Token Coherence:

Enabling Faster Multiprocessors by Decoupling Performance and Correctness

Milo M. K. Martin, Mark D. Hill, David A. Wood

Wisconsin Multifacet Project http://www.cs.wisc.edu/multifacet/ University of Wisconsin—Madison

Refs: ISCA'03, Martin Thesis, Martin Interview Talk

(C) 2003 Milo Martin











Hardware Cache-Coherence Problem Goal: provide a consistent view of memory Permissions in each cache per block One read/write -orMany readers Cache coherence protocols

- Distributed & complex
- Correctness critical
- Performance critical

slide 9

• Races: the main source of complexity - Requests for the **same block** at the **same time**

Token Coherence – Milo Marti



Token Coherence – Milo Mar

- · Motivation & background
- Problem: resolving races
 - Traditional solutions
 Performance overheads
- Solution: Token Coherence (specific)
- Evaluation
- Token Coherence (generalized)
- Conclusions

slide 10



















Token Coher























- · As before:
 - Broadcast with direct responses (like snooping)
 - Use unordered interconnect (like directory)
- Track tokens for safety

Reissue requests as needed

- Needed due to racing requests (uncommon)
- Timeout to detect failed completion
 - Wait twice average miss latencySmall hardware overhead
- (Ignore starvation for a moment)

















More Depth in ISCA '03 & Martin Thesis

- Traffic optimizations
 - Transfer tokens without data
 - Add an "owner" token
 - Upgrade (Read-only to read/write)
 - "Exclusive Clean" State
- Note: no silent read-only replacements
 Worst case: 10% traffic overhead
- Encoding tokens in memory

 Using ECC bits

slide 39

lido 41

- Reduce read/modify/writes with token cache

Outline Motivation & background

Token Coherence – Milo M

- Problem: resolving races
- Solution: Token Coherence (specific)
- Evaluation
- Token Coherence (generalized)
- Conclusions

Evaluation Goal

- Non-goal: exact speedup numbers
 - Many parameters and assumptions
 - Key parameter: **16 processors**
- Goal: Four Questions
 - 1. Are races rare?
 - 2. Can Token Coherence outperform Snooping?
 - 3. Can Token Coherence outperform Directories?
 - 4. Is broadcast overhead reasonable?

Quantitative evidence for qualitative behavior

Token Coherence – Milo Mart

Token Coherence – Milo Martin

OLTP	SPECjbb	Apache

Q1: F (p	Reissue	d Reques all misses)	sts				
Outcome	OLTP	SPECjbb	Apache				
Not Reissued	98%	98%	96%				
Reissued Once	2%	2%	3%				
Reissued > 1	0.4%	0.3%	0.7%				
Persistent Requests (Reissued > 4)	0.2%	0.1%	0.3%				
Y	Yes; races are rare						
slide 43	slide 43 Token Coherence – Milo Martin						

























Performance Protocols Opportunities - Aggressively target the common case - Requests are just "hints" to move data & tokens · Correctness substrate may ignore hints Robust - Can't cause "correctness" violations - A null or random protocol is correct - Rely on correctness substrate

Verifiability & Complexity

- · Divide and conquer complexity
 - Formal verification is future work - Difficult to quantify, but promising
 - E.g. simple replacements (no handshake)
- · Strong invariants - Locally enforced with tokens
 - Response-centric; independent of requests
 - Prevent data-corruption bugs
- Explicit forward progress Simple mechanism

slide 57

elido 50

Further innovation → no correctness worries

Outline

Token Coherence – Milo M

Token Coherence – Milo M

- Motivation & background
- · Problem: resolving races
- Solution: Token Coherence (specific)
- Evaluation
- Token Coherence (generalized)
- · Other research & future directions
- Conclusions

slide 58

Token Coherence – Milo N

Token Coherence – Milo I



- Coherence protocols
 - Destination-set prediction [ISCA '03]
 - Bandwidth adaptive coherence [HPCA '02]
 - Timestamp-based ordering [ASPLOS '00]
- Correctness & specification
 - Value prediction & consistency [Micro '01]
 - Protocol specification [TPDS '02]
 - SLICC: A domain specific language for protocol specification
- Improving hardware availability [ISCA '02]
- Compiler/hardware interaction [Micro '97]

Tackling complexity through decoupling - Simpler, faster, cheaper, more robust • Other hardware examples

Future Directions: Servers

- Hierarchical multiprocessors (CMPs)
- Multiprocessor interconnection networks
- (e.g., domain specific TCP/IP)
- Simplify processor design with "correctness" checker

Token Coh

 Possible software example (for collaboration) - Content distribution networks, secure systems

Future Directions: Clients

- Mobile, wireless, embedded, sensors – Not PCs
- Energy consumption, cost, performance – Applying techniques from servers to clients
 - Exploit explicit parallelism
- Example: Intel Xscale

slide 61

- 1 GHz \rightarrow 200 MHz reduces energy used by 30x
- 5 x 200 MHz in parallel, use 1/6th the energy
- Same research approach:
 - Identify emerging workloads, trends, problems
 - Build tools, characterize, find innovate solutions

Token Coherence – Milo Marti

Conclusions

- Token Coherence (generalization)
- Decouple correctness from performance
 - Correctness substrate
 - Tokens for safety
 - Persistent requests for forward progress
 - Performance protocol for performance
- Token Coherence (broadcast version)

 Low sharing-miss latency (no indirection)
 - Requires no interconnect ordering
 - Can be faster than current alternatives
- Enables further protocol innovation

Token Coherence – Milo Marti