

AN INVESTIGATION OF FACTORS INFLUENCING
ALGORITHM SELECTION FOR HIGH DIMENSIONAL
CONTINUOUS OPTIMISATION PROBLEMS

KEVIN GRAHAM



Doctor of Philosophy

Institute of Computing Science and Mathematics

University of Stirling

November 2019

DECLARATION

I, Kevin Graham, hereby declare that all material in this thesis is my own original work except where reference has been given to other authors. The work presented here has not been submitted for the award or partial award of any other degree at the University of Stirling nor any other institute.

Stirling, November 2019

Kevin Graham

ABSTRACT

The problem of algorithm selection is of great importance to the optimisation community, with a number of publications present in the Body-of-Knowledge. This importance stems from the consequences of the No-Free-Lunch Theorem which states that there cannot exist a single algorithm capable of solving all possible problems. However, despite this importance, the algorithm selection problem has of yet failed to gain widespread attention . In particular, little to no work in this area has been carried out with a focus on large-scale optimisation; a field quickly gaining momentum in line with advancements and influence of big data processing. As such, it is not as yet clear as to what factors, if any, influence the selection of algorithms for very high-dimensional problems (> 1000) - and it is entirely possible that algorithms that may not work well in lower dimensions may in fact work well in much higher dimensional spaces and vice-versa. This work therefore aims to begin addressing this knowledge gap by investigating some of these influencing factors for some common metaheuristic variants.

To this end, typical parameters native to several metaheuristic algorithms are firstly tuned using the state-of-the-art automatic parameter tuner, SMAC. Tuning produces separate parameter configurations of each metaheuristic for each of a set of continuous benchmark functions; specifically, for every algorithm-function pairing, configurations are found for each dimensionality of the function from a geometrically increasing scale (from 2 to 1500 dimensions). The nature of this tuning is therefore highly computationally expensive necessitating the use of SMAC. Using these sets of parameter configurations, a vast amount of performance data relating to the large-scale optimisation of our benchmark suite by each metaheuristic was subsequently generated.

From the generated data and its analysis, several behaviours presented by the metaheuristics as applied to large-scale optimisation have been identified and discussed. Further, this thesis provides a concise review of the relevant literature for the consumption of other researchers looking to progress in this area in addition to the large volume of data produced, relevant to the large-scale optimisation of our benchmark suite by the applied set of common metaheuristics.

All work presented in this thesis was funded by EPSRC grant: EP/J017515/1 through the DAASE project ¹.

¹ <http://daase.cs.ucl.ac.uk/>

ACKNOWLEDGMENTS

“Let not your mind run on what you lack as much as on what you already have”

– *Marcus Aurelius*

“Why do we fall, Bruce? So we can learn to pick ourselves up”

– *Thomas Wayne (Batman Begins 2005)*

I am dedicating this thesis to my loving wife Tracy, my sons Carson and Callan and my dad Michael Graham.

My dad sadly passed away from an asbestos related disease at the beginning of my 2nd year - a particularly difficult time for family and myself. Of all the occasions when things got tough during the course of my studies, one of the things that kept me going was his pride that I had gotten to where I was; not through luck or innate ability, but like he did his whole life - through hard work and dedication. Being a family-centric man, I still believe that hard work on its own could not have brought me to this point, but the love and support from my dad and my family throughout my life. We all miss his presence, humour and love but his memory will live on through his wife - my mum, sons and grandchildren for years to come.

My PhD years were not all filled with grief however, I married my wife Tracy and 3 years later we had our first son, Carson Michael Graham. I simply would not be here writing this now if not for the unconditional love and support Tracy has shown me over the years - extending far beyond PhD studies. Whenever I needed to vent, talk things over or was in dire need of mental support I could always count on her to be there for me - even when I needed to be ‘rescued’ from Nottingham University when I fell ill with severe appendicitis. I really could not have done any of this were it not for her.

I make no secret of the fact about how many times - and how close I came - to throwing in the towel. The illogical side of my psyche had me convinced that the PhD was ‘cursed’ - I had lost multiple supervisors (through work opportunities), family members, had several failed/abandoned projects, suffered severe illness and was getting nowhere fast despite my efforts. That’s until my son Carson was born, and then my son Callan Thomas Graham just short of 2 years later. Every day since then I vowed that I would show them how to not give

up, to keep fighting through obstacles in your way in order to become the best version of yourself and that getting knocked down only teaches us how to get back up - I didn't want their first lesson in life to be to stop moving forward.

And so my sons, this is for you, your mummy and your grandad - Always remember that you are braver than you believe, stronger than you seem, smarter than you think and more loved than you'll ever know.

Finally, I would like to thank my current supervisors Prof. Leslie Smith, Dr. Alexander (Sandy) Brownlee and Dr. David Cairns for all their help and support over the last two years - without whom I could not have reached this point. Also, I'd like to give thanks to a former supervisor - Dr John Woodward, who never gave up trying to organise me, even up to the day he took up his new appointment. Also, I'd like to thank my fellow PhD students for all the years of interesting discussion and laughter. In particular, I'd like to thank Jason Adair (now Dr. Adair) for all the help and support (both personally and work-based), chats and putting up with my long venting sessions on the not-so-rare occasion - and of whom will always be considered as a friend. Of course, I cannot forget my mum who has always shown me love and support throughout my life and who taught ME that getting knocked down only teaches us how to get back up.

LIST OF PUBLICATIONS

[42] Kevin Graham, Jerry Swan, and Simon Martin. The Blackboard Pattern for Metaheuristics. In *Proceedings of the Companion Publication of the 2015 Annual Conference on Genetic and Evolutionary Computation*, pages 1265-1267, 2015.

[41] Kevin Graham and Leslie Smith. Comparing Hyper-heuristics with Blackboard Systems. *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, pages 1141-1145, 2017.

The publications listed here are related to previous research directions no longer being followed and are therefore not relevant to this thesis.

CONTENTS

i INTRODUCTION	1
1 INTRODUCTION TO THE THESIS	2
1.1 Introduction	2
1.2 Structure of the Thesis	3
ii BACKGROUND & LITERATURE REVIEW	5
2 CHAPTER 2: OPTIMISATION	6
2.1 Introduction	6
2.2 General Optimisation	7
2.2.1 Terminology	7
2.2.2 A Formal Definition of an Optimisation Problem	9
2.2.3 Problem Difficulty	9
2.3 Classical Optimisation Methods	11
2.3.1 Strong Methods	12
2.3.2 Weak Methods	16
2.4 Metaheuristic Optimisation	17
2.4.1 What is a Metaheuristic?	17
2.4.2 Exploration vs. Exploitation in Metaheuristics	19
2.4.3 Problems Addressed by Metaheuristics	20
2.4.4 Categorisations of Metaheuristics	20
2.4.5 General Concepts	22
2.5 Large-Scale Global Optimisation	31
2.5.1 Cooperative Co-Evolutionary Algorithm	34
3 CHAPTER 3: CONTINUOUS OPTIMISATION BENCHMARK FUNCTIONS	42
3.1 Introduction	42
3.2 What Makes a Function Difficult to Optimise	42
3.3 Constructing a Custom Function Suite Versus Using Existing Benchmark Function Suites	46
3.4 Functions	48
3.4.1 Summary of Functions	49
4 CHAPTER 4: METAHEURISTIC APPROACHES USED IN THIS THESIS	50

4.1	Introduction	50
4.1.1	Exposing Algorithm Parameters	51
4.1.2	Solution Representation	52
4.1.3	Neighbourhood Function	52
4.2	Hill Climbing Algorithms	54
4.2.1	Common Variations	55
4.2.2	Global Optimisation with Hill Climbing	58
4.2.3	Random Mutation Hill Climb: Implementation	61
4.3	Genetic Algorithms (GA)	65
4.3.1	The Basic GA Procedure	67
4.3.2	Genetic Operators	69
4.3.3	Steady State GA vs. Generational GA	78
4.3.4	Steady State Genetic Algorithm: Implementation	79
4.4	Particle Swarm Optimisation (PSO)	85
4.4.1	Neighbourhood Topologies	85
4.4.2	Parameters of PSO	87
4.4.3	Boundary-Handling in PSO	90
4.4.4	Particle Swarm Optimisation (PSO): Implementation	92
4.5	Simulated Annealing (SA)	98
4.5.1	Annealing in Condensed Matter Physics	99
4.5.2	Acceptance Probability	99
4.5.3	Cooling Schedules	100
4.5.4	Simulated Annealing (SA): Implementation	101
4.6	Differential Evolution (DE)	107
4.6.1	The Differential Evolution Process	107
4.6.2	Variants of DE	110
4.6.3	Differential Evolution (DE): Implementation	110
4.7	Covariance Matrix Adaption Evolutionary Strategy (CMA-ES)	115
4.7.1	Algorithm Overview	115
4.7.2	CMA-ES for Large-scale Optimisation	116
4.7.3	Covariance Matrix Adaption Evolution Strategy (CMA-ES): Implementation	117
5	CHAPTER 5: THE PROBLEMS OF ALGORITHM SELECTION AND CONFIGURATION	120
5.1	Introduction	120
5.2	Relationships Between Algorithm Selection and Other Areas of Research	121
5.2.1	Meta-learning	121

5.3	Automatic Parameter Tuning	125
5.3.1	Model-free vs. Model-based Methods	126
5.4	SMAC: Sequential Model-based Algorithm Configuration	129
5.4.1	The SMAC Tuning Procedure	129
iii	METHOD	132
6	CHAPTER 6: ALGORITHM PARAMETER CONFIGURATION	133
6.1	Introduction	133
6.2	Parameter Tuning Method	134
6.2.1	Parameter Identification and Range Selection	135
6.2.2	Conducting the Tuning Procedures	136
6.2.3	Generation of Tuned Performance Data	136
6.2.4	Obtaining a Final Parameter Configuration Set	137
6.3	Improving the Generality of Empirical Results and Observations	138
6.4	Result of Parameter Tuning with SMAC	139
iv	RESULTS AND CONCLUSIONS	141
7	CHAPTER 7: EMPIRICAL RESULTS	142
7.1	Algorithm Tuning Results	142
7.1.1	Comparing SMAC and Irace Tuning Performance	143
7.2	Function Evaluation Budget Comparison	145
7.2.1	Effects Observed when Tuning with 10,000 and 50,000 Function Evaluation Budgets	145
7.3	Performance Comparison Between Algorithm Implementations	154
7.3.1	Random Mutation Hill Climbing (RMHC)	154
7.3.2	Changes to Rank Ordering of Performance at Larger Scales	155
7.4	Conclusions and Discussion	164
7.4.1	Tuning Evaluation Budget Comparisons	164
7.4.2	Performance Comparison Between Algorithm Implementations	166
7.4.3	Rank Ordering Changes at Larger Problem Scales	166
7.5	Contributions	166
7.6	Future Work	167
v	APPENDICES	1
A	APPENDIX A: 10,000 AND 50,000 EVALUATION COMPARISON PLOTS	2
B	APPENDIX B: 10K & 50K EVALUATION TUNING BUDGET DESCRIPTIVE STATISTICS	32

C	APPENDIX C: SMALL TO LARGE-SCALE METAHEURISTIC COMPARISON MATERIAL	51
C.1	Metaheuristic Comparison Plots	52
C.1.1	Comparison of Approaches Tuned at the 10,000 Evaluation Budget (With Q1 & Q3 Error Bars)	52
C.1.2	Comparison of Approaches Tuned at the 10,000 Evaluation Budget (No Error Bars)	55
C.1.3	Comparison of Approaches Tuned at the 50,000 Evaluation Budget (With Q1 & Q3 Error Bars)	58
C.1.4	Comparison of Approaches Tuned at the 50,000 Evaluation Budget (No Error Bars)	61
C.2	Graphical Metaheuristic Rankings	64
C.2.1	Rank Ordering By Benchmark Function	64
C.2.2	Rank Ordering: Mean Aggregate and Pooled Mean Over Subsets of the Benchmark Suite	67
C.3	Algorithm Ranking Tables	69
C.3.1	Algorithm Rank Ordering By Benchmark Function	69
C.3.2	Mean Aggregated Algorithm Rank Ordering By All Problems and Problems Classified by Problem Feature Investigated	79
D	APPENDIX D: PLOTS SHOWING THE RESULTS OF TUNING EACH ALGORITHM AT A 10,000 AND 50,000 EVALUATION BUDGET COMPARED TO THE CORRESPONDING UNTUNED ALGORITHM	83
D.1	SMAC Parameter Tuning Results	84
D.2	Irace vs. SMAC Parameter Tuning Results	114
E	APPENDIX E: DEFINITIONS OF BENCHMARK FUNCTIONS IN THE TEST SUITE	132
E.1	Ackley Function no.1	132
E.2	Alpine Function no.1	132
E.3	Bent Cigar Function	133
E.4	Brown Function	134
E.5	Chung-Reynolds Function	135
E.6	Deflected Corrugated Spring Function	136
E.7	Exponential Function	137
E.8	Griewank Function	138
E.9	Inverted Cosine Wave Function	139
E.10	Levy Function	140
E.11	Qing Function	142

E.12	Rastrigin Function	143
E.13	Rosenbrock Function	144
E.14	Schwefel Function no.26	146
E.15	Sphere Function	147
E.16	Sum of Different Powers Function	148
E.17	Sum Squares Function	148
F	SUPPLEMENTARY EXPERIMENT MATERIAL	150
F.1	Extended PSO Experiment Results	150

LIST OF FIGURES

Figure 3.1	Barchart Summarising Prevalence of Problems With Certain Features Present in Each Benchmark Test Suite	47
Figure 4.1	Effects of inappropriate fixed step sizes	53
Figure 4.2	Definition of a neighbourhood of a solution sol using the stepsize parameter and the stochastic generation of a neighbouring solution sol_2	64
Figure 4.3	Basic operation of a typical genetic algorithm implementation	68
Figure 4.4	Roulette Wheel Selection example showing a single selection over a population of four individuals. Each individual, A,B,C and D represent a 32.7%, 12.7%, 3.6% and 51% portion of the roulette wheel respectively. Random point r selects the portion representing chromosome D similar to a hypothetical roulette ball	70
Figure 4.5	Tournament selection, where $t = 8$ (tournament size)	72
Figure 4.6	Uniform Crossover	73
Figure 4.7	One-point Crossover	73
Figure 4.8	Uniform Mutation: where x is transformed to x'	75
Figure 4.9	Velocity and Position Update in PSO	88
Figure 4.10	Production of a mutant vector \vec{V}_i from the scaled difference vector $(\vec{X}_{r_2} - \vec{X}_{r_3})$	109
Figure 5.1	Taken from [58], This figure shows two steps of SMBO (EGO) for the optimization of a 1D function. The true function is shown as a solid line, and the circles denote our observations. The dotted line represents the mean prediction of a noise-free Gaussian process model (the DACE model), with the grey area representing its uncertainty. Expected improvement, EI, (scaled by the authors for visualisation) is shown by a dashed line.	128
Figure 6.1	Example of a parameter configuration file compiled from the individual problem-dimension configurations obtained from SMAC	137
Figure 6.2	Result Obtained from DE Implementation (tuning at 50,000 function evaluations) vs. Untuned Version	140
Figure 7.1	Examples of Successful Tunings of RMHC and SSGA Using SMAC on Ackleys Function and Alpine n.1	143

Figure 7.2	Example of GA Continuing to Progress Beyond 10,000 Function Evaluations for Ackley’s Function	147
Figure 7.3	SSGA Performance Results Comparison (Irace Tuned) Example - 10,000 and 50,000 Tuning Evaluation Budgets	148
Figure 7.4	SSGA: Divergence of Performance between 10,000 and 50,000 Tuning Evaluation Budgets of SMAC and Irace	149
Figure 7.5	Example showing divergence between parameter tuning evaluation budgets for PSO (Alpine n.1 Function) at dimensions 11-16 and re-convergence occurring at dimension 128	150
Figure 7.6	Plots showing individual values for each tuned parameter from our original 5 parameter configurations sets obtained from SMAC	151
Figure 7.7	Comparison of Performance (Median Solution Quality of 20 Independent Repeats) Against the Chung-Reynolds Function Between Original and Extended Scale PSO Experiments	153
Figure 7.8	Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Benchmark Functions (where $n = 5$)	157
Figure 7.9	Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Uni-modal Functions (where $n = 5$)	159
Figure 7.10	Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Multi-modal Functions (where $n = 5$)	159
Figure 7.11	Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Separable Functions (where $n = 5$)	162
Figure 7.12	Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Non-separable Functions (where $n = 5$)	162
Figure E.1	Ackley Function no.1 in its 2-dimensional form	132
Figure E.2	Alpine Function no.1 in its 2-dimensional form	133
Figure E.3	Bent Cigar Function in its 2-dimensional form	134
Figure E.4	Brown Function in its 2-dimensional form: (a) and (b) shows the function’s full domain, and (c) and (d) shows the function with domain limited to $-1 \leq x_i \leq 1$	135
Figure E.5	Chung-Reynolds Function in its 2-dimensional form	136

Figure E.6	Deflected Corrugated Spring Function in its 2-dimensional form	137
Figure E.7	Exponential Function in its 2-dimensional form	138
Figure E.8	Griewank Function in its 2-dimensional form	139
Figure E.9	Inverted Cosine Wave Function in its 2-dimensional form	140
Figure E.10	Levy Function in its 2-dimensional form	141
Figure E.11	Qing Function in its 2-dimensional form: (a) and (b) shows the function's full domain, and (c) and (d) shows the function with domain limited to $-2 \leq x_i \leq 2$	143
Figure E.12	Rastrigin Function in its 2-dimensional form	144
Figure E.13	Rosenbrock Function in its 2-dimensional form: (a) and (b) shows the function's full domain, and (c) and (d) shows the function with domain limited to $-2.5 \leq x_i \leq 2.5$	145
Figure E.14	Schwefel Function in its 2-dimensional form	146
Figure E.15	Sphere Function in its 2-dimensional form	147
Figure E.16	Sum of Different Powers Function in its 2-dimensional form	148
Figure E.17	Sum Squares Function in its 2-dimensional form	149

LIST OF TABLES

Table 3.1	Summary of the Proportion of Problems With Certain Features Present in Various Existing Test Suites Versus the Test Suite Constructed for this Study	47
Table 3.2	Summary of n-Dimensional Benchmark Function Suite	49
Table 7.1	Plot summary of problem dimensionalities where algorithm performance between 10k and 50k evaluations diverge in favour of tuning at 50k function evaluations	146
Table 7.2	Pooled Mean Rank Table Over All Functions	157
Table 7.3	Table of Pooled Mean Ranks Over All Uni-modal Functions	160
Table 7.4	Table of Pooled Mean Ranks Over All Multi-modal Functions	160
Table 7.5	Table of Pooled Mean Ranks Over All Separable Functions	163
Table 7.6	Table of Pooled Mean Ranks Over All Non-separable Functions	163
Table B.1	sep-CMA-ES Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (10,000 Evaluation Tuning)	33
Table B.2	sep-CMA-ES Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (50,000 Evaluation Tuning)	33
Table B.3	sep-CMA-ES Descriptive Statistics: Exponential to Rastrigin Functions (10,000 Evaluation Tuning)	34
Table B.4	sep-CMA-ES Descriptive Statistics: Exponential to Rastrigin Functions (50,000 Evaluation Tuning)	34
Table B.5	sep-CMA-ES Descriptive Statistics: Rosenbrock to Sum Squares Functions (10,000 Evaluation Tuning)	35
Table B.6	sep-CMA-ES Descriptive Statistics: Rosenbrock to Sum Squares Functions (50,000 Evaluation Tuning)	35
Table B.7	DE Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (10,000 Evaluation Tuning)	36
Table B.8	DE Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (50,000 Evaluation Tuning)	36
Table B.9	DE Descriptive Statistics: Exponential to Rastrigin Functions (10,000 Evaluation Tuning)	37

Table B.10	DE Descriptive Statistics: Exponential to Rastrigin Functions (50,000 Evaluation Tuning)	37
Table B.11	DE Descriptive Statistics: Rosenbrock to Sum Squares Functions (10,000 Evaluation Tuning)	38
Table B.12	DE Descriptive Statistics: Rosenbrock to Sum Squares Functions (50,000 Evaluation Tuning)	38
Table B.13	GA Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (10,000 Evaluation Tuning)	39
Table B.14	GA Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (50,000 Evaluation Tuning)	39
Table B.15	GA Descriptive Statistics: Exponential to Rastrigin Functions (10,000 Evaluation Tuning)	40
Table B.16	GA Descriptive Statistics: Exponential to Rastrigin Functions (50,000 Evaluation Tuning)	40
Table B.17	GA Descriptive Statistics: Rosenbrock to Sum Squares Functions (10,000 Evaluation Tuning)	41
Table B.18	GA Descriptive Statistics: Rosenbrock to Sum Squares Functions (50,000 Evaluation Tuning)	41
Table B.19	PSO Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (10,000 Evaluation Tuning)	42
Table B.20	PSO Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (50,000 Evaluation Tuning)	42
Table B.21	PSO Descriptive Statistics: Exponential to Rastrigin Functions (10,000 Evaluation Tuning)	43
Table B.22	PSO Descriptive Statistics: Exponential to Rastrigin Functions (50,000 Evaluation Tuning)	43
Table B.23	PSO Descriptive Statistics: Rosenbrock to Sum Squares Functions (10,000 Evaluation Tuning)	44
Table B.24	PSO Descriptive Statistics: Rosenbrock to Sum Squares Functions (50,000 Evaluation Tuning)	44
Table B.25	RMHC Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (10,000 Evaluation Tuning)	45
Table B.26	RMHC Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (50,000 Evaluation Tuning)	45
Table B.27	RMHC Descriptive Statistics: Exponential to Rastrigin Functions (10,000 Evaluation Tuning)	46

Table B.28	RMHC Descriptive Statistics: Exponential to Rastrigin Functions (50,000 Evaluation Tuning)	46
Table B.29	RMHC Descriptive Statistics: Rosenbrock to Sum Squares Function (10,000 Evaluation Tuning)	47
Table B.30	RMHC Descriptive Statistics: Rosenbrock to Sum Squares Function (50,000 Evaluation Tuning)	47
Table B.31	SA Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (10,000 Evaluation Tuning)	48
Table B.32	SA Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (50,000 Evaluation Tuning)	48
Table B.33	SA Descriptive Statistics: Exponential to Rastrigin Functions (10,000 Evaluation Tuning)	49
Table B.34	SA Descriptive Statistics: Exponential to Rastrigin Functions (50,000 Evaluation Tuning)	49
Table B.35	SA Descriptive Statistics: Rosenbrock to Sum Squares Functions (10,000 Evaluation Tuning)	50
Table B.36	SA Descriptive Statistics: Rosenbrock to Sum Squares Functions (50,000 Evaluation Tuning)	50
Table C.1	Algorithm Rankings for Ackley Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	70
Table C.2	Algorithm Rankings for Alpine n.1 Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	70
Table C.3	Algorithm Rankings for Bent Cigar Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	71
Table C.4	Algorithm Rankings for Brown Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	71
Table C.5	Algorithm Rankings for Chung-Reynolds Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	72
Table C.6	Algorithm Rankings for Deflected Corrugated Spring Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	72
Table C.7	Algorithm Rankings for Exponential Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	73
Table C.8	Algorithm Rankings for Griewank Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	73
Table C.9	Algorithm Rankings for Inverted Cosine Wave Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	74

Table C.10	Algorithm Rankings for Levy Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	74
Table C.11	Algorithm Rankings for Qing Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	75
Table C.12	Algorithm Rankings for Rastrigin Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	75
Table C.13	Algorithm Rankings for Rosenbrock Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	76
Table C.14	Algorithm Rankings for Schwefel Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	76
Table C.15	Algorithm Rankings for Sphere Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	77
Table C.16	Algorithm Rankings for Sum of Different Powers Function over 2 to 1448 Dimensions (1=Best, 5=Worst)	77
Table C.17	Algorithm Rankings for Sum Squares Function over 2 to 1448 Dimen- sions (1=Best, 5=Worst)	78
Table C.18	Mean Aggregated Algorithm Rankings Over All Benchmark Functions	80
Table C.19	Mean Aggregated Algorithm Rankings Over All Uni-modal Functions	80
Table C.20	Mean Aggregated Algorithm Rankings Over All Multi-modal Functions	81
Table C.21	Mean Aggregated Algorithm Rankings Over All Separable Functions .	81
Table C.22	Mean Aggregated Algorithm Rankings Over All Non-separable Functions	82

LIST OF ACRONYMS

acp	Algorithm Configuration Problem
apt	Automatic Parameter Tuning
asp	Automatic Selection Problem
bcd	Binary Coded Decimal
ccd	Central Composite Design
cdiw	Chaotic Descending Inertial Weight
cma-es	Covariance Matrix Adaption - Evolution Strategy
criw	Chaotic Random Inertial Weight
de	Differential Evolution
ga	Genetic Algorithm
gga	Generational Genetic Algorithm
ils	Iterated Local Search
nfl	no free lunch (theorem)
pso	Particle Swarm Optimisation
psosa	Particle Swarm Optimisation Simulated Annealing
pso-tvac	Particle Swarm Optimisation - Time Varying Acceleration Coefficients
rmhc	Random Mutation Hill Climbing
rws	Roulette Wheel Selection
sa	Simulated Annealing
sga	Simple Genetic Algorithm
shc	Stochastic Hill Climbing
smac	Sequential Model-based Algorithm Configuration

ssga Steady-State Genetic Algorithm

sus Stochastic Universal Sampling

vns Variable Neighbourhood Search

Part I

INTRODUCTION

INTRODUCTION TO THE THESIS

1.1 INTRODUCTION

Human beings are a species of optimisers. Throughout history, advances in civilisation can be said to have been built, almost solely, on our ability to optimise the various processes essential to our survival. From our early ancestors gradually improving on their tools for hunting or gathering to developing more efficient means to automatically mass produce items of necessity during the industrial period, we owe much of our current, and relatively comfortable, state of being to our insatiable desire to 'find something better'.

As our mathematics has developed, an ever increasing number of processes of increasing complexity are being modelled. Ever more complex models from the physical sciences, natural sciences and human sciences have provided us with new and interesting insights and new optimisation problems to solve. From molecular biology, for example, we have the problem of protein folding...

However, many of these real world problems can have very large search spaces; that is, the number of possible solutions to the problem can be incredibly large - increasing exponentially with the dimensionality of the problem. This makes it difficult for more traditional and exact optimisation methods, e.g., exhaustive search, to effectively search through all these possible solutions in a time that is acceptable. In terms of optimisation in continuous spaces; that is, where the problem variables are from the set of real numbers, these traditional approaches often make use of function derivatives. However, many continuous problems have search spaces that are inherently non-differentiable and thus inhibiting the use of these strategies.

For exactly these reasons, heuristic based optimisation methods - able to effectively search within arbitrarily large search spaces in a reasonable length of time - have become increasingly popular with several thousand associated publications even in the last year alone (2017). However, a fundamental theorem of any search based optimisation strategy - the No Free Lunch Theorem (NFL) - states that no single search algorithm, heuristic or otherwise, can be shown to exhibit higher than average performance when all possible problems are considered.

The consequences of the NFL therefore pose a problem to the users and developers of search based optimisation algorithms; specifically, what algorithm or algorithms will perform best on a given problem or set of problems of a certain type?

This problem, dubbed the *Algorithm Selection Problem (ASP)* by Rice [112] in 1975, has been of some interest to researchers however despite its importance it has failed to gain widespread attention in the research community. Furthermore, of the works that do exist, none were found to address the issue of algorithm selection for large-scale optimisation problems - those of very high dimensionality. This is an important omission, and one we attempt to begin addressing - at least partially - because the exponential increase in search space volume that invariably results from increases to problem dimensionality means that the effectiveness - at least in terms of solution quality - of all search algorithms, including heuristic algorithms, diminishes with dimensionality. Some strategies may however remain reasonably effective on certain problems at higher dimensions - where their rate of decreasing performance with dimensionality is slower - but no data is yet available to even attempt to make such discriminations.

1.2 STRUCTURE OF THE THESIS

We begin in Chapter 2 our review of the available subject matter. In this chapter we cover the subject of optimisation along with descriptions of some more ‘classical’ algorithmic approaches to addressing continuous optimisation problems.

In Chapter 3, we provide a description of our selected suite of 17 benchmark functions used throughout the experiments in this thesis. Along with the descriptions, we also provide visual (2D) representations of each function we generated through sampling of the functions implemented. A review of the metaheuristic search-based optimisation algorithms being used in these studies is then provided in Chapter 4.

We then briefly cover both the problem of algorithm selection and automatic parameter tuning (APT), in particular of the former - our selected tuning method, Sequential Model-based Algorithm Configuration (SMAC) in Chapter 5.

Beginning our method part with Chapter 6, attention is firstly given to the metaheuristic implementations developed for use in our experiments, including, for each approach, a list

and discussion of the various ‘tunable parameters’ configured by our selected APT described in Chapter 5. To supplement this description, and to provide more specificity to our use of SMAC, Chapter 7 provides a brief description of the steps taken to tune each metaheuristic for each of our benchmark functions at each dimensionality.

Finally, in Chapter 8, we blend our experimental methodology with individual findings and discussions and for each, provide our conclusions. These conclusions are then summarised at the end of the chapter along with a discussion of future work directions.

Part II

BACKGROUND & LITERATURE REVIEW

CHAPTER 2: OPTIMISATION

2.1 INTRODUCTION

The task of optimisation has become one of the main staples of modern day science, finding utility in almost every scientific field of enquiry. In fact, Schwefel in [118] states:

“There is scarcely a modern journal, whether of engineering, economics or the social sciences, in which the concept ‘optimization’ is missing from the subject index.”

Given our long relationship with optimisation and its many successes, having been paramount to our continued advancement as a species for thousands of years, this is not too surprising.

However, as problems have become increasingly more complex, compounded with the fact that in many cases we do not know the mathematical formulation of our problem, we are less able to rely on some of our more ‘classical’ or exact methods to optimisation. This has led to the popularisation of stochastic optimisation methods, specifically those techniques such as metaheuristics and hyper-heuristics, that are more able to deal with these intractable problems in a reasonable time scale albeit with no guarantee of exact optimality. Therefore, metaheuristics and other stochastic approaches have their place in situations where a ‘good enough’ solution is acceptable or no known exact algorithm can find the optimal solutions in reasonable runtimes.

This chapter begins by discussing some of the background and terminology of optimisation in general, including: a definition of optimisation, an overview of how problem difficulty is defined and what this means to optimisation on the whole and a comparison of the main categories of optimisation: discrete vs. continuous and local vs. global optimisation. Next we present an overview of some of the more classical and exact methods to optimisation and how these approaches are less able to perform efficiently when problems being approached become larger and more complex. This leads us give a relatively concise discussion on the topic of metaheuristics, specifically how they are applied to help solve continuous optimisation problems.

2.2 GENERAL OPTIMISATION

2.2.1 Terminology

Here we present some terminology useful to further discussions in this chapter.

2.2.1.1 Objective Function

An objective function is an encoding of the goal of some problem as a mathematical function or computational model used to measure the quality of a given solution in respect to the problem. It essentially represents the interface between an algorithm and the real problem and plays a crucial role in the guidance of a search algorithm through a given search space - without it the algorithm would have no clue as to how the solutions produced are performing. The problem goal, or goals in the case of multi-objective problems, can be defined as either being a minimisation goal, where the objective is to produce solutions with low values with respect to the objective function, or a maximisation problem where higher valued solutions are sought. When indirect solution representations are used within the algorithm - those that do not yet represent a solution in their current form i.e., they are encoded so must first be converted (decoded) into a form that can be directly evaluated by the objective function [17].

2.2.1.2 Feasible vs. Infeasible Solutions

Feasible solutions are all solutions within the search space of the problem which satisfy the constraints on the problems parameters - specifically, the hard constraints (see below) [17]. Conversely, infeasible solutions all violate at least one of the problems hard constraints.

2.2.1.3 Hard vs. Soft Constraints

Hard constraints placed on a problem are conditions which have to be satisfied in order for a solution to be feasible [17]. On the other hand, soft constraints are those which we would like to have satisfied but which are not essential; so a solution violating one or more soft-constraints would still be considered feasible [17]. A common way of handling soft constraints when searching for a solution to an optimisation problem is to have solutions incur a penalty if a soft constraint is not satisfied; possibly weighted by desirability. Therefore, the objective function of many problems are represented by the summation of penalties for the soft constraints [17]. This means that solutions that do not violate the soft constraints will be considered as being of a higher quality [17].

2.2.1.4 *Local vs. Global Optimisation*

The difference between local and global optimisation lies in the scope of the optimisation task. Local optimisation looks to optimise within a local feasible region [97], that is, to discover a local minimiser (or optimum) - the smallest objective value within some feasible neighbourhood in the problem domain. A neighbourhood function N defines the neighbourhood of a solution $s \in S$ as a mapping $N : S \rightarrow 2^S$ which assigns to each s a set of solutions $N(s) \subset S$ [135]. A solution $s' \in N(s)$ is referred to as a neighbour of s [135]. Therefore, given a domain S and a feasible neighbourhood of the domain $N(s) \subset S$, local optimisation seeks to find a local optimum $x^* \in N$ such that $f(x^*) \leq f(x) \forall x \in N$ [135]. Global optimisation on the other hand is concerned with discovering the smallest objective value or values over the entire feasible domain of a function - a global optimum, also often referred to as the 'best solution'. As before, given a domain S , global optimisation seeks to find a global optimum $s^* \in S$ such that $f(s^*) \leq f(s) \forall s \in S$ [135]. The main outcome of optimisation is therefore to discover a global optimum s_* of which many such solutions may exist [135]. Where many global optima may be present (it is not always known), the outcome may be defined as attempting to discover all of the global optima in order to generate alternative choices of solution [135].

2.2.1.5 *Best Solution*

Using a term like the best solution implies that there must be more than one solution, all of which are considered to have differing levels of value [50]. The meaning of best, in terms of quality, is quite often depends on the actual problem, the solution method and any allowed tolerances [50]. Some types of problem have exact answers where the term 'best' has an absolute definition [50] - in optimisation terms, 'best' in these instances refers to a problem's single global optimum. Other problems however can have several global optima, so 'best' in these situations is more of a relative term [50]

2.2.1.6 *Discrete vs. Continuous Optimisation*

Discrete optimisation involves the optimisation of problems which are characterised by a finite number of states accepted by their input parameters; on the other hand, continuous problems are typically those that take their solutions from an uncountably infinite set of solutions whose parameters are members of the set of real values - of which there are no discontinuities [97]. However, given that digital computer technologies cannot provide the kind of precision required of a 'true' continuous solution, the set of possible solutions will be far less than uncountably infinite albeit still intractably large.

Although the practical optimisation of both types of problem can be approached using similar methods - at least in terms of inexact heuristic methods - there are nuances found for each type of problem that need to be taken into account in order to optimise effectively. For example, continuous optimisation problems are considered easier to solve than discrete problems as their typical smoothness make it possible to work out the functions behaviour at all points close to a point x by using objective and constraint information [97]. On the other hand, the behaviour of the objective function may change drastically between points - which are deemed close by some measure of distance [97].

2.2.2 A Formal Definition of an Optimisation Problem

Burke and Kendall in [17] define optimisation as trying to find the best solution possible from amongst all possible solutions. They state further that optimisation can therefore be considered as the task of modelling the problem to be solved as a mathematical evaluation (or objective) function representing the quality of a solution and then search through the space of all possible solutions in order to find one that either minimises or maximises this function [17].

The canonical form for an optimisation problems is usually stated as the minimisation of the objective function subject to constraints placed on its variables and expressed as [97, 6]:

$$\min f(x) : x \in \mathbb{R}^n \quad \text{subject to} \quad \begin{cases} g_j(x) \geq 0, & j = \{1, 2, \dots, J\} \\ h_k(x) = 0, & k = \{1, 2, \dots, K\} \end{cases} \quad (2.1)$$

$$\text{such that: } x_i^{(L)} \leq x_i \leq x_i^{(U)}, \quad i = \{1, 2, \dots, n\}$$

where $f(x)$ is the function being optimised, $g_j(x)$ are J inequality constraints and $h_k(x)$ represent the K equality constraints [6].

2.2.3 Problem Difficulty

Michalewicz and Fogel in [90] state several different reasons why real-world optimisations are difficult to solve effectively:

- The search space is too vast to allow the use of an exhaustive search

- The problem is so complex that in order to determine any answer whatsoever requires the use of such simplified models of the problem domain that any solution found is essentially useless
- The objective function is noisy or varies over time therefore requiring that an entire series of solutions are required as opposed to a single solution
- The search space is so heavily constrained that even finding one feasible solution to the problem is difficult

2.2.3.1 Complexity and P vs. NP

Complexity refers to the study of how difficult an optimisation problem is to solve [17]. Here, problems are classified according to the properties of optimisation algorithms where roughly speaking, a problem is classified as being ‘hard’ if there exists no known fast solver and ‘easy’ otherwise [29].

The first key concept in computational complexity is that of problem size which is rooted in the dimensionality of a given problem - that is, the number of problem parameters - and the size of the set of values of which each of the problem parameters can be set [29].

The second key concept relates to algorithms rather than to the problems themselves, that of the running-time [29]. Running-time is the number of “elementary” operations required of an algorithm before it terminates when running against a given problem - in general, the intuition of computational complexity is that larger problems will require more computational time to solve, although this is not always true [29]. The best known definitions of the ‘hardness’ of a problem is a relationship between the size of the problem being solved to the worst-case runtime of the algorithm to be used to solve it. This is encapsulated as a formula that defines an upper bound on runtime of the algorithm as a function of the problem size [29]. Basically, shorter running times are expected when this formula is a polynomial or longer when the formula is ‘super-polynomial’ such as exponential [29].

The final key concept of computational complexity relates to the concept that it is possible to transform one problem into another through an appropriate mapping which may or may not be reversible. This is referred to as *problem reduction* [29].

These key concepts can now be described more formally as [29]:

1. A problem can be said to belong to the class P (polynomial) if there exists at least one algorithm capable of solving it in polynomial time [29]
2. A problem can be said to belong to class NP (Non-deterministic polynomial) when it can be solved by some algorithm, specifically a non-deterministic turing machine, with no runtime guarantees, and its solution can be tested in polynomial time [29]. P can be considered as a subset on NP since a solver with a polynomial runtime can also be used to test solutions in polynomial time [29]
3. A problem is said to belong to the class NP-complete when it belongs to the class NP and another problem in NP can be reduced to this particular problem by an algorithm that runs in polynomial time [29]
4. A problem belongs to the class NP-hard if (i) it is at least as difficult as any problem in NP-complete, but (ii) where solutions cannot be necessarily tested in polynomial time - the canonical example being the halting problem [29]

One of the 'grand' challenges in complexity theory is to provide a proof that $P = NP$ or conversely that $P \neq NP$, in other words, is the class P and the class NP in fact the same class - or not, in the latter case [29]. It should be noted that the proof required of the latter case will need to be acquired through the application of complex mathematics, but for the former, 'simply' showing that an algorithm exists that can solve at least problem of class NP-complete would provide enough of a proof [29].

2.3 CLASSICAL OPTIMISATION METHODS

Before continuing on further to discuss stochastic optimisation, it is worth outlining some of the existing 'classical' mathematical optimisation techniques. Wehrens and Buydens [142] categorise the following methods under two headings: *strong* methods and *weak* methods.

Strong methods make certain assumptions about the structure of the solutions space; if the assumptions made are correct then strong methods are both fast and reliable - however, if they are incorrect the methods will never manage to find the global optimum even through increased repetitions [142]. They are often used in the final part of the optimisation of difficult problems where the region of the global optimum has been found but an exact solution has yet to be discovered [142]. Additionally, they find utility where problem dimensionality is low [142].

Weak problems, on the other hand, make very little - if any - assumptions about the search space but at a cost to their effectiveness and/or performance and should only be used when there is no other option [142]. Each essentially samples the search space in the hope of discovering good solutions quickly [142]. The quintessential example of a weak method is that of exhaustive enumeration of the search space, a technique whose use is infeasible for all but the smallest of problem instances [142]. They tend to make use of a random component in generating solution instances, given the fact they also use no assumptions - repeatability then tends to be low as repeated runs of these approaches typically produce very different results [142].

Between these two classes lies an *intermediate* class encapsulating approaches such as metaheuristics and other stochastic algorithms. Metaheuristics will be discussed later in this chapter in section 2.4.5.3 and specific methods described in Chapter 4.

2.3.1 Strong Methods

2.3.1.1 Gradient-based Optimisation Methods

GRADIENT ASCENT The first gradient-based method that will be discussed is that of Gradient Ascent - or Descent if dealing with minimisation. Here, the basic idea is to find the slope of the function to be optimised, from the current point, and move up the hill towards the maxima at its peak [82]. Progress continues until the slope of the function reaches zero - which may or may not have occurred around a function maxima. Since progress cannot continue where the slope equals zero, it is possible for gradient ascent to converge elsewhere, such as: (i) at minima of the function and (ii) at *saddle points*. The gradient ascent method does not require the computation or knowledge of $f(x)$, but it does however make the assumption that the first derivative $f'(x)$ can be known and calculated [82]. The procedure for this method is quite straightforward for both the 1-dimensional and the n-dimensions cases.

In 1-dimension, a random initial point is selected. We then continually add to it a portion of its slope i.e., $x \leftarrow x + \alpha f'(x)$ - where α is a small positive value [82]. This process continues until the slope reaches zero and therefore x is unable to change further [82]. For the n-dimensional case, the slope at the current point is simply replaced with the gradient so that $\bar{x} \leftarrow \bar{x} + \alpha \nabla f(\bar{x})$. There are a couple of known problems with gradient descent, one has already been discussed - that of search progress being unable to continue when a slope not found at an optima equals zero. The other issue is speed of convergence. The reason behind this second issue is that as the slope approaches zero at the maximum, the value for x will overshoot the peak - landing

on the other side of the hill - and may do this many times before reaching the maximum. The cause of this is due to the step size α being based solely on the slope of the function at the current position. If a slope is very steep, α will be large even if it is not required. One way to deal with this is to tune the value of α to the particular problem instance being tackled. Another way, if the second derivative of the function, $f''(x)$, can also be calculated, would be to abandon the gradient ascent method altogether, opting instead for Newtons Method, a variation of gradient ascent described next.

NEWTON'S METHOD This variation on gradient ascent/descent involves the use of the second derivative, i.e. $x \leftarrow x - \alpha \frac{f'(x)}{f''(x)}$ of the function such that α is dampened as the algorithm approaches a slope of zero [82]. This means that the algorithm will now converge to any kind of zero slope - maxima, minima, saddle points and to points of inflection [82]. Additionally, the multidimensional case of the second derivative is not as straightforward to compute as the gradient $\nabla f(\bar{x})$ was for the first derivative, but is instead a multi-dimensional Hessian matrix $H_f(\bar{x})$ comprised of partial second derivatives over each dimension [82]. This extra complexity is compounded by the fact that newtons method divides by the second derivative, requiring that inverse hessian matrix be calculated [82].

Since the second derivative is being used, the algorithm maintains all the information required to identify whether it has reached a maximum solution (or minimum point in the case of minimisation) as opposed to other points of zero slope discussed previously. This is due to the fact that the second derivative will have negative value when a maximum point is met and positive otherwise [82].

The use of the functions second derivative still does not solve the overall problem of local search methods however, in that it can still converge to local optima rather than the global optima [82]. One simple method of constructing a global search algorithm from both gradient ascent and newtons method is to place these algorithms in a loop such that as a local optima is discovered, the algorithm may be restarted from a random position in the search space - keeping track of the best optimum found so far - in the hope of finding better optima which hopefully lead to the global optima [82]. This is essentially the same approach taken for derivative-free hill climbing algorithms when looking to search for global optima - see Chapter 4 Section 4.2.2.1 for more details on using restarts for global optimisation.

2.3.1.2 *Response Surface Methods*

Originating from the area of experimental design, response surface methods make the assumption that the response surface (objective function landscape / fitness landscape) can be parameterised under a simple functions containing a single optimum [142]; that is, these

methods aim to estimate the response surface by choosing parameters in an intelligent way [142]. A common implementation is known as a central composite design (CCD) which makes use of 2^N points (factorial points) [70] combined with $2N$ points (axial points) [70] and one central points [142, 70].

2.3.1.3 Simplex Methods

The Nelder-Mead Simplex algorithm, or simplex algorithm, was introduced in 1965 by Nelder and Mead in [94] for multidimensional unconstrained optimisation problems. The approach achieves this without derivative information - using only the returned objective function values at given points - and do not attempt to approximate a gradient at any point [122]. From its name, the Nelder-Mead Simplex algorithm is simplex based ¹ Simplex-based algorithms begin with an initial working simplex composed of $n + 1$ solution points considered as the vertices of the simplex and a corresponding set of objective function values $f_j := f(x_j) : j = \{1, 2, \dots, n\}$ [122]. It is a requirement of the method that the initial vertices do not lie in the same hyperplane, that is, the simplex should be non-degenerate [122].

To execute the algorithm, a sequence of transformation operators are applied to the simplex which are determined by testing one or more test points and comparing these with the objective values of the vertices [122]. The general procedure can be stated simply as [122]

1. Generate an initial working simplex
2. While termination criteria not satisfied:
 - a) calculate the information relevant to testing termination criteria
 - b) transform the current working simplex
3. Return the objective value of the best vertex of the final simplex

There are four fundamental transformation operators, each controlled by a parameter: reflection α , contraction β , expansion γ and shrink δ [122]. Each of the four parameters α , β , γ and δ must satisfy the following constraints [122]:

- $0 < \alpha < \gamma > 1, 0 < \beta < 1, 0 < \delta < 1$

The most common values of these parameters that satisfy the constraints are: $\alpha = 1$, $\beta = 0.5$, $\gamma = 2$ and $\delta = 0.5$ [122].

¹ A simplex S in a search space R^n is defined as the convex hull of $n + 1$ vertices in R^n [122]

TRANSFORMATION STEPS Transformation of the current simplex follows three steps [123, 122]:

1. determine the worst, second-worst and best vertices in the simplex: h , s and l respectively as in:

$$f_h = \max_j f_j, \quad f_s = \max_{j \neq h} f_j, \quad f_l = \min_{j \neq h} f_j$$

2. A new centroid is computed for the best side - opposite from the worst vertex - by:

$$c = \frac{1}{n} \sum_{\substack{j=0 \\ j \neq h}}^n x_j$$

3. Calculate the new working simplex by:

a) First attempt to replace the worst vertex with a new vertex through reflection...

- If the new vertex x_r is better than the second best vertex x_s but is not better than the best vertex x_l then simply accept the new vertex, replacing x_h and continue on to the next iteration
- If the new reflected vertex x_r is actually better than the best vertex x_l , attempt to explore further by expanding x_r to x_e by:

$$x_c = c + \gamma(x_r - c)$$
- The vertex x_h is replaced by the better of the vertices x_r and x_e

b) If the reflected point x_r was found to be worse than the second best vertex x_s then this suggests that the region of x_r is not a promising direction so we perform contraction on the simplex by: $x_c = c + \beta(x_h - c)$

- if the new contracted vertex x_c is better than x_h then x_h is replaced by x_c and move to the next iteration, that is: $x_h \leftarrow x_c$
- If the contracted vertex x_c is worse than x_h we must resort to the final transformation, the shrink operator which transforms the entire simplex by maintaining the current best point x_l and modifying all others relative to it. The j^{th} new point is calculated using:

$$x_j = x_l + \delta(x_j - x_l)$$

Termination is then determined by the following criteria:

1. Iteration budget has been exhausted
2. The minimum size for the simplex has been reached
3. An acceptable objective value for one of the vertices has been reached

The main advantage of simplex methods according to [122] is that it succeeded in achieving a good minimisation of the objective function over a number of numerical trials and does so over a relatively low number of objective function evaluations. However, a disadvantage given is that “numerical breakdown” of the algorithm occurs in practice resulting in requiring an “enormous” number of iteration and making very little progress toward the minimum - even when the simplex is nowhere near the minimum [122].

2.3.2 *Weak Methods*

2.3.2.1 *Random Search*

This very simple approach involves sampling the search space at random, maintaining the best found and terminating when the time budget is exceeded [142]. As such, this is not an approach that is widely used [142] given the multitude of other more successful methods available - of which some are close in terms of their simplicity (See Hill Climbing Approaches, Chapter 4 - Section 4.2). In fact, the only class of problem that this approach is expected to perform as well as any other is that of needle-in-the-haystack problems - where the entire search space is completely ‘flat’ except for a single point, and so no information involving the proximity of the optimum from the current location can be obtained by any existing algorithm [142].

2.3.2.2 *Sampling Methods*

A somewhat more feasible weak method is to sample the search space using a grid of a pre-determined size [142] - aptly referred to as *grid search* methods. When promising regions are identified, grids with smaller spacings may be placed in these locations in order to perform a more fine grained search [142]. These approaches suffer more than most from the curse of dimensionality (Chapter 3 - Section 3.2) where as the number of problem dimensions increases the number of samples required to form the grids increases exponentially [142]. Another disadvantage is that since points between the samples are not known so good solutions and regions smaller than the spacing of the grid may be missed [142].

2.4 METAHEURISTIC OPTIMISATION

2.4.1 *What is a Metaheuristic?*

The name metaheuristic is derived from a combination of the Greek prefix ‘meta’, meaning beyond - in terms of high-level and ‘heuristic’ from the greek word ‘heuriskein’ meaning ‘to search’ [39]. This combination can therefore translated roughly to ‘high-level search’.

Although this etymology may provide us with a loose idea of what it may mean to be a metaheuristic, unfortunately, no one definition of metaheuristic available in the literature seems to capture everything that it means to be a metaheuristic and so several will be presented in this section for comparison.

As stated by Burke and Kendall in [17], Glover first used the term metaheuristics in 1986 in [39] and defines it as²:

“A meta-heuristic refers to a master strategy that guides and modifies other heuristics to produce solutions beyond those that are normally generated in a quest for local optimality. The heuristics guided by such a meta-strategy may be high level procedures or may embody nothing more than a description of available moves for transforming one solution into another, together with an associated evaluation rule”

Stutzle in [131], found in [11], offers another quite explanatory definition:

“Metaheuristics are typically high-level strategies which guide an underlying, more problem specific heuristics, to increase their performance. The main goal is to avoid the disadvantages of iterative improvement and, in particular, multiple descent by allowing the local search to escape from local optima. This is achieved by either allowing worsening moves or generating new starting solutions for the local search in a more “intelligent” way than just providing random initial solutions. Many of the methods can be interpreted as introducing a bias such that high quality solutions are produced quickly. This bias can be of various forms and can be cast as descent bias (based on the objective function), memory bias (based on previously made decisions) or experience bias (based on prior performance). Many of the metaheuristic approaches rely on probabilistic decisions made during the

² there are other earlier papers making use of the term metaheuristic

search. But, the main difference to pure random search is that in these algorithms randomness is not used blindly but in an intelligent, biased form.”

A less formal definition of metaheuristics, as defined by Sean Luke in [82], is that metaheuristics comprise a general class of optimisation algorithms - forming the main sub-field of stochastic optimisation - that make use of some level of randomness in order to find the solution to ‘hard’ problems [82].

Blum and Roli in [11] provide a summary of the fundamental properties in which they feel characterise metaheuristics:

- “Metaheuristics are strategies that “guide” the search process”
- “The goal is to efficiently explore the search space in order to find (near-)optimum solutions”
- “Techniques which constitute metaheuristic algorithms range from simple local search procedures to complex learning processes”
- “Metaheuristic algorithms are approximate and usually non-deterministic”
- “They may incorporate mechanisms to avoid getting trapped in confined areas of the search space”
- “The basic concepts of metaheuristics permit an abstract level description”
- “Metaheuristics are not problem specific”
- “Metaheuristics may make use of domain-specific knowledge in the form of heuristics that are controlled by the upper level strategy”
- “Today’s more advanced metaheuristics use search experience (embodied in some form of memory) to guide the search”

They continue to then offer their own short definition of a metaheuristic, as [11]:

“In short we could say that metaheuristics are high level strategies for exploring search spaces by using different methods. Of great importance hereby is that a dynamic balance is given between diversification and intensification . . .”

For our purposes, we define a metaheuristic to be an algorithm making use of some form of randomness, comprised of one or more high-level heuristics whose task it is to drive and control a set of lower-level search heuristics in their quest to discover ‘good enough’ solutions to ‘hard’ optimisation problems.

2.4.2 *Exploration vs. Exploitation in Metaheuristics*

A common theme in all metaheuristics is that of the trade-off between exploration and exploitation [6] - also referred to as diversification and intensification respectively [11]. Exploration refers to the force in a metaheuristic, provided by its operators, that diversifies solutions maintained by the algorithm from the search space [6]. This divergent behaviour of the algorithm encourages a global search behaviour [6]. Exploitation on the other hand refers to how well the metaheuristic operators are able to use information from previously discovered solutions in order to focus the search on specific regions of the search space [6]. Metaheuristics exist on a spectrum depending on how much of each property is displayed i.e., some metaheuristic approaches can be more exploitative than explorative and vice-versa [6].

On the extreme ends of the scale, an example of a purely explorative search can be observed in the trivial random search algorithm which involves randomly selecting solutions in the search space for a given number of iterations [6]. Examples of purely exploitative search are hill climbing algorithms (Chapter 4 - Section 4.2) that perform small incremental steps in the space while new improvements are discovered or no further steps can be taken without finding a worsening solution [6]. Most metaheuristics, however, offer the ability to tune these aspects through the parameters used by their operators [6].

The exploitative and explorative nature of metaheuristics, along with their effective use of randomness, give them several advantages over classical (exact) optimisation methods. Bandaru and Deb illustrate this in [6] with the following list:

- Metaheuristics are able to find ‘good enough’ solutions to computationally ‘easy’ problems with a large input complexity [6]
- They can find ‘good enough’ solutions for NP-hard problems (see Section 2.2.3.1) [6]
- As opposed to most classical optimisation approaches, metaheuristics make use - and indeed, do not require - gradient information about the objective function and can therefore be applied to non-analytic, black-box or simulation-based objective functions [6]
- Most metaheuristic algorithms are capable of escaping from local optima [6]
- Due to the above ability to escape local optima, metaheuristics can better deal with uncertainties related to the objectives [6]

- Multiple objectives can be handled by most metaheuristics given minor algorithmic modifications [6]

2.4.3 *Problems Addressed by Metaheuristics*

The types of problem addressed by metaheuristics are those where

- The size of the solution space is large enough to make the use of exhaustive or exact algorithms infeasible,
- There is no known deterministic algorithm to solve the problem at all or within a reasonable time frame, and
- Problems where there is little knowledge of the domain that can be used to construct an optimal solution[82].
- A near-optimal solution to a problem is acceptable as opposed to the searching for the exact optimal solution [6]. Metaheuristics do not come with any guarantees of optimality.

2.4.4 *Categorisations of Metaheuristics*

Several different categorisations of metaheuristics, based on a variety of properties, have been made over the years. Unfortunately, there is little consensus as to a single scheme for categorising metaheuristics from the list given in the next subsections, they are often used interchangeably in the literature and even combined for increased specificity.

2.4.4.1 *Nature-inspired vs. Non-nature Inspired Metaheuristics*

One method used to classify metaheuristics is by basing this classification on the conceptual origins of the approach [11]. In metaheuristic research, these origins roughly fall into the category of either being inspired by natural processes (Genetic Algorithms (GA) and Particle Swarm Optimisation (PSO)) or “artificial” (Tabu Search and ILS) in the sense that nature was not queried for insight. Blum and Roli [11] do not consider this classification very meaningful for two reasons. The first is that many ‘hybrid’ approaches do not fit cleanly into either class - and in fact, may fit both simultaneously [11]. The second reason given is that it is often difficult to attribute even some non-hybrid metaheuristics to either class, for example, the use of memory in Tabu search could also be considered as being nature inspired [11].

2.4.4.2 *Single-solution vs. Population-based Metaheuristics*

Another characteristic that may be used for classifying metaheuristics is whether one or more solutions are being utilised at once [11]. Metaheuristic algorithms that operate only on a single solution are often referred to as *trajectory methods* which include the set of local search algorithms such as: Iterated Local Search (ILS), Simulated Annealing (SA) and Tabu Search [11]. What all these algorithms have in common is their property of describing the path, or trajectory, through the search space [11]. In contrast, population based metaheuristics such as: Particle Swarm Optimisation (PSO), Genetic Algorithms (GA) and Differential Evolution (DE), focus on the descriptions of the evolution of multiple, often related, solutions in the search space [11]. Another way of referring to these classifications is as *Local vs. Global Optimisation Metaheuristics* - typically single solution metaheuristics are only employed for local search whereas population-based approaches are able to search for a globally optimum solution to a given problem.

2.4.4.3 *Dynamic vs. Static Objective Function*

Metaheuristics can also be classified by the way in which they might make use of the objective function [11]. Some approaches use the objective function in the same way it is provided - as a 'black-box' - where others such as guided local search make modifications to the objective function as the search progresses [11]. The logic behind this approach is to enable the algorithm to escape from local optima by changing the search landscape itself [11].

2.4.4.4 *One vs. Many Neighbourhood Structures*

The majority of metaheuristics make use of only a single neighbourhood structure throughout the search process. What this means is that the topology of the search landscape does not change [11]. Others, the canonical example being Variable Neighbourhood Search (VNS), make use of a set of neighbourhood topologies allowing for the algorithm to switch between them in order to diversify the search [11].

2.4.4.5 *Memory vs. Memoryless Metaheuristics*

Another classification scheme is to focus on whether or not the metaheuristic makes use of search history / trajectory information - these days considered to be a fundamental component of successful metaheuristics [11]. Memoryless algorithms, such as Simulated Annealing (SA), carry out a Markov process - where the information used to decide where the search should explore next is taken exclusively from the current solution [11]. For Memory algorithms, a further delineation can be found in the various ways that an algorithm may make use of

memory - usually considered between the categories of short or long-term memory [11]. Short-term memory approaches typically maintain only information about the previous few iterations, the last few solutions found and decisions made [11]. Long-term memory approaches on the other hand make use of a large set of various 'synthetic' parameters about the search [11].

2.4.5 *General Concepts*

2.4.5.1 *Solution Representations Used in Metaheuristics*

In order to develop solutions to a given problem; that is, to provide input to the objective function representing the problem, a suitable solution representation for use in optimising the problem in the context of metaheuristics must be selected. There are several drawbacks and advantages to the use of each type of solution representation discussed here, however the selection of a representation is often heavily influenced by the problem being solved. Any change to solution representation will often require a subsequent change to the objective function and in effect, the nature of the search space will also change. For example, it may be tempting to use a representation simply because it feels more natural and intuitive but some search spaces based on a representation (and its corresponding objective function) may well be more rugged, sparse, deceptive or of a larger scale than other spaces defined by other representations and functions. The solution representation also determines how the search operators might create neighbourhoods and how the objective function will be evaluated; for example, a representation requiring heavy or complex conversion by the objective function will require more time and computational resources. Objective functions are considered to be one of the main bottlenecks in search, so it is always prudent to minimise computation costs here where possible. All in all, it is important to keep the nature of the resulting search space in mind when selecting a representation as the efficiency of search is very much determined by it. In the following paragraphs, since we work primarily with a vector-based representation³, we outline several of the more common vector-based representations and how they are typically used.

(A) **BINARY REPRESENTATION** The simplest of the common representations used, the *binary representation*, is one in which each solution is held as a vector of bits where each bit or sub-sequence of bits represents some parameter of the problem being solved [29]. For a given problem, a practitioner should: (i) ensure that it is clear how the binary string is to be

³ By vector, we mean a one-dimensional array of a fixed length, as per Luke in [82]

interpreted, that is, how it is to be decoded into a meaningful solution and (ii) in ensuring (i), that each possible binary vector encodes a valid (but not necessarily feasible) solution to the problem given and conversely that each possible solution can be encoded [29]. For some problem types, the binary representation is a natural one - especially those considering boolean decisions (yes/no) [29, 111, 135] - thus ensuring the above considerations are met. An canonical example is a *0/1 Knapsack Problem* where a solution can be encoded as simply as whether or not an item is included in the knapsack or not - a binary decision [29, 111, 135]. It is also common that a binary vector be used to represent other data types such as integer or real-valued numbers; for example, a bit string of length $l = 80$ can be used to encode ten integer values using 8 bits each or five 16 bit real-valued numbers [29] - a scheme known as Binary Coded Decimal (BCD) [143]. BCD represents each decimal digit as a string of bits (4 or more bits). So a string of length 80 could code 20 decimal digits (each length 4), or 10 decimal digits (each length 8), or 5 decimal digits (each length 16). The sign is usually coded as another single bit. However, these schemes tend to have the disadvantage that consecutive decimal numbers are considered neighbours in the context of decimal space but not in the context of binary space [143], that is, the Hamming distance between two consecutive decimal numbers, for example integers, often does not equal 1 [29] forming what are known as Hamming Cliffs: where adjacent integers are represented by complimentary bit vectors and therefore share few or no bits in common [143, 82]. This problem can be helped by a binary to decimal mapping variation, *Gray Coding*, where consecutive integers will always have a Hamming Distance of 1 [29, 143].

(B) INTEGER REPRESENTATION Used in situations where individual solution parameters can naturally take one of n discrete values, the Integer representation has been used for solutions representations, such as: of unordered sets i.e., $\{1 = \text{red}, 2 = \text{green}, 3 = \text{blue}\}$ or the definition of a discrete metric space [82]. When designing an integer-based solution representation of a problem, it is worth considering if there are any natural relationships between all the possible values that a solution parameter can have [29]. This is often more obvious for ordinal parameters, such as for the integers themselves, than for cardinal parameters - such as compass directions that do not appear to have a natural ordering [29].

(C) REAL-VALUED REPRESENTATION Usually, the most natural encoding for parameters that take their values from a continuous rather than discrete space is a real-valued or floating-point representation [29]. Given that precision on current digital computer technology is limited, it is technically more accurate to refer to this representation as a floating-point representation [29]; however, for consistency we will continue with using the term *real-valued*

for the remainder of this thesis. An example of a natural class of problem for a real-valued representation is that of continuous function optimisation, where the parameters of a solution represent a continuous point within the metric space of a given continuous function. When recombination operators for real-valued representations are considered, it is possible to generalise many of the operators used for discrete spaces to real-valued representations [29], for example: one-point, n-point crossover and uniform crossover (Chapter 4, Section 4.3.2.2); however these *discrete recombination operators* - as they are for discrete representations - are unable to introduce new values for the parameters thus leaving this task to the mutation operators [29]. Other options are to use recombination operators specifically suited to real-valued solutions, ones that are able to generate new parameter values [29]. Operators of this type that generate new values between two solutions to be recombined are known as *intermediate* or *arithmetic* recombination operators [29]. Alternatively operators that generate new parameter values close to - and can exceed - one of the solutions are known as *blend* recombination operators [29]. Unlike recombination, mutation operators for real-valued solutions cannot be extended naturally from discrete representations and so specialised mutation operators have been developed; one such example is that of *uniform* mutation (Chapter 4, Section 4.3.2.3).

2.4.5.2 *Metaheuristics and the 'No Free Lunch Theorem'*

In the fields of optimisation and metaheuristics, the no free lunch theorem (NFL) states that [146, 145]:

When averaged across all possible problems, the performance of all algorithms is equal, irrespective of the choice of evaluation criteria.

There are a couple of consequences to this statement. Firstly, since all algorithms are equally good on average across all problems, this suggests there must be some subset of problems on which an algorithm performs well and another subset on which it performs poorly⁴. Secondly, it suggests that it is not correct to state that one algorithm is better than any other and that comparisons between algorithms must take into account the types of problems the algorithm performs best on.

There have been attempts to develop techniques that are able to circumvent the no free lunch theorem to various degree. Two notable examples are that of hybrid metaheuristics and hyper-heuristics, both of which are beyond the scope of this thesis.

⁴ of course, this is an over-simplification made for the purposes of illustration

One caveat of the NFL is that, at the moment, there is some contention on whether or not the NFL is only valid for discrete problem domains. For most problems being solved computationally on a Von Neumann architecture this detail should not present a problem due to the fact that we make use of a finite representation with precision bounded by internal representation. For more information and discussion on both sides of the argument about this detail of NFL, readers can refer to both [4, 141].

2.4.5.3 *Comparisons Between Metaheuristic Approaches*

Comparison between metaheuristic techniques has been said to be more difficult than for other types of algorithm [121]. However, comparing metaheuristics to one another is an important step in metaheuristic development in order to determine where a new approach stands in relation to others in terms of various performance measures a suitability to certain situations. Since the NFL theorem prohibits any implication that one approach is in general ‘better’ than any other, one cannot make such comparisons based solely on solution quality over a set of problem instances. As such, there are several important considerations that should be taken into account in order to effectively compare metaheuristics to one another in a meaningful way - without violating the implications of the NFL. A discussion of these considerations appear in the next few sub-sections based exclusively on the excellent presentation of the topic provided by Silberholz and Golden in [121].

(A) TESTBEDS

(A.i) The Use of Existing Testbeds In the majority of cases, it is of benefit - and actually considered to be the ideal case [121] - to make use of pre-existing testbeds available from the metaheuristic literature [121]. Not only does this allow comparisons with other metaheuristics to be carried out on a “by-instance” basis [121], but where algorithm results are published, it also removes the need to either obtain or reproduce algorithm source code from the literature in order to compare against the same or a different testbed - of which it will likely provide a different set of results. However, these existing testbeds can sometimes be insufficient or their use is not possible [121] - examples include:

- The existing testbeds have not been made available to other researchers or the problem instances have been generated randomly [121], for which either or both the function used to generate the instances or the seeds used for the pseudo-random number generator are not published or available.

- The testbeds are too small either in terms of available problem instances or of problem instance size [121].
- A new problem is being addressed by the metaheuristics being compared, therefore, there will be no existing testbed available for use [121]

When the available testbeds are insufficient for use in comparing metaheuristics available, there may be little option other than to develop new testbeds.

(A.ii) Creating a New Testbed Although making use of existing testbeds is the preferred option when comparing metaheuristics, there are several situations where there is little option except for developing a new testbed - some of which have already been discussed in the previous sub-section. This section will discuss several factors that should be taken into consideration when developing a new testbed. The first goal in developing a new testbed is that of ensuring problem instances included in the suite effectively mimic real-world situations [121]. Although artificial problems have their uses when developing and testing a metaheuristic, being artificial, many if not most of these fail to fully represent the important characteristics of a corresponding real-life problem instance. Since metaheuristics find their utility being applied to real-world problems as opposed to 'toy' problems it is crucial during the analysis of metaheuristics that the behaviour of algorithms on these kinds of problem is represented in results.

Another goal in creating an effective testbed is the provision of problem instances of various types and levels of difficulty [121]. In terms of types of problem, every metaheuristic has a 'niche' set of problems on which it performs very well and others where it does not perform as well. This observation is a consequence of the 'No Free Lunch Theorem' (Section 2.4.5.2) which states that when averaged over all possible problems, every search method, including metaheuristics, show the same performance - this means that there has to be a subset of all problems where a given metaheuristic performs best. It is therefore important when analysing and comparing metaheuristics that a wide range of problem types are represented in the testbed in order to identify those types of problem where a metaheuristic performs best. In terms of the levels of difficulty found among the test instances, several factors are involved in determining what makes a problem difficult to solve, which are discussed in Chapter 3, Section 3.2. When applied to real-world settings, metaheuristics will likely be targeted towards optimising objective functions which exhibit some of these factors - when analysing metaheuristics, representing all of these factors in the testbed ensures a more accurate depiction of overall performance.

Silberholz and Golden also suggest that large problem instances must be represented in the testbed [121]. One reason for this is that exact methods are often able to run in a reasonable

length of time whilst also providing the guarantee of returning the optimal solution [121]. Metaheuristics are not required for these types of problem, they are most useful for those problems where the optimal solution cannot be found so readily. As metaheuristics show their utility being applied to these more intractable problems, it is important to ensure that they are represented in any new testbed.

(A.iii) Making the New Testbeds Available After creation of a new testbed, it is important to make sure it is easily accessible. Providing a testbed for others not only ensures it is widely used but also allows other researchers to make their own comparisons easily [121]. According to Silberholz and Golden in [121], one easy way to make a testbed available is to provide a simple generating function for the problem instances in the testbed. Another approach is simply to publish the testbed [121] or make it available through an author provided webpage or web service. Of these last suggestions, publishing will likely be the preferred option as it relies less on other researchers simply 'stumbling' upon the webpage and has the benefits of having the testbed go through the peer review process.

(B) CLASSIFICATION OF THE PROBLEM INSTANCES Proper classification of the problem instances present in the testbed is crucial to the effective analysis of metaheuristic techniques [121]. As discussed previously, since the no free lunch theorem suggests that every search algorithm has its niche set of problems on which it performs well, this niche set can only be discovered and subsequently reported for a new metaheuristic when problems in the testbed are appropriately separated into classes. According to Silberholz and Golden, state that each of these classifications should be:

- Noted prior to any experimentation
- Discussed individually in terms of a metaheuristics performance

The authors note further that the proper classification of real-world instances, and subsequent analyses, is of particular importance as it can help algorithm developers in industry, with a certain type of dataset, to decide on the most suitable approach to use [121].

(C) PARAMETERS In terms of comparing the actual algorithms themselves, it is useful to consider the factors surrounding the parameters included in these algorithms. Here, these are defined as the variables belonging to an algorithm that are involved in modifying the behaviour of the algorithm. For example, common parameters in a GA can include: population size, mutation rate and tournament size (for tournament selection). Silberholz and Golden note that in a metaheuristic, parameters can be set statically, such as using a fixed population

size in a genetic algorithm, or based on the problem instance - a population size of $5\sqrt{n}$ for instance, where n is the number of nodes in the instance [121].⁵ Four factors which will be discussed here are: the relation of the number of parameters to the simplicity of an algorithm, tuning and visualising the parameter space, parameter interactions and fair testing of metaheuristics in the context of parameters.

(C.i) Parameters and Algorithmic Simplicity When making comparisons between different metaheuristics in terms of solution qualities, it is useful to consider the complexity of the approach. For example, given two approaches which produce similar results, the simpler of two approaches is the superior algorithm [121]. A few reasons behind this superiority include: (i) being simple to reimplement in an industrial setting, (ii) being simpler to reimplement by researchers and (iii) being simpler to explain and analyse [121]. A number of reasonable metrics exist including the number of lines of source code required to implement the algorithm or the number of steps of pseudocode used to describe it [121]. However, both of these approaches rely too heavily on the programming language used to implement an algorithm, the individual style of the programmer and the level of detail in terms of pseudocode descriptions [121]. Silberholz and Golden prescribe a more meaningful metric - that of the number of parameters used in an approach [121]. Although this metric does not provide an estimation of programming complexity in terms of actually implementing an algorithm, it does however give a good estimation of how difficult the algorithms behaviour will be to understand and tune according to a given application.

(C.ii) Visualising and Tuning the Parameter Space The consequence of the parameter metric prescribed by Silberholz and Golden in the previous section is that metaheuristics making use of many parameters are considered more complex than those with only a few [121]. The main source of this complexity is that the effort required to understand and tune the parameters becomes far greater as the number of parameters considered increases [121] - the size of the parameter space in fact increases exponentially as more parameters are added. For example, if we consider a brute-force approach to tuning the parameters in a metaheuristic, tuning involves testing the performance of a set of m parameter values for each of the n parameters under consideration [121]. If only three values for each of the parameters were tested i.e. 3^n different configurations, a configuration set would be obtained of size: 9 for a metaheuristic with only 2 parameters, 2187 configurations for a metaheuristic with 7 parameters and an intractable 1,853,020,188,851,841 configurations if 32 parameters are used [121]. It is true, however, that there are more intelligent methods used to search for good

⁵ Although, other algorithms in the literature can do adjust parameters dynamically based on the progress of the search

parameter settings - automatic parameter tuners such as SMAC (see Chapter 5 - Section 5.4) and Irace are two such examples - but the exponential growth of the parameter space still poses problems to these methods [121]. In addition to these issues, is that of visualisation of the parameter space as the amount of parameters in an algorithm increases [121]. Summarising an example from Silberholz and Golden [121], for an algorithm with only two - or even three - parameters, the parameter space can be quite easily visualised with many standard 2D visualisation techniques, e.g., a scatter plot where the axes each represent one of the parameters and the colour or size of the points can represent the cost of the parameter setting at a given point. In the same way, a space of three parameters may be viewed on a 3D scatter plot, however there will likely be some difficulty in viewing interior points unless various 'slices' through the visualisation are taken. Although Silberholz and Golden make note of a potential visualisation for 4D parameters spaces; such as in [56], they are quick to point out that these types of approaches in no way makes 4-dimensions as intuitively visualised as the 2D and 3D cases above increasing the visualisation difficulty [121].

(C.iii) Interactions Between Parameters Another complexity caused by large parameter sets in a metaheuristic is that these have a tendency to produce complex parameter interactions [121]. Silberholz and Golden state that from the point of view of optimisation, the interaction of parameters implies that the optimisation of individual parameters or even small subsets of parameters will become increasingly ineffective as the number of parameters considered is increased [121]. The authors however cite an instance where non-trivial parameter interactions were observed in a genetic algorithm containing only 3 parameters. The implication of this is that to an extent it is often difficult to avoid all parameter interactions with the exception of metaheuristics making use of only 1 parameter or no parameters at all [121] - where no interaction can occur.

(C.iv) Fair Testing Concerning Parameters The other concern when comparing metaheuristics, apart from simplicity, is that of fairness when tuning parameters [121]. For example, if one algorithm is over-tuned to the entire set of problem instances of which it is being tested, that can produce unfair comparisons; alternatively, if only a representative subset of the problems are used to tune parameters and the remainder to compare the metaheuristics, this would represent a much fairer comparison strategy [121].

(D) COMPARING SOLUTION QUALITY Since metaheuristics are designed to produce solutions of high quality, comparisons between metaheuristics in the context of solution quality performance is amongst the most important [121].

(D.i) Quality Metrics In order to compare the quality of solutions produced by two or more metaheuristics, a suitable objective metric is required - the best metric being the deviation of the produced solutions from optimal [121]. Comparisons made over separate problem instances are, in general, weaker than those made over the same problem; this is because different instances will have different structures in their search landscapes and more than likely different values for the optimum [121]. This metric is clearly only applicable to those problems of which the optimum value is known. For problems where the optimum is not known, several alternative metrics exist. One popular metric is deviation from best known solution [121], where the best known solution published in literature can be used in place of the true optimum for determining performance. It is also possible to use this metric when no best known solution has been published [121] - possibly through the generation of a solution using an algorithm known to perform well on the particular class of problem being addressed - however, comparisons made without the true optimum will be less meaningful [121].

(D.ii) Fairness of Performance Comparison As mentioned briefly in Section 2.4.5.3, it is usually prudent to make use of existing test beds when comparing metaheuristics. In particular, the fairness of such comparisons against original implementations can be ensured, as no re-implementation of algorithms needs to be carried out (where source code is unavailable and results published) and the inadvertent introduction of implementation errors in the test bed itself is unlikely compared to creating a new test bed e.g., when using re-implementations of some problem instances from other test sets. In the context of this thesis, testbeds which can facilitate the fair comparison between metaheuristics used to solve continuous problems include those provided for black-box optimisation competitions such as the Black-Box Optimisation Benchmarking (BBOB) competition [32, 47], where the test set used is available as the *Comparing Continuous Optimisers (COCO)* platform described in [49] as well as the various benchmark sets provided for the *IEEE Congress on Evolutionary Computations (CEC)* competitions such as CEC-2005 [133].

The usual approach to fair performance comparison in the literature is using rank-based analysis. This is most commonly carried out by:

1. Applying the set of algorithms being compared to a suite of problems to obtain a 'raw' performance measure (whether a fitness or time-based measure)
2. Using the raw performance data to generate a rank ordering (best to worst performing) of the algorithms against each of the problems in the set
3. Aggregating the ranks generated for all problems for each algorithm to arrive at a *consensus* ranking against the set as a whole, as per Mersmann et al. [88]

Typically, these ranks will then be presented in tabular form in order to analyse the rankings, however in [88], Mersmann et al. make use of parallel coordinate plots to improve comprehensibility and ease of analysis. In this thesis, both approaches are used for algorithm comparison, however for our plots the axes are inverted to provide easier visualisation of results.

In terms of fairness of comparison, by normalising the raw performances as ranks any bias in comparison from the nature of the problems, e.g., different perceived difficulty levels, are eliminated as only the ordering of the algorithms in the context of the selected performance measure is considered. Subsequent statistical tests will therefore be far less influenced by an unbalanced weighting given to any one data point. However, since ranking certainly involves loss of information, they are far less sensitive to large variances in the data, in this case the variances between performances obtained by algorithms, and thus do not readily show large differences in performance which may be an important factor in a particular study. One possible benefit to this resistance to noise in the data is that it more clearly highlights the differences of real interest, cutting through the stochastic noise inherent to the data obtained from metaheuristics.

2.5 LARGE-SCALE GLOBAL OPTIMISATION

A Large Scale Optimisation Problem (LSOP) is characterised as having many hundreds or thousands of dimensions. These kinds of problem are considered harder to solve than lower dimensional problems and algorithms that show good performance in lower dimensions can perform poorly when used against a scaled version of the same problem [18]. The nature of problems of this scale differ hugely from lower dimensional problems as each increase of dimensionality causes an exponential increase in domain volume [18] - often referred to as the 'curse of dimensionality' (See 'Dimensionality' - Chapter 3, Section 3.2). For example, a 1-dimensional search space containing 100 possible solutions, where one solution can be considered as being the global minimum, will require 100 function evaluations in the worst-case. Scaling the same problem to 2-dimensions has $100^2 = 10000$ possible solutions but extending this even to a 50-dimensional problem would place the global minimum as a single point in $100^{50} = 1 \times 10^{100}$ (a googol) feasible solutions [18]. Further to the scaling of the search domain, dimensionality can also affect structural features of the search space itself. An example from Caraffini et al. in [18] states the fact that a unit sphere in 3-dimensions has a surface area of $S_2 = 4\pi$ and a volume of $V_3 = \frac{4}{3}\pi$. If we were then to extend this unit sphere into n-dimensional space, it can be proved that the ratio between surface area and volume

is $\frac{1}{n}$; meaning that as the dimensionality increases, more and more points can be found on the surface as the volume becomes comparatively smaller [18]. This greatly increases the size of the neighbourhoods around the optimal solutions, in turn greatly increasing the time for conventional metaheuristics to solve the problem.

Computationally, the effective solution of LSOPs can increase the cost significantly; due not only to the large increase in the search space, but also by the large function evaluation budgets required to find high performing solutions [18]. For example, an optimisation algorithm that is capable of solving a 10-dimensional problem within 50,000 function evaluations would require $50,000^{40} = 9.1 \times 10^{187}$ [18] which far outweighs not only the number of atoms estimated in the observable universe (10^{78} to 10^{82}) but also the number of seconds since the universe began $4.361170769 \times 10^{17}$. Clearly, maintaining search coverage in this way simply isn't tractable.

An often forgotten consideration is that of the parameter configurations of the algorithms being used against LSOPs. For example, Caraffini et al. [18] describe a common population scaling strategy whereby the population size should increase exponentially with dimensionality in order to maintain search coverage; however, they quickly point out that, as with function evaluations budgets, this strategy quickly becomes impractical [18]. Using the same example as before, if a population-based optimisation algorithm with a population size of 30 was successful against the 10-dimensional problem, to maintain a consistent solution density in 50-dimensional space, the population size would have to contain $30^{40} \approx 1.22 \times 10^{59}$ individuals [18]. Metaheuristic approaches are therefore only able to cover a small portion of large-scale search spaces since population sizes of this magnitude are impractical, and often impossible, to actually use [18].

From a similar standpoint, consideration can also be given towards the computational budget afforded to the optimisation algorithm i.e., commonly the number of objective function evaluations. Using a similar example to the above, an algorithm which is originally tasked with optimising a 10-dimensional problem using a budget of 50,000 function evaluations, would require - for the same problem in only 50-dimensions - $50000^{40} \approx 9.1 \times 10^{187}$ function evaluations in order to cover the same proportion of the search landscape [18]. Therefore, although it is usual practice when optimising problems with relatively few dimensions (certainly < 50) that the budget be scaled proportionally to the dimensionality, this strategy simply won't scale to LSOPs. However, real-world LSOPs exist and as with smaller real-world problems, often need to be optimised effectively and efficiently in order to produce useful solutions, so research has been focussed on the development of algorithms and approaches

aimed specifically at solving LSOPs whilst at the same time trying to overcome the various difficulties described above [18]. Caraffini et al. roughly categorise the methods for tackling LSOPs as follows [18]:

- Methods which exploit promising search trajectories intensively
 - Perhaps quite unintuitively, these methods effectively abandon the search for global optimality and instead intensively exploit promising regions already explored in order to find high performing solutions. Perhaps just as unintuitively, one popular method implementing this strategy involves the use of only a small population or by using a population which diminishes during the search [18]. Doing this allows one to control the amount of exploration carried out as there is less diversification present in the population, this in turn allows the population to begin converging towards one of the promising regions discovered early on in the search - with minimum exploration beyond this region. One other way to implement this approach is to combine highly exploitative local search algorithms with other algorithms and population-based structures [18]. Specifically, the local search algorithms used tend to exploit promising regions by perturbing a solution in an axis-parallel manner i.e., only one parameter of the solution is modified at any given time. Indeed, we found that the neighbourhood function utilised in the local search algorithms implemented for this study and making use of this axis-parallel strategy (see Chapter 4, Section 4.1.3) produced far better performance (in terms of solution quality) than allowing steps in any and all axes in a single iteration. Further, and echoed by Caraffini et al. [18], the modification of the Co-variance Matrix Adaption Evolution Strategy (CMA-ES) for separable problems *sep-CMA-ES* (detailed in Chapter 4, Section 4.7) implements this kind of axis-parallel trajectory through use of a diagonal rather than a full co-variance matrix in its updates. This variation of CMA-ES has been shown to be particularly promising within high-dimensional spaces, not only due to its use of an axis-parallel trajectory, but also, it is more suitable to large-scale problems as the calculation of a diagonal co-variance matrix requires far fewer computational resources and thus can be calculated far more quickly than a full matrix.
- Methods which use some mechanism for decomposing the search space
 - An implementation of this strategy is realised in an approach known as Cooperative Co-evolution, where LSOPs are decomposed into several smaller sub-problems which can each be solved independently of each other by a set of co-evolving

populations (one for each sub-problem). A more detailed discussion of Cooperative Co-evolution is presented in the next section (Section 2.5.1).

A common link between these two strategies for tackling large-scale problems is that each, in its own way, is achieving improvement by the exploitation of search directions [18]. Caraffini et al. go further, with the assertion that this means that “every modern metaheuristic for LSOPs gives up the search for the global optimum and simply tries to enhance as much as possible upon an initial sampling” [18]. However, it could be argued that an approach such as cooperative co-evolution does not in fact abandon the search for optimality, but instead is actively seeking it through the combination of simpler (and hopefully optimal) sub-solutions. Further, Ma et al. [83] point out that having multiple co-evolving populations helps the search as a whole maintain diversity; a property widely held to facilitate an effective global search. Also, one could state this same assertion about any metaheuristic algorithm, including those not specifically targeted at LSOPs, where metaheuristics as a general concept each work to enhance the performance of one or more initial states. Such algorithms, with the exception of simple local search algorithms, although not guaranteed to discover the optimal solution, do appear to have this goal in mind; which is particularly true of population-based approaches where there is a general aim to explore as much of the space as possible over the course of the search.

I do however agree with the authors of [83] that this is certainly true of approaches categorised by the first strategy; local search algorithms - as their name suggests - do not have a strong focus towards global optimality and small populations will reduce the amount of exploration of the space, which can more easily be interpreted as abandoning the idea of achieving global optimality than for approaches such as cooperative co-evolution.

2.5.1 Cooperative Co-Evolutionary Algorithm

Similarly to other evolutionary approaches, a Cooperative Co-Evolution Algorithm (CCEA) finds its inspiration from theories of natural evolution. The predominant evolutionary concept adopted by CCEA is that of *mutualism*, a facet of co-evolution, where the success of two or more species or sub-populations is partially dependent on shared mutually beneficial relationships [83]; in essence, the species can be said to be evolving together (co-evolution) with the mutual goal of continued survival. Three interactions between species have been observed to occur in nature [83]: (i) Resource-Resource Interactions - where each population/species trade

consumable products necessary for each of their survival, (ii) Resource-Service Interactions - where a consumable product is provided in turn for conducting a useful function e.g., bees consume pollen in trade for dispersing pollen to other places, and (iii) Service-Service Interactions - when each species provides a function useful to the other, e.g., a pro-biotic bacterial species provides services to the immune system of humans (as well as having other functions) and the human provides a safe, warm and nutrient-rich environment for the bacteria to live and procreate.

Looking closely at the history of our own species, it is easy to see that humans have themselves benefit from mutualistic relationships with other species as well as between other human groups. Consider our relationships with domesticated animals, particularly dogs, whereby humans benefit from positive emotional attachment as well as work provided by dogs e.g., herding of cattle and sheep, hunting, pest control, protection etc.; and humans reciprocate by providing a home, protection, food, warmth etc. Amongst ourselves, we have, and do, benefit from trade and cooperation between different countries and cultures as well as many other beneficial interactions. Henrik Valeur, a Danish architect who introduced the concept of co-evolution to the architectural world considers that; “As we become more and more interconnected and interdependent, human development is no longer a matter of the evolution of individual groups of people but rather a matter of the co-evolution of all people”[137].

CCEA, then, works by applying a simplified version of this natural concept in terms of co-evolving sub-populations of solutions to a given problem. Particularly useful against large-scale problems, each independent sub-population evolves towards the solution to a non-overlapping sub-problem of lower dimensionality. An overall solution can then be obtained through the combination of sub-solutions found by the separately evolving sub-populations [18, 83]. As such, the evaluation of any given individual in any given sub-population requires cooperation with all other sub-populations [83]. Specifically, the fitness of an individual from a sub-population is calculated in terms of complete solutions in which the individual is participating [83]. It is chiefly this divide-and-conquer decomposition approach that places CCEA at an advantage over more traditional evolutionary algorithms [83].

According to Ma et al. [83], CCEA can be said to have four main advantages stemming from its decomposition strategy:

1. Decomposition of the problem into sub-problems allows CCEA to be readily parallelised in order to speed up the optimisation process

2. Since each sub-problem is tackled by a separate sub-population, CCEA can maintain good solution diversity - a difficulty often plaguing other evolutionary algorithms
3. By decomposing a system into “sub-modules”, the system as a whole can be made more robust against the errors and failings of any one module; increasing reusability in dynamic environments
4. If decomposition is carried out effectively, the ‘curse of dimensionality’ (Chapter 3, Section 3.2) can be somewhat alleviated allowing for generally better performance when solving LSGO problems when compared to more traditional EA

Despite the power that can be displayed by CCEA, a difficulty often faced when using this approach is deciding how to go about decomposing a problem into its various sub-problems. Clearly, this problem is trivial if the problem is known to be fully additively or multiplicatively separable (see Chapter 3, Section 3.2) where there is minimal interaction between variables; problems can then be naively separated out as $k = \frac{d}{n}$ instances of the problem where d is the problem dimensionality and n is the number of dimensions per sub-problem.⁶ It is also possible, mostly for separable problems, to have $k = 1$, where each of the d sub-populations evolve solutions for sub-problems of a single dimension. In fact, this was the strategy employed in the first CCEA, Cooperative Co-evolutionary Genetic Algorithm (CCGA), by Potter and De Jong [107]. Complete solutions to separable problems can then be combined through the summation, multiplication or some other reconstitution of sub-solutions. However, a significant disadvantage of CCEA is that as sub-populations are added, the total number of evaluations of the objective function available to CCEA also needs to increase proportionally, since each population conducts its own set of function evaluations in determining the fitness of solutions to its sub-problem. Therefore, depending on the scale of the problem to be solved, it may be wise to consider larger sub-problems in order to reduce computation time.

Several ‘flavours’ of CCEA can be found in the literature, which include: the first cooperative co-evolutionary approach - CCGA - proposed by Potter and De Jong in [107], Cooperative Co-evolution with Differential Evolution (DECC) by Yang et al. [150] and Cooperative Co-evolution based on PSO (CCPSO) by Li and Yao [75] to name but a few. However, for the remainder of this section we will not focus on these variations on an individual basis, but instead form our discussions around one or another of the two basic CCEA algorithmic

⁶ Of course, this is a simplified view of problem decomposition for CCEA in general, sub-problems do not have to be of the same dimensionality as one another and can thus be of variable size. Indeed, for problems that are only partially-separable, one decomposition method may be to evolve smaller separable parts of the problem independently from a comparatively larger non-separable part(s)

frameworks outlined by Popovici et al. in [106] as many CCEAs are a variation on one of these frameworks or the other. Before discussing these however, I present a brief discussion of problem decomposition for CCEA.

For a more detailed discourse of cooperative co-evolutionary algorithms, three of the main sources cited in this section can be recommended: The survey paper by Ma et al. [83], focussing specifically of function optimisation [107] by Potter and De Jong and in particular [106] by Popovici et al. which provides a relatively thorough treatment of the subject matter.

2.5.1.1 *Problem Decomposition*

Many decomposition strategies have been developed since the inception of CCEA, including the simpler of those as mentioned previously, however, an overarching issue that arises when performing problem decomposition using any such strategy, is how to decompose problems where variable interactions are present [83]. Indeed, it is exactly this issue that motivates the search for more effective decomposition strategies - since without such interactions i.e., if every problem was fully separable, then a simple strategy using a sub-problem size of $k = 1$ would not have a more effective substitute.

In the ideal case then, the decomposition should be conducted with a focus on any possible variable interactions with a view to minimise such interactions; for example, a problem that can arise if interacting variables do not form part of the same sub-problem, is that the CCEA can easily become trapped in a pseudo-minimum [83]⁷. By keeping all interacting components together, these sub-problems can be solved separately without negatively affecting the solution to the other (separable) sub-problems.

Decomposition methods range from those which make use of Static Variable Grouping or Random Variable Grouping, to somewhat more complex methods such as linkage learning-based variable grouping, domain knowledge-based variable grouping and overlap and hierarchical grouping [83]. However for brevity, we will only discuss the first two of these here.

(A) **STATIC VARIABLE GROUPING** Static variable grouping methods make no attempt to discern any interdependence amongst variables, but instead simply decompose high-

⁷ Ma et al. note that a pseudo-minimum does not refer to a local minimum of the original (full) problem, but is in fact a minimum that is created by incorrect decomposition. A difficulty that likely plagued Potter and De Jongs early CCGA which indeed performed much better against separable problems than against non-separable problems

dimensional problems into a series of low-dimensional sub-problems where the grouping remains fixed throughout the execution of the algorithm [83]. The original CCEA by Potter and De Jong (CCGA) [107] is an example of static variable grouping, where the grouping is held as n 1-dimensional sub-problems and each searched by a separate sub-population [83]. As mentioned previously, this algorithm did not perform well when faced with non-separable problems. To overcome this problem, several other static grouping methods were proposed. First, in [124], a pool of static decompositions, where each could potentially vary in size from a single variable sub-problem to a full length problem, were used to produce new solutions through random combination. Another decomposition method which statically decomposed a problem into m s -dimensional sub-problems, where $n = m \times s$, was introduced by van den Bergh and Engelbrecht in [138] which was shown to perform better than Potter and De Jongs 1-dimensional decomposition on several benchmark functions [83]. A later method, very similar to this, by Cao et al. [152], made use of a sequential sliding window to decompose problems and was compared to several random variable grouping methods, however, it was not compared against the method of van den Bergh and Engelbrecht or any other static variable grouping method for that matter.

(B) **RANDOM VARIABLE GROUPING** Random Variable Grouping (RG) is often employed to partially alleviate the problem faced by static variable grouping, whereby since groups are fixed throughout the optimisation process, if two interacting problem variables are placed into different groups on initialisation of the algorithm, there can be no way of them to be part of the same group at all during execution [83]. It does this by randomly selecting variable for each sub-problem [83], that is, the grouping of variables is changed throughout the course of optimisation. Two schemes for selecting the size of the sub-problems in RG are: (i) using a *fixed* size or (ii) *dynamically* changing the size throughout the optimisation process [83]; each of these will be discussed separately.

(A.i) Random Variable Grouping with a Fixed Sub-problem Size Here the number of sub-problems m is fixed and in each co-evolutionary cycle, randomly selects a set of variables s for each sub-problem such that $n = m \times s$. In this way, RG looks to improve the probability of combining any interacting variables into the same sub-problems [83]. A proof presented in [151], concludes that there is a relatively high probability that two variables which interact will be part of the same sub-problem during at least two co-evolutionary cycles, and thus further conclude that with little domain knowledge, RG is quite effective at capturing variable interactions. On the other hand, if the number of interacting variables is large, the probability of finding these variables together in the same sub-problem, for at least one co-evolutionary

cycle, remains very low [99]. The first CCEA using random grouping, and using fixed sub-problem sizes - was proposed by Yang et al in [150] implemented as part of their cooperative co-evolution with differential evolution (DECC-I) approach, where results showed that, particularly for non-separable problems, it outperformed CCEA using a static grouping approach. Later, the authors improved upon this algorithm with the development of DECC-G which makes use of adaptive weighting strategy during its co-adaptation process; it was found that DECC-G outperformed standard DE as well as DECC using a static grouping scheme [151].

(A.ii) Random Variable Grouping with Dynamic Sub-problem Size According to Ma et al. [83] one of the main disadvantages faced when using static grouping and random grouping with a fixed sub-problem size is that in order to correctly set the sub-problem size s , one must have some prior knowledge of the problem domain in question. For example, a small s is suitable if the problems being addressed are separable in nature, where a larger s is more suitable for use against non-separable problems so as to improve the probability that interacting variables will be grouped together [83]. Therefore, having an approach which dynamically tunes the value of s to the domain is desirable [83]. In [150], additionally to DECC-I, Yang et al. DECC-II in which given a predefined range, randomly tuned the value of s on each cycle [83]. The probability of selecting s from $S = \{s_1, \dots, s_t\}$ is calculated based on its recent performance [83]. An approach presented in [57] (CCDECD) shows the applicability of the use of heuristic rules based on knowledge of the domain for the modification of s [83]. Several studies including [99, 100, 76] make use of a fixed s until no further improvements to solution fitness can be achieved at which time a new value of s is selected uniformly at random from $S = \{s_1, \dots, s_t\}$ [83]. Additionally to this, in [103] the value for s is decreased gradually as the search progresses.

2.5.1.2 Basic Algorithmic Frameworks for CCEA

In this section we describe two general frameworks which underpin most CCEA implementations, namely; the single population CCEA and the multi-population CCEA. In the single population scheme, individuals interact with other individuals from the same population; whereas, for the multiple population scheme, individuals can interact with individuals from one or more different populations. Pseudocode for each of these schemes, derived from [106], is given in Algorithm 1 and 2 respectively.

Algorithm 1 Single Population CCEA

```
1: procedure SP-CCEA
2:   population ← INITIALISE
3:   evaluators ← select evaluators from the population
4:   Evaluate each individual ∈ population by interacting with evaluators
5:   while terminate == false do
6:     parents ← select parents from population
7:     children ← produce children from parents using variational operators (x-
      over/mutation)
8:     evaluators ← select evaluators from parents + children
9:     Evaluate each individual ∈ children by interacting with evaluators
10:    Select survivors to be passed on to the next generation
11:  end while
12:  return solution
13: end procedure
```

Algorithm 2 Multiple Population CCEA

```
1: procedure MP-CCEA
2:   for each ( $pop \in populations$ ) do
3:      $pop \leftarrow INITIALISE$ 
4:      $evaluators \leftarrow \text{select } evaluators \text{ from } (populations - pop)$ 
5:     Evaluate each  $individual \in pop$  by interacting with  $evaluators$ 
6:   end for
7:   while  $terminate == false$  do
8:     for each ( $pop \in populations$ ) do
9:        $parents \leftarrow \text{select } parents \text{ from } pop$ 
10:       $children \leftarrow \text{produce } children \text{ from } parents \text{ using variational operators (x-}$ 
      over/mutation)
11:       $evaluators \leftarrow \text{select } evaluators \text{ from } (populations - pop)$ 
12:      Evaluate each  $individual \in children$  by interacting with  $evaluators$ 
13:      Select survivors to be passed on to the next generation
14:    end for
15:  end while
16:  return solution
17: end procedure
```

A few details have been omitted from both algorithms shown here. Firstly, it is common to find CCEA implementations which make use of a so called *archive*, a form of search memory maintained over generations i.e., throughout the whole evolutionary process [106]. Also, there is no indication about which specific interactions should be evaluated, how the results of any interactions can be combined in order to provide individuals a measure of fitness or how communication between populations is achieved when multiple populations are considered [106].

CHAPTER 3: CONTINUOUS OPTIMISATION BENCHMARK FUNCTIONS

3.1 INTRODUCTION

Firstly in this chapter, some of the main factors affecting continuous problem difficulty are discussed. In the next section, a short discussion of the decision to construct a new benchmark suite is provided. Finally in this chapter, we present a summary of the functions included in the new suite. Full descriptions of all functions included in the suite can be found in Appendix E.

The primary source used for the selection of benchmark functions in this suite was the highly cited review by Jamil and Yang [63]; cross-checked for the commonality of selected functions with other continuous benchmark function sets (see Chapter 3 Section 3.3). Others (i.e., Bent Cigar, Deflected Corrugated Spring, Inverted Cosine Wave, Levy, Rastrigin, Schwefel and Sum of Different Powers Functions) were found in various secondary sources [73, 93, 134, 36, 51]. One other reason for using secondary sources was simply to ensure that the function definitions and generated plots (see Appendix E) were correct and that corresponding definitions and plots were consistent between sources. However, another reason was that in order to construct a balanced suite of generally well known functions, the sets of functions such as those found in other benchmark suites, e.g., [133, 32, 47], could not be referenced in their entirety due to over/under representation of problems with certain features (see section 3.3 for discussion). Therefore, other - perhaps less well known functions - had to be selected where their definitions and outputs could be easily and reliably validated. Functions with discrepancies between their descriptions and output from several sources were therefore not used in order to help ease the replicability of the study being undertaken.

3.2 WHAT MAKES A FUNCTION DIFFICULT TO OPTIMISE

Several factors influence the difficulty in optimising an objective function [63]:

1. **Ruggedness** of the fitness landscape

2. **Dimensionality** of the objective function
3. **Separability** of the objective function
4. **Global Basin Structure**
5. **Variable Scaling**
6. **Search Space Homogeneity**
7. **Basin Size Homogeneity**
8. **Contrast Between Local and Global Optima**
9. **Plateau Sizes**

Each of these will be discussed in turn over the following few paragraphs.

RUGGEDNESS The ruggedness, or modality, of a function can either refer to: (i) the number of local optima, present in the landscape, or (ii) the number of global optima present in the landscape. Weise defines ruggedness simply as “From a simplified point of view, ruggedness is multi-modality plus steep ascends or descends in the fitness landscape”. We focus on multi-modality in the local sense so we use the term multi-modality to this effect. If a search algorithm encounters these optima in the search for global optimality, and becomes trapped in (unable to escape) that region of space, this can have a negative impact on search progress [63]. Therefore, functions with a large number of optima (multi-modal functions) tend to be more difficult to search than functions with few optima (also multi-modal) or a single optima (Uni-modal functions). Local search strategies will tend to perform poorly on multi-modal functions, demanding the use of a global strategy.

DIMENSIONALITY When considering objective functions of multiple dimensions, optimisation practitioners often run up against the so-called ‘curse of dimensionality’, an idea first highlighted by Bellman in 1961 in [8]. The curse of dimensionality refers to the phenomena where an exponential increase of volume in a function space occurs as new dimensions are considered. If for instance a 1D function was to be sampled evenly 20 times in the range $[0, 1]$, we would of course have a sample of size 20^1 . However to retain the same density of sample points in the 2D case, in terms of distance between consecutive points, we would require a sample of size $20^2 = 400$, $20^3 = 8000$ for the 3D case and a massive $20^4 = 160,000$ for a 4D sample. It is clear to see that for dimensions even higher than this, and likely including a smaller granularity of samples per dimension, that adequately sampling the available space quickly becomes intractable. The obvious consequence for search and optimisation is that

algorithms need to search over increasing numbers of possible solutions in the search for optimality - causing a rapid deterioration of solution quality obtained within a given time budget and the speed of convergence to the optimum solution when no budget is used. However, another less obvious consequence is that a search methodology that is found to be useful for problems in smaller dimensions may turn out to be almost useless when considering the same problems in higher dimensions. Therefore, given that real-world optimisation problems tend to be defined in highly n -dimensional space, optimisation methods that work well within some reasonably high range of dimensions are incredibly valuable.

SEPARABILITY A separable function is defined as a function of n parameters that can be rewritten as the product of n functions of one parameter each. That is, they can be easily partitioned into several sub-problems each with lower dimensionality and greater ease of searching [63, 12]. Combining the solutions to each sub-problem constructs a final solution.

There are two types of separable function related to the nature of this final solution construction: (i) additively separable functions and (ii) multiplicatively separable functions. Given a function F of n parameters $\{x_1, \dots, x_n\}$, a function can be described as additively separable if there exist a set of functions $\{f_1, \dots, f_n\}$ of one parameter, such that: $F(x_1, \dots, x_n) = f(x_1) + \dots + f(x_n)$. An example of such a function is that of the Rastrigin function (see section E.12). Similarly, given a function F of n parameters $\{x_1, \dots, x_n\}$, a function can be described as multiplicatively separable if there exist a set of functions $\{f_1, \dots, f_n\}$ of one parameter, such that: $F(x_1, \dots, x_n) = f(x_1)f(x_2) \dots f(x_n)$. The Alpine Function no.2 [63] (not included in our benchmark suite) can be described as being multiplicatively separable. In general non-separable functions tend to be harder to optimise than separable ones. The main reason for this is due to the fact that a separable function of n independent parameters can be optimised by n independent optimisation processes on each parameter [63].

Optimisation methods that can take advantage of the separable nature of functions tend to perform well. In order to increase the difficulty of functions for the purposes of benchmarking new and existing optimisation methods, and to consequently drive the development of ever more effective techniques, researchers have been making use of testbeds where the normally separable problem instances have been made non-separable through coordinate rotation by use of a rotation matrix [12, 114]. We do not consider functions that have undergone coordinate rotation in this thesis since (i) there are many inherently non-separable benchmark functions available, and (ii) for our purposes, it was useful to look at performances of algorithms against separable problems as well as non-separable as no comparative studies exist, to our knowledge, that show the behaviour of different metaheuristic algorithms for large-scale separable problems. Further, multiplicatively separable functions are currently beyond the

scope of the present study and so only functions that are additively separable are included in our benchmark suite. The reason for this is that due to the \prod term present in these functions, some solution values returned would eventually exceed the maximum value of the internal representation when searching above some threshold dimensionality. For fairness and in order to avoid the introduction of bugs in our algorithms - if modifying the internal representation - the decision was made to avoid the use of this type of function.

GLOBAL BASIN STRUCTURE For a local optimum \hat{x} , a *basin* or *basin of attraction* consists of the set of configurations from which a gradient walk, or steepest descent, will reach \hat{x} [12]. Therefore, since every basin contains at least the local optimum itself then every basin can be considered as being non-empty [12]. The global basin is a similar structure to local basins, however at a larger granularity instead considering the set local basins able to reach the global optimum by means of a suitable basin-hopping algorithm. For example, a highly multi-modal problem such as Rastrigin's Function (Section E.12) forms a clearly defined parabola of local minima and maxima when viewed from a distance tending towards the global optimum [12]; thus giving the impression that the non-linear Rastrigin function is convex. Problems without a global basin structure are considered harder to solve as the global optima can appear at any point in the search space with no overall guidance as to where it is located [12]. Some algorithms such as ILS and other basin-hopping approaches can take advantage of global basin structures in order to converge on the global optimum.

VARIABLE SCALING Variable scaling refers to the difference in scale with respect to the search space dimensional bounds. This means that it would be necessary to perform smaller steps in some dimensions and larger steps in others [12]. Even when the bounds on the search space are uniform across all dimensions, the problem may still behave very differently in certain dimensions than for others [12]. Algorithms such as CMA-ES (section 4.7) are particularly suited to this type of problem [12].

SEARCH SPACE HOMOGENEITY Many of the standard benchmark functions developed for testing the effectiveness of optimisation approaches are often designed with a homogeneous search space structure - often defined by an individual simple formula [12]. However, in the real world it is rare to find problems with such structures; instead, problems tend to have far more complex structure [12]. To combat this issue, many benchmark sets consist of functions that are composed of several functions blended into one such that the search space of the resulting function has different characteristics in different areas [12] - much more like how real-world problems might appear.

BASIN SIZE HOMOGENEITY The size of global basin is known to play a role in determining the hardness of a problem [12]. Torn et al. in [136] emphasise this fact by stating that if the basin is large the global optimum is easier to detect and the problem is easier to solve than for those where the global basin is far smaller. However, problems can also be caused in cases where many local optima basins exist. For example, many algorithms for use with multi-modal functions, such as niching methods, make the general assumption that basin sizes are mostly similar and as a result use appropriate distances to differentiate between the basins; of course, this means that if the variance of basin sizes is large then these techniques will perform poorly [136].

CONTRAST BETWEEN LOCAL AND GLOBAL OPTIMA This problem characteristic refers to the quality, or height, differences between the local and global optima when compared to the average quality of the search space as a whole. If this contrast is high then good locations in the search space are more easily detected [136].

PLATEAU SIZES Plateaus make optimisation problems harder as they do not provide any directional information required by trajectory-oriented approaches - large plateaus effectively partitioning the search space making it difficult to transition from one area to another [136].

3.3 CONSTRUCTING A CUSTOM FUNCTION SUITE VERSUS USING EXISTING BENCHMARK FUNCTION SUITES

As there is a need for a balanced suite of benchmark functions for this study, i.e., a balance of the features under investigation, one was constructed by sampling from the large quantity of existing - and widely used - functions. Previous attempts have been made to construct similar standard benchmark suites, such as: for continuous optimisation competitions like CEC-2005, CEC-2010 and CEC-2013. As stated in Chapter 2, Section 2.4.5.3, it is recommended that one should try to make use of existing and widely used benchmark sets, such as these, rather than constructing a separate suite.

However, the need for a balanced suite of functions restricted the direct use of the existing test suites; although many of the functions present in these suites have been included. There were several reasons for this. Firstly, the current focus of the continuous global optimisation community is to find approaches capable of solving problems which are resistant to current methods [21]. As such, and since effective algorithms currently exist for solving classes of

problem with uni-modal, low dimensional and separable characteristics, these are rather under-represented in these suites compared to problems with more resistant characteristics, such as: multi-modal, high-dimensional and non-separable [21]. Table 3.1 provides a summary of the proportion of problems with certain features present in these suites as compared to the suite constructed for this study. The proportions are also presented graphically in Fig. 3.1 for convenience. The Black-Box Optimisation Benchmark (BBOB) benchmark suite is not presented, due to the difficulty in verifying the problem features possessed by the benchmark functions used; specifically, the features are not always explicitly stated in the published suite descriptors [32, 47]. However, Clerc in [21], from which I also verify the correctness of some of the feature proportions found in these other suites for Table 3.1, provides data suggesting that this set is also too unbalanced for our purposes - in favour of multi-modal and non-separable functions.

Table 3.1: Summary of the Proportion of Problems With Certain Features Present in Various Existing Test Suites Versus the Test Suite Constructed for this Study

Feature / Suite	CEC-2005	CEC-2010 (LSGO)	CEC-2013	CEC-2013 (LSGO)	Bespoke Suite
<i>Uni-modal</i>	5/25 (20%)	8/20 (40%)	5/28 (18%)	7/15 (47%)	10/17 (59%)
<i>Multi-modal</i>	20/25 (80%)	12/20 (60%)	22/28 (82%)	8/15 (53%)	7/17 (41%)
<i>Separable</i>	2/25 (8%)	3/20 (15%)	4/28 (15%)	4/15 (27%)	8/17 (47%)
<i>Non-separable</i>	22/25 (88%)	2/20 (10%)	24/28 (85%)	3/15 (20%)	9/17 (53%)
<i>Partially-separable</i>	1/25 (4%)	15/20 (75%)	0/28 (0%)	8/15 (53%)	0/17 (0%)
<i>Less than 10 Dimensions</i>	0/25 (0%)	0/20 (0%)*	0/28 (0%)	0/20 (0%)*	17/17 (100%)

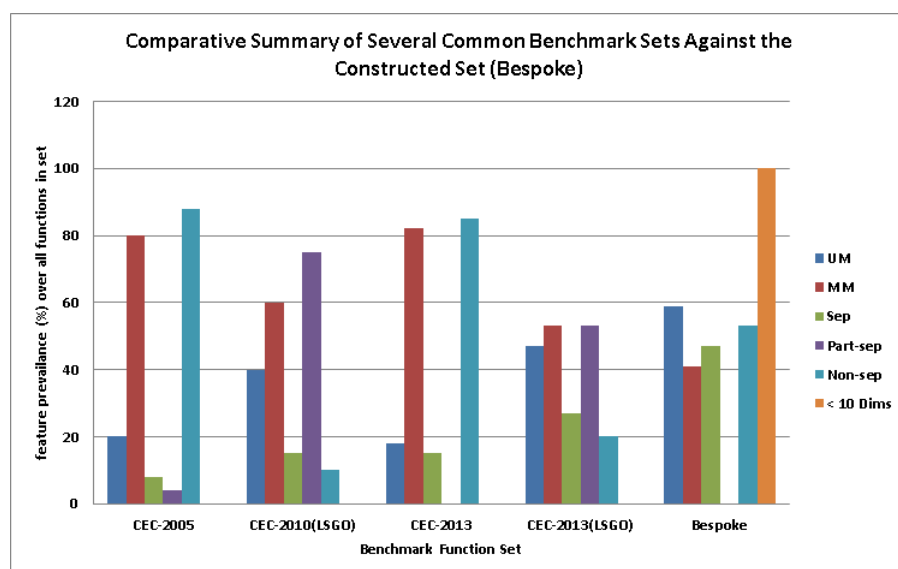


Figure 3.1: Barchart Summarising Prevalence of Problems With Certain Features Present in Each Benchmark Test Suite

From Table 3.1 and Fig. 3.1, we can see that the bespoke function suite used for this thesis is indeed more evenly distributed in terms of the problem featured of interest (i.e., separable/non-separable and uni-modal/multi-modal). Additionally, as we explore how different metaheuristics behave at different problem scales, it is required that low dimensional versions of our benchmark functions can be evaluated during our investigations. It can be clearly seen that every one of the published standard benchmark suites fails to meet this requirement, where for the ‘bespoke’ suite, any scale of instance for each function is possible.

3.4 FUNCTIONS

Appendix E provides details and definitions of the standard continuous benchmark functions we implemented and used in this thesis, where all included 2D function plots were generated using statistics package *R* through sampling the functions by a uniform step size in each dimension.

In order to eliminate positional bias in terms of the selected initial solutions when benchmarking algorithms, existing suites are commonly constructed from functions in which the coordinate system has been shifted with respect to the base function definition (e.g. ‘CEC-2005 Special Session on Real-Parameter Optimisation’ function suite [133]). In this way, the location of optima (including the global optima) can be modified whilst preserving the features and structure of the function. With respect to the experiments in this thesis, the issue of positional bias is addressed by selecting an initial solution in the search space uniformly at random for each individual run of an algorithm. Therefore, the use of shifted functions for the purposes of preventing positional bias is unnecessary.

The following section provides a table summarising key details of the functions included in the suite, including: domain of each function, location of the global minimum, value of the global minimum and a note of the features of interest (separability and modality) provided by each function.

3.4.1 Summary of Functions

Table 3.2: Summary of n-Dimensional Benchmark Function Suite

Function Name	Domain	Location of Global Minima/ Minimum	Global Minima Value(s)
Ackley Function	$-32.768 \leq x_i \leq 32.768$	$x_i^* = 0$	$f(x^*) = 0$
Alpine Function no.1	$-10 \leq x_i \leq 10$	$x_i^* = 0$	$f(x^*) = 0$
Bent Cigar Function	$-100 \leq x_i \leq 100$	$x_i^* = 0$	$f(x^*) = 0$
Brown Function	$-1 \leq x_i \leq 4$	$x_i^* = 0$	$f(x^*) = 0$
Chung-Reynolds Function	$-100 \leq x_i \leq 100$	$x_i^* = 0$	$f(x^*) = 0$
Deflected Corrugated Spring Function	$0 \leq x_i \leq 2\alpha^*$	$x_i^* = \alpha$	$f(x^*) = -1$
Exponential Function	$-1 \leq x_i \leq 1$	$x_i^* = 0$	$f(x^*) = -1$
Griewank Function	$-600 \leq x_i \leq 600$	$x_i^* = 0$	$f(x^*) = 0$
Inverted Cosine Wave Function	$-5 \leq x_i \leq 5$	$x_i^* = 0$	$f(x^*) = -n + 1$
Levy Function	$-10 \leq x_i \leq 10$	$x_i^* = 1$	$f(x^*) = 0$
Qing Function	$-500 \leq x_i \leq 500$	$x_i^* = \pm\sqrt{i}$	$f(x^*) = 0$
Rastrigin Function	$-5.12 \leq x_i \leq 5.12$	$x_i^* = 0$	$f(x^*) = 0$
Rosenbrock Function	$-30 \leq x_i \leq 30$	$x_i^* = 1$	$f(x^*) = 0$
Schwefel Function	$-500 \leq x_i \leq 500$	$x_i^* = 420.9687$	$f(x^*) = 0$
Sphere Function	$-5.12 \leq x_i \leq 5.12$	$x_i^* = 0$	$f(x^*) = 0$
Sum of Different Powers Function	$-1 \leq x_i \leq 1$	$x_i^* = 0$	$f(x^*) = 0$
Sum Squares Function	$-10 \leq x_i \leq 10$	$x_i^* = 0$	$f(x^*) = 0$

4.1 INTRODUCTION

In this chapter we give some background to the metaheuristic approaches implemented for this study. We cover the essentials of each technique and include descriptions of common operators and algorithmic variations found in the literature. When discussing several different operator implementations for a metaheuristic, the operator described last in the corresponding section was used in our implementation of the metaheuristic. For clarity, the operators for each algorithm are given as follows:

- Random Mutation Hill Climb (RMHC)
 - Neighbourhood function (described in Section 4.1.3)
- Steady-state Genetic Algorithm (SSGA)
 - Selection: Tournament Selection
 - Crossover: 1-Point Crossover
 - Mutation: Single Gene Mutation (as described in Section 4.3.2.3)
 - Replacement: Replace-Worst

Following the background of each metaheuristic algorithm, we outline the metaheuristic implementations used throughout our experiments along with a brief outline of the implementation specifics - highlighting important design decisions. We also provide pseudocode representing the actual JavaTM source code we produced or sourced externally. Finally, we describe the ‘tunable parameters’ for each metaheuristic optimised in the automatic parameter tuning phase (Chapter 6) along with their value ranges - as used by SMAC - and justification for their selection.

In the remainder of this section, we must provide some general details that affect multiple implementations, namely: the solution representation used consistently throughout all implementations and the neighbourhood function utilised by our single solution methods - RMHC and SA.

4.1.1 *Exposing Algorithm Parameters*

The selection of the parameters to be tuned by the selected automatic parameter tuning method - SMAC (see Chapter 5 Section 5.4) - for each of the algorithms used in this study were selected as being the most common parameters being tuned by practitioners as described in published works. Of course, more parameters for some of the algorithms could have been exposed beyond those commonly found in the literature; for example, the selection of the most appropriate boundary handling schemes for PSO or the crossover or mutation operator used by the GA. However, there are a couple of reasons for not including these kinds of parameters in the tuning procedures.

Firstly, parameters of the type given in the example above go beyond simply modifying the behaviour of the algorithms through the changing of values; selecting one setting over another would change the internal mechanics of the algorithm itself, in effect producing a different algorithm altogether. As tuning takes place for this study for each problem-dimension pairing and any comparison between algorithms in also in this context, the returned results and analyses could be considered invalid as essentially different algorithms are being compared - even when considering plots of a single problem with dimensionality along the x-axis.

The second reason exposing more parameters was problematic is that of the scale of the data generation that was to be conducted; specifically, the correspondingly large amount of time and computational cost required. As will be discussed later in Chapter 6 Section 6.1, the number of independent tuning procedures for each algorithm-problem-dimension triplet amounted to 1700. With the worst-case maximum number of runs of the target algorithm being carried out by SMAC during these processes set to 1000 - the recommended default number of runs; the wall clock time required to complete even a single set of parameter configurations for one algorithm (i.e., each problem-dimension pairing) grew enough on a single machine to warrant parallelisation of these independent tuning processes on our departments computer cluster. Even then, finding a full set of configurations of an algorithm would take upwards of a day and worst of all, tuning sep-CMA-ES with SMAC - considering its higher computational cost (due to calculating the diagonal covariance matrices) - would take almost 3 days to find a complete set of configurations. With these numbers in mind, since SMAC - like all automated parameter configuration tools - is still bound by the 'curse of dimensionality', a parameter space volume increasing exponentially with each new parameter considered would call for a larger maximum number of target algorithm runs to ensure SMAC is still able to find effective configurations. Further, since SMAC is probabilistic in nature, more repeats of the entire tuning process would have to be carried out to improve the likelihood

and therefore confidence that SMAC has - to the best of its ability - found the most effective configurations for each algorithm.

4.1.2 *Solution Representation*

For each of the implemented metaheuristics, a fixed double solution representation was used. Each parameter p in a solution of size d , where d is the dimensionality of the problem, can take on a value from the interval $[0.0, 1.0] \subset \mathcal{R}$ which is then later scaled by the objective function i.e., the benchmark problem, to fall within the interval $[UB, LB] \subset \mathcal{R}$ where LB and UB are the lower and upper bounds for the given problem dimension respectively. There is a couple of reasons to having the solutions encoded in this way. Firstly, the implementations are independent from the problem being solved; allowing the implementations to be reused on different problems - should time allow - e.g., real-world problems, without requiring significant modification therefore avoiding the possibility of programming error that may introduce noise in further comparisons. Secondly, the search space of the encoded representation and the decoded solution are equivalent, that is, the representation does not change the nature of the search landscape and its associated features. The encoded solutions are decoded within the objective functions using the following:

$$x'_i = LB + x_i \times (UB - LB) \quad (4.1)$$

Where x_i represents an encoded parameter i of an encoded solution representation x and x'_i is the decoded (scaled) parameter. Once scaled, the solution can then be used directly in evaluation of the objective function.

4.1.3 *Neighbourhood Function*

For those algorithms that require the use of a neighbourhood for generating new solutions in the search space from a given point, namely: random mutation hill climbing and simulated annealing, a common neighbourhood function is used. If working within discrete spaces - e.g., if the problems were combinatorial in nature - a neighbourhood can be quite easily defined as being a move to the next (or previous) discrete state in one or more dimensions. However, given that optimisation takes place in a continuous context in these studies - i.e., the notion of a next discrete step in any given direction is meaningless due to continuous nature of values in any given direction - steps within the space have to be defined in terms of a 'step size' variable. It is common to use a fixed step size tuned to a given optimisation problem, however this strategy has several disadvantages (see Fig. 4.1) and so, dynamic step sizes that scale with

search progress and/or solution state are sometimes used. When using a fixed step size, it is critical to the performance of the search algorithm that an appropriate value be selected as it imposes a structure (as viewed by the algorithm) to the underlying search landscape. A value that is too large can make the landscape appear flatter to the algorithm making it harder to traverse the landscape as well as hampering progress towards local optima - including global optima, and in general moving the algorithms behaviour closer to that of random search. On the other hand, if the step size value is too small, it may be more difficult to escape from local optima and by utilising small moves the region of global optimality may never be reached with an allocated computational budget, i.e., the algorithm is too slow to converge. These issues with fixed step sizes can be better illustrated through example as shown in Fig. 4.1 where we assume an approach that can accept worsening solutions such as SA.

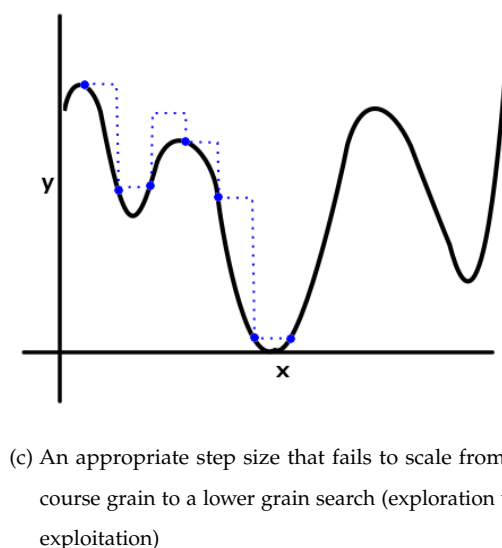
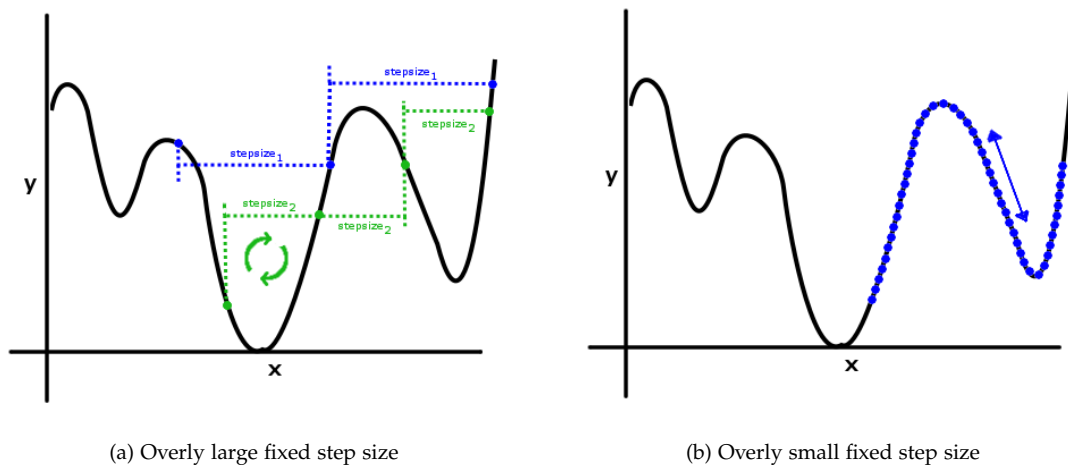


Figure 4.1: Effects of inappropriate fixed step sizes

In Fig.4.1a(a), an inappropriately large value for the step size can cause the search to 'skip over' local optima as with $stepsize_1$ or end up becoming trapped within the basin of a local optimum unable to reach the minimum like $stepsize_2$. In Fig. 4.1b(b) on the other hand, a step size which is too small makes very slow progress towards the minimum. Even with a suitably appropriate step size for the problem being solved it can still affect the quality of the final solution returned - even if the global basin has been discovered - this can be seen in Fig.4.1c(c). Although this is less likely to cause problems similar to the step size in Fig.4.1b(b) - due to the scale of the basin of attraction in relation to the step size - this example shows that there is still a possibility of this occurring as the search nears optimality.

Several methods exist for dynamically selecting step sizes and have been found to be quite successful albeit far more computationally expensive; however, for brevity we will not discuss them further.

However, another method exists, and is the one we use for our algorithms in this study, namely, the use of stochastic step sizes. We use this here as it combines the benefit of low computational cost when using fixed step sizes with some ability to scale appropriately as the search progresses, as found in dynamic step methods. In order to avoid the problems of overly large steps sizes - as described above - a step size is not simply selected before execution, bounded only by the boundary constraints of the problem, but rather a step size is selected uniformly at randomly on each iteration of the algorithm from a neighbourhood defined by a user supplied parameter and centred at the current solution. If the neighbourhood exceeds any bounds of the problem then the violating bound of the neighbourhood is set to that of the violated bound. All of this means that the neighbourhood can be larger than allowable when compared to one when using a strictly fixed step size, giving some of the benefits of a higher grained exploration of the search space, but since the defined neighbourhood essentially contain all possible neighbours perpendicular to the current solution within its bounds, a stochastic selection from among these neighbours can allow for a finer-grained search.

4.2 HILL CLIMBING ALGORITHMS

Hill climbing algorithms are a family of local search metaheuristics originally developed for searching within discrete combinatorial search spaces - although the algorithm can be adapted to work with a variety of different solution encodings [16]. As with many local search methods, hill climbing algorithms tend to maintain only a single individual throughout the entire search

procedure [115]. Although there are many varieties, many of which were developed to solve issues found with other hill climbing variations, the general functioning of these algorithms are quite similar and tend to differ mainly in how candidate solutions are sampled [96]. Hill climbing, or descent if dealing with minimisation, therefore attempts to discover solutions of better quality by iteratively selecting a member of the immediate local neighbourhood and testing whether or not there is an improvement - accepting an improving solution as the new current solution and rejecting worsening solutions [96]. Given this description, hill climbing methodologies can be said to be composed of two main phases: (i) a step function, that defines the neighbourhood of the current solution to be searched and (ii) an acceptance strategy determining what solutions are to be accepted as a new current solution. Although we discuss hill climbing algorithms here, these phases also apply to most other local search methods.

4.2.1 *Common Variations*

4.2.1.1 *Simple Hill Climbing Algorithm*

Simple Hill Climbing, also referred to as Next-best Neighbour and Next-descent/ascent Hill Climbing, is a deterministic search algorithm which involves the generation of neighbours sequentially one at a time, by the neighbour operator(s), and evaluating their objective value. If any improvement is found over the current solution, the neighbour becomes the new current solution without evaluating any more solutions of the previous current solution [35]. Pseudocode of simple hill climbing is shown in Algorithm 3 derived from the description of the approach found in [35].

Algorithm 3 Pseudocode of Simple Hill Climbing

```
1: procedure HILLCLIMB
2:   currentSolution  $\leftarrow$  RANDOMSOLUTION
3:   terminategetsfalse

4:   while terminate == false do
5:     candidateSolution  $\leftarrow$  NEXTNEIGHBOUR(currentSolution)  $\triangleright$  Get next reachable
       neighbour
6:     if (SCORE(candidateSolution)  $\leq$  SCORE(currentSolution)) then
7:       currentSolution  $\leftarrow$  candidateSolution
8:     end if
9:     if (currentSolution.score is good enough) then
10:      terminate = true
11:    end if
12:  end while
13:  return currentSolution  $\triangleright$  return best solution discovered
14: end procedure
```

4.2.1.2 Steepest Descent Hill Climbing Algorithm

In this variation on simple hill climbing, all possible neighbours of the current solution are considered - at least when considering discrete representations - and the best neighbour solution is selected to compete with the current solution [89]. If the best neighbour is found to be better than the current solution, then neighbour replaces it and the whole process begins anew [89]. If no neighbour is found that is better than the current solution then the algorithm has converged on a local optimum and the algorithm can either terminate here or the search can be restarted from a random position in the search space - keeping note, but making no further use of, the best solution discovered during the previous search [89].

Steepest descent can be extended to continuous domains, as with all hill climbing approaches, by introducing the notion of a step size parameter which defines a discrete neighbourhood for the current solution. For course, one must be careful of the value selected for the step size. If the size is too small the neighbourhoods can become very large - increasing computational cost, the speed of algorithm will be bounded by the small step size and once a local optimum is discovered the step size will be too small to jump over to other 'hill/troughs' [82]. If on the other hand the step size is too large, the algorithm may skip over local optima

and lose all information related to promising search trajectory; moreover, the algorithm will have a difficult time converging to a local optimum [82]. The consequences of the selected step sizes are best illustrated by Fig. 4.1 in Chapter 6.

Pseudocode of steepest descent hill climbing is shown in Algorithm 4 derived from the pseudocode found in [89]. This pseudocode includes random restarts which makes it more akin to a global optimiser such as those found in section 4.2.2.

Algorithm 4 Pseudocode of Steepest Descent Hill Climbing

```

1: procedure SD-HILLCLIMB
2:    $t \leftarrow 0$ 
3:    $overallBest \leftarrow INITIALISE$ 

4:   while ( $t \neq t.MAX$ ) do
5:      $foundLocal \leftarrow FALSE$ 
6:      $currentSolution \leftarrow RANDOMSOLUTION$ 
7:     while ( $\neg foundLocal$ ) do
8:        $neighbours \leftarrow GETALLNEIGHBOURS(currentSolution)$ 
9:        $candidateSolution \leftarrow GETBEST(neighbours)$ 
10:      if ( $SCORE(candidateSolution) \leq SCORE(currentSolution)$ ) then
11:         $currentSolution \leftarrow candidateSolution$ 
12:      else
13:         $foundLocal \leftarrow TRUE$ 
14:      end if
15:    end while
16:     $t \leftarrow t + 1$ 
17:    if ( $currentSolution \leq overallBest$ ) then
18:       $overallBest \leftarrow currentSolution$ 
19:    end if
20:  end while
21:  return  $overallBest$  ▷ return best solution discovered
22: end procedure

```

4.2.1.3 Stochastic Hill Climbing Algorithm

Stochastic Hill Climbing SHC is a local optimisation algorithm and is an extension of deterministic hill climbing approaches such as simple hill climbing (first-best neighbour) and steepest

ascent/descent hill climbing [16]. One of the most popular implementations of SHC, and the one utilised for this study, is that of Random Mutation Hill Climbing (RMHC) [16] introduced by Forrest and Mitchell in [35] where it was used in conjunction with other hill climbing approaches to investigate the behaviour of genetic algorithms (section 4.3). Pseudocode of a stochastic hill climbing approach - specifically the RMHC algorithm by Forrest and Mitchell as derived from [16] - is given in Algorithm 5.

Algorithm 5 Pseudocode of Stochastic Hill Climbing (RMHC)

```

1: procedure RMHC
2:   currentSolution  $\leftarrow$  RANDOMSOLUTION

3:   while terminate == false do
4:     candidateSolution  $\leftarrow$  RANDOMNEIGHBOUR(currentSolution)
5:     if (SCORE(candidateSolution)  $\leq$  SCORE(currentSolution)) then
6:       currentSolution  $\leftarrow$  candidateSolution
7:     end if
8:   end while
9:   return currentSolution ▷ return best solution discovered
10: end procedure

```

For RMHC, the call to RandomNeighbour mutates a single locus of the current solution at random [16], in effect taking a single dimensional step. This contrasts with simple hill climbing, and is in fact the only difference, which generates each neighbour of the current solution in turn until one is found that is improving.

4.2.2 Global Optimisation with Hill Climbing

As hill climbing algorithms are local methods that maintain only a single solution at a time, when performing global optimisation, they can easily become trapped in local optima. Extension of the basic algorithms to search within multiple local neighbourhoods provides the needed diversification to avoid this problem. Approaches such as: random-restart hill climbing, local beam search and stochastic beam search can provide diversification through repetition of the search process - accomplished by various means each with different properties. Simulated annealing can also be considered as a global variation of hill climbing utilising non-deterministic search, however, this is covered separately in section 4.5. The next few subsections will cover these other global hill climbing approaches in more detail.

4.2.2.1 *Random-restart Hill Climbing*

Another approach to improving on the quality of solution discovered by local search techniques such as hill climbing is to simply restart the search from a different randomly selected initial point in the search space. This is the basis of random restart hill climbing, the motivation being the hope that the global optimum will eventually be found through trial and error. Again there are several variations, each determined either by how a restart of the algorithm is triggered - after a given length of time or until the algorithm terminates naturally - or whether a fixed number of restarts are used or not. These variations are not mutually exclusive and can be combined. An obvious disadvantage of this approach is that it requires more time and computational resources than a single repeat of the standard hill climbing algorithm used, increasing linearly with the number of restarts required. Another problem with this approach is that as the problem instance size increases, random sampling has a much reduced chance of finding a point in the search space that would lead the algorithm to converge on a solution of higher quality [80] - this chance reducing further as the current solution quality approaches optimality. Therefore, reaching these higher quality solutions requires more of a biased sampling which can be facilitated by stochastic search [80]. Despite its issues, random-restart hill climbing is fairly easy to implement, requiring only one or two additional parameters, an additional loop surrounding the main hill climbing algorithm and some very minor modifications e.g. keeping track of the overall best solution globally rather than returning the best overall solution from a single local repeat.

4.2.2.2 *Local Beam Search*

Instead of making use of only a single solution as in local hill climbing algorithms, local beam search maintains a population of k solutions. Here, k solutions are randomly generated over the search space and all members of their neighbourhoods are generated [127]. The k members over all neighbourhoods are selected and the algorithm begins a new iteration. Although this might appear as though several repeats of hill climbing are being performed concurrently, the repeats run independently and no information is shared between search processes; on the other hand, local beam search does share information between localities by considering all neighbourhood members to become part of the newly maintained set of size k [127]. In this way, areas of the search space that do not hold much promise can be rejected in order to redirect search efforts to more promising regions. This however leads to the disadvantage that over time the maintained set of k solutions can still converge to a local optimum as weaker but more diverse solutions and regions are quickly rejected.

4.2.2.3 *Stochastic Beam Search*

To address the disadvantage of local beam search eventually becoming stuck in local optima, stochastic beam search does not select the best k members from the whole list of k neighbourhood members, but instead selects k members randomly with a probability of selection given by an increasing function of a members quality [127].

4.2.2.4 *Iterated Local Search (ILS)*

The main basis of iterated local search (ILS) is that the search should not focus on the entire set of all possible feasible solutions S , but instead, the resulting solutions returned from some underlying heuristic - usually local search - in other words, the far smaller set S^* of locally optimal solutions s^* [81]. The search behaviour can therefore be categorised as a markov chain of locally optimal solutions returned from this underlying heuristic, leading to solutions better than those that can be found through random repeated trials - as with random restarts [80]. In practice, local search methods such as hill climbing have been the most commonly used search heuristic in ILS, but in fact, any optimisation algorithm can be used even if this is deterministic rather than stochastic [80, 81].

Lourenco, Martin and Stützl in [80] and [81] state two factors that determine whether an algorithm is an instance of ILS:

1. There must only be a single chain of solutions being followed - this then excludes the use of any population-based heuristics
2. The search for improving solutions must be taking place within a reduced search space defined by the output of the underlying search heuristic

ILS requires a mechanism for moving within the space S^* of local optima. Unlike random restart this mechanism - a perturbation of the current s^* - provides a sampling biased by increasingly better local optima. However, as Blum and Roli point out in [11], it is - in most cases - infeasible to construct such a neighbourhood structure as S^* ; therefore the trajectory between local optima has to be performed without introducing the notion of a neighbourhood structure explicitly [11].

However, and as discussed in section 4.2.2.1, the bias provided by ILS will still help alleviate the issue faced by the random restart method where random sampling has a decreasing chance of discovering new points in the space that would lead to convergence on a new s^* . The perturbation of the current local optima should neither be too large or too small

[80]. Too large and the new intermediate solution s' will be random, we lose the useful bias and the algorithm will degenerate into a random restart-like approach [80, 11]. If conversely, the perturbation is too small, the search will often return back to s^* meaning that little of the space of S^* will be investigated [80, 11].

4.2.3 *Random Mutation Hill Climb: Implementation*

4.2.3.1 *Details and Description*

As detailed above, a random mutation hill climbing algorithm does not generate all improving neighbourhood solutions, as with steepest ascent/descent hill climbing - but rather 'mutates' a single solution parameter selected uniformly at random and also mutating the selected parameter by uniform mutation (Section 4.3.2.3) - and accepts the solution as the new current solution if and only if the new solution is improving. Furthermore, the improving neighbour is not accepted by a probability based on the magnitude of the improvement but is instead accepted immediately as the 'first choice' generated. As with the simulated annealing implementation (Chapter 4 - Section 4.5) this implementation makes use of the neighbourhood function described in Section 4.1.3. As such, the implementation only uses a single tunable parameter - *stepSize* - that defines the neighbourhood of a solution.

4.2.3.2 Pseudocode

Algorithm 6 Pseudocode of RMHC Implementation

```
1: procedure RMHC(stepValue, pSize) ▷ Main loop procedure
2:   problemSize ← pSize
3:   stepSize ← stepValue
4:   S ← GENERATEINITIALSOLUTION(problemSize)
5:   Sscore ← EVALUATESOLUTION(S)

6:   while (terminate == false) do
7:     N ← GENERATENEIGHBOUR(S, problemSize, stepSize) ▷ get random neighbour of
       S
8:     Nscore ← EVALUATESOLUTION(N)

9:     if (Nscore < Sscore) then
10:      S ← N
11:      Sscore ← Nscore
12:     end if
13:   end while

14:   return S ▷ return best solution discovered
15: end procedure

16: procedure GENERATEINITIALSOLUTION(problemSize)
17:   solution ← init
18:   for (i = 0 to problemSize-1 by 1) do
19:     solution[i] ← uniform random float ∈ [0.0, 1.0]
20:   end for
21:   return solution
22: end procedure

23: procedure GENERATENEIGHBOUR(solution, problemSize, stepSize)
24:   neighbour ← solution
25:   r ← uniform random float ∈ [0.0, 1.0]
```

```
26: boundMin  $\leftarrow$  neighbour[i] - stepSize
27: boundMax  $\leftarrow$  neighbour[i] + stepSize

28: if (boundMin < 0.0) then                                 $\triangleright$  Boundary handling
29:     boundMin  $\leftarrow$  0.0
30: else
31:     if (boundMin > 1.0) then
32:         boundMin  $\leftarrow$  1.0
33:     end if
34: end if

35: boundRange  $\leftarrow$  boundMax - boundMin
36: r  $\leftarrow$  uniform random float  $\in$  [0.0, 1.0]
37: step  $\leftarrow$  boundRange * r
38: neighbour[i]  $\leftarrow$  boundMin + step

39: return neighbour
40: end procedure
```

In this implementation, as with all the implementations described in this section, the procedure *evaluateSolution* is provided by the benchmark problem suite implementation presented in Chapter 3 of which all implementations make use. Therefore, *evaluateSolution* is not expanded in the pseudocode; however, the solution is decoded and evaluated by the benchmark suite in the manner described in section 4.1.2.

4.2.3.3 Tunable Parameters

Only one parameter is used that can affect the performance of the algorithm, specifically, *step size*. As discussed in section 4.1.3, this parameter defines the maximum radius of an axis-aligned hyperplane along each dimension of the search space centred on the current solution location. A neighbouring point - perpendicular to the current solution - from within the hyperplane is then generated stochastically by uniformly selecting a new value for one of the solution parameters. The result is a stochastically variable neighbourhood with the ability to scale with search granularity.

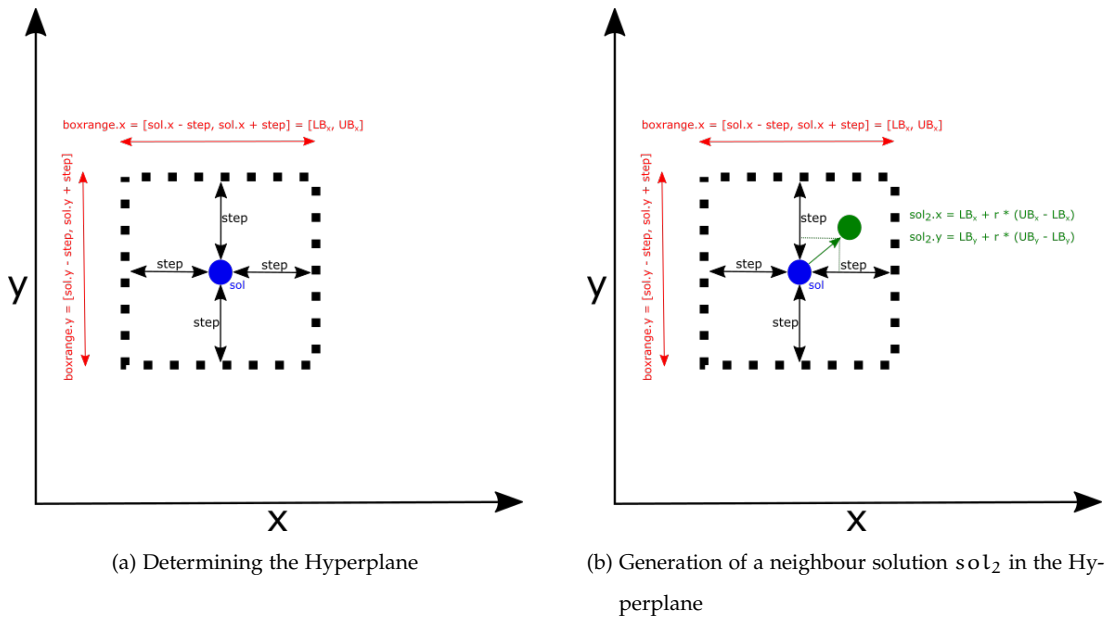


Figure 4.2: Definition of a neighbourhood of a solution sol using the stepsize parameter and the stochastic generation of a neighbouring solution sol_2

In the example shown in Fig. 4.1.3, a neighbourhood is defined by firstly determining its lower and upper bounds in each dimension i of the problem (Fig. 4.2a) - where $i \in [1, \dots, d]$ and d is the problem dimensionality. Since the box is to be centred on the current solution S , the lower bound LB_i and upper bound UB_i of each box dimension is calculated by:

$$LB_i = S_i - stepSize \tag{4.2}$$

$$UB_i = S_i + stepSize$$

where S_i is the position of the solution S in the search space in the i^{th} dimension. To then generate a neighbouring solution N , a different uniform pseudorandom value $r_i \in [0.0, 1.0]$ is generated in a single random dimension i and is applied to the current solution in the range of LB_i and UB_i as follows:

$$N_i = LB_i + r_i \times (UB_i - LB_i) \tag{4.3}$$

Through initial experimentation using the neighbourhood function as shown in Fig. 4.1.3, generating neighbours along a single dimension only i.e., directly perpendicular, was found to be a more effective strategy than those generated through perturbing along several dimensions. Thus, we opted to use the single dimensional version for each of the algorithms making use of a neighbourhood function.

4.3 GENETIC ALGORITHMS (GA)

Genetic algorithms are a population-based evolutionary algorithm, developed and introduced by Holland in the 1960's. Genetic algorithms (GA) are inspired by natural evolution and genetics [17]. In fact, Holland specifies in [54] that genetic algorithms are based on the classic view of chromosomes as a sequence of genes; a view that eventually led to the development of mathematical genetics founded by Fischer in [34]¹, which provide formulae that describes the rate of spread of given genes throughout a population. A GA is therefore a more generalised and computer executable equivalent of the formulation used by Fischer in mathematical genetics [54]. A typical implementation of a genetic algorithm might consist of the following:

1. A population
 - A set of candidate solutions to a given problem being solved
2. An Objective Function or "Fitness"
 - A measure of quality used to determine good solutions from not so good solutions in the population
3. A Selection Operator
 - A strategy of determining which individuals in a population, at a given time, are allowed to perform crossover and which ones can be carried across to the next generation. There are many such strategies.
 - Solutions which are among the fittest in a given population are able to bypass the selection procedure employed, whereby a small portion of the very best solutions in the population are brought brought across to the next generation, typically without undergoing any changes - mimicking the natural concept of 'survival of the fittest'. This is known as elitism or elitist selection.
4. A Crossover Operator

¹ First published in 1930

- Also referred to as a recombination operator, this is a method of combining the existing genetic information from two 'parent' solutions together in order to create one or more 'child' solutions. In search terms, crossover represents the main source of 'exploitation' during a search i.e., along with selection operators, it facilitates the convergence of the population towards a local optimum
- Typical methods of combination include: uniform crossover, one-point crossover and two-point crossover. There are however a wide selection of possible alternatives.

5. A Mutation Operator

- A means of ensuring that the population contains enough diversity so as to be able to better cover the space of possible solutions. Mutation in a search is the main source of 'exploration' i.e., it is considered a divergence operator which works against selection and crossover to help prevent premature convergence by allowing the search to escape local optima and explore the global search space.
- Some selection strategies, such as elitism, help create less diversity in a population by only accepting better solutions as parents and replacing the worst solutions with their children - meaning that the population will eventually converge to sharing a common set of traits with the side-effect of causing the crossover operator to become ineffective.

The main generalisations of a GA from mathematical genetics, as discussed by Holland in [55] and [54], consist of:

1. A concern with the interaction of genes on a chromosome rather than assume that genetic alleles can act independently from the others [55, 54]
2. An enlargement of the set of available genetic operators that can be used, over and above the mutation operator used almost exclusively in mathematical genetics. [55, 54]

From the first generalisation, the fitness function becomes complex and non-linear which quite often cannot be approximated usefully through the summation of the effects from different genes [55]. The second generalisation places more emphasis on the genetic operators, in particular on crossover [55]. Holland considers crossover to be essential to the success of a GA since crossover is an operation that is carried out in all mating organisms in nature as opposed to mutation whereby the mutation of a given gene takes place at a far lower frequency (less than 1 in a million individuals) [55]. The effect of crossover within a GA will be discussed more fully in section 4.3.2.2.

4.3.1 The Basic GA Procedure

A basic instance of GA is illustrated by pseudocode given in Algorithm 7 and the process flowchart shown in Fig. 4.3 and can be implemented as follows:

1. A population of candidate solutions is maintained by the algorithm, each one representing a potential solution to the given problem. Each individual solution is encoded as a number of variables, each representing some feature of the problem being optimised, whose values are taken from some predetermined alphabet - for example, binary, integer and real numbers - (Alg. 7: line 2)²
2. Each individual in the population is assigned a fitness score representing the quality of the solution, i.e., how well it solves the problem according to the objective function, which ultimately determines whether or not the solution will be selected for crossover - (Alg. 7: line 4)
3. A number of solutions are selected according to some selection strategy. These strategies can often be categorised as being 'elitist' or 'non-elitist', however, in practice selection strategies tend to fall within a spectrum with selection strategies being described as more or less elitist than others. A common strategy that can be said to fall at the centre of this spectrum is tournament selection described in section 4.3.2.1 - (Alg. 7: line 6)
4. The selected solutions are recombined, by application of a crossover operator, to create one or more candidate (child) solutions that somehow inherit components from multiple parents - (Alg. 7: line 7)
5. Given some probability, typically very small, each generated child solution may undergo mutation whereby one or more of its genes is modified in some way - (Alg 7: line 8)
6. Child solutions are then replaced into the population either by replacing other - often less fit - individuals or by generating a new population using the child solutions - (Alg 7: line 9)

In this procedure, steps 2 to 5 represent the main GA loop (Alg 7: line 3) and are therefore repeated until a user specified termination condition is satisfied.

² The individual solutions in a GA are designed to be analogous to natural genetic structures; therefore, chromosomes (solutions) are composed of a number of genes (representing problem variables/parameters. This terminology amongst others from evolutionary genetics is commonplace in the study of genetic algorithms and some other evolutionary algorithms.

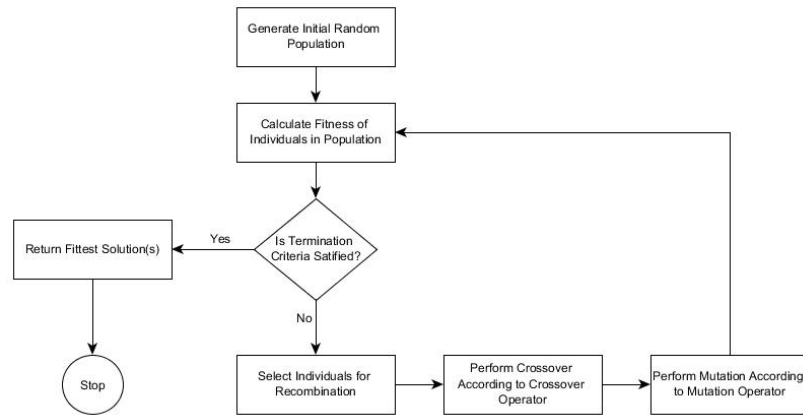


Figure 4.3: Basic operation of a typical genetic algorithm implementation

Algorithm 7 Pseudocode of a Basic Genetic Algorithm

```

1: procedure GA
2:   population ← INITIALISEPOPULATION
3:    $S_{best} \leftarrow \emptyset$ 
4:   terminate ← false

5:   while terminate == false do
6:     EVALUATEPOPULATION(population)
7:      $S_{best} \leftarrow$  GETBEST(population)

8:     parents ← SELECTPARENTS(population, numberOfParents)
9:     children ← CROSSOVER(parents, numberOfChildren)
10:    children ← MUTATION(children)
11:    population ← REPLACE(population, children)
12:    if ( $S_{best}$ .score is good enough) then
13:      terminate ← true
14:    end if
15:  end while
16:  return  $S_{best}$  ▷ return best solution discovered

17: end procedure

```

4.3.2 Genetic Operators

GAs work on three simple rules of genomic evolution analogous to those observed in natural evolutionary genetics; referred to as the genetic operators of the algorithm. Each operator, namely: selection, crossover and mutation [29], has a specific role to play within the evolutionary process and work in conjunction to maintain the course of improvement within the population over time as well as maintaining the stability of the system. Since we focus on continuous (real-valued) problems in our study, all genetic operators will be discussed in terms of a real-valued encoding scheme.

4.3.2.1 Selection

The responsibility of the selection operator in a GA is to choose a number of solutions from the population to become parents in order to produce offspring that will make up the next generation [29]. In this way, it acts as the force that increases the mean quality of the solutions in the population [29]. Parent selection in evolutionary algorithms such as GA is typically a probabilistic procedure - one in which the solutions of higher quality have more chance of being selected as parents [29]. However due to this probabilistic nature of selection, lower quality solutions typically have a small - but positive - opportunity of being selected; this helps to maintain diversity in the population otherwise the search can quickly become trapped in local optima [29]. Over the following paragraphs, we briefly describe the selection operator used in the GA studied for this thesis - namely, Tournament Selection - however, firstly, a couple of other common selection operators found in the literature are presented for comparison purposes: Fitness Proportionate Selection and Rank-based Selection.

(A) **FITNESS PROPORTIONATE SELECTION** Each solution in the population becomes a parent with a probability proportionate to its fitness [77]. Higher quality solutions have more chance of being selected to propagate their information to the next generation than those of lower quality. Therefore this strategy applies a selection pressure to those individuals in the population with higher fitness - to drive further progress. Two fitness proportional selection methods are given in the literature, namely: Roulette Wheel Selection and Stochastic Universal Sampling (SUS) - with the latter being developed to solve limitations found for the former - however, for brevity, we will only be discussing Roulette Wheel Selection here.

(A.i) Roulette Wheel Selection In Roulette Wheel Selection RWS, a probability of selection, proportionate to fitness value, is first determined for each solution in the population given by the equation [77]:

$$p_i = \frac{f(x_i)}{\sum_{j=1}^n f(x_j)} \quad (4.4)$$

where, $f(x_i)$ is the fitness of the solution x_i and n is the total number of solutions in the population. This equation provides, for each solution, a real valued probability in the range $[0, 1]$. Considering each of these probabilities as corresponding to a portion of a roulette wheel, a random number from the range $[0, 1]$ is generated representing a roulette ball landing in a portion of our hypothetical roulette wheel, selecting an individual for inclusion into the mating pool [22]. To gain an intuition of roulette wheel selection, an example is given in Fig. 4.4.

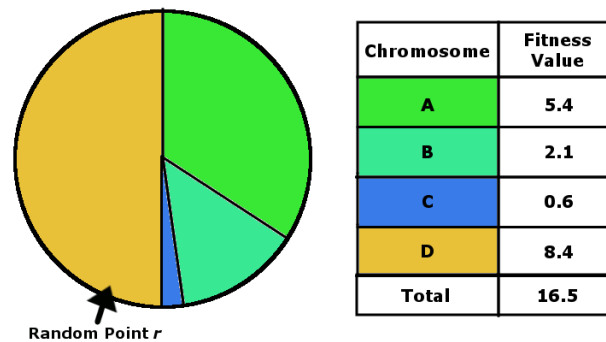


Figure 4.4: Roulette Wheel Selection example showing a single selection over a population of four individuals. Each individual, A,B,C and D represent a 32.7%, 12.7%, 3.6% and 51% portion of the roulette wheel respectively. Random point r selects the portion representing chromosome D similar to a hypothetical roulette ball

In this example, solution D has the highest probability of being selected, with solution C being the lowest. This is a useful feature of roulette wheel selection; although the solutions of higher fitness are more likely to be selected, it does not exclude the possibility that some solutions of lower fitness may be selected to contribute to the next generation, thus preserving some diversity in the population. The sum total of all the fitnesses t is calculated and a random point r is selected in the range $[0, t]$. The fitness of each individual in the population is subtracted from the value of r until $r < 0$. The individual whose fitness causes r to become negative is selected and the process is repeated for the total number of selections to be made. There is also nothing preventing the same individual to be chosen more than once and this can indeed happen. The repeated subtraction step in the above example, equivalent to a linear search over the population elements gives this implementation of roulette wheel selection a time complexity of $O(n)$; if, however, an implementation were to make use of binary search,

this can be reduced to $O(\log n)$ [77]. Another possibility is to make use of the implementation outlined in [77] whereby “stochastic acceptance” is used to select an individual with a typical time complexity of $O(1)$ depending on the distribution of fitness weights.

Although a popular selection operator, roulette wheel selection carries some limitations. Firstly, the operator in its current form is not suitable for minimisation problems or problems where fitnesses can take on negative values - although, both issues may be overcome with appropriate fitness scaling or transformation [5]. Another limitation is that since the likely variance of fitnesses in the population is initially high, the fittest solution can quickly take over the population leading to premature convergence [38, 5].

(B) TOURNAMENT SELECTION In Tournament Selection, a competition, or tournament, is held between n solutions chosen randomly from the current population. The winner of each tournament is allowed to pass through to the intermediate population - either with or without replacement ³[38]. Usually the size of each tournament is kept small, this is due to the increasing percentage of diversity invariably lost as the size of tournament groups increase; likely resulting in premature convergence [38]. An example of tournament selection is shown in Fig. 4.5.

Tournament selection is popular amongst GA practitioners and has become the main selection method used in GA for several reasons [82]:

1. Unlike fitness proportionate selection, tournament selection is readily compatible with minimisation problems, fitness functions that may return negative values and generally any other particulars of the fitness function [82]
2. The selection pressure is easily controlled by changing the size of the tournament [82]; larger sizes make it harder for weaker solutions to compete as the chance of much better solutions being included in a tournament increases with group size
3. It provides a means of automatic fitness rescaling in that solutions of similar fitness near optimality are not considered as equal - as would be the case with fitness proportionate selection and SUS [82]

In terms of tournament sizes, a typical tournament size is $t = 2$ [82] however for certain representations it is common to be a bit more selective by using a larger tournament size [82]. It is also worth noting that $t = 1$ is equivalent to random selection, however, if $t \geq p$, where p

³ With replacement means that a selected individual can be selected more than once from the population meaning that subsequent tournaments have the possibility of containing a given individual from a previous tournament. Conversely, without replacement means that once an individual has taken part in a tournament it is no longer allowed to take part in subsequent ones.

is the population size, the probability of the fittest solution in the population being part of the tournament approaches 1.0 and so effectively selects the fittest solution on each tournament [82].

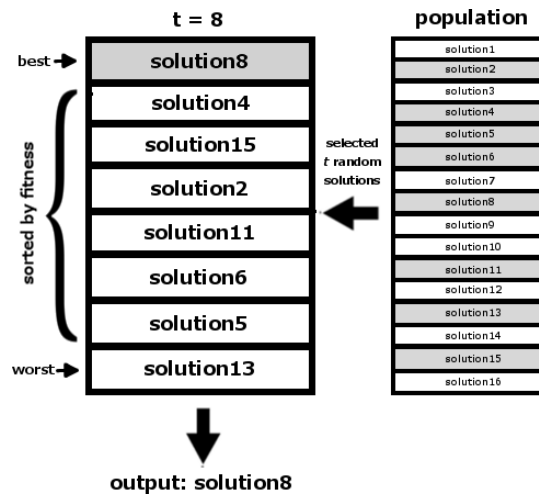


Figure 4.5: Tournament selection, where $t = 8$ (tournament size)

4.3.2.2 Crossover (Recombination)

Crossover is used to combine existing solutions into new, hopefully better, candidate solutions [29] and represents the main source of exploitation in the algorithm [104] - as such, they encourage the eventual convergence of the GA to a optimum solution. Many techniques exist for recombining parent solutions whose use is heavily reliant on the properties of the solution representation used; e.g., the specific encoding of the solution chromosomes and whether such solutions are of fixed or variable length. Several of the more common crossover operators will now be described.

(A) **UNIFORM CROSSOVER** Here, each parent solution is able to contribute genes to child solutions at the gene level rather than at the segment level - as is the case with one and two-point crossover [29]. For certain problems, this added flexibility can far outweigh any disadvantages stemming from the destruction or dismantling of whole building blocks. A probability is defined, often referred to as the mixing ratio (usually set to 0.5), which determines which parent will contribute each corresponding gene value to the child solution [29] - and by doing so determining an approximate ratio of genetic material taken from each parent. For example, supposing the crossover ratio of Fig. 4.6 below is 0.5, then it follows that roughly half of the genes from each parent will contribute to the resulting child solution.

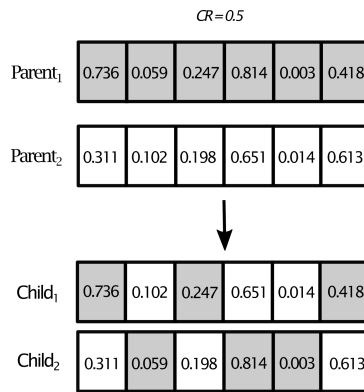


Figure 4.6: Uniform Crossover

(B) ONE-POINT CROSSOVER A single crossover point within a solution chromosome is selected at random [116]. When solutions in a GA are of a fixed length, a single crossover point common to both parents is used [65, 92]. Usually, two child solutions are produced from two parent solutions [116] however, one child can also be produced (see below). Given a fixed-length scheme, two variations of crossover can occur and are illustrated in Figs. 4.7a and 4.7b:

1. All genes after the crossover point in both parent solutions are swapped over between the two, resulting in two child solutions
2. All genes before the crossover point are copied from the first parent and all genes above the crossover points are copied from the second parent resulting in a single child solution

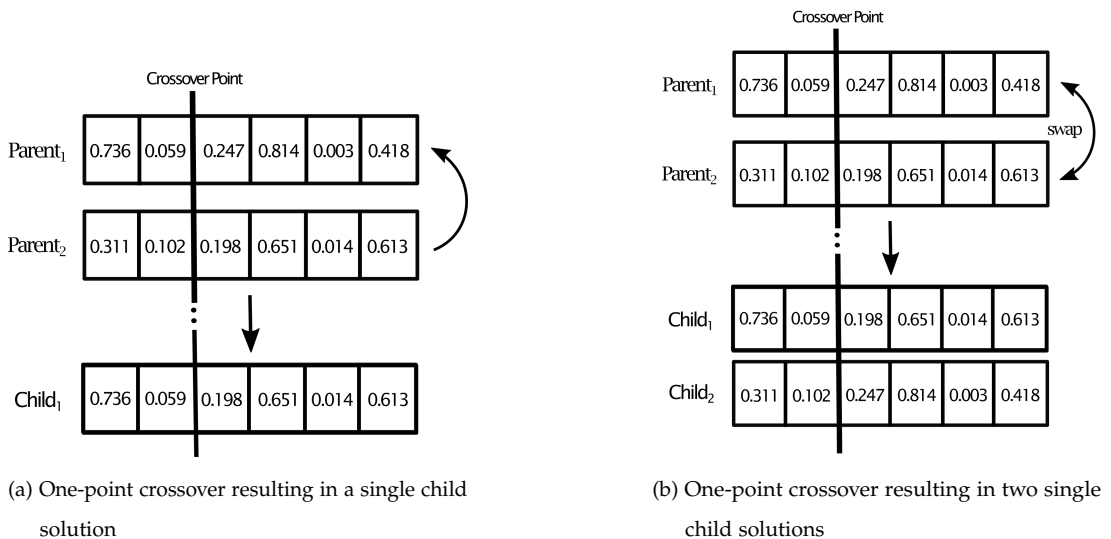


Figure 4.7: One-point Crossover

4.3.2.3 Mutation

Although the two genetic operators discussed so far serve an important purpose in genetic algorithms, an algorithm implemented with only these two operators would simply perform an organisation of the genes of various solutions in order to construct, hopefully, better solutions [82]. This can place an upper limit to the quality of the solutions produced, as only gene values available to the population on initialisation could be arranged to make solutions. The problem here is that the population lacks diversity, without this the population would simply converge upon a local optimum with no means of escape [82]. A solution, and the solution adopted in genetic algorithms, is to allow the application of a third operator, mutation, which also mimics the corresponding process found in nature. Mutations are often found to be unhelpful, such as, the mutations of one or more genes causing anything from colour blindness to Turner syndrome. The crux of the matter however is that without this process in our own genes, although able to mix between individuals, human beings would be unchanged from our first appearance (aside from cross breeding with other primate species) and so would not have been able to evolve into the species we are today. Mutation is utilised in genetic algorithms for the same reason; without mutation new gene values are not created and, as mentioned previously, genes are simply swapped between the populous to create new children - whose possible optimisation potential is highly constrained to the available genes [82].

Mutation makes use of user-defined parameter - the mutation rate P_m - that determines which gene(s) will undergo mutation. As far back as 1989, a parametric study was carried out suggesting that a rate calculated as: $P_m = \frac{1}{n}$ represents an optimal mutation rate for many cases [111, 26]. Although this finding is in reference to binary-coded GA, our own studies show that improved performance can be found using this scheme with real-valued GA. Typically, mutation is carried out over the entire solution string, i.e., every gene has a chance to be mutated, however schemes where only one gene is mutated each time also exist such as the one described in [26].

A few common mutation operators common to real-parameter genetic algorithms will now be discussed along with a short note about mutation operators for binary representations.

(A) UNIFORM MUTATION This mutation operator, generates a single child solution given a single parent solution. It selects a gene at random, replacing its value with one selected uniformly at random between the upper and lower bounds specified for that particular gene

[128]. More formally; given a solution vector x , the operator generates a child x' by selecting a gene $i \in \{1, \dots, n\}$ at random from the vector $x = \{x_1, \dots, x_i, \dots, x_n\}$ in order to generate a mutated vector $x' = \{x_1, \dots, x'_i, \dots, x_n\}$, where x'_i is a uniformly distributed random value selected from the range $(l(x_i), u(x_i))$ where $l(x_i)$ and $u(x_i)$ represent the lower and upper bound of x_i . The operator can play an important role in the beginning stages of the search process, where solutions should be allowed to explore more freely [91].

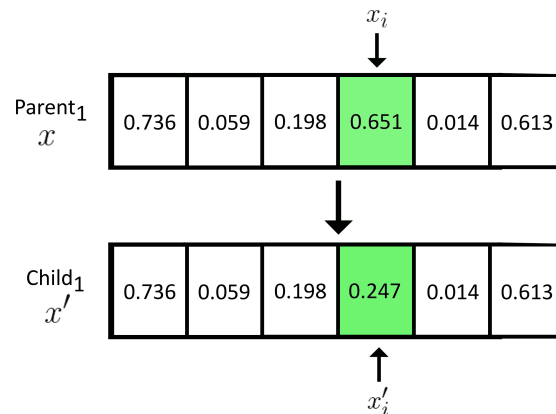


Figure 4.8: Uniform Mutation: where x is transformed to x'

(B) SINGLE GENE MUTATION In [26], the authors describe the use of a mutation operator for real-valued GA where a single gene in a solution is selected and mutated using polynomial mutation. In this thesis, we present - and in fact make use of - an alternative to this strategy where like the operator described above, the gene to be mutated is selected randomly (actually, uniform random) however, an additional random decision is made before this step that determines whether to mutate the solution at all. This additional random element of the operator is added by generating a uniform random number $r \in [0.0, 1.0]$ and comparing with the user supplied mutation rate P_m . If $r \leq P_m$, mutation is carried out on the solution as described previously and the solution is left in its un-mutated state otherwise. The actual mutation is carried out using uniform mutation. Over a number of trials with various mutation operators and replacement strategies, this operator coupled with worst-replacement (see Section 4.3.2.4) was found to be greatly superior to the other variants trialled over the majority of our selected benchmark suite (Chapter 3) and comparable over the remaining functions.

4.3.2.4 Replacement

The replacement process is somewhat similar to that of parent/survivor selection, where selected solutions are distinguished based on quality [29], but is instead concerned with

the selection of replacement slots for child solutions to enter into the population - and thus, is used at a different stage in evolution. Where they differ however, is that, in practice, selection tends to be stochastic whereas replacement is usually deterministic [29] - although, replacement schemes such as 'reverse tournament selection' and random replacement with an element of stochasticity show that replacement strategies can exist on a spectrum between deterministic and non-deterministic approaches. Replacement is a process that is mainly relevant in steady-state genetic algorithms (section 4.3.3) where the entire parent population is not being replaced by the child population, but some replacement strategies can be used with generational genetic algorithms such as: age-based replacement where $\mu = \lambda$, $(\mu + \lambda)$ selection and (μ, λ) selection. Given that the size of GA populations, steady state or otherwise, tend to remain constant, selections therefore have to be made as to which existing solutions will be discarded to make room for newly generated solutions. Replacement selection, as mentioned, is often based on fitness values - with bias often given to those of higher quality - however, the concept of solution age is also frequently used [29]. Therefore, replacement strategies can broadly be categorised as: (i) fitness-based replacement and (ii) age-based replacement [29]; it is worth noting however that some fitness-based replacement strategies can also take age into account as well. Since replacement selection is so similar to parent selection, in theory any of the selection operators discussed in section 4.3.2.1 can be used for replacement [29], however, there have been several special replacement operators developed and now commonly in use by practitioners [29]. The following paragraphs outline some of the more commonly used fitness-based and age-based replacement schemes.

(A) FITNESS-BASED REPLACEMENT STRATEGIES

(A.i) Replace Worst Originally used in the GENITOR steady-state GA [144], this replacement strategy selects the λ worst solutions in the population to be replaced [29] - where λ refers to the number of child solutions that have been produced. Removing the worst solutions from the population in this way will lead to fast improvements to the mean fitness of the population but can also quickly lead to premature convergence [29] as much needed diversity is purged from the population, reducing the takeover time of the best solution and significantly increasing innovation time. Therefore, this replacement strategy is mostly used in situations where the population is large and/or the detection and pruning of duplicate solutions is being carried out [29].

(A.ii) $(\mu + \lambda)$ Selection This replacement strategy comes from evolution strategies and in general refers to the merging of the parent population μ with the child population λ which

is then ranked according to fitness. The best μ individuals are then chosen to form the next generation [29].

(A.iii) (μ, λ) Selection Also originating from evolution strategies, this replacement strategy is made up of both an age component and a fitness component where typically, $\lambda > \mu$ children are generated from a set of μ parents [29]. The age component of the strategy means that the entire parent population will be deleted so that no individual lasts beyond a single generation [29]. The remaining λ population are then ranked according to their fitness and the best μ solutions go on to form the next generation [29]. Eiben and Smith note that in evolution strategies, this replacement operator is often preferred over $(\mu + \lambda)$ as: (i) since all parents are discarded in (μ, λ) it could allow escape from small local optima and can therefore be useful in multimodal search spaces containing many local optima [29], and (ii) if the fitness function is non-stationary i.e., it changes over time along with the location of the global optimum, $(\mu + \lambda)$ selection preserves those solutions relevant to an outdated optima and so is not able to track the new optima as well [29].

(B) AGE-BASED REPLACEMENT STRATEGIES The general concept of age-based replacement is that the fitnesses of solutions are not used in determining which solutions will be replaced and instead designed in such a way as to have each solution exist for the same number of algorithm iterations [29]. Eiben and Smith point out that since age-based replacement is not reliant on fitness values, there is the possibility that the mean or best fitness of a given iteration will be lower than the preceding iteration [29], but point out that, although counter-intuitive, this may actually be of some benefit should the population be centred near to a local optimum and if a decrease in the mean or best fitness between iterations does not happen too often [29].

(B.i) Random Replacement The naive Random replacement involves selecting, uniformly at random, a solution to be deleted from the population. At first glance it may not appear as if random replacement belongs to the class of age-based replacement strategies, or fitness based strategies for that matter, however, Eiben and Smith [29] suggest that this strategy has the same mean effect, meaning that on average individual solutions live for μ iterations - where μ is the size of the parent population - giving us some concept, however probabilistic, of age. Eiben and Smith further state that the random replacement strategy is not recommended since experiments by De Jong and Sarma in [24] between a steady-state GA using random replacement and a comparable generational GA showed that the steady state GA exhibited more variance in performance [29]. A reason for this was given by Smith and Vavak in [125]

stating that random replacement is far more likely to lose the best solution in the population than ‘delete-oldest’ replacement [29] - a comparable age-based replacement strategy.

(B.ii) Delete-oldest Replacement This operator represents the steady state GA equivalent of the simple genetic algorithm (SGA) age-based replacement strategy - where the number of offspring equals the number of parents, i.e., ($\mu = \lambda$) and so the child population entirely replaces the parent population so that no one individual survives for more than a single generation [29]. In a steady state GA where the child population is smaller than the parent population ($\lambda < \mu$) age-based replacement instead forms a first-in-first-out (FIFO) queue [29]. As such, the first λ solutions at the head of the queue will then be replaced on each generation.

4.3.3 *Steady State GA vs. Generational GA*

The first occurrence of a steady state genetic algorithm was that of the GENITOR algorithm [144] developed by Whitley and Kauth in 1988. The main difference between a generational GA (GGA) and a Steady State GA (SSGA) such as GENITOR is how the offspring in each iteration enter the population. In generational GA, new offspring are created from a subset of the individuals in a population through the mechanisms of crossover and mutation and form a new population which replaces the old population [140]. In SSGA, typically a single child solution is generated which replaces a single individual in the current population chosen by a replacement strategy. SSGA is computationally equivalent to GGA but since only one individual is replaced in each iteration, it typically takes more iterations (generations) for the algorithm to converge; however, the computational costs associated with genetic operations in SSGA are far lower and SSGA itself is much simplified over GGA [113].

In GGA, if the new population is to include some members of the previous population this is determined by a variable G often referred to as the ‘generation gap’ [24] as originally defined from a set of empirical studies of overlapping generations by De Jong in [25]. However, around the same time, Holland was performing analyses on two reproductive schemes: (i) R_1 , replaced a single individual in the current population at random with the child of a single parent using fitness-proportionate selection - a likely precursor to SSGA developed and popularised by Whitley in [144] - and (ii) R_d in which the expected number of offspring from parents selected deterministically were produced which replaced the entire current population [24].

Since a member of the population is to be replaced by a child solution on each iteration, another step is required in each iteration; namely, the selection of a solution to be replaced. No such step is required in GGA as the whole population is typically replaced, at least for

those implementations with no overlapping generations. As with other genetic operators there are many strategies for selecting a replacement slot. Common methods include: random replacement and worst replacement. These will be given more attention in section 4.3.2.4.

The SSGA model has been widely studied and applied in various contexts [29]. In [140] an SSGA and GGA were compared in terms of their ability to work in non-stationary environments, that is, their ability to adapt to the position or a dynamic optimum. An SSGA using a “delete the oldest” replacement strategy was shown to outperform an equivalent GGA in terms of both online and offline optimisation. The authors attribute the success of SSGAs for these applications to the fact that unlike GGAs, since an SSGA is incremental and a newly generated child solution is immediately considered as part of the mating pool from the next iteration, the algorithm is more able to more quickly move towards the optimum at a much earlier stage in optimisation [140].

4.3.4 *Steady State Genetic Algorithm: Implementation*

4.3.4.1 *Details and Description*

The genetic algorithm implementation chosen is that of an SSGA (Section 4.3.3) that replaces only a single solution in the population with a single child solution generated from the recombination of two parents in any given iteration. The selection of parents in this implementation is achieved by application of ‘tournament selection’ which are subsequently recombined using one-point crossover to form the child solution. Mutation is accomplished through the use of a single-gene uniform crossover (Chapter 4 - Section 4.3.2.3); as the solution representation used (as in section 4.1.2) bounds each gene by the interval [0.0, 1.0]. If mutation of a solution is to be performed, a randomly selected gene from a mutation operation is set to a uniform random value selected from this interval. This type of mutation was decided upon empirically through performance comparison with some more typical operators; such as those which are applied to the entire chromosome, and was found to perform best on average across the full benchmark suite ⁴.

A steady-state GA was selected due to its reduced computational complexity - both in terms of runtime and memory. It is well known that one of the most expensive steps in generational GA is that of population generation; *therefore*, if generational GA was tuned

⁴ The algorithm implementation was identical for each variant; the only variable being modified in the experiment was the choice of mutation operator

using our parameter tuning method (Chapter 7), the generation of parameter configuration sets at high dimensionalities by SMAC would likely take far longer than it did for an SSGA.

4.3.4.2 Pseudocode

Algorithm 8 Pseudocode of SSGA Implementation

```
1: procedure GA(pSize, mutationRate, tSizeProportion)
2:   popSize  $\leftarrow$  pSize
3:   mutation  $\leftarrow$  mutationRate
4:   population  $\leftarrow$  INITPOPULATION(popSize)
5:   tSizeProp  $\leftarrow$  tSizeProportion    ▷ Proportion of population to use in tournaments
6:   tSize  $\leftarrow$  2 + FLOOR(popSize * tSizeProp)
7:   Sbest  $\leftarrow$  GETBEST(population)
8:   while terminate == false do                                ▷ Main GA loop procedure
9:     EVALUATEPOPULATION(population)
10:    parent1  $\leftarrow$  TOURNAMENTSELECT(tSize)
11:    parent2  $\leftarrow$  TOURNAMENTSELECT(tSize)
12:    child  $\leftarrow$  ONEPOINTCROSSOVER(parent1, parent2)
13:    child  $\leftarrow$  UNIFORMMUTATION(child, mutationRate)

14:    replace  $\leftarrow$  GETWORST(pop)
15:    population[replace]  $\leftarrow$  child

16:    Sbest  $\leftarrow$  GETBEST(population)
17:  end while
18:  return Sbest                                                ▷ return best solution discovered
19: end procedure
```

```

20: procedure TOURNAMENTSELECT(tournamentSize)
21:   tournamentSet  $\leftarrow$  random set of tournamentSize solutions from pop
22:   best  $\leftarrow$  tournamentSet[0]
23:   for (i=1 to tournamentSize-1 by 1) do
24:     if (tournamentSet[i] < best) then
25:       best  $\leftarrow$  tournamentSet[i]
26:     end if
27:   end for
28:   return best
29: end procedure

30: procedure ONEPOINTCROSSOVER(parent1, parent2)
31:   r  $\leftarrow$  uniform random float  $\in$  [0.0, 1.0]
32:   chromosomeIndex  $\leftarrow$   $\lfloor r * \text{parent1.length} \rfloor$ 
33:   for (i = 0 to parent1.length-1 by 1) do
34:     if (i < chromosomeIndex) then
35:       child[i]  $\leftarrow$  parent1[i]
36:     else
37:       child[i]  $\leftarrow$  parent2[i]
38:     end if
39:   end for
40:   return child
41: end procedure

42: procedure UNIFORMMUTATION(child, mutationRate)
43:   r  $\leftarrow$  uniform random float  $\in$  [0.0, 1.0]
44:   if (r  $\leq$  mutationRate) then
45:     r  $\leftarrow$  next uniform random float  $\in$  [0.0, 1.0]
46:     geneIndex  $\leftarrow$   $\lfloor r * \text{child.length} \rfloor$ 
47:     child [geneIndex]  $\leftarrow$  r
48:   end if
49:   return child
50: end procedure

```

```

51: procedure GETWORST
52:   worst ← 0
53:   for (i = 1 to popSize by 1) do
54:     if (population[i].score > population[worst].score) then
55:       worst ← i
56:     end if
57:   end for
58:   return worst
59: end procedure

60: procedure GETBEST
61:   best ← 0
62:   for (i = 1 to popSize by 1) do
63:     if (population[i].score < population[worst].score) then
64:       best ← i
65:     end if
66:   end for
67:   return best
68: end procedure
69: procedure EVALUATEPOPULATION(population)
70: end procedure

```

4.3.4.3 Tunable Parameters

The parameters used to control the performance of the algorithm in this implementation are as follows:

- Population Size
 - Type:- integer $\in [2, 250]$
 - Used by the initialization procedure to define the size of the population that will remain fixed throughout execution of the algorithm. A lower bound for this parameter of 2 is the minimum population size that would allow the genetic algorithm to be referred to as such; specifically, a population size of only one would essentially represent regular stochastic hill climbing (as crossover would not be possible). The upper bound was set as to fall in line with some of the other population based approaches represented here, as there seemed to be no relatively

general agreement to a reasonable range for population size in genetic algorithms. The default value for this parameter as used initially in SMAC is 100.

- Tournament Size Percent

- Type:- float $\in [0.01, 0.7]$

- Defines the size of the tournament groups used for parent selection as a proportion of the main population i.e., how many randomly selected solutions from the population are chosen to take part in the tournament. This is equivalent to setting an explicit number of individuals from the population - as is more typical - but provides prevention against the tournament size becoming larger than the main population when SMAC is actively tuning both tournament size and population size independently. In this implementation of the operator, n tournaments are run in order to select n parent solutions; thus, two tournaments are carried out to produce two parents - one from each tournament. This parameter also allows for the adjustment of selection pressure; a large group size will reduce the chances of weaker solutions being selected for recombination with the inverse being true of a smaller group.

The upper bound on the value, of 0.7, was selected arbitrarily for the purposes of the automatic parameter tuning performed by SMAC. This value was judged to be a reasonable upper limit in order to prevent premature convergence resulting from the repeated selection of the very best solutions. Additionally, should an upper bound of 1.0 have been allowed, any parameter tuning process that produced a value for this parameter that lay very close to this bound would have caused the tournament selection procedure to be roughly equivalent (depending on the actual proportion of the population involved) to selecting the best solution in the entire population. Since our tournament selection creates tournaments with replacement and separate tournaments are carried out to select each parent, not only will the selection result in the same solution being selected far more often, but also that - without the action of mutation - the entire population would be quickly replaced with exact copies of a single 'super-individual'. Mutation would indeed continue to be applied however, becoming the primary means of search progress, making the GA process equivalent to a random search. It is clear that avoiding such problems makes comparison of this SSGA implementation with other approaches used in this study more straightforward and meaningful.

- Mutation Rate

- Type:- float $\in [0.01, 1.0]$
- Defines the rate, and therefore the probability, of a randomly selected gene being mutated or not. In other words, if a uniform random number generated for each mutation procedure is less than the mutation rate, then the gene is mutated. Given our implementation makes use of 'single gene' uniform mutation operator (Chapter 4 - Section 4.3.2.3); the effect on mutation is different than for other more common mutation schemes. Usually, lower values of this parameter results in less frequent mutation operations and in turn, results in slower overall convergence towards the global optimum with more risk of premature convergence to local optima - due to imbalance with crossover operations. Higher values on the other hand will result in faster convergence but will invariably produce too many perturbations to high quality solutions to progress further towards optimality. For this operator, better performance was found over more common operators where the rate of mutation is much higher - typically closer to 1.0. I hypothesise that the operator attempts to behave more like other typical operators - where the mutation rate give each gene in a solution the chance to mutate - but there appears to be some benefit to having no mutation at all. There is credence to this hypothesis in that early Evolution Strategy (EA) implementations made use of only a mutation operator - where crossover was omitted [89]. A mutation only GA variant was used successfully in [72] and was selected as it was found to be fast and reliable; however this was in the context of binary search spaces. We intend to investigate this further in terms of continuous search spaces in future work.

4.4 PARTICLE SWARM OPTIMISATION (PSO)

Particle swarm optimisation (PSO) is a population based metaheuristic introduced by Kennedy and Eberhart in [67] and is heavily influence by the computational models and rules for coordination developed at the time that simulated the flocking behaviour of birds and the schooling of fish [67]. The PSO algorithm maintains a population of solutions (particles), known as the swarm, each representing a point in n-dimensional space [7]. Initialising at random locations in the search space, the particles traverse the landscape in order to discover the global optimum of the given objective function [105]. This is in a way similar to a swarm of organisms in nature foraging for food or a new nesting location, the particle swarm in PSO makes use of the objective function - akin to biological feedback from an environment - to determine the quality or the amount of food at a given location; the swarm searches the space for the location containing the best solution - or most amount of food [105].

In models of flocking birds, the movement of each individual is determined both by (i) the influence of it's own perception of the environment and satisfaction of its goals and (ii) the influence of social pressures such as: separation from other individuals in the swarm, alignment to the average velocity of the flock and cohesion where the individual moves towards the average position of its neighbours. Similarly, in PSO the movement of a particle depends on its own velocity, good solutions or locations previous discovered by the particle and the good solutions discovered by its neighbouring particles [33]. This particle behaviour is modelled in PSO by having each particle maintain both the location of the best solution discovered so far by the particle itself - along with a corresponding objective value - and the location of the best solution discovered by all particles in its immediate neighbourhood [105].

4.4.1 *Neighbourhood Topologies*

The choice of neighbourhood topology used by a PSO algorithm can play a significant role in its overall behaviour. Investigations, such as in [66] by Kennedy, show that a particles neighbourhood topology interacts directly with the objective function of the problem being optimised, and further, that it significantly affects the performance of the swarm where the particular effect it has is dependent on the problem [87]. Therefore, where some topologies are well suited to certain functions, other functions can present problems to these same topologies [66, 87].

Several topologies common to PSO implementations include [87]:

1. Ring Topology

- Also referred to as *lbest* (or local best) topology, here each particle in the swarm is only affected by the best performance of its k immediately adjacent neighbours. A common case is to have $lbest = 2$.

2. Star Topology

- In this topology, social influence passes through only one central individual/hub, commonly selected at random, and is itself influenced by the rest of the population [68, 87]. On each iteration, every particle in the swarm moves towards only the central particle, and the central particle directs itself towards the best performing particle in the neighbourhood. This topology can be used to control the possibility of premature convergence, as since the population remains isolated from each other directly, the central particle - constantly adjusting its trajectory towards the best particle - effectively slows the transmission of good solutions throughout the population and thus maintaining diversity.

3. Von Neumann / Toroidal Topology

- Here, each particle is connected to several others, specifically to the particles immediately above, below and on each side of it and any edges are wrapped [68].

4. Fully Connected Topology

- Also referred to as a *gbest* (or global best) topology or *full* topology, this neighbourhood scheme prescribes the sharing of information amongst all particles in a population; that is, each particle is directly connected to one another. Here, all particles in the population are influenced and directed by the best solution found by any particle in the whole population. Kennedy and Mendes in [68] observed that fully connected (*gbest*) populations had a tendency of converging more quickly than ring (*lbest*) populations and further, are more prone to local optima convergence. This however does not prevent it from being one of the most widely used topologies [87].

In the PSO algorithm used for this thesis, a fully connected topology is implemented since: (i) no additional computation or neighbourhood link management needs to be performed to make use of the topology, which would otherwise extend the required runtime of our already very expensive parameter tuning processes, and (ii) it allows the PSO implementation

to more closely resemble the neighbourhood structures of the other population-based metaheuristics being investigated; where operators typically choose members of the population at random (i.e., fully connected). This allows for a fairer comparison to be conducted between these approaches based on the merits of algorithmic approach rather than on the choice of neighbourhood structure.

In the PSO literature, such as in [48], it is usual to find the term *informants* used in place of ‘neighbour’ as used throughout the rest of the thesis as it more clearly represents the role of individuals and their neighbours ‘informing’ the direction of search. However, for consistency with the other metaheuristics described, the terms neighbour and neighbourhood will be used.

4.4.2 Parameters of PSO

4.4.2.1 Particle Position, Velocity and Related Parameters

Every particle in the swarm maintains for each dimension of the problem both its position within the search space and its velocity [105]. Velocity is added to the current position of the particle to move it to a new position of the search space - representing a particle step [105]. One of the first identified issues related to PSO was that of swarm explosion, where the swarm would gradually begin to drift apart and more regularly violate bounds on the search space [13]. This goes against the desire for a convergent algorithm. The solution proposed by researchers was to ‘clamp’ the velocity to a lower and upper bound specified by a range $[-V_{\max}, \dots, V_{\max}]$ [105] where V_{\max} is the maximum velocity threshold. By clamping velocity, uncontrolled increases in magnitude from a particles current location can be avoided. Whenever the velocity of a particle in any dimension exceeds $\pm V_{\max}$ it is set to the value of the violated bound as in[30]:

$$v_{id}(t+1) = \begin{cases} V_{\max} & \text{if } v_{id}(t+1) > V_{\max} \\ -V_{\max} & \text{if } v_{id}(t+1) < -V_{\max} \end{cases} \quad (4.5)$$

As with all tunable parameters in metaheuristics, and as per NFL (Chapter 2 Section 2.4.5.2), there is not a single optimum value of V_{\max} to maximise performance in all conceivable problem space and so time should be taken to find a suitable value for each problem space being addressed [105].

The velocity is updated on each iteration of the algorithm based on the three components: (i) inertia weight ω , (ii) the cognitive component c_1 and (iii) the social component c_2 [10].

The inertia weight is used to specify how the current velocity influences the new velocity and also controls the scope of the search [10]. A higher value of ω favours exploration and lower values favour exploitation [10]; it is important to determine an appropriate balance, as with all metaheuristics, between these two actions and so an inertia weight should be carefully selected to avoid either slow convergence or premature convergence. In a similar way, the cognitive and social parameters specify how much social and cognitive influence contributes to the new velocity [30]. When $c_1 < c_2$ the search is more biased towards global influences and biased towards personal influence when $c_1 > c_2$ [30]. Also, and as with ω , higher and lower values of c_1 and c_2 favour more exploration and exploitation respectively [30] and is determined by the velocity update formula [10]:

$$v_{id}(t+1) = \omega \times v_{id}(t) + c_1 \times r_1 \times (p_{id} - x_{id}) + c_2 \times r_2 \times (p_{gd} - x_{id}) \quad (4.6)$$

where v_{id} is the velocity of a particle at iteration t of the algorithm, c_1 and c_2 represent the cognitive and social components respectively and ω is the inertia weight. Two random parameters r_1 and r_2 are also used to randomly modify the influence of the social and cognitive parameters have on any given iteration [30, 10]. An illustration of a single velocity and position update is shown in Fig. 4.9.

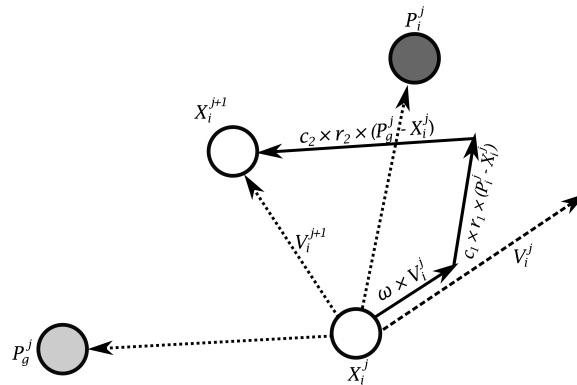


Figure 4.9: Velocity and Position Update in PSO

4.4.2.2 Inertia Weight

The motivation for introducing an inertia weight ω into the velocity equation was to gain more control over exploration vs. exploitation which would eliminate the need for the V_{max} parameter [7]. The concept of inertia weight was first proposed by Shi and Eberhart in [119] when they introduced a fixed inertia weight [7]. They stated that inertia weight plays a role in balancing global vs. local search (exploration vs. exploitation); a low value of inertia weight encourages a more local search of the available space where larger values would facilitate a more global search [119, 7]. As well as simply selecting a fixed inertia weight for a given

problem, as above, many other methods that adaptively control the value of ω have been proposed in the literature; some of which will be discussed here.

A simple adaptive method is to linearly decrease value of ω so as to shift the mode of the search from being more exploratory early on in the search process to more exploitative in later stages [119, 120]. Xin et al. in [148] calculate the value of ω as follows [148]:

$$\omega = (\omega_s - \omega_e)(t_{\max} - t)/t_{\max} + \omega_e \quad (4.7)$$

where t_{\max} represents the total number of algorithm iterations allowed, t is the current iteration count and ω_s and ω_e represent the initial and final values for the inertia weight respectively [148] - bounding the weight to the range $[\omega_e, \omega_s]$. A weight decreasing linearly from $\omega_s = 0.9$ to $\omega_e = 0.4$ is shown to improve performance [120, 148].

Eberhart and Shi proposed a random inertia in [28], where a value of ω is selected randomly from the range $[0.5, 1.0]$ and by the equation [28]:

$$\omega = 0.5 + \frac{\text{rand}[0, 1]}{2} \quad (4.8)$$

where $\text{rand}[0, 1]$ is a uniformly distributed random number between 0 and 1. Thus, the mean weight will be ≈ 0.75 . It was used to address the issue with using linear decreasing inertia weight when tracking non-linear dynamic systems, where it is not predictable whether more exploration or more exploitation would be of benefit at any given time during execution [28]. They report that convergence rate is increased in early the stages of execution [28], although, it is not reported whether this strategy is of any use for static systems.

Based on linear decreasing inertial weight and random inertia weight, Feng et al. presented two strategies that used a chaotic mapping to set the inertia weight referred to as the Chaotic Descending Inertia Weight (CDIW) and the Chaotic Random Inertia Weight (CRIW) [31]. They report "outstanding" performance of CRIW in comparison to random inertial weight due to the alternation between rough search and "minute" search throughout the evolutionary process [31].

Al-Hassan et al. in [2] developed an optimised PSO algorithm known as PSOSA that makes use of simulated annealing for optimising the inertia weight and was tested against an urban planning problem [7].

4.4.2.3 *The Cognitive and Social Parameters*

The cognitive and social parameters also referred to as the acceleration coefficients, are represented by c_1 and c_2 respectively and as mentioned earlier, each determines how much either cognitive (personal performance) influence or social (global performance) influence contributes to the velocity update equation. It is often the case that these parameters are fixed throughout execution of the algorithm with typical, and widely used, values of $c_1 = c_2 = 2$ [132]. However, a couple of adaptive variants have been identified.

Ratnaweera et al. in [110] presented a time-varying acceleration coefficient scheme (PSO-TVAC) with the goal of enhancing global search at the beginning of the search process and improving convergence to the global optima near the end of the search by adjusting the values of c_1 and c_2 over time - specifically, by reducing the cognitive component and increasing the social component [110]. With a large value of the cognitive component and small value of the social component at the beginning of the search, the search will favour global exploration of the available space. Conversely, a small value of the cognitive component and a large value of the social component will better facilitate local exploitation of the global best near the end of the search [110].

4.4.3 *Boundary-Handling in PSO*

Due to the nature of the displacement update of individual particles in PSO, there is always the possibility that particles can overshoot any bounds placed on the search space and end up in infeasible regions if not adequately caught and remedied. Left unchecked, individual particles and indeed most of the population (the global search itself) can become lost outside the boundaries resulting either with: (i) infeasible solutions to the given problem, if the search converges on a solution in the infeasible regions and violating particles are allowed to be evaluated, or (ii) potential for sub-optimal solutions to be returned, if the search eventually returns to converge on a solution within bounds or if out-of-bound solutions are not evaluated. Also, real-world applications often search in spaces that are high-dimensional and bounded in order to retain the semantics of the parameters [19]. Therefore, detecting and repairing boundary violations during the search process becomes an essential consideration. Additional to this, the method chosen to deal with boundary violations has also been shown to be crucial to the performance of PSO - where some methods may “paralyze” a PSO searching in high-dimensional and complex search spaces [19].

Despite this being an important consideration for constrained optimisation problems, particularly in the case of real-world high-dimensional spaces, how a PSO should handle boundary violations has been given little mention in the literature. A mathematical proof, and empirical results using common velocity initialisation techniques, by Helwig and Wanka [52] showed that for a PSO algorithm in high-dimensional spaces:

“Uniform velocity initialization causes all particles to leave the search space with overwhelming probability (w.o.p.) with respect to the search space dimensionality n ... Examples demonstrated that this probability rapidly approaches 1 if the search space dimensionality is increased. [52]”

It is further noted that none of the initialisation approaches used could fully avoid this eventuality, the results suggest that half-diff velocity initialisation has the fewest drawbacks - thus demonstrating the importance of boundary handling in influencing search performance [52].

Chu et al. later provided further empirical results to strengthen this proof in [19] where the average fraction of particles “flying outside” the boundary in the first iteration, from 100 independent runs of PSO against two benchmark functions at increasing dimensionality, were recorded. The results showed a clear trend that this fraction increased significantly as the dimensionality of each benchmark function increased. Also in this same study, Chu et al. compared the ability of three basic boundary-handling schemes: random, reflecting and absorbing, in handling and preventing boundary violations and how this translates to algorithm performance over the two original benchmark functions. These results showed that both the random and absorbing schemes had the ability to hamper the searches performed with PSO when compared with the reflecting approach.

4.4.3.1 *Boundary-Handling Approaches*

This section covers several common boundary-handling schemes in the literature, including those studied by Chu et al. in [19].

(I) **RANDOM** In this scheme, if a particle violates the boundary of one of its parameters, a value chosen randomly from a uniform distribution, between the lower and upper bounds of the offending parameter, is assigned as the new value [19]. That is, if X_c represents the the location of the current infeasible particle variable and LB and UB denote the lower and upper bounds of the variable respectively, then a new feasible location can be chosen as $X_c \in [LB, UB]$ [101].

(II) **REFLECTING** In this boundary-handling scheme, any violation to a parameter boundary is ‘reflected’ back inside the feasible region as if having ricocheted off a hard surface [149]. This is done by relocating the particle at the violated boundary and the sign of the velocity component in the offending dimension is flipped [149]

(III) **DAMPING** The scheme works similarly to reflection in that the particle rebounds back into the feasible region; however, a randomly generated damping factor from the interval $[0, 1]$ is used to reduce the velocity component in that dimension which, as with reflecting, has its sign changed [149]. Also as with the reflecting scheme, this reflection is carried out after having been absorbed to the violated boundary [149].

(IV) **ABSORBING** Here, a particle violating a parameter boundary is set to the boundary value it violated. For example, if a particle would pass across the upper bound for a given parameter, it would instead be placed at the upper boundary and its velocity at that dimension zeroed [149]. The name of the approach stems from the metaphor of a particle leaving a parameter boundary and being ‘absorbed’ back. As before, if we consider X_c to be the location of the current infeasible particle variable, then X_c is set to the violated boundary as follows [101]:

$$X_c = \begin{cases} LB & \text{if } X_c < LB \\ UB & \text{if } X_c > UB \end{cases} \quad (4.9)$$

Even though the velocity is zeroed, the particles memory of personal best and global best locations will reorient the particle inside the feasible search space [149].

4.4.4 Particle Swarm Optimisation (PSO): Implementation

4.4.4.1 Details and Description

This implementation was derived from code sourced from [85] and [86] by Gandhi Manalu and was modified in order to work within our experimental framework. The main omission from this code however, which had to be implemented, was that of handling boundary violations (see Section 4.4.3) so particles were simply allowed to move beyond the bounds on the search space and therefore assuming that enough particles would remain within bounds in order to find the global optimum. This approach is suitable if it is known that there are no points or regions in the infeasible space with a better objective value than that of the feasible global optima - as is the case with the example benchmark function provided with the source code

cited above - as any bound violating particles would eventually, and likely rather quickly, be drawn back within bounds through the effect of the cognitive and social components of the velocity formula (See section 4.4.2.1). Otherwise, any infeasible solution measuring better performance than the true global optima will likely cause the search to converge somewhere in infeasible space.

Since we do not know the nature of the search spaces beyond the defined bounds on our benchmark function set, it was necessary to use a boundary handling strategy to avoid any possibility of objective performances being reported from infeasible regions. For this we opted to use a simple absorbing scheme (Section 4.4.3) - one in which the velocity of the violating parameter is also set to zero. During our studies, we found that this scheme outperformed other variants tested, namely: (i) absorbing and truncating the velocity to 30% of its original, (ii) Random replacement within the bounds and a random update to the velocity and (iii) a dampened reflect scheme - where the solution was reflected back into the bounds at 20% of the distance violated and the velocity was set to zero. A 'full' reflecting scheme was tried initially - where 100% of the violated distance was used and velocity zeroed - however this lead to particles which were reflected outside the opposing boundary. Although we could have dealt with this with multiple reflects (within a while loop) until the particle lay within the bounds, it was felt that since there was no mention found in the literature relating to this issue with the reflecting scheme that it would not be conducive to a 'common' implementation of the algorithm.

4.4.4.2 Pseudocode

Algorithm 9 Pseudocode of PSO Implementation

```

1: procedure PSO(iWeight, cogCoeff, socialCoeff, probSize, popSize)           ▷ Main loop
2:   problemSize  $\leftarrow$  probSize
3:   swarm  $\leftarrow$   $\emptyset$ 
4:   swarmSize  $\leftarrow$  popSize
5:   w  $\leftarrow$  iWeight                                                       ▷ Inertia Weight
6:   c1  $\leftarrow$  cogCoeff                                                     ▷ Cognitive Acceleration Coefficient
7:   c2  $\leftarrow$  socialCoeff                                                 ▷ Social Acceleration Coefficient
8:   Vmax  $\leftarrow$  maximum allowed value
9:   fitnessValueList  $\leftarrow$  INSTANTIATE(swarmSize) ▷ Current list of all particle fitnesses
10:  pBestLocations  $\leftarrow$  INSTANTIATE(swarmSize) ▷ List of particles best locations so far
11:  pBestScores  $\leftarrow$   $\emptyset$                                              ▷ Best objective scores of all particles so far
12:  gBestScore  $\leftarrow$   $\emptyset$                                              ▷ Objective score of global best
13:  gBestLocation  $\leftarrow$   $\emptyset$                                          ▷ Location of global best

14:  INITIALISESWARM
15:  UPDATEFITNESSLIST

16:  while (terminate == false) do                                         ▷ Main Loop
17:    for (i = 0 to swarmSize - 1 by 1) do                                  ▷ Update pBest Lists
18:      pBestScores[i]  $\leftarrow$  fitnessValueList[i]
19:      pBestLocations[i]  $\leftarrow$  swarm[i].location
20:    end for

21:    bestParticleIndex  $\leftarrow$  GETMININDEX(fitnessValueList)
22:    if (fitnessValueList[bestParticleIndex] < gBestScore) then       ▷ Update
    gBestScore
23:      gBestScore  $\leftarrow$  fitnessValueList[bestParticleIndex]
24:      gBestLocation = swarm[bestParticleIndex].location
25:    end if
26:    for (i = 0 to swarmSize - 1 by 1) do
27:      r1  $\leftarrow$  randomfloat  $\in$  [0.0, 1.0]
28:      r2  $\leftarrow$  randomfloat  $\in$  [0.0, 1.0]
29:      particle  $\leftarrow$  swarm[i]

```

```

30:     newVel ← INITIALISE()           ▷ New list of velocities for the particle
31:     newLoc ← INITIALISE()         ▷ New location of the particle
32:     for (j = 0 to problemSize by 1) do
33:         /*Update the velocity and position of the current particle */
34:         newVel[j] ← (w * particle.velocity[j]) + (r1) * c1 *
(pBestLocations[i].location[j] – particle.location[j]) + (r2 * c2) * (gBestLocation –
particle.location[j])
35:         newLoc[j] ← particle.location[j] + newVel[j]
36:         /*Boundary Handling: Absorbing and Zero Velocity */
37:         if (newLoc[j] < LB) then
38:             newLoc[j] ← LB
39:             newVel[j] ← 0.0
40:         else
41:             if (newLoc[j] > UB) then
42:                 newLoc[j] ← UB
43:                 newVel[j] ← 0.0
44:             end if
45:         end if
46:     end for
47:     particle.velocity ← newVel
48:     particle.location ← newLoc
49: end for
50:     UPDATEFITNESSLIST
51: end while
52:     return gBestLocation and gBestScore
53: end procedure

54: procedure INITIALISESWARM
55:     for (i = 0 in swarmSize – 1 by 1) do
56:         newParticle ← INITIALISE
57:         /* Randomise the location in the space defined by the representation */
58:         newLoc ← INSTANTIATE()
59:         newVel ← INSTANTIATE()
60:         for (j = 0 in problemSize – 1) do
61:             newLoc[j] ← uniform random float ∈ [LB, UB]
62:             newVel[j] ← uniform random float ∈ [–Vmax, Vmax]
63:         end for

```

```

64:     newParticle.location ← newLoc
65:     newParticle.velocity ← newVel
66:     swarm[i] ← newParticle
67:   end for
68: end procedure

69: procedure UPDATEFITNESSLIST
70:   for i = 0 in swarmSize - 1 by 1 do
71:     fitnessValueList[i] ← SCORESOLUTION(swarm[i])
72:   end for
73: end procedure

74: procedure GETMININDEX(fitnessList)
75:   index ← 0
76:   minValue ← fitnessList[0]

77:   for (i = 0 in fitnessList.length by 1) do
78:     if (fitnessList[i] < minValue) then           ▷ For a minimisation problem
79:       index ← i
80:       minValue ← fitnessList[i]
81:     end if
82:   end for
83:   return index
84: end procedure

```

4.4.4.3 Tunable Parameters

The parameters used to control the performance of the algorithm in this implementation are the following:

- Population Size
 - Type:- integer $\in [2, 250]$
 - Used by the initialization procedure to define the size of the population that will remain fixed throughout execution of the algorithm. The default population size initially trialled by SMAC is set to 50
- Inertia Weight

- Type:- float $\in [0.4, 1.0]$
 - The inertia weight is used to scale the effect of the previous velocity of a particle in a given dimension. A lower and upper bound of 0.5 and 1.0 respectively was suggested by Eberhart and Shi in [28] and so we allowed SMAC to explore a little below this value. Also as per [28] the default value used by SMAC in the tuning process was 0.5.
- Cognitive and Social Acceleration Coefficients (c_1 and c_2)
 - Type:- float $\in [1.0, 2.0]$
 - These parameters are responsible for scaling the influence of the location of a particles 'best ever' solution and the location of the global best solution in the velocity update formula respectively. The upper bound for these parameters was determined by Clerc [20], where $c_1 = c_2 = 2.05$ and rounded to 2. Suganathan in [132] also used a fixed value of 2 for both coefficients. The default value is also 2.0.

4.5 SIMULATED ANNEALING (SA)

Simulated Annealing (SA), developed by Kirkpatrick et al., is a general purpose algorithm for the global optimisation of continuous objective functions[27]⁵ which takes its inspiration from the process of metallurgical annealing, and is implemented as an extension of the Metropolis-Hastings algorithm - a Monte-Carlo algorithm originally used to draw sample states from the probability distributions of thermodynamic systems. SA is essentially a descent method for global optimisation - modified by random ascent moves [27] - in which the search for solutions mimics the controlled heating and cooling of annealing in order to find the global optima of an objective function. The slow cooling, often referred to as the cooling schedule, is used to decrease the probability of accepting worse solutions nearer the end of the search process, and conversely, allowing for greater exploration of the solution space in the early stages. SA has been shown to produce solutions that lie near to the global optimum in a computation time within a polynomial upper bound; these findings have further been shown to be independent of initial conditions [27].

The basic idea behind SA, is to have the heuristic generate, uniformly at random, a move from a set of neighbouring states N and probabilistically decide on whether or not to accept the new move - determined by the acceptance probability function. By making use of the temperature, updated by the cooling schedule as the search progresses, the acceptance probability function will accept more worsening moves at the beginning of the search - maximising exploration of the space - and gradually reducing to stochastic descent as the temperature approaches zero. Therefore, SA is composed of two stochastic processes: one for the generation of candidate moves and another of the acceptance of candidate moves [27].

In the next few sections, the underlying concepts behind simulated annealing will be discussed. Firstly, we gain some intuition into annealing within the field of condensed matter physics and how this motivated the development of the Metropolis-Hastings algorithm - the backbone of the complete SA algorithm. Next, SA itself is described in comparison to physical annealing and how the Metropolis-Hastings algorithm is used in SA to imitate annealing for solving optimisation problems.

⁵ It is also a popular monte-carlo approach for any kind of optimisation including discrete optimisation such as combinatorial optimisation [27]

4.5.1 *Annealing in Condensed Matter Physics*

Although largely beyond the scope of this section, a short description of the physical annealing process will aid further discussion and understanding of the simulated annealing algorithm. In condensed matter physics, annealing is a thermal process for allowing solids in a heat bath to obtain low energy states [1]. This process can be stated in the following two steps [1]:

1. Melt the solid in a heat bath by increasing the temperature to a maximum value
2. Slowly and carefully decrease the temperature of the heat bath until the particles of the solid rearrange themselves at it's ground state

As described by Aarts et al. in [1], in the liquid phase of the material, that is, when the atoms are in an excited state, the atoms arrange themselves randomly. Conversely, at the ground state of the material the atoms form in a highly structured lattice in which the energy value is minimal [1]. The authors continue to explain that the ground state of the material can only be reached if the maximum temperature is high enough and the cooling sufficiently slow.

4.5.2 *Acceptance Probability*

When the neighbouring solution selected is a worsening move, the probability P of accepting the solution is calculated from a function of the difference between the solution scores ΔE (between the current solution and the neighbour solution) and the current value of the temperature parameter T , as in [82]:

$$P = e^{\frac{-\Delta E}{T}} \quad (4.10)$$

The difference ΔE is calculated as the fitness of the neighbour score minus the fitness of the current solution score [82]:

$$\Delta E = f(x') - f(x) \quad (4.11)$$

From equation 4.10, the probability that a worsening move is accepted is larger at a higher value of T and smaller at lower values of T . At very low values of T , the simulated annealing algorithm behaves more like a hill climbing algorithm [117] - accepting only those moves which improve on the current solution. Also, as the value of ΔE increases, i.e., the difference between the current and neighbouring solution scores, the probability of acceptance at a given T decreases exponentially [117].

4.5.3 Cooling Schedules

The cooling schedule, also known as the annealing schedule, in SA is used to determine a finite sequence of value for T over time as the algorithm progresses and a finite number of steps taken at each value of T [1]. A cooling schedule is fully defined by [1]:

- An *initial value* of T
- A *minimum value* of T
- A *decrement function* used to reduce the value of T over time
- The *number of steps* to be taken at each temperature

There are many cooling strategies proposed in the current SA literature however a common method, first proposed by Kirkpatrick in [69] is simply to decrease the temperature geometrically by a fixed constant $\alpha \in [0.0, \dots, 1.0]$ [1] and is thus known as the geometric cooling schedule. This scheme falls under the category of static cooling schedules [1], where the rate of cooling is not influenced by the state of the search; thus, the strategy remains fixed throughout the search process and must be specified before execution of the algorithm [53]. Conversely, adaptive (or dynamic) cooling schedules adjust the rate of temperature decrement from information obtained during execution of the algorithm [53].

In [130], Strenski and Kirkpatrick attempted to categorise optimal cooling schedules by studying SA over a set of very small problem instances over a finite number of iterations using three common schedules: geometric, linear and inverse-logarithmic [130, 53]. Their results suggested that optimal schedules are not those which decrease monotonically [130, 53] - that is, they do not only decrease with time but may remain stable or increase. The experiments also show that there is very little difference, in terms of performance, between linear and geometric schedules [130, 53]. Furthermore, the authors observe that “excessively high” initial temperatures do not affect geometric schedules [130, 53].

For our experiments, a static schedule implementation (Linear Cooling) is used with the purpose of producing performance results closer to that of a baseline typical SA implementation. Dynamic schedules have been shown to improve the performance (in terms of solution quality) over some static schedules - therefore, their use would present an atypical, and thus unfair, viewpoint of SA when compared to typical implementations of other metaheuristic approaches. The following subsections outline a few of the more common static schedules; as we make no use of dynamic schedules, these are not covered.

4.5.3.1 Geometric Cooling Schedule

Also referred to as the exponential cooling schedule and first proposed by Kirkpatrick et al. in [69] during their early work on cooling schedules [1], the geometric cooling schedule is one of the simplest of the static schedules but is still often in use in many practical applications [1]. It is categorised by the use of a positive real-valued constant α with a value less than but close to 1.0; typically set between 0.8 and 0.99 [1]. The temperature decrement function is then given by [98]:

$$T_t = T_0 \alpha^t \quad (4.12)$$

where T_t represents the temperature at time/cycle t and T_0 is the starting temperature.

4.5.3.2 Logarithmic Cooling Schedule

Of some special theoretical importance, the logarithmic cooling schedule first proposed by Geman and Geman in [37], makes use of the following decrement function [98]:

$$T_t = \frac{c}{\log(c + d)} \quad (4.13)$$

where d is typically set to a value of one. It has been shown theoretically in [43], “that when c is set greater than or equal to the largest energy barrier in the problem, this schedule will lead the system to the global minimum state in the limit of infinite time”[98]. However its use is impractical due to its extremely slow rate of temperature decrease [98].

4.5.3.3 Linear Cooling Schedule

In the linear schedule, another commonly used static schedule [98], the value of the temperature parameter is decreased by a constant amount by the decrement function [98]:

$$T_t = T_0 - \eta t \quad (4.14)$$

4.5.4 Simulated Annealing (SA): Implementation

4.5.4.1 Details and Description

The main detail defining this implementation of SA is that it makes use of a linear cooling schedule as opposed to the more common geometric schedule. This decision was made to account for the automatic tuning of the algorithm and the fact that comparisons would be limited to a maximum number of function evaluations. When tuning SA making use of a geometric schedule, the automatic parameter tuner would often produce a value for α

that reduced the temperature too quickly resulting in a large portion of the search being carried out at the minimum temperature - often very close to zero. Since these portions of the search processes are almost equivalent to a hill climb (descent) algorithm, it was felt that the results from the algorithm would not fully represent the performance achievable by SA when compared to hill climbing. As such, a linear schedule was selected that would not use parameters tuned by our parameter tuner but instead used a constant decrement calculated as a function of: (i) the maximum number of function evaluations E_{max} to be carried out before termination, (ii) the number of iterations to be carried out at each temperature T_{iters} and (iii) the initial and minimum values of T - T_{init} and T_{min} . Specifically, the decrement d was calculated as follows:

$$d = \frac{T_{diff}}{T_{states}} \tag{4.15}$$

where: $T_{diff} = T_{init} - T_{min}$

and: $T_{states} = E_{max}/T_{iters}$

T_{diff} represents the difference between the initial temperature and the minimum temperature and T_{states} is the total number of temperature states. This function results in a value for d that adapts to the other parameters tuned automatically by our selected tuner so as to reach T_{min} on the final round of iterations. The linear cooling schedule is also used during the tuning of the remaining parameters, so although we have removed the ability to tune the cooling schedule itself, parameter configurations will be sought that maximise the performance of the algorithm when using this scheme.

4.5.4.2 Pseudocode

Algorithm 10 Pseudocode of SA Implementation

```

1: procedure SA(stepValue, pSize, startT, minT, iterT  $\alpha$ )           ▷ Main loop procedure
2:   stepSize  $\leftarrow$  stepValue
3:   problemSize  $\leftarrow$  pSize
4:   T  $\leftarrow$  startT
5:   Tmin  $\leftarrow$  minT
6:   Tdec  $\leftarrow$  GETTDECREMENT
7:   iterPerT  $\leftarrow$  iterT

8:   Sbest  $\leftarrow$  null                                           ▷ Overall best solution discovered
9:   Sbest.score  $\leftarrow$  null

10:  S  $\leftarrow$  GENERATEINITIALSOLUTION(problemSize)
11:  Sscore  $\leftarrow$  EVALUATESOLUTION(S)

12:  while (terminate == false) do
13:    for (i = 0 to iterPerT-1 by 1) do
14:      N  $\leftarrow$  GENERATENEIGHBOUR(S, problemSize, stepSize) ▷ random neighbour of
      S
15:      Nscore  $\leftarrow$  EVALUATESOLUTION(N)

16:      if (Nscore < Sscore) then
17:        S  $\leftarrow$  N
18:        Sscore  $\leftarrow$  Nscore
19:        if (Sscore < Sbest.score) then
20:          Sbest  $\leftarrow$  S
21:          Sbest.score  $\leftarrow$  Sscore
22:        end if
23:      else
24:        r  $\leftarrow$  uniform random float  $\in$  [0.0, 1.0]
25:        if (CALCULATEACCEPTANCE(Sscore, Nscore, T) < r) then
26:          S  $\leftarrow$  N
27:          Sscore  $\leftarrow$  Nscore
28:          if (Sscore < Sbest.score) then
29:            Sbest  $\leftarrow$  S

```

```

30:          $S_{\text{best.score}} \leftarrow S_{\text{score}}$ 
31:     end if
32: end if
33: end if
34: end for
35:  $T \leftarrow T - T_{\text{dec}}$ 
36: if ( $T < T_{\text{min}}$ ) then
37:      $T \leftarrow T_{\text{min}}$ 
38: end if
39: end while

40: return S ▷ return best solution discovered
41: end procedure

42: procedure GENERATEINITIALSOLUTION(problemSize)
43:     solution  $\leftarrow$  init
44:     for ( $i = 0$  to  $\text{problemSize} - 1$  by 1) do
45:         solution[i]  $\leftarrow$  uniform random float  $\in [0.0, 1.0]$ 
46:     end for
47:     return solution
48: end procedure

49: procedure GENERATENEIGHBOUR(solution, problemSize, stepSize)
50:     neighbour  $\leftarrow$  solution
51:      $r \leftarrow$  uniform random float  $\in [0.0, 1.0]$ 

52:     boundMin  $\leftarrow$  neighbour[i] - stepSize
53:     boundMax  $\leftarrow$  neighbour[i] + stepSize

54:     if (boundMin < 0.0) then ▷ Boundary handling
55:         boundMin  $\leftarrow$  0.0
56:     else
57:         if (boundMin > 1.0) then
58:             boundMin  $\leftarrow$  1.0
59:         end if
60:     end if

```

```

61:   boundRange ← boundMax – boundMin
62:   r ← uniform random float ∈ [0.0, 1.0]
63:   step ← boundRange * r
64:   neighbour[i] ← boundMin + step

65:   return neighbour
66: end procedure

67: procedure CALCULATEACCEPTANCE( $S_{score}$ ,  $N_{score}$ , T)
68:   return  $\exp((S_{score} - N_{score})/T)$ 
69: end procedure

70: procedure GETTDECREMENT
71:   Tdiff ← T – Tmin
72:   totalT – iters ← maxEvals/iterPerT
73:   return Tdiff/totalT – iters
74: end procedure

```

4.5.4.3 Tunable Parameters

The parameters used to control the performance of the algorithm in this implementation are as follows:

- Step Size
 - See Sections 4.1.3 and 4.2.3.3 for a complete description
- Initial T
 - Type:- float ∈ [0.5, 200.0]
 - Defines the initial temperature from which to start execution of the algorithm. The default value used initially by SMAC is 1.0
- Min T
 - Type:- float ∈ [0.00001, 0.45]
 - Defines the minimum temperature of the algorithm where the annealing schedule should stop reducing the temperature. In regular implementations, the minimum temperature would represent a stopping criterion. However, since we require each algorithm to run for a set number of objective function evaluations, if the minimum

temperature is reached before the evaluation budget is exhausted, the algorithm will perform at that temperature for the remainder of the allocated budget

- Iterations Per T
 - Type :- integer $\in [10, 500]$
 - Defines the number of iterations performed by the algorithm at each temperature step. The default value used initially by SMAC during tuning is 100

4.6 DIFFERENTIAL EVOLUTION (DE)

Differential Evolution (DE) is a population based metaheuristic from the class of evolutionary algorithms first introduced by Storn and Price in [129] and was designed specifically for the optimisation of real-valued, real parameter functions. In a similar way to other evolutionary algorithms such as GA, DE generates new solutions through mutation and crossover. Mutation in DE creates a new solution by adding a scaled difference between two randomly chosen solutions from the population to a third random solution [95, 129]. The resulting solution is referred to as the donor solution and will be crossed over with another randomly preselected (parent) solution referred to as the target solution to produce a new candidate solution referred to as the trial vector [129, 23]. If the trial solution possesses a better objective score than the target solution the target solution will be replaced by the trial solution in the next generation - this final step is referred to as selection [129].

This solution generation scheme differentiates DE from other evolutionary algorithms that typically generate new solutions in a more probabilistic manner through the action of probabilistic crossover and mutation operators [95]. Another differentiation is that DE prescribes the replacement of child solutions only if that solution is better than its parent [95]; GA, for example, instead replaces each child solution, even if that child creates a worsening move in respect to the population.

A typical DE implementation can now be said to consist of four fundamental steps: (i) population initialisation, (ii) mutation, (iii) crossover and (iv) selection, each of these will be covered in more detail in the next sub section.

4.6.1 *The Differential Evolution Process*

In this section each of the four fundamental steps in differential evolution will be discussed in order of occurrence in the evolutionary process.

4.6.1.1 *Initialisation of the Population*

DE begins with a randomly initialised population of NP n-dimensional real-valued solution vectors which represents a candidate solution to the n-dimensional problem being tackled [23]. For each parameter in a solution vector, it is often the case that it is to be bound within a given range beyond which the solution will become infeasible. Storn and Price [129] note

that random population generation should cover the whole search space and so prescribe the use of a uniform probability distribution for all random decisions in DE, including those involved in solution generation. This is also noted by Das and Suganthan [23]. Therefore, a given solution will be constrained by the minimum and maximum bounds of its parameters, as in [23]:

$$\vec{X} = \{[x_{1,\min}, x_{1,\max}], [\dots], [x_{n,\min}, x_{n,\max}]\} \quad (4.16)$$

So then the j^{th} parameter of the i^{th} solution vector may be instantiated as follows [23]:

$$x_{j,i} = x_{j,\min} + r_{i,j} \in [0, 1] \cdot (x_{j,\max} - x_{j,\min}) \quad (4.17)$$

where $r_{i,j}$ is a uniformly distributed real value from the interval $[0, 1]$ and is generated individually for each parameter j of solution vector i [23].

4.6.1.2 Mutation

As mentioned briefly earlier, for each target vector of index i , where $i = \{1, 2, \dots, NP\}$ and NP is the population size, the simplest form of mutation in DE is to create a donor vector from three randomly selected individuals $\vec{X}_{r_1^i}$, $\vec{X}_{r_2^i}$ and $\vec{X}_{r_3^i}$ in the population [23]. Where r_1^i , r_2^i and r_3^i are indices of the population in the range $[1, NP]$ which are to be mutually exclusive as well as being distinct from the target vector index i [129, 23]. This condition clearly holds only if the population size $NP \geq 4$ [129]. A constant $F \in [0, 2]$ [129], which we refer to as the *differential weight*, scales the difference between any two of the three vectors and this scaled difference is added to the third [129, 23]. Thus, F is used to control the magnitude of the differential variation [129]. The mutation function for generating a mutant vector \vec{V}_i is as follows:

$$\vec{V}_i = \vec{X}_{r_1} + F \cdot (\vec{X}_{r_2} - \vec{X}_{r_3}) \quad (4.18)$$

A graphical example of the production of a mutant vector using differential mutation is shown in Fig. 4.10 below.

As shown in this example, the mutation step begins with the selection of three random position vectors (indicated in black) are first selected from the population. Next, the difference vector is calculated between the two vectors \vec{X}_{r_2} and \vec{X}_{r_3} (represented by the longest red arrow); this difference is calculated by the bracketed sub-term in the second term of equation 4.18. Finally, the difference vector is weighted by $F \in [0, 2]$ and added to remaining random vector \vec{X}_{r_1} resulting in a new mutant position vector (the shorter red arrow) to be recombined with the parent vector in the crossover step. It is worth briefly mentioning that the labelling of

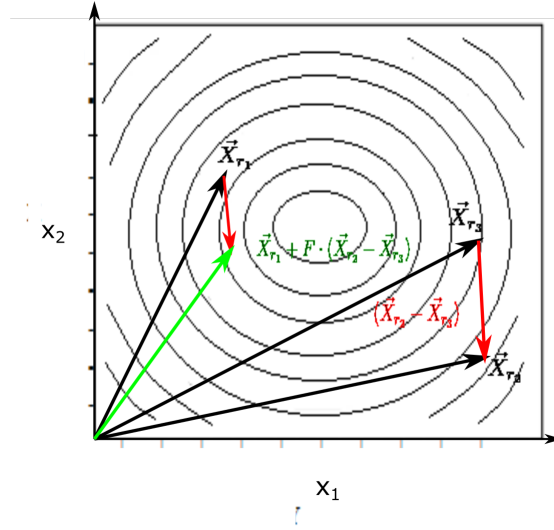


Figure 4.10: Production of a mutant vector \vec{V}_i from the scaled difference vector $(\vec{X}_{r_2} - \vec{X}_{r_3})$

the three vectors X_{r_1} , X_{r_2} and X_{r_3} is down to the order in which they have been selected from the population, thus, the choice of vectors to be used for the different terms of the mutation equation are arbitrary and so any other ordering of vectors in the equation should not impact the performance of the algorithm.

4.6.1.3 Recombination

A recombination step is introduced that is designed to improve the diversity in the resulting trial vectors, where a trial vector [129]:

$$\vec{U}_i = [u_{1,i}, u_{2,i}, \dots, u_{n,i}] \quad (4.19)$$

is generated according to [129]:

$$u_{ji} = \begin{cases} v_{ji,G+1} & \text{if } (\text{rand}_u \in [0, 1] \leq \text{CR}) \vee j = \text{rand}_n \in [0, n] \\ x_{ji,G} & \text{if } (\text{rand}_u \in [0, 1] > \text{CR}) \wedge j \neq \text{rand}_n \in [0, n] \end{cases} \quad (4.20)$$

where, $j \in [0, n]$

where, $\text{CR} \in [0, 1]$ is a user supplied constant referred to as the *crossover weight* which determines an approximate rate in which parameters will be donated from the mutant vector, rand_u is a uniformly random number - regenerated for every j^{th} parameter of the i^{th} vector [23]. rand_n is randomly chosen parameter index that ensured that at least one parameter value in \vec{V}_i is recombined with \vec{U}_i and is generated once for each vector in each generation [23]. This type of crossover represents the binomial (uniform) crossover described in Section 4.6.2 below [23].

4.6.1.4 Selection

In order to maintain the size of the population, a selection operator is applied to decide which of the target vector \vec{X}_i or \vec{U}_i is carried over to the next generation G , i.e., $G = G + 1$ [23]. For a minimisation problem, this selection is determined by:

$$\vec{X}_{i,G+1} = \begin{cases} \vec{U}_{i,G} & \text{if } f(\vec{U}_{i,G}) \leq f(\vec{X}_{i,G}) \\ \vec{X}_{i,G} & \text{if } f(\vec{U}_{i,G}) > f(\vec{X}_{i,G}) \end{cases} \quad (4.21)$$

where, $f(\vec{X})$ is the objective function being minimised [23]. In short, this means that the trial vector will replace the target vector in the subsequent population iff returns an equal or better objective function value and the target vector is kept otherwise [23]. As noted in Section 4.6 this is one way DE differs from other evolutionary algorithms, where the replacement of only solutions of greater or equal value means that the mean objective value of the population can only increase or stay the same but never degenerates [23]

4.6.2 Variants of DE

Apart from the commonly used DE scheme described thus far, there are as many as 10 different variants proposed by Price et al. in [129, 108]. Different DE strategies are represented by the notation $DE/x/y/z$ where [129, 108]:

1. x represents the base vector that is to be mutated
2. y denotes the number of difference vectors to be used, and
3. z specifies the crossover method to be used either *bin* (independent binomial) or *exp* (exponential)

The strategy presented above is described as $DE/rand/1/bin$, as: (i) the base vector is randomly selected, (ii) 1 vector difference is added to this base vector and (iii) the number of parameters provided by the mutant vector closely follows a binomial distribution [129, 108].

4.6.3 Differential Evolution (DE): Implementation

4.6.3.1 Details and Description

This implementation was based loosely on the differential evolution implementation source code found at [84] and pseudocode found in [129] and represents a $DE/rand/1/bin$ approach as described in Section 4.6.2. As such the implementation randomly selects a base vector (the

vector to be mutated), one difference vector is added to the base vector and the number of mutant parameters recombined with the base vector is roughly binomial.

4.6.3.2 Pseudocode

Algorithm 11 Pseudocode of DE Implementation

```
1: procedure DIFFERENTIALEVOLUTION(dWeight, xProb, popSize, pSize)      ▷ Main loop
2:   problemSize ← pSize
3:   populationSize ← popSize
4:   diffWeight ← dWeight
5:   xoverProb ← xProb
6:   population ← init()
7:   S ← GENERATEPOPULATION(problemSize)
8:   while (terminate == false) do
9:     xIndex ← uniform random integer ∈ [0, populationSize − 1]
10:    aIndex, bIndex, cIndex ← null
11:    do
12:      aIndex ← uniform random integer ∈ [0, populationSize − 1]
13:    while (aIndex == xIndex)
14:    do
15:      bIndex ← uniform random integer ∈ [0, populationSize − 1]
16:    while (bIndex == xIndex || bIndex == aIndex)
17:    do
18:      cIndex ← uniform random integer ∈ [0, populationSize − 1]
19:    while (cIndex == xIndex || cIndex == bIndex || cIndex == aIndex)
20:    randIndex ← uniform random integer ∈ [0, pSize − 1]
21:    parent ← population[xIndex]
22:    candidate ← parent                                          ▷ clone
23:    individualA ← population[aIndex]
24:    individualB ← population[bIndex]
25:    individualC ← population[cIndex]
26:    for (i = 0 to pSize by 1) do
27:      r ← uniform random float ∈ [0.0, 1.0]
28:      if (i == randIndex || r < xoverProb) then
29:        candidate[i] ← individualA[i] + diffWeight × (individualB[i] −
        individualC[i])
30:      end if
31:    end for
```

```

32:   for (j = 0 to pSize by 1) do           ▷ Boundary handling: absorbing to LB or UB
33:       if (candidate[j] < 0) then
34:           candidate[j] = 0.0
35:       else
36:           if (candidate[j] > 1) then
37:               candidate[j] = 1.0
38:           end if
39:       end if
40:   end for

41:   parent.score ← EVALUATESOLUTION(parent)
42:   candidate.score ← EVALUATESOLUTION(candidate)
43:   if (candidate.score < parent.score) then
44:       population[xIndex] ← candidate
45:   end if
46: end while
47: return GETBEST(population)
48: end procedure

```

4.6.3.3 Tunable Parameters

The parameters used to control the performance of the algorithm in this implementation are as follows:

- Population Size
 - Type:- integer $\in [2, 250]$
 - Used by the initialization procedure to define the size of the population that will remain fixed throughout execution of the algorithm. The default population size initially trialled by SMAC is 20
- Differential Weight
 - Type:- float $\in [0.01, 1.0]$
 - Used in the mutation step to scale the difference between two of the three randomly selected mutation vectors which are to be added to the third random vector. The default value of this parameter and therefore the first one attempted by SMAC is 0.5. In [23], the authors state that this parameter has a value that typically lies in the range $[0.4, 1.0]$ so this default seemed an appropriate starting point.

- Crossover Rate
 - Type: float $\in [0.01, 1.0]$
 - Used in the binomial crossover procedure, this parameter represents the probability (or rate) that parameters from the mutant vector is passed on to the trial vector. Further, Das and Suganthan [23] note that the use of crossover rate only aids in approximating the true probability of a crossover occurring.

CMA-ES is an extension of the well-known Evolution Strategies, and was mostly developed by Nikolaus Hansen and Ann Auger [82]. It targets difficult black-box optimisation problems, which may be non-linear and non-convex [45].

The framework is highly complex so we focus here only on the high-level aspects and parameters we consider for variation. Additionally, whilst variants intended for multi-objective optimisation and elitist variants have been proposed [45] this section will only deal with single objective optimisation and non-elitist selection. A more in depth treatment of this approach can be found in the excellent tutorials by Hansen [45] and Luke [82].

4.7.1 Algorithm Overview

The classic version of CMA-ES, more formally referred to as $(\mu/\mu_w, \lambda)$ CMA-ES, contains several components intended to improve its performance. Further from [82], we will begin with a description of a basic implementation without inclusion of these components before discussing these components further. CMA-ES generates new solutions by sampling a multivariate Normal distribution, specified by a mean vector \vec{m} of length n and an $n \times n$ covariance matrix C [82].

During execution of the basic CMA-ES loop, the algorithm samples the μ highest quality solutions from its current distribution, evaluates these solutions and uses these to update the distribution [82]. Sampling and updating the covariance matrix to reflect these solution vectors C is achieved more easily through eigendecomposition of C into BDB^T , which effectively splits the matrix into two separate matrices B and D [82].

In CMA-ES, new solutions are sampled randomly from the current distribution - where the actual sampling is scaled by a value σ , the step size parameter, which for all intents and purposes acts as the algorithms mutation rate; the larger the value of σ , the further spread out the samples are relative to the current distribution [82].

Given the components: C , \vec{m} and σ , a point P_i in the solution space is represented by three different but equivalent forms, aimed at making the mathematics easier to perform[82]. Specifically, a solution is given by: $P_i = \{\vec{x}^{(i)}, \vec{y}^{(i)}, \vec{z}^{(i)}\}$; these are defined as:

1. $\vec{z}^{(i)}$: an n-dimensional real-valued solution which has been sampled at the multivariate Normal distribution $N(\vec{0}, I)$ where the mean $\vec{0}$ is the origin (at zero) with I - the identity matrix - as the $n \times n$ covariance matrix [82]
2. $\vec{y}^{(i)}$: is the point $\vec{z}^{(i)}$ that has been transformed into the space where the distribution currently is, i.e., $N(\vec{0}, I)$. This is achieved through the use of the equation $\vec{y}^{(i)} = C^{\frac{1}{2}} \vec{z}^{(i)} = BDB^T \vec{z}^{(i)} = BD^T \vec{z}^{(i)}$ [82]
3. $\vec{x}^{(i)}$: this is the point $\vec{y}^{(i)}$ scaled by the step size σ and translated by \vec{m} i.e., $\vec{x}^{(i)} = \vec{m} + \sigma \vec{y}^{(i)}$ [82]

Remaining in the context of this simple implementation, the distribution can be updated by completely rebuilding the distribution to fit with the new samples, for example by firstly updating the mean as in [82]:

$$\vec{m} \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} \vec{x}^{(i)} \quad (4.22)$$

and update the covariance matrix by [82]:

$$C \leftarrow \frac{1}{\mu-1} \sum_{i=1}^{\mu} (\vec{x}^{(i)} - \vec{m})(\vec{x}^{(i)} - \vec{m})^T \quad (4.23)$$

However, given that the fitnesses for our set of \vec{x} have already been calculated so we could use their rank order instead of simply relying on truncation selection. This is done by sorting the truncated population such that fitter solutions appear first [82]. Weights are then calculated for each member of the ranked truncated population to give each more or less weight that could impact the outcome of \vec{m} and C . These should be assigned in such a way as to have all weights sum to 1.0 [82]. Using these weights the distribution mean \vec{m} and the covariance matrix C can be updated by the following [82]:

$$\vec{m} \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} \omega_i \vec{x}^{(i)} \quad (4.24)$$

and update the covariance matrix by [82]:

$$C \leftarrow \frac{1}{\mu-1} \sum_{i=1}^{\mu} \omega_i (\vec{x}^{(i)} - \vec{m})(\vec{x}^{(i)} - \vec{m})^T \quad (4.25)$$

However, and as mentioned previously, CMA-ES makes use of several components in order to improve the performance of the algorithm [82].

4.7.2 CMA-ES for Large-scale Optimisation

It is well documented that due to the high computational cost associated with calculation of the co-variance matrices in particular, that CMA-ES in its original form is unsuitable for use on

large-scale optimisation problems [114]. Indeed, from our own experiences with the algorithm, by the time we attempted to solve for 1000D problems, runtime was already exceeding 2-3 hours for a single repetition. Given that there was also a need to tune CMA-ES as we did for our other approaches (See Chapter 6) where up to 1000 repeats of the algorithm would be performed, using the full (original) approach was not feasible.

A partial solution to the above problem comes from a small modification to the approach *sep-CMA-ES* - as introduced by Rose and Hansen in [114] and provided with our implementation (see Section 4.7.3) - where instead of calculating full co-variance matrices, only diagonal matrices are computed; reducing the time and space complexity from quadratic (as we ourselves were observing) to linear - so effectively reducing model complexity and increasing the learning rate [114]. This modification however comes with reduced generality to the types of problem of which it is suitable; specifically, by reducing the co-variance matrices to diagonal (identity) matrices the degrees of freedom possessed by the original implementation that provided rotational invariance - allowing it to successfully solve many non-separable problems - is much reduced meaning that dependencies between solution variables are not captured and coordinates are sampled independently [114]. Consequently, some decreased performance should be expected from this variant [114]. Ros and Hansen report that for separable problems, *sep-CMA-ES* “significantly outperforms” regular CMA-ES with the added benefit that the time and space scale-up is linear [114]. Further, for low to moderate levels of non-separability *sep-CMA-ES* still outperforms CMA-ES for all dimensionalities by almost a factor of $\frac{n}{10} + 1$ [114]. However, for larger levels of non-separability, the authors report an advance to *sep-CMA-ES* only on problem dimensionalities greater than 100 [114]. Lastly, for a single fully non-separable problem instance (the block-ellipsoid function) regular CMA-ES was always found to be significantly superior [114].

However despite the possible loss of performance over our benchmark suite, of which there are several non-separable problems, we continue to use *sep-CMA-ES* for our experiments for its significant - and in practice, noticeable - reduction in time complexity.

4.7.3 Covariance Matrix Adaption Evolution Strategy (CMA-ES): Implementation

4.7.3.1 Details and Description

The CMA-ES implementation is that of a Java library sourced from [44] to work within our experimental framework. This implementation of the algorithm, using its default settings,

utilises a full covariance matrix for its updates, which for all benchmark functions at much higher dimensions, (≥ 1000), had a runtime in the region of 2-3+ hours for a single repeat of the algorithm. As such, this setup was unsuitable for our purposes where dimensions greater than 1000 were being probed and 20 repeats at each dimension of a problem were necessary. Furthermore, when tuning the algorithm, up to 1000 repetitions of the algorithm would be carried out by SMAC in search for parameter configurations, so this represented a greater obstacle to the approach.

However, the developers site, from which the source code was found, noted an inbuilt setting which if toggled would instruct CMA-ES to use only diagonal covariance matrices and speed up the approach from quadratic to linear time as described in the relevant article pertaining to this development, *sep-CMA-ES*, by Ros and Hansen [114]. They state however that this modification to the original approach makes *sep-CMA-ES* less effective on most non-separable objective functions, however it was found to vastly outperform the original implementation on separable functions. We found that *sep-CMA-ES* was still competitive with our other implementations and so still included in our algorithm set.

4.7.3.2 *Tunable Parameters*

The parameters used to control the performance of the algorithm in this implementation are as follows:

- Population Size
 - Type:- integer $\in [6, 150]$
 - As with the other population based approaches described, this parameter determines the number of solutions maintained by the algorithm at any one time. The lower bound on this parameter of 6 was determined to be suitable since the implementation as written defaulted to a population size λ of 6 for a 2-dimensional problem through the function: $\lambda = (4 + 3) \times \log_e(n)$, where n is the problem dimensionality. The upper bound was determined to be lower than that of the other population based approaches to further reduce the overall runtime since this placed a smaller upper bound on the number of diagonal covariance matrices that would need to be calculated.
- Initial Sigma (Step-size)
 - Type:- float $\in [0.2, 30.0]$
 - Represents the initial standard deviation used in the parameterised Normal distribution modelled by the algorithm and thus can be regarded as the initial step

size. Again, the default σ used by the algorithm was 0.2 so therefore we made use of it also as the parameters lower bound. The upper bound was determined by studying an article [46] by Hansen and Kern, which investigated the modification of this parameter in the context of several standard continuous benchmarks - some of which are included in our own suite. The authors concluded that a value of half the range between the lower and upper bound of the search space represented a reasonable setting; stating that a “considerable impact” on the performance against multi-modal functions is observed if initial step size is set too low [46]. However, since we are investigating with several benchmark functions a reasonable intermediate upper bound had to be selected. Given that one of our functions (specifically, the Griewank function) had parameter bounds of $[-600, 600]$, half of this would present too many possible combinations for SMAC to enable it to find an effective set of parameter configurations.

CHAPTER 5: THE PROBLEMS OF ALGORITHM SELECTION AND CONFIGURATION

5.1 INTRODUCTION

The algorithm selection problem (ASP) first identified and discussed by Rice [112] is the problem of selecting the most appropriate algorithm for solving a given problem and has become increasingly relevant to optimisation researchers due to the realisation that no one algorithm can be developed that can work effectively against all problems - as per the NFL (Chapter 2 - section 2.4.5.2). Rice in [112] presented a formal abstract model that aimed to help explore this problem [139]:

1. The space of problems P [139]
2. The feature space F of all measurable characteristics of the problems contained in p as calculated by a feature extraction procedure f [139]
3. The space of algorithms A [139] that can be used against problems in P
4. The performance space Y containing a mapping between each algorithm $a \in A$ to a set of performance metrics [139]

Following Rice [112] and Vanschoren [139] the ASP can then be defined as follows:

“For a given problem instance $x \in P$ with features $f(x) \in F$, find a mapping $S(f(x))$ into algorithm space A , such that the selected algorithm $\alpha \in A$ maximises the performance mapping $y(\alpha(x)) \in Y$ ”

Where S is a selection mapping method responsible for selecting the appropriate algorithm α given features $f(x)$ [139].

Simply put, the algorithm selection problem involves selecting the most appropriate algorithm from a suite of algorithms that best solves a given problem instance.

The various algorithms used to solve optimisation problems, in terms of both deterministic and stochastic algorithms, are often highly parameterised, sometimes involving tens or

hundreds of parameters ¹. Finding an optimal set of values for these settings, hereafter referred to as a *parameter configuration*, can greatly improve the performance of a target algorithm against a given set of problem instances but this can often be an intractable goal - especially when large numbers of parameters and many possible values for each parameter are concerned.

The rationale for tuning the metaheuristic algorithms used for our purposes is to provide an appropriate set of parameter configurations for each such that all later experiments comparing the results of the metaheuristics can be made in the context of well performing algorithms for each problem instance at each given dimension. By tuning individually for each dimensionality of a problem we can obtain an appropriately 'good' parameter configuration for a given dimensionality - meaning that comparisons between other algorithms for the same problem instance at the same dimensionality can be made more meaningfully. All experimentation with the metaheuristic approaches in this thesis requires that each approach gives the best possible performance, in terms of average solution quality, for each given problem instance. This is very different from requiring that all algorithms perform equally well on all problem instances as prohibited by the no-free-lunch theorem. When comparing metaheuristic approaches, the no-free-lunch theorem falsifies any claim that one approach is better than any other when averaged over all problems, this also carries the implication that for some problems an approach may simply perform badly when compared to others - regardless of the settings of any tunable parameters. However, there will exist, for all approaches with at least one tunable parameter, one or more parameter configurations that facilitate the best performance on a given problem in relation to all other configurations.

5.2 RELATIONSHIPS BETWEEN ALGORITHM SELECTION AND OTHER AREAS OF RESEARCH

5.2.1 *Meta-learning*

Most related to the problem of algorithm selection is that of *meta-learning*, a term originating from the field of educational psychology. Meta-learning has been described as "*being in aware and taking control of one's own learning*" [9] or alternatively, *learning to learn*. More concisely, Lemke et al. considers meta-learning as being an understanding and subsequent adaptation of learning itself - on a higher level - rather than simply about gaining "subject knowledge". An individual can then introspect about their approach to leaning and adapt it depending on

¹ a parameter here refers to a variable that is part of an algorithm that can be tuned to alter the behaviour and/or performance of the algorithm

the given task requirements [74]. As used in the context of machine learning, the concept of meta-learning is still closely based on this same description. However, here the term *subject knowledge* will translate to base learning, which refers to the acquisition of experience of a specific learning task [74]. In [14], Brazdil et al. explain the distinction between base-learning and meta-learning, as being categorised mainly by the scope of the level of their adaptations; specifically, base-learning is concerned only with the accumulation of experience on a specified task, where meta-learning on the other hand focusses on gathering experience over a number of applications of a learning system [14].

In this section we will follow a brief outline of meta-learning; specifically on its relationship to algorithm selection in particular but also introducing specific methods widely held to be instances of meta-learning as well as a comparison in terms of the how the constituent components of meta-learning are represented, or not represented, in each.

5.2.1.1 *Definition*

So close is the relationship between meta-learning and algorithm selection, meta-learning as a concept can be traced back directly to the formalisation of the algorithm selection problem by Rice [112]. However, the actual use of the term only appeared in the literature for machine learning in the 1990s [74]. Since then, many definitions of meta-learning have appeared in the literature, most sharing the common understanding - as noted by Lemke et al. [74] - that "*meta-learning becomes meta-learning by looking at different problems, domains, tasks or contexts or simply past experience*". This understanding is most easily seen in the definition presented by Brazdil et al. in [14] stating that; "*meta-learning is the study of principled methods that exploit metaknowledge to obtain efficient models and solutions by adapting machine learning and data mining processes*". The authors explain that the term *metaknowledge* refers to the knowledge obtained from past learning experiences, either from past experiences on the same data or using data obtained from another problem domain [14]. Therefore, it would be reasonable to conclude from this definition that the main prerequisite for any meta-learning approach is the use of at least one form of metaknowledge. From the common understanding found in the various definitions, Lemke et al. propose their own more formalised definition - provided here verbatim [74]:

1. A meta-learning system must include a learning sub-system, which adapts with experience
2. Experience is gained by extracting metaknowledge extracted
 - a) ...in a previous learning episode on a single dataset, and/or

b) ...from different domains or problems

A further concept in meta-learning, also noted by Lemke et al. [74] is that of bias, here, referring to a set of assumptions which influence the choice of hypotheses that could explain the data. Brazdil et al. [14] refer to two types of bias that could impart influence in a meta-learning system. Firstly, *declarative bias* which relates to the representation given to the space of hypotheses, for example, where hypotheses are represented exclusively by neural networks; and *procedural bias*, affecting the precedence of hypotheses, for example, preferring hypotheses providing shorter runtimes over others [74, 14]. The assumption given this notion is that the bias in the base-learners remains fixed, where for meta-learning the system attempts to select the most appropriate bias dynamically [74].

5.2.1.2 Algorithm Recommendation / Selection

Out of the many forms of 'meta-learning' discussed in the literature, including but not limited to ensemble methods and dynamic bias selection, algorithm recommendation is the only form which is applicable to each point from the definition outlined above according to Lemke et al.[74]. Here the interest to the metalearner is the relationship between the characteristics of the data and the algorithm performance with the end goal of being able to select one or more algorithms from a set of algorithms suitable for a given problem ² [74]. Since it is generally infeasible to select amongst all the possible alternative algorithms available for a given problem in a trial and error manner, meta-learning can therefore prove useful in recommending an algorithm to an end-user or in automatically weighting algorithms according to their suitability [74].

Vanschoren [139] also points out that it is not simply the algorithms themselves that allows the performance of the same algorithm to vary over different datasets, but also that these algorithms can have different parameter configurations. Therefore it would be reasonable to suggest that an algorithm with various different parameter settings can be considered as representing different algorithms altogether [139]. The author however suggests that the study of this subject and its effects should remain separate.

In machine learning the traditional application of algorithm selection is classification; however, efforts have been made to generalise the concepts to other areas of application such as: regression, sorting, constraint satisfaction and optimisation [126, 74]. For the latter,

² Note that this description of algorithm selection in the meta-learning context is equivalent to that defined by Rice as described in Section 5.1

optimisation, one such application example by Pavelski et al. [102] is that of meta-learning optimisation for the flowshop problem through the use of decision trees. Specifically, the work proposed a meta-learning approach for knowledge discovery that operated on the performance data from four different metaheuristic algorithms whilst they solved several flowshop problem instances. As well as being able to recommend the most suitable metaheuristic for each problem instance, it is also capable of suggesting well-suited parameter configurations [102]. Although the authors suggest that for algorithm recommendation in this particular case study, there is “a lot of room for improvement”, their results were deemed promising as the rules induced indicated that some metaheuristic parameters were more preferable than others [102].

5.2.1.3 *Considerations for Using a Meta-learning Approach*

This section focusses on several practical decisions which have to be made before applying meta-learning to a problem, which can include: the selection of a meta-learning algorithm, the choice of appropriate metaknowledge and the issue of implementing and maintaining *metadatabases* [74].

A) PREREQUISITES It is usually beneficial to consider whether or not a particular methodology is appropriate for use at all against a given problem or application before effort is expended. Meta-learning is no different, and as such can not be considered as a “magic cure” to any and all machine learning problems [74]. As Brazdil points out in [14], it is important that the metafeatures extracted from a problem domain are representative of that domain; a failure to ensure this will prevent meta-learning from identifying whether or not other domains are similar [74]. Similarly, attempting to apply meta-learning to new problems which have never been tackled in the past (or rarely so) will not be able to take advantage of past experiences which would improve its predictions [14, 74]. Other factors influencing the applicability of meta-learning also include: (i) the possibility that performance estimations - being reliant on the accuracy and effectiveness of being able to estimate true performance - may be too unreliable and (ii) different metafeatures might be applicable to each dataset [74].

B) META-LEARNING ALGORITHMS The selection of a particular meta-learning algorithm is very problem specific, in general, Lemke et al. [74] suggest that traditional algorithms used in classification can be very successful for meta-learning algorithm selection and could include: meta-decision trees, neural networks or support vector machines (SVM) to name but a few. The application of regression algorithms however remains less popular [74].

C) **META-KNOWLEDGE** As discussed in section 5.2.1.1, Brazdil et al. [14] define meta-knowledge as being derived from the application of a learning system. A common type of metaknowledge is the performance of algorithms in the context of certain problem domains, which is then linked to characteristics of the task [74].

In terms of extracting metaknowledge from the data, the most intuitive include statistical or information-theoretic features e.g., according to Brazdil [14], for classification problems we could use the number of classes or features, ratio of examples to features and the average entropy of a class/classes [74]. For different areas of application however, features can be very different [74].

As an alternative to focussing on the data only for extracting metaknowledge, information about individual algorithms and how they produced a solution to the problem could also be used such as their predicted confidence [74]. One method of achieving this is by building a model that is both easy to generate and train and interrogating its properties. Another is an approach known as landmarking and roughly involves using the performance of simple algorithms in order to describe a problem and then correlating this information with the performance of far more complex algorithms [74].

It is also worth noting that as with all learning problems, meta-learning is also subject to the curse-of-dimensionality (Chapter 3, Section 3.2), which is ‘traditionally’ solved by selecting only a subset of the most relevant features - as would be done with regular feature selection of which meta-feature selection is not foundationally different [74].

5.3 AUTOMATIC PARAMETER TUNING

The problem can be laid out more formally as the *algorithm configuration problem* (ACP) [59, 58], where, given a target algorithm with tunable parameters A , a set or distribution of problem instances I and a cost metric c the goal is to find a parameter configuration(s) of A that minimises the cost c on I [59, 58].

If we were to choose to tune exhaustively on an approach with p tunable parameter settings each with j discrete values, a total number of dimensions d (according to the geometric progression) over a set of benchmark functions \mathcal{F} with cardinality $|\mathcal{F}|$, we would need to test $j^p \times d \times |\mathcal{F}|$ individual parameter combinations. For example, an approach with only 2 tunable parameters with 10 discrete values each, considering 10 different dimensions and a

set of benchmark functions with cardinality 10 then $10^2 \times 10 \times 10 = 10,000$ separate tunings would need to be carried out to ensure the best parameter combination for each scenario. It follows that with more tunable parameters, an increased number of parameter values, more dimensions and over a larger set of benchmark functions, this approach quickly becomes intractable.

For this reason many automatic parameter tuners (APTs) have been developed that can produce a reasonable, albeit approximate, set of parameter configurations without the need to test every possible configuration of the target algorithm.³ Several notable APTs are described and made use of in the literature, including model-free methods such as: F-Race, iRace, ParamILS and ROAR along with model-based techniques such as SMAC; each with their own set of benefits and drawbacks. In terms of optimising the parameters of the algorithms used in this study we opted to make use of the SMAC automatic configuration tool described in [59, 58].

5.3.1 *Model-free vs. Model-based Methods*

The existing APTs differ in whether or not explicit models are used to represent the dependence of the performance of a target algorithm on particular parameter configurations [59]. Model-free approaches, studied as far back as the 1990s, have the advantages of being able to be used ‘out-of-the-box’ and are relatively simple when compared to model-based methods [59, 58]. The various model-free approaches can be separated into those which focus only on optimising numerical parameters (i.e., integer or real-valued) such as F-race and those which can also be used within categorical domains (i.e., discrete and unordered) such as ParamILS [59, 58].

Model-based methods, on the other hand, has been cited as a promising path towards the next generation of automatic algorithm configuration methods [59, 58]. A well known branch of this class of methods is the Sequential Model-Based Optimisation (SMBO) framework, with benefits ranging from:

³ However, in our case, automatic parameter tuning only improves the value of j^P meaning that, given the same example as above, the total number of individual tuning operations, if the APT only tested half of the possible configurations, would be reduced from 10,000 to $0.5(10^2) \times 10 \times 10 = 5,000$ tuning operations - the number of dimensions to tune against as well as the number of benchmark functions is not affected by the use of an APT.

1. Being able to interpolate performances between between parameter settings that have been observed [59, 58] ⁴
2. Able to extrapolate from what is known to unseen regions of the parameter space [59, 58]
3. And, can be used to provide quantitative data regarding the importance of both individual parameters and interactions between parameters [59, 58]

SMBO has however been shown to have a few limitations - on account of being derived from the statistics literature on black-box optimisation - that make it an inappropriate choice for automated algorithm configuration. Amongst other limitations, these include: a reliance of computationally expensive models and being focussed heavily on deterministic algorithms [59, 58]. However, work continues to overcome these and other limitations of the framework - such examples include the development of methods such as Random Online Aggressive Racing (ROAR) and Sequential Model-based Algorithm Configuration (SMAC).

In simple terms, SMBO works by iterating between fitting models based on observed configurations and using these to select promising configurations [59]. Specifically, SMBO methods construct regression models from sampled configurations that are used to predict performance and then use these model in the optimisation process. According to Hutter at al. [59, 58]:

“In the context of parameter configuration, the model is fitted to training set $\{(\theta_1, o_1), \dots, (\theta_n, o_n)\}$ where parameter setting $\theta_i = (\theta_i, 1, \dots, \theta_i, d)$ is a complete instantiation of the target algorithm’s d parameters and o_i is the target algorithm’s observed performance when run with configuration θ_i . Given a new configuration θ_{n+1} , the model tries to predict its performance o_{n+1} .”

As mentioned above, SMBO was derived from the statistics literature on global black-box function optimisation. As such, a simple example of its basic workflow as provided by Hutter et al. in [59, 58], considers a rough approximation of the algorithm of the most notable example in this literature [59, 58], Efficient Global Optimisation or (EGO) [64]. Here, they consider a deterministic algorithm A with a single continuous parameter x and runtime as a 1-dimensional function of x represented by the solid line shown in Fig.5.1a, SMBO looks to optimise (minimise) the runtime of algorithm A by searching for a suitable value for x . Also note in Fig.5.1a that SMBO is first initialised by running algorithm A with the parameter

⁴ interpolate here is used in the mathematical sense, where performances of new parameter settings within the range of a set of known parameter setting performances (data points) can be constructed.

values indicated by circles [59, 58]. Next, a Gaussian process (GP) regression model⁵ [109] is fit to the sampled parameter value data [59, 58]. The black dotted line in Fig.5.1a and b shows the predictive mean of the GP regression model that has been trained on the parameter value samples given by the circles; where the shaded region surrounding this line represents the uncertainty in prediction - growing as distance increases from the training data [59]. SMBO then uses this performance data in order to predict and then select a promising configuration for the next run of algorithm A [59, 58]. Promising parameter configurations are predicted to perform well and/or be found in regions where the model is still uncertain [59]. These expectations are both encapsulated by the Expected Improvement (EI), the dashed lines in Fig.5.1, which has a larger value in areas of low predictive mean and high predictive variance [59, 58] (representing high performing configurations and found in uncertain areas respectively). Expected Improvement (EI) therefore provides a automatic tradeoff between focussing on promising parts of the parameter space (exploitation) and finding out more information about more unknown regions (exploration) [59, 58].

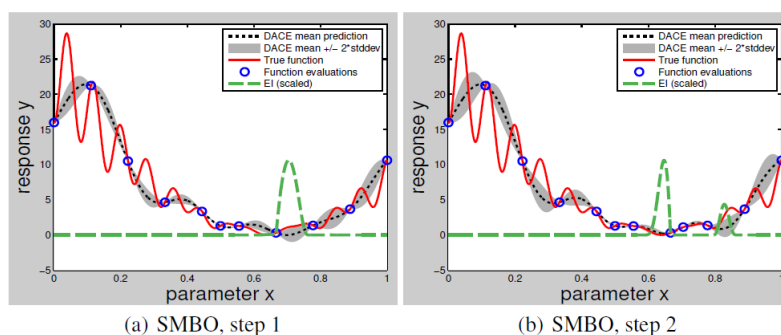


Figure 5.1: Taken from [58], This figure shows two steps of SMBO (EGO) for the optimization of a 1D function. The true function is shown as a solid line, and the circles denote our observations. The dotted line represents the mean prediction of a noise-free Gaussian process model (the DACE model), with the grey area representing its uncertainty. Expected improvement, EI, (scaled by the authors for visualisation) is shown by a dashed line.

An exact formula for EI, is introduced in [61] as $E[I_{exp}]$ for log-transformed costs and used for this example by Hutter et al. [59, 58], is defined as:

$$EI(\theta) := E[I_{exp}(\theta)] = f_{\min} \Phi(v) - e^{\frac{1}{2} \sigma_{\theta}^2 + \mu_{\theta}} \cdot \Phi(v - \sigma_{\theta}), \quad (5.1)$$

where, μ_{θ} and σ_{θ}^2 represent the predictive mean and variance of the log-transformed cost obtained by a configuration θ respectively. Further, Φ denotes the cumulative distribution

⁵ A generalisation of a gaussian distribution, a gaussian process is defined by its mean function and covariance function (a gaussian distribution is defined by its mean vector and covariance matrix). Therefore, a Gaussian distribution is over vectors, where a Gaussian process is over functions. [109, 71, 40]

function (CDF) of a standard Normal distribution and f_{\min} is the empirical mean performance of the incumbent configuration θ_{inc} [59, 58].

5.4 SMAC: SEQUENTIAL MODEL-BASED ALGORITHM CONFIGURATION

SMAC is an instantiation of the general Sequential Model-Based Optimization (SMBO) framework and an extension of a previous approach known as Randomised Online Aggressive Racing (ROAR) which is also an instantiation of the general framework. This section discusses what makes SMAC different from ROAR and the basic iterative procedure by which it constructs and samples from its models in order to make predictions about suitable parameter configurations for a given target algorithm.

5.4.1 *The SMAC Tuning Procedure*

5.4.1.1 *Random Online Aggressive Racing (ROAR)*

Before continuing on to discuss SMAC specifically, it is first necessary to briefly cover its precursor, Random Online Aggressive Racing (ROAR), as SMAC makes use of the aggressive racing element of this approach. ROAR is a model-free instantiation of the general SMBO framework and works by firstly selecting new parameter configurations of the target algorithm uniformly at random which it iteratively compares to the current best configuration, known as the *incumbent*, using a form of aggressive racing.

This racing is applied by a generalisation of the SMBO frameworks *intensify* procedure, adapted by the SMAC developers to be able to configure algorithms to multiple problem instances, so as to solve the additional problem - introduced by multiple instances - of deciding which instances to use on each run. The original *intensify* procedure as defined in the SMBO framework is only responsible for: (i) determining how many evaluations with each candidate configuration are to be performed, and crucially, (ii) when confidence in a configuration is strong enough to accept it as the new incumbent configuration.

5.4.1.2 *Basic SMAC Procedure*

As hinted in the previous section, SMAC makes several generalisations which solve some of the problems relating to ROAR, namely: (i) use of models that allow for categorical parameters and (ii) use of models that allow for sets of problem instances to be used [59, 58]. The most commonly used model in SMAC is a random forest model, which supports these

generalisations, and are covered in section 5.4.1.3. Along with the updated *intensify* procedure from ROAR, a summary of the basic iterative SMAC procedure can then be defined as follows [60]:

1. Construct a random forest model to predict algorithm performance over instances
2. Use the constructed model to select new promising configurations
 - a) Uses the expected improvement criterion (EI) which combines predicted mean with prediction uncertainty (section 5.3.1)
 - b) Finds configuration with the highest expected improvement via local search
3. Compare each new configuration to the current 'best' configuration (the incumbent)
 - a) Uses the same aggressive racing strategy as for ROAR
 - b) Save all data generated from runs to construct a model in the next iteration

5.4.1.3 *Random Forest Models*

The default model used within SMAC is Random Forests, an ensemble learning method used as a standard model in machine learning for regression and classification [59]. From a basic standpoint, random forests are essentially collections of regression trees. Similar in nature to decision trees, regression trees make use of real values at their leaves instead of class labels - in the context of SMAC these real values represent the performance of the target algorithm [59]. Regression trees have been used in the past to model heuristic performance both in terms of runtime and quality of solutions and are known to perform well for categorical input data [59]. Random forests also share the benefit of being able to model performance as runtime and solution quality, but additionally, allowing the uncertainty of given predictions to be quantified [59]. In addition to these benefits, further benefits stated by Breiman and Cutler [15] make it clear why the developers of SMAC would choose random forests as their default model:

- They can be run efficiently on large data sets
- They can handle thousands of input variables without variable deletion
- They are able to give estimates of what variables are important in the classification
- They create internal and unbiased estimates of the generalisation error as the building of the forest progresses
- They possess an effective approach to estimating missing data and can additionally maintain accuracy even when a large proportion of the data is missing

- They have methods for balancing error in class population unbalanced datasets
- The forests generated can be saved and reused on other data
- Prototypes can be developed which give information regarding the relationship between variables and a generated classification
- They are able to compute proximities between pairs of cases to be used in the contexts of clustering, outlier detection or giving more interesting views of the data through scaling
- The capabilities of the previous point can be extended to unlabelled data
- They offer an experimental approach to detecting variable interactions

In terms of computational costs, random forests are also fast, Breiman and Cutler claim that with a dataset containing 500 cases and 100 variables, an implementation of a random forest managed to produce 100 trees in 11 minutes using an 800MHz processor [15]. They also state that the main memory requirements come from the storage of the datasets themselves along with 3 arrays each with the same dimensionality as the data but if proximities are calculated the storage requirements grow as the number of cases multiplied by the number of trees in the forest [15].

Part III

METHOD

CHAPTER 6: ALGORITHM PARAMETER CONFIGURATION

6.1 INTRODUCTION

Before conducting any experimentation using the metaheuristic approaches outlined in Chapter 4, it was first necessary to conduct parameter tuning processes on each approach. It is considered standard practice when using stochastic approaches such as metaheuristics to tune the parameters of the utilised approach in order to produce the best possible performance against a particular problem instance or set of instances being addressed. This tuning aimed to obtain configurations that provided the best possible performance for each algorithm - based on a measure of solution quality - against a specific problem at a specific dimensionality.

From the point of view of comparing the performances of the algorithms in our experiments, it would be unwise to leave our implementations in an untuned state as any findings would carry lower confidence since there would be no way of knowing if the implementations could have performed better. For example, if one approach was found to perform better on some problem at higher dimensionalities than some other approach we could not state this with high confidence since we do not know how much better each approach can become with appropriate tuning. In this scenario, the better approach may already be close to an optimally performing state whereas the other approach may have much performance to be gained if tuned. If both approaches are tuned, comparisons like this can more accurately reflect the full potential of each approach and as such lead to more trustworthy conclusions.

It is sometimes the case that a practitioner will tune their approaches against several different problem instances (possibly of different sizes) in order to generate an average set of parameter configurations that, it is hoped, will provide good performance against all similar instances. For our purposes, this type of tuning of our implementations would cause inconsistent performance across problems and their various dimensionalities. Similarly to the disadvantages of having the implementations remain untuned, here we would not know with any confidence to which problems and to which dimensionalities the approach has been tuned to address better than others.

To eliminate the potential issues discussed above we tuned our implementations on a per-instance-per-dimension basis; that is, the implementations would be tuned against each benchmark problem at each dimensionality independently so as to obtain the best performance possible from an approach against a single problem at one specific dimensionality. This would result in far fairer comparisons; if one approach is compared with another and is found to perform better at higher dimensionalities for a given problem, this result carries a higher confidence level since there is no tuning that would allow the worse approach to perform better. Furthermore, this approach to tuning allows us to more easily and more justifiably generalise any conclusions to higher levels such as, to other approaches of the same class or to problem types.

To reduce the computational cost of tuning each metaheuristic approach when considering each subsequent dimension x_i as a linear progression (e.g., $x_i \in \{2, 3, 4, \dots, n\}$) the dimensions to be considered instead follow a geometric progression, specifically $x_i = \lfloor x_{i-1} \sqrt{2} \rfloor$, where $x_i \in \mathbb{N} : \mathbb{N} \leq n$. To illustrate this point, considering we make use of our 17 standard benchmark functions \mathcal{F} with cardinality $|\mathcal{F}|$ (see Chapter 3, Section 3.4) up to a maximum dimensionality of 1500, we would need to carry out $d - 1 \times |\mathcal{F}| = 25,483$ separate tuning procedures if we were to consider a linear increment of 1 to the dimensionality of each problem ¹. This problem would be compounded further when we consider that multiple repeats of each of these would have to be performed to offset the stochasticity inherent in our selected APT (see Section 6.2 and Chapter 5, Section 5.4 for details) - this would require 127,415 tuning procedures if we conduct each tuning procedure five times. Conversely, using a geometric progression of dimensionality results in a set of dimensions d_{gp} with cardinality $|d_{gp}| = 20$ meaning we only need to perform a more manageable $|d_{gp}| \times |\mathcal{F}| = 340$ tuning procedures, equating to 1700 when conducting each procedure five times.

6.2 PARAMETER TUNING METHOD

We tuned our metaheuristic implementations using the state-of-the-art automatic parameter tuner SMAC. For each metaheuristic approach, the tuning process consisted of the following steps:

1. Identified tunable parameters and appropriate value ranges. These were specified in SMAC parameter files to be used by the SMAC tuning procedures.

¹ a dimensionality of 1 is not considered so we only consider 1499 out of the 5000 dimensionalities

2. Compiled a set of SMAC instance files each containing an algorithm recognisable string identifying a benchmark function and a specific dimensionality to tune against.
3. Compiled a set of scenario files to be used by SMAC which among other things addresses the instance file to be used and the path to an algorithms executable code.
4. Conducted five independent tuning procedures for each scenario file in order to offset the possibility of obtaining less optimal parameter configurations due to the stochastic noise that is invariably part of random processes such as those used by SMAC
5. Obtained tuned approach performance data by conducting a series of repeated runs using each of the five independent configuration sets against their corresponding problem-dimension pairing.
6. Using the tuned performance data above, produced a 'final' parameter configuration set containing the best configurations from each of the five independent problem-dimension configurations

Each of these steps will be covered in more detail in the following subsections.

6.2.1 *Parameter Identification and Range Selection*

For each metaheuristic approach implemented, a suitable set of SMAC parameter files needed to be produced. Identification of the 'tunable' parameters was straightforward as any parameter of an approach can be considered tunable when its value can affect the behaviour of the algorithm [121] and thus its performance outcome. Since all user specified variable parameters used during an algorithms execution have an affect on algorithm behaviour these are the ones tuned to obtain better performance. The parameter file content for each of our implementations showing the identified parameters and ranges will be made available via the university's digital repository in due course.

On the other hand, selection of an appropriate range for these parameters was not as straightforward. Care needed to be taken when selecting the ranges of the parameters used in conjunction within the algorithm; as SMAC, having no sense of correspondence between parameters, can potentially return an infeasible value for one parameter in respect to the other. One such example is that of the parameters used in the GA tournament selection, *population size* and *tournament size*. When originally tuned by SMAC, both parameters were stated as integer types each representing a number of individuals. Since SMAC has no knowledge of each use within the algorithm and it will always return the best configuration found

during tuning - even if that means the operator in question become equivalent to another, we found in early tuning attempts that tournament size (of type integer) would often return with a larger value than population size suggesting a tournament consisting of the whole population. Of course, given a GA tournament with replacement consisting of the whole population, it becomes increasingly likely that the best solution in the population will take part in all tournaments - although small there is still a non-zero probability that this will not be the case on every iteration. This was overcome by making tournament size a parameter representing a proportion of the main population, which allowed us to prevent this problem for any population size SMAC produced.

6.2.2 *Conducting the Tuning Procedures*

It is suggested that when using SMAC that multiple independent tuning procedures are carried out to better ensure a high quality parameter configuration output due to the inherent stochasticity used within the SMAC implementation. This being considered, five independent tuning procedures were carried out per SMAC scenario/instance file (problem-dimension pair) for each metaheuristic approach implemented. We make use of 17 standard benchmark functions (Chapter 3) and a rounded geometric progression of problem dimensionality - resulting in 20 different dimensionalities ranging from [0, 1448]. Across all problem-dimension pairs for each approach, we carried out a total of $17 \times 20 \times 5 = 1700$ independent SMAC tuning processes run by utilising the departments CPU cluster. It was determined that the dimensionality increase of each problem being addressed by our approaches was to follow a geometric progression. The reason for this was mainly to reduce computational cost, since even tuning with a linear progression using dimensionality increases of size 10 in the range [0, 1500] would need 12,750 independent tunings.

6.2.3 *Generation of Tuned Performance Data*

Having obtained a set of independent parameter configurations, five for each problem-dimension pairing, the algorithms were initialised with each configuration in order to generate the required performance data. Each algorithm was repeated for 20 independent runs with each parameter configuration relating to the current benchmark problem with a fixed maximum function evaluation budget of 50,000. In practice, it was far simpler to compile the tuning configurations obtained from SMAC into five comma separated value (CSV) files per problem with each line representing one of the five tuned parameter configurations at a given

dimensionality. These files would then be passed to the algorithm on execution which would carry out 20 independent repeats for each configuration. Fig. 6.1 shows an example of one of these compiled CSV files.

```

dim=2, c1='2.0', c2='2.0', inertiaRate='0.7', popSize='50'
dim=3, c1='2.0', c2='2.0', inertiaRate='0.7', popSize='50'
dim=4, c1='1.1278175019901626', c2='1.4833886999205852', inertiaRate='0.4401729629396143', popSize='165'
dim=6, c1='1.1262434695182078', c2='1.2571467889597245', inertiaRate='0.7353011476808933', popSize='146'
dim=8, c1='1.005954467241296', c2='1.0160879984589266', inertiaRate='0.7582652478158005', popSize='108'
dim=11, c1='1.2062138052663433', c2='1.1213648961513178', inertiaRate='0.7171198456046894', popSize='221'
dim=16, c1='1.0489475842909564', c2='1.2809857205555781', inertiaRate='0.6766207861718936', popSize='182'
dim=23, c1='1.1062310957599955', c2='1.4050099636591442', inertiaRate='0.6677417805705328', popSize='236'
dim=32, c1='1.1678393383411656', c2='1.0314616634985572', inertiaRate='0.72912131489805', popSize='223'
dim=45, c1='1.9267431996625533', c2='1.1399060916013448', inertiaRate='0.6763259760109199', popSize='182'
dim=64, c1='1.1768606688287502', c2='1.0669707080618969', inertiaRate='0.7029910688358594', popSize='220'
dim=91, c1='1.8342116750479402', c2='1.0010597879025378', inertiaRate='0.6846094516554043', popSize='207'
dim=128, c1='1.4988992860205215', c2='1.0018270758408594', inertiaRate='0.6769109230043653', popSize='239'
dim=181, c1='1.5236922574245115', c2='1.0028616429469575', inertiaRate='0.5961094011782737', popSize='236'
dim=256, c1='1.6325529288481204', c2='1.0430337303902295', inertiaRate='0.672410269291469', popSize='222'
dim=362, c1='1.6168078198323692', c2='1.0191106976778248', inertiaRate='0.5904451836113005', popSize='248'
dim=512, c1='1.7186775640269385', c2='1.1847251742530893', inertiaRate='0.6013429922694741', popSize='227'
dim=724, c1='1.691634360228445', c2='1.0156918360808054', inertiaRate='0.6243383375260789', popSize='250'
dim=1024, c1='1.9002598650777534', c2='1.0130073290318302', inertiaRate='0.5759667214321772', popSize='245'
dim=1448, c1='1.9049669817841566', c2='1.0584854069921459', inertiaRate='0.5782790957502374', popSize='220'

```

Figure 6.1: Example of a parameter configuration file compiled from the individual problem-dimension configurations obtained from SMAC

For each dimension of the current problem, a separate CSV file was produced as output at different stages of the search process - in 5000 evaluation increments - containing the objective value of the best solution at that stage of the search for each of the 20 independent repeats carried out. Therefore, the output CSV files for 50,000 function evaluations would represent final solutions to the problem instances given our budget.

6.2.4 Obtaining a Final Parameter Configuration Set

Given the output CSV files at 50,000 function evaluations from our implementations described above, the median of the 20 independent repeats was used from each to determine the best parameter configuration for each dimension-problem pairing. The median was used in this case as the distribution of the data contained in most of the output CSV files was not Normally distributed; that is, it tended to be either very left or right skewed. Therefore, in the majority of cases the mean of the repeat data would not accurately convey the central tendency of the performance data from an algorithm for a given dimension-problem pair. This however is not unusual for data obtained from stochastic approaches such as metaheuristics.

To calculate the median from our algorithm outputs and subsequently compile a 'final' set of parameter configuration files for use in our experiments an R script was written. This would read in the set of five CSV parameter configuration files for each problem - from our five independent SMAC tuning procedures - and for each dimension of the problem, generate the median solution performance from the output CSV files associated with each of the five SMAC tuning procedures. Where a median was found to be lowest out of the five for a given

dimension, the corresponding line of the associated parameter configuration CSV file is copied into a new file.

6.3 IMPROVING THE GENERALITY OF EMPIRICAL RESULTS AND OBSERVATIONS

In order to better generalise my observations, results and conclusions derived from our tuned metaheuristic algorithms as well as to provide stronger evidence that observed relationships are not simply produced from biases that might be inherent to SMAC; a repetition of the tuning procedure was carried out using another popular automatic parameter tuning approach. The selected alternative, Irace [79], was used to tune our algorithms on a representative subset (7 functions) of the benchmark function suite. Although other widely used APTs exist, such as ParamILS, Irace is capable of handling real-valued parameters as well as discrete parameters, whereas ParamILS is designed only for use with discrete parameters.

Irace, developed by López-Ibáñez et al. [79] makes use of an iterated racing approach to algorithm configuration and consists of three distinct steps: (i) sampling of new configurations based on a particular distribution, (ii) selects the most promising configurations from the samples through the use of racing and (iii) updating the distribution used for sampling in order to bias further sampling towards promising configurations [79]. A racing procedure typically starts with a finite number of candidate configurations. At each step of a race, the candidate configurations are evaluated on a single problem instance. After a number of steps, those configurations which are shown to perform statistically worse are rejected and the race continues with those configurations which survive [79].

The Irace implementation used (Version 3.3) was produced by the developers as a CRAN library for the statistics package *R*. An R script, all required input files and command line scripts (.bat files for Microsoft Windows command line and a shell script (.sh) for UNIX-based systems) were constructed using templates provided as part of the source files of same CRAN library - available from [78].

As with SMAC, the maximum number of single algorithm executions allowed during tuning was set to 1000, in keeping with the execution budget afforded for SMAC tuning processes, and all other required settings were kept as default. Also, as Irace also has a stochastic element, five independent tuning procedures were carried out for each Algorithm-

Problem-Dimension triplet in order to better insure against unsuccessful tuning procedures (see Section 6.2.2).

In an effort to make the comparison between the two tuning approaches as fair as possible, the allowable parameters ranges used in tuning each algorithm were identical to those used by SMAC. This removes the potential of any difference in tuning performance observed between the two approaches being down simply to one tuning method having to search over larger parameter spaces than the other, i.e., if the difference between the lower and upper bounds are larger.

6.4 RESULT OF PARAMETER TUNING WITH SMAC

For each of the algorithms tuned with SMAC, the median performance over all function-dimension pairings (compared to the untuned median performance) was improved significantly; particularly when searching at higher dimensionalities. The line plot in Fig. 6.2 shows the median performance results - aggregated from 20 repeats - of our DE implementation against the Ackley function (Chapter 3) over all dimensionalities. The red line depicts the median performances of the untuned algorithm and the black line is the median performances of the algorithm using the final configuration set. For this particular result, no significant difference in the performance between the tuned and untuned versions was observed up to 91 dimensions, however, from the 128 dimensional problem instance, the median performance difference increases noticeably. Similar trends were found true for each of the algorithms implemented, where the tuned version was found to perform significantly better on all benchmark functions.

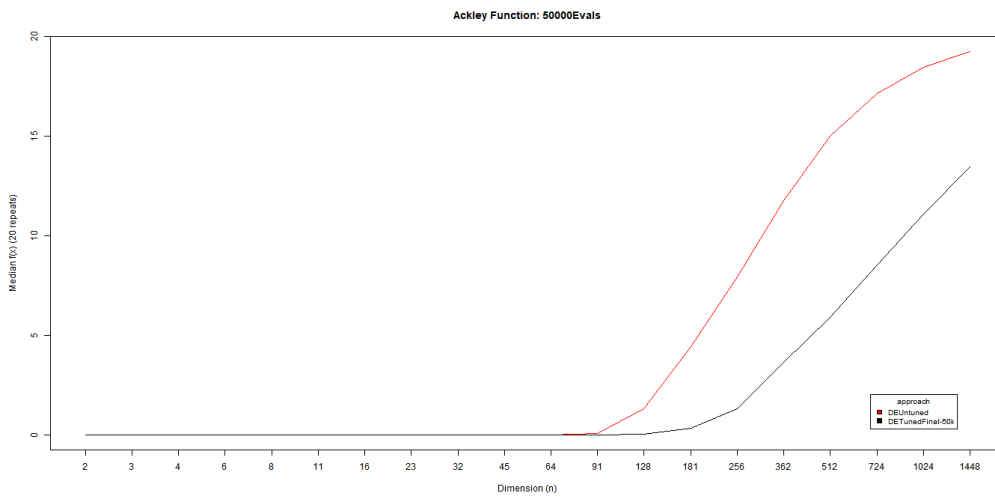


Figure 6.2: Result Obtained from DE Implementation (tuning at 50,000 function evaluations) vs. Untuned Version

Part IV

RESULTS AND CONCLUSIONS

CHAPTER 7: EMPIRICAL RESULTS

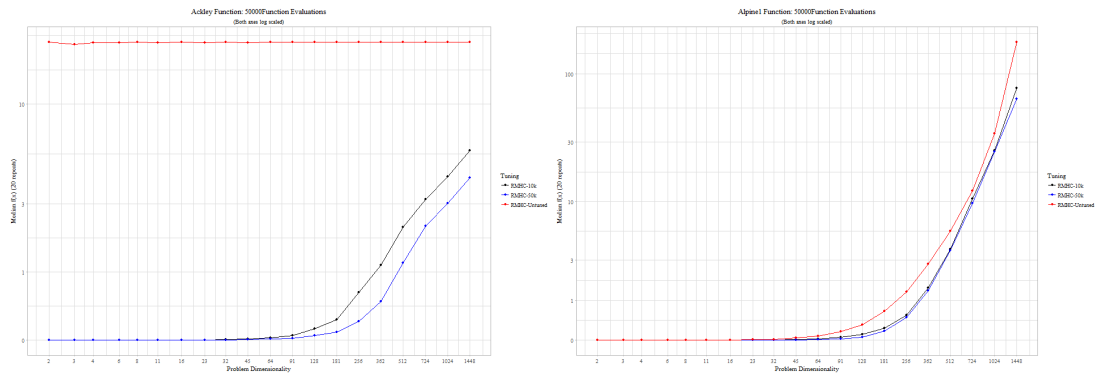
7.1 ALGORITHM TUNING RESULTS

In this section we outline the results obtained through our selected tuning procedure as compared to the same algorithms in their untuned states as well as present findings from our analysis related to parameter tuning at differing function evaluation budgets.

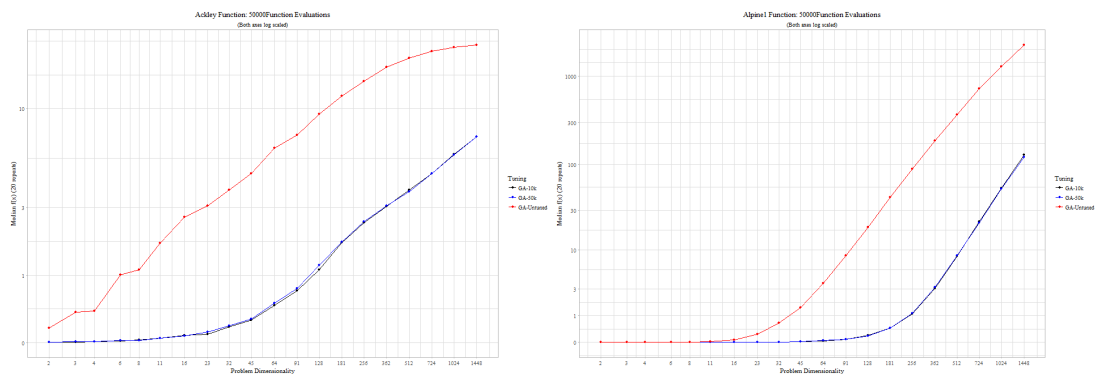
This first set of results represents a ‘sanity-check’ for further experimentation, specifically, I compare the metaheuristic approaches in the context of the best tuning found by the selected parameter tuner (SMAC), therefore, it was sensible to ensure that the performance of the algorithms was indeed being improved beyond the baseline ‘untuned’ equivalent algorithms. Non-improvement may have indicated a problem/bug in the tuner which would have to be remedied before any data generation and comparison could be performed.

Before continuing on to presenting the tuning results, I must clarify what is meant by an ‘untuned’ algorithm in this context. This is where an algorithm has been initialised with a parameter configuration obtained either from: (i) common settings found in literature sources or (ii) where no suitable settings could be found in the literature based on specific design choices - such as, when using a different parameter representation to work better with SMAC e.g., GA Tournament Size vs Tournament Size Percent (See Chapter 4 Section 4.3.4.3) - a sensible value for the parameter was used instead. In short, ‘untuned’ here refers to an algorithm that has a single set of parameter settings that have not been selected in the context of our benchmark suite and therefore, are not necessarily appropriate or effective at obtaining high performance on this problem set and at larger problem scales.

Here we present several plots showing that SMAC indeed effectively tuned all algorithms often well beyond their baseline untuned performance - however, the full plot listing showing the tuning results for all problems and algorithms can be found in Appendix D section D.1.



(a) RMHC Tuned to Ackley Function at 50,000 GA Evaluations (b) RMHC Tuned to Alpine n.1 Function at 50,000 GA Evaluations



(c) SSGA Tuned to Ackley Function at 50,000 GA Evaluations (d) SSGA Tuned to Alpine n.1 Function at 50,000 GA Evaluations

Figure 7.1: Examples of Successful Tunings of RMHC and SSGA Using SMAC on Ackleys Function and Alpine n.1

7.1.1 Comparing SMAC and Irace Tuning Performance

To better generalise my observations, results and conclusions and generally provide stronger evidence that observed relationships are not simply produced from biases present in SMAC; a repetition of the tuning procedure, over a representative subset of my benchmark set (7 functions), was carried out using the automatic parameter tuning approach, Irace (See Chapter 6, Section 6.3 for a discussion on this method). To provide a fair comparison, Irace was provided with the same initial conditions as were true for SMAC where possible¹, specifically these include:

¹ Where these conditions were not individual to a particular approach and implementation thereof

- The tuning algorithms were both afforded a maximum budget of 1000 algorithm executions to use during the tuning procedures. This was kept constant as the problem size was increased ²
- The ranges for each algorithms parameter set were consistent between tuning methods
- The exact same target algorithm files and benchmark function suite file (utilised as executable Java .jar files) were used by both SMAC and Irace
- As with SMAC, the final set of configurations for each algorithm were taken as the best of five independent tuning procedures (determined individually for each problem-dimension tuple)

The results of tuning each algorithm using Irace in comparison to tuning with SMAC are presented in Appendix D, section D.2.

From the results, it can be easily verified that, in general, Irace found parameter configurations that produce results which follow a similar trend to those of the algorithms tuned with SMAC. However, the similarities between the two tuning methods is limited only to this single observation. Over the representative subset of functions utilised by the Irace experiments, SMAC was shown consistently to have produced far better parameter configurations which again, in general, produced a far more stable set of performance values (i.e., little variance between independent repeats of an algorithm, illustrated by the interquartile range bars on the plots in Appendix D.2) than the corresponding performances obtained with Irace configurations. Given that this is the case, and also that the initial conditions were kept consistent between SMAC and Irace where possible, then SMAC has shown itself to be the more reliable, consistent and robust to noise of the two approaches, at least against this particular suite of problem instances. Additionally, these results also show SMAC as being less prone to converging to locally optimal parameter settings - prematurely or otherwise - as a result of its inherent stochasticity.

² In general for SMAC, the search for configurations converges before the maximum number of algorithm executions is reached - even for the maximum problem size of 1448 - so not increasing this budget as problem size increases should not present a problem in terms of observing how an algorithms real performance decreases as problem size increases. Further, for SMAC at least, the configurations returned should therefore be reasonably close to optimal for that tuning method.

7.2 FUNCTION EVALUATION BUDGET COMPARISON

Here we presents several findings related to the tuning of our metaheuristic implementations over two distinct evaluation budgets.

Plots showing the median objective value (raw) returned from 20 algorithm repeats vs. problem dimensionality for each of the benchmark problems and metaheuristics is provided in Appendix A. For each of these plots, each line represents an approach either tuned using a 10,000 or 50,000 function evaluation budget and error bars represent the lower and upper quartiles of the data at each point - Q1 and Q3 respectively. Additionally, the plots have been generated at 5000 evaluation increments for the same set of algorithm repeats, i.e., output was generated for each algorithm repeat when the number of evaluations: $(\text{eval} \bmod 5000) == 0$. For brevity, we will form our discussions here in terms of only one or two Functions from our suite (Chapter 3) and present select plots from Appendix A as appropriate.

Additionally, Appendix B provides supplementary descriptive statistics tables for each of the plots in Appendix A.

7.2.1 *Effects Observed when Tuning with 10,000 and 50,000 Function Evaluation Budgets*

We compared our implementations on their ability to produce quality solutions when tuned with two different function evaluation budgets. Each set of resulting parameter configurations generated using these two budget schemes were then applied against our benchmark function suite - running with a budget of 50,000 objective function evaluations. The goal of this experiment was to see whether algorithms should be tuned at the same number of evaluations of which it would be expected to perform - hereafter referred to as the *target evaluations* - or if, in some instances, it were possible to obtain better or equal performance with a smaller evaluation budget.

In general, we found that, when configuring the parameters for the metaheuristics, it is better to use an evaluation budget at least equal to the target evaluations when solving problems of higher dimensionality. Specifically, at lower dimensionalities (typically < 91D), using either tuning evaluation budget will produce comparable median solution quality; and at higher dimensionalities (typically > 91D), the performance obtained from each scheme begins to diverge in favour of tuning at the higher evaluation budget - here 50,000 evaluations.

However, in a relatively high number of cases no such divergence occurs, i.e., there is no difference when tuning algorithms at the target evaluation budget or a much lower budget. Table 7.1 summarises this divergent behaviour as observed in the comparison plots in Appendix A. In this table, *no-div* indicates where very little or no divergence of solution quality performance is found between tuning evaluation budget schemes, *UM* and *MM* indicates whether a function is uni-modal or multi-modal respectively and *S* and *NS* indicate whether or not a function is separable (either multiplicatively or additively). Additionally, PSO is intentionally omitted from the table as this will be treated separately.

Table 7.1: Plot summary of problem dimensionalities where algorithm performance between 10k and 50k evaluations diverge in favour of tuning at 50k function evaluations

Function/Alg.	CMA-ES	DE	GA	RMHC	SA
Ackley (MM,NS)	no-div	91	no-div	91	no-div
Alpine n.1 (MM,S)	181	91	no-div	1448	no-div
Bent Cigar (UM,S)	no-div	91	no-div	32	45
Brown (UM,NS)	no-div	1448	no-div	256	128
Chung-Reynolds (UM,PS)	1448	128	no-div	181	91
Deflected Corrugated Spring (MM,NS)	512	91	no-div	no-div	1448
Exponential (MM,NS)	1448	181	1448	512	no-div
Griewank (MM,NS)	no-div	91	no-div	181	64
Inverted Cosine Wave (MM,NS)	no-div	91	no-div	no-div	no-div
Levy (MM,NS)	181	128	no-div	362	128
Qing (MM,S)	no-div	64	no-div	64	64
Rastrigin (MM,NS)	181	16	no-div	181	181
Rosenbrock (UM,NS)	no-div	64	no-div	362	128
Schwefel (MM,S)	n/a	16	no-div	no-div	no-div
Sphere (UM,S)	1448	128	no-div	362	91
Sum of Different Powers (UM,S)	no-div	no-div	no-div	724	no-div
Sum Squares (UM,S)	no-div	128	no-div	128	64
Mean	771	178	1448	349	221
No Divergence Count	9 of 17	1 of 17	16 of 17	3 of 17	6 of 17

The first thing to notice in Table 7.1 is that for 16 out of the 17 benchmark functions, when using our GA implementation no significant difference can be observed between tuning at 10,000 or 50,000 evaluations. This finding suggests that when using our GA (steady-state) for continuous function optimisation that the approach may be tuned at a significantly lower number of evaluations than the required target evaluation budget. Since automatic parameter tuning, despite not being an exhaustive process, can still be a computationally costly operation - and will likely require a lot of wall clock time to complete - this finding potentially eases this restriction that may prohibit APT for other approaches, allowing GA to be tuned faster. However, one caveat is that this in no way means that the target evaluation budget can also be lowered; looking at the the plots in Appendix A it is clear to see that the algorithm has not yet converged by 10,000 evaluations for the vast majority of our benchmark suite, as progress towards better solutions continues up until our 50,000 target evaluation budget - and in some cases appears as though this progress would continue given a larger budget. Fig. 7.2 shows our GAs performance on Ackley's function for 1000, 10000, 30000 and 50000 evaluations illustrating this observation more clearly.

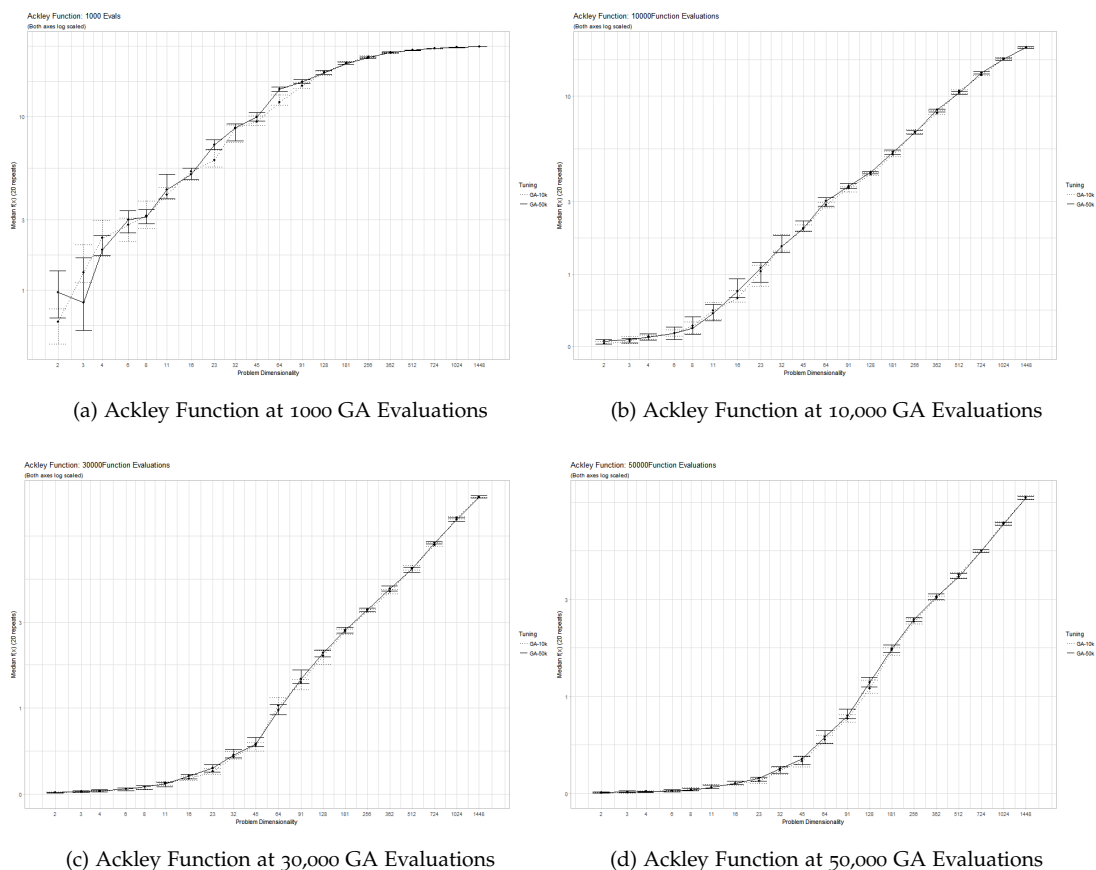
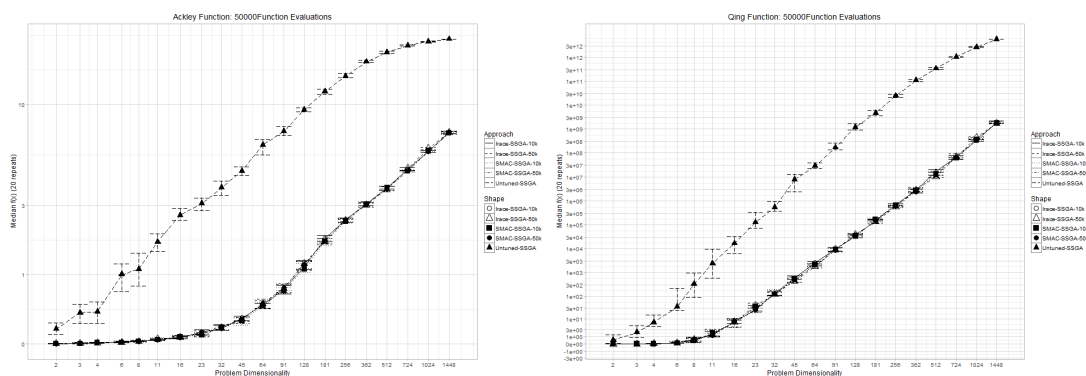


Figure 7.2: Example of GA Continuing to Progress Beyond 10,000 Function Evaluations for Ackley's Function

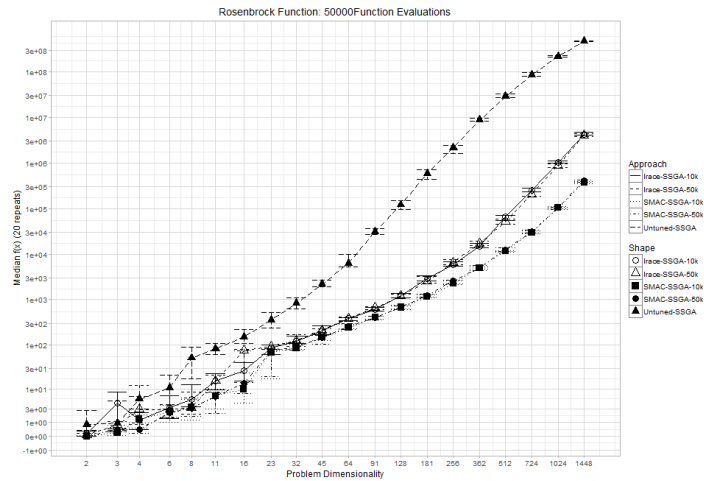
The performance results for SSGA tuned by Irace, across a representative subset of the function suite, provides further evidence for the generality of this result. As with the algorithm tuned by SMAC, the performances for both the 10,000 evaluation and 50,000 evaluation tuning budget strategies are for all intents and purposes identical, with no significant divergence between the results. Fig. 7.4 illustrates this observation for both Ackley and Qing functions.



(a) Ackley Function at 50,000 SSGA Function Evaluations (b) Qing Function at 50,000 SSGA Function Evaluations

Figure 7.3: SSGA Performance Results Comparison (Irace Tuned) Example - 10,000 and 50,000 Tuning Evaluation Budgets

Moreover, with the exception of the Rosenbrock function, the SSGA tuned with Irace produces performances that are almost identical to those of the SSGA algorithm tuned using SMAC across the representative function subset used for the Irace experiments. For Rosenbrock, we find that the configurations found by SMAC slightly outperform those of Irace. Despite this performance improvement only being ‘slight’, the interquartile ranges for each tuning method (SMAC vs. Irace) are very small and diverge cleanly with no overlap.



(a) Rosenbrock Function at 50,000 SSGA Function Evaluations

Figure 7.4: SSGA: Divergence of Performance between 10,000 and 50,000 Tuning Evaluation Budgets of SMAC and Irace

The Irace results for the remainder of the benchmark function subset can be found in Appendix D Section D.2.

The next thing we can observe from Table 7.1 is that DE appears to be the most sensitive - of all the implementations - to the number of evaluations used to tune the approach. Further, when compared with the other approaches not yet discussed, DE requires to be tuned at the target evaluation budget of 50,000 evaluations far sooner, in terms of problem size - with a mean point of divergence occurring at 178 (if we consider browns function an outlier this is decreased further to 93).

As mentioned previously, we deal with PSO here separately as it behaves quite differently from the other approaches discussed so far. What can be observed in general here is not simply a divergence of performance at a certain dimensionality when using different tuning evaluation budgets or even observing no such divergence at all. For around $\frac{2}{3}$ of our function suite, it can be seen that no performance deviation occurs between budgets at lower dimensions, but after a significant divergence of performance at mid ranged dimensions when using different budgets, a re-convergence occurs in higher dimensions that is here interpreted as there being no difference between using lower or higher tuning budgets. Also, against many of these functions, this divergence and re-convergence of performance can occur several times throughout an experiment. This effect can be seen clearly for Alpine no.1 function shown in Fig. 7.5.

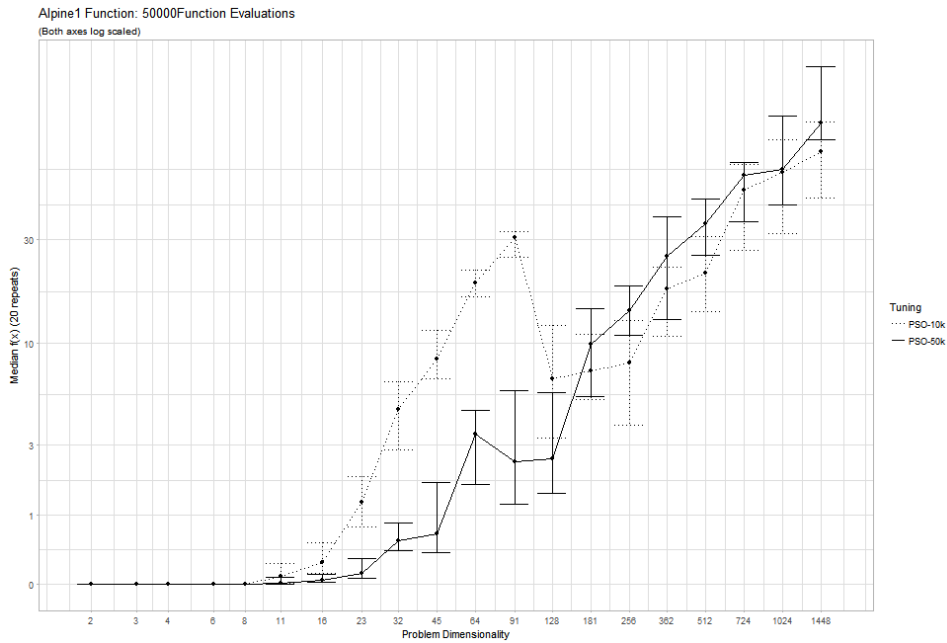
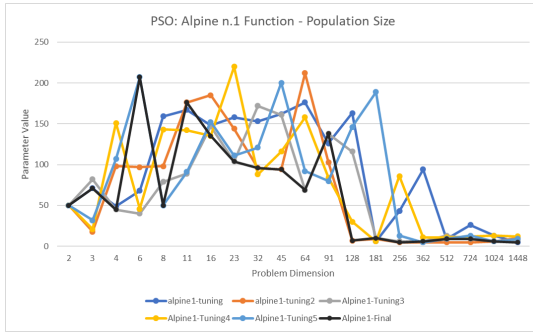
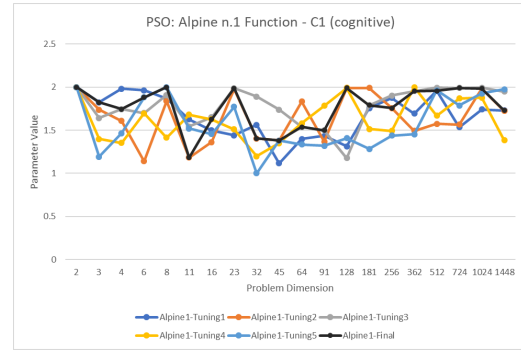


Figure 7.5: Example showing divergence between parameter tuning evaluation budgets for PSO (Alpine n.1 Function) at dimensions 11-16 and re-convergence occurring at dimension 128

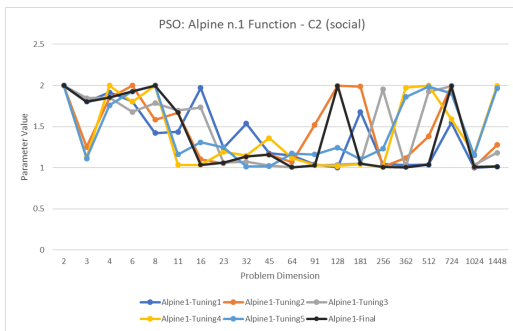
Through some analysis of the original five parameter configuration set returned for PSO, we feel that this behaviour can be attributed to a dramatic change in population size: (i) a large increase from 10s of individuals to 100s when the lines representing tuned evaluation budgets diverge, and (ii) a large decrease in population size causing re-convergence of the budget strategies. Fig. 7.6, shows plots of each parameter used to tune our PSO implementation, where each line represents the parameter values from one of our original five configuration sets tuned with a 10,000 evaluation budget (from which we selected our final configuration). We have also included in each a line representing the parameter values of the final parameter configuration. We make use of the 10,000 budget configuration files as in Fig. 7.5 the 50,000 evaluation budget line is far more stable and so is less likely to include any ‘strange’ parameter settings around within the dimensionalities of interest.



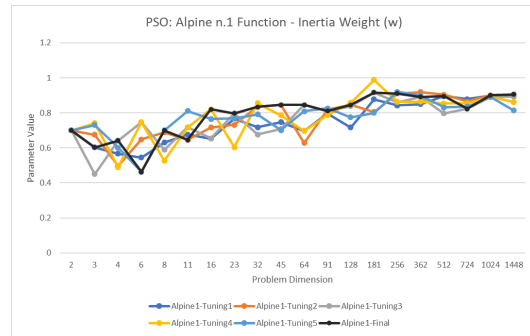
(a) PSO: Plot of Population Sizes from the Five Original Parameter Configuration Sets and the Final Configuration Set



(b) PSO: Plot of c_1 (Cognitive Component) from the Five Original Parameter Configuration Sets and the Final Configuration Set



(c) PSO: Plot of c_2 (Social Component) from the Five Original Parameter Configuration Sets and the Final Configuration Set



(d) PSO: Plot of Inertia Weight ω from the Five Original Parameter Configuration Sets and the Final Configuration Set

Figure 7.6: Plots showing individual values for each tuned parameter from our original 5 parameter configurations sets obtained from SMAC

From the plots in Fig. 7.6, it is clear to see in the population size plot that a dramatic increase of population size corresponds to a similarly dramatic decrease in solution quality obtained from PSO at the same dimensionality (Fig. 7.5). The reverse is also true when the lines in Fig. 7.5 re-converge at a problem dimensionality of 128, in that a corresponding increase of the solution quality occurs. This effect is not only limited to the Alpine n.1 function, the remaining function plots in Appendix A in which we observe this effect³ also show similar reliance on the population size parameter. Clearly, there is a threshold above which the presence of a large population presents a problem to our PSO implementation and conversely, below this threshold there appears to be no ideal value for population size at all.

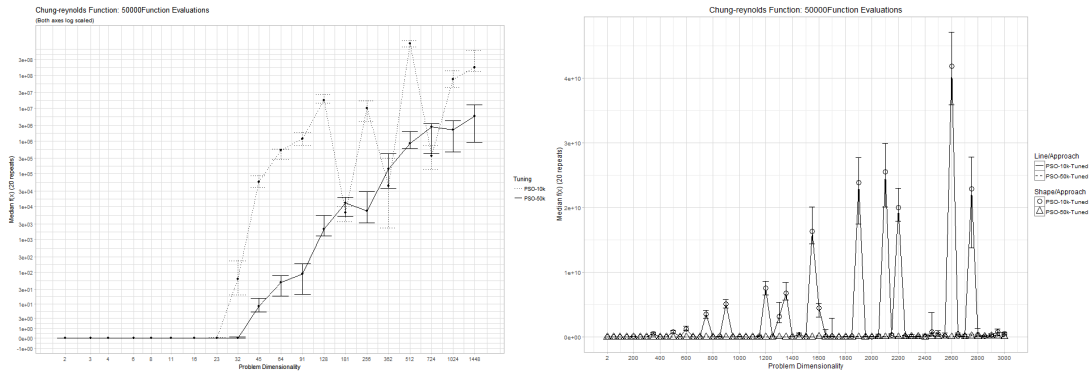
³ Ackley, Bent Cigar, Brown, Chung-Reynolds, Exponential (however less pronounced), Griewank, Inverted Cosine Wave, Levy, Qing, Rastrigin, Sphere (less pronounced) and Sum Squares

Also from the plots in Fig. 7.6, it can be seen that an apparent negative correlation between population size and the social acceleration coefficient (c_2) exists. A suggested explanation of this effect could be due to the sharing of components between particles i.e., in cases where the population size is smaller, the social coefficient may have to be larger in order to maintain a consistent level of exploitation.

I felt that it would be worth exploring this behaviour further through an extension of the tuning experiment at a larger scale, specifically by: (i) doubling the maximum number of dimensions of each problem to 3000 dimensions in order to see if this behaviour continues beyond the original 1500-D limit, and, (ii) increasing the dimensionality of the problem linearly - by increments of 50 - in place of the usual geometric progression - to observe how this behaviour presents itself at smaller dimensional granularities. The results of this experiment can be found in Appendix F Section F.1.

The results show that the effect observed for the functions in my original PSO experiment is also represented in this extended version for the same set of functions, with the same $\frac{1}{3}$ of functions exhibiting normal behaviour i.e., performance divergence with very little to no re-convergence. Further, the effect does present itself as being cyclic in nature - even at these extended problem scales and smaller granularity. As an example, in Fig. 7.7, we see that the performances between tuning budget strategies diverge and re-converge multiple times, and that each consecutive divergence increases in magnitude as the dimensionality increases. ⁴

⁴ For some of the other functions included in Appendix F that display this behaviour, such as: Ackleys Function and Alpine Function no.1, the high frequency of the divergences coupled with the similar baseline performances for both tuning budget schemes does not facilitate clear comparison, however, on close examination the same behaviour can also be observed for these functions. The majority of the comparisons however provide clear indications of where divergence is occurring.



(a) PSO Performance Against Chung-Reynolds Function Over a Geometric Progression of Dimensionalities from 2 to 1448 Dimensions (Original Experiment) after 50,000 Evaluations
 (b) PSO Performance Against Chung-Reynolds Function Over a Linear Progression of Dimensionalities in Increments of 50 Dimensions (Extended Experiment) after 50,000 Evaluations

Figure 7.7: Comparison of Performance (Median Solution Quality of 20 Independent Repeats) Against the Chung-Reynolds Function Between Original and Extended Scale PSO Experiments

Of course, this exponential decrease in performance could easily be attributed to the increased difficulty of finding good solutions within exponentially increasing search space volumes (curse-of-dimensionality). However, the root cause for the occurrence of this cycle of behaviour is still largely unclear. One possibility is that where the performance divergence occurs, the underlying structure of the function at that particular dimensionality becomes such that the automatic parameter tuner finds it more difficult in earlier stages of the search to identify a parameter configuration providing a significant improvement in performance in the shorter search scale of 10,000 function evaluations. As suggested of the original PSO experiment, from SMACs perspective, no configuration is found to be more effective than any other. In effect, this would mean that the same behaviour is largely unseen for the 50,000 budget strategy, as we observe, where PSO is given enough time during tuning to explore the potential of the various configurations provided to it.

Finally, the performance results of PSO when using the parameter configurations obtained through Irace do not clearly show that the observed effect can be further generalised to alternative tuning methods other than SMAC; the corresponding PSO performances when tuned using Irace (Appendix D, Section D.2) are too erratic to reliably differentiate between noise produced by more poorly performing parameter configurations and actual instances of the same effect discussed thus far. However, comparing the noise between different scales of the same problem for PSO to the other algorithms configured using Irace, we see that the configurations for PSO produces the most inconsistent performance data. It may therefore

be hypothesised that what we are observing with the Irace results for PSO are exaggerated instances of the original effect, caused by the increased difficulty and inconsistently in finding reliable configurations when compared with SMAC.

7.3 PERFORMANCE COMPARISON BETWEEN ALGORITHM IMPLEMENTATIONS

7.3.1 *Random Mutation Hill Climbing (RMHC)*

The first, rather unexpected, observation that can be made from the algorithm comparison plots of Appendix C Sections C.1.4 and C.1.4 is that our Random Mutation Hill Climbing (RMHC) implementation regularly performs competitively with far more complex approaches. It places third - albeit jointly in some cases - in 6 of the 17 benchmark functions, second in 4 of the functions and first position in 7 functions (beating every other approach outright in 3 out of the 7 - including the far more complex sep-CMA-ES).

This observation, although relatively trivial, actually provides some evidence contrary to the widely held belief that hill climbing algorithms, particularly of the most basic stochastic variations, are not viewed as very useful when compared to more powerful evolutionary algorithms [115]. Rosete-Suarez et al. in [115] describe this belief as being a bias likely introduced due to: (i) a lack of rigour in terms of defining the ‘true’ complexity of problems and (ii) the situation that it is rare to find comparisons against basic stochastic hill climbing variations. In fact, Rosete-Suarez et al. also offer a multitude of examples from the literature where simple stochastic hill climbing variants outperform some Evolutionary Algorithms in various contexts [115]. This, conclusion should not be new to the search and optimisation community - we know from the NFL theorem [146] that all search methods perform equally well when averaged over all possible search spaces; the obvious consequence of which is Hill Climbing Algorithms, as a valid search strategy, must have a niche set of functions for which it is ideally suited [115].

When compared to the Irace performance results, we see that Irace - although being afforded similar initial conditions (e.g. total number of algorithm executions) - was not as successful in tuning RMHC as SMAC. Therefore, depending on whether SMAC or Irace was used for the tuning procedure, the performance results will not likely be consistent with the observations above; especially if such comparisons are made between the remaining metaheuristics that have also been tuned through the use of Irace.

7.3.2 *Changes to Rank Ordering of Performance at Larger Scales*

Here we make use of tables of ranks and ranking plots to investigate the scalability of our metaheuristic implementations given a fixed budget of function evaluations (50,000) as problem scale increases. The full set of results can be found in Appendix C, which includes plots showing the raw performance medians with and without IQR error bars (Sections Sections C.1.4 and C.1.4) as well as providing the rank ordering plots and tables (Sections C.2 and C.3) respectively. However, for the purposes of discussion, selected results will be presented in this section. Additionally, in order to gain some insight into the general robustness and effectiveness of the metaheuristic implementations over our benchmark suite as dimensionality increases, and in effect, an idea of the comparative performance when averaged over functions possessing features of interest for this study i.e., multi-modality, uni-modality, full separability and full non-separability; plots presenting the mean aggregate rank ordering of each algorithm over several subsets of the benchmark suite will be discussed. Due to space restrictions, the tables corresponding to these plots can be found in Appendix C Section C.3.2. Additionally, and in order to more easily discuss the general trends in terms of scalability, results in this section are often presented as pooled mean rankings over dimensionality ranges rather than individual dimensions e.g., for: $\{2, 3, 4, 6, 8\}D$, $\{11, 16, 23, 32, 45\}D$ and so on.

For both of these types of plots, care should be taken in the interpretation of the aggregate ranks as there are a couple of reasons an algorithm can appear to improve as the scale of problems increase. Specifically, an improvement of the overall movement of rank over time of a given algorithm can be attributed to: (i) the rank changing in relation to the worsening of other algorithms ranks where the performance at larger scale is only better comparatively, and (ii) a genuine improvement in the performance of the algorithm resulting in better ranking. Depending on the situation however, the interpretation may matter more or less, for example, when simply determining which algorithms scale better than others in comparative terms, whether improved scaling is due to better performance or not matters little. On the other hand, when looking to identify algorithms which are more successful - performance-wise - as the scale increases, interpretation of these plots will be more troublesome⁵. In the plots found in Appendix C Sections C.1.3 and C.1.4, we can observe only a few isolated instances where the median performance is found to be better than that achieved at the previous dimensionality scale, where afterwards the expected trend of decreasing performance recovers. Therefore for our purposes in this section, any perceived improvement to rank order as dimensionality

⁵ This situation was always very unlikely to transpire due to the effects of the curse of dimensionality. It is mentioned here only to draw attention to the possibility of misinterpreting this kind of rank visualisation

increases can be safely assumed to attributed only to the rank reduction of other algorithms (point (i) above).

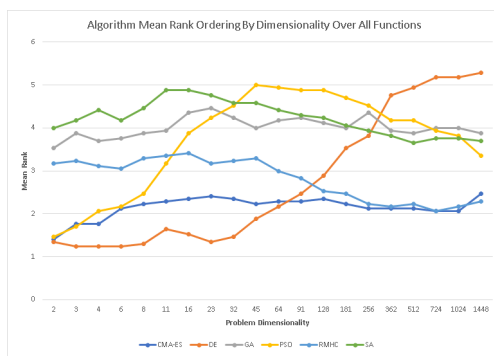
To reiterate the corresponding hypothesis, we expected that the ordering (and therefore the performance) between algorithms would change as problem dimensionality increases. In providing evidence supporting this hypothesis, we hoped to further find out: (i) whether there is a pattern, in terms of the presence of certain problem features, as to what kind of problems and/or at which point an algorithm performs well or ceases to perform well, (ii) the scale of performance degradation for each algorithm in terms of dimensionality increase, and (iii) the comparative scalability of the algorithms given (ii).

7.3.2.1 *Rank Ordering Over the Entire Benchmark Suite*

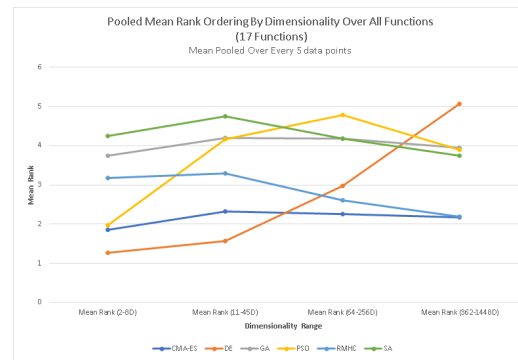
Before exploring the results in the context of problems possessing specific features, it is worth discussing the comparative results in terms of the entire benchmark suite i.e., the general performance of algorithms when all four of the investigated features are present. It is often the case, particularly for real-world problems, that little is known about the underlying structure of a particular problem or problem class; algorithm selection then must be performed from a more general 'higher level' context, such as in terms of a given application domain e.g., the space of symbolic regression problems, of a certain size. Algorithm selection at these higher granularities will almost certainly produce less meaningful recommendations compared to those made in the context of a more defined set of problems, however they can still usefully give some idea as to the likelihood that one algorithm is more suitable for use than another. For example in Fig. 7.8 and Table 7.2, the mean aggregate rank (between approximately 1.5 - 2.5) of CMA-ES over all scales provides some confidence - given no deeper investigation - that most of the comparative rankings are good (around ranks 1 or 2) for functions which can include: unimodal, multimodal, separable and non-separable features, and further that it remains robust at larger scales. On the other hand, at scales of above 128 dimensions for this set of problems, DE would likely be unsuitable. Thus, performing a general comparison such as this allows one to focus in on promising recommendations for further analysis if and when more is known about the structure of the specific problem being tackled, whilst simultaneously being able to discard algorithms which are more unlikely to be useful.

Table 7.2: Pooled Mean Rank Table Over All Functions

Alg./Dim. Range	Mean Rank (2-8D)	Mean Rank (11-45D)	Mean Rank (64-256D)	Mean Rank (362-1448D)
<i>CMA-ES</i>	1.86	2.33	2.26	2.16
<i>DE</i>	1.27	1.58	2.98	5.07
<i>GA</i>	3.75	4.20	4.18	3.94
<i>PSO</i>	1.98	4.16	4.79	3.89
<i>RMHC</i>	3.18	3.29	2.61	2.19
<i>SA</i>	4.25	4.74	4.19	3.74



Mean Aggregate Rankings



Pooled Mean Rankings

Figure 7.8: Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Benchmark Functions (where n = 5)

From the plots in Fig. 7.8, we can immediately make observations about the behaviour of both PSO and DE in comparison to the other algorithms as the scale of the problems in the benchmark suite increase. Firstly, DE is shown to be the best approach to use on average on our benchmark suite until the problem scale reaches 64 dimensions, after which it sharply decreases in performance in comparison to the other approaches, becoming the worst performing approach at a scale of 362 dimensions of which it does not recover. This observation is also mirrored clearly in the pooled mean plot albeit with less precision. Similarly to the observation made of PSO in section 7.2.1, here PSO performs relatively well at the lower dimensionalities of the benchmark suite becoming significantly worse at the low-mid to mid-high dimensionality ranges, to the point of ranking as the worst approach on average at 45 dimensions. Again, and where this observation is similar to that made in section 7.2.1, is that from 45 dimensions onwards, PSO becomes more efficient in comparison to the other

algorithms - settling at an average rank of approximately 3.5 at a problem scale of 1448 dimensions, beaten only by CMA-ES and RMHC. It is unclear if this behaviour is a direct effect of the main behavioural observation discussed in section 7.2.1, as the cyclic and sharp decreases in performance which would readily explain the similarity of behaviour here were related only to parameter configurations from a tuning strategy making use of 10,000 function evaluations, where the ranks discussed in this section are related to the performance of algorithms tuned with a budget of 50,000 evaluations. However, what can be observed from the 50,000 function evaluation plots for PSO in Appendix A is that for each of the benchmark functions where the cyclic behaviour can be observed for PSO tuned to 10,000 evaluations, a corresponding but often far less pronounced decrease - occurring at similar dimensionalities - can also be seen for PSO tuned at 50,000 evaluations which also recovers along with the 10,000 evaluation tuned PSO. Looking at the comparison plots in Appendix C, sections C.1.3 and C.1.4, we can see the effect of these small decreases in performance on rank ordering more clearly - particularly for the plots concerning: Bent Cigar, Chung-Reynolds, Rastrigin and Sum Squares. Here, where there is a small decrease in performance around the mid ranges for PSO, it is often enough to worsen the standing of the approach in comparison with the others - regularly falling into the worst performing algorithm at these points. As dimensionality increases, and since the drop off in performance often presents itself as being more gradual for PSO when compared to other approaches (i.e., is more robust to scale), PSO is eventually ranked better than those approaches which do not scale as well. The effect of these small decreases in performance on the ranking of PSO is also viewed clearly on the rank orderings on individual functions in Appendix C Section C.2.1.

From a more general standpoint, from the plots in Fig. 7.8 we can observe only slight changes in ordering amongst the remaining approaches - at least in the context of aggregated ranks. The ranks of SA and GA switch in favour of SA beyond 181 dimensions, this does not reoccur higher scales. The difference in average ranking between SA and GA, although larger at small to small-medium scales (2-64 dimensions) decreases at higher scales to that of being negligible. Additionally, the gradual improvement to the rank of both these approaches as dimensionality increases beyond 16 dimensions seems likely to be an effect originating mainly from the sharp decline of the ranks of both DE and PSO as discussed previously rather than an improvement in the performance of SA and GA at larger scales. Both RMHC and CMA-ES are seen to be the best algorithms to use on this benchmark suite on average, where the typical mean rank ordering is in the range of $\approx 1.6 - 3.5$. At larger scales, the mean ranks of both algorithms coalesce at around 2.1 - 2.5, with the only ordering change occurring

at 1448 dimensions in favour of RMHC. Otherwise, CMA-ES maintains its rank as the best algorithm to use on this suite overall.

7.3.2.2 Rank Ordering Over the Set of Uni-modal Versus Multi-modal Functions

Here and in the next section, the ranking of algorithms against sets of problems possessing complementary pairs of features will be discussed. Fig. 7.9 presents both the mean aggregate ranks and pooled mean ranks over several dimensionality ranges (as in the previous section) over all uni-modal problems in the benchmark suite. Fig. 7.10 similarly presents the rank orderings over all multi-modal problems. The pooled mean ranking tables for uni-modal and multi-modal problems are provided as Tables 7.3 and 7.4 respectively.

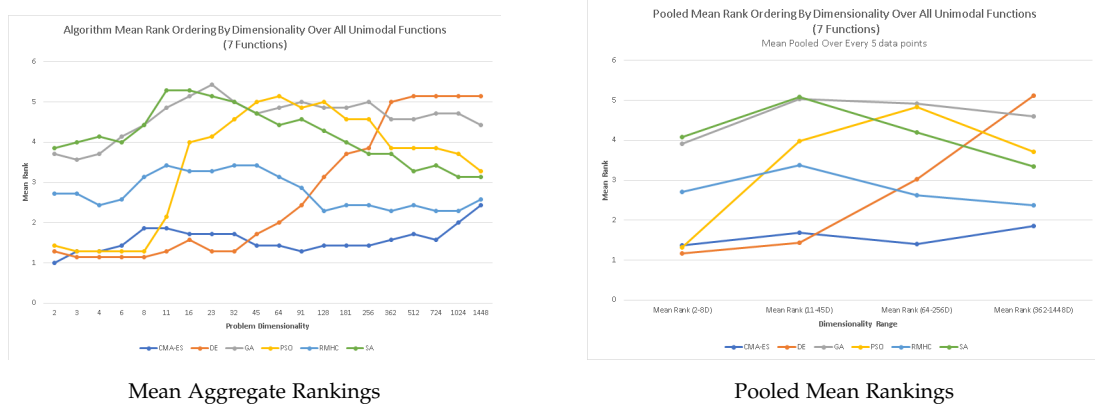


Figure 7.9: Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Uni-modal Functions (where n = 5)



Figure 7.10: Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Multi-modal Functions (where n = 5)

Table 7.3: Table of Pooled Mean Ranks Over All Uni-modal Functions

Alg./Dim. Range	Mean Rank (2-8D)	Mean Rank (11-45D)	Mean Rank (64-256D)	Mean Rank (362-1448D)
<i>CMA-ES</i>	1.37	1.69	1.40	1.86
<i>DE</i>	1.17	1.43	3.03	5.11
<i>GA</i>	3.91	5.03	4.91	4.60
<i>PSO</i>	1.31	3.97	4.83	3.71
<i>RMHC</i>	2.71	3.37	2.63	2.37
<i>SA</i>	4.09	5.09	4.20	3.34

Table 7.4: Table of Pooled Mean Ranks Over All Multi-modal Functions

Alg./Dim. Range	Mean Rank (2-8D)	Mean Rank (11-45D)	Mean Rank (64-256D)	Mean Rank (362-1448D)
<i>CMA-ES</i>	2.2	2.78	2.86	2.38
<i>DE</i>	1.34	1.68	2.94	5.04
<i>GA</i>	3.64	3.62	3.66	3.48
<i>PSO</i>	2.44	4.3	4.76	4.02
<i>RMHC</i>	3.5	3.24	2.6	2.06
<i>SA</i>	4.36	4.5	4.18	4.02

From this set of results, we again see evidence of the poor comparative scalability of DE and PSO on high dimensional problems previously observed when all benchmark functions were considered in the last section. Over both sets of problems, the mean ranking trends for DE appear very similar showing little discernible difference in rankings on average against each set. However, despite the scalability issues at higher dimensions, DE ranks best overall for both uni-modal and multi-modal problems at lower scales. For uni-modal problems, the comparative effectiveness of DE arises only up until a scale of 32 dimensions after which the rankings of the algorithm rapidly decline placing as the least effective algorithm at a scale of 362 dimensions onwards. DE however does fair better comparatively with others against the multi-modal set, maintaining the best average ranking until a scale of 91 dimensions where, as for its rankings against uni-modal problems, there is a sharp deterioration in rank, again becoming the worst performing algorithm at the 362 dimensional scale. As with DE, there appears to be little difference between the trends of PSO against each problem set, again suggesting that problems possessing either of the two features do not significantly affect the scalability of the approach. Additionally, there is little variation between the PSO ranks which have been mean aggregated across all benchmark functions (Fig. 7.8) and those aggregated by either of the problem features. For uni-modal problems however, PSO does maintain a good averaged rank compared to the other algorithms over a larger range of dimensions albeit negligible extending only so far as 8 dimensional problems before seeing a sharp decline in its rankings - again falling in line with those presented for all benchmark problems (Fig. 7.8) and multi-modal problems (Fig. 7.10). The rank ordering between DE and PSO switches in favour of PSO at a scale of 256 dimensions for uni-modal problems and 362 dimensions for multi-modal problems.

As with the rank results in the context of the entire benchmark suite, the visualisations similarities of mean rankings between distinct pairs of algorithms - where the ranks fall closely to one another but more distant from either algorithm comprising the other pair. Again, this pairing between algorithms with similar scalability relative to the other algorithms can also be observed in the results over the whole suite of problems (Fig. 7.8). However, in the context of the set of multi-modal problems, the distinct pairings become more evident as the scale of the problems increase as at low scales, RMHC and GA rank more similarly to each other than with the remaining algorithms.

Although changes to the rank ordering of the pairing comprised of GA and SA and the pairing of RMHC and CMA-ES are negligible for the set of uni-modal problems, as well as for GA and SA against the set of multi-modal functions, a change of rank order does occur

between CMA-ES and RMHC at a scale of 64 dimensions for the multi-modal set after which the mean rank of RMHC remains better than that of CMA-ES at the remaining scales.

7.3.2.3 Rank Ordering Over the Set of Separable Versus Non-separable Functions

In this section, the mean aggregate rankings and pooled mean rankings of the algorithms against the set of separable problems and the set of non-separable problems are presented (Figs. 7.11 and 7.12 respectively). Tables 7.5 and 7.6 additionally provide the pooled mean rankings for both sets of problems.

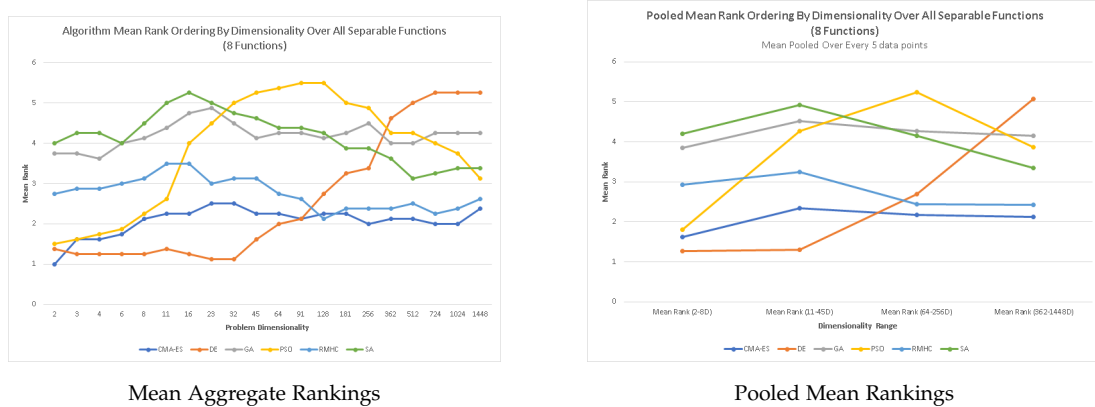


Figure 7.11: Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Separable Functions (where $n = 5$)

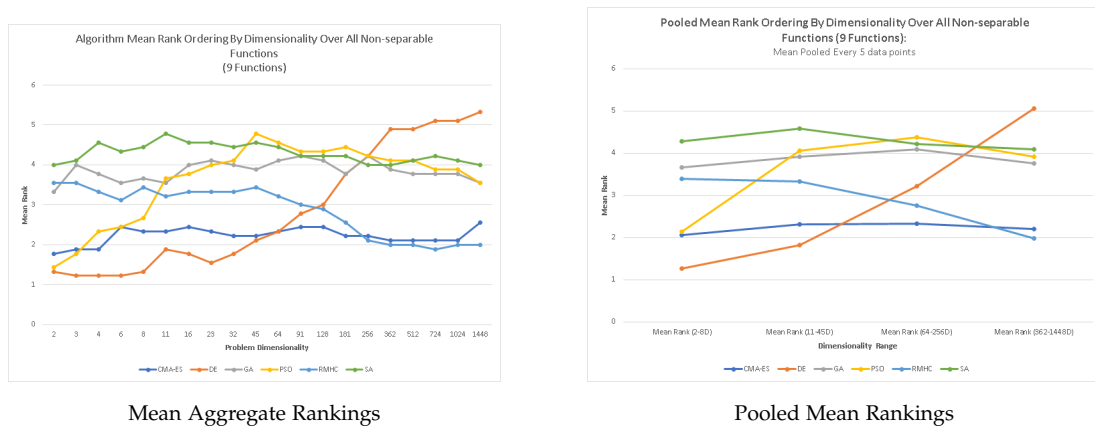


Figure 7.12: Mean Aggregate Ranking for Each Dimensionality and Pooled Mean Rank Over Every n Data Points Over All Non-separable Functions (where $n = 5$)

As with the ranking results over uni-modal and multi-modal problems and for those over the benchmark suite as a whole, a similar pattern of rapid decline in scalability relative to the other approaches again emerges for both DE and PSO. Where the ranking trend for DE

Table 7.5: Table of Pooled Mean Ranks Over All Separable Functions

Alg./Dim. Range	Mean Rank (2-8D)	Mean Rank (11-45D)	Mean Rank (64-256D)	Mean Rank (362-1448D)
<i>CMA-ES</i>	1.625	2.35	2.175	2.125
<i>DE</i>	1.275	1.3	2.7	5.075
<i>GA</i>	3.85	4.525	4.275	4.15
<i>PSO</i>	1.8	4.275	5.25	3.875
<i>RMHC</i>	2.925	3.25	2.45	2.425
<i>SA</i>	4.2	4.925	4.15	3.35

Table 7.6: Table of Pooled Mean Ranks Over All Non-separable Functions

Alg./Dim. Range	Mean Rank (2-8D)	Mean Rank (11-45D)	Mean Rank (64-256D)	Mean Rank (362-1448D)
<i>CMA-ES</i>	2.07	2.31	2.33	2.20
<i>DE</i>	1.27	1.82	3.22	5.07
<i>GA</i>	3.67	3.91	4.09	3.76
<i>PSO</i>	2.13	4.07	4.38	3.91
<i>RMHC</i>	3.40	3.33	2.76	1.98
<i>SA</i>	4.29	4.58	4.22	4.09

over both separable and non-separable problems largely matches those for sets of problems discussed previously; PSO in the context the set of non-separable problems, although showing the same pattern of rapid worsening of rank before improving at larger scales, does not show the same level of decline in ranking as we have seen over other problem subsets. Where previously, in the worst case, PSO achieves a mean aggregate ranking of between 5 - 5.5, this time around the medium scale problems PSO achieves a worst case mean ranking of approximately 4.7. Although this still places PSO as the worst performing at this scale, it is still far closer to the rankings of SA and GA which we have not previously observed over other sets of problems.

The close pairings of the rankings of the remaining algorithms previously discussed in relation to other problem sets remains the same for both separable and non separable problems. In terms of the rank order of the members of these pairings, a change occurs at a problem scale of 128 dimensions in favour of SA when considering only separable problems with no such change in ordering between CMA-ES and RMHC. On the non-separable set however, there is no change to rank order between SA and GA, despite achieving closely related mean rankings, but a change occurs for RMHC and CMA-ES at a scale of 256 dimensions.

7.4 CONCLUSIONS AND DISCUSSION

7.4.1 *Tuning Evaluation Budget Comparisons*

Here we firstly presented a finding that suggests it is generally better, given no other indicators, to use an evaluation budget at least equal to the target evaluations when solving large-scale optimisation problems. This conclusion runs contrary to a current SMAC recommendation that one should use a budget of $10 \times D$ evaluations, where D is the problem dimensionality [62]. For the largest number of dimensions we experimented with (1448D) this equates to a budget recommendation of around 14,480 function evaluations. Further in this study, a suggestion that this should be further extended to $100 \times D$ evaluations [62]. The latter suggestion we feel is 'overkill' and only aids in worsening the situation in relation to the runtime of tuning processes - where common sense would dictate that actually, since an algorithm tuned for this number of evaluations targeting a lower budget would expect more time to explore the available space, that algorithm performance would worsen. We have observed a similar effect from the plots in Appendix A where an algorithm is tuned at a 50,000 evaluation budget, but when the algorithm reaches only 10,000 evaluations during experiments the performance

is shown to be significantly worse (compared to results from the 10,000 evaluation budget) over the various dimension thresholds of all our functions, with the exception of: Deflected Corrugated Spring, Inverted-Cosine-Wave and Rastrigin which showed little or no difference between the two budgets at 10,000 evaluations.

We also found no appreciable difference between tuning at either tuning evaluation budget for our implementation of GA (an SSGA). We conclude from our findings here that when using an SSGA making use of the same operators and solution representation for continuous function optimisation as we have here, the approach can be tuned at a significantly lower number of evaluations than the required target evaluation budget. Automatic parameter tuning can often be computationally expensive so this finding potentially eases this restriction that may prohibit APT for the other approaches we tested, allowing GA to be tuned faster.

Next we presented our findings relating to the initial observation that our implementation of PSO behaves differently from the others when tuned for our two evaluation budgets (10,000 and 50,000). In general, what we observed was that for all but $\approx \frac{1}{4}$ of our function suite, no deviation occurs between budgets at lower dimensions and after a significant divergence of performance in mid ranged dimensions when using different budgets a re-convergence occurs in higher dimensions. We initially interpreted these observations as there being no difference between using lower or higher tuning budgets. After some investigation we concluded further that the effect could be attributed to there being a threshold number of dimensions above which a larger population causes our implementation some difficulty. Similar past results in the field of Genetic Programming in [3], suggested that small population sizes were better at solving some problems. Further, evidence also shows that small population sizes can sometimes lead the search to the optimum solutions in fewer objective function evaluations - for certain problems - when compared to the use of larger sizes [147]. This may explain why a dramatic decrease in population size for the algorithm using the 10,000 evaluation parameter set produces comparable results to the alternate scheme (50,000 evaluation budget); SMAC could have determined that a smaller population size (increasing exploitation effects around any decent local minimum) would be more successful where the scale of the problem instances reached a point where a more equal balance between exploration and exploitation could not do better than a highly exploitative search. Indeed, in the large scale optimisation literature, small populations are used with LSOPs to 'brute-force' promising search directions [18], since full coverage of a very large search space is intractable using reasonable population sizes and/or function evaluation budgets. However, we currently have no explanation or

reason as to why SMAC found a high population size - which invariably produced worse solutions - to be the best configuration.

7.4.2 *Performance Comparison Between Algorithm Implementations*

Here we made the moderately surprising observation that our Random Mutation Hill Climbing (RMHC) algorithm regularly performs competitively with far more complex approaches. We conclude through further analysis of the available literature that results such as those discussed here oppose widely held opinions that simple algorithm like Stochastic Hill Climbing - of which RMHC is a variation - are not very useful when compared to more complex evolutionary algorithms. A contribution here is that although it is rare to find comparisons between more powerful EAs and simple hill climbing algorithms, we have now produced data for this kind of comparison, for large-scale optimisation, which can be made available to the wider research community through the University's Digital Repository.

7.4.3 *Rank Ordering Changes at Larger Problem Scales*

In general we can conclude that when viewed as the aggregate rank over several functions from the benchmark set, much of the noise from small intermittent re-orderings can be removed leaving only the most meaningful changes to rank order thus reflecting a clearer picture of comparative performance between the algorithms covered. Such meaningful transitions between the algorithms are therefore few, but occur around the medium to large problem scales as hypothesised and persist over the remaining (higher) scales. In terms of ranked performance over subsets of the benchmark suite, i.e., subsets delineated by the problem features: modality (uni-modal vs. multi-modal) and separability (separable vs. non-separable), there is not sufficient reason to conclude that any of the problem features considered pose any kind of differing levels of difficulty for any particular one of the metaheuristic implementations compared. Despite generally localised differences to the mean aggregate rankings for different problem subsets, these when compared by problem scale do not show much variation.

7.5 CONTRIBUTIONS

- We have produced a concise review of the subject area for future consumption by other researchers looking to make progress in this area

- We have generated a vast amount of data relating to the large-scale optimisation of several standard continuous function benchmarks found in the literature that can be made available to the wider community through the university's digital repository in due course.
- We have identified several behaviours of common metaheuristics as applied to optimisation problems with many hundreds of dimensions. Filling some of the gaps in our understanding of how the optimisation approaches, we may assume behave in certain ways, actually do behave when presented with much harder problems than have they have typically been investigated against.
- We have identified an effect in PSO in which for the majority of our benchmark set behaves unusually in the presence of large population sizes - an effect we felt to be counter-intuitive. Through further investigation we hope to further explain the causes of this effect and be better able to predict the type of functions where it is likely we will observe this effect occurring, and further, research into the possible benefits or disadvantages such a situation may present
- Results show that the SSGA used in this study can be tuned at a significantly lower evaluation budget than that of the intended target budget without affecting performance (in terms of solution quality). This finding carries the implication that the observed requirement of tuning at the target evaluation budget in other algorithms can be largely alleviated for SSGA - allowing the algorithm to be tuned faster.

7.6 FUTURE WORK

Based on some of the behaviours we are beginning to observe in even the most common of metaheuristic implementations, there is of course scope to continue this line of research - not only to develop more evidence to support our findings here but also to investigate further into the possible reasons - an overlap into the realm of 'Explanatory Metaheuristics' - for behavioural changes observed at the kind of scale we have been dealing with.

In particular, the discovery of a threshold above which large population sizes in PSO are unhelpful warrants much further investigation; as an understanding of how and when this effect will present itself may inspire improvements to the current state of affairs in regards to particle swarm optimisation research and state-of-the-art PSO algorithm performance as well as aiming to generate insight into how parameters interact within metaheuristic algorithms in general - a very interesting prospect quickly gaining momentum in the research community.

BIBLIOGRAPHY

- [1] Emile Aarts, Jan Korst, and Wil Michiels. Simulated annealing. *Search Methodologies: Introductory Tutorials in Optimization and Decision Support Techniques*, page 187, 2006.
- [2] W. Al-Hassan, M. B. Fayek, and S. I. Shaheen. Psosa: An optimized particle swarm technique for solving the urban planning problem. In *2006 International Conference on Computer Engineering and Systems*, pages 401–405, Nov 2006. doi: 10.1109/ICCES.2006.320481.
- [3] D. A. Ashlock, K. M. Bryden, and S. Corns. Small population effects and hybridization. In *2008 IEEE Congress on Evolutionary Computation (IEEE World Congress on Computational Intelligence)*, pages 2637–2643, June 2008. doi: 10.1109/CEC.2008.4631152.
- [4] Anne Auger and Olivier Teytaud. Continuous lunches are free! In *Proceedings of the 9th annual conference on Genetic and evolutionary computation*, pages 916–922. ACM, 2007.
- [5] Thomas Bäck, David B Fogel, and Zbigniew Michalewicz. *Evolutionary Computation 1: Basic algorithms and Operators*, volume 1. CRC press, 2000.
- [6] Sunith Bandaru and Kalyanmoy Deb. Metaheuristic techniques. *Decision Sciences: Theory and Practice*, pages 693–750, 2016.
- [7] Jagdish Chand Bansal, PK Singh, Mukesh Saraswat, Abhishek Verma, Shimpi Singh Jadon, and Ajith Abraham. Inertia weight strategies in particle swarm optimization. In *Nature and Biologically Inspired Computing (NaBIC), 2011 Third World Congress on*, pages 633–640. IEEE, 2011.
- [8] Richard E Bellman. *Adaptive Control Processes: a Guided Tour*, volume 2045. Princeton University Press, 1961.
- [9] John B Biggs. The role of metalearning in study processes. *British journal of educational psychology*, 55(3):185–212, 1985.
- [10] James Blondin. Particle swarm optimization: A tutorial. Available from: http://cs.armstrong.edu/saad/csci8100/psa_tutorial.pdf, 2009.
- [11] Christian Blum and Andrea Roli. Metaheuristics in combinatorial optimization: Overview and conceptual comparison. *ACM computing surveys (CSUR)*, 35(3):268–308, 2003.

- [12] Yossi Borenstein and Alberto Moraglio. *Theory and Principled Methods for the Design of Metaheuristics*. Springer, 2014.
- [13] Daniel Bratton and James Kennedy. Defining a standard for particle swarm optimization. In *Swarm Intelligence Symposium, 2007. SIS 2007. IEEE*, pages 120–127. IEEE, 2007.
- [14] Pavel Brazdil, Christophe Giraud Carrier, Carlos Soares, and Ricardo Vilalta. *Metalearning: Applications to Data Mining*. Springer Science & Business Media, 2008.
- [15] Leo Breiman and Adele Cutler. Random forests. URL https://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm#remarks.
- [16] Jason Brownlee. *Clever Algorithms: Nature-inspired Programming Recipes*. Jason Brownlee, 2011.
- [17] Edmund K Burke, Graham Kendall, et al. *Search methodologies*. Springer, 2005.
- [18] Fabio Caraffini, Ferrante Neri, and Giovanni Iacca. Large scale problems in practice: The effect of dimensionality on the interaction among variables. In Giovanni Squillero and Kevin Sim, editors, *Applications of Evolutionary Computation*, pages 636–652. Springer International Publishing, 2017. ISBN 978-3-319-55849-3.
- [19] Wei Chu, Xiaogang Gao, and Soroosh Sorooshian. Handling boundary constraints for particle swarm optimization in high-dimensional search space. *Information Sciences*, 181(20):4569–4581, October 2011. ISSN 0020-0255. doi: 10.1016/j.ins.2010.11.030. URL <http://dx.doi.org/10.1016/j.ins.2010.11.030>.
- [20] Maurice Clerc. The swarm and the queen: Towards a deterministic and adaptive particle swarm optimization. In *Evolutionary Computation, 1999. CEC 99. Proceedings of the 1999 Congress on*, volume 3, pages 1951–1957. IEEE, 1999.
- [21] Maurice Clerc. *Guided randomness in optimization*, volume 1. John Wiley & Sons, 2015.
- [22] David A Coley. *An introduction to genetic algorithms for scientists and engineers*. World Scientific Publishing Company, 1999.
- [23] Swagatam Das and Ponnuthurai Nagaratnam Suganthan. Differential evolution: A survey of the state-of-the-art. *IEEE transactions on evolutionary computation*, 15(1):4–31, 2011.
- [24] Kenneth A De Jong and Jayshree Sarma. Generation gaps revisited. In *Foundations of Genetic Algorithms*, volume 2, pages 19–28. Elsevier, 1993.

- [25] Kenneth Alan De Jong. *An Analysis of the Behavior of a Class of Genetic Adaptive Systems*. PhD thesis, MI, USA, 1975. AAI7609381.
- [26] Kalyanmoy Deb and Debayan Deb. Analysing mutation schemes for real-parameter genetic algorithms. *International Journal of Artificial Intelligence and Soft Computing*, 4(1): 1–28, 2014.
- [27] Ke-Lin Du and MNS Swamy. Simulated annealing. In *Search and Optimization by Metaheuristics*, pages 29–36. Springer, 2016.
- [28] Russell C Eberhart and Yuhui Shi. Tracking and optimizing dynamic systems with particle swarms. In *Evolutionary Computation, 2001. Proceedings of the 2001 Congress on*, volume 1, pages 94–100. IEEE, 2001.
- [29] Agoston E Eiben, James E Smith, et al. *Introduction to Evolutionary Computing*, volume 53. Springer, 2003.
- [30] Andries P. Engelbrecht. *Computational Intelligence: An Introduction*. Wiley Publishing, 2nd edition, 2007. ISBN 0470035617.
- [31] Yong Feng, Gui-Fa Teng, Ai-Xin Wang, and Yong-Mei Yao. Chaotic inertia weight in particle swarm optimization. In *Innovative Computing, Information and Control, 2007. ICICIC'07. Second International Conference on*, pages 475–475. IEEE, 2007.
- [32] Steffen Finck, Nikolaus Hansen, Raymond Ros, and Anne Auger. Real-parameter black-box optimization benchmarking 2009: Presentation of the noiseless functions. Technical report, Citeseer, 2010.
- [33] Len Fisher. *The Perfect Swarm: The Science of Complexity in Everyday Life*. Basic Books, 2009.
- [34] Ronald Aylmer Fisher. *The genetical theory of natural selection: a complete variorum edition*. Oxford University Press, 1999.
- [35] Stephanie Forrest and Melanie Mitchell. Relative building-block fitness and the building-block hypothesis. In *Foundations of genetic algorithms*, volume 2, pages 109–126. Elsevier, 1993.
- [36] Andrea Gavana. Test functions index - ampgo 0.1.0 documentation. URL http://infinity77.net/global_optimization/test_functions.html#test-functions-index.

- [37] Stuart Geman and Donald Geman. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. In *Readings in Computer Vision*, pages 564–584. Elsevier, 1987.
- [38] Robert Ghanea-Hercock. *Applied Evolutionary Algorithms in Java*. Springer Science & Business Media, 2013.
- [39] F. Glover and K. Sörensen. Metaheuristics. *Scholarpedia*, 10(4):6532, 2015. doi: 10.4249/scholarpedia.6532. revision #149834.
- [40] Jochen Gortler, Rebecca Kehlbeck, and Oliver Deussen. Visual exploration of gaussian processes, 2019. URL <https://distill.pub/2019/visual-exploration-gaussian-processes/>.
- [41] Kevin Graham and Leslie Smith. Comparing hyper-heuristics with blackboard systems. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, pages 1141–1145, 2017.
- [42] Kevin Graham, Jerry Swan, and Simon Martin. The ‘blackboard pattern’ for metaheuristics. In *Proceedings of the Companion Publication of the 2015 Annual Conference on Genetic and Evolutionary Computation*, pages 1265–1267, 2015.
- [43] Bruce Hajek. Cooling schedules for optimal annealing. *Mathematics of operations research*, 13(2):311–329, 1988.
- [44] Nikolaus Hansen. Cma evolution strategy source code, 2011. URL http://cma.gforge.inria.fr/cmaes_sourcecode_page.html.
- [45] Nikolaus Hansen. The cma evolution strategy: A tutorial. 2016.
- [46] Nikolaus Hansen and Stefan Kern. Evaluating the cma evolution strategy on multimodal test functions. In *International Conference on Parallel Problem Solving from Nature*, pages 282–291. Springer, 2004.
- [47] Nikolaus Hansen, Steffen Finck, Raymond Ros, and Anne Auger. Real-parameter black-box optimization benchmarking 2009: Presentation of the noiseless functions. Technical report, 2010.
- [48] Nikolaus Hansen, Raymond Ros, Nikolas Mauny, Marc Schoenauer, and Anne Auger. Impacts of invariance in search: When cma-es and pso face ill-conditioned and non-separable problems. *Applied Soft Computing*, 11(8):5755–5769, 2011.

- [49] Nikolaus Hansen, Anne Auger, Olaf Mersmann, Tea Tusar, and Dimo Brockhoff. Coco: A platform for comparing continuous optimizers in a black-box setting. *arXiv preprint arXiv:1603.08785*, 2016.
- [50] Randy L Haupt and Sue Ellen Haupt. *Practical genetic algorithms*, volume 2. Wiley New York, 2004.
- [51] Abdel-Rahman Hedar. Global optimisation test problems. URL http://www-optima.amp.i.kyoto-u.ac.jp/member/student/hedar/Hedar_files/TestG0.htm.
- [52] Sabine Helwig and Rolf Wanka. Theoretical analysis of initial particle swarm behavior. In *International conference on parallel problem solving from nature*, pages 889–898. Springer, 2008.
- [53] Darrall Henderson, Sheldon H Jacobson, and Alan W Johnson. The theory and practice of simulated annealing. In *Handbook of metaheuristics*, pages 287–319. Springer, 2003.
- [54] J. H. Holland. Genetic algorithms. *Scholarpedia*, 7(12):1482, 2012. doi: 10.4249/scholarpedia.1482. revision #128222.
- [55] John H. Holland. *Hidden Order: How Adaptation Builds Complexity*. Addison Wesley Longman Publishing Co., Inc., Redwood City, CA, USA, 1995. ISBN 0-201-40793-0.
- [56] S. Hollasch. *Four-space Visualization of 4D Objects*. PhD thesis, 1991.
- [57] Qiang Huang, Thomas White, Guanbo Jia, Mirco Musolesi, Nil Turan, Ke Tang, Shan He, John K Heath, and Xin Yao. Community detection using cooperative co-evolutionary differential evolution. In *International Conference on Parallel Problem Solving from Nature*, pages 235–244, 2012.
- [58] F. Hutter, H. H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algorithm configuration (extended version). Technical Report TR-2010-10, University of British Columbia, Department of Computer Science, 2010. Available online: <http://www.cs.ubc.ca/~hutter/papers/10-TR-SMAC.pdf>.
- [59] F. Hutter, H.H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *Proc. of LION-5*, pages 507–523, 2011.
- [60] F. Hutter, H.H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algorithm configuration (slide presentation), 2011.
- [61] Frank Hutter, Holger H Hoos, Kevin Leyton-Brown, and Kevin P Murphy. An experimental investigation of model-based parameter optimisation: Spo and beyond. In

- Proceedings of the 11th Annual conference on Genetic and evolutionary computation*, pages 271–278. ACM, 2009.
- [62] Frank Hutter, Holger H. Hoos, and Kevin Leyton-Brown. An evaluation of sequential model-based optimization for expensive blackbox functions. In *GECCO*, 2013.
- [63] Momin Jamil and Xin-She Yang. A literature survey of benchmark functions for global optimisation problems. *International Journal of Mathematical Modelling and Numerical Optimisation*, 4(2):150–194, 2013.
- [64] Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 13(4):455–492, 1998.
- [65] Yilmaz Kaya, Murat Uyar, et al. A novel crossover operator for genetic algorithms: Ring crossover. *arXiv preprint arXiv:1105.0355*, 2011.
- [66] J. Kennedy. Small worlds and mega-minds: Effects of neighborhood topology on particle swarm performance. In *Proceedings of the 1999 Congress on Evolutionary Computation-CEC99 (Cat. No. 99TH8406)*, volume 3, pages 1931–1938 Vol. 3, July 1999. doi: 10.1109/CEC.1999.785509.
- [67] J. Kennedy and R. Eberhart. Particle swarm optimization. In *Neural Networks, 1995. Proceedings., IEEE International Conference on*, volume 4, pages 1942–1948 vol.4, Nov 1995. doi: 10.1109/ICNN.1995.488968.
- [68] James Kennedy and Rui Mendes. Population structure and particle swarm performance. In *Proceedings of the 2002 Congress on Evolutionary Computation. CEC’02 (Cat. No. 02TH8600)*, volume 2, pages 1671–1676. IEEE, 2002.
- [69] Scott Kirkpatrick, C Daniel Gelatt, and Mario P Vecchi. Optimization by simulated annealing. *science*, 220(4598):671–680, 1983.
- [70] Jack PC Kleijnen. Response surface methodology. In *Handbook of simulation optimization*, pages 81–104. Springer, 2015.
- [71] Oscar Knagg. An intuitive guide to gaussian processes, 2015. URL <https://towardsdatascience.com/an-intuitive-guide-to-gaussian-processes-ec2f0b45c71d>.
- [72] Johannes W Kruisselbrink, Rui Li, Edgar Reehuis, Jeroen Eggermont, and Thomas Bäck. On the log-normal self-adaptation of the mutation rate in binary search spaces. In *Proceedings of the 13th annual conference on Genetic and evolutionary computation*, pages 893–900. ACM, 2011.

- [73] Manuel Laguna and Rafael Martí. Experimental testing of advanced scatter search designs for global optimization of multimodal functions. *Journal of Global Optimization*, 33(2):235–255, 2005.
- [74] Christiane Lemke, Marcin Budka, and Bogdan Gabrys. Metalearning: a survey of trends and technologies. *Artificial intelligence review*, 44(1):117–130, 2015.
- [75] Xiaodong Li and Xin Yao. Tackling high dimensional nonseparable optimisation problems by cooperatively coevolving particle swarms. In *2009 IEEE Congress on Evolutionary Computation*, pages 1546–1553, 2009.
- [76] Xiaodong Li and Xin Yao. Cooperatively coevolving particle swarms for large scale optimization. *IEEE Transactions on Evolutionary Computation*, 16(2):210–224, 2011.
- [77] Adam Lipowski and Dorota Lipowska. Roulette-wheel selection via stochastic acceptance. *Physica A: Statistical Mechanics and its Applications*, 391(6):2193–2196, 2012.
- [78] Manuel López-Ibáñez, Jérémie Dubois-Lacoste, Leslie Pérez Cáceres, Thomas Stützle, and Mauro Birattari. Iterated racing for automatic algorithm configuration [r package irace version 3.3]. URL <https://cran.r-project.org/package=irace>.
- [79] Manuel López-Ibáñez, Jérémie Dubois-Lacoste, Leslie Pérez Cáceres, Mauro Birattari, and Thomas Stützle. The irace package: Iterated racing for automatic algorithm configuration. *Operations Research Perspectives*, 3:43–58, 2016.
- [80] Helena R Lourenço, Olivier C Martin, and Thomas Stützle. Iterated local search. In *Handbook of Metaheuristics*, pages 320–353. Springer, 2003.
- [81] Helena R Lourenço, Olivier C Martin, and Thomas Stützle. Iterated local search: Framework and applications. In *Handbook of metaheuristics*, pages 363–397. Springer, 2010.
- [82] Sean Luke. *Essentials of metaheuristics*, volume 113. Lulu Raleigh, 2009.
- [83] X. Ma, X. Li, Q. Zhang, K. Tang, Z. Liang, W. Xie, and Z. Zhu. A survey on cooperative co-evolutionary algorithms. *IEEE Transactions on Evolutionary Computation*, 23(3):421–441, June 2019. ISSN 1089-778X. doi: 10.1109/TEVC.2018.2868770.
- [84] Sarim Mahmood, Scott Rosenquist, and Bryce Benn. skarjoko/differential-evolution, 2013. URL <https://github.com/skarjoko/differential-evolution/blob/master/DifferentialEvolution.java>.

- [85] Gandhi Manalu. Particle swarm optimization: Sample code using java, 2010. URL <https://gandhim.wordpress.com/2010/04/04/particle-swarm-optimization-pso-sample-code-using-java/>.
- [86] Gandhi Manalu. Particle swarm optimization: Sample code using java, 2014. URL <https://github.com/therealmanalu/pso-example-java>.
- [87] Angelina Jane Reyes Medina, Gregorio Toscano Pulido, and José Gabriel Ramírez-Torres. A comparative study of neighborhood topologies for particle swarm optimizers. In *IJCCI*, pages 152–159, 2009.
- [88] Olaf Mersmann, Mike Preuss, Heike Trautmann, Bernd Bischl, and Claus Weihs. Analyzing the bbob results by means of benchmarking concepts. *Evolutionary computation*, 23(1):161–185, 2015.
- [89] Zbigniew Michalewicz. *Genetic Algorithms + Data Structures = Evolution Programs (3rd Ed.)*. Springer-Verlag, London, UK, UK, 1996. ISBN 3-540-60676-9.
- [90] Zbigniew Michalewicz and David B Fogel. *How to solve it: modern heuristics*. Springer Science & Business Media, 2013.
- [91] Zbigniew Michalewicz, Thomas Logan, and Swarnalatha Swaminathan. Evolutionary operators for continuous convex parameter spaces. In *Proceedings of the 3rd Annual conference on Evolutionary Programming*, pages 84–97. World Scientific, 1994.
- [92] Melanie Mitchell. *An introduction to genetic algorithms*. MIT press, 1998.
- [93] Marcin Molga and Czesław Smutnicki. Test functions for optimization needs. *Test functions for optimization needs*, 101, 2005.
- [94] John A Nelder and Roger Mead. A simplex method for function minimization. *The Computer Journal*, 7(4):308–313, 1965.
- [95] Ferrante Neri and Ville Tirronen. Recent advances in differential evolution: a survey and experimental analysis. *Artificial Intelligence Review*, 33(1-2):61–106, 2010.
- [96] Alexander G Nikolaev and Sheldon H Jacobson. Simulated annealing. In *Handbook of Metaheuristics*, pages 1–39. Springer, 2010.
- [97] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, New York, NY, USA, 2nd edition, 2006.
- [98] Yaghout Nourani and Bjarne Andresen. A comparison of simulated annealing cooling strategies. *Journal of Physics A: Mathematical and General*, 31(41):8373, 1998.

- [99] Mohammad Nabi Omidvar, Xiaodong Li, Zhenyu Yang, and Xin Yao. Cooperative co-evolution for large scale optimization through more frequent random grouping. In *IEEE Congress on Evolutionary Computation*, pages 1–8, 2010.
- [100] Mohammad Nabi Omidvar, Xiaodong Li, and Xin Yao. Cooperative co-evolution with delta grouping for large scale non-separable function optimisation. In *IEEE Congress on Evolutionary Computation*, pages 1–8, 2010.
- [101] Nikhil Padhye, Kalyanmoy Deb, and Pulkrit Mittal. Boundary handling approaches in particle swarm optimization. In *Proceedings of Seventh International Conference on Bio-Inspired Computing: Theories and Applications (BIC-TA 2012)*, pages 287–298. Springer, 2013.
- [102] Lucas Pavelski, Myriam Delgado, and Marie-Eleonore Kessaci. Meta-learning for optimization: A case study on the flowshop problem using decision trees. In *2018 IEEE Congress on Evolutionary Computation (CEC)*, pages 1–8. IEEE, 2018.
- [103] Peter Frank Perroni, Daniel Weingaertner, and Myriam Regattieri Delgado. Automated iterative partitioning for cooperatively coevolving particle swarms in large scale optimization. In *2015 Brazilian Conference on Intelligent Systems (BRACIS)*, pages 19–24, 2015.
- [104] Stjepan Picek, Domagoj Jakobovic, and Marin Golub. On the recombination operator in the real-coded genetic algorithms. In *Evolutionary Computation (CEC), 2013 IEEE Congress on*, pages 3103–3110. IEEE, 2013.
- [105] Riccardo Poli, James Kennedy, and Tim Blackwell. Particle swarm optimization. *Swarm intelligence*, 1(1):33–57, 2007.
- [106] Elena Popovici, Anthony Bucci, R. Paul Wiegand, and Edwin D. De Jong. *Coevolutionary Principles*, pages 987–1033. Springer, Berlin, Heidelberg, 2012. ISBN 978-3-540-92910-9. doi: 10.1007/978-3-540-92910-9_31. URL https://doi.org/10.1007/978-3-540-92910-9_31.
- [107] Mitchell A Potter and Kenneth A De Jong. A cooperative coevolutionary approach to function optimization. In *International Conference on Parallel Problem Solving from Nature*, pages 249–257, 1994.
- [108] Kenneth Price, Rainer M Storn, and Jouni A Lampinen. *Differential Evolution: a Practical Approach to Global Optimization*. Springer Science & Business Media, 2006.

- [109] Carl Edward Rasmussen. Gaussian processes in machine learning. In *Summer School on Machine Learning*, pages 63–71. Springer, 2003.
- [110] Asanga Ratnaweera, Saman K Halgamuge, and Harry C Watson. Self-organizing hierarchical particle swarm optimizer with time-varying acceleration coefficients. *IEEE Transactions on evolutionary computation*, 8(3):240–255, 2004.
- [111] Colin Reeves. Genetic algorithms. In *Handbook of metaheuristics*, pages 55–82. Springer, 2003.
- [112] John R Rice. The algorithm selection problem. In *Advances in computers*, volume 15, pages 65–118. Elsevier, 1976.
- [113] Alex Rogers and Adam Prügel-Bennett. Modelling the dynamics of a steady-state genetic algorithm. *Foundations of genetic algorithms*, 5:57–68, 1999.
- [114] Raymond Ros and Nikolaus Hansen. A simple modification in cma-es achieving linear time and space complexity. In *International Conference on Parallel Problem Solving from Nature*, pages 296–305. Springer, 2008.
- [115] Alejandro Rosete-Suárez, Alberto Ochoa-Rodríguez, and Michele Sebag. Automatic graph drawing and stochastic hill climbing. In *Proceedings of the 1st Annual Conference on Genetic and Evolutionary Computation - Volume 2, GECCO'99*, pages 1699–1706, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc. ISBN 1-55860-611-4. URL <http://dl.acm.org/citation.cfm?id=2934046.2934177>.
- [116] Franz Rothlauf. *Design of Modern Heuristics: Principles and Application*. Springer Science & Business Media, 2011.
- [117] Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Malaysia; Pearson Education Limited,, 2016.
- [118] Hans-Paul Schwefel. *Numerical Optimization of Computer Models*. John Wiley & Sons, Inc., New York, NY, USA, 1981. ISBN 0471099880.
- [119] Yuhui Shi and Russell Eberhart. A modified particle swarm optimizer. In *Evolutionary Computation Proceedings, 1998. IEEE World Congress on Computational Intelligence., The 1998 IEEE International Conference on*, pages 69–73. IEEE, 1998.
- [120] Yuhui Shi and Russell C Eberhart. Empirical study of particle swarm optimization. In *Evolutionary computation, 1999. CEC 99. Proceedings of the 1999 congress on*, volume 3, pages 1945–1950. IEEE, 1999.

- [121] John Silberholz and Bruce Golden. *Comparison of Metaheuristics*, pages 625–640. Springer US, Boston, MA, 2010. ISBN 978-1-4419-1665-5. doi: 10.1007/978-1-4419-1665-5_21. URL https://doi.org/10.1007/978-1-4419-1665-5_21.
- [122] S. Singer and J. Nelder. Nelder-Mead algorithm. *Scholarpedia*, 4(7):2928, 2009. doi: 10.4249/scholarpedia.2928. revision #91557.
- [123] Sanja Singer and Saša Singer. Complexity analysis of nelder-mead search iterations. In *Proceedings of the 1. Conference on Applied Mathematics and Computation, Dubrovnik, Croatia*, pages 185–196. PMF–Matematički odjel, Zagreb, 1999.
- [124] Jim Smith and Terence C. Fogarty. An adaptive poly-parental recombination strategy. In *AISB Workshop on Evolutionary Computing*, pages 48–61, 1995.
- [125] Jim Smith and Frank Vavak. Replacement strategies in steady state genetic algorithms: Static environments. *Foundations of genetic algorithms*, 5:219–233, 1999.
- [126] Kate A Smith-Miles. Cross-disciplinary perspectives on meta-learning for algorithm selection. *ACM Computing Surveys (CSUR)*, 41(1):6, 2009.
- [127] Padhraic Smyth. Local search and optimization, 2007.
- [128] Nitasha Soni and Tapas Kumar. Study of various mutation operators in genetic algorithms. *International Journal of Computer Science and Information Technologies (IJCSIT)*, 3, 2014. ISSN 0975-9646.
- [129] Rainer Storn and Kenneth Price. Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces. *Journal of Global Optimization*, 11(4): 341–359, 1997.
- [130] Philip N Strenske and Scott Kirkpatrick. Analysis of finite length annealing schedules. *Algorithmica*, 6(1-6):346–366, 1991.
- [131] T Stützle. *Local Search Algorithms for Combinatorial Problems*. PhD thesis, 1998.
- [132] Ponnuthurai N Suganthan. Particle swarm optimiser with neighbourhood operator. In *Evolutionary Computation, 1999. CEC 99. Proceedings of the 1999 Congress on*, volume 3, pages 1958–1962. IEEE, 1999.
- [133] Ponnuthurai N Suganthan, Nikolaus Hansen, Jing J Liang, Kalyanmoy Deb, Ying-Ping Chen, Anne Auger, and Santosh Tiwari. Problem definitions and evaluation criteria for the cec 2005 special session on real-parameter optimization. *KanGAL report*, page 2005, 2005.

- [134] Sonja Surjanovic and Derek Bingham. Virtual library of simulation experiments: Test functions and datasets, 2013. URL <https://www.sfu.ca/~ssurjano/>.
- [135] El-Ghazali Talbi. *Metaheuristics: from design to implementation*, volume 74. John Wiley & Sons, 2009.
- [136] Aimo Törn, Montaz M Ali, and Sami Viitanen. Stochastic global optimization: Problem classes and solution techniques. *Journal of Global Optimization*, 14(4):437–447, 1999.
- [137] Henrik Valeur. *India: the Urban Transition - A Case Study of Development Urbanism*. The Architectural Publisher, 2014.
- [138] Frans van den Bergh and Andries Petrus Engelbrecht. A cooperative approach to particle swarm optimization. 8(3):225–239, 2004.
- [139] Joaquin Vanschoren. *Understanding Machine Learning Performance With Experiment Databases*. PhD thesis, 2010.
- [140] F. Vavak and T. C. Fogarty. A comparative study of steady state and generational genetic algorithms for use in nonstationary environments. In Terence C. Fogarty, editor, *Evolutionary Computing*, pages 297–304, Berlin, Heidelberg, 1996. Springer Berlin Heidelberg. ISBN 978-3-540-70671-7.
- [141] Michael D Vose. Continuous lunches are not free. *arXiv preprint arXiv:1507.00581*, 2015.
- [142] Ron Wehrens and Lutgarde MC Buydens. Classical and nonclassical optimization methods. *Encyclopedia of Analytical Chemistry*.
- [143] Darrell Whitley. Genetic algorithms and evolutionary computing. *Van Nostrand's Scientific Encyclopedia*, 2002.
- [144] Darrell Whitley and Joan Kauth. GENITOR: A Different Genetic Algorithm. In *Proceedings of the 1988 Rocky Mountain Conference on Artificial Intelligence*, pages 118–130. Computer Science Department, Colorado State University, 1988.
- [145] David H Wolpert and William G Macready. No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1):67–82, 1997.
- [146] David H Wolpert, William G Macready, et al. No free lunch theorems for search. Technical report, Technical Report SFI-TR-95-02-010, Santa Fe Institute, 1995.
- [147] Jonathan Wright and Ali Alajmi. Efficient genetic algorithm sets for optimizing constrained building design problem. *International Journal of Sustainable Built Environment*, 5(1):123–131, 2016.

- [148] Jianbin Xin, Guimin Chen, and Yubao Hai. A particle swarm optimizer with multi-stage linearly-decreasing inertia weight. In *Computational Sciences and Optimization, 2009. CSO 2009. International Joint Conference on*, volume 1, pages 505–508. IEEE, 2009.
- [149] S. Xu and Y. Rahmat-Samii. Boundary conditions in particle swarm optimization revisited. *IEEE Transactions on Antennas and Propagation*, 55(3):760–765, March 2007. ISSN 0018-926X. doi: 10.1109/TAP.2007.891562.
- [150] Zhenyu Yang, Ke Tang, and Xin Yao. Differential evolution for high-dimensional function optimisation. In *2007 IEEE Congress on Evolutionary Computation*, pages 3523–3530, 2007.
- [151] Zhenyu Yang, Ke Tang, and Xin Yao. Large scale evolutionary optimisation using cooperative coevolution. *Information Sciences*, 178(15):2985–2999, 2008.
- [152] z. Cao, L. Wang, Y. Shi, X. Hei, X. Rong, Q. Jiang, and H. Li. An effective cooperative coevolution framework integrating global and local search for large scale optimization problems. In *2015 IEEE Congress on Evolutionary Computation (CEC)*, pages 1986–1993, May 2015. doi: 10.1109/CEC.2015.7257129.

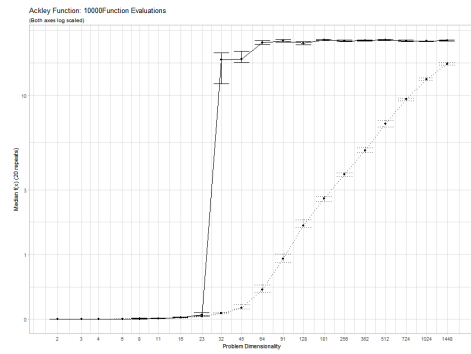
Part V

APPENDICES

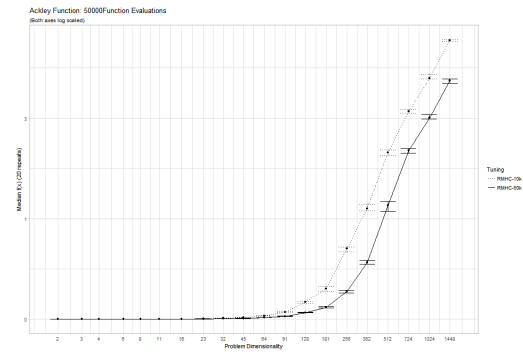


APPENDIX A: 10,000 AND 50,000 EVALUATION COMPARISON PLOTS

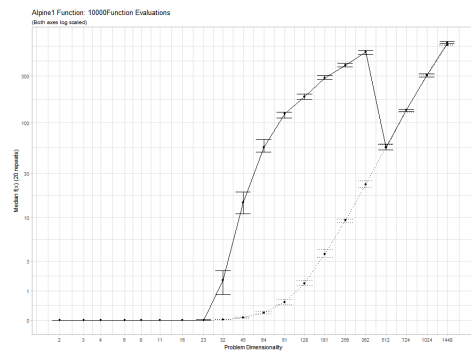
RANDOM MUTATION HILL CLIMBING (RMHC)



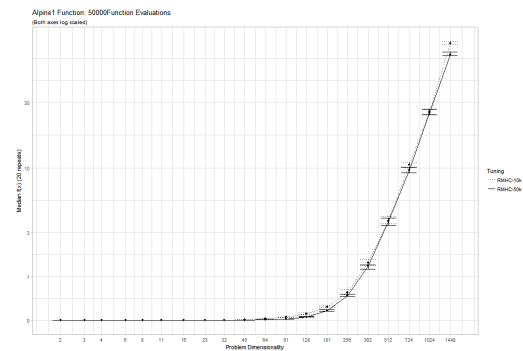
Ackley Function at 10,000 Evaluations



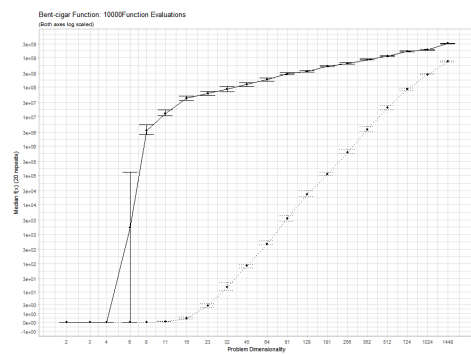
Ackley Function at 50,000 Evaluations



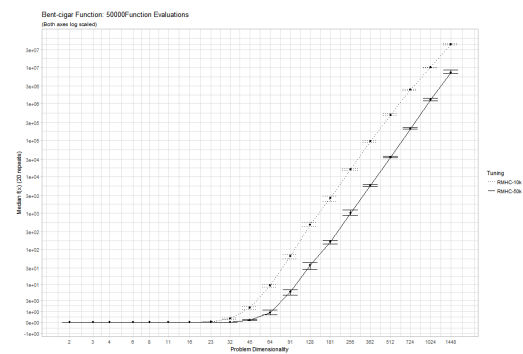
Alpine no.1 Function at 10,000 Evaluations



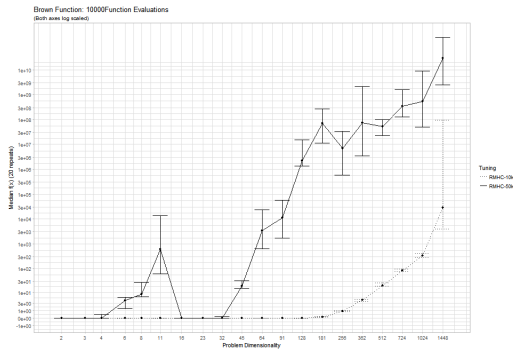
Alpine no.1 Function at 50,000 Evaluations



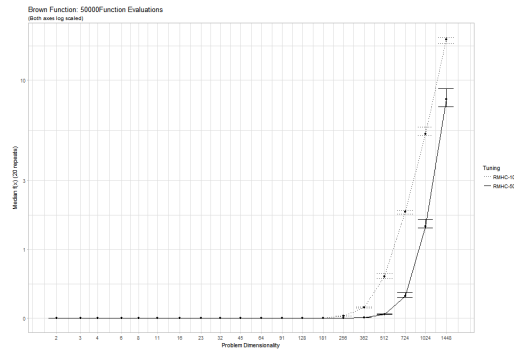
Bent Cigar Function at 10,000 Evaluations



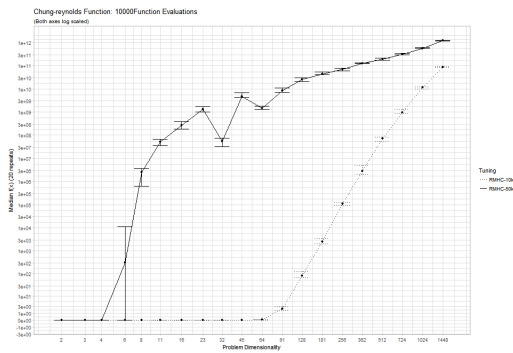
Bent Cigar Function at 50,000 Evaluations



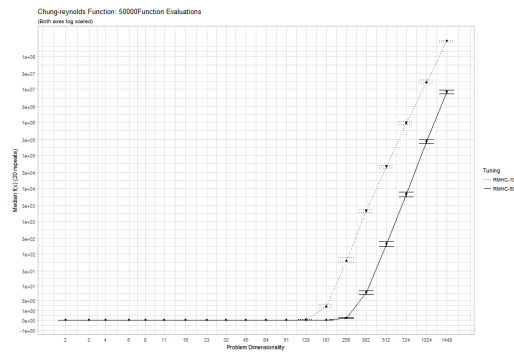
Brown Function at 10,000 Evaluations



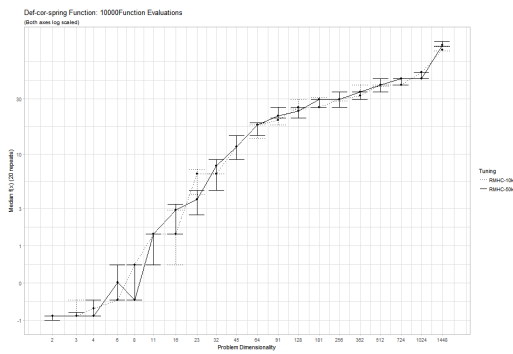
Brown Function at 50,000 Evaluations



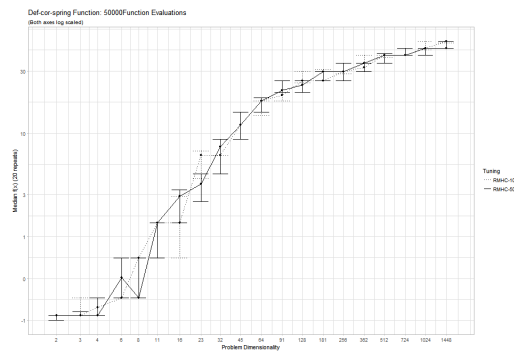
Chung-Reynolds Function at 10,000 Evaluations



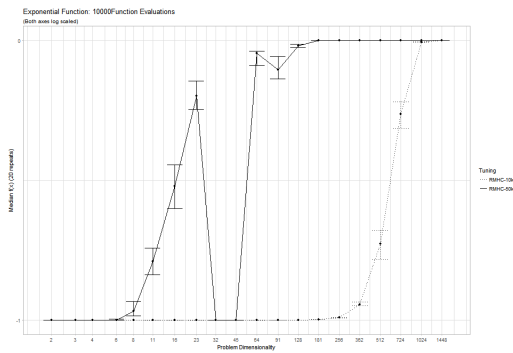
Chung-Reynolds Function at 50,000 Evaluations



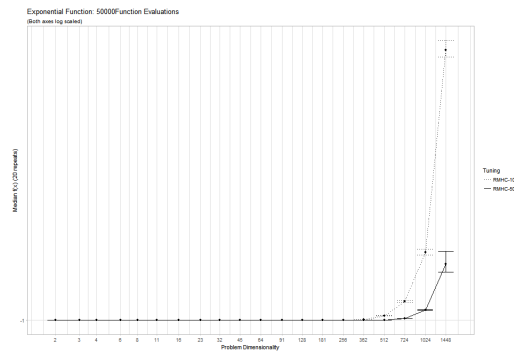
Deflected Corrugated Spring Function at 10,000 Evaluations



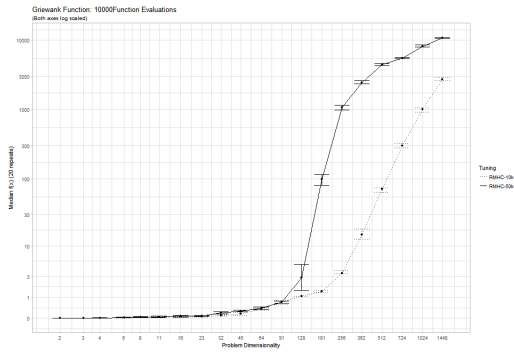
Deflected Corrugated Spring Function at 50,000 Evaluations



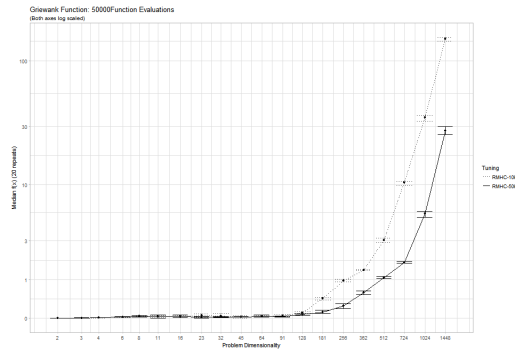
Exponential Function at 10,000 Evaluations



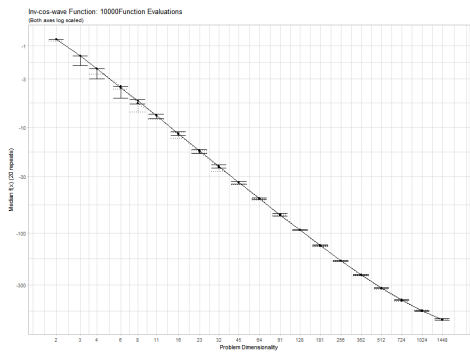
Exponential Function at 50,000 Evaluations



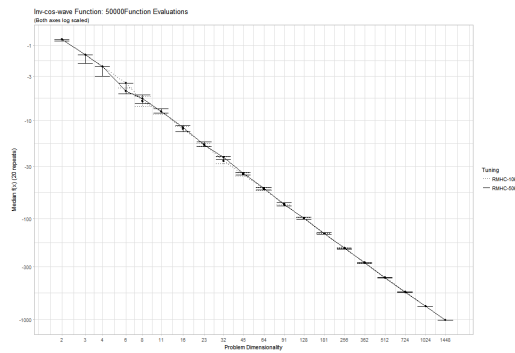
Griewank Function at 10,000 Evaluations



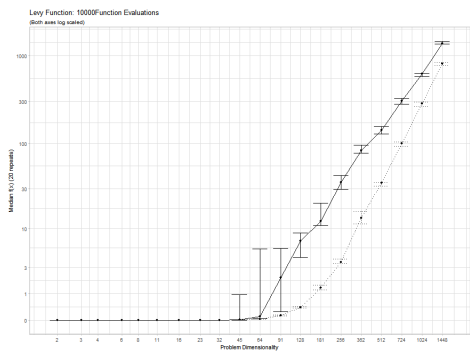
Griewank Function at 50,000 Evaluations



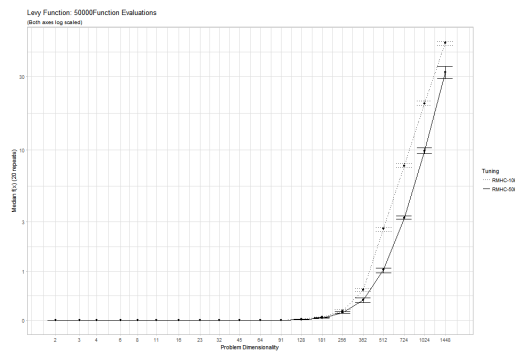
Inverted Cosine Wave Function at 10,000 Evaluations



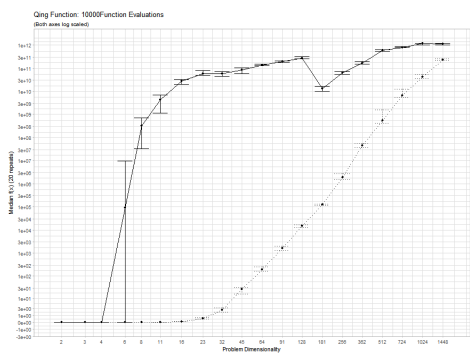
Inverted Cosine Wave Function at 50,000 Evaluations



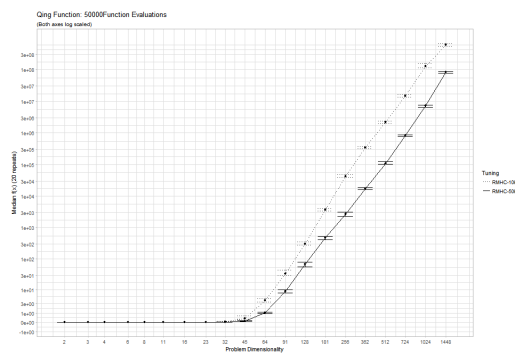
Levy Function at 10,000 Evaluations



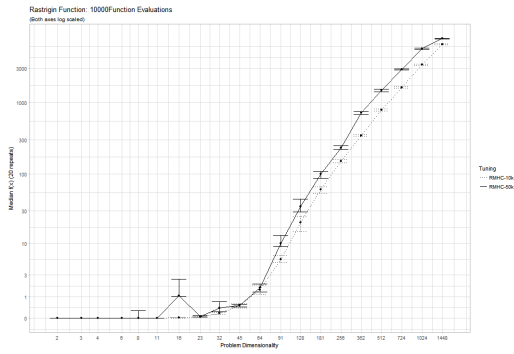
Levy Function at 50,000 Evaluations



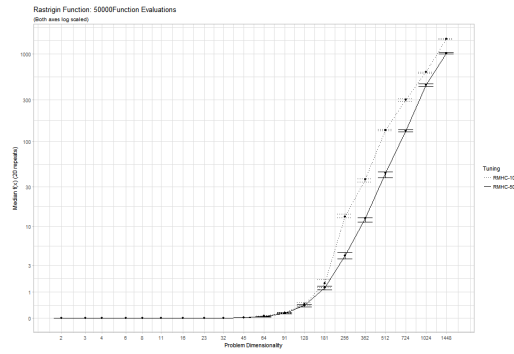
Qing Function at 10,000 Evaluations



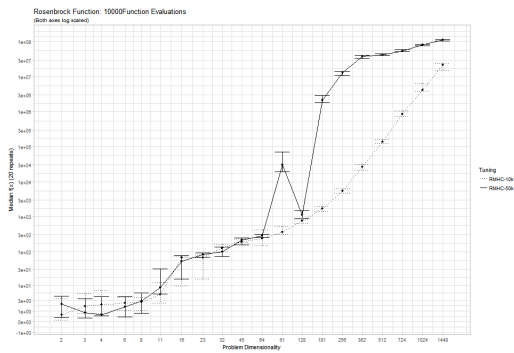
Qing Function at 50,000 Evaluations



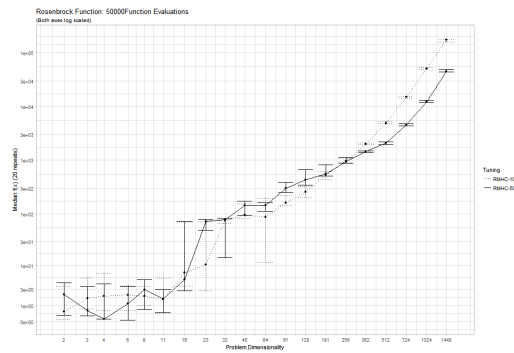
Rastrigin Function at 10,000 Evaluations



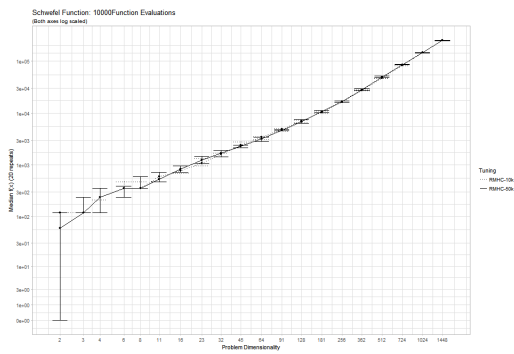
Rastrigin Function at 50,000 Evaluations



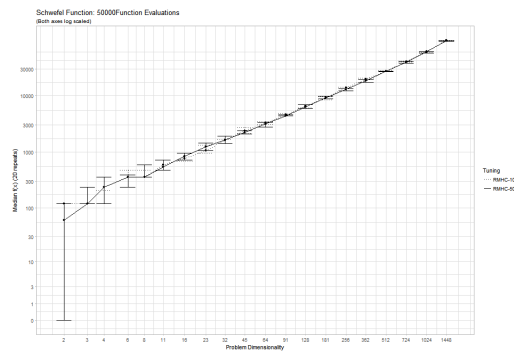
Rosenbrock Function at 10,000 Evaluations



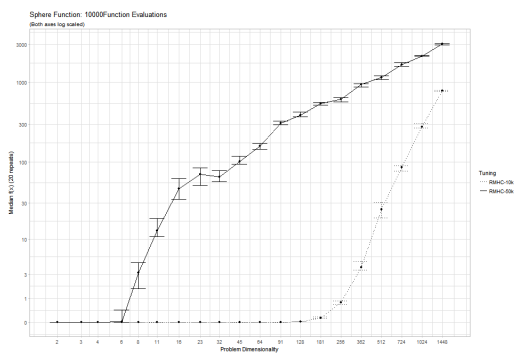
Rosenbrock Function at 50,000 Evaluations



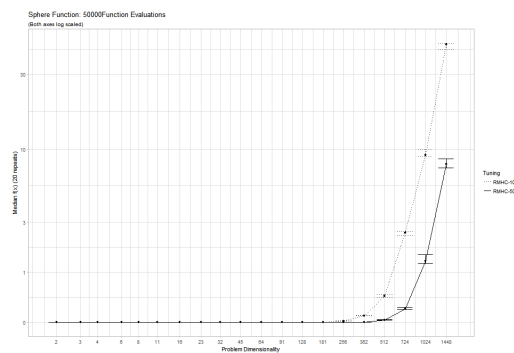
Schwefel Function at 10,000 Evaluations



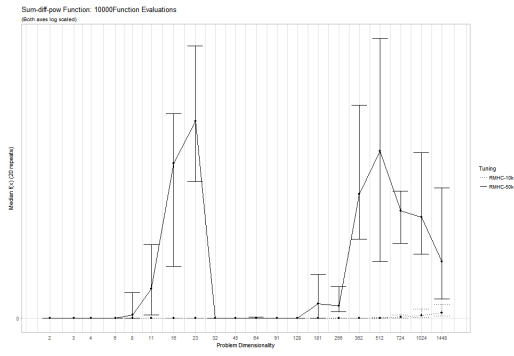
Schwefel Function at 50,000 Evaluations



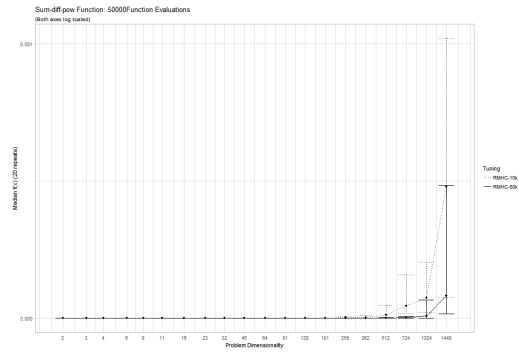
Sphere Function at 10,000 Evaluations



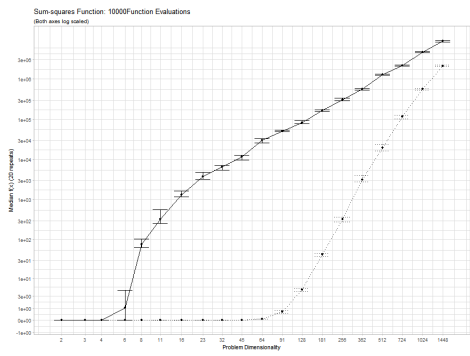
Sphere Function at 50,000 Evaluations



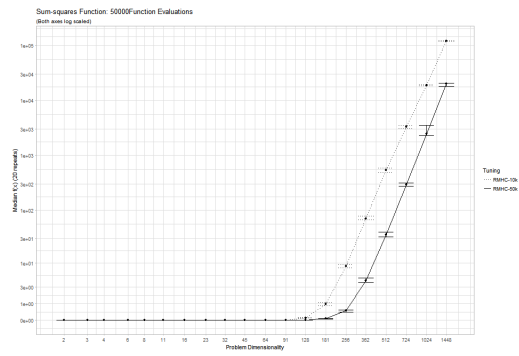
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

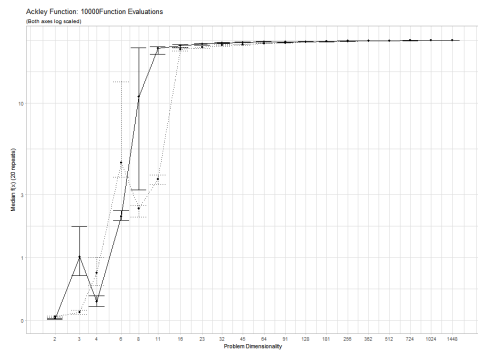


Sum Squares Function at 10,000 Evaluations

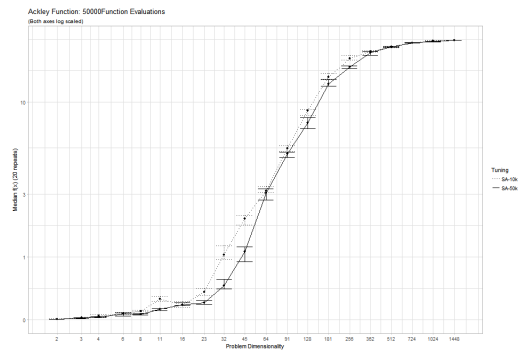


Sum Squares Function at 50,000 Evaluations

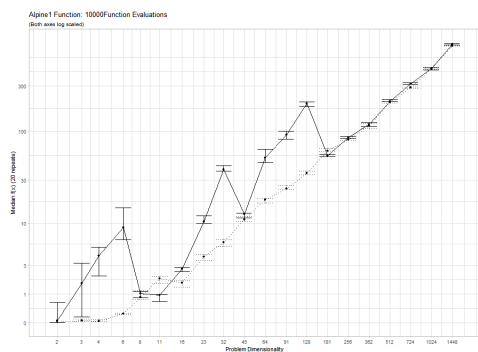
SIMULATED ANNEALING (SA)



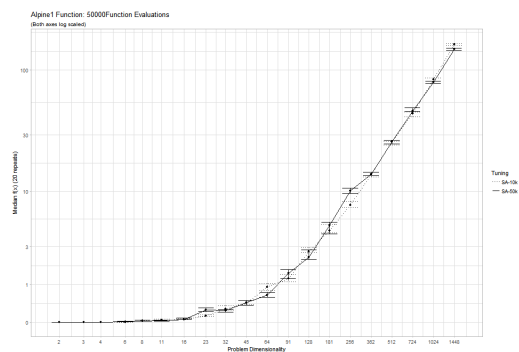
Ackley Function at 10,000 Evaluations



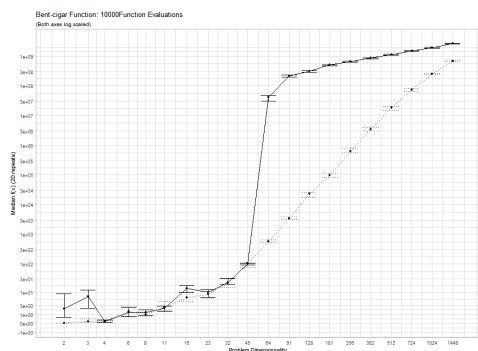
Ackley Function at 50,000 Evaluations



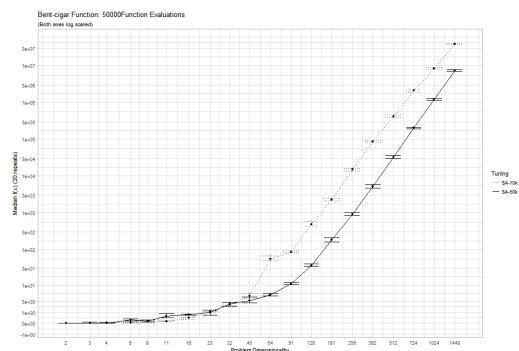
Alpine no.1 Function at 10,000 Evaluations



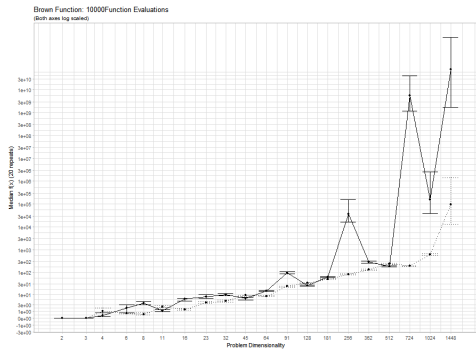
Alpine no.1 Function at 50,000 Evaluations



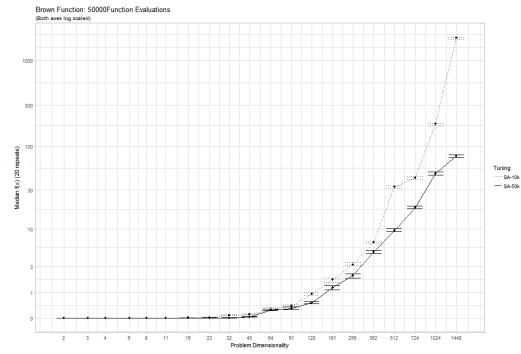
Bent Cigar Function at 10,000 Evaluations



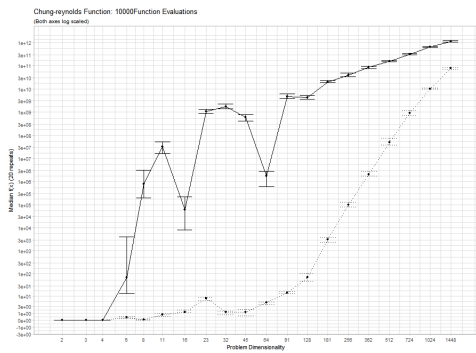
Bent Cigar Function at 50,000 Evaluations



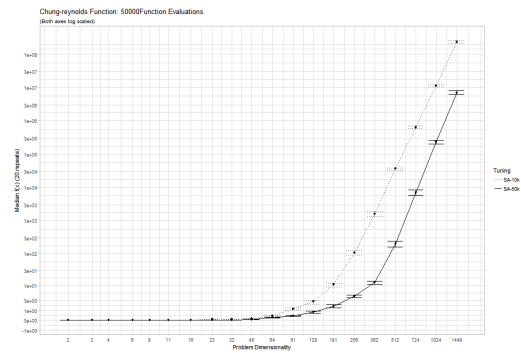
Brown Function at 10,000 Evaluations



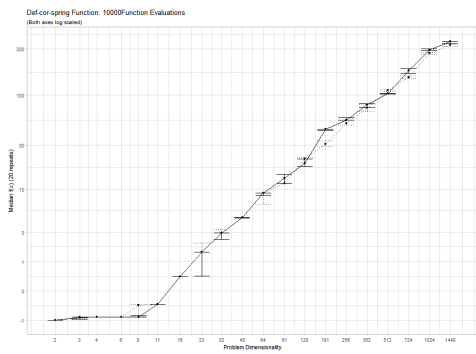
Brown Function at 50,000 Evaluations



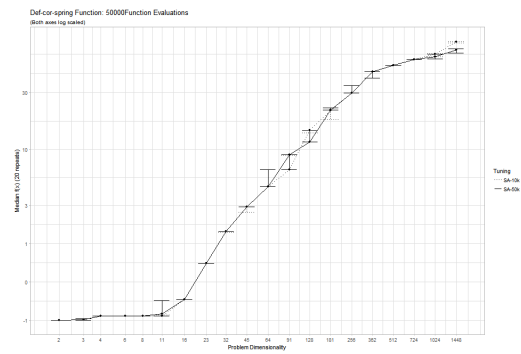
Chung-Reynolds Function at 10,000 Evaluations



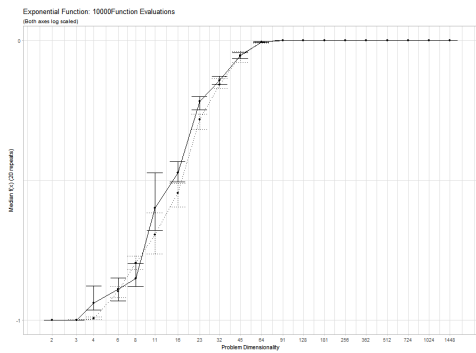
Chung-Reynolds Function at 50,000 Evaluations



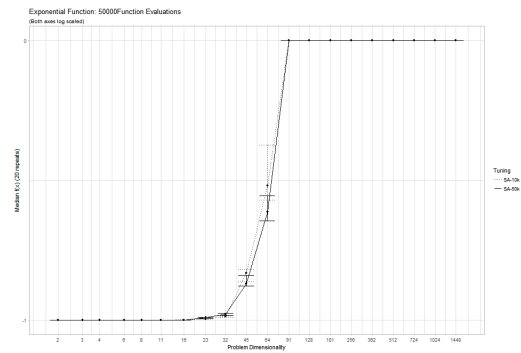
Deflected Corrugated Spring Function at 10,000 Evaluations



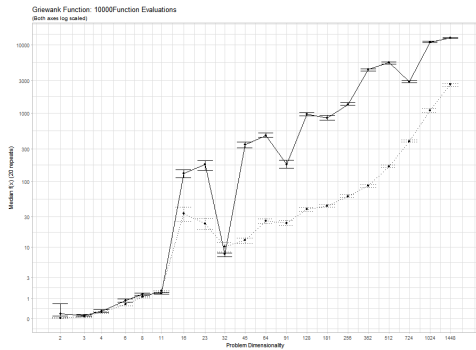
Deflected Corrugated Spring Function at 50,000 Evaluations



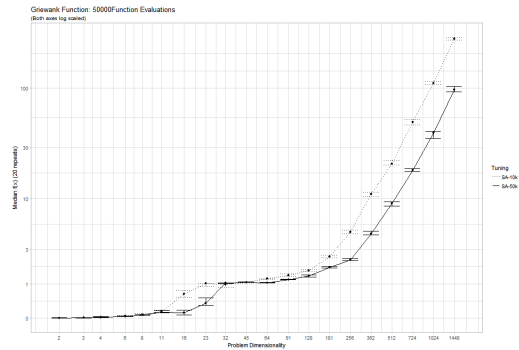
Exponential Function at 10,000 Evaluations



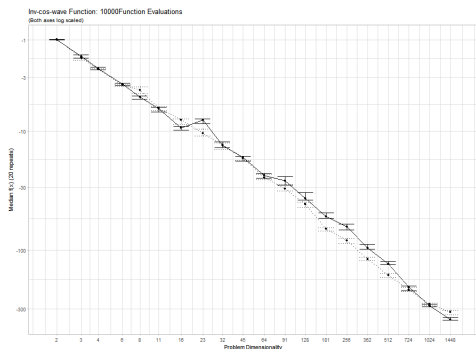
Exponential Function at 50,000 Evaluations



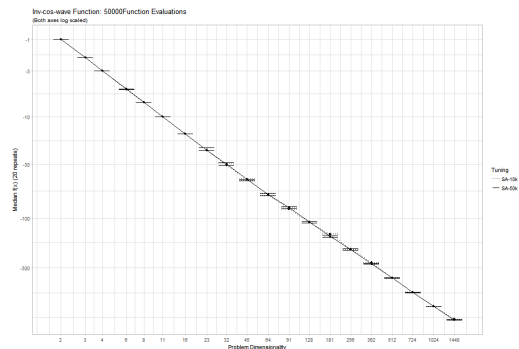
Griewank Function at 10,000 Evaluations



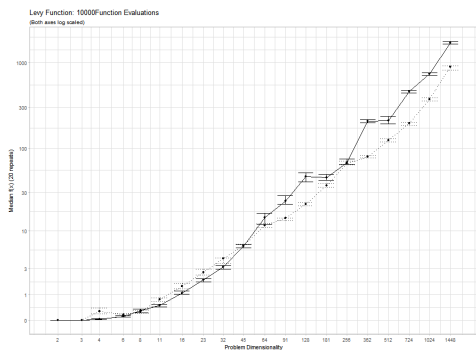
Griewank Function at 50,000 Evaluations



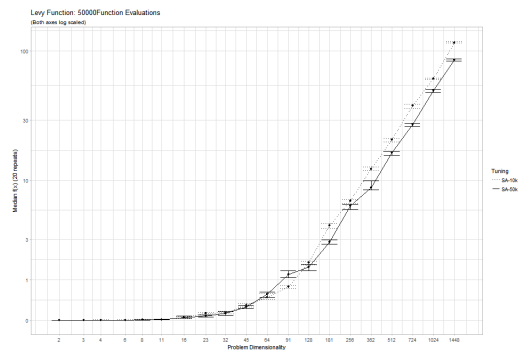
Inverted Cosine Wave Function at 10,000 Evaluations



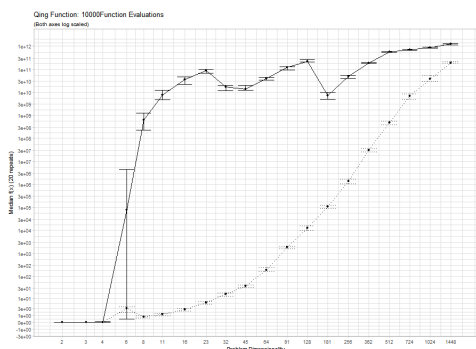
Inverted Cosine Wave Function at 50,000 Evaluations



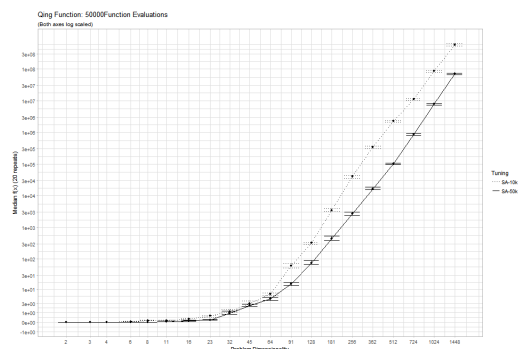
Levy Function at 10,000 Evaluations



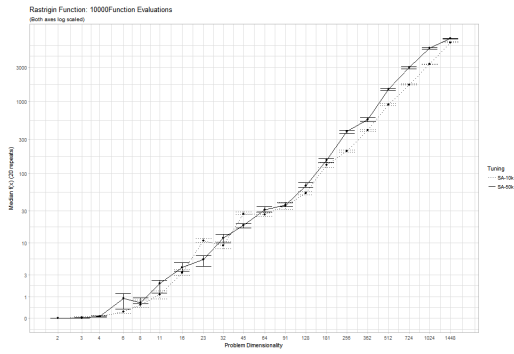
Levy Function at 50,000 Evaluations



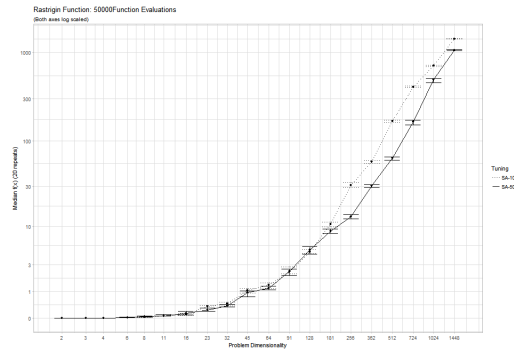
Qing Function at 10,000 Evaluations



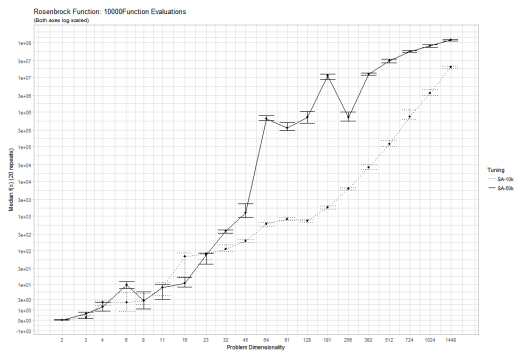
Qing Function at 50,000 Evaluations



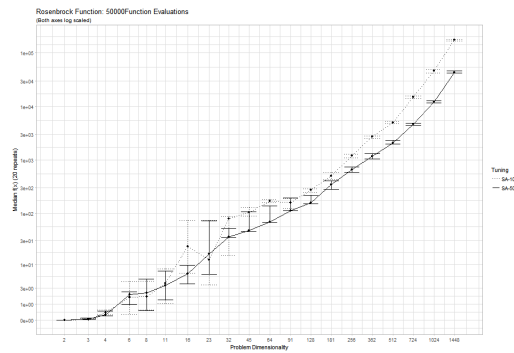
Rastrigin Function at 10,000 Evaluations



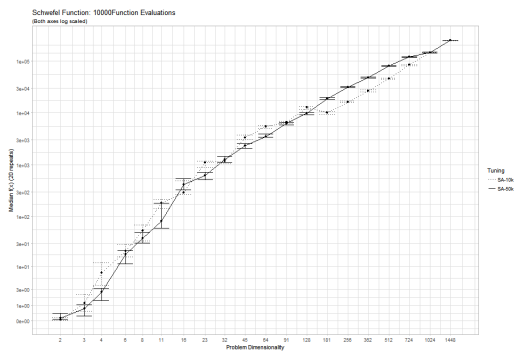
Rastrigin Function at 50,000 Evaluations



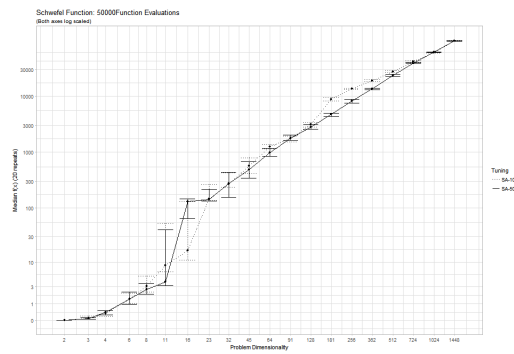
Rosenbrock Function at 10,000 Evaluations



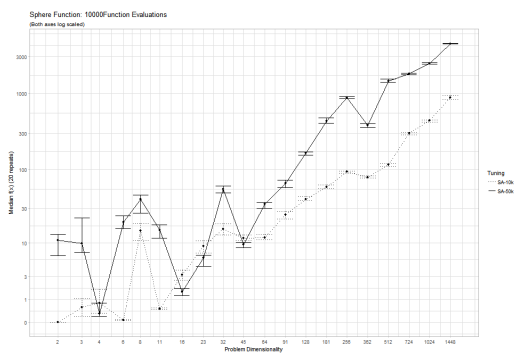
Rosenbrock Function at 50,000 Evaluations



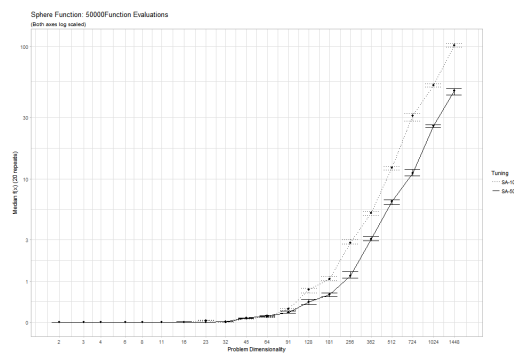
Schwefel Function at 10,000 Evaluations



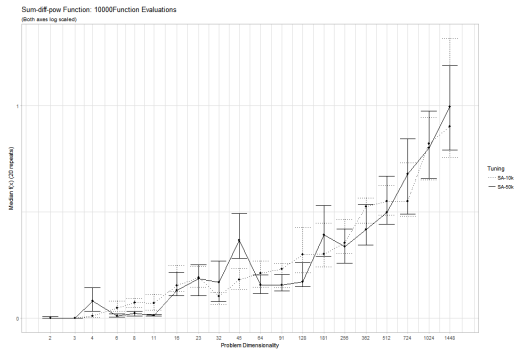
Schwefel Function at 50,000 Evaluations



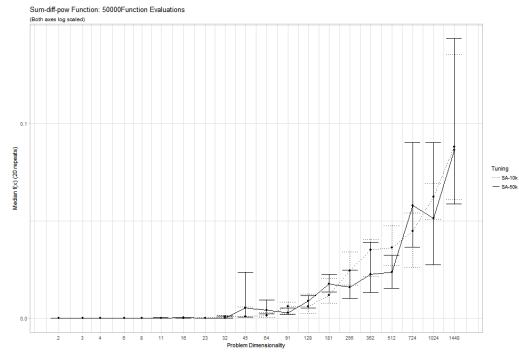
Sphere Function at 10,000 Evaluations



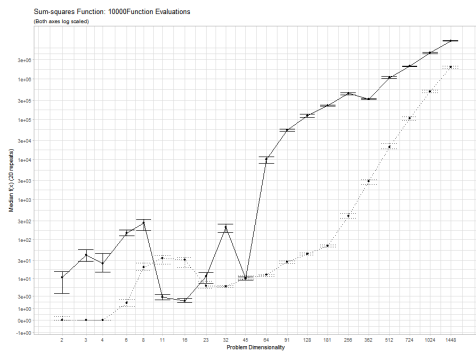
Sphere Function at 50,000 Evaluations



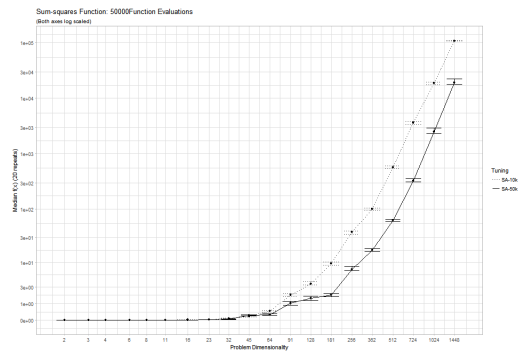
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

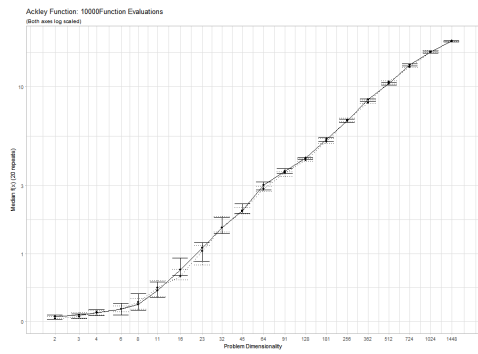


Sum Squares Function at 10,000 Evaluations

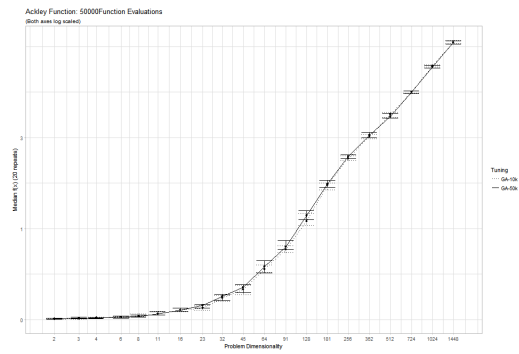


Sum Squares Function at 50,000 Evaluations

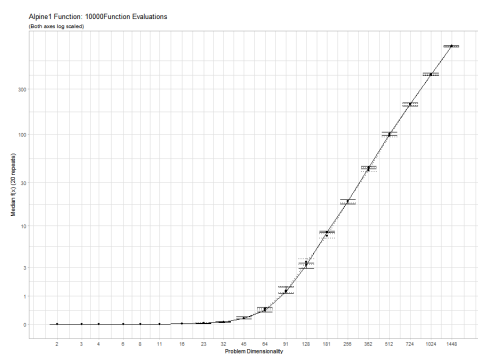
STEADY STATE GENETIC ALGORITHM (SSGA)



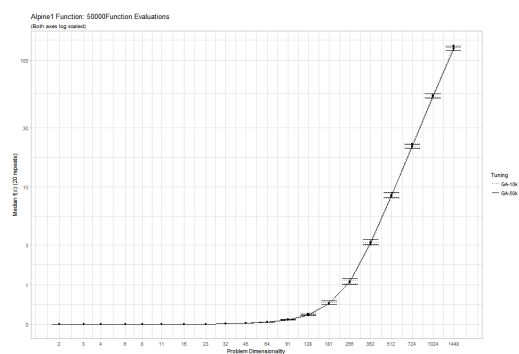
Ackley Function at 10,000 Evaluations



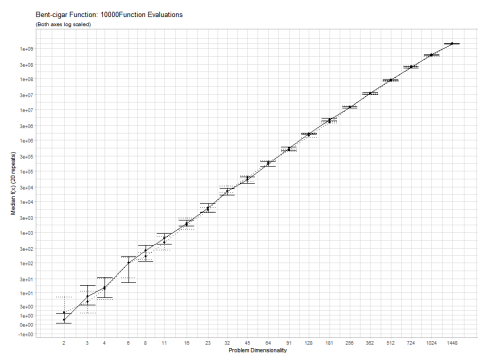
Ackley Function at 50,000 Evaluations



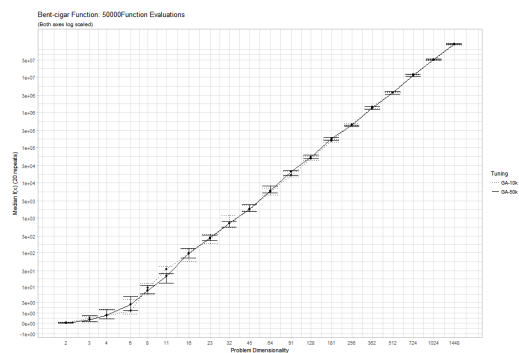
Alpine no.1 Function at 10,000 Evaluations



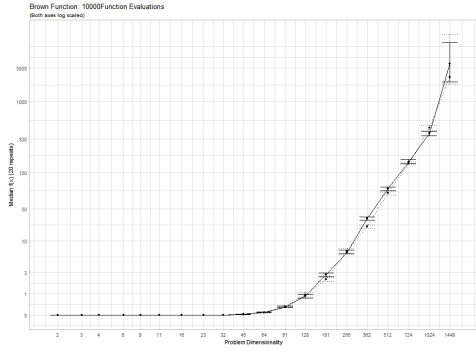
Alpine no.1 Function at 50,000 Evaluations



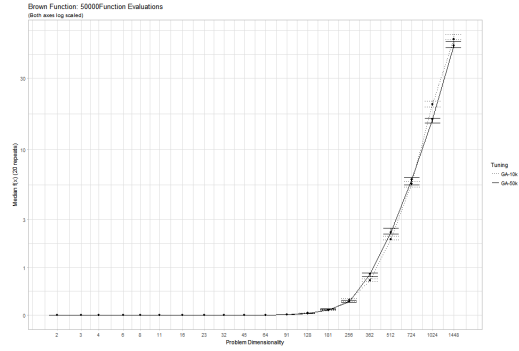
Bent Cigar Function at 10,000 Evaluations



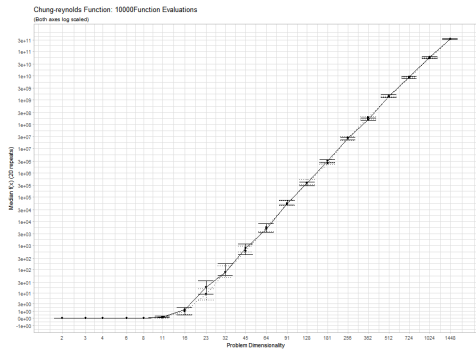
Bent Cigar Function at 50,000 Evaluations



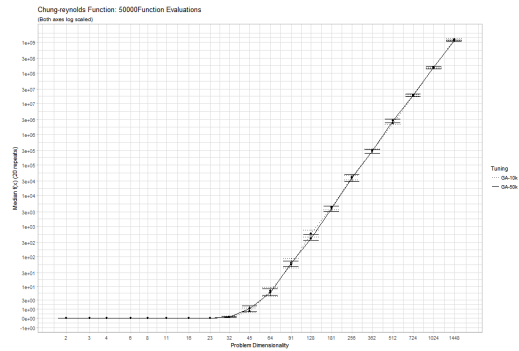
Brown Function at 10,000 Evaluations



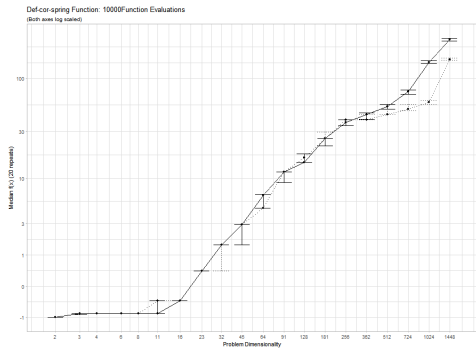
Brown Function at 50,000 Evaluations



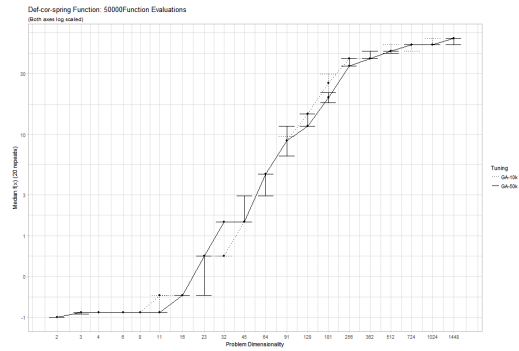
Chung-Reynolds Function at 10,000 Evaluations



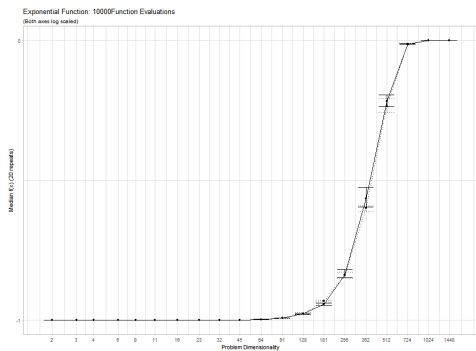
Chung-Reynolds Function at 50,000 Evaluations



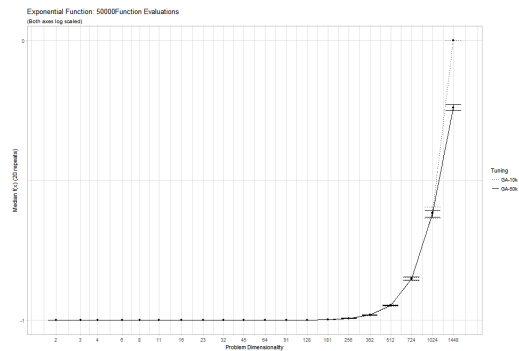
Deflected Corrugated Spring Function at 10,000 Evaluations



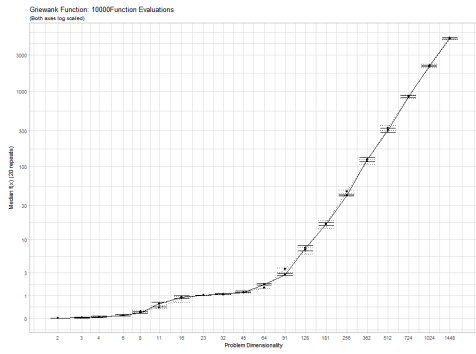
Deflected Corrugated Spring Function at 50,000 Evaluations



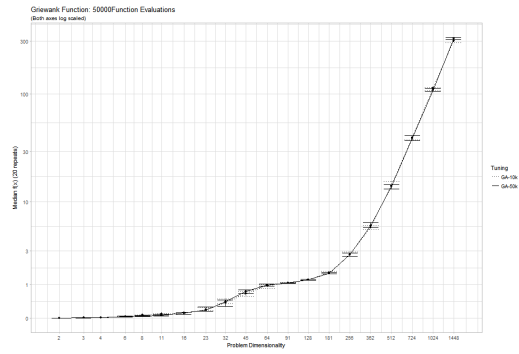
Exponential Function at 10,000 Evaluations



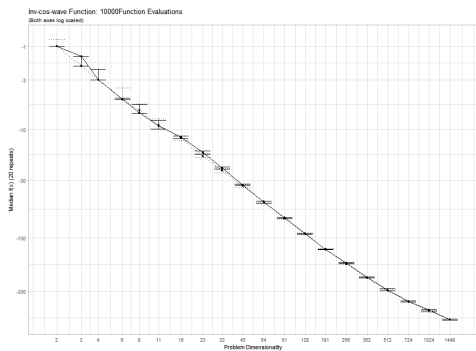
Exponential Function at 50,000 Evaluations



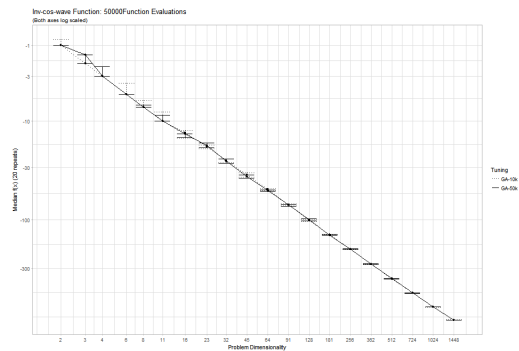
Griewank Function at 10,000 Evaluations



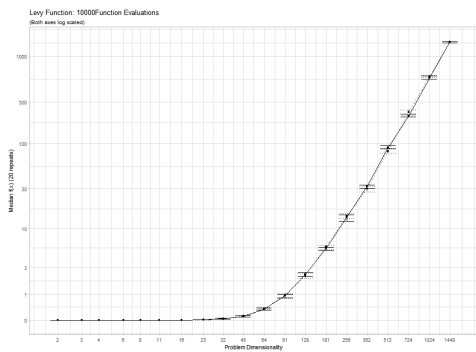
Griewank Function at 50,000 Evaluations



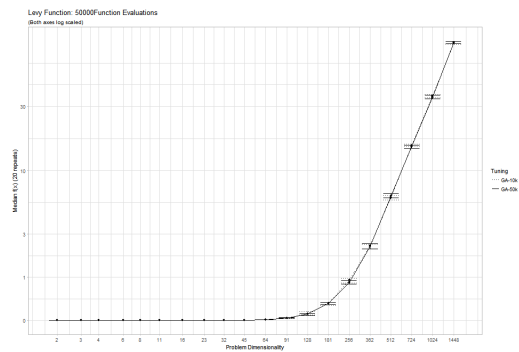
Inverted Cosine Wave Function at 10,000 Evaluations



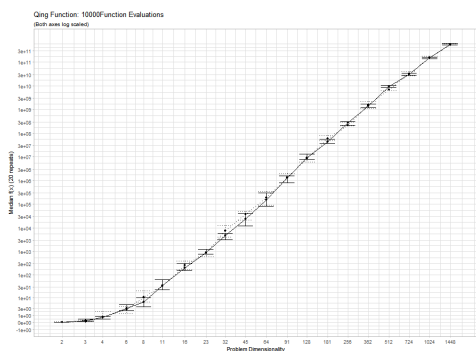
Inverted Cosine Wave Function at 50,000 Evaluations



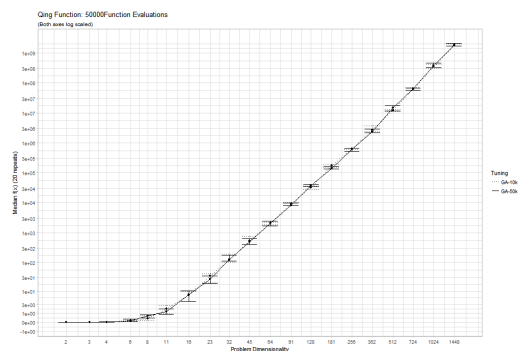
Levy Function at 10,000 Evaluations



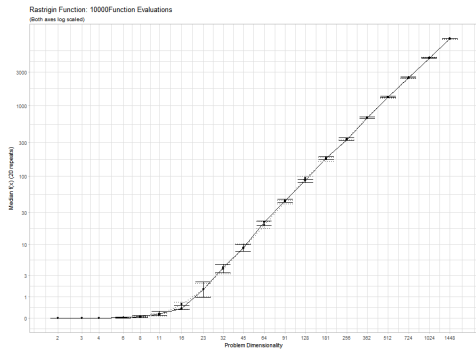
Levy Function at 50,000 Evaluations



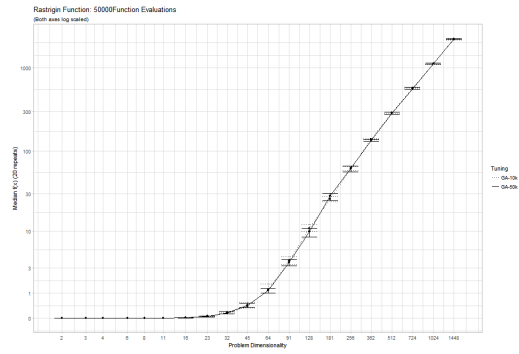
Qing Function at 10,000 Evaluations



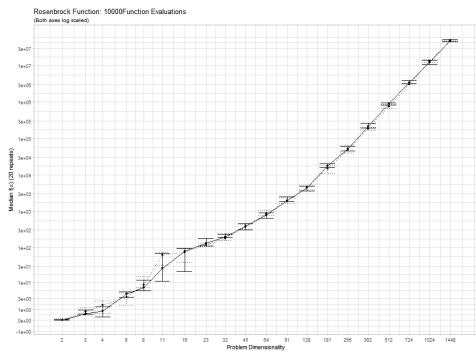
Qing Function at 50,000 Evaluations



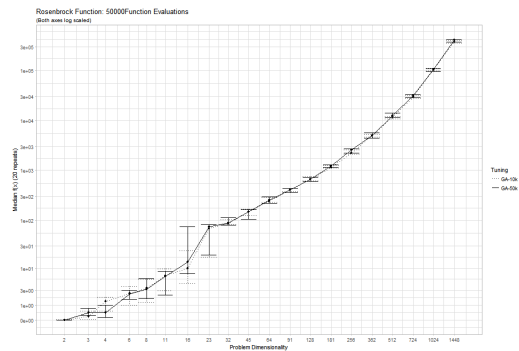
Rastrigin Function at 10,000 Evaluations



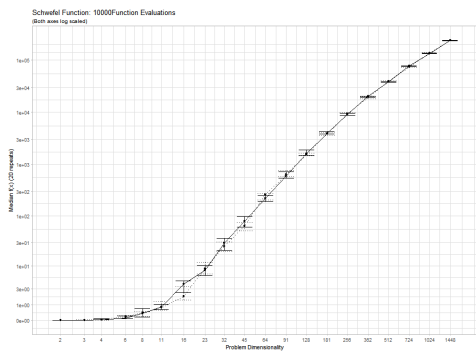
Rastrigin Function at 50,000 Evaluations



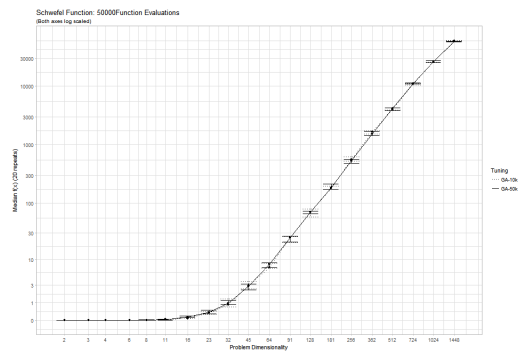
Rosenbrock Function at 10,000 Evaluations



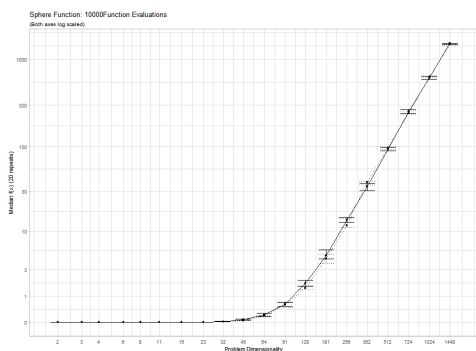
Rosenbrock Function at 50,000 Evaluations



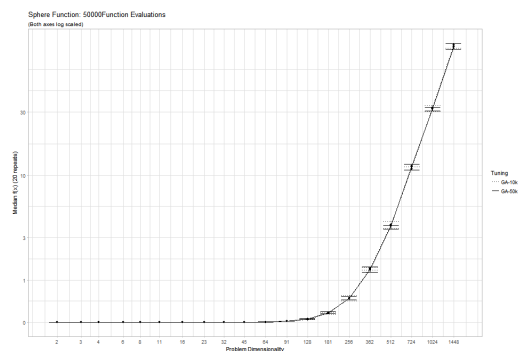
Schwefel Function at 10,000 Evaluations



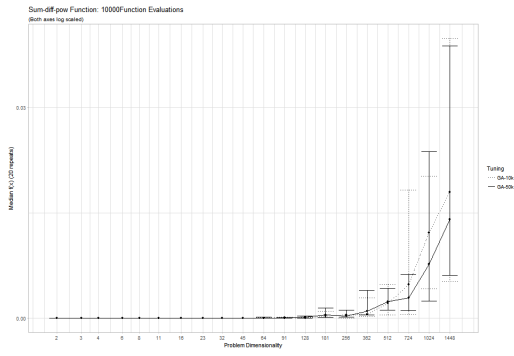
Schwefel Function at 50,000 Evaluations



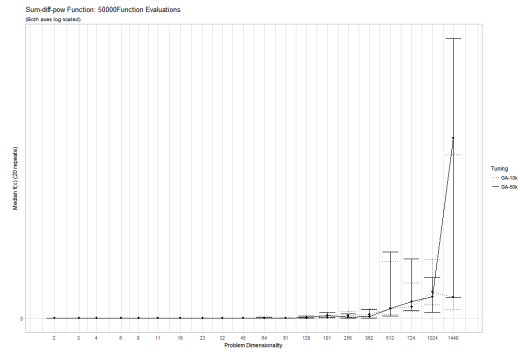
Sphere Function at 10,000 Evaluations



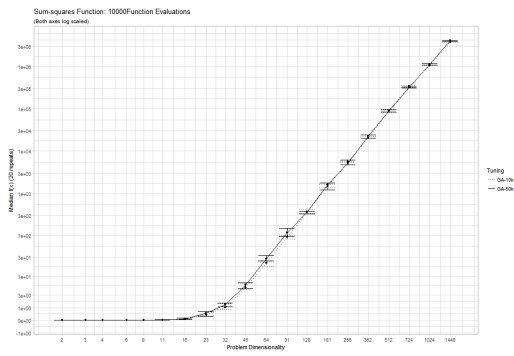
Sphere Function at 50,000 Evaluations



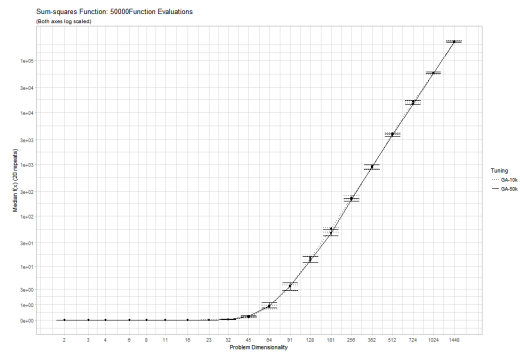
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

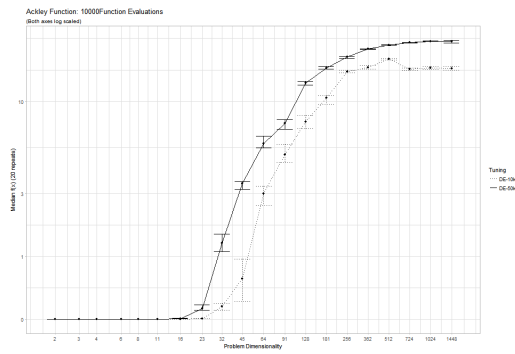


Sum Squares Function at 10,000 Evaluations

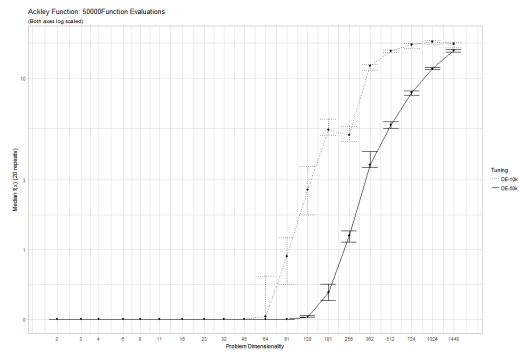


Sum Squares Function at 50,000 Evaluations

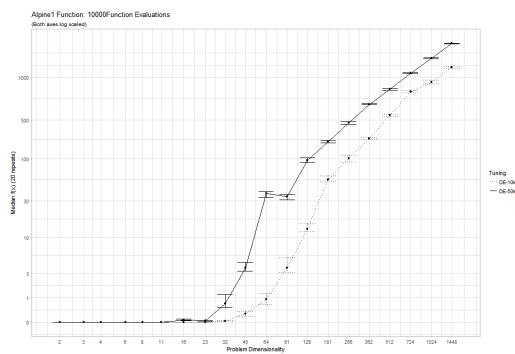
DIFFERENTIAL EVOLUTION (DE)



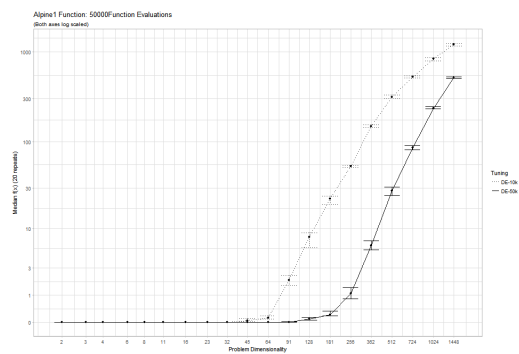
Ackley Function at 10,000 Evaluations



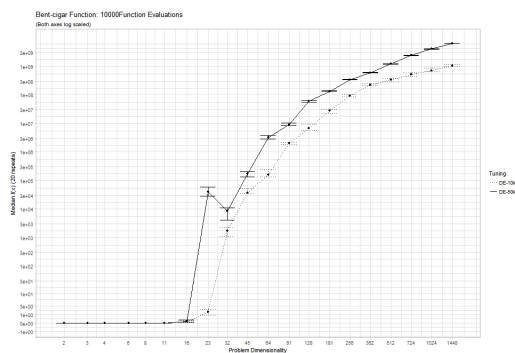
Ackley Function at 50,000 Evaluations



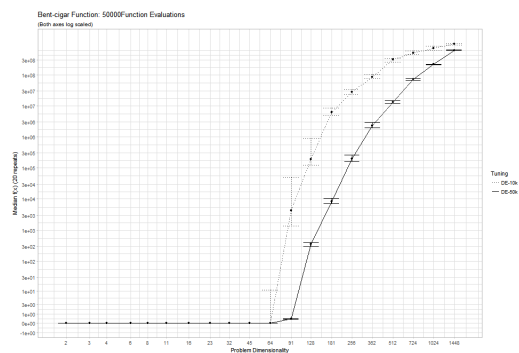
Alpine no.1 Function at 10,000 Evaluations



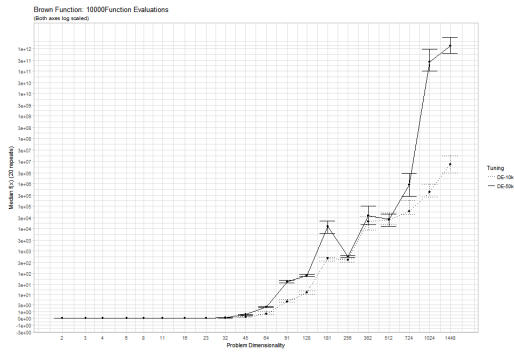
Alpine no.1 Function at 50,000 Evaluations



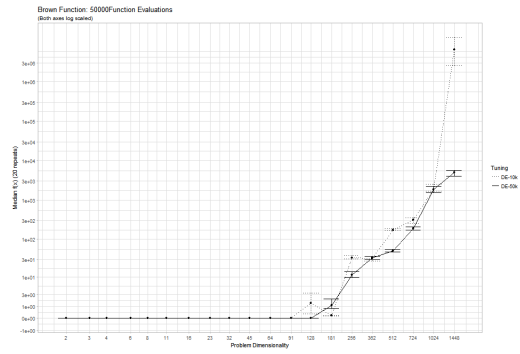
Bent Cigar Function at 10,000 Evaluations



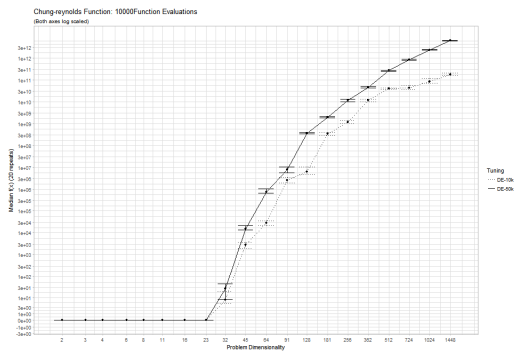
Bent Cigar Function at 50,000 Evaluations



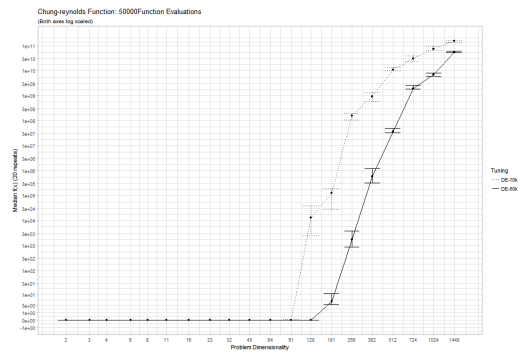
Brown Function at 10,000 Evaluations



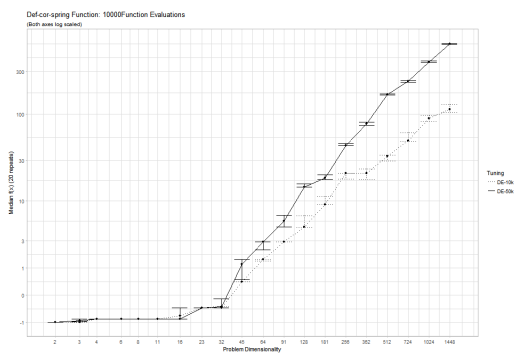
Brown Function at 50,000 Evaluations



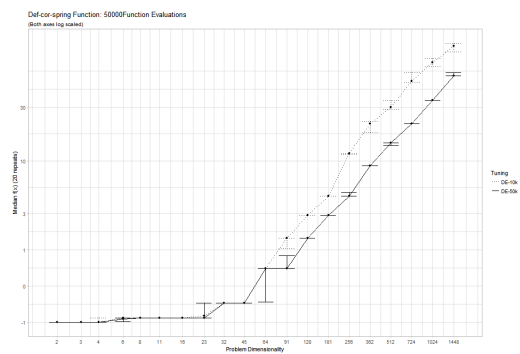
Chung-Reynolds Function at 10,000 Evaluations



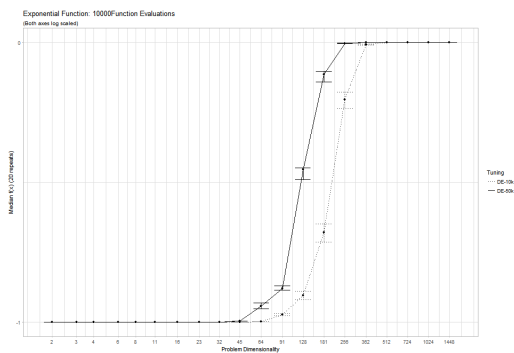
Chung-Reynolds Function at 50,000 Evaluations



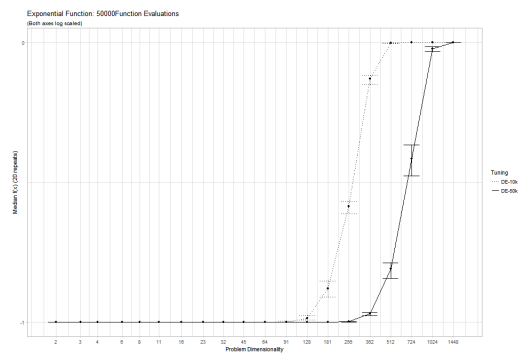
Deflected Corrugated Spring Function at 10,000 Evaluations



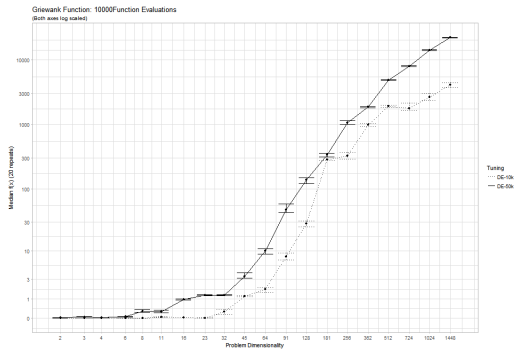
Deflected Corrugated Spring Function at 50,000 Evaluations



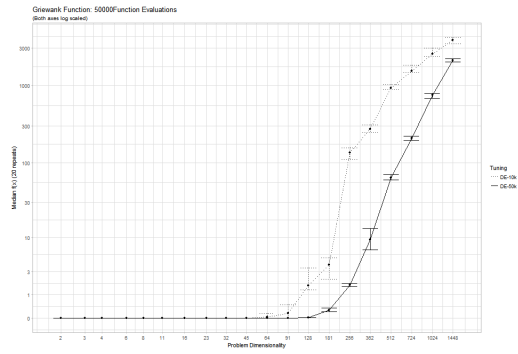
Exponential Function at 10,000 Evaluations



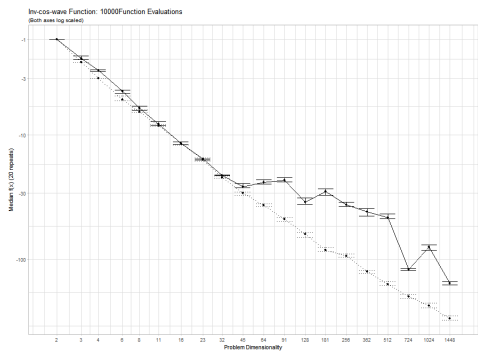
Exponential Function at 50,000 Evaluations



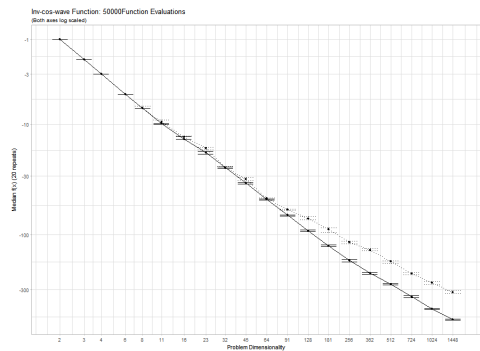
Griewank Function at 10,000 Evaluations



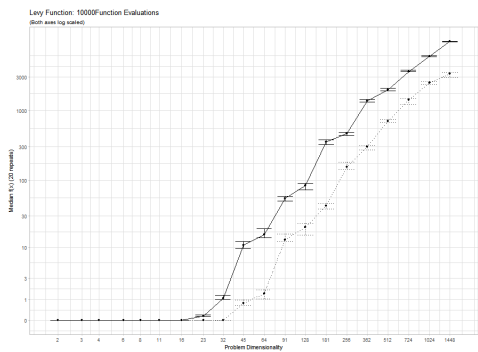
Griewank Function at 50,000 Evaluations



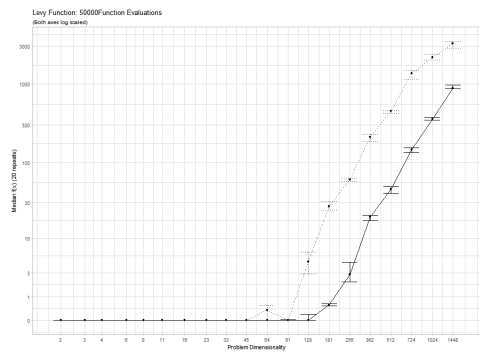
Inverted Cosine Wave Function at 10,000 Evaluations



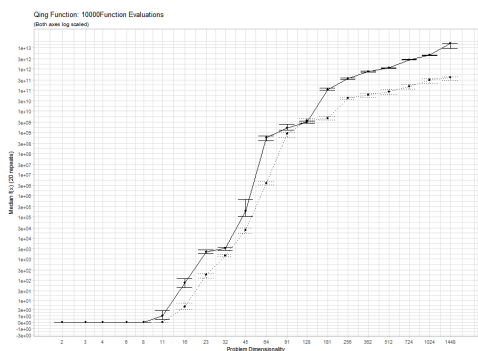
Inverted Cosine Wave Function at 50,000 Evaluations



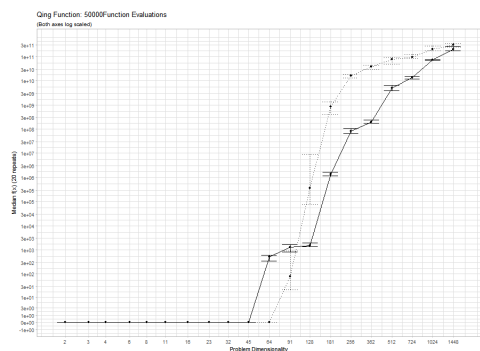
Levy Function at 10,000 Evaluations



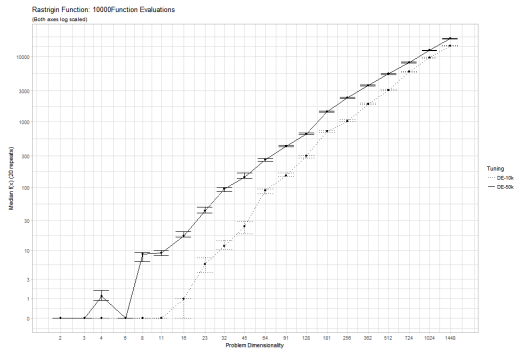
Levy Function at 50,000 Evaluations



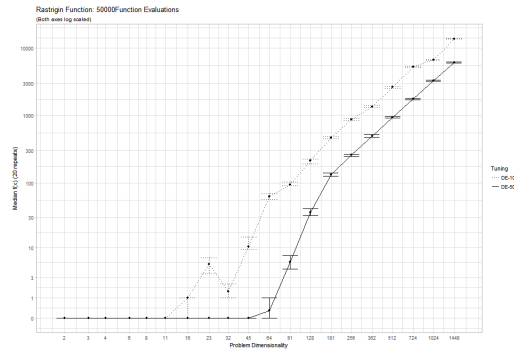
Qing Function at 10,000 Evaluations



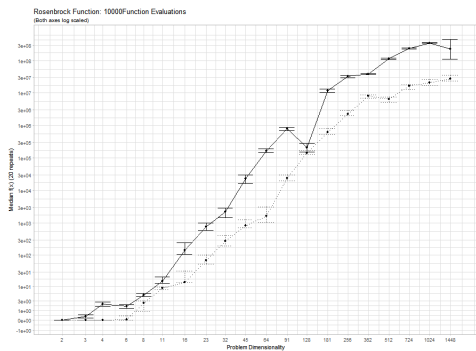
Qing Function at 50,000 Evaluations



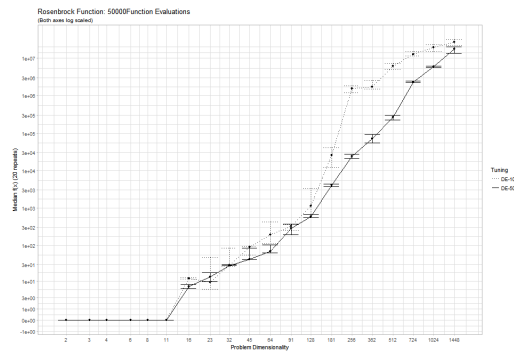
Rastrigin Function at 10,000 Evaluations



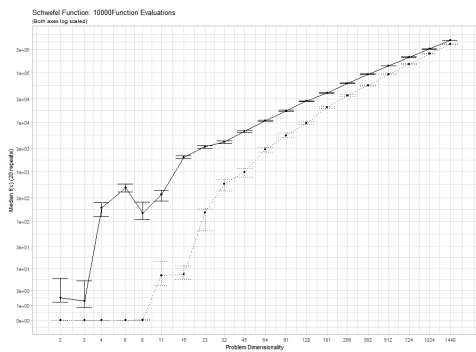
Rastrigin Function at 50,000 Evaluations



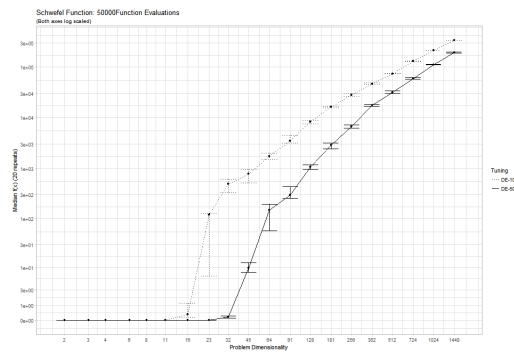
Rosenbrock Function at 10,000 Evaluations



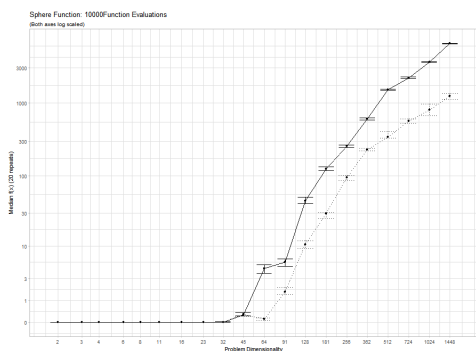
Rosenbrock Function at 50,000 Evaluations



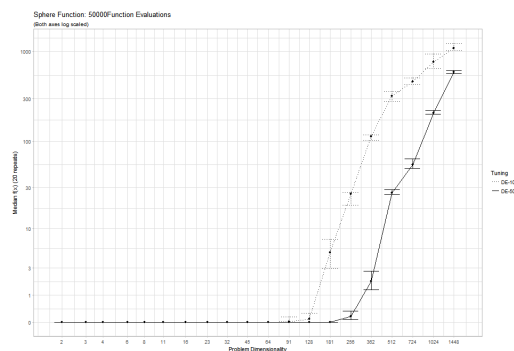
Schwefel Function at 10,000 Evaluations



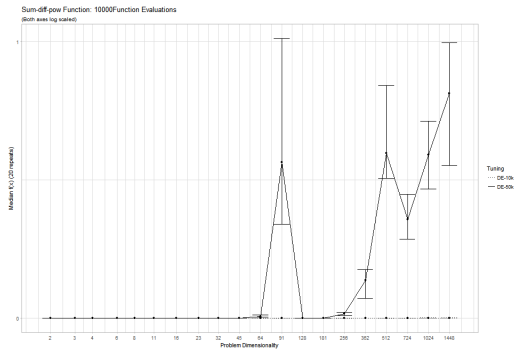
Schwefel Function at 50,000 Evaluations



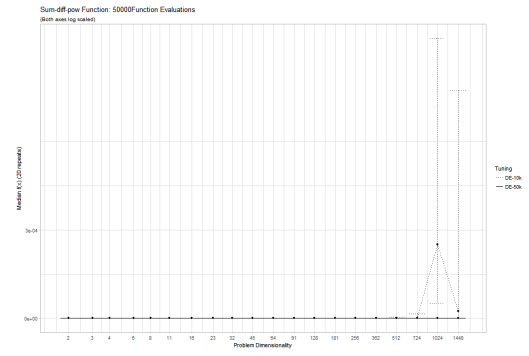
Sphere Function at 10,000 Evaluations



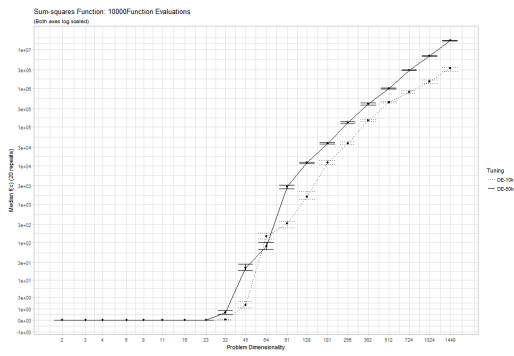
Sphere Function at 50,000 Evaluations



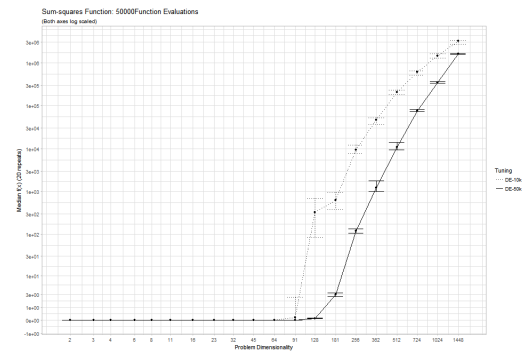
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

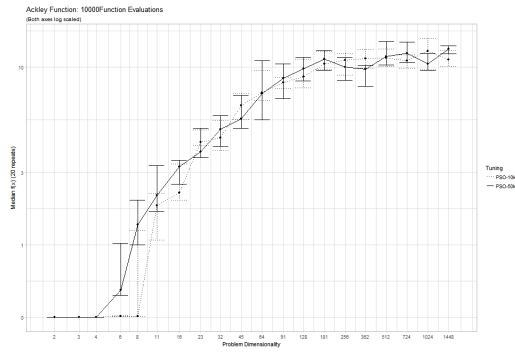


Sum Squares Function at 10,000 Evaluations

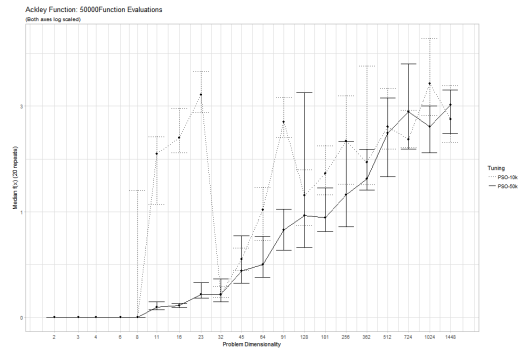


Sum Squares Function at 50,000 Evaluations

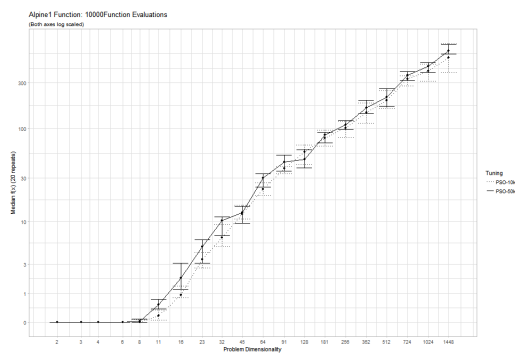
PARTICLE SWARM OPTIMISATION



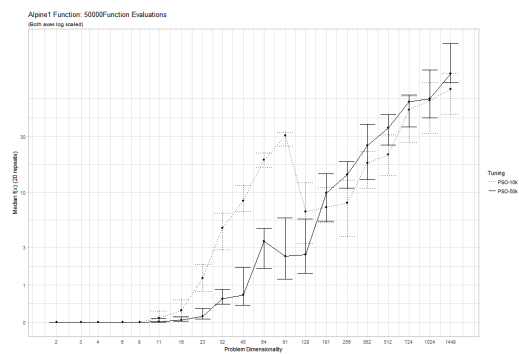
Ackley Function at 10,000 Evaluations



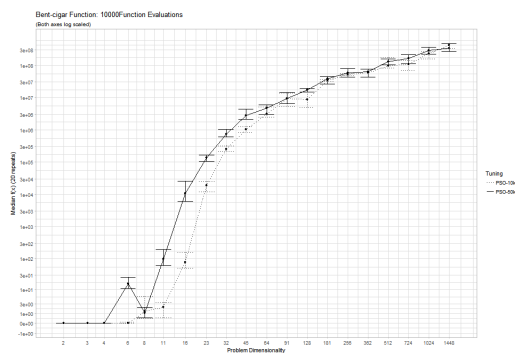
Ackley Function at 50,000 Evaluations



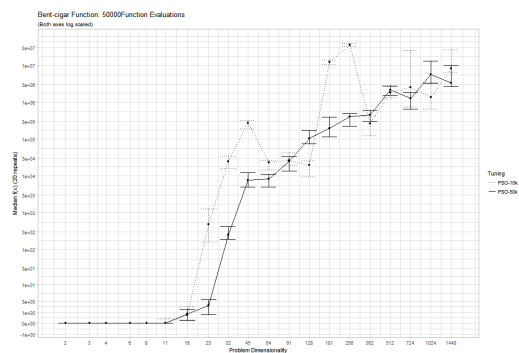
Alpine no.1 Function at 10,000 Evaluations



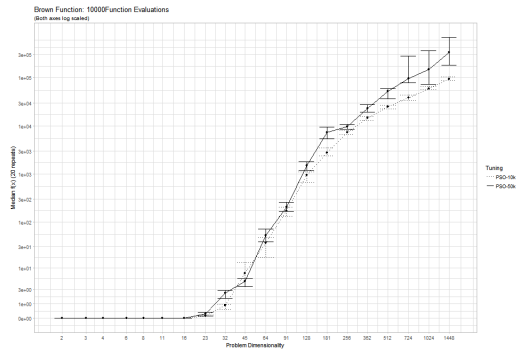
Alpine no.1 Function at 50,000 Evaluations



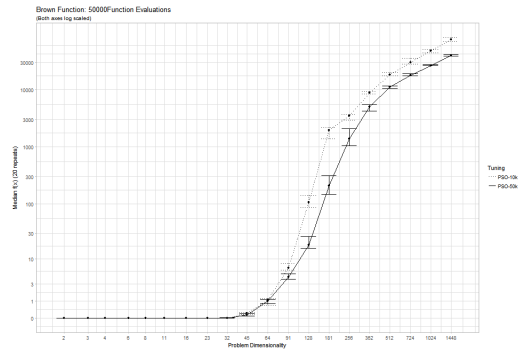
Bent Cigar Function at 10,000 Evaluations



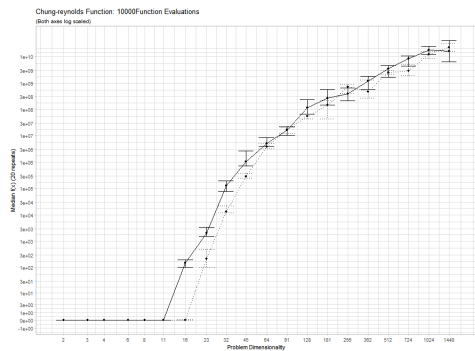
Bent Cigar Function at 50,000 Evaluations



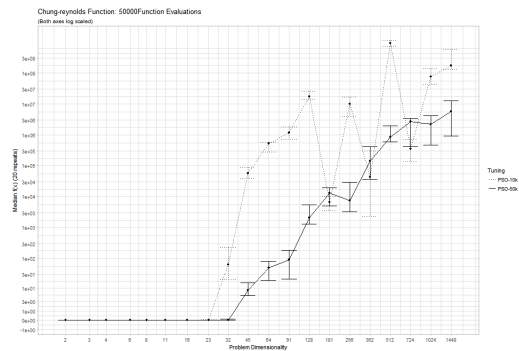
Brown Function at 10,000 Evaluations



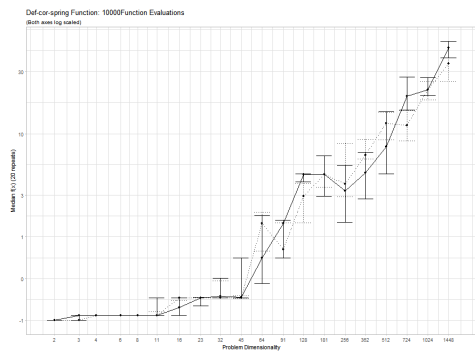
Brown Function at 50,000 Evaluations



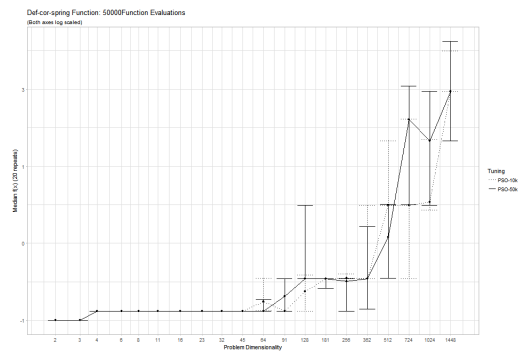
Chung-Reynolds Function at 10,000 Evaluations



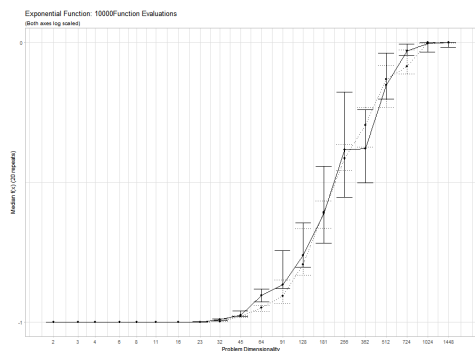
Chung-Reynolds Function at 50,000 Evaluations



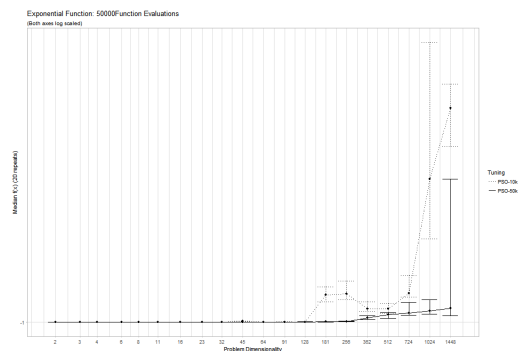
Deflected Corrugated Spring Function at 10,000 Evaluations



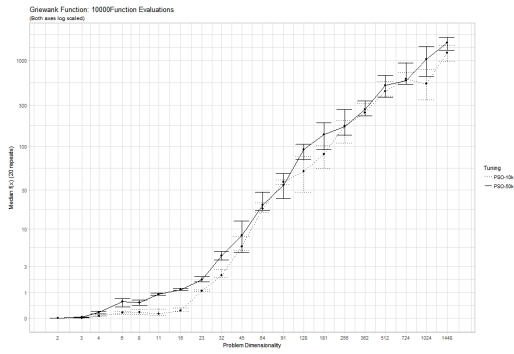
Deflected Corrugated Spring Function at 50,000 Evaluations



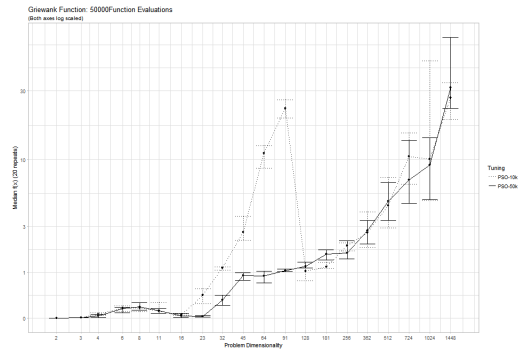
Exponential Function at 10,000 Evaluations



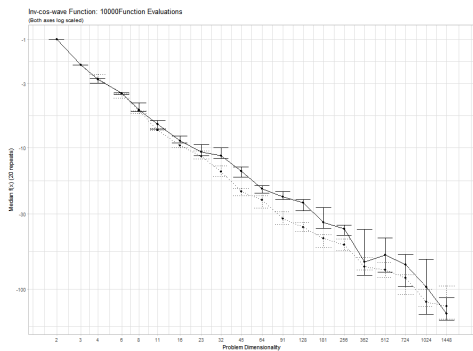
Exponential Function at 50,000 Evaluations



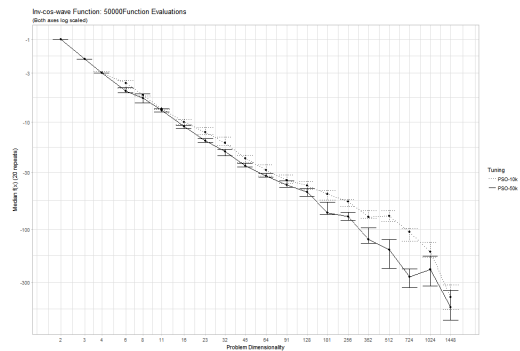
Griewank Function at 10,000 Evaluations



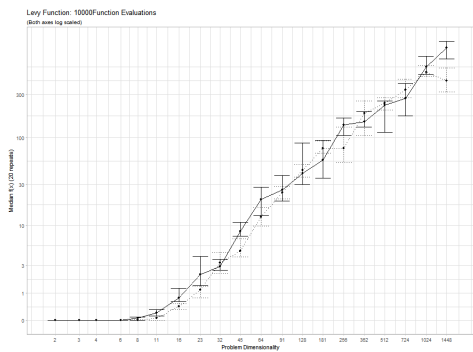
Griewank Function at 50,000 Evaluations



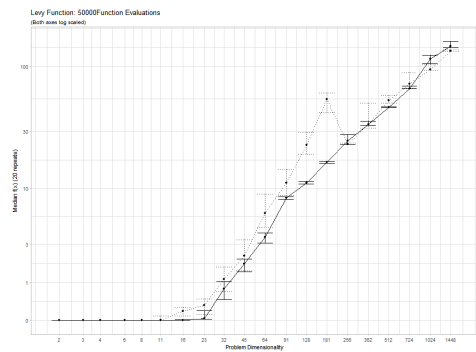
Inverted Cosine Wave Function at 10,000 Evaluations



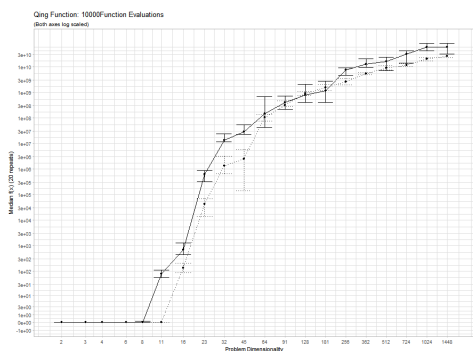
Inverted Cosine Wave Function at 50,000 Evaluations



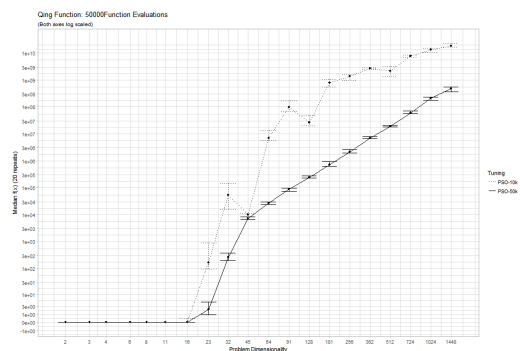
Levy Function at 10,000 Evaluations



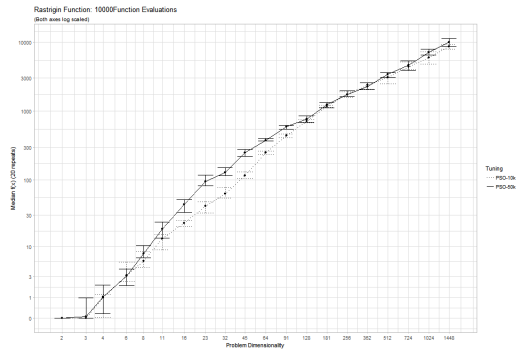
Levy Function at 50,000 Evaluations



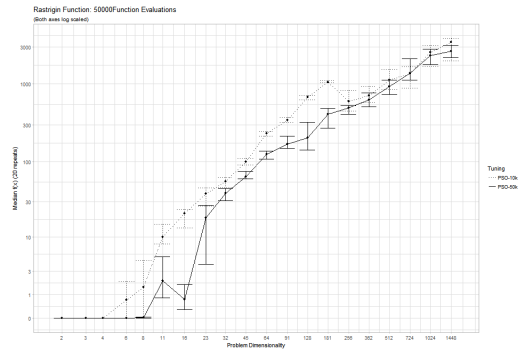
Qing Function at 10,000 Evaluations



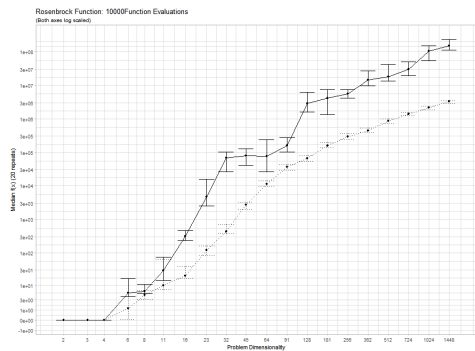
Qing Function at 50,000 Evaluations



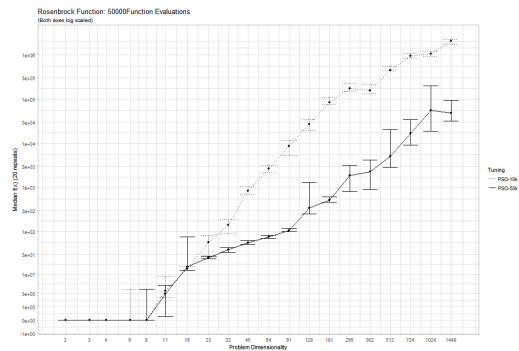
Rastrigin Function at 10,000 Evaluations



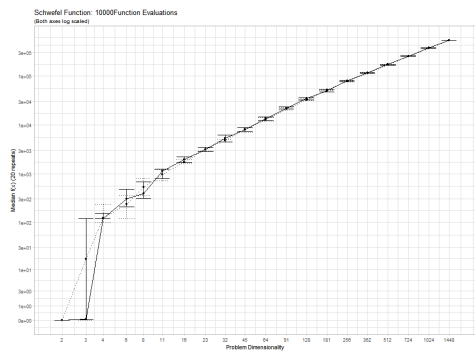
Rastrigin Function at 50,000 Evaluations



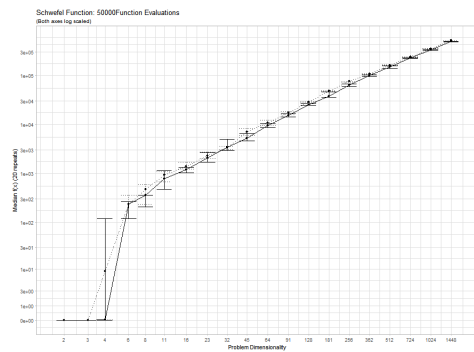
Rosenbrock Function at 10,000 Evaluations



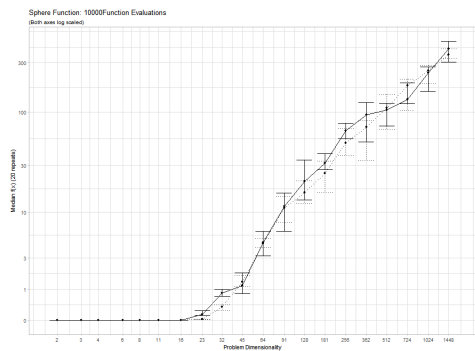
Rosenbrock Function at 50,000 Evaluations



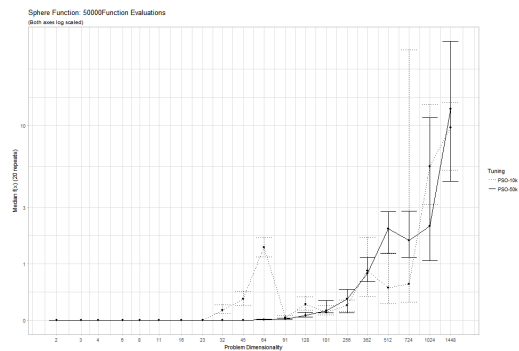
Schwefel Function at 10,000 Evaluations



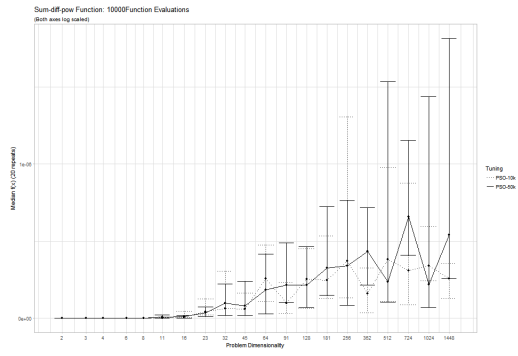
Schwefel Function at 50,000 Evaluations



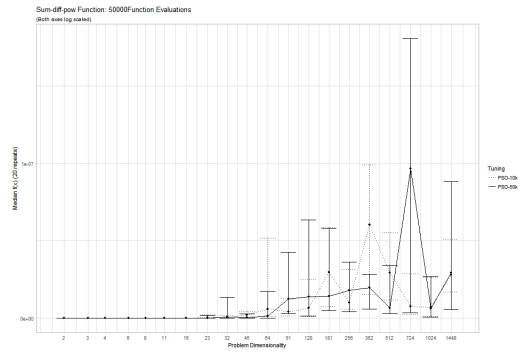
Sphere Function at 10,000 Evaluations



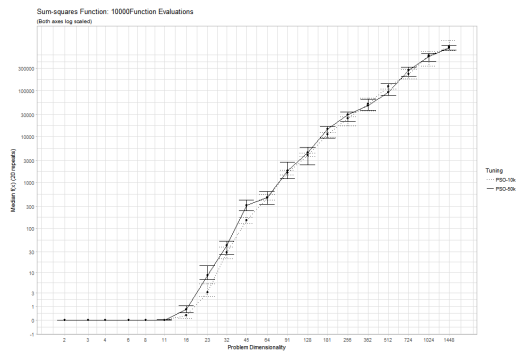
Sphere Function at 50,000 Evaluations



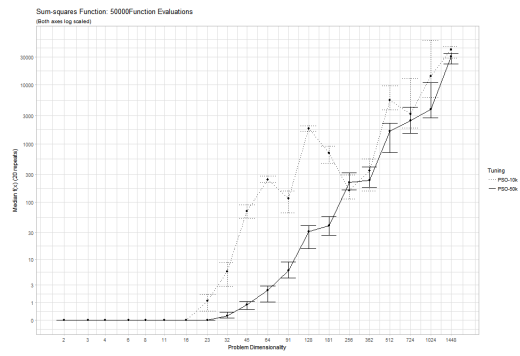
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

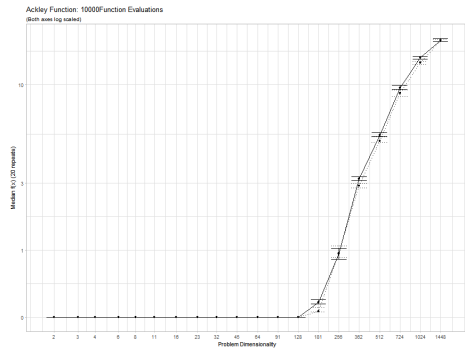


Sum Squares Function at 10,000 Evaluations

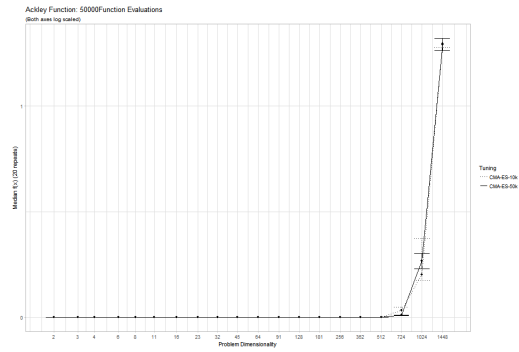


Sum Squares Function at 50,000 Evaluations

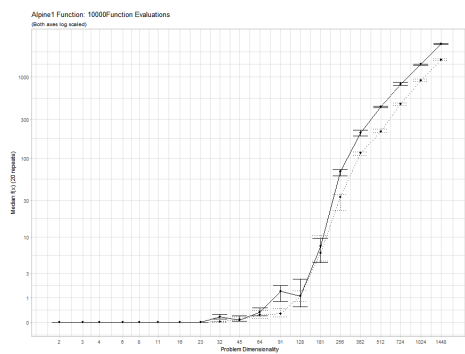
COVARIANCE MATRIX ADAPTION EVOLUTIONARY STRATEGY (CMA-ES)



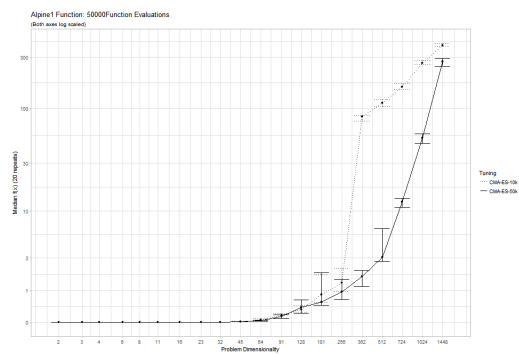
Ackley Function at 10,000 Evaluations



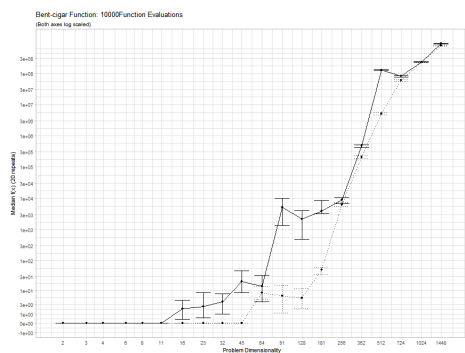
Ackley Function at 50,000 Evaluations



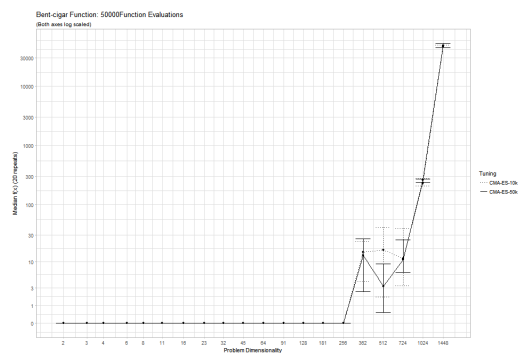
Alpine no.1 Function at 10,000 Evaluations



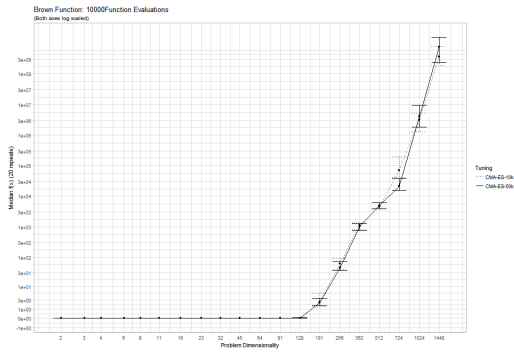
Alpine no.1 Function at 50,000 Evaluations



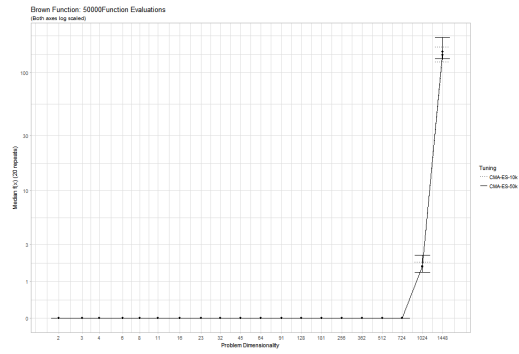
Bent Cigar Function at 10,000 Evaluations



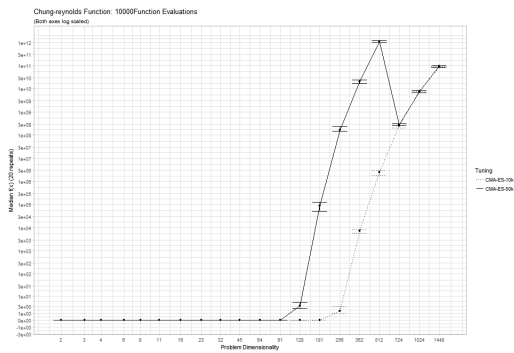
Bent Cigar Function at 50,000 Evaluations



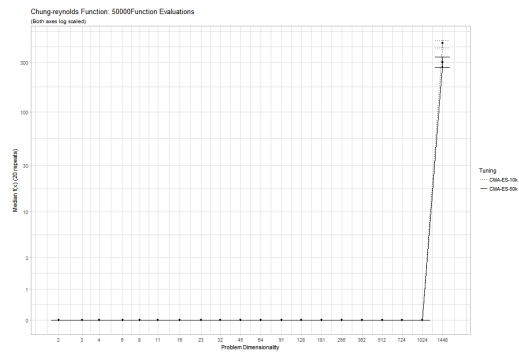
Brown Function at 10,000 Evaluations



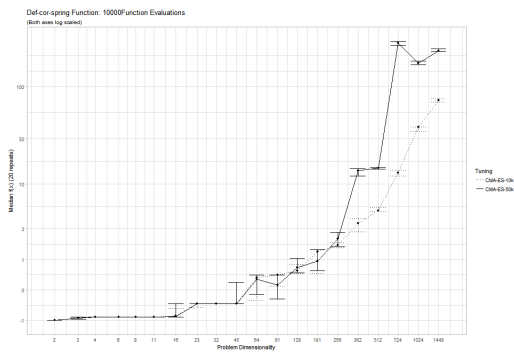
Brown Function at 50,000 Evaluations



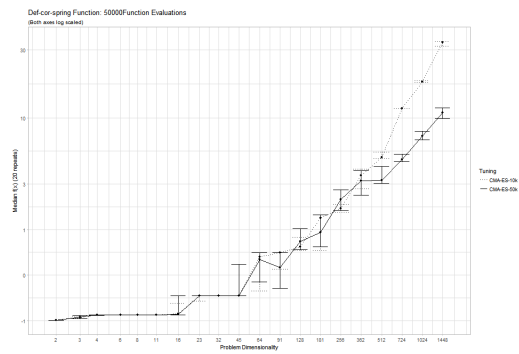
Chung-Reynolds Function at 10,000 Evaluations



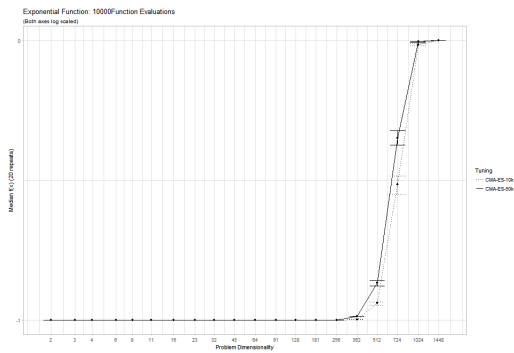
Chung-Reynolds Function at 50,000 Evaluations



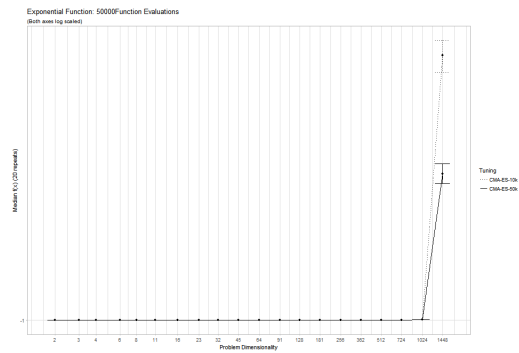
Deflected Corrugated Spring Function at 10,000 Evaluations



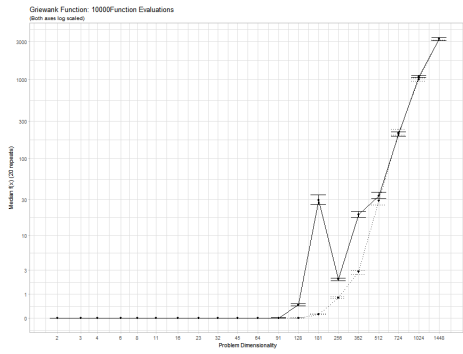
Deflected Corrugated Spring Function at 50,000 Evaluations



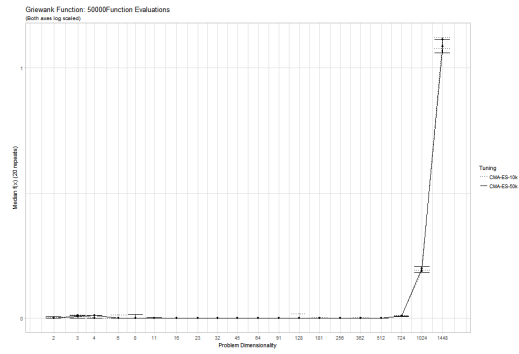
Exponential Function at 10,000 Evaluations



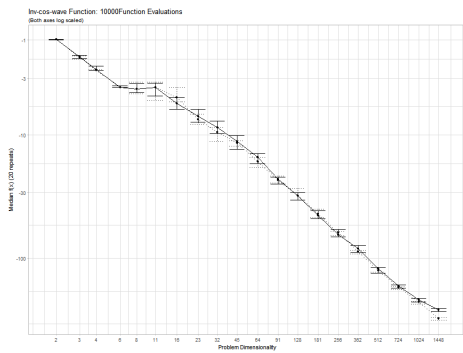
Exponential Function at 50,000 Evaluations



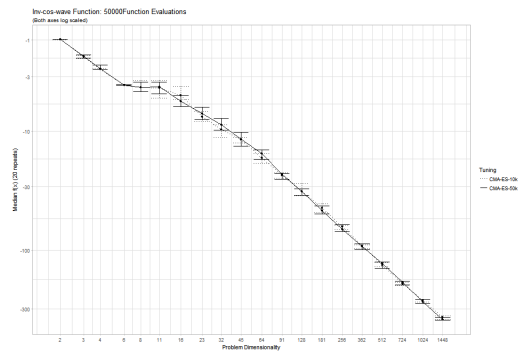
Griewank Function at 10,000 Evaluations



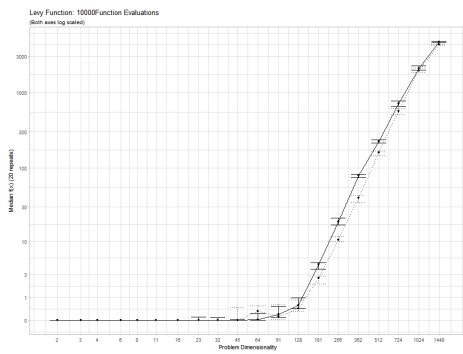
Griewank Function at 50,000 Evaluations



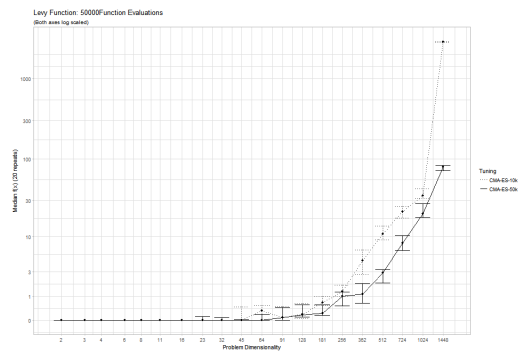
Inverted Cosine Wave Function at 10,000 Evaluations



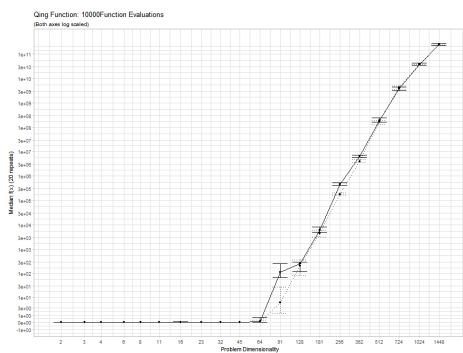
Inverted Cosine Wave Function at 50,000 Evaluations



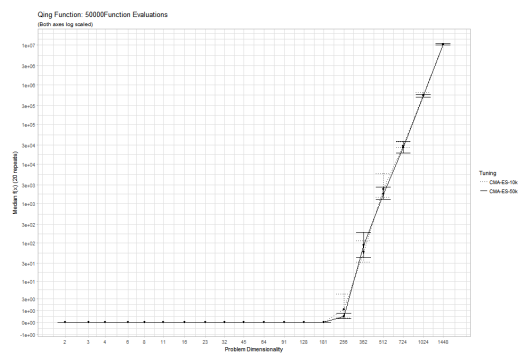
Levy Function at 10,000 Evaluations



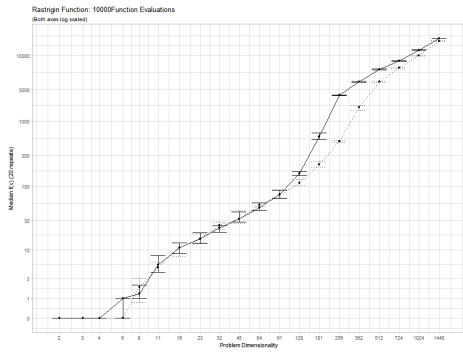
Levy Function at 50,000 Evaluations



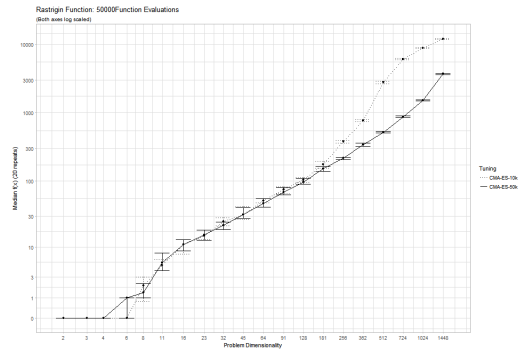
Qing Function at 10,000 Evaluations



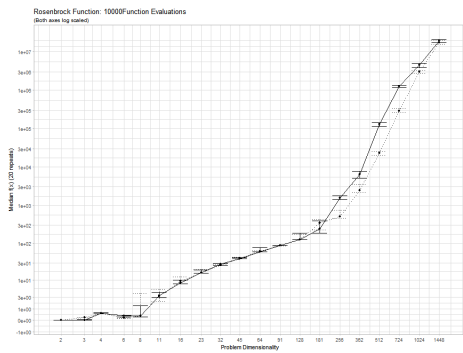
Qing Function at 50,000 Evaluations



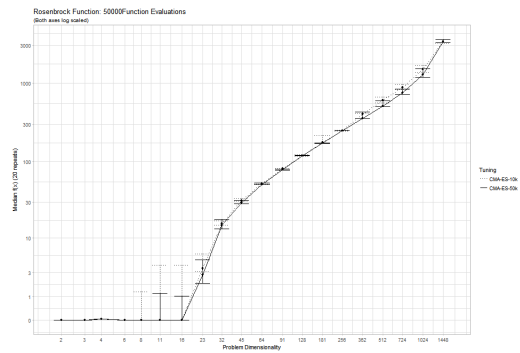
Rastrigin Function at 10,000 Evaluations



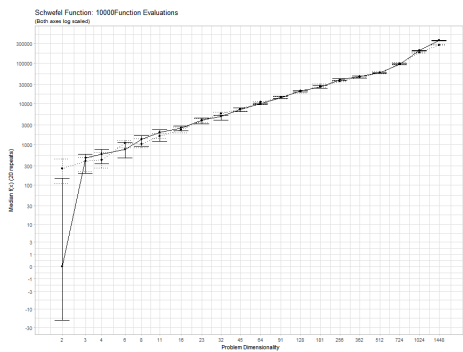
Rastrigin Function at 50,000 Evaluations



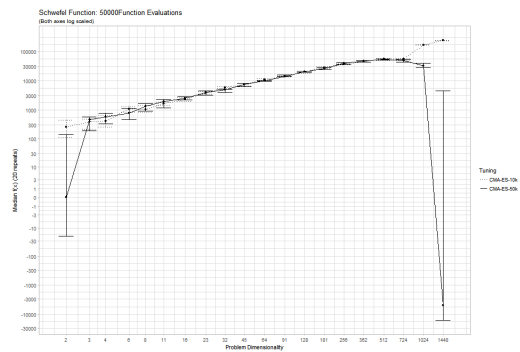
Rosenbrock Function at 10,000 Evaluations



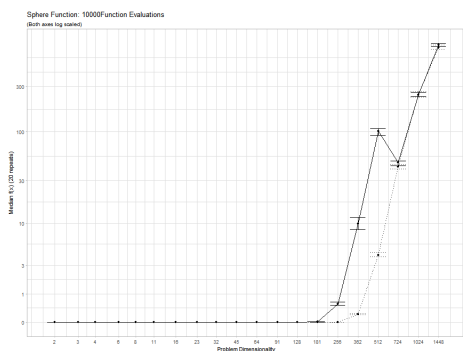
Rosenbrock Function at 50,000 Evaluations



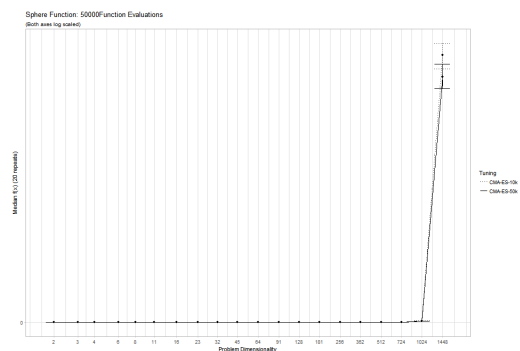
Schwefel Function at 10,000 Evaluations



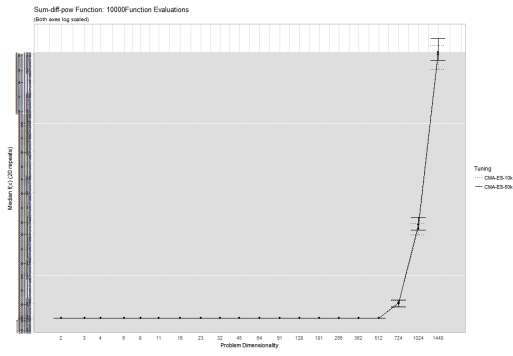
Schwefel Function at 50,000 Evaluations



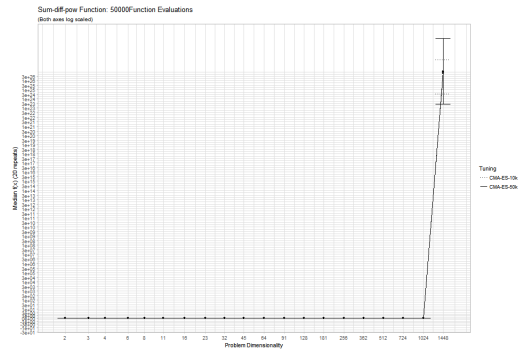
Sphere Function at 10,000 Evaluations



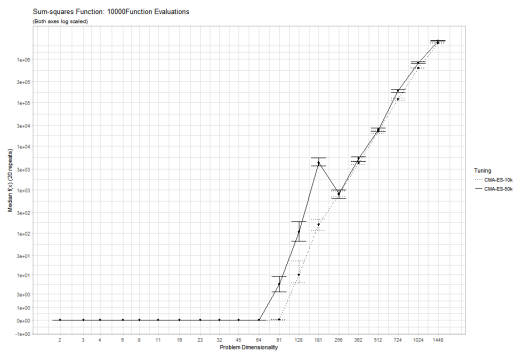
Sphere Function at 50,000 Evaluations



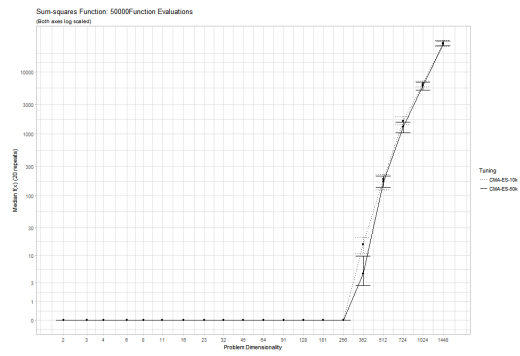
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations



Sum Squares Function at 10,000 Evaluations



Sum Squares Function at 50,000 Evaluations

APPENDIX B: 10K & 50K EVALUATION TUNING BUDGET
DESCRIPTIVE STATISTICS

Table B.15: GA Descriptive Statistics: Exponential to Rastrigin Functions (10,000 Evaluation Tuning)

Dims.	Exponential			Griewank			Inverted Cosine Wave			Levy			Qing			Rastrigin		
	no. evals			no. evals			no. evals			no. evals			no. evals			no. evals		
	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3
2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
3	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
4	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
6	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
8	-1	-1	-0.9999	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
11	-0.9999	-0.9999	-0.9999	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
16	-0.9996	-0.9987	-0.9994	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
23	-0.9983	-0.9987	-0.9979	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999
31	-0.9995	-0.9981	-0.9985	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996
46	-0.9981	-0.9985	-0.9981	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988
64	-0.9958	-0.9919	-0.9919	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968	-0.9968
91	-0.9847	-0.9928	-0.9923	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961
128	-0.9684	-0.9727	-0.9721	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868	-0.9868
184	-0.9629	-0.9721	-0.9721	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882	-0.9882
256	-0.9784	-0.9811	-0.9811	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929	-0.9929
364	-0.9517	-0.9526	-0.9526	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748	-0.9748
512	-0.9176	-0.9324	-0.9324	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682	-0.9682
744	-0.9091	-0.9124	-0.9124	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697	-0.9697
1024	-0.8156	-0.8156	-0.8156	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381	-0.9381
1408	-0.6853	-0.6853	-0.6853	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154	-0.8154

Table B.16: GA Descriptive Statistics: Exponential to Rastrigin Functions (50,000 Evaluation Tuning)

Dims.	Exponential			Griewank			Inverted Cosine Wave			Levy			Qing			Rastrigin		
	no. evals			no. evals			no. evals			no. evals			no. evals			no. evals		
	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3
2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
3	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
4	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
6	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
8	-0.9999	-1	-0.9999	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
11	-0.9996	-0.9987	-0.9993	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
16	-0.9984	-0.9986	-0.9986	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999
23	-0.9954	-0.9961	-0.9945	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988	-0.9988
31	-0.9987	-0.9988	-0.9984	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995	-0.9995
46	-0.9952	-0.9945	-0.9948	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987
64	-0.9888	-0.9873	-0.9879	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987
91	-0.9644	-0.9724	-0.9729	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987	-0.9987
128	-0.9328	-0.9438	-0.9438	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971	-0.9971
184	-0.9171	-0.9296	-0.9296	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961	-0.9961
256	-0.8717	-0.8798	-0.8798	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939	-0.9939
364	-0.8191	-0.8289	-0.8289	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927	-0.9927
512	-0.7631	-0.7689	-0.7689	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914	-0.9914
744	-0.6922	-0.7014	-0.7014	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897	-0.9897
1024	-0.6155	-0.6155	-0.6155	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876	-0.9876
1408	-0.5213	-0.5213	-0.5213	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812	-0.9812

Table B.17: GA Descriptive Statistics: Rosenbrock to Sum Squares Functions (10,000 Evaluation Tuning)

Dims.	Rosenbrock										Sphere										Sum of Different Powers										Sum Squares									
	10,000 Evaluations		50,000 Evaluations		100,000 Evaluations		Schwefel		50,000 Evaluations		100,000 Evaluations		10,000 Evaluations		50,000 Evaluations		100,000 Evaluations		10,000 Evaluations		50,000 Evaluations		100,000 Evaluations		10,000 Evaluations		50,000 Evaluations		100,000 Evaluations		10,000 Evaluations		50,000 Evaluations							
	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3							
2	0.02528	0.01924	0.12805	0.00704	0.0014	0.26099	0.00292	0.00837	0.00659	0.00774	0.28505	0.00053	3.69E-06	8.1E-07	6.53E-07	1.03E-07	1.75E-08	2.35E-07	4.31E-08	1.1E-08	1.59E-07	1.3E-09	3.06E-10	6.19E-09	1.34E-05	1.03E-05	2.09E-05	4.16E-07	9.3E-06	1.87E-06	2.82E-06	7.61E-07	6E-06							
3	0.72047	0.59294	1.19049	0.19077	0.02689	0.52749	0.00578	0.03687	0.02139	0.00489	0.00025	0.00183	6.83E-06	2.1E-06	1.03E-05	2.88E-07	1.2E-07	4.94E-07	1.03E-07	1.81E-08	1.93E-07	1.2E-09	4.54E-10	4.2E-09	9.24E-05	3.8E-05	0.000134	0.000266	9.5E-06	2.77E-06	1.1E-05									
4	1.63691	0.55146	2.39305	1.36192	0.52123	1.93391	0.03318	0.01595	0.00358	0.00732	0.00071	0.00334	2.66E-05	6.4E-06	5.07E-05	6.34E-07	3.4E-07	1.43E-06	1.62E-07	7.01E-08	4.4E-07	3.98E-09	5.94E-10	1.24E-08	0.000266	9.5E-06	0.000266	9.5E-06	5.8E-06	1.33E-05	1.33E-05									
6	3.37236	1.51759	4.71083	2.41043	1.02891	3.82719	0.19646	0.07473	0.02628	0.03032	0.00071	0.00334	6.00E-05	4.3E-05	9.16E-05	6.34E-06	3.1E-06	9.16E-06	1.41E-07	8.08E-08	4.4E-07	2.95E-08	2.11E-09	5.77E-08	0.001335	0.000902	0.001335	5.2E-05	3.6E-05	9.3E-05	9.3E-05									
8	8.48635	6.89698	15.21876	3.36448	1.26971	5.59599	0.40964	0.19136	0.07948	0.00823	0.00094	0.00341	0.00283	0.00094	0.00341	1.43E-05	9.42E-06	1.93E-05	3.38E-07	1.69E-07	3.94E-07	2.7E-08	3.3E-09	4.9E-08	0.005235	0.003499	0.005235	0.00171	8.4E-05	0.00021	0.00021									
11	62.32526	31.28997	74.071	6.73340	2.91898	9.80385	0.78951	0.36289	1.27643	0.04794	0.02266	0.05113	0.00779	0.00349	0.01092	2E-05	1.34E-05	4.12E-05	6.94E-07	1.1E-07	1.81E-06	3.5E-08	7.24E-09	7.72E-08	0.001874	0.011746	0.002983	0.000773	0.000494	0.001033	0.001033									
12	174.11967	97.01949	97.70042	4.888156	44.3368	23.9881	0.92726	1.44724	3.39707	0.0811	0.06232	0.21294	0.00381	0.0028	0.00428	0.00019	6.42E-05	0.00024	1.38E-06	2.88E-07	5.63E-06	5.1E-08	8.8E-09	1.08E-07	0.00585	0.00585	0.00585	0.00154	0.00154	0.00154	0.00154									
23	1.169447	1.051494	2.38928	68.22008	17.34984	80.81342	8.92761	6.85796	12.10966	0.33826	0.27932	0.93214	0.00731	0.00936	0.00777	0.00027	0.00026	0.00039	4.11E-06	4.4E-07	1.31E-05	1.7E-07	4.4E-08	3.87E-07	0.002885	0.014847	0.002885	0.001544	0.001544	0.001544	0.001544									
32	1.774997	1.551924	2.38928	86.27265	80.46126	12.83956	2.56997	19.91987	31.39683	0.85473	0.69345	1.34743	0.00839	0.01579	0.02841	0.00063	0.00053	0.000816	1.95E-06	2.16E-06	2.49E-05	2.01E-07	6.99E-08	4.86E-06	21.31085	18.09797	23.28736	0.81652	0.672371	0.95922	0.95922									
45	3.071316	3.137452	4.272211	1.473368	12.49446	16.99816	6.432316	51.79084	76.26882	2.947021	2.473912	3.666669	0.05494	0.051942	0.069793	0.00803	0.00118	0.002314	6.21E-06	2.64E-06	3.44E-05	1.89E-06	1.11E-07	3.42E-06	4.82E-04	4.489955	6.365666	0.144087	0.122747	0.191202	0.191202									
64	8.67755	7.293335	10.85001	51.47321	67.28888	62.61995	15.61433	143.45	165.3375	69.84972	57.79794	7.977394	1.40093	1.38022	1.66887	0.00339	0.04975	0.05925	4.93E-05	1.62E-05	0.000123	3.29E-05	7.41E-06	6.88E-05	56.7455	54.9795	57.8211	2.73797	21.40086	25.32681	25.32681									
91	1.97881	1.653803	2.298381	4.93809	3.759192	4.377244	6.001599	5.809822	7.437214	2.464011	2.021165	2.564497	0.62847	0.48719	0.700443	0.021791	0.07875	0.024827	3.96E-05	7.27E-07	0.00026	4.43E-07	5.63E-06	1.14E-05	3.714734	3.123212	4.860699	1.47987	12.0375	16.9931	16.9931									
128	4.475084	3.87997	5.147321	1.55549	1.108545	1.253713	3.956486	3.944005	4.113655	1.809955	1.713967	1.972746	4.269995	3.70881	4.68113	1.40247	1.34104	1.46266	0.00566	0.00047	0.000273	9.8E-06	2.64E-06	2.53E-05	1.46E-05	9.01139	8.184493	9.792434	3.97159	37.46787	47.7675	47.7675								
256	5.93487	4.01877	6.97636	2.236161	2.119873	2.554026	9.32177	8.90347	9.86269	5.57997	5.14149	6.209793	1.153482	1.115385	1.44159	0.44897	0.416682	0.524236	0.00471	4.8E-05	0.001123	3.06E-05	7.41E-06	6.88E-05	56.7455	54.9795	57.8211	2.73797	21.40086	25.32681	25.32681									
368	1.866284	1.678248	1.99977	4.66679	4.60223	5.29198	1.93375	1.816647	2.04961	1.665486	1.522826	1.73649	3.84797	3.58339	3.91631	0.82977	0.816682	0.945436	0.00471	4.8E-05	0.001123	3.06E-05	7.41E-06	6.88E-05	56.7455	54.9795	57.8211	2.73797	21.40086	25.32681	25.32681									
512	7.529539	6.85414	8.94381	1.168891	1.028274	1.201916	3.992666	3.756942	4.018416	4.026338	3.899438	4.270114	9.58486	9.33909	10.31439	3.897918	3.61895	4.20977	0.00259	0.00042	0.000273	9.8E-06	2.64E-06	2.53E-05	1.46E-05	9.01139	8.184493	9.792434	3.97159	37.46787	47.7675	47.7675								
724	3.088964	2.979269	3.159919	2.978157	2.754433	3.18679	4.79167	7.328355	7.91356	1.118401	1.118401	1.118401	4.22072	2.98756	2.55333	11.27116	10.99258	11.72117	0.00475	0.00056	0.000304	0.00013	6.51E-05	0.00039	3.97884	3.8643	3.7935	1.99657	15.9325	16.9931	16.9931									
1024	1.326623	1.245190	1.447943	1.04471	0.951288	1.111086	1.371269	1.340318	1.88974	2.65861	2.591963	2.76847	6.27102	6.07297	6.441198	3.149729	3.030458	3.314484	0.02026	0.00439	0.00211	0.00024	0.00039	0.000593	1.065913	1.04059	1.065913	5.71127	5.04226	5.90175	5.90175									
1448	4.824258	4.166949	4.921079	3.722858	3.491481	3.99782.8	2.409653	2.358462	2.417614	6.017679	5.966393	6.120498	4.64431	4.42491	4.90488	9.52725	9.294918	9.46888	0.00786	0.00576	0.0032	0.00021	8.53E-05	0.000647	3.95823	3.91896	3.95823	2.31264	2.12732	2.42541	2.42541									

Table B.18: GA Descriptive Statistics: Rosenbrock to Sum Squares Functions (50,000 Evaluation Tuning)

Dims.	Rosenbrock										Sphere										Sum of Different Powers										Sum Squares									
	10,000 Evaluations		50,000 Evaluations		100,000 Evaluations		Schwefel		50,000 Evaluations		100,000 Evaluations		10,000 Evaluations		50,000 Evaluations		100,000 Evaluations		10,000 Evaluations		50,000 Evaluations		100,000 Evaluations		10,000 Evaluations		50,000 Evaluations		100,000 Evaluations		10,000 Evaluations		50,000 Evaluations							
	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3							
2	0.017267	0.006635	0.059116	0.00293	0.00097	0.012168	0.003415	0.00481	0.005297	8.3E-05	4.8E-05	0.00092	1.1E-06	2.67E-07	3.91E-06	1.43E-07	5.67E-08	3.4E-07	4.7E-08	6.08E-09	1.2E-07	6.0E-10	2.69E-10	3.97E-09	1.56E-05	2.7E-06	3.39E-07	3.3E-07	1.81E-07	8.6E-07	8.6E-07									
3	0.537161	0.411496	0.914106	0.463273	0.227982	0.71719	0.010266	0.006055	0.016734	0.000379	0.000152	0.00071	9.47E-06	3.1E-06	1.31E-05	3.14E-07	2.07E-07	6.97E-07	1.69E-07	2.01E-08	2.82E-07	4.3E-09	7.93E-10	1.58E-08	4.32E-05	2.98E-05	0.00021	3.54E-06	2.13E-06	4.17E-06	4.17E-06									
4	0.83147	0.534605	1.33995	0.46793	0.14901	0.990386	0.02747	0.020996	0.07169	0.001776	0.00026	0.003433	3.54E-05	2.37E-05	7.23E-05	1.44E-06	6.69E-07	2.05E-06	2.48E-07	5.11E-08	8.84E-07	7.43E-09	1.23E-09	2.77E-08	0.00154	0.00013	0.00029	6.75E-06	2.77E-06	1.1E-05	1.1E-05									
6	4.26393	3.3332	5.954793	2.352184	1.66653	3.095987	0.110683	0.081042	0.240887	0.009931	0.0021	0.007393	0.00132	7.01E-05	0.000188	3.72E-06	1.8E-06	6.11E-06	1.88E-07	8.79E-08	4.14E-07	5.5E-09	9.93E-10	9.93E-09	0.001353	0.000339	0.002369	5.8E-05	3.98E-05	9.5E-05	9.5E-05									
8	6.842301	5.49856	11.22564	3.182453	1.782366	5.8391	0.369568	0.171969	0.201353	0.01157	0.00756	0.020356	0.002023	0.00069	0.000348	1.23E-05	8.32E-06	1.6E-05	1.32E-06	1.8E-07	3.64E-06	2.62E-08	5.59E-09	8.54E-08	0.0047	0.000937	0.008185	0.000185	0.000185	0.000185	0.000185									
11	2.597931	1.10088	6.8583	6.58131	2.173116	8.54278	8.08423	5.99811	1.045966	0.027601	0.03149	0.049776	0.001968	0.00174	0.00087	9.53E-05	7.99E-05	0.00017	8.2E-07	2.34E-07	1.51E-06	2.53E-08	7.8E-09	7.65E-08	0.02173	0.013319	0.01331	0.00056	0.00049	0.00049	0.00049									
23	1.324265	1.118662	1.799186	7.296168	19.27156	8.88839	8.62278	6.26956	10.4975	0.355448	0.33033	0.494161	0.00771	0.00494	0.00988	0.00029	0.00039	0.00037	1.46E-06	6.98E-07	1E-05	1.95E-07	1.47E-08	5.5E-07	0.45696	0.241335	0.61545	0.01376	0.00072	0.00939	0.00939									
32	1.90673	1.883138	2.359797	8.87967	78.4404	11.53158	29.4569	20.99724	36.85242	0.9739	0.88144	1.172481	0.025964	0.021225	0.02996	0.000825	0.00076	0.00109	1.69E-06	1.93E-06	1.6E-05	5.6E-08	2.14E-08	2.03E-07	1.333654	1.0431	1.594143	0.04318	0.03904	0.06647	0.06647									
45	3.96467	3.04701	4.28195	1.489692																																				

Table B.25: RMHC Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (10,000 Evaluation Tuning)

Dims.	Ackley										Brent-Cigar										Brown										Chung-Keays										Deflected Corrugated Spring									
	5000 Evaluations					10000 Evaluations					15000 Evaluations					20000 Evaluations					25000 Evaluations					30000 Evaluations					35000 Evaluations					40000 Evaluations					45000 Evaluations					50000 Evaluations				
	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3								
2	0.00074	0.00037	0.00017	0.00043	8.54e-05	0.00043	2.4e-05	1.88e-05	3.18e-05	5.92e-06	3.92e-06	1.44e-05	1.51e-06	8.83e-06	3.28e-06	2.42e-07	3.35e-07	3.35e-07	4.42e-12	1.12e-12	1.42e-11	1.97e-14	2.49e-14	3.23e-13	6.64e-16	9.92e-16	6.42e-17	2.92e-20	1.28e-21	8.61e-21	-0.8334	-0.9997	-0.8334	-0.8334	-1	-0.8334														
3	0.00147	0.00075	0.00038	0.00028	0.00028	0.00028	8.88e-05	6.82e-05	0.00031	1.19e-05	5.01e-06	2.72e-05	2.18e-05	8.24e-06	0.00048	7.11e-07	3.05e-07	2.31e-06	4.27e-11	1.37e-11	4.49e-10	3.97e-11	2.17e-11	8.06e-13	6.64e-16	2.92e-16	2.4e-15	8.08e-20	2.87e-20	4.13e-19	-0.8334	-0.8334	-0.8334	-0.8334	-0.8334	-0.8334														
4	0.00231	0.00121	0.00059	0.00047	0.00047	0.00047	0.00013	5.95e-05	0.00036	3.18e-05	2.26e-05	4.57e-05	0.00017	7.15e-05	0.00083	2.33e-06	1.27e-06	3.27e-06	5.87e-09	1.67e-09	1.07e-08	2.72e-10	1.92e-10	2.62e-10	2.72e-14	1.92e-15	1.35e-16	1.62e-19	1.34e-19	1.34e-19	-0.6853	-0.8334	-0.7336	-0.6853	-0.8334	-0.7336														
6	0.00337	0.00231	0.00147	0.00098	0.00098	0.00098	0.00032	0.00022	0.00042	5.79e-05	4.33e-05	9.22e-05	0.00019	0.00022	0.00088	3.28e-05	1.94e-05	6.44e-05	2.88e-09	2.88e-09	1.44e-09	3.92e-11	1.44e-11	4.43e-11	1.35e-15	2.81e-16	1.73e-16	4.81e-16	1.47e-16	1.73e-16	-0.7336	-0.7336	0.49922	-0.7336	-0.7336	0.49922														
8	0.00436	0.00297	0.00179	0.00116	0.00116	0.00116	0.00041	0.00028	0.00048	8.92e-05	6.27e-05	0.00013	0.00033	0.00022	0.00089	0.00058	7.40e-05	0.00013	1.94e-09	2.42e-09	4.97e-09	2.81e-10	1.44e-10	1.44e-10	4.35e-14	1.35e-14	1.24e-14	1.35e-14	1.24e-14	1.35e-14	-0.49922	-0.7336	0.49922	-0.49922	-0.7336	0.49922														
11	0.00584	0.00396	0.00249	0.00159	0.00159	0.00159	0.00054	0.00036	0.00054	0.00021	0.00014	0.00025	0.00043	0.00029	0.00085	0.00068	0.00019	0.00029	6.97e-09	8.24e-09	1.97e-08	9.27e-09	6.02e-09	1.29e-08	3.68e-08	1.92e-08	2.4e-08	1.81e-11	8.82e-12	9.27e-12	1.90499	0.49922	1.065	1.90499	0.49922	1.065														
16	0.00749	0.00492	0.00308	0.00199	0.00199	0.00199	0.00064	0.00041	0.00064	0.00024	0.00014	0.00024	0.00043	0.00029	0.00085	0.00068	0.00019	0.00029	6.97e-09	8.24e-09	1.97e-08	9.27e-09	6.02e-09	1.29e-08	3.68e-08	1.92e-08	2.4e-08	1.81e-11	8.82e-12	9.27e-12	1.90499	0.49922	1.065	1.90499	0.49922	1.065														
23	0.01353	0.00908	0.00561	0.00356	0.00356	0.00356	0.00107	0.00067	0.00107	0.00042	0.00024	0.00042	0.00064	0.00041	0.00089	0.00073	0.00023	0.00041	1.07e-08	1.29e-08	3.23e-08	1.51e-08	7.92e-09	1.44e-08	4.35e-08	2.42e-08	3.12e-06	1.44e-08	8.42e-08	6.67277	4.26859	7.29396	6.67277	4.26859	7.29396															
32	0.02129	0.01409	0.00864	0.00529	0.00529	0.00529	0.00159	0.00107	0.00159	0.00073	0.00042	0.00073	0.00107	0.00067	0.00138	0.00107	0.00038	0.00067	1.51e-08	1.81e-08	4.81e-08	2.27e-08	1.24e-07	2.42e-07	3.92e-07	3.92e-07	3.92e-07	3.92e-07	3.92e-07	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
45	0.03208	0.02129	0.01353	0.00864	0.00864	0.00864	0.00249	0.00166	0.00249	0.00107	0.00067	0.00107	0.00159	0.00107	0.00166	0.00138	0.00048	0.00089	2.27e-08	2.81e-08	7.29e-08	3.23e-08	1.66e-07	3.23e-07	5.15e-07	5.15e-07	5.15e-07	5.15e-07	5.15e-07	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
64	0.05029	0.03208	0.02129	0.01353	0.01353	0.01353	0.00396	0.00261	0.00396	0.00200	0.00121	0.00200	0.00281	0.00200	0.00320	0.00261	0.00089	0.00138	3.23e-08	4.07e-08	1.07e-07	4.81e-08	2.42e-07	4.81e-07	7.29e-07	7.29e-07	7.29e-07	7.29e-07	7.29e-07	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
91	0.07492	0.05029	0.03208	0.02129	0.02129	0.02129	0.00584	0.00396	0.00584	0.00308	0.00199	0.00308	0.00436	0.00308	0.00436	0.00356	0.00138	0.00200	4.81e-08	5.92e-08	1.51e-07	6.67e-08	3.23e-07	6.67e-07	1.07e-06	1.07e-06	1.07e-06	1.07e-06	1.07e-06	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
128	0.10749	0.07492	0.05029	0.03208	0.03208	0.03208	0.00864	0.00584	0.00864	0.00492	0.00308	0.00492	0.00684	0.00492	0.00684	0.00492	0.00166	0.00249	6.67e-08	8.24e-08	2.12e-07	8.24e-08	4.07e-07	8.24e-07	1.29e-06	1.29e-06	1.29e-06	1.29e-06	1.29e-06	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
181	0.15129	0.10749	0.07492	0.05029	0.05029	0.05029	0.01249	0.00864	0.01249	0.00684	0.00436	0.00684	0.00908	0.00684	0.00908	0.00684	0.00249	0.00396	9.27e-08	1.16e-07	3.08e-07	1.16e-07	5.15e-07	1.16e-06	1.81e-06	1.81e-06	1.81e-06	1.81e-06	1.81e-06	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
249	0.20749	0.15129	0.10749	0.07492	0.07492	0.07492	0.01812	0.01249	0.01812	0.01074	0.00684	0.01074	0.01409	0.01074	0.01409	0.01074	0.00396	0.00584	1.29e-07	1.51e-07	4.07e-07	1.51e-07	7.29e-07	1.51e-06	2.27e-06	2.27e-06	2.27e-06	2.27e-06	2.27e-06	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
368	0.28129	0.20749	0.15129	0.10749	0.10749	0.10749	0.0261	0.01812	0.0261	0.0166	0.01074	0.0166	0.02129	0.0166	0.02129	0.0166	0.00584	0.00864	1.81e-07	2.12e-07	5.92e-07	2.12e-07	1.07e-06	2.12e-06	3.23e-06	3.23e-06	3.23e-06	3.23e-06	3.23e-06	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
512	0.37492	0.28129	0.20749	0.15129	0.15129	0.15129	0.0396	0.0261	0.0396	0.0249	0.0166	0.0249	0.03208	0.0249	0.03208	0.0249	0.00864	0.01249	2.42e-07	2.81e-07	7.29e-07	2.81e-07	1.407e-06	2.81e-06	4.07e-06	4.07e-06	4.07e-06	4.07e-06	4.07e-06	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
724	0.49279	0.37492	0.28129	0.20749	0.20749	0.20749	0.0584	0.0396	0.0584	0.0356	0.0249	0.0356	0.0436	0.0356	0.0436	0.0356	0.01249	0.01812	3.23e-07	3.67e-07	1.07e-06	3.67e-07	1.81e-06	3.67e-06	5.15e-06	5.15e-06	5.15e-06	5.15e-06	5.15e-06	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
1024	0.64279	0.49279	0.37492	0.28129	0.28129	0.28129	0.0864	0.0584	0.0864	0.0492	0.03208	0.0492	0.0612	0.0492	0.0612	0.0492	0.01812	0.0261	4.81e-07	5.52e-07	1.51e-06	5.52e-07	2.81e-06	5.52e-06	7.29e-06	7.29e-06	7.29e-06	7.29e-06	7.29e-06	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															
1448	0.84279	0.64279	0.49279	0.37492	0.37492	0.37492	0.1249	0.0864	0.1249	0.0684	0.0436	0.0684	0.0864	0.0684	0.0864	0.0684	0.0261	0.0396	7.29e-07	8.24e-07	2.12e-06	8.24e-07	4.07e-06	8.24e-06	1.07e-05	1.07e-05	1.07e-05	1.07e-05	1.07e-05	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08	1.81e-08															

Table B.26: RMHC Descriptive Statistics: Ackley to Deflected Corrugated Spring Functions (50,000 Evaluation Tuning)

Dims.	Ackley										Brent-Cigar										Brown										Chung-Keays										Deflected Corrugated Spring									
	5000 Evaluations					10000 Evaluations					15000 Evaluations					20000 Evaluations					25000 Evaluations					30000 Evaluations					35000 Evaluations					40000 Evaluations					45000 Evaluations					50000 Evaluations				
	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3	Min	Q1	Q3	Median	Q1	Q3								
2	0.0013	0.00078	0.00043	0.00061	0.00043	0.00043	1.92e-05	1.07e-05	4.13e-05	4.37e-06	2.72e-06	8.34e-06	2.71e-06	1.44e-06	8.88e-06	1.51e-07	8.88e-08	3.55e-07	3.08e-12	1.81e-12	7.92e-12	3.23e-14	4.49e-14	1.54e-13	2.27e-17	9.81e-18	2.27e-17	1.47e-19	2.87e-20	3.55e-19	-0.8334	-0.9998	-0.8334	-0.8334	-1	-0.8334														
3	0.00167	0.00087	0.00047	0.00032	0.00032	0.00032	8.24e-05	4.49e-05	0.00017	1.71e-05	8.41e-06	2.31e-05	2.72e-05	1.27e-05	0.00013	3.77e-07	2.81e-07	2.49e-06	4.27e-11	1.37e-11	4.49e-10	3.97e-11	2.17e-11	8.06e-13	6.64e-16	2.92e-16	2.4e-15	8.08e-20	2.87e-20	4.13e-19	-0.8334	-0.8334	-0.7284	-0.8334	-0.8334	-0.7284														
4	0.00241	0.00121	0.00059	0.00037	0.00037	0.00037	0.00013	5.95e-05	0.00036	3.18e-05	2.26e-05	4.57e-05	0.00017	7.15e-05	0.00083	2.33e-06	1.27e-06	3.27e-06	5.87e-09	1.67e-09	1.07e-08	2.72e-10	1.92e-10	2.62e-10	2.72e-14	1.92e-15	1.35e-16	1.62e-19	1.34e-19	1.34e-19	-0.6853	-0.8334	-0.7336	-0.6853	-0.8334	-0.7336														
6	0.00337	0.00231	0.00147	0.00098	0.00098	0.00098	0.00032	0.00022	0.00042	5.79e-05	4.33e-05	9.22e-05	0.00019	0.00022	0.00088	3.28e-05	1.94e-05	6.44e-05	2.88e-09	2.88e-09	1.44e-09	3.92e-11	1.44e-11	4.43e-11	1.35e-15	2.81e-16	1.73e-16	4.81e-16	1.47e-16	1.73e-16	-0.7336	-0.7336	0.49922	-0.7336	-0.7336	0.49922														
8	0.00436	0.00297	0.00179	0.00116	0.00116	0.00116	0.00041	0.00028	0.00048	8.92e-05	6.27e-05	0.00013	0.00033	0.00022																																				

Table B.27: RMHC Descriptive Statistics: Exponential to Rastrigin Functions (10,000 Evaluation Tuning)

Dims.	Exponential			Circus			Inverted Circus Wave			Levy			Qing			Rastrigin		
	no. evals	Q1	Q3	no. evals	Q1	Q3	no. evals	Q1	Q3	no. evals	Q1	Q3	no. evals	Q1	Q3	no. evals	Q1	Q3
	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3
2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
3	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
4	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
6	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
8	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
11	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
16	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
23	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
34	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
46	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
64	-0.9999	-0.9999	-0.9998	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
91	-0.9993	-0.9994	-0.9992	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
128	-0.996	-0.9948	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999	-0.9999
184	-0.99748	-0.99774	-0.9971	-0.99985	-0.99985	-0.99984	-1.00000	-1.00000	-1.00000	-1.00000	-1.00000	-1.00000	-1.00000	-1.00000	-1.00000	-1.00000	-1.00000	-1.00000
256	-0.9842	-0.9772	-0.9855	-0.9968	-0.9968	-0.9964	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996	-0.9996
364	-0.9329	-0.9267	-0.9318	-0.9886	-0.9885	-0.9881	-0.9951	-0.9951	-0.9951	-0.9951	-0.9951	-0.9951	-0.9951	-0.9951	-0.9951	-0.9951	-0.9951	-0.9951
512	-0.9423	-0.9248	-0.9408	-0.9841	-0.9841	-0.9842	-0.9902	-0.9902	-0.9902	-0.9902	-0.9902	-0.9902	-0.9902	-0.9902	-0.9902	-0.9902	-0.9902	-0.9902
744	-0.929	-0.94436	-0.9595	-0.9535	-0.9535	-0.9536	-0.9586	-0.9586	-0.9586	-0.9586	-0.9586	-0.9586	-0.9586	-0.9586	-0.9586	-0.9586	-0.9586	-0.9586
1044	-0.9045	-0.9053	-0.9024	-0.9459	-0.9459	-0.9456	-0.9486	-0.9486	-0.9486	-0.9486	-0.9486	-0.9486	-0.9486	-0.9486	-0.9486	-0.9486	-0.9486	-0.9486
1448	-0.8507	-0.8507	-0.8507	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251	-0.94251

Table B.28: RMHC Descriptive Statistics: Exponential to Rastrigin Functions (50,000 Evaluation Tuning)

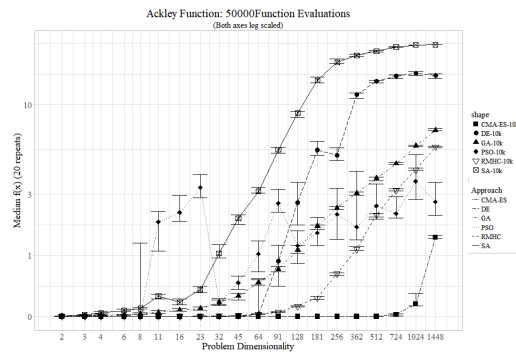
Dims.	Exponential			Circus			Inverted Circus Wave			Levy			Qing			Rastrigin		
	no. evals	Q1	Q3	no. evals	Q1	Q3	no. evals	Q1	Q3	no. evals	Q1	Q3	no. evals	Q1	Q3	no. evals	Q1	Q3
	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3	Median	Q1	Q3
2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
3	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
4	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
6	-0.9994	-1	-0.9949	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
11	-0.9562	-0.9245	-0.9575	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
16	-0.92784	-0.92922	-0.9272	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
23	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
34	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
46	-0.9436	-0.9463	-0.9466	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
64	-0.924	-0.9247	-0.923	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
91	-0.9781	-0.9375	-0.9224	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
128	-0.9413	-0.9268	-0.9319	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
184	-0.9001	-0.9002	-0.9106	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
256	-0.9106	-0.9106	-0.9106	-0.9998	-0.9998	-0.9998	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000
364	-0.9106	-0.9106	-0.9106	-0.9998	-0.9998	-0.9998	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000
512	-0.9106	-0.9106	-0.9106	-0.9998	-0.9998	-0.9998	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000
744	-0.9106	-0.9106	-0.9106	-0.9998	-0.9998	-0.9998	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000
1044	-0.9106	-0.9106	-0.9106	-0.9998	-0.9998	-0.9998	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000
1448	-0.9106	-0.9106	-0.9106	-0.9998	-0.9998	-0.9998	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000	-1.0000



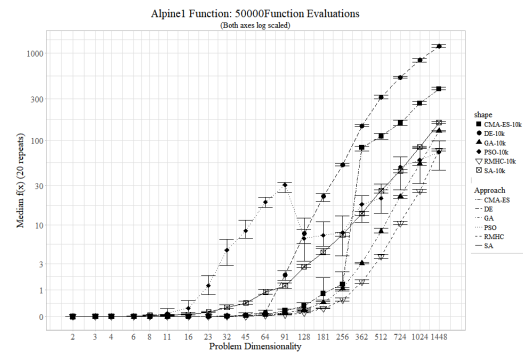
APPENDIX C: SMALL TO LARGE-SCALE METAHEURISTIC COMPARISON MATERIAL

C.1 METAHEURISTIC COMPARISON PLOTS

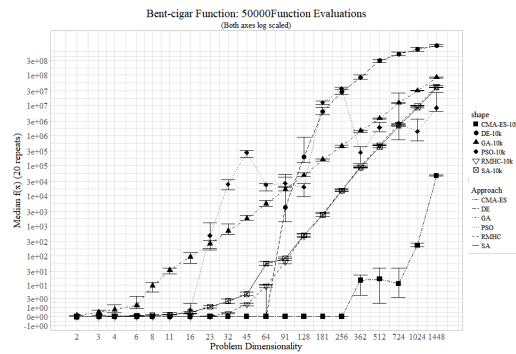
C.1.1 Comparison of Approaches Tuned at the 10,000 Evaluation Budget (With Q_1 & Q_3 Error Bars)



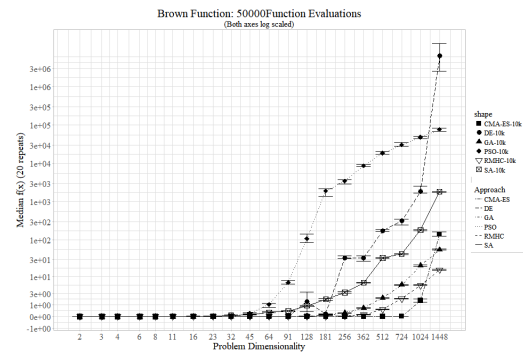
Ackley Function at 50,000 Evaluations



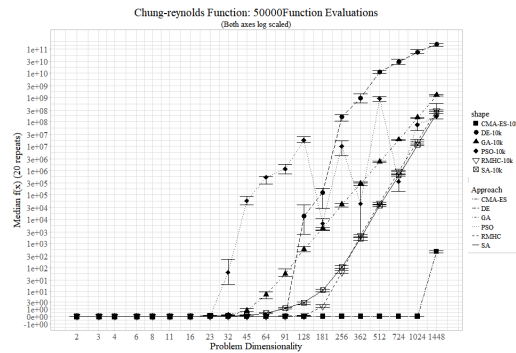
Alpine n.1 Function at 50,000 Evaluations



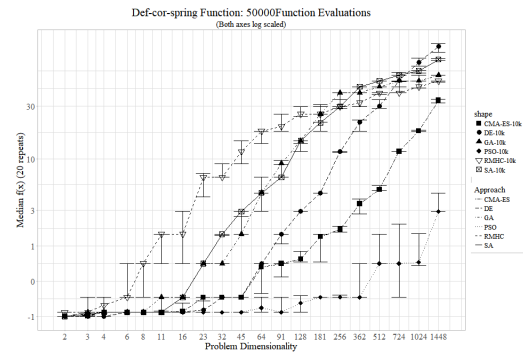
Bent Cigar Function at 50,000 Evaluations



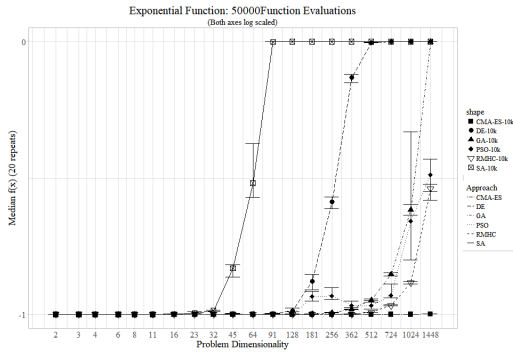
Brown Function at 50,000 Evaluations



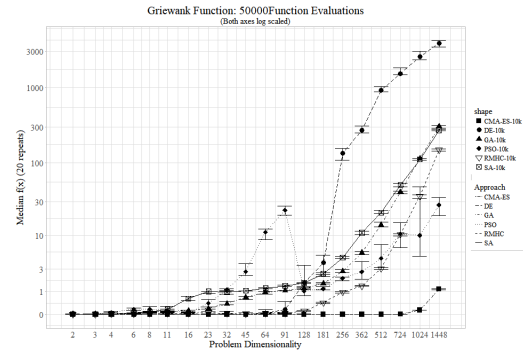
Chung-Reynolds Function at 50,000 Evaluations



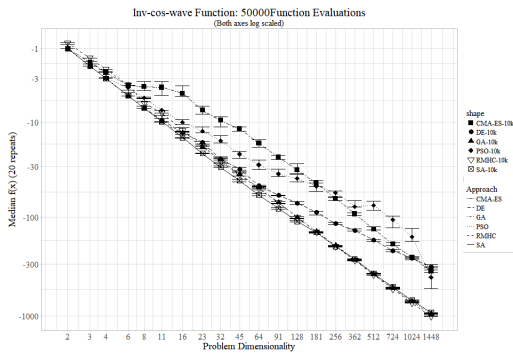
Deflected Corrugated Spring Function at 50,000 Evaluations



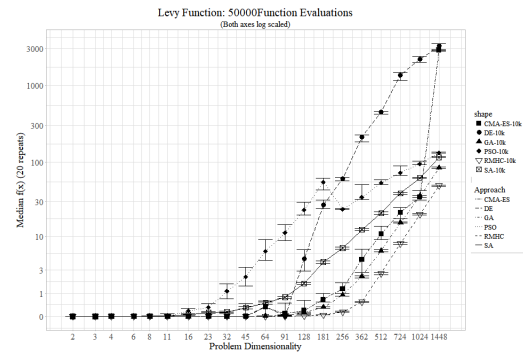
Exponential Function at 50,000 Evaluations



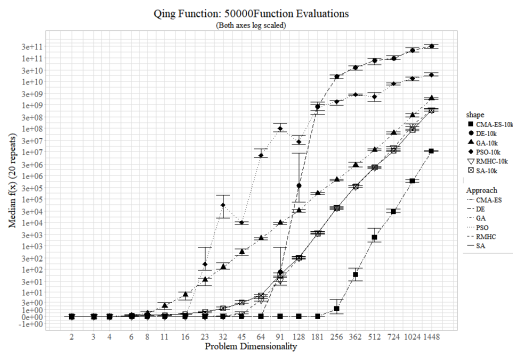
Griewank Function at 50,000 Evaluations



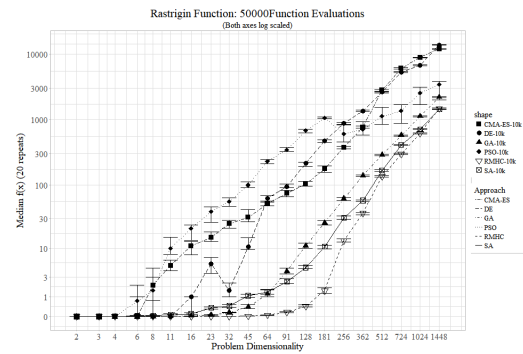
Inverted Cosine Wave Function at 50,000 Evaluations



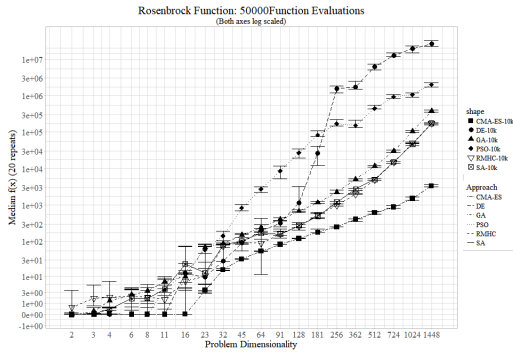
Levy Function at 50,000 Evaluations



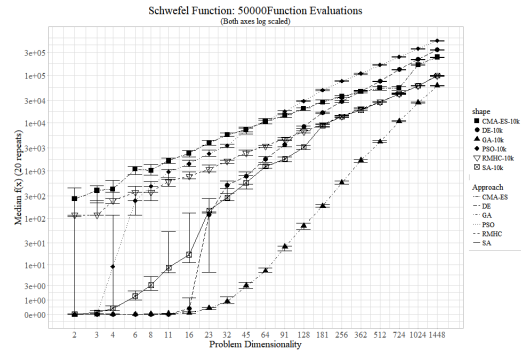
Qing Function at 50,000 Evaluations



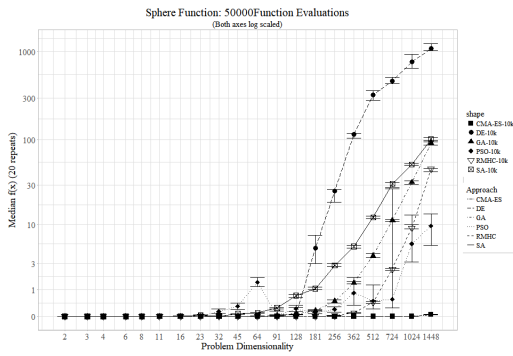
Rastrigin Function at 50,000 Evaluations



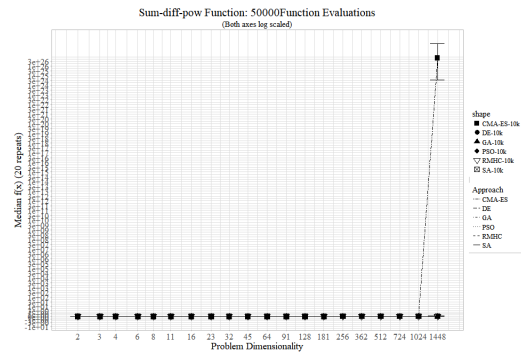
Rosenbrock Function at 50,000 Evaluations



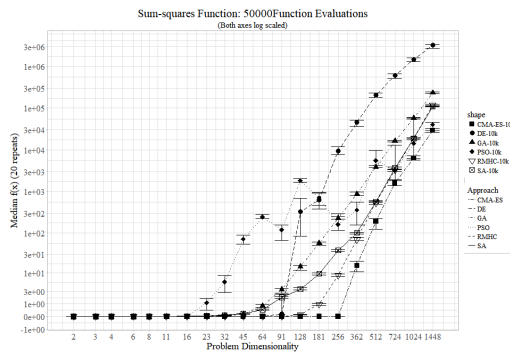
Schwefel Function at 50,000 Evaluations



Sphere Function at 50,000 Evaluations

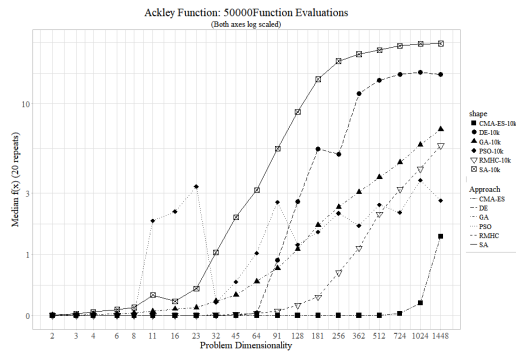


Sum of Different Powers Function at 50,000 Evaluations

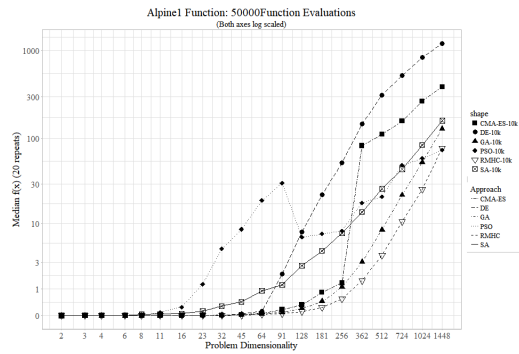


Sum Squares Function at 50,000 Evaluations

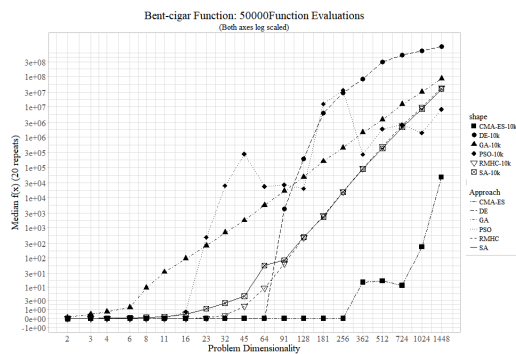
C.1.2 Comparison of Approaches Tuned at the 10,000 Evaluation Budget (No Error Bars)



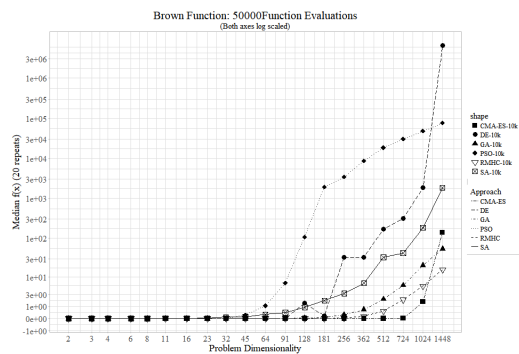
Ackley Function at 50,000 Evaluations



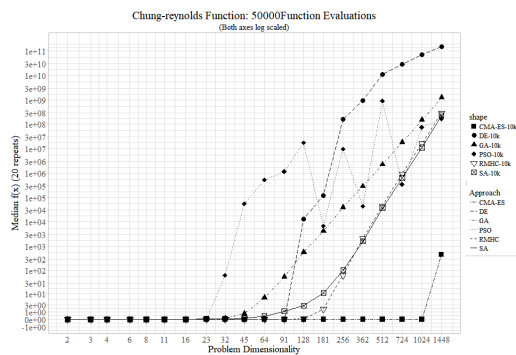
Alpine n.1 Function at 50,000 Evaluations



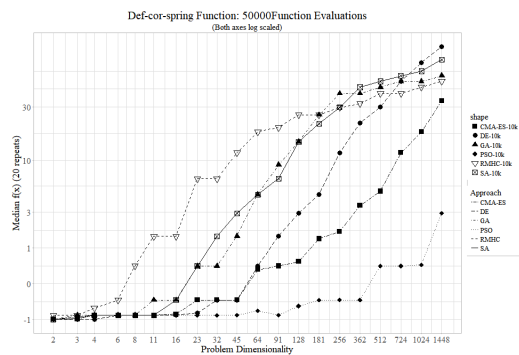
Bent Cigar Function at 50,000 Evaluations



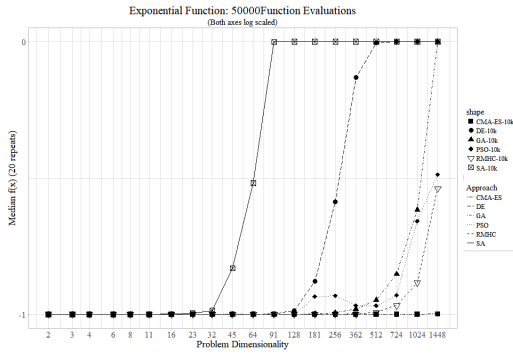
Brown Function at 50,000 Evaluations



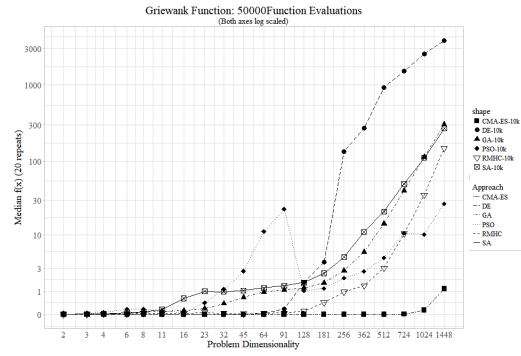
Chung-Reynolds Function at 50,000 Evaluations



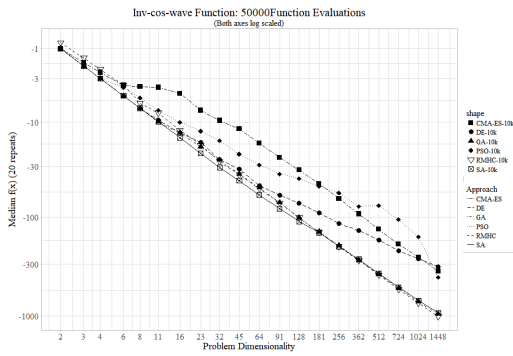
Deflected Corrugated Spring Function at 50,000 Evaluations



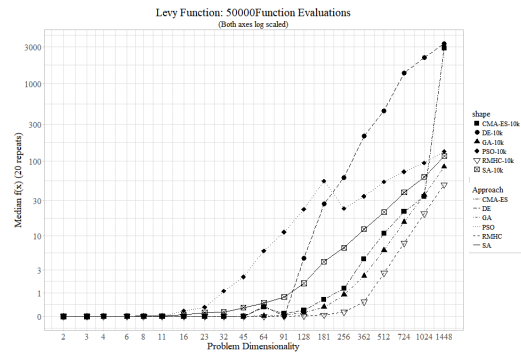
Exponential Function at 50,000 Evaluations



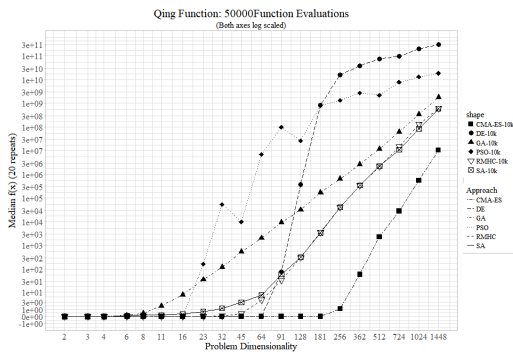
Griewank Function at 50,000 Evaluations



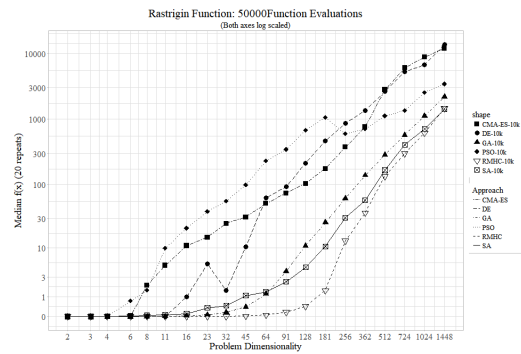
Inverted Cosine Wave Function at 50,000 Evaluations



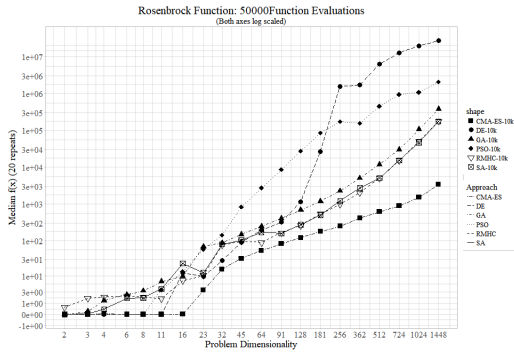
Levy Function at 50,000 Evaluations



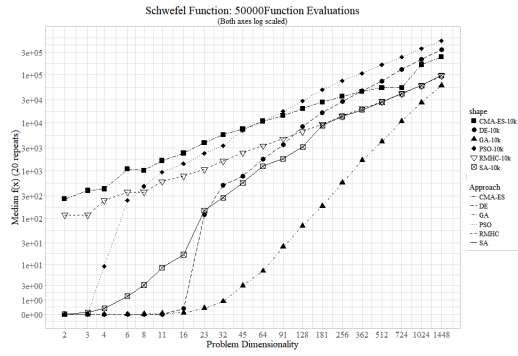
Qing Function at 50,000 Evaluations



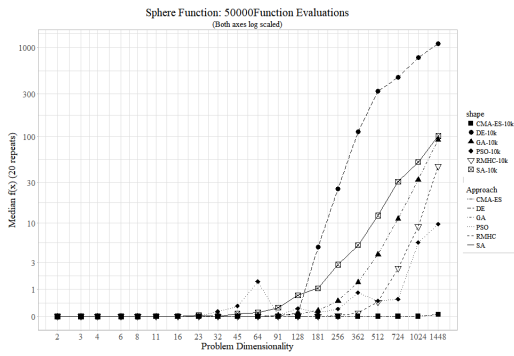
Rastrigin Function at 50,000 Evaluations



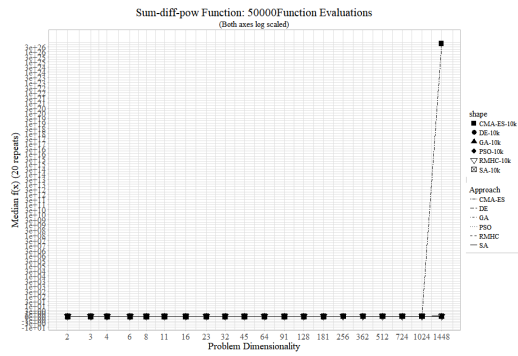
Rosenbrock Function at 50,000 Evaluations



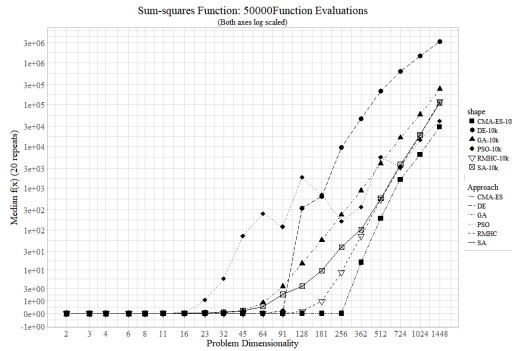
Schwefel Function at 50,000 Evaluations



Sphere Function at 50,000 Evaluations

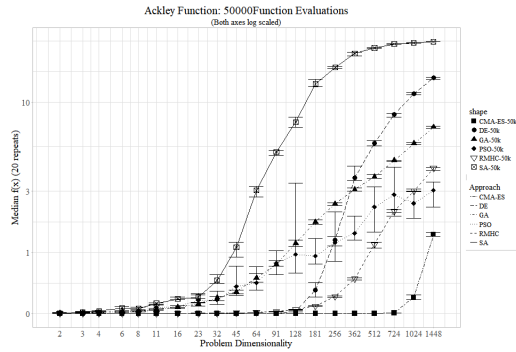


Sum of Different Powers Function at 50,000 Evaluations

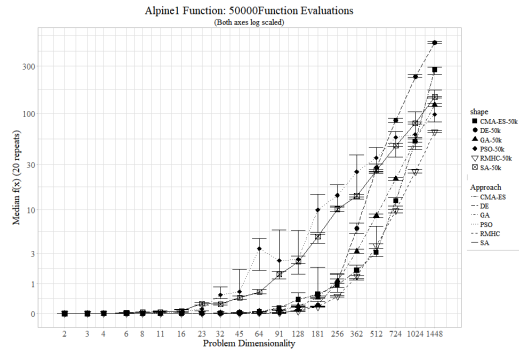


Sum Squares Function at 50,000 Evaluations

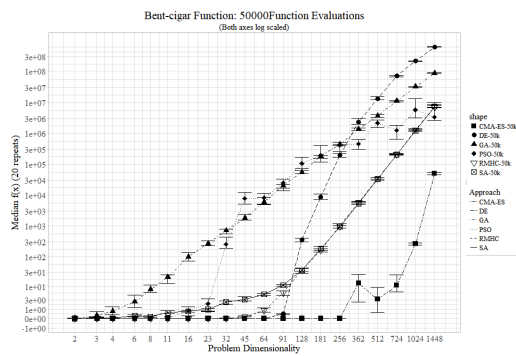
C.1.3 Comparison of Approaches Tuned at the 50,000 Evaluation Budget (With Q1 & Q3 Error Bars)



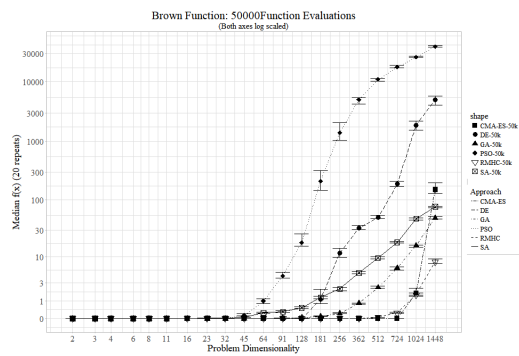
Ackley Function at 50,000 Evaluations



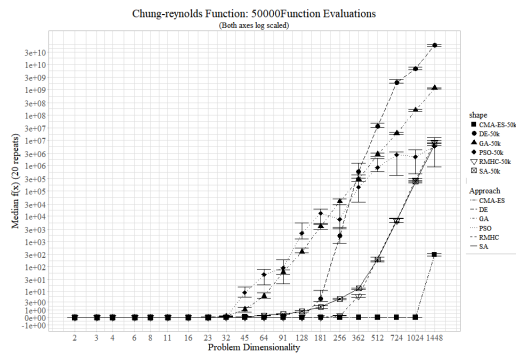
Alpine n.1 Function at 50,000 Evaluations



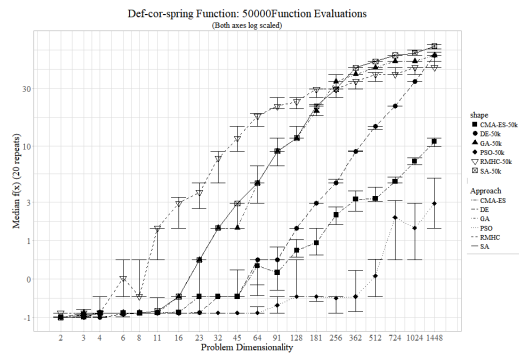
Bent Cigar Function at 50,000 Evaluations



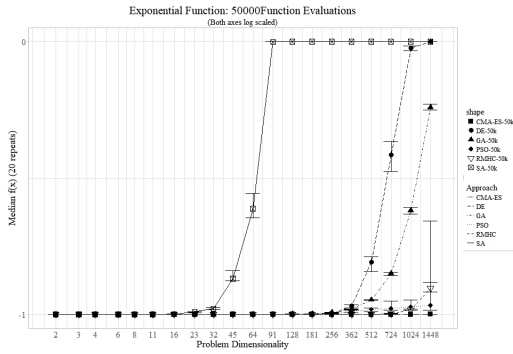
Brown Function at 50,000 Evaluations



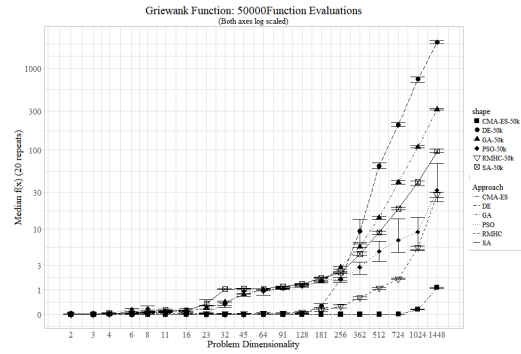
Chung-Reynolds Function at 50,000 Evaluations



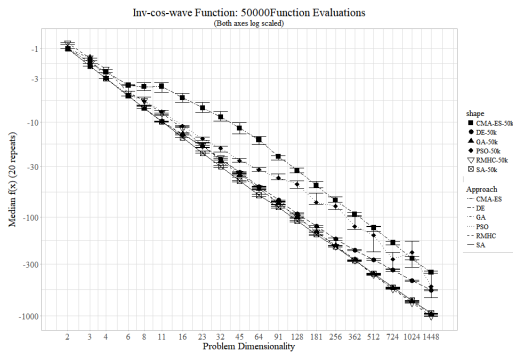
Deflected Corrugated Spring Function at 50,000 Evaluations



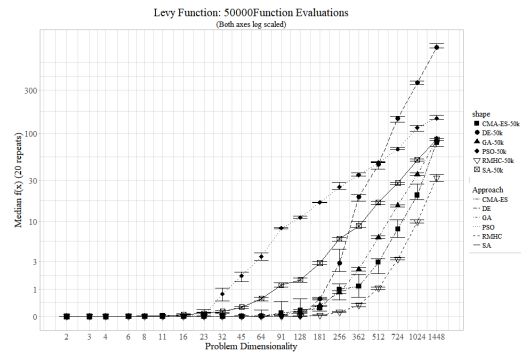
Exponential Function at 50,000 Evaluations



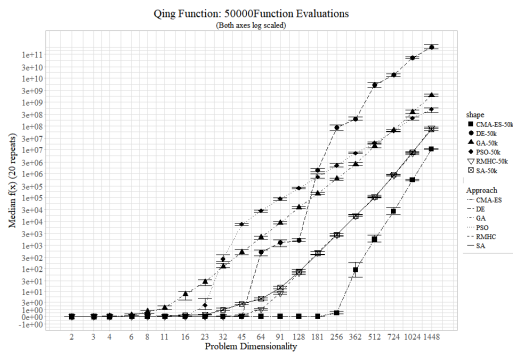
Griewank Function at 50,000 Evaluations



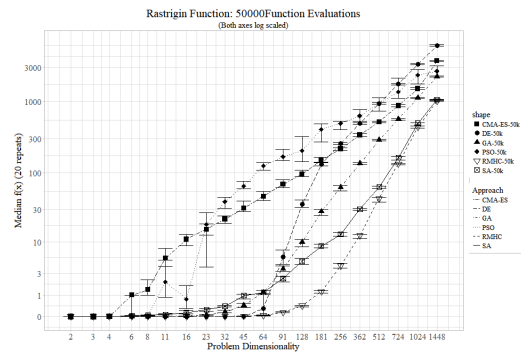
Inverted Cosine Wave Function at 50,000 Evaluations



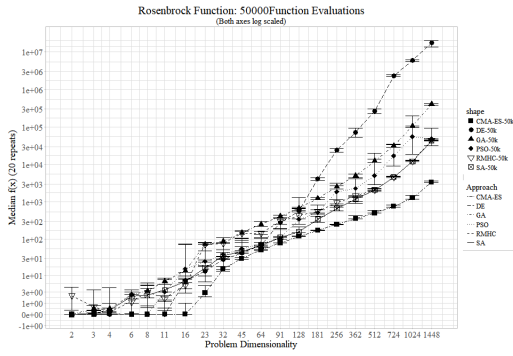
Levy Function at 50,000 Evaluations



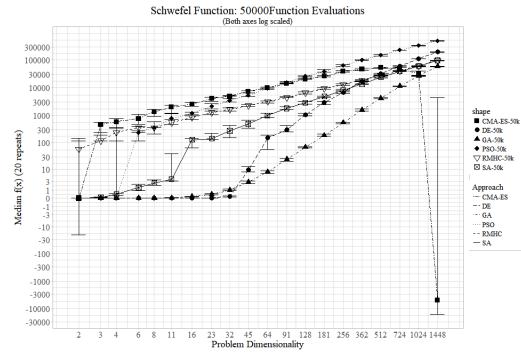
Qing Function at 50,000 Evaluations



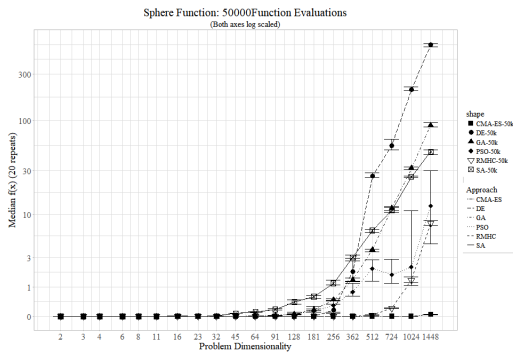
Rastrigin Function at 50,000 Evaluations



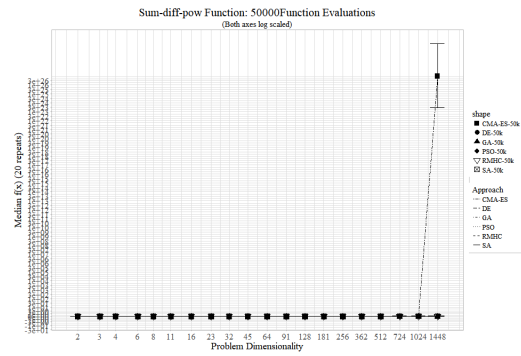
Rosenbrock Function at 50,000 Evaluations



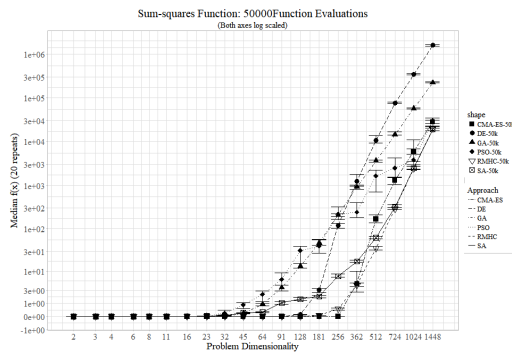
Schwefel Function at 50,000 Evaluations



Sphere Function at 50,000 Evaluations

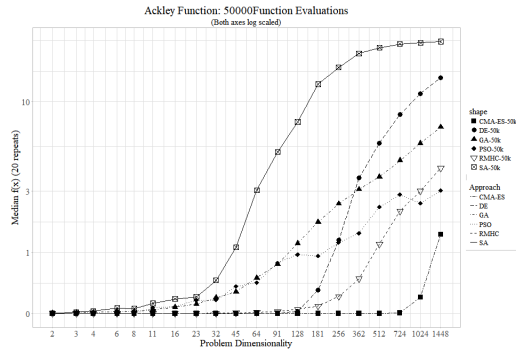


Sum of Different Powers Function at 50,000 Evaluations

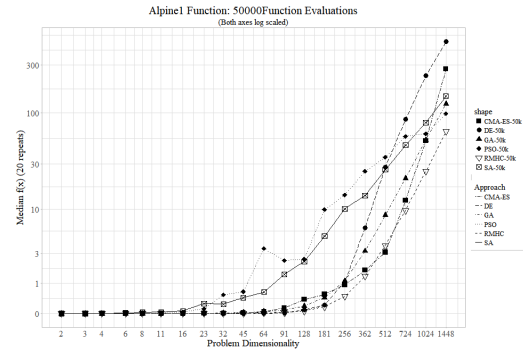


Sum Squares Function at 50,000 Evaluations

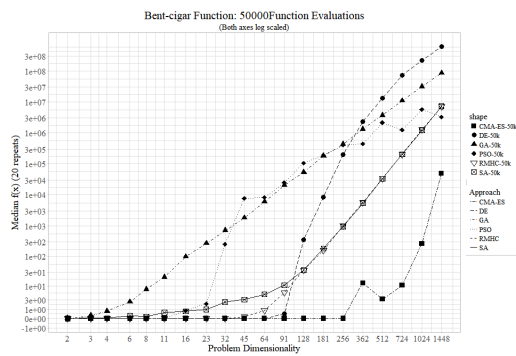
C.1.4 Comparison of Approaches Tuned at the 50,000 Evaluation Budget (No Error Bars)



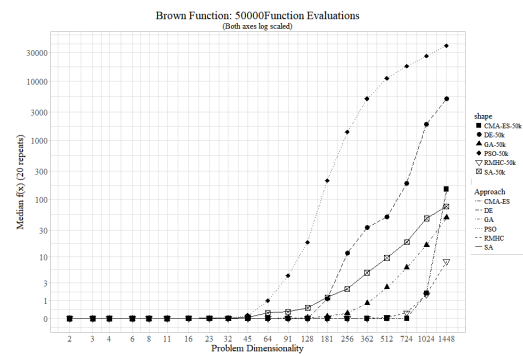
Ackley Function at 50,000 Evaluations



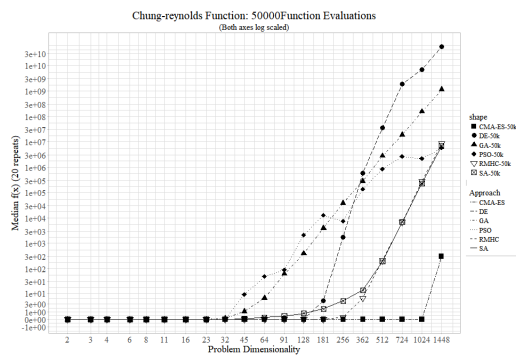
Alpine n.1 Function at 50,000 Evaluations



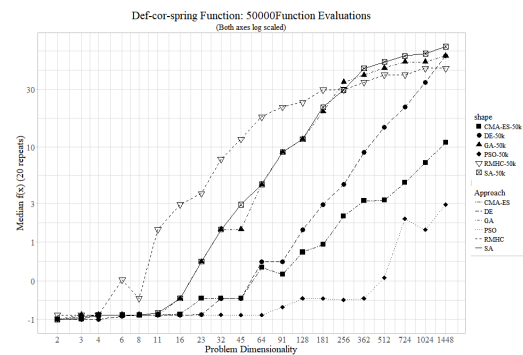
Bent Cigar Function at 50,000 Evaluations



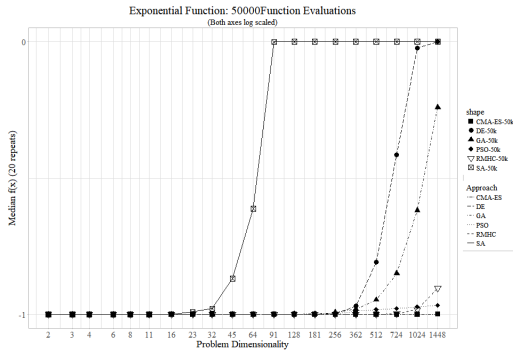
Brown Function at 50,000 Evaluations



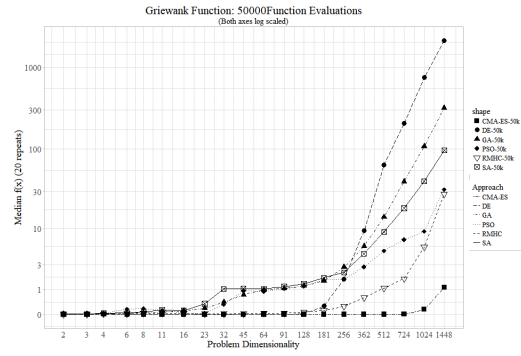
Chung-Reynolds Function at 50,000 Evaluations



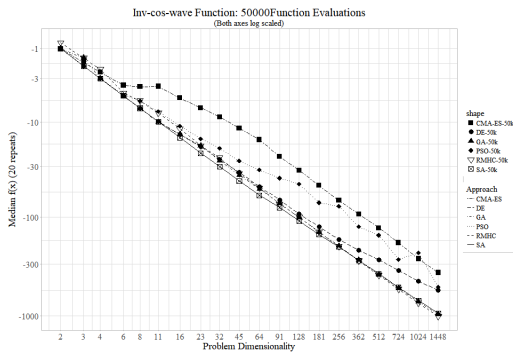
Deflected Corrugated Spring Function at 50,000 Evaluations



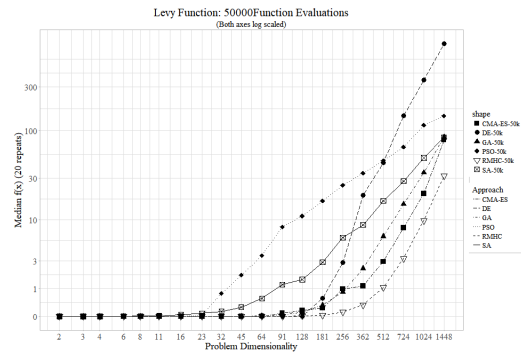
Exponential Function at 50,000 Evaluations



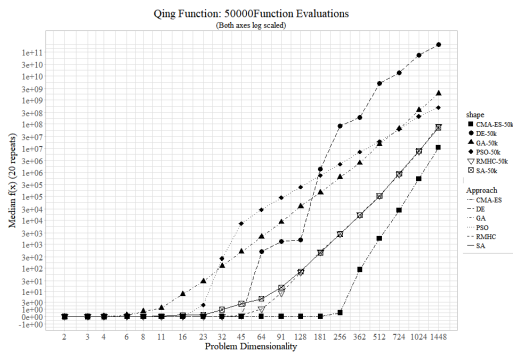
Griewank Function at 50,000 Evaluations



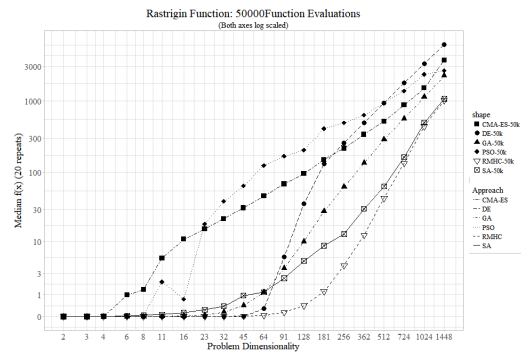
Inverted Cosine Wave Function at 50,000 Evaluations



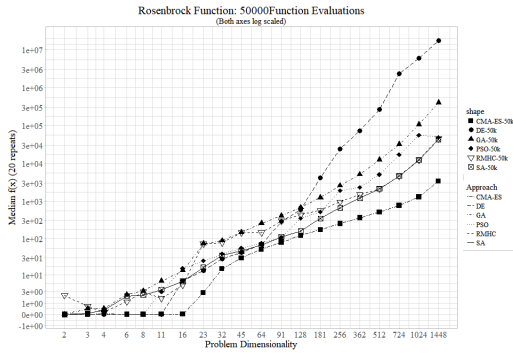
Levy Function at 50,000 Evaluations



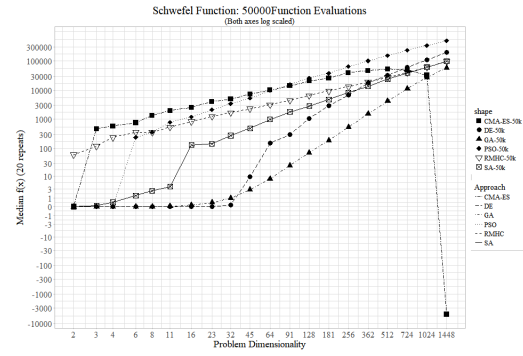
Qing Function at 50,000 Evaluations



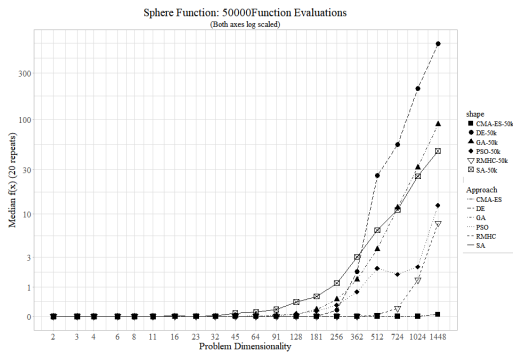
Rastrigin Function at 50,000 Evaluations



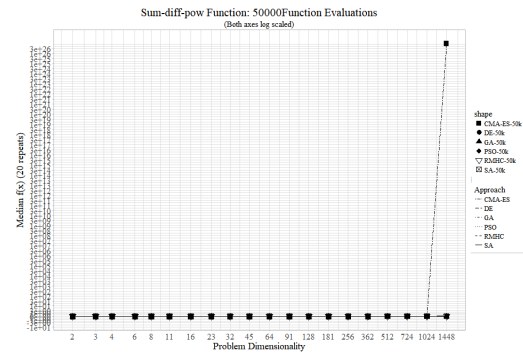
Rosenbrock Function at 50,000 Evaluations



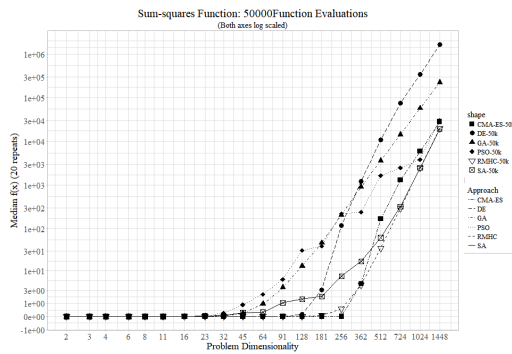
Schwefel Function at 50,000 Evaluations



Sphere Function at 50,000 Evaluations



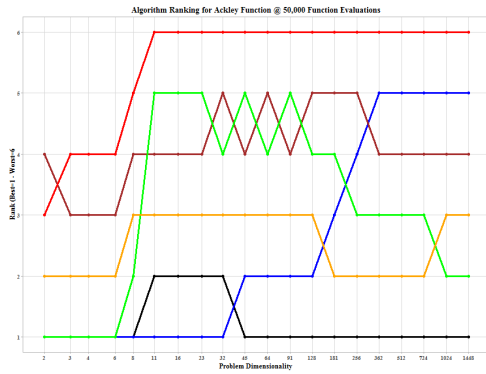
Sum of Different Powers Function at 50,000 Evaluations



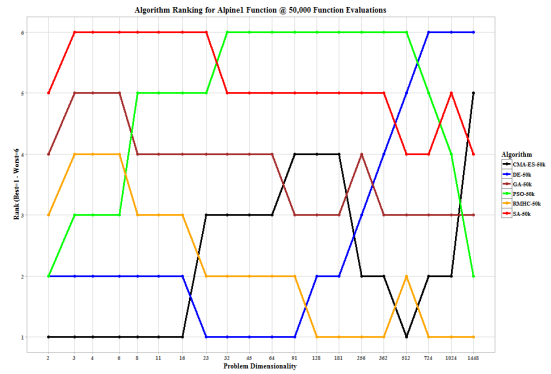
Sum Squares Function at 50,000 Evaluations

C.2 GRAPHICAL METAHEURISTIC RANKINGS

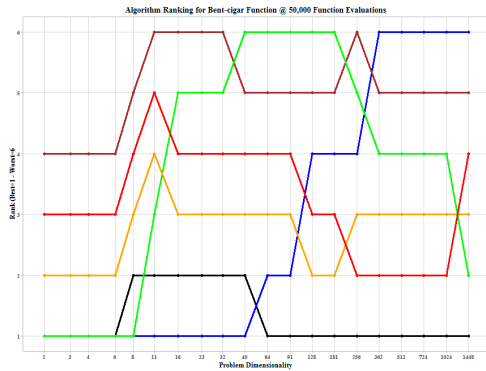
C.2.1 Rank Ordering By Benchmark Function



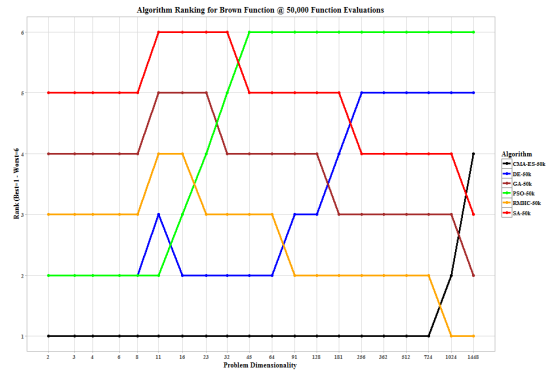
Rankings for Ackley Function



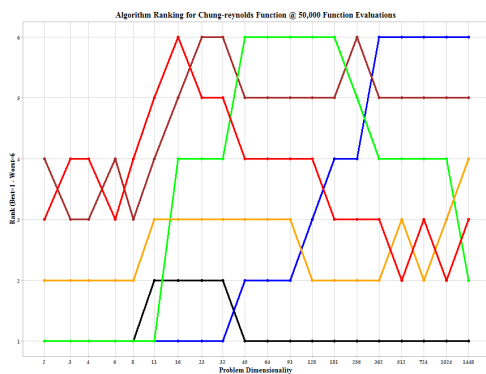
Rankings for Alpine no.1 Function



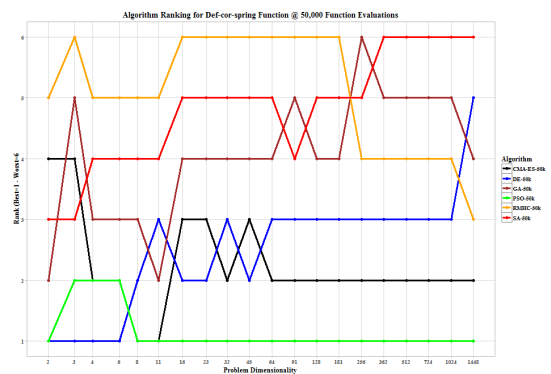
Rankings for Bent Cigar Function



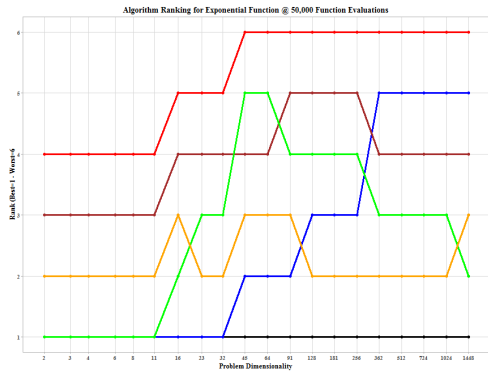
Rankings for Brown Function



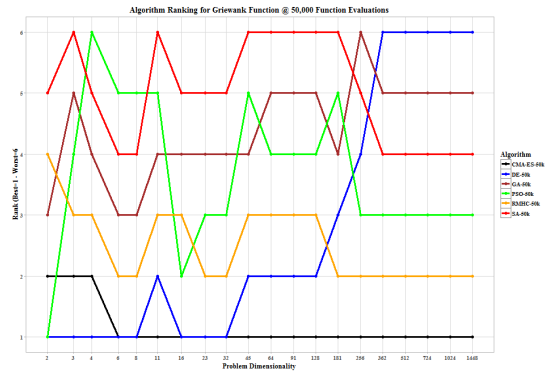
Rankings for Chung Reynolds Function



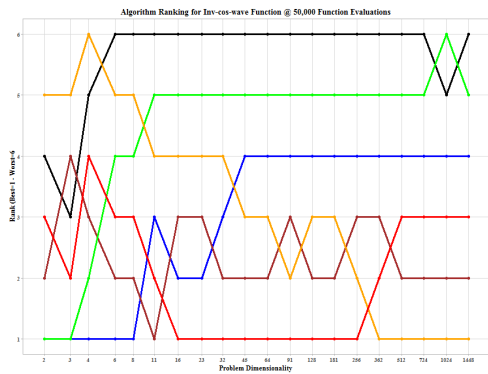
Rankings for Deflected Corrugated Spring Function



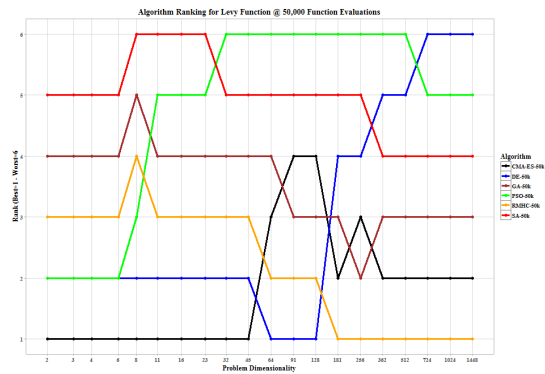
Rankings for Exponential Function



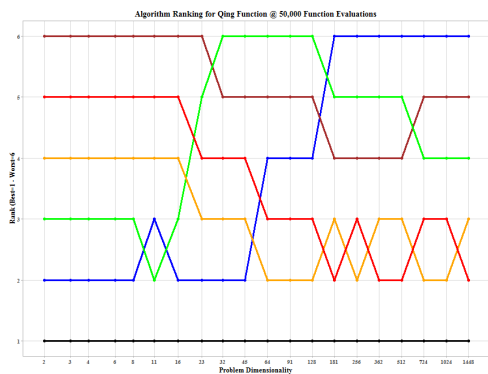
Rankings for Griewank Function



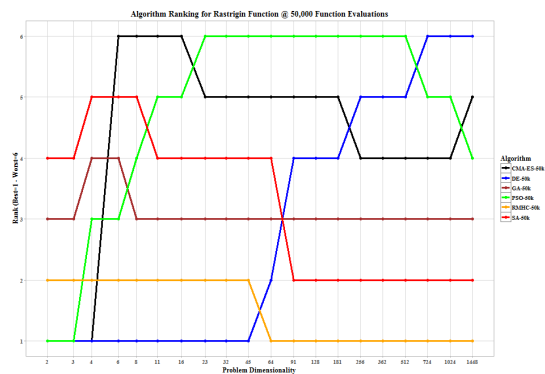
Rankings for Inverted Cosine Wave Function



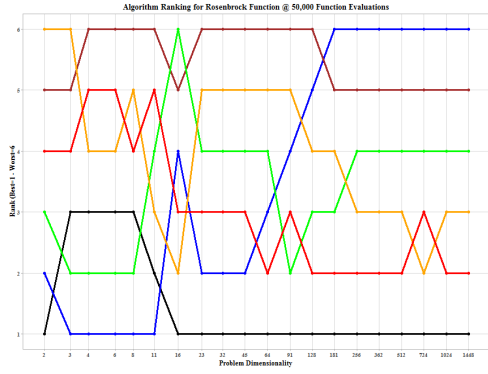
Rankings for Levy Function



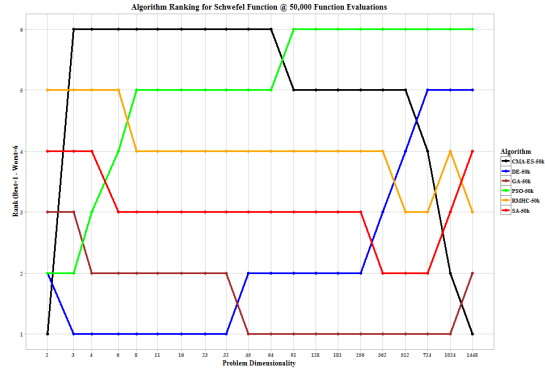
Rankings for Qing Function



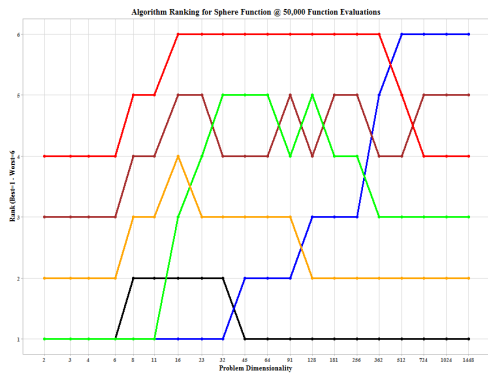
Rankings for Rastrigin Function



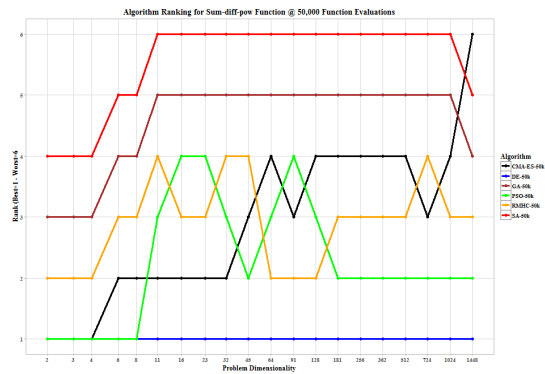
Rankings for Rosenbrock Function



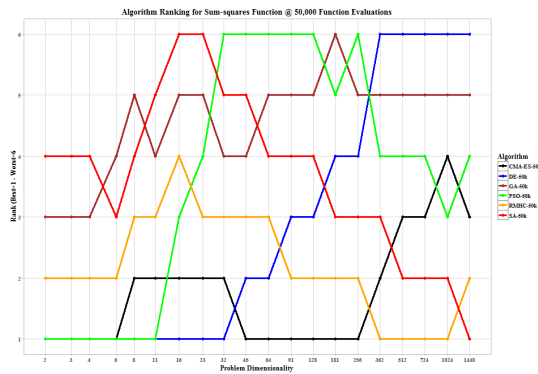
Rankings for Schwefel Function



Rankings for Sphere Function



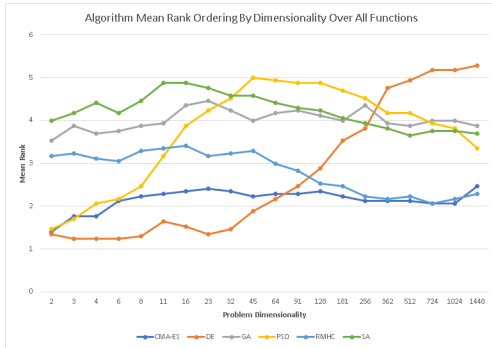
Rankings for Sum of Different Powers Function



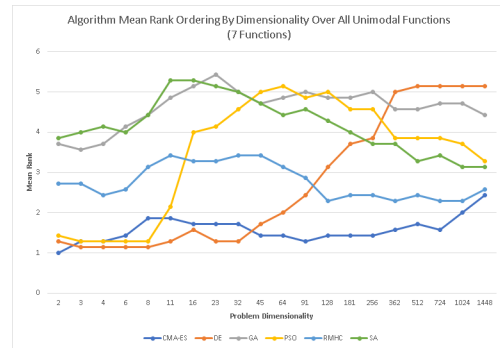
Rankings for Sum Squares Function

C.2.2 Rank Ordering: Mean Aggregate and Pooled Mean Over Subsets of the Benchmark Suite

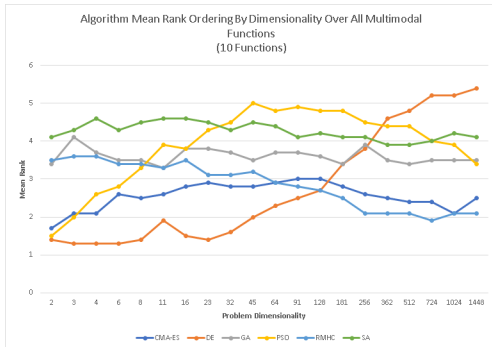
C.2.2.1 Mean Aggregated Rank Ordering Over All Dimensionalities



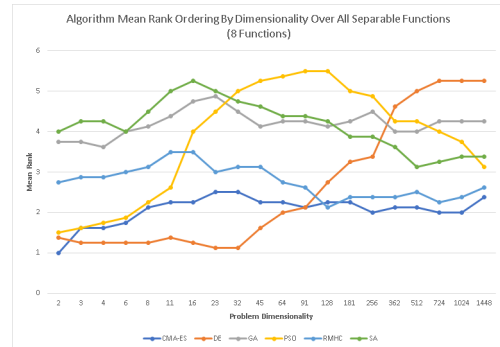
Mean Aggregate Rankings Over Entire Benchmark Suite



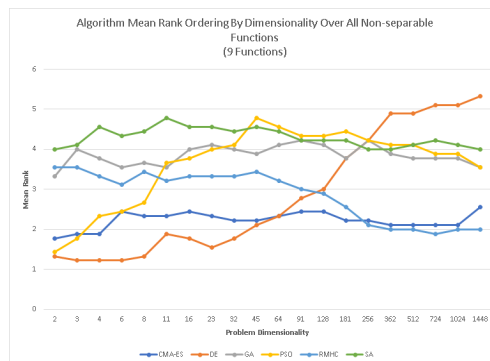
Mean Aggregate Rankings Over All Uni-modal Problems



Mean Aggregate Rankings Over All Multi-modal Problems

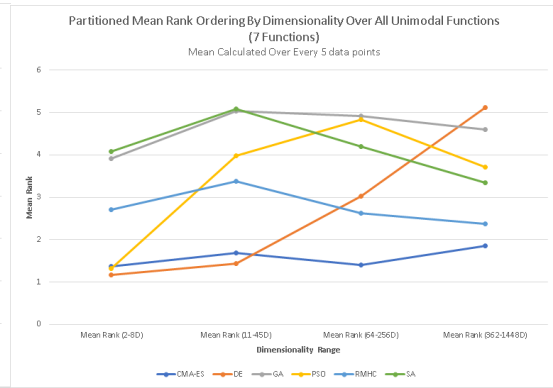
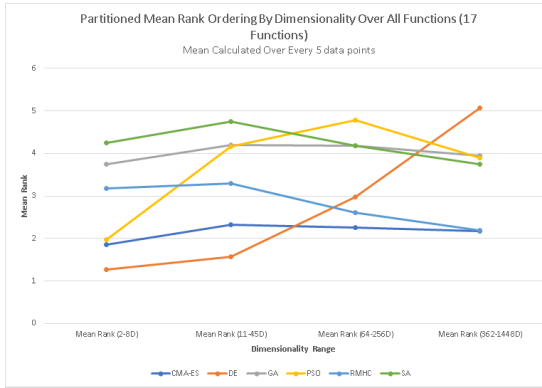


Mean Aggregate Rankings Over All Separable Problems



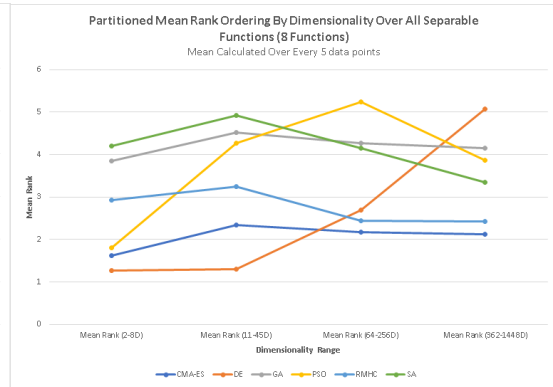
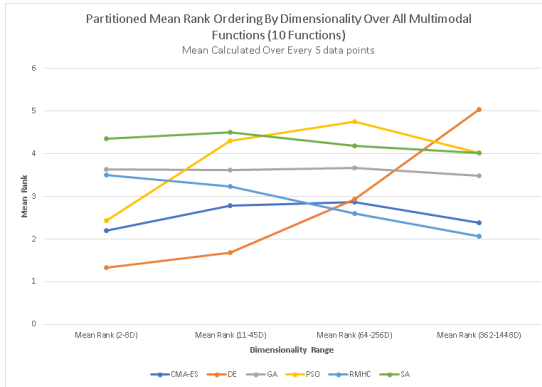
Mean Aggregate Rankings Over All Non-separable Problems

C.2.2.2 Pooled Mean Rank Ordering Over Several Dimension Ranges



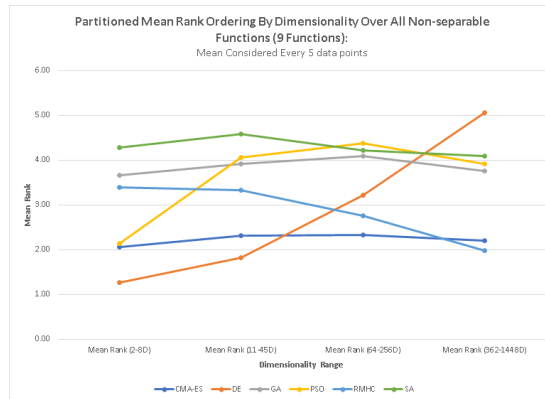
Pooled Mean Rankings Over Entire Benchmark Suite

Pooled Mean Rankings Over All Uni-modal Problems



Pooled Mean Rankings Over All Multi-modal Problems

Pooled Mean Rankings Over All Separable Problems



Pooled Mean Rankings Over All Non-separable Problems

C.3 ALGORITHM RANKING TABLES

C.3.1 *Algorithm Rank Ordering By Benchmark Function*

Table C.1: Algorithm Rankings for Ackley Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	1	2	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1.2
DE	1	1	1	1	1	1	1	1	1	2	2	2	2	3	4	5	5	5	5	5	5	2.45
GA	4	3	3	3	4	4	4	4	5	4	5	4	5	5	5	4	4	4	4	4	4	4.1
PSO	1	1	1	1	2	5	5	5	4	5	4	5	4	4	3	3	3	3	2	2	2	3.15
RMHC	2	2	2	2	3	3	3	3	3	3	3	3	3	2	2	2	2	2	3	3	3	2.55
SA	3	4	4	4	5	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	5.5

Table C.2: Algorithm Rankings for Alpine n.1 Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	1	1	1	3	3	3	3	4	4	4	2	2	1	2	2	5	5	2.25
DE	2	2	2	2	2	2	2	1	1	1	1	1	2	2	3	4	5	6	6	6	6	2.65
GA	4	5	5	5	4	4	4	4	4	4	4	3	3	3	4	3	3	3	3	3	3	3.75
PSO	2	3	3	3	5	5	5	5	6	6	6	6	6	6	6	6	6	5	4	2	2	4.8
RMHC	3	4	4	4	3	3	3	2	2	2	2	2	1	1	1	1	2	1	1	1	1	2.15
SA	5	6	6	6	6	6	6	6	5	5	5	5	5	5	5	5	4	4	5	4	4	5.2

Table C.3: Algorithm Rankings for Bent Cigar Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1.3
DE	1	1	1	1	1	1	1	1	1	1	2	2	4	4	4	6	6	6	6	6	6	2.8
GA	4	4	4	4	5	6	6	6	6	5	5	5	5	5	6	5	5	5	5	5	5	5.05
PSO	1	1	1	1	1	3	5	5	5	6	6	6	6	6	5	4	4	4	4	4	2	3.8
RMHC	2	2	2	2	3	4	3	3	3	3	3	3	2	2	3	3	3	3	3	3	3	2.75
SA	3	3	3	3	4	5	4	4	4	4	4	4	3	3	2	2	2	2	2	2	4	3.25

Table C.4: Algorithm Rankings for Brown Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1.2
DE	2	2	2	2	2	3	2	2	2	2	2	3	3	4	5	5	5	5	5	5	5	3.15
GA	4	4	4	4	4	5	5	5	4	4	4	4	4	3	3	3	3	3	3	2	2	3.75
PSO	2	2	2	2	2	2	3	4	5	6	6	6	6	6	6	6	6	6	6	6	6	4.5
RMHC	3	3	3	3	3	4	4	3	3	3	3	2	2	2	2	2	2	2	1	1	1	2.55
SA	5	5	5	5	5	6	6	6	6	5	5	5	5	5	4	4	4	4	4	3	3	4.85

Table C.5: Algorithm Rankings for Chung-Reynolds Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	1	2	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1.2
DE	1	1	1	1	1	1	1	1	1	2	2	2	3	4	4	6	6	6	6	6	6	2.8
GA	4	3	3	4	3	4	5	6	6	5	5	5	5	5	6	5	5	5	5	5	5	4.7
PSO	1	1	1	1	1	1	4	4	4	6	6	6	6	6	5	4	4	4	4	4	2	3.55
RMHC	2	2	2	2	2	3	3	3	3	3	3	3	2	2	2	2	3	2	3	4	4	2.55
SA	3	4	4	3	4	5	6	5	5	4	4	4	4	3	3	3	2	3	2	3	3	3.7

Table C.6: Algorithm Rankings for Deflected Corrugated Spring Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	4	4	2	2	1	1	3	3	2	3	2	2	2	2	2	2	2	2	2	2	2	2.25
DE	1	1	1	1	2	3	2	2	3	2	3	3	3	3	3	3	3	3	3	3	5	2.5
GA	2	5	3	3	3	2	4	4	4	4	4	5	4	4	6	5	5	5	5	4	4	4.05
PSO	1	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1.15
RMHC	5	6	5	5	5	5	6	6	6	6	6	6	6	6	4	4	4	4	4	3	3	5.1
SA	3	3	4	4	4	4	5	5	5	5	5	4	5	5	5	6	6	6	6	6	6	4.8

Table C.7: Algorithm Rankings for Exponential Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
DE	1	1	1	1	1	1	1	1	1	2	2	2	3	3	3	5	5	5	5	5	5	2.45
GA	3	3	3	3	3	3	4	4	4	4	4	5	5	5	5	4	4	4	4	4	4	3.9
PSO	1	1	1	1	1	1	2	3	3	5	5	4	4	4	4	3	3	3	3	2	2	2.7
RMHC	2	2	2	2	2	2	3	2	2	3	3	3	2	2	2	2	2	2	2	2	3	2.25
SA	4	4	4	4	4	4	5	5	5	6	6	6	6	6	6	6	6	6	6	6	6	5.25

Table C.8: Algorithm Rankings for Griewank Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1.15
DE	1	1	1	1	1	2	1	1	1	2	2	2	2	3	4	6	6	6	6	6	6	2.75
GA	3	5	4	3	3	4	4	4	4	4	5	5	5	4	6	5	5	5	5	5	5	4.4
PSO	1	4	6	5	5	5	2	3	3	5	4	4	4	5	3	3	3	3	3	3	3	3.7
RMHC	4	3	3	2	2	3	3	2	2	3	3	3	3	2	2	2	2	2	2	2	2	2.5
SA	5	6	5	4	4	6	5	5	5	6	6	6	6	6	5	4	4	4	4	4	4	5

Table C.9: Algorithm Rankings for Inverted Cosine Wave Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	4	3	5	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	5	6	6	5.65
DE	1	1	1	1	1	3	2	2	3	4	4	4	4	4	4	4	4	4	4	4	4	2.95
GA	2	4	3	2	2	1	3	3	2	2	2	3	2	2	3	3	2	2	2	2	2	2.35
PSO	1	1	2	4	4	5	5	5	5	5	5	5	5	5	5	5	5	5	6	5	5	4.4
RMHC	5	5	6	5	5	4	4	4	4	3	3	2	3	3	2	1	1	1	1	1	1	3.15
SA	3	2	4	3	3	2	1	1	1	1	1	1	1	1	1	2	3	3	3	3	3	2

Table C.10: Algorithm Rankings for Levy Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	1	1	1	1	1	1	3	4	4	2	3	2	2	2	2	2	2	1.8
DE	2	2	2	2	2	2	2	2	2	2	1	1	1	4	4	5	5	6	6	6	6	2.95
GA	4	4	4	4	5	4	4	4	4	4	4	3	3	3	2	3	3	3	3	3	3	3.55
PSO	2	2	2	2	3	5	5	5	6	6	6	6	6	6	6	6	6	5	5	5	5	4.75
RMHC	3	3	3	3	4	3	3	3	3	3	2	2	2	1	1	1	1	1	1	1	1	2.2
SA	5	5	5	5	6	6	6	6	5	5	5	5	5	5	5	4	4	4	4	4	4	4.95

Table C.11: Algorithm Rankings for Qing Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
DE	2	2	2	2	2	3	2	2	2	2	4	4	4	6	6	6	6	6	6	6	6	3.75
GA	6	6	6	6	6	6	6	6	5	5	5	5	5	4	4	4	4	5	5	5	5	5.2
PSO	3	3	3	3	3	2	3	5	6	6	6	6	6	5	5	5	5	4	4	4	4	4.35
RMHC	4	4	4	4	4	4	4	3	3	3	2	2	2	3	2	3	3	2	2	3	3	3.05
SA	5	5	5	5	5	5	5	4	4	4	3	3	3	2	3	2	2	3	3	2	2	3.65

Table C.12: Algorithm Rankings for Rastrigin Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	6	6	6	6	5	5	5	5	5	5	5	4	4	4	4	4	5	5	4.35
DE	1	1	1	1	1	1	1	1	1	1	2	4	4	4	5	5	5	6	6	6	6	2.85
GA	3	3	4	4	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3.1
PSO	1	1	3	3	4	5	5	6	6	6	6	6	6	6	6	6	6	5	5	4	4	4.8
RMHC	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1.5
SA	4	4	5	5	5	4	4	4	4	4	4	2	2	2	2	2	2	2	2	2	2	3.25

Table C.13: Algorithm Rankings for Rosenbrock Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	3	3	3	3	2	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1.45
DE	2	1	1	1	1	1	4	2	2	2	3	4	5	6	6	6	6	6	6	6	6	3.55
GA	5	5	6	6	6	6	5	6	6	6	6	6	6	5	5	5	5	5	5	5	5	5.5
PSO	3	2	2	2	2	4	6	4	4	4	4	2	3	3	4	4	4	4	4	4	4	3.45
RMHC	6	6	4	4	5	3	2	5	5	5	5	5	4	4	3	3	3	2	3	3	3	4
SA	4	4	5	5	4	5	3	3	3	3	2	3	2	2	2	2	2	3	2	2	2	3.05

Table C.14: Algorithm Rankings for Schwefel Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	6	6	6	6	6	6	6	6	6	6	5	5	5	5	5	5	4	2	1	1	4.9
DE	2	1	1	1	1	1	1	1	1	2	2	2	2	2	2	3	4	5	5	5	5	2.2
GA	3	3	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	2	2	1.6
PSO	2	2	3	4	5	5	5	5	5	5	5	6	6	6	6	6	6	6	6	6	6	5
RMHC	5	5	5	5	4	4	4	4	4	4	4	4	4	4	4	4	3	3	4	3	3	4.05
SA	4	4	4	3	3	3	3	3	3	3	3	3	3	3	3	2	2	2	3	4	4	3.05

Table C.15: Algorithm Rankings for Sphere Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	1	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1.25
DE	1	1	1	1	1	1	1	1	1	2	2	2	3	3	3	5	6	6	6	6	6	2.65
GA	3	3	3	3	4	4	5	5	4	4	4	5	4	5	5	4	4	5	5	5	5	4.2
PSO	1	1	1	1	1	1	3	4	5	5	5	4	5	4	4	3	3	3	3	3	3	3
RMHC	2	2	2	2	3	3	4	3	3	3	3	3	2	2	2	2	2	2	2	2	2	2.45
SA	4	4	4	4	5	5	6	6	6	6	6	6	6	6	6	6	5	4	4	4	4	5.15

Table C.16: Algorithm Rankings for Sum of Different Powers Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>	
Sep-CMA-ES	1	1	1	2	2	2	2	2	2	3	4	3	4	4	4	4	4	3	4	4	6	2.9
DE	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
GA	3	3	3	4	4	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	4	4.55
PSO	1	1	1	1	1	3	4	4	3	2	3	4	3	2	2	2	2	2	2	2	2	2.25
RMHC	2	2	2	3	3	4	3	3	4	4	2	2	2	3	3	3	3	4	3	3	3	2.9
SA	4	4	4	5	5	6	6	6	6	6	6	6	6	6	6	6	6	6	6	5	5	5.55

Table C.17: Algorithm Rankings for Sum Squares Function over 2 to 1448 Dimensions (1=Best, 5=Worst)

<i>Alg. / Dim.</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean</i>
Sep-CMA-ES	1	1	1	1	2	2	2	2	2	1	1	1	1	1	1	2	3	3	4	3	1.75
DE	1	1	1	1	1	1	1	1	1	2	2	3	3	4	4	6	6	6	6	6	2.85
GA	3	3	3	4	5	4	5	5	4	4	5	5	5	6	5	5	5	5	5	5	4.55
PSO	1	1	1	1	1	1	3	4	6	6	6	6	6	5	6	4	4	4	3	4	3.65
RMHC	2	2	2	2	3	3	4	3	3	3	3	2	2	2	2	1	1	1	1	2	2.2
SA	4	4	4	3	4	5	6	6	5	5	4	4	4	3	3	3	2	2	2	1	3.7

c.3.2 Mean Aggregated Algorithm Rank Ordering By All Problems and Problems Classified by Problem Feature Investigated

Table C.18: Mean Aggregated Algorithm Rankings Over All Benchmark Functions

<i>Alg.</i> / <i>Dim</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean Rank</i>	
CMA-ES	1.41	1.76	1.76	2.12	2.24	2.29	2.35	2.41	2.35	2.24	2.29	2.29	2.35	2.24	2.12	2.12	2.12	2.06	2.06	2.06	2.47	2.15
DE	1.35	1.24	1.24	1.24	1.29	1.65	1.53	1.35	1.47	1.88	2.18	2.47	2.88	3.53	3.82	4.76	4.94	5.18	5.18	5.29	5.29	2.72
GA	3.53	3.88	3.71	3.76	3.88	3.94	4.35	4.47	4.24	4.00	4.18	4.24	4.12	4.00	4.35	3.94	3.88	4.00	4.00	3.88	3.88	4.02
PSO	1.47	1.71	2.06	2.18	2.47	3.18	3.88	4.24	4.53	5.00	4.94	4.88	4.88	4.71	4.53	4.18	4.18	3.94	3.82	3.35	3.35	3.71
RMHC	3.18	3.24	3.12	3.06	3.29	3.35	3.41	3.18	3.24	3.29	3.00	2.82	2.53	2.47	2.24	2.18	2.24	2.06	2.18	2.29	2.29	2.82
SA	4.00	4.18	4.41	4.18	4.47	4.88	4.88	4.76	4.59	4.59	4.41	4.29	4.24	4.06	3.94	3.82	3.65	3.76	3.76	3.71	3.71	4.23

Table C.19: Mean Aggregated Algorithm Rankings Over All Uni-modal Functions

<i>Alg.</i> / <i>Dim</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean Rank</i>	
CMA-ES	1.00	1.29	1.29	1.43	1.86	1.86	1.71	1.71	1.71	1.43	1.43	1.29	1.43	1.43	1.43	1.57	1.71	1.57	2.00	2.43	2.43	1.58
DE	1.29	1.14	1.14	1.14	1.14	1.29	1.57	1.29	1.29	1.71	2.00	2.43	3.14	3.71	3.86	5.00	5.14	5.14	5.14	5.14	5.14	2.69
GA	3.71	3.57	3.71	4.14	4.43	4.86	5.14	5.43	5.00	4.71	4.86	5.00	4.86	4.86	5.00	4.57	4.57	4.71	4.71	4.43	4.43	4.61
PSO	1.43	1.29	1.29	1.29	1.29	2.14	4.00	4.14	4.57	5.00	5.14	4.86	5.00	4.57	4.57	3.86	3.86	3.71	3.71	3.29	3.29	3.46
RMHC	2.71	2.71	2.43	2.57	3.14	3.43	3.29	3.29	3.43	3.43	3.14	2.86	2.29	2.43	2.43	2.29	2.43	2.29	2.29	2.57	2.57	2.77
SA	3.86	4.00	4.14	4.00	4.43	5.29	5.29	5.14	5.00	4.71	4.43	4.57	4.29	4.00	3.71	3.71	3.29	3.43	3.14	3.14	3.14	4.18

Table C.20: Mean Aggregated Algorithm Rankings Over All Multi-modal Functions

<i>Alg.</i> / <i>Dim</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean Rank</i>
CMA-ES	1.70	2.10	2.10	2.60	2.50	2.60	2.80	2.90	2.80	2.80	2.90	3.00	3.00	2.80	2.60	2.50	2.40	2.40	2.10	2.50	2.56
DE	1.40	1.30	1.30	1.30	1.40	1.90	1.50	1.40	1.60	2.00	2.30	2.50	2.70	3.40	3.80	4.60	4.80	5.20	5.20	5.40	2.75
GA	3.40	4.10	3.70	3.50	3.50	3.30	3.80	3.80	3.70	3.50	3.70	3.70	3.60	3.40	3.90	3.50	3.40	3.50	3.50	3.50	3.60
PSO	1.50	2.00	2.60	2.80	3.30	3.90	3.80	4.30	4.50	5.00	4.80	4.90	4.80	4.80	4.50	4.40	4.40	4.00	3.90	3.40	3.88
RMHC	3.50	3.60	3.60	3.40	3.40	3.30	3.50	3.10	3.10	3.20	2.90	2.80	2.70	2.50	2.10	2.10	2.10	1.90	2.10	2.10	2.85
SA	4.10	4.30	4.60	4.30	4.50	4.60	4.60	4.50	4.30	4.50	4.40	4.10	4.20	4.10	4.10	3.90	3.90	4.00	4.20	4.10	4.27

Table C.21: Mean Aggregated Algorithm Rankings Over All Separable Functions

<i>Alg.</i> / <i>Dim</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean Rank</i>
CMA-ES	1.00	1.63	1.63	1.75	2.13	2.25	2.25	2.50	2.50	2.25	2.25	2.13	2.25	2.25	2.00	2.13	2.13	2.00	2.00	2.38	2.07
DE	1.38	1.25	1.25	1.25	1.25	1.38	1.25	1.13	1.13	1.63	2.00	2.13	2.75	3.25	3.38	4.63	5.00	5.25	5.25	5.25	2.59
GA	3.75	3.75	3.63	4.00	4.13	4.38	4.75	4.88	4.50	4.13	4.25	4.25	4.13	4.25	4.50	4.00	4.00	4.25	4.25	4.25	4.20
PSO	1.50	1.63	1.75	1.88	2.25	2.63	4.00	4.50	5.00	5.25	5.38	5.50	5.50	5.00	4.88	4.25	4.25	4.00	3.75	3.13	3.80
RMHC	2.75	2.88	2.88	3.00	3.13	3.50	3.50	3.00	3.13	3.13	2.75	2.63	2.13	2.38	2.38	2.38	2.50	2.25	2.38	2.63	2.76
SA	4.00	4.25	4.25	4.00	4.50	5.00	5.25	5.00	4.75	4.63	4.38	4.38	4.25	3.88	3.88	3.63	3.13	3.25	3.38	3.38	4.16

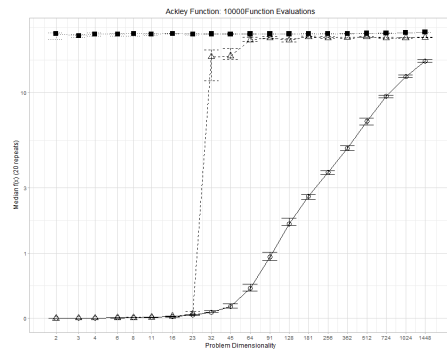
Table C.22: Mean Aggregated Algorithm Rankings Over All Non-separable Functions

<i>Alg.</i>	<i>Dim</i>	2	3	4	6	8	11	16	23	32	45	64	91	128	181	256	362	512	724	1024	1448	<i>Mean Rank</i>	
CMA-ES		1.78	1.89	1.89	2.44	2.33	2.33	2.44	2.33	2.22	2.22	2.33	2.44	2.44	2.22	2.22	2.11	2.11	2.11	2.11	2.11	2.56	2.23
DE		1.33	1.22	1.22	1.22	1.33	1.89	1.78	1.56	1.78	2.11	2.33	2.78	3.00	3.78	4.22	4.89	4.89	5.11	5.11	5.33	2.84	
GA		3.33	4.00	3.78	3.56	3.67	3.56	4.00	4.11	4.00	3.89	4.11	4.22	4.11	3.78	4.22	3.89	3.78	3.78	3.78	3.56	3.86	
PSO		1.44	1.78	2.33	2.44	2.67	3.67	3.78	4.00	4.11	4.78	4.56	4.33	4.44	4.44	4.22	4.11	4.11	3.89	3.89	3.56	3.62	
RMHC		3.56	3.56	3.33	3.11	3.44	3.22	3.33	3.33	3.33	3.44	3.22	3.00	2.89	2.56	2.11	2.00	2.00	1.89	2.00	2.00	2.87	
SA		4.00	4.11	4.56	4.33	4.44	4.78	4.56	4.56	4.44	4.56	4.44	4.22	4.22	4.22	4.00	4.00	4.11	4.22	4.11	4.00	4.29	

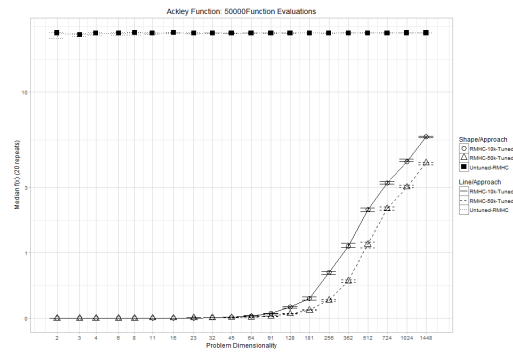
APPENDIX D: PLOTS SHOWING THE RESULTS OF TUNING EACH ALGORITHM AT A 10,000 AND 50,000 EVALUATION BUDGET COMPARED TO THE CORRESPONDING UNTUNED ALGORITHM

D.1 SMAC PARAMETER TUNING RESULTS

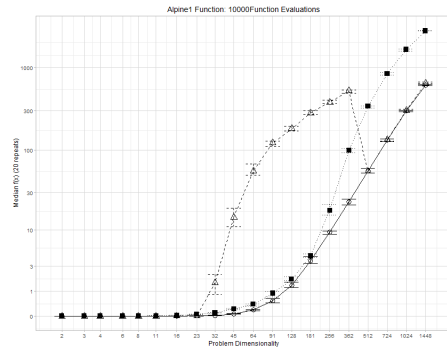
Random Mutation Hill Climbing (RMHC)



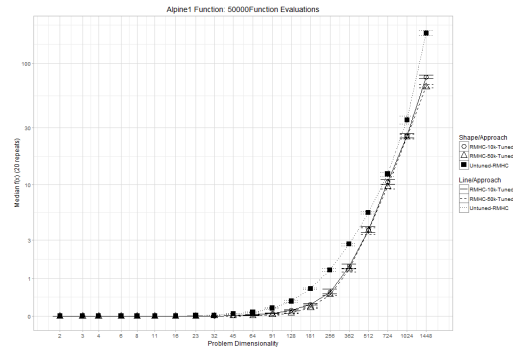
Ackley Function at 10,000 Evaluations



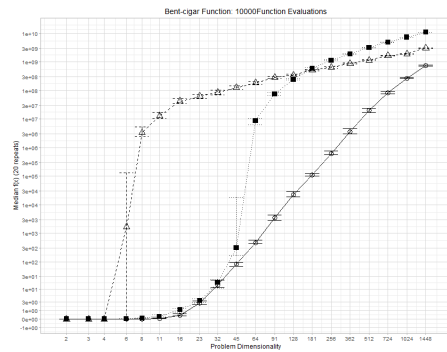
Ackley Function at 50,000 Evaluations



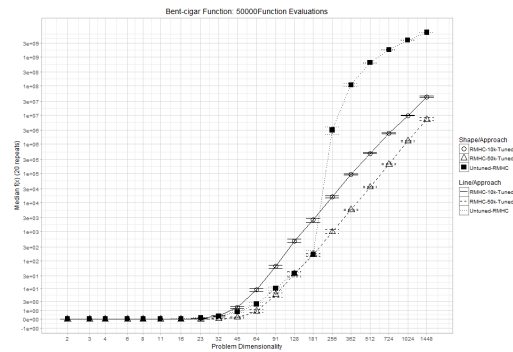
Alpine no.1 Function at 10,000 Evaluations



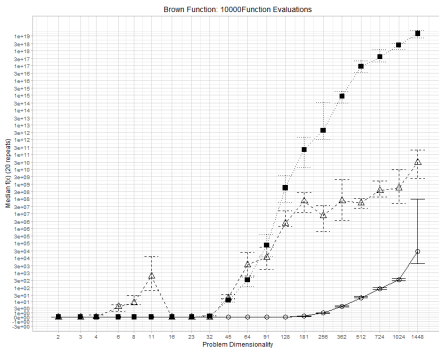
Alpine no.1 Function at 50,000 Evaluations



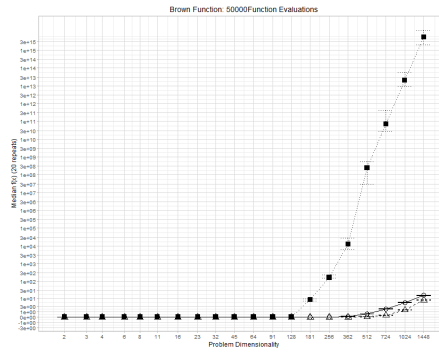
Bent Cigar Function at 10,000 Evaluations



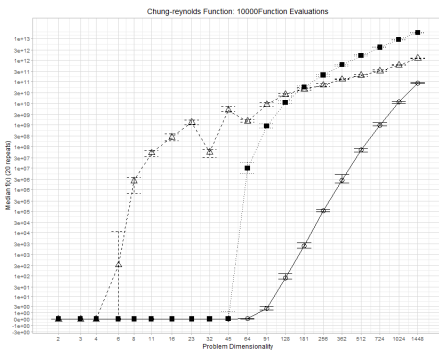
Bent Cigar Function at 50,000 Evaluations



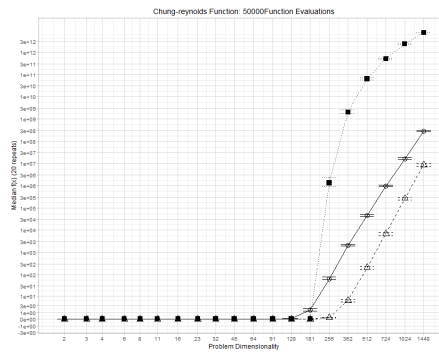
Brown Function at 10,000 Evaluations



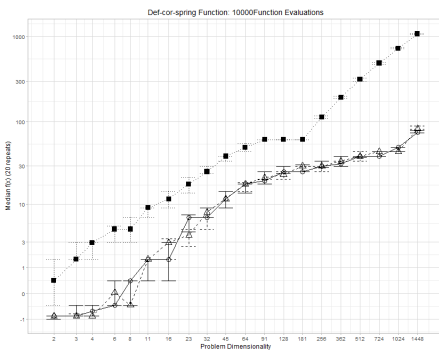
Brown Function at 50,000 Evaluations



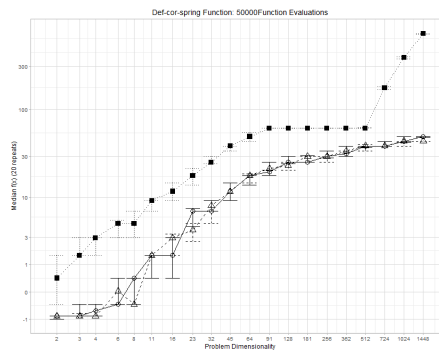
Chung-Reynolds Function at 10,000 Evaluations



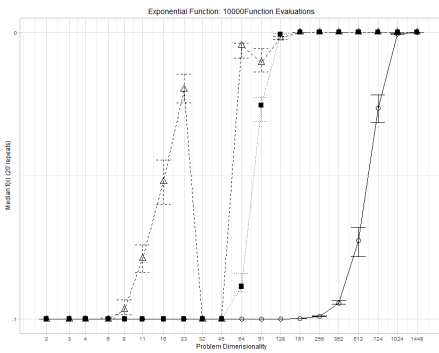
Chung-Reynolds Function at 50,000 Evaluations



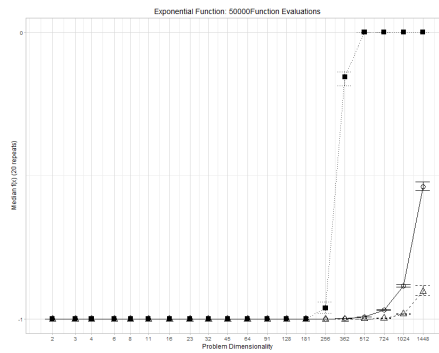
Deflected Corrugated Spring Function at 10,000 Evaluations



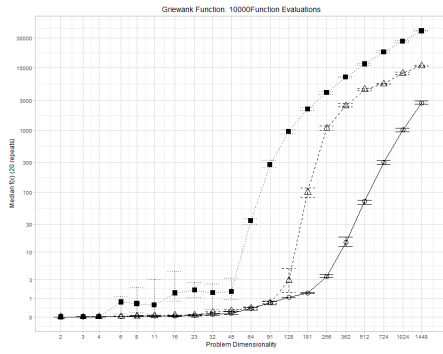
Deflected Corrugated Spring Function at 50,000 Evaluations



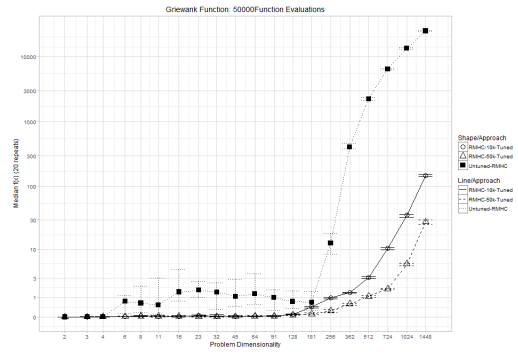
Exponential Function at 10,000 Evaluations



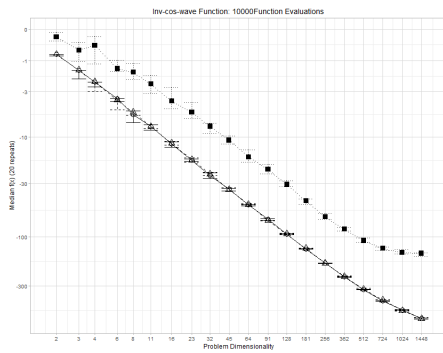
Exponential Function at 50,000 Evaluations



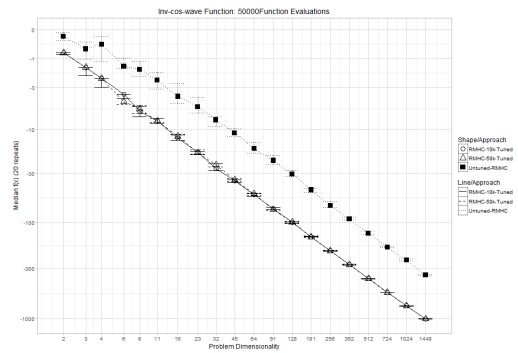
Griewank Function at 10,000 Evaluations



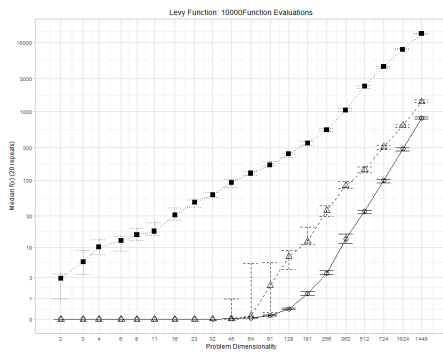
Griewank Function at 50,000 Evaluations



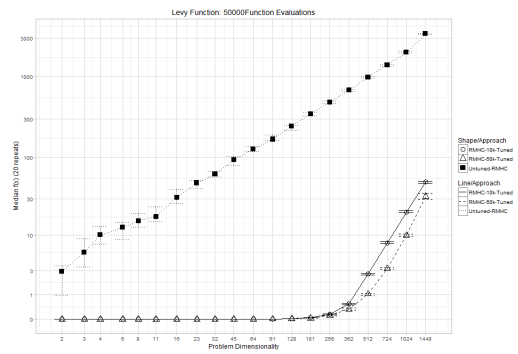
Inverted Cosine Wave Function at 10,000 Evaluations



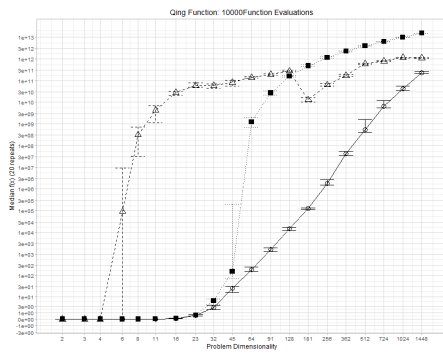
Inverted Cosine Wave Function at 50,000 Evaluations



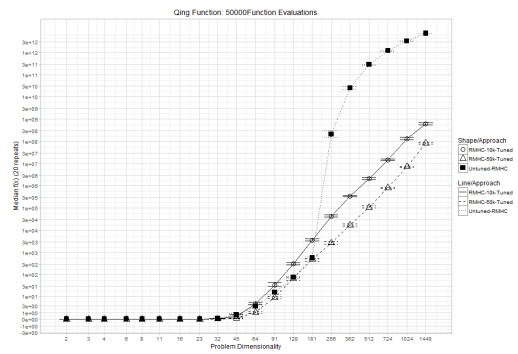
Levy Function at 10,000 Evaluations



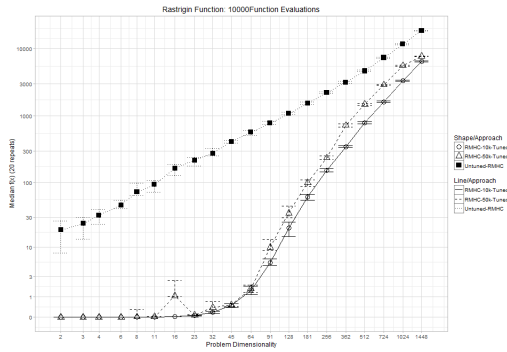
Levy Function at 50,000 Evaluations



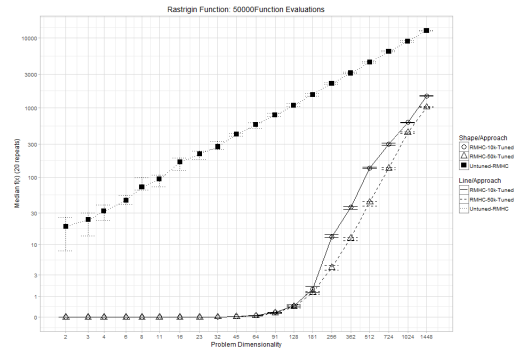
Qing Function at 10,000 Evaluations



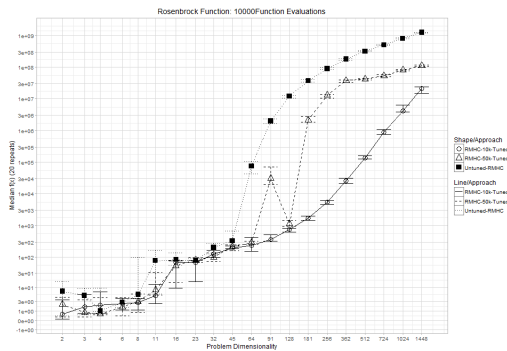
Qing Function at 50,000 Evaluations



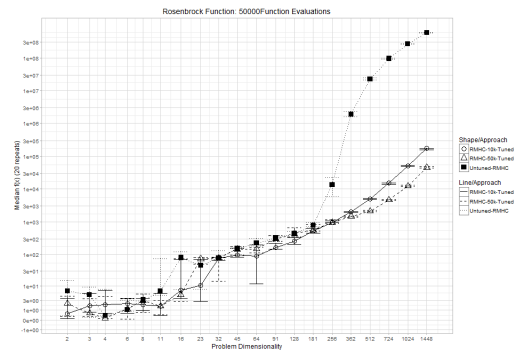
Rastrigin Function at 10,000 Evaluations



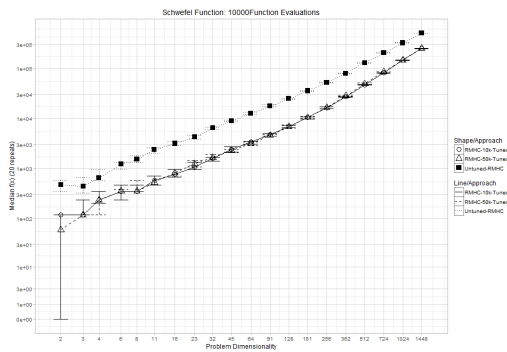
Rastrigin Function at 50,000 Evaluations



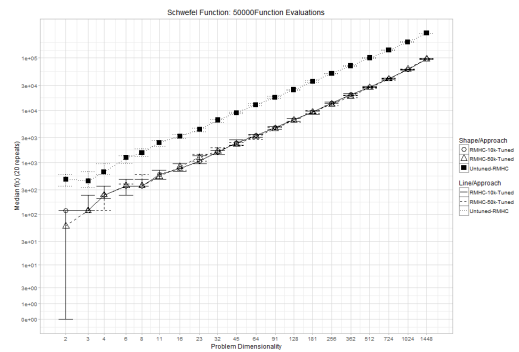
Rosenbrock Function at 10,000 Evaluations



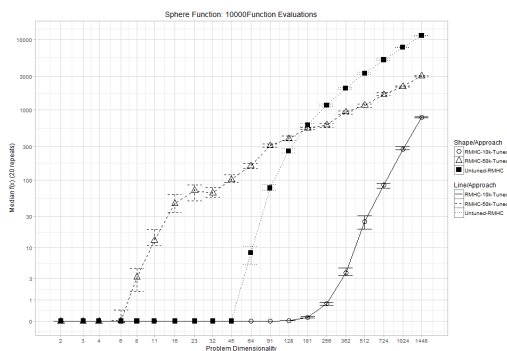
Rosenbrock Function at 50,000 Evaluations



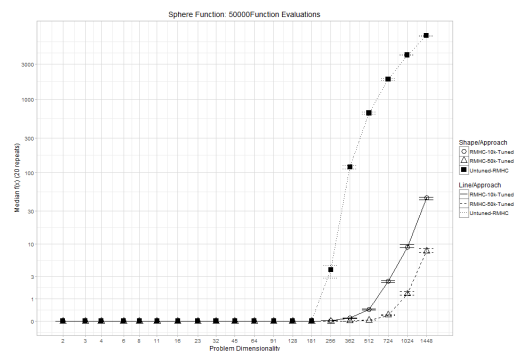
Schwefel Function at 10,000 Evaluations



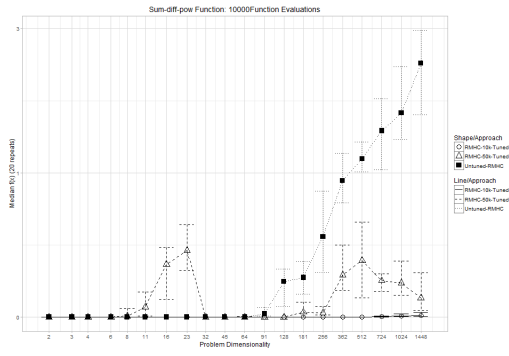
Schwefel Function at 50,000 Evaluations



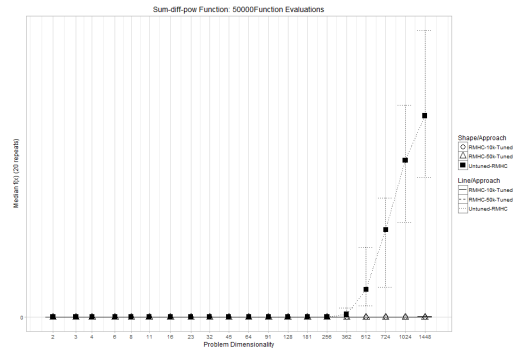
Sphere Function at 10,000 Evaluations



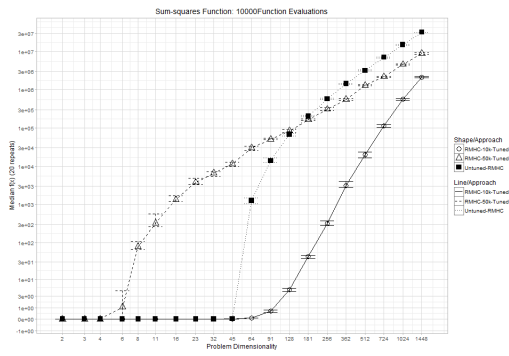
Sphere Function at 50,000 Evaluations



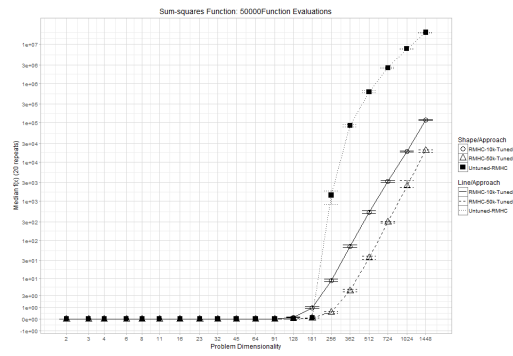
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

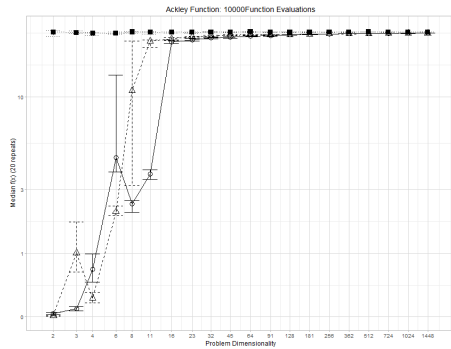


Sum Squares Function at 10,000 Evaluations

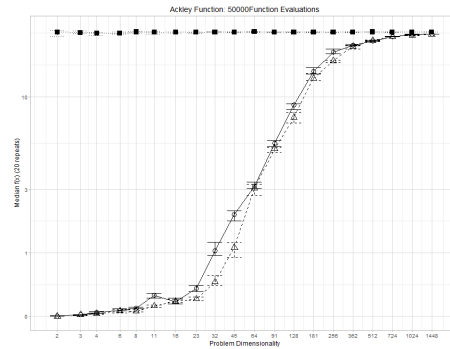


Sum Squares Function at 50,000 Evaluations

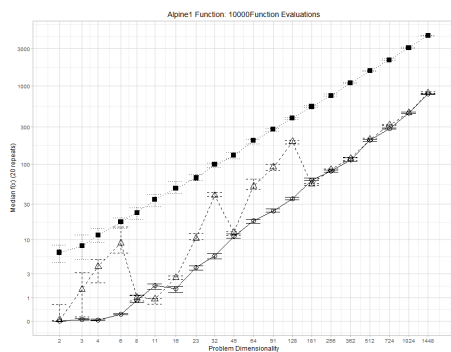
Simulated Annealing (SA)



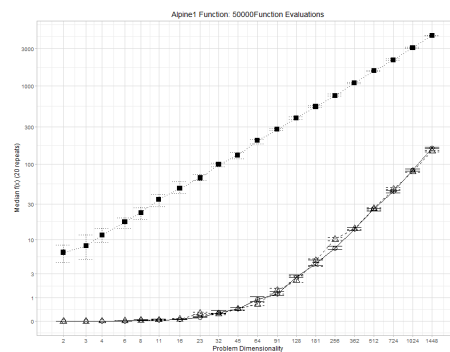
Ackley Function at 10,000 Evaluations



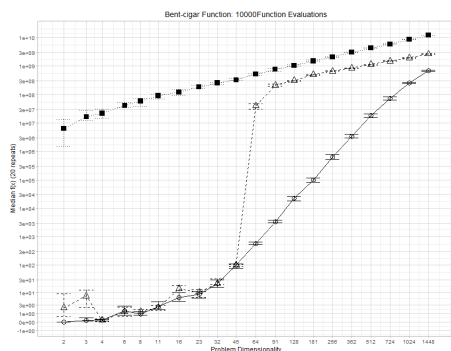
Ackley Function at 50,000 Evaluations



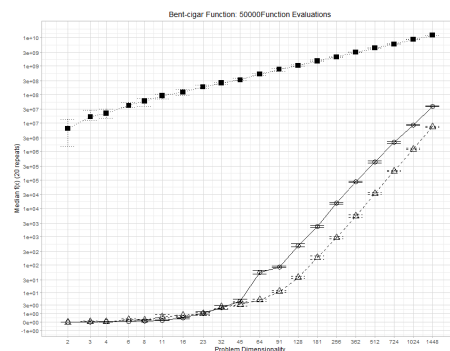
Alpine no.1 Function at 10,000 Evaluations



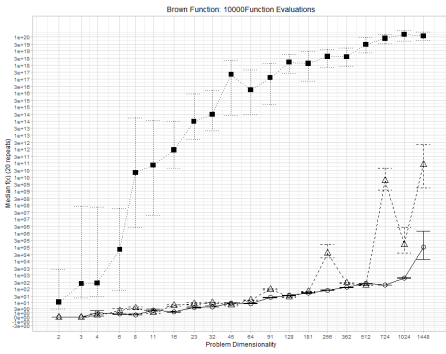
Alpine no.1 Function at 50,000 Evaluations



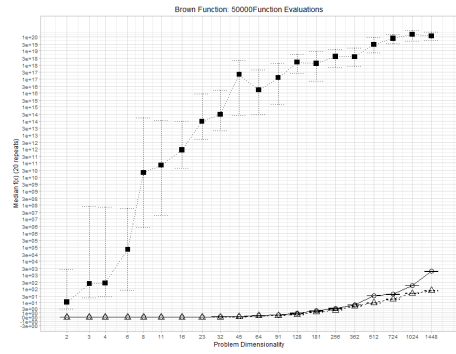
Bent Cigar Function at 10,000 Evaluations



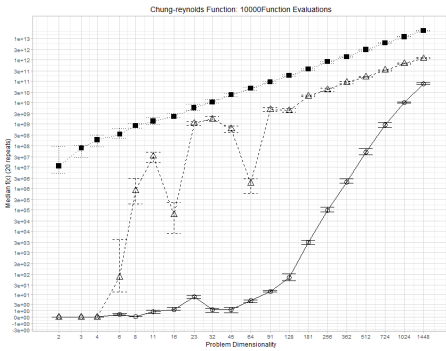
Bent Cigar Function at 50,000 Evaluations



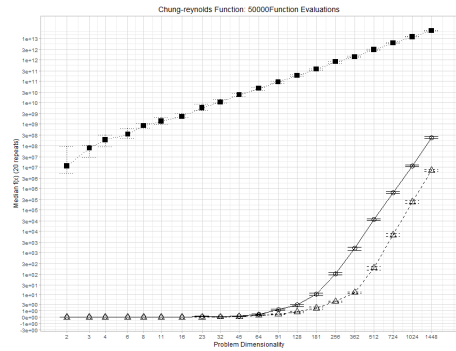
Brown Function at 10,000 Evaluations



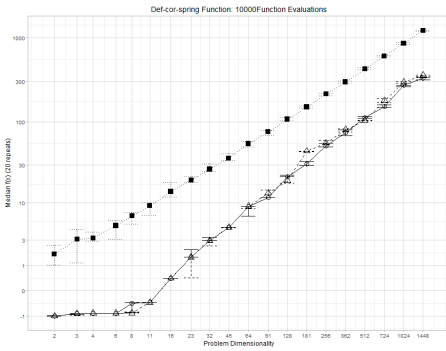
Brown Function at 50,000 Evaluations



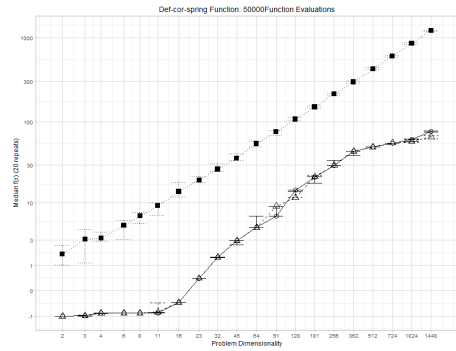
Chung-Reynolds Function at 10,000 Evaluations



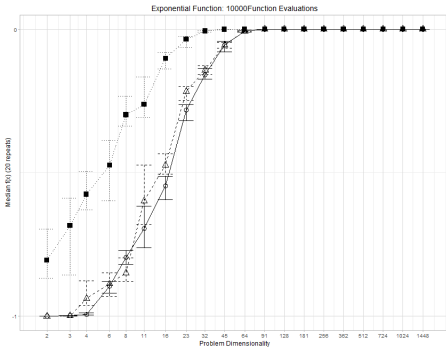
Chung-Reynolds Function at 50,000 Evaluations



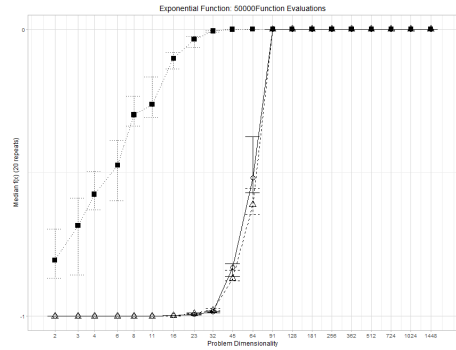
Deflected Corrugated Spring Function at 10,000 Evaluations



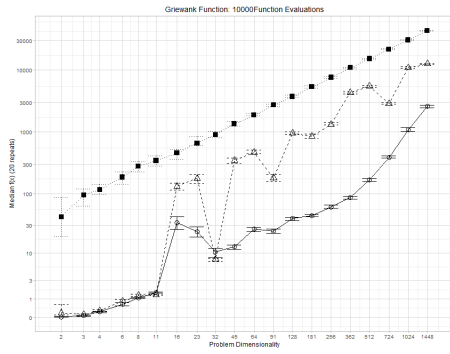
Deflected Corrugated Spring Function at 50,000 Evaluations



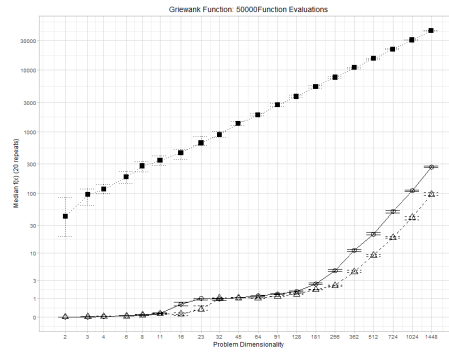
Exponential Function at 10,000 Evaluations



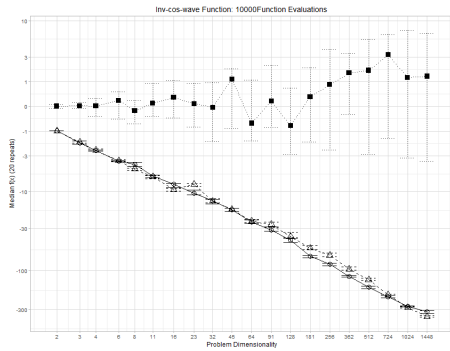
Exponential Function at 50,000 Evaluations



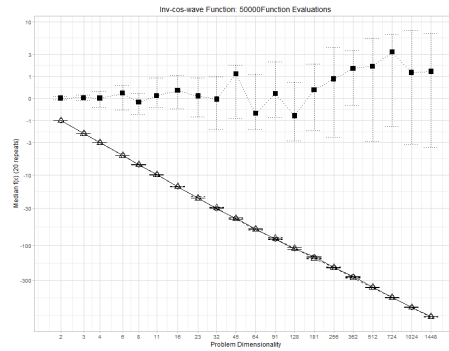
Griewank Function at 10,000 Evaluations



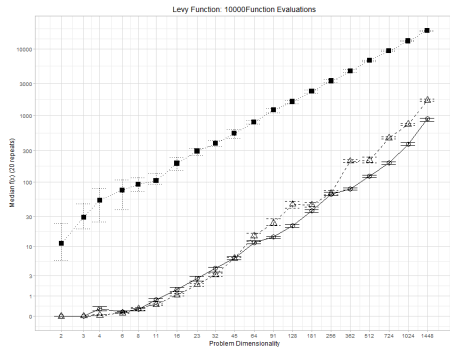
Griewank Function at 50,000 Evaluations



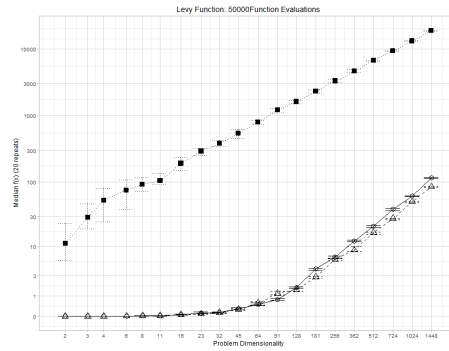
Inverted Cosine Wave Function at 10,000 Evaluations



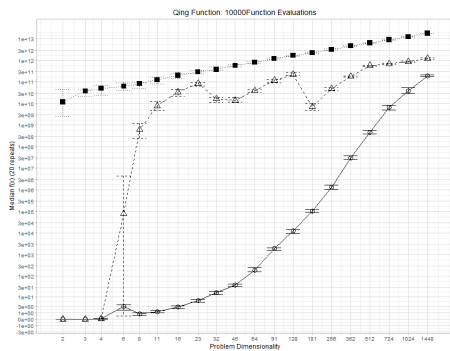
Inverted Cosine Wave Function at 50,000 Evaluations



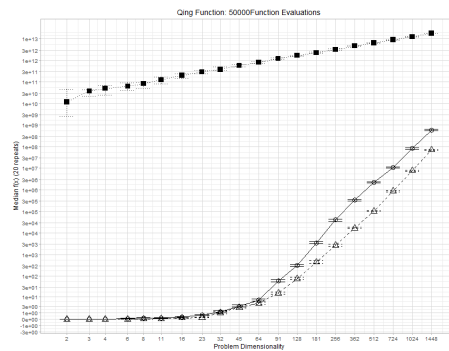
Levy Function at 10,000 Evaluations



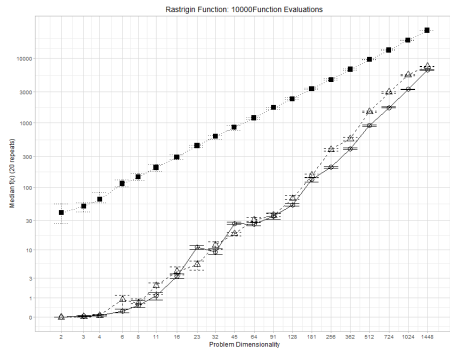
Levy Function at 50,000 Evaluations



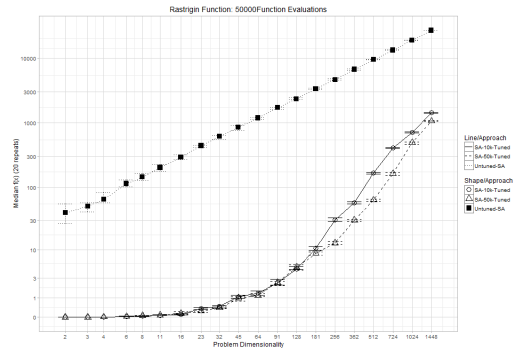
Qing Function at 10,000 Evaluations



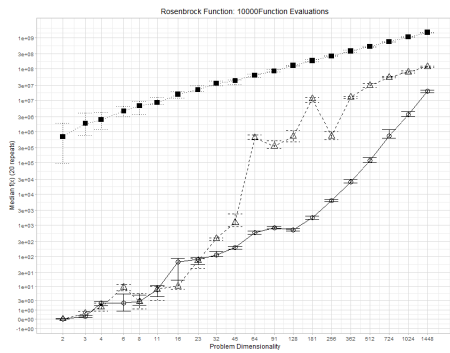
Qing Function at 50,000 Evaluations



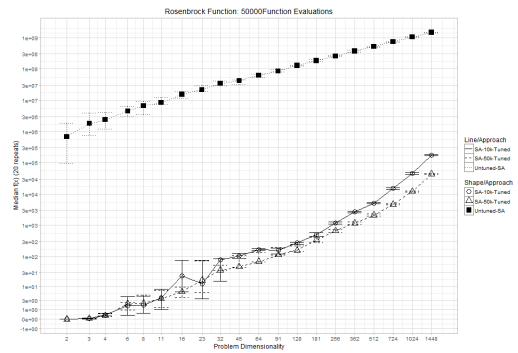
Rastrigin Function at 10,000 Evaluations



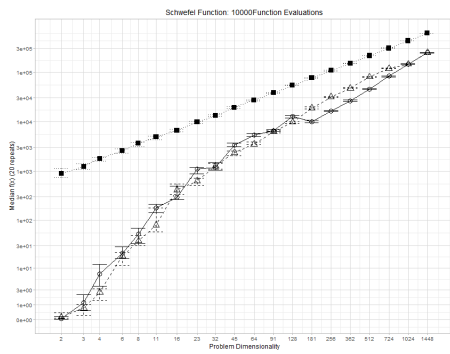
Rastrigin Function at 50,000 Evaluations



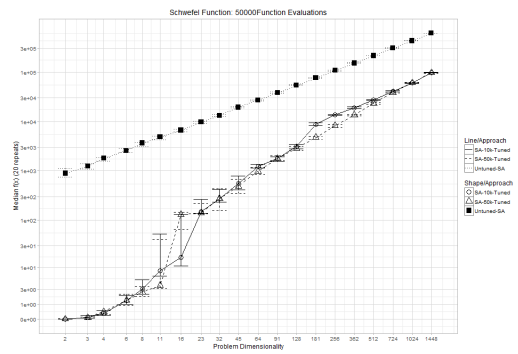
Rosenbrock Function at 10,000 Evaluations



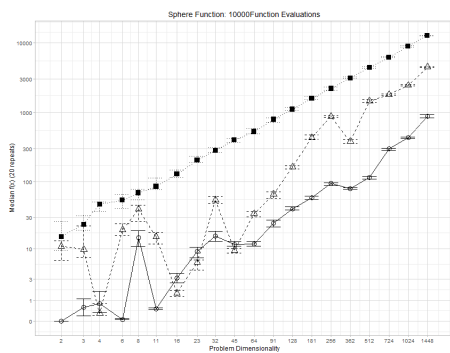
Rosenbrock Function at 50,000 Evaluations



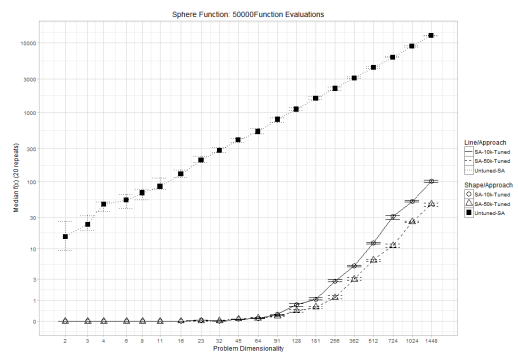
Schwefel Function at 10,000 Evaluations



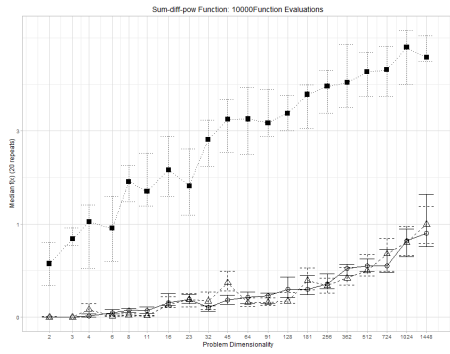
Schwefel Function at 50,000 Evaluations



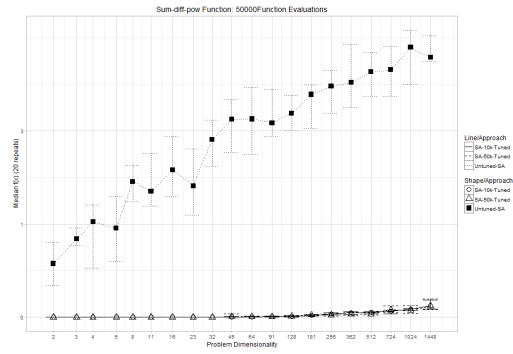
Sphere Function at 10,000 Evaluations



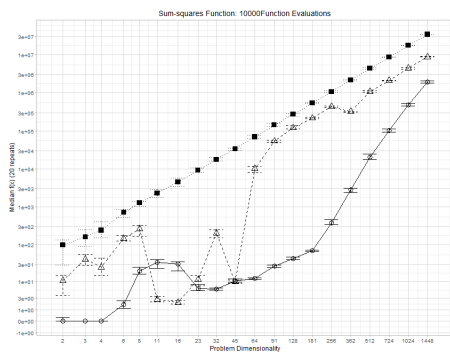
Sphere Function at 50,000 Evaluations



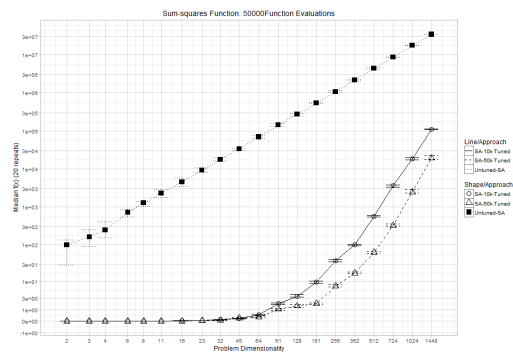
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

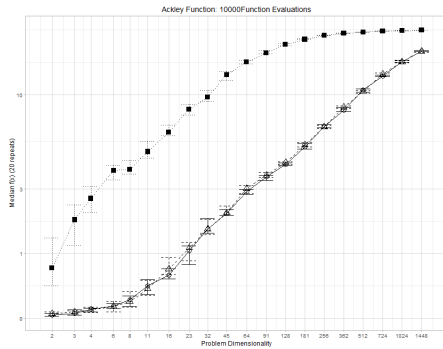


Sum Squares Function at 10,000 Evaluations

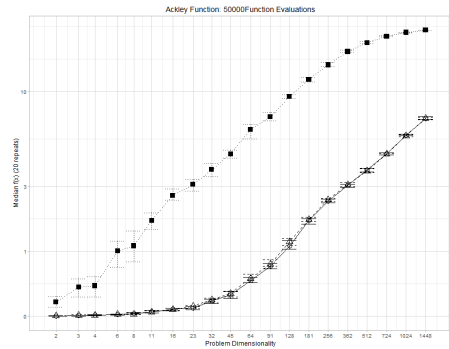


Sum Squares Function at 50,000 Evaluations

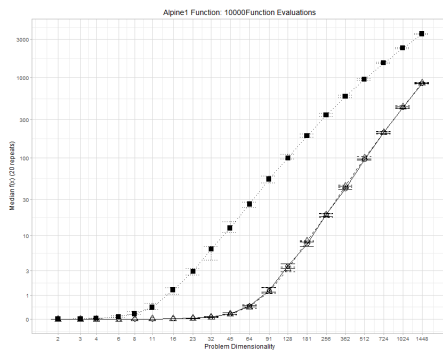
Steady State Genetic Algorithm (SSGA)



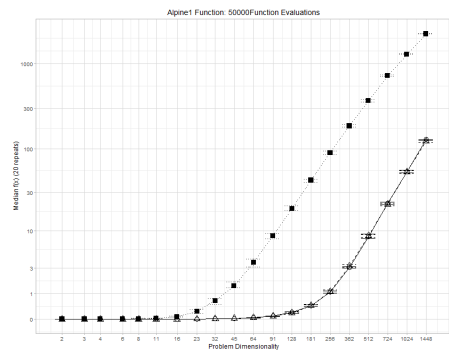
Ackley Function at 10,000 Evaluations



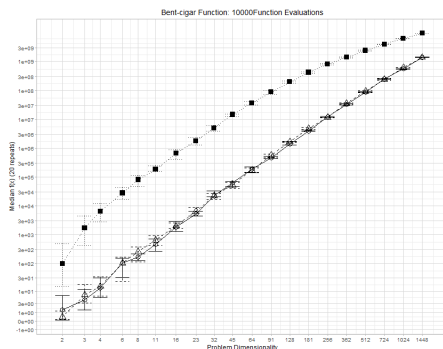
Ackley Function at 50,000 Evaluations



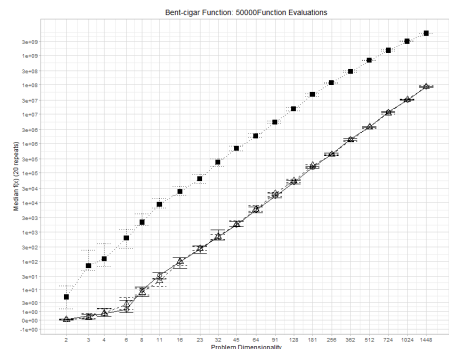
Alpine no.1 Function at 10,000 Evaluations



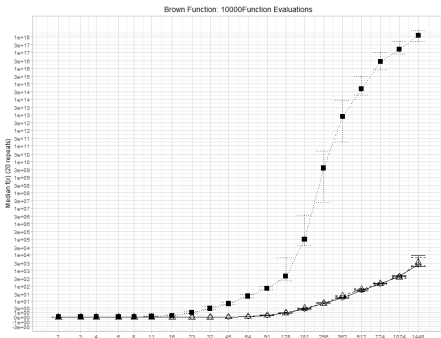
Alpine no.1 Function at 50,000 Evaluations



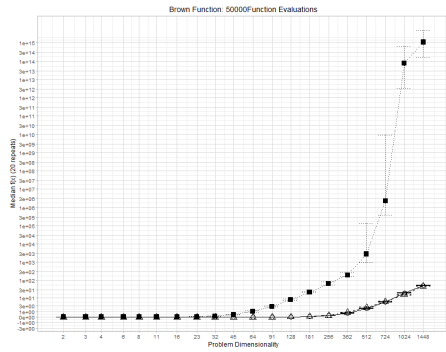
Bent Cigar Function at 10,000 Evaluations



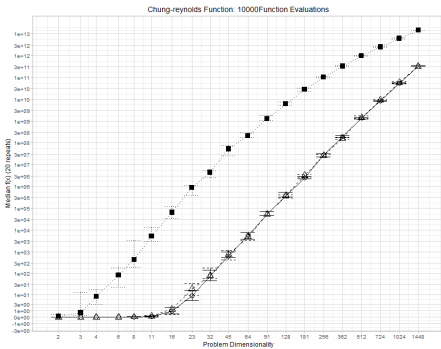
Bent Cigar Function at 50,000 Evaluations



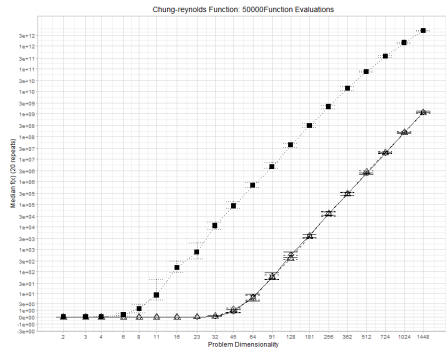
Brown Function at 10,000 Evaluations



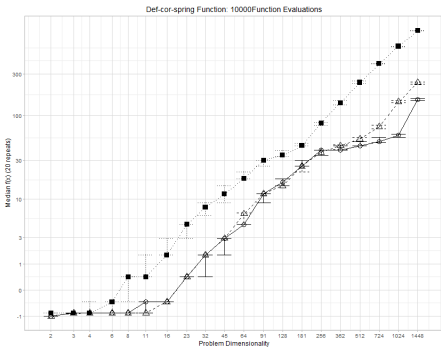
Brown Function at 50,000 Evaluations



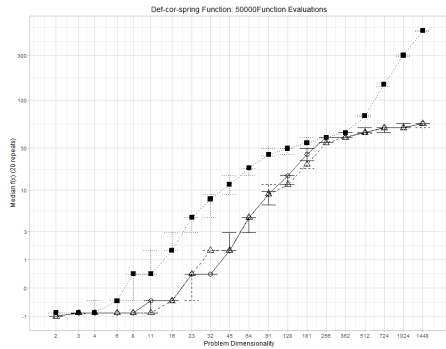
Chung-Reynolds Function at 10,000 Evaluations



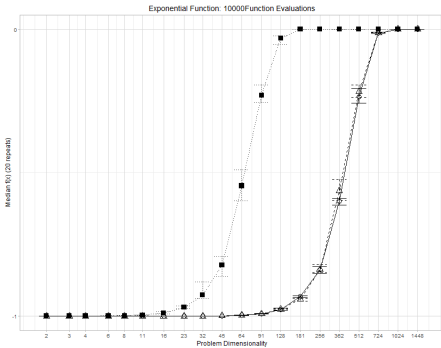
Chung-Reynolds Function at 50,000 Evaluations



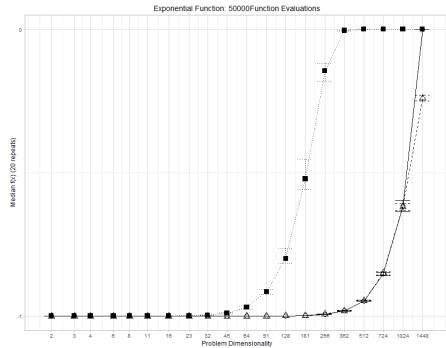
Deflected Corrugated Spring Function at 10,000 Evaluations



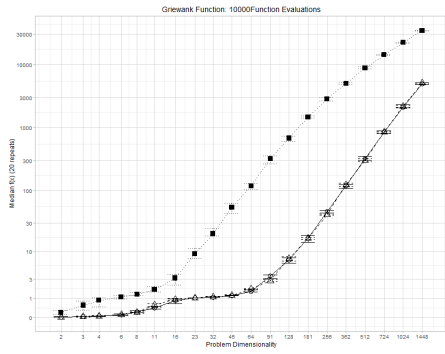
Deflected Corrugated Spring Function at 50,000 Evaluations



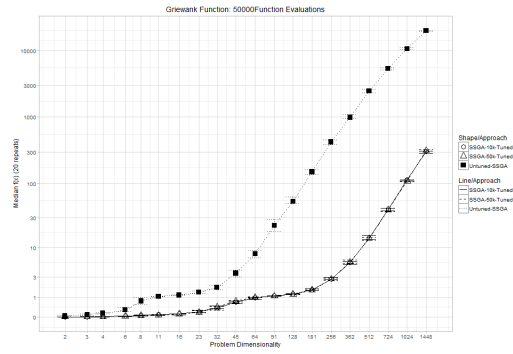
Exponential Function at 10,000 Evaluations



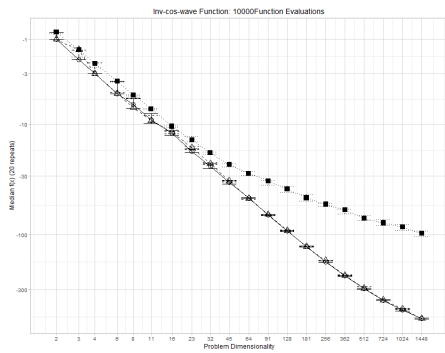
Exponential Function at 50,000 Evaluations



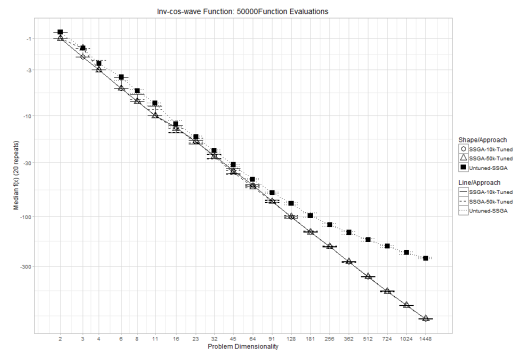
Griewank Function at 10,000 Evaluations



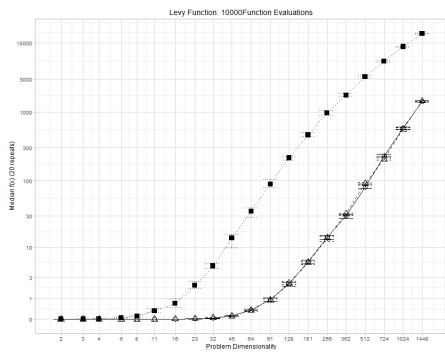
Griewank Function at 50,000 Evaluations



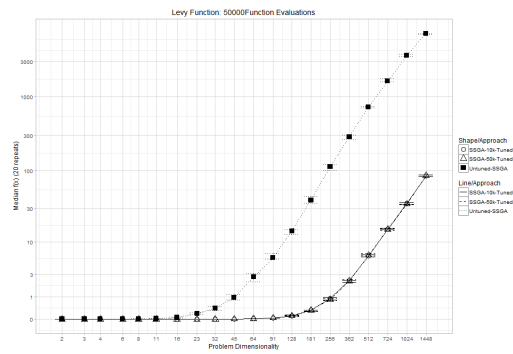
Inverted Cosine Wave Function at 10,000 Evaluations



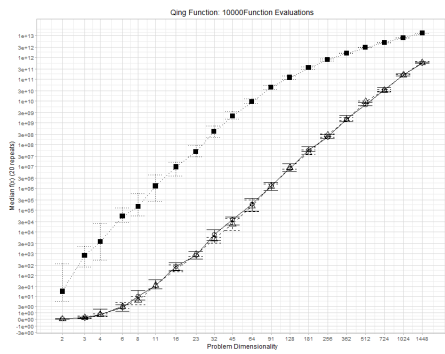
Inverted Cosine Wave Function at 50,000 Evaluations



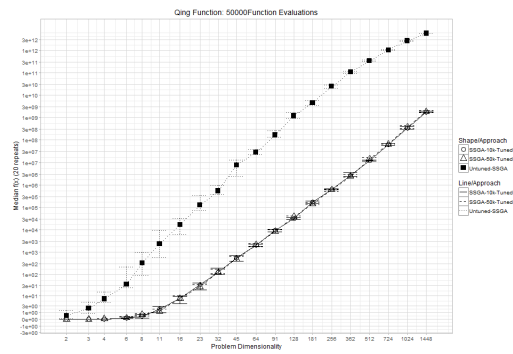
Levy Function at 10,000 Evaluations



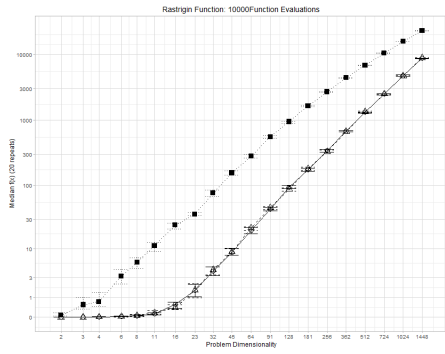
Levy Function at 50,000 Evaluations



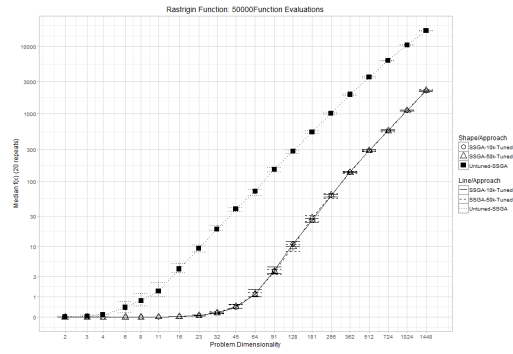
Qing Function at 10,000 Evaluations



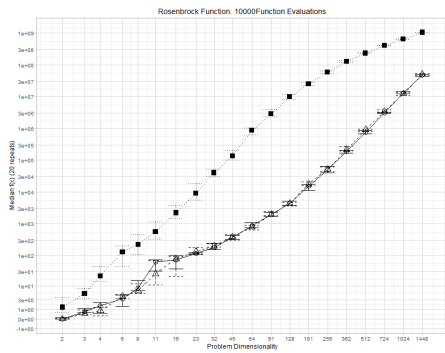
Qing Function at 50,000 Evaluations



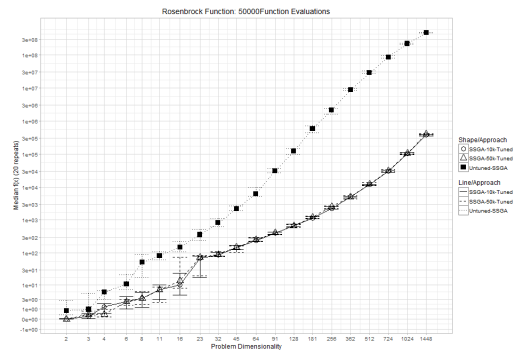
Rastrigin Function at 10,000 Evaluations



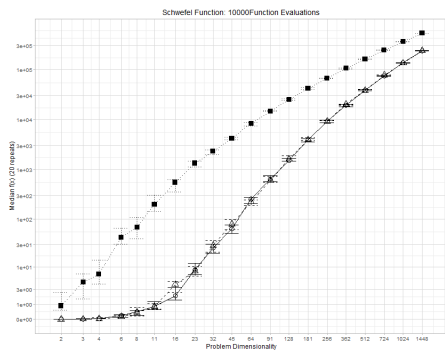
Rastrigin Function at 50,000 Evaluations



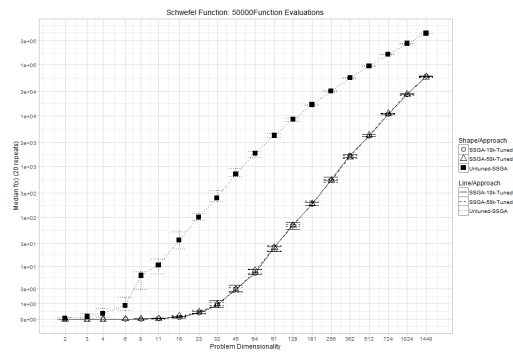
Rosenbrock Function at 10,000 Evaluations



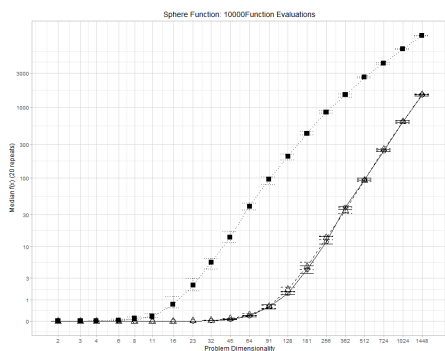
Rosenbrock Function at 50,000 Evaluations



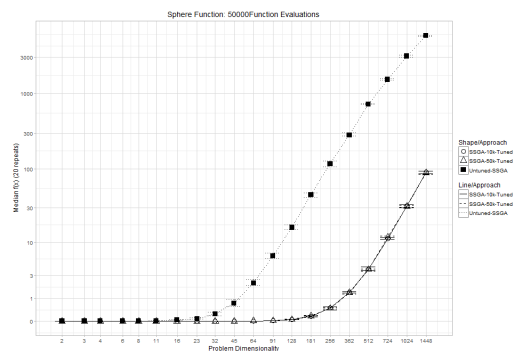
Schwefel Function at 10,000 Evaluations



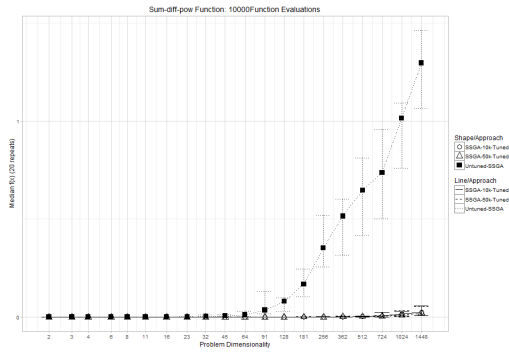
Schwefel Function at 50,000 Evaluations



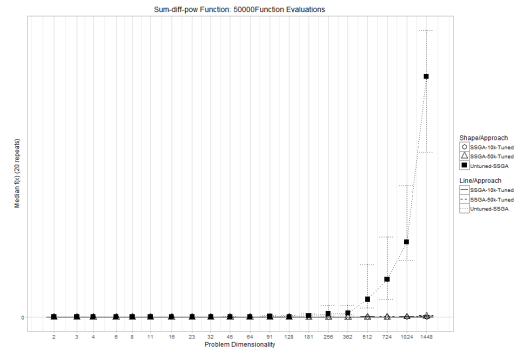
Sphere Function at 10,000 Evaluations



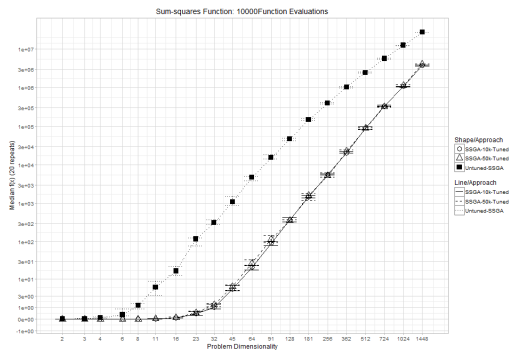
Sphere Function at 50,000 Evaluations



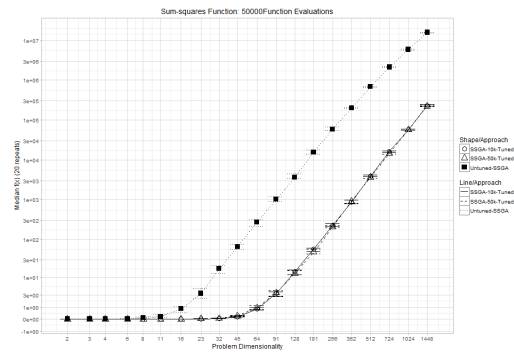
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

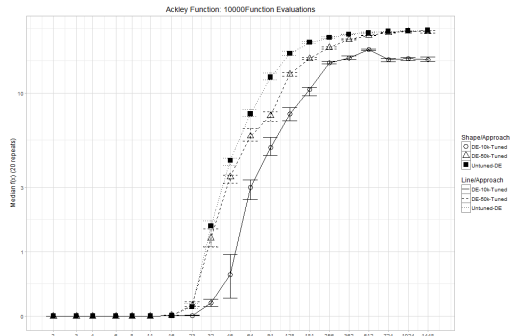


Sum Squares Function at 10,000 Evaluations

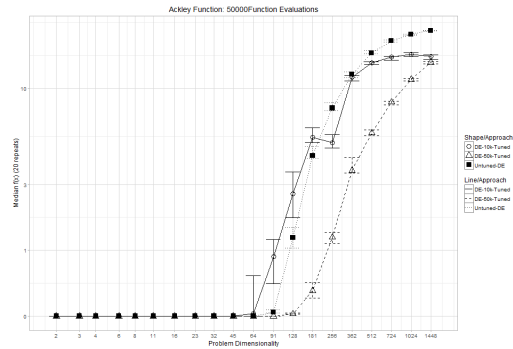


Sum Squares Function at 50,000 Evaluations

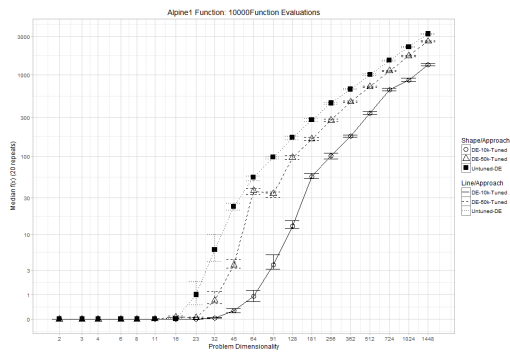
Differential Evolution (DE)



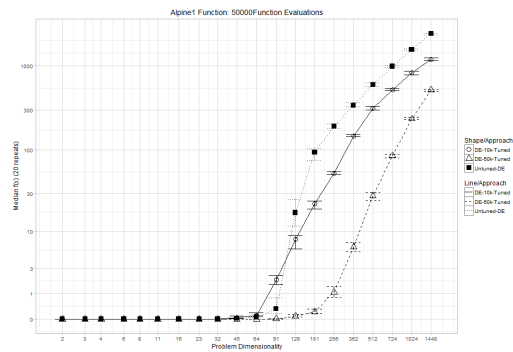
Ackley Function at 10,000 Evaluations



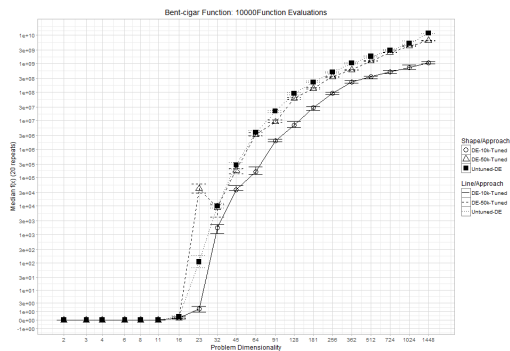
Ackley Function at 50,000 Evaluations



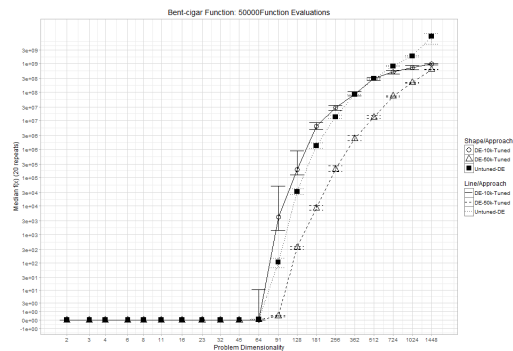
Alpine no.1 Function at 10,000 Evaluations



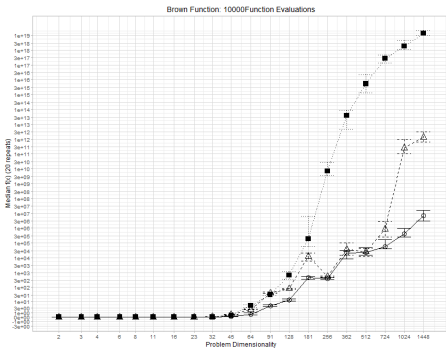
Alpine no.1 Function at 50,000 Evaluations



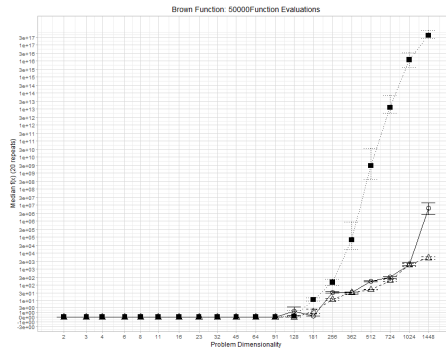
Bent Cigar Function at 10,000 Evaluations



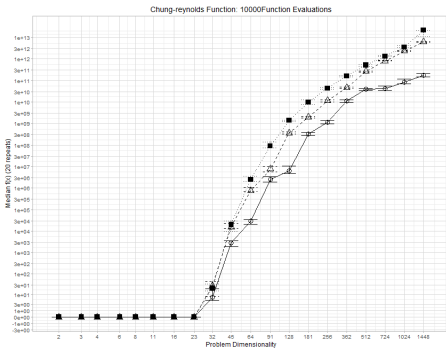
Bent Cigar Function at 50,000 Evaluations



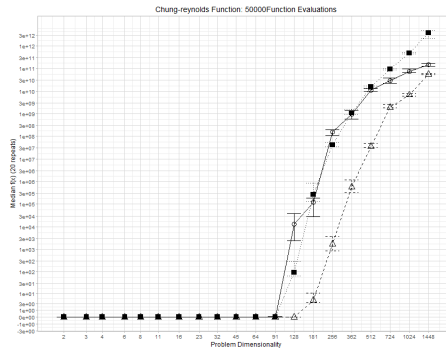
Brown Function at 10,000 Evaluations



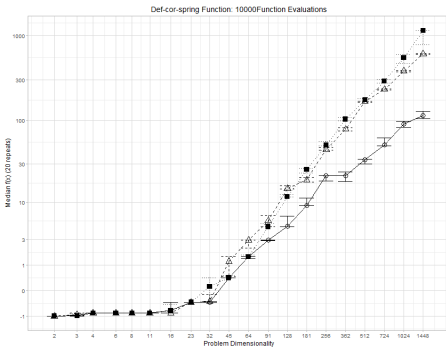
Brown Function at 50,000 Evaluations



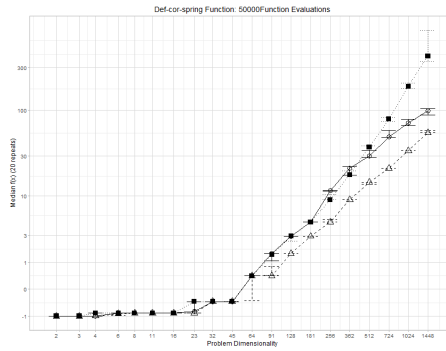
Chung-Reynolds Function at 10,000 Evaluations



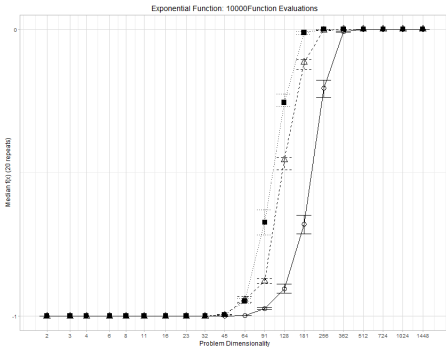
Chung-Reynolds Function at 50,000 Evaluations



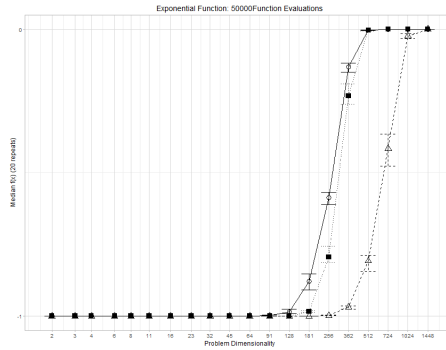
Deflected Corrugated Spring Function at 10,000 Evaluations



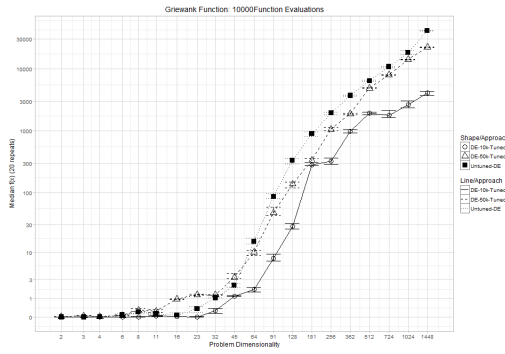
Deflected Corrugated Spring Function at 50,000 Evaluations



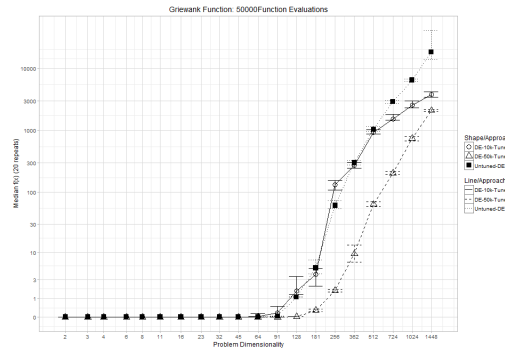
Exponential Function at 10,000 Evaluations



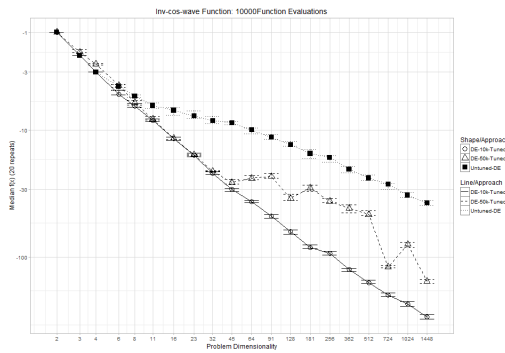
Exponential Function at 50,000 Evaluations



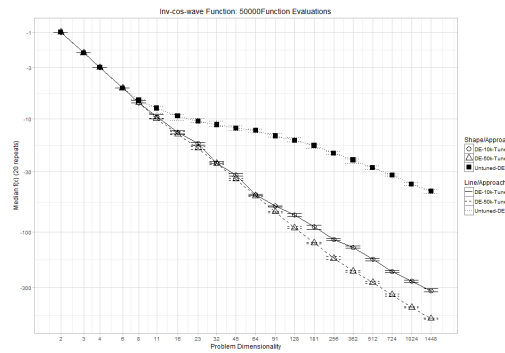
Griewank Function at 10,000 Evaluations



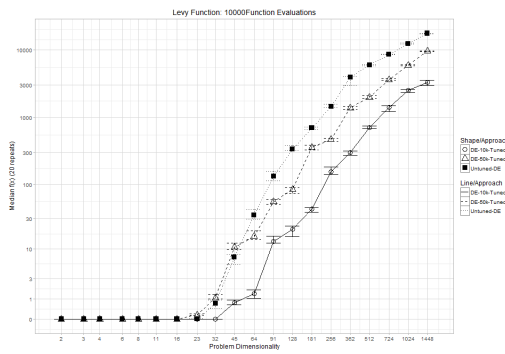
Griewank Function at 50,000 Evaluations



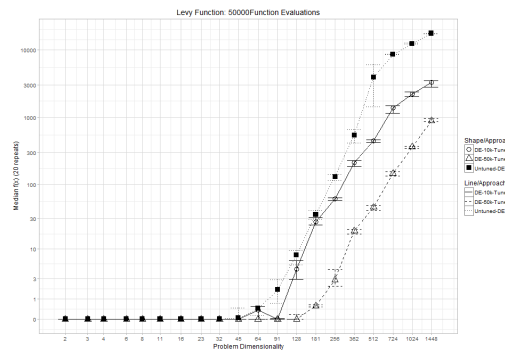
Inverted Cosine Wave Function at 10,000 Evaluations



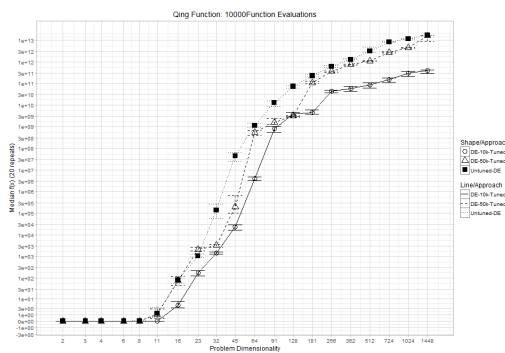
Inverted Cosine Wave Function at 50,000 Evaluations



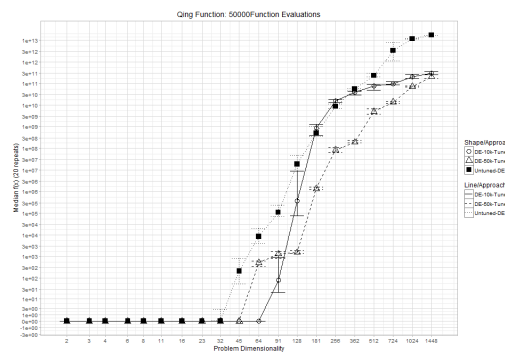
Levy Function at 10,000 Evaluations



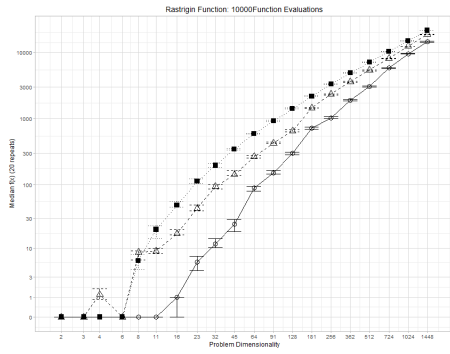
Levy Function at 50,000 Evaluations



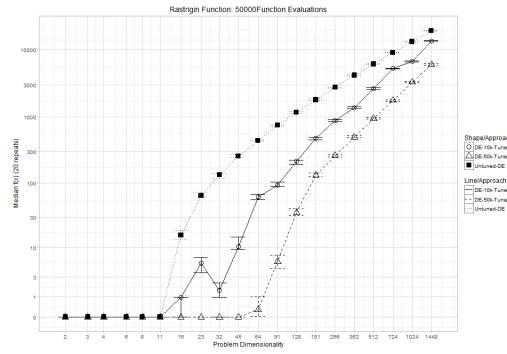
Qing Function at 10,000 Evaluations



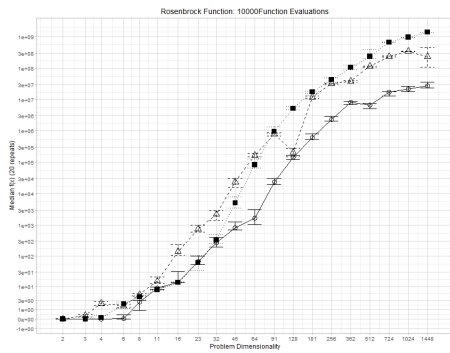
Qing Function at 50,000 Evaluations



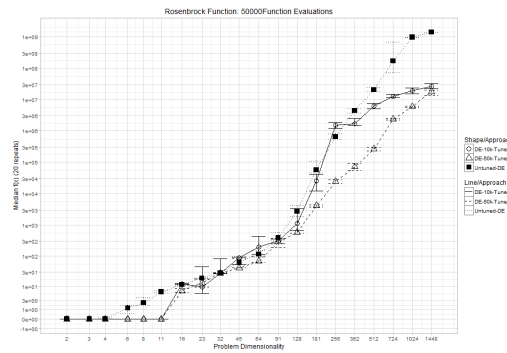
Rastrigin Function at 10,000 Evaluations



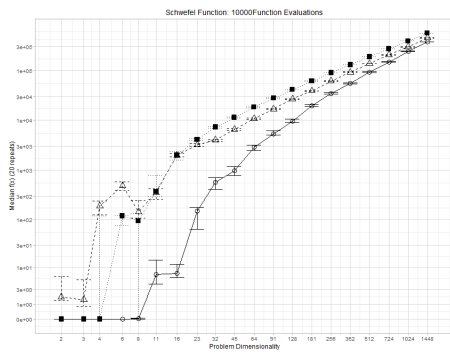
Rastrigin Function at 50,000 Evaluations



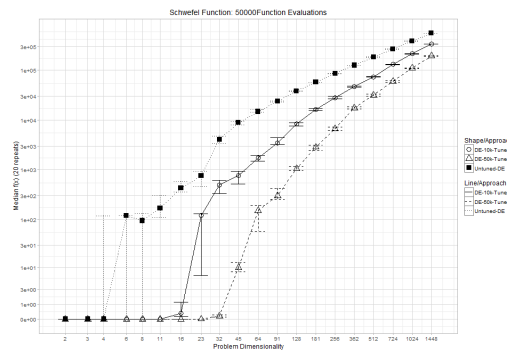
Rosenbrock Function at 10,000 Evaluations



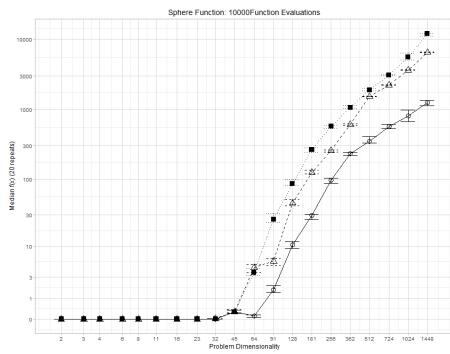
Rosenbrock Function at 50,000 Evaluations



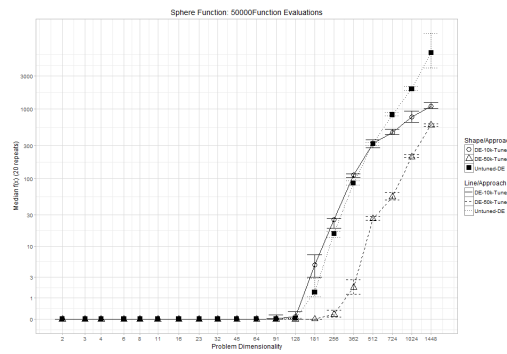
Schwefel Function at 10,000 Evaluations



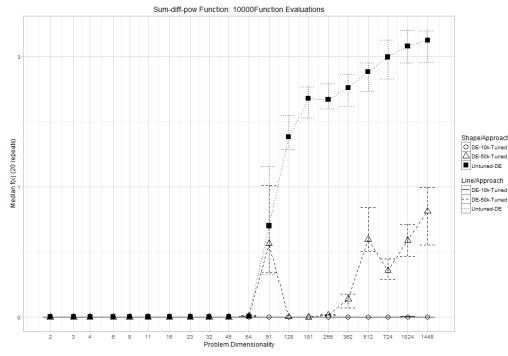
Schwefel Function at 50,000 Evaluations



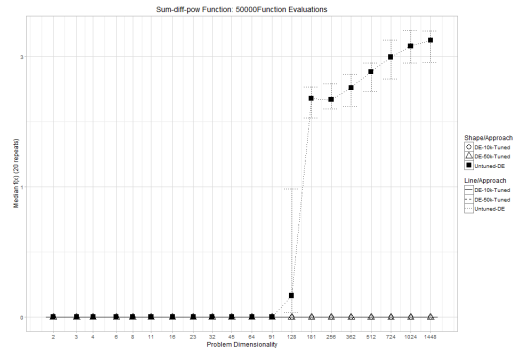
Sphere Function at 10,000 Evaluations



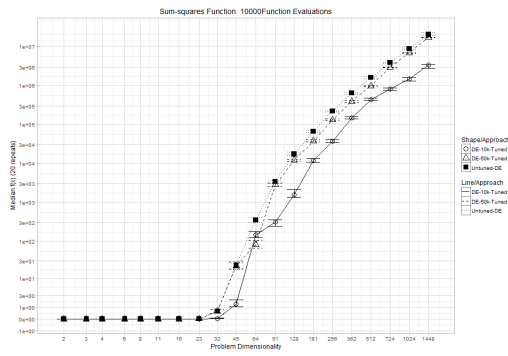
Sphere Function at 50,000 Evaluations



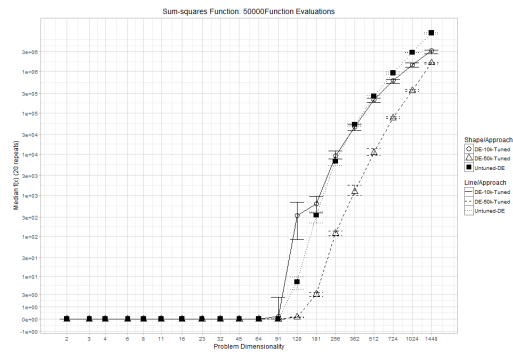
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

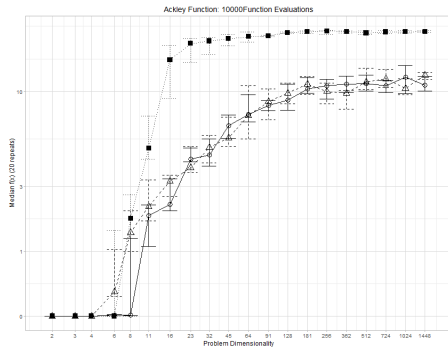


Sum Squares Function at 10,000 Evaluations

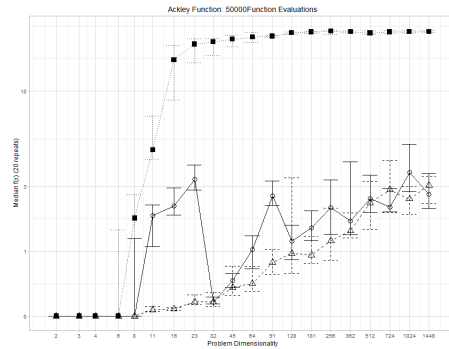


Sum Squares Function at 50,000 Evaluations

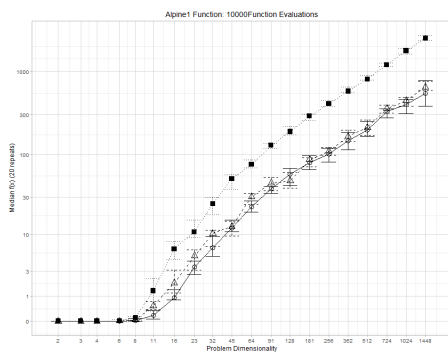
Particle Swarm Optimisation



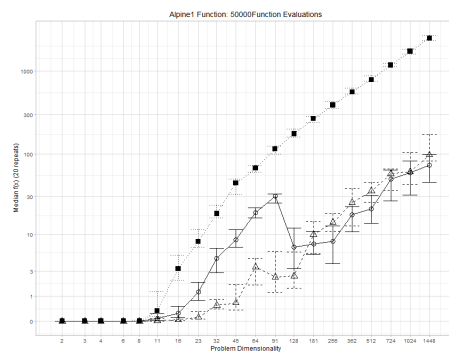
Ackley Function at 10,000 Evaluations



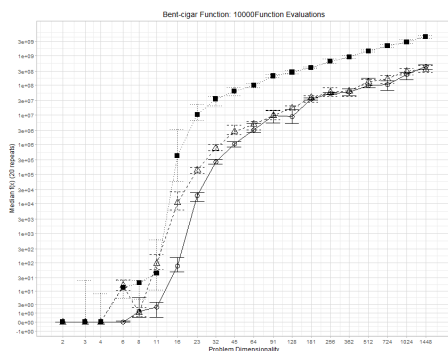
Ackley Function at 50,000 Evaluations



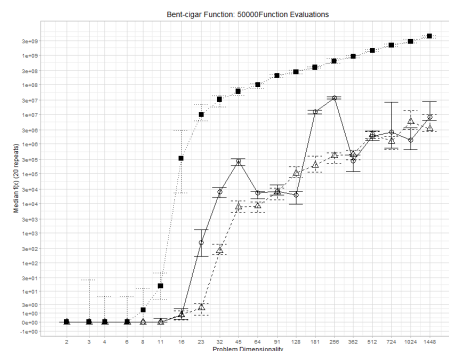
Alpine no.1 Function at 10,000 Evaluations



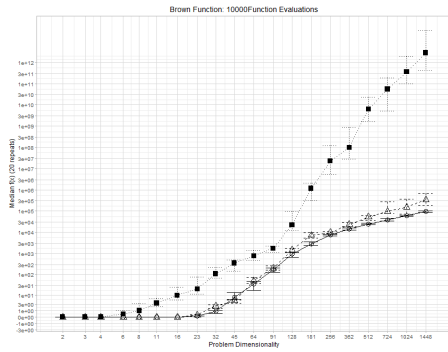
Alpine no.1 Function at 50,000 Evaluations



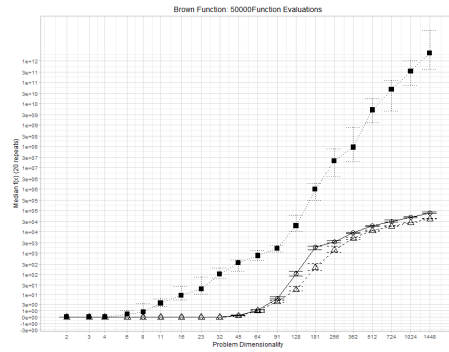
Bent Cigar Function at 10,000 Evaluations



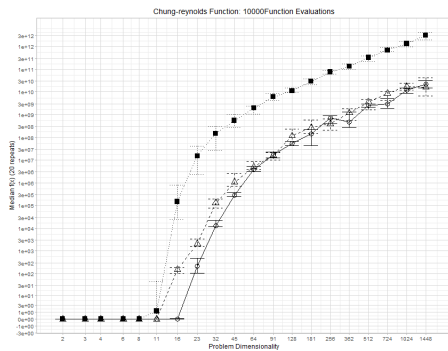
Bent Cigar Function at 50,000 Evaluations



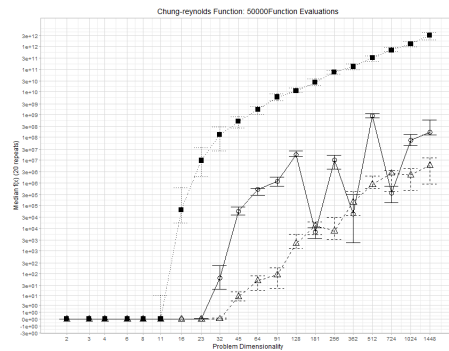
Brown Function at 10,000 Evaluations



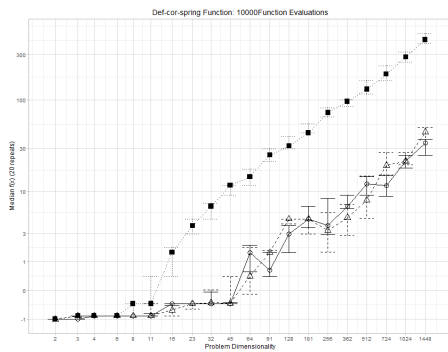
Brown Function at 50,000 Evaluations



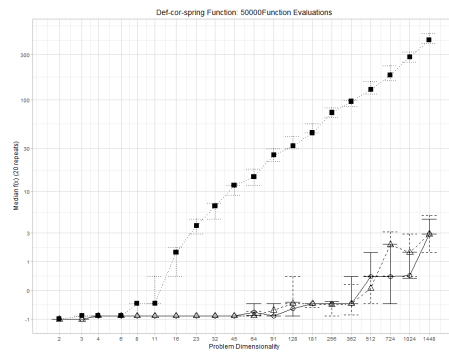
Chung-Reynolds Function at 10,000 Evaluations



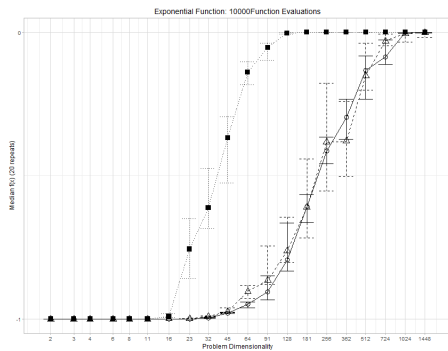
Chung-Reynolds Function at 50,000 Evaluations



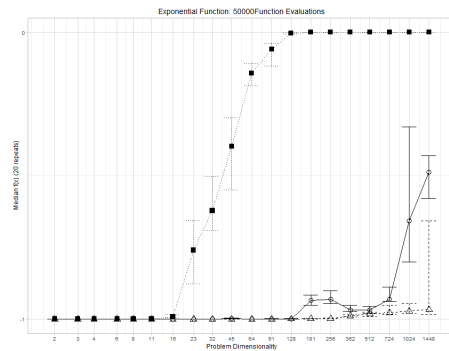
Deflected Corrugated Spring Function at 10,000 Evaluations



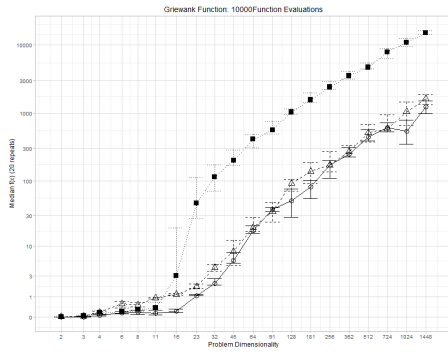
Deflected Corrugated Spring Function at 50,000 Evaluations



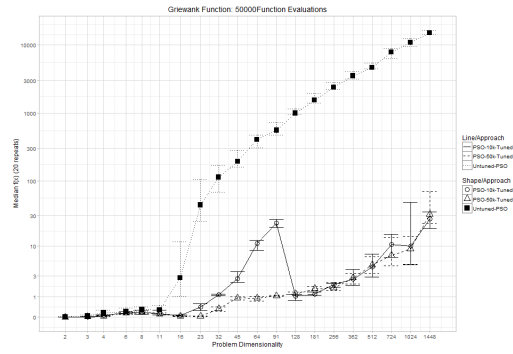
Exponential Function at 10,000 Evaluations



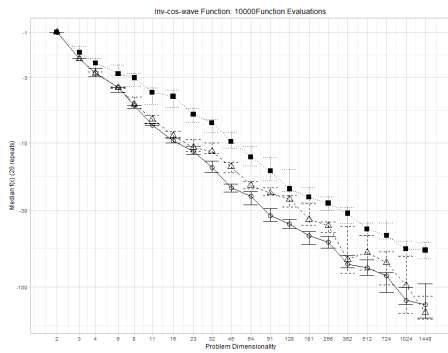
Exponential Function at 50,000 Evaluations



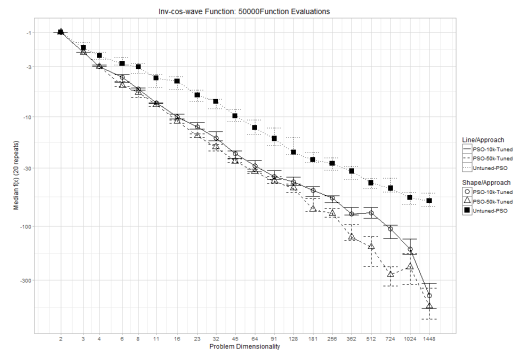
Griewank Function at 10,000 Evaluations



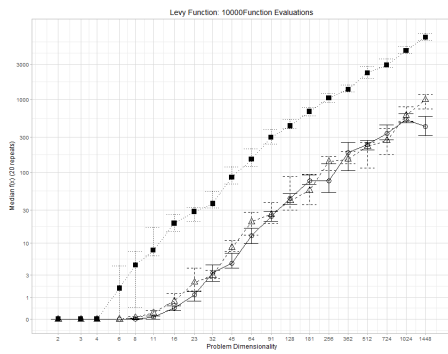
Griewank Function at 50,000 Evaluations



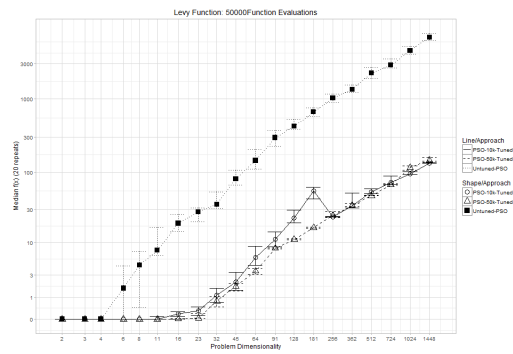
Inverted Cosine Wave Function at 10,000 Evaluations



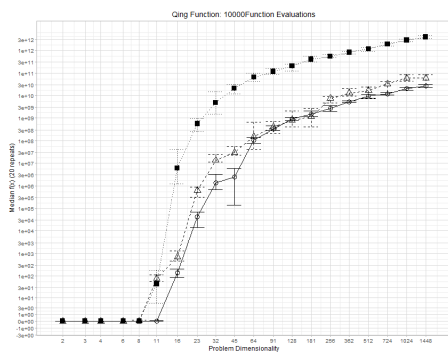
Inverted Cosine Wave Function at 50,000 Evaluations



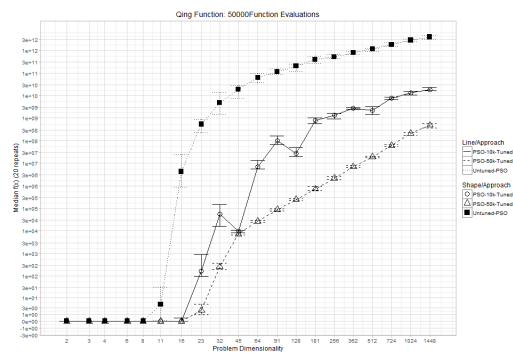
Levy Function at 10,000 Evaluations



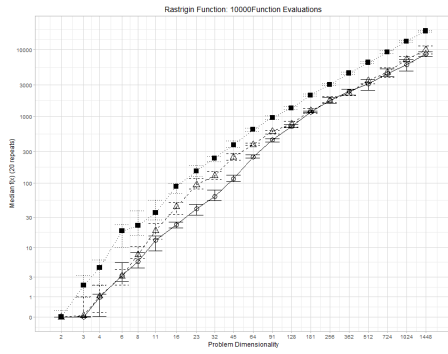
Levy Function at 50,000 Evaluations



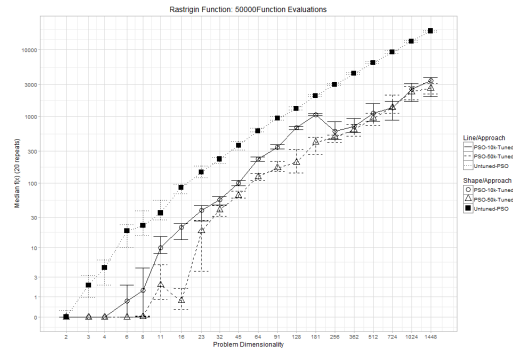
Qing Function at 10,000 Evaluations



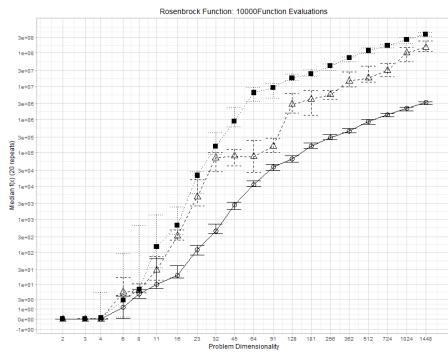
Qing Function at 50,000 Evaluations



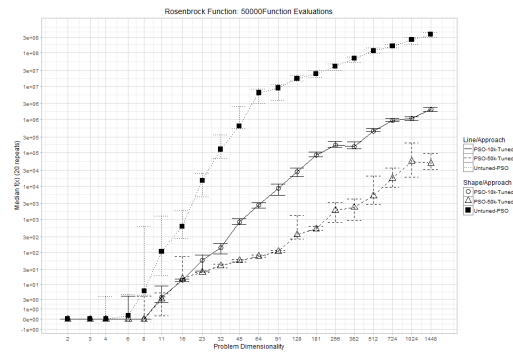
Rastrigin Function at 10,000 Evaluations



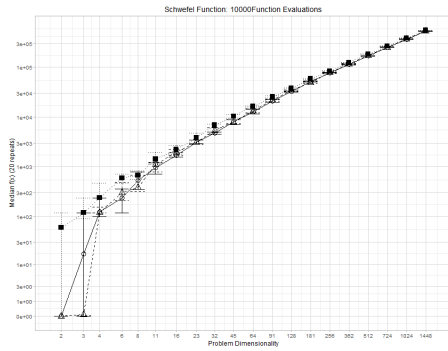
Rastrigin Function at 50,000 Evaluations



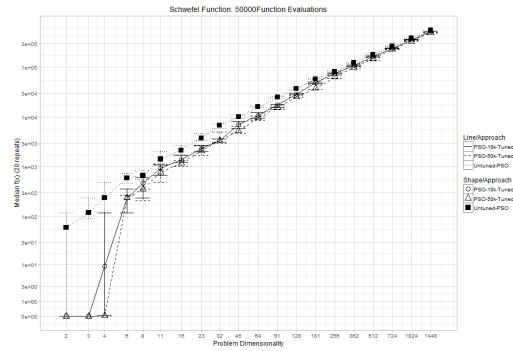
Rosenbrock Function at 10,000 Evaluations



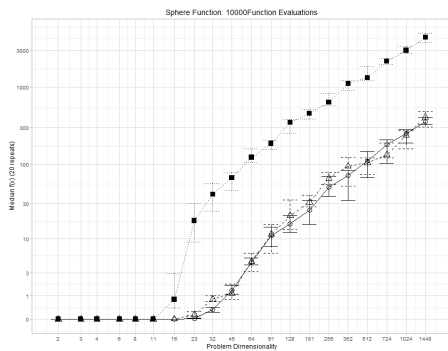
Rosenbrock Function at 50,000 Evaluations



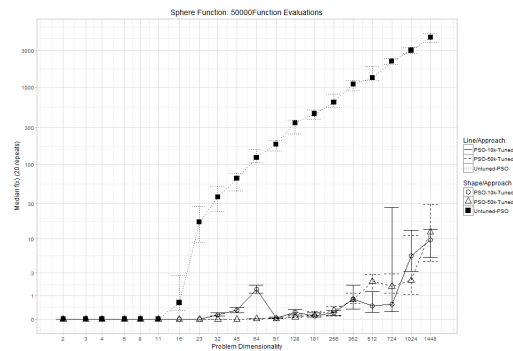
Schwefel Function at 10,000 Evaluations



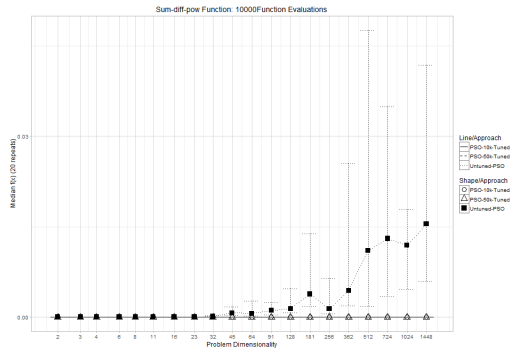
Schwefel Function at 50,000 Evaluations



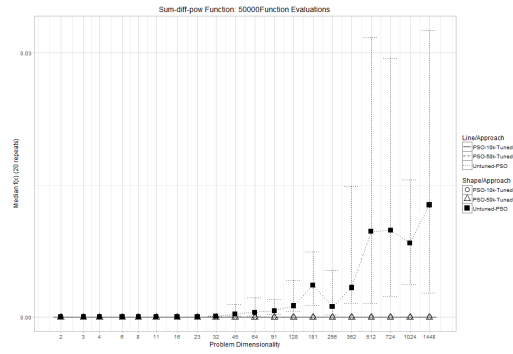
Sphere Function at 10,000 Evaluations



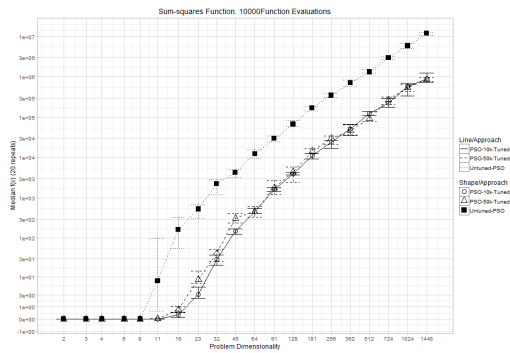
Sphere Function at 50,000 Evaluations



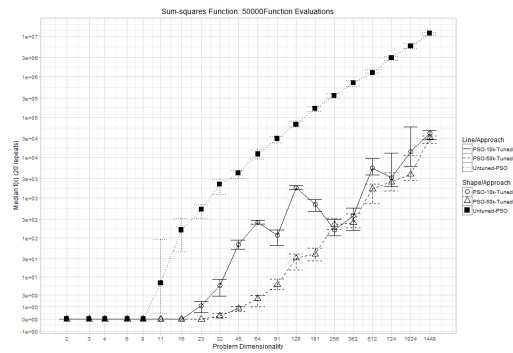
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations

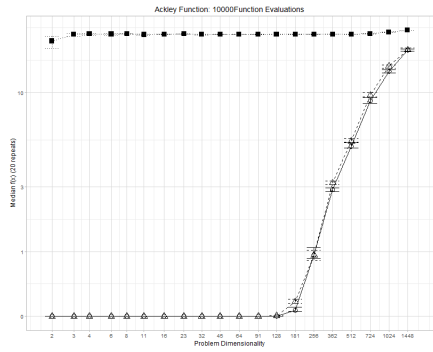


Sum Squares Function at 10,000 Evaluations

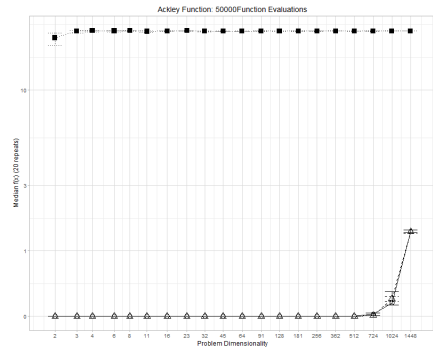


Sum Squares Function at 50,000 Evaluations

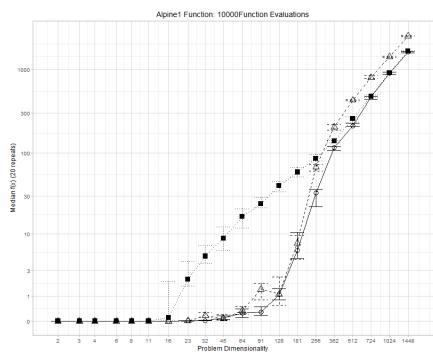
Covariance Matrix Adaption Evolutionary Strategy (CMA-ES)



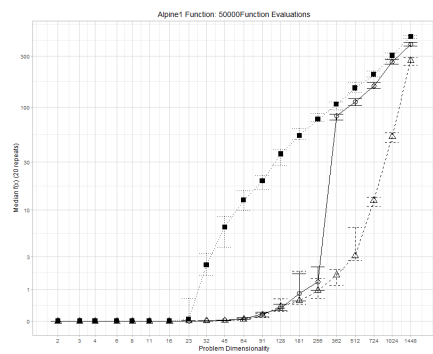
Ackley Function at 10,000 Evaluations



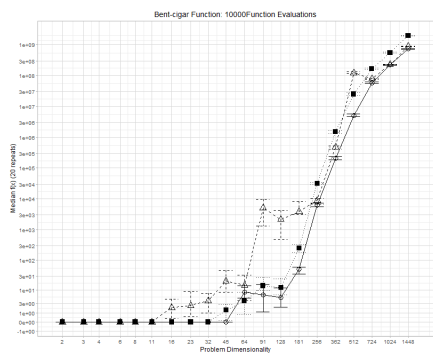
Ackley Function at 50,000 Evaluations



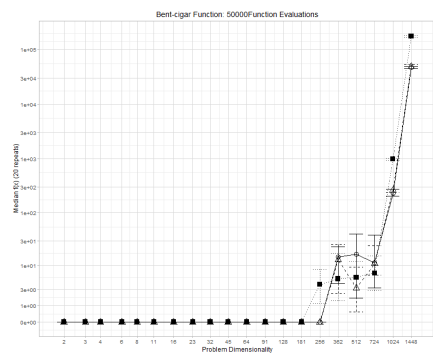
Alpine no.1 Function at 10,000 Evaluations



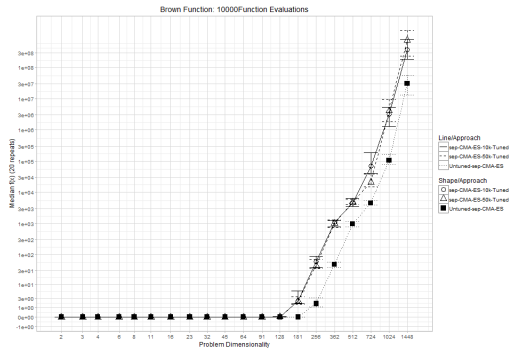
Alpine no.1 Function at 50,000 Evaluations



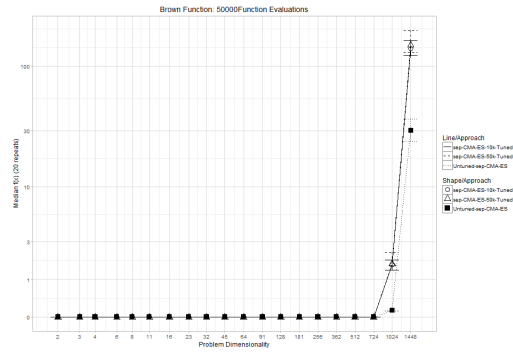
Bent Cigar Function at 10,000 Evaluations



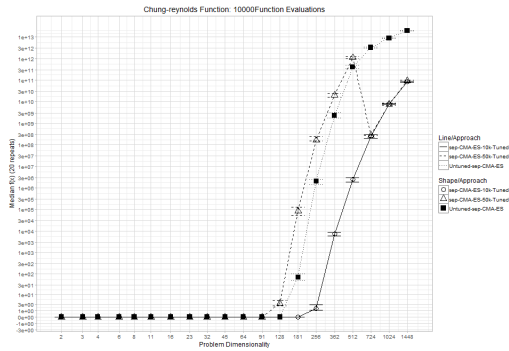
Bent Cigar Function at 50,000 Evaluations



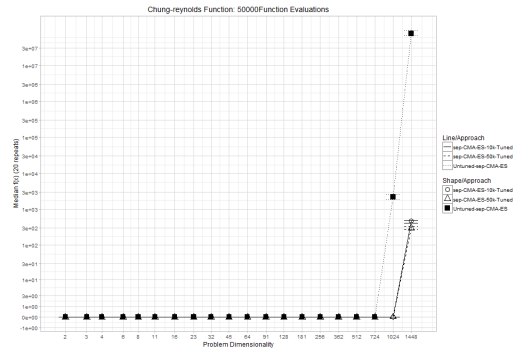
Brown Function at 10,000 Evaluations



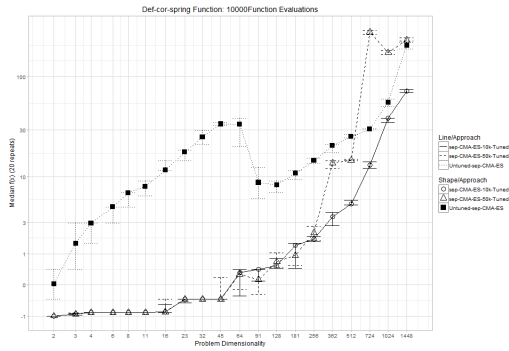
Brown Function at 50,000 Evaluations



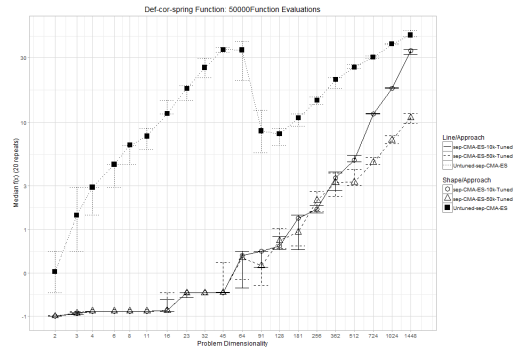
Chung-Reynolds Function at 10,000 Evaluations



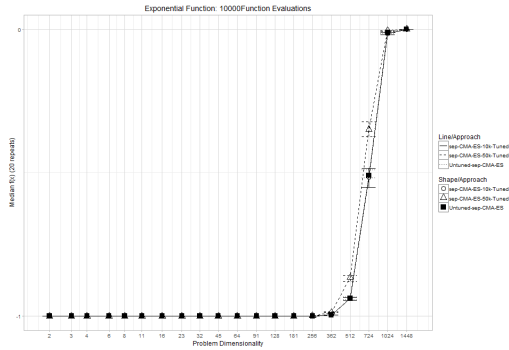
Chung-Reynolds Function at 50,000 Evaluations



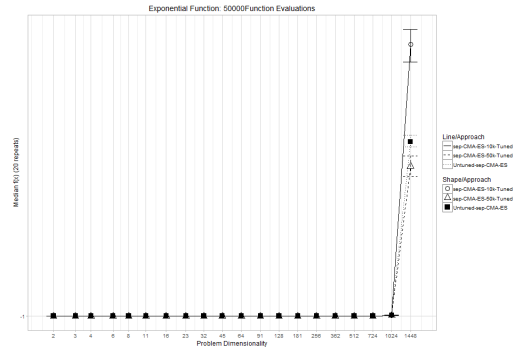
Deflected Corrugated Spring Function at 10,000 Evaluations



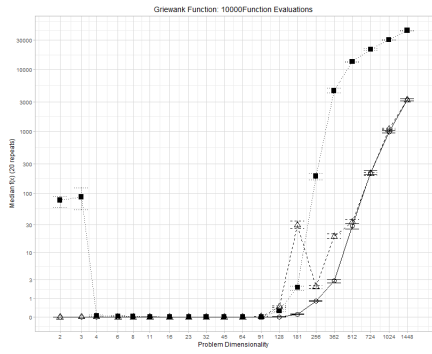
Deflected Corrugated Spring Function at 50,000 Evaluations



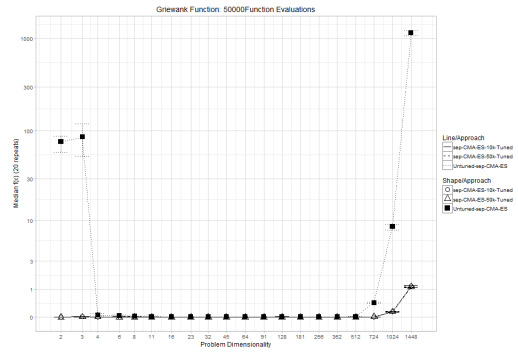
Exponential Function at 10,000 Evaluations



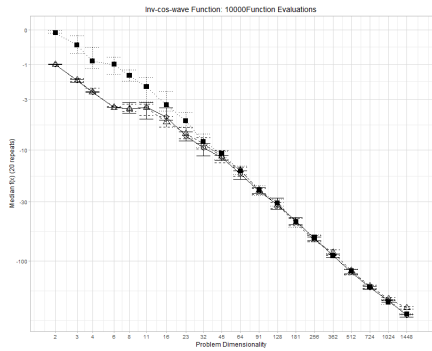
Exponential Function at 50,000 Evaluations



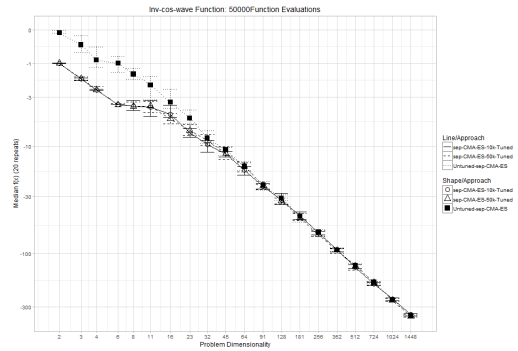
Griewank Function at 10,000 Evaluations



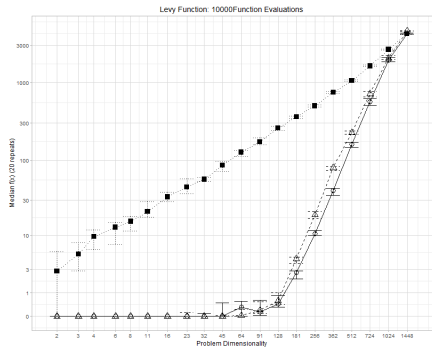
Griewank Function at 50,000 Evaluations



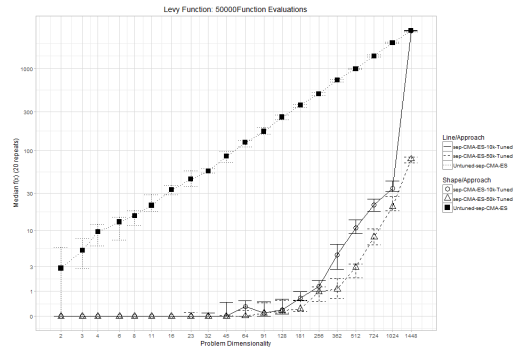
Inverted Cosine Wave Function at 10,000 Evaluations



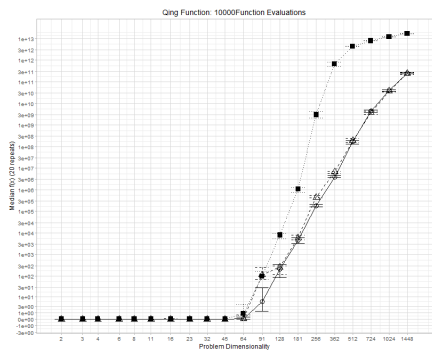
Inverted Cosine Wave Function at 50,000 Evaluations



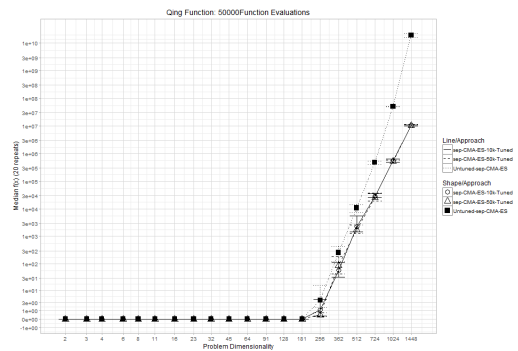
Levy Function at 10,000 Evaluations



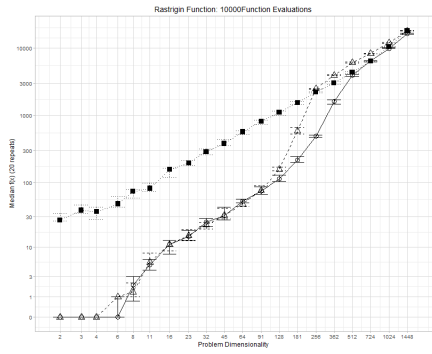
Levy Function at 50,000 Evaluations



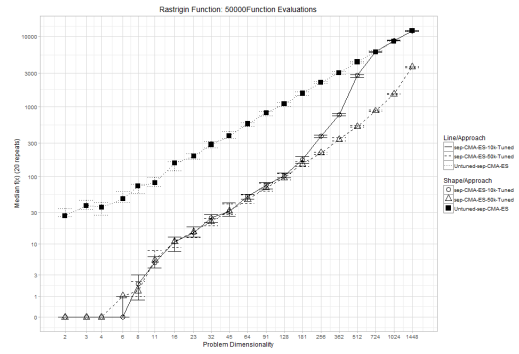
Qing Function at 10,000 Evaluations



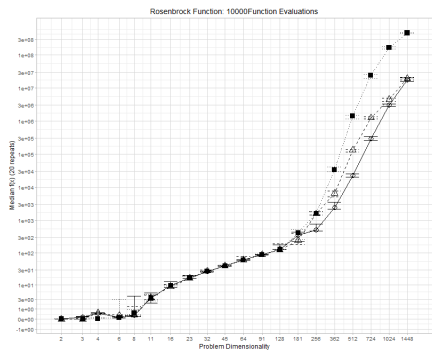
Qing Function at 50,000 Evaluations



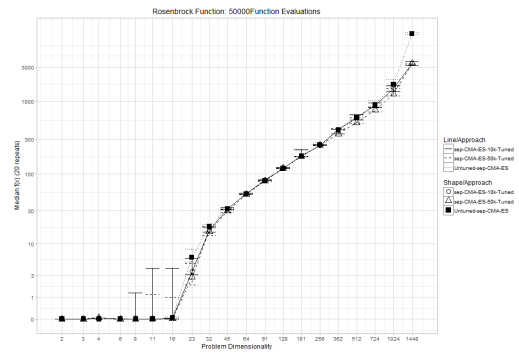
Rastrigin Function at 10,000 Evaluations



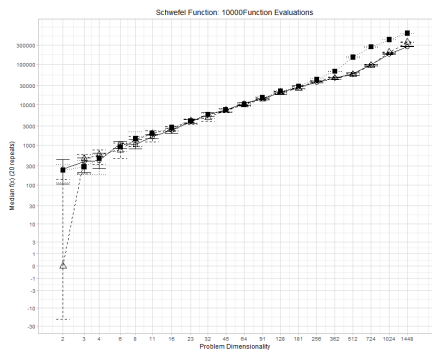
Rastrigin Function at 50,000 Evaluations



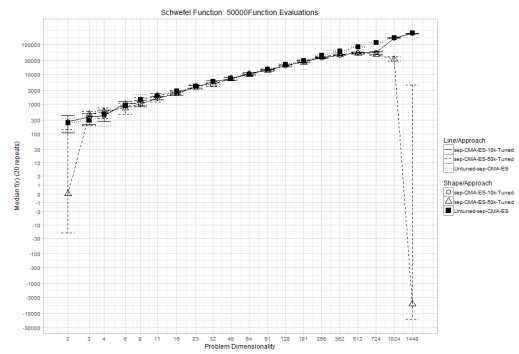
Rosenbrock Function at 10,000 Evaluations



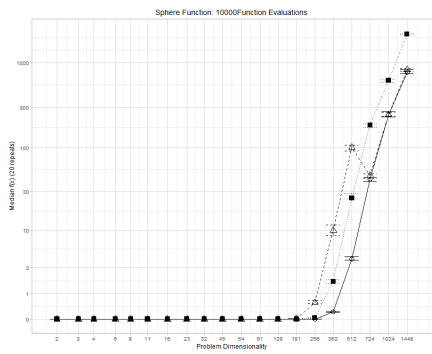
Rosenbrock Function at 50,000 Evaluations



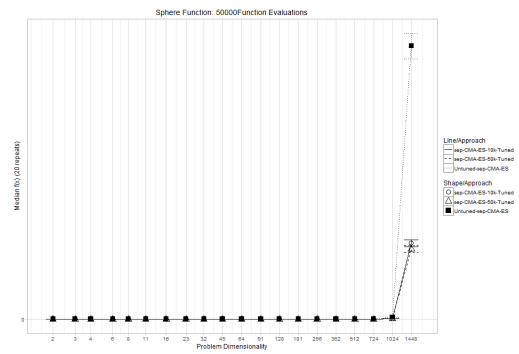
Schwefel Function at 10,000 Evaluations



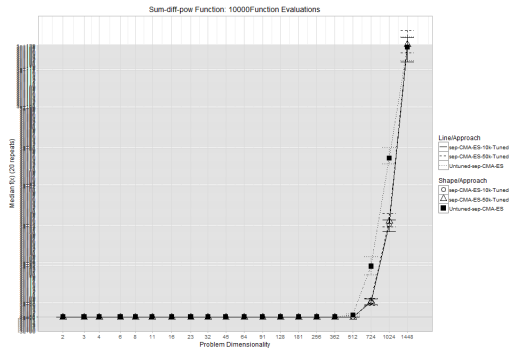
Schwefel Function at 50,000 Evaluations



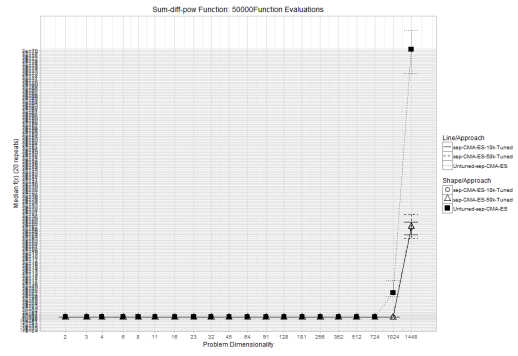
Sphere Function at 10,000 Evaluations



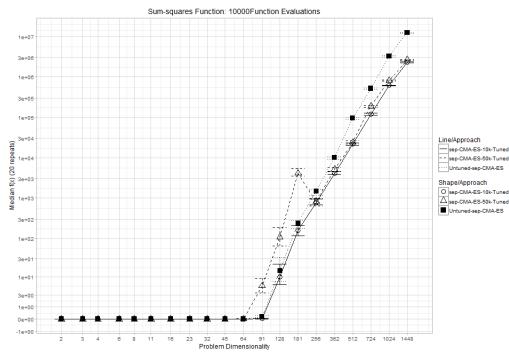
Sphere Function at 50,000 Evaluations



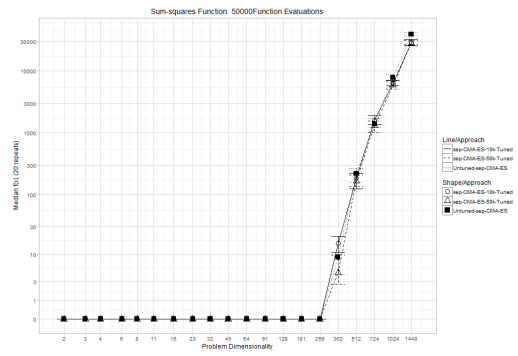
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations



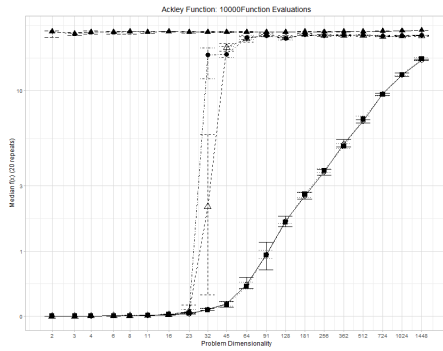
Sum Squares Function at 10,000 Evaluations



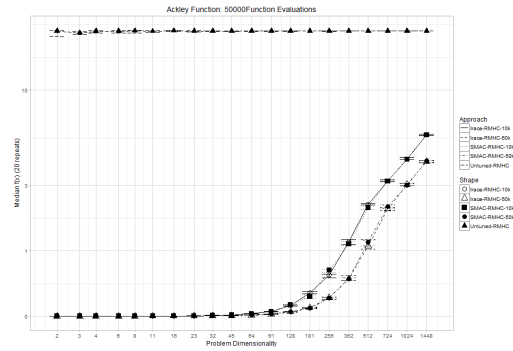
Sum Squares Function at 50,000 Evaluations

D.2 IRACE VS. SMAC PARAMETER TUNING RESULTS

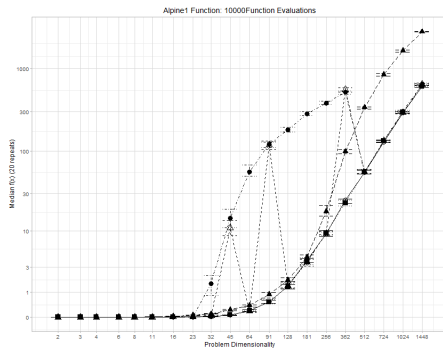
Random Mutation Hill Climbing (RMHC)



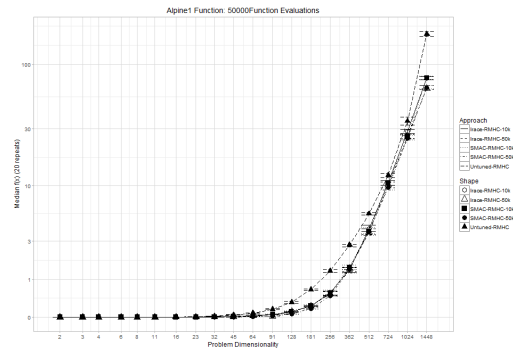
Ackley Function at 10,000 Evaluations



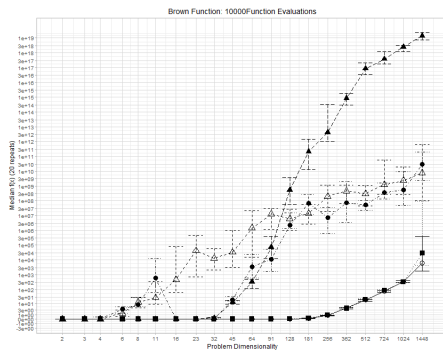
Ackley Function at 50,000 Evaluations



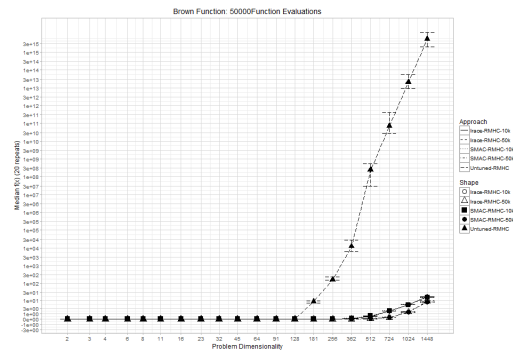
Alpine no.1 Function at 10,000 Evaluations



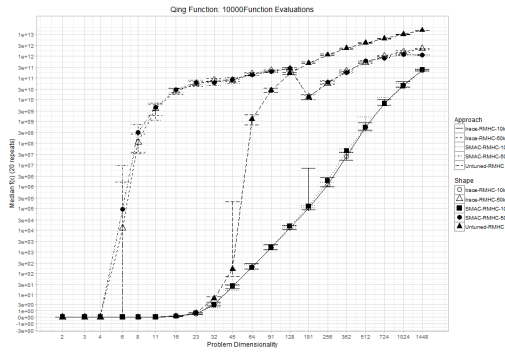
Alpine no.1 Function at 50,000 Evaluations



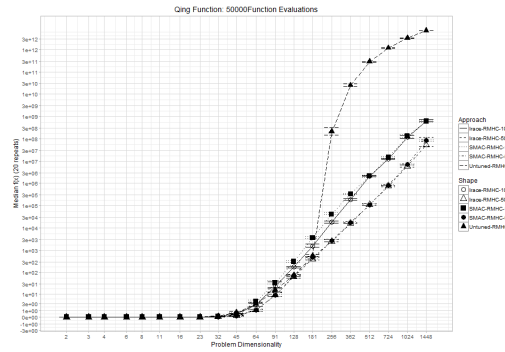
Brown Function at 10,000 Evaluations



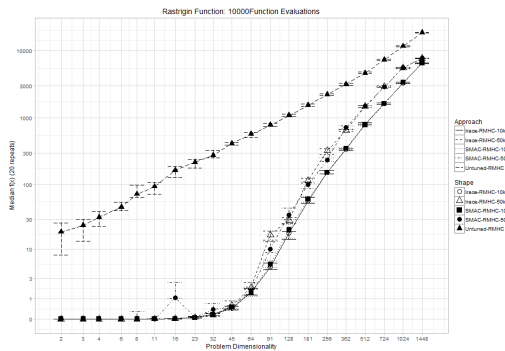
Brown Function at 50,000 Evaluations



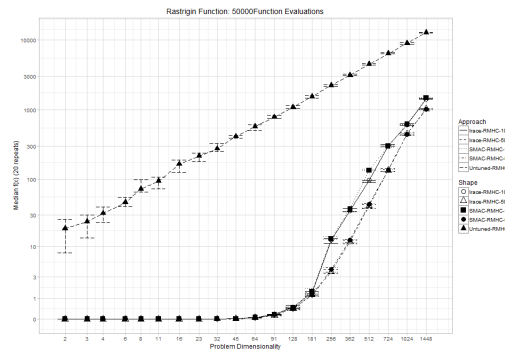
Qing Function at 10,000 Evaluations



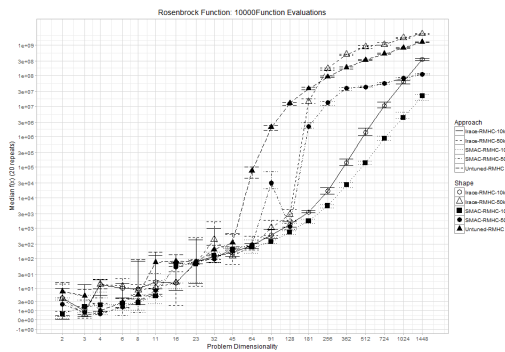
Qing Function at 50,000 Evaluations



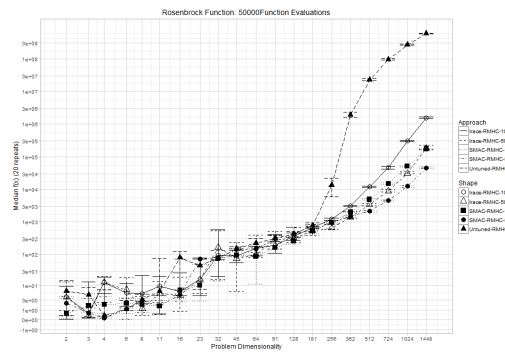
Rastrigin Function at 10,000 Evaluations



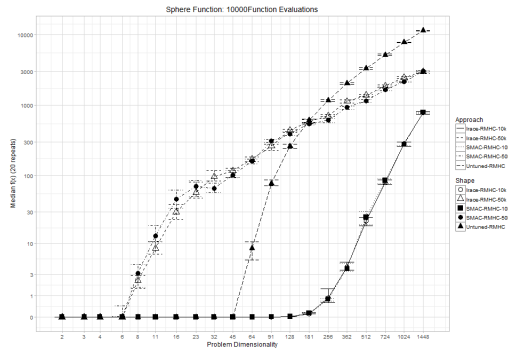
Rastrigin Function at 50,000 Evaluations



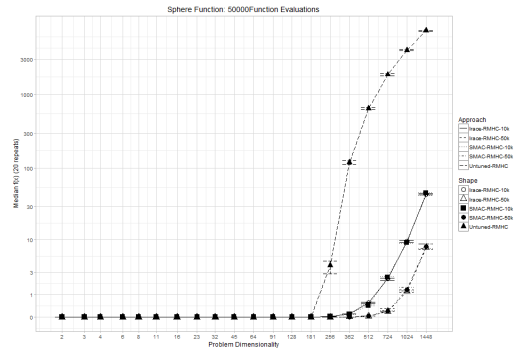
Rosenbrock Function at 10,000 Evaluations



Rosenbrock Function at 50,000 Evaluations

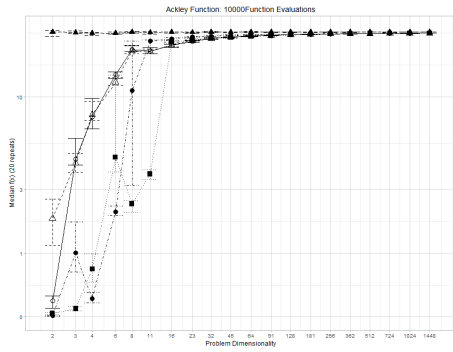


Sphere Function at 10,000 Evaluations

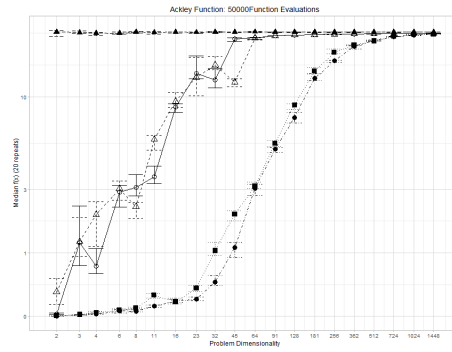


Sphere Function at 50,000 Evaluations

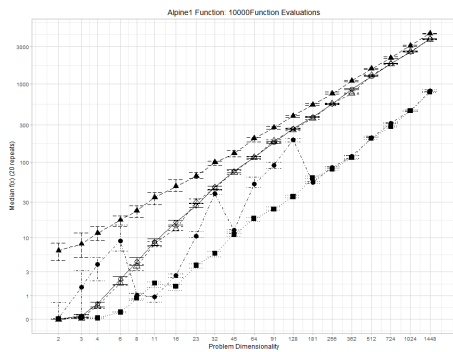
Simulated Annealing (SA)



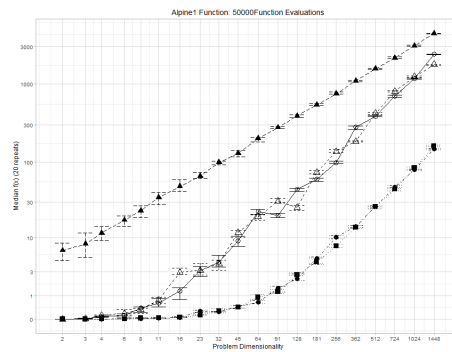
Ackley Function at 10,000 Evaluations



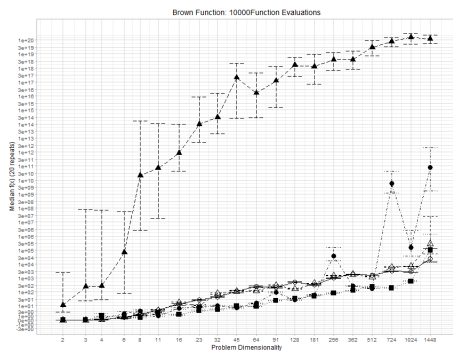
Ackley Function at 50,000 Evaluations



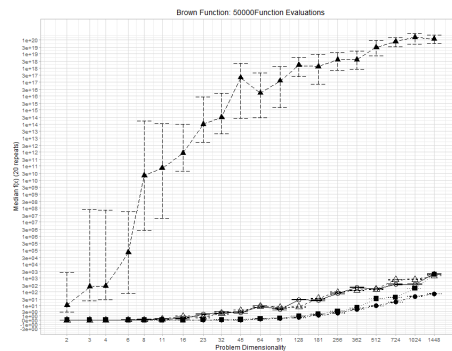
Alpine no.1 Function at 10,000 Evaluations



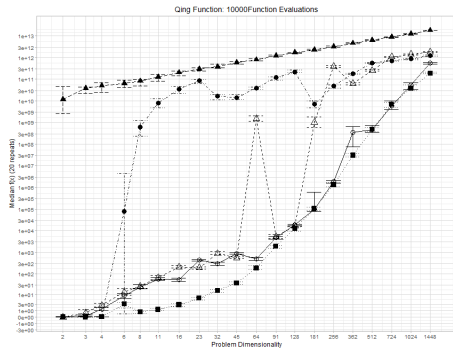
Alpine no.1 Function at 50,000 Evaluations



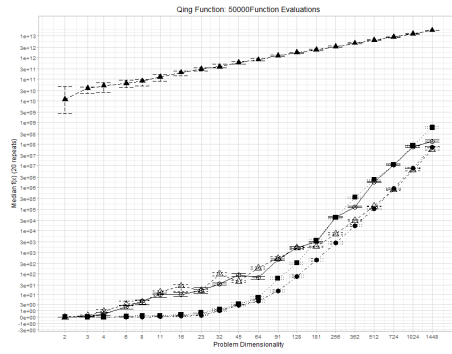
Brown Function at 10,000 Evaluations



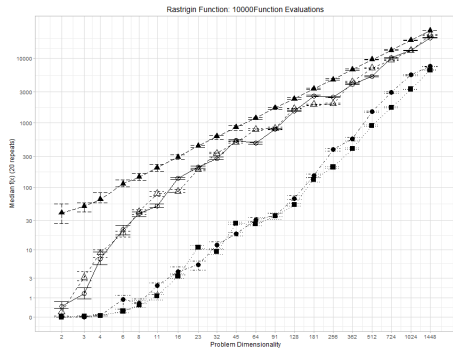
Brown Function at 50,000 Evaluations



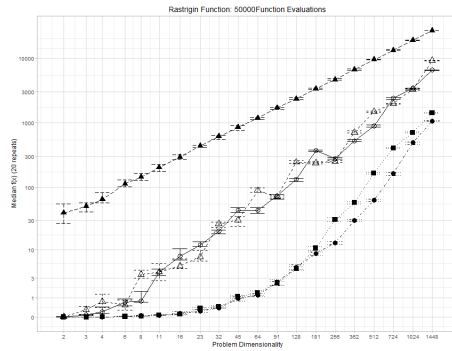
Qing Function at 10,000 Evaluations



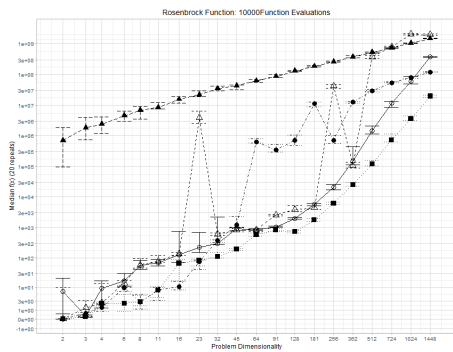
Qing Function at 50,000 Evaluations



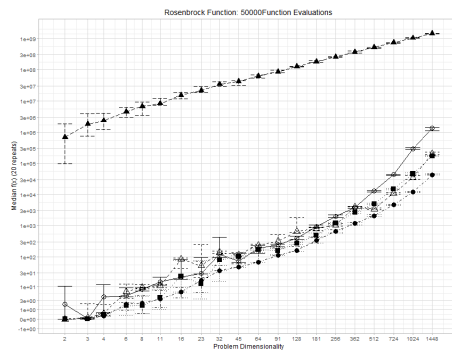
Rastrigin Function at 10,000 Evaluations



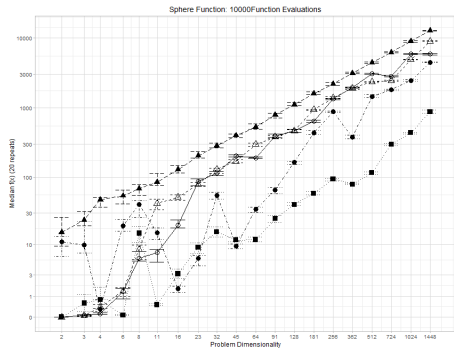
Rastrigin Function at 50,000 Evaluations



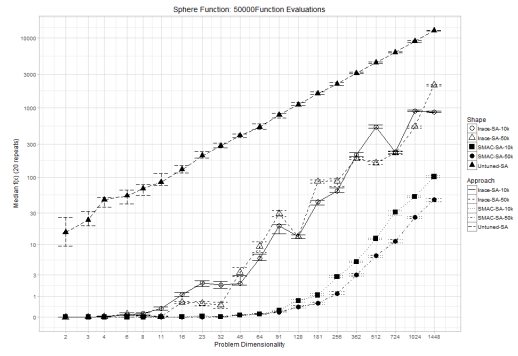
Rosenbrock Function at 10,000 Evaluations



Rosenbrock Function at 50,000 Evaluations

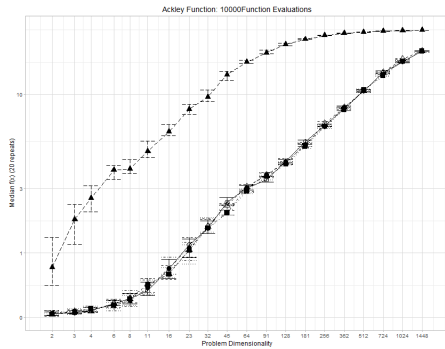


Sphere Function at 10,000 Evaluations

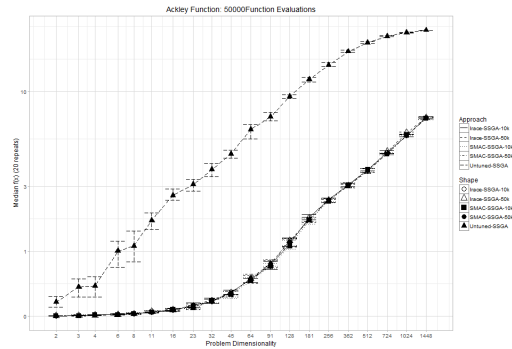


Sphere Function at 50,000 Evaluations

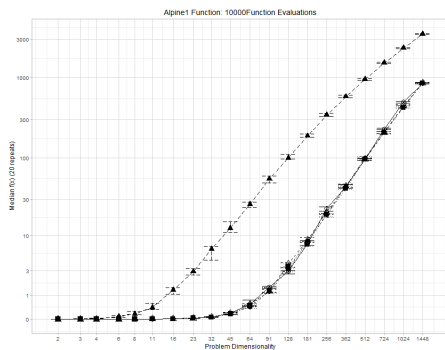
Steady State Genetic Algorithm (SSGA)



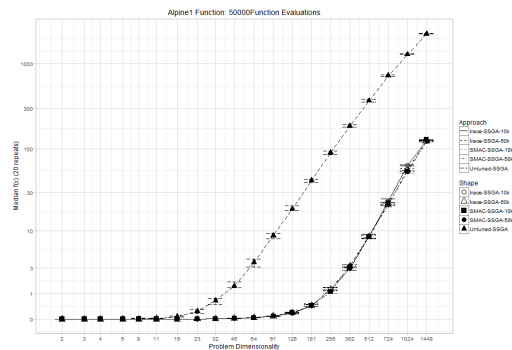
Ackley Function at 10,000 Evaluations



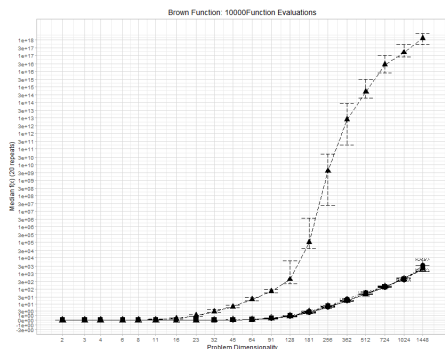
Ackley Function at 50,000 Evaluations



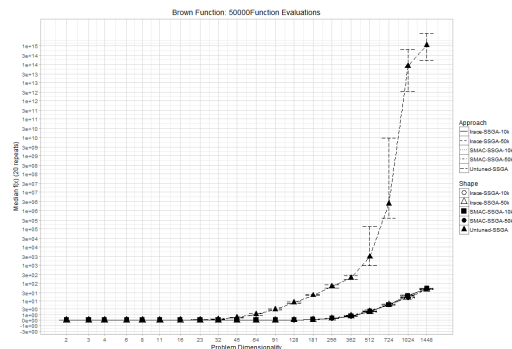
Alpine no.1 Function at 10,000 Evaluations



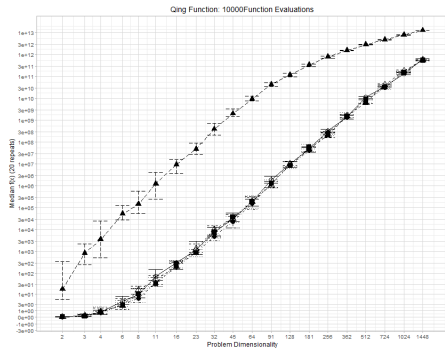
Alpine no.1 Function at 50,000 Evaluations



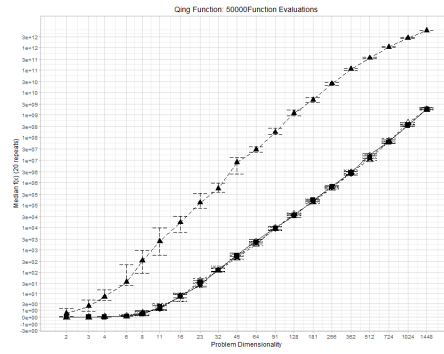
Brown Function at 10,000 Evaluations



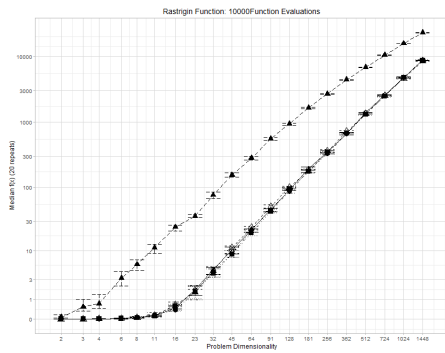
Brown Function at 50,000 Evaluations



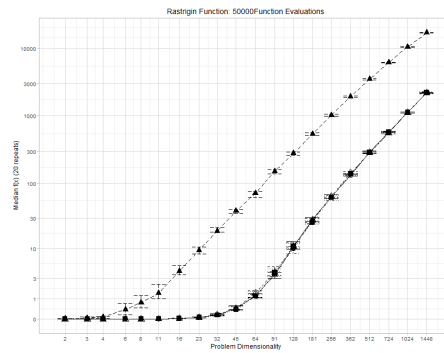
Qing Function at 10,000 Evaluations



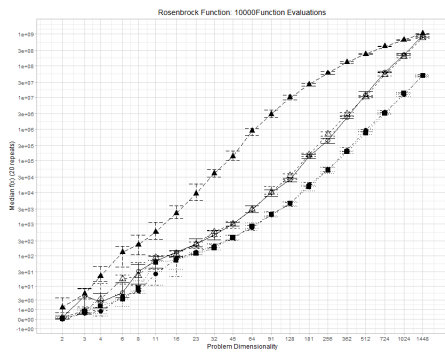
Qing Function at 50,000 Evaluations



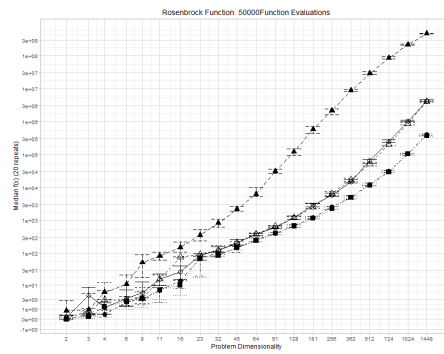
Rastrigin Function at 10,000 Evaluations



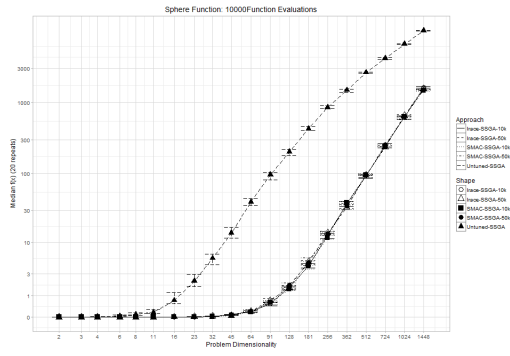
Rastrigin Function at 50,000 Evaluations



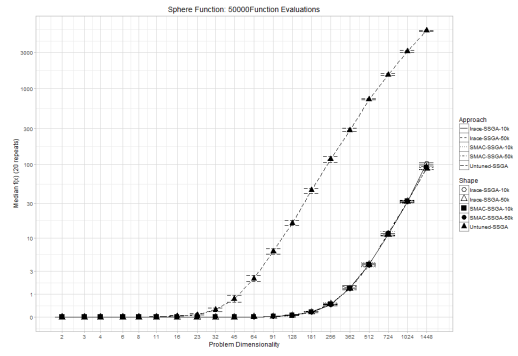
Rosenbrock Function at 10,000 Evaluations



Rosenbrock Function at 50,000 Evaluations

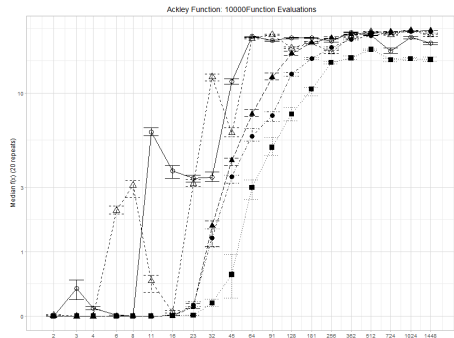


Sphere Function at 10,000 Evaluations

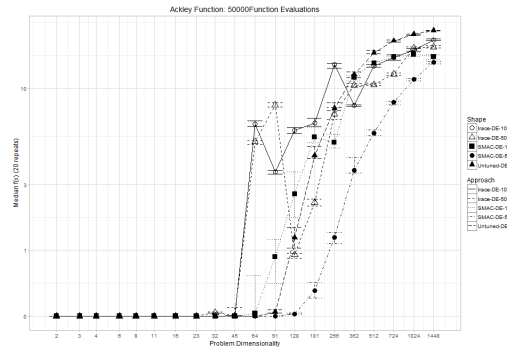


Sphere Function at 50,000 Evaluations

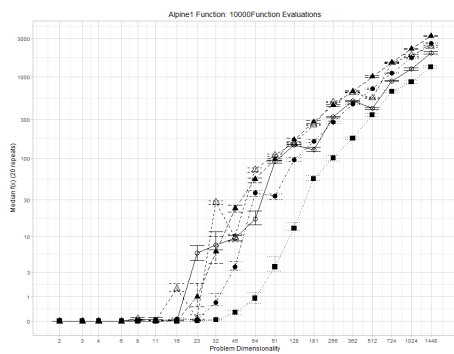
Differential Evolution (DE)



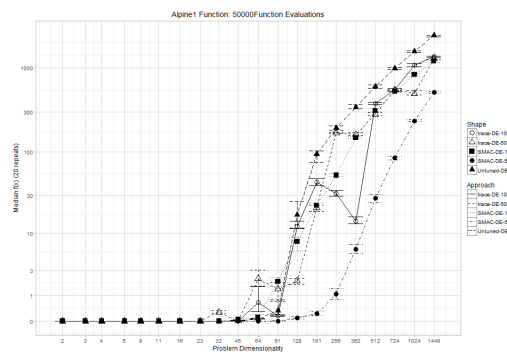
Ackley Function at 10,000 Evaluations



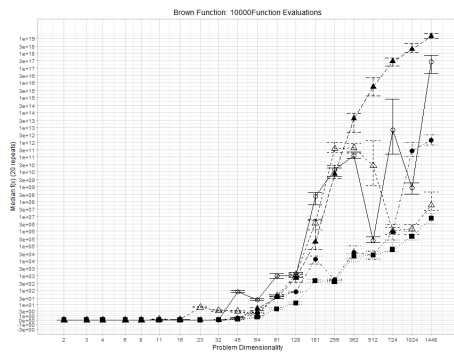
Ackley Function at 50,000 Evaluations



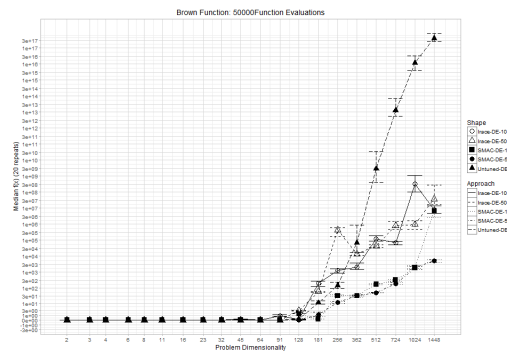
Alpine no.1 Function at 10,000 Evaluations



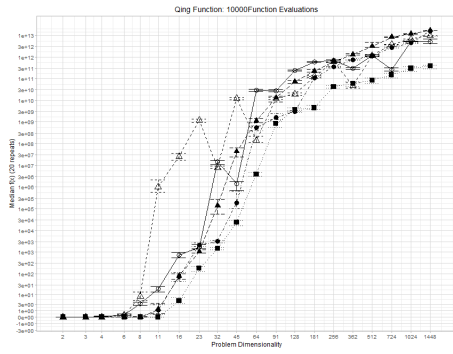
Alpine no.1 Function at 50,000 Evaluations



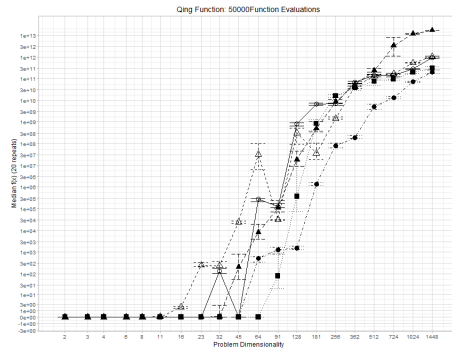
Brown Function at 10,000 Evaluations



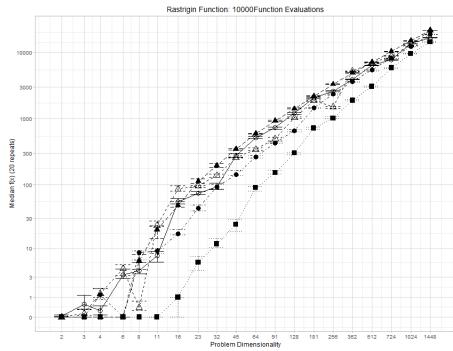
Brown Function at 50,000 Evaluations



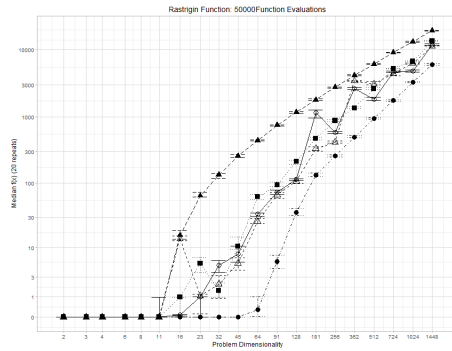
Qing Function at 10,000 Evaluations



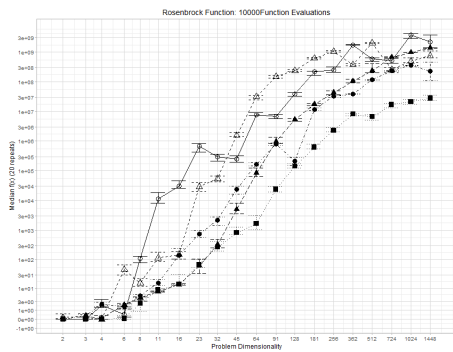
Qing Function at 50,000 Evaluations



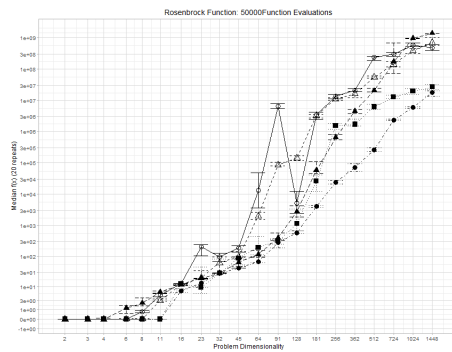
Rastrigin Function at 10,000 Evaluations



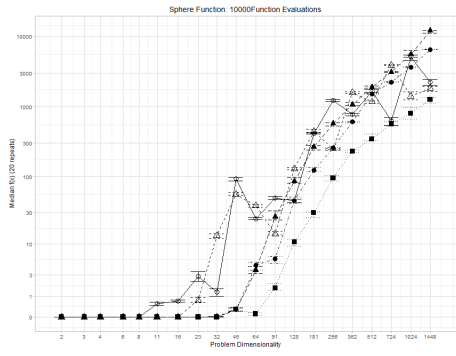
Rastrigin Function at 50,000 Evaluations



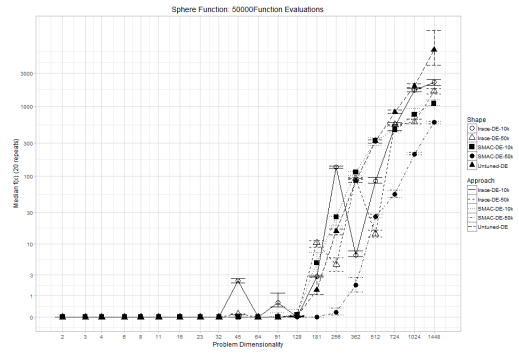
Rosenbrock Function at 10,000 Evaluations



Rosenbrock Function at 50,000 Evaluations

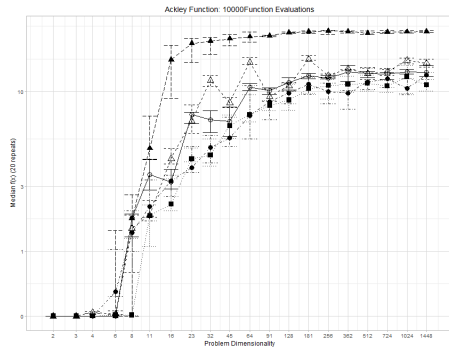


Sphere Function at 10,000 Evaluations

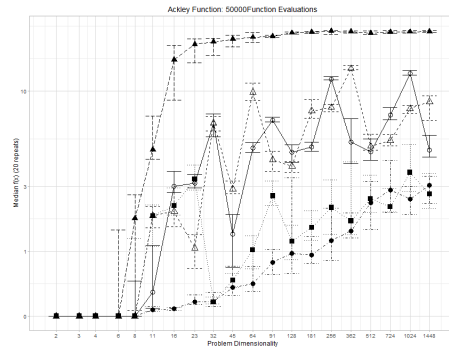


Sphere Function at 50,000 Evaluations

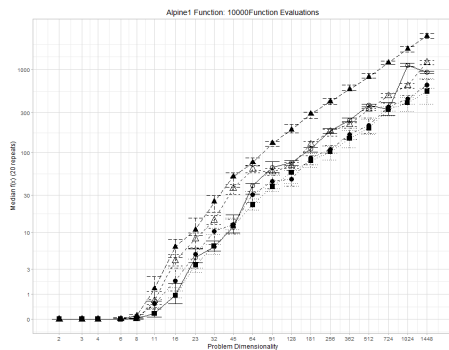
Particle Swarm Optimisation (PSO)



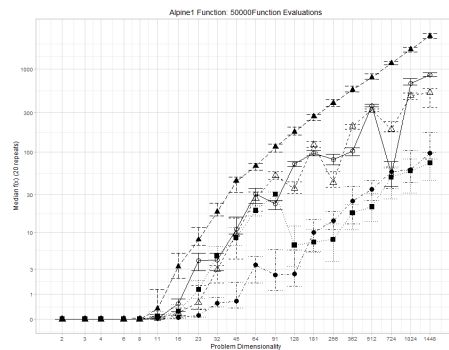
Ackley Function at 10,000 Evaluations



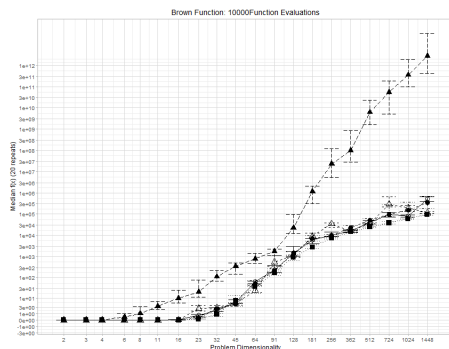
Ackley Function at 50,000 Evaluations



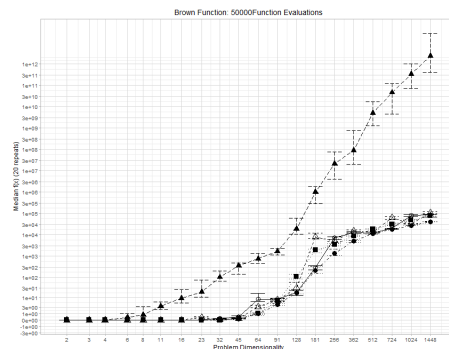
Alpine no.1 Function at 10,000 Evaluations



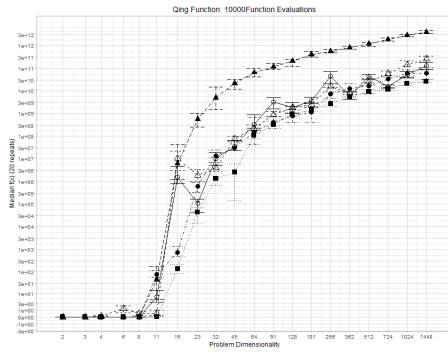
Alpine no.1 Function at 50,000 Evaluations



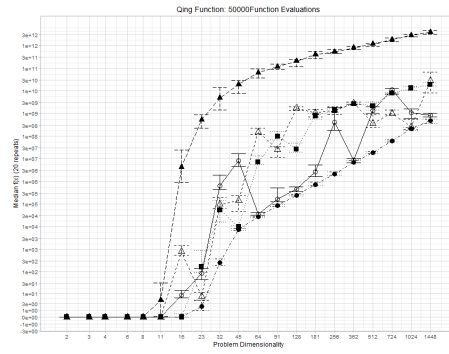
Brown Function at 10,000 Evaluations



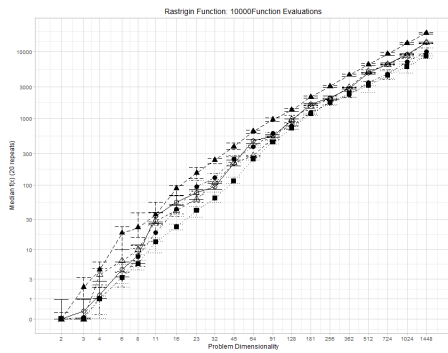
Brown Function at 50,000 Evaluations



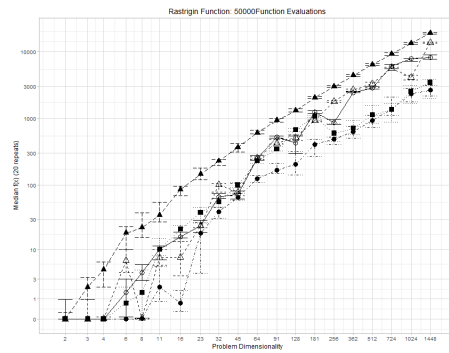
Qing Function at 10,000 Evaluations



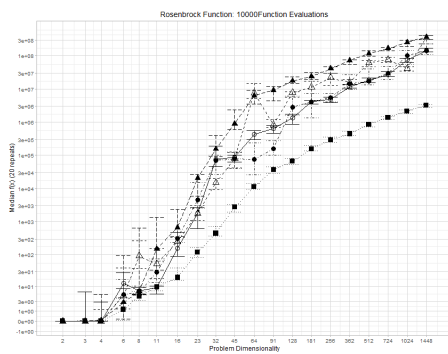
Qing Function at 50,000 Evaluations



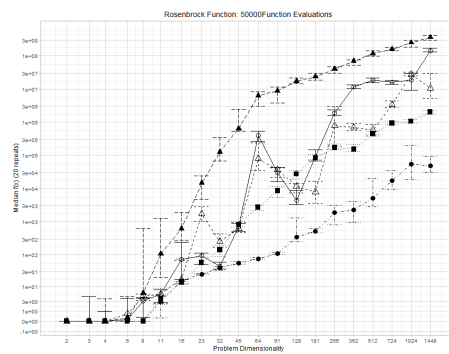
Rastrigin Function at 10,000 Evaluations



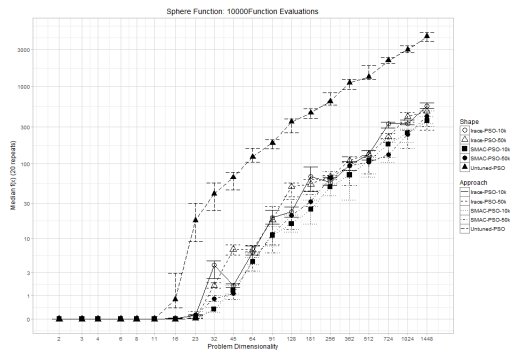
Rastrigin Function at 50,000 Evaluations



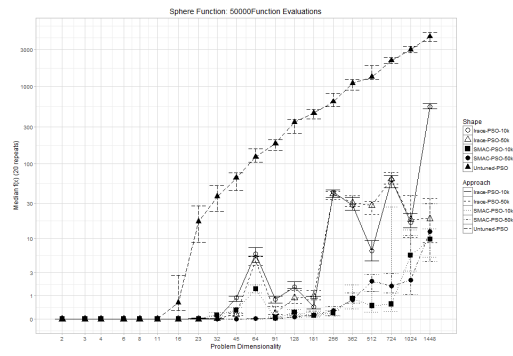
Rosenbrock Function at 10,000 Evaluations



Rosenbrock Function at 50,000 Evaluations

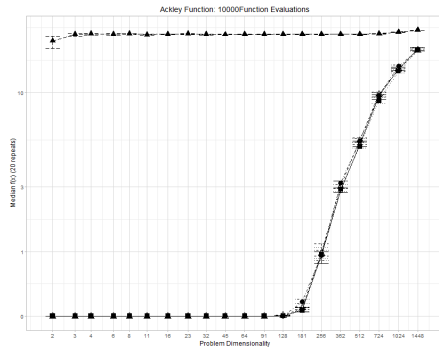


Sphere Function at 10,000 Evaluations

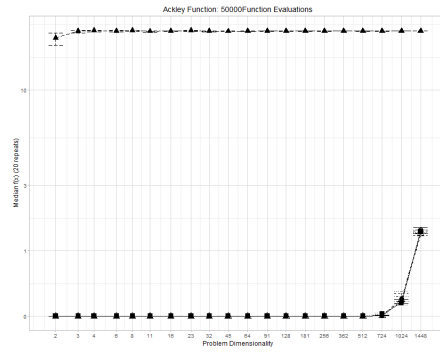


Sphere Function at 50,000 Evaluations

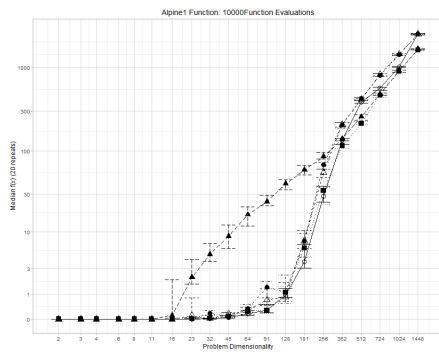
Covariance Matrix Adaption Evolutionary Strategy (CMA-ES)



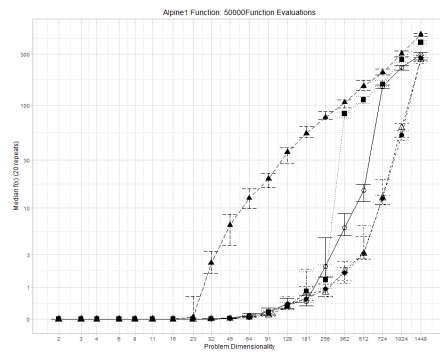
Ackley Function at 10,000 Evaluations



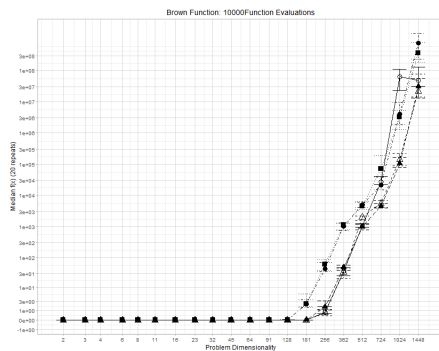
Ackley Function at 50,000 Evaluations



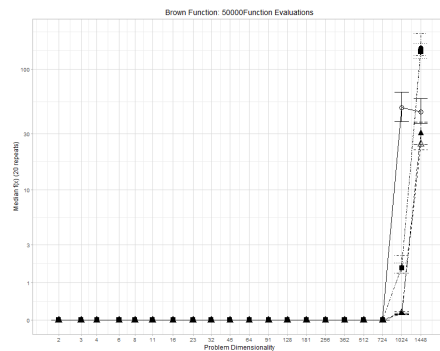
Alpine no.1 Function at 10,000 Evaluations



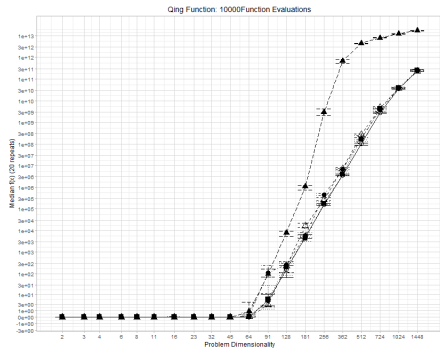
Alpine no.1 Function at 50,000 Evaluations



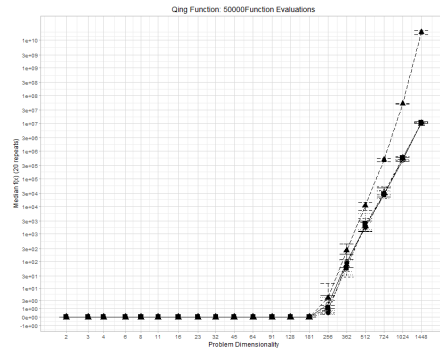
Brown Function at 10,000 Evaluations



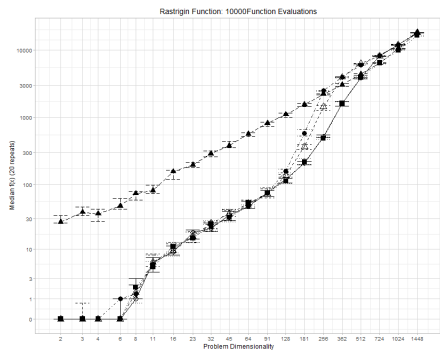
Brown Function at 50,000 Evaluations



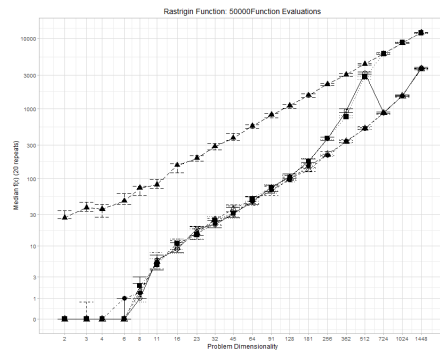
Qing Function at 10,000 Evaluations



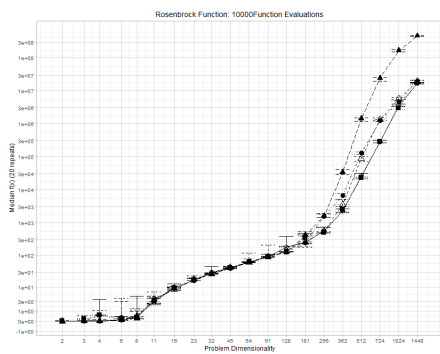
Qing Function at 50,000 Evaluations



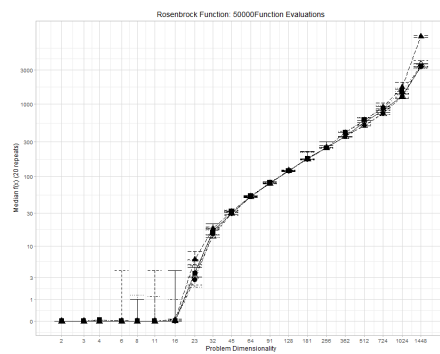
Rastrigin Function at 10,000 Evaluations



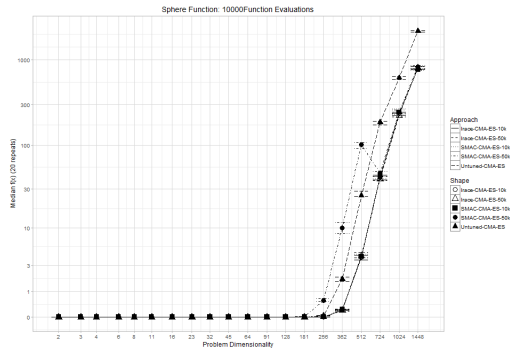
Rastrigin Function at 50,000 Evaluations



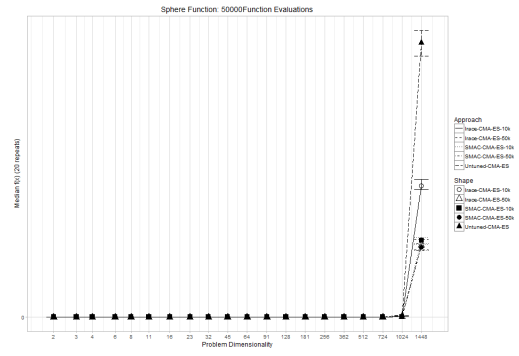
Rosenbrock Function at 10,000 Evaluations



Rosenbrock Function at 50,000 Evaluations



Sphere Function at 10,000 Evaluations



Sphere Function at 50,000 Evaluations

APPENDIX E: DEFINITIONS OF BENCHMARK FUNCTIONS IN THE TEST SUITE

E.1 ACKLEY FUNCTION NO.1

$$f(x) = -20 \exp \left(-0.2 \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2} \right) - \exp \left(\frac{1}{n} \sum_{i=1}^n \cos(2\pi x_i) \right) + 20 + \exp(1) \quad (\text{E.1})$$

Subject to domain, $-32.768 \leq x_i \leq 32.768$ for all $i \in \{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.

Figure E.1a and E.1b shows a 3D and contour plot representing Ackley function no.1 in 2-dimensional space.

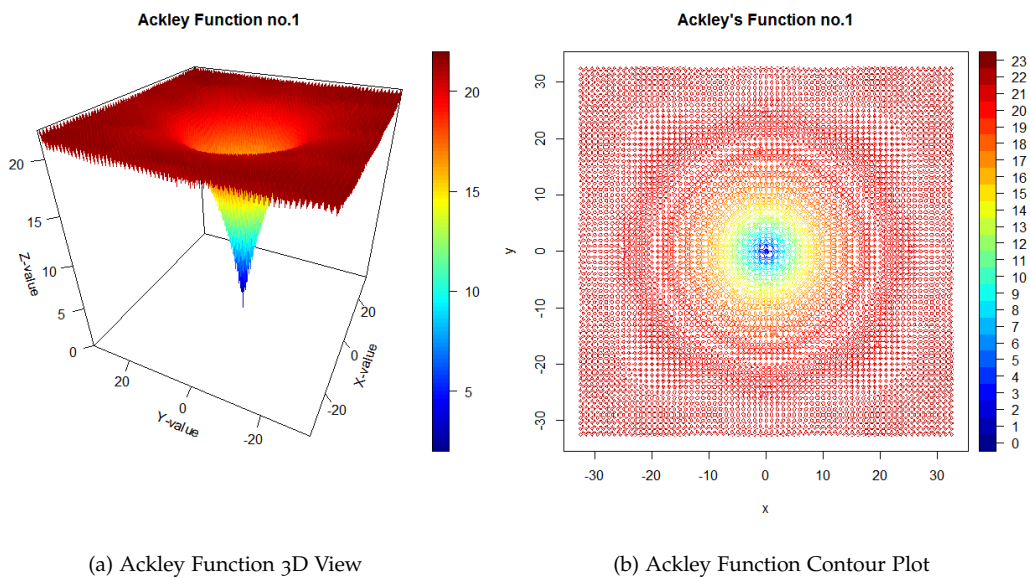


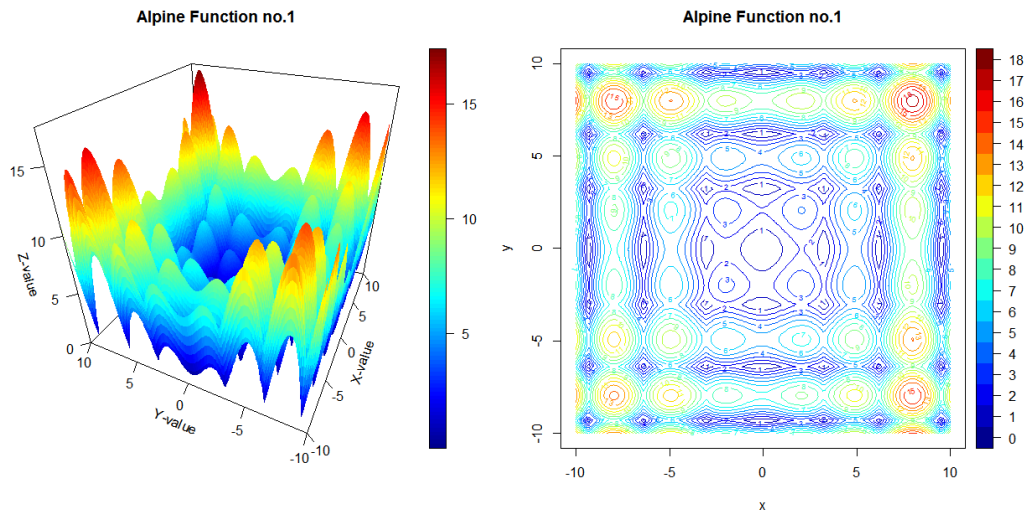
Figure E.1: Ackley Function no.1 in its 2-dimensional form

E.2 ALPINE FUNCTION NO.1

$$f(x) = \sum_{i=1}^n |x_i \sin(x_i) + 0.1x_i| \quad (\text{E.2})$$

Subject to domain, $-10 \leq x_i \leq 10$ for all $i \in \{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value 0.

Figure E.2a and E.2b shows a 3D and contour plot representing Alpine function no.1 in 2-dimensional space.



(a) Alpine Function no.1 3D View

(b) Alpine Function no.1 Contour Plot

Figure E.2: Alpine Function no.1 in its 2-dimensional form

E.3 BENT CIGAR FUNCTION

$$f(x) = x_1^2 + 10^6 \sum_{i=2}^n x_i^2 \tag{E.3}$$

Subject to domain, $-100 \leq x_i \leq 100$ for all $i \in \{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.

Figure E.3a and E.3b shows a 3D and contour plot representing Bent Cigar function in 2-dimensional space.

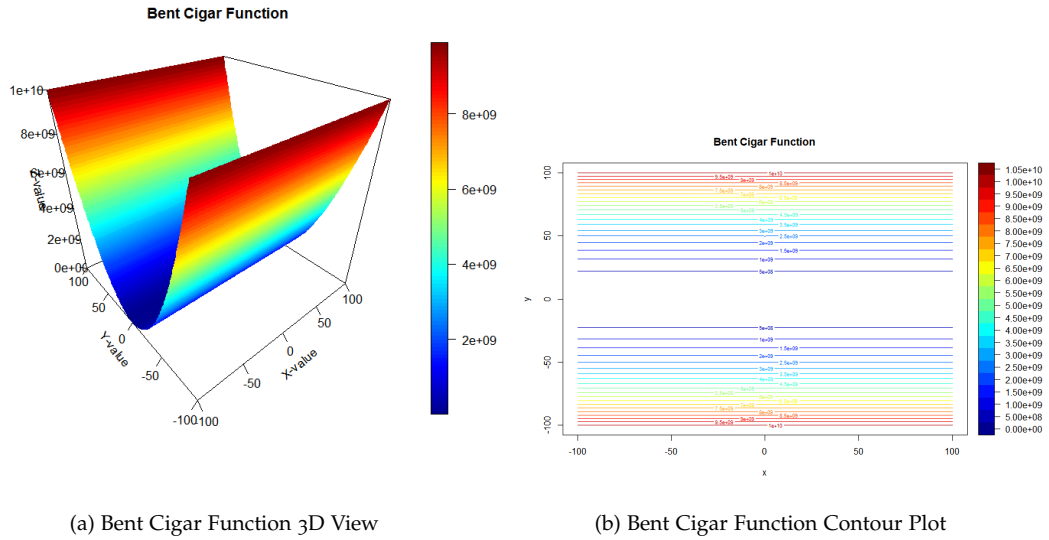


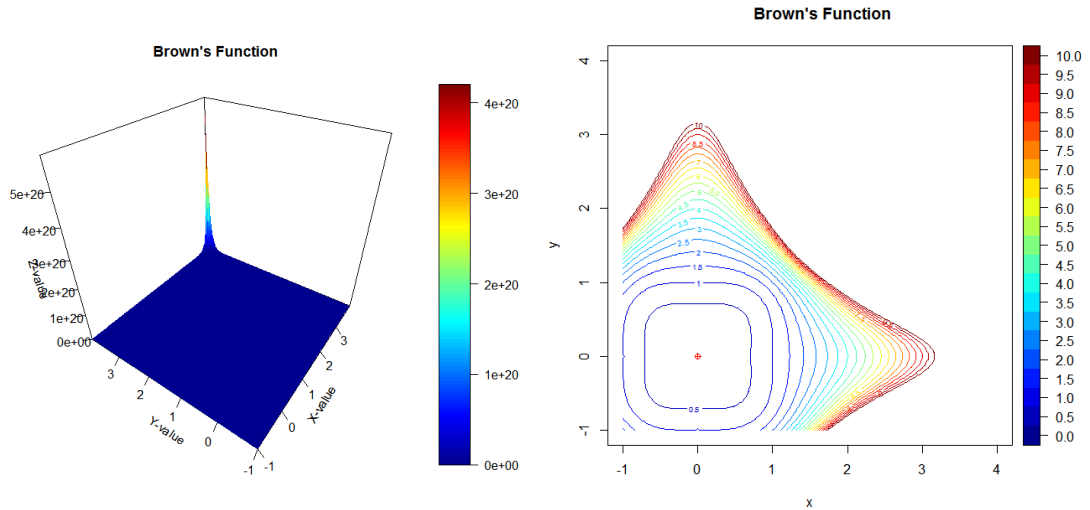
Figure E.3: Bent Cigar Function in its 2-dimensional form

E.4 BROWN FUNCTION

$$f(x) = \sum_{i=1}^{n-1} (x_i^2)^{x_{i+1}^2+1} + (x_{i+1}^2)^{x_i^2+1} \quad (E.4)$$

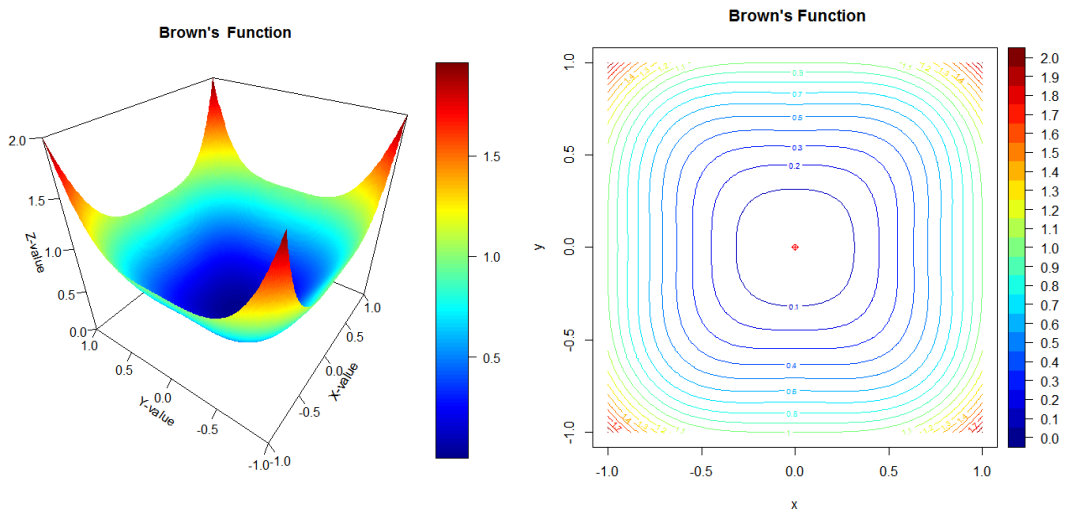
Subject to domain, $-1 \leq x_i \leq 4$ for all $i \in \{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.

Figures E.4a and E.4b shows a 3D and contour plot representing the Brown function in 2-dimensional space, subject to the full sized domain. To clearly show the function structure near the global optimum, Figures E.4c and E.4d represents the function with domain restricted to $-1 \leq x_i \leq 1$



(a) Brown Function 3D View (Full Domain)

(b) Brown Function Contour Plot (Full Domain)



(c) Brown Function Contour Plot (Limited Domain)

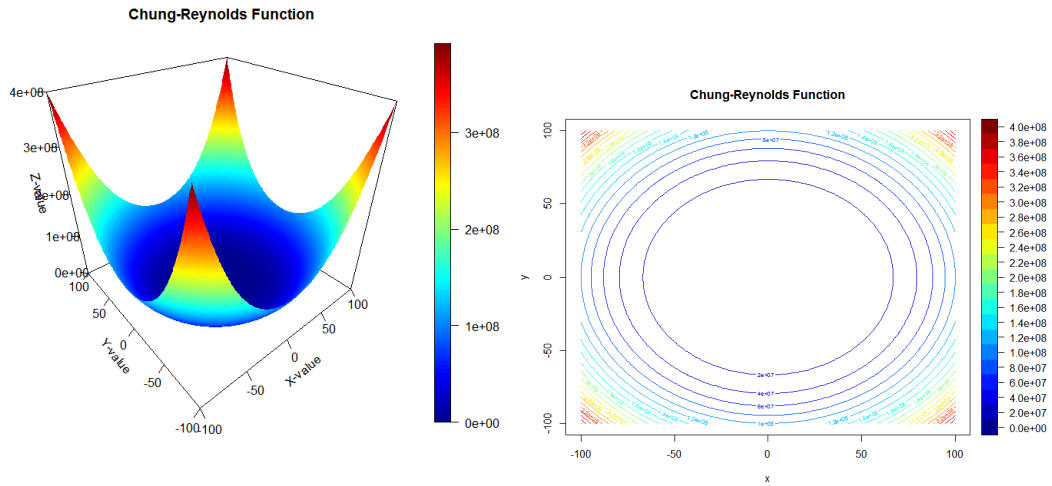
(d) Brown Function Contour Plot (Limited Domain)

Figure E.4: Brown Function in its 2-dimensional form: (a) and (b) shows the function's full domain, and (c) and (d) shows the function with domain limited to $-1 \leq x_i \leq 1$

E.5 CHUNG-REYNOLDS FUNCTION

$$f(x) = \left(\sum_{i=1}^n x_i^2 \right)^2 \quad (\text{E.5})$$

Subject to domain, $-100 \leq x_i \leq 100$ for all $i \{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.



(a) Chung-Reynolds Function 3D View

(b) Chung-Reynolds Function Contour Plot

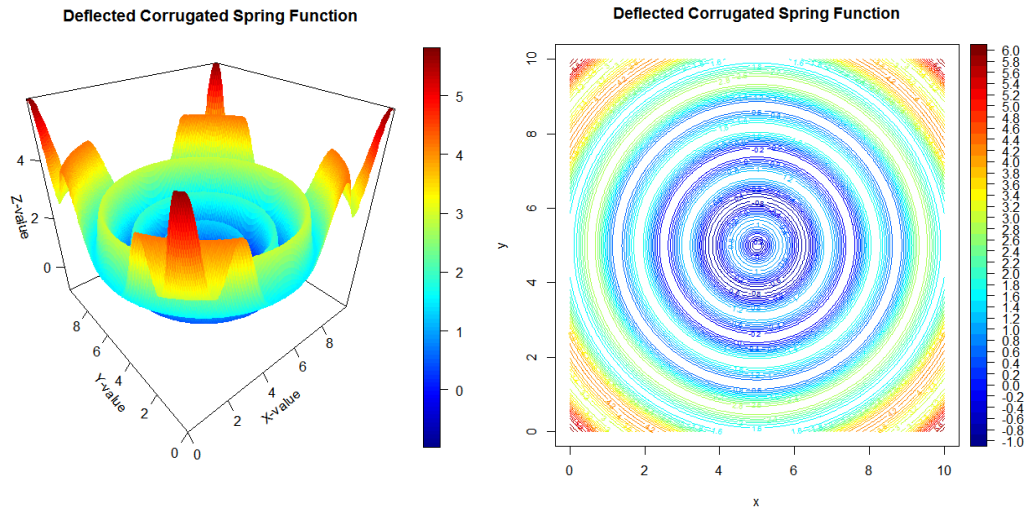
Figure E.5: Chung-Reynolds Function in its 2-dimensional form

E.6 DEFLECTED CORRUGATED SPRING FUNCTION

$$f(x) = 0.1 \sum_{i=1}^n \left[(x_i - \alpha)^2 - \cos \left(K \sqrt{\sum_{i=1}^n (x_i - \alpha)^2} \right) \right] \quad (\text{E.6})$$

The function is subject to domain, $0 \leq x_i \leq 2\alpha$ for all $i \in \{1, 2, \dots, n\}$. and the global minima is found at the point $x_i^* = \alpha$ with value $f(x^*) = -1$. The variables α and K are both set equal to 5.

Figures E.6a and E.6b shows a 3D and contour plot representing the Deflected Corrugated Spring function in 2-dimensional space.



(a) Deflected Corrugated Spring Function 3D View (b) Deflected Corrugated Spring Function Contour Plot

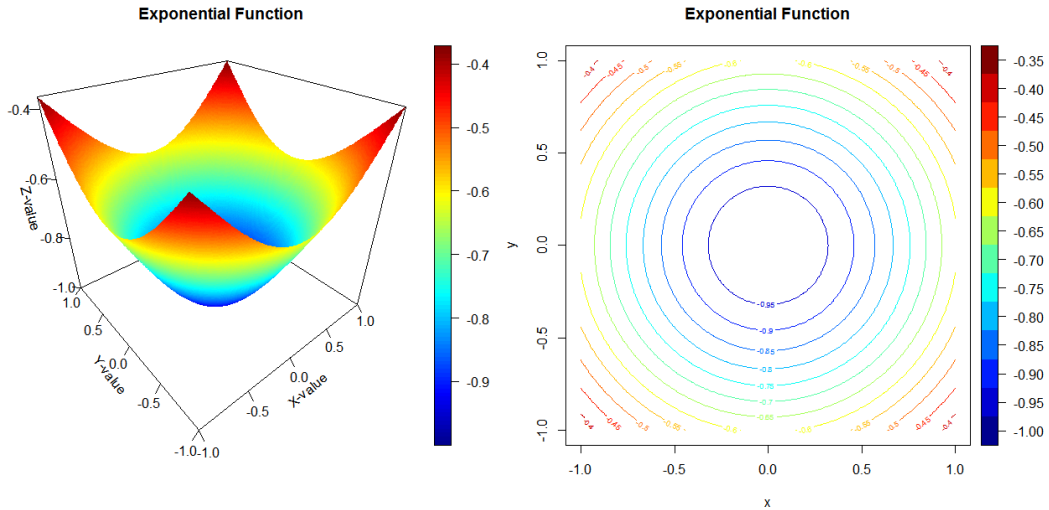
Figure E.6: Deflected Corrugated Spring Function in its 2-dimensional form

E.7 EXPONENTIAL FUNCTION

$$f(x) = -\exp\left(-0.5 \sum_{i=1}^n x_i^2\right) \quad (\text{E.7})$$

Subject to domain, $-1 \leq x_i \leq 1$ for all $i\{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = -1$. This optimum of this function is it's maxima not it's minima - to be used for minimisation purposed the result of the function is negated - i.e., the minima becomes the negative maxima.

Figures E.7a and E.7b shows a 3D and contour plot representing the Exponential function in 2-dimensional space.



(a) Exponential Function 3D View

(b) Exponential Function Contour Plot

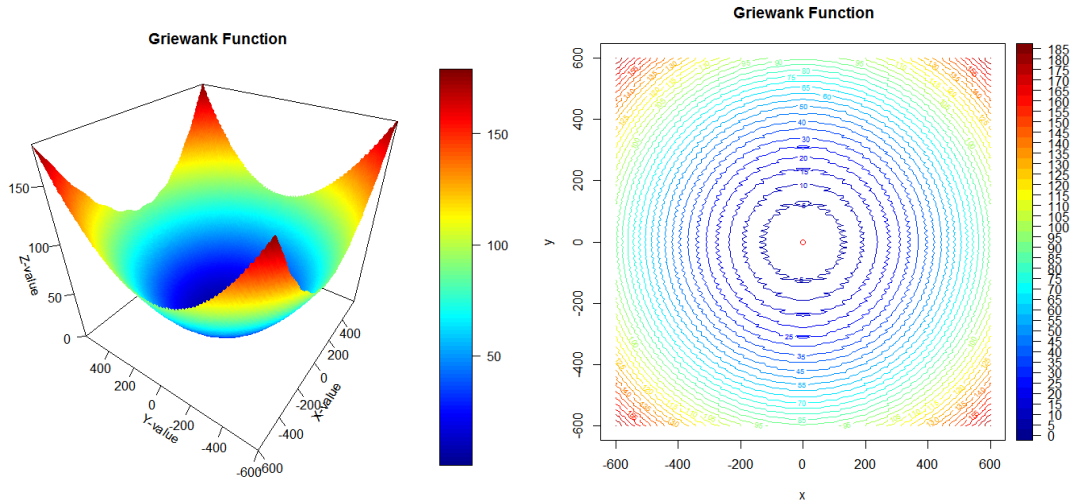
Figure E.7: Exponential Function in its 2-dimensional form

E.8 GRIEWANK FUNCTION

$$f(x) = \frac{1}{4000} \sum_{i=1}^n x_i^2 - \prod_{i=1}^n \cos\left(\frac{x_i}{\sqrt{i}}\right) + 1 \quad (\text{E.8})$$

Subject to domain, $-600 \leq x_i \leq 600$ for all $i\{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.

Figures E.8a and E.8b shows a 3D and contour plot representing the Griewank function in 2-dimensional space.



(a) Griewank Function 3D View

(b) Griewank Function Contour Plot

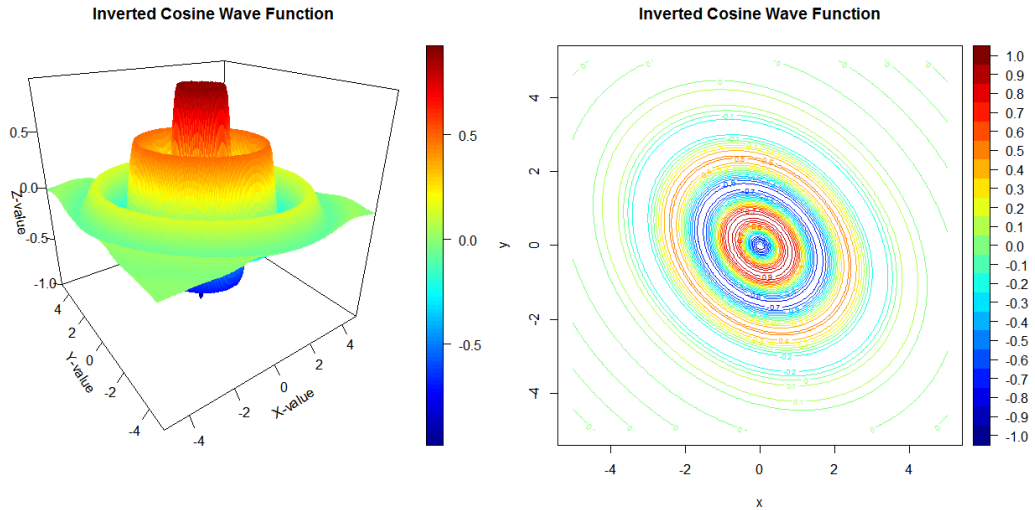
Figure E.8: Griewank Function in its 2-dimensional form

E.9 INVERTED COSINE WAVE FUNCTION

$$f(x) = \sum_{i=1}^{n-1} \left\{ e^{\left[\frac{-(x_i^2 + x_{i+1}^2 + 0.5x_i x_{i+1})}{8} \right]} \cos \left(4 \times \sqrt{x_i^2 + x_{i+1}^2 + 0.5x_i x_{i+1}} \right) \right\} \quad (E.9)$$

Subject to domain, $-5 \leq x_i \leq 5$ for all $i \{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = -n + 1$.

Figures E.9a and E.9b shows a 3D and contour plot representing the Inverted Cosine Wave function in 2-dimensional space.



(a) Inverted Cosine Wave Function 3D View

(b) Inverted Cosine Wave Function Contour Plot

Figure E.9: Inverted Cosine Wave Function in its 2-dimensional form

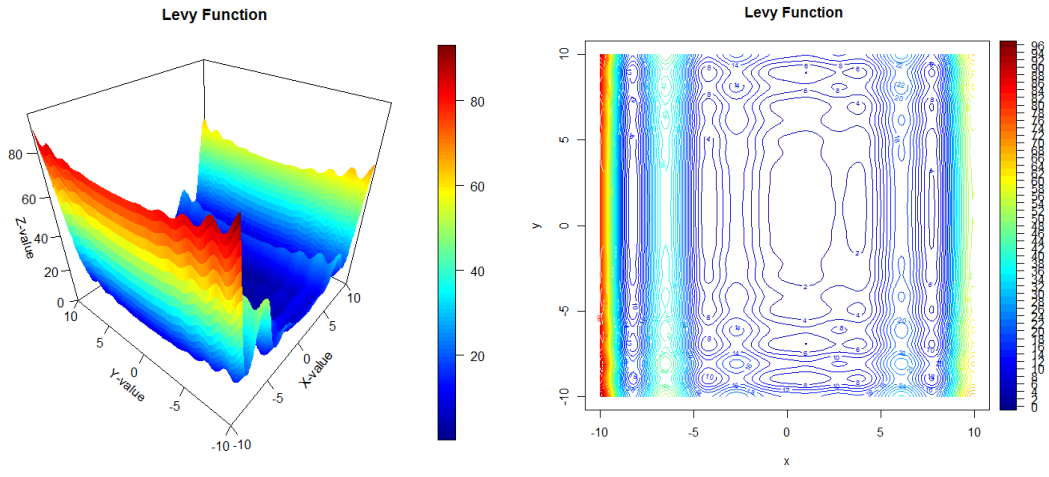
E.10 LEVY FUNCTION

$$f(x) = \sin^2(\pi w_1) + \sum_{i=1}^{n-1} (w_i - 1)^2 \left[1 + 10 \sin^2(\pi w_i + 1) \right] + (w_n - 1)^2 \left[1 + \sin^2(2\pi w_n) \right], \quad (\text{E.10})$$

$$\text{where } w_i = 1 + \frac{x_i - 1}{4} \text{ for all } i\{1, 2, \dots, n\}$$

Subject to domain, $-10 \leq x_i \leq 10$ for all $i\{1, 2, \dots, n\}$. The global minima is found at the point $x_i^* = 1$ with value $f(x^*) = 0$.

Figures E.10a and E.10b shows a 3D and contour plot representing the Levy function in 2-dimensional space.



(a) Levy Function 3D View

(b) Levy Function Contour Plot

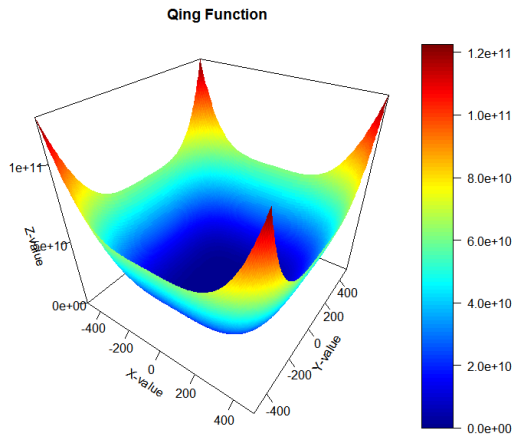
Figure E.10: Levy Function in its 2-dimensional form

E.11 QING FUNCTION

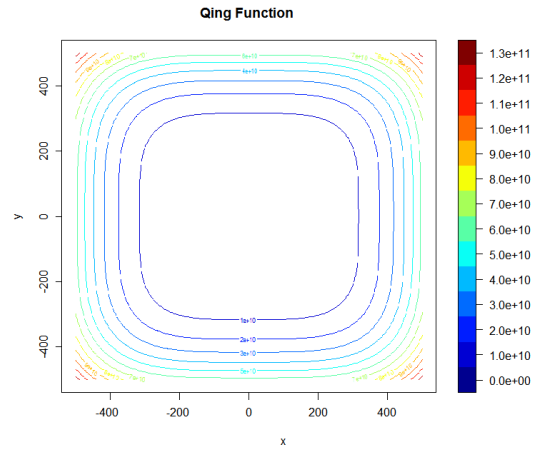
$$f(x) = \sum_{i=1}^n (x_i^2 - i)^2 \quad (\text{E.11})$$

Subject to domain, $-500 \leq x_i \leq 500$ for all $i \in \{1, 2, \dots, n\}$. There are 2^n global minima found at points $x_i^* = \pm\sqrt{i}$, where $i = 1 \dots n$, with value $f(x^*) = 0$.

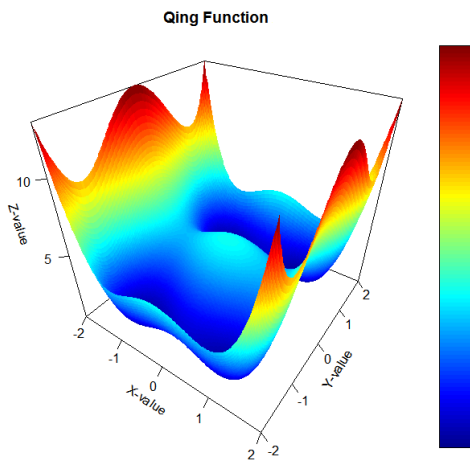
Figures E.11a and E.11b shows a 3D and contour plot representing the Qing function in 2-dimensional space, subject to the full sized domain. To clearly show the function structure near the global optima, Figures E.11c and E.11d represents the function with domain restricted to $-2 \leq x_i \leq 2$



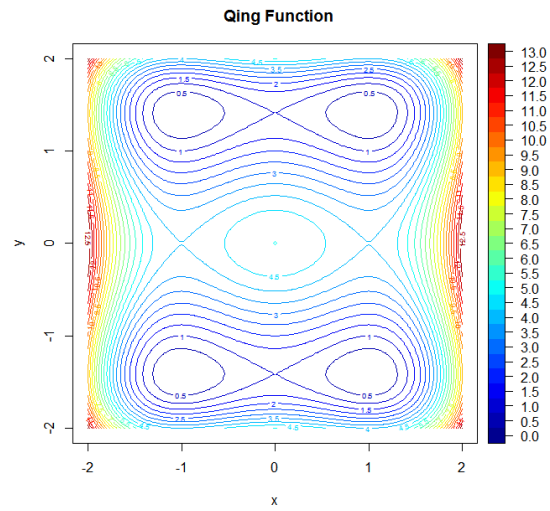
(a) Qing Function 3D View (Full Domain)



(b) Qing Function Contour Plot (Full Domain)



(c) Qing Function Contour Plot (Limited Domain)



(d) Qing Function Contour Plot (Limited Domain)

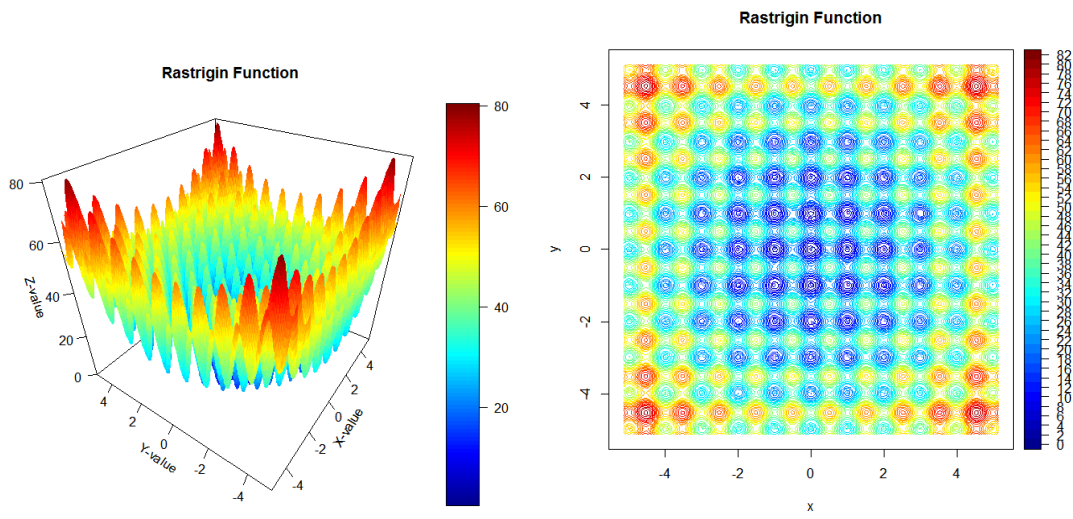
Figure E.11: Qing Function in its 2-dimensional form: (a) and (b) shows the function's full domain, and (c) and (d) shows the function with domain limited to $-2 \leq x_i \leq 2$

E.12 RASTRIGIN FUNCTION

$$f(x) = 10n \sum_{i=1}^n \left[x_i^2 - 10 \cos(2\pi x_i) \right] \quad (\text{E.12})$$

Subject to domain, $-5.12 \leq x_i \leq 5.12$ for all $i\{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.

Figures E.12a and E.12b shows a 3D and contour plot representing the Rastrigin function in 2-dimensional space.



(a) Rastrigin Function 3D View

(b) Rastrigin Function Contour Plot

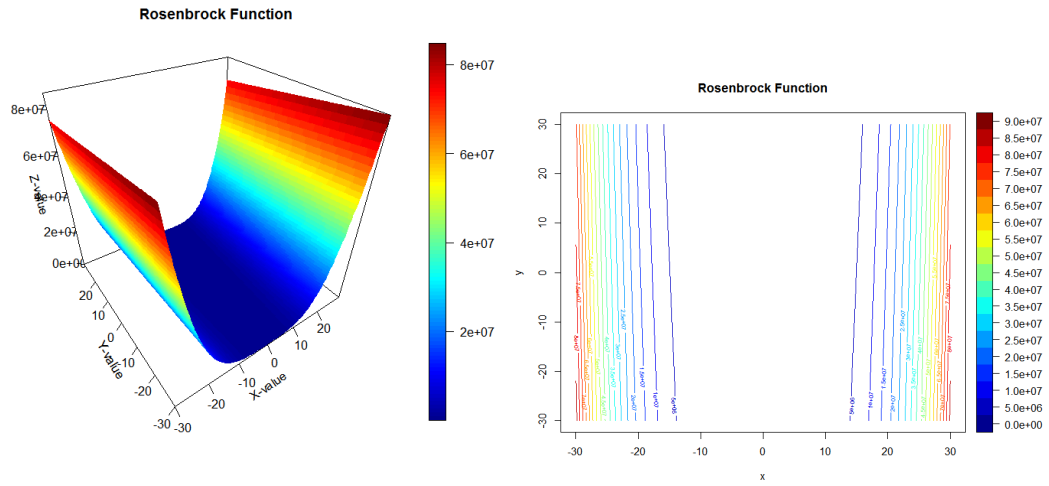
Figure E.12: Rastrigin Function in its 2-dimensional form

E.13 ROSENBROCK FUNCTION

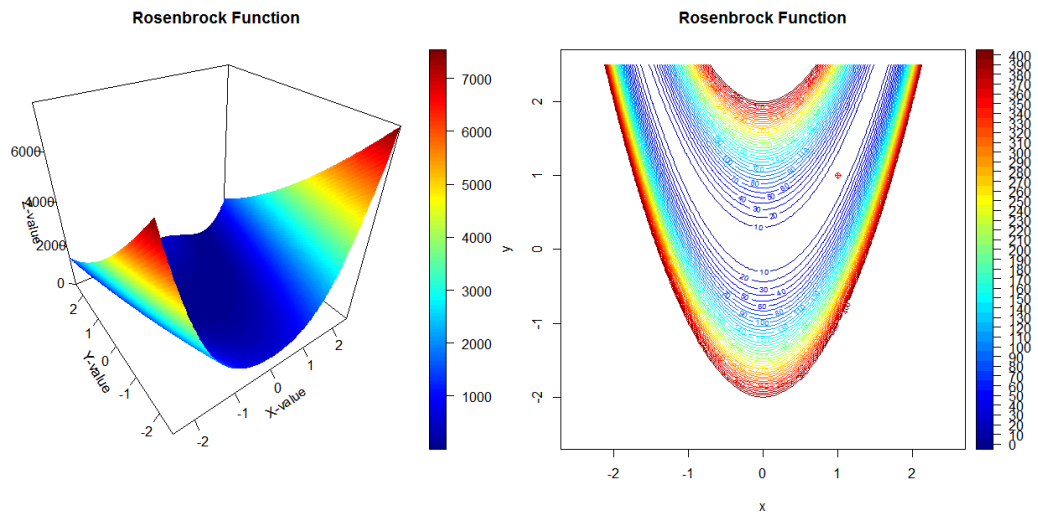
$$f(x) = \sum_{i=1}^{n-1} \left[100 (x_{i+1} - x_i^2)^2 + (x_i - 1)^2 \right] \quad (\text{E.13})$$

Subject to domain, $-30 \leq x_i \leq 30$ for all $i \{1, 2, \dots, n\}$. The global minima is found at the point $x_i^* = 1$ with value $f(x^*) = 0$.

Figures E.13a and E.13b shows a 3D and contour plot representing the Rosenbrock function in 2-dimensional space, subject to the full sized domain. To clearly show the function structure near the global optima, Figures E.13c and E.13d represents the function with domain restricted to $-2.5 \leq x_i \leq 2.5$



(a) Rosenbrock Function 3D View (Full Domain) (b) Rosenbrock Function Contour Plot (Full Domain)



(c) Rosenbrock Function Contour Plot (Limited Domain) (d) Rosenbrock Function Contour Plot (Limited Domain)

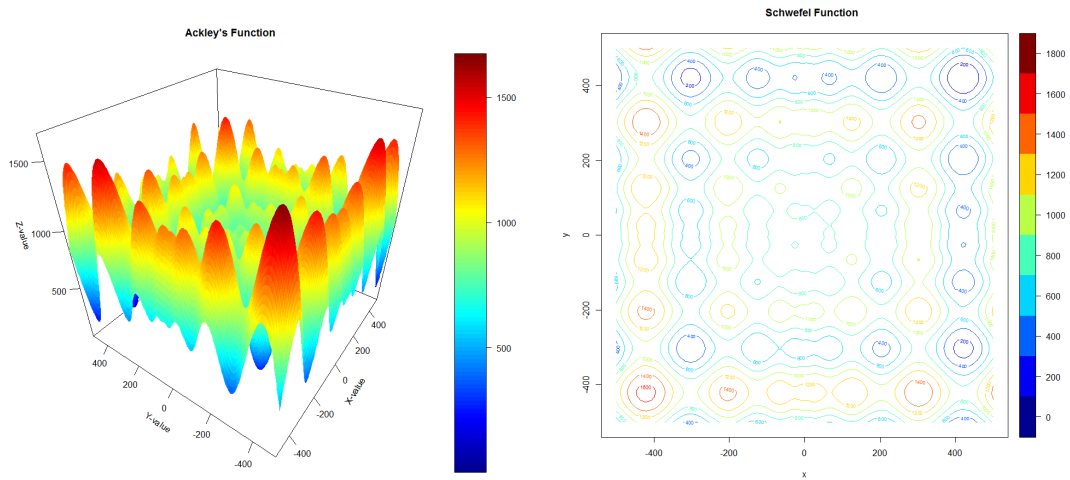
Figure E.13: Rosenbrock Function in its 2-dimensional form: (a) and (b) shows the function's full domain, and (c) and (d) shows the function with domain limited to $-2.5 \leq x_i \leq 2.5$

E.14 SCHWEFEL FUNCTION NO.26

$$f(x) = 418.9829n - \sum_{i=1}^n x_i \sin(\sqrt{|x_i|}) \quad (\text{E.14})$$

Subject to domain, $-500 \leq x_i \leq 500$ for all $i\{1,2,\dots,n\}$. The global minima is found at the point $x_i^* = 420.9687$ with value $f(x^*) = 0$.

Figures E.14a and E.14b shows a 3D and contour plot representing the Schwefel function in 2-dimensional space.



(a) Schwefel Function 3D View

(b) Schwefel Function Contour Plot

Figure E.14: Schwefel Function in its 2-dimensional form

E.15 SPHERE FUNCTION

$$f(x) = \sum_{i=1}^n x_i^2 \tag{E.15}$$

Subject to domain, $-5.12 \leq x_i \leq 5.12$ for all $i \in \{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.

Figures E.15a and E.15b shows a 3D and contour plot representing the Sphere function in 2-dimensional space.

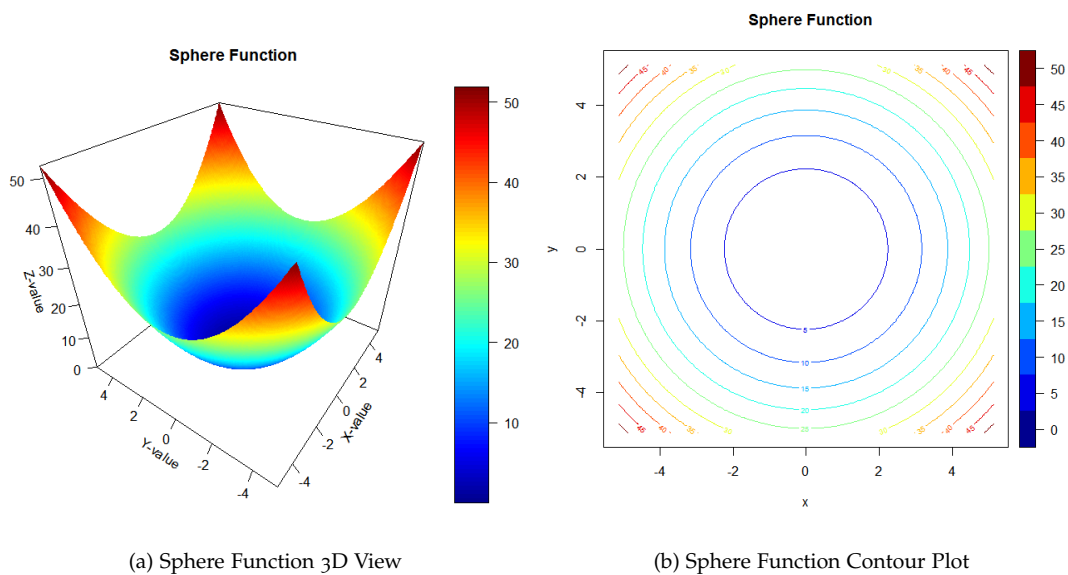


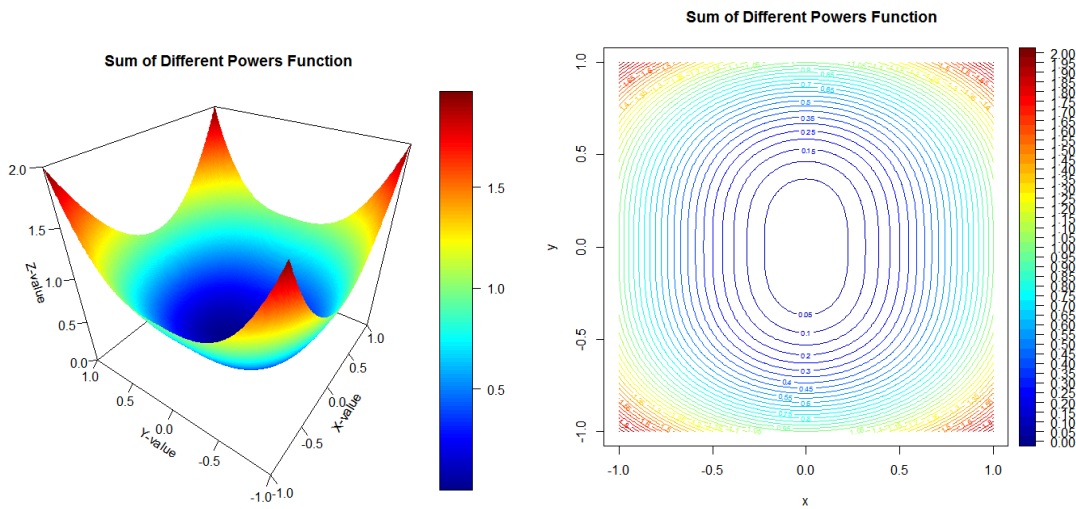
Figure E.15: Sphere Function in its 2-dimensional form

E.16 SUM OF DIFFERENT POWERS FUNCTION

$$f(x) = \sum_{i=1}^n |x_i|^{i+1} \tag{E.16}$$

Subject to domain, $-1 \leq x_i \leq 1$ for all $i\{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.

Figures E.16a and E.16b shows a 3D and contour plot representing the Sum of Different Powers function in 2-dimensional space.



(a) Sum of Different Powers Function 3D View

(b) Sum of Different Powers Function Contour Plot

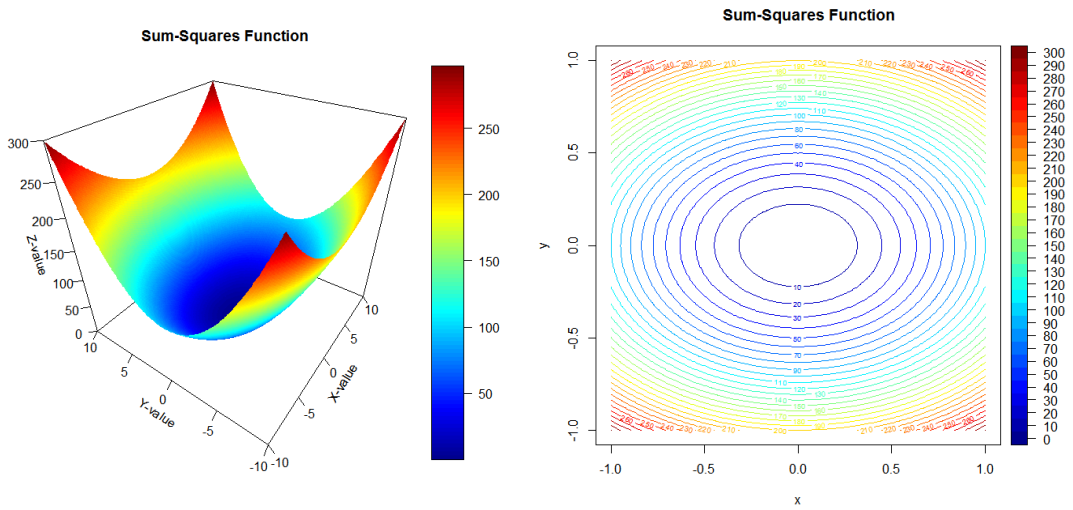
Figure E.16: Sum of Different Powers Function in its 2-dimensional form

E.17 SUM SQUARES FUNCTION

$$f(x) = \sum_{i=1}^n ix_i^2 \tag{E.17}$$

Subject to domain, $-10 \leq x_i \leq 10$ for all $i\{1, 2, \dots, n\}$. The global minima is found at the origin $x_i^* = 0$ with value $f(x^*) = 0$.

Figures E.17a and E.17b shows a 3D and contour plot representing the Sum Squares function in 2-dimensional space.



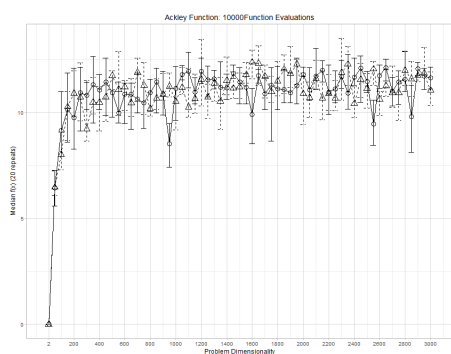
(a) Sum Squares Function 3D View

(b) Sum Squares Function Contour Plot

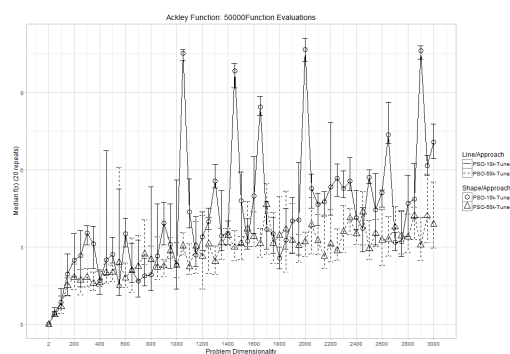
Figure E.17: Sum Squares Function in its 2-dimensional form

SUPPLEMENTARY EXPERIMENT MATERIAL

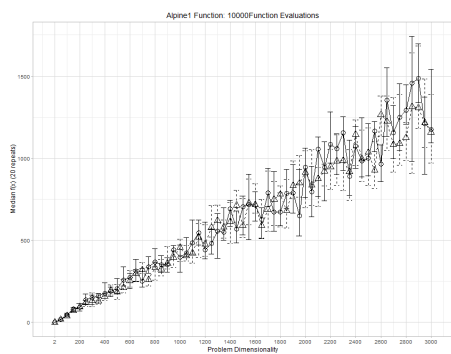
F.1 EXTENDED PSO EXPERIMENT RESULTS



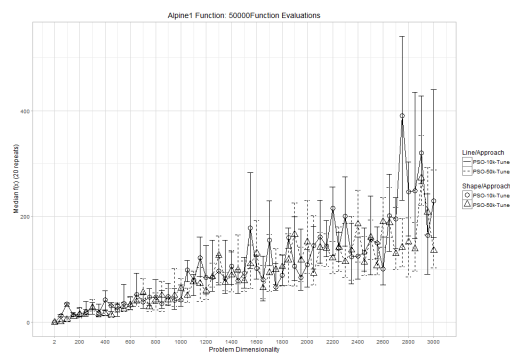
Ackley Function at 10,000 Evaluations



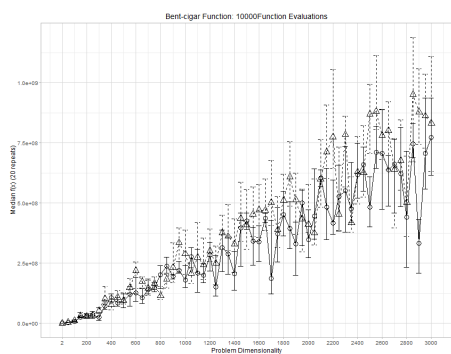
Ackley Function at 50,000 Evaluations



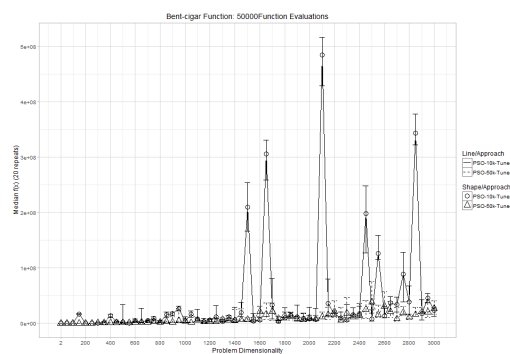
Alpine no.1 Function at 10,000 Evaluations



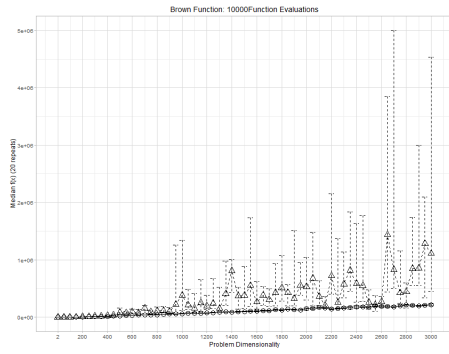
Alpine no.1 Function at 50,000 Evaluations



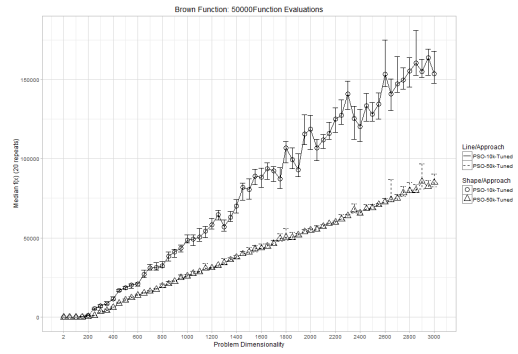
Bent Cigar Function at 10,000 Evaluations



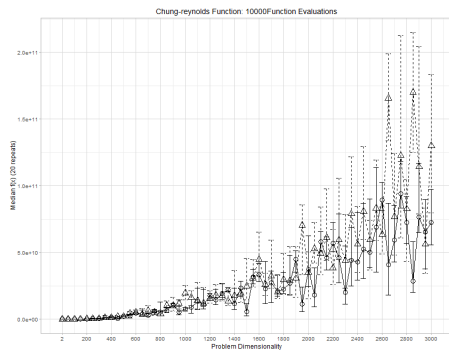
Bent Cigar Function at 50,000 Evaluations



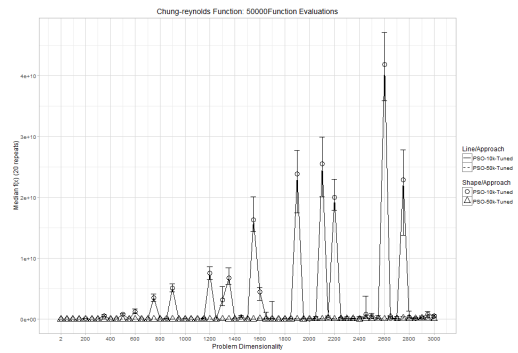
Brown Function at 10,000 Evaluations



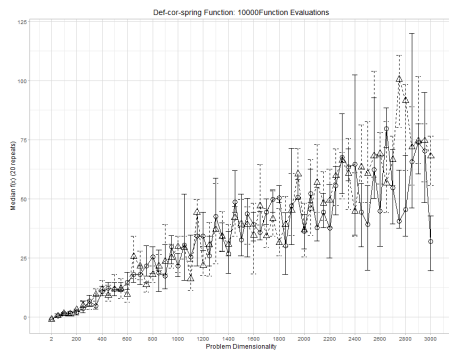
Brown Function at 50,000 Evaluations



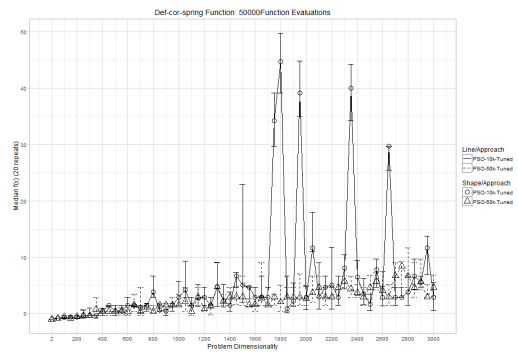
Chung-Reynolds Function at 10,000 Evaluations



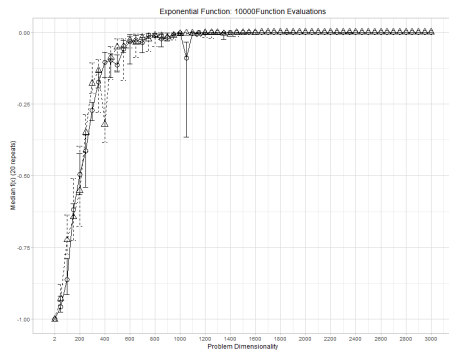
Chung-Reynolds Function at 50,000 Evaluations



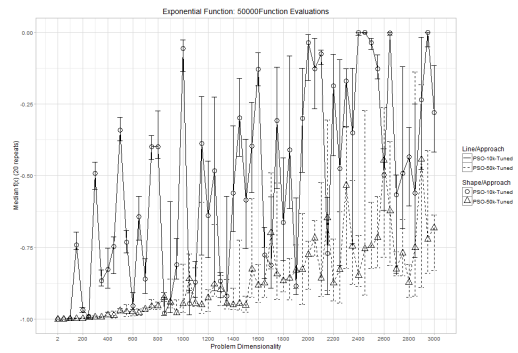
Deflected Corrugated Spring Function at 10,000 Evaluations



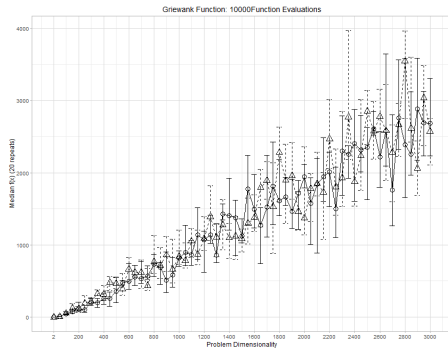
Deflected Corrugated Spring Function at 50,000 Evaluations



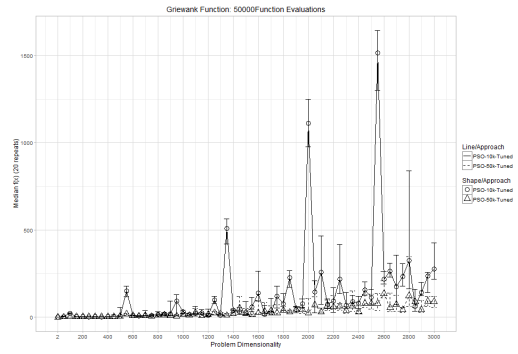
Exponential Function at 10,000 Evaluations



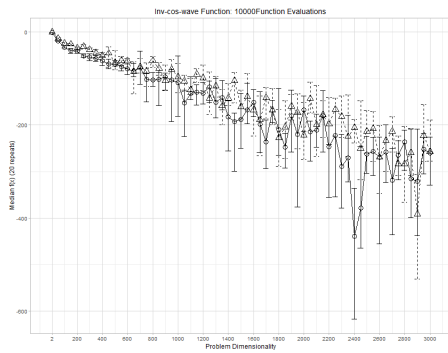
Exponential Function at 50,000 Evaluations



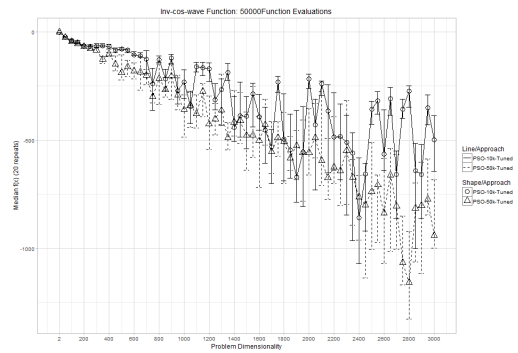
Griewank Function at 10,000 Evaluations



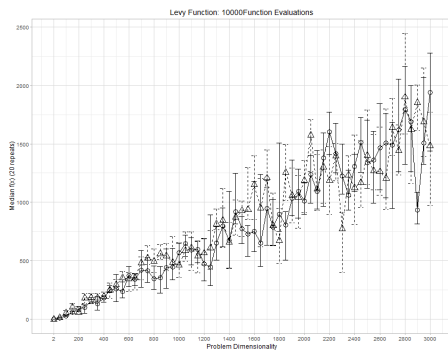
Griewank Function at 50,000 Evaluations



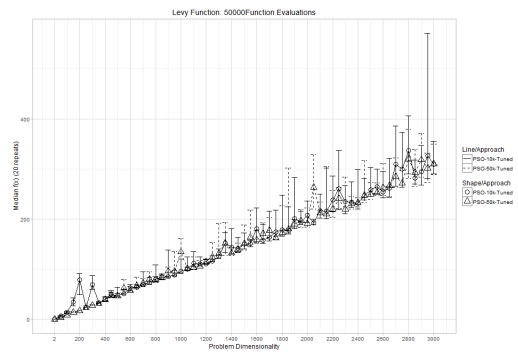
Inverted Cosine Wave Function at 10,000 Evaluations



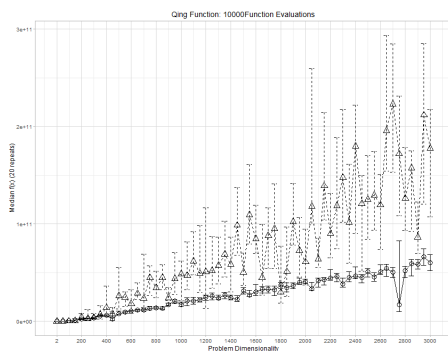
Inverted Cosine Wave Function at 50,000 Evaluations



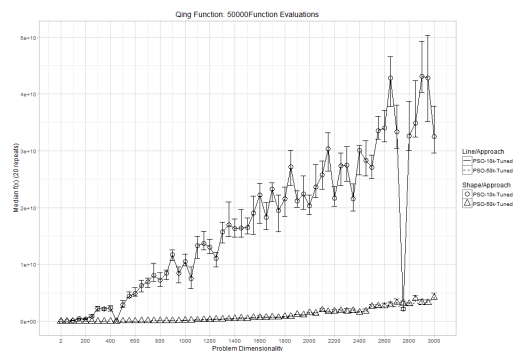
Levy Function at 10,000 Evaluations



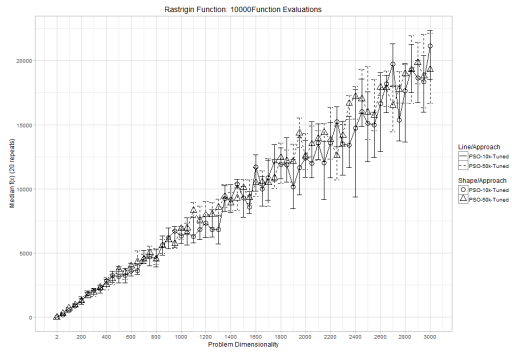
Levy Function at 50,000 Evaluations



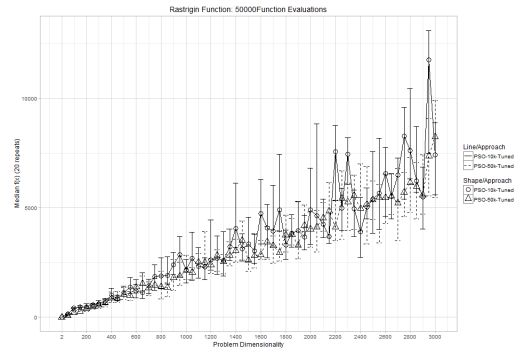
Qing Function at 10,000 Evaluations



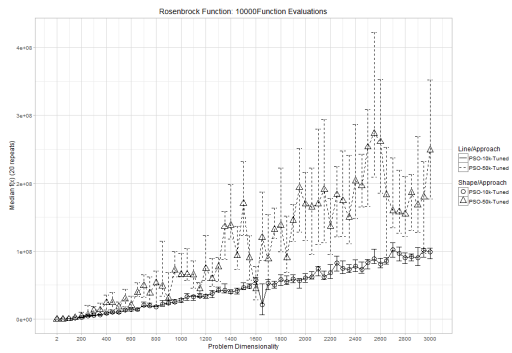
Qing Function at 50,000 Evaluations



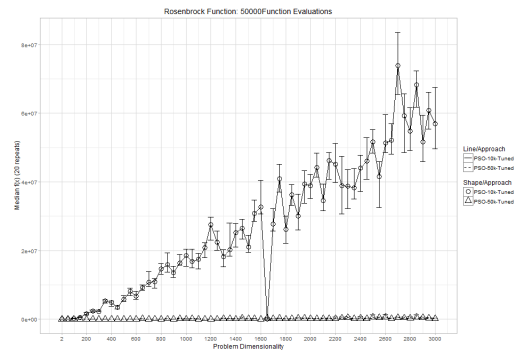
Rastrigin Function at 10,000 Evaluations



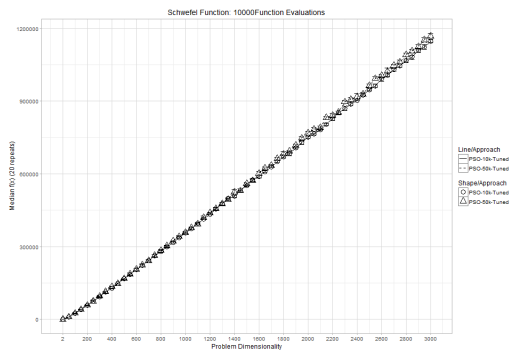
Rastrigin Function at 50,000 Evaluations



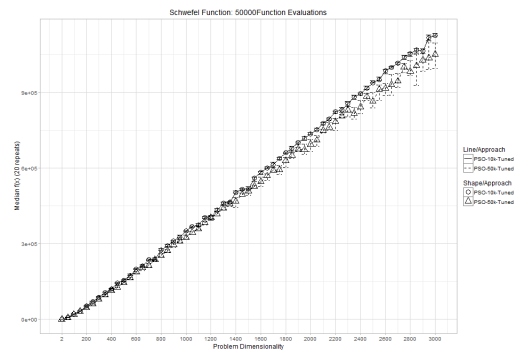
Rosenbrock Function at 10,000 Evaluations



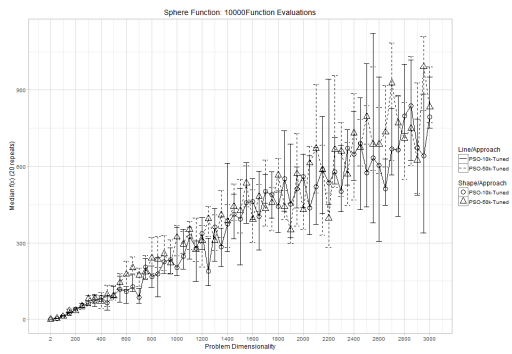
Rosenbrock Function at 50,000 Evaluations



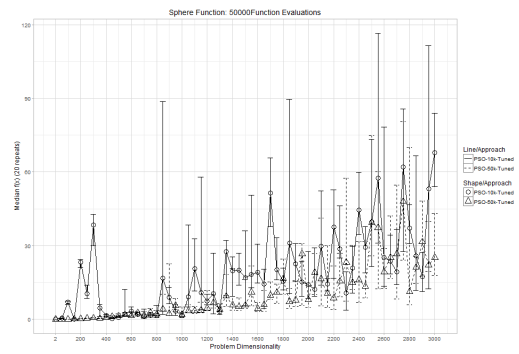
Schwefel Function at 10,000 Evaluations



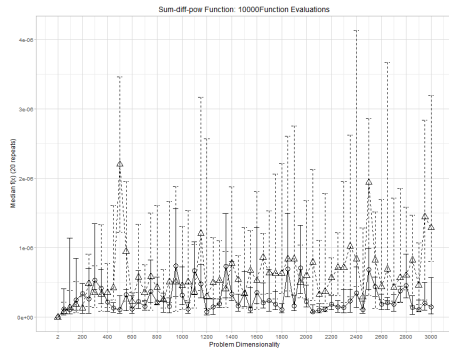
Schwefel Function at 50,000 Evaluations



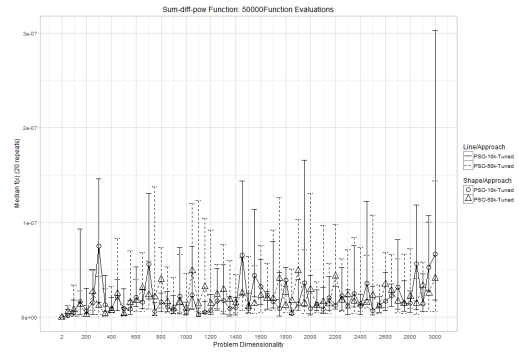
Sphere Function at 10,000 Evaluations



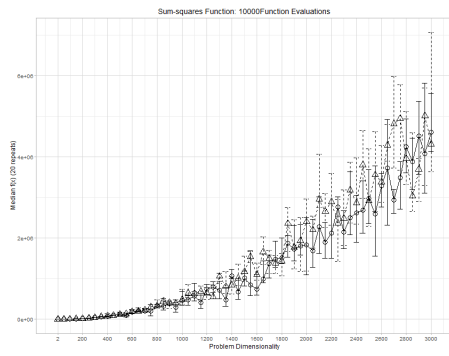
Sphere Function at 50,000 Evaluations



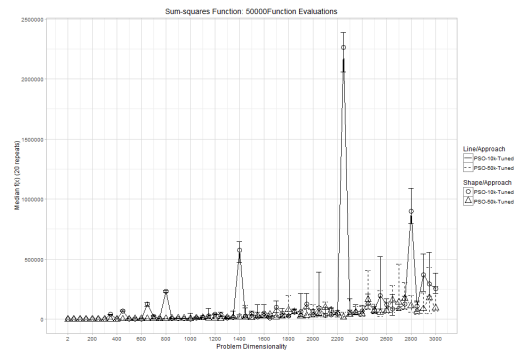
Sum of Different Powers Function at 10,000 Evaluations



Sum of Different Powers Function at 50,000 Evaluations



Sum Squares Function at 10,000 Evaluations



Sum Squares Function at 50,000 Evaluations