memory for large-scale systems

Prakhar Jawre

Outline

• Overview

Current Technologies

Challenges

Conclusion





Rapid growth in digital data!



Rapid growth in digital data!

From Computing Centric to Data Centric



Rapid growth in digital data!

From Computing Centric to Data Centric

Data growth in 10 year(by 2015) = 100X

Today's Data Centers

- Per-vendor Layer
- Product Based
- Near neighbor optimization
- Big vendors(Google, Amazon etc)
- Can do Cross Layer Optimization but,
 - Limited to service of Interest
 - Limited to Extend(Software)
 - Closed Technologies



Outline

• Overview

Current Technologies

Challenges

Conclusion

In-Memory Computing for high Performance



In-Memory Computing for high Performance

Why: Tight Latency Constraints How: Place data in main memory



In-Memory Computing for high Performance

Why: Tight Latency Constraints How: Place data in main memory



Need large memory pool to accommodate all data

Nodes Frequently Access Remote Memory



Nodes Frequently Access Remote Memory



Nodes Frequently Access Remote Memory



Need fast access to both small and large objects in remote memory

Rack-Scale Systems: Fast Access to Big Memory















Rack-Scale Systems: Fast Access to Big Memory

- Vast memory pool in small form factor
- On-chip integrated, cache-coherent NIs
 - Examples: Scale-Out NUMA [Novakovic et al., '14]



Large memory capacity, low latency, high bandwidth

Rack-Scale Systems: Fast Access to Big Memory

- Vast memory pool in small form factor
- On-chip integrated, cache-coherent NIs
 - Examples: Scale-Out NUMA [Novakovic et al., '14]
- High-performance inter-node interconnect
 - NUMA environment
 - Low latency for fine-grained transfers
 - High bandwidth for bulk transfers



Large memory capacity, low latency, high bandwidth

Remote Direct Memory Access (RDMA)

- Frequent remote memory accesses within datacenter
- Networking traditionally expensive



Remote Direct Memory Access (RDMA)

- Frequent remote memory accesses within datacenter
- Networking traditionally expensive



Remote Direct Memory Access (RDMA)

- Frequent remote memory accesses within datacenter
- Networking traditionally expensive
- RDMA technology
- Fast one-sided operations
- Destination CPU not involved



Outline

- Overview
- Current Technologies
- Challenges
 - Atomic Object Read
 - Network Interface for Many-Core Chip
- Proposed Solutions
- Conclusion

Atomicity limited to single cache line in one sided ops



Atomicity limited to single cache line in one sided ops



Atomicity limited to single cache line in one sided ops



Atomicity limited to single cache line in one sided ops



Atomicity limited to single cache line in one sided ops

Software construct for object atomicity
Embedded MetaData



□ SW check expensive=>upto 50% of remote object read



















Can we do this using hardware?






















Speculatively read version and data \rightarrow remove serial read latency



Speculatively read version and data \rightarrow remove serial read latency



Speculatively read version and data \rightarrow remove serial read latency

Simple hardware, atomicity w/ zero latency overhead



Up to 2x latency & throughput benefit for distributed object stores



 Up to 2x latency & throughput benefit for distributed object stores

Outline

• Overview

Current Technologies

Challenges

Atomic Object Read

Network Interface for Many-Core Chip

Conclusion

Interwork Interface Design for Manycore Chips

In I placement & design is key for remote access performance Obvious NI designs suffer from poor latency or bandwidth

Contributions

In the second state of the second

Goal: "Low-latency, high-bandwidth NI design for manycore chips"

RDMA-like Queue-Pair (QP) model

Cores and NIs communicate through *cacheable memorymapped* queues



RDMA-like Queue-Pair (QP) model

Cores and NIs communicate through *cacheable memorymapped* queues



RDMA-like Queue-Pair (QP) model

Cores and NIs communicate through *cacheable memorymapped* queues



RDMA-like Queue-Pair (QP) model

Cores and NIs communicate through *cacheable memorymapped* queues



RDMA-like Queue-Pair (QP) model

Cores and NIs communicate through cacheable memorymapped queues



Implications of Manycore Chips on Remote Access

More Cores \rightarrow Higher request rate NI capabilities have to match chip's communication demands

- Approaches \rightarrow
- Scale NI across edge
- $\Box \rightarrow Close to network pins$
- → Access to NOC's full bisection bandwidth
- $\beth \rightarrow \bot$ Large average core to NI distance $ar{arepsilon}$

- 1
 (
 1
 (
 1
 (

 1
 (
 1
 (
 1
 (
 1

 1
 (
 1
 (
 1
 (
 1

 1
 (
 1
 (
 1
 (
 1

 1
 (
 1
 (
 1
 (
 1

 1
 (
 1
 (
 1
 (
 1

 1
 (
 1
 (
 1
 (
)
- Collocate NI logic per core, to localize all interactions
- $\Box \rightarrow$ Transparent to coherence mechanism
- $\Box \rightarrow$ No modifications to core's IP block





















QP interactions account for up to 50% of end-to-end latency












Per-Core NI 4-Cache-Block Remote Read



Minimized latency but bandwidth misuse for large requests











Data handling









Both low-latency small transfers and high-bandwidth large transfers





Outline

• Overview

Current Technologies

Challenges

Conclusion

Conclusion

- □ Software construct for object atomicity incur heavy overheads.
- Need for hardware solutions to atomic object reads.
- □ NI design for manycore chips is crucial to remote memory access.
- □ Fast core-NI interaction critical to achieve low latency
- □ Large transfers best handled at the edge
- □ Split NI: low-latency, high-bandwidth remote memory access

Reference

- Alexandros Daglis, Stanko Novaković, Edouard Bugnion, Babak Falsafi, and Boris Grot. 2015. Manycore network interfaces for in-memory rack-scale computing. SIGARCH Comput. Archit. News 43, 3 (June 2015)
- Daglis, A, Ustiugov, D, Novakovic, S, Bugnion, E, Falsafi, B & Grot, B 2016, SABRes: Atomic Object Reads for In-Memory Rack-Scale Computing. in In Proceedings of the 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO 2016)
- https://parsa.epfl.ch/~falsafi/talks/specialized-server.pdf