

Boosted Regression Trees for Predictive Auto-Tuning

First Last
Affiliation line 1
Affiliation line 2
anon@mail.com

First Last
Affiliation line 1
Affiliation line 2
anon@mail.com

First Last
Affiliation line 1
Affiliation line 2
anon@mail.com

ABSTRACT

Auto-tuning of parameteric, instrumented code templates has proven to be a highly effective technique for optimizing a GPU kernel for a particular computation on particular hardware ?. However, a drawback of this approach is that exhaustive, or even thorough, auto-tuning requires compiling many kernels and calling each one many times, and this process is slow. Furthermore, it is desirable to provide users with unified library abstraction boundaries for operations such as image filtering and matrix multiplication, even those these operations actually correspond to a large set of potential problem configurations with a wide variety of memory access patterns and computational bottlenecks. How can we draw on data from previous empirical auto-tuning of related problems on related hardware to make a just-in-time implementation decision for a novel problem? This paper presents a machine learning approach to auto-tuning, in which features of the current hardware platform, the kernel configuration and the problem instance are passed to a regression model (boosted regression trees) which predicts how much faster this kernel will be than a reference baseline. Combinatorial optimization strategies for auto-tuning that would normally require evaluating large number of kernel configurations on the real hardware, are made orders of magnitude faster by evaluating the surrogate regression model instead. We validate our approach using the filterbank correlation kernel described in Pinto and Cox [2012], where we find that 0.1 seconds (XXX) of hill climbing on the regression model (which we dub “*predictive* auto-tuning”) can achieve an average of 95% of the speed-up brought by 120 seconds of empirical auto-tuning XXX. Our approach is not specific to filterbank correlation, nor even to GPU kernel auto-tuning, and can be applied to almost any templated-code optimization problem, spanning a wide variety of problem types, kernel types, and platforms.

XXX: use 300% boost instead of 3x

1. INTRODUCTION

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$10.00.

Due to power consumption and heat dissipation concerns, scientific applications have shifted from computing platforms where performance had been primarily driven by rises in the clock frequency of a single “heavy-weight” processor (with complex out-of-order control and cache structures) to a platform with ever increasing numbers of “light-weight” cores. Interestingly, this shift is now not only relevant to computational sciences but to the development of all computer systems: from ubiquitous consumer-facing devices (e.g. phones) to high-end computer farms for web-scale applications (e.g. social networks).

Although the future lies in low-power multi-core hardware designs, the field lacks consensus on exactly how the different subsystems (memory, communication and computation) should be efficiently integrated, modeled and programmed. These systems have exhibited varying degrees of memory hierarchy and multi-threading complexity and, as a consequence, they have been increasingly relying on flexible but low-level software-controlled cache management and parallelism [Asanovic et al., 2006] in order to better control and understand the various trade-offs among performance, reliability, energy efficiency, production costs, etc. This evolution has profoundly altered the landscape of application development: programmers are now facing a wide diversity of low-level architectural features to write high-performance *and* portable code.

1.1 Motivation

In this rapidly evolving landscape, the construction of general development tools and libraries that fully utilize system resources remains a daunting task. Even within specialized architectures from the *same* vendor, like NVIDIA’s Graphics Processing Units (GPUs) and the Compute Unified Device Architecture (CUDA) [Nickolls et al., 2008, NVIDIA, 2011], many developers default to massive amounts of manual labor to optimize CUDA code to specific input domains. In addition, hand-tuning rarely generalizes well to new hardware generations or different input domains, and it can also be error-prone or far from optimal. One of the reason is that kernels can produce staggeringly large optimization spaces [Datta et al., 2008]. The problem is further compounded by the fact that these spaces can be highly discontinuous [Ryoo et al., 2008], difficult to explore and quasi-optimal solutions lie at the edge of “performance cliffs” induced by device-specific hard constraints (e.g. register file size and shared memory size).

1.2 Auto-Tuning

One strategy for addressing these challenges is to use one of a variety of automatic methods known collectively as “auto-tuning”. Two major auto-tuning approaches have emerged in the extensive literature covering the subject (see surveys in [Vuduc et al., 2001, Demmel et al., 2005, Vuduc et al., 2005, Williams, 2008, Datta et al., 2008, Cavazos, 2008, Li et al., 2009, Park et al., 2011]): analytical model-driven optimization and empirical optimization [Yotov et al., 2003].

The model-driven optimization approach uses analytical abstractions to model the hardware architectures, the possible code transformations and their complex interactions. Even though highly-accurate analytical models are generally intractable to build, this approach has been quite successful in the past at accelerating serial code with simplified but general abstractions. However, large speed-ups for parallel code require more accurate high-dimensional models and since this approach is bound by the quality and scalability of its abstraction, it has been less suited for highly-specialized kernels. This approach has been dominant in the compiler community, and as a result, it has generally been applied at compile-time where important run-time characteristics such as input domains are missing. These drawbacks render the model-driven optimization approach less attractive for high-performance library developers.

The empirical optimization approach, on the contrary, seeks to find the best performing code configuration by automatically generating many versions of a parametrized kernel and benchmarking them on the actual hardware (possibly at runtime, when contextual information about the hardware and software stack is the richest). This method directly optimizes the metric(s) of interest and not surrogates. One of its main advantage is that it allows any metric to be optimized without loss of generality. It is indeed possible to formulate the problem as a multi-objective optimization and minimize both speed *and* power consumption [Rahman et al., 2011], a feat that renders the analytical model-driven approach even more difficult. Due to its flexibility, empirical auto-tuning has been successfully applied to build high-performance domain-specific libraries including dense linear algebra [Clint Whaley et al., 2001, Bilmes et al., 1997], sparse linear algebra [Vuduc et al., 2005], signal processing [?], sorting [Li et al., 2004], general stencil operations [Kamil et al., 2010], etc.

Even though such a “black-box” approach does not try to build an explicit model per se, it could, if care were taken, help to greatly accelerate our understanding of the complex interactions between the set of optimizations, their parameters and the hardware architectural features (XXX: that’s what we are trying to do w/ ML-based, decision trees, etc.).

The empirical approach is very sensitive to the choice of instrumented optimizations *and* to the search method. The size of the search space is so large that the current best empirical auto-tuners can only consider highly-specialized functions with a limited set of code transformations and compiler options, on very few input domains [Ganapathi et al., 2009]. Although searching for good code configurations in this highly-discontinuous space can be made embarrassingly parallel, it remains a very difficult and prohibitively expensive combinatorial optimization problem as many variants of the code must be generated, compiled, and benchmarked on specific input domains with meaningful statistics (that may require multiple runs). Consequently, most proposed methods prune the space with hard-coded heuristics that offer

little generalization guarantees. This has been the main drawback of the empirical approach compared to the model-driven approach where good code configurations can be directly derived from the analytical model.

To alleviate this major weakness, it is intuitively appealing to combine the two approaches by first constraining the search space with an analytical model and then exploring the reduced space empirically [Chen et al., 2005, Li et al., 2009]. Unfortunately, this hybrid approach is still bound by the quality of the analytical model, which remains hard to build by hand.

In this paper, we propose to *learn the model* using non-linear regression modelling techniques instead of constructing a model manually. By learning the model, one can hope to achieve elements of the best of both approaches: XXX.

Various statistical prediction techniques have been applied with success at compile-time for general programs on various CPU architectures [Monsifrot et al., 2002, Stephenson et al., 2003, Yotov et al., 2003, Kulkarni et al., 2004, Cooper et al., 2005, Franke et al., 2005, Hutter et al., 2006, Cavazos et al., 2007, Cavazos, 2008, Hartono et al., 2009, Park et al., 2011, Fursin et al., 2008]. Relative to this work, our contribution is to show how to do fast predictive auto-tuning that satisfies the requirements to: (a) handle the variety of recent multi-core architectures like GPUs [Schaa and Kaeli, 2009], (b) provide high-performance domain-specific libraries [Nukada and Matsuoka, 2009, Li et al., 2009, Kamil et al., 2010], (c) that select good implementations at run-time [Klöckner et al., 2011, Pinto and Cox, 2012], and (d) for the full input domain of a library routine [Liu et al., 2009, Grauer-Gray and Cavazos, 2011].

1.3 Case Study

Current languages and software engineering practices have conspired to make it difficult to provide high-performance libraries for general mathematical operations on GPU devices. An important approach in implementing high-performance GPU codes has been auto-tuning, which is fundamentally slow. Auto-tuning requires running a large number of sub-compilations and representative computations to *measure* which of several possible implementations is fastest. General-purpose libraries such as CUBLAS and CUFFT hide this overhead by auto-tuning offline and using intimate knowledge of the mathematical problem being solved as well as the kernel implementations to heuristically choose a good computation strategy quickly when a user asks for e.g. a matrix product of particular arguments. Designing and testing libraries of this form is time-consuming, difficult, and provides little reusable infrastructure to assist in the implementation of different computations. For instance, even if one has already developed an auto-tuned library for dense matrix algebra, one must basically start from scratch to develop an auto-tuned library for stencil computations or sparse matrix operations.

This work shows that auto-tuning can be accelerated by orders of magnitude by using a regression model built offline as a surrogate for actual computations on the real hardware. The general form of an auto-tuning based library routine is illustrated in Figure 1 (top). We will use filterbank correlation (Section 1.4) as our running example, but the approach is both general and powerful. An auto-tuning based routine must operate on three sets of variables:

A: task description (argument shapes, physical layout)

| | |
|---|---|
| AUTOTUNE_EMPIRICAL (<i>shapes, strides</i>) | |
| 1 | $a \leftarrow \text{TaskFeatures}(\text{shapes}, \text{strides})$ |
| 2 | $c \leftarrow \text{PlatformFeatures}()$ |
| 3 | $b^* \leftarrow \text{argmin}_{b \in \mathcal{B}} \text{MeasureTime}(a, b, c) \triangleright \text{slow}$ |
| 4 | return b^* |
| AUTOTUNE_PREDICTIVE (<i>shapes, strides</i>) | |
| 1 | $a \leftarrow \text{TaskFeatures}(\text{shapes}, \text{strides})$ |
| 2 | $c \leftarrow \text{PlatformFeatures}()$ |
| 3 | $f \leftarrow \text{LoadTimingModel}()$ |
| 4 | $b^* \leftarrow \text{argmin}_{b \in \mathcal{B}} f(a, b, c) \triangleright \text{fast}$ |
| 5 | return b^* |

Figure 1: Pseudo-code template for on-demand auto-tuning. Empirical auto-tuning (above) is inevitably slow because dynamically-generated code must be compiled and run on a number of actual-size inputs. Predictive auto-tuning (below) can be orders of magnitude faster. We show that it can also be accurate.

\mathcal{B} : implementation description (auto-tuning parameters)

C : platform description (capabilities, micro-benchmarks)

The hypothetical auto-tuning routine described at the top of Figure 1 might take many minutes or hours to perform the argmin at step 3 (during which time it computes the desired result many times!) so it would be all but unusable as a normal library routine. However, the form of the auto-tuning routine suggests the potential for enormous acceleration: if only there were a fast (even approximate) surrogate for the costly *MeasureTime*(\cdot) function, then the argmin could be done in a fraction of a second and the routine could be used normally (Figure 1, bottom).

Our contribution of the present work is to show that boosted regression tree models, a powerful machine learning technique for function approximation, can meet the requirements of this application. We call this approach *predictive* auto-tuning to contrast it with the standard measurement-based approach which we will call *empirical* auto-tuning. We show that the kernel from the GCG chapter can be modeled on a variety of hardware and get good performance compared to empirical auto-tuning.

1.4 Filterbank Correlation

Filterbank correlation is a simple spatial image filtering operation that is an important subroutine in many image processing applications. It has a relatively high arithmetic intensity which makes it a natural fit for GPU platforms [Pinto and Cox, 2012]. XXX cite some computer vision papers.

Mathematically, we define filterbank correlation in terms of an image x and a filterbank f . The image x has R rows, C columns, and D channels (e.g. color channels) that we call its *depth*. We index x like $\mathbf{x}[i, j, d]$ where $0 \leq i < R$, $0 \leq j < C$, and $0 \leq d < D$. The filterbank f has F filters that are like little images: each has a height H , a width W , and D channels. We will restrict ourselves to what are called *valid* correlations, in which the image is larger in both rows and columns than the filters. The result of filterbank

correlation of x with f is an image-like array z with $R-H+1$ rows, $C-W+1$ columns, and depth F , whose elements are defined according to Equation 1:

$$\mathbf{z}[r, c, k] = \sum_{w=0}^{W-1} \sum_{h=0}^{H-1} \sum_{d=0}^{D-1} \mathbf{x}[r+h, c+h, d] \mathbf{f}[k, h, w, d]. \quad (1)$$

In terms of floating point operations, a filterbank correlation requires the inner sums to be computed for each output pixel, yielding the quantity in Eq. 2:

$$\text{FLOPS} = 2FHW D(R-H+1)(C-W+1) \quad (2)$$

The multiplicative factor of 2 arises because we must first multiply an element of x with an element of f and then add the result to an element of z .

The memory transfer requirements of filterbank correlation are more difficult to quantify. Assuming three kinds of non-register memory – constant, shared, and global – and assuming optimistically that the entire filterbank fits into the GPU’s constant memory, then we can establish a lower bound (Eq. 3) on the amount of memory that must be moved in order to store the computed result to global memory starting from arguments in global memory:

$$\begin{aligned} \text{Bytes} &= 4RCD \\ &\quad + 4FHW D \\ &\quad + 4(R-H+1)(C-W+1)F. \end{aligned} \quad (3)$$

In short, we must read the filterbank and image once, and store the result.

The arithmetic intensity of filterbank correlation, assuming our lower bound on memory transfers is therefore approximately

$$\text{intensity} \approx \frac{FDHW}{2(D+F)}, \quad (4)$$

for images that are large relative to filters. Each F output writes corresponds to approximately D input reads and F inner products between DHW elements.

The high potential for arithmetic intensity makes the GPU an ideal platform for computing filterbank correlations, and and filterbank correlation is used extensively in image and video processing, where it is often a computational bottleneck. One might expect then, that it would be easy to implement a library providing this operation as a simple function that takes pointers and strides for x , f , and z and performs the computation. However, as shown in Pinto and Cox [2012] and numerous articles on related stencil operations XXX, it is challenging to provide an implementation or even an implementation strategy that provides satisfactory performance across the range of inputs (shapes, physical layouts) that occur in typical usage. Kamil et al. [2009] summarize a related situation related to general stencil computations in their abstract: “Although the auto-tuning strategy has been successfully applied to libraries, generalized stencil kernels are not amenable to packaging as libraries.”

XXX Datta [2009]

This paper contributes a model-based approach to auto-tuning in which the model is determined almost entirely using machine learning rather than domain knowledge. We show that even though our timing model has no built-in knowledge of GPU devices, or our kernel, or filterbank correlation, it is able to provide fast and accurate timing esti-

mates. We show that a model built by boosting regression trees is fast and accurate enough to be used profitably in the algorithm schema shown at the bottom of Figure 1, which permits doing on-demand auto-tuning for novel (not already auto-tuned) argument combinations. We call this approach *predictive auto-tuning*.

The paper is organized as follows: Section 2 describes the boosted regression tree model and the procedure for fitting it to empirical timing data. Section 3 describes the sort of kernel we employ for our benchmarking, Section ?? presents the results of our benchmarking experiments, which compare a reference implementation to a) empirical auto-tuning over a domain-specific grid, b) empirical auto-tuning over a hill-climbing search, and c) predictive autotuning. Section ?? summarizes previous and related work involving machine learning with performance auto-tuning. Section ?? summarizes our findings and outlines directions for future work.

2. BOOSTED REGRESSION TREES

A regression tree is a piece-wise constant function from one vector space to another, that works by recursively subdividing the input space into constant regions [??]. They are widely used in statistics and data-mining applications because the fitting algorithm is quick and reliable, and the form of the tree can provide insight into the relevant input variables. We use a standard fitting procedure, which takes a set of $(x, y) \in \mathbb{R}^k \times \mathbb{R}$ pairs and constructs a tree with a low mean squared error. To construct each node of a regression tree, we sort the set \mathcal{D} of (x, y) pairs along each of the k features to find the best partitioning $f_{i,\gamma}$ of the input space along feature i at point γ (Eqs. 5, 6).

$$f_{i,\gamma}(x) = \begin{cases} \alpha & \text{if } x_i < \gamma \\ \beta & \text{if } x_i \geq \gamma \end{cases} \quad (5)$$

$$i^*, \gamma^* = \underset{i, \gamma}{\operatorname{argmin}} \hat{\mathbb{E}} [(y - f_{i,\gamma}(x))^2] \quad (6)$$

One disadvantage of the regression tree is that it does not make full use of broad patterns in the data – each partition formed by the fitting procedure is fit independently in the recursive training procedure, so it is impossible for the model to extract more than one bit of information from each training partition (XXX unclear). This disadvantage is mitigated to a significant extent by the practice of *boosting* ?.

Boosting is an iterative procedure for constructing an *ensemble* of regression trees that is coordinated to fit training examples as accurately as possible. In a recent empirical study of a range of machine learning regressino problems, boosted decision trees were found to be among the best and easiest models to apply ?. On each boosting iteration, a regression tree is fit to the residual error remaining after all previously-fit models have made their predictions. There are essentially three parameters that control the boosted regression tree training procedure: 1) the depth of tree constructed on each boosting iteration, 2) the minimum number of examples to allow at a regression tree leaf, and 3) the number of trees constructed by boosting. We did not attempt a systematic study of the effect of these variables on performance. We chose a maximum depth of 4, a minimum number of examples of 10, and 100 iterations of boosting.

3. GPU IMPLEMENTATION

The strategy we use for computing filterbank correlation on the GPU using CUDA follows Pinto and Cox [2012]. The overall strategy is to load the filterbank into constant memory, which is relatively fast and visible to all threads, and then launch a grid of blocks that tiles the output image. Each thread computes $4 \times n_output_4s$ channels for some column and row of z . Each block of threads computes $4 \times n_output_4s$ channels for a subrectangle of the output image (z). When there are more than $4 \times n_output_4s$ channels in z , or if the filterbank is too large to fit into constant memory, then multiple kernel executions perform the full computation. Our approach permits splitting the filterbank along the number-of-filters dimension (F) and the height dimension (H). All the filterbanks in our study are small enough that at least one row of a single filter can fit into constant memory. Pseudo-code for the kernel is given in Figure 2.

```

THREAD_FBCORR( $gX, cF, gZ$ )
1  shared  $sX \leftarrow$  all channels of region ( $\beta$ ) of  $gX$ 
2   $x, y \leftarrow$  position of this thread in output image
3  __syncthreads()
4   $v[0 : N] \leftarrow 0$ , for  $N = 4 \times n\_output\_4s$ 
5  for  $d \leftarrow 0$  to  $D$ ,
6    for  $h \leftarrow 0$  to  $H/n\_filter\_r$ ,
7      for  $w \leftarrow 0$  to  $W$ ,
8         $u \leftarrow sX[x + h, y + w, d]$ 
9        for  $n \leftarrow 0$  to  $n\_output\_4s - 1$ ,
10          $v[n] \leftarrow v[n] + cF[n, h, w, d]$ 
11        for  $n \leftarrow 0$  to  $n\_output\_4s - 1$ ,
12          $gZ[x][y][4n:4n+n] += v[4n:4n+n]$ , (float4)

```

Figure 2: Kernel pseudo-code for filterbank correlation. Input gX is a pointer to x in global memory, input cF is a pointer to f in either constant or texture memory, and output gZ is a pointer to z in global memory. Each block of threads modifies $4 \times n_output_4s$ channels of a rectangle (called β in code listing) within z . A grid of blocks covers all rows and columns of z . Multiple calls can be used to apply all filters of a large filterbank f to x .

The kernel is parametrized by 10 parameters:

$\text{block_h} \in (4, 8, 16, 32, 64, 128)$
 $\text{block_w} \in (4, 8, 16, 32, 64, 128)$
 $\text{n_filter_r} \in (1, 2)$
 $\text{n_output_4s} \in (\text{all}, 1, 2)$
 $\text{spill} \in (\text{False}, \text{True})$
 $\text{imul_fast} \in (\text{False}, \text{True})$
 $\text{pad_shared} \in (\text{False}, \text{True})$
 $\text{use_tex1d} \in (\text{False}, \text{True})$
 $\text{maxrreg} \in (8, 16, 20, 24, 28, 32, \infty)$
 $\text{fast_math} \in (\text{False}, \text{True})$

The block height (“block_h”) and block width (“block_w”) parameters control the number of threads that run within

each block. Each kernel call loads some number of filter rows (“`n_filter_r`”) into constant memory and processes the correlation of the image with just those rows, incrementing the output buffer. Each thread can compute several output elements at once, in multiples (“`n_output_4s`”) of 4; this increases the efficiency of each thread, but can lead to lower occupancy. Registers are a precious commodity on the GPU, and this kernel accumulates elements of v in registers. The “spill” parameter controls whether the current thread’s output position in gZ is stored in a register (faster access) or in shared memory (frees up a register). The “`imul_fast`” parameter controls whether integer multiplication is done in 24-bit (True) or 32-bit (False) precision. The “`pad_shared`” parameter controls whether the sX shared memory buffer is padded, which wastes space in shared memory but reduces bank conflicts. The “`use_texld`” parameter controls whether the image is loaded into shared memory with global pointer dereferences or texture fetches. The “`maxrreg`” and “`fast_math`” parameters are passed to the nvcc compiler to limit the number of registers available to each thread, and to enable XXX, respectively.

When the entire filterbank does not fit into the GPU’s constant memory, P passes are necessary to compute all of z , where

$$P = \frac{FH}{4 \cdot n_output_4s \cdot n_filter_r}.$$

In such cases, the number of bytes moved to and from global memory is much higher than the theoretical lower limit.

$$\begin{aligned} \text{Bytes} = & 4RCDP \\ & + 4FHW D \\ & + 8(R - H + 1)(C - W + 1)FP. \end{aligned}$$

These passes make the I/O requirements increase quadratically in F and H . At the same time, the total number of floating-point operations (Eq. 2) is quadratic in H and W . In our experiments, we only considered square filters so in our setting the total number of flops is proportional to H^4 .

XXX Arithmetic density with passes.

Critically: The arithmetic intensity, shared storage, and register requirements of this kernel change significantly and in a complicated platform-dependent way with the argument parameters (R, C, D, F, H, W) and with the implementation parameters, especially “`block_w`”, “`block_h`”, “`n_output_4s`” and “`n_filter_rows`.”

4. EXPERIMENT SETUP

Recall from the introduction (Eq. ??) that auto-tuning can be seen as a conditional optimization problem in which we seek an implementation ($b \in \mathcal{B}$) that minimizes runtime or some other scalar-valued cost function for given arguments ($a \in \mathcal{A}$) on a particular platform ($c \in \mathcal{C}$). In order to perform predictive auto-tuning with a regression model, it is necessary to characterize these three types of variables with *features*. We describe the arguments to a filterbank correlation with the 6-tuple (R, C, D, F, H, W) . We randomly sampled arguments (uniformly) from the following product

space:

$$\begin{aligned} R = C & \in \{256, 512, 1024, 2048, 4096\} \\ H = W & \in \{3, 5, 7, 9, 11\} \\ D & \in \{1, 4, 8, 16, 32, 64, 128, 256\} \\ F & \in \{1, 4, 8, 16, 32, 64, 128, 256\} \end{aligned}$$

A library implementation of this operation would ideally support all image and filter sizes as well as variations due to strided memory layouts. In such a setting would be useful to characterize the arguments with features such as whether the inputs are Fortran-style contiguous, C-style contiguous, or row-padded to various byte alignments. These additional options would make our approach of automatic auto-tuning even more important, because there would be a greater variety in the kinds of computations and memory transfers to perform. Our experiments consider a somewhat simplified setting in which the arguments are always stored with depth channels being contiguous in memory, followed by columns, then rows, and then filters having the largest stride.

The product space in our study includes 1600 argument combinations, but we restricted our experiments to correlations that represented between 1 and 50 gigaflops of arithmetic. Smaller problems do not fully utilize GPU hardware and are handled equally well by many kernel settings. Larger take so long to evaluate that there is negligible inefficiency in implementing them via multiple calls with smaller images and fewer filters. With the experiments searched an argument space included 602 configurations with between 1 and 50 gigaflops.

For the implementation features b , we directly used the integer and binary values (`block_w`, `block_h`, etc.) that parametrized the kernels. We did not use platform features (c) in our experiments. We leave the investigation of cross-platform predictive auto-tuning for future work.

4.1 Search Algorithms

The core of any online auto-tuning algorithm (Fig. ??) is the search procedure used to explore the space of implementations. Our implementation space includes 16000 (XXX) points, too many to search exhaustively for every input configuration. ? recommends a particular hyper-parameter implementation based on empirical autotuning on several platforms, that achieves good performance on a variety of GPUs from older-generation cards such as the 8600GT all the way to current-generation flagship cards such as the GTX 580 and C2070. We call this configuration the *reference* implementation.

Additionally, ? advocate a particular grid search over what was estimated to be the most relevant part of the configuration space. This grid iterates over all combinations of XXX, but leaves the value of XXX at. In our experiments, we call this algorithm the *grid* search procedure. The grid included 72 points in addition to the reference implementation, for a total of 73 points. XXX: isn’t 73 prime?? explain this so-called “grid”.

In this work, we also consider a third search algorithm based on a local search heuristic. This *hill-climbing* (HC) algorithm starts from the reference implementation and re-samples each of the parameters of the current best implementation randomly with probability 0.25 (keeping the current best setting with probability 0.75). On each hill-climbing iteration, if the speed of the newly sampled point is greater

than the previous point, then it becomes the current point. We show results for search variants HC25, HC50, and HC75, which correspond to hill-climbing for 25, 50, and 75 iterations respectively.

4.2 Platforms and Implementation

This kernel was implemented in the meta-programming style advocated in Pinto and Cox [2012] in Python using Cheetah for string processing [?] and PyCUDA [Klöckner et al., 2009] for dynamic kernel compilation and interfacing with CUDA.

We evaluated our approach on several GPU devices: XXX.

XXX Table showing theoretical peak FLOPS, theoretical peak bandwidth, core count, shared memory size, register size, etc.

5. RESULTS

Figure 11 shows the effectiveness of empirical auto-tuning in this setting. Averaging across the range of problem configurations in our study, implementations discovered by empirical auto-tuning are on average about 50% (XXX check) faster than the reference implementation. XXX: recalculate with simpler mean speedup formula. On the GTX 295 speedups over the reference implementation ranged from XXX to XXX with mean XXX. On the GTX 480 speedups over the reference implementation ranged from XXX to XXX with mean XXX. XXX: Consider table for these numbers. XXX: Put in some timings of mixed & matched all-stars to show that indeed the reference was a good general baseline, and in fact the auto-tuned models are special-purpose.

Figure 11 also shows that random hill-climbing is at least competitive and often better than the hand-chosen grid approach of Pinto and Cox [2012]. Figure ?? shows how the number of hill-climbing iterations affects the quality of the auto-tuned model. The quicker HC25 approach is clearly inferior but the quality of the 75 iteration search (HC75) is statistically similar to the shorter HC50, suggesting that XXX. Still, HC75 is not always better than the grid, so evidently neither search algorithm is perfect.

The previous figures establish that auto-tuning on a problem-by-problem basis can achieve good performance improvements over a high-quality reference baseline, but they do not show how long it takes to find these implementations. A search of 75 iterations took an average of about 120 seconds XXX on a fast computer, because it requires compiling 75 CUDA kernels, and evaluating up to 750 filterbank correlations to get reliable timing information for each candidate implementation.

Our contrib

TL;DR: Hill-climbing is a good way to search the implementation space, can search the full space and do about as well as grid search in the more restricted space.

5.1 Model-based Autotuning

On the GTX 295, a single regression-tree fit to the log-speed-multiplier (LSM) function (XXX: how many problem configurations, how much vs. what kind of training data) achieves an average Spearman correlation of 0.9 on held-out test data (XXX: what kind of test data).

Now: the big test is that if we use the genetic search on the model instead of the original function, do we get good performance?

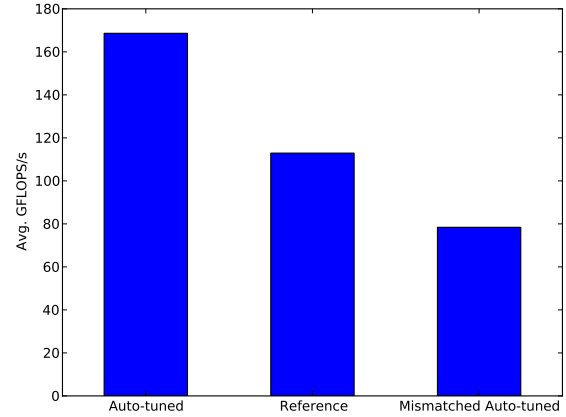


Figure 3: Different arguments call for different kernels: *left* is the average argument-specific empirical auto-tuned performance across 100 random argument configurations, *middle* is the average speed of our reference implementation, *right* is the average speed of kernels auto-tuned for different argument configurations than the one being tested. For 38/100 random configurations, the kernel auto-tuned for another problem could not even run on the GTX 295 hardware. These points contributed a speed of 0, bringing down the average much lower than the reference. Good performance across a variety of inputs requires *input-dependent auto-tuning*.

6. DISCUSSION

We’re characterizing hardware by a 1-of-N feature vector, simply describing which hardware is the current hardware. To make better use of auto-tuning data from other platforms, it would be more useful to have precise and descriptive features such as: what compute capability is present, how many cores are active, what is the bandwidth between the various kinds of memory, and how much of various kinds of memory is present. With these features a model-assisted auto-tuning approach might be able to make very good guesses on hardware for which no auto-tuning has ever been done.

In a complete implementation of Eq. 1, there would be parameters related to the physical layout (e.g. strides) of x , f and z arrays in addition to the mathematical parameters of heights and widths and so on, so the total number of filterbank correlation computations that a user might be able to demand is astronomical. A typical coping mechanism would be to cast an arbitrary problem configuration into a more standardized form, such as by copying inputs into aligned memory buffers with appropriate padding, and then choosing an appropriate blocking strategy for the computation. At that point, auto-tuning efforts can be focussed on the kernel for each blocked strategy. However, on GPU hardware, the cost of aligning the inputs can be relatively large. In future work we hope to apply our model-based technique to auto-tune in full problem configuration space, so that we can optimize as much of the computation as possible within the predictive auto-tuning framework (XXX name).

Platform Space: XXX

XXX: GCG chapter pg 13 notes that among the grid

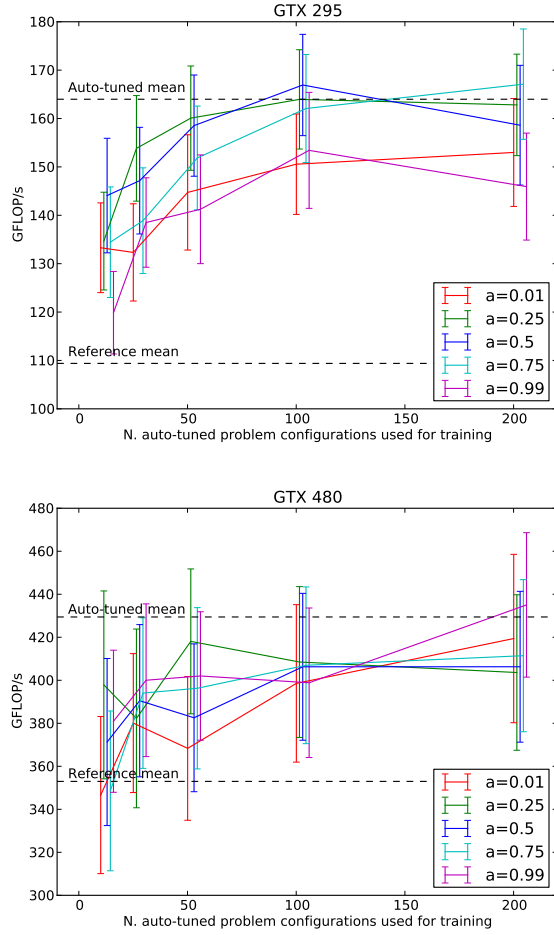


Figure 4: The effect of the invalid configuration score (a) and training set size on simulated auto-tuning. All candidates timed during the grid and hill-climbing search procedures were used as training points, so the training set sizes ranged from an average of 1500 (10 problem configurations) to 30 thousand (200 problem configurations). Training on 10 or 25 configurations was useful (higher mean than the reference), but not as useful as training on 50 or more configurations. The results on the GTX 295 suggest that moderate values of a between 0.25 and 0.75 might be best, but a had no significant effect on the GTX 480.

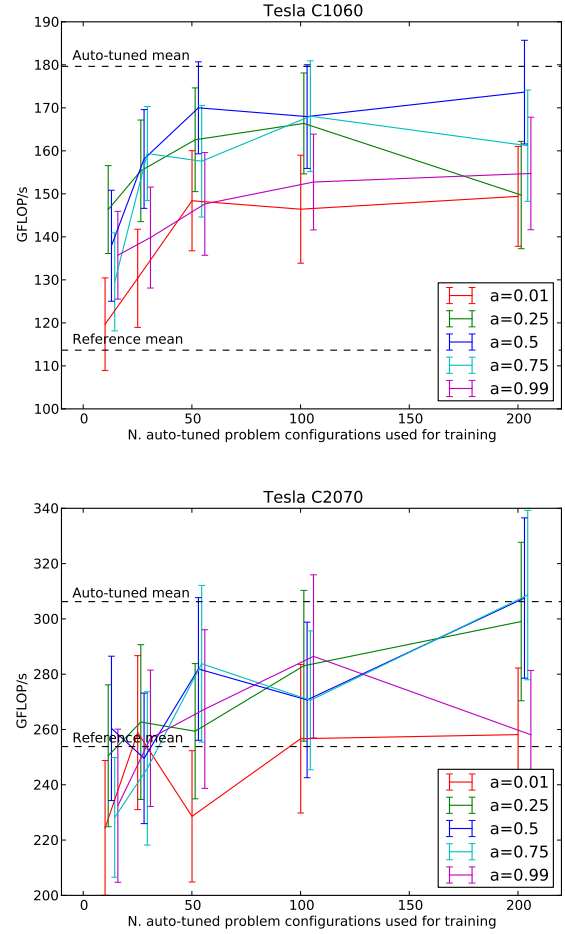


Figure 5: more timing on munctional0

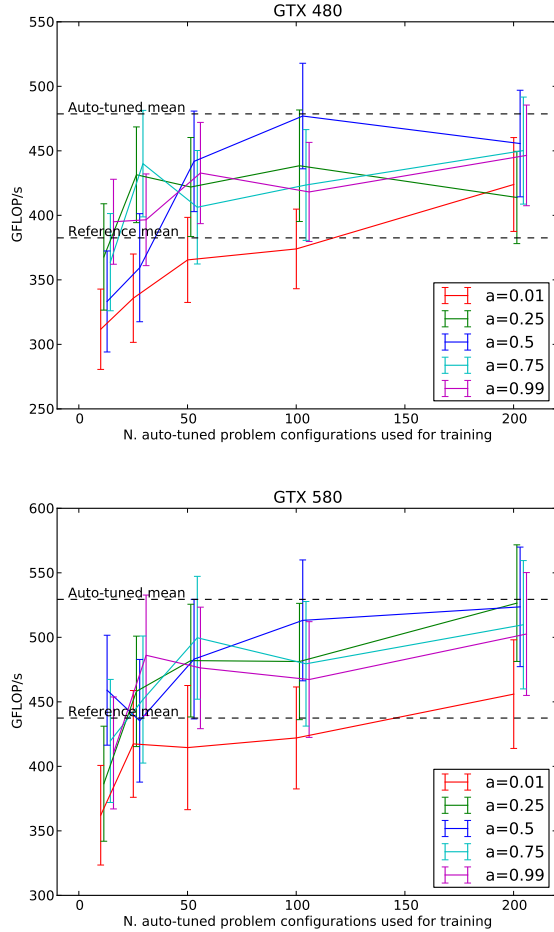


Figure 6: more timing on munctional0

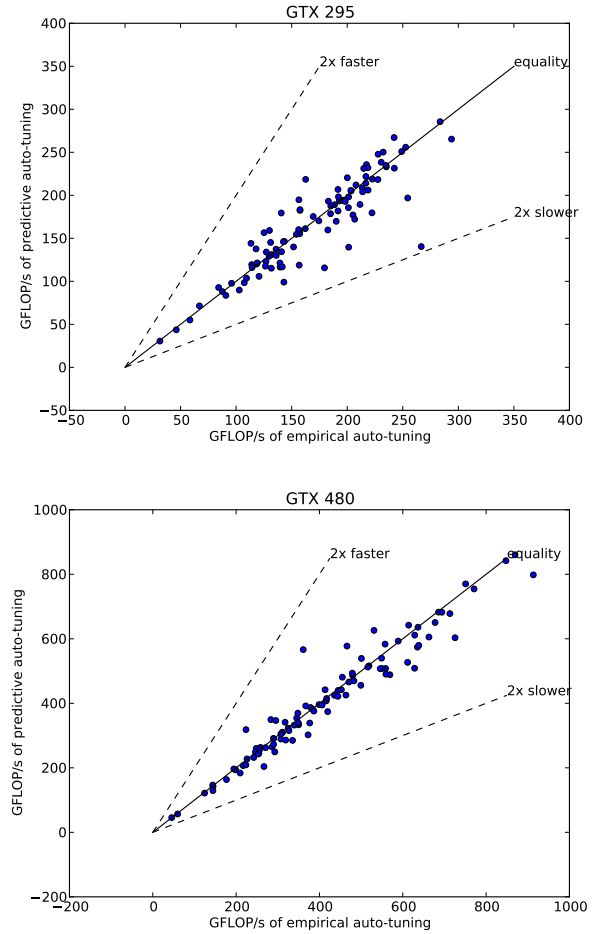


Figure 7: Computation speed for novel problem configurations using predictive vs. empirical code specialization. The scatterplots are roughly symmetric about the axis of equality (main diagonal) with occasional outliers, indicating that for most problem configurations the predictive approach gives as good a solution as the empirical approach. The predictive approach typically requires about 0.1 seconds to suggest an implementation, whereas the empirical approach requires about 2 minutes.

search, the best parameter setting is different for the various platforms.

XXX: so far we only tried 3 cards... so we're using a 1-of-N feature.

One might expect then, that it would be easy to implement a library providing this operation – this hypothetical library could have quite a simple interface following the spirit of FFTW or a BLAS routine:

```
gpu_fbcorr(  
    img_shape, filters_shape,  
    img, filters, output)
```

However, as shown in Pinto and Cox [2012] and numerous articles on related stencil operations XXX, it is challenging to provide an implementation or even an implementation strategy that provides satisfactory performance across the range of inputs (shapes, physical layouts) that occur even in typical usage.

Current state-of-the art is to do empirical auto-tuning. XXX

Increasingly, machine learning and statistical inference strategies are finding application in compiler technology. XXX

machine learning

This paper shows that a kernel with several boosted

There is a natural function from problem space \times implementation space \times platform space to runtime: how much wall time elapses on the given platform when solving the given problem with the given implementation.

Auto-tuning is a family of empirical techniques for finding the implementation that minimizes that runtime function, when the platform and the problem configuration are given.

XXX: some math, or possibly a picture with the three boxes that was on Pinto's whiteboard.

This paper is about different ways of doing the minimization. For a particular problem (filterbank correlation) we look at a parameterized kernel implementation (from GCG) and a few hardware platforms, and compare random search, grid search, and a hill-climbing strategy with a good, but generic, reference implementation. Then we show that we can model the LSM function with a regression tree well enough to usefully perform the argmin of Equation ?? on the model rather than the actual system. Whereas it may take several minutes to perform a hill-climbing or grid search in the original system, a hill-climbing search in the model requires a small fraction of a second. Model-based autotuning makes it possible to simulate auto-tuning quickly and accurately enough to be used within a single library call of the computational routine.

Best arithmetic density

A bandwidth analysis that takes the blocking structure into account would reveal that larger H and W require each block to read a larger input region, and store it in a larger shared memory region, which reduces the potential for high occupancy.

Filterbank correlation represents a difficult code-optimization problem because of the number of interacting constraints: a large value of P reduces arithmetic density, but a small value of P requires that each thread do more work, and makes global memory access latency harder to hide. (XXX: is there a better reason?)

The kernel used to perform filterbank correlation was parametrized pages 111–122. IEEE, 2005.

by 10 parameters, some of which were binary and others of which were integer-valued. The full configuration space included 12000 kernels. XXX: how to make sense of these parameters without code listing?

XXX: What to all the parameters mean? XXX: Point to github / GCG for full code listing.

XXX: how many configurations are in the configuration space

XXX: what was the reference kernel, and how was it chosen?

XXX: analysis of READ traffic, WRITE traffic, DATA re-use and CACHEing strategies. FLOPS per write etc.

XXX: Why not use FFT: convolution in spectral domain?

XXXX: Talk about memory transfer requirements vs. speed vs. copy...

7. FUTURE WORK

more features (this is the hardest, and most manual part, may require domain-specific knowledge, but not necessarily, e.g. microbenchmarking) * features derived from microbenchmarking [Wong et al., 2010] * features derived from performance limiter analysis [Micikevicius] * more instrumentation * more data-layout * more threading strategies to lower latency

online distributed auto-tuning with learning, especially w/ HT-based experiments ? our method is indeed very well suited for parallel exploring of multiple configurations for learning (in a cluster, e.g. in HT context) this may get us *more* than linear speed ups for better accuracy, especially when more features will be present (see above)

interpret the learned model and confront the performance against theoretical (or approximated) models of mem/com/compute, arithmetic intensity, etc. to understand performance limiters, the hardware, etc.

References

- K. Asanovic, R. Bodik, B. Catanzaro, J. Gebis, P. Husbands, K. Keutzer, D. Patterson, W. Plishker, J. Shalf, S. Williams, et al. The landscape of parallel computing research: A view from Berkeley. *EECS Department University of California Berkeley Tech Rep UCBECS2006183*, 18(UCB/EECS-2006-183), 2006.
- J. Bilmes, K. Asanovic, C. Chin, and J. Demmel. Optimizing matrix multiply using phipac: a portable, high-performance, ansi c coding methodology. In *Proceedings of the 11th international conference on Supercomputing*, pages 340–347. ACM, 1997.
- J. Cavazos. Intelligent compilers. In *Cluster Computing, 2008 IEEE International Conference on*, pages 360–368. IEEE, 2008.
- J. Cavazos, G. Fursin, F. Agakov, E. Bonilla, M. O'Boyle, and O. Temam. Rapidly selecting good compiler optimizations using performance counters. In *Code Generation and Optimization, 2007. CGO'07. International Symposium on*, pages 185–197. IEEE, 2007.
- C. Chen, J. Chame, and M. Hall. Combining models and guided empirical search to optimize for multiple levels of the memory hierarchy. In *Code Generation and Optimization, 2005. CGO 2005. International Symposium on*,

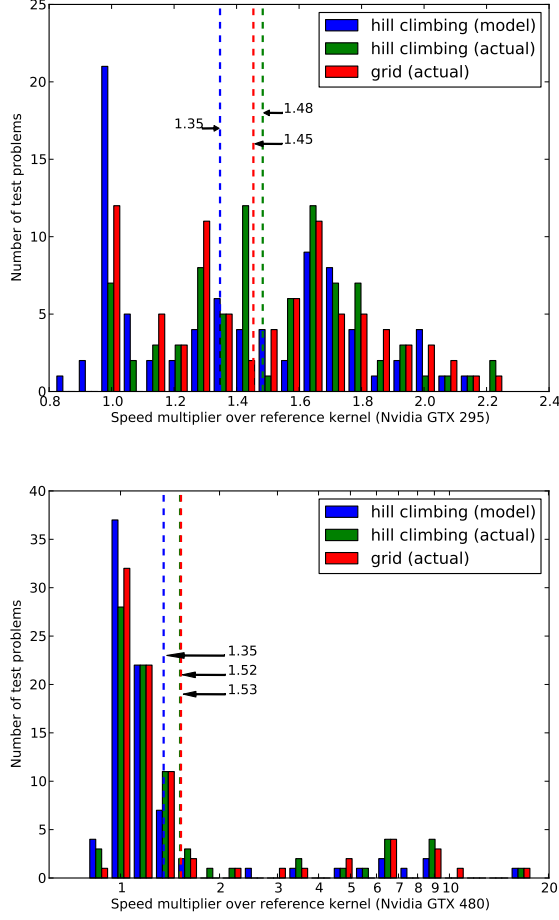


Figure 8: Speedup of various auto-tuning strategies over the reference implementation. Each strategy was evaluated the same set of 83 problem configurations, which was disjoint from the problem configurations used to build the model for model-based hill-climbing. The vertical dashed lines are positioned at the geometric mean speedup (XXX?) of each strategy. The grid and hill-climbing approaches tested 73 and 75 configurations respectively, and took 45?? XXX seconds on average. The model-based approach tested 0 configuration evaluations, and took 0.05 seconds on average, even with a naive Python implementation.

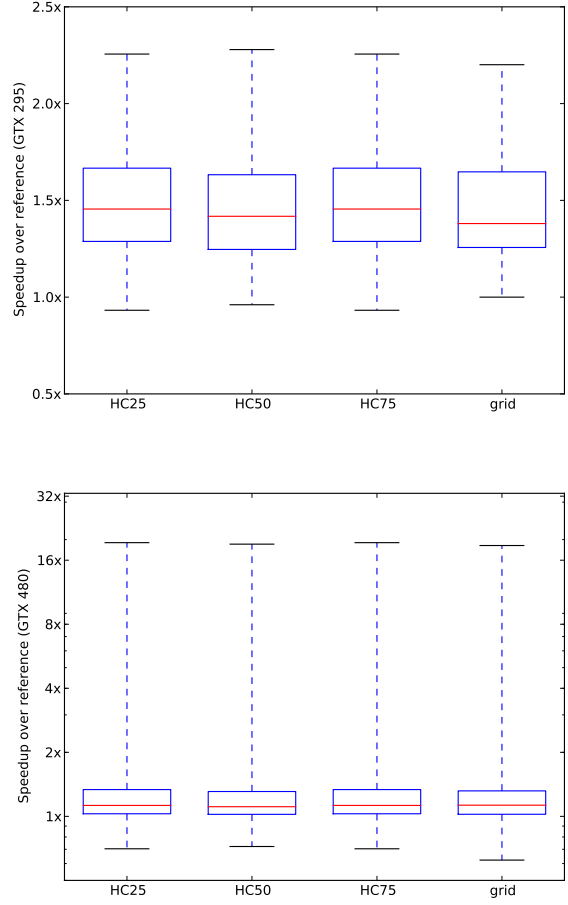


Figure 9: The speedup of hill-climbing (HC) and grid algorithms for empirical autotuning

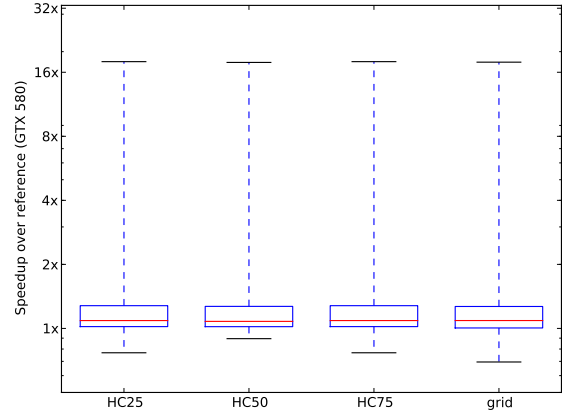


Figure 10: The speedup of hill-climbing (HC) and grid algorithms for empirical autotuning

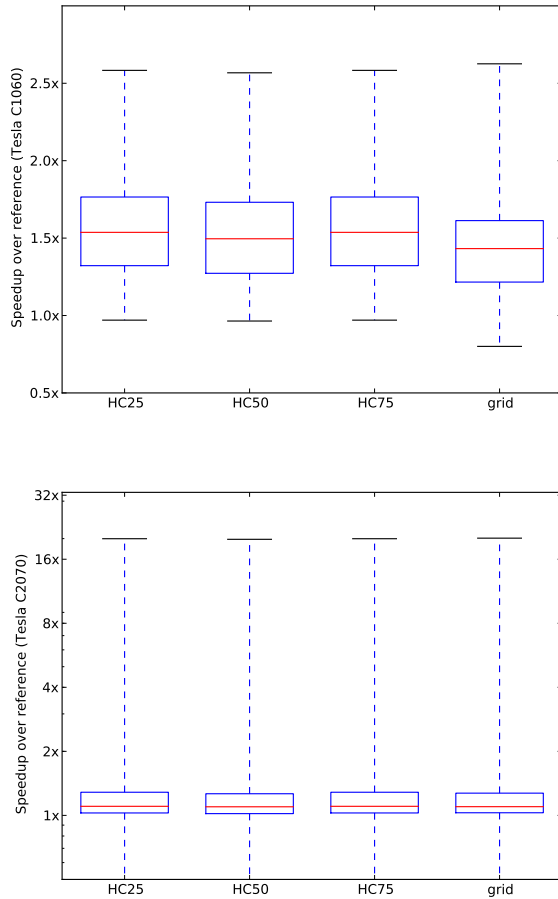


Figure 11: The speedup of hill-climbing (HC) and grid algorithms for empirical autotuning

- R. Clint Whaley, A. Petitet, and J. Dongarra. Automated empirical optimizations of software and the atlas project. *Parallel Computing*, 27(1-2):3–35, 2001.
- K. Cooper, A. Grosul, T. Harvey, S. Reeves, D. Subramanian, L. Torczon, and T. Waterman. Acme: adaptive compilation made efficient. In *ACM SIGPLAN Notices*, volume 40, pages 69–77. ACM, 2005.
- K. Datta. *Auto-tuning Stencil Codes for Cache-Based Multicore Platforms*. PhD thesis, Computer Science Division, U.C. Berkeley, Dec. 2009.
- K. Datta, M. Murphy, V. Volkov, S. Williams, J. Carter, L. Oliker, D. Patterson, J. Shalf, and K. Yelick. Stencil computation optimization and auto-tuning on state-of-the-art multicore architectures. In *Proceedings of the 2008 ACM/IEEE conference on Supercomputing*, page 4. IEEE Press, 2008.
- J. Demmel, J. Dongarra, V. Eijkhout, E. Fuentes, A. Petitet, R. Vuduc, R. Whaley, and K. Yelick. Self-adapting linear algebra algorithms and software. *Proceedings of the IEEE*, 93(2):293–312, 2005.
- B. Franke, M. O’Boyle, J. Thomson, and G. Fursin. Probabilistic source-level optimisation of embedded programs. In *ACM SIGPLAN Notices*, volume 40, pages 78–86. ACM, 2005.
- G. Fursin, C. Miranda, O. Temam, M. Namolaru, E. Yom-Tov, A. Zaks, B. Mendelson, E. Bonilla, J. Thomson, H. Leather, et al. Milepost gcc: machine learning based research compiler. 2008.
- A. Ganapathi, K. Datta, A. Fox, and D. Patterson. A case for machine learning to optimize multicore performance. In *Proceedings of the First USENIX conference on Hot topics in parallelism*, pages 1–1. USENIX Association, 2009.
- S. Grauer-Gray and J. Cavazos. Optimizing and auto-tuning belief propagation on the gpu. *Languages and Compilers for Parallel Computing*, pages 121–135, 2011.
- 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*, pages 1–11. IEEE, 2009.
- F. Hutter, Y. Hamadi, H. Hoos, and K. Leyton-Brown. Performance prediction and automated tuning of randomized and parametric algorithms. *Principles and Practice of Constraint Programming-CP 2006*, pages 213–228, 2006.
- S. Kamil, C. Chan, S. Williams, L. Oliker, J. Shalf, M. Howison, and P. E. W. Bethel. A generalized framework for auto-tuning stencil computations. In *Cray User Group Conference*, May 2009.
- S. Kamil, C. Chan, L. Oliker, J. Shalf, and S. Williams. An auto-tuning framework for parallel multicore stencil computations. In *Parallel & Distributed Processing (IPDPS), 2010 IEEE International Symposium on*, pages 1–12. IEEE, 2010.

- A. Klöckner, N. Pinto, Y. Lee, B. C. Catanzaro, P. Ivanov, and A. Fasih. PyCUDA: GPU run-time code generation for high-performance computing. arXiv, 2009.
- A. Klöckner, N. Pinto, Y. Lee, B. Catanzaro, P. Ivanov, and A. Fasih. Pycuda and pyopencl: A scripting-based approach to gpu run-time code generation. *Parallel Computing*, 2011.
- P. Kulkarni, S. Hines, J. Hiser, D. Whalley, J. Davidson, and D. Jones. Fast searches for effective optimization phase sequences. In *ACM SIGPLAN Notices*, volume 39, pages 171–182. ACM, 2004.
- X. Li, M. Garzarán, and D. Padua. A dynamically tuned sorting library. In *Code Generation and Optimization, 2004. CGO 2004. International Symposium on*, pages 111–122. IEEE, 2004.
- Y. Li, J. Dongarra, and S. Tomov. A note on auto-tuning gemm for gpus. *Computational Science–ICCS 2009*, pages 884–892, 2009.
- Y. Liu, E. Zhang, and X. Shen. A cross-input adaptive framework for gpu program optimizations. In *Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on*, pages 1–10. IEEE, 2009.
- P. Micikevicius. Analysis-driven optimization. In *GPU Technology Conference*. NVIDIA, 2010.
- A. Monsifrot, F. Bodin, and R. Quiniou. A machine learning approach to automatic production of compiler heuristics. *Artificial Intelligence: Methodology, Systems, and Applications*, pages 389–409, 2002.
- J. Nickolls, I. Buck, M. Garland, and K. Skadron. Scalable parallel programming with cuda. *Queue*, 6(2):40–53, 2008.
- A. Nukada and S. Matsuoka. Auto-tuning 3-d fft library for cuda gpus. In *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, page 30. ACM, 2009.
- NVIDIA. Compute unified device architecture (cuda) programming guide. 2011.
- E. Park, S. Kulkarni, and J. Cavazos. An evaluation of different modeling techniques for iterative compilation. In *Proceedings of the 14th international conference on Compilers, architectures and synthesis for embedded systems*, pages 65–74. ACM, 2011.
- N. Pinto and D. D. Cox. GPU metaprogramming: A case study in biologically inspired machine vision. In *GPU Computing Gems*, volume 2. Morgan Kauffmann, 2012. to appear.
- S. Rahman, J. Guo, and Q. Yi. Automated empirical tuning of scientific codes for performance and power consumption. In *Proceedings of the 6th International Conference on High Performance and Embedded Architectures and Compilers*, pages 107–116. ACM, 2011.
- S. Ryoo, C. Rodrigues, S. Stone, S. Baghsorkhi, S. Ueng, J. Stratton, and W. Hwu. Program optimization space pruning for a multithreaded gpu. In *Proceedings of the 6th annual IEEE/ACM international symposium on Code generation and optimization*, pages 195–204. ACM, 2008.
- D. Schaa and D. Kaeli. Exploring the multiple-gpu design space. In *Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on*, pages 1–12. IEEE, 2009.
- M. Stephenson, S. Amarasinghe, M. Martin, and U. O’Reilly. Meta optimization: Improving compiler heuristics with machine learning. *ACM SIGPLAN Notices*, 38(5):77–90, 2003.
- R. Vuduc, J. Demmel, and J. Bilmes. Statistical models for automatic performance tuning. *Computational Science—ICCS 2001*, pages 117–126, 2001.
- R. Vuduc, J. Demmel, and K. Yelick. Oski: A library of automatically tuned sparse matrix kernels. In *Journal of Physics: Conference Series*, volume 16, page 521. IOP Publishing, 2005.
- S. Williams. *Auto-tuning performance on multicore computers*. ProQuest, 2008.
- H. Wong, M. Papadopoulou, M. Sadooghi-Alvandi, and A. Moshovos. Demystifying gpu microarchitecture through microbenchmarking. In *Performance Analysis of Systems & Software (ISPASS), 2010 IEEE International Symposium on*, pages 235–246. IEEE, 2010.
- K. Yotov, X. Li, G. Ren, M. Cibulskis, G. DeJong, M. Garzaran, D. Padua, K. Pingali, P. Stodghill, and P. Wu. A comparison of empirical and model-driven optimization. In *ACM SIGPLAN Notices*, volume 38, pages 63–76. ACM, 2003.