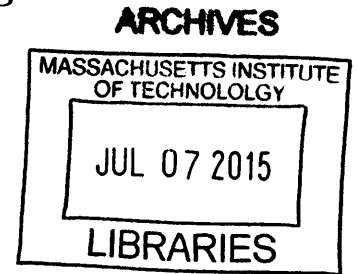# Reducing Data Movement in Multicore Chips with Computation and Data Co-Scheduling

by

Po-An Tsai

B.S. in Electrical Engineering
National Taiwan University, 2012

Submitted to the Department of Electrical Engineering and Computer
Science
in partial fulfillment of the requirements for the degree of

Master of Science in Electrical Engineering and Computer Science

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2015

**Signature redacted**

Author . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Department of Electrical Engineering and Computer Science
May 19, 2015

**Signature redacted**

Certified by . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Daniel Sanchez
Assistant Professor of Electrical Engineering and Computer Science
Thesis Supervisor

**Signature redacted**

Accepted by . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Leslie A. Kolodziejski
Chairman, Department Committee on Graduate Students

# Reducing Data Movement in Multicore Chips with Computation and Data Co-Scheduling

by

Po-An Tsai

## Abstract

Energy efficiency is the main limitation to the performance of parallel systems. Current architectures often focus on making cores more efficient. However, data movement is much more costly than basic compute operations. For example, at 28 nm, a main memory access is 100× slower and consumes 1000× the energy of a floating-point operation, and moving 64 bytes across a 16-core processor is 50× slower and consumes 20× the energy. Without a drastic reduction in data movement, memory accesses and communication costs will limit the scalability of future computing systems. Conventional hardware-only and software-only techniques miss many opportunities to reduce data movement.

This thesis presents computation and data co-scheduling (CDCS), a technique that jointly performs computation and data placement to reduce both on-chip and off-chip data movement. CDCS integrates hardware and software techniques: Hardware lets software control data mapping to physically distributed caches, and software uses this support to periodically reconfigure the chip, minimizing data movement. On a simulated 64-core system, CDCS outperforms a standard last-level cache by 46% on average (up to 76%) in weighted speedup, reduces both on-chip network traffic (by 11×) and off-chip traffic (by 23%), and saves 36% of system energy.

# Acknowledgments

I would like to thank my advisor Professor Daniel Sanchez for being such an excellent and amazing advisor. Daniel is a knowledgeable and approachable advisor. I enjoyed discussing ideas and early results of this project with him over the past two years. Daniel is also encouraging and provides me insightful advice even outside research. I am grateful that I can have Daniel as my advisor; without his guidance and consistent help, this thesis would not have been possible.

I would also like to thank my group mate, Nathan Beckmann, who was instrumental in this project. The whole project is based on his earlier research, and he contributed to the implementation and algorithms in this project. As a senior student, Nathan mentors me in many ways, from coding to work-life balance, and also divides the blame due to my own mistakes with me. I am lucky to have a chance to work with him.

Our group also supports me and gives useful opinions to this project. In particular, I would like to thank Professor Joel Emer, Mark, Harshad, and Suvinay for their helpful feedback on this manuscript, as well as my presentations. I would also like to thank Alin, Albert, Ilia, Sarah, and the CSAIL community. Hanging out with them is always lots of fun. Moreover I would like to thank Yi-Min, my long-time friend and roommate, as well as the ROCSA community for relieving my homesickness from time to time.

Last but not the least, I also want to thank my parents, brother, and sister for providing support and opportunity for me to pursue graduate studies overseas. I am fortunate to have their support to focus on my research freely. Finally, I want to thank my girlfriend, Man-Chi, for always being nice, patient, and supportive.

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1 Motivation

Energy efficiency is the main limitation to performance in future computing systems. Traditional architectures have focused on optimizing computation, neglecting data movement. However, without significant improvements, data movement will become a major limitation: A main memory access is 100× slower and consumes 1000× more energy than a floating point instruction; moving 64 bytes across a 16-core multicore is 50× slower and consumes 20× more energy [14]. Thus, this work focuses techniques to minimize data movement. To systematically reduce data movement, future systems require changes in both hardware and software. This thesis thus proposes a *data-centric* view of computation that minimizes data movement instead of maximizing core utilization.

The cache hierarchy is one of the common methods to reduce data movement in current CMPs [14, 21], but the trend towards many simpler and specialized cores further constrains the energy and latency of cache accesses [14]. To address this problem, cache architectures are becoming increasingly non-uniform (NUCA [33]), providing fast access to physically close banks, and slower access to far-away banks. For systems to scale efficiently, data must be close to the computation that uses it to minimize on-chip traffic. This requires keeping cached data in banks close to threads (to minimize on-chip data movement), while judiciously allocating cache capacity among threads (to minimize off-chip data movement).

Ideally, threads and data should be allocated and placed *jointly* across the chip to account for competing tradeoffs and avoid hotspots. Nevertheless, prior work has not addressed this problem in a unified way. On the one hand, dynamic and partitioned NUCA techniques [2,3,4,8,10,11,20,27,41,50,62] allocate cache space among threads, and then place data close to the threads that use it. However, these techniques ignore thread placement, which can have a large impact on data movement. On the other hand, thread placement techniques mainly focus on non-uniform memory architectures (NUMA) [7,15,28,56,58,63] reduce the distance to data, and use policies, such as clustering, that do not translate well to NUCA. Applying these techniques separately leaves significant performance on the table.

## 1.2    Contribution: computation and data co-scheduling

The focus of this work is to investigate and design hardware and software mechanisms for data-centric computing. We treat the hardware NUCA cache as a resource pool with different capacity and latency tradeoffs, and build *virtual caches* from these resources that minimizes the application's data movement. This design takes a cross-layer approach, in which hardware monitors applications' performance and virtualizes the cache, and software uses a simple but accurate performance model to place thread and find the best virtual cache hierarchy. A cross-layer approach is a must for data-centric computation: hardware alone cannot optimize the full system, while software alone cannot control the cache without help from hardware. We believe this design will improve performance and energy efficiency substantially over conventional, hardware-only caches with low overheads.

With the cross-layer method in mind, this thesis proposes CDCS, a technique that manages cache capacity and schedules threads at the same time to reduce data movement. By taking both data allocation and access intensity into account, a prototype of CDCS is proposed to jointly place threads and data across CMP tiles, depicted in Fig. 1-1. CDCS builds on a partitioned NUCA baseline and reconfigures the system periodically. It takes the cross-layer approach, in which hardware monitors applications' performance and provides flexible configuration technique, and software decides how to place threads and data. Since the joint placement is a complex, multi-

Figure 1-1: Example configuration of CDCS system with five running applications

dimensional optimization problem, this thesis proposes a heuristic-based, hierarchical optimization method to solve the problem efficiently. The proposed design works on arbitrary mixes of single- and multi-threaded processes and reconfigures the system with small overheads.

## 1.3    Thesis Structure

This thesis is organized as follows: Chapter 2 discusses prior work related to computation and data co-scheduling. Chapter 3 introduces the baseline architecture used in both CDCS. Chapter 4 cover the design and implementation of CDCS. Chapter 5 presents the methodology used in the thesis. Chapter 6 shows the evaluation for CDCS with different benchmarks and that those proposed techniques improve performance and energy efficiency. Chapter 7 concludes the thesis.

# Chapter 2

# Related work

This section discusses prior work related to computation and data co-scheduling. First, we discuss related work about data placement in multicore last-level caches (LLCs). Next, we present a case study that compares different NUCA schemes and shows that thread placement significantly affects performance. Finally, we present prior work on thread placement in different contexts.

## 2.1  Data placement techniques to reduce data movement

**Non-uniform cache architectures**  NUCA techniques [33] reduce access latency of distributed caches, and are concerned with data placement, but do not place threads or divide cache capacity among them. Static NUCA (S-NUCA) [33] spreads data across banks with a fixed line-bank mapping, and exposes a variable bank latency. Commercial CMPs often use S-NUCA [34]. Dynamic NUCA (D-NUCA) schemes adaptively place data close to the requesting core [2,3,8,10,11,20,27,41,50,62] using a mix of placement, migration, and replication techniques. Placement and migration bring lines close to cores that use them, possibly introducing capacity contention between cores depending on thread placement. Replication makes multiple copies of frequently used lines, reducing latency for widely read-shared lines (e.g., hot code) at the expense of some capacity loss.

Most D-NUCA designs build on a *private-cache baseline*, where each NUCA bank

is treated as a private cache. All banks are under a coherence protocol, which makes such schemes either hard to scale (in snoopy protocols) or require large directories that incur significant area, energy, latency, and complexity overheads (in directory-based protocols). To avoid these costs, some D-NUCA schemes instead build on a *shared-cache baseline*: banks are not under a coherence protocol, and virtual memory is used to place data. Cho and Jin [11] use page coloring and NUCA-aware allocation to map pages to banks. R-NUCA [20] specializes placement and replication policies for different data classes (instructions, private data, and shared data), on a per-page basis, outperforming prior D-NUCA schemes. Shared-baseline schemes are cheaper, as LLC data does not need coherence, and have a simpler lookup mechanism. However, remapping data is expensive as it requires page copies and invalidations.

Conventional D-NUCA techniques are chiefly concerned with data placement, and do not explicitly allocate capacity. This incurs unnecessary misses and degrades performance [4] and cannot achieve isolation or quality of service.

**Partitioned shared caches**  Partitioning enables software to explicitly allocate cache space among threads or cores, but it is not concerned with data or thread placement. Partitioning can be beneficial because applications vary widely in how well they use the cache. Cache arrays can support multiple partitions with small modifications [9, 37, 39, 52, 55, 61]. Software can then set these sizes to maximize throughput [51], or to achieve fairness [43], isolation and prioritization [12, 18, 31], and security [46]. Unfortunately, partitioned caches scale poorly because they do not optimize placement. Moreover, they allocate capacity to cores, which works for single-threaded mixes, but incorrectly accounts for shared data in multi-threaded workloads.

**Partitioned NUCA**  Recent work has developed techniques to perform spatial partitioning of NUCA caches. These schemes jointly consider data allocation and placement, reaping the benefits of NUCA and partitioned caches. However, they do not consider thread placement. Virtual Hierarchies rely on a logical two-level directory to partition the cache [40], but they only allocate full banks, double directory overheads, and make misses slower. CloudCache [35] implements virtual private

18

caches that can span multiple banks, but allocates capacity to cores, needs a directory, and uses broadcasts, making it hard to scale. Jigsaw [4] is a shared-baseline NUCA with partitionable banks and single-lookup accesses. Jigsaw lets software divide the distributed cache in finely-sized *virtual caches*, place them in different banks, and map pages to different virtual caches. Using utility monitors [51], an OS-level software runtime periodically gathers the miss curve of each virtual cache and co-optimizes data allocation and placement. Since Jigsaw let software control data placement, techniques proposed in this project will build on Jigsaw to leverage its data placement technique.

## 2.2   Case study: Tradeoffs in thread and data placement

To explore the effect of thread placement on different NUCA schemes, we simulate a 36-core CMP running a specific mix. The CMP is a scaled-down version of the 64-core chip in Fig. 3-1, with 6×6 tiles. Each tile has a simple 2-way OOO core and a 512 KB LLC bank.

We run a mix of single- and multi-threaded workloads. From single-threaded SPEC CPU2006, we run six instances of omnet (labeled O1-O6) and 14 instances of milc (M1-M14). From multi-threaded SPEC OMP2012, we run two instances of ilbdc (labeled I1 and I2) with eight threads each. We choose this mix because it illustrates the effects of thread and data placement—chapter 6 uses a comprehensive set of benchmarks.

Fig. 2-1 shows how thread and data are placed across the chip under different schemes. Each square represents a tile. The label on each tile denotes the thread scheduled in the tile's core (labeled by benchmark as discussed before). The colors on each tile show a breakdown of the data in the tile's bank. Each process uses the same color for its threads and data. For example, in Fig. 2-1b, the upper-leftmost tile has thread O1 (colored blue) and its data (also colored blue); data from O1 also occupies parts of the top row of banks (portions in blue).

Fig. 2-1a shows the thread and data breakdown under R-NUCA when applications

19

**Legend**

Tile (1 core+LLC bank) →
Thread running on core → I1
LLC data breakdown →

**Example**
I1 (ilbdc) thread
Data from O1 and O2

ilbdc — Threads I1 x 8, I2 x 8; Data

omnet — Threads O1 O2 O3 O4 O5 O6; Data

milc — Threads M1 M2 ... M14; Data

(a) R-NUCA

(b) Jigsaw+Clustered

(c) Jigsaw+Random

(d) CDCS

Figure 2-1: Case study: 36-tile CMP with a mix of single- and multi-threaded workloads (omnet×6, milc×14, 8-thread ilbdc×2) under different NUCA organizations and thread placement schemes. Threads are labeled and data is colored by process.

are grouped by type (e.g., the six copies of omnet are in the top-left corner). R-NUCA maps thread-private data to each thread's local bank, resulting in very low latency. Banks also have some of the shared data from the multithreaded processes I1 and I2

20

Figure 2-2: Application miss curves.

Table 2.1: Per-app and weighted speedups for the mix studied.

|  | omnet | ilbdc | milc | WSpdp |
|---|---|---|---|---|
| R-NUCA | 1.09 | 0.99 | 1.15 | 1.08 |
| Jigsaw+Cl | 2.88 | 1.40 | 1.21 | 1.48 |
| Jigsaw+Rnd | 3.99 | 1.20 | 1.21 | 1.47 |
| CDCS | 4.00 | 1.40 | 1.20 | 1.56 |

(shown hatched), because R-NUCA spreads shared data across the chip. Finally, code pages are mapped to different banks using rotational interleaving, though this is not visible in this mix because apps have small code footprints. These policies excel at reducing LLC access latency to private data vs. an S-NUCA cache. This helps `milc` and `omnet`, as shown in Fig. 2.1. Overall, R-NUCA speeds up this mix by 8% over S-NUCA.

In R-NUCA, other thread placements would make little difference for this mix, as most capacity is used for either thread-private data, which is confined to the local bank, or shared data, which is spread out across the chip. But R-NUCA does not use capacity efficiently in this mix. Fig. 2-2 shows why, giving the miss curves of each app. Each miss curve shows the misses per kilo-instruction (MPKI) that each process incurs as a function of LLC space (in MB). `omnet` is very memory-intensive, and suffers 85 MPKI below 2.5 MB. However, over 2.5 MB, its data fits in the cache and misses turn into hits. `ilbdc` is less intensive and has a smaller footprint of 512 KB. Finally, `milc`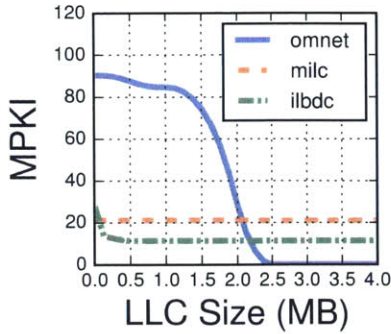 gets no cache hits no matter how much capacity it is given—it is a streaming application. In R-NUCA, `omnet` and `milc` apps get less than 512 KB, which does not benefit them, and `ilbdc` apps use more capacity than they need.

Jigsaw uses capacity more efficiently, giving 2.5 MB to each instance of `omnet`, 512 KB to each `ilbdc` (8 threads), and near-zero capacity to each `milc`. Fig. 2-1b shows how Jigsaw tries to place data close to the threads that use it. By using partitioning, Jigsaw can share banks among multiple types of data without introducing capacity interference. However, the `omnet` threads in the corner heavily contend for capacity of neighboring banks, and their data is placed farther away than if they were spread out.

Clearly, *when capacity is managed efficiently, thread placement has a large impact on capacity contention and achievable latency*. Nevertheless, because omnet's data now fits in the cache, its performance vastly improves, by 2.88× over S-NUCA (its AMAT improves from 15.2 to 3.7 cycles, and its IPC improves from 0.22 to 0.61). ilbdc is also faster, because its shared data is placed close by instead of across the chip; and because omnet does not consume memory bandwidth anymore, milc instances have more of it and speed up moderately (Fig. 2.1). Overall, Jigsaw speeds up this mix by 48% over S-NUCA.

Fig. 2-1c shows the effect of randomizing thread placement to spread capacity contention among the chip. omnet instances now have their data in neighboring banks (1.2 hops on average, instead of 3.2 hops in Fig. 2-1b) and enjoy a 3.99× speedup over S-NUCA. Unfortunately, ilbdc's threads are spread further, and its performance suffers relative to clustering threads (Fig. 2.1). This shows why one policy does not fit all: depending on capacity contention and sharing behavior, apps prefer different placements. Specializing policies for single- and multithreaded apps would only be a partial solution, since multithreaded apps with large per-thread footprints and little sharing also benefit from spreading.

Finally, Fig. 2-1d shows how CDCS handles this mix. CDCS spreads omnet instances across the chip, avoiding capacity contention, but clusters ilbdc instances across their shared data. CDCS achieves a 4× speedup for omnet and a 40% speedup for ilbdc. CDCS speeds up this mix by 56%.

In summary, this case study shows that partitioned NUCA schemes use capacity more effectively and improve performance, but they are sensitive to thread placement, as threads in neighboring tiles can aggressively contend for capacity. This presents an opportunity to perform smart thread placement, but fixed policies have clear shortcomings.

## 2.3 Thread placement to reduce data movement

**NUMA thread scheduling** Prior work has studied NUMA-aware thread scheduling and migration in multi-socket systems. Tam et al. [56] profile which threads have frequent sharing and place them in the same socket. DINO [7] clusters single-threaded

22

processes to equalize memory intensity, places these clusters in different sockets, and migrates pages along with threads. While DINO directly considers the relationship between thread and data placement, it does not address applications with data sharing. Linux will consider swapping two threads across NUMA nodes instead of migrating their remote pages [13]. Multi-socket NUMA systems have simple off-chip topologies (e.g., dance-hall), and thread clustering is all-or-nothing proposition: threads either are in the same socket and enjoy fast sharing, or they are not and suffer slow shared data accesses. In contrast, on-chip networks (e.g., meshes) have more continuous latency profiles, which make the spatial mapping of thread to cores more lenient but computationally harder. In the context of on-chip NUMA systems, Tumanov et al. [58] and Das et al. [15] profile memory accesses and schedule high-intensity threads closer to their memory controller. As in multi-socket systems, the goal of these schemes is to minimize network latency while balancing controller load. These NUMA schemes focus on equalizing memory bandwidth, whereas in this project, proper cache allocations cause capacity contention to be the main constraint on thread placement in NUCA.

**Cache contention-aware thread scheduling** Rather than focus on bandwidth limitation in NUMA, this thesis proposes to use thread placement to address cache contention problem in multicore system. Among few prior work that shares same motivation, CRUISE [28] is the closest work to the technique proposed in the thesis. CRUISE schedules single-threaded apps in CMPs with multiple fixed-size last-level caches, each shared by multiple cores and unpartitioned. CRUISE takes a classification-based approach, dividing apps in thrashing, fitting, friendly, and insensitive, and applies fixed scheduling policies to each class (spreading some classes among LLCs, using others as filler, etc.). CRUISE bin-packs apps into fixed-size caches, but partitioned NUCA schemes provide flexibly sized virtual caches that can span multiple banks. It is unclear how CRUISE's policies and classification would apply to NUCA. CRUISE improves on DI [63], which profiles miss rates and schedules apps across chips to balance intensive and non-intensive apps. Both schemes only consider single-threaded apps, and have to contend with the lack of partitioned LLCs. In contrast to CRUISE, CDCS, the first technique proposed in this thesis, addresses capacity contention problem for NUCA within single socket, and is applicable to both single- and multi-threaded

applications.

**Static schedulers** Prior work has used high-quality *static mapping* techniques to spatially decompose regular parallel problems. For example, integer linear programming is useful in spatial architectures [44], and graph partitioning has been extensively used in stream scheduling [48, 49]. While some of these techniques can be applied to the thread and data placement problem, they are too expensive to use dynamically.

## 2.4  Summary

All the prior work mentioned in the previous two sections tries to reduce data movement by either placing threads close to their data or placing data close threads. CDCS, instead, aims to consider thread and data placement at the same time and thus is different from prior work. Moreover, papers relate to thread and data co-scheduling deal with balancing the memory requirements of multiple co-scheduled threads over multiple intervals. In contrast, CDCS deals with the spatial scheduling and resource management aspects of large-scale, non-uniform, non-overcommitted multi-cores.

# Chapter 3

# Baseline architecture

CDCS builds on Jigsaw [4], a partitioned NUCA technique. Fig. 3-1 shows the tiled CMP architecture and the hardware additions of Jigsaw. Each tile has a core and a slice of the LLC. An on-chip network of arbitrary topology connects tiles, and memory controllers are at the edges. Pages are interleaved across memory controllers, as in Tilera and Knights Corner chips [6].

## 3.1 Virtual caches

Jigsaw lets software divide each cache bank in multiple partitions and uses Vantage [52] to efficiently partition banks at cache-line granularity. Collections of bank partitions are ganged and exposed to software as a single VC (originally called a *share* in Jigsaw). This allows software to define many VCs cheaply (several per thread), and to finely size and place them among banks.

## 3.2 Mapping data to VCs

Unlike other D-NUCAs, in Jigsaw lines do not migrate in response to accesses. Instead, between reconfigurations, each line can only reside in a single LLC bank. Jigsaw maps data to VCs using the virtual memory subsystem, similar to R-NUCA [20]. Each page table entry is tagged with a VC id. On an L2 miss, Jigsaw uses the line address and its VC id to determine the bank and bank partition that the line maps to.

**64-tile CMP**

Mem / IO

NoC (Mesh)    Tile

Mem / IO

Mem / IO

Mem / IO

**Tile Organization**

◻ Modified structures
◻ New/added structures

**Jigsaw L3 Bank**

Bank partitioning HW

Monitoring HW | Bulk inv HW

NoC Router

L2 ↔ STB

L1I | L1D

Core | TLBs

**Share-bank Translation Buffer**

Share Id (from TLB)        Address (from L1 miss)

2706                          0x5CA1AB1E
                                    H

STB Entry

4 entries, associative,
exception on miss

Bank/Part 0          Bank/Part N-1
1/3 | 3/5 | 1/3  ...  0/8

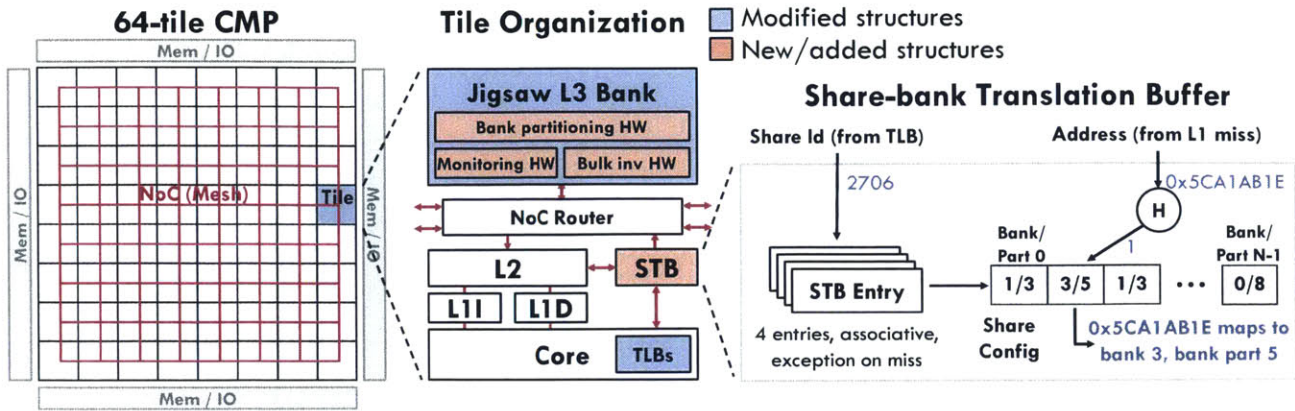Share Config    0x5CA1AB1E maps to
                bank 3, bank part 5

Figure 3-1: Target CMP, with tile configuration and microarchitectural additions introduced for Jigsaw. Jigsaw gangs portions of banks into virtual caches, and uses the VTB to find the bank and bank partition to use on each access.

The *virtual-cache translation buffer* (VTB), shown in Fig. 3-1, determines the bank and bank partition for each access. The VTB stores the configuration of all VCs that the running thread can access [4]. In the implementation of Jigsaw, it is a 3-entry lookup table, as each thread only accesses 3 VCs (as explained below).

Each VTB entry contains a VC descriptor, which consists of an array of $N$ bank and bank partition ids (in our implementation, $N = 64$ buckets). As shown in Fig. 3-1, to find the bank and bank partition ids, the address is hashed, and the hash value (between 0 and $N - 1$) selects the bucket. Hashing allows spreading accesses across the VC's bank partitions in proportion to their capacities, which makes them behave as a cache of their aggregate size. For example, if a VC consists of two bank partitions $A$ of 1 MB and $B$ of 3 MB, by setting array elements 0–15 in the VC descriptor to $A$ and elements 16–63 to $B$, $B$ receives $3\times$ more accesses than $A$. In this case, $A$ and $B$ behave like a 4 MB VC [4,5].

Periodically (e.g., every 25 ms), software changes the configuration of some or all VCs, changing both their bank partitions and sizes. The OS recomputes all VC descriptors based on the new data placement, and cores coordinate via inter-processor interrupts to update the VTB entries simultaneously. The two-level translation of pages to VCs and VCs to bank partitions allows Jigsaw to be more responsive and take more drastic reconfigurations than prior shared-baseline D-NUCAs: reconfigurations simply require changing the VC descriptors, and software need not copy pages or alter page table entries.

26

## 3.3 Types of VCs

Jigsaw's OS-level runtime creates one thread-private VC per thread, one per-process VC for each process, and a global VC. Data accessed by a single thread is mapped to its thread-private VC, data accessed by multiple threads in the same process is mapped to the per-process VC, and data used by multiple processes is mapped to the global VC. Pages can be reclassified to a different VC efficiently [4] (e.g., when a page in a per-thread VC is accessed by another thread, it is remapped to the per-process VC), though in steady-state this happens rarely.

## 3.4 Summary

Jigsaw provides a flexible mechanism to place data across the chip. It enables CDCS to freely decide where to place data. Therefore, the proposed technique can focus on how to leverage this mechanism to further perform data and thread co-scheduling and optimize the system.

# Chapter 4

# CDCS: computation and data co-scheduling

## 4.1 Overview

CDCS uses a combination of hardware and software techniques to perform joint thread and data co-scheduling. Fig. 4-1 gives an overview of the steps involved. Novel, scalable performance monitors sample the miss curves of each virtual cache. OS runtime periodically reads these miss curves and uses them to jointly place VCs and threads using a 4-step procedure. Finally, this runtime uses hardware support to move cache lines to their new locations. This hardware addresses overheads that would hinder CDCS performance, especially on large systems.

Although CDCS leverages the data placement mechanism from Jigsaw, all aspects of this optimization process differ from Jigsaw [4]. Jigsaw uses a simple runtime that sizes VCs obliviously to their latency, places them greedily, and does not place threads. Jigsaw also uses conventional utility monitors [51] that do not scale to large caches, and reconfigurations require long pauses while banks invalidate data, which adds jitter.

All of the techniques are *topology-agnostic*; our algorithms accept arbitrary distance vectors between tiles. However, to make the discussion concrete, the following discussion assumes a mesh topology.
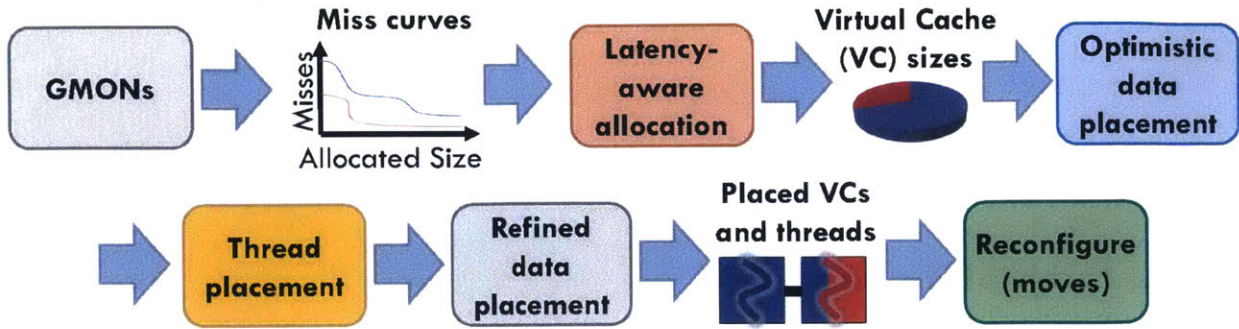
Figure 4-1: Overview of CDCS's periodic reconfiguration procedure

## 4.2 A simple cost model for thread and data placement

As discussed in Sec. 2.3, classification-based heuristics are hard to apply to NUCA. Instead, CDCS uses a simple analytical cost model that captures the effects of different placements on *total memory access latency*, and uses it to find a low-cost solution. This latency is better analyzed as the sum of on-chip (L2 to LLC) and off-chip (LLC to memory) latencies.

**Off-chip latency**  Assume a system with $T$ threads and $D$ VCs. Each thread $t$ accesses VC $d$ at a rate $a_{t,d}$ (e.g., $50\,\mathrm{K}$ accesses in $10\,\mathrm{ms}$). If VC $d$ is allocated $s_d$ lines in the cache, its miss ratio is $M_d(s_d)$ (e.g., 10% of accesses miss). Then the total off-chip access latency is:

$$\text{Off-chip latency} = \sum_{t=1}^{T} \sum_{d=1}^{D} a_{t,d} \times M_d(s_d) \times \text{MemLatency} \tag{4.1}$$

where MemLatency is the latency of a memory access. This includes network latency, and relies on the average distance of all cores to memory controllers being the same (Fig. 3-1). Extending CDCS to NUMA would require modeling a variable main memory latency in Eq. 4.1.

**On-chip latency**  Each VC's capacity allocation $s_d$ consists of portions of the $N$ banks on chip, so that $s_d = \sum_{b=1}^{N} s_{d,b}$. The capacity of each bank $B$ constrains allocations, so that $B = \sum_{d=1}^{D} s_{d,b}$. Limited bank capacities sometimes force data to

30

be further away from the threads that access it. Because the VTB spreads accesses across banks in proportion to capacity, the number of accesses from thread $t$ to bank $b$ is $\alpha_{t,b} = \sum_{d=1}^{D} \frac{s_{d,b}}{s_d} \times a_{t,d}$. If thread $t$ is placed in a core $c_t$ and the network distance between two tiles $t_1$ and $t_2$ is $D(t_1, t_2)$, then the on-chip latency is:

$$\text{On-chip latency} = \sum_{t=1}^{T} \sum_{b=1}^{N} \alpha_{t,b} \times D(c_t, b) \tag{4.2}$$

## 4.3 Overview of CDCS optimization

With this cost model, the computation and data co-scheduling problem is to choose the $c_t$ (thread placement) and $s_{t,b}$ (VC size and data placement) that minimize total latency, subject to the given constraints. However, finding the optimal solution is NP-hard [24, 47], and different factors are intertwined. For example, the size of VCs and the thread placement affect how close data can be placed to the threads that use it. Therefore, CDCS takes a multi-step approach to disentangle these interdependencies. CDCS first adopts optimistic assumptions about the contention introduced by thread and data placement, and gradually refines them to produce the final placement. Specifically, this thesis proposes to reconfigure the system with four steps:

1. *Latency-aware allocation* divides capacity among VCs assuming all threads that access each VC are placed as close as possible and data is compactly placed (no capacity contention).

2. *Optimistic contention-aware VC placement* places VCs among banks to avoid capacity contention. This step produces a rough picture of where data should be in the chip.

3. *Thread placement* uses the optimistic VC placement to place threads close to the VCs they access. For example, this step places single-threaded applications close to the center of mass of their data, and clusters threads that shares data around their shared data.

4. *Refined VC placement* improves on the previous data placement to, now that thread locations are known, place data closer to minimize on-chip latency. For example, a thread that accesses its data intensely may swap allocations with a less intensive thread to bring its data closer; while this increases latency for the other thread,

31

overall it is beneficial.

By considering data placement twice (steps 2 and 4), CDCS disentangles the circular relationship between thread and data placement and therefore performs efficient thread and data co-placement. This design can find high-quality solution for the latency-minimization problem within a couple of milliseconds. With the solution, CDCS reconfigures the system to minimize potential data movement and further improve energy efficiency and performance of multicore processors.



Figure 4-2: Access latency vs. capacity allocation.



Figure 4-3: Optimistic uncontended virtual cache placement.



(a) Partial optimistic data placement

(b) Estimating contention for VC

(c) VC placed near least contended tile

Figure 4-4: Optimistic, contention-aware virtual cache placement.

## 4.4    Latency-aware capacity allocation

Cache misses generally decrease with cache capacity, but not all threads see equal benefits from additional capacity. For example, streaming applications access memory

frequently and consume most of the capacity in unpartitioned caches, but they have no reuse and do not benefit from this capacity. Prior work has partitioned cache capacity to reduce cache misses [4,51], i.e. off-chip latency. However, it is well-known that larger caches take longer to access [22,23,57]. Most prior partitioning work has targeted fixed-size LLCs with constant latency. But capacity allocation in NUCA caches provides an opportunity to also reduce on-chip latency: if an application sees little reduction in misses from a larger VC, the additional network latency to access it can negate the benefits of having fewer misses.

In summary, larger allocations have two competing effects: decreasing off-chip latency and increasing on-chip latency. This is illustrated in Fig. 4-2, which shows the average memory access latency to one VC (e.g., a thread's private data). Fig. 4-2 breaks latency into its off- and on-chip components, and shows that there is a "sweet spot" that minimizes total latency.

This has two important consequences. First, unlike in other D-NUCA schemes, it is sometimes better to leave cache capacity *unused*. Second, incorporating on-chip latency changes the curve's shape and also its marginal utility [51], leading to different cache allocations even when all capacity is used.

CDCS allocates capacity from *total memory latency curves* (the sum of Eq. 4.1 and Eq. 4.2) instead of miss curves. These curves estimate the total memory latency from accessing one type of data over some interval, say the next 25 ms. However, Eq. 4.2 requires knowing the thread and data placements, which are unknown at the first step of reconfiguration. CDCS instead uses an optimistic on-chip latency curve, found by the *minimum-latency, uncontended placement*, and compactly place the VC around the center of the chip and computing the resulting average latency. For example, Fig. 4-3 shows the optimistic placement of an 8.2-bank VC accessed by a single thread, with an average distance of 1.27 hops. Optimistic data placement gives a lower bound on the on-chip latency curve. CDCS computes the off-chip latency curve by multiplying the miss curve by the main memory access latency, and add both curves to produce the total memory latency curve.

With this simplification, CDCS uses the Peekahead optimization algorithm [4] to efficiently find the sizes of all VCs that minimize latency. While these allocations account for on-chip latency, they generally underestimate it due to capacity contention.

Nevertheless, this scheme works well because the next steps are effective at limiting contention.

## 4.5 Optimistic contention-aware VC placement

Once VC sizes are known, CDCS first finds a rough picture of how data should be placed around the chip to avoid placing large VCs close to each other. The main goal of this step is to *inform thread placement* by avoiding VC placements that produce high capacity contention.

To this end, CDCS sort VCs by size and place the largest ones first. Intuitively, this works well because larger VCs can cause more contention, while small VCs can fit in a fraction of a bank and cause little contention. For each VC, the algorithm iterates over all the banks, and chooses the bank that yields the least contention with already-placed VCs as the center of mass of the current VC. To make this search efficient, CDCS approximate contention by keeping a running tally of *claimed capacity* in each bank, and relax capacity constraints, allowing VCs to claim more capacity than is available at each bank. With $N$ banks and $D$ VCs, the algorithm runs in $\mathcal{O}(N \cdot D)$.

Fig. 4-4 shows an example of optimistic contention-aware VC placement at work. Fig. 4-4a shows claimed capacity after two VCs have been placed. Fig. 4-4b shows the contention for the next VC at the center of the mesh (hatched), where the uncontended placement is a cross. Contention is approximated as the claimed capacity in the banks covered by the hatched area—or 3.6 in this case. To place a single VC, CDCS computes the contention from placing the VC centered around that tile. CDCS then place the VC around the tile that had the lowest contention, updating the claimed capacity accordingly. For instance, Fig. 4-4c shows the final placement for the third VC in our example.

Optimistic data placement consists of a triply nested loop over VCs, tiles, and neighboring banks. The run-time of this algorithm is nonetheless $\mathcal{O}(N \cdot T)$, where $N$ and $T$ are the number of cores and threads. The critical insight is that its impossible for $\mathcal{O}(T)$ threads to each require $\mathcal{O}(N)$ banks to place their data, since that would imply a total allocation greatly exceeding cache capacity. In fact across all threads,

the total allocation equals $N$ banks exactly. Thus while placing a large VC can take $\mathcal{O}(N \cdot T)$ time, the vast majority of VCs only check their local bank and are placed in $\mathcal{O}(N)$. Summing over all VCs, the total number of tiles checked is thus $\mathcal{O}(N \cdot T)$, and this computation overhead is evaluated in Sec. 6.3.

## 4.6 Thread placement

Given the previous data placement, CDCS tries to place threads closest to the center of mass of their accesses. Recall that each thread accesses multiple VCs, so this center of mass is computed by weighting the centers of mass of each VC by the thread's accesses to that VC. Placing the thread at this center of mass minimizes its on-chip latency (Eq. 4.2).

Unfortunately, threads sometimes have the same centers of mass. To break ties, CDCS places threads in descending *intensity-capacity product* (sum of VC accesses $\times$ VC size for each VC accessed). Intuitively, this order prioritizes threads for which low on-chip latency is important, and for which VCs are hard to move.

For multithreaded workloads, this approach clusters shared-heavy threads around their shared VC, and spreads private-heavy threads to be close to their private VCs. Should threads access private and shared data with similar intensities, CDCS places threads relatively close to their shared VC but does not tightly cluster them, avoiding capacity contention among their private VCs.

## 4.7 Refined VC placement

Finally, CDCS performs a round of detailed VC placement to reduce the distance between threads and their data.

CDCS first simply round-robins VCs, placing capacity as close to threads as possible without violating capacity constraints. This greedy scheme, which was used in Jigsaw [4], is a reasonable starting point, but produces sub-optimal placements. For example, a thread's private VC always gets space in its local bank, regardless of the thread's memory intensity. Also, shared VCs can often be moved at little or no cost to make room for data that is more sensitive to placement. This is because moving
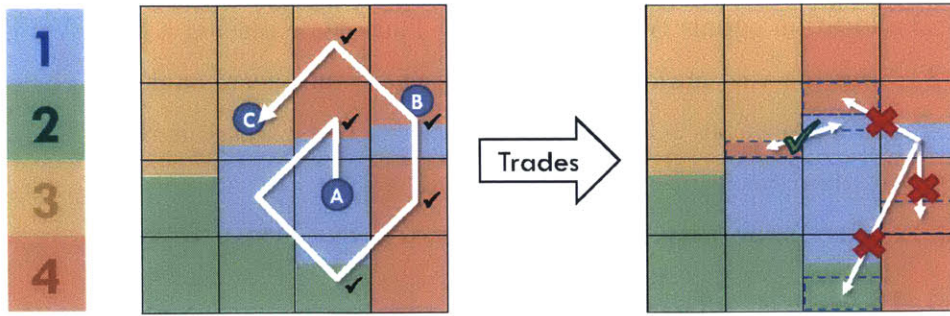
Figure 4-5: Trading data placement: Starting from a simple initial placement, VCs trade capacity to move their data closer. Only trades that reduce total latency are permitted.

shared data farther away from one accessing thread often moves it closer to another.

Furthermore, unlike in previous steps, it is straightforward to compute the effects of moving data, since there is an initial placement to compare against. CDCS therefore looks for beneficial trades between pairs of VCs after the initial, greedy placement. Specifically, CDCS computes the latency change from trading capacity between $VC_1$ at bank $b_1$ and $VC_2$ at bank $b_2$ using Eq. 4.2. The change in latency for $VC_1$ is:

$$\Delta \text{Latency} = \frac{\text{Accesses}}{\text{Capacity}} \times \left( D(VC_1, b_1) - D(VC_1, b_2) \right)$$

The first factor is $VC_1$'s accesses per byte of allocated capacity. Multiplying by this factor accounts for the number of accesses that are affected by moving capacity, which varies between VCs. The equation for $VC_2$ is similar, and the net effect of the trade is their sum. If the net effect is negative (lower latency is better), then the VCs swap bank capacity.

Naïvely enumerating all possible trades is prohibitively expensive, however. Instead, CDCS performs a bounded search by iterating over all VCs: Each VC spirals outward from its center of mass, iterating over banks in order of distance while trying to move its data closer. At each bank $b$ along the outward spiral, if the VC has not claimed all of $b$'s capacity then it adds $b$ to a list of desirable banks. These are the banks it will try to trade into later. Next, the VC tries to move its data placed in $b$ (if any) closer by iterating over closer, desirable banks and offering trades with VCs that have data in these banks. If the trades are beneficial, they are performed. The spiral terminates when the VC has seen all of its data, since no farther banks will allow it to move any

36

data closer.

Fig. 4-5 illustrates this for an example CMP with four VCs and some initial data placement. CDCS spirals outward starting from VC1's center of mass at bank A, and terminate at VC1's farthest data at bank C. Desirable banks are marked with black checks on the left of Fig. 4-5. CDCS only attempts a few trades, shown on the right side of Fig. 4-5. At bank B, VC1's data is two hops away, so CDCS tries to trade it to any closer, marked bank. For illustration, suppose none of the trades are beneficial, so the data does not move. This repeats at bank C, but suppose the first trade is now beneficial. VC1 and VC4 trade capacity, moving VC1's data one hop closer.

This approach gives every VC a chance to improve its placement. Since any beneficial trade must benefit one party, it would discover all beneficial trades. However, for efficiency, in CDCS each VC trades only once, since CDCS has empirically found this discovers most trades. Finally, this scheme incurs negligible overheads, as shown in Sec. 6.3. Limiting the search space in this way greatly improves performance, by $10\times$ over naïve enumeration at 64 cores, because (similar to thread placement) CDCS only considers a few nearby banks for the vast majority of VCs and even large VCs only consider a few trades. Computation overhead for this scheme is negligible, as in Sec. 6.3.

These techniques are cheap and effective. We also experimented with more expensive approaches commonly used in placement problems: integer linear programming, simulated annealing, and graph partitioning. Sec. 6.3 shows that they yield minor gains and are too expensive to be used online.

## 4.8   Monitoring large caches

Monitoring miss curves in large CMPs is challenging. To allocate capacity efficiently, it is important to manage it in small chunks (e.g., the size of the L1s) so that it isn't over-allocated where it produces little benefit. This is crucial for VCs with small working sets, which see large gains from a small size and no benefit beyond. Yet, miss curves that cover the full LLC are needed because a few VCs may benefit from taking most capacity. These two requirements—fine granularity and large coverage—are problematic for existing monitors.
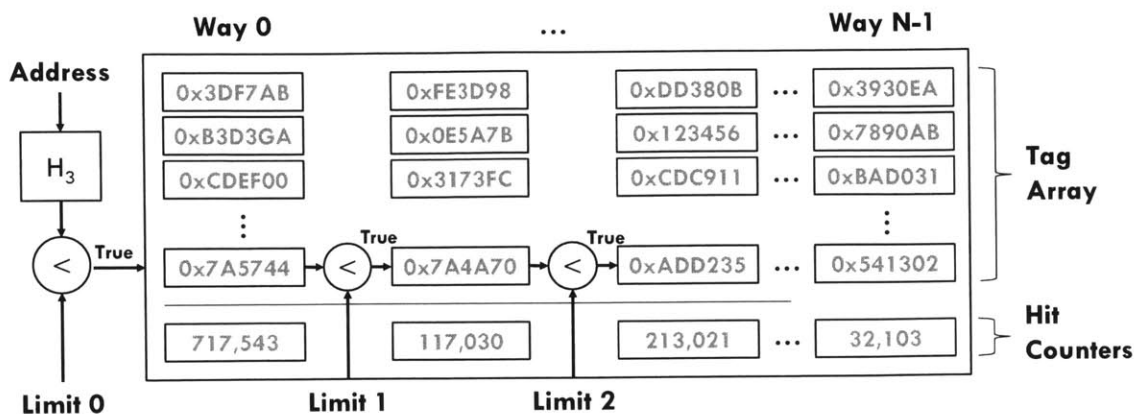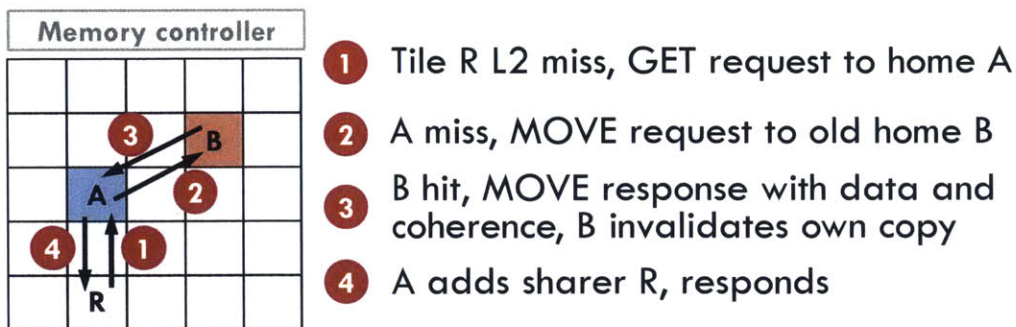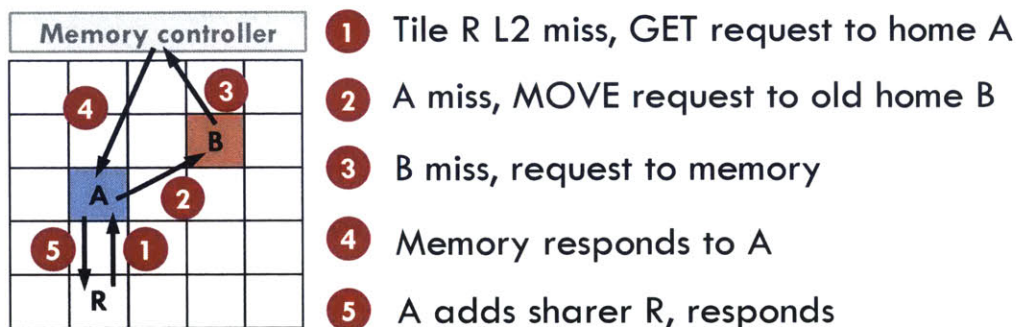
Figure 4-6: GMONs enhance UMONs with varying sampling rate across ways, controlled with per-way *limit registers*.

Conventional cache partitioning techniques use utility monitors (UMONs) [51] to monitor a fraction of sets, counting hits at each way to gather miss curves. UMONs monitor a fixed cache capacity per way, and would require a prohibitively large associativity to achieve both fine detail and large coverage. Specifically, in an UMON with $W$ ways, each way models $1/W$ of LLC capacity. With a 32 MB LLC (chapter 5, Table 5.1) if we want to allocate capacity in 64 KB chunks, a conventional UMON needs 512 ways to have enough resolution. This is expensive to implement, even for infrequently used monitors.

Instead, a novel monitor, called a *geometric monitor* (GMON) is developed. GMONs need fewer ways—64 in our evaluation—to model capacities from 64 KB up to 32 MB. This is possible because GMONs vary the sampling rate across ways, giving both fine detail for small allocations and large coverage, while using many fewer ways than conventional UMONs. Fig. 4-6 shows this design: A GMON consists of small set-associative, tag-only array. Instead of storing address tags, GMON tags store 16-bit hashed addresses. GMONs also have a *limit register* per way. The limit registers progressively decrease across ways, and are used to filter out a fraction of lines per way as follows. In a conventional UMON, when an address is inserted or moved up to the first way, all other tags are moved to the next way. (This requires potentially shifting as many tags as ways, but only a small fraction of accesses are monitored, so the energy impact is small.) In a GMON, on each move, the hash value of the tag is checked against the way's limit register. If the value exceeds the limit register, the tag is discarded instead of moved to the next way, and the process

38

(a) Demand move, old bank hit.



(b) Demand move, old bank miss.

Figure 4-7: Messages and protocol used on incremental reconfigurations: demand moves when old bank hits or misses.

terminates. Discarding lines achieves a variable, decreasing sampling rate per way, so GMONs model an increasing capacity per way [4, 5, 32]. This increases sampling rate geometrically across ways, and since the size of the modeled cache is inversely proportional to the sampling rate, the coverage increases geometrically as well.

Limit registers are set to decrease the sampling rate by a factor $\gamma < 1$, so the sampling rate at way $w$ is $k_w = \gamma^w$ less than at way zero, and then choose $\gamma$ to cover the full cache capacity. For example, with a 32 MB LLC, a 64-way GMON with $\gamma \approx 0.95$ covers the full cache while having the first way model 64 KB. Modeled capacity per way grows by 26×, from 0.125 to 3.3 banks. This means GMONs miss curves are sparse, with high resolution at small sizes, and reduced resolution at large sizes. We find these GMONs work as well as the impractical UMONs described above (Sec. 6.3).

# 4.9 Incremental reconfigurations

Moving threads and data reduces steady-state network latency, but requires more drastic reconfigurations. Jigsaw reconfigures through bulk invalidations: all bank controllers walk the tag array and invalidate lines that should be mapped somewhere else, which requires pausing cores for tens to hundreds of thousands of cycles [4]. This is simple, but pauses, extra writebacks, and misses hurt performance. Because CDCS reconfigures more frequently, this approach can significantly degrade performance. With few hardware modifications, it is possible to spatially reconfigure the cache incrementally, without pausing cores, by moving lines instead of invalidating them.

The key idea is to, upon a reconfiguration, temporarily treat the cache as a two-level hierarchy. We add a *shadow VC descriptor* to each VTB entry, as shown in Fig. 3-1. Upon reconfiguration, each core copies the VC descriptors into the shadow descriptors, and updates the normal VC descriptors with the new configuration. When the shadow descriptors are active, the VTB finds both the current and previous banks and bank partitions for the line, and, if they are different, sends the old bank id along with its request. Fig. 4-7 illustrates this protocol. If the current bank misses, it forwards the request to the old bank instead of memory. If the old bank hits, it sends both the line and its coherence state to the new bank, invalidating its own copy. This moves the line to its new location. We call this a *demand move* (Fig. 4-7a). If the old bank misses, it forwards the request to the memory controller (Fig. 4-7b). Demand moves have no races because all requests follow the same path through both virtual levels, i.e. they access the same sequence of cache banks. If a second request for the same line arrives at the new bank, it is queued at the MSHR that's currently being used to serve the first request.

Demand moves quickly migrate frequently used lines to their new locations, but the old banks must still be checked until all lines have moved. This adds latency to cache accesses. To limit this cost, banks walk the array in the background and incrementally invalidate lines whose location has changed. Unlike bulk invalidations, these *background invalidations* are not on the critical path and can proceed at a comparatively slow rate. For example, by scanning one set every 200 cycles, background invalidations finish in 100 Kcycles. Background invalidations begin after a short period (e.g., 50 Kcycles)

to allow frequently-accessed lines to migrate via demand moves. After banks walk the entire array, cores stop using the shadow VTB descriptors and resume normal operation.

In addition to background invalidations, this thesis also experimented with background moves, i.e. having banks send lines to their new locationths instead of invalidating them. However, it was found that background moves and background invalidations performed similarly—most of the benefit comes from not pausing cores as is done in bulk invalidations. CDCS prefers background invalidations because they are simpler: background moves require additional state at every bank (viz., *where* the line needs to be moved, not just that its location has changed), and require a more sophisticated protocol (as there can be races between demand and background moves).

Overall, by taking invalidations off the critical path, CDCS can reconfigure the cache frequently without pausing cores or invalidating frequently-accessed data. As Sec. 6.3 shows, background invalidations narrow the performance gap with an idealized architecture that moves lines immediately to their destination banks, and are faster than using bulk invalidations.

## 4.10 Overheads and applicability

### 4.10.1 Hardware overheads

Implementing CDCS as described imposes small overheads that the achieved system-wide performance and energy savings compensate for:

- Each bank is partitioned. With 512 KB banks and 64-byte lines, Vantage adds 8 KB of state per bank to support 64 bank partitions [4] (each tag needs a 6-bit partition id and each bank needs 256 bits of per-partition state).

- Each tile's VTB is 588 bytes: 576 bytes for 6 VC descriptors (3 normal + 3 shadow) and 12 bytes for the 3 tags.

- CDCS uses 4 monitors per bank. Each GMON has 1024 tags and 64 ways. Each tag is a 16-bit hash value (it does not store full addresses, since rare false positives are fine for monitoring purposes). Each way has a 16-bit limit register. This yields 2.1 KB monitors, and 8.4 KB overhead per tile.

41

This requires 17.1 KB of state per tile (2.9% of the space devoted to the tile's bank) and simple logic. Overheads are similar to prior partitioning-based schemes [4, 51].

Unlike Jigsaw [4], CDCS places each VC's monitor in a fixed location on the chip to avoid clearing or migrating monitor state. Since cores already hash lines to index the VTB, CDCS store the GMON location at the VTB. For full LLC coverage with $\gamma = 0.95$ and 64 cores, CDCS sample every $64^{th}$ access. Monitoring is off the critical path, so this has negligible impact on performance (traffic is roughly $1/64^{th}$ of an S-NUCA cache).

## 4.10.2 Software overheads

Periodically (every 25 ms in our implementation), a software runtime wakes up on core 0 and performs the steps in Fig. 4-1. Reconfiguration steps are at most quadratic on the number of tiles (Sec. 4.5), and most are simpler. Software overheads are small, 0.2% of system cycles, and are detailed in Sec. 6.3 and Table 6.1.

## 4.10.3 CDCS on other NUCA schemes

If partitioned banks are not desirable, CDCS can be used as-is with non-partitioned NUCA schemes [11, 27, 40]. To allow more VCs than threads, it is possible to use several smaller banks per tile (e.g., 4×128 KB), and size and place VCs at bank granularity. This would eliminate partitioning overheads, but would make VTBs larger and force coarser allocations. We evaluate the effects of such a configuration in Sec. 6.3. CDCS could also be used in spilling D-NUCAs [50], though the cost model (Sec. 4.2) would need to change to account for the multiple cache and directory lookups.

# Chapter 5

# Methodology

To evaluate proposed techniques, in this thesis, we use zsim, an open-source micro-architectural x86-64 simulator [53] for modeling. Zsim is a parallel simulator and adapts several novel techniques to be fast, scalable, and accurate. Zsim is suitable for this project because of its speed and scalability: simulation of thousands of cores at speeds of 100-1000s of MIPS.

## 5.1   Modeled System

We perform microarchitectural, execution-driven simulation using zsim [53], and model a 64-tile CMP connected by an $8 \times 8$ mesh NoC with 8 memory controllers at the edges. Each tile has one lean 2-way OOO core similar to Silvermont [29] and a 3-level cache hierarchy with parameters shown in Table 5.1. This system is similar to Knights Landing [26]. We use McPAT 1.1 [36] to derive the area and energy numbers of chip components (cores, caches, NoC, and memory controller) at 22 nm, and Micron DDR3L datasheets [42] for main memory. This system is implementable in $408\,\mathrm{mm}^2$ with typical power consumption of 80-130 W in our workloads, consistent with area and power of scaled Silvermont-based systems [26, 29].

## 5.2   Schemes

We compare CDCS with Jigsaw, R-NUCA, and S-NUCA organizations. R-NUCA and Jigsaw are implemented as proposed. R-NUCA uses 4-way rotational interleaving

| | |
|---|---|
| **Cores** | 64 cores, x86-64 ISA, 2 GHz, Silvermont-like OOO [29]: 8B-wide ifetch; 2-level bpred with 512×10-bit BHSRs + 1024×2-bit PHT, 2-way decode/issue/rename/commit, 32-entry IQ and ROB, 10-entry LQ, 16-entry SQ |
| **L1 caches** | 32 KB, 8-way set-associative, split D/I, 3-cycle latency |
| **L2 caches** | 128 KB private per-core, 8-way set-associative, inclusive, 6-cycle latency |
| **L3 cache** | 512 KB per tile, 4-way 52-candidate zcache, inclusive, 9 cycles, S-NUCA/R-NUCA/Jigsaw/CDCS |
| **Coherence protocol** | MESI, 64 B lines, in-cache directory, no silent drops; sequential consistency |
| **Global NoC** | 8×8 mesh, 128-bit flits and links, X-Y routing, 3-cycle pipelined routers, 1-cycle links |
| **Memory controllers** | 8 MCUs, 1 channel/MCU, 120 cycles zero-load latency, 12.8 GB/s per channel |

Table 5.1: Configuration of the simulated 64-core CMP.

and page-based reclassification [20]. CDCS and Jigsaw use 64-way, 1 Kline GMONs from Sec. 4.8, and reconfigure every 25 ms.

# 5.3 Workloads

We simulate mixes of single and multithreaded workloads, with a methodology similar to prior work [4, 51, 52]. We simulate single-threaded mixes of SPEC CPU2006 applications. We use the 16 SPEC CPU2006 applications with $\geq 5$ L2 MPKI: bzip2, gcc, bwaves, mcf, milc, zeusmp, cactusADM, leslie3d, calculix, GemsFDTD, libquantum, lbm, astar, omnet, sphinx3, and xalancbmk. We simulate mixes of 1–64 random applications. We fast-forward all applications for 20 billion instructions. We use a fixed-work methodology and equalize sample lengths to avoid sample imbalance, as in FIESTA [25]: We simulate each application alone, and measure how many instructions it executes in 1 billion cycles, $I_i$. Each experiment then runs the full mix until all applications execute at least $I_i$ instructions, and consider only the first $I_i$ instructions of each application to report performance. This ensures that each mix runs for at least 1 billion cycles.

We simulate multithreaded mixes of SPEC OMP2012 workloads. Since IPC is not a valid measure of work in multithreaded workloads [1], we instrument each application with heartbeats that report global progress (e.g., when each timestep or

44

transaction finishes). For each application, we find the smallest number of heartbeats that complete in over 1 billion cycles from the start of the parallel region when running alone. This is the region of interest (ROI). We then run the mixes by fast-forwarding all applications to the start of their parallel regions, and running the full mix until all applications complete their ROI.

We report weighted speedup over the S-NUCA baseline, which accounts for throughput and fairness [51, 54]. To achieve statistically significant results, we introduce small amounts of non-determinism as in [1], and perform enough runs to achieve 95% confidence intervals $\leq 1\%$.
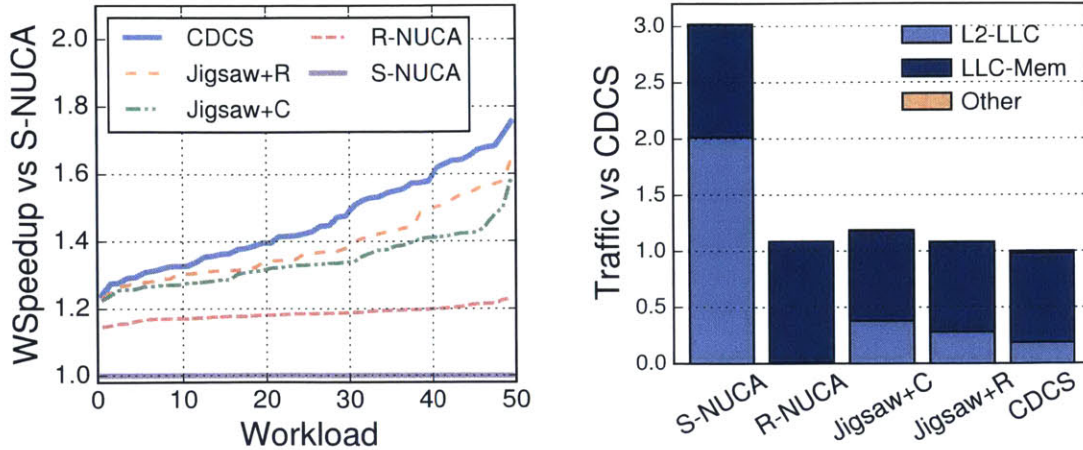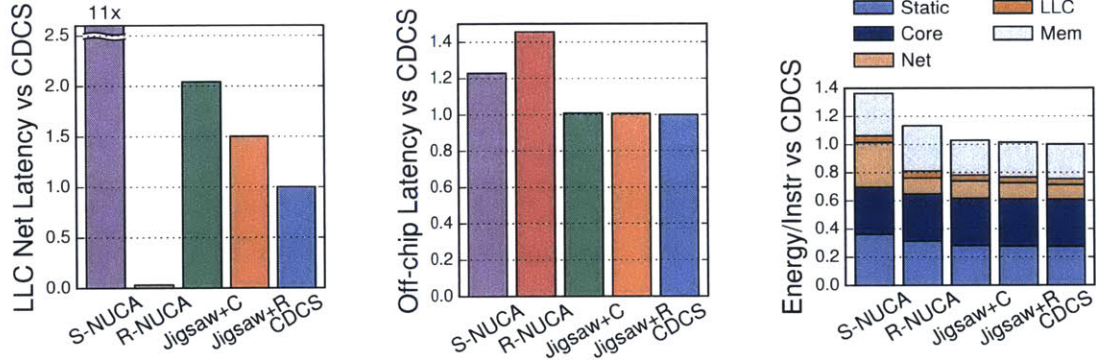
# Chapter 6

# Evaluation

## 6.1 Single-threaded mixes

Fig. 6-1a shows the distribution of weighted speedups that S-NUCA, R-NUCA, Jigsaw, and CDCS achieve in 50 mixes of 64 randomly-chosen, memory-intensive SPEC CPU2006 applications. We find that S-NUCA and R-NUCA are insensitive to thread placement (performance changes by $\leq 1\%$): S-NUCA because it spreads accesses among banks, and R-NUCA because its policies cause little contention. Therefore, we report results for both with a *random* scheduler, where threads are placed randomly at initialization, and stay pinned. We report Jigsaw results with two schedulers: *random* (Jigsaw+R) and *clustered* (Jigsaw+C). As we will see, neither choice is better in general—different mixes prefer one over the other. Each line shows the weighted speedup of a single scheme over the S-NUCA baseline, sorted along workload mixes ($x$-axis) by improvement (inverse CDF).

Fig. 6-1a shows that CDCS significantly improves system performance, achieving 46% gmean weighted speedup and up to 76%. Jigsaw+R achieves 38% gmean weighted speedup and up to 64%, Jigsaw+C achieves 34% gmean weighted speedup and up to 59%, and R-NUCA achieves 18% gmean weighted speedup and up to 23%. Jigsaw+C shows near-pathological behavior, as different instances of the same benchmark are placed close by, introducing capacity contention and hurting latency when they get large VCs. Jigsaw+R avoids this behavior and performs better, but CDCS avoids capacity contention much more effectively and attains higher speedups across all

47

(a) Distribution of weighted speedups over S-NUCA.

(b) Breakdown of avg. on-chip network traffic per instruction.



(c) Avg. on-chip latency on LLC accesses.

(d) Avg. off-chip latency (due to LLC misses).

(e) Breakdown of avg. energy per instruction.

Figure 6-1: Evaluation of S-NUCA, R-NUCA, Jigsaw, and CDCS across 50 mixes of 64 SPEC CPU2006 apps on a 64-core CMP.

mixes. CDCS and Jigsaw widely outperform R-NUCA, as R-NUCA does not manage capacity efficiently in heterogeneous workload mixes.

Fig. 6-1 gives more insight into these differences. Fig. 6-1c shows the average network latency incurred by LLC accesses across all mixes (Eq. 4.2), normalized to CDCS, while Fig. 6-1d compares off-chip latency (Eq. 4.1). S-NUCA incurs $11\times$ more on-chip network latency than CDCS on L2-LLC accesses, and 23% more off-chip latency. R-NUCA classifies most pages as private and maps them to the nearest bank, so its network latency for LLC accesses is negligible. However, the lack of capacity management degrades off-chip latency by 46% over CDCS. Jigsaw+C, Jigsaw+R and CDCS achieve similar off-chip latency, but Jigsaw+C and Jigsaw+R have $2\times$ and 51% higher on-chip network latency for LLC accesses than CDCS.

Fig. 6-1b compares the network traffic of different schemes, measured in flits, and split in L2-to-LLC and LLC-to-memory traffic. S-NUCA incurs 3× more traffic than CDCS, most of it due to LLC accesses. For other schemes, traffic due to LLC misses dominates, because requests are interleaved across memory controllers and take several hops. We could combine these schemes with NUMA-aware techniques [15, 16, 38, 56, 58, 59] to further reduce this traffic. Though not explicitly optimizing for it, CDCS achieves the lowest traffic.

Because CDCS improves performance and reduces network and memory traffic, it reduces energy as well. Fig. 6-1e shows the average energy per instruction of different organizations. Static energy (including chip and DRAM) decreases with higher performance, as each instruction takes fewer cycles. S-NUCA spends significant energy on network traversals, but other schemes make it a minor overhead; and R-NUCA is penalized by its more frequent memory accesses. Overall, Jigsaw+C, Jigsaw+R and CDCS reduce energy by 33%, 34% and 36% over S-NUCA, respectively.

CDCS benefits apps with large cache-fitting footprints, such as omnet, xalanc, and sphinx3, the most. They require multi-bank VCs to work well, and benefit from lower access latencies. Apps with smaller footprints benefit from the lower contention, but their speedups are moderate.

Fig. 6-2a shows how each of the proposed techniques in CDCS improves performance when applied to Jigsaw+R individually. We show results for latency-aware allocation (+L), thread placement (+T), and refined data placement (+D); +LTD is CDCS. Since cache capacity is scarce, latency-aware allocation helps little, whereas thread and data placement achieve significant, compounding benefits.

**Under-committed systems** Fig. 6-3 shows the weighted speedups achieved by S-NUCA, R-NUCA, Jigsaw+R, Jigsaw+C, and CDCS when the 64-core CMP is under-committed: each set of bars shows the gmean weighted speedup when running 50 mixes with an increasing number of single-threaded applications per mix, from 1 to 64. Besides characterizing these schemes on CMPs running at low utilization (e.g., due to limited power or parallelism), this scenario is similar to introducing a varying number of non-intensive benchmarks, for which LLC performance is a second-order effect.
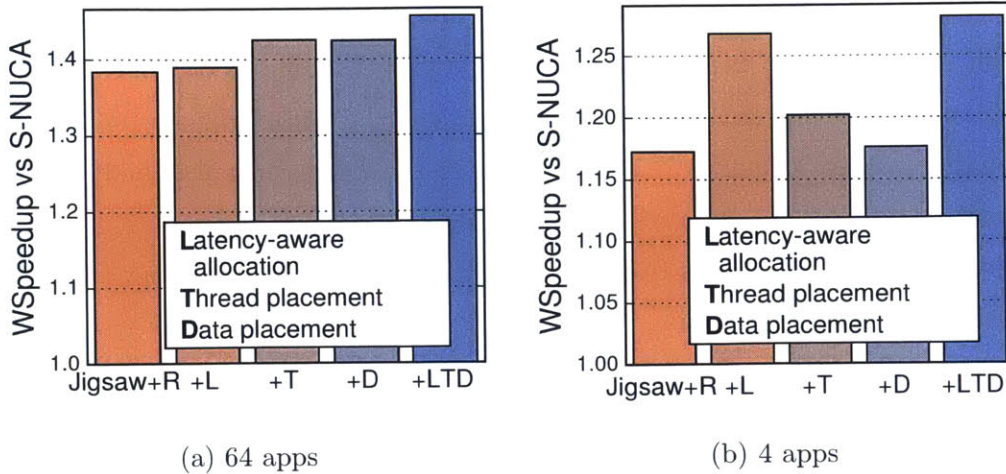
(a) 64 apps                (b) 4 apps

Figure 6-2: Factor analysis of CDCS with 50 mixes of 64 and 4 SPEC CPU2006 applications on a 64-core CMP.
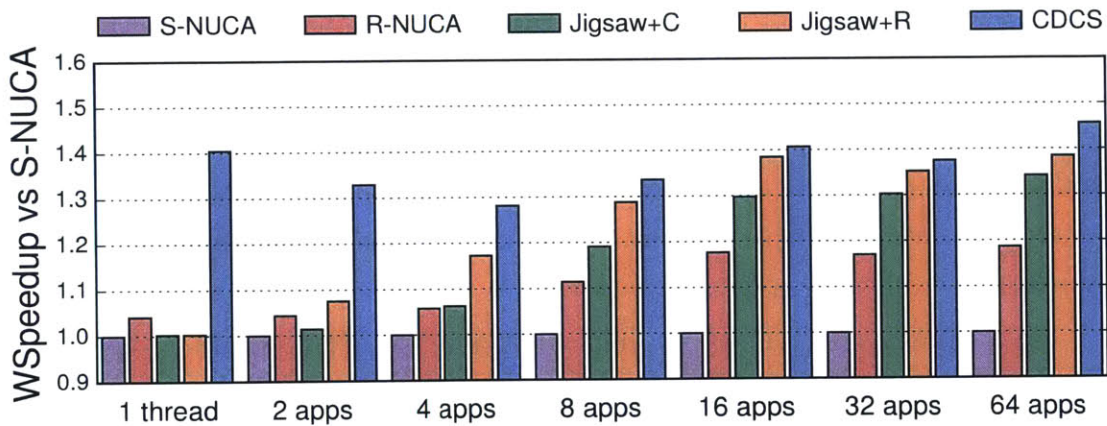


Figure 6-3: Weighted speedups for 50 mixes of 1, 2, 4, 6, 8, 16, 32, and 64 SPEC CPU2006 applications on a 64-core CMP.

Fig. 6-3 shows that CDCS maintains high weighted speedups throughout the whole range, while Jigsaw+R and Jigsaw+C work poorly on 1–8 app mixes. To see why, Fig. 6-4 shows the weighted speedup distribution and network traffic breakdown for the 4-app case. On-chip latency (L2-LLC) dominates Jigsaw's latency. In these mixes, cache capacity is plentiful, so large VC allocations hurt on-chip latency more than they help off-chip latency. CDCS's latency-aware allocation avoids using banks when detrimental, and yields most of the speedup in the 4-app mixes, as shown in Fig. 6-2b. At 4 apps, CDCS achieves 28% gmean weighted speedup, while Jigsaw+R sees 17% and Jigsaw+C sees 6%. Overall, latency-aware allocation becomes more important as capacity becomes plentiful.
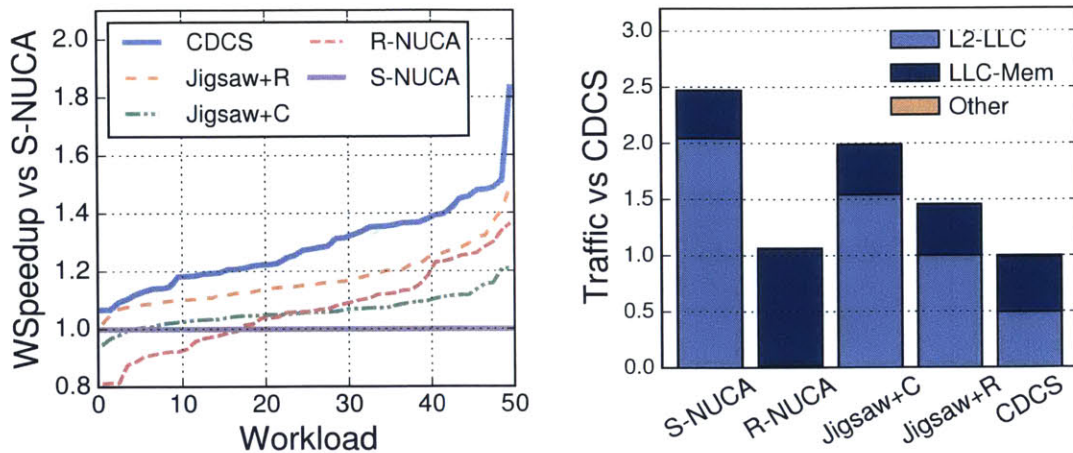
Figure 6-4: Weighted speedup distribution and traffic breakdown of 50 mixes of 4 SPEC CPU2006 apps on a 64-core CMP.



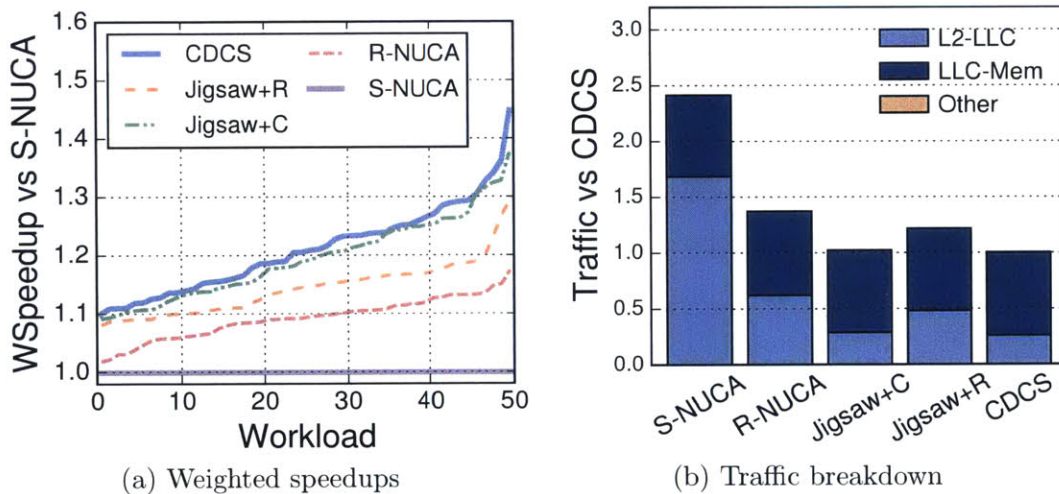(a) Weighted speedups       (b) Traffic breakdown

Figure 6-5: Weighted speedup distribution and traffic breakdown of 50 mixes of eight 8-thread SPEC OMP2012 apps on a 64-core CMP.

## 6.2 Multithreaded mixes

Fig. 6-5a shows the distribution of weighted speedups for 50 mixes of eight 8-thread SPEC OMP2012 applications (64 threads total) running on the 64-core CMP. CDCS achieves gmean weighted speedup of 21%. Jigsaw+R achieves 14%, Jigsaw+C achieves 19%, and R-NUCA achieves 9%. Trends are reversed: on multi-threaded benchmarks, Jigsaw works better with clustered thread placement than with random (S-NUCA and R-NUCA are still insensitive). CDCS sees smaller benefits over Jigsaw+C. Fig. 6-5b shows that they get about the same network traffic, while others are noticeably worse.

Fig. 6-6a shows the distribution of weighted speedups with under-committed
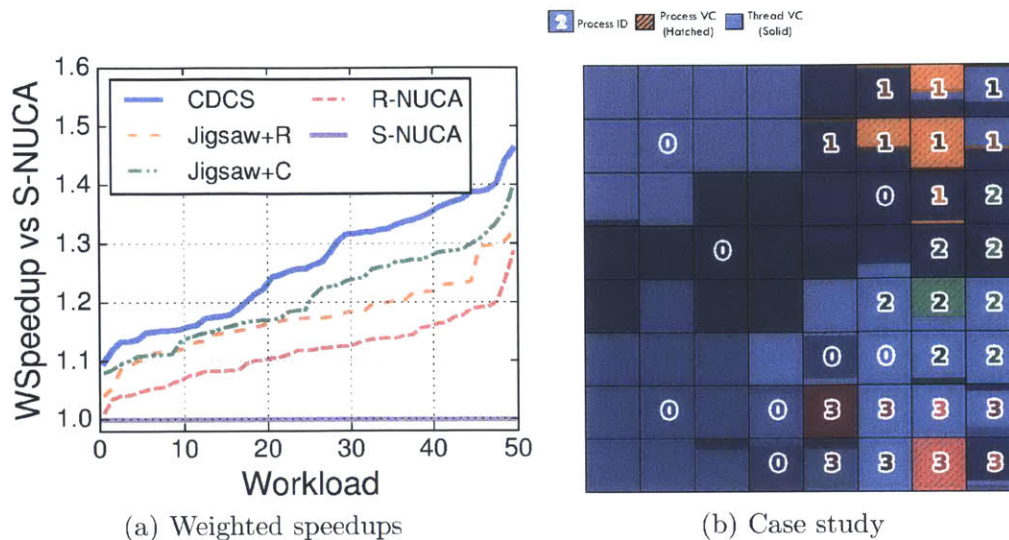
(a) Weighted speedups  (b) Case study

Figure 6-6: Weighted speedups for 50 mixes of four 8-thread SPEC OMP2012 apps (32 threads total) on a 64-core CMP, and case study with private-heavy and shared-heavy apps.

system running mixes of four 8-thread applications. CDCS increases its advantage over Jigsaw+C, as it has more freedom to place threads. CDCS dynamically clusters or spreads each process as the context demands: shared-heavy processes are clustered, and private-heavy processes are spread out. Fig. 6-6b illustrates this behavior by showing the thread and data placement of a specific mix, where one of the apps, mgrid (process P0), is private and intensive, and the others, md (P1), ilbdc (P2), and nab (P3) access mostly shared data. CDCS gives most capacity to mgrid, spreads its threads over the CMP, and tightly clusters P1–3 around their shared data.

From the results of Sec. 6.1 and Sec. 6.2, we can see that Jigsaw+R and Jigsaw+C help different types of programs, but no option is best in general. Yet by jointly placing threads and data, CDCS always provides the highest performance across all mixes. Thus, beyond improving performance, CDCS provides an important advantage in *guarding against pathological behavior incurred by fixed policies*.

## 6.3 CDCS analysis

**Reconfiguration overheads**  Table 6.1 shows the CPU cycles spent, on average, in each of the steps of the reconfiguration procedure. Overheads are negligible: each reconfiguration consumes a mere 0.2% of system cycles. Sparse GMON curves improve

| Threads / Cores | 16 / 16 | 16 / 64 | 64 / 64 |
|---|---|---|---|
| Capacity allocation (Mcycles) | 0.30 | 0.30 | 1.20 |
| Thread placement (Mcycles) | 0.29 | 0.80 | 3.44 |
| Data placement (Mcycles) | 0.13 | 0.36 | 1.85 |
| Total runtime (Mcycles) | 0.72 | 1.46 | 6.49 |
| Overhead @ 25 ms (%) | 0.09 | 0.05 | 0.20 |

Table 6.1: CDCS runtime analysis. Avg Mcycles per invocation of each reconfiguration step, total runtime, and relative overhead.

Peekahead's runtime, taking 1.2 Mcycles at 64 cores instead of the 7.6 Mcycles it would require with 512-way UMONs [4]. Although thread and data placement have quadratic runtime, they are practical even at thousands of cores (1.2% projected overhead at 1024 cores).

**Alternative thread and data placement schemes** We have considered more computationally expensive alternatives for thread and data placement. First, we explored using integer linear programming (ILP) to produce the best achievable data placement. We formulate the ILP problem by minimizing Eq. 4.2 subject to the bank capacity and VC allocation constraints, and solve it in Gurobi [19]. ILP data placement improves weighted speedup by 0.5% over CDCS on 64-app mixes. However, Gurobi takes about 219 Mcycles to solve 64-cores, far too long to be practical. We also formulated the joint thread and data placement ILP problem, but Gurobi takes at best tens of minutes to find the solution and frequently does not converge.

Since using ILP for thread placement is infeasible, we have implemented a simulated annealing [60] thread placer, which tries 5000 rounds of thread swaps to find a high-quality solution. This thread placer is only 0.6% better than CDCS on 64-app runs, and is too costly (6.3 billion cycles per run).

We also explored using METIS [30], a graph partitioning tool, to jointly place threads and data. We were unable to outperform CDCS. We observe that graph partitioning methods recursively divide threads and data into equal-sized partitions of the chip, splitting around the center of the chip first. CDCS, by contrast, often clusters one application around the center of the chip to minimize latency. In trace-driven runs, graph partitioning increases network latency by 2.5% over CDCS.
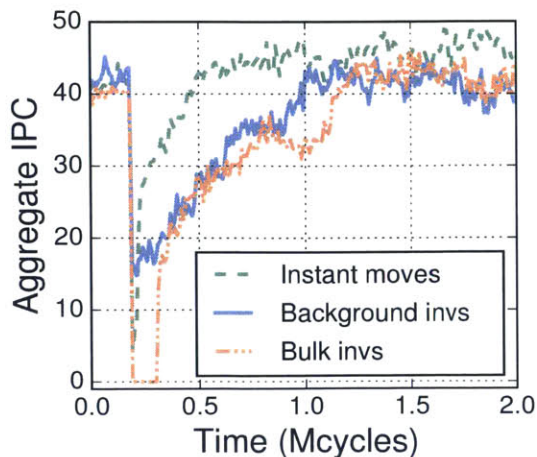
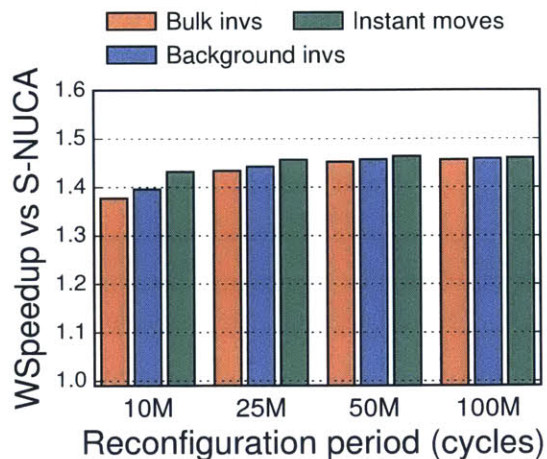Figure 6-7: IPC throughput of a 64-core CMP with various data movement schemes during one reconfiguration.

Figure 6-8: Weighted speedup of 64-app mixes for various data movement schemes vs. reconfiguration period.

**Geometric monitors** 1K-line, 64-way GMONs match the performance of 256-way UMONs. UMONs lose performance below 256 ways because of their poor resolution: 64-way UMONs degrade performance by 3% on 64-app mixes. In contrast, unrealistically large 16K-line, 1K-way UMONs are only 1.1% better than 64-way GMONs.

**Reconfiguration schemes** We evaluate several LLC reconfiguration schemes: demand moves plus background invalidations (as in CDCS), bulk invalidations (as in Jigsaw), and idealized, instant moves. The main benefit of demand moves is avoiding global pauses, which take 114 Kcycles on average, and up to 230 Kcycles. While this is a 0.23% overhead if reconfigurations are performed every 50 Mcycles (25 ms), many applications cannot tolerate such pauses [17, 45]. Fig. 6-7 shows a trace of aggregate IPC across all 64 cores during one representative reconfiguration. This trace focuses on a small time interval to show how performance changes right after a reconfiguration, which happens at 200 Kcycles. By serving lines with demand moves, CDCS prevents pauses and achieves smooth reconfigurations, while bulk invalidations pause the chip for 100 Kcycles in this case. Besides pauses, bulk invalidations add misses and hurt performance. With 64 apps (Fig. 6-1), misses are already frequent and per-thread capacity is scarce, so the average slowdown is 0.5%. With 4 apps (Fig. 6-4), VC allocations are larger and threads take longer to warm up the LLC, so the slowdown

54

is 1.4%. Note that since SPEC CPU2006 is stable for long phases, these results may underestimate overheads for apps with more time-varying behavior. Fig. 6-8 compares the weighted speedups of different schemes when reconfiguration intervals increase from 10 Mcycles to 100 Mcycles. CDCS outperforms bulk invalidations, though differences diminish as reconfiguration interval increases.

**Bank-partitioned NUCA**  CDCS can be used without fine-grained partitioning (Sec. 4.10). With the parameters in Table 5.1 but 4 smaller banks per tile, CDCS achieves 36% gmean weighted speedup (up to 49%) over S-NUCA in 64-app mixes, vs. 46% gmean with partitioned banks. This difference is mainly due to coarser-grain capacity allocations, as CDCS allocates full banks in this case.

# Chapter 7

# Conclusion

Data movement limits the performance and energy efficiency of multicore chips. This thesis has presented computation and data co-scheduling (CDCS), a joint hardware and software technique to reduce data movement. CDCS hardware provides efficient monitoring and reconfiguration techniques, and CDCS software optimizes the full system with simple performance model. Results shown that the proposed technique successfully reduces data movement and thus improves performance and energy efficiency. In ongoing and future work, we plan to extend CDCS to manage the whole memory hierarchy and handle heterogeneous memory resources.

# Bibliography

[1] A. Alameldeen and D. Wood, "IPC considered harmful for multiprocessor work-loads," *IEEE Micro*, vol. 26, no. 4, 2006.

[2] B. Beckmann, M. Marty, and D. Wood, "ASR: Adaptive selective replication for CMP caches," in *Proc. MICRO-39*, 2006.

[3] B. Beckmann and D. Wood, "Managing wire delay in large chip-multiprocessor caches," in *Proc. MICRO-37*, 2004.

[4] N. Beckmann and D. Sanchez, "Jigsaw: Scalable Software-Defined Caches," in *Proc. PACT-22*, 2013.

[5] ——, "Talus: A Simple Way to Remove Cliffs in Cache Performance," in *Proc. of the 21th intl. symp. on High Performance Computer Architecture*, 2015.

[6] S. Bell, B. Edwards, J. Amann, R. Conlin, K. Joyce, V. Leung, J. MacKay, M. Reif, L. Bao, J. Brown, M. Mattina, C.-C. Miao, C. Ramey, D. Wentzlaff, W. Anderson, E. Berger, N. Fairbanks, D. Khan, F. Montenegro, J. Stickney, and J. Zook, "TILE64 processor: A 64-core SoC with mesh interconnect," in *Proc. ISSCC*, 2008.

[7] S. Blagodurov, S. Zhuravlev, M. Dashti, and A. Fedorova, "A case for NUMA-aware contention management on multicore systems," in *Proc. USENIX ATC*, 2011.

[8] J. Chang and G. Sohi, "Cooperative caching for chip multiprocessors," in *Proc. ISCA-33*, 2006.

[9] D. Chiou, P. Jain, L. Rudolph, and S. Devadas, "Application-specific memory management for embedded systems using software-controlled caches," in *Proc. DAC-37*, 2000.

[10] Z. Chishti, M. Powell, and T. Vijaykumar, "Optimizing replication, communication, and capacity allocation in cmps," in *ISCA-32*, 2005.

[11] S. Cho and L. Jin, "Managing distributed, shared L2 caches through OS-level page allocation," in *Proc. MICRO-39*, 2006.

[12] H. Cook, M. Moreto, S. Bird, K. Dao, D. A. Patterson, and K. Asanovic, "A hardware evaluation of cache partitioning to improve utilization and energy-efficiency while preserving responsiveness," in *ISCA-40*, 2013.

[13] J. Corbet, "AutoNUMA: the other approach to NUMA scheduling," in *LWN*, 2012. [Online]. Available: https://lwn.net/Articles/488709/

[14] W. J. Dally, "GPU Computing: To Exascale and Beyond," in *Supercomputing '10, Plenary Talk*, 2010.

[15] R. Das, R. Ausavarungnirun, O. Mutlu, A. Kumar, and M. Azimi, "Application-to-core mapping policies to reduce memory system interference in multi-core systems," in *Proc. HPCA-19*, 2013.

[16] M. Dashti, A. Fedorova, J. Funston, F. Gaud, R. Lachaize, B. Lepers, V. Quema, and M. Roth, "Traffic management: a holistic approach to memory placement on NUMA systems," in *Proc. ASPLOS-18*, 2013.

[17] J. Dean and L. Barroso, "The Tail at Scale," *CACM*, vol. 56, 2013.

[18] F. Guo, Y. Solihin, L. Zhao, and R. Iyer, "A framework for providing quality of service in chip multi-processors," in *Proc. MICRO-40*, 2007.

[19] Gurobi, "Gurobi optimizer reference manual version 5.6," 2013. [Online]. Available: http://www.gurobi.com

[20] N. Hardavellas, M. Ferdman, B. Falsafi, and A. Ailamaki, "Reactive NUCA: near-optimal block placement and replication in distributed caches," in *Proc. ISCA-36*, 2009.

[21] ——, "Toward dark silicon in servers," *IEEE Micro*, vol. 31, no. 4, 2011.

[22] N. Hardavellas, I. Pandis, R. Johnson, and N. Mancheril, "Database Servers on Chip Multiprocessors: Limitations and Opportunities," in *Proc. CIDR*, 2007.

[23] J. L. Hennessy and D. A. Patterson, *Computer Architecture: A Quantitative Approach (5th ed.)*. Morgan Kaufmann, 2011.

[24] H. J. Herrmann, "Geometrical cluster growth models and kinetic gelation," *Physics Reports*, vol. 136, no. 3, pp. 153–224, 1986.

[25] A. Hilton, N. Eswaran, and A. Roth, "FIESTA: A sample-balanced multi-program workload methodology," in *Proc. MoBS*, 2009.

[26] Intel, "Knights Landing: Next Generation Intel Xeon Phi," in *SC Presentation*, 2013.

[27] J. Jaehyuk Huh, C. Changkyu Kim, H. Shafi, L. Lixin Zhang, D. Burger, and S. Keckler, "A NUCA substrate for flexible CMP cache sharing," *IEEE Trans. Par. Dist. Sys.*, vol. 18, no. 8, 2007.

[28] A. Jaleel, H. H. Najaf-Abadi, S. Subramaniam, S. C. Steely, and J. Emer, "CRUISE: Cache replacement and utility-aware scheduling," in *Proc. ASPLOS*, 2012.

[29] D. Kanter, "Silvermont, Intels Low Power Architecture," 2013. [Online]. Available: http://www.realworldtech.com/silvermont/

[30] G. Karypis and V. Kumar, "A fast and high quality multilevel scheme for partitioning irregular graphs," *SIAM Journal on scientific Computing*, vol. 20, no. 1, 1998.

[31] H. Kasture and D. Sanchez, "Ubik: Efficient Cache Sharing with Strict QoS for Latency-Critical Workloads," in *Proc. ASPLOS-19*, 2014.

[32] R. E. Kessler, M. D. Hill, and D. A. Wood, "A comparison of trace-sampling techniques for multi-megabyte caches," *Computers, IEEE Transactions on*, vol. 43, no. 6, 1994.

[33] C. Kim, D. Burger, and S. Keckler, "An adaptive, non-uniform cache structure for wire-delay dominated on-chip caches," in *ASPLOS*, 2002.

[34] N. Kurd, S. Bhamidipati, C. Mozak, J. Miller, T. Wilson, M. Nemani, and M. Chowdhury, "Westmere: A family of 32nm IA processors," in *Proc. ISSCC*, 2010.

[35] H. Lee, S. Cho, and B. R. Childers, "CloudCache: Expanding and shrinking private caches," in *Proc. HPCA-17*, 2011.

[36] S. Li, J. H. Ahn, R. Strong, J. Brockman, D. Tullsen, and N. Jouppi, "McPAT: an integrated power, area, and timing modeling framework for multicore and manycore architectures," in *MICRO-42*, 2009.

[37] J. Lin, Q. Lu, X. Ding, Z. Zhang, X. Zhang, and P. Sadayappan, "Gaining insights into multicore cache partitioning: Bridging the gap between simulation and real systems," in *Proc. HPCA-14*, 2008.

[38] Z. Majo and T. R. Gross, "Memory system performance in a NUMA multicore multiprocessor," in *Proc. ISMM*, 2011.

[39] R. Manikantan, K. Rajan, and R. Govindarajan, "Probabilistic shared cache management (PriSM)," in *Proc. ISCA-39*, 2012.

[40] M. Marty and M. Hill, "Virtual hierarchies to support server consolidation," in *Proc. ISCA-34*, 2007.

[41] J. Merino, V. Puente, and J. Gregorio, "ESP-NUCA: A low-cost adaptive non-uniform cache architecture," in *Proc. HPCA-16*, 2010.

[42] Micron, "1.35V DDR3L power calculator (4Gb x16 chips)," 2013.

[43] M. Moreto, F. J. Cazorla, A. Ramirez, R. Sakellariou, and M. Valero, "FlexDCP: A QoS framework for CMP architectures," *ACM SIGOPS Operating Systems Review*, vol. 43, no. 2, 2009.

[44] T. Nowatzki, M. Sartin-Tarm, L. D. Carli, K. Sankaralingam, C. Estan, and B. Robatmili, "A general constraint-centric scheduling framework for spatial architectures," in *Proc. PLDI-34*, 2013.

[45] J. Ousterhout, P. Agrawal, D. Erickson, C. Kozyrakis, J. Leverich, D. Mazières, S. Mitra, A. Narayanan, G. Parulkar, M. Rosenblum, S. M. Rumble, E. Stratmann, and R. Stutsman, "The case for RAMClouds: scalable high-performance storage entirely in DRAM," *ACM SIGOPS Operating Systems Review*, vol. 43, no. 4, 2010.

[46] D. Page, "Partitioned cache architecture as a side-channel defence mechanism," *IACR Cryptology ePrint archive*, no. 2005/280, 2005.

[47] P. N. Parakh, R. B. Brown, and K. A. Sakallah, "Congestion driven quadratic placement," in *Proc. DAC-35*, 1998.

[48] J. Park and W. J. Dally, "Buffer-space efficient and deadlock-free scheduling of stream applications on multi-core architectures," in *Proc. SPAA-22*, 2010.

[49] F. Pellegrini and J. Roman, "SCOTCH: A software package for static mapping by dual recursive bipartitioning of process and architecture graphs," in *Proc. HPCN*, 1996.

[50] M. Qureshi, "Adaptive Spill-Receive for Robust High-Performance Caching in CMPs," in *Proc. HPCA-10*, 2009.

[51] M. Qureshi and Y. Patt, "Utility-based cache partitioning: A low-overhead, high-performance, runtime mechanism to partition shared caches," in *Proc. MICRO-39*, 2006.

[52] D. Sanchez and C. Kozyrakis, "Vantage: Scalable and Efficient Fine-Grain Cache Partitioning," in *Proc. ISCA-38*, 2011.

[53] ——, "ZSim: Fast and Accurate Microarchitectural Simulation of Thousand-Core Systems," in *ISCA-40*, 2013.

[54] A. Snavely and D. M. Tullsen, "Symbiotic jobscheduling for a simultaneous multithreading processor," in *Proc. ASPLOS-8*, 2000.

[55] D. Tam, R. Azimi, L. Soares, and M. Stumm, "Managing shared l2 caches on multicore systems in software," in *WIOSCA*, 2007.

[56] D. Tam, R. Azimi, and M. Stumm, "Thread clustering: sharing-aware scheduling on smp-cmp-smt multiprocessors," in *Proc. Eurosys*, 2007.

[57] S. Thoziyoor, N. Muralimanohar, J. H. Ahn, and N. P. Jouppi, "CACTI 5.1," HP Labs, Tech. Rep. HPL-2008-20, 2008.

[58] A. Tumanov, J. Wise, O. Mutlu, and G. R. Ganger, "Asymmetry-aware execution placement on manycore chips," in *SFMA-3 workshop*, 2013.

[59] B. Verghese, S. Devine, A. Gupta, and M. Rosenblum, "Operating system support for improving data locality on CC-NUMA compute servers," in *Proc. ASPLOS*, 1996.

[60] D. Wong, H. W. Leong, and C. L. Liu, *Simulated annealing for VLSI design.* Kluwer Academic Publishers, 1988.

[61] C. Wu and M. Martonosi, "A Comparison of Capacity Management Schemes for Shared CMP Caches," in *WDDD-7*, 2008.

[62] M. Zhang and K. Asanovic, "Victim replication: Maximizing capacity while hiding wire delay in tiled chip multiprocessors," in *ISCA*, 2005.

[63] S. Zhuravlev, S. Blagodurov, and A. Fedorova, "Addressing shared resource contention in multicore processors via scheduling," in *Proc. ASPLOS*, 2010.