



June 2003

Token Coherence: Decoupling Performance and Correctness

Milo Martin

University of Pennsylvania, milom@cis.upenn.edu

Mark D. Hill

University of Wisconsin

David A. Wood

University of Wisconsin

Follow this and additional works at: https://repository.upenn.edu/cis_papers

Recommended Citation

Milo Martin, Mark D. Hill, and David A. Wood, "Token Coherence: Decoupling Performance and Correctness", . June 2003.

Copyright 2003 IEEE. Reprinted from *Proceedings of the 30th Annual International Symposium on Computer Architecture (ISCA'03)*, pages 182-193.

This material is posted here with permission of the IEEE. Such permission of the IEEE does not in any way imply IEEE endorsement of any of the University of Pennsylvania's products or services. Internal or personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution must be obtained from the IEEE by writing to pubs-permissions@ieee.org. By choosing to view this document, you agree to all provisions of the copyright laws protecting it.

At the time of publication, author Milo M.K. Martin was affiliated with the University of Michigan. Currently, November 2006, he is a faculty member in the Department of Computer and Information Science at the University of Pennsylvania.

This paper is posted at ScholarlyCommons. https://repository.upenn.edu/cis_papers/264
For more information, please contact repository@pobox.upenn.edu.

Token Coherence: Decoupling Performance and Correctness

Abstract

Many future shared-memory multiprocessor servers will both target commercial workloads and use highly-integrated "glueless" designs. Implementing low-latency cache coherence in these systems is difficult, because traditional approaches either add indirection for common cache-to-cache misses (directory protocols) or require a totally-ordered interconnect (traditional snooping protocols).

Unfortunately, totally-ordered interconnects are difficult to implement in glueless designs. An ideal coherence protocol would avoid indirections and interconnect ordering; however, such an approach introduces numerous protocol races that are difficult to resolve. We propose a new coherence framework to enable such protocols by separating performance from correctness. A performance protocol can optimize for the common case (i.e., absence of races) and rely on the underlying correctness substrate to resolve races, provide safety, and prevent starvation. We call the combination *Token Coherence*, since it explicitly exchanges and counts tokens to control coherence permissions.

This paper develops *TokenB*, a specific Token Coherence performance protocol that allows a glueless multiprocessor to both exploit a low-latency unordered interconnect (like directory protocols) and avoid indirection (like snooping protocols). Simulations using commercial workloads show that our new protocol can significantly outperform traditional snooping and directory protocols.

Comments

Copyright 2003 IEEE. Reprinted from *Proceedings of the 30th Annual International Symposium on Computer Architecture (ISCA'03)*, pages 182-193.

This material is posted here with permission of the IEEE. Such permission of the IEEE does not in any way imply IEEE endorsement of any of the University of Pennsylvania's products or services. Internal or personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution must be obtained from the IEEE by writing to pubs-permissions@ieee.org. By choosing to view this document, you agree to all provisions of the copyright laws protecting it.

At the time of publication, author Milo M.K. Martin was affiliated with the University of Michigan. Currently, November 2006, he is a faculty member in the Department of Computer and Information Science at the University of Pennsylvania.

Token Coherence: Decoupling Performance and Correctness

Milo M. K. Martin, Mark D. Hill, and David A. Wood
Computer Sciences Department
University of Wisconsin-Madison
{milo, markhill, david}@cs.wisc.edu

Abstract

Many future shared-memory multiprocessor servers will both target commercial workloads and use highly-integrated “glueless” designs. Implementing low-latency cache coherence in these systems is difficult, because traditional approaches either add indirection for common cache-to-cache misses (directory protocols) or require a totally-ordered interconnect (traditional snooping protocols). Unfortunately, totally-ordered interconnects are difficult to implement in glueless designs. An ideal coherence protocol would avoid indirections and interconnect ordering; however, such an approach introduces numerous protocol races that are difficult to resolve.

We propose a new coherence framework to enable such protocols by separating performance from correctness. A performance protocol can optimize for the common case (i.e., absence of races) and rely on the underlying correctness substrate to resolve races, provide safety, and prevent starvation. We call the combination Token Coherence, since it explicitly exchanges and counts tokens to control coherence permissions.

This paper develops TokenB, a specific Token Coherence performance protocol that allows a glueless multiprocessor to both exploit a low-latency unordered interconnect (like directory protocols) and avoid indirection (like snooping protocols). Simulations using commercial workloads show that our new protocol can significantly outperform traditional snooping and directory protocols.

1 Introduction

The performance and cost of database and web servers is important because the services they provide are becoming increasingly part of our daily lives. Many of these servers are shared-memory multiprocessors. In our view, workload and technology trends point toward a new design space that provides opportunities to improve performance and cost of these multiprocessor servers.

This work is supported in part by the National Science Foundation (EIA-9971256, EIA-0205286, CDA-9623632, and CCR-0105721), a Norm Koo Graduate Fellowship and an IBM Graduate Fellowship (Martin), two Wisconsin Romnes Fellowships (Hill and Wood), Spanish Secretaría de Estado de Educación y Universidades (Hill sabbatical), and donations from Intel Corporation, IBM, and Sun Microsystems.

Workload trends. Since many commercial workloads exhibit abundant thread-level parallelism, using multiple processors is an attractive approach for increasing their performance. To efficiently support the frequent communication and synchronization in these workloads [8], servers should optimize the latency of cache-to-cache misses (*i.e.*, those misses—often caused by accessing shared data—that require data to move directly between caches).

To reduce cache-to-cache miss latency, many multiprocessor servers use snooping cache coherence. Bus-based snooping protocols exploit totally-ordered broadcasts (*i.e.*, all processors are guaranteed to observe all broadcasts in the same order) to both satisfy cache-to-cache misses directly and resolve protocol races. Protocol races can occur, for example, when two processors request the same block at the same time.

To overcome the increasingly difficult challenge of scaling the bandwidth of shared-wire buses [13], some recent snooping designs broadcast requests on a “virtual bus” created with an indirect switched interconnect (*e.g.*, Figure 1a). These interconnects can provide higher bandwidth than buses, at the cost of additional switch chips. The use of broadcast limits snooping’s scalability, but small- to medium-sized snooping-based multiprocessors (4-16 processors) suffice for many workloads, because larger services tend to cluster machines to increase both throughput and availability.

Technology trends. The increasing number of transistors per chip predicted by Moore’s Law has encouraged and will continue to encourage more integrated designs, making “glue” logic (*e.g.*, discrete switch chips) less desirable. Many current and future systems will integrate processor(s), cache(s), coherence logic, switch logic, and memory controller(s) on a single die (*e.g.*, Alpha 21364 [32] and AMD’s Hammer [5]). Directly connecting these highly-integrated nodes leads to a high-bandwidth, low-cost, low-latency “glueless” interconnect (*e.g.*, Figure 1b).

These glueless interconnects are fast but do not easily provide the virtual bus behavior required by traditional snooping protocols. Instead, most such systems use directory protocols, which provide coherence without requiring broadcast or a totally-ordered interconnect. These systems maintain a directory at the home node (*i.e.*, memory) that

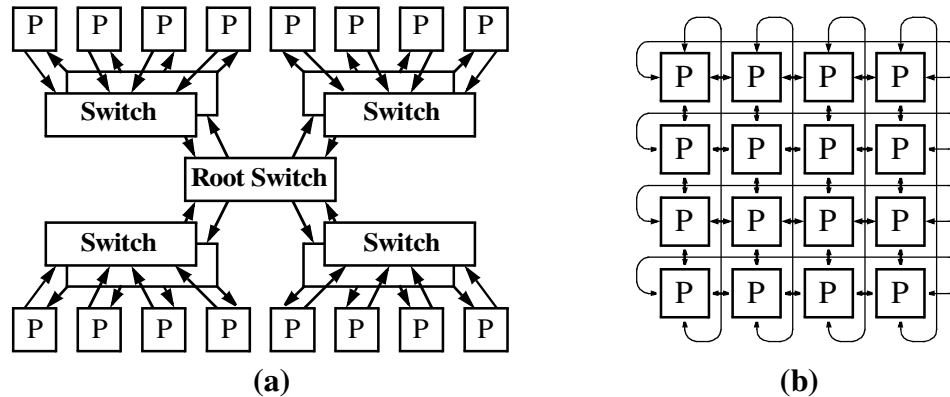


Figure 1. (a) 16-processor two-level tree interconnect and (b) 16-processor (4x4) bi-directional torus interconnect. The boxes marked “P” represent highly-integrated nodes that include a processor, caches, memory controller, and coherence controllers. The indirect broadcast tree uses discrete switches, while the torus is a directly connected interconnect. In this example, the torus has lower latency (two vs. four chip crossings on average) and does not require any glue chips; however, unlike the indirect tree, the torus provides no request total order, making it unsuitable for traditional snooping.

resolves some protocol races by ordering requests on a per-cache-block basis. Unfortunately, traditional directory protocols must first send all requests to the home node, adding an indirection to the critical path of cache-to-cache misses—a poor match for commercial workloads.

Our approach. Ideally, a coherence protocol would both avoid indirection latency for cache-to-cache misses (like snooping protocols) and not require any interconnect ordering (like directory protocols). One obvious approach is to directly send broadcasts on an unordered interconnect. This general approach has not been used, however, because it suffers from numerous race cases that are difficult to make correct (as discussed in Section 2).

Token Coherence. Rather than abandoning this fast approach, we use it to make the common case fast, but we back it up with a substrate that ensures correctness. To this end, we propose *Token Coherence*, which has two parts: a correctness substrate and a performance protocol.

- **Correctness substrate:** The substrate (as described in Section 3) provides a foundation for building correct coherence protocols on unordered interconnects by ensuring safety (*i.e.*, guaranteeing all reads and writes are coherent) and starvation avoidance (*i.e.*, guaranteeing all reads and writes are eventually completed). The substrate ensures safety by (1) associating a fixed number of *tokens* with each logical block of shared memory, and (2) ensuring that a processor may read a cache block only if it holds at least one of the block’s tokens, and it may write a cache block only if it holds all of the block’s tokens (allowing for a single writer or many readers, but not both). Tokens are held with copies of the block in caches and memory and exchanged using coherence messages. The substrate provides starvation freedom via *persistent requests*, which a processor invokes when it detects possible starvation. Persistent requests always succeed in obtaining data and tokens—even when races

occur—because once activated they persist in forwarding data and tokens until the request is satisfied.

- **Performance protocol:** Performance protocols (as described in Section 4.1) use *transient requests* as “hints” to direct the correctness substrate to send data and tokens to the requesting processor. In the common case, a transient request succeeds in obtaining the requested data and tokens. However, transient requests may fail to complete, principally due to races. Since the correctness substrate prevents starvation (via persistent requests) and guarantees safety (via token counting), performance protocol bugs and various races may hurt performance, but they cannot affect correctness.

Since Token Coherence never speculatively modifies memory state, it is not a speculative execution technique, and it requires no rollback or recovery mechanism.

TokenB. We target medium-sized glueless multiprocessors with unordered interconnects using a specific performance protocol called *Token-Coherence-using-Broadcast* or *TokenB* (Section 4.2). In *TokenB*, processors broadcast transient requests and respond like a traditional MOSI snooping protocol. *TokenB* directly locates the data in the common case of no races, allowing for low-latency cache-to-cache misses. When transient requests fail (due to races), the protocol reissues them until the processor times out and invokes a persistent request to prevent starvation.

For selected commercial workloads on a full-system simulation of a 16-processor system (described in Section 5), we find that (1) reissued and persistent requests are rare (3.0% and 0.2% of requests, respectively), (2) *TokenB* is faster than traditional snooping, because it allows use of an unordered interconnect (15-28% faster), (3) *TokenB* is faster than a directory protocol, because it avoids directory indirections (17-54%), and (4) a directory protocol uses less bandwidth than *TokenB* (21-25% for 16 processors), but this additional bandwidth may not be a significant

problem for the high-bandwidth glueless interconnects that will be common in future systems. We present these and other results in Section 6.

While *TokenB* provides an attractive alternative for small- to medium-sized systems, Token Coherence is a general coherence framework that enables the creation of other performance protocols (described in Section 7) that can reduce traffic for larger systems, use prediction to push data, and support hierarchy with low complexity.

2 A Motivating Example Race

In this section, we first present a simple coherence race to illustrate that naively sending requests without an ordering point is incorrect. We then review how the race is handled by traditional snooping protocols (by ordering requests in the interconnect) and directory protocols (by ordering requests at the home node). Finally, we forecast how Token Coherence handles this and other races.

Fast (but incorrect) approach. Invalidation-based coherence protocols provide the illusion of a single memory shared by all processors by allowing either a single writer or many readers, but not both at the same time. A fast way to obtain read or write permission to a cache block would allow requesters to broadcast requests to all other processors over a low-latency unordered interconnect. Doing this naively, however, may allow a processor to erroneously read stale (or incoherent) data when processors act on racing requests in different orders.

Consider the example illustrated in Figure 2a, in which processor P_0 desires read/write access to the block (*i.e.*, MOESI [41] state *modified* or *M*), and processor P_1 desires read-only access (*i.e.*, state *shared* or *S*). P_0 broadcasts its request at time ①, which the interconnect promptly delivers to P_1 at time ② but belatedly to memory at time ⑥ (*e.g.*, due to contention on the unordered interconnect). Processor P_1 handles P_0 's request but takes no action other than an invalidation acknowledgment at time ③, because it lacks a valid copy. Later, at time ④, P_1 issues its request, which the memory quickly satisfies at time ⑤. Finally, P_0 's delayed request arrives at memory at time ⑥, and the memory satisfies the request at time ⑦. After receiving both responses, P_0 believes (erroneously) that it has a writable copy of the block, but P_1 still holds a read-only copy. If this situation arises, the memory consistency model—the definition of memory access correctness in a multiprocessor system—may be violated.

Snooping protocols. Traditional split-transaction snooping protocols resolve this example race and other races by relying on a totally-ordered interconnect—a virtual bus—to provide a total order of all requests. This ordering ensures that all processors (and memories) observe requests in the same order (including their own requests, to establish their place in the total order). In our example

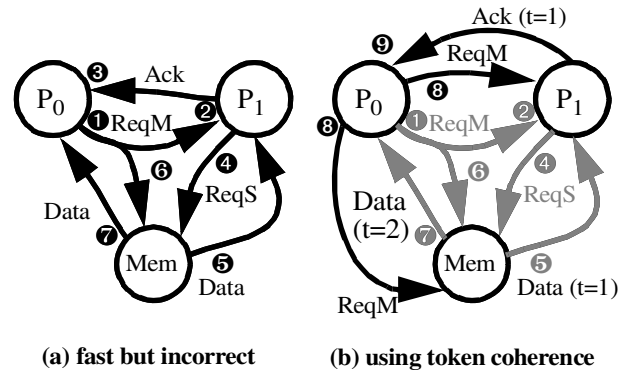


Figure 2. Example Race. A request for shared (ReqS) racing with a request for modified (ReqM).

race, request ordering was not consistent with a total order. P_1 observed its request as occurring *after* P_0 's request, while the memory observed P_1 's request *before* P_0 's request. A total order would guarantee correct operation because either P_0 's invalidation would have occurred before P_1 's request (and thus P_0 would transition to read-only and respond to P_1 's request), or P_0 's request would have arrived at P_1 after P_1 's request (invalidating P_1 's shared copy of the block). Unfortunately, interconnects that enforce a total order may have higher latency or cost.

Directory protocols. Directory protocols resolve this example race and other races without an ordered interconnect by providing a per-block ordering point at the directory. A directory protocol prevents this example race by (1) relying on the directory controller to determine which request will be satisfied first, (2) using forwarded requests, invalidations, and explicit acknowledgements to enforce ordering, and (3) possibly blocking, queuing, or negatively acknowledging requests in some cases. In our example, P_1 's request arrives at the home memory/directory first, and the memory provides data. P_0 's request arrives later, and the home forwards an invalidation message to P_1 . P_0 's request completes after receiving an acknowledgment from P_1 and data from the home, knowing that no other copies are present in the system. Unfortunately, the cost of this solution is an added level of indirection on the critical path of cache-to-cache misses.

Token Coherence. Token Coherence allows races to occur but provides correct behavior in all cases with a *correctness substrate*. This substrate ensures that processors only read and write coherent blocks appropriately (safety) and that processors eventually obtain a needed block (starvation avoidance). *Performance protocols* seek to make the common case fast with requests (or hints) for data movement that do not always succeed due to races. After describing the correctness substrate (in Section 3) and performance protocols (in Section 4), we revisit this example race in the context of Token Coherence (in Section 4.2).

3 Correctness Substrate

The correctness substrate uses token counting to enforce safety (do no harm), and it uses persistent requests to prevent starvation (do some good). These two mechanisms allow the substrate to move data around the system without concern for order or races, allowing processors to only read or write the block as appropriate, but still ensures that a request for a block will eventually succeed in all cases.

3.1 Enforcing Safety via Token Counting

The correctness substrate uses tokens to ensure safety without requiring indirection or a totally-ordered interconnect. The system associates a fixed number of *tokens* with each block of shared memory, and a processor is only allowed to read a cache block when it holds at least one token, or write a cache block when holding all tokens.

During system initialization, the system assigns each block T tokens (where T is at least as large as the number of processors). Initially, the block's home memory module holds all tokens for a block. Later, tokens are held also by processor caches and coherence messages. Tokens and data are allowed to move throughout the system as long as the substrate maintains these four invariants:

- **Invariant #1:** At all times, each block has T tokens in the system.
- **Invariant #2:** A processor can write a block only if it holds all T tokens for that block.
- **Invariant #3:** A processor can read a block only if it holds at least one token for that block.
- **Invariant #4:** If a coherence message contains one or more tokens, it must contain data.

Invariant #1 ensures that the substrate never creates or destroys tokens. Invariants #2 and #3 ensure that a processor will not write the block while another processor is reading it. Adding invariant #4 ensures that processors holding tokens always have a valid copy of the data block. In more familiar terms, token possession maps directly to traditional coherence states: holding all T tokens is *modified* (M); one to $T-1$ tokens is *shared* (S); and no tokens is *invalid* (I).

The token-based correctness substrate enforces these invariants directly by counting tokens. The substrate maintains these invariants by induction; the invariants hold for the initial system state, and all movements of data and tokens preserve the invariants. Thus, safety is ensured without reasoning about the interactions among non-stable protocol states, data responses, acknowledgment messages, interconnect ordering, or system hierarchy.

Token Coherence enforces a memory consistency model [4]—the definition of correctness for multiprocessor systems—in a manner similar to directory protocols. The above guarantee of a “single writer” or “multiple readers with no writer” is the same property provided by tradi-

tional invalidation-based directory protocols. For example, the MIPS R10k processors [43] in the Origin 2000 [23] use this guarantee to provide sequential consistency, even without a global ordering point¹. The Origin protocol uses explicit invalidation acknowledgments to provide the above guarantee. We provide the same guarantee by explicitly tracking tokens for each block. As with any coherence scheme, the processors are also intimately involved in enforcing the memory consistency model.

Optimized token counting. An issue with the invariants above is that data must always travel with tokens, even when gathering tokens from shared copies. To avoid this bandwidth inefficiency, the substrate actually allows tokens to be transferred without data (similar to the data-less invalidation acknowledgment messages in a directory protocol). To enable this optimization, the substrate distinguishes a separate *owner token*, adds a data valid bit (distinct from the traditional tag valid bit), and maintains the following four invariants (changes in *italics*):

- **Invariant #1'**: At all times, each block has T tokens in the system, *one of which is the owner token*.
- **Invariant #2'**: A processor can write a block only if it holds all T tokens for that block.
- **Invariant #3'**: A processor can read a block only if it holds at least one token for that block *and has valid data*.
- **Invariant #4'**: If a coherence message contains *the owner token*, it must contain data.

Invariants #1', #2' and #3' continue to provide safety. Invariant #4' allows coherence messages with non-owner tokens to omit data, but it still requires that messages with the owner token contain data (to prevent all processors from simultaneously discarding data). Possession of the owner token but not all other tokens maps to the familiar MOESI state *owned* (O). System components (processors and the home memory) maintain a valid bit, to allow components to receive and hold non-owner tokens without valid data. A component sets the valid bit when a message with data and at least one token arrives, and a component clears the valid bit when it no longer holds any tokens.

Tokens are held in processor caches (*e.g.*, part of tag state), memory (*e.g.*, encoded in ECC bits [17, 32, 34]), and coherence messages². Since we do not track which processors hold them but only count tokens, tokens can be stored in $2 + \lceil \log_2 T \rceil$ bits (valid bit, owner-token bit, and non-owner token count). For example, encoding 64 tokens with 64-byte blocks adds one byte of storage (1.6% overhead).

1. The Origin protocol uses a directory to serialize some requests *for the same block*; however, since memory consistency involves the ordering relationship between different memory locations [4], using a distributed directory is not alone sufficient to implement a memory consistency.

2. Like most coherence protocols (*e.g.*, [11, 23, 32, 42]), we assume the interconnect provides reliable message delivery.

Finally, there is important freedom in what the invariants do *not* specify. While our invariants restrict the data and token content of coherence messages, the invariants do not restrict *when* or *to whom* the substrate can send coherence messages. For example, to evict a block (and thus tokens) from a cache, the processor simply sends all its tokens (and data if the message includes the owner token) to the memory. Likewise, anytime a processor receives a message carrying tokens (with or without data), it can either choose to accept it (*e.g.*, if there is space in the cache) or redirect it to memory (using another virtual network to avoid deadlock). We use this freedom in three additional ways. First, we define persistent requests to prevent starvation (Section 3.2). Second, we define transient requests that allow a performance protocol to send “hints” to inform the substrate to which processor it should send data and tokens (Section 4.1). Third, this freedom enables many performance protocols (Section 4.2 and Section 7).

3.2 Avoiding Starvation via Persistent Requests

The correctness substrate provides persistent requests to prevent starvation. A processor invokes a persistent request whenever it detects possible starvation (*e.g.*, it has failed to complete a cache miss within a timeout period). Since processors should only infrequently resort to persistent requests, persistent requests must be correct but not necessarily fast. The substrate uses persistent requests to prevent starvation by performing the following steps:

- When a processor detects possible starvation, it initiates a persistent request.
- The substrate *activates* at most one persistent request per block.
- System nodes remember all *activated* persistent requests and forward all tokens for the block—those tokens currently present and received in the future—to the initiator of the request.
- When the initiator has sufficient tokens, it performs a memory operation (*e.g.*, a load or store instruction) and *deactivates* its persistent request.

To guarantee starvation freedom, the system must provide a fair mechanism for activating persistent requests.

Implementation. Processors invoke a persistent request when a cache miss has not been satisfied within ten average miss times. The correctness substrate implements persistent requests with a simple *arbiter* state machine at each home memory module. The substrate directs persistent requests to the home node of the requested block. Requests may queue in a dedicated virtual network or at the home node. The arbiter state machine activates at most one request by informing all nodes. Each node responds with an acknowledgement (to avoid races) and remembers all active persistent requests using a hardware table. This table contains an 8-byte entry per arbiter (*i.e.*, per home memory module). For example, a 64-node system requires

only a 512-byte table at each node. While a persistent request is active, nodes must forward all tokens (and data, if they have the owner token) to the requester. The node will also forward tokens (and data) that arrive later, because the request persists until the requester explicitly deactivates it. Once the requester is satisfied, it sends a message to the arbiter at the home memory module to deactivate the request. The arbiter deactivates the request by informing all nodes, who delete the entry from their table and send an acknowledgement (again, to eliminate races). Figure 3 shows the general operation of our implementation of the correctness substrate.

While this implementation of activating persistent requests is sufficient for our experiments, we are currently developing a distributed arbitration scheme that efficiently transfers highly-contended blocks directly between contending processors.

4 Performance Protocols

This section first discusses performance protocol requirements and then presents *TokenB*, a performance protocol targeting medium-sized glueless multiprocessors.

4.1 Obligations and Opportunities

Obligations. Performance protocols have no obligations, because the processors and correctness substrate ensure correctness. A null or random performance protocol would perform poorly but not incorrectly. Therefore, performance protocols may aggressively seek performance without concern for corner-case errors.

Opportunities via Transient Requests. One way in which performance protocols seek high performance is by specifying a policy for using *transient requests*. Transient requests are fast, unordered “hint” requests sent to one or more nodes that often succeed, but may fail to obtain a readable or writable block due to races, insufficient recipients, or being ignored by the correctness substrate. Performance protocols can detect when a request has not (yet) succeeded, because the requester has not obtained sufficient tokens (*i.e.*, one token to read the block, and all tokens to write it). Performance protocols may reissue transient requests or do nothing (since the processor will eventually timeout and issue a persistent request).

A performance protocol also specifies a policy for how system components respond to transient requests. If a transient request for a shared block encounters data with all tokens, for example, a performance protocol can specify whether the substrate should reply with the data and one token or the data and all tokens (much like a migratory sharing optimization [12, 40]). Active persistent requests always override performance protocol policies to prevent starvation.

A good performance protocol will use transient requests to quickly satisfy most cache misses. Returning to the exam-

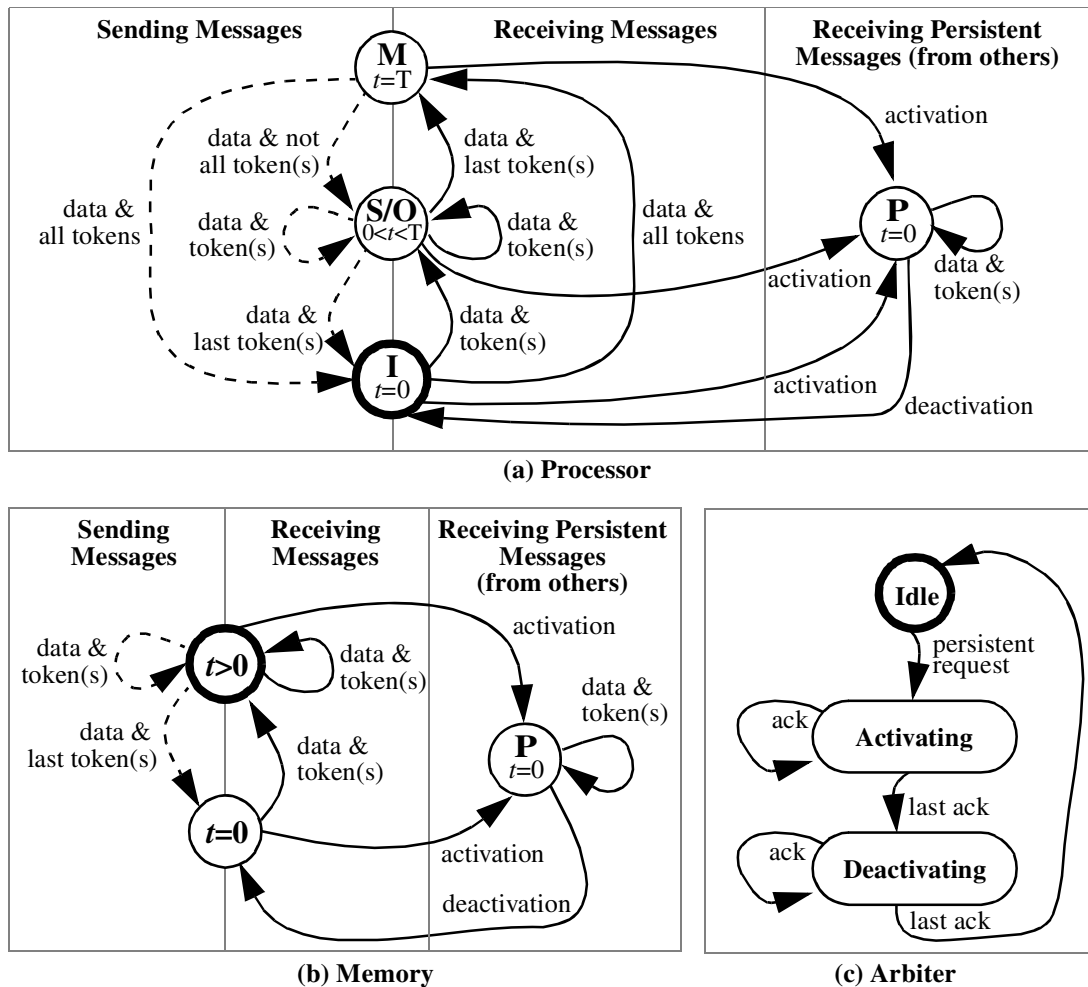


Figure 3. Correctness substrate state transitions for the (a) processor, (b) memory, and (c) persistent request arbiter. As a simplification, the figure shows only tokens sent with data. The symbol t represents the current token count, and T represents all the tokens. Solid arcs are transitions in response to incoming messages. Dashed arcs are transitions a performance protocol (Section 4) can invoke at any time (e.g., when receiving a transient request). The “P” states occur when a node receives a another processor’s persistent request from the arbiter. Each processor must also remember its own persistent request, not explicitly shown in this figure. The initial states are emphasized with thick borders.

ple from Section 2 (illustrated in Figure 2b), both processors could broadcast transient requests. Even though the requests race, frequently both processor’s misses would be satisfied. In other cases, one or both may not succeed (detected by insufficient tokens and a timeout). When this occurs, the performance protocol can reissue those transient requests. In the worst case, one or both processors may time out and issue persistent requests.

4.2 TokenB: Targeting Glueless Multiprocessors

The *Token-Coherence-using-Broadcast (TokenB)* performance protocol uses three policies to target glueless multiprocessors with high-bandwidth unordered interconnects.

Issuing transient requests. Processors broadcast all transient requests. This policy works well for moderate-sized systems where interconnect bandwidth is plentiful and when racing requests are rare.

Responding to transient requests. Components (processors and the home memory) respond³ to transient requests as they would in most MOSI protocols. A component with no tokens (state I) ignores all requests. A component with only non-owner tokens (state S) ignores shared requests, but on an exclusive request it sends all its tokens in a data-less message (like an invalidation acknowledgment in a directory protocol). A component with the owner token but not all other tokens (state O) sends the data with one token (usually not the owner token) on a shared request, and it sends the data and all its tokens on an exclusive request. A component with all the tokens (state M) responds the same way as a component in state O , with the exception given in the next paragraph.

3. Technically, the performance protocol asks the correctness substrate to respond on its behalf.

To optimize for common migratory sharing patterns, we implement a well-known optimization for migratory data [12, 40]. If a processor with all tokens (state M) has written the block, and it receives a shared request, it provides read/write permission to the block by responding with the data and all tokens (instead of the data and one token). We also implement an analogous optimization in all other protocols we compare against in the evaluation.

Reissuing transient requests. If a transient request has not completed after a reasonable interval, we reissue the transient request. We continue to reissue transient requests until the processor invokes the persistent request mechanism (approximately 4 times). We use both a small randomized exponential backoff (much like ethernet) and twice the recent average miss latency to calculate the re-issue timeout. This policy adapts to the average miss latency of the system (to avoid reissuing too soon), but it also quickly reissues requests that do not succeed due to occasional races. Since races are rare, on average only 3.0% of all misses are reissued even once (for our workloads and simulation assumptions, described next).

Example. Returning to the example race in Section 2 (Figure 2b), the block has three tokens that are initially held by memory. P_1 received one token in the response at time ⑤, allowing it to read the block. Due to the race, P_0 only received two tokens in the response at time ⑦ but requires all three before it can write the block. After the specified timeout interval, *TokenB* reissues P_0 's request at time ⑧. P_1 responds with the missing token at time ⑨, allowing P_0 to finally complete its request.

5 Evaluation Methods

To evaluate Token Coherence, we simulate a multiprocessor server running commercial workloads using multiple interconnection networks and coherence protocols. Our target system is a 16-processor SPARC v9 system with highly integrated nodes that each include a pipelined dynamically scheduled processor, two levels of cache, coherence protocol controllers, and a memory controller for part of the globally shared memory. The system implements sequential consistency using invalidation-based cache coherence and an aggressive, speculative processor implementation [18, 43].

Our benchmarks consist of three commercial workloads: an online transaction processing workload (OLTP), a static web serving workload (Apache), and a Java middleware workload (SPECjbb). We refer interested readers to Alameldeen *et al.* [6] for a more detailed description and characterization of these workloads.

We selected a number of coherence protocols, interconnection networks, latencies, bandwidths, cache sizes, and other structure sizes. Table 1 lists the system parameters for both the memory system and the processors, chosen to

Table 1. Target System Parameters

Coherent Memory System	
split L1 I & D caches	128kBytes, 4-way, 2ns latency
unified L2 cache	4MBytes, 4-way, 6ns latency
cache block size	64 Bytes
DRAM/dir. latency	80ns (2 GBytes of DRAM)
memory/dir. controllers	6ns latency
network link bandwidth	3.2 GBytes/sec
network link latency	15ns (incl. wire, sync. & route)
Dynamically Scheduled Processors	
clock frequency	1 Ghz
reorder buffer/scheduler	128/64 entries
pipeline width	4-wide fetch & issue
pipeline stages	11
direct branch predictor	1kBytes YAGS
indirect branch predictor	64 entry (cascaded)
return address stack	64 entry

approximate the published parameters of systems like the Alpha 21364 [32]. The coherence protocols and interconnection networks are described below.

5.1 Coherence Protocols

We compare target systems using four distinct MOSI coherence protocols. They all implement the previously described migratory sharing optimization that improves the performance of all the protocols. All request, acknowledgment, invalidation, and dataless token messages are 8 bytes in size (including the 40+ bit physical address and token count if needed); data messages include this 8 byte header and 64 bytes of data. We compare the *TokenB* protocol (described in Section 4) with three other coherence protocols:

Snooping. We based our traditional snooping protocol on a modern protocol [11], and we added additional non-stable states to relax synchronous timing requirements. To avoid the complexity and latency of a snoop response combining tree to implement the “owner” signal, the protocol uses a single bit in memory to determine when the memory should respond to requests [16].

Directory. We use a standard full-map directory protocol inspired by the Origin 2000 [23] and Alpha 21364 [32]. The protocol requires no ordering in the interconnect and does not use negative acknowledgments (nacks) or retries, but it does queue requests at the directory controller in some cases. The base system stores the directory state in the main memory DRAM [17, 32, 34], but we also evaluate systems with “perfect” directory caches by simulating a zero cycle directory access latency.

Hammer. We use a reverse-engineered *approximation* of AMD’s Hammer protocol [5] to represent a class of recent systems whose protocols are not described in the academic literature (*e.g.*, Intel’s E8870 Scalability Port [7], IBM’s

Power4 [42] and xSeries Summit [10] systems). The protocol targets small systems (where broadcast is acceptable) with unordered interconnects (where traditional snooping is not possible), while avoiding directory state overhead and directory access latency. In this protocol, a processor first sends its request to a home node to be queued behind other requests to the same block. In parallel with the memory access, the home node broadcasts the request to all nodes who each respond to the requester with data or an acknowledgment. Finally, the requester sends a message to unblock the home node. By avoiding a directory lookup, this protocol has lower latency for cache-to-cache misses than a standard directory protocol, but it still requires indirection through the home node.

5.2 Interconnection Networks

We selected two interconnects with high-speed point-to-point links: an ordered “virtual bus” pipelined broadcast tree (sufficient for traditional snooping) and an unordered torus. We do not consider shared-wire (multi-drop) buses, because designing high-speed buses is increasingly difficult due to electrical issues [13, section 3.4.1]. We selected the link bandwidth of 3.2 GBytes/sec (4-byte wide links at 800 Mhz) and latency of 15 ns based on descriptions of current systems (e.g., the Alpha 21364 [32] and AMD’s hammer [5]). Messages are multiplexed over a single shared interconnect using virtual networks and channels, and broadcast messages use bandwidth-efficient tree-based multicast routing [14, section 5.5].

Tree (Figure 1a). For our totally-ordered interconnect, we use a two-level hierarchy of switches to form a pipelined broadcast tree with a fan-out of four, resulting in a message latency of four link crossings. This tree obtains the total order required for traditional snooping by using a single switch at the root. To reduce the number of pins per switch, a 16-processor system using this topology has nine switches (four incoming switches, four outgoing switches, and a single root switch).

Torus (Figure 1b). For our unordered interconnect, we use a two-dimensional, bidirectional torus like that used in the Alpha 21364 [32]. A torus has reasonable latency and bisection bandwidth, especially for small to mid-sized systems. For 16-processor systems, this interconnect has an average message latency of two link crossings.

5.3 Simulation Methods

We simulate our target systems with the Simics full-system multiprocessor simulator [26], and we extend Simics with a processor and memory hierarchy model to compute execution times [6]. Simics is a system-level architectural simulator developed by Virtutech AB that can run unmodified commercial applications and operating systems. Simics is a functional simulator only, but it provides an interface to support our detailed timing simulation. We use TFSim [30], configured as described in Table 1, to model

Table 2. Overhead due to reissued requests

Workload	Percentage of Misses			
	Not Reissued	Reissued Once	Reissued > Once	Persistent Requests
Apache	95.75%	3.25%	0.71%	0.29%
OLTP	97.57%	1.79%	0.43%	0.21%
SPECjbb	97.60%	2.03%	0.30%	0.07%
Average	96.97%	2.36%	0.48%	0.19%

superscalar processor cores that are dynamically scheduled, exploit speculative execution, and generate multiple outstanding coherence requests. Our detailed memory hierarchy simulator models the latency and bandwidth of the interconnects described above, and it also captures timing races and all state transitions (including non-stable states) of the coherence protocols. All workloads were warmed up and checkpointed to avoid system cold-start effects, and we ensure that caches are warm by restoring the cache contents captured as part of our checkpoint creation process. To address the variability in commercial workloads, we adopt the approach of simulating each design point multiple times with small, pseudo-random perturbations of request latencies to cause alternative operating system scheduling paths in our otherwise deterministic simulations [6]. Error bars in our runtime results represent one standard deviation from the mean in each direction.

6 Evaluation via Five Questions

We present evidence that Token Coherence can improve performance via five questions.

Question #1: Can the number of reissued and persistent requests be small? Answer: Yes; on average for our workloads, 97% of *TokenB*’s cache misses are issued only once. Since reissued requests are slower and consume more bandwidth than misses that succeed on the first attempt, reissued requests must be uncommon for *TokenB* to perform well. Races are rare in our workloads, because—even though synchronization and sharing are common—multiple processors rarely access the same data simultaneously due to the large amount of shared data. Table 2 shows the percentage of all *TokenB* misses that are not reissued, reissued once, reissued more than once, and that eventually use persistent requests. For our workloads, on average only 3.0% of cache misses are issued more than once and only 0.2% resort to persistent requests. (Table 2 shows *Torus* interconnect results, but *Tree* results, not shown, are similar.)

Question #2: Can *TokenB* outperform Snooping? Answer: Yes; with the same interconnect, *TokenB* and *Snooping* perform similarly for our workloads; however, by exploiting the lower-latency unordered *Torus*, *TokenB* on the *Torus* is faster than *Snooping* on the *Tree* intercon-

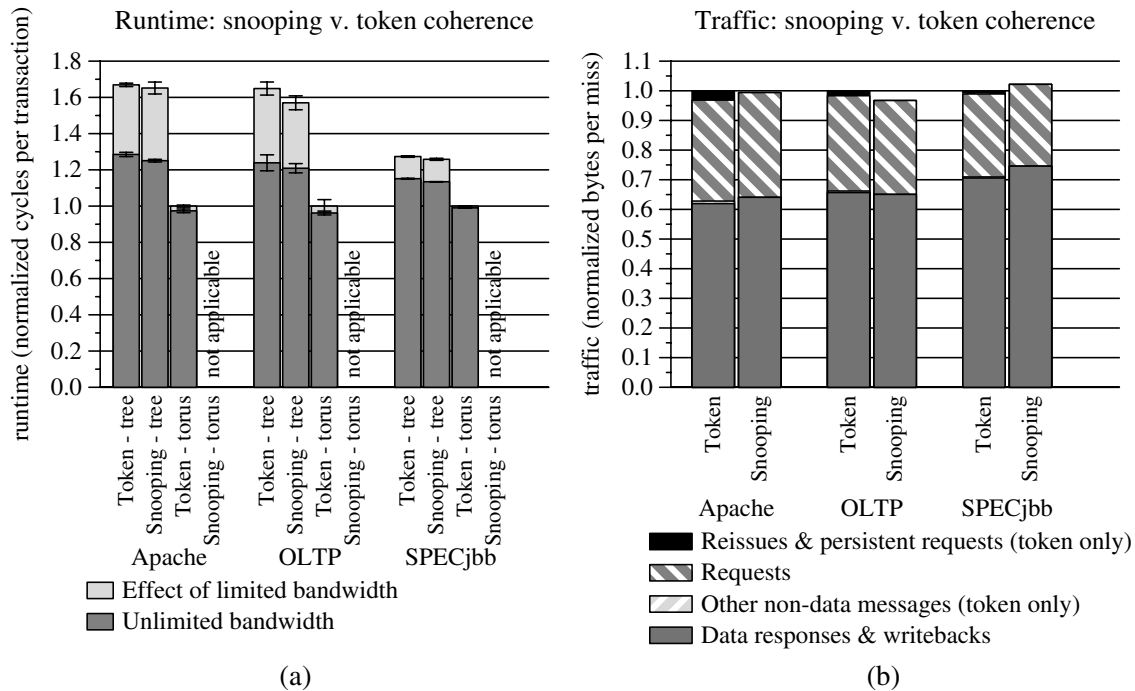


Figure 4. *Snooping v. TokenB*: runtime and traffic

nect (15-28% faster). Figure 4a shows the normalized runtime (smaller is better) of *TokenB* on the *Tree* and *Torus* interconnects and *Snooping* on the *Tree* interconnect. *Snooping* on the *Torus* is not applicable, because the *Torus* does not provide the required total order of requests. The dark grey bar shows the runtime when the bandwidth is changed from 3.2 GB/s to unlimited. Figure 4b shows the traffic in normalized average bytes per miss.

On the *Tree* interconnect, due to *TokenB*'s occasionally reissued requests, *Snooping* is slightly faster than *TokenB* (1-5% and 1-3%) with both limited and unlimited bandwidth, respectively (Figure 4a), and both protocols use approximately the same interconnect bandwidth (Figure 4b). However, since *Snooping* requires a totally-ordered interconnect, only *TokenB* can exploit a lower-latency unordered interconnect. Thus, by using the *Torus*, *TokenB* is 26-65% faster than *Snooping* on *Tree* with limited bandwidth links, and 15-28% faster with unlimited bandwidth links. This speedup results from (1) lower latency for all misses (cache-to-cache or otherwise) due to lower average interconnect latency, and (2) lower contention in *Torus* (by avoiding *Tree*'s central-root bottleneck).

Question #3: Can *TokenB* outperform *Directory* and *Hammer*? Answer: Yes; by removing the latency of indirection through the home node from the critical path of cache-to-cache misses, *TokenB* is faster than both *Directory* and *Hammer* (17-54% and 8-29% faster, respectively). Figure 5a shows the normalized runtime (smaller is better) for *TokenB*, *Hammer*, and *Directory* on the *Torus* interconnect with 3.2 GB/second links (the relative performances on *Tree*, not shown, are similar). The light grey

bars illustrate the small increase in runtime due to limited bandwidth in the interconnect. The grey striped bar for *Directory* illustrates the runtime increase due to the DRAM directory lookup latency.

TokenB is faster than *Directory* and *Hammer* by (1) avoiding the third interconnect traversal for cache-to-cache misses, (2) avoiding the directory lookup latency (*Directory* only), and (3) removing blocking states in the memory controller. Even if the directory lookup latency is reduced to zero (to approximate a fast SRAM directory or directory cache), shown by disregarding the grey striped bar in Figure 5a, *TokenB* is still faster than *Directory* by 6-18%. *Hammer* is 7-17% faster than *Directory* by avoiding the directory lookup latency (but not the third interconnect traversal), but *Directory* with the zero-cycle directory access latency is 2-9% faster than *Hammer* due to contention in the interconnect. The performance impact of *TokenB*'s additional traffic is negligible, because (1) the interconnect has sufficient bandwidth due to high-speed point-to-point links, and (2) the additional traffic of *TokenB* is moderate, discussed next.

Question #4: How does *TokenB*'s traffic compare to *Directory* and *Hammer*? Answer: *TokenB* generates less interconnect traffic than *Hammer*, but a moderate amount more than *Directory* (*Hammer* uses 79-90% more traffic than *TokenB*, *Directory* uses 21-25% less than *TokenB*). Figure 5b shows a traffic breakdown in normalized bytes per miss (smaller is better) for *TokenB*, *Hammer*, and *Directory*. The extra traffic of *TokenB* over *Directory* is not as large as one might expect, because (1) both protocols send a similar number of 72-byte data messages (81% of

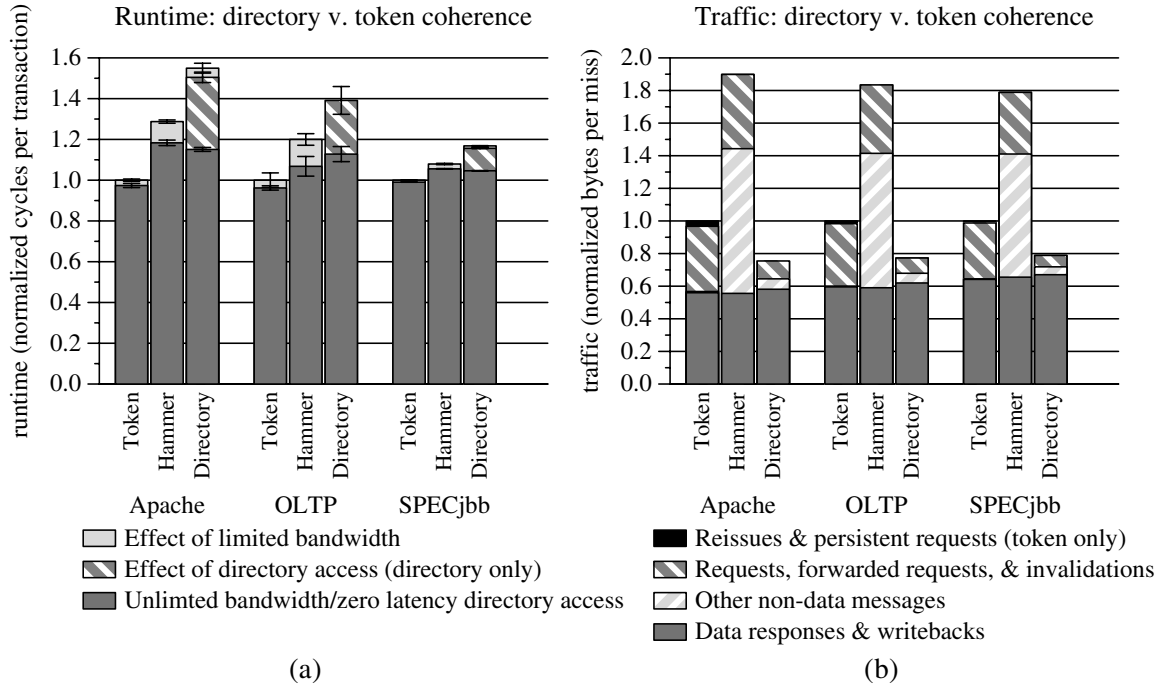


Figure 5. *Directory and Hammer v. TokenB: runtime and traffic*

Directory's traffic on average), (2) request messages are small (8 bytes), and (3) *Torus* supports broadcast tree routing (as stated in Section 5.2). *Hammer*, which targets smaller systems, uses much more bandwidth than *TokenB* or *Directory*, because every processor acknowledges each request (shown by the light grey striped segment).

Question #5: Can the *TokenB* protocol scale to an unlimited number of processors? Answer: No; *TokenB* relies on broadcast, limiting its scalability. However, *Token Coherence* is not limited to always broadcasting. *TokenB* is more scalable than *Hammer*, because *Hammer* uses broadcast and many acknowledgment messages. *TokenB* is less scalable than *Directory*, because *Directory* avoids broadcast. However, *TokenB* can perform well for perhaps 32 or 64 processors if bandwidth is abundant (by using high-bandwidth links [20] and coherence controllers with high throughput [33, 36] and low power consumption [31]). Experiments (not shown) using a simple micro-benchmark indicate that, for a 64 processor system, *TokenB* uses twice the interconnect bandwidth of *Directory*⁴. However, *TokenB* is a poor choice for larger or more bandwidth-limited systems. For this reason, the next section discusses other potential performance protocols.

7 Other Performance Protocol Opportunities

Token Coherence enables many performance protocols beyond the broadcast-always *TokenB* protocol. Furthermore, since its correctness substrate guarantees safety and prevents starvation, performance protocol designers can innovate without fear of corner-case correctness errors.

4. The additional cost of tree-based broadcast on *Torus* grows as $\Theta(\sqrt{n})$.

Reducing traffic. We can reduce request traffic (by not broadcasting transient requests) in several ways. First, we can reduce the traffic to directory protocol-like amounts by constructing a directory-like performance protocol. Processors first send transient requests to the home node, and the home redirects the request to likely sharers and/or the owner by using a “soft state” directory [25]. Second, bandwidth-adaptive techniques would allow a system to dynamically adapt between *TokenB* and this directory-like mode, providing high performance for multiple system sizes and workloads [29]. Third, *Token Coherence* can use destination-set prediction [2, 3, 9, 27] to achieve the performance of broadcast while using less bandwidth by predicting a subset of processors to which to send requests. Previously, these proposals required complicated protocols or protocol extensions. By multicasting transient requests, *Token Coherence* provides a simpler implementation of these proposals, while eliminating the totally-ordered interconnect required by some proposals [9, 27] and complex races in other proposals [2, 3, 9, 27].

Predictive push. The decoupling of correctness and performance provides an opportunity to reduce the number of cache misses by predictively pushing data between system components. This predictive transfer of data can be triggered by a coherence protocol predictor [1, 21, 35], by software (e.g., the KSR1’s “poststore” [37] and DASH’s “deliver” [24]), or by allowing the memory to push data into processor caches. Since *Token Coherence* allows data and tokens to be transferred between system components without affecting correctness, these schemes are easily implemented correctly as part of a performance protocol.

Hierarchical system support. Token Coherence can also accelerate hierarchical systems, an increasingly important concern with the rise of chip multiprocessors (CMPs, *e.g.*, IBM's Power4 [42]). Power4 uses extra protocol states to allow neighboring processors to respond with data, reducing traffic and average miss latency. A Token Coherence performance protocol could achieve this more simply by granting extra tokens to requesters, and allowing those processors to respond with data and tokens to neighboring processors. Other hierarchical systems connect smaller snooping based modules into larger systems (*e.g.*, [7, 10, 24]). Token Coherence may allow for a single protocol to more simply achieve the latency and bandwidth characteristics of these hierarchical systems, without requiring the complexity of two distinct protocols and the bridge logic between them.

8 Related Work

Protocol and interconnect co-design. Timestamp Snooping [28] adds ordering sufficient for traditional snooping to an unordered interconnect by using timestamps and reordering requests at the interconnect end points. Other approaches eschew virtual buses by using rings or a hierarchy of rings [15, 37, 42] or race-free interconnects [22]. Token Coherence addresses similar problems, but instead uses a new coherence protocol on an unordered interconnect to remove indirection in the common case.

Coherence protocols. Acacio *et al.* separately target read-miss latency [2] and write-miss latency [3] by augmenting a directory protocol with support for predicting current holders of the block. In many cases, the system grants permission without indirection, but in other cases, prediction is not allowed, requiring normal directory-based request ordering and directory indirection. In contrast, we introduce a simpler, unified approach that allows for correct direct communication in all cases, only resorting to a slower mechanism for starvation prevention. Shen *et al.* [38] use term rewriting rules to create a coherence protocol that allows operations from any of several sub-protocols, forming a trivially-correct hybrid protocol. We similarly provide a correctness guarantee, but we use tokens to remove ordering overheads from the common case. The Dir₁SW directory protocol [19] keeps a count of sharers at the memory for detecting the correct use of check-in/check-out annotations, and Stenstrom [39] proposes systems that use limited directory state and state in caches to track sharers. Token Coherence uses similar state in the memory and caches to count tokens, but it uses tokens to enforce high-level invariants that avoid directory indirection in the common case. Li and Hudak [25] explore a protocol in which each node tracks a probable owner, allowing requests to quickly find the current owner of the line. This approach could be used to improve a performance protocol or a persistent request mechanism.

9 Conclusions

To enable low-latency cache-to-cache misses on unordered interconnects, this paper introduces Token Coherence. Token Coherence resolves protocol races without indirection or a totally-ordered interconnect by decoupling coherence into a correctness substrate and a performance protocol. The correctness substrate guarantees correct transfer and access to blocks by tracking tokens, and it prevents starvation using persistent requests. Free from the burden of correctness, the performance protocol directly requests blocks without concern for races. We introduced *TokenB*, a specific performance protocol based on broadcasting transient requests and reissuing requests when occasional races occur. *TokenB* can outperform traditional snooping by using low-latency, unordered interconnects. *TokenB* outperforms directory protocols by avoiding cache-to-cache miss indirections, while using only a moderate amount of additional bandwidth.

By decoupling performance and correctness, Token Coherence may be an appealing framework for attacking other multiprocessor performance and design complexity problems. Future performance protocols may reduce request bandwidth via destination-set prediction, reduce miss frequency via predictive push, and gracefully handle hierarchical systems. By using the substrate to ensure correctness, these optimizations can be implemented with little impact on system complexity.

Acknowledgments

We thank Virtutech AB, the Wisconsin Condor group, and the Wisconsin Computer Systems Lab for their help and support. We thank Alaa Alameldeen, Allan Baum, Adam Butts, Joel Emer, Kouros Gharachorloo, Anders Landin, Alvin Lebeck, Carl Mauer, Kevin Moore, Shubu Mukherjee, Amir Roth, Dan Sorin, Craig Zilles, the Wisconsin Multifacet group, and the Wisconsin Computer Architecture Affiliates for their comments on this work.

References

- [1] H. Abdel-Shafi, J. Hall, S. V. Adve, and V. S. Adve. An Evaluation of Fine-Grain Producer-Initiated Communication in Cache-Coherent Multiprocessors. In *Proceedings of the Third IEEE Symposium on High-Performance Computer Architecture*, Feb. 1997.
- [2] M. E. Acacio, J. González, J. M. García, and J. Duato. Owner Prediction for Accelerating Cache-to-Cache Transfers in a cc-NUMA Architecture. In *Proceedings of SC2002*, Nov. 2002.
- [3] M. E. Acacio, J. González, J. M. García, and J. Duato. The Use of Prediction for Accelerating Upgrade Misses in cc-NUMA Multiprocessors. In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, pages 155–164, Sept. 2002.
- [4] S. V. Adve and K. Gharachorloo. Shared Memory Consistency Models: A Tutorial. *IEEE Computer*, 29(12):66–76, Dec. 1996.
- [5] A. Ahmed, P. Conway, B. Hughes, and F. Weber. AMD Opteron Shared Memory MP Systems. In *Proceedings of the 14th HotChips Symposium*, Aug. 2002. http://www.hotchips.org/archive/hc14/program/28_AMD_Hammer_MP_HC_v8.pdf.

- [6] A. R. Alameldeen, M. M. K. Martin, C. J. Mauer, K. E. Moore, M. Xu, D. J. Sorin, M. D. Hill, and D. A. Wood. Simulating a \$2M Commercial Server on a \$2K PC. *IEEE Computer*, 36(2):50–57, Feb. 2003.
- [7] M. Azimi, F. Briggs, M. Cekleov, M. Khare, A. Kumar, and L. P. Looi. Scalability Port: A Coherent Interface for Shared Memory Multiprocessors. In *Proceedings of the 10th Hot Interconnects Symposium*, pages 65–70, Aug. 2002.
- [8] L. A. Barroso, K. Gharachorloo, and E. Bugnion. Memory System Characterization of Commercial Workloads. In *Proceedings of the 25th Annual International Symposium on Computer Architecture*, pages 3–14, June 1998.
- [9] E. E. Bilir, R. M. Dickson, Y. Hu, M. Plakal, D. J. Sorin, M. D. Hill, and D. A. Wood. Multicast Snooping: A New Coherence Method Using a Multicast Address Network. In *Proceedings of the 26th Annual International Symposium on Computer Architecture*, pages 294–304, May 1999.
- [10] J. M. Borkenhagen, R. D. Hoover, and K. M. Valk. EXA Cache/Scalability Controllers. In *IBM Enterprise X-Architecture Technology: Reaching the Summit*, pages 37–50. International Business Machines, 2002.
- [11] A. Charlesworth. Starfire: Extending the SMP Envelope. *IEEE Micro*, 18(1):39–49, Jan/Feb 1998.
- [12] A. L. Cox and R. J. Fowler. Adaptive Cache Coherency for Detecting Migratory Shared Data. In *Proceedings of the 20th Annual International Symposium on Computer Architecture*, May 1993.
- [13] W. J. Dally and J. W. Poulton. *Digital Systems Engineering*. Cambridge University Press, 1998.
- [14] J. Duato, S. Yalamanchili, and L. Ni. *Interconnection Networks: An Engineering Approach*. Morgan Kaufmann, revised edition, 2003.
- [15] K. Farkas, Z. Vranesic, and M. Stumm. Scalable Cache Consistency for Hierarchically Structured Multiprocessors. *The Journal of Supercomputing*, 8(4), 1995.
- [16] S. J. Frank. Tightly Coupled Multiprocessor System Speeds Memory-access Times. *Electronics*, 57(1):164–169, Jan. 1984.
- [17] K. Gharachorloo, L. A. Barroso, and A. Nowatzyk. Efficient ECC-Based Directory Implementations for Scalable Multiprocessors. In *Proceedings of the 12th Symposium on Computer Architecture and High-Performance Computing (SBAC-PAD 2000)*, Oct. 2000.
- [18] K. Gharachorloo, A. Gupta, and J. Hennessy. Two Techniques to Enhance the Performance of Memory Consistency Models. In *Proceedings of the International Conference on Parallel Processing*, volume I, pages 355–364, Aug. 1991.
- [19] M. D. Hill, J. R. Larus, S. K. Reinhardt, and D. A. Wood. Cooperative Shared Memory: Software and Hardware for Scalable Multiprocessor. *ACM Transactions on Computer Systems*, 11(4):300–318, Nov. 1993.
- [20] M. Horowitz, C.-K. K. Yang, and S. Sidiropoulos. High-Speed Electrical Signaling: Overview and Limitations. *IEEE Micro*, 18(1), January/February 1998.
- [21] D. A. Koufaty, X. Chen, D. K. Poulsen, and J. Torrellas. Data Forwarding in Scalable Shared-Memory Multiprocessors. In *Proceedings of the 1995 International Conference on Supercomputing*, July 1995.
- [22] A. Landin, E. Hagersten, and S. Haridi. Race-Free Interconnection Networks and Multiprocessor Consistency. In *Proceedings of the 18th Annual International Symposium on Computer Architecture*, May 1991.
- [23] J. Laudon and D. Lenoski. The SGI Origin: A ccNUMA Highly Scalable Server. In *Proceedings of the 24th Annual International Symposium on Computer Architecture*, pages 241–251, June 1997.
- [24] D. Lenoski, J. Laudon, K. Gharachorloo, W.-D. Weber, A. Gupta, J. Hennessy, M. Horowitz, and M. Lam. The Stanford DASH Multiprocessor. *IEEE Computer*, 25(3):63–79, Mar. 1992.
- [25] K. Li and P. Hudak. Memory Coherence in Shared Virtual Memory Systems. *ACM Transactions on Computer Systems*, 7(4):321–359, 1989.
- [26] P. S. Magnusson et al. Simics: A Full System Simulation Platform. *IEEE Computer*, 35(2):50–58, Feb. 2002.
- [27] M. M. K. Martin, P. J. Harper, D. J. Sorin, M. D. Hill, and D. A. Wood. Using Destination-Set Prediction to Improve the Latency/Bandwidth Tradeoff in Shared Memory Multiprocessors. In *Proceedings of the 30th Annual International Symposium on Computer Architecture*, June 2003.
- [28] M. M. K. Martin, D. J. Sorin, A. Ailamaki, A. R. Alameldeen, R. M. Dickson, C. J. Mauer, K. E. Moore, M. Plakal, M. D. Hill, and D. A. Wood. Timestamp Snooping: An Approach for Extending SMPs. In *Proceedings of the Ninth International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 25–36, Nov. 2000.
- [29] M. M. K. Martin, D. J. Sorin, M. D. Hill, and D. A. Wood. Bandwidth Adaptive Snooping. In *Proceedings of the Eighth IEEE Symposium on High-Performance Computer Architecture*, Feb. 2002.
- [30] C. J. Mauer, M. D. Hill, and D. A. Wood. Full System Timing-First Simulation. In *Proceedings of the 2002 ACM Sigmetrics Conference on Measurement and Modeling of Computer Systems*, pages 108–116, June 2002.
- [31] A. Moshovos, G. Memik, B. Falsafi, and A. Choudhary. JETTY: Filtering Snoops for Reduced Power Consumption in SMP Servers. In *Proceedings of the Seventh IEEE Symposium on High-Performance Computer Architecture*, Jan. 2001.
- [32] S. S. Mukherjee, P. Bannon, S. Lang, A. Spink, and D. Webb. The Alpha 21364 Network Architecture. In *Proceedings of the 9th Hot Interconnects Symposium*, Aug. 2001.
- [33] A. K. Nanda, A.-T. Nguyen, M. M. Michael, and D. J. Joseph. High-Throughput Coherence Controllers. In *Proceedings of the Sixth IEEE Symposium on High-Performance Computer Architecture*, Jan. 2000.
- [34] A. Nowatzyk, G. Aybay, M. Browne, E. Kelly, and M. Parkin. The S3.mp Scalable Shared Memory Multiprocessor. In *Proceedings of the International Conference on Parallel Processing*, volume I, pages 1–10, Aug. 1995.
- [35] D. Poulsen and P.-C. Yew. Data Prefetching and Data Forwarding in Shared-Memory Multiprocessors. In *Proceedings of the International Conference on Parallel Processing*, volume II, pages 296–280, Aug. 1994.
- [36] I. Pragaspathy and B. Falsafi. Address Partitioning in DSM Clusters with Parallel Coherence Controllers. In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, Oct. 2000.
- [37] E. Rosti, E. Smirni, T. Wagner, A. Apon, and L. Dowdy. The KSR1: Experimentation and Modeling of Poststore. In *Proceedings of the 1993 ACM Sigmetrics Conference on Measurement and Modeling of Computer Systems*, pages 74–85, May 1993.
- [38] X. Shen, Arvind, and L. Rudolph. CACHET: An Adaptive Cache Coherence Protocol for Distributed Shared-Memory Systems. In *Proceedings of the 1999 International Conference on Supercomputing*, pages 135–144, June 1998.
- [39] P. Stenström. A Cache Consistency Protocol for Multiprocessors with Multistage Networks. In *Proceedings of the 16th Annual International Symposium on Computer Architecture*, May 1989.
- [40] P. Stenström, M. Brorsson, and L. Sandberg. Adaptive Cache Coherence Protocol Optimized for Migratory Sharing. In *Proceedings of the 20th Annual International Symposium on Computer Architecture*, pages 109–118, May 1993.
- [41] P. Sweazey and A. J. Smith. A Class of Compatible Cache Consistency Protocols and their Support by the IEEE Futurebus. In *Proceedings of the 13th Annual International Symposium on Computer Architecture*, pages 414–423, June 1986.
- [42] J. M. Tandler, S. Dodson, S. Fields, H. Le, and B. Sinharoy. POWER4 System Microarchitecture. IBM Server Group Whitepaper, Oct. 2001.
- [43] K. C. Yeager. The MIPS R10000 Superscalar Microprocessor. *IEEE Micro*, 16(2):28–40, Apr. 1996.