# EFFICIENT AND SCALABLE CACHE COHERENCE FOR MANY-CORE ARCHITECTURES

## Alberto Ros

Research assistant
Computer Engineering Department
Technical University of Valencia
aros@gap.upv.es

Adjunct professor
Computer Engineering Department
University of Murcia
a.ros@ditec.um.es

Uppsala, May 27, 2011

## OUTLINE

**1** INTRODUCTION
- Challenges in many-core computing

**2** CACHE COHERENCE PROTOCOLS
- Direct coherence (DiCo)
- Coherence deactivation
- Synchronous coherence

**3** MEMORY HIERARCHY ORGANIZATION
- Replacement policies for shared caches
- Indexing policies for shared caches
- Impact of NUCA mapping policies on directory scalability

**4** CONCLUSIONS

## OUTLINE

TRENDS

- The increasing number of transistors per chip can be used to obtain more performance.

## TRENDS

- The increasing number of transistors per chip can be used to obtain more performance.

### EXPLOITING ILP

- Very complex core
- Small improvements

## TRENDS

- The increasing number of transistors per chip can be used to obtain more performance.

### EXPLOITING ILP

- Very complex core
- Small improvements

$\Longrightarrow$

### EXPLOITING TLP

- Many simple cores
- Programming effort

# TRENDS

- The increasing number of transistors per chip can be used to obtain more performance.

EXPLOITING ILP
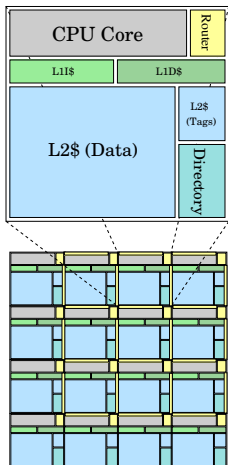- Very complex core
- Small improvements

$\Longrightarrow$

EXPLOITING TLP
- Many simple cores
- Programming effort

- Chip Multiprocessors (CMPs) constitute the new trend for increasing performance.
- Tiled CMPs are a scalable alternative for building CMPs.
  - Designed as arrays of replicated tiles.
  - Cores connected through a direct network.

# TILED-CMPS
ARCHITECTURE



- Each tile contains:
    - A processing core.
    - A private L1 cache (both instruction and data caches).
    - A shared or private L2 cache bank, and a directory.
    - A network interface (router).
- All tiles are connected through a scalable point-to-point interconnection network.

# CACHE COHERENCE PROBLEM

- Most parallel software in the commercial market relies on a shared-memory programming model.
- The presence of private caches requires to keep coherence among data stored in them.
- Solution ⇒ Keep cache coherence in hardware.
- Problem ⇒ Cache coherence protocols introduce extra overhead:
  - In terms of execution time.
  - In terms of area requirements.
  - In terms of power consumption.
  - In terms of designing and verification time.

# CACHE COHERENCE PROBLEM



- Directory-based cache coherence protocols constitute the most scalable alternative.
  - But they have some inefficiencies and constraints:

  1. Scalability of the directory structure.
  2. Indirection to the home node.
  3. Large verification time.

# MEMORY HIERARCHY ORGANIZATION
## SHARED VS. PRIVATE LASTL-LEVEL (L2) CACHE ORGANIZATION



### PRIVATE ORGANIZATION

- ☺ L2 hits have short latencies (local accesses).
- ☹ Blocks potentially replicated in multiple L2 banks.
- ☹ Load balancing problems.

# MEMORY HIERARCHY ORGANIZATION
### SHARED VS. PRIVATE LASTL-LEVEL (L2) CACHE ORGANIZATION



### PRIVATE ORGANIZATION

- ☺ L2 hits have short latencies (local accesses).
- ☹ Blocks potentially replicated in multiple L2 banks.
- ☹ Load balancing problems.

### SHARED ORGANIZATION (NUCA ARCHITECTURE)

- ☺ Better use of the aggregate L2 cache capacity.
- ☹ Long latencies when compared to a private L2 design.
  - The access latency to the L2 depends on where the requested block is mapped.

# SHARED CACHE ORGANIZATION CHALLENGES



- Tiled-CMPs distribute the shared last-level cache among the different tiles (Non Uniform Cache Access or NUCA architecture).
  - The access latency to the last-level cache depends on where the requested block is mapped.
  - Blocks requested by different threads competing for the same resources.

  4. Reduce long access latencies.
  5. Manage conflicting data requests from different threads.

## OUTLINE

## DIRECT COHERENCE
MOTIVATION

### REMEMBER

Directory protocols introduce indirection in the critical path of cache misses.

- This indirection impacts on applications' performance.

# DIRECT COHERENCE
MOTIVATION

---

### REMEMBER

Directory protocols introduce indirection in the critical path of cache misses.

- This indirection impacts on applications' performance.

---

- Token protocols have been proposed to avoid the indirection problem.
    - But they rely on broadcasting requests to all nodes.
    - They are not scalable in terms of network-traffic.
- An ideal protocol should avoid indirection while keeping traffic requirements low.



---

# DIRECT COHERENCE
## INDIRECTION PROBLEM



### CACHE-TO-CACHE TRANSFER IN DIRECTORY-BASED PROTOCOLS

Cache miss        R

# DIRECT COHERENCE
## INDIRECTION PROBLEM



CACHE-TO-CACHE TRANSFER IN DIRECTORY-BASED PROTOCOLS

Cache miss    R

Why?
- To order requests
- To get directory information
- To provide main memory storage

1 GetS

H & D

# DIRECT COHERENCE
## INDIRECTION PROBLEM

# DIRECT COHERENCE
INDIRECTION PROBLEM

# DIRECT COHERENCE
## THE ROLES

DIRECTORY

# DIRECT COHERENCE
## THE ROLES

DIRECTORY                                    DIRECT COHERENCE

# DIRECT COHERENCE
## THE ROLES

DIRECTORY                                      DIRECT COHERENCE

## DIRECT COHERENCE
THE ROLES

# DIRECT COHERENCE
## THE ROLES

DIRECTORY

DIRECT COHERENCE

## DIRECT COHERENCE
CHANGES IN THE STRUCTURE OF TILES

- This distribution of the roles in direct coherence implies changes in the structure of each tile.

# DIRECT COHERENCE
## CHANGES IN THE STRUCTURE OF TILES

- T ........... nplies
c

> **L1D$**: Adds Sharing Information
> - Every owner cache must keep track of the sharers to keep coherence.
> - This field replaces the directory structure.

# DIRECT COHERENCE
CHANGES IN THE STRUCTURE OF TILES

> **L1C\$**: L1 Coherence Cache
> - Each requesting cache stores the identity of the owner for some memory blocks.
> - This information is used to directly send the requests to the corresponding owner cache.

- This ... ies
  cha...

# DIRECT COHERENCE
## CHANGES IN THE STRUCTURE OF TILES

> **L2C$**: L2 Coherence Cache
> ① Each home tile needs to store the identity of the owner cache of each one of its blocks.
> ② This information is accessed when the requestor is not able to locate the owner cache.

- This di…
  chang…

# DIRECT COHERENCE
## BEHAVIOR: CACHE-TO-CACHE READ MISS



DIRECT COHERENCE

R

# DIRECT COHERENCE
## BEHAVIOR: CACHE-TO-CACHE READ MISS

# DIRECT COHERENCE
## BEHAVIOR: CACHE-TO-CACHE READ MISS

# DIRECT COHERENCE
## BEHAVIOR: CACHE-TO-CACHE READ MISS



- The critical path of the miss is reduced from three to two hops.
- The number of coherence messages is halved.
- The waiting time at the home tile is removed.

# DIRECT COHERENCE
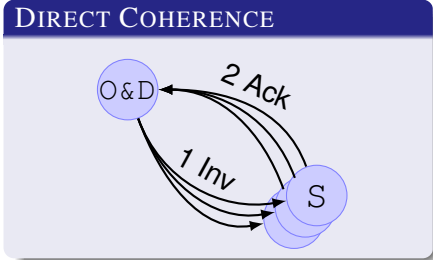BEHAVIOR: UPGRADE IN OWNER
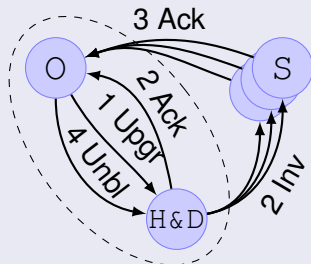


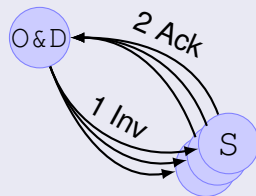DIRECT COHERENCE

O & D

# DIRECT COHERENCE
## BEHAVIOR: UPGRADE IN OWNER

# DIRECT COHERENCE
BEHAVIOR: UPGRADE IN OWNER



DIRECT COHERENCE

# DIRECT COHERENCE
## BEHAVIOR: UPGRADE IN OWNER



- The critical path of the miss is reduced from three to two hops.
- The number of coherence messages is also reduced.

# DIRECT COHERENCE
UPDATING THE L2 COHERENCE CACHE

- The L2C\$ must keep the identity of the current owner cache for each block allocated in any L1 data cache.
  - This information is accessed when the requestor is not able to locate the owner cache.
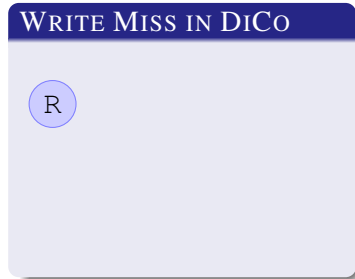- The L2C\$ is notified on every owner change through control messages.

# DIRECT COHERENCE
UPDATING THE L2 COHERENCE CACHE

- The L2C$ must keep the identity of the current owner cache for each block allocated in any L1 data cache.
  - This information is accessed when the requestor is not able to locate the owner cache.

- The L2C$ is notified on every owner change through control messages.

<div style="border:1px solid; padding:10px;">

### WRITE MISS IN DICO

R

</div>

# DIRECT COHERENCE
UPDATING THE L2 COHERENCE CACHE

- The L2C\$ must keep the identity of the current owner cache for each block allocated in any L1 data cache.
  - This information is accessed when the requestor is not able to locate the owner cache.
- The L2C\$ is notified on every owner change through control messages.
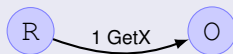
### WRITE MISS IN DICO



R    1 GetX    O

# DIRECT COHERENCE
UPDATING THE L2 COHERENCE CACHE

- The L2C$ must keep the identity of the current owner cache for each block allocated in any L1 data cache.
  - This information is accessed when the requestor is not able to locate the owner cache.

- The L2C$ is notified on every owner change through control messages.

### WRITE MISS IN DICO

# DIRECT COHERENCE
UPDATING THE L2 COHERENCE CACHE

- The L2C\$ must keep the identity of the current owner cache for each block allocated in any L1 data cache.
  - This information is accessed when the requestor is not able to locate the owner cache.

- The L2C\$ is notified on every owner change through control messages.
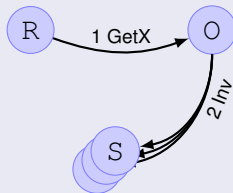
## WRITE MISS IN DICO

# DIRECT COHERENCE
UPDATING THE L2 COHERENCE CACHE

- The L2C\$ must keep the identity of the current owner cache for each block allocated in any L1 data cache.
  - This information is accessed when the requestor is not able to locate the owner cache.

- The L2C\$ is notified on every owner change through control messages.

- These messages should be processed by the L2C\$ in the very same order in which they were generated.
  - To ensure this, the L2C\$ sends an ACK message to the new owner when it receives a change owner message.



### WRITE MISS IN DICO

# DIRECT COHERENCE
## UPDATING THE L2 COHERENCE CACHE

- The L2C\$ must keep the identity of the current owner cache for each block allocated in any L1 data cache.
  - This information is accessed when the requestor is not able to locate the owner cache.

- The L2C\$ is notified on every owner change through control messages.
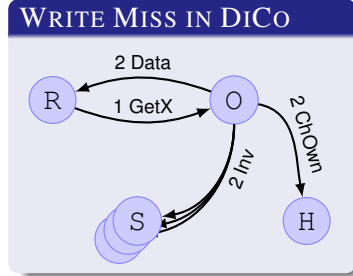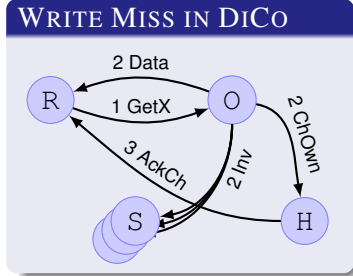
- These messages should be processed by the L2C\$ in the very same order in which they were generated.
  - To ensure this, the L2C\$ sends an ACK message to the new owner when it receives a change owner message.
  - Until this message is not received by the owner node, it could use the block but cannot give the onwership to another cache.

### WRITE MISS IN DICO

# DIRECT COHERENCE
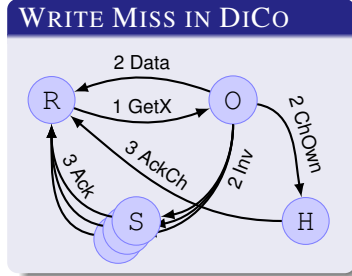UPDATING THE L2 COHERENCE CACHE

- The L2C\$ must keep the identity of the current owner cache for each block allocated in any L1 data cache.
  - This information is accessed when the requestor is not able to locate the owner cache.

- The L2C\$ is notified on every owner change through control messages.

- These messages should be processed by the L2C\$ in the very same order in which they were generated.
  - To ensure this, the L2C\$ sends an ACK message to the new owner when it receives a change owner message.
  - Until this message is not received by the owner node, it could use the block but cannot give the onwership to another cache.



### WRITE MISS IN DICO

## DIRECT COHERENCE
UPDATING THE L1 COHERENCE CACHE

- Base: information about the last core that invalidated or
  provided each block is kept in the L1C$.
  - Extra messages are not needed.
  - In some cases this information is not enough to obtain
    accurate predictions.

# DIRECT COHERENCE
## UPDATING THE L1 COHERENCE CACHE

- Base: information about the last core that invalidated or provided each block is kept in the L1C$.
  - Extra messages are not needed.
  - In some cases this information is not enough to obtain accurate predictions.
- Hints: control messages update the L1C$.
  - More accurate predictions.
  - Area and network traffic overhead.

### FREQUENT SHARERS (**FS**)

- Area: Duplicated sharing information.
- Network: Hints sent on each owner change.

### ADDRESS SIGNATURES (**AS**)

- Area: Two address signatures.
- Network: Hints filtering.

# DIRECT COHERENCE
## UPDATING THE L1 COHERENCE CACHE

- **Base**: information about the last core that invalidated or provided each block is kept in the L1C$.
  - Extra messages are not needed.
  - In some cases this information is not enough to obtain accurate predictions.
- **Hints**: control messages update the L1C$.
  - More accurate predictions.
  - Area and network traffic overhead.

---

### FREQUENT SHARERS (**FS**)

- Area: Duplicated sharing information.
- Network: Hints sent on each owner change.

### ADDRESS SIGNATURES (**AS**)

- Area: Two address signatures.
- Network: Hints filtering.

---

- **Oracle**: the requestor always knows the identity of the current owner.

# DIRECT COHERENCE
EVALUATION



TRAFFIC-INDIRECTION TRADE-OFF

- *Directory* introduces indirection in the critical path of cache misses.
- *Token* generates high levels of network traffic.
- *DiCo-Base* reduces traffic even compared to *Directory*, but the indirection avoidance is limited.
- *DiCo-Hints* policies slightly increase traffic compared to *DiCo-Base* and successfully avoid indirection.

# DIRECT COHERENCE
EVALUATION

## APPLICATIONS' EXECUTION TIME



- *DiCo-Hints AS* reduces execution time compared to *Directory* (9%) and *Token* (8%).

# DIRECT COHERENCE
TRAFFIC-AREA TRADE-OFF IN DICO

- We have obtained a good trade-off between execution time and network traffic.
- However, the area requirements of *DiCo* do not scale with the number of cores.
- There are other protocols that scale better in terms of area.





### CLASSIFICATION OF PROTOCOLS

|  | Traditional | Indirection-aware |
|---|---|---|
| Traffic-intensive | Hammer | Token |
| Area-demanding | Directory | DiCo |

# DIRECT COHERENCE
## TRAFFIC-AREA TRADE-OFF IN DICO

- Extra structures for keeping coherence:
  - L1C\$: One pointer to the predicted owner $\Rightarrow O(log_2 n)$
  - L2C\$: One pointer to the current owner $\Rightarrow O(log_2 n)$
  - Sharing information (L1 and L2): One bit per tile $\Rightarrow O(n)$
    - This structure compromises scalability.
- Solution: To use compressed sharing codes.
- Advantage of DiCo: The owner tile keeps cache coherence, so the first sharer (i.e., the owner) is always known.
  - Read misses do not need to check the sharing code field, so the compressed sharing code employed do not affect them.
  - Reduces network traffic compared to broadcast-based protocols even when the sharing information field is removed.

# DIRECT COHERENCE
COMPRESSED SHARING CODES

### SHARING CODES EVALUATED

| Protocol | Sharing Code | Bits L1 cache and L2 cache | Bits L1C\$ and L2C\$ | Order |
|----------|--------------|----------------------------|----------------------|-------|
| DiCo-FM | Full-map | $n$ | $log_2 n$ | $O(n)$ |
| DiCo-CV-K | Coarse vector | $\frac{n}{K}$ | $log_2 n$ | $O(n)$ |
| DiCo-LP-P | Limited pointers | $1 + P \times log_2 n$ | $log_2 n$ | $O(log_2 n)$ |
| DiCo-BT | Binary Tree | $\lceil log_2(1 + log_2 n) \rceil$ | $log_2 n$ | $O(log_2 n)$ |
| DiCo-NoSC | None | $0$ | $log_2 n$ | $O(log_2 n)$ |

- We evaluate the *DiCo-Hints AS* policy.
- *DiCo-FM* is the previously evaluated *DiCo-Hints AS* policy.

# DIRECT COHERENCE
EVALUATION



TRAFFIC-AREA TRADE-OFF

- *Hammer* and *Token* are traffic-intensive.
- *Directory* and *DiCo-FM* are area-demanding.
- *DiCo-BT* achieves a good compromise.
- *DiCo-NoSC* also achieves a good compromise without modifying the data caches.

# DIRECT COHERENCE
EVALUATION

## OVERALL TRADE-OFF



Execution time (normalized)

Network traffic (normalized)

Area required (mm$^2$)

# DIRECT COHERENCE
EVALUATION

### OVERALL TRADE-OFF

# DIRECT COHERENCE
EVALUATION

## OVERALL TRADE-OFF

# DIRECT COHERENCE
EVALUATION

## OVERALL TRADE-OFF

# DIRECT COHERENCE
EVALUATION

## OVERALL TRADE-OFF

# DIRECT COHERENCE
EVALUATION

## OVERALL TRADE-OFF



*DiCo-BT* obtains a very good compromise among execution time, network traffic and area requirements.

# DIRECT COHERENCE
## CONCLUSIONS

- Direct coherence protocols:
  - Do not rely on broadcasting requests.
  - Avoid the indirection for most cache misses.
  - Work well with compressed sharing codes.
- The following improvements have been obtained by DiCo-FM (Hints AS):
  - Execution time: 9% compared to *Directory* and 8% compared to *Token*.
  - Network traffic: 37% compared to *Token* and a slightly increase compared to *Directory*.
- *DiCo-BT* and *DiCo-NoSC* obtain a good trade-off among execution time, network traffic and area requirements.

# DIRECT COHERENCE
PUBLICATIONS

### INTERNATIONAL CONFERENCES

- **A. Ros**, M. E. Acacio and J. M. García, *"Direct Coherence: Bringing Together Performance and Scalability in Shared-Memory Multiprocessors"*. **HiPC'07**.
- **A. Ros**, M. E. Acacio and J. M. García, *"DiCo-CMP: Efficient Cache Coherency in Tiled CMP Architectures"*. **IPDPS'08**.
- **A. Ros**, M. E. Acacio and J. M. García, *"Dealing with Traffic-Area Trade-Off in Direct Coherence Protocols for Many-Core CMPs"*. **APPT'09**.

### INTERNATIONAL JOURNALS

- **A. Ros**, M. E. Acacio and J. M. García, *"A Direct Coherence Protocol for Many-Core Chip Multiprocessors"*. **TPDS**, Dec 2010.

### BOOK CHAPTERS

- **A. Ros**, M. E. Acacio and J. M. García, *"Cache Coherence Protocols for Many-Core CMPs"*. Parallel and Distributed Computing.

# DIRECT COHERENCE
CURRENT WORK

- Heterogeneous networks:
  - Network provided with fast and low-power links.
  - Non-critical messages can be sent by low-power links.
  - DiCo increases the number of non-critical messages: hints.



NON-CRITICAL TRAFFIC

# DIRECT COHERENCE
CURRENT WORK

- Heterogeneous networks:
  - Network provided with fast and low-power links.
  - Non-critical messages can be sent by low-power links.
  - DiCo increases the number of non-critical messages: hints.



NON-CRITICAL TRAFFIC

- Server consolidation or multiprogrammed workloads:
  - Several virtual machines (VM) in a CMP.
  - Home nodes can map anywhere.
  - Owner nodes will likely be in the same VM.

## OUTLINE

## COHERENCE DEACTIVATION
MOTIVATION

### REMEMBER

Directory protocols are the most scalable alternative for keeping cache coherence.

- But the area requirements of the directory structure could become prohibitive for large-scale multiprocessors.

# COHERENCE DEACTIVATION
## MOTIVATION

> ### REMEMBER
> Directory protocols are the most scalable alternative for keeping cache coherence.
>
> - But the area requirements of the directory structure could become prohibitive for large-scale multiprocessors.

- Directory caches accelerate the access to the coherence information and reduce directory overhead with respect to a memory directory but...
  - ...directory cache evictions cause the invalidation of cached data, resulting in performance degradation.

# COHERENCE DEACTIVATION
## MOTIVATION

- Is it necessary to keep cache coherence for all referenced blocks?
  - Private blocks will never be incoherent!
    - 75% of referenced blocks (on average).

- If we do not maintain directory information for these blocks we can save a lot of directory storage.

# COHERENCE DEACTIVATION
PROPOSAL

- We propose a mechanism that:
    - Classifies memory blocks into private and shared (How?).
    - Deactivates the coherence protocol for private blocks
        - i.e., do not keep track of them.
- A block-grain classification would require significant storage resources.
    - Blocks are classified at page granularity.
    - The operating system detects when a page (initially considered private) becomes shared (minimal OS overhead).
        - Performed upon TLB misses: state stored in the page table.
        - A coherence recovery mechanism is necessary to restore block's coherence status.
    - Collaboration between hardware and operating system.

# COHERENCE DEACTIVATION
EXAMPLE



REQUESTS FOR PRIVATE AND SHARED BLOCKS

# COHERENCE DEACTIVATION
## THE RECOVERY MECHANISM

- When a private page becomes shared, it is necessary:
  1. Make coherent directory caches with cached blocks ⇒ recovery mechanism.
     - Triggered by OS during TLB miss resolution (critical section)
  2. From this point on, keep track of the blocks in this page.
- To main options:
  - Flushing-based recovery: evicts blocks in the page being recovery from processor caches.
  - Updated-based recovery: updates directory caches with the information about cached blocks.

# COHERENCE DEACTIVATION
FLUSHING-BASED RECOVERY MECHANISM

## COHERENCE DEACTIVATION
UPDATING-BASED RECOVERY MECHANISM

### EXAMPLE OF THE UPDATING-BASED RECOVERY

# COHERENCE DEACTIVATION
## ADVANTAGES

- The amount of directory information required to maintain coherence is reduced.
  - Reduce directory cache evictions to improve performance.
  - Reduce directory cache size while keeping performance.
- Request for private blocks do not need to access the directory structure.
  - Savings in both cache miss latency and power consumption.

# COHERENCE DEACTIVATION
## RESULTS: DETECTED BLOCKS

- Private blocks (75%) $\Rightarrow$ Detected private blocks (57%).
- 57% memory blocks are not tracked

## PRIVATE BLOCKS VS. DETECTED BLOCKS

# COHERENCE DEACTIVATION
## RESULTS: CACHE MISSES

- 5C cache miss classification.
- 75% coverage misses avoided (35% overall).

### NORMALIZED NUMBER OF PROCESSOR CACHE MISSES

# COHERENCE DEACTIVATION
RESULTS: EXECUTION TIME

- Similar runtime for flushing and updating recovery mechanisms.
- With same directory size ⇒ performance improvement: 15%.
- With same performance ⇒ directory cache 8 times smaller.



NORMALIZED RUNTIME

# COHERENCE DEACTIVATION
## RESULTS: ENERGY CONSUMPTION

- Reduced by 40% when keeping directory cache size.



NORMALIZED DYNAMIC ENERGY CONSUMPTION

## COHERENCE DEACTIVATION
### FUTURE WORK AND PUBLICATIONS

- Publications:

#### INTERNATIONAL CONFERENCES

- B. Cuesta, **A. Ros**, M. E. Gómez, A. Robles, and J. Duato, *"Increasing the Effectiveness of Directory Caches by Deactivating Coherence for Private Memory Blocks"*. **ISCA'11**.

#### INTERNATIONAL JOURNALS

- B. Cuesta, **A. Ros**, M. E. Gómez, A. Robles, and J. Duato, *"Increasing the Effectiveness of Directory Caches by Avoiding the Tracking of Non-Coherent Memory Blocks"*. Submitted to **TC**.

# COHERENCE DEACTIVATION
## FUTURE WORK AND PUBLICATIONS

- Publications:

### INTERNATIONAL CONFERENCES

- B. Cuesta, **A. Ros**, M. E. Gómez, A. Robles, and J. Duato, *"Increasing the Effectiveness of Directory Caches by Deactivating Coherence for Private Memory Blocks"*. **ISCA'11**.

### INTERNATIONAL JOURNALS

- B. Cuesta, **A. Ros**, M. E. Gómez, A. Robles, and J. Duato, *"Increasing the Effectiveness of Directory Caches by Avoiding the Tracking of Non-Coherent Memory Blocks"*. Submitted to **TC**.

- Future work:
  - Thread migration can reduce the number of non-coherent blocks.
  - A page-grained classification misclassifies about 18% of blocks.
    - Blocks detected as coherent are actually non-coherent.
  - Possible solutions: Hardware mechanisms or modifications in the programing language?

## OUTLINE

## SYNCHRONOUS COHERENCE
MOTIVATION

### REMEMBER

- The verification of a cache coherence protocol is very time-consuming and tedious.

# SYNCHRONOUS COHERENCE
MOTIVATION

## REMEMBER

- The verification of a cache coherence protocol is very time-consuming and tedious.

- The more complex the coherence protocol is, the more verification time is required.
- The appearance of race conditions makes even harder the protocol verification.
- Some authors reduce protocol races by relying on atomic transitions [1].
- Another approach: simple request-response protocols.

## REFERENCES

[1] D. Vantrease, M. H. Lipasti, and N. Binkert, *"Atomic Coherence: Leveraging Nanophotonics to Build Race-Free Cache Coherence Protocols"*. **HPCA'10**.

# SYNCHRONOUS COHERENCE
## REQUEST-RESPONSE PROTOCOLS

- A request-response protocol does not forwards requests to other nodes (2-hop protocol).
    - The requester issues a message to the home node.
    - The home node directly responds with a copy of the request block.
- What happens with dirty cached copies?
    - Write-through caches? $\Rightarrow$ Not very efficient.
    - Solution: time-based cache coherence protocols (synchronous coherence).
        - A global clock is needed $\Rightarrow$ use of global lines [2].
        - Block stored in cache will have expiration date!
        - When a cached block expires it will be invalidated, performing a writeback in case the block is dirty.

## REFERENCES

[2] R. T. Chang, N. Talwalkar, C. P. Yue, and S. S. Wong, *"Near Speed-of-Light Signaling Over On-Chip Electrical Interconnects"*. IEEE Journal of Solid-State Circuits, 2003.

# SYNCHRONOUS COHERENCE
EXAMPLE



REQUEST-RESPONSE PROTOCOL
WITH EXPIRATION DATE FOR CACHED BLOCKS

# SYNCHRONOUS COHERENCE
EXAMPLE



REQUEST-RESPONSE PROTOCOL
WITH EXPIRATION DATE FOR CACHED BLOCKS

Cycle 0                                         Directory

Cache 0        Cache 1        Cache n

...

GetS

# SYNCHRONOUS COHERENCE
EXAMPLE

- Directory does not keep list of sharers but expiration date.



REQUEST-RESPONSE PROTOCOL
WITH EXPIRATION DATE FOR CACHED BLOCKS

# SYNCHRONOUS COHERENCE
EXAMPLE

- Directory does not keep list of sharers but expiration date.

REQUEST-RESPONSE PROTOCOL
WITH EXPIRATION DATE FOR CACHED BLOCKS

# SYNCHRONOUS COHERENCE
EXAMPLE

- Directory does not keep list of sharers but expiration date.

REQUEST-RESPONSE PROTOCOL
WITH EXPIRATION DATE FOR CACHED BLOCKS

# SYNCHRONOUS COHERENCE
### EXAMPLE

- Directory does not keep list of sharers but expiration date.
- GetX transaction waits until the block expires.

# SYNCHRONOUS COHERENCE
## EXAMPLE

- Directory does not keep list of sharers but expiration date.
- GetX transaction waits until the block expires.
- Memory sends the block to the requester and a new expiration date is assigned.

### REQUEST-RESPONSE PROTOCOL
### WITH EXPIRATION DATE FOR CACHED BLOCKS

## OUTLINE

## INTRODUCTION

- There are several challenges to address for the memory hierarchy organization of a CMP.
  - Thread Balancing problems.
    - Imbalance in time: Some threads arrive to a barrier before the other ones $\Rightarrow$ Can increase execution time.
    - Imbalance in storage: The working set of threads also varies $\Rightarrow$ Can increase cache misses (off-chip accesses).
  - Conflict misses.
    - Reduce last level conflict misses also can save off-chip accesses.
  - Long access latency to NUCA banks.
    - Several authors address this problem but they do not care about directory scalability.

## OUTLINE

# REPLACEMENT POLICIES FOR SHARED CACHES
## MOTIVATION

- In parallel applications, some threads arrive to a barrier before the other ones.
- The first threads arriving to a barrier start a busy waiting.
  - This consumes extra power.
- Some authors propose to save power consumption by slowing down faster threads (e.g., reducing processor frequency) [3,4].
  - This saves power but it does not improve execution time.
- Another approach: A thread-aware replacement policy.

## REFERENCES

[3]    J. Li, J. F. Martínez, and M. C. Huang, *"The Thrifty Barrier: Energy-Aware Synchronization in Shared-Memory Multiprocessors"*. **HPCA'04**.

[4]    Q. Cai, J. González, R. Rakvic, G. Magklis, P. Chaparro, and A. González, *"Meeting Points: Using Thread Criticality to Adapt Multicore Hardware to Parallel Regions"*. **PACT'08**.

# REPLACEMENT POLICIES FOR SHARED CACHES
## A THREAD-AWARE REPLACEMENT POLICY

- Sets in a shared cache hold blocks from different threads.
- A smart policy can be implemented:
    - Avoid evictions of blocks accessed by slower threads, or widely shared.
    - Evicts private blocks accessed by faster threads.
- Directory caches already store information about which processors hold the blocks.
- Since slower threads are accelerated, the final execution time can be reduced.
- Another option for balancing threads ⇒ Lock priorities.
    - Give more priority for lock acquisition to slower threads.

## OUTLINE

| Alberto Ros | **Research lines** | Uppsala, May 27, 2011 | **52 / 63** |

# INDEXING POLICIES FOR SHARED CACHES
MOTIVATION

- Memory references are not often distributed across cache sets.
  - Some sets exhibit large miss ratios, while other are underutilized.
- This causes the appearance of conflict misses.
  - Can be reduced by increasing associativity.
    - But this would increase power consumption and access latency.
- Misses in the shared last-level cache cause expensive off-chip accesses.
  - Some authors reduce conflict misses by reallocating blocks to underutilized sets [5].
  - Another approach: adaptive selection of index bits.

## REFERENCES

[5]  D. Rolán, B. B. Fraguela, and R Doallo, *"Adaptive Line Placement with the Set Balancing Cache"*. **MICRO'09**.

# INDEXING POLICIES FOR SHARED CACHES
## MOTIVATION EXAMPLE

- If we carefully chose the address bits for indexing the cache, a better set balancing can be obtained.
- Why choose other bits apart from the least significant bits (LSB)?
    - Example 1: stride memory access pattern.
    - Example 2: parallel applications (shared array partitioned).
    - Example 3: second level caches (L1 filter accesses to contiguous blocks).

### LSB INDEXING VS. OTHER BITS INDEXING



| | LSB indexing | | | | Other-bits indexing | |
|---|---|---|---|---|---|---|

Memory addresses requested — LSB indexing — Memory addresses requested — Other−bits indexing

4 3 2 1 0 ◄— Bit position    DM cache

```
0 0 0 0 0  ──►  0 0 0  ┐            set 0
0 0 1 0 0  ──►  1 0 0  │            set 1
0 1 0 0 0  ──►  0 0 0  │            set 2
0 1 1 0 0  ──►  1 0 0  │            set 3
1 0 0 0 0  ──►  0 0 0  │            set 4
1 0 1 0 0  ──►  1 0 0  │            set 5
1 1 0 0 0  ──►  0 0 0  │            set 6
1 1 1 0 0  ──►  1 0 0  │            set 7
```

4 3 2 1 0 ◄— Bit position    DM cache

```
0 0 0 0 0  ──►  0 0 0  ──►  set 0
0 0 1 0 0  ──►  0 0 1  ──►  set 1
0 1 0 0 0  ──►  0 1 0  ──►  set 2
0 1 1 0 0  ──►  0 1 1  ──►  set 3
1 0 0 0 0  ──►  1 0 0  ──►  set 4
1 0 1 0 0  ──►  1 0 1  ──►  set 5
1 1 0 0 0  ──►  1 1 0  ──►  set 6
1 1 1 0 0  ──►  1 1 1  ──►  set 7
```

■ Bits selected to form the index    □ Unused set    □ Conflict−free set    ■ Conflicting set

## INDEXING POLICIES FOR SHARED CACHES
CURRENT AND FUTURE WORK

- Current work:
  - A first approach for direct-mapped first-level caches.
  - Submitted to **PACT'11**.
- Future work:
  - Private caches in CMPs.
  - Set-associative caches.
  - Last-level cache in CMPs.
    - Where misses cause expensive off-chip accesses.
    - Thread balancing?

## OUTLINE

# NUCA MAPPING AND DIRECTORY SCALABILITY
MOTIVATION

### REMEMBER

In NUCA (Non-Uniform Cache Architecture) caches, the access latency depends on where the requested block is mapped (home bank).

# NUCA MAPPING AND DIRECTORY SCALABILITY
MOTIVATION

### REMEMBER

In NUCA (Non-Uniform Cache Architecture) caches, the access latency depends on where the requested block is mapped (home bank).

- This mapping is commonly performed by taking some bits from the block address leading to a Round-Robin mapping.
    - The *Round-Robin* mapping does not care about the distance between requesting cores and home banks $\Rightarrow$ long access latency.
- A First-Touch mapping policy can lessen this latency.
    - But can cause imbalance among cache banks $\Rightarrow$ high cache miss rate.

# NUCA MAPPING AND DIRECTORY SCALABILITY
## PREVIOUS WORK

- Several authors have study the trade-off between low miss rate and low access time in NUCA caches [6,7].
  - But these works does not care about directory scalability.
  - They are based on OS allocation policies at page granularity...
    - ...which can affect directory scalability.

## REFERENCES

[6]  N. Hardavellas, M. Ferdman, B. Falsafi, and A. Ailamaki, *"Reactive NUCA: Near-optimal block placement and replication in distributed caches"*. **ISCA'09**.

[7]  **A. Ros**, M. Cintra, M. E. Acacio and J. M. García, *"Distance-Aware Round-Robin Mapping for Large NUCA Caches"*. **HiPC'09**.

# NUCA MAPPING AND DIRECTORY SCALABILITY
## A NEW METRIC: DIRECTORY SCALABILITY

- A directory cache based on duplicated tags can perfectly scale (in size) up to a certain number of nodes [8].
  - This number of nodes corresponds to the number of private cache sets.
    - Commonly, first-level caches have between 128 and 512 sets.
  - Constraint: the mapping of memory blocks to home banks must be done at fine granularity (i.e., block granularity).



TRADE-OFF DIAGRAM

LLC hit rate

Locality

## REFERENCES

[8] **A. Ros**, M. E. Acacio and J. M. García, *"A Scalable Organization for Distributed Directories"*. **JSA**, Mar, 2010.

# NUCA MAPPING AND DIRECTORY SCALABILITY
## A NEW METRIC: DIRECTORY SCALABILITY

- A directory cache based on duplicated tags can perfectly scale (in size) up to a certain number of nodes [8].
  - This number of nodes corresponds to the number of private cache sets.
    - Commonly, first-level caches have between 128 and 512 sets.
  - Constraint: the mapping of memory blocks to home banks must be done at fine granularity (i.e., block granularity).



TRADE-OFF DIAGRAM

LLC hit rate

RR

[6,7]

FT

Locality

### REFERENCES

[8]  **A. Ros**, M. E. Acacio and J. M. García, *"A Scalable Organization for Distributed Directories"*. **JSA**, Mar, 2010.

# NUCA MAPPING AND DIRECTORY SCALABILITY
A NEW METRIC: DIRECTORY SCALABILITY

- A directory cache based on duplicated tags can perfectly scale (in size) up to a certain number of nodes [8].
  - This number of nodes corresponds to the number of private cache sets.
    - Commonly, first-level caches have between 128 and 512 sets.
  - Constraint: the mapping of memory blocks to home banks must be done at fine granularity (i.e., block granularity).



TRADE-OFF DIAGRAM

LLC hit rate, [6,7], RR, [8], FT, Locality, Scalability

REFERENCES

[8]  **A. Ros**, M. E. Acacio and J. M. García, *"A Scalable Organization for Distributed Directories"*. **JSA**, Mar, 2010.

# NUCA MAPPING AND DIRECTORY SCALABILITY
## A NEW METRIC: DIRECTORY SCALABILITY

- A directory cache based on duplicated tags can perfectly scale (in size) up to a certain number of nodes [8].
  - This number of nodes corresponds to the number of private cache sets.
    - Commonly, first-level caches have between 128 and 512 sets.
  - Constraint: the mapping of memory blocks to home banks must be done at fine granularity (i.e., block granularity).



TRADE-OFF DIAGRAM

### REFERENCES

[8]   **A. Ros**, M. E. Acacio and J. M. García, *"A Scalable Organization for Distributed Directories"*. **JSA**, Mar, 2010.

# NUCA MAPPING AND DIRECTORY SCALABILITY
## POSSIBLE SOLUTION

- Directory scalability requires block-grained interleaving.
- Low latency and miss rate approaches employ page-grained interleaving.
- Possible solution ⇒ Decoupling directory information and data blocks.
  - Read requests may not require directory information.
    - So they can be sent to the data home bank.
  - Upgrade requests do not require data.
    - So they can be sent to the directory home bank.
  - Only few requests will require both data and directory information.

## OUTLINE

## CONCLUSIONS

- Currently, working on several proposals for improve both coherence protocols and cache hierarchy.
- Publications in high-level conferences and journals.
- Future work on both fields seems promising.
- Collaboration between hardware and operating system is very effective.
- Collaboration at other levels: compiler and programming languages is also interesting.

# EFFICIENT AND SCALABLE CACHE COHERENCE FOR MANY-CORE ARCHITECTURES

### Alberto Ros

Research assistant
Computer Engineering Department
Technical University of Valencia
aros@gap.upv.es

Adjunct professor
Computer Engineering Department
University of Murcia
a.ros@ditec.um.es

Uppsala, May 27, 2011