

Adaptive 3D Cache Architecture for Manycores

Eric Guthmuller CEA-Leti, DACLE, LISAN UPMC, LIP6, ALSoC



GDR SoC SiP - 2012-11-15

Introduction

- Number of cores in computing ICs is exploding
- Applications use more and more memory amount
- Cores need a fast access to the memory

=> Memory Wall





Introduction

Memory bandwidth per core in Nvidia GPUs over time:



source: http://denalimemoryreport.com



Outline

- Context and motivations
 - The cube law
 - 3D technologies
- TSAR manycore architecture
 - Architecture
 - Cache Coherency
- Adaptive distributed 3D cache
 - TSAR 3D L3
 - Adaptivity
 - Architecture
- Results

- Tile allocation
- Fault tolerance
- Hardware implementation

Outline

- Context and motivations
 - The cube law
 - 3D technologies
- TSAR manycore architecture
 - Architecture
 - Cache Coherency
- Adaptive distributed 3D cache
 - TSAR 3D L3
 - Adaptivity
 - Architecture
- Results

- Tile allocation
- Fault tolerance
- Hardware implementation



Context and Motivations - Cube Law

- The cache miss rate is a S^x function where S is the size of the cache (see [1])
- x is between 0.36 and 0.62 on commercial workloads
 - => Mean of 0.5

leti&li/t ĽP



Context and Motivations - Cube Law

- With no sharing between processors: memory accesses are multiplied by the number of cores
- To keep constant the external memory bandwidth when the number of cores is multiplied by N, the total cache size must be multiplied by N³
- With 16% of data sharing (mean for PARSEC benchmarks [1]):

 $=> (1-0.16)^{3*}N^{3} = 0.6^{*}N^{3}$ (with N big)

leti&list ID

Manycores need very big caches and/or fast memory
 => Go to 3D caches



Context and Motivations - 3D Stack

- Memory on top of processors
- Multiple layers of logic => Through Silicon Vias (TSV)
- Example of stacking:





Context and Motivations - 3D TSV

Different TSV technologies (lined / filled)



Trench AR 20, AR 10, Ø10x100µm X700'

AR= Aspect Ratio



Adaptive 3D Cache Architecture for Manycores – GDR SOCSIP – Eric Guthmuller | 15 November 2012 9 © CEA. All rights reserved

5x100µm

Source: Patrick Leduc; LETI

Context and Motivations - 3D Bumps



© CEA. All rights reserved

Context and Motivations - Conclusion

- Use dense 3D on-chip memory to build a n+1 level cache
- 3D memory distributed amongst the circuit with many 3D links
 - => High bandwidth and fault tolerance
- Goal: stackable and reusable architecture to adapt memory quantity to the application needs, including fault tolerance

leti&list ID



Outline

- Context and motivations
 - The cube law
 - 3D technologies
- TSAR manycore architecture
 - Architecture
 - Cache Coherency
- Adaptive distributed 3D cache
 - TSAR 3D L3
 - Adaptivity
 - Architecture
- Results

- Tile allocation
- Fault tolerance
- Hardware implementation



TSAR Manycore [2] - Architecture

leti&li/t LP

MP²SOC with a 2D-mesh NoC (DSPIN) => GALS



Adaptive 3D Cache Architecture for Manycores – GDR SOCSIP – Eric Guthmuller | 15 November 2012 | 13 © CEA. All rights reserved

TSAR Manycore - Architecture

leti&li/t LP







TSAR Manycore - Cache Coherency

- Thanks to the scalable bandwidth provided by the NoC technology, the TSAR architecture use a write-through policy between L1 and L2 caches
 - ⇒ the L1 cache architecture and the coherency protocol are simplified
 - ⇒ the cache coherency is ensured by the L2 cache (Home Agent)
- As snooping does not scale with a large number of cores, TSAR relies on a L2 cache directory-based approach

⇒ Each L2 cache tracks all the copies stored in the L1 caches in order to update them

leti&list UP



TSAR Manycore - Cache Coherency

- TSAR implements the DHCCP protocol (Distributed Hybrid Cache Coherence Protocol):
 - multicast / update when the number of copies is smaller than a predefined threshold
 - broadcast / invalidate when the number of copies is larger than a predefined threshold
- The DHCCP has been analyzed, from the point of view of deadlock prevention. Three types of traffic have been identified:
 - Direct read/write transactions (L1 caches => L2 caches)
 - Coherency transactions (L2 caches <=> L1 caches)
 - External memory accesses (L2 caches => Ext. RAM)

6 communication channels

TSAR Manycore - Cache Coherency

Write-through:

- Simple coherency protocol
- Simple hardware
- BUT more network traffic
- Caches need to be inclusive

Write-back:

- Network traffic reduced
- Automatic replication and distribution
- Cache storage better exploited
- BUT complex protocol
- Problem of false sharing



Outline

- Context and motivations
 - The cube law
 - 3D technologies
- TSAR manycore architecture
 - Architecture
 - Cache Coherency
- Adaptive distributed 3D cache
 - TSAR 3D L3
 - Adaptivity
 - Architecture
- Results

- Tile allocation
- Fault tolerance
- Hardware implementation



3D Cache - TSAR 3D L3

leti&li/t LP



3D Cache - TSAR 3D L3

leti&li/t LP







leti&li/t LP

3D Cache - TSAR 3D L3

Proposed 3D architecture: regular 3D mesh

Hypothesis: 1:1 Proc tile/Cache tile





3D Cache - TSAR 3D L3

leti&li/t LP

Proposed 3D architecture: regular 3D mesh





3D Cache - Adaptivity

leti&list LP

Allocation of cache tiles to memory segments



3D Cache - Adaptivity

Allocation of cache tiles to memory segments



- Different cache quantities for different memory segments
- Shared cache tiles: the most accessed memory segment takes the most part of the cache tile
- Conclusion: 2 levels of adaptability:

- Programmable by the soft (OS) at runtime
- Very dynamic when cache tiles are shared

Cache controllers integrated in the computing tier



Selects the tile coordinates and its cache set to use



Cache controllers integrated in the computing tier



Selects the tile coordinates and its cache set to use



• Four 3D 64 bits DSPIN NoC channels:

L2->L3 (req+rsp) + L3->Ext Mem (req + rsp)

leti&li/t LP





Adaptive 3D Cache Architecture for Manycores – GDR SOCSIP – Eric Guthmuller | 15 November 2012 | 30 © CEA. All rights reserved



- Low frequency high throughput tile
- 1 MB cache tile



Cache Tile



Outline

- Context and motivations
 - The cube law
 - 3D technologies
- TSAR manycore architecture
 - Architecture
 - Cache Coherency
- Adaptive distributed 3D cache
 - TSAR 3D L3
 - Adaptivity
 - Architecture
- Results

- Tile allocation
- Fault tolerance
- Hardware implementation





16 tiles of 4 cores

leti&list LP

- 16x1MB 3D cache tiles per memory tier
- L3 tiles frequency == Processor frequency

© CFA. All rights reserved

Results: Benchmark platform

- Bandwidth to External Memory: we emulate an equivalent 1024 cores circuit
- Splash 2 multithreaded benchmarks on MutekH [5] OS:
 - Radix: a radix sorting algorithm
 - FFT: a distributed FFT
 - LU: a matrix factorization algorithm
 - Ocean: a sea current solver



First case: Private caches



- Two applications, 32 cores each:
 - Memory intensive (red)
 - Computational (grey)
- 1st case: private 3D cache tiles

leti&list LP

Adaptive 3D Cache Architecture for Manycores – GDR SOCSIP – Eric Guthmuller | 15 November 2012 | 35 © CEA. All rights reserved

ALLOCATION !

Second case: Shared caches



2nd case: shared 3D cache tiles

leti&list UP

- Competition between the two applications on cache utilization
 - => The most demanding application will occupy more cache space



Results: Tiles Allocation

- We study execution time for shared caches wrt. private caches
- Mix of Splash FFT and Splash LU applications for test cases

The most memory intensive application performs better with shared caches

Execution Time with shared caches



=> Quality of Service vs Best Effort

leti&li/t ĽP



Total traffic to memory with shared caches

external memory is reduced by up to 55%

Results: Tiles allocation

App Set	Memory Intensive App	Computational App	L3 cache
1	FFT 2 ¹⁸	FFT 2 ¹⁴	1 tier (16MB)
2	FFT 2 ²⁰	FFT 2 ¹⁸	4 tiers (64MB)
3	LU 1024	LU 512	1 tier (16MB)
4	FFT 2 ¹⁸	LU 512	1 tier (16 MB)
5	FFT 2 ¹⁸	LU 1024	1 tier (16 MB)
6	FFT 2 ²⁰	LU 512	4 tiers (64MB)



Adaptive 3D Cache Architecture for Manycores – GDR SOCSIP – Eric Guthmuller | 15 November 2012 | 39 © CEA. All rights reserved



We consider only faulty cache tiles

First order algorithm

leti&li/t ĽP

- A list of faulty tiles in each cache access controller and a decision algorithm
- Two possible algorithms if we select a faulty tile



V	Faulty	Good
1	(x _f ,y _f ,z _f)	(x_g, y_g, z_g)
•	•••	•••

Second order algorithm

Results: Fault Tolerance

- We compare the two algorithms while we introduce faults
- 1 tier 3D cache (16 tiles) with FFT 2¹⁸
 - ⇒ 1 faulty tile: 13.5% execution time increase for first order algo versus 2.9% for second order
- 4 tiers 3D cache (64 tiles) with FFT 2²⁰

leti&list

- ⇒ 1 faulty tile: 1.8% execution time increase for first order algo versus 0% for second order
- ⇒ 4 faulty tiles: 9.3% execution time increase for first order algo versus 2.7% for second order



Outline

- Context and motivations
 - The cube law
 - 3D technologies
- TSAR manycore architecture
 - Architecture
 - Cache Coherency
- Adaptive distributed 3D cache
 - TSAR 3D L3
 - Adaptivity
 - Architecture
- Results

- Tile allocation
- Fault tolerance
- Hardware implementation

Hardware Implementation - Synthesis

- Synthesis of a 1MB cache tile in 65nm with SRAM memories
 - \Rightarrow Extrapolation to 32nm and eDRAM
- 3 TSV sizes : 10μm (WideIO), 5μm and 2μm
- Power TSVs: according to the Wide IO example \Rightarrow 68 TSVs, 10µm diam. (272 for 5µm and 1700 for 2µm)
- Signal TSVs: 284 for the 3D NoCs (worst case tile)
 - L2/L3 NoC

leti&list ID

L3/ExtMem NoC

Hardware Implementation - Synthesis

- Synthesis of a 1MB cache tile in 65nm with the two 3D NoCs
- TSV arrays:
 - 3D NoC signals
 - Memory tier power supply
- Control:
 - Directory
 - Logic
 - NoC routers
- Data:

leti&li/t LP

Cached data



Conclusion

- How many cores in a manycore architecture without saturating memory access ?
 - TSAR, NUMA cache coherent manycore
 - 4 x 2133 MHz DDR3 interfaces
 - Constant computing density (GFLOPS/mm²)
 - HPC application



Conclusion

- The memory bandwidth is the limiting factor for manycore architectures
- 3D big caches are a solution
 - Even if they have only a S^{0.5} impact on memory accesses
- We propose a low power and high performance 3D cache architecture with
 - Adaptivity

- Fault tolerance
- A 28nm hardware implementation has been performed with low 3D overhead



References

- [1] B. M. Rogers, A. Krishna, G. B. Bell, K. Vu, X. Jiang, et Y. Solihin, "Scaling the bandwidth wall: challenges in and avenues for CMP scaling," SIGARCH Comput. Archit. News, vol. 37, n

 371–382, juin 2009.
- [2] TSAR architecture: https://www-soc.lip6.fr/trac/tsar/
- [3] Guthmuller, E.; Miro-Panades, I.; Greiner, A.; , "Adaptive Stackable 3D Cache Architecture for Manycores," VLSI (ISVLSI), 2012 IEEE Computer Society Annual Symposium on , vol., no., pp.39-44, 19-21 Aug. 2012
- [4] SoClib project: http://www.soclib.fr
- [5] MutekH Operating System: http://www.mutekh.org

leti

LABORATOIRE D'ÉLECTRONIQUE ET DE TECHNOLOGIES DE L'INFORMATION

list

LABORATOIRE D'INTEGRATIO DES SYSTÈMES ET DES TECHNOLOGIES



Thank you

