# High-Performance Matrix Multiplication Algorithms for Architectures with Hierarchical Memories *

John A. Gunnels[†]     Greg M. Henry[‡]     Robert A. van de Geijn[§]

June 20, 2001

### Abstract

During the last half-decade, a number of research efforts have centered around the development of software for generating automatically tuned matrix multiplication kernels. These include the PHiPAC project and the ATLAS project. The software products of both projects employ brute force to search a parameter space for blockings that accommodate multiple levels of memory hierarchy. We take a different approach. Using a simple model of hierarchical memories we employ mathematics to determine a locally-optimal strategy for blocking matrices. The theoretical results show that, depending on the shape of the matrices involved, different strategies are locally-optimal. Rather than determining a blocking strategy at library generation time, the theoretical results show that ideally one should pursue a heuristic that allows the blocking strategy to be determined dynamically at run-time as a function of the shapes of the operands. When the resulting family of algorithms is combined with a highly optimized inner-kernel for a small matrix multiplication, the approach yields performance that is superior to that of methods that automatically tune such kernels. Preliminary results, for the Intel Pentium (R) III processor, support the theoretical insights.

## 1   Introduction

Research in the development of linear algebra libraries has recently shifted to the automatic generation and optimization of the matrix multiplication kernels. The idea is that many linear algebra operations can be implemented in terms of matrix multiplication [3, 11, 7] and that thus it is this operation that should be highly optimized on different platforms. Since the coding effort is considerable, especially when multiple layers of cache are involved, the general concensus is that this process should be automated.

In this paper, we develop a theoretical framework that (1) suggests a formula for the block sizes that should be used at each level of the memory hierarchy, and (2) restricts the possible loop orderings to a specific family of algorithms for matrix multiplication. Together, we show how to use these results to build highly optimized matrix multiplication implementations that utilize the caches in a locally-optimal fashion. The results could be equally well used to limit the search space that must be examined by packages that automatically tune such kernels.

The current pursuit of highly optimized matrix kernels achieved by coding in a high-level programming language started with the implementation of the FORTRAN implementation of Basic linear Algebra Subprograms (BLAS) [5] for the IBM Power2 [1]. Subsequently, the PHiPAC project [4] at UC-Berkeley demonstrated that high-performance matrix multiplication kernels can be written in C and that code generators could be used to automatically generate many different blockings, allowing automatic tuning. Next, the Automatically Tuned Linear Algebra Software (ATLAS) project [12] at the University of Tennessee extended the ideas developed as part of the PHiPAC project by reducing the kernel that is called once matrices are massaged to be in the L1 cache into one specific case: $C = A^T B + \beta C$ for small matrices $A$, $B$, and $C$ and reducing the space searched for optimal blockings. Furthermore it marketed the methodology allowing it to gain wide-spread acceptance and igniting the current craze in the linear algebra community towards automatically tuned libraries. Finally, there has been a considerable recent interest in recursive algorithms and recursive data structures. The idea here is that by recursively partitioning the operands blocks that fit in the different levels of the caches will automatically be encountered [9]. By storing matrices recursively, blocks that are encountered during the execution of the recursive algorithms will be in contiguous memory [2, 8, 10].

Other work closely related to this topic is discussed in other papers presented as part of this session of the conference.

## 2 Notation and Terminology

### 2.1 Special cases of matrix multiplication

The general form of a matrix multiply is $C \leftarrow \alpha A B + \beta C$ where $C$ is $m \times n$, $A$ is $m \times k$, and $B$ is $k \times n$. We will use the following terminology when referring to a matrix multiply when two dimensions are large and one is small:

| | Condition | Shape |
|---|---|---|
| Matrix-panel multiply | $n$ is small |  |
| Panel-matrix multiply | $m$ is small |  |
| Panel-panel multiply | $k$ is small |  |

The following observation will become key to understanding concepts encountered in the rest of the paper: Partition $X = \left( \begin{array}{c|c|c} X_1 & \cdots & X_{N_X} \end{array} \right) = \left( \begin{array}{c} \hat{X}_1 \\ \hline \vdots \\ \hline \hat{X}_{M_X} \end{array} \right)$ for $X \in \{A, B, C\}$, where $C_j$ is $m \times n_j$, $\hat{C}_i$ is $m_i \times n$, $A_p$ is $m \times k_p$, $\hat{A}_i$ is $m_i \times k$, $B_j$ is $k \times n_j$, and $\hat{B}_p$ is $k_p \times n$. Then $C \leftarrow AB + C$ can be achieved by

| multiple matrix-panel multiplies: | $C_j \leftarrow AB_j + C_j$ for $j = 1, \ldots, N_C$ | $C_1 C_2 C_3 \;+= \; \boxed{A} \; B_1 B_1 B_1$ |
|---|---|---|
| multiple panel-matrix multiplies: | $\hat{C}_i \leftarrow \hat{A}_i B + \hat{C}_i$ for $i = 1, \ldots, M_C$ | $\begin{matrix} \hat{C}_1 \\ \hat{C}_2 \\ \hat{C}_3 \end{matrix} \;+= \; \begin{matrix} \hat{A}_1 \\ \hat{A}_2 \\ \hat{A}_3 \end{matrix} \; \boxed{B}$ |
| multiple panel-panel multiplies | $C \leftarrow \sum_p^{N_A} A_p \hat{B}_p + C$ | $\boxed{C} \;+= \; A_1 A_2 A_3 \; \begin{matrix} \hat{B}_1 \\ \hat{B}_2 \\ \hat{B}_3 \end{matrix}$ |

## 2.2 A cost model for hierarchical memories

The memory hierarchy of a modern microprocessor is often viewed as a pyramid: At the top of the pyramid, there are the processor registers, with extremely fast access. At the bottom, there are disks and even slower media. As one goes down the pyramid, while the cost of memory decreases, the amount of memory increases along with the time required to access that that memory.

We will model the above-mentioned hierarchy naively as follows: (1) The memory hierarchy consists of $H$ levels, indexed $0, \ldots, H-1$. Level 0 corresponds to the registers. We will often denote the $i$th level by $L_i$. Notice that on a typical current architecture $L_1$ and $L_2$ correspond the level 1 and level 2 data caches and $L_3$ corresponds to RAM. (2) Level $h$ of the memory hierarchy can store $S_h$ floating point numbers. Generally $S_0 \leq S_1 \leq \cdots \leq S_{H-1}$. (3) Loading a floating point number stored in level $h+1$ to level $h$ costs time $\rho_h$. We will assume that $\rho_0 < \rho_1 < \cdots < \rho_{H-1}$. (4) Storing a floating point number from level $h$ to level $h+1$ costs time $\sigma_h$. We will assume that $\sigma_0 < \sigma_1 < \cdots < \sigma_{H-1}$. (5) If $m_h \times n_h$ matrix $C$, $m_h \times k_h$ matrix $A$, and $k_h \times n_h$ matrix $B$ are all stored in level $h$ of the memory hierarchy then forming $C \leftarrow AB + C$ costs time $2m_h n_h k_h \gamma_h$. (Notice that $\gamma_h$ will depend on $m_h$, $n_h$, and $k_h$).

# 3 Building-blocks for matrix multiplication

Consider the matrix multiplication $C \leftarrow AB + C$ where $m_{h+1} \times n_{h+1}$ matrix $C$, $m_{h+1} \times k_{h+1}$ matrix $A$, and $k_{h+1} \times n_{h+1}$ matrix $B$ are all stored in $L_{h+1}$. Let us assume that somehow an efficient matrix multiplication kernel exists for matrices stored in $L_h$. In this section, we develop three distinct approaches for matrix multiplication kernels for matrices stored in $L_{h+1}$.

Partition

$$(1) \quad C = \left( \begin{array}{c|c|c} C_{11} & \cdots & C_{1N} \\ \hline \vdots & & \vdots \\ \hline C_{M1} & \cdots & C_{MN} \end{array} \right), A = \left( \begin{array}{c|c|c} A_{11} & \cdots & A_{1K} \\ \hline \vdots & & \vdots \\ \hline A_{M1} & \cdots & A_{MK} \end{array} \right), \text{ and } B = \left( \begin{array}{c|c|c} B_{11} & \cdots & B_{1N} \\ \hline \vdots & & \vdots \\ \hline B_{K1} & \cdots & B_{KN} \end{array} \right)$$

where $C_{ij}$ is $m_h \times n_h$, $A_{ip}$ is $m_h \times k_h$, and $B_{pj}$ is $k_h \times n_h$. The objective of the game will be to determine optimal $m_h$, $n_h$, and $k_h$.

## 3.1 Multiple panel-panel multiplies in $L_h$

Noting that $C_{ij} \leftarrow \sum_{p=1}^{K} A_{ip} B_{pj} + C_{ij}$, let us consider the algorithm in Fig. 1 for computing the matrix multiplication. In that figure the costs of the various operations are shown to the right. The order of the outer-most loops is irrelevant to the analysis.

**Algorithm 1**

```
for j = 1, ..., N
    for i = 1, ..., M
        Load C_ij from L_{h+1} to L_h.              m_h n_h ρ_h
        for p = 1, ..., K
            Load A_ip from L_{h+1} to L_h.          m_h k_h ρ_h
            Load B_pj from L_{h+1} to L_h.          k_h n_h ρ_h
            Update C_ij ← A_ip B_pj + C_ij          2 m_h n_h k_h γ_h
        endfor
        Store C_ij from L_h to L_{h+1}              m_h n_h σ_h
    endfor
endfor
```

Figure 1: Multiple panel-panel multiply based blocked matrix multiplication.

The cost for updating $C$ is given by

$$m_{h+1}n_{h+1}(\rho_h + \sigma_h) + m_{h+1}n_{h+1}k_{h+1}\frac{\rho_h}{n_h} + m_{h+1}n_{h+1}k_{h+1}\frac{\rho_h}{m_h} + 2m_{h+1}n_{h+1}k_{h+1}\gamma_h$$

Since it also equals $2m_{h+1}n_{h+1}k_{h+1}$, solving for $\gamma_{h+1}$, the effective cost per floating point operation at level $L_{h+1}$, yields

$$\gamma_{h+1}^{PP} = \frac{\rho_h + \sigma_h}{2k_{h+1}} + \frac{\rho_h}{2n_h} + \frac{\rho_h}{2m_h} + \gamma_h$$

The question now is how to find the $m_h$, $n_h$, and $k_h$ that minimize $\gamma_{h+1}$ under the constraint that $C_{ij}$, $A_{ik}$ and $B_{kj}$ all fit in $L_h$, i.e., $m_h n_h + m_h k_h + n_h k_h \leq S_h$. The smaller $k_h$, the more space in $L_h$ can be dedicated to $C_{ij}$ and thus the smaller the fractions $\rho_h/m_h$ and $\rho_h/n_h$ can be made. A good strategy is thus to let essentially all of $L_h$ be dedicated to $C_{ij}$, i.e., $m_h n_h \approx S_h$. The minimum is then attained when essentially $m_h \approx n_h \approx \sqrt{S_h}$.

Notice that it suffices to have $m_{h+1} = m_h$ or $n_{h+1} = n_h$ for the above cost of $\gamma_{h+1}$ to be achieved. Thus, the above already for the special cases

$$(2) \qquad \begin{pmatrix} C_{11} \\ \hline \vdots \\ \hline C_{M1} \end{pmatrix} += \begin{pmatrix} A_{11} | \cdots | A_{1K} \\ \hline \vdots \quad | \quad | \quad \vdots \\ \hline A_{M1} | \cdots | A_{MK} \end{pmatrix} \begin{pmatrix} B_{11} \\ \hline \vdots \\ \hline B_{K1} \end{pmatrix}$$

$$(3) \quad \left( C_{11} \| \cdots \| C_{1N} \right) += \left( A_{11} | \cdots | A_{1K} \right) \begin{pmatrix} B_{11} \| \cdots \| B_{1N} \\ \hline \vdots \quad \| \quad \| \quad \vdots \\ \hline B_{K1} \| \cdots \| B_{KN} \end{pmatrix}$$

Here the distance between single/thin lines is $k_h$ and between double/thick lines $m_h = n_h$, where $k_h$ is much smaller than $m_h$ and $n_h$.

The inner-most loop in Alg. 1 implements multiple panel-panel multiplies since $k_h$ is small relative to $m_h$ and $n_h$. Hence the name of this section.

## 3.2 Multiple matrix-panel multiplies in $L_h$

Moving the loops over $l$ and $i$ to the outside we obtain the algorithm in Fig. 2(left). Performing an

**Algorithm 2**
for $p = 1, \ldots, K$
    for $i = 1, \ldots, M$
        Load $A_{ip}$ from $L_{h+1}$ to $L_h$.
        for $j = 1, \ldots, N$
            Load $C_{ij}$ from $L_{h+1}$ to $L_h$.
            Load $B_{pj}$ from $L_{h+1}$ to $L_h$.
            Update $C_{ij} \leftarrow A_{ip}B_{pj} + C_{ij}$
            Store $C_{ij}$ from $L_h$ to $L_{h+1}$
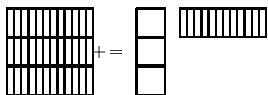        endfor
    endfor
endfor

**Algorithm 3**
for $j = 1, \ldots, N$
    for $p = 1, \ldots, K$
        Load $B_{pj}$ from $L_{h+1}$ to $L_h$.
        for $i = 1, \ldots, M$
            Load $C_{ij}$ from $L_{h+1}$ to $L_h$.
            Load $A_{ip}$ from $L_{h+1}$ to $L_h$.
            Update $C_{ij} \leftarrow A_{ip}B_{pj} + C_{ij}$
            Store $C_{ij}$ from $L_h$ to $L_{h+1}$
        endfor
    endfor
endfor

Figure 2: Multiple matrix-panel (left) and panel-matrix (right) multiply based blocked matrix multiplication.

analysis similar to that given in Section 3.1 the effective cost of a floating point operation is now given by

$$\text{(4)} \qquad \gamma_{h+1}^{MP} = \frac{\rho_h}{2n_{h+1}} + \frac{\rho_h + \sigma_h}{2k_h} + \frac{\rho_h}{2m_h} + \gamma_h$$

Again, the question is how to find the $m_h$, $n_h$, and $k_h$ that minimize $\gamma_{h+1}$ under the constraint that $C_{ij}$, $A_{ik}$ and $B_{kj}$ all fit in $L_h$, i.e., $m_h n_h + m_h k_h + n_h k_h \leq S_h$. Note that the smaller $n_h$, the more space in $L_h$ can be dedicated to $A_{il}$ and thus the smaller the fractions $(\rho_h + \sigma_h)/2k_h$ and $\rho_h/2m_h$ can be made. A good strategy is thus to let essentially all of $L_h$ be dedicated to $A_{il}$, i.e., $m_h k_h \approx S_h$. The minimum is then attained when essentially $m_h \approx k_h \approx \sqrt{S_h}$.

Notice that it suffices to have $m_{h+1} = m_h$ or $k_{h+1} = k_h$ for the above cost of $\gamma_{h+1}$ to be achieved. In other words, the above holds for the special cases

$$\text{(5)} \qquad \begin{pmatrix} C_{11} | \cdots | C_{1N} \\ \hline \vdots \quad | \quad \vdots \\ \hline C_{M1} | \cdots | C_{MN} \end{pmatrix} += \begin{pmatrix} A_{11} \\ \hline \vdots \\ \hline A_{M1} \end{pmatrix} \begin{pmatrix} B_{11} | \cdots | B_{1N} \end{pmatrix}$$

$$\text{(6)} \quad \begin{pmatrix} C_{11} \| \cdots \| C_{1N} \end{pmatrix} += \begin{pmatrix} A_{11} | \cdots | A_{1K} \end{pmatrix} \begin{pmatrix} B_{11} \| \cdots \| B_{1N} \\ \hline \vdots \quad \| \quad \vdots \\ \hline B_{K1} \| \cdots \| B_{KN} \end{pmatrix}$$

The inner-most loop in Alg. 2 implements multiple matrix-panel multiplies since $n_h$ is small relative to $m_h$ and $k_h$. Thus the name of this section.

### 3.3 Multiple panel-matrix multiplies in $L_h$

Finally, moving the loops over $p$ and $j$ to the outside we obtain the algorithm given in Fig. 2(right). This time, the effective cost of a floating point operation is given by

$$\text{(7)} \qquad \gamma_{h+1}^{PM} = \frac{\rho_h}{2m_{h+1}} + \frac{\rho_h + \sigma_h}{2k_h} + \frac{\rho_h}{2n_h} + \gamma_h$$

Again, the question is how to find the $m_h$, $n_h$, and $k_h$ that minimize $\gamma_{h+1}$ under the constraint that $C_{ij}$, $A_{ik}$ and $B_{kj}$ all fit in $L_h$, i.e., $m_h n_h + m_h k_h + n_h k_h \leq S_h$. Note that the smaller $m_h$, the

more space in $L_h$ can be dedicated to $B_{pj}$ and thus the smaller the fractions $(\rho_h + \sigma_h)/2k_h$ and $\rho_h/2n_h$ can be made. A good strategy in this case is to dedicate essentially all of $L_h$ to $B_{pj}$, i.e., $n_h k_h \approx S_h$. The minimum is then attained when essentially $n_h \approx k_h \approx \sqrt{S_h}$.

Notice that it suffices to have $n_{h+1} = n_h$ and/or $k_{h+1} = k_h$ for the above cost of $\gamma_{h+1}$ to be achieved. In other words, the above holds for the special cases

$$(8) \qquad \begin{pmatrix} C_{11} & \cdots & C_{1N} \\ \vdots & & \vdots \\ C_{M1} & \cdots & C_{MN} \end{pmatrix} + = \begin{pmatrix} A_{11} \\ \vdots \\ A_{M1} \end{pmatrix} \begin{pmatrix} B_{11} & \cdots & B_{1N} \end{pmatrix}$$

$$(9) \qquad \begin{pmatrix} C_{11} \\ \vdots \\ C_{M1} \end{pmatrix} + = \begin{pmatrix} A_{11} & \cdots & A_{1K} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MK} \end{pmatrix} \begin{pmatrix} B_{11} \\ \vdots \\ B_{K1} \end{pmatrix}$$

an observation that will become important later.

## 3.4 Summary

The conclusions to draw from Sections 2.1 and 3.1–3.3 are: (1) There are three shapes of matrix multiplication that one expects to encounter at each level of the memory hierarchy: panel-panel, matrix-panel, and panel-matrix multiplication. (2) If one such shape is encountered at $L_{h+1}$, a locally-optimal approach to utilizing $L_h$ will perform multiple instances with one of the other two shapes. (3) Given that multiple instances of a given shape are to be performed, the strategy is to move a submatrix of one of the three operands into $L_h$ (we will call this the resident matrix in $L_h$), filling most of that layer, and to amortize the cost of this data movement by streaming submatrices from the other operands from $L_{h+1}$ to $L_h$.

Interestingly enough, the shapes discussed are exactly those that we encountered when studying a class of matrix multiplication algorithms on distributed memory architectures [6]. This is not surprising, since distributed memory is just another layer in the memory hierarchy.

# 4    A Family of Algorithms

We now show how to turn the observations made in the previous section into a practical implementation.

High-performance implementations of matrix multiplication typically start with an "inner-kernel". This kernel carefully orchestrates the movement of data in and out of the registers and the computation under the assumption that one or more of the operands are in the L1 cache. For our implementation on the Intel Pentium (R) III processor, the inner-kernel performs the operation $C = A^T B + \beta C$ where $64 \times 8$ matrix $A$ is kept in the L1 cache. Matrices $B$ and $C$ have a large number of columns, which we view as multiple-panels, with each panel of width one. Thus, our inner-kernel performs a multiple matrix-panel multiply (MMP) with a transposed resident matrix $A$. The technical reasons why this particular shape was selected go beyond the scope of this paper.

While it may appear that we thus only have one of the three kernels for operation in the L1 cache, notice that for the submatrices with which we compute at that level one can instead compute $C^T = B^T A + C^T$, reversing the role of $A$ and $B$. This simple observation allows us to claim that we also have an inner-kernel that performs a multiple panel-matrix multiply (MPM).

Let us introduce a naming convention for a family of algorithms that perform the discussed algorithms at different levels of the memory hierarchy:
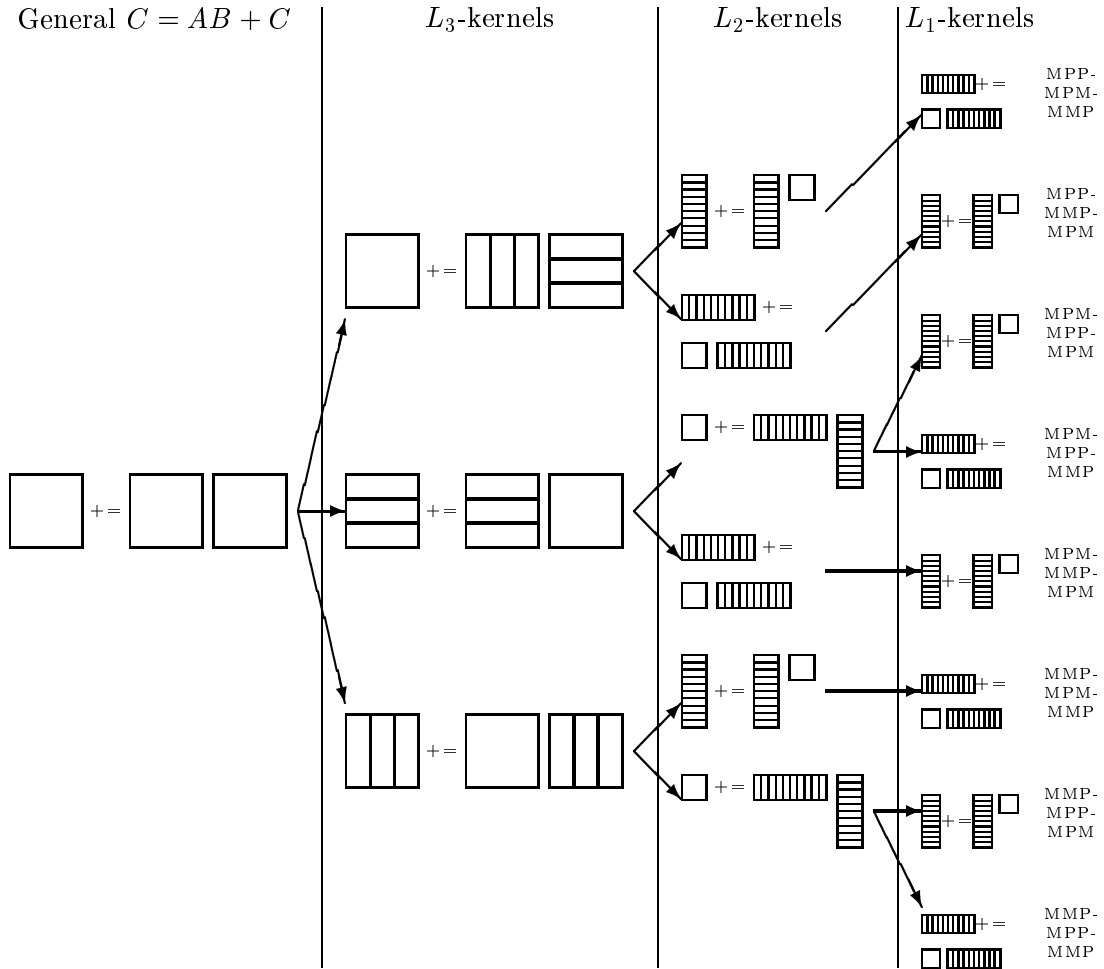
Figure 3: Possible algorithms for matrices in memory level $L_3$ given all $L_2$-kernels.

$$\text{<kernel at } L_3\text{>-<kernel at } L_2\text{>-<kernel at } L_1\text{>}.$$

For example MPP-MPM-MMP will indicate that the $L_3$-kernel uses multiple panel-panel multiplies, calls the $L_2$-kernel that uses multiple matrix-panel multiplies, which in turn calls the $L_1$-kernel that uses multiple panel-matrix multiplies. Given the constraint that only two of the possible three kernel algorithms are implemented at $L_1$, the tree of algorithms in Fig. 3 can be implemented.

## 5    Performance

In this section, we report performance attained by the different algorithms. Performance is reported by the rate of computations attained, in millions of floating point operations per second (MFLOPS/sec). For the usual matrix dimensions $m$, $n$, and $k$, we use the operation count $2mnk$ for the matrix multiplication. We tested performance of the operation $C = C - AB$ ($\alpha = -1$ and $\beta = 1$) since this is the case most frequently encountered when matrix multiplication is used in libraries such as LAPACK.

We report performance on an Intel Pentium (R) III (650 MHz) processor with a 16 Kbyte L1