

Program Autotuning with Cloud Computing and OpenTuner

Pedro Bruel

phrb@ime.usp.br

Alfredo Goldman

gold@ime.usp.br

Daniel Batista

batista@ime.usp.br

Instituto de Matemática e Estatística (IME)

Universidade de São Paulo (USP)

R. do Matão, 1010 – Vila Universitária, São Paulo – SP, 05508-090

Abstract

This research project presents a proposal for a conference paper. The objective is to extend the OpenTuner autotuning framework to leverage the cloud computing resources from Google Compute Engine. An overview of the autotuning research area is presented, supporting the novelty and contribution potential of this proposal. The objectives and current work are discussed in detail, and a research schedule is presented.

1 Introduction

The program autotuning problem fits in the framework of the Algorithm Selection Problem, introduced by Rice in 1976 [16]. The objective of an autotuner is to select the best algorithm, or algorithm configuration, for each instance of a problem. Algorithms or configurations are selected according to performance metrics such as the time to solve the problem instance, the accuracy of the solution and the energy consumed. The set of all possible algorithms and configurations that solve a problem define a *search space*. Guided by the performance metrics, various optimization techniques search this space for the algorithm or configuration that best solves the problem.

Autotuners can specialize in domains such as matrix multiplication [3], dense [19] or sparse [18] matrix linear algebra, and parallel programming [13]. Other autotuning frameworks provide more general tools for the representation and search of program configurations, enabling the implementation of autotuners for different problem domains [2, 12].

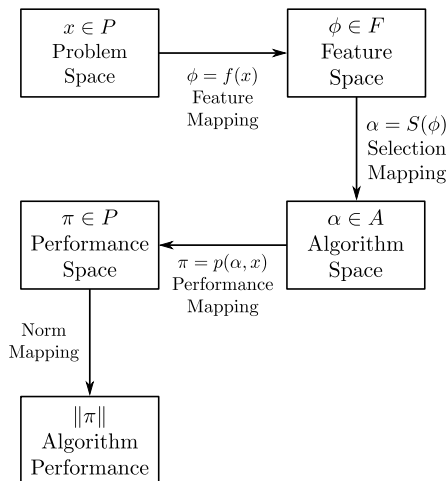


Figure 1: Rice’s Algorithm Selection Framework [16].

Figure 1 shows the Algorithm Selection Problem framework. Bougeret *et al.* [5] proved that the Algorithm Selection Problem is NP-complete when calculating static distributions of algorithms in parallel machines. Guo [8] proved the problem is undecidable in the general case.

This research proposal consists in the implementation of an extension to the OpenTuner framework [2] that will enable it to leverage cloud computing resources. The remaining of this document is organized as follows. Section 2 discusses related work and the OpenTuner framework. Section 3 presents a detailed discussion of the objectives and reports the current state of the research. Section 4 discusses the methods and benchmarks that will be used to test and validate the implementations. Section 5 presents a preliminary schedule. Section 6 concludes the proposal.

2 Related Work

Rice’s conceptual framework [16] formed the foundation of autotuners in various problem domains. In 1997, the PHiPAC system [3] 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. Whalley *et al.* [19] introduced the ATLAS project, that optimizes dense matrix multiply

routines. The OSKI [18] library provides automatically tuned kernels for sparse matrices. The FFTW [7] library provides tuned C subroutines for computing the Discrete Fourier Transform. In an effort to provide a common representation of multiple parallel programming models, the INSIEME compiler project [13] implements abstractions for OpenMP, MPI and OpenCL, and generates optimized parallel code for heterogeneous multi-core architectures.

Some autotuning systems provide generic tools that enable the implementation of autotuners in various domains. PetaBricks [1] is a language, compiler and autotuner that introduces abstractions, such as the `either...or` construct, 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 [2] provides ensembles of techniques that search spaces of program configurations. Bosboom *et al.* and Eliahu use OpenTuner to implement a domain specific language for data-flow programming [4] and a framework for recursive parallel algorithm optimization [6].

In a progression of papers [9, 10, 11], Gupta *et al.* provide experimental evaluations of the application of cloud computing to high performance computing, describing which kind of applications has the greatest potential to benefit from cloud computing. Their work highlights small and medium scale projects as the main beneficiaries of cloud computing resources.

2.1 OpenTuner

OpenTuner search spaces are defined by *Configurations*, that are composed of *Parameter* of various types. Each type has restricted bounds and manipulation functions that enable the exploration of the search space. OpenTuner implements ensembles of optimization techniques that perform well in different problem domains. The framework uses *meta-techniques* to coordinate the distribution of resources between techniques. Results found during search are shared through a database. An OpenTuner application can implement its own search techniques and meta-techniques, making the ensemble more robust. The source code is available¹ under the MIT License.

¹Hosted at GitHub: github.com/jansel/opentuner

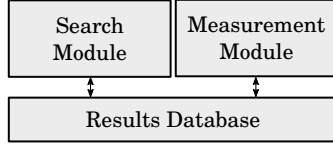


Figure 2: Simplified OpenTuner Architecture.

Figure 2 shows a high-level view OpenTuner’s architecture. Measurement and searching are done in separate modules, whose main classes are called *drivers*. The search driver requests measurements by registering configurations to the database. The measurement driver reads those configurations and writes back the desired results. Currently, the measurements are performed sequentially.

OpenTuner implements optimization techniques such as the Nelder-Mead [15] simplex method and Simulated Annealing [14]. A resource sharing mechanism, called *meta-technique*, aims to take advantage of the strengths of each technique by balancing the exploitation of a technique that has produced good results in the past and the exploration of unused and possibly better ones.

3 Objectives

The objectives of this research project are the following:

- Implement an extension to the OpenTuner measurement driver, the *MeasurementClient*, that requests virtual machines in the Google Compute Engine to produce measurements of configurations;
- Implement a *MeasurementServer* that runs in the virtual machines. The server receives requests for results, computes them, and send them to the client;
- Implement a Google Compute Engine Interface, to manage the communication between *MeasurementClient* and *MeasurementServer*;
- Define an application protocol;
- Normalize the results obtained from virtual machines to the local machine;

- Evaluate the performance of the implementation, using the benchmark that validated OpenTuner [2].

The remaining of this section describes the implementation objectives in detail, and reports the current state of the research. Section 4 describes the performance evaluations and the benchmark.

3.1 Measurement Server and Client

The interactions between the local and virtual machines will follow the client-server model. The local machine will run a measurement client, that will request results from various measurement servers running in virtual machines hosted at the Google Compute Engine.

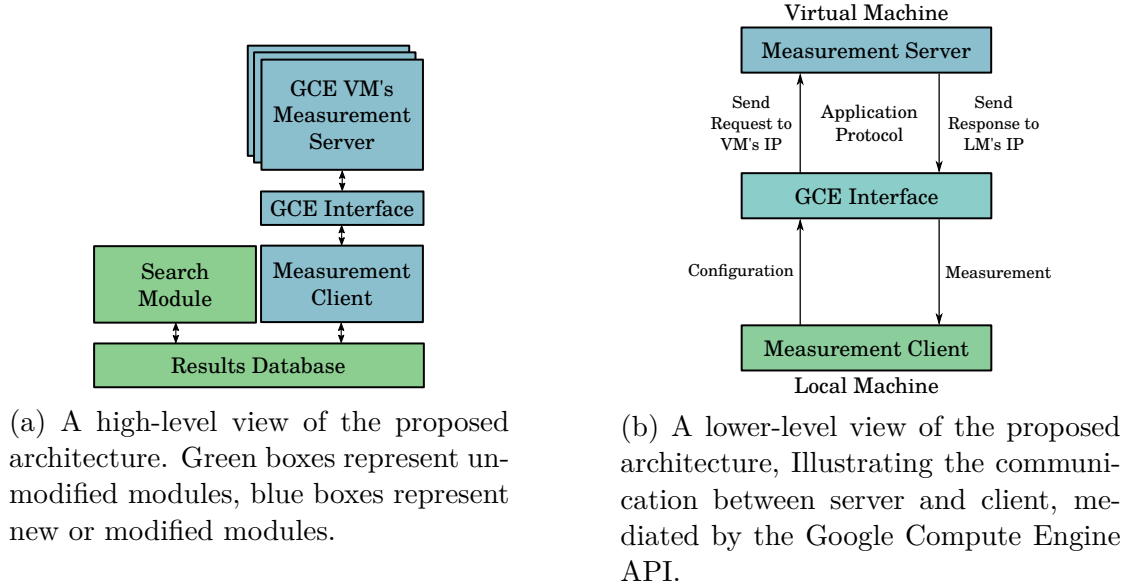


Figure 3: Proposed OpenTuner Architecture, using the Google Compute Engine (GCE) API.

Figure 3a shows the proposed architecture of an OpenTuner application running the measurement client and communicating with the measurement servers. Green boxes in the figure represent OpenTuner modules that will not be modified, and blue boxes represent new or modified modules.

Figure 3b shows, on a lower level of abstraction, the interactions between the measurement client and servers. A Google Compute Engine interface and a simple application protocol will be devised to mediate the requests and responses. The configurations will be sent to the server, that will run the correspondent program and respond with the measurement.

OpenTuner controls the execution flow of an application with the `main` function of the `TuningRunMain` class. This function sets up the database and the search and measurement modules. It then calls the `main` function of the search driver, which runs the main loop of the application. The search driver generates configurations to be tested and saves them to the database. It then calls the `process_all` function of the measurement driver and blocks until the function returns.

```
class MeasurementClient(MeasurementDriver):
    """
    Reads DesiredResults and requests VMs
    in the cloud to compute Results.
    """
    def __init__(self,
                  measurement_interface,
                  input_manager,
                  **kwargs):
        super(MeasurementDriver, self).__init__(**kwargs)
        """
        Instantiate and configure the VM Measurement
        Servers using the GCE interface.
        """

    def process_all(self):
        """
        Process all results in the
        database, by sending requests to
        VMs in the cloud and waiting
        for responses.
        """
```

Listing 1: Proposed modifications to `MeasurementDriver`.

The `process_all` function is able to compile programs in parallel, but the measurements are done sequentially. Listing 1 shows the functions of the measurement driver that will be modified to enable the `MeasurementClient` to process results with Google Compute Engine resources.

During its initialization the measurement client will initialize and configure the virtual machines, storing each measurement server's IP in the GCE interface. When the search driver makes requests for results, the `process_all` function will

route them to the servers via the GCE interface. The interface will call the appropriate GCE Python API functions, and wait for the responses.

3.2 Normalizing the Results

The autotuner will optimize programs for a machine that will typically have an architectural specification different from the machines in the cloud. A normalization technique must be devised that enables the results found in the virtual machines to be valid for the local machine. Four preliminary approaches to this problem are discussed in the following. The best approach will be experimentally determined, and could be combination of the approaches described here.

Use OpenTuner to Model Performance

Another autotuner could be implemented to optimize parameters of a simple performance model, that would associate a configuration's measurement and the virtual machine that produced it with a conversion function that transposes performance results to the target architecture.

Compose Ensembles of Virtual Machines

The cloud application could be composed of virtual machines with different architectures. The final performance measurement for a configuration would be built from some combination of the results obtained in these different virtual machines.

Simulate the Target Architecture

The target machine could be modeled by an architecture simulator such as *zsim* [17], a simulator for multi-core architectures available² under the GNU General Public License. Using a simulator would solve the normalization problem but introduce other problems, such as the simulator's accuracy and performance.

²Hosted at GitHub: <https://github.com/s5z/zsim>

Run the Autotuner in the Cloud

Finally, the normalization problem could be sidestepped, at least in initial stages of research, by running the servers and clients in the cloud using the same kind of virtual machine.

3.3 Using the Google Compute Engine

Simple experiments with the Google Compute Engine have been performed. A project was started and utility functions were implemented, such as adding and removing virtual machines, configuring firewall access rules, configuring virtual machines with a startup script and obtaining a virtual machine's IP.

For testing purposes, each virtual machine is setup with a *Debian 7* image and starts running a simple Python TCP echo server at the **8080** port. Listing 2 shows the code for the server.

The utility functions and the measurement server's and client's code are available³ under the GNU General Public License.

```
#!/usr/bin/env python

import socket

TCP_IP = ''
TCP_PORT = 8080
BUFFER_SIZE = 1024

s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
s.bind((TCP_IP, TCP_PORT))
s.listen(1)

conn, addr = s.accept()

while 1:
    data = conn.recv(BUFFER_SIZE)
    if not data: break
    conn.send(data)
```

Listing 2: A simple echo server in Python.

³All code is hosted at GitHub:
github.com/phrb/measurement-server
github.com/phrb/autotuning-gce

4 Experiments

A benchmark of applications will be composed to compare the performance of the OpenTuner framework with and without the proposed modifications. An ideal benchmark would comprise the applications used to validate OpenTuner [2]. The normalization techniques devised for the results from virtual machines will be compared using the same benchmark.

We expect that the experiments provide insight into the situations when using cloud computing resources for autotuning is beneficial, and into the application and efficiency of the result normalization techniques.

5 Research Schedule

Table 1 presents a first version of the research schedule following the assignment’s delivery dates. The schedule spreads the tasks so that portions of the work predicted to be harder are assigned to longer time frames. The delivery dates for harder tasks are also assigned to later in the project’s timeline.

AP2 & PR 25/09	AR1 16/10	AP3 & AR2 13/11	AR3 25/11
Paper proposal.	First version of the paper.	Second version of the paper.	
10-minute presentation.	Measurement Client implementation.	20-minute presentation.	Final version of the paper.
Preliminary experiments at GCE.	GCE Interface implementation.	Result normalization techniques.	Experiments, results and analysis.
High and low-level views of the implementation.	Measurement Server implementation.	Benchmark choice.	
		Application protocol full specification.	

Table 1: Task scheduling for the project.

6 Conclusion

This proposal presented a research project aiming to extend the OpenTuner autotuning framework enabling it to leverage the cloud computing resources from Google Compute Engine. We believe that distributing measurements in the cloud will considerably speedup autotuning for every problem domain. We propose four approaches to solve the result normalization problem which would enable transposing the results obtained in virtual machines to a local machine. The benchmark for the experiments will be detailed during the development of the research. It will be composed preferably by problems that evidence strengths and weaknesses of the cloud environment.

References

- [1] Jason Ansel, Cy Chan, Yee Lok Wong, Marek Olszewski, Qin Zhao, Alan Edelman, and Saman Amarasinghe. Petabricks: A language and compiler for algorithmic choice. In *ACM SIGPLAN Conference on Programming Language Design and Implementation*, Dublin, Ireland, Jun 2009.
- [2] Jason Ansel, Shoaib Kamil, Kalyan Veeramachaneni, Jonathan Ragan-Kelley, Jeffrey Bosboom, Una-May O'Reilly, and Saman Amarasinghe. Opentuner: An extensible framework for program autotuning. In *Proceedings of the 23rd international conference on Parallel architectures and compilation*, pages 303–316. ACM, 2014.
- [3] Jeff Bilmes, Krste Asanovic, Chee-Whye Chin, and Jim Demmel. Optimizing matrix multiply using phipac: A portable, high-performance, ansi c coding methodology. In *Proceedings of the 11th International Conference on Supercomputing*, ICS '97, pages 340–347, New York, NY, USA, 1997. ACM. ISBN 0-89791-902-5.
- [4] Jeffrey Bosboom, Sumanaruban Rajadurai, Weng-Fai Wong, and Saman Amarasinghe. Streamjit: a commensal compiler for high-performance stream programming. In *Proceedings of the 2014 ACM International Conference*

- on Object Oriented Programming Systems Languages & Applications*, pages 177–195. ACM, 2014.
- [5] Marin Bougeret, P-F Dutot, Alfredo Goldman, Yanik Ngoko, and Denis Trystram. Combining multiple heuristics on discrete resources. In *Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on*, pages 1–8. IEEE, 2009.
 - [6] David Eliahu, Omer Spillinger, Armando Fox, and James Demmel. Frpa: A framework for recursive parallel algorithms. Master’s thesis, EECS Department, University of California, Berkeley, May 2015.
 - [7] Matteo Frigo and Steven 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*, volume 3, pages 1381–1384. IEEE, 1998.
 - [8] Haipeng Guo. *Algorithm selection for sorting and probabilistic inference: a machine learning-based approach*. PhD thesis, Citeseer, 2003.
 - [9] Abhishek Gupta, Laxmikant V Kalé, Dejan S Milojicic, Paolo Faraboschi, Richard Kaufmann, Verdi March, Filippo Gioachin, Chun Hui Suen, and Bu-Sung Lee. Exploring the performance and mapping of hpc applications to platforms in the cloud. In *Proceedings of the 21st international symposium on High-Performance Parallel and Distributed Computing*, pages 121–122. ACM, 2012.
 - [10] Arpan Gupta, Laxmikant V Kale, Filippo Gioachin, Victor March, Chun Hui Suen, Bu-Sung Lee, Paolo Faraboschi, Richard Kaufmann, and Dejan Milojicic. The who, what, why, and how of high performance computing in the cloud. In *Cloud Computing Technology and Science (CloudCom), 2013 IEEE 5th International Conference on*, volume 1, pages 306–314. IEEE, 2013.
 - [11] Arpan Gupta, Paolo Faraboschi, Filippo Gioachin, Laxmikant V Kale, Richard Kaufmann, Bu-Sung Lee, Victor March, Dejan Milojicic, and Chun Hui Suen. Evaluating and improving the performance and scheduling of hpc applications in cloud. 2014.

- [12] Frank Hutter, Holger H Hoos, Kevin Leyton-Brown, and Thomas Stützle. Paramils: an automatic algorithm configuration framework. *Journal of Artificial Intelligence Research*, 36(1):267–306, 2009.
- [13] Herbert Jordan, Peter Thoman, Juan J Durillo, Sara Pellegrini, Philipp Gschwandtner, Thomas Fahringer, and Hans Moritsch. A multi-objective auto-tuning framework for parallel codes. In *High Performance Computing, Networking, Storage and Analysis (SC), 2012 International Conference for*, pages 1–12. IEEE, 2012.
- [14] Scott Kirkpatrick, C Daniel Gelatt, Mario P Vecchi, et al. Optimization by simulated annealing. *science*, 220(4598):671–680, 1983.
- [15] John A Nelder and Roger Mead. A simplex method for function minimization. *The computer journal*, 7(4):308–313, 1965.
- [16] John R Rice. The algorithm selection problem. 1976.
- [17] Daniel Sanchez and Christos Kozyrakis. Zsim: fast and accurate microarchitectural simulation of thousand-core systems. In *ACM SIGARCH Computer Architecture News*, volume 41, pages 475–486. ACM, 2013.
- [18] Richard Vuduc, James W Demmel, and Katherine A Yelick. Oski: A library of automatically tuned sparse matrix kernels. In *Journal of Physics: Conference Series*, volume 16, page 521. IOP Publishing, 2005.
- [19] R. Clint Whaley and Jack J. Dongarra. Automatically tuned linear algebra software. In *Proceedings of the 1998 ACM/IEEE Conference on Supercomputing*, SC '98, pages 1–27, Washington, DC, USA, 1998. IEEE Computer Society. ISBN 0-89791-984-X.