

Towards Transparent and Parsimonious Methods for Automatic Performance Tuning

Pedro Bruel

April 2, 2021

Contents

Contents	2
Introduction	9
1 Introduction	9
2 Hardware Design Trends and Consequences for Code Optimization	12
2.1 Historical Hardware Design Trends	14
2.2 Consequences for Code Optimization: Autotuning Loop Nest Optimization . .	17
2.3 Autotuning Approaches and Search Spaces	21
I Optimization Methods for Autotuning	25
3 Notation and Search Spaces	26
3.1 Notation	26
3.2 Search Spaces	28
4 Methods for Function Minimization	32
4.1 Methods Based on Derivatives	32
4.1.1 Gradient Descent	33
4.1.2 Newton's Method	34
4.2 Stochastic Methods	37
4.2.1 Single-State Methods: Random Walk and Simulated Annealing	37
4.2.2 Population-Based Methods: Genetic Algorithms and Particle Swarm Optimization	39

4.3	Summary	42
5	Learning: Building Surrogates	44
5.1	Linear Regression	44
5.1.1	Fitting the Model: The Ordinary Least Squares Estimator	45
5.1.2	Assessing the Quality of Fit of Linear Model Surrogates	47
5.1.3	Inference: Interpreting Significance with ANOVA	50
5.1.4	Linear Models: Interpretable but Biased Surrogates	51
5.2	Gaussian Process Regression	51
5.2.1	Fitting the Model: Posterior Distributions over Functions	54
5.2.2	Assessing the Quality of Fit of Gaussian Process Surrogates	57
5.2.3	Inference: Interpreting Significance with Sobol Indices	60
5.2.4	Gaussian Processes: Flexible but Hard to Interpret Surrogates	60
5.3	Summary	60
6	Design of Experiments	62
6.1	Estimating Linear Effects	62
6.1.1	2-Level Factorial Designs	63
6.1.2	Interactions and Fractional Factorial Designs	68
6.2	Parsimony for Linear Effects with Screening Designs	69
6.3	Increasing Flexibility: Optimal Design	72
6.4	Exploration with Space-Filling Designs	75
6.5	Comparing Sampling Methods	76
6.6	Summary	77
7	Online Learning	79
7.1	Independent Bandits	79
7.2	Expected Improvement for Gaussian Process Regression	79
7.3	Reinforcement Learning	79
7.4	Summary	79
8	Optimization Methods for Autotuning: An Overview	80

Contents

II Applications	82
9 Research Methodology: Efforts for Reproducible Science	83
9.1 Reproducible Research	83
9.2 Computational Documents with Emacs and Org mode	83
9.3 Versioning for Code, Text, and Data	83
10 Compiler Parameters for CUDA Kernels	84
10.1 General-Purpose Computing on NVIDIA GPUs	85
10.1.1 NVIDIA GPU Micro-Architecture	85
10.1.2 Compute Unified Device Architecture (CUDA)	86
10.1.3 GPU Performance Models and Autotuning	86
10.2 Autotuner and Search Space for the NVCC Compiler	86
10.2.1 An Autotuner for CUDA Parameters using OpenTuner	87
10.2.2 Search Space for CUDA Parameters	87
10.3 Target GPUs and Kernels	87
10.3.1 Target GPU Architectures	88
10.3.2 Benchmark of CUDA Algorithms	88
10.4 Performance Improvements and Parameter Clustering Attempt	89
10.4.1 Performance Improvements	90
10.4.2 CUDA Compiler Autotuner Performance	91
10.4.3 Clustering Parameters found by Stochastic Experiments	93
10.5 Assessing Parameter Significance with Screening	93
10.6 Summary	97
11 High-Level Synthesis Kernels for FPGAs	98
11.1 Autotuning High-Level Synthesis for FPGAs	100
11.1.1 Tools for HLS	100
11.1.2 Autotuning for FPGAs	100
11.2 Autotuner and Search Space for the LegUp HLS Compiler	101
11.2.1 Autotuner	101
11.2.2 High-Level Synthesis Parameters	102

11.2.3 HLS Autotuning Metrics	102
11.3 Target Optimization Scenarios and HLS Kernels	103
11.3.1 Optimization Scenarios	103
11.3.2 Kernels	104
11.3.3 Experiments	104
11.4 Performance Improvements using Stochastic Methods	105
11.5 Summary	107
12 An OpenCL Laplacian Kernel	109
12.1 The Laplacian Kernel	110
12.1.1 BOAST Code and the OpenCL Kernel	110
12.2 A Transparent and Parsimonious Approach to Autotuning using Optimal Design	114
12.3 Building a Performance Model	115
12.3.1 Search Space	115
12.3.2 Modeling the Impact of Each Factor	116
Linear Terms	116
Linear plus Inverse Terms	116
The Complete Initial Model	117
12.4 Looking at a Single DLMT Run	117
12.5 Evaluation of Optimization Methods	120
12.6 Summary	122
13 SPAPT Kernels and Orio	124
13.1 Choosing a Design of Experiments Method for Autotuning	124
13.2 The SPAPT Benchmark Suite	125
13.3 Performance Improvements using Optimal Design	126
13.4 Identifying Significant Factors for the <i>bicg</i> Kernel	128
Removing Cubic Terms	130
Reusing Data	130
Quantile Regression	130
Running more Steps	130
13.4.1 Removing Cubic Terms	130

Contents

13.4.2 Reusing Data from All Steps	134
13.4.3 Summary: Tuning SPAPT Kernels with Optimal Design	136
13.5 Autotuning SPAPT kernels with Gaussian Process Regression and Expected Improvement	136
13.5.1 Peak Performance for the <i>DGEMV</i> kernel	139
13.6 Summary	140
14 Mixed-Precision Quantization for <i>ResNet50</i>	142
14.1 Autotuning Bit Precision for Convolutional Neural Networks	143
14.1.1 Further Applications of Autotuning to Neural Networks	144
14.2 ResNet50 and ImageNet	144
14.2.1 ResNet50 Architecture	144
14.2.2 The ImageNet Dataset	146
14.3 Search Space and Objective Function	146
14.4 Optimization Methods	147
14.4.1 Random and Sobol Sampling	147
14.4.2 Gaussian Process Regression with Expected Improvement	147
Respecting the Size Constraint with Sobol Samples	148
Pseudocode for the GPR Method	148
14.4.3 Reinforcement Learning: The Baseline Method	148
Respecting the Size Constraint	148
14.5 Performance Evaluation and Comparison	150
14.5.1 Optimizing <i>Top5</i> for Constrained Total Weight	150
Comparing Policy Variability on Different Samples	155
14.6 Measuring Performance and Detecting Significance for GPR	155
14.6.1 Performance of the GPR Method	156
14.6.2 Detecting Significance with Sobol Indices	156
14.6.3 Perspectives for Further Exploration	157
Running GPR with Data from Past Experiments	157
Optimizing Accuracy and Weight	158
14.7 Summary	160

Conclusion	162
15 Conclusion	162
Bibliography	163
List of Tables	177
List of Figures	179

Introduction

Chapter 1

Introduction

Computer performance has sustained exponential increases over the last half century, despite current physical limits on hardware design. Software optimization has performed an increasingly substantial role on performance improvement over the last 15 years, requiring the exploration of larger and more complex search spaces than ever before.

Autotuning methods are one approach to tackle performance optimization of complex search spaces, enabling exploitation of existing relationships between program parameters and performance. We can derive autotuning methods from well-established statistics, although their usage is not common in or standardized for autotuning domains.

Initial work on this thesis studied the effectiveness of classical and standard search heuristics, such as Simulated Annealing, on autotuning problems. The first target autotuning domain was the set of parameters of a compiler for CUDA programs. The search heuristics for this case study were implemented using the OpenTuner framework [1], and consisted of an ensemble of search heuristics coordinated by a Multi-Armed Bandit algorithm. The autotuner searched for a set of compilation parameters that optimized 17 heterogeneous GPU kernels, from a set of approximately 10^{23} possible combinations of all parameters. With 1.5h autotuning runs we have achieved up to 4 \times speedup in comparison with the CUDA compiler's high-level optimizations. The compilation and execution times of programs in this autotuning domain are relatively fast, and were in the order of a few seconds to a minute. Since measurement costs are relatively small, search heuristics could find good optimizations using as many measurements as needed. A detailed description of this work is available in our paper [2] published in the *Concurrency and Computation: Practice and Experience* journal.

The next case study was developed in collaboration with *Hewlett-Packard Enterprise*, and consisted of applying the same heuristics-based autotuning approach to the configuration of parameters involved in the generation of FPGA hardware specification from source code in the C language, a process called *High-Level Synthesis* (HLS). The main difference from our work with GPU compiler parameters was the time to obtain the hardware specification, which could be in the order of hours for a single kernel.

In this more complex scenario, we achieved up to 2 \times improvements for different hardware metrics using conventional search algorithms. These results were obtained in a simple HLS benchmark, for which compilation times were in the order of minutes. The search space was composed of approximately 10^{123} possible configurations, which is much larger than the search space in our previous work with GPUs. Search space size and the larger measurement cost meant that we did not expect the heuristics-based approach to have the same effectiveness as in the GPU compiler case study. This work was published [3] at the 2017 *IEEE International Conference on ReConfigurable Computing and FPGAs*.

Approaches using classical machine learning and optimization techniques would not scale to industrial-level HLS, where each compilation can take hours to complete. Search space properties also increase the complexity of the problem, in particular its structure composed of binary, factorial and continuous variables with potentially complex interactions. Our results on autotuning HLS for FPGAs corroborate the conclusion that the empirical autotuning of expensive-to-evaluate functions, such as those that appear on the autotuning of HLS, require a more parsimonious and transparent approach, that can potentially be achieved using the DoE methodology. The next section describes our work on applying the DoE methodology to autotuning.

The main contribution of this thesis is a strategy to apply the Design of Experiments methodology to autotuning problems. The strategy is based on linear regression and its extensions, Analysis of Variance (ANOVA), and Optimal Design. The strategy requires the formulation of initial assumptions about the target autotuning problem, which are refined with data collected by efficiently selected experiments. The main objectives are to identify relationships between parameters, suggesting regions for further experimentation in a transparent way, that is, in a way that is supported by statistical tests of significance. The effectiveness of the proposed strategy, and its ability to explain optimizations it finds, are evaluated on autotuning problems in the source-to-source transformation domain. The initial stages of this work resulted in a publication on IEEE/ACM CCGrid [4].

Further efforts in this direction were dedicated to describe precisely what can be learned about search spaces from the application of the methodology, and to refine the differentiation of the approach for the sometimes conflicting objectives of model assessment and prediction. These discussions are presented in Chapter 12.

Because the Design of Experiments methodology requires the specification of a class of initial performance models, the methodology can sometimes achieve worse prediction capabilities when there is considerable uncertainty on initial assumptions about the underlying relationships. Chapter 12 and describe the application to autotuning of Gaussian Process Regression, an approach that trades some explanatory power for a much larger and more flexible class of underlying models. We also evaluate the performance of this approach on the source-to-source transformation problem from Chapter 12, and on larger autotuning problem on the mixed-precision quantization of *Convolutional Neural Network* (CNN) layers, in Chapter 14.

Describe thesis structure

Chapter 2

Hardware Design Trends and Consequences for Code Optimization

High Performance Computing has been a cornerstone of scientific and industrial progress for at least five decades. By paying the cost of increased complexity, software and hardware engineering advances continue to overcome several challenges on the way of the sustained performance improvements observed during the last fifty years. A consequence of this mounting complexity is that reaching the theoretical peak hardware performance for a given program requires not only expert knowledge of specific hardware architectures, but also mastery of programming models and languages for parallel and distributed computing.

If we state performance optimization problems as *search* or *learning* problems, by converting implementation and configuration choices to *parameters* which might affect performance, we can draw from and adapt proven methods from search, mathematical optimization, and statistical learning. The effectiveness of these adapted methods on performance optimization problems varies greatly, and hinges on practical and mathematical properties of the problem and the corresponding *search space*. The application of such methods to the automation of performance tuning for specific hardware, under a set of *constraints*, is named *autotuning*.

Improving performance also relies on gathering application-specific knowledge, which entails extensive experimental costs since, with the exception of linear algebra routines, theoretical peak performance is not always a reachable comparison baseline. When adapting methods for autotuning we must face challenges emerging from practical properties, such as restricted time and cost budgets, constraints on feasible parameter values, and the need to mix *categorical*, *continuous*, and *discrete* parameters. To achieve useful results we must also choose methods that make hypotheses compatible with problem search spaces, such as the existence of *discoverable*, or at least *exploitable*, relationships between parameters and performance. Choosing an autotuning method requires balancing the exploration of a problem, that is, seeking to discover and explain relationships between parameters and performance, and

the exploitation of known or discovered relationships, seeking only to find the best possible performance.

Search algorithms based on machine learning heuristics are not the best candidates for autotuning domains where measurements are lengthy and costly, such as compiling industrial-level FPGA programs, because these algorithms rely on the availability of a large number of measurements. They also assume good optimizations are reachable from a starting position, and that tendencies observed locally in the search space are exploitable. These assumptions are not usually true in common autotuning domains, as shown in the work of Seymour *et al.* [5].

Autotuning search spaces also usually have non-linear constraints and undefined regions, which are also expected to decrease the effectiveness of search based on heuristics and machine learning. An additional downside to heuristics- and machine learning-based search is that, usually, optimization choices cannot be explained, and knowledge gained during optimization is not reusable. The main contribution of this thesis is to study how to overcome the reliance on these assumptions about search spaces, and the lack of explainable optimizations, from the point of view of a *Design of Experiments* (DOE), or *Experimental Design*, methodology to autotuning.

One of the first detailed descriptions and mathematical treatment of DOE was presented by Ronald Fisher [6] in his 1937 book *The Design of Experiments*, where he discussed principles of experimentation, latin square sampling and factorial designs. Later books such as the ones from Jain [7], Montgomery [8] and Box *et al.* [9] present comprehensive and detailed foundations. Techniques based on DOE are *parsimonious* because they allow decreasing the number of measurements required to determine certain relationships between parameters and metrics, and are *transparent* because parameter selections and configurations can be justified by the results of statistical tests.

In DOE terminology, a *design* is a plan for executing a series of measurements, or *experiments*, whose objective is to identify relationships between *factors* and *responses*. While factors and responses can refer to different concrete entities in other domains, in computer experiments factors can be configuration parameters for algorithms and compilers, for example, and responses can be the execution time or memory consumption of a program.

Designs can serve diverse purposes, from identifying the most significant factors for performance, to fitting analytical performance models for the response. The field of DOE encompasses the mathematical formalization of the construction of experimental designs. More practical works in the field present algorithms to generate designs with different objectives and restrictions.

The contributions of this thesis are strategies to apply to program autotuning the DOE methodology, and *Gaussian Process Regression* [10]. This thesis presents background and a high-level view of the theoretical foundations of each method, and detailed discussions of the challenges involved in specializing the general definitions of search heuristics and statistical

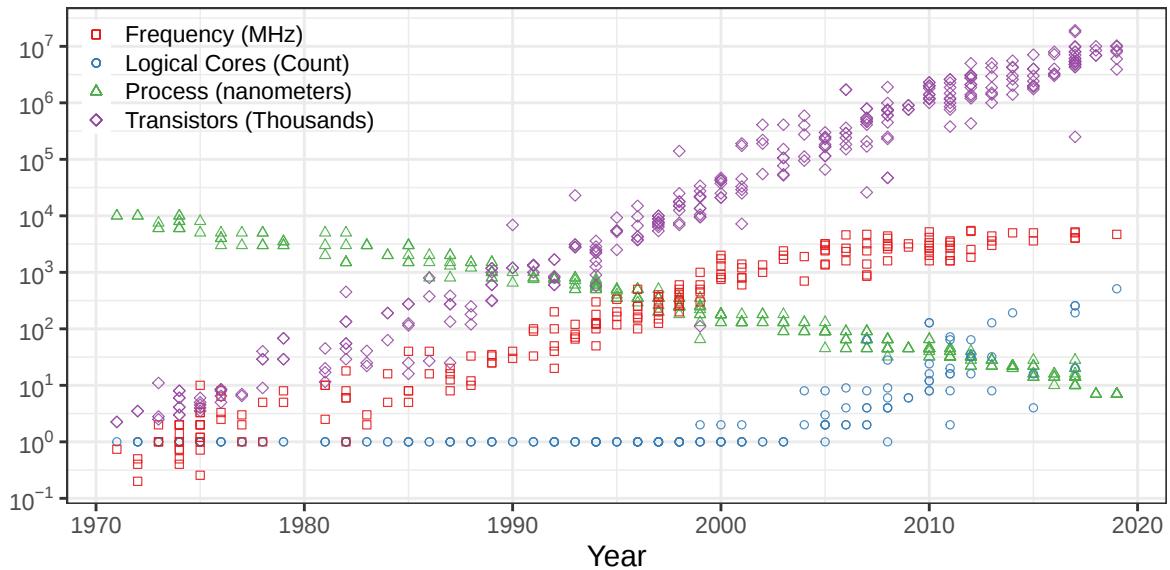


Figure 2.1: 49 years of microprocessor data, highlighting the sustained exponential increases and reductions on transistor counts and fabrication processes, the stagnation of frequency scaling around 2005, and one solution found for it, the simultaneous exponential increase on logical core count. Data from Wikipedia [13, 14]

learning methods to different autotuning problems, as well as what can be *learned* about specific autotuning search spaces, and how that acquired knowledge can be leveraged for further optimization.

This chapter aims to substantiate the claim that autotuning methods have a fundamental role to play on the future of program performance optimization, arguing that the value and the difficulty of the efforts to carefully tune software became more apparent ever since advances in hardware stopped leading to effortless performance improvements, at least from the programmer's perspective. The following sections discuss the historical context for the changes in trends on computer architecture, and characterize the search spaces found when optimizing performance on different domains.

2.1 Historical Hardware Design Trends

The physical constraints imposed by technological advances on circuit design were evident since the first vacuum tube computers that already spanned entire floors, such as the ENIAC in 1945 [11]. The practical and economical need to fit more computing power into real estate is one force for innovation in hardware design that spans its history, and is echoed in modern supercomputers, such as the *Summit* from *Oak Ridge National Laboratory* [12], which spans an entire room.

Figure 2.1 highlights the unrelenting and so far successful pursuit of smaller transistor fabrication processes, and the resulting capability to fit more computing power on a fixed

2.1. Historical Hardware Design Trends

chip area. This trend was already observed in integrated circuits by Gordon Moore *et al.* in 1965 [15], who also postulated its continuity. The performance improvements produced by the design efforts to make Moore's forecast a self-fulfilling prophecy were boosted until around 2005 by the performance gained from increases in circuit frequency.

Robert Dennard *et al.* remarked in 1974 [16] that smaller transistors, in part because they generate shorter circuit delays, decrease the energy required to power a circuit and enable an increase in operation frequency without breaking power usage constraints. This scaling effect, named *Dennard's scaling*, is hindered primarily by leakage current, caused by quantum tunneling effects in small transistors. Figure 2.1 shows a marked stagnation on frequency increase after around 2005, as transistors crossed the 10²nm fabrication process. It was expected that leakage due to tunneling would limit frequency scaling strongly, even before the transistor fabrication process reached 10nm [17].

Current hardware is now past the effects of Dennard's scaling. The increase in logical cores around 2015 can be interpreted as preparation for and mitigation of the end of frequency scaling, and ushered in an age of multicore scaling. Still, in order to meet power consumption constraints, up to half of a multicore processor could have to be powered down, at all times. This phenomenon is named *Dark Silicon* [18], and presents significant challenges to current hardware designers and programmers [19, 20, 21].

The *Top500* [22] list gathers information about commercially available supercomputers, and ranks them by performance on the *LINPACK* benchmark [23]. Figure 2.2 shows the peak theoretical performance *RPeak*, and the maximum performance achieved on the LINPACK benchmark *RMax*, in *Tflops/s*, for the top-ranked supercomputers on TOP500. Despite the smaller performance gains from hardware design that are to be expected for post-Dennard's scaling processors, the increase in computer performance has sustained an exponential climb, sustained mostly by software improvements.

Although *hardware accelerators* such as GPUs and FPGAs, have also helped to support exponential performance increases, their use is not an escape from the fundamental scaling constraints imposed by current semiconductor design. Figure 2.3 shows the increase in processor and accelerator core count on the top-ranked supercomputers on Top500. Half of the top-ranked supercomputers in the last decade had accelerator cores and, of those, all had around ten times more accelerator than processor cores. The apparent stagnation of core count in top-ranked supercomputers, even considering accelerators, highlights the crucial impact software optimization has on performance.

Advances in hardware design are currently not capable of providing performance improvements via frequency scaling without dissipating more power than the processor was designed to support, which violates power constraints and risks damaging the circuit. From the programmer's perspective, effortless performance improvements from hardware have not been expected for quite some time, and the key to sustaining historical trends in performance scaling has lied in accelerators, parallel and distributed programming libraries, and fine tun-

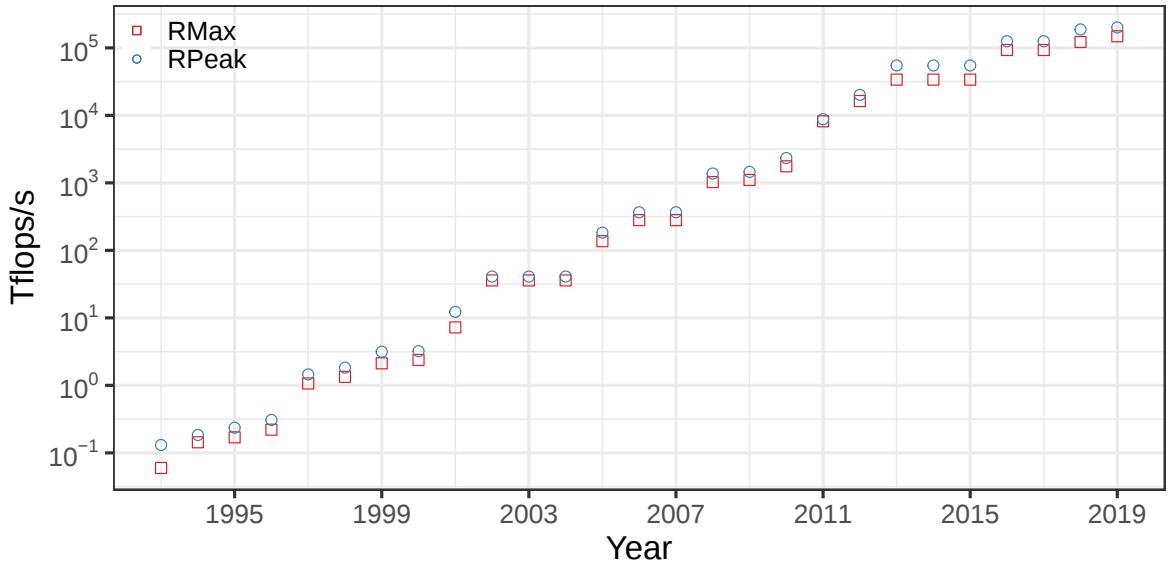


Figure 2.2: Sustained exponential increase of theoretical *RPeak* and achieved *RMax* performance for the supercomputer ranked 1st on TOP500 [22]

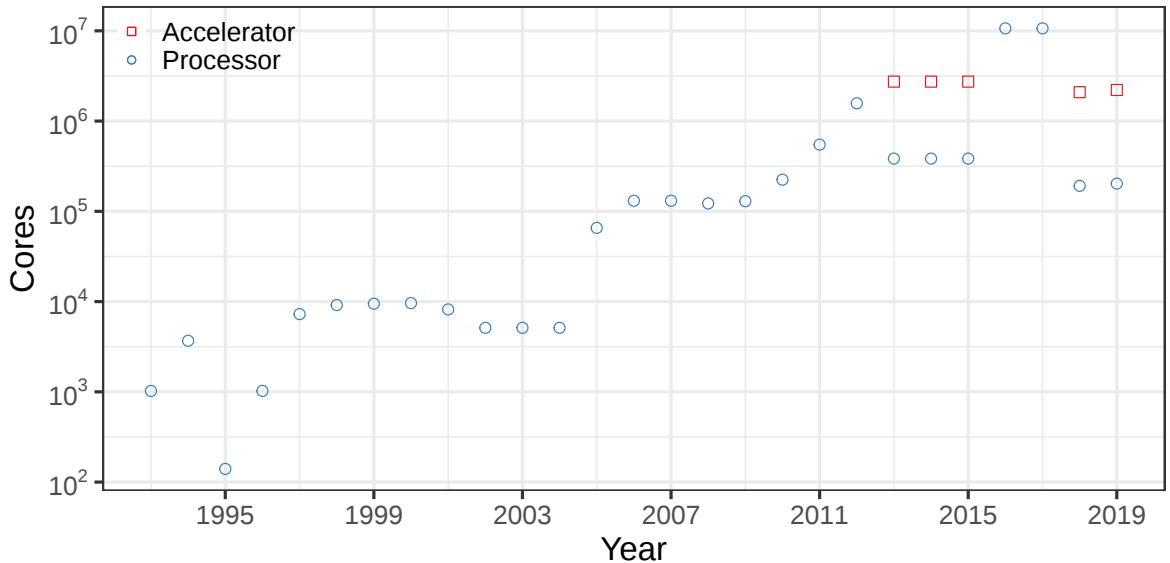


Figure 2.3: Processor and accelerator core count in supercomputers ranked 1st on TOP500 [22]. Core count trends for supercomputers are not necessarily bound to processor trends observed on Figure 2.1.

ing of several stages of the software stack, from instruction selection to the layout of neural networks.

The problem of optimizing software for performance presents its own challenges. The search spaces that emerge from autotuning problems grow quickly to a size for which it would take a prohibitive amount of time to determine the best configuration by exhaustively evaluating all possibilities. Although this means we must seek to decrease the amount of

possibilities, by restricting allowed parameter values, or dropping parameters completely, it is often unclear how to decide which parameters should be restricted or dropped. The next sections introduce a simple autotuning problem, present an overview of the magnitude of the dimension of autotuning search spaces, and briefly introduce the methods commonly used to explore search spaces, some of which are discussed in detail in Chapter I.

2.2 Consequences for Code Optimization: Autotuning Loop Nest Optimization

Algorithms for linear algebra problems are fundamental to scientific computing and statistics. Therefore, decreasing the execution time of algorithms such as general matrix multiplication (GEMM) [24], and others from the original BLAS [25], is an interesting and well motivated example, that we will use to introduce the autotuning problem.

One way to improve the performance of such linear algebra programs is to exploit cache locality by reordering and organizing loop iterations, using source code transformation methods such as loop *tiling*, or *blocking*, and *unrolling*. We will now briefly describe loop tiling and unrolling for a simple linear algebra problem. After, we will discuss an autotuning search space for blocking and unrolling applied to GEMM, and how these transformations generate a relatively large and complex search space, which we can explore using autotuning methods.

Figure 2.4 shows three versions of code in the C language that, given three square matrices A , B , and C , computes $C = C + A + B^T$. The first optimization we can make is to preemptively load to cache, or *prefetch*, as many as possible of the elements we know will be needed at any given iteration, as is shown in Figure 2.4a. The shaded elements on the top row of Figure 2.5 represent the elements that could be prefetched in iterations of Figure 2.4a.

Since C matrices are stored in *row-major order*, each access of an element of B forces loading the next row elements, even if we explicitly prefetch a column of B . Since we are accessing B in a *column-major order*, the prefetched row elements would not be used until we reached the corresponding column. Therefore, the next column elements will have to be loaded at each iteration, considerably slowing down the computation.

We can solve this problem by reordering memory accesses to request only prefetched elements. It suffices to adequately split loop indices into blocks, as shown in Figure 2.4b. Now, memory accesses are performed in *tiles*, as shown on the bottom row of Figure 2.5. If blocks are correctly sized to fit in cache, we can improve performance by explicitly prefetching each tile. After blocking, we can still improve performance by *unrolling* loop iterations, which forces register usage and helps the compiler to identify regions that can be *vectorized*. A conceptual implementation of loop unrolling is shown in Figure 2.4c.

Looking at the loop nest optimization problem from the autotuning perspective, the two *parameters* that emerge from the implementations are the *block size*, which controls the stride,

```

int N = 256;
float A[N][N], B[N][N], C[N][N];
int i, j;
// Initialize A, B, C
for(i = 0; i < N; i++){
    // Load line i of A to fast memory
    for(j = 0; j < N; j++){
        // Load C[i][j] to fast memory
        // Load column j of B to fast memory
        C[i][j] += A[i][j] + B[j][i];
        // Write C[i][j] to main memory
    }
}
    
```

(a) Regular implementation

```

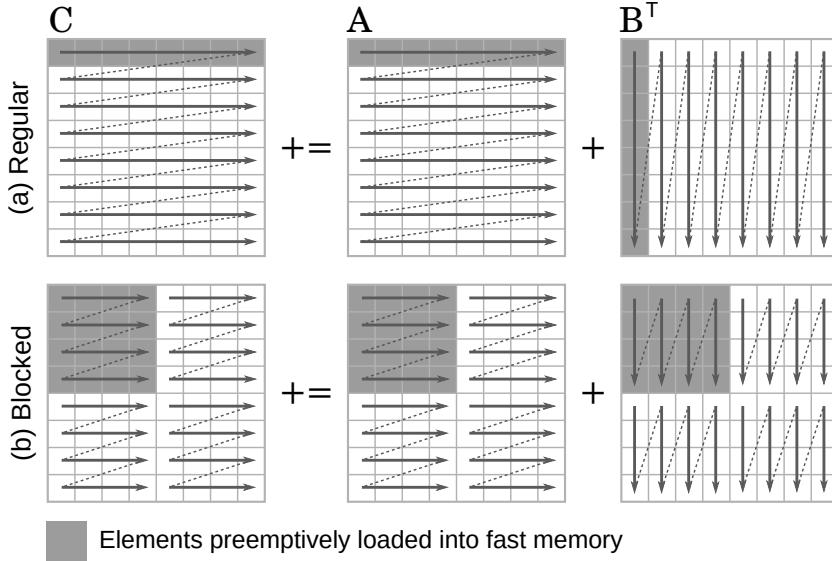
int N = 256;
int B_size = 4;
int A[N][N], B[N][N], C[N][N];
int i, j, x, y;
// Initialize A, B, C
for(i = 0; i < N; i += B_size){
    for(j = 0; j < N; j += B_size){
        // Load block (i, j) of C to fast memory
        // Load block (i, j) of A to fast memory
        // Load block (j, i) of B to fast memory
        for(x = i; x < min(i + B_size, N); x++){
            for(y = j; y < min(j + B_size, N); y++){
                C[x][y] += A[x][y] + B[y][x];
            }
        }
        // Write block (i, j) of C to main memory
    }
}
    
```

(b) Blocked, or tiled

```

int N = 256;
int B_size = 4;
int A[N][N], B[N][N], C[N][N];
int i, j, k;
// Initialize A, B, C
for(i = 0; i < N; i += B_size){
    for(j = 0; j < N; j += B_size){
        // Load block (i, j) of C to fast memory
        // Load block (i, j) of A to fast memory
        // Load block (j, i) of B to fast memory
        C[i + 0][j + 0] += A[i + 0][j] * B[i][j + 0];
        C[i + 0][j + 1] += A[i + 0][j] * B[i][j + 1];
        // Unroll the remaining 12 iterations
        C[i + Bsize - 1][j + B_size - 2] += A[i + Bsize - 1][j] * B[i][j + B_size - 2];
        C[i + Bsize - 1][j + B_size - 1] += A[i + Bsize - 1][j] * B[i][j + B_size - 1];
        // Write block (i, j) of C to main memory
    }
}
    
```

(c) Tiled and unrolled

 Figure 2.4: Loop nest optimizations for $C = C + A + B^T$, in C

 Figure 2.5: Access patterns for matrices in $C = C + A + B^T$, with loop nest optimizations. Panel (a) shows the access order of a regular implementation, and panel (b) shows the effect of loop tiling, or blocking

and the *unrolling factor*, which controls the number of unrolled iterations. Larger block sizes are desirable, because we want to avoid extra comparisons, but blocks should be small enough

2.2. Consequences for Code Optimization: Autotuning Loop Nest Optimization

```

int N = 256;
float A[N][N], B[N][N], C[N][N];
int i, j, k;
// Initialize A, B, C
for(i = 0; i < N; i++){
    // Load line i of A to fast memory
    for(j = 0; j < N; j++){
        // Load C[i][j] to fast memory
        // Load column j of B to fast memory
        for(k = 0; k < N; k++){
            C[i][j] += A[i][k] * B[k][j];
        }
        // Write C[i][j] to main memory
    }
}

```

(a) Regular implementation

```

int N = 256;
int B_size = 4;
int A[N][N], B[N][N], C[N][N];
int i, j, k;
// Initialize A, B, C
for(i = 0; i < N; i += B_size){
    for(j = 0; j < N; j += B_size){
        // Load block (i, j) of C to fast memory
        for(k = 0; k < N; k++){
            // Load block (i, k) of A to fast memory
            // Load block (k, j) of B to fast memory
            for(x = i; x < min(i + B_size, N); x++){
                for(y = j; y < min(j + B_size, N); y++){
                    C[x][y] += A[x][k] * B[k][y];
                }
            }
        }
        // Write block (i, j) of C to main memory
    }
}

```

(b) Blocked, or tiled

```

int N = 256;
int B_size = 4;
int A[N][N], B[N][N], C[N][N];
int i, j, k;
// Initialize A, B, C
for(i = 0; i < N; i += B_size){
    for(j = 0; j < N; j += B_size){
        // Load block (i, j) of C to fast memory
        for(k = 0; k < N; k++){
            // Load block (i, k) of A to fast memory
            // Load block (k, j) of B to fast memory
            C[i + 0][j + 0] += A[i + 0][k] * B[k][j + 0];
            C[i + 0][j + 1] += A[i + 0][k] * B[k][j + 1];
            // Unroll the remaining 12 iterations
            C[i + Bsize - 1][j + B_size - 2] += A[i + Bsize - 1][k] * B[k][j + B_size - 2];
            C[i + Bsize - 1][j + B_size - 1] += A[i + Bsize - 1][k] * B[k][j + B_size - 1];
        }
        // Write block (i, j) of C to main memory
    }
}

```

(c) Tiled and unrolled

Figure 2.6: Loop nest optimizations for GEMM, in C

to ensure access to as few as possible out-of-cache elements. Likewise, the unrolling factor should be large, to leverage vectorization and available registers, but not so large that it forces memory to the stack.

The values of block size and unrolling factor that optimize performance will depend on the cache hierarchy, register layout, and vectorization capabilities of the target processor, but also on the memory access pattern of the target algorithm. In addition to finding the best values for each parameter independently, an autotuner must ideally aim to account for the *interactions* between parameters, that is, for the fact that the best value for each parameter might also depend on the value chosen for the other.

The next loop optimization example comes from Seymour *et al.* [5], and considers 128 blocking and unrolling values, in the interval [0, 127], for the GEMM algorithm. The three panels of Figure 2.6 show conceptual implementations of loop blocking and unrolling for GEMM in C. A block size of zero results in the implementation from Figure 2.6a, and an unrolling factor of zero performs a single iteration per condition check.

It is straightforward to change the block size of the implementations from Figure 2.6, but the unrolling factor is not exposed as a parameter. To test different unrolling values we need

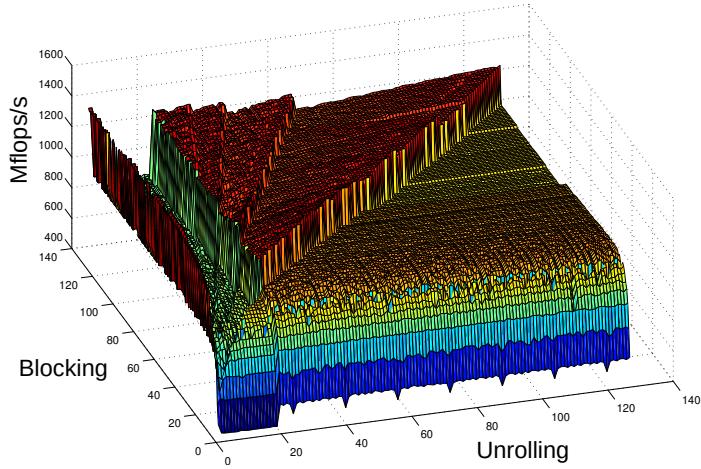


Figure 2.7: An exhaustively measured search space, defined by loop blocking and unrolling parameters, for a sequential GEMM kernel. Reproduced from Seymour *et al.* [5]

to generate new versions of the source code with different numbers of unrolled iterations. We can do that with code generators or with *source-to-source transformation* tools [26, 27, 28]. It is often necessary to modify the program we wish to optimize in order to provide a configuration interface and expose its implicit parameters. Once we are able to control the block size and the loop unrolling factor, we determine the target search space by choosing the values to be explored.

In this example, the search space is defined by the $128^2 = 16384$ possible combinations of blocking and unrolling values. The performance of each combination in the search space, shown in *Mflops/s* in Figure 2.7, was measured for a sequential GEMM implementation, using square matrices of size 400 [5]. We can represent this autotuning search space as a 3D *landscape*, since we have two configurable parameters and a single target performance metric. In this setting, the objective is to find the *highest* point, since the objective is to *maximize* Mflops/s, although usually the performance metric is transformed so that the objective is its *minimization*.

On a first look, there seems to be no apparent global search space structure in the landscape on Figure 2.7, but local features jump to the eyes, such as the “valley” across all block sizes for low unrolling factors, the “ramp” across all unrolling factors for low block sizes, and the series of jagged “plateaus” across the middle regions, with ridges for identical or divisible block sizes and unrolling factors. A careful look reveals also that there is a curvature along the unrolling factor axis. Also of note is the abundance in this landscape of *local minima*, that is, points with relatively good performance, surrounded by points with worse performance. By exhaustively evaluating all possibilities, the original study determined that the best performance on this program was achieved with a block size of 80 and an unrolling factor of 2.

In this conceptual example, all $\approx 1.64 \times 10^4$ configurations were exhaustively evaluated, but it is impossible to do so in most settings where autotuning methods are useful. The

next section provides a perspective of the autotuning domains and methods employed in current research, presenting a selection of search spaces and discussing the trends that can be observed on search space size, targeted HPC domains, and chosen optimization methods.

2.3 Autotuning Approaches and Search Spaces

Autotuning methods have been used to improve performance in an increasingly large variety of domains, from the earlier applications to linear algebra subprograms, to the now ubiquitous construction and configuration of neural networks, to the configuration of the increasingly relevant tools for the re-configurable hardware of FPGAs. In this setting, it is not far-fetched to establish a link between the continued increases in performance and hardware complexity, that we discussed previously in this chapter, to the increases in dimension and size of the autotuning problems that we can now tackle.

Figure 2.8 presents search space dimension, measured as the number of parameters involved, and size, measured as the number of possible parameter combinations, for a selection of search spaces from 14 autotuning domains. Precise information about search space characterization is often missing from works on autotuning methods and applications. The characterization of most of the search spaces in Figure 2.8 was obtained directly from the text of the corresponding published paper, but for some it was necessary to extract characterizations from the available source code. Still, it was impossible to obtain detailed descriptions of search spaces for many of the published works on autotuning methods and applications, and in that way the sample shown in this section is biased, because it contains only information on works that provided it.

The left hand panel of Figure 2.8 shows search spaces with up to 60 parameters. The over-representation of search spaces for linear algebra domains in this sample stands out on the left hand panel, but the domain is not present on the remaining portion of the sample, shown on the right hand panel. The largest search spaces for which we were able to find information on published work are defined for the domains of neural network configuration, High-Level Synthesis for FPGAs, compiler parameters, and domain-specific languages.

None of the largest search spaces in this sample, that is, the ones outside the zoomed area of the left hand panel, come from works earlier than 2009. The sustained performance improvements we discussed previously have enabled and pushed autotuning research toward progressively larger problems, which has also been done to most research areas. Increased computing power has made it feasible, or at least tolerable, to apply search heuristics and statistical learning methods to find program configurations that improve performance.

It is straightforward to produce an extremely large autotuning search space. Compilers have hundreds of binary flags that can be considered for selection, generating a large set of combinations. Despite that, regarding performance improvements, it is likely that most configuration parameters will have a small impact, that is, that only a handful of parameters



Figure 2.8: Dimension and search space size for autotuning problems from 14 domains [29, 1, 30, 31, 32, 4, 33, 2, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 5] The left panel shows a zoomed view of the right panel

are responsible for changes in performance. Search spaces are often much more restrictive than the one we discussed in Section 2.2. Autotuning problem definitions usually come with *constraints* on parameter values and limited *experimental budgets*, and *runtime* failures for some configurations are often unpredictable. In this context, finding configurations that improve performance and determining the subset of *significant parameters* are considerable challenges.

Search heuristics, such as methods based on genetic algorithms and gradient descent, are a natural way to tackle these challenges because they consist of procedures for exploiting existing and unknown relationships between parameters and performance without making or requiring *explicit hypotheses* about the problem. Despite that, most commonly used heuristics make *implicit hypotheses* about search spaces which are not always verified, such as assuming that good configurations are *reachable* from a random starting point.

Autotuning methods that make explicit hypotheses about the target program, such as methods based on Design of Experiments, require some initial knowledge, or willingness to make assumptions, about underlying relationships, and are harder to adapt to constrained scenarios, but have the potential to produce *explainable* optimizations. In general, methods based on Machine Learning have enough flexibility to perform well in complex search spaces and make few assumptions about problems, but usually provide little, if any, that can be used to explain optimization choices or derive relationships between parameters and performance.

Table 2.1 lists some autotuning systems, their target domains, and the employed method, ordered by publication date. Some systems that did not provide detailed search space descriptions and could not be included in Figure 2.8, especially some of the earlier work, provided enough information to categorize their autotuning methods. In contrast, many more recent

Table 2.1: Autotuning methods used by a sample of systems, in different domains, ordered by publishing year. Methods were classified as either Search Heuristics (SH), Machine Learning (ML), or more precisely when the originating work provided detailed information. Earlier work favored employing Search Heuristics, which are less prominent in recent work, which favors methods based on Machine Learning.

System	Domain	Method	Year
PhiPAC [46]	Linear Algebra	SH (Exhaustive)	1997
ATLAS [47]	Linear Algebra	SH (Exhaustive)	1998
FFTW [48]	Digital Signal Processing	SH (Exhaustive)	1998
Active Harmony [49]	Domain-Specific Language	SH	2002
OSKI [50]	Linear Algebra	SH	2005
Seymour, K. <i>et al.</i> [5]	Linear Algebra	SH	2008
PRO [37]	Linear Algebra	SH	2009
ParamILS [38]	Combinatorial Auctions	SH	2009
PetaBricks [28]	Domain-Specific Language	SH (Genetic Algorithm)	2009
MILEPOST GCC [51]	Compiler Parameters	ML	2011
Orio [29]	Linear Algebra	ML (Decision Trees)	2012
pOSKI [30]	Linear Algebra	SH	2012
INSIEME [52]	Compiler Parameters	SH (Genetic Algorithm)	2012
OpenTuner [1]	Compiler Parameters	SH	2014
Lgen [53]	Linear Algebra	SH	2014
OPAL [54]	Parallel Computing	SH	2014
Mametjanov, A. <i>et al.</i> [34]	High-Level Synthesis	ML (Decision Trees)	2015
CLTune [55]	Parallel Computing	SH	2015
Guerreirro, J. <i>et al.</i> [56]	Parallel Computing	SH	2015
Collective Mind [57]	Compiler Parameters	ML	2015
Abdelfattah, A. <i>et al.</i> [35]	Linear Algebra	SH (Exhaustive)	2016
TANGRAM [58]	Domain-Specific Language	SH	2016
MASE-BDI [59]	Environmental Land Change	SH	2016
Xu, C. <i>et al.</i> [36]	High-Level Synthesis	SH	2017
Apollo [60]	Parallel Computing	ML (Decision Trees)	2017
DeepHyper [32]	Neural Networks	ML (Decision Trees)	2018
Tuzov, I. <i>et al.</i> [40]	High-Level Synthesis	Design of Experiments	2018
Periscope [42]	Compiler Parameters	SH	2018
SynTunSys [41]	High-Level Synthesis	SH	2019
Kwon, J. <i>et al.</i> [43]	High-Level Synthesis	ML	2019
FancyTuner [44]	Compiler Parameters	SH	2019
Ol'ha, J. <i>et al.</i> [45]	Parallel Computing	Sensitivity Analysis	2019
Petrovic, F. <i>et al.</i> [31]	Linear Algebra	SH	2020
Chu, Y. <i>et al.</i> [39]	Search/MVWCP	SH	2020

works, especially those using methods based on Machine Learning, did not provide specific method information. Earlier work often deals with search spaces small enough to exhaustively evaluate, and using search heuristics to optimize linear algebra programs is the most prominent category of earlier work in this sample. Later autotuning work target more varied domains, with the most prominent domains in this sample being parallel computing, compiler

parameters, and High-Level Synthesis. Systems using methods based on Machine Learning become more common on later work than systems using heuristics.

Chapter 4 provides more detailed definitions and discussions of the applicability, effectiveness, and explanatory power of stochastic autotuning methods based on search heuristics. The remainder of this chapter details the contributions of this thesis and the structure of this document.

Third Chapter, or next Section: Proposal of this thesis, list of contributions

Part I

Optimization Methods for Autotuning

Chapter 3

Notation and Search Spaces

The following chapters present the optimization methods that we applied to autotuning problems in different domains. Each chapter presents a group of methods, briefly discussing each method in the group and their underlying hypotheses. The objective of each chapter is to provide high-level descriptions of optimization methods, presenting clear definitions tied to the autotuning context.

Each chapter concludes with a discussion of the applicability of each group of methods to autotuning problems. The methods we discuss have significant differences but employ basic concepts. We will use the same mathematical notation to discuss all methods when possible. The remainder of this chapter presents common basic concepts and the associated notation to be used in subsequent chapters, which are summarized in Table 3.1. Completing this introduction to optimization in the context of autotuning, the chapter ends with a discussion of common search space properties.

3.1 Notation

We will call *optimization* the minimization of a real-valued function with a single vector input. For a function $f : \mathcal{X} \rightarrow \mathbb{R}$, we wish to find the input vector $\mathbf{x}^* = [x_1 \dots x_p]^\top$ in the *parameter space*, or *search space* \mathcal{X} for which $f(\mathbf{x}^*)$ is the smallest, compared to all other $\mathbf{x} \in \mathcal{X}$. The function f represents, in the autotuning context, the *performance metric* we wish to optimize, such as the execution time of some application, and the parameter space \mathcal{X} represents the set of possible *configurations* we can explore, such as compiler flags. Therefore, we define optimization in the autotuning context as finding the configuration that minimizes the target performance metric. For the sake of simplicity, we assume we can use the opposite of the performance metrics that should be maximized.

As an example of an autotuning problem, consider optimizing the choice of flags for a compiler with $p = |\mathbf{x}|$ flags, where a configuration $\mathbf{x} = [x_1 \dots x_p]^\top$ consists of a vector of

p boolean values, denoting whether each flag x_1, \dots, x_p is turned on for the compilation of a specific application. To find the compiler configuration that generates the binary with the smallest execution time, we conduct a set of $n = |\mathbf{X}|$ experiments, chosen according to some criterion, generating the $n \times p$ experimental design matrix $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^\top$. Each experiment consists of compiling the target application using the specified compiler configuration, and measuring the execution time of the resulting binary.

In this example, evaluating $f(\mathbf{x})$ involves generating the compiler configuration corresponding to the vector \mathbf{x} of selected flags. This involves writing a shell command or a configuration file, running the configured compiler, checking for compilation errors, measuring the execution time of the resulting binary, and verifying the correctness of its output.

In practice, we may never be able to observe the *true value* of $f(\mathbf{x})$. In fact, empirical tests of this nature are always subject to unknown or uncontrollable effects, and to inherent imprecision in measurement. In practice, we settle for observing $y = f(\mathbf{x}) + \varepsilon$, where ε encapsulates all unknown and uncontrollable effects, as well as the measurement error. Returning to the compiler flag example, suppose that we could conduct $n = |\mathcal{X}| = 2^p$ experiments, measuring the performance of the binaries generated with all possible flag combinations. With such experimental design we would obtain the measurements

$$\mathbf{y} = [y_i = f(\mathbf{x}_i) + \varepsilon_i, i = 1, \dots, 2^p]^\top. \quad (3.1)$$

The measurement y_i is an *estimate* of $f(\mathbf{x}_i)$, with error ε_i . If the error is reasonably small, an estimate of the *global optimum* \mathbf{x}^* in this example is the \mathbf{x}_i that produces the binary with the smallest estimated execution time y^* , the smallest $y_i \in \mathbf{y}$.

Assuming we are capable of cheaply evaluating f for a large set of experiments \mathbf{X} , and that we are not interested in building statistical models for the performance of our application, we can directly optimize f using stochastic descent methods, or gradient- and hessian-based methods if f is suitably convex and differentiable. These function minimization methods are discussed in Chapter 4.

If we are not capable of directly measuring f , if it is unreasonably expensive or time-consuming to do so, or if constructing statistical performance models for our application is of crucial importance, we can employ the surrogate-based methods discussed in Chapters 5, 6, and 7. These methods use different strategies to construct a *surrogate model*

$$\hat{f}_\theta : \mathcal{X} \rightarrow \mathbb{R}, \text{ with } \theta \in \Theta, \quad (3.2)$$

where θ is a parameter vector usually estimated from measurements (\mathbf{X}, \mathbf{y}) . The function $\theta : \mathcal{X}^n \times \mathbb{R}^n \rightarrow \Theta$ represents the *estimation process*, that uses the *observations* $\mathbf{y} \in \mathbb{R}^n$ from an experimental design $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^\top$. The *parameter vector* $\theta(\mathbf{X}, \mathbf{y})$ in the *parameter space* Θ will be used to compute the estimate.

The constructed surrogate model $\hat{f}_{\theta}(\mathbf{x}, \mathbf{y})$ can be used as a tool to attempt to describe and optimize the underlying real function f , provided we are able to construct a useful estimate of $\theta(\mathbf{X}, \mathbf{y})$. We will discuss in Chapters 5 and 6 methods that use the m individual parameter estimates $\theta(\mathbf{X}, \mathbf{y}) = [\theta_1 \dots \theta_m]^\top$ to assess the significance of specific factors.

If we use the Ordinary Least Squares (OLS) estimator for θ in our compiler flag example, the parameter estimates $\theta_2, \dots, \theta_m$ could correspond to each one of the p flags. In this case, θ_1 is the estimate of the intercept of the linear model, and $m = p + 1$. Optimization methods using the Bayesian inference framework, such as Gaussian Process Regression, which we discuss in Chapter 5, associate a probability distribution to $\theta(\mathbf{X}, \mathbf{y})$, which propagates with \hat{f}_{θ} and can be exploited in our optimization context.

Chapter 5 differentiates between parametric methods which make the hypothesis that the number $m = |\theta|$ of estimated parameters is finite and often interpretable, and nonparametric methods, which operate in parameter spaces of infinite dimension.

Table 3.1 summarizes the notation and concepts we have discussed so far, tying those concepts to the compiler flag example we used in the discussion. The notation and the basic concepts we have described in this section, although referring to abstract entities, enable a uniform discussion of different optimization methods in the next chapters.

Table 3.1: Summary of the notation, concepts, and examples discussed in this chapter, and common to the autotuning methods discussed in further chapters

Symbol	Concept	Example
\mathcal{X}	Search space	All possible compiler flag assignments
$\mathbf{x} = [x_1 \dots x_p]^\top \in \mathcal{X}$	Input variable	A specific flag assignment
$\mathbf{X} = [x_1 \dots x_n]^\top \subseteq \mathcal{X}$	Experimental design	A set of flag assignments
$p = \mathbf{x} $	Search space dimension	The number of flags to assign
$n = \mathbf{X} $	Number of experiments	Size of the set of flags to compile and measure
$f : \mathcal{X} \rightarrow \mathbb{R}$	Function to minimize	Performance metric, such as the execution time of a binary
$y_i = f(\mathbf{x}_i) + \varepsilon_i$	Observable quantity	Execution time with flags \mathbf{x}_i , with error ε_i
$\mathbf{y} = [y_1 \dots y_n]^\top \in \mathbb{R}^n$	Observations	List of execution times for all flags
$\mathbf{x}^* : y^* = f(\mathbf{x}^*) \leq f(\mathbf{x}), \forall \mathbf{x} \in \mathcal{X}$	Global optimum	Flag assignment with smallest execution time
Θ	Parameter space	All possible values of the coefficients of a linear model
$\theta(\mathbf{X}, \mathbf{y}) \in \Theta$	Parameter vector	OLS estimate of the coefficients of a linear model
$m = \theta $	Number of parameters	Number of OLS coefficient estimates
$\hat{f}_{\theta} : \mathcal{X} \rightarrow \mathbb{R}$	Surrogate for f	A linear model fit used for <i>predictions</i>
$\hat{y} = \hat{f}_{\theta}(\mathbf{x})$	Estimate of f at \mathbf{x}	A prediction of execution time for a flag assignment

Before moving on to the descriptions of derivative-based and stochastic function minimization methods and their application to autotuning problems, we discuss in the next section some of the properties of search spaces that are relevant for both autotuning and mathematical optimization.

3.2 Search Spaces

Consider a more abstract optimization problem than the compiler flag selection from the last section, consisting of finding the global optimum of the paraboloid surface defined by the



Figure 3.1: Contour plots and slices through the global optimum, marked with a $\textcolor{red}{\times}$, for search spaces defined by variations of the Booth function. Panels (a) and (d) correspond to Equation 3.3, panels (b) and (e) to Equation 3.4, and panels (c) and (f) to Equation 3.5.

Booth function,

$$y_a = f(\mathbf{x} = [x_1, x_2]^\top) = (x_1 + 2x_2 - 7)^2 + (2x_1 + x_2 - 5)^2, \quad x_1, x_2 \in [-10, 10]. \quad (3.3)$$

In our notation, the search space for this example is $\mathcal{X} = (x_1, x_2) \in \mathbb{R}^2$, $x_1, x_2 \in [-10, 10]$, with global optimum $y^* = f(\mathbf{x}^* = [1, 3]) = 0$.

Besides the search space defined by the observations y_a , we will consider search spaces for two variations

$$y_b = f(\mathbf{x}) + \varepsilon, \quad \text{with } \varepsilon \sim \mathcal{N}(0, \sigma^2), \quad \text{and} \quad (3.4)$$

$$y_c = f(\mathbf{x}) + \varepsilon, \quad \text{with } x_1 \in [3, -6]. \quad (3.5)$$

The underlying objective function in this example has a closed-form expression, but in the context of our applications we consider that we can never observe the true $f(\mathbf{x})$, even in ideal experimental conditions. In that sense, the observations y_b closer to a real application, and incorporate the unknown effects and measurement errors to which the underlying objective function is subject, represented by the normally distributed random variable ε , with mean 0 and variance σ^2 . There are often algorithmic, theoretical, or practical *constraints* on the allowed

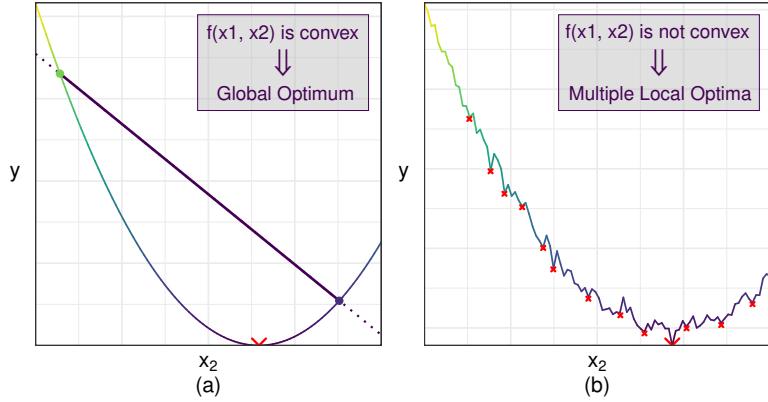


Figure 3.2: Illustrating the relationship between convexity of a function over a compact set and the presence of local minima. Panels (a) and (b) match the functions on the same panels of Figure 3.1

combinations of parameters of a given objective function. The observations y_c represent this scenario by incorporating constraints on the parameter x_1 .

Panels (a), (b), and (c) of Figure 3.1 show contour plots for our search space variations in Equations 3.3, 3.4, and 3.5 respectively. The global optimum in each variation is represented by a red cross, and its location changes between scenarios because of ε . Panels (d), (e), and (f) show slices of panels (a), (b), and (c), respectively, that pass through the global optimum $\mathbf{x}^* = (x_1^*, x_2^*)$ for fixed $x_1 = x_1^*$.

The noise-free example shown in panel (a) has no local optima, and its smooth surface can be quickly navigated by the derivative-based methods discussed in Chapter 4. Such methods aim to follow the direction of *greatest descent* in the neighborhood of a given point. In a contour plot, this direction is always orthogonal to the contour lines. The panel (a) from Figure 3.2 illustrates the *convexity* of our noise-free function. Informally, a line segment connecting any two points in the graph of a convex function will not cross the graph of the function. Convexity of a function over a compact set implies the existence of a single global optimum, whereas lack of convexity implies the existence of *local optima*, as happens in our noisy functions, and is highlighted in panel (b) of Figure 3.2.

Local optima are by definition surrounded by higher values of the objective function, and can thus trap optimization methods that do not plan for such situations, such as a naive implementation of a derivative-based method. Adapting and randomizing the step size of a derivative-based method is one way to deal with functions with many local optima, such as the ones on panels (b) and (c) of Figure 3.1.

Objective functions including constraints can present much harder problems to optimization methods, if certain conditions are met. For example, on panel (c) of Figure 3.1, we have constraints that cut contour lines in such a way that prevents attempts to move inside the feasible space in the direction of greatest descent. A derivative-based method would have

to drastically decrease its step size upon reaching the constraint border, and coast along the border in small steps until it finds a more appealing direction of descent.

This more abstract and simple example aimed to illustrate that it can be non-trivial to find the global optimum of the simplest of search spaces, if we consider the significant challenges introduced by unknown effects and measurement error. The additional challenges introduced by the time cost to obtain measurements, which are discussed in Chapter 6, guided the selection of the optimization methods studied in this thesis. The next chapter discusses derivative-based and stochastic methods for function minimization.

Chapter 4

Methods for Function Minimization

This chapter aims to present the intuition guiding the construction of derivative-based and stochastic methods for function minimization. We discuss the key hypotheses of these groups of methods, and for which autotuning problems they can be most effective.

We will put aside for the moment the idea of using observations to estimate a parameter vector and construct a surrogate function $\hat{f}(\mathbf{x}, \theta(\mathbf{X}, \mathbf{y}))$. This idea will return in later chapters. The methods discussed in this chapter do not construct a surrogate function, and thus attempt to directly optimize the objective function $f(\mathbf{x})$. In this sense, because they need to know how to evaluate it during optimization, these methods make the hypothesis that the objective function is known. The effectiveness of methods based on derivatives requires additional properties of $f(\mathbf{x})$ to be known or estimable, such as its first and second order derivatives, which imposes additional constraints on objective functions.

Evaluating derivatives to determine the next best step or using heuristics to explore a search space cannot be done parsimoniously, because a large number of function evaluations is required to estimate derivatives when closed-forms are unknown, and to explore a search space in the expectation of leveraging unknown structure. We will discuss Design of Experiments in Chapter 6, and present optimization strategies for situations where the cost of evaluating the objective function is prohibitive.

Methods based on derivatives can be powerful, provided their strong hypothesis are respected. We now briefly define and discuss these methods and their application to autotuning.

4.1 Methods Based on Derivatives

The derivatives of a function f at a point \mathbf{x} provide information about the values of f in a neighborhood of \mathbf{x} . It is straightforward to construct optimization methods that use this local information, although iteratively leveraging it requires closed-form expressions for the derivatives of f , or estimates obtained by evaluating f at the neighborhoods of each point.

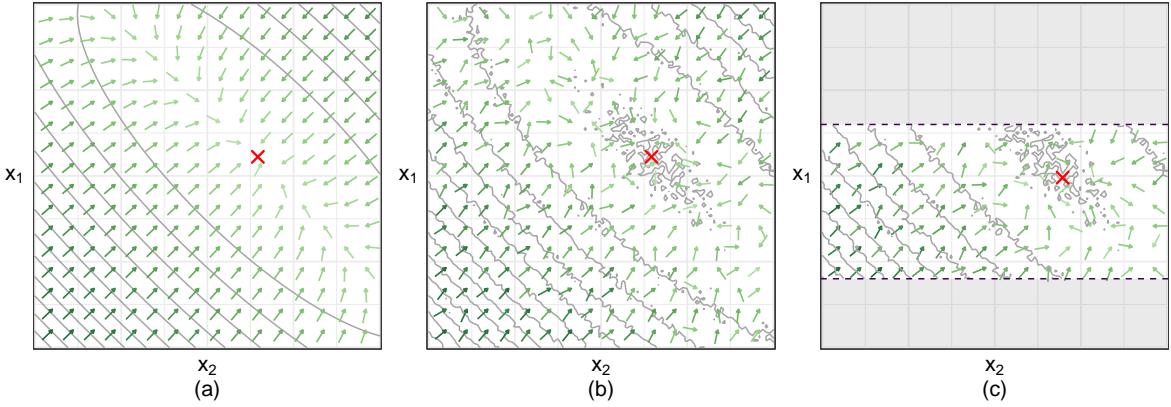


Figure 4.1: Contour plots and direction of greatest descent $-\nabla f(\mathbf{x})$, for search spaces defined by variations of the Booth function. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. The global optimum is marked with a $\textcolor{red}{\times}$. To aid visualization, vector magnitude was encoded by color intensity, so that darker vectors have larger magnitude. The gradient along the function's basin is near zero.

This section discusses gradient descent and Newton’s method, optimization methods using first and second derivatives of f respectively, which for convex functions quickly converge to the global optimum. We will use examples with Booth’s function to discuss how noise and uncertainty impact these methods.

4.1.1 Gradient Descent

The gradient $\nabla f(\mathbf{x})$ of a function $f : \mathcal{X} \rightarrow \mathbb{R}$ at point $\mathbf{x} = [x_1 \dots x_n]^\top$ is defined as

$$\nabla f(\mathbf{x} = [x_1 \dots x_n]^\top) = [f'_{x_1}(\mathbf{x}) \dots f'_{x_n}(\mathbf{x})]^\top, \quad (4.1)$$

where $f'_{x_i}(\mathbf{x})$ is the partial derivative of f at \mathbf{x} , with respect to variable x_i . The vector $\nabla f(\mathbf{x}_i)$ points to the direction of *greatest ascent* in which, from the perspective of $f(\mathbf{x}_i)$, the value of f increases the most.

The gradient descent method is one of the simplest ways to leverage derivative information for optimization. It consists in moving iteratively in the direction of *greatest descent*, opposite the gradient, from a starting point \mathbf{x}_1 . If we follow the opposite of the gradient at \mathbf{x}_1 for additional points $\mathbf{x}_2, \dots, \mathbf{x}_n$, each iteration is written

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \alpha_k \nabla f(\mathbf{x}_{k-1}), \quad k = 2, \dots, n, \quad (4.2)$$

where α_k is the step size at iteration k . The step size for each iteration can be a parameter fixed at the beginning of optimization, but the best α_k can alternatively be determined by searching along the direction of greatest descent.



Figure 4.2: Representation of paths taken by the gradient descent method, with adaptive choice of α_k , on the search spaces defined by variations of the Booth function. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. Contour plots and direction of greatest descent $-\nabla f(\mathbf{x})$ are also shown, and the global optimum is marked with a \times .

Figure 4.1 shows the opposites of the gradients of the three variations of Booth's function, described by Equations 3.3, 3.4, and 3.5, in panels (a), (b), (c), respectively. Figure 4.2 shows, in equally marked panels, the paths taken by a gradient descent algorithm where α_k is chosen at each step, according to the values of f in the neighborhood x_{k-1} .

The gradient descent method iterates along the direction of greatest descent at each point and easily reaches the optimum on the search space of panel (a), unless we make an unlucky choice of α_k . The descent paths on panels (b) and (c) are not so straightforward since the several local minima, represented by the crests and loops on the contour lines, trap the descent path if α_k is not carefully chosen.

The situation is thornier in panel (c), where an unlucky choice of \mathbf{x}_1 or α_k throws the descent path against the top constraint border, forcing the method to zigzag along the border in short steps. This happens in this particular situation in panel (c) because all gradient information guides the descent across the constraint border, but the method cannot cross it. Gradient descent has a harder time on panels (b) and (c) even upon reaching the basin where the optimum lies, because gradient information there is also conflicting due to noise ε . Gradient descent gets stuck in our example paths, but restarting strategies picking new \mathbf{x}_1 or α_k could help escaping the local minima along the basin, as is shown in Figure 4.3.

4.1.2 Newton's Method

With Newton's method we can improve upon the intuition of descending along the opposite of the gradient of f by using the second partial derivatives $f''_{\mathbf{x}_{k-1}}$ to approximate f in the neighborhood of \mathbf{x}_{k-1} using its second Taylor polynomial

$$f(\mathbf{x}) \approx f(\mathbf{x}_{k-1}) + \nabla f(\mathbf{x}_{k-1})^\top (\mathbf{x} - \mathbf{x}_{k-1}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_{k-1})^\top \mathbf{H} f(\mathbf{x}_{k-1}) (\mathbf{x} - \mathbf{x}_{k-1}), \quad (4.3)$$



Figure 4.3: Representation of paths taken by the gradient descent method with 4 restarts, with adaptive choice of α_k , on the search spaces defined by variations of the Booth function. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. Contour plots and direction of greatest descent $-\nabla f(\mathbf{x})$ are also shown, and the global optimum is marked with a \times .

where \mathbf{x} is in the neighborhood of \mathbf{x}_{k-1} , and $\mathbf{H}f(\mathbf{x}_{k-1})$ denotes the Hessian of f , a square matrix of second derivatives of f , with elements

$$(\mathbf{H}f(\mathbf{x}_{k-1}))_{i,j} = f''_{x_i, x_j}(\mathbf{x}_{k-1}). \quad (4.4)$$

We are not going to consider the approximation of f by the second Taylor polynomial to be an estimation process in the statistical sense, because it does not involve dealing with measurement or modeling error.

The second Taylor polynomial uses information about the partial derivatives of f at \mathbf{x}_{k-1} to produce an approximation of f for points \mathbf{x} around \mathbf{x}_{k-1} . If we compute the gradient of this approximation polynomial and set it to zero, we obtain the next point \mathbf{x}_k , as well as the iterative step of Newton's method. Starting at \mathbf{x}_1 , for points $\mathbf{x}_2, \dots, \mathbf{x}_n$, we have

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \mathbf{H}f(\mathbf{x}_{k-1})\nabla f(\mathbf{x}_{k-1}). \quad (4.5)$$

Provided the strong hypotheses of convexity and differentiability are respected, derivative-based methods are extremely effective. In particular, Newton's method converges to the global optimum in a single step if f is quadratic and $\mathbf{H}f$ is positive definite. This happens in the following example with Booth's function

$$y = f(\mathbf{x} = [x_1, x_2]^\top) = (x + 2y - 7)^2 + (2x + y - 5)^2, \quad x_1, x_2 \in [-10, 10], \quad (4.6)$$

adapted from Kochenderfer and Wheeler [61].

If we start with $\mathbf{x}_1 = [9, 8]^\top$ and plug

$$\nabla f(\mathbf{x}_1 = [9, 8]^\top) = [10 \cdot 9 + 8 \cdot 8 - 34, 8 \cdot 9 + 10 \cdot 8 - 38]^\top = [120, 144]^\top, \quad (4.7)$$

and

$$\mathbf{H}f(\mathbf{x}_1 = [9, 8]^\top) = \begin{bmatrix} 10 & 8 \\ 8 & 10 \end{bmatrix} \quad (4.8)$$

into the Newton's method update step in Equation 4.5, we reach $\mathbf{x}^* = [1, 3]^\top$ in the next step

$$\mathbf{x}_2 = \begin{bmatrix} 9 \\ 8 \end{bmatrix} - \begin{bmatrix} 10 & 8 \\ 8 & 10 \end{bmatrix}^{-1} \begin{bmatrix} 120 \\ 144 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \mathbf{x}^*. \quad (4.9)$$

Function minimization methods work very well if the objective function respects their strong hypotheses of convexity and differentiability. Unfortunately, autotuning problems seldom fulfill the conditions necessary for the application of such methods.

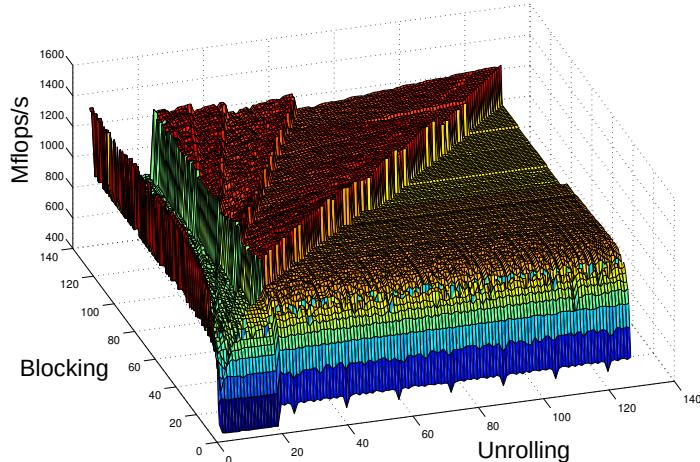


Figure 4.4: An exhaustively measured search space, defined by loop blocking and unrolling parameters, for a sequential GEMM kernel. Reproduced from Seymour *et al.* [5]

Typical autotuning search spaces, such as the one we encountered in Chapter 2, which is reproduced in Figure 4.4, clearly present considerable challenges to derivative-based methods, due to the abundance of local minima, valleys, and ridges. Since we have closed-form expression for the search space in Figure 4.4, we would have to perform a considerable number of evaluations of f in order to estimate its derivative at each step. Although gradient descent and other derivative-based methods can be effective and are historically important, measuring f is usually not cheap for autotuning search spaces. Therefore, minimizing experimentation cost is also a strong concern.

Before discussing methods to construct surrogate models in Chapter 5, and how to use such surrogates to minimize experimental cost on Chapter 6, we will relax the convexity

and differentiability requirements on objective functions and discuss stochastic methods for function minimization, and their applicability to autotuning.

4.2 Stochastic Methods

In this section we discuss some stochastic methods for function minimization that drop the convexity and differentiability requirements of gradient-based methods, becoming applicable to a wider range of autotuning problems at the cost of providing no convergence guarantees. In fact, these methods have no clearly stated hypotheses and are based on search space exploration heuristics. Often the best possible understanding of how these heuristics work comes from the intuition and motivation behind their definition, and from the analysis of empirical tests.

There are multiple ways to categorize heuristics. Our choice was to make a distinction between single-state and population-based methods. In summary, single-state methods have update rules mapping a single point \mathbf{x}_k to a point \mathbf{x}_{k+1} , while rules for population-based methods map a population $\mathbf{P}_k = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ to a population \mathbf{P}_{k+1} , which may retain, combine, and modify elements from \mathbf{P} .

We will first discuss single-state methods, building up to Simulated Annealing from Random Walk, then we discuss Genetic Algorithms and Particle Swarm Optimization, representing widely used population-based methods.

4.2.1 Single-State Methods: Random Walk and Simulated Annealing

Random Sampling is arguably the simplest exploration heuristic, consisting of picking uncorrelated samples from a uniform distribution over the search space \mathcal{X} . There is no guarantee of finding local or global minima, but the chances of improving over a starting point increase if the objective function has many local minima. Despite its simplicity, Random Sampling can be surprisingly effective.

A simple way to derive a single-state heuristic from Random Sampling is to take correlated samples, so that each sample lies in a neighborhood of its predecessor, which is called Random Walk. The neighborhood $N(\mathbf{x})$ of a point is the set of points within distance d from \mathbf{x} . One way to define it is

$$N(\mathbf{x}) = \{\mathbf{x}_i \in \mathcal{X} : \mathbf{x}_i \neq \mathbf{x}, \|\mathbf{x}_i - \mathbf{x}\|^2 \leq d\}. \quad (4.10)$$

As long as it is possible to compute distances between the elements of \mathcal{X} , we will be able to construct the neighborhood of a point and employ the stochastic methods we discuss in this section.

A random walk of length n starting at \mathbf{x}_1 produces a sequence where each point $\mathbf{x}_{k>1}$ is a random variable with uniform distribution over $N(\mathbf{x}_{k-1})$. There is no guarantee that we will ever find a better point with respect to the objective function than the one we started with. A straightforward extension of Random Walk is to pick at each step $k > 1$ the first point \mathbf{x}_k we come across in $N(\mathbf{x}_{k-1})$ for which $f(\mathbf{x}_k) < f(\mathbf{x}_{k-1})$. This greedy strategy would require measuring the value of f for possibly many elements of $N(\mathbf{x}_{k-1})$ but it ensures that we will only move toward a better point.

If we are willing to pay the cost to measure all the points in $N(\mathbf{x}_{k-1})$ we can choose the \mathbf{x}_k^* that brings the best improvement, for which $f(\mathbf{x}_k^*) < f(\mathbf{x})$ for all $\mathbf{x} \in N(\mathbf{x}_{k-1})$. This best improvement strategy always moves to the best point in a neighborhood, but it can still get stuck in local minima if the current point is already the best one in its neighborhood. Adapting the distance that defines a neighborhood can help escape local minima, but the best improvement strategy still requires measuring the entire neighborhood of a point.

Simulated Annealing is a probabilistic improvement heuristic inspired by the process of annealing, where temperature is carefully controlled to first agitate a material's crystalline structure with higher temperature, and then settle it into more desirable configurations. Adapting this idea to optimization, Simulated Annealing makes a compromise between a greedy approach to exploration and a random walk. At each step k , we pick a random uniformly distributed $\mathbf{x}_k \in N(\mathbf{x}_{k-1})$ and move to it if $f(\mathbf{x}_k) < f(\mathbf{x}_{k-1})$. In contrast to a greedy approach, we also move if $f(\mathbf{x}_k) \geq f(\mathbf{x}_{k-1})$ with probability p .

In analogy to the real annealing process, p starts high to enable exploration of the search space and then decreases to force a descent towards an optimum. All through the process, a nonzero value of p permits the heuristic to leave a local minimum. The probability p_k of moving to a worse point at iteration k is

$$p_k = \exp\left(-\frac{f(\mathbf{x}_k) - f(\mathbf{x}_{k-1})}{t_k}\right), \quad (4.11)$$

where t_k is the temperature at iteration k , and can follow different decaying rates. For a starting temperature t_1 the logarithmic annealing schedule is

$$t_k = \frac{t_1}{\log(k + 1)}. \quad (4.12)$$

Figure 4.5 shows the paths taken by Simulated Annealing on search spaces defined by variations of Booth's function. In panel (a) we can see the effects of the annealing schedule, which enables large detours to worse points early on, but forces descent on later iterations. Since this search space has no local minima, the descent eventually reaches the optimum. In panel (b) the higher initial temperature also allows escaping the many local minima found during exploration, but as the path approaches the global optimum with the lower temperature of later iterations it gets trapped by one the local minima around the global optimum.



Figure 4.5: Representation of paths taken by the Simulated Annealing method on the search spaces defined by variations of the Booth function. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. Contour plots and direction of greatest descent are also shown, and the global optimum is marked with a \times .

Likewise, encountering a constraint border early on in panel (c) is not a challenge for Simulated Annealing, since it can bounce back toward worse points, but on later iterations the method is forced to wander around the border like gradient descent, also getting trapped at local minima.

Single-state methods for optimization provide heuristics for exploring a search space, usually based on the neighborhoods of the points that compose a path. The cost of measuring the objective function f for an entire neighborhood is prohibitive in large dimensional search spaces, but methods such as Simulated Annealing provide strategies for escaping local minima without completely evaluating a neighborhood. Despite that, the local nature of single-state methods means that the final results are heavily reliant on the starting point. Restarting and performing parallel searches are among the strategies to reduce dependence on initial choices, as is the idea of using a population of points, which we discuss next.

4.2.2 Population-Based Methods: Genetic Algorithms and Particle Swarm Optimization

Instead of progressing from a single starting point \mathbf{x}_1 toward a final state \mathbf{x}_n , we can use heuristics for moving a starting set of points, or population, $\mathbf{P}_1 = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ toward a final population \mathbf{P}_n . For simplicity, we keep the population size constant. The inspiration for population-based methods for function minimization comes in general from observations of processes such as evolution by natural selection and animal group behavior.

The intuition behind population-based methods is that the points in a population would provide variability that, when combined in specific ways, would eventually lead to better values of the objective function f . Genetic Algorithms, in analogy to the process of evolution by natural selection, select the fittest points in a population for generating offspring.



Figure 4.6: Some ways of producing offspring from two parents with binary chromosomes. Crossover splits parent chromosomes and combine the resulting pieces. In general, pieces from multiple splits can be combined. Mutations are introduced randomly and correspond to flipping bits on binary chromosomes.

Individuals can be selected according to multiple metrics, aiming to produce the best possible combinations in an iteration but also to maintain population variability. Since we do not know how mutations and chromosome combinations might impact the fitness of an individual, keeping worse individuals during optimization could pay off later.

The new population is generated by combining the chromosomes of each parent, using strategies that are also inspired by natural processes such as mutation and crossover. To be able to perform these operations, chromosomes must be encoded in a suitable representation. Individuals in the search space defined by compiler flags, for example, a binary array indicating whether each flag is used could be a suitable representation. Figure 4.6 shows what the mutation and crossover operations could look like on a binary encoding.

Genetic Algorithms have the potential to explore a search space more globally, simultaneously maintaining several populations distributed over a search space, and have also the potential of escaping local minima by mutation and combination of chromosomes of different individuals. Figure 4.7 shows a representation of the paths a population in a Genetic Algorithm could take while searching for the global minimum on three variations of Booth’s function. Each region marked by dashed lines represents the spread of a generation on a given optimization step, and previous individuals are marked by hollow points. The final generation is marked by the filled points. In contrast to Gradient Descent and Simulated Annealing, a Genetic Algorithm do not seek to measure and minimize local properties of f , and consequently its behavior would be less impacted by the noisy scenarios on panels (b) and (c). Since the population can spread across the search space, it could still be possible to end the process with individuals in different local minima, which is represented in panel (b).

Differential Evolution presents an alternative to the strategies of mutations and crossover, represented in Figure 4.6. Consider a population \mathbf{P}_{k-1} , and individuals $\mathbf{x}_a, \mathbf{x}_b, \mathbf{x}_c$ uniformly distributed over the population at step $k - 1$. The update step in Differential Evolution sets the parameters, or components, of an offspring \mathbf{x}_k to a corresponding component of \mathbf{x}_a or $\mathbf{x}_d = w \cdot (\mathbf{x}_b - \mathbf{x}_c)$, according to

$$\mathbf{x}_{k,i} = \begin{cases} \mathbf{x}_{d,i} & \text{if } i = j, \text{ or with probability } p \\ \mathbf{x}_{a,i} & \text{otherwise} \end{cases}, \quad (4.13)$$

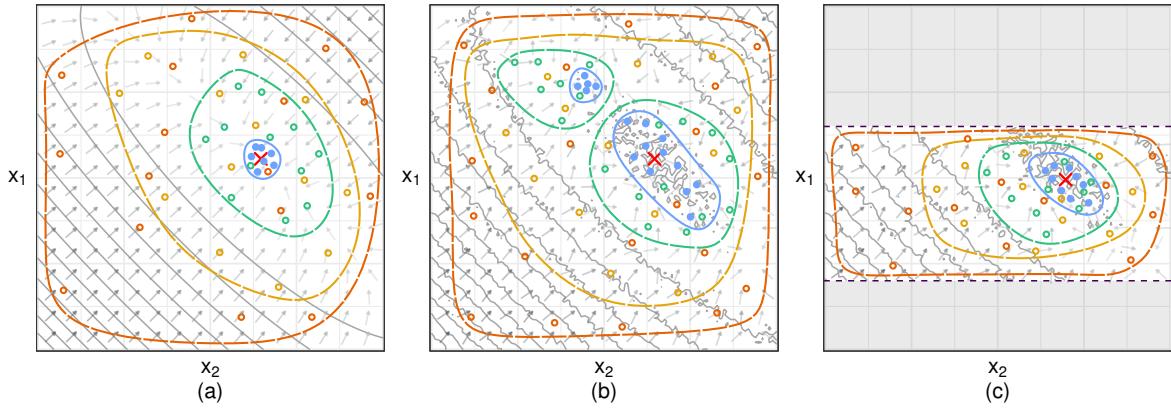


Figure 4.7: Color-coded generation snapshots of a Genetic Algorithm on the search spaces defined by variations of the Booth function. Hollow points and dashed lines mark members of previous populations and the regions they covered, while filled points and complete lines mark the final population. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. Contour plots and direction of greatest descent are also shown, and the global optimum is marked with a \times .

where j is a dimension favored for updating picked at random, and w is a weight given to the difference between \mathbf{x}_b and \mathbf{x}_c . The offspring \mathbf{x}_k replaces \mathbf{x}_a in the population at step k if $f(\mathbf{x}_k) < f(\mathbf{x}_a)$.

A different approach to leveraging a population for optimization is to think of it as an analogy for a swarm, in a method called Particle Swarm Optimization. Each individual \mathbf{x}_k in the swarm keeps track of its velocity \mathbf{v}_k . The velocity is updated at each step, and points to a combination of \mathbf{x}^{best} , the best position found by the swarm so far, and \mathbf{x}_k^{best} , the individual's personal best. At each step, individual \mathbf{x}_k updates its position to

$$\mathbf{x}'_k = \mathbf{x}_k + \mathbf{v}_k, \quad (4.14)$$

and its velocity \mathbf{v}_k according to

$$\mathbf{v}'_k = \alpha_1 \mathbf{v}_k + \alpha_2 (\mathbf{x}^{best} - \mathbf{x}_k) + \alpha_3 (\mathbf{x}_k^{best} - \mathbf{x}_k), \quad (4.15)$$

where α_1 , α_2 , and α_3 are chosen beforehand. The intuition behind this method is that the momentum of each particle toward the best points found so far would accelerate convergence and allow escaping from local minima.

In whichever way we choose to combine individuals in population-based methods, these heuristics require extensive exploration and, consequently, abundant evaluation of the objective function. By making virtually no hypotheses these methods inspired by natural processes become applicable to a much broader range of search spaces than derivative-based methods, while abdicating from convergence guarantees. Population-based methods perform a more global optimization, reliant on the initial distribution of individuals over the search space to provide variability.

4.3 Summary

The methods for function minimization we have discussed in this chapter do not strive to be parsimonious. To be effective they require many estimates of values of f and may require additional information, such as ∇f and $\mathbf{H}f$. The strong hypotheses of derivative-based methods restrict their applicability by invalidating convergence guarantees. Stochastic methods have no such guarantees to begin with, and thus require costly exploration.

Figure 4.8 shows representations of optimization trajectories for Booth’s function made by Gradient Descent with restarts, on panels (a), (b), and (c), Simulated Annealing, on panels (d), (e), and (f), and a Genetic Algorithm, on panels (g), (h), and (i). Restarting optimization from a different point, seen on panels (a), (b), and (c), is a common and widely used technique to reduce reliance on starting conditions.

Despite the high cost of exploration and strong hypotheses, methods for function minimization are used for autotuning and can indeed achieve interesting results in certain problems. We will present our results with these methods on Chapters 10 and 11, where we also review their application to autotuning problems in different domains. The explorations performed by these methods are not structured in a way that favors statistical analysis. We will postpone the discussion of how to obtain well-structured experimental data until Chapter 6, and in the following chapter we will discuss learning methods that enable building surrogate models \hat{f} and identifying relationships between parameters \mathbf{X} in the search space \mathcal{X} and observations of the objective function f .



Figure 4.8: Representation of paths taken by the gradient descent method, with adaptive choice of α_k , on the search spaces defined by variations of the Booth function. Panel groups (a,d,g), (b,e,h), and (c,f,i) correspond to Equations 3.3, 3.4, and 3.5 respectively. Panel groups (a,b,c), (d,e,f), and (g,h,i) correspond to Gradient Descent with restarts, Simulated Annealing, and a Genetic Algorithm, respectively. Contour plots and direction of greatest descent $-\nabla f(\mathbf{x})$ are also shown, and the global optimum is marked with a X .

Chapter 5

Learning: Building Surrogates

This chapter discusses two methods for building surrogate models. Linear Regression is a parametric method capable of modeling a restricted class of surrogates, for which it is relatively simple to interpret significance. Gaussian Process Regression is a flexible nonparametric method capable of modeling a large class of surrogates, for which interpreting significance is possible, but costly.

The surrogate $\hat{f}_\theta : \mathcal{X} \rightarrow \mathbb{R}$ depends on the parameter vector $\theta(\mathbf{X}, \mathbf{y})$ in the parameter space Θ constructed using the pair (\mathbf{X}, \mathbf{y}) , where \mathbf{X} is a set of experiments, and \mathbf{y} the corresponding vector of observations of the objective function f . The process of using the available data to construct a surrogate model is called *learning* and involves fitting the surrogate model, which can also be called *training*. The class of surrogates we can fit depends on each method's definition of Θ .

In this chapter we assume that an initial pair of (\mathbf{X}, \mathbf{y}) is given, and that we can no longer observe $f(\mathbf{x})$ for new \mathbf{x} outside of \mathbf{X} . In this sense, we do not know f but we can still make hypotheses about it in order to construct surrogate models. For the two methods we discuss, we present the hypotheses embedded in each parameter space Θ and how each method fits a surrogate. We also discuss strategies to evaluate the quality of fit of a surrogate, and to interpret parameter significance.

We assume for now that $\mathbf{X} \sim \text{Uniform}(\mathcal{X})$, and we discuss how we can construct better distributions of experiments over a search space on Chapter 6. We will now discuss Linear Regression, a method to build surrogates using linear models on the parameters θ .

5.1 Linear Regression

We will build a surrogate \hat{f}_θ for a function $f : \mathcal{X} \rightarrow \mathbb{R}$ using a fixed experimental design \mathbf{X} and observations \mathbf{y} . We make the hypothesis that f is a linear model on θ , with error ε ,

written

$$f(\mathbf{X}) = \mathbf{X}\theta(\mathbf{X}, \mathbf{y}) + \varepsilon. \quad (5.1)$$

The experimental design \mathbf{X} is an $n \times p$ matrix of column vectors, where each element $\mathbf{x}_1, \dots, \mathbf{x}_n$ is an experiment. We do not know how well the linear model hypothesis represents the true f , but we assume that each experiment in \mathbf{X} was run, producing the response vector $\mathbf{y} = [y_1 \dots y_n]^\top$, where each element is an observation of f subject to measurement error.

Using \mathbf{y} we can construct parameter vectors θ in the parameter space

$$\Theta = \{(\theta_0, \dots, \theta_{m-1}) : \theta_0, \dots, \theta_{m-1} \in \mathbb{R}\}. \quad (5.2)$$

For linear models we can construct the optimal parameter vector $\hat{\theta}$ using the Ordinary Least Squares (OLS) estimator, which we discuss in the next section, and use it to write the linear model surrogate

$$\hat{f}_\theta(\mathbf{X}) = \mathbf{X}\hat{\theta}(\mathbf{X}, \mathbf{y}). \quad (5.3)$$

The vector $\hat{\mathbf{y}}$ produced by evaluating $\hat{f}_\theta(\mathbf{X})$ is called the prediction of the surrogate model for experiments \mathbf{X} , and the distance between \mathbf{y} and $\hat{\mathbf{y}}$ can be used to optimize the parameter vector and estimate the quality of fit of the surrogate.

A key advantage of a surrogate model is that we can use \hat{f}_θ to estimate the value of f for a new design \mathbf{X}' without evaluating $f(\mathbf{X}')$. For well chosen \mathbf{X} and well fitted θ , these estimates can be accurate and useful. We will discuss how to choose experiments on Chapter 6, and how to evaluate the quality of fit for linear models on Section 5.1.2.

5.1.1 Fitting the Model: The Ordinary Least Squares Estimator

Remove reference to H , add dot prod. equal zero

For simplicity, we will for now use θ to refer to $\theta(\mathbf{X}, \mathbf{y})$. The sum of the squared differences between f and \hat{f}_θ for all points in \mathbf{X} is the squared model error

$$\|\varepsilon\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 = (\mathbf{y} - \mathbf{X}\theta)^\top(\mathbf{y} - \mathbf{X}\theta), \quad (5.4)$$

which is a quadratic function of θ . We can therefore differentiate it with respect to θ and set it to zero to obtain the OLS estimator

$$\hat{\theta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (5.5)$$

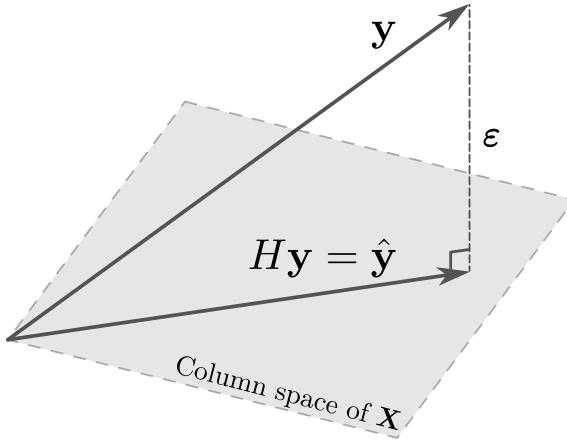


Figure 5.1: In Linear Regression, the prediction vector $\hat{\mathbf{y}}$ is the orthogonal projection of the observations vector \mathbf{y} into the vector space spanned by the columns of \mathbf{X} .

provided $\mathbf{X}^\top \mathbf{X}$ is invertible. The variance of the OLS estimator is written

$$\text{Var}(\hat{\theta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2, \quad (5.6)$$

assuming the observations \mathbf{y} are uncorrelated and homoscedastic, with constant variance σ^2 . Since we assume in this chapter that \mathbf{X} is fixed, the variance of $\hat{\theta}$ shows explicitly how uncertainty in measurements due to error propagates to the linear model surrogate, and we can use this uncertainty to compute confidence intervals for the surrogate's predictions. We will use the fact that the variance of $\hat{\theta}$ depends on $(\mathbf{X}^\top \mathbf{X})^{-1}$ when we discuss Optimal Design on Section 6.3.

An interesting interpretation of the OLS estimator is that it approximates the observation vector \mathbf{y} by its orthogonal projection into the vector space spanned by the columns of \mathbf{X} . Substituting $\hat{\theta}$ on the surrogate model 5.3, we obtain the projection matrix

$$H = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top. \quad (5.7)$$

The difference between \mathbf{y} and its projection $H\mathbf{y} = \hat{\mathbf{y}}$ is the model error $\boldsymbol{\varepsilon}$, orthogonal to the column space of \mathbf{X} , as represented in Figure 5.1. The OLS estimator can also be derived as a special case of the Maximum Likelihood Estimator with $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mu, \sigma^2)$.

We assumed that the design matrix \mathbf{X} was used directly in the surrogate model 5.3, but this limits the models we can represent with $\theta_{0,\dots,m-1}$ to those with $m = p$ linear terms on each factor $\mathbf{x}_{1,\dots,p}$, without an intercept term. More generally, we can have a linear model surrogate with $m \neq p$ terms. In a straightforward expansion of the linear model, we can obtain m model terms from an $n \times p$ design matrix \mathbf{X} by using a basis function set $\mathcal{H} = \{h_{0,\dots,m-1} : \mathcal{X} \rightarrow \mathbb{R}\}$,

generating the $n \times m$ model matrix

$$\mathcal{M} \left(\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}_{n \times p} \right) = \begin{bmatrix} h_0(\mathbf{x}_1) & \dots & h_{m-1}(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ h_0(\mathbf{x}_n) & \dots & h_{m-1}(\mathbf{x}_n) \end{bmatrix}_{n \times m}. \quad (5.8)$$

The OLS estimator does not change if we write the linear model surrogate in Equation 5.3 with $\mathcal{M}(\mathbf{X})$. Unless we need to discuss the underlying model terms, we assume from now on that \mathbf{X} represents a suitable model matrix $\mathcal{M}(\mathbf{X})$. We will now discuss how to evaluate the quality of fit of a linear model surrogate.

5.1.2 Assessing the Quality of Fit of Linear Model Surrogates

We will construct three linear model surrogates for the Booth function and use the model fits to discuss the assessment of model quality. Figure 5.2 shows the three surrogate model fits, constructed using the same design \mathbf{X} with $n = 10$ random uniformly distributed measurements of the Booth function

$$f(\mathbf{x} = [x_1, x_2]^\top) = (x_1 + 2x_2 - 7)^2 + (2x_1 + x_2 - 5)^2, \quad x_1, x_2 \in [-10, 10], \quad (5.9)$$

subject to measurement error ε . Panel (a) shows the surface produced by measuring f in its entire domain. The 10 random measurements composing \mathbf{X} are highlighted. Panels (b), (c), and (d) show the three surrogates, which used the basis function sets

$$\mathcal{H}_{(b)} = \{ h_0(\mathbf{x}) = 1, h_1(\mathbf{x}) = x_1, h_2(\mathbf{x}) = x_2 \}, \quad (5.10)$$

$$\mathcal{H}_{(c)} = h_{(b)} \cup \{ h_3(\mathbf{x}) = x_1^2, h_4(\mathbf{x}) = x_2^2 \}, \text{ and} \quad (5.11)$$

$$\mathcal{H}_{(d)} = h_{(c)} \cup \{ h_5(\mathbf{x}) = x_1 x_2 \}. \quad (5.12)$$

Panels (b), (c), and (d) also show the training, testing, and true Mean Squared Error (MSE) for each surrogate on the design \mathbf{X} , written

$$MSE_{\hat{f}_\theta(\mathbf{x})} = \frac{1}{n} \sum_{i=1}^n \left((f(\mathbf{x}_i) + \varepsilon) - \hat{f}_\theta(\mathbf{x}_i) \right)^2, \quad (5.13)$$

$$MSE_{\hat{f}_\theta} = \frac{1}{|\mathcal{X}|} \sum_{i=1}^{|\mathcal{X}|} \left((f(\mathbf{x}_i) + \varepsilon) - \hat{f}_\theta(\mathbf{x}_i) \right)^2, \text{ and} \quad (5.14)$$

$$MSE_f = \frac{1}{|\mathcal{X}|} \sum_{i=1}^{|\mathcal{X}|} \left(f(\mathbf{x}_i) - \hat{f}_\theta(\mathbf{x}_i) \right)^2. \quad (5.15)$$

The training error $MSE_{\hat{f}_\theta(\mathbf{X})}$ can be computed with the design points and the fitted model. In real applications we can sometimes compute the testing error $MSE_{\hat{f}_\theta}$ for the entire search space, but usually we settle for the error on a testing set distinct from \mathbf{X} . We can almost never compute the true error MSE_f of our surrogate, but in our toy example we can use the Booth function without the error term to get a sense of how well our models generalize. As we could expect, since the true f in our example can be represented by the basis functions on set $\mathcal{H}_{(d)}$, the surrogate from panel (c) has the smallest MSE and the generalizes the best, although it still differs from the true f .

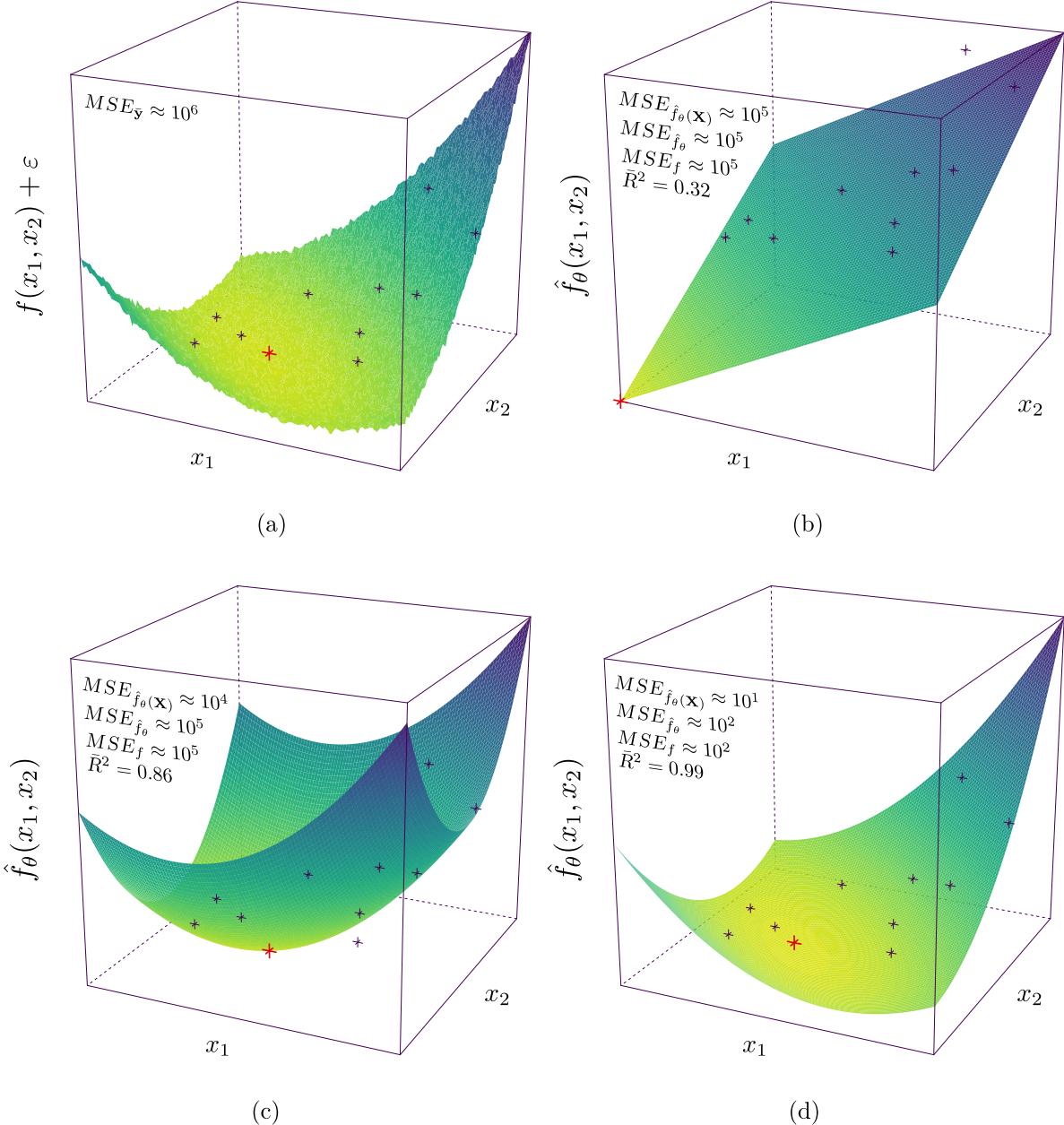


Figure 5.2: Three linear model surrogates for the Booth function, with a \times marking the global optimum and best surrogate predictions. The fixed experimental design \mathbf{X} used to fit all surrogates is marked by \times 's. Panel (a) shows noisy measurements of Booth's function, panels (b), (c), and (d) show surrogate predictions for models fit with basis functions sets from Equations 5.10, 5.11, and 5.12, respectively

We can use the mean \bar{y} of the observations of experiments in \mathbf{X} as a surrogate, for which the MSE is written

$$MSE_{\bar{y}} = \frac{1}{n} \sum_{i=1}^n ((f(\mathbf{x}_i) + \varepsilon) - \bar{y})^2, \quad (5.16)$$

shown in panel (a) from Figure 5.2. We can use this surrogate to compute the adjusted coefficient of determination

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-p-1} \cdot \frac{MSE_{\hat{f}_\theta(\mathbf{x})}}{MSE_{\bar{y}}} \right), \quad (5.17)$$

which compares the squared error of a given linear model with the error of the model that predicts the mean of all observations, adjusted by the flexibility introduced by new parameters.

In general, we can measure the complexity of a surrogate model by the number of parameters that we need to estimate when fitting it to experimental data. The bias of a surrogate, for a fixed point \mathbf{x}_0 , is the expected value of the distance between the surrogate's prediction and $f(\mathbf{x}_0)$, over surrogate fits using a large number of different experimental designs. The squared bias is written

$$\text{Bias}^2 (\hat{f}_\theta(\mathbf{x}_0)) = \left(E [\hat{f}_\theta(\mathbf{x}_0)] - f(\mathbf{x}_0) \right)^2, \quad (5.18)$$

because $E[f(\mathbf{x}_0)] = f(\mathbf{x}_0)$. As we increase the flexibility of a model by adding parameters we allow the surrogate model to better approximate design points, reducing its bias. Concomitantly, we increase the number of different parameter vectors that can describe the experimental data, increasing the surrogate's variance

$$\text{Var} (\hat{f}_\theta(\mathbf{x}_0)) = E \left[\hat{f}_\theta(\mathbf{x}_0) - E [\hat{f}_\theta(\mathbf{x}_0)] \right]^2. \quad (5.19)$$

The total model error for the prediction of \mathbf{x}_0 still has to factor in the irreducible error ε associated with the measurements of f , and is written

$$\text{Error} (\mathbf{x}_0) = \text{Var}(\varepsilon) + \text{Bias}^2 (\hat{f}_\theta(\mathbf{x}_0)) + \text{Var} (\hat{f}_\theta(\mathbf{x}_0)). \quad (5.20)$$

Increasing surrogate complexity reduces bias but increases variance, and this trade-off is central to selecting and assessing the quality of surrogate models, however we define the parameter space Θ .

5.1.3 Inference: Interpreting Significance with ANOVA

The variance of the OLS estimator enables us to compute confidence intervals and p -values for the effects of each model term, or factor, for a single surrogate model fit. In a frequentist interpretation a 95% confidence interval for the estimate of a mean of a factor's effect is

interpreted as the interval that would contain 95% of our estimates, were we to repeat the estimation multiple times. The p -value of a factor effect's estimate is interpreted in frequentist inference as the probability of observing an effect at least as large as what was observed, if the factor's true effect is zero. When they can be computed, confidence intervals are in general more useful than p -values for judging the accuracy of a factor effect's estimate, because they are explicitly defined in the context of the magnitude of that estimate.

Analysis of Variance (ANOVA) is a more refined statistical tool for significance testing, able to estimate relative factor significance. The steps of an ANOVA test are grouping the observations \mathbf{y} by factors and factor levels, computing separate group means, and testing the significance of the differences between group means with an F -test.

The ANOVA test can be understood as a special case of the linear model we have discussed in this section, in which case its formal hypotheses are the same as the linear model's, that is, that the observations \mathbf{y} are uncorrelated, the residuals are normally distributed, and the variances of each group are homoscedastic.

Running an ANOVA test as a special case of the linear model consists of running F -tests for multiple models $\mathbf{y} = \mathbf{X}\theta$, with specially constructed model matrices of indicator variables for group and group interaction membership. A detailed description of ANOVA in relation to linear models can be found Chapter 6 of Dobson *et al.* [62], among other reference texts [63, 64].

5.1.4 Linear Models: Interpretable but Biased Surrogates

Linear regression can be successful in learning and interpreting relationships between factors and an objective function f . If the underlying functions have complex structure that cannot be sufficiently well represented by a finite number of parameters, linear models might not be useful beyond identifying the strongest factor effects. The following section will discuss Gaussian Process Regression, a nonparametric approach that fits a model by conditioning a probability distribution over functions.

5.2 Gaussian Process Regression

Similarly to the surrogate model we built using linear models, in this section we will fit a surrogate $\hat{f}_\theta : X \rightarrow \mathbb{R}$ to a fixed experimental design \mathbf{X} and observations \mathbf{y} . We make the hypothesis that the observations y_1, \dots, y_n are normally distributed. This is a reasonable hypothesis to make, especially after our discussion on linear models, where our hypothesis over f was

$$f(\mathbf{X}) = \mathbf{X}\theta + \varepsilon, \quad (5.21)$$

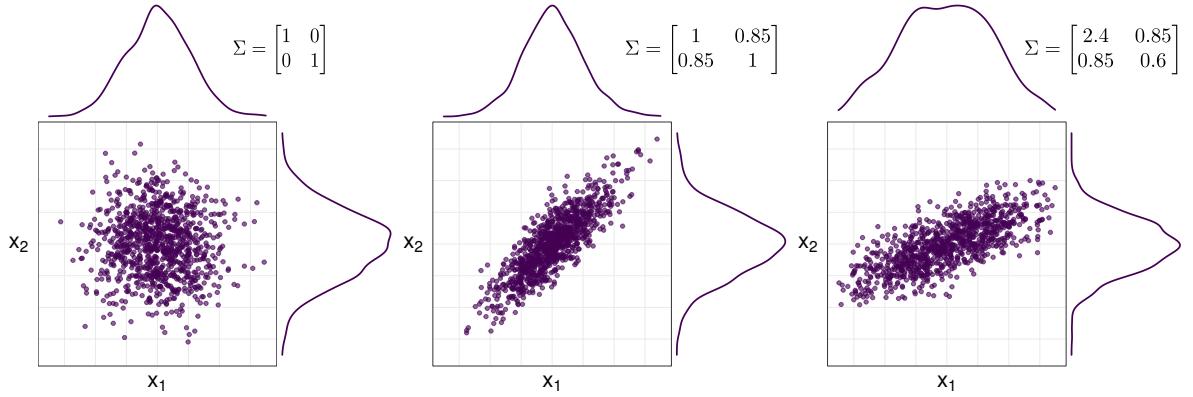


Figure 5.3: Effects of three covariance matrices on a multivariate normal distribution with mean vector $\mu = [0 \ 0]^\top$ and covariance matrix Σ as shown on the upper right corner of each plot

with $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, which makes $f(\mathbf{X})$ a multivariate random variable, written

$$f(\mathbf{X}) \sim \mathcal{N}(\mathbf{X}\theta, \sigma^2 \mathbf{I}), \quad (5.22)$$

where \mathbf{I} is the $n \times p$ identity matrix.

Another way to state this key idea is to say that we make the hypothesis that f is a Gaussian Process, that is, that it belongs to the class of models containing all the functions whose values on any set of experiments \mathbf{X} can be represented by a single sample of a multivariate normal distribution, with dimension $n = |\mathbf{X}|$. We can write this hypothesis as a prior probability distribution

$$f(\mathbf{X}) = [y_1 \ \dots \ y_n]^\top \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0), \quad (5.23)$$

where the mean vector $\boldsymbol{\mu}_0$ is usually a vector of n zeros and the covariance matrix $\boldsymbol{\Sigma}_0$ is computed using a kernel and depends on \mathbf{X} .

This prior is more general than the linear model hypothesis in Equation 5.22, but it generalizes even further. In Chapter 6 of *Gaussian Processes for Machine Learning* [10], Rasmussen and Williams describe relationships and equivalences of Gaussian Processes to other methods such as Splines, Support Vector Machines, and OLS estimation. They also describe how Neural Networks can be represented by specific Gaussian Processes in Section 4.2.3.

In contrast to the linear model, the definition and exploration of a parameter space Θ for a Gaussian Process surrogate is done indirectly, by controlling kernel parameters and noise amplitudes. A first approach to define the parameter space could be explicitly searching for the covariance matrix and mean vector that best fit the data, where the size of the parameter space would increase as more training data becomes available. For a given design \mathbf{X} of size n ,

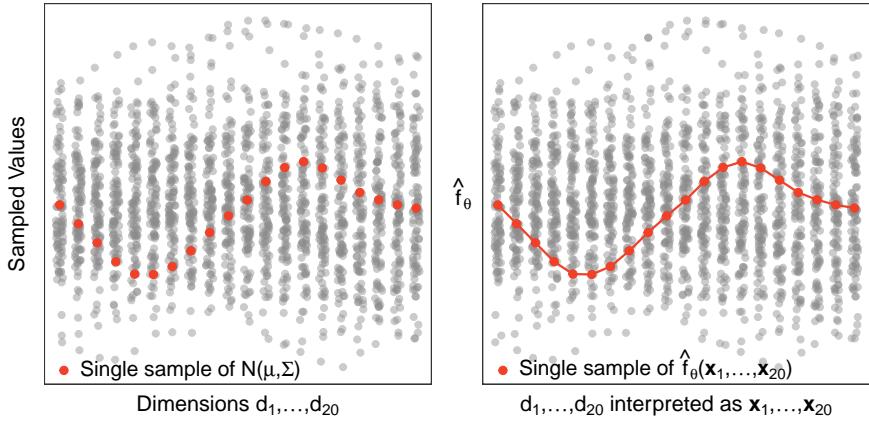


Figure 5.4: Reinterpreting the unrolled dimensions of 100 samples of a 20 dimension multivariate normal, on the left panel, to obtain 100 samples of functions evaluated on 20 different input points, on the right panel

this parameter space would be

$$\Theta_0 = \{ (\mu_{\mathbf{X}, \mathbf{y}}, \Sigma_{\mathbf{X}}), \mu_{\mathbf{X}, \mathbf{y}} \in \mathbb{R}^n, \Sigma_{\mathbf{X}} \in \mathbb{R}^{n \times n} \}. \quad (5.24)$$

Note that the mean vector depends on the observations \mathbf{y} , while the covariance matrix depends only on \mathbf{X} . In a multivariate Gaussian, the diagonal of the covariance matrix contains the variances associated with its dimensions, and the off-diagonal elements contain the covariances between pairs of dimensions. Figure 5.3 shows the impact of three covariance matrices, which must always be symmetric and positive semi-definite.

Instead of looking for the best parameter vector in Θ_0 , we will condition the prior Gaussian distribution to the observed data, obtaining a posterior distribution. The conditioned distribution representing our surrogate model is also a multivariate normal written

$$\hat{f}_{\theta}(\mathbf{X}) \sim f(\mathbf{X}) | \mathbf{X}, \mathbf{y}. \quad (5.25)$$

We will see how to compute this posterior in the next section. Note that it would require an infinite parameter vector to fit a Gaussian Process for all points in a search space \mathcal{X} consisting of a single real number. In this sense, Gaussian Process Regression is a nonparametric method.

To build an intuitive understanding of Gaussian Processes it can be helpful to think of them as a reinterpretation of the dimensions of a sample of a multivariate normal distribution. If we take 100 samples of a 20-dimension Gaussian and unroll each dimension into a single axis, we end up with the left panel of Figure 5.4. Each column of points contains the 100 values on each dimension of our samples, and each dimension is correlated according to the distribution's covariance matrix. The values for one arbitrary sample are marked in larger red dots.

Table 5.1: Expressions for the covariance functions, or kernels, shown in Figure 5.5. The variables v and l are kernel parameters that can themselves be estimated. The Matérn kernel depends on the gamma function Γ and on the Bessel function of the second kind K_v . We refer the reader to Chapter 4 of Rasmussen and Williams [10] for detailed definitions and discussions

Kernel	Expressions
Exponential	$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\ \mathbf{x} - \mathbf{x}'\ }{l}\right)$
Squared Exponential	$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\ \mathbf{x} - \mathbf{x}'\ ^2}{2l^2}\right)$
Matérn	$k(\mathbf{x}, \mathbf{x}') = \frac{1}{2^{v-1}\Gamma(v)} \left(\frac{\sqrt{2v}}{l}\ \mathbf{x} - \mathbf{x}'\ \right)^v K_v\left(\frac{\sqrt{2v}}{l}\ \mathbf{x} - \mathbf{x}'\ \right)$

The right panel of Figure 5.4 shows the same data in the left panel but we now interpret the values of each dimension d_i , for all 100 samples, as a distribution of values for the surrogate \hat{f}_θ evaluated at point \mathbf{x}_i . Each $\hat{f}_\theta(\mathbf{x}_i)$ is correlated to other values of the surrogate function, with covariance given by the distribution's covariance matrix. In this example we used the Matérn kernel, discussed in the next section, to compute the covariance matrix. This prior over functions directly estimates a mean and its associated variance, for each value of the surrogate. We will now discuss how to fit a Gaussian Process to observed data, generating predictions of means and variances conditioned to observations.

5.2.1 Fitting the Model: Posterior Distributions over Functions

Before fitting our surrogate, we must compute the covariance matrix of the prior Gaussian distribution from Equation 5.23 using a covariance function, or kernel, $K : \mathcal{X}^2 \rightarrow \mathbb{R}$. For any pair $(\mathbf{x}, \mathbf{x}') \in \mathcal{X}^2$, the kernel determines how strong the covariance between $\hat{f}_\theta(\mathbf{x})$ and $\hat{f}_\theta(\mathbf{x}')$ should be, based on the distance $\|\mathbf{x} - \mathbf{x}'\|$. Figure 5.5 shows four exponential kernels, also called radial basis functions, whose formulas are shown in Table 5.1.

The actual parameter space Θ over which optimization is performed when fitting a Gaussian Process surrogate is composed by the parameters of the chosen covariance kernel, called hyperparameters. Typical hyperparameters are the Section 5.1 of Rasmussen and Williams [10], . Explicitly listing the hyperparameters for a Gaussian Process fit using the exponential kernel from Table 5.1 would result in

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|}{l}\right), \quad (5.26)$$

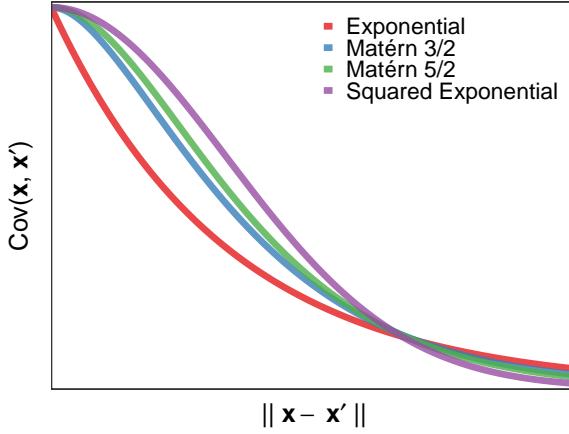


Figure 5.5: Covariance of points $(\mathbf{x}, \mathbf{x}')$ according to four covariance functions based on the distance $\|\mathbf{x} - \mathbf{x}'\|$. Expressions for each kernel are shown in Table 5.1, and Matérn kernels use parameters $v_1, v_2 = \{\frac{3}{2}, \frac{5}{2}\}$

and the parameter space Θ would be defined as

$$\Theta = \{ \sigma_f, l \in \mathbb{R} \}. \quad (5.27)$$

We can then use the conditioned posterior distribution, discussed below, to determine the best specific hyperparameter values by minimizing the cross-validated mean squared error, for example, which we discuss in Section 5.2.2, or by maximizing the posterior likelihood.

The posterior distribution \hat{f}_θ is computed by conditioning the prior distribution from Equation 5.23 to observed data $(\mathbf{X}_k, \mathbf{y}_k)$, obtaining the distribution in Equation 5.25, which is also a Gaussian distribution, and can be written

$$\hat{f}_\theta(\mathbf{X}_k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (5.28)$$

where the conditioned vector mean, of size $|\mathbf{X}_k|$, is

$$\boldsymbol{\mu}_k = \mathbf{K}(\mathbf{X}_k, \mathbf{X}_{k-1}) \mathbf{K}(\mathbf{X}_{k-1}, \mathbf{X}_{k-1})^{-1} \mathbf{y}_{k-1}^\top, \quad (5.29)$$

and the conditioned $|\mathbf{X}_k| \times |\mathbf{X}_k|$ covariance matrix is written

$$\boldsymbol{\Sigma}_k = \mathbf{K}(\mathbf{X}_k, \mathbf{X}_k) - \mathbf{K}(\mathbf{X}_k, \mathbf{X}_{k-1}) \mathbf{K}(\mathbf{X}_{k-1}, \mathbf{X}_{k-1})^{-1} \mathbf{K}(\mathbf{X}_{k-1}, \mathbf{X}_k). \quad (5.30)$$

The function \mathbf{K} in Equations 5.29 and 5.30 produces the covariance matrices corresponding to applying one of the covariance kernels $k(\mathbf{x}, \mathbf{x}')$ from Table 5.1 to all pairs of points on the input designs for \mathbf{K} . The complete text of Rasmussen and Williams [10], the Chapter 6 of Bishop [65], and the Chapter 15 of Kochenderfer and Wheeler [61] are among the texts that present detailed discussions and derivations of important properties of Gaussian Process Regression.

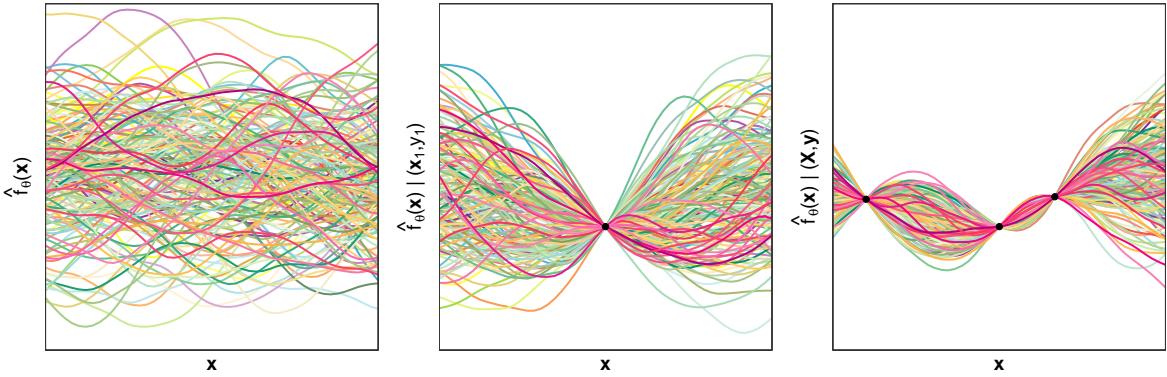


Figure 5.6: Fitting a Gaussian Process to three noise-free observations. The left panel shows 300 samples from a Gaussian prior, using the Matérn kernel to compute the covariance matrix. The center and right panels show 300 samples from the posterior distributions conditioned by one, then two more, successive noise-free observations

The left panel of Figure 5.6 shows 300 sampled functions from the prior in Equation 5.23, for a 50-dimension Gaussian distribution, with covariance matrix Σ_0 given by the Matérn kernel from Table 5.1, with $v = \frac{5}{2}$. The starting mean vector μ_0 is zero. We see that, for all radial basis functions shown in Figure 5.5, the covariance between inputs decreases as the distance between inputs increases, approaching zero. Using different kernels we can control properties of sampled functions, such as smoothness and periodicity. The center panel of Figure 5.6 shows samples from the conditioned posterior after a single observation (\mathbf{x}, y) , computed using Equations 5.28, 5.29, and 5.30. The right panel shows samples from the posterior after observing two additional observations. Note that all 300 functions sampled from the conditioned distributions pass exactly through the observations on the center and rightmost panels.

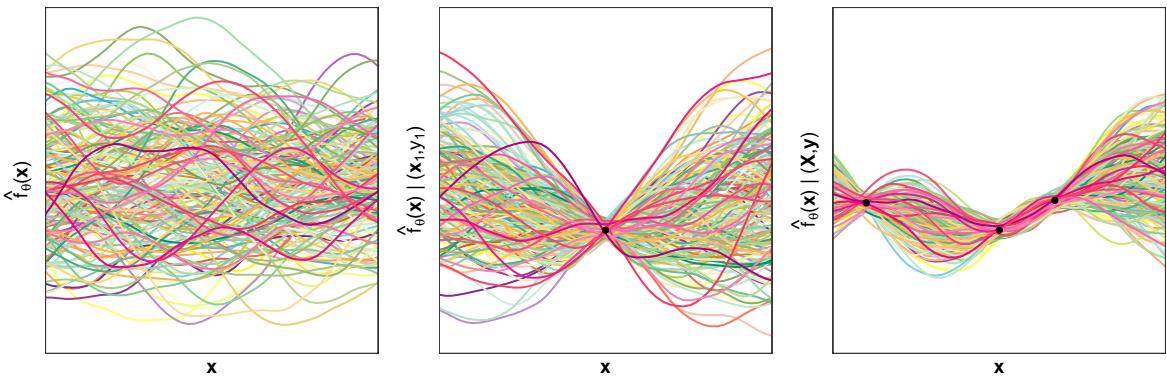


Figure 5.7: Fitting a Gaussian Process to three noisy observations, in the same conditions and with the same panel structure in Figure 5.6

Although in specific autotuning applications involving deterministic or extremely fast processes we can observe consistent measurements and produce useful models with fits that assume noiseless measurements, such as in Figure 5.6, we can also produce a Gaussian

Process fit that incorporates the uncertainty from noisy data. Figure 5.7 shows 300 prior samples, in the left panel, and results of conditioning the prior distribution to one and then three observations, under the assumption that the underlying objective function is subject to measurement error $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, as described in Equation 5.21. In this scenario, we write the mean vector and covariance matrix for Equation 5.28 as

$$\mu_k = \mathbf{K}(\mathbf{X}_k, \mathbf{X}_{k-1}) (\mathbf{K}(\mathbf{X}_{k-1}, \mathbf{X}_{k-1})^{-1} + \sigma^2 \mathbf{I}) \mathbf{y}_{k-1}^\top, \quad (5.31)$$

and

$$\Sigma_k = \mathbf{K}(\mathbf{X}_k, \mathbf{X}_k) - \mathbf{K}(\mathbf{X}_k, \mathbf{X}_{k-1}) (\mathbf{K}(\mathbf{X}_{k-1}, \mathbf{X}_{k-1})^{-1} + \sigma^2 \mathbf{I}) \mathbf{K}(\mathbf{X}_{k-1}, \mathbf{X}_k). \quad (5.32)$$

Before discussing quality of fit assessment metrics for Gaussian Process surrogates, we will discuss how to incorporate hypotheses over the search space to fitted surrogates using basis functions, which we described when we discussed linear regression. Trend functions can be added to the surrogate's predicted mean to leverage underlying trends in data, that should be followed on prediction regions far from measurements.

Figure 5.8 shows Gaussian Process surrogates fitted with and without model trends. Black circles represent six measurements of a single-input function, and each line represents a model trend. We are considering noise-free measurements in these fits, so all surrogates agree on the predictions at the measurements. We can see that predictions for input between or sufficiently far from measurements present stronger influence from the underlying trend. A trend can help leverage previous knowledge of the relationships between factors and the values of the objective function, especially when the cost of measuring a single point is expensive, or when new experiments cannot be performed. However, if we can choose which measurements to perform we can use space-filling designs, which we discuss on Chapter 6, to decrease or sometimes remove the need for a model trend.

The fits and predictions shown in Figure 5.8 were computed using the *DiceKriging R* package [66]. We refer the reader to the detailed descriptions of trend functions and their application to Gaussian Process Regression presented in Section 2.1 of the package's accompanying paper [67].

As we did for the linear model surrogates discussed earlier in this chapter, in the next section we will assess the quality of fit of Gaussian Process surrogates, and quantify their prediction error on a testing set outside of the measurements used for fitting.

5.2.2 Assessing the Quality of Fit of Gaussian Process Surrogates

We will construct three Gaussian Process surrogates for the Booth function and use the model fits to discuss the assessment of model quality. Figure 5.9 shows the three surrogate fits, constructed using the same design \mathbf{X} with $n = 10$ random uniformly distributed measurements

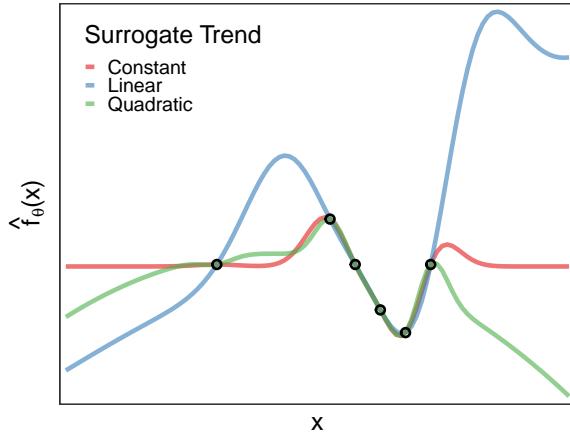


Figure 5.8: Gaussian Process surrogates using three different model trends, fit to six noise-free observations of a single-input objective function, marked with black circles

of the Booth function

$$f(\mathbf{x} = [x_1, x_2]^\top) = (x_1 + 2x_2 - 7)^2 + (2x_1 + x_2 - 5)^2, \quad x_1, x_2 \in [-10, 10], \quad (5.33)$$

subject to measurement error ε . Panel (a) shows the surface produced by measuring f in its entire domain and the MSE of the constant mean predictor $MSE_{\bar{y}}$. The 10 random measurements in \mathbf{X} are highlighted. Panels (b), (c), and (d) show the three surrogates, which used model trends using the basis function sets described by Equations 5.10, 5.11, and 5.12. Panels (b), (c), and (d) show the testing $MSE_{\hat{f}_\theta}$ and true MSE_f prediction errors for each surrogate, described by Equations 5.14 and 5.15.

We have assumed noisy measurements when fitting the Gaussian Process surrogates in this example. In this setting, we could compute the training prediction error $MSE_{\hat{f}_\theta(\mathbf{X})}$ as described by Equation 5.13, but this would not work for a surrogate assuming noise-free measurements. Since all surrogates would interpolate the observations, the prediction error for points in the training set would always be zero. Still, analogously, in our example we also controlled the variances representing the measurement errors, as described by Equations 5.31 and 5.32.

The strategy we used for computing test prediction error for linear models would not be accurate for this example, and we used a Leave-One-Out (LOO) cross-validation strategy to compute the testing error $MSE_{\hat{f}_\theta}^{LOO}(\mathbf{X})$. The strategy consists of computing the mean of all $MSE_{\hat{f}_\theta(\mathbf{X})}$ for different surrogates, with the same model trend, fit to the testing sets generated by removing one distinct training point at a time, for all training points. This strategy also works to compute the testing prediction error for noise-free Gaussian Process fits.

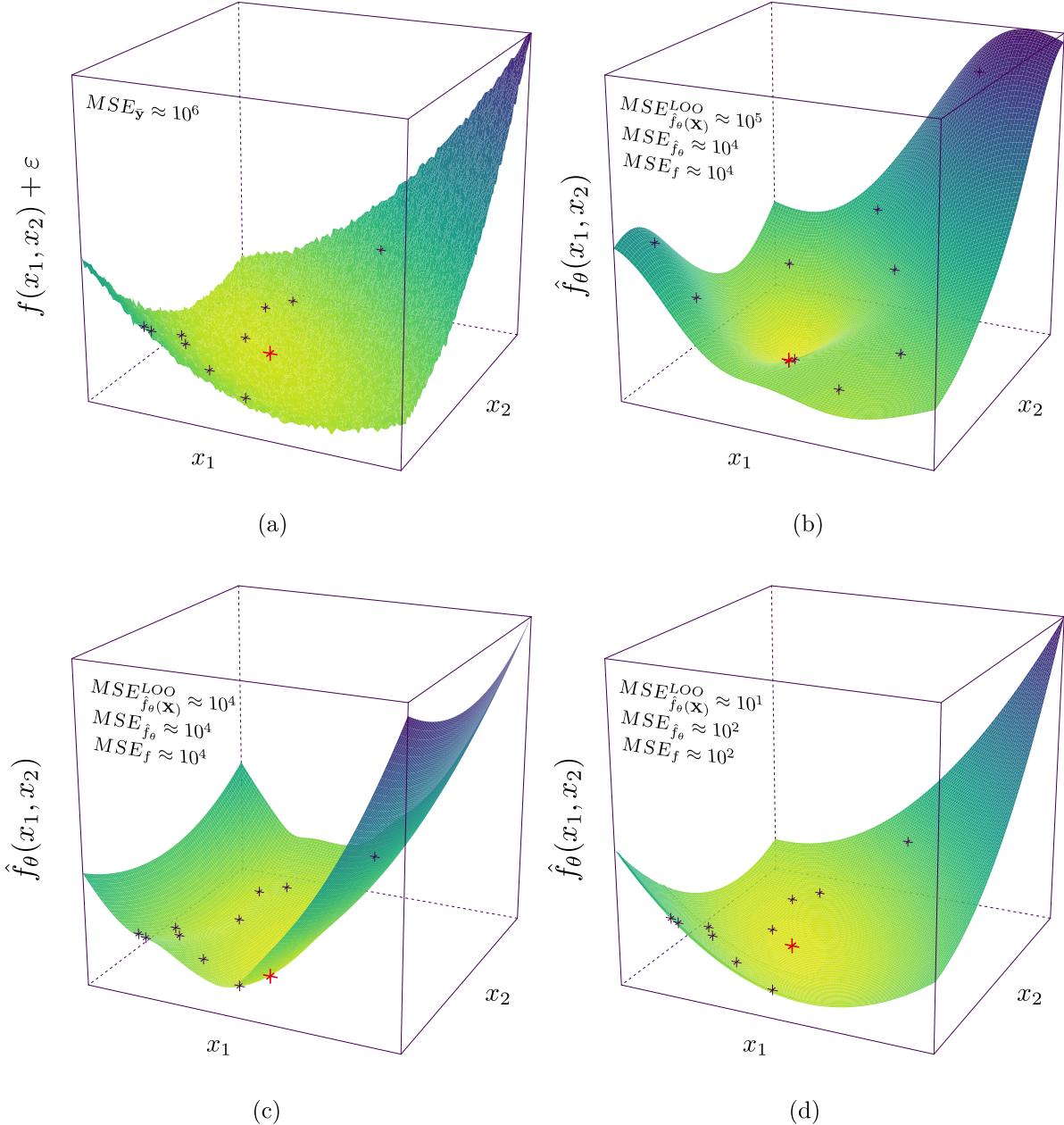


Figure 5.9: Three Gaussian Process Regression surrogates with noisy fits for the Booth function, with a \times marking the global optimum and best surrogate predictions. The fixed experimental design \mathbf{X} used to fit all surrogates is marked by \times s. Panel (a) shows noisy measurements of Booth's function, panels (b), (c), and (d) show surrogate predictions for models fit with linear, quadratic, and quadratic plus interactions trends, respectively

5.2.3 Inference: Interpreting Significance with Sobol Indices

Sobol indices [68] are a variance-based global sensitivity analysis method, which computes the relative importance of the input factors $\mathbf{x}_1, \dots, \mathbf{x}_p$ of an objective function f by decomposing the variance observed in measurements of f . The first-order Sobol index S_i for each factor \mathbf{x}_i represents the normalized variance of the expected value of f given a fixed factor value, written

$$S_i = \frac{V_{\mathbf{x}_i} (E_{\mathbf{X} \setminus \mathbf{x}_i} (f(\mathbf{X}) | \mathbf{x}_i))}{V(f(\mathbf{X}))}. \quad (5.34)$$

The variance of a function can be further decomposed to produce $2^p - 1$ indexes that measure the sensitivity of factor interactions of up to the p^{th} level. Additionally, we can aggregate the effects of the interactions of a factor and compute total-effect indices, in order to decrease the total number of indices to compute.

Sobol sensitivity indices can be estimated using $N(p + 2)$ evaluations of the objective function, where N controls the accuracy of the estimate and is typically in the order of thousands. The Monte Carlo method described in Section 4.1.6 of Saltelli *et al.* [69] determines how to construct the experiment matrices containing all samples, and computes the sensitivity indices using the estimators described in the book and later reviewed and summarized by the authors [70].

5.2.4 Gaussian Processes: Flexible but Hard to Interpret Surrogates

Gaussian Process Regression is a nonparametric learning method with low bias and high variance, that can produce surrogate models capable of describing observations with essentially arbitrary precision, without entirely committing to structure outside of observations. Although we can attempt to interpret the results produced by a Gaussian Process fit with Sobol indices, it is not always possible or viable to collect sufficient data to reach accurate results. Recent work interprets Gaussian Process fits by identifying the contribution of different covariance kernels to the fitted model. Such approaches help constructing useful models in high-dimensional spaces, and are valuable to experimental research, such as the studies we present in this thesis.

Add reference to PhD thesis from Rasmussen's student

5.3 Summary

This chapter presented Linear and Gaussian Process Regression, two large classes of surrogate model construction methods. From the point of view of an experimenter, a practical difference between these two classes lies in their balance of interpretability and generality.

Linear models are easier to interpret, they have high bias and low variance, and can be used to test hypotheses about search space structure, or leverage known relationships. Gaussian Process Regression is hard to interpret, has low bias and high variance, and can be used to model search spaces for which it is hard to elaborate clear hypotheses, or where simple structure does not exist to be exploited.

In this chapter we assumed that the design matrix \mathbf{X} is a given, or is sampled from a fixed set of observations. In the next chapter we will discuss methods of Design of Experiments that enable the construction of carefully chosen design matrices for different purposes, such as testing hypotheses, improving model fits, or minimizing the objective function.

Chapter 6

Design of Experiments

This chapter presents some elements of Design of Experiments (DOE). Complementing the learning methods we discussed in the previous chapter, DOE methods enable choosing the experiments on a design matrix \mathbf{X} according to different needs and criteria. A well-constructed design enables testing hypotheses about the objective function, which can lead either to a useful model or to a set of questions to ask and hypotheses to test.

This chapter is organized as follows. Section 6.1 introduces key concepts and presents 2-level Full and Fractional designs. Section 6.2 presents Screening designs, an efficient way to detect linear effects. Section 6.3 discusses Optimal Designs, which enable constructing efficient designs for more complex hypotheses. Section 6.4 presents Space-filling designs, useful for search space exploration and more complex learning methods. Section 6.5 compares different sampling strategies, and Section 6.6 summarizes the discussion and concludes the chapter.

6.1 Estimating Linear Effects

We have been using Design of Experiments terminology in this thesis, where the vector \mathbf{y} with observations of an objective function f in an input set \mathbf{X} is called a *response*. The input set is an $n \times p$ matrix called a *design*, where each row is an *experiment* with p *factors*. Choosing an adequate design depends on characteristics of the factors, of the objective function, and of the statistical analysis or learning methods we wish to use.

In this chapter we discuss methods for constructing experimental designs with the assumption that the analysis of the results will be made using linear models. This section will present the difficulties of estimating coefficients for linear model terms in settings with high levels of error or noise, heteroscedasticity, and interactions between factors.

6.1.1 2-Level Factorial Designs

Figure 6.1 shows measurements of an objective function whose input is a single real number factor $x_1 \in [0, 1]$. Each of the six panels represents an experiment where we fit a linear model to a limited experimental budget of two observations of f , represented by blue points. Hollow points in each panel represent a much larger sample of values of f , to aid the visualization of each experimental condition.

Our objective in this example is to distribute our 2-point budget so as to produce a linear model fit for f which will guide future experiments in the correct direction for minimizing it. The direction in which we would explore if we followed each model fit is represented by dashed vertical red lines. With such a budget restriction we will not have sufficient data to estimate the confidence intervals for the predictions of our linear model. We did this in this example so we can discuss the impact of the choice of design points in a single prediction, and we later look at an experiment with more measurements.

The first row of Figure 6.1 shows measurements subject to small measurement error. In these conditions we can produce a fairly good fit regardless of how close we pick our two experiments. This experiment with small noise reveals that the “true” shape of f is linear, minimized by $x_1 = 0$.

The experiments in the second row have larger measurement error. In real situations this could happen if we were measuring f in a shared machine, for example, while in the first row we might have had exclusive access. In this setting, picking experiments too close to each other can invert the direction in which the model tells us to explore, and we will be better off spending our two measurements at the interval extremes.

The last row shows an experimental setting where measurement error depends on the value of x_1 , violating a core hypothesis of linear models. Although we could still fit a linear model indicating the right exploration direction in this scenario, chances are high that we won’t, even if we pick extreme points. Heteroscedastic error can be produced in real conditions by resource management overhead, for example, and the simplest way to try and deal with it in practice is to fit a linear model to the logarithm of the response.

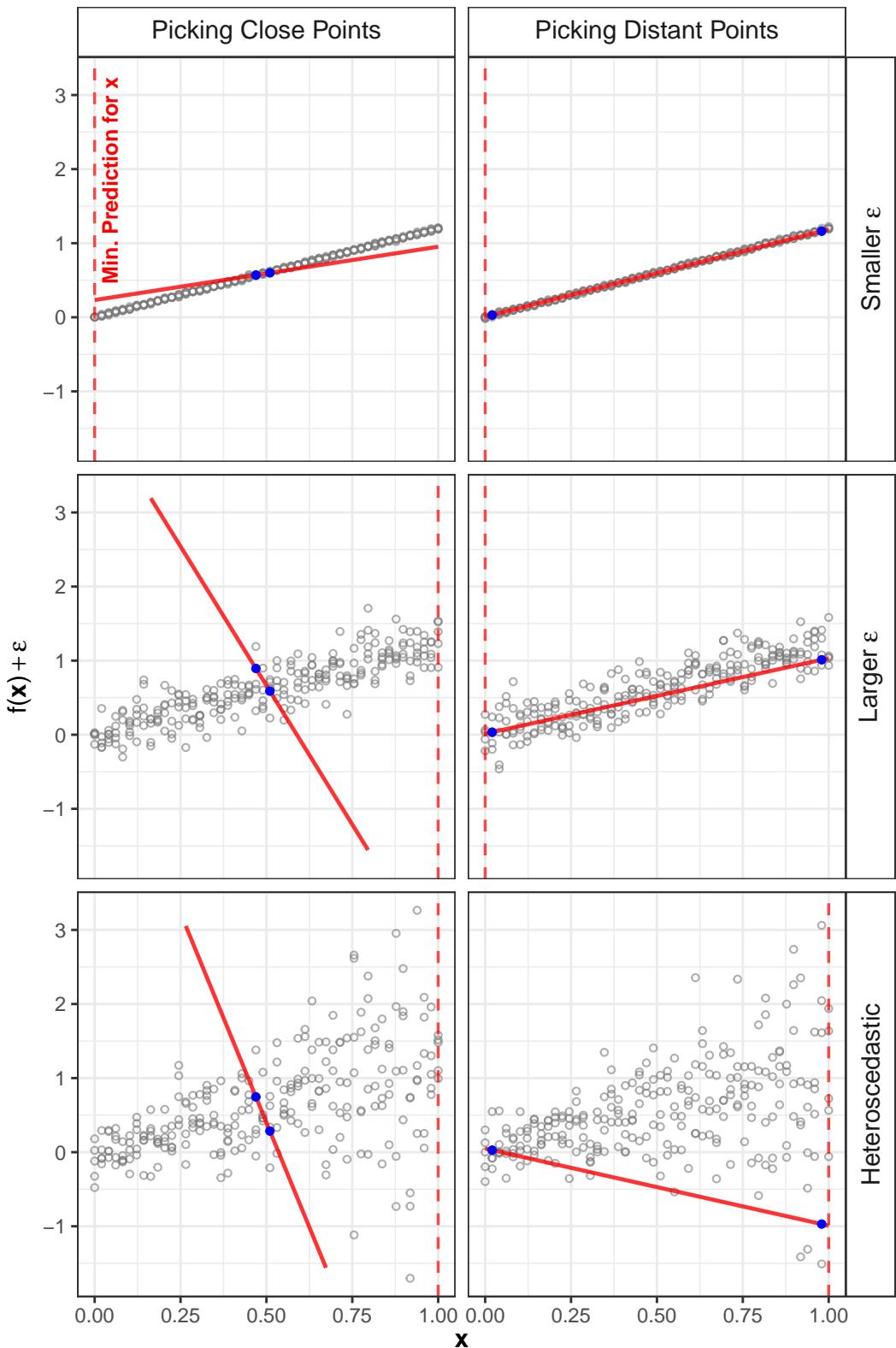


Figure 6.1: Fitting linear models using only two measurements, in experimental settings with different properties of the error ε

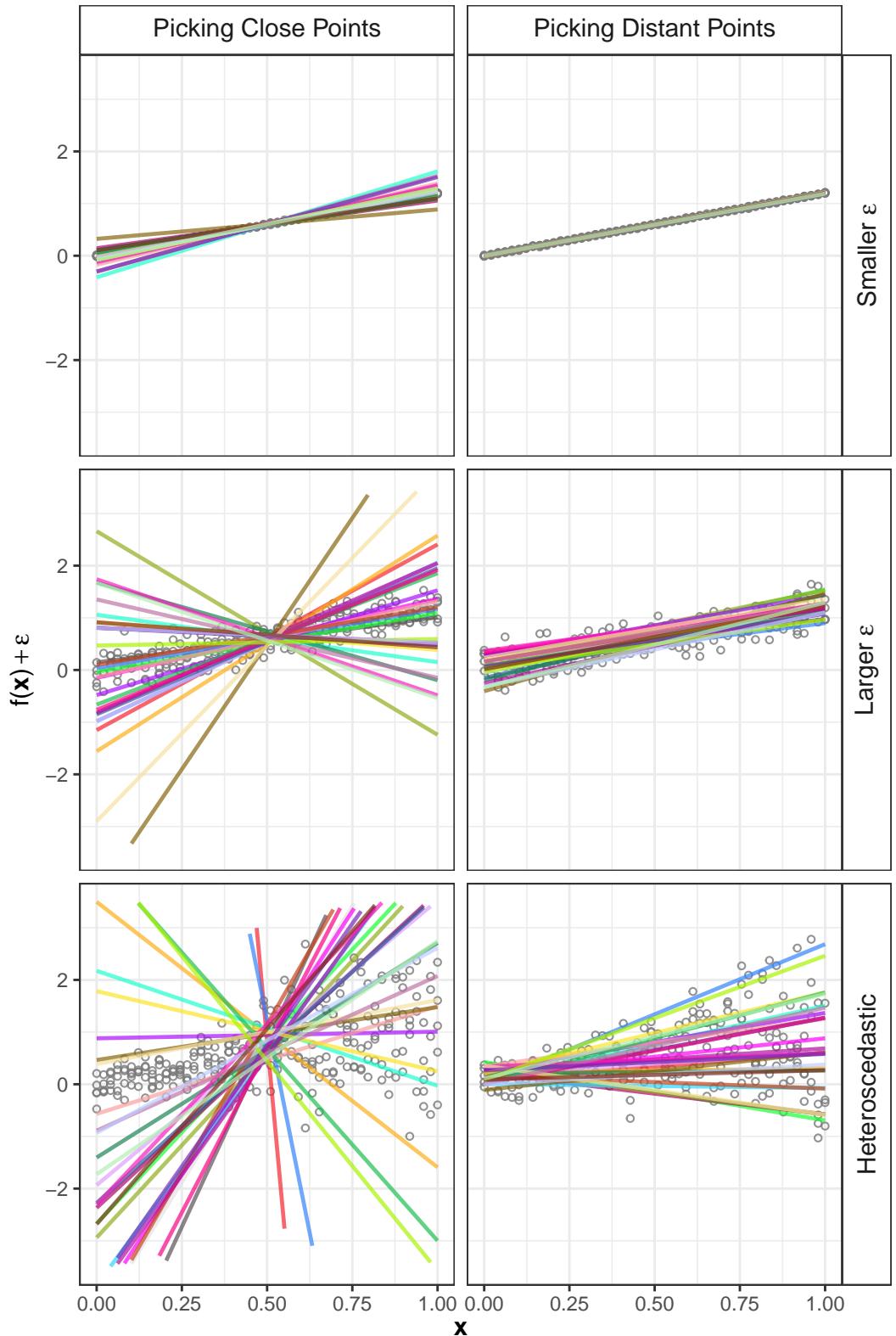


Figure 6.2: Fitting linear models using 30 distinct random sets of two measurements, in experimental settings with different properties of the error ε

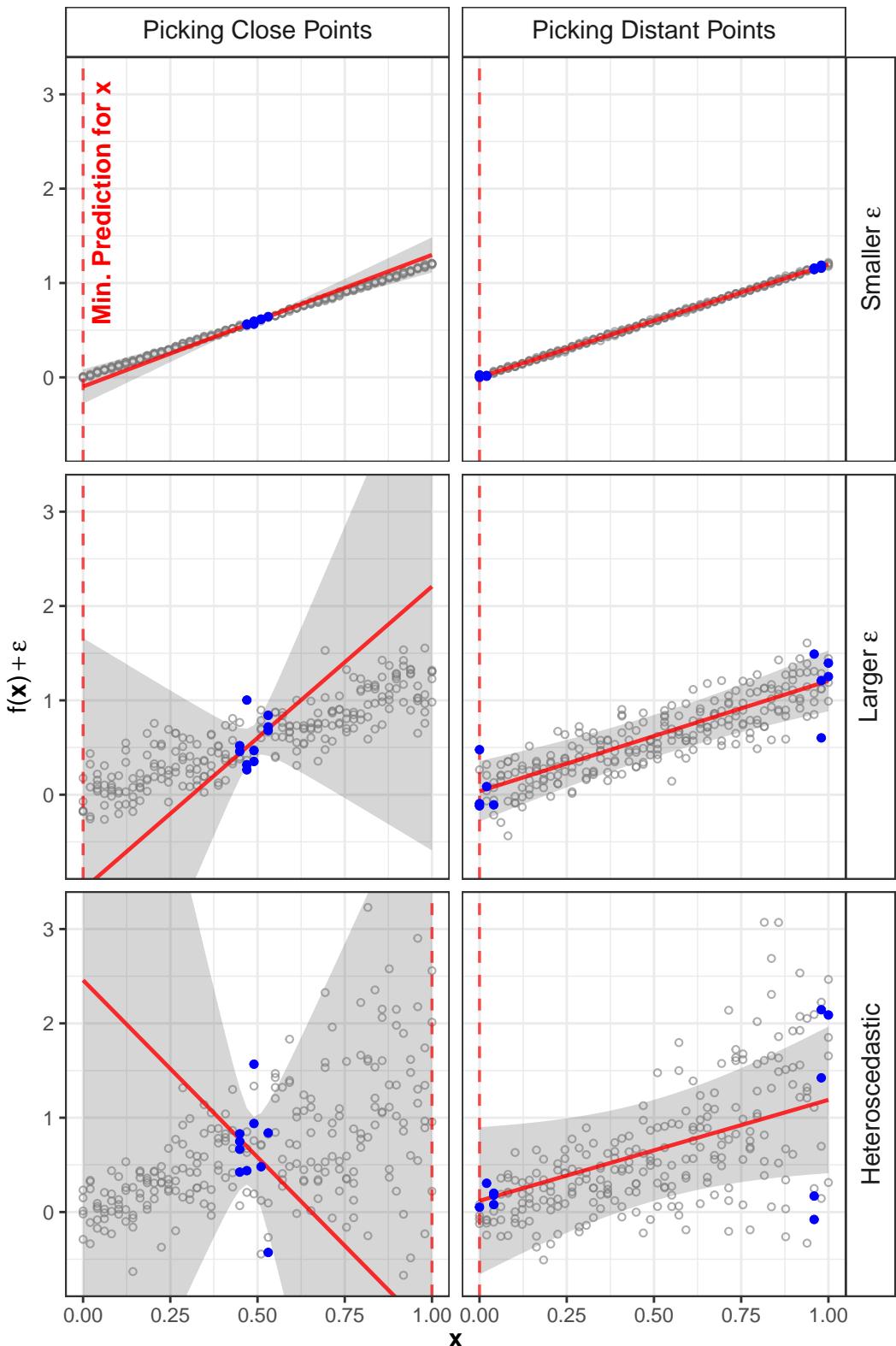


Figure 6.3: Modeling linear effects in different experimental settings

Table 6.1: A 2-level full factorial design for 3 factors, with response $y = f(x_1, x_2, x_3) + \varepsilon$

x_1	x_2	x_3	y
-1	-1	1	1.82
1	1	1	-3.30
1	1	-1	0.01
-1	-1	-1	5.36
1	-1	1	0.09
1	-1	-1	3.29
-1	1	-1	-1.07
-1	1	1	-2.46

In general, picking far-apart points will help estimating linear effects. This principle keeps working in higher-dimension search spaces, and is the motivation for 2-level full factorial designs. The first step to construct such a design is to choose extreme *levels* for each of the p factors, as we did in the right column of Figure 6.1, and then *encode* these levels, mapping their original values to the $[-1, 1]$ interval. The design is produced by listing all 2^p combinations of factor levels.

Table 6.1 shows a 2-level full factorial design for a search space with $p = 3$ factors and $2^p = 8$ experiments. We added a fourth column to the experimental design matrix, representing the 8 measurements. To produce a linear model surrogate \hat{f}_θ , it suffices to estimate θ as we did in Section 5.1, by computing

$$\hat{\theta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (6.1)$$

where \mathbf{X} contains the x_1, x_2, x_3 columns from the design matrix, and \mathbf{y} consists of the y column. With this design we are able to estimate eight model coefficients, one for the intercept, or mean, term, three for the linear effects of each factor, three for the two-way interactions, and one for the three-way interaction $x_1 x_2 x_3$.

If we fit models to distinct sets of two points picked at random, considering distant and close sets separately, we would get much more variance in the predictions of the models using close point sets, as we see in Figure 6.2. In other words, the linear model lacks robustness when designs are not constructed properly, especially for noisy or heteroscedastic measurements.

We can fit more robust linear models if we allow an expansion to our experimental budget and evaluate eight additional points in each experiment. Figure 6.3 shows the resulting model fits, this time accompanied by 95% confidence intervals for predictions. We observe great improvements in the heteroscedastic scenario due to the extra samples of f we were able to get, but we see the confidence intervals are very large. This example allows pointing out that the *uncertainty* in our model fits, which is proportional to the narrowness of the confidence intervals, is related to the distribution of experiments and to the variance associated with measuring f .

The previous practical examples show that the choice of the design matrix \mathbf{X} impacts the variance of the estimator of θ in the model fit in Equation 6.1. We will present a more precise statement of this interpretation in Section 6.3, and in the next section we will discuss the impact of interactions on our estimates, and present *fractional* factorial designs.

6.1.2 Interactions and Fractional Factorial Designs

Figure 6.4 shows two different objective functions of two input factors x_1 and x_2 , which can take real values in the interval $[0, 1]$. The top and bottom panels show the measures from a full factorial design using the interval extremes of both factors. We see the impact on the measured values of f when we change x_1 from zero to one, while keeping x_2 fixed to zero, represented by the red line, or fixed to one, represented by the blue line.

The function on the top panel has no interaction between x_1 and x_2 , and we see that changing x_2 from zero to one increases the value of the measured f , but does not change the slope of the line. These parallel lines can be represented in a linear model with two separate coefficients for x_1 and x_2 . The bottom panel shows measurements using the same design, evaluated for an objective function where an interaction between x_1 and x_2 is present. There is still an additive effect, since the two lines start at different values, but now the slopes do not match. In a linear model we can represent this interaction as a multiplicative effect, that is, in addition to separate coefficients for each factor there is a third coefficient for the multiplication $x_1 x_2$.

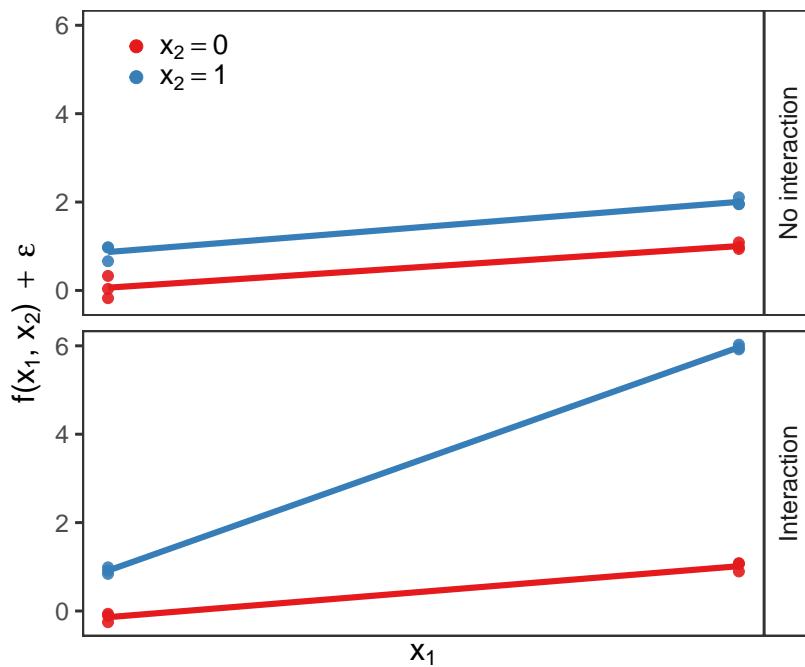


Figure 6.4: Effect of interactions in the slopes for extreme values of factors

Table 6.2: A 2-level fractional factorial design of resolution 3 for 5 factors, with response $y = f(x_1, x_2, x_3, x_4, x_5) + \varepsilon$. The main effects for x_4 and x_5 are aliased by, or confounded with, the interaction effects x_1x_2 and x_1x_3

x_1	x_2	x_3	$x_4 = x_1x_2$	$x_5 = x_1x_3$	y
1	-1	1	-1	1	5.10
1	1	1	1	1	4.47
1	-1	-1	-1	-1	0.54
-1	-1	-1	1	1	-0.79
-1	-1	1	1	-1	-2.39
1	1	-1	1	-1	-1.70
-1	1	-1	-1	1	-1.78
-1	1	1	-1	-1	-1.97

All interactions can be estimated with full factorial designs, but the cost of evaluating all combinations of extreme levels is prohibitive for most high-dimensional problems. It is common practice, at least on the initial stages of the statistical modeling of a problem, to consider that higher order interactions, such as those between three or more factors, are not of practical significance. This hypothesis supports the construction and analysis of *fractional* factorial designs.

A fractional factorial design for p factors consists of a full factorial design for $p - k$ factors, augmented with k columns produced by multiplying two or more of the original $p - k$ columns. Table 6.2 shows a fractional factorial design for $p = 5$ and $k = 2$, said to have *resolution 3*, with $2^{5-2} = 8$ experiments.

We can use this design to estimate the six coefficients we need for the for the main effects of each factor plus the intercept, but the estimates for the main effects of x_4 and x_5 will be multiplied by the interaction effects of the combined columns. Effects that are indistinguishable because of the way the design is constructed are said to be aliased by, or confounded with, each other. We will only be able to accurately estimate main effects in this setting if the interaction confounded with a main effect has small or nonexistent impact on the response.

Fractional factorial designs enable reducing the cost of estimating main effects, but the quality of the estimates relies on weak factor interactions. In the next section we will see how we can reduce the experimental costs even more, while simultaneously solving the confounding problem, by carefully constructing orthogonal arrays.

6.2 Parsimony for Linear Effects with Screening Designs

Screening designs identify parsimoniously the main effects of 2-level factors in the initial stages of studying a problem. While interactions are still not considered at this stage, identifying main effects early enables focusing on a smaller set of factors on subsequent

experiments. A specially efficient design construction method for screening designs was presented by Plackett and Burman [71] in 1946, and is available in the *FrF2* package [72] of the *R* language [73].

Despite having strong restrictions on the number of factors they support, *Plackett-Burman* designs enable the identification of main effects of p factors with $n = p + 1$ experiments. Factors may have many levels, but Plackett-Burman designs can only be constructed for 2-level factors. Therefore, like we did for fractional and full factorial designs, before constructing a Plackett-Burman design we must identify and encode *high* and *low* levels for each factor.

This section presents an example of *screening*, an interesting first stage on the sequential approach to building a linear model for a given objective function. The first step in the sequential approach is to build a design for the factors we wish to study. Then, we run experiments and analyze the results to find *significant* factors, using a statistical analysis method such as ANOVA. Next, we *restrict* the target search space by fixing significant factors to their best predicted levels, and start the process again. There are many caveats to the effectiveness of this approach, as we will discuss in Part II. In the next section we will go beyond screening, and discuss how we can explore more complex linear models using Optimal Design.

Because we are only looking at the two extreme levels for each factor, we must start from the hypothesis of linear relationships between factors and the response. This is always the case when using Plackett-Burman designs and two-level screening in general. We present how the screening approach works the example that follows. Suppose we wish to minimize an objective function written

$$f(\mathbf{x}) = -1.5x_1 + 1.3x_3 + 3.1x_5 + -1.4x_7 + 1.35x_8^2 + 1.6x_3x_5 + \varepsilon, \quad (6.2)$$

where x_1, \dots, x_8 are real numbers in the $[-1, 1]$ interval. For the purpose of this example, suppose the measurement of each point $f(\mathbf{x})$ is very expensive. Note that factors x_2, x_4 , and x_6 have no contribution to the response, and we can think of the error term ε as representing not only measurement noise, but our uncertainty regarding the model. Constructing a Plackett-Burman design will minimize the experiments needed to identify significant linear effects, and the analysis of this design will enable decreasing the dimension of the search space of interest.

Our initial modeling hypothesis for the screening experiment is written

$$\hat{f}_\theta(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\theta} + \varepsilon, \quad (6.3)$$

and Table 6.3 shows the Plackett-Burman design we generated for testing it. The design consists of an orthogonal matrix with twelve rows and eleven columns, containing the extreme levels of factors x_1, \dots, x_8 plus three “dummy” factors d_1, d_2, d_3 , completing the columns

Table 6.3: Randomized Plackett-Burman design for factors x_1, \dots, x_8 , with experiments $\mathbf{x}_1, \dots, \mathbf{x}_{12}$, and added “dummy” factors d_1, d_2, d_3 . The response y is subject to measurement error ε

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	d_1	d_2	d_3	$y = f(\mathbf{x}) + \varepsilon$
1	1	-1	1	1	1	-1	-1	-1	1	-1	8.00
1	1	-1	-1	-1	1	-1	1	1	-1	1	8.45
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	4.67
1	-1	1	1	1	-1	-1	-1	1	-1	1	12.48
-1	1	1	1	-1	-1	-1	1	-1	1	1	11.93
-1	-1	1	-1	1	1	-1	1	1	1	-1	14.27
1	-1	1	1	-1	1	1	1	-1	-1	-1	11.70
-1	-1	-1	1	-1	1	1	-1	1	1	1	2.96
1	1	1	-1	-1	-1	1	-1	1	1	-1	7.53
1	-1	-1	-1	1	-1	1	1	-1	1	1	10.69
-1	1	1	-1	1	1	1	-1	-1	-1	1	8.85
-1	1	-1	1	1	-1	1	1	1	-1	-1	10.39

needed to construct a Plackett-Burman design for 8 factors [71]. The final design matrix also contains a column for the computed response, subject to measurement error.

We fit a linear model corresponding to our initial hypothesis in Equation 6.3, and identify the most significant factors with an ANOVA test. The test results are shown in Table 6.4, where the *significance* of each factor is interpreted from $p(> F)$ values. We interpret this p -value as the probability that we would observe a value of the F statistic at least as large as what we did if the corresponding factor did not have an impact on the response. This approach is subject to the usual caveats regarding p -values, and it does not consider the practical significance of an effect.

Table 6.4 uses a “*”, as is convention in R , to represent significance levels for each factor. The table shows that factors x_1, x_3, x_5, x_7 , and x_8 have at least two “*” of significance, corresponding to a p -value smaller than 0.01. For the purpose of this example, this will be sufficient reason to include them in our linear model for the next step. We decide as well to discard factors x_2, x_4 , and x_6 from the next model, due to their low significance.

Table 6.4: ANOVA table for the fit of the initial model, with default significance levels from R

Term	Sum of Squares	F-value	$p(> F)$	“Significance”
x1	2.78	35.91	9.3×10^{-3}	**
x2	0.22	2.79	1.9×10^{-1}	
x3	38.90	501.80	1.9×10^{-4}	***
x4	0.75	9.67	5.3×10^{-2}	.
x5	25.35	327.08	3.7×10^{-4}	***
x6	0.99	12.83	3.7×10^{-2}	*
x7	4.93	63.59	4.1×10^{-3}	**
x8	43.84	565.55	1.6×10^{-4}	***

Table 6.5: Comparison between the response y predicted by our linear model fit and the true global optimum. Factors that were used in the model are in bold

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	y
<i>Global Optimum</i>	1.00	–	-1.00	–	-1.00	–	1.00	-0.40	-8.998
<i>Screening Prediction</i>	-1.00	–	-1.00	–	-1.00	–	1.00	-1.00	-1.258

Now, we fit a linear model with factors x_1 , x_3 , x_5 , x_7 , and x_8 , using the measurements of y we obtained when running our design in Table 6.3, and use the estimated model coefficients to predict the levels of each factor that minimize the real response. We can do that because these factors are numerical, even though only discrete values are allowed in the screening design.

We now proceed to the prediction step, where the objective is to identify the factor levels that minimize our fitted model, without running any new experiments. This can be done in practice by generating a large sample of experiments, predicting the value of f for each one using the estimated coefficients, and picking the experiment with the smallest predicted response. Table 6.5 compares the prediction for the best value of f , made with our linear model with significant factors, with the actual global optimum of f .

Using only 12 measurements and a simple linear model, the predicted best value of f was around nine times larger than the global optimum. Despite that, note that our model predicted the correct levels for x_3 , x_5 , and x_7 . The linear model predicted wrong levels for x_1 , perhaps due to this factor's interaction with x_3 , and for x_8 . It would be impossible to predict the correct level for x_8 using only linear terms since a quadratic term composes the true formula of f , shown in Equation 6.2.

We can improve upon this result if we introduce some information about the problem and use a more flexible design construction technique. Next, we will discuss the construction of efficient designs using problem-specific formulas and continue the optimization of our example.

6.3 Increasing Flexibility: Optimal Design

The application of DOE to autotuning problems requires design construction techniques that support factors of arbitrary types and number of levels. Autotuning problems typically combine factors such as binary flags, integer and floating point numerical values, and unordered enumerations of abstract values. Because Plackett-Burman designs only support 2-level factors, we had to restrict factor levels to interval extremities in the example of the previous section. We have seen that this restriction makes it impossible to test hypotheses using quadratic terms. We will continue to explore the search space of our example by using *D-Optimal designs*, which increase the number of levels we can efficiently screen for and enables testing more complex linear model hypotheses.

Like we did for screening, to construct a D-Optimal design we must first determine initial modeling hypotheses, which can be done based on previous experiments and analyses, or on knowledge about the problem. This initial model can be more complex than in the screening example, using any basis function set for linear models, as we discussed in Section 5.1.1, including polynomial terms of higher orders, interactions, and other non-linear relationships between factors.

We can then use *exchange* algorithms to construct an optimal design for our selected hypothesis, optimizing different criteria related to the trace, determinant, condition number, and other properties of the design's model matrix. We refer the reader to Chapters 9 through 12 of Atkinson *et al.* [74] for discussions on Optimal Design theory, the general equivalence between design optimality criteria, and the construction of *D-Optimal* designs using the *KL Exchange* algorithm. In the following examples, and in our work with optimal designs for autotuning, we have used D-Optimal designs constructed with Fedorov's algorithm [75], a particular case of KL Exchange, implemented in the *AlgDesign R* package [76].

Exchange algorithms for optimal designs build a design model matrix \mathbf{X} by selecting experiments $\mathbf{x}_1, \dots, \mathbf{x}_n$ in a search space \mathcal{X} that optimize a given design metric. The *D-Criterion* is one such metric, which measures the variance of the estimator $\hat{\theta}$ for the coefficients of a linear model, written

$$\text{Var}(\hat{\theta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2, \quad (6.4)$$

where σ^2 is the variance of measurements of the objective function. The D-Criterion $D(\mathbf{X})$ of an $n \times p$ design model matrix is a measure of the determinant of the information matrix related to the linear model we wish to test, and is written

$$D(\mathbf{X}) = \det\left(\frac{\mathbf{X}^\top \mathbf{X}}{p}\right)^{\frac{1}{n}}. \quad (6.5)$$

The higher the value of the D-Criterion, the smaller will be the variance of estimators. The search for the best design model matrix is done by swapping experiments from the current candidate design with experiments from a pool of possible experiments until some stopping criterion is met.

Returning to our screening example from the previous section, suppose that, in addition to using our previous screening results, we decide to consult an expert in the problem represented by our objective function, in Equation 6.2. The expert mentions that factor x_1 typically interacts strongly with factor x_3 , and also that we should replace include a quadratic term for x_8 , because there are some management overheads related to it. In Chapter 12 we will see a real example of the impact of such expert-guided models on the effectiveness of autotuning.

Using this improved performance model, built with a screening experiment and expert knowledge, we can construct a D-Optimal design using Fedorov's algorithm. We have much

Chapter 6. Design of Experiments

Table 6.6: D-Optimal design with 12 experiments, constructed for the surrogate $\hat{f}_\theta \sim x_1 + x_3 + x_5 + x_7 + x_8 + x_8^2 + x_1x_3$, with corresponding measurements

x_1	x_3	x_5	x_7	x_8	$y = f(\mathbf{x}) + \varepsilon$
-1.0	-1.0	-1.0	-1.0	-1.0	1.85
1.0	1.0	1.0	-1.0	-1.0	7.36
-1.0	1.0	-1.0	1.0	-1.0	-2.22
1.0	-1.0	1.0	1.0	-1.0	-1.54
1.0	-1.0	-1.0	-1.0	0.0	-6.30
-1.0	1.0	1.0	-1.0	0.0	5.55
-1.0	-1.0	-1.0	1.0	0.0	-2.92
1.0	1.0	1.0	1.0	0.0	2.86
1.0	-1.0	-1.0	-1.0	1.0	-4.78
-1.0	1.0	1.0	-1.0	1.0	7.03
1.0	1.0	-1.0	1.0	1.0	-1.60
-1.0	-1.0	1.0	1.0	1.0	4.39

more flexibility than screening when using exchange algorithms, being able to choose more model terms and also how many experiments we would like to run. We still need sufficient degrees of freedom to fit our new model, and we use Fedorov's algorithm to produce the 12 experiments shown in Table 6.6. Note that the design includes levels $-1, 0$, and 1 for factor x_8 . Because we have added a quadratic term, the optimal design will contain factor levels other than the extremes.

Table 6.7 shows the results of the new ANOVA test, where all factors lie below our chosen p -value threshold of 0.01. Note that if we consider that x_8^2 is significant we also include the corresponding linear term. We can then fit this model using the results of the optimal design we constructed. Table 6.8 show how close the coefficients estimated by this model fit are to the true model coefficients, from Equation 6.2. This example illustrates that the DOE methodology can help experimenters produce accurate models while using few measurements, provided the modeling hypotheses are correct.

Table 6.9 compares the true global optimum to the predictions made by the initial screening model from the previous section, and to the improved model fit to the experiments from the optimal design. Using screening, D-Optimal designs, and domain knowledge to guide modeling hypotheses, we predicted a set of factor levels whose response lies within 10% of

Table 6.7: ANOVA table for the fit of the improved model, with default significance levels from R

Term	Sum of Squares	F-value	$p(> F)$	"Significance"
x_1	26.10	701.66	1.2×10^{-5}	***
x_3	66.73	1794.21	1.9×10^{-6}	***
x_5	97.13	2611.61	8.8×10^{-7}	***
x_7	11.51	309.40	6.1×10^{-5}	***
x_8	0.02	0.57	4.9×10^{-1}	
x_8^2	6.11	164.28	2.1×10^{-4}	***
x_1x_3	30.75	826.70	8.7×10^{-6}	***

6.4. Exploration with Space-Filling Designs

Table 6.8: Model fit comparing real and estimated coefficients , with significance measurements as standard in *R*

	True Coefficient	Estimate	t-value	$p(> t)$	"Significance"
Intercept	0.00	-0.20	-2.11	1.0×10^{-1}	
x_1	-1.50	-1.47	-26.49	1.2×10^{-5}	***
x_3	1.30	1.35	22.90	2.2×10^{-5}	***
x_5	3.10	3.02	51.10	8.8×10^{-7}	***
x_7	-1.40	-1.55	-26.17	1.3×10^{-5}	***
x_8	0.00	-0.05	-0.76	4.9×10^{-1}	
x_8^2	1.35	1.51	12.82	2.1×10^{-4}	***
x_1x_3	1.60	1.70	28.75	8.7×10^{-6}	***

the global optimum, performing 24 measurements of f . We did that by first reducing the problem's dimension when we eliminated insignificant factors in the screening step. Then, we performed a more careful exploration of this new subspace, aided by domain knowledge provided by an expert. Note that we could have reused some of the 12 experiments from the previous step to reduce even further the size of the new design.

Because of the carefully constructed experiments, we are able to *explain* the improvements we obtained in each step of the process using results of statistical tests. If we needed to optimize a similar problem, we could start by *exploiting* the relationships we have discovered, but we could also attempt to *explore* new hypotheses for the factors we did not include at first. This balance between exploitation and exploration will be discussed further in Chapter 7.

In the absence of reasonable hypotheses about the search space, we can still construct designs that favor a balanced exploration, enabling tests for more complex hypotheses and increasing the effectiveness of nonparametric optimization methods. The next section presents some methods for constructing such designs.

6.4 Exploration with Space-Filling Designs

For large integers n , most of the volume of an n -sphere lies in a thin outer shell. Therefore, in search spaces with a large number of p factors, such as the ones we find in some autotuning search spaces, uniform random samples concentrate in the extreme levels, which can make it impossible to test more complex linear models and to properly evaluate lack of fit.

Table 6.9: Comparison between the response y predicted by our screening and optimal design model fits and the true global optimum. Factors that were used in the models are in bold

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	y
<i>Global Optimum</i>	1.0	–	-1.0	–	-1.0	–	1.0	-0.4	-8.998
<i>Screening Prediction</i>	-1.0	–	-1.0	–	-1.0	–	1.0	-1.0	-1.258
<i>Optimal Design Prediction</i>	1.0	–	-1.0	–	-1.0	–	1.0	0.0	-8.894

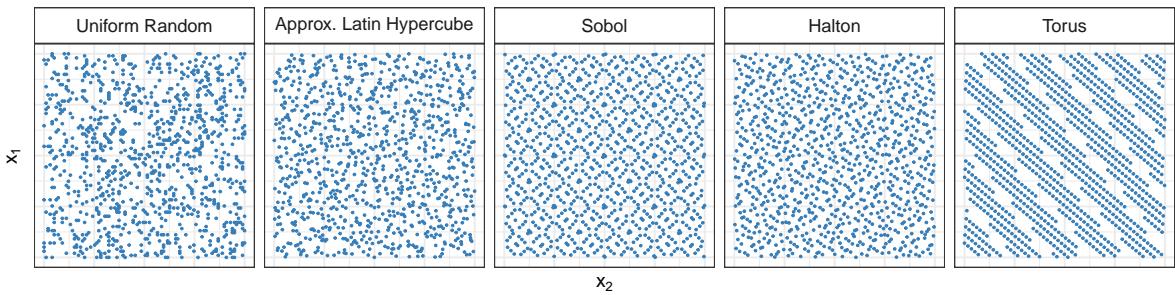


Figure 6.5: Comparison of 1000 samples obtained with different sampling methods. The *Latin Hypercube* design was produced with the *DiceDesign R* package [80], and was optimized after an initial sampling, with Simulated Annealing. Sobol, Halton, and Torus are low-discrepancy sequences from the *randtoolbox R* package [79].

Space-filling designs are capable of generating better-distributed samples, even for high-dimensional search spaces, enabling more precise measurements of model quality. These designs also help to improve the accuracy of acquisition functions for nonparametric optimization methods, some of which we discuss in Chapter 7.

Figure 6.5 compares 5 sampling methods in a two-dimensional search space. Even in low dimension, we can see how uniform random sampling produces clusters of points and leaves large regions unsampled. *Latin Hypercube Sampling* [81] (LHS) is one kind of space-filling design, constructed by partitioning the search space and attempting to distribute samples evenly across partitions. It is still hard to cover high-dimensional spaces, and it is common to have to optimize an LHS design with function minimization methods.

Discrepancy is one criterion typically used to evaluate how evenly a sample is distributed. For a given region of a search space, discrepancy can be described intuitively as the ratio between the number of samples that lie in the region and the size of the region. In a *low-discrepancy* design, any arbitrary small region of a high-dimensional search space contains a proportionately small number of samples. Sobol, who is also responsible for sensitivity indices [68] from the previous chapter, and Halton introduced deterministic *quasirandom* low-discrepancy sequences [77, 78] that are extensively used. Figure 6.5 also shows these sequences, implemented in the *randtoolbox R* package [79], which also implements a *Torus* low-discrepancy sampler.

6.5 Comparing Sampling Methods

Figure 6.6 shows the contour of the search space defined by the Booth function and its exploration by six different sampling methods, restricted to 30 measurements. Here, a uniform random sample finds points close to the global optimum, even though measurements are not particularly distributed and there is no attempt to exploit search space structure. An LHS sample covers the search space more evenly, but still does not exploit structure. The Sobol

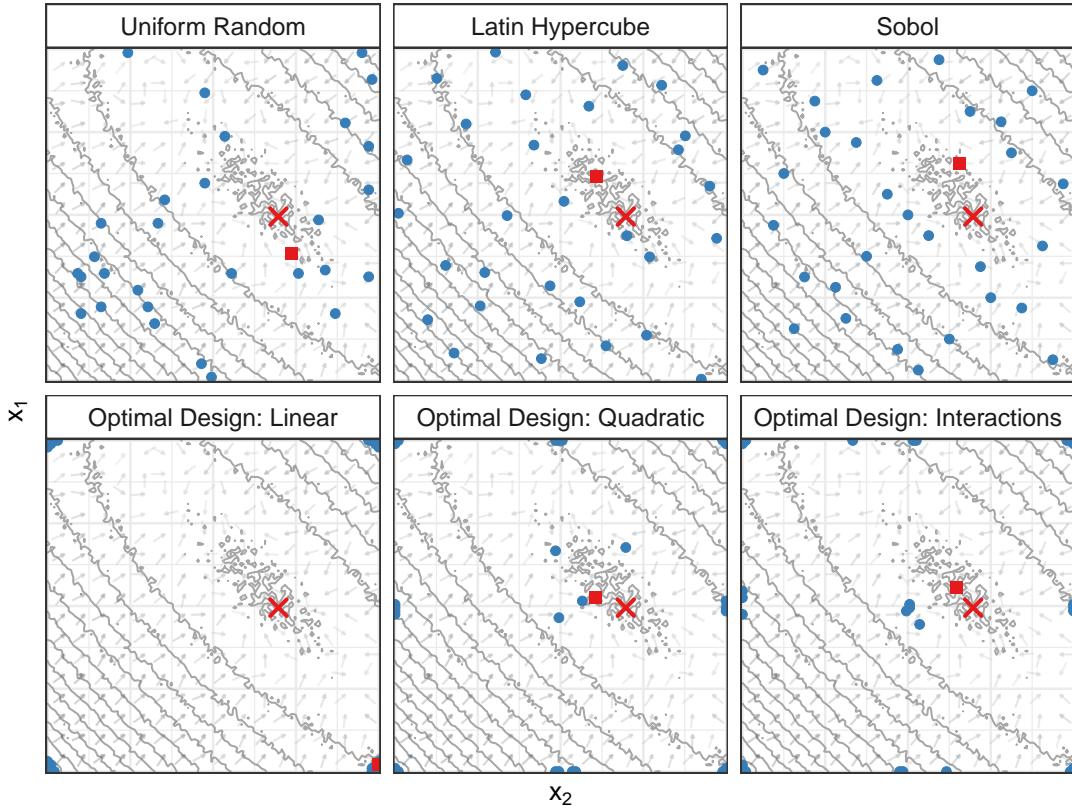


Figure 6.6: Exploring the search space for Booth’s function with different methods, using a fixed budget of 30 measurements. A \times marks the global optimum, a \blacksquare marks the best point found by each strategy, and \bullet s mark measurements

sample is much more useful than the others for such a small space, but we can clearly see the impact of the optimal designs built for different models.

The D-optimal design for a model with linear terms samples the corners of the search space, and adding quadratic and interaction terms to the model adds samples to the inner regions. Predictions for models with these added terms are closer to the global optimum, but the selection of the experiments included in the designs is still made over an initial sample of the search space. Therefore, combining optimal design and space-filling designs should always be considered.

6.6 Summary

- Designs to obtain good-quality parameter estimates
 - Screening and D-opt for LM
 - SFD for GP
- Designs to test the model quality (lack of fit)

Chapter 6. Design of Experiments

- SFD fo LM
- If model based, parameter significance and estimation can be used to reduce dimension and guide the optimization.
- With this DOE approach, we have a clear separation between the sampling phase and the interpretation phase. But what if no parameter really appears significant anymore?
- Paley construction [82]

Chapter 7

Online Learning

TODO: Fill in the sketched section below

7.1 Independent Bandits

- simple (discrete choice, optimize $regret = \sum_t R_t$)
- UCB

7.2 Expected Improvement for Gaussian Process Regression

- (continuous choice, optimize $EI = \max_t R_t$)
- Also mention GP-UCB and contrast with EGO. There are also variants for Linear Model (LinUCB)

7.3 Reinforcement Learning

- Unsupervised learning
- Actor-critic model
- Main motivation for adding this section is to contextualize results with GPR and Neural Nets. Could move to last Application chapter

7.4 Summary

- These methods seamlessly mix exploration and exploitation but the overall objective function is generally the regret, which makes sense for a self-optimizing system (e.g. facebook) but not in an autotuning context (where EI is more meaningful)

Chapter 8

Optimization Methods for Autotuning: An Overview

TODO: Revisit semi-conclusions of each chapter, wrap up this part, hook to Applications part

- use glassbox (DOE based) approach to perform the optimization, always try to interpret the results
- 2 big methods based on different exploration/exploitation strategy:
 1. Evaluate parameter significance and reduce dimension (two phases, iterative)
 2. Expected Improvement (first a general exploration phase with a SFD, then seamlessly mix exploration and exploitation)

Possibly combinations of both approaches could be used back and forth depending on the specific information we learn on the use case.

- In this thesis, we evaluate these DOE-based approaches for several autotuning use cases and try to compare them with approaches that had been previously proposed for these use cases.

TODO: Adapt this paragraph

The effectiveness of stochastic descent methods on autotuning can be limited [5, 83, 84], between other factors, by underlying hypotheses about the search space, such as the reachability of the global optimum and the smoothness of search space surfaces, which are frequently not respected. The derivation of relationships between parameters and performance from search heuristic optimizations is greatly hindered, if not rendered impossible, by the biased way these methods explore parameters. Some parametric learning methods, such as Design

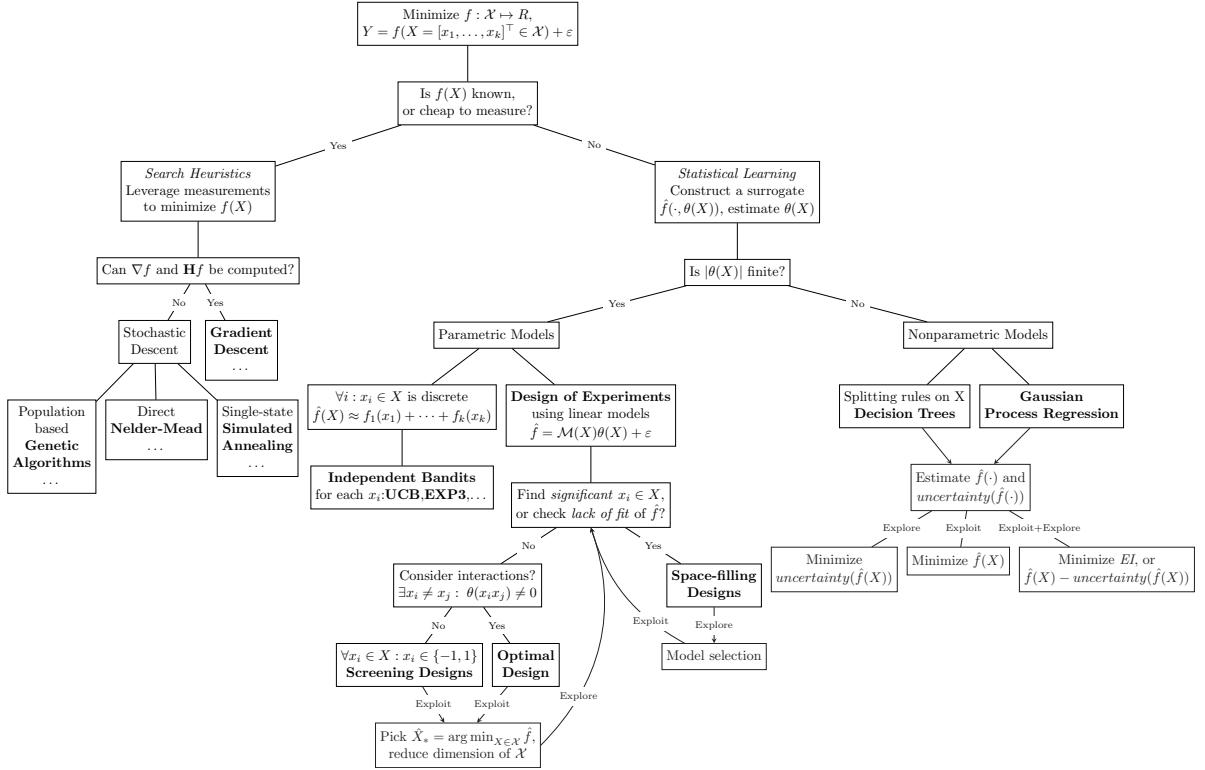


Figure 8.1: Some autotuning methods

of Experiments, are not widely applied to autotuning. These methods perform structured parameter exploration, and can be used to build and validate performance models, generating transparent and cost-effective optimizations [34, 4]. Other methods from the parametric family are more widely used, such as Bandit Algorithms [36]. Nonparametric learning methods, such as Decision Trees [85] and Gaussian Process Regression [86], are able to reduce model bias greatly, at the expense of increased prediction variance. Figure 8.1 categorizes some autotuning methods according to some of the key hypotheses and branching questions underlying each method.

Part II

Applications

Chapter 9

Research Methodology: Efforts for Reproducible Science

Sketch:

- Tools and such
- Explain difficulty of finding a needle in a haystack:
 - how to know whether we found the optimal value ?
 - how to know how far we are from the optimal value ?
 - how to know whether there is anything to find ?
 - how to know whether the geometry hypothesis we make are sound in an unknown space ?
 - ...

9.1 Reproducible Research

Conclusions derived from experimental data can only be considered validated after it is possible to reproduce them under independent experimental conditions. This principle has guided experimental research and scientific progress. Experimental researchers on computer science are in a singular position to promote and practice reproducible research, because some computational experiments can be monitored, registered, and repeated with precision and control that are simply impossible in other experimental fields such as chemistry and biology.

9.2 Computational Documents with Emacs and Org mode

9.3 Versioning for Code, Text, and Data

Chapter 10

Compiler Parameters for CUDA Kernels

A Graphics Processing Unit (GPU) is a parallel computing coprocessor specialized in accelerating vector operations such as graphics rendering. General Purpose computing on GPUs (GPGPU) depends on accessible programming interfaces for languages such as C and Python that enable the use of GPUs in different parallel computing domains. The Compute Unified Device Architecture (CUDA) is a GPGPU toolchain introduced by the NVIDIA corporation.

In this study we measured the performance improvements obtained with an ensemble of stochastic methods for function minimization applied to the parameters of the CUDA compiler. We implemented an autotuner using the OpenTuner framework [1] and used it to search for the compilation parameters that optimize the performance of 17 heterogeneous GPU kernels, 12 of which are from the Rodinia Benchmark Suite [87]. We used 3 different NVIDIA GPUs in the experiments, the Tesla K40, the GTX 980, and the GTX 750. We published the results of this study, which we summarize in this chapter, at the *Concurrency and Computation: Practice and Experience* journal [2].

The optimizations we found often beat compiler high-level optimization options, such as the `-O2` flag. The autotuner found compilation options that produced speedups, in relation to `-O2`, of around 2 times for the *Gaussian Elimination* problem from the Rodinia Benchmark Suite, 2 times for *Heart Wall* problem, also from Rodinia, and 4 times for one of the matrix multiplication algorithms we implemented. We observed that the compilation parameters that optimize an algorithm for a given GPU architecture would not always achieve the same performance in different hardware, which was a reasonable expectation.

Figure 10.1 shows a representation of the autotuner implemented in this study, and illustrates the time scale of our experiments. The autotuner, in light blue, receives as input a GPU kernel, a target GPU, input data for the kernel, and a search space composed of NVCC flags. After running an ensemble of stochastic search methods, the autotuner outputs a flag

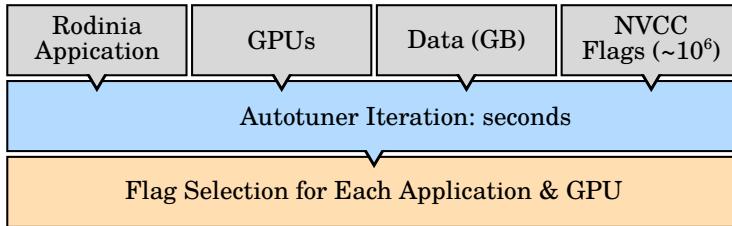


Figure 10.1: Autotuner representation and time scale of the experiments

selection. The compilation of the kernels we targeted in this study takes few seconds, so it was possible to test thousands of flag combinations per hour using a sequential autotuner.

The rest of this chapter is organized as follows. Section 10.1 discusses NVIDIA GPU architecture and programming toolchain, and presents related work on GPU performance tuning and modeling. Section 10.2 presents the search space of CUDA parameters and discusses the autotuner we implemented. Section 10.3 presents GPUs and the kernels whose compilation we attempted to optimize. Section 10.4 presents and discusses the performance improvements achieved by the autotuner, and attempts to identify significant flags using clustering. Section 10.5 was written at a later moment, after studying Experimental Design methods, and revisits this problem with a screening experiment. Finally, Section 10.6 summarizes the discussion and presents perspectives for future work.

10.1 General-Purpose Computing on NVIDIA GPUs

10.1.1 NVIDIA GPU Micro-Architecture

NVIDIA GPU architectures have multiple asynchronous and parallel Streaming Multiprocessors (SMs) which contain Scalar Processors (SPs), Special Function Units (SFUs) and load/store units. Each group of 32 parallel threads scheduled by an SM, or *warp*, is able to read from memory concurrently. NVIDIA architectures vary in a large number of features, such as number of cores, registers, SFUs, load/store units, on-chip and cache memory sizes, processor clock frequency, memory bandwidth, unified memory spaces and dynamic kernel launches. Those differences are summarized in the Compute Capability (C.C.) of an NVIDIA GPU.

The hierarchical memory of an NVIDIA GPU contains global and shared portions. Global memory is big, off-chip, has a high latency and can be accessed by all threads of the kernel. Shared memory is small, on-chip, has a low-latency and can be accessed only by threads in a same SM. Each SM has its own shared L1 cache, and new architectures have coherent global L2 caches. Optimizing thread accesses to different memory levels is essential to achieve good performance.

10.1.2 Compute Unified Device Architecture (CUDA)

The CUDA programming model and platform enables the use of NVIDIA GPUs for scientific and general purpose computation. A single *main* thread runs in the CPU, launching and managing computations on the GPU. Data for the computations has to be transferred from the *main memory* to the GPU's memory. Multiple computations launched by the main thread, or *kernels*, can run asynchronously and concurrently. If the threads from a same warp must execute different instructions the CUDA compiler must generate code to branch the execution correctly, making the program lose performance due to this *warp divergence*.

The CUDA language extends C and provides a multi-step compiler, called NVCC, that translates CUDA code to Parallel Thread Execution code (PTX). NVCC uses the host's C++ compiler in several compilation steps, and also to generate code to be executed in the host. The final binary generated by NVCC contains code for the GPU and the host. When PTX code is loaded by a kernel at runtime, it is compiled to binary code by the host's device driver. This binary code can be executed in the device's processing cores, and is architecture-specific. The target architecture can be specified using NVCC parameters.

10.1.3 GPU Performance Models and Autotuning

The accuracy of a GPU performance model is subject to low level elements such as instruction pipeline usage and small cache hierarchies. A GPU's performance approaches its peak when the instruction pipeline is saturated, but becomes unpredictable when the pipeline is under-utilized [88, 89]. Considering the effects of small cache hierarchies [90, 91] and memory-access divergence [92, 93] is also critical to a GPU performance model.

Guo and Wang [94] introduce a framework for autotuning the number of threads and the sizes of blocks and warps used by the CUDA compiler for sparse matrix and vector multiplication GPU applications. Li *et al.* [95] discuss the performance of autotuning techniques in highly-tuned GPU General Matrix to Matrix Multiplication (GEMMs) routines, highlighting the difficulty in developing optimized code for new GPU architectures. Grauer-Gray *et al.* [96] autotune an optimization space of GPU kernels focusing on tiling, loop permutation, unrolling, and parallelization. Chaparala *et al.* [97] autotune GPU-accelerated Quadratic Assignment Problem solvers. Sedaghati *et al.* [98] build a decision model for the selection of sparse matrix representations in GPUs.

10.2 Autotuner and Search Space for the NVCC Compiler

This section discusses our autotuner implementation and the search space of NVCC compiler parameters.

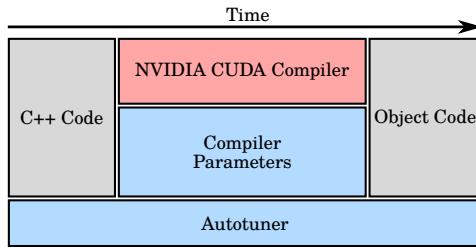


Figure 10.2: Simplified view of NVCC compilation

10.2.1 An Autotuner for CUDA Parameters using OpenTuner

The autotuner was implemented using the OpenTuner framework [1], and used a multi-armed bandit algorithm that aims to maximize the successes across the autotuning process. A success is defined as finding a program configuration that improves upon the best performance found so far. The stochastic method which finds such a configuration gets its score increased for that tuning session. A bandit algorithm, called *MABAUC* in the paper, for Multi-Armed Bandit Area Under the Curve, considers a sliding window covering a given number of measurements to compute the scores of methods in an ensemble, and to decide which method, or “arm”, to play next. The ensemble of methods we used was composed by implementations of the Nelder-Mead algorithm and three variations of genetic algorithms. Figure 10.2 shows a simplified version of the steps necessary to generate the object code that will be measured later. The code for our autotuner and all the experiments and results is available [99] under the GNU General Public License.

10.2.2 Search Space for CUDA Parameters

Table 10.1 presents a subset of the CUDA configuration parameters [100] used in this study. Parameters can target different compilation steps, namely the *PTX* optimizing assembler, the *NVLINK* linker, and the *NVCC* compiler. We compared the performance of programs generated by tuned parameters with the standard compiler optimizations *opt-level* equal to 0, 1, 2, 3. Different optimization levels could also be selected during tuning. We did not use compiler options that target the host linker or the library manager since they do not affect performance. The size of the search space defined by all possible combinations of flags is in the order of 10^6 , making hand-optimization or exhaustive searches unwieldy.

10.3 Target GPUs and Kernels

This section presents the GPU testbed, the algorithm benchmark, the autotuner implementation and its search space.

Table 10.1: Description of flags in the search space

Flag	Description
no-align-double	Specifies that <code>malign-double</code> should not be passed as a compiler argument on 32-bit platforms. Step: NVCC
use_fast_math	Uses the fast math library, implies <code>ftz=true</code> , <code>prec-div=false</code> , <code>prec-sqrt=false</code> and <code>fmad=true</code> . Step: NVCC
gpu-architecture	Specifies the NVIDIA virtual GPU architecture for which the CUDA input files must be compiled. Step: NVCC Values: <code>sm_20</code> , <code>sm_21</code> , <code>sm_30</code> , <code>sm_32</code> , <code>sm_35</code> , <code>sm_50</code> , <code>sm_52</code>
relocatable-device-code	Enables the generation of relocatable device code. If disabled, executable device code is generated. Relocatable device code must be linked before it can be executed. Step: NVCC
ftz	Controls single-precision denormals support. <code>ftz=true</code> flushes denormal values to zero and <code>ftz=false</code> preserves denormal values. Step: NVCC
prec-div	Controls single-precision floating-point division and reciprocals. <code>prec-div=true</code> enables the IEEE round-to-nearest mode and <code>prec-div=false</code> enables the fast approximation mode. Step: NVCC
prec-sqrt	Controls single-precision floating-point square root. <code>prec-sqrt=true</code> enables the IEEE round-to-nearest mode and <code>prec-sqrt=false</code> enables the fast approximation mode. Step: NVCC
def-load-cache	Default cache modifier on global/generic load. Step: PTX Values: <code>ca</code> , <code>cg</code> , <code>cv</code> , <code>cs</code>
opt-level	Specifies high-level optimizations. Step: PTX Values: <code>0 - 3</code>
fmad	Enables the contraction of floating-point multiplies and adds/subtracts into floating-point multiply-add operations (FMAD, FFMA, or DFMA). Step: PTX
allow-expensive-optimizations	Enables the compiler to perform expensive optimizations using maximum available resources (memory and compile-time). If unspecified, default behavior is to enable this feature for optimization level <code>>O2</code> . Step: PTX
maxrregcount	Specifies the maximum number of registers that GPU functions can use. Step: PTX Values: <code>16 - 64</code>
preserve-relocs	Makes the PTX assembler generate relocatable references for variables and preserve relocations generated for them in the linked executable. Step: NVLINK

10.3.1 Target GPU Architectures

To be able to show that different GPUs require different options to improve performance, and that it is possible to achieve speedups in different hardware, we wanted to tune our benchmark for different NVIDIA architectures. Table 10.2 summarizes the hardware characteristics of the three GPUs.

10.3.2 Benchmark of CUDA Algorithms

We composed a benchmark with 17 heterogeneous GPU kernels. The benchmark contains 4 code optimizations for *matrix multiplication* counted as a single kernel, 1 vector addition problem, 1 solution for the *maximum sub-array problem* [101], 2 *sorting* algorithms and 12 kernels from the Rodinia Benchmark Suite [87].

Table 10.2: Hardware specifications of the target GPU architectures

Model	C.C.	Global Memory	Bus	Bandwidth	L2	Cores/SM	Clock
Tesla-K40	3.5	12 GB	384-bit	276.5 GB/s	1.5 MB	2880/15	745 Mhz
GTX-750	5.0	1 GB	128-bit	86.4 GB/s	2 MB	512/4	1110 Mhz
GTX-980	5.2	4 GB	256-bit	224.3 GB/s	2 MB	2048/16	1216 Mhz

Table 10.3 shows the Rodinia kernels contained in our benchmark and the corresponding three-letter code. The other kernels in the benchmark were the following CUDA implementations:

- Matrix multiplications using:
 - Global memory with non-coalesced accesses (MMU)
 - Global memory with coalesced accesses (MMG)
 - Shared memory with non-coalesced accesses to global memory (MSU)
 - Shared memory with coalesced accesses to global memory (MMS)
- Simple vector addition algorithm (VAD)
- Solution for the Maximum Sub-Array Problem (MSA) [103, 101]
- Quicksort (QKS) and Bitonicsort (BTN)

10.4 Performance Improvements and Parameter Clustering Attempt

This section presents the speedups achieved for all algorithms in the benchmark, highlights the most significant speedups, and discusses the performance and accuracy of the autotuner.

Table 10.3: Rodinia [87] kernels used in the experiments

Kernel	Berkeley Dwarf[102]	Domain
B+Tree (BPT)	Graph Traversal	Search
Back Propagation (BCK)	Unstructured Grid	Pattern Recognition
Breadth-First Search (BFS)	Graph Traversal	Graph Algorithms
Gaussian Elimination (GAU)	Dense Linear Algebra	Linear Algebra
Heart Wall (HWL)	Structured Grid	Medical Imaging
Hot Spot (HOT)	Structured Grid	Physics Simulation
K-Means (KMN)	Dense Linear Algebra	Data Mining
LavaMD (LMD)	N-Body	Molecular Dynamics
LU Decomposition (LUD)	Dense Linear Algebra	Linear Algebra
Myocyte (MYO)	Structured Grid	Biological Simulation
Needleman-Wunsch (NDL)	Dynamic Programming	Bioinformatics
Path Finder (PTF)	Dynamic Programming	Grid Traversal

10.4.1 Performance Improvements

Figure 10.3 compares the means of 30 measurements of the tuned results for each kernel to the high-level optimization *opt-level* set to 2. We see that that the autotuned solution for the Heart Wall problem (HWL) in the *Tesla K40* achieved over 2 times speedup in comparison with the high-level CUDA optimizations. The tuned kernel for the *GTX-980* reached almost 2.5 times speedup for Gaussian Elimination (GAU). We also found smaller speedups for most kernels, such as a 10% speedup for Path Finder (PTF) in *Tesla K40*, and around 10% speedup for Myocyte (MYO) on *GTX-750*.

Figure 10.4 summarizes the results for the other CUDA kernels we implemented. The autotuner did not improve upon the high-level optimizations for BTN, VAD, and QKS in any of the GPUs in the testbed, but it found solutions that achieved speedups for at least one GPU for the other kernels.

We were not able to determine the hardware characteristics that impacted performance. We believe that the Maxwell GPUs, the *GTX-980* and *GTX-750*, had differing results from the *Tesla K40* because they are consumer grade GPUs, producing less precise results, and configured with different default optimizations. The similarities between the compute capabilities of the *GTX* GPUs could also explain the observed differences from the K40. The overall greater computing power of the *GTX-980* could explain its differing results from the *GTX-750*, since the *GTX 980* has a greater number of Cores, SMs/Cores, Bandwidth, Bus, Clock and Global Memory.

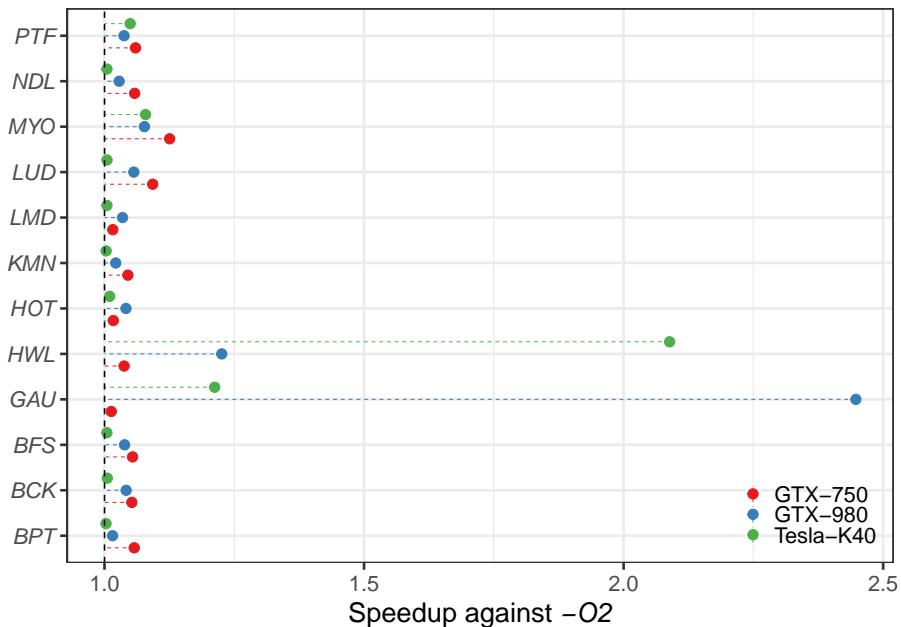


Figure 10.3: Mean speedup over 30 repetitions of the tuned solutions for Rodinia kernels versus the *opt-level* = 2 baseline

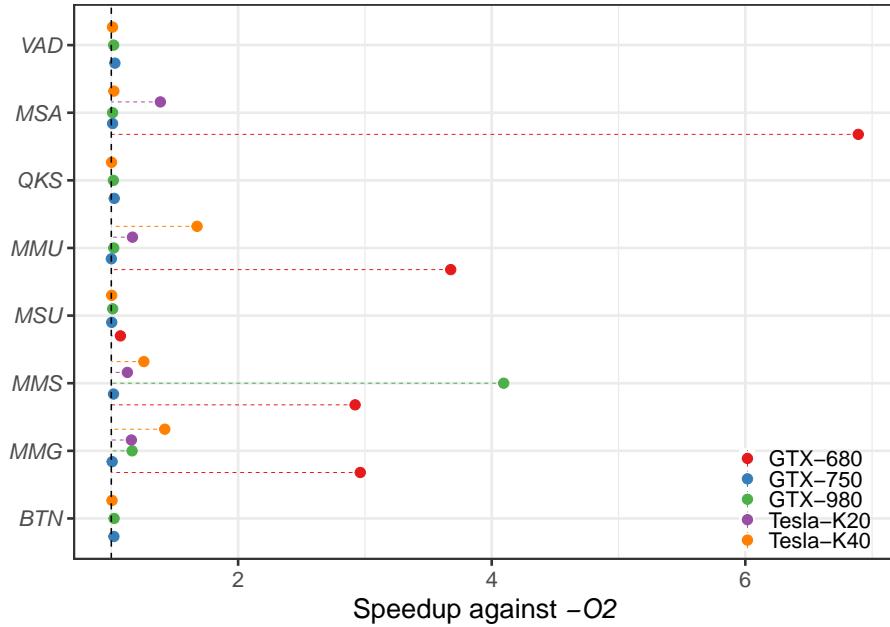


Figure 10.4: Mean speedup over 30 repetitions of the tuned solutions for kernels we implemented versus the *opt-level* = 2 baseline

10.4.2 CUDA Compiler Autotuner Performance

This section presents an assessment of the autotuner’s performance. Figures 10.5 and 10.6 present the speedup of the best solution found across tuning time. Points in each graph represent the performance, in the *y*-axis, of the best configuration found at the corresponding tuning time, shown in the *x*-axis. The leftmost point in each graph represents the performance of a configuration chosen at random by the autotuner. Each subsequent point represents the performance of the best configuration found across optimization. Note that the autotuner is able to quickly improve upon the initial random configuration, but the rate of improvement also decays quickly. The duration of all tuning runs was two hours, or 7200 seconds. The rightmost point in each graph represents the performance of the last improving configuration found by the autotuner.

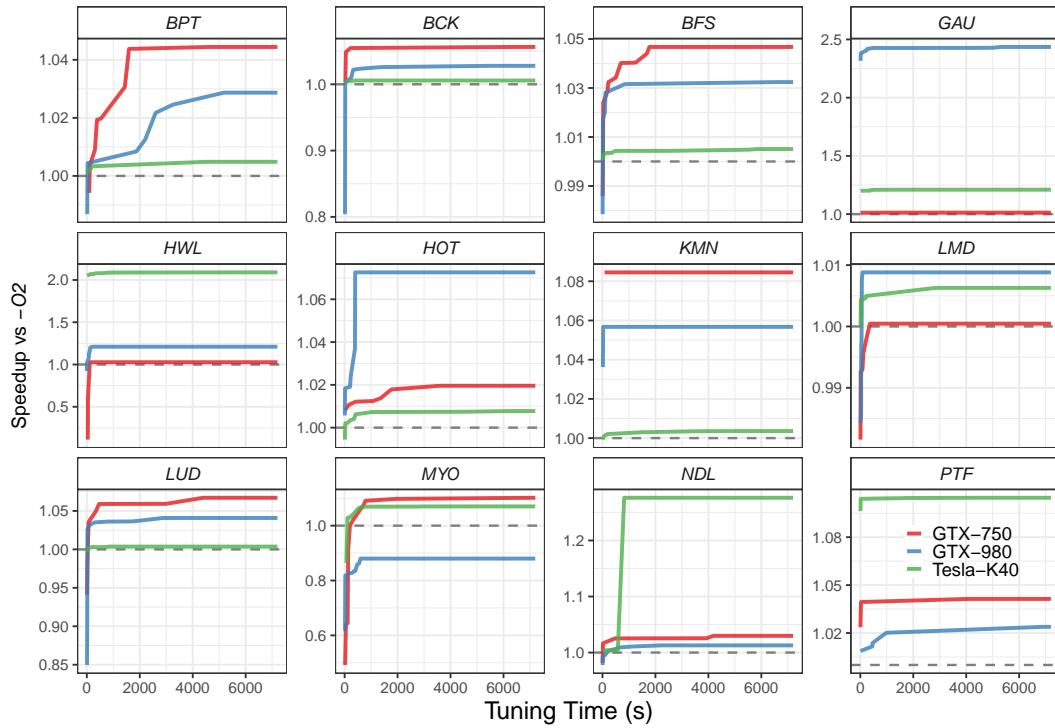


Figure 10.5: Mean speedup over 30 repetitions of the tuned solutions for Rodinia kernels versus the $opt-level = 2$ baseline, across two hours of tuning. Notice the difference in the y -axis scales for each panel

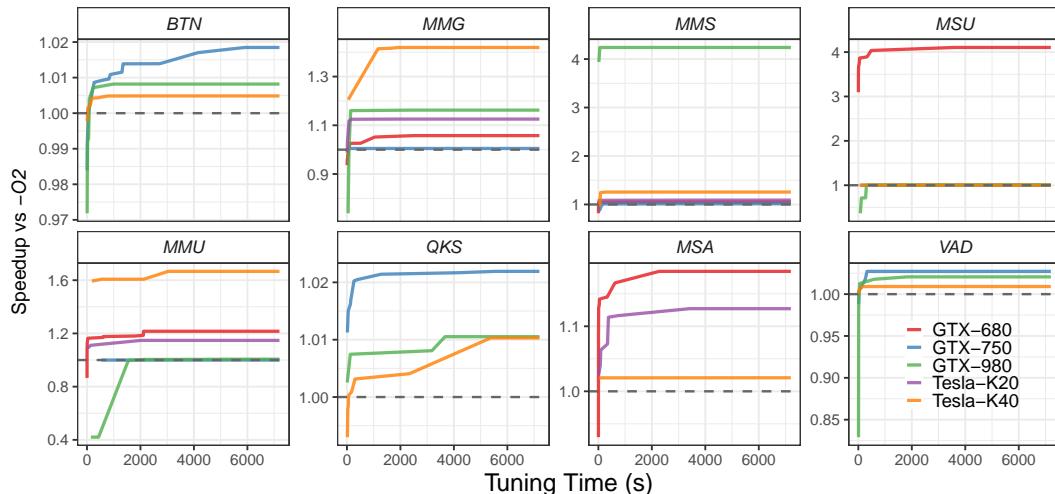


Figure 10.6: Mean speedup over 30 repetitions of the tuned solutions for the kernels we implemented versus the $opt-level = 2$ baseline, across two hours of tuning. Notice the difference in the y -axis scales for each panel

Flag	Cluster 0 (17%)	Cluster 1 (83%)
no-align-double	on	on
use_fast_math	on	on
preserve-relocs	off	off
relocatable-device-code	true	false
ftz	true	true
prec-div	true	false
prec-sqrt	true	true
fmad	false	false
allow-expensive-optimizations	false	true
gpu-architecture	sm_20	sm_50
def-load-cache	cv	ca
opt-level	1	3
maxrregcount	42	44.6

Table 10.4: Parameter clusters for all Rodinia problems in the GTX 750

10.4.3 Clustering Parameters found by Stochastic Experiments

We attempted to associate compilation parameters to kernels and GPUs using WEKA clustering algorithms [104]. Although we could not find significant relations for most kernels we detected that the *ftz=true* in MMS and the Compute Capabilities 3.0, 5.0 and 5.2 in GAU caused the speedups observed in the GTX 980 for these kernels. Table 10.4 shows clusters obtained with the K-means WEKA algorithm for autotuned parameter sets for the Rodinia Benchmark in the GTX 750. Unlike most clusters found for all GPUs and problems, these clusters did not contain an equal number of instances. The difficulty of finding associations between compiler optimizations, kernels and GPUs justifies the autotuning of compiler parameters in the general case.

10.5 Assessing Parameter Significance with Screening

This study was done in 2015, the first year of this thesis. Since at that time we did not have experience with the methods for statistical modeling which we later studied, and since the measurement time for each kernel was sufficiently small, we did not deal with the problem of building a performance model for the CUDA compiler in this study, and originally attempted to find good compiler configurations using stochastic methods for function minimization. Since then, we performed a new experiment using screening and targeting the same problem, in an attempt to identify the flags responsible for any improvements we observe. In this section, before closing the chapter, we will present and discuss this new screening experiment.

The data in this section has not yet been published at the time of the writing of this thesis. We designed and ran screening experiments for the CUDA compiler, adding the parameters described in Table 10.5 to the search space defined in Table 10.1. We no longer had access to the GPUs from the testbed described in the previous section, and we targeted the two GPUs described in Table 10.6.

Table 10.5: Flags added to the search space from Table 10.1

Flag	Description
force-store-cache	Force specified cache modifier on global/generic store. Step: PTX Values: cs, cg
force-load-cache	Replaces def-load-cache, force specified cache modifier on global/- generic load. Step: PTX Values: cs, cg
optimize	Specify optimization level for host code. Step: NVCC Values: 2, 3

For this experiment we picked *Heart Wall* (HWL), *Gaussian Elimination* (GAU), and *Needleman-Wunsch* (NDL), the three Rodinia [87] kernels for which we found the largest speedups in the last study. Table 10.3 gives more details on each kernel. We constructed a Plackett-Burman design for the 15 parameters in the search space, defining a minimum and maximum value for each parameter. For categorical parameters we fixed two arbitrary levels, for the optimization level parameters we picked the equivalents of -O2 and -O3 for both host and device code, and for the parameter that controls the number of available registers, we picked the largest number and the midpoint.

Table 10.6: Hardware specifications of the GPU architectures targeted in the screening experiments

Model	C.C.	Global Memory	Bus	Bandwidth	L2	Cores	Clock
Quadro M1200	5.0	4 GB	128-bit	80.19 GB/s	2 MB	640	1148 Mhz
Titan X	5.2	12 GB	384-bit	336.5 GB/s	3 MB	3072	1075 Mhz

The Plackett-Burman design for 15 2-level factors contains 20 experiments, and we performed 20 repetitions of each parameter assignment for each of the three kernels. The top and bottom panels of Figure 10.7 show for each target GPU the means and 95% confidence intervals for the coefficients of a linear model fit to the screening measurements, which represent the main effect estimates for each factor.

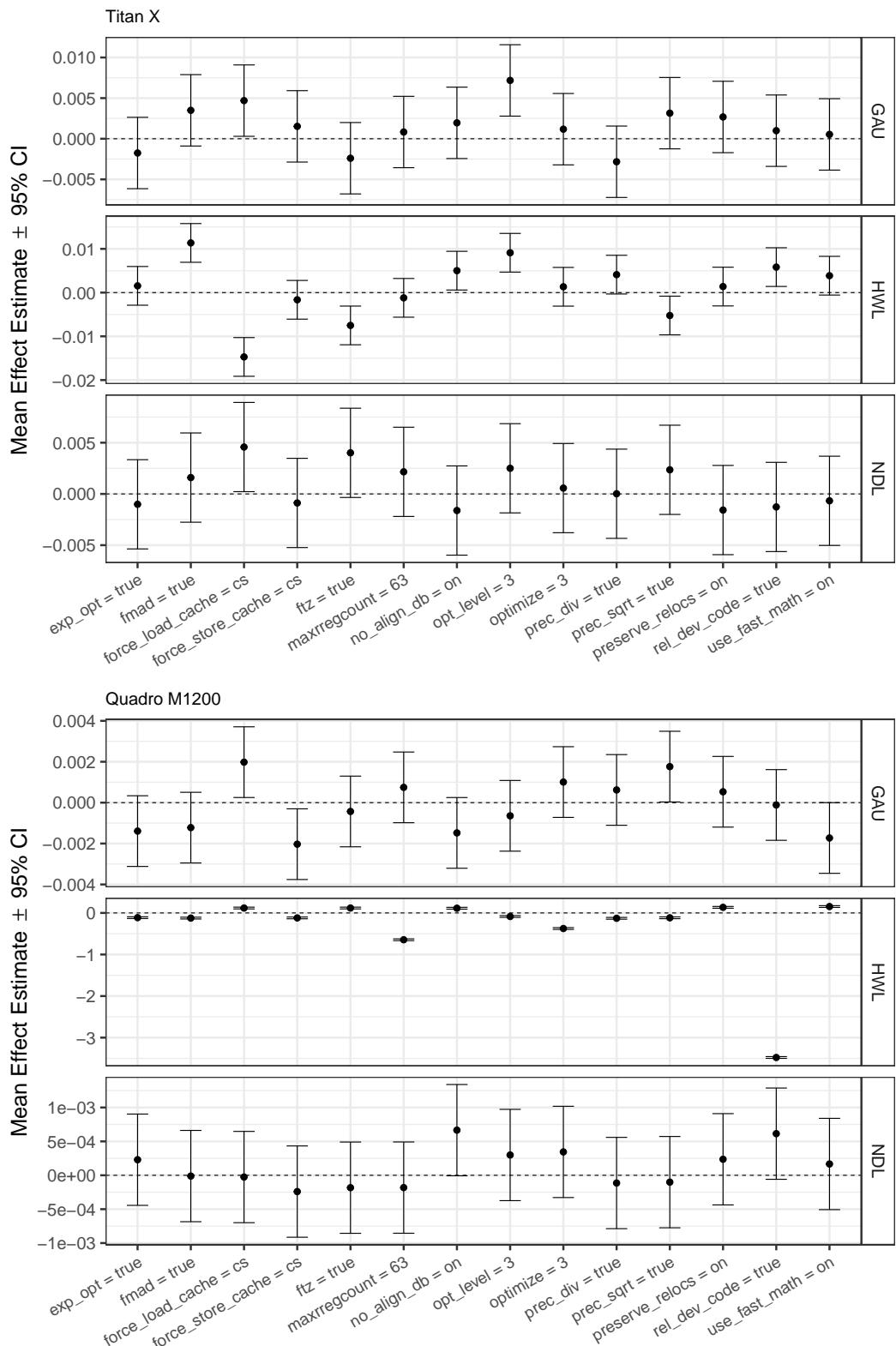


Figure 10.7: Main effects estimates and 95% confidence intervals for the *high level* of the 15 factors in the Plackett-Burman design, for target kernels measured on the Titan X and Quadro M1200 GPUs

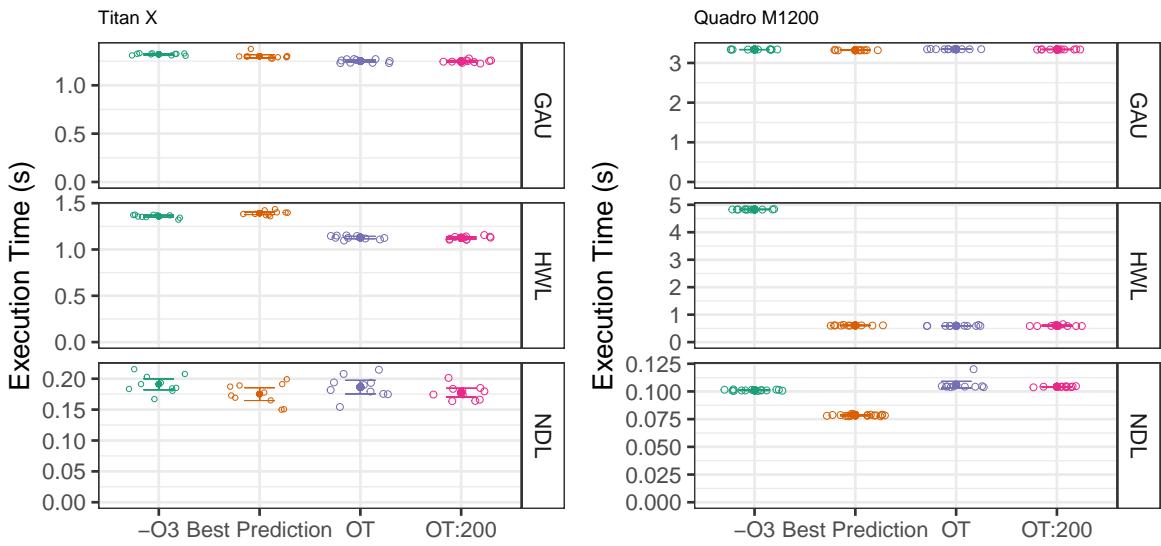


Figure 10.8: Execution time of baseline and model-predicted NVCC flag configurations for three Rodinia kernels on the Titan X and Quadro M1200 GPUs. Circles mark each of the 20 measurements of each flag configuration, filled dots and whiskers mark the mean and 95% confidence intervals

Note that because the Plackett-Burman design matrix is orthogonal the confidence intervals for the effect mean estimates all have the same size. In opposition to the exploration of stochastic methods, this well-balanced design allows the detection of factor main effects, which are mostly small, except in the HWL kernel on the Quadro GPU, where two factors present large effects. The linear models we built for each kernel using screening results do not account for interactions between factors, and thus their predictions favor turning on the factors whose mean estimates are negative and have confidence intervals that do not cross zero.

We generated a large set of points in the CUDA parameters search space and picked the points for which the models predict the smallest execution time. Figure 10.8 shows the performances of the points picked by the models and the baseline for comparison, for each Rodinia kernel. The baseline for comparison consisted of using only the PTX-stage *opt-level*=3 high-level compiler optimization. There was no significant difference in performance between second and third levels of *opt-level* in these experiments.

The left panel shows results for the Titan X GPU, where we found small speedups for the GAU and NDL kernels, but small slowdowns for HWL. The right panel shows results for the Quadro M1200 GPU, where we found statistically significant but small speedups for GAU and NDL, and 8 times speedup for the HWL kernel. The speedups obtained using screening are highlighted in Figure 10.9.

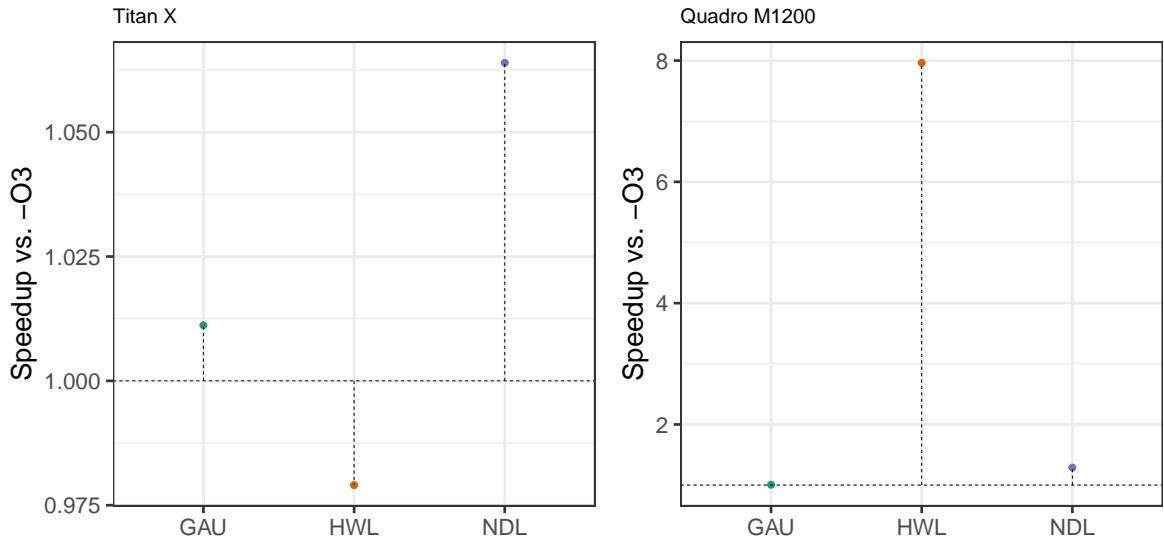


Figure 10.9: Mean speedups found using the predictions of the screening model, for three Rodinia kernels on the Titan X and Quadro M1200 GPUs. The horizontal dashed lines mark the PTX-stage `opt-level=3` comparison baseline

10.6 Summary

This study provides evidence that it is possible to improve the performance of GPU kernels with by selecting CUDA compiler parameters and flags. The inconclusive clustering attempts for the data obtained from heuristic and stochastic explorations of the search spaces emphasize the importance of well-designed experiment. Later experiments with screening reproduced the interesting speedups found for some kernels in different GPUs, while simultaneously allowing the identification of the flags responsible for performance improvements. Future work in this direction will explore the impact of adding host compiler optimizations and kernel-specific parameters. In this larger search space it would be more interesting to employ the more flexible Optimal Design methods we discuss in Chapters 12 and 13.

Chapter 11

High-Level Synthesis Kernels for FPGAs

A Field-Programmable Gate Array (FPGA) is a reprogrammable circuit that enables writing specialized hardware specifications for a variety of applications without changing the underlying chip. In recent years High Performance Computing tasks that require low power and latency have increasingly switched to FPGAs. Hardware specifications are written using low-level languages such as Verilog and VHDL, creating a challenge for software engineers to leverage FPGA capabilities. Developing strategies to decrease the effort required to program FPGAs is becoming more relevant with the overspread adaptation of FPGAs for data centers [105, 106, 107, 108, 109, 110], with direct vendor backing [111, 112].

Essential support for software engineers can be provided by High-Level Synthesis (HLS), where hardware descriptions are generated from high-level code. HLS compilers intend to lower the complexity of hardware design and have become increasingly valuable as part of the FPGA design workflow, with support from vendor HLS tools [113, 114] for C/C++ and OpenCL. The benefits of higher-level abstractions often come with the cost of decreased performance, making FPGAs less viable as accelerators. Thus, optimizing HLS still requires domain expertise and exhaustive or manual exploration of design spaces and configurations.

High-Level Synthesis is a challenging problem, and a common strategy for its solution involves the divide-and-conquer approach [115]. The most important sub-problems to solve are *scheduling*, where operations are assigned to specific clock cycles, and *binding*, where operations are assigned to specific hardware functional units, which can be shared between operations. LegUp [116] was an initially open-source HLS tool, that has since gone paid and closed-source [117], and was implemented as a compiler pass for the LLVM Compiler Infrastructure [118]. LegUp receives code in LLVM’s intermediate representation as input and produces as output a hardware description in Verilog, exposing configuration parameters of its HLS process, which are set with a configuration file.

This chapter presents our implementation of an autotuner for LegUp HLS parameters, also using the OpenTuner framework, which we published at ReConFig in 2017 [3]. The autotuner we implemented in this study targeted 8 hardware metrics obtained from Altera Quartus, for applications of the CHStone HLS benchmark suite [119] targeting the Intel StratixV FPGA. The program whose configurations are explored in this study is the LegUp HLS compiler, and the search space is composed of approximately 10^{126} possible combinations of HLS parameters.

One of the obstacles we faced was the impossibility of making accurate predictions for hardware metrics from a set of HLS parameters, or even for the generated Verilog, requiring that we run the autotuner using the metrics reported by the lengthier process of complete hardware synthesis instead. To mitigate the extra time cost, we implemented a virtualized deployment workflow which enabled launching several distributed tuning runs at the same time. We present data showing that the autotuner using stochastic methods from OpenTuner found optimized HLS parameters for CHStone applications that decreased the *Weighted Normalized Sum* (WNS) of hardware metrics by up to 21.5%, in relation to the default LegUp configuration.

It takes a considerable time to synthesize hardware from specifications in hardware description languages, ranging from minutes to hours. The process of generating those specifications from high-level C code is much faster in comparison, taking only seconds. Section 11.2 describes in more detail the High-Level Synthesis process.

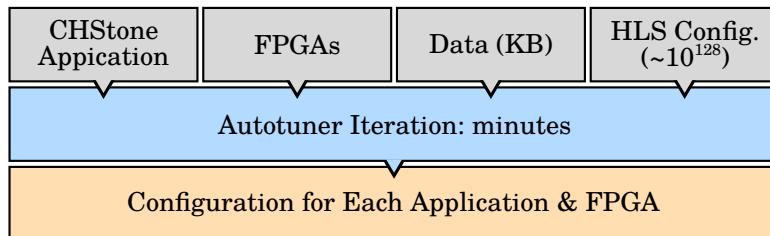


Figure 11.1: Autotuner representation and time scale of the experiments

Figure 11.1 shows a representation of our autotuner for this experiment, and highlights the time scales involved. The autotuner, represented by the light blue box, receives as input a CHStone kernel, a target FPGA, input data for the kernel in the order of Kilobytes and a search space composed of LegUp’s HLS parameters. The autotuner outputs an HLS configuration for each kernel and FPGA.

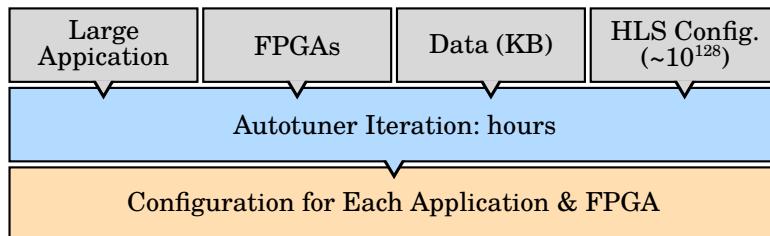


Figure 11.2: Autotuner representation and time scale of more complex FPGA applications

The kernels from CHStone presented further in Table 11.3 are very simple benchmark kernels, and generating hardware for them takes time in the order of minutes. Generating hardware for more complex FPGA applications can take several hours. Figure 11.2 shows a representation and time scale for an autotuner that targets complex FPGA applications. Each iteration now takes several hours, and limits the configurations that can be tested per hour. We worked around this limitation in this experiment by implementing a virtualized autotuner using Docker containers, which enabled parallel measurements of different configurations.

The rest of this chapter is organized as follows. Section 11.1 discusses the background on HLS tools and autotuning. Section 11.2 presents the search space defined by LegUp HLS parameters and the autotuner we implemented using OpenTuner. Section 11.3 introduces the optimization scenarios balancing different hardware metrics we targeted in this study, discusses the HLS kernels we optimized, and presents the settings in which the experiments were performed. Section 11.4 presents and discusses the results on each scenario. Finally, Section 11.5 summarizes the chapter and discusses future work.

11.1 Autotuning High-Level Synthesis for FPGAs

In this section we discuss background work related to HLS tools, and autotuning for FPGAs.

11.1.1 Tools for HLS

Various research and vendor tools for High-Level Synthesis have been developed [113, 114]. Villareal *et al.* [120] implemented extensions to the Riverside Optimizing Compiler for Configurable Circuits (ROCCC), which also uses the LLVM compiler infrastructure, to add support for generating VHDL from C code. Implemented within GCC, GAUT [121] is an open-source HLS tool for generating VHDL from C/C++ code. Other HLS tools such as Mitrion [122], Impulse [123] and Handel [124] also generate hardware descriptions from C code. We refer the reader to the survey from Nane *et al.* [125] for a comprehensive analysis of recent approaches to HLS.

11.1.2 Autotuning for FPGAs

Recent work studies autotuning approaches for FPGA compilation. Xu *et al.* [36] uses distributed OpenTuner instances to optimize the compilation flow from hardware description to bitstream. They optimize configuration parameters from the Verilog-to-Routing (VTR) toolflow [126] and target frequency, wall-clock time and logic utilization. Huang *et al.* [127] study the effect of LLVM pass ordering and application in LegUp’s HLS process, demonstrating the complexity of the search space and the difficulty of its exhaustive exploration. They exhaustively explore a subset of LLVM passes and target logic utilization, execution cycles,

frequency, and wall-clock time. Mamejanov {et al.} [34] propose a machine-learning-based approach to tune design parameters for performance and power consumption. Nabi and Vanderbauwhede [128] present a model for performance and resource utilization for designs based on an intermediate representation.

11.2 Autotuner and Search Space for the LegUp HLS Compiler

This section describes our autotuner implementation, the LegUp HLS parameters selected for tuning, and the autotuning metrics used to measure the quality of HLS configurations.

11.2.1 Autotuner

We implemented our autotuner with OpenTuner [1], using ensembles of search techniques to find an optimized selection of LegUp [129] HLS parameters, according to our cost function, for 11 of the CHStone [119] kernels.

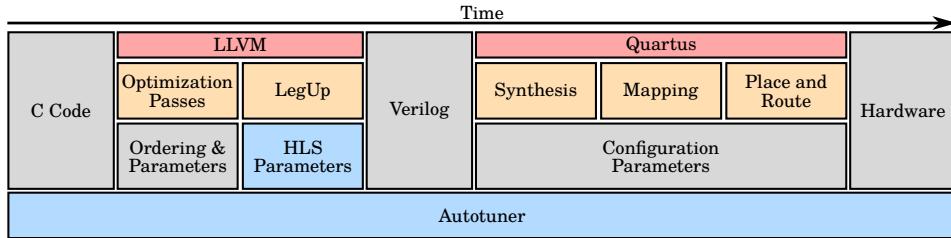


Figure 11.3: High-Level Synthesis compilation process. The autotuner search space at the HLS stage is highlighted in blue

Figure 11.3 shows the steps to generate a hardware description from C code. It also shows the Quartus steps to generate bitstreams from hardware descriptions and to obtain the hardware metrics we targeted. Our autotuner used LegUp’s HLS parameters as the search space, but it completed the hardware generation process to obtain metrics from Quartus, as represented by the blue boxes in Figure 11.3.

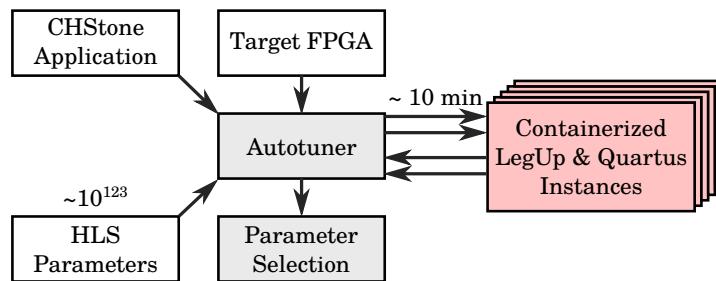


Figure 11.4: Autotuner Setup

Figure 11.4 shows our setup using Docker containers running LegUp and Quartus. This virtualization setup enabled portable dependency management and can be used to run exper-

iments in distributed environments. The arrows coming from the autotuner to the containers represent the flow of new configurations generated by search techniques, and the arrows coming from the containers to the autotuner represent the flow of measurements for a set of parameters. For CHStone kernels measurements take approximately 10 minutes to complete, and the majority of this time is spent in Quartus's synthesis, mapping, and place and route steps.

11.2.2 High-Level Synthesis Parameters

We selected an extensive set of LegUp High-Level Synthesis parameters, shown partially in Table 11.1. Each parameter in the first two rows of Table 11.1 has an 8, 16, 32 and 64 bit variant. *Operation Latency* parameters define the number of clock cycles required to complete a given operation when compiled with LegUp. *Zero-latency* operations can be performed in a single clock cycle. *Resource Constraint* parameters define the number of times a given operation can be performed in a clock cycle. *Boolean or Multi-Valued* parameters are used to set various advanced configurations. For example, the *enable_{patternsharing}* parameter can be set to enable resource sharing for patterns of computational operators, as is described by Hadjis *et al.* [130]. For a complete list and description of each parameter, please refer to LegUp's official documentation [131].

Type	Parameters
Operation Latency	<code>altpf_[divide, truncate, fptosi, add, subtract, multiply, extend, sitofp], unsigned_[multiply, divide, add, modulus], signed_[modulus, divide, multiply, add, comp_[o, u]], [local_mem, mem]_dual_port, reg</code>
Resource Constraint	<code>signed_[divide, multiply, modulus, add], altpf_[multiply, add, subtract, di- vide], unsigned_[modulus, multiply, add, divide], [shared_mem, mem]_dual_port</code>
Boolean or Multi-value	<code>pattern_share_[add, shift, sub, bitops], sdc_[multipump, no_chaining, prior- ity], pipeline_[resource_sharing, all], ps_[min_size, min_width, max_size, bit_diff_threshold], mb_[minimize_hw, max_back_passes], no_roms, multi- plier_no_chain, dont_chain_get_elem_ptr, clock_period, no_loop_pipeline, incremental_sdc, disable_reg_sharing, set_combine_basicblock, en- able_pattern_sharing, multipumping, dual_port_binding, modulo_scheduler, explicit_lpm_mults</code>

Table 11.1: Subset of All Autotuned LegUP HLS Parameters

11.2.3 HLS Autotuning Metrics

To obtain values for hardware metrics we needed to perform the synthesis, mapping and place and route steps. We used Quartus to do so, and selected 8 hardware metrics reported by Quartus to compose our cost or fitness function. From the fitter summary we obtained 6 metrics. *Logic Utilization (LUT)* measures the number of logic elements and is composed of Adaptive Look-Up Table (ALUTs), memory ALUTs, logic registers or dedicated logic registers. The *Registers (Regs.)*, *Virtual Pins (Pins)*, *Block Memory Bits (Blocks)*, *RAM Blocks (BRAM)* and *DSP Blocks (DSP)* metrics measure the usage of the resources indicated by their names.

From the timing analysis we obtained the *Cycles* and *FMax* metrics, used to compute the *Wall-Clock Time* metric. This metric composed the cost function, but *Cycles* and *FMax* were not individually used. We chose to do that because all of our hardware metrics needed to be minimized except for *FMax*, and computing *Wall-Clock Time* instead solved that restriction. The *Wall-Clock Time* wct is computed by $wct = Cycles \times (\alpha/FMax)$, where $\alpha = 10^6$ because *FMax* is reported in MHz.

The objective function used by the autotuner to evaluate sets of HLS parameters is written

$$f(\mathbf{m}, \mathbf{w}) = \frac{\sum_{m \in \mathbf{m}} w \left(\frac{m}{m^0} \right)}{\sum_{w \in \mathbf{w}} w}, \quad (11.1)$$

and computes a *Weighted Normalized Sum* (WNS) of the measured metrics $m \in \mathbf{m}$, where \mathbf{m} is a vector with measurements of the 8 hardware metrics described previously. Each weight $w \in \mathbf{w}$ corresponds to one of the scenarios in Table 11.2. A value is computed for each metric m in relation to an initial default value m^0 , measured for each metric in a given kernel, using LegUp's default configuration. For a given set measurement vector \mathbf{x} , $f(\mathbf{x}, \mathbf{w}) = 1.0$ means that there was no improvement relative to the default HLS configuration.

11.3 Target Optimization Scenarios and HLS Kernels

This section describes the optimization scenarios, the CHStone kernels, and the experimental settings.

11.3.1 Optimization Scenarios

Table 11.2 shows the assigned weights in our 4 optimization scenarios. The *Area*-targeting scenario assigns low weights to wall-clock time metrics. The *Performance and Latency* scenario assigns high weights to wall-clock time metrics and also to the number of registers used. The *Performance* scenario assigns low weights to area metrics and cycles, assigning a high weight only to frequency. The *Balanced* scenario assigns the same weight to every metric. The weights assigned to the metrics that do not appear on Table 11.2 are always 1. The weights are integers and powers of 2.

We compared results when starting from a *Default* configuration with the results when starting at a *Random* set of parameters. The default configuration for the StratixV was provided by LegUp and the comparison was performed in the *Balanced* optimization scenario.

Table 11.2: Weights for Optimization Scenarios
(*High* = 8, *Medium* = 4, *Low* = 2)

Metric	Area	Perf. and Lat	Performance	Balanced
<i>LUT</i>	High	Low	Low	Medium
<i>Registers</i>	High	High	Medium	Medium
<i>BRAMs</i>	High	Low	Low	Medium
<i>DSPs</i>	High	Low	Low	Medium
<i>FMax</i>	Low	High	High	Medium
<i>Cycles</i>	Low	High	Low	Medium

11.3.2 Kernels

To test and validate our autotuner we used 11 kernels from the CHStone HLS benchmark suite [119]. CHStone kernels are implemented in the C language and contain inputs and previously computed outputs, allowing for correctness checks to be performed for all kernels.

Table 11.3: Autotuned CHStone Kernels

Kernel	Short Description
blowfish	Symmetric-key block cypher
aes	Advanced Encryption Algorithm (AES)
adpcm	Adaptive Differential Pulse Code Modulation dec. and enc.
sha	Secure Hash Algorithm (SHA)
motion	Motion vector decoding from MPEG-2
mips	Simplified MIPS processor
gsm	Predictive coding analysis of systems for mobile comms.
dfsin	Sine function for double-precision floating-point numbers
dfmul	Double-precision floating-point multiplication
dfdiv	Double-precision floating-point division
dfadd	Double-precision floating-point addition

Table 11.3 provides short descriptions of the 11 CHStone kernels we used. We were not able to compile the *jpeg* CHStone kernel, so did not use it. All experiments targeted the *Intel StratixV 5SGXEA7N2F45C2* FPGA.

11.3.3 Experiments

We performed 10 tuning runs of 1.5h for each kernel. Section 11.4 presents the mean relative improvements for each kernel and individual metric. The code needed to run the experiments and generate the figures, as well as the implementation of the autotuner and all data we generated, is open and hosted at GitHub [132].

The experimental settings included Docker for virtualization and reproducibility, *LegUp v4.0*, *Quartus Prime Standard Edition v16.0*, and *CHStone* [119]. All experiments were performed on a machine with two *Intel Xeon CPU E5-2699 v3* with 18 *x86_64* cores each, and 503GB of

RAM. The instructions and the code to reproduce the software experimental environment are open and hosted at GitHub [133].

11.4 Performance Improvements using Stochastic Methods

This section presents summaries of the results from 10 autotuning runs of $1.5h$ in the scenarios from Table 11.2. Results are presented in *heatmaps* where each row has one of the 11 CHStone kernels in Table 11.3 and each column has one of the 8 hardware metrics and their *Weighted Normalized Sum* (WNS) as described in Section 11.2.3.

Cells on heatmaps show the ratio of tuned to initial values of a hardware metric in a CHStone kernel, averaged over 10 autotuning runs. The objective of the autotuner is to minimize all hardware metrics, except for *FMax*, whose inverse is minimized. Cell values less than 1.0 always mark an improvement on a given metric. Darker blue squares in the following heatmaps mark improvements, and darker red squares mark worse values in relation to the starting point.

	<i>aes</i>	0.79	0.91	1.00	0.56	1.00	0.47	1.00	1.12
<i>adperm</i>	0.68	1.13	1.00	0.54	1.00	0.56	0.60	0.98	
<i>sha</i>	1.00	1.03	1.00	0.82	1.00	0.55	1.00	0.89	
<i>motions</i>	1.02	1.00	0.60	0.85	1.00	0.57	1.00	0.94	
<i>mips</i>	1.00	0.93	1.00	0.44	1.00	0.45	0.98	1.24	
<i>gsm</i>	0.83	1.17	1.00	0.48	1.00	0.56	0.52	0.99	
<i>dfsin</i>	0.79	0.97	1.00	0.61	1.00	0.60	1.49	1.53	
<i>dfmul</i>	1.00	1.06	0.90	0.47	0.90	0.47	1.31	1.10	
<i>dfdiv</i>	0.83	1.07	0.80	0.73	0.80	0.65	1.32	1.49	
<i>dfadd</i>	1.00	0.94	1.00	0.82	1.00	0.71	1.00	0.93	
<i>blowfish</i>	—	—	—	—	—	—	—	—	
	<i>LUTs</i>	<i>Pins</i>	<i>BRAM</i>	<i>Regs</i>	<i>Blocks</i>	<i>Cycles</i>	<i>DSP</i>	<i>FMax</i>	

Figure 11.5: Comparison of the absolute values for Random and Default starting points in the Balanced scenario

Figure 11.5 compares the ratios of absolute values for each hardware metric for *Default* and *Random* starts, in the *Balanced* scenario. Cell values less than 1.0 mean that the *Default* start achieved smaller absolute values than the *Random* start. Cells with “—” mean that the *Default* start could not find a set of HLS parameters that produced valid output during any of the $1.5h$ tuning runs. The *Default* start found better values for most metrics.

The *Random* start found better values for *DSP*, *Pins* and *FMax* for some kernels. For example, it found values 49% smaller, 6% smaller and 53% larger for *DSP*, *Pins* and *FMax*, respectively, for the *dfsin* kernel. The *Default* start found better values for *Regs* and *Cycles* for all kernels. For example, it found values 53% smaller for *Regs* and *Cycles* for the *dfmul* kernel, and 56% and 55% smaller for *Regs* and *Cycles*, respectively, for the *mips* kernel.

The *Random* start found worst values in most cases because of the size of the search space and the stochastic nature of the explorations performed by OpenTuner methods. In such scenarios, knowing and leveraging a reasonably good starting configuration for a given

kernel is extremely important. The remaining results in this Section used one of the *Default* starting configurations provided by LegUp, specifically targeted to Stratix V boards.

Figure 11.6 shows the results for the *Balanced* scenario. These results are the baseline for evaluating the autotuner in other scenarios, since all metrics had the same weight. The optimization target was the *Weighted Normalized Sum* (WNS) of hardware metrics, but we were also interested in changes in other metrics, as their relative weights changed. In the *Balanced* scenario we expected to see smaller improvements of WNS due to the competition of concurrent improvements on every metric.

The autotuner found values of WNS 16% smaller for the *adpcm* and *dfdiv* kernels, and 15% smaller for *dfmul*. Even for the *Balanced* scenario it is possible to see that some metrics decreased while others decreased consistently over the 10 tuning runs. *FMax* and *DSP* had the larger improvements for most kernels, for example, 51% greater *FMax* in *adpcm* and 69% smaller *DSP* in *dfmul*. *Cycles*, *Regs* and *Pins* had the worst results in this scenario, with 34% larger *Cycles* in *dfdiv*, 15% larger *Regs* in *dfdiv* and 17% larger *Pins* in *gsm*. Other metrics had smaller improvements or no improvements at all in most kernels.

	aes	0.94	0.90	1.00	1.00	0.97	1.00	0.98	1.00	0.76
adpcm		0.84	0.73	1.13	1.00	0.91	1.00	1.05	0.63	0.49
sha		0.98	1.00	1.03	1.00	0.96	1.00	0.86	1.00	0.97
motion		0.98	0.97	1.00	1.00	0.99	1.00	0.95	1.00	0.95
mips		0.95	1.00	1.07	1.00	0.90	1.00	1.01	0.80	0.89
gsm		0.95	0.95	1.17	1.00	0.99	1.00	0.95	0.49	1.08
dfsim		0.93	1.03	1.03	1.00	1.05	1.00	1.07	0.65	0.73
dfmul		0.85	1.00	1.13	0.90	0.88	0.90	1.06	0.31	0.76
dfdiv		0.84	1.00	1.07	0.80	1.15	0.80	1.34	0.41	0.57
dfadd		1.00	1.00	1.00	1.00	0.98	1.00	0.98	1.00	0.97
blowfish		0.99	1.00	1.00	1.00	0.98	1.00	0.98	1.00	0.99
WNS										
LUTs										
Pins										
BRAM										
Regs										
Blocks										
Cycles										
DSP										
FMax										

Figure 11.6: Relative improvement for all metrics in the *Balanced* scenario

	aes	0.96	0.92	1.00	1.00	0.93	1.00	0.97	1.00	0.86
adpcm		0.83	0.78	1.11	1.00	0.96	1.00	1.04	0.63	0.52
sha		0.96	1.00	1.06	1.00	0.84	1.00	0.83	1.00	1.08
motion		0.97	0.94	1.00	1.00	0.97	1.00	0.92	1.00	0.95
mips		0.91	1.00	1.06	1.00	0.89	1.00	1.00	0.72	0.84
gsm		0.89	0.83	1.06	1.00	0.86	1.00	0.89	0.82	1.14
dfsim		0.94	1.00	1.06	1.00	1.00	1.00	1.02	0.76	0.80
dfmul		0.82	1.00	1.06	1.00	0.90	1.00	1.21	0.27	0.86
dfdiv		0.77	1.00	1.22	0.67	1.03	0.67	1.62	0.28	0.77
dfadd		0.99	1.00	1.11	1.00	0.94	1.00	1.00	1.00	1.04
blowfish		0.99	1.00	1.00	1.00	0.97	1.00	0.91	1.00	1.01
WNS										
LUTs										
Pins										
BRAM										
Regs										
Blocks										
Cycles										
DSP										
FMax										

Figure 11.7: Relative improvement for all metrics in the *Area* scenario

Figure 11.7 shows the results for the *Area* scenario. We believe that the greater coherence of optimization objectives is responsible for the greater improvements of WNS in the following scenarios. The autotuner found values of WNS 23% smaller for *dfdiv*, 18% smaller for *dfmul*, and smaller values overall in comparison with the *Balanced* scenario. Regarding individual metrics, the values for *FMax* were worse overall, with 14% smaller *FMax* in *gsm* and 62% greater *Cycles*, for example. As expected for this scenario, metrics related to area had better improvements than in the *Balanced* scenario, with 73% and 72% smaller *DSP* for *dfmul* and *dfdiv* respectively, 33% smaller *Blocks* and *BRAM* in *dfdiv* and smaller values overall for *Regs* and *LUTs*.

Figure 11.8 shows the results for the *Performance* scenario. The autotuner found values of WNS 23% smaller for *dfmul*, 19% smaller for *dfdiv*, and smaller values overall than in the *Balanced* scenario. *FMax* was the only metric with a *High* weight in this scenario, so most metrics had improvements close overall to the *Balanced* scenario. The values for *FMax* were

best overall, with better improvements in most kernels. For example, 41%, 30%, 44% and 37% greater *FMax* in *dfdiv*, *dfmul*, *dfsinc* and *aes* respectively.

	<i>aes</i>	0.82	0.75	1.00	1.00	0.92	1.00	1.05	1.00	0.63
<i>adpcm</i>		0.84	0.94	1.17	1.00	0.96	1.00	0.94	0.69	0.74
<i>sha</i>		0.92	1.00	1.06	1.00	0.92	1.00	0.88	1.00	0.96
<i>motion</i>		0.99	1.00	1.00	1.01	1.00	1.00	1.00	1.00	0.98
<i>mips</i>		0.90	1.00	1.06	1.00	0.93	1.00	1.03	0.83	0.80
<i>gsm</i>		0.95	0.92	1.00	1.00	0.95	1.00	0.90	1.00	1.02
<i>dfsinc</i>		0.87	1.11	1.06	1.00	1.27	1.00	1.25	0.53	0.56
<i>dfmul</i>		0.77	1.00	1.17	0.83	0.85	0.83	1.04	0.21	0.70
<i>dfdiv</i>		0.81	1.00	1.06	1.00	1.03	1.00	1.25	0.50	0.59
<i>dfadd</i>		0.97	1.00	1.00	1.00	0.97	1.00	1.00	1.00	0.94
<i>blowfish</i>		0.94	1.00	1.00	1.00	0.98	1.00	0.91	1.00	0.95
	WNS	LUTs	Pins	BRAM	Regs	Blocks	Cycles	DSP	FMax	

Figure 11.8: Relative improvement for all metrics in the *Performance* scenario

Figure 11.9 shows the results for the *Performance and Latency* scenario. The autotuner found values of WNS 24% smaller for *adpcm*, 18% smaller for *dfdiv*, and smaller values overall than in the *Balanced* scenario. *Regs*, *Cycles* and *FMax* had higher weights in this scenario, and also better improvements overall. For example, 16% and 15% smaller *Regs* in *dfdiv* and *sha* respectively, 23% and 11% smaller *Cycles* in *sha* and *aes* respectively, and 53% greater *FMax* in *adpcm*. Although *FMax* had the worst improvements in relation to the *Balanced* scenario, the *Wall-Clock Time* was still decreased by the smaller values of *Cycles*.

Figure 11.10 summarizes the average improvements on WNS in the 4 scenarios over 10 runs. Only the *Weighted Normalized Sum* of metrics directly guided optimization. With the exception of *dfadd* in the *Balanced* scenario, the autotuner decreased WNS for all kernels in all scenarios by 10% on average, and up to 24% for *adpcm* in the *Performance and Latency* scenario. The figure also shows the average decreases for each scenario.

	Balanced	Area	Performance	Perf. & Lat.
<i>aes</i>	0.94	0.96	0.82	0.83
<i>adpcm</i>	0.84	0.83	0.84	0.76
<i>sha</i>	0.98	0.96	0.92	0.88
<i>motion</i>	0.98	0.97	0.99	0.98
<i>mips</i>	0.95	0.91	0.90	0.95
<i>gsm</i>	0.95	0.89	0.95	0.92
<i>dfsinc</i>	0.93	0.94	0.87	0.97
<i>dfmul</i>	0.85	0.82	0.77	0.86
<i>dfdiv</i>	0.84	0.77	0.81	0.82
<i>dfadd</i>	1.00	0.99	0.97	0.98
<i>blowfish</i>	0.99	0.99	0.94	0.96
Average	0.93	0.91	0.89	0.90

Figure 11.10: Relative improvement for WNS in all scenarios

11.5 Summary

This study highlights the importance of starting positions for an autotuner based on stochastic methods. This becomes more relevant as the size of the search space and the number

of targeted metrics increase. The flexibility of our virtualized approach is evidenced by the results for different optimization scenarios. Improvements in WNS increased when higher weights were assigned to metrics that express a coherent objective such as area, performance and latency. The improvements of metrics related to those objectives also increased.

Kernels with large measurement time still presented a challenge for the autotuning approach we used in this problem, even with distributed measurements, because the stochastic methods implemented in OpenTuner reach better results if longer explorations of the search space are performed. A Design of Experiments approach such as the one we later developed and applied to other problems, as described in Chapter 12, would be ideal for this type of problem. Since the HLS toolchain used in this study became proprietary software and FPGA programming is still a very closed domain, requiring paid licenses and still expensive hardware, we did not have the opportunity to revisit this study with the Design of Experiments methods we later studied.

Future work in this direction will study the impact of different starting points on the final tuned values in each optimization scenario, for example we could start tuning for *Performance* at the best autotuned *Area* value. We expected that starting positions tailored for each target kernel will enable the autotuner to find better WNS values faster. We will also apply this autotuning methodology to HLS tools that enable a fast prediction of metric values. These tools will enable the exploration of the trade-off between prediction accuracy and the time to measure an HLS configuration.

Chapter 12

An OpenCL Laplacian Kernel

This chapter introduces our Experimental Design approach to autotuning, and presents an application where we continued the exploration, started by Masnada [134], of performance optimization methods for the search space defined by the parameters of a Laplacian GPU kernel. The kernel was implemented using BOAST [26], a framework that enables writing and optimizing HPC applications using metaprogramming, and is publicly hosted on GitHub [135]. We targeted the *Nvidia K40c* GPU, and the objective function we minimized was the time to compute each pixel. We evaluated the performance of nine optimization methods for this kernel, including sampling strategies, stochastic methods, variations on linear regression, our Design of Experiments approach, and Gaussian Process Regression. Our transparent and parsimonious Experimental Design approach achieved the most consistent results across the methods we tested, which motivated searching for a more comprehensive set of applications in which we could evaluate the performance of our approach. This more comprehensive study is described in Chapter 13. The initial steps of the studies presented in this chapter and the next were published together in the CCGRID conference [4].

The remainder of this chapter is organized as follows. Section 12.1 presents the target kernel and the associated search space exposed by the BOAST implementation. Section 12.2 presents our adaptation of a sequential approach to autotuning, starting with modeling assumptions about the search space and iteratively refining models based on ANOVA tests. Section 12.3 discusses our initial modeling hypotheses. Section 12.4 looks at a single run of our approach, linking the factor levels that were chosen at each step to ANOVA p -values, and comparing to the levels of the global optimum. Section 12.5 evaluates the performance of each method. Finally, Section 12.6 the discussion and concludes the chapter.

Figure 12.1: Edge-detection effect of a *Laplacian of Gaussian* filter

12.1 The Laplacian Kernel

The Laplacian Δf of a function f with inputs x and y is written

$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}, \quad (12.1)$$

and can be interpreted as providing a measure of how much the values of f in a neighborhood of a point (x, y) deviate, on average, from the value of $f(x, y)$. The discrete Laplacian operator is commonly used in image processing algorithms such as edge detection, where it is typically composed with smoothing filters. In the discrete setting the Laplacian of an image can be computed by a convolution filter, where a single kernel encodes the smoothing and Laplacian filters. Figure 12.1 shows the effect of applying a *Laplacian of Gaussian* convolution filter to detect edges in a picture.

12.1.1 BOAST Code and the OpenCL Kernel

Our Experimental Design approach can be applied to any autotuning domain that expresses optimization as a search problem, but the performance evaluations we present in this thesis were obtained in the domain of source-to-source transformation. Several frameworks, compilers, and autotuners provide tools to generate and optimize architecture-specific code [27, 26, 37, 136, 28]. We used BOAST [26] to generate code for GPUs by generating an OpenCL Laplacian kernel, optimizing parameters controlling vectorization, loop transformations, and data structure size and copying.

Figure 12.1.1 shows a Laplacian kernel naively written in C. In this code, the RGB channels of an image are stored a single *int* array, representing a color image. A three by three convolution kernel is applied to each color channel separately.

```

void kernel(int width, int height, uint8_t *source, uint8_t *destination){
    int i, j, d, ii, jj, tmp;
    for (j = 1; j < height - 1; j++) {
        for (i = 1; i < width - 1; i++) {
            // Process RGB components separately
            for (d = 0; d < 3; d++) {
                tmp = 0;
                for(jj = -1; jj < 1; jj++) {
                    for(ii = -1; ii < 1; ii++) {
                        if(ii == 0 && jj == 0) {
                            tmp += 9 * source[d + (3 * (i + (width * j))]];
                        } else {
                            tmp -= source[d + (3 * ((i + ii) + (width * (j + jj))))];
                        }
                    }
                }
                // Clamp to valid pixel values
                destination[d + (3 * (i + (width * j)))] = tmp < 0 ? 0 :
                    (tmp > 255 ? 255 : tmp);
            }
        }
    }
}

```

Figure 12.2: A CPU Laplacian kernel written in C

For comparison, the BOAST Ruby script for generating a parameterized Laplacian kernel is partially shown in Figure 12.1.1. A sample of the resulting OpenCL kernel is shown in Figure 12.1.1. The difference in complexity between the naive implementation in Figure 12.1.1 and the parameterized generated version helps motivate the usage of automated optimization methods for autotuning.

Chapter 12. An OpenCL Laplacian Kernel

```

def laplacian(options)
  default_options = {:x_component_number => 1,
                     :vector_length => 1,
                     :y_component_number => 1,
                     :temporary_size => 2,
                     :vector_recompute => false,
                     :load_overlap => false}
  # [...] (Omitted)
  # Beginning of the generated procedure, based on the baseline
  p = Procedure("kernel", [psrc, pdst, width, height]) {
    # [...] (Omitted)
    # Using vector_length to compute the number of vectors
    vector_number = (x_component_number.to_f/vector_length).ceil
    total_x_size = vector_recompute ? vector_number *
                               vector_length : x_component_number
    # [...] (Omitted)
    # Using y_component_number to compute offsets
    y_offset = y_component_number + 1
    # Controlling load_overlap
    if not load_overlap then
      total_load_window = total_x_size + 6
      tempload = []
      ranges = split_in_ranges(total_load_window, vector_length)
      ranges.each { |r|
        tempload.push( Int("tempload#{r.begin}#{r.end}", :size => 1,
                           :vector_length => (r.end - r.begin + 1), :signed => false) )
      }
      decl *(tempload)
    else
      tempnn = (0..2).collect { |v_i|
        (0...vector_number).collect { |x_i|
          (0...(y_component_number+2)).collect { |y_i|
            Int("temp#{x_i}#{v_i}#{y_i}", :size => 1,
                :vector_length => vector_length, :signed => false) }
        }
      }
      decl *(tempnn.flatten)
    end
    # [...] (Omitted)
    # Using temporary_size to allocate data
    tempcnn = (0..2).collect { |i|
      (0...vector_number).collect { |x_i|
        (0...(y_component_number+2)).collect { |y_i|
          Int("tempcnn#{x_i}#{v_i}#{y_i}", :size => temporary_size,
              :vector_length => vector_length)
        }
      }
    }
    decl *(tempcnn.flatten)
    # [...] (Omitted)
    # Generate kernel and clamp to valid pixel values
    (0...vector_number).each { |v_i|
      (0...y_component_number).each { |y_i|
        pr rescnn[v_i][y_i] === - tempcnn[0][v_i][y_i] -
          tempcnn[1][v_i][y_i] -
          tempcnn[2][v_i][y_i] -
          tempcnn[0][v_i][y_i + 1] +
          tempcnn[1][v_i][y_i + 1] *
          "#{temp_type}9" -
          tempcnn[2][v_i][y_i + 1] -
          tempcnn[0][v_i][y_i + 2] -
          tempcnn[1][v_i][y_i + 2] -
          tempcnn[2][v_i][y_i + 2]
        pr resnn[v_i][y_i] === clamp(rescnn[v_i][y_i],
                                      "#{temp_type}0",
                                      "#{temp_type}255",
                                      :returns => rescnn[v_i][y_i])
      }
    }
    # [...] (Omitted)
  }
  # Variable p contains the complete generated OpenCL code
  pr p
  k.procedure = p
  return k
end

```

Figure 12.3: Excerpts of the BOAST code, in Ruby, that generates the Laplacian OpenCL kernel. The code is publicly hosted at GitHub [137]

```

__kernel void kernel(const __global uchar * psrc, __global uchar * pdst,
                    const int width, const int height){
    int y, x, w;
    x = (get_global_id(0)) * (1);
    y = (get_global_id(1)) * (3);
    w = (width) * (3);
    x = (x < 3 ? 3 : (x > w - (11) ? w - (11) : x));
    y = (y < 1 ? 1 : (y > height - (4) ? height - (4) : y));
    uchar8 temppload0_7, temppload8_11, temppload12_13,
            res00, res01, res02, tempc000;
    int8 tempc001, tempc002, tempc003, tempc004, tempc010, tempc011, tempc012,
        tempc013, tempc014, tempc020, tempc021, tempc022, tempc023, tempc024,
        resc00, resc01, resc02;
    temppload0_7 = vload8(0, &psrc[x + -3 + (w) * (y + -1)]);
    temppload8_11 = vload4(0, &psrc[x + 5 + (w) * (y + -1)]);
    temppload12_13 = vload2(0, &psrc[x + 9 + (w) * (y + -1)]);
    tempc000 = convert_int8( (uchar8)(temppload0_7.s01234567) );
    tempc010 = convert_int8( (uchar8)(temppload0_7.s345,temppload0_7.s67,temppload8_11.s012) );
    tempc020 = convert_int8( (uchar8)(temppload0_7.s67,temppload8_11.s0123,
                                         temppload12_13.s01) );
    temppload0_7 = vload8(0, &psrc[x + -3 + (w) * (y + 0)]);
    temppload8_11 = vload4(0, &psrc[x + 5 + (w) * (y + 0)]);
    temppload12_13 = vload2(0, &psrc[x + 9 + (w) * (y + 0)]);
    tempc001 = convert_int8( (uchar8)(temppload0_7.s01234567) );
    tempc011 = convert_int8( (uchar8)(temppload0_7.s345,temppload0_7.s67,temppload8_11.s012) );
    tempc021 = convert_int8( (uchar8)(temppload0_7.s67,temppload8_11.s0123,
                                         temppload12_13.s01) );
    temppload0_7 = vload8(0, &psrc[x + -3 + (w) * (y + 1)]);
    temppload8_11 = vload4(0, &psrc[x + 5 + (w) * (y + 1)]);
    temppload12_13 = vload2(0, &psrc[x + 9 + (w) * (y + 1)]);
    tempc002 = convert_int8( (uchar8)(temppload0_7.s01234567) );
    tempc012 = convert_int8( (uchar8)(temppload0_7.s345,temppload0_7.s67,temppload8_11.s012) );
    tempc022 = convert_int8( (uchar8)(temppload0_7.s67,temppload8_11.s0123,
                                         temppload12_13.s01) );
    temppload0_7 = vload8(0, &psrc[x + -3 + (w) * (y + 2)]);
    temppload8_11 = vload4(0, &psrc[x + 5 + (w) * (y + 2)]);
    temppload12_13 = vload2(0, &psrc[x + 9 + (w) * (y + 2)]);
    tempc003 = convert_int8( (uchar8)(temppload0_7.s01234567) );
    tempc013 = convert_int8( (uchar8)(temppload0_7.s345,temppload0_7.s67,
                                         temppload8_11.s012) );
    tempc023 = convert_int8( (uchar8)(temppload0_7.s67,temppload8_11.s0123,
                                         temppload12_13.s01) );
    temppload0_7 = vload8(0, &psrc[x + -3 + (w) * (y + 3)]);
    temppload8_11 = vload4(0, &psrc[x + 5 + (w) * (y + 3)]);
    temppload12_13 = vload2(0, &psrc[x + 9 + (w) * (y + 3)]);
    tempc004 = convert_int8( (uchar8)(temppload0_7.s01234567) );
    tempc014 = convert_int8( (uchar8)(temppload0_7.s345,temppload0_7.s67,temppload8_11.s012) );
    tempc024 = convert_int8( (uchar8)(temppload0_7.s67,temppload8_11.s0123,
                                         temppload12_13.s01) );
    // Kernel computation
    resc00 = -(tempc000) - (tempc010) - (tempc020) - (tempc001) + (tempc011) *
        ((int)9) - (tempc021) - (tempc002) - (tempc012) - (tempc022);
    res00 = convert_uchar8( clamp(resc00, (int)0, (int)255) );
    resc01 = -(tempc001) - (tempc011) - (tempc021) - (tempc002) + (tempc012) *
        ((int)9) - (tempc022) - (tempc003) - (tempc013) - (tempc023);
    res01 = convert_uchar8( clamp(resc01, (int)0, (int)255) );
    resc02 = -(tempc002) - (tempc012) - (tempc022) - (tempc003) + (tempc013) *
        ((int)9) - (tempc023) - (tempc004) - (tempc014) - (tempc024);
    res02 = convert_uchar8( clamp(resc02, (int)0, (int)255) );
    vstore8( res00, 0, &pdst[x + 0 + (w) * (y + 0)] );
    vstore8( res01, 0, &pdst[x + 0 + (w) * (y + 1)] );
    vstore8( res02, 0, &pdst[x + 0 + (w) * (y + 2)] );
}

```

Figure 12.4: A sample OpenCL Laplacian kernel generated by BOAST

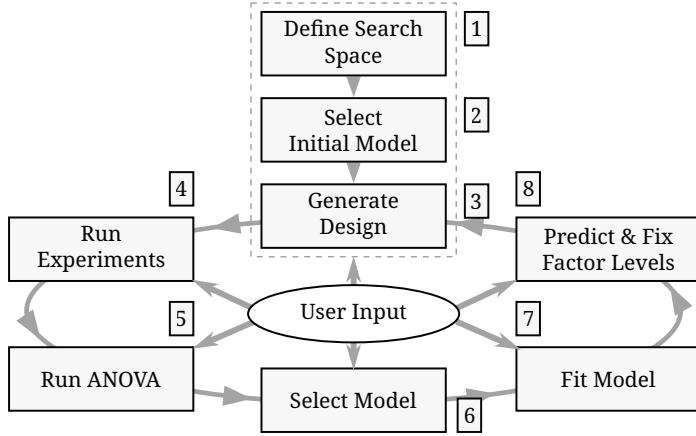


Figure 12.5: Overview of the DOE approach to autotuning we implemented

12.2 A Transparent and Parsimonious Approach to Autotuning using Optimal Design

In this section we discuss in detail our iterative, transparent, and parsimonious Experimental Design approach to autotuning. Figure 12.5 presents an overview. In step 1 we define the factors and levels that compose the search space of the target problem, in step 2 we select an initial performance model, and in step 3 we generate an experimental design. We run the experiments in step 4 and then, as we discuss in the next section, we identify significant factors with an ANOVA test in step 5. This enables selecting and fitting a new performance model in steps 6 and 7. The new model is used in step 8 for predicting levels for each significant factor. We then go back to step 3, generating a new design for the new problem subspace with the remaining factors. Informed decisions made by the user at each step guide the outcome of each iteration.

Step 1 of our approach is to define target factors and which of their levels are worth exploring. Then, the user must select an initial performance model in step 2. Compilers typically expose many 2-level factors in the form of configuration flags, and the performance model for a single flag can only be a linear term, since there are only 2 values to measure. Interactions between flags and numerical factors such as block sizes in CUDA programs or loop unrolling amounts are also common. Deciding which levels to include for these kinds of factors requires more careful analysis. For example, if we suspect the performance model has a quadratic term for a certain factor, the design should include at least three factor levels. The ordering between the levels of other compiler parameters, such as $-O(0,1,2,3)$, is not obviously translated to a number. Factors like these are named *categorical*, and must be treated differently when constructing designs in step 3 and analyzing results in step 5.

After the design is constructed in step 3, we run each selected experiment in step 4. This step can run in parallel since experiments are independent. Not all target programs run successfully in their entire input range, making runtime failures common in this step. The user

can decide whether to construct a new design using the successfully completed experiments or to continue to the analysis step if enough experiments succeed.

After running the ANOVA test in step 5, the user should apply domain knowledge to analyze the ANOVA table and determine which factors are significant. Certain factors might not appear significant and should not be included in the regression model. Selecting the model after the ANOVA test in step 6 also benefits from domain knowledge.

A central assumption of ANOVA is the *homoscedasticity* of the response, which can be interpreted as requiring the observed error on measurements to be independent of factor levels and of the number of measurements. Fortunately, there are statistical tests and corrections for lack of homoscedasticity. Our approach uses the homoscedasticity check and correction by power transformations from the `car` package [138] of the R language.

We fit the selected model to our design's data in step 7, and use the fitted model in step 8 to find levels that minimize the response. The choice of the method used to find these levels depends on factor types and on the complexity of the model and search space. If factors have discrete levels, neighborhood exploration might be needed to find levels that minimize the response around predicted levels. Constraints might put predicted levels on an undefined or invalid region on the search space. This presents challenge, because the borders of valid regions would have to be explored.

In step 8 we also fix factor levels to those predicted to achieve best performance. The user can also decide the level of trust placed on the prediction at this step, by keeping other levels available. In step 8 we perform a reduction of problem dimension by eliminating factors and decreasing the size of the search space. If we identified significant parameters correctly, we will have restricted further search to better regions.

In this chapter and in Chapter 13 we evaluate the performance of this Design of Experiments approach to autotuning in different applications. The next section discusses the implementation of a Laplacian kernel for GPUs and the resulting search space.

12.3 Building a Performance Model

This section describes the search space exposed by the BOAST Laplacian kernel and our initial modeling hypotheses.

12.3.1 Search Space

The factors and levels defining the search space are listed in Table 12.1. The complete search space contains 1.9×10^5 configurations, but removing invalid configurations, or configurations that fail at runtime, yields a search space with 2.3×10^4 configurations. The valid search space

Table 12.1: Parameters of the Laplacian Kernel

Factor	Levels	Short Description
<i>vector_length</i>	$2^0, \dots, 2^4$	Size of vectors
<i>load_overlap</i>	<i>true, false</i>	Load overlaps in vectorization
<i>temporary_size</i>	2, 4	Byte size of temporary data
<i>elements_number</i>	$1, \dots, 24$	Size of equal data splits
<i>y_component_number</i>	$1, \dots, 6$	Loop tile size
<i>threads_number</i>	$2^5, \dots, 2^{10}$	Size of thread groups
<i>lws_y</i>	$2^0, \dots, 2^{10}$	Block size in <i>y</i> dimension

took 154 hours to be completely evaluated on an *Intel Xeon E5-2630v2* CPU, with *gcc* version 4.8.3 and *Nvidia* driver version 340.32.

12.3.2 Modeling the Impact of Each Factor

This section briefly describes each factor and our modeling hypotheses regarding their impact on performance.

Linear Terms

The effects of the following factors were modeled with linear terms, whether because they are categorical binary factors or to attempt to exploit the simplest relationship under uncertainty of the effects.

- The *vector_length* factor controls the vector size used during computation. Architectures supporting vectorization provide speedups by saving instruction decoding time, but vectors must be properly sized. Although NVIDIA GPUs did not support vectorization at the time of this study, vectorization can still impact performance via cache effects
- The *load_overlap* factor is binary, and its effect is consequently modeled with a linear term. The parameter encodes the choice of whether to save memory load instructions by overlapping vectors in memory

Linear plus Inverse Terms

The effects of the following factors were modeled with a linear term plus an inverse term. In our initial model we assumed that these factors had a linear effect on performance, but also that a more complex relationship existed due to different kinds of overhead. In all cases, we modeled the expected overhead as an inverse term.

- The *elements_number* factor controls the size of image portion that will be processed by each thread, and using smaller portions implies using more threads. As more threads

are used we expect an improvement on performance, but also an overhead due to extra memory loads and to the cost of managing threads

- The lws_y factor controls the size of a thread block, in the y -axis. We did not have assumptions about the behavior of this factor, other than a possible overhead due to thread management, so we attempted linear plus inverse terms
- The $y_component_number$ factor controls tiling on the y -axis of the target image, inside each thread. The optimal tiling size is a compromise between fitting a tile in cache and providing enough prefetched memory to not slow computation down
- The $threads_number$ parameter controls the size of an OpenCL thread work group. Threads in the same group share data and can be scheduled together. Using more smaller groups can improve performance from better scheduling, but imply in more management overhead

The Complete Initial Model

Putting together the terms for all factors, the complete initial model for the time to compute a single pixel is written

$$\begin{aligned}
 time_per_pixel \sim & y_component_number + \frac{1}{y_component_number} + \\
 & temporary_size + vector_length + load_overlap + \\
 & lws_y + \frac{1}{lws_y} + elements_number + \frac{1}{elements_number} + \\
 & threads_number + \frac{1}{threads_number},
 \end{aligned} \tag{12.2}$$

where coefficients were omitted. This model was pruned at each iteration of our method, fixing significant factors to their best predicted level.

12.4 Looking at a Single DLMT Run

This section presents a more detailed look at the iterative process performed during the optimization of the Laplacian kernel by our method. Table 12.2 shows the ANOVA tests performed at each step in Figure 12.6, highlighting the factors identified to be significant at each step, which were subsequently fixed to their best predicted levels. Note that the size of the model being tested decreases at each step, as factors are fixed. These progressive restrictions to slices of the search space enable the ANOVA test to detect differences in means that were previously unclear.

Factors were chosen based on ANOVA tests where factors with p -values below a threshold of $p < 0.05$ were considered significant. The threshold for significance must be adjusted

Table 12.2: ANOVA tests at each step. Red lines mark model terms that were considered significant, with $p < 0.05$, and fixed in a given step

Step	Term	Sum Sq.	F-value	p(>F)
1 st	<i>y_component_number</i>	2.1×10^{-18}	7.3×10^{-1}	4.1×10^{-1}
	<i>1/y_component_number</i>	4.4×10^{-18}	1.6×10^0	2.4×10^{-1}
	<i>vector_length</i>	1.3×10^{-17}	4.4×10^0	4.7×10^{-2}
	<i>lws_y</i>	6.9×10^{-17}	2.4×10^1	3.5×10^{-4}
	<i>1/lws_y</i>	1.8×10^{-17}	6.2×10^0	2.8×10^{-2}
	<i>load_overlap</i>	9.1×10^{-20}	3.2×10^{-2}	8.6×10^{-1}
	<i>temporary_size</i>	7.1×10^{-18}	2.5×10^0	1.4×10^{-1}
	<i>elements_number</i>	3.1×10^{-19}	1.1×10^{-1}	7.5×10^{-1}
	<i>1/elements_number</i>	1.3×10^{-18}	4.4×10^{-1}	5.2×10^{-1}
	<i>threads_number</i>	7.2×10^{-18}	2.5×10^0	1.4×10^{-1}
	<i>1/threads_number</i>	4.3×10^{-18}	1.5×10^0	2.4×10^{-1}
2 nd	<i>y_component_number</i>	1.2×10^{-19}	2.1×10^1	1.4×10^{-3}
	<i>1/y_component_number</i>	1.4×10^{-20}	2.4×10^0	1.5×10^{-1}
	<i>load_overlap</i>	4.1×10^{-21}	7.3×10^{-1}	4.1×10^{-1}
	<i>temporary_size</i>	1.4×10^{-21}	2.6×10^{-1}	6.2×10^{-1}
	<i>elements_number</i>	6.0×10^{-22}	1.1×10^{-1}	7.5×10^{-1}
	<i>1/elements_number</i>	2.7×10^{-21}	4.8×10^{-1}	5.0×10^{-1}
3 rd	<i>threads_number</i>	7.2×10^{-21}	1.3×10^0	2.9×10^{-1}
	<i>1/threads_number</i>	2.9×10^{-20}	5.1×10^0	4.0×10^{-2}
	<i>load_overlap</i>	7.4×10^{-25}	3.8×10^0	1.1×10^{-1}
	<i>temporary_size</i>	1.1×10^{-22}	5.7×10^2	2.4×10^{-1}
	<i>elements_number</i>	9.3×10^{-22}	4.7×10^3	1.2×10^{-8}
	<i>1/elements_number</i>	3.1×10^{-22}	1.6×10^3	1.9×10^{-7}

to better fit the target problem, and could be dispensed with altogether, according to the desired degree of automation and trust in the initial modeling assumptions. All experiments with this method presented in this thesis were completely automated using fixed significance thresholds.

Figure 12.6 shows the factors that were eliminated in each of the 3 steps performed by DLMT, and compares each factor's predicted best level with the level for that factor in the global optimum.

We see that the chose factor level matches the level on the global optimum in all factor fixing steps, except for the number of threads. This seems to have a small effect in the final results, but could be fixed in this specific case by dropping the inverse term, since the level on the global optimum is the highest one possible. This is an example of how this transparent approach can help learning about the problem begin optimized by challenging initial modeling assumptions. In this case, perhaps thread group management is a larger overhead than expected, and it is better to have larger groups.

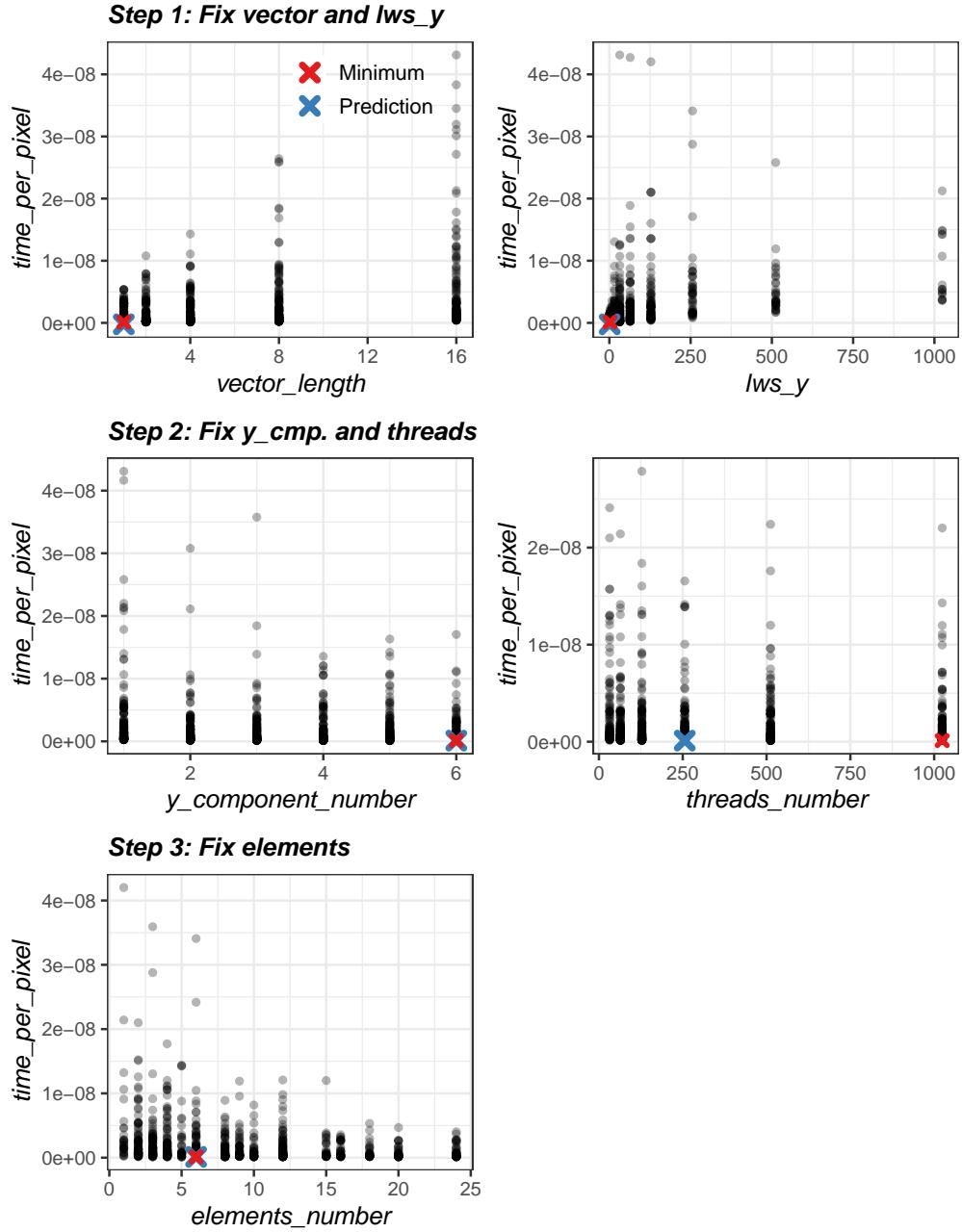


Figure 12.6: Overview of a single DLMT run. Each panel shows a view of the completely evaluated search space, from the perspective of one of the factors eliminated in a given step. A **X** marks the predicted best level of each factor, and a **x** marks the factor level on the global optimum.

12.5 Evaluation of Optimization Methods

This section presents a performance evaluation of the approach described in Section 12.2, comparing this approach to eight other optimization methods. The performance model described in the previous section was used by the Iterative Linear Model (LM) and the Quantile Regression (QR) methods, as well as by our approach, named *D-Optimal Designs, with Linear Model and heteroscedasticity correction Transform* (DLMT). The LM method is almost identical to our approach, described in Chapter 12, but it uses a fixed-size random sample of the search space instead of generating D-Optimal designs. Likewise, the QR method used random samples and Quantile Regression, in an attempt to decrease the impact of noisy outliers in the final fit. Additionally, we compared the performance of our approach with nine other methods, namely uniform Random Sampling (RS), Latin Hypercube Sampling (LHS), Greedy Search (GS), Greedy Search with Restart (GSR), and a Genetic Algorithm (GA). At a later date, after the CCGRID publication, we added to this performance comparison an implementation of Gaussian Process Regression with Expected Improvement acquisition function (GPR). Each method performed at most 125 measurements over 1000 repetitions, without user intervention.

Since we measured the entire valid search space for this kernel, we could use the *slowdown* relative to the global optimum to compare the performances of each method. Table 12.3 shows the mean, minimum, and maximum slowdowns in comparison to the global optimum for each method. It also shows the mean and maximum budget used by each method. Figure 12.7 presents histograms with the count of the slowdowns found by each of the 1000 repetitions. Arrows point the maximum slowdown found by each method. Note that maximum slowdown of the GS method was left out of range to help the comparison between the other method.

All methods performed relatively well in this kernel, with only GS not being able to find slowdowns smaller than 4 times in some of the runs. As expected, other search algorithms had results similar to RS. LM was able to find slowdowns close to the global optimum on most runs, but some of the runs could not find slowdowns smaller than 4 times. Our approach reached a slowdown of 1% from the global optimum in all of the 1000 runs while using at most fewer than half of the allotted budget.

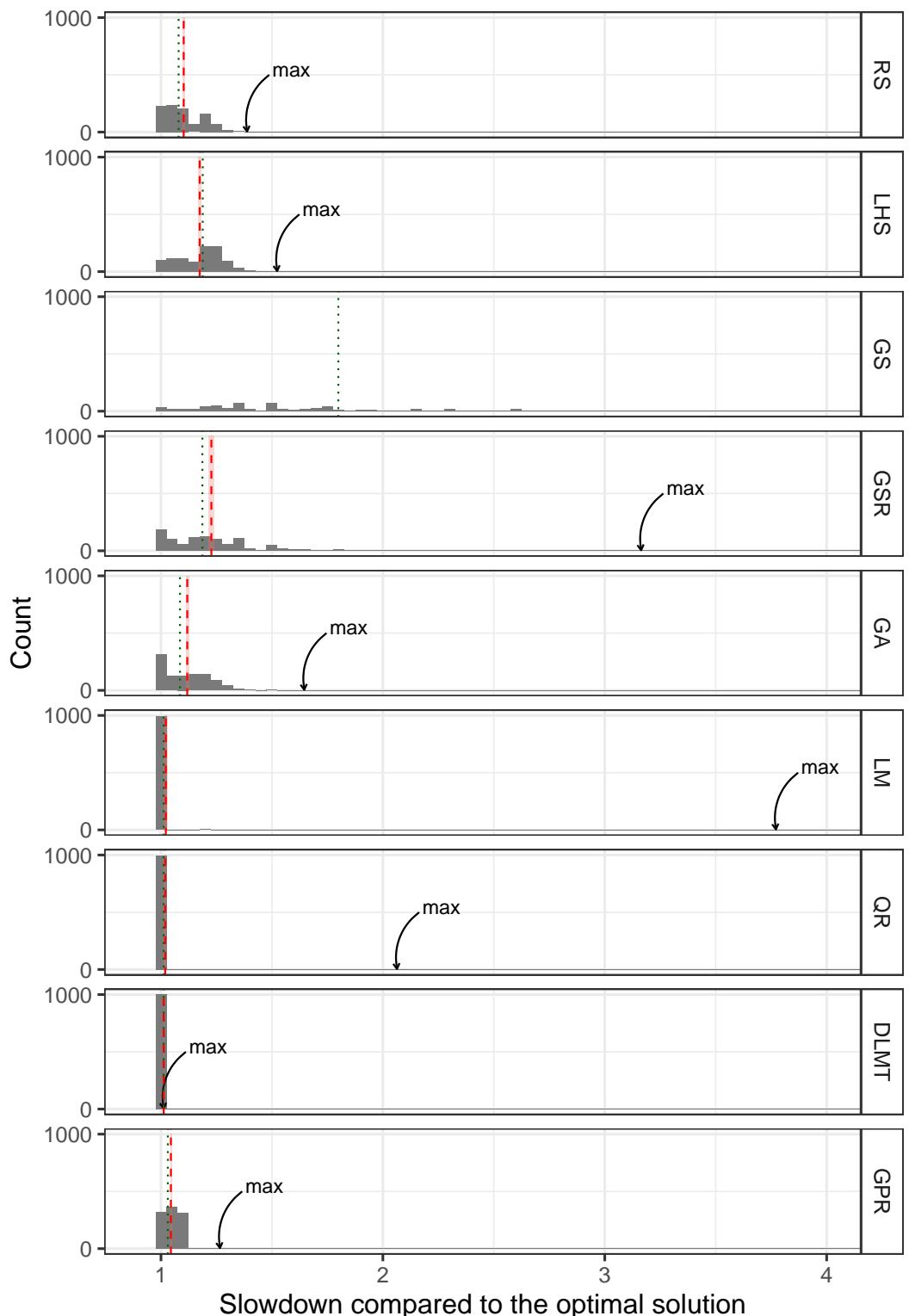


Figure 12.7: Distribution of slowdowns in relation to the global optimum for 7 optimization methods on the Laplacian Kernel, using a budget of 125 points over 1000 repetitions

We implemented a simple approach for the prediction step in this problem, choosing the best value of our fitted models on the complete set of valid level combinations. This was possible for this problem since all valid combinations were known. For problems where the search space is too large to be generated, we would have to either adapt this step and run the prediction on a sample.

This kernel provided ideal conditions for using our approach, where the performance model is approximately known and the complete valid search space is small enough to be used for prediction. The global optimum also appears to not be isolated in a region of points with bad performance, since our approach was able to exploit search space geometry. In Chapter 13 we present a more comprehensive performance evaluation of our approach in a larger benchmark suite.

12.6 Summary

This chapter described a transparent and parsimonious approach to autotuning, based on D-Optimal Designs, ANOVA tests, and Linear Regression. We dealt with the heteroscedasticity that could violate modeling assumptions by using standard detection tests and transformation procedures. This approach produced optimizations consistently within 1% of the global optimum of an OpenCL Laplacian kernel run in an NVIDIA GPU, using half of the allotted measurement budget for optimization.

We demonstrated how the choices made during optimization are completely transparent, in such a way that it is straightforward to envision how users could interfere with and learn from the optimization process by, for example, selecting specific factors that should be treated differently, and changing initial modeling assumptions based on the results produced.

The development, testing, and tweaking of this approach relied heavily on the availability of a completely evaluated search space, with known global optimum. Having access to this sort of search space is not possible in real applications, as we will see in the next two

Table 12.3: Slowdown and budget used by 7 optimization methods on the Laplacian Kernel, using a budget of 125 points with 1000 repetitions

Method	Slowdown			Budget	
	Mean	Min.	Max.	Mean	Max.
Random Sampling (RS)	1.10	1.00	1.39	120.00	120
Latin Hypercube Sampling (LHS)	1.17	1.00	1.52	98.92	125
Greedy Search (GS)	6.46	1.00	124.76	22.17	106
Greedy Search w. Restart (GSR)	1.23	1.00	3.16	120.00	120
Genetic Algorithm (GA)	1.12	1.00	1.65	120.00	120
Linear Model (LM)	1.02	1.01	3.77	119.00	119
Quantile Regression (QR)	1.02	1.01	2.06	119.00	119
D-Opt., Linear Model w. Transform (DLMT)	1.01	1.01	1.01	54.84	56
Gaussian Process Regression w. EI (GPR)	1.04	1.01	1.27	71.45	120

chapters, but it is also unfortunately uncommon in the autotuning literature. We hope the studies presented in this thesis, and the one in this chapter in particular, help arguing for and justifying the elaboration and evaluation of well-defined and completely-evaluated search spaces for autotuning, composed of high-interest computing kernels, that could be used for the development and comparison of optimization methods for autotuning.

Chapter 13

SPAPT Kernels and Orio

This chapter presents an application of Optimal Design and Gaussian Process Regression to the optimization of source-to-source transformation kernels from the SPAPT [29] benchmark suite.

The chapter is organized as follows. Section 13.1 discusses our choice of using D-Optimal designs, considering its applicability to autotuning among other Experimental Design methods. Section 13.2 presents the CPU kernels to which we applied our approach, and Section 13.3 discusses the results. The work presented up to Section 13.3 was published at CCGrid [4] in 2019. The studies we performed later were not yet published at the writing of this thesis. In Section 13.4 we explore the optimization of one of the target kernels and identify the factors and levels that were responsible for the observed performance improvements. In Section 13.5 we discuss the application of Gaussian Process Regression with Expected Improvement to the same kernel. Finally, Section 13.6 summarizes the chapter and discusses future work.

13.1 Choosing a Design of Experiments Method for Autotuning

Our application of the DOE methodology requires support for factors of different types and numbers of levels, such as binary flags, integer and floating point numerical values and enumerations of categorical values. We also need designs that minimize the number of experiments needed for identifying the most relevant factors of a problem, since at this moment we are not interested in a precise analytical model. The design construction techniques that fit these requirements are limited. The first DOE approach we studied was *screening* which, despite being extremely parsimonious, does not have enough flexibility to explore the search spaces which we wished to study.

Despite being parsimonious, screening designs are extremely restrictive. Algorithms for constructing Optimal Designs can adapt to our requirements, being able to mix factors of different types and optimizing a starting designs. Before settling on D-Optimal designs, which

we use in this study, we explored other design construction techniques such as extensions to Plackett-Burman [71] designs that use random factor levels, the *contractive replacement* technique presented by Addelman-Kempthorne [139], and the *direct generation* algorithm by Grömping and Fontana [140]. These techniques have strong requirements on design size and level mixing, so we opted for a more flexible technique that would enable exploring a more comprehensive class of autotuning problems.

The modified screening method provides a strategy to use factors with more than two levels by, instead of encoding fixed high and low levels, sampling two levels uniformly, each time the factor is included in a Plackett-Burman design. This approach has the advantage of a small design size with good main effect estimation capability, but it is still not capable of estimating interactions.

The *contractive replacement* method starts with a large 2-level design and generates mixed-level designs by re-encoding and replacing pairs of columns with a new column for a multi-level factor. The contractive replacement method presented by Addelman-Kempthorne [139] is a strategy of this kind. In addition to small design size and good estimation capability, their method maintains the orthogonality of starting designs, although it places strong requirements on initial designs, such that only orthogonal 2-level matrices can be contracted with their method.

The *direct generation* algorithm introduced by Grömping and Fontana [140] enables the direct generation of multi-level designs that minimize the *Generalized Minimum Aberration* [141] optimality criterion by solving mixed integer problems. We did not pursue this method because of the limitations it imposes on the size and the shape of the designs that can be generated, and also because it relied on proprietary MIP solvers that we did not have access to.

Weighting flexibility, effectiveness, parsimony, and the cost and availability of algorithmic construction, we picked *D-Optimal Designs* among the methods that fulfilled the requirements for application to autotuning problems. In particular, we used the KL-exchange algorithm [142], which enables mixing categorical and numerical factors in the same design, while biasing sampling according to the performance model we wish to explore. This enables the exploitation of global search space structures, if we use the right model. We can safely optimize for the *D* optimality criterion among other criteria without loosing quality of designs [143].

13.2 The SPAPT Benchmark Suite

The *Search Problems in Automatic Performance Tuning* (SPAPT) [29] benchmark suite provides parametrized CPU kernels from different HPC domains. The kernels shown in Table 13.1 are implemented using the code annotation and transformation tools provided by Orio [27]. Search space sizes are larger than in the Laplacian Kernel example. Kernel factors are either

Table 13.1: Kernels from the SPAPT benchmark used in this evaluation

Kernel	Operation	Factors	Size
atax	Matrix transp. & vector mult.	18	2.6×10^{16}
dgemv3	Scalar, vector & matrix mult.	49	3.8×10^{36}
gemver	Vector mult. & matrix add.	24	2.6×10^{22}
gesummv	Scalar, vector, & matrix mult.	11	5.3×10^9
hessian	Hessian computation	9	3.7×10^7
mm	Matrix multiplication	13	1.2×10^{12}
mvt	Matrix vector product & transp.	12	1.1×10^9
tensor	Tensor matrix mult.	20	1.2×10^{19}
trmm	Triangular matrix operations	25	3.7×10^{23}
bicg	Subkernel of BiCGStab	13	3.2×10^{11}
lu	LU decomposition	14	9.6×10^{12}
adi	Matrix sub., mult., & div.	20	6.0×10^{15}
jacobi	1-D Jacobi computation	11	5.3×10^9
seidel	Matrix factorization	15	1.3×10^{14}
stencil3d	3-D stencil computation	29	9.7×10^{27}
correlation	Correlation computation	21	4.5×10^{17}

integers in an interval, such as loop unrolling and register tiling amounts, or binary flags that control parallelization and vectorization.

We used the Random Sampling (RS) implementation available in Orio and integrated an implementation of our approach (DLMT) to the system. We omitted the other Orio algorithms because other studies using SPAPT kernels [83, 84] showed that their performance is similar to RS regarding budget usage. The global minima are not known for any of the problems, and search spaces are too large to allow complete measurements.

13.3 Performance Improvements using Optimal Design

We used the performance of each SPAPT kernel compiled with the `gcc -O3` flag, with no code transformations, as a *baseline* for computing the *speedups* achieved by each strategy. We performed 10 autotuning repetitions for each kernel using RS and DLMT, using a budget of *at most* 400 measurements. DLMT was allowed to perform only 4 of the iterations shown in Figure 12.5. Experiments were performed using Grid5000 [144], on *Debian Jessie*, using an *Intel Xeon E5-2630v3* CPU and `gcc` version 6.3.0.

The time to measure each kernel varied from a few seconds to up to 20 minutes. In testing, some transformations caused the compiler to enter an internal optimization process that did not stop for over 12 hours. We did not study why these cases delayed for so long, and implemented an execution timeout of 20 minutes, considering cases that took longer than that to compile to be runtime failures.

We automated factor elimination based on ANOVA tests so that a comprehensive evaluation could be performed. We also did not tailor initial performance models, which were the same for all kernels. Initial models had a linear term for each factor with two or more levels, plus quadratic and cubic terms for factors with sufficient levels. Although automation and

13.3. Performance Improvements using Optimal Design

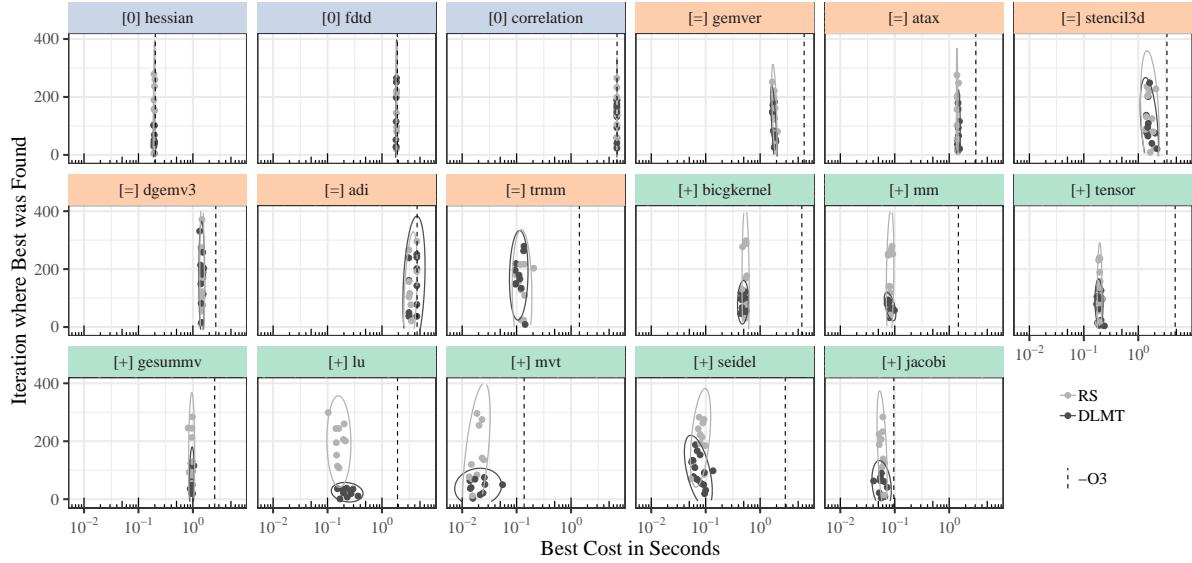


Figure 13.1: Cost of best points found on each run, and the iteration where they were found. RS and DLMT found no speedups with similar budgets for kernels marked with “[0]” and *blue* headers, and similar speedups with similar budgets for kernels marked with “[=]” and *orange* headers. DLMT found similar speedups using smaller budgets for kernels marked with “[+]” *green* headers. Ellipses delimit an estimate of where 95% of the underlying distribution lies

identical initial models might have limited the improvements at each step of our application, our results show that it still succeeded in decreasing the budget needed to find significant speedups for some kernels.

Figure 13.1 presents the *speedup* found by each run of RS and DLMT, plotted against the algorithm *iteration* where that speedup was found. We divided the kernels into 3 groups according to the results. The group where no algorithm found any speedups contains 3 kernels and is marked with “[0]” and *blue* headers. The group where both algorithms found similar speedups, in similar iterations, contains 6 kernels and is marked with “[=]” and *orange* headers. The group where DLMT found similar speedups using a significantly smaller budget than RS contains 8 kernels and is marked with “[+]” and *green* headers. Ellipses delimit an estimate of where 95% of the underlying distribution lies, and a dashed line marks the -03 baseline. In comparison to RS, our approach significantly decreased the average number of iterations needed to find speedups for the 8 kernels in the green group.

Figure 13.2 shows the search space exploration performed by RS and DLMT. It uses the same color groups as Figure 13.1, and shows the distribution of the speedups that were found during all repetitions of the experiments. Histogram areas corresponding to DLMT are usually smaller because it always stopped at 4 iterations, while RS always performed 400 measurements. This is particularly visible in *lu*, *mvt*, and *jacobi*. We also observe that the quantity of configurations with high speedups found by DLMT is higher, even for kernels on the orange group. This is noticeable in *gemver*, *bicgkernel*, *mm* and *tensor*, and means that

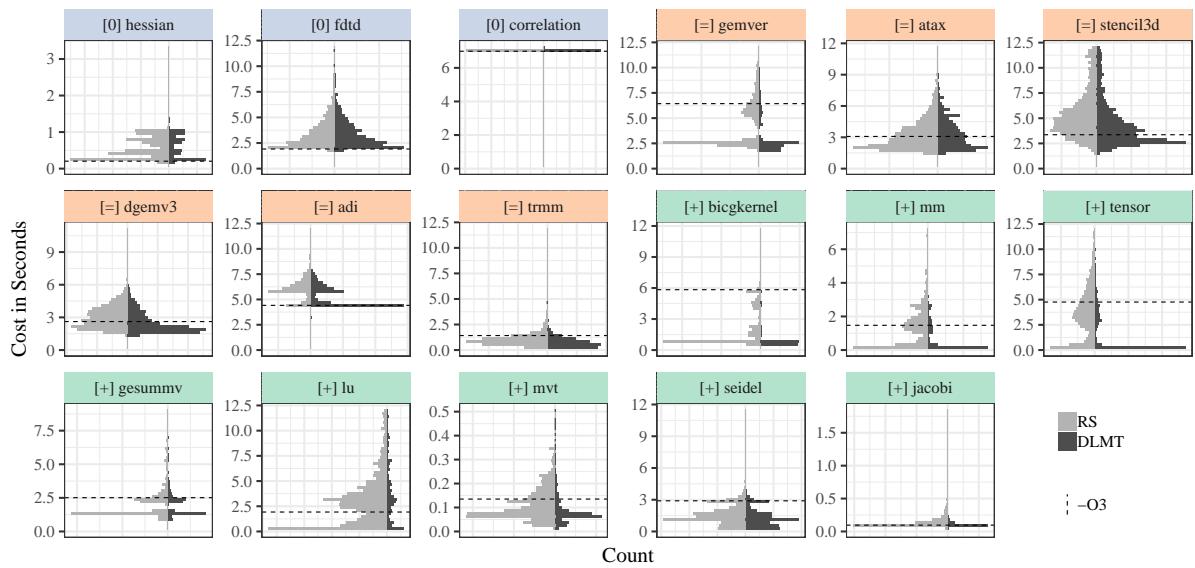


Figure 13.2: Histograms of explored search spaces, showing the real count of measured configurations. Kernels are grouped in the same way as in Figure 13.1. DLMT spent fewer measurements than RS in configurations with smaller speedups or with slowdowns, even for kernels in the orange group. DLMT also spent more time exploring configurations with larger speedups

DLMT spent less of the budget exploring configurations with small speedups or slowdowns, in comparison with RS.

Analyzing the significant performance parameters identified by our automated approach for every kernel, we were able to identify interesting relationships between parameters and performance. In *bicgkernel*, for example, DLMT identified a linear relationship for OpenMP and scalar replacement optimizations, and quadratic relationships between register and cache tiling, and loop unrolling. This is an example of the transparency in the optimization process that can be achieved with a DOE approach.

Our approach used a generic initial performance model for all kernels, but since it iteratively eliminates factors and model terms based on ANOVA tests, it was still able to exploit global search space structures for kernels in the orange and green groups. Even in this automated setting, the results with SPAPT kernels illustrate the ability our approach has to reduce the budget needed to find good speedups by efficiently exploring search spaces.

13.4 Identifying Significant Factors for the *bicg* Kernel

The discussion on the following sections was not yet published at the time of writing of this thesis. We continued to study the application of Design of Experiments methods to the optimization of SPAPT kernels, modifying the original DLMT algorithm by removing cubic terms to simplify the target performance model, fixing binary parameters to look for other

13.4. Identifying Significant Factors for the *bicg* Kernel

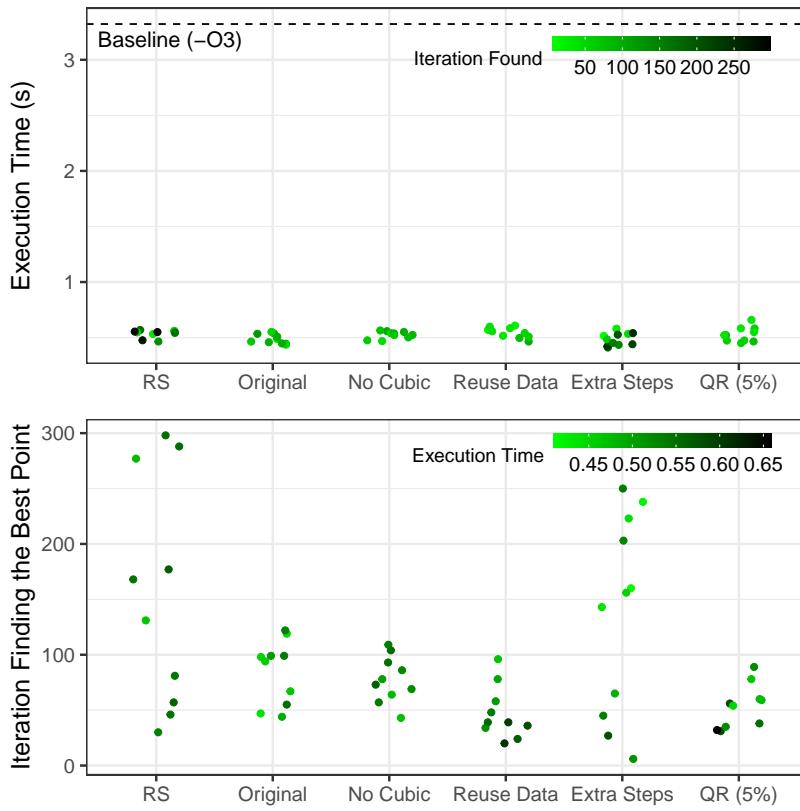


Figure 13.3: Summary of our DLMT results, compared to uniform Random Sampling (RS), showing the configurations with smallest speedups, found in 10 independent runs. The top panel compares execution time, with color-encoded iterations, and the bottom panel compares iterations where configurations were found, with color-encoded execution times.

significant effects, leveraging data collected in previous designs, running more optimization steps, and assessing the performance of quantile regression as an alternative to the linear model. In this study we sought to expose and analyze the inner workings of our approach.

In a non-automated setting, it becomes more clear how users could interfere and guide search space restriction and exploration, by intervening at each ANOVA step, for example. We started with the *bicg* SPAPT kernel because our approach achieved equivalent solutions using a smaller budget than a random sampling approach.

Figure 13.3 shows the execution times of the best kernel configurations, and the corresponding iterations where each of these points was found. Results are shown for Random Sampling, for our DLMT implementation from the CCGRID paper, and for four variations that we later explored, attempting to interpret and improve the optimization process. The *y*-axis on the top panel shows the execution time of the best points found in 10 distinct runs, and lighter colors mark points found with fewer iterations. Similarly, the *y*-axis on the bottom panel shows the iterations where points were found for the same data, with lighter colors marking points with lower execution times.

The next sections discuss incremental modifications to the algorithm published at CCGRID, which were performed to help identify significant factors and attempt to find better speedups. Although we did not significantly improve the speedups found by our method, we were able to present a deeper look into the inner workings of this white-box autotuning approach. The next section describe in detail the two incremental modifications below:

Removing Cubic Terms We initially used third-degree polynomials for the performance models of all kernels. Cubic terms were never eliminated when following significance levels reported by ANOVA, so we decided to drop cubic terms and fit a quadratic polynomial, which reduced the necessary design size and produced equivalent results.

Reusing Data In our initial approach, we performed ANOVA in each iteration using only the data from the subspaces to which we restricted exploration. This can increase flexibility because the models are free to vary outside the restricted regions, but this also wastes experimental data. In this modification we reused all measurements performed during optimization, leveraging all measurements regardless of their position in the search space, which increases the benefits of multiple runs.

The two other modifications shown in Figure 13.3 are:

Quantile Regression We identified a substantial amount of noise on the performance measurements of non-binary parameters. Models accounting for factor interactions could help, since there is a good chance that interactions are responsible for the noise. In this study, we decided to try Quantile Regression, which is a modification of Linear Regression where it is possible to weight points based on the quantiles they belong to, allowing the model to fit to a different quantile, opposition to fitting to the mean. We explored multiple quantiles, and Figure 13.3 shows results using the 5% quantile.

Running more Steps This is less of a modification, and involved simply running the original approach for twice the number of steps.

13.4.1 Removing Cubic Terms

The data in the figures shown next were obtained with a new set of 10 repetitions of the *bicg* kernel experiment from Figures 13.1 and 13.2, but using a performance model with only linear and quadratic terms. The results are similar to the ones from the CCGRID paper, but our approach found slightly better configurations slightly faster. Random sampling also found a much better configuration much faster than before in one experiment. Figure 13.4 shows the *count of model terms* that were eliminated at each color-coded DLMT step. As before, we see that *OMP* and *SCR* were the most eliminated factors, especially on the first step.

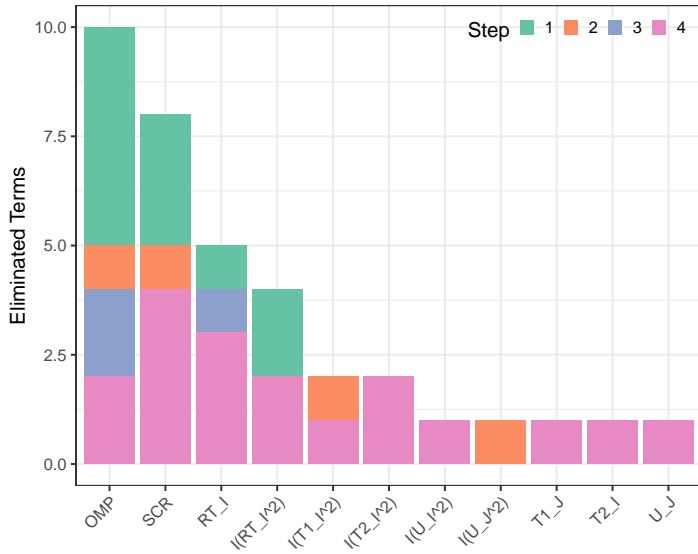


Figure 13.4: Count of the model terms eliminated in each of the 4 steps of 10 independent runs, without cubic terms

Figure 13.5 shows the factors that were identified as significant by ANOVA, with significance threshold of 0.01, at each DLMT step. Identifying a factor removes the corresponding parameter from the model, regardless of the model term that was identified. Figures are grouped by each of the 10 experiments. Figure headers identify each run, and correspond to the names of the Grid5000 machines where experiments were run.

It is interesting that only *parasilo-19*, *parasilo-20*, and *parasilo-9* eliminated any factor other than *OMP* and *SCR* on the first step where any factor was eliminated, and also that those runs fixed the factor *RT_I* to the same value. We can also see that *OMP* seems to be the parameter behind the most extreme changes in the execution time of tested configurations. This is specially clear at *parasilo-13* and *parasilo-10*, where the explored configurations have relatively high execution time until after step 3, where *OMP* is fixed. A slight worsening of execution times, an increase, that is, can be seen at the fourth step at *parasilo-11*, after *RT_I* was fixed. The minimum execution time seems to be higher than in the third step.

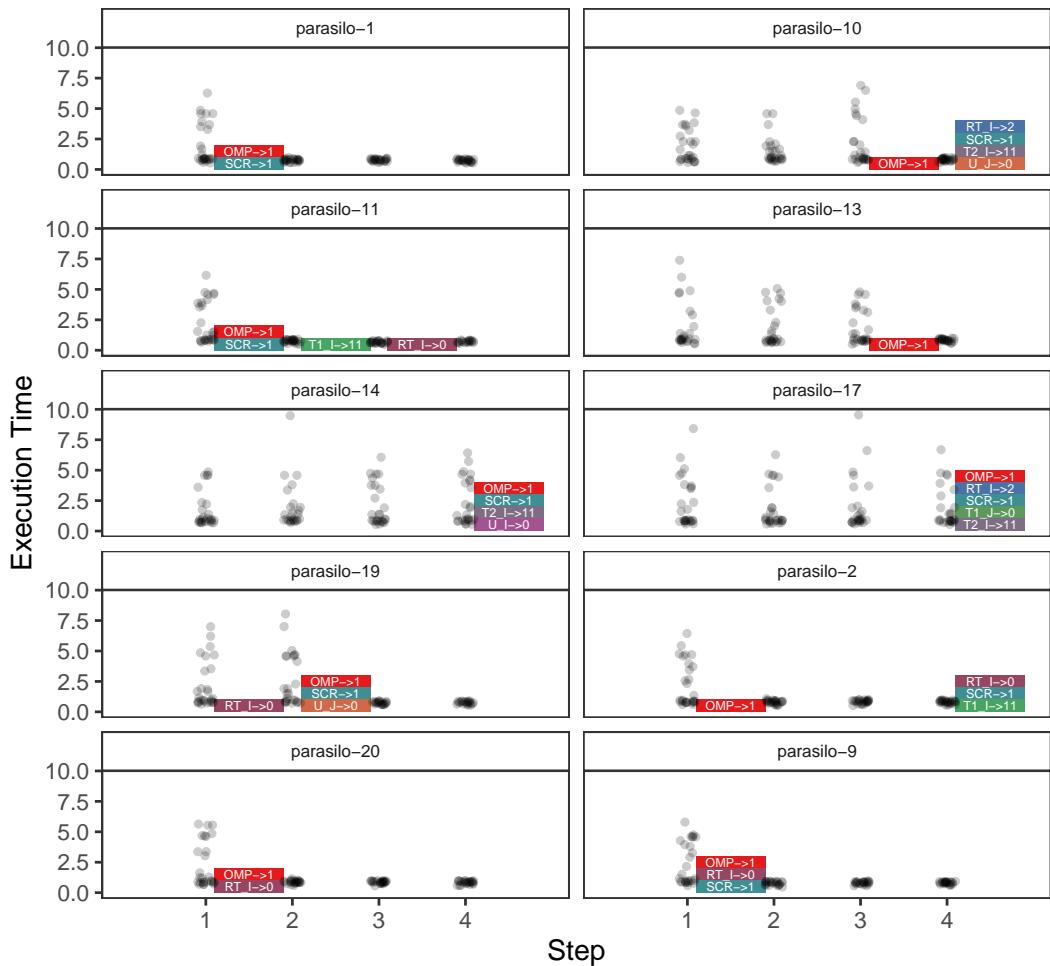


Figure 13.5: Measured execution time of all points in the designs constructed at each step, where panel headers mark the hostname of the machines used in each of the 10 separate experiments

13.4. Identifying Significant Factors for the *bicg* Kernel

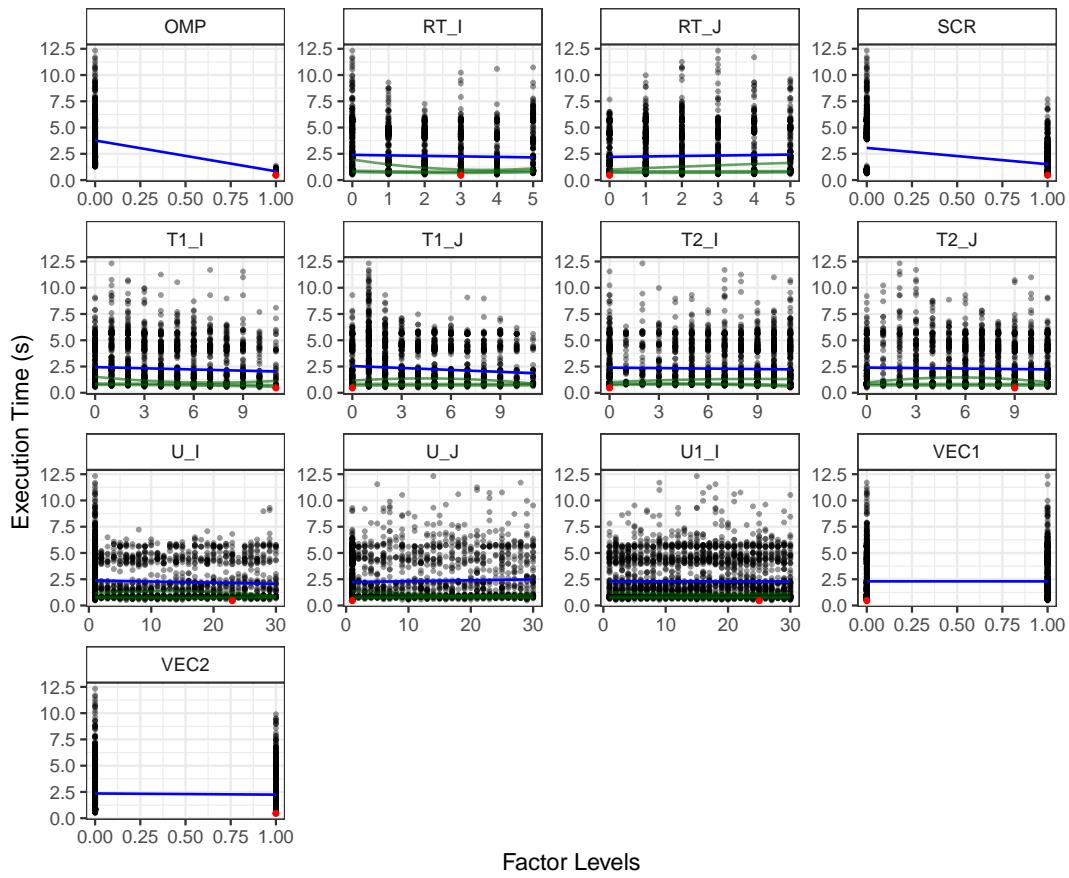


Figure 13.6: Fitting linear models and quadratic quantile regression, with $\tau = 0.05, 0.25$, and 0.5 , to separate factors for 10 uniform random sampling runs

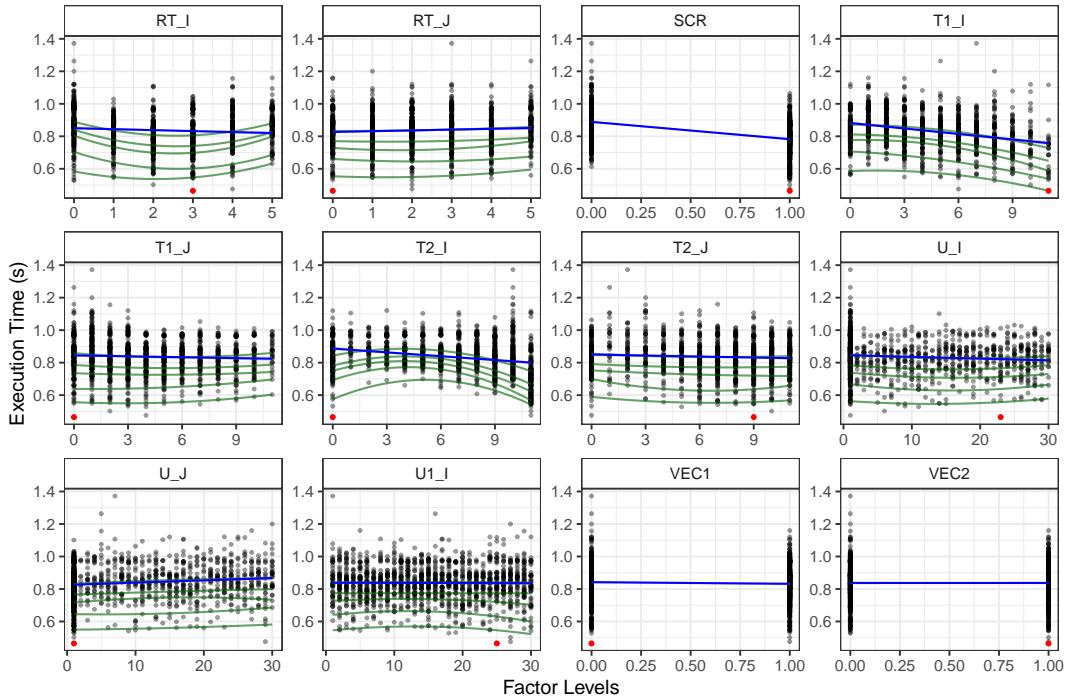


Figure 13.7: Same as above, for 10 uniform sampling runs with *OMP* fixed to *on*

Figures 13.6 and 13.7 contain all points measured in 10 uniform random sampling runs. Each panel the same execution time data in the y -axis, viewed from the perspective of the changes in a single factor, shown in the x -axis. The blue lines in both figures represents a linear model fit, using only linear terms for each factor separately. The green lines denote three quantile regression fits using quadratic model terms for each factor separately. We used the 5%, 25%, and 50% quantiles. Using the 50% quantile is equivalent to linear regression. Red dots mark the values of each factor in the point with the best performance across the 10 runs.

The OMP parameter controls whether OpenMP is used to run multithreaded execution, and has clearly the largest impact when its level is 1.0, which encodes the binary level corresponding to turning parallelization on. All measurements larger than 1.4 seconds happened when the parallelization flag was off, as can be seen on Figure 13.7, which shows all data collected in 10 uniform random sampling runs with the OMP factor fixed on the *on* level. The SCR parameter also has a strong effect. This parameter controls *scalar replacement*, which consists of moving array references to the stack, preferably into registers.

Fixing OMP to *on* makes the effects of other factors more pronounced and easier to detect using ANOVA. Most of the best factor levels in this example were near the limits of the search space. The RT_I factor controls tiling of one of the loops in *bicg*, and its best level was exceptionally at the middle of the range.

We can also verify that the optimal design approach is restricting the exploration of the search space to regions with faster configurations. Figure 13.8 compares the performance predicted by the fitted model, represented by the green line, with the performance of the best point found at each design, represented by the blue line. The red line marks the best point found so far.

It is only at the fourth step of *parasilo-17* that the point with the best predicted performance was better than the best point on the design for that step, while also being better than the best point found so far. Although we seem to be effectively restricting the search space with our exploration, which is evidenced by the improvement that occurs as steps progress and by the best points being found inside designs, the models fit using experiment data are not able to improve the current best on the majority of cases. This is an indication that the model is not capable of describing the relationships between factors and response, perhaps because of the presence of interactions between parameters.

13.4.2 Reusing Data from All Steps

Significant changes were performed on the initial DLMT implementation. We decided that there was no good reason to not reuse the data obtained from evaluating designs at each step, and the various samples of the search space taken at different points. Now, all evaluated experiments compose a single growing design used by *aov* to identify the best

13.4. Identifying Significant Factors for the *bicg* Kernel

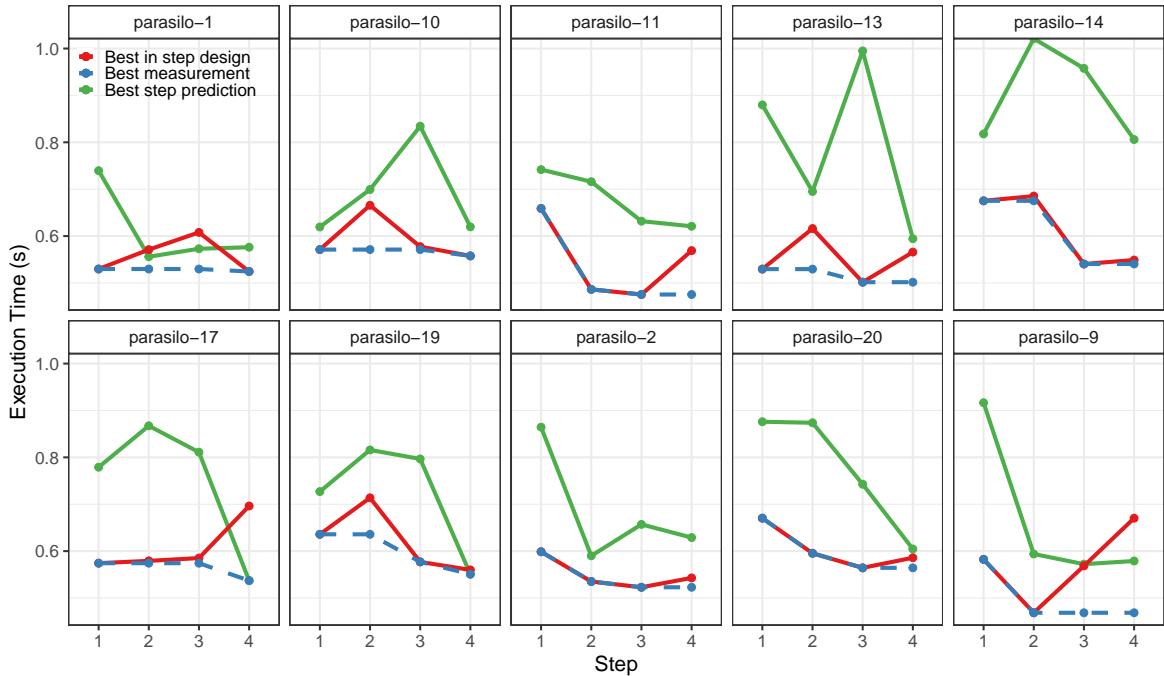


Figure 13.8: Performance of best points, step-by-step, by experiment

factors, and all samples from the search space compose a single data set used by *optFederov* to select experiments. The data set is pruned in both *aov* and *lm* analyses, to guarantee only experiments with the correct levels of fixed factors are used. This is crucial for both analyses, since having different fixed values of factors that are not in the model would imply that we would have inexplicable variance in the data set.

Using all experimental data on *aov* is interesting because it is always worth it to consider additional information on the factors that are currently being studied. On one hand, it might not allow for enough flexibility when we consider regression only on a small restricted subspace, because points outside the subspace would impact regression, and we would be interested in the significance of the factor inside the subspace at the moment. On the other hand, using all data available makes sense because we are also interested in exploring any global structures of the search space, and keeping points from other subspaces would increase the chance of “catching” the significance of a factor globally.

Using all sampled space, across all steps, as a candidate for *optFederov* has no downsides, provided we prune the search space to account for current constraints on factor levels fixed on previous steps. We increase the size of the available set of configurations that can compose a new design each time we sample a new subspace. This would hopefully improve the quality of designs produced as we progress.

13.4.3 Summary: Tuning SPAPT Kernels with Optimal Design

We could not significantly improve upon the speedups found by the uniform random sampling strategy, although we did find good speedups using a significantly smaller budget. This happens mainly because the DLMT approach quickly detects the large effect of the *OMP* binary parameter. Subsequent parameters that are detected as significant have much smaller effects, so consequently the performance does not improve by much.

We can conclude that the SPAPT search space structures defined by kernel parameters exposed and modified by Orio cannot be fully exploited by models based on second or third degree polynomials, even with well designed experiments. We could not determine whether kernel configurations with better performance than the ones we found exist in these search spaces. Additionally, the strongly biased models we used could be responsible for not finding better configurations. We attempted to find better configurations using a much more flexible approach, based on Gaussian Process Regression and Expected Improvement, which we discuss next. We also compared the performance of the *dgemv* SPAPT kernel with the peak theoretical performance for the *Xeon 2630v3* processor, indicating that better kernel configurations might exist.

13.5 Autotuning SPAPT kernels with Gaussian Process Regression and Expected Improvement

We used Gaussian Process Regression with an Expected Improvement acquisition function, denoted GPR, to try and find *bicg* kernel configurations improving upon those found by the uniform Random Sampling (RS) method. We found statistically but not practically significant improvements upon RS.

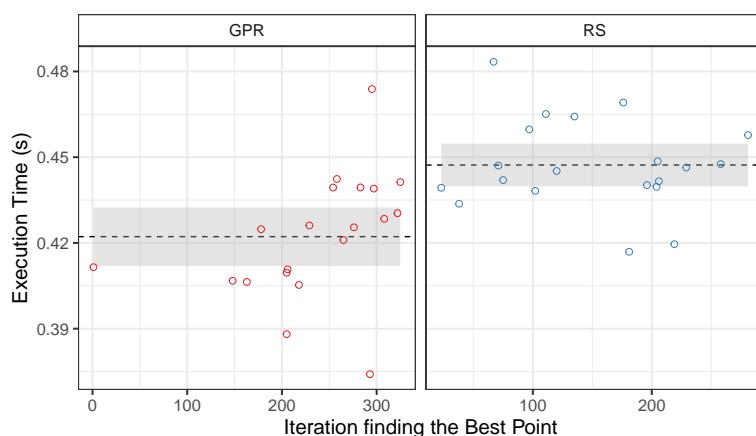


Figure 13.9: Best *bicg* configurations and iterations that found them, for Gaussian Process Regression with Expected Improvement acquisition function (GPR), and uniform Random Sampling (RS), with a budget of 400 measurements.

We used two *R* packages to construct our GPR algorithm [67]. We used the *DiceKriging* package to fit Gaussian Processes using *Matérn* covariance functions, and the *DiceOptim* package to compute the Expected Improvement metric at each step. Figure 13.9 compares the performance of the best points found by GPR and RS, across 20 independent tuning runs. The figure also shows the iteration finding the best point at each run. Dashed lines mark the estimate of the mean execution time of the best points across the 20 experiments, and the gray bands show the 95% confidence intervals for the mean estimate.

The GPR method found statistically significant better configurations than RS, but the practical difference is less than 0.1 seconds. Since Gaussian Process Regression is a much more flexible method than our previous Design of Experiments approach, not finding expressively better configurations hints that it might not be possible to improve this kernel much more, in comparison to RS.

Figures 13.10 and 13.11 show the progression of all configurations evaluated during each of the 20 runs of GPR and RS, respectively. We can see that after an initial Sobol sample the GPR method quickly finds regions with better configurations. The RS method, as expected, continues to randomly explore the search space.

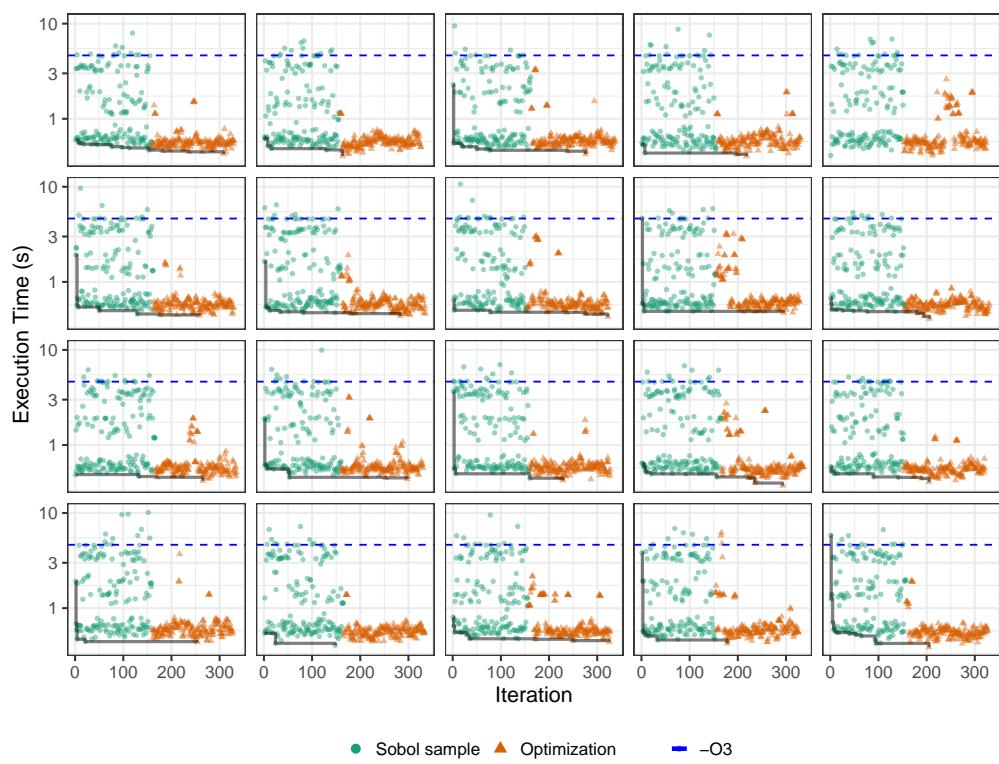


Figure 13.10: Execution time of points measured by GPR along iterations, with pareto border in red

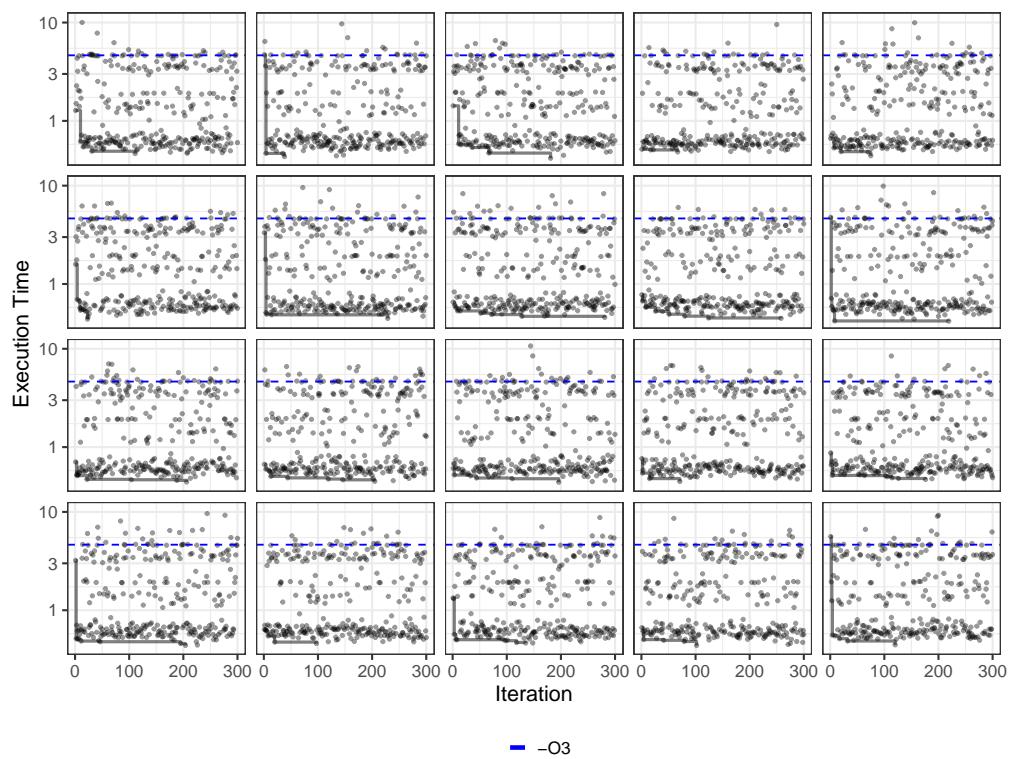


Figure 13.11: Execution time of points measured by RS along iterations



Figure 13.12: Best $dgemv3$ configurations and iterations that found them, for Gaussian Process Regression with Expected Improvement acquisition function (GPR), and uniform Random Sampling (RS), with a budget of 400 measurements.

13.5.1 Peak Performance for the DGEMV kernel

We also performed GPR experiments with the $dgemv3$ SPAPT kernel which, since we can compute its theoretical performance, allowed determining that the best kernel configuration we have found is still around 20 times slower than the theoretical peak, despite being around two times faster than the $-O3$ flag. We did not perform a detailed analysis of the SPAPT kernel, which could find out whether more performance can be obtained from this kernel implementation.

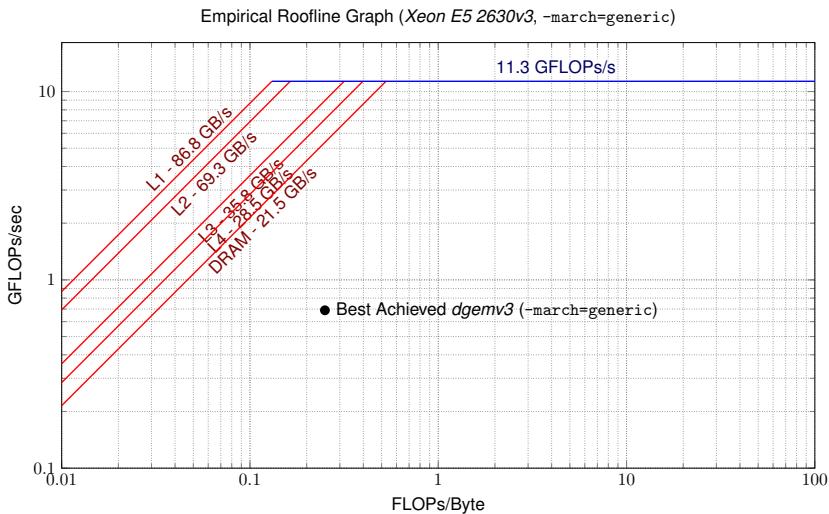


Figure 13.13: Empirical Roofline graph for the Xeon E5-2630v3 processor [145, 146], showing the best point we found during all experiments with the memory-bound $dgemv3$ SPAPT kernel

Figure 13.12 shows the best kernel configurations found across 20 repetitions of the GPR and RS methods. GPR could not find statistically or practically significant speedups in relation to RS for this kernel, although one specific GPR run found an outlier configuration. Figures 13.13 and 13.14 show the empirical Roofline [147] graphs for the processor targeted by our

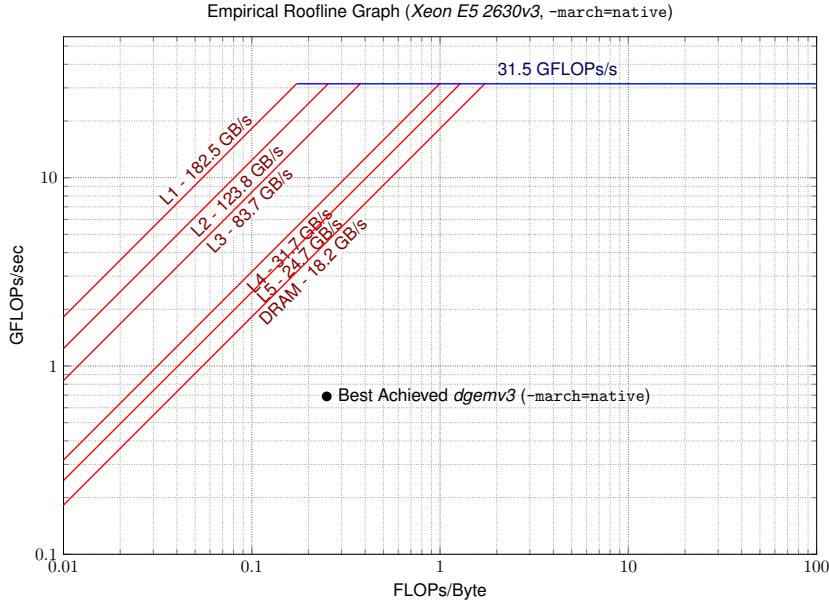


Figure 13.14: Empirical Roofline graph for the *Xeon E5-2630v3* processor [145, 146], showing the best point we found during all experiments with the memory-bound *dgemv3* SPAPT kernel

experiments, which was the same from our CCGRID paper, using *GCC -march=generic,native* options. The figure also shows the performance that was achieved by the best configuration we found for the *dgemv3* SPAPT kernel, the red outlier in Figure 13.12.

13.6 Summary

We have implemented a Design of Experiments approach to autotuning based on ANOVA and D-Optimal Designs. We have applied this approach to a set of configurable kernels from the SPAPT benchmark suite and, although our approach found significant speedups in relation to an *-O3* compiler baseline, we could not significantly improve upon the results found by an uniform Random Sampling algorithm using the same exploratory budget. Other experiments with the same kernels also showed than uniform random sampling [83, 84] performs well, especially when budgets are short.

We could identify and eliminate significant factors using our approach, but the largest effect was consistently that of the binary flags, especially *OMP*, which controlled parallelization. We could consistently detect significant albeit small effects for a few other numerical parameters.

The polynomial models coupled with optimal design construction have the potential for low prediction variance, but tend to have strong biases. We explored a Gaussian Process Regression approach which is much more complex, implying in high variance, but has small bias. We found small improvements upon Random Sampling with this new approach, but it is still unclear if Orio configurations are capable of further improving the kernel configurations

we already found. Analyzing the Roofline theoretical model for our target CPU hints that tweaks to the kernel implementation might be needed to improve the current performance.

A possible future direction for exploring the potential of this Design of Experiments approach could involve developing a lower level set of configurable kernels, where performance measurements are always tied to a theoretical peak or to a known global optimum. As we discussed in Chapter 12, knowing the global optimum helps evaluating the performance of an autotuning approach.

Chapter 14

Mixed-Precision Quantization for *ResNet50*

This chapter presents an application of Gaussian Process Regression with Expected Improvement acquisition function (GPR) to bit precision tuning for a Convolutional Neural Network (CNN). We compare GPR to a Reinforcement Learning approach which uses an actor-critic model (RL) [148], to a uniform random sampling algorithm, and to a low-discrepancy sampler using Sobol sequences. The objective function we attempted to optimize was the network accuracy over the ImageNet dataset [149, 150], subject to a limit imposed on network size. Our GPR approach achieved comparable results to the baseline RL approach, but the random samplers also performed well, especially the low-discrepancy Sobol sequence.

This work started during a two-month stay at *Hewlett Packard Enterprise* (HPE), former *HP Labs*, in Palo Alto, California. The problem of tuning mixed-precision quantization for CNNs is of interest to HPE researchers because it produces smaller and more energy-efficient networks, which are especially suited to embedded and custom-built hardware. We expanded the initial study in the last year of this thesis, still in collaboration with HPE researchers, but this work is still unpublished at the time of writing.

The remainder of this chapter is organized as follows. Section 14.1 describes the *HAQ* framework [148] for mixed-precision bit quantization, which we used to change the bitwidth of network layers, and which implements the baseline RL approach. Section 14.2 presents the *ResNet50* network used in our experiments and the ImageNet dataset. Section 14.4 discusses the implementation details of the methods we compared, and Section 14.5 evaluates their performance with respect to network accuracy. Section 14.6 presents a more detailed assessment of the GPR method we implemented. Section 4.3 summarizes the discussion and concludes the chapter.

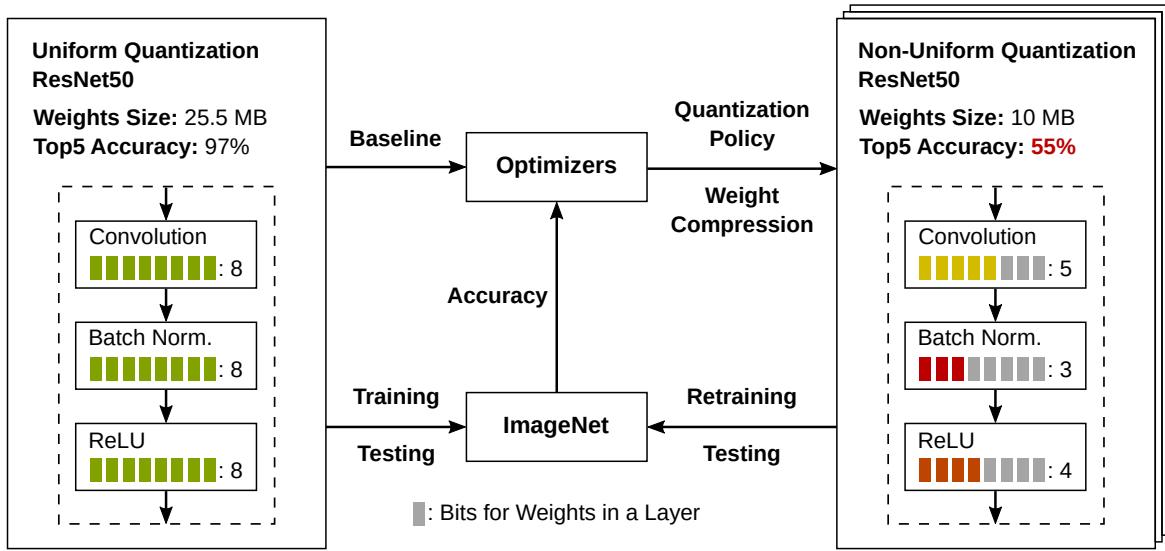


Figure 14.1: Optimizing nonuniform weight quantization policies for the *ResNet50* network, keeping total weight size below 10 MB, and attempting to keep accuracy. The optimization methods we compared were a Sobol Sampler, Reinforcement Learning [148], and Gaussian Process Regression with Expected Improvement

14.1 Autotuning Bit Precision for Convolutional Neural Networks

A Neural Network is composed of interconnected layers of "neurons". Each neuron can be interpreted as an agent that receives *weighted* inputs from connections with other neurons, combines them, and produces an output if the combined input crosses an *activation* threshold. The process of *quantization* of a network consists in choosing, according to a performance metric, the best bit precision for data representing weights and activation thresholds in each layer. Quantizing a network consequently determines the bit precision of the computations performed during training and inference. Quantization is increasingly relevant in Neural Network research [148, 151, 152, 153, 154, 155] because it reduces energy and space requirements, enabling deployment in embedded, mobile, or custom-engineered hardware.

Figure 14.1 shows the bit precision tuning experiments we present in this chapter. We used different methods to optimize quantization policies for the layer groups of the PyTorch [156] implementation of a *Residual Network* [157] with 50 residual blocks, called *ResNet50*. We used the HAQ framework [148] to quantize each layer with a bitwidth in the [1, 8] interval, producing a smaller network that was retrained on NVIDIA GPUs. We ensured that quantization policies respected size constraints, keeping total weight size under 10 MB, and allowing the comparison of our results to a reproduction of the approach of the original paper [148].

14.1.1 Further Applications of Autotuning to Neural Networks

Autotuning methods can be applied to optimize networks targeting different stages and structures. Online Gradient Descent coupled with backpropagation is a well established approach to network training [158, 159], but fine tuning of the network architecture and its parameters is also a promising target for the application of optimization methods.

A recent survey and taxonomy by Talbi [160] discusses work on *Neural Architecture Search* (NAS), and *Hyperparameter Tuning*. Neural Architecture Search consists of using optimization methods to find the modifications of the layers and connections of a network that best adapt to resource constraints or hardware architecture features. Hyperparameter Tuning refers to the optimization of any network parameters that have an impact on performance without necessarily requiring structure modification.

14.2 ResNet50 and ImageNet

14.2.1 ResNet50 Architecture

Residual Networks, or *ResNets*, were introduced by He *et al.* [157], and the intuition behind the idea is to add the original layer input, or *residual*, to the output of a stack of convolution layers. Keeping this residual component in a layer output enables better global optimization of the network. Figure 14.2 shows the *PyTorch* implementation of a 50-block Residual Network, called *ResNet50*.

We generated this visualization of *ResNet50*'s structure with the *HiddenLayer* library [161], adding the underlying structure of *Bottleneck Residual Blocks* in panel (a), and of *Residual Blocks* in panel (b). Bottleneck blocks serve the same purpose as Residual blocks, adding an input component to layer output, but are favored in larger networks because they are computationally less expensive.

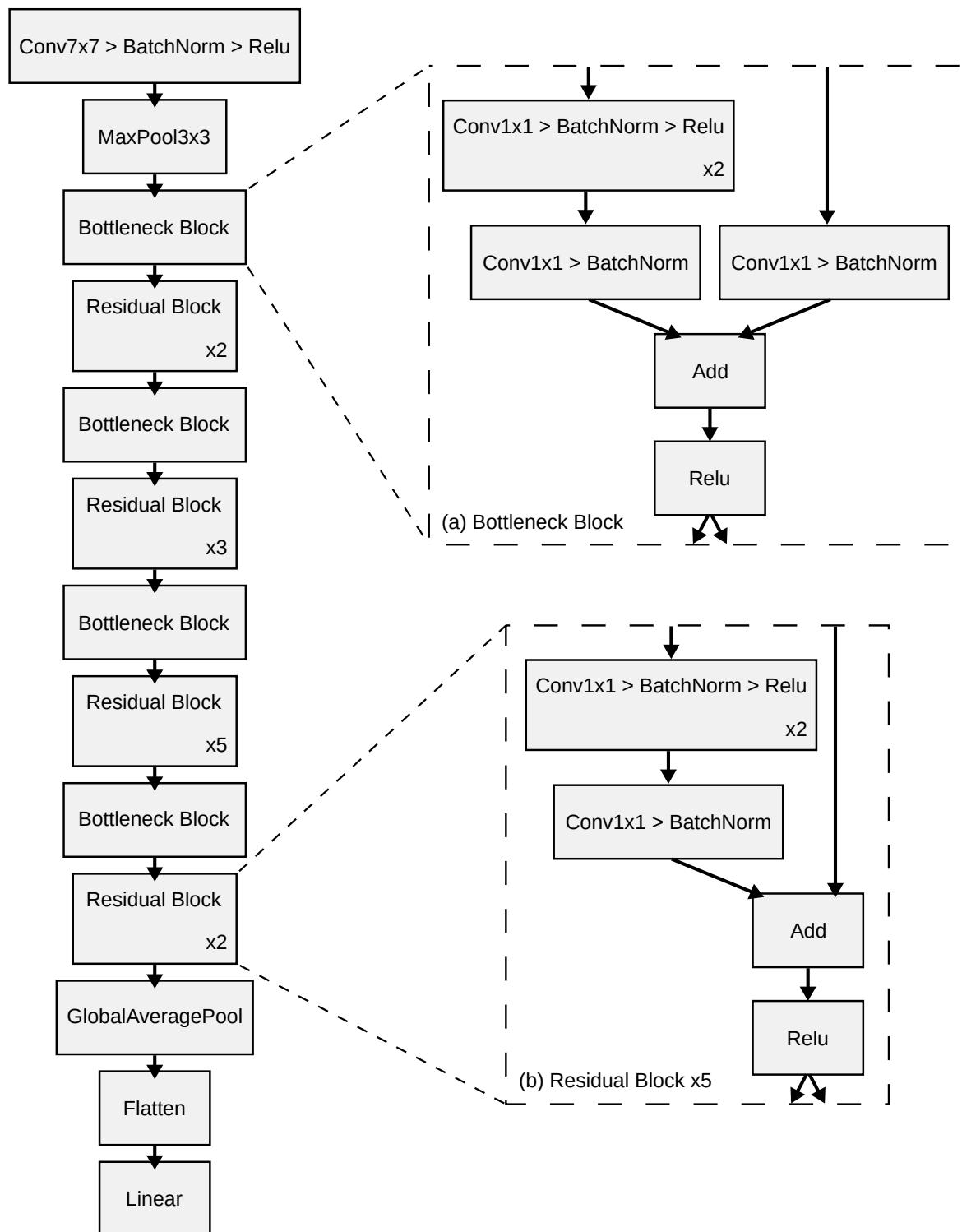
Figure 14.2: ResNet50 architecture, implemented in *pytorch*



Figure 14.3: Images sampled from some of the categories in *ImageNet*, adapted from Deng *et al.* [149]

14.2.2 The ImageNet Dataset

In this study we were interested in maintaining the accuracy of *ResNet50* on the *ImageNet* dataset as we decreased the bit precision for a set of network layers. *ImageNet* [149] is a massive image dataset, with 3.2 million labeled and semantically categorized images. Figure 14.3 shows a sample of the dataset. *ImageNet* is widely used to validate works on image recognition, and the dataset version we obtained for the work presented in this chapter came from a series of visual recognition challenges [150].

In each repetition the experiments discussed in Section 14.5 we used distinct subset of 30 thousand uniformly sampled images, of which 20 thousand were used for training and 10 thousand for validation.

14.3 Search Space and Objective Function

The search space we explored consisted of 54 convolutional layers, that is, 54 component layers of the 50 residual blocks in the *PyTorch ResNet50* implementation. We optimized only the bit precision for weights, since the code that quantized activation precision was only released by the authors after this study was completed.

Choosing between a bit precision in the $[1, 8]$ interval for the 54 quantizable layers generates a search space with $8^{54} \approx 10^{49}$ configurations. This search space is relatively large, although not the largest we studied in this thesis. Completely evaluating it is unfeasible, considering we would have to retrain and validate the quantized network on our 30 thousand image sample of *ImageNet*, which took tens of minutes for each quantization on an NVIDIA K80 GPU.

The accuracy of a network is typically measured by its *Top1* and *Top5* accuracy measures. Both metrics measure the ratio between right and wrong network predictions in the testing dataset, but the *Top1* accuracy considers as a prediction only the class with the highest predicted probability, while the *Top5* metric considers the classes with the top five highest

predicted probabilities. For example, if a network predicts a picture of a Chihuahua is in the "Muffin" class with 85% probability and in the "Dog" class with 80% probability, this would be considered a miss in *Top1* accuracy but a hit in *Top5*.

In this work we were interested in minimizing the distance between the 8-bit uniformly quantized *Top5* accuracy and the *Top5* for a quantized network whose total weight size was under 10 MB. We targeted *Top5* because, from the perspective of the HPE researchers we collaborated with, it was more interesting that the network was made to fit on the smallest possible device, perhaps at the cost of some *Top1* accuracy, because then it would be cheap to run a larger set of predictions. In the end, suggesting the two metrics are linked, optimizing for *Top5* produced *Top1* results comparable to the ones found by the original paper [148] and by our replication of their results in our settings.

We also performed experiments that added the total weight size to the objective function and dropped the 10 MB restriction, but we found it hard to compare the results we produced with the original paper, so the main discussion that follows will regard the performance of optimization methods on the *Top5* distance minimization.

14.4 Optimization Methods

This section describes the optimization methods we applied to the quantization policy optimization for *ResNet50*. The performance comparison between methods is presented in the next section.

14.4.1 Random and Sobol Sampling

Although practically impossible, in order to properly explore the *ResNet50* quantization search space it would be ideal know the global optimum. In the absence of this knowledge, it is useful to compare the performance of the optimization methods we wish to use to simple uniform random sampling. Since in high dimensional spaces almost all uniformly sampled points would lie in the extremes of the search space, we also compare the performance of methods with a sampling algorithm using Sobol sequences, that produce low-discrepancy, space-filling, samples.

14.4.2 Gaussian Process Regression with Expected Improvement

We used the same two *R* packages we used in Chapters 13 and 12 to implement the GPR algorithm [67]. The *DiceKriging* package provided Gaussian Processes Regression fits using the *Matérn* covariance kernel, and the *DiceOptim* package provided functions to compute the Expected Improvement (EI) metric at each step, as described in Section 7.2. Additionally, we used the *future.apply* package [162] to run EI computations in parallel.

Respecting the Size Constraint with Sobol Samples

We used the implementation of the Sobol sequence low-discrepancy sampling algorithm from the *randtoolbox* [79] R package. Low-discrepancy sampling and space-filling designs were discussed in Section 6.4. To ensure all samples respected the size constraint, we progressively doubled the sample size until enough valid configurations were produced.

Pseudocode for the GPR Method

Figure 14.4 shows the pseudocode for our adaptation of GPR to the bit precision tuning problem, including the constrained Sobol sampling algorithm. The algorithm starts with a smaller valid Sobol sample which is completely evaluated and used to fit a Gaussian Process. This model fit enables computing the Expected Improvement for a much larger valid Sobol sample.

We pick the three configurations with the highest EI from this sample and compute additionally the EI for a sampled neighborhood around these configurations. The neighborhood is also generated using the Sobol sampling algorithm.

We then measure the unique configuration with the highest EI over this larger sample, augmented with a targeted neighborhood, and keep a backlog of the 200 other configurations with the highest EI in the sample. The EI of the configurations in the backlog is constantly reevaluated, and the backlog is used if the highest EI in the current sample is smaller than the EI of any configuration in the backlog.

14.4.3 Reinforcement Learning: The Baseline Method

We reused the Reinforcement Learning (RL) *Python* implementation provided by the authors of the *HAQ* paper [148]. Reinforcement Learning using the *actor-critic* model was discussed in Section 7.3, and the implementation in *HAQ* employed additionally the *Deep Deterministic Policy Gradient* method from Lillicrap *et al.* [163]. The cumulative regret to minimize is the distance of the quantized accuracy from the initial 8-bit uniform quantization, across all tested quantization policies.

Respecting the Size Constraint

To respect the constraint imposed on the sum of network weights, the original paper implementation used a *round-robin* algorithms to decrease the weight bitwidth for each layer sequentially, starting with the first, until the total sum crossed the 10 MB threshold. After analyzing our results with RL and other methods, our hypothesis for this implementation choice is that it encodes the implicit knowledge that the bit precision of the initial layers of the *PyTorch ResNet50* implementation has a smaller impact on accuracy.

```

1: function FILTER_SOBOl(dimension, target_size, max_weight)
2:   sample_size  $\leftarrow 10^6$ 
3:   samples  $\leftarrow \emptyset$ 

4:   while SIZE(samples)  $\leq$  target_size do
5:     samples  $\leftarrow$  SOBOl(dimension, sample_size)
6:     samples  $\leftarrow$  FILTER(samples,  $\leq$  max_weight)
7:     sample_size  $\leftarrow 2 * \text{sample\_size}$ 
8:   end while

9:   return samples
10: end function

11: function GPR(budget, dimension, starting_size, sample_size, max_weight)
12:   starting_sample  $\leftarrow$  FILTER_SOBOl(dimension, starting_size, max_weight)
13:   training_set  $\leftarrow$  MEASURE(starting_sample)
14:   measurements  $\leftarrow$  SIZE(training_set)
15:   EI_backlog  $\leftarrow \emptyset$ 

16:   while measurements  $\leq$  budget do
17:     model  $\leftarrow$  FitGP(training_set)

18:     new_sample  $\leftarrow$  FILTER_SOBOl(dimension, sample_size, max_weight)
19:     candidates  $\leftarrow$  COMPUTEI(model, new_sample)
20:     neighbors  $\leftarrow$  COMPUTEI(model, NEIGHBORHOOD(candidates)))

21:     new_measurement  $\leftarrow$  MaxEI(candidates  $\cup$  neighbors  $\cup$  EI_backlog)
22:     EI_backlog  $\leftarrow (\text{EI\_backlog} \cup \text{candidates} \cup \text{neighbors}) \setminus \text{new\_measurement}$ 

23:     training_set  $\leftarrow$  training_set  $\cup$  MEASURE(new_measurement)
24:     measurements  $\leftarrow$  SIZE(training_set)
25:   end while

26:   return MaxTop5(training_set)
27: end function

```

Figure 14.4: Pseudocode implementation of the GPR method using Expected Improvement (EI), and of the constrained sampler using Sobol sequences. In our experiments, variables were $\text{budget} \leftarrow 245$, $\text{dimension} \leftarrow 54$, $\text{starting_size} \leftarrow 64$, $\text{sample_size} \leftarrow 3240$, and $\text{max_weight} \leftarrow 10 \text{ MB}$

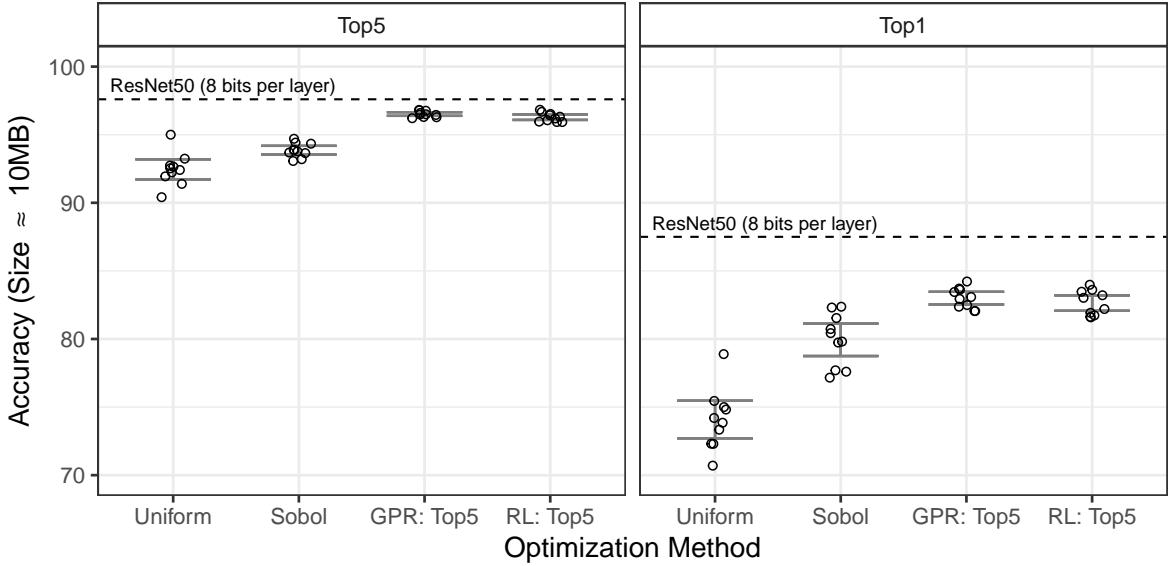


Figure 14.5: 10 repetitions of the 4 methods we compared, using a budget of 245 measurements, with 95% confidence intervals for estimates of the mean

14.5 Performance Evaluation and Comparison

This section presents our evaluation of the performance of our implementation of Gaussian Process Regression with Expected Improvement (GPR) and of the baseline Reinforcement Learning (RL) [148] on the optimization of quantization policies for *ResNet50*, targeting the *ImageNet* dataset. We compare the performance of GPR and RL with a uniform random sampling method and with a low-discrepancy space-filling sampler using Sobol sequences.

The simple samplers performed interestingly well in comparison with more complex approaches, a pattern we have already observed in many of the studies presented in this thesis. This time, the more complex methods were able to improve upon the results of the simpler approaches. Under the same budget, our GPR implementation achieved accuracy statistically indistinguishable from the original RL implementation. In comparison with the RL algorithm, we argue that our approach produces more consistent quantization policies between repetitions, and that it is more extensible and perhaps more robust to changes in the objective function.

14.5.1 Optimizing *Top5* for Constrained Total Weight

Figure 14.5 compares the best quantization policies found by each optimization method across 10 repetitions, with a fixed budget. The total *ResNet50* weight size was constrained to be no larger than 10 MB. Comparing the 95% confidence intervals for the estimates of the mean best configuration for each method we see that GPR and RL have similar results, but that the Sobol sampler produces policies with good accuracy. Figure 14.6 shows the accuracy

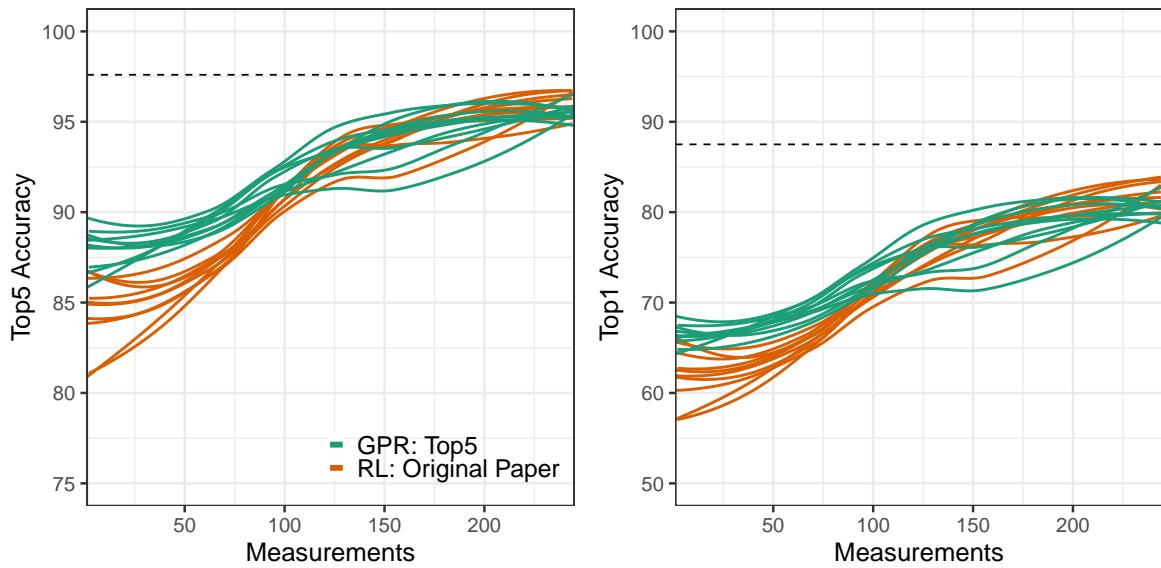


Figure 14.6: Best quantization policy found by each method during optimization, across 10 repetitions

metrics of the best configurations found at each iteration, for the 10 repetitions of GPR and RL.

Figures 14.7 and 14.8 present four columns with more detailed looks into the quantization policies explored by GPR and RL and the corresponding network configuration, in the first column, total weight size, in the second column, and *Top5* and *Top1* accuracy metrics in the last two columns.

The first columns of both Figures show stacked *ResNet50* layers in the *y*-axis, in similar fashion to what is presented in Figure 14.2. The *x*-axis, as in all columns, shows each of the 245 measurements. Each square represents a layer at a given time, and is colored with a gradient between red for 1-bit and green for 8-bit quantization. Initial blocks of measurements for all experiments consist of 64 *warm-up* measurements, generating visually distinct starting regions in all images.

The three horizontal dashed lines in the second column, from top to bottom, mark the size of *ResNet50* with all layers quantized with 8, 2, and 1 bit for all layers. Likewise for the third and fourth columns, the horizontal line marks the accuracy of the 8-bit quantized network. The vertical lines in all columns mark the iteration that found the point with the highest *Top5* accuracy, for each repetition.

Figure 14.7 shows that the policies arrived at by RL have different quantization distributions. Experiments 3, 5, 6, and 10 have clearly defined red bands in similar regions, corresponding to layers that were quantized with fewer bits. Fainter bands in the same places can be seen in the other experiments, but we can identify other features as well, such as green bands or more uniform orange policies. The original RL implementation used random

sampling and round-robin decrease of weights to obtain the warm-up set, while our GPR implementation used Sobol sequences and filtering, resulting in different-looking starting stages.

Figure 14.8 shows clear red bands in the same top layers, and less variable policy choices than RL across experiments. We see regions that abruptly change from darker to light greens in experiments 2, 4, 9, and 10, which could be due to the algorithm using a configuration on the EI backlog and picking a new region of the search space to explore.

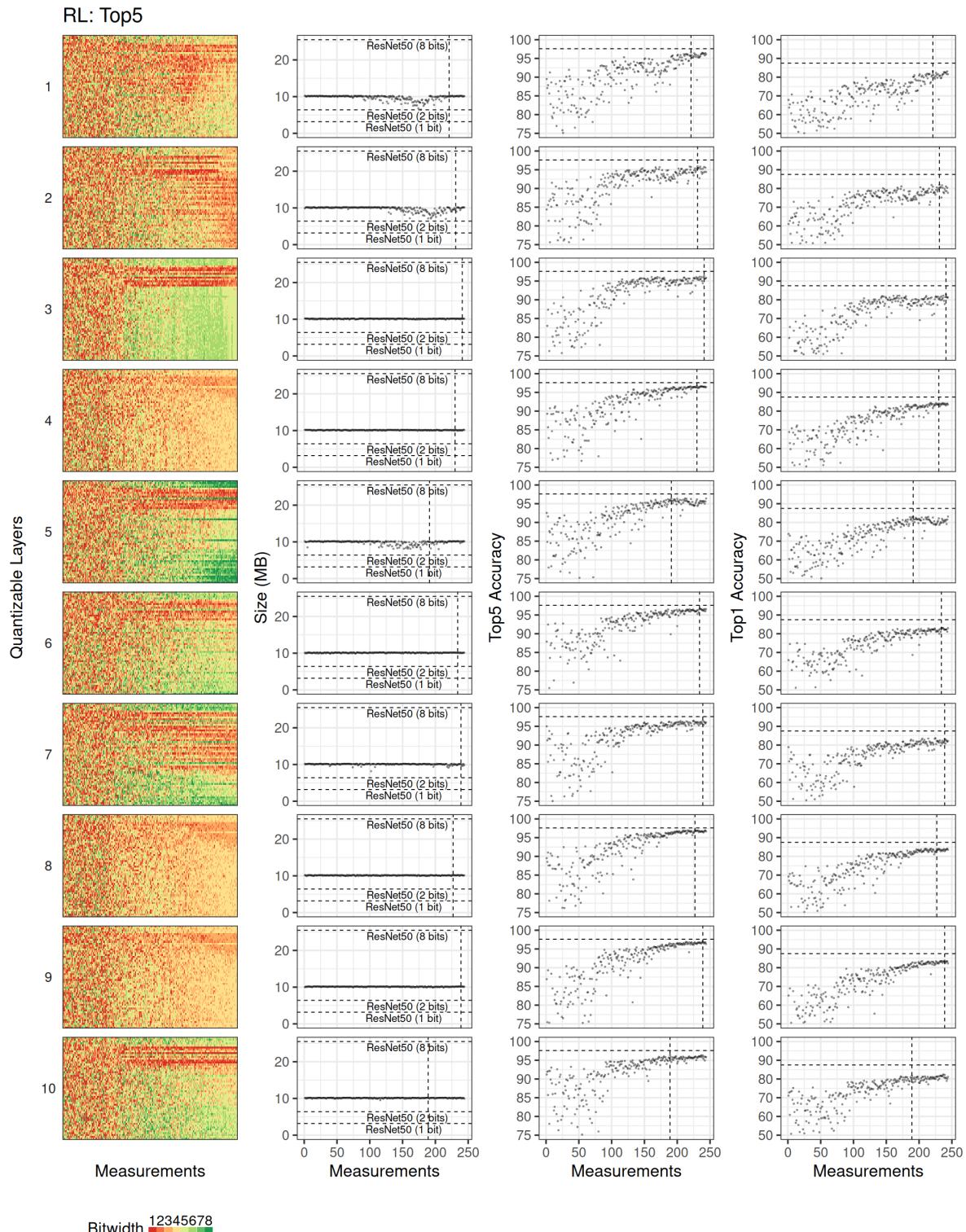


Figure 14.7: Results with the baseline Reinforcement Learning method

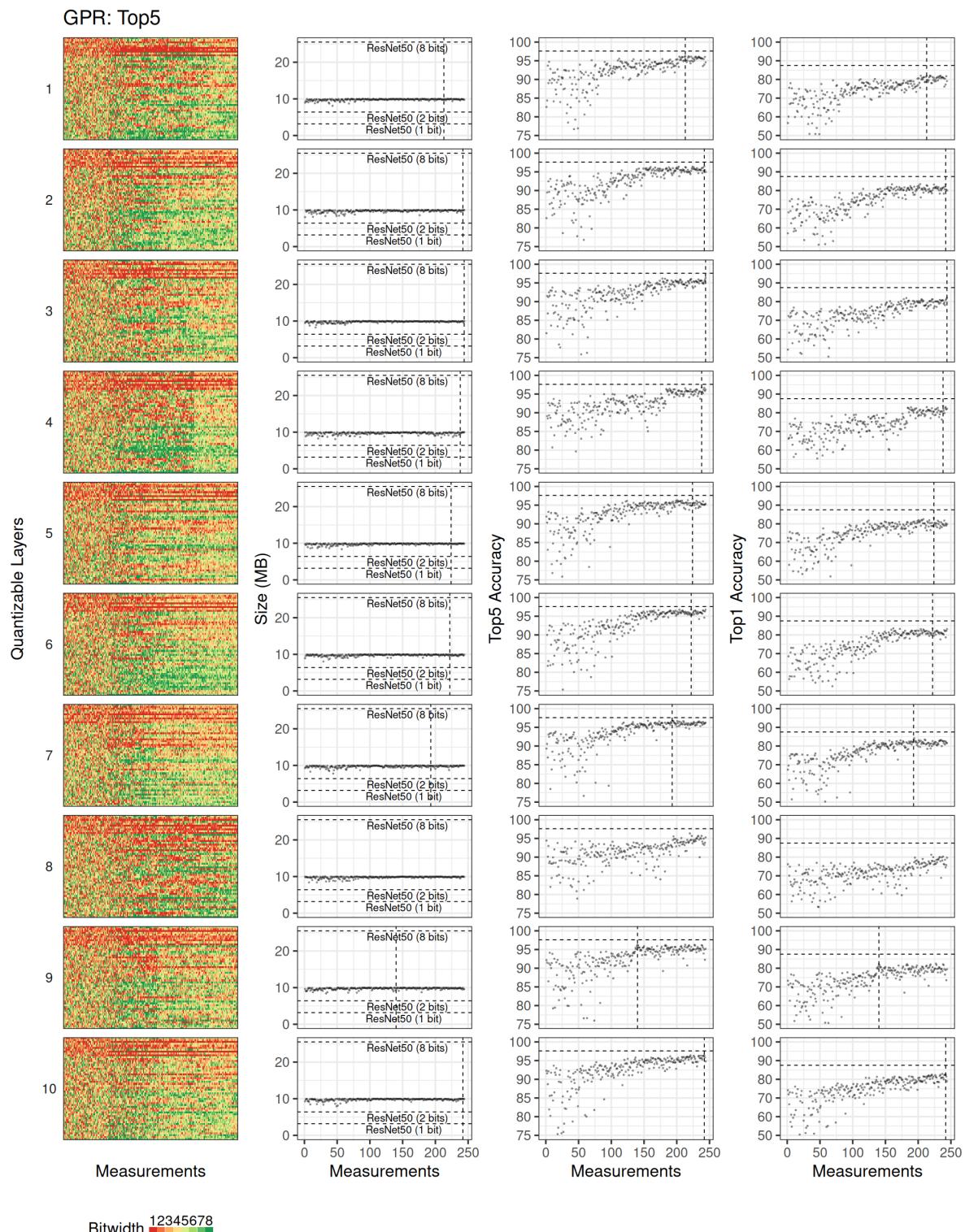


Figure 14.8: Results with Gaussian Process Regression

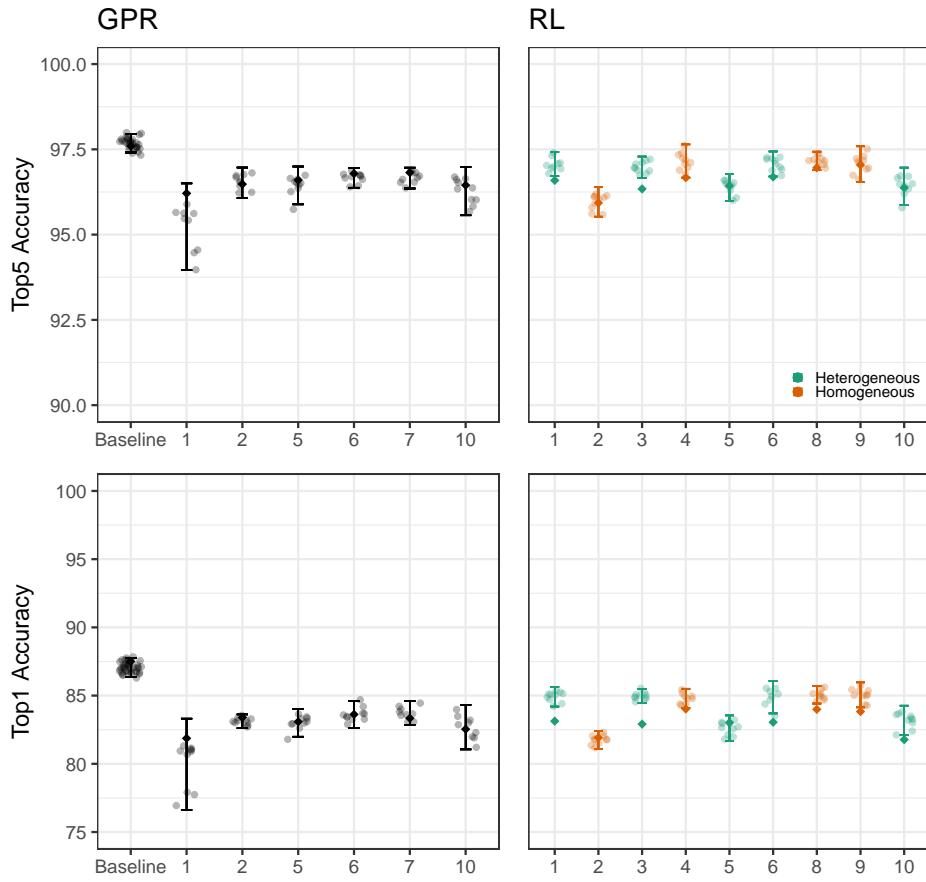


Figure 14.9: Running optimized quantization policies with new *ImageNet* samples, with 95% confidence intervals for the mean estimate. The baseline is the uniform 8-bit *ResNet50* quantization

Comparing Policy Variability on Different Samples

Figure 14.9 shows the result of running some of the policies produced by GPR, and all produced by RL, in new random *ImageNet* samples. We wanted to check if the “homogeneous” RL experiments, that is, the ones without clear red bands, had a different behavior when pitted against new datasets. We did not see any difference between RL experiments, but we were surprised to find that RL experiments were consistently better, and outside of the confidence intervals, for new datasets, although differences are not practically significant. Experiment numbering is the same from Figures 14.7 and 14.8.

14.6 Measuring Performance and Detecting Significance for GPR

In this section we analyze the performance of our GPR implementation, discuss our attempts to detect the significance of the quantization impact of each layer on *Top5* accuracy, and present perspectives on future explorations of GPR for CNN quantization.

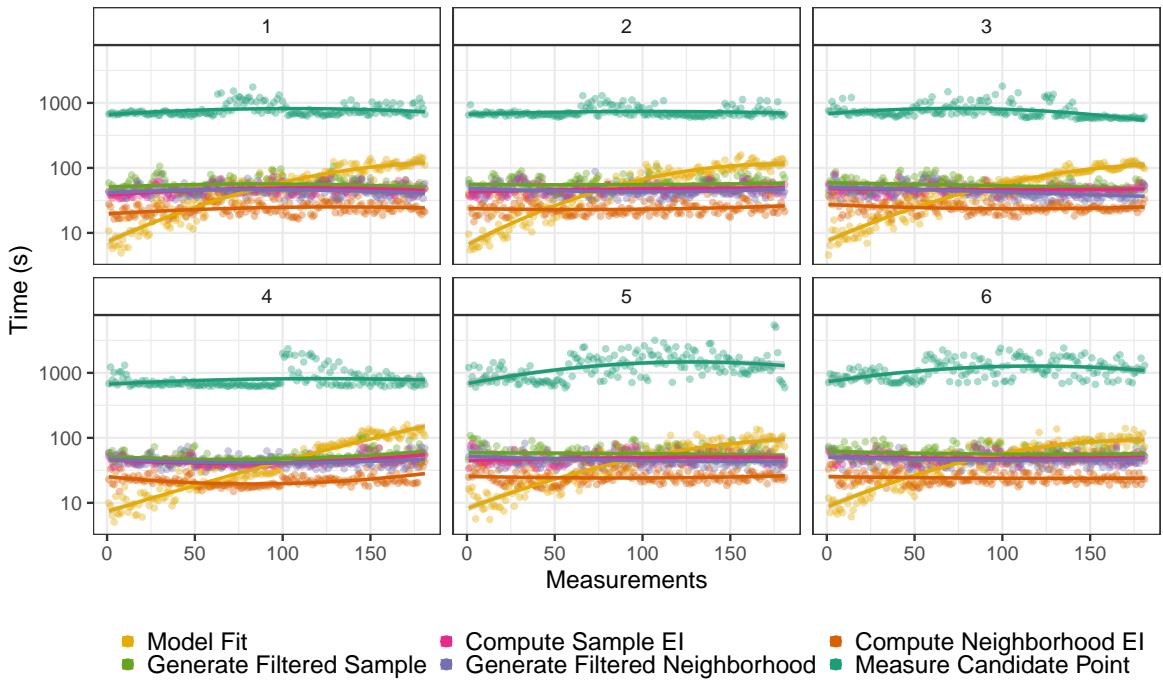


Figure 14.10: Execution time of the different steps of our GPR implementation. The initial sample measurement is not shown, and the y -axis is logarithmic

14.6.1 Performance of the GPR Method

Figure 14.10 shows the performance of six executions of the GPR algorithm described in Figure 14.4, broken down by step. Most steps of the algorithm have fixed computational cost, since Sobol samples and neighborhoods have fixed sizes across the optimization process, but fitting the model takes more time as the size of training data grows. The accuracy measurement time for a quantization policy on an *ImageNet* sample depends on the policy. We did not explore the reasons why this happens, but we suspect that training a network using operations with less precision is faster.

14.6.2 Detecting Significance with Sobol Indices

We used the *sensitivity R* package [164] to compute Sobol indices, using Gaussian Process surrogate models, for the impact on *Top5* of quantizing each *ResNet50* layer. Figure 14.11 shows the indices computed for the ensemble of Sobol samples we collected, in the left panel, and for the Sobol samples plus the search space exploration performed by GPR, in the right panel.

It is expected that the first and last layers have the most impact, since they determine the precision of all computations that follow, and the precision of the final output, respectively. The confidence intervals for the means of the other indices in the Sobol sample all cross zero but we know that we can impact accuracy by changing them, so the interpretation is unclear.

Likewise, it is unclear how to interpret that all indices seem to have an equally large impact when the explored search spaces are added on the right panel, since we see clear bands of strongly quantized layers on the GPR results.

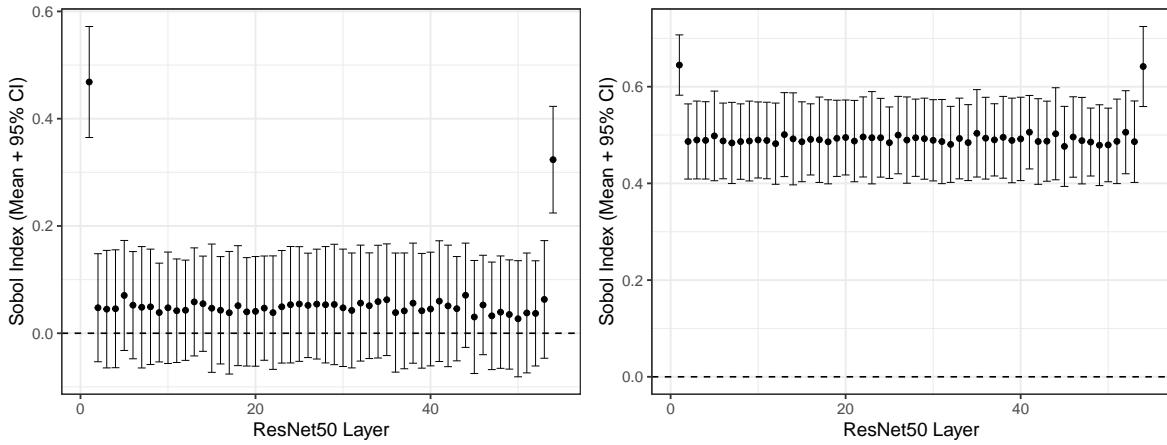


Figure 14.11: Sobol indices computed for the impact on *Top5* accuracy of the quantization of each *ResNet50* layer using only Sobol samples, in the left panel, and Sobol samples plus the search spaces explored by GPR, in the right panel

14.6.3 Perspectives for Further Exploration

Running GPR with Data from Past Experiments

An interesting feature of Gaussian Process Regression is that, due to its stability, we can leverage past executions by keeping only a *.csv* file with the previous measurements. It is more expensive to relaunch GPR with each added past exploration, since we need to fit a model for more data, but the way in which new experiments are picked does not change.

The top row of Figure 14.12 shows a GPR instance that was launched with the ensemble of data from all past GPR experiments. On the last two columns we can see each peak in accuracy, corresponding to the end of each experiment from Figure 14.8. The bottom row shows the same GPR instance, but focuses only on the experiments picked after fitting the model with the data ensemble.

We see that all newly picked configurations have around 95% *Top5* accuracy, and that accuracy seems to continue to increase. A long-running optimizer could be continuously fed with new experiments as tests are made, even during prototyping, and could be used to produce configurations similar to the best known ones, which could help understanding the impact of each controllable parameter.

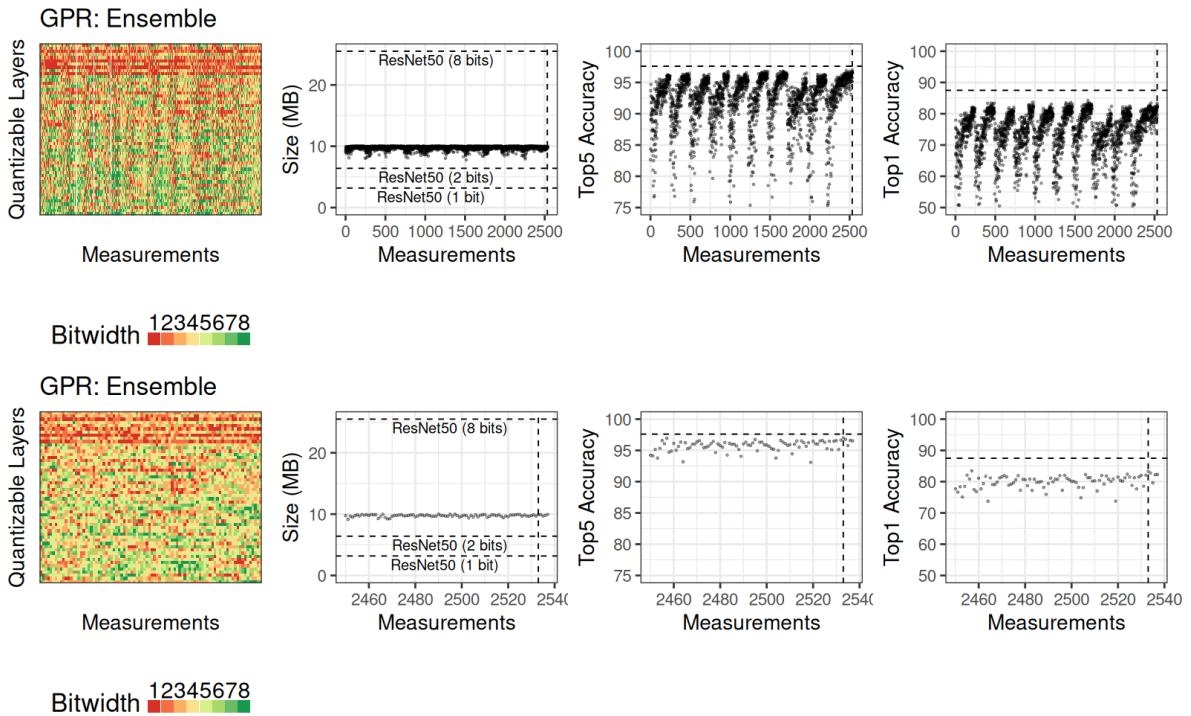


Figure 14.12: Running GPR with the gathered ensemble of GPR experiments. The bottom row shows only the new measurements performed after fitting a Gaussian Process with the ensemble of experiments, and the top row shows all the data that were used

Optimizing Accuracy and Weight

We ran preliminary experiments with adding the total weight size as a component of the performance metric for GPR and RL, but this made it unclear how to compare the GPR results to the original RL implementation, because to add the new metric we removed the round-robin weight decrease used in the original paper.

Figures 14.14 and 14.13 show the results for GPR and RL, respectively. We see that this expanded problem is a much harder task for both methods, but RL is especially impacted by it. This might be a good reason for the original paper's usage of the round-robin procedure for weight decrease. Experiments 2, 3, and 4 with the RL method show quirky features that we are not able to explain, but we hypothesize that some emerging restarting strategy has taken place, especially for experiment 4. Experiment 2 is the only one that seems to behave as expected.

The experiments with GPR were much more stable. We show two of them in Figure 14.14, and we can see red bands that resemble the ones from Figure 14.8. Although this modified version does not achieve the same *Top5* accuracy, it does find interesting configurations producing quantization policies equivalent to using less than 2 bits per layer, without letting *Top5* go below 80%. This hints that policies exist that generate much smaller networks with a manageable compromise in accuracy.

14.6. Measuring Performance and Detecting Significance for GPR

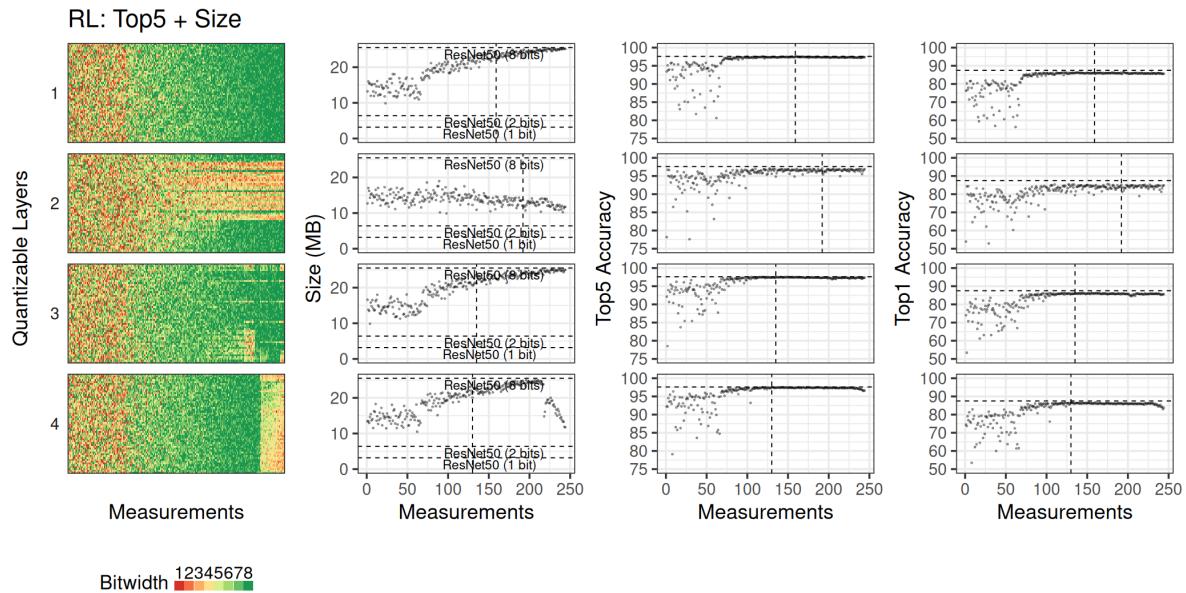


Figure 14.13: Optimizing accuracy and weight with the modified baseline Reinforcement Learning method

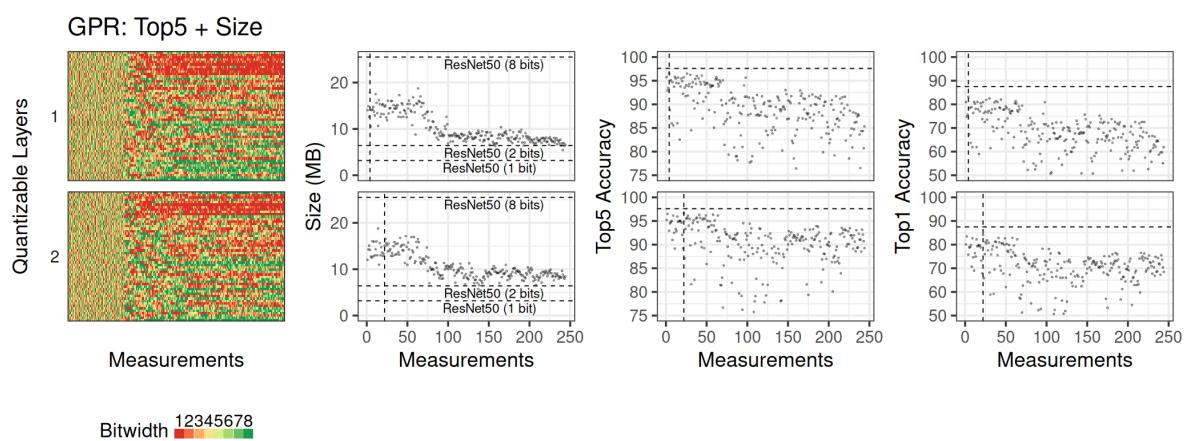


Figure 14.14: Optimizing accuracy and weight with Gaussian Process Regression with EI.

14.7 Summary

This chapter presented an application of Gaussian Process Regression with Expected Improvement to the problem of finding the smallest mixed-precision quantization policies for the layers of a neural network, subject to keeping the accuracy comparable to the 8-bit per layer network version, and the total weight size below 10 MB.

We compared our GPR approach with a baseline Reinforcement Learning method, and with uniform and space-filling samplers. The sampling methods performed relatively well, but worse than the more complex approaches. Our implementation produced policies with comparable accuracy to the Reinforcement Learning approach, with the advantage of more consistent explorations.

Although our attempts to estimate Sobol indices were inconclusive, we believe that the choices made by the GPR method are overall easier to interpret and understand, in comparison with the results produced by the RL.

The GPR method was also capable of leveraging past executions to produce policies with similar accuracy. Additionally, we argue that GPR was more robust to the addition of the total weight size to the objective function. Accurately estimating Sobol indices, performing multi-objective GPR optimization, and learning more about high-performing configurations using past data are possibly interesting paths that we could point to during this chapter, but that we do not explore in this thesis.

Conclusion

Chapter 15

Conclusion

Bibliography

- [1] J. Ansel, S. Kamil, K. Veeramachaneni, J. Ragan-Kelley, J. Bosboom, U.-M. O'Reilly, and S. Amarasinghe, "Opentuner: An extensible framework for program autotuning," in *Proceedings of the 23rd international conference on Parallel architectures and compilation*. ACM, 2014, pp. 303–316.
- [2] P. Bruel, M. Amarís, and A. Goldman, "Autotuning CUDA compiler parameters for heterogeneous applications using the OpenTuner framework," *Concurrency and Computation: Practice and Experience*, pp. e3973–n/a, 2017.
- [3] P. Bruel, A. Goldman, S. R. Chalamalasetti, and D. Milojicic, "Autotuning High-Level Synthesis for FPGAs Using OpenTuner and LegUp," in *International Conference on Reconfigurable Computing and FPGAs (ReConFig)*, 2017.
- [4] P. Bruel, S. Quinito Masnada, B. Videau, A. Legrand, J.-M. Vincent, and A. Goldman, "Autotuning under Tight Budget Constraints: A Transparent Design of Experiments Approach," in *The 19th Annual IEEE/ACM International Symposium in Cluster, Cloud, and Grid Computing (CCGrid 2019)*. IEEE/ACM, 2019.
- [5] K. Seymour, H. You, and J. Dongarra, "A comparison of search heuristics for empirical code optimization." in *CLUSTER*, 2008, pp. 421–429.
- [6] R. A. Fisher, *The design of experiments*. Oliver And Boyd; Edinburgh; London, 1937.
- [7] P. N. D. Bukh, *The art of computer systems performance analysis, techniques for experimental design, measurement, simulation and modeling*. JSTOR, 1992.
- [8] D. C. Montgomery, *Design and analysis of experiments*. John wiley & sons, 2017.
- [9] G. E. Box, J. S. Hunter, and W. G. Hunter, *Statistics for experimenters: design, innovation, and discovery*. Wiley-Interscience New York, 2005, vol. 2.
- [10] C. K. Williams and C. E. Rasmussen, *Gaussian processes for machine learning*. MIT press Cambridge, MA, 2006, vol. 2, no. 3.
- [11] P. E. Ceruzzi, E. Paul *et al.*, *A history of modern computing*. MIT press, 2003.

Bibliography

- [12] O. R. N. Laboratory, "Summit – Oak Ridge Leadership Computing Facility," <https://www.olcf.ornl.gov/olcf-resources/compute-systems/summit>, 2020, [Online; accessed 08-Jun-2020].
- [13] Wikipedia, "Transistor count," https://en.wikipedia.org/wiki/Transistor_count, 2020, [Online; accessed 08-Jun-2020].
- [14] ——, "Microprocessor Chronology," https://en.wikipedia.org/wiki/Microprocessor_chronology, 2020, [Online; accessed 08-Jun-2020].
- [15] G. E. Moore *et al.*, "Cramming more components onto integrated circuits," 1965.
- [16] R. H. Dennard, F. H. Gaensslen, H.-N. Yu, V. L. Rideout, E. Bassous, and A. R. LeBlanc, "Design of ion-implanted mosfet's with very small physical dimensions," *IEEE Journal of solid-state circuits*, vol. 9, no. 5, pp. 256–268, 1974.
- [17] D. J. Frank, R. H. Dennard, E. Nowak, P. M. Solomon, Y. Taur, and H.-S. P. Wong, "Device scaling limits of si mosfets and their application dependencies," *Proceedings of the IEEE*, vol. 89, no. 3, pp. 259–288, 2001.
- [18] H. Esmaeilzadeh, E. Blem, R. S. Amant, K. Sankaralingam, and D. Burger, "Dark silicon and the end of multicore scaling," in *2011 38th Annual international symposium on computer architecture (ISCA)*. IEEE, 2011, pp. 365–376.
- [19] S. Venkataramani, S. T. Chakradhar, K. Roy, and A. Raghunathan, "Approximate computing and the quest for computing efficiency," in *2015 52nd ACM/EDAC/IEEE Design Automation Conference (DAC)*. IEEE, 2015, pp. 1–6.
- [20] H.-Y. Cheng, J. Zhan, J. Zhao, Y. Xie, J. Sampson, and M. J. Irwin, "Core vs. uncore: The heart of darkness," in *2015 52nd ACM/EDAC/IEEE Design Automation Conference (DAC)*. IEEE, 2015, pp. 1–6.
- [21] J. Henkel, H. Khdr, S. Pagani, and M. Shafique, "New trends in dark silicon," in *2015 52nd ACM/EDAC/IEEE Design Automation Conference (Dac)*. IEEE, 2015, pp. 1–6.
- [22] Top500, "Top500 List," <https://www.top500.org/>, 2020, [Online; accessed 08-Jun-2020].
- [23] J. J. Dongarra, P. Luszczek, and A. Petitet, "The LINPACK benchmark: past, present and future," *Concurrency and Computation: practice and experience*, vol. 15, no. 9, pp. 803–820, 2003.
- [24] J. J. Dongarra, J. Du Croz, S. Hammarling, and I. S. Duff, "A set of level 3 basic linear algebra subprograms," *ACM Transactions on Mathematical Software (TOMS)*, vol. 16, no. 1, pp. 1–17, 1990.
- [25] C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh, "Basic linear algebra subprograms for fortran usage," *ACM Transactions on Mathematical Software (TOMS)*, vol. 5, no. 3, pp. 308–323, 1979.

- [26] B. Videau, K. Pouget, L. Genovese, T. Deutsch, D. Komatitsch, F. Despres, and J.-F. Méhaut, "BOAST: A metaprogramming framework to produce portable and efficient computing kernels for hpc applications," *The International Journal of High Performance Computing Applications*, p. 1094342017718068, 2017.
- [27] A. Hartono, B. Norris, and P. Sadayappan, "Annotation-based empirical performance tuning using Orio," in *Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on.* IEEE, 2009, pp. 1–11.
- [28] J. Ansel, C. Chan, Y. L. Wong, M. Olszewski, Q. Zhao, A. Edelman, and S. Amarasinghe, "Petabricks: a language and compiler for algorithmic choice," *ACM Sigplan Notices*, vol. 44, no. 6, pp. 38–49, 2009.
- [29] P. Balaprakash, S. M. Wild, and B. Norris, "SPAPT: Search problems in automatic performance tuning," *Procedia Computer Science*, vol. 9, pp. 1959–1968, 2012.
- [30] J.-H. Byun, R. Lin, K. A. Yelick, and J. Demmel, "Autotuning sparse matrix-vector multiplication for multicore," *EECS, UC Berkeley, Tech. Rep*, 2012.
- [31] F. Petrovič, D. Střelák, J. Hozzová, J. Ol'ha, R. Trembecký, S. Benkner, and J. Filipovič, "A benchmark set of highly-efficient cuda and opencl kernels and its dynamic autotuning with kernel tuning toolkit," *Future Generation Computer Systems*, 2020.
- [32] P. Balaprakash, M. Salim, T. Uram, V. Vishwanath, and S. Wild, "Deephypert: Asynchronous hyperparameter search for deep neural networks," in *2018 IEEE 25th International Conference on High Performance Computing (HiPC).* IEEE, 2018, pp. 42–51.
- [33] P. Bruel, M. Amarís, and A. Goldman, "Autotuning gpu compiler parameters using opentuner," in *Proceedings of the WSCAD (Simpósio em Sistemas Computacionais de Alto Desempenho)*, 2015.
- [34] A. Mametjanov, P. Balaprakash, C. Choudary, P. D. Hovland, S. M. Wild, and G. Sabin, "Autotuning fpga design parameters for performance and power," in *Field-Programmable Custom Computing Machines (FCCM), 2015 IEEE 23rd Annual International Symposium on.* IEEE, 2015, pp. 84–91.
- [35] A. Abdelfattah, A. Haidar, S. Tomov, and J. Dongarra, "Performance, design, and autotuning of batched gemm for gpus," in *International Conference on High Performance Computing.* Springer, 2016, pp. 21–38.
- [36] C. Xu, G. Liu, R. Zhao, S. Yang, G. Luo, and Z. Zhang, "A parallel bandit-based approach for autotuning fpga compilation," in *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays.* ACM, 2017, pp. 157–166.
- [37] A. Tiwari, C. Chen, J. Chame, M. Hall, and J. K. Hollingsworth, "A scalable auto-tuning framework for compiler optimization," in *Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on.* IEEE, 2009, pp. 1–12.

Bibliography

- [38] F. Hutter, H. H. Hoos, K. Leyton-Brown, and T. Stützle, "ParamILS: an automatic algorithm configuration framework," *Journal of Artificial Intelligence Research*, vol. 36, no. 1, pp. 267–306, 2009.
- [39] Y. Chu, C. Luo, H. H. Hoos, Q. Lin, and H. You, "Improving the performance of stochastic local search for maximum vertex weight clique problem using programming by optimization," *arXiv*, pp. arXiv–2002, 2020.
- [40] I. Tuzov, D. de Andrés, and J.-C. Ruiz, "Tuning synthesis flags to optimize implementation goals: Performance and robustness of the leon3 processor as a case study," *Journal of Parallel and Distributed Computing*, vol. 112, pp. 84–96, 2018.
- [41] M. M. Ziegler, H.-Y. Liu, G. Gristede, B. Owens, R. Nigaglioni, J. Kwon, and L. P. Carloni, "Syntunsys: A synthesis parameter autotuning system for optimizing high-performance processors," in *Machine Learning in VLSI Computer-Aided Design*. Springer, 2019, pp. 539–570.
- [42] M. Gerndt, S. Benkner, E. César, C. Navarrete, E. Bajrovic, J. Dokulil, C. Guillén, R. Mijakovic, and A. Sikora, "A multi-aspect online tuning framework for hpc applications," *Software Quality Journal*, vol. 26, no. 3, pp. 1063–1096, 2018.
- [43] J. Kwon, M. M. Ziegler, and L. P. Carloni, "A learning-based recommender system for autotuning design flows of industrial high-performance processors," in *2019 56th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2019, pp. 1–6.
- [44] T. Wang, N. Jain, D. Beckingsale, D. Boehme, F. Mueller, and T. Gamblin, "Funcytuner: Auto-tuning scientific applications with per-loop compilation," in *Proceedings of the 48th International Conference on Parallel Processing*, 2019, pp. 1–10.
- [45] J. Ol'ha, J. Hozzová, J. Fousek, and J. Filipovič, "Exploiting historical data: pruning autotuning spaces and estimating the number of tuning steps," in *European Conference on Parallel Processing*. Springer, 2019, pp. 295–307.
- [46] J. Bilmes, K. Asanovic, C.-W. Chin, and J. Demmel, "Optimizing matrix multiply using PHiPAC: a portable, high-performance, ansi c coding methodology," in *Proceedings of International Conference on Supercomputing, Vienna, Austria*, 1997.
- [47] J. J. Dongarra and C. R. Whaley, "Automatically tuned linear algebra software (ATLAS)," *Proceedings of SC*, vol. 98, 1998.
- [48] M. Frigo and S. G. Johnson, "FFTW: An adaptive software architecture for the fft," in *Acoustics, Speech and Signal Processing, 1998. Proceedings of the 1998 IEEE International Conference on*, vol. 3. IEEE, 1998, pp. 1381–1384.
- [49] C. Tăpuş, I.-H. Chung, J. K. Hollingsworth *et al.*, "Active harmony: Towards automated performance tuning," in *Proceedings of the 2002 ACM/IEEE conference on Supercomputing*. IEEE Computer Society Press, 2002, pp. 1–11.

- [50] R. Vuduc, J. W. Demmel, and K. A. Yelick, "OSKI: A library of automatically tuned sparse matrix kernels," in *Journal of Physics: Conference Series*, vol. 16, no. 1. IOP Publishing, 2005, p. 521.
- [51] G. Fursin, Y. Kashnikov, A. W. Memon, Z. Chamski, O. Temam, M. Namolaru, E. Yom-Tov, B. Mendelson, A. Zaks, E. Courtois *et al.*, "Milepost gcc: Machine learning enabled self-tuning compiler," *International journal of parallel programming*, vol. 39, no. 3, pp. 296–327, 2011.
- [52] H. Jordan, P. Thoman, J. J. Durillo, S. Pellegrini, P. Gschwandtner, T. Fahringer, and H. Moritsch, "A multi-objective auto-tuning framework for parallel codes," in *High Performance Computing, Networking, Storage and Analysis (SC), 2012 International Conference for*. IEEE, 2012, pp. 1–12.
- [53] D. G. Spampinato and M. Püschel, "A basic linear algebra compiler," in *Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization*. ACM, 2014, p. 23.
- [54] C. Audet, K.-C. Dang, and D. Orban, "Optimization of algorithms with opal," *Mathematical Programming Computation*, vol. 6, no. 3, pp. 233–254, 2014.
- [55] C. Nugteren and V. Codreanu, "Cltune: A generic auto-tuner for opencl kernels," in *Embedded Multicore/Many-core Systems-on-Chip (MCSoC), 2015 IEEE 9th International Symposium on*. IEEE, 2015, pp. 195–202.
- [56] J. Guerreiro, A. Ilic, N. Roma, and P. Tomás, "Multi-kernel auto-tuning on gpus: Performance and energy-aware optimization," in *Parallel, Distributed and Network-Based Processing (PDP), 2015 23rd Euromicro International Conference on*. IEEE, 2015, pp. 438–445.
- [57] G. Fursin, A. Memon, C. Guillou, and A. Lokhmotov, "Collective mind, part ii: Towards performance-and cost-aware software engineering as a natural science," *arXiv preprint arXiv:1506.06256*, 2015.
- [58] L.-W. Chang, I. El Hajj, C. Rodrigues, J. Gómez-Luna, and W.-m. Hwu, "Efficient kernel synthesis for performance portable programming," in *Microarchitecture (MICRO), 2016 49th Annual IEEE/ACM International Symposium on*. IEEE, 2016.
- [59] C. G. Coelho, C. G. Abreu, R. M. Ramos, A. H. Mendes, G. Teodoro, and C. G. Ralha, "Mase-bdi: agent-based simulator for environmental land change with efficient and parallel auto-tuning," *Applied Intelligence*, vol. 45, no. 3, pp. 904–922, 2016.
- [60] D. Beckingsale, O. Pearce, I. Laguna, and T. Gamblin, "Apollo: Reusable models for fast, dynamic tuning of input-dependent code," in *The 31th IEEE International Parallel and Distributed Processing Symposium*, 2017.

Bibliography

- [61] M. J. Kochenderfer and T. A. Wheeler, *Algorithms for optimization*. Mit Press, 2019.
- [62] A. J. Dobson and A. G. Barnett, *An introduction to generalized linear models*. CRC press, 2018.
- [63] A. Agresti, *Foundations of linear and generalized linear models*. John Wiley & Sons, 2015.
- [64] G. Seber, *The Linear Model and Hypothesis: A General Unifying Theory*. Springer, 2015.
- [65] C. M. Bishop, *Pattern recognition and machine learning*. springer, 2006.
- [66] Rousant, Olivier and Ginsbourger, David and Deville, Yves and Chevalier, Clement and Richet, Yann, "DiceKriging." [Online]. Available: <https://cran.r-project.org/web/packages/DiceKriging/index.html>
- [67] O. Rousant, D. Ginsbourger, and Y. Deville, "DiceKriging, DiceOptim: Two R Packages for the Analysis of Computer Experiments by Kriging-Based Metamodeling and Optimization," vol. 51, no. 1, 2012, pp. 1–55.
- [68] I. M. Sobol', "Sensitivity estimates for nonlinear mathematical models," *Math. Model. Comput. Exp.*, vol. 1, no. 4, pp. 407–414, 1993.
- [69] A. Saltelli, M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli, M. Saisana, and S. Tarantola, *Global sensitivity analysis: the primer*. John Wiley & Sons, 2008.
- [70] A. Saltelli, P. Annoni, I. Azzini, F. Campolongo, M. Ratto, and S. Tarantola, "Variance based sensitivity analysis of model output. design and estimator for the total sensitivity index," *Computer physics communications*, vol. 181, no. 2, pp. 259–270, 2010.
- [71] R. L. Plackett and J. P. Burman, "The design of optimum multifactorial experiments," *Biometrika*, vol. 33, no. 4, pp. 305–325, 1946.
- [72] U. Grömping, "R package FrF2 for creating and analyzing fractional factorial 2-level designs," *Journal of Statistical Software*, vol. 56, no. 1, pp. 1–56, 2014. [Online]. Available: <http://www.jstatsoft.org/v56/i01/>
- [73] R Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2018. [Online]. Available: <https://www.R-project.org/>
- [74] A. Atkinson, A. Donev, and R. Tobias, *Optimum experimental designs, with SAS*. Oxford University Press, 2007, vol. 34.
- [75] V. V. Fedorov, *Theory of optimal experiments*. Elsevier, 1972.
- [76] B. Wheeler, *AlgDesign: Algorithmic Experimental Design*, 2014, R package version 1.1-7.3. [Online]. Available: <https://CRAN.R-project.org/package=AlgDesign>

- [77] I. M. Sobol', "On the distribution of points in a cube and the approximate evaluation of integrals," *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, vol. 7, no. 4, pp. 784–802, 1967.
- [78] J. H. Halton, "Algorithm 247: Radical-inverse quasi-random point sequence," *Communications of the ACM*, vol. 7, no. 12, pp. 701–702, 1964.
- [79] D. Christophe and S. Petr, *randtoolbox: Generating and Testing Random Numbers*, 2020, r package version 1.30.1.
- [80] D. Dupuy, C. Helbert, and J. Franco, "DiceDesign and DiceEval: Two R packages for design and analysis of computer experiments," *Journal of Statistical Software*, vol. 65, no. 11, pp. 1–38, 2015. [Online]. Available: <https://www.jstatsoft.org/v65/i11/>
- [81] R. Carnell, *lhs: Latin Hypercube Samples*, 2018, R package version 0.16. [Online]. Available: <https://CRAN.R-project.org/package=lhs>
- [82] R. E. Paley, "On orthogonal matrices," *Journal of Mathematics and Physics*, vol. 12, no. 1-4, pp. 311–320, 1933.
- [83] P. Balaprakash, S. M. Wild, and P. D. Hovland, "Can search algorithms save large-scale automatic performance tuning?" in *ICCS*, 2011, pp. 2136–2145.
- [84] ——, "An experimental study of global and local search algorithms in empirical performance tuning," in *International Conference on High Performance Computing for Computational Science*. Springer, 2012, pp. 261–269.
- [85] P. Balaprakash, A. Tiwari, S. M. Wild, and P. D. Hovland, "AutoMOMML: Automatic Multi-objective Modeling with Machine Learning," in *High Performance Computing: 31st International Conference, ISC High Performance 2016, Frankfurt, Germany, June 19-23, 2016, Proceedings*, M. J. Kunkel, P. Balaji, and J. Dongarra, Eds. Springer International Publishing, 2016, pp. 219–239.
- [86] M. Parsa, A. Ankit, A. Ziabari, and K. Roy, "PABO: Pseudo Agent-Based Multi-Objective Bayesian Hyperparameter Optimization for Efficient Neural Accelerator Design," in *2019 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 2019, pp. 1–8.
- [87] S. Che, M. Boyer, J. Meng, D. Tarjan, J. W. Sheaffer, S.-H. Lee, and K. Skadron, "Rodinia: A benchmark suite for heterogeneous computing," in *Workload Characterization, 2009. IISWC 2009. IEEE International Symposium on*. IEEE, 2009, pp. 44–54.
- [88] Y. Zhang and J. D. Owens, "A quantitative performance analysis model for gpu architectures," in *High Performance Computer Architecture (HPCA), 2011 IEEE 17th International Symposium on*. IEEE, 2011, pp. 382–393.

Bibliography

- [89] M. Amaris, D. Cordeiro, A. Goldman, and R. Y. de Camargo, "A simple bsp-based model to predict execution time in gpu applications," in *High Performance Computing (HiPC), 2015 IEEE 22nd International Conference on.* IEEE, 2015, pp. 285–294.
- [90] T. T. Dao, J. Kim, S. Seo, B. Egger, and J. Lee, "A performance model for gpus with caches," *IEEE Transactions on Parallel and Distributed Systems*, vol. 26, no. 7, pp. 1800–1813, 2015.
- [91] J. Picchi and W. Zhang, "Impact of l2 cache locking on gpu performance," in *SoutheastCon 2015.* IEEE, 2015, pp. 1–4.
- [92] D. Sampaio, R. M. d. Souza, S. Collange, and F. M. Q. Pereira, "Divergence analysis," *ACM Transactions on Programming Languages and Systems (TOPLAS)*, vol. 35, no. 4, p. 13, 2013.
- [93] S. S. Baghsorkhi, M. Delahaye, S. J. Patel, W. D. Gropp, and W.-m. W. Hwu, "An adaptive performance modeling tool for gpu architectures," in *ACM Sigplan Notices*, vol. 45, no. 5. ACM, 2010, pp. 105–114.
- [94] P. Guo and L. Wang, "Auto-tuning cuda parameters for sparse matrix-vector multiplication on gpus," in *Computational and Information Sciences (ICCIS), 2010 International Conference on.* IEEE, 2010, pp. 1154–1157.
- [95] Y. Li, J. Dongarra, and S. Tomov, "A note on auto-tuning gemm for gpus," in *Computational Science–ICCS 2009.* Springer, 2009, pp. 884–892.
- [96] S. Grauer-Gray, L. Xu, R. Searles, S. Ayala-Somayajula, and J. Cavazos, "Auto-tuning a high-level language targeted to gpu codes," in *Innovative Parallel Computing (InPar), 2012.* IEEE, 2012, pp. 1–10.
- [97] A. Chaparala, C. Novoa, and A. Qasem, "Autotuning gpu-accelerated qap solvers for power and performance," in *High Performance Computing and Communications (HPCC), 2015 IEEE 7th International Symposium on Cyberspace Safety and Security (CSS), 2015 IEEE 12th International Conference on Embedded Software and Systems (ICESS), 2015 IEEE 17th International Conference on.* IEEE, 2015, pp. 78–83.
- [98] N. Sedaghati, T. Mu, L.-N. Pouchet, S. Parthasarathy, and P. Sadayappan, "Automatic selection of sparse matrix representation on gpus," in *Proceedings of the 29th ACM on International Conference on Supercomputing.* ACM, 2015, pp. 99–108.
- [99] "GitHub Repository for Autotuning CUDA Compiler Parameters for Heterogeneous Applications using the OpenTuner Framework," <https://github.com/phrb/gpu-autotuning>, accessed: 18 January 2021.
- [100] "CUDA Compiler Driver," <https://docs.nvidia.com/cuda/cuda-compiler-driver-nvcc/>, accessed: 18 January 2021.

- [101] C. S. Ferreira, R. Y. Camargo, and S. W. Song, "A parallel maximum subarray algorithm on gpus," in *Computer Architecture and High Performance Computing Workshop (SBAC-PADW), 2014 International Symposium on*. IEEE, 2014, pp. 12–17.
- [102] K. Asanovic, R. Bodik, J. Demmel, T. Keaveny, K. Keutzer, J. Kubiatowicz, N. Morgan, D. Patterson, K. Sen, J. Wawrzynek *et al.*, "A view of the parallel computing landscape," *Communications of the ACM*, vol. 52, no. 10, pp. 56–67, 2009.
- [103] C. E. Alves, E. N. Cáceres, and S. W. Song, "Bsp/cgm algorithms for maximum subsequence and maximum subarray," *PVM/MPI*, vol. 3241, pp. 139–146, 2004.
- [104] G. Holmes, A. Donkin, and I. H. Witten, "Weka: A machine learning workbench," in *Intelligent Information Systems, 1994. Proceedings of the 1994 Second Australian and New Zealand Conference on*. IEEE, 1994, pp. 357–361.
- [105] A. M. Caulfield, E. S. Chung, A. Putnam, H. Angepat, J. Fowers, M. Haselman, S. Heil, M. Humphrey, P. Kaur, J.-Y. Kim *et al.*, "A cloud-scale acceleration architecture," in *Microarchitecture (MICRO), 2016 49th Annual IEEE/ACM International Symposium on*. IEEE, 2016, pp. 1–13.
- [106] X. Yu, Y. Wang, J. Miao, E. Wu, H. Zhang, Y. Meng, B. Zhang, B. Min, D. Chen, and J. Gao, "A data-center fpga acceleration platform for convolutional neural networks," in *2019 29th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2019, pp. 151–158.
- [107] N. Tarafdar, T. Lin, E. Fukuda, H. Bannazadeh, A. Leon-Garcia, and P. Chow, "Enabling flexible network fpga clusters in a heterogeneous cloud data center," in *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2017, pp. 237–246.
- [108] R. Shu, P. Cheng, G. Chen, Z. Guo, L. Qu, Y. Xiong, D. Chiou, and T. Moscibroda, "Direct universal access: Making data center resources available to {FPGA},," in *16th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 19)*, 2019, pp. 127–140.
- [109] J. Zhang, Y. Xiong, N. Xu, R. Shu, B. Li, P. Cheng, G. Chen, and T. Moscibroda, "The feniks fpga operating system for cloud computing," in *Proceedings of the 8th Asia-Pacific Workshop on Systems*, 2017, pp. 1–7.
- [110] N. Tarafdar, N. Eskandari, V. Sharma, C. Lo, and P. Chow, "Galapagos: A full stack approach to fpga integration in the cloud," *IEEE Micro*, vol. 38, no. 6, pp. 18–24, 2018.
- [111] "Unleash Your Data Center for Intel Xeon CPU with FPGAs, howpublished = <https://www.intel.com/content/www/us/en/data-center/products/programmable/overview.html>, note = Accessed: 25 January 2020."

Bibliography

- [112] "Adaptable Acceleration for the Modern Data Center, howpublished = <https://www.xilinx.com/applications/data-center.html>, note = Accessed: 25 January 2020."
- [113] D. Singh, "Implementing fpga design with the opencl standard," *Altera whitepaper*, 2011.
- [114] T. Feist, "Vivado design suite," *White Paper*, vol. 5, 2012.
- [115] P. Coussy, D. D. Gajski, M. Meredith, and A. Takach, "An introduction to high-level synthesis," *IEEE Design & Test of Computers*, vol. 26, no. 4, pp. 8–17, 2009.
- [116] A. Canis, J. Choi, M. Aldham, V. Zhang, A. Kammoona, T. Czajkowski, S. D. Brown, and J. H. Anderson, "Legup: An open-source high-level synthesis tool for fpga-based processor/accelerator systems," *ACM Transactions on Embedded Computing Systems (TECS)*, vol. 13, no. 2, p. 24, 2013.
- [117] "LegUp Computing, now closed-source, paid software," <https://www.legupcomputing.com/>, accessed: 26 January 2021.
- [118] C. Lattner and V. Adve, "LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation," in *Proceedings of the 2004 International Symposium on Code Generation and Optimization (CGO'04)*, Palo Alto, California, Mar 2004.
- [119] Y. Hara, H. Tomiyama, S. Honda, H. Takada, and K. Ishii, "Chstone: A benchmark program suite for practical c-based high-level synthesis," in *Circuits and Systems, 2008. ISCAS 2008. IEEE International Symposium on*. IEEE, 2008, pp. 1192–1195.
- [120] J. Villarreal, A. Park, W. Najjar, and R. Halstead, "Designing modular hardware accelerators in c with roccc 2.0," in *Field-Programmable Custom Computing Machines (FCCM), 2010 18th IEEE Annual International Symposium on*. IEEE, 2010, pp. 127–134.
- [121] P. Coussy, G. Lhairech-Lebreton, D. Heller, and E. Martin, "Gaut a free and open source high-level synthesis tool," in *IEEE DATE*, 2010.
- [122] V. V. Kindratenko, R. J. Brunner, and A. D. Myers, "Mitrion-c application development on sgi altix 350/rc100," in *Field-Programmable Custom Computing Machines, 2007. FCCM 2007. 15th Annual IEEE Symposium on*. IEEE, 2007, pp. 239–250.
- [123] A. Antola, M. D. Santambrogio, M. Fracassi, P. Gotti, and C. Sandionigi, "A novel hardware/software codesign methodology based on dynamic reconfiguration with impulse c and codeveloper," in *Programmable Logic, 2007. SPL'07. 2007 3rd Southern Conference on*. IEEE, 2007, pp. 221–224.
- [124] S. Loo, B. E. Wells, N. Freije, and J. Kulick, "Handel-c for rapid prototyping of vlsi coprocessors for real time systems," in *System Theory, 2002. Proceedings of the Thirty-Fourth Southeastern Symposium on*. IEEE, 2002, pp. 6–10.

- [125] R. Nane, V. M. Sima, C. Pilato, J. Choi, B. Fort, A. Canis, Y. T. Chen, H. Hsiao, S. Brown, F. Ferrandi, J. Anderson, and K. Bertels, "A survey and evaluation of fpga high-level synthesis tools," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. PP, no. 99, pp. 1–1, 2016.
- [126] J. Luu, J. Goeders, M. Wainberg, A. Somerville, T. Yu, K. Nasartschuk, M. Nasr, S. Wang, T. Liu, N. Ahmed *et al.*, "Vtr 7.0: Next generation architecture and cad system for fpgas," *ACM Transactions on Reconfigurable Technology and Systems (TRETS)*, vol. 7, no. 2, p. 6, 2014.
- [127] Q. Huang, R. Lian, A. Canis, J. Choi, R. Xi, N. Calagar, S. Brown, and J. Anderson, "The effect of compiler optimizations on high-level synthesis-generated hardware," *ACM Transactions on Reconfigurable Technology and Systems (TRETS)*, vol. 8, no. 3, p. 14, 2015.
- [128] S. W. Nabi and W. Vanderbauwhede, "A fast and accurate cost model for fpga design space exploration in hpc applications," in *Parallel and Distributed Processing Symposium Workshops, 2016 IEEE International*. IEEE, 2016, pp. 114–123.
- [129] A. C. Canis, "Legup: Open-source high-level synthesis research framework," Ph.D. dissertation, University of Toronto, 2015.
- [130] S. Hadjis, A. Canis, J. H. Anderson, J. Choi, K. Nam, S. Brown, and T. Czajkowski, "Impact of fpga architecture on resource sharing in high-level synthesis," in *Proceedings of the ACM/SIGDA international symposium on Field Programmable Gate Arrays*. ACM, 2012, pp. 111–114.
- [131] "LegUp Documentation, now closed-source, paid software," <https://www.legupcomputing.com/main/resource>, accessed: 18 January 2021.
- [132] "GitHub Repository for Autotuning High-Level Synthesis for FPGAs using OpenTuner and LegUp," <https://github.com/phrb/legup-tuner>, accessed: 18 January 2021.
- [133] "GitHub Repository for the dockerfile used in Autotuning High-Level Synthesis for FPGAs using OpenTuner and LegUp," <https://github.com/phrb/legup-dockerfile>, accessed: 18 January 2021.
- [134] S. Q. Masnada, "Semi-Automatic Performance Optimization of HPC Kernels," Master's thesis, Université Grenoble Alpes, Jun. 2016. [Online]. Available: <https://hal.inria.fr/hal-01579422>
- [135] "BOAST Source Code," <https://github.com/Nanosim-LIG/boast>, accessed: 15 March 2021.
- [136] Q. Yi, K. Seymour, H. You, R. Vuduc, and D. Quinlan, "POET: Parameterized optimizations for empirical tuning," in *Parallel and Distributed Processing Symposium, 2007. IPDPS 2007. IEEE International*. IEEE, 2007, pp. 1–8.

Bibliography

- [137] “BOAST Laplacian Sample GPU Kernel,” <https://github.com/Nanosim-LIG/boast/blob/master/sample/Laplacian.rb>, accessed: 15 March 2021.
- [138] J. Fox and S. Weisberg, *An R Companion to Applied Regression*, 2nd ed. Thousand Oaks CA: Sage, 2011. [Online]. Available: <http://socserv.socsci.mcmaster.ca/jfox/Books/Companion>
- [139] S. Addelman and O. Kempthorne, “Some main-effect plans and orthogonal arrays of strength two,” *The Annals of Mathematical Statistics*, pp. 1167–1176, 1961.
- [140] U. Grömping and R. Fontana, “An algorithm for generating good mixed level factorial designs,” Beuth University of Applied Sciences, Berlin, Tech. Rep., 2018.
- [141] L.-Y. Deng and B. Tang, “Generalized resolution and minimum aberration criteria for plackett-burman and other nonregular factorial designs,” *Statistica Sinica*, pp. 1071–1082, 1999.
- [142] A. C. Atkinson and A. N. Donev, “The construction of exact d-optimum experimental designs with application to blocking response surface designs,” *Biometrika*, vol. 76, no. 3, pp. 515–526, 1989.
- [143] J. Kiefer, “General equivalence theory for optimum designs (approximate theory),” *The annals of Statistics*, pp. 849–879, 1974.
- [144] D. Balouek, A. Carpen Amarie, G. Charrier, F. Desprez, E. Jeannot, E. Jeanvoine, A. Lèbre, D. Margery, N. Niclausse, L. Nussbaum, O. Richard, C. Pérez, F. Quesnel, C. Rohr, and L. Sarzyniec, “Adding virtualization capabilities to the Grid’5000 testbed,” in *Cloud Computing and Services Science*, ser. Communications in Computer and Information Science, I. I. Ivanov, M. van Sinderen, F. Leymann, and T. Shan, Eds. Springer International Publishing, 2013, vol. 367, pp. 3–20.
- [145] “APP Metrics for Intel Xeon Processors,” <https://www.intel.com/content/dam/support/us/en/documents/processors/APP-for-Intel-Xeon-Processors.pdf>, accessed: 04 March 2021.
- [146] “Intel Xeon Processor E5-2630 v3,” <https://ark.intel.com/content/www/us/en/ark/products/83356/intel-xeon-processor-e5-2630-v3-20m-cache-2-40-ghz.html>, accessed: 04 March 2021.
- [147] S. Williams, A. Waterman, and D. Patterson, “Roofline: an insightful visual performance model for multicore architectures,” *Communications of the ACM*, vol. 52, no. 4, pp. 65–76, 2009.
- [148] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han, “HAQ: Hardware-Aware Automated Quantization With Mixed Precision,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8612–8620.

- [149] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [150] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [151] A. Elthakeb, P. Pilligundla, F. Mireshghallah, A. Yazdanbakhsh, S. Gao, and H. Esmaeilzadeh, "ReLeQ: An automatic reinforcement learning approach for deep quantization of neural networks," in *NeurIPS ML for Systems workshop, 2018*, 2019.
- [152] C. Zhu, S. Han, H. Mao, and W. J. Dally, "Trained ternary quantization," *arXiv preprint arXiv:1612.01064*, 2016.
- [153] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [154] A. Mishra, E. Nurvitadhi, J. J. Cook, and D. Marr, "Wrpn: Wide reduced-precision networks," *arXiv preprint arXiv:1709.01134*, 2017.
- [155] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou, "Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients," *arXiv preprint arXiv:1606.06160*, 2016.
- [156] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "PyTorch: An Imperative Style, High-Performance Deep Learning Library," in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [157] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [158] D. R. Wilson and T. R. Martinez, "The general inefficiency of batch training for gradient descent learning," *Neural networks*, vol. 16, no. 10, pp. 1429–1451, 2003.
- [159] Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, "Efficient backprop," in *Neural networks: Tricks of the trade*. Springer, 2012, pp. 9–48.
- [160] E.-G. Talbi, "Automated design of deep neural networks: A survey and unified taxonomy," *ACM Computing Surveys (CSUR)*, vol. 54, no. 2, pp. 1–37, 2021.

Bibliography

- [161] "HiddenLayer: A lightweight library for neural network graphs and training metrics for PyTorch, Tensorflow, and Keras," <https://github.com/waleedka/hiddenlayer>, accessed: 24 March 2021.
- [162] H. Bengtsson, "A unifying framework for parallel and distributed processing in r using futures," aug 2020. [Online]. Available: <https://arxiv.org/abs/2008.00553>
- [163] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [164] "sensitivity: Global Sensitivity Analysis of Model Outputs," <https://cran.r-project.org/web/packages/sensitivity/index.html>, accessed: 25 March 2021.

List of Tables

2.1	Autotuning methods used by a sample of systems, in different domains, ordered by publishing year. Methods were classified as either Search Heuristics (SH), Machine Learning (ML), or more precisely when the originating work provided detailed information. Earlier work favored employing Search Heuristics, which are less prominent in recent work, which favors methods based on Machine Learning.	23
3.1	Summary of the notation, concepts, and examples discussed in this chapter, and common to the autotuning methods discussed in further chapters	28
5.1	Expressions for the covariance functions, or kernels, shown in Figure 5.5. The variables v and l are kernel parameters that can themselves be estimated. The Matérn kernel depends on the gamma function Γ and on the Bessel function of the second kind K_v . We refer the reader to Chapter 4 of Rasmussen and Williams [10] for detailed definitions and discussions	54
6.1	A 2-level full factorial design for 3 factors, with response $y = f(x_1, x_2, x_3) + \varepsilon$	67
6.2	A 2-level fractional factorial design of resolution 3 for 5 factors, with response $y = f(x_1, x_2, x_3, x_4, x_5) + \varepsilon$. The main effects for x_4 and x_5 are aliased by, or confounded with, the interaction effects x_1x_2 and x_1x_3	69
6.3	Randomized Plackett-Burman design for factors x_1, \dots, x_8 , with experiments $\mathbf{x}_1, \dots, \mathbf{x}_{12}$, and added “dummy” factors d_1, d_2, d_3 . The response y is subject to measurement error ε	71
6.4	ANOVA table for the fit of the initial model, with default significance levels from R	71
6.5	Comparison between the response y predicted by our linear model fit and the true global optimum. Factors that were used in the model are in bold	72
6.6	D-Optimal design with 12 experiments, constructed for the surrogate $\hat{f}_\theta \sim x_1 + x_3 + x_5 + x_7 + x_8 + x_8^2 + x_1x_3$, with corresponding measurements	74

List of Tables

6.7	ANOVA table for the fit of the improved model, with default significance levels from <i>R</i>	74
6.8	Model fit comparing real and estimated coefficients , with significance measurements as standard in <i>R</i>	75
6.9	Comparison between the response y predicted by our screening and optimal design model fits and the true global optimum. Factors that were used in the models are in bold	75
10.1	Description of flags in the search space	88
10.2	Hardware specifications of the target GPU architectures	89
10.3	Rodinia [87] kernels used in the experiments	89
10.4	Parameter clusters for all Rodinia problems in the GTX 750	93
10.5	Flags added to the search space from Table 10.1	94
10.6	Hardware specifications of the GPU architectures targeted in the screening experiments	94
11.1	Subset of All Autotuned LegUP HLS Parameters	102
11.2	Weights for Optimization Scenarios (<i>High</i> = 8, <i>Medium</i> = 4, <i>Low</i> = 2)	104
11.3	Autotuned CHStone Kernels	104
12.1	Parameters of the Laplacian Kernel	116
12.2	ANOVA tests at each step. Red lines mark model terms that were considered significant, with $p < 0.05$, and fixed in a given step	118
12.3	Slowdown and budget used by 7 optimization methods on the Laplacian Kernel, using a budget of 125 points with 1000 repetitions	122
13.1	Kernels from the SPAPT benchmark used in this evaluation	126

List of Figures

2.1	49 years of microprocessor data, highlighting the sustained exponential increases and reductions on transistor counts and fabrication processes, the stagnation of frequency scaling around 2005, and one solution found for it, the simultaneous exponential increase on logical core count. Data from Wikipedia [13, 14]	14
2.2	Sustained exponential increase of theoretical <i>RPeak</i> and achieved <i>RMax</i> performance for the supercomputer ranked 1 st on TOP500 [22]	16
2.3	Processor and accelerator core count in supercomputers ranked 1 st on TOP500 [22]. Core count trends for supercomputers are not necessarily bound to processor trends observed on Figure 2.1.	16
2.4	Loop nest optimizations for $C = C + A + B^\top$, in C	18
2.5	Access patterns for matrices in $C = C + A + B^\top$, with loop nest optimizations. Panel (a) shows the access order of a regular implementation, and panel (b) shows the effect of loop <i>tiling</i> , or <i>blocking</i>	18
2.6	Loop nest optimizations for GEMM, in C	19
2.7	An exhaustively measured search space, defined by loop blocking and unrolling parameters, for a sequential GEMM kernel. Reproduced from Seymour <i>et al.</i> [5]	20
2.8	Dimension and search space size for autotuning problems from 14 domains [29, 1, 30, 31, 32, 4, 33, 2, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 5] The left panel shows a zoomed view of the right panel	22
3.1	Contour plots and slices through the global optimum, marked with a \times , for search spaces defined by variations of the Booth function. Panels (a) and (d) correspond to Equation 3.3, panels (b) and (e) to Equation 3.4, and panels (c) and (f) to Equation 3.5.	29
3.2	Illustrating the relationship between convexity of a function over a compact set and the presence of local minima. Panels (a) and (b) match the functions on the same panels of Figure 3.1	30

List of Figures

- 4.1 Contour plots and direction of greatest descent $-\nabla f(\mathbf{x})$, for search spaces defined by variations of the Booth function. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. The global optimum is marked with a \times . To aid visualization, vector magnitude was encoded by color intensity, so that darker vectors have larger magnitude. The gradient along the function's basin is near zero. 33
- 4.2 Representation of paths taken by the gradient descent method, with adaptive choice of α_k , on the search spaces defined by variations of the Booth function. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. Contour plots and direction of greatest descent $-\nabla f(\mathbf{x})$ are also shown, and the global optimum is marked with a \times 34
- 4.3 Representation of paths taken by the gradient descent method with 4 restarts, with adaptive choice of α_k , on the search spaces defined by variations of the Booth function. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. Contour plots and direction of greatest descent $-\nabla f(\mathbf{x})$ are also shown, and the global optimum is marked with a \times 35
- 4.4 An exhaustively measured search space, defined by loop blocking and unrolling parameters, for a sequential GEMM kernel. Reproduced from Seymour *et al.* [5] 36
- 4.5 Representation of paths taken by the Simulated Annealing method on the search spaces defined by variations of the Booth function. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. Contour plots and direction of greatest descent are also shown, and the global optimum is marked with a \times 39
- 4.6 Some ways of producing offspring from two parents with binary chromosomes. Crossover splits parent chromosomes and combine the resulting pieces. In general, pieces from multiple splits can be combined. Mutations are introduced randomly and correspond to flipping bits on binary chromosomes. 40
- 4.7 Color-coded generation snapshots of a Genetic Algorithm on the search spaces defined by variations of the Booth function. Hollow points and dashed lines mark members of previous populations and the regions they covered, while filled points and complete lines mark the final population. Panels (a), (b), and (c) correspond to Equations 3.3, 3.4, and 3.5 respectively. Contour plots and direction of greatest descent are also shown, and the global optimum is marked with a \times 41

4.8	Representation of paths taken by the gradient descent method, with adaptive choice of α_k , on the search spaces defined by variations of the Booth function. Panel groups (a,d,g), (b,e,h), and (c,f,i) correspond to Equations 3.3, 3.4, and 3.5 respectively. Panel groups (a,b,c), (d,e,f), and (g,h,i) correspond to Gradient Descent with restarts, Simulated Annealing, and a Genetic Algorithm, respectively. Contour plots and direction of greatest descent $-\nabla f(\mathbf{x})$ are also shown, and the global optimum is marked with a \times	43
5.1	In Linear Regression, the prediction vector $\hat{\mathbf{y}}$ is the orthogonal projection of the observations vector \mathbf{y} into the vector space spanned by the columns of \mathbf{X}	46
5.2	Three linear model surrogates for the Booth function, with a \times marking the global optimum and best surrogate predictions. The fixed experimental design \mathbf{X} used to fit all surrogates is marked by \times s. Panel (a) shows noisy measurements of Booth's function, panels (b), (c), and (d) show surrogate predictions for models fit with basis functions sets from Equations 5.10, 5.11, and 5.12, respectively	49
5.3	Effects of three covariance matrices on a multivariate normal distribution with mean vector $\mu = [0 \ 0]^\top$ and covariance matrix Σ as shown on the upper right corner of each plot	52
5.4	Reinterpreting the unrolled dimensions of 100 samples of a 20 dimension multivariate normal, on the left panel, to obtain 100 samples of functions evaluated on 20 different input points, on the right panel	53
5.5	Covariance of points $(\mathbf{x}, \mathbf{x}')$ according to four covariance functions based on the distance $\ \mathbf{x} - \mathbf{x}'\ $. Expressions for each kernel are shown in Table 5.1, and Matérn kernels use parameters $v_1, v_2 = \{\frac{3}{2}, \frac{5}{2}\}$	55
5.6	Fitting a Gaussian Process to three noise-free observations. The left panel shows 300 samples from a Gaussian prior, using the Matérn kernel to compute the covariance matrix. The center and right panels show 300 samples from the posterior distributions conditioned by one, then two more, successive noise-free observations	56
5.7	Fitting a Gaussian Process to three noisy observations, in the same conditions and with the same panel structure in Figure 5.6	56
5.8	Gaussian Process surrogates using three different model trends, fit to six noise-free observations of a single-input objective function, marked with black circles	58

List of Figures

5.9 Three Gaussian Process Regression surrogates with noisy fits for the Booth function, with a \times marking the global optimum and best surrogate predictions. The fixed experimental design X used to fit all surrogates is marked by \times s. Panel (a) shows noisy measurements of Booth's function, panels (b), (c), and (d) show surrogate predictions for models fit with linear, quadratic, and quadratic plus interactions trends, respectively	59
6.1 Fitting linear models using only two measurements, in experimental settings with different properties of the error ε	64
6.2 Fitting linear models using 30 distinct random sets of two measurements, in experimental settings with different properties of the error ε	65
6.3 Modeling linear effects in different experimental settings	66
6.4 Effect of interactions in the slopes for extreme values of factors	68
6.5 Comparison of 1000 samples obtained with different sampling methods. The <i>Latin Hypercube</i> design was produced with the <i>DiceDesign R</i> package [80], and was optimized after an initial sampling, with Simulated Annealing. Sobol, Halton, and Torus are low-discrepancy sequences from the <i>randtoolbox R</i> package [79].	76
6.6 Exploring the search space for Booth's function with different methods, using a fixed budget of 30 measurements. A \times marks the global optimum, a \blacksquare marks the best point found by each strategy, and \bullet s mark measurements	77
8.1 Some autotuning methods	81
10.1 Autotuner representation and time scale of the experiments	85
10.2 Simplified view of NVCC compilation	87
10.3 Mean speedup over 30 repetitions of the tuned solutions for Rodinia kernels versus the <i>opt-level</i> = 2 baseline	90
10.4 Mean speedup over 30 repetitions of the tuned solutions for kernels we implemented versus the <i>opt-level</i> = 2 baseline	91
10.5 Mean speedup over 30 repetitions of the tuned solutions for Rodinia kernels versus the <i>opt-level</i> = 2 baseline, across two hours of tuning. Notice the difference in the <i>y</i> -axis scales for each panel	92
10.6 Mean speedup over 30 repetitions of the tuned solutions for the kernels we implemented versus the <i>opt-level</i> = 2 baseline, across two hours of tuning. Notice the difference in the <i>y</i> -axis scales for each panel	92

10.7 Main effects estimates and 95% confidence intervals for the <i>high level</i> of the 15 factors in the Plackett-Burman design, for target kernels measured on the Titan X and Quadro M1200 GPUs	95
10.8 Execution time of baseline and model-predicted NVCC flag configurations for three Rodinia kernels on the Titan X and Quadro M1200 GPUs. Circles mark each of the 20 measurements of each flag configuration, filled dots and whiskers mark the mean and 95% confidence intervals	96
10.9 Mean speedups found using the predictions of the screening model, for three Rodinia kernels on the Titan X and Quadro M1200 GPUs. The horizontal dashed lines mark the PTX-stage <code>opt-level=3</code> comparison baseline	97
11.1 Autotuner representation and time scale of the experiments	99
11.2 Autotuner representation and time scale of more complex FPGA applications .	99
11.3 High-Level Synthesis compilation process. The autotuner search space at the HLS stage is highlighted in blue	101
11.4 Autotuner Setup	101
11.5 Comparison of the absolute values for Random and Default starting points in the Balanced scenario	105
11.6 Relative improvement for all metrics in the <i>Balanced</i> scenario	106
11.7 Relative improvement for all metrics in the <i>Area</i> scenario	106
11.8 Relative improvement for all metrics in the <i>Performance</i> scenario	107
11.9 Relative improvement for all metrics in the <i>Performance and Latency</i> scenario . .	107
11.10 Relative improvement for WNS in all scenarios	107
12.1 Edge-detection effect of a <i>Laplacian of Gaussian</i> filter	110
12.2 A CPU Laplacian kernel written in C	111
12.3 Excerpts of the BOAST code, in Ruby, that generates the Laplacian OpenCL kernel. The code is publicly hosted at GitHub [137]	112
12.4 A sample OpenCL Laplacian kernel generated by BOAST	113
12.5 Overview of the DOE approach to autotuning we implemented	114
12.6 Overview of a single DLMT run. Each panel shows a view of the completely evaluated search space, from the perspective of one of the factors eliminated in a given step. A x marks the predicted best level of each factor, and a X marks the factor level on the global optimum.	119

List of Figures

12.7 Distribution of slowdowns in relation to the global optimum for 7 optimization methods on the Laplacian Kernel, using a budget of 125 points over 1000 repetitions	121
13.1 Cost of best points found on each run, and the iteration where they were found. RS and DLMT found no speedups with similar budgets for kernels marked with “[0]” and <i>blue</i> headers, and similar speedups with similar budgets for kernels marked with “[=]” and <i>orange</i> headers. DLMT found similar speedups using smaller budgets for kernels marked with “[+]” <i>green</i> headers. Ellipses delimit an estimate of where 95% of the underlying distribution lies	127
13.2 Histograms of explored search spaces, showing the real count of measured configurations. Kernels are grouped in the same way as in Figure 13.1. DLMT spent fewer measurements than RS in configurations with smaller speedups or with slowdowns, even for kernels in the orange group. DLMT also spent more time exploring configurations with larger speedups	128
13.3 Summary of our DLMT results, compared to uniform Random Sampling (RS), showing the configurations with smallest speedups, found in 10 independent runs. The top panel compares execution time, with color-encoded iterations, and the bottom panel compares iterations where configurations were found, with color-encoded execution times.	129
13.4 Count of the model terms eliminated in each of the 4 steps of 10 independent runs, without cubic terms	131
13.5 Measured execution time of all points in the designs constructed at each step, where panel headers mark the hostname of the machines used in each of the 10 separate experiments	132
13.6 Fitting linear models and quadratic quantile regression, with $\tau = 0.05, 0.25$, and 0.5, to separate factors for 10 uniform random sampling runs	133
13.7 Same as above, for 10 uniform sampling runs with OMP fixed to <i>on</i>	133
13.8 Performance of best points, step-by-step, by experiment	135
13.9 Best <i>bicg</i> configurations and iterations that found them, for Gaussian Process Regression with Expected Improvement acquisition function (GPR), and uniform Random Sampling (RS), with a budget of 400 measurements.	136
13.10 Execution time of points measured by GPR along iterations, with pareto border in red	138
13.11 Execution time of points measured by RS along iterations	138

13.12 Best <i>dgemv3</i> configurations and iterations that found them, for Gaussian Process Regression with Expected Improvement acquisition function (GPR), and uniform Random Sampling (RS), with a budget of 400 measurements.	139
13.13 Empirical Roofline graph for the <i>Xeon E5-2630v3</i> processor [145, 146], showing the best point we found during all experiments with the memory-bound <i>dgemv3</i> SPAPT kernel	139
13.14 Empirical Roofline graph for the <i>Xeon E5-2630v3</i> processor [145, 146], showing the best point we found during all experiments with the memory-bound <i>dgemv3</i> SPAPT kernel	140
14.1 Optimizing nonuniform weight quantization policies for the <i>ResNet50</i> network, keeping total weight size below 10 MB, and attempting to keep accuracy. The optimization methods we compared were a Sobol Sampler, Reinforcement Learning [148], and Gaussian Process Regression with Expected Improvement	143
14.2 <i>ResNet50</i> architecture, implemented in <i>pytorch</i>	145
14.3 Images sampled from some of the categories in <i>ImageNet</i> , adapted from Deng <i>et al.</i> [149]	146
14.4 Pseudocode implementations of the GPR method using Expected Improvement (EI), and of the constrained sampler using Sobol sequences	149
14.5 10 repetitions of the 4 methods we compared, using a budget of 245 measurements, with 95% confidence intervals for estimates of the mean	150
14.6 Best quantization policy found by each method during optimization, across 10 repetitions	151
14.7 Results with the baseline Reinforcement Learning method	153
14.8 Results with Gaussian Process Regression	154
14.9 Running optimized quantization policies with new <i>ImageNet</i> samples, with 95% confidence intervals for the mean estimate. The baseline is the uniform 8-bit <i>ResNet50</i> quantization	155
14.10 Execution time of the different steps of our GPR implementation. The initial sample measurement is not shown, and the <i>y</i> -axis is logarithmic	156
14.11 Sobol indices computed for the impact on <i>Top5</i> accuracy of the quantization of each <i>ResNet50</i> layer using only Sobol samples, in the left panel, and Sobol samples plus the search spaces explored by GPR, in the right panel	157
	185

List of Figures

14.12Running GPR with the gathered ensemble of GPR experiments. The bottom row shows only the new measurements performed after fitting a Gaussian Process with the ensemble of experiments, and the top row shows all the data that were used	158
14.13Optimizing accuracy and weight with the modified baseline Reinforcement Learning method	159
14.14Optimizing accuracy and weight with Gaussian Process Regression with EI. . .	159