

Analytical Modeling is Enough for High Performance BLIS

TZE MENG LOW, The University of Texas at Austin
FRANCISCO D. IGUAL, Universidad Complutense de Madrid
TYLER M. SMITH, The University of Texas at Austin
ENRIQUE S. QUINTANA-ORTI, Universidad Jaume I

We show how the BLAS-like Library Instantiation Software (BLIS) framework, which provides a more detailed layering of the GotoBLAS (now maintained as OpenBLAS) implementation, allows one to analytically determine optimal tuning parameters for high-end instantiations of the matrix-matrix multiplication. This is of both practical and scientific importance, as it greatly reduces the development effort required for the implementation of the level-3 BLAS while also advancing our understanding of how hierarchically layered memories interact with high performance software. This allows the community to move on from valuable engineering solutions (empirically autotuning) to scientific understanding (analytical insight).

Categories and Subject Descriptors: G.4 [**Mathematical Software**]: Efficiency

General Terms: Algorithms, Performance

Additional Key Words and Phrases: linear algebra, libraries, high-performance, matrix multiplication, analytical modeling

ACM Reference Format:

ACM V, N, Article A (January YYYY), 19 pages.
DOI: <http://dx.doi.org/10.1145/0000000.0000000>

1. INTRODUCTION

The field of dense linear algebra (DLA) was among the first to realize, understand, and promote that standardizing (*de facto*) an interface to fundamental primitives supports portable high performance. For almost four decades, those primitives have been the *Basic Linear Algebra Subprograms* (BLAS) [Lawson et al. 1979; Dongarra et al. 1988; Dongarra et al. 1990]. The expectation was and still remains that some expert will provide to the community a high performance implementation of the BLAS every time a new architecture reaches the market. For this purpose, many hardware vendors currently enroll a numerical libraries group and distribute well-maintained libraries (e.g., Intel’s MKL [Intel 2015], AMD’s ACML [AMD 2015], and IBM’s ESSL [IBM 2015] libraries), while over the years we have enjoyed a number of alternative Open Source solutions (e.g., GotoBLAS [Goto and van de Geijn 2008b; 2008a], OpenBLAS [OpenBLAS 2015], ATLAS [Whaley and Dongarra 1998]). Nevertheless, these solutions all require considerable effort in time and labor for each supported target architecture.

In order to reduce the exertion of developing high performance implementations, ATLAS and its predecessor PHiPAC [Bilmes et al. 1997b] introduced autotuning to the field of high performance computing. The fundamental rationale behind these two libraries is that optimizing is too difficult and time-consuming for a human expert and that autogeneration in combination with autotuning is the answer.

Permission to make digital or hard copies of part or all 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 show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.
© YYYY ACM 0000-0000/YYYY/01-ARTA \$15.00
DOI: <http://dx.doi.org/10.1145/0000000.0000000>

The work in [Yotov et al. 2005] shows that the ATLAS approach to optimizing the general matrix-matrix product (GEMM) can be, by and large, analytically modelled. Thus, it reveals that autotuning is unnecessary for the operation that has been touted by the autotuning community as *the* example of the success of autotuning. The problem with that work ([Yotov et al. 2005]) is that the ATLAS approach to optimizing GEMM had been previously shown to be suboptimal, first in theory [Gunnels et al. 2001] and then in practice [Goto and van de Geijn 2008b]. Furthermore, ATLAS leverages an inner kernel optimized by a human expert, which still involves a substantial manual encoding. Precisely, the presence of this black-box kernel makes analytical modeling problematic for ATLAS, since some of the tuning parameters are hidden inside.

GotoBLAS (now supported as OpenBLAS) also builds on a substantial inner kernel that is implemented by a human, turning the analysis equally difficult without a complete understanding of the complex details embedded into the inner kernel. (Until recently, no paper that detailed the intricacies of the GotoBLAS inner kernel existed.)

BLIS (*BLAS-like Library Instantiation Software*) [Van Zee and van de Geijn 2014] is a new framework for instantiating the BLAS. A key benefit of BLIS is that it is a productivity multiplier, as the framework allows developers to rapidly unfold new high-performance implementations of BLAS and BLAS-like operations on current and future architectures [Van Zee et al. 2014]. From the implementation point of view, BLIS modularizes the approach underlying GotoBLAS2 (now adopted also by many other implementations) to isolate a *micro-kernel* that performs GEMM upon which all the level-3 BLAS functionality can be built. Thus, for a given architecture, the programmer only needs to develop an efficient micro-kernel for BLIS, and adjust a few key parameter values¹ that optimize the loops around the micro-kernel, to automatically instantiate all the level-3 BLAS. Importantly, BLIS features a layering that naturally exposes the parameters that need to be tuned.

While BLIS simplifies the implementation of BLAS(-like) operations, a critical step to optimize BLIS for a target architecture is for the developer to identify the specific parameter values for both the micro-kernel and the loops around it. Conventional wisdom would dictate that empirical search must be applied. In this paper, we adopt an alternative model-driven approach to analytically identify the optimal parameter values for BLIS. In particular, we apply an analysis of the algorithm and architecture, similar to that performed in [Yotov et al. 2005], to determine the data usage pattern of the algorithm ingrained within the BLIS framework, and to build an analytical model based on hardware features in modern processor architectures.

The specific contributions of this work include:

- **An analysis of the best known algorithm.** We address the algorithm underlying BLIS (and therefore GotoBLAS and OpenBLAS) instead of the algorithm in ATLAS. This is relevant because it has been shown that, on most architectures, a hand-coded BLIS/GotoBLAS/OpenBLAS implementation (almost) always yields higher performance than an implementation automatically generated via ATLAS [Van Zee et al. 2014], even if ATLAS starts with a hand-coded inner kernel.
- **A more modern model.** We accommodate a processor model that is more representative of modern architectures than that adopted in [Yotov et al. 2005]. Concretely, our model includes hardware features such as a vector register file and a memory hierarchy with multiple levels of set-associative data caches. This considerably improves upon the analytical model in [Yotov et al. 2005] which considered one level only of fully associative cache and ignored vector instructions.

¹The micro-kernel itself is characterized by three additional parameters, which the programmer has to consider when implementing it.

- **A more comprehensive model.** Our analytical model is more comprehensive in that it includes the parameter values that also characterize the micro-kernel. These are values that a developer would otherwise have to determine empirically. This is important because the best implementations provided by ATLAS often involve loops around a hand-coded (black-box) kernel. Since the internals of the hand-coded kernel are not known, the model in [Yotov et al. 2005] was not able to identify the globally optimal parameter values.
- **Competitive with expert-tuned implementations.** Unlike previous work comparing the performance attained against auto-generated implementations in C obtained by ATLAS [Yotov et al. 2005; Kelefouras et al. 2014], which are typically not competitive with those that use hand-coded inner kernels (so-called “ATLAS unleashed” in [Yotov et al. 2005]), this paper shows that analytically-obtained parameter values can yield performance that is competitive with manually-tuned implementations that achieve among best-in-class performance. Hence, this paper provides an answer to the question posed in [Yotov et al. 2005] —i.e. *whether empirical search is really necessary in this context*—, by demonstrating that *analytical modeling suffices for high performance BLIS implementations* which then shows that careful layering combined with analytical modeling is competitive with GotoBLAS, OpenBLAS, and vendor BLAS (since other papers show BLIS to be competitive with all these implementations).

In this paper, we restrict ourselves to the case of a single-threaded implementation of GEMM in double precision. How to determine optimal parameter values for a multithreaded implementation is orthogonal to this paper, and is at least partially answered in [Smith et al. 2014].

2. APPROACHES TO IDENTIFY THE OPTIMAL PARAMETER VALUES FOR GEMM

It has been argued that empirical search is the only way to obtain highly optimized implementations for DLA operations [Demmel et al. 2005; Bilmes et al. 1997a; Whaley and Dongarra 1998], and an increasing number of recent projects (Build-To-Order BLAS [Belter et al. 2010] and AuGEM [Wang et al. 2013]) now adopt empirical search to identify optimal parameter values for DLA algorithms.

The problem with empirical-based approaches is that they unleash a walloping search space, due to the combination of a large number of possible values for a substantial set of parameters. Therefore, even with the help of heuristics, generating, executing, and timing all variants of the program, with their unique combinations of parameter values, requires a considerable amount of time (sometimes measured in days). In addition, empirical search has to be performed on the target machine architecture. This implies that a high performance implementation of BLAS for a new machine architecture cannot be developed until the programmer is granted access to the target machine, and then it still requires a significant amount of time. This is especially critical on state-of-the-art embedded architectures, where cross-compiling processes are necessary, greatly complicating or even invalidating in practice the application of empirical-based approaches.

A refinement of the empirical approach is iterative optimization [Knijnenburg et al. 2002; Kisuki et al. 2000; 2000]. As the name suggests, instead of a single run that explores the entire search space, multiple optimization passes are executed. There exist different techniques to perform iterative optimization, but they all share the following similarities: An initial sampling run is performed while a monitoring tool captures specific performance information such as the number of cache misses. Using the captured information and the application of some heuristics, the search space is trimmed/refined, and new program variants (with new sets of parameter values) are generated for further exploration in subsequent passes.

Although iterative optimization represents an improvement over empirical search, the information obtained, e.g. from the performance counters, may not be sufficient to limit the

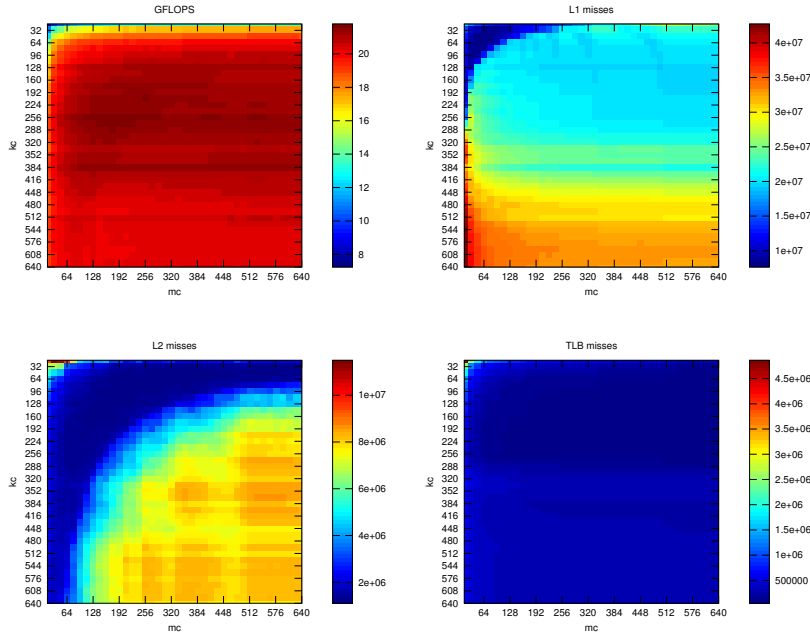


Fig. 1: Performance metrics when empirically testing different parameter values. Clockwise from top-left: GFLOPS, and misses in the L1, TLB and L2 caches. A heuristic for reducing the search space is to find the parameter values that minimize the cache misses and TLB misses. From the above data, this still leaves a relatively large space that needs to be searched.

search space. Consider the graphs in Figure 1, which report the GFLOPS (i.e., billions of floating-point arithmetic operations, or flops, per second), and L1, L2 and TLB cache misses observed for the execution of a tuned implementation of `dgemm` from BLIS on a single core of an Intel Xeon E3-1220 processor (3.1 GHz), while varying two of the optimization parameters only (m_c and k_c , to be introduced later). The top-left plot illustrates the challenge for empirical search: generate and evaluate all different combinations of an exponentially growing search space. (BLIS, e.g., needs to optimize 5–6 different parameters, which is a reduced quantity when compared with other solutions including ATLAS.) In contrast to this, iterative optimization runs a coarse-grained grid of experiments that can help to identify contour lines which roughly delimit the area which simultaneously minimizes the amount of cache misses for all three caches (L1, L2, and TLB). However, even armed with that information, we can observe from the three cache miss plots that such region still comprises an enormous number of cases that nonetheless have to be individually evaluated to finally reveal the optimal combination for the architecture ($m_c=96$ and $k_c=256$ for this particular case).

A third approach to identify the optimal parameter values is to use analytical models as advocated by some in the compiler community. Concretely, simple analytical models have been previously developed in the compiler domain in order to determine optimizing parameter values for tile sizes and unrolling factors [Bacon et al. 1994; Yotov et al. 2005; Kelefouras et al. 2014]. These models are based on the data access pattern and exploit a fair knowledge

of the hardware features commonly found in modern processor architectures. Previous work has shown that parameter values derived analytically can deliver performance rates that are similar to those attained using values determined through empirical search [Yotov et al. 2005]. More importantly, the time it takes to analytically identify the optimal parameter combination is often a fraction of that required by empirical search. Let us distinguish this from autotuning by calling it auto-*fine*-tuning. We adopt this in our quest for the parameter values that optimize the algorithm in BLIS, with the aim to make even auto-*fine*-tuning unnecessary for many architectures.

3. EXPERT IMPLEMENTATION OF GEMM

The GEMM algorithm embedded within the BLIS framework is described in [Van Zee and van de Geijn 2014], and the same approach is instantiated in BLAS libraries optimized by DLA experts such as GotoBLAS [Goto and van de Geijn 2008b; 2008a] and its successor, OpenBLAS [OpenBLAS 2015]. Importantly, BLIS exposes three loops within the inner kernel used by the GotoBLAS and OpenBLAS, which then facilitates the analytical determination of the parameters. For completeness, we provide a brief overview of the algorithm next. In addition, we identify the data usage pattern and highlight the parameters that characterize the algorithm.

3.1. An expert's implementation of GEMM

Without loss of generality, we consider the special case of the matrix-matrix multiplication, $C := AB + C$, where the sizes of A , B , C are $m \times k$, $k \times n$, and $m \times n$, respectively. In the following elaboration, we will consider different partitionings of the m , n and k -dimensions of the problem. For simplicity, when considering a (vertical or horizontal) partitioning of one of the problem dimensions, say q , into panels/blocks of size (length or width) r , we will assume that q is an integer multiple of r .

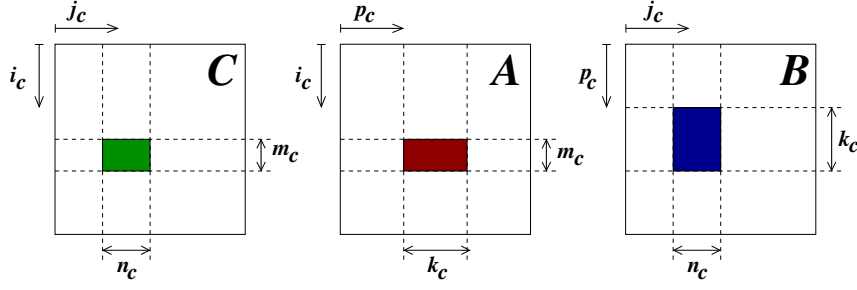
BLIS implements GEMM as three nested loops around a *macro-kernel*², plus two packing routines. The macro-kernel is implemented in terms of two additional loops around a *micro-kernel*. The micro-kernel is a loop around a rank-1 (i.e., outer product) update, and it is typically implemented in assembly or with vector intrinsics. In BLIS, the remaining five loops are implemented in C.

Pseudo-code for the GEMM algorithm is given in Figure 2. The outermost loop (Loop 1, indexed in the figure by j_c) traverses the n -dimension of the problem, partitioning both C and B into column panels of width n_c . The second loop (indexed by p_c) processes the k -dimension, partitioning A into column panels of width k_c and the current panel of B into blocks of length k_c . The third loop (indexed by i_c) iterates over the m -dimension, yielding partitionings of the current column panels of C and A into blocks of length m_c . The following loops comprise the macro-kernel. Let us define $C_c = C(i_c : i_c + m_c - 1, j_c : j_c + n_c - 1)$. Loop 4 (indexed by j_r) traverses the n -dimension, partitioning C_c and the packed block B_c into *micro-panels* of width n_r . Loop 5 (indexed by i_r) then progresses along the m -dimension, partitioning the current micro-panel of C_c into *micro-tiles* of length m_r and the packed block A_c into micro-panels of the same length. The innermost loop (Loop 6, indexed by p_r), inside the micro-kernel, computes the product of a micro-panel of A_c and micro-panel of B_c as a sequence of k_c rank-1 updates, accumulating the result to a micro-tile of C_c , which can be described mathematically as

$$C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1) += A_c(i_r : i_r + m_r - 1, 0 : k_c - 1) \cdot B_c(0 : k_c - 1, j_r : j_r + n_r - 1).$$

At this point it is worth emphasizing the difference between C_c and A_c, B_c . While the former is just a notation artifact, introduced to ease the presentation of the algorithm, the latter two correspond to actual buffers that are involved in data copies.

²The macro-kernel is also known as the inner kernel in GotoBLAS.



```

Loop 1  for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
Loop 2  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
         $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$  // Pack into  $B_c$ 
Loop 3  for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
         $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$  // Pack into  $A_c$ 
Loop 4  for  $j_r = 0, \dots, n_c - 1$  in steps of  $n_r$  // Macro-kernel
Loop 5  for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
Loop 6  for  $p_r = 0, \dots, k_c - 1$  in steps of 1 // Micro-kernel
         $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
         $+= A_c(i_r : i_r + m_r - 1, p_r)$ 
         $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 
        endfor
    endfor
endfor
endfor
endfor
endfor

```

Fig. 2: High performance implementation of GEMM by an expert.

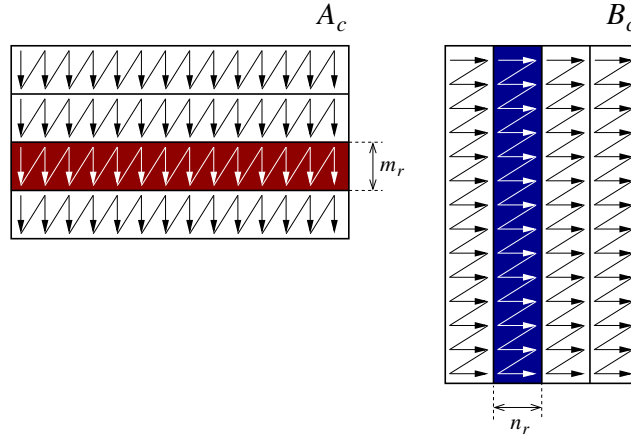


Fig. 3: Packing in the BLIS and GotoBLAS implementations of GEMM.

3.2. The importance of packing

It is generally known that accessing consecutive memory locations (also known as in-stride access) is usually faster than non-consecutive memory accesses. By packing elements of A in Loop 3 and B in Loop 2 in a special manner into A_c and B_c , respectively, we can ensure that the elements of these two matrices will be read with in-stride access inside the micro-kernel.

Loop	Reused data	Size of reused data	Reuse factor
6	Micro-tile, C_r	$m_r \times n_r$	k_c
5	Micro-panel, B_r	$k_c \times n_r$	m_c/m_r
4	Packed block A_c	$m_c \times k_c$	n_c/n_r
3	Packed block B_c	$k_c \times n_c$	m/m_c
2	Column panel of C	$m \times n_c$	k/k_c
1	Matrix A	$m \times k$	n/n_c

Table I: Analysis of data reuse performed in the different loops of the BLIS implementation of GEMM.

Concretely, each $m_c \times k_c$ block of A is packed into A_c . Elements are organized as row micro-panels of size $m_r \times k_c$, and within each micro-panel of A_c , the elements are stored in column-major order. Each $k_c \times n_c$ block of B is packed into B_c . In this case, elements are packed into column micro-panels of size $k_c \times n_r$, and each column micro-panel is stored into row-major order. The reorganization of the entries of A and B into blocks of A_c and B_c with the packing layout illustrated in Figure 3 ensures that these elements are accessed with unit stride when used to update a micro-tile of C . Packing A_c and B_c also has the additional benefit of aligning data from A_c and B_c to cache lines boundaries and page boundaries. This enables to use instructions for accessing aligned data, which are typically faster than their non-aligned counterparts.

A third benefit of packing is that data from A_c and B_c are preloaded into certain cache levels of the memory hierarchy. This reduces the time required to access the elements of A_c and B_c when using them to update a micro-tile of C . Since A_c is packed in Loop 3, after B_c has been packed in Loop 2, the elements of A_c will likely be higher up in the memory hierarchy (i.e. closer to the registers) than those of B_c . By carefully picking the sizes for A_c and B_c , the exact location of A_c and B_c within the memory hierarchy can be determined.

3.3. Data usage pattern

Consider Loop 6. This innermost loop updates a micro-tile of C , say C_r , via a series of k_c rank-1 updates involving a micro-panel A_r , from the packed matrix A_c , and a micro-panel B_r , from the packed matrix B_c . At each iteration, m_r elements from A_r and n_r elements from B_r are multiplied to compute $m_r n_r$ intermediate results that will be accumulated into the micro-tile of C . Between two iterations of Loop 6, different elements from A_r and B_r are used, but the results from the rank-1 updates are accumulated into the same micro-tile. Hence, the micro-tile C_r is the data reused between iterations (*Reused data*) in Loop 6. In addition, because k_c rank-1 updates are performed each time Loop 6 is executed, C_r is said to have a reuse factor of k_c times. Since the columns of A_r and the rows of B_r are used exactly once each time through Loop 6, there is no reuse of these elements.

Identifying pieces of data that are reused is important because it tells us which data portions should be kept in cache in order to allow fast access and reduce cache trashing. An analysis similar to that with Loop 6, performed on the remaining loops to identify the data that is reused in each one of them, the size of the reused data, and the number of times the data is being reused, yields the results summarized in Table I.

Notice that the size of the reused data becomes smaller as the loop depth increases. This nicely matches the memory hierarchy, where faster memory layers (caches) feature lower capacity than slower caches. As such, GEMM maps smaller reused data to faster layers in the memory hierarchy, as depicted in Figure 4.

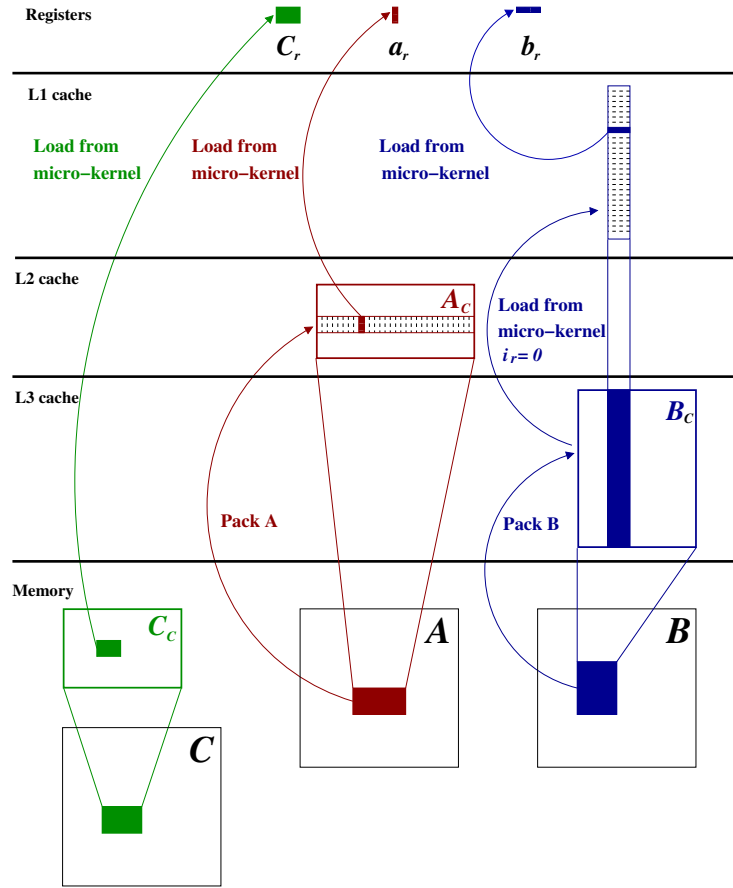


Fig. 4: Packing in the BLIS and GotoBLAS implementations of GEMM.

3.4. Characteristic parameters that drive performance

As described in the previous subsection and captured in Table I, the GEMM algorithm that underlies BLIS is characterized by the following five parameters:

$$m_c, k_c, n_c, m_r \text{ and } n_r,$$

which correspond to the block/panel/tile size for each of the loops around the micro-kernel. In addition, this set of parameters also determines the dimension of the reused data and the reuse factor of each of the loops.

Loop 6 is defined by three of the five parameters (m_r, n_r and k_c), while Loops 3, 4 and 5 are characterized by four of the five parameters. In addition, three of the four parameters that characterize Loops 3, 4 and 5 carry beyond to the next inner loop as well. Now, when the parameters that characterize Loop 5 have been identified, the only unknown parameter for Loop 4 is n_c . This observation suggests that the parameters should be identified from the innermost loop outwards. In the next section, we leverage this observation to optimize this collection of parameters analytically.

4. MAPPING ALGORITHM TO ARCHITECTURE

An effective mapping of the micro-kernel and the loops around it to the architecture is fundamental to attain high performance with the BLIS approach. This means that we need to identify values for the characteristic parameters that are tuned for the target machine architecture. As discussed in the previous section, we start by mapping the micro-kernel to the architecture and work our way outwards in order to identify the optimal values for the different characteristic parameters.

4.1. Architecture model

To develop an optimization model for mapping the GEMM algorithm onto an architecture, we first need to define a model for our target architecture. We make the following assumptions regarding our hypothetical processor:

- **Load/store architecture and vector registers.** Data have to be loaded into the processor registers before computation can be performed with them. A total of N_{REG} vector registers exist, where each can hold N_{VEC} elements (floating-point numbers) of size S_{DATA} . (In a machine with no vector registers, $N_{\text{VEC}} = 1$.) In addition, we assume that memory instructions can be issued in parallel with floating-point arithmetic instructions.
- **Vector instructions.** The processor floating-point arithmetic units have a throughput of N_{VFMA} vector (or SIMD) *fused multiply-add instructions* (VFMA) per clock cycle (i.e., $2N_{\text{VEC}}N_{\text{VFMA}}$ flops per cycle). A single VFMA instruction combines N_{VEC} regular (non-vector) instructions and produces N_{VEC} scalar results. Furthermore, each VFMA has a latency of L_{VFMA} cycles, which is the minimum number of cycles between the issuance of two dependent consecutive VFMA instructions. On an architecture without VFMA instructions, L_{VFMA} is computed by adding the latencies for a multiply and an addition instruction.
- **Caches.** All data caches are set-associative, and each cache level L_i is characterized by four parameters as follows:
 - C_{L_i} : size of cache line,
 - W_{L_i} : associativity degree,
 - S_{L_i} : size, and
 - N_{L_i} : number of sets,

where $S_{L_i} = N_{L_i}C_{L_i}W_{L_i}$. A fully-associative cache can be modelled by setting $N_{L_i} = 1$ and $W_{L_i} = S_{L_i}/C_{L_i}$.

The replacement policy for all caches is least-recently-used (LRU) [Hennessy and Patterson 2003]. For simplicity, we assume that the size of a cache line is the same for all cache levels.

4.2. Parameters for the body of the inner-most loop: m_r and n_r

Recall that the micro-kernel is characterized by three parameters: m_r , n_r and k_c , where the former two determine the size of the micro-tile of C that is reused across every iteration. In addition, these two parameters also define the number of elements of A_r and B_r that are involved in the rank-1 update at each iteration of Loop 6. In this subsection, we discuss how m_r and n_r can be identified analytically.

4.2.1. Strategy for finding m_r and n_r . The main strategy behind our approach to identify m_r, n_r is to choose them “large enough” so that no stalls due to the combination of dependencies and latencies are introduced in the floating-point pipelines during the repeated updates of the micro-tile C_r . In addition, the smallest values of m_r, n_r which satisfy this condition should be chosen because this implies that C_r occupies the minimum number of registers, releasing more registers for the entries of A_r, B_r . This in turn enables larger loop unrolling factors and more aggressive data preloading [Hennessy and Patterson 2003] to reduce loop overhead.

4.2.2. Latency of instructions. Consider the k_c successive updates of C_r occurring in Loop 6. In each of the k_c iterations, each element of C_r is updated by accumulating the appropriate intermediate result to it (obtained from multiplying corresponding elements of A_r and B_r). This update can be achieved with a VFMA instruction for each element of C_r in each iteration.

Now, recall that there is a latency of L_{VFMA} cycles between the issuance of a VFMA instruction and the issuance of another dependent VFMA instruction. This implies that L_{VFMA} cycles must lapse between two successive updates to the same element of C_r . During those L_{VFMA} cycles, a minimum of $N_{VFMA} L_{VFMA}$ VFMA instructions must be issued without introducing a stall in the floating-point pipelines. This means that, in order to avoid introducing stalls in the pipelines, at least $N_{VFMA} L_{VFMA} N_{VEC}$ output elements must be computed. Therefore, the size of the micro-tile C_r must satisfy:

$$m_r n_r \geq N_{VEC} L_{VFMA} N_{VFMA}. \quad (1)$$

4.2.3. Maximizing Register Usage. Ideally, m_r should be equal to n_r as this maximizes the ratio of computation to data movement during the update of C_r in Loop 6 ($2m_r n_r k_c$ flops vs $2m_r n_r + m_r k_c + n_r k_c$ memory operations). However, there are two reasons why this is not always possible in practice. First, m_r and n_r have to be integer values as they correspond to the dimensions of the micro-tile C_r . Second, it is desirable for either m_r or n_r to be integer multiples of N_{VEC} . As the number of registers is limited in any architecture, it is necessary to maximize their use, and choosing m_r (or n_r) to be an integer multiple of N_{VEC} will ensure that the registers are filled with elements from A_r , B_r and C_r such that data not used in a particular iteration is not kept in registers.

Therefore, in order to satisfy this criterion as well as (1), m_r and n_r are computed as follows:

$$m_r = \left\lceil \frac{\sqrt{N_{VEC} L_{VFMA} N_{VFMA}}}{N_{VEC}} \right\rceil N_{VEC} \quad (2)$$

and

$$n_r = \left\lceil \frac{N_{VEC} L_{VFMA} N_{VFMA}}{m_r} \right\rceil. \quad (3)$$

The astute reader will recognize that one could have computed n_r before computing m_r . Doing this is analogous to swapping the values of m_r and n_r , and also yields a pair of values that avoid the introduction of stalls in the pipeline. We choose the values of m_r and n_r that maximize the value of k_c , which will be discussed in the following section.

4.3. Parameters for the remaining loops: k_c , m_c and n_c

After analytically deriving optimal values for m_r, n_r , we discuss next how to proceed for k_c, m_c and n_c .

4.3.1. Strategy for identifying k_c, m_c and n_c . We recall that k_c, m_c and n_c (together with n_r) define the dimensions of the reused data for Loops 3 to 5; see Table I. Since the reused blocks are mapped to the different layers of the memory hierarchy, this implies that there is a natural upper bound on k_c, m_c and n_c , imposed by the sizes of the caches. In addition, because the reused data should ideally be kept in cache between iterations, the cache replacement policy and the cache organization impose further restrictions on the optimal values for k_c, m_c and n_c .

In the remainder of this section, we describe our analytically model for identifying values for these characteristic parameters. For clarity, the discussion will focus on identifying the optimal value for k_c . The values for m_c and n_c can be derived in a similar manner.

4.3.2. Keeping B_r in the L1 cache. Recall that the micro-panel B_r is reused in every iteration of Loop 5. As such, it is desirable for B_r to remain in the L1 cache while Loop 5 is being

executed. Since the L1 cache implements an LRU replacement policy, conventional wisdom suggests choosing the sizes of the different parts of the matrices such that data required for two consecutive iterations of the loop are kept in cache. This implies that the cache should be filled with B_r , two micro-panels of A_c , and two micro-tiles of C_c . The problem with this approach is that keeping two micro-panels of A_c and two micro-tiles of C_c consequently reduces the size of the micro-panel B_r that fits into the L1 cache, which means that the value k_c is smaller, hence decreasing the amount of data that could possibly be reused. In addition, k_c is the number of iterations for Loop 6, and reducing this value also means less opportunities to amortize data movement with enough computation in the micro-kernel.

Instead, we make the following observations:

- (1) At each iteration of Loop 5, a micro-panel A_r from A_c and the reused micro-panel B_r are accessed exactly once. This implies that the micro-panel that is loaded first will contain the least-recently-used elements. Therefore, because of the LRU cache replacement policy, as long as B_r is loaded after A_r , the entries from A_r will be evicted from the cache before those from B_r .
- (2) Since each micro-panel A_r is used exactly once per iteration, it is advantageous to overwrite the entries of the micro-panel A_r^{prev} which was used in the previous iteration with the corresponding entries of the new A_r^{next} that will be used in the present iteration. Doing so potentially allows a larger B_r to fit into the cache, hence increasing the amount of data being reused.

4.3.3. Evicting A_r^{prev} from cache. We assume that within the micro-kernel (Loop 6), the elements from B_r are loaded after those of A_r . From the second observation above, a strategy to keep a large micro-panel B_r in cache is to evict the old micro-panel A_r^{prev} from the cache, loaded in the previous iteration of Loop 5, by replacing its entries with those of the new micro-panel A_r^{next} to be used in the current iteration of the loop. To ensure this, we need to enforce that the same entries of all micro-panels of A_c are mapped to the same cache sets. Since the L1 cache comprises N_{L1} sets, then the memory addresses that are $N_{\text{L1}}C_{\text{L1}}$ bytes apart will be mapped to the same set in the L1 cache. This implies that the corresponding elements of consecutive micro-panels of A_c must lie in memory an integer multiple (say C_{A_r}) of $N_{\text{L1}}C_{\text{L1}}$ bytes apart.

Recall that consecutive micro-panels of A_c are packed contiguously in memory, and each micro-panel of A_c contains exactly $m_r \times k_c$ elements. This means that the distance between the same entries of two consecutive micro-panels of A_c must be $m_r k_c S_{\text{DATA}}$ bytes. Therefore, A_r^{prev} will be replaced by A_r^{next} only if

$$m_r k_c S_{\text{DATA}} = C_{A_r} N_{\text{L1}} C_{\text{L1}}.$$

Rearranging the above expression yields the following expression for k_c ,

$$k_c = \frac{C_{A_r} N_{\text{L1}} C_{\text{L1}}}{m_r S_{\text{DATA}}}, \quad (4)$$

which ensures that a newly read micro-panel of A_c is mapped to the same set as the existing micro-panel of A_c . Thus, solving for k_c in (4) is equivalent to finding C_{A_r} , since the remaining variables in the equation are hardware parameters.

4.3.4. Finding C_{A_r} . The astute reader will recognize that C_{A_r} is the number of cache lines taken up by a micro-panel A_r in each set of the L1 cache. Similarly, we can define C_{B_r} as the number of cache lines in each set dedicated to the micro-panel B_r . In order for B_r to remain in a W_{L1} -associative cache, it is necessary that

$$C_{A_r} + C_{B_r} \leq W_{\text{L1}}.$$

In practice, the number of cache lines filled with elements of A_r and B_r has to be strictly less than the degree of associativity of the cache. This is because at least one cache line

must be used when the entries of a micro-tile C_r are loaded into the registers. Since C_r is not packed, it is possible that its entries are loaded into the same set as the entries of A_r and B_r . Hence, we update the previous expression to account for the loading of C_r as follows

$$C_{A_r} + C_{B_r} \leq W_{L1} - 1. \quad (5)$$

As the micro-panel of B_c is packed into $n_r k_c S_{\text{DATA}}$ contiguous memory, we can compute C_{B_r} as follows:

$$C_{B_r} = \left\lceil \frac{n_r k_c S_{\text{DATA}}}{N_{L1} C_{L1}} \right\rceil = \left\lceil \frac{n_r}{m_r} C_{A_r} \right\rceil.$$

Therefore, replacing this expression in the inequality in (5), we get

$$C_{A_r} \leq \left\lfloor \frac{W_{L1} - 1}{1 + \frac{n_r}{m_r}} \right\rfloor,$$

which suggests choosing C_{A_r} as large as possible and then allows us to solve for k_c .

4.3.5. Two-way set associativity. Our value for k_c was chosen to enforce A_r^{next} to be loaded into those cache lines previously occupied by A_r^{prev} . To achieve this effect, we made two assumptions:

- Each set in the cache has to be filled equally with the micro-panel of A_r ; i.e. the number of cache lines in each set containing elements from A_r is the same across all sets. This assumption ensures that the corresponding elements of different micro-panels will be assigned to the same cache set.
- One cache line in each set of the cache is reserved for the elements of C_r so that loading them will not evict the micro-panel B_r that already resides in cache.

The problem with a two-way set associative cache is that satisfying these two assumptions would imply that there remain no available cache lines to hold the elements of B_r , precisely the block that we want to keep in cache.

However, if the size of A_r is $N_{L1} C_{L1} / k$, where N_{L1} is an integer multiple of k , then the micro-panel of A_r that is loaded in the $(k + 1)$ -th iteration will be mapped to the same cache sets as the first micro-panel of A_r . When $k = 2$, this is identical to keeping two iterations worth of data in the cache, which ensures that the micro-panel of B_r is kept in cache. Any larger value of k decreases the size of the micro-panel of A_r , which implies that k_c is reduced. Therefore,

$$m_r k_c S_{\text{DATA}} = \frac{N_{L1} C_{L1}}{2},$$

which implies that k_c is given by the formula:

$$k_c = \frac{N_{L1} C_{L1}}{2 m_r S_{\text{DATA}}} \quad (6)$$

when the cache is 2-way set associative.

5. VALIDATION

In this section, we compare the optimal values derived via our analytical model with those employed by implementations that were either manually tuned by experts or empirically tuned. Unlike previous work [Yotov et al. 2005; Kelefouras et al. 2014] that compared performance against an implementation based on ATLAS, we chose to compare our parameter values against the parameter values from manually-optimized implementations using the BLIS framework. As the algorithms used are identical, any deviation is the result of a difference in the parameter values.

Architecture	N_{VEC}	L_{VFMA}	N_{VFMA}	Analytically derived		Expert BLIS		Expert OpenBLAS	
				m_r	n_r	m_r	n_r	m_r	n_r
Intel Dunnington	2	5+3	1	4	4	4	4	4	4
Intel SandyBridge	4	5+3	1	8	4	8	4	8	4
AMD Kaveri	2	6	2	4	6	4	6	8	2
TI C6678	2	3+4	1	4	4	4	4	-	-

Table II: Comparison between analytical values for m_r , n_r and empirical values chosen by experts.

For this experiment, we chose the following mixture of traditional (x86) architectures and low-power processors: Intel Dunnington (X7660), Intel SandyBridge (E3-1220), AMD Kaveri (A10-7850K), and Texas Instruments C6678. This sample includes both dated (end-of-life) and recent architectures in order to evaluate the accuracy and validity of the model, while taking into account different trends in processor architecture design. We consider double precision real arithmetic to obtain our parameter values.

5.1. Evaluating the model for m_r and n_r

Table II lists the machine parameters necessary for computing m_r and n_r , the optimal values analytically derived for m_r and n_r using the model, and the values of m_r and n_r chosen by the expert when implementing a micro-kernel using the BLIS framework. In addition we include the expert’s values for the micro-kernel underlying OpenBLAS. The instruction latency for the different architectures were obtained from the vendors’ instruction set/optimization manuals. In these cases where the architecture did not include an actual VFMA instruction, we computed L_{VFMA} by adding the latencies for a floating-point multiply and a floating-point addition.

Note that for all the architectures our analytical values match those chosen by the BLIS experts who were not aware of this parallel research. Furthermore, the values for m_r and n_r used in OpenBLAS also match our analytical values in two out of three cases³. We note with interest that OpenBLAS utilized an 8×2 micro-kernel for the AMD Kaveri. While it differs from our analytical 4×6 micro-kernel, we note that AuGEM [Wang et al. 2013], an empirical search tool developed by the authors of OpenBLAS in order to automatically generate the micro-kernel, generated a micro-kernel that operates on a 6×4 micro-tile. This suggests that the original 8×2 micro-kernel currently used by OpenBLAS may not be optimal for the AMD architecture.

The results provide evidence that our analytical model for m_r and n_r is reasonably robust across a variety of architectures.

5.2. Evaluating the model for k_c and m_c

Table III presents the machine parameters to derive k_c and m_c using the model, the optimal analytical values, and the empirical values adopted in BLIS after a manual optimization process (inside parenthesis). We do not report results for n_c because most of these processor architectures do not include an L3 cache, which means that n_c is, for all practical purposes, redundant. For the SandyBridge processor, there was minimal variation in performance when n_c was modified.

Again, our analytical model yields similar values if not identical for both k_c and m_c . The analytical model offered the same values the expert picked when optimizing for k_c .

³There exist no OpenBLAS implementations for the remaining architectures.

Architecture	S_{L1} (Kbytes)	W_{L1}	N_{L1}	k_c	S_{L2} (Kbytes)	W_{L2}	N_{L2}	m_c
Intel Dunnington	32	8	64	256 (256)	3,072	12	4,096	384 (384)
Intel SandyBridge	32	8	64	256 (256)	256	8	512	96 (96)
AMD Kaveri	16	4	64	128 (128)	2,048	16	2,048	1,792 (1,368)
TI C6678	32	4	256	256 (256)	512	4	2,048	128 (128)

Table III: Comparison between analytical values for k_c , m_c and empirical (manually-optimized) values from existing implementations of BLIS (inside parenthesis).

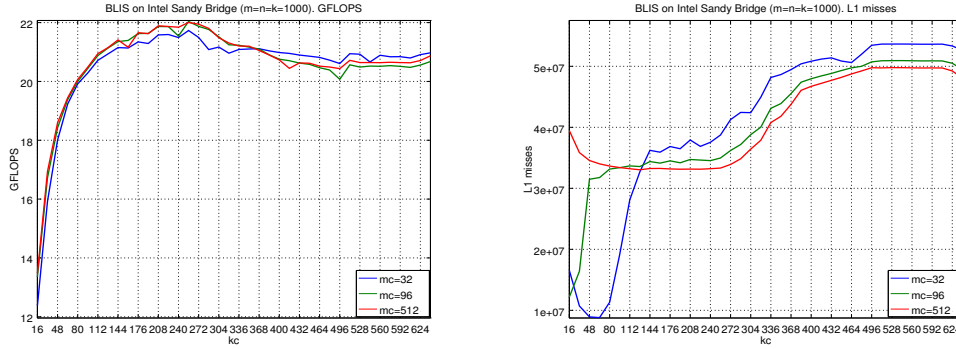


Fig. 5: GFLOPS and L1 cache misses (left and right plots, respectively) on the Intel SandyBridge, for a matrix multiplication with all three operands square of dimension $m = n = k = 1,000$, different values of k_c ; $m_c = 32, 96$ (optimal) and 512; and the remaining three parameters set to the optimal analytical/empirical values for this architecture ($m_r = 8, n_r = 4, n_c = 4,096$).

We expected to encounter more variation between the expert-chosen values and the manually-tuned ones for the m_c parameter. This is because most L2 caches are unified (i.e., they contain both instructions and data), which makes predicting the cache replacement behavior of the L2 cache more difficult, as it depends on both the amount of data and instructions in that level of the cache. Nonetheless, we note that with the exception of the AMD Kaveri, the values computed by our analytical model for m_c were similar if not identical to those chosen by the expert.

5.3. Experimental validation

In this subsection we aim to assess the accuracy of the analytical model in more detail. Our goal is twofold: first, to analyze the behavior of BLIS from the performance and cache miss ratio perspectives on an architecture for which the analytical model exactly predicts the empirical optimal values for m_c and k_c (Intel SandyBridge); and second, to quantify and explain the deviations in performance for an architecture for which the analytical values do not match the observed ones (AMD Kaveri).

5.3.1. Empirical validation of the analytical model. The Intel SandyBridge case. We experimentally measure the behaviour of the L1 cache misses/GFLOPS rate against the parameter k_c .

Figure 5 shows the result of this evaluation on the Intel SandyBridge, for three different values of m_c : the analytical/empirical optimal ($m_c = 96$) and two values, one below and one above the optimal ($m_c = 32$ and $m_c = 512$, respectively). Let us relate the size of the blocks in the L1 cache to the performance lines and the *compulsory* (also known as cold), *conflict* and *capacity* misses [Hennessy and Patterson 2003]. Consider the m_c -optimal case (green line) first. We can clearly identify three ranges or intervals in the results for the L1 cache misses (right-hand side plot):

- $k_c < 48$ (i.e., $k_c = 16$ or 32). For these two small values, a complete block A_c ($m_c \times k_c$) fits into the L1 cache together with a single micro-panel B_r ($k_c \times n_r$). In this scenario, the volume of conflict/capacity misses is negligible, and most of L1 cache misses are compulsory, occurring when A_c is packed by the corresponding routine and each B_r is loaded from the micro-kernel. (The number of L1 cache misses during the packing of B_c is negligible as well.) We emphasize that no L1 cache misses are expected in the access to any of the micro-panels A_r from the micro-kernel, as they are all part of A_c and this block already resides in the L1 cache after it is packed. The repeated use of the micro-panels A_r from within Loop 4 explains why this yields such a low rate of L1 cache misses. On the other hand, the low value of k_c does not allow to hide the cost of transferring the data from the cache, and produces the low GFLOPS rate for these two particular values.
- $48 \leq k_c \leq 256$. This range of values for k_c ensures that A_r and B_r respectively occupy 16 and 8 KBytes at most (i.e., for the largest value of the range, $k_c = 256$). Therefore, no capacity misses are to be expected in the L1 cache. Furthermore, because of the 8-way set-associative configuration of the L1 cache in this architecture, the careful aligned packing of BLIS will ensure that A_r occupies up to $C_{A_r} = 4$ lines per set of the L1 cache, B_r up to $C_{B_r} = 2$ lines/set, and the remaining 2 lines/set are left “empty” for entries of C_c . Therefore, no significant amount of conflict misses is to be expected. The L1 cache misses are mostly compulsory, and originate when the micro-panels A_r and B_r are loaded from within the micro-kernel. (This time, the volume of L1 cache misses during the packing of both A_c and B_c is negligible.) Indeed, most of the misses come from the load of A_r , as the block A_c does not fit into the L1 cache, and has to be repeatedly accessed for each iteration of Loop 4. We also note that the number of misses is almost constant for the full range of values. This is because A_c is streamed by the micro-kernel, from the L2 cache through the L1 cache, approximately n_c/n_r times within Loop 4, independently of the actual value of k_c . From the point of view of performance, the best option occurs at $k_c = 256$, because this is in the range of value that match the architecture configuration, and is the largest one so that it better amortizes the overheads of packing and overlaps data transfers from the L1/L2 cache with computation.
- $k_c > 256$. As k_c exceeds the analytical optimal, conflict misses first and capacity misses later dominate resulting in lower performance.

Three intervals can be distinguished as well for the suboptimal case $m_c = 32$ (blue line): $k_c < 144$, $144 \leq k_c \leq 256$ and $k_c > 256$. The first interval has the same cause as in the optimal scenario above (i.e., $k_c < 48$ and $m_c = 96$), but given that the value of m_c is now $3\times$ smaller than the optimal, the value of k_c that yields a block A_c that fits into the L1 cache is multiplied by 3, yielding the upper threshold $k_c < 144$ for the first interval. Now, given that the second interval is defined by the dimensions of the micro-panels A_r and B_r that fit into the L1 cache, it should not be surprising that the upper threshold is defined by $k_c = 256$ when $m_c = 32$ as well, since m_c plays no role in the dimension of these two micro-panels. Thirdly, the higher amount of L1 cache misses of this suboptimal case when compared to the optimal one, with $k_c \geq 144$ is explained because the number of times that A_c is repacked and repeatedly accessed is proportional to m/m_c .

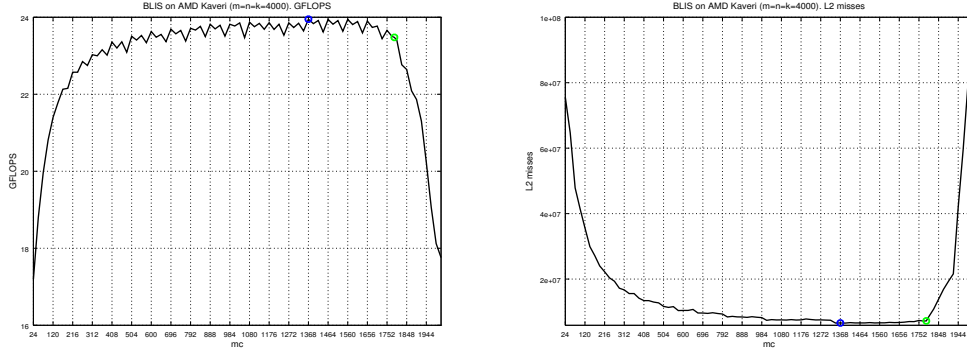


Fig. 6: GFLOPS and L2 cache misses (left and right plots, respectively) on the AMD Kaveri , for a matrix multiplication with all three operands square of dimension $m = n = k = 4,000$, the optimal $k_c = 128$, and varying m_c in steps of 24. The blue and green markers respectively identify the empirical optimal value $m_c^E = 1,368$ and the analytical suboptimal $\tilde{m}_c^A = 1,788$.

The same comments of the previous case apply to the suboptimal case $m_c = 512$ (red line), with the only exception of the high number of L1 cache misses when k_c is very small. The reason is that, for this large value of m_c , the packing of A_c produces this effect.

5.3.2. Quantification of the analytical model deviations. The AMD Kaveri case. Let us analyze the discrepancy between the analytical and empirical optimal values for m_c on the AMD Kaveri. In particular, the model fixes the analytical optimal at $m_c^A = 1,792$, while the experimentation found the empirical optimal at $m_c^E = 1,368$. Unfortunately, BLIS enforces that m_c is an integer multiple of both m_r and n_r (on this architecture, 4 and 6 respectively) so that the analytical optimal cannot be experimentally tested. To avoid exceeding the capacity of the cache, in the following analysis we will thus consider the analytical suboptimal value at $\tilde{m}_c^A = 1,788$, which corresponds to the closest integer multiple of both 4 and 6 smaller than m_c^A .

Figure 6 relates the performance (GFLOPS) and L2 cache misses on the AMD Kaveri, using the optimal value $k_c = 128$, for a wide range of values for m_c that include both m_c^E and \tilde{m}_c^A . The first point to note in the right-hand side plot is that the analytical suboptimal correctly determines the dimension from which the number of L2 cache misses initiates an exponential growth due to capacity constraints (specifically, this occurs from $m_c \geq 1,800 > m_c^A > \tilde{m}_c^A$). Compared with the empirical optimal, the case \tilde{m}_c^A increases the L2 cache misses by 8.85% which, in principle, seems a large value but cannot be appreciated in the plot due to the scale. However, relative to the maximum number of cache misses in the plot (i.e, normalized against the misses for $m_c = 2,016$), the misses for \tilde{m}_c^A only incur an increase of 0.78% with respect to the figure for m_c^E . With these parameters, the slightly larger number of L2 cache misses results in a small decrease of performance: from 23.95 GFLOPS for m_c^E to 23.48 GFLOPS for \tilde{m}_c^A , i.e. -1.97% .

Let us consider now the empirical suboptimal $\tilde{m}_c^E = 1,656$, which corresponds to the largest possible value of m_c that still is within the “noise” of the best performance. From a practical point of view, the difference in performance between m_c^E and \tilde{m}_c^E is negligible (23.946 GFLOPS for the latter, which represents a decrease of -0.017% with respect to m_c^E and can be considered to be in the noise level). Taking into account the dimensions of A_c , for \tilde{m}_c^E this block occupies 1,656 KBytes, i.e. 85.55% of the L2 cache or, in this 16-way associative cache, close to 13 lines of each 16-line set (the exact value is 12.93). We can observe that this is not far from the analytical optimal m_c^A , which dictates that

A_c should occupy 87.50% of the L2 cache, or 14 lines of each 16-line set. The conclusion from this analysis is that the discrepancies between the analytical values \tilde{m}_c^A , m_c^A and the empirical values \tilde{m}_c^E , m_c^E seem to be due to problems with the alignment of A_c in the 16-way associative cache, probably related to the packing of A_c producing unexpected conflict misses, and not to a deviation in the analytical model. A second source of alignment problems may be that, for this architecture, $n_r (= 6)$ is not a power of 2. Therefore, it may occur that each micro-panel of B_c does not overwrite the previous micro-panel of this block, in turn causing conflict misses since the streaming of this matrix does not overwrite the micro-panel of B in the L2 cache.

6. CONCLUSION AND FUTURE DIRECTIONS

We have developed an analytical model that yields optimal parameter values for the GEMM algorithm that underlies BLIS as well as other manually-tuned high performance dense linear algebra libraries. The key to our approach is the recognition that many of the parameters that characterize BLIS are bounded by hardware features. By understanding how the algorithm interacts with the processor architecture, we can choose parameter values that leverage, more aggressively, hardware features such as the multi-layered cache memory of the system.

We compare the values obtained via our analytical approach with the values from manually-optimized implementations and report that they are similar if not identical. Unlike similar work that compares with ATLAS, this is the first paper, to the best of our knowledge, that demonstrates that an analytical approach to identifying parameter values can be competitive with best-in-class, expert-optimized implementations. This demonstrates that a high performance implementation for the GEMM algorithm can be obtained without relying on a costly empirical search.

We believe that more can be done in terms of enhancing the analytical model. One possible direction we are exploring is to extend the analytical model to more complicated linear algebra operations such as those in LAPACK [Anderson et al. 1999]. In general these operations, e.g. matrix factorizations, are implemented as blocked algorithms, and the GEMM operation is often a sub-operation in the loop body. It would be interesting to determine how to identify optimal block sizes, given that the parameters that favor performance for the GEMM operation are known.

A second direction to enhance the analytical model is to incorporate data movement considerations. Currently, the parameters for the micro-kernel (m_r and n_r) are determined by only considering the latency of the floating-point arithmetic instructions. However, this makes the assumption that bandwidth is large enough. On low-power systems, this assumption may no longer hold. In such scenario, m_r and n_r may have to be larger, so that the time for computation is sufficiently long to hide the time it takes to load the next elements of A_r and B_r into the processor registers.

Finally, our current analytical model assumes double precision arithmetic. With complex arithmetic, a micro-tile of the same size incurs four times as many flops, and involves twice as much data. These changes necessarily mean that the formulas in this paper have to be updated. Nonetheless, we believe that a similar analysis of the algorithm and the hardware features yields an analytical model for a micro-kernel that operates with complex arithmetic.

A question becomes whether our analysis can be used to better design future architectures. This question is at least partially answered in [Pedram et al. 2012b; Pedram et al. 2012a], which examines how to design specialized hardware (both compute core and entire processor) for linear algebra computation. The models used for such purpose have much in common with our model for determining the parameter values for the micro-kernel. This suggests that our model for determining the parameter values for loops around the micro-kernel can potentially be leveraged to either determine the ideal cache size and/or cache replacement policy.

Acknowledgments

This research was sponsored in part by NSF grants ACI-1148125/1340293 and CCF-0917167.

Enrique S. Quintana-Ortí was supported by project TIN2011-23283 of the *Ministerio de Ciencia e Innovación* and FEDER. Francisco D. Igual was supported by project TIN2012-32180 of the *Ministerio de Ciencia e Innovación*. This work was partially performed during their visit to The University of Texas at Austin (UT), funded by the JTO visitor applications programme from the Institute for Computational Engineering and Sciences (ICES) at UT.

REFERENCES

- AMD. 2015. AMD Core Math Library. <http://developer.amd.com/tools-and-sdks/cpu-development/amd-core-math-library-acml/>. (2015).
- E. Anderson, Z. Bai, C. Bischof, L. S. Blackford, J. Demmel, Jack J. Dongarra, J. Du Croz, S. Hammarling, A. Greenbaum, A. McKenney, and D. Sorensen. 1999. *LAPACK Users' guide* (3rd ed.). SIAM.
- David F. Bacon, Susan L. Graham, and Oliver J. Sharp. 1994. Compiler Transformations for High-performance Computing. *ACM Comput. Surv.* 26, 4 (Dec. 1994), 345–420. DOI: <http://dx.doi.org/10.1145/197405.197406>
- G. Belter, J. G. Siek, I. Karlin, and E. R. Jessup. 2010. Automatic Generation of Tiled and Parallel Linear Algebra Routines: A partitioning framework for the BTO Compiler. In *Proceedings of the Fifth International Workshop on Automatic Performance Tuning (iWAPT10)*. 1–15.
- Jeff Bilmes, Krste Asanovic, Chee-Whye Chin, and Jim Demmel. 1997a. Optimizing Matrix Multiply Using PHiPAC: A Portable, High-performance, ANSI C Coding Methodology. In *Proceedings of the 11th International Conference on Supercomputing (ICS '97)*. ACM, New York, NY, USA, 340–347. DOI: <http://dx.doi.org/10.1145/263580.263662>
- Jeff Bilmes, Krste Asanović, Chee whye Chin, and Jim Demmel. 1997b. Optimizing Matrix Multiply using PHiPAC: a Portable, High-Performance, ANSI C Coding Methodology. In *Proceedings of International Conference on Supercomputing*. Vienna, Austria.
- Jim Demmel, Jack Dongarra, Victor Eijkhout, Erika Fuentes, Antoine Petit, Rich Vuduc, R. Clint Whaley, and Katherine Yelick. 2005. Self adapting linear algebra algorithms and software. In *Proceedings of the IEEE*. 2005.
- Jack J. Dongarra, Jeremy Du Croz, Sven Hammarling, and Iain Duff. 1990. A Set of Level 3 Basic Linear Algebra Subprograms. *ACM Trans. Math. Soft.* 16, 1 (March 1990), 1–17.
- Jack J. Dongarra, Jeremy Du Croz, Sven Hammarling, and Richard J. Hanson. 1988. An Extended Set of FORTRAN Basic Linear Algebra Subprograms. *ACM Trans. Math. Soft.* 14, 1 (March 1988), 1–17.
- Kazushige Goto and Robert van de Geijn. 2008a. High Performance Implementation of the Level-3 BLAS. *ACM Trans. Math. Software* 35, 1 (July 2008), 4:1–4:14. <http://doi.acm.org/10.1145/1377603.1377607>
- Kazushige Goto and Robert A. van de Geijn. 2008b. Anatomy of a High-Performance Matrix Multiplication. *ACM Trans. Math. Software* 34, 3 (May 2008), 12:1–12:25. <http://doi.acm.org/10.1145/1356052.1356053>
- John A. Gunnels, Greg M. Henry, and Robert A. van de Geijn. 2001. A Family of High-Performance Matrix Multiplication Algorithms. In *Computational Science - ICCS 2001, Part I (Lecture Notes in Computer Science 2073)*, Vassil N. Alexandrov, Jack J. Dongarra, Benjoe A. Juliano, René S. Renner, and C.J. Kenneth Tan (Eds.). Springer-Verlag, 51–60.
- John L. Hennessy and David A. Patterson. 2003. *Computer Architecture: A Quantitative Approach*. Morgan Kaufmann Pub., San Francisco.
- IBM. 2015. Engineering and Scientific Subroutine Library. <http://www-03.ibm.com/systems/power/software/essl/>. (2015).
- Intel. 2015. Math Kernel Library. <https://software.intel.com/en-us/intel-mkl>. (2015).
- Vasilios Kelefouras, Angeliki Kritikakou, and Costas Goutis. 2014. A MatrixMatrix Multiplication methodology for single/multi-core architectures using SIMD. *The Journal of Supercomputing* (2014), 1–23. DOI: <http://dx.doi.org/10.1007/s11227-014-1098-9>
- T. Kisuki, P.M.W. Knijnenburg, M.F.P. O'Boyle, and H. A. G. Wijshoff. 2000. Iterative Compilation in Program Optimization. (2000).
- Peter M. W. Knijnenburg, Toru Kisuki, and Michael F. P. O'Boyle. 2002. Iterative Compilation.. In *Embedded Processor Design Challenges (Lecture Notes in Computer Science)*, Ed F. Depretere, Jrgen Teich, and Stamatis Vassiliadis (Eds.), Vol. 2268. Springer, 171–187. <http://dblp.uni-trier.de/db/conf/samos/samos2002.html#KnijnenburgKO02>

- C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. 1979. Basic Linear Algebra Subprograms for Fortran Usage. *ACM Trans. Math. Soft.* 5, 3 (Sept. 1979), 308–323.
- OpenBLAS 2015. <http://www.openblas.net>. (2015).
- Ardavan Pedram, Andreas Gerstlauer, and Robert A. van de Geijn. 2012a. On the Efficiency of Register File versus Broadcast Interconnect for Collective Communications in Data-Parallel Hardware Accelerators. *Computer Architecture and High Performance Computing (SBAC-PAD), 2012 24th International Symposium on* (October 2012).
- Ardavan Pedram, Robert A. van de Geijn, and Andreas Gerstlauer. 2012b. Codesign Tradeoffs for High-Performance, Low-Power Linear Algebra Architectures. *IEEE Trans. Comput.* 61 (December 2012), 1724–1736. DOI:<http://dx.doi.org/10.1109/TC.2012.132>
- Tyler M. Smith, Robert van de Geijn, Mikhail Smelyanskiy, Jeff R. Hammond, and Field G. Van Zee. 2014. Anatomy of High-Performance Many-Threaded Matrix Multiplication. In *IPDPS '14: Proceedings of the International Parallel and Distributed Processing Symposium*. To appear.
- Field G. Van Zee, Tyler Smith, Bryan Marker, Tze Meng Low, Robert A. van de Geijn, Francisco D. Igual, Mikhail Smelyanskiy, Xianyi Zhang, Michael Kistler, Vernon Austel, John Gunnels, and Lee Killough. 2014. The BLIS Framework: Experiments in Portability. *ACM Trans. Math. Soft.* (2014). In review.
- Field G. Van Zee and Robert A. van de Geijn. 2014. BLIS: A Framework for Generating BLAS-like Libraries. *ACM Trans. Math. Soft.* (2014). To appear.
- Qian Wang, Xianyi Zhang, Yunquan Zhang, and Qing Yi. 2013. AUGEM: Automatically Generate High Performance Dense Linear Algebra Kernels on x86 CPUs. In *Proceedings of SC13: International Conference for High Performance Computing, Networking, Storage and Analysis (SC '13)*. ACM, New York, NY, USA, Article 25, 12 pages. DOI:<http://dx.doi.org/10.1145/2503210.2503219>
- R. Clint Whaley and Jack J. Dongarra. 1998. Automatically Tuned Linear Algebra Software. In *Proceedings of SC'98*.
- Kamen Yotov, Xiaoming Li, María Jesús Garzarán, David Padua, Keshav Pingali, and Paul Stodghill. 2005. Is Search Really Necessary to Generate High-Performance BLAS? *Proceedings of the IEEE, special issue on "Program Generation, Optimization, and Adaptation"* 93, 2 (2005).