

INTERCONNECTION NETWORKS

Mahdi Nazm Bojnordi

Assistant Professor

School of Computing

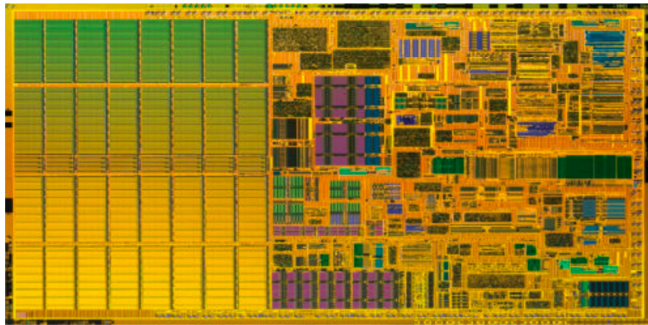
University of Utah

Overview

- Upcoming deadline
 - ▣ Feb.3rd: project group formation
 - ▣ No groups have sent me emails!
- This lecture
 - ▣ Cache interconnects
 - ▣ Basics of the interconnection networks
 - ▣ Network topologies
 - ▣ Flow control

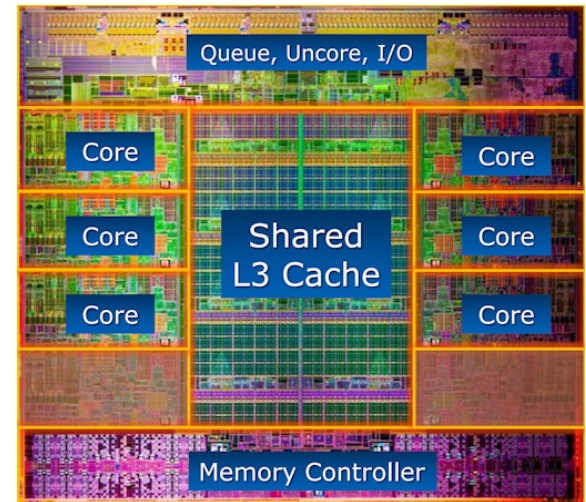
Where Interconnects Are Used?

- About 60% of the dynamic power in modern microprocessors is dissipated in on-chip interconnects



- **Analysis subject: Processor, 0.13 μm**
- **77 million transistors, die size of 88 mm^2**
- **Data sources (AF, Capacitance, Length)**
- **Excluded: L2 cache, global clock, analog units**

[Magen'04]

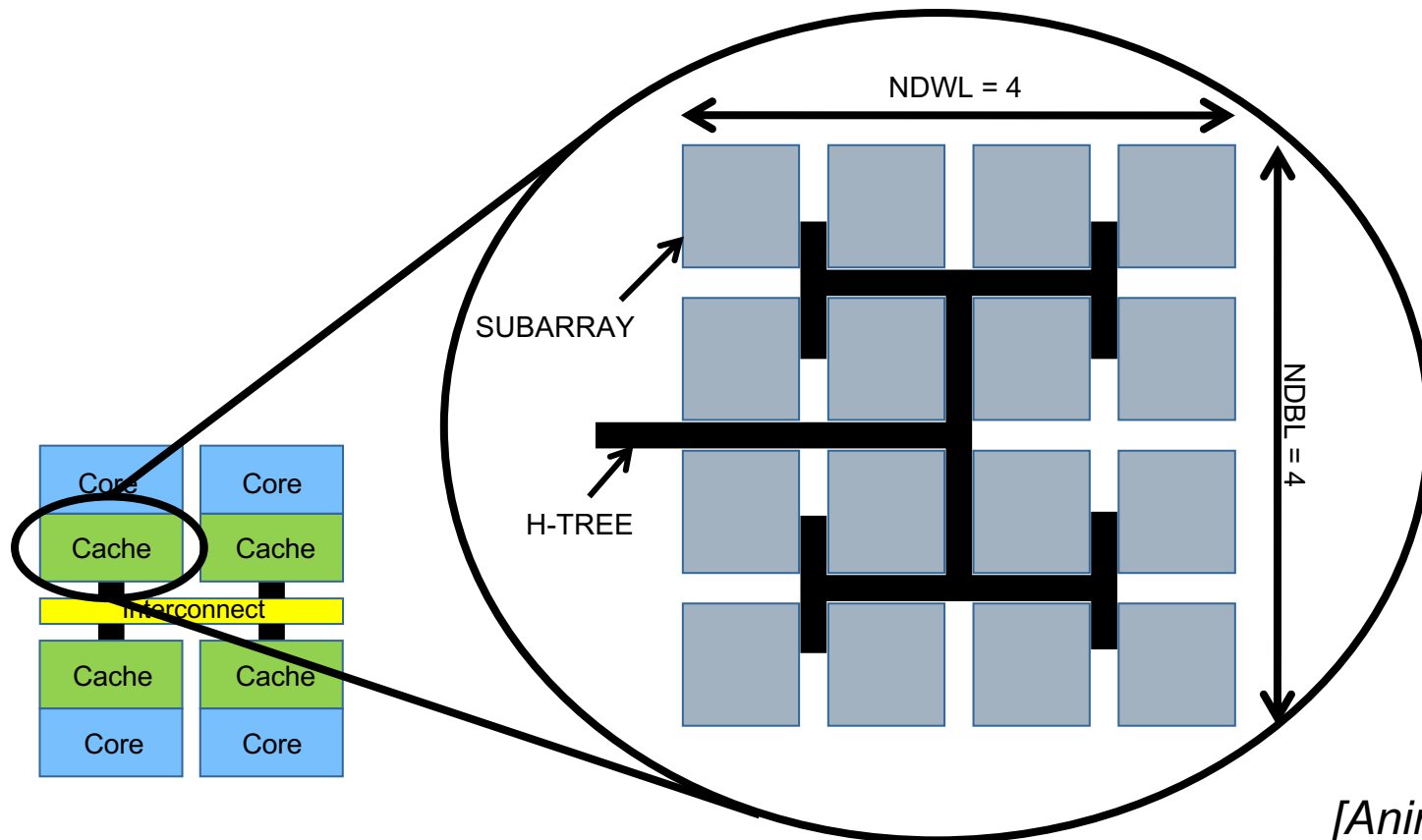


- **Six processor cores**
 - **8MB Last level cache**
- [Intel Core i7]*

Cache Interconnect Optimizations

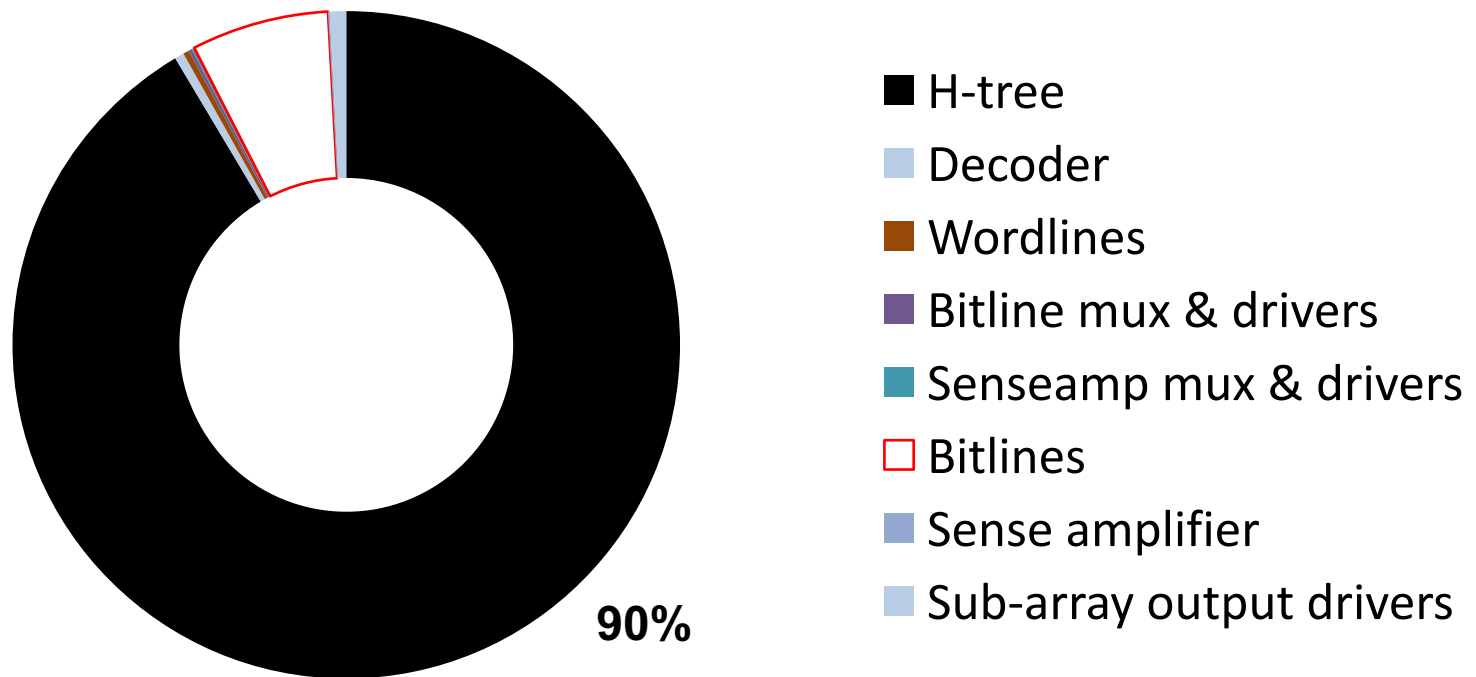
Large Cache Organization

- Fewer subarrays gives increased area efficiency, but larger delay due to longer wordlines/bitlines



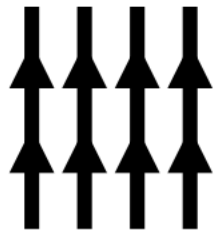
Large Cache Energy Consumption

- H-tree is clearly the dominant component of energy consumption

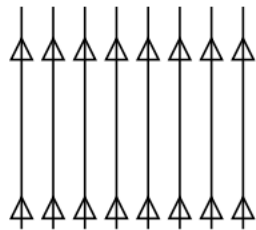


Heterogeneous Interconnects

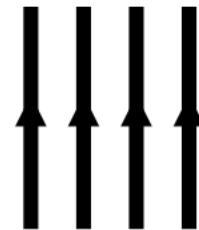
- A global wire management at the microarchitecture level
- A heterogeneous interconnect that is comprised of wires with varying **latency**, **bandwidth**, and **energy** characteristics



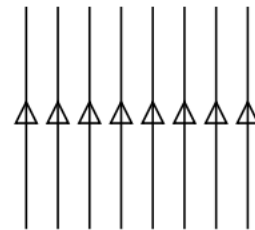
Delay Optimized



Bandwidth Optimized



Power Optimized



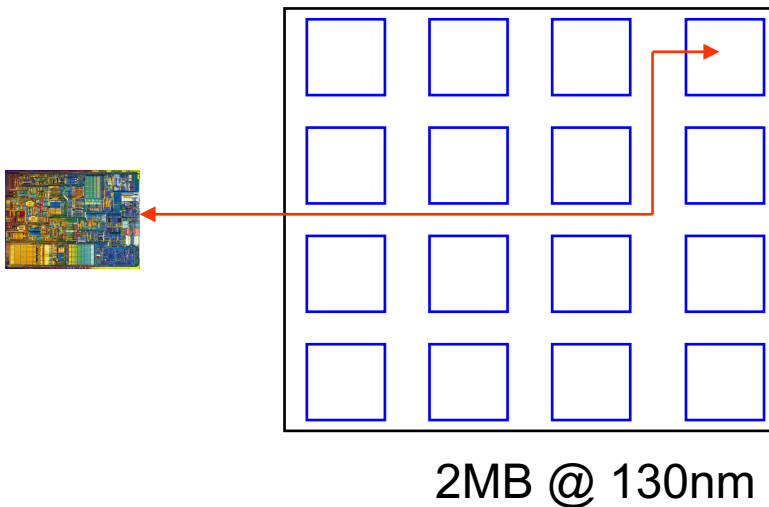
Power and Bandwidth Optimized

Heterogeneous Interconnects

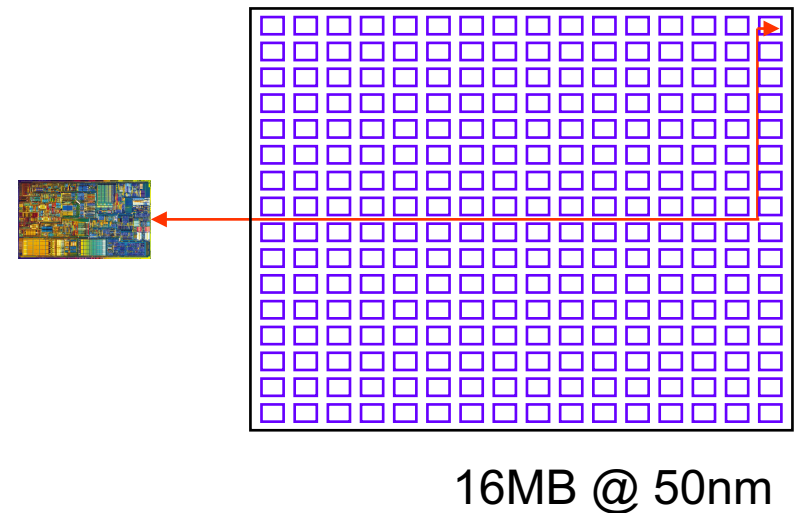
- Better energy-efficiency for a dynamically scheduled partitioned architecture
 - ▣ ED^2 is reduced by 11%
- A low-latency low-bandwidth network can be effectively used to hide wire latencies and improve performance
- A high-bandwidth low-energy network and an instruction assignment heuristic are effective at reducing contention cycles and total processor energy.

Non-Uniform Cache Architecture

- NUCA optimizes energy and time based on the proximity of the cache blocks to the cache controller.



Bank Access time = 3 cycles
Interconnect delay = 8 cycles

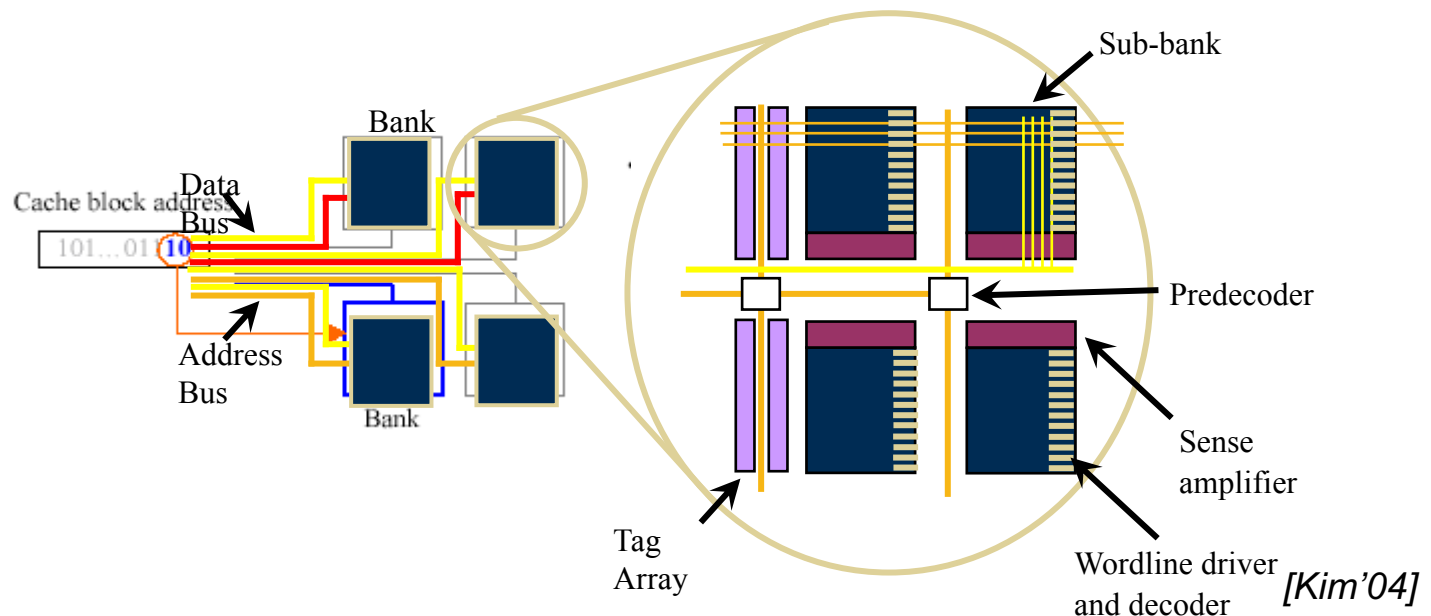


Bank Access time = 3 cycles
Interconnect delay = 44 cycles

Non-Uniform Cache Architecture

□ S-NUCA-1

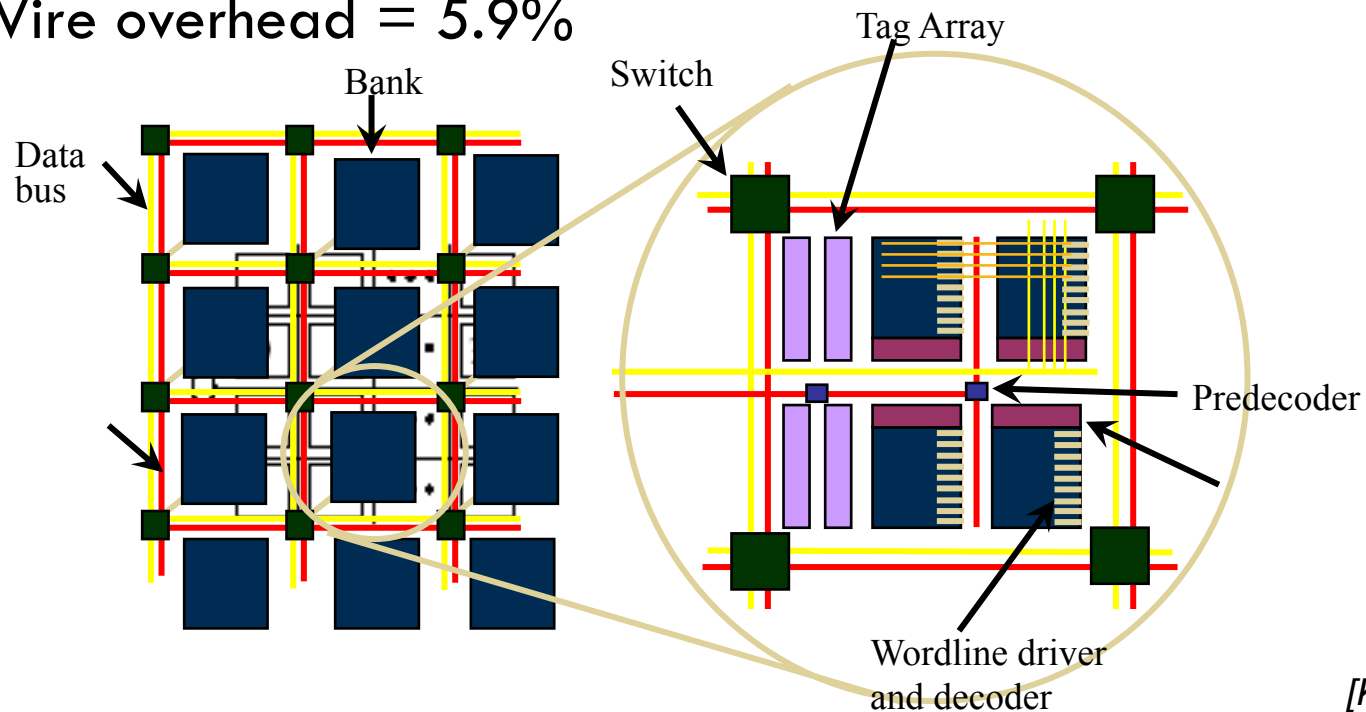
- Use private per-bank channel
- Each bank has its distinct access latency
- Statically decide data location for its given address
- Average access latency = 34.2 cycles
- Wire overhead = 20.9% → an issue



Non-Uniform Cache Architecture

□ S-NUCA-2

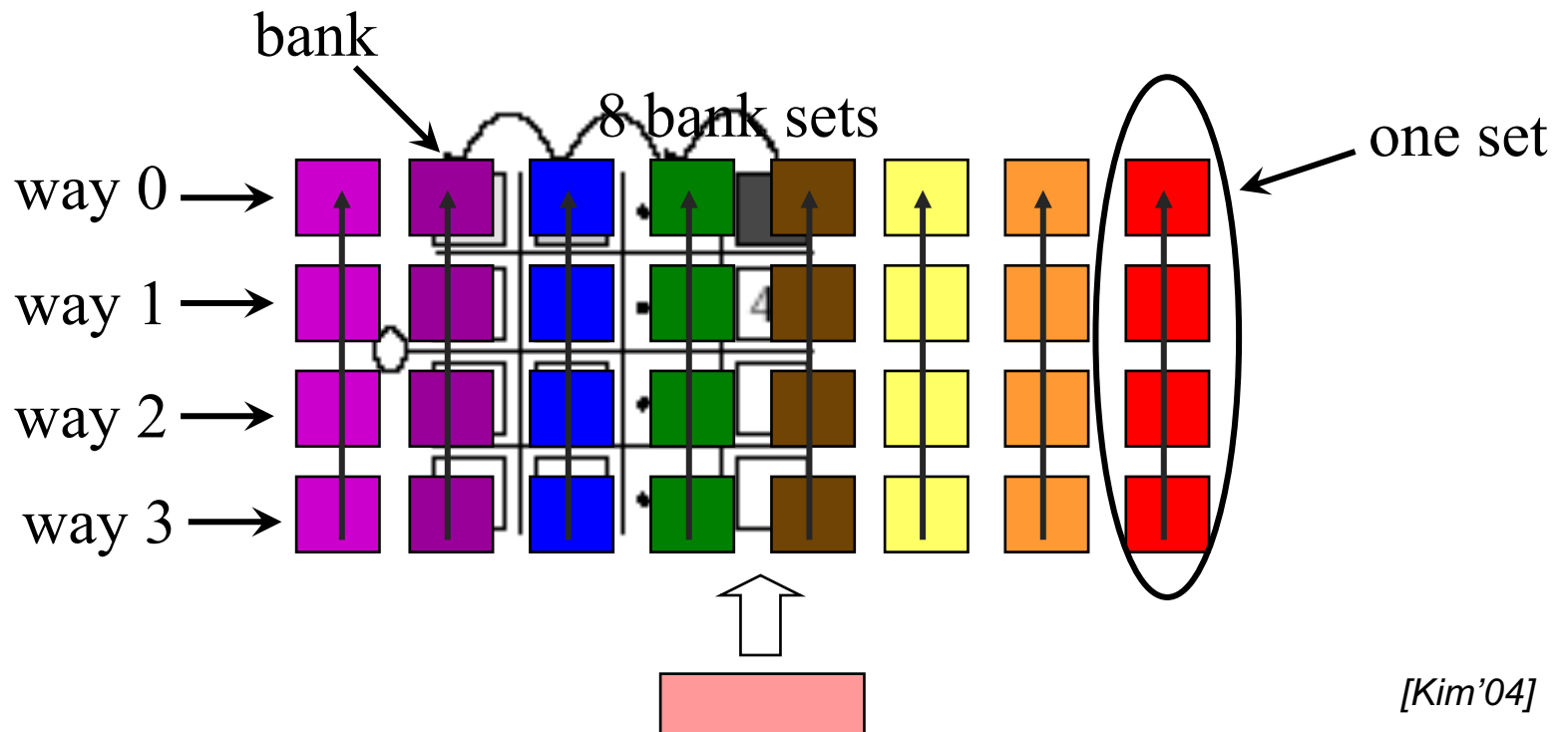
- Use a 2D switched network to alleviate wire area overhead
- Average access latency = 24.2 cycles
- Wire overhead = 5.9%



Non-Uniform Cache Architecture

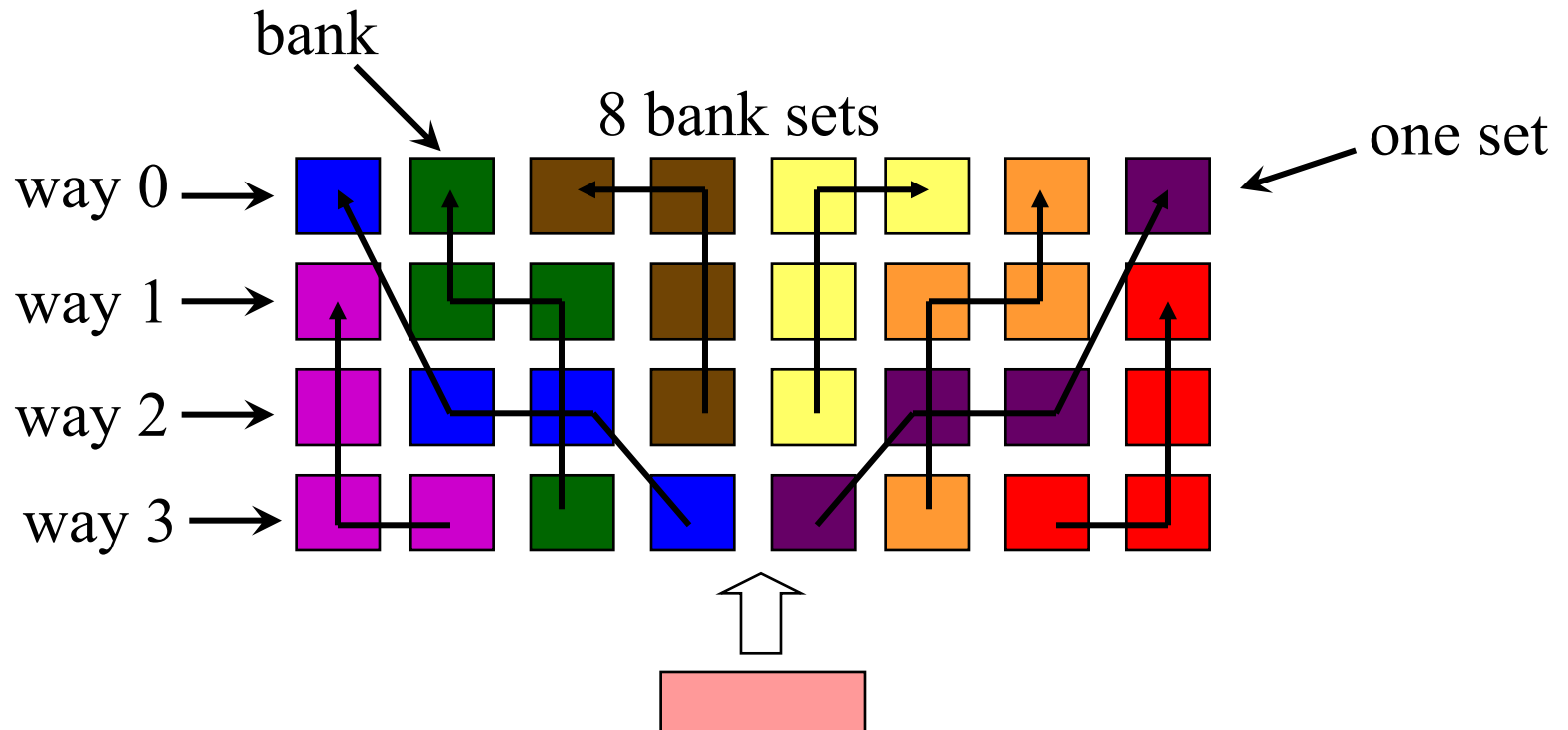
- Dynamic NUCA

- ▣ Data can dynamically migrate
- ▣ Move frequently used cache lines closer to CPU



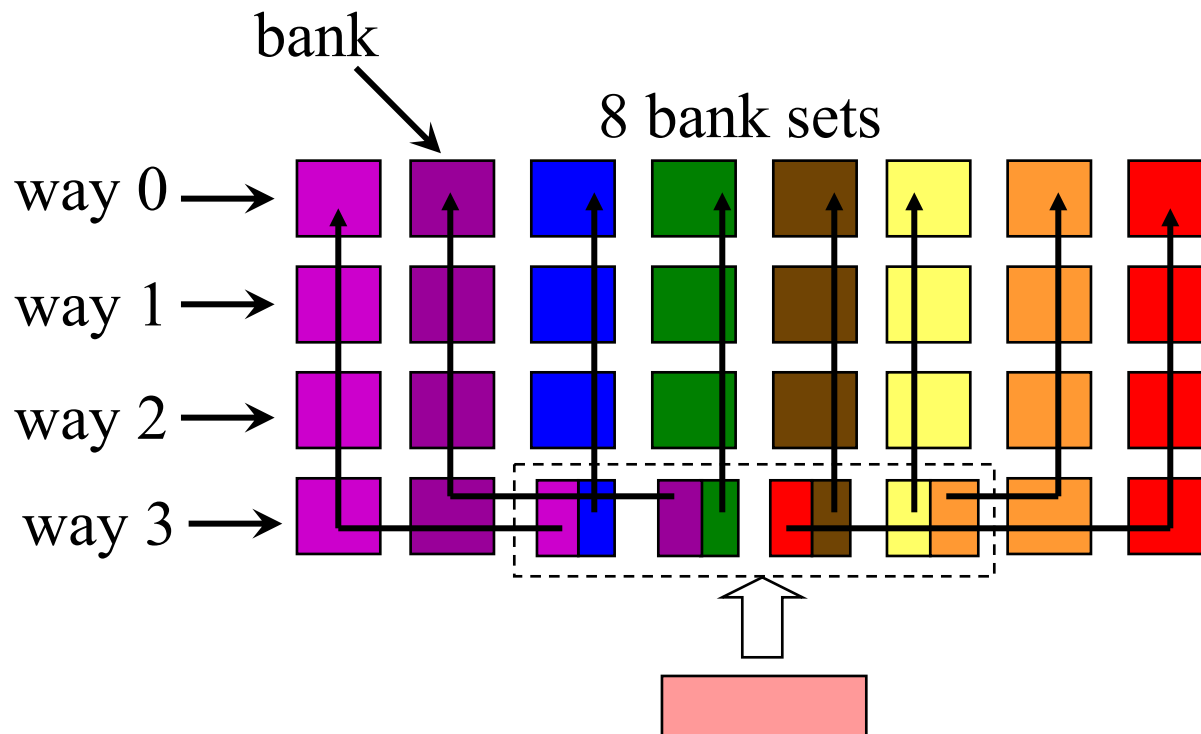
Non-Uniform Cache Architecture

- Fair mapping
 - ▣ Average access time across all bank sets are equal



Non-Uniform Cache Architecture

- Shared mapping
 - ▣ Sharing the closet banks for farther banks

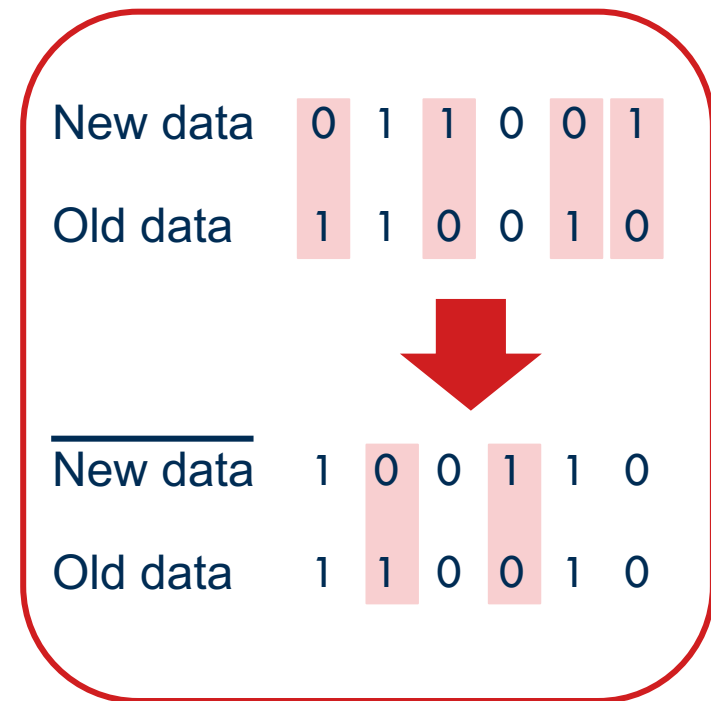


Encoding Based Optimizations

Cache Interconnect Optimizations

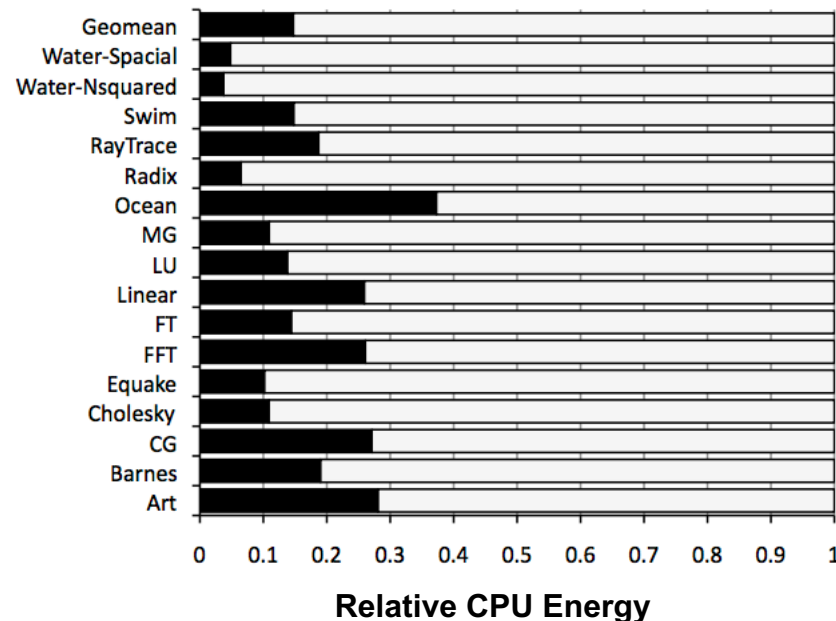
- Bus invert coding transfers either the data or its complement to minimize the number of bit flips on the bus.

$$P_{\text{switching}} = \alpha C V_{DD}^2 f$$



Time-Based Data Transfer

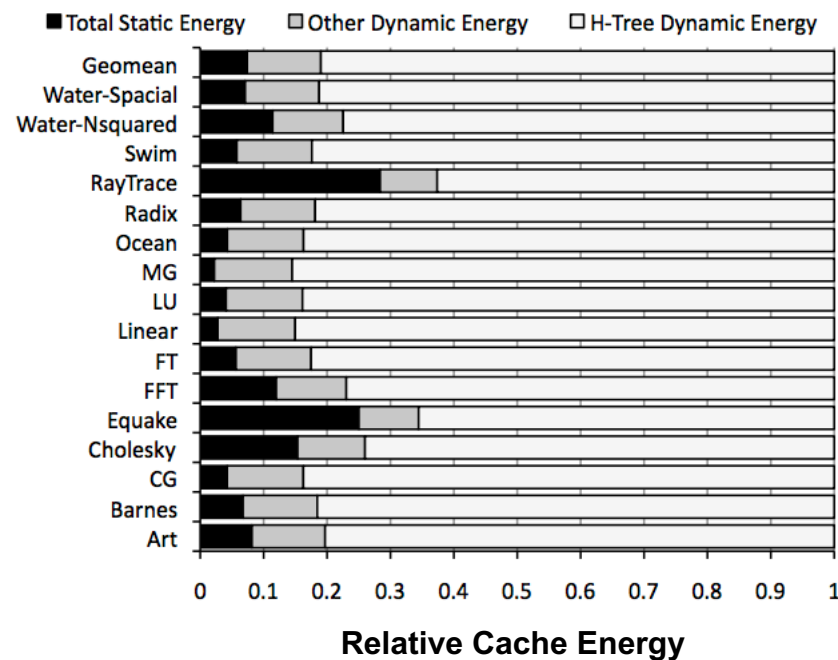
- The percentage of processor energy expended on an 8MB cache when running a set of parallel applications on a Sun Niagara-like multicore processor



[Bojnordi'13]

Time-Based Data Transfer

- Communication over the long, capacitive H-tree interconnect is the dominant source of energy consumption (80% on average) in the L2 cache

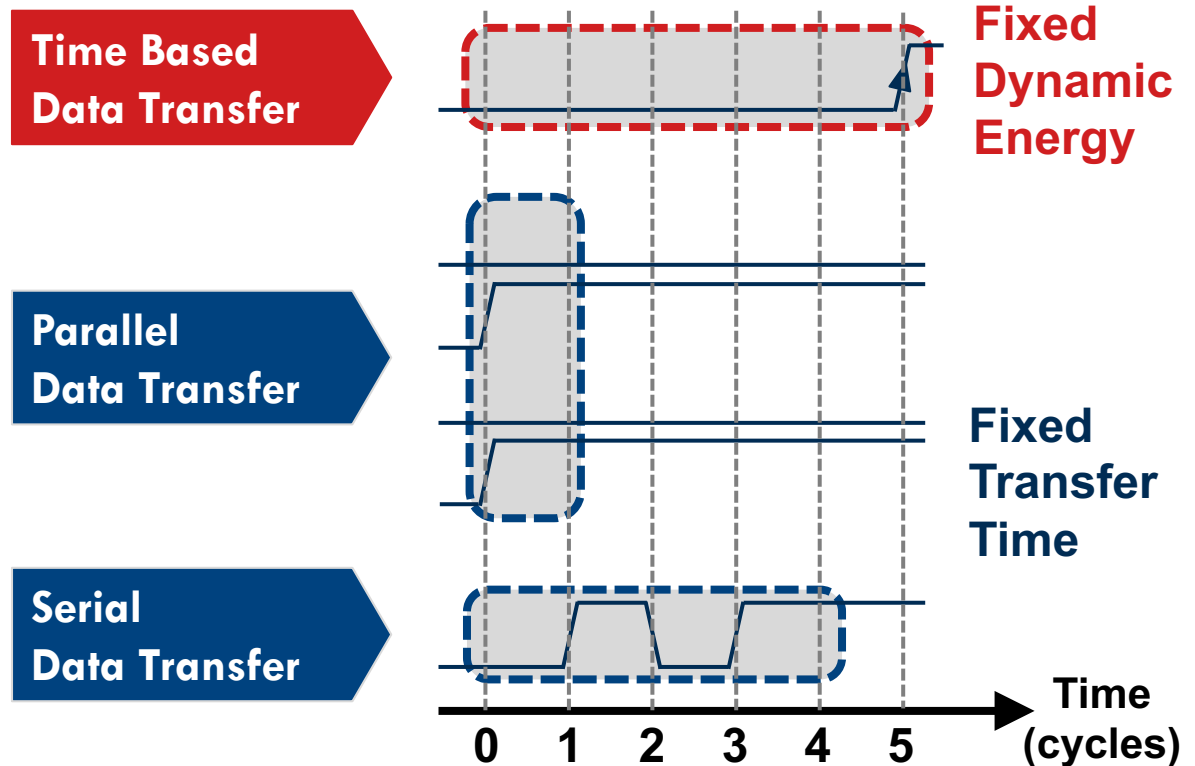


[Bojnordi'13]

Time-Based Data Transfer

Key idea: represent information by the number of clock cycles between two consecutive pulses to reduce interconnect activity factor.

Example: transmitting the value 5



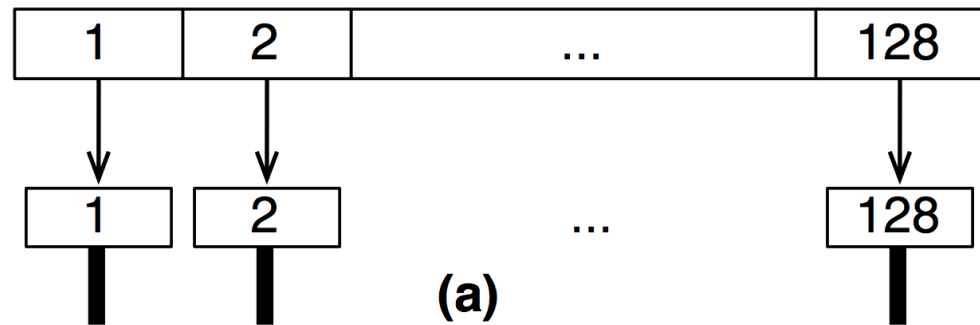
[Bojnordi'13]

Time-Based Data Transfer

- Cache blocks are partitioned into small, contiguous **chunks**.

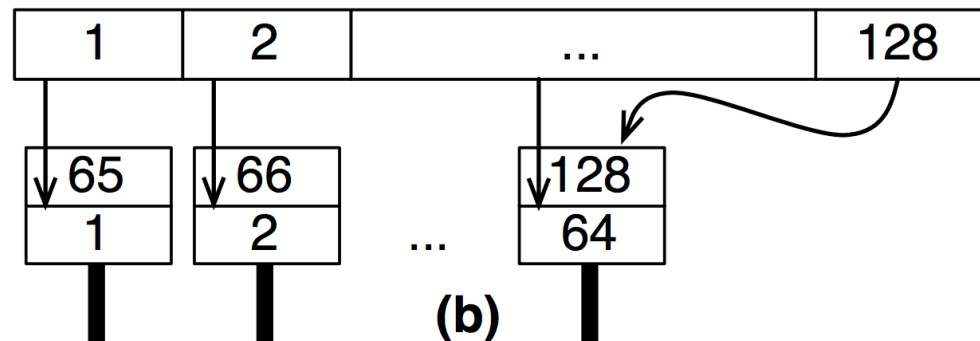
Cache Block
Partitioned into Chunks

Communication Wires
and FIFO Queues



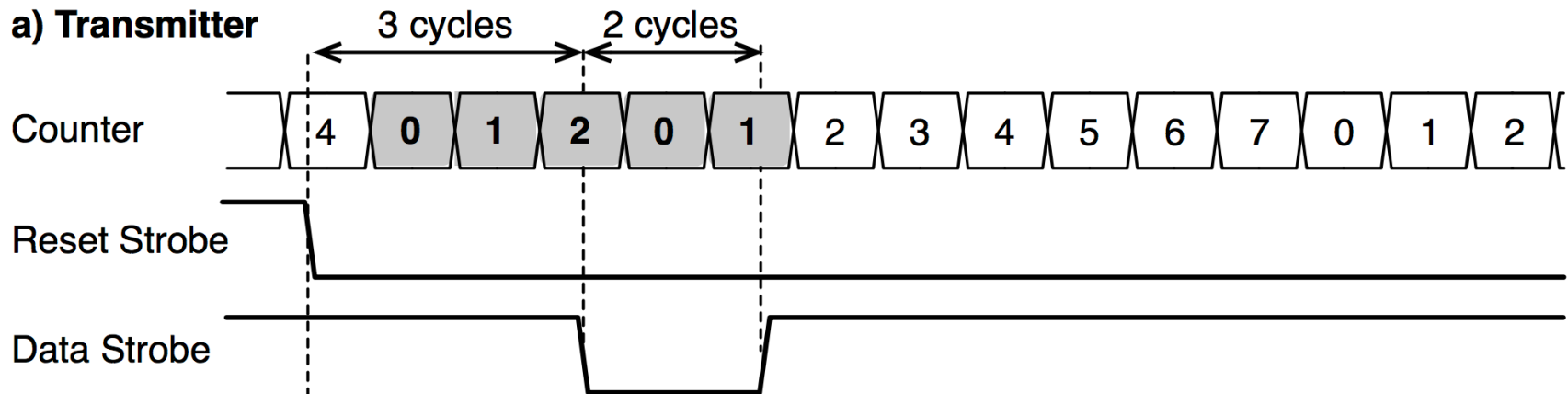
Cache Block
Partitioned into Chunks

Communication Wires
and FIFO Queues



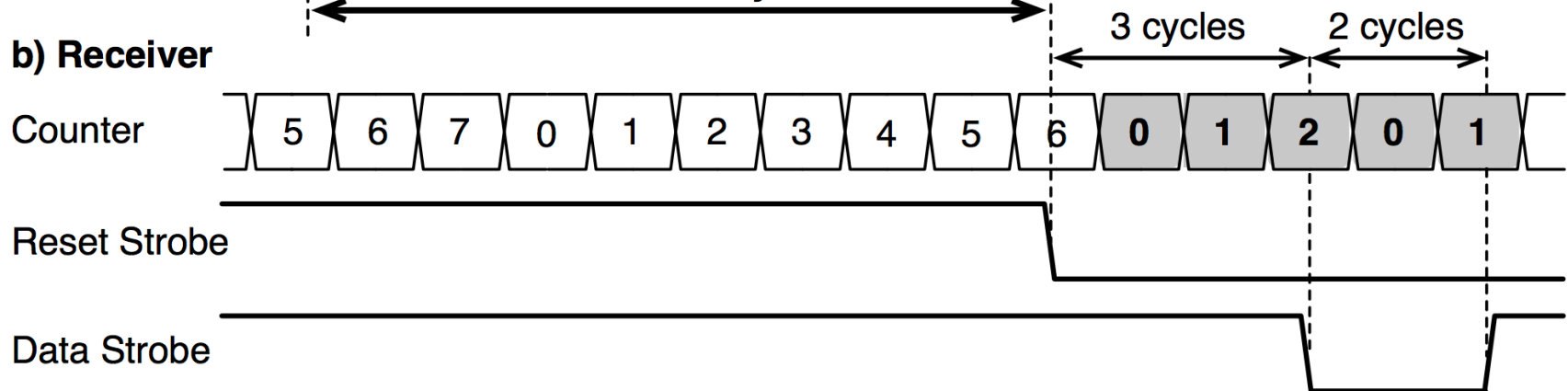
Time-Based Data Transfer

a) Transmitter



Wire Delay

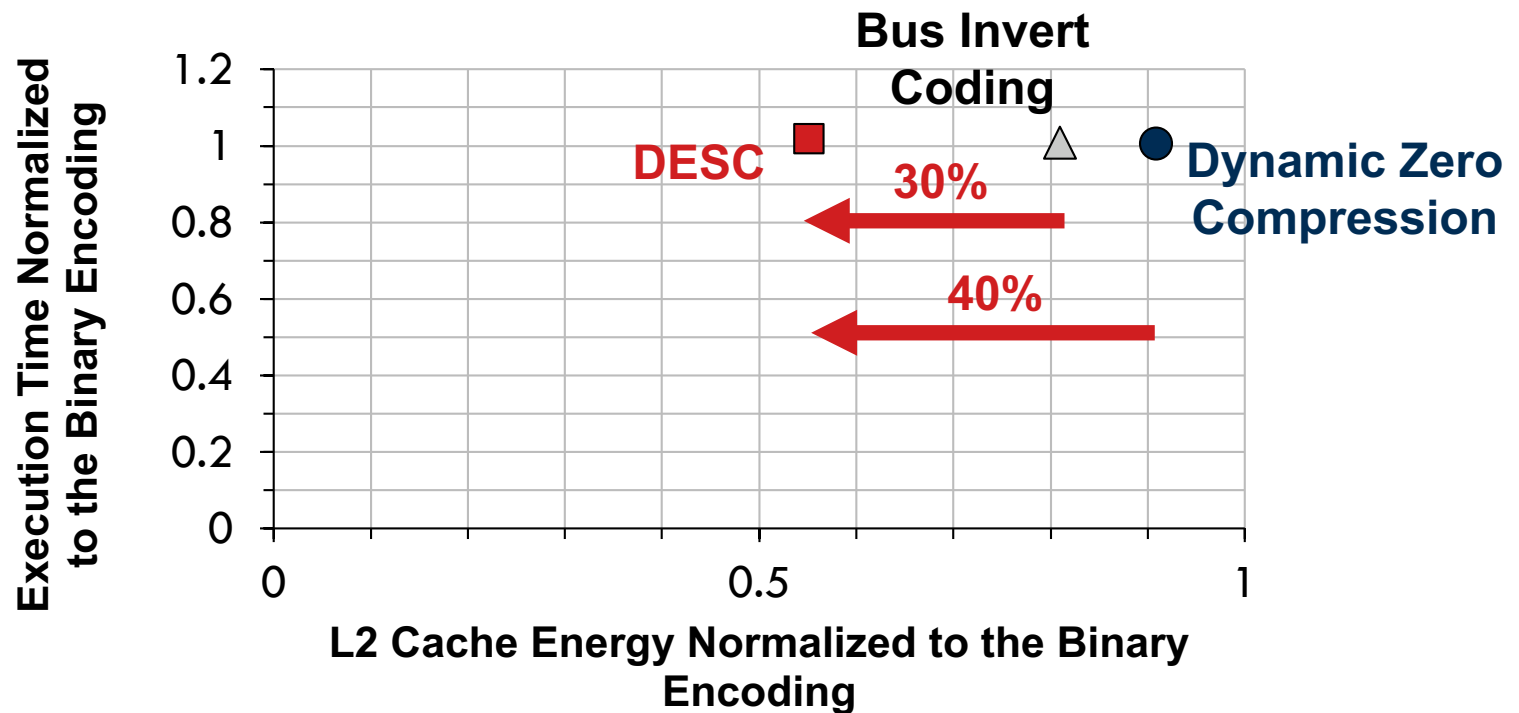
b) Receiver



[Bojnordi'13]

Time-Based Data Transfer

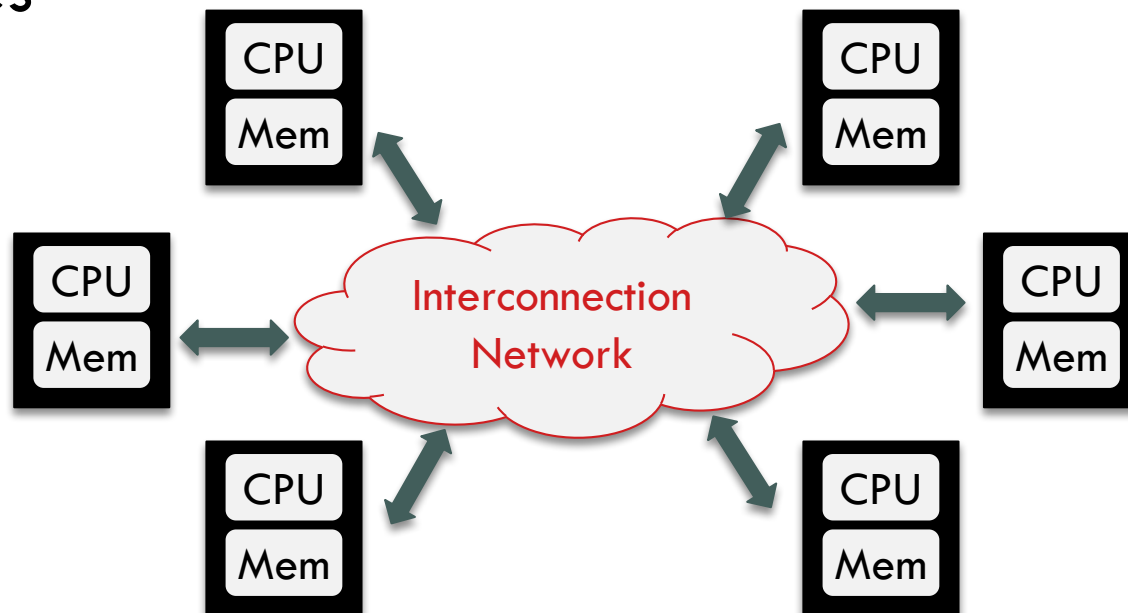
- L2 cache energy is reduced by 1.8x at the cost of less than 2% increase in the execution time.



Interconnection Networks

Interconnection Networks

- Goal: transfer maximum amount of information with the minimum time and power
- Connects processors, memories, caches, and I/O devices



Types of Interconnection Networks

- Four domains based on number and proximity of devices
 - ▣ On-chip networks (OCN or NOC)
 - Microarchitectural elements: cores, caches, reg. files, etc.
 - ▣ System/storage area networks (SAN)
 - Computer subsystems: storage, processor, IO device, etc.
 - ▣ Local area networks (LAN)
 - Autonomous computer systems: desktop computers etc.
 - ▣ Wide area networks (WAN)
 - Interconnected computers distributed across the globe

Basics of Interconnection Networks

- Network topology
 - ▣ How to wire switches and nodes in the network
- Routing algorithm
 - ▣ How to transfer a message from source to destination
- Flow control
 - ▣ How to control the flow messages within the network