

Towards a Theory of Cache-Efficient Algorithms ^{*}

Sandeep Sen[†]

Siddhartha Chatterjee[‡]

Neeraj Dumir[§]

Abstract

We describe a model that enables us to analyze the running time of an algorithm in a computer with a memory hierarchy with limited associativity, in terms of various cache parameters. Our model, an extension of Aggarwal and Vitter's I/O model, enables us to establish useful relationships between the cache complexity and the I/O complexity of computations. As a corollary, we obtain cache-optimal algorithms for some fundamental problems like sorting, FFT, and an important subclass of permutations in the single-level cache model. We also show that ignoring associativity concerns could lead to inferior performance, by analyzing the average-case cache behavior of mergesort. We further extend our model to multiple levels of cache with limited associativity and present optimal algorithms for matrix transpose and sorting. Our techniques may be used for systematic exploitation of the memory hierarchy starting from the algorithm design stage, and dealing with the hitherto unresolved problem of limited associativity.

1 Introduction

Models of computation are essential for abstracting the complexity of real machines and enabling the design and analysis of algorithms. The widely-used RAM model owes its longevity and usefulness to its simplicity and robustness. While it is far removed from the complexities of any physical computing device, it successfully predicts the relative performance of algorithms based on an abstract notion of operation counts.

The RAM model assumes a flat memory address space with unit-cost access to any memory location. With the increasing use of caches in modern machines, this assumption grows less justifiable. On modern computers, the running time of a program is as much a function of operation count as of its cache reference pattern. A result of this growing divergence between model and reality is that operation count alone is not always a true predictor of the running time of a program, and manifests itself in anomalies such as a matrix multiplication algorithm demonstrating $O(n^5)$ running time instead of the expected $O(n^3)$ behavior predicted by the RAM model [5]. Such shortcomings of the RAM model motivate us to seek an alternative model that more realistically models the presence of a memory hierarchy. In this paper, we address the issue of better and systematic utilization of caches starting from the algorithm design stage.

^{*}Some of the results in this appeared in a preliminary form in the Proceedings of the Eleventh ACM-SIAM Symposium on Discrete Algorithms 2000 [29].

This work is supported in part by DARPA Grant DABT63-98-1-0001, NSF Grants CDA-97-2637 and CDA-95-12356, The University of North Carolina at Chapel Hill, Duke University, and an equipment donation through Intel Corporation's Technology for Education 2000 Program. The views and conclusions contained herein are those of the authors and should not be interpreted as representing the official policies or endorsements, either expressed or implied, of DARPA or the U.S. Government.

[†]Department of Computer Science and Engineering, IIT Delhi, New Delhi 110016, India. E-mail: ssen@cse.iitd.ernet.in. Part of the work was done when the author was a visiting faculty in the Department of Computer Science, University of North Carolina, Chapel Hill, NC 27599-3175, USA.

[‡]Department of Computer Science, University of North Carolina, Chapel Hill, NC 27599-3175, USA. E-mail: sc@cs.unc.edu.

[§]Department of Computer Science and Engineering, IIT Delhi, New Delhi 110016, India.

A challenge in coming up with a good model is achieving a balance between abstraction and fidelity, so as not to make the model unwieldy for theoretical analysis or simplistic to the point of lack of predictiveness. Memory hierarchy models used by computer architects to design caches have numerous parameters and suffer from the first shortcoming [1, 26]. Early algorithmic work in this area focussed on a two-layered memory model[21]—a very large capacity memory with slow access time (secondary memory) and a limited size faster memory (internal memory). All computation is performed on elements in the internal memory and there is no restriction on placement of elements in the internal memory (fully associative).

The focus of this paper is on the interaction between main memory and *cache*, which is the first level of memory hierarchy once the address is provided by the CPU. The structure of a single level hierarchy of cache memory is adequately characterized by the following three parameters: **Associativity**, **Block size**, and **Capacity**. Capacity and block size are in units of the minimum memory access size (usually one byte). A cache can hold a maximum of C bytes. However, due to physical constraints, the cache is divided into *cache frames* of size B that contain B contiguous bytes of memory—called a *memory block*. The associativity A specifies the number of different frames in which a memory block can reside. If a block can reside in any frame (i.e., $A = \frac{C}{B}$), the cache is said to be *fully associative*; if $A = 1$, the cache is said to be *direct-mapped*; otherwise, the cache is *A-way set associative*.

For a given memory access, the hardware inspects the cache to determine if the corresponding memory element is resident in the cache. This is accomplished by using an indexing function to locate the appropriate set of cache frames that may contain the memory block. If the memory block is not resident, a cache miss is said to occur. From an architectural standpoint, cache misses can be classified into one of three classes [20].

- A *compulsory miss* (also called a *cold miss*) is one that is caused by referencing a previously unreference memory block. Eliminating a compulsory miss requires prefetching the data, either by an explicit prefetch operation or by placing more data items in a single memory block.
- A reference that is not a compulsory miss but misses in a fully-associative cache with LRU replacement is classified as a *capacity miss*. Capacity misses are caused by referencing more memory blocks than can fit in the cache. Restructuring the program to re-use blocks while they are in cache can reduce capacity misses.
- A reference that is not a compulsory miss that hits in a fully-associative cache but misses in an A -way set-associative cache is classified as a *conflict miss*. A conflict miss to block X indicates that block X has been referenced in the recent past, since it is contained in the fully-associative cache, but at least A other memory blocks that map to the same cache set have been accessed since the last reference to block X. Eliminating conflict misses requires transforming the program to change either the memory allocation and/or layout of the two arrays (so that contemporaneous accesses do not compete for the same sets) or the manner in which the arrays are accessed.

Conflict misses pose an additional challenge in designing efficient algorithms in the cache. This class of misses is not present in the I/O models, where the mapping between internal and external memory is fully associative.

Existing memory hierarchy models [4, 2, 3, 5] do not model certain salient features of caches, notably the lack of full associativity in address mapping and the lack of explicit control over data movement and replacement. Unfortunately, these small differences are malign in the effect.¹ The *conflict misses* that they introduce make analysis of algorithms much more difficult [16]. Carter and Gatlin [9] conclude a recent paper saying

¹See the discussion in [9] on a simple matrix transpose program.

What is needed next is a study of “messy details” not modeled by UMH (particularly cache associativity) that are important to the performance of the remaining steps of the FFT algorithm.

In the first part of this paper, we develop a two-level memory hierarchy model to capture the interaction between cache and main memory. Our model is a simple extension of the two-level I/O model that Aggarwal and Vitter [4] proposed for analyzing external memory algorithms. However, it captures three additional constraints of caches: lower miss penalties; lack of full associativity in address mapping; and lack of explicit program control over data movement. The work in this paper shows that the constraint imposed by limited associativity can be tackled quite elegantly, allowing us to extend the results of the I/O model to the cache model very efficiently.

Most modern architectures have a memory hierarchy consisting of multiple cache levels. In the second half of this paper, we extend the two-level cache model to a multi-level cache model.

The remainder of this paper is organized as follows. Section 2 surveys related work. Section 3 defines our cache model and establishes an efficient emulation scheme between the I/O model and our cache model. As direct corollaries of the emulation scheme, we obtain cache-optimal algorithms for several fundamental problems such as sorting, FFT, and an important class of permutations. Section 4 illustrates the importance of the emulation scheme by demonstrating that a direct (*i.e.*, bypassing the emulation) implementation of an I/O-optimal sorting algorithm (multiway mergesort) is provably inferior, even in the average case, in the cache model. Section 5 describes a natural extension of our model to multiple levels of caches. We present an algorithm for transposing a matrix in the multi-level cache model that attains optimal performance in the presence of any number of levels of cache memory. Our algorithm is not cache-oblivious, *i.e.*, we do make explicit use of the sizes of the cache at various levels. Next, we show that with some simple modifications, the funnel-sort algorithm of Frigo *et al.* attains optimal performance in a single level (direct mapped) cache in an oblivious sense, *i.e.*, without prior knowledge of memory parameters. Finally, Section 6 presents conclusions, possible refinements to the model, and directions for future work.

2 Related work

The I/O model assumes that most of the data resides on disk and has to be transferred to main memory to do any processing. Because of the tremendous difference in speeds, it ignores the cost of internal processing and counts only the number of I/Os. Floyd [15] originally defined a formal model and proved tight bounds on the number of I/Os required to transpose a matrix using two pages of internal memory. Hong and Kung [21] extended this model and studied the I/O complexity of FFT when the internal memory size is bounded by M . Aggarwal and Vitter [4] further refined the model by incorporating an additional parameter B , the number of (contiguous) elements transferred in a single I/O operation. They gave upper and lower bounds on the number of I/Os for several fundamental problems including sorting, selection, matrix transposition, and FFT. Following their work, researchers have designed I/O-optimal algorithms for fundamental problems in graph theory [13] and computational geometry [19].

Researchers have also modeled multiple levels of memory hierarchy. Aggarwal *et al.* [2] defined the *Hierarchical Memory Model* (HMM) that assigns a function $f(x)$ to accessing location x in the memory, where f is a monotonically increasing function. This can be regarded as a continuous analog of the multi-level hierarchy. Aggarwal *et al.* [3] added the capability of block transfer to the HMM, which enabled them to obtain faster algorithms. Alpern *et al.* [5] described the *Uniform Memory Hierarchy* (UMH) model, where the access costs differ in discrete steps. Very recently, Frigo *et al.* [18] presented an alternate strategy of algorithm design on these models which has the added advantage that explicit values of parameters related to different levels of the memory hierarchy are not required. Bilardi and Peserico [8] investigate further the

complexity of designing algorithms without the knowledge architectural parameters.² Other attempts were directed towards extracting better performance by parallel memory hierarchies [32, 33, 14], where several blocks could be transferred simultaneously.

Ladner *et al.* [23] describe a stochastic model for performance analysis in cache. Our work is different in nature, as we follow a more traditional worst-case analysis. Our analysis of sorting in Section 4 provides a better theoretical basis for some of the experimental work of LaMarca and Ladner [25].

To the best of our knowledge, the only other paper that addresses the problem of limited associativity in cache is recent work of Mehlhorn and Sanders[27]. They show that for a class of algorithms based on merging multiple sequences, the I/O algorithms can be made nearly optimal by use of a simple randomized shift technique. The emulation theorem in Section 3 of this paper not only provides a deterministic solution for the same class of algorithms, but also works for a very general situation. The results in [27] are nevertheless interesting from the perspective of implementation.

3 The cache model

The (two-level) I/O model of Aggarwal and Vitter [4] captures the interaction between a slow (secondary) memory of infinite capacity and a fast (primary) memory of limited capacity. It is characterized by two parameters: M , the capacity of the fast memory; and B , the size of data transfers between slow and fast memories. Such data movement operations are called *I/O operations* or *block transfers*. The use of the model is meaningful when the problem size $N \gg M$.

The I/O model contains the following further assumptions.

1. A datum can be used in a computation iff it is present in fast memory. All data initially resides in slow memory. Data can be transferred between slow and fast memory (in either direction) by I/O operations.
2. Since the latency for accessing slow memory is very high, the average cost of transfer per element can be reduced by transferring a block of B elements at little additional cost. This may not be as useful as it may seem at first sight, since these B elements are not arbitrary, but are contiguous in memory. The onus is on the programmer to use all the elements, as traditional RAM algorithms are not designed for such restricted memory access patterns. We denote the map from a memory address to its block address by \mathbb{B} . The internal memory can hold at least three blocks, i.e., $M \geq 3 \cdot B$.
3. The computation cost is ignored in comparison to the cost of an I/O operation. This is justified by the high access latency of slow memory.
4. A block of data from slow memory can be placed in any block of fast memory.
5. I/O operations are explicit in the algorithm.

The goal of algorithm design in this model is to minimize the number of I/O operations.

We adopt much of the framework of the I/O model in developing a cache model to capture the interactions between cache and main memory. In this case, the cache is the fast memory, while main memory is the slow memory. Assumptions 1 and 2 of the I/O model continue to hold in our cache model. However, assumptions 3–5 are no longer valid and need to be replaced as follows.

- The difference between the access times of slow and fast memory is considerably smaller than in the I/O model, namely a factor of 5–100 rather than factor of 10000. We will use L to denote the

²However, none of these models address the problem of limited associativity in cache.

normalized cache latency. This cost function assigns a cost of 1 for accessing an element in cache and L for accessing an element in the main memory. This way, we also account for the computation in cache.

- Main memory blocks are mapped into cache sets using a *fixed* and pre-determined mapping function that is implemented in hardware. Typically, this is a modulo mapping based on the low-order address bits. However, the results of this section will hold as long as there is a *fixed* address mapping function that distributes the main memory evenly in the cache. We denote this mapping from main memory blocks to cache sets by \mathbb{S} . We will occasionally slightly abuse this notation and apply \mathbb{S} directly to a memory address x rather than to $\mathbb{B}(x)$.
- The cache is not visible to the programmer (not even at the assembly level). When a program issues a reference to a memory location x , an *image* (copy) of the main memory block $b = \mathbb{B}(x)$ is brought into the cache set $\mathbb{S}(b)$ if it is not already present there. The block b continues to reside in cache until it is evicted by another block b' that is mapped to the same cache set (*i.e.*, $\mathbb{S}(b) = \mathbb{S}(b')$). In other words, a cache set c contains the latest memory block referenced that is mapped to this set.

To summarize, we use the notation $\mathfrak{C}(M, B, L)$ to denote our three-parameter cache model, and the notation $\mathfrak{J}(M, B)$ to denote the I/O model with parameters M and B . We will use n and m to denote N/B and M/B respectively. The assumptions of our cache model parallel those of the I/O model, except as noted above.³ The goal of algorithm design in the cache model is to minimize *running time*, defined as the number of cache accesses plus L times the number of main memory accesses.

3.1 Emulating I/O algorithms

The differences between the two models listed above would appear to frustrate any efforts to naively map an I/O algorithm to the cache model, given that we neither have the control nor the flexibility of the I/O model. Our main result in this section establishes a connection between the I/O model and the cache model using a very simple emulation scheme.

Theorem 3.1 (Emulation Theorem) *An algorithm A in $\mathfrak{J}(M, B)$ using T block transfers and I processing time can be converted to an equivalent algorithm A^c in $\mathfrak{C}(M, B, L)$ that runs in $O(I + (L + B) \cdot T)$ steps. The memory requirement of A^c is an additional $m + 2$ blocks beyond that of A .*

Proof: Note that I is usually not accounted for in the I/O model, but we will keep track of the internal memory computation done in A in our emulation. The idea behind the emulation is as follows. We will mimic the behavior of the I/O algorithm A in the cache model, using an array Buf of m blocks to play the role of the fast memory. We will view the main memory in the cache model as an array Mem of B -element blocks. Although Buf is also part of the memory, we are using different notations to make their roles explicit in this proof. Likewise, we will view the cache as an array of sets and denote the i th set by $C[i]$.

As discussed above, we do not have explicit control on the contents of the cache locations. However, we can control the memory access pattern through a level of indirection so as to maintain a 1-1 correspondence between Buf and the cache. Wlog, we assume that \mathbb{S} maps block i of Buf to cache set $C[i]$ for $i \in [1, m]$.

We divide the I/O algorithm into rounds, where in each round, the I/O algorithm A transfers a block between the slow memory and the fast memory and (possibly) does some computations. The cache algorithm A^c transfers the same blocks between Mem and Buf and then does the identical computations in Buf . Figure 1 formally describes the procedure. Note that the B elements must be explicitly copied in the cache model.

³Frigo *et al.* [18] independently arrive at a very similar parameterization of their model.

Round t of the emulation

<i>I/O Algorithm A</i>	<i>Cache Emulation A^c</i>
1. Transfer block b_t from slow memory to block a_t of the fast memory	1. Copy contents of the B locations of $Mem[b_t]$ into $Buf[a_t]$
2. Perform computations in fast memory	2. Perform identical computations in Buf

Figure 1: The emulation scheme used in the proof of Theorem 3.1.

It must be obvious that the final outcome of algorithm A^c is the same as algorithm A . The more interesting issue is the cost of the emulation.

A block of size B is transferred into cache if its image does not exist in the cache at the time of reference. The invariant that we try to maintain at the end of each round is that there is a 1-1 correspondence between Buf and C . This will ensure that all the I operations are done within the cache at minimal cost.

Assume that we have maintained the above invariant at the end of round $t - 1$. In round t , we transfer block $Mem[b_t]$ into $Buf[a_t]$. Accessing the memory block $Mem[b_t]$ will displace the existing block in cache set $C[q]$, where $q = \mathbb{S}(b_t)$. From the invariant, we know that the block displaced from $C[q]$ is $Buf[q]$, which must be restored to cache to restore the invariant. We can bring it back by a single memory reference and charge this to the round t itself, which is L . (Actually it will be brought back during the subsequent reference, so the previous step is only to simplify the accounting.)

The cost of copying $Mem[b_t]$ to $Buf[a_t]$ is $L + B$ assuming that $Mem[b_t]$ and $Buf[a_t]$ are not mapped to the same cache set ($\mathbb{S}(b_t) \neq \mathbb{S}(a_t)$). Otherwise it will cause alternate cache misses (*thrashing*) of the blocks $Mem[b_t]$ and $Buf[a_t]$ leading to $L \cdot B$ steps for copying. This can be prevented by transferring through an intermediate memory block $Mem[Y]$ such that $\mathbb{S}(Y) \neq \mathbb{S}(b_t)$. Having two such intermediate buffers that map to distinct cache sets would suffice in all cases. So, we first transfer $Mem[b_t]$ to $Mem[Y]$ followed by $Mem[Y]$ to $Buf[j]$. The first copying has cost $2L + B$ since both blocks must be fetched from main memory. The second transfer is between blocks, one of which is present in the cache, so it has cost $L + B$. To this we must also add cost L for restoring the block of Buf that was mapped to the same cache set as $Mem[Y]$. So, the total cost of the *safe* method is $4L + 2B$.

The internal processing remains identical. If I_t denotes the internal processing cost of step t , the total cost of the emulation is at most $\sum_{t=1}^T (I_t + 2(L + B) + 2L) = I + 4L \cdot T + 2B \cdot T$. \square

Remark 1

- A possible alternative to using intermediate memory-resident buffers to avoid thrashing is to use registers, since register access is much faster. In particular, if we have B registers, then we can save two extra memory accesses, bringing down the emulation cost to $2L + 2B$.
- We can make the emulation somewhat simpler by using a randomized mapping scheme. That is, if we choose the starting location of array Buf randomly, then the probability that $Mem[b_t]$ and $Buf[a_t]$ have the same image is $1/M$. So the expected emulation cost is $I + 2L \cdot T + (B + (LB)/M) \cdot T$ without using any intermediate copying.
- The basic idea of copying data into contiguous memory locations to reduce interference misses has been exploited before in some specific contexts like matrix multiplication [24] and bit-reversal permutation [9]. Theorem 3.1 unifies these previous results within a common framework.

The term $O(B \cdot T)$ is subsumed by $O(I)$ if computation is done on at least a constant fraction of the elements in the block transferred by the I/O algorithm. This is usually the case for efficient I/O algorithms. We will call such I/O algorithms *block-efficient*.

Corollary 3.2 *A block-efficient I/O algorithm for $\mathfrak{J}(M, B)$ that uses T block transfers and I processing can be emulated in $\mathfrak{C}(M, B, L)$ in $O(I + L \cdot T)$ steps.*

Remark 2 The algorithms for sorting, FFT, matrix transposition, and matrix multiplication described in Aggarwal and Vitter [4] are block-efficient.

3.2 Extension to set-associative cache

The trend in modern memory architectures is to allow limited flexibility in the address mapping between memory blocks and cache sets. The k -way set-associative cache has the property that a memory block can reside in any (one) of k cache frames. Thus, $k = 1$ corresponds to the direct-mapped cache we have considered so far, while $k = m$ corresponds to a fully associative cache. Values of k for data caches are generally small, usually in the range 1–4.

If all the k sets are occupied, a replacement policy like LRU is used (by the hardware) to find an assignment for the referenced block. The emulation technique of the previous section would extend to this scenario easily if we had explicit control on the replacement. This not being the case, we shall tackle it indirectly by making use of an useful property of LRU that Frigo *et al.* [18] exploited in the context of designing cache-oblivious algorithms for a fully associative cache.

Lemma 3.1 (Sleator-Tarjan[30]) *For any sequence s , F_{LRU} , the number of misses incurred by LRU with cache size n_{LRU} is no more than $(n_{LRU}/(n_{LRU} - n_{OPT} + 1))F_{OPT}$, where F_{OPT} is the minimum number of misses by an optimal replacement strategy with cache size n_{OPT} .*

We use this lemma in the following way. We run the emulation technique for only half the cache size, *i.e.*, we choose the buffer to be of size $m/2$, such that for every k cache lines in a set, we have only $k/2$ buffer blocks. From Lemma 3.1, we know that the number of misses in each cache set is no more than twice the optimal, which is in turn bounded by the number of misses incurred by the I/O algorithm.

Theorem 3.3 (Generalized Emulation Theorem) *An algorithm A in $\mathfrak{J}(M/2, B)$ using T block transfers and I processing time can be converted to an equivalent algorithm \mathcal{A} in the k -way set-associative cache model with parameters M, B, L that runs in $O(I + (L + B) \cdot T)$ steps. The memory requirement of \mathcal{A} is an additional $m/2 + 2$ blocks beyond that of A .*

3.3 The cache complexity of sorting and other problems

Aggarwal and Vitter [4] prove the following lower bound for sorting and FFT in the I/O model.

Lemma 3.2 ([4]) *The average-case and the worst-case number of I/O's required for sorting N records and for computing the N -input FFT graph in $\mathfrak{J}(M, B)$ is $\Omega\left(\frac{N}{B} \frac{\log(1+N/B)}{\log(1+M/B)}\right)$.*

Theorem 3.4 *The lower bound for sorting in $\mathfrak{C}(M, B, L)$ is $\Omega(N \log N + L \frac{N}{B} \frac{\log N/B}{\log M/B})$.*

Proof: Any lower bound in the number of block transfers in $\mathfrak{J}(M, B)$ carries over to $\mathfrak{C}(M, B, L)$. Since the lower bound is the maximum of the lower bound on number of comparisons and the bound in Lemma 3.2, the theorem follows by dividing the sum of the two terms by 2. \square

Theorem 3.5 In $\mathfrak{C}(M, B, L)$, N numbers can be sorted in $O(N \log N + L \cdot \frac{N}{B} \cdot \frac{\log N/B}{\log M/B})$ steps and this is optimal.

Proof: The M/B -way mergesort algorithm described in Aggarwal and Vitter [4] has an I/O complexity of $O(\frac{N}{B} \frac{\log N/B}{\log M/B})$. The processing time involves maintaining a heap of size M/B and $O(\log M/B)$ per output element. For N elements, the number of phases is $\frac{\log N}{\log M/B}$, so the total processing time is $O(N \log N)$. From Corollary 3.2, and Remark 2, the cost of this algorithm in the cache model is $O(N \log N + L \cdot \frac{N}{B} \cdot \frac{\log N/B}{\log M/B})$. Optimality follows from Theorem 3.4. \square

Remark 3 The M/B -way distribution sort (multiway quicksort) also has the same upper bound.

We can prove a similar result for FFT computations.

Theorem 3.6 The FFT of N numbers can be computed in $O(N \log N + L \cdot \frac{N \log N/B}{B \log M/B})$ in $\mathfrak{C}(M, B, L)$.

Remark 4 The FFTW algorithm [17] is optimal only for $B = 1$. Barve [6] has independently obtained a similar result.

The class of Bit Matrix Multiply Complement (BMMC) permutations include many important permutations like matrix transposition and bit reversal. Combining the work of Cormen *et al.* [14] with our emulation scheme, we obtain the following result.

Theorem 3.7 The class of BMMC permutations for N elements can be achieved in $\Theta\left(N + L \cdot \frac{N}{B} \frac{\log M}{\log(M/B)}\right)$ steps in $\mathfrak{C}(M, B, L)$.

Remark 5 Many known geometric [13] and graph algorithms [19] in the I/O model, such as convex hull and graph connectivity, can be transformed optimally into the cache model.

4 Average-case performance of mergesort in the cache model

In this section, we analyze the average-case performance of k -way mergesort in the cache model. Of the three classes of misses described in Section 1, we note that compulsory misses are unavoidable and that capacity misses are minimized while designing algorithms for the I/O model. We are therefore interested in bounding the number of conflict misses for a straightforward implementation of the I/O-optimal k -way mergesort algorithm. It is easy to construct a worst-case input permutation where there will be a conflict miss for every input element (a cyclic distribution suffices), so the average case is more interesting.

We assume that s cache sets are available for the leading blocks of the k runs S_1, \dots, S_k . In other words, we ignore the misses caused by heap operations (or equivalently ensure that the heap area in the cache does not overlap with the runs).

We create a random instance of the input as follows. Consider the sequence $\{1, \dots, N\}$, and distribute the elements of this sequence to runs by traversing the sequence in increasing order and assigning element i to run S_j with probability $1/k$. From the nature of our construction, each run S_i is sorted. We denote j -th element of S_i as $S_{i,j}$. The expected number of elements in any run S_i is N/k .

During the k -way merge, the leading blocks are critical in the sense that the heap is built on the *leading element* of every sequence S_i . The leading element of a sequence is the smallest element that has not been added to the merged (output) sequence. The *leading block* is the cache line containing the leading element. Let b_i denote the leading block of run S_i . *Conflict* can occur when the leading blocks of different sequences

are mapped to the same cache set. In particular, a *conflict miss* occurs for element $S_{i,j+1}$ when there is at least one element $x \in b_k$, for some $k \neq i$, such that $S_{i,j} < x < S_{i,j+1}$ and $\mathbb{S}(b_i) = \mathbb{S}(b_k)$. (We do not count conflict misses for the first element in the leading block, *i.e.*, $S_{i,j}$ and $S_{i,j+1}$ must belong to the same block, but we will not be very strict about this in our calculations.)

Let p_i denote the probability of conflict for element $i \in [1, N]$. Using indicator random variables X_i to count the conflict miss for element i , the total number of conflict misses $X = \sum_i X_i$. It follows that the expected number of conflict misses $E[X] = \sum_i E[X_i] = \sum_i p_i$. In the remaining section we will try to estimate a lower bound on p_i for i large enough to validate the following assumption.

A1 The cache sets of the leading blocks, $\mathbb{S}(b_i)$, are randomly distributed in cache sets $1, \dots, s$ independent of the other sorted runs. Moreover, the exact position of the leading element within the leading block is also uniformly distributed in positions $\{1, \dots, sB\}$.

Remark 6 A recent variation of the mergesort algorithm (see [7]) actually satisfies **A1** by its very nature. So, the present analysis is directly applicable to its average-case performance in cache. A similar observation was made independently by Sanders [27] who obtained upper-bounds for mergesort for a set associative cache.

From our previous discussion and the definition of a conflict miss, we would like to compute the probability of the following event.

E1 For some i, j , for all elements x , such that $S_{i,j} < x < S_{i,j+1}$, $\mathbb{S}(x) \neq \mathbb{S}(S_{i,j})$.

In other words, none of the leading blocks of the sorted sequences S_j , $j \neq i$, conflicts with b_i . The probability of the complement of this event (*i.e.*, $\Pr[\overline{E1}]$) is the probability that we want to estimate. We will compute an upper bound on $\Pr[E1]$, under the assumption **A1**, thus deriving a lower bound on $\Pr[\overline{E1}]$.

Lemma 4.1 For $k/s > \epsilon$, $\Pr[E1] < 1 - \delta$, where ϵ and δ are positive constants (dependent only on s and k but not on n or B).

Proof: See Appendix A. □

Thus we can state the main result of this section as follows.

Theorem 4.1 The expected number of conflict misses in a random input for doing a k -way merge in an s -set direct-mapped cache, where k is $\Omega(s)$, is $\Omega(N)$, where N is the total number of elements in all the k sequences. Therefore the (ordinary I/O-optimal) M/B -way mergesort in an M/B -set cache will exhibit $O(N \frac{\log N/B}{\log M/B})$ cache misses which is asymptotically larger than the optimal value of $O(\frac{N}{B} \frac{\log N/B}{\log M/B})$.

Proof: The probability of conflict misses is $\Omega(1)$ when k is $\Omega(s)$. Therefore the expected total number of conflict misses is $\Omega(N)$ for N elements. The I/O-optimal mergesort uses M/B -way merging at each of the $\frac{\log N/B}{\log M/B}$ levels, hence the second part of the theorem follows. □

Remark 7 Intuitively, by choosing $k \ll s$, we can minimize the probability of conflict misses resulting in an increased number of merge phases (and hence running time). This underlines the critical role of conflict misses *vis-a-vis* capacity misses that forces us to use only a small fraction of the available cache. Recently, Sanders [27] has shown that by choosing k to be $O(\frac{M}{B^{1+1/\alpha}})$ in an α -way set associative cache with a modified version of mergesort of [7], the expected number of conflict misses per phase can be bounded by $O(N/B)$. In comparison, the use of the emulation theorem guarantees minimal worst-case conflict misses while making good use of cache.

5 The Multi-level Cache Model

Most modern architectures have a memory hierarchy consisting of multiple levels of cache. Consider two cache levels \mathcal{L}_1 and \mathcal{L}_2 preceding main memory, with \mathcal{L}_1 being faster and smaller. The operation of the memory hierarchy in this case is as follows. The memory location being referenced is first looked up in \mathcal{L}_1 . If it is not present in \mathcal{L}_1 , then it is searched for in \mathcal{L}_2 (these can be overlapped with appropriate hardware support). If the item is not present in \mathcal{L}_1 but it is in \mathcal{L}_2 , then it is brought into \mathcal{L}_1 . In case that it is not in \mathcal{L}_2 , then a cache line is brought in from main memory into \mathcal{L}_2 and into \mathcal{L}_1 . The size of cache line brought into \mathcal{L}_2 (denoted by B_2) is usually no smaller than the one brought into \mathcal{L}_1 (denoted by B_1). The expectation is that the more frequently used items will remain in the faster cache.

The Multi-level Cache Model is an extension to multiple cache levels of the previously introduced Cache Model. Let \mathcal{L}_i denote the i -th level of cache memory. The parameters involved here are the problem size N , the size of \mathcal{L}_i which is denoted by M_i , the frame size (unit of allocation) of \mathcal{L}_i denoted by B_i and the latency factor l_i . If a data item is present in the \mathcal{L}_i , then it is present in \mathcal{L}_j for all $j \geq i$ (sometimes referred to as the **inclusion property**). If it is not present in \mathcal{L}_i , then the cost for a miss is l_i plus the cost of fetching it from \mathcal{L}_{i+1} (if it is present in \mathcal{L}_{i+1} , then this cost is zero). For convenience, the latency factor l_i is the ratio of time taken on a miss from the i -th level to the amount of time taken for a unit operation.

Figure 2 shows the memory mapping for a two-level cache architecture. The shaded part of main memory is of size B_1 and therefore occupies only a part of a line of the \mathcal{L}_2 cache which is of size B_2 . There is a natural generalization of the memory mapping to multiple levels of cache.

We make the following assumptions in this section, which are consistent with existing architectures.

A1. For all i , B_i, L_i are powers of 2.

A2. $2B_i \leq B_{i+1}$ and the number of Cache Lines $L_i \leq L_{i+1}$.

A3. $B_k \leq L_1$ and $4B_k \leq B_1L_1$ (i.e. $B_1 \geq 4$) where \mathcal{L}_k is the largest and slowest cache. This implies that

$$L_i \cdot B_i \geq B_k \cdot B_i \quad (1)$$

This will be useful for the analysis of the algorithms and are sometimes termed as *tall cache* in reference to the aspect ratio.

5.1 Matrix Transpose

In this section, we provide an approach for transposing a matrix in the Multi-level Cache Model.

The trivial lower bound for matrix transposition of an $N \times N$ matrix in the multi-level cache hierarchy is clearly the time to scan N^2 elements, namely,

$$\Omega\left(\sum_i \frac{N^2}{B_i} l_i\right)$$

where

B_i is the number of elements in one cache line in \mathcal{L}_i cache; L_i is the number of cache lines in \mathcal{L}_i cache, which is $\frac{M_i}{B_i}$; and l_i is the latency for \mathcal{L}_i cache.

Our algorithm uses a more general form of the emulation theorem to get the submatrices to fit into cache in a regular fashion. The work in this section shows that it is possible to handle the constraints imposed by limited associativity even in a multi-level cache model.

We subdivide the matrix into $B_k \times B_k$ submatrices. Thus we get $\lceil n/B_k \rceil \times \lceil n/B_k \rceil$ submatrices from an $n \times n$ submatrix.

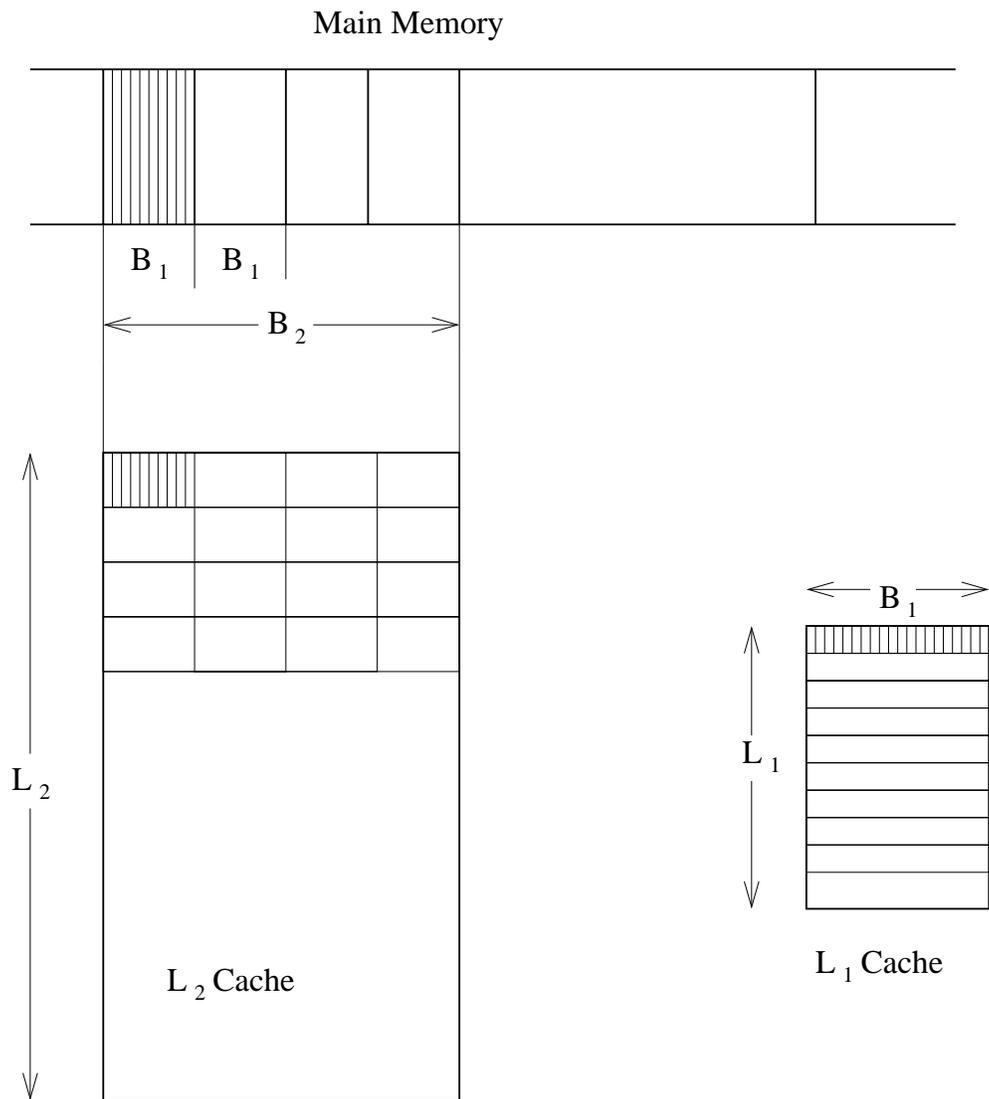


Figure 2: Memory mapping in a two-level cache hierarchy

$$A = \begin{pmatrix} a_1 & a_2 & \dots & \dots & a_n \\ a_{n+1} & a_{n+2} & \dots & \dots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{n^2-n+1} & \dots & \dots & \dots & a_{n^2} \end{pmatrix} = \begin{pmatrix} A_1 & A_2 & \dots & A_{n/B} \\ \vdots & \vdots & \vdots & \vdots \\ A_{n^2-nB/B} & \dots & \dots & A_{n^2/B^2} \end{pmatrix}$$

Note that the submatrices in the last row and column need not be square as one side may have $\leq B$ rows or columns.

Let $m = n/B$ then

$$A^T = \begin{pmatrix} A_1^T & A_{m+1}^T & \dots & \dots & A_{m^2-m+1}^T \\ A_2^T & A_{m+2}^T & \dots & \dots & A_{2m}^T \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ A_m^T & \dots & \dots & \dots & A_{m^2}^T \end{pmatrix}$$

For simplicity, we describe the algorithm as transposing a square matrix A in another matrix B , i.e. $B = A^T$. The main procedure is **Rec.Trans**(A, B, s), where A is transposed into B by dividing A and B into s^2 submatrices and then recursively transposing the sub-matrices. Let $A_{i,j}$ ($B_{i,j}$) denote the submatrices for $1 \leq i, j \leq s$. Then $B = A^T$ can be computed as **Rec.Trans**($A_{i,j}, B_{j,i}, s'$) for all i, j and some appropriate s' which depends on B_k and B_{k-1} . In general, if t_k, t_{k-1}, \dots, t_1 denote the values of s' at the $1, 2 \dots$ level of recursion, then $t_i = B_{i+1}/B_i$. If the submatrices are $B_1 \times B_1$ (base case), then perform the transpose exchange of the symmetric submatrices directly. We perform matrix transpose as follows, which is similar to the familiar recursive transpose algorithm.

1. Subdivide the matrix as shown into $B_k \times B_k$ submatrices.
2. Move the symmetric submatrices to contiguous memory locations.
3. **Rec.Trans**($A_{i,j}, B_{j,i}, B_k/B_{k-1}$).
4. Write back the $B_k \times B_k$ submatrices to original locations.

In the following subsections we analyze the data movement of this algorithm to bound the number of cache misses at various levels.

5.2 Moving a submatrix to contiguous locations

To move a submatrix we will move it cache line by cache line. By choice of size of submatrices ($B_k \times B_k$) each row will be an array of size B_k , but the rows themselves may be far apart.

Lemma 5.1 *If two memory blocks x and y of size B_k are aligned in \mathcal{L}_k -cache map to the same cache set in \mathcal{L}_i -cache for some $1 \leq i \leq k$, then x and y map to the same set in each \mathcal{L}_j -cache for all $1 \leq j \leq i$.*

Proof: If x and y map to the same cache set in \mathcal{L}_i cache then their i -th level memory block numbers (to be denoted by $b^i(x)$ and $b^i(y)$) differ by a multiple of L_i . Let $b^i(x) - b^i(y) = \alpha L_i$. Since $L_j | L_i$ (both are powers of two), $b^i(x) - b^i(y) = \beta L_j$ where $\beta = \alpha \cdot L_i / L_j$. Let x', y' be the *corresponding* sub-blocks of x and y at the j -th level. Then their block numbers $b^j(x'), b^j(y')$ differ by $B_i/B_j \cdot \beta \cdot L_j$, i.e., a multiple of L_j as $B_j | B_i$. Note that blocks are aligned across different levels of cache. Therefore x and y also collide in \mathcal{L}_j . \square

Corollary 5.1 *If two blocks of size B_k that are aligned in \mathcal{L}_k -cache do not conflict in level i they do not conflict in any level j for all $i \leq j \leq k$.*

Theorem 5.2 *There is an algorithm which moves a set of blocks of size B_k (where there are k levels of cache with block size B_i for each $1 \leq i \leq k$) into a contiguous area in main memory in*

$$O\left(\sum \frac{N}{B_i} l_i\right)$$

where N is the total data moved and l_i is the cost of a cache miss for the i^{th} level of cache.

Proof: Let the set of blocks of size B_k be I (we are assuming that the blocks are aligned). Let the target block in the contiguous area for each block $i \in I$ be in the corresponding set J where each block $j \in J$ is also aligned with a cache line in \mathcal{L}_k Cache.

Let block a map to $R_{b,a}$, $b = \{1, 2, \dots, k\}$ where $R_{b,a}$ denote the set of cache lines in the \mathcal{L}_b -cache. (Since a is of size B_k , it will occupy several blocks in lower levels of cache.)

Let the i^{th} block map to set $R_{k,i}$ of the \mathcal{L}_k Cache. Let the target block j map to set $R_{k,j}$. In the worst case, $R_{k,j}$ is equal to $R_{k,i}$. Thus in this case the line $R_{k,i}$ has to be moved to a temporary block say x (mapped to $R_{k,x}$) and then moved back to $R_{k,j}$. We choose x such that $R_{1,x}$ and $R_{1,i}$ do not conflict and also $R_{1,x}$ and $R_{1,j}$ do not conflict. Such a choice of x is always possible because our temporary storage area X of size $4B_k$ has at least 4 lines of \mathcal{L}_k -cache (i and j will take up two blocks of \mathcal{L}_k -cache, thus leaving at least one block free to be used as temporary storage). *This is why we have the assumption that $4B_k \leq B_1 L_1$.* That is, by dividing the \mathcal{L}_1 -cache into $B_1 L_1 / B_k$ zones, there is always a zone free for x .

For convenience of analysis, we maintain the invariant that X is always in \mathcal{L}_k -cache. By application of the previous corollary on our choice of x (such that $R_{1,i} \neq R_{1,x} \neq R_{1,j}$) we also have $R_{a,i} \neq R_{a,x} \neq R_{a,j}$ for all $1 \leq a \leq k$. Thus we can move i to x and x to j without any conflict misses. The number of cache misses involved is three for each level—one for getting the i^{th} block, one for writing the j^{th} block, and one to maintain the invariant since we have to touch the line displaced by i . Thus we get a factor of 3.

Thus the cost of this process is

$$3\left(\sum \frac{N}{B_i} l_i\right)$$

where N is the amount of data moved. □

Remark 8 For blocks I that are not aligned in \mathcal{L}_k Cache, the constant would increase to 4 since we would need to bring up to 2 cache lines for each $i \in I$. The rest of the proof would remain the same.

Corollary 5.3 *A $B_k \times B_k$ submatrix can be moved into contiguous locations in the memory in $O\left(\sum_{i=1}^{i=k} \frac{B_k^2}{B_i} l_i\right)$ time in a computer that has k levels of (direct-mapped) cache.*

This follows from the preceding discussion. We allocate memory say C of size $B_k \times B_k$ for placing the submatrix and memory, say, X of size $4B_k$ for temporary storage and keep both these areas distinct.

Remark 9 If we have set associativity (≥ 2) in all levels of cache then we do not need an intermediate buffer x as line i and j can both reside in cache simultaneously and movement from one to the other will not cause thrashing. Thus the constant will come down to two. Since at any point in time we will only be dealing with two cache lines and will not need the lines i or j once we have read or written to them the replacement policy of the cache does not affect our algorithm.

Remark 10 If the capacity of the register file is greater than the size of the cache line (B_k) of the outermost cache level (\mathcal{L}_k) then we can move data without worrying about collision by copying from line i to registers and then from registers to line j . Thus even in this case the constant will come down to two.

Once we have the submatrices in contiguous locations we perform the transpose as follows. For each of the submatrices we divide the $B_r \times B_r$ submatrix (say S) in level \mathcal{L}_r (for $2 \leq r \leq k$) further into $B_{r-1} \times B_{r-1}$ size submatrices as before. Each $B_{r-1} \times B_{r-1}$ size subsubmatrix fits into \mathcal{L}_{r-1} cache completely (since $B_{r-1} \cdot B_{r-1} \leq B_{r-1} \cdot B_k \leq B_{r-1} \cdot L_{r-1}$ from equation (1)). Let $B_r/B_{r-1} = k_r$.

Thus we have the submatrices as

$$\begin{pmatrix} S_{1,1} & S_{1,2} & \dots & S_{1,k_r} \\ \vdots & \vdots & \vdots & \vdots \\ S_{k_r,1} & \dots & \dots & S_{k_r,k_r} \end{pmatrix}$$

So we perform matrix transpose of each $S_{i,j}$ in place without incurring any misses as it resides completely inside the cache. Once we have transposed each $S_{i,j}$ we exchange $S_{i,j}$ with $S_{j,i}$. We will show that $S_{i,j}$ and $S_{j,i}$ can not conflict in \mathcal{L}_{r-1} -cache for $i \neq j$.

Since $S_{i,j}$ and $S_{j,i}$ lie in different parts of the \mathcal{L}_r -cache lines, they will map to different cache sets in the \mathcal{L}_{r-1} -cache. The rows of $S_{i,j}$ and $S_{j,i}$ correspond to $(iB_{r-1} + a_1)k_r + j$ and $(jB_{r-1} + a_2)k_r + i$ where $a_1, a_2 \in \{1, 2, \dots, B_{r-1}\}$ and

$$B_r/B_{r-1} = k_r.$$

If these conflict then

$$(iB_{r-1} + a_1)k_r + j \equiv (jB_{r-1} + a_2)k_r + i \pmod{L_{r-1}}.$$

Since $B_{r-1} = 2^u$ and $B_r = 2^v$ and $\mathcal{L}_{r-1} = 2^w$ (all powers of two)

$$k_r = 2^{v-u}$$

Therefore k_r divides \mathcal{L}_{r-1} (because $k_r = B_r/B_{r-1} < B_r \leq L_{r-1}$). Hence

$$j \equiv i \pmod{k_r}.$$

Since $i, j \leq k_r$ the above implies

$$i = j.$$

Note that $S_{i,i}$'s do not have to be exchanged. Thus, we have shown that a $B_r \times B_r$ matrix can be divided into $B_{r-1} \times B_{r-1}$ which completely fits into \mathcal{L}_{r-1} -cache. Moreover, the symmetric sub-matrices do not interfere with each other. The same argument can be extended to any $B_j \times B_j$ submatrix for $j < r$. Applying this recursively we end up dividing the $B_k \times B_k$ size matrix in \mathcal{L}_k -cache to $B_1 \times B_1$ sized submatrices in \mathcal{L}_1 -cache, which can then be transposed and exchanged easily. From the preceding discussion, the corresponding submatrices do not interfere in any level of the cache.

(Note that even though we keep subdividing the matrix at every cache level recursively and claim that we then have the submatrices in cache and can take the transpose and exchange them, the actual movement, i.e., transpose and exchange happens only at the \mathcal{L}_1 -cache level, where the submatrices are of size $B_1 \times B_1$.)

The time taken by this operation is

$$\sum \frac{N^2}{B_i} l_i.$$

This is because each $S_{i,j}$ and $S_{j,i}$ pair (such that $i \neq j$) has to be brought into \mathcal{L}_{r-1} Cache only once for transposing and exchanging of $B_1 \times B_1$ submatrices. Similarly, at any level of cache, a block from the

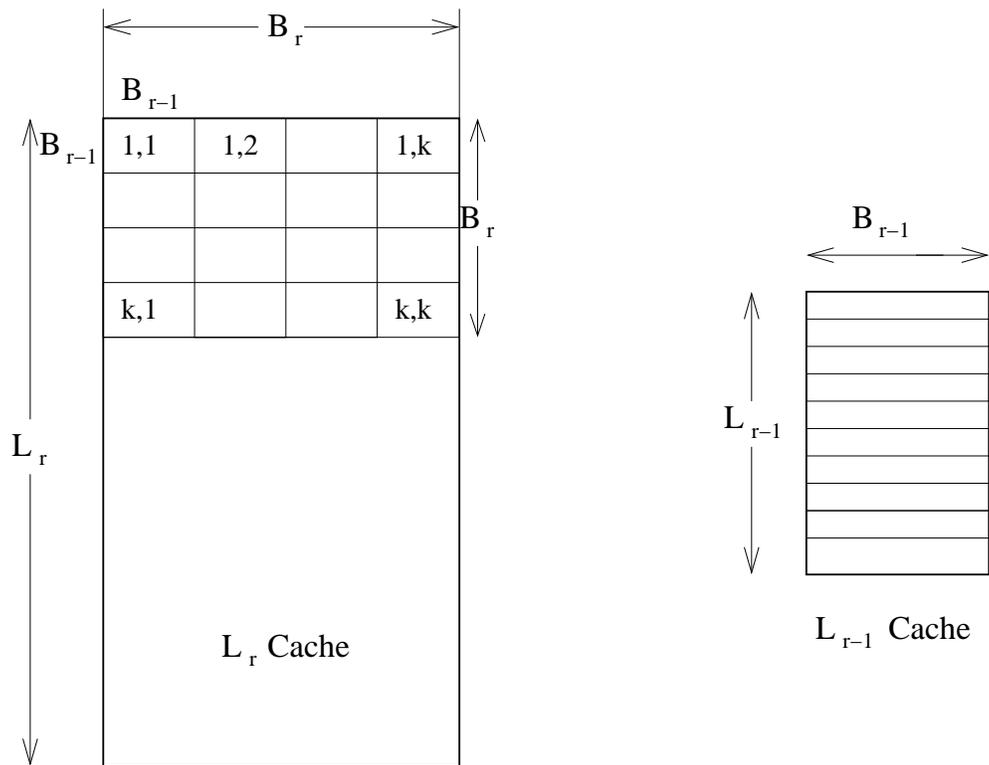


Figure 3: Positions of symmetric submatrices in Cache

matrix is brought in only once. The sequence of the recursive calls ensures that each cache line is used completely as we move from sub-matrix to sub-matrix.

Finally, we move the transposed symmetric submatrices of size $B_k \times B_k$ to their location in memory, *i.e.*, reverse the process of bringing in blocks of size B_k from random locations to a contiguous block. This procedure is exactly the same as in Theorem 5.2 in the previous section that has the constant 3.

Remark 11 The above constant of 3 for writing back the matrix to an appropriate location depends on the assumption that we can keep the two symmetric submatrices of size $B_k \times B_k$ in contiguous locations at the same time. This would allow us to exchange the matrices during the write back stage. If we are restricted to a contiguous temporary space of size $B_k \times B_k$ only, then we will have to move the data twice, incurring the cost twice.

Remark 12 Even though in the above analysis we have always assumed a square matrix of size $N \times N$ the algorithm works correctly without any change for transposing a matrix of size $M \times N$ if we are transposing a matrix A and storing it in B . This is because the same analysis of subdividing into submatrices of size $B_k \times B_k$ and transposing still holds. However if we want to transpose a $M \times N$ matrix in place then the algorithm fails because the location to write back to would not be obvious and the approach used here would fail.

Theorem 5.4 *The algorithm for matrix transpose runs in*

$$O\left(\sum_{i=1}^{i=k} \frac{N^2}{B_i} l_i\right) + O(N^2)$$

steps in a computer that has k levels of direct-mapped cache.

If we have temporary storage space of size $2B_k \times B_k + 4B_k$ and assume block alignment of all submatrices then the constant is 7. This includes 3 for initial movement to contiguous location, 1 for transposing the symmetric submatrices of size $B_k \times B_k$ and 3 for writing back the transposed submatrix to its original location. Note that the constant is independent of the number of levels of cache.

Remark 13 Even if we have set associativity (≥ 2) in any level of cache the analysis goes through as before (though the constants will come down for data copying to contiguous locations). For the transposing and exchange of symmetric submatrices the set associativity will not come into play because we need a line only once in the cache and are using only 2 lines at a given time. So either LRU or even FIFO replacement policy would only evict a line that we have already finished using.

5.3 Sorting in multiple levels

We first consider a restriction of the model described above where data cannot be transferred simultaneously across non-consecutive cache levels. We use C_i to denote $\sum_{j=1}^{j=i} M_j$.

Theorem 5.5 *The lower bound for sorting in the restricted multi-level cache model is $\Omega(N \log N + \sum_{i=1}^k l_i \cdot \frac{N \log N/B_i}{B_i \log C_i/B_i})$.*

Proof: The proof of Aggarwal and Vitter can be modified to disregard block transfers that merely rearrange data in the external memory. Then it can be applied separately to each cache level, noting that the data transfer in the higher levels do not contribute for any given level. \square

These lower bounds are in the same spirit as those of Vitter and Nodine [32] (for the S-UMH model) and Savage [28], that is, the lower bounds do not capture the simultaneous interaction of the different levels.

If we remove this restriction, then the following can be proved along similar lines as Theorem 3.4.

Lemma 5.2 *The lower bound for sorting in the multi-level cache model is*

$$\Omega\left(\max_{i=1}^k \left\{ N \log N, \ell_i \cdot \frac{N \cdot \log N / B_i}{B_i \log C_i / B_i} \right\}\right).$$

□

This bound appears weak if k is large. To rectify this, we observe the following. Across each cache boundary, the minimum number of I/Os follow from Aggarwal and Vitter's arguments. The difficulty arises in the multi-level model as a block transfer in level i propagates in all levels $j < i$ although the block sizes are different. The minimum number of I/Os from (the highest) level k remains unaffected, namely, $\frac{N}{B_k} \frac{\log N / B_k}{\log C_k / B_k}$. For level $k - 1$, we will subtract this number from the lower bound of $\frac{N}{B_{k-1}} \frac{\log N / B_{k-1}}{\log C_{k-1} / B_{k-1}}$. Continuing in this fashion, we obtain the following lower bound.

Theorem 5.6 *The lower bound for sorting in the multi-level cache model is*

$$\Omega\left(N \log N + \sum_{i=1}^k \ell_i \cdot \left(\frac{N \cdot \log N / B_i}{B_i \log C_i / B_i} - \left(\sum_{j=i+1}^k \frac{N \cdot \log N / B_j}{B_j \log C_j / B_j} \right) \right)\right).$$

□

If we further assume that $\frac{C_i}{C_{i-1}} \geq \frac{B_i}{B_{i-1}} \geq 3$, we obtain a relatively simple expression that resembles Theorem 5.5. Note that the consecutive terms in the expression in the second summation of the previous lemma decrease by a factor of 3.

Corollary 5.7 *The lower bound for sorting in the multi-level cache model with geometrically decreasing cache sizes and cache lines is $\Omega(N \log N + \frac{1}{2} \sum_{i=1}^k \ell_i \cdot \frac{N \cdot \log N / B_i}{B_i \log C_i / B_i})$.* □

Theorem 5.8 *In a multi-level cache, where the B_i blocks are composed of B_{i-1} blocks, we can sort in expected time $O\left(N \log N + \left(\frac{\log N / B_1}{\log M_1 / B_1}\right) \cdot \sum_{i=1}^k \ell_i \cdot \frac{N}{B_i}\right)$.*

Proof: We perform a M_1/B_1 -way mergesort using the variation proposed by Barve *et al.* [7] in the context of parallel disk I/Os. The main idea is to shift each sorted stream cyclically by a random amount R_i for the i th stream. If $R_i \in [0, M_k - 1]$, then the leading element is in any of the cache sets with equal likelihood. Like Barve *et al.* [7], we divide the merging into phases where a phase outputs m elements, where m is the merge degree. In the previous section we counted the number of conflict misses for the input streams, since we could exploit symmetry based on the random input. It is difficult to extend the previous arguments to a worst case input. However, it can be shown easily that if $\frac{m}{s} < \frac{1}{m^s}$ (where s is the number of cache sets), the expected number of conflict misses is $O(1)$ in each phase. So the total expected number of cache misses is $O(N/B_i)$ in the level i cache for all $1 \leq i \leq k$.

The cost of writing a block of size B_1 from level k is spread across several levels. The cost of transferring B_k/B_1 blocks of size B_1 from level k is $\ell_k + \ell_{k-1} \frac{B_k}{B_{k-1}} + \ell_{k-2} \frac{B_k}{B_{k-1}} \frac{B_{k-1}}{B_{k-2}} + \dots + \ell_1 \frac{B_k}{B_1}$. Amortizing this cost over B_k/B_1 transfers gives us the required result. Recall that $O\left(N/B_1 \left(\frac{\log N / B_1}{\log M_1 / B_1}\right)\right)$ B_1 block transfers suffice for $(M_1/B_1)^{1/3}$ -way mergesort. □

Remark 14 This bound is reasonably close to that of Corollary 5.7 if we ignore constant factors. Extending this to the more general emulation scheme of Theorem 3.1 is not immediate as we require the block transfers across various cache boundaries to have a nice pattern, namely the *sub-block* property. This is satisfied by the mergesort and quicksort and a number of other algorithms but cannot be assumed in general.

5.4 Cache-oblivious sorting

In this section, we will focus on a two-level cache model that has limited associativity. One of the *cache-oblivious* algorithms presented by Frigo *et al.* [18] is the funnel sort algorithm. They showed that the algorithm is optimal in the I/O model (which is fully associative). However it is not clear whether the optimality holds in the cache model. In this section, we show that, with some simple modification, funnel sort is optimal even in the direct-mapped cache model.

The funnel sort algorithm can be described as follows.

- Split the input into $n^{1/3}$ contiguous arrays of size $n^{2/3}$ and sort these arrays recursively.
- Merge the $n^{1/3}$ sorted sequences using a $n^{1/3}$ -merger, where a k -merger works as follows.

A k -merger operates by recursively merging sorted sequences. Unlike mergesort, a k -merger stops working on a merging sub-problem when the merged output sequence becomes “long enough” and resumes working on another merging sub-problem (see Figure 4).

INVARIANT The invocation of a k -merger outputs the first k^3 elements of the sorted sequence obtained by merging the k input sequences.

BASE CASE $k = 2$ producing $k^3 = 8$ elements whenever invoked.

NOTE The intermediate buffers are twice the size of the output obtained by a $k^{1/2}$ merger.

To output k^3 elements, the k -merger is invoked $k^{3/2}$ times. Before each invocation the k -merger fills each buffer that is less than half full so that every buffer has at least $k^{3/2}$ elements—the number of elements to be merged in that invocation.

Frigo *et al.* [18] have shown that the above algorithm (that does not make explicit use of the various memory-size parameters) is optimal in the I/O model. However, the I/O model does not account for conflict misses since it assumes full associativity. This could be a degrading influence in the presence of limited associativity (in particular direct-mapping).

5.4.1 Structure of k -merger

It is sufficient to get a bound on cache misses in the cache model since the bounds for capacity misses in the cache model are the same as the bounds shown in the I/O model.

Let us get an idea of what the structure of a k -merger looks like by looking at a 16-merger (see Figure 5). A k -merger, unrolled, consists of 2-mergers arranged in a tree-like fashion. Since the number of 2-mergers gets halved at each level and the initial input sequences are k in number there are $\lg k$ levels.

Lemma 5.3 *If the buffers are randomly placed and the starting position is also randomly chosen (since the buffers are cyclic this is easy to do) the probability of conflict misses is maximized if the buffers are less than one cache line long.*

The worst case for conflict misses occurs when the buffers are less than one cache line in size. This is because if the buffers collide then all data that goes through them will thrash. If however the size of the buffers were greater than one cache line then even if some two elements collide the probability of future collisions would depend upon the data input or the relative movement of data in the two buffers. The

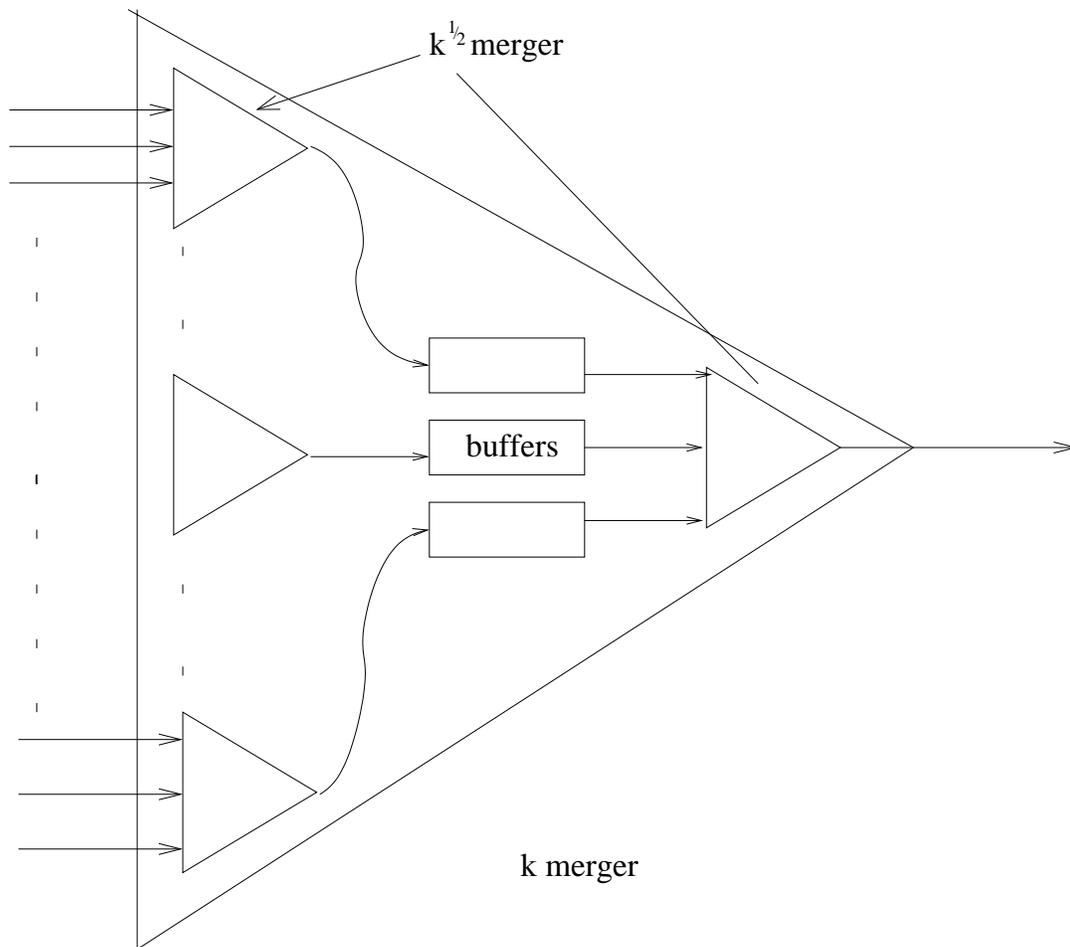


Figure 4: Recursive definition of a k -merger in terms of $k^{1/2}$ -mergers

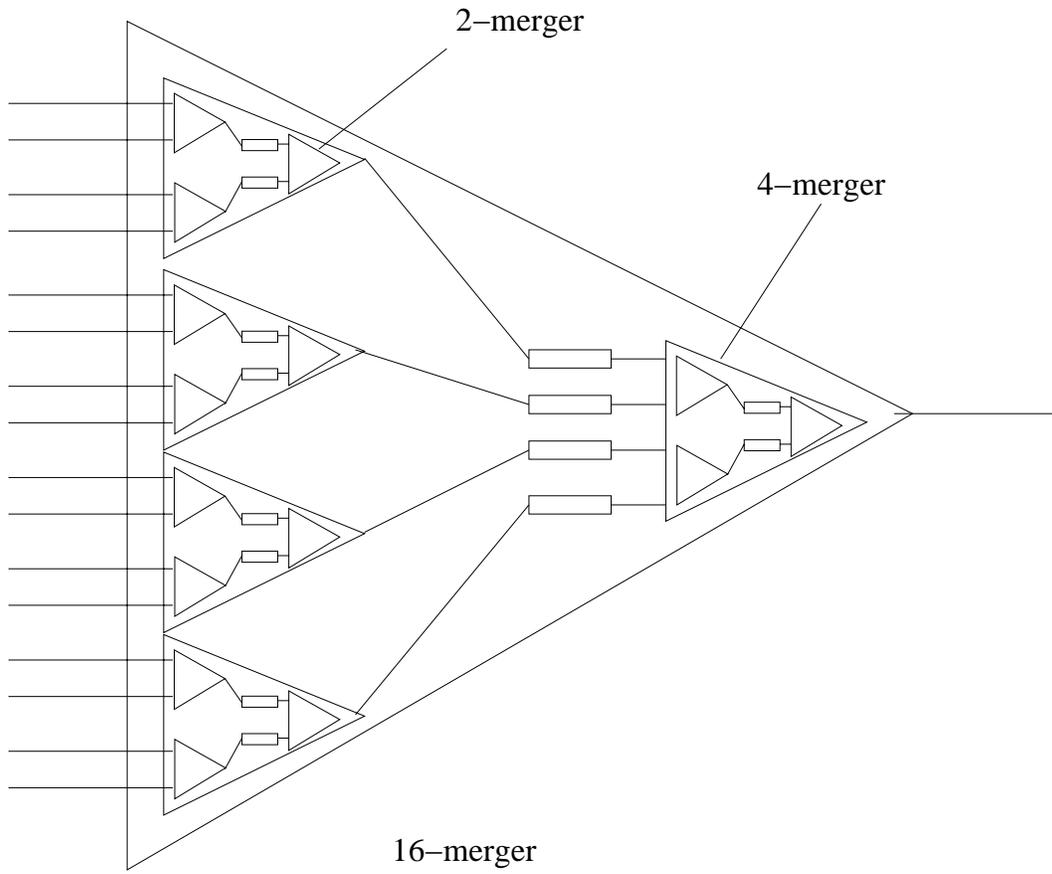


Figure 5: Expansion of a 16-merger into 2-mergers

probability of conflict miss is maximized when the buffers are less than one cache line. Then probability of conflict is $1/m$, where m is equal to the cache size M divided by the cache line size B , *i.e.*, the number of cache lines.

5.4.2 Bounding conflict misses

The analysis for compulsory and capacity misses goes through without change from the I/O model to the cache model. Thus, funnel sort is optimal in the cache model if the conflict misses can be bounded by

$$\frac{N}{B} \times \frac{\log N/B}{\log M/B}$$

Lemma 5.4 *If the cache is 3-way or more set associative, there will be no conflict misses for a 2-way merger.*

Proof: The two input buffers and the output buffer, even if they map to the same cache set can reside simultaneously in the cache. Since at any stage only one 2-merger is active there will be no conflict misses at all and the cache misses will only be in the form of capacity or compulsory misses. \square

5.4.3 Direct-Mapped case

For an input of size N , a $N^{1/3}$ -merger is created. The number of levels in such a merger is $\log N^{1/3}$ (*i.e.*, the number of levels of the tree in the unrolled merger). Every element that travels through the $N^{1/3}$ -merger sees $\log N^{1/3}$ 2-mergers (see Figure 6). For an element passing through a 2-merger there are 3 buffers that could collide. We *charge* an element for a conflict miss if it is swapped out of the cache before it passes to the output buffer or collides with the output buffer when it is being output. So the expected number of collisions is 3C_2 times the probability of collision between any two buffers (two input and one output). Thus the expected number of collisions for a single element passing through a 2-merger is ${}^3C_2 \times 1/m \leq 3/m$ where $m = M/B$.

If $x_{i,j}$ is the probability of a cache miss for element i in level j then summing over all elements and all levels we get

$$\begin{aligned} E \left(\sum_{i=1}^N \sum_{j=1}^{N^{1/3}} x_{i,j} \right) &= \sum_{i=1}^N \sum_{j=1}^{\log N^{1/3}} E(x_{i,j}) \\ &\leq \sum_{i=1}^N \sum_{j=1}^{\log N^{1/3}} \frac{3}{m} = \frac{3N}{m} \times \log N^{1/3} \\ &= O \left(\frac{N}{m} \times \log N \right) \end{aligned}$$

Lemma 5.5 *The expected performance of funnel sort is optimal in the direct-mapped cache model if $\log \frac{M}{B} \leq \frac{M}{B^2 \log B}$. It is also optimal for a 3-way associative cache.*

Proof: If M and B are such that

$$\log \frac{M}{B} \leq \frac{M}{B^2 \log B}$$

we have the total number of conflict misses

$$\frac{N \log N}{m} = \frac{N \log N}{B \log B \frac{M}{B^2 \log B}} \leq \frac{N}{B} \times \frac{\log N/B}{\log M/B}$$

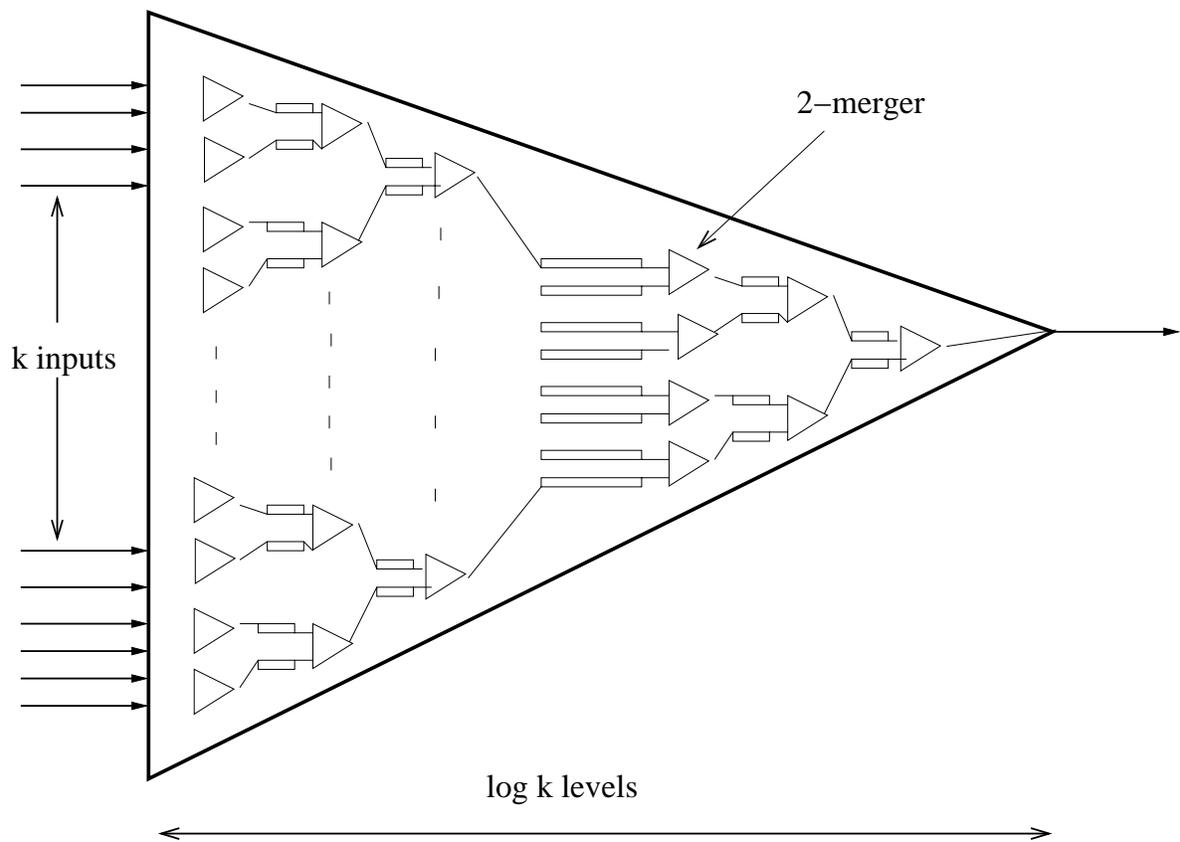


Figure 6: A k -merger expanded out into 2-mergers

Note that the condition is satisfied for $M \geq B^{2+\epsilon}$ for any fixed $\epsilon > 0$ which is similar to the *tall-cache* assumption made by Frigo *et al.*.

The set associative case is proved by Lemma 5.4. □

The same analysis is applicable between successive levels \mathcal{L}_i and \mathcal{L}_{i+1} of a multi-level cache model. This yields an optimal algorithm for sorting in the multilevel cache model.

Theorem 5.9 *In a multi-level cache model, the number of cache misses at level \mathcal{L}_i in the funnel sort algorithm can be bounded by $\frac{N \log(N/B_i)}{B_i \log(M_i/B_i)}$.*

This bound matches the lower bound of Lemma 5.5 within a constant factor, which makes it an optimal algorithm when simultaneous transfers are not allowed across multiple levels.

6 Conclusions

We have presented a cache model for designing and analyzing algorithms. Our model, while closely related to the I/O model of Aggarwal and Vitter, incorporates three additional salient features of cache: lower miss penalty, limited associativity, and lack of direct program control over data movement. We have established an emulation scheme that allows us to systematically convert an I/O-efficient algorithm into a cache-efficient algorithm. This emulation provides a generic starting point for cache-conscious algorithm design; it may be possible to further improve cache performance by problem-specific techniques to control interference misses. We have also demonstrated the relevance of the emulation scheme by demonstrating that a direct mapping of an I/O-efficient algorithm does not guarantee a cache-efficient algorithm. Finally, we have extended our basic cache model to multiple cache levels.

Our single-level cache model is based on a blocking direct-mapped cache that does not distinguish between reads and writes. Modeling a non-blocking cache or distinguishing between reads and writes would appear to require queuing-theoretic extensions and does not appear to be appropriate at the algorithm design stage. The *translation lookaside buffer* or TLB is another important cache in real systems that caches virtual-to-physical address translations. Its peculiar aspect ratio and high miss penalty raise different concerns for algorithm design. Our preliminary experiments with certain permutation problems suggests that TLBs are important to model and can contribute significantly to program running times.

We have begun to implement some of these algorithms to validate the theory on real machines, and also using cache simulation tools like *fast-cache*, *ATOM*, or *cprof*. Preliminary observations indicate that our predictions are more accurate with respect to miss ratios than actual running times (see [12]). We have traced a number of possible reasons for this. First, because the cache miss latencies are not astronomical, it is important to keep track of the constant factors. An algorithmic variation that guarantees lack of conflict misses at the expense of doubling the number of memory references may turn out to be slower than the original algorithm. Second, our preliminary experiments with certain permutation problems suggests that TLBs are important to model and can contribute significantly to program running times. Third, several low-level details hidden by the compiler related to instruction scheduling, array address computations, and alignment of data structures in memory can significantly influence running times. As argued earlier, these factors are more appropriate to tackle at the level of implementation than algorithm design.

Several of the cache problems we observe can be traced to the simple array layout schemes used in current programming languages. It has shown elsewhere [10, 11, 31] that nonlinear array layout schemes based on quadrant-based decomposition are better suited for hierarchical memory systems. Further study of such array layouts is a promising direction for future research.

Acknowledgments

We are grateful to Alvin Lebeck for valuable discussions related to present and future trends of different aspects of memory hierarchy design. We would like to acknowledge Rakesh Barve for discussions related to sorting, FFTW, and BRP. The first author would also like to thank Jeff Vitter for his comments on an earlier draft of this paper.

References

- [1] A. Agarwal, M. Horowitz, and J. Hennessy. An analytical cache model. *ACM Trans. Comput. Syst.*, 7(2):184–215, May 1989.
- [2] A. Aggarwal, B. Alpern, A. Chandra, and M. Snir. A model for hierarchical memory. In *Proceedings of ACM Symposium on Theory of Computing*, pages 305–314, 1987.
- [3] A. Aggarwal, A. Chandra, and M. Snir. Hierarchical memory with block transfer. In *Proceedings of IEEE Foundations of Computer Science*, pages 204–216, 1987.
- [4] A. Aggarwal and J. Vitter. The input/output complexity of sorting and related problems. *Commun. ACM*, 31(5):1118–1127, 1988.
- [5] B. Alpern, L. Carter, E. Feig, and T. Selker. The uniform memory hierarchy model of computation. *Algorithmica*, 12(2):72–109, 1994.
- [6] R. Barve. Private communication.
- [7] R. Barve, E. Grove, and J. Vitter. Simple randomized mergesort on parallel disks. *Parallel Computing*, 23(4):109–118, 1997. A preliminary version appeared in SPAA 96.
- [8] G. Bilardi and E. Peserico. Efficient portability across memory hierarchies, 2000. Unpublished manuscript.
- [9] L. Carter and K. Gatlin. Towards an optimal bit-reversal permutation program. In *Proceeding of IEEE Foundations of Computer Science*, 1998.
- [10] S. Chatterjee, V. V. Jain, A. R. Lebeck, S. Mundhra, and M. Thottethodi. Nonlinear array layouts for hierarchical memory systems. In *Proceedings of the 1999 ACM International Conference on Supercomputing*, pages 444–453, Rhodes, Greece, June 1999.
- [11] S. Chatterjee, A. R. Lebeck, P. K. Patnala, and M. Thottethodi. Recursive array layouts and fast parallel matrix multiplication. In *Proceedings of Eleventh Annual ACM Symposium on Parallel Algorithms and Architectures*, pages 222–231, Saint-Malo, France, June 1999.
- [12] S. Chatterjee and S. Sen. Cache-efficient matrix transposition. In *Proceedings of HPCA-6*, pages 195–205, Toulouse, France, Jan. 2000.
- [13] Y. Chiang, M. Goodrich, E. Grove, R. Tamassia, D. Vengroff, and J. Vitter. External memory graph algorithms. In *Proceedings of the ACM-SIAM Symposium of Discrete Algorithms*, pages 139–149, 1995.
- [14] T. H. Cormen, T. Sundquist, and L. F. Wisniewski. Asymptotically tight bounds for performing BMMC permutations on parallel disk systems. *SIAM Journal of Computing*, 28(1):105–136, 1999.
- [15] R. Floyd. Permuting information in idealized two-level storage. In R. E. Miller and J. W. Thatcher, editors, *Complexity of Computer Computations*, pages 105–109. Plenum Press, New York, NY, 1972.
- [16] C. Fricker, O. Temam, and W. Jalby. Influence of cross-interference on blocked loops: A case study with matrix-vector multiply. *ACM Trans. Prog. Lang. Syst.*, 17(4):561–575, July 1995.
- [17] M. Frigo and S. G. Johnson. FFTW: An adaptive software architecture for the FFT. In *Proceedings of ICASSP’98*, volume 3, page 1381, Seattle, WA, 1998. IEEE.
- [18] M. Frigo, C. E. Leiserson, H. Prokop, and S. Ramachandran. Cache-oblivious algorithms. In *Proceedings of the 40th Annual Symposium on Foundations of Computer Science (FOCS ’99)*, New York, NY, Oct. 1999.
- [19] M. Goodrich, J. Tsay, D. Vengroff, and J. Vitter. External memory computational geometry. In *Proceeding of IEEE Foundations of Computer Science*, pages 714–723, 1993.
- [20] M. D. Hill and A. J. Smith. Evaluating associativity in CPU caches. *IEEE Trans. Comput.*, C-38(12):1612–1630, Dec. 1989.
- [21] J. Hong and H. Kung. I/O complexity: The red blue pebble game. In *Proceedings of ACM Symposium on Theory of Computing*, 1981.

- [22] A. Kamath, R. Motwani, K. Palem, and P. Spirakis. Tail bounds for occupancy and the satisfiability threshold conjecture. In *Proceeding of IEEE Foundations of Computer Science*, pages 592–603, 1994.
- [23] R. Ladner, J. Fix, and A. LaMarca. Cache performance analysis of algorithms. In *Proceedings of the ACM-SIAM Symposium of Discrete Algorithms*, 1999.
- [24] M. S. Lam, E. E. Rothberg, and M. E. Wolf. The cache performance and optimizations of blocked algorithms. In *Proceedings of the Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 63–74, Apr. 1991.
- [25] A. LaMarca and R. Ladner. The influence of cache on the performance of sorting. In *Proceedings of the ACM-SIAM Symposium of Discrete Algorithms*, pages 370–379, 1997.
- [26] S. A. Przybylski. *Cache and Memory Hierarchy Design: A Performance-Directed Approach*. Morgan Kaufmann Publishers, San Mateo, CA, 1990.
- [27] P. Sanders. Accessing multiple sequences through set associative caches. In *Proceedings of ICALP*, 1999. A more recent version by Mehlhorn and Sanders was communicated to the authors in Dec 1999.
- [28] J. Savage. Extending the Hong-Kung model to memory hierarchies. In *Proceedings of COCOON*, volume LNCS 959, pages 270–281. Springer Verlag, 1995.
- [29] S. Sen and S. Chatterjee. Towards a theory of cache-efficient algorithms. In *Proceedings of the Symposium on Discrete Algorithms*, 2000.
- [30] D. Sleator and R. Tarjan. Amortized efficiency of list update and paging rules. *Commun. ACM*, 28(2):202–208, 1985.
- [31] M. Thottethodi, S. Chatterjee, and A. R. Lebeck. Tuning Strassen’s matrix multiplication for memory efficiency. In *Proceedings of SC98 (CD-ROM)*, Orlando, FL, Nov. 1998. Available from <http://www.supercomp.org/sc98>.
- [32] J. Vitter and M. Nodine. Large scale sorting in uniform memory hierarchies. *Journal of Parallel and Distributed Computing*, 17:107–114, 1993.
- [33] J. Vitter and E. Shriver. Algorithms for parallel memory I: Two-level memories. *Algorithmica*, 12(2):110–147, 1994.

A Approximating probability of conflict

Let μ be the number of elements between $S_{i,j}$ and $S_{i,j+1}$, *i.e.*, one less than the difference in ranks of $S_{i,j}$ and $S_{i,j+1}$. (μ may be 0, which guarantees event E1.) Let E_m denote the event that $\mu = m$. Then $\Pr[E1] = \sum_m \Pr[E1 \cap E_m]$, since E_m ’s are disjoint. For each m , $\Pr[E1 \cap E_m] = \Pr[E1|E_m] \cdot \Pr[E_m]$. The events E_m correspond to a geometric distribution, *i.e.*,

$$\Pr[E_m] = \Pr[\mu = m] = \frac{1}{k} \left(1 - \frac{1}{k}\right)^m. \quad (2)$$

To compute $\Pr[E1|E_m]$, we further subdivide the event into cases about how the m numbers are distributed into the sets $S_j, j \neq i$. Wlog, let $i = 1$ to keep notations simple. Let m_2, \dots, m_k denote the case that m_j numbers belong to sequence S_j ($\sum_j m_j = m$). We need to estimate the probability that for sequence S_j, b_j does not conflict with $\mathbb{S}(b_1)$ (recall that we have fixed $i = 1$) during the course that m_j elements arrive in S_j . This can happen only if $\mathbb{S}(b_j)$ (the cache set position of the leading block of S_j right after element $S_{1,t}$) does not lie roughly $\lceil m_i/B \rceil$ blocks from $\mathbb{S}(b_1)$. From assumption A1 and some careful counting this is $1 - \frac{m_i - 1 + B}{sB}$ for $m_j \geq 1$. For $m_j = 0$, this probability is 1 since no elements go into S_j and hence there is no conflict.⁴ These events are independent from our assumption A1 and hence these can be multiplied. The probability for a fixed partition m_2, \dots, m_k is the multinomial $\frac{m!}{m_2! \dots m_k!} \cdot \left(\frac{1}{k-1}\right)^m$ (m is partitioned into $k - 1$ parts). Therefore we can write the following expression for $\Pr[E1|E_m]$.

$$\Pr[E1|E_m] = \sum_{m_2 + \dots + m_k = m} \frac{m!}{m_2! \dots m_k!} \cdot \left(\frac{1}{k-1}\right)^m \prod_{m_i \neq 0} \left(1 - \frac{m_j - 1 + B}{sB}\right) \quad (3)$$

⁴The reader will soon realize that this case leads to some non-trivial calculations.

In the remainder of this section, we will obtain an upper bound on the right hand side of equation (3). Let $nz(m_2, \dots, m_k)$ denote the number of j s for which $m_j \neq 0$ (non-zero partitions). Then equation (3) can be rewritten as the following inequality.

$$\Pr[E1|E_m] \leq \sum_{m_2+\dots+m_k=m} \frac{m!}{m_2! \dots m_k!} \cdot \left(\frac{1}{k-1}\right)^m \left(1 - \frac{1}{s}\right)^{nz(m_2 \dots m_k)} \quad (4)$$

since $\left(1 - \frac{m_j-1+B}{sB}\right) \leq \left(1 - \frac{1}{s}\right)$ for $m_j \geq 1$. In other words, the right side is the expected value of $\left(1 - \frac{1}{s}\right)^{NZ(m, k-1)}$, where $NZ(m, k-1)$ denotes the number of non-empty bins when m balls are thrown into $k-1$ bins. Using equation (2) and the preceding discussion, we can write down an upper bound for the (unconditional) probability of $E1$ as

$$\sum_{m=0}^{\infty} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m \cdot E \left[\left(1 - \frac{1}{s}\right)^{NZ(m, k-1)} \right] \quad (5)$$

We use known sharp concentration bounds for the occupancy problem to obtain the following approximation for the expression (5) in terms of s and k .

Theorem A.1 ([22]) *Let $r = m/n$, and Y be the number of empty bins when m balls are thrown randomly into n bins. Then*

$$E[Y] = n \left(1 - \frac{1}{m}\right)^m \sim ne^{-r}$$

and for $\lambda > 0$

$$\Pr[|Y - E[Y]| \geq \lambda] \leq 2 \exp\left(-\frac{\lambda^2(n-1)/2}{n^2 - \mu^2}\right).$$

□

Corollary A.2 *Let NZ be the number of non-empty bins when m balls are thrown into k bins. Then*

$$E[NZ] = k(1 - e^{-m/k})$$

and

$$\Pr[|NZ - E[NZ]| \geq \alpha \sqrt{2k \log k}] \leq 1/k^\alpha.$$

□

So in equation (4), $E\left[\left(1 - \frac{1}{s}\right)^{NZ(m, k-1)}\right]$ can be bounded by

$$1/k^\alpha \left(1 - \frac{1}{s}\right) + \left(1 - \frac{1}{s}\right)^{k(1 - e^{-m/k} - \alpha \sqrt{2k \log k}/k)} \quad (6)$$

for any α and $m \geq 1$.

Proof: (of Lemma 4.1): We will split up the summation of (5) into two parts, namely, $m \leq e/2 \cdot k$ and $m > e/2 \cdot k$. One can obtain better approximations by refining the partitions, but our objective here is to demonstrate the existence of ϵ and δ and not necessarily obtain the best values.

$$\begin{aligned} \sum_{m=0}^{\infty} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m \cdot E\left[\left(1 - \frac{1}{s}\right)^{NZ(m,k-1)}\right] &= \sum_{m=0}^{ek/2k} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m \cdot E\left[\left(1 - \frac{1}{s}\right)^{NZ(m,k-1)}\right] \\ &+ \sum_{m=ek/2+1}^{\infty} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m \cdot E\left[\left(1 - \frac{1}{s}\right)^{NZ(m,k-1)}\right] \quad (7) \end{aligned}$$

The first term can be upper bounded by

$$\sum_{m=0}^{ek/2} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m$$

which is $\sim 1 - \frac{1}{e^{e/2}} \sim 0.74$.

The second term can be bounded using equation (6) using $\alpha \geq 2$.

$$\begin{aligned} \sum_{m=ek/2+1}^{\infty} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m \cdot E\left[\left(1 - \frac{1}{s}\right)^{NZ(m,k-1)}\right] &\leq \sum_{m=ek/2+1}^{\infty} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m \cdot 1/k^2 (1 - 1/s) \\ &+ \sum_{m=ek/2+1}^{\infty} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m \cdot \left(1 - \frac{1}{s}\right)^{k(1-e^{-m/k} - \alpha\sqrt{2k \log k/k})} \quad (8) \end{aligned}$$

The first term of the previous equation is less than $1/k$ and the second term can be bounded by

$$\sum_{m=ek/2+1}^{\infty} \frac{1}{k} \left(1 - \frac{1}{k}\right)^m \cdot \left(1 - \frac{1}{s}\right)^{0.25k}$$

for sufficiently large k ($k > 80$ suffices). This can be bounded by $\sim 0.25e^{-0.25k/s}$, so equation (8) can be bounded by $1/k + 0.25e^{-0.25k/s}$. Adding this to the first term of equation (7), we obtain an upper bound of $0.75 + 0.25e^{-0.25k/s}$ for $k > 100$. Subtracting this from 1 gives us $\frac{1-e^{-0.25k/s}}{4}$, i.e., $\delta \geq \frac{1-e^{-0.25k/s}}{4}$. \square