

Computer Architecture: Interconnects (Part I)

Prof. Onur Mutlu
Carnegie Mellon University

Memory Interference and QoS Lectures

- These slides are from a lecture delivered at Bogazici University (June 10, 2013)
- Videos:
 - <https://www.youtube.com/playlist?list=PLVngZ7BemHHV6N0ejHhwOfLwTr8Q-UKXj>
 - <https://www.youtube.com/watch?v=rKRFlpavGs8&list=PLVngZ7BemHHV6N0ejHhwOfLwTr8Q-UKXj&index=3>
 - <https://www.youtube.com/watch?v=go34f7v9KIs&list=PLVngZ7BemHHV6N0ejHhwOfLwTr8Q-UKXj&index=4>
 - <https://www.youtube.com/watch?v=nWDoUqoZY4s&list=PLVngZ7BemHHV6N0ejHhwOfLwTr8Q-UKXj&index=5>

Multi-Core Architectures and Shared Resource Management

Lecture 3: Interconnects

Prof. Onur Mutlu

<http://www.ece.cmu.edu/~omutlu>

onur@cmu.edu

Bogazici University

June 10, 2013

SAFARI Carnegie Mellon

Last Lecture

- Wrap up Asymmetric Multi-Core Systems
 - Handling Private Data Locality
 - Asymmetry Everywhere
- Resource Sharing vs. Partitioning
- Cache Design for Multi-core Architectures
 - MLP-aware Cache Replacement
 - The Evicted-Address Filter Cache
 - Base-Delta-Immediate Compression
 - Linearly Compressed Pages
 - Utility Based Cache Partitioning
 - Fair Shared Cache Partitining
 - Page Coloring Based Cache Partitioning

Agenda for Today

- Interconnect design for multi-core systems
- (Prefetcher design for multi-core systems)
- (Data Parallelism and GPUs)

Readings for Lecture June 6 (Lecture 1.1)

- Required – Symmetric and Asymmetric Multi-Core Systems
 - Mutlu et al., “Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-order Processors,” HPCA 2003, IEEE Micro 2003.
 - Suleman et al., “Accelerating Critical Section Execution with Asymmetric Multi-Core Architectures,” ASPLOS 2009, IEEE Micro 2010.
 - Suleman et al., “Data Marshaling for Multi-Core Architectures,” ISCA 2010, IEEE Micro 2011.
 - Joao et al., “Bottleneck Identification and Scheduling for Multithreaded Applications,” ASPLOS 2012.
 - Joao et al., “Utility-Based Acceleration of Multithreaded Applications on Asymmetric CMPs,” ISCA 2013.
- Recommended
 - Amdahl, “Validity of the single processor approach to achieving large scale computing capabilities,” AFIPS 1967.
 - Olukotun et al., “The Case for a Single-Chip Multiprocessor,” ASPLOS 1996.
 - Mutlu et al., “Techniques for Efficient Processing in Runahead Execution Engines,” ISCA 2005, IEEE Micro 2006.

Videos for Lecture June 6 (Lecture 1.1)

■ Runahead Execution

- <http://www.youtube.com/watch?v=z8YpjqXQJIA&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=28>

■ Multiprocessors

- Basics: http://www.youtube.com/watch?v=7ozCK_Mgxfk&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=31
- Correctness and Coherence: <http://www.youtube.com/watch?v=U-VZKMgItDM&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=32>
- Heterogeneous Multi-Core: <http://www.youtube.com/watch?v=r6r2NJxj3kI&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=34>

Readings for Lecture June 7 (Lecture 1.2)

- Required – Caches in Multi-Core
 - Qureshi et al., “A Case for MLP-Aware Cache Replacement,” ISCA 2005.
 - Seshadri et al., “The Evicted-Address Filter: A Unified Mechanism to Address both Cache Pollution and Thrashing,” PACT 2012.
 - Pekhimenko et al., “Base-Delta-Immediate Compression: Practical Data Compression for On-Chip Caches,” PACT 2012.
 - Pekhimenko et al., “Linearly Compressed Pages: A Main Memory Compression Framework with Low Complexity and Low Latency,” SAFARI Technical Report 2013.
- Recommended
 - Qureshi et al., “Utility-Based Cache Partitioning: A Low-Overhead, High-Performance, Runtime Mechanism to Partition Shared Caches,” MICRO 2006.

Videos for Lecture June 7 (Lecture 1.2)

- Cache basics:

- <http://www.youtube.com/watch?v=TpMdBrM1hVc&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=23>

- Advanced caches:

- <http://www.youtube.com/watch?v=TboaFbjTd-E&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=24>

Readings for Lecture June 10 (Lecture 1.3)

■ Required – Interconnects in Multi-Core

- Moscibroda and Mutlu, “A Case for Bufferless Routing in On-Chip Networks,” ISCA 2009.
- Fallin et al., “CHIPPER: A Low-Complexity Bufferless Deflection Router,” HPCA 2011.
- Fallin et al., “MinBD: Minimally-Buffered Deflection Routing for Energy-Efficient Interconnect,” NOCS 2012.
- Das et al., “Application-Aware Prioritization Mechanisms for On-Chip Networks,” MICRO 2009.
- Das et al., “Aergia: Exploiting Packet Latency Slack in On-Chip Networks,” ISCA 2010, IEEE Micro 2011.

■ Recommended

- Grot et al. “Preemptive Virtual Clock: A Flexible, Efficient, and Cost-effective QOS Scheme for Networks-on-Chip,” MICRO 2009.
- Grot et al., “Kilo-NOC: A Heterogeneous Network-on-Chip Architecture for Scalability and Service Guarantees,” ISCA 2011, IEEE Micro 2012.

More Readings for Lecture 1.3

- Studies of congestion and congestion control in on-chip vs. internet-like networks
- George Nychis, Chris Fallin, Thomas Moscibroda, Onur Mutlu, and Srinivasan Seshan,
"On-Chip Networks from a Networking Perspective: Congestion and Scalability in Many-core Interconnects"
Proceedings of the 2012 ACM SIGCOMM Conference (SIGCOMM), Helsinki, Finland, August 2012. Slides (pptx)
- George Nychis, Chris Fallin, Thomas Moscibroda, and Onur Mutlu,
"Next Generation On-Chip Networks: What Kind of Congestion Control Do We Need?"
Proceedings of the 9th ACM Workshop on Hot Topics in Networks (HOTNETS), Monterey, CA, October 2010. Slides (ppt) (key)

Videos for Lecture June 10 (Lecture 1.3)

■ Interconnects

- <http://www.youtube.com/watch?v=6xEpbFVgnf8&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=33>

■ GPUs and SIMD processing

- Vector/array processing basics: <http://www.youtube.com/watch?v=f-XL4BNRoBA&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=15>
- GPUs versus other execution models:
<http://www.youtube.com/watch?v=dl5TZ4-oao0&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=19>
- GPUs in more detail:
<http://www.youtube.com/watch?v=vr5hbSkb1Eg&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=20>

Readings for Prefetching

■ Prefetching

- Srinath et al., "Feedback Directed Prefetching: Improving the Performance and Bandwidth-Efficiency of Hardware Prefetchers," HPCA 2007.
- Ebrahimi et al., "Coordinated Control of Multiple Prefetchers in Multi-Core Systems," MICRO 2009.
- Ebrahimi et al., "Techniques for Bandwidth-Efficient Prefetching of Linked Data Structures in Hybrid Prefetching Systems," HPCA 2009.
- Ebrahimi et al., "Prefetch-Aware Shared Resource Management for Multi-Core Systems," ISCA 2011.
- Lee et al., "Prefetch-Aware DRAM Controllers," MICRO 2008.

■ Recommended

- Lee et al., "Improving Memory Bank-Level Parallelism in the Presence of Prefetching," MICRO 2009.

Videos for Prefetching

■ Prefetching

- <http://www.youtube.com/watch?v=IIkIwiNNI0c&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=29>
- <http://www.youtube.com/watch?v=yapQavK6LUk&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=30>

Readings for GPUs and SIMD

- GPUs and SIMD processing
 - Narasiman et al., “Improving GPU Performance via Large Warps and Two-Level Warp Scheduling,” MICRO 2011.
 - Jog et al., “OWL: Cooperative Thread Array Aware Scheduling Techniques for Improving GPGPU Performance,” ASPLOS 2013.
 - Jog et al., “Orchestrated Scheduling and Prefetching for GPGPUs,” ISCA 2013.
 - Lindholm et al., “NVIDIA Tesla: A Unified Graphics and Computing Architecture,” IEEE Micro 2008.
 - Fung et al., “Dynamic Warp Formation and Scheduling for Efficient GPU Control Flow,” MICRO 2007.

Videos for GPUs and SIMD

- GPUs and SIMD processing

- Vector/array processing basics: <http://www.youtube.com/watch?v=f-XL4BNRoBA&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=15>
- GPUs versus other execution models: <http://www.youtube.com/watch?v=dl5TZ4-oao0&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=19>
- GPUs in more detail: <http://www.youtube.com/watch?v=vr5hbSkb1Eg&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=20>

Online Lectures and More Information

■ Online Computer Architecture Lectures

- <http://www.youtube.com/playlist?list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ>

■ Online Computer Architecture Courses

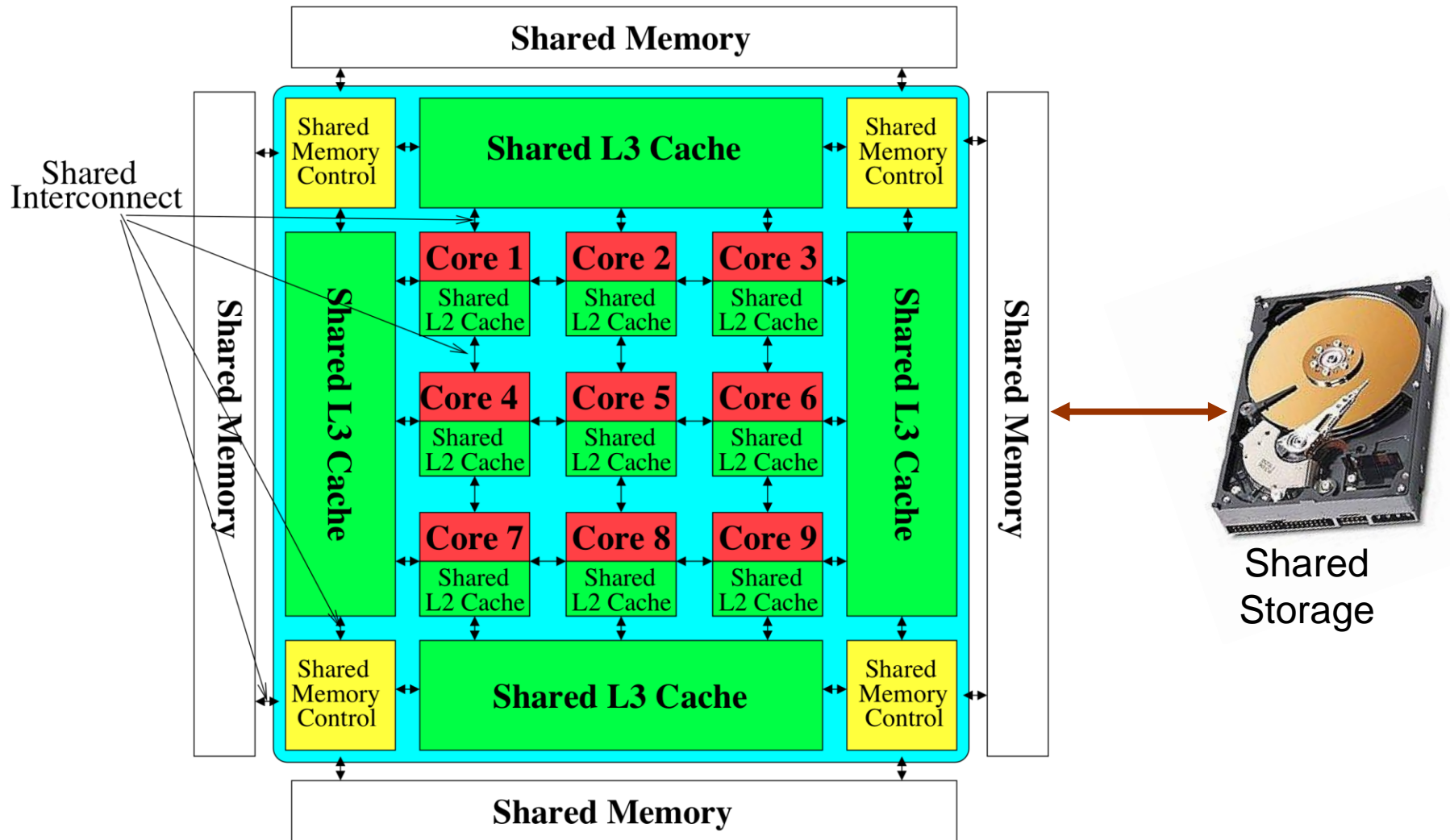
- Intro: <http://www.ece.cmu.edu/~ece447/s13/doku.php>
- Advanced: <http://www.ece.cmu.edu/~ece740/f11/doku.php>
- Advanced: <http://www.ece.cmu.edu/~ece742/doku.php>

■ Recent Research Papers

- <http://users.ece.cmu.edu/~omutlu/projects.htm>
- <http://scholar.google.com/citations?user=7XyGUGkAAAAJ&hl=en>

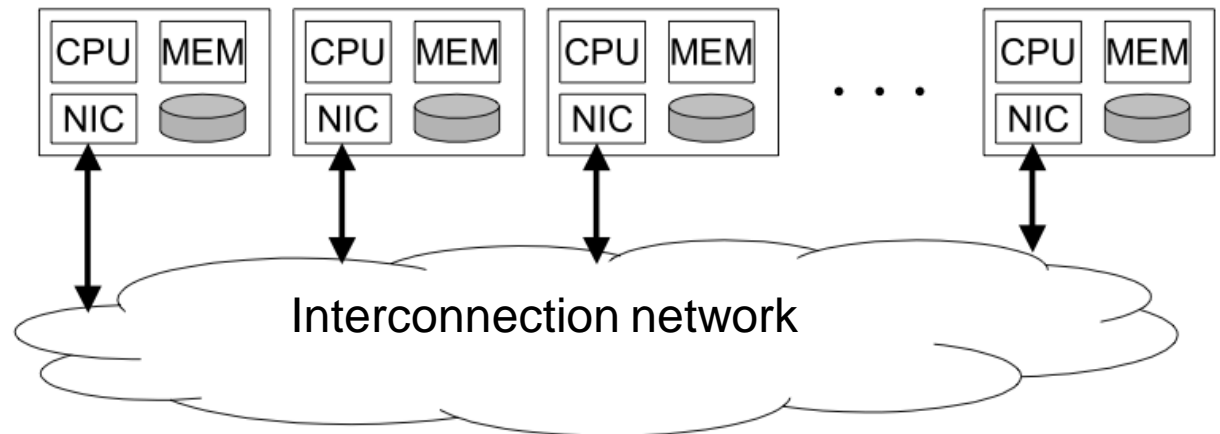
Interconnect Basics

Interconnect in a Multi-Core System



Where Is Interconnect Used?

- To connect components
- Many examples
 - Processors and processors
 - Processors and memories (banks)
 - Processors and caches (banks)
 - Caches and caches
 - I/O devices



Why Is It Important?

- Affects the scalability of the system
 - How large of a system can you build?
 - How easily can you add more processors?
- Affects performance and energy efficiency
 - How fast can processors, caches, and memory communicate?
 - How long are the latencies to memory?
 - How much energy is spent on communication?

Interconnection Network Basics

- Topology
 - Specifies the way switches are wired
 - Affects routing, reliability, throughput, latency, building ease
- Routing (algorithm)
 - How does a message get from source to destination
 - Static or adaptive
- Buffering and Flow Control
 - What do we store within the network?
 - Entire packets, parts of packets, etc?
 - How do we throttle during oversubscription?
 - Tightly coupled with routing strategy

Topology

- Bus (simplest)
- Point-to-point connections (ideal and most costly)
- Crossbar (less costly)
- Ring
- Tree
- Omega
- Hypercube
- Mesh
- Torus
- Butterfly
- ...

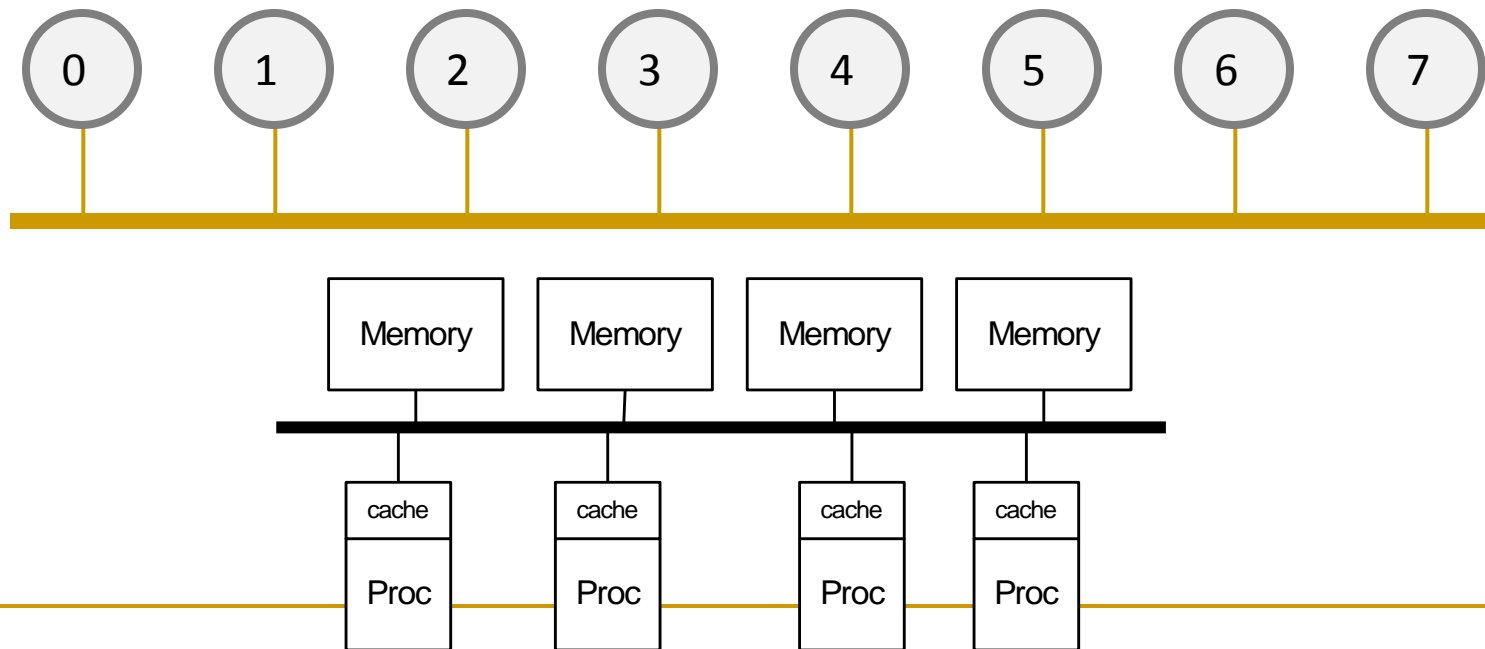
Metrics to Evaluate Interconnect Topology

- Cost
- Latency (in hops, in nanoseconds)
- Contention

- Many others exist you should think about
 - Energy
 - Bandwidth
 - Overall system performance

Bus

- + Simple
- + Cost effective for a small number of nodes
- + Easy to implement coherence (snooping and serialization)
- Not scalable to large number of nodes (limited bandwidth, electrical loading → reduced frequency)
- High contention → fast saturation

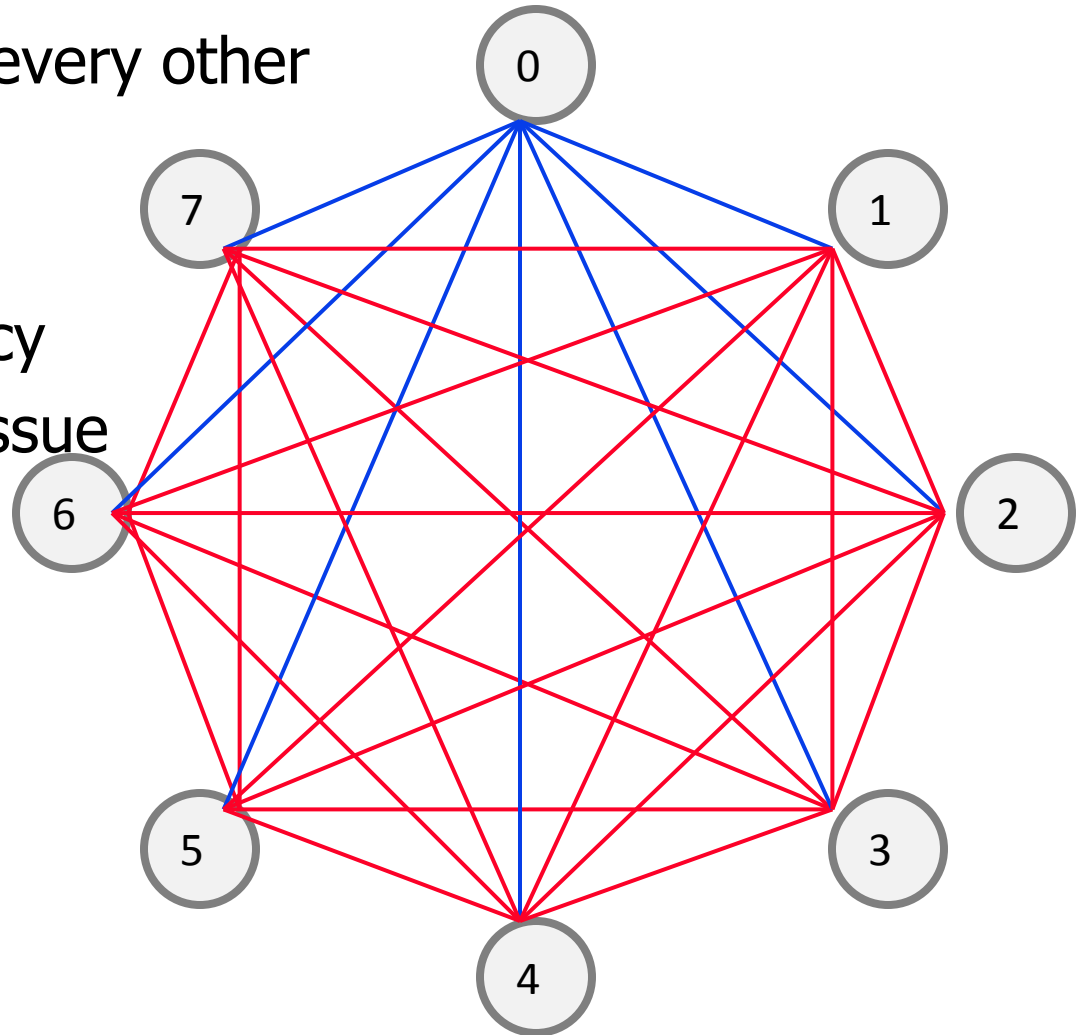


Point-to-Point

Every node connected to every other

- + Lowest contention
- + Potentially lowest latency
- + Ideal, if cost is not an issue

- Highest cost
 - $O(N)$ connections/ports per node
 - $O(N^2)$ links
- Not scalable
- How to lay out on chip?



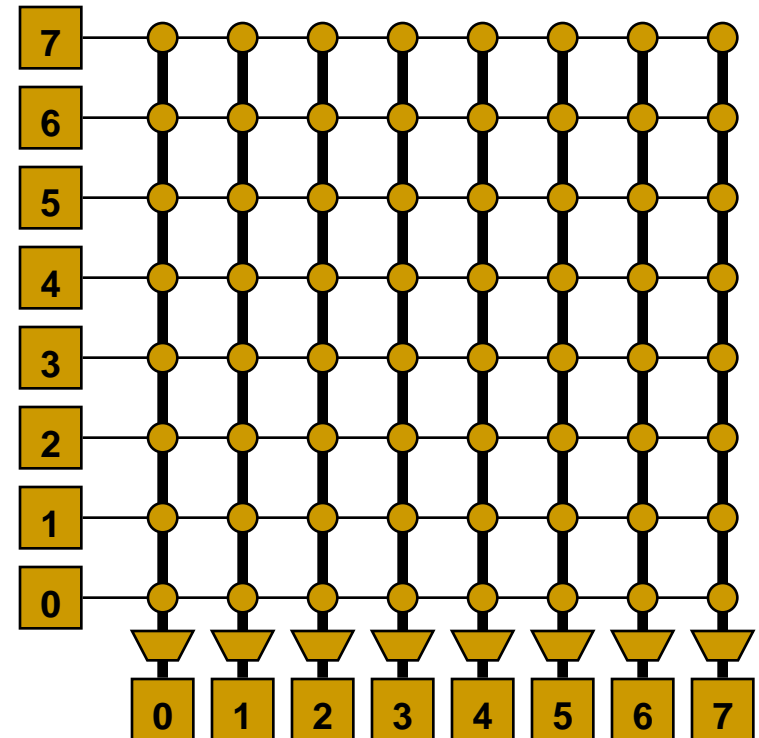
Crossbar

- Every node connected to every other (non-blocking) except one can be using the connection at any given time
- Enables concurrent sends to non-conflicting destinations
- Good for small number of nodes

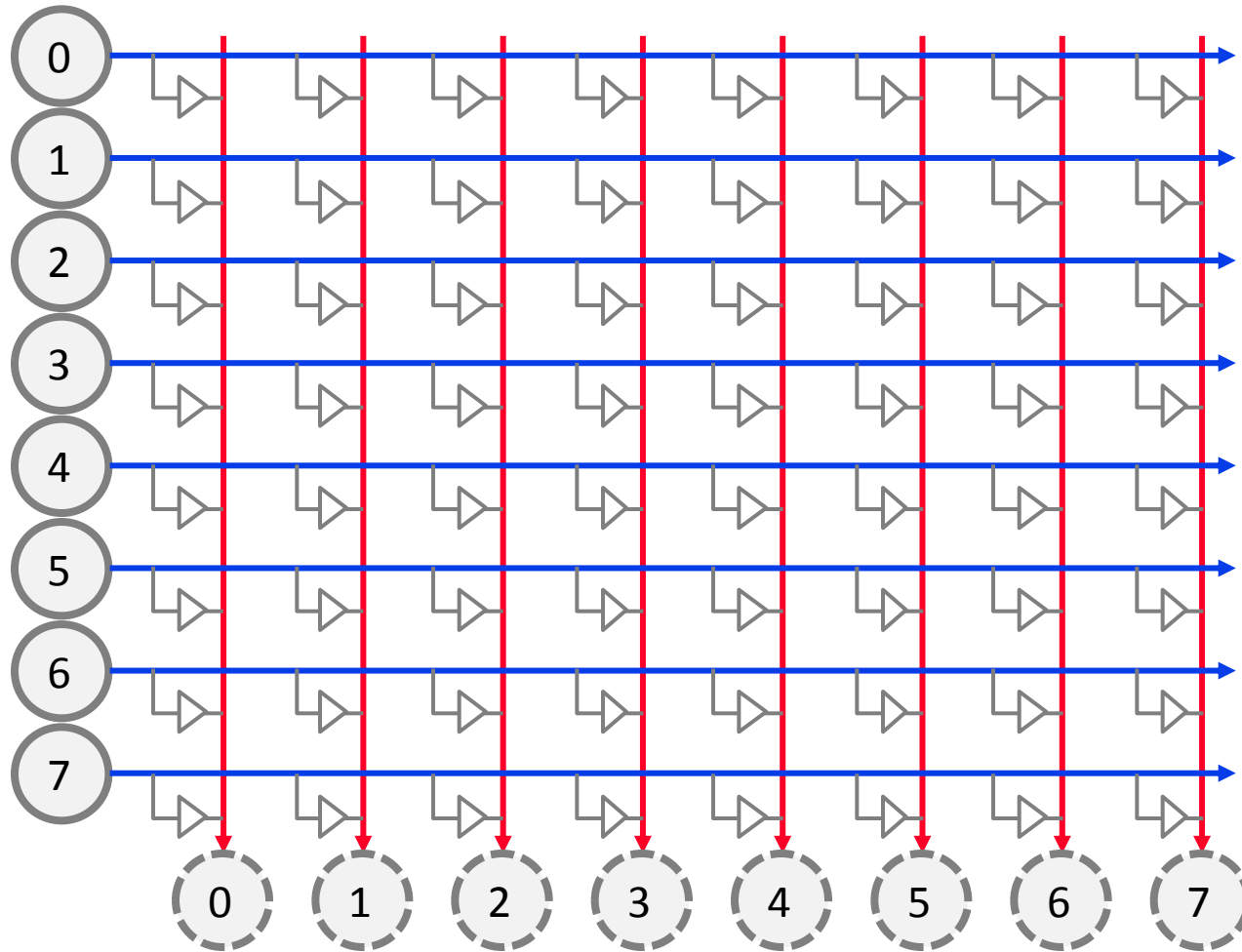
- + Low latency and high throughput
- Expensive
- Not scalable $\rightarrow O(N^2)$ cost
- Difficult to arbitrate as N increases

Used in core-to-cache-bank networks in

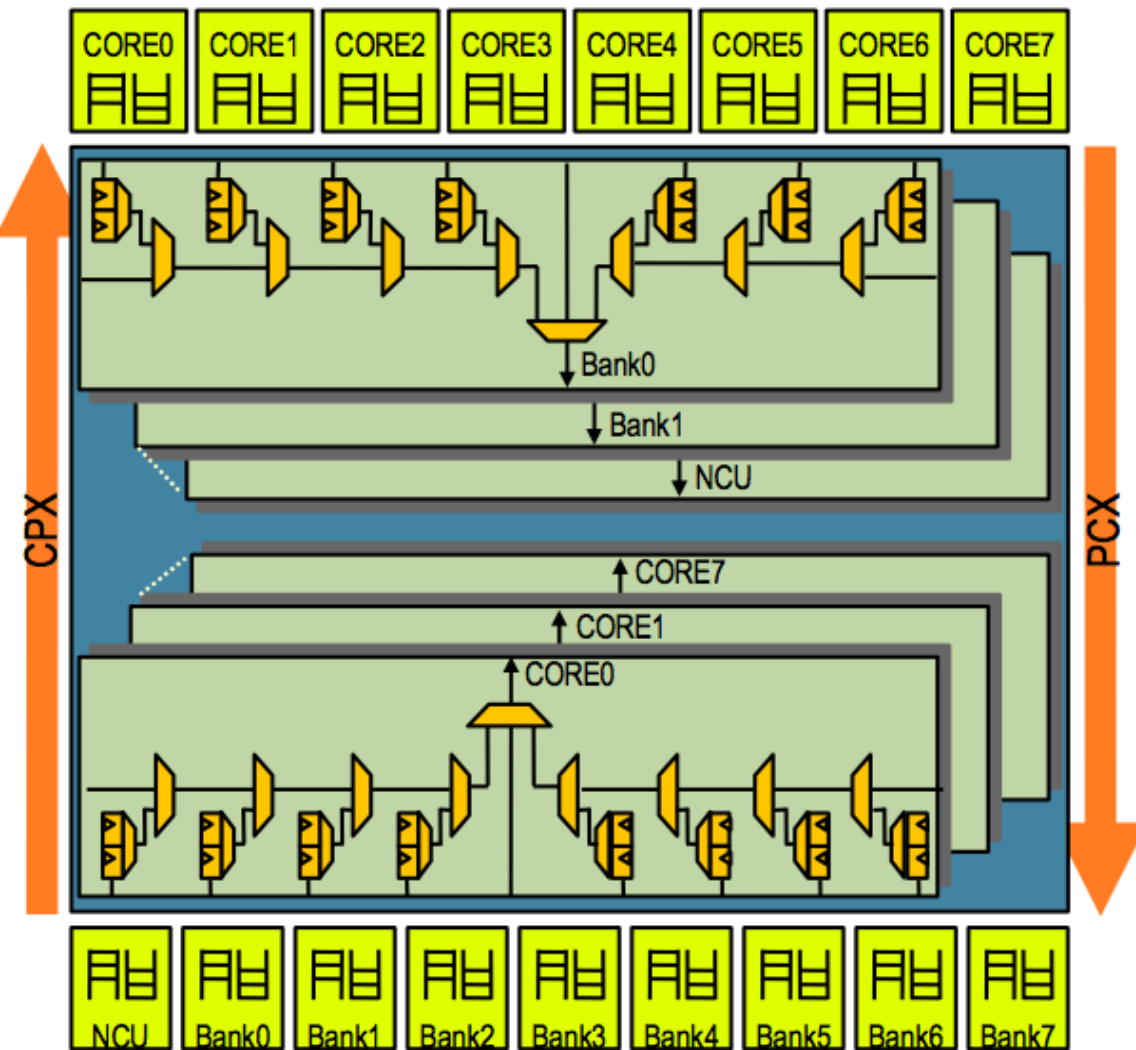
- IBM POWER5
- Sun Niagara I/II



Another Crossbar Design

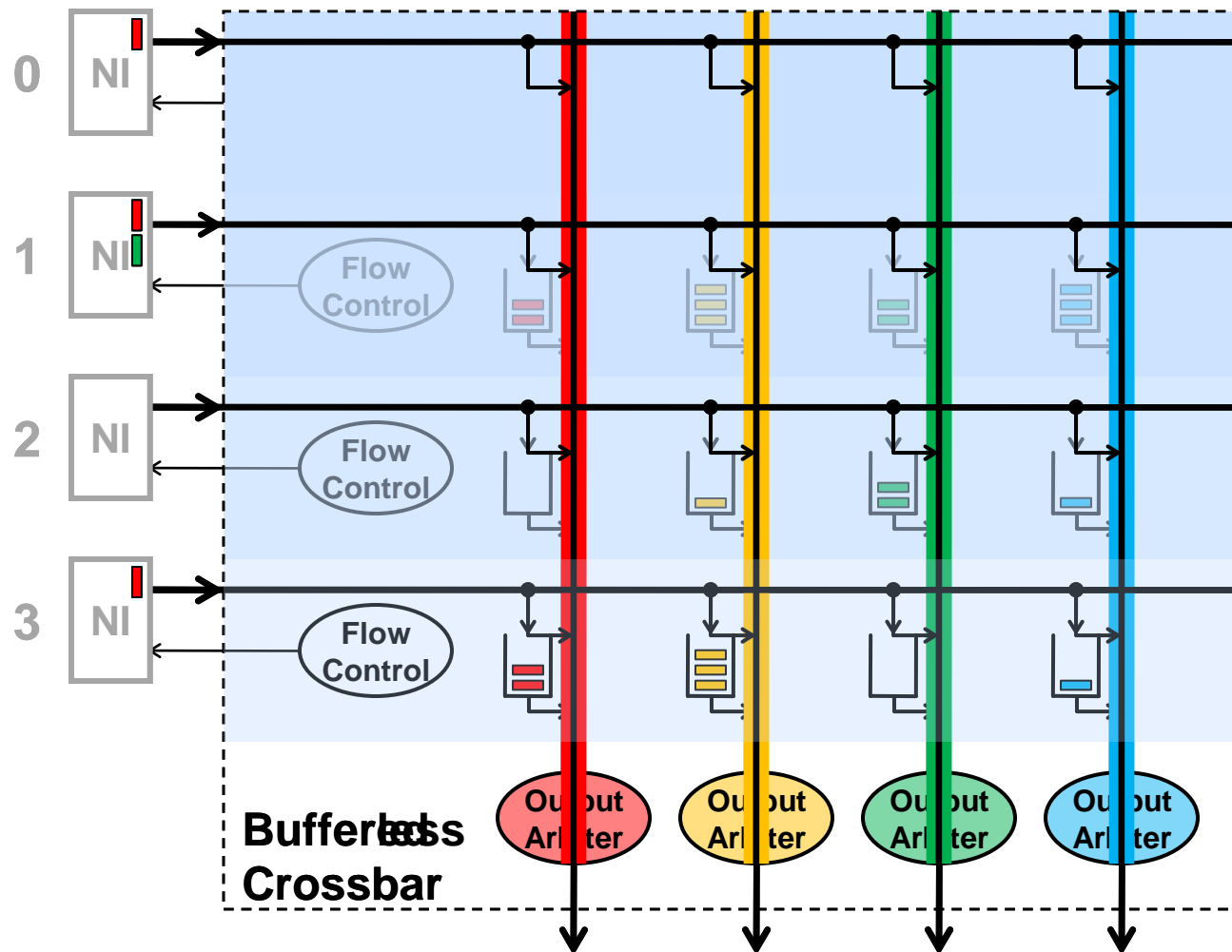


Sun UltraSPARC T2 Core-to-Cache Crossbar



- High bandwidth interface between 8 cores and 8 L2 banks & NCU
- 4-stage pipeline: req, arbitration, selection, transmission
- 2-deep queue for each src/dest pair to hold data transfer request

Buffered Crossbar



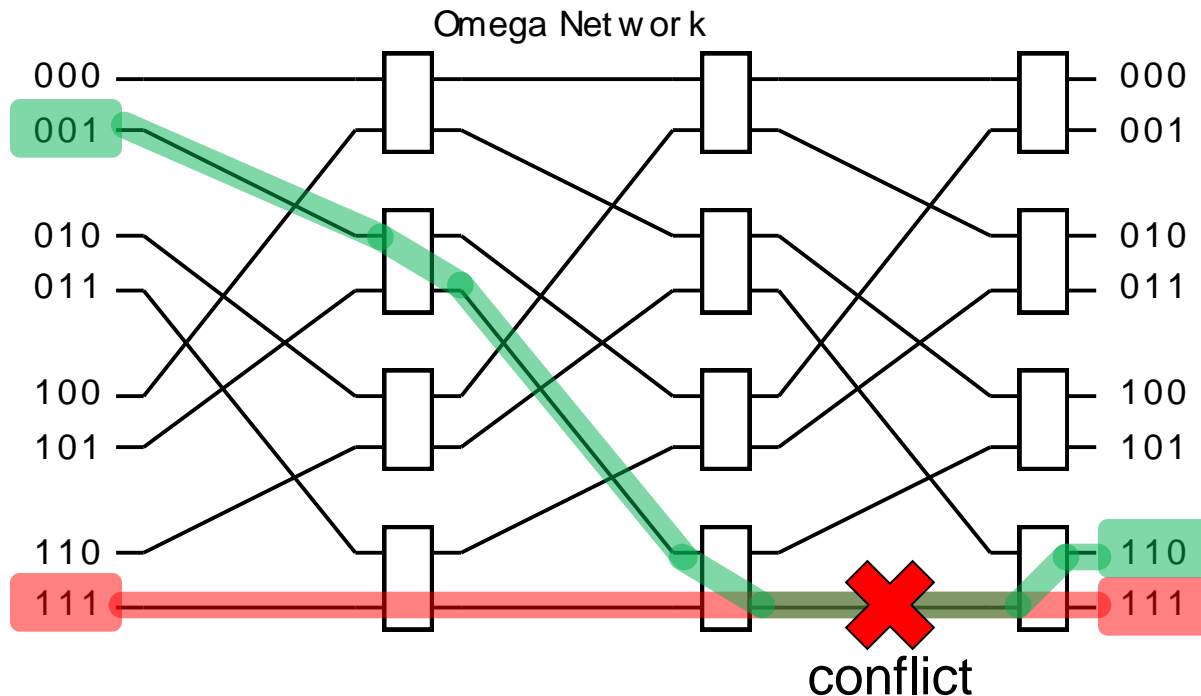
- + Simpler arbitration/scheduling
- + Efficient support for variable-size packets
- Requires N^2 buffers

Can We Get Lower Cost than A Crossbar?

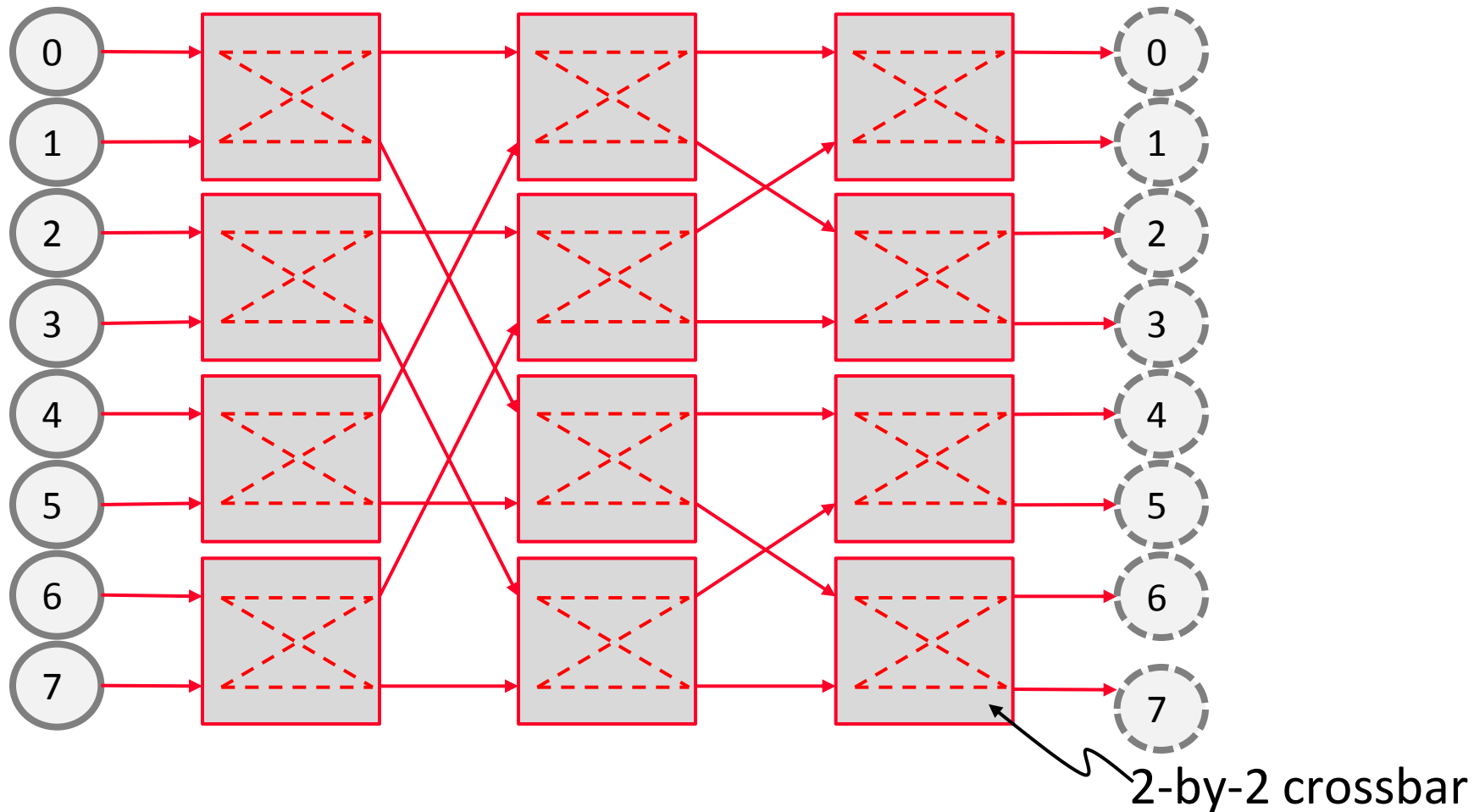
- Yet still have low contention?
- Idea: Multistage networks

Multistage Logarithmic Networks

- Idea: Indirect networks with multiple layers of switches between terminals/nodes
- Cost: $O(N \log N)$, Latency: $O(\log N)$
- Many variations (Omega, Butterfly, Benes, Banyan, ...)
- **Omega Network:**

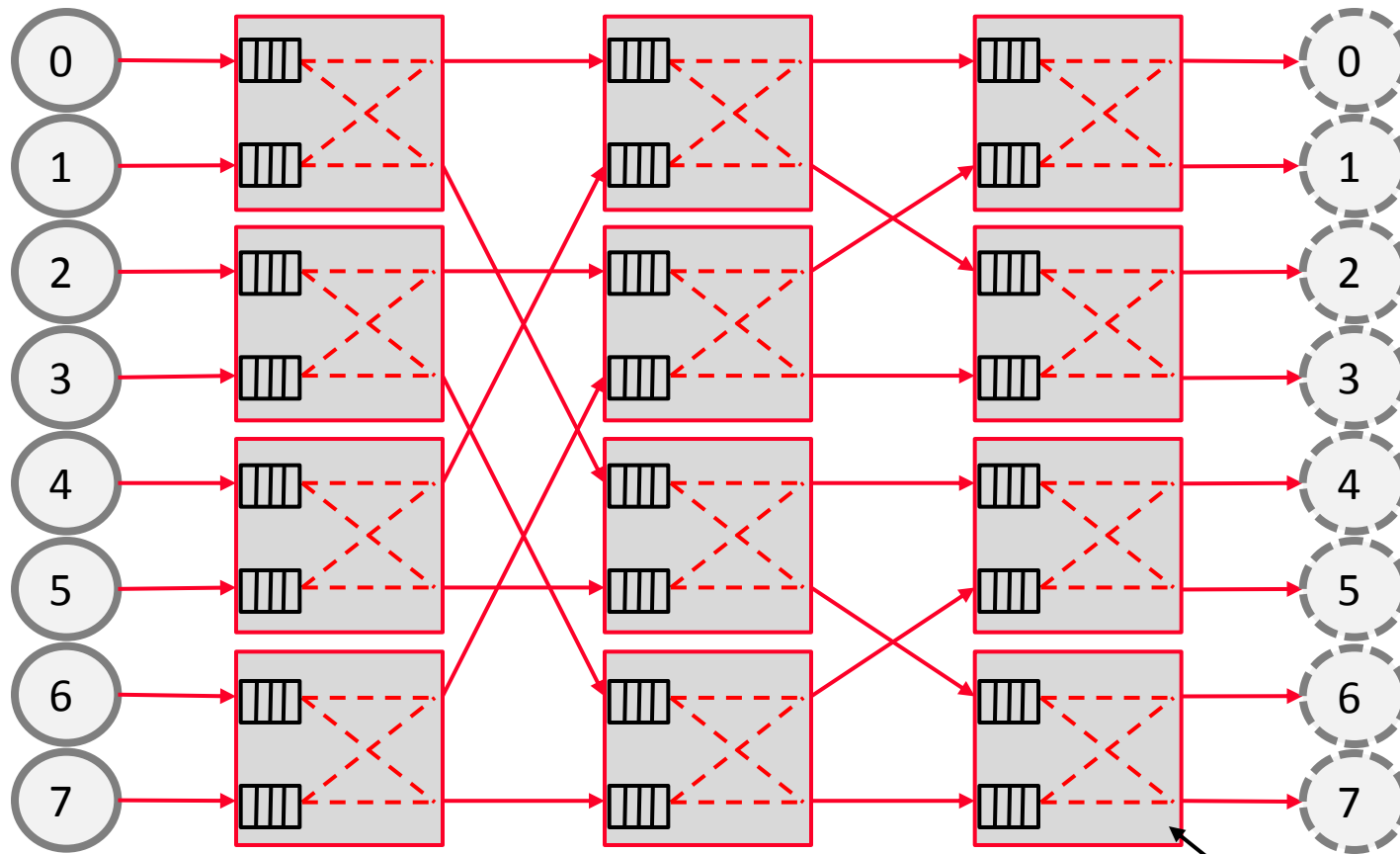


Multistage Circuit Switched



- More restrictions on feasible concurrent Tx-Rx pairs
- But more scalable than crossbar in cost, e.g., $O(N \log N)$ for Butterfly

Multistage Packet Switched



- Packets "hop" from router to router, pending availability of the next-required switch and buffer

Aside: Circuit vs. Packet Switching

- **Circuit switching** sets up full path
 - Establish route then send data
 - (no one else can use those links)
 - + faster arbitration
 - setting up and bringing down links takes time
- **Packet switching** routes per packet
 - Route each packet individually (possibly via different paths)
 - if link is free, any packet can use it
 - potentially slower --- must dynamically switch
 - + no setup, bring down time
 - + more flexible, does not underutilize links

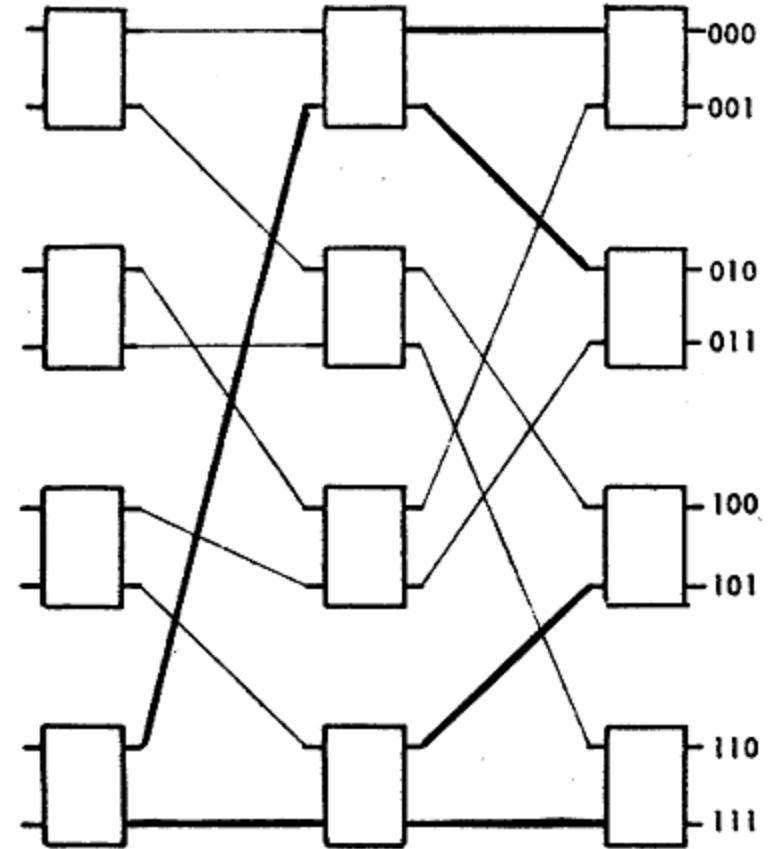
Switching vs. Topology

- Circuit/packet switching choice independent of topology
- It is a higher-level protocol on how a message gets sent to a destination

- However, some topologies are more amenable to circuit vs. packet switching

Another Example: Delta Network

- Single path from source to destination
- Does not support all possible permutations
- Proposed to replace costly crossbars as processor-memory interconnect
- Janak H. Patel , “[Processor-Memory Interconnections for Multiprocessors](#),” ISCA 1979.



8x8 Delta network

Another Example: Omega Network

- Single path from source to destination
- All stages are the same
- Used in NYU Ultracomputer
- Gottlieb et al. “[The NYU Ultracomputer-designing MIMD, shared-memory parallel machine,](#)” ISCA 1982.

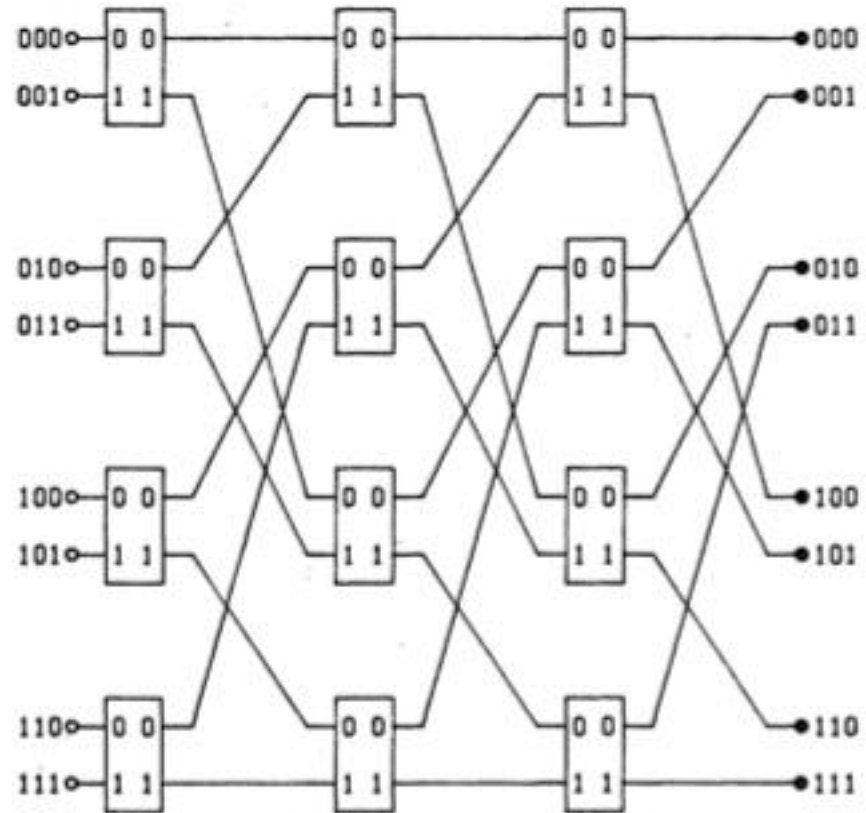
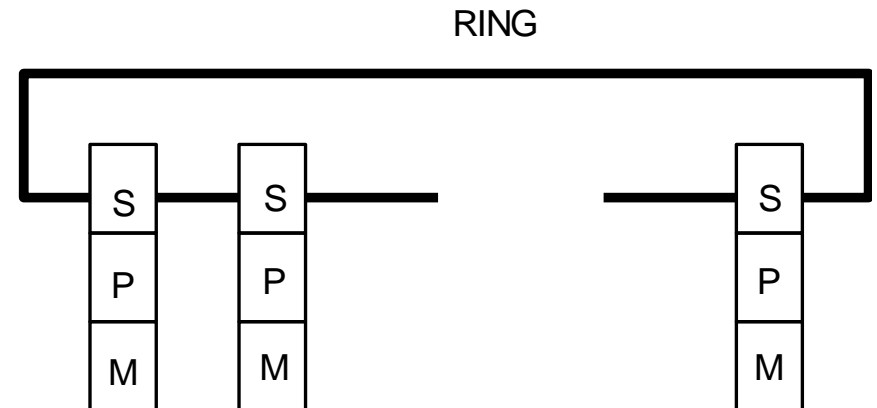


Fig. 2. Omega-network ($N = 8$).

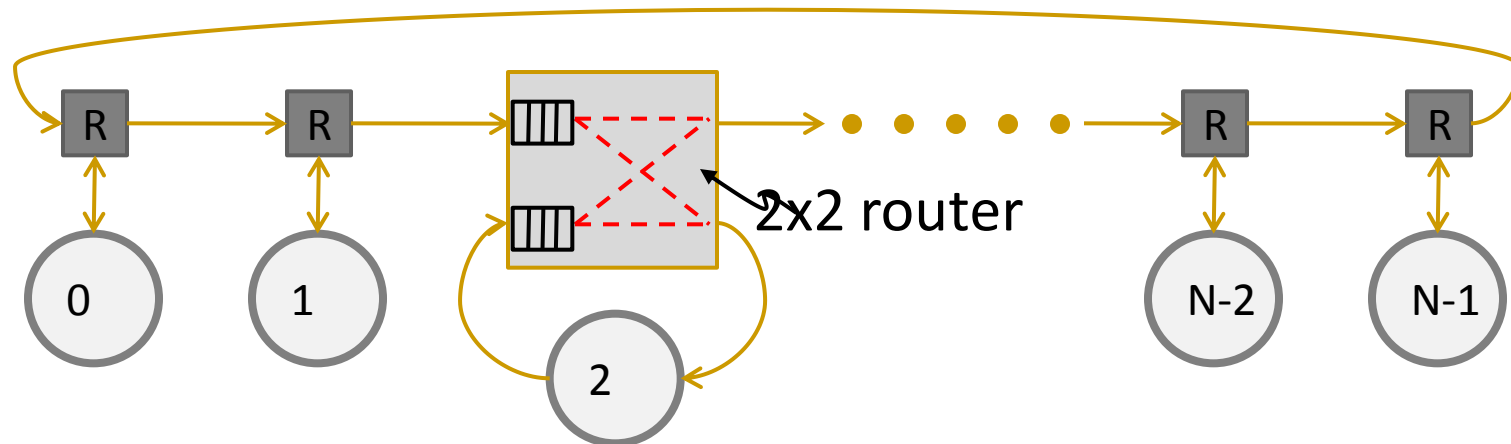
Ring

- + Cheap: $O(N)$ cost
- High latency: $O(N)$
- Not easy to scale
 - Bisection bandwidth remains constant

Used in Intel Haswell, Intel Larrabee, IBM Cell, many commercial systems today



Unidirectional Ring

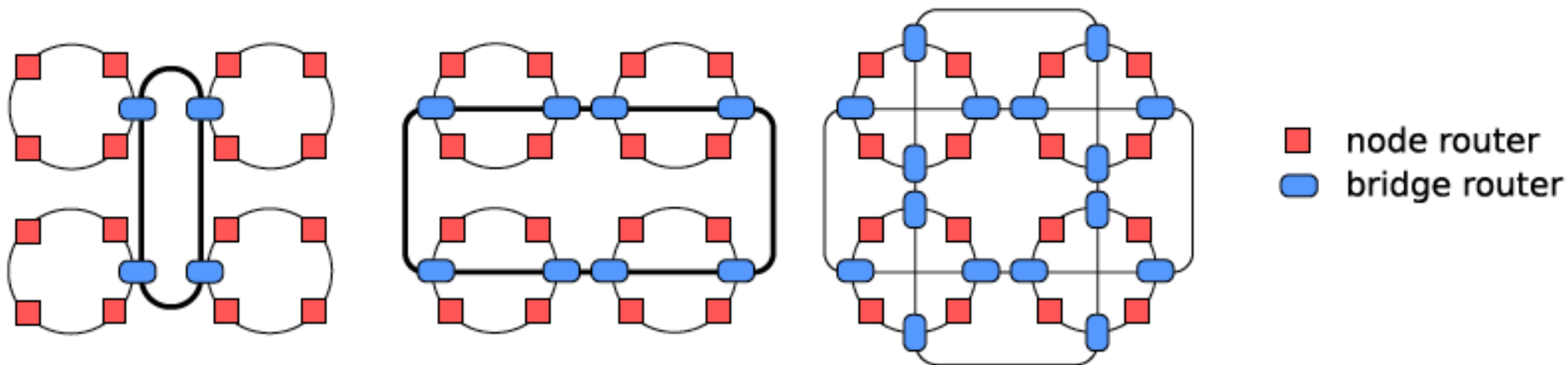


- Simple topology and implementation
 - Reasonable performance if N and performance needs (bandwidth & latency) still moderately low
 - $O(N)$ cost
 - $N/2$ average hops; latency depends on utilization

Bidirectional Rings

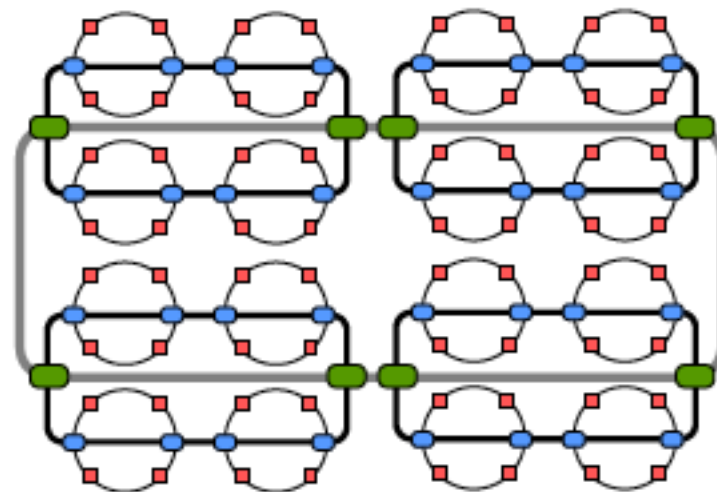
- + Reduces latency
- + Improves scalability
- Slightly more complex injection policy (need to select which ring to inject a packet into)

Hierarchical Rings



(a) 4-, 8-, and 16-bridge hierarchical ring topologies.

- + More scalable
- + Lower latency
- More complex



(b) Three-level hierarchy (8x8).

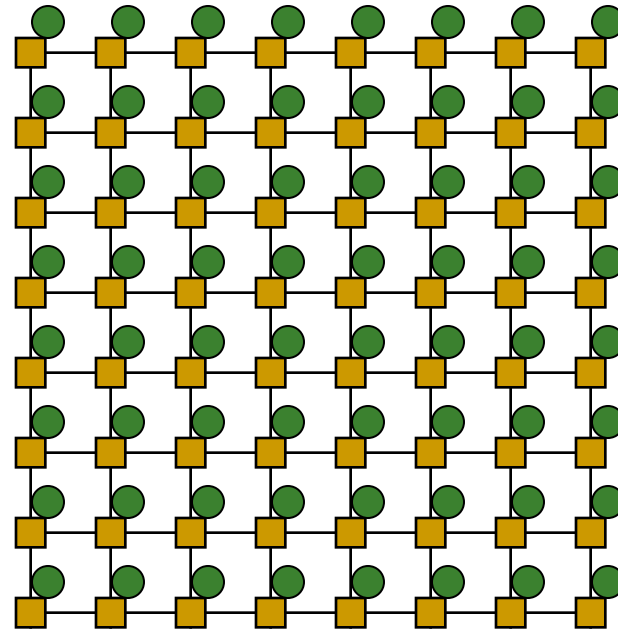
More on Hierarchical Rings

- Chris Fallin, Xiangyao Yu, Kevin Chang, Rachata Ausavarungnirun, Greg Nazario, Reetuparna Das, and Onur Mutlu,
"HiRD: A Low-Complexity, Energy-Efficient Hierarchical Ring Interconnect"
SAFARI Technical Report, TR-SAFARI-2012-004, Carnegie Mellon University, December 2012.
- Discusses the design and implementation of a mostly-bufferless hierarchical ring

Mesh

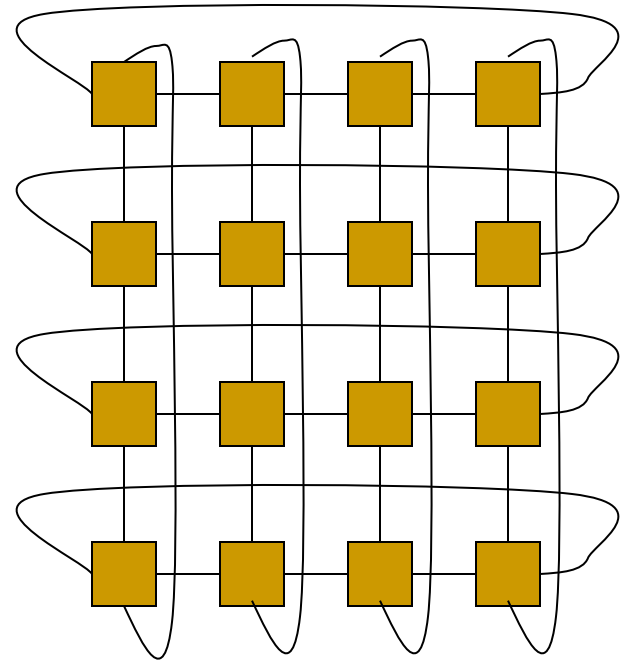
- $O(N)$ cost
- Average latency: $O(\sqrt{N})$
- Easy to layout on-chip: regular and equal-length links
- Path diversity: many ways to get from one node to another

- Used in Tileria 100-core
- And many on-chip network prototypes



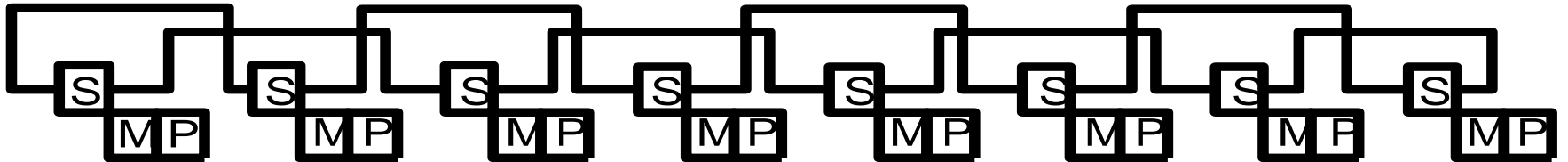
Torus

- Mesh is not symmetric on edges: performance very sensitive to placement of task on edge vs. middle
 - Torus avoids this problem
- + Higher path diversity (and bisection bandwidth) than mesh
- Higher cost
 - Harder to lay out on-chip
 - Unequal link lengths



Torus, continued

- Weave nodes to make inter-node latencies \sim constant



Trees

Planar, hierarchical topology

Latency: $O(\log N)$

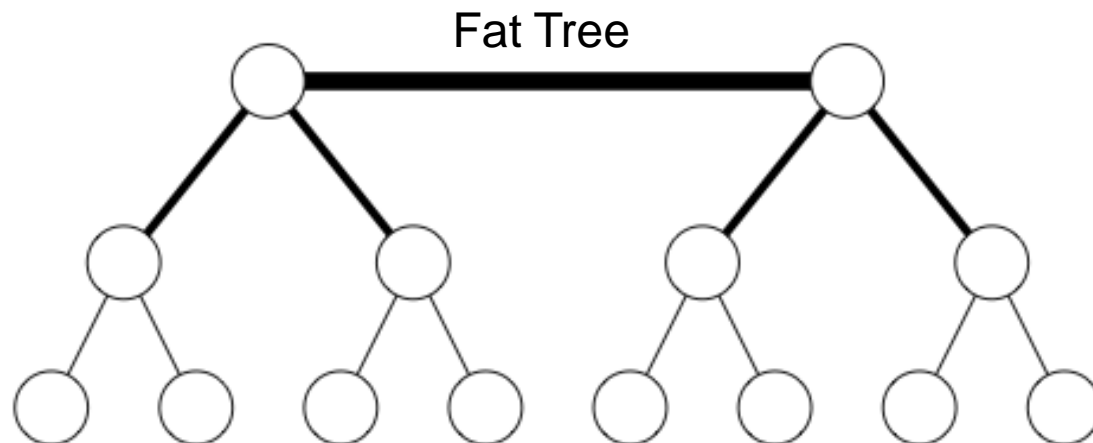
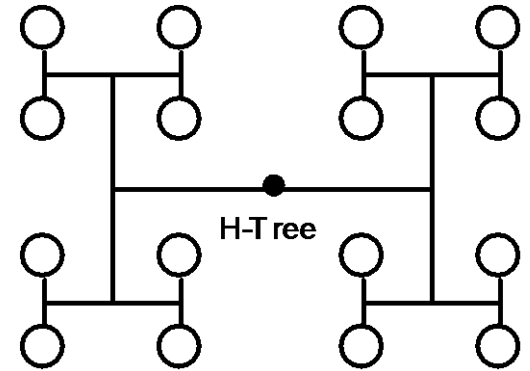
Good for local traffic

+ Cheap: $O(N)$ cost

+ Easy to Layout

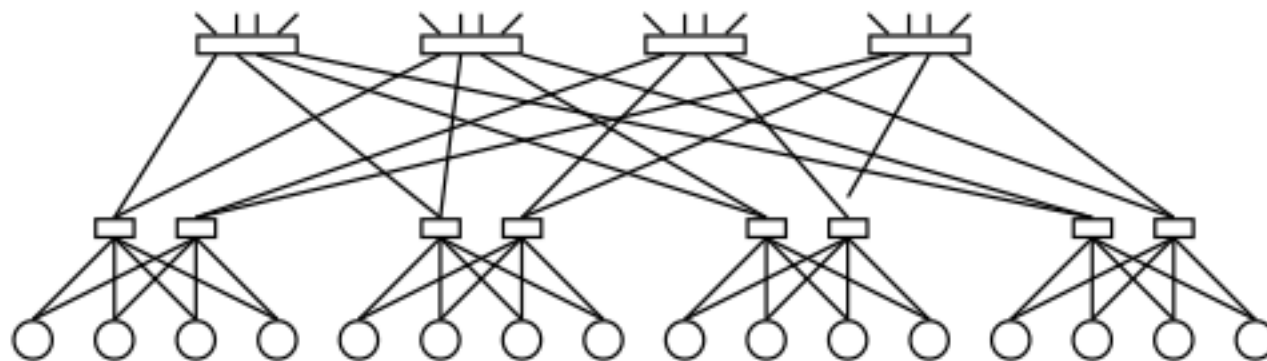
- Root can become a bottleneck

Fat trees avoid this problem (CM-5)



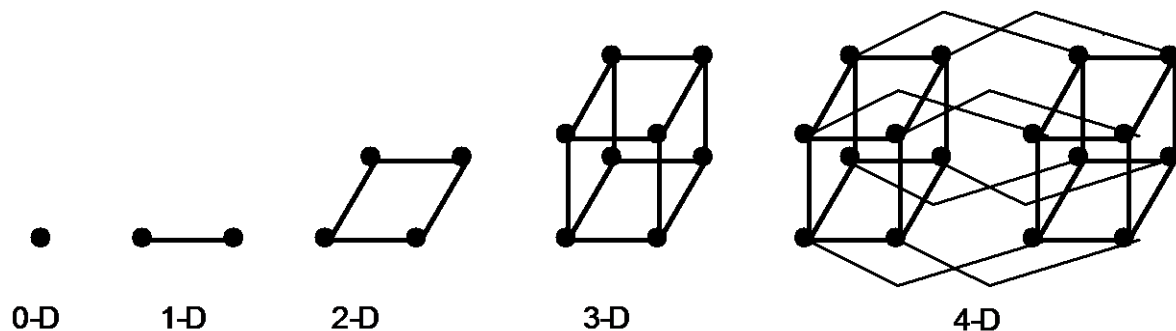
CM-5 Fat Tree

- Fat tree based on 4x2 switches
- Randomized routing on the way up
- Combining, multicast, reduction operators supported in hardware
 - Thinking Machines Corp., “[The Connection Machine CM-5 Technical Summary](#),” Jan. 1992.

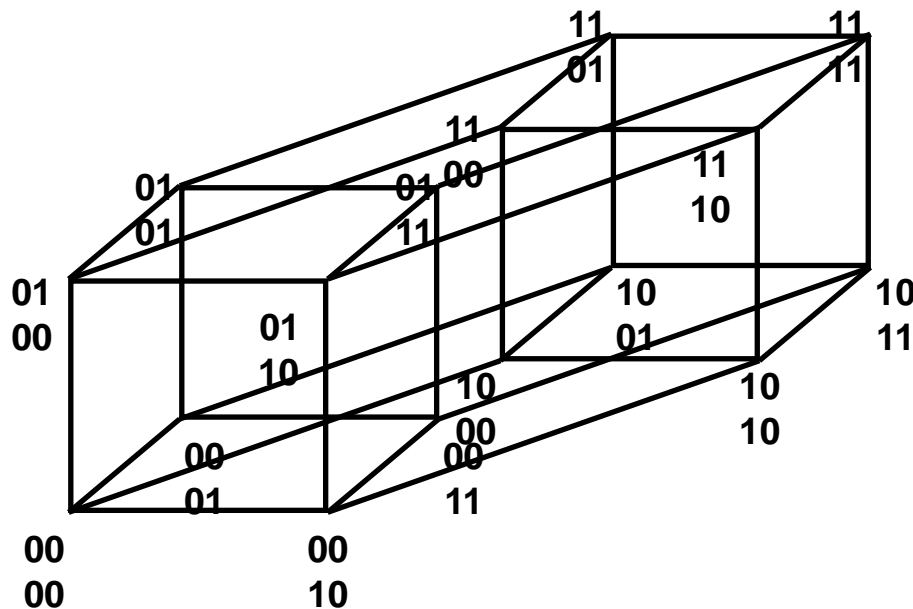


CM-5 Thinned Fat Tree

Hypercube

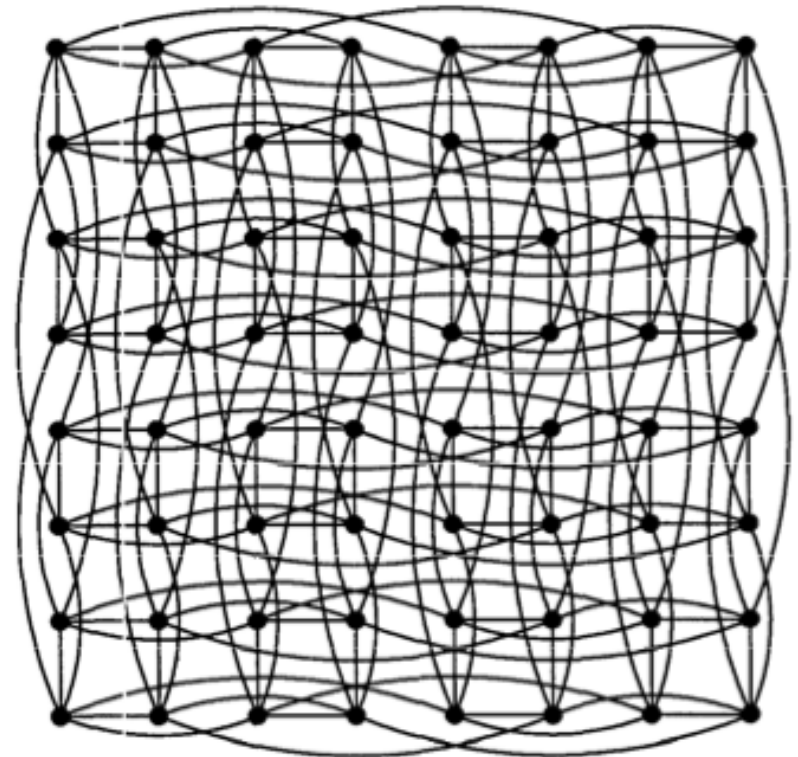
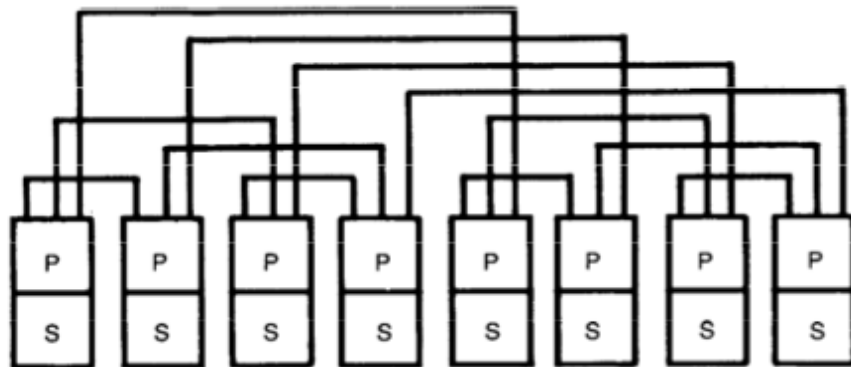


- Latency: $O(\log N)$
 - Radix: $O(\log N)$
 - #links: $O(N \log N)$
- + Low latency
- Hard to lay out in 2D/3D



Caltech Cosmic Cube

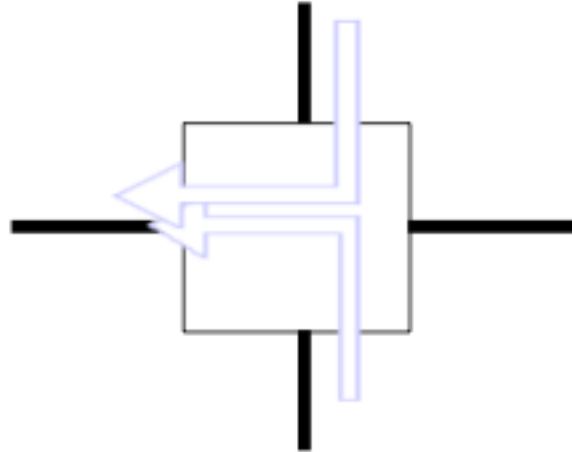
- 64-node message passing machine
- Seitz, “The Cosmic Cube,” CACM 1985.



A hypercube connects $N = 2^n$ small computers, called nodes, through point-to-point communication channels in the Cosmic Cube. Shown here is a two-dimensional projection of a six-dimensional hypercube, or binary 6-cube, which corresponds to a 64-node machine.

FIGURE 1. A Hypercube (also known as a binary cube or a Boolean n -cube)

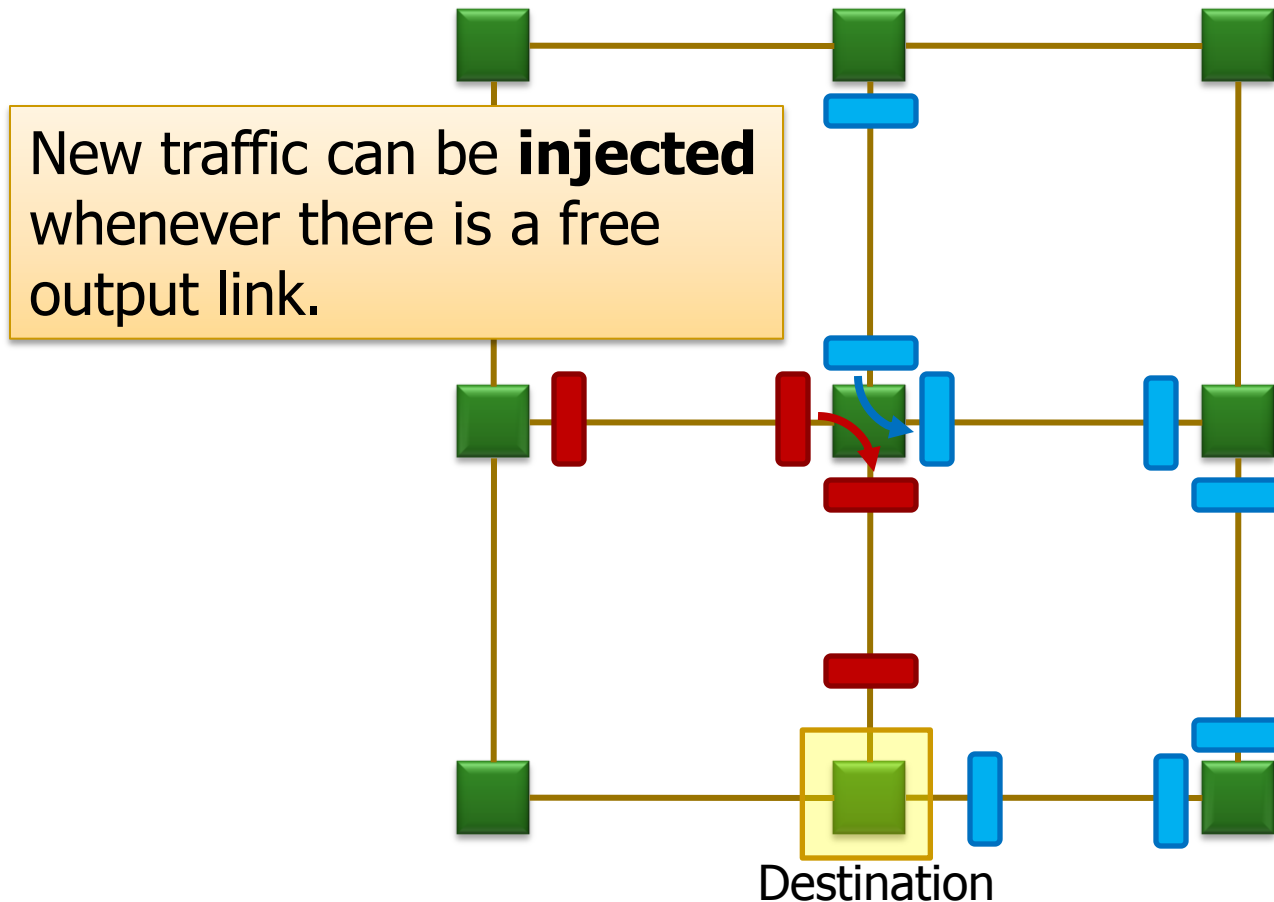
Handling Contention



- Two packets trying to use the same link at the same time
- What do you do?
 - Buffer one
 - Drop one
 - Misroute one (deflection)
- Tradeoffs?

Bufferless Deflection Routing

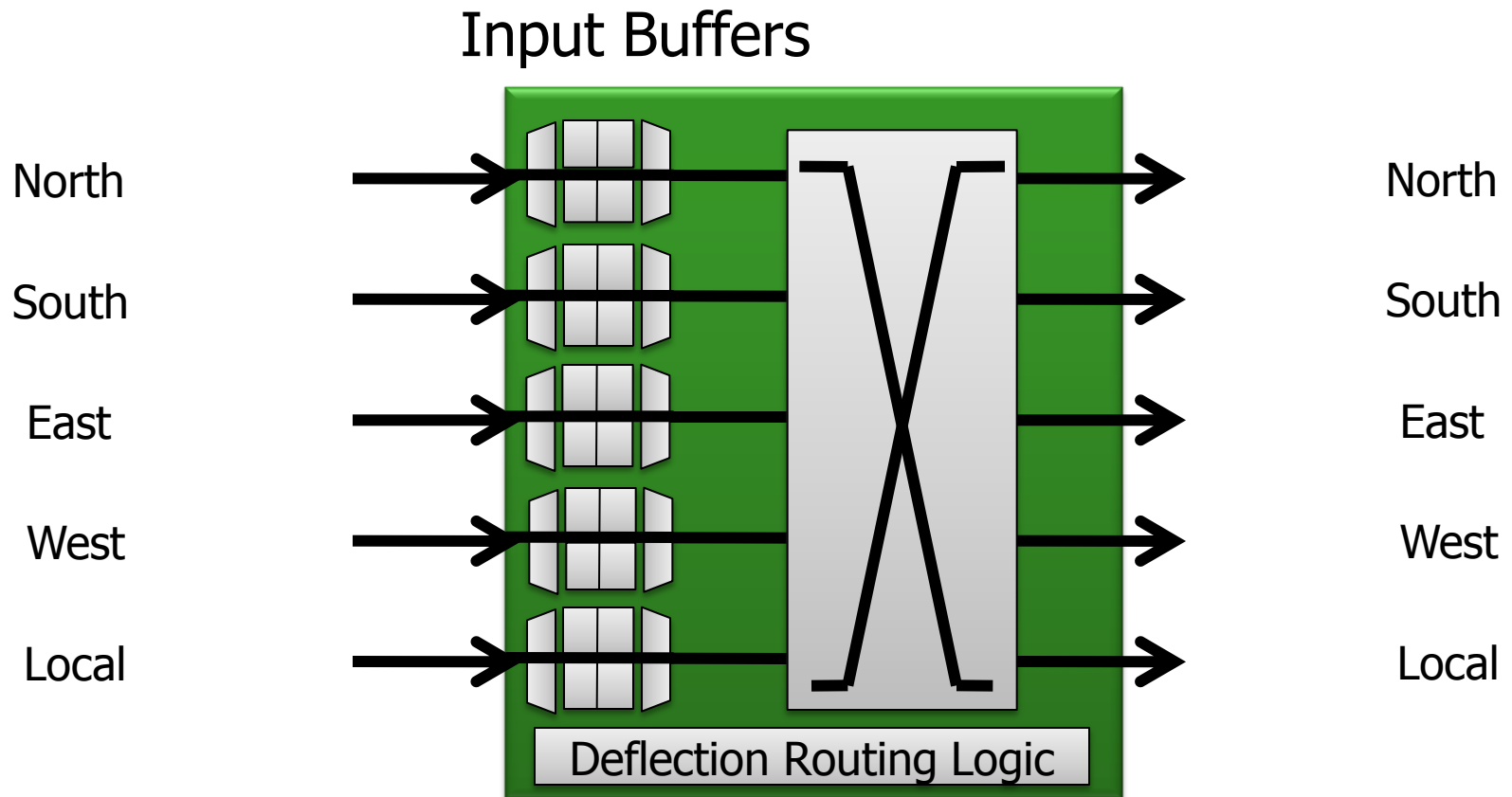
- **Key idea:** Packets are never buffered in the network. When two packets contend for the same link, one is **deflected**.¹



¹Baran, "On Distributed Communication Networks." RAND Tech. Report., 1962 / IEEE Trans.Comm., 1964.52

Bufferless Deflection Routing

- Input buffers are eliminated: flits are buffered in **pipeline latches** and on **network links**



Routing Algorithm

- Types
 - **Deterministic:** always chooses the same path for a communicating source-destination pair
 - **Oblivious:** chooses different paths, without considering network state
 - **Adaptive:** can choose different paths, adapting to the state of the network

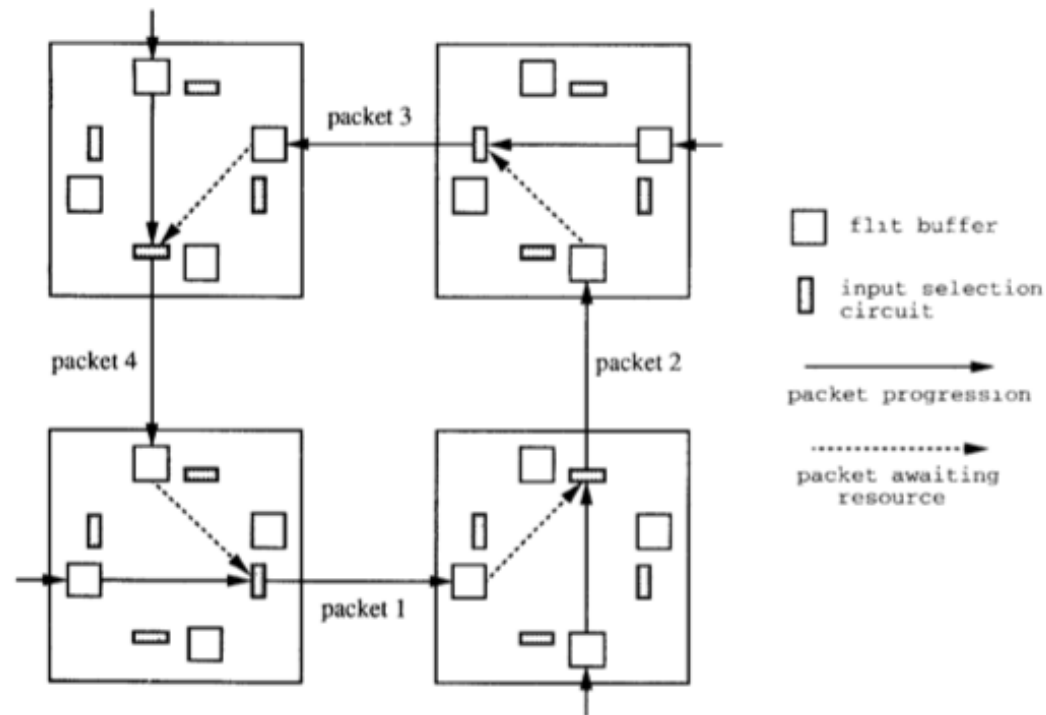
- How to adapt
 - Local/global feedback
 - Minimal or non-minimal paths

Deterministic Routing

- All packets between the same (source, dest) pair take the same path
 - Dimension-order routing
 - E.g., XY routing (used in Cray T3D, and many on-chip networks)
 - First traverse dimension X, then traverse dimension Y
- + Simple
- + Deadlock freedom (no cycles in resource allocation)
- Could lead to high contention
- Does not exploit path diversity

Deadlock

- No forward progress
- Caused by circular dependencies on resources
- Each packet waits for a buffer occupied by another packet downstream



Handling Deadlock

- Avoid cycles in routing
 - Dimension order routing
 - Cannot build a circular dependency
 - Restrict the “turns” each packet can take

- Avoid deadlock by adding more buffering (escape paths)

- Detect and break deadlock
 - Preemption of buffers

Turn Model to Avoid Deadlock

■ Idea

- Analyze directions in which packets can turn in the network
- Determine the cycles that such turns can form
- Prohibit just enough turns to break possible cycles

- Glass and Ni, “[The Turn Model for Adaptive Routing](#),” ISCA 1992.

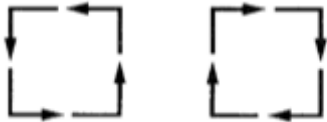


FIG. 2. The possible turns and simple cycles in a two-dimensional mesh.



FIG. 3. The four turns allowed by the xy routing algorithm.

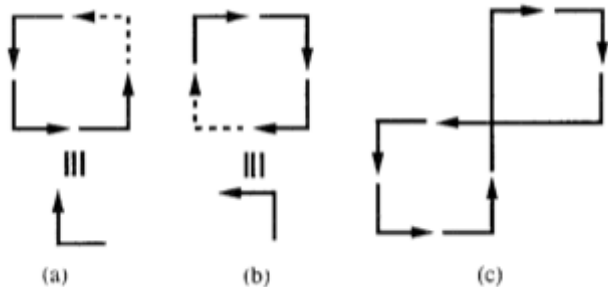


FIG. 4. Six turns that complete the cycles and allow deadlock.

Oblivious Routing: Valiant's Algorithm

- An example of oblivious algorithm
 - Goal: Balance network load
 - Idea: Randomly choose an intermediate destination, route to it first, then route from there to destination
 - Between source-intermediate and intermediate-dest, can use dimension order routing
- + Randomizes/balances network load
- Non minimal (packet latency can increase)
- Optimizations:
 - Do this on high load
 - Restrict the intermediate node to be close (in the same quadrant)

Adaptive Routing

■ Minimal adaptive

- Router uses network state (e.g., downstream buffer occupancy) to pick which “productive” output port to send a packet to
- Productive output port: port that gets the packet closer to its destination

+ Aware of local congestion

- Minimality restricts achievable link utilization (load balance)

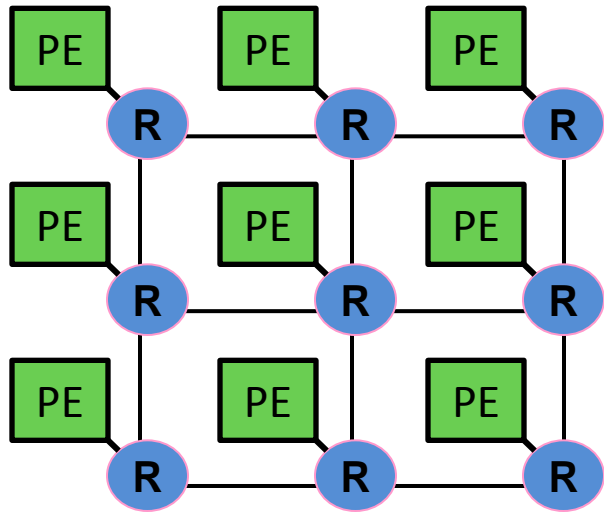
■ Non-minimal (fully) adaptive

- “Misroute” packets to non-productive output ports based on network state

+ Can achieve better network utilization and load balance

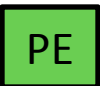
- Need to guarantee livelock freedom

On-Chip Networks

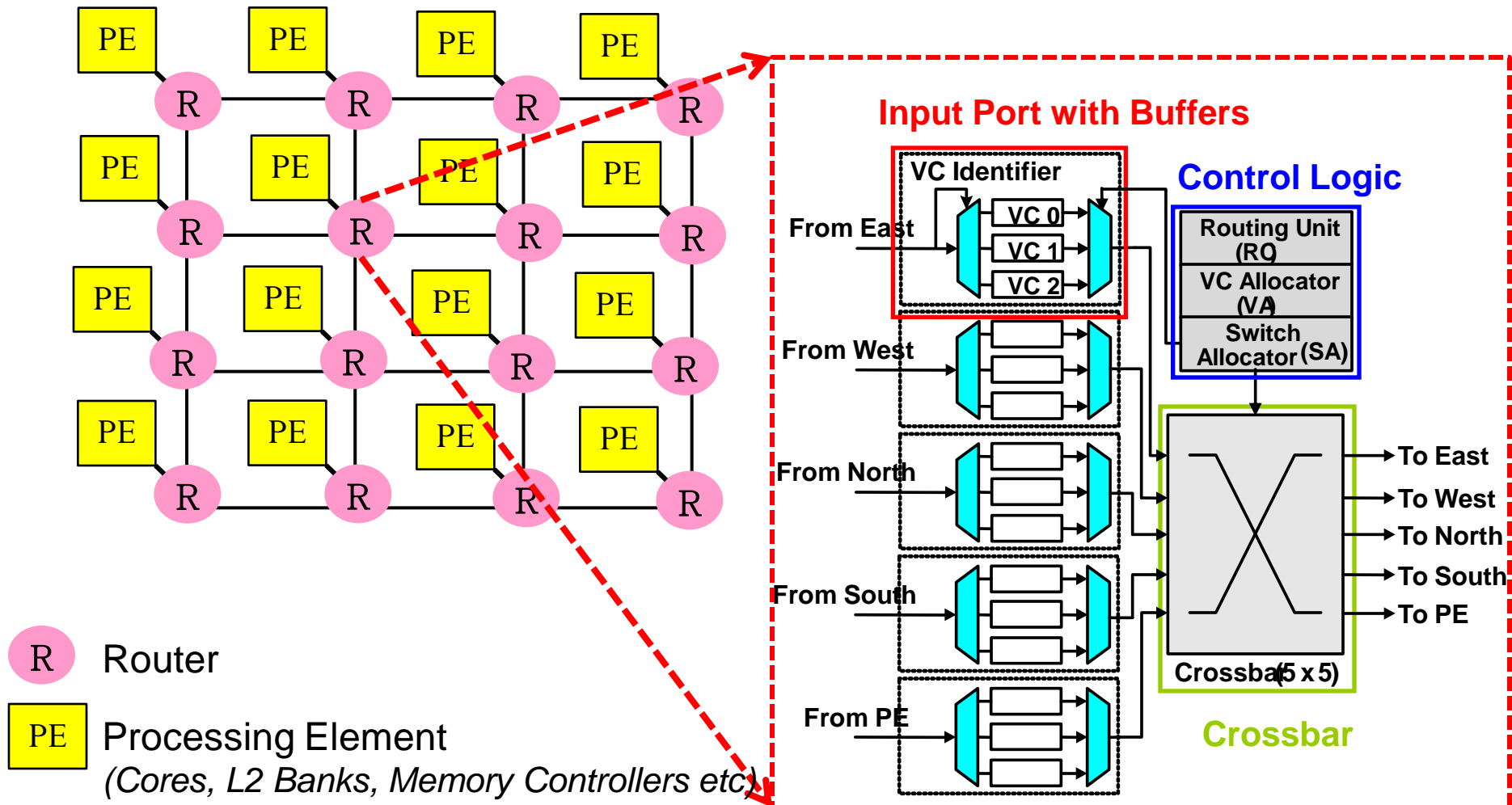


- Connect **cores, caches, memory controllers, etc**
 - Buses and crossbars are not scalable
- **Packet switched**
- **2D mesh:** Most commonly used topology
- Primarily serve **cache misses** and **memory requests**

 Router

 Processing Element
(Cores, L2 Banks, Memory Controllers, etc)

On-chip Networks



On-Chip vs. Off-Chip Interconnects

- On-chip advantages
 - Low latency between cores
 - No pin constraints
 - Rich wiring resources
 - Very high bandwidth
 - Simpler coordination

- On-chip constraints/disadvantages
 - 2D substrate limits implementable topologies
 - Energy/power consumption a key concern
 - Complex algorithms undesirable
 - Logic area constrains use of wiring resources

On-Chip vs. Off-Chip Interconnects (II)

- Cost
 - Off-chip: Channels, pins, connectors, cables
 - On-chip: Cost is storage and switches (wires are plentiful)
 - Leads to networks with many wide channels, few buffers
- Channel characteristics
 - On chip short distance → low latency
 - On chip RC lines → need repeaters every 1-2mm
 - Can put logic in repeaters
- Workloads
 - Multi-core cache traffic vs. supercomputer interconnect traffic

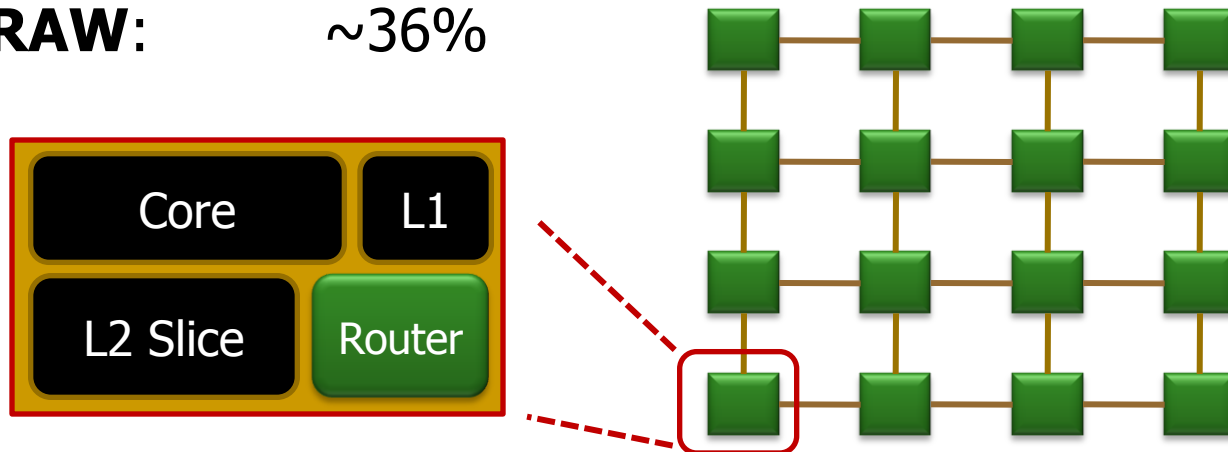
Motivation for Efficient Interconnect

- In many-core chips, on-chip interconnect (NoC) consumes **significant power**

Intel Terascale: ~28% of chip power

Intel SCC: ~10%

MIT RAW: ~36%



- Recent work¹ uses **bufferless deflection routing** to reduce power and die area

¹Moscibroda and Mutlu, "A Case for Bufferless Deflection Routing in On-Chip Networks." ISCA 2009.

Research in Interconnects

Research Topics in Interconnects

- Plenty of topics in interconnection networks. Examples:
- **Energy/power** efficient and proportional design
- **Reducing Complexity**: Simplified router and protocol designs
- **Adaptivity**: Ability to adapt to different access patterns
- **QoS and performance isolation**
 - Reducing and controlling interference, admission control
- **Co-design of NoCs with other shared resources**
 - End-to-end performance, QoS, power/energy optimization
- **Scalable topologies** to many cores, heterogeneous systems
- Fault tolerance
- Request prioritization, priority inversion, coherence, ...
- New technologies (optical, 3D)

One Example: Packet Scheduling

- Which packet to choose for a given output port?
 - Router needs to prioritize between competing flits
 - Which input port?
 - Which virtual channel?
 - Which application's packet?
- Common strategies
 - Round robin across virtual channels
 - Oldest packet first (or an approximation)
 - Prioritize some virtual channels over others
- Better policies in a multi-core environment
 - Use application characteristics
 - Minimize energy

Computer Architecture: Interconnects (Part I)

Prof. Onur Mutlu
Carnegie Mellon University

Bufferless Routing

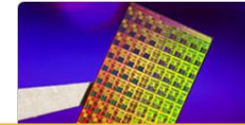
Thomas Moscibroda and Onur Mutlu,
"A Case for Bufferless Routing in On-Chip Networks"
*Proceedings of the 36th International Symposium on Computer
Architecture (ISCA)*, pages 196-207, Austin, TX, June 2009.
Slides (pptx)

On-Chip Networks (NoC)

- Connect **cores, caches, memory controllers**, etc...

- Examples:

- Intel 80-core Terascale chip
- MIT RAW



- Design goals in NoC

- High throughput
- Fairness between nodes
- Low complexity
- Power, low

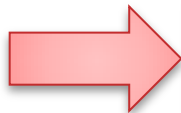


Energy/Power in On-Chip Networks

- Power is a **key constraint** in the design of high-performance processors
- NoCs consume substantial portion of system power
 - **~30%** in Intel 80-core Terascale [IEEE Micro' 07]
 - **~40%** in MIT RAW Chip [ISCA' 04]
- NoCs estimated to consume **100s of Watts** [Borkar, DAC' 07]

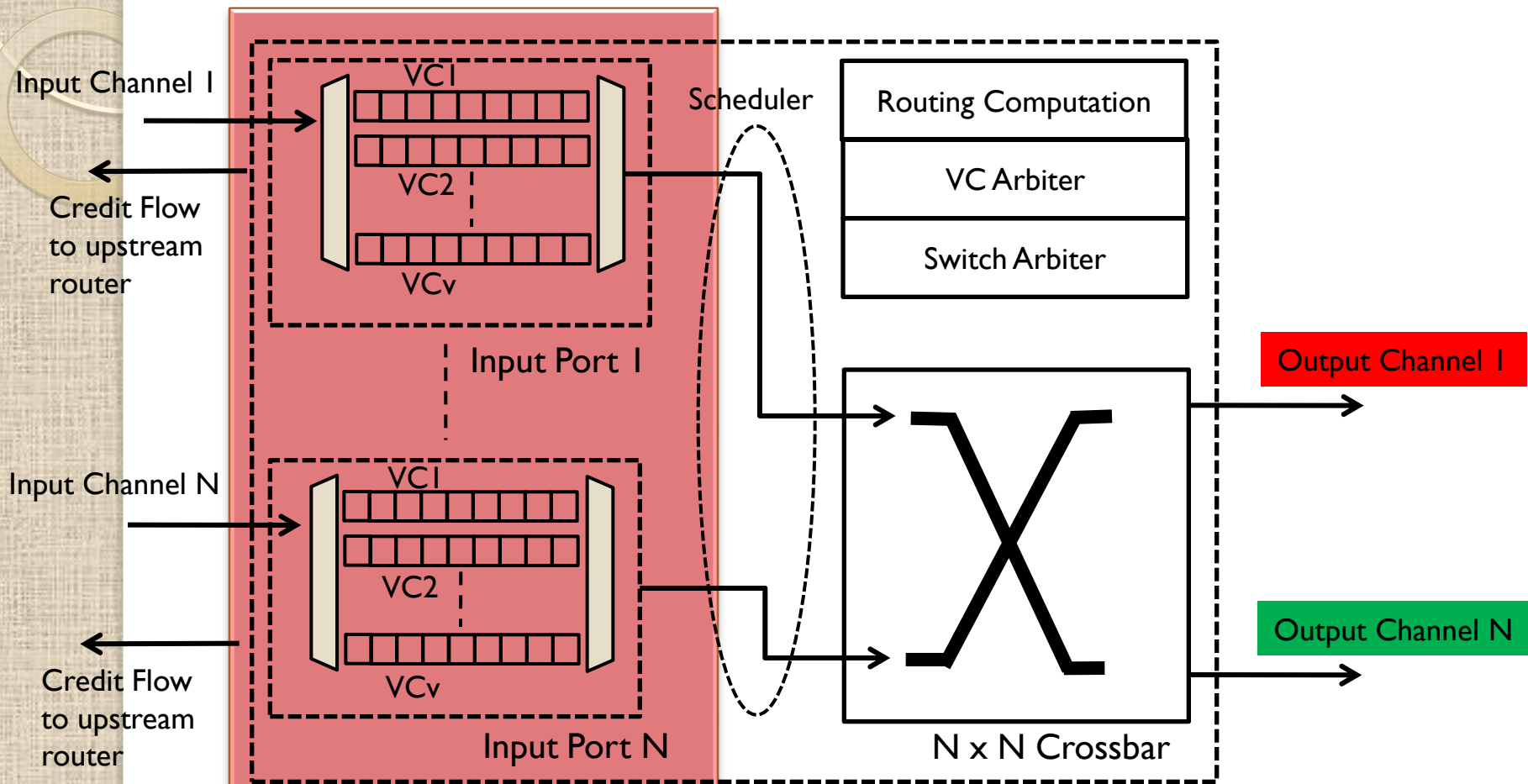
Current NoC Approaches

- Existing approaches differ in numerous ways:
 - Network topology [Kim et al, ISCA' 07, Kim et al, ISCA' 08 etc]
 - Flow control [Michelogiannakis et al, HPCA' 09, Kumar et al, MICRO' 08, etc]
 - Virtual Channels [Nicolopoulos et al, MICRO' 06, etc]
 - QoS & fairness mechanisms [Lee et al, ISCA' 08, etc]
 - Routing algorithms [Singh et al, CAL' 04]
 - Router architecture [Park et al, ISCA' 08]
 - Broadcast, Multicast [Jerger et al, ISCA' 08, Rodrigo et al, MICRO' 08]



Existing work assumes existence of buffers in routers!

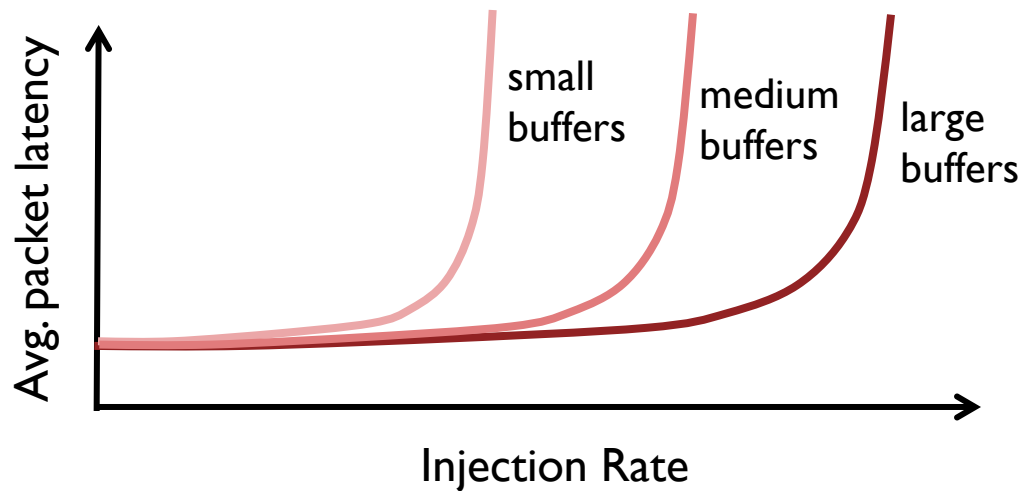
A Typical Router




Buffers are integral part of existing NoC Routers




Buffers in NoC Routers

- Buffers are necessary for high network throughput → buffers increase total available bandwidth in network



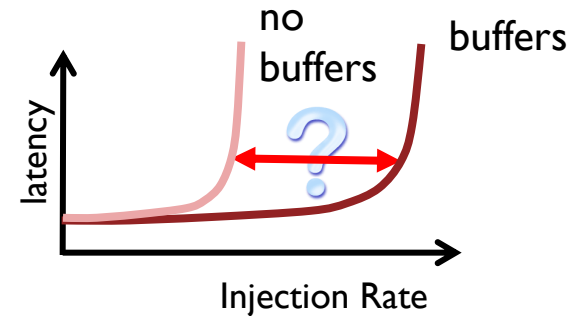
Buffers in NoC Routers

- Buffers are necessary for high network throughput 
→ buffers increase total available bandwidth in network
-

- Buffers consume significant **energy/power**
 - **Dynamic energy** when read/write 
 - **Static energy** even when not occupied
- Buffers add **complexity** and **latency** 
 - Logic for buffer management
 - Virtual channel allocation
 - Credit-based flow control
- Buffers require significant **chip area** 
 - E.g., in TRIPS prototype chip, input buffers occupy 75% of total on-chip network area [Gratz et al, ICCD' 06]

Going Bufferless...?

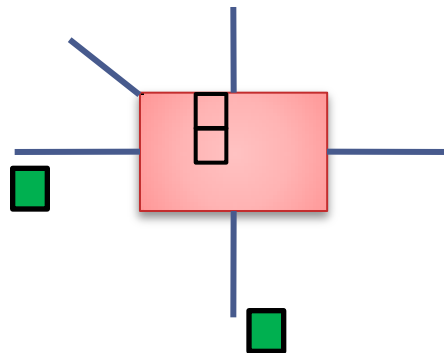
- **How much throughput do we lose?**
→ How is latency affected?



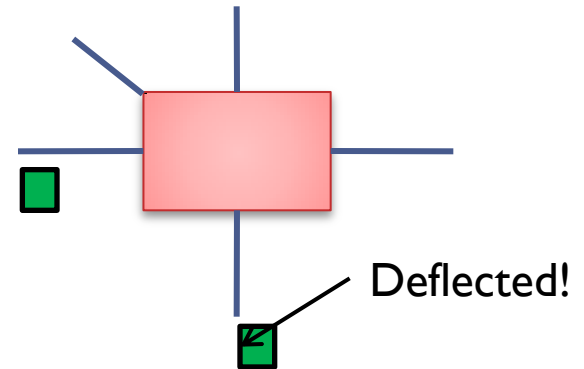
- Up to what **injection rates** can we use bufferless routing?
→ Are there **realistic scenarios** in which NoC is operated at injection rates below the threshold?
- Can we achieve **energy reduction**?
→ If so, how much...?
- Can we reduce **area, complexity**, etc...?

BLESS: Bufferless Routing

- Always forward *all* incoming flits to some output port
- If no productive direction is available, send to another direction
- → packet is deflected
- → **Hot-potato routing** [Baran' 64, etc]

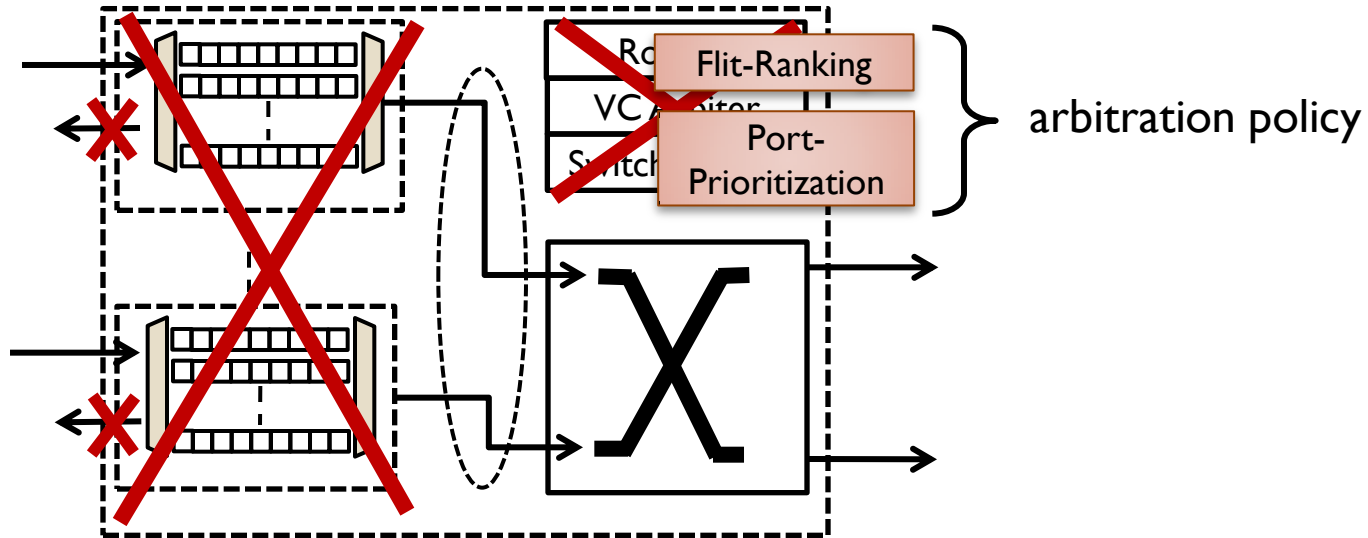


Buffered



BLESS

BLESS: Bufferless Routing



Flit-Ranking

1. Create a ranking over all incoming flits

Port-Prioritization

2. For a given flit in this ranking, find the best free output-port
Apply to each flit in order of ranking

FLIT-BLESS: Flit-Level Routing

- Each flit is routed independently.
- **Oldest-first arbitration** (other policies evaluated in paper)

Flit-Ranking	1. Oldest-first ranking
Port-Prioritization	2. Assign flit to productive port, if possible. Otherwise, assign to non-productive port.

- **Network Topology:**
 - Can be applied to most topologies (Mesh, Torus, Hypercube, Trees, ...)
 - 1) #output ports , #input ports at every router
 - 2) every router is reachable from every other router
- **Flow Control & Injection Policy:**
 - Completely **local**, inject whenever input port is free
- **Absence of Deadlocks:** every flit is always moving
- **Absence of Livelocks:** with oldest-first ranking

BLESS with Buffers

- BLESS without buffers is extreme end of a continuum
- BLESS can be integrated with buffers
 - FLIT-BLESS with Buffers
 - WORM-BLESS with Buffers
- Whenever a buffer is full, it's first flit becomes **must-schedule**
- **must-schedule** flits must be deflected if necessary

See paper for details...

BLESS: Advantages & Disadvantages

Advantages

- No buffers
- Purely local flow control
- Simplicity
 - no credit-flows
 - no virtual channels
 - simplified router design
- No deadlocks, livelocks
- Adaptivity
 - packets are deflected around congested areas!
- Router latency reduction
- Area savings

Disadvantages

- Increased latency
- Reduced bandwidth
- Increased buffering at receiver
- Header information at each flit
- Oldest-first arbitration complex
- QoS becomes difficult

Impact on energy...?



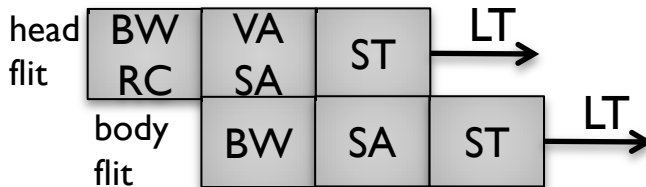
Reduction of Router Latency

- BLESS gets rid of input buffers and virtual channels

BW: Buffer Write
 RC: Route Computation
 VA: Virtual Channel Allocation
 SA: Switch Allocation
 ST: Switch Traversal
 LT: Link Traversal
 LA LT: Link Traversal of Lookahead

[Dally, Towles '04]

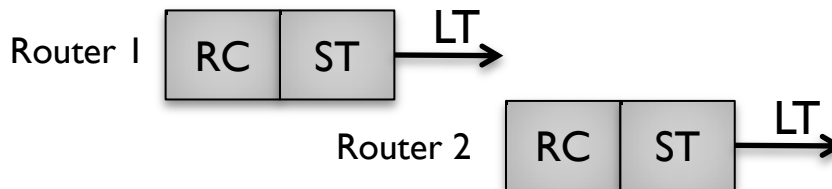
Baseline Router (speculative)



Router Latency = 3

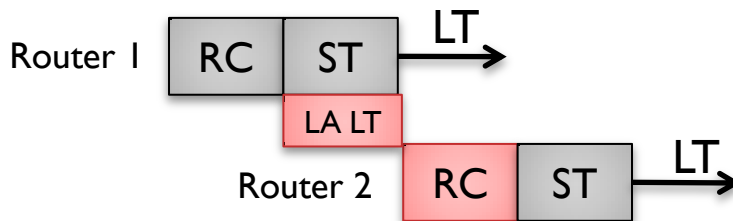
Can be improved to 2.

BLESS Router (standard)



Router Latency = 2

BLESS Router (optimized)



Router Latency = 1

our evaluations

BLESS: Advantages & Disadvantages

Advantages

- No buffers
- Purely local flow control
- Simplicity
 - no credit-flows
 - no virtual channels
 - simplified router design
- No deadlocks, livelocks
- Adaptivity
 - packets are deflected around congested areas!
- Router latency reduction
- Area savings

Disadvantages

- Increased latency
- Reduced bandwidth
- Increased buffering at receiver
- Header information at each flit

Extensive evaluations in the paper!

Impact on energy...?



Evaluation Methodology

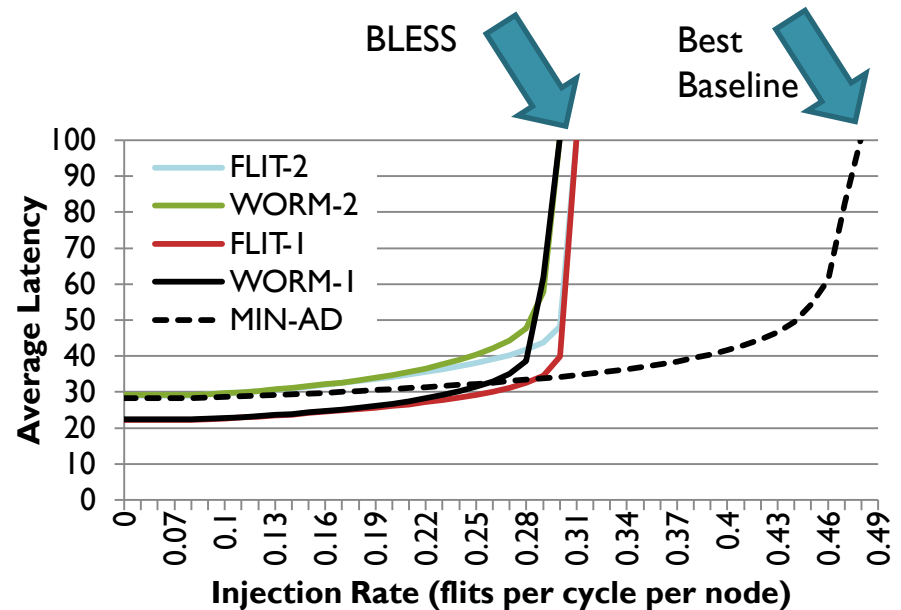
- 2D mesh network, router latency is 2 cycles
 - 4x4, 8 core, 8 L2 cache banks (each node is a core or an L2 bank)
 - 4x4, 16 core, 16 L2 cache banks (each node is a core and an L2 bank)
 - 8x8, 16 core, 64 L2 cache banks (each node is L2 bank and may be a core)
 - 128-bit wide links, 4-flit data packets, 1-flit address packets
 - For baseline configuration: 4 VCs per physical input port, 1 packet deep
- Benchmarks
 - Multiprogrammed SPEC CPU2006 and Windows Desktop applications
 - Heterogeneous and homogenous application mixes
 - Synthetic traffic patterns: UR, Transpose, Tornado, Bit Complement
- x86 processor model based on Intel Pentium M
 - 2 GHz processor, 128-entry instruction window
 - 64Kbyte private L1 caches
 - Total 16Mbyte shared L2 caches; 16 MSHRs per bank
 - DRAM model based on Micron DDR2-800

Evaluation Methodology


- **Energy model** provided by Orion simulator [MICRO' 02]
 - 70nm technology, 2 GHz routers at $1.0V_{dd}$
- For BLESS, we model
 - Additional energy to transmit header information
 - Additional buffers needed on the receiver side
 - Additional logic to reorder flits of individual packets at receiver
- We partition network energy into **buffer energy**, **router energy**, and **link energy**, each having **static** and **dynamic** components.
- Comparisons against non-adaptive and aggressive adaptive buffered routing algorithms (**DO**, **MIN-AD**, **ROMM**)

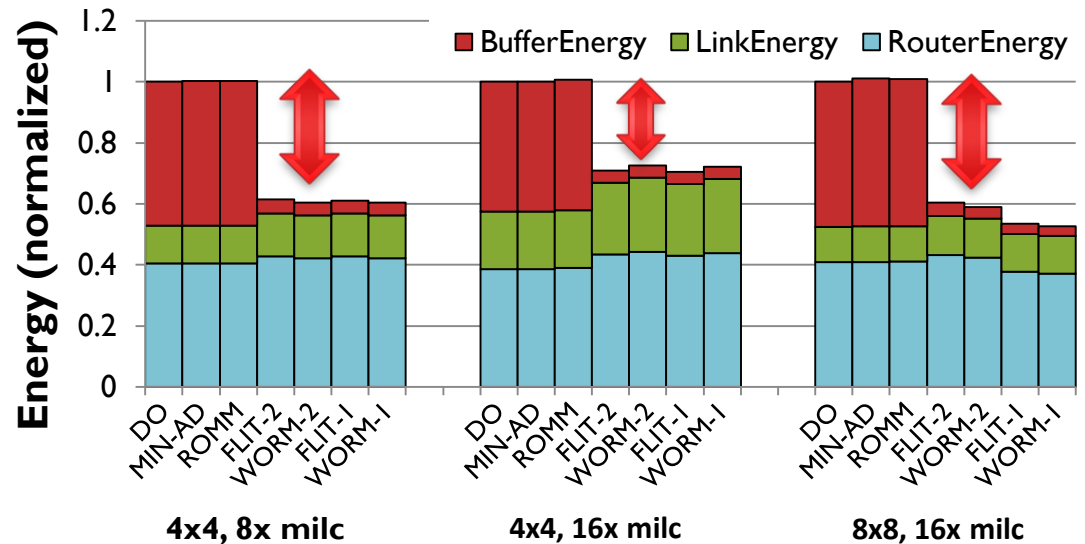
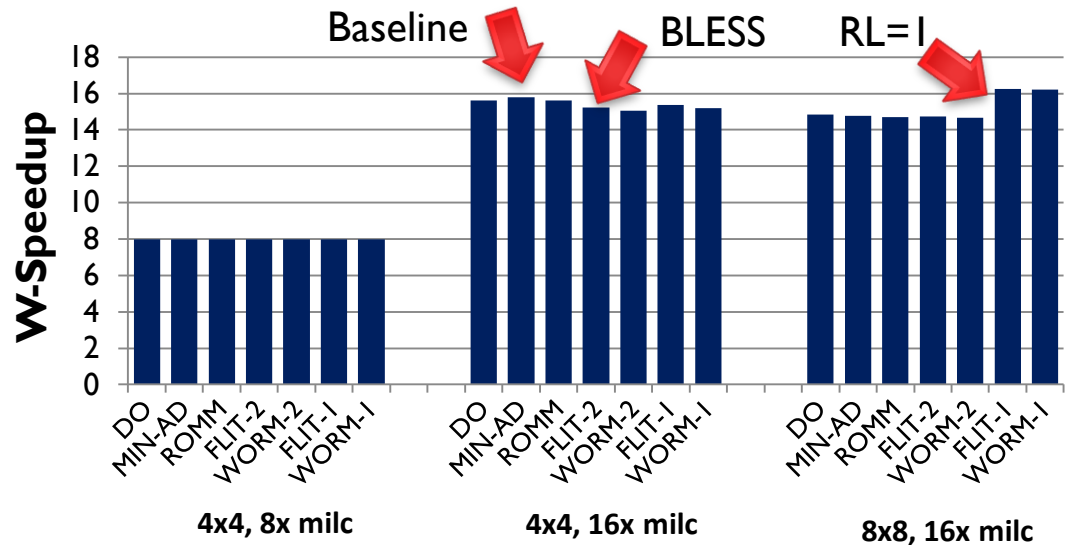
Evaluation – Synthetic Traces

- First, the bad news 😊
- Uniform random injection
- BLESS has significantly lower saturation throughput compared to buffered baseline.



Evaluation – Homogenous Case Study

- **milc** benchmarks (moderately intensive)
- **Perfect caches!**
- Very little performance degradation with BLESS (less than 4% in dense network)
- With router latency 1, BLESS can even outperform baseline (by ~10%)
- **Significant energy improvements (almost 40%)** 



Evaluation – Homogenous Case Study

- milc benchmarks

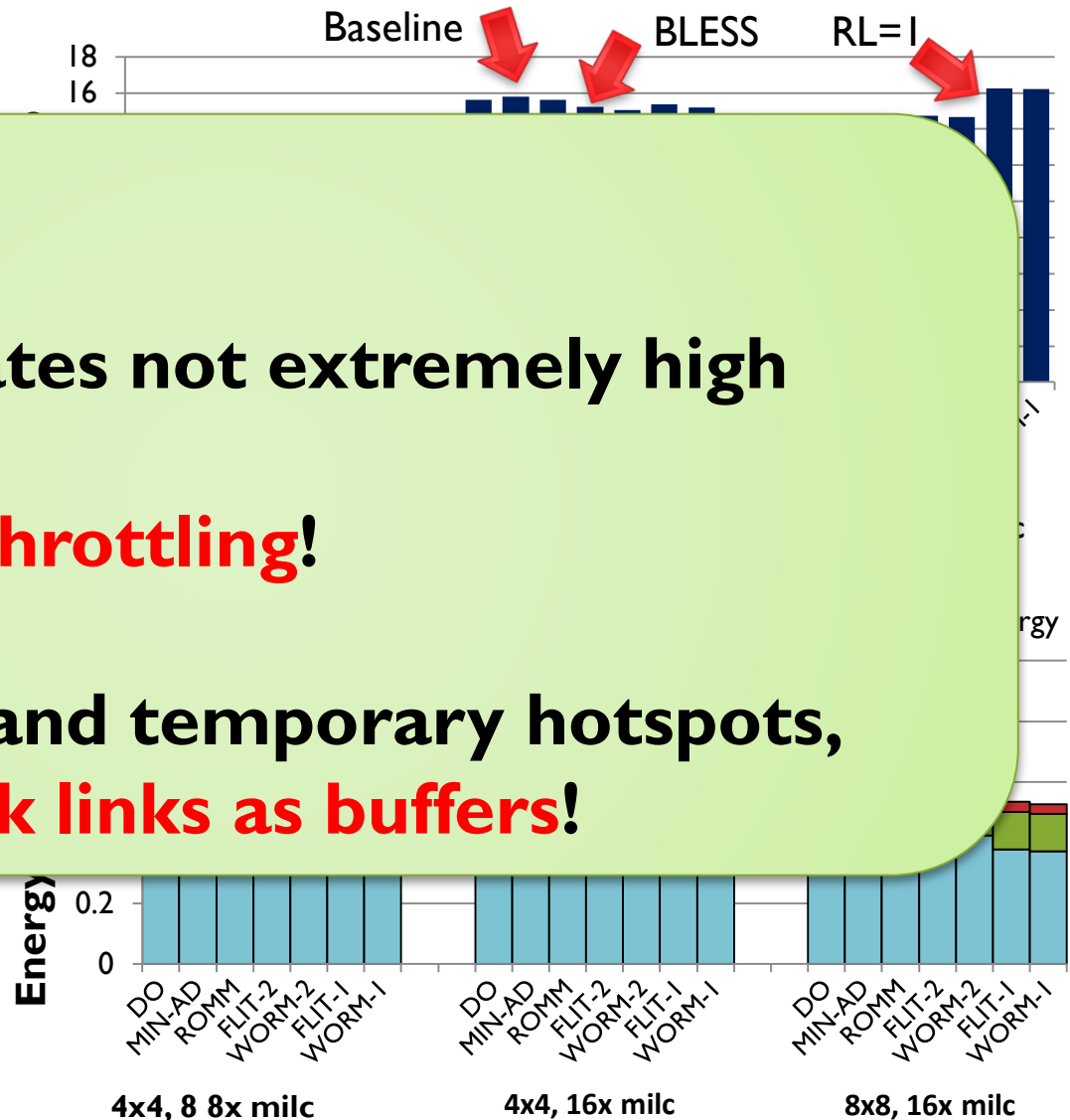
Observations:

1) Injection rates not extremely high on average

→ **self-throttling!**

2) For bursts and temporary hotspots, **use network links as buffers!**

- Significant energy improvements (almost 40%)



Evaluation – Further Results

See paper for details...

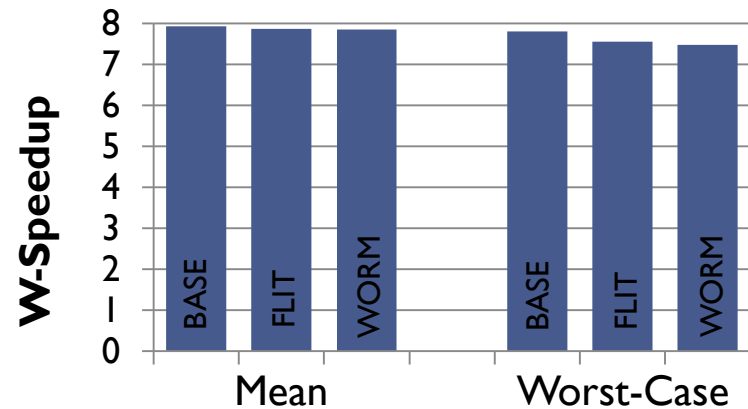
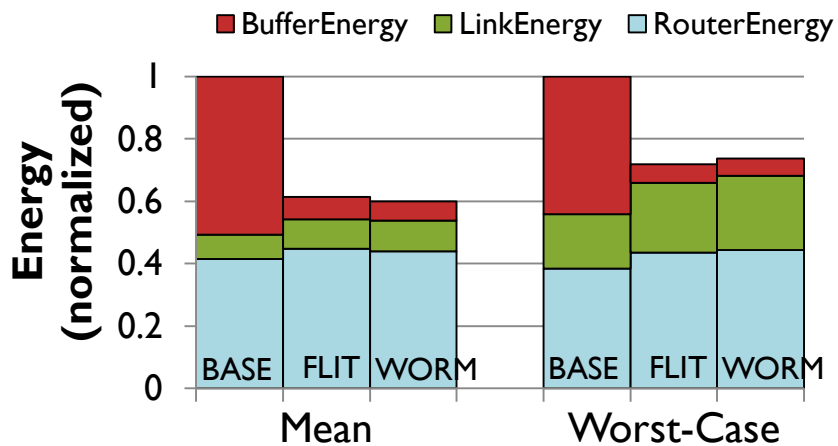
- BLESS increases **buffer requirement** at receiver by at most 2x
→ overall, energy is still reduced
- Impact of **memory latency**
→ with **real caches**, very **little slowdown!** (at most 1.5%)
- **Heterogeneous application mixes**
(we evaluate several mixes of intensive and non-intensive applications)
→ little performance degradation
→ significant energy savings in all cases
→ no significant increase in **unfairness** across different applications
- **Area savings:** ~60% of network area can be saved!



Evaluation – Aggregate Results

- Aggregate results over all 29 applications

Sparse Network	Perfect L2		Realistic L2	
	Average	Worst-Case	Average	Worst-Case
Δ Network Energy	-39.4%	-28.1%	-46.4%	-41.0%
Δ System Performance	-0.5%	-3.2%	-0.15%	-0.55%



Evaluation – Aggregate Results

- Aggregate results over all 29 applications

Sparse Network	Perfect L2		Realistic L2	
	Average	Worst-Case	Average	Worst-Case
Δ Network Energy	-39.4%	-28.1%	-46.4%	-41.0%
Δ System Performance	-0.5%	-3.2%	-0.15%	-0.55%

Dense Network	Perfect L2		Realistic L2	
	Average	Worst-Case	Average	Worst-Case
Δ Network Energy	-32.8%	-14.0%	-42.5%	-33.7%
Δ System Performance	-3.6%	-17.1%	-0.7%	-1.5%

BLESS Conclusions

- For a very wide range of applications and network settings, buffers are not needed in NoC
 - Significant energy savings (32% even in dense networks and perfect caches)
 - Area-savings of 60%
 - Simplified router and network design (flow control, etc...)
 - Performance slowdown is minimal (can even increase!)

➤ A strong case for a **rethinking of NoC design!**

- **Future research:**
 - Support for quality of service, different traffic classes, energy-management, etc...

CHIPPER: A Low-complexity Bufferless Deflection Router

Chris Fallin, Chris Craik, and Onur Mutlu,

"CHIPPER: A Low-Complexity Bufferless Deflection Router"

Proceedings of the 17th International Symposium on High-Performance Computer Architecture (HPCA), pages 144-155, San Antonio, TX, February 2011. Slides (pptx)

SAFARI Carnegie Mellon

Motivation

- Recent work has proposed **bufferless deflection routing** (BLESS [Moscibroda, ISCA 2009])
 - **Energy savings:** ~40% in total NoC energy
 - **Area reduction:** ~40% in total NoC area
 - **Minimal performance loss:** ~4% on average
 - **Unfortunately: unaddressed complexities in router**
 - ➔ long critical path, large reassembly buffers
- **Goal:** obtain these benefits while simplifying the router in order to **make bufferless NoCs practical.**

Problems that Bufferless Routers Must Solve

1. Must provide **livelock freedom**

→ A packet should not be deflected forever

2. Must **reassemble packets** upon arrival

Flit: atomic routing unit

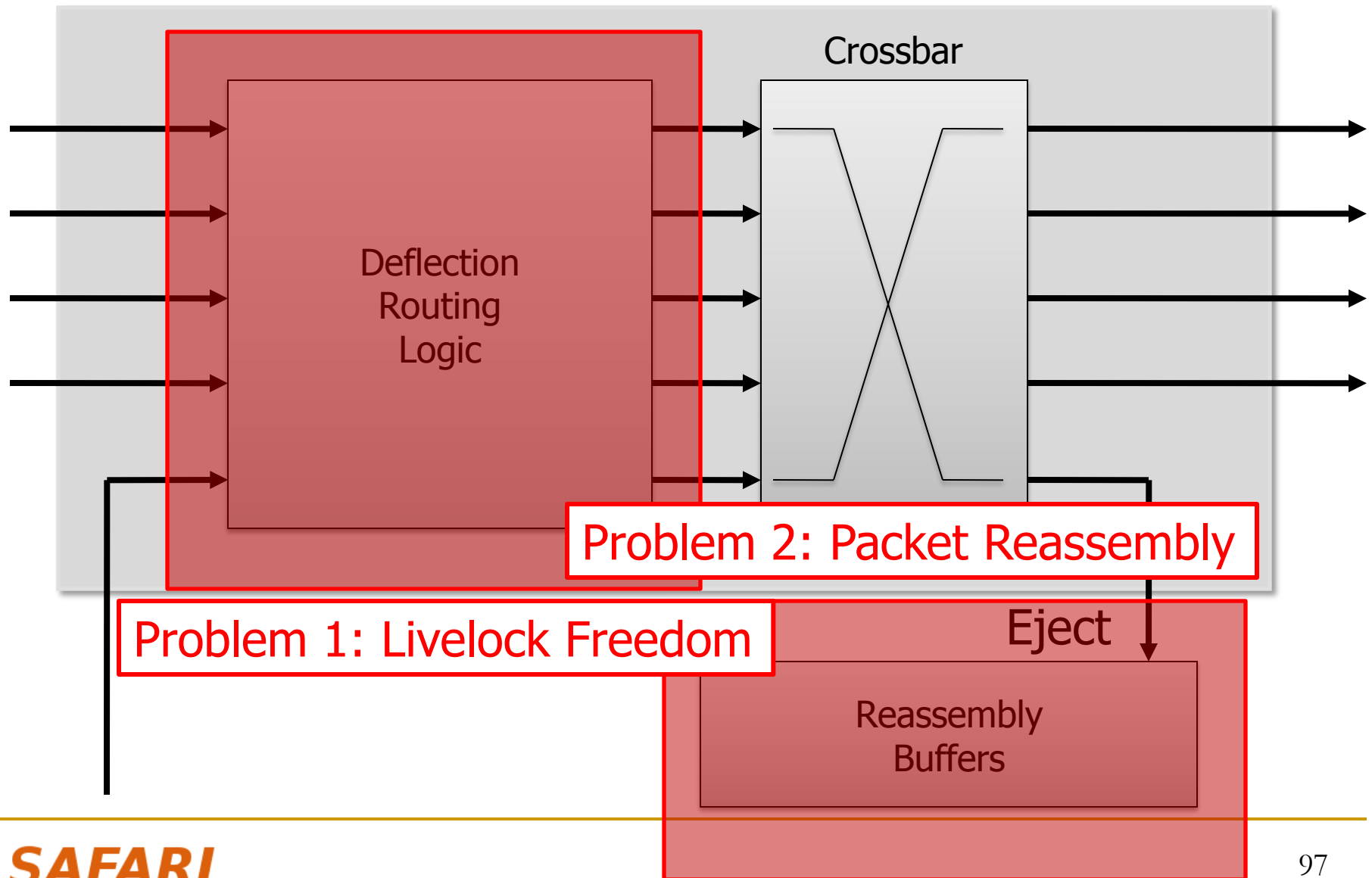


Packet: one or multiple flits



0 1 2 3

A Bufferless Router: A High-Level View



Complexity in Bufferless Deflection Routers

1. Must provide livelock freedom

Flits are sorted by age, then assigned in age order to output ports

→ **43% longer critical path than buffered router**

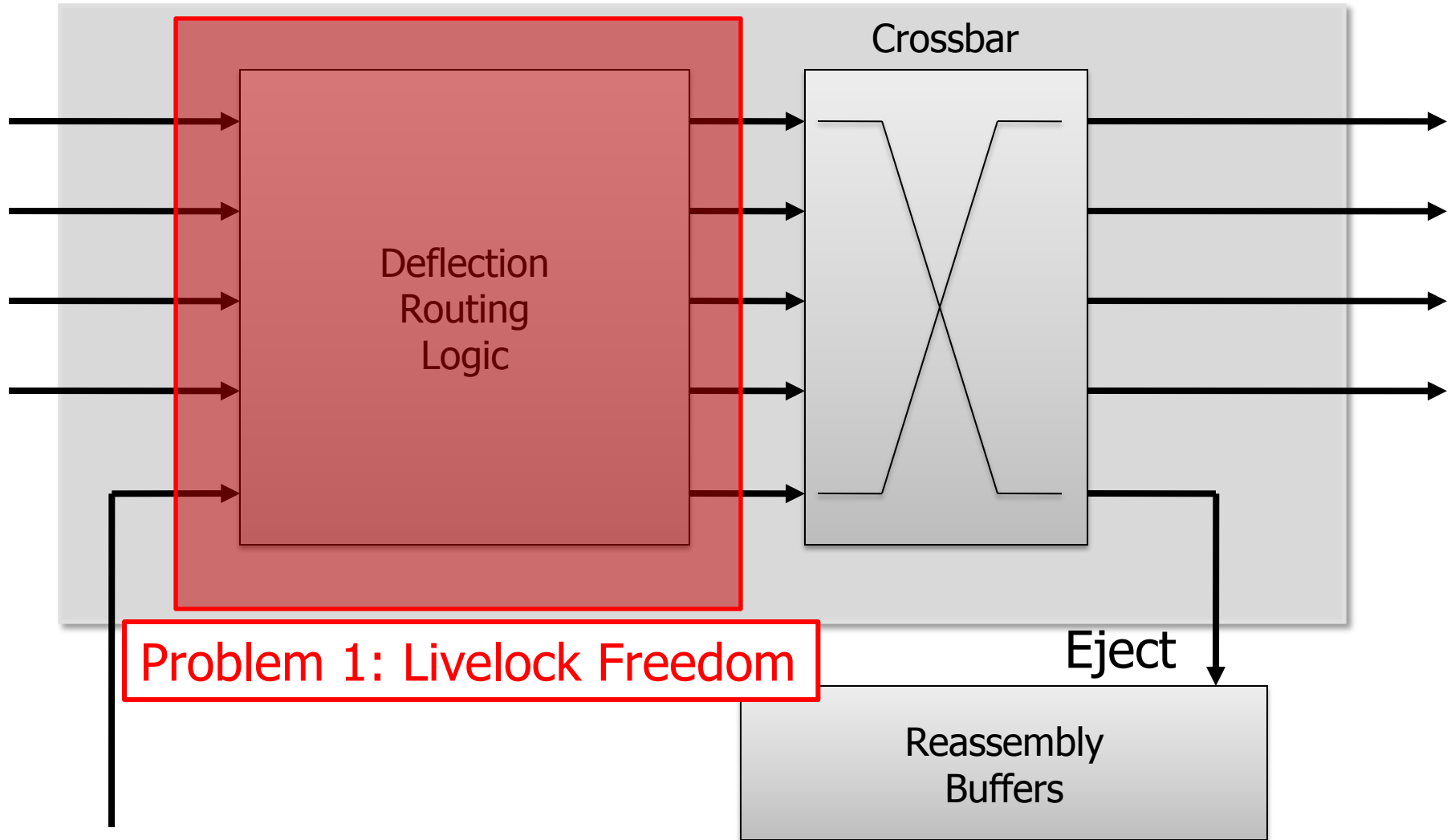
2. Must reassemble packets upon arrival

Reassembly buffers must be sized for worst case

→ **4KB per node**

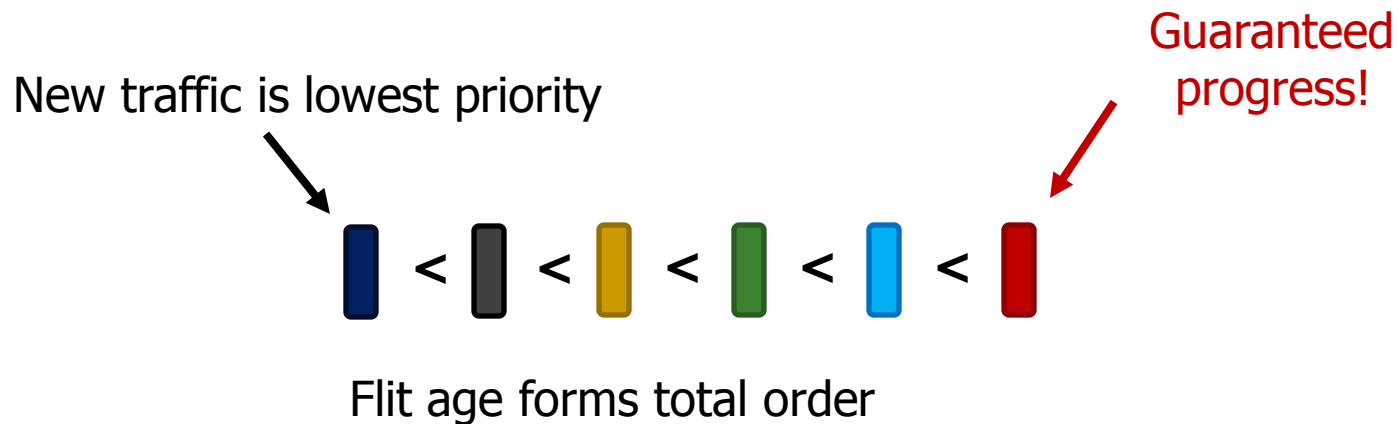
(8x8, 64-byte cache block)

Problem 1: Livelock Freedom



Livelock Freedom in Previous Work

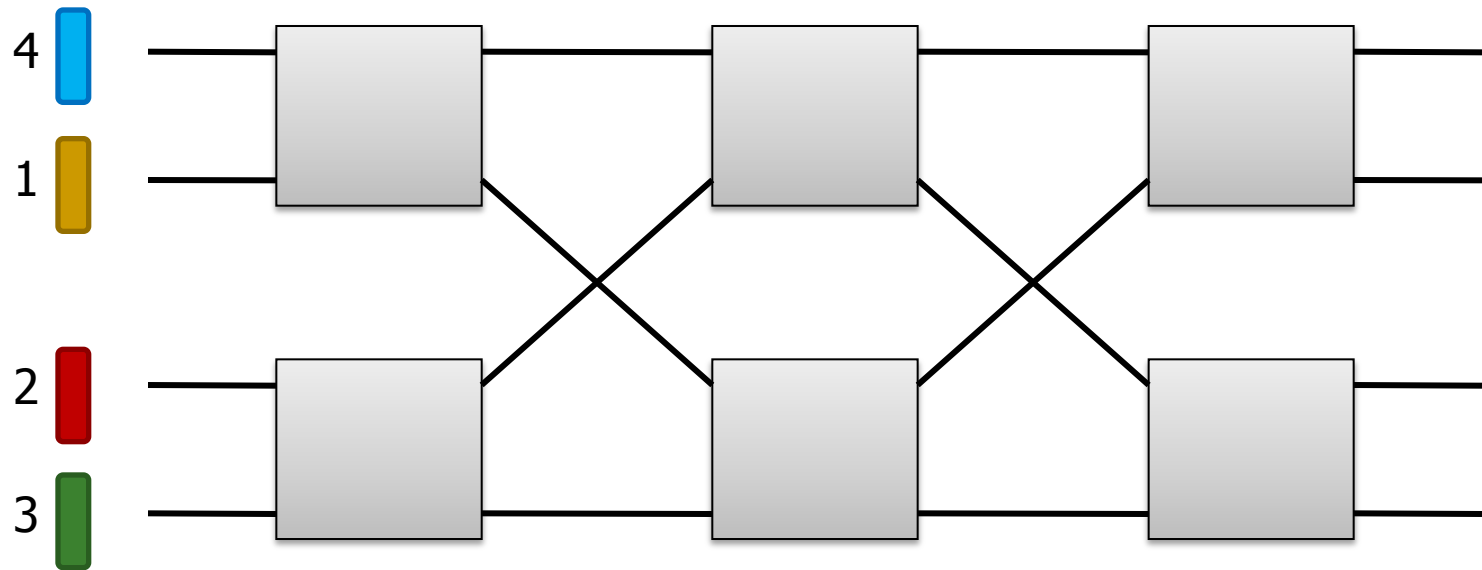
- What stops a flit from deflecting forever?
- All flits are **timestamped**
- **Oldest flits** are assigned their desired ports
- **Total order among flits**



- But what is the **cost** of this?

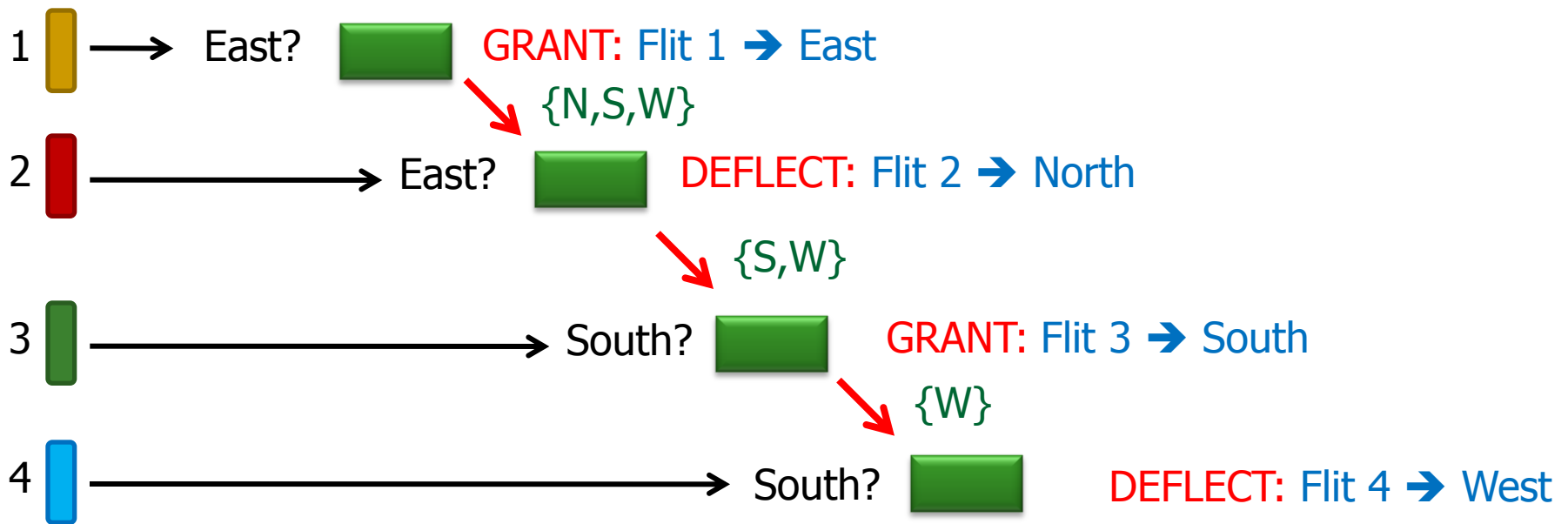
Age-Based Priorities are Expensive: Sorting

- Router must **sort flits by age**: long-latency sort network
 - **Three comparator stages** for 4 flits



Age-Based Priorities Are Expensive: Allocation

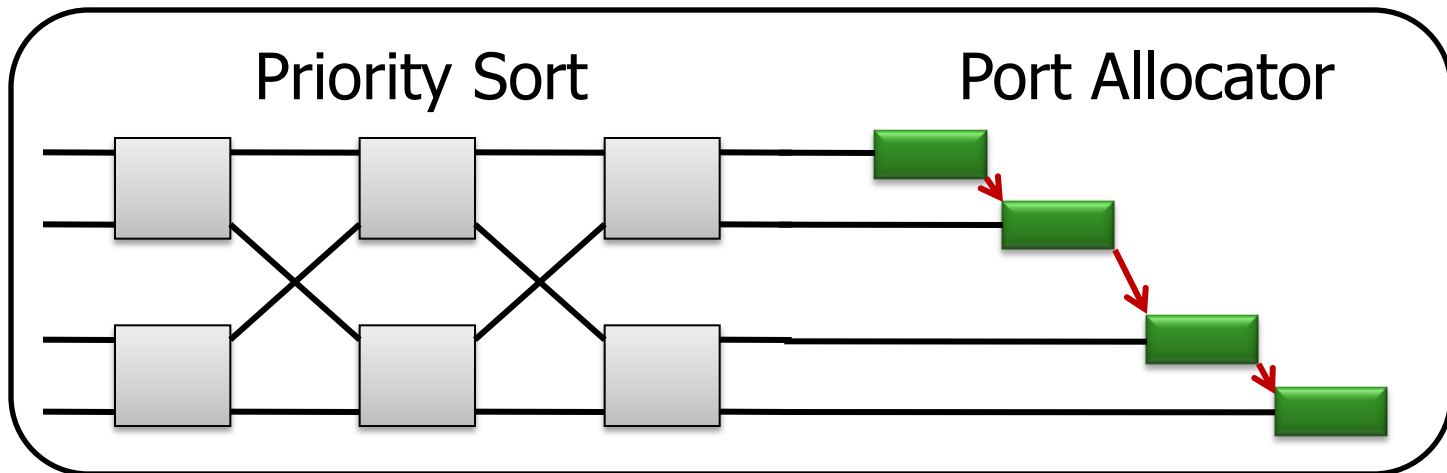
- After sorting, flits assigned to output ports in priority order
- Port assignment of younger flits depends on that of older flits
 - **sequential dependence** in the port allocator



Age-Ordered Flits

Age-Based Priorities Are Expensive

- Overall, **deflection routing logic** based on **Oldest-First** has a **43% longer critical path** than a buffered router



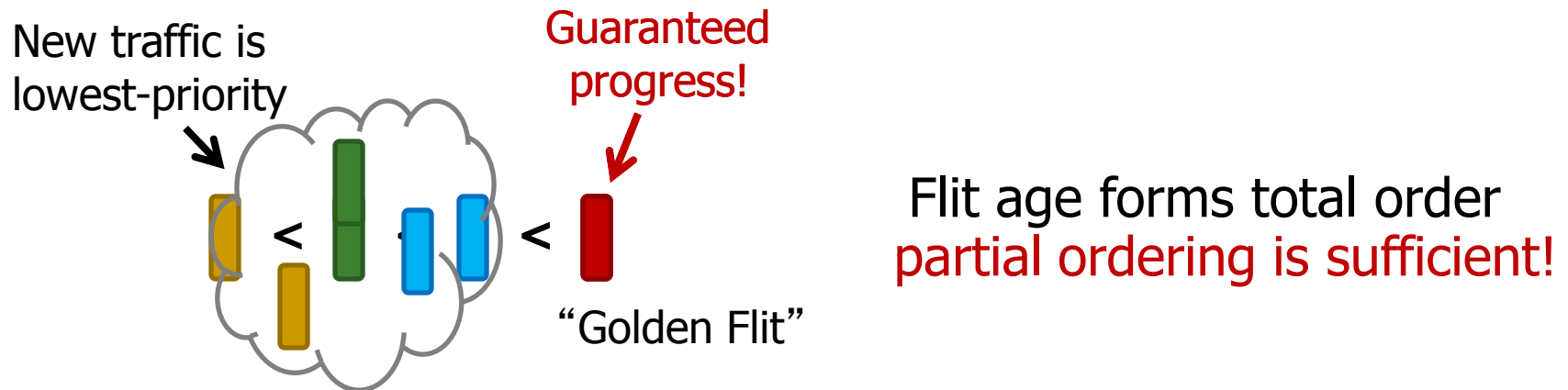
- Question: is there a cheaper way to route while guaranteeing livelock-freedom?

Solution: Golden Packet for Livelock Freedom

- What is *really necessary* for livelock freedom?

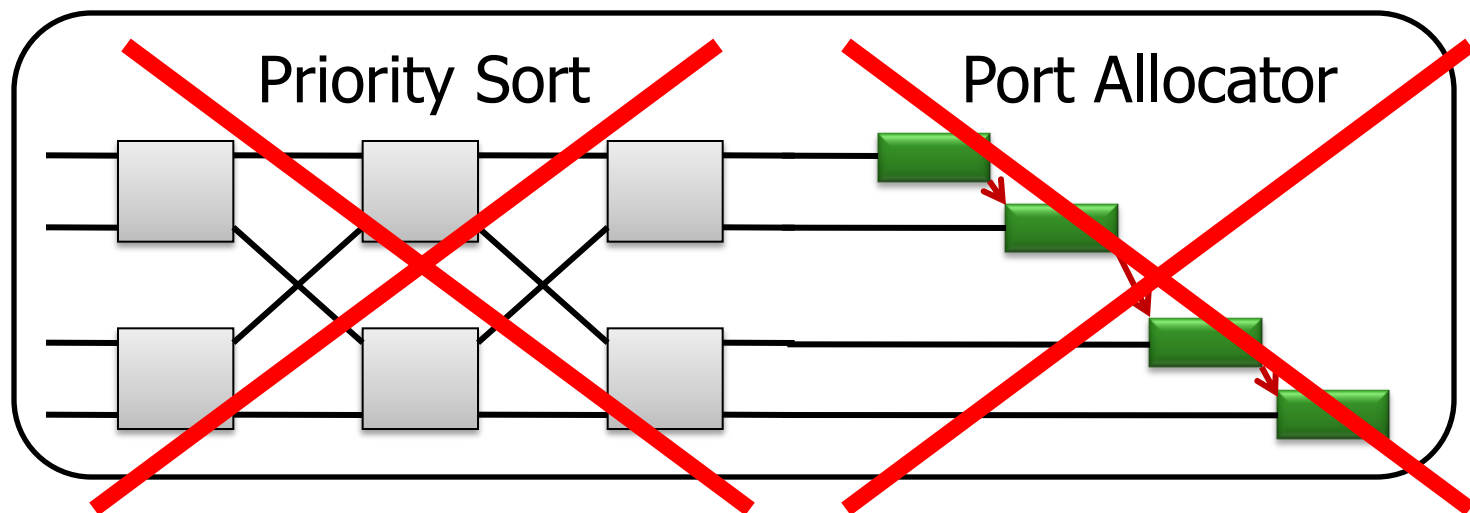
Key Insight: No total order. it is enough to:

1. Pick one flit to **prioritize until arrival**
2. Ensure any flit is **eventually** picked



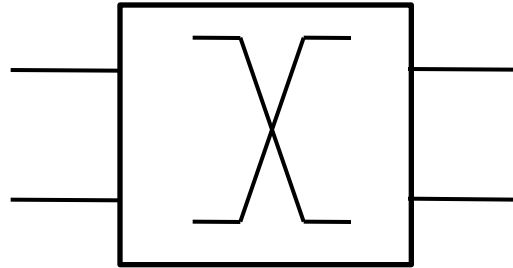
What Does Golden Flit Routing Require?

- Only **need** to properly route the Golden Flit
- **First Insight:** no need for full sort
- **Second Insight:** no need for sequential allocation



Golden Flit Routing With Two Inputs

- Let's route the Golden Flit in a two-input router first

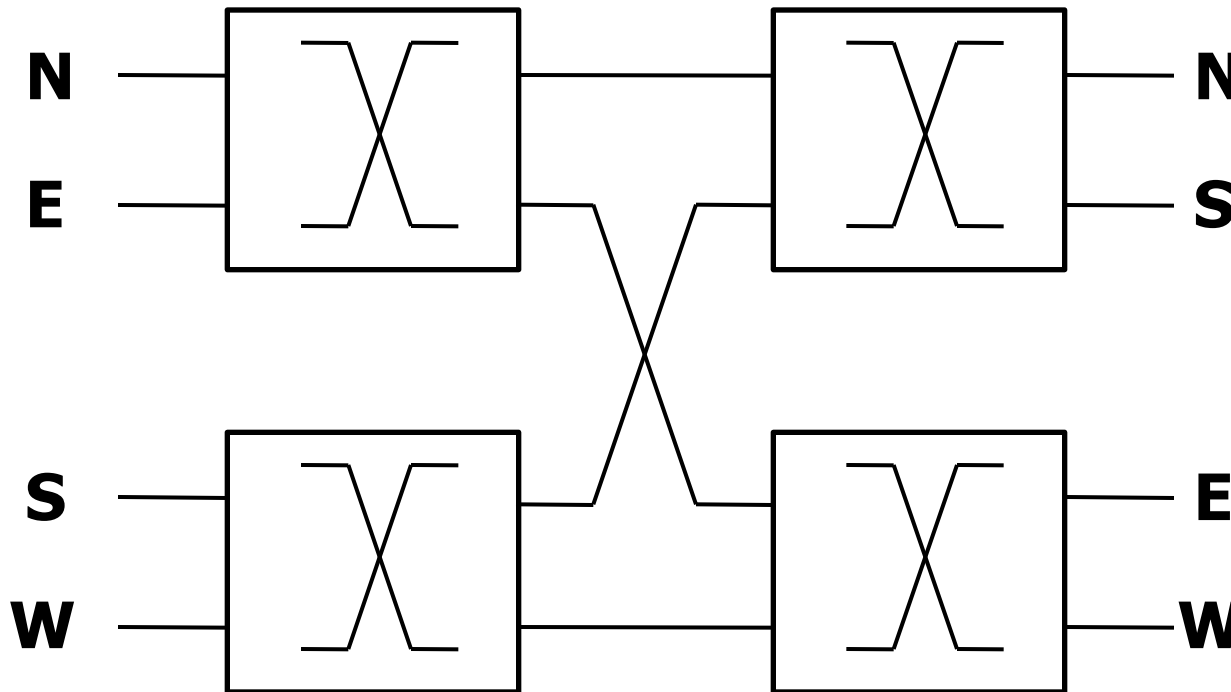


- **Step 1:** pick a “winning” flit: Golden Flit, else random
- **Step 2:** steer the winning flit to its desired output and deflect other flit

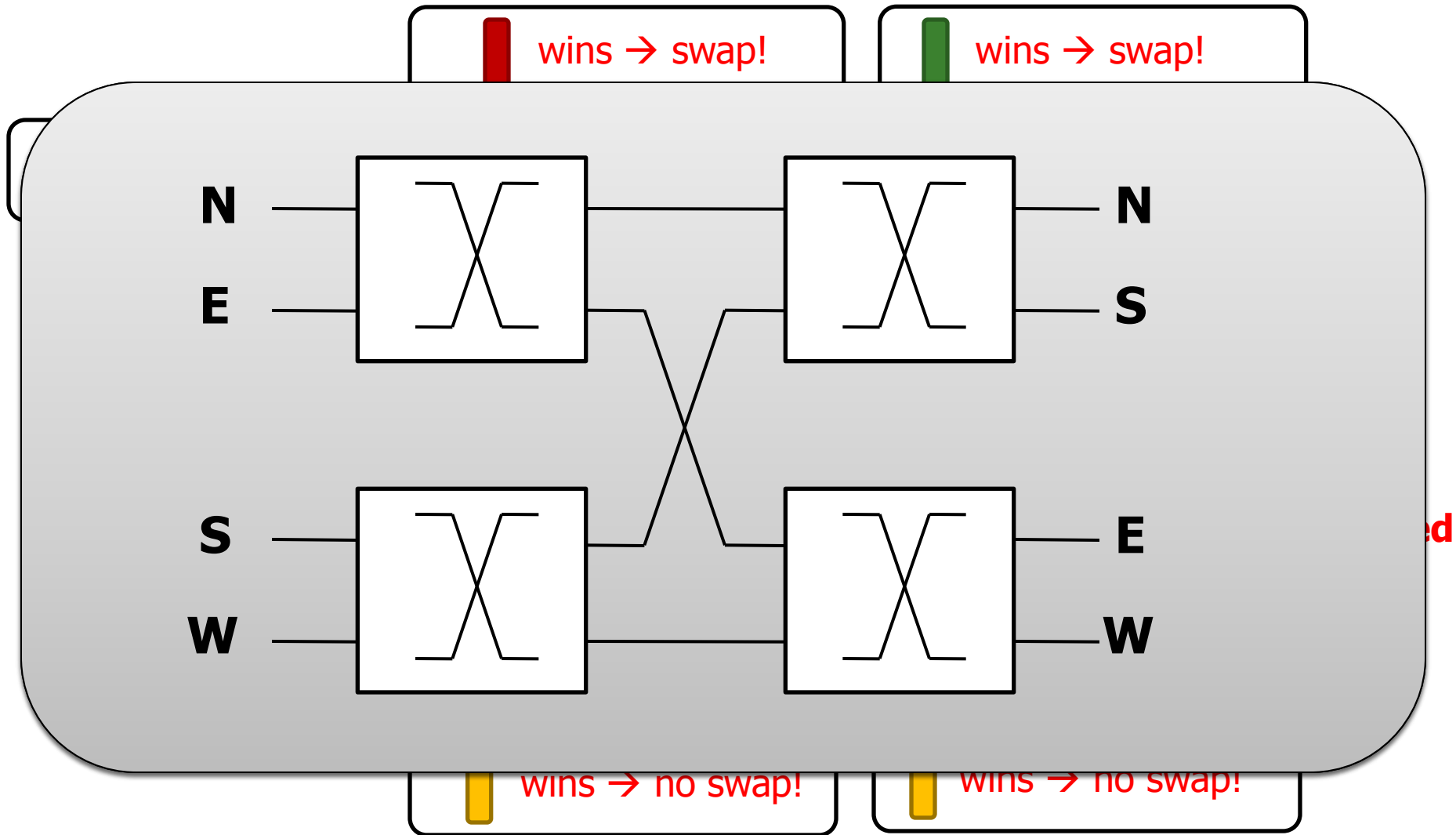
→ Golden Flit always routes toward destination

Golden Flit Routing with Four Inputs

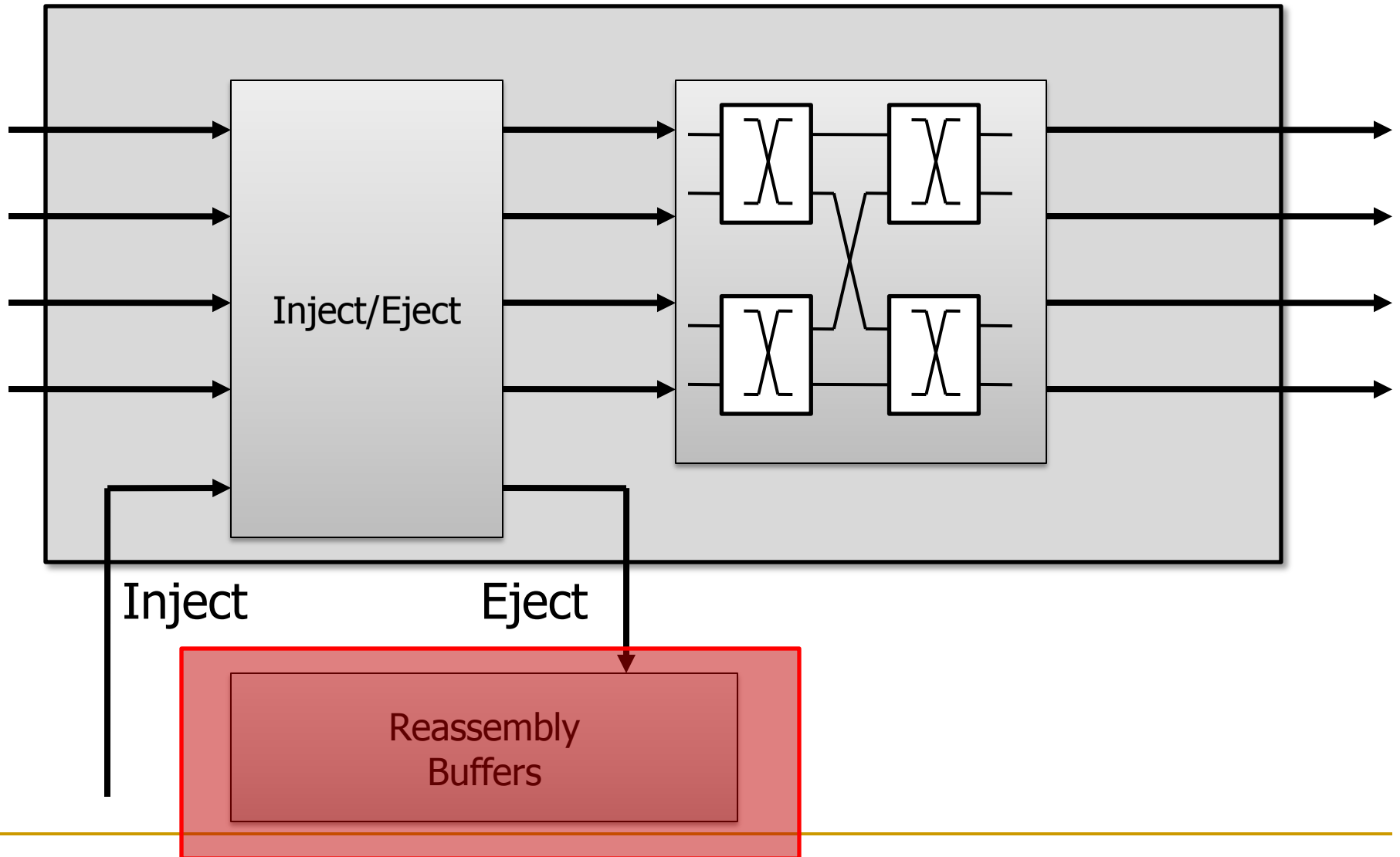
- Each block makes decisions **independently!**
 - **Deflection is a distributed decision**



Permutation Network Operation

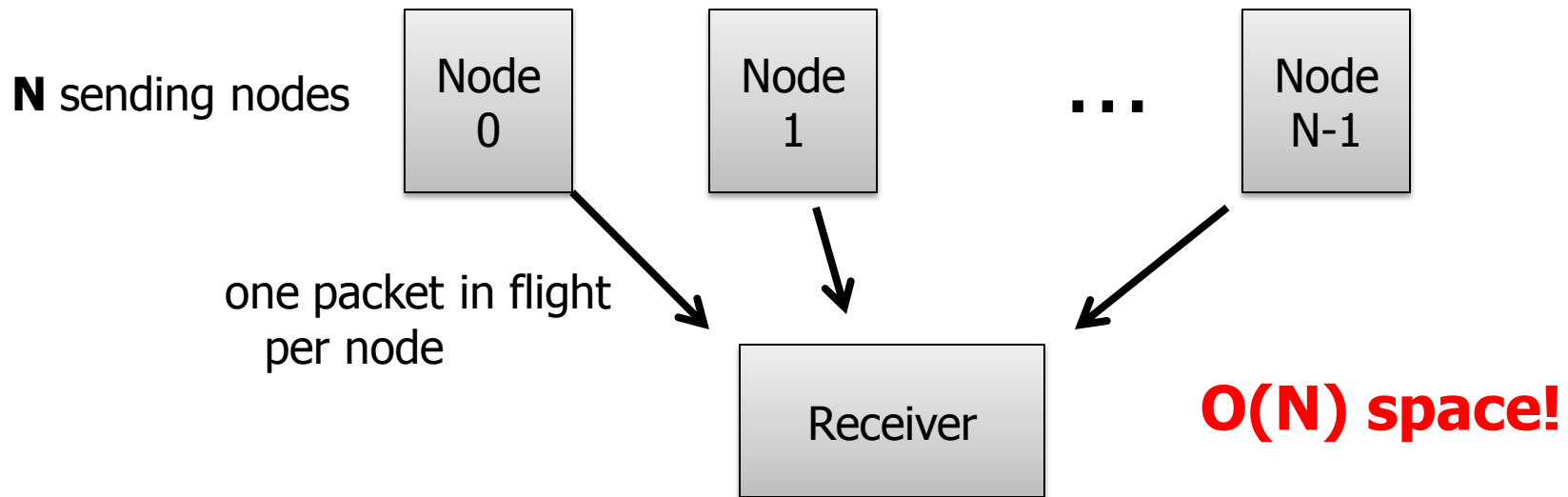


Problem 2: Packet Reassembly



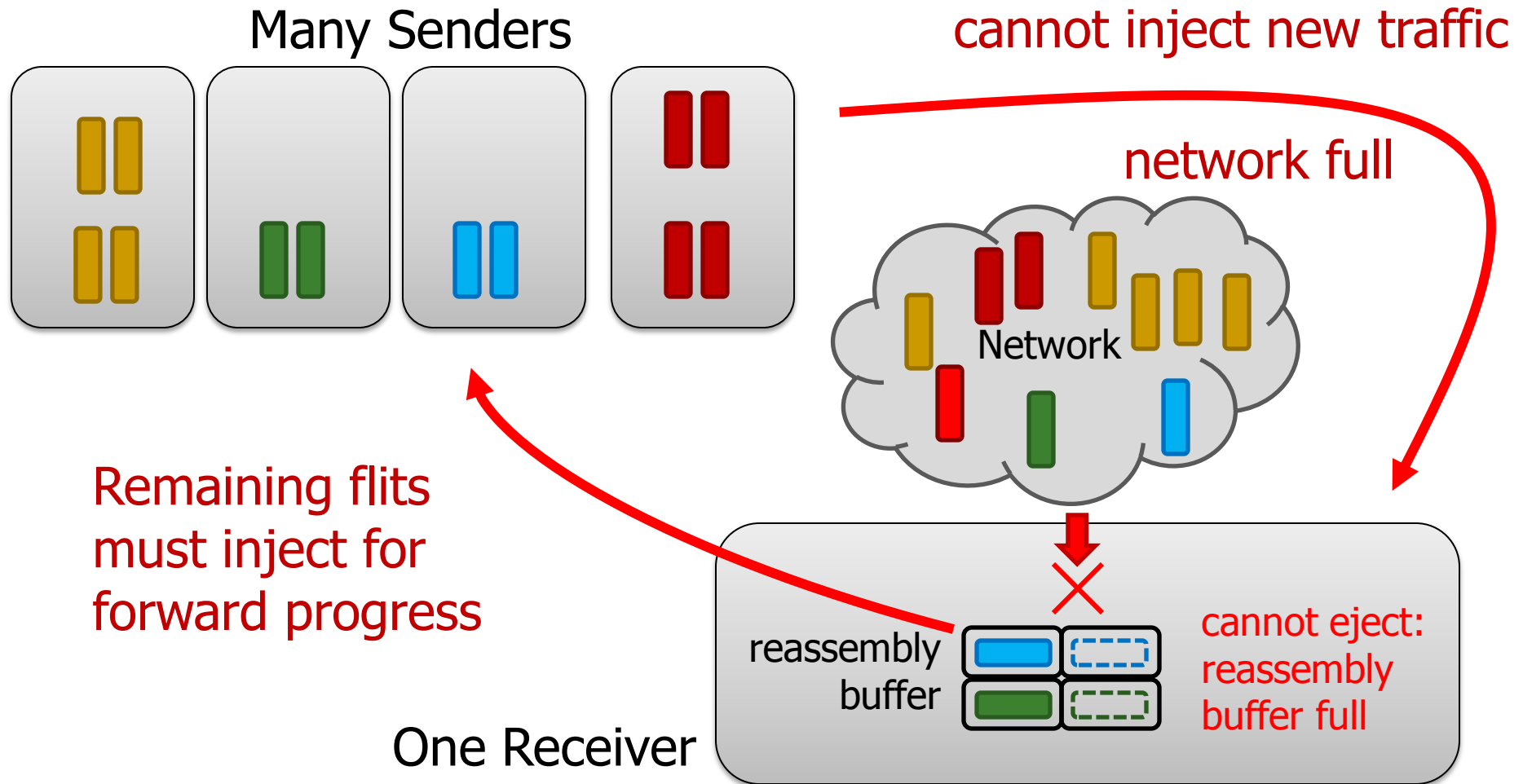
Reassembly Buffers are Large

- **Worst case:** every node sends a packet to one receiver
- Why can't we **make reassembly buffers smaller?**



Small Reassembly Buffers Cause Deadlock

- What happens when reassembly buffer is too small?

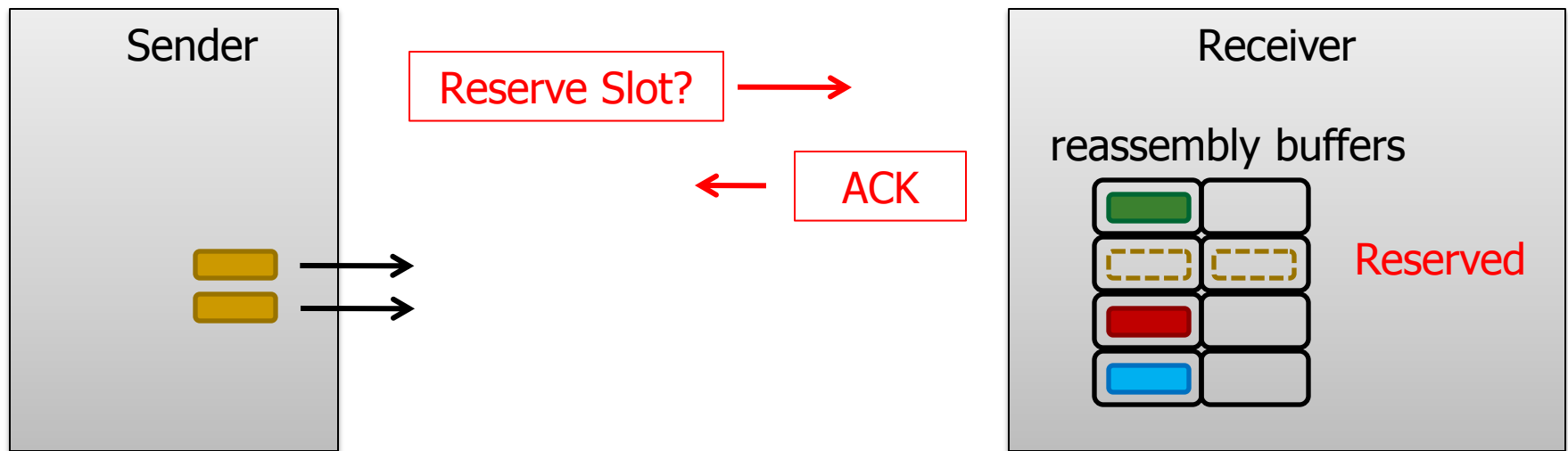


Reserve Space to Avoid Deadlock?

- What if every sender **asks permission** from the receiver before it sends?

→ **adds additional delay** to every request

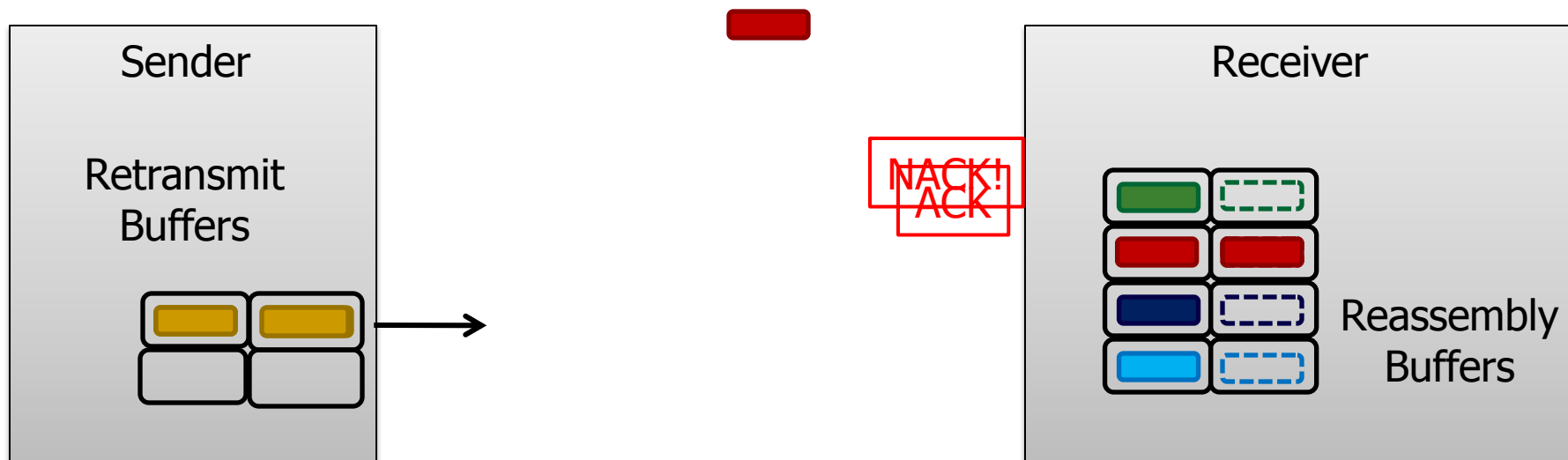
- 1. Reserve Slot
- 2. ACK
- 3. Send Packet



Escaping Deadlock with Retransmissions

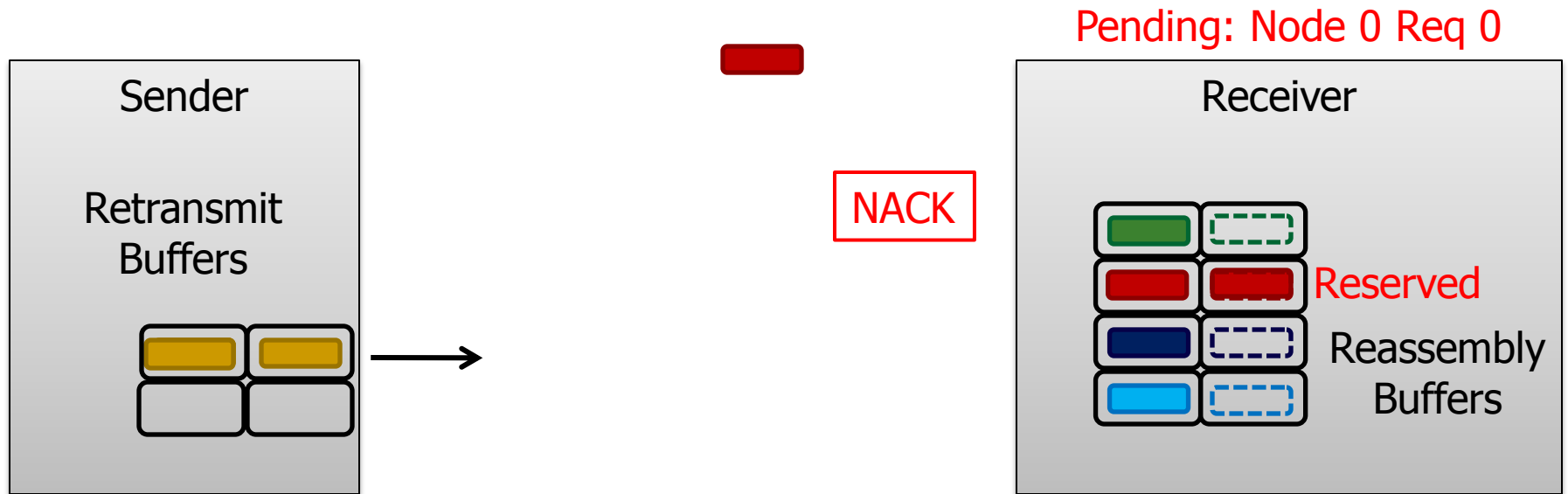
- Sender is optimistic instead: assume buffer is free
 - If not, receiver **drops** and NACKs; sender **retransmits**

- **no additional delay in best case**
 - **transmit buffering overhead for all packets**
 - **potentially many retransmits**
1. Send (2 flits)
 2. Drop, NACK
 3. Other packet completes
 4. Retransmit packet
 5. ACK
 6. Sender frees data

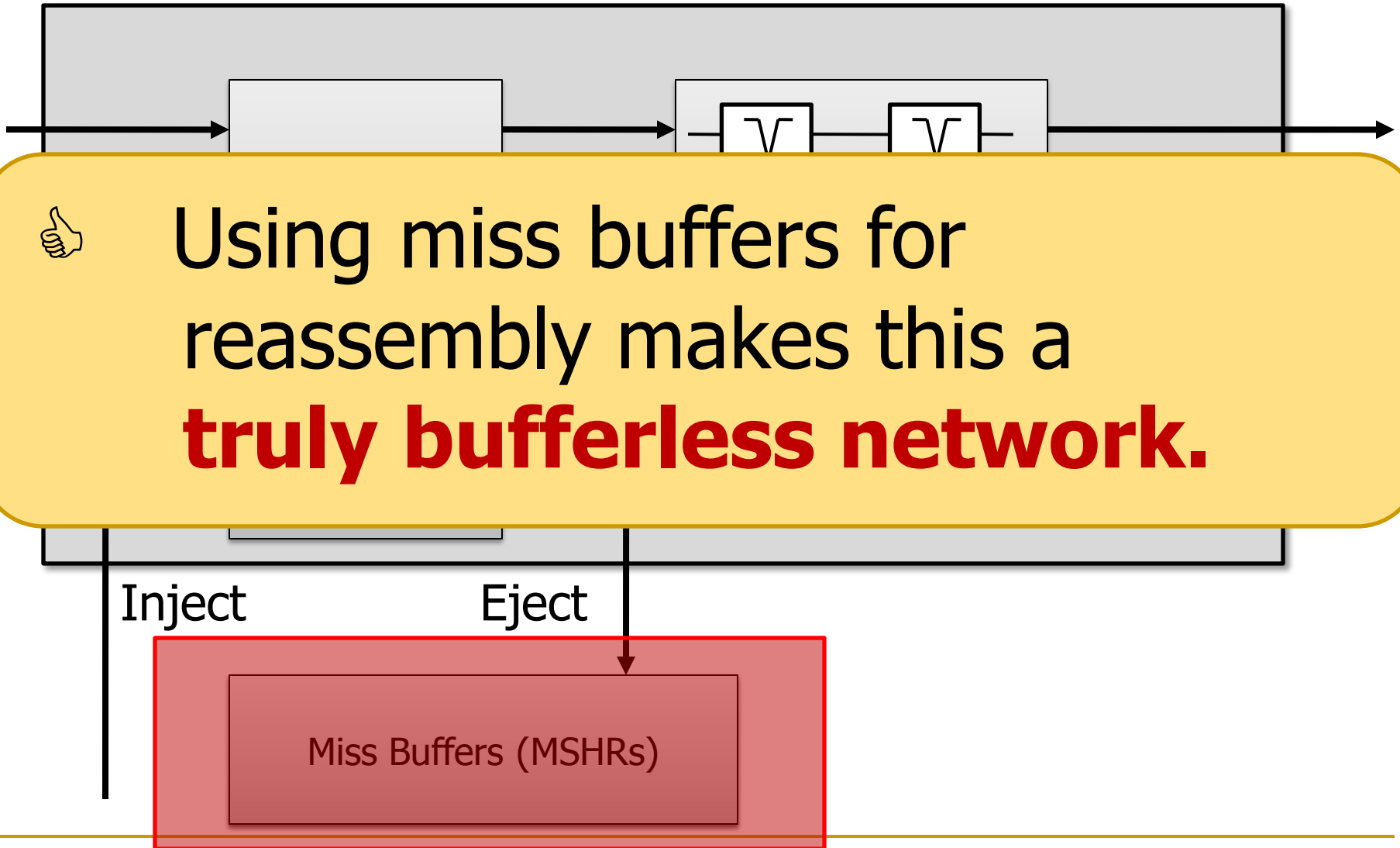


Solution: Retransmitting Only Once

- **Key Idea:** Retransmit only **when space becomes available.**
 - Receiver **drops packet** if full; **notes** which packet it drops
 - When space frees up, receiver **reserves space** so retransmit is successful
 - Receiver notifies sender to retransmit

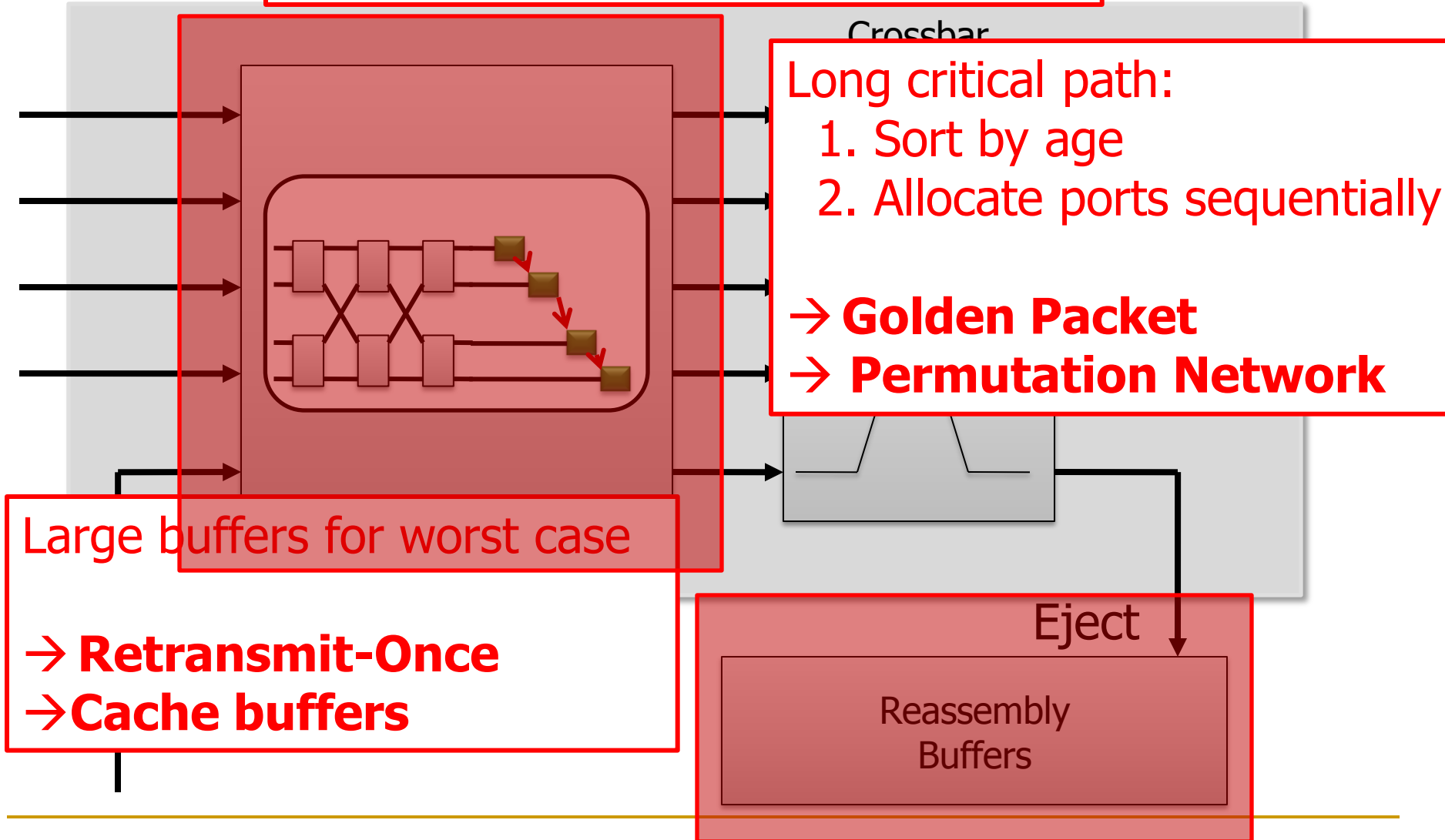


Using MSHRs as Reassembly Buffers

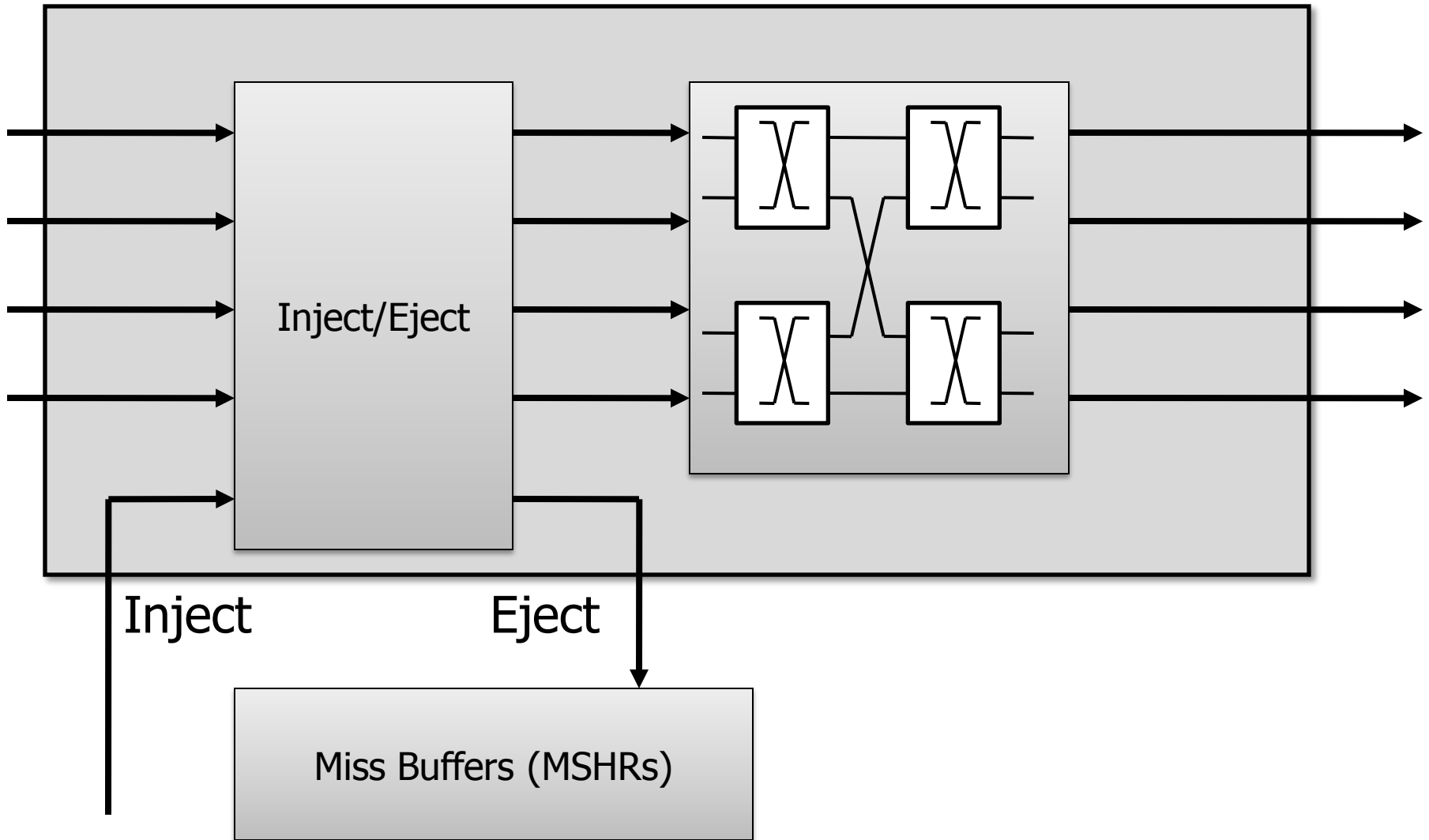


CHIPPER: Cheap Interconnect Partially-Permuting Router

Baseline Bufferless Deflection Router



CHIPPER: Cheap Interconnect Partially-Permuting Router



EVALUATION

Methodology

- **Multiprogrammed** workloads: CPU2006, server, desktop
 - 8x8 (64 cores), 39 homogeneous and 10 mixed sets
- **Multithreaded** workloads: SPLASH-2, 16 threads
 - 4x4 (16 cores), 5 applications
- **System configuration**
 - **Buffered** baseline: 2-cycle router, 4 VCs/channel, 8 flits/VC
 - **Bufferless** baseline: 2-cycle latency, FLIT-BLESS
 - Instruction-trace driven, closed-loop, 128-entry OoO window
 - 64KB L1, **perfect L2 (stresses interconnect)**, XOR mapping

Methodology

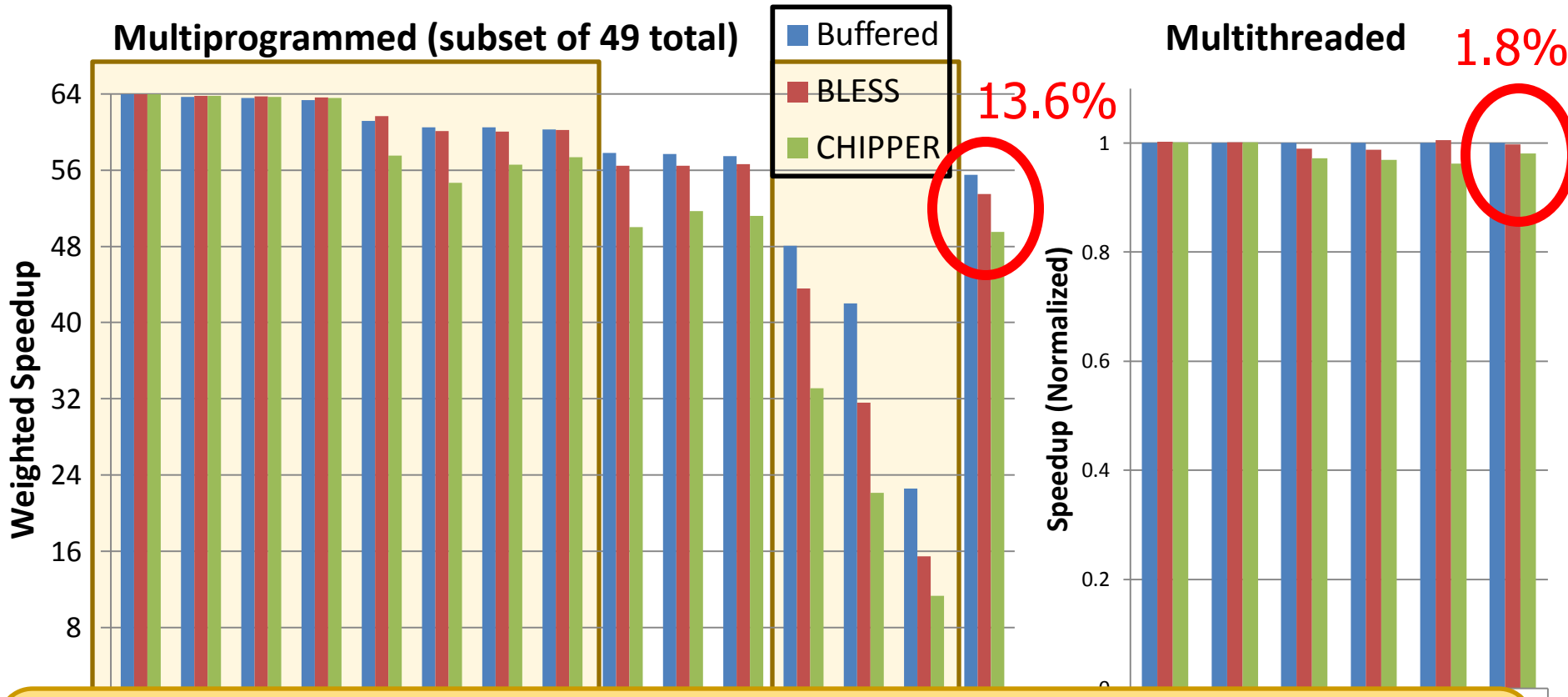
- **Hardware modeling**

- Verilog models for CHIPPER, BLESS, buffered logic
 - Synthesized with commercial 65nm library
- ORION for crossbar, buffers and links

- **Power**

- Static and dynamic power from hardware models
- Based on event counts in cycle-accurate simulations

Results: Performance Degradation



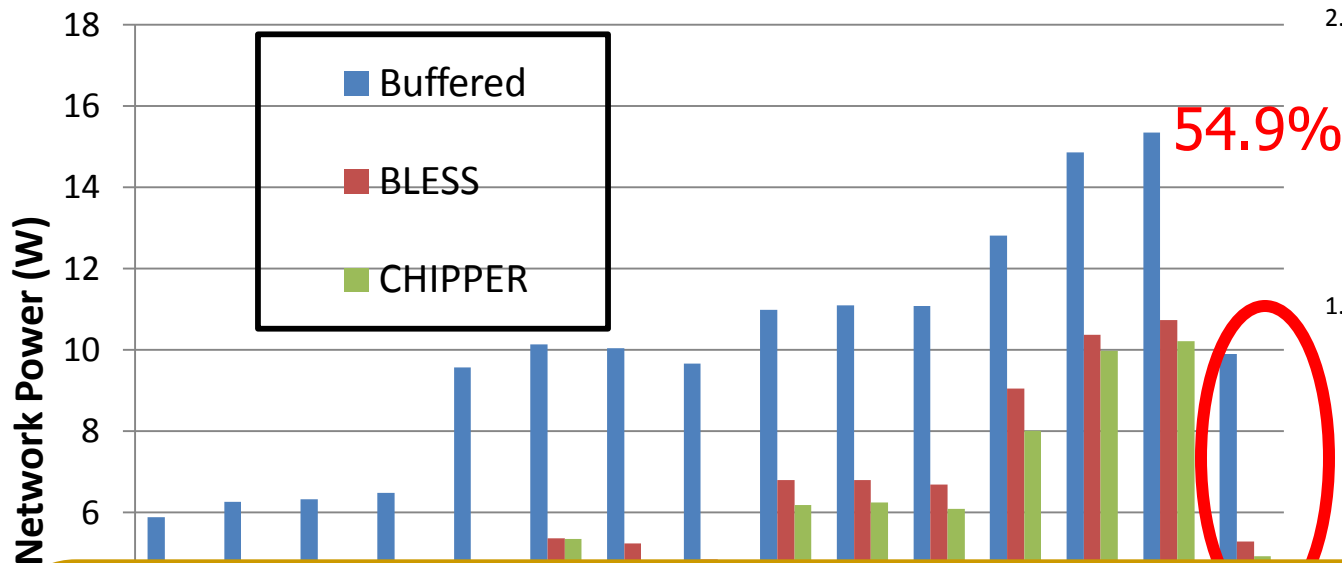
Minimal loss for low-to-medium-intensity workloads

5.0%

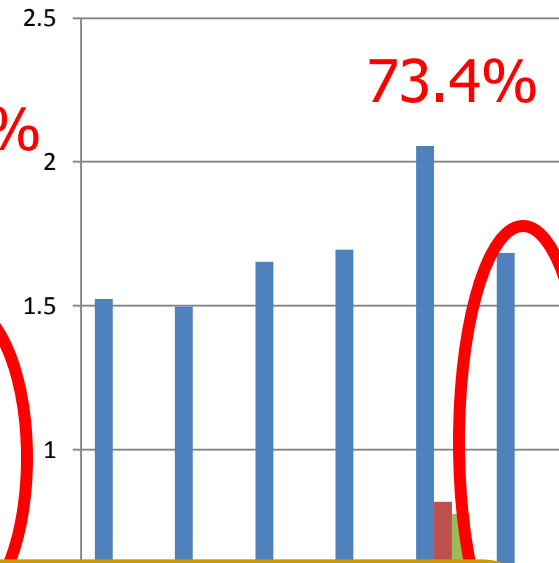
49.8%^A

Results: Power Reduction

Multiprogrammed (subset of 49 total)



Multithreaded



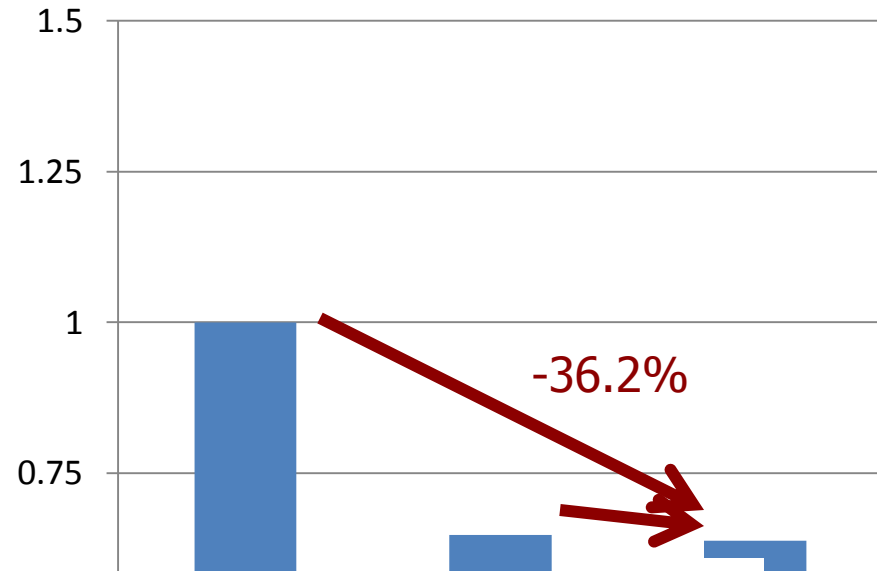
Removing buffers → majority of power savings



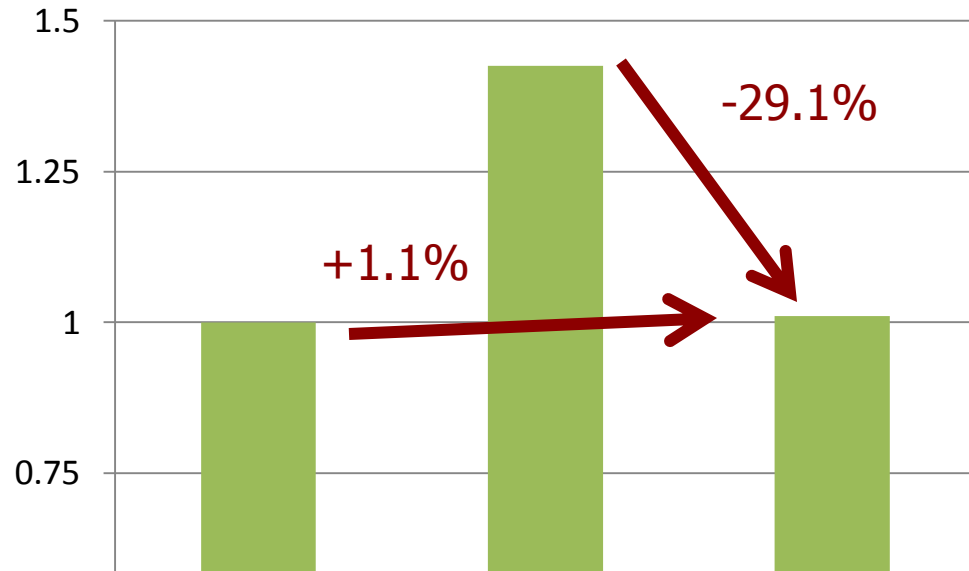
Slight savings from BLESS to CHIPPER

Results: Area and Critical Path Reduction

Normalized Router Area



Normalized Critical Path



CHIPPER maintains area savings of BLESS



Critical path **becomes competitive** to buffered

Conclusions

- Two key issues in **bufferless deflection routing**
 - livelock freedom and packet reassembly
- Bufferless deflection routers were high-complexity and impractical
 - **Oldest-first prioritization** → long critical path in router
 - **No end-to-end flow control for reassembly** → prone to deadlock with reasonably-sized reassembly buffers
- CHIPPER is a new, practical bufferless deflection router
 - **Golden packet prioritization** → short critical path in router
 - **Retransmit-once protocol** → deadlock-free packet reassembly
 - **Cache miss buffers as reassembly buffers** → truly bufferless network
- CHIPPER frequency comparable to buffered routers at much lower area and power cost, and minimal performance loss

MinBD: Minimally-Buffered Deflection Routing for Energy-Efficient Interconnect

Chris Fallin, Greg Nazario, Xiangyao Yu, Kevin Chang, Rachata Ausavarungnirun, and Onur Mutlu,

"MinBD: Minimally-Buffered Deflection Routing for Energy-Efficient Interconnect"

Proceedings of the 6th ACM/IEEE International Symposium on Networks on Chip (NOCS), Lyngby, Denmark, May 2012. Slides (pptx) (pdf)

SAFARI Carnegie Mellon University

Bufferless Deflection Routing

- **Key idea:** Packets are never buffered in the network. When two packets contend for the same link, one is **deflected**.
- Removing **buffers** yields significant benefits
 - Reduces **power** (CHIPPER: reduces NoC power by 55%)
 - Reduces **die area** (CHIPPER: reduces NoC area by 36%)
- But, at **high network utilization** (load), bufferless deflection routing causes **unnecessary link & router traversals**
 - Reduces **network throughput** and application performance
 - Increases **dynamic power**
- **Goal:** Improve **high-load performance** of low-cost deflection networks by reducing the deflection rate.

Outline: This Talk

- **Motivation**
- **Background:** Bufferless Deflection Routing
- **MinBD:** Reducing Deflections
 - Addressing Link Contention
 - Addressing the Ejection Bottleneck
 - Improving Deflection Arbitration
- **Results**
- **Conclusions**

Outline: This Talk

- **Motivation**
- **Background:** Bufferless Deflection Routing
- **MinBD:** Reducing Deflections
 - Addressing Link Contention
 - Addressing the Ejection Bottleneck
 - Improving Deflection Arbitration
- **Results**
- **Conclusions**

Issues in Bufferless Deflection Routing

- **Correctness:** Deliver all packets without **livelock**
 - **CHIPPER¹: Golden Packet**
 - **Globally prioritize one packet** until delivered
- **Correctness:** Reassemble packets without **deadlock**
 - **CHIPPER¹: Retransmit-Once**
- **Performance:** Avoid performance degradation at **high load**
 - **MinBD**

Key Performance Issues

- 1. Link contention:** no buffers to hold traffic → any link contention causes a deflection
→ use side buffers
- 2. Ejection bottleneck:** only one flit can eject per router per cycle → simultaneous arrival causes deflection
→ eject up to 2 flits/cycle
- 3. Deflection arbitration:** practical (fast) deflection arbiters deflect unnecessarily
→ new priority scheme (silver flit)

Outline: This Talk

- **Motivation**
- **Background:** Bufferless Deflection Routing
- **MinBD:** Reducing Deflections
 - Addressing Link Contention
 - Addressing the Ejection Bottleneck
 - Improving Deflection Arbitration
- **Results**
- **Conclusions**

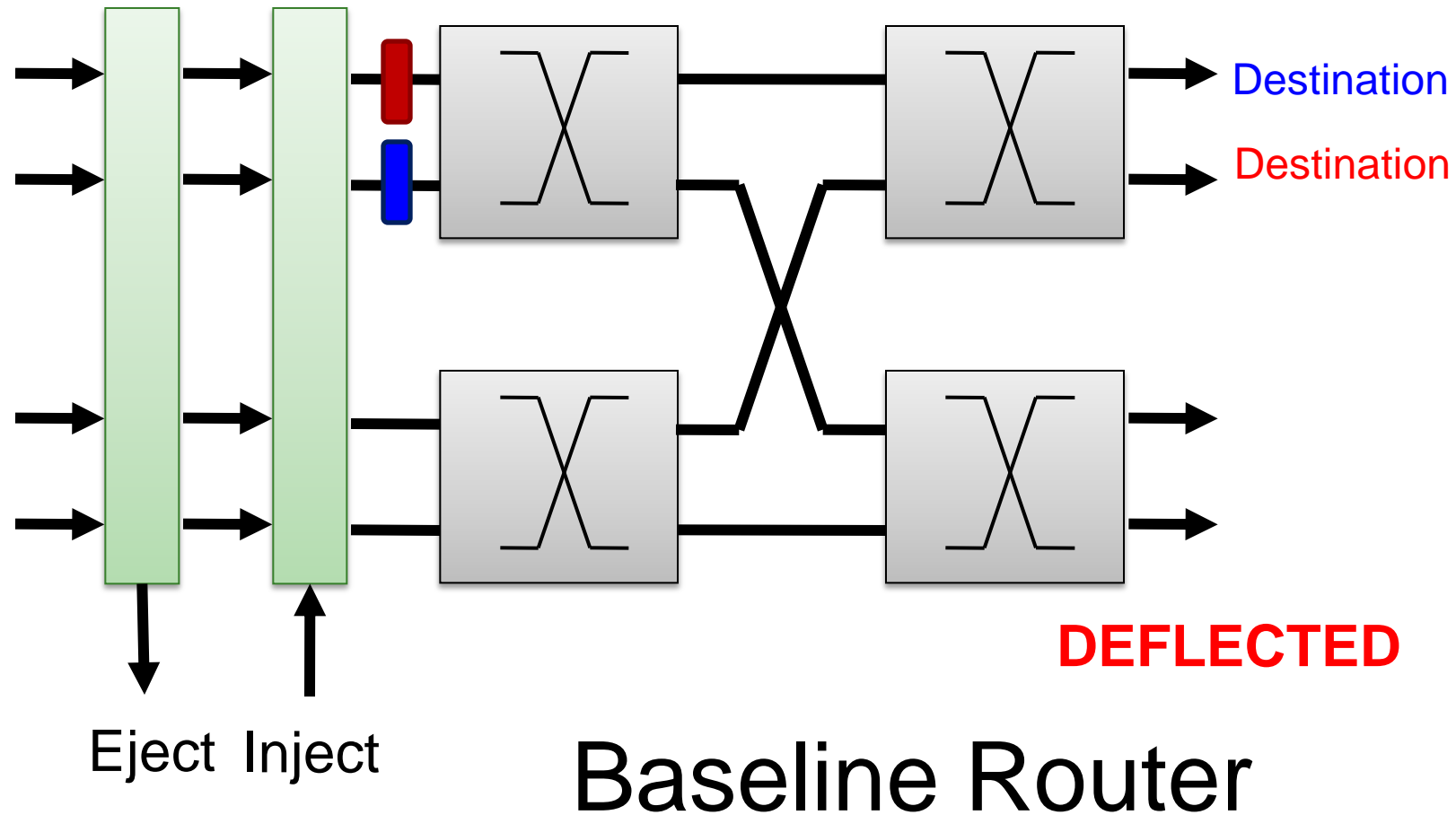
Outline: This Talk

- **Motivation**
- **Background:** Bufferless Deflection Routing
- **MinBD:** Reducing Deflections
 - Addressing Link Contention
 - Addressing the Ejection Bottleneck
 - Improving Deflection Arbitration
- **Results**
- **Conclusions**

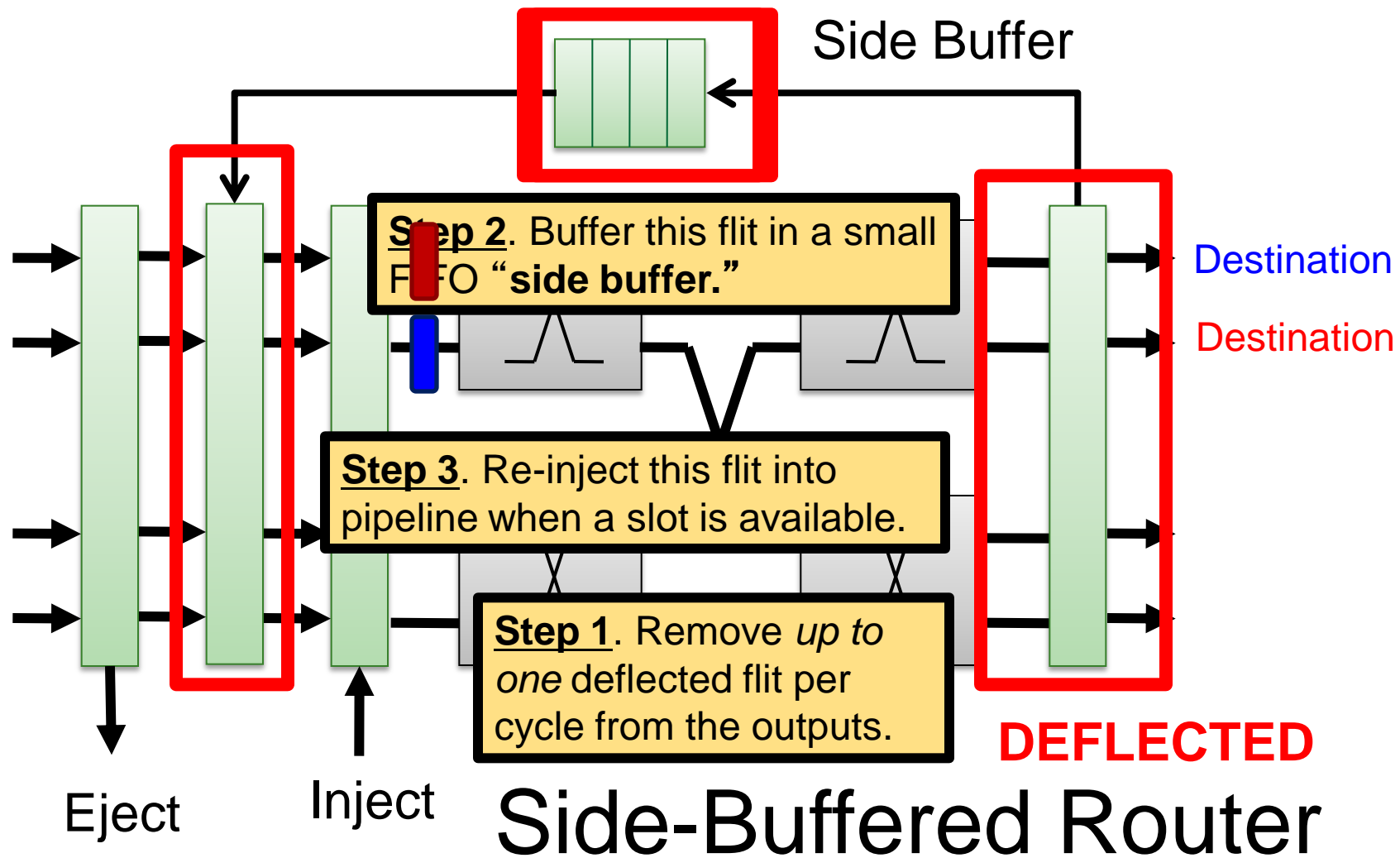
Addressing Link Contention

- **Problem 1:** Any link contention causes a deflection
- **Buffering** a flit can avoid deflection on contention
- But, **input buffers** are expensive:
 - All flits are buffered on every hop → **high dynamic energy**
 - Large buffers necessary → **high static energy** and **large area**
- **Key Idea 1:** add a **small buffer** to a bufferless deflection router to buffer **only** flits that **would have been deflected**

How to Buffer Deflected Flits



How to Buffer Deflected Flits



Why Could A Side Buffer Work Well?

- Buffer some flits and deflect other flits at **per-flit level**
 - Relative to **bufferless routers**, **deflection rate reduces** (need not deflect all contending flits)
 - 4-flit buffer reduces deflection rate by **39%**
 - Relative to **buffered routers**, **buffer is more efficiently used** (need not buffer all flits)
 - similar performance with **25%** of buffer space

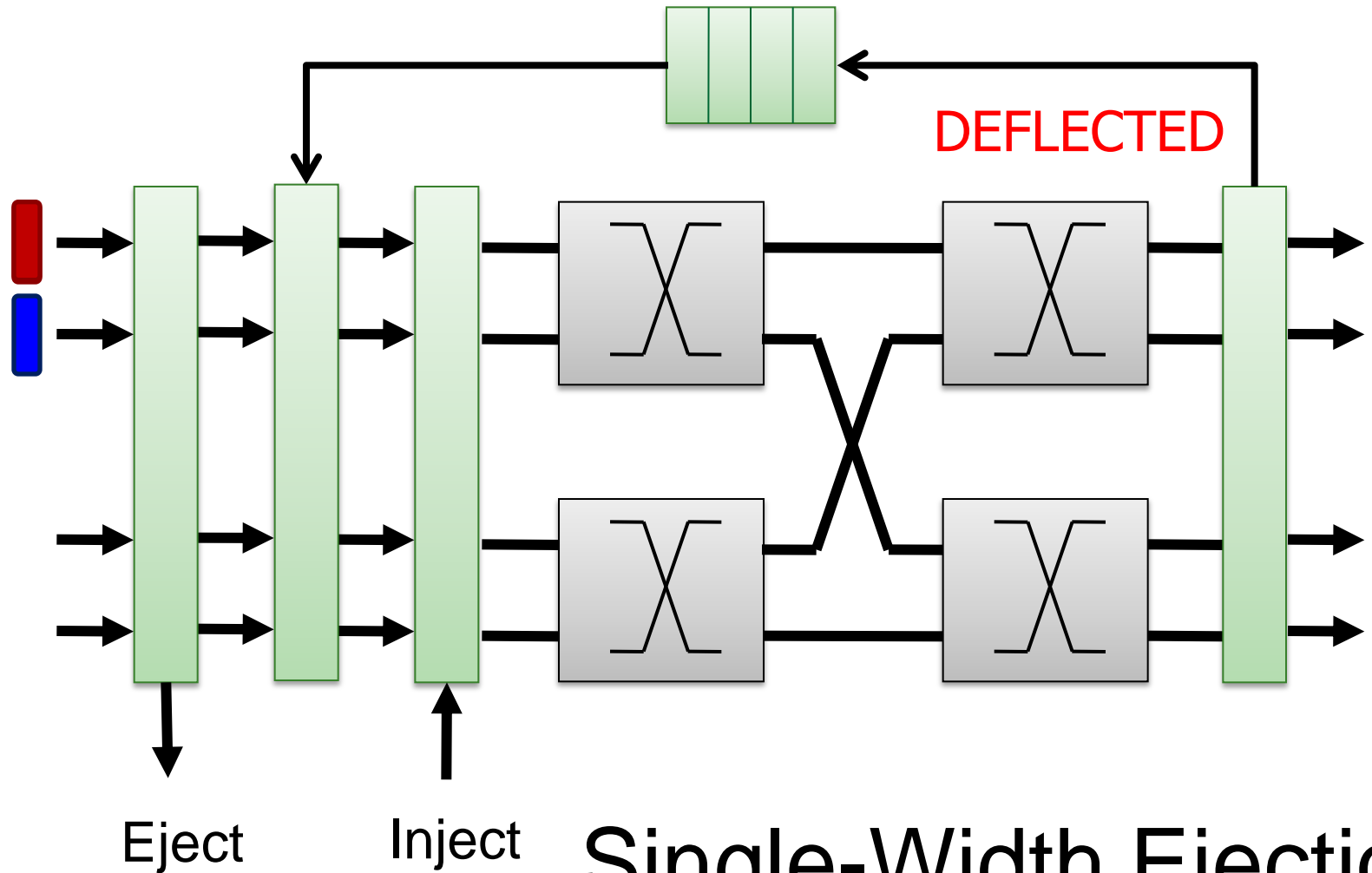
Outline: This Talk

- **Motivation**
- **Background:** Bufferless Deflection Routing
- **MinBD:** Reducing Deflections
 - Addressing Link Contention
 - Addressing the Ejection Bottleneck
 - Improving Deflection Arbitration
- **Results**
- **Conclusions**

Addressing the Ejection Bottleneck

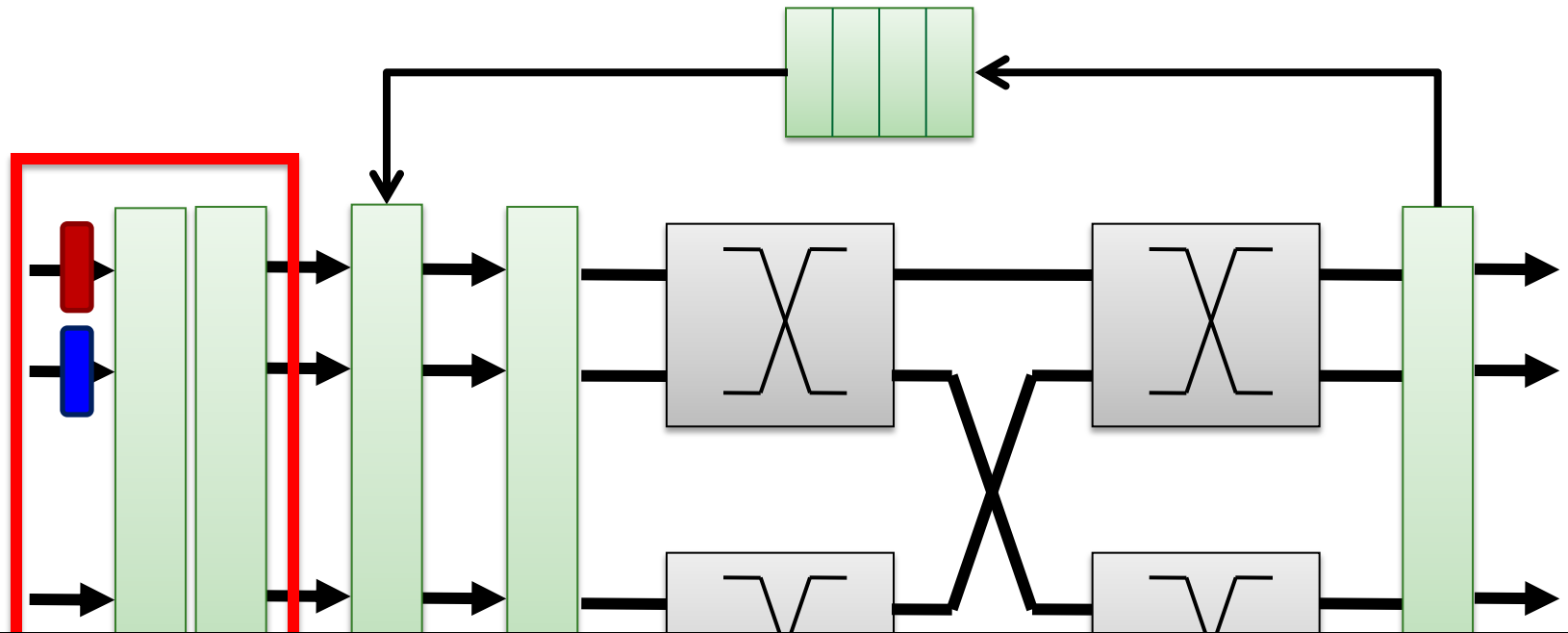
- **Problem 2:** Flits deflect unnecessarily because only one flit can **eject** per router per cycle
- In 20% of all ejections, ≥ 2 flits could have ejected
 - all but one flit must **deflect and try again**
 - these deflected flits cause additional contention
- Ejection width of 2 flits/cycle reduces **deflection rate 21%**
- **Key idea 2:** Reduce deflections due to a single-flit ejection port by allowing **two flits** to eject per cycle

Addressing the Ejection Bottleneck



Single-Width Ejection

Addressing the Ejection Bottleneck



For fair comparison, **baseline routers** have dual-width ejection for perf. (not power/area)

Eject

Inject

Dual-Width Ejection

Outline: This Talk

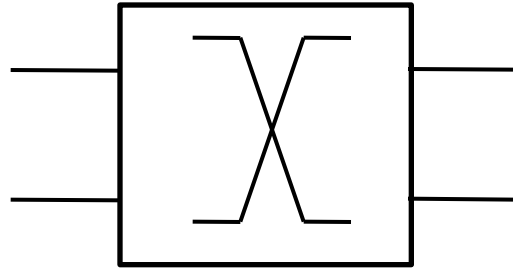
- **Motivation**
- **Background:** Bufferless Deflection Routing
- **MinBD:** Reducing Deflections
 - Addressing Link Contention
 - Addressing the Ejection Bottleneck
 - Improving Deflection Arbitration
- **Results**
- **Conclusions**

Improving Deflection Arbitration

- **Problem 3:** Deflections occur unnecessarily because fast arbiters must use simple priority schemes
- Age-based priorities (several past works): full priority order gives fewer deflections, but requires **slow arbiters**
- State-of-the-art deflection arbitration (Golden Packet & two-stage permutation network)
 - Prioritize one packet globally (**ensure forward progress**)
 - Arbitrate other flits randomly (**fast critical path**)
- Random common case leads to **uncoordinated arbitration**

Fast Deflection Routing Implementation

- Let's route in a two-input router first:

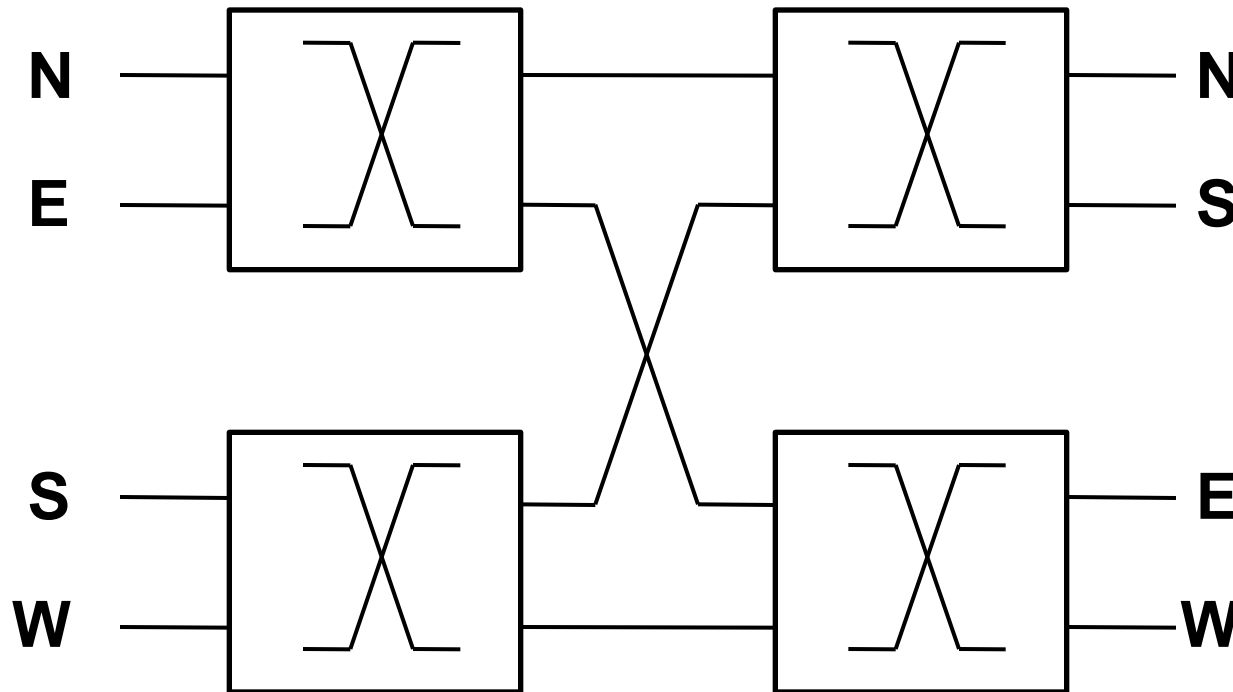


- **Step 1:** pick a “winning” flit (Golden Packet, else random)
- **Step 2:** steer the winning flit to its desired output and deflect other flit

→ Highest-priority flit always routes to destination

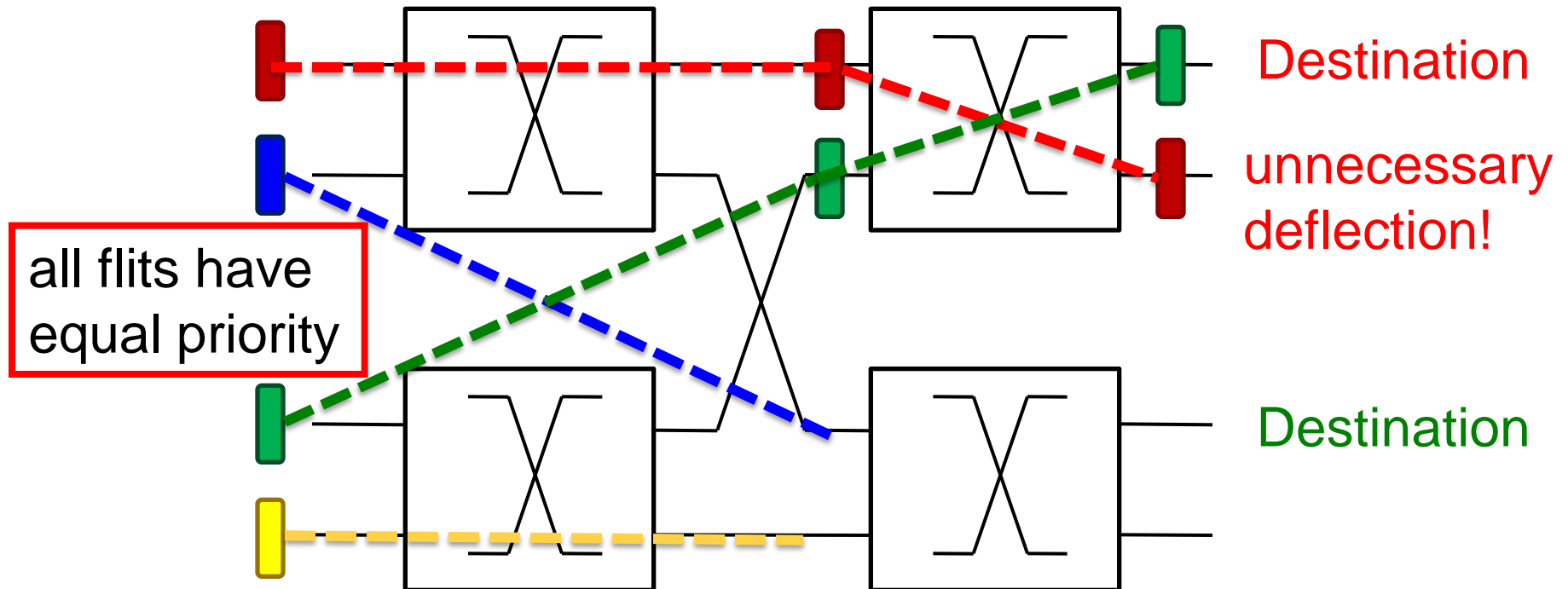
Fast Deflection Routing with Four Inputs

- Each block makes decisions **independently**
 - **Deflection is a distributed decision**



Unnecessary Deflections in Fast Arbiters

- How does lack of coordination cause unnecessary deflections?
 - No flit is golden (pseudorandom arbitration)
 - Red flit wins at first stage
 - Green flit loses at first stage (must be deflected now)
 - Red flit loses at second stage; Red and Green are deflected



Improving Deflection Arbitration

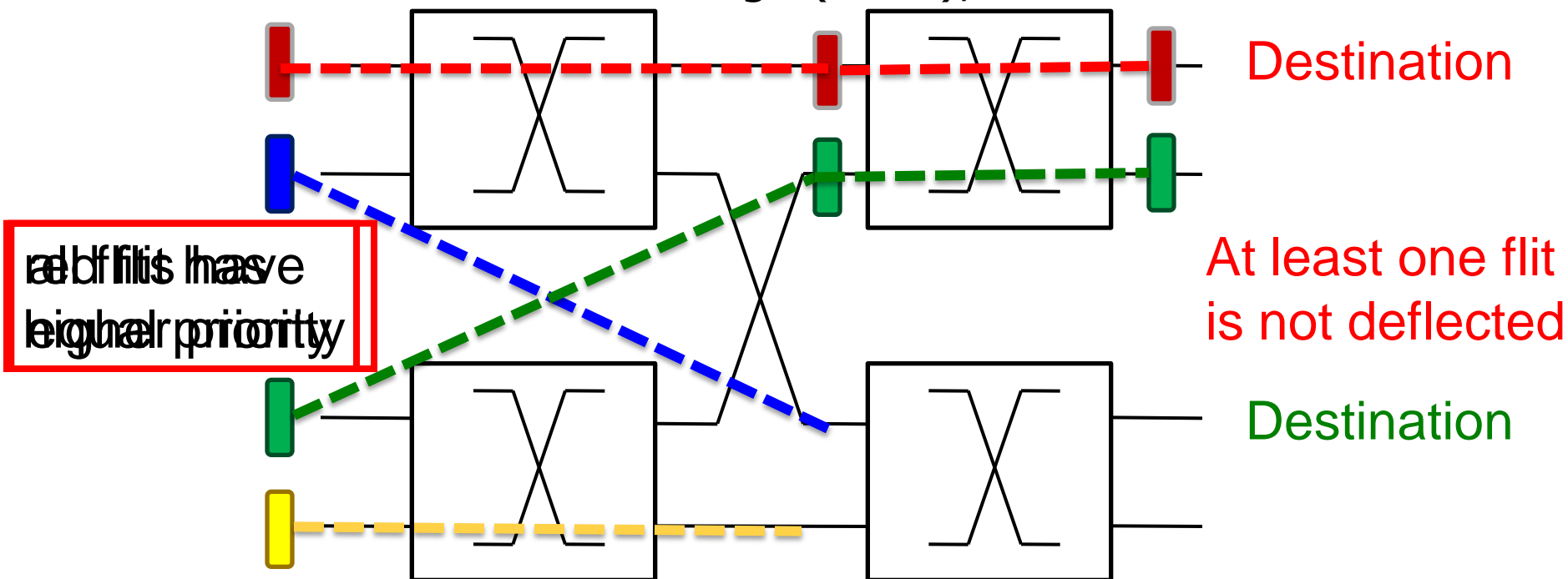
- **Key idea 3: Add a priority level** and prioritize one flit to ensure at least **one flit is not deflected in each cycle**

- **Highest priority:** one **Golden Packet** in network
 - Chosen in static round-robin schedule
 - **Ensures correctness**

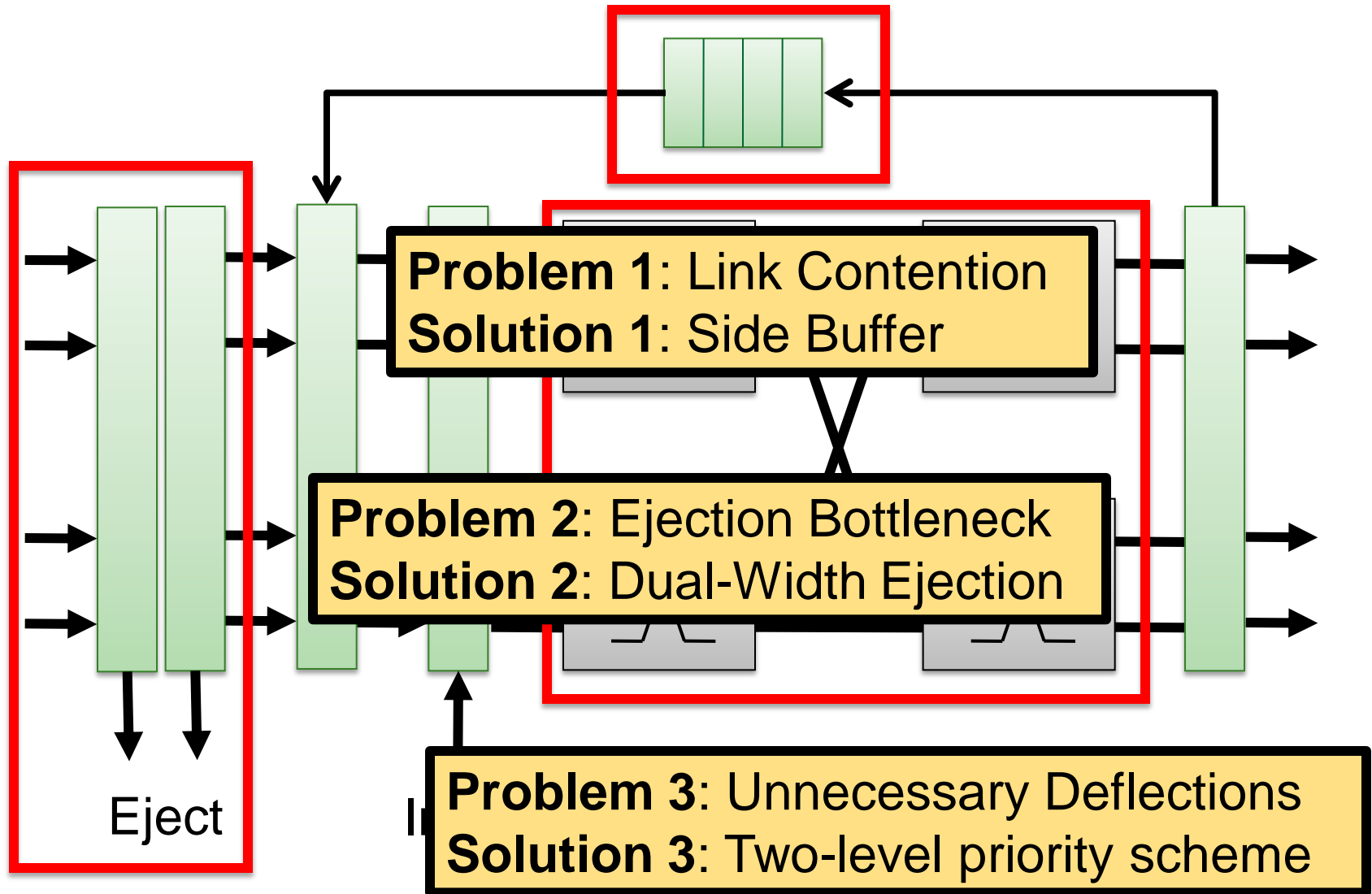
- **Next-highest priority:** one **silver flit** per router per cycle
 - Chosen pseudo-randomly & local to one router
 - **Enhances performance**

Adding A Silver Flit

- Randomly picking a silver flit ensures **one flit is not deflected**
 - No flit is golden but Red flit is silver
 - Red flit wins at first stage (silver)
 - Green flit is deflected at first stage
 - Red flit wins at second stage (silver); not deflected



Minimally-Buffered Deflection Router



Outline: This Talk

- **Motivation**
- **Background:** Bufferless Deflection Routing
- **MinBD:** Reducing Deflections
 - Addressing Link Contention
 - Addressing the Ejection Bottleneck
 - Improving Deflection Arbitration

Outline: This Talk

- **Motivation**
- **Background:** Bufferless Deflection Routing
- **MinBD:** Reducing Deflections
 - Addressing Link Contention
 - Addressing the Ejection Bottleneck
 - Improving Deflection Arbitration
- **Results**
- **Conclusions**

Methodology: Simulated System

■ **Chip Multiprocessor Simulation**

- ❑ **64-core** and **16-core** models
- ❑ **Closed-loop** core/cache/NoC cycle-level model
- ❑ Directory cache coherence protocol (SGI Origin-based)
- ❑ 64KB L1, perfect L2 (stresses interconnect), XOR-mapping
- ❑ Performance metric: **Weighted Speedup**
(similar conclusions from network-level latency)
- ❑ Workloads: multiprogrammed SPEC CPU2006
 - 75 randomly-chosen workloads
 - Binned into network-load categories by average injection rate

Methodology: Routers and Network

- **Input-buffered** virtual-channel router
 - 8 VCs, 8 flits/VC [**Buffered(8,8)**]: large buffered router
 - 4 VCs, 4 flits/VC [**Buffered(4,4)**]: typical buffered router
 - 4 VCs, 1 flit/VC [**Buffered(4,1)**]: smallest deadlock-free router
 - All power-of-2 buffer sizes up to (8, 8) for perf/power sweep
- **Bufferless deflection** router: **CHIPPER**¹
- **Bufferless-buffered hybrid** router: **AFC**²
 - Has input buffers and deflection routing logic
 - Performs coarse-grained (multi-cycle) mode switching
- **Common parameters**
 - 2-cycle router latency, 1-cycle link latency
 - 2D-mesh topology (16-node: 4x4; 64-node: 8x8)
 - Dual ejection assumed for baseline routers (for perf. only)

¹Fallin et al., “CHIPPER: A Low-complexity Bufferless Deflection Router”, HPCA 2011.

²Jafri et al., “Adaptive Flow Control for Robust Performance and Energy”, MICRO 2010.

Methodology: Power, Die Area, Crit. Path

■ **Hardware modeling**

- Verilog models for CHIPPER, MinBD, buffered control logic
 - Synthesized with commercial 65nm library
- ORION 2.0 for datapath: crossbar, muxes, buffers and links

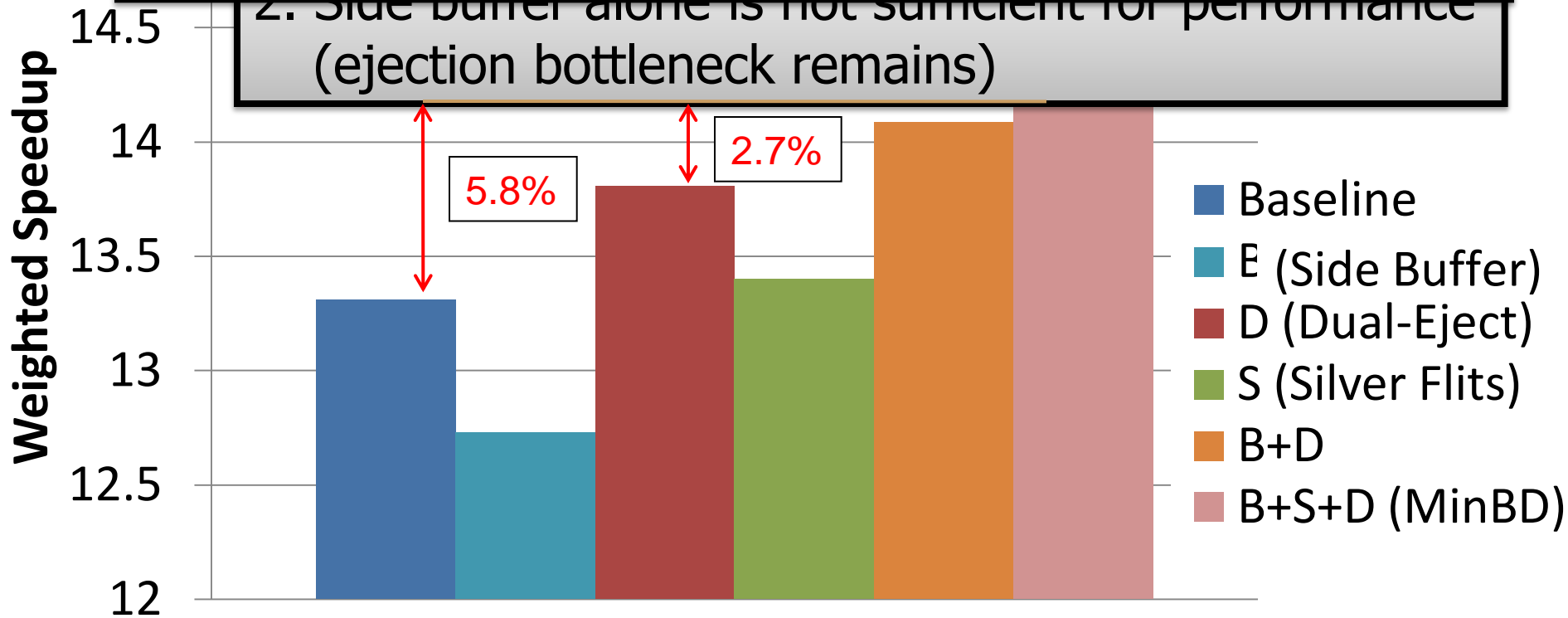
■ **Power**

- Static and dynamic power from hardware models
- Based on event counts in cycle-accurate simulations
- Broken down into buffer, link, other

Reduced Deflections & Improved Perf.

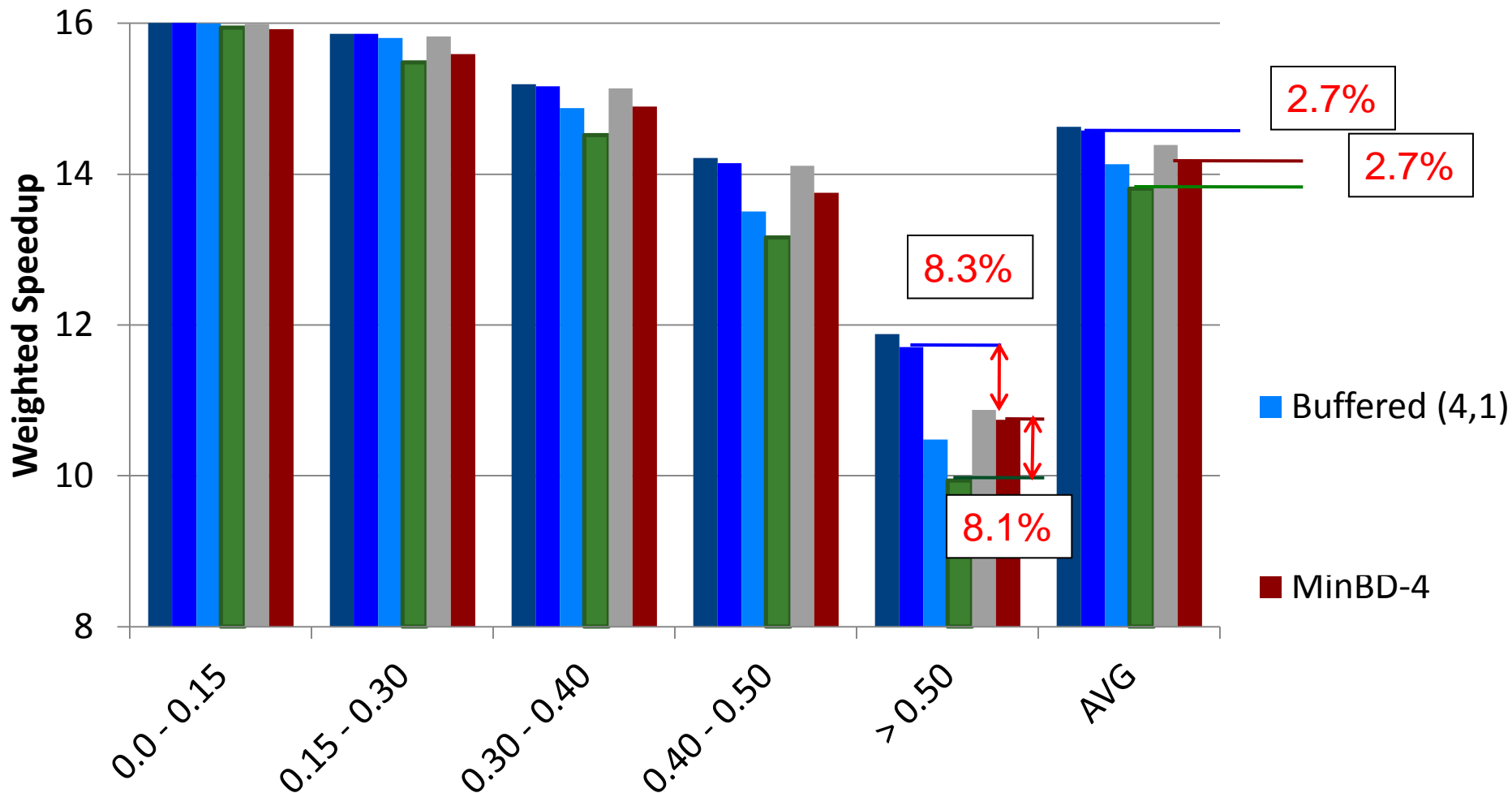
3. Overall, **5.8%** over baseline, **2.7%** over dual-eject by reducing deflections **64%** / **54%**

2. Side buffer alone is not sufficient for performance (ejection bottleneck remains)



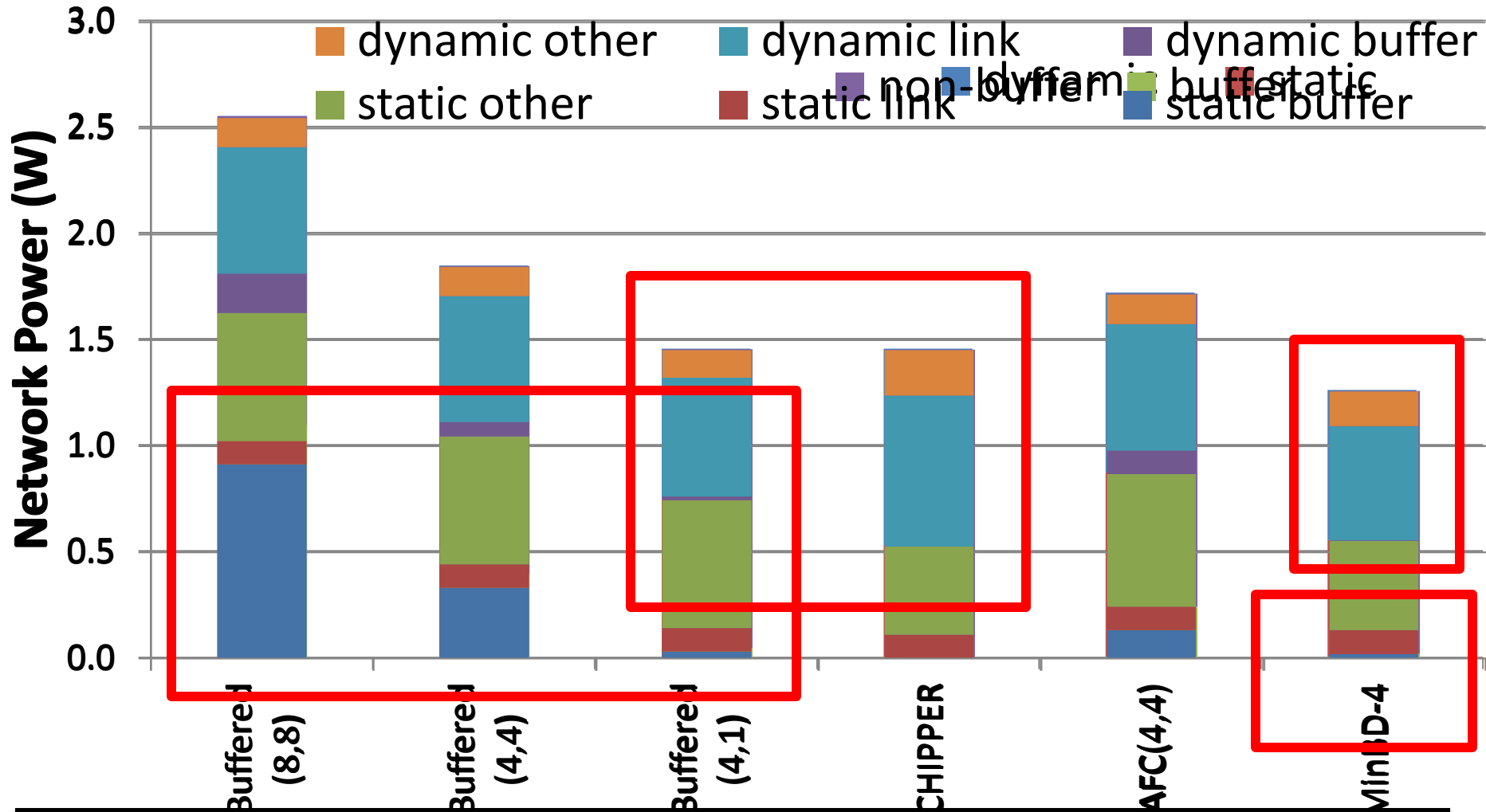
Deflection Rate	28%	17%	22%	27%	11%	10%
-----------------	-----	-----	-----	-----	-----	-----

Overall Performance Results



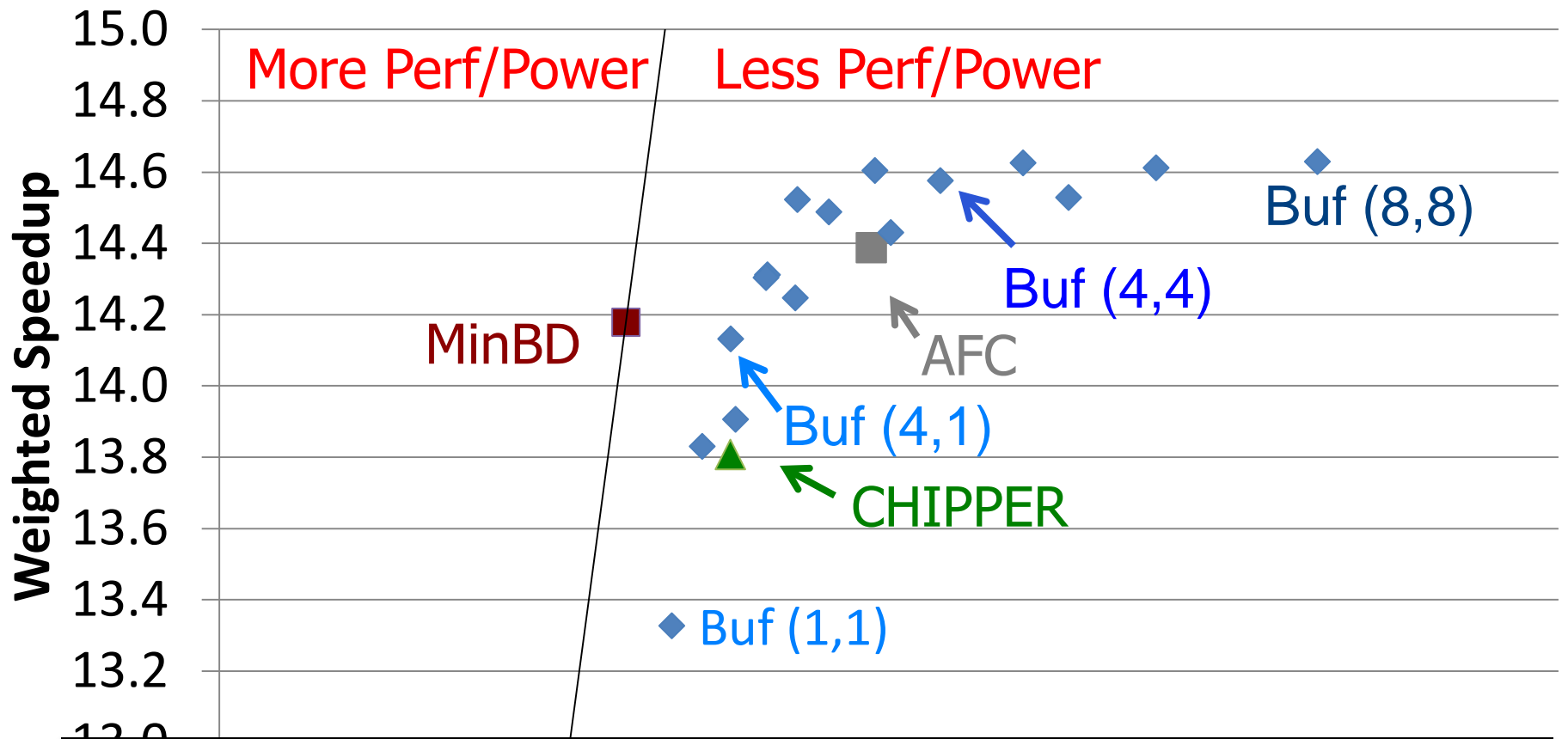
- Similar perf. to Buffered (4,1) @ 25% of buffering space
- Within **2.7%** of Buffered (4,4) (**8.3%** at high load)

Overall Power Results



- Dynamic power increases with deflection routing
- Dynamic power reduces in MinBD relative to CHIPPER

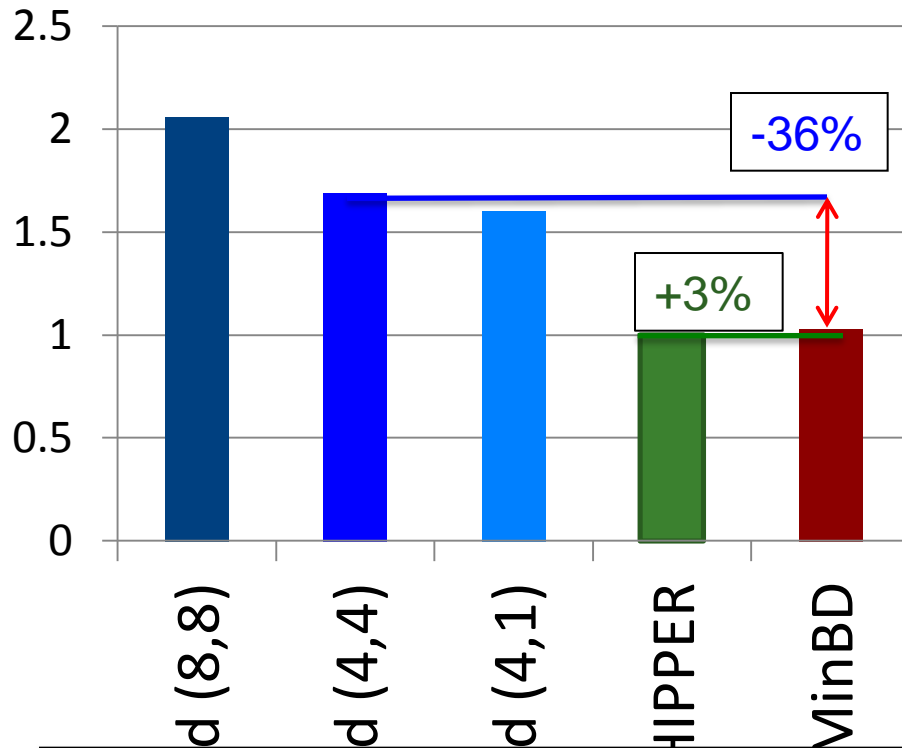
Performance-Power Spectrum



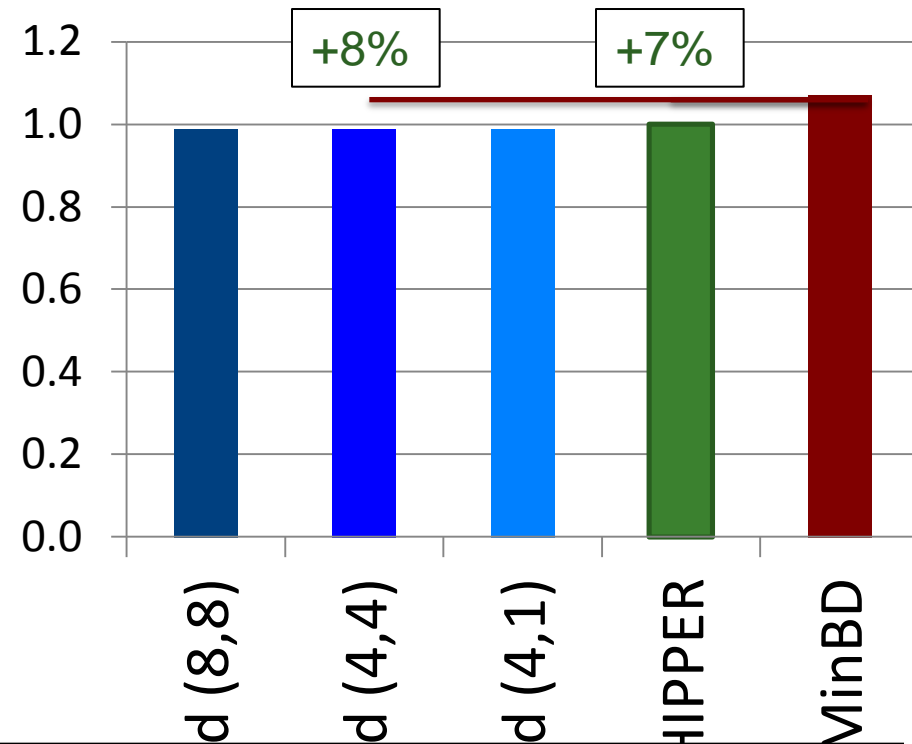
- Most **energy-efficient** (perf/watt) of any evaluated network router design

Die Area and Critical Path

Normalized Die Area



Normalized Critical Path



- Only **3%** area increase over CHIPPER (4-flit buffer)
- Increases by **7%** over CHIPPER, **8%** over Buffered (4,4)

Conclusions

- Bufferless deflection routing offers **reduced power & area**
- But, high deflection rate hurts **performance at high load**

- **MinBD** (Minimally-Buffered Deflection Router) introduces:
 - **Side buffer** to hold **only** flits that would have been deflected
 - **Dual-width ejection** to address ejection bottleneck
 - **Two-level prioritization** to avoid unnecessary deflections

- MinBD yields **reduced power (31%) & reduced area (36%)** relative to **buffered** routers
- MinBD yields **improved performance (8.1% at high load)** relative to **bufferless** routers → closes half of perf. gap

- MinBD has the **best energy efficiency** of all evaluated designs with **competitive performance**

More Readings

- **Studies of congestion and congestion control in on-chip vs. internet-like networks**
- George Nychis, Chris Fallin, Thomas Moscibroda, Onur Mutlu, and Srinivasan Seshan,
"On-Chip Networks from a Networking Perspective: Congestion and Scalability in Many-core Interconnects"
Proceedings of the 2012 ACM SIGCOMM Conference (SIGCOMM), Helsinki, Finland, August 2012. Slides (pptx)
- George Nychis, Chris Fallin, Thomas Moscibroda, and Onur Mutlu,
"Next Generation On-Chip Networks: What Kind of Congestion Control Do We Need?"
Proceedings of the 9th ACM Workshop on Hot Topics in Networks (HOTNETS), Monterey, CA, October 2010. Slides (ppt) (key)

HAT: Heterogeneous Adaptive Throttling for On-Chip Networks

Kevin Chang, Rachata Ausavarungnirun, Chris Fallin, and Onur Mutlu,

"HAT: Heterogeneous Adaptive Throttling for On-Chip Networks"

Proceedings of the 24th International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD), New York, NY, October 2012. Slides
(pptx) (pdf)

Carnegie Mellon University

SAFARI

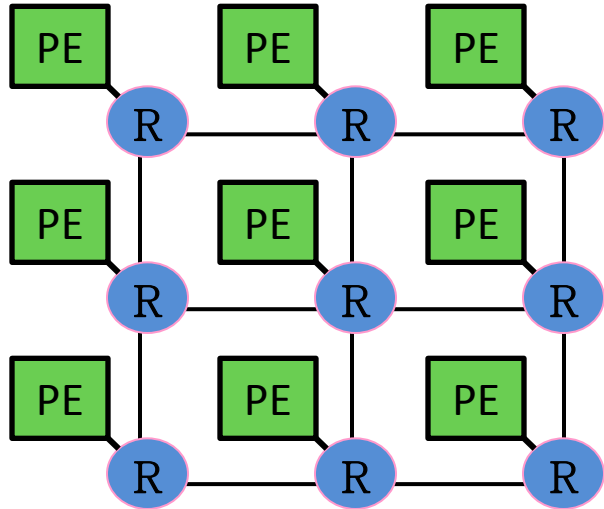
Executive Summary

- **Problem**: Packets contend in on-chip networks (NoCs), causing congestion, thus reducing performance
- **Observations**:
 - 1) Some applications are more sensitive to network latency than others
 - 2) Applications must be throttled differently to achieve peak performance
- **Key Idea: Heterogeneous Adaptive Throttling (HAT)**
 - 1) Application-aware source throttling
 - 2) Network-load-aware throttling rate adjustment
- **Result**: Improves performance and energy efficiency over state-of-the-art source throttling policies


Outline

- **Background and Motivation**
- Mechanism
- Prior Works
- Results

On-Chip Networks

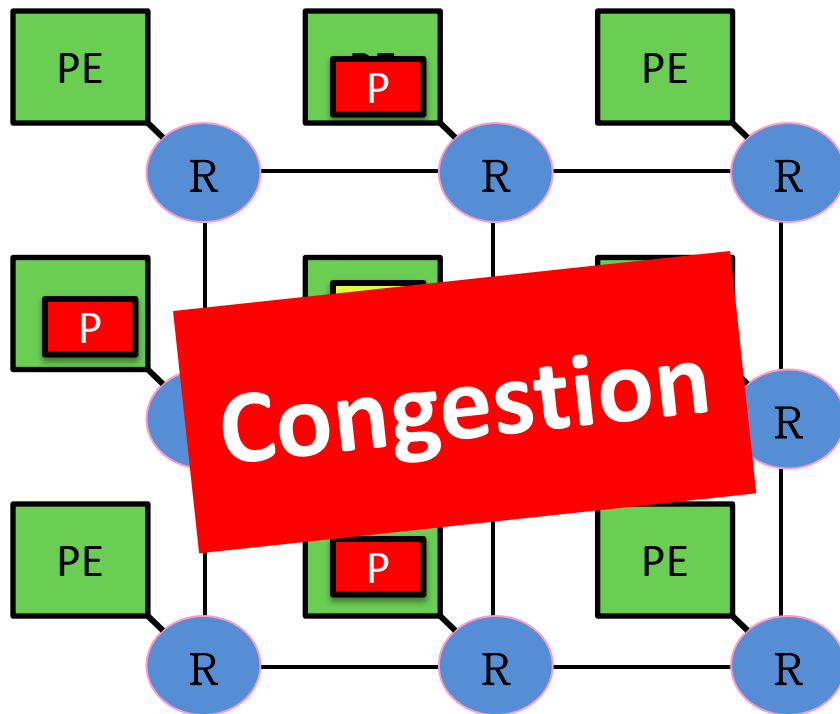


 Router

 Processing Element
(Cores, L2 Banks, Memory Controllers, etc)

- Connect **cores, caches, memory controllers, etc**
- **Packet switched**
- **2D mesh:** Most commonly used topology
- Primarily serve **cache misses** and **memory requests**
- **Router designs**
 - Buffered: **Input buffers** to hold contending packets
 - Bufferless: **Misroute (deflect)** contending packets

Network Congestion Reduces Performance



Limited shared resources
(buffers and links)

- **Design constraints: power, chip area, and timing**

Network congestion:

↓ Network throughput

↓ Application performance

Goal

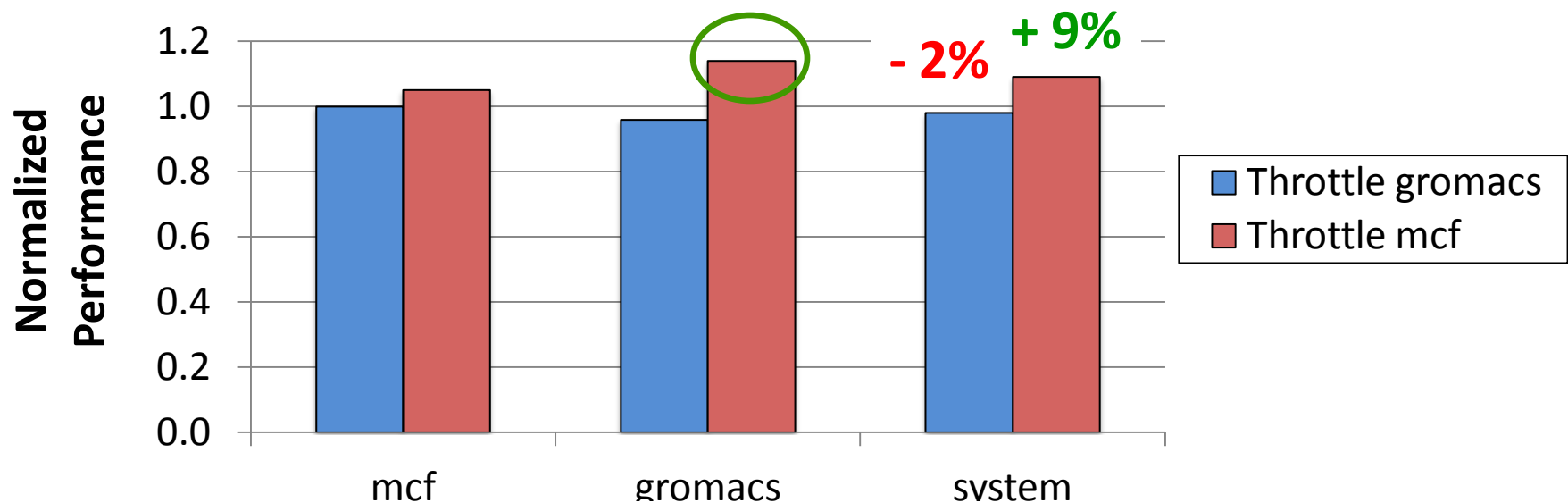
- **Improve performance in a highly congested NoC**
- Reducing network load decreases network congestion, hence improves performance
- **Approach: source throttling to reduce network load**
 - Temporarily delay new traffic injection
- **Naïve mechanism: throttle every single node**

Key Observation #1

Different applications respond differently to changes in **network latency**

gromacs: network-**non**-intensive

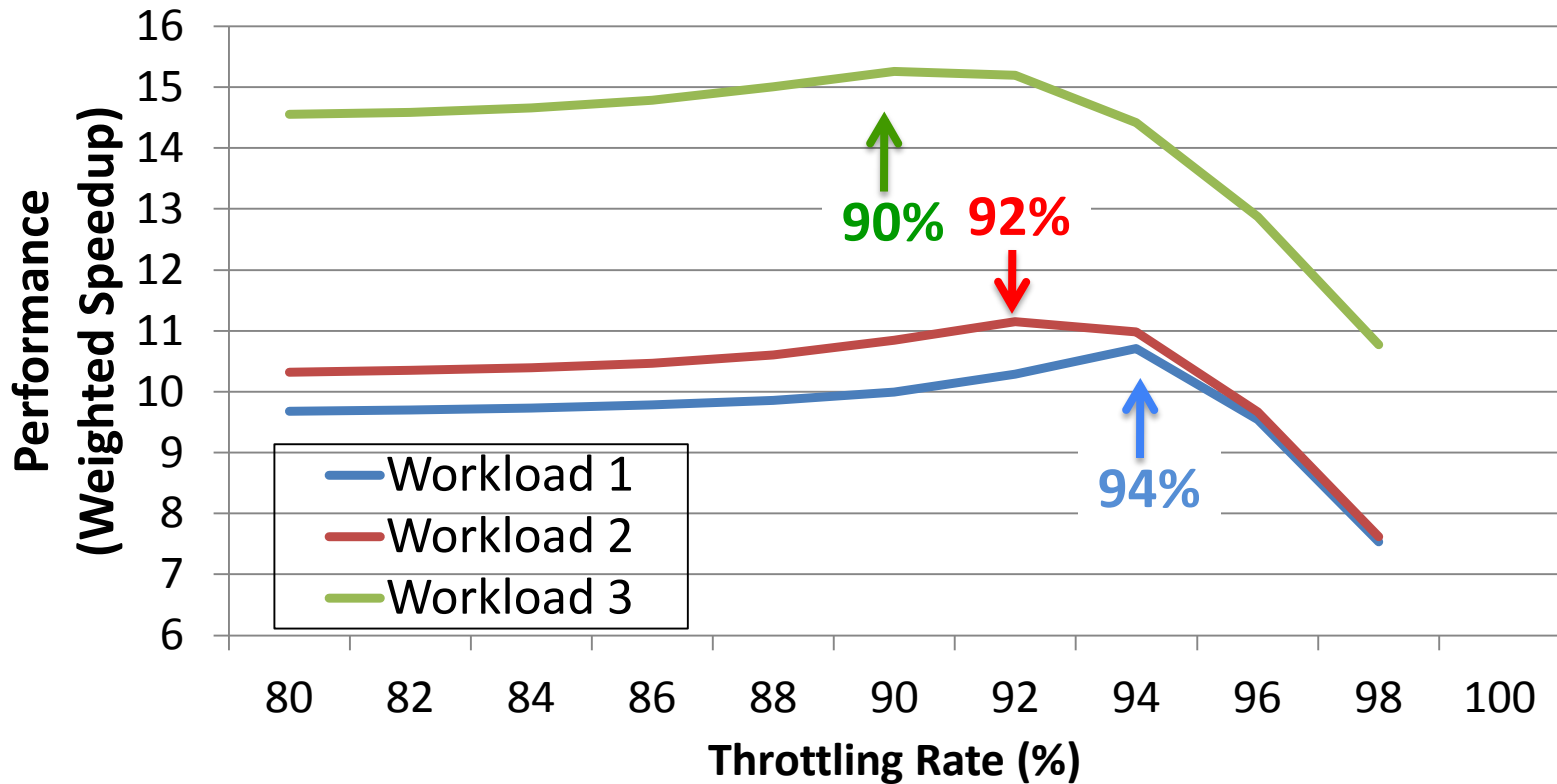
mcf: network-**intensive**



Throttling **network-intensive** applications benefits system performance more

Key Observation #2

Different workloads achieve peak performance at different throttling rates



Dynamically adjusting throttling rate yields better performance than a single static rate

Outline

- Background and Motivation
- **Mechanism**
- Prior Works
- Results

Heterogeneous Adaptive Throttling (HAT)

1. Application-aware throttling:

Throttle **network-intensive** applications that interfere with **network-non-intensive** applications

2. Network-load-aware throttling rate adjustment:

Dynamically adjusts throttling rate to adapt to different workloads

Heterogeneous Adaptive Throttling (HAT)

1. Application-aware throttling:

Throttle **network-intensive** applications that interfere with **network-non-intensive** applications

2. Network-load-aware throttling rate adjustment:

Dynamically adjusts throttling rate to adapt to different workloads

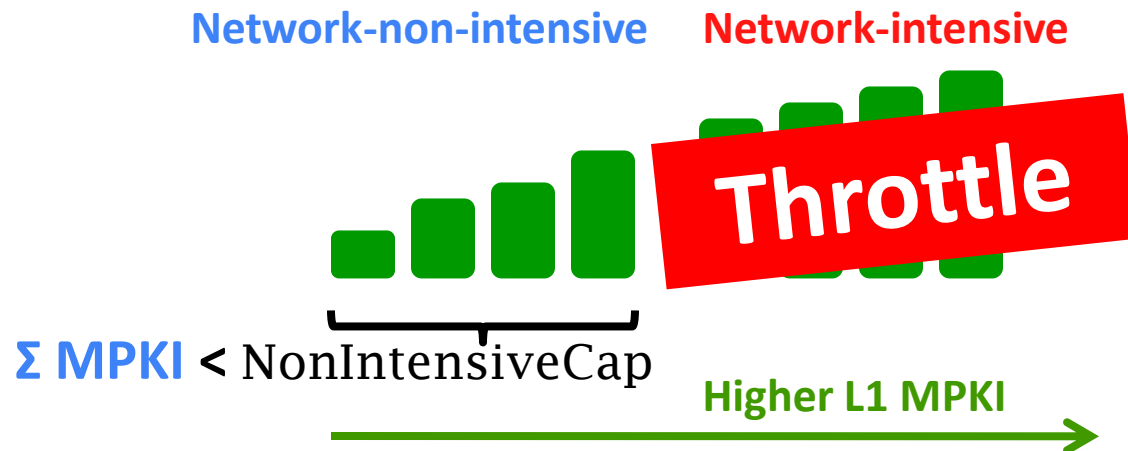
Application-Aware Throttling

1. Measure Network Intensity

Use **L1 MPKI** (misses per thousand instructions) to estimate network intensity

2. Classify Application

Sort applications by L1 MPKI



3. Throttle network-intensive applications

Heterogeneous Adaptive Throttling (HAT)

1. Application-aware throttling:

Throttle **network-intensive** applications that interfere with **network-non-intensive** applications

2. Network-load-aware throttling rate adjustment:

Dynamically adjusts throttling rate to adapt to different workloads

Dynamic Throttling Rate Adjustment

- For a given **network design**, peak performance tends to occur at a **fixed network load point**
- **Dynamically** adjust throttling rate to achieve that network load point

Dynamic Throttling Rate Adjustment

- **Goal:** maintain network load at a peak performance point

1. Measure network load

2. Compare and adjust throttling rate

If **network load** > **peak point**:

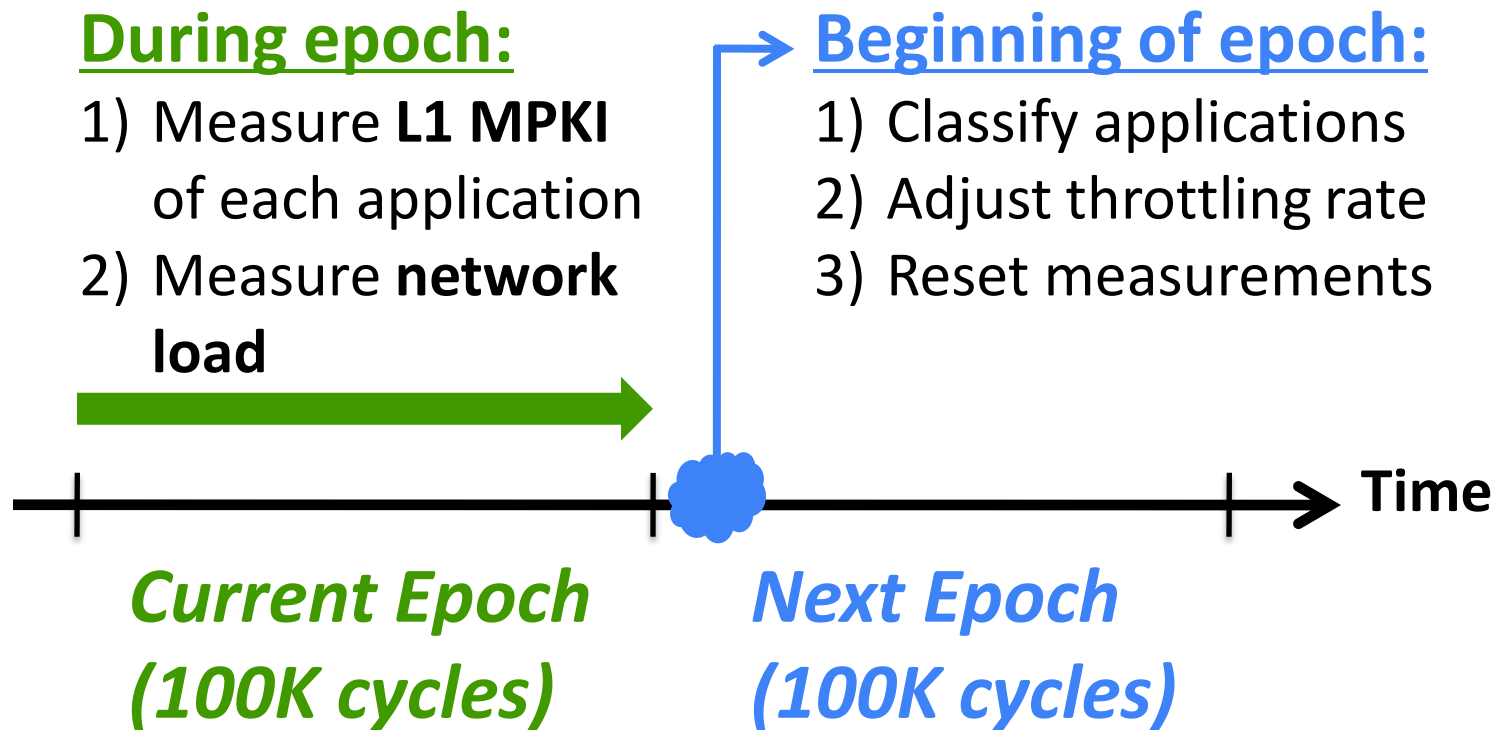
Increase throttling rate

elif **network load** ≤ **peak point**:

Decrease throttling rate

Epoch-Based Operation

- Continuous **HAT** operation is expensive
- **Solution:** performs **HAT** at epoch granularity



Outline

- Background and Motivation
- Mechanism
- **Prior Works**
- Results

Prior Source Throttling Works

- **Source throttling for bufferless NoCs**

[Nychis+ Hotnets'10, SIGCOMM'12]

- Application-aware throttling based on starvation rate
- **Does not adaptively adjust throttling rate**
- “Heterogeneous Throttling”

- **Source throttling off-chip buffered networks**

[Thottethodi+ HPCA'01]

- Dynamically trigger throttling based on fraction of buffer occupancy
- **Not application-aware: fully block packet injections of every node**
- “Self-tuned Throttling”

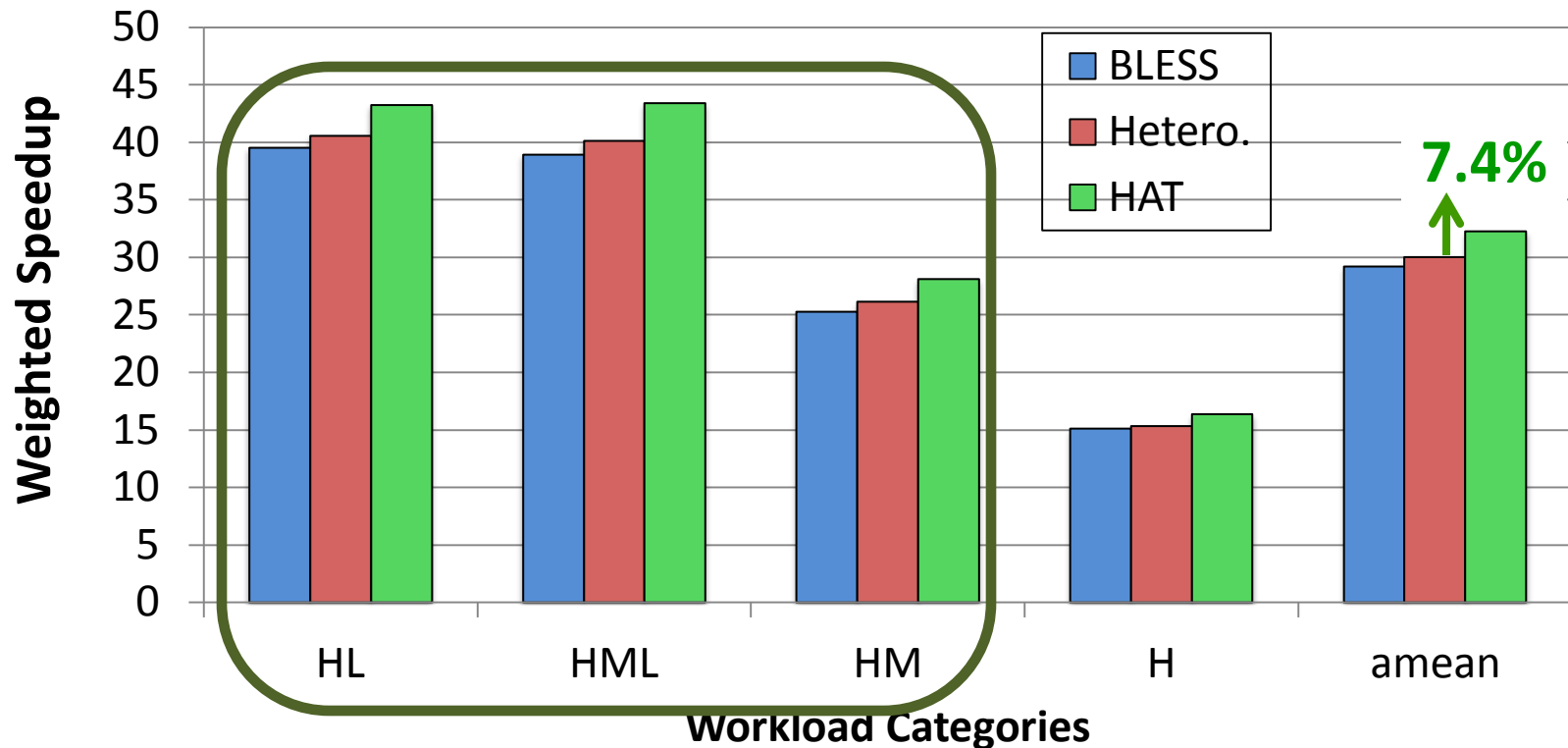
Outline

- Background and Motivation
- Mechanism
- Prior Works
- **Results**

Methodology

- **Chip Multiprocessor Simulator**
 - **64-node** multi-core systems with a **2D-mesh topology**
 - Closed-loop core/cache/NoC cycle-level model
 - 64KB L1, perfect L2 (always hits to stress NoC)
- **Router Designs**
 - **Virtual-channel buffered** router: 4 VCs, 4 flits/VC [Dally+ IEEE TPDS'92]
 - **Bufferless deflection** routers: **BLESS** [Moscibroda+ ISCA'09]
- **Workloads**
 - 60 multi-core workloads: SPEC CPU2006 benchmarks
 - Categorized based on their network intensity
 - Low/**M**edium/**H**igh intensity categories
- **Metrics:** Weighted Speedup (perf.), perf./Watt (energy eff.), and maximum slowdown (fairness)

Performance: Bufferless NoC (BLESS)

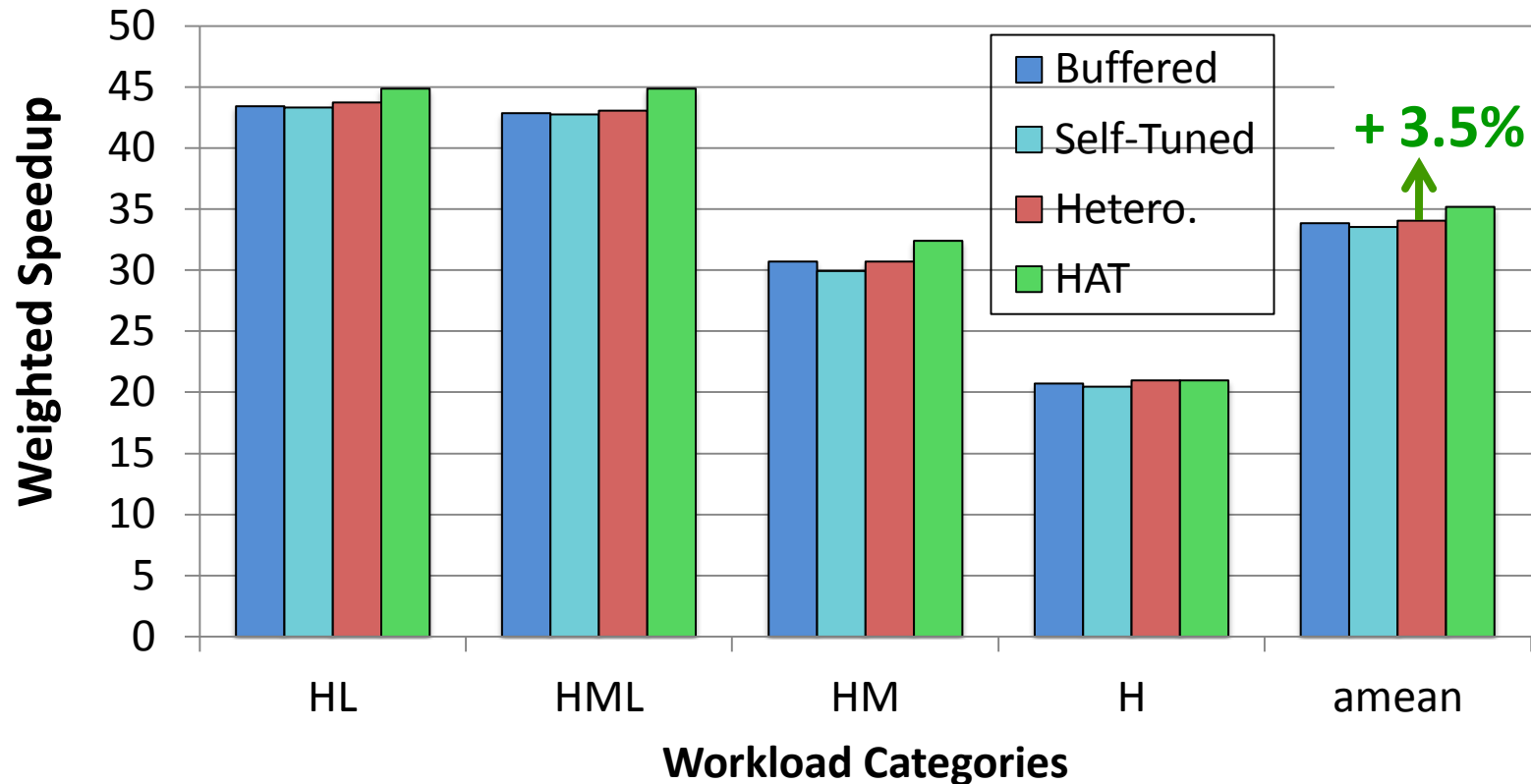


HAT provides better performance improvement than past work

Highest improvement on **heterogeneous** workload mixes

- **L** and **M** are more **sensitive** to network latency

Performance: Buffered NoC



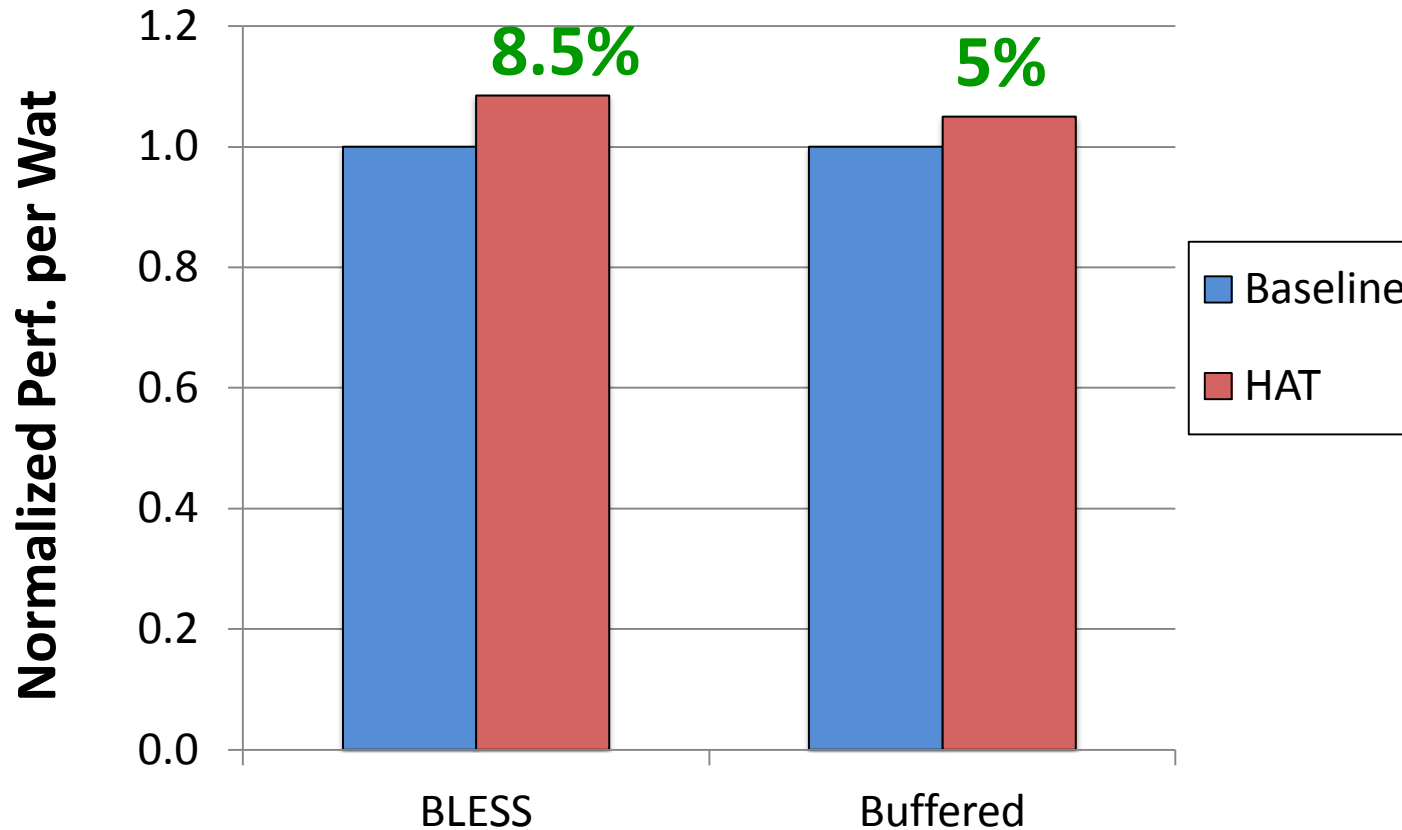
Congestion is much lower in Buffered NoC, but **HAT** still provides performance benefit

Application Fairness



HAT provides better fairness than prior works

Network Energy Efficiency



HAT increases energy efficiency by reducing congestion

Other Results in Paper

- Performance on **CHIPPER**
- Performance on **multithreaded** workloads
- Parameters sensitivity sweep of **HAT**

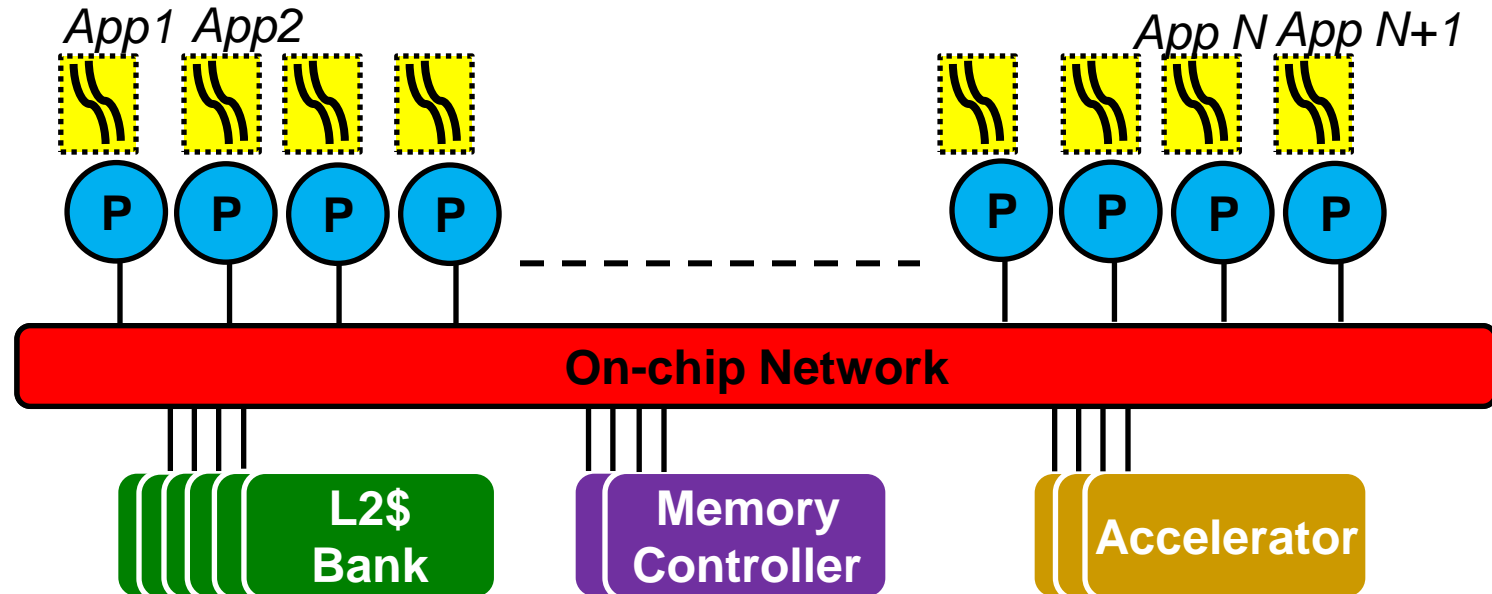
Conclusion

- **Problem**: Packets contend in on-chip networks (NoCs), causing congestion, thus reducing performance
- **Observations**:
 - 1) Some applications are more sensitive to network latency than others
 - 2) Applications must be throttled differently to achieve peak performance
- **Key Idea**: **Heterogeneous Adaptive Throttling (HAT)**
 - 1) Application-aware source throttling
 - 2) Network-load-aware throttling rate adjustment
- **Result**: Improves performance and energy efficiency over state-of-the-art source throttling policies

Application-Aware Packet Scheduling

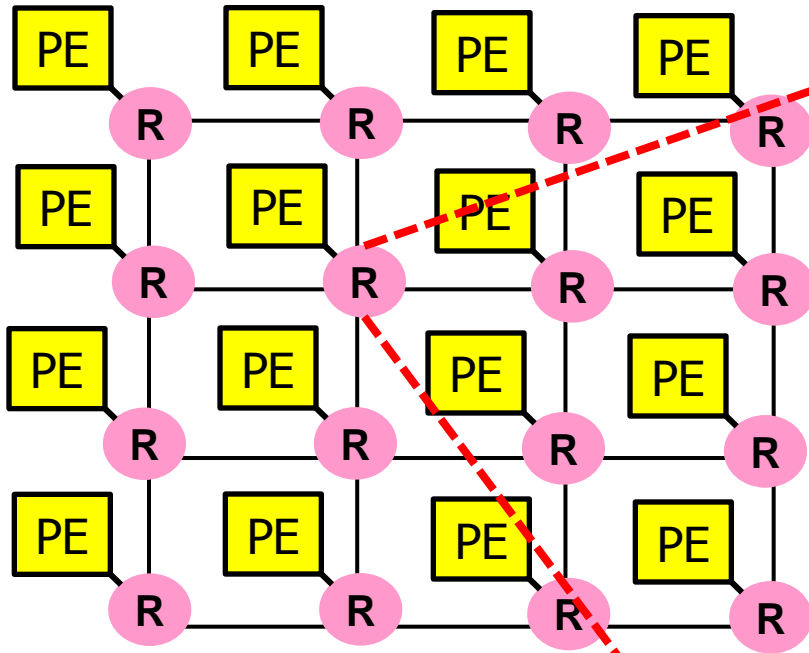
Reetuparna Das, Onur Mutlu, Thomas Moscibroda, and Chita R. Das,
"Application-Aware Prioritization Mechanisms for On-Chip Networks"
Proceedings of the 42nd International Symposium on Microarchitecture
(MICRO), pages 280-291, New York, NY, December 2009. Slides (pptx)

The Problem: Packet Scheduling

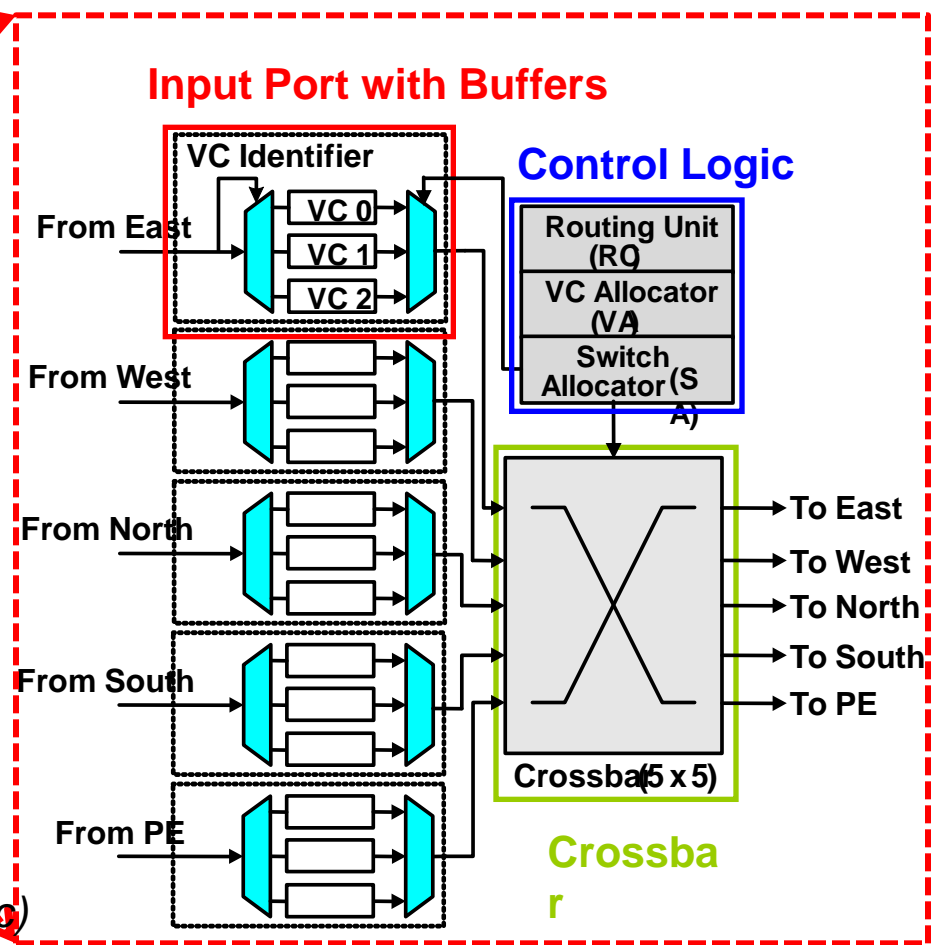


On-chip Network is a critical resource shared by multiple applications

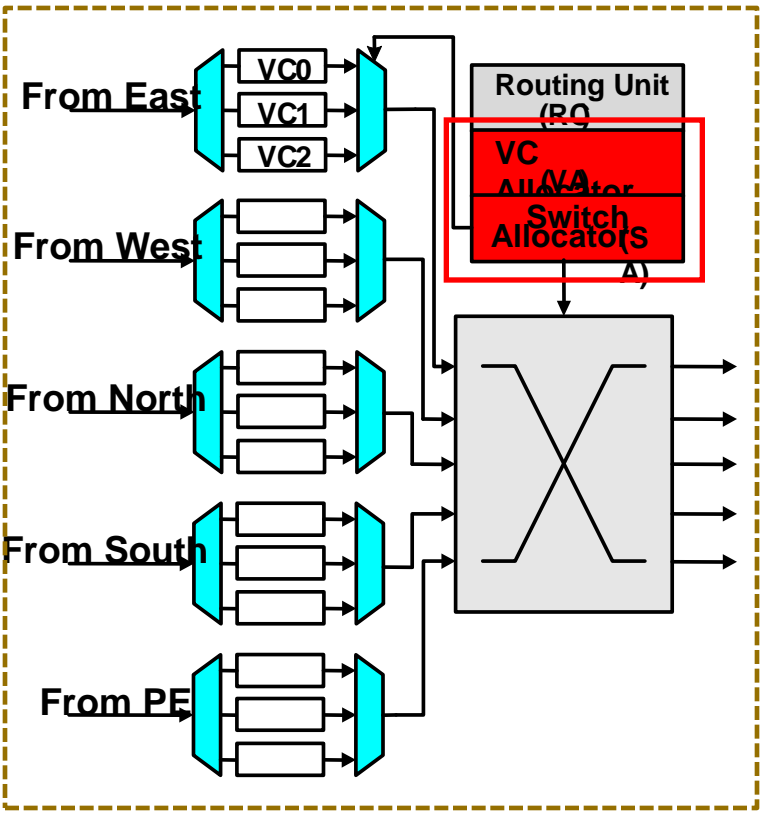
The Problem: Packet Scheduling



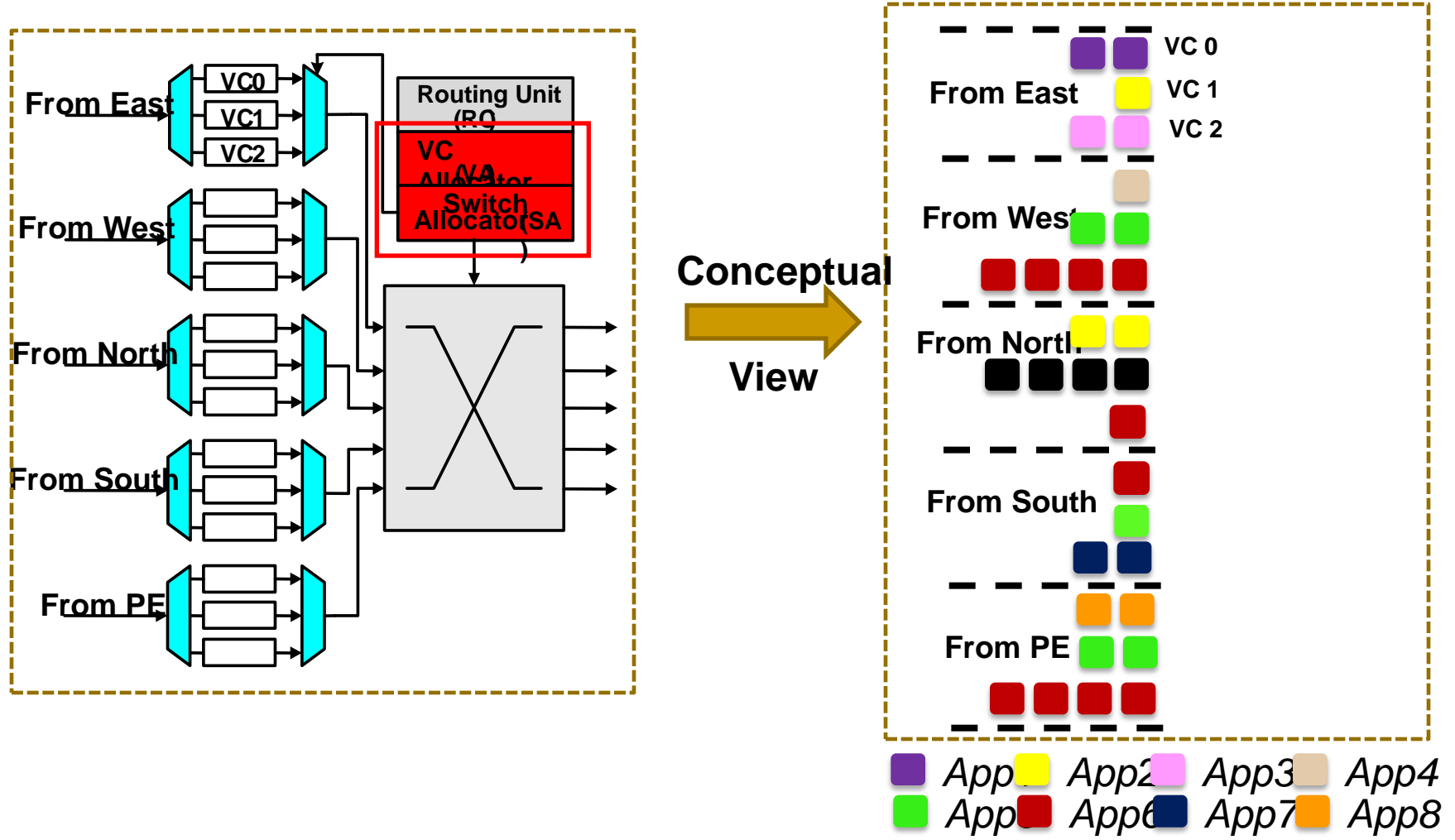
R Routers
PE Processing Element
(Cores, L2 Banks, Memory Controllers etc)



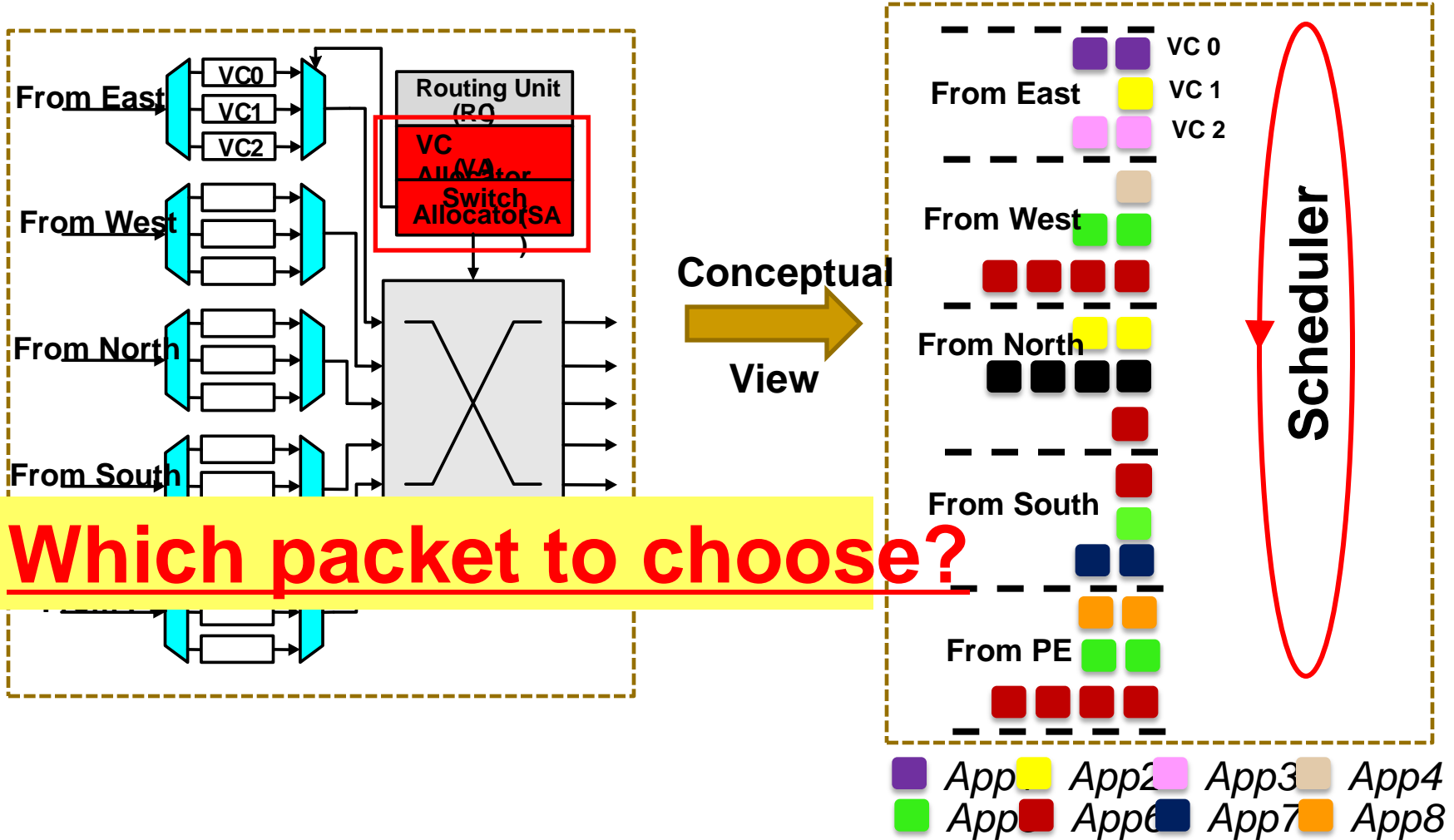
The Problem: Packet Scheduling



The Problem: Packet Scheduling



The Problem: Packet Scheduling



The Problem: Packet Scheduling

- Existing scheduling policies
 - Round Robin
 - Age
- Problem 1: **Local** to a router
 - Lead to contradictory decision making between routers: packets from one application may be prioritized at one router, to be delayed at next.
- Problem 2: **Application oblivious**
 - Treat all applications **packets equally**
 - But applications are heterogeneous
- **Solution**: Application-aware global scheduling policies.

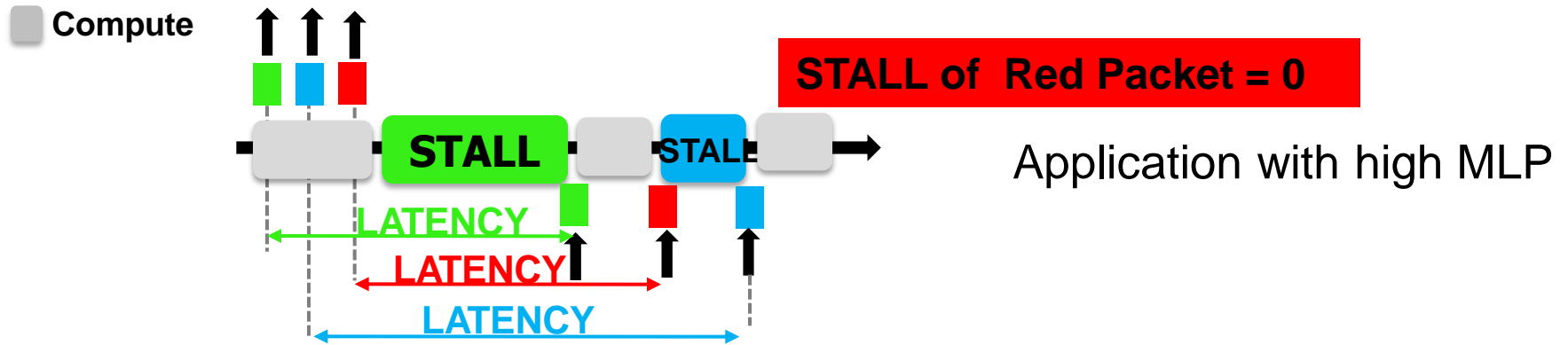
Motivation: Stall-Time Criticality

- Applications are **not homogenous**
- Applications have different **criticality** with respect to the **network**
 - Some applications are network latency sensitive
 - Some applications are network latency tolerant
- Application's **Stall Time Criticality (STC)** can be measured by its average network stall time per packet (**i.e. NST/packet**)
 - **Network Stall Time (NST)** is number of cycles the processor stalls waiting for network transactions to complete

Motivation: Stall-Time Criticality

- Why do applications have different network stall time criticality (STC)?
 - Memory Level Parallelism (MLP)
 - **Lower MLP leads to higher criticality**
 - Shortest Job First Principle (SJF)
 - **Lower network load leads to higher criticality**

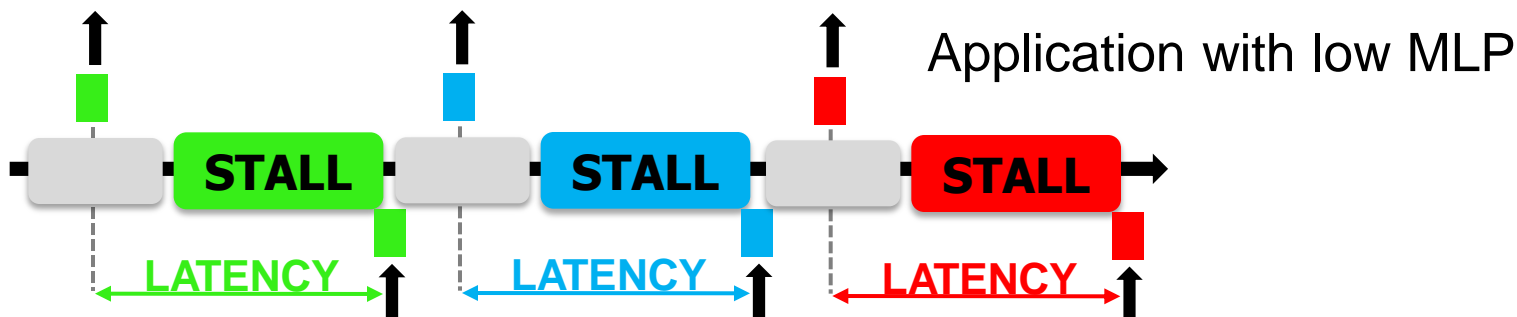
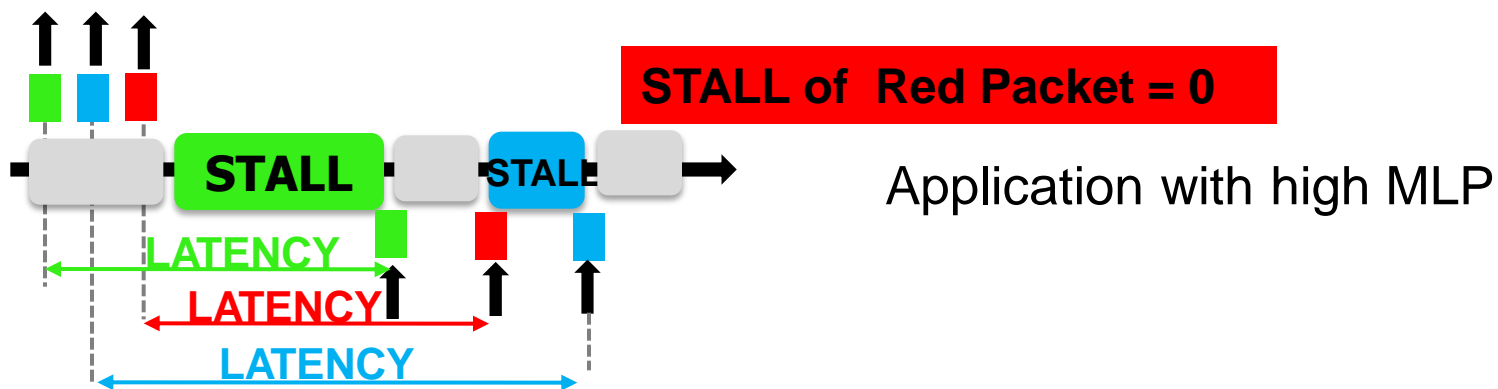
STC Principle 1: MLP



- Observation 1: **Packet Latency \neq Network Stall Time**

STC Principle 1: MLP

■ Compute



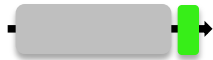
- Observation 1: **Packet Latency \neq Network Stall Time**
- Observation 2: A low MLP application's packets have higher criticality than a high MLP application's

STC Principle 2: Shortest-Job-First

Light Application

Heavy Application

Running ALONE



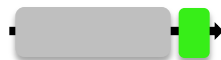
Baseline (RR) Scheduling



4X network slow down

1.3X network slow down

SJF Scheduling



1.2X network slow down

1.6X network slow down

Overall system throughput (weighted speedup) increases by 34%

Solution: Application-Aware Policies

- Idea
 - Identify critical applications (i.e. network sensitive applications) and prioritize their packets in each router.
- Key components of scheduling policy:
 - Application Ranking
 - Packet Batching
- Propose low-hardware complexity solution

Component 1: Ranking

- Ranking distinguishes applications based on Stall Time Criticality (STC)
- Periodically **rank** applications based on STC
- Explored many **heuristics** for estimating STC
 - Heuristic based on **outermost private cache Misses Per Instruction (L1-MPI)** is the most effective
 - **Low L1-MPI => high STC => higher rank**
- Why Misses Per Instruction (L1-MPI)?
 - Easy to Compute (low complexity)
 - Stable Metric (unaffected by interference in network)

Component 1 : How to Rank?

- Execution time is divided into fixed “ranking intervals”
 - Ranking interval is 350,000 cycles
- At the end of an interval, each core calculates their L1-MPI and sends it to the **Central Decision Logic (CDL)**
 - CDL is located in the central node of mesh
- CDL forms a rank order and sends back its rank to each core
 - Two control packets per core every ranking interval
- Ranking order is a “partial order”

- Rank formation is **not** on the **critical path**
 - Ranking interval is significantly longer than rank computation time
 - Cores use older rank values until new ranking is available

Component 2: Batching

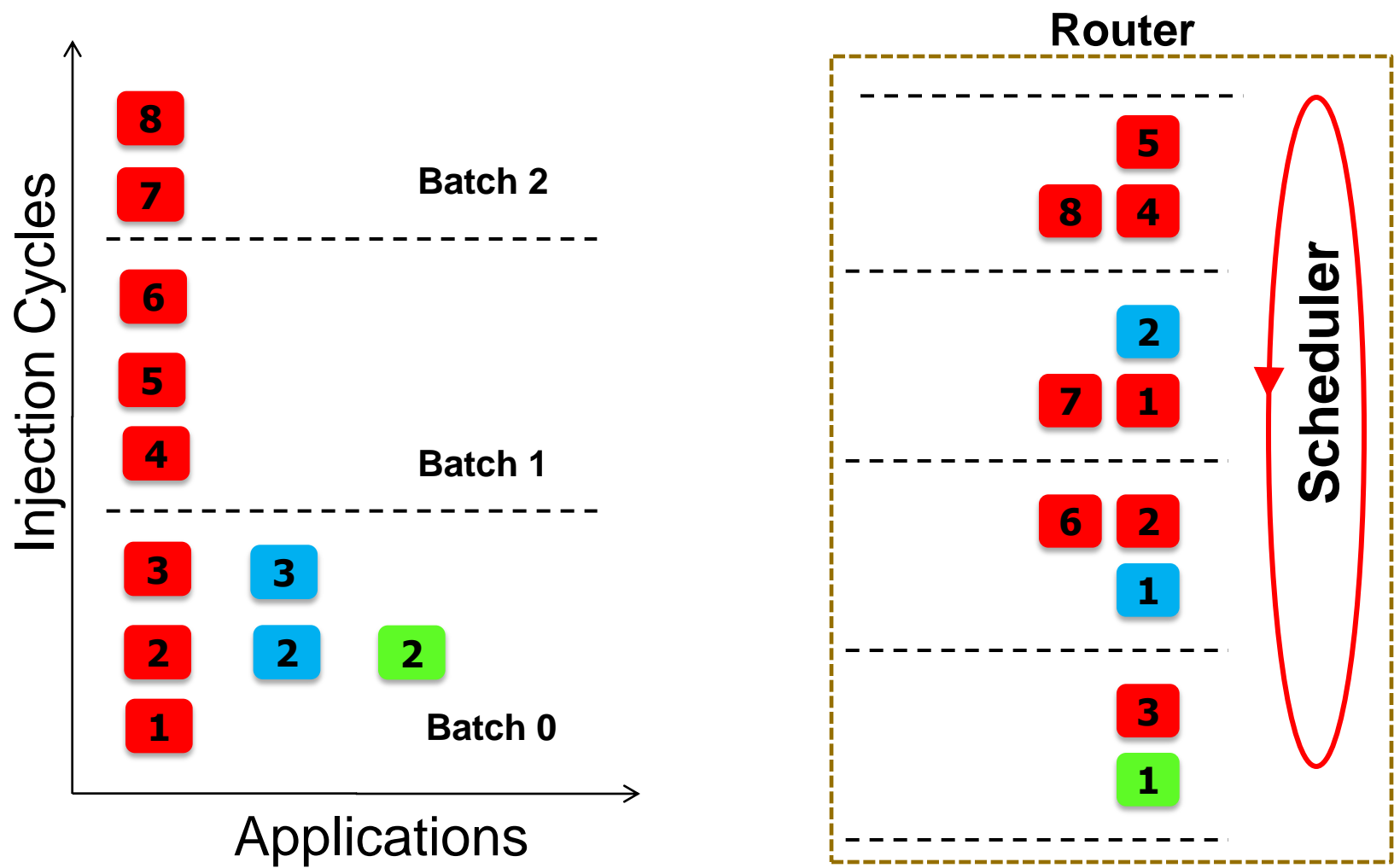
- Problem: **Starvation**
 - Prioritizing a higher ranked application can lead to starvation of lower ranked application
- Solution: **Packet Batching**
 - Network packets are grouped into finite sized batches
 - **Packets of older batches are prioritized over younger batches**
- **Time-Based Batching**
 - New batches are formed in a periodic, synchronous manner across all nodes in the network, every T cycles

Putting it all together: STC Scheduling

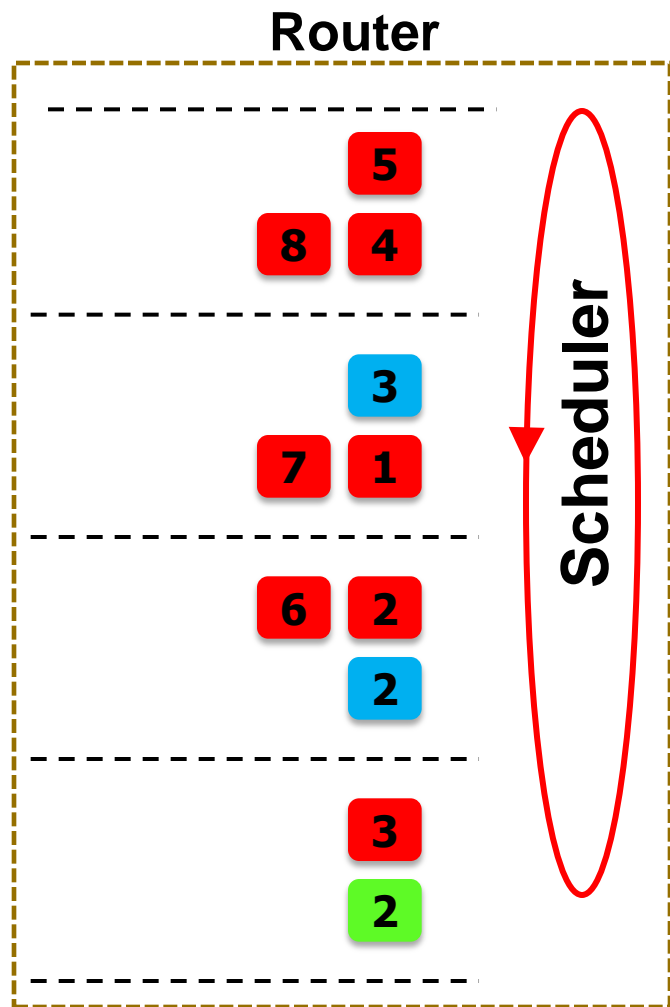
Policy

- Before injecting a packet into the network, it is tagged with
 - Batch ID (*3 bits*)
 - Rank ID (*3 bits*)
- Three tier priority structure at routers
 - **Oldest batch first** (*prevent starvation*)
 - **Highest rank first** (*maximize performance*)
 - **Local Round-Robin** (*final tie breaker*)
- Simple hardware support: priority arbiters
- **Global coordinated scheduling**
 - Ranking order and batching order are same across all routers

STC Scheduling Example

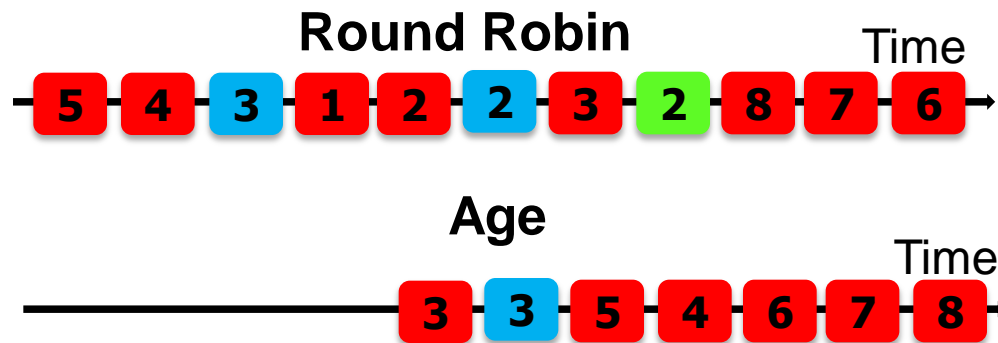
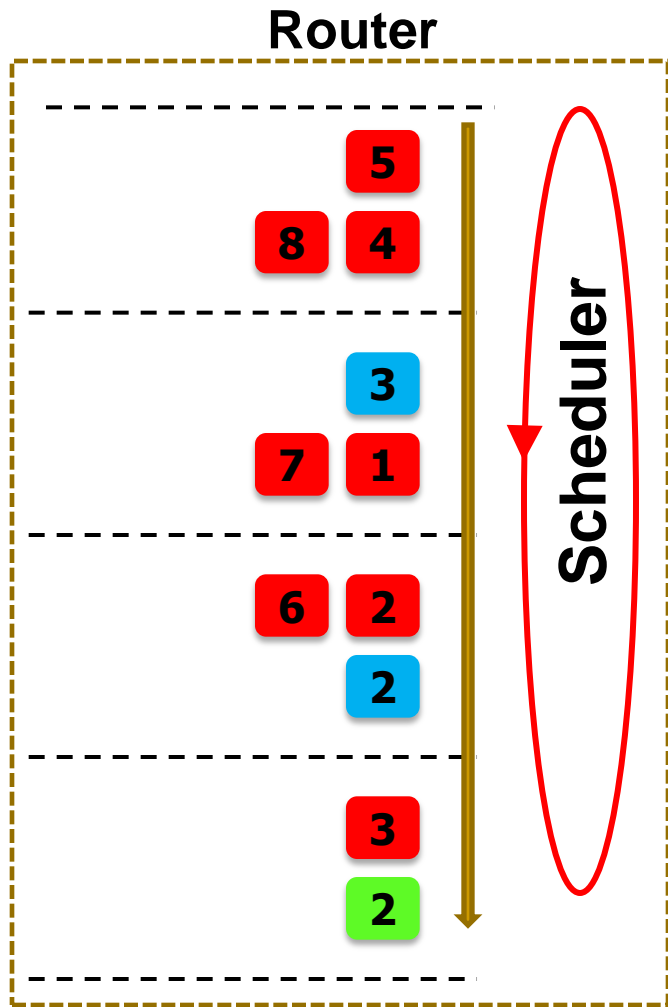


STC Scheduling Example



	STALL CYCLES			Avg
RR	8	6	11	8.3
Age				
STC				

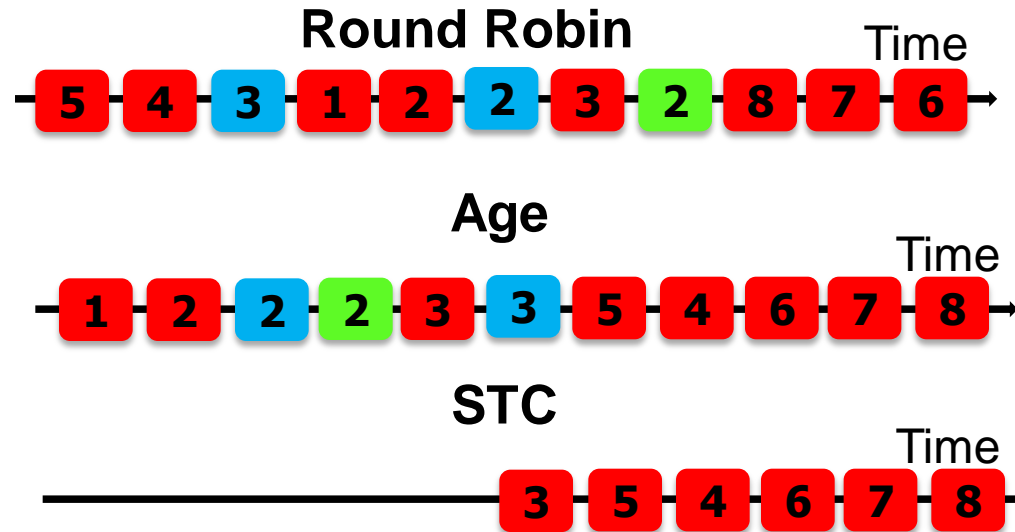
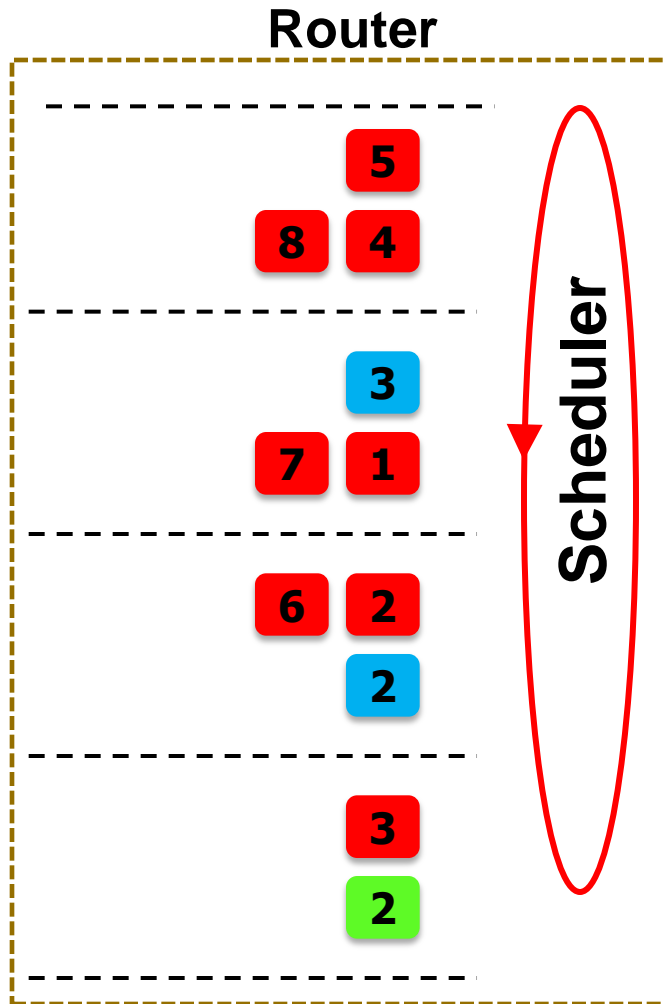
STC Scheduling Example



	STALL CYCLES			Avg
RR	8	6	11	8.3
Age	4	6	11	7.0
STC				

STC Scheduling Example

Rank order ■ > ■ > ■



STALL CYCLES				Avg
RR	8	6	11	8.3
Age	4	6	11	7.0
STC	1	3	11	5.0

STC Evaluation Methodology

- 64-core system
 - x86 processor model based on Intel Pentium M
 - 2 GHz processor, 128-entry instruction window
 - 32KB private L1 and 1MB per core shared L2 caches, 32 miss buffers
 - 4GB DRAM, 320 cycle access latency, 4 on-chip DRAM controllers
- Detailed Network-on-Chip model
 - 2-stage routers (with speculation and look ahead routing)
 - Wormhole switching (8 flit data packets)
 - Virtual channel flow control (6 VCs, 5 flit buffer depth)
 - 8x8 Mesh (128 bit bi-directional channels)
- Benchmarks
 - Multiprogrammed scientific, server, desktop workloads (**35 applications**)
 - **96 workload combinations**

Comparison to Previous Policies

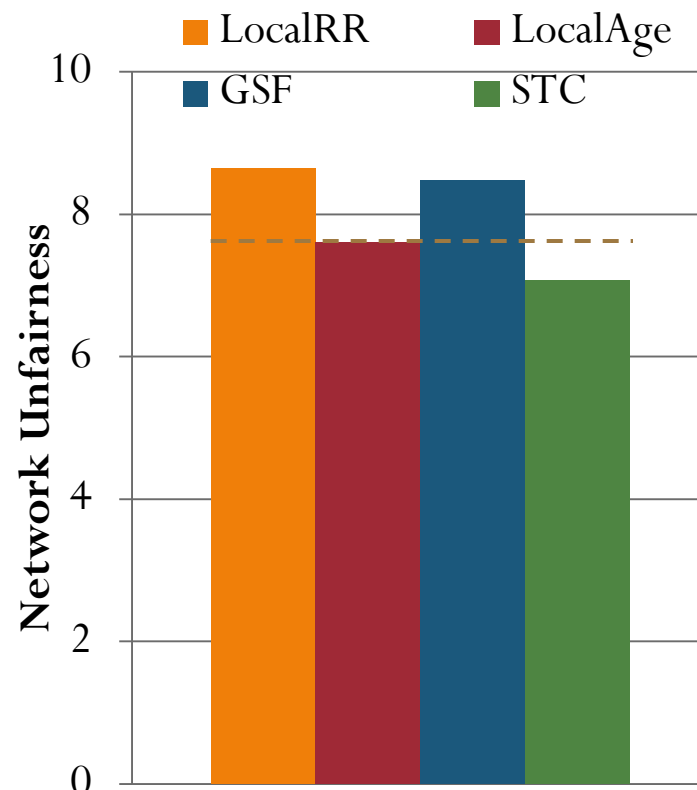
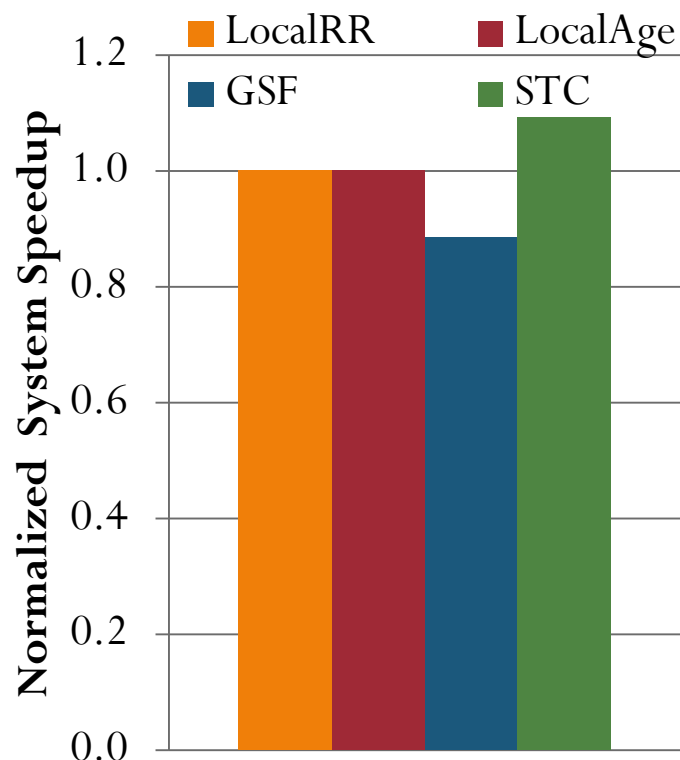
■ **Round Robin & Age (Oldest-First)**

- Local and application oblivious
- Age is biased towards heavy applications
 - heavy applications flood the network
 - higher likelihood of an older packet being from heavy application

■ **Globally Synchronized Frames (GSF)** [Lee et al., ISCA 2008]

- Provides **bandwidth fairness** at the expense of **system performance**
- Penalizes heavy and bursty applications
 - Each application gets equal and fixed quota of flits (credits) in each batch.
 - Heavy application quickly run out of credits after injecting into all active batches & stalls until oldest batch completes and frees up fresh credits.
 - Underutilization of network resources

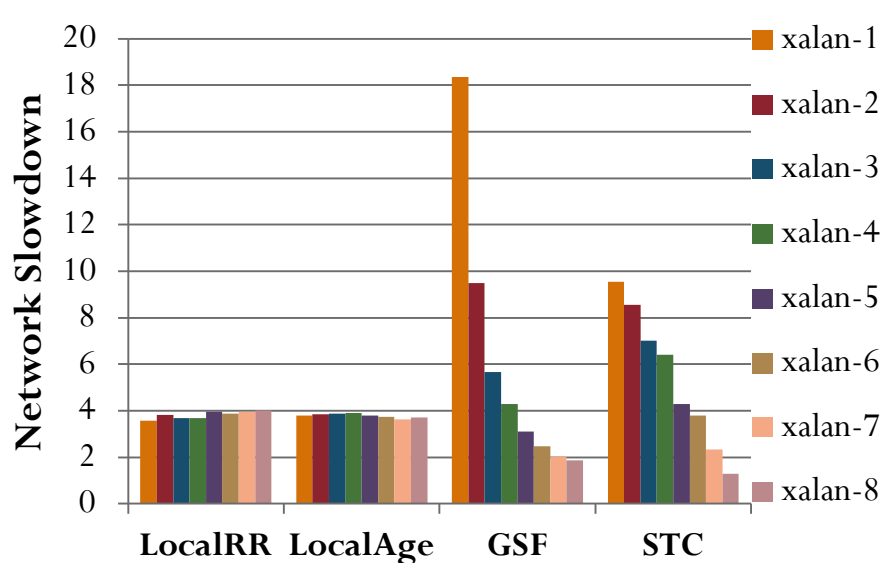
STC System Performance and Fairness



- 9.1% improvement in weighted speedup over the best existing policy (averaged across 96 workloads)

Enforcing Operating System Priorities

- Existing policies cannot enforce operating system (OS) assigned priorities in Network-on-Chip
- Proposed framework can enforce OS assigned priorities
 - Weight of applications => Ranking of applications
 - Configurable batching interval based on application weight



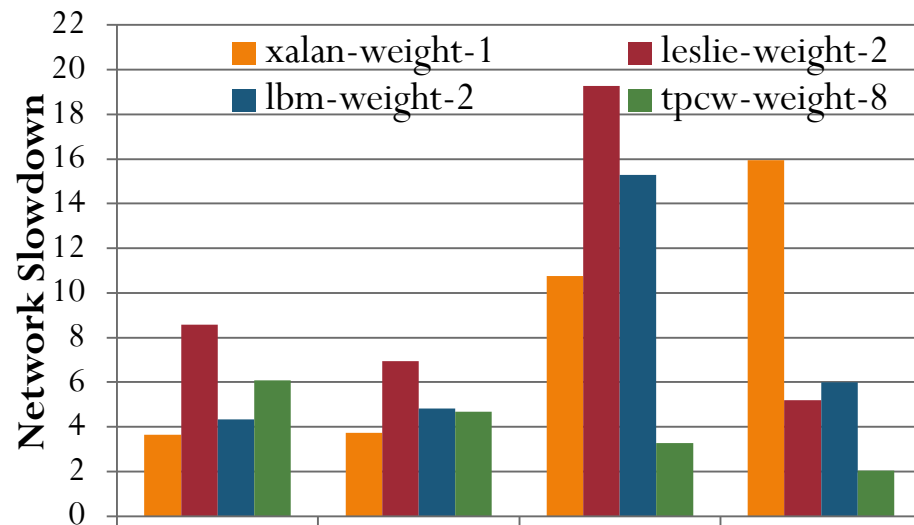
W. Speedup

0.49

0.49

0.46

0.52



W. Speedup

0.46

0.44

0.27

0.43

Application Aware Packet Scheduling: Summary

- Packet scheduling policies critically impact performance and fairness of NoCs
- Existing packet scheduling policies are **local** and **application oblivious**
- STC is a new, **global, application-aware approach to packet scheduling** in NoCs
 - **Ranking:** differentiates applications based on their criticality
 - **Batching:** avoids starvation due to rank-based prioritization
- Proposed framework
 - provides **higher system performance and fairness** than existing policies
 - can **enforce OS assigned priorities** in network-on-chip

Slack-Driven Packet Scheduling

Reetuparna Das, Onur Mutlu, Thomas Moscibroda, and Chita R. Das,
"Aergia: Exploiting Packet Latency Slack in On-Chip Networks"
Proceedings of the 37th International Symposium on Computer Architecture
(ISCA), pages 106-116, Saint-Malo, France, June 2010. Slides (pptx)

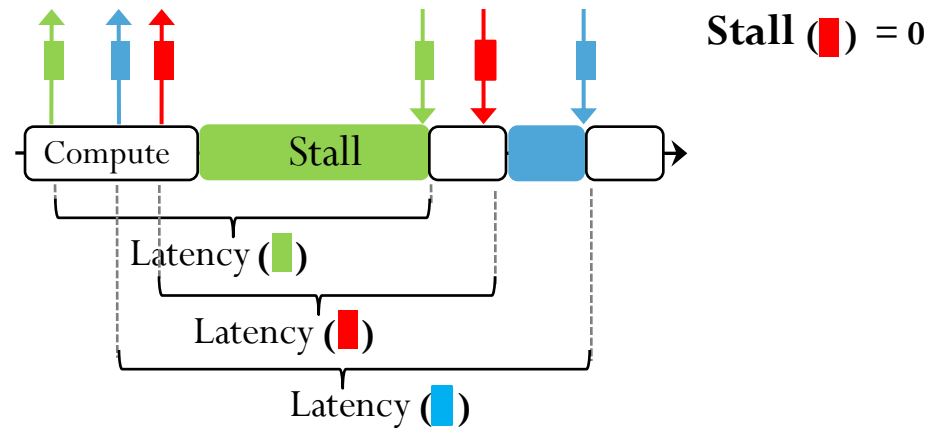
Packet Scheduling in NoC

- Existing scheduling policies
 - Round robin
 - Age
- Problem
 - Treat all packets equally
 - Application-oblivious
- Packets have **different criticality**
 - Packet is critical if latency of a packet affects application's performance
 - Different criticality due to memory level parallelism (MLP)

All packets are not the same...!!!



MLP Principle



Packet Latency \neq Network Stall Time

Different Packets have different criticality due to MLP

Criticality(green) > Criticality(blue) > Criticality(red)

Outline

- Introduction
 - Packet Scheduling
 - Memory Level Parallelism
- Aergia
 - Concept of Slack
 - Estimating Slack
- Evaluation
- Conclusion

What is Aérgia?



- Aérgia is the spirit of laziness in Greek mythology
- Some packets can afford to **slack!**

Outline

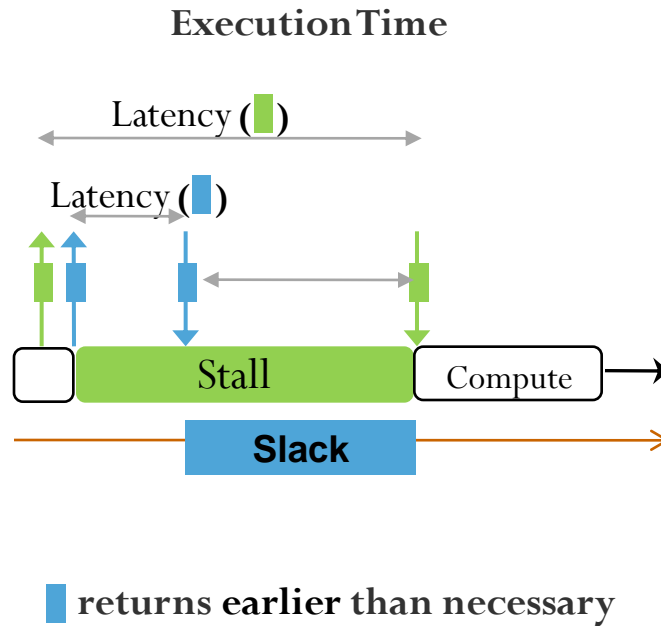
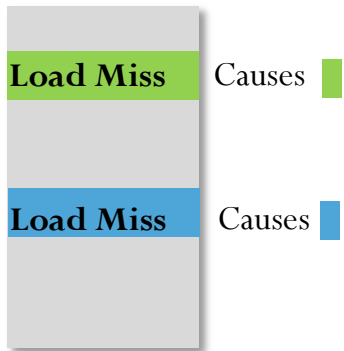
- Introduction
 - Packet Scheduling
 - Memory Level Parallelism
- Aergia
 - Concept of Slack
 - Estimating Slack
- Evaluation
- Conclusion

Slack of Packets

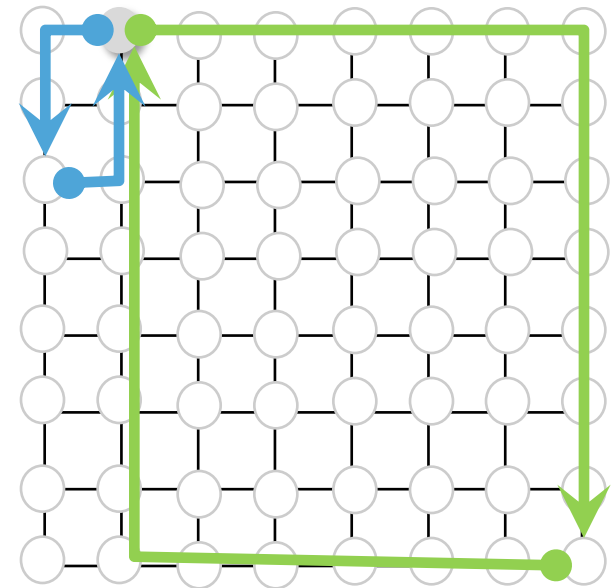
- What is slack of a packet?
 - Slack of a packet is number of cycles it can be delayed in a router without (significantly) reducing application's performance
 - **Local network slack**
- Source of slack: Memory-Level Parallelism (MLP)
 - Latency of an application's packet hidden from application due to **overlap** with latency of pending cache miss requests
- Prioritize packets with **lower slack**

Concept of Slack

Instruction Window



Network-on-Chip

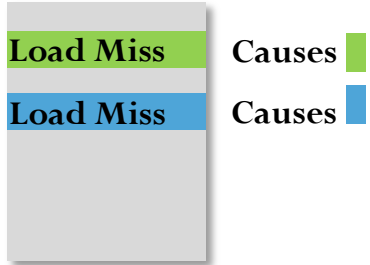


$$\text{Slack (■)} = \text{Latency (■)} - \text{Latency (■)} = 26 - 6 = 20 \text{ hops}$$

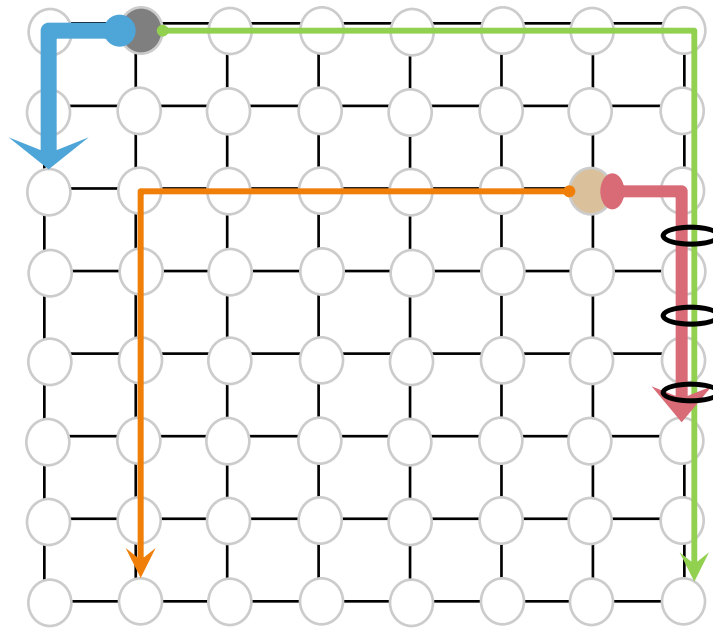
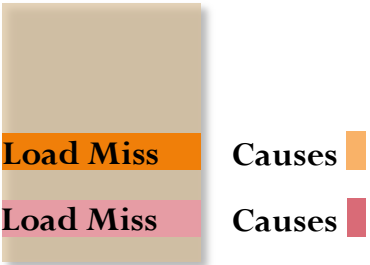
Packet(■) can be delayed for available slack cycles without reducing performance!

Prioritizing using Slack

Core A



Core B



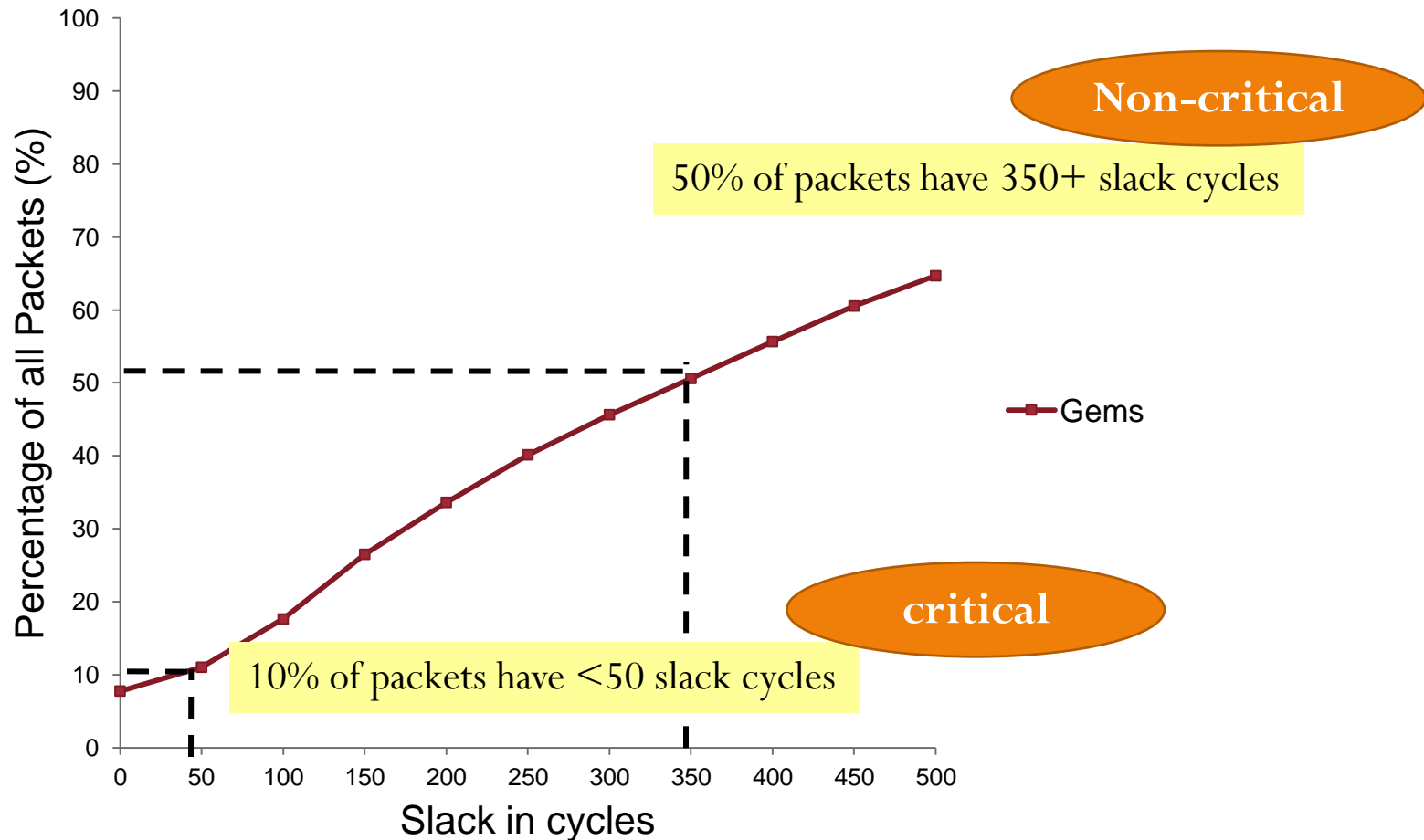
Packet	Latency	Slack
■	13 hops	0 hops
■	3 hops	10 hops

○ Interference at 3 hops

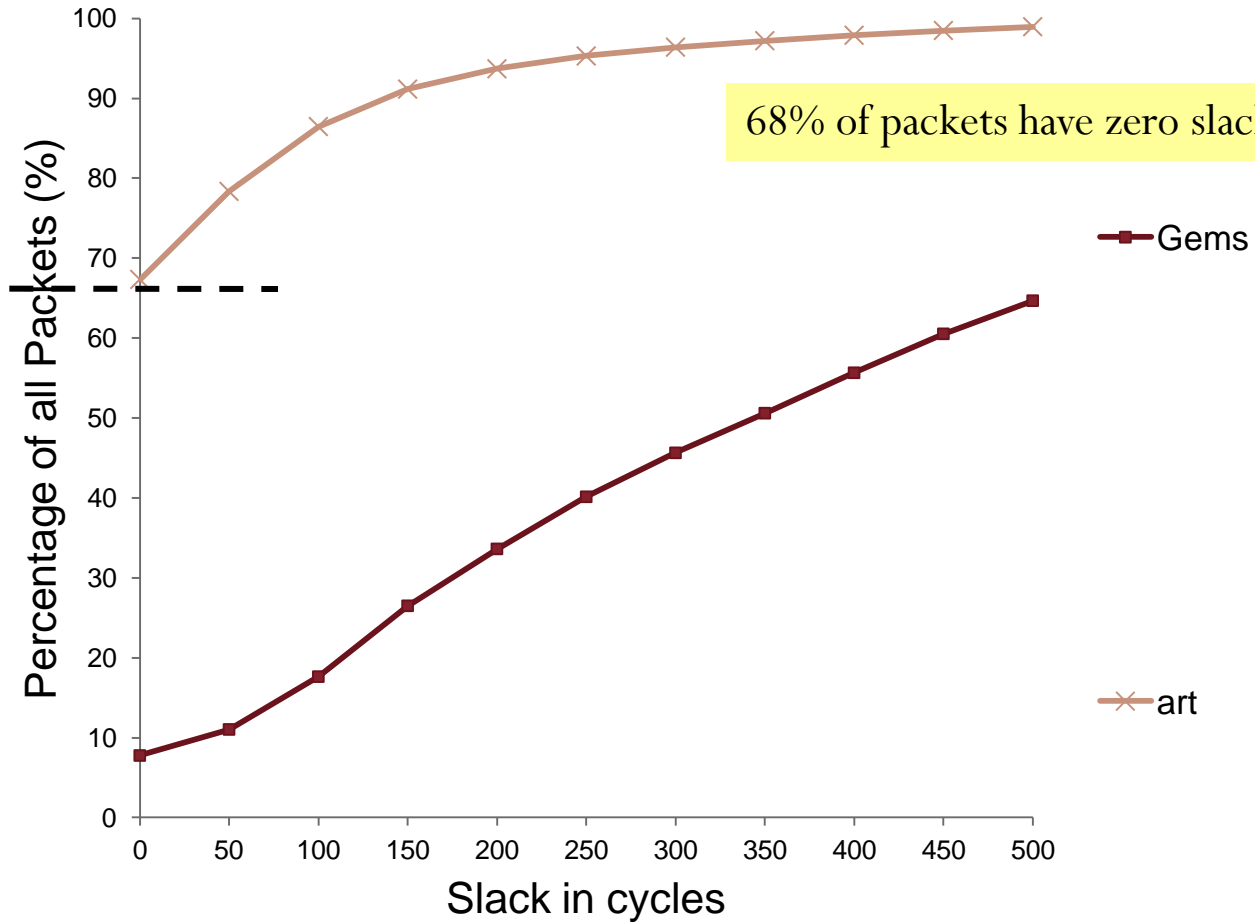
Slack(■) > Slack (■)

Prioritize ■

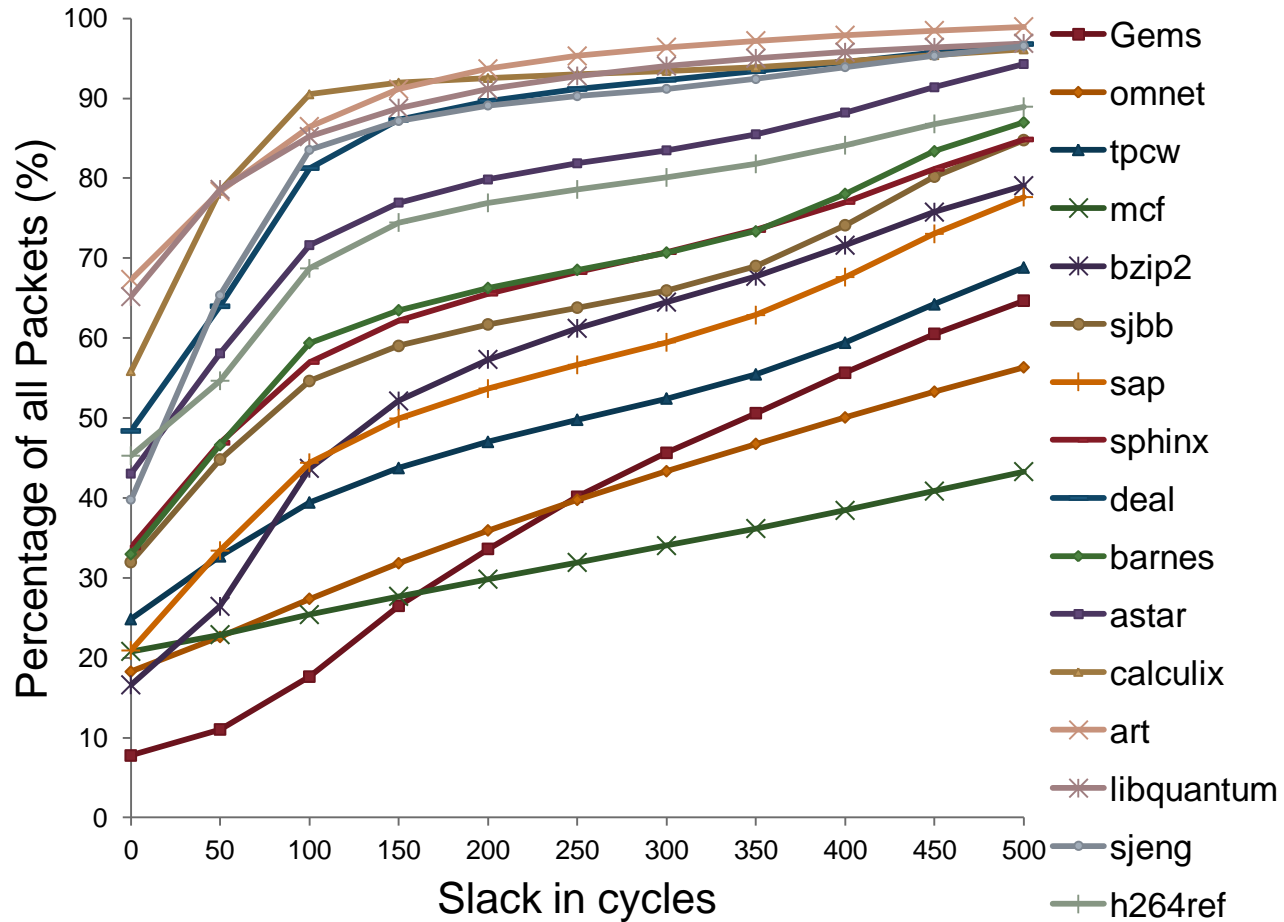
Slack in Applications



Slack in Applications



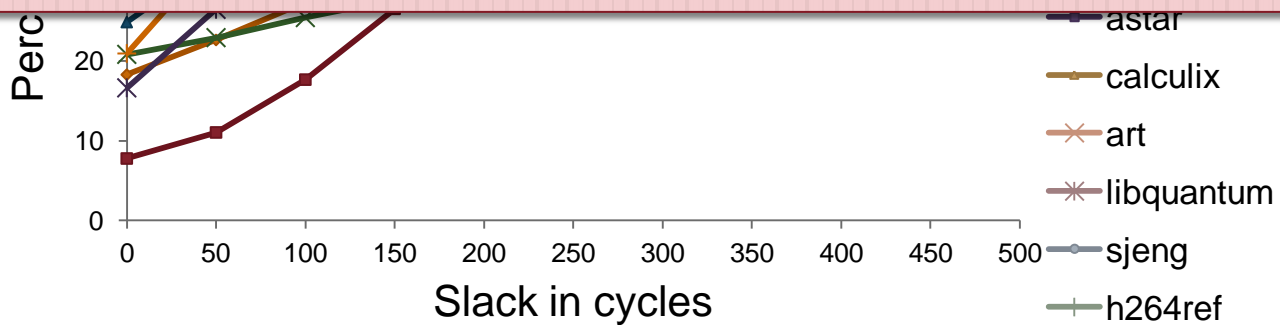
Diversity in Slack



Diversity in Slack



Slack varies **between** packets of **different** applications



Slack varies **between** packets of a **single** application



Outline

- Introduction
 - Packet Scheduling
 - Memory Level Parallelism
- Aergia
 - Concept of Slack
 - Estimating Slack
- Evaluation
- Conclusion

Estimating Slack Priority

Slack (P) = Max (Latencies of P's Predecessors) – Latency of P

Predecessors(P) are the packets of outstanding cache miss requests when P is issued

- Packet latencies not known when issued
- Predicting latency of any packet Q
 - Higher latency if Q corresponds to an L2 miss
 - Higher latency if Q has to travel farther number of hops

Estimating Slack Priority

- Slack of P = Maximum Predecessor Latency – Latency of P

- Slack(P) =

PredL2 (2 bits)	MyL2 (1 bit)	HopEstimate (2 bits)
--------------------	-----------------	-------------------------

PredL2: Set if any predecessor packet is servicing L2 miss

MyL2: Set if P is NOT servicing an L2 miss

HopEstimate: Max (# of hops of Predecessors) – hops of P

Estimating Slack Priority

- How to predict L2 hit or miss at core?
 - *Global Branch Predictor* based L2 Miss Predictor
 - Use Pattern History Table and 2-bit saturating counters
 - *Threshold* based L2 Miss Predictor
 - If #L2 misses in “M” misses \geq “T” threshold then next load is a L2 miss.
- Number of miss predecessors?
 - List of outstanding L2 Misses
- Hops estimate?
 - Hops $\Rightarrow \Delta X + \Delta Y$ distance
 - Use predecessor list to calculate slack hop estimate

Starvation Avoidance

- Problem: **Starvation**
 - Prioritizing packets can lead to starvation of lower priority packets
- Solution: **Time-Based Packet Batching**
 - New batches are formed at every T cycles
 - Packets of older batches are prioritized over younger batches

Putting it all together

- Tag header of the packet with priority bits before injection

Priority (P) =

Batch (3 bits)	PredL2 (2 bits)	MyL2 (1 bit)	HopEstimate (2 bits)
-------------------	--------------------	-----------------	-------------------------

- Priority(P)?
 - P's batch *(highest priority)*
 - P's Slack
 - Local Round-Robin *(final tie breaker)*

Outline

- Introduction
 - Packet Scheduling
 - Memory Level Parallelism
- Aergia
 - Concept of Slack
 - Estimating Slack
- Evaluation
- Conclusion

Evaluation Methodology

- 64-core system
 - x86 processor model based on Intel Pentium M
 - 2 GHz processor, 128-entry instruction window
 - 32KB private L1 and 1MB per core shared L2 caches, 32 miss buffers
 - 4GB DRAM, 320 cycle access latency, 4 on-chip DRAM controllers
- Detailed Network-on-Chip model
 - 2-stage routers (with speculation and look ahead routing)
 - Wormhole switching (8 flit data packets)
 - Virtual channel flow control (6 VCs, 5 flit buffer depth)
 - 8x8 Mesh (128 bit bi-directional channels)
- Benchmarks
 - Multiprogrammed scientific, server, desktop workloads (35 applications)
 - 96 workload combinations

Qualitative Comparison

- **Round Robin & Age**

- Local and application oblivious
- Age is biased towards heavy applications

- **Globally Synchronized Frames (GSF)**

[Lee et al., ISCA 2008]

- Provides **bandwidth fairness** at the expense of **system performance**
- Penalizes heavy and bursty applications

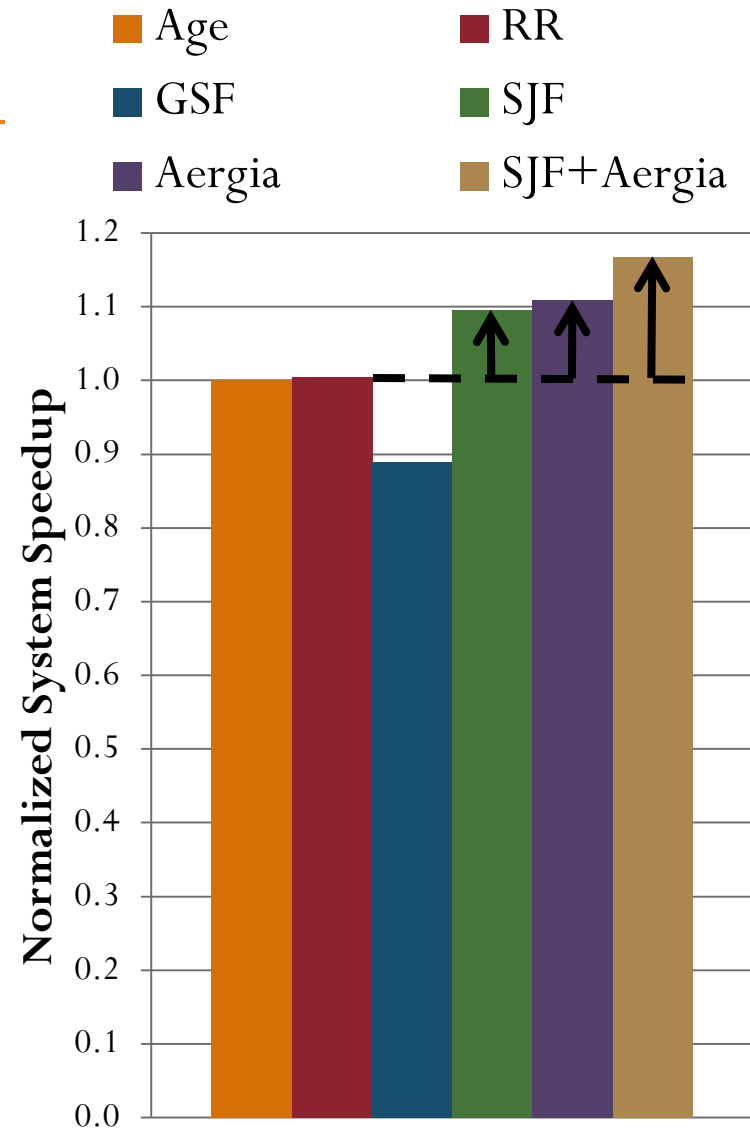
- **Application-Aware Prioritization Policies (SJF)**

[Das et al., MICRO 2009]

- **Shortest-Job-First Principle**
- Packet scheduling policies which prioritize network sensitive applications which inject lower load

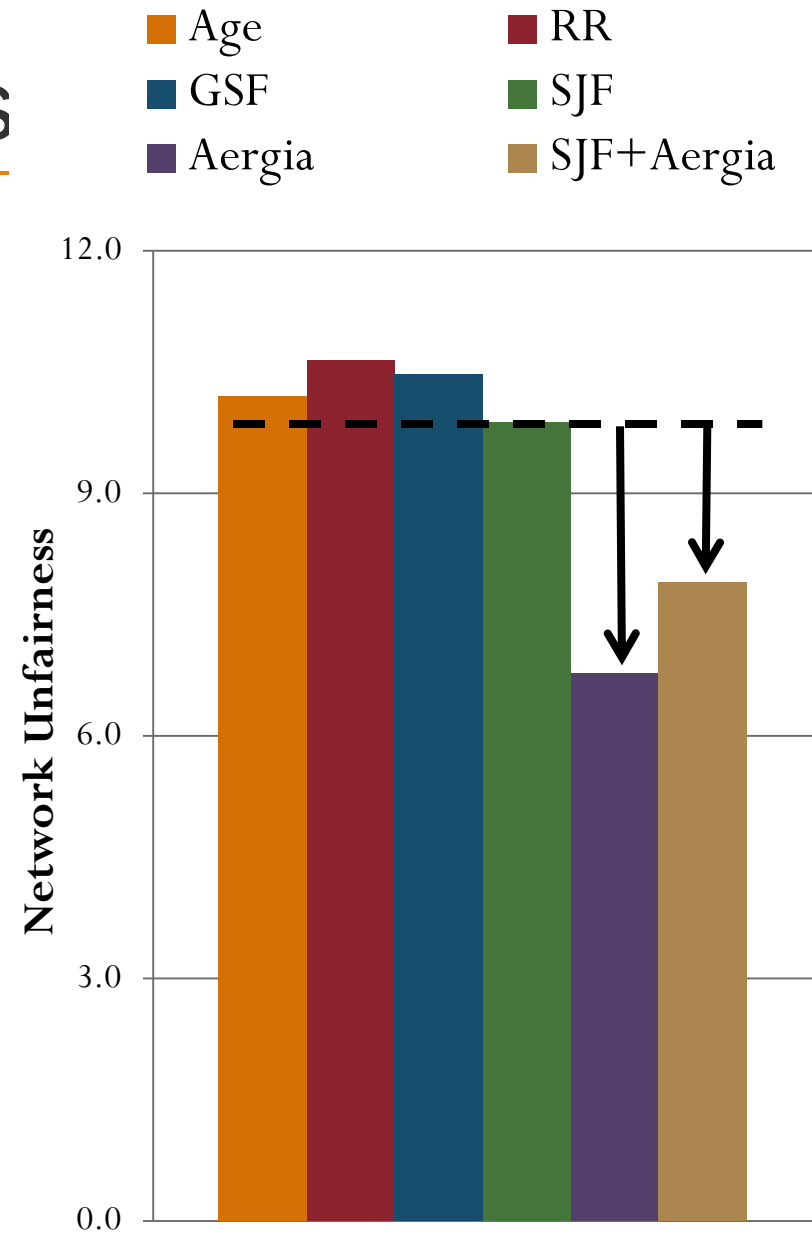
System Performance

- SJF provides 8.9% improvement in weighted speedup
- Aergia improves system throughput by 10.3%
- Aergia+SJF improves system throughput by 16.1%



Network Unfairness

- SJF does not imbalance network fairness
- Aergia improves network unfairness by 1.5X
- SJF+Aergia improves network unfairness by 1.3X



Conclusions & Future Directions

- Packets have different criticality, yet existing packet scheduling policies **treat all packets equally**
- We propose a new approach to packet scheduling in NoCs
 - We define **Slack** as a key measure that characterizes the relative importance of a packet.
 - We propose **Aérgia** a novel architecture to accelerate low slack critical packets
- Result
 - Improves system performance: 16.1%
 - Improves network fairness: 30.8%

Express-Cube Topologies

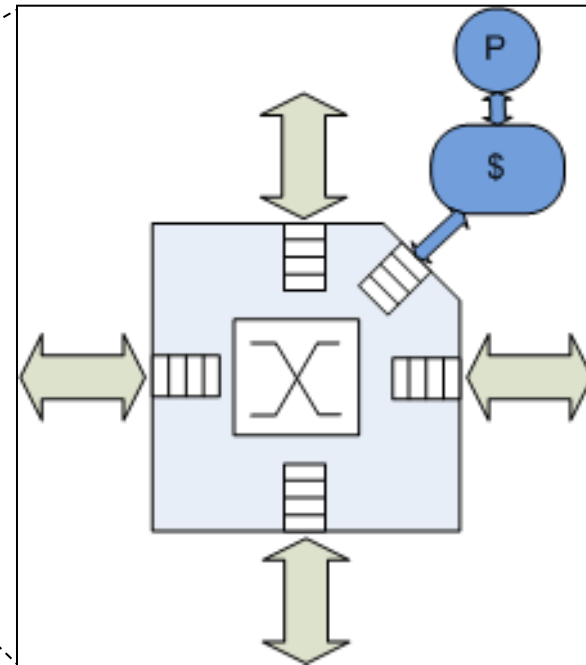
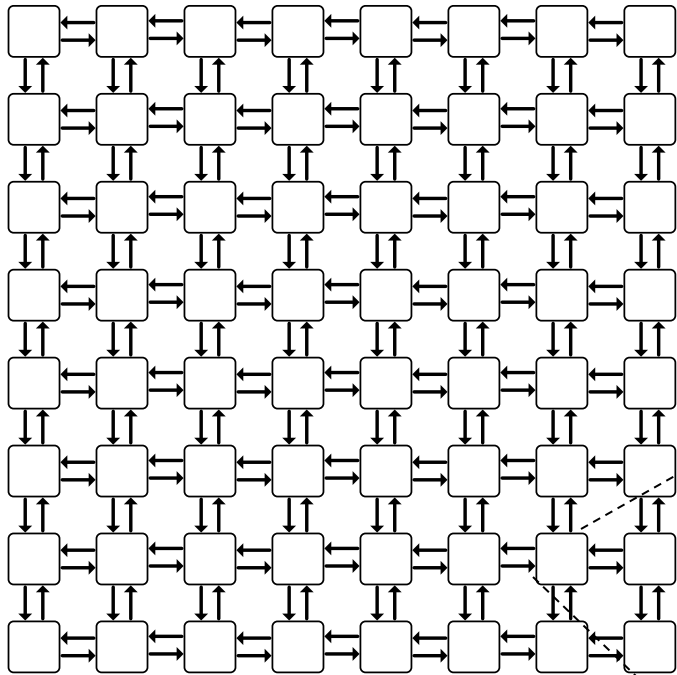
Boris Grot, Joel Hestness, Stephen W. Keckler, and Onur Mutlu,

"Express Cube Topologies for On-Chip Interconnects"

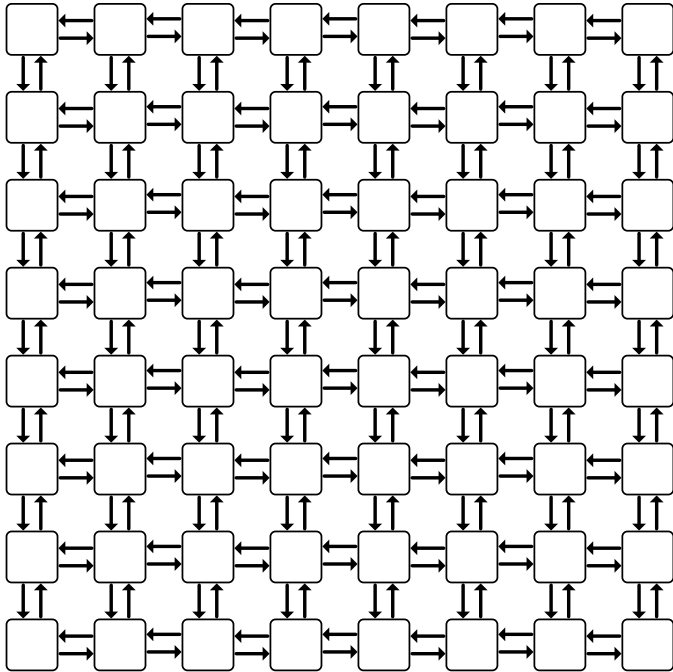
Proceedings of the 15th International Symposium on High-Performance Computer Architecture (HPCA), pages 163-174, Raleigh, NC, February 2009.

Slides (ppt)

2-D Mesh

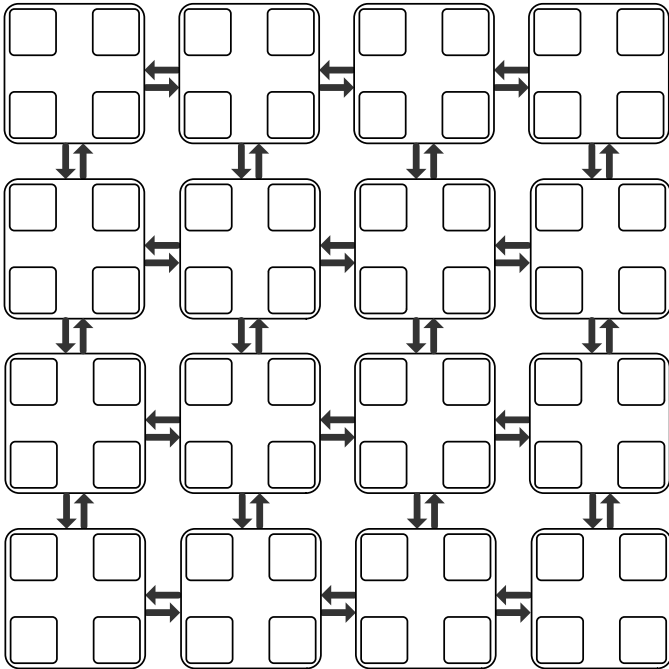


2-D Mesh



- Pros
 - Low design & layout complexity
 - Simple, fast routers
- Cons
 - Large diameter
 - Energy & latency impact

Concentration *(Balfour & Dally, ICS '06)*



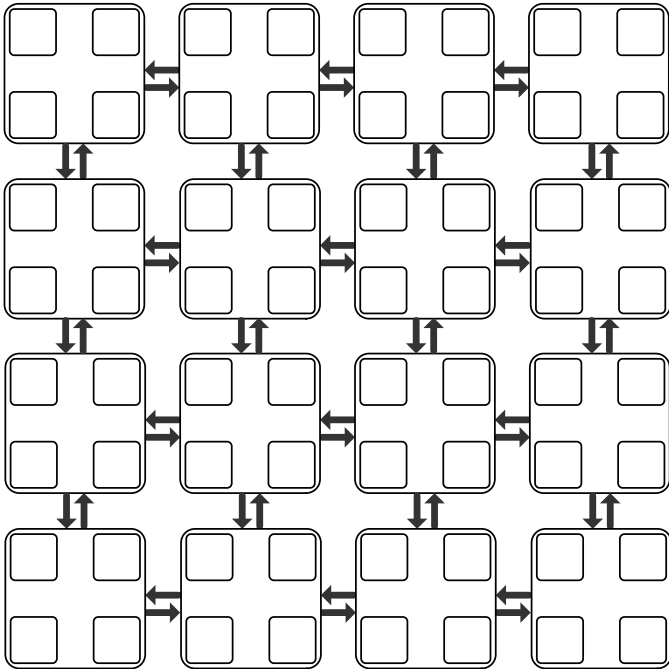
□ Pros

- Multiple *terminals* attached to a router node
- Fast nearest-neighbor communication via the crossbar
- Hop count reduction proportional to *concentration* degree

□ Cons

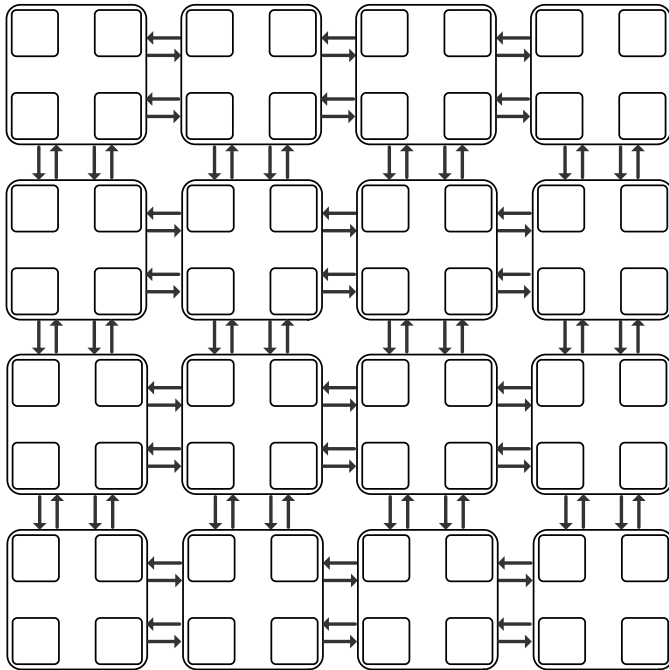
- Benefits limited by crossbar complexity

Concentration



- Side-effects
 - Fewer channels
 - Greater channel width

Replication

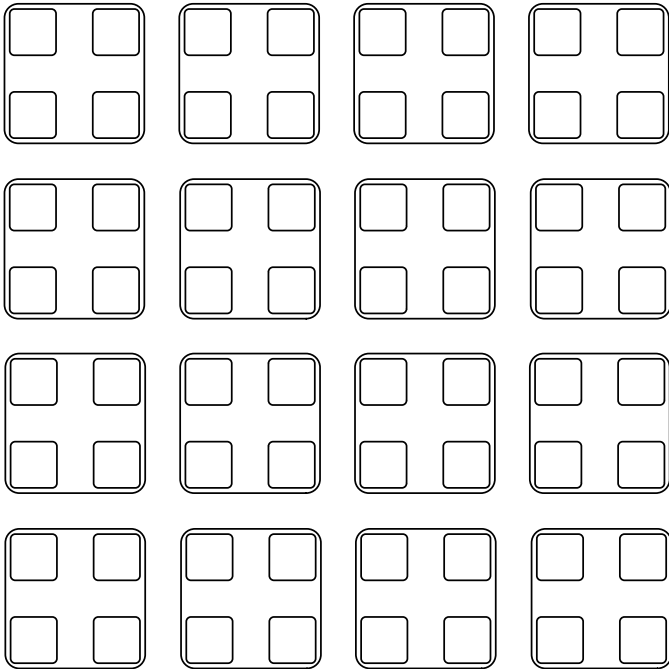


CMesh-X2

□ Benefits

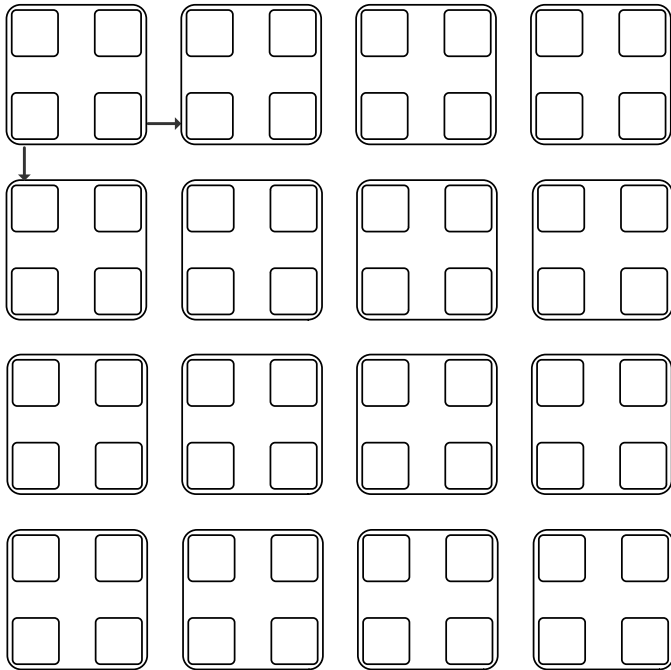
- Restores bisection channel count
- Restores channel width
- Reduced crossbar complexity

Flattened Butterfly *(Kim et al., Micro '07)*

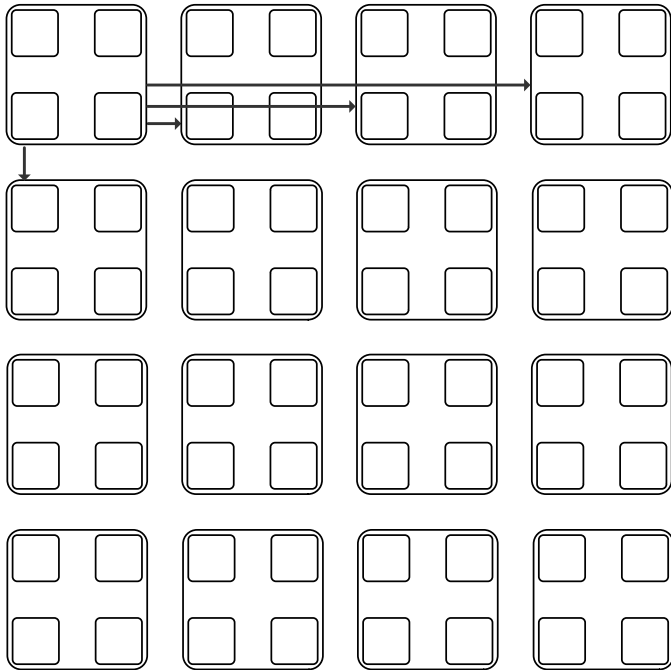


- ▣ Objectives:
 - Improve connectivity
 - Exploit the wire budget

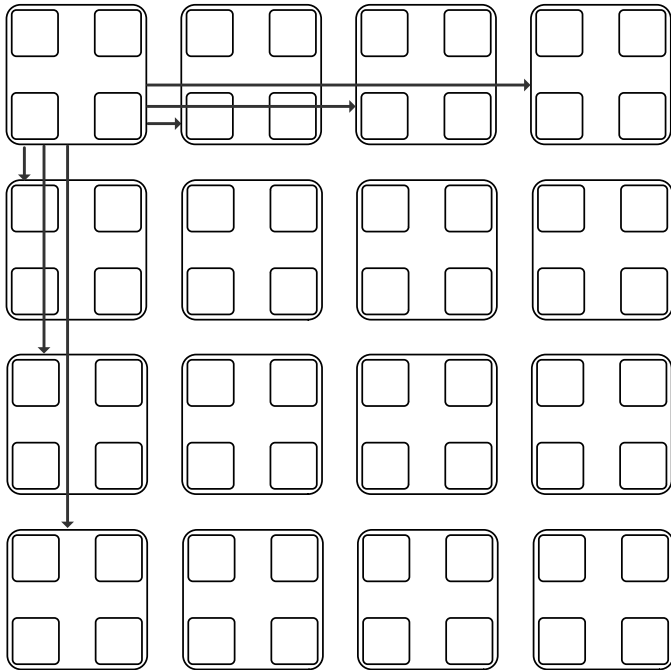
Flattened Butterfly *(Kim et al., Micro '07)*



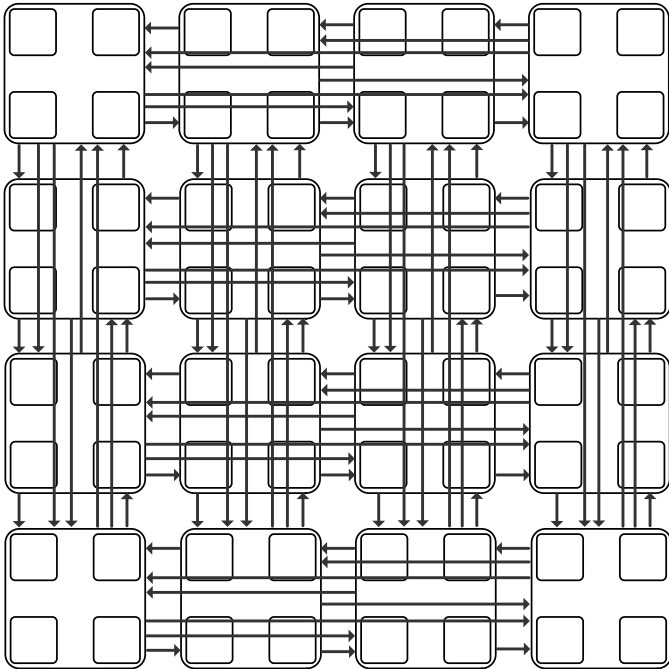
Flattened Butterfly *(Kim et al., Micro '07)*



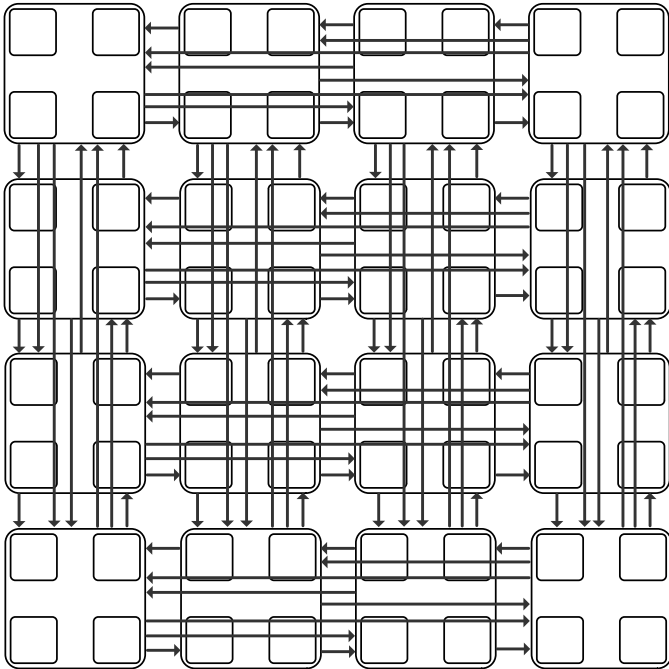
Flattened Butterfly *(Kim et al., Micro '07)*



Flattened Butterfly *(Kim et al., Micro '07)*

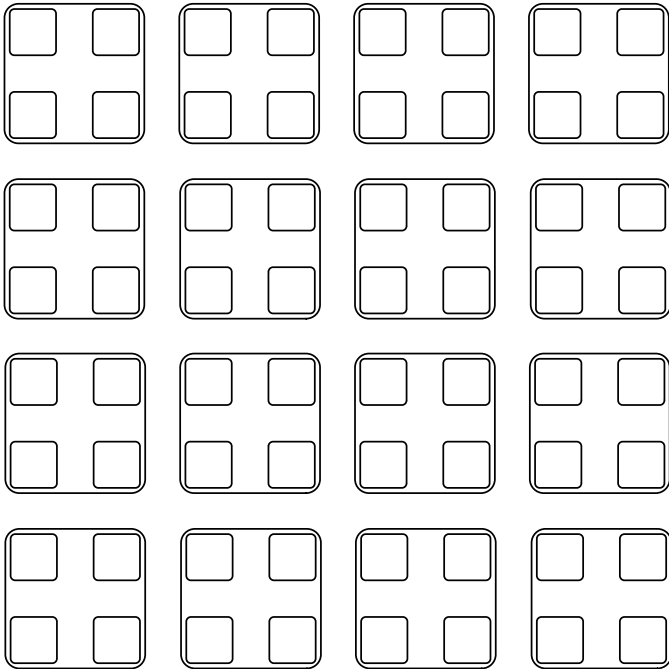


Flattened Butterfly (*Kim et al., Micro '07*)



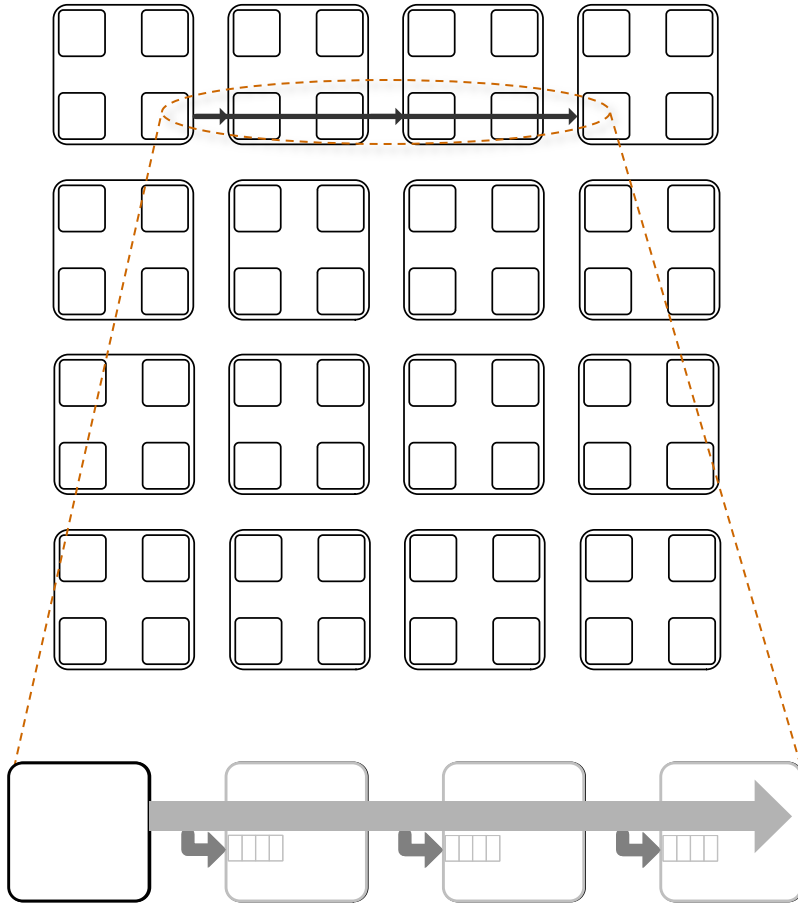
- Pros
 - Excellent connectivity
 - Low diameter: 2 hops
- Cons
 - High channel count: $k^2/2$ per row/column
 - Low channel utilization
 - Increased control (arbitration) complexity

Multidrop Express Channels (MECS)

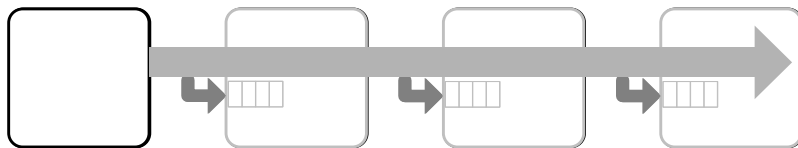
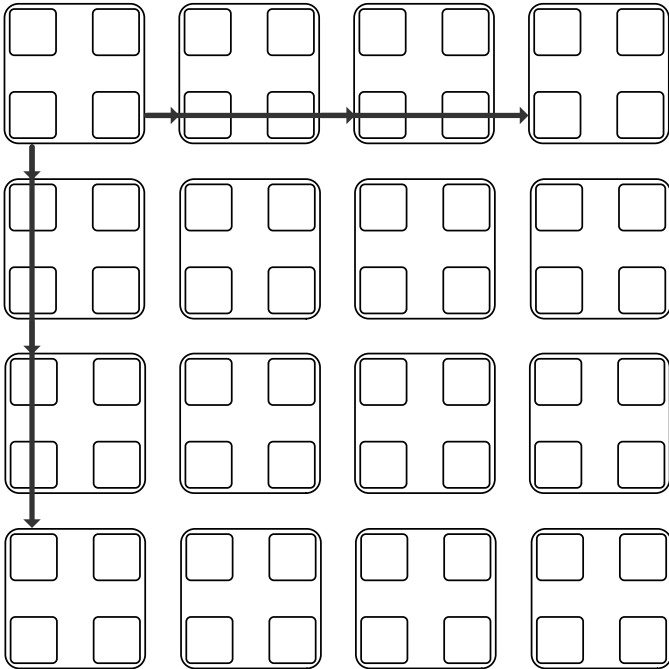


- Objectives:
 - Connectivity
 - More scalable channel count
 - Better channel utilization

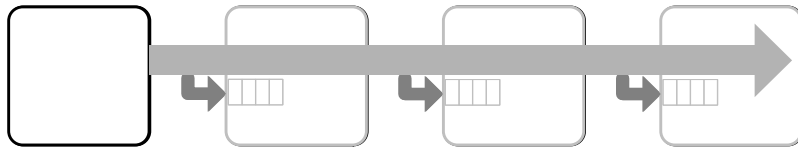
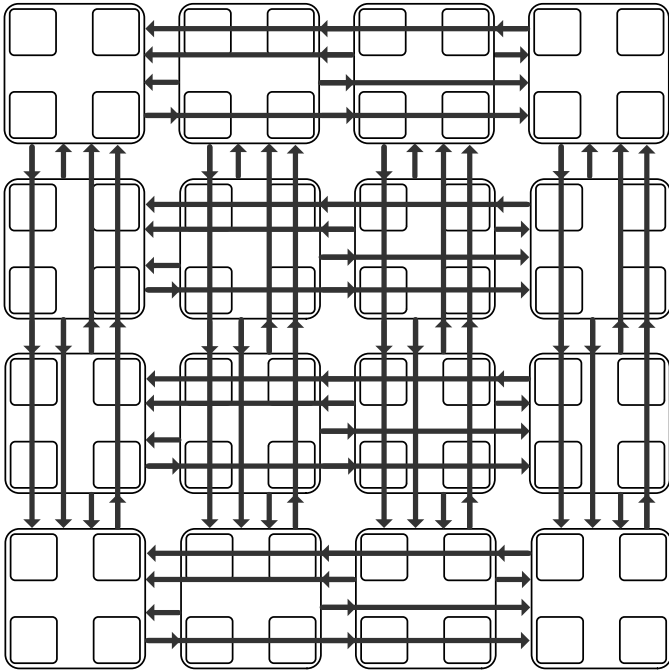
Multidrop Express Channels (MECS)



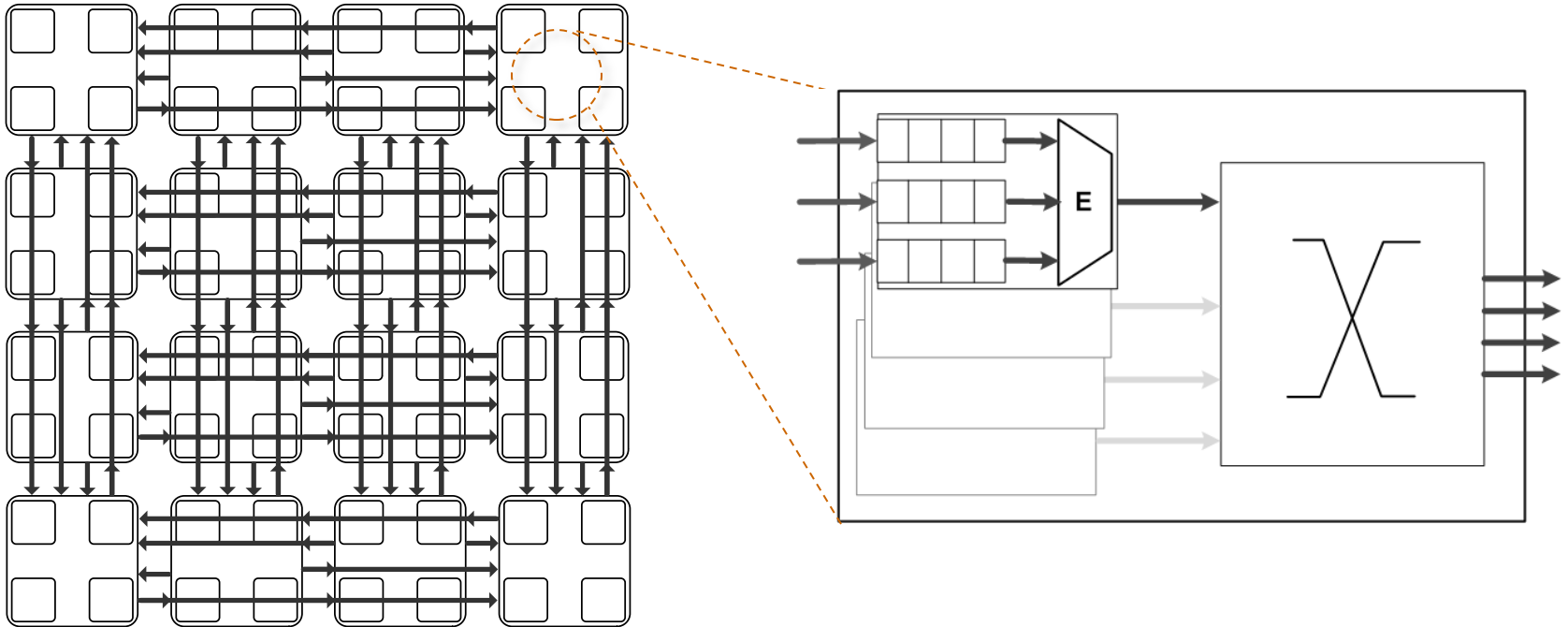
Multidrop Express Channels (MECS)



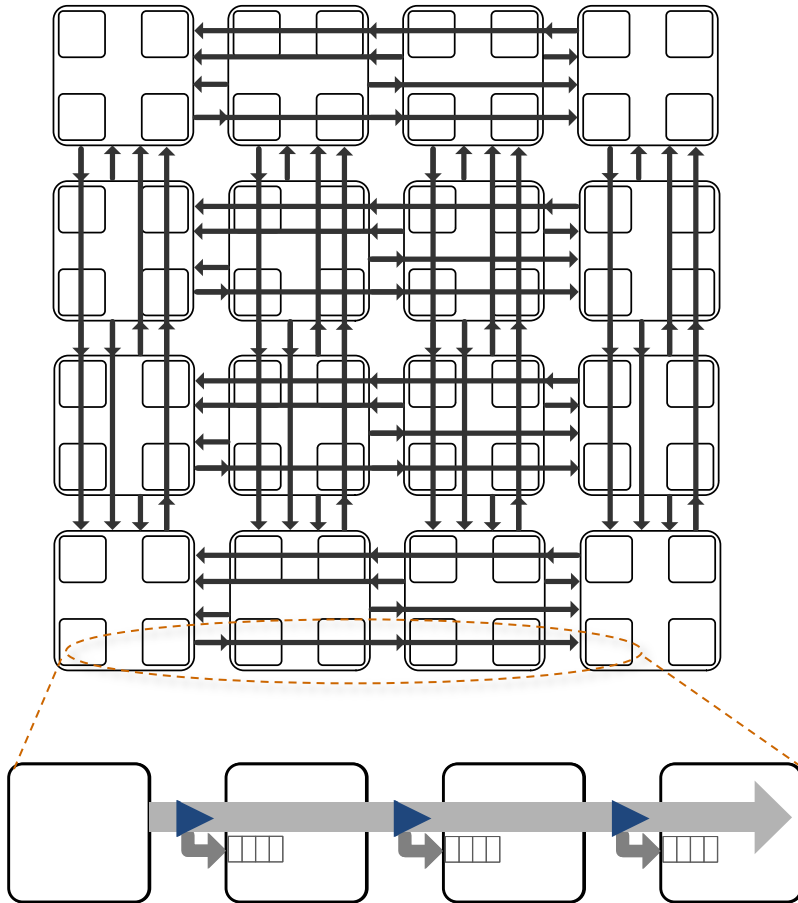
Multidrop Express Channels (MECS)



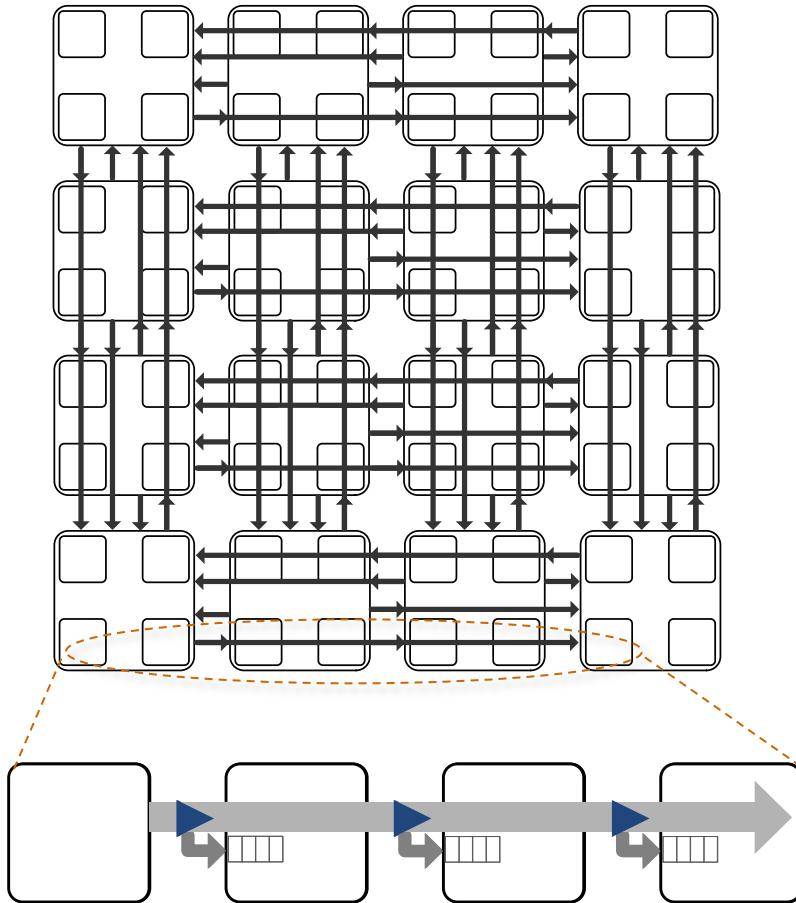
Multidrop Express Channels (MECS)



Multidrop Express Channels (MECS)

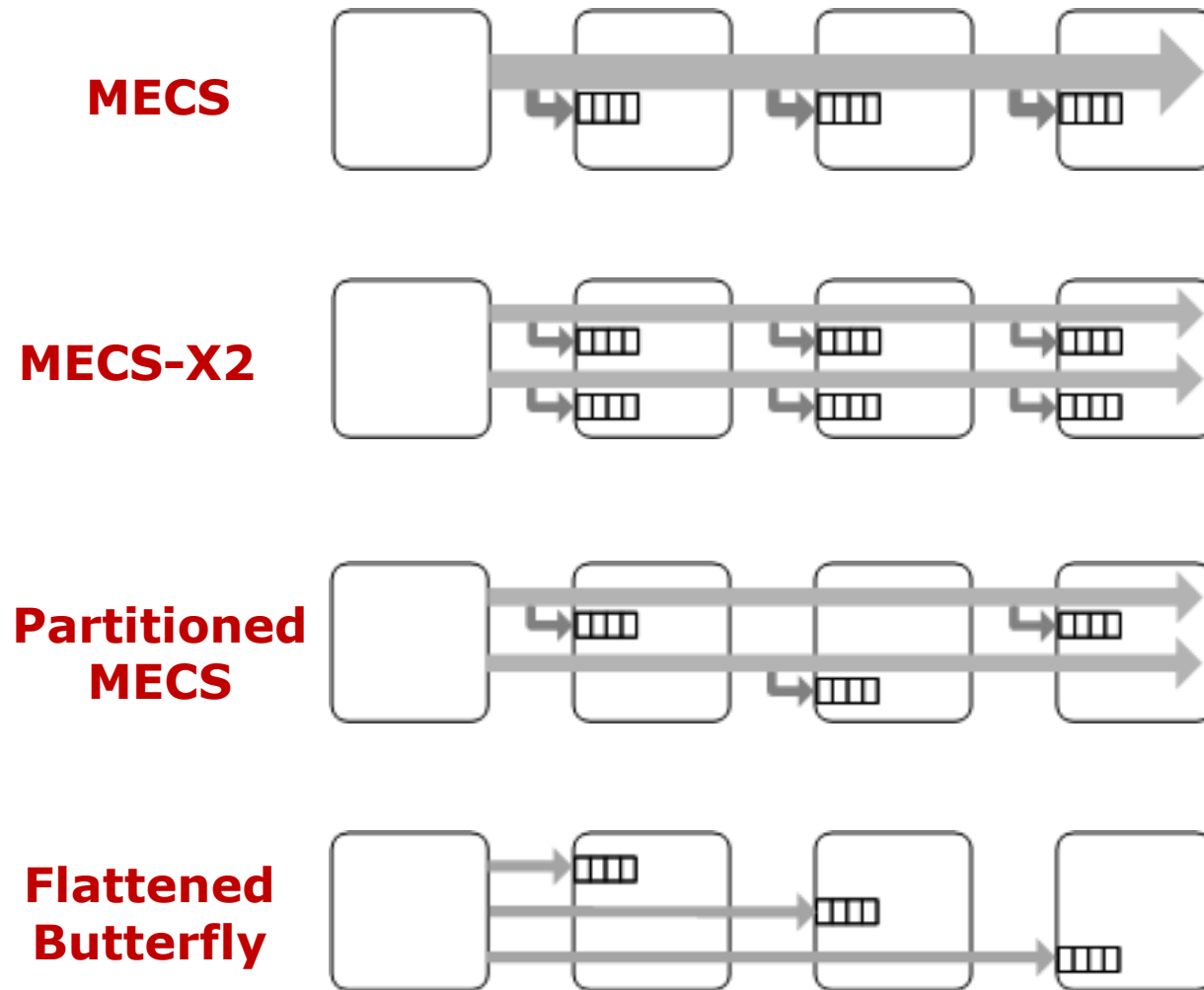


Multidrop Express Channels (MECS)



- Pros
 - One-to-many topology
 - Low diameter: 2 hops
 - k channels row/column
 - Asymmetric
- Cons
 - Asymmetric
 - Increased control (arbitration) complexity

Partitioning: a GEC Example



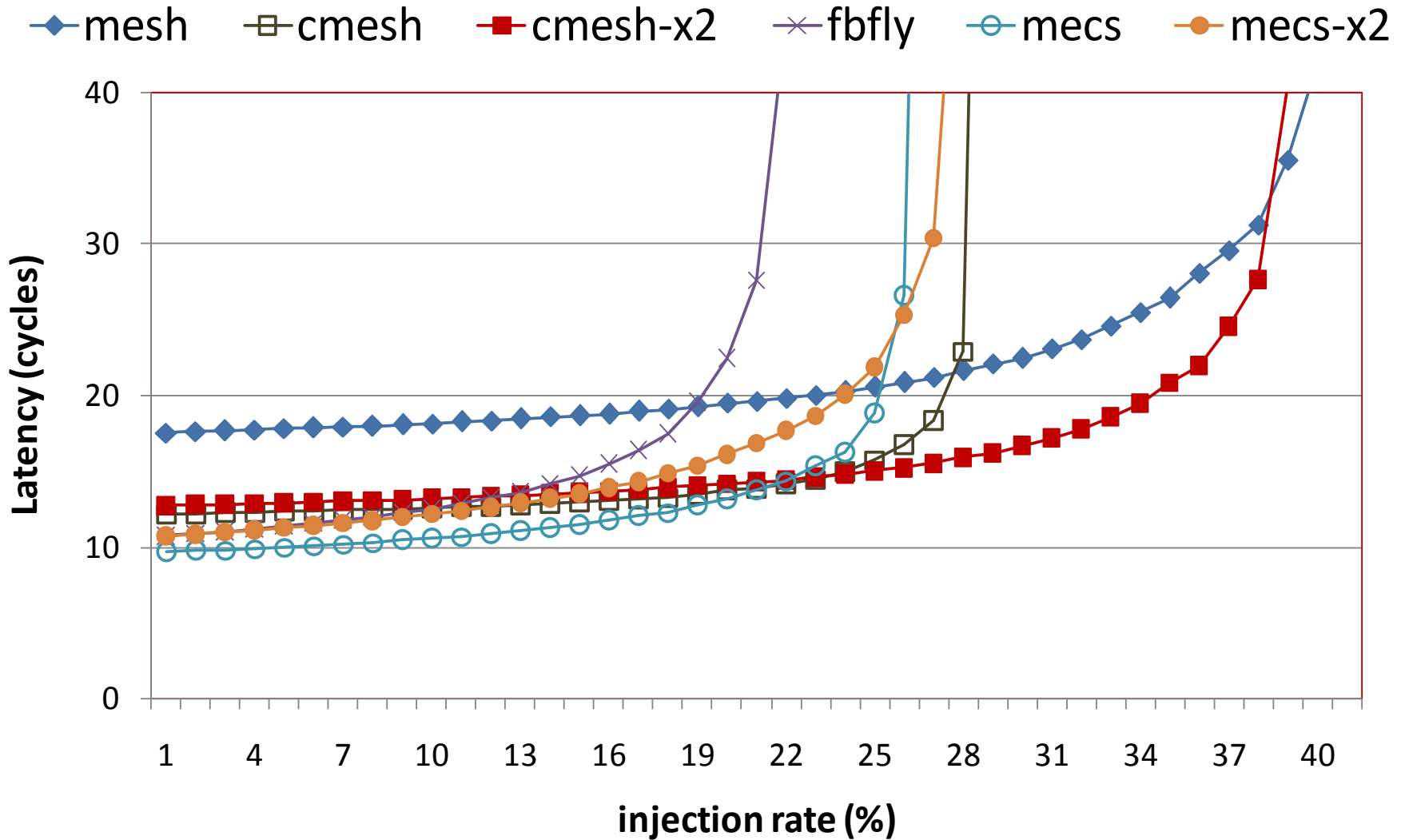
Analytical Comparison

	CMesh		FBfly		MECS	
Network Size	64	256	64	256	64	256
Radix (conctr' d)	4	8	4	8	4	8
Diameter	6	14	2	2	2	2
Channel count	2	2	8	32	4	8
Channel width	576	1152	144	72	288	288
Router inputs	4	4	6	14	6	14
Router outputs	4	4	6	14	4	4

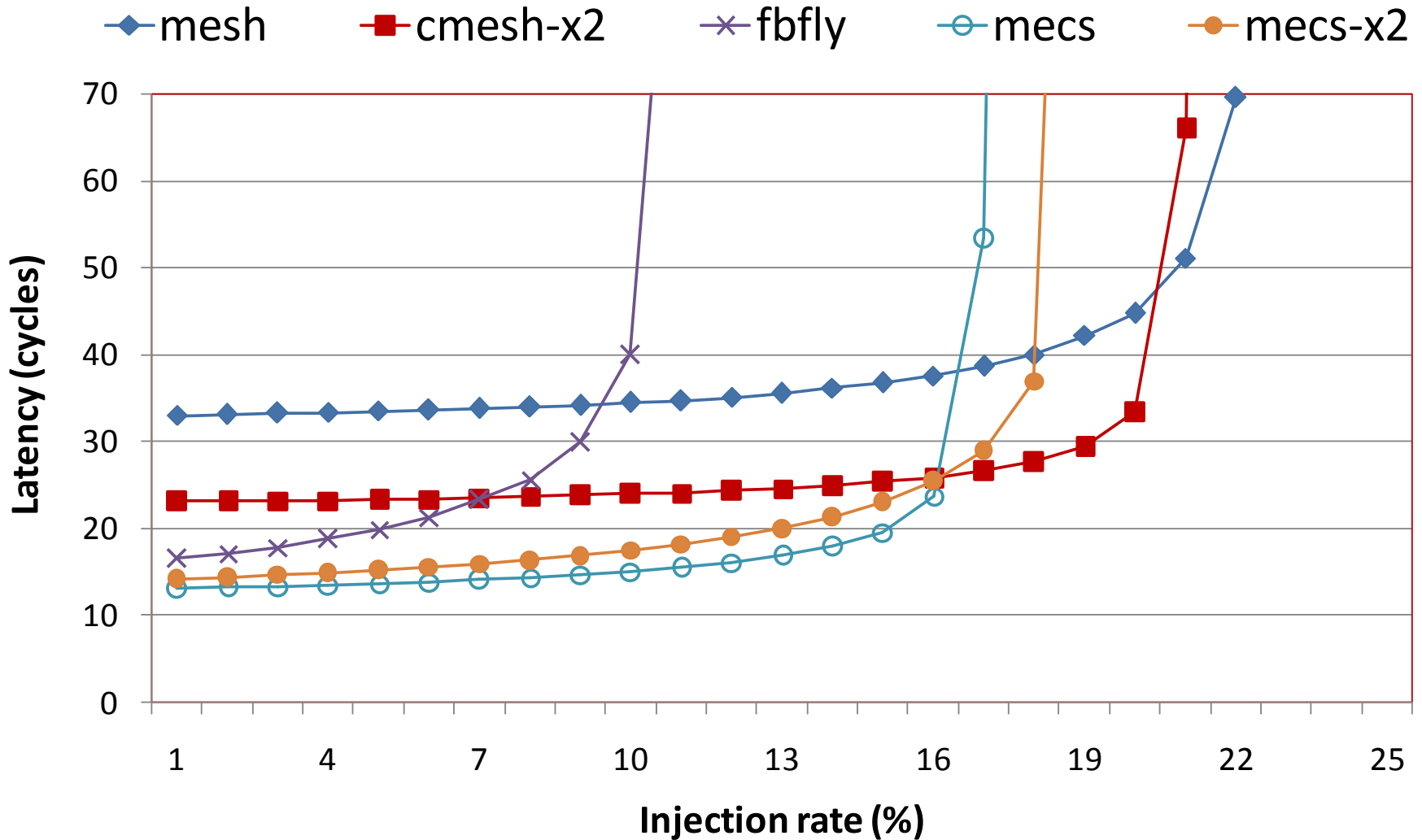
Experimental Methodology

Topologies	Mesh, CMesh, CMesh-X2, FBFly, MECS, MECS-X2
Network sizes	64 & 256 terminals
Routing	DOR, adaptive
Messages	64 & 576 bits
Synthetic traffic	Uniform random, bit complement, transpose, self-similar
PARSEC benchmarks	Blackscholes, Bodytrack, Canneal, Ferret, Fluidanimate, Freqmine, Vip, x264
Full-system config	M5 simulator, Alpha ISA, 64 000 cores
Energy evaluation	Orion + CACTI 6

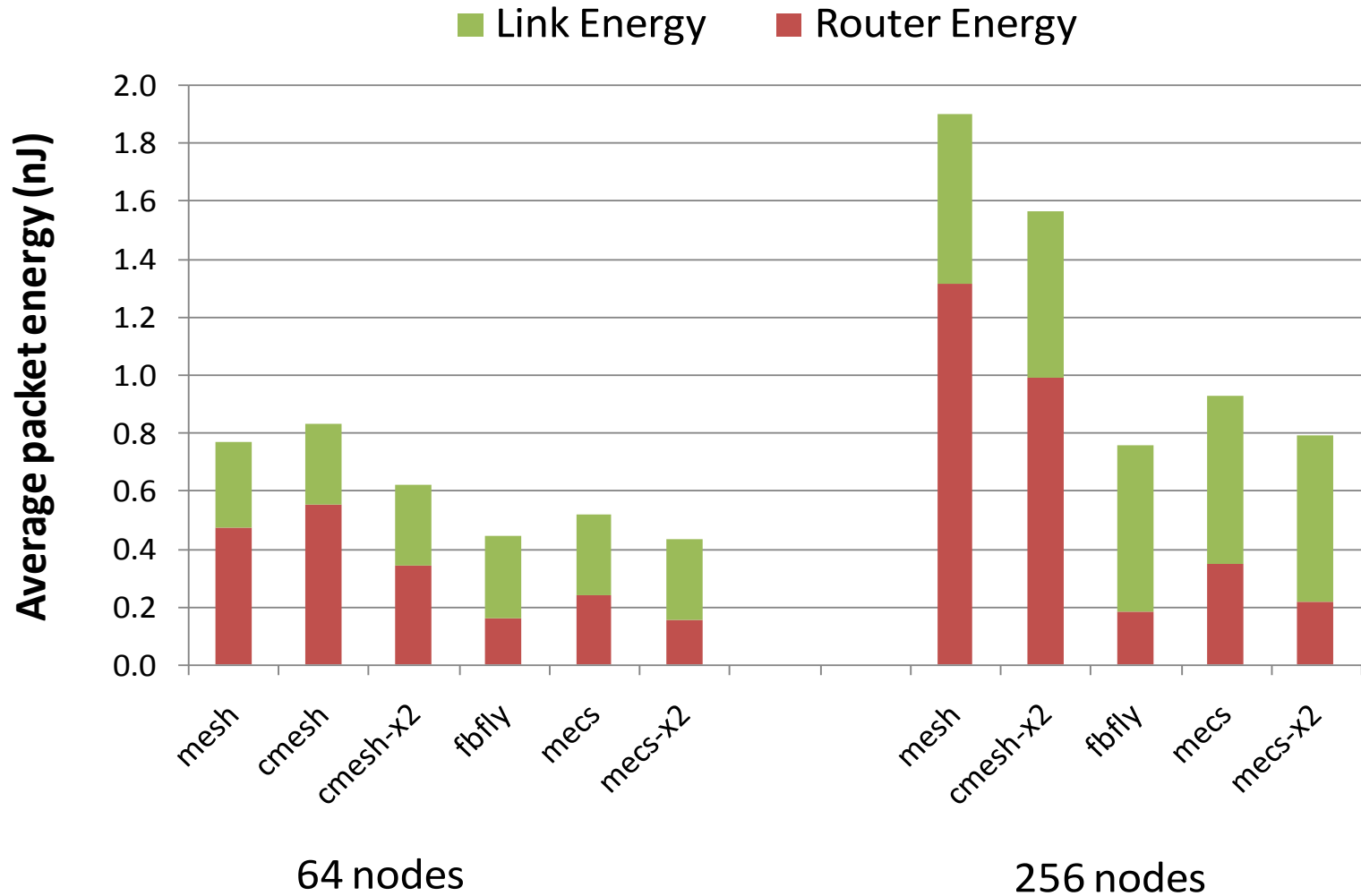
64 nodes: Uniform Random



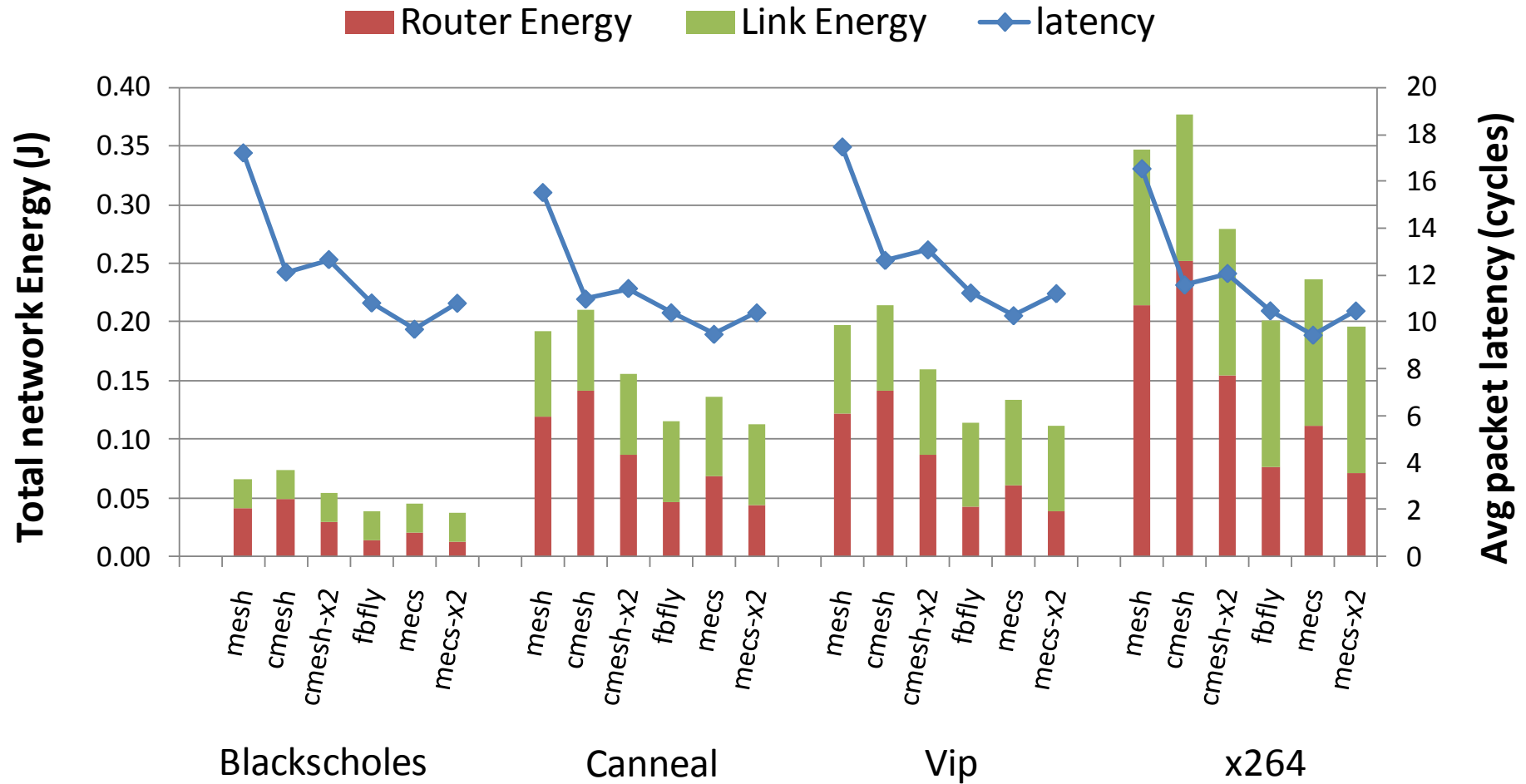
256 nodes: Uniform Random



Energy (100K pkts, Uniform Random)



64 Nodes: PARSEC



Summary

□ MECS

- A new one-to-many topology
- Good fit for planar substrates
- Excellent connectivity
- Effective wire utilization

□ Generalized Express Cubes

- Framework & taxonomy for NOC topologies
- Extension of the k -ary n -cube model
- Useful for understanding and exploring on-chip interconnect options
- Future: expand & formalize

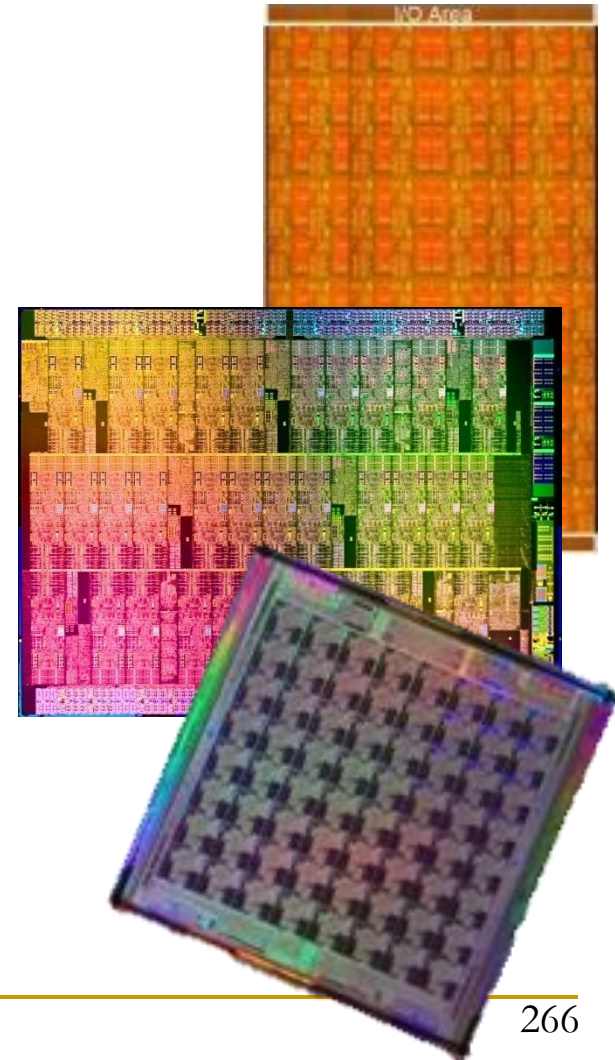
Kilo-NoC: Topology-Aware QoS

Boris Grot, Joel Hestness, Stephen W. Keckler, and Onur Mutlu,
"Kilo-NOC: A Heterogeneous Network-on-Chip Architecture for Scalability and Service Guarantees"

Proceedings of the 38th International Symposium on Computer Architecture (ISCA), San Jose, CA, June 2011. Slides (pptx)

Motivation

- Extreme-scale chip-level integration
 - Cores
 - Cache banks
 - Accelerators
 - I/O logic
 - Network-on-chip (NOC)
- 10-100 cores today
- 1000+ assets in the near future



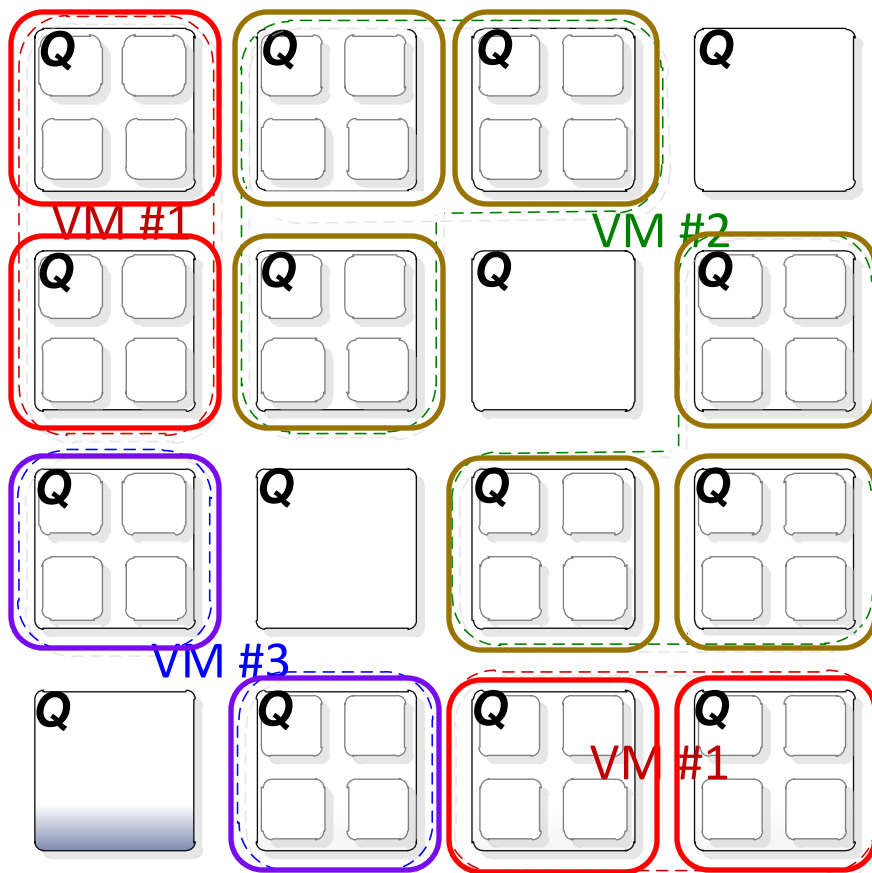
Kilo-NOC requirements

- High efficiency
 - Area
 - Energy
- Good performance
- Strong service guarantees (QoS)




Topology-Aware QoS

- Problem: QoS support in each router is expensive (in terms of buffering, arbitration, bookkeeping)
 - E.g., Grot et al., "Preemptive Virtual Clock: A Flexible, Efficient, and Cost-effective QOS Scheme for Networks-on-Chip," MICRO 2009.
- Goal: Provide QoS guarantees at low area and power cost
- Idea:
 - Isolate shared resources in a region of the network, support QoS within that area
 - Design the topology so that applications can access the region without interference

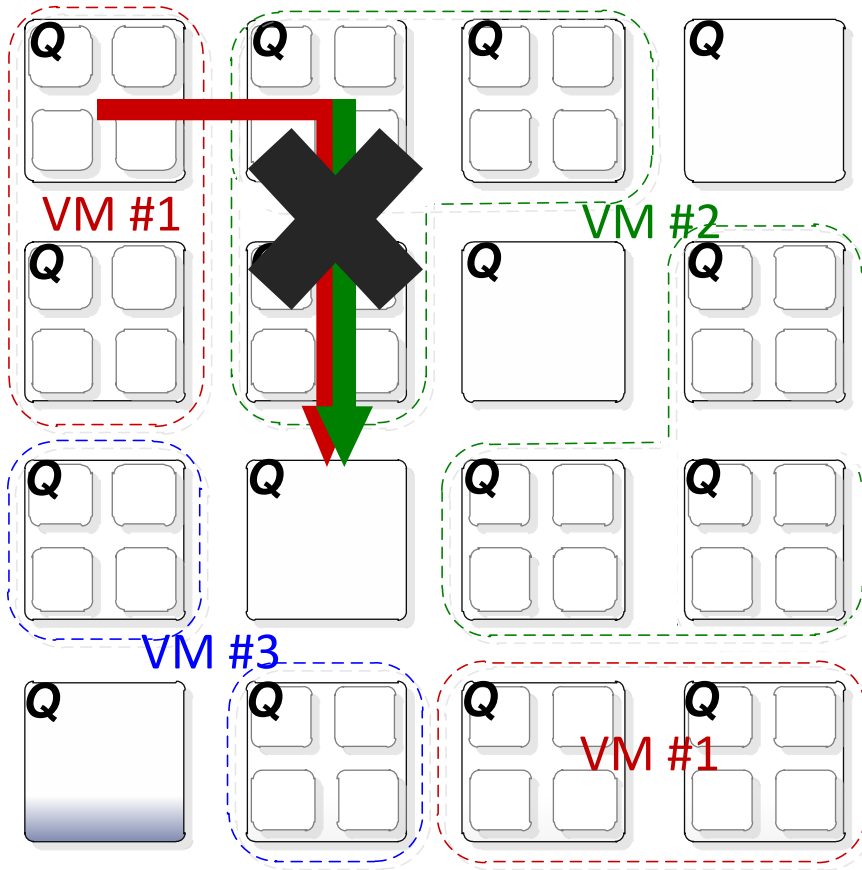
Baseline QOS-enabled CMP



Multiple VMs
sharing a die

-  Shared resources
(e.g., memory controllers)
-  VM-private resources
(cores, caches)
-  QOS-enabled router

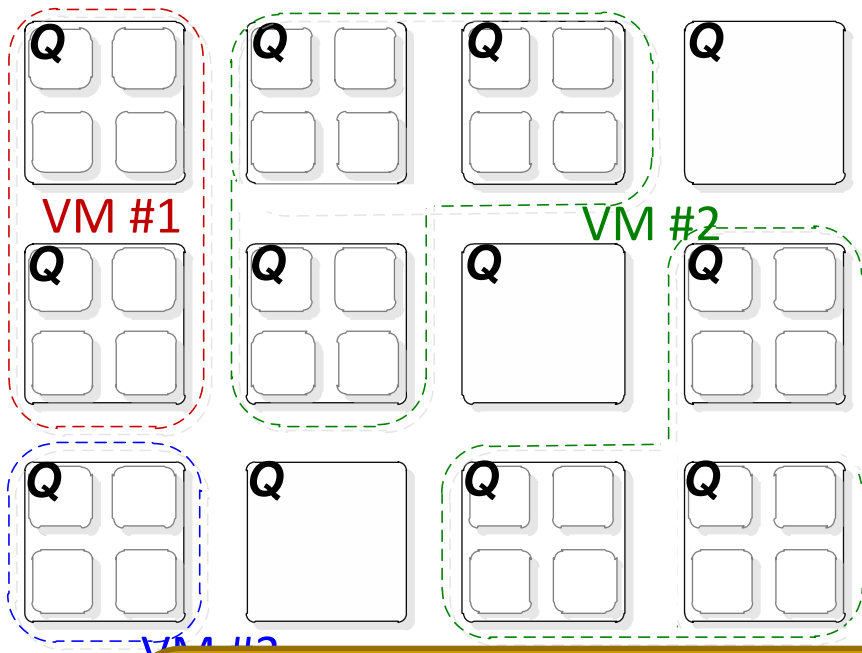
Conventional NOC QoS



Contention scenarios:

- Shared resources
 - memory access
- Intra-VM traffic
 - shared cache access
- Inter-VM traffic
 - VM page sharing

Conventional NOC QOS



Contention scenarios:

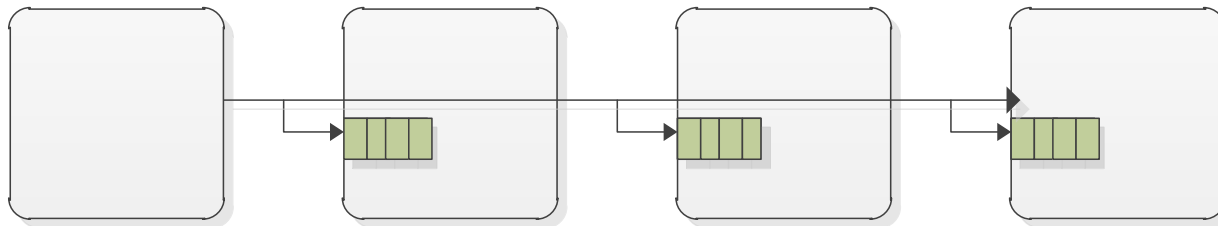
- Shared resources
 - memory access
- Intra-VM traffic
 - shared cache access
- Inter-VM traffic
 - VM page sharing

Network-wide guarantees *without* network-wide QOS support

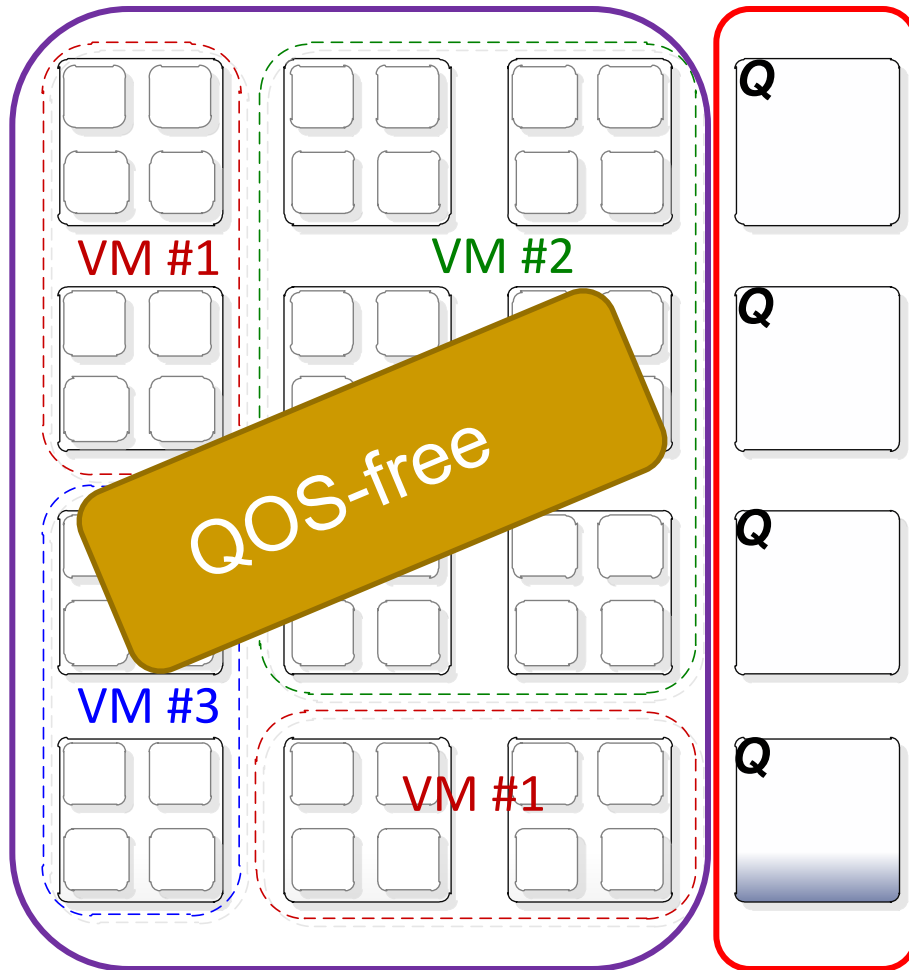
Kilo-NOC QOS

- Insight: leverage rich network connectivity
 - Naturally reduce interference among flows
 - Limit the extent of hardware QOS support
- Requires a low-diameter topology
 - This work: Multidrop Express Channels (MECS)

*Grot et al., HPCA
2009*

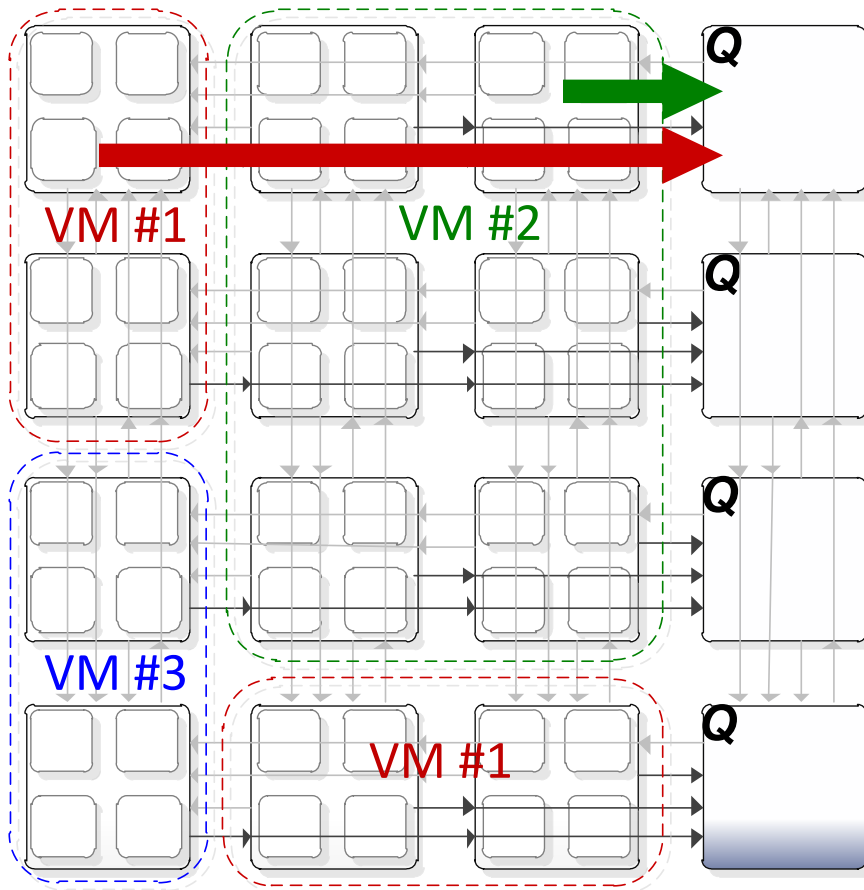


Topology-Aware QOS



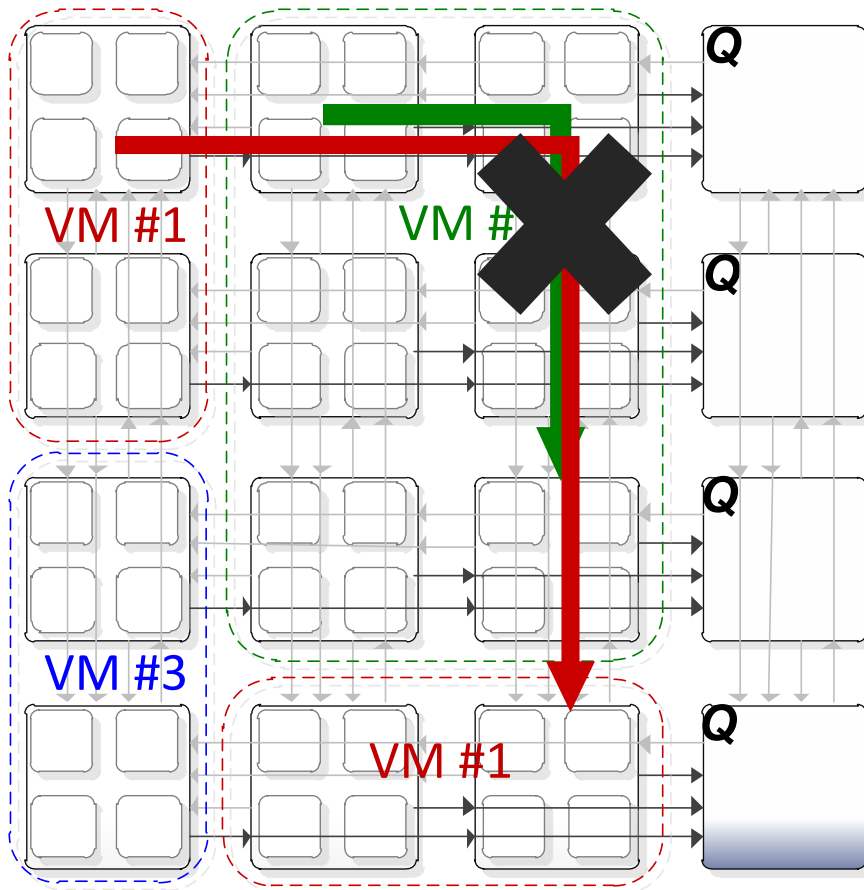
- Dedicated, QOS-enabled regions
 - Rest of die: QOS-free
- Richly-connected topology
 - Traffic isolation
- Special routing rules
 - Manage interference

Topology-Aware QoS



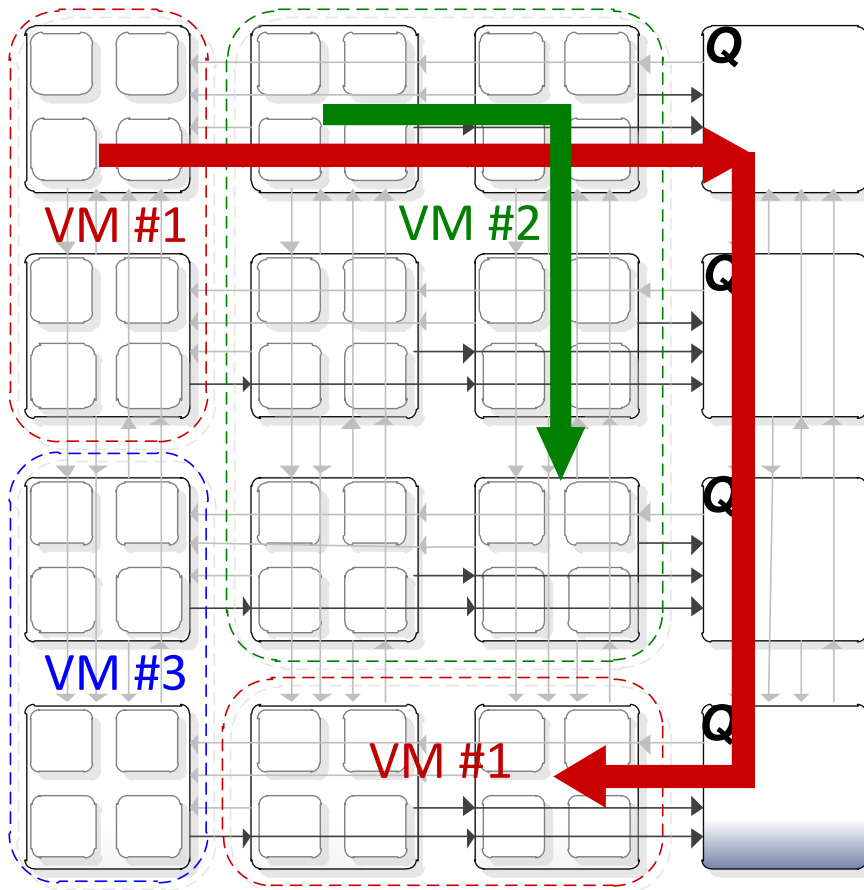
- Dedicated, QoS-enabled regions
 - Rest of die: QoS-free
- Richly-connected topology
 - Traffic isolation
- Special routing rules
 - Manage interference

Topology-Aware QoS



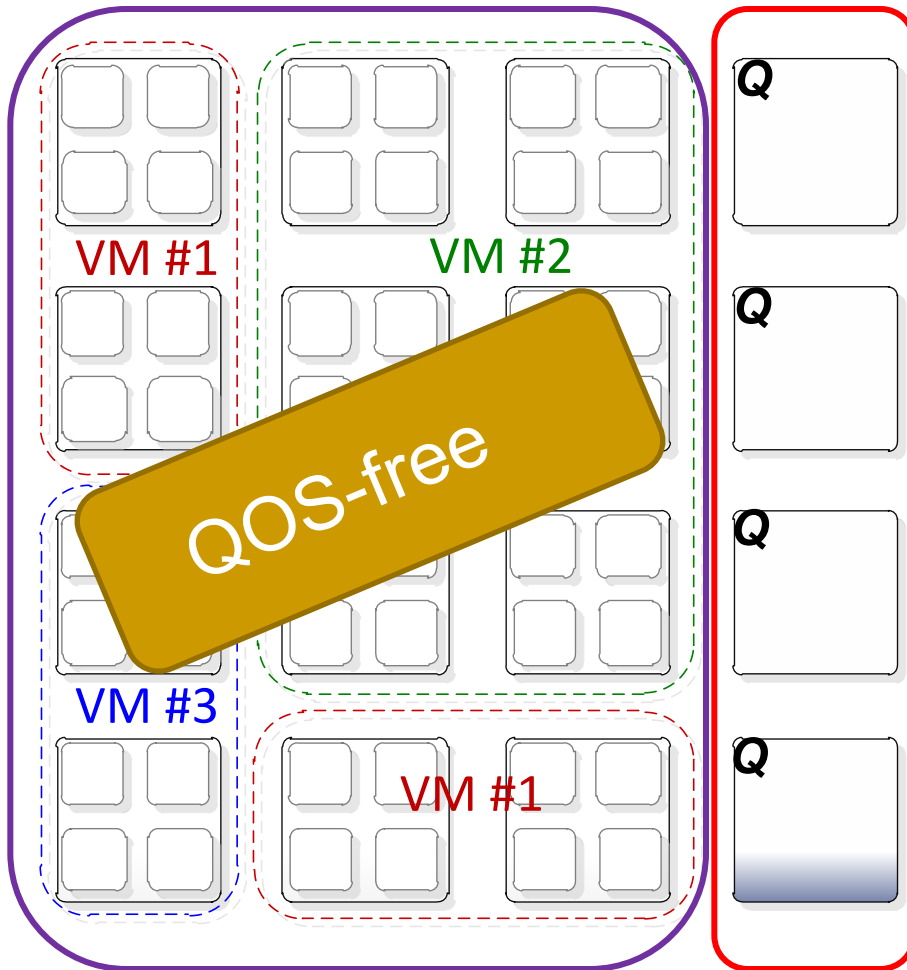
- Dedicated, QoS-enabled regions
 - Rest of die: QoS-free
- Richly-connected topology
 - Traffic isolation
- Special routing rules
 - Manage interference

Topology-Aware QoS



- Dedicated, QoS-enabled regions
 - Rest of die: QoS-free
- Richly-connected topology
 - Traffic isolation
- Special routing rules
 - Manage interference

Kilo-NOC view

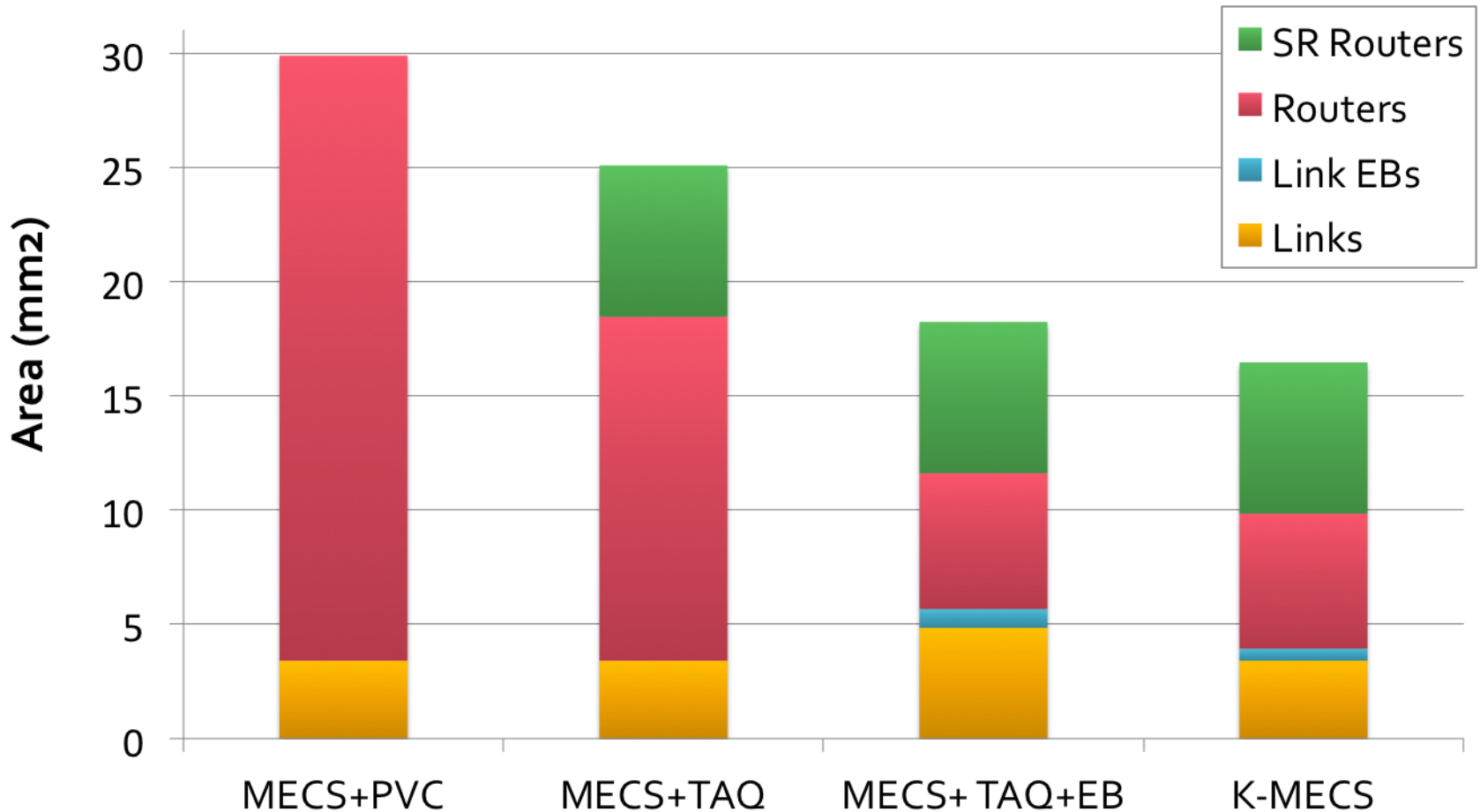


- Topology-aware QOS support
 - Limit QOS complexity to a fraction of the die
- Optimized flow control
 - Reduce buffer requirements in QOS-free regions

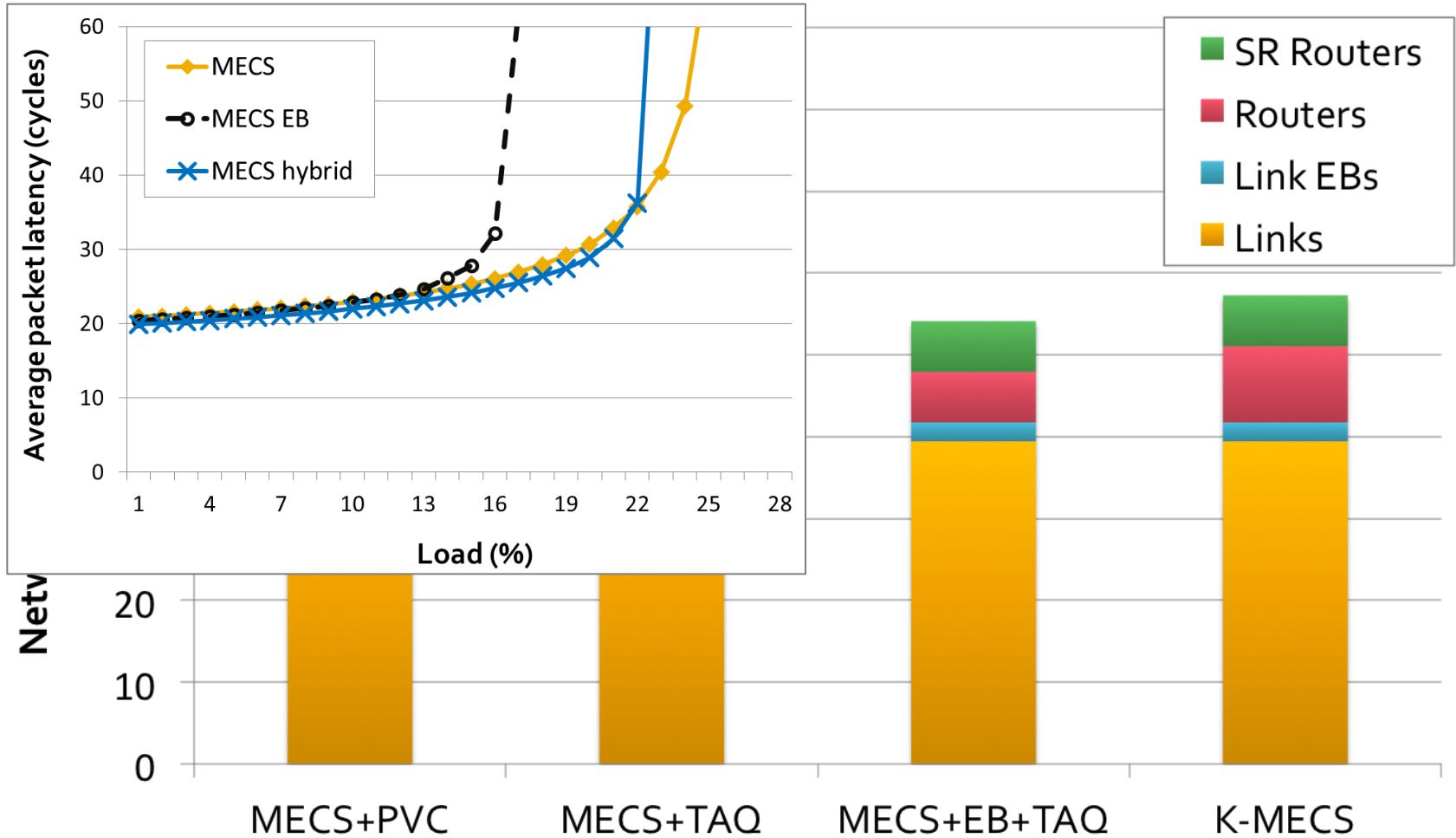
Evaluation Methodology

Parameter	Value
Technology	15 nm
Vdd	0.7 V
System	1024 tiles: 256 concentrated nodes (64 shared resources)
Networks:	
MECS+PVC	VC flow control, QOS support (PVC) at each node
MECS+TAQ	VC flow control, QOS support only in shared regions
MECS+TAQ+EB	EB flow control outside of SRs, Separate <i>Request</i> and <i>Reply</i> networks
K-MECS	Proposed organization: TAQ + hybrid flow control

Area comparison



Energy comparison



Summary

Kilo-NOC: a heterogeneous NOC architecture for kilo-node substrates

- Topology-aware QOS
 - Limits QOS support to a fraction of the die
 - Leverages low-diameter topologies
 - Improves NOC area- and energy-efficiency
 - Provides strong guarantees