

Predicting Memory-Access Cost Based on Data-Access Patterns

Surendra Byna

Xian-He Sun

William Gropp

Rajeev Thakur

Computer Science Department
Illinois Institute of Technology
Chicago, IL 60616
{renbyna, sun}@iit.edu

Math. and Comp. Science Division
Argonne National Laboratory
Argonne, IL 60439
{gropp, thakur}@mcs.anl.gov

Abstract

The expanding processor-memory performance gap cause several scientific applications spend a major part of their execution time stalled in accessing memory. Loops are the basic blocks, where most of these data accesses take place. Memory performance improvement relies in optimizing these loop accesses. Loop transformations (e.g. loop unrolling, loop tiling) and array restructuring optimizations improve this performance by increasing the locality of memory accesses. However, to find the best optimization parameters dynamically at runtime, we need a fast and simple analytical model to predict the memory access cost. Most of the existing models are complex and impractical to be integrated in the automatic tuning systems. In this paper, we propose a simple, fast and reasonably accurate model that is capable of predicting the memory access cost based on a wide range of data access patterns. We classify the data access patterns that appear in the nested loops of scientific applications. We also discuss our applying this model in optimizing the performance of derived datatypes in Message Passing Interface (MPI) implementation.

Keywords: Memory access cost, data access patterns, computer architecture, memory hierarchy, performance prediction, memory access cost prediction.

1. Introduction

Immense research effort has been spent on reducing the performance gap between processor and memory. Processor speed continues to increase every year. CPU clock frequency is doubling every 18 months complying with Moore's Law. On the other hand main memory (DRAM) speeds haven't increased enough to catch up with the processor speed. This performance gap has been increasing for the last 20 years [Patt96] and the trend appears to continue in the near future.

Advanced hierarchical memories that include cache memories at various levels are developed to bridge this gap. A cache memory works on the principle of spatial and temporal locality [Smit82]. However, there are many applications that lack locality in accessing the memory. These applications spend a major fraction of execution time waiting for data accesses. In other words, cache memories are exploited better if the cached blocks of data are reused extensively before other cache blocks replace them.

Transforming and reordering the memory accesses improve application performance [Mcki96, Kand99]. As loops are the basic blocks, where most of time is spent in applications, various loop optimization techniques have been developed to enhance the memory hierarchy utilization. Loop transformations (loop unrolling, loop fusion, loop interchange, loop reversal and loop tiling) are some of the most effective loop optimizations.

Some advanced compilers utilize these optimization techniques in various levels to improve application performance. However, compilers alone are not sufficient to achieve the best possible optimization [Byna03]. Optimizations that are implemented manually, achieve better performance. But superior manual optimizations require extensive knowledge of the hardware architecture and data access patterns. The developer needs to be aware of efficient optimization techniques to be applied in the right place. It is also required to automate the optimization process to obtain standard and consistent performance.

Performance prediction of memory access cost is required to automate the optimizations. Currently there are few automatic tuning software tools. One of the most popular tools of optimization is Automatically Tuned Linear Algebra Software (ATLAS) [Whal01]. This tool runs subroutines multiple times to obtain the best optimization parameters by a trial and error method. A prediction model can remove these multiple runs and be extended to optimize more than just linear algebra subroutines. When the prediction model is simple and fast, it can facilitate to perform the optimization dynamically, at runtime, based on the data access pattern and available memory hierarchy.

Many performance prediction models are available to help programmers in estimating the cost of memory. Copious research effort has been spent in this area to develop accurate cache performance models. But most of them [Chat00, SSen00, Jaco96] lack generality. They are complex, and are bounded to a few algorithms or data access patterns. Jacob [Jaco96] extracts address traces from the code, which requires execution of the program, and consumes a lot of time if an optimization has to be applied. We base our prediction model on various access patterns, which are parameterized. This helps in predicting the memory cost with very small complexity and skips the costly process of tracing the references every time a loop parameter is changed. Chatterjee et al. [Chat01] studies the exact analysis of cache misses based on the polyhedral model, which is very complex. The Cache Miss Equations model (CME) [Ghos99] is the least costly performance model to our knowledge. However, this model also requires tracing the references to create the reuse vectors and solve cache miss equations. These models are accurate but expensive, and are better choices for static analysis of cache behavior. Our model fits better in choosing the optimization parameters dynamically at runtime than CMEs.

Our model also focuses on wide range of data access patterns with multiple array variables. Most of the other cache analysis models hold good results for a specific algorithm [Chat00, SSen00] and fall short in acquiring generality.

The rest of this paper is organized as follows. Section 2 classifies various data access patterns that are used in most of the scientific applications. Section 3 discusses about the hierarchical memory system and the parameters of it. In Section 4, we propose the memory access cost analysis and prediction equations. Section 5 provides experimental verification, and section 6 discusses the possible application of our model and current projects. Section 7 concludes with further objectives.

2. Data Access Patterns

Loops and arrays are fundamental structures of most numerical and scientific applications [Peak98]. A major share of the execution time of these applications is spent in loops, accessing

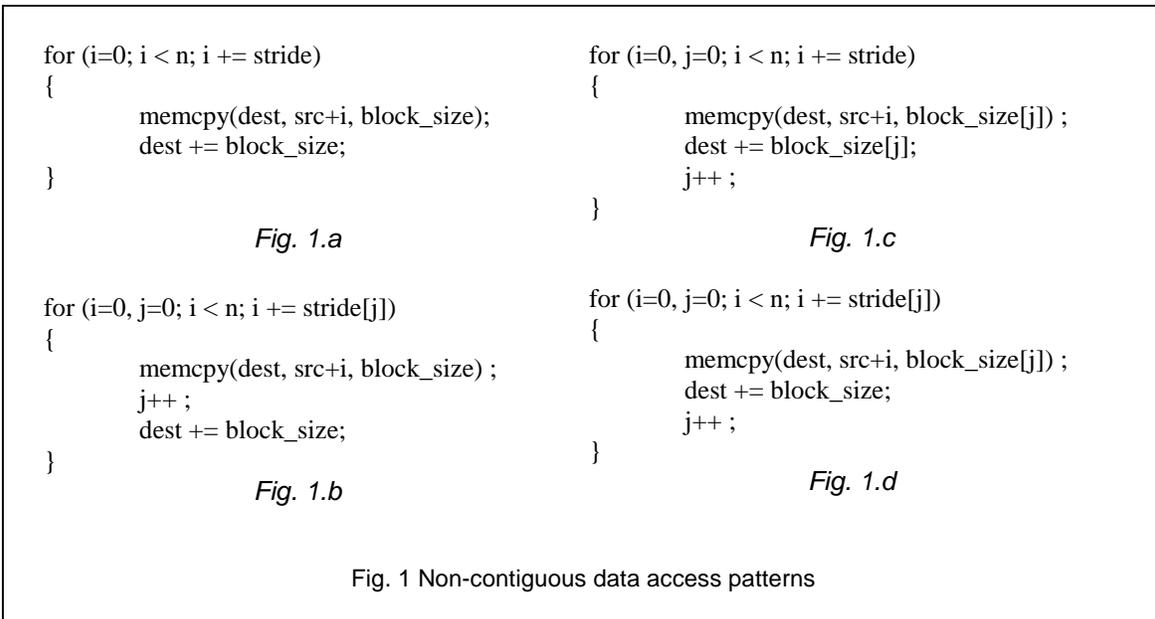
data from arrays. Loop variables are incremented or decremented to find the reference of an array. Analyzing these reference patterns of array accesses is needed to find out the hotspots and to optimize the performance by reorganizing these memory references.

Data access patterns are classified based on the stride between successive accesses. Modal model of memory [Mitic01] categorize data accesses as constant, strided and non-monotonic modes. Yan et al. [Ryan02] classify memory access patterns into three types: migratory, group and unpredictable patterns.

We classify data access patterns in scientific applications as constant, contiguous and non-contiguous. Non-contiguous pattern is further divided into four patterns. This classification is based on the size of data blocks accessed with each reference and their successive strides. Stride is the distance between the previous reference and current reference.

Constant accesses are those where the same data block is accessed repeatedly i.e., stride is equal to zero. In this type of accesses, once the data block is loaded into the cache, further accesses would not cause any other cache misses.

Contiguous access pattern is where the stride between successive accesses is equal to the size of datatype. These are divided further as fixed length block accesses and variable length block



accesses. Fixed length block accesses refer to the same datatype in consecutive references.

Non-contiguous access pattern is where the stride of next reference is greater than the size of currently accessed datatype. These can be further divided as follows:

- a. *Fixed length block, with fixed stride*: Stride is similar through out the access pattern. As shown Fig. 1.a., a block with a size of constant *block_size* is copied into *dest* from *src*. The next block is copied from *src+stride* to *dest+block_size*, i.e. array *src* is being accessed non-contiguously with a fixed stride and array *dest* is being accessed contiguously.
- b. *Fixed length block, with varying stride*: The stride varies between each access. (Fig. 1.b.)
- c. *Variable size block, with fixed stride*: Accessing different or varying size datatypes, where the strides of accesses are similar. (Fig. 1.c.)
- d. *Variable size block, with variable stride*: Accessing different or varying size datatypes, where the strides of accesses are varying. (Fig. 1.d.)

In Section 4, we predict the memory access cost for these data access patterns.

3. Model Parameters

To bridge the gap between processor and memory performance, modern computer architectures include multiple levels of memory hierarchies that consist of cache memory and TLB. The cache memories are inserted either on the die of the processor, or outside between the CPU and main memory. The size as well as the latency increases to access a level of cache as it is placed further from the CPU. It is common to place first level of cache (L1) on the processor. Recently processors are including second level cache (L2) also on the die. A cache of virtual memory mappings is stored on a Translation Look-aside Buffer (TLB). In this section, we discuss the details of memory hierarchy parameters, which are used in developing the prediction model.

A Cache memory is characterized by its size, line size and associativity. Cache size (C) represents its capacity in bytes. Caches are organized in cache lines. The data transfer between

two levels of memory hierarchy is done based on the size of these lines. When a cache miss occurs, a block of data of size equal to cache line size (L) is fetched from the next level of cache or memory. This property conforms to the spatial locality. Associativity of a cache helps in deciding how many places are there to place a cache line. For a *direct mapped* cache (1-set associativity), a fetched cache line can be placed in one of the (C/L) positions. If another cache line is mapped to the same location, the currently residing cache line has to be replaced. Higher associative caches provide more than one place to map a cache block. In a *fully associative* cache, a cache block can be placed anywhere in the cache. But *fully associative* caches are costly. It is

C_k	Cache size at k^{th} level cache of memory hierarchy
L_k	Cache line size at k^{th} level cache of memory hierarchy
A_k	Associativity of k^{th} level cache of memory hierarchy
M_k	Number of cache misses at k^{th} level cache of memory hierarchy
$M_{(k,i)}^c$	Number of cache misses at k^{th} level cache of memory hierarchy in accessing i^{th} array, contiguously.
$M_{(k,i)}^n$	Number of cache misses at k^{th} level cache of memory hierarchy in accessing i^{th} variable, non-contiguously.
T_s	Number of page table entries (PTE) in TLB
P_s	Page size of each PTE
A_T	TLB associativity
M	Number of cache levels in memory hierarchy

Table 1. Memory hierarchy parameters

$R_{(k,i)}^c$	Number of contiguous references of i^{th} array at cache level k of the memory hierarchy.
$R_{(k,i)}^n$	Number of non-contiguous references of i^{th} array at cache level k of the memory hierarchy.
W_i	Fixed size of the data block being accessed in i^{th} array.
S_i	Fixed stride of accessing i^{th} array non-contiguously.
W_i^c	Variable size of the data block being contiguously accessed in i^{th} array.
$W_{(i,j)}^n$	Variable size of j^{th} data block being non-contiguously accessed in i^{th} array.
$S_{(i,j)}$	Variable stride of the j^{th} data block being contiguously accessed in i^{th} array. (stride, in stride signature)
D	Size of working set

Table 2. Data access parameters

common for architectures to have 2-way to-8-way associative caches. In an *8-way associative* cache, a block can be placed in any of the 8 positions in each set of 8 lines.

We treat the TLB as a level of memory hierarchy. Its parameters are page size P (similar to cache line size of a cache) and the capacity. The capacity of a TLB is the amount of memory page mapping it can store and is equal to *number of page entries* multiplied by page size.

Table 1 summarizes the memory hierarchy parameters. Subscript i of a parameter signifies the level of that cache/TLB in the hierarchy of memory. Cache memory at level i has three properties: its size in bytes (C_i), cache line size (L_i) and its associativity (A_i). TLB is represented with the number of page table entries (T_s), page size (P_s) and its associativity (A_T). M refers the total number of cache levels.

Cache misses are classified into three types [Mark87]. *Compulsory misses*: Misses that occur when the CPU accesses a data block for the first time. *Capacity misses*: Misses that occur due to insufficient cache to hold an entire range of data that is being accessed at a time. This range is also called working-set. *Conflict misses*: These misses occur when more than one cache block maps to the same location even though the existing cache block still is being used in the near future. Cache misses at level i are represented with M_i . $M_{(k,i)}^c$ refers to the number of cache misses at level k of memory hierarchy, in accessing i^{th} array (variable) contiguously. If it is being accessed non-contiguously, it is represented by $M_{(k,i)}^n$.

Data access pattern parameters are shown in Table 2. The subscript i represents the i^{th} array being accessed. The parameters $R_{(k,i)}^c$ and $R_{(k,i)}^n$ represent the number of contiguous and non-contiguous references separately. S_i is the fixed stride in accessing the i^{th} array and $S_{(i,j)}$ is the variable stride. D is the working set size and W_i represents the block (word) size of the i^{th} array.

4. Memory Access Cost Prediction

Our goal is to predict the memory access cost of a basic block of loop with any type of data access patterns discussed in section 2, and for multiple data array variables. We assume LRU replacement policy for cache and TLB. We assume that the memory hierarchy is following inclusive property. The total cost of accessing memory includes the access time and the miss penalties of these levels in the hierarchy. If there are k levels of cache memory and one level TLB [Patt96],

$$\begin{aligned}
 \text{Total Memory cost} = & (\text{Number of TLB hits} * \text{Time to access TLB}) + (\text{Number of TLB misses} * \\
 & \text{TLB miss penalty}) + (\text{Number of } L_1 \text{ hits}) * (\text{Time to access } L_1) + (L_1 \\
 & \text{misses} * L_1 \text{ penalty}) + (L_2 \text{ misses} * L_2 \text{ penalty}) + \dots + (L_k \text{ misses} * L_k \\
 & \text{penalty})
 \end{aligned} \tag{4.1}$$

To predict this cost, we have to find the number of cache hits/misses at each level and TLB hit rate. We predict the cache and TLB misses based on the access pattern.

Assuming that there are M levels of cache, the total miss penalty due to cache misses is the sum of miss penalty at each level.

$$T_m = \sum_{k=1}^M (M_k * T_k) - \alpha \tag{4.2}$$

where M_k is the total number of cache misses and T_k is the miss penalty at level k cache. α is the overlapping the cache misses with prefetching and other OS optimizations.

Consider that there are m array variables accessed contiguously and n array variables accessed non-contiguously, the total number of misses at cache level k is the sum of misses caused in accessing contiguously accessed arrays and those of non-contiguously accessed arrays.

$$M_k = \sum_{i=1}^m M_{(k,i)}^c + \sum_{i=1}^n M_{(k,i)}^n \tag{4.3}$$

$M_{(k,i)}^c$ is the number of cache misses at k^{th} level cache of memory hierarchy in accessing i^{th} variable, contiguously. $M_{(k,i)}^n$ is the number of cache misses at k^{th} level cache of memory hierarchy in accessing i^{th} variable, non-contiguously.

Now we count the number of cache misses based on the data access pattern.

Constant access pattern: In this type of accesses, once a word is loaded into the cache, the following accesses to the same word cause no extra cache misses. If the word size of a variable is W_i and there are the number of cache misses is equal to $\lceil (W_i / L_k) \rceil$. (4.4)

Contiguous access pattern: In this pattern the stride between successive accesses is the same as data size. All the cache misses caused in this pattern are compulsory misses. Each reference fetches a cache line into the cache. Cache line contains more than one word of data. This is to assure spatial locality property of using cache memory. If the cache line size is more than the data type accessed, the next reference utilizes the prefetched data from the cache. Each reference causes $\lceil (W_i / L_k) \rceil$ misses, i.e. if the word size is more than cache line size, then it causes more than one miss, otherwise just one miss occurs for every $\lceil (L_k / W_i) \rceil$ references. If there are n references, the number of cache misses caused at cache level k in accessing variable i is:

$$M_{(k,i)}^c = \lceil n * (W_i / L_k) \rceil \quad (4.5)$$

If i^{th} variable has $R_{(k,i)}^c$ references, the number of cache misses is:

$$M_{(k,i)}^c = \left\lceil R_{(k,i)}^c * \left(\frac{W_i^c}{L_k} \right) \right\rceil \quad (4.6)$$

where W_i^c is size of the data block being contiguously accessed in i^{th} variable.

The number of cache references at level k ($R_{(k,i)}^c$) is the number of cache misses at the lower level cache, i.e. $R_{(k,i)}^c = M_{(k-1,j)}^c$. (4.7)

Non-contiguous access patterns: As described in section 2, there are four main types of access patterns. These patterns are classified based on the variability of stride and data block size. The occurrence of cache misses is categorized into four regions based on the working set size, similar to Saavedra and Smith [Saav95]. First region is the one where all the working set fits in the cache. As long as the working set size is less than the cache size, the total data fits into the cache. All the cache misses are compulsory misses. This number is equal to $M_{(k,i)}^n = \lceil n * (W_i / L_k) \rceil$ at level k of memory hierarchy in accessing i^{th} variable non-contiguously. This number is the same for all types of non-contiguous access patterns.

When size of the data working set exceeds the cache size, three regions of memory operations are defined. The first region is when the stride (S) is between 1 and cache line size ($1 < S \leq L_k$). The second region is $L_k < S \leq D / A_k$, where A_k is the associativity of k^{th} level cache of memory hierarchy. The third region is $D / A_k < S \leq D / 2$. In this last case, although the $D > C_k$, only $D / S < A_k$ amount of data is needed for access. In the last region the number of references mapping to a single set is less than the set associativity. Thus, only compulsory misses are caused in the third access pattern, i.e.

$$M_{(k,i)}^n = \lceil n * (W_i / L_k) \rceil \quad (4.8)$$

where n is the number of data accesses. Thus, we set our focus on the first two regions to count the number of cache misses.

First we find the cache misses for a fixed size of data block accesses of one variable, with a fixed stride.

If the stride (fixed) is less than the cache line size, one cache miss occurs for (L_k / S) references. If there are n references, the number of cache misses is: $n * (S / L_k)$, where n is the number of data accesses.

If the stride (fixed) is more than the cache line size, each reference causes $(\max(\lceil W_i / L_k \rceil, 1))$ cache misses, i.e. each access causes one miss when the word size is less than L_k . If word size is more than L_k , each reference causes $\lceil W_i / L_k \rceil$ misses. If there are n references, the number of cache misses is equal to $n * (\max(\lceil W_i / L_k \rceil, 1))$.

$$M_{(k,i)}^n = R_{(k,i)}^n * (\max(\lceil W_i / L_k \rceil, 1)) \quad (4.9)$$

For variable stride with fixed size block accesses, the cache misses have to be counted for each stride. If the stride is less than L_k , it does not cause a cache miss as the pre-fetched line of data is reused. The number of cache misses is:

$$M_{(k,i)}^n = \left(\sum_{j=1}^{R_{(k,i)}^n} \lfloor \min((S_{(i,j)} / L_k), 0) \rfloor \right) * (\max(\lceil W_{(i,j)}^n / L_k \rceil, 1)) \quad (4.10)$$

$M_{(k,i)}^c$ is the number of cache misses at k^{th} level cache of memory hierarchy in accessing i^{th} variable, non-contiguously, $S_{(i,j)}$ is variable stride of the data block being contiguously accessed in i^{th} variable. $R_{(k,i)}^n$ is the number of non-contiguous references of i^{th} array at cache level k of the memory hierarchy. $W_{(i,j)}^n$ is the size of j^{th} data block being non-contiguously accessed in i^{th} variable. In this pattern when the stride $S_{(i,j)}$ is less than the cache line size, we assume that the cache line has already been fetched into the cache. However when this stride is causing to fetch a new cache line, then this formula misses to count that cache miss. This can be corrected by maintaining the history of cache line that has been fetched recently.

The number of cache references at level k ($R_{(k,i)}^n$) is the number of cache misses at the lower level cache, i.e. $R_{(k,i)}^n = M_{(k-1,j)}^n$. (4.11)

For fixed or variable stride with variable size block accesses, the cache misses have to be counted for each block size. In this case, we assume that the stride is always more than L_k . If the

data block size is less than L_k , it does not cause a cache miss as the pre-fetched line of data is reused. The number of cache misses is:

$$M_{(k,i)}^n = \sum_{j=1}^{R_{(k,i)}^n} (\max(\lceil W_{(i,j)}^n / L_k \rceil, 1)) \quad (4.12)$$

Refer Table 3 for a summary of formulae to calculate the cache misses for all data access patterns. Using (4.3) total number of cache misses in accessing contiguous and non-contiguous data is calculated. Formula 4.2 gives the total memory access cost.

Data access pattern		Number of cache misses
Constant		$M_{(k,i)}^c = \lceil W_i^c / L_k \rceil$
Contiguous		$M_{(k,i)}^c = \lceil R_{(k,i)}^c * (W_i^c / L_k) \rceil$
Non-contiguous ($D < C_k$)		$M_{(k,i)}^n = \lceil R_{(k,i)}^n * (W_i^n / L_k) \rceil$
Non-contiguous ($D > C_k$)	$1 < S \leq L_k$	$M_{(k,i)}^n = \lceil R_{(k,i)}^n * (S / L_k) \rceil$
	$L_k < S \leq D / A_k$	$M_{(k,i)}^n = R_{(k,i)}^n * (\max(\lceil W_i^n / L_k \rceil, 1))$
	Variable stride, fixed data block size	$M_{(k,i)}^n = (\sum_{j=1}^{R_{(k,i)}^n} \lfloor \min((S_{(i,j)} / L_k), 0) \rfloor) * (\max(\lceil W_{(i,j)}^n / L_k \rceil, 1))$
	Variable stride ($L_k < S \leq D / A_k$), variable data block size	$M_{(k,i)}^n = \sum_{j=1}^{R_{(k,i)}^n} (\max(\lceil W_{(i,j)}^n / L_k \rceil, 1))$
	$D / A_k < S < D / 2$	$M_{(k,i)}^n = \lceil R_{(k,i)}^n * (W_i^n / L_k) \rceil$

Table 3. Number of cache misses for all data access patterns

5. Performance Verification

This section presents performance measurements to verify the predicted memory access cost with the measured cost on various architectures. We measure the performance of loops with all the data access patterns mentioned above and compare that performance with the predicted performance.

We took the measurements on a Sun Solaris based cluster called *Sunwulf*, which is located at the Scalable Computer Software Lab of Illinois Institute of Technology. *Sunwulf* is composed of a four-processor E450 server and 63 high-end workstations. We run our experiments on one of the nodes. Each node is a SUN Blade-100 workstation with one UltraSparc-IIe, 500MHz CPU. The L1 cache is 16KB, with a 16-byte cache line size. The L2 cache has a capacity of 8MB and its line size is 64 bytes. It also has a TLB with 4KB page size and 48 entries. We used a microbenchmark to find the average access time and miss penalty of each level of memory hierarchy. This is similar to the microbenchmark proposed by Saavedra and Smith [Saav95].

Another platform we used for experiments is a 32-node Beowulf, located at University of South Carolina. Each node consists of 933MHz, Pentium III processor. It has 16 KB L1 cache and 256KB L2 cache. Both these caches are on the die, and the average penalty for load misses is measured as 7 cycles and 70 cycles for L1 and L2 respectively. We chose these processors, as they apply inclusive property in the memory hierarchy with less aggressive pre-fetching.

We used the loops similar to Fig.1 and measured the time to execute those loops. In all these loops, two array variables are accessed with different access patterns. We chose these loops since many applications contain loop blocks where the data accesses are similar to the access patterns discussed above in section 2. We can apply the same prediction model for any number of arrays. Execution time of these loops contains only the data access cost, without any computation cost. We used pointer-to-pointer copy to avoid the cost of memcopy. In these experiments we ran many iterations of the program to find the minimum cost. We also flushed the cache after measuring the time for an iteration to replace any cache blocks that are reusable. We compiled these programs using *gcc 3.0* and padded the arrays to avoid any cache thrashing. The comparison of predicted cost and measured memory access cost is presented in the following paragraphs. The memory access cost is presented as a ratio of execution time to the number of memory references. This normalization is done to fit all the data into the graph. The performance is better for lower values.

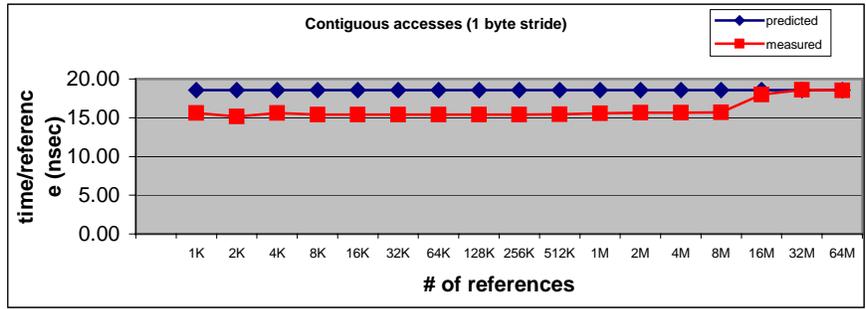


Fig. 2.a.

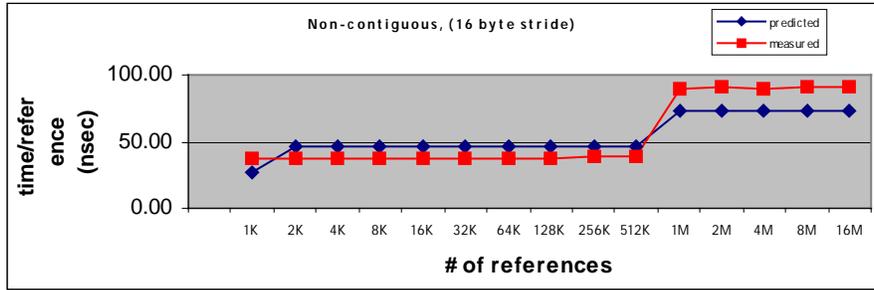


Fig. 2.b.

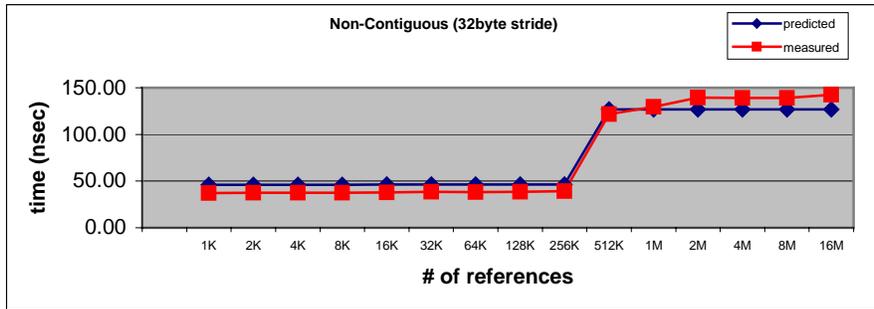


Fig. 2.c.

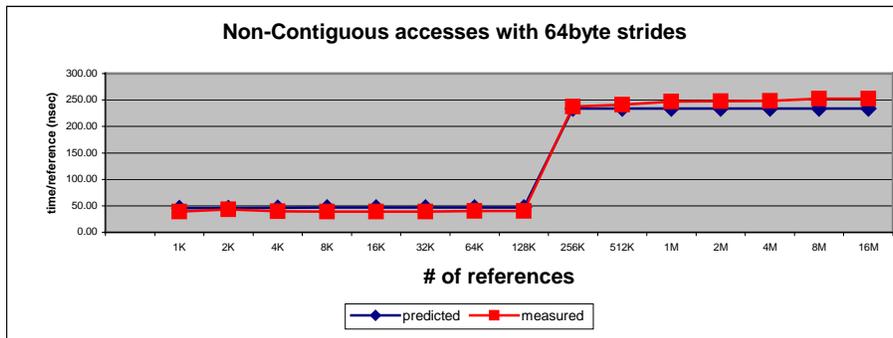


Fig. 2.d.

Fig.2. Comparison of measured and predicted memory access cost. The access patterns are: 4.a. Contiguous data access (word size: 1byte, stride: 1byte). 4.b. Non-contiguous data access with fixed word size and stride (word size: 8 bytes, stride: 16 bytes), 4.c. Non-contiguous data access with fixed word size and stride (word size: 8 bytes, stride: 32 bytes) 4.d. Non-contiguous data access with fixed word size and stride (word size: 8 bytes, stride: 64 bytes)

Fig. 2 and Fig. 3 compare the predicted memory access cost with measured cost in running the loops in various data access patterns explained in section 2 (Fig.1) on Sunwulf cluster. For

contiguous data accesses (Fig. 2.a.), the predicted cost is constant per reference. The prediction error reduced as the number of references increased. The error was mainly due to the approach of counting cache misses pessimistically without taking prefetching into consideration. The prediction error was below 20% for small data and below 4% for large data with this data access pattern.

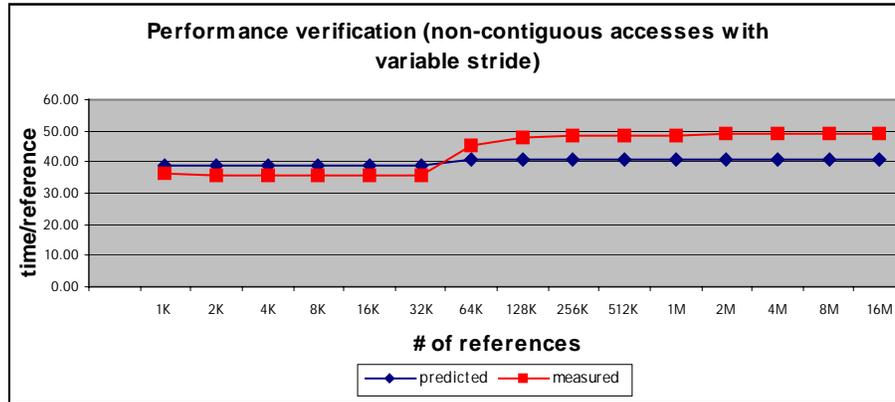


Fig.3. Comparison of measured and predicted memory access cost for non-contiguous data access with fixed word size and variable strides (word size: 8 bytes, stride varies from 1 to 128 bytes periodically)

To test the non-contiguous access pattern performance we used three sizes of fixed strides (16bytes, 32 bytes and 64 bytes) that are equal to L1 cache line size, more than L1 line size and that of equal to L2 line size. For non-contiguous accesses, with stride equal to L1 cache line size, the prediction error reduced as the data size increase. It can be seen from Fig 2.b and Fig 2.c, that the utilization of caches are more effective when the data size is less than L2 cache size. Overall the error is below 20% in most of the cases. For the remaining two non-contiguous access patterns with fixed strides, the prediction error is below 10% for larger data sizes.

For non-contiguous access pattern with variable strides, we initialized an array that contains strides of accesses. Prediction cost of this access pattern contains the cost of accessing non-contiguous arrays as well as the cost of accessing the array of strides. The prediction error is below 15% (Fig 3). This error is caused by missing some of the cache misses in non-contiguous accesses, which requires maintaining the history of the length of cache lines that are already been

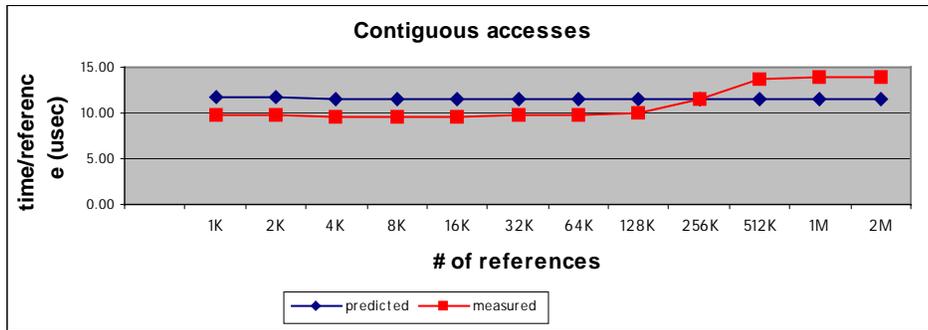


Fig. 4.a.

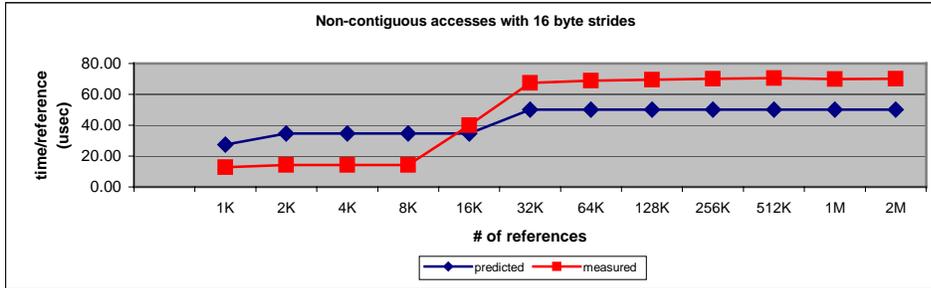


Fig. 4.b.

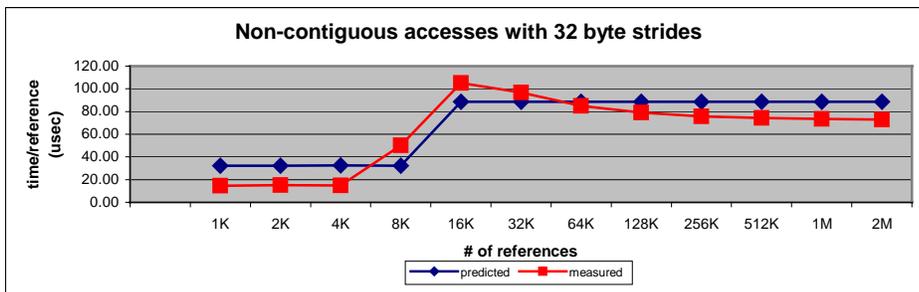


Fig. 4.c.

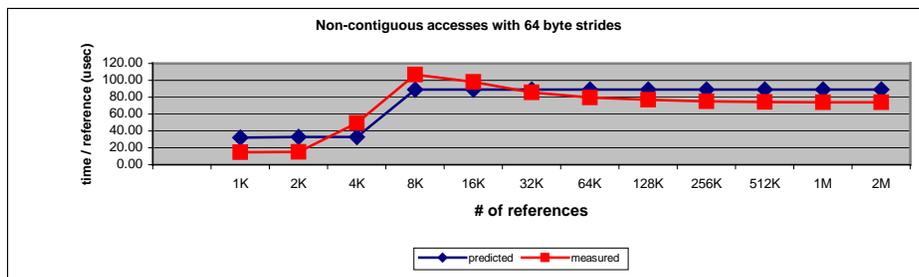


Fig. 4.d.

Fig.4. Comparison of measured and predicted memory access cost on Pentium III processor. The access patterns are: 6.a. Contiguous data access (word size: 1byte, stride: 1byte). 6.b. Non-contiguous data access with fixed word size and stride (word size: 8 bytes, stride: 16 bytes), 6.c. Non-contiguous data access with fixed word size and stride (word size: 8 bytes, stride: 32 bytes) 6.d. Non-contiguous data access with fixed word size and stride (word size: 8 bytes, stride: 64 bytes)

fetched into the cache. Another reason for prediction error for all these access patterns is that we are using average miss penalties, which may not be accurate.

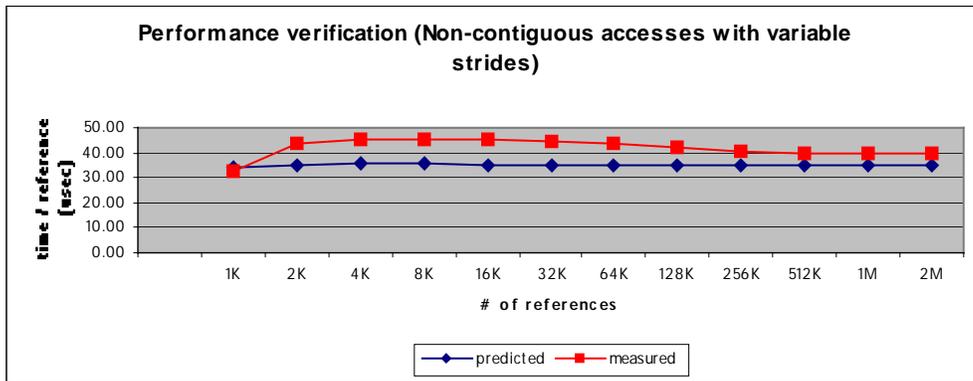


Fig.5. Comparison of measured and predicted memory access cost for non-contiguous data access with fixed word size and variable strides (word size: 8 bytes, stride varies from 1 to 128 bytes periodically)

We observe the similar results on Pentium III processor on Beowulf cluster (Fig 4 and 5). The prediction error is slightly high for small data sizes where the prefetching of this processor is effective. As the working set size increase, the L2 misses increase and the prediction error is below 20% in these cases.

We also verified the performance of the loops in NAS Parallel benchmarks that are performing matrix transpose operation. We have measured the performance two variations of matrix transpose algorithms from NAS Parallel benchmarks' Fast Fourier Transform program. The first algorithm is a simple matrix transpose of copying rows of one matrix to columns of another matrix. The second algorithm uses cache-blocking optimization to improve the

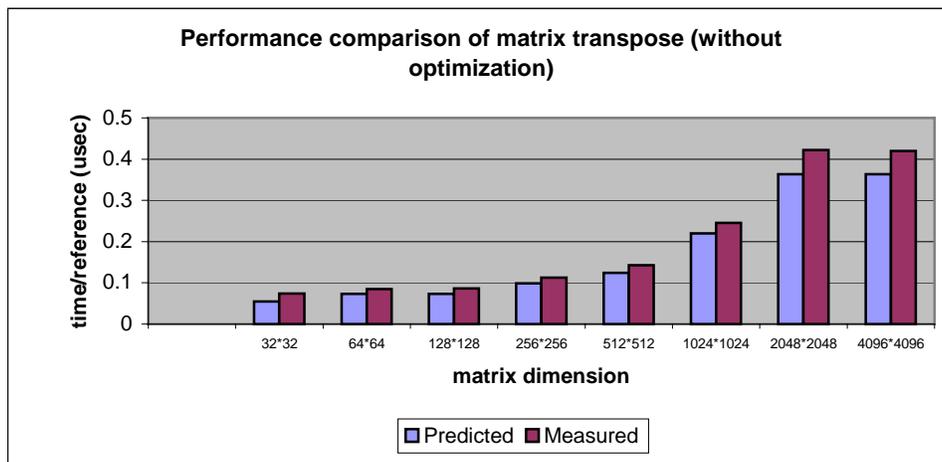


Fig.6. Comparison of measured and predicted memory access cost Matrix transpose algorithm without cache blocking optimization

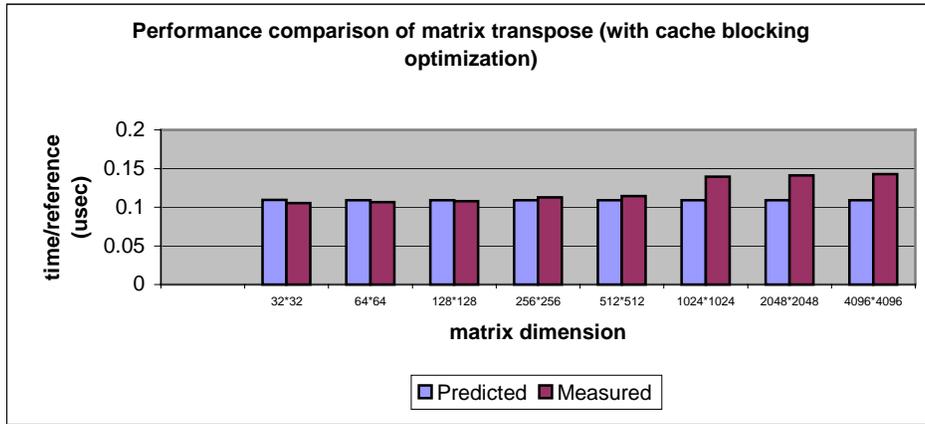


Fig.7. Comparison of measured and predicted memory access cost Matrix transpose algorithm with cache blocking optimization

performance. Both algorithms fit into the data access patterns explained in section 2. The data working set of the first algorithm increases with the dimension of the matrices. Due to the row major ordering of arrays (in C or column major ordering in Fortran), one matrix is accessed contiguously and the other is accessed non-contiguously with fixed stride. The second algorithm makes sure that a block of data is fully utilized before replacing it from the cache. In this algorithm, the two matrices are accessed non-contiguously with fixed strides. However, as the whole data block is reused before it is being replaced, and we chose the block size such that it fits into the cache, the number of cache misses is very less compared to the unoptimized version of matrix transpose. These experiments are performed on Sun UltraSparc IIe processor node.

As expected, the performance (time/reference) increases as the data size increases for the unoptimized transpose algorithm (Fig 6). Predicted values of performance are slightly different from the measured values. The error is around 13%. In the second algorithm, the performance is improved for the transpose algorithm due to the cache-blocking optimization (Fig 7). The performance error was below 5% for most of the data sizes, but increased for large data sizes. This is mainly due the increase in average time per memory reference for the large data sizes.

6. An Application of the model

Parallel communication models such as LogP [Cull96] focus on network communication, with limited consideration of memory communication. Recently, the LogP model was extended to

incorporate memory-communication cost. The *memory-LogP* model formally characterizes the memory-communication cost under four parameters: l : the effective latency, defined as the length of time the processor is engaged in transmission or reception of a message due to the influence of data size (D) and distribution also called as strides (S), $l=f(D,S)$; o : the overhead, defined as the length of time the processor is engaged in transmission or reception of an ideally distributed (contiguous) message (during this time, the processor cannot perform other operations); g : the gap, defined as the minimum time interval between consecutive message receptions at the processor (the reciprocal of g corresponds to the available per processor bandwidth for a given implementation of data transfer on a given system); and P : the number of processor/memory modules (point-to-point communication in the memory hierarchy implies $P=1$). Detailed information about the memory-LogP model can be found in [Kirk03].

The memory-communication cost for sending a data segment depends on architectural parameters, such as cache capacity, and code characteristics, such as data distribution, as explained in the *memory-LogP* model. In general, the overall communication cost includes data-collection overhead, the cost of data copying to the network buffer, the cost of data forwarding to the receiver (network-communication cost), and other costs added by the middleware implementation. When data distribution in memory is noncontiguous, the data is typically collected into a contiguous buffer before being copied to the network buffer. This process adds extra buffering overhead to the overall communication cost and is implementation dependent. The memory access cost predicted in this paper is a part of the latency (l) parameter of the *memory-logP* model.

Currently we apply this model in improving the performance of MPI derived datatypes by optimizing the memory access cost [Byna03]. The MPI Standard [MPI298] supports derived datatypes, which allow users to describe noncontiguous memory layout and communicate noncontiguous data with a single communication function. This feature enables an MPI implementation to optimize the transfer of noncontiguous data. In practice, however, few MPI

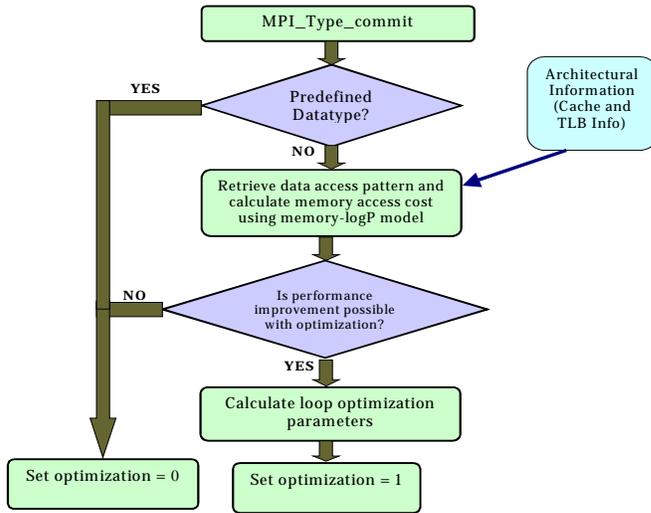


Fig. 8. Predicting the performance of memory access cost and improving the performance of MPI derived datatypes

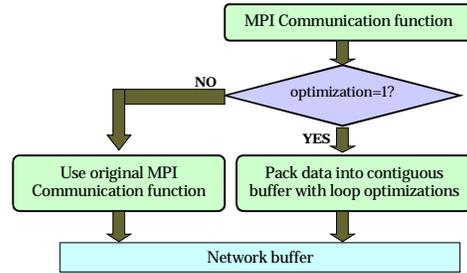


Fig. 9. Packing the data to improve the performance of MPI derived datatypes

implementations provide derived datatypes in a way that performs better than what the user can achieve by manually packing data into a contiguous buffer and then calling an MPI function. Memory access cost has been the reason for this performance bottleneck. We use *memory-logP* model to predict this cost and apply memory access optimization techniques to improve the performance (see Fig 8). When an MPI derived datatype is defined, we convert it into a parameterized loop that fits into one of the data access patterns we described earlier. This loop data structure contains the values of working set size, depth of nested loop, strides, number of iterations etc. We predict the memory access cost using these parameters. We then reorder the loop to optimize the data accesses (e.g. loop tiling, loop interchange etc.) and find the pattern with the least memory access cost. The prediction overhead is very low as our model is using simple formulae.

When the communication function is called, and where optimization is possible, we pack the data using the proposed loop optimization parameters (Fig 9). The packed data is sent to the receiver and the same technique is applied to unpack the data if necessary. Initial results [Byna03]

suggest that this technique is promising and achieving better performance on the advanced architectures such as SGI Origin 2000 and IBM's Blue Horizon.

7. Conclusion

Loop transformations and loop access reordering techniques improve the memory access performance. To obtain these loop optimization parameters, a simple, fast and accurate memory access cost prediction model is necessary. This improves the standard of application level optimizations and reduces the burden on the programmers to learn the rapidly improving processor and computer architecture technology. Towards achieving this goal, in this paper we proposed an analytical model to predict the memory access cost based on the data access patterns. We first classified the most common data access patterns in scientific computing applications. We then proposed a model to predict the memory access cost. We verified this model with measurements and showed that this model is practical. The accuracy of our model is reasonable given its simplicity. We also applied this model to matrix transpose routines in Fast Fourier Transform program of NAS benchmarks, which was implemented in different memory access patterns.

Our model is simple, effective, and easy to be incorporated into memory cost tuning tools, where optimization parameters are to be found at runtime. The prediction errors of 10% to 20% exist, they are reasonably accurate in making optimization decisions. We are currently utilizing this model to improve the performance of MPI derived datatypes, by optimizing the memory access cost. This cost prediction is a part of our memory-logP model, which emphasizes the importance of memory communication performance in point-to-point communication. Our model is practical because of its simplicity. We are able to fit this easily into any optimization library to choose optimization parameters dynamically at runtime. This is not possible with the existing models due to their complexity.

We plan to extend this work in various aspects. We will extend this model to include external and internal conflict misses. We will broaden this model for replacement policies other than LRU, such as FIFO, LFU, MRU, MFU etc. We plan to incorporate this model in a automatic performance tuning system that improves the application performance by optimizing the memory access cost.

Acknowledgments

This work was supported by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Advanced Scientific Computing Research, Office of Science, U.S. Department of Energy, under Contract W-31-109-ENG-38, and in part by a grant from the Office of Advanced Simulation and Computing, National Nuclear Security Administration, U.S. Department of Energy. The U.S. Government retains for itself, and others acting on its behalf, a paid-up, nonexclusive, irrevocable worldwide license in said article to reproduce, prepare derivative works, distribute copies to the public, and perform publicly and display publicly, by or on behalf of the Government.

References

- [Byna03] Surendra Byna, William Gropp, Xian-He Sun, and Rajeev Thakur, "*Improving the Performance of MPI Derived Datatypes by Optimizing Memory-Access Cost*," IEEE International Conference on Cluster Computing, 2003, Hong Kong, December 2003
- [Chat00] S. Chatterjee and S. Sen, "*Cache-Efficient Matrix Transposition*", Proceedings of the 6th International Symposium on High-Performance Computer Architecture, Toulouse, France, January 2000, pages 195-205.
- [Chat01] S. Chatterjee, E. Parker, P. J. Hanlon, and A. R. Lebeck, "Exact Analysis of the Cache Behavior of Nested Loops", Proceedings of the ACM SIGPLAN 2001 Conference on Programming Language Design and Implementation, Snowbird, UT, June 2001

- [Cull96] D.E. Culler et al., “*LogP: A Practical Model of parallel computation*”, Communications of the ACM, vol. 39, pp. 78-85, 1996
- [Ghos99] Somnath Ghosh , Margaret Martonosi , Sharad Malik, “*Cache miss equations: a compiler framework for analyzing and tuning memory behavior*”, ACM Transactions on Programming Languages and Systems (TOPLAS), v.21 n.4, p.703-746, July 1999
- [Jaco96] B.L. Jacob, “*An analytical model for designing memory hierarchies*”, IEEE Transaction on Computers, volume 45, pp. 83-105, 1996.
- [Kand99] M. Kandemir, J. Ramanujam and A. Choudhary, “*Cache Locality by a Combination of Loop and Data Transformations,*” IEEE Transactions on Computers (TC) 48(2): 159–167, February 1999.
- [Kirk03] Kirk W. Cameron, Xian-He Sun, “*Quantifying Locality Effect in Data Access Delay: Memory logP,*” in Proceedings of the 17th International Parallel and Distributed Processing Symposium (IPDPS '03), April 2003.
- [Mark87] Hill, M.D. “*Aspects of cache memory and instruction buffer performance*”, Ph.D. Thesis, University of California, Berkeley, 1987.
- [Mcki96] K.S. McKinley, S.Carr, and C.W. Tseng, “*Improving data locality with loop transformations*”, ACM TOPLAS, 18(4): 424-453. July 1996
- [Mitc01] N. Mitchell, L. Carter, and J. Ferrante, “*A modal model of memory*”. In V.N.Alexandrov, J.J. Dongarra, Computer Science. Springer, May 28-30, 2001
- [MPI298] Message Passing Interface Forum, “*MPI-2: A message passing interface standard*”, High Performance Computing Applications, 12(1-2):1-299, 1998
- [Patt96] D. A. Patterson and J. L. Hennessy, “*Computer Architecture: A quantitative approach*”, 2nd edition. San Francisco, CA: Morgan Kaufmann Publishers, 1996.
- [Peak98] Y. Paek, J. Hoeflinger, and D. Padua, “*Simplification of Array Access Patterns for Compiler Optimizations*”. In Proceedings of the ACM SIGPLAN 98 Conference on Programming Language Design and Implementation, June 1998. 23

- [Ryan02] Rong Yan and Seth C, “*Goldstein Mobile Memory: Improving Memory Locality in Very Large Reconfigurable Fabrics*”, FCCM '02, Napa Valley, CA, April 2002
- [Saav95] Rafael H. Saavedra and Alan Jay Smith, “*Measuring Cache and {TLB} Performance and Their Effect on Benchmark Runtimes*”, IEEE Transactions on Computers, Volume: 44, number: 10, p1223-1235, 1995.
- [Smit82] A.J. Smith, “*Cache Memories*”, Computing Surveys, 14(3), p.473, September 1982
- [SSen00] S. Sen and S. Chatterjee, “*Towards a theory of Cache efficient algorithms*”, SODA, 2000
- [Whal01] R. Clint Whaley, Antoine Petitet, and Jack Dongarra, “*Automated Empirical Optimizations of Software and the ATLAS Project*”, Parallel Computing, Volume 27, Numbers 1-2, pp 3-25, 2001