

Comparing Memory Systems for Chip Multiprocessors

Jacob Leverich
Amin Firoozshahian

Hideho Arakida
Mark Horowitz

Alex Solomatnikov
Christos Kozyrakis

Computer Systems Laboratory
Stanford University

{leverich, arakida, sols, aminf13, horowitz, kozyraki}@stanford.edu

ABSTRACT

There are two basic models for the on-chip memory in CMP systems: *hardware-managed coherent caches* and *software-managed streaming memory*. This paper performs a direct comparison of the two models under the same set of assumptions about technology, area, and computational capabilities. The goal is to quantify how and when they differ in terms of performance, energy consumption, bandwidth requirements, and latency tolerance for general-purpose CMPs. We demonstrate that for data-parallel applications, the cache-based and streaming models perform and scale equally well. For certain applications with little data reuse, streaming scales better due to better bandwidth use and macroscopic software prefetching. However, the introduction of techniques such as hardware prefetching and non-allocating stores to the cache-based model eliminates the streaming advantage. Overall, our results indicate that there is not sufficient advantage in building streaming memory systems where all on-chip memory structures are explicitly managed. On the other hand, we show that streaming at the programming model level is particularly beneficial, even with the cache-based model, as it enhances locality and creates opportunities for bandwidth optimizations. Moreover, we observe that stream programming is actually easier with the cache-based model because the hardware guarantees correct, best-effort execution even when the programmer cannot fully regularize an application's code.

Categories and Subject Descriptors: B.3.2 [Memory Structures]: Design Styles; C.1.2 [Processor Architectures]: Multiple Data Stream Architectures (Multiprocessors); D.1.3 [Programming Techniques]: Concurrent Programming

General Terms: Performance, Design

Keywords: Chip multiprocessors, coherent caches, streaming memory, parallel programming, locality optimizations

1. INTRODUCTION

The scaling limitations of uniprocessors [2] have led to an industry-wide turn towards chip multiprocessor (CMP) systems. CMPs are becoming ubiquitous in all computing domains. Unlike uniprocessors, which have a dominant, well-understood model

for on-chip memory structures, there is no widespread agreement on the memory model for CMP designs. The choice of memory model can significantly impact efficiency in terms of performance, energy consumption, and scalability. Moreover, it is closely coupled with the choice of parallel programming model, which in turn affects ease of use. While it is possible to map any programming model to any memory model, it is typically more efficient if the programming model builds upon the basic principles of the memory model.

Similar to larger parallel systems [8], there are two basic memory models for contemporary CMP systems: *hardware-managed, implicitly-addressed, coherent caches* and *software-managed, explicitly-addressed, local memories* (also called *streaming memory*). With the cache-based model, all on-chip storage is used for private and shared caches that are kept coherent by hardware. The advantage is that the hardware provides best-effort locality and communication management, even when the access and sharing patterns are difficult to statically analyze. With the streaming memory model, part of the on-chip storage is organized as independently addressable structures. Explicit accesses and DMA transfers are needed to move data to and from off-chip memory or between two on-chip structures. The advantage of streaming memory is that it provides software with full flexibility on locality and communication management in terms of addressing, granularity, and replacement policy. Since communication is explicit, it can also be proactive, unlike the mostly reactive behavior of the cache-based model. Hence, streaming allows software to exploit producer-consumer locality, avoid redundant write-backs for temporary results, manage selective data replication, and perform application-specific caching and macroscopic prefetching. Streaming eliminates the communication overhead and hardware complexity of the coordination protocol needed for cache coherence. On the other hand, it introduces software complexity, since either the programmer or the compiler must explicitly manage locality and communication.

Traditional desktop and enterprise applications are difficult to analyze and favor cache-based systems. In contrast, many upcoming applications from the multimedia, graphics, physical simulation, and scientific computing domains are being targeted by both cache-based [4, 45] and streaming [3, 39, 9, 16, 32] systems. The debate is particularly interesting vis-à-vis the latest game consoles. The CMPs for the Xbox360 and PlayStation 3 differ dramatically in their on-chip memory model, as Xbox360's Xenon processor is a cache-based CMP and PlayStation 3's Cell processor is a streaming memory CMP. Hence, it is interesting to evaluate if streaming provides specific benefits to this domain, or if using a cache-based approach across all application domains should be preferable.

The goal of this paper is to compare the efficiency of the two memory models under the same set of assumptions about tech-

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

ISCA'07, June 9–13, 2007, San Diego, California, USA.

Copyright 2007 ACM 978-1-59593-706-3/07/0006 ...\$5.00.

nology, area, and computational capabilities. Specifically, we are interested in answering the following questions: How do the two models compare in terms of overall *performance* and *energy consumption*? How does the comparison change as we *scale* the number or compute throughput of the processor cores? How sensitive is each model to *bandwidth* or *latency variations*? We believe that such a direct comparison will provide valuable information for the CMP architecture debate and generate some guidelines for the development of future systems.

The major conclusions from our comparison are:

- For data-parallel applications with abundant data reuse, the two models perform and scale equally well. Caches are as effective as software-managed memories at capturing locality and reducing average memory access time. For some of these applications, streaming has an energy advantage of 10% to 25% over write-allocate caches because it avoids superfluous refills on output data streams. Using a no-write-allocate policy for output data in the cache-based system reduces the streaming advantage.
- For applications without significant data reuse, macroscopic prefetching (double-buffering) provides streaming memory systems with a performance advantage when we scale the number and computational capabilities of the cores. The use of hardware prefetching with the cache-based model eliminates the streaming advantage for some latency-bound applications. There are also cases where streaming performs worse, such as when it requires redundant copying of data or extra computation in order to manage local stores.
- Our results indicate that a pure streaming memory model is not sufficiently advantageous at the memory system level. With the addition of prefetching and non-allocating writes, the cache-based model provides similar performance, energy, and bandwidth behavior. On the other hand, we found that “streaming” at the programming model level is very important, even with the cache-based model. Properly blocking an application’s working set, exposing producer-consumer locality, and identifying output-only data leads to significant efficiency gains. Moreover, stream programming leads to code that requires coarser-grain and lower-frequency coherence and consistency in a cache-based system. This observation will be increasingly relevant as CMPs scale to much larger numbers of cores.
- Finally, we observe that stream programming is actually easier with the cache-based model because the hardware guarantees correct, best-effort execution, even when the programmer cannot fully regularize an application’s code. With the streaming memory model, the software must orchestrate locality and communication perfectly, even for irregular codes.

The rest of the paper is organized as follows. Section 2 summarizes the two memory models and discusses their advantages and drawbacks. Section 3 presents the architectural framework for our comparison and Section 4 describes the experimental methodology. Section 5 analyzes our evaluation results. Section 6 focuses on streaming at the programming level and its interactions with CMP architecture. Section 7 addresses limitations in our methodology. Finally, Section 8 reviews related work and Section 9 concludes the paper.

		Communication	
		HW	SW
Locality	HW	Coherent Cache-based	Incoherent Cache-based
	SW	(impractical)	Streaming Memory

Table 1. The design space for on-chip memory for CMPs. This work focuses on the two highlighted options: coherent cache-based and streaming memory. The third practical option, incoherent cache-based, is briefly discussed in Section 7.

2. ON-CHIP MEMORY MODELS FOR CHIP MULTIPROCESSORS

General-purpose systems are designed around a *computational model* and a *memory model*. The computational model defines the organization of execution resources and register files and may follow some combination of the the superscalar, VLIW, vector, or SIMD approaches. The memory model defines the organization of on-chip and off-chip memories, as well as the communication mechanisms between the various computation and memory units. The two models are linked, and an efficient system is carefully designed along both dimensions. However, we believe that the on-chip memory model creates more interesting challenges for CMPs. First, it is the one that changes most dramatically compared to uniprocessor designs. Second, the memory model has broader implications on both software and hardware. Assuming we can move the data close to the execution units in an efficient manner, it is not difficult to select the proper computational model based on the type(s) of parallelism available in the computation kernels.

For both the cache-based and streaming models, certain aspects of the on-chip memory system are set by VLSI constraints such as wire delay [18]. Every node in the CMP system is directly associated with a limited amount of storage (first-level memory) that can be accessed within a small number of cycles. Nodes communicate by exchanging packets over an on-chip network that can range from a simple bus to a hierarchical structure. Additional memory structures (second-level memory) are also connected to the network fabric. The system scales with technology by increasing the number of nodes, the size of the network, and the capacity of second-level storage. Eventually, the scale of a CMP design may be limited by off-chip bandwidth or energy consumption [20].

Although they are under the same VLSI constraints, the cache-based and streaming models differ significantly in the way they manage data locality and inter-processor communication. As shown in Table 1, the cache-based model relies on hardware mechanisms for both, while the streaming model delegates management to software. The rest of this section overviews the protocol and operation for each memory model for CMP systems.

2.1 Coherent Cache-based Model

With the cache-based model [8], the only directly addressable storage is the off-chip memory. All on-chip storage is used for caches with hardware management of locality and communication. As cores perform memory accesses, the caches attempt to capture the application’s working set by fetching or replacing data at the granularity of cache blocks. The cores communicate implicitly through loads and stores to the single memory image. Since many caches may store copies of a specific address, it is necessary to query multiple caches on load and store requests and potentially invalidate entries to keep caches coherent. A coherence protocol, such as MESI, minimizes the cases when remote cache

lookups are necessary. Remote lookups can be distributed through a broadcast mechanism or by first consulting a directory structure. The cores synchronize using atomic operations such as compare-and-swap. Cache-based systems must also provide event ordering guarantees within and across cores following some consistency model [1]. Caching, coherence, synchronization, and consistency are implemented in hardware.

Coherent caching techniques were developed for board-level and cabinet-level systems (SMPs and DSMs), for which communication latency ranges from tens to hundreds of cycles. In CMPs, coherence signals travel within one chip, where latency is much lower and bandwidth is much higher. Consequently, even algorithms with non-trivial amounts of communication and synchronization can scale reasonably well. Moreover, the efficient design points for coherent caching in CMPs are likely to be different from those for SMP and DSM systems.

2.2 Streaming Memory Model

With streaming, the local memory for data in each core is a separately addressable structure, called a scratch-pad, local store, or stream register file. We adopt the term *local store* in this work. Software is responsible for managing locality and communication across local stores. Software has full flexibility in placing frequently accessed data in local stores with respect to location, granularity, replacement policy, allocation policy, and even the number of copies. For applications with statically analyzable access patterns, software can exploit this flexibility to minimize communication and overlap it with useful computation in the best possible application-specific way. Data are communicated between local stores or to and from off-chip memory using explicit DMA transfers. The cores can access their local stores as FIFO queues or as randomly indexed structures [23]. The streaming model requires DMA engines, but no other special hardware support for coherence or consistency.

Streaming is essentially message-passing applied to CMP designs, though there are some important differences compared to conventional message-passing for clusters and massively parallel systems [8]. First, communication is always orchestrated at the user-level, and its overhead is low. Next, messages are exchanged at the first level of the memory hierarchy, not the last one. Since the communicating cores are on the same chip, communication latencies are low and bandwidth is high. Finally, software manages both the communication between cores and the communication between a core and off-chip memory.

The discussion above separates the two memory models from the selection of a computation model. The streaming memory model is general and has already been used with VLIW/SIMD systems [3], RISC cores [39], vector processors [30], and even DSPs [32]. The cache-based model can be used with any of these computation models as well.

2.3 Qualitative Comparison

When considering the memory system alone, we find several ways in which cache-based and streaming memory systems may differ: bandwidth utilization, latency tolerance, performance, energy consumption, and cost. This section summarizes each of these differences in a qualitative manner.

Off-chip Bandwidth & Local Memory Utilization: The cache-based model leads to bandwidth and storage capacity waste on sparsely strided memory accesses. In the absence of spatial locality, manipulating data at the granularity of wide cache lines is wasteful. Streaming memory systems, by virtue of strided scatter and gather DMA transfers, can use the minimum memory channel bandwidth

necessary to deliver data, and also compact the data within the local store. Note, however, that memory and interconnect channels are typically optimized for block transfers and may not be bandwidth efficient for strided or scatter/gather accesses.

Caching can also waste off-chip bandwidth on unnecessary refills for output data. Because caches often use write-allocate policies, store misses force memory reads before the data are overwritten in the cache. If an application has disjoint input and output streams, the refills may waste a significant percentage of bandwidth. Similarly, caching can waste bandwidth on write-backs of dead temporary data. A streaming system does not suffer from these problems, as the output and temporary buffers are managed explicitly. Output data are sent off-chip without refills, and dead temporary data can be ignored, as they are not mapped to off-chip addresses. To mitigate the refill problem, cache-based systems can use a no-write-allocate policy. In this case, it is necessary to group store data in write buffers before forwarding them to memory in order to avoid wasting bandwidth on narrow writes [4]. Another approach is to use cache control instructions, such as “Prepare For Store,” [34] that instruct the cache to allocate a cache line but avoid retrieval of the old values from memory. Similarly, temporary data can be marked invalid at the end of a computation [7, 43]. In any case, software must determine when to use these mechanisms.

Streaming systems may also waste bandwidth and storage capacity on programs with statically unpredictable, irregular data patterns. A streaming system can sometimes cope with these patterns by fetching a superset of the needed input data. Alternatively, at the cost of enduring long latencies, the system could use a DMA transfer to collect required data on demand from main memory before each computational task. For programs that operate on overlapping blocks or graph structures with multiple references to the same data, the streaming system may naively re-fetch data. This can be avoided through increased address generation complexity or software caching. Finally, for applications that fetch a block and update some of its elements in-place, a streaming system will often write back the whole block to memory at the end of the computation, even if some data were not updated. In contrast to all of these scenarios, cache-based systems perform load and store accesses on demand, and hence only move cache lines as required. They may even search for copies of the required data in other on-chip caches before going off-chip.

Latency Tolerance: The cache-based model is traditionally reactive, meaning that a miss must occur before a fetch is triggered. Memory latency can be hidden using hardware prefetching techniques, which detect repetitive access patterns and issue memory accesses ahead of time, or proactive software prefetching. In practice, the DMA transfers in the streaming memory model are an efficient and accurate form of software prefetching. They can hide a significant amount of latency, especially if double-buffering is used. Unlike hardware prefetching, which requires a few misses before a pattern is detected (microscopic view), a DMA access can start arbitrarily early (macroscopic view) and can capture both regular and irregular (scatter/gather) accesses.

Performance: From the discussion so far, one can conclude that streaming memory may have a performance or scaling advantage for regular applications, due to potentially better latency tolerance and better usage of off-chip bandwidth or local storage. These advantages are important only if latency, bandwidth, or local storage capacity are significant bottlenecks to begin with. For example, reducing the number of misses is unimportant for a computationally intensive application that already has very good locality. In the event that an application is bandwidth-bound, latency tolerance measures will be ineffective. A drawback for streaming, even with

regular code, is that it often has to execute additional instructions to set up DMA transfers. For applications with unpredictable data access patterns or control flow, a streaming system may execute a significantly higher number of instructions than that of a cache-based system in order to produce predictable patterns or to use the local store to emulate a software cache.

Energy Consumption: Any performance advantage also translates to an energy advantage, as it allows us to turn off the system early or scale down its power supply and clock frequency. Streaming accesses to the first-level storage eliminate the energy overhead of caches (tag access and tag comparison). The cache-based model consumes additional energy for on-chip coherence traffic, snoop requests or directory lookups. Moreover, efficient use of the available off-chip bandwidth by either of the two models (through fewer transfers or messages) reduces the energy consumption by the interconnect network and main memory.

Complexity & Cost: It is difficult to make accurate estimates of hardware cost without comparable implementations. The hardware for the cache-based model is generally more complex to design and verify, as coherence, synchronization, consistency, and prefetching interact in subtle ways. Still, reuse across server, desktop, and embedded CMP designs can significantly reduce such costs. On the other hand, streaming passes the complexity to software, the compiler, and/or the programmer. For applications in the synchronous data-flow, DSP, or dense matrix domains, it is often straight forward to express a streaming algorithm. For other applications, it is non-trivial, and a single good algorithm is often a research contribution in itself [11, 14]. Finally, complexity and cost must also be considered with scaling in mind. A memory model has an advantage if it allows efficient use of more cores without the need for disproportional increases in bandwidth or some other resource in the memory system.

3. CMP ARCHITECTURE FRAMEWORK

We compare the two models using the CMP architecture shown in Figure 1. There are numerous design parameters in a CMP system, and evaluating the two models under all possible combinations is infeasible. Hence, we start with a baseline system that represents a practical design point and vary only the key parameters that interact significantly with the on-chip memory system.

3.1 Baseline Architecture

Our CMP design is based on in-order processors similar to Piranha [5], RAW [39], Ultrasparc T1 [25], and Xbox360 [4]. Such CMPs have been shown to be efficient for multimedia, communications, and throughput computing workloads as they provide high compute density without the area and power overhead of out-of-order processors [10]. We use the Tensilica Xtensa LX, 3-way VLIW core [22]. Tensilica cores have been used in several embedded CMP designs, including the 188-core Cisco Metro router chip [12]. We also performed experiments with a single-issue Xtensa core that produced similar results for the memory model comparison. The VLIW core is consistently faster by 1.6x to 2x for the applications we studied. The core has 3 slots per instruction, with up to 2 slots for floating-point operations and up to 1 slot for loads and stores. Due to time constraints, we do not use the Tensilica SIMD extensions at this point. Nevertheless, Section 5.3 includes an experiment that evaluates the efficiency of each model with processors that provide higher computational throughput. Each processor has a 16-KByte, 2-way set-associative instruction cache. The fixed amount of local data storage is used differently in each memory model.

We explore systems with 1 to 16 cores using a hierarchical interconnect similar to that suggested by [26]. We group cores in clusters of four with a wide, bidirectional bus (*local network*) providing the necessary interconnect. The cluster structure allows for fast communication between neighboring cores. If threads are mapped intelligently, the intra-cluster bus will handle most of the core-to-core communication. A global crossbar connects the clusters to the second-level storage. There is buffering at all interfaces to tolerate arbitration latencies and ensure efficient bandwidth use at each level. The hierarchical interconnect provides sufficient bandwidth, while avoiding the bottlenecks of long wires and centralized arbitration [26]. Finally, the secondary storage communicates to off-chip memory through some number of memory channels.

Table 2 presents the parameters of the CMP system. We vary the number of cores, the core clock frequency, the available off-chip bandwidth, and the degree of hardware prefetching. We keep the capacity of the first-level data-storage constant.

3.2 Cache-based Implementation

For the cache-based model, the first-level data storage in each core is organized as a 32-KByte, 2-way set-associative data cache. The second-level cache is a 512-KByte, 16-way set-associative cache. Both caches use a write-back, write-allocate policy. Coherence is maintained using the MESI write-invalidate protocol. Coherence requests are propagated in steps over the hierarchical interconnect. First, they are broadcast to other processors within the cluster. If the cache-line is not available or the request cannot be satisfied within one cluster (e.g., upgrade request), it is broadcast to all other clusters as well. Snooping requests from other cores occupy the data cache for one cycle, forcing the core to stall if it tries to do a load/store access in the same cycle. Each core includes a store-buffer that allows loads to bypass store misses. As a result, the consistency model is weak. Since the processors are in-order, it is easy to provide sufficient MSHRs for the maximum possible number of concurrent misses.

Each core additionally includes a hardware stream-based prefetch engine that places data directly in the L1 cache. Modeled after the tagged prefetcher described in [41], the prefetcher keeps a history of the last 8 cache misses for identifying sequential accesses, runs a configurable number of cache lines ahead of the latest cache miss, and tracks 4 separate access streams. Our experiments include hardware prefetching only when explicitly stated.

We use POSIX threads to manually parallelize and tune applications [27]. The applications used are regular and use locks to implement efficient task-queues and barriers to synchronize SPMD code. Higher-level programming models, such as OpenMP, are also applicable to these applications.

3.3 Streaming Implementation

For the streaming model, the first-level data storage in each core is split between a 24-KByte local store and an 8-KByte cache. The small cache is used for stack data and global variables. It is particularly useful for the sequential portions of the code and helps simplify the programming and compilation of the streaming portion as well. The 24-KByte local store is indexed as a random access memory. Our implementation also provides hardware support for FIFO accesses to the local store, but we did not use this feature with any of our applications. Each core has a DMA engine that supports sequential, strided, and indexed transfers, as well as command queuing. At any point, the DMA engine may have up to 16 32-byte outstanding accesses.

The local store is effectively smaller than a 24-KByte cache, since it has no tag or control bits overhead. We do not raise the

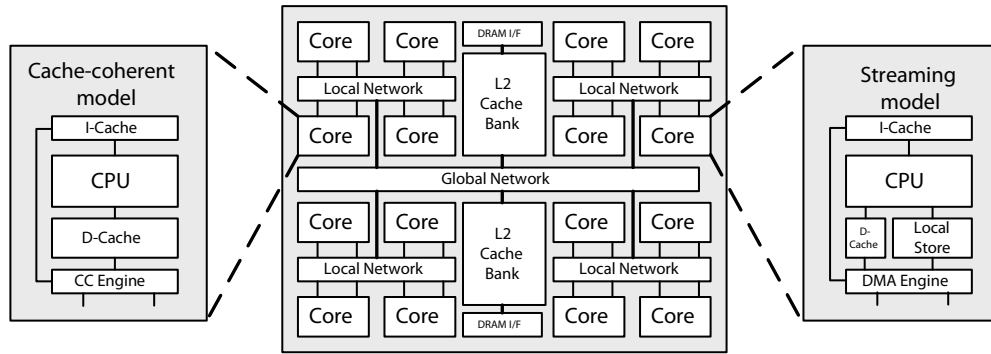


Figure 1. The architecture of the CMP system with up to 16 processors. The core organizations for the cache-based and streaming models are shown on the left and right side respectively.

size of the local store, since the small increment (2 KBytes) does not make a difference for our applications. Still, the energy consumption model accounts correctly for the reduced capacity. The secondary storage is again organized as a 16-way set-associative cache. L2 caches are useful with stream processors, as they capture long-term reuse patterns and avoid expensive accesses to main memory [36, 9]. The L2 cache avoids refills on write misses when DMA transfers overwrite entire lines.

We developed streaming code using a simple library for DMA transfers within threaded C code. We manually applied the proper blocking factor and double-buffering in order to overlap DMA transfers with useful computation. We also run multiple computational kernels on each data block to benefit from producer-consumer locality without additional memory accesses or write-backs for intermediate results. Higher-level stream programming models should be applicable to most of our applications [15, 13]. In some cases, the DMA library uses a scheduling thread that queues pending transfers. We avoid taking up a whole core for this thread by multiplexing it with an application thread. The performance impact is negligible.

4. METHODOLOGY

4.1 Simulation and Energy Modeling

We used Tensilica’s modeling tools to construct a CMP simulator for both memory models [40]. The simulator captures all stall and contention events in the core pipeline and the memory system. Table 2 summarizes the major system characteristics and the parameters we varied for this study. The default values are shown in bold. The applications were compiled with Tensilica’s optimizing compiler at the -O3 optimization level. We fast-forward over the initialization for each application but simulate the rest to completion, excluding 179.art for which we measure 10 invocations of the `train_match` function.

We also developed an energy model for the architecture in a 90nm CMOS process (1.0V power supply). For the cores, the model combines usage statistics (instruction mix, functional unit utilization, stalls and idle cycles, etc.) with energy data from the layout of actual Tensilica designs at 600MHz in 90nm. The energy consumed by on-chip memory structures is calculated using CACTI 4.1 [38], which includes a leakage power model and improved circuit models compared to CACTI 3. Interconnect energy is calculated based on our measured activity statistics and scaled power measurements from [19]. The energy consumption for off-chip DRAM is derived from DRAMsim [42]. We model the effect of leakage and clock gating on energy at all levels of the model.

4.2 Applications

Table 3 presents the set of applications used for this study. They represent applications from the multimedia, graphics, physical simulation, DSP, and data management domains. Such applications have been used to evaluate and motivate the development of streaming architectures. MPEG-2, H.264, Raytracer, JPEG, and Stereo Depth Extraction are compute-intensive applications and show exceptionally good cache performance despite their large datasets. They exhibit good spatial or temporal locality and have enough computation per data element to amortize the penalty for any misses. FEM is a scientific application, but has about the same compute intensity as multimedia applications. The remaining applications—Bitonic Sort, Merge Sort, FIR, and 179.art—perform a relatively small computation on each input element. They require considerably higher off-chip bandwidth and are sensitive to memory latency.

We manually optimized both versions of each application to eliminate bottlenecks and schedule its parallelism in the best possible way. Whenever appropriate, we applied the same data-locality optimizations (i.e. blocking, producer-consumer, etc.) to both models. In Section 6, we explore the impact of data-locality optimizations. The following is a brief description of how each application was parallelized.

MPEG-2 and *H.264* are parallelized at the macroblock level. Both dynamically assign macroblocks to cores using a task queue. Macroblocks are entirely data-parallel in MPEG-2. In H.264, we schedule the processing of dependent macroblocks so as to minimize the length of the critical execution path. With the CIF resolution video frames we encode for this study, the macroblock parallelism available in H.264 is limited. *Stereo Depth Extraction* is parallelized by dividing input frames into 32x32 blocks and statically assigning them to processors.

KD-tree Raytracer is parallelized across camera rays. We assign rays to processors in chunks to improve locality. Our streaming version reads the KD-tree from the cache instead of streaming it with a DMA controller. *JPEG Encode* and *JPEG Decode* are parallelized across input images, in a manner similar to that done by an image thumbnail browser. Note that Encode reads a lot of data but outputs little; Decode behaves in the opposite way. The *Finite Element Method (FEM)* is parallelized across mesh cells. The *FIR filter* has 16 taps and is parallelized across long strips of samples. *179.art* is parallelized across F1 neurons; this application is composed of several data-parallel vector operations and reductions between which we place barriers.

Merge Sort and *Bitonic Sort* are parallelized across sub-arrays of a large input array. The processors first sort chunks of 4096 keys in

	Cache-Coherent Model (CC)	Streaming Model (STR)
Core	1, 2, 4, 8, or 16 Tensilica LX cores, 3-way VLIW, 7-stage pipeline 800MHz , 1.6GHz, 3.2GHz, or 6.4GHz clock frequency 2 FPU, 2 integer units, 1 load/store unit	
I-cache	16KB, 2-way associative, 32-byte blocks, 1 port	
Data Storage	32KB, 2-way associative cache 32-byte blocks, 1 port, MESI Hardware stream prefetcher	24KB local store, 1 port 8KB, 2-way associative cache, 32-byte blocks, 1 port DMA engine with 16 outstanding accesses
Local Network	bidirectional bus, 32 bytes wide, 2 cycle latency (after arbitration)	
Global Crossbar	1 input and output port per cluster or L2 bank, 16 bytes wide, 2.5ns latency (pipelined)	
L2-cache	512KB , 16-way set associative, 1 port, 2.2ns access latency, non-inclusive	
DRAM	One memory channel at 1.6GB/s, 3.2GB/s , 6.4GB/s, or 12.8GB/s; 70ns random access latency	

Table 2. Parameters for the CMP system. For parameters that vary, we denote the default value in bold. Latencies are for a 90nm CMOS process.

Application	Input Dataset	L1 D-Miss Rate	L2 D-Miss Rate	Instr. per L1 D-Miss	Cycles per L2 D-Miss	Off-chip B/W
MPEG-2 Encoder [28]	10 CIF frames (Foreman)	0.58%	85.3%	324.8	135.4	292.4 MB/s
H.264 Encoder	10 CIF frames (Foreman)	0.06%	30.8%	3705.5	4225.9	10.8 MB/s
KD-tree Raytracer	128x128, 16371 triangles	1.06%	98.9%	256.3	654.6	45.1 MB/s
JPEG Encoder [21]	128 PPMs of various sizes	0.40%	72.9%	577.1	84.2	402.2 MB/s
JPEG Decoder [21]	128 JPGs of various sizes	0.58%	76.2%	352.9	44.9	1059.2 MB/s
Stereo Depth Extraction	3 CIF image pairs	0.03%	46.1%	8662.5	3995.3	11.4 MB/s
2D FEM	5006 cell mesh, 7663 edges	0.60%	86.2%	368.8	55.5	587.9 MB/s
FIR filter	2 ²⁰ 32-bit samples	0.63%	99.8%	214.6	20.4	1839.1 MB/s
179.art	SPEC reference dataset	1.79%	7.4%	150.1	230.9	227.7 MB/s
Bitonic Sort	2 ¹⁹ 32-bit keys (2 MB)	2.22%	98.2%	140.9	26.1	1594.2 MB/s
Merge Sort	2 ¹⁹ 32-bit keys (2 MB)	3.98%	99.7%	71.1	33.7	1167.8 MB/s

Table 3. Memory characteristics of the applications measured on the cache-based model using 16 cores running at 800MHz.

parallel using quicksort. Then, sorted chunks are merged or sorted until the full array is sorted. Merge Sort gradually reduces in parallelism as it progress, whereas Bitonic Sort retains full parallelism for its duration. Merge Sort alternates writing output sublists to two buffer arrays, while Bitonic Sort operates on the list *in situ*.

5. EVALUATION

Our evaluation starts with a comparison of the streaming system to the baseline cache-based system without prefetching or other enhancements (Sections 5.1 and 5.2). We then study the bandwidth consumption and latency tolerance of the two systems (Section 5.3). We conclude by evaluating means to enhance the performance of caching systems (Sections 5.4 and 5.5).

5.1 Performance Comparison

Figure 2 presents the execution time for the coherent cache-based (CC) and streaming (STR) models as we vary the number of 800 MHz cores from 2 to 16. We normalize to the execution time of the sequential run with the cache-based system. The components of execution time are: useful execution (including fetch and non-memory pipeline stalls), synchronization (lock, barrier, wait for DMA), and stalls for data (caches). Lower bars are better. The cache-based results assume no hardware prefetching.

For 7 out of 11 applications (MPEG-2, H.264, Depth, Raytracing, FEM, JPEG Encode and Decode), the two models perform almost identically for all processor counts. These programs perform a significant computation on each data element fetched and can be classified as compute-bound. Both caches and local stores capture their locality patterns equally well.

The remaining applications (179.art, FIR, Merge Sort, and Bitonic Sort) are data-bound and reveal some interesting differences between the two models. The cache-based versions stall regularly

due to cache misses. Streaming versions eliminate many of these stalls using double-buffering (macroscopic prefetching). This is not the case for Bitonic Sort, because off-chip bandwidth is saturated at high processor counts. Bitonic Sort is an in-place sorting algorithm, and it is often the case that sublists are moderately in-order and elements don't need to be swapped, and consequently don't need to be written back. The cache-based system naturally discovers this behavior, while the streaming memory system writes the unmodified data back to main memory anyway. H.264 and Merge Sort have synchronization stalls with both models due to limited parallelism.

There are subtle differences in the useful cycles of some applications. FIR executes 14% more instructions in the streaming model than the caching model because of the management of DMA transfers (128 elements per transfer). In the streaming Merge Sort, the inner loop executes extra comparisons to check if an output buffer is full and needs to be drained to main memory, whereas the cache-based variant freely writes data sequentially. Even though double-buffering eliminates all data stalls, the application runs longer because of its higher instruction count. The streaming H.264 takes advantage of some boundary-condition optimizations that proved difficult in the cache-based variant. This resulted in a slight reduction in instruction count when streaming. MPEG-2 suffers a moderate number of instruction cache misses due the cache's limited size.

Overall, Figure 2 shows that neither model has a consistent performance advantage. Even without prefetching, the cache-based model performs similarly to the streaming one, and sometimes it is actually faster (Bitonic Sort for 16 cores). Both models use data efficiently, and the fundamental parallelism available in these applications is not affected by how data are moved. Although the differences in performance are small, there is a larger variation in

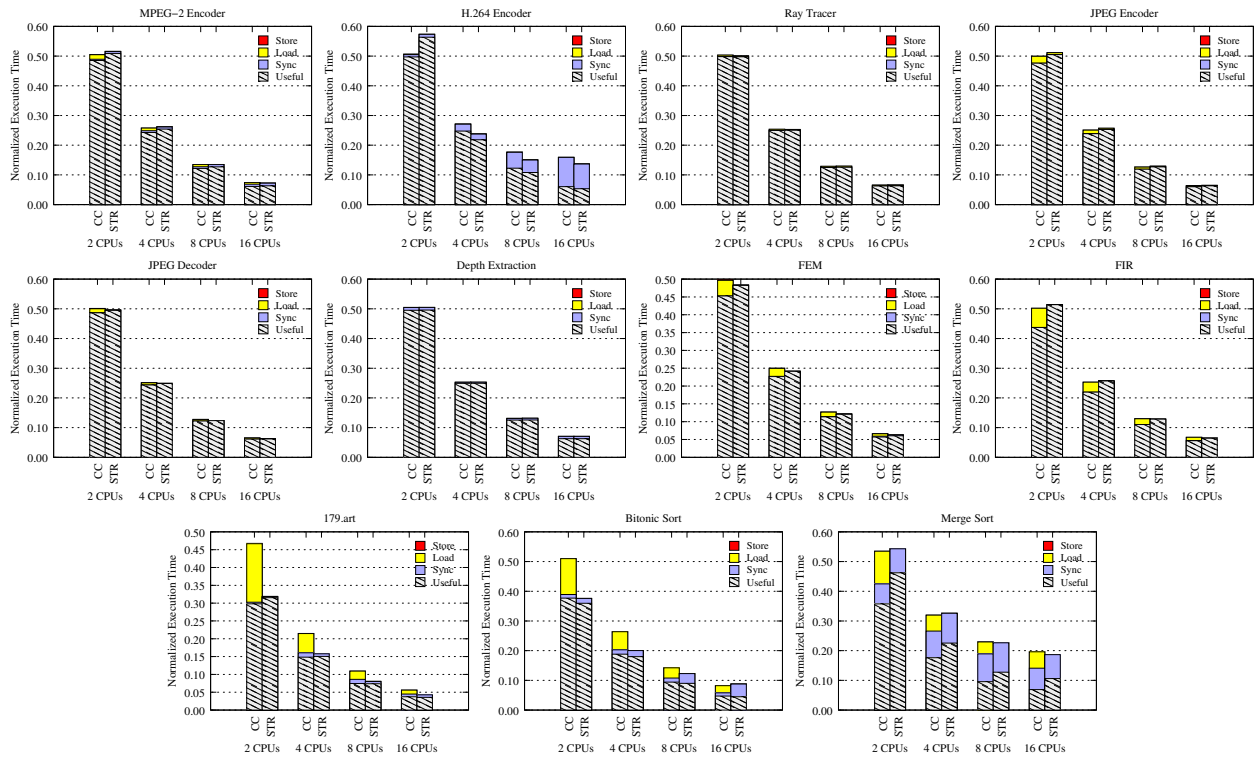


Figure 2. Execution times for the two memory models as the number of cores is increased, normalized to a single caching core.

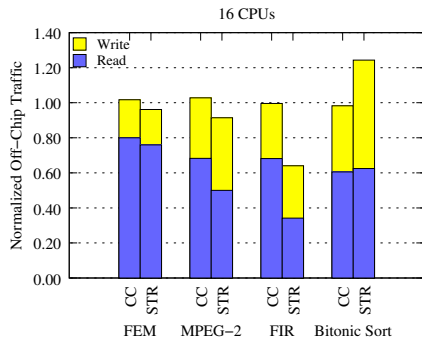


Figure 3. Off-chip traffic for the cache-based and streaming systems with 16 CPUs, normalized to a single caching core.

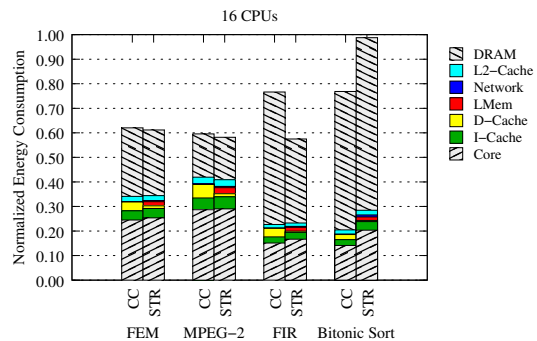


Figure 4. Energy consumption for the cache-based and streaming systems with 16 CPUs, normalized to a single caching core.

off-chip bandwidth utilization of the two models. Figure 3 shows that each model has an advantage in some situations.

5.2 Energy Comparison

Figure 4 presents energy consumption for the two models running FEM, MPEG-2, FIR, and Bitonic Sort. We normalize to the energy consumption of a single caching core for each application. Each bar indicates the energy consumed by the cores, the caches and local stores, the on-chip network, the second-level cache, and the main memory. The numbers include both static and dynamic power. Lower bars are better. In contrast to performance scaling, energy consumption does not always improve with more cores, since the amount of hardware used to run the application increases.

For 5 out of 11 applications (JPEG Encode, JPEG Decode, FIR, 179.art, and Merge Sort), streaming consistently consumes less energy than cache-coherence, typically 10% to 25%. The energy differential in nearly every case comes from the DRAM system. This

can be observed in the correlation between Figures 3 and 4. Specifically, the streaming applications typically transfer fewer bytes from main memory, often through the elimination of superfluous refills for output-only data. The opposite is true for our streaming Bitonic Sort, which tends to communicate more data with main memory than the caching version due to the write-back of unmodified data. For applications where there is little bandwidth difference between the two models (such as FEM) or the computational intensity is very high (such as Depth), the difference in energy consumption is insignificant.

We expected to see a greater difference between the local store and L1 data cache, but it never materialized. Since our applications are data-parallel and rarely share data, the energy cost of an average cache miss is dominated by the off-chip DRAM access rather than the modest tag broadcast and lookup. Hence, the per-access energy savings by eliminating tag lookups in the streaming system made little impact on the total energy footprint of the system.

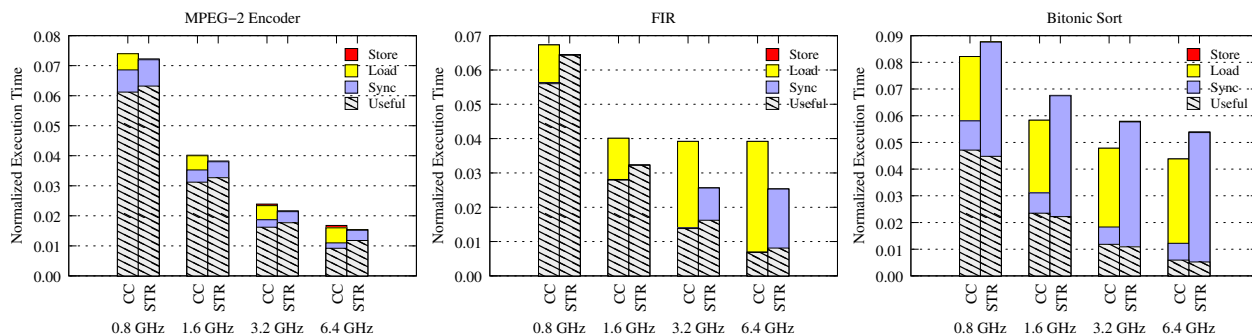


Figure 5. Normalized execution time as the computation rate of processor cores is increased (16 cores).

5.3 Increased Computational Throughput

Up to now, the results assume 800 MHz cores, which are reasonable for embedded CMPs for consumer applications. To explore the efficiency of the two memory models as the computational throughput of the processor is increased, we vary the clock frequency of the cores while keeping constant the bandwidth and latency in the on-chip networks, L2 cache, and off-chip memory. In some sense, the higher clock frequencies tell us in general what would happen with more powerful processors that use SIMD units, out-of-order schemes, higher clock frequency, or a combination. For example, the 6.4 GHz configuration can be representative of the performance of an 800 MHz processor that uses 4- to 8-wide SIMD instructions. The experiment was performed with 16 cores to stress scaling and increase the system’s sensitivity to *both* memory latency and memory bandwidth.

Applications with significant data reuse, such as H.264 and Depth, show no sensitivity to this experiment and perform equally well on both systems. Figure 5 shows the results for some of the applications that are affected by computational scaling. These applications fall into one of two categories: bandwidth-sensitive or latency-sensitive. Latency-sensitive programs, like MPEG-2 Encoding, perform a relative large degree of computation between off-chip memory accesses (hundreds of instructions). While the higher core frequency shortens these computations, it does not reduce the amount of time (in nanoseconds, not cycles) required to fetch the data in between computations. The macroscopic prefetching in the streaming system can tolerate a significant percentage of the memory latency. Hence at 6.4 GHz, the streaming MPEG-2 Encoder is 9% faster.

Bandwidth-sensitive applications, like FIR and Bitonic Sort, eventually saturate the available off-chip bandwidth. Beyond that point, further increases in computational throughput do not improve overall performance. For FIR, the cache-based system saturates before the streaming system due to the superfluous refills on store misses to output-only data. At the highest computational throughput, the streaming system performs 36% faster. For Bitonic Sort, the streaming version saturates first, since it performs more writes than the cache-based version (as described in 5.1). This gives the cache-based version a 19% performance advantage.

5.4 Mitigating Latency and Bandwidth Issues

The previous section indicates that when a large number of cores with high computational throughput are used, the cache-based model faces latency and bandwidth issues with certain applications. To characterize these inefficiencies, we performed experiments with increased off-chip bandwidth and hardware prefetching.

Figure 6 shows the impact of increasing the available off-chip bandwidth for FIR. This can be achieved by using higher frequency

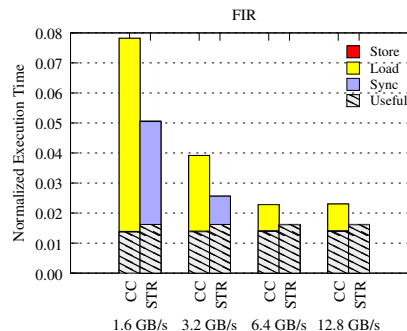


Figure 6. The effect of increased off-chip bandwidth on FIR performance. Measured on 16 cores at 3.2 GHz.

DRAM (e.g. DDR2, DDR3, GDDR) or multiple memory channels. With more bandwidth available, the effect of superfluous refills is significantly reduced, and the cache-based system performs nearly as well as the streaming one. When hardware prefetching is introduced at 12.8 GB/s, load stalls are reduced to 3% of the total execution time. However, the additional off-chip bandwidth does not close the energy gap for this application. An energy-efficient solution for the cache-based system is to use a non-allocating write policy, which we explore in Section 5.5.

For Merge Sort and 179.art (Figure 7), hardware prefetching significantly improves the latency tolerance of the cache-based systems; data stalls are virtually eliminated. This is not to say that we never observed data stalls—at 16 cores, the cache-based Merge Sort saturates the memory channel due to superfluous refills—but that a small degree of prefetching is sufficient to hide over 200 cycles of memory latency.

5.5 Mitigating Superfluous Refills

For some applications, the cache-based system uses more off-chip bandwidth (and consequently energy) because of superfluous refills for output-only data. This disadvantage can be addressed by using a non-write-allocate policy for output-only data streams. We mimic non-allocating stores by using an instruction similar to the MIPS32 “Prepare for Store” (PFS) instruction [34]. PFS allocates and validates a cache line without refilling it. The results for FIR, Mergesort, and MPEG-2 are shown in Figure 8. For each application, the elimination of superfluous refills brings the memory traffic and energy consumption of the cache-based model into parity with the streaming model. For MPEG-2, the memory traffic due to write misses was reduced 56% compared to the cache-based application without PFS.

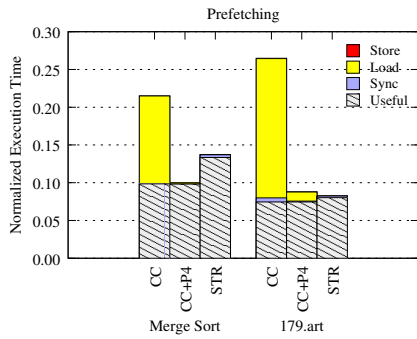


Figure 7. The effect of hardware prefetching on performance. P4 refers to the prefetch depth of 4. Measured on 2 cores at 3.2 GHz with a 12.8 GB/s memory channel.

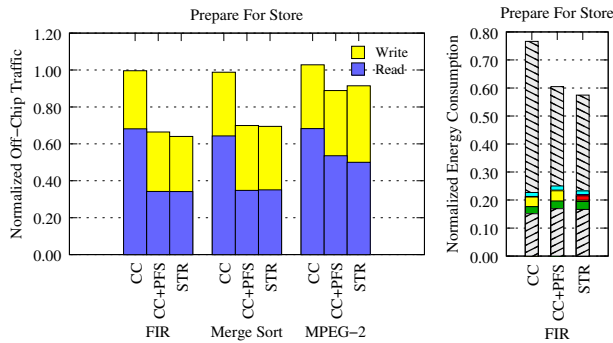


Figure 8. The effect of “Prepare for Store” (PFS) instructions on the off-chip traffic for the cache-based system, normalized to a single caching core. Also shown is energy consumption for FIR with 16 cores at 800 MHz. See Figure 4 for the second graph’s legend.

Note that a full hardware implementation of a non-write-allocate cache policy, along with the necessary write-gathering buffer, might perform better than PFS as it would also eliminate cache *replacements* due to output-only data.

6. STREAMING AS A PROGRAMMING MODEL

Our evaluation so far shows that the two memory models lead to similar performance and scaling. It is important to remember that we took advantage of streaming optimizations, such as blocking and locality-aware scheduling, on both memory models. To illustrate this, it is educational to look at stream programming and its implications for CMP architecture.

Stream programming models encourage programmers to think about the locality, data movement, and storage capacity issues in their applications [15, 13]. While they do not necessarily require the programmer to manage these issues, the programmer structures the application in such a way that it is easy to reason about them. This exposed nature of a stream program is vitally important for streaming architectures, as it enables software or compiler management of data locality and asynchronous communication with architecturally visible on-chip memories. Despite its affinity for stream architectures, we find that stream programming is beneficial for cache-based architectures as well.

Figure 9 shows the importance of streaming optimizations in the cache-based MPEG-2 Encoder. The original parallel code from [28] performs an application kernel on a whole video frame before the next kernel is invoked (i.e. Motion Estimation, DCT, Quantization,

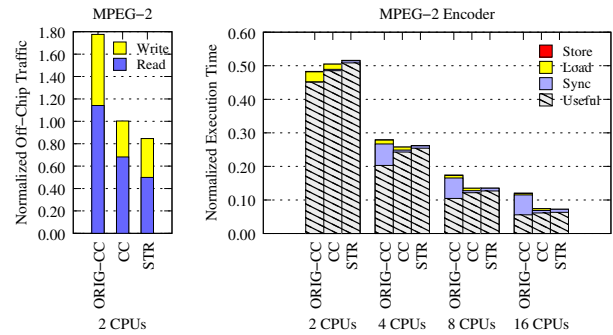


Figure 9. The effect of stream programming optimizations on the off-chip bandwidth and performance of MPEG-2 at 800 MHz.

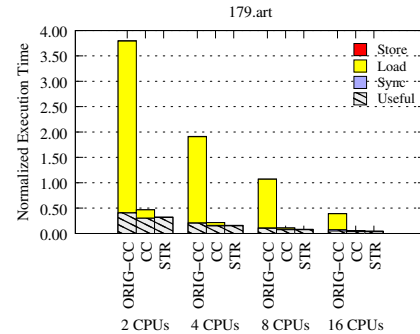


Figure 10. The effect of stream programming optimizations on the performance of 179.art at 800 MHz.

etc.). We restructured this code by hoisting the inner loops of several tasks into a single outer loop that calls each task in turn. In the optimized version, we execute all tasks on a block of a frame before moving to the next block. This also allowed us to condense a large temporary array into a small stack variable. The improved producer-consumer locality reduced write-backs from L1 caches by 60%. Data stalls are reduced by 41% at 6.4 GHz, even without prefetching. Furthermore, improving the parallel efficiency of the application became a simple matter of scheduling a single data-parallel loop, which alone is responsible for a 40% performance improvement at 16 cores. However, instruction cache misses are notably increased in the streaming-optimized code.

For 179.art, we reorganized the main data structure in the cache-based version in the same way as we did for the streaming code. We were also able to replace several large temporary vectors with scalar values by merging several loops. These optimizations reduced the sparseness of 179.art’s memory access pattern, improved both temporal and spatial locality, and allowed us to use prefetching effectively. As shown in Figure 10, the impact on performance is dramatic, even at small core counts (7x speedup).

Overall, we observed performance, bandwidth, and energy benefits whenever stream programming optimizations were applied to our cache-based applications. This is not a novel result, since it is well-known that locality optimizations, such as blocking and loop fusion [29], increase computational intensity and cache efficiency. However, stream programming models encourage users to write code that explicitly exposes an application’s parallelism and data-access pattern, more often allowing such optimizations.

Our experience is that that stream programming is actually *easier* with the cache-based model rather than the streaming model. With streaming memory, the programmer or the compiler must orchestrate all data movement and positioning exactly right in order for

the program to operate correctly and fast. This can be burdensome for irregular access patterns (overlapping blocks, search structures, unpredictable or data-dependent patterns, etc.), or for accesses that do not affect an application’s performance. It can lead to additional instructions and memory references that reduce or eliminate streaming hardware’s other advantages. With cache-based hardware, stream programming is just an issue of performance optimization. Even if the algorithm is not blocked exactly right, the caches will provide best-effort locality and communication management. Hence, the programmer or the compiler can focus on the most promising and most regular data structures instead of managing all data structures in a program.

Moreover, stream programming can address some of the coherence and consistency challenges when scaling cache-based CMPs to large numbers of cores. Since a streaming application typically operates in a data-parallel fashion on a sequence of data, there is little short-term communication or synchronization between processors. Communication is only necessary when processors move from one independent set of input/output blocks to the next or reach a cycle in an application’s data-flow graph. This observation may be increasingly important as CMPs grow, as it implies that less aggressive, coarser-grain, or lower-frequency mechanisms can be employed to keep caches coherent.

7. DISCUSSION AND LIMITATIONS

It is important to recognize that our study has limitations. Our experiments focus on CMPs with 1-16 cores and uniform memory access. Some of the conclusions may not generalize to larger-scale CMPs with non-uniform memory access (NUMA). A large-scale, cache-based CMP, programmed in a locality-oblivious way, will undoubtedly suffer stalls due to long memory delays or excessive on-chip coherence traffic. We observe that the stream programming model may be able to address both limitations; it exposes the flow of data early enough that they can be prefetched, and motivates a far coarser-grained, lower-frequency coherence model.

Towards the goal of designing large CMPs that are still easy to program, a hybrid memory system that combines caching and software-managed memory structures can mitigate efficiency challenges without exacerbating the difficulty of software development. For example, bulk transfer primitives for cache-based systems could enable more efficient macroscopic prefetching. Conversely, small, potentially incoherent caches in streaming memory systems could vastly simplify the use of static data structures with abundant temporal locality. An industry example of a hybrid system is the NVIDIA G80 GPU [6] which, in addition to cached access to global memory, includes a small scratch-pad for application-specific locality and communication optimizations.

Besides the limits of scalability, we did not consider architectures that expose the streaming model all the way to the register file [39] or applications without abundant data parallelism. We also did not consider changes to the pipeline of our cores, since that is precisely what makes it difficult to evaluate existing streaming memory processors compared to cache-based processors. Finally, our study was performed using general-purpose CMPs. A comparison between the two memory models for specialized CMPs remains an open issue. Despite these limitations, we believe this study’s conclusions are important in terms of understanding the actual behavior of CMP memory system and motivating future research and development.

8. RELATED WORK

Several architectures [3, 39, 16, 32] use streaming hardware with multimedia and scientific code in order to get performance and

energy benefits from software-managed memory hierarchies and regular control flow. There are also corresponding proposals for stream programming languages and compiler optimizations [15, 13]. Such tools can reduce the burden on the programmer for explicit locality and communication management. In parallel, there are significant efforts to enhance cache-based systems with traffic filters [35], replacement hints [43], or prefetching hints [44]. These enhancements target the same access patterns that streaming memory systems benefit from.

To the best of our knowledge, this is the first direct comparison of the two memory models for CMP systems under a unified set of assumptions. Jayasena [23] compared a stream register file to a single-level cache for a SIMD processor. He found that the stream register file provides performance and bandwidth advantages for applications with significant producer-consumer locality. Loghi and Poncino compared hardware cache coherence to not caching shared data at all for embedded CMPs with on-chip main memory [31]. The ALP report [28] evaluates multimedia codes on CMPs with streaming support. However, for all but one benchmark, streaming implied the use of enhanced SIMD instructions, not software managed memory hierarchies. Suh *et al.* [37] compared a streaming SIMD processor, a streaming CMP chip, a vector design, and a superscalar processor for DSP kernels. However, the four systems varied vastly at all levels, hence it is difficult to compare memory models directly. There are several proposals for configurable or hybrid memory systems [33, 36, 23, 28]. In such systems, a level in the memory hierarchy can be configured as a cache or as a local store depending on an application’s needs. Gummaraju and Rosenblum have shown benefits from a hybrid architecture that uses stream programming on a cache-based superscalar design for scientific code [17]. Our work supports this approach as we show that cache-based memory systems can be as efficient as streaming memory systems, but could benefit in terms of bandwidth consumption and latency tolerance from stream programming.

Our study follows the example of papers that compared shared memory and message passing for multi-chip systems [24, 8].

9. CONCLUSIONS

The choice of the on-chip memory model has far-reaching implications for CMP systems. In this paper, we performed a direct comparison of two basic models: coherent caches and streaming memory.

We conclude that for the majority of applications, both models perform and scale equally well. For some applications without significant data reuse, streaming has a performance and energy advantage when we scale the number and computational capabilities of the cores. However, the efficiency gap can be bridged by introducing prefetching and non-allocating write policies to cache-coherent systems. We also found some applications for which streaming scales worse than caching due to the redundancies introduced in order to make the code regular enough for streaming.

Through the adoption of *stream programming* methodologies, which encourage blocking, macroscopic prefetching, and locality aware task scheduling, cache-based systems are equally as efficient as streaming memory systems. This indicates that there is not a sufficient advantage in building general-purpose cores that follow a pure streaming model, where all local memory structures and all data streams are explicitly managed. We also observed that stream programming is actually easier when targeting cache-based systems rather than streaming memory systems, and that it may be beneficial in scaling coherence and consistency for caches to larger systems.

10. ACKNOWLEDGEMENTS

We sincerely thank Bill Mark, Mattan Erez, and Ron Ho for their invaluable comments on early versions of this paper. We are also grateful for the insightful critique of this work by the reviewers. This work was partially funded by DARPA under contract number F29601-03-2-0117. Jacob Leverich is supported by a David Cheriton Stanford Graduate Fellowship.

11. REFERENCES

- [1] S. V. Adve and K. Gharachorloo. Shared Memory Consistency Models: A Tutorial. *IEEE Computer*, 29(12):66–76, Dec. 1996.
- [2] V. Agarwal, M. S. Hrishikesh, S. W. Keckler, and D. Burger. Clock Rate versus IPC: the End of the Road for Conventional Microarchitectures. In *Proceedings of the 27th Intl. Symp. on Computer Architecture*, June 2000.
- [3] J. Ahn et al. Evaluating the Imagine Stream Architecture. In *Proceedings of the 31st Intl. Symp. on Computer Architecture*, May 2004.
- [4] J. Andrews and N. Backer. Xbox360 System Architecture. In *Conf. Record of Hot Chips 17*, Stanford, CA, Aug. 2005.
- [5] L. A. Barroso et al. Piranha: A Scalable Architecture Based on Single-Chip Multiprocessing. In *Proceedings of the 27th Intl. Symp. on Computer Architecture*, Vancouver, Canada, June 2000.
- [6] I. Buck. GPU Computing: Programming a Massively Parallel Processor, Mar. 2005. Keynote presentation at the International Symposium on Code Generation and Optimization, San Jose, CA.
- [7] T. Chiueh. A Generational Algorithm to Multiprocessor Cache Coherence. In *International Conference on Parallel Processing*, pages 20–24, Oct. 1993.
- [8] D. Culler, J. P. Singh, and A. Gupta. *Parallel Computer Architecture: A Hardware/Software Approach*. Morgan Kaufman, 1999.
- [9] W. Dally et al. Merrimac: Supercomputing with Streams. In *Proceedings of the 2003 Conf. on Supercomputing*, Nov. 2003.
- [10] J. D. Davis, J. Laudon, and K. Olukotun. Maximizing CMP Throughput with Mediocre Cores. In *Proceedings of the 14th Intl. Conf. on Parallel Architectures and Compilation Techniques*, Sept. 2005.
- [11] M. Drake, H. Hoffmann, R. Rabbah, and S. Amarasinghe. MPEG-2 Decoding in a Stream Programming Language. In *Proceedings of the 20th IEEE International Parallel & Distributed Processing Symposium, Rhodes Island (IPDPS)*, Apr. 2006.
- [12] W. Eatherton. The Push of Network Processing to the Top of the Pyramid, Oct. 2005. Keynote presentation at the Symposium on Architectures for Networking and Communication Systems, Princeton, NJ.
- [13] K. Fatahalian et al. Sequoia: Programming The Memory Hierarchy. In *Supercomputing Conference*, Nov. 2006.
- [14] T. Foley and J. Sugerman. KD-Tree Acceleration Structures for a GPU Raytracer. In *Proceedings of the Graphics Hardware Conf.*, July 2005.
- [15] M. I. Gordon et al. A Stream Compiler for Communication-exposed Architectures. In *Proceedings of the 10th Intl. Conf. on Architectural Support for Programming Languages and Operating Systems*, Oct. 2002.
- [16] M. Gschwind et al. A Novel SIMD Architecture for the Cell Heterogeneous Chip-Multiprocessor. In *Conf. Record of Hot Chips 17*, Stanford, CA, Aug. 2005.
- [17] J. Gummaraju and M. Rosenblum. Stream Programming on General-Purpose Processors. In *Proceedings of the 38th Intl. Symp. on Microarchitecture*, Nov. 2005.
- [18] R. Ho, K. Mai, and M. Horowitz. The Future of Wires. *Proceedings of the IEEE*, 89(4), Apr. 2001.
- [19] R. Ho, K. Mai, and M. Horowitz. Efficient On-chip Global Interconnects, June 2003.
- [20] M. Horowitz and W. Dally. How Scaling Will Change Processor Architecture. In *International Solid-State Circuits Conference*, pages 132–133, Feb. 2004.
- [21] Independent JPEG Group. JG’s JPEG Software Release 6b, 1998.
- [22] D. Jani, G. Ezer, and J. Kim. Long Words and Wide Ports: Reinventing the Configurable Processor. In *Conf. Record of Hot Chips 16*, Stanford, CA, Aug. 2004.
- [23] N. Jayasena. *Memory Hierarchy Design for Stream Computing*. PhD thesis, Stanford University, 2005.
- [24] A. C. Klaiber and H. M. Levy. A Comparison of Message Passing and Shared Memory Architectures for Data Parallel Programs. In *Proceedings of the 21th Intl. Symp. on Computer Architecture*, Apr. 1994.
- [25] P. Kongetira. A 32-way Multithreaded Sparc Processor. In *Conf. Record of Hot Chips 16*, Stanford, CA, Aug. 2004.
- [26] R. Kumar, V. Zyuban, and D. M. Tullsen. Interconnections in Multi-Core Architectures: Understanding Mechanisms, Overheads and Scaling. In *Proceedings of the 32nd Intl. Symp. on Computer Architecture*, June 2005.
- [27] B. Lewis and D. J. Berg. *Multithreaded Programming with Pthreads*. Prentice Hall, 1998.
- [28] M. Li et al. ALP: Efficient Support for All Levels of Parallelism for Complex Media Applications. Technical Report UIUCDCS-R-2005-2605, UIUC CS, July 2005.
- [29] A. W. Lim, S.-W. Liao, and M. S. Lam. Blocking and Array Contraction Across Arbitrarily Nested Loops Using Affine Partitioning. *ACM SIGPLAN Notices*, 36(7):103–112, July 2001.
- [30] Y. Lin. A Programmable Vector Coprocessor Architecture for Wireless Applications. In *Proceedings of the 3rd Workshop on Application Specific Processors*, Sept. 2004.
- [31] M. Loghi and M. Pncino. Exploring Energy/Performance Tradeoffs in Shared Memory MPSoCs: Snoop-Based Cache Coherence vs. Software Solutions. In *Proceedings of the Design Automation and Test in Europe Conf.*, Mar. 2005.
- [32] E. Machnicki. Ultra High Performance Scalable DSP Family for Multimedia. In *Conf. Record of Hot Chips 17*, Stanford, CA, Aug. 2005.
- [33] K. Mai et al. Smart Memories: a Modular Reconfigurable Architecture. In *Proceedings of the 27th Intl. Symp. on Computer Architecture*, June 2000.
- [34] MIPS32 Architecture For Programmers Volume II: The MIPS32 Instruction Set. MIPS Technologies, Inc., 2001.
- [35] A. Moshovos. RegionScout: Exploiting Coarse Grain Sharing in Snoop-Based Coherence. In *Proceedings of the 32nd Intl. Symp. on Computer Architecture*, June 2005.
- [36] K. Sankaralingam. TRIPS: A Polymorphous Architecture for Exploiting ILP, TLP, and DLP. *ACM Trans. Archit. Code Optim.*, 1(1):62–93, Mar. 2004.
- [37] J. Suh et al. A Performance Analysis of PIM, Stream Processing, and Tiled Processing on Memory-intensive Signal Processing Kernels. In *Proceedings of the 30th Intl. Symp. on Computer Architecture*, June 2003.
- [38] D. Tarjan, S. Thoziyoor, and N. P. Jouppi. CACTI 4.0. Technical Report HPL-2006-86, HP Labs, 2006.
- [39] M. Taylor et al. Evaluation of the Raw Microprocessor: An Exposed-Wire-Delay Architecture for ILP and Streams. In *Proceedings of the 31st Intl. Symp. on Computer Architecture*, May 2004.
- [40] Tensilica Software Tools. <http://www.tensilica.com/products/software.htm>.
- [41] S. P. VanderWiel and D. J. Lilja. Data Prefetch Mechanisms. *ACM Computing Surveys*, 32(2):174–199, 2000.
- [42] D. Wang et al. DRAMsim: A Memory-System Simulator. *SIGARCH Computer Architecture News*, 33(4), 2005.
- [43] Z. Wang et al. Using the Compiler to Improve Cache Replacement Decisions. In *Proceedings of the Conf. on Parallel Architectures and Compilation Techniques*, Sept. 2002.
- [44] Z. Wang et al. Guided Region Prefetching: A Cooperative Hardware/Software Approach. In *Proceedings of the 30th Intl. Symp. on Computer Architecture*, June 2003.
- [45] T.-Y. Yeh. The Low-Power High-Performance Architecture of the PWRficient Processor Family. In *Conf. Record of Hot Chips 17*, Stanford, CA, Aug. 2005.