The Best of All Worlds: Improving Predictability at the Performance of Conventional Coherence with No Protocol Modifications

Salah Hessien, <u>Mohamed Hassan</u>







Data is a Key in all modern applications

ilot tor iters	Air Data Rreference Unit	Actuator Control Electronics		\mathcal{T}
P	F	P		//
-0-			200	
, _				
ft Inform nagem System	Flap/Slat Electronics U mation ent	Proximity nit Electronics Unit		





Common Approach

• Adopts an independent-task model \rightarrow No communication

• Enforcing complete isolation between tasks. • At the shared cache: strict cache partitioning and coloring • At the DRAM: bank privatization





- May result in a poor memory or cache utilization
 - e.g.: a task has conflict misses, while other partitions may remain underutilized
- Does not scale with increasing number of cores
 - e.g.: number of PEs \leq number of DRAM banks
- Not viable in emerging systems due to increased functionality and massive data

Common Approach





Data Sharing









Simpler timing analysis
 Hardware changes
 Long execution time

Solution: No caching of shared data [Hardy et al., RTSS'09] [Lesage et al., RTNS'10] [Bansal et al., arXiv'19] [Chisholm et al., RTSS'16]



Solution:

[Hardy et al., RTSS'09] [Lesage et al., RTNS'10] [Bansal et al., arXiv'19] [Chisholm et al., RTSS'16]

No caching of shared data



The mainstream solution is to provide shared memory and prevent incoherence through a hardware cache coherence protocol, making caches functionally invisible to software.

Coherence is the norm in COTS platforms

DOI:10.1145/2209249.2209269

On-chip hardware coherence can scale gracefully as the number of cores increases.

BY MILO M.K. MARTIN, MARK D. HILL, AND DANIEL J. SORIN

Why On-Chip Cache **Coherence Is** Here to Stay

SHARED MEMORY IS the dominant low-level communication paradigm in today's mainstream multicore processors. In a shared-memory system, the (processor) cores communicate via loads and stores to a shared address space. The cores use caches to reduce the average memory latency and memory traffic. Caches are thus beneficial, but private caches lead to the possibility of cache incoherence. The mainstream solution is to provide shared memory and prevent incoherence through a hardware cache coherence protocol, making caches functionally invisible to software. The incoherence problem and basic hardware coherence solution are outlined in the sidebar, "The Problem of Incoherence," page 86.

Cache-coherent shared memory is provided by mainstream servers, desktops, laptops, and mobile devices and is available from all major vendors, including AMD, ARM, IBM, Intel, and Oracle (Sun).

78 COMMUNICATIONS OF THE ACM | JULY 2012 | VOL. 55 | NO.

Data Sharing





Heterogeneous compute requires coherency

- Flexible heterogeneous architecture
- · Blend compute and acceleration for target solution
- · Fast, reliable transport to shared memory
- Maximize throughput, minimize latency
- Accelerate SoC deployment
- · IP designed, optimized and validated for systems
- CAMP12016 41



Coherence is the Industry's Choice





Coherency: The New Normal in SoCs Anush Mohandass anush@netspeedsystems.com





oday's SoCs include a mix of CPU cores, computing clusters, GPUs and other computing resources and specialized accelerators.

Getting heterogeneous processors to communicate efficiently is a daunting design challenge. A popular approach is to use high-performance and power-efficient shared-memory communication and a sophisticated on-chip cache-coherent

interconnect. This presentation will introduce a new technology that automates the architecture design process, supports CHI and ACE in one design, and uses advanced machinelearning algorithms to create an optimal pre-verified cache-coherent solution.

Data Sharing







Coherence is the Industry's Choice

autonomous driving requirements are mandating the simultaneous use of multiple types of processing units to efficiently execute sophisticated image processing, sensor fusion, and machine learning/AI algorithms. This presentation introduces **new** coherency platform technology that enables the integration of heterogeneous cache coherent hardware accelerators and CPUs, using a mixture of ARM ACE, CHI, and CHI Issue B protocols, into systems that meet both the requirements of high compute performance and ISO 26262-compliant functional safety.









* = Other-GetM, Other-GetS, or Other-PutM. ** = Other-GetM or Other-GetS.

Coherence is A complex Machinery

Data Sharing



- •High performance arbitration \checkmark
 - FCFS, Split-Tx, Priority-Based, ...
- •Shared data is entirely supported through a hardware cache coherence protocol 💜
- •Not Predictable X



State-of-the-art



- High performance arbitration 🗡
 - FCFS, Split-Tx, Priority-Based, ...
- •Shared data is entirely supported through a hardware cache coherence protocol 💜
- •Not Predictable X



State-of-the-art



- •High performance arbitration \checkmark
 - FCFS, Split-Tx, Priority-Based, ...
- •Shared data is entirely supported through a hardware cache coherence protocol 💙
- •Not Predictable X



State-of-the-art









State-of-the-art

Data Sharing





State-of-the-art

Data Sharing



- High performance arbitration
 - FCFS, Split-Tx, Priority-Based, ...
- •Shared data is entirely supported through a hardware cache coherence protocol 💜
- •Not Predictable 🗡





Data-Aware Arbitration:

- Builds on traditional arbitration
- Shared data is supported
- •BUT requires coherence protocol modifications
- Coherence leads to good performance
- •Also achives predictability, but with significant V latency bounds
- •PMSI [RTAS'17], CARP [RTSS'19], HourGlass [Arxiv'18], PENDULUM [RTSS'19]

State-of-the-art

Traditional Real-Time:

- Predictable by design
 - TDM, RR, WRR, HKK, ...
- Does not support shared data, put aside coherence









Benefit of Coherence: Up to 3x performance

[RTAS'17] Mohamed Hassan, Anirudh M. Kaushik, Hiren Patel, "Predictable Cache Coherence for Multi-Core Real-Time Systems"





CO cannot receive A from Mem, it has to wait for C1 to WB first

- •Assume same example from before
- •But, let all 3 cores request to access same cache line A
- •Assume initially A is modified by C1 in its private cache
- •Let as before, our request under analysis is St(A) from C2

Problem: Coherence effect on WC





CO cannot receive A from Mem, it has to wait for C1 to WB first

C1 and C2 has to wait for CO to conform to the coherence order

Problem: Coherence effect on WC





C1 performs the WB

Problem: Coherence effect on WC

Only then, CO receives data and finishes its request





wait	wait	WB(A)	wait	wait		Rx(A)	wait
	t+5	500		t+65	0		t+8

In next period, C0 performs the WB

Only then, C1 receives data and finishes its request





wait	wait	WB(A)	wait	wait		Rx(A)	wait		WB(A)	wait		
	t+5	500		t+6.	50		t+8	300		t·	+950	
		Fina	ally, C	1 per	form	s the '	WB	Af	terwa	rds,	C2	receiv

and finishes its request





wait	wait	WB(A)	wait	wait	Rx(A)	wait	WB(A)	wait	
	t+5	500		t+650		t+800	0	t+950	
Coherence Latency									







	wait	wait	WB(A)	wait	wait	Rx(A)	wait	WB(A)	wait		
		t+5	500		t+650		t+80	00	t+95	50	
Coherence Latency											









- Adopting Traditional Arbitration for coherence results in coupling coherence requests with responses
 - This results in huge WCLs of PMSI
 - It also requires the **coherence modifications** to conform to this arbitration, while maintaining predictable behavior
- Can we achieve predictable coherence without imposing coherence modifications?
 - Yes \rightarrow by deploying a predictable split-Tx bus

PISCOT Main Idea







- High performance arbitration
 - FCFS, Split-Tx, Priority-Based, ...
- •Shared data is entirely supported through a hardware cache coherence protocol 💜
- •Not Predictable X

PISCOT

- Predictable split-transaction bus
- Shared data is supported
- With **NO** protocol modifications



Also achives predictability, with tight latency bounds

PISCOT vs. State-of-the-art



Traditional Real-Time:

- Predictable by design
 - TDM, RR, WRR, HKK, ...
- Does not support shared data, put aside coherence

Data-Aware Arbitration:

- Builds on traditional arbitration
- Shared data is supported
- •BUT requires coherence protocol modifications
- •Coherence leads to good performance
- •Also achives predictability, but with significant latency bounds
- •PMSI [RTAS'17], CARP [RTSS'19], HourGlass [Arxiv'18], PENDULUM [RTSS'19]









Unlike COTS systems, the Request Bus Arbiter is a work-conserving TDM arbiter



PISCOT High-Level Architecture

(1) PISCOT implements a splittransaction bus:

- A Request Bus is responsible for broadcasting the coherence messages initiating memory requests, while
- A Response Bus transfers data as a response to these requests.
 Both buses operate fully in parallel and synchronize using a Service Queue





(1) Split-Bus PISCOT implements a split-transaction bus: A Request Bus is responsible for broadcasting the coherence messages initiating memory requests, while A Response Bus transfers data as a response to these requests.



PISCOT High-Level Architecture

PISCOT

Shared Memory



- A Request Bus is responsible for broadcasting the coherence messages initiating memory requests, while
- A Response Bus transfers data as a response to these requests.

(2) Request Arbiter Unlike COTS systems, the Request Bus Arbiter is a work-conserving TDM arbiter → To achieve **PREDICTABILITY**



PISCOT High-Level Architecture





- A Request Bus is responsible for broadcasting the coherence messages initiating memory requests, while
- A Response Bus transfers data as a response to these requests.

(2) Request Arbiter Unlike COTS systems, the Request Bus Arbiter is a work-conserving TDM arbiter \rightarrow To achieve **PREDICTABILITY**



PISCOT High-Level Architecture

(3) **Response Arbiter**

The Response Bus's arbiter is FCFS:

serves requests based on their arrival time on the Service Queue









- A Request Bus is responsible for broadcasting the coherence messages initiating memory requests, while
- A Response Bus transfers data as a response to these requests.

(2) Request Arbiter Unlike COTS systems, the Request Bus Arbiter is a work-conserving TDM arbiter \rightarrow To achieve **PREDICTABILITY**



PISCOT High-Level Architecture

(3) Response Arbiter

The Response Bus's arbiter is FCFS:

serves requests based on their arrival time on the Service Queue

(4) Service Queue

Both buses operate fully in parallel and synchronize using a Service Queue







- A Request Bus is responsible for broadcasting the coherence messages initiating memory requests, while
- A Response Bus transfers data as a response to these requests.

(2) Request Arbiter Unlike COTS systems, the Request Bus Arbiter is a work-conserving TDM arbiter \rightarrow To achieve **PREDICTABILITY**



PISCOT High-Level Architecture

Shared Memory

(3) Response Arbiter

The Response Bus's arbiter is FCFS:

serves requests based on their arrival time on the Service Queue

(4) Service Queue

Both buses operate fully in parallel and synchronize using a Service Queue

(5) Bounding Interference

- Allows cores to issue multiple \bullet outstanding requests.
- However, to limit coherence \bullet interference, it only services at most one request from any given core at a time.

PISCOT





C	1→Mem: WB(A)	Mem→C0: Send(A)		
t+8	t+	58 t+	+108	

Service Queue			CO: St(A)	C1: PutM(A)
			1	



EXAMPLE

- •our request under analysis is St(A) from C2
- •Assume initially A is modified by C1 in its private cache
- •all 3 cores request to access same cache line A
- Assume same example from before









Service Queue		C1: St(A)	C0: PutM(A)	C0: St(A)	C1: PutM(A)

C1	→Mem: M WB(A) S	em→C0: Send(A)	
t+8	t+58	t+108	t+208



EXAMPLE

- •our request under analysis is St(A) from C2
- •Assume initially A is modified by C1 in its private cache
- •all 3 cores request to access same cache line A
- Assume same example from before











Sanvica Oulous	C2:	C1:	C1:	C0:	C0:	C1:
Service Queue	St(A)	PutM(A)	St(A)	PutM(A)	St(A)	PutM(A)
						· · · · ·

	C1→Mem WB(A)	: Mem→C Send(A	CO: \)	
t	+8	t+58	t+108	t+208



EXAMPLE

- •our request under analysis is St(A) from C2
- •Assume initially A is modified by C1 in its private cache
- •all 3 cores request to access same cache line A
- Assume same example from before











Samica Outour	C2:	C1:	C1:	C0:	C0:	C1:
Service Queue	 St(A)	PutM(A)	St(A)	PutM(A)	St(A)	PutM(A)

			C1→Me WB(A)	m: Men) Ser	n→C0: nd(A)	
			t+8	t+58	t+108	t+208
A: IM ^{ad}	A: IM ^d	A: IM ^d I	A: IM ^d I		A: M->I	
A: M	A: I	A: IM ^d	A: IM ^d I		A: IM ^d I	
A: IM ^{ad}	A: IM ^{ad}	A: IM ^{ad}	A: IM ^d		A: IM ^d	
	t+8	t+12	t+16		t+108	

EXAMPLE

- •our request under analysis is St(A) from C2
- •Assume initially A is modified by C1 in its private cache
- •all 3 cores request to access same cache line A
- •Assume same example from before











	C2:	C1:	C1:	C0:	C0:	C1:
Scivice Queue	St(A)	PutM(A)	St(A)	PutM(A)	St(A)	PutM(A)

			C1→Mem: WB(A)	Mem→C0: Send(A)	C0→Mem: WB(A)	Mem→C1: Send(A)	
			t+8 t-	+58 t+1	08	t+2	08
A: IM ^{ad}	A: IM ^d	A: IM ^d I	A: IM ^d I	A: M	->	A:	
A: M	A: I	A: IM ^d	A: IM ^d I	A: IN	۸d۱	A: M	->
A: IM ^{ad}	A: IM ^{ad}	A: IM ^{ad}	A: IM ^d	A: IN	√lq	A: IN	Md
	t+8	t+12	t+16	t+10)8	t+20)8

EXAMPLE

- •our request under analysis is St(A) from C2
- •Assume initially A is modified by C1 in its private cache
- •all 3 cores request to access same cache line A
- Assume same example from before











Service Queue	C2: C1:	C1: C0: C0:	C1:
	St(A) PutM(A)	St(A) PutM(A) St(A)	PutM(A)
C1→Mem: Mem→C0:	CO \rightarrow Mem: Mem \rightarrow C1: C1 \rightarrow I	Mem: Mem→C2:	
WB(A) Send(A)	WB(A) Send(A) WB	B(A) Send(A)	
t+8 t+58 t+10	18 t+208	t+308	

			C1→Mem WB(A)	n: Mem→C0: Send(A)	C0→Mem: WB(A)	Mem→C1: C1 Send(A) \
			t+8	t+58 t+1	08	t+208
A: IM ^{ad}	A: IM ^d	A: IM ^d I	A: IM ^d I	A: M	->	A: I
A: M	A: I	A: IM ^d	A: IM ^d I	A: IN	۸d۱	A: M->I
A: IM ^{ad}	A: IM ^{ad}	A: IM ^{ad}	A: IM ^d		Mq	A: IM ^d
	t+8	t+12	t+16	t+1(08	t+208



- •Assume same example from before
- •all 3 cores request to access same cache line A
- Assume initially A is modified by C1 in its private cache
- •our request under analysis is St(A) from C2

EXAMPLE









EXAMPLE









Total Latency of 308 cycles compared to 1150 for PMSI

Generally, improves coherence interference from αN^2 to αN





Predictability and Latency Bounds

- Four-core system and figures are for SPLASH-3 and synthetics workloads.
- PISCOT shows more than 4× improvement in both the analytical and observed WCL compared to PMSI

Evaluation

https://gitlab.com/FanosLab/piscot-and-msi-split-bus





EVALUATION



Predictability and Latency Bounds

- Four-core system (Fig is for Ocean BM)
- PISCOT is able to avoid the unpredictability/ variability behavior of conventional MSI

Evaluation https://gitlab.com/FanosLab/piscot-and-msi-split-bus





400









Average-Case Performance

- Slowdown compared to conventional MSI
- PMSI's slowdown is 2× on average (and up to 4.3×) across all benchmarks.
 - coupling of coherence and data transfer
- PISCOT achieves comparable results with slowdown in the range of 1% – 4%

Evaluation https://gitlab.com/FanosLab/piscot-and-msi-split-bus







- High performance arbitration
 - FCFS, Split-Tx, Priority-Based, ..
- •Shared data is entirely supported through a hardware cache coherence protocol 👐
- •Not Predictable 🗙

PISCOT

- Predictable split-transaction bus
- •Shared data is supported **V**
- •With NO protocol modifications
- •Achieves better performance than existing coherence solutions 📢
- •Also achives predictability, with tight latency bounds

- Predictable by design V
- coherence 🗙

Data-Aware Arbitration:

- Builds on traditional arbitration
- •Shared data is supported
- BUT requires coherence protocol modifications
- •Coherence leads to good performance
- •Also achives predictability, but with significant latency bounds
- •PMSI [RTAS'17], CARP [RTSS'19], HourGlass [Arxiv'18], PENDULUM [RTSS'19]

PISCOT vs. State-of-the-art

Summary

Traditional Real-Time: • TDM, RR, WRR, HRR, ... • Does not support shared data, put aside







Conventional COTS: •High performance arbitration •FCFS, Split-Tx, Priority-Based, .. •Shared data is entirely supported through

 Shared data is entirely supported throu hardware cache coherence protocol
 Not Predictable

PISCOT

Predictable split-transaction bus
Shared data is supported
With *NO* protocol modifications
Achieves better performance than existing coherent solutions
Also achives predictability, with tight latency bounds

 Does not support shared data, put aside coherence
 Coherence
 Coherence
 Coherence
 Coherence
 Coherence
 Coherence
 Coherence leads to good performance
 Coherence leads to good performance
 Also achives predictability, but with significan
 PMSI [RTAS'17], CARP [RTSS'19], Hour
 PENDULUM (RTSS'19]

Traditional Real-Time:

Predictable by design 💜

• TDM, RR, WRR, HRR, .

PISCOT vs. State-of-the-art

(1) Split-Bus

PISCOT implements a split-transaction bus:

- A Request Bus is responsible for broadcasting the coherence messages initiating memory requests, while
- A Response Bus transfers data as a response to these requests.

(2) Request Arbi
 Unlike COTS system
 Bus Arbiter is a wo
 TDM arbiter
 → To achieve PREL



PISCOT High-Level Architec

Summary

iter	(3) Response Arbiter
ork-conserving	The Response Bus's arbiter is FCFS:
	 serves requests based on their arrival
	time on the Service Queue
DICIADILITI	
	(4) Service Queue
	Both buses operate fully in parallel
	and synchronize using a Service
>	Queue
nor	(5) Bounding Interference
Aer	 Allows cores to issue multiple
2 p	outstanding requests.
are	 However, to limit coherence
Sh	interference, it only services at
	most one request from any given
	core at a time.
ture	





46

•High performance arbitration 🗸 • FCFS, Split-Tx, Priority-Based, •Shared data is entirely supported through a hardware cache coherence protocol 📢 •Not Predictable 🗙

PISCOT

- Predictable split-transaction bus
- •Shared data is supported 💜
- •With NO protocol modifications
- •Achieves better performance than existing coherence solutions √
- •Also achives predictability, with tight latency bounds

Traditional Real-Time:

- Predictable by design
- TDM, RR, WRR, HRR, ...
- Does not support shared data, put aside coherence X

Data-Aware Arbitration:

- Builds on traditional arbitration
- •Shared data is supported
- •BUT requires coherence protocol modifications 🗸
- •Coherence leads to good performance
- •Also achives predictability, but with significant latency bounds
- •PMSI [RTAS'17], CARP [RTSS'19], HourGlass [Arxiv'18], PENDULUM [RTSS'19]

PISCOT vs. State-of-the-art

Data Sharing





Summary





•High performance arbitration 🗸 •FCFS, Split-Tx, Priority-Based, •Shared data is entirely supported through a hardware cache coherence protocol 📢 •Not Predictable 🗙

PISCOT

- Predictable split-transaction bus
- •Shared data is supported
- •With NO protocol modifications
- •Achieves better performance than existing coherence solutions √
- •Also achives predictability, with tight latency bounds

Traditional Real-Time:

- Predictable by design V
- TDM, RR, WRR, HRR, ...
- Does not support shared data, put aside coherence X

Data-Aware Arbitration:

- Builds on traditional arbitration
- Shared data is supported
- •BUT requires coherence protocol modifications 🗸
- •Coherence leads to good performance
- •Also achives predictability, but with significant latency bounds
- •PMSI [RTAS'17], CARP [RTSS'19], HourGlass [Arxiv'18], PENDULUM [RTSS'19]

PISCOT vs. State-of-the-art

Data Sharing

(1) Split-Bus



PISCOT High-Level Architecture



Summary



(b) PISCOT





PISCOT

solutions 📈

•High performance arbitration 🗸 •FCFS, Split-Tx, Priority-Based, •Shared data is entirely supported through a hardware cache coherence protocol 📢 •Not Predictable 🗙

• Predictable split-transaction bus

•Shared data is supported

•With NO protocol modifications

Traditional Real-Time:

- Predictable by design
- TDM, RR, WRR, HRR, ...
- Does not support shared data, put aside coherence X

Data-Aware Arbitration:

- Builds on traditional arbitration
- •Shared data is supported
- •BUT requires coherence protocol modifications 🗸
- Coherence leads to good performance
- •Also achives predictability, but with significant latency bounds
- •PMSI [RTAS'17], CARP [RTSS'19], HourGlass [Arxiv'18], PENDULUM [RTSS'19]

PISCOT vs. State-of-the-art

•Achieves better performance than existing coherence

•Also achives predictability, with tight latency bounds

Data Sharing



PISCOT High-Level Architecture



Summary



(b) PISCOT

SUMMARY

