# Computer Architecture for the Next Millenium

November 1, 1999

William J. Dally

Computer Systems Laboratory

Stanford University

billd@csl.stanford.edu

# Outline
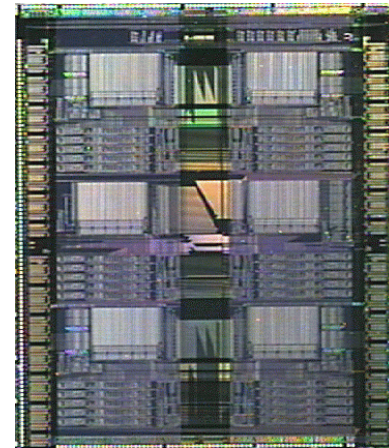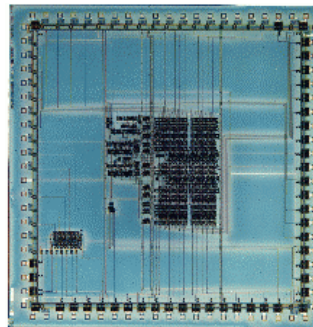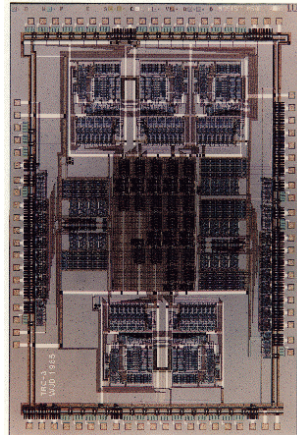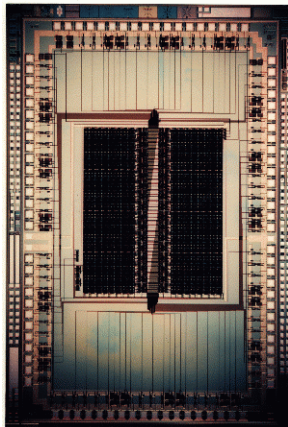
- The Stanford Concurrent VLSI Architecture Group
- Forces acting on computer architecture
  - applications (media)
  - technology (wire-limited)
  - techniques (explicit parallelism)
- Example: register organization
  - distributed register files
- *Imagine* a stream processor
  - 20GFLOPS on a $0.5cm^2$ chip
- Tremendous opportunities and challenges for computer architecture in the next millenium
  - its not a *mature* field yet

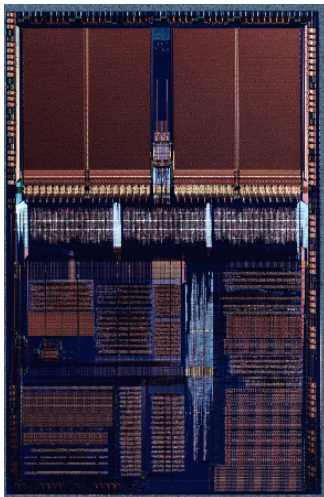# The Concurrent VLSI Architecture Group

- Architecture and design technology for VLSI
- Routing chips
  - Torus Routing Chip, Network Design Frame, Reliable Router
  - Basis for Intel, Cray/SGI, Mercury, Avici network chips
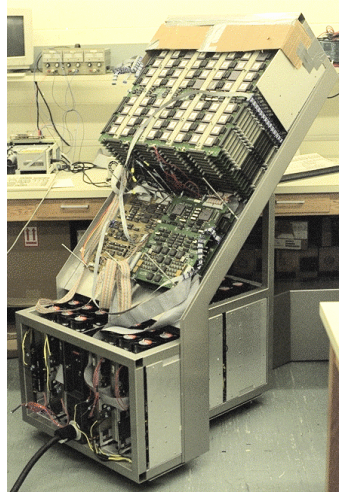
# Parallel computer systems

- J-Machine (MDP) led to Cray T3D/T3E
- M-Machine (MAP)
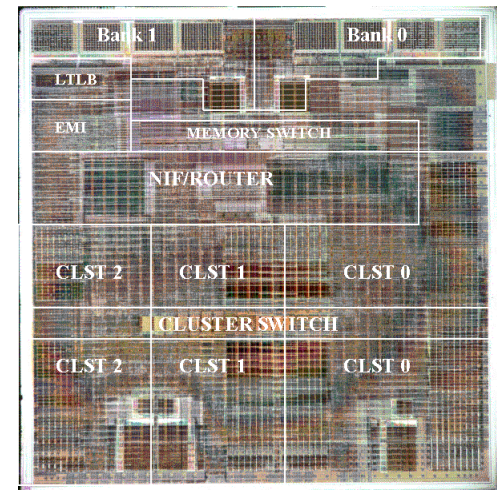  - Fast messaging, scalable processing nodes, scalable memory architecture

MDP Chip          J-Machine          Cray T3D          MAP Chip
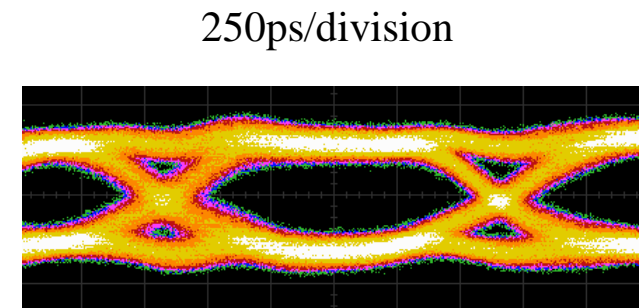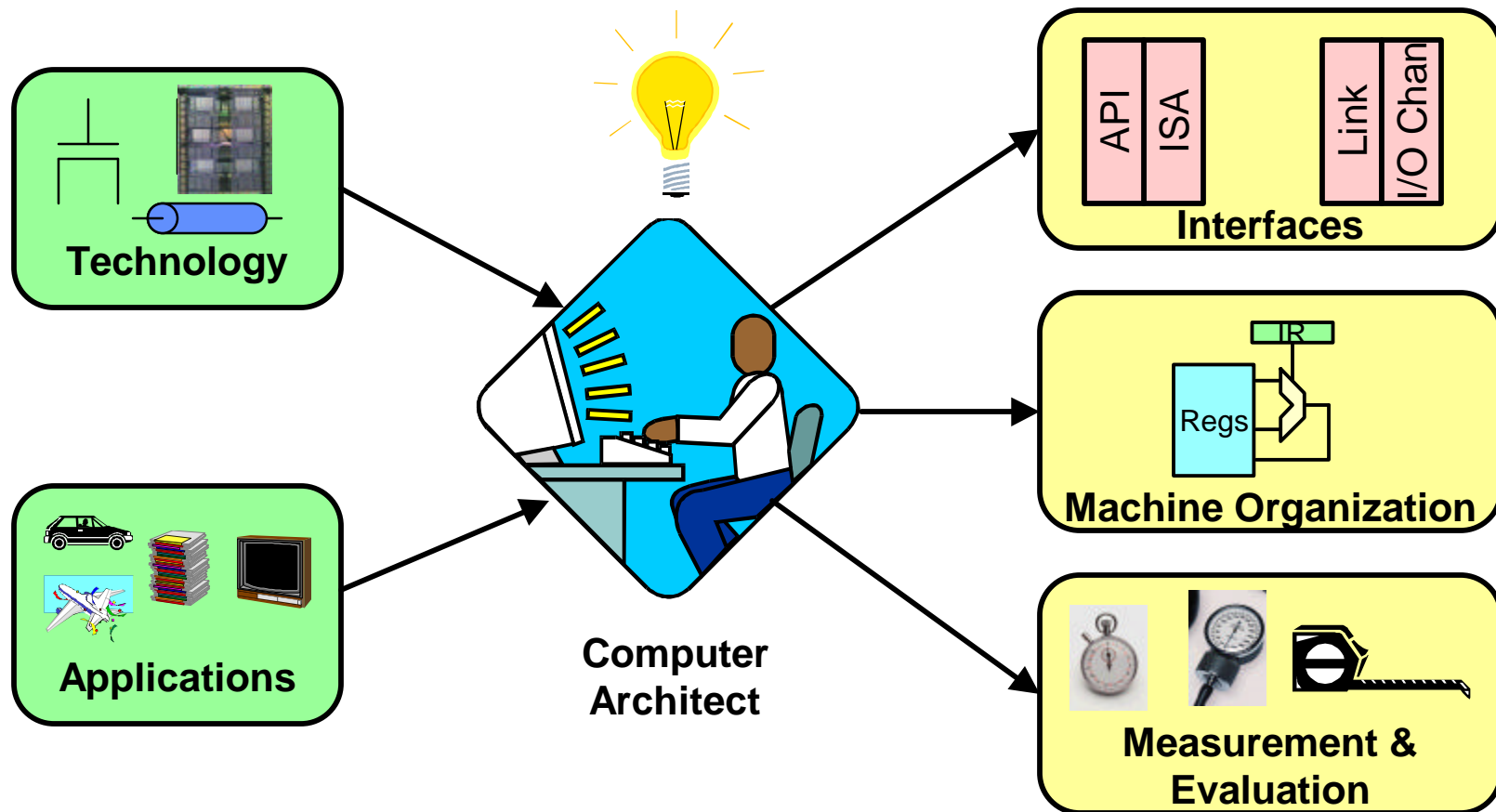
# Design technology

- ## Off-chip I/O
  - Simultaneous bidirectional signaling, 1989
    - now used by Intel and Hitachi
  - High-speed signalling
    - 4Gb/s in 0.6µm CMOS, Equalization, 1995

- ## On-Chip Signalling
  - Low-voltage on-chip signalling
  - Low-skew clock distribution

- ## Synchronization
  - Mesochronous, Plesiochronous
  - Self-Timed Design

250ps/division



4Gb/s CMOS I/O

# What is Computer Architecture?



**Technology**

**Applications**

**Computer Architect**

**Interfaces**

API | ISA | Link | I/O Chan

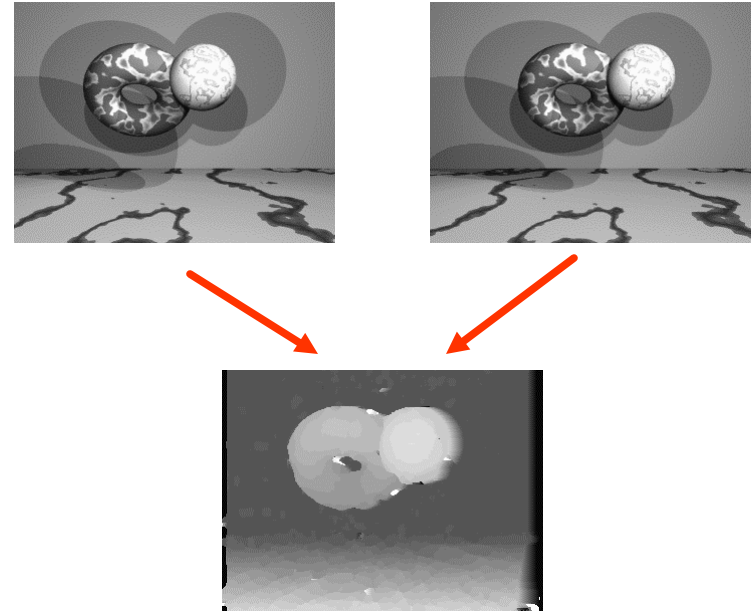**Machine Organization**

IR

Regs

**Measurement & Evaluation**

# Forces Acting on Architecture

- Applications - shifting towards *media* applications dealing with *streams* of low-precision samples
  - – video, graphics, audio, DSL modems, cellular base stations
- Technology - becoming *wire-limited*
  - – power and delay dominated by communication, not arithmetic
  - – global structures: register files and instruction issue don't scale
- Technique - Micro-architecture - ILP has been *mined out*
  - – to the point of diminishing returns on squeezing performance from sequential code
  - – explicit parallelism (data parallelism and thread-level parallelism) required to continue scaling performance

# Applications

- ## Little locality of reference
  - read each pixel once
  - often non-unit stride
  - but there is producer-consumer locality
- ## Very high arithmetic intensity
  - 100s of arithmetic operations per memory reference
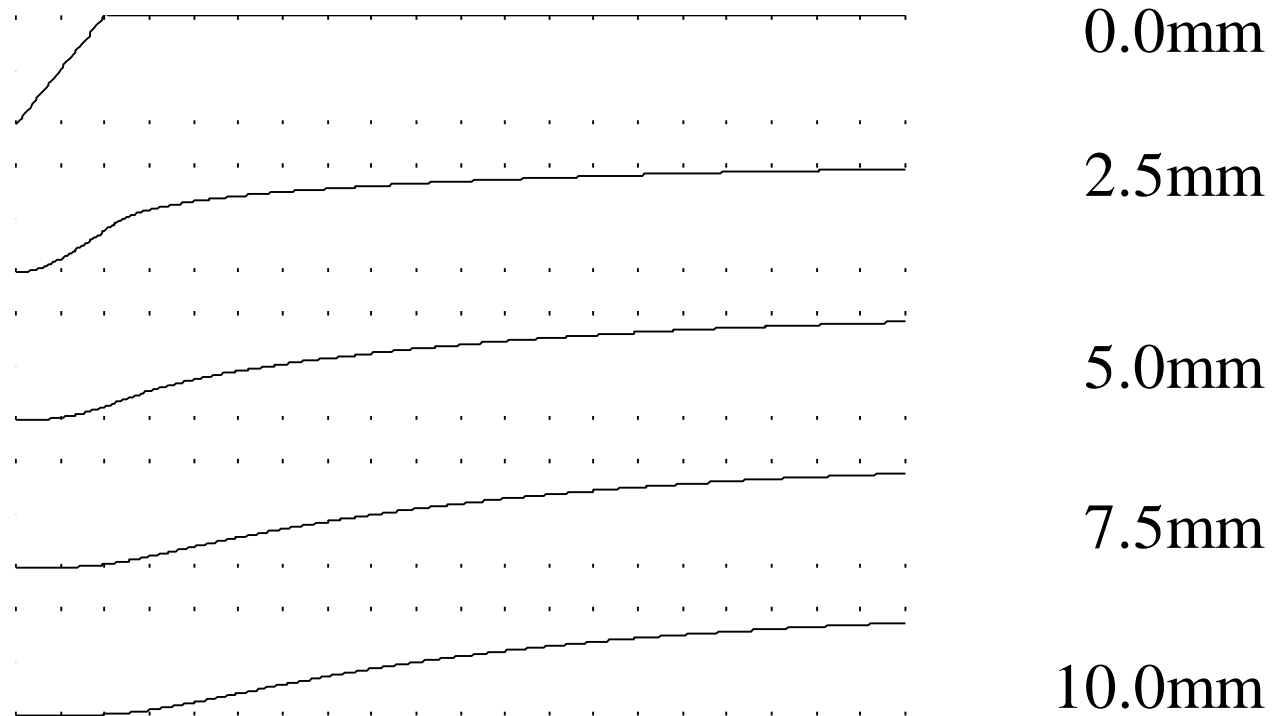- ## Dominated by low-precision (16-bit) integer operations
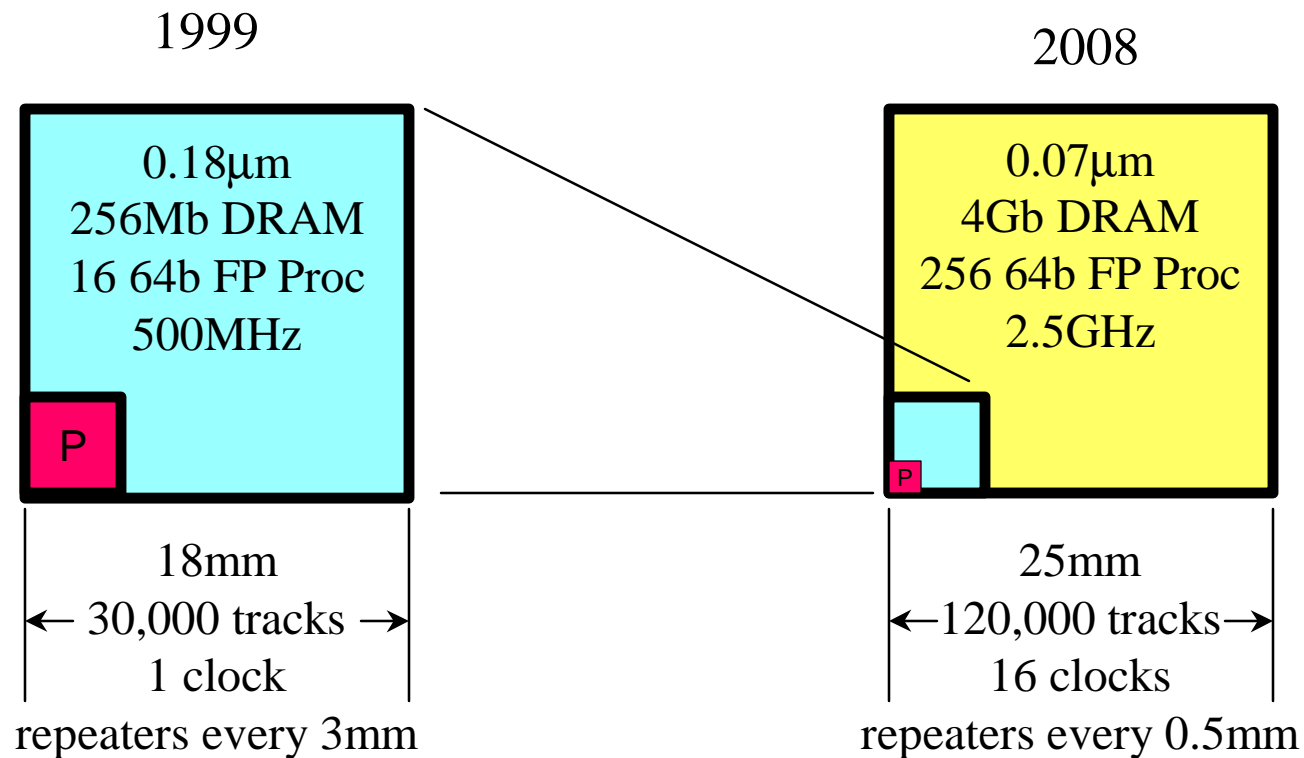
# Wires Are Becoming Like Wet Noodles

Minimum width
wire in an 0.35μm
process

0.0mm

2.5mm

5.0mm

7.5mm

10.0mm

# Technology scaling makes communication *the* scarce resource

1999

0.18μm
256Mb DRAM
16 64b FP Proc
500MHz

P

18mm
← 30,000 tracks →
1 clock
repeaters every 3mm

2008

0.07μm
4Gb DRAM
256 64b FP Proc
2.5GHz

P

25mm
←120,000 tracks→
16 clocks
repeaters every 0.5mm

# Care and Feeding of ALUs



Instruction
Bandwidth

Data
Bandwidth

Instr.
Cache

IP

IR

Regs

'Feeding' Structure Dwarfs ALU

# What Does This Say About Architecture?

- Tremendous opportunities
  - Media problems have lots of parallelism and locality
  - VLSI technology enables 100s of ALUs per chip (1000s soon)
    - (in 0.18um 0.1mm$^2$ per integer adder, 0.5mm$^2$ per FP adder)
- Challenging problems
  - Locality - global structures won't work
  - Explicit parallelism - ILP won't keep 100 ALUs busy
  - Memory - streaming applications don't cache well
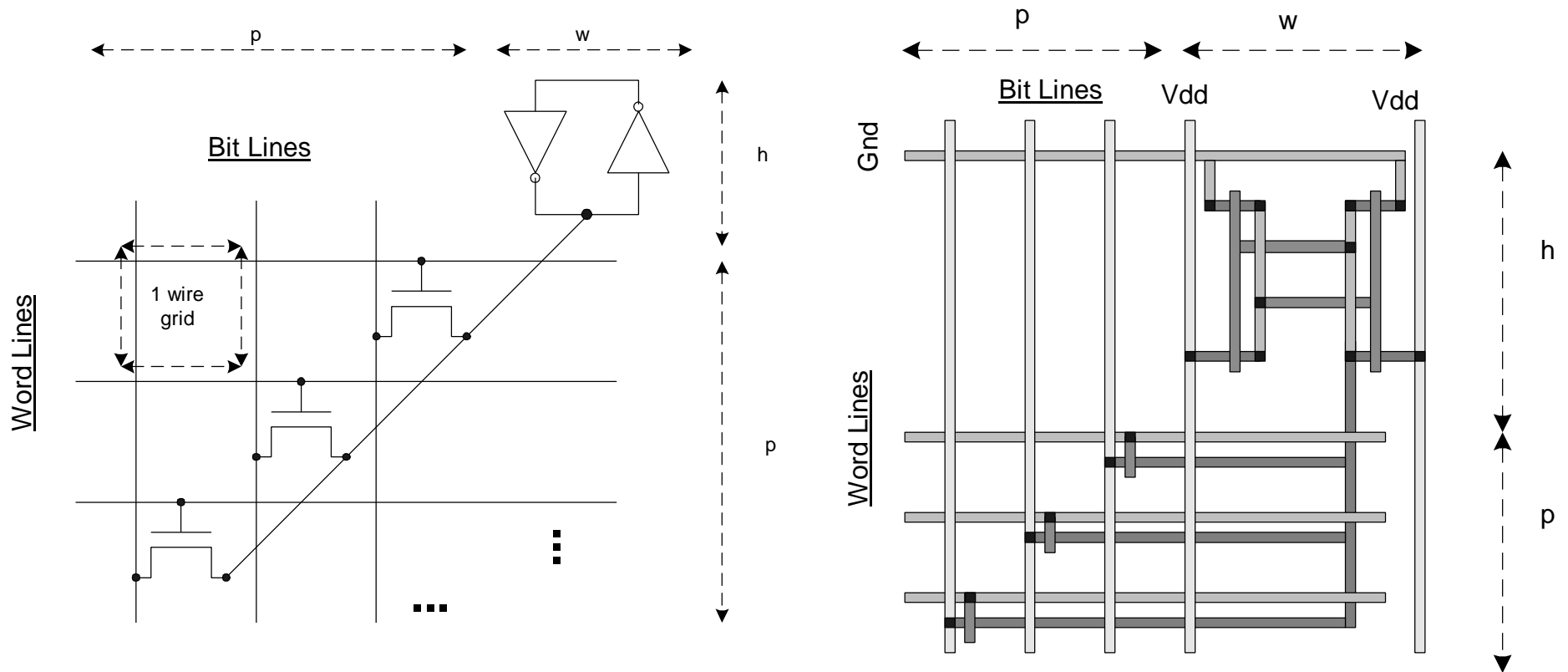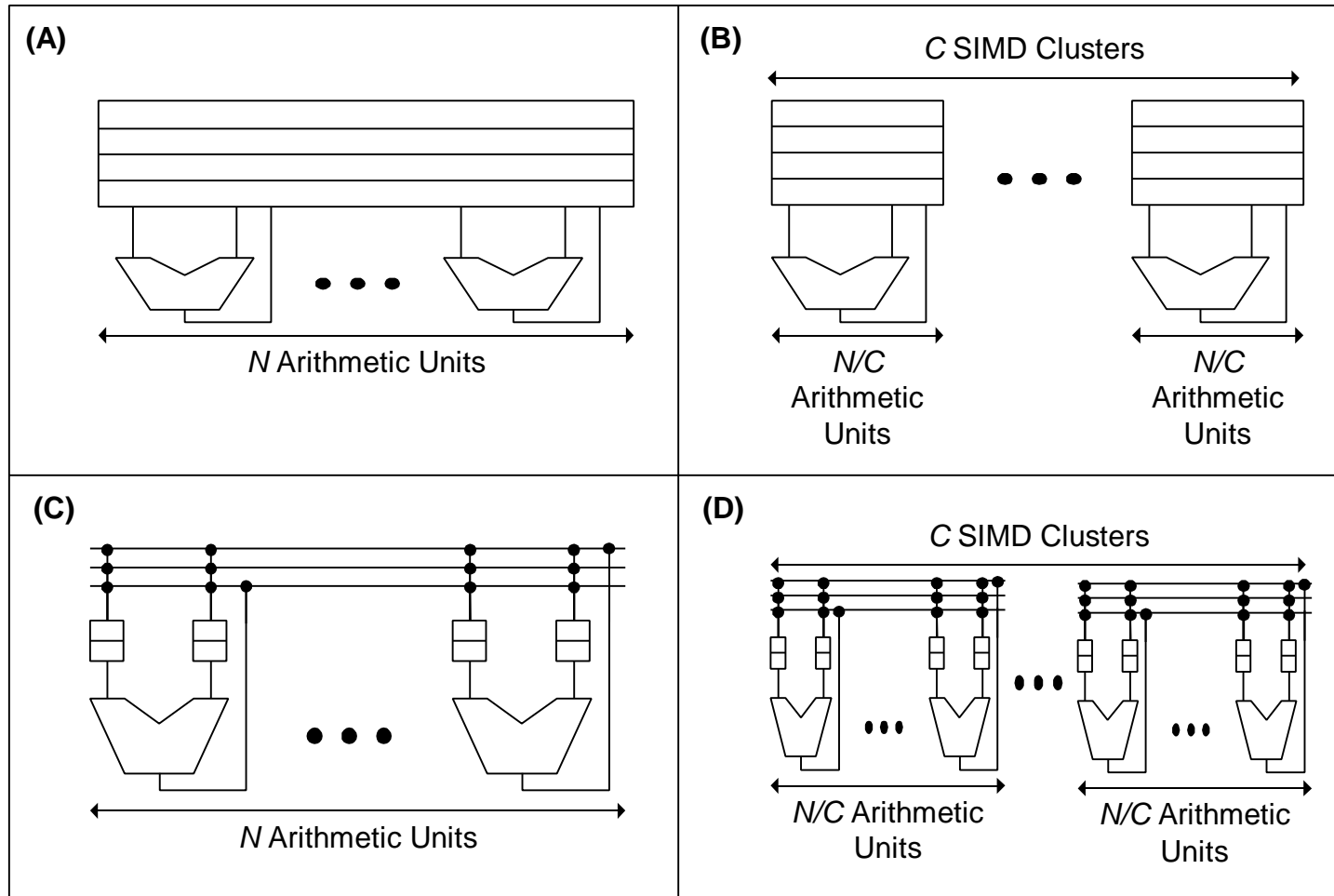- Its time to try some new approaches

# Example
# Register File Organization

- Register files serve two functions:
  - Short term storage for intermediate results
  - Communication between multiple function units

- Global register files don't scale well as N, number of ALUs increases
  - Need more registers to hold more results (grows with N)
  - Need more ports to connect all of the units (grows with $N^2$)
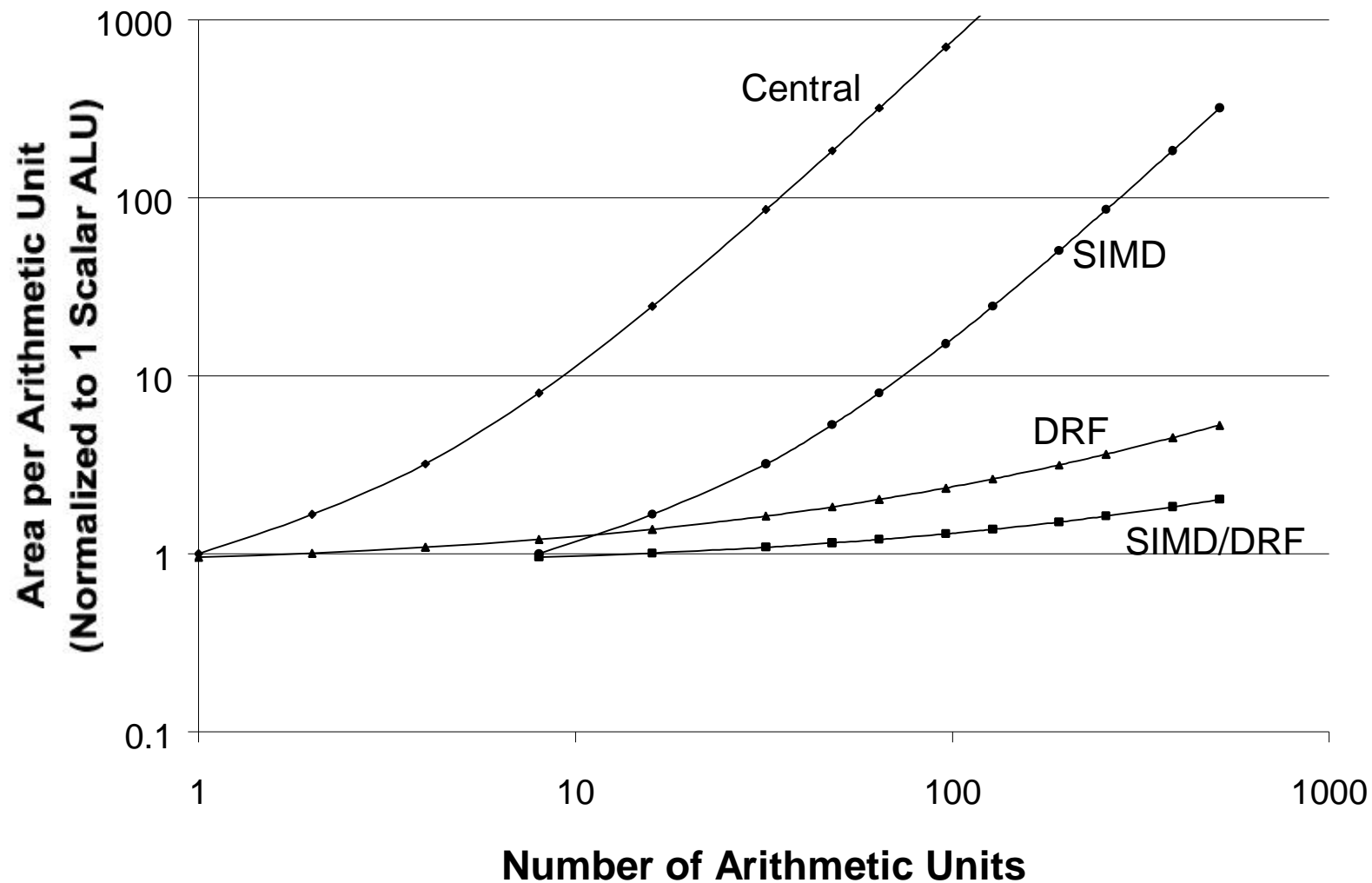
# Register Cells are Mostly Switch
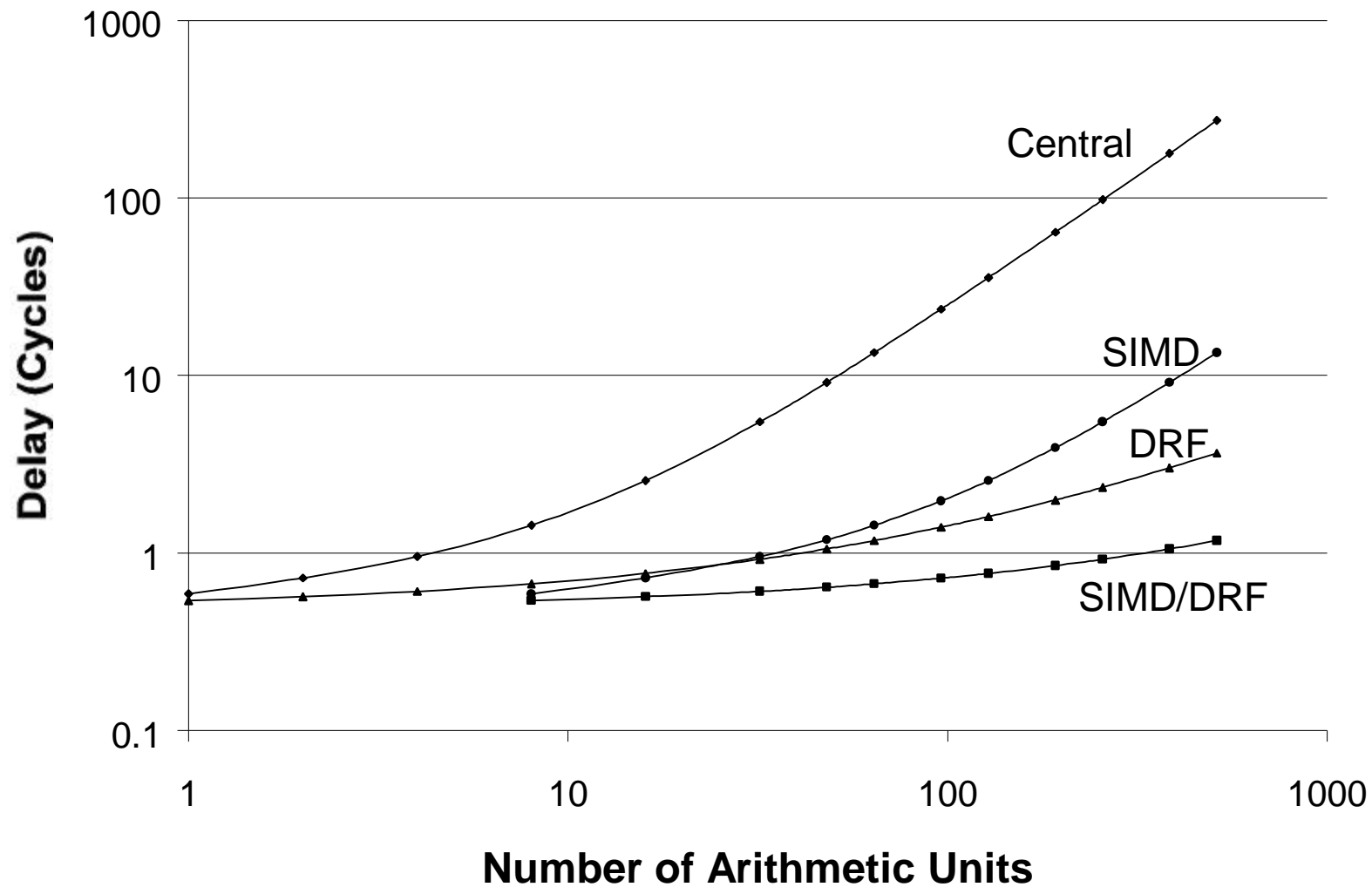
# Register Architecture for 'wide' Processors

**(A)**

N Arithmetic Units

**(B)** C SIMD Clusters

N/C Arithmetic Units

N/C Arithmetic Units

**(C)**

N Arithmetic Units

**(D)** C SIMD Clusters

N/C Arithmetic Units
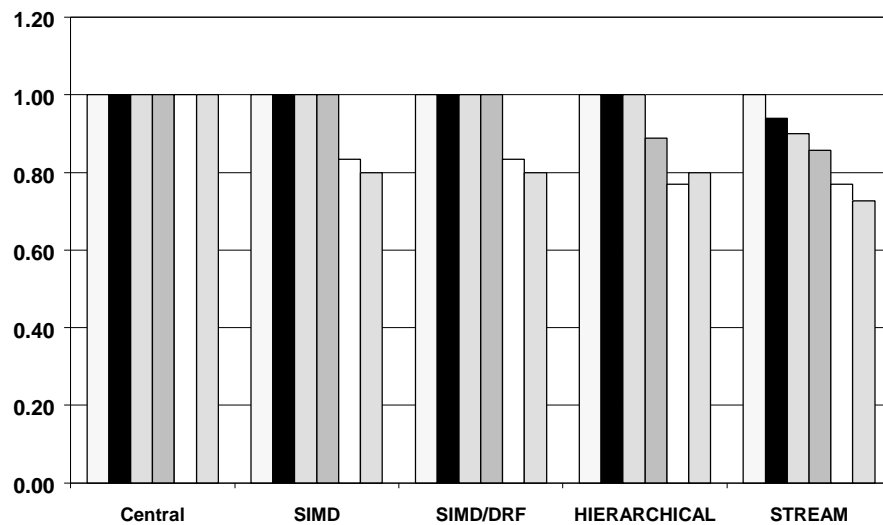
N/C Arithmetic Units

# Area of Register Organizations

# Delay of Register Organizations
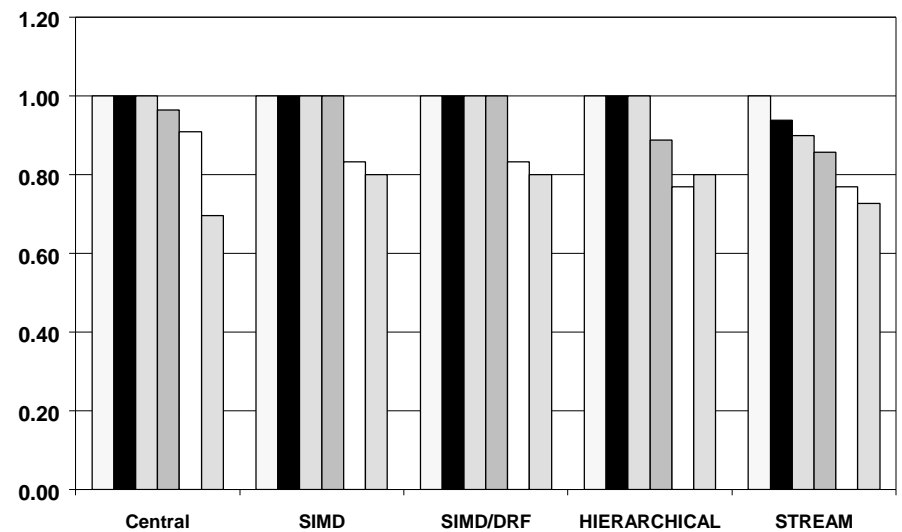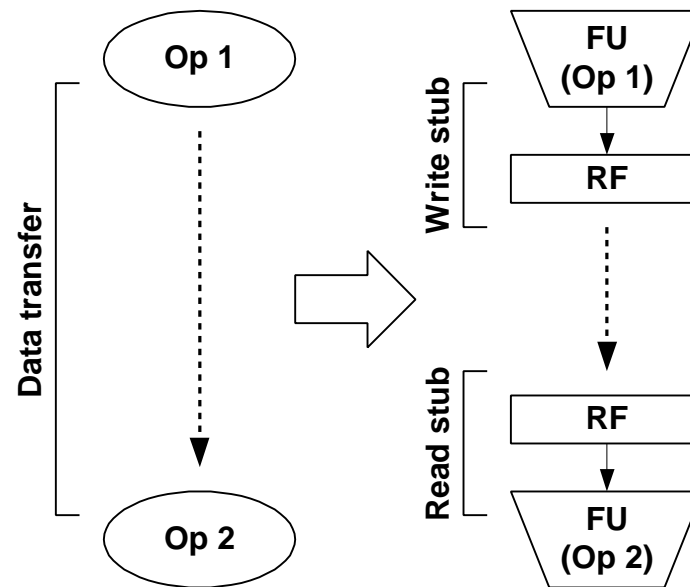
# Performance of Register Organizations
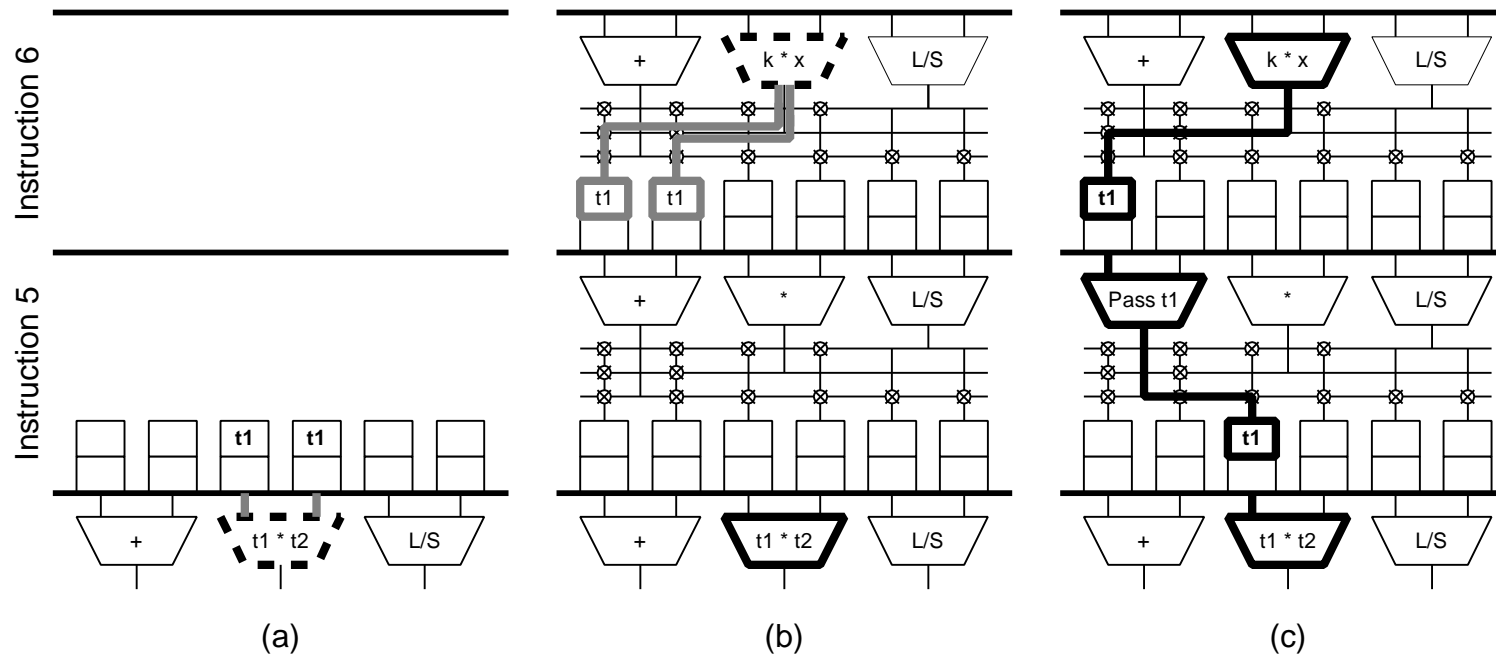


**(A) Raw Performance**

**(B) Performance with Latency**

# Stubs Abstract the Communication Between Operations

# A Communication Example



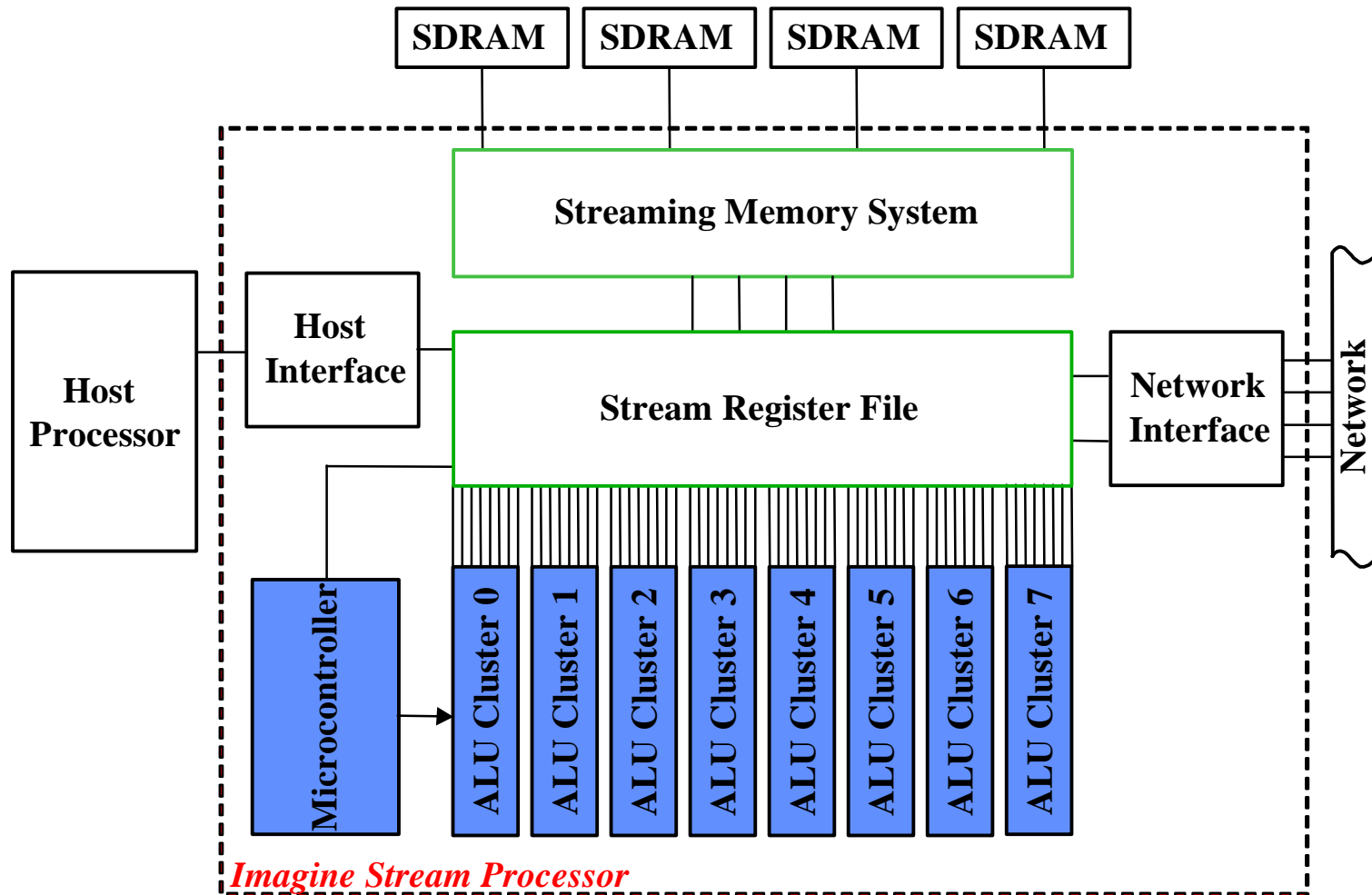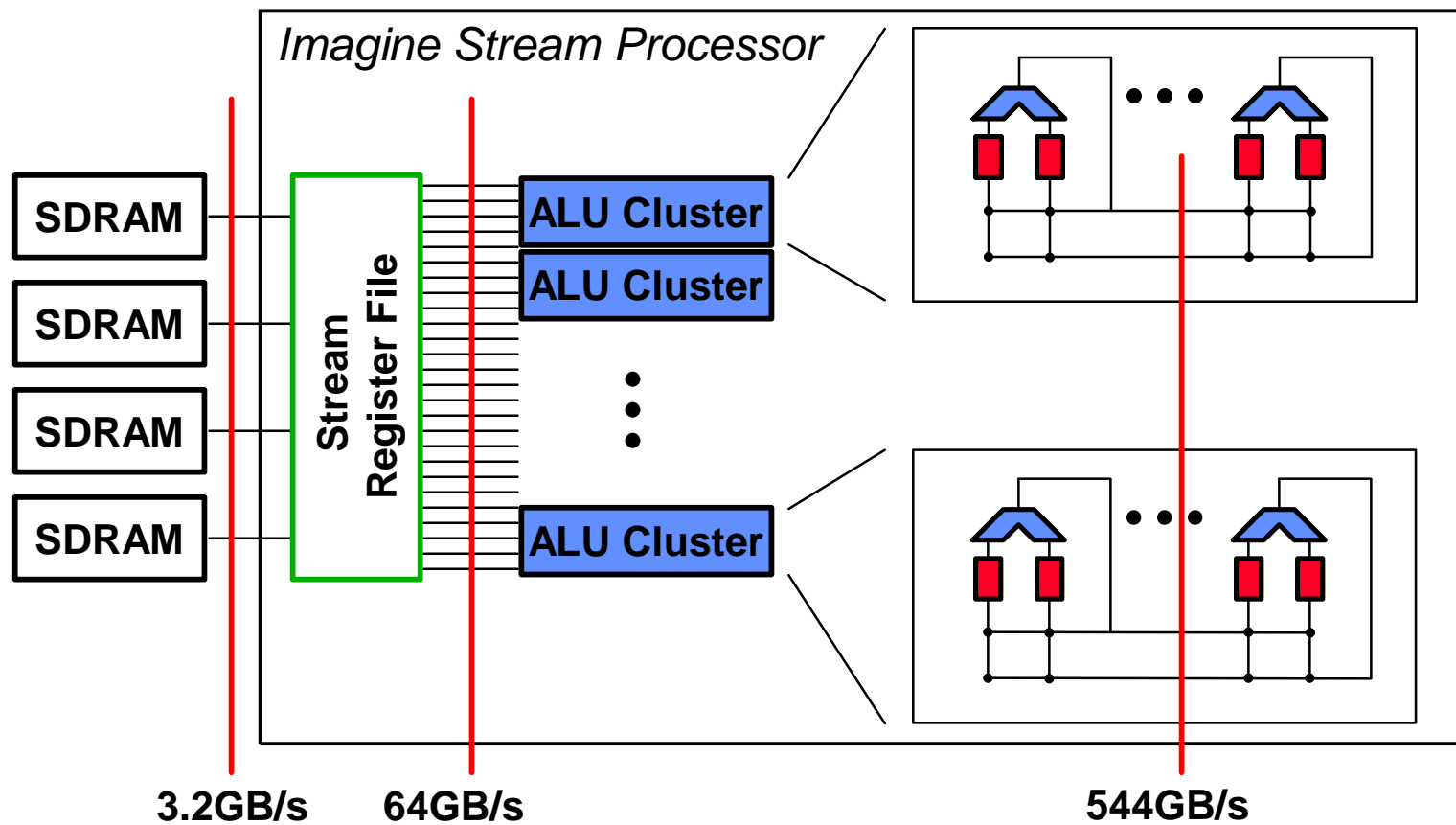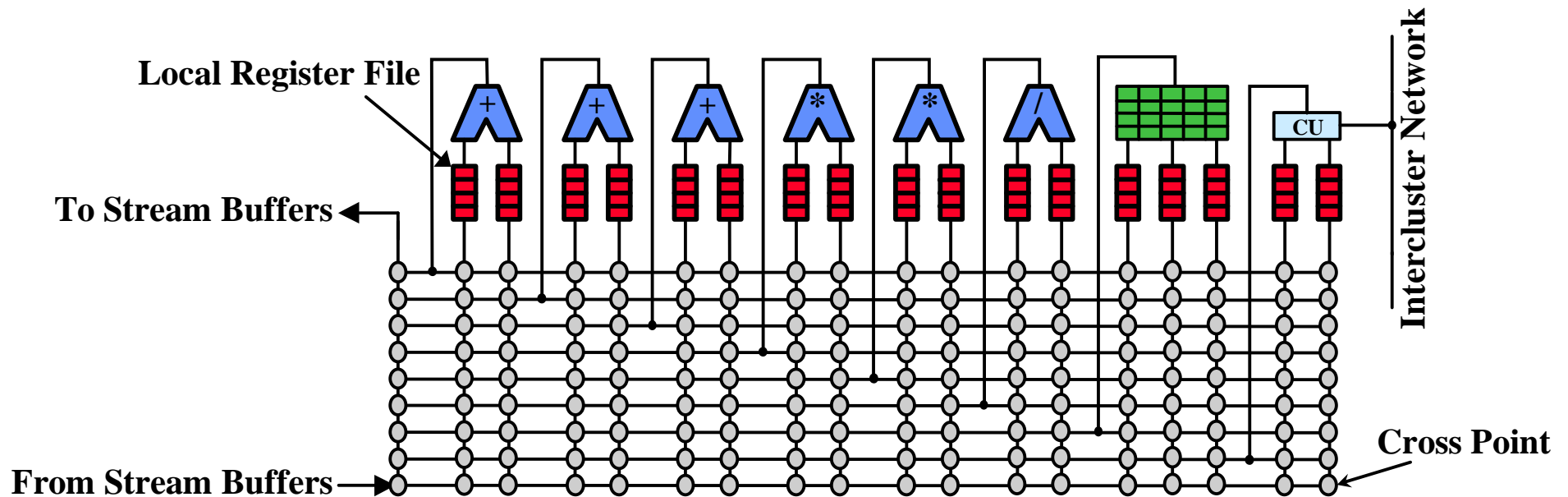(a)          (b)          (c)

# The Imagine Stream Processor

# Data Bandwidth Hierarchy



**3.2GB/s**  **64GB/s**  **544GB/s**

# Cluster Architecture



- VLIW organization with shared control
- Local register files provide high data bandwidth

# Imagine is a *Stream Processor*

- Instructions are Load, Store, and Operate
  - operands are streams
  - also Send and Receive for multiple-imagine systems
- Operate performs a *compound stream operation*
  - read elements from input streams
  - perform a local computation
  - append elements to output streams
  - repeat until input stream is consumed
  - (e.g., triangle transform)
- Order of magnitude less global register bandwidth than a vector processor

# Triangle Rendering

Memory Bandwidth    Register Bandwidth

**Input Data**

**Image Depth & Color**

Transform
Shade
Project/ Cull
Span Setup
Process Span
Sort
Compact
Z-Composite

Triangle Records
Triangle Records
Shaded Triangle Records
Projected Triangle Records
Span Records
Fragment Records
Fragment Records
Image Buffer Indices
Pixel Depth & Color
Pixel Depth & Color
Pixel Depth & Color

word
record

# Bandwidth Demands

## Transform Kernel

| References (per Δ) | Stream | Scalar | | Vector | |
|---|---|---|---|---|---|
| Memory | 5.5 | 117 | (21.3) | 48 | (8.7) |
| Global RF | 48 | 624 | (13.0) | 261 | (5.4) |
| Local RF | 372 | N/A | | N/A | |

# Data Parallelism is easier than ILP

| Kernel | 1 to 8 Cluster Speedup |
|---|---|
| FFT (1024) | 6.4 |
| DCT (8x8) | 7.8 |
| Blockwarp (8x8) | 7.2 |
| Transform ($\Delta$) | 8.0 |
| Harmonic Mean | 7.3 |

# Conventional Approaches to
# Data-Dependent Conditional Execution



Y    x>0    N

B

C

J

K

**Data-Dependent
Branch**

A

x>0

Y

B

C

Whoops

J

K

Speculative
Loss
D x W
~100s

A

y=(x>0)

B    if y

J    if ~y

C    if y

K    if ~y
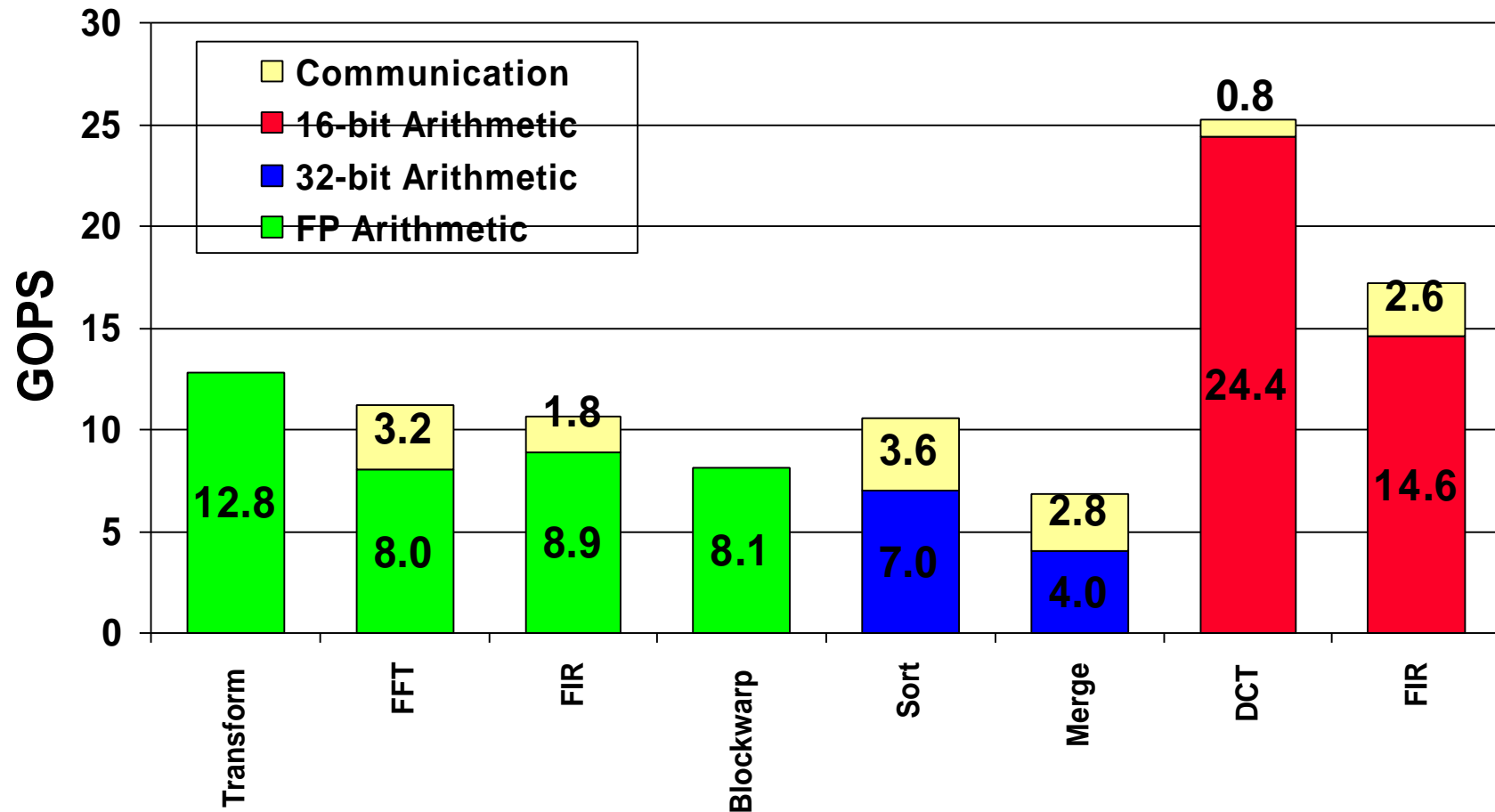
Exponentially
Decreasing
Duty Factor

# Zero-Cost Conditionals

- Most Approaches to Conditional Operations are Costly
  - Branching control flow - dead issue slots on mispredicted branches
  - Predication (SIMD select, masked vectors) - large fraction of execution 'opportunities' go idle.
- Conditional Streams
  - append an *element* to an output stream depending on a *case* variable.
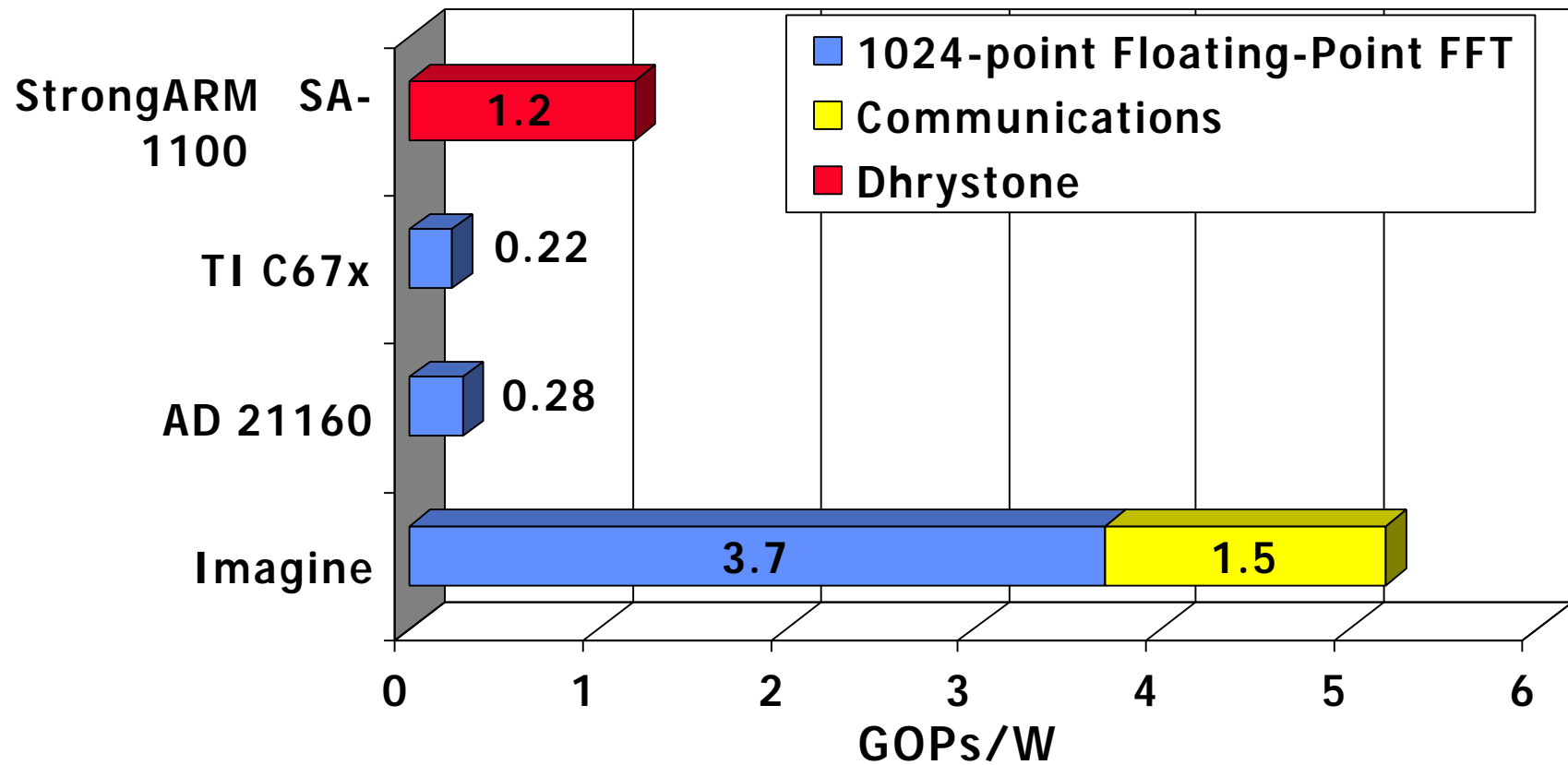
Value Stream

Output Stream

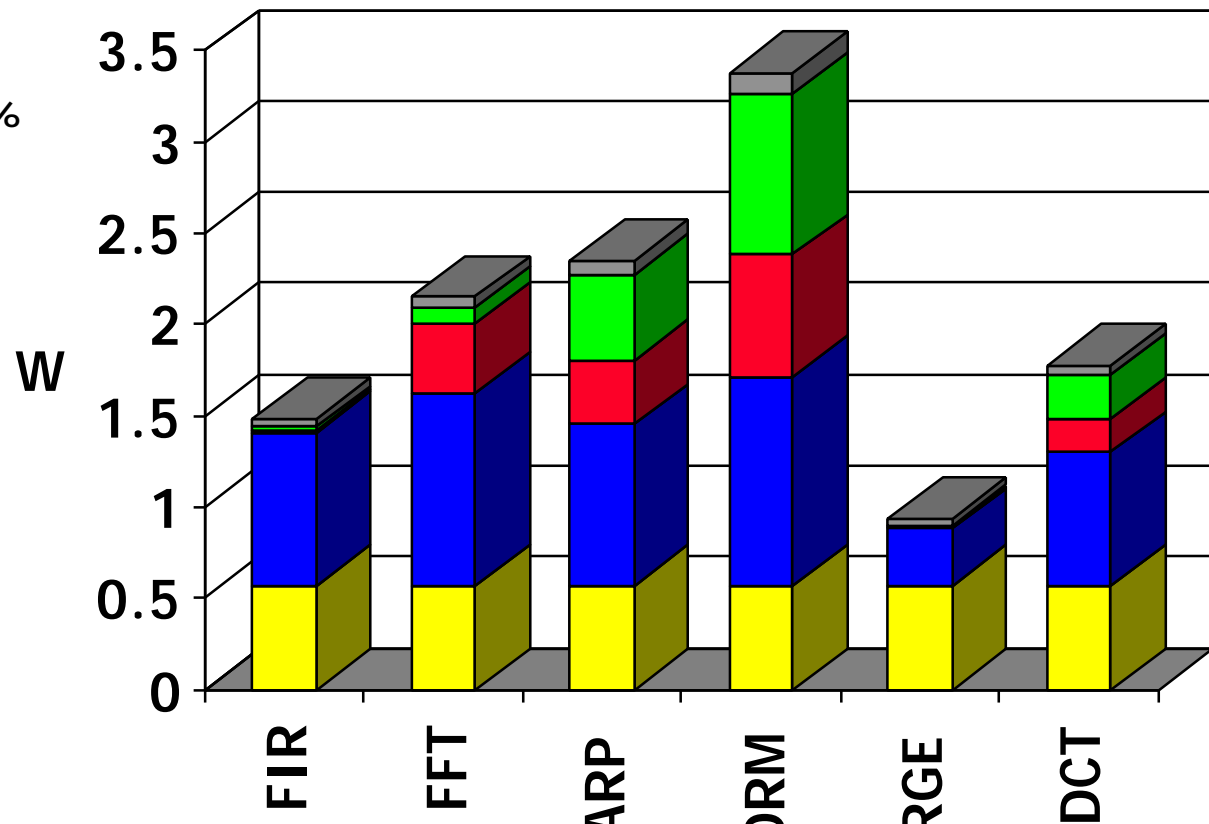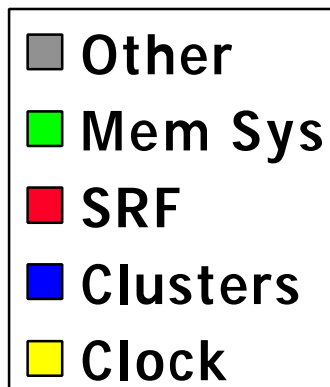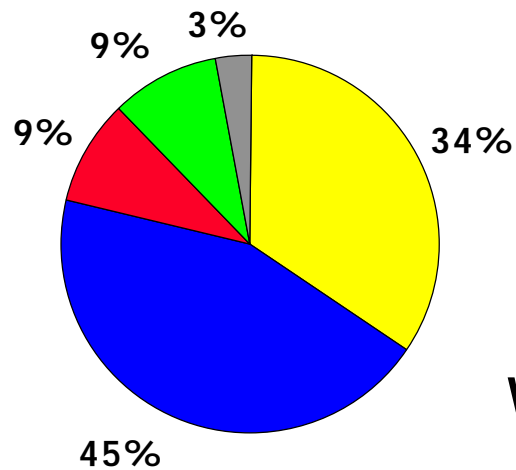Case Stream {0,1}

# Sustainable Performance

# Power Comparison
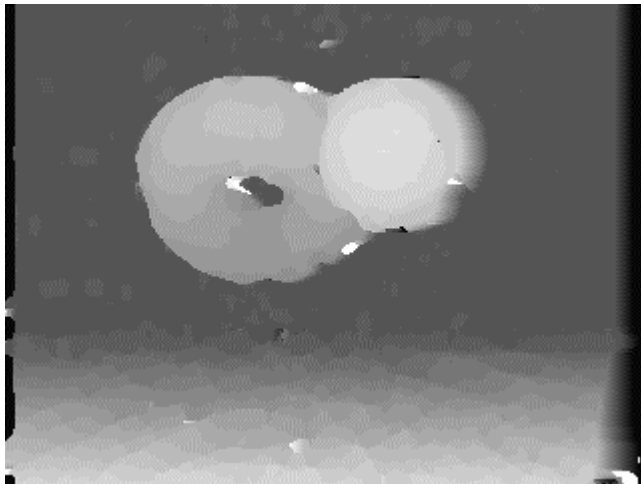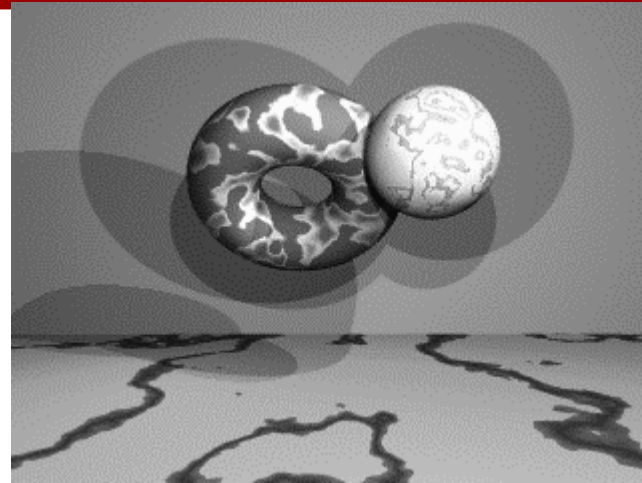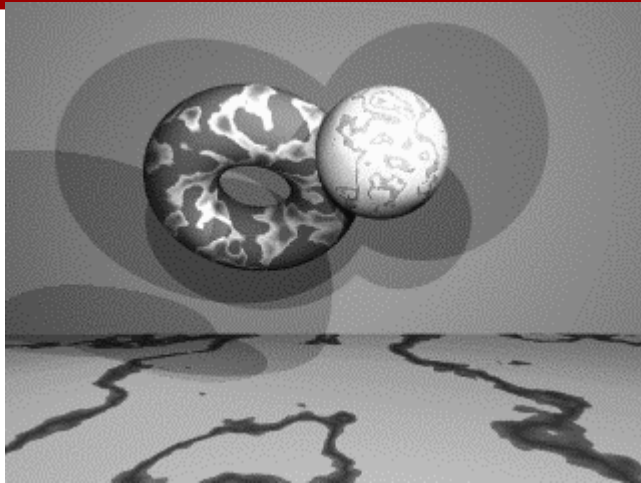


-Source: Web Pages of Intel, TI, and Analog Devices

# Power and Performance

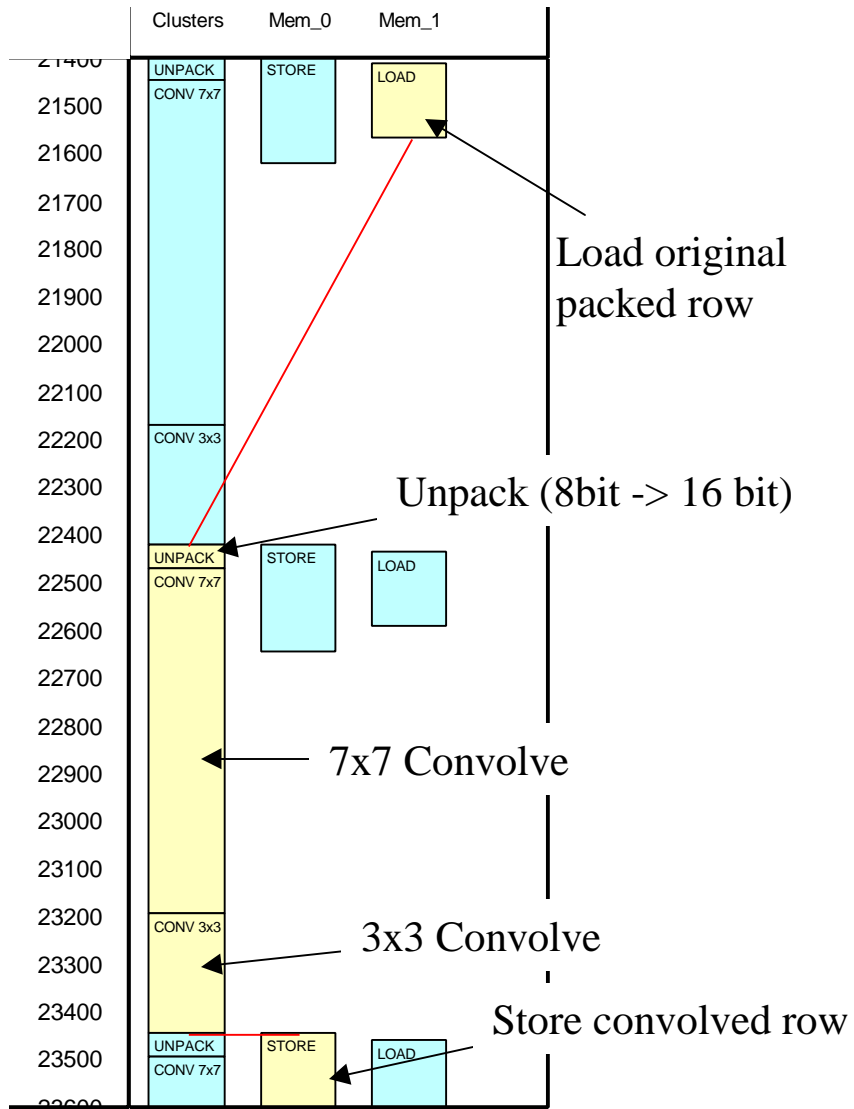# A Look Inside an Application
## Stereo Depth Extraction
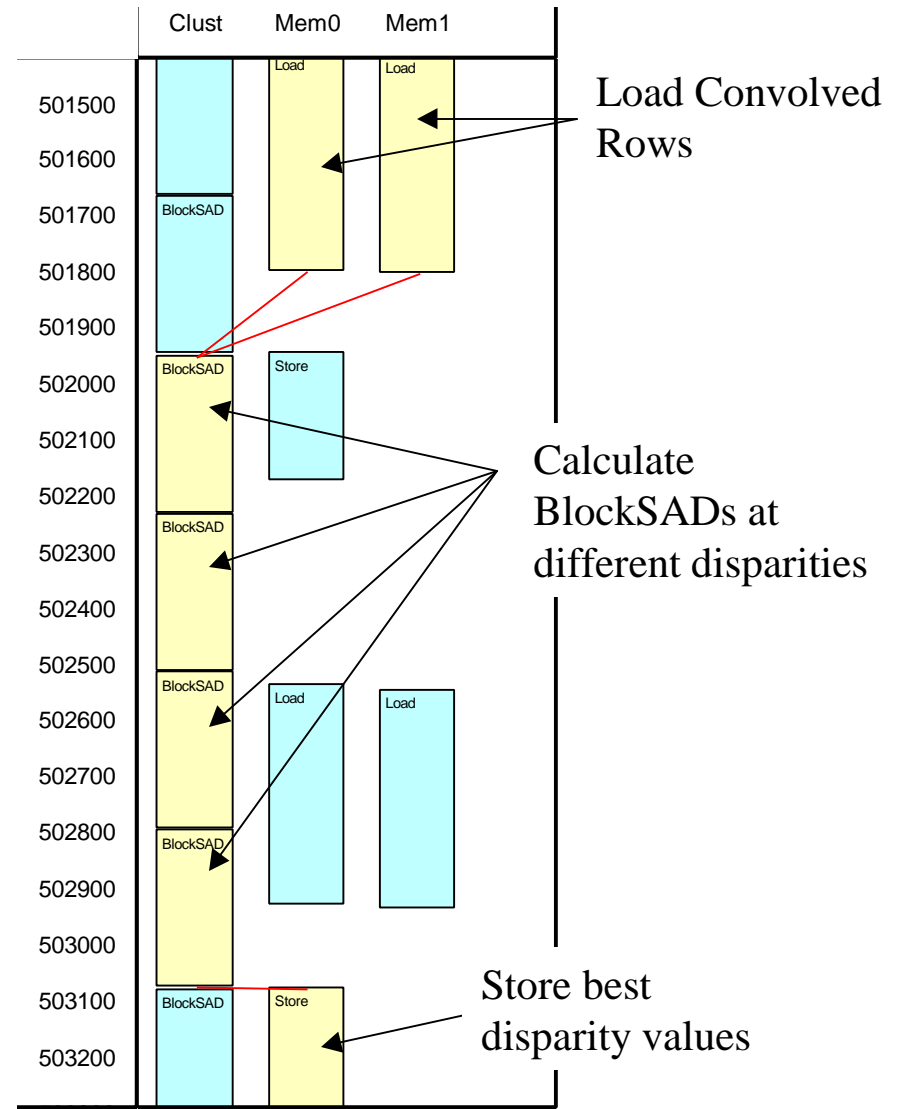


- 320x240 8-bit grayscale images
- 30 disparity search
- 220 frames/second
- 12.7 GOPS
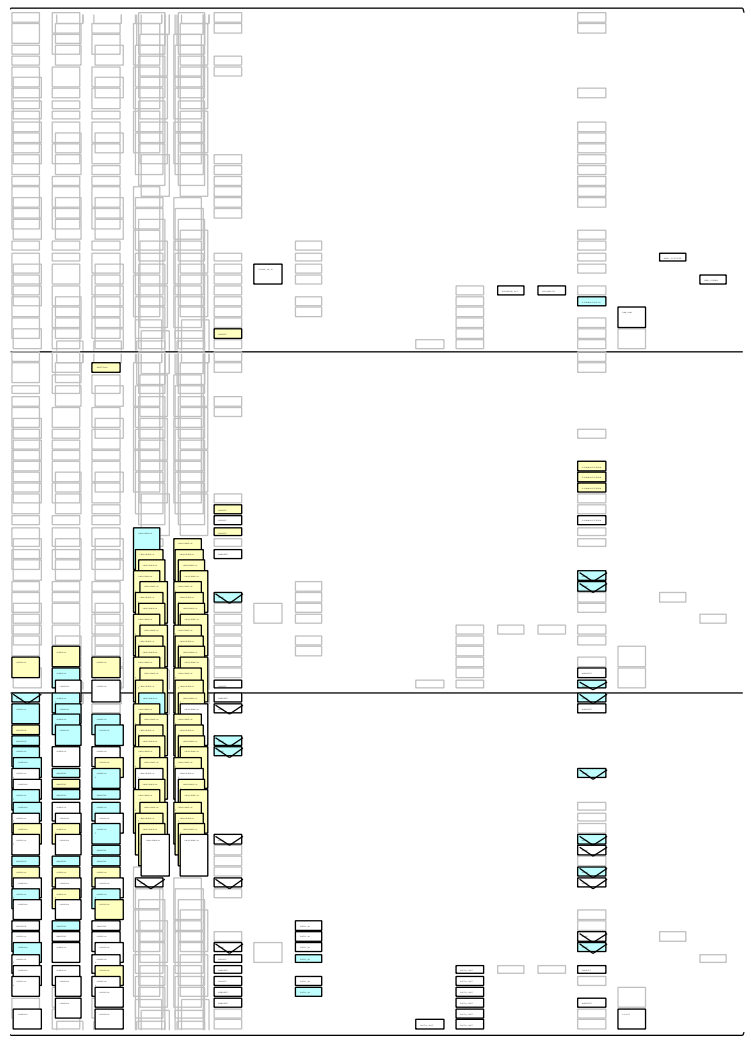- 5.7 GOPS/W

# Stereo Depth Extractor

## Convolutions

| Clusters | Mem_0 | Mem_1 |
|---|---|---|

- 21400 — UNPACK / CONV 7x7 | STORE | LOAD
- 21500
- 21600
- 21700
- 21800
- 21900
- 22000
- 22100
- 22200 — CONV 3x3
- 22300
- 22400 — UNPACK / CONV 7x7 | STORE | LOAD
- 22500
- 22600
- 22700
- 22800
- 22900
- 23000
- 23100
- 23200 — CONV 3x3
- 23300
- 23400
- 23500 — UNPACK / CONV 7x7 | STORE | LOAD
- 23600

Load original packed row

Unpack (8bit -> 16 bit)

7x7 Convolve

3x3 Convolve

Store convolved row

## Disparity Search

| Clust | Mem0 | Mem1 |
|---|---|---|

- 501500 — | Load | Load
- 501600
- 501700 — BlockSAD
- 501800
- 501900
- 502000 — BlockSAD | Store
- 502100
- 502200
- 502300 — BlockSAD
- 502400
- 502500
- 502600 — BlockSAD | Load | Load
- 502700
- 502800 — BlockSAD
- 502900
- 503000
- 503100 — BlockSAD | Store
- 503200

Load Convolved Rows

Calculate BlockSADs at different disparities

Store best disparity values

# 7x7 Convolve Kernel

# Imagine Summary

- Imagine operates on *streams* of records
  - simplifies programming
  - exposes locality and concurrency
- Compound stream operations
  - perform a subroutine on each stream element
  - reduces global register bandwidth
- Bandwidth hierarchy
  - use bandwidth where its inexpensive
  - distributed and hierarchical register organization
- Conditional stream operations
  - sort elements into homogeneous streams
  - avoid predication or speculation

# Computer Architecture for the Next Millenium

- Applications and technology are changing
  - *media* applications process streams of low-precision samples
  - wires dominate gates
- ILP is at the point of diminishing returns
- Tremendous opportunities for new architectures
  - new applications have *lots* of parallelism and locality
  - modern technology can build chips with 100s of ALUs (32b FP) 1000s in the near future
- The challenge is to develop architectures
  - that can harness this potential performance
  - in a way that can be easily programmed
- Stream processing is one approach, there are many others. We need to start exploring them