

Application-to-Core Mapping Policies to Reduce Interference in On-Chip Networks

Reetuparna Das[§] Onur Mutlu[†] Akhilesh Kumar[§] Mani Azimi[§]
§Intel Labs †Carnegie Mellon University

SAFARI Technical Report No. 2011-001

May 25, 2011

Abstract

Future many-core processors are likely to concurrently execute a large number of diverse applications. How these applications are mapped to cores largely determines the interference between these applications in critical shared resources such as the network-on-chip. In this paper, we propose application-to-core mapping policies to reduce the contention in network-on-chip and memory controller resources and hence improve overall system performance. The key ideas of our policies are to: 1) map network-latency-sensitive applications to separate node clusters in the network from network-bandwidth-intensive applications such that the former makes fast progress without heavy interference from the latter, 2) map those applications that benefit more from being closer to the memory controllers close to these resources. Contrary to the conventional wisdom of balancing network or memory load across the network-on-chip and controllers, we observe that it is also important to ensure that applications that are more sensitive to network latency experience little interference from applications that are network-bandwidth-intensive, even at the cost of load imbalance.

We evaluate the proposed application-to-core mapping policies on a 60-core system with an 8x8 mesh NoC using a suite of 35 diverse applications. Averaged over 128 randomly generated multiprogrammed workloads, the final proposed policy improves system throughput by 16.7% in terms of weighted speedup over a state-of-the-art baseline, while also reducing system unfairness by 22.4% and average interconnect power consumption by 52.3%.

1 Introduction

Landscape of computing is undergoing a sea change with the advent of multi-core processors. Parallelism is now all pervasive, from cloud services to laptops to hand-held devices. Unlike high-performance computing environments, where a single parallel application runs on all cores, a large fraction of future computing systems are expected to run many diverse applications competing for shared resources. Thus, in the era of many-core processors, managing shared resources among co-scheduled interfering applications is one of the most fundamental challenges we face. The Network-on-Chip (NoC) is a *critical shared resource* and its effective utilization is essential for improving overall system performance and fairness.

We observe that in a large many-core processor, which processor core is selected to execute an application could have a significant impact on system performance. Performance of an application critically depends how its network packets *interfere* with other applications' packets in the interconnect and on how

far away it is scheduled from shared resources such as memory controllers. Hence, the core selection policy has to be aware of the spatial geometry of a many-core processor as well as applications' interference characteristics.

While prior research [21, 27, 30] has tackled the problem of how to map tasks/threads *within* an application, the interference behavior *between* applications in the NoC is less well understood. Unfortunately, current operating systems are unaware of the network topology and application interference characteristics at any instant of time, and employ naive methods while mapping applications to cores. For instance, Linux 2.6.x [1] assigns a static numbering to cores and chooses the numerically smallest numbered core when allocating an idle core to an application. This leads to an application-to-core mapping that is oblivious to application characteristics and inter-application interference. This causes two major problems we aim to solve in this paper. First, applications that interfere significantly with each other can get mapped to closeby cores, thereby causing significant interference to each other in both the network-on-chip and the memory controllers, reducing overall system performance. Second, an application may benefit significantly from being mapped to a core close to a shared resource (e.g., a memory controller), yet it can be mapped far away from that resource (while another application that does not benefit from being close to the resource is mapped closeby to the resource), again reducing system performance. If, on the other hand, the operating system took into account the interconnect topology and application performance/interference characteristics when deciding where to map applications in a network-on-chip, it can 1) significantly reduce the destructive network and memory interference between applications (by mapping them far away from each other), 2) map close to memory controllers those applications that benefit the most from being close to memory controllers, thereby improving overall system performance.

In this paper, we develop intelligent application-to-core mapping policies to reduce inter-application interference in the NoC and thus improve system performance. Our policies are built upon two major observations. First, we observe that some applications are more *sensitive* to interference than others: network-latency-sensitive applications slow down more significantly when interfered with than others. Thus, system performance can be improved by separating network-sensitive applications from aggressive applications which have high demand for network bandwidth. To allow this separation of applications, we partition the network into clusters using memory data placement techniques, develop heuristics to estimate each application's network sensitivity, and devise algorithms that use these estimates to distribute applications to clusters. While partitioning applications between clusters to reduce interference, our algorithms also try to ensure that network load is balanced among clusters as much as possible.

Second, we observe that some applications benefit significantly more from being placed close to memory controllers than others: an application that is both memory-intensive and network-latency-sensitive gains more performance from being close to a memory controller than one that does not have either of these properties. Thus, system performance can be improved by mapping such applications to cores that are close to memory controllers. To this end, we develop heuristics to identify such applications dynamically and devise a new algorithm that maps applications to cores within each cluster based on each application's performance sensitivity to distance from the memory controller.

We make the following new contributions in this paper:

- We develop a novel core assignment algorithm for applications in a network-on-chip based many-core system, which aims to minimize destructive inter-application network interference and thereby maximize overall system performance. To our knowledge, this is the first work that takes into account interference characteristics of applications in the on-chip network when assigning cores to applications.
- We develop new insights on application interference in NoC based systems, which form the foundation of our algorithm. In particular, we demonstrate that 1) mapping network-sensitive applications to parts of

the network such that they do not interfere with network-intensive applications and 2) mapping memory-intensive and network-sensitive applications close to the memory controller can significantly improve performance. Contrary to conventional wisdom [12, 18, 22], we show that sometimes reducing network load balance to isolate network-sensitive applications can be beneficial.

- We demonstrate that intelligent application-to-core mappings can save network energy significantly. By separating applications into clusters and placing memory-intensive applications closer to the memory controllers, our techniques reduce the communication distance and hence communication energy of applications that are responsible for the highest fraction of overall network load.
- We extensively evaluate the proposed application-to-core mapping policies on a 60-core CMP with an 8x8 mesh NoC using a suite of 35 diverse applications. Averaged over 128 randomly generated multiprogrammed workloads, our final proposed policy improves system throughput by 16.7% in terms of weighted speedup over a state-of-the-art baseline, while also reducing application-level unfairness by 22.4% and NoC power consumption by 52.3%. We also show that our policies are complementary to application-aware network packet prioritization techniques [13] and can be combined with them to further improve system performance.

2 Background

In this section, we provide a brief background on NoC architectures. For an in-depth introduction to NoCs, we refer the reader to [12].

Router: A generic NoC router architecture is illustrated in Figure 1. The router has P input and P output channels/ports; typically $P = 5$ for a 2D mesh, one from each direction and one from the network interface (NI). The Routing Computation unit, RC, is responsible for determining the next router and the virtual channel within the next router for each packet. The Virtual channel Arbitration unit (VA) arbitrates amongst all packets requesting access to the same VCs and decides on winners. The Switch Arbitration unit (SA) arbitrates amongst all VCs requesting access to the crossbar and grants permission to the winning packets/flits. The winners are then able to traverse the crossbar and are placed on the output links. The baseline packet scheduling policy we model is round-robin arbitration across all virtual channels [12], but we also show that our proposal works well with recently proposed application-aware packet scheduling policies [13].

Network Transactions: In the many-core processor architecture we study in this paper, the NoC connects the core nodes (a CPU core with its private caches) and the on-chip memory controllers. Figure 2 shows the layout of the many core processor with a 8x8 mesh. All tiles have a router. The memory controllers (triangles) are placed in the corner tiles. All other tiles have a CPU core, private L1 and private L2 cache. The core nodes send request packets to on-chip memory controller nodes via the NoC and receive response

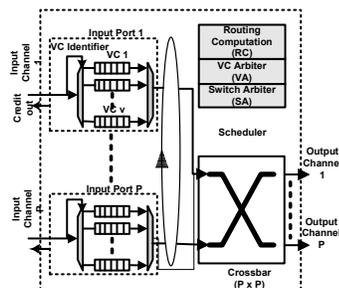


Figure 1: Generic NoC router

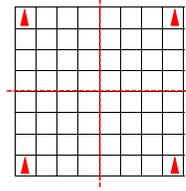


Figure 2: An 8x8 mesh network-on-chip divided into four clusters. Each cluster has a memory controller (triangle) tile.

data packets from the memory controller, once data has been fetched from off-chip DRAMs.

Interference in NoC: Each packet spends at least 2-4 cycles at each router depending on the number of stages in the router pipeline. In addition to the router pipeline stages, a packet can spend *many cycles* waiting in a router, competing for buffers or switch with other packets. Thus, an application's packet may be blocked in the network due to interference from other applications' packets. While its packets are buffered in remote routers, the application (running on the core node) stalls waiting for its packets to return.

3 Motivation

A many-core processor with n cores can run n concurrent applications. Each of these applications can be mapped to any of n cores. Thus, there can be $n!$ possible mappings. From the interconnect perspective, a application-to-core mapping can determine the degree of interference of an application with other applications in the NoC as well as how well the overall network load is balanced. The application-to-core mapping also determines how the application's memory accesses are distributed among the memory controllers. These different factors can lead to variation in performance. For example, Figure 3 shows the system performance for 576 different application-to-core mappings for the same workload (detailed system configuration is given in Section 6). The workload consists of 10 copies each of applications `gcc`, `barnes`, `soplex`, `lbm`, `milc` and `leslie` running together. The Y-axis shows the system performance in terms of weighted speedup of the different mappings (higher is better). Each blue dot represents a different mapping. It can be seen that the best possible mapping provides 1.6X higher system performance than the worst possible mapping. **Our goal** is to devise new policies that can find a application-to-core mapping to maximize system performance.

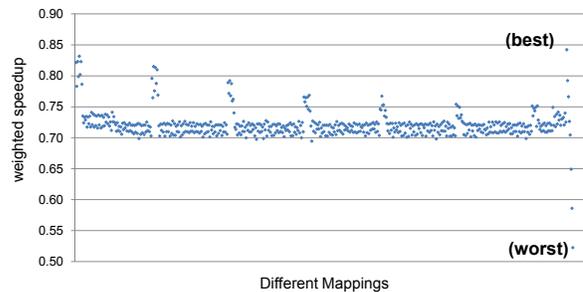


Figure 3: Performance of 576 different application to core mappings for one multiprogrammed workload

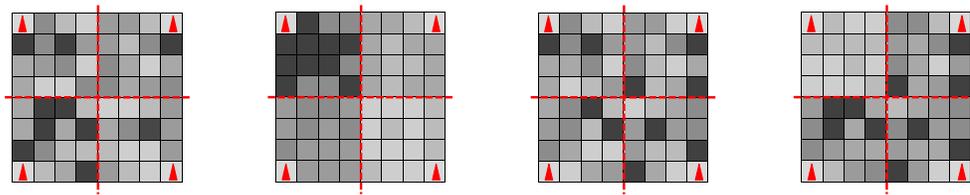


Figure 4: Inter-Cluster Application-to-Core Mapping Examples: (a) Random (RND) (b) Imbalanced (IMBL) (c) Balanced (BL) (d) Balanced with Reduce Interference (BLRI)

4 Application-to-Core Mapping Policies

Our goal is to improve system throughput by designing efficient application-to-core mapping policies that reduce interference in on-chip-network based multi-core systems. To enable this, we partition the cores into clusters and develop two novel techniques: 1) algorithms to distribute applications between clusters, and 2) algorithms to map applications to cores within a cluster.

4.1 Cluster Formation

What is a cluster? We define a cluster as a sub-network such that the majority of network traffic originating in the sub-network can be constrained within the sub-network. *Clustering factor* is defined as the percentile of accesses that can be constrained within the cluster. As a first step, we form clusters to partition network traffic between clusters. *Forming clusters ensures that applications mapped to different clusters interfere minimally with each other.* Figure 2 shows a many-core processor with cores organized in an 8x8 mesh on-chip network. We partition the tiles of the on-chip network such that each cluster has one memory controller (the *home memory controller* of that cluster). The dotted lines show the division of the on-chip network into four clusters. In addition to reducing interference among applications, clustering also improves communication locality since the applications in the cluster communicate mainly with the memory controller in the cluster. This has two positive effects: 1) it reduces overall congestion in the network, 2) it reduces the average distance packets need to traverse (to get to the memory controller), thereby reducing packet latency and network energy consumed per packet.

How to enforce clustering for memory accesses? Clustering can be achieved by mapping physical pages requested by cores to memory controllers in an appropriate manner. Typically, physical pages (or even cache blocks) are *interleaved* among memory controllers such that adjacent pages (or cache blocks) are mapped to different memory controllers [37, 32, 25]. To enable clustering, page allocation and replacement policies should be modified such that data requested by a core is opportunistically mapped to the *home memory controller (home MC)* of the core. To achieve this, we slightly modify the commonly-used CLOCK [23] page allocation and replacement algorithm to what we call the *cluster-CLOCK* algorithm.

When a page fault occurs and free pages exist, the operating system gives *preference* to free pages belonging to the *home MC* of a requesting core when allocating the new page to the requesting core. If no free pages belonging to the *home MC* exist, a free page from another MC is allocated. When a page fault occurs and no free pages exist, *preference* is given to a page belonging to the *home MC*, while finding the replacement page candidate. We look N pages beyond the default candidate found by CLOCK [23] algorithm to find a page which belongs to *home MC*¹. If unsuccessful to find a replacement candidate belonging to *home MC* in N pages beyond the default candidate, the algorithm simply selects the default

¹We use N=512 in our experiments, a value empirically determined to maximize the possibility of finding a page in home MC while minimizing the overhead of searching for one.

candidate for replacement.

The above modifications ensure that the new page replacement policy does not perturb significantly the existing replacement order, and at the same time achieves the effect of clustering opportunistically. Note that these modifications to virtual memory management (for both page allocation and page replacement) *do not enforce a static partitioning of DRAM memory capacity*; they only *bias* the page replacement policy such that it likely allocates pages to a core from the core’s home memory controller.

How to enforce clustering for cache accesses? Clustering can be enforced for cache accesses for shared cache architectures by slightly modifying state-of-the-art cooperative caching techniques [7, 34] or page coloring techniques [9] such that cache accesses of a core remain within the shared cache slices within the core’s cluster. Such caching techniques have been shown to improve system performance by trading-off cache capacity for communication locality. Since our focus is to reduce interference for NoCs, for this work we assume private caches which automatically enforces clustering by restricting the cache accesses from an application to the private cache co-located in the same tile as the core running the application.

4.2 Mapping Policy between Clusters

In this subsection, we devise clustering algorithms that decide *to which cluster* an application should be mapped to. The key idea is to map the network-sensitive applications to a separate cluster such that they suffer minimum interference. While doing so, we also try to ensure that overall network load is balanced between clusters as much as possible.

Before we develop our algorithms, we briefly describe two possible inter-cluster mappings, which we will contrast our algorithms to the baseline. Figure 4 illustrates four different example inter-cluster mappings for a particular abstract workload. The former two, (a) and (b), are baseline mappings and the latter two, (c) and (d), are possible results of our algorithms. Each core tile in the figure is shaded according to network intensity of the application running on it; a darker tile corresponds to an application with a higher private cache miss rate, i.e. Misses per Thousand Instructions (MPKI), running on the core in the tile. Figure 4 (a) shows a possible random mapping of applications to clusters (called RND), which is the baseline mapping policy in existing general-purpose systems. Figure 4 (b) shows a mapping where applications are completely imbalanced between clusters (called IMBL). The upper left cluster is heavily loaded while the lower right cluster has very low load. An imbalanced inter-cluster mapping can severely degrade system performance because of poor utilization of aggregate available bandwidth in the NoC and in the off-chip memory channels. We investigate this policy for solely to provide an evaluation of how much performance can degrade with an undesirable mapping, which can possibly happen with existing policies.

We develop two inter-cluster mapping policies. The goal of the *Balanced (BL)* policy, an example of which is shown in Figure 4 (c), is to balance load between clusters. The goal of the *Balanced with Reduced Interference (BLRI)* policy, an example of which is shown in Figure 4 (d), is to protect interference-sensitive applications from other applications by assigning them their own cluster (the top left cluster in the figure) while trying to keep load balance in the rest of the network as much as possible.

4.2.1 Balanced Mapping (BL): Mapping to Balance Load between Clusters

The first goal of our inter-cluster mapping policy is to balance the load among clusters such that there is better overall utilization of network and memory bandwidth. To form an application-to-cluster mapping singularly optimized to balance load, we start with the list of applications *sorted with respect to their network intensity*: $\{A_0, A_1, \dots, A_{n-1}\}$ where A_i is i^{th} highest intensity application and n is number of cores

in the chip. Network intensity is measured in terms of injection rate into the network, which is quantified by private cache Misses per Thousand Instructions (MPKI), collected periodically at run-time or provided statically before the application starts. The basic idea is to map consecutively-ranked applications to consecutive clusters in a round robin fashion. Hence, applications are mapped to clusters as follows: $C_i = \{A_i, A_{i+k}, A_{i+2k}, A_{i+3k}, \dots\}$ where C_i is the set of applications mapped to the i^{th} cluster, and k is the number of clusters. Figure 4 (c) shows an example mapping after applying the inter-cluster BL mapping algorithm to a workload mix.

4.2.2 Reduced Interference Mapping (RI): Mapping to Reduce Interference for Sensitive Applications

We observe that some applications are more sensitive to interference in the network compared to other applications. In other words, the relative slowdown these applications experience due to the same amount of interference in the shared NoC is higher than that other applications experience. Our insight is that if we can separate such sensitive applications from other network-intensive ones by mapping them to their own cluster, we can improve performance significantly because the slowdowns of the most sensitive applications can be kept under control. To achieve this, our mechanism: 1) identifies sensitive applications 2) maps the sensitive applications to their own cluster *only if* there are enough such applications *and* the overall load in the network is high enough to cause congestion.

How to Identify Sensitive Applications We characterize applications to investigate their relative sensitivity. Our studies show that interference-sensitive applications have two main characteristics. First, they have low memory level parallelism (MLP) [16, 31]: such applications are in general more sensitive to interference because any delay for the application’s network packet likely results in extra stalls, as there is little or no overlap of packet latencies. Second, they inject enough load into the network for the network interference to make a difference in their execution time. In other words, applications with very low network intensity are not sensitive because their performance does not significantly change due to extra delay in the network.

We use two metrics to identify interference-sensitive applications. We find that Stall Cycles Per Miss (STPM) metric correlates with memory level parallelism (MLP). STPM is the average number of cycles for which a core is stalled because it is waiting for a cache miss packet to return from the network. Relative STPM is an application’s STPM value normalized to the minimum STPM among all applications to be mapped. Applications with high relative STPM are likely to have relatively low MLP. Such applications are classified as sensitive only if they inject enough load into network, i.e., if their private cache misses per thousand instructions (MPKI) is greater than a threshold. Algorithm 1 formally summarizes how our technique categorizes applications as sensitive.

How to Decide Whether or Not to Allocate a Cluster for Sensitive Applications After identifying sensitive applications, our technique tests if a separate cluster should be allocated for them. This cluster is called the $RI_{cluster}$, which stands for *Reduced-Interference Cluster*. There are three conditions that need to be satisfied for this cluster to be formed:

- First, there have to be enough sensitive applications to fill at least $R\%$ of the cores in a cluster. *This condition ensures that there are enough sensitive applications such that their separation from others actually reduces interference significantly.* We empirically found $R = 75$ is a good threshold.
- Second, the entire workload should exert a large amount of pressure on the network. We found that allocating interference-sensitive applications to their own cluster makes sense only for workloads that have a mixture of interference-sensitive applications and network-intensive (high-MPKI) applications

Algorithm 1: Algorithm to identify sensitive applications

Input: $Applications = \{A_0, A_1, \dots, A_{n-1}\}$ such that
 $\forall i \text{ } STPM(A_i) \geq STPM(A_{i+1})$
 $Thresh_{MPI_{Low}}, Thresh_{MPI_{High}}, Thresh_{SensitivityRatio}$

```

1: begin
2:    $SensitiveApplications \leftarrow \emptyset$ 
   // The STPM of all applications ( $A_0$  to  $A_{n-1}$ ) is normalized to  $MinSTPM$ 
3:    $MinSTPM \leftarrow \text{minimumof}STPM(Applications)$ 
4:   for  $i \leftarrow 0$  to  $n - 1$  do
5:     if  $NetworkDemand(A_i) > Thresh_{MPI_{Low}}$  and  $NetworkDemand(A_i) < Thresh_{MPI_{High}}$  and
6:      $STPM(A_i)/MinSTPM \geq Thresh_{SensitivityRatio}$  then
7:        $SensitiveApplications \leftarrow SensitiveApplications \cup A_i$ 
8:     end
9:   end
10: end

```

that can cause severe interference by pushing the network towards saturation. *Therefore, we consider forming a separate cluster only for very intensive workloads that are likely to saturate the network.* As a result, our algorithm considers forming an $RI_{cluster}$ if the aggregate bandwidth demand of the entire workload is higher than 1500 MPKI.²

- Third, the aggregate MPKI of the $RI_{cluster}$ should be small enough so that interference-sensitive applications mapped to it do not significantly slow down each other. If separating applications to an $RI_{cluster}$ ends up causing too much interference within the $RI_{cluster}$, this would defeat the purpose of forming the $RI_{cluster}$ in the first place. To avoid this problem, our algorithm does not form an $RI_{cluster}$ if the aggregate MPKI of $RI_{cluster}$ exceeds the bandwidth capacity of any NoC channel.³

If these three criteria are not satisfied, Balanced Load (BL) algorithm is used to perform mapping in all clusters, without forming a separate $RI_{cluster}$.

How to Map Sensitive Applications to Their Own Cluster Algorithm 2 (Reduced Interference) illustrates how to form a separate cluster ($RI_{cluster}$) for sensitive applications. The goal of the Reduced Interference (RI) algorithm is to fill the $RI_{cluster}$ with as many sensitive applications as possible, as long as the aggregate MPKI of $RI_{cluster}$ does not exceed the capacity of any NoC channel. The problem of choosing p (p = number of cores in a cluster) sensitive applications that have an aggregate MPKI less than an upper bound, while maximizing aggregate sensitivity of $RI_{cluster}$ can be easily shown to be equivalent to the *0-1 knapsack problem* [10]. We use a simple solution described in [10] to the knapsack problem to choose p sensitive applications. In case there are fewer sensitive applications than p , we pack the empty cores in the $RI_{cluster}$ with the *insensitive applications that are the least network-intensive*.⁴

Forming More than One $RI_{cluster}$ If the number of sensitive applications identified by Algorithm 1 is more than $2p$, then our technique forms two $RI_{cluster}$ s.

²Each memory controller feeds the cores with two outgoing NoC channels (each channel @32 GB/s, total 64GB/s). The off-chip memory channels are matched to the network channel capacity (4 channels @16 GB/s, total 64GB/s). The total bandwidth demand of 64 GB/s translates to an MPKI of 500 MPKI (assuming 64 byte cache lines, throughput demand of one memory access per cycle and core frequency of 2 GHz). We consider forming an RI cluster if the bandwidth demand is $R=0.75$ times the cluster capacity of 64 GB/s (or 375 MPKI). With 4 clusters, the total demand of the workload should be $4*375=1500$ MPKI. We empirically validated this threshold.

³Each NoC channel has a capacity of 32 GB/s. Assuming 64 byte cache lines and throughput demand of one memory access per cycle at core frequency of 2 GHz, 32 GB/s translates to 250 MPKI. Thus the aggregate MPKI of $RI_{cluster}$ should be less than 250 MPKI ($Thresh_{MPKIRI} = 250$ MPKI).

⁴Note that this solution does not have high overhead our algorithm is invoked at long time intervals at the granularity of OS time quanta.

Algorithm 2: Reduced-Interference (RI) Mapping: Algorithm to form the $RI_{cluster}$ for Sensitive Applications

Input: $SensitiveApplications = \{S_0, S_1, \dots, S_{k-1}\}$
 $FillerApplications = \{F_0, F_1, \dots, F_{l-1}\}$ such that
 $\forall i \ NetworkDemand(F_i) \leq NetworkDemand(F_{i+1})$ and $F_i \notin SensitiveApplications$
 $Thresh_{MPKI-RI} (= 250MPKI), p = number\ of\ cores\ in\ a\ cluster$

```

1: begin
2:    $RI_{cluster} \leftarrow \emptyset$   $MaxWeight \leftarrow Thresh_{MPKI-RI}$ 
3:   for  $i \leftarrow 0$  to  $k - 1$  do
4:      $Weight_i \leftarrow MPKI(S_i)$ 
5:      $Value_i \leftarrow STPM(S_i)/MinSTPM$ 
6:   end
7:    $RI_{cluster} \leftarrow \mathbf{KnapSack}(Weight, Value, k, MaxWeight)$ 
8:   for  $i \leftarrow 0$  to  $p - |RI_{cluster}| - 1$  do
9:      $RI_{cluster} \leftarrow RI_{cluster} \cup F_i$ 
10:  end
11: end

```

Once the $RI_{cluster}$ has been formed, the rest of the applications are mapped to the remaining clusters using the BL algorithm. **We call this final inter-cluster mapping algorithm, which combines RI and BL algorithms, as the BLRI algorithm.** Figure 4 (d) shows an example mapping after applying the inter-cluster BLRI mapping algorithm for a workload mix.

4.3 Mapping Policy within a Cluster

After mapping each application to a cluster, a question remains: *which core within a cluster* should an application be mapped to? Figure 5 shows different possible intra-cluster mappings for a single cluster. Figure 5 (a) depicts a random intra-cluster mapping; this is not the best intra-cluster mapping as it is agnostic to application characteristics. We observe that 1) a memory-intensive application benefits more from being placed closer to memory controller than other applications because it demands faster communication with the memory controller (See Figure 6), 2) an interference-sensitive application benefits more from being placed closer to memory controller (especially if it is memory intensive) than other applications, 3) mapping a memory-intensive application close to the memory controller reduces the network interference it causes to less network-intensive applications due to its frequent communication with the controller because its network packets travel much shorter distances. Figure 6 empirically supports the first observation by showing that applications with higher MPKI (shown in graphs to the left) gain more in terms of IPC performance when mapped closer to the memory controller than applications with lower MPKI (shown in graphs to the right).

Our intra-cluster mapping exploits these insights. It differentiates applications based on *both* their 1) *network/memory demand* (i.e. rate of injection of packets) measured as MPKI, and 2) *sensitivity* to network latency measured as *Stall Time per Miss (STPM)* at the core. It then computes a metric, stall time per thousand instructions, $STPKI = MPKI * STPM$ for each application, and sorts applications based on the value of this metric. Applications with higher $STPKI$ are assigned to cores closer to the memory controller. To achieve this, the algorithm maps applications *radially in concentric circles around the memory controller* in sorted order, starting from the application with the highest $STPKI$. Algorithm 3 shows this process. Figure 5 (c) shows an example resulting mapping within a cluster. Darker (inner and closer) tiles represent heavy and sensitive applications and lighter (outer and farther) tiles represent lower intensity

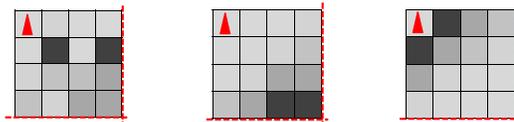


Figure 5: Intra-Cluster Application-to-Core Mapping Examples: (a) Random (b) Inverse Radial (c) Radial applications with low sensitivity. As an example of a contrasting policy which we will evaluate, Figure 5 (b) shows the resulting mapping if we perform *inverse-radial* mapping based on the sorted application order.

Algorithm 3: Intra-Cluster Mapping: Radial Mapping Algorithm

Input: $Applications = \{A_0, A_1, \dots, A_{c-1}\}$ such that
 $\forall i \text{ } MPKI(A_i) * STPM(A_i) \geq MPKI(A_{i+1}) * STPM(A_{i+1})$
 $Cores = \{N_0, N_1, \dots, N_{c-1}\}$ such that
 $\forall i \text{ } Distance(N_i, HomeMemoryController) \leq Distance(N_{i+1}, HomeMemoryController)$

- 1: **for** $i \leftarrow 0$ **to** $c - 1$ **do**
 - 2: | $ApplicationMap(N_i) \leftarrow A_i$
 - 3: **end**
-

4.4 Putting It All Together: Our Application-to-Core (A2C) Mapping Algorithm

Our final algorithm consists of three steps combining the above algorithms. First, cores are clustered into subnetworks using the cluster-CLOCK page mapping algorithm (Section 4.1). Second, the BLRI algorithm is invoked to map applications to clusters (Section 4.2.2). In other words, the algorithm attempts to allocate a separate cluster to interference-sensitive applications (Section 4.2.2), if possible, and distributes the applications to remaining clusters to balance load (Section 4.2.1). Third, after applications are assigned to clusters, the applications are mapped to cores within the clusters by invoking the intra-cluster radial mapping algorithm (Section 4.3). **We call this final algorithm, which combines clustering, BLRI and radial algorithms, as the A2C mapping.**

4.5 Enforcing System-Level Priorities in the NoC

So far, our mechanisms were targeted to improve system throughput and energy efficiency. We implicitly assumed that all applications have equal system-level priority. In a real system, the system software (the OS or virtual machine monitor) may want to assign priorities (or, weights) to applications in order to convey that some applications are more/less important than others. We seamlessly modify our mapping policies to incorporate system-assigned application weights. First the load balancing algorithm (Section 4.2.1) is invoked to distribute applications between clusters as we already described. When performing intra-cluster mapping, the system software takes into account system-level weights of applications to sort the applications instead of purely using the STPKI metric to do the sorting. This can be done by sorting the applications based on 1) purely their system-level weights or 2) their STPKI scaled with weights. The applications are mapped radially around the memory controllers within the clusters, with higher ranked ones placed closer to the memory controller.

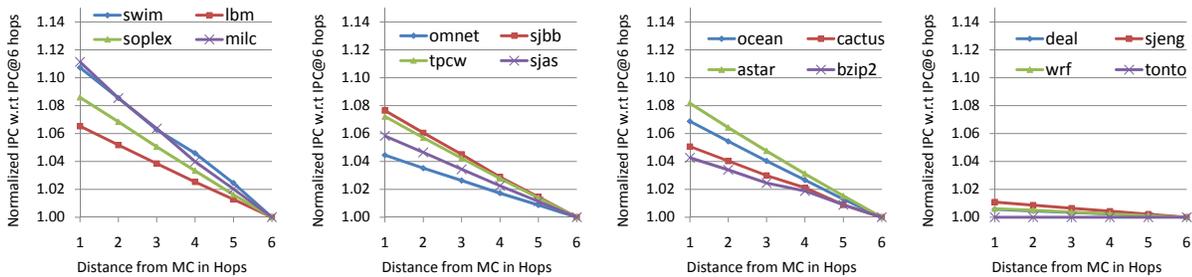


Figure 6: Application performance with distance from memory controller (MC) when running *alone*. Each figure has four benchmarks each with (a) MPKI range from 25.4 to 29.6 (b) MPKI range from 8.5 to 10.6 (c) MPKI range from 3.2 to 7.2 (d) MPKI range from 0.26 to 0.40. The figures show 1) speedups of higher-MPKI applications are higher if placed closer to MC and 2) some applications with similar or lower MPKI, show high speedups when placed closer to MC because they are more sensitive to network latency.

5 Enforcement of Application-to-Core Mapping

5.1 Profiling

The proposed mapping policies assume knowledge of two metrics: a) network demand in terms of MPKI and b) sensitivity in terms of stall cycles per miss (STPM). These metrics can be either, pre-computed for applications *a priori*, or measured online during a profiling phase. We evaluate both scenarios, where metrics are known *a priori* (static A2C) and when metrics are measured online (dynamic A2C). For dynamic A2C, we profile the workload for 10 million instructions (profiling phase) and then compute mappings that are enforced for 300 million instructions (enforcement phase). The profiling phase and enforcement phases are repeated periodically. The profiling to determine MPKI requires two hardware counters in the core: 1) instruction counter and 2) L2 miss counter. The profiling to determine STPM requires one additional hardware counter at the core which is incremented every cycle the oldest instruction cannot be retired because it is waiting for an L2 miss. Note that the A2C technique requires only a relative ordering among applications and hence quantizing applications to classes based on the above metrics is sufficient.

5.2 Operating System and Firmware Support

Our proposal requires support from the operating system. First, the operating system page allocation and replacement routine is modified to enforce clustering, as described in Section 4.1. Second, the A2C algorithm can be integrated as part of the operating system task scheduler. If this is the case, the OS scheduler allocates cores to applications based on the optimized mapping computed by the A2C algorithm. The complexity of the algorithm is relatively modest and we found its time overhead is negligible since the algorithm needs to be invoked very infrequently (e.g., every OS time quantum). Alternatively, the A2C algorithm can be implemented as part of the firmware of a multi-core processor. There is already thrust towards firmware support for multi-core processors to manage power [35, 20] and thermal issues [17, 11], support application migration, and provide fault tolerance [2]. If integrated into the firmware, our techniques can be either exposed to the operating system or completely managed by the hardware.

5.3 Adapting to Dynamic Runtime Environment

The runtime environment of a manycore processor will be dynamic with continuous flow of incoming programs (process creation), outgoing programs (process completion), and context switches. Thus, it is hard to

predict a priori which set of applications will run simultaneously as a workload on the manycore processor. Our application-to-core mapping techniques have the capability to adapt to such dynamic scenarios via two key elements. First, the ability to determine *online the application characteristics* (Section 5.1). Second, since application-to-core mapping of an application can change between different execution phases, we migrate the applications between cores to enforce new mappings. We discuss the costs of application migration in the next subsection.

5.4 Migration Costs

A new application-to-core mapping may require migration of an application from one core to another. We can split the cost associated with application migration into four parts: 1) A constant cost due to operating system bookkeeping to facilitate migration. This cost is negligible because both the cores are managed by a single unified operating system. Thus, unlike process migration in MP systems where a process is migrated between processors managed by different operating systems, in this case, minimal system state needs to be accounted for. For example, the file handling, memory management, IO, network socket state, etc are shared between the cores due to the single operating system image and need not be saved or restored; 2) A constant cost (in terms of bytes) of transferring the application’s architectural context (including registers) to the new core; 3) A variable cache warmup cost due to cache misses incurred after transferring the application to a new core. We quantify this cost and show that averaged over the entire execution phase, this cost is negligibly small across all benchmarks (see Section 7.7); and 4) A variable cost due to potential reduction in *clustering factor*⁵. This cost is incurred only when we migrate applications between clusters and, after migration, the application continues to access pages mapped to its old cluster. We quantify this cost as well in our performance evaluation for all our workloads (see Section 7.7). Our evaluations faithfully account for all of these four types of migration costs.

Processor Pipeline	2 GHz processor, 128-entry instruction window
Fetch/Exec/Commit width	2 instructions per cycle in each core; only 1 can be a memory operation
Memory Management	4KB physical and virtual pages, 512 entry TLBs, CLOCK page allocation and replacement
L1 Caches	32KB per-core (private), 4-way set associative, 64B block size, 2-cycle latency, write-back, split I/D caches, 32 MSHRs
L2 Caches	256KB per core (private), 16-way set associative, 64B block size, 6-cycle bank latency, 32 MSHRs
Main Memory	4GB DRAM, up to 16 outstanding requests per-core, 160 cycle access, 4 DDR Channels at 16GB/s 4 on-chip Memory Controllers.
Network Router	2-stage wormhole switched, virtual channel flow control, 4 VC’s per Port, 4 flit buffer depth, 4 flits per data packet, 1 flit per address packet.
Network Interface	16 FIFO buffer queues with 4 flit depth
Network Topology	8x8 mesh, 128 bit bi-directional links (32GB/s).

Table 1: Baseline Processor, Cache, Memory, and Network Configuration

6 Methodology

6.1 Experimental Setup

We evaluate our techniques using an instruction-trace-driven, cycle-level x86 CMP simulator. The functional frontend of the simulator is based on Pin dynamic binary instrumentation tool [29], which is used to collect instruction traces from applications, which are then fed to the core models that model the execution and timing of each instruction.

⁵Clustering factor is defined as the percentile of accesses that are constrained within the cluster.

Functional Simulations for Page Fault Rates: We run 500 million instructions per core (totally 30 billion instructions across 60 cores) from the simulation fast forward point obtained from [33] to evaluate cluster-CLOCK and CLOCK page replacement algorithms.
Performance Simulations for Static A2C mapping: To have tractable simulation time, we choose a smaller representative window of instructions, obtained by profiling each benchmark, from the representative execution phase obtained from [33]. All our experiments study multi-programmed workloads, where each core runs a separate benchmark. We simulate 10 million instructions per core, corresponding to at least 600 million instructions across 60 cores.
Performance Simulations for Dynamic A2C mapping: To evaluate the dynamic application-to-core mapping faithfully, we need longer simulations that model at least one dynamic profiling phase and one enforcement phase (as explained in Section 5.1). We simulate an entire profiling+enforcement phase (300 million instructions) per benchmark per core, corresponding to at least 18.6 billion instructions across 60 cores.

Table 2: Simulation Methodology

Table 1 provides the configuration of our baseline, which consists of 60 cores and 4 memory controllers connected by a 2D, 8x8 Mesh NoC. Each core is modeled as an out-of-order execution core with a limited instruction window and limited buffers. The memory hierarchy uses a two-level directory-based MESI cache coherence protocol. Each core has a private write-back L1 cache and private L2 cache. The network connects the core tiles and memory controller tiles. The system we model is self-throttling as real systems are: if the miss buffers are full the core cannot inject more packets into the network. Each router uses a state-of-the-art two-stage microarchitecture. We use the deterministic X-Y routing algorithm, finite input buffering, wormhole switching, and virtual-channel flow control. We use the Orion power model for estimating the router power [36].

We also implemented a detailed functional model for virtual memory management to study page access and page fault behavior of our workloads. The baseline page allocation and replacement policy is CLOCK [23]. The modified page replacement and allocation policy, cluster-CLOCK, looks ahead 512 pages beyond the first replacement candidate to potentially find a replacement page belonging to home memory controller.

The parameters used for our A2C algorithm are: $Thresh_{MPILow} = 5$ MPKI, and $Thresh_{MPIHigh} = 25$ MPKI, $Thresh_{SensitivityRatio} = 5$ and $Thresh_{MPKI-RI} = 250$ MPKI. These parameters are empirically determined but not tuned. The constant cost for OS book keeping while migrating applications is assumed to be 50K cycles. The migrating applications write and read 128 bytes to/from the memory to save and restore their register contexts.

6.2 Evaluation Metrics

Our evaluation uses several metrics. We measure **system performance** in terms of average **weighted speedup** [14], a commonly used multi-program performance metric, which is the average of the sum of slowdowns of each application compared to when it is run alone on the same system. We measure **system fairness** in terms of the maximum slowdown any application experiences in the system.

$$(Average) \text{ Weighted Speedup} = \frac{1}{NumThreads} \times \sum_{i=1}^{NumThreads} \frac{IPC_i^{shared}}{IPC_i^{alone}} \quad (1)$$

$$UnfairnessIndex = \max_i \frac{IPC_i^{shared}}{IPC_i^{alone}} \quad (2)$$

IPC_{alone} is the IPC of the application when run alone on our baseline system. We also report **IPC throughput**.

$$IPC \text{ Throughput} = \sum_{i=1}^{NumThreads} IPC_i \quad (3)$$

6.3 Workloads and Simulation Methodology

Table 2 describes the three different types of simulations we run to evaluate our mechanisms. Due to simulation time limitations, we evaluate the execution time effects of our static and dynamic mechanisms without taking into account the effect on page faults. We evaluate the effect of our proposal on page fault rate separately using functional simulation, showing that our proposal reduces page fault rate (Sec. 7.2).

We use a diverse set of multiprogrammed application workloads comprising scientific, commercial, and desktop applications. In total, we study 35 applications, including SPEC CPU2006 benchmarks, applications from SPLASH-2 and SpecOMP benchmark suites, and four commercial workloads traces (*sap*, *tpcw*, *sjobb*, *sjas*). We choose representative execution phases using PinPoints [33] for all our workloads except commercial traces, which were collected over Intel servers. Figure 12 (b) lists each application and includes results showing the MPKI of each application on our baseline system.

Multiprogrammed Workloads and Categories: All our results are across 128 randomly generated workloads. Each workload consists of 10 copies each of 6 applications randomly picked from our suite of 35 applications. The 128 workloads are divided into four subcategories of 32 workloads each: 1) MPKI500: relatively less network-intensive workloads with aggregate MPKI less than 500, 2) MPKI1000: aggregate MPKI is between 500-1000, 3) MPKI1500: aggregate MPKI is between 1000-1500, 4) MPKI2000: relatively more network-intensive workloads with aggregate MPKI between 1500-2000.

7 Performance Evaluation

7.1 Overall Results for A2C Algorithm

We first show the overall results of our final Application-to-Core Mapping algorithm (A2C). We evaluate three systems: 1) the baseline system with random mapping of applications to cores (BASE), which is representative of existing systems, 2) our enhanced system which uses our clustering and modified CLOCK algorithm (described in Section 4.1) but still uses random mapping of applications to cores (CLUSTER+RND), 3) our final system which uses our combined A2C algorithm (summarized in Section 4.4), which consists of clustering, inter-cluster BLRI mapping, and intra-cluster radial mapping algorithms (A2C).

Figure 7 (a) and (b) respectively show system performance (higher is better) and system unfairness (lower is better) of the three systems. Solely using clustering (CLUSTER+RND) improves weighted speedup by 9.3% over the baseline (BASE). A2C improves weighted speedup by 16.7% over the baseline, while also reducing unfairness by 22%.⁶

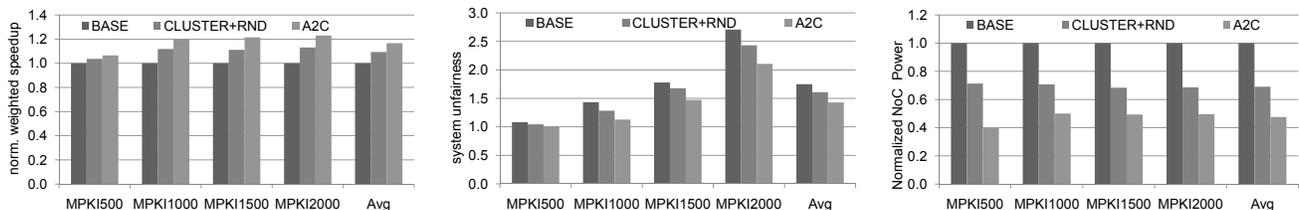


Figure 7: (a) System performance (b) system unfairness and (c) interconnect power of the A2C algorithm for 128 workloads

⁶We do not show graphs for instruction throughput due to space limitations. Clustering alone (CLUSTER+RND) improves IPC throughput by 7.0% over the baseline, while A2C improves IPC by 14.0%.

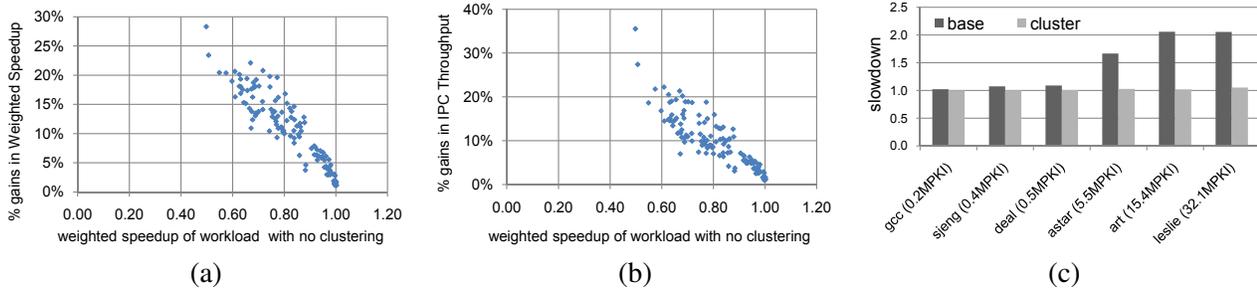


Figure 8: Performance gain of clustering across 128 workloads (a) Weighted speedup (b) IPC throughput (c) A Case Study

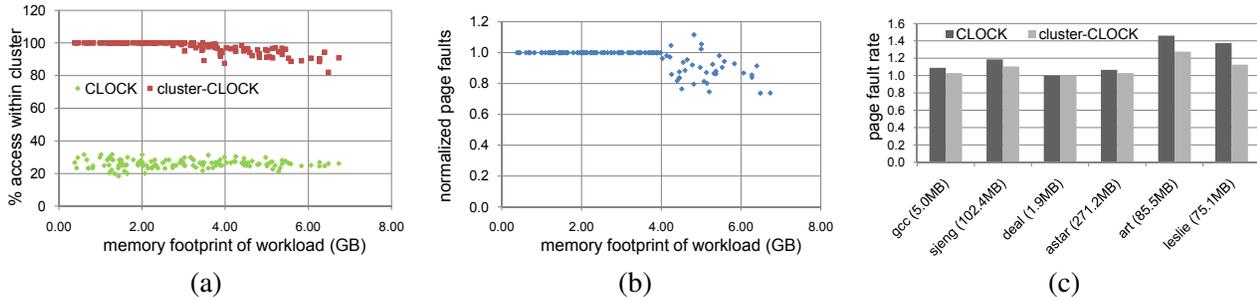


Figure 9: Page access behavior across 128 multiprogrammed workloads (a) Clustering factor (b) Page Faults (c) A Case Study

Interconnect Power Figure 7 (c) shows the normalized interconnect average power consumption (lower is better). Clustering only reduces power consumption by 31.2%; A2C mapping reduces power consumption by 52.3% over baseline (BASE). The clustering of applications to memory controllers reduces the average hop count significantly, reducing the energy spent in moving data over the interconnect. Using inter-cluster and intra-cluster mapping further reduces hop count and power consumption by ensuring that network-intensive applications get mapped close to the memory controllers and network load is balanced across controllers after clustering.

In the next three sections, we analyze the benefits and tradeoffs of each component of A2C.

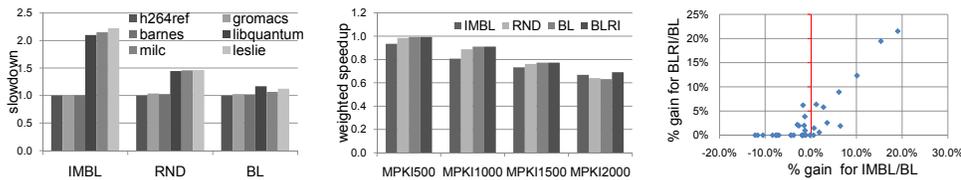


Figure 10: (a) Case study analyzing BL mapping (b) Average results for different inter-cluster mappings across 128 workloads (c) Average results for BLRI mapping for the MPKI2000 workload category

7.2 Analysis of Clustering and Cluster-CLOCK

The goal of clustering is to reduce interference between applications mapped to different clusters. Averaged across 128 workloads, clustering improves system throughput by 9.3% in terms of weighted speedup and 8.0% in terms of IPC throughput.

Figure 8 (a) and (b) plot the gains in weighted speedup and IPC throughput due to clustering for each workload against the baseline weighted speedup of the workload. A lower baseline weighted speedup indi-

cates that average slowdown of applications are higher and hence contention/interference is high between applications in the baseline. The figures show that performance gains due to clustering are significantly higher for workloads with lower weighted speedup (i.e., higher slowdowns due to interference). This is intuitive because the main goal of clustering is to reduce interference between applications. We conclude that the benefits of clustering are higher when interference is higher in the network.

To provide detailed insight, we zoom in on one workload as a case study in Figure 8 (c). This workload consists of 10 copies each of applications `gcc`, `sjeng`, `deal`, `astar`, `art`, and `leslie` running together on the many-core processor. The Y-axis measures the slowdowns (lower is better) of individual applications compared to when run alone. We observe that clustering 1) reduces slowdown of all applications because it reduces interference and reduces their latency of access to the memory controller, 2) provides larger benefits to more network-intensive (higher MPKI) applications because they have significantly higher slowdown in the baseline system and separating these applications' accesses via clustering reduces this large slowdown they cause to each other.

Analysis of the cluster-CLOCK page replacement algorithm: To enforce clustering, we have developed the cluster-CLOCK algorithm (Sec. 4.1) which modifies the default page allocation and page replacement policies. The results in Figure 9 quantify the effect of cluster-CLOCK across 128 workloads. Figure 9 (a) plots the clustering factor with the baseline policy (CLOCK) and our new policy (cluster-CLOCK). Recall that the clustering factor is the percentage of all accesses that are serviced by the home memory controller. On average, cluster-CLOCK improves the clustering factor from 26.0% to 97.4%, thereby reducing interference among applications.

Figure 9 (b) shows the normalized page fault rate of cluster-CLOCK for each workload (Y axis) versus the memory footprint of the workload (X axis). A lower relative page fault rate indicates that cluster-CLOCK reduces the page fault rate compared to the baseline. We observe that cluster-CLOCK 1) does not affect the page fault rate for workloads with small memory footprint, 2) in general reduces the page fault rate for workloads with large memory footprint. On average, cluster-CLOCK reduces the page fault rate by 4.1% over 128 workloads. This is a side effect of cluster-CLOCK since the algorithm is not designed to reduce page faults. Yet, it reduces page faults because it happens to make better page replacement decisions than CLOCK (i.e., replace pages that are less likely to be reused soon) by reducing the interference between applications in physical memory space: by biasing replacement decisions to be made within each memory controller as much as possible, applications mapped to different controllers interfere less with each other in the physical memory space. As a result, applications with lower page locality disturb applications with higher page locality less, improving page fault rate. Note that our execution time results do not include this effect of reduced page faults due to simulation speed limitations.

To illustrate this behavior, we focus on one workload as a case study in Figure 9 (c), which depicts the page fault rate in terms of page faults incurred per unique page accessed by each application with CLOCK and cluster-CLOCK. Applications `art` and `leslie` have higher page fault rate but we found that they have good locality in page access. On the other hand, `astar` also has high footprint but low locality in page access. When these applications run together using the CLOCK algorithm, `astar`'s pages contend with `art` and `leslie`'s pages in the entire physical memory space, causing those pages to be evicted from physical memory. On the other hand, if cluster-CLOCK is used, and `astar` is mapped to a different cluster from `art` and `leslie`, the likelihood that `astar`'s pages replace `art` and `leslie`'s pages reduces significantly because cluster-CLOCK attempts to replace a page from the home memory controller `astar` is assigned to instead of any page in the physical memory space. Hence, cluster-CLOCK can reduce page fault rates by likely localizing page replacement and thereby limiting as much as possible page-level interference among applications to pages assigned to a single memory controller.

7.3 Analysis of Inter-Cluster Mapping

We study the effect of inter-cluster load balancing algorithms described in Section 4.2. For these evaluations, all other parts of the A2C algorithm is kept the same: we employ clustering and radial intra-cluster mapping. We first show the benefits of balanced mapping (BL), then show when and why imbalanced mapping (IMBL) can be beneficial, and later evaluate our BLRI algorithm, which aims to achieve the benefits of both balance and imbalance.

BL Mapping Figure 10 (a) shows application slowdowns in an example workload consisting of 10 copies each of `h264ref`, `gromacs`, `barnes`, `libquantum`, `milc` and `leslie` applications running together. The former three are very network-insensitive and exert very low load and the latter three are both network-intensive and network-sensitive. A completely imbalanced (IMBL) inter-cluster mapping (described in Section 4.2) severely slows down the latter three network-intensive and network-sensitive applications because they get mapped to the same clusters, causing significant interference to each other, whereas the former three applications do not utilize the bandwidth available in their clusters. A random (RND) mapping (which is our baseline) still slows down the same applications, albeit less, by providing better balance in interference. Our balanced (BL) mapping algorithm, which distributes load in a balanced way among all clusters provides the best speedup (19.7% over IMBL and 5.9% over RND) by reducing the slowdown of all applications because it does not overload any single cluster.

IMBL Mapping Figure 10 (b) shows average weighted speedups across 128 workloads, categorized by the overall intensity of the workloads. When the overall workload intensity is not too high (i.e., less than 1500 MPKI), balanced mapping (BL) provides significantly higher performance than IMBL and RND by reducing and balancing interference. However, when the workload intensity is very high (i.e., greater than 1500 MPKI), BL performs worse than IMBL mapping. The reason is that IMBL mapping isolates network-non-intensive (low MPKI) applications from network-intensive (high MPKI) ones by placing them into separate clusters. When the network load is very high, such isolation significantly improves the performance of network-non-intensive applications without significantly degrading the performance of intensive ones by reducing interference between the two types of applications.

The main takeaway is that when network load is very high, relatively less intensive applications' progress gets slowed down *too* significantly because other applications keep injecting interfering requests. On the other hand, when the network load is not as high, less-intensive applications can still make fast progress because interference in the network is less. As a result, if network load is very high, it is more beneficial to separate the accesses of non-intensive applications from others by placing them into separate clusters, thereby allowing their fast progress. However, such separation is not beneficial for non-intensive applications and harmful for performance if the network load is not high: it causes wasted bandwidth in some clusters and too much interference in others. This observation motivates our BLRI algorithm (which creates a separate cluster for non-intensive applications only when the network load is high), which we analyze next.

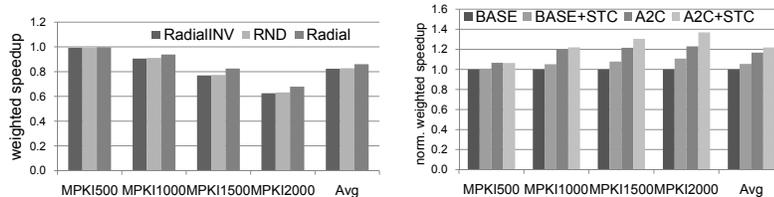


Figure 11: (a) Average results for radial mapping across 128 workloads (b) Performance comparison of A2C mapping and application-aware prioritization (STC) for 128 workloads

Cache Size	256KB	512KB	1MB	Number of MCs	4 MC	8 MC
Performance Gain	16.7%	13.6%	12.1%	Performance Gain	16.7%	17.9%

Table 3: Sensitivity to last level per-core cache size and number of memory controllers

BLRI Mapping Figure 10 (c) shows the speedup achieved by BLRI mapping over BL for all 32 workloads in MPKI2000 category (recall that BLRI is not invoked for workloads with aggregate MPKI of less than 1500) on the Y axis against the speedup achieved for the same workload by IMBL over BL. We make two observations. First, for workloads where imbalanced mapping (IMBL) improves performance over balanced mapping (BL), shown in the right half of the plot, BLRI also significantly improves performance over BL. Second, for workloads where imbalance reduces performance (left half of the plot), BLRI either improves performance or does not affect performance. *We conclude that BLRI achieves the best of both worlds (load balance and imbalance) by isolating those applications that would most benefit from imbalance and performing load balancing for the remaining ones.*

7.4 Effect of Intra-Cluster Mapping

We analyze the effect of intra-cluster mapping policy, after applications are assigned to clusters using the BLRI inter-cluster policy. We examine three different intra-cluster mapping algorithms: 1) Radial: our proposed radial mapping described in Section 4.2, 2) RND: Cores in a cluster are assigned randomly to applications, 3) RadialINV: This is the inverse of our radial algorithm; those applications that would benefit least from being close to the memory controller (i.e., those with low STPKI) are mapped closest to the memory controller. Figure 11 (a) shows the average weighted speedup of 128 workloads with BLRI inter-cluster mapping and different intra-cluster mappings. The radial intra-cluster mapping provides 0.4%, 3.0%, 6.8%, 7.3% for MPKI500, MPKI1000, MPKI1500, MPKI2000 category workloads over RND intra-cluster mapping. Radial mapping is the best mapping for all workloads; RadialINV is the worst. We conclude that our metric and algorithm for identifying and deciding which workloads to map close to the memory controller is effective.

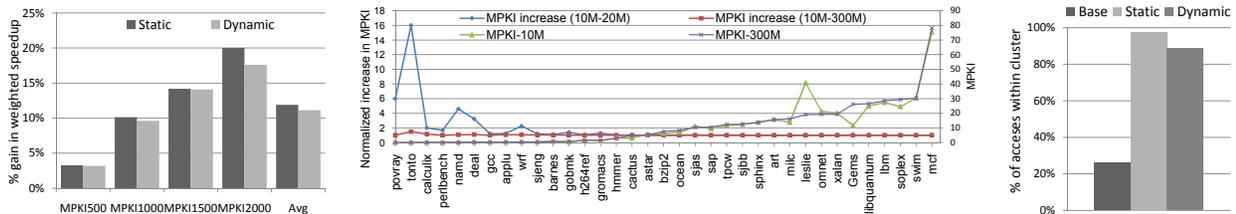


Figure 12: (a) Performance of static and dynamic A2C for 8 workloads (b) Increase in MPKI due to migration for individual applications (c) Clustering factor for 128 workloads

7.5 A2C vs Application-Aware Prioritization

We compare the benefits of A2C mapping to application-aware prioritization in the network to show that our interference-aware mapping mechanisms are orthogonal to interference-aware packet scheduling in the NoC.

Das et al. [13] proposed an application-aware arbitration policy (called STC) to accelerate network-sensitive applications. The key idea is to rank applications at regular intervals based on their network intensity (outermost private cache MPKI), and prioritize packets of non-intensive applications over packets of intensive applications. We compare Application-to-Core mapping policies to application-aware prioritiza-

tion policies because both techniques have similar goals: to effectively handle inter-application interference in NoC. *However, both techniques take different approaches towards the same goal.* STC tries to make the best decision when interference happens in the NoC by efficient packet scheduling in routers while A2C tries to reduce interference by mapping applications to separate clusters and controllers.

Figure 11 (b) shows that STC⁷ is orthogonal to our proposed A2C technique and its benefits are additive to A2C. STC prioritization improves performance by 5.5% over the baseline whereas A2C mapping improves performance by 16.7% over the baseline. When STC and A2C are combined together, overall performance improvement is 21.9% over the baseline, greater than the improvement provided by either alone. STC improves performance when used with A2C mapping because it prioritizes non-intensive applications (shorter jobs) within a cluster in a coordinated manner, ensuring all the routers act in tandem. In contrast, the baseline round-robin arbitration policy is uncoordinated: it causes one application to be prioritized in one router only to be delayed in the next, resulting in a slowdown for all applications. This coordinated packet scheduling effect of STC is orthogonal to the benefits of A2C. Hence, we conclude that our mapping mechanisms interact synergistically with application-aware packet scheduling.

7.6 Sensitivity Studies

We evaluated the sensitivity of A2C to last-level cache size and number of memory controllers. Table 3 shows A2C's performance improvement over baseline on 128 multiprogrammed workloads. As cache size increases, performance benefit of A2C decreases because traffic and hence interference in the NoC reduces. However, A2C still provides 12.1% performance improvement even when total on-chip cache size is 60MB (i.e., 1MB/core). As we vary the number of memory controllers from 4-8, A2C's performance benefit stays similar because more controllers 1) on the one hand enable more fine-grained clustering, improving communication locality and reducing interference, 2) on the other hand can reduce clustering factor when application behavior changes. We conclude that our proposed techniques provide consistent performance improvements when we vary major relevant system parameters.

7.7 Effect of Dynamic A2C Mapping

Our evaluation assumed so far static mapping is formed with pre-runtime knowledge of application characteristics. We now evaluate the dynamic A2C scheme, described in Section 5.1. We use a profiling phase of 10 million instructions, after which the operating system forms a new application-to-core mapping and enforces it for the whole phase (300 million instructions). An application can migrate at the end of the dynamic profiling interval after each phase.

Figure 12 (c) compares the performance improvement achieved by static and dynamic A2C schemes for eight workloads, which consist of two workloads from each MPKI category, over a baseline that employs clustering but uses random application-to-core mapping.⁸ The performance of dynamic A2C is close to that of static A2C (within 1% on average) for these eight workloads. However, static A2C performs better for the two MPKI2000 workloads. We found this is because the BLRI scheme (which determines sensitive applications online and forms a separate cluster for them) requires re-mapping at more fine-grained execution

⁷The default parameters used for STC [13] are: 1) ranking levels $R = 8$, 2) batching levels $B = 8$, 3) ranking interval = 350,000 cycles, 4) batching interval = 16,000 cycles, 5) BCIP packet sent every $U = 4000$ cycles.

⁸We show only eight workloads because we were severely limited by simulation time for evaluation and analysis of dynamic scenarios. The performance simulations for dynamic mapping took us more than a month for different workloads because the minimum length of cycle-level simulation had to be a whole phase of 300 million instructions per application (totally 18 billion instructions).

phases. Unfortunately, given the simulation time constraints we could not fully explore the best thresholds for the dynamic scheme. We conclude that the dynamic A2C scheme provides significant performance improvements, even with untuned parameters.

Analysis In Section 5, we qualitatively discussed the overheads of application migration. In this section, we first quantify the increases in cache misses when an application migrates from one core to another. We then quantify the reduction in clustering factor due to migrations. Overall, these results provide quantitative insight into why the dynamic A2C algorithm works.

Figure 12 (b) analyzes the MPKI of the 35 applications during the profiling phase (MPKI-10M) and the enforcement phase (MPKI-300M). It also analyzes the increase in the MPKI due to migration to another core right during the 10M instruction interval right after the migration happens (MPKI increase (10-20M)) and during the entire enforcement phase (MPKI increase (300M)). The left Y-axis is the normalized MPKI of the application when it is migrated to another core compared to the MPKI when it is running alone (a value of 1 means the MPKI of the application does not change after migration). Benchmarks on the X axis are sorted from lightest (lowest baseline MPKI) to heaviest (highest baseline MPKI) from left to right. We make several key observations:

- The MPKI in the profiling phase (MPKI-10M) correlates well with the MPKI in the enforcement phase (MPKI-300M), indicating why dynamic profiling can be effective.
- MPKI increase within 10M instructions after migration is negligible for high-intensity workloads, but significant for low-intensity workloads. However, since these benchmarks have very low MPKI to begin with, their execution time is not significantly affected.
- MPKI increase during the entire phase after migration is negligible for almost all workloads. This increase is 3% on average and again observed mainly in applications with low intensity. These results show that cache migration cost of migrating an application to another core is minimal over the enforcement phase of the new mapping.

The clustering factor (i.e., the ratio of memory accesses that are serviced by the home memory controller) is also affected by application migration. The clustering factor may potentially decrease, if, after migration, an application continues to access the pages mapped to its old cluster. Our dynamic algorithm minimizes the number of inter-cluster application migrations by placing applications with similar MPKI in equivalence classes for A2C mapping computation (we omit the detailed explanation of this optimization due to space limitations).

Figure 12 (c) shows the clustering factor for our 128 workloads with 1) baseline RND mapping, 2) static A2C mapping, and 3) dynamic A2C mapping. The clustering factor reduces from 97.4% with static A2C to 89.0% with dynamic A2C due to inter-cluster application migrations. However, dynamic A2C mapping still provides a very large improvement in clustering factor compared to the baseline mapping. We conclude that both dynamic and static versions of A2C are effective: static A2C is desirable when application information is profiled before runtime, but dynamic A2C does not cause significant overhead and achieves similar performance.

8 Related Work

To our knowledge, this is the first work that tackles the problem of how to map applications to cores in a network-on-chip to minimize inter-application interference. We briefly describe the most closely related previous work in this section.

Thread-to-Core Mapping: Prior works have proposed thread-to-core mapping to improve locality of communication by placing frequently communicating threads/tasks closer to each other [27, 30] in parallel ap-

plications and for shared cache management [8]. As such, their goal was to reduce inter-thread or inter-task communication. Our techniques solve a different problem: inter-application interference. As such, our goals are to 1) reduce interference in the network between *different* independent applications and 2) reduce contention at the memory controllers.

Reducing Interference in NoCs: Recent works [26, 19, 13] propose prioritization and packet scheduling policies to provide quality of service or improve application-level throughput in the NoC. These works are orthogonal to our mechanism since they perform packet scheduling. We have already shown that our proposal works synergistically with application-aware prioritization [13].

Memory Controllers and Page Allocation: Awasthi et al. [4] have explored page allocation and page migration in the context of multiple memory controllers in a multi-core processor with the goal of balancing load between memory controllers and improving DRAM row buffer locality. The problem we solve is different and our techniques are complementary to theirs. First, our goal is to reduce interference between applications in the interconnect and the memory controllers. Second, our mechanism solves the problem of mapping applications to the cores; [4] does not solve this problem. We believe the page migration techniques proposed in [4] can be employed to reduce the costs of migration in our dynamic application-to-core mapping policies.

Page allocation and migration has been explored extensively to improve locality of access within a single application executing on a NUMA multiprocessor (e.g., [6, 15]). These works do not aim to reduce interference between multiple applications sharing the memory system, and hence do not solve the problem we are aiming to solve. In addition, these works do not consider how to map applications to cores in a network-on-chip based system. Page coloring techniques have been employed in caches to reduce cache conflict misses [24, 5] and to improve cache locality in shared caches (e.g., [9]). These works solve an orthogonal problem and can be combined with our techniques.

Abts et al. [3] recently explore placement of memory controllers in a multi-core processor to minimize channel load. Memory controller placement is orthogonal to our mechanisms.

Thread Migration: Thread migration has been explored to manage power [35, 20], thermal hotspots [17, 11] or to exploit heterogeneity in multi-core processors [28]. We use thread migration to enable dynamic application-to-core mappings to reduce interference in NoCs.

9 Conclusion

We presented application-to-core mapping policies that largely reduce inter-application interference in network-on-chip based multi-core systems. We have shown that by intelligently mapping applications to cores in a manner that is aware of applications' characteristics, significant increases in system performance, fairness, and energy-efficiency can be achieved at the same time. Our proposed algorithm, A2C, achieves this by using two key principles: 1) mapping network-latency-sensitive applications to separate node clusters in the network from network-bandwidth-intensive applications such that the former makes fast progress without heavy interference from the latter, 2) mapping those applications that benefit more from being closer to the memory controllers close to these resources.

We have extensively evaluated our mechanism and compared it both qualitatively and quantitatively to different application mapping and packet prioritization policies. Our main results show that: 1) averaged over 128 multiprogrammed workload mixes on a 60-core 8x8-mesh system, our proposed A2C improves system throughput by 16.7%, while also reducing system unfairness by 22.4% and interconnect power consumption by 52.3%, 2) A2C is orthogonal to application-aware packet prioritization techniques. We conclude that the proposed approach can be an effective way of improving overall system throughput, fairness,

and power-efficiency and therefore can be a promising way to exploit the non-uniform structure of network-on-chip-based multi-core systems.

References

- [1] “Linux source code (version 2.6.39),” <http://www.kernel.org>.
- [2] H. Abdel-Shafi, E. Speight, and J. K. Bennett, “Efficient user-level thread migration and checkpointing on windows nt clusters,” in *Usenix (3rd Windows NT Symposium)*, 1999.
- [3] D. Abts, N. D. E. Jerger, J. Kim, D. Gibson, and M. H. Lipasti, “Achieving predictable performance through better memory controller placement in many-core cmps.” in *ISCA-36*, 2009.
- [4] M. Awasthi, D. Nellans, K. Sudan, R. Balasubramonian, and A. Davis, “Handling the problems and opportunities posed by multiple on-chip memory controllers,” in *PACT-19*, 2010.
- [5] B. N. Bershad, D. Lee, T. H. Romer, and J. B. Chen, “Avoiding conflict misses dynamically in large direct-mapped caches,” in *ASPLOS-VI*, 1994.
- [6] R. Chandra, S. Devine, B. Verghese, A. Gupta, and M. Rosenblum, “Scheduling and page migration for multiprocessor compute servers,” in *ASPLOS-VI*, 1994.
- [7] J. Chang and G. S. Sohi, “Cooperative caching for chip multiprocessors,” in *ISCA-33*, 2006.
- [8] S. Chen, P. B. Gibbons, M. Kozuch, V. Liaskovitis, A. Ailamaki, G. E. Blelloch, B. Falsafi, L. Fix, N. Hardavellas, T. C. Mowry, and C. Wilkerson, “Scheduling threads for constructive cache sharing on cmps,” in *SPAA-19*, 2007.
- [9] S. Cho and L. Jin, “Managing distributed, shared L2 caches through OS-level page allocation,” in *MICRO-39*, 2006.
- [10] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. MIT Press, 2009.
- [11] A. K. Coskun, R. Strong, D. M. Tullsen, and T. Simunic Rosing, “Evaluating the impact of job scheduling and power management on processor lifetime for chip multiprocessors,” in *SIGMETRICS-11*, 2009.
- [12] W. J. Dally and B. Towles, *Principles and Practices of Interconnection Networks*. Morgan Kaufmann, 2003.
- [13] R. Das, O. Mutlu, T. Moscibroda, and C. Das, “Application-Aware Prioritization Mechanisms for On-Chip Networks,” in *MICRO-42*, 2009.
- [14] S. Eyerhan and L. Eeckhout, “System-level performance metrics for multiprogram workloads,” *IEEE Micro*, May-June 2008.
- [15] B. Falsafi and D. A. Wood, “Reactive numa: a design for unifying s-coma and cc-numa,” in *ISCA-24*, 1997.

- [16] A. Glew, “MLP Yes! ILP No! Memory Level Parallelism, or, Why I No Longer Worry About IPC,” in *ASPLOS Wild and Crazy Ideas Session*, 1998.
- [17] M. Goma, M. D. Powell, and T. N. Vijaykumar, “Heat-and-run: leveraging SMT and CMP to manage power density through the operating system,” in *ASPLOS-XI*, 2004.
- [18] P. Gratz, B. Grot, and S. W. Keckler, “Regional congestion awareness for load balance in networks-on-chip,” in *HPCA-14*, 2008.
- [19] B. Grot, S. W. Keckler, and O. Mutlu, “Preemptive Virtual Clock: A Flexible, Efficient, and Cost-effective QOS Scheme for Networks-on-a-Chip,” in *MICRO-42*, 2009.
- [20] S. Heo, K. Barr, and K. Asanović, “Reducing power density through activity migration,” in *ISLPED*, 2003.
- [21] J. Hu and R. Marculescu, “Energy-aware mapping for tile-based NoC architectures under performance constraints,” in *ASPDAC*, 2003.
- [22] —, “Dyad: smart routing for networks-on-chip,” in *DAC-41*, 2004.
- [23] F. C. Jiang and X. Zhang, “CLOCK-Pro: an effective improvement of the CLOCK replacement,” in *USENIX*, 2005.
- [24] R. E. Kessler and M. D. Hill, “Page placement algorithms for large real-indexed caches,” *ACM Transactions on Computer Systems*, 1992.
- [25] Y. Kim, D. Han, O. Mutlu, and M. Harchol-Balter, “ATLAS: A scalable and high-performance scheduling algorithm for multiple memory controllers,” in *HPCA-16*, 2010.
- [26] J. W. Lee, M. C. Ng, and K. Asanovic, “Globally-Synchronized Frames for Guaranteed Quality-of-Service in On-Chip Networks,” in *ISCA-35*, 2008.
- [27] T. Lei and S. Kumar, “A two-step genetic algorithm for mapping task graphs to a network on chip architecture,” in *Euromicro Symposium on Digital Systems Design*, 2003.
- [28] T. Li, P. Brett, R. Knauerhase, D. Koufaty, D. Reddy, and S. Hahn, “Operating system support for overlapping-isa heterogeneous multi-core architectures,” in *HPCA-16*, 2010.
- [29] C. K. Luk *et al.*, “Pin: building customized program analysis tools with dynamic instrumentation,” in *PLDI*, 2005.
- [30] S. Murali and G. D. Micheli, “Bandwidth-constrained mapping of cores onto NoC architectures,” in *DATE*, 2004.
- [31] O. Mutlu, H. Kim, and Y. N. Patt, “Efficient runahead execution: Power-efficient memory latency tolerance,” *IEEE Micro*, 2006.
- [32] C. Natarajan, B. Christenson, and F. Briggs, “A study of performance impact of memory controller features in multi-processor server environment,” in *WMPI-3*, 2004.
- [33] H. Patil, R. Cohn, M. Charney, R. Kapoor, A. Sun, and A. Karunanidhi, “Pinpointing representative portions of large intel itanium programs with dynamic instrumentation,” in *MICRO-37*, 2004.

- [34] M. K. Qureshi, “Adaptive spill-recvive for robust high-performance caching in cmps,” in *HPCA-15*, 2009.
- [35] K. K. Rangan, G.-Y. Wei, and D. Brooks, “Thread motion: fine-grained power management for multi-core systems,” in *ISCA-36*, 2009.
- [36] H.-S. Wang, X. Zhu, L.-S. Peh, and S. Malik, “Orion: A power-performance simulator for interconnection networks,” in *MICRO-35*, 2002.
- [37] Z. Zhang, Z. Zhu, and X. Zhang, “A permutation-based page interleaving scheme to reduce row-buffer conflicts and exploit data locality,” in *MICRO-33*, 2000.