

Traffic Management: A Holistic Approach to Memory Placement on NUMA Systems

ABSTRACT

NUMA systems are characterized by Non-Uniform Memory Access times, where accessing data in a remote node takes longer than in a local node. NUMA hardware has been built since the late 80's, and the operating systems designed for it optimized for access locality. They placed threads and memory pages on the system's nodes so as to reduce the number of remote accesses. We discovered that contrary to older systems, remote wire delays have been substantially reduced on modern hardware, and so remote access costs *per se* are not the main concern for performance. Instead, *congestion on memory controllers and interconnects*, caused by memory traffic from data-intensive applications, hurts performance a lot more. Because of that, thread and memory placement algorithms must be completely redesigned to target traffic congestion, which requires different techniques than optimizing locality. We propose *Carrefour*, an algorithm that addresses this goal. We implemented *Carrefour* in Linux and obtained performance improvements of up to $2.5\times$ relative to the default kernel, as well as significant improvements compared to NUMA-aware patchsets available for Linux. *Carrefour* never hurts performance by more than 6% when memory placement cannot be improved. In this paper we describe the design of *Carrefour*, the challenges of implementing it on modern hardware, and draw insights on hardware support in future NUMA systems.

Categories and Subject Descriptors

D.4.1 [OPERATING SYSTEMS]: Process Management—*scheduling*

General Terms

Algorithms, Performance, Measurement

Keywords

NUMA, operating systems, multicore, scheduling

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

ASPLOS 2013 Houston, Texas USA

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

1. INTRODUCTION

Modern servers are built from several *processor nodes*, each containing a multicore CPU and a local DRAM serviced by one or more memory controllers (see Figure 1). The nodes are connected into a single cache-coherent system by a high-speed *interconnect*. Physical address space is globally shared, so all cores can transparently access memory in all nodes. Accesses to a local node go through a local memory controller; accesses to remote nodes must traverse the interconnect and access a remote controller. Remote accesses typically take longer than local ones, giving these systems a property of Non-Uniform Memory Access time (NUMA).

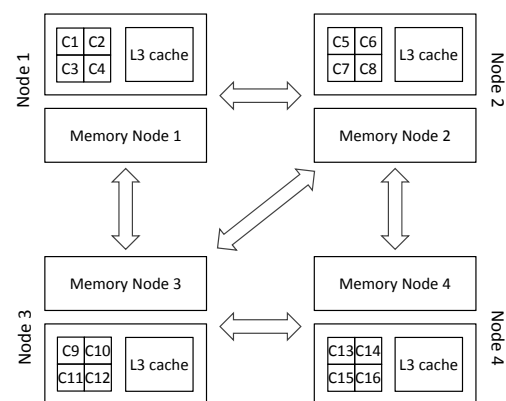


Figure 1: A modern NUMA system, with four nodes and four cores per node. At the time of the writing, NUMA systems are built with up to 8 nodes and 10 cores per node.

It is well understood that optimal performance on NUMA systems can be achieved only if we place threads and their memory in consideration of the system's physical layout. For instance, previous work on NUMA-aware thread and memory placement focused on maximizing locality of accesses, that is, placing threads and memory pages such that data accesses are satisfied from a local node whenever possible. That was done to avoid very high costs of remote memory accesses. Contrary to insights from previous work, we discover that on modern NUMA systems remote wire delays (that is, delays resulting from traversing a greater physical distance to reach a remote node) are *not* the most important source of performance overhead. On the other hand, *congestion on interconnect links and in memory controllers*, which results from high volume of data flowing across the

system, can dramatically hurt performance. This motivates the design of new NUMA-aware memory placement policies.

To make these statements concrete, consider the following facts. On NUMA systems circa 1990s, the time to access data from a remote node took 4-10 times longer than from a local node [31]. On NUMA systems that are built today, remote wire delays add at most 30% to the cost of a memory access [7]. For most programs, this latency differential alone would not have a substantial impact on performance. However, fast modern CPUs are able to generate memory requests at very high rates. Massive data traffic creates congestion in memory controller queues and on interconnects. When this happens, memory access latencies can become as large as 1000 cycles (!), from a normal latency of only around 200. Such a dramatic increase in latencies can slow down data-intensive applications by more than a factor of 2 \times . Fortunately, high latencies can be avoided or substantially reduced if we carefully place threads and memory pages on nodes so as to avoid traffic congestion.

In response to the changes in hardware bottlenecks, we approach the problem of thread and memory placement on NUMA systems from an entirely new perspective. We look at it as the problem of *traffic management*. Our algorithm, called *Carrefour*¹, places threads and memory so as to avoid traffic hotspots and prevent congestion in memory controllers and on interconnect links. This is akin to traffic management in the context of city planning: popular residential and business hubs must be placed so as to avoid congestion on the roads leading to these destinations.

The mechanisms used in our algorithm: e.g., migration and replication of physical memory pages, are well understood, but the algorithm itself is new. Our algorithm makes decisions based on global observations of traffic congestion. Previous algorithms optimized only for locality, not congestion, and relied only on local information (e.g., access pattern of individual pages) when making decisions.

Implementing an effective NUMA-aware algorithm on modern systems presents several challenges. Modern systems do not have the same performance monitoring hardware that was present (or assumed) on earlier systems. Existing instruction sampling hardware cannot gather the profiling data needed for the algorithm with the desired accuracy and speed. We had to navigate around this problem in our design. Furthermore, the memory latencies that we are optimizing are lower than on older systems, so we can tolerate less overhead in the algorithm.

We implemented *Carrefour* in Linux and evaluated it with several data-centric applications: k-means clustering, face recognition, map/reduce, and others. *Carrefour* improves performance of these applications, with the largest gain of 2.5 \times speedup. When memory placement cannot be improved *Carrefour* never hurts performance by more than 6%. Existing NUMA-aware patches for the Linux kernel perform less reliably and in general fall short of improvements achieved with *Carrefour*.

In addition to addressing NUMA-aware memory placement, *Carrefour* includes important ideas proposed in other multicore scheduling algorithms, such as contention-aware thread placement [32, 5] and sharing-aware thread clustering [14, 29]. As a result, *Carrefour* was designed to be a complete scheduling algorithm for multicore systems.

¹*Carrefour*— (French) intersection, crossroads.

2. TRAFFIC CONGESTION ON MODERN NUMA SYSTEMS

In this section, we demonstrate that the effects of traffic congestion are more substantial than those of wire delays, and motivate why memory placement algorithms must be redesigned to consider them. To that end, we report data from two sets of experiments. In the first set, our goal is to measure purely the effects of wire delays. We run applications in two configurations: *local-memory* and *remote-memory*. Under *local-memory*, the thread and its data are co-located on the same NUMA node; under *remote-memory*, the thread runs on a different node than its data. To ensure that wire delay is the dominant performance factor, we had to avoid congestion on memory controllers and interconnects, so we run one application at a time and use only one thread in each application. We do not include applications with CPU utilization less than 30%, because memory performance is not their main bottleneck. The experiments are run on a system described in Section 4 (Machine A – illustrated in Figure 1). We use applications from the NAS, PARSEC and Map/reduce Metis suites, also described in Section 4.

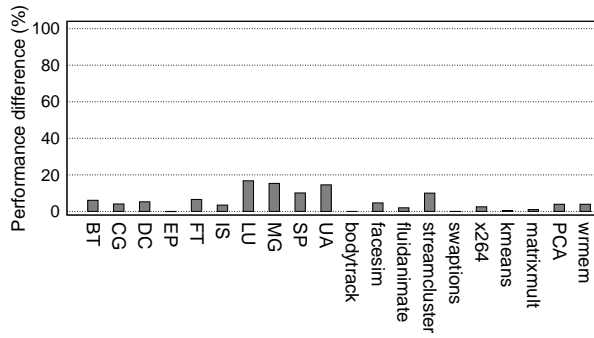
Figure 2(a) shows relative completion time under *remote-memory* vs. *local-memory* configuration. The performance degrades by at most 20% under *remote-memory*, which is consistent with at most 30% difference in local-vs-remote memory latencies measured in microbenchmarks [7].

In the second set of experiments, we want to observe traffic congestion, so we run each application with as many threads as there are cores. We demonstrate how performance varies under two memory placement policies on Linux, as they induce different degrees of traffic congestion. We compare *First-touch* (F) - the default policy where the memory pages are placed on the node where they are first accessed, and *Interleaving* (I) - where memory pages are spread evenly across all nodes. Although these are not the only possible and not necessarily the best policies, comparing them illustrates the salient effects of traffic congestion.

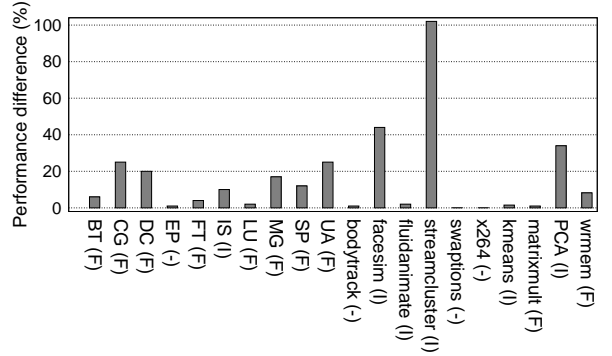
Figure 2(b) shows the absolute difference in completion time achieved under *first-touch* and *interleaving*. The policy that performed the best is indicated in parenthesis next to the application name; a "-" shown when the application performs equally well with either policy. We observe that the differences are often much larger than what we can expect from wire delays alone. For a k-means clustering application *Streamcluster* from PARSEC, the performance varies by a factor of 2 \times depending on the memory placement policy!

To illustrate that these differences are due to traffic congestion, we show Table 1 some supporting data for the two most affected applications, *Streamcluster* and *PCA* (a map/reduce application)². **Local access ratio** is the percent of all memory accesses sourced from a local node; **Memory latency** is the average number of cycles to satisfy a memory request from any node; **Memory controller imbalance** is the standard deviation of the load across all memory controllers, expressed as percent of the mean. Load is measured as the number requests per time unit; **Average interconnect (IC) usage** shows the utilized interconnect bandwidth as percent of total, averaged across all links, the **imbalance** shows the standard deviation of utilization across the links

²We are unable to present the same data for all applications due to space constraints, but the conclusions reached from their measurements are qualitatively similar.



(a) Perf. difference for single-thread versions of applications between local and remote memory configurations.



(b) Absolute perf. difference for multi-thread versions of applications between First-touch (F) and interleaving (I).

Figure 2: Performance difference of applications depending on the thread and memory configuration.

	<i>Streamcluster</i>		<i>PCA</i>	
	Best (I)	Worst (F)	Best (I)	Worst (F)
Local access ratio	25%	25%	25%	33%
Memory latency	471	1169	480	665
Mem-ctrl. imbalance	7%	200%	4%	154%
IC: imbalance, (avg)	21% (58%)	86% (32%)	19% (47%)	64% (30%)
L3MPKI	16.61	16.35	6.99	7.1
IPC	0.30	0.15	0.52	0.37

Table 1: Traffic congestion effects

as percent of mean utilization; *L3MPKI* is the number of last-level (L3) cache misses per thousand instructions; *IPC* is the number of instructions per cycle.

The data in Table 1 leads to several curious observations. First, we see that locality of memory accesses either does not change regardless of the memory management policy, or *decreases* under the better performing policy. For *Streamcluster*, most of the memory pages happen to be placed on a single node under *first-touch* (because a single thread initializes them at the beginning of the program). Under *interleaving* the pages are spread across all nodes, but since the threads access data from all four nodes, the overall access ratio is about the same in both configurations. For *PCA*, interleaving decreases the local access ratio and yet *increases* performance. So the first surprising conclusion is that the better-performing policy *does not* necessarily improve access locality!

And yet, the IPC substantially improves (2× for *Streamcluster* and 41% for *PCA*), and not because of changes in the L3 miss rate (L1 and L2 cache miss rates also remained unchanged). The explanation emerges if we look at the memory latency. Under interleaving, memory latency reduces by a factor of 2.48 for *Streamcluster* and 1.39 for *PCA*. The question is, *what is responsible for memory latency improvements?* It turns out that interleaving *dramatically reduces memory controller and interconnect congestion* by alleviating the load imbalance and mitigating traffic hotspots. Rows

5, 6 in Table 1 show significant reductions in imbalance under interleaving, and Figure 3 illustrates these effects visually for *Streamcluster*. So even without improving locality – we even *reduce* it for *PCA*(!), we are able to substantially improve performance. And yet, existing NUMA-aware algorithms disregarded traffic congestion, optimizing for locality only. Our work addresses this shortcoming.

Although the two selected applications performed significantly better under interleaving, this does not mean that interleaving is the only desired policy on modern NUMA hardware. In fact, as Figure 2(b) shows, many NAS applications fared a lot worse with interleaving. In the process of designing the algorithm we learned that an arsenal of techniques – interleaving, page replication and co-location (placing the pages on the same node as the threads that access them) – must be judiciously applied to different parts of the address space depending on global traffic conditions and page access patterns, and further combined with methods like sharing-aware thread clustering. So the challenge in designing a good algorithm is understanding when to apply each technique, while navigating around the challenges of obtaining accurate performance data and limiting the overhead.

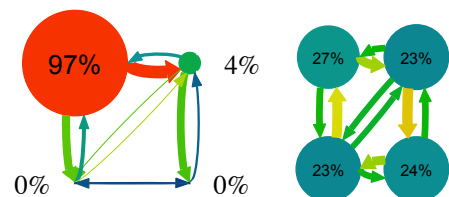


Figure 3: Traffic imbalance under first-touch (left) and interleaving (right) for *streamcluster*. Nodes and links bearing the majority of the traffic are shown proportionately larger in size and in brighter colours. The figure is drawn to scale.

3. DESIGN AND IMPLEMENTATION

We begin by describing the mechanisms composing the algorithm: *page co-location*, *interleaving*, *replication* and *thread clustering*. Then we explain how they fit together.

3.1 The mechanisms

Page co-location is when we re-locate the physical page to the same node as the thread that accesses it. Co-location works well for pages that are accessed by a single thread or by threads co-located on the same node.

Page interleaving is about evenly distributing physical pages across nodes. Interleaving is useful when we have imbalance on memory controllers and interconnect links, and when pages are accessed by many threads. Operating systems usually provide an interleaving allocation policy, but only give an option to enable or disable it globally for the entire application. We found that interleaving works best when judiciously applied to parts of the address space that will benefit from it.

Page replication is about placing a copy of a page on each memory node. Replication distributes the pressure across memory controllers, alleviating traffic hotspots. An added bonus is eliminating remote accesses on replicated pages. When done right, replication can bring very large performance improvements. Unfortunately, replication also has costs. Since we keep multiple copies of the same page, we must synchronize their contents, which is like running a cache coherency protocol in software. The costs can be very significant if there is a lot of fine-grained read/write sharing. Another potential source of overhead is the synchronization of page tables. Since modern hardware walks page tables automatically, page tables themselves must be replicated across nodes and kept in sync. For workloads with many page faults, page table synchronization has a large overhead. Finally, replication increases the memory footprint, so it should be avoided for workloads with very large working sets, for fear of introducing additional hard page faults.

Thread clustering is about co-locating threads that share data on the same node, to the extent possible³. Two thread clustering algorithms have been described in previous work [14, 29]. Tam’s algorithm [29] always clustered data-sharing threads. Kamali’s algorithm [14] balanced between the benefits of co-operative sharing and contention, and only clustered threads when that would not create excessive contention for shared resources. We confirmed the importance of balancing the two goals in our experiments, so our algorithm includes Kamali’s variant. Since thread clustering is well documented and understood, we did not evaluate it in our implementation at this point⁴, but our algorithm assumes a thread placement phase that is carried out before any memory placement decisions.

3.2 The Algorithm

Our memory management algorithm has three components: *measurement*, *global decisions* and *page-local decisions*. The measurement component continuously gathers various metrics (Table 2) that later drive page placement decisions. Global and per-application metrics are collected using hardware counters with very low overhead. Per-page statistics are collected via instruction-based sampling (IBS), which can introduce significant overheads at high sampling rates. Section 3.3 describes how we keep the overheads at bay. Global decisions look at system-wide traffic congestion

³We never want to sacrifice CPU load balance in favour of clustering.

⁴We plan to add this to the final version of the paper.

Global statistics	
MC-IMB	Memory controller imbalance (as defined in Section 2)
LAR	Local access ratio (as defined in Section 2)
MAPTU	Memory (DRAM) accesses per time unit (microsecond)
Per-application statistics	
MRR	Memory read ratio. Fraction of DRAM accesses that are reads
IPC	Instructions per cycle
Per-page statistics	
Number of accesses	The number of sampled data loads that fell in that page
Access type	Read-only or read-write

Table 2: Statistics collected for the algorithm.

and workload properties to determine what mechanisms to use. Page-local decisions examine access patterns of individual pages to decide their fate.

3.2.1 Global Decisions

The global decision-making process is outlined in Figure 4. First of all, we decide whether to enable *Carrefour*. We only want to run *Carrefour* if memory accesses appear to be the bottleneck, otherwise there is no point in incurring the sampling overhead. This decision is driven by two parameters: memory access rate (MAPTU) and the IPC. A high MAPTU indicates large memory traffic, while a low IPC typically means that the CPU spends lots of time stalled on DRAM accesses. In our implementation, we enable *Carrefour* if the memory access rate (MAPTU) exceeds 50 and the IPC of at least one application is below 0.7.

Second, we decide whether it is worthwhile to use replication. Replication is likely to help when we have sufficient physical RAM capacity for replicated pages, when the workload accesses its data set in a read-only or read-mostly fashion, and when we do not have an extremely high page fault rate. (Section 3.3 elaborates why replication has high overheads when the read-write ratio is low or the page fault rate is high.) As a result, this decision is driven by three parameters: free physical RAM, read-write ratio, and page fault rate. We enable replication if free RAM is at least $1 - \frac{1}{NUM_NODES}$ of the total. This configuration is crucial for applications with very large memory footprints.

With regards to the read-write ratio, ideally, we want to compute it for each individual page, so as to avoid replicating pages that are frequently written. However, due to overheads of IBS, as we explain in Section 3.3, obtaining per-page memory read ratio is very difficult, so instead we use a per-application MRR, which we can obtain inexpensively via hardware counters. We disable replication if the MRR is below 95%⁵.

Finally, we disable replication for workloads that generate more than 500 page faults per second. In general we found that performance is not very sensitive to this parameter, as only the workloads that produce thousands of page faults per second suffer high overheads under replication.

⁵MRR is approximated as fraction of L1 refills from DRAM in modified state, because there is not a hardware counter that provides this quantity precisely per core, as opposed to per-node.

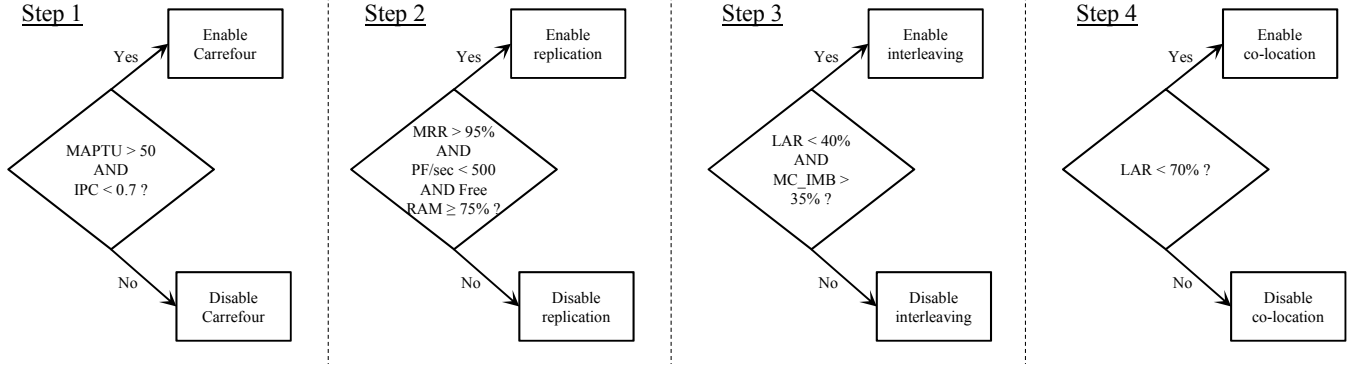


Figure 4: Global decisions in *Carrefour*.

Third, we decide whether or not to use interleaving. Interleaving is useful if we have a large memory controller imbalance and a low local-access ratio. We enable interleaving if the local access ratio (LAR) is below 40% and the memory controller imbalance is above 35%.

Finally, we decide whether or not to enable co-location. Co-location will be triggered only for pages that are accessed from a single node, and so it will not exacerbate the imbalance if memory-intensive threads are evenly spread across nodes, which is what the thread-clustering step would do. Therefore, we enable co-location if the local access rate is slightly less than ideal ($LAR < 70\%$).

3.2.2 Page-local Decisions

Carrefour makes page-local decisions depending on the mechanisms enabled: e.g., pages are only considered for replication if replication is enabled for that application. The following explanation assumes that all three mechanisms are enabled.

To decide the fate of each page, we need at least two memory-access samples for that page. If the page was accessed from only a single node we migrate it to that node. We do not migrate the thread, because we assume that thread clustering, performed before memory placement, already made good thread placement decisions. If the page was accessed from two or more nodes in a read-only fashion, we replicate it. Otherwise we mark it for interleaving. To decide where to move a page marked for interleaving, we use two probabilities: $P_{migrate}$ and P_{node} . $P_{migrate}$ determines the likelihood of migrating the page away from the current node. $P_{migrate}$ is the MAPTU of the current node as fraction of MAPTU on all nodes, so the higher the load on the current node relative to others, the higher the chance that we will migrate a page. P_{node} gives us the probability of migrating a page to a particular node, and it is the opposite of $P_{migrate}$ for that node, so *Carrefour* will migrate the page to the least loaded node.

3.3 Implementation

We implemented *Carrefour* in the Linux kernel 3.2.1. *Carrefour* runs continuously, alternating the measurement and decision-making phases with periodicity of one second. To a large extent, *Carrefour* relies on well-understood mechanisms in the Linux kernel, such as physical page migration. The non-trivial aspects of the implementation were under-

standing how to accomplish fast and accurate sampling of memory accesses and navigating around the overheads of replication.

3.3.1 Fast and accurate memory access sampling

A crucial goal of the algorithm is to quickly and accurately detect memory pages that cause the most DRAM accesses, and accurately estimate the read/write ratio of those pages. To that end, we used Instruction-Based Sampling (IBS) – hardware-supported sampling of instructions available in AMD processors⁶. IBS can be configured to deliver instruction samples at a desired interval, e.g., after expiration of a certain number of cycles or micro-ops. Each sample contains detailed information about the sampled instruction, such as the address of the accessed data (if the instruction is a load or a store), whether or not it missed in the cache and how long it took to fetch the data. Unfortunately, every delivered sample generates an interrupt, so processing samples at a high rate becomes very costly. Other systems that used IBS relied on it for off-line profiling [16, 26], so they could tolerate much higher overhead.

After experimenting with IBS, we found that the sampling interval of 130,000 cycles incurs a reasonable overhead of less than 5% for most programs we tested. Our initial decision was to filter out all the samples that did not generate a DRAM access. However, we found that the resulting number of samples was extremely low. Even very memory-intensive workloads access DRAM only a few times for every thousand instructions. That, combined with a low IBS sampling frequency, gave us the sampling rate of less than one hundred thousandth of a percent, and made it very difficult to generate a sufficient number of samples.

The second problem with considering only instructions that access DRAM is that we miss a lot of the accesses generated by the hardware prefetcher. These accesses are not part of any instruction so they will not be tagged by IBS. For prefetch-intensive applications, we obtain a very small number of samples and a very distorted read-write ratio. For example, at least five applications that we tested had a very high read-write ratio (close to 95%) when hardware prefetch requests were included, but showed a ratio of only about around 5% when prefetch requests were excluded.

As a result, our replication mechanism performed very

⁶PEBS – Precise Event-Based Sampling, is a similar feature in Intel processors

poorly, slowing down *streamcluster*, the application that stands to gain the most from replication, by a factor of 4, while manual replication *improved* its performance by more than 2.5 \times . We were not replicating enough pages, simply because we had too few samples, and were unable to identify frequently written pages, because of inaccurate read-write ratio. The latter caused our replication mechanism to perform a lot of page synchronization (more discussion below) and caused significant overhead.

To address this problem, we used two solutions. First, is the adaptive sampling rate. When the program begins to run, we sample it at a relatively high rate of 1/65,000 cycles. If after this measurement phase we take fewer than ten actions in the algorithm – an action is considered any change in page placement – we switch to a much lower rate of 1/260,000 cycles. Otherwise we continue sampling at the high rate.

The second solution was, when filtering IBS samples, to retain not just the data samples that accessed the DRAM, but those that hit in the first-level cache as well. First-level cache loads include accesses to prefetched data, so we avoid prefetcher-related inaccuracy. Considering cache accesses can introduce “noise” in the data, because we could be sampling pages that never access DRAM. On the other hand, *Carrefour* is only activated for memory-intensive applications, and for them there is a higher correlation between the accesses that hit in the cache and those that go to DRAM.

With the two solutions combined, we were able to detect all the pages that are worth replicating for *streamcluster*⁷.

However, even though performance became much better (we were able to speed up *streamcluster* by 26% relative to the default kernel), we were still far from the “ideal” manual replication, which sped it up by more than 2.5 \times . To approach ideal performance, we had to mitigate the overheads of replication, which we describe next.

3.3.2 Replication

Replication has three perils. First, it increases the workload’s memory footprint and can generate additional hard page faults (when we have to swap pages to disk). As explained earlier, we avoided this overhead by conservatively disabling replication for workloads with large working sets. Second, the kernel must synchronize the contents of replicated pages when they are written. This involves a physical page copy and is very costly. Third, the kernel must synchronize the contents of page tables. To illustrate the source of these overheads and explain how we mitigated them, we describe the relevant details of our implementation.

In Linux, a process address space is represented by a `mm_struct`, which keeps track of valid address space segments and holds a pointer to the page table, which is stored as a hierarchical array. Since modern hardware walks page tables automatically, we cannot modify the structure of the page table to point to several physical locations (one for each node) for a given virtual page. Instead, we must maintain a separate copy of the page table for each node and synchronize the page tables when they are modified, even for virtual pages that are not replicated. Our implementation

maintains a separate copy of `mm_struct` for each node, called a *slave mm*, as well as the *master mm*. The master serves as a synchronization point for the slaves. Whenever a core generates a page fault for a previously unmapped page, it must install a mapping in the slave `mm` corresponding to its node, as well as on the master. Cores on other nodes will then find the mapping on the master if they fault on the same page. An alternative is to implement a distributed synchronization mechanism, but this would incur a greater complexity, so we did not consider it.

Installing the mapping on the master requires holding a write lock. Lock contention can become a problem for workloads that allocate lots of memory and install lots of new mappings into page tables. We observed this to be the case with one map/reduce workload. That is why, to avoid synchronization overhead, we disable replication for workloads with an extremely high soft page fault rate.

When a page is replicated, we create a physical copy on every memory node that runs threads from the corresponding application. We install a different virtual-to-physical translation in each node’s page table. We write-protect the replicated page, so when any node writes that page we receive a page protection fault. To handle this fault, we read-protect the page on all nodes except the faulting one, and enable writing on the faulting node. If another node accesses that page, we must copy the new version of the page to that node, enable the page for reading and protect it from writing.

We refer to all these actions needed to keep the page synchronized as *page collapses*. Collapses are extremely costly, and would occur if we replicate a page that is write-shared. With limited capabilities of IBS, we are unable to detect read/write ratio on individual pages with sufficient accuracy. Collapses are so costly that even with very infrequent writes (e.g., one in 1000 accesses), we could have a high collapse overhead, but with IBS we are unable to obtain a sufficient number of samples. That is why we look at application-wide MRR and disable replication for the entire application if the MRR is low. This optimization was crucial to avoid the overhead of replication. Furthermore, we monitor collapse statistics of individual pages and disable replication for any page that generated more than five collapses.

With these optimizations, including those described in Section 3.3.1, we were able to avoid replication costs for the applications we tested and approached within 10% the performance of manual replication for *streamcluster*.

4. EVALUATION

In this section, we study the performance of *Carrefour* on a set of benchmarks. The main questions that we address are the following ones: (i) *how does Carrefour impact the performance of applications, including those that cannot benefit from its heuristics?*, (ii) *how does Carrefour compare against existing heuristics for modern NUMA hardware?*, (iii) *how does Carrefour leverage the different memory placement mechanisms?*, and (iv) *what is the overhead of Carrefour?*.

To assess the performance of *Carrefour*, we compare it against three other configurations: *Linux* – a standard Linux kernel with the default first-touch memory allocation policy, *Manual interleaving* – a standard Linux kernel with the interleaving policy manually enabled for the application, and *AutoNUMA* – the recent Linux patchset with a NUMA-

⁷*Streamcluster* holds shared data in a single large array, so it is trivial to detect which data is worth replicating and implement a manual solution to use as the performance upper-bound.

aware memory management policy [2]⁸.

The rest of the section is organized as follows: we first describe our experimental testbed. Next, we study single-application scenarios, followed by workloads with multiple co-scheduled applications. We then detail the overhead of *Carrefour* and conclude with a discussion on additional hardware support that would improve *Carrefour*'s operation.

4.1 Testbed

We used two different machines for the experiments:

Machine A has four 2.3GHz AMD Opteron 8385 processors with 4 cores in each (16 cores in total) and 64GB of RAM. It features 4 memory nodes (i.e., 4 cores and 16GB of RAM per node) interconnected with HyperTransport 1.0 links.

Machine B has two 1.7GHz AMD Opteron 6164 HE processors with 12 cores in each (24 cores in total) and 48GB of RAM. It features 4 memory nodes (i.e., 6 cores and 12GB of RAM per node) interconnected with HyperTransport 3.0 links.

All the single-application experiments (i.e., from section 4.2) were performed on machine A and the multi-application experiments were performed on machine B. We actually ran each experiment on both machines and obtained qualitatively similar results but, due to space constraints we cannot show all the numbers.

We used the Linux kernel v3.2.1 for all experiments, except for the AutoNUMA configuration, where we used v3.5.0.

We used the following set of applications: the PARSEC benchmark suite v2.1 [25], the FaceRec facial recognition engine v5.0 [10], the Metis MapReduce benchmark suite [21] and the NAS parallel benchmark suite v3.3 [23]. PARSEC applications run with the native workload. From the available workloads in NAS we chose those with the running time of at least ten seconds. We excluded applications whose CPU utilization was below 33%, because they were not affected by memory management policies. We also excluded applications using shared memory across processes (as opposed to threads) or memory-mapped files, because our replication mechanism does not yet support this behaviour. For *FaceRec*, we used two kinds of workloads: a short-running one and a long running one (named *FaceRecLong* in the remainder of the paper). The reason why we present two different workloads is that we found out that the behavior of *FaceRec* is drastically impacted by the size of the workload that is used. Each experiment was run at least three times. Overall, unless explicitly mentioned, we observed a standard deviation between 1% and 2%.

4.2 Single-application workloads

Performance comparison. Figures 5 and 6 show the execution time of the benchmarks under all configurations. For readability, we normalize the execution times to default Linux.

We can make two main observations. First, *Carrefour* almost systematically outperforms default Linux, AutoNUMA, and Manual interleaving, often quite substantially. For instance, when running *Streamcluster* or *FaceRecLong*, we observe that *Carrefour* is about 20% faster than Manual interleaving, about 40% faster than AutoNUMA, and about

⁸Note that we also ran experiments with the NUMAsched patchset for Linux [30]. However, the obtained results were significantly poorer than the ones of AutoNUMA.

twice as fast as default Linux. Second, we observe that, unlike other techniques, *Carrefour* never performs significantly worse than default Linux: the maximum performance degradation over Linux is below 3%, with the exception of *UA* from NAS, which slows down by 6%. In contrast, AutoNUMA and Manual Interleaving cause performance degradations of up to 40% and 25% respectively.

In order to understand the reasons for performance of different policies, we study in detail a set of applications whose performance is affected the most by *Carrefour*. To that end, we present several metrics: the load imbalance on memory controllers (Figure 7(a)), the load imbalance on interconnect links (Figure 7(b)), the average memory latency (Figure 8(a)) and the local access ratio (Figure 8(b)).

We draw the following observations. First, *Carrefour* much better balances the load on both memory controllers and interconnect links than Linux and AutoNUMA. Not surprisingly, Manual interleaving is also very good at balancing the load. Nevertheless, we observe in Figure 8(a) that *Carrefour* induces lower average memory latencies than Manual interleaving, which explains its better performance. To understand why *Carrefour* reduces memory latencies we refer to Figure 8(b), which shows that *Carrefour* not only balances the load on memory controllers and interconnect links, but also often induces a much higher ratio of local memory accesses than other techniques. This is a consequence of *Carrefour*'s judiciously applying the rights techniques (interleaving, replication or co-location) in places where they are beneficial. Interleaving mostly balances the load; replication and co-location in addition to balancing the load improve the local access ratio.

Looking inside *Carrefour*. To better understand the behavior of *Carrefour*, we show in Table 3 the number of replicated pages, the number of interleaved pages, and the number of co-located pages for the chosen benchmarks. These numbers provide a better insight into how *Carrefour* manages the memory. We observe that *all* three memory placement mechanisms are in use, but for some applications a single technique is dominant. *Streamcluster* and *SP* rely almost exclusively on a single technique (replication and co-location respectively). *Facesim*, *FaceRec*, *FaceRecLong* and *PCA* rely on two or more.

A legitimate question we can ask is whether *Carrefour* always selects the best technique. In Table 4, we report the performance improvement over Linux obtained when running a full-fledged *Carrefour* and when running a reduced version of *Carrefour* enabling only one technique at a time. We observe that *Carrefour* systematically selects the best technique. It is also interesting to remark how different techniques work together and how the numbers provided here echo those in Table 3. *Streamcluster* achieves the best performance with a single technique: replication, just as *SP* relies solely on co-location. The remaining benchmarks, and especially *PCA* need a range of techniques to obtain the ultimate performance.

4.3 Multi-application workloads

In this section, we study how *Carrefour* behaves in the context of workloads with multiple applications that are co-scheduled on the same machine (Machine B). The goal is to assess that *Carrefour* is able to work on complex access patterns and to make the distinction between the diverse

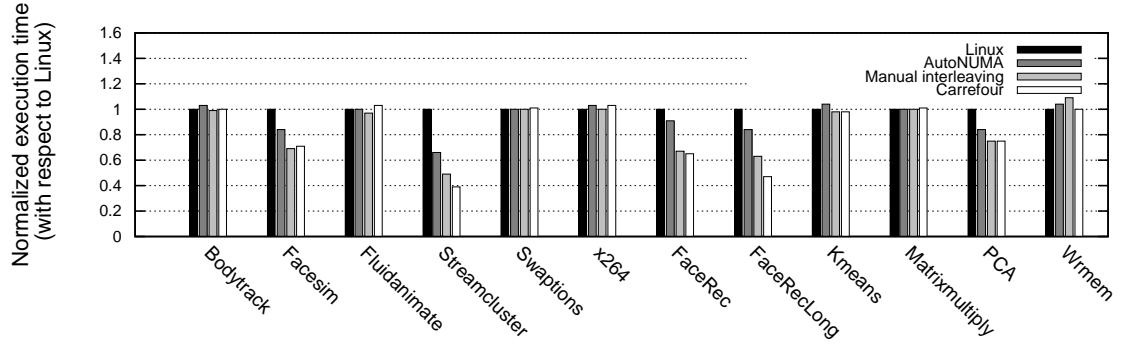


Figure 5: Parsec/Metis: AutoNUMA, Manual interleaving and Carrefour vs. Default Linux.

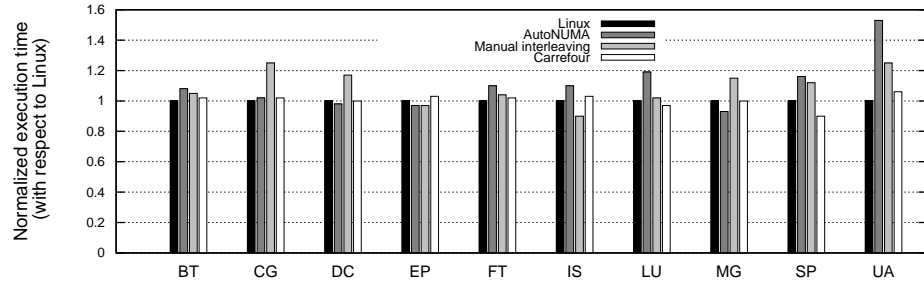


Figure 6: NAS: AutoNUMA, Manual interleaving and Carrefour vs. Default Linux.

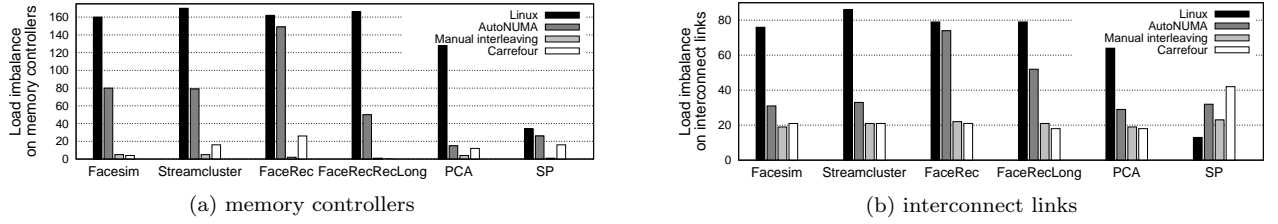


Figure 7: Load imbalance for selected single-application benchmarks.

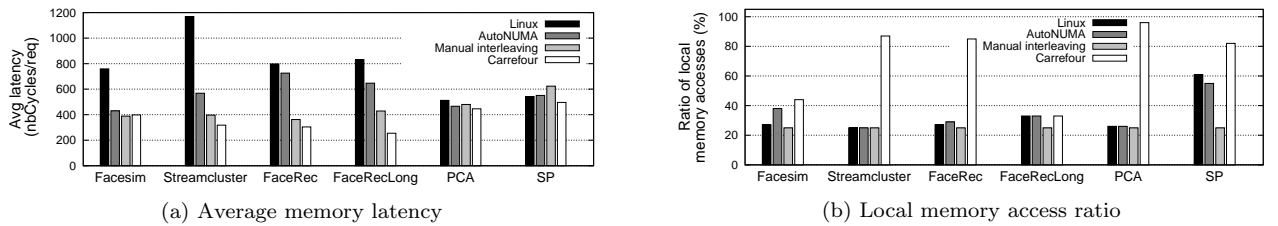


Figure 8: DRAM latency and locality for selected single-application benchmarks.

	Nb replicated pages	Nb interleaved pages	Nb migrated pages
FaceRec	3.9k	4	3.3k
FaceRecLong	3.9k	5	1.2k
Facesim	737	1.4k	12k
PCA	0	2.9k	81.3k
SP	0	0	2.1k
Streamcluster	25k	3	3k

Table 3: Number of memory pages that are replicated, interleaved and co-located.

	<i>Carrefour</i>	Replication	Interleaving	Co-location
FaceRec	35%	33%	22%	5%
FaceRecLong	53%	52%	32%	17%
Facesim	29%	12%	25%	27%
PCA	25%	0%	19%	8%
SP	10%	0%	0%	10%
Streamcluster	61%	61%	47%	34%

Table 4: Performance improvement (percentage of reduction in execution time) over Linux when running *Carrefour* and the three different techniques individually.

requirements of different applications. We consider several scenarios based on some of the applications previously studied in Section 4.2: (i) *FaceRec* + *PCA*, (ii) *Streamcluster* + *MG*, (iii) *PCA* + *MG*. Each application runs with 12 threads. We chose these scenarios because they exhibit interesting patterns, which require combining several memory placement techniques in order to achieve good performance.

Figure 9 shows, for each workload, the execution times of each application, comparing the different configurations: Linux, AutoNUMA, Manual interleaving and *Carrefour*. As in 4.2, we normalize the execution times to default Linux. We observe that, as in the case of single-application workloads, *Carrefour* outperforms Linux and AutoNUMA by up to 35% and 30% respectively. Manual interleaving and *Carrefour* yield close results. Manual interleaving outperforms *Carrefour* by 6% for *PCA*+*FaceRec*, but *Carrefour* outperforms interleaving by 3%—11% for the other workloads.

The reason why Manual interleaving performs relatively well in these scenarios is because there is a lot less inter-domain data sharing, and as a result a lot less traffic congestion. Hence there are fewer “problems” that need to be fixed. With two applications on four domains, each application uses two domains on average, so there is a lot less cross-domain traffic than in the single-application case. Hence there is a smaller discrepancy between the performance of different memory management algorithms.

To explain the results, we show the same detailed metrics as in Section 4.2: load imbalance on memory controllers, load imbalance on interconnect links, average DRAM latency and ratio of local DRAM accesses in Figures 10, 11, 12 and 13 respectively. The metrics are aggregated for the two applications of each workload. As was the case with the single-application workloads, we see that *Carrefour* systematically improves the latency of the studied workloads for two reasons: a more balanced load on memory controllers and interconnect links as well as an improved locality for DRAM accesses.

Finally, we show in Table 5 the performance improvement for each application when the effects of only one of the techniques are enabled. We observe, as we did in single-application experiments, that the multi-applications workloads perform better with an arsenal of techniques used in *Carrefour* rather than with any single technique alone.

4.4 Overhead

Carrefour incurs CPU and memory overhead. CPU overhead is a result of periodic IBS profiling. Memory overhead is a result of allocating data structures to keep track of profiling data. Memory overhead is negligible: e.g., 15MB on Machine B with 48GB of RAM. *Carrefour*’s data structures are pre-allocated on startup to avoid memory allocation during the runtime. We limit the number of profiled pages to 30,000

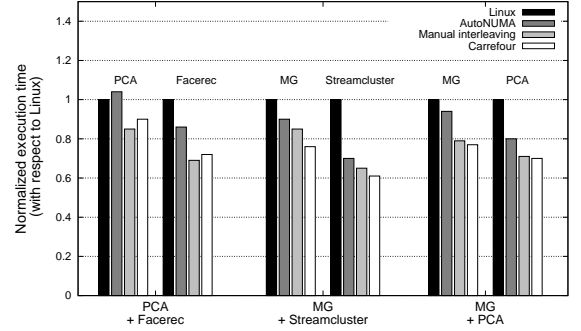


Figure 9: Multi-application workloads: AutoNuma, Manual interleaving and Carrefour vs. Default Linux.

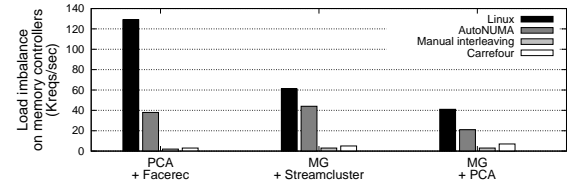


Figure 10: Multi-application workloads: load imbalance on memory controllers.

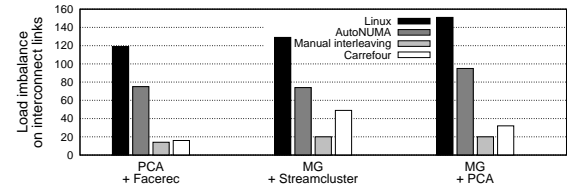


Figure 11: Multi-application workloads: load imbalance on interconnect links.

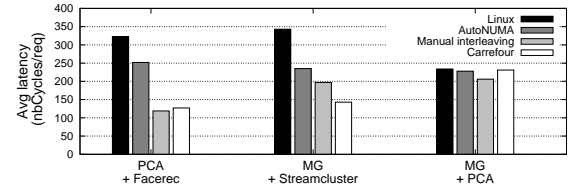


Figure 12: Multi-application workloads: average memory latency.

	Carrefour	Replication	Interleaving	Co-location
FaceRec + PCA	20% / 29%	3% / 0%	18% / 29%	18% / 1%
MG + Streamcluster	24% / 38%	19% / 39%	23% / 35%	24% / 37%
MG + PCA	27% / 30%	-3% / 11%	19% / 29%	15% / 25%

Table 5: Multi-application workloads: performance improvement (percentage of reduction in execution time) over Linux when running Carrefour and the three different techniques individually.

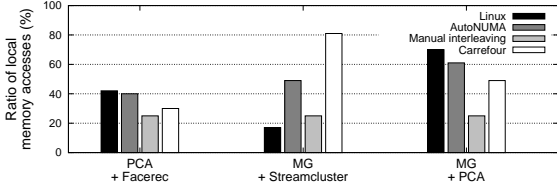


Figure 13: Multi-application workloads: local memory access ratio.

to avoid the cost of managing dynamically sized structures.

To measure CPU overhead we compared performance of *Carrefour* with Linux on those applications where *Carrefour* does not yield any performance benefits, but performs the profiling. We observed the overhead between 0.2% and 3.2%. Adaptive sampling rate in *Carrefour* is crucial to keeping this overhead low.

Replication can introduce significant overhead if we perform a lot of collapses. A single collapse costs a few hundred microseconds when it occurs in isolation. Parallel collapses can take a few milliseconds because of lock contention. That is why it is crucial to avoid collapses and other synchronization events by disabling replication for write-intensive and page-fault-heavy workloads, as done in *Carrefour*.

4.5 Discussion: hardware support

We have shown that a traffic management system like *Carrefour* can bring significant performance benefits. However, the challenge in building *Carrefour* was the need to navigate around the limitations of the performance monitoring units of our hardware as well as the costs of replicating pages. In this section, we draw some insights on the features that could be integrated into future machines in order to further mitigate the overhead and improve accuracy, efficiency and performance of traffic management algorithms.

First, *Carrefour* would benefit from hardware profiling mechanisms that sample memory accesses with high precision and low overhead. For instance, it would be useful to have a profiling mechanism that accumulates and aggregates page access statistics in an internal buffer before triggering an interrupt. In this regard, the AMD Lightweight Profiling [1] facility seems a promising evolution of profiling hardware⁹, but we believe the hardware should go even further, and not only accumulate the samples but be configured to aggregate them according to user needs, to reduce the number of interrupts even further.

Second, *Carrefour* would benefit from dedicated hardware support for memory replication. We believe that there should be interfaces allowing the operating system to indi-

cate to the processor which pages to replicate. The processor would then be in charge of replicating the pages on the nodes accessing it and maintaining consistency between the various replicas (in the same way as it maintains consistency for cache lines). Given that maintaining consistency between frequently written pages is costly, we believe that such processors should also be able to trigger an interrupt when a page is too frequently written. The OS would then decide to keep the page replicated or to revert the replication decision.

This hardware support can be made a lot more scalable than cache coherency protocols, because it is not automatic, but controlled by the OS, which, armed with better hardware profiling, will only invoke it for pages that perform very little write sharing. So the actual synchronization protocol would be triggered infrequently.

5. RELATED WORK

In this section, we explain how *Carrefour* relates to different works on multicore systems. First, we review systems aimed at maximizing data locality. Second, we contrast *Carrefour* with previous contention-aware systems. Third, we consider application-level techniques to mitigate contention on data-sharing applications. Finally, we discuss traffic characterization observations for modern NUMA systems.

Locality-driven optimizations.

NUMA-aware thread and memory management policies were proposed for earlier research systems [6, 9, 17, 31] as well as in commercial OS. Their main difference from our work is that their goal was to optimize *locality*. However, on modern systems, the main performance problems are due to traffic congestion. Our algorithm is the first one that meets the goal of mitigating traffic congestion. Among the above-mentioned works, the one most related to *Carrefour* is the system by Verghese et al. [31] for early cache-coherent NUMA machines, which leverages page replication and migration mechanisms. Their system relies on assumptions about hardware support that do not hold on currently available machines (e.g., precise per-page access statistics). Thus, *Carrefour*'s logic is more involved, as it is more difficult to amortize the costs of the monitoring and memory page placement mechanisms. The authors noticed that locality-driven optimizations could, as a side effect, reduce the overall contention in the system. However, their system does not systematically address contention issues. For instance, shared written pages are not taken into account, whereas *Carrefour* uses memory interleaving techniques when there is contention on such pages. Moreover, the load on memory controllers is ignored when making page replication/migration decisions.

Similarly to earlier NUMA-aware policies, Solaris and Linux focus primarily on co-location of threads and data, but to the disadvantage of data-sharing workloads, replication is not supported. Linux provides the option to interleave parts

⁹Unfortunately, Lightweight Profiling is only available on recent AMD processors (AMD Bulldozer) and we were not able to evaluate it in this work. We plan to study it in the future.

of the address space across all memory nodes, but the decision when to invoke the interleaving is left to the programmer or the administrator. Solaris supports the notion of a home load group, such that the thread's memory is always allocated in its home group and the thread is preferentially scheduled in its home group. This, again, favours locality, but does not necessarily address traffic congestion.

The recent AutoNUMA patches for Linux also implement locality-driven optimizations, along two main heuristics. First, threads migrate toward nodes holding the majority of the pages accessed by these threads. Memory residence is determined by page fault statistics. Second, pages are periodically unmapped from a process address space and, upon the next page fault, migrated to the requesting node. As shown in the evaluation section, this approach yields irregular results. We attribute this limitation to the following sources of overhead: local thread/page migration decisions that do not take data sharing patterns nor access frequencies into account (thus leading to page bouncing or useless migrations) nor MC/interconnect load (thus leading to memory load imbalance/congestion), continuous overhead due to the scanning/unmapping of page-table entries and the corresponding soft page faults. In contrast, *Carrefour* makes global data placement decisions based on precise traffic patterns and adjusts the monitoring overhead based on the observed contention level.

Locality-driven optimizations for data-sharing applications were addressed in a study that dynamically identified data-sharing thread groups and co-located them on the same node [29]. However, that solution was for a system with a centralized (UMA) memory architecture. Thus, it only studied the benefits of thread placement for improved cache locality and did not address the placement of memory pages on multiple memory nodes.

Contention management techniques.

Several recent studies addressed contention issues in the memory hierarchy. Some of these works were designed for UMA systems [15, 20, 32] and are inefficient on NUMA systems because they fail to address or even accentuate issues such as remote access latencies and contention on memory controllers and on the interconnect links [5]. Other works have been specifically designed for NUMA systems but only partially address contention issues. The N-Mass thread placement algorithm [18] attempts to achieve good DRAM locality while avoiding cache contention. However, it does not address contention issues at the level of memory controllers and interconnect links. Two studies [3, 19] have shown the importance of taking memory controller congestion into account for data placement decisions, but they did not provide a complete solution to address multi-level resource contention. The most comprehensive work to date on NUMA-aware contention management is the DINO scheduler [5], which spreads memory intensive threads across memory domains and accordingly migrates the corresponding memory pages. However, DINO does not address workloads with data sharing between threads or processes, which require identifying per-page memory access patterns and making the appropriate data placement decisions.

No-sharing application design principle.

When introducing a new resource-management policy in the OS, it is worth asking whether a similar or better effect

could be achieved by restructuring the application. In our context, it is important to consider the so-called *no-sharing* principle of application design. The key idea behind no-sharing is that the data must be partitioned or replicated between memory nodes, and a thread needing to access data in a different domain than its own either migrates to the target domain or asks the thread running in that domain to perform the work on its behalf, instead of fetching the data over the memory channels [4, 8, 13, 22, 24, 27, 28]. While the no-sharing architecture was primarily motivated by the need to avoid locking, it could similarly help reduce the amount of traffic sent across the interconnect, and thus alleviate the traffic congestion problem.

Unfortunately, no-sharing architectures are not a universal remedy. First of all, they trade-off data accesses for messages or thread migration; the trade-off is only worth making if the size of the data used in a single operation is much larger than the size of the message or the state of the migrated thread [8]. Second, adopting a no-sharing architecture often requires very significant changes to the application (and to the OS, if the application is OS-intensive). A good illustration of the potential challenges can be gleaned from two studies that converted a database system to the no-sharing design. The first study took the path of least resistance and simply replicated and/or partitioned the database among domains, adding a message-routing layer on top [27]. While this worked well for small read-mostly workloads, for large workloads replication had very significant memory overhead (unacceptable because of increased paging), and partitioning required *a priori* knowledge of query-to-data mapping, which is not a reasonable assumption in a general case. A solution that overcame these limitations, DORA [24], required a very significant restructuring of the database system, which could easily amount to millions of dollars in development costs for a commercial database.

Our goal was to address scenarios where adopting a no-sharing architecture is not feasible either for technical reasons or for practical considerations. Providing an OS-level, rather than an application-level, solution allows us to address many applications at once. Understanding the limitations of the OS-level solution and determining what optimizations can be done only at the level of an application is an open research question. Although we did provide some comparison of *Carrefour* with application-level techniques, a deeper analysis of this problem would be worthwhile.

Traffic characterization on modern systems.

A recent study characterized the performance of emerging "scale-out" workloads on modern hardware [12]. The authors observed that there is little inter-thread data sharing, and the utilization of the off-chip memory bandwidth is low. As a result, they argue that memory bandwidth on existing processors is unnecessarily high. Our findings do not agree with this observation. Although it is also true that the workloads we consider perform very little fine-grained data sharing, they still stress the cross-chip interconnect, because they access a large working set, which is spread across the entire NUMA memory space. The authors of the scale-out study reported a very low bandwidth utilization (<10%), even for database workloads. In contrast, our measurements show a utilization as low as 45% in some cases. These differences could be because the authors of [12] used a system with only two chips and 12 cores in their experiments. On larger

systems, more threads are making requests to remote memories and so there is greater pressure on bandwidth. Further, we found that low bandwidth utilization is not necessarily a healthy symptom. In our experiments, performance dropped steeply even as bandwidth utilization went from 10% to 30%. The interconnect became the bottleneck even at fraction of the total available bandwidth! The reason is that memory requests are not spread evenly in time; they are bursty. Burstiness causes requests to clash on the link even if the overall bandwidth is not exceeded. In summary, we conclude that contrary to the suggestion made in [12], it is too early to reduce the bandwidth of cross-chip interconnects on large multicore systems, especially if they are used for running large data-centric workloads.

6. CONCLUSION

We presented *Carrefour*, a new memory management algorithm for NUMA systems that manages traffic on memory controllers and interconnect links. Other NUMA-aware memory management policies optimized locality of memory accesses only, and their performance falls short of *Carrefour*. *Carrefour*'s design was motivated by bottleneck shift in modern hardware, where remote wire delays play a lot smaller role in performance than traffic congestion.

System design principles embodied in *Carrefour* are important not only for today's systems, but also for future hardware. The amount of memory bandwidth per core is projected to decrease in the future [8], so managing interconnect congestion will be as crucial as ever. Furthermore, *energy cost* of remote accesses remains higher than that of local accesses [11], so a solution that reduces remote latencies will be the cornerstone of energy-efficient system design.

7. REFERENCES

- [1] AMD64 Technology Lightweight Profiling Specification, Aug. 2010. http://support.amd.com/us/Processor_TechDocs/43724.pdf.
- [2] AutoNUMA: the other approach to NUMA scheduling. LWN.net, Mar. 2012. <http://lwn.net/Articles/488709/>.
- [3] M. Awasthi, D. Nellans, K. Sudan, et al. Handling the problems and opportunities posed by multiple on-chip memory controllers. In *PACT*, 2010.
- [4] A. Baumann, P. Barham, P.-E. Dagand, et al. The Multikernel: A New OS Architecture for Scalable Multicore Systems. In *SOSP*, 2009.
- [5] S. Blagodurov, S. Zhuravlev, M. Dashti, et al. A case for numa-aware contention management on multicore systems. In *USENIX ATC*, 2011.
- [6] W. Bolosky, R. Fitzgerald, and M. Scott. Simple but Effective Techniques for NUMA Memory Management. In *SOSP*, 1989.
- [7] S. Boyd-Wickizer, H. Chen, R. Chen, et al. Corey: an operating system for many cores. In *OSDI*, 2008.
- [8] S. Boyd-Wickizer, M. F. Kaashoek, R. Morris, et al. A software approach to unifying multicore caches. Technical Report MIT-CSAIL-TR-2011-032, 2011.
- [9] T. Brecht. On the Importance of Parallel Application Placement in NUMA Multiprocessors. In *USENIX SEDMS*, 1993.
- [10] CSU Face Identification Evaluation System. <http://www.cs.colostate.edu/evalfacerec/index10.php>.
- [11] B. Dally. Power, programmability, and granularity: The challenges of exascale computing. <http://techtalks.tv/talks/54110>.
- [12] M. Ferdman, A. Adileh, O. Kocberber, et al. Clearing the clouds: A study of emerging scale-out workloads on modern hardware. In *ASPLOS*, 2012.
- [13] B. Gamsa, O. Krieger, and M. Stumm. Tornado: Maximizing Locality and Concurrency in a Shared Memory Multiprocessor Operating System. In *OSDI*, 1999.
- [14] A. Kamali. Sharing aware scheduling on multicore systems. In *MSc Thesis, Simon Fraser Univ.*, 2010.
- [15] R. Knauerhase, P. Brett, B. Hohlt, et al. Using OS Observations to Improve Performance in Multicore Systems. *IEEE Micro*, 28(3):pp. 54–66, 2008.
- [16] R. Lachaize, B. Lepers, and V. Quéma. Memprof: A memory profiler for numa multicore systems. In *USENIX ATC*, 2012.
- [17] R. P. LaRowe, Jr., C. S. Ellis, et al. Evaluation of NUMA Memory Management Through Modeling and Measurements. *IEEE TPDS*, 3:686–701, 1991.
- [18] Z. Majo and T. R. Gross. Memory management in numa multicore systems: Trapped between cache contention and interconnect overhead. In *ISMM*, 2011.
- [19] Z. Majo and T. R. Gross. Memory system performance in a numa multicore multiprocessor. In *SYSTOR*, 2011.
- [20] A. Merkel, J. Stoess, and F. Bellosa. Resource-Conscious Scheduling for Energy Efficiency on Multicore Processors. In *EuroSys*, 2010.
- [21] Metis MapReduce Library. <http://pdos.csail.mit.edu/metis/>.
- [22] Z. Metreveli, N. Zeldovich, and F. Kaashoek. Cphash: a cache-partitioned hash table. In *PPoPP*, 2012.
- [23] NAS Parallel Benchmarks. <http://www.nas.nasa.gov/publications/npb.html>.
- [24] I. Pandis, R. Johnson, N. Hardavellas, et al. Data-oriented transaction execution. *Proc. VLDB Endow.*, 3:928–939, September 2010.
- [25] PARSEC Benchmark Suite. <http://parsec.cs.princeton.edu/>.
- [26] A. Pesterev, N. Zeldovich, and R. T. Morris. Locating cache performance bottlenecks using data profiling. In *EuroSys*, 2010.
- [27] T.-I. Salomie, I. E. Subasu, J. Giceva, et al. Database engines on multicores, why parallelize when you can distribute? In *EuroSys*, 2011.
- [28] X. Song, H. Chen, R. Chen, et al. A case for scaling applications to many-core with os clustering. In *EuroSys*, 2011.
- [29] D. Tam, R. Azimi, and M. Stumm. Thread Clustering: Sharing-Aware Scheduling on SMP-CMP-SMT Multiprocessors. In *EuroSys*, 2007.
- [30] Toward better NUMA scheduling. Linux Weekly News, Mar. 2012. <http://lwn.net/Articles/486858/>.
- [31] B. Verghese, S. Devine, A. Gupta, et al. Operating system support for improving data locality on cc-numa compute servers. In *ASPLOS*, 1996.
- [32] S. Zhuravlev, S. Blagodurov, and A. Fedorova. Addressing Contention on Multicore Processors via Scheduling. In *ASPLOS*, 2010.