

Issues in VLSI interconnect

Avinoam Kolodny
Technion – Israel Institute of Technology

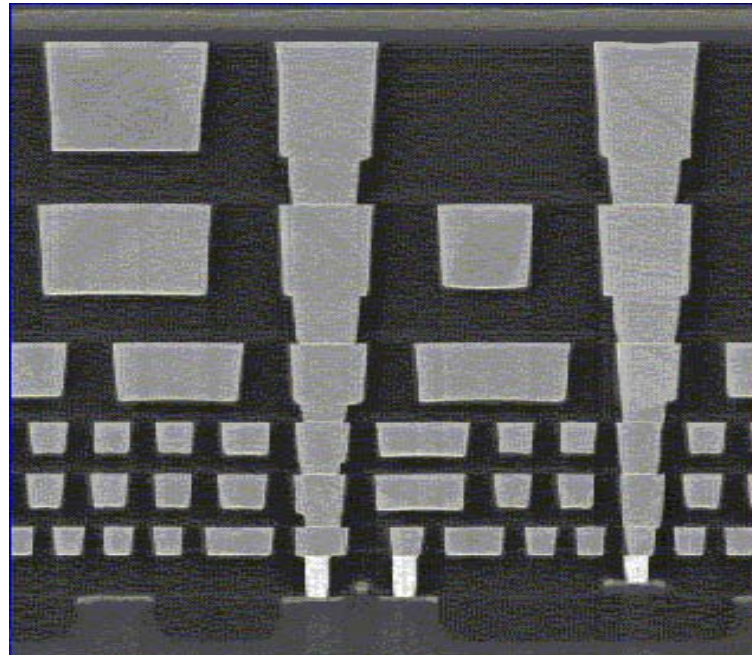
February 2004



In nanometer CMOS

Interconnect becomes critical

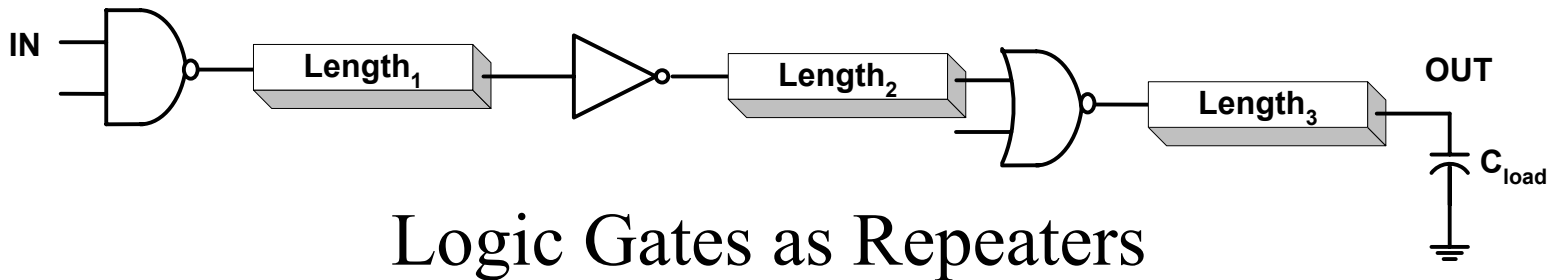
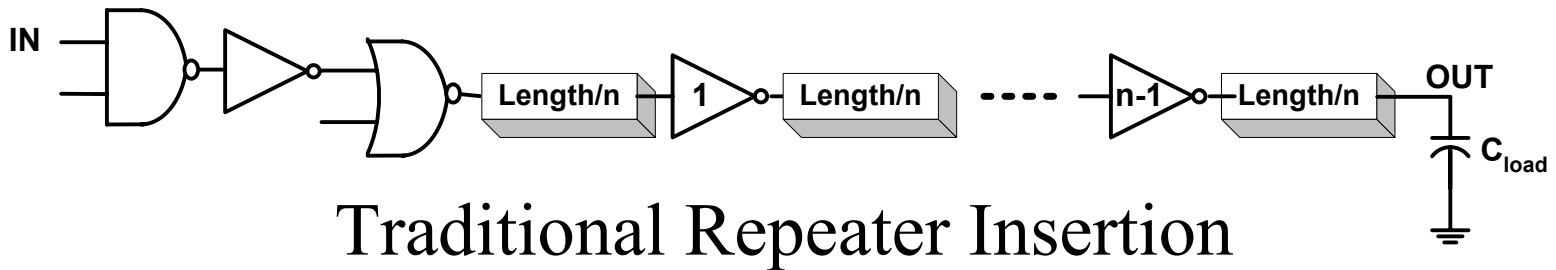
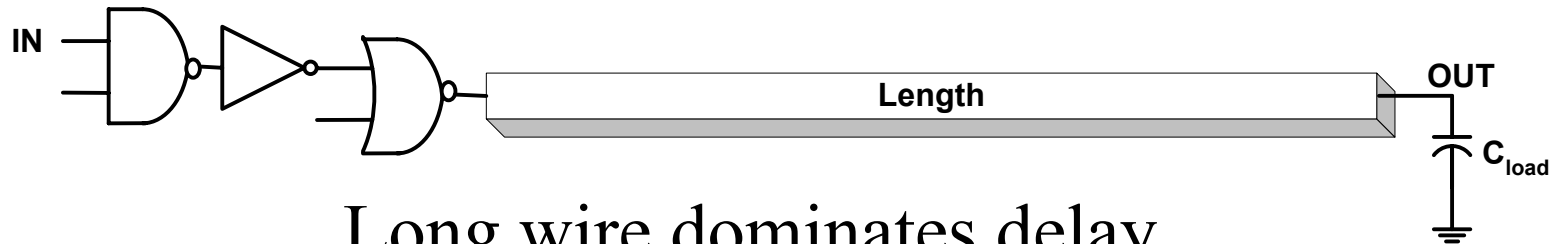
- Interconnect can dominate
 - Timing
 - Power
 - Noise
 - Layout density
 - Design effort



0.13 μm cross-section, source - Intel



Logic Gates as Repeaters (LGR)



- + *Timing improvement similar to traditional Repeater Insertion*
- + *No additional devices – area and power efficient*

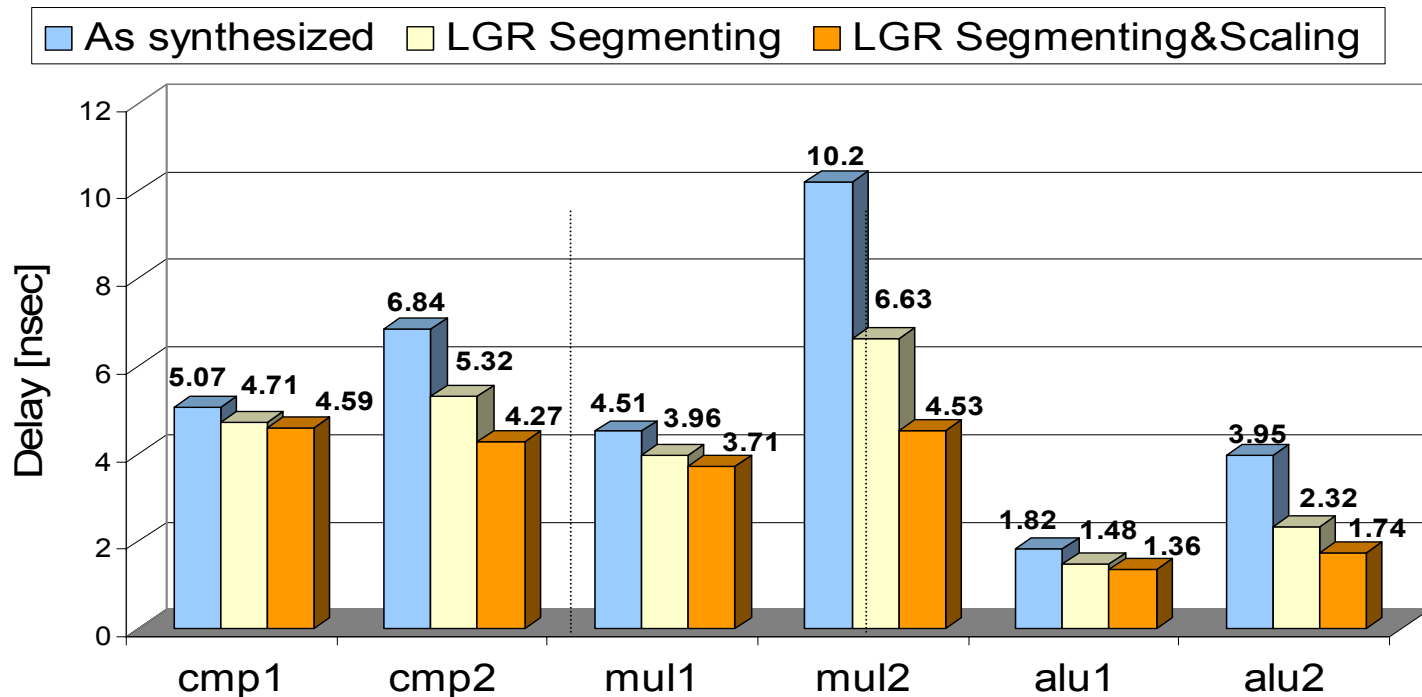


LGR in Physical Synthesis

Stage 1: Test circuits Implementation using commercial tools

Stage 2: Post-Layout extraction of critical path

Stage 3: LGR optimization of critical path

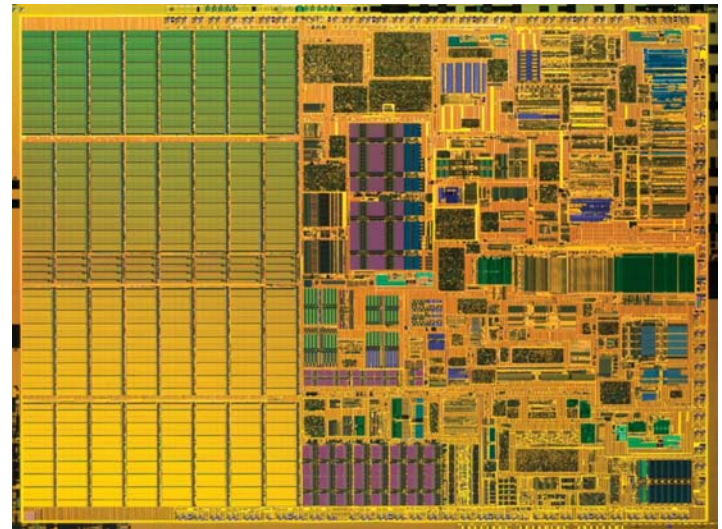
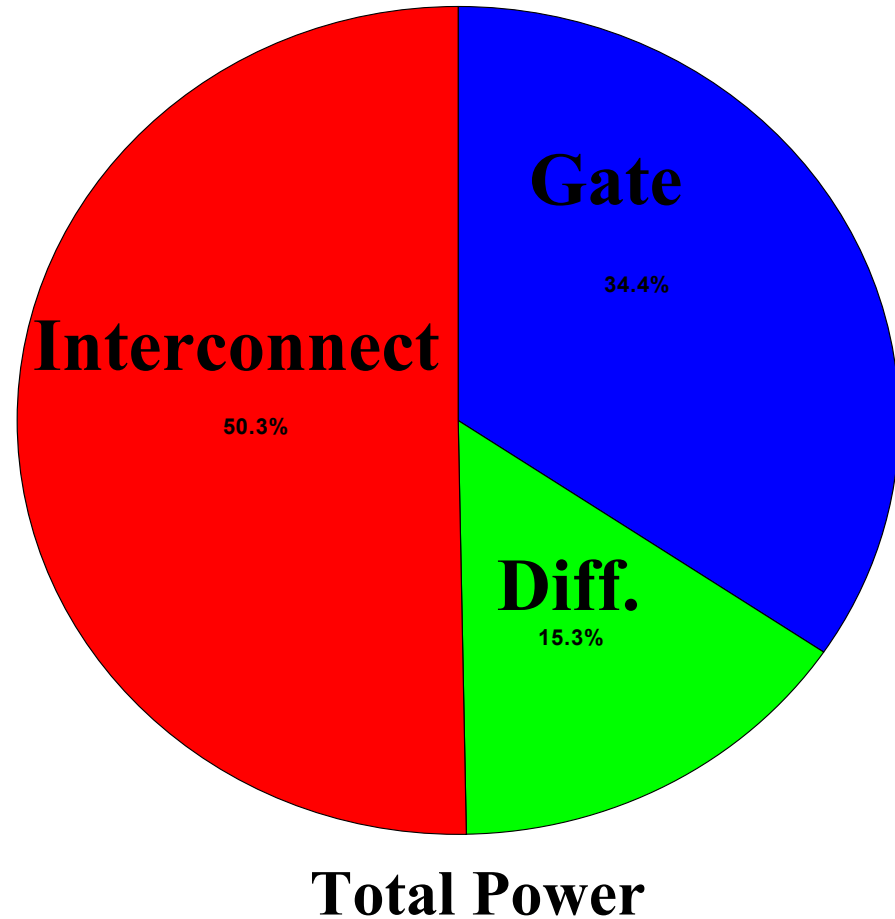


* A. Morgenshtein, M. Moreinis, I. Wagner and A. Kolodny,
“Logic gates as Repeaters”, IFIP VLSI-SoC , December 2003.



Dynamic Power Dissipation in a microprocessor

- 0.13 μm CMOS
- 77 million transistors
- Die size: 88 mm²

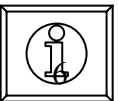
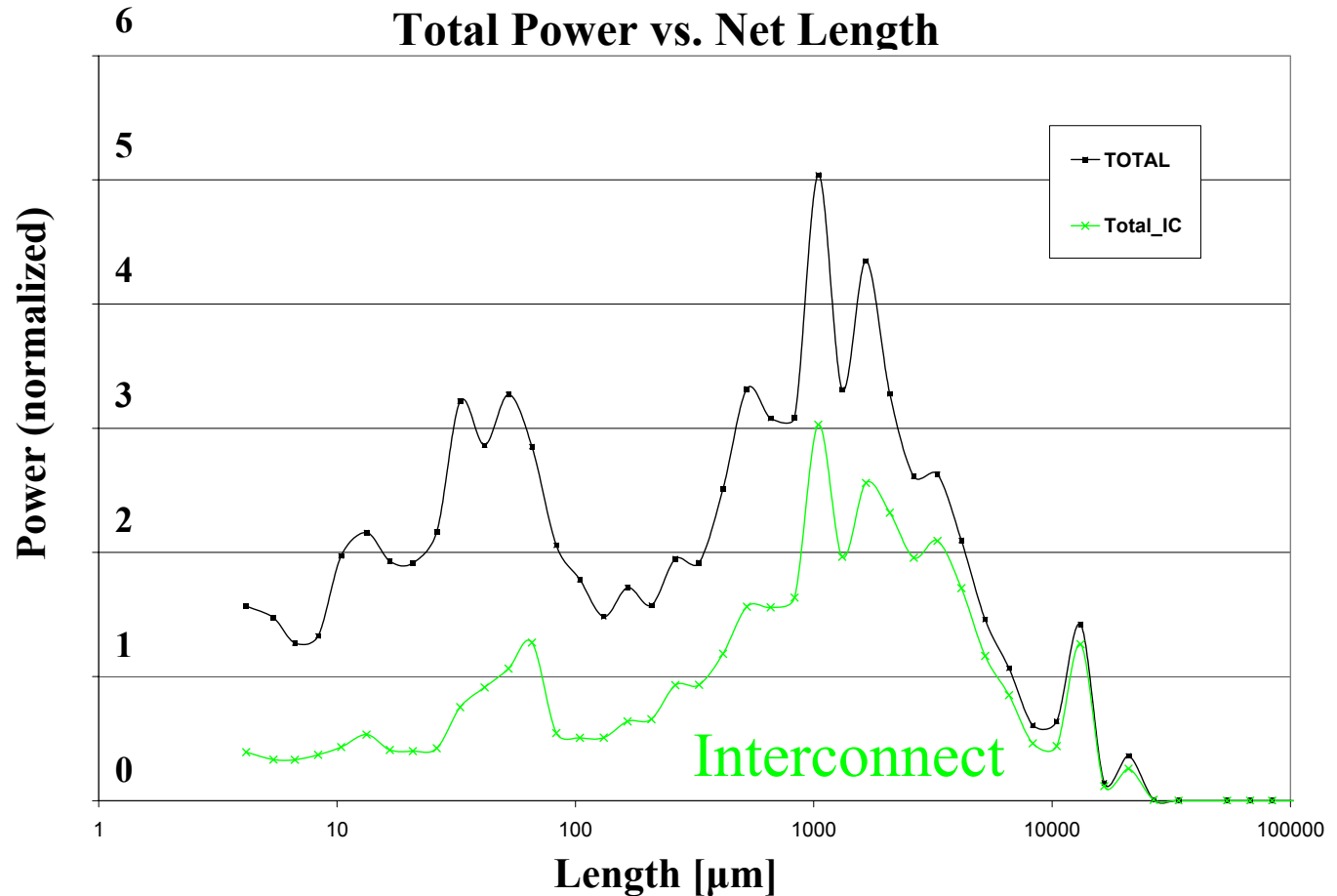


* N. Magen, A. Kolodny, U. Weiser and N. Shamir, “Interconnect-Related Energy dissipation in a Low-Power Microprocessor”, Proc. International System Level Interconnect Prediction workshop (SLIP) , February 2004.



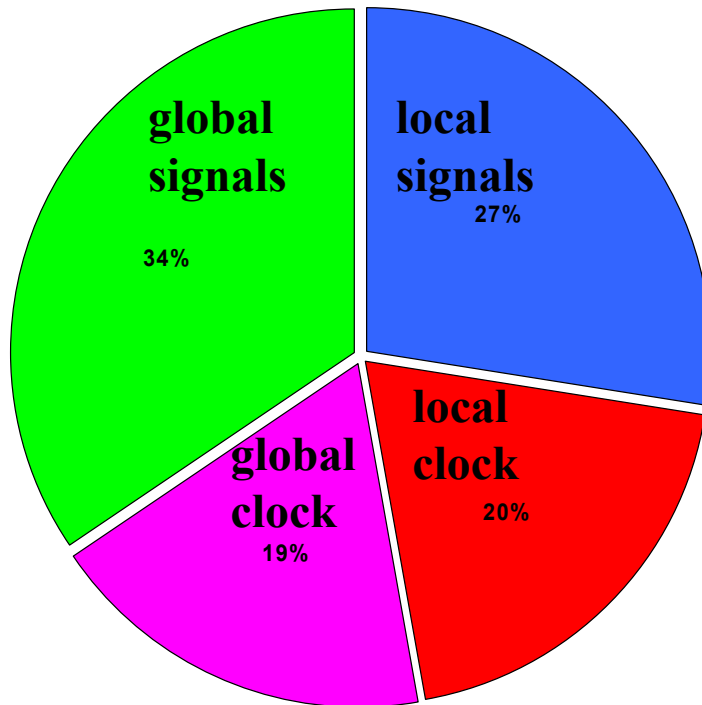
Power Distribution by Wirelength

- Total Dynamic Power
- Global clock – not included
- Local signals nets = 66%
- Global signals nets = 34%



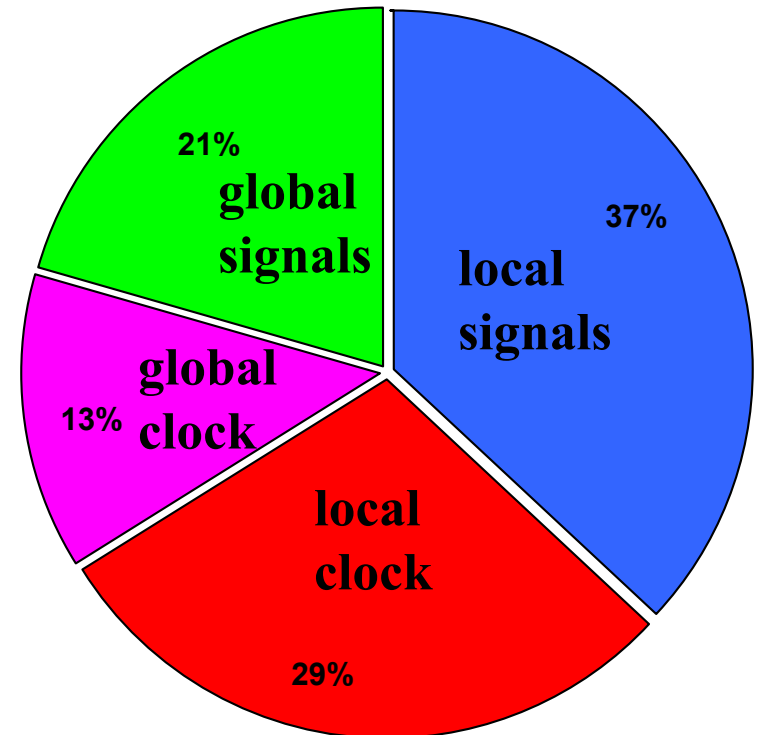
Power in local and global signals

- Interconnect consumes 50% of total dynamic power
- Clock power ~40% (of Interconnect and total)
- 90% of power consumed by 10% of nets
- Interconnect design is NOT power-aware !



Interconnect power

(Interconnect only)

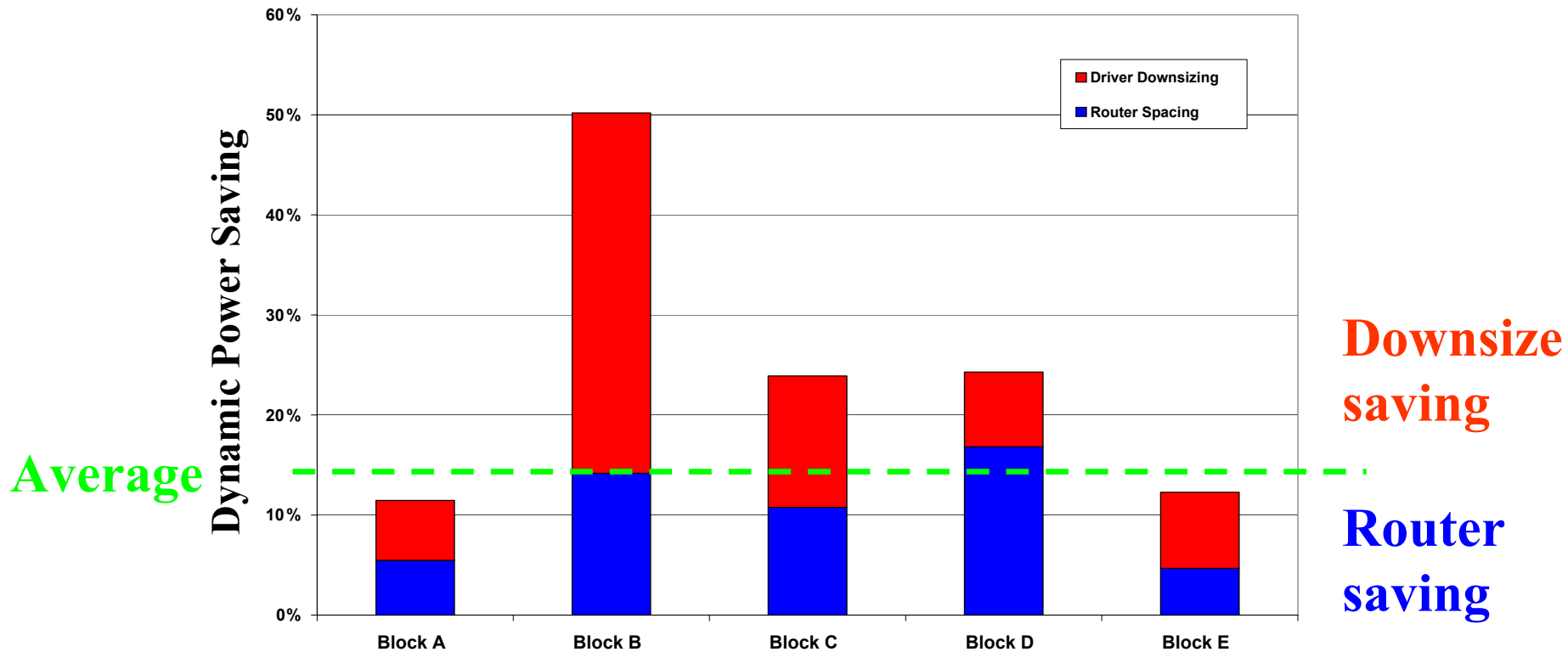


Total dynamic power

(Gate, Diffusion and Interconnect)



Savings by power-aware router

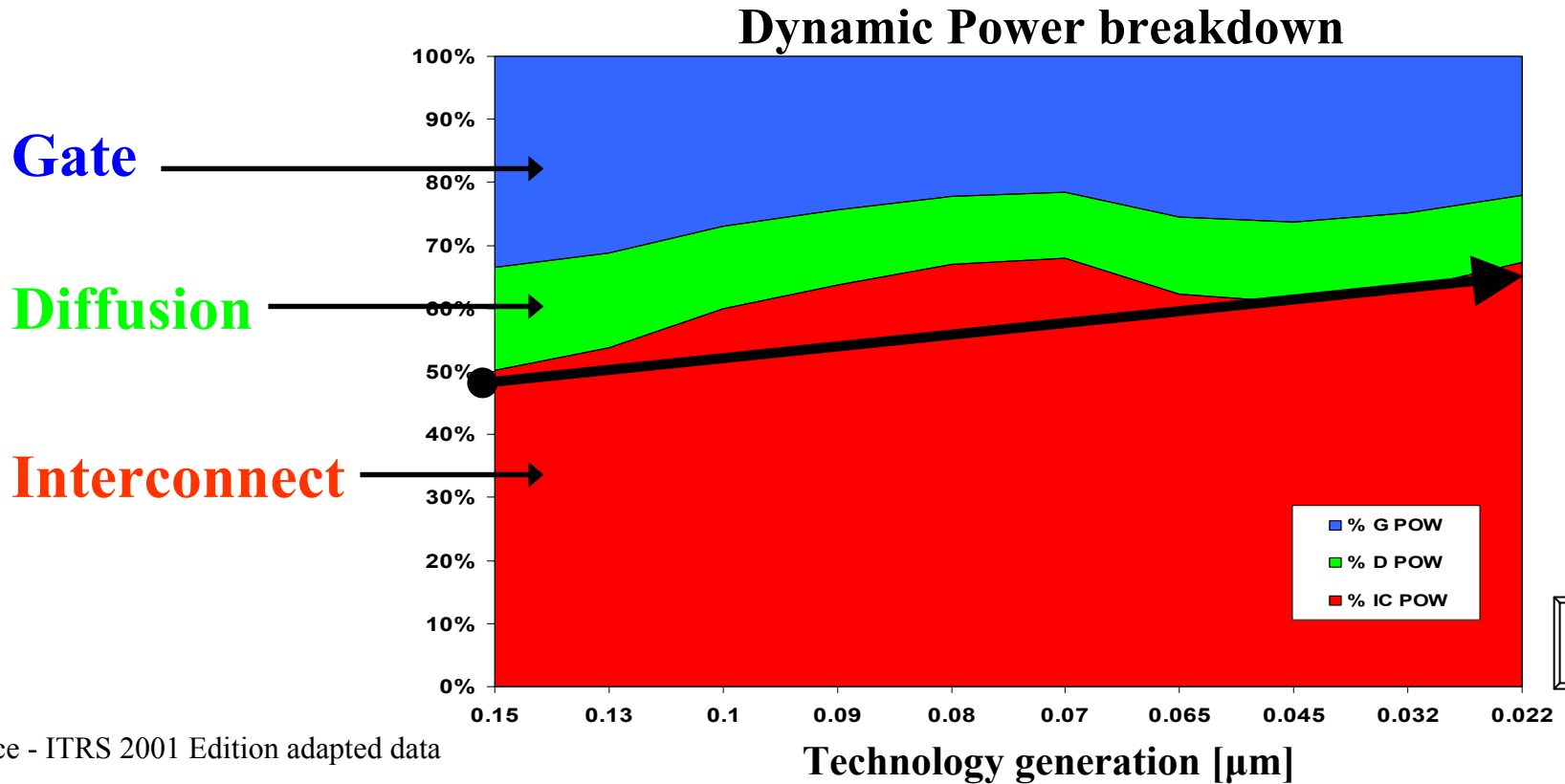


Average saving results: 14.3% for ASIC blocks ¹

1 - Estimated based on clock interconnect power



Future of Interconnect Power



Source - ITRS 2001 Edition adapted data

Interconnect power grows to 65%-80% within 5 years !





Networks-on-Chip (NoC)

Faculty: Israel Cidon, Ran Ginosar, Avinoam Kolody

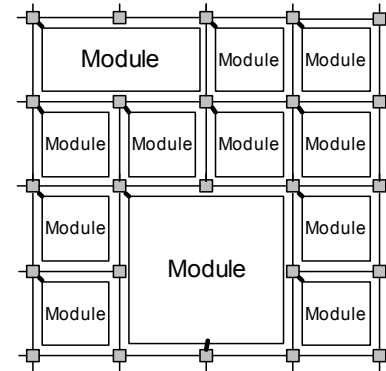
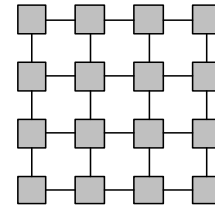
Students: Evgeny Bolotin, Walter Isaskhar

Arkadiy Morgenshtein, Reuven Dobkin



Why SoC needs NoC?

- Scalability
 - Traditional chip interconnections are not scalable
- Optimized NoC links
 - Interconnect dominates delay, noise and power in nanometer technologies
- Synchronization
 - Delay of global wire longer than a clock cycle
 - GALS Systems
- Modularity and reuse
 - SoC design productivity

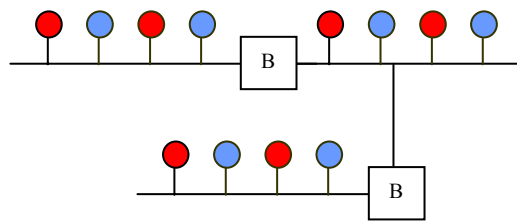


Goal: Cost-effective network that satisfies on-chip QoS requirements

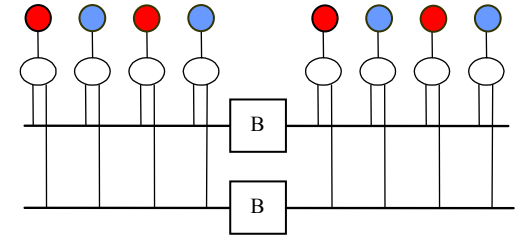


Traditional Solution-Bus

Segmented Bus



Multi-Level Segmented Bus



Original bus features:

- One transaction at a time
- Central Arbiter
- Limited bandwidth
- Synchronous
- Low cost **Is it still?**

New features:

- Versatile bus architectures
- Pipelining capability
- Burst transfer
- Split transactions
- Transaction preemption and resume
- Transaction reordering...



NoC Cost Scalability vs. Alternative Solutions

<i>Arch</i>	<i>Total Area</i>	<i>Power Dissipation</i>	<i>Operating Frequency</i>
<i>NS-Bus</i>	$O(n^3 \sqrt{n})$	$O(n\sqrt{n})$	$O\left(\frac{1}{n^2}\right)$
<i>S-Bus</i>	$O(n^2 \sqrt{n})$	$O(n\sqrt{n})$	$O\left(\frac{1}{n}\right)$
<i>NoC</i>	$O(n)$	$O(n)$	$O(1)$
<i>PTP</i>	$O(n^2 \sqrt{n})$	$O(n\sqrt{n})$	$O\left(\frac{1}{n}\right)$

* E. Bolotin, I. Cidon, R. Ginosar and A. Kolodny, “Cost Considerations in Network-on-Chip”, INTEGRATION – the VLSI journal, 2003)



Solution – Network on Chip

Networks are preferred over buses:

- Scalability
- Higher bandwidth
- Concurrency, effective spatial reuse of resources
- Higher levels of abstraction
- Modularity - Design Productivity Improvement



Network on Chip Requirements

- Different QoS must be supported
 - Bandwidth
 - Latency
- Distributed deadlock free routing
- Distributed congestion/flow control
- Low VLSI Cost



QNoC: QoS NoC

Define Service Levels (SLs):

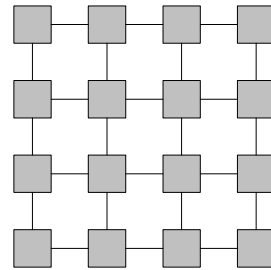
- *Signaling*
 - *Real-Time*
 - *Read/Write (RD/WR)*
 - *Block-Transfer*
- ✓ Different QoS for each SL



* E. Bolotin, I. Cidon, R. Ginosar and A. Kolodny., “QNoC: QoS architecture and design process for Network on Chip”, JSA special issue on NOC, December 2003

QNoC Architecture

- Mesh Topology



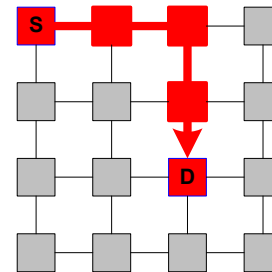
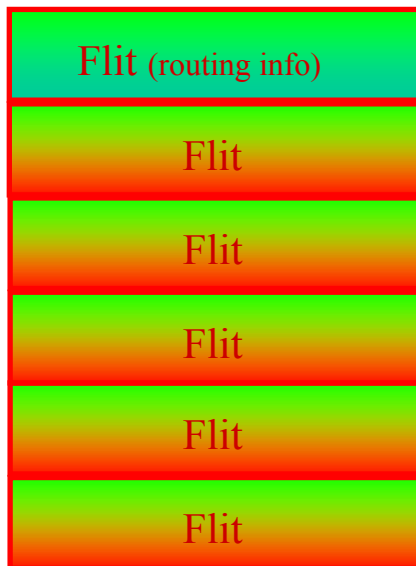
- Fixed shortest path routing (X-Y)
 - ✓ Simple Router (no tables, simple logic)
 - ✓ Power efficient communication
 - ✓ No deadlock scenario



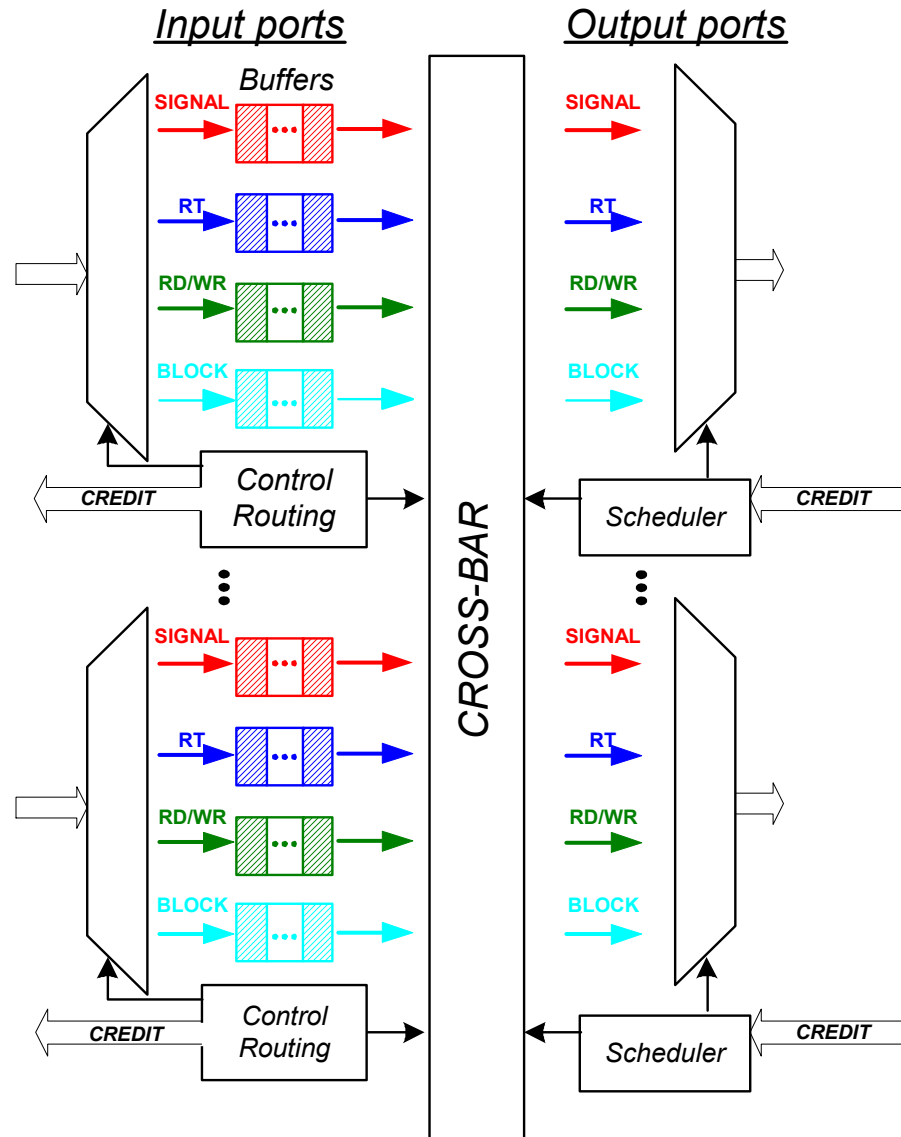
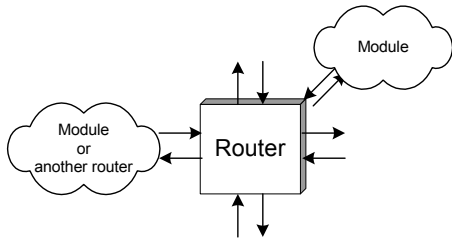
QNoC Architecture

- Wormhole Routing
 - ✓ For reduced buffering

Wormhole Packet:



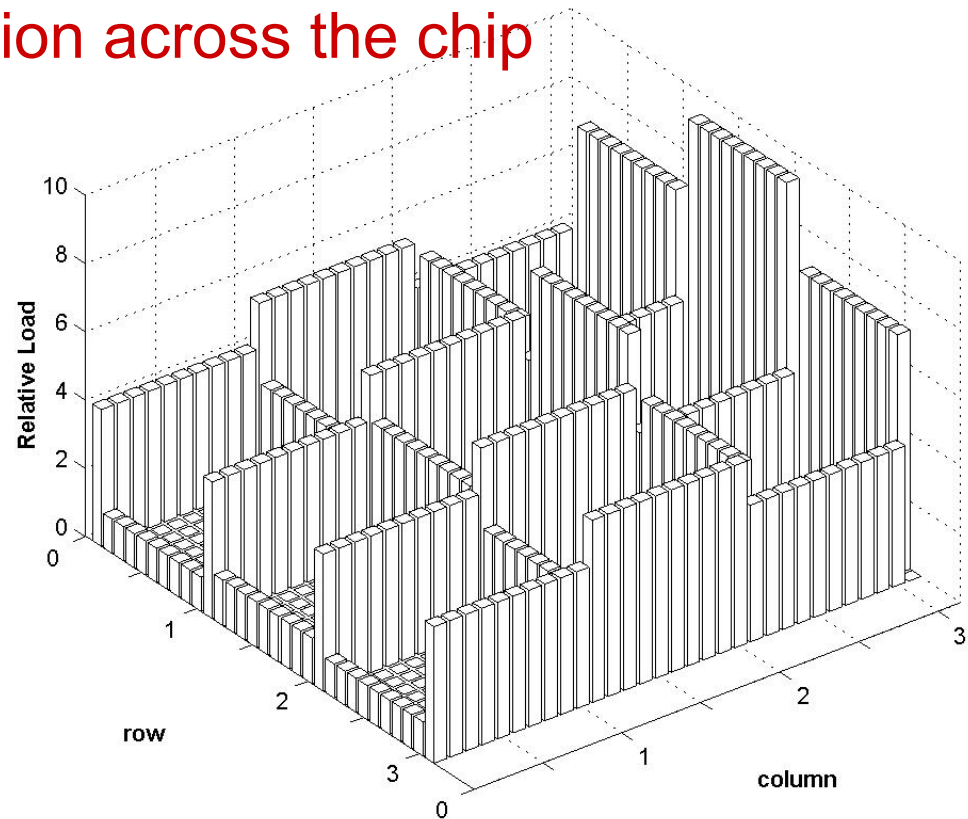
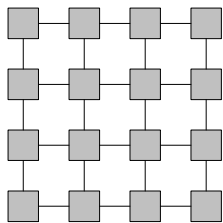
QNoC Wormhole Router



QNoC Design Process - Optimization

- Trim Unnecessary Resources
- Adjust each link capacity according to its load
 - ✓ Equal link utilization across the chip

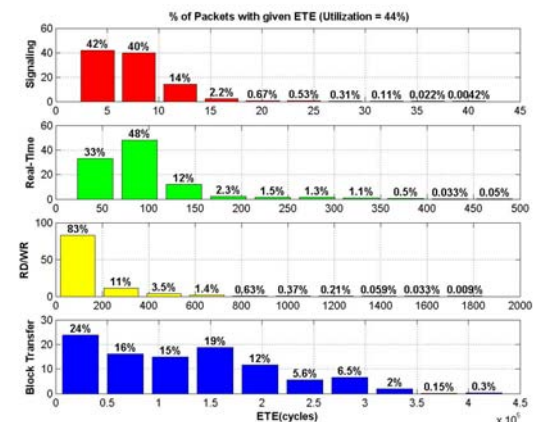
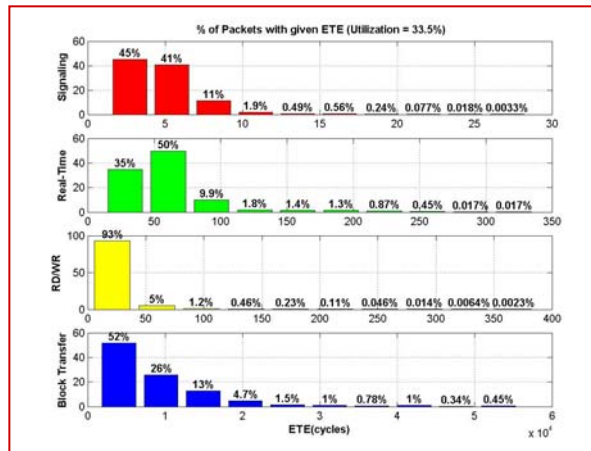
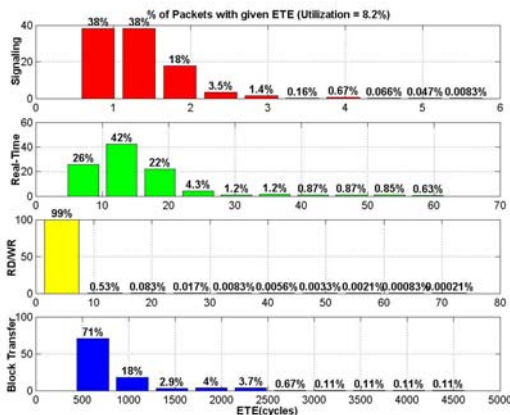
Example: (Uniform mesh)



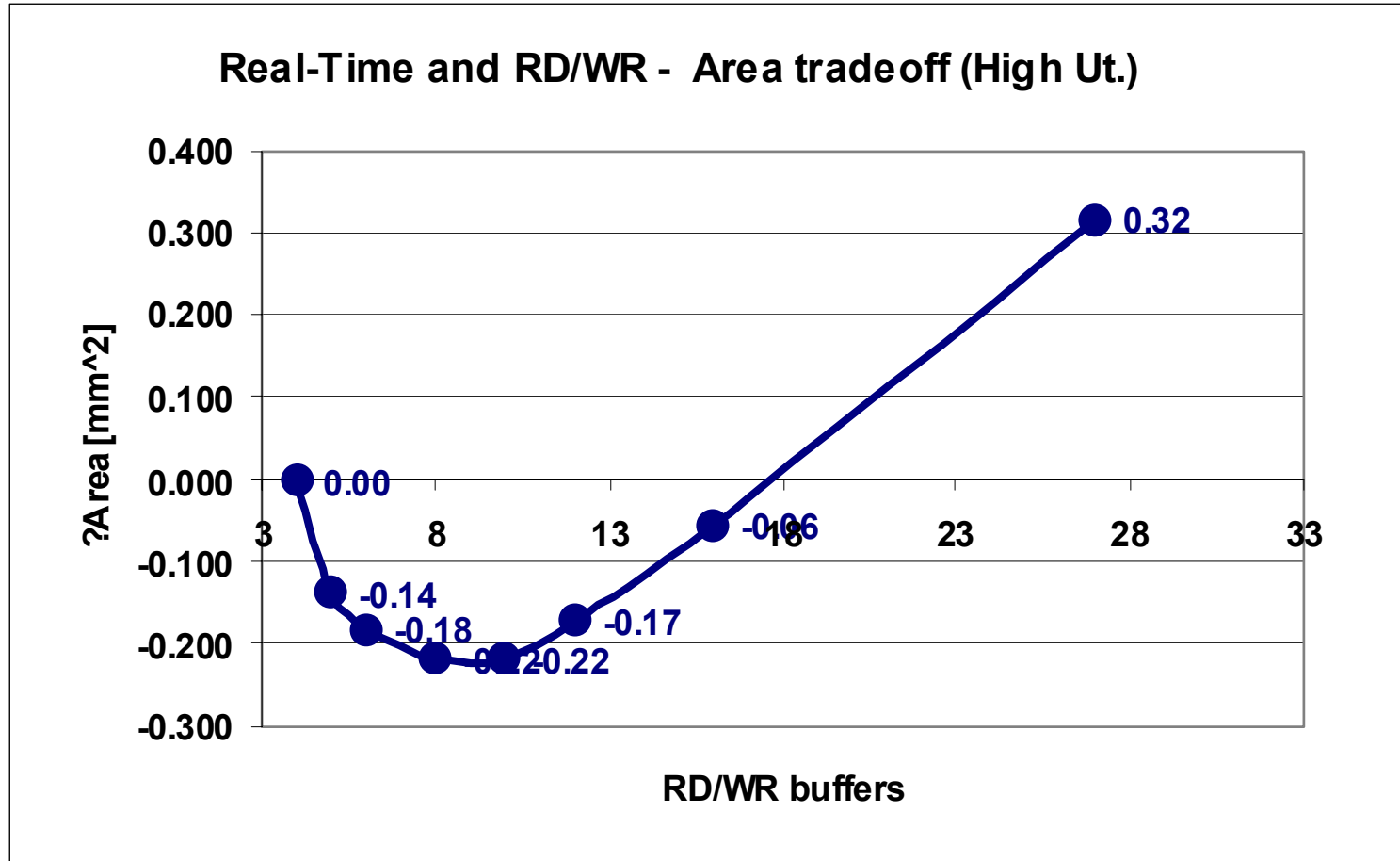
Non-Uniform Traffic - Observations

Various network BW allocations were tried in order to reach desired QoS requirements:

Allocated Link BW [Gbps]	Average Link Utilization [%]	Packet ETE delay of packets [ns or cycles]			
		Signaling (99.9%)	Real-Time (99.9%)	RD/WR (99%)	Block-Transfer (99%)
2752Gbps	8.2	5ns	60ns	20ns	4.5 μ s
1376Gbps	16.5	10ns	120ns	50ns	13 μ s
688Gbps	33.5	20ns	270ns	120ns	45 μ s
459Gbps	44	35ns	400ns	1.3 μ s	350 μ s



Cost tradeoffs



* E. Bolotin, I. Cidon, R. Ginosar and A. Kolodny, “Cost Considerations in Network-on-Chip”, INTEGRATION – the VLSI journal, 2003)



Future Directions

- Module placement
- Link Optimization
- Interfaces
- Physical design
- Error correction
- Synchronization
- NoC for reconfigurable chips
- Design and validation methods
-

