# Προηγμένη Αρχιτεκτονική Υπολογιστών

## Non-Uniform Cache Architectures

Νεκτάριος Κοζύρης & Διονύσης Πνευματικάτος

{nkoziris,pnevmati}@cslab.ece.ntua.gr

Διαφάνειες από τον Ανδρέα Μόσχοβο, University of Toronto

8ο εξάμηνο ΣΗΜΜΥ – Ακαδημαϊκό Έτος: 2019-20

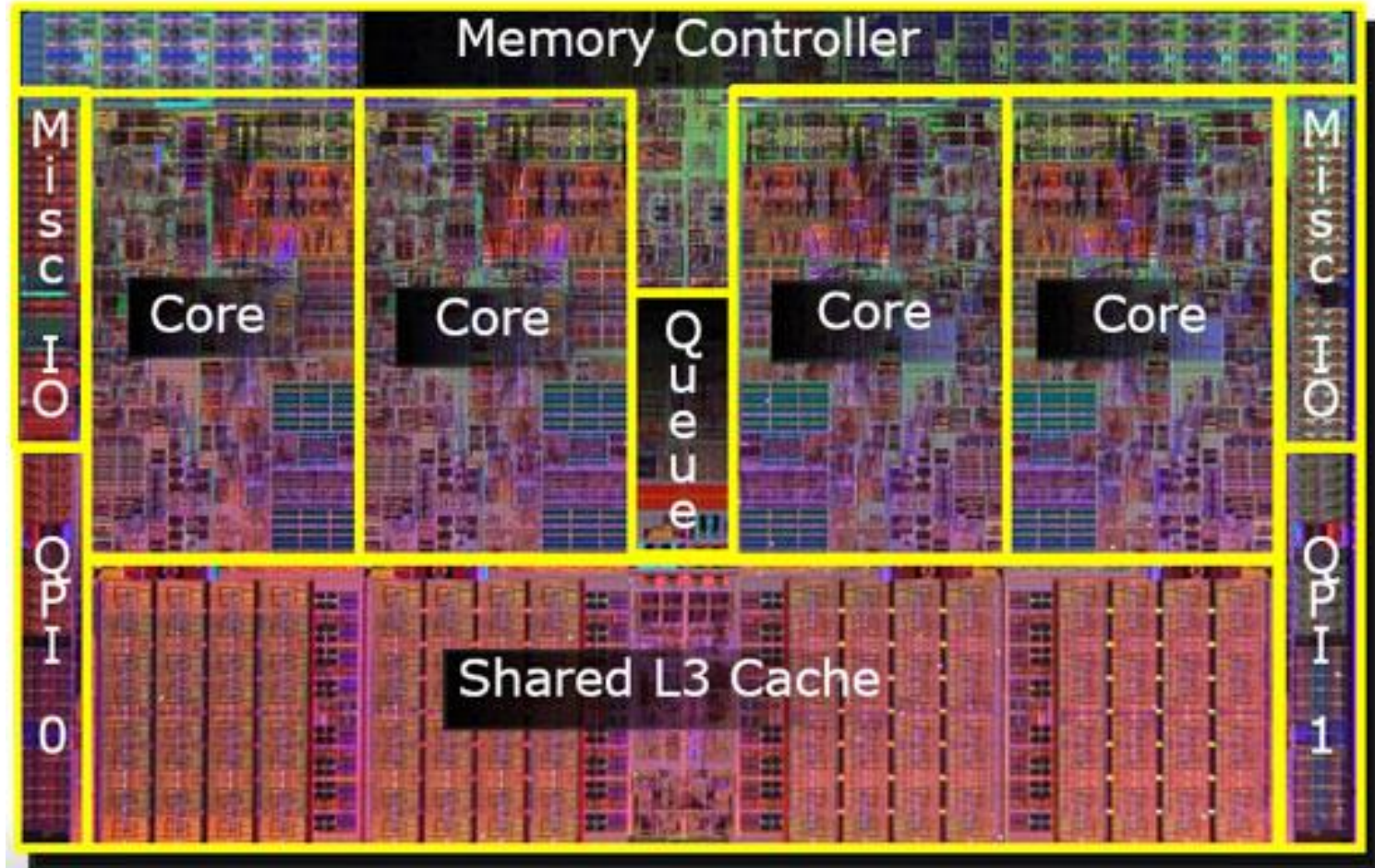http://www.cslab.ece.ntua.gr/courses/advcomparch/

- Sun Niagara T1



From http://jinsatoh.jp/ennui/archives/2006/03/opensparc.html

- Intel i7  (Nehalem)

Memory Controller

Misc IO

Misc IO

QPI 0

QPI 1

Core

Core

Queue

Core

Core

Shared L3 Cache

- AMD Shanghai



From http://www.chiparchitect.com
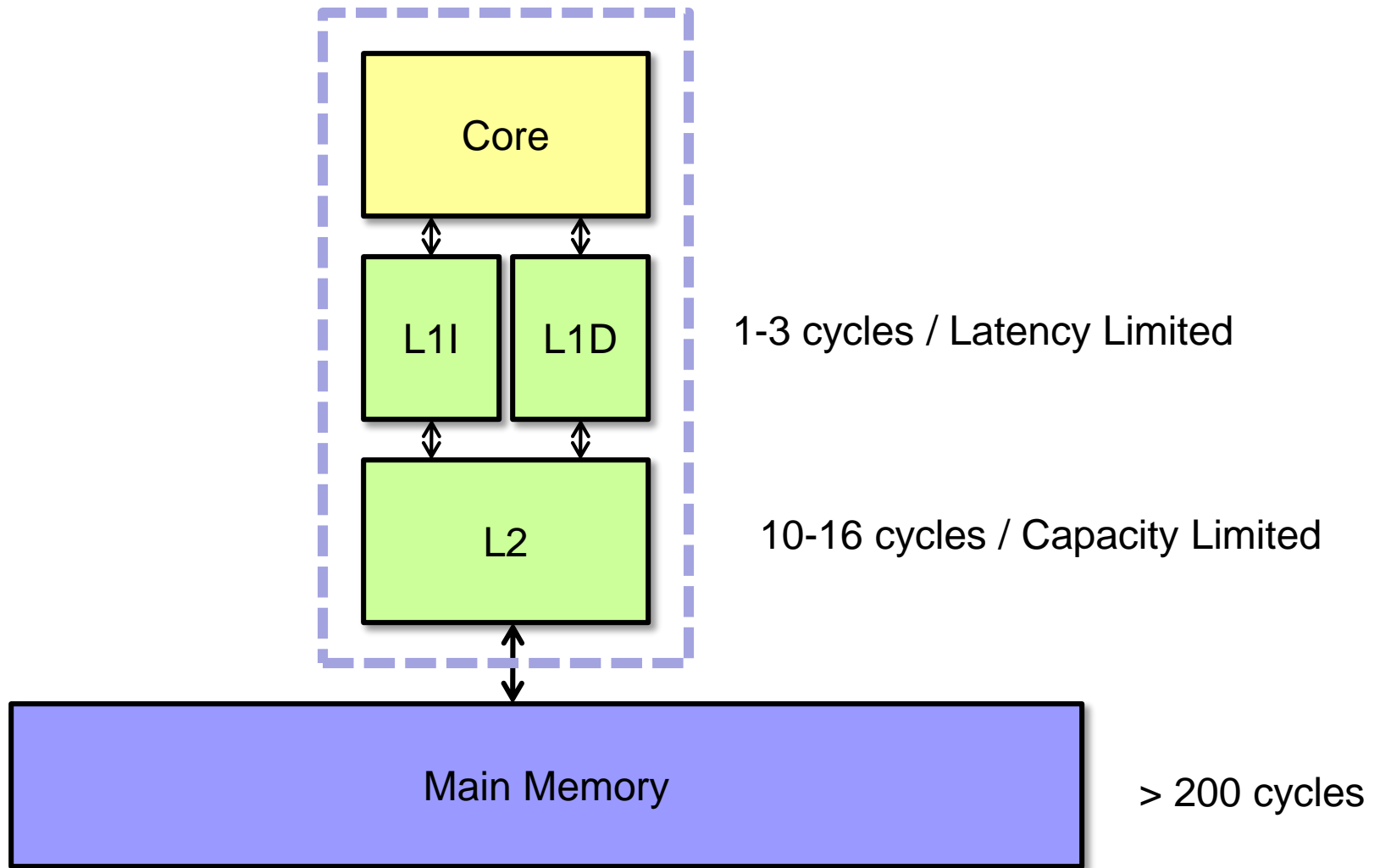
# Modern Processors Have Lots of Cores and Large Caches

- IBM Power 5
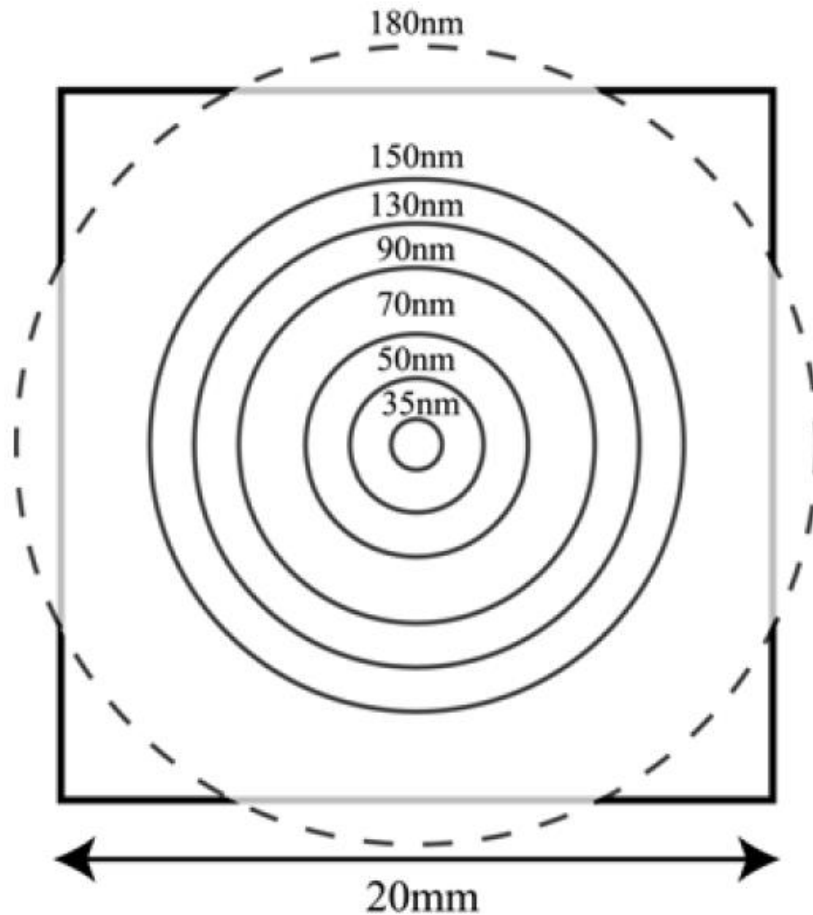
- Helps with Performance and Energy

  - Find graph with perfect vs. realistic memory system

# What Cache Design Used to be About



Core

L1I  L1D    1-3 cycles / Latency Limited

L2    10-16 cycles / Capacity Limited

Main Memory    > 200 cycles

- L2: Worst Latency == Best Latency

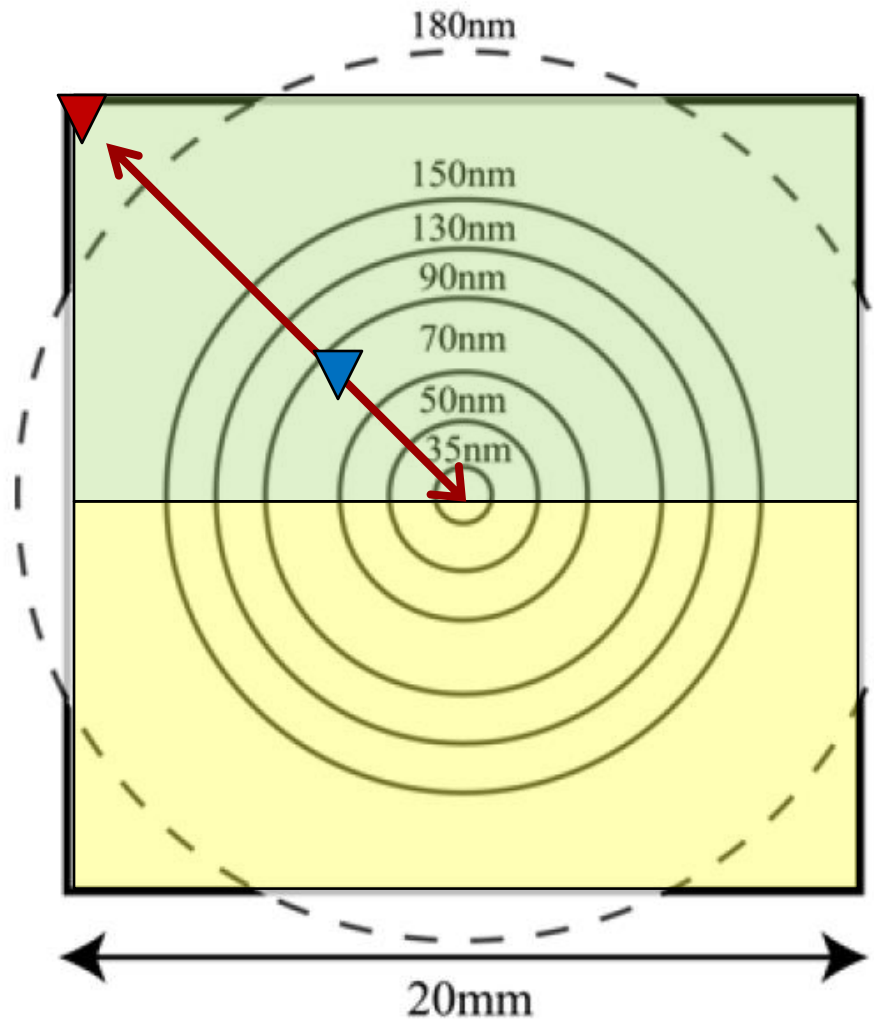- **Key Decision: What to keep in each cache level**

ISSCC 2003

**9.6    A Wire-Delay Scalable Microprocessor Architecture for High Performance Systems**

Stephen W. Keckler[1], Doug Burger[1], Charles R. Moore[1], Ramadass Nagarajan[1], Karthikeyan Sankaralingam[1], Vikas Agarwal[2], M.S. Hrishikesh[2], Nitya Ranganathan[1], Premkishore Shivakumar[1]

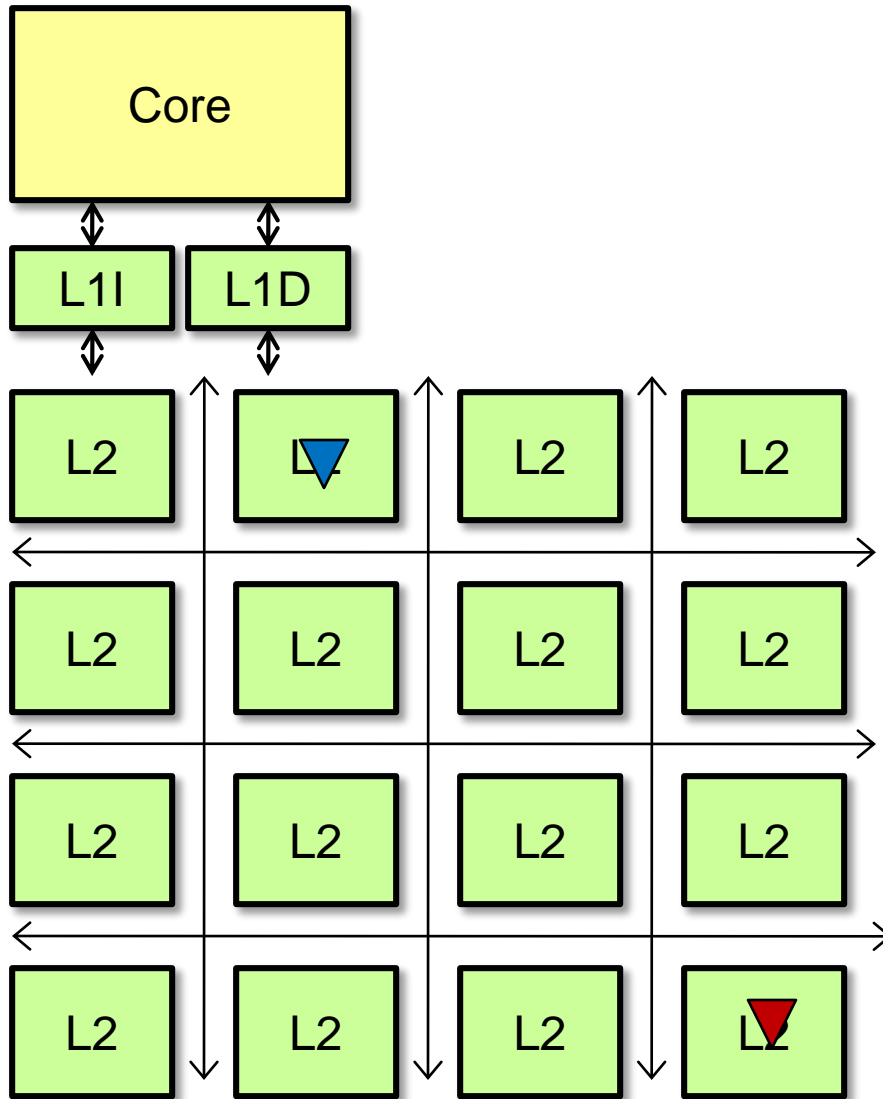**Figure 9.6.2: Projected fraction of chip reachable in one cycle with an 8FO4 clock period.**

Figure 9.6.2: Projected fraction of chip reachable in one cycle with an 8FO4 clock period.

- **Where something is matters**

- **More time for longer distances**
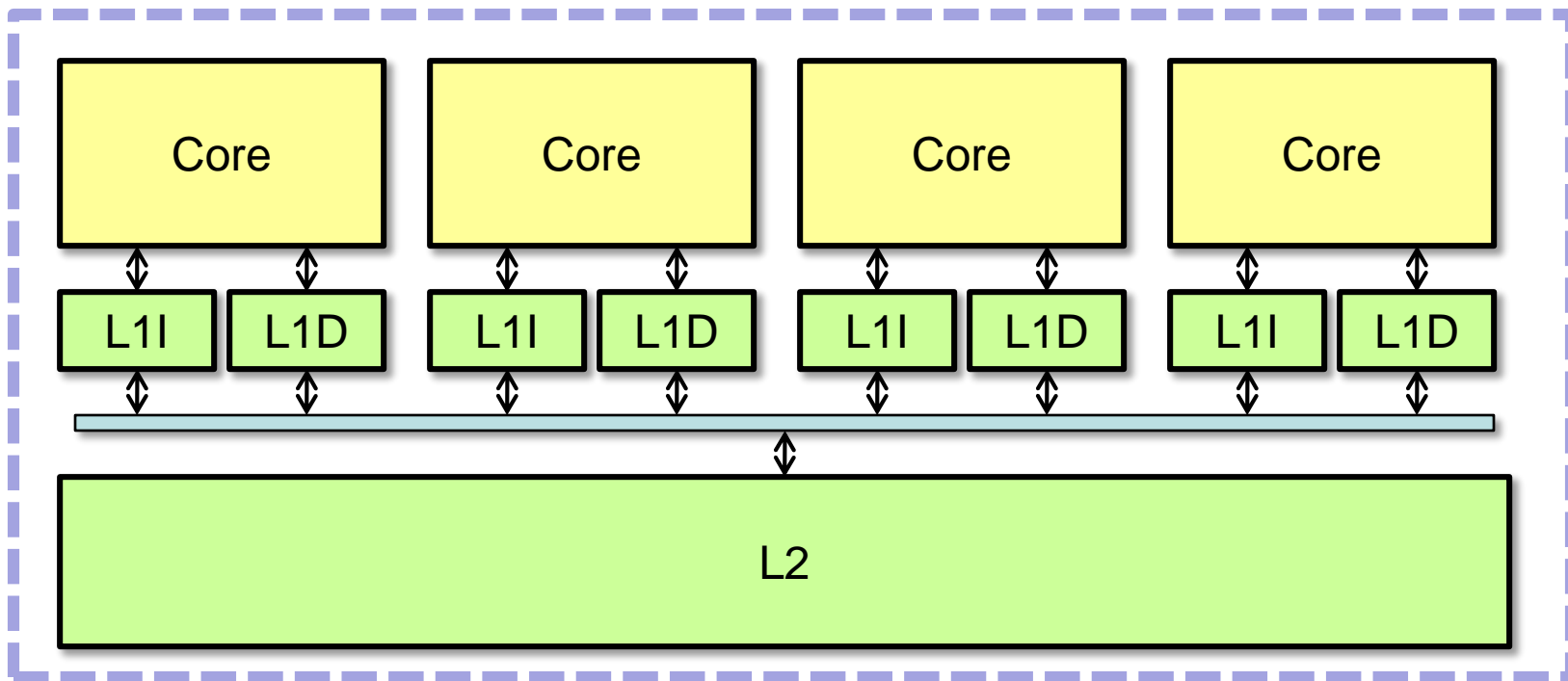
- Tiled Cache
- Variable Latency
- Closer tiles = Faster

- **Key Decisions:**
  - **Not only *what* to cache**
  - **Also *where* to cache**

- Initial Research focused on Uniprocessors

- Data Migration Policies
  - When to move data among tiles
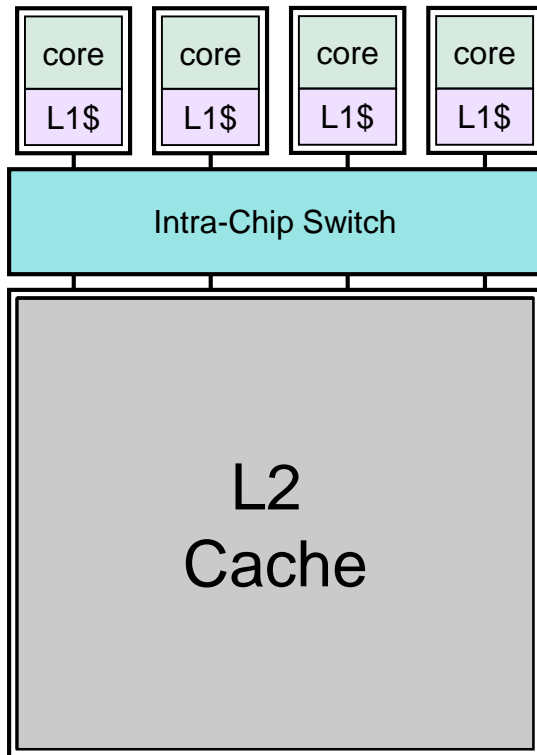
- L-NUCA: Fine-Grained NUCA

# Another Development: Chip Multiprocessors



- Easily utilize on-chip transistors

- Naturally exploit thread-level parallelism

- Dramatically reduce design complexity

- Future CMPs will have more processor cores

- Future CMPs will have more cache

Text from Michael Zhang & Krste Asanovic, MIT

# Initial Chip Multiprocessor Designs



| core | core | core | core |
| L1$ | L1$ | L1$ | L1$ |

Intra-Chip Switch

L2
Cache

*A 4-node CMP with a
large L2 cache*

- *Layout*: *"Dance-Hall"*
  - *Core + L1 cache*
  - *L2 cache*

- *Small L1 cache:* Very low access latency

- *Large L2 cache*

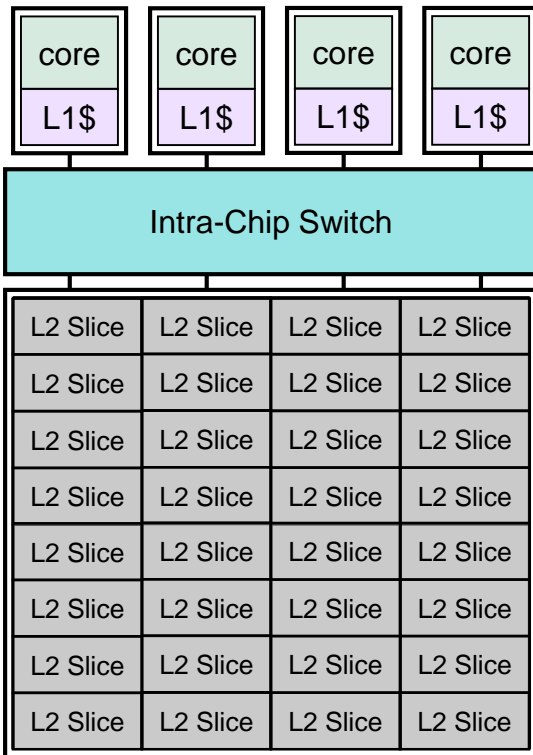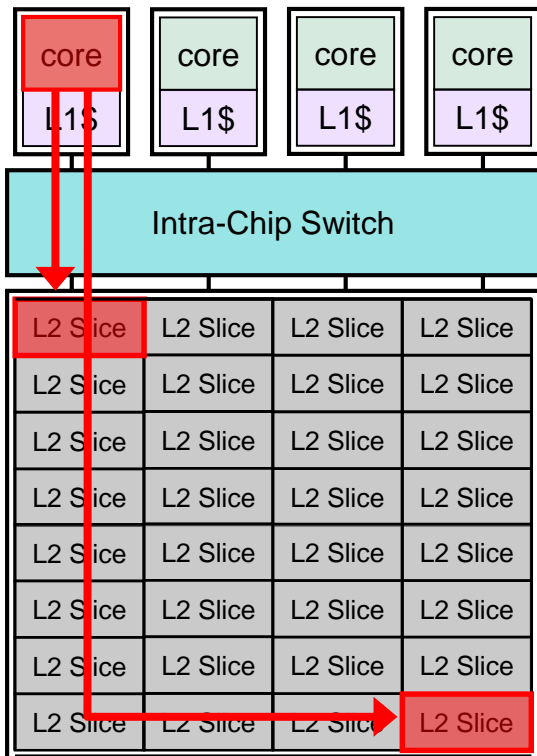Slide from Michael Zhang & Krste Asanovic, MIT

# Chip Multiprocessor w/ Large Caches



A 4-node CMP with a large L2 cache

- *Layout*: "Dance-Hall"
  - *Core + L1 cache*
  - *L2 cache*

- *Small L1 cache:* Very low access latency

- *Large L2 cache:* Divided into slices to minimize access latency and power usage

Slide from Michael Zhang & Krste Asanovic, MIT

# Chip Multiprocessors + NUCA



*A 4-node CMP with a large L2 cache*

- *Current:* Caches are designed with (long) uniform access latency for the worst case:
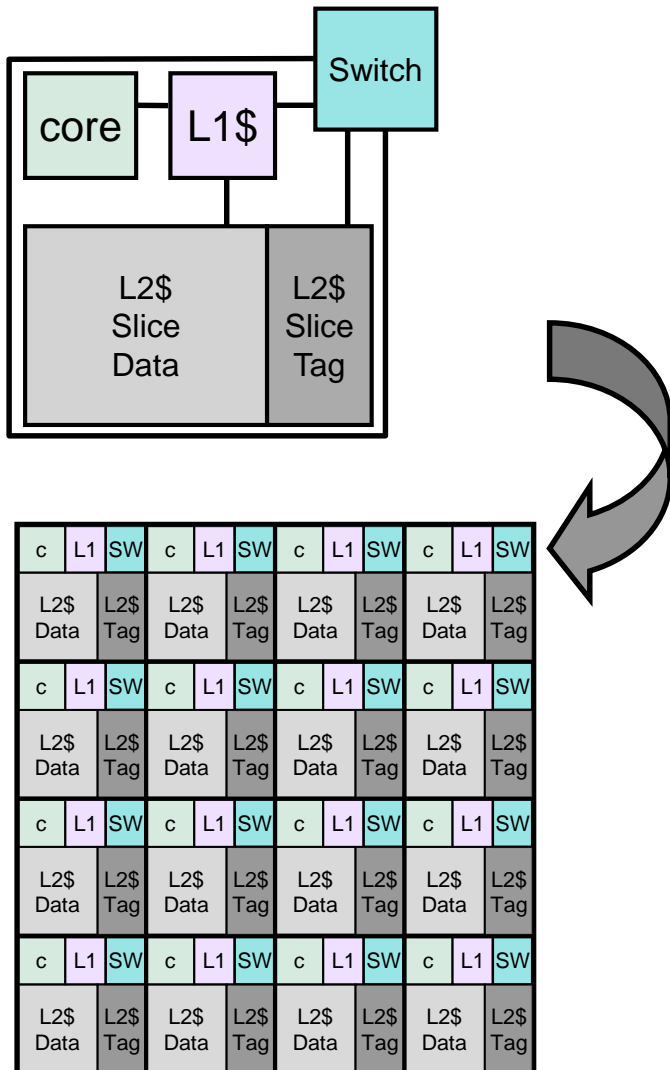
  ***Best Latency == Worst Latency***

- *Future:* Must design with non-uniform access latencies depending on the on-die location of the data:

  ***Best Latency << Worst Latency***

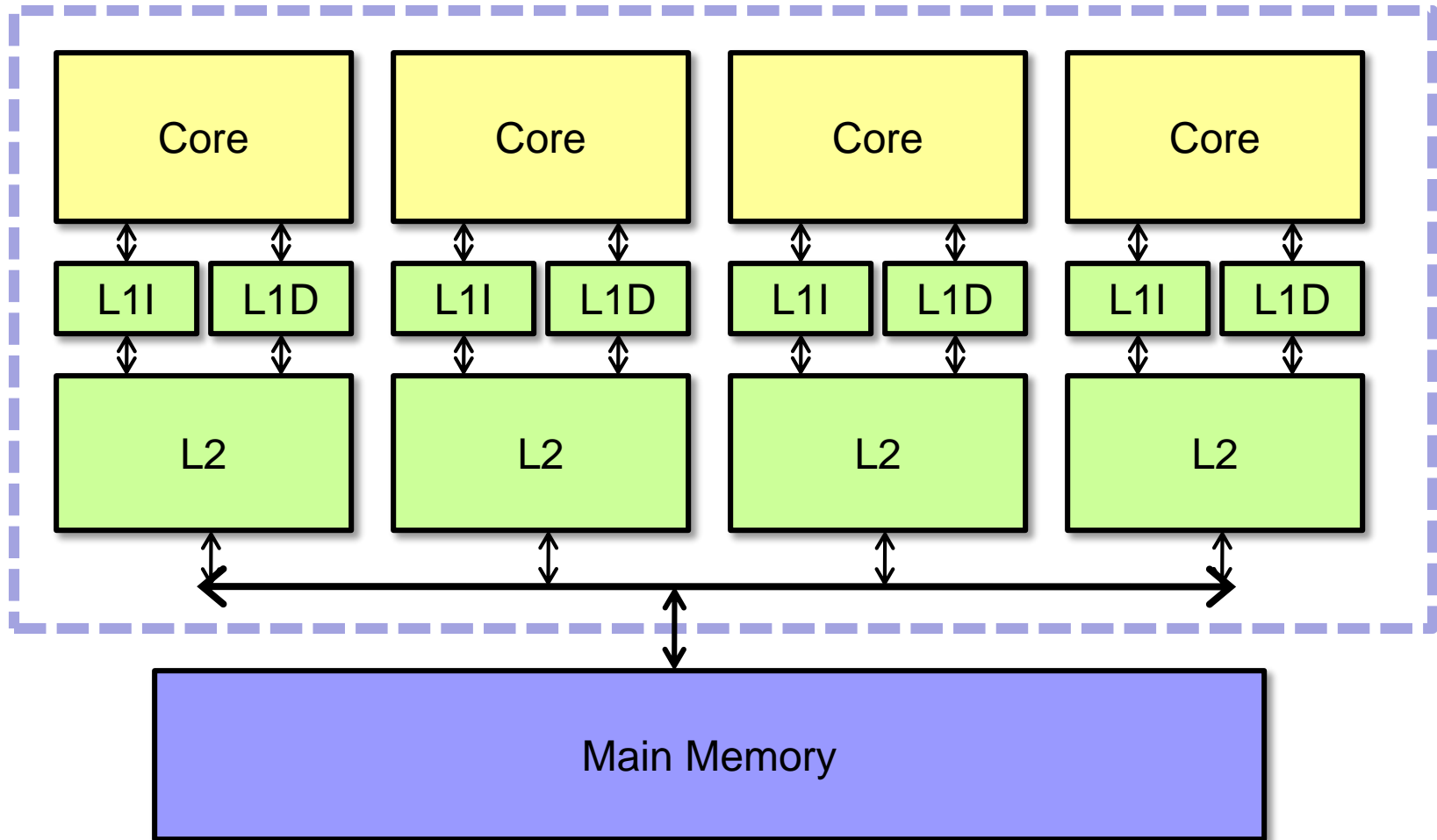- *Challenge:* How to minimize average cache access latency:

  ***Average Latency → Best Latency***

Slide from Michael Zhang & Krste Asanovic, MIT
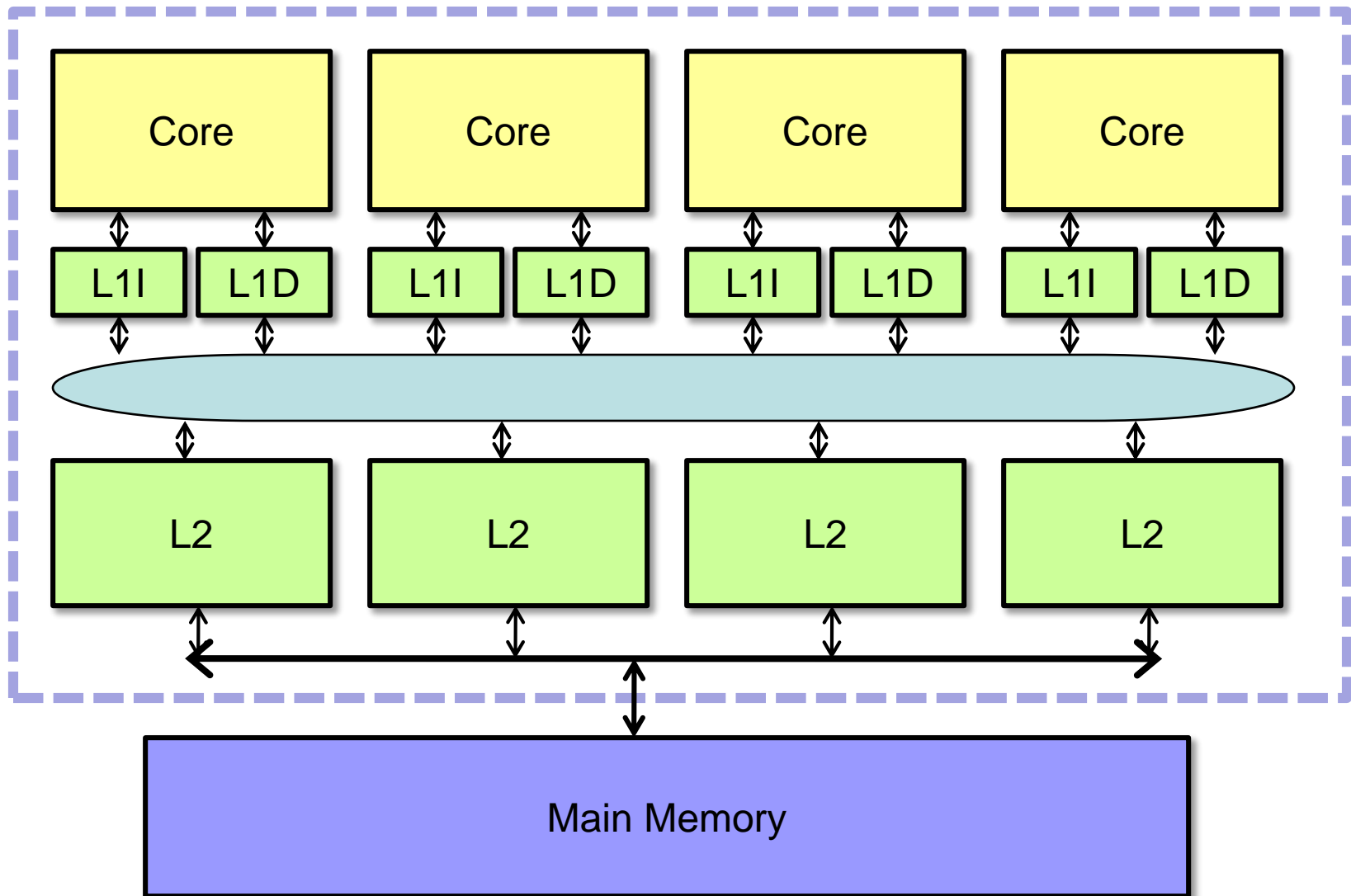
# Tiled Chip Multiprocessors

- **Tiled CMPs for *Scalability***
  - *Minimal redesign effort*
  - *Use directory-based protocol for scalability*

- **Managing the L2s to minimize the effective access latency**
  - *Keep data close to the requestors*
  - *Keep data on-chip*

Slide from Michael Zhang & Krste Asanovic, MIT

- + Low Latency
- - Fixed allocation

# Option #2: Shared Caches



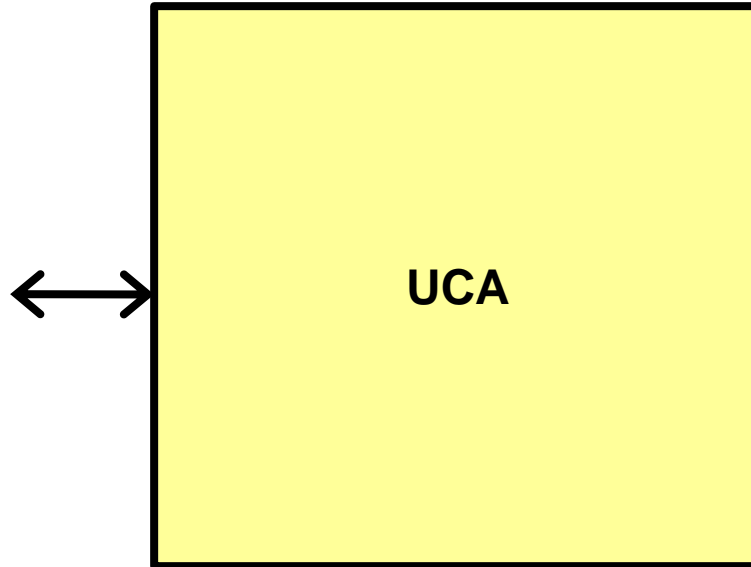- Higher, variable latency
- One Core can use all of the cache

- Get the best of both worlds
  - Low Latency of Private Caches
  - Capacity Adaptability of Shared Caches

# NUCA: A Non-Uniform Cache Access Architecture for Wire-Delay Dominated On-Chip Caches

**Changkyu Kim, D.C. Burger, and S.W. Keckler,**

10th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-X), October, 2002.
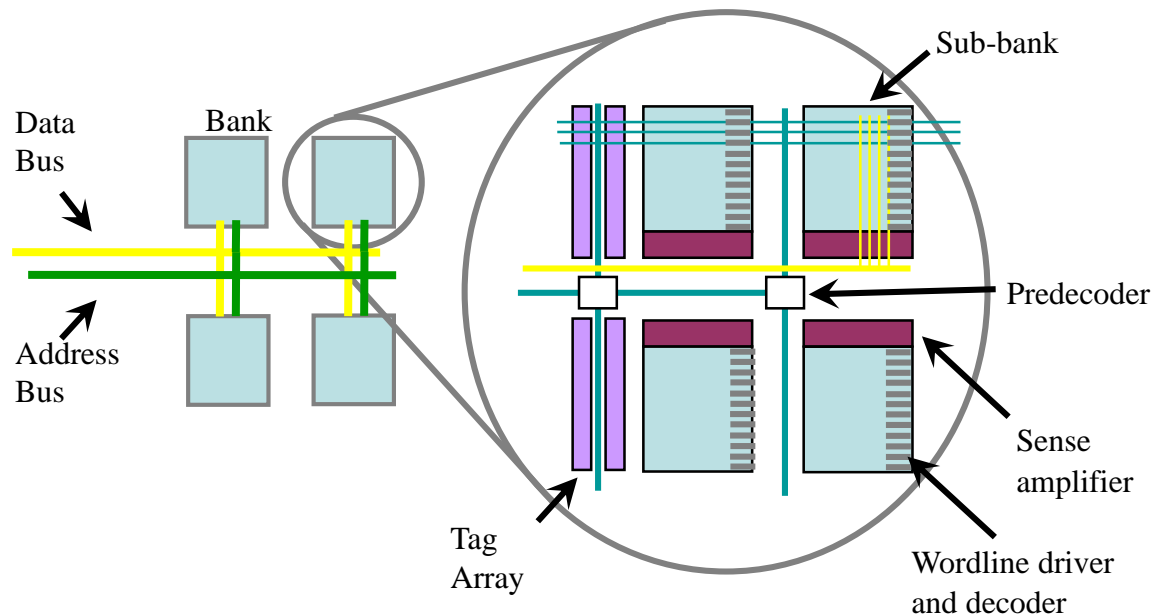
- **UCA:** Uniform Access Cache



- **Best Latency = Worst Latency**
  - **Time to access the farthest possible bank**

- **Partitioned in Banks**



Data Bus

Bank

Address Bus
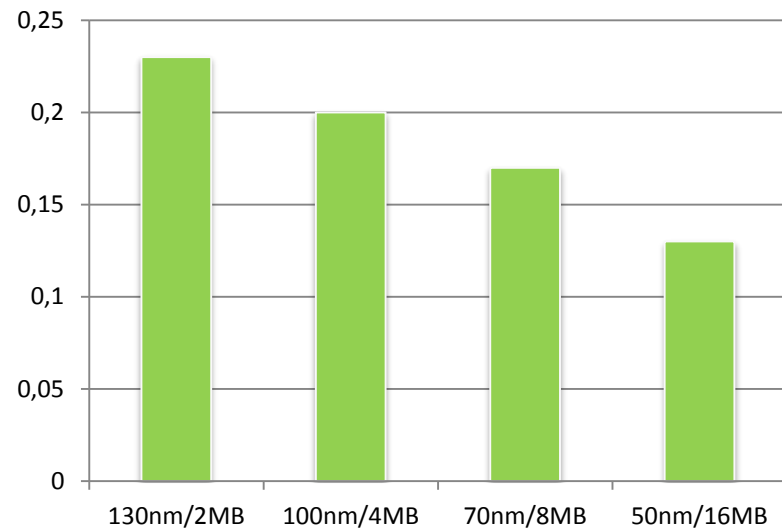
Sub-bank

Predecoder

Sense amplifier

Tag Array

Wordline driver and decoder

- Conceptually a single address and a single data bus
  - Pipelining can increase throughput
- See **CACTI** tool:
  - http://www.hpl.hp.com/research/cacti/
  - http://quid.hpl.hp.com:9081/cacti/
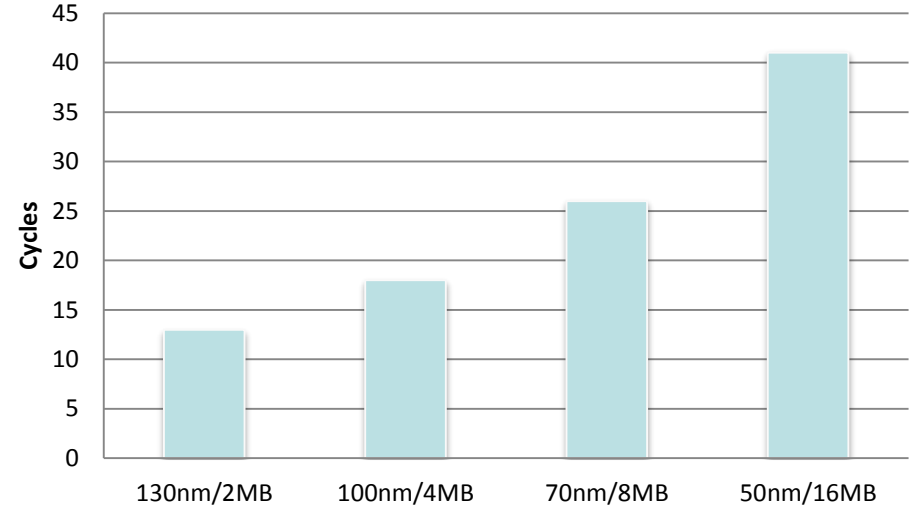
# Experimental Methodology

- SPEC CPU 2000

- Sim-Alpha

- CACTI

- 8 FO4 cycle time

- 132 cycles to main memory

- Skip and execute a sample

- Technology Nodes
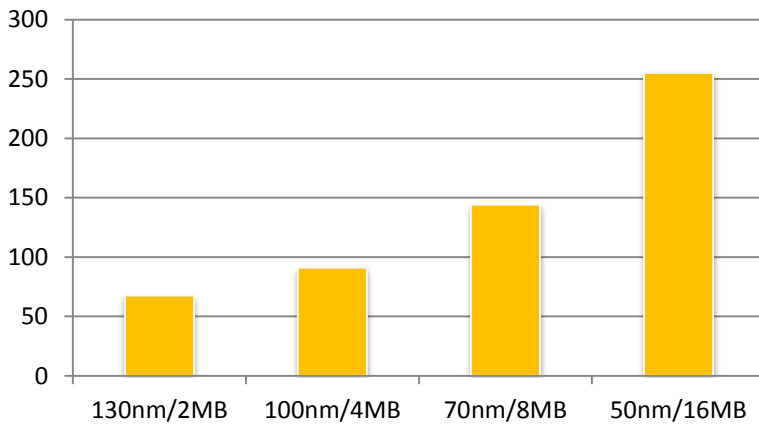  - 130nm, 100nm, 70nm, 50nm

# UCA Scaling – 130nm to 50nm



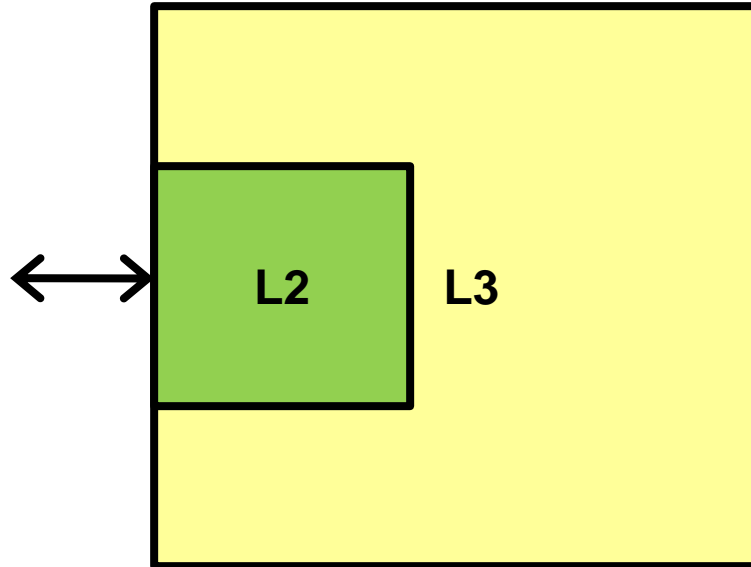**Relative Latency and Performance Degrade as Technology Improves**

- Loaded Latency: Contention
  - Bank
  - Channel
    - Bank may be free but path to it is not

- Conventional Hierarchy



- Common Usage:

  - **Serial-Access** for Energy and Bandwidth Reduction

- This paper:

  - **Parallel Access**

  - Prove that even then their design is better

# ML-UCA Evaluation

**Latency**



- 130nm/512K/2MB
- 100nm/512K/4M
- 70nm/1M/8M
- 50nm/1M/16M

■ Unloaded L2   ■ Unloaded L3

**IPC**



- 130nm/512K/2MB
- 100nm/512K/4M
- 70nm/1M/8M
- 50nm/1M/16M

- Better than UCA

- Performance Saturates at 70nm

- No benefit from larger cache at 50nm

- Aggressively banked

- Multiple parallel accesses

- Static NUCA with per bank set busses



| Tag | Set | Offset |
|-----|-----|--------|

**Bank Set**

- Use private **per bank set** channel
- Each bank has its distinct access latency
- A given address maps to a given bank set
  - Lower bits of **block address**

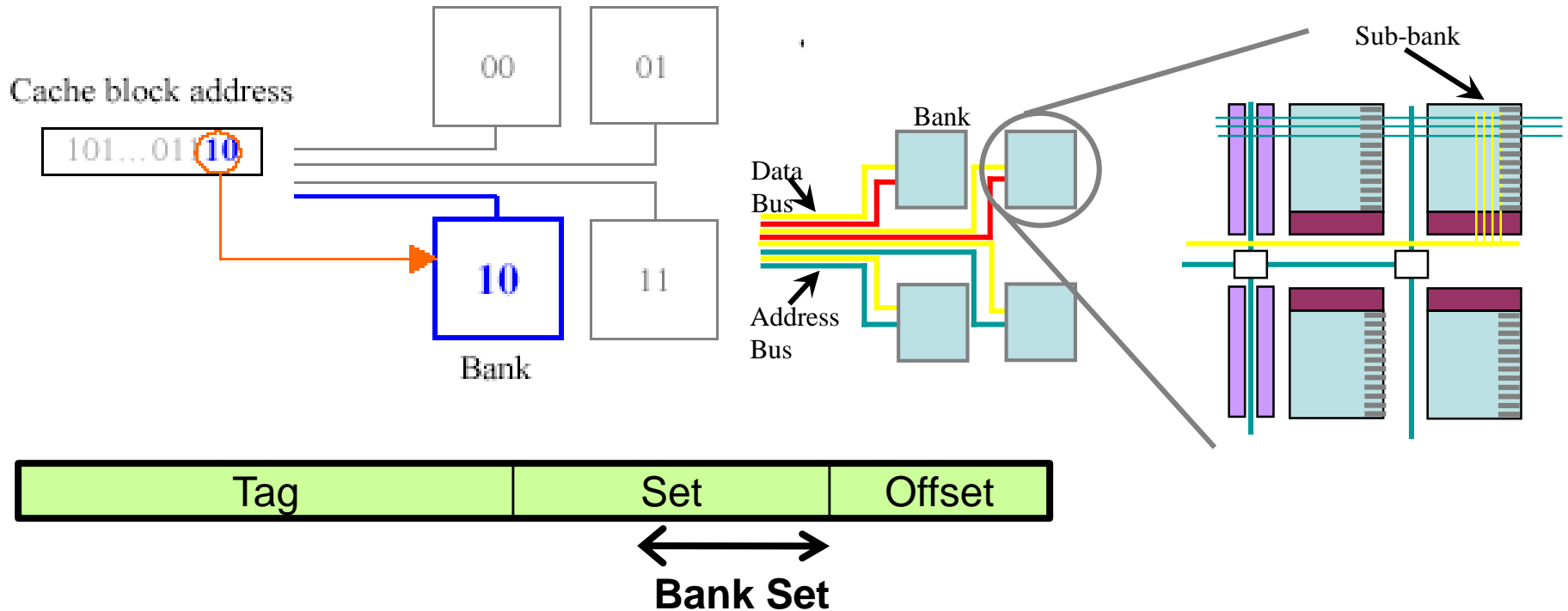- How fast can we initiate requests?
  - If c = scheduler delay
- Conservative **/Realistic:**
  - Bank + 2 x interconnect + c



- Aggressive **/ Unrealistic**:
  - Bank + c

- What is the optimal number of bank sets?
  - Exhaustive evaluation of all options
  - Which gives the highest IPC

# S-NUCA-1 Latency Variability



- Variability increases for finer technologies
- Number of banks does not increase beyond 4M
  - Overhead of additional channels
  - Banks become larger and slower

**Aggr. Loaded**



- Better than ML-UCA

# S-NUCA-1: IPC Performance

**IPC S1**



- Per bank channels become an overhead
- Prevent finer partitioning @70nm or smaller

- Use a 2-D Mesh P2P interconnect



- Wire overhead much lower:

  - S1: 20.9% vs. S2: 5.9% at 50nm and 32banks

- Reduces contention

- 128-bit bi-directional links

# S-NUCA2 vs. S-NUCA1 Unloaded Latency



- S-NUCA2 almost always better

# S-NUCA2 vs. S-NUCA-1 IPC Performance



Aggressive scheduling for S-NUCA1
Channel fully pipelined

Legend: IPC S1, IPC S2

X-axis: 130nm/2M/16, 100nm/4M/32, 70nm/8M/32, 50nm/16M/32

- S2 better than S1

- Data can dynamically migrate

- Move frequently used cache lines closer to CPU

**One way of each set in fast d-group**; compete within set

Cache blocks "screened" for fast placement

- **Where can a block map to?**



- **Simple Mapping**
- All 4 ways of each bank set need to be searched
- Farther bank sets → longer access

bank

8 bank sets

one set

way 0 →

way 1 →

way 2 →

way 3 →

- **Fair Mapping**
- Average access times across all bank sets are equal

- **Shared Mapping**
- Sharing the closest banks → every set has some fast storage
- If n bank sets share a bank then all banks must be n-way set associative

- **Where is a block?**

- **Incremental Search**
  - Search in order

- **Multicast**
  - Search all of them in parallel

- **Partitioned Multicast**
  - Search groups of them in parallel



way 0
way 1
way 2
way 3

- ## Tags are distributed
  - May search many banks before finding a block
  - Farthest bank determines miss determination latency

- ## Solution: Centralized Partial Tags
  - Keep a few bits of all tags (e.g., 6) at the cache controller
  - If no match → Bank doesn't have the block
  - If match → Must access the bank to find out

**Partial Tags** → R.E. Kessler, R. Jooss, A. Lebeck, and M.D. Hill. Inexpensive implementations of set-associativity. In *Proceedings of the 16th Annual International Symposium on Computer Architecture,* pages 131–139, May 1989.

- SS-Performance:
  - Partial Tags and Banks accessed in parallel
  - Early Miss Determination
  - Go to main memory if no match
  - Reduces latency for misses

- SS-Energy:
  - Partial Tags first
  - Banks only on potential match
  - Saves energy
  - Increases Delay

- Want data that will be accessed to be close

- Use LRU?

  - Bad idea: must shift all others



- Generational Promotion

  - Move to next bank

  - Swap with another block

- Where to place a new block coming from memory?

- Closest Bank?

  - May force another important block to move away

- Farthest Bank?

  - Takes several accesses before block comes close

- A new block must replace an older block **victim**

- What happens to the victim?

- **Zero Copy**
  - Get's dropped completely

- **One Copy**
  - Moved away to a slower bank (next bank)

- DN-BEST
  - Shared Mapping
  - SS Energy
    - Balance performance and access account/energy
    - Maximum performance is 3% higher
  - Insert at tail
    - Insert at head → reduces avg. latency but increases misses
  - Promote on hit
    - No major differences with other polices

# Baseline D-NUCA

- Simple Mapping
- Multicast Search
- One-Bank Promotion on Hit
- Replace from the slowest bank

# D-NUCA Unloaded Latency

# Performance Comparison



UPPER = all hits are in the closest bank 3 cycle latency

- D-NUCA and S-NUCA2 scale well
- D-NUCA outperforms all other designs
- ML-UCA saturates – UCA Degrades

# Distance Associativity for High-Performance Non-Uniform Cache Architectures

**Zeshan Chishti, Michael D Powell, and T. N. Vijaykumar**

Slides mostly directly from the authors' conference presentation

Large Cache Design
- L2/L3 growing (e.g., 3 MB in Itanium II)
- Wire-delay becoming dominant in access time

Conventional large-cache
- Many subarrays => wide range of access times
- Uniform cache access => access-time of *slowest* subarray
- Oblivious to access-frequency of data

Want often-accessed data faster: improve access time

Pioneered Non-Uniform Cache Architecture

Access time: Divides cache into many distance-groups

- D-group closer to core => faster access time

Data Mapping: conventional

- Set determined by block index; each set has n-ways

Within a set, place frequently-accessed data in fast d-group

- Place blocks in farthest way; bubble closer if needed

Processor core

fast

slow

tag    data

way 0

way 1

way n-2

way n-1

d-group

One way of each set in fast d-group; compete within set

Cache blocks "screened" for fast placement

Processor core



fast

tag    data            way 0

                       way 1

slow                   way n-2

                       way n-1

                       d-group

One way of each set in fast d-group; compete within set
Cache blocks "screened" for fast placement

Processor core

fast

slow

tag | data | way 0

way 1

way n-2

way n-1

d-group

One way of each set in fast d-group; compete within set

Cache blocks "screened" for fast placement

One way of each set in fast d-group; compete within set

Cache blocks "screened" for fast placement

Processor core

fast

slow

| tag | |
|-----|---|
| | way 0 |

| | |
|-----|---|
| | way 1 |

| | |
|-----|---|
| | way n-2 |

| | |
|-----|---|
| | way n-1 |

d-group

One way of each set in fast d-group; compete within set

Cache blocks "screened" for fast placement

Processor core

fast

slow

tag data

way 0

way 1

way n-2

way n-1

d-group

Want to change restriction; more flexible data-placement

Artificial coupling between s-a way # and d-group
- Only one way in each set can be in fastest d-group
    - Hot sets have > 1 frequently-accessed way
    - Hot sets can place only one way in fastest d-group

Swapping of blocks is bandwidth- and energy-hungry
- D-NUCA uses a switched network for fast swaps

# Common Large-cache Techniques

Sequential Tag-Data: e.g., Alpha 21164 L2, Itanium II L3

- Access tag first, and then access only matching data
- Saves energy compared to parallel access

Data Layout: Itanium II L3

- Spread a block over many subarrays (e.g., 135 in Itanium II)
- For area efficiency and hard- and soft-error tolerance

These issues are important for large caches

Key observation:

- sequential tag-data => indirection through tag array
- Data may be located anywhere

**Distance Associativity:**

Decouple tag and data => flexible mapping for sets

Any # of ways of a hot set can be in fastest d-group

**NuRAPID** cache: Non-uniform access with Replacement And Placement usIng Distance associativity

Benefits:

More accesses to faster d-groups

Fewer swaps => less energy, less bandwidth

But:

More tags + pointers are needed

# Outline

- Overview


- NuRAPID Mapping and Placement


- NuRAPID Replacement


- NuRAPID layout


- Results


- Conclusion

Distance-Associative Mapping:
- decouple tag from data using forward pointer
- Tag access returns forward pointer, data location

Placement: data block can be placed anywhere
- Initially place all data in fastest d-group
- Small risk of displacing often-accessed block

## Tag array

## Data Arrays

frame #

set #

Way-0

Way-(n-1)

$A_{tag}, grp_0, frm_1$

0
1
2
3

· · ·

forward pointer

A

fast

d-group 0

0
1
k

0
1
k

d-group 1

0
1
k

d-group 2

slow

All blocks initially placed in fastest d-group

## Tag array

## Data Arrays

frame #

set #   Way-0                Way-(n-1)

$A_{tag}, grp_0, frm_1$          $B_{tag}, grp_0, frm_k$

0
1
2
3

0
1
k

A

B

0
1
k

0
1
k

d-group 0

d-group 1

d-group 2

fast

slow

Multiple blocks from one set in same d-group

Tag array

Data Arrays

frame #

set #   Way-0          Way-(n-1)

0   $A_{tag}, grp_0, frm_1$        $B_{tag}, grp_0, frm_k$
1
2   . . .
3
$C_{tag}, grp_2, frm_0$        $D_{tag}, grp_1, frm_1$

0
1
k

A
B

d-group 0

fast

0
1
k

D

d-group 1

0
1
k

C

d-group 2

slow

No coupling between tag and data mapping

- Overview

- NuRAPID Mapping and Placement

- NuRAPID Replacement

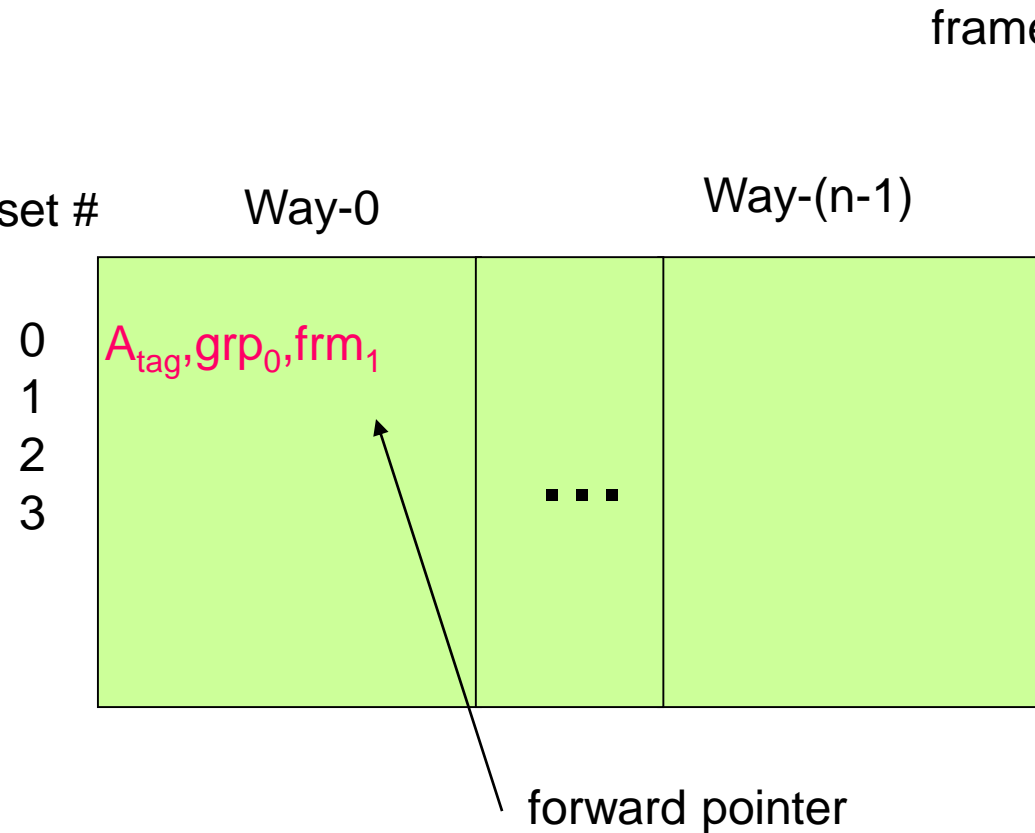- NuRAPID layout

- Results

- Conclusion

Two forms of replacement:

- Data Replacement: Like conventional
  - Evicts blocks from cache due to tag-array limits

- Distance Replacement: Moving blocks among d-groups
  - Determines which block to demote from a d-group
  - Decoupled from data replacement
  - No blocks evicted
    - Blocks are swapped

Tag array

Data Arrays

frame #

set #    Way-0                    Way-(n-1)

$Z_{tag}, grp_1, frm_k$

...

$B_{tag}, grp_0, frm_1$

0
1

0
1
k

B

d-group 0

fast

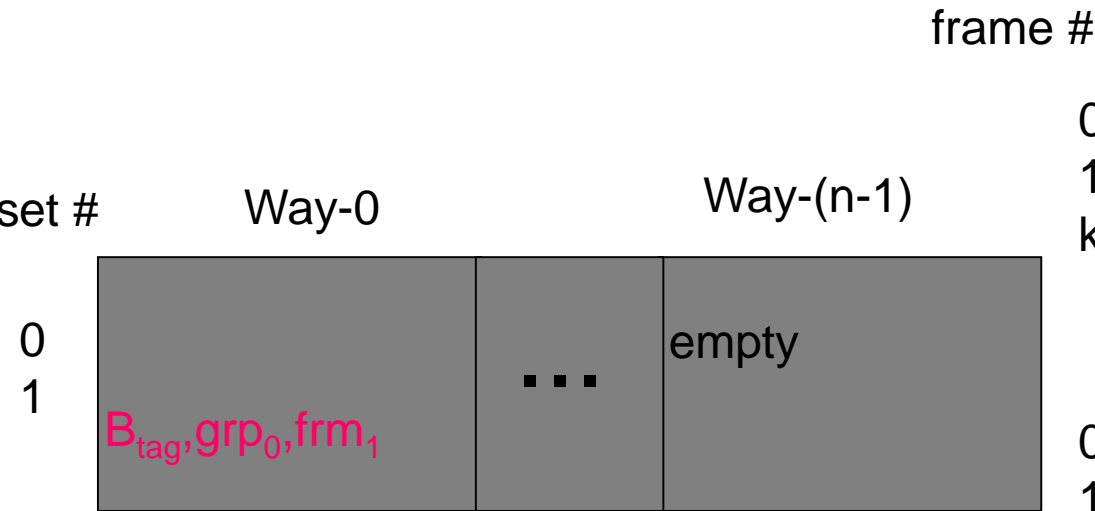0
1
k

Z

d-group 1

0
1
k

d-group 2

slow

**Place new block, A, in set 0.
Space must be created in
the tag set:**

Data-Replace **Z**

Z may not be in the target d-group

## Tag array

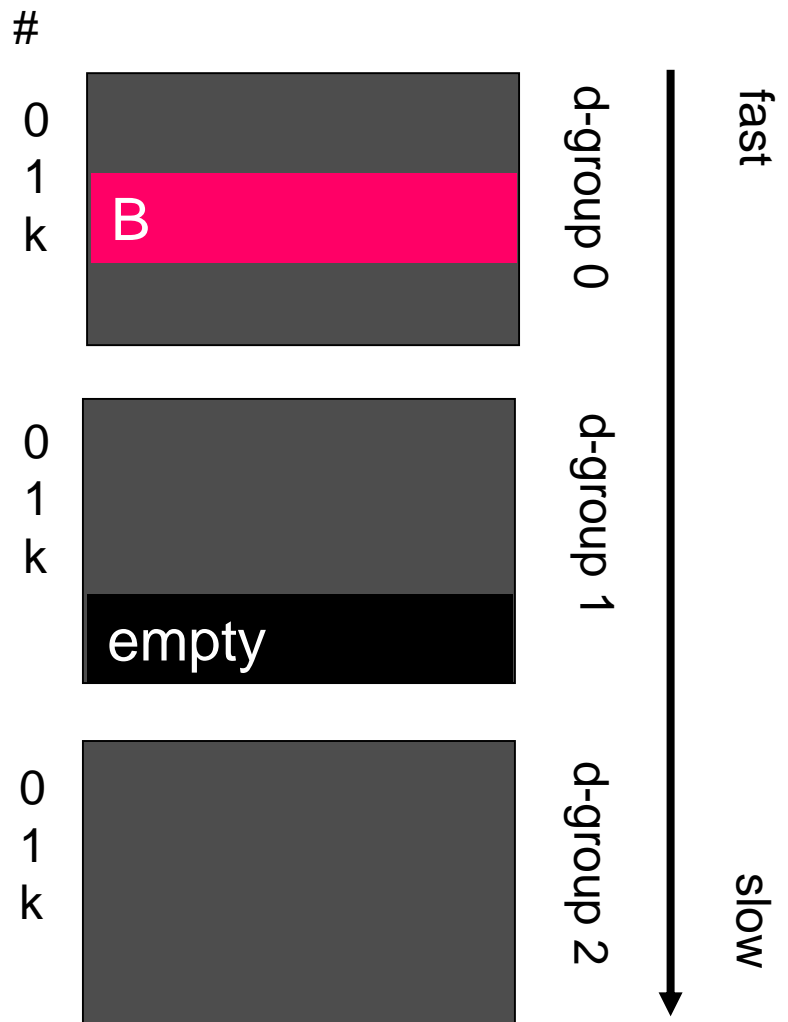## Data Arrays

frame #

set #   Way-0   Way-(n-1)



$B_{tag}$,$grp_0$,$frm_1$   empty

Place new block, A, in set 0.

Data-Replace Z

0
1
k
d-group 0   fast

B

0
1
k
d-group 1

empty

0
1
k
d-group 2   slow

reverse pointer

Tag array

Data Arrays

frame #

fast

d-group 0

0
1
k

B, $set_1$ $way_0$

set #

Way-0

Way-(n-1)

0
1

$A_{tag}$

...

$B_{tag}, grp_0, frm_1$

d-group 1

0
1
k

empty

d-group 2

0
1
k

slow

Place $A_{tag}$, in set 0.
**Must create an empty data block**
B is selected to demote. Use reverse-pointer to locate $B_{tag}$

# Tag array

# Data Arrays

frame #

set #   Way-0                           Way-(n-1)

Way-0                           $A_{tag}$

0
1
$B_{tag}, grp_1, frm_k$

0
1
k
empty

d-group 0
fast

0
1
k

d-group 1

B, $set_1$ $way_0$

B is demoted to **empty** frame.
$B_{tag}$ updated
There was an empty frame
because Z was evicted
This may not always be the
case

0
1
k

d-group 2
slow

## Tag array

## Data Arrays
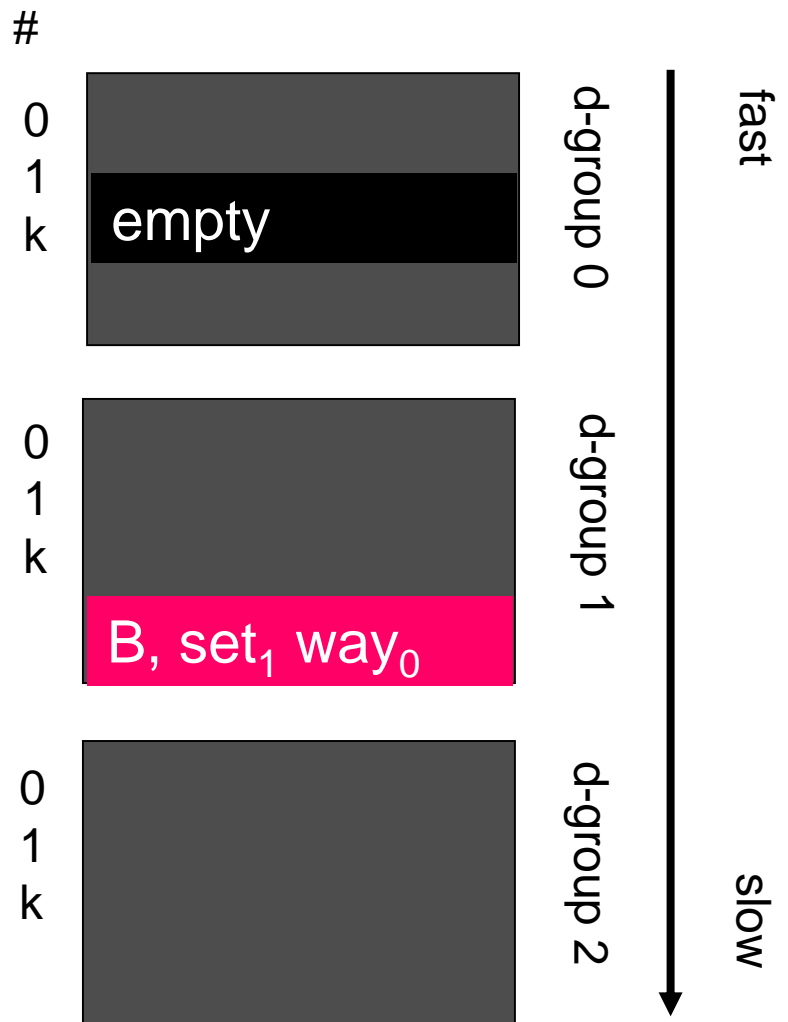
frame #

set #

Way-0

Way-$(n-1)$

. . .

$A_{tag}, grp_0, frm_1$

$B_{tag}, grp_1, frm_k$

0
1

0
1
k

A, $set_0$ $way_{n-1}$

d-group 0

0
1
k

B, $set_1$ $way_0$

d-group 1

0
1
k

d-group 2

fast

slow

A is placed in d-group 0
pointers updated

Always empty block for demotion for dist.-replacement
- May require multiple demotions to find it

Example showed only demotion
- Block could get stuck in slow d-group
- Solution: Promote upon access (see paper)

How to choose block for demotion?
- Ideal: LRU-group
- LRU hard.  We show random OK (see paper)
  - Promotions fix errors made by random

# Outline

- Overview

- NuRAPID Mapping and Placement

- NuRAPID Replacement

- NuRAPID layout

- Results

- Conclusion

## Key: Conventional caches spread block over subarrays

+ Splits the "decoding" into the address decoder and muxes at the output of the subarrays

e.g., <u>5-to-1 decoder + 2 2-to-1 muxes</u> better than <u>10-to-1 decoder</u>

?? 9-to-1 decoder ??

+ more flexibility to deal with defects

+ more tolerant to transient errors

- Non-uniform cache: can spread over only one d-group
  - So all bits in a block have same access time

**Small d-groups** (e.g., 64KB of 4 16-KB subarrays)

- Fine granularity of access times
- Blocks spread over few subarrays

**Large d-groups** (e.g., 2 MB of 128 16-KB subarrays)

- Coarse granularity of access times
- Blocks spread over many subarrays

Large d-groups superior for spreading data

# Outline

- Overview

- NuRAPID Mapping and Placement

- NuRAPID Replacement

- NuRAPID layout

- Results

- Conclusion

- 64 KB, 2-way L1s.  8 MSHRs on d-cache

NuRAPID: 8 MB, 8-way, 1-port, no banking
- 4 d-groups (14-, 18-, 36-, 44- cycles)
- 8 d-groups (12-, 19-, 20-, . . . 49- cycles) shown in paper

Compare to:
- **BASE:**
  - 1 MB, 8-way L2 (11-cycles) + 8-MB, 8-way L3 (43-cycles)
- 8 MB, 16-way D-NUCA (4 – 31 cycles)
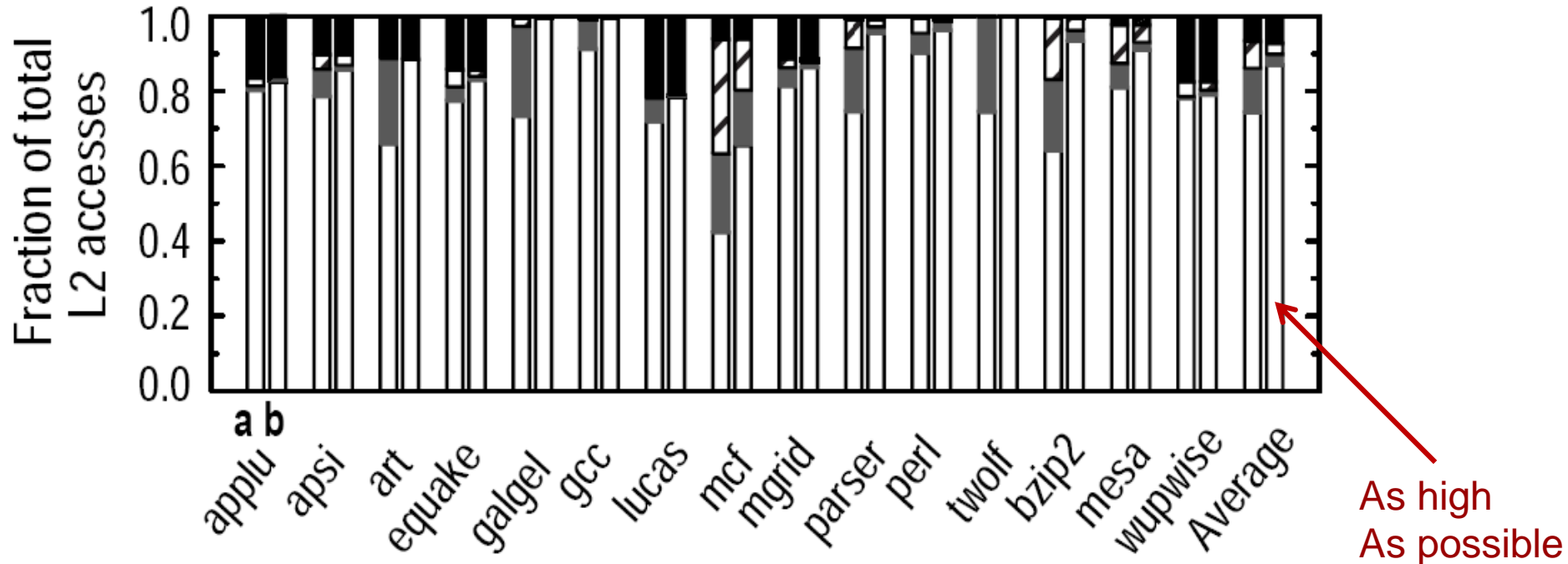  - Multi-banked, infinite-bandwidth interconnect

FIGURE 4: Distribution of group accesses for set-associative and distance-associative placement.

**a**: D-NUCA    **b**: 4-d-group NuRAPID    **c**: 8-d-group NuRAPID

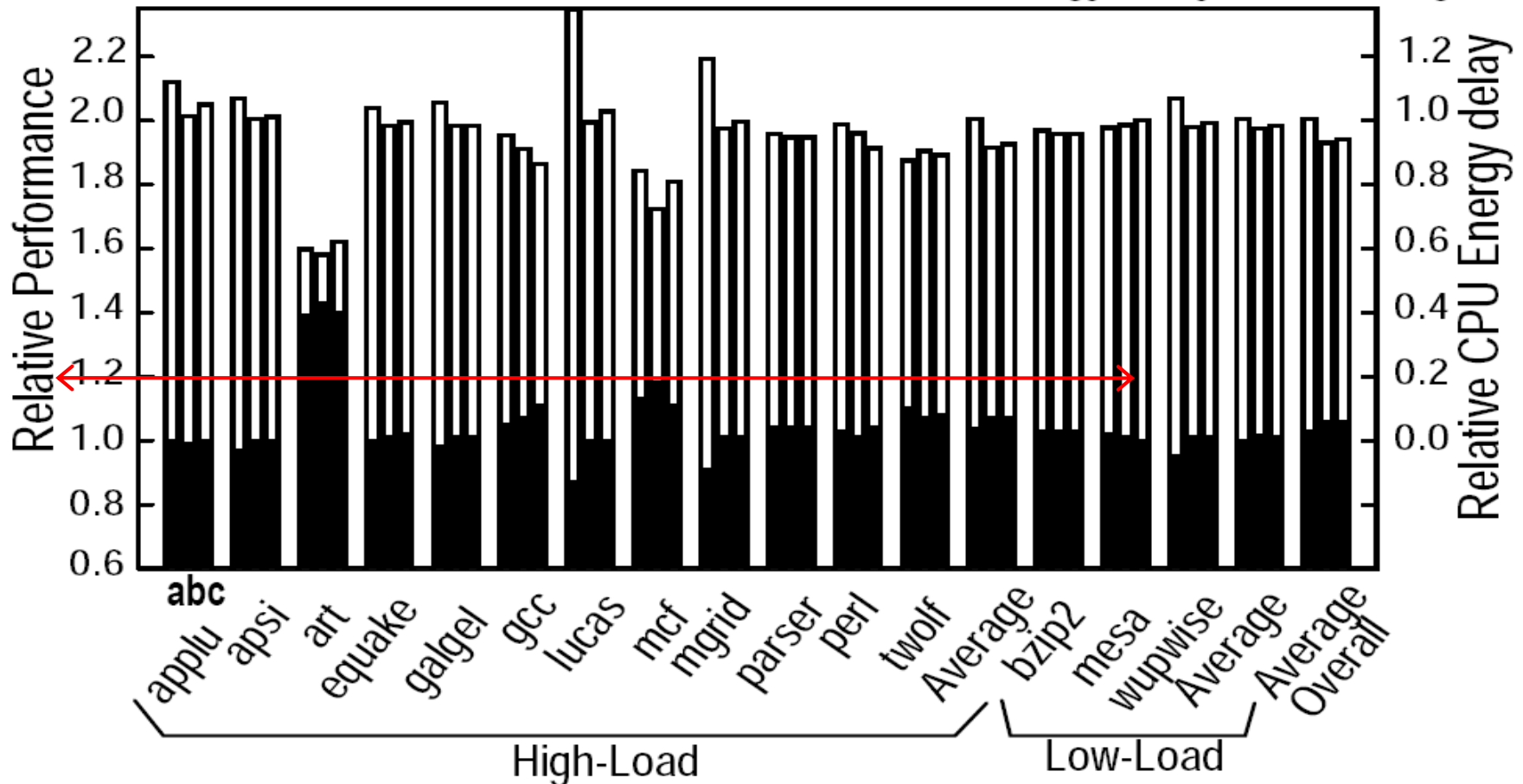■ Relative Performance (scale on left)  □ Relative energy delay (scale on right)

**FIGURE 9:  Performance and energy-delay comparison of NuRAPID and D-NUCA.**
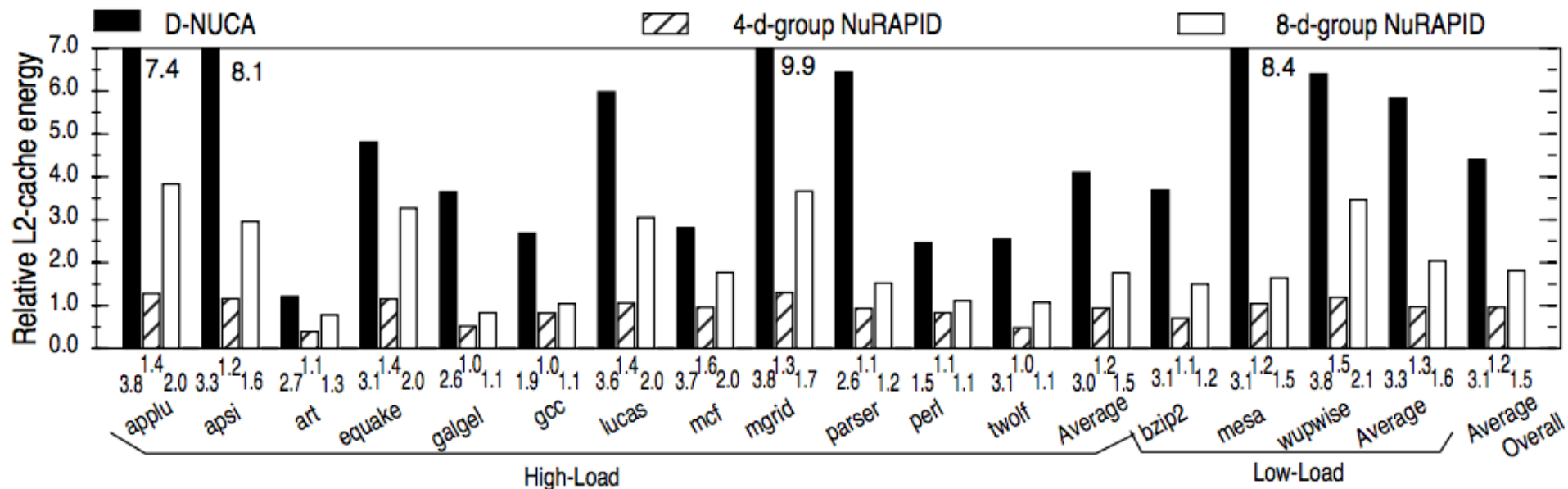
**3.0% better than D-NUCA and up to 15% better**

FIGURE 10: Cache energy comparison of NuRAPID and D-NUCA.

**Energy effects are much more significant**

- NuRAPID

  - leverage seq. tag-data

  - flexible placement, replacement for non-uniform cache

  - Achieves *7% overall processor E-D* savings

    over conventional cache, D-NUCA
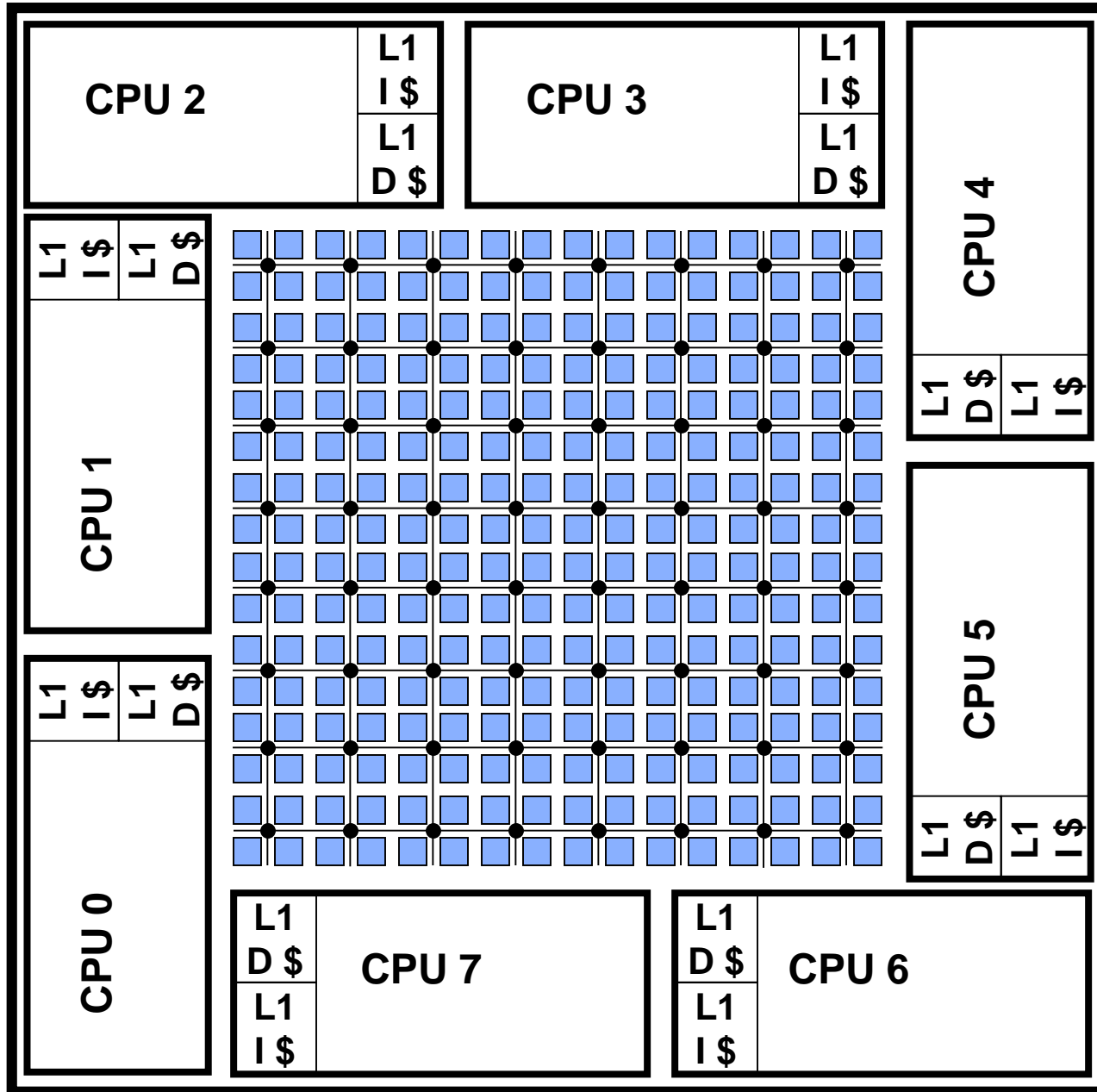
  - Reduces L2 energy by 77% over D-NUCA

NuRAPID an important design for wire-delay dominated caches

# **Managing Wire Delay in Large CMP Caches**

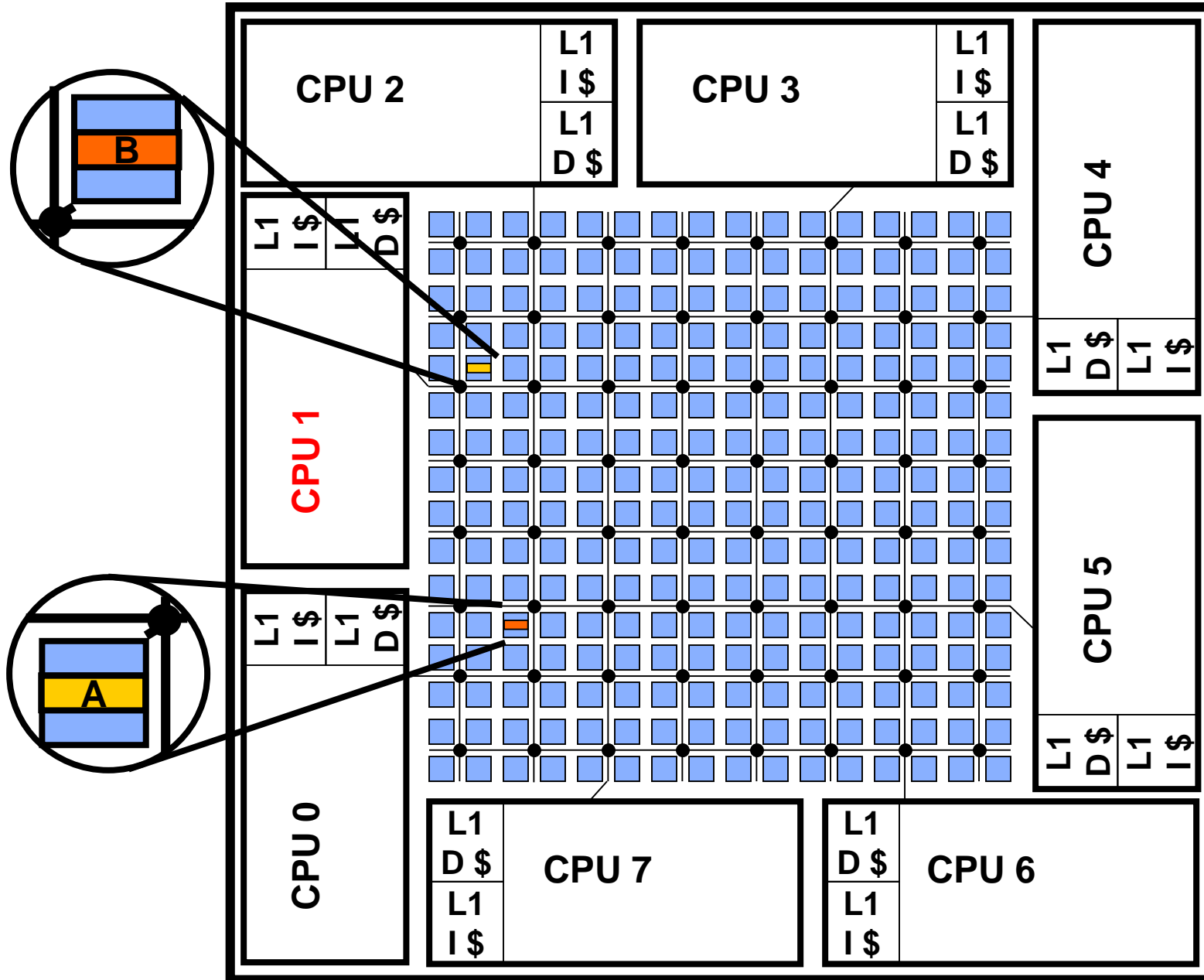**Bradford M. Beckmann and David A. Wood**

Multifacet Project
University of Wisconsin-Madison
MICRO 2004

- Managing wire delay in shared CMP caches
- Three techniques extended to CMPs
    1. **On-chip Strided Prefetching** (not in talk – see paper)
        – Scientific workloads: **10%** average reduction
        – Commercial workloads: **3%** average reduction
    2. **Cache Block Migration** (e.g. D-NUCA)
        – Block sharing limits average reduction to **3%**
        – Dependence on difficult to implement smart search
    3. **On-chip Transmission Lines** (e.g. TLC)
        – Reduce runtime by **8%** on average
        – Bandwidth contention accounts for **26%** of L2 hit latency
- **Combining techniques**
    + Potentially alleviates isolated deficiencies
        – Up to **19%** reduction vs. baseline
    – Implementation complexity

- Global interconnect and CMP trends

- **Latency Management Techniques**

- Evaluation

  - Methodology

  - **Block Migration**: CMP-DNUCA

  - **Transmission Lines**: CMP-TLC
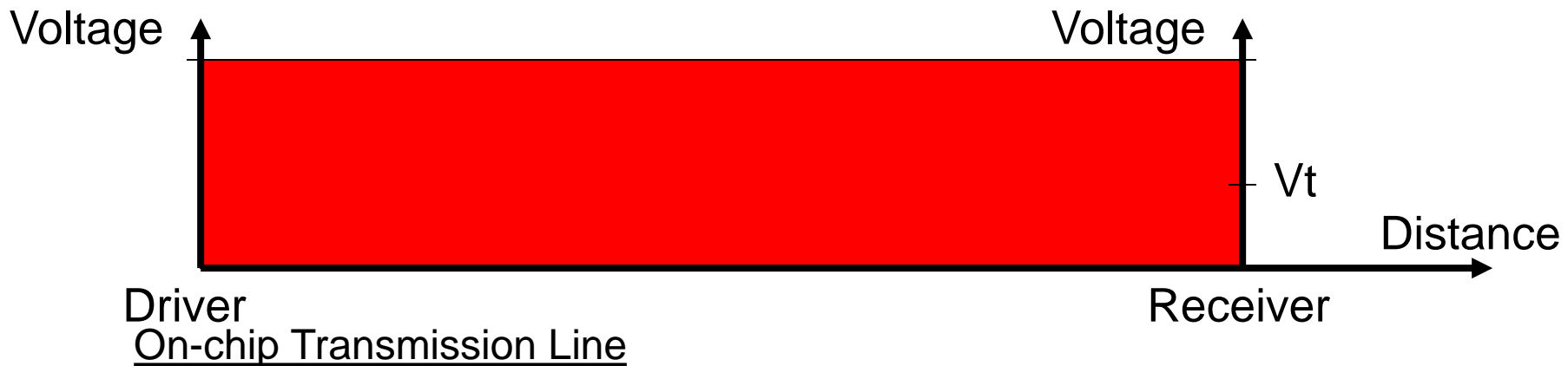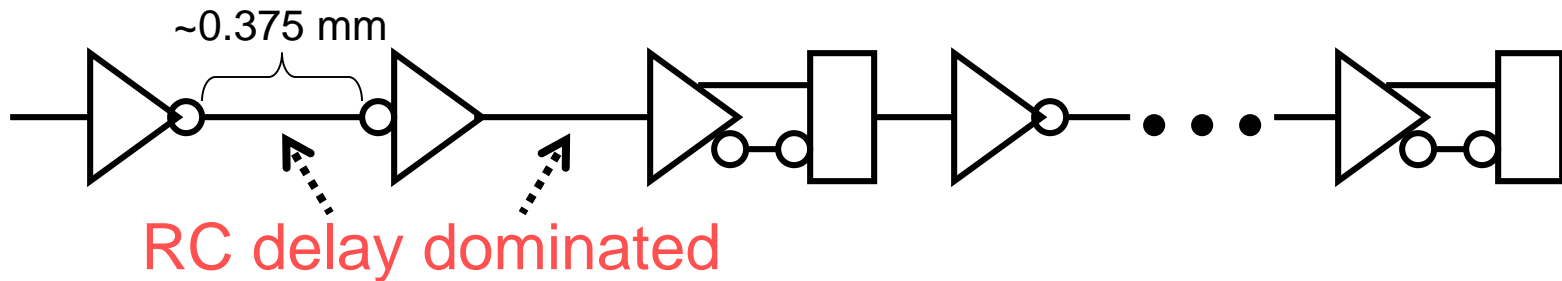
  - **Combination**: CMP-Hybrid

# On-chip Transmission Lines

- Similar to contemporary off-chip communication
- Provides a different latency / bandwidth tradeoff
- Wires behave more "transmission-line" like as frequency increases
  - Utilize transmission line qualities to our advantage
  - No repeaters – route directly over large structures
  - **~10x lower latency across long distances**
- Limitations
  - Requires thick wires and dielectric spacing
  - Increases manufacturing cost
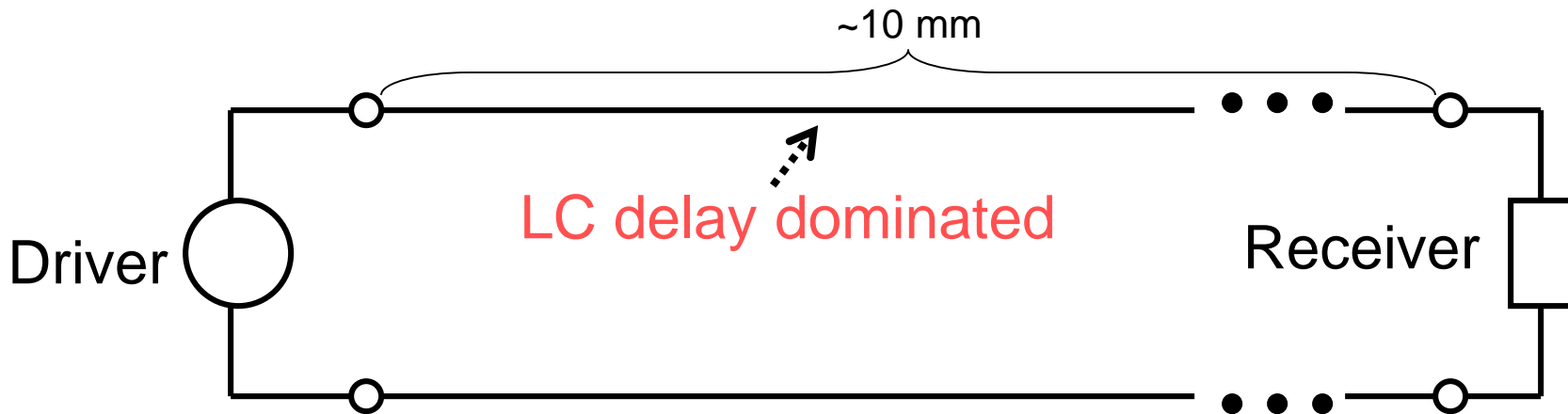- See "TLC: Transmission Line Caches" Beckman, Wood, MICRO'03

Conventional Global RC Wire



Voltage                                    Voltage

                                                    Vt

                                              Distance

Driver                                    Receiver

On-chip Transmission Line



Voltage                                    Voltage

                                                    Vt

                                              Distance

Driver                                    Receiver

MICRO '03 - TLC: Transmission Line Caches

Conventional Global RC Wire

~0.375 mm

RC delay dominated

On-chip Transmission Line

~10 mm

LC delay dominated

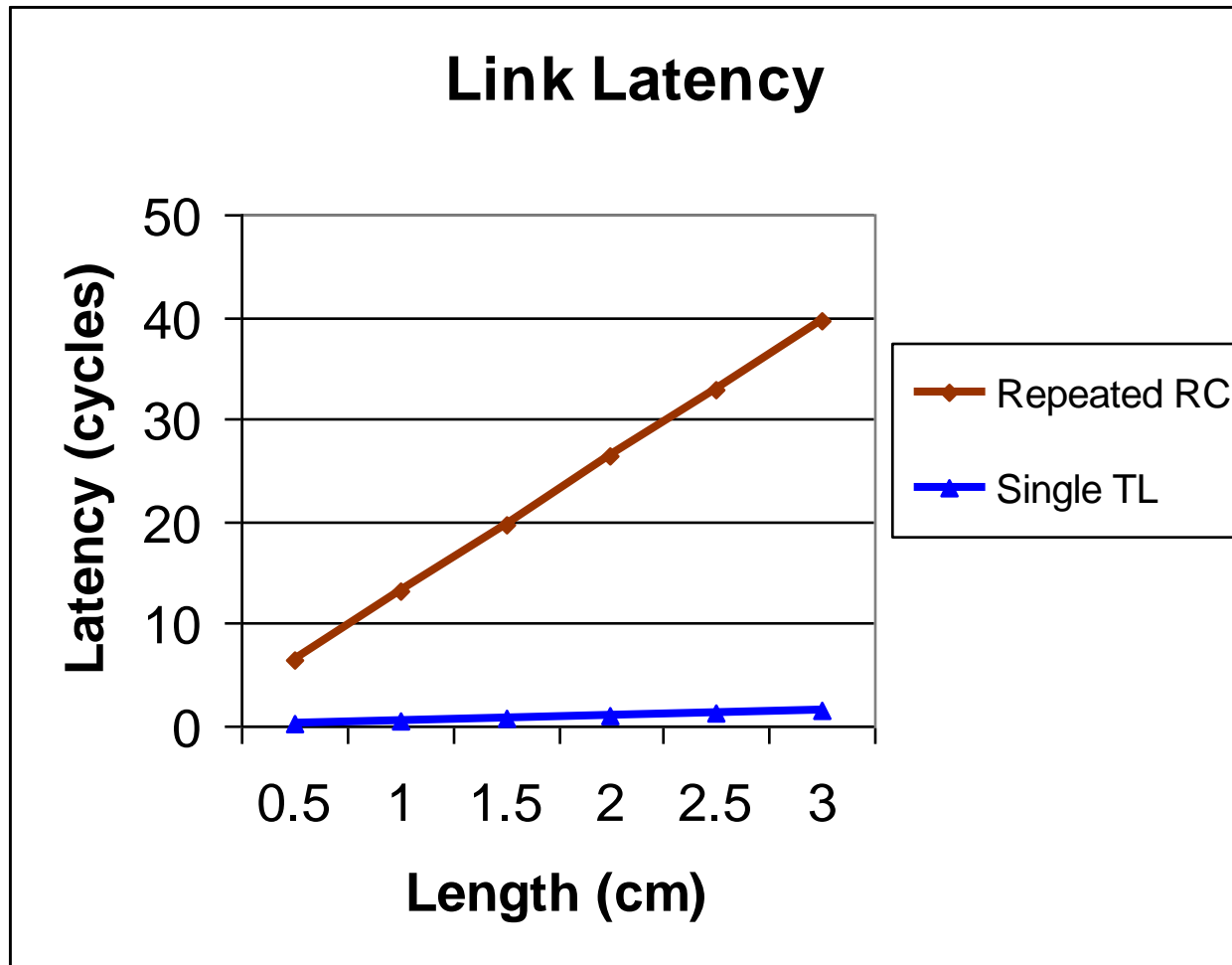Driver

Receiver

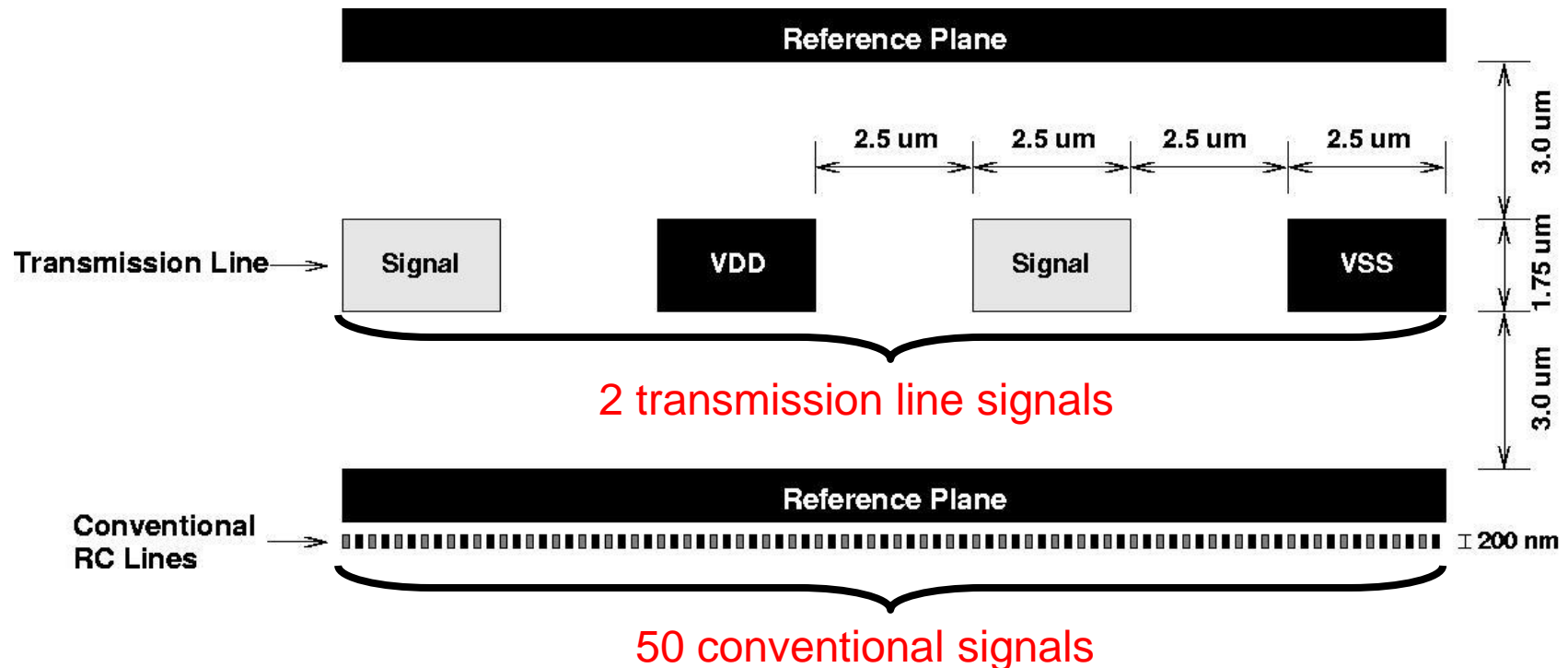MICRO '03 - TLC: Transmission Line
Caches

- Why now? → 2010 technology

  – Relative RC delay ↑

  – Improve latency by 10x or more

- What are their limitations?

  – Require thick wires and dielectric spacing

  – Increase wafer cost

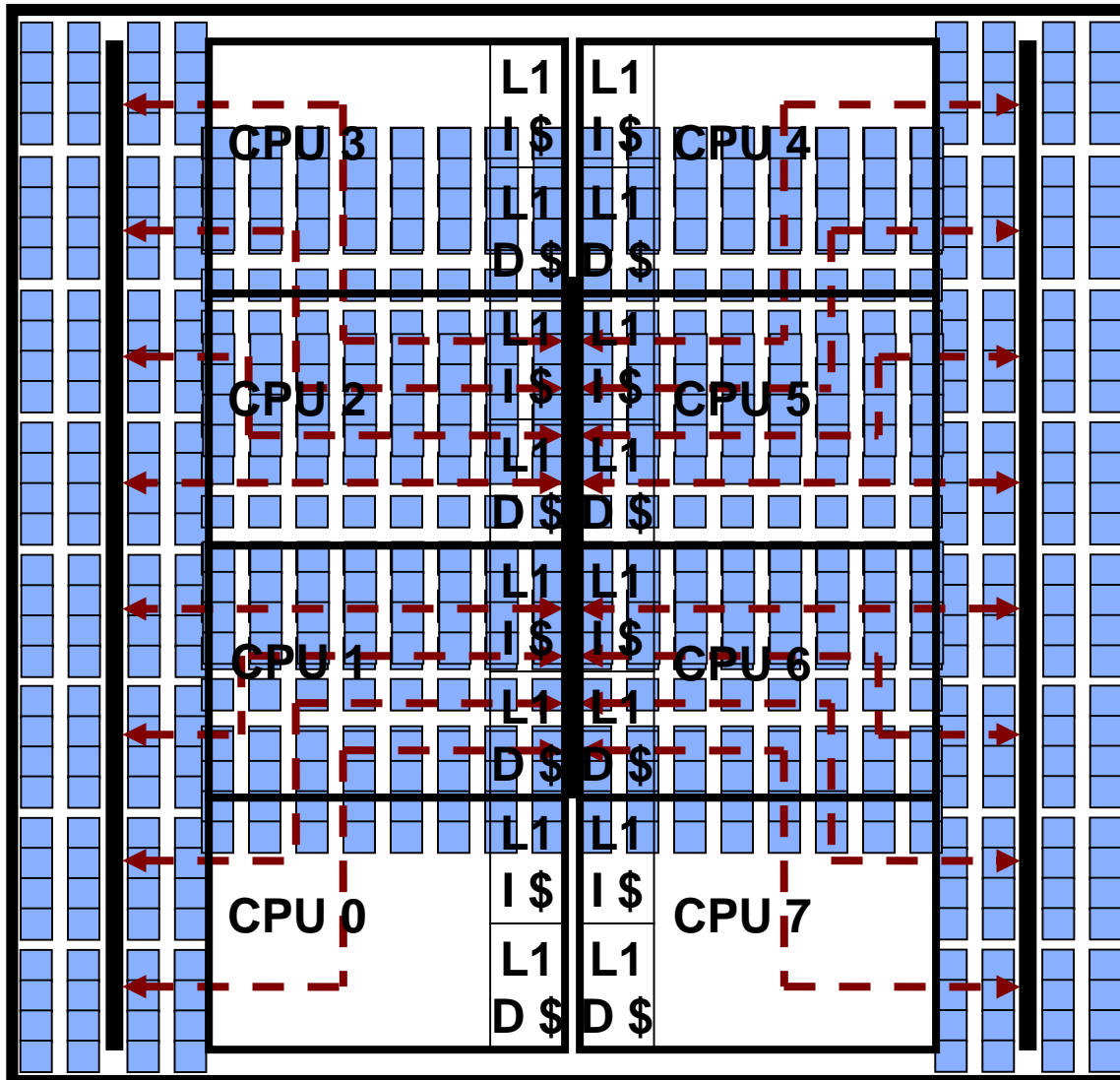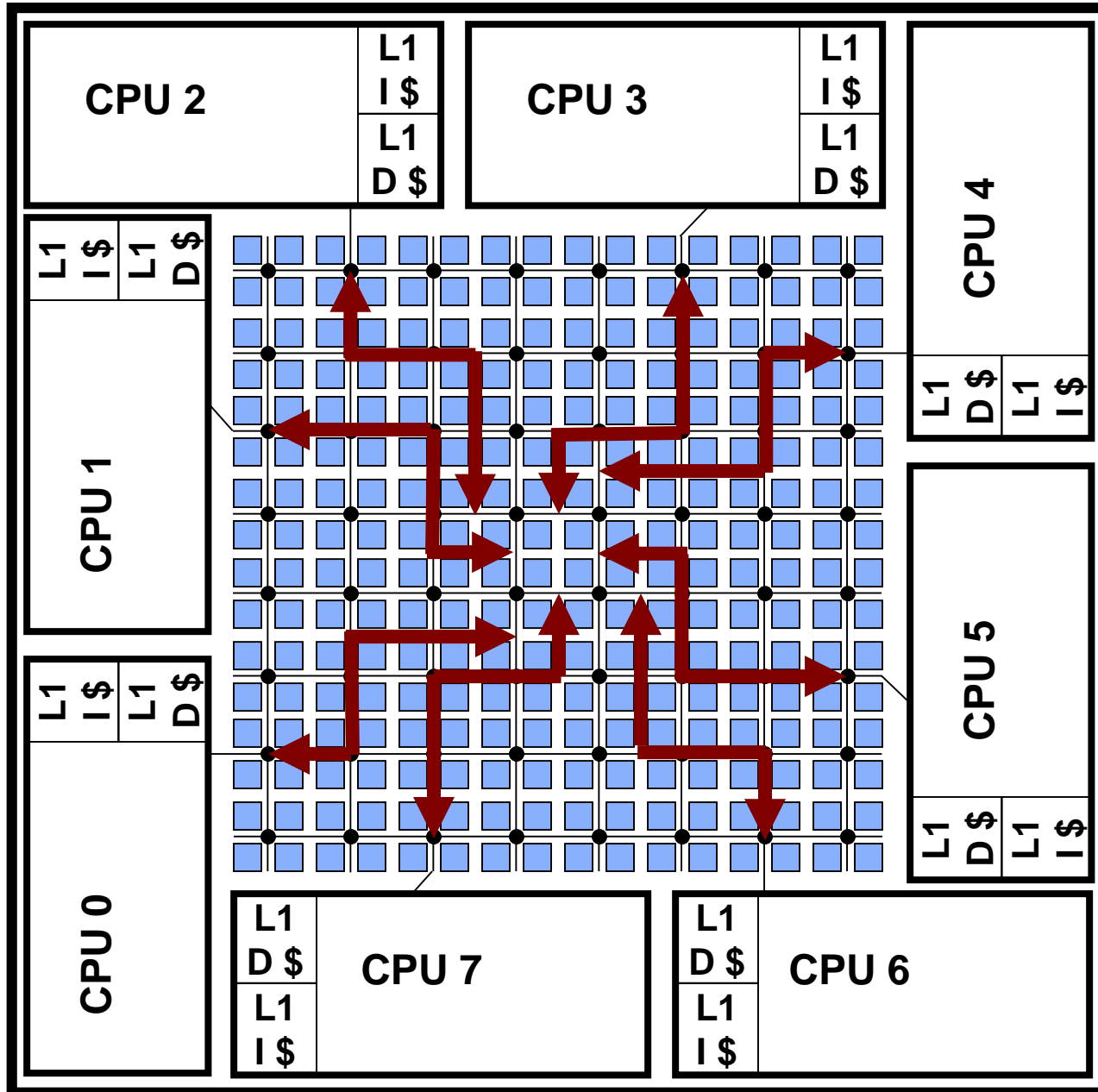Presents a different Latency/Bandwidth Tradeoff

MICRO '03 - TLC: Transmission Line
Caches

**Link Latency**

*Latency (cycles)* vs *Length (cm)*

Legend: Repeated RC, Single TL

MICRO '03 - TLC: Transmission Line Caches

# Bandwidth Comparison



**Key observation**
- Transmission lines – route ***over*** large structures
- Conventional wires – substrate area & vias for repeaters

MICRO '03 - TLC: Transmission Line
Caches

CPU 3

CPU 2

CPU 1

CPU 0

L1 I $

L1 D $

L1 I $

L1 D $

L1 I $

L1 D $

L1 I $

L1 D $

L1 I $

L1 D $

L1 I $

L1 D $

L1 I $

L1 D $

L1 I $

L1 D $

CPU 4

CPU 5

CPU 6

CPU 7

**16 8-byte links**

CPU 2

L1
I $
L1
D $

CPU 3

L1
I $
L1
D $

CPU 4

L1
I $
L1
D $

CPU 1

L1
D $
L1
I $

CPU 5

CPU 0

L1
I $
L1
D $

CPU 7

L1
D $
L1
I $

L1
D $
L1
I $

CPU 6

L1
D $
L1
I $

8
32-byte
links

- Global interconnect and CMP trends

- Latency Management Techniques

- **Evaluation**

  - **Methodology**

  - **Block Migration**: CMP-DNUCA

  - **Transmission Lines**: CMP-TLC

  - **Combination**: CMP-Hybrid

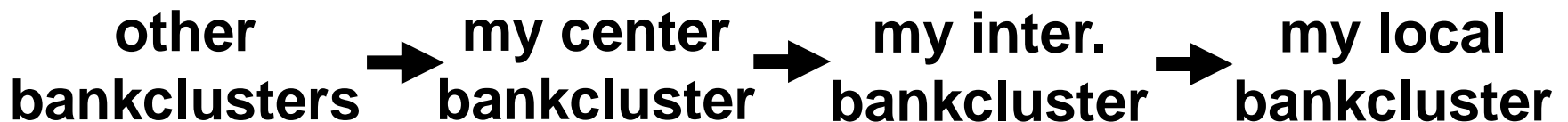Managing Wire Delay in Large CMP Caches

- Full system simulation
  - Simics
  - Timing model extensions
    - Out-of-order processor
    - Memory system

- Workloads
  - **Commercial**
    - apache, **jbb**, **otlp**, zeus
  - **Scientific**
    - *Splash*: barnes & **ocean**
    - *SpecOMP*: **apsi** & fma3d

Managing Wire Delay in Large CMP Caches

# System Parameters

| Memory System | | Dynamically Scheduled Processor | |
|---|---|---|---|
| L1 I & D caches | 64 KB, 2-way, 3 cycles | **Clock frequency** | **10 GHz** |
| **Unified L2 cache** | **16 MB, 256x64 KB, 16-way, 6 cycle bank access** | Reorder buffer / scheduler | 128 / 64 entries |
| L1 / L2 cache block size | 64 Bytes | Pipeline width | 4-wide fetch & issue |
| Memory latency | 260 cycles | Pipeline stages | 30 |
| Memory bandwidth | 320 GB/s | Direct branch predictor | 3.5 KB YAGS |
| Memory size | 4 GB of DRAM | Return address stack | 64 entries |
| Outstanding memory request / CPU | 16 | Indirect branch predictor | 256 entries (cascaded) |

Beckmann & Wood    Managing Wire Delay in Large CMP Caches

- Global interconnect and CMP trends

- Latency Management Techniques

- **Evaluation**

  - Methodology

  - **Block Migration: CMP-DNUCA**

  - **Transmission Lines**: CMP-TLC

  - **Combination**: CMP-Hybrid

Managing Wire Delay in Large CMP Caches

# Hit Distribution: Grayscale Shading



Managing Wire Delay in Large CMP Caches

- Migration policy
  - **Gradual** movement
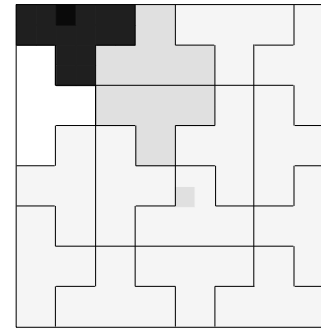  - Increases local hits and reduces distant hits

**other bankclusters** → **my center bankcluster** → **my inter. bankcluster** → **my local bankcluster**
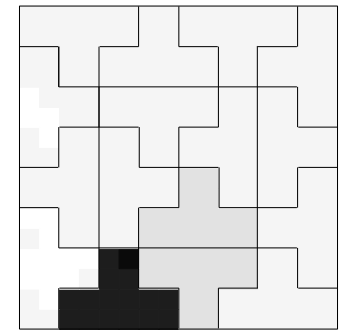
Managing Wire Delay in Large CMP Caches
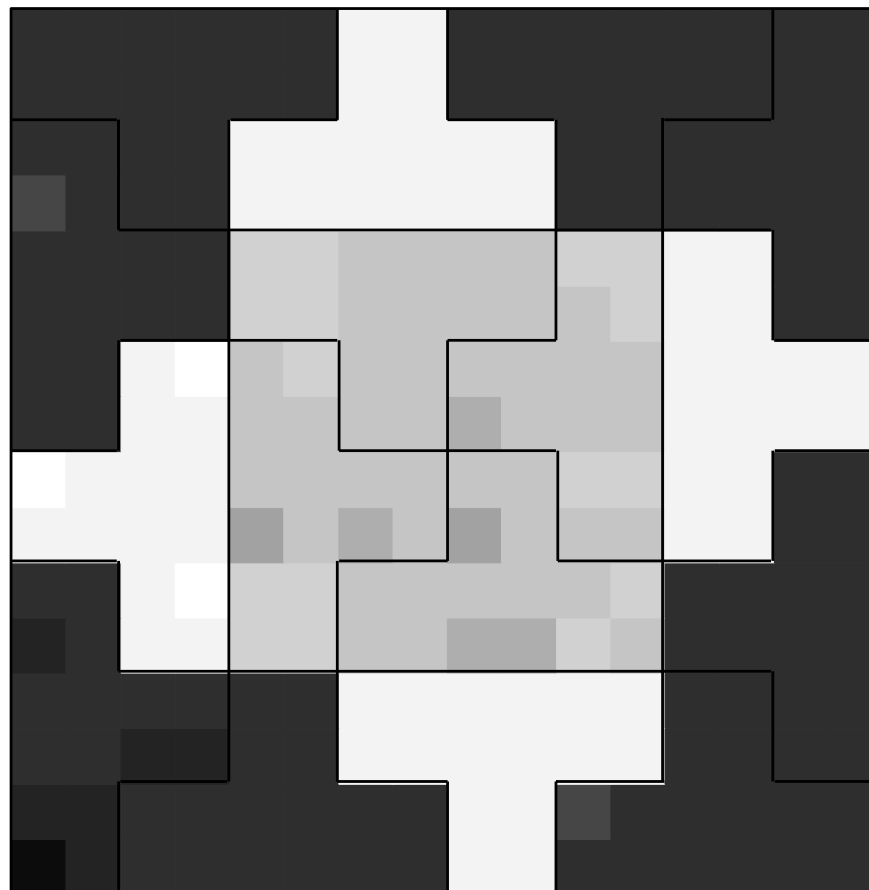
CPU 0　CPU 1　CPU 2　CPU 3
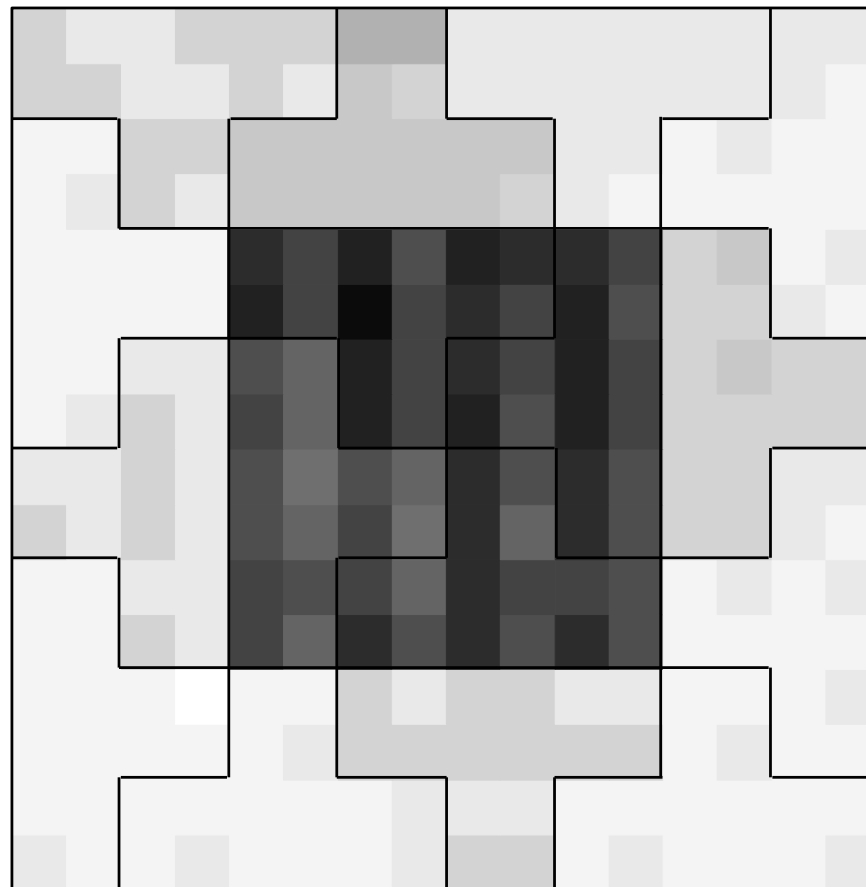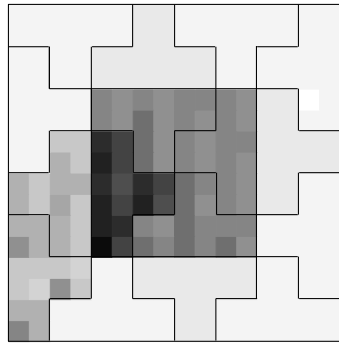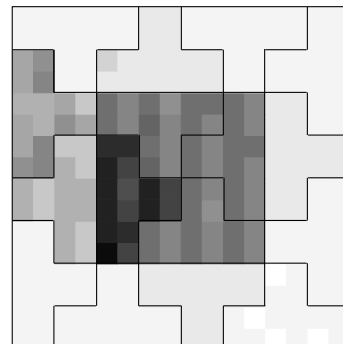
CPU 4　CPU 5　CPU 6　CPU 7

Managing Wire Delay in Large CMP Caches

**Block migration successfully <span style="color:red">separates</span> the data sets**

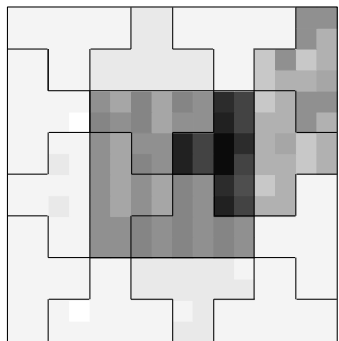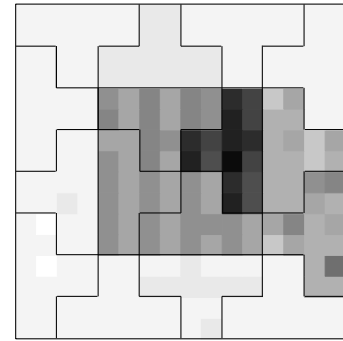**CPU 0**   **CPU 1**   **CPU 2**   **CPU 3**
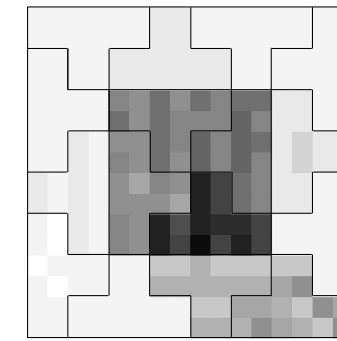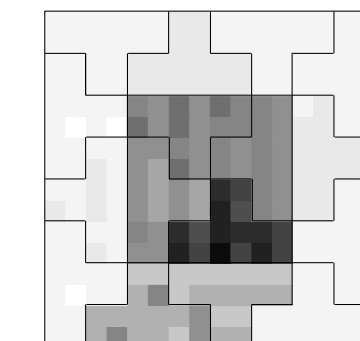
**CPU 4**   **CPU 5**   **CPU 6**   **CPU 7**

**Hit Clustering: Most L2 hits satisfied by the center banks**

- Search policy
  - Uniprocessor DNUCA solution: partial tags
    - Quick summary of the L2 tag state at the CPU
    - **No known practical implementation for CMPs**
      - Size impact of multiple partial tags
      - Coherence between block migrations and partial tag state
  - CMP-DNUCA solution: two-phase search
    - 1st phase: CPU's local, inter., & 4 center banks
    - 2nd phase: remaining 10 banks
    - **Slow 2nd phase hits and L2 misses**

Managing Wire Delay in Large CMP Caches

Managing Wire Delay in Large CMP Caches

- Limited success
  - Ocean successfully splits
    - Regular scientific workload – **little sharing**
  - OLTP congregates in the center
    - Commercial workload – **significant sharing**
- Smart search mechanism
  - Necessary for performance improvement
  - **No known implementations**
  - Upper bound – perfect search

Managing Wire Delay in Large CMP Caches

- Global interconnect and CMP trends

- Latency Management Techniques

- **Evaluation**

    - Methodology

    - **Block Migration**: CMP-DNUCA

    - **Transmission Lines: CMP-TLC**

    - **Combination: CMP-Hybrid**

Managing Wire Delay in Large CMP Caches

Managing Wire Delay in Large CMP Caches

Transmission lines improve **L2 hit** and **L2 miss** latency

Managing Wire Delay in Large CMP Caches

- Individual Latency Management Techniques

  – Strided Prefetching: **subset of misses**

  – Cache Block Migration: **sharing impedes migration**

  – On-chip Transmission Lines: **limited bandwidth**

- Combination: CMP-Hybrid

  – Potentially alleviates bottlenecks

  – Disadvantages

    - Relies on smart-search mechanism
    - Manufacturing cost of transmission lines

Managing Wire Delay in Large CMP Caches

- Initial NUCA designs → Uniprocessors
  - NUCA:
    - Centralized Partial Tag Array
  - NuRAPID:
    - Decouples Tag and Data Placement
    - More overhead
  - L-NUCA
    - Fine-Grain NUCA close to the core
  - Beckman & Wood:
    - Move Data Close to User
    - Two-Phase Multicast Search
    - Gradual Migration
    - Scientific: Data mostly "private" → move close / fast
    - Commercial: Data mostly "shared" → moves in the center / "slow"

- **Beckman & Wood:**
  - Move Data Close to User
  - Two-Phase Multicast Search
  - Gradual Migration
  - Scientific: Data mostly "private" $\rightarrow$ move close / fast
  - Commercial: Data mostly "shared" $\rightarrow$ moves in the center / "slow"

- **CMP-NuRapid:**
  - Per core, L2 tag array
    - Area overhead
    - Tag coherence