

Multicore-Aware Reuse Distance Analysis *

Derek L. Schuff, Benjamin S. Parsons, and Vijay S. Pai
Purdue University
West Lafayette, IN 47907

E-mail: dschuff@purdue.edu, bspardon@purdue.edu, vpai@purdue.edu

Abstract

This paper presents and validates methods to extend reuse distance analysis of application locality characteristics to shared-memory multicore platforms by accounting for invalidation-based cache-coherence and inter-core cache sharing. Existing reuse distance analysis methods track the number of distinct addresses referenced between reuses of the same address by a given thread, but do not model the effects of data references by other threads. This paper shows several methods to keep reuse stacks consistent so that they account for invalidations and cache sharing, either as references arise in a simulated execution or at synchronization points. These methods are evaluated against a Simics-based coherent cache simulator running several OpenMP and transaction-based benchmarks. The results show that adding multicore-awareness substantially improves the ability of reuse distance analysis to model cache behavior, reducing the error in miss ratio prediction (relative to cache simulation for a specific cache size) by an average of 69% for per-core caches and an average of 84% for shared caches.

1 Introduction

Reuse distance analysis (also called stack distance analysis) is a method for characterizing application memory access behavior without being tied to a specific memory hierarchy [22]. The reuse distance of a given reference to address x is the count of the number of distinct addresses that have been referenced since the last access of address x (or ∞ if the address has not been referenced before). Conceptually, this is implemented by tracking the depth of the access in a stack, though more efficient implementations also exist [3, 5, 23, 30]. Forming a cumulative histogram of reuse distances gives a prediction of cache hit ratio; in particular, the hit ratio of a cache with N one-word blocks organized in a single LRU set would correspond to the fraction of references with reuse distance less than N . A single run of reuse distance analysis can thus predict the behavior of an application over a variety of cache sizes, but does not account for real cache constraints like associativity, block size, or replacement policy details. Despite those limitations, reuse distance analysis has been used and validated for predicting cache performance (even for low-associativity caches), modeling program locality, providing cache hints during code generation, and for feeding models that restructure code and data to improve locality [7, 8, 13, 14, 20, 21, 27, 32].

When considering multicore systems, however, existing reuse distance analysis methods are insufficient, as their abstract notion of access ordering by counting only makes sense within a single reference stream. Considering events within only a single reference stream ignores inter-core data communication that may impact or even dominate parallel program performance. These inter-core interactions depend largely on the cache configuration, which may consist of private caches or

*This work is supported in part by the National Science Foundation under Grant Nos. CCF-0532448, CNS-0532452, and CCF-0621457.

shared caches. Private caches are typically kept coherent using invalidations; shared caches can allow fast inter-thread communication when more than one core sharing the cache accesses the same block (inter-core prefetching), and can also use cache capacity more efficiently by reducing the number of replicated copies of widely-read data. Reuse distances are supposed to represent a measure of locality, with shorter distances being more likely to hit and longer ones less likely. However, if one thread has written to the same address between two reuses of that address by another thread with a different cache, an invalidation will take place and the access will miss regardless of how short the reuse distance is. Similarly, if a thread fetches an address into a shared cache, another thread will see that access as a hit even if it has never referenced that data before. Thus, reuse distance alone does not consistently correspond to locality in multicore systems.

There is thus a fundamental tradeoff between completely machine-independent metrics and those that include details of the machine such as the organization of the memory hierarchy. Machine independent metrics are advantageous in that they require less information about machine parameters, and data from a single run may be more broadly applicable, but at the cost of some accuracy. Although reuse distance is intended to be largely machine-independent, applying it to multithreaded programs requires striking a balance: adding enough machine awareness to accurately capture locality characteristics but not so much as to diminish its simplicity and broad applicability.

This paper enables reuse distance analysis for parallel systems by modeling the cache sharing configuration of the multicore processors used. Although the analysis methods are independent of cache sizes, they do require an understanding of which (if any) cores share a cache. In practice, real systems include private L1 caches even if they have a shared L2 or L3; this paper only considers the configuration of the outermost level of on-chip cache since that determines off-chip bandwidth consumption and typically has the most pronounced impact on performance [10]. Systems with per-core caches use private reuse distance stacks, but must guarantee that addresses are removed from the reuse distance stack when they would be invalidated in a real system; a reuse after an invalidation is just as “non-local” as a cold miss, and is thus treated as such. Since the precise timing of invalidations varies from system to system, the system here considers different points to remove an address from reuse distance stacks in the system as a result of a write by one thread: immediately at the time of the write, lazily at the following synchronization point, or at the preceding synchronization point (in an oracular fashion). The lazy option represents the longest delay option for those invalidations, while the oracular option represents the earliest possible time for invalidations to propagate and for those blocks to become free for use by other data.

For systems with shared caches, the most important issue is to make sure that accesses by one thread become promptly available in cache for other threads to access. This is done by using a shared reuse distance stack that sees all references from all threads and measures reuse distances in terms of total number of distinct addresses referenced by all threads since the last reuse by any thread. This paper also considers caches shared by a subset of the cores, such as pairwise-shared caches in a quad-core processor (e.g., the Intel Core 2 Quad Processor [15]).

To test these models, all analyses described are implemented in a Simics fast functional execution model that switches

between cores every instruction, although the basic ideas would also apply to trace-driven or detailed execution-driven simulation. Both private and shared-cache models are evaluated using applications coded in the OpenMP shared-memory API and the TL2 software transactional memory system [12, 24]. Although the analyses described here apply to any model of parallelism, the lazy and oracular invalidation methods depend on explicit identification of synchronization points such as those used in data-race-free programs [1]. This paper’s framework currently implements those methods only for the barriers at the end of OpenMP parallel regions.

The experimental results show that adding multicore-awareness to reuse distance analysis substantially improves locality modeling when compared to various cache configurations for a simulated quad-core system. When comparing the predictions from reuse distance analysis against actual cache simulation results for 1 MB, 2 MB, and 4 MB total cache sizes split across 4 per-core caches, adding multicore-awareness reduces the error in miss ratio prediction by an average of 69% across 13 applications. This result arises because the unaware method does not account for the fact that a reuse will not take place if another thread has issued an intervening write. Comparing against simulation results for shared caches of the same total cache size, multicore-awareness reduces the errors in predicted miss ratios by an average of 84%. This result arises because the unaware method does not account for the positive effects of inter-thread reuse when considering limited cache sizes.

2 Background on Reuse Distance

The reuse distance seen by a memory reference is defined as the number of distinct data elements or memory locations referenced since the last reference to the same location. The data elements can be pages, cache lines, individual words, or even instructions, yielding a generalized machine-independent measure of program locality. A complete histogram of reuse distances can be used to predict hit ratios for any size cache; in particular the hit ratio of a fully associative cache with N elements using LRU replacement is the fraction of references with a reuse distance less than N . Reuse distance thus enables simulation of multiple cache configurations in one simulation pass. Beyls and D’Hollander showed that reuse distance prediction is accurate even for caches of lower associativity [7].

Reuse distance analysis was first introduced by Mattson et al. in 1970 in the context of virtual memory pages as “LRU stack processing”, a way to model different storage configurations in one pass [22]. The measurement is implemented by keeping a record of all previously-accessed addresses in a stack. On each reference, the stack is searched for the address and the distance is recorded as the number of addresses above it. The address is then moved to the top of the stack. Mattson used a list-based stack which required $O(NM)$ time, where N is the length of execution and M is the size of the program data. Since then there have been many improvements in reuse distance implementations, such as representing the trace as a vector or using a tree instead of a list to implement the stack [3, 5, 23, 30]. Approximate methods have also been developed, which trade accuracy for space efficiency. Ding and Zhong introduced dynamic tree compression to reduce space requirements to $O(\log M)$ while bounding relative or absolute error as desired [13].

3 Multicore-Aware Reuse Distance Models

Reuse distance is a model of program locality, which determines cache performance. Its fundamental limitation with respect to multicore systems is that cache performance for a given thread is affected by the actions of other threads. Because these interactions can have a significant impact on performance, they should not be ignored by models that predict performance. Since intra-thread locality is still the primary factor affecting cache performance, the basic model of reuse distance is the same; however, it must be extended to account for inter-thread interactions because behaviors such as coherence and sharing are fundamental elements in multicore systems. The modifications to the reuse distance model depend on the cache configuration used in a multicore system. Caches may be private to each processor core or shared among two or more cores; each policy has a different effect on reuse distance since it has a different impact on coherence and sharing. On the other hand, issues such as cache size or associativity do not have any additional impact on multicore systems, so the existing reuse distance models may remain independent of these constraints.

The discussion below makes several assumptions. First, the threads in the program of interest are assumed to be compute-bound and the only load on the system. This assumption allows reuse distance analysis to ignore the effects of context switches, OS behavior, and timesharing, and thereby follows the traditional methods for reuse distance analysis. This assumption is valid for the workloads studied in this paper. Second, this work assumes that threads are tied to specific cores so that thread migration has no impact; this is likely to be the case as most real systems support some form of affinity scheduling [17, 31]. Third, this work treats the cache as though it were just a single-level corresponding to the last level of cache (L2 or L3 today). This assumption allows this work to consider only fast hits and slow misses, rather than hits at multiple levels of cache. This assumption is not required by reuse distance analysis, as it models locality abstractly and can simultaneously model multiple levels of cache [9]. Rather, it is just a simplification for providing more concise metrics. Additionally, this is based on the observation that misses from the last level are generally the most critical in shaping performance [10].

3.1 Modeling private caches

When caches are private to each processor core in a shared-memory system, a coherence protocol is employed to prevent caches from containing stale data. Most coherence protocols are write-invalidate, in which a write to a cache block in one cache invalidates copies of that block in other caches. If a subsequent reference is made to the recently invalidated block, it will miss in the cache no matter how recently the core accessed it. (This is called a coherence miss.)

Per-core caches with invalidation-based coherence can be modeled by using per-thread reuse distance stacks for which a write to any address removes that address from all other stacks which contain it, leaving it only in the stack belonging to the thread which wrote it. Then a subsequent reference to that address by some other thread will have infinite reuse distance, as though it had never been seen before — just like a cold miss. When a block is invalidated from a cache, it leaves behind a free slot that can be re-used for the next block of data that maps to that set. Blocks that have been evicted will not go back into

the cache, but the slot can be filled without evicting any other blocks. In a fully-associative cache, any block can be mapped to this slot. To model this behavior, blocks that are invalidated are left in the stack but marked as invalid, so an access deeper in the stack has the same reuse distance. Then when the deeper block is brought to the top, the slot is removed (filled in).

Invalidation can be performed on every write, as with real caches. However the order and interleaving of references varies with the machine characteristics and the particular execution of the program. In a data-race-free (DRF) program, any data invalidated from thread t 's cache has been written by some other thread; thus t will not access it until at least the next synchronization point. Similarly, any data referenced by t will not be written by other threads, and thus will not be invalidated until the next synchronization. Thus, invalidations can be propagated anytime between the previous synchronization and the next one while still representing a possible execution of the DRF program.

This paper investigates 3 different timings for the invalidations. The first, *eager*, mimics the behavior of a cache, performing invalidations for each write reference as soon as it occurs during simulation. In the second, *lazy* invalidation, each thread accesses its own reuse stack as normal, but write references are also buffered. All invalidations caused by writes in the buffer take place at the next synchronization point. The third, *oracular* invalidation, uses future knowledge of the write references for the upcoming parallel region and invalidates the addresses that will be written before the region begins. This is implemented by buffering all references between synchronization points, then performing all invalidations related to the buffered references, and then actually accessing the thread's private reuse stack for the buffered references. The lazy and oracular schemes test the effect of invalidations freeing space in the cache to be used by other data. Since the DRF property guarantees that a thread will not reference an address in the same parallel region in which it is invalidated, these addresses just take up space until they are invalidated and removed. The lazy invalidation scheme leaves them in the stack as long as possible, which should increase the overall reuse distance and thus represent a worst case. The oracular invalidation scheme removes these addresses as soon as possible, thus freeing up their space early and reducing overall reuse distance; it represents a best case. Both lazy and oracular are currently implemented only for OpenMP fork-join synchronization; other benchmarks are tested using only eager. These approaches, could, however, be extended to other DRF programs with well-identified synchronization points.

Accounting for false sharing. False sharing is an additional concern in write-invalidate protocols for caches with larger block sizes. Although most tests in this study use a single-word (8-byte) cache block size, some tests also consider a larger and more realistic block size of 64 bytes. False sharing is dependent on machine configuration and timing, but can be explicitly detected by extending the reuse distance analysis in the following way: each reuse stack entry (representing one block accessed by a particular thread) is extended to include information on which words in the block have actually been accessed. If the block is subsequently invalidated by another thread this information can be compared to the word causing the invalidation to determine whether true or false sharing has occurred.

3.2 Modeling shared caches

When threads run on cores that share a cache, the locality of each thread is affected by the behavior of the others in several ways. First, one thread accessing data will effectively prefetch that data for all others with the same shared cache. Second, all threads sharing a cache can use just one copy of a given widely-read data element, thus reducing unnecessary replication and using the cache capacity more efficiently. Third, different threads may have different working set sizes and thus require different partitions of total cache capacity, possibly freeing up space for others, or taking space away from them.

Reuse distance can be measured across all references from all threads; in this case the distance seen by one thread's references are affected by the data use patterns of the other threads. This type of analysis can be enabled by directly using a shared reuse stack.

Additionally, caches may have hierarchical sharing. For example, a 4-core system could have 2 caches, each shared by 2 cores, which use a coherence protocol between them (e.g., the Intel Core 2 Quad Processor [15]). In this case, any of the private techniques could be used across caches.

4 Experimental Methodology

Simulation Platform. The reuse distance stack implementation is based on Sugumar and Abraham's splay tree code as distributed with SimpleScalar 4.0 [18, 30]. The reuse distance stacks were implemented as a memory hierarchy module in the Virtutech Simics full-system simulator, which executes all code from the Linux kernel and system daemons [19]. However, only userspace references from the benchmark were fed to the stacks. In addition the stacks were tied to the thread rather than the processor to eliminate any effects of thread migration. Timing of the instructions and memory references was not modeled; Simics' default timing model of one instruction per cycle was used for most tests, with round-robin switching among processors on each cycle. Some additional tests were also performed with different switching frequencies to determine if the results were sensitive to that particular assumption.

Workloads. This system is evaluated with thirteen benchmarks: six from SPEC OMP2001, three from the OpenMP implementation of the NAS Parallel Benchmarks version 3.3, and four from the STAMP transactional memory benchmark suite [4, 11, 16]. They include `312.swim`, `316.applu`, `318.galgel`, `320.equake`, `324.apsi`, and `326.gafort` from SPEC OMP; `CG.W`, `FT.W`, and `MG.W` from NAS; and `genome`, `intruder`, `kmeans`, and `vacation` from STAMP. All benchmarks were run with 4 threads. The "test" data input was used for SPEC benchmarks and the "W" input size for NAS benchmarks. High-contention input parameters were used for `kmeans` and `vacation`. For the OpenMP benchmarks under the lazy and oracular invalidation models, the library itself was modified to notify the simulator at the end of each parallel region and thus allow the simulator to ignore the implementation of barriers between parallel regions. For the STAMP benchmarks, all transactions are implemented in software using the TL2-x86 STM system distributed with STAMP [12].

Validation against simulated caches. This study also validates the accuracy of reuse distance as a predictor of cache

performance by comparing against Simics’ coherent cache simulation module called *g-cache*. The caches operated on the same reference stream as the reuse distance stacks. The stacks were then used to predict hits and misses at each reference for 3 specific cache sizes by comparing the reuse distance for the reference with the size of the simulated cache. The base tests use caches with single-word (8 byte) lines and 8-way set associativity; as discussed in Section 3, some sensitivity tests also consider 64 byte cache block sizes.

The cache simulations model configurations with private per-thread caches with MSI coherence, with a single cache shared among all 4 threads, and with a pair of caches, each of which is shared by 2 threads, and which use MSI coherence between them. The private cache experiments use 256KB, 512KB and 1MB sizes. The shared cache experiments use 1MB, 2MB, and 4MB sizes, so as to correspond to the same total cache size used in the private cache tests. Similarly, the pairwise shared cache tests used 512KB, 1MB, and 2MB sizes.

5 Experimental Results

5.1 Private Cache Configurations

Reuse distance results. Figure 1 shows the cumulative reuse distance function for each application using unaware per-thread reuse distance analysis without considering multicore effects and using the eager, lazy, and oracular multicore-aware reuse distance analysis methods discussed in Section 3. The X axis represents a given reuse distance value x , while the Y axis represents the percentage of application data accesses that exhibit a reuse distance of x or less under each model. Because these applications are dominated by double-precision floating-point accesses, addresses are considered at 64-bit granularity. A particular (x, y) point on these curves thus predicts a cache hit ratio y for a fully-associative LRU cache of size x 64-bit words.

In most cases, there is little discernible difference in predicted hit ratio between unaware and multicore-aware reuse distance measurement until the right end where reuse distance is the longest, and the plots are zoomed to show these differences. The fact that the differences primarily appear for large reuse distances stems from traditional working set analysis. Reuse distance plots often have one or more knees and plateaus, which correspond to different working sets. These working sets can be seen clearly, for example, in the `equake` benchmark. In data-parallel applications the largest datasets would be the ones most commonly divided among the cores in the system, so any differences between traditional and multicore-aware reuse distance would be seen at this end of the plot. In addition, with smaller caches and lower overall hit ratios, capacity misses increase, causing coherence misses to become a smaller fraction of overall misses, mitigating their effects. The STAMP benchmarks do not have the clean data distribution or the well-defined (and often long) periods between synchronizations that the OpenMP benchmarks do, and write more frequently to actively shared data. Thus, a higher fraction of the misses are coherence misses even at small cache sizes, causing unaware and invalidating schemes to diverge in the plots even at small reuse distances.

The three multicore-aware reuse models generally track fairly close to each other, with the only truly substantial differences seen in `316.applu` for distances between 8000 and 18000. The difference between the lazy and oracular methods is a function of the amount of data invalidated from a particular stack during an OpenMP parallel region (which could not have been accessed by this stack during this region due to the DRF property) and how recently that invalidated data was used compared to the data that actually was accessed on this stack during this region: invalidations of data that was last used less recently than the data of this region will not impact the reuse distance for the data of this region since the invalidations occur deeper on the stack. The larger discrepancies in this application indicate that a large amount of recently-used data is being invalidated.

Impact on performance prediction. Comparing the multicore-aware models to the unaware model shows that for most applications, the multicore-aware cumulative reuse level reaches a peak at least a few percent below that of the unaware case. The multicore-aware versions never reach 100% reuse because coherence-related misses behave just like cold misses, with an infinite reuse distance. When seen on a plot as hit ratio variations of just a few or even fractions of a percent, these differences may seem insignificant. However the memory hierarchy is such an important factor in system performance because the difference in stall time between a hit and a miss is so great. Consequently, the difference between a 98% hit ratio (2% miss ratio) and a 99% hit ratio (1% miss ratio) leads to a performance impact far greater than 1%, and indeed often closer to doubling.

To illustrate the way these hit ratios affect execution times, Figure 2 shows a simple prediction of application execution time computed as the number of instructions plus a constant miss penalty for each cache miss at a given cache size, for a subset of the benchmarks. On a real system, the cost of a miss or other instruction would be affected by issues such as instruction-level parallelism, memory-level parallelism, data dependence depth, and bandwidth contention [25, 26, 29]. The figure considers several values of relative miss penalty to capture a range of reasonable possibilities in a real system. Since the three invalidation schemes yield similar results, only the eager invalidation model is shown.

The difference between the times for Unaware and Eager (solid and dotted lines of the same shade) become clearer here and are most evident among large cache sizes with low miss ratios, where coherence misses become a significant fraction of all misses. With the largest cache sizes and miss penalty, accounting for invalidations results in execution times up to nearly 9 times longer, with an average of 2.1 times over all miss penalties and all benchmarks studied (including those not shown).

Validating against simulated caches. Table 1 shows a comparison of miss ratios predicted by the various reuse distance analysis schemes with miss ratios determined by simulated caches using MSI coherence. Reuse distances for each thread are used to predict the miss ratio for each of the cache sizes on each of the threads, and then averaged together. Miss ratios are used instead of hit ratios because as illustrated by Figure 2, when hit ratios are high, the performance correlates better numerically with the miss ratio number than the hit ratio. The predicted and simulated cache miss ratios for cache sizes of

Benchmark	Percent error			
	Unaware	Eager	Lazy	Oracular
316.applu	91.55	6.58	10.30	8.78
324.apsi	82.31	4.76	5.82	5.82
CG.W	2.89	0.93	2.89	2.89
320.quake	1.58	1.65	1.65	1.65
FT.W	5.59	5.58	5.60	5.60
326.gafort	99.89	0.37	4.57	3.71
318.galgel	79.57	5.80	6.02	6.03
genome	8.37	3.68	N/A	N/A
intruder	20.15	0.14	N/A	N/A
kmeans	45.58	0.70	N/A	N/A
MG.W	4.62	4.66	4.65	4.66
312.swim	3.72	3.72	3.72	3.72
vacation	8.37	1.25	N/A	N/A
Average	34.94	3.06	5.04	4.76

Table 1: Percent error of each reuse distance method compared to simulated caches. As discussed in Section 3, Lazy and Oracular are currently only implemented for OpenMP benchmarks.

256 KB, 512 KB, and 1 MB are compared, and the reported prediction error for each application represents the average error when modeling these three cache sizes.

The prediction error for the unaware non-invalidating method varies from 1.58 for *quake* to more than 99% for *gafort*, with an average of 35%. The accuracy of Unaware correlates strongly with overall miss ratio, especially for the OpenMP benchmarks. For example, *applu*, *galgel*, *apsi*, and *gafort* have miss ratios ranging from 1.2-4.0% for the smallest cache size, and their error for Unaware is greater than 79%. The remaining OpenMP benchmarks have miss ratios ranging from 11-51%, with error of less than 5%. The STAMP benchmarks fall in between, with an average of 6.9% miss ratio and 21% error for Unaware. In these cases where the overall number of misses is small, the proportion of coherence misses, which are not accounted for by Unaware, is much higher and its accuracy is reduced. In contrast, the multicore-aware methods improve accuracy dramatically, with Eager, Lazy, and Oracular all seeing 7% or less inaccuracy for all benchmarks. The average reduction in error is 69% for Eager, and about 50% for Lazy and Oracular. Thus, it is important to add multicore-awareness to reuse distance analysis, but the precise method for doing so is not as important for these applications.

5.2 Shared Cache Configurations

Figure 3 shows the cumulative reuse distance function for each application using a single shared reuse distance stack (solid black line) and pairwise shared stacks that invalidate each other using the eager method (dotted line). The unaware reuse distance with no invalidation or sharing is shown for comparison. The lines are shifted for the unaware and pairwise shared stacks such that the x value represents the total cache size in all cases (4x the size of the unaware caches and 2x the size of the pairwise caches).

Benchmark	Percent error			
	Shared cache		Pairwise shared	
	Unaware	Aware	Unaware	Aware
316.applu	66.05	7.73	84.69	5.04
324.apsi	12478.76	0.09	68.08	6.09
CG.W	2.24	0.05	1.62	0.17
320.equake	3.43	0.80	2.91	1.81
FT.W	5.87	3.61	5.24	0.59
326.gafort	123.13	0.00	99.80	0.26
318.galgel	435.79	45.87	60.62	8.79
genome	23.20	0.88	7.34	1.88
intruder	28.85	0.34	7.99	0.26
kmeans	0.97	0.00	35.62	1.09
MG.W	5.43	4.04	5.78	5.93
312.swim	3.94	3.74	3.83	3.73
vacation	73.16	2.73	8.07	2.48
Average	1019.29	5.37	30.12	2.93

Table 2: Percent error of each reuse distance method compared to a simulated shared or pairwise shared cache configuration

When reuse distances are shortened by effects such as inter-core prefetching or reduced replication, the shared plots reflect higher fractions of accesses less than a given distance (and thus a higher predicted hit ratio for a cache of that size). For most benchmarks (especially the OpenMP benchmarks), this corresponding distinction between unaware and multicore-aware models is seen over a wider range of sizes than with the private caches. Though not easily seen on all plots, the line for the pairwise shared stacks dips below those of the shared stacks in the upper right corner and flattens out at a hit ratio less than 100% due to invalidations between the shared caches.

On seven benchmarks (all four STAMP benchmarks plus *CG.W*, *gafort*, and *MG.W*), the predicted hit ratio for the shared cache is actually lower than for private caches when using very small cache sizes such as 4 words. This illustrates on a small scale the ability of one thread to monopolize shared cache resources as mentioned in Section 3.2. In such a small cache, a run of 4 consecutive references from one thread can evict all the lines from the other threads, resulting in worse performance than if each thread had one line to itself.

Table 2 shows a comparison of miss ratios predicted by the reuse distance analysis schemes with miss ratios determined by simulated shared caches and pairwise-shared caches. For the unaware method, the predictions are made for individual caches of size $1/4$ the size of the shared cache or $1/2$ the size of the pairwise-shared caches, so that the total cache size is the same. Reuse distances for each thread are used to predict the miss ratio for each of the cache sizes on each of the threads (for the shared stack there is only one stack for all threads), and then averaged together.

Focusing first on the single shared cache, there are two main sources of error when attempting to model this with independent unaware reuse stacks. The first is that they will indicate capacity misses that will not be seen by a shared stack; this is due to duplication of data reducing the effective aggregate capacity, and to accesses by other cores that keep shared

blocks in the cache when they would have been evicted from a private cache. The second is that accurately modeling inter-core prefetching will also cause the shared stacks to report fewer cold misses. These effects represent benefits to using a shared-cache configuration.

Both of these effects are strongly intensified when the working set is small compared to the total size of the shared cache. This can be seen in the inaccuracy of the unaware method on the `apsi`, `galgel`, `applu`, and `gafort` benchmarks. All 4 have very low simulated miss ratios (less than 0.14% in the smallest simulated cache size). The particularly extreme numbers seen in `apsi` and `galgel` are due to reductions in capacity misses; the working set fits entirely into a 1MB cache, so the shared 1MB cache correctly predicts the extremely low miss ratio, but individual caches of 256KB do not. Conversely, for `applu` and `gafort` the errors are due mostly to a discrepancy in cold misses; they too have very low miss ratios, so the difference in cold misses is significant. For the rest of the benchmarks, which do not have such small working sets, the breakdown is more mixed, but slightly favors the impact of reduced capacity misses.

Overall, the prediction accuracy for shared caches is improved by 84% by using the multicore-aware method. Even if the 4 benchmarks with extremely low miss ratios are ignored, the improvement is substantial, averaging 80%.

The `galgel` benchmark has significant error even with the shared stack; this benchmark had significant numbers of conflict misses in the shared configuration that did not appear in the private configuration.

Looking next at the pairwise-shared cache results, the accuracy on most of the benchmarks is quite similar to that of the private invalidating stacks. Those where errors are higher tend to be those that also had higher error in the shared-cache configuration. The overall effect of adding multicore-awareness is a 76% improvement in modeling accuracy (69% if not counting the benchmarks with low miss ratios). Detailed results also show that effective reuse distance modeling for pairwise-shared caches requires pairwise-shared stacks; basic Eager or Shared models see large errors here.

5.3 Sensitivity Analysis

Besides the basic tests described above, additional experiments considered the sensitivity of our methods to number of cycles between simulator context switches and cache block size.

The base tests switch between simulating all of the processors in round-robin order on each instruction. Since this is unlikely to accurately represent a real execution, additional experiments were run with context switching every 100, 1000, 5000, and 10000 instructions. In all cases, all of these switching times caused less than 3% difference in accuracy for any of the reuse distance models tested: namely, the multicore-unaware models continued to show large and highly variable errors while the aware models had much smaller and more bounded errors. Most tests saw less than 1% difference between the various switching times.

Although the base tests use single-word cache lines, simulations were also run for all benchmarks under all sharing configurations with 64-byte line sizes (to more accurately reflect the designs of modern caches and to detect false sharing),

and with 4-byte line sizes for `genome` and `kmeans` (the only benchmarks dominated by 4-byte instead of 8-byte accesses). The results were very similar, with average accuracy and improvement within 5% of the numbers reported above.

6 Related Work

Section 2 covers the work most closely related to this paper. This section covers other related efforts.

Several works have used reuse distance analysis to model locality in broader ways. For example, Marin and Mellor-Crummey extend reuse distance analysis to identify the source, destination, and scope of a data reuse pattern, providing a list of suggested transformations such as data-splitting and loop fusion that can address misses generated in various ways [21]. Ding and Zhong present a scheme to detect access patterns based on approximate reuse analysis results [13]. Zhong et al. use those access patterns to predict cache behavior for all possible data sizes using a limited number of actual observations [32]. Shen et al. show that reuse-distance in programs exhibits phase-dependent behavior and that information-processing methods (such as wavelet transformations) can be used to understand the time-frequency characteristics of reuse [27]. Other efforts have focused on per-instruction characterization of reuse distance [14, 20]. All of these analysis and transformation methods considered the behavior of a single reference stream; the strategies therein could be combined with the methods from this paper to provide a better understanding of locality in multicore programs.

Berg et al. present a statistical model of multiprocessor caches with full associativity and random replacement [6]. That work models coherence misses probabilistically using a definition of reuse distance that counts all intervening references, whether distinct or not. In contrast, our work models coherence misses directly, considers how synchronization points may be used to propagate coherence information, and performs reuse distance analysis using the more commonly-studied metric of stack distance. These distinctions makes their work well-suited for cache performance evaluation but less so for application locality analysis.

Shi et al. present an analytical model of data replication and a method to simulate multiple caches in a CMP in a single pass [28]. Their method uses a list-based reuse stack to model private and shared caches and to account for data replication in the shared cache, but instead of calculating the actual reuse distance for each access, it simply accumulates the appropriate hit and miss counters for the various caches. Therefore it is unsuitable for uses other than simulation.

Agarwal and Gupta study shared memory reference patterns independent of a specific memory hierarchy by marking a reference as a *ping* if it accesses a block that was last accessed by a different core and as a *cling* otherwise [2]. The number of clings between pings on a block indicates how often that block remains associated with the same core (*processor locality*), the number of references (distinct or not) between clings at a given address represents temporal locality, and successive references by the same core to nearby addresses represent spatial locality. Their per-address approach aims primarily to explore cache coherence alternatives, whereas this work aims to model how application reuse characteristics interact with the memory hierarchy's sharing configuration to shape overall performance.

7 Conclusions

Reuse distance analysis has previously proven to be a valuable tool for cache modeling, application locality prediction, and locality-aware code transformations [8, 13, 14, 20, 21, 22, 27, 32]. This paper extends reuse distance analysis to the multicore domain by exploring methods to make the reuse stack account for inter-thread interactions, such as invalidations for per-core caches or inter-core prefetching for shared caches. The results show that modeling these interactions can lead to reuse distance analysis that is much more accurate in predicting application locality characteristics. The net effect is to substantially improve the robustness and applicability of reuse distance analysis in the face of evolving multicore hardware.

References

- [1] S. V. Adve and M. D. Hill. Weak ordering—a new definition. In *ISCA '90: Proceedings of the 17th annual international symposium on Computer Architecture*, pages 2–14, New York, NY, USA, 1990. ACM.
- [2] A. Agarwal and A. Gupta. Memory-reference characteristics of multiprocessor applications under mach. In *SIGMETRICS '88: Proceedings of the 1988 ACM SIGMETRICS conference on Measurement and modeling of computer systems*, pages 215–225, New York, NY, USA, 1988. ACM.
- [3] G. Almási, C. Caşcaval, and D. A. Padua. Calculating stack distances efficiently. *SIGPLAN Not.*, 38(2 supplement):37–43, 2003.
- [4] V. Aslot, R. Eigenmann, G. Gaertner, W. B. Jones, and B. Parady. Specomp: A new benchmark suite for measuring parallel computer performance. In *Workshop on OpenMP Applications and Tools*, pages 1–10, 2001.
- [5] B. T. Bennett and V. J. Kruskal. LRU stack processing. *IBM Journal of Research and Development*, 19:353–357, 1975.
- [6] E. Berg, H. Zeffer, and E. Hagersten. A statistical multiprocessor cache model. *Performance Analysis of Systems and Software, 2006 IEEE International Symposium on*, pages 89–99, March 2006.
- [7] K. Beyls and E. D'Hollander. Reuse distance as a metric for cache behavior. In *IASTED conference on Parallel and Distributed Computing and Systems 2001 (PDCS01)*, pages 617–662, 2001.
- [8] K. Beyls and E. D'Hollander. Generating cache hints for improved program efficiency. *Journal of Systems Architecture*, 51(4):223–250, 4 2005.
- [9] K. Beyls and E. H. D'hollander. Reuse distance-based cache hint selection. In *In Proceedings of the 8th International Euro-Par Conference*, pages 265–274, 2002.
- [10] D. Bhandarkar and J. Ding. Performance characterization of the pentium®pro processor. In *HPCA '97: Proceedings of the 3rd IEEE Symposium on High-Performance Computer Architecture*, page 288, Washington, DC, USA, 1997. IEEE Computer Society.
- [11] C. Cao Minh, J. Chung, C. Kozyrakis, and K. Olukotun. STAMP: Stanford transactional applications for multi-processing. In *IISWC '08: Proceedings of The IEEE International Symposium on Workload Characterization*, September 2008.
- [12] D. Dice, O. Shalev, and N. Shavit. Transactional Locking II. In *20th International Symposium on Distributed Computing*, pages 194–208, 2006.
- [13] C. Ding and Y. Zhong. Predicting whole-program locality through reuse distance analysis. In *PLDI '03: Proceedings of the ACM SIGPLAN 2003 conference on Programming language design and implementation*, pages 245–257, New York, NY, USA, 2003. ACM.
- [14] C. Fang, S. Carr, S. Önder, and Z. Wang. Reuse-distance-based miss-rate prediction on a per instruction basis. In *In Proceedings of the first ACM SIGPLAN Workshop on Memory System Performance*, pages 60–68, 2004.
- [15] Intel Corporation. Intel core 2 extreme quad-core processor qx6000 sequence and intel core 2 quad processor q6000 sequence. Datasheet, Aug. 2007.
- [16] H. Jin, M. Frumkin, and J. Yan. The OpenMP implementation of NAS parallel benchmarks and its performance.
- [17] M. T. Jones. Inside the Linux scheduler. *IBM DeveloperWorks*, June 2006.
- [18] E. Larson and S. Chatterjee. Mase: A novel infrastructure for detailed microarchitectural modeling. In *Proceedings of the 2001 International Symposium on Performance Analysis of Systems and Software*, pages 1–9, 2001.
- [19] P. S. Magnusson, M. Christensson, J. Eskilson, D. Forsgren, G. Hallberg, J. Hogberg, F. Larsson, A. Moestedt, and B. Werner. Simics: A full system simulation platform. *IEEE Computer*, 35(2):50–58, Feb. 2002.

- [20] G. Marin and J. Mellor-Crummey. Cross-architecture performance predictions for scientific applications using parameterized models. In *SIGMETRICS '04/Performance '04: Proceedings of the joint international conference on Measurement and modeling of computer systems*, pages 2–13, New York, NY, USA, 2004. ACM.
- [21] G. Marin and J. M. Mellor-Crummey. Pinpointing and exploiting opportunities for enhancing data reuse. In *IEEE International Symposium on Performance Analysis of Systems and Software, ISPASS 2008*, pages 115–126, Apr. 2008.
- [22] R. L. Mattson, J. Gececi, D. R. Slutz, and I. L. Traiger. Evaluation Techniques for Storage Hierarchies. *IBM Systems Journal*, 9(2):78–117, 1970.
- [23] F. Olken. Efficient methods for calculating the success function of fixed space replacement policies. Technical Report LBL-12370, Lawrence Berkeley Laboratory, 1981.
- [24] OpenMP Architecture Review Board. *OpenMP Application Program Interface: Version 3.0*, May 2008.
- [25] V. S. Pai, P. Ranganathan, and S. V. Adve. The impact of instruction-level parallelism on multiprocessor performance and simulation methodology. In *HPCA '97: Proceedings of the 3rd IEEE Symposium on High-Performance Computer Architecture*, page 72, Washington, DC, USA, 1997. IEEE Computer Society.
- [26] M. K. Qureshi, D. N. Lynch, O. Mutlu, and Y. N. Patt. A Case for MLP-Aware Cache Replacement. In *ISCA '06: Proceedings of the 33rd annual international symposium on Computer Architecture*, pages 167–178, Washington, DC, USA, 2006. IEEE Computer Society.
- [27] X. Shen, Y. Zhong, and C. Ding. Locality phase prediction. In *ASPLOS-XI: Proceedings of the 11th international conference on Architectural support for programming languages and operating systems*, pages 165–176, New York, NY, USA, 2004. ACM.
- [28] X. Shi, F. Su, J.-K. Peir, Y. Xia, and Z. Yang. Modeling and single-pass simulation of cmp cache capacity and accessibility. In *Proc. IEEE International Symposium on Performance Analysis of Systems & Software ISPASS 2007*, pages 126–135, 2007.
- [29] S. T. Srinivasan and A. R. Lebeck. Load latency tolerance in dynamically scheduled processors. In *MICRO 31: Proceedings of the 31st annual ACM/IEEE international symposium on Microarchitecture*, pages 148–159, Los Alamitos, CA, USA, 1998. IEEE Computer Society Press.
- [30] R. A. Sugumar and S. G. Abraham. Multi-configuration simulation algorithms for the evaluation of computer architecture designs. Technical report, University of Michigan, 1993.
- [31] R. Vaswani and J. Zahorjan. The Implications of Cache Affinity on Processor Scheduling on Multiprogrammed, Shared Memory Multiprocessors. In *Proceedings of the 13th Symposium on Operating Systems Principles*, pages 26–40, Oct. 1991.
- [32] Y. Zhong, S. G. Dropsho, and C. Ding. Miss rate prediction across all program inputs. In *In Proceedings of the 12th International Conference on Parallel Architectures and Compilation Techniques*, pages 91–101. IEEE Computer Society, 2003.

A Detailed Cache Results

Figure 4 expands upon the results reported in Section 5, specifically considering the prediction error in cache miss ratio between the reuse distance analysis models and the actual execution-driven simulation for various cache sizes. The x axis is cache size and the y axis is the percent error, plotted on a log scale (except for a linear scale on equake, MG, and swim, all of which have low errors). Some of the bars are too small to be visible. Some benchmarks have very small errors in all cases; for example MG, equake and swim have errors of less than 2% for all configurations. Conversely, applu, apsi, gafort and galgel have errors close to 100% for the Unaware methods and much lower for the aware method, with the rest of the 13 benchmarks somewhere in between. Nine of the benchmarks never have more than 10% error for any size cache when using the multicore-aware methods. Aplu and FT see slightly larger errors for 256 K caches, and Apsi and Galgel see slightly larger errors for 512 K caches: in both cases, these are caused by conflict misses that arise only at certain cache sizes.

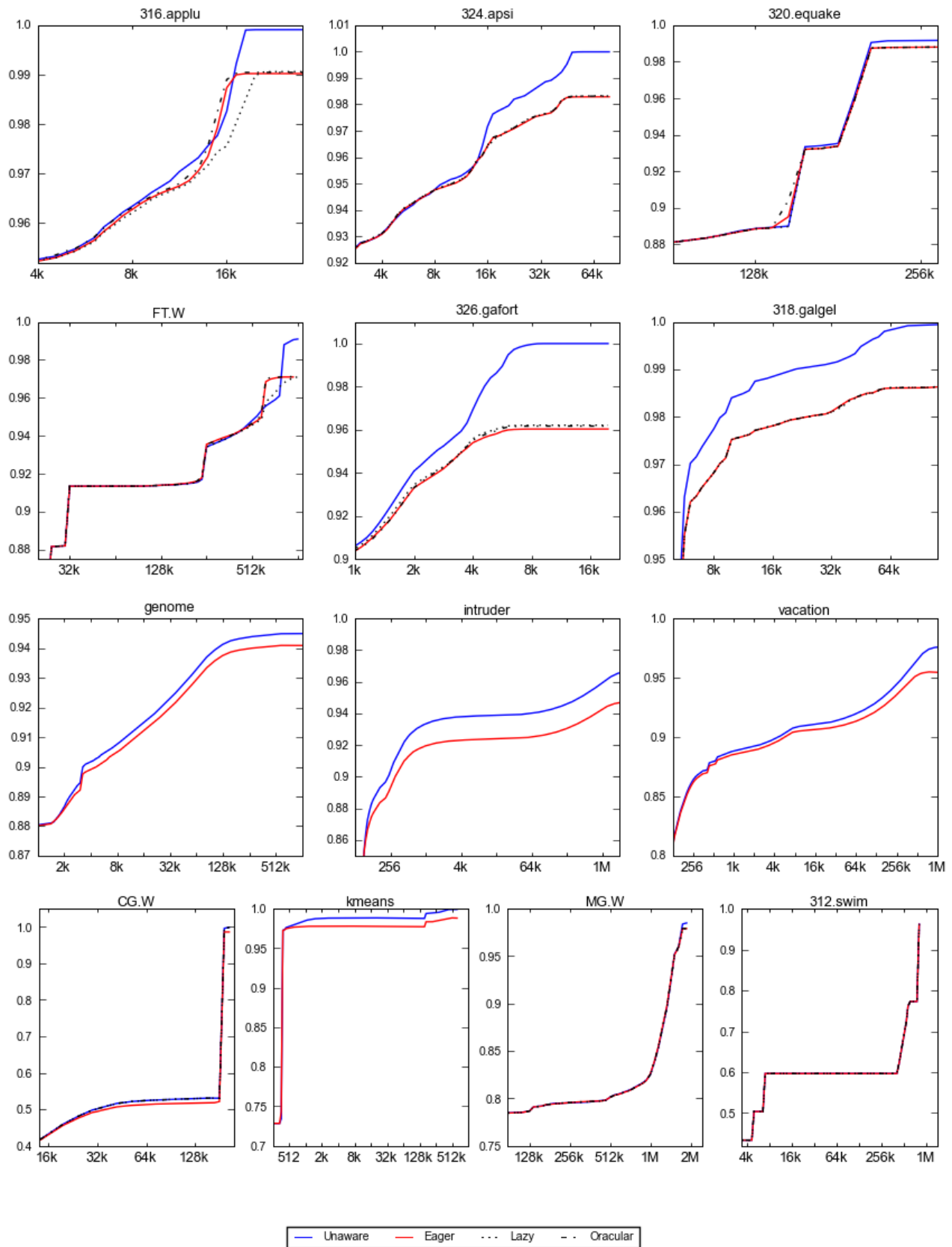


Figure 1: CDFs of reuse distance using private per-thread stacks. X axis is reuse distance in 64-bit words, Y axis is fraction of references with reuse distance less than x .

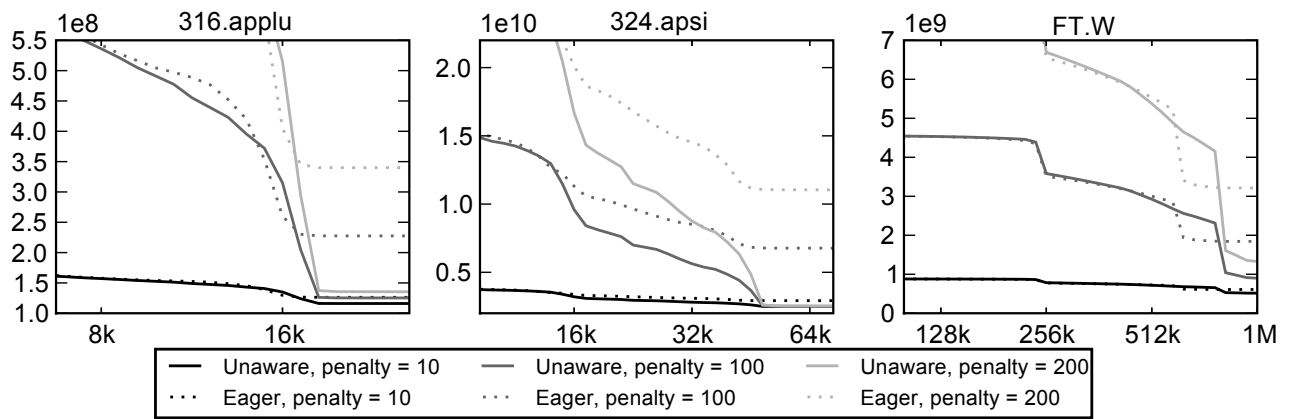


Figure 2: The simple model of execution time shown here just considers a unit cycle count for all ordinary instructions and cache hits, along with varying values for cache miss penalty as shown. The X axis represents the fully-associative LRU cache size in 64-bit words, and the Y axis represents the execution time in cycles predicted by this model for multicore-unaware and aware reuse distance methods.

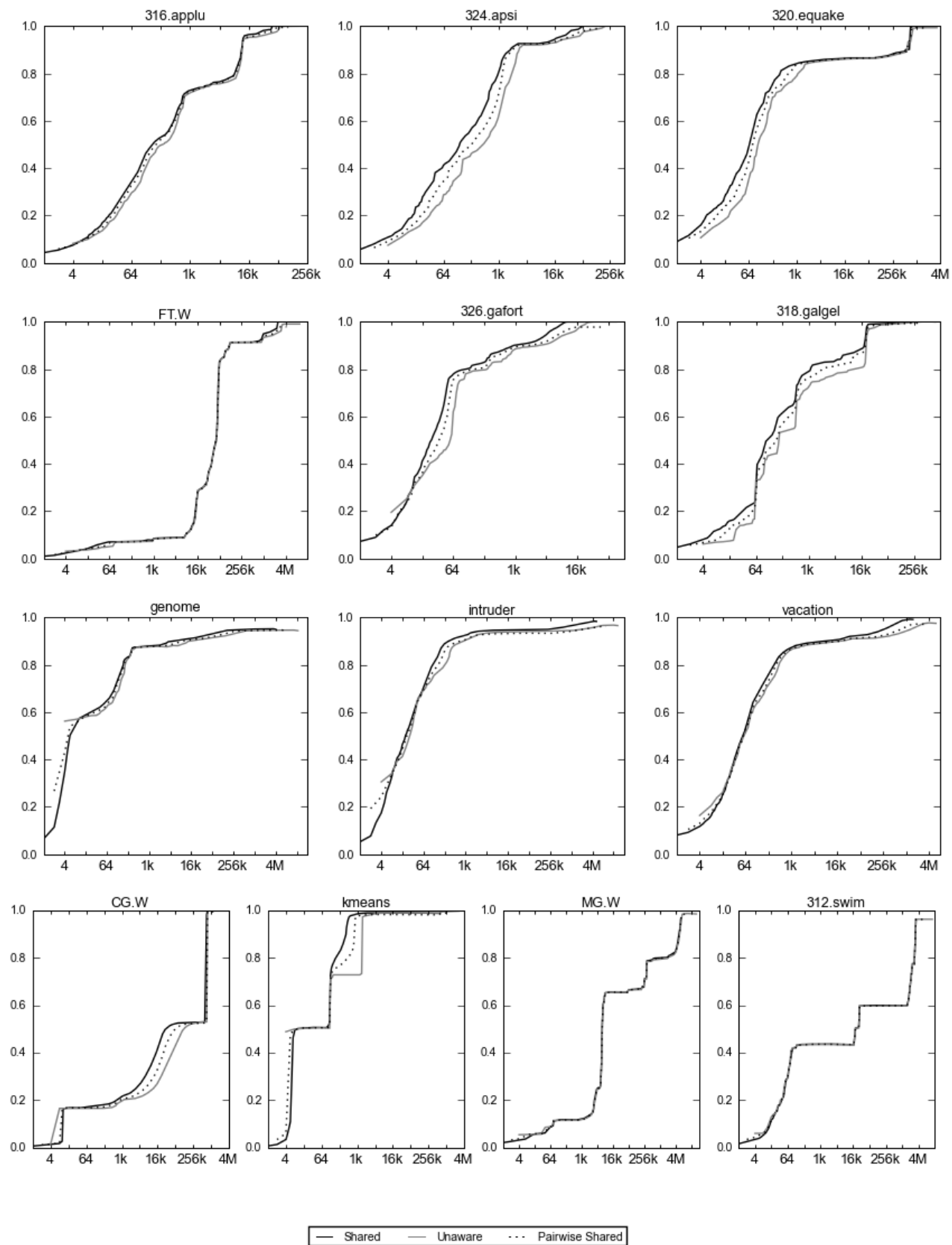


Figure 3: CDFs of reuse distance using a single shared or pairwise shared stacks. X axis is reuse distance in 64-bit words, Y axis is fraction of references with reuse distance less than x .



Figure 4: Detailed breakdown of percent error for all sizes of private caches for Unaware and Aware methods