

# A Design of Experiments Approach to Autotuning under Tight Budget Constraints

Pedro Bruel<sup>\*†</sup>, Arnaud Legrand<sup>\*</sup>, Brice Videau<sup>\*</sup> and Alfredo Goldman<sup>†</sup>

<sup>\*</sup>University of Grenoble Alpes, CNRS, INRIA, LIG - Grenoble, France

Email: {arnaud.legrand, brice.videau}@imag.fr

<sup>†</sup>University of São Paulo - São Paulo, Brazil

Email: {phrb, gold}@ime.usp.br

## Abstract—Abstract

### I. INTRODUCTION

Optimizing code for objectives such as performance and power consumption is fundamental to the success and cost effectiveness of industrial and scientific endeavors in High Performance Computing. A considerable amount of highly specialized time and effort is spent in porting and optimizing code for GPUs, FPGAs and other hardware accelerators. Experts are also needed to leverage bleeding edge software improvements in compilers, languages, libraries and frameworks. The automatic configuration and optimization of High Performance Computing applications, or *autotuning*, is a technique effective in decreasing the cost and time needed to adopt efficient hardware and software. Typical targets for autotuning include algorithm selection, source-to-source transformations and compiler configuration.

Autotuning can be studied as a search problem, where the objective is to minimize single or multiple software of hardware metrics. The exploration of the search spaces defined by configurations and optimizations present interesting challenges to search strategies. These search spaces grow exponentially with the number of considered configuration parameters and their possible values. They are also difficult to extensively explore due to the often prohibitive costs of hardware utilization and program compilation and execution times. Developing autotuning strategies capable of producing good optimizations while minimizing resource utilization is therefore essential. The capability of acquiring knowledge about an optimization problem is also a desired feature of an autotuning strategy, since this knowledge can decrease the cost of subsequent optimizations of the same application or for the same hardware.

It is common and usually effective to use search meta-heuristics such as genetic algorithms and simulated annealing in autotuning. These strategies usually attempt to exploit local properties and are not capable of fully exploiting global search space structures. They are also not much more effective in comparison with a naive uniform random sample of the search space [1], [2], and usually rely on a large number of measurements and frequent restarts to achieve good performance improvements. Search strategies based on gradient descent also are commonly used in autotuning and rely on a large number

of measurements. Their effectiveness diminishes additionally in search spaces with complex local structures. Completely automated machine learning autotuning strategies are effective in building models for predicting important optimization parameters, but still rely on a sizable data set for training. Large data sets are fundamental to strategies based on machine learning since they select models from a generally very large class.

Search strategies based on meta-heuristics, gradient descent and machine learning require a large number of measurements to be effective, and are usually incapable of providing knowledge about search spaces to users. At the end of each autotuning session it is difficult to decide if and where further exploration is warranted, and impossible to know which parameters are responsible for the observed improvements. After exploring a search space, it is impossible to confidently deduce its global properties since its was explored with unknown biases.

In this paper we propose an autotuning strategy that leverages existing expert and approximate knowledge about a problem in the form of a performance model, and refines this initial model iteratively using empirical performance evaluations, statistical analysis and user input. Our strategy puts a heavy weight on decreasing the costs of autotuning by using efficient Design of Experiments strategies to minimize the number of experiments needed to find good optimizations. Each optimization iteration uses *Analysis of Variance* (ANOVA) to help identify the relative significance of each configurable parameter to the performance observations. An architecture- and problem-specific performance model is built iteratively and with user input, enabling informed decisions on which regions of the search space are worth exploring.

We present the performance of our approach on a Laplacian Kernel for GPUs where the search space, global optimum and performance model approximation are known. The experimental budget on this application were tightly constrained. The speedups achieved and the budget utilization of our approach on this setting motivated a more comprehensive performance evaluation. We chose the *Search Problems in Automatic Performance Tuning* (SPAPT) [3] benchmark for this evaluation, where our approach was able to find speedups of over 50× for some SPAPT applications, finding speedups better than random sampling in some scenarios. Despite using

generic performance models for every SPAPT application, our approach was able to significantly decrease the budget used to find performance improvements.

The rest of this paper is organized as follows. Section II presents related work on source-to-source transformation, which is the main optimization target in SPAPT problems, on autotuning systems and on search space exploration strategies. Section III presents a detailed description of the implementation of our approach and its background. It discusses the Design of Experiments concepts we incorporate, and the ANOVA and linear regression algorithms we use in analysis steps. Section IV presents our results with the GPU Laplacian Kernel and the SPAPT benchmark. Section V discusses our conclusions and future work.

## II. BACKGROUND

### A. Source-to-source Transformation

#### B. Autotuning

Rice’s conceptual framework [4] formed the foundation of autotuners in various problem domains. In 1997, the PHiPAC system [5] used code generators and search scripts to automatically generate high performance code for matrix multiplication. Since then, systems tackled different domains with a diversity of strategies. Dongarra *et al.* [6] introduced the ATLAS project, that optimizes dense matrix multiplication routines. The OSKI [7] library provides automatically tuned kernels for sparse matrices. The FFTW [8] library provides tuned C subroutines for computing the Discrete Fourier Transform. Periscope [9] is a distributed online autotuner for parallel systems and single-node performance. In an effort to provide a common representation of multiple parallel programming models, the INSIEME compiler project [10] implements abstractions for OpenMP, MPI and OpenCL, and generates optimized parallel code for heterogeneous multi-core architectures.

Some systems provide generic tools that enable the implementation of autotuners in various domains. PetaBricks [11] is a language, compiler and autotuner that introduces abstractions that enable programmers to define multiple algorithms for the same problem. The ParamILS framework [12] applies stochastic local search methods for algorithm configuration and parameter tuning. The OpenTuner framework [13] provides ensembles of techniques that search spaces of program configurations.

### C. Search Space Exploration Strategies

## III. APPLYING DESIGN OF EXPERIMENTS TO AUTOTUNING

An *experimental design* is a plan for executing a series of experiments whose objective is to identify the 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.

Experimental designs can be used with various objectives, from identifying the most important factors to building an analytical model for the response. The field of Design of Experiments 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.

One of the first detailed descriptions and mathematical treatment of Design of Experiments was presented by Ronald Fisher [14] 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 [15], Montgomery [16] and Box *et al.* [17] present comprehensive and detailed foundations.

### A. Design Construction Techniques

Our application of Design of Experiments requires support for factors of different types and number of possible values, such as binary flags, integer and floating point numerical values and unordered enumerations of abstract values. We also need designs that minimize the number of experiments needed for identifying the most relevant factors in a given problem, since at this moment we are not interesting in a precise analytical model.

The design construction techniques that fit these requirements are limited. In the *2-level screening with random level sampling* technique, factors with more than two unordered levels are sampled at two random levels. This enables using small design such as the Plackett-Burman screening design. Advantages are the small design size and good estimation capability for main effects. Disadvantages are the incapability of estimating interactions, but mainly the lack of information regarding the response for levels not selected in the initial screening.

In *contractive replacement*, an initial 2-Level design is used to generate mixed-level designs by re-encoding columns into a new single column representing a multi-level factor. The contractive replacement of Addelman-Kempthorne is a strategy of this kind. Advantages are also small design sizes and good estimation capability of main effects. Additionally, the contractive replacement technique keeps orthogonality of designs. Disadvantages are the requirements on the initial designs. Not all 2-Level designs can be contracted with those methods if orthogonality is desired.

The *direct generation* algorithm presented by Grömping and Fontana [18] enables the generation of multi-level designs with the Generalized Minimum Aberration optimality criterion by solving mixed integer problems. Advantages are the direct generation of multi-level designs and the optimality criteria. Disadvantages are the use of proprietary MIP solvers and the limitations on the size and shape of the designs that can be generated.

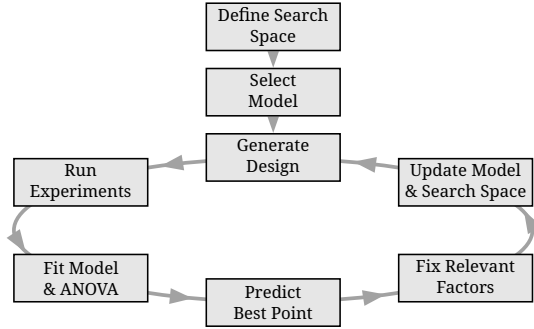
### B. D-Optimal Designs

The best candidate we have found so far are *D-Optimal* designs. Considering that we are going to analyse the results of

an experiments plan, the *D-Efficiency* of a design is inversely proportional to the *geometric mean* of the *eigenvalues* of the plan's *covariance matrix*. A D-Optimal design has the best D-Efficiency. Our current approach is based on D-Optimal designs.

### C. ANOVA

#### D. The DLMT Strategy



## IV. PERFORMANCE EVALUATIONS

### A. Example on a GPU Laplacian Kernel

### B. Results on the SPAPT Benchmark

#### 1) Experimental Methodology:

## V. CONCLUSION

## ACKNOWLEDGMENT

## REFERENCES

- [1] K. Seymour, H. You, and J. Dongarra, "A comparison of search heuristics for empirical code optimization," in *CLUSTER*, 2008, pp. 421–429.
- [2] P. M. Knijnenburg, T. Kisuki, and M. F. O'Boyle, "Combined selection of tile sizes and unroll factors using iterative compilation," *The Journal of Supercomputing*, vol. 24, no. 1, pp. 43–67, 2003.
- [3] P. Balaprakash, S. M. Wild, and B. Norris, "Spapt: Search problems in automatic performance tuning," *Procedia Computer Science*, vol. 9, pp. 1959–1968, 2012.
- [4] J. R. Rice, "The algorithm selection problem," in *Advances in Computers* 15, 1976, pp. 65–118.
- [5] 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.
- [6] J. J. Dongarra and C. R. Whaley, "Automatically tuned linear algebra software (atlas)," *Proceedings of SC*, vol. 98, 1998.
- [7] 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.
- [8] 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.
- [9] M. Gerndt and M. Ott, "Automatic performance analysis with periscope," *Concurrency and Computation: Practice and Experience*, vol. 22, no. 6, pp. 736–748, 2010.
- [10] 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.
- [11] 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, 2009, vol. 44, no. 6.

- [12] 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.
- [13] 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.
- [14] R. A. Fisher, *The design of experiments*. Oliver And Boyd; Edinburgh; London, 1937.
- [15] P. N. D. Bukh, *The art of computer systems performance analysis, techniques for experimental design, measurement, simulation and modeling*. JSTOR, 1992.
- [16] D. C. Montgomery, *Design and analysis of experiments*. John wiley & sons, 2017.
- [17] G. E. Box, J. S. Hunter, and W. G. Hunter, *Statistics for experimenters: design, innovation, and discovery*. Wiley-Interscience New York, 2005, vol. 2.
- [18] U. Grömping and R. Fontana, "An algorithm for generating good mixed level factorial designs," Beuth University of Applied Sciences, Berlin, Tech. Rep., 2018.