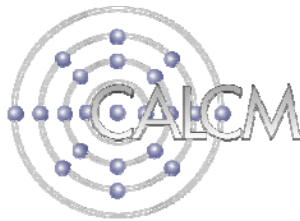


# Reactive NUCA: Near-Optimal Block Placement and Replication in Distributed Caches

Nikos Hardavellas

Michael Ferdman, Babak Falsafi, Anastasia Ailamaki  
Carnegie Mellon and EPFL



**Databases**  
@ Carnegie Mellon



CALCM Computer Architecture Lab
Carnegie Mellon

## Data Placement in Distributed Caches

cache slice

core	core	core	core	core	core	core	core
L2	L2	L2	L2	L2	L2	L2	L2
core	core	core	core	core	core	core	core
L2	L2	L2	L2	L2	L2	L2	L2
core	core	core	core	core	core	core	core
L2	L2	L2	L2	L2	L2	L2	L2
core	core	core	core	core	core	core	core
L2	L2	L2	L2	L2	L2	L2	L2

➡ Data placement determines performance

➡ Goal: place data on chip close to where they are used

**Databases**  
@ Carnegie Mellon
2
© 2009 Hardavellas

**CALCM** Computer Architecture Lab
Carnegie Mellon

## Prior Work

- Several proposals for CMP cache management
  - ASR, cooperative caching, victim replication, CMP-NuRapid, D-NUCA
- ...but suffer from shortcomings
  - complex, high-latency lookup/coherence
  - don't scale
  - lower effective cache capacity
  - optimize only for subset of accesses

We need:

➔ Simple, scalable mechanism for fast access to all data

Databases @ Carnegie Mellon
3
© 2009 Hardavellas

**CALCM** Computer Architecture Lab
Carnegie Mellon

## Our Proposal: Reactive NUCA

- Cache accesses can be classified at run-time
  - Each class amenable to different placement
- Per-class block placement
  - Simple, scalable, transparent
  - No need for HW coherence mechanisms at LLC
  - Avg. speedup of 6% & 14% over shared & private
    - Up to 32% speedup
    - -5% on avg. from ideal cache organization
- Rotational Interleaving
  - Data replication and fast single-probe lookup

Databases @ Carnegie Mellon
4
© 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

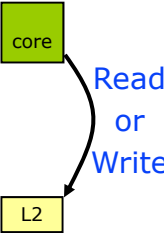
## Outline

- Introduction
- Access Classification and Block Placement
- Reactive NUCA Mechanisms
- Evaluation
- Conclusion

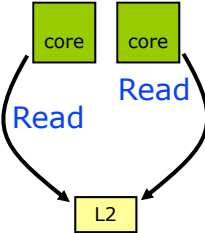
Databases @ Carnegie Mellon 5 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

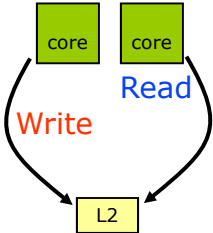
## Terminology: Data Types



Private



Shared  
Read-Only



Shared  
Read-Write

Databases @ Carnegie Mellon 6 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Conventional Multicore Caches

**Shared**

**Private**

**Addr-interleave blocks**

- + High effective capacity
- Slow access

**Each block cached locally**

- + Fast access (local)
- Low capacity (replicas)
- Coherence: via indirection (distributed directory)

➔ We want: high capacity (shared) + fast access (priv.)

**Databases @ Carnegie Mellon** 7 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Where to Place the Data?

- Close to where they are used!
- Accessed by single core: migrate locally
- Accessed by many cores: replicate (?)
  - If read-only, replication is OK
  - If read-write, coherence a problem
    - Low reuse: evenly distribute across sharers

**Databases @ Carnegie Mellon** 8 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Methodology

**Flexus:** Full-system cycle-accurate timing simulation

### Workloads

- OLTP: TPC-C 3.0 100 WH
  - IBM DB2 v8
  - Oracle 10g
- DSS: TPC-H Qry 6, 8, 13
  - IBM DB2 v8
- SPECweb99 on Apache 2.0
- Multiprogrammed: Spec2K
- Scientific: em3d

### Model Parameters

- Tiled, LLC = L2
- Server/Scientific wrkld.
  - 16-cores, 1MB/core
- Multi-programmed wrkld.
  - 8-cores, 3MB/core
- OoO, 2GHz, 96-entry ROB
- Folded 2D-torus
  - 2-cycle router
  - 1-cycle link
- 45ns memory

Databases @ Carnegie Mellon
9
© 2009 Hardavellas

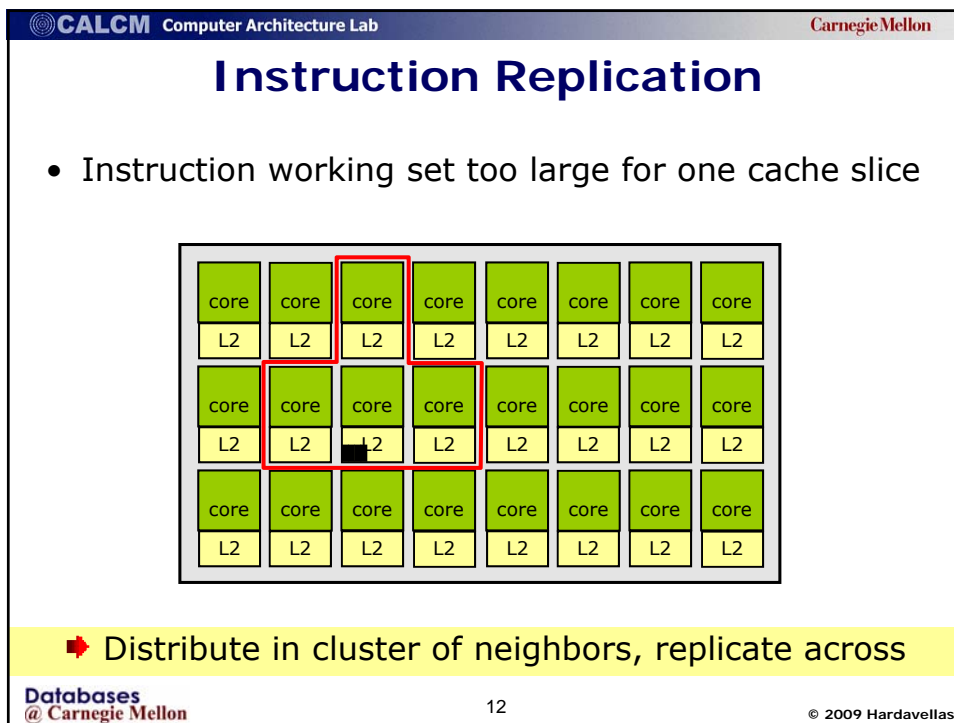
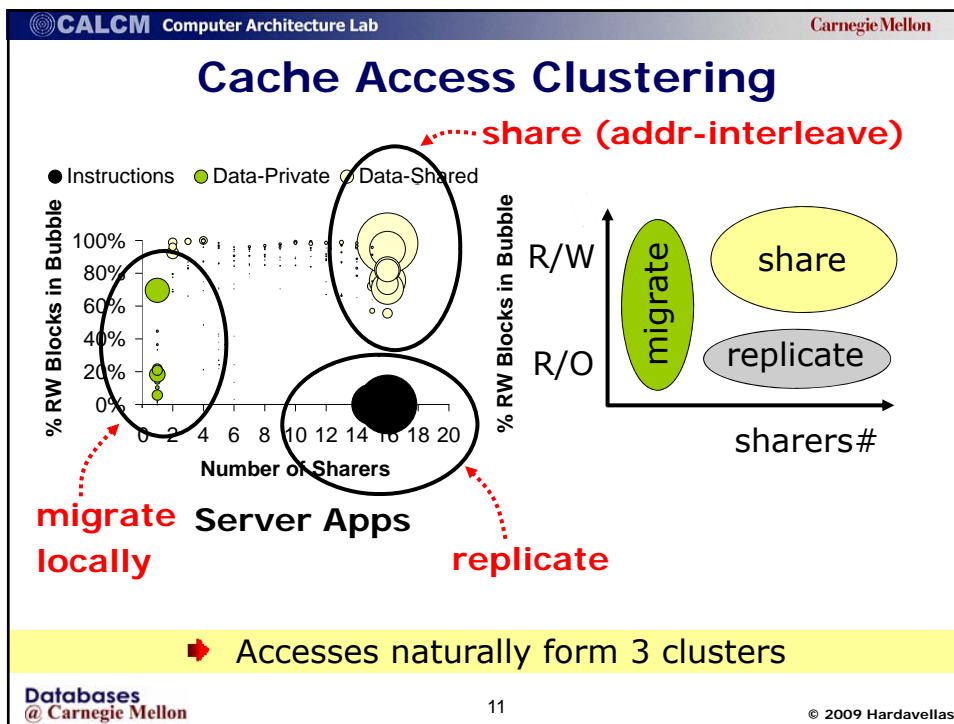
**CALCM** Computer Architecture Lab Carnegie Mellon

## Cache Access Classification Example

- Each bubble: cache blocks shared by x cores
- Size of bubble proportional to % L2 accesses
- y axis: % blocks in bubble that are read-write

● Instructions   ● Data-Private   ○ Data-Shared

Databases @ Carnegie Mellon 10 © 2009 Hardavellas



**CALCM** Computer Architecture Lab Carnegie Mellon

## Outline

- Introduction
- Access Classification and Block Placement
- Reactive NUCA Mechanisms
- Evaluation
- Conclusion

Databases @ Carnegie Mellon 13 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Rotational Interleaving

RID → +1

↓

+log<sub>2</sub>(k)

0	1	2	3	0	1	2	3
2	3	0	1	2	3	0	1
0	1	2	3	0	1	2	3
2	3	0	1	2	3	0	1

size-4 clusters:  
local slice + 3 neighbors

- ◆ Fast access (nearest-neighbor, simple lookup)
- ◆ Balance access latency with capacity constraints
- ◆ Equal capacity pressure at overlapped slices

Databases @ Carnegie Mellon 14 © Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Rotational Interleaving

RID → +1

↓

+log<sub>2</sub>(k)

0	1	2	3	0	1	2	3
2	3	0	1	2	3	0	1
0	1	2	3	0	1	2	3
2	3	0	1	2	3	0	1

PC: 0xfa480

$$\text{Destination} = (\text{Addr} + \overline{\text{RID}} + 1) \& (n - 1)$$

- ◆ Fast access (nearest-neighbor, simple lookup)
- ◆ Balance access latency with capacity constraints
- ◆ Equal capacity pressure at overlapped slices

Databases @ Carnegie Mellon
15
© Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Rotational Interleaving

RID → +1

↓

+log<sub>2</sub>(k)

0	1	2	3	0	1	2	3
2	3	0	1	2	3	0	1
0	1	2	3	0	1	2	3
2	3	0	1	2	3	0	1

each slice caches the same blocks on behalf of any cluster

- ◆ Fast access (nearest-neighbor, simple lookup)
- ◆ Balance access latency with capacity constraints
- ◆ Equal capacity pressure at overlapped slices

Databases @ Carnegie Mellon
16
© Hardavellas



**CALCM** Computer Architecture Lab Carnegie Mellon

## Classification Mechanisms

- Instructions classification: all accesses from L1-I
- Per-page classification for data: at TLB miss
  - Utilize OS page table & TLB for book-keeping info
  - Page classification is accurate (<0.5% error)

On 1<sup>st</sup> access

Core i Ld A  
L2 TLB Miss

OS  
A: Private to "i"

On access by another core

Ld A Core j  
L2 TLB Miss

OS  
~~A: Private to "i"~~  
A: Shared

Databases @ Carnegie Mellon 17 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Coherence: No Need for HW Mechanisms at LLC

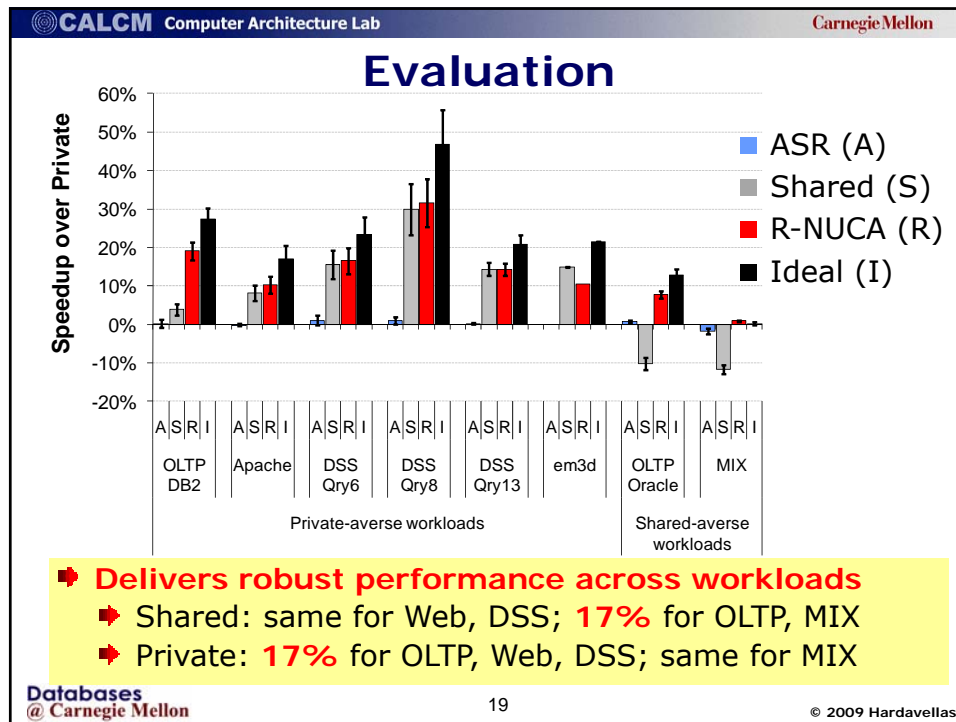
- Reactive NUCA placement guarantee
  - Each R/W datum in unique & known location

Shared data: addr-interleave

Private data: local slice

➔ Fast access, eliminates HW overhead


Databases @ Carnegie Mellon 18 © 2009 Hardavellas



- CALCM Computer Architecture Lab** **Carnegie Mellon**
- ## Conclusions
- Reactive NUCA: near-optimal block placement and replication in distributed caches**
- Cache accesses can be classified at run-time
    - Each class amenable to different placement
  - Reactive NUCA: placement of each class
    - Simple, scalable, low-overhead, transparent
    - Obviates HW coherence mechanisms for LLC
  - Rotational Interleaving
    - Replication + fast lookup (neighbors, single probe)
  - Robust performance across server workloads
    - Near-optimal placement (-5% avg. from ideal)
- Databases @ Carnegie Mellon** 20 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Questions ?



### Impetus Group

Computer Architecture Lab (CALCM)  
Carnegie Mellon University  
[www.ece.cmu.edu/CALCM/](http://www.ece.cmu.edu/CALCM/)

Flexus full-system simulator available at  
[www.ece.cmu.edu/~simflex/](http://www.ece.cmu.edu/~simflex/)

**Databases @ Carnegie Mellon** 21 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Backup Slides

### Classification and Lookup

**Databases @ Carnegie Mellon** 22 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Classification Granularity: OS Page

- Per-block classification: high area/latency overhead
- Per-page classification (utilize OS page table)
  - Core accesses it for every access anyway (TLB)
  - Page classification is accurate (<0.5% error)

TLB entry: 

P/S	vpage	ppage
-----	-------	-------

} 1 bit

Page table entry: 

P/S/I	L2 id	vpage	ppage
-------	-------	-------	-------

} 2 bits    } log(n)

➤ Page granularity allows simple + practical HW

**Databases @ Carnegie Mellon** 23 © 2009 Hardavellas

**CALCM** Computer Architecture Lab Carnegie Mellon

## Data Class Bookkeeping

- **private data:** place in local L2 slice
 

Page table entry: 

P	L2 id	vpage	ppage
---	-------	-------	-------

TLB entry: 

P	vpage	ppage
---	-------	-------
- **shared data:** place in aggregate L2 (addr interleave)
 

Page table entry: 

S	L2 id	vpage	ppage
---	-------	-------	-------

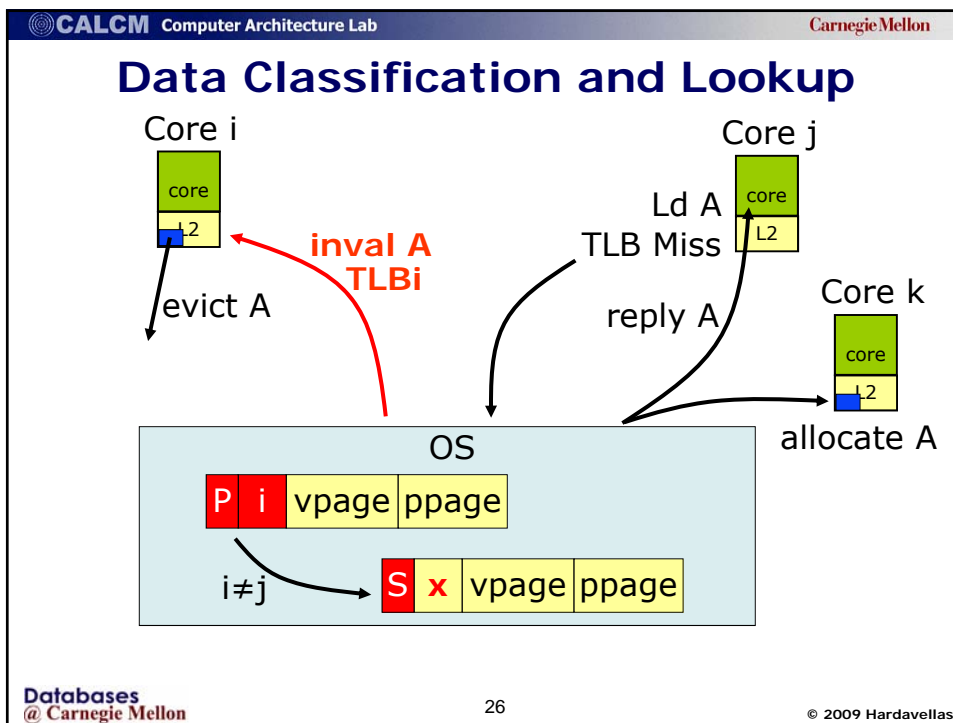
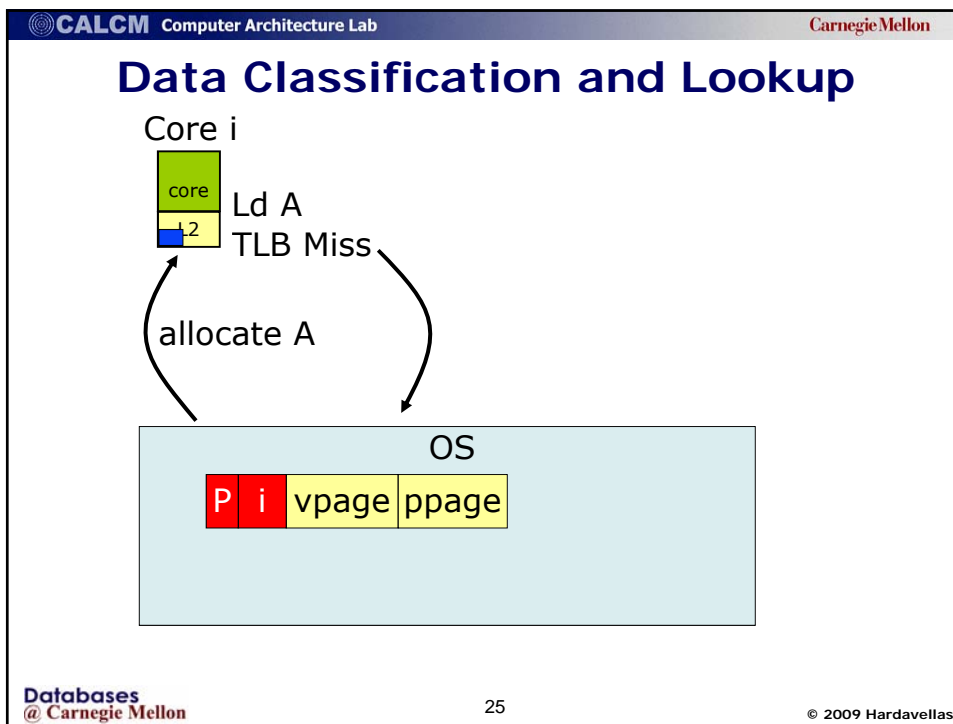
TLB entry: 

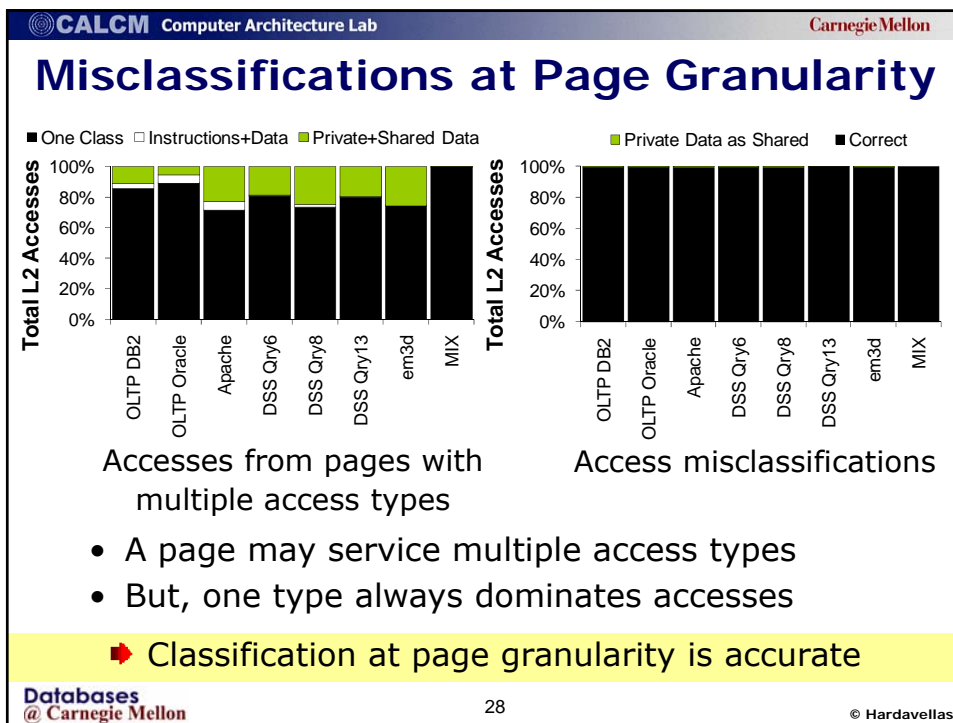
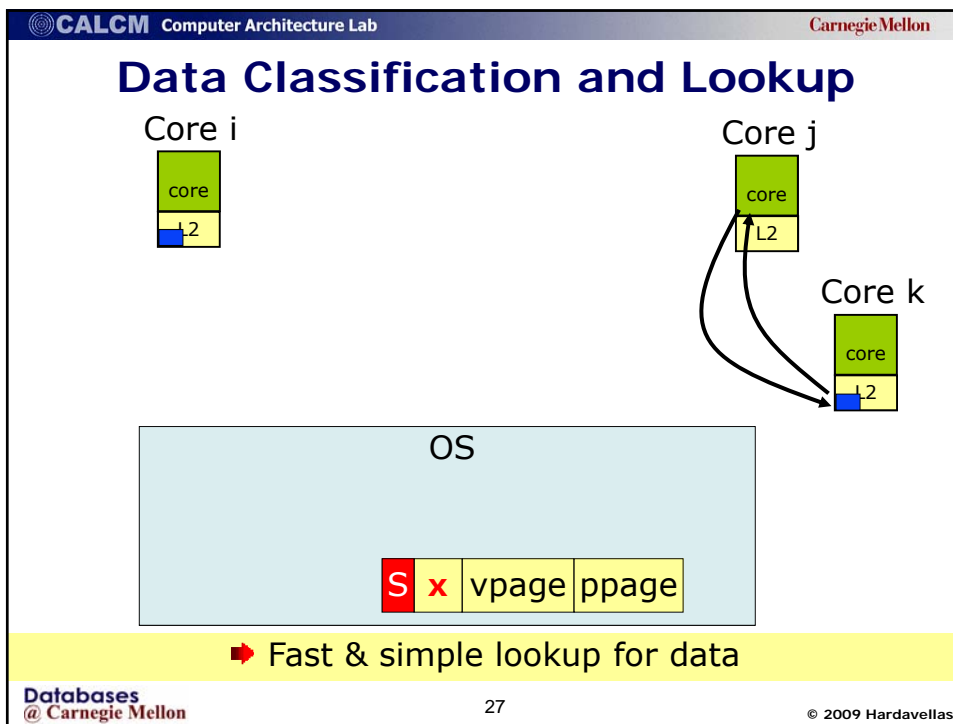
S	vpage	ppage
---	-------	-------

Physical Addr.: 

tag	L2 id	cache index	offset
-----	-------	-------------	--------

**Databases @ Carnegie Mellon** 24 © Hardavellas





CALCM Computer Architecture Lab Carnegie Mellon

# Backup Slides

## Placement

Databases @ Carnegie Mellon 29 © 2009 Hardavellas

CALCM Computer Architecture Lab Carnegie Mellon

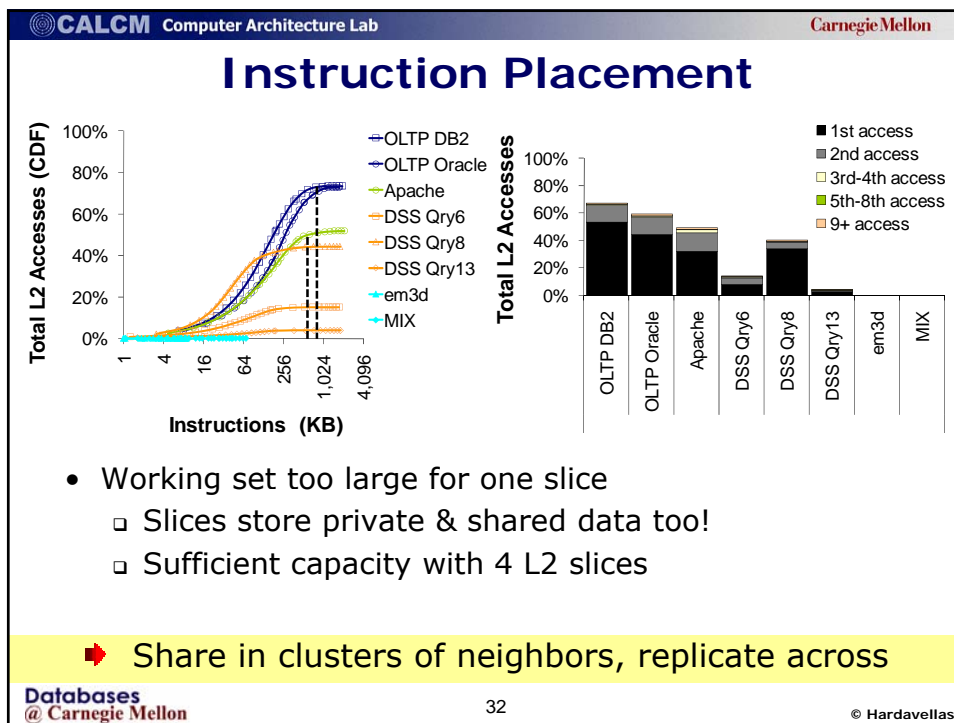
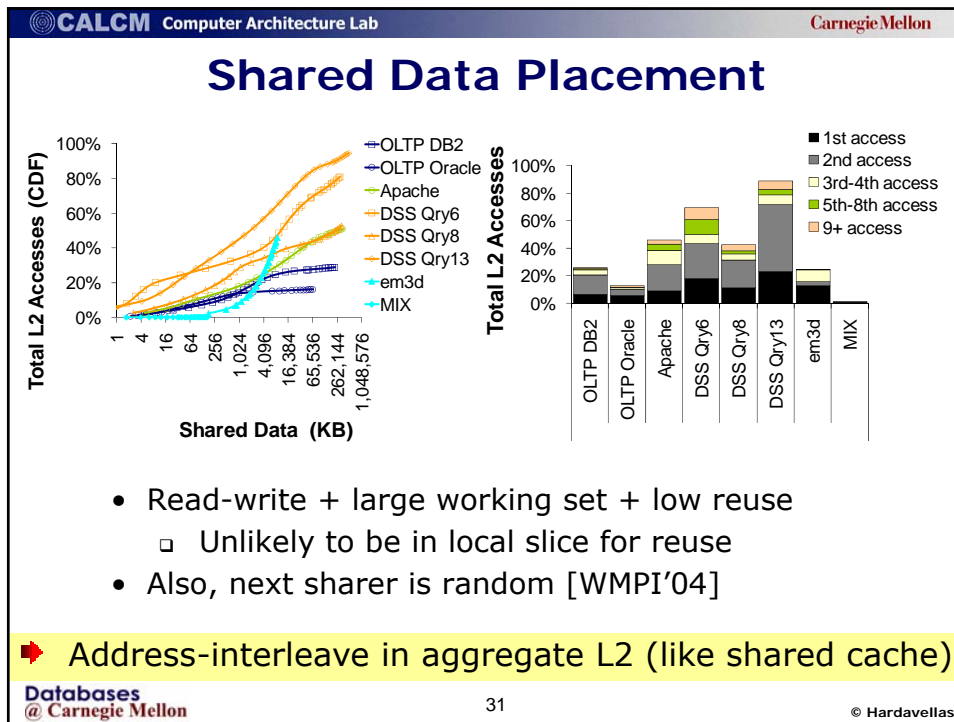
## Private Data Placement

Private Data (KB)	MIX (%)	em3d (%)	DSS Qry13 (%)	DSS Qry8 (%)	DSS Qry6 (%)	Apache (%)	OLTP Oracle (%)	OLTP DB2 (%)
1	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0
64	0	0	0	0	0	0	0	0
256	0	0	0	0	0	0	0	0
1,024	0	0	0	0	0	0	0	0
4,096	0	0	0	0	0	0	0	0
16,384	0	0	0	0	0	0	0	0
65,536	0	0	0	0	0	0	0	0
262,144	0	0	0	0	0	0	0	0
1,048,576	100	100	100	100	100	100	100	100

- Spill to neighbors if working set too large?
  - NO!!! Each core runs similar threads

➡ Store in local L2 slice (like in private cache)

Databases @ Carnegie Mellon 30 © Hardavellas





CALCM Computer Architecture Lab Carnegie Mellon

# Backup Slides

## Detailed Evaluation

Databases @ Carnegie Mellon 33 © 2009 Hardavellas

CALCM Computer Architecture Lab Carnegie Mellon

# Cache Accesses Breakdown

Workload	Instructions (%)	Data-Private (%)	Data-Shared-RW (%)	Data-Shared-RO (%)
OLTP DB2	~68	~10	~15	~7
OLTP Oracle	~60	~25	~10	~5
Apache	~50	~5	~30	~15
DSS Qry6	~15	~15	~60	~10
DSS Qry8	~40	~15	~40	~5
DSS Qry13	~5	~5	~85	~5
em3d	~0	~75	~20	~5
MIX	~0	~95	~5	~0

■ Instructions   ■ Data-Private   ■ Data-Shared-RW   ■ Data-Shared-RO

Total L2 Accesses

OLTP DB2   OLTP Oracle   Apache   DSS Qry6   DSS Qry8   DSS Qry13   em3d   MIX

➡ Instr. + shared-RW dominate server workloads  
 ➡ Private dominate scientific/MIX

Databases @ Carnegie Mellon 34 © Hardavellas

