# SCALING DISTRIBUTED CACHE HIERARCHIES THROUGH COMPUTATION AND DATA CO-SCHEDULING
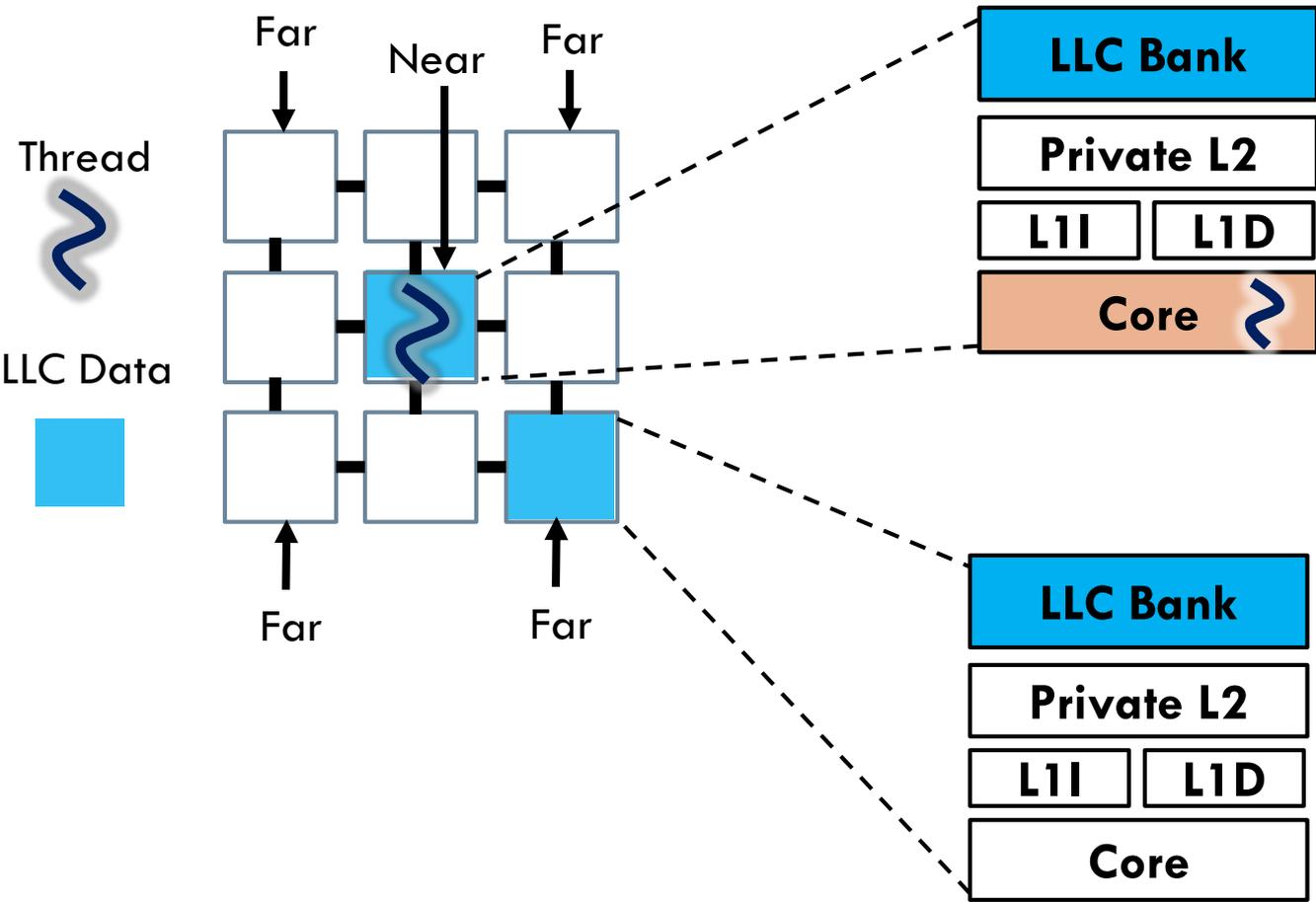
NATHAN BECKMANN, **PO-AN TSAI** AND DANIEL SANCHEZ

MIT CSAIL

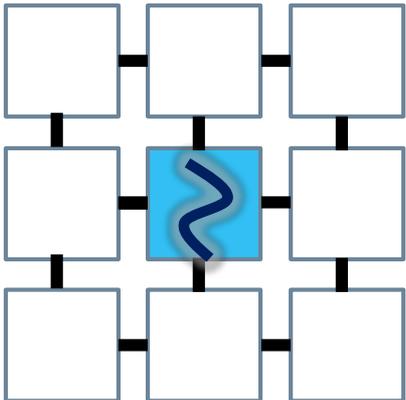**Massachusetts Institute of Technology**

CSAIL

# Executive Summary

Far    Near    Far

Thread

LLC Data

Far      Far

| LLC Bank |
| --- |
| Private L2 |
| L1I    L1D |
| Core |

| LLC Bank |
| --- |
| Private L2 |
| L1I    L1D |
| Core |

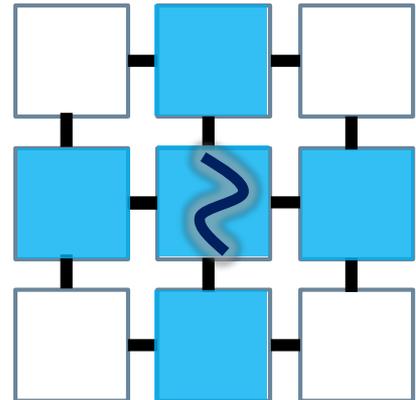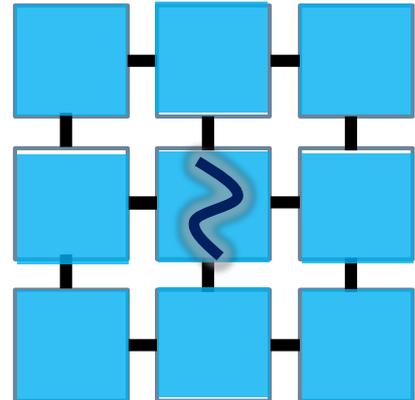# Executive Summary

Thread

LLC Data

Many misses
Low hit latency
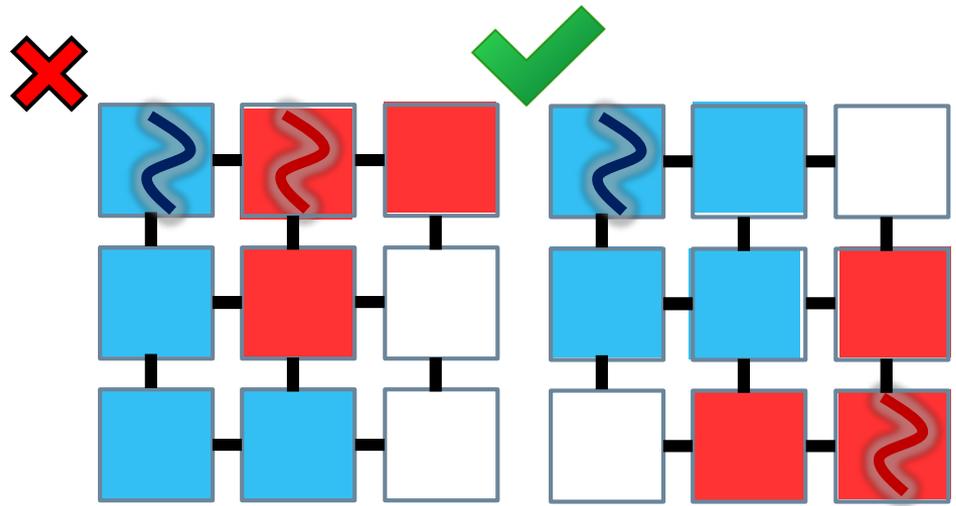
Moderate misses
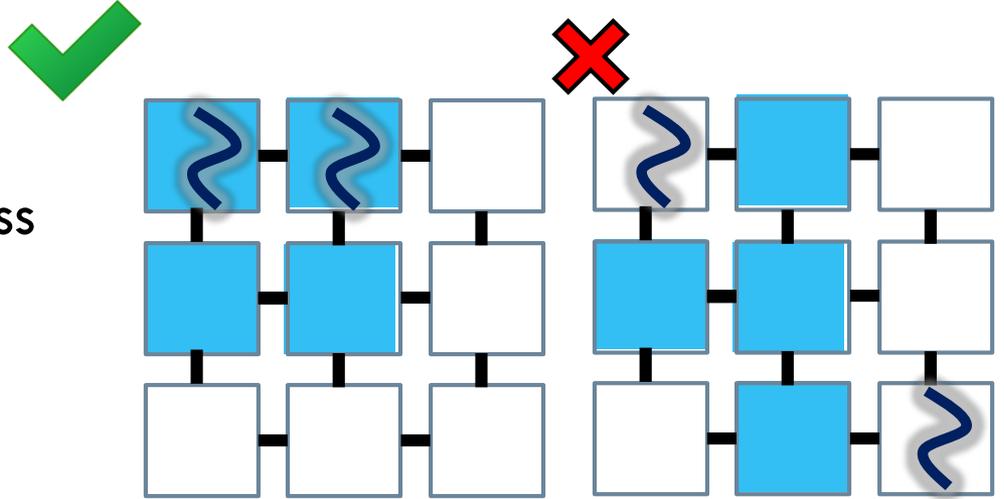Medium hit latency

Few misses
High hit latency

More capacity does not always mean better performance

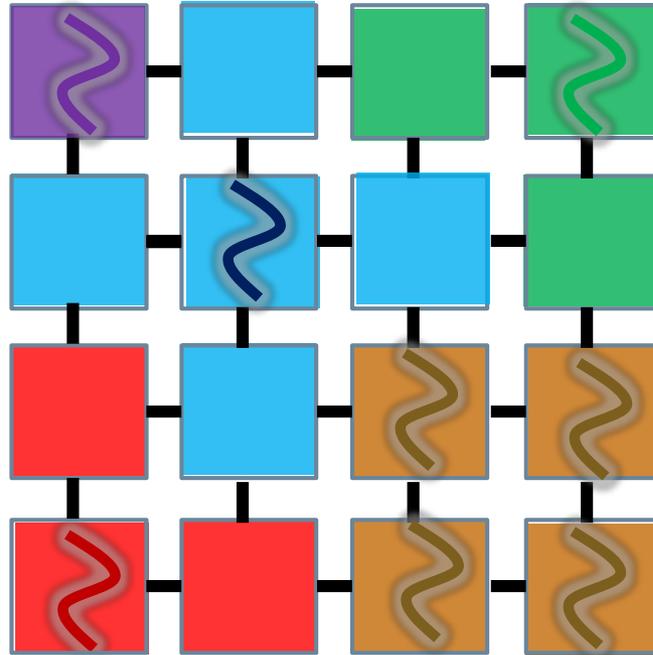# Executive Summary

Threads access
different data

Threads access
same data

# Executive Summary

- CDCS jointly places threads and data to reduce data movement



  - Improves performance by 46% on average and by up to 76%
  - Saves 36% of system energy
  - Uses low-overhead algorithms that perform within 1% of impractical, idealized solutions
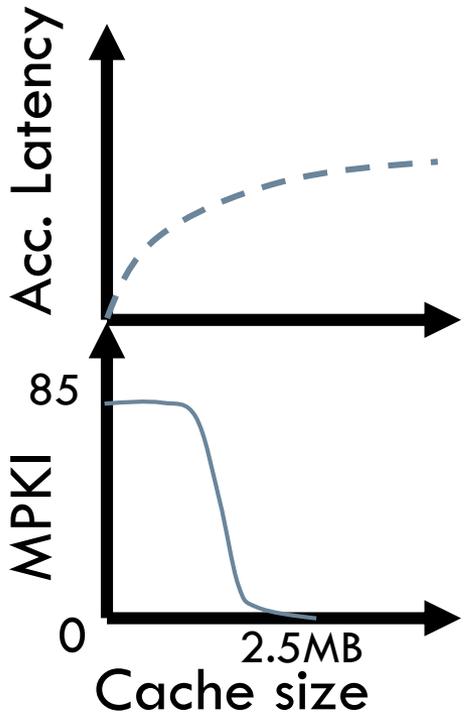
# Agenda

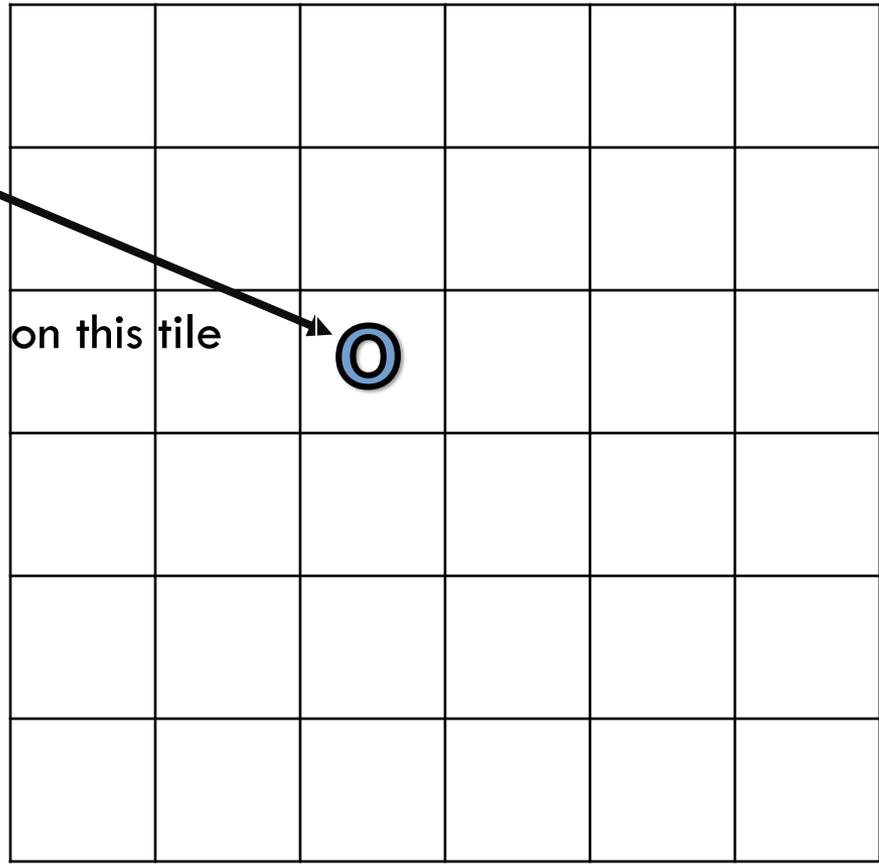- Background


- CDCS Design


- Evaluation

# Capacity vs latency

6x6 mesh, 18MB NUCA
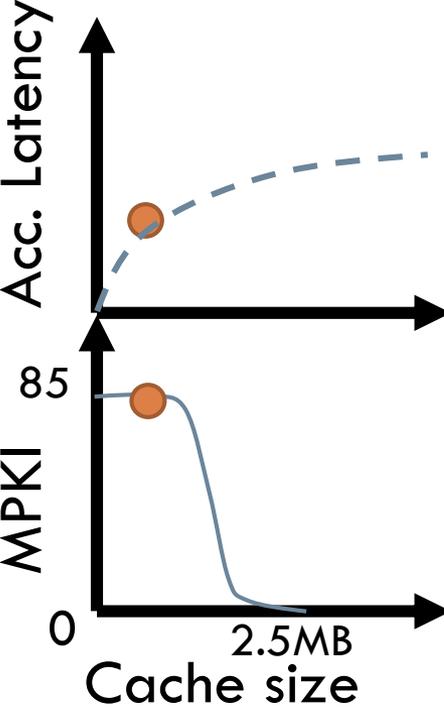
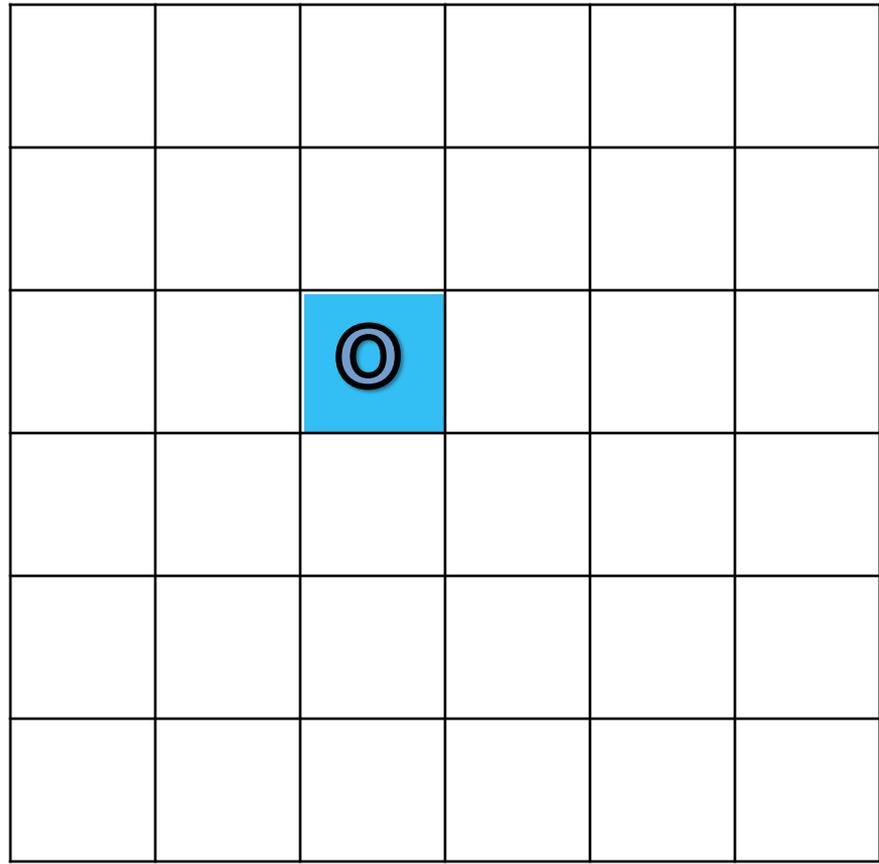App: 471.omnetpp
from SPECCPU2006

Thread running on this tile

**Acc. Latency**

**MPKI**

85

0

2.5MB

**Cache size**

# Capacity vs latency

App: 471.omnetpp
from SPECCPU2006

Place data in local bank

Acc. Latency

MPKI

85

0

2.5MB

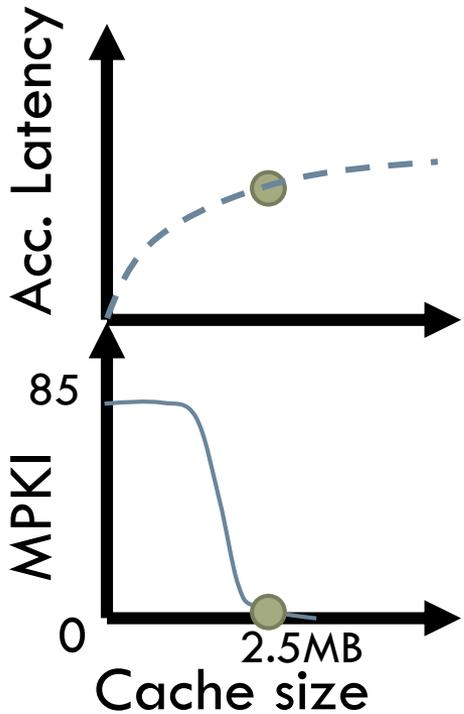Cache size

# Capacity vs latency

Use closest banks that just fit working set

App: 471.omnetpp
from SPECCPU2006

3.7x speedup

| Banks | Perf |
|-------|------|
| 1 | 1x |
| | |
| | |

Acc. Latency

MPKI

85

0

2.5MB

Cache size

# Capacity vs latency

Place data across the chip

App: 471.omnetpp
from SPECCPU2006



2.4x speedup

| Banks | Perf |
|-------|------|
| 1 | 1x |
| 5 | 3.7x |

Acc. Latency

MPKI

85

0

2.5MB

Cache size

In NUCA, using more capacity than needed is detrimental!

# Thread placement matters

App: 471.omnetpp
471.omnetpp
471.omnetpp
471.omnetpp
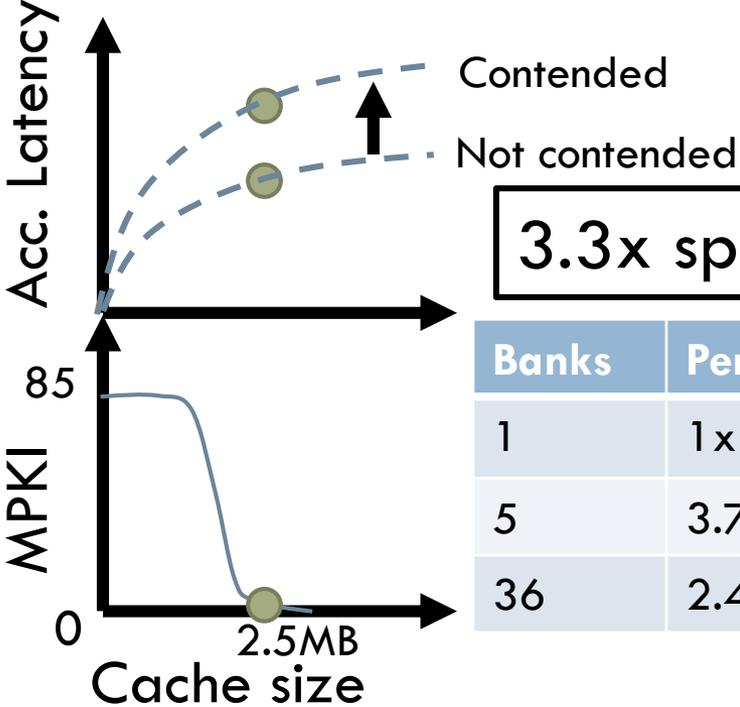
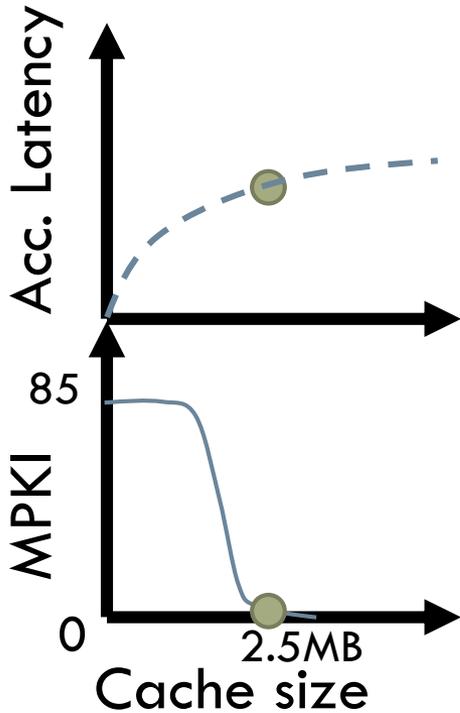**Capacity contention** changes achievable access latency!

Contended

Not contended

3.3x speedup each

| Banks | Perf |
|-------|------|
| 1 | 1x |
| 5 | 3.7x |
| 36 | 2.4x |

Acc. Latency

MPKI

85

0

2.5MB

Cache size

# Thread placement matters

App: 471.omnetpp
471.omnetpp
471.omnetpp
471.omnetpp

Spread out threads



3.7x speedup

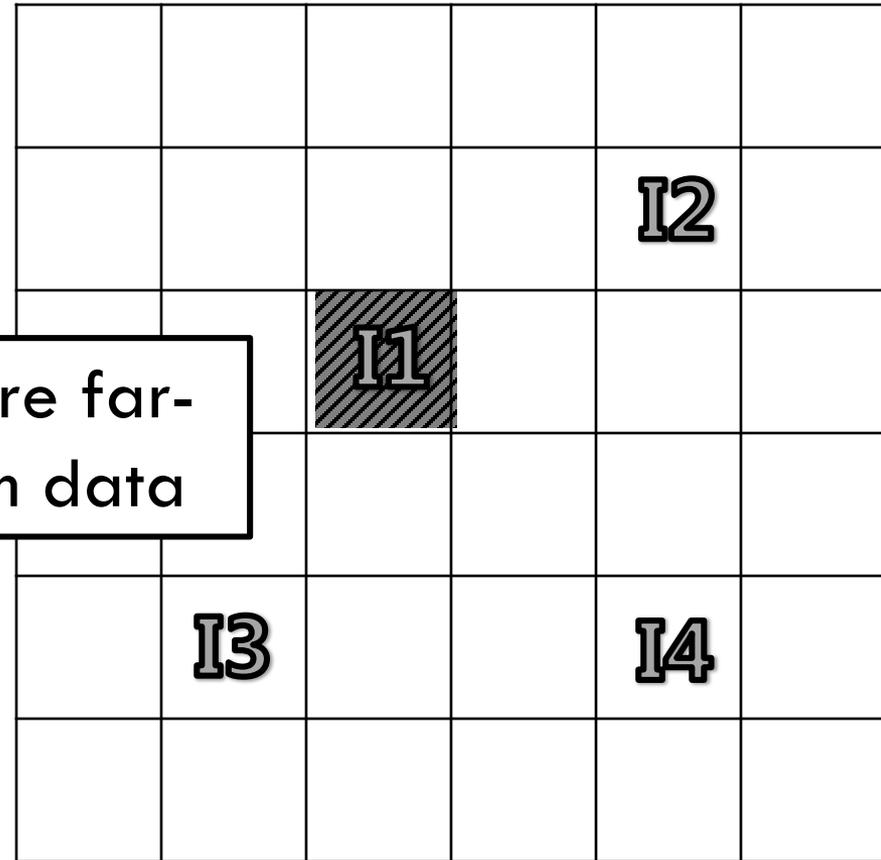| Banks | Perf |
|-------|------|
| 1 | 1x |
| 5 | 3.7x |
| 36 | 2.4x |

85

0    2.5MB
Cache size

# Thread placement matters

App: 4-thread `360.ilbdc` from SPECOMP2012

Spread out threads

Acc. Latency

MPKI

85

0

0.5MB

Cache size

Threads are far-away from data

I1

I2

I3

I4

# Thread placement matters

App: 4-thread 360.ilbdc
from SPECOMP2012

Cluster threads

Acc. Latency

Spread

Clustered

1.4x speedup

I3  I1  I2

I4

85

MPKI

One scheduler does not fit
all applications!

0

0.5MB

Cache size

# Dynamic NUCA

Mix:

4 x 471.omnetpp

4-thread 360.ilbdc

D-NUCA

✔ Place data

✘ Control capacity

✘ Place threads

R-NUCA [Hardavellas'09]

# Partitioned NUCA

Mix:

4 x 471.omnetpp

4-thread 360.ilbdc

Jigsaw [Beckmann'13]

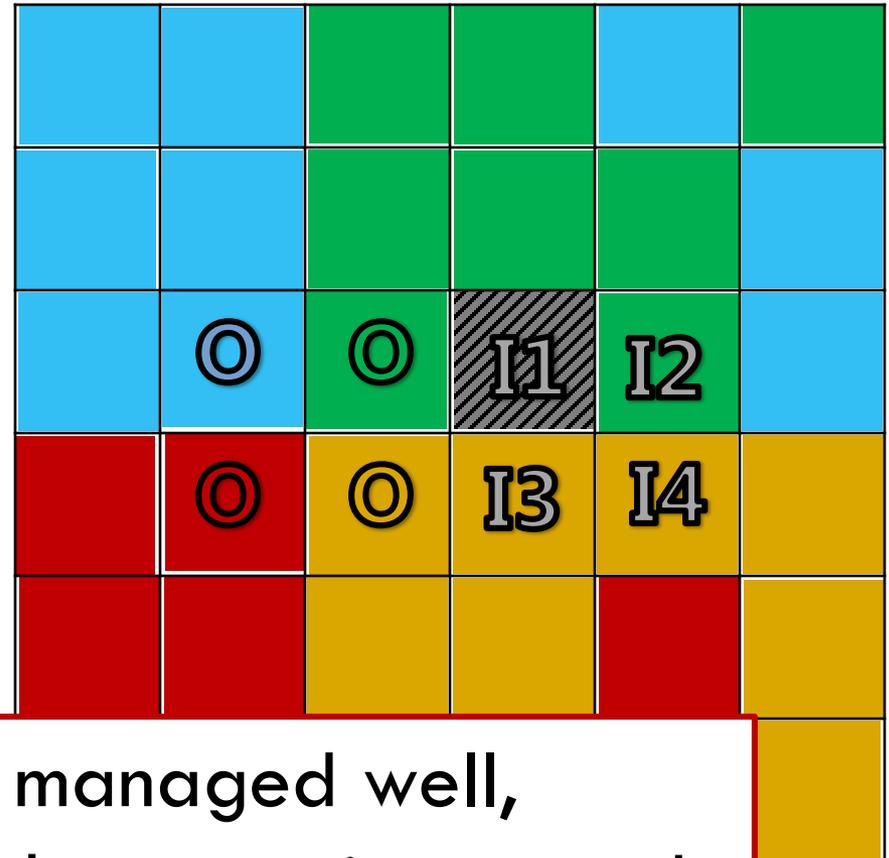|  | D-NUCA | Partitioned NUCA |
| --- | --- | --- |
| Place data | ✔ | ✔ |
| Control capacity | ✘ | ✔ |
|  | ✘ | ✘ |

When capacity is managed well, thread placement becomes important!
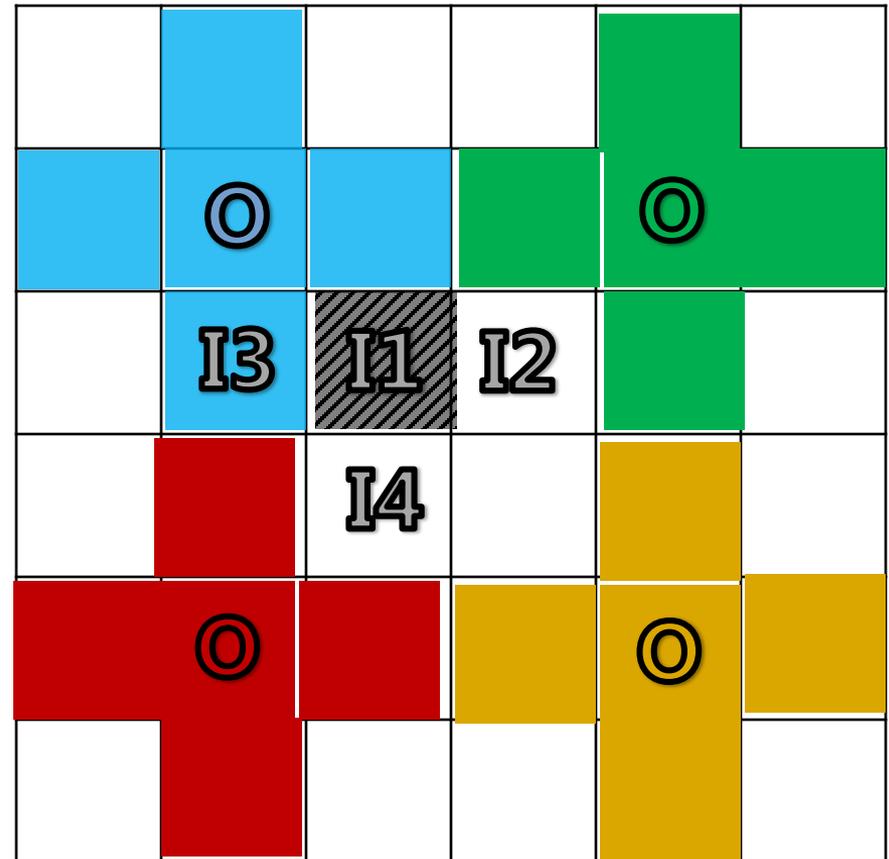
# CDCS

Mix:

4 x 471.omnetpp

4-thread 360.ilbdc

| D-NUCA | Partitioned NUCA | CDCS | |
|--------|------------------|------|--|
| ✔ | ✔ | ✔ | Place data |
| ✘ | ✔ | ✔ | Control capacity |
| ✘ | ✘ | ✔ | Place threads |

**CDCS**

# Agenda

- Background

- CDCS Design
  - Operation
  - Optimization

- Evaluation

# CDCS Overview



Hardware ——————

Software – – – – –

Steady-state Operation
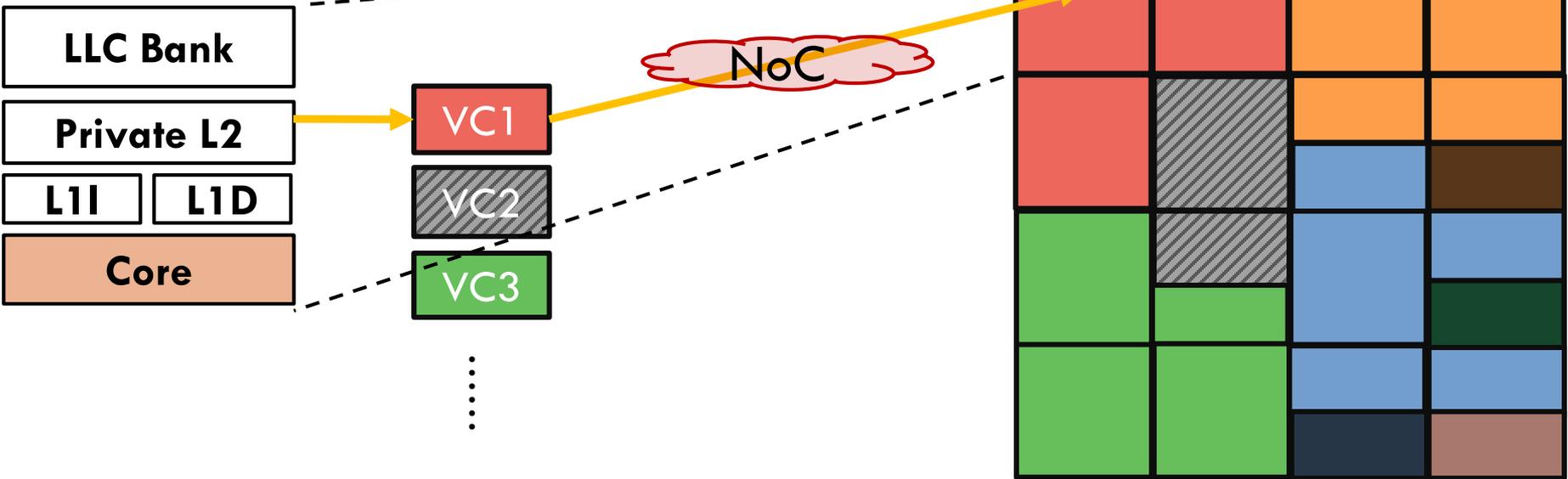
Place threads & data

Sample accesses

Optimization

Monitoring
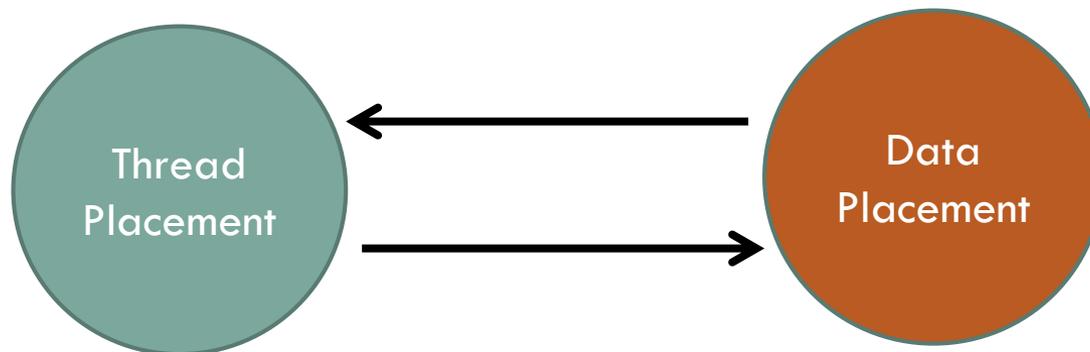
Miss curves

# Operation

□ Group partitions from different banks to create *virtual caches (VCs)*
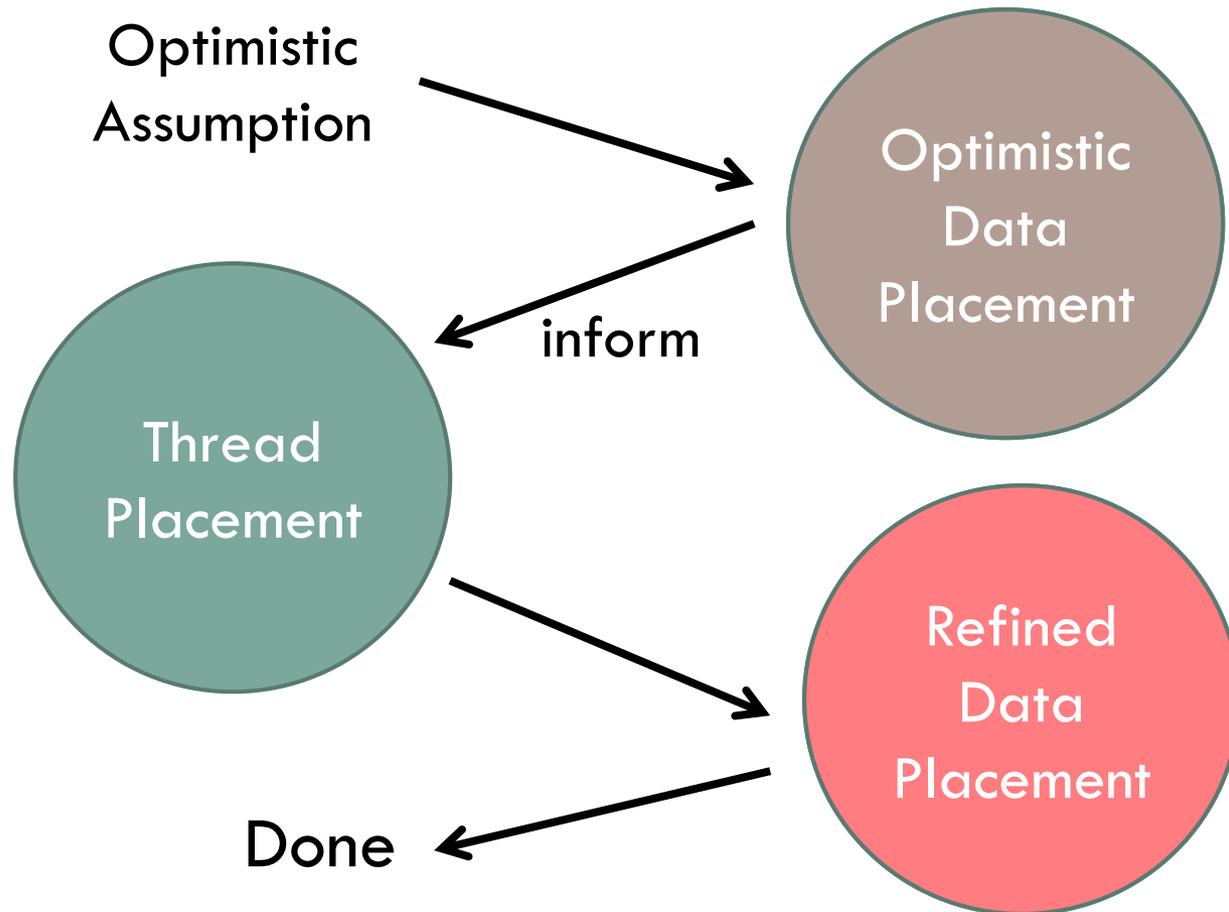
□ Similar to Jigsaw [Beckmann'13]

4x4 mesh NUCA LLC

# Optimization

- Minimize sum of on-chip latency and off-chip latency by deciding:
  - Thread placement
  - Virtual cache capacity
  - Virtual cache data placement

- It's an NP-hard problem
  - Thread and data placement are interrelated
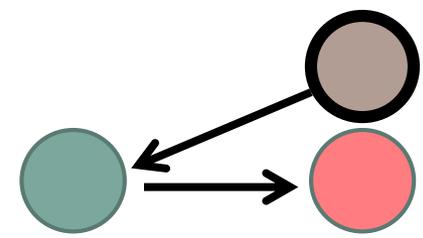  - Similar to VLSI place & route, HPC cluster scheduling
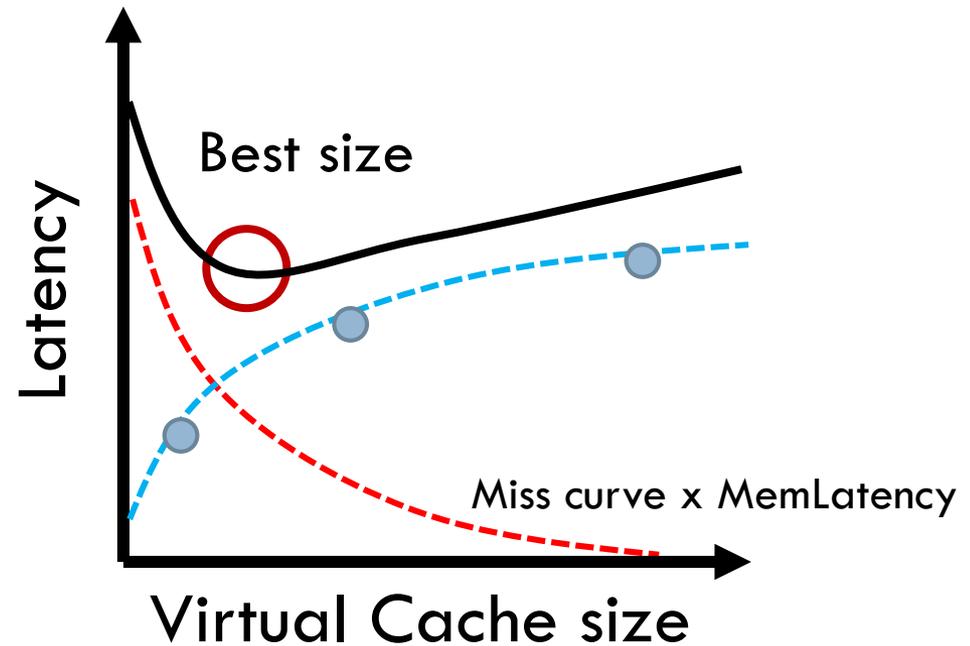
# Insight: Decouple the dependency



Optimistic
Assumption

Optimistic
Data
Placement

inform

Thread
Placement

Refined
Data
Placement

Done

By placing data twice, CDCS disentangles the dependencies

# Latency-aware allocation

□ Assume no contention



Best size

Latency

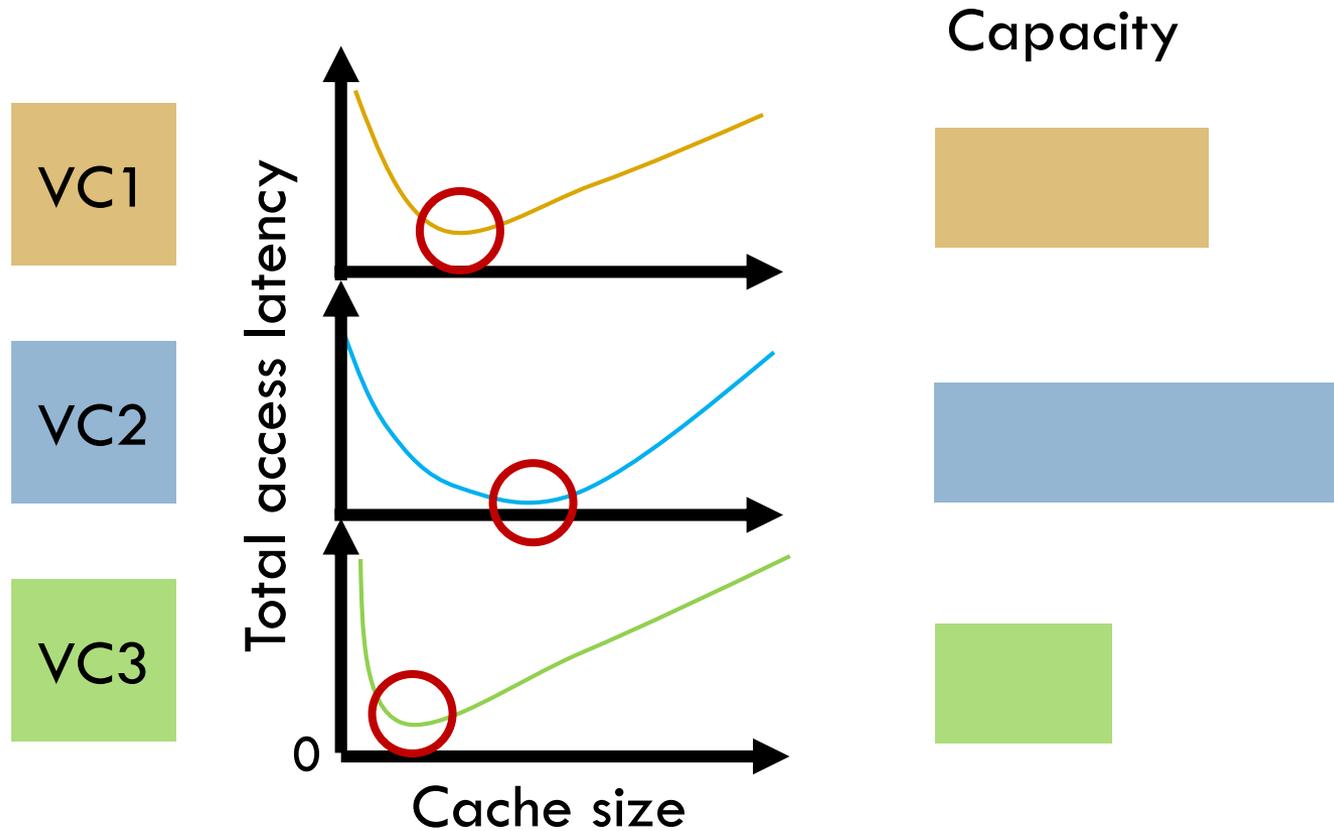Miss curve x MemLatency

Virtual Cache size

On-chip latency      – – –
Off-chip latency      – – –
Total latency      ——

# Latency-aware allocation

- Use total latency curve to partition cache among VCs

Capacity

VC1

VC2

VC3

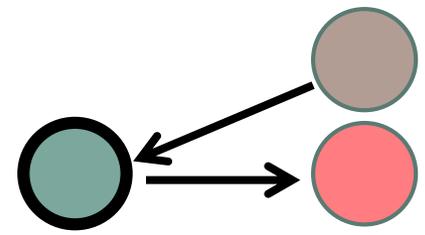Total access latency

Cache size

0

# Optimistic VC placement

□ Place VC as compactly as possible



Estimating contention
of every bank for VC

VC placed around
least-contended tile

# Thread placement

□ Place threads at *center of mass of their accesses*

# Refined VC placement

Greedily place VC close to thread first

Move/trade cache lines between VCs

# Scalable reconfiguration & monitoring

□ Incremental reconfiguration

  ◻ Allows chip to reconfigure smoothly, without pausing cores

□ Geometric monitor

  ◻ Monitors large LLC with low overhead

See paper for details

# Agenda

- Background

- CDCS design

- Evaluation
  - Methodology
  - Performance
  - Sensitivity

# Methodology

- Systems:
  - 64-core, 512KB/L3 bank
  - OOO cores (Silvermont-like)
  - 8x8 Mesh network
  - Similar to Knights Landing

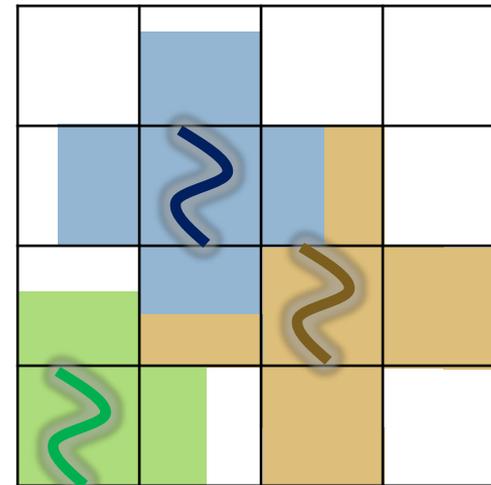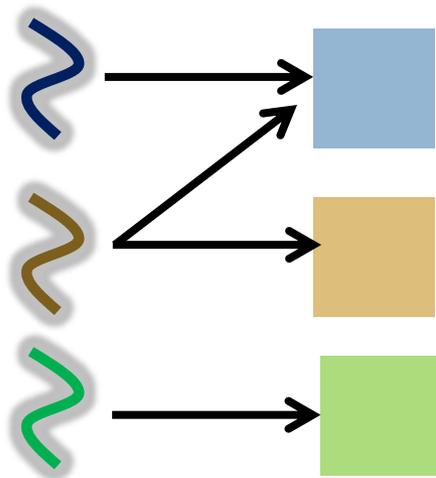64-core; 8x8 mesh network

| | |
|---|---|
| **Mem / IO** | |

| L3 Bank |
|---|
| L2 |

| L1I | L1D |
|---|---|

| OOO Core |
|---|

**Mem / IO**

- Zsim [Sanchez'13]: Pin-based simulator

- Workloads: SPEC CPU2006, SPEC OMP2012

# Methodology

☐ Schemes

    ❑ S-NUCA (baseline) with clustering thread scheduler

    ❑ R-NUCA with clustering thread scheduler

    ❑ Jigsaw

        ■ Jigsaw+C: Jigsaw with clustering thread scheduler

        ■ Jigsaw+R: Jigsaw with random thread scheduler

    ❑ CDCS

| | D-NUCA | Partitioned NUCA | CDCS | |
|---|---|---|---|---|
| | ✔ | ✔ | ✔ | Place data |
| | ✖ | ✔ | ✔ | Control capacity |
| | ✖ | ✖ | ✔ | Place threads |

# Multi-programmed mixes

Workloads that do not share data

CDCS avoids capacity contention more effectively than random scheduler

# Multi-threaded mixes

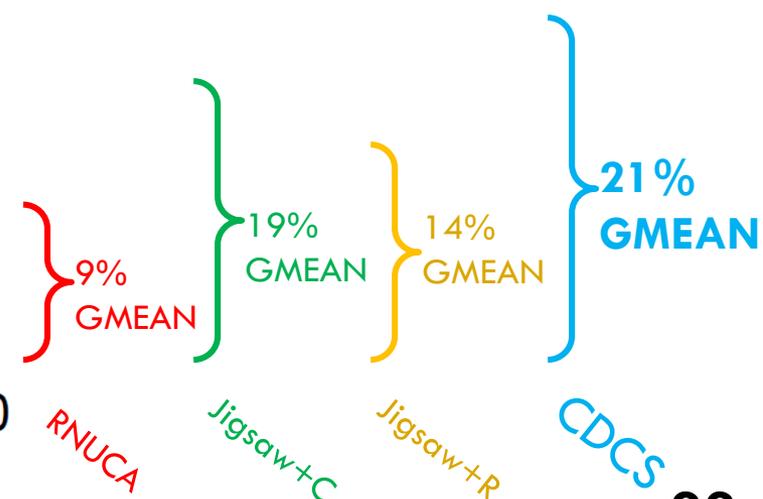Workloads that share data



CDCS guards against **pathological behavior** incurred by fixed thread scheduling policies

Clustering is better now

RNUCA — 9% GMEAN

Jigsaw+C — 19% GMEAN
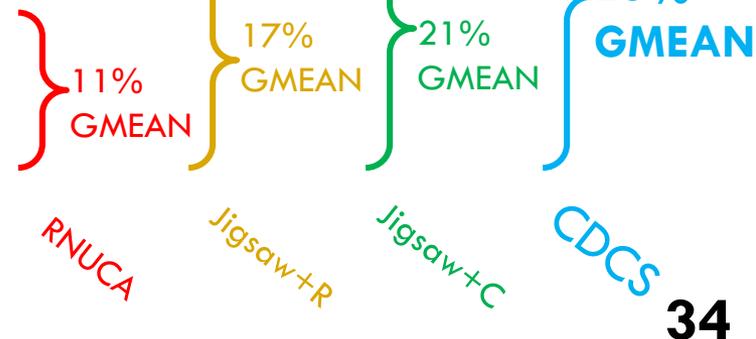
Jigsaw+R — 14% GMEAN

CDCS — **21% GMEAN**

# Undercommitted multi-threaded mixes

- SPECOMP mixes using half of the cores



With more flexibility, CDCS dynamically clusters or spreads out threads

11% GMEAN RNUCA

17% GMEAN Jigsaw+R

21% GMEAN Jigsaw+C

26% GMEAN CDCS

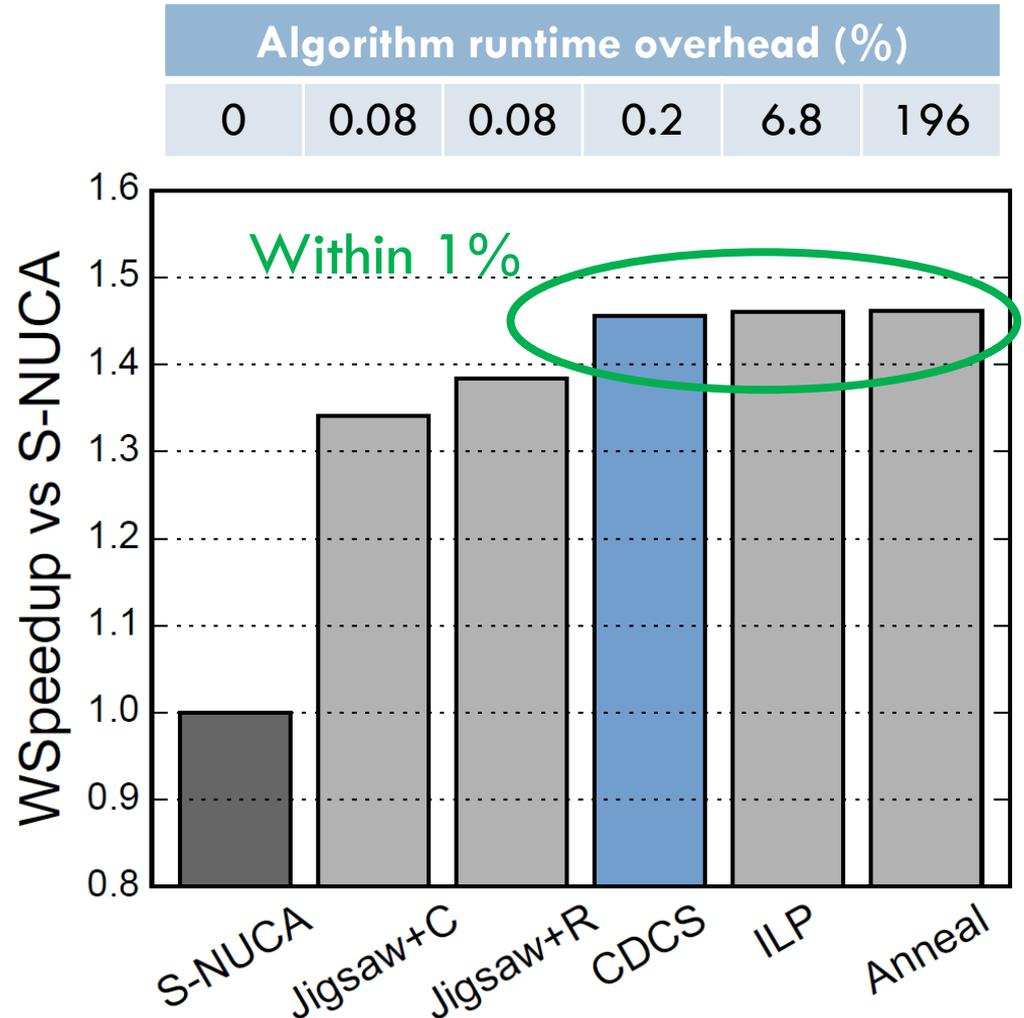# CDCS vs idealized algorithms

☐ Integer Linear Programming (ILP)

☐ Simulated annealing

Multi-programmed mixes

**Same** for multi-threaded mixes

| Algorithm runtime overhead (%) | | | | | |
|---|---|---|---|---|---|
| 0 | 0.08 | 0.08 | 0.2 | 6.8 | 196 |



Within 1%

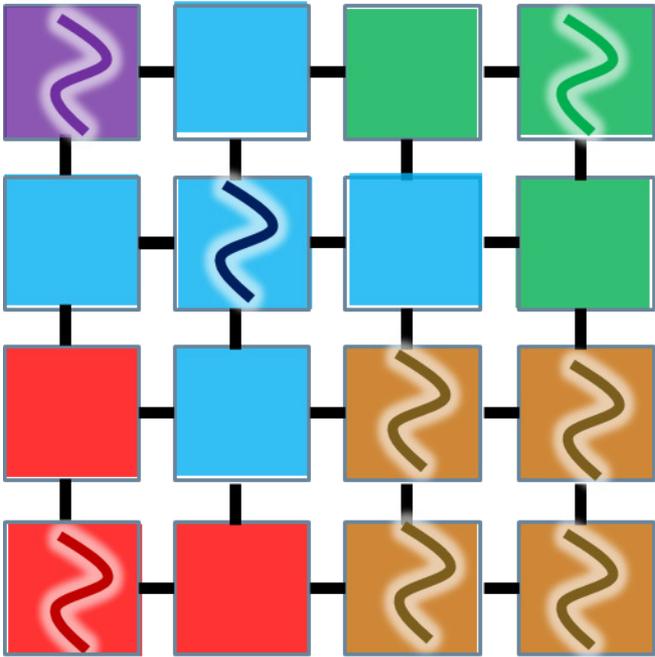WSpeedup vs S-NUCA — S-NUCA, Jigsaw+C, Jigsaw+R, CDCS, ILP, Anneal

# See paper for additional results

- Under-committed system

- Traffic breakdown

- Energy breakdown

- Factor analysis

- Other sensitivity studies

  - Reconfiguration interval sweep

  - Incremental reconfiguration IPC trace

# Conclusions

- Thread placement has a large impact on NUCA performance when capacity is well managed

- CDCS reduces the distance to data through joint thread and data placement

- CDCS outperforms state-of-the-art NUCA techniques with different thread scheduling policies and prevents pathological behavior of fixed policies

# QUESTIONS