

New Directions in Computer Architecture

David A. Patterson

`http://cs.berkeley.edu/~patterson/talks`

`patterson@cs.berkeley.edu`

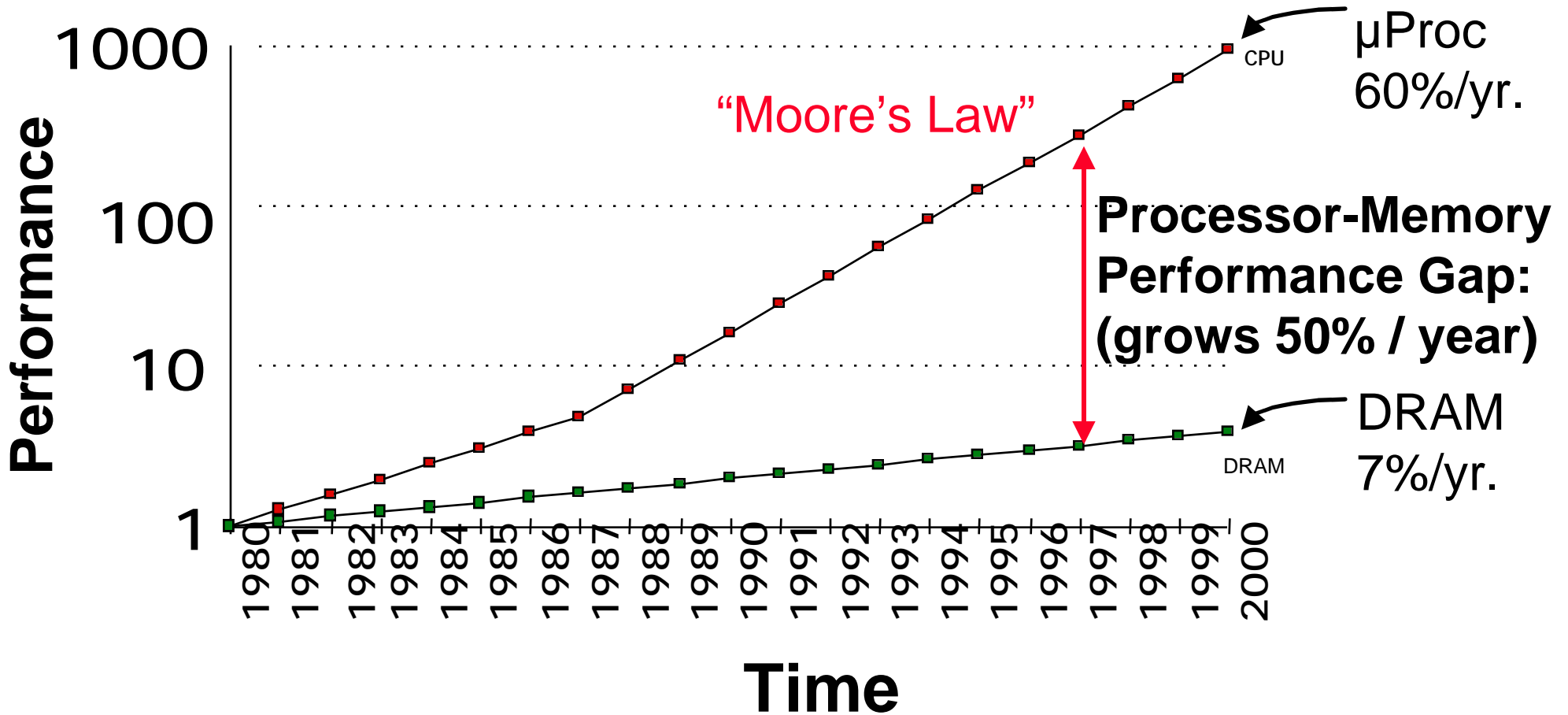
EECS, University of California

Berkeley, CA 94720-1776

Outline

- Desktop/Server Microprocessor State of the Art
- Mobile Multimedia Computing as New Direction
- A New Architecture for Mobile Multimedia Computing
- A New Technology for Mobile Multimedia Computing
- Berkeley's Mobile Multimedia Microprocessor
- Radical Bonus Application
- Challenges & Potential Industrial Impact

Processor-DRAM Gap (latency)



Processor-Memory Performance Gap “Tax”

Processor	% Area (<i>cost</i>)	%Transistors (<i>power</i>)
■ Alpha 21164	37%	77%
■ StrongArm SA110	61%	94%
■ Pentium Pro	64%	88%
– 2 dies per package: Proc/I\$/D\$ + L2\$		
■ Caches have no inherent value, only try to close performance gap		

Today's Situation: Microprocessor

- Microprocessor-DRAM performance gap
 - time of a full cache miss in instructions executed
 - 1st Alpha (7000): $340 \text{ ns} / 5.0 \text{ ns} = 68 \text{ clks} \times 2$ or 136
 - 2nd Alpha (8400): $266 \text{ ns} / 3.3 \text{ ns} = 80 \text{ clks} \times 4$ or 320
 - 3rd Alpha (t.b.d.): $180 \text{ ns} / 1.7 \text{ ns} = 108 \text{ clks} \times 6$ or 648
 - 1/2X latency x 3X clock rate x 3X Instr/clock 5X
- Benchmarks: SPEC, TPC-C, TPC-D
 - Benchmark highest optimization, ship lowest optimization?
 - Applications of past to design computers of future?

Today's Situation: Microprocessor

MIPS MPUs	R5000	R10000	10k/5k
■ Clock Rate	200 MHz	195 MHz	1.0x
■ On-Chip Caches	32K/32K	32K/32K	1.0x
■ Instructions/Cycle	1(+ FP)	4	4.0x
■ Pipe stages	5	5-7	1.2x
■ Model	In-order	Out-of-order	---
■ Die Size (mm ²)	84	298	3.5x
– without cache, TLB	32	205	6.3x
■ Development (man yr.)	60	300	5.0x
■ SPECint_base95	5.7	8.8	1.6x

Challenge for Future Microprocessors

- “...wires are not keeping pace with scaling of other features. ... In fact, for CMOS processes below 0.25 micron ... *an unacceptably small percentage of the die will be reachable during a single clock cycle.*”
- *“Architectures that require long-distance, rapid interaction will not scale well ...”*
 - “Will Physical Scalability Sabotage Performance Gains?” Matzke, *IEEE Computer* (9/97)

Billion Transistor Architectures and “Stationary Computer” Metrics

	SS++	Trace	SMT	CMP	IA-64	RAW
SPEC Int	+	+	+	=	+	=
SPEC FP	+	+	+	+	+	=
TPC (DataBse)	=	=	+	+	=	-
SW Effort	+	+	=	=	=	-
Design Scal.	-	=	-	=	=	=
Physical Design Complexity	-	=	-	=	=	+

(See *IEEE Computer* (9/97), Special Issue on
Billion Transistor Microprocessors)

Desktop/Server State of the Art

- Primary focus of architecture research last 15 years
- Processor performance doubling / 18 months
 - assuming SPEC compiler optimization levels
- Growing MPU-DRAM performance gap & tax
- Cost \$200-\$500/chip, power whatever can cool
 - 10X cost, 10X power => 2X integer performance?
- Desktop apps slow at rate processors speedup?
- Consolidation of stationary computer industry?

PA-RISC

MIPS

PowerPC

Alpha

SPARC

IA-64

Outline

- Desktop/Server Microprocessor State of the Art
- Mobile Multimedia Computing as New Direction
- A New Architecture for Mobile Multimedia Computing
- A New Technology for Mobile Multimedia Computing
- Berkeley's Mobile Multimedia Microprocessor
- Radical Bonus Application
- Challenges & Potential Industrial Impact

Intelligent PDA (2003?)

Pilot PDA

- + gameboy, cell phone, radio, timer, camera, TV remote, am/fm radio, garage door opener, ...
- + Wireless data (WWW)
- + Speech, vision recog.
- + Voice output for conversations



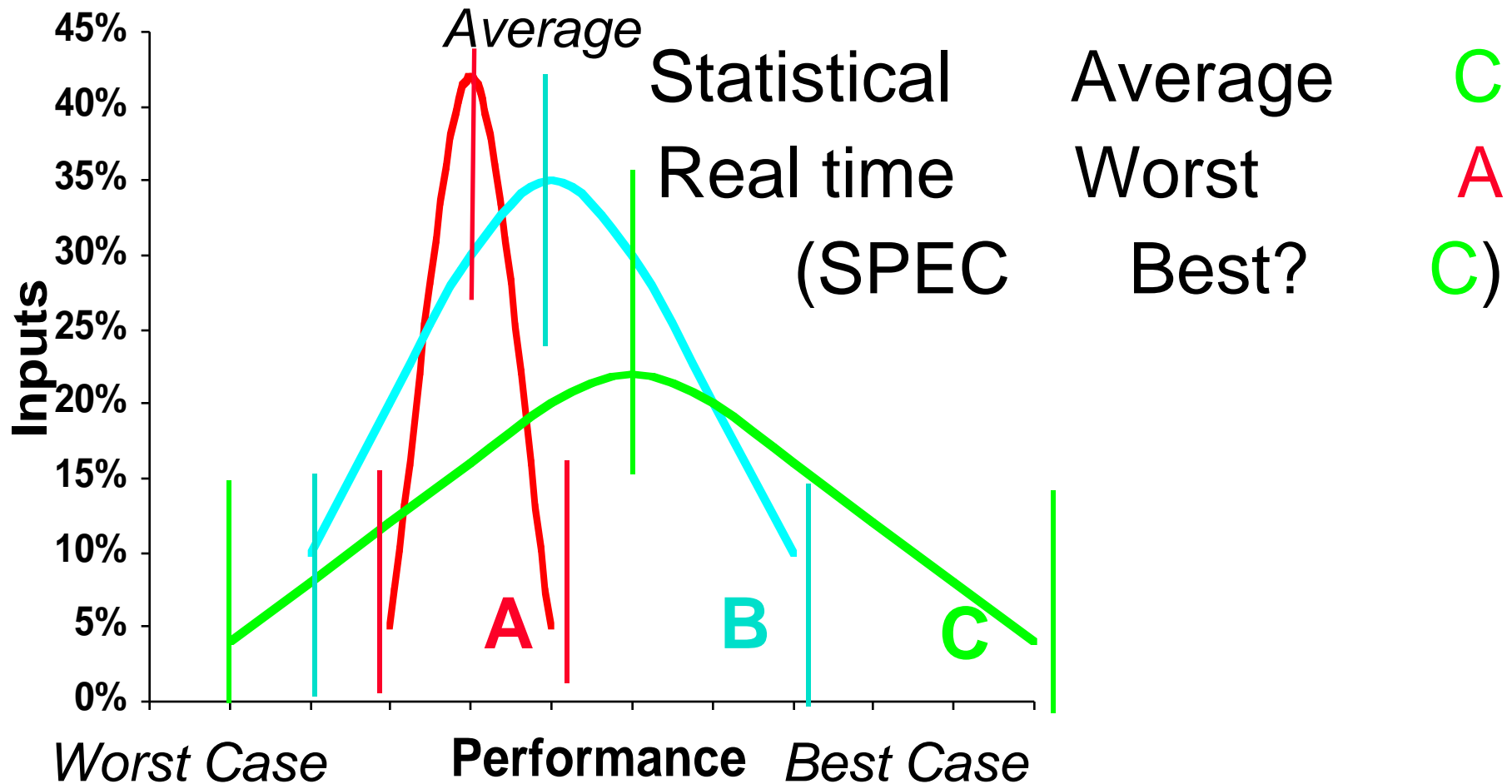
- Speech control of all devices
- Vision to see surroundings, scan documents, read bar code, measure room, ...

New Architecture Directions

- “...media processing will become the dominant force in computer arch. & microprocessor design.”
- “... new media-rich applications... involve significant real-time processing of continuous media streams, and make heavy use of vectors of packed 8-, 16-, and 32-bit integer and Fl. Pt.”
- Needs include real-time response, continuous media data types (no temporal locality), fine grain parallelism, coarse grain parallelism, memory BW
 - “How Multimedia Workloads Will Change Processor Design”, Diefendorff & Dubey, *IEEE Computer* (9/97)

Which is Faster?

Statistical v. Real time v. SPEC



Billion Transistor Architectures and “Mobile Multimedia” Metrics

SS++ Trace SMT CMP IA-64 RAW

Design Scal.	-	=	-	=	=	=
Energy/power	-	-	-	=	=	-
Code Size	=	=	=	=	-	=
Real-time	-	-	=	=	=	=
Cont. Data	=	=	=	=	=	=
Memory BW	=	=	=	=	=	=
Fine-grain Par.	=	=	=	=	=	+
Coarse-gr.Par.	=	=	+	+	=	+

Outline

- Desktop/Server Microprocessor State of the Art
- Mobile Multimedia Computing as New Direction
- A New Architecture for Mobile Multimedia Computing
- A New Technology for Mobile Multimedia Computing
- Berkeley's Mobile Multimedia Microprocessor
- Radical Bonus Application
- Challenges & Potential Industrial Impact

Potential Multimedia Architecture

- “New” model: VSIW=Very Short Instruction Word!
 - Compact: Describe N operations with 1 short instruct.
 - Predictable (real-time) perf. vs. statistical perf. (cache)
 - Multimedia ready: choose $N*64b$, $2N*32b$, $4N*16b$
 - Easy to get high performance; N operations:
 - » are independent
 - » use same functional unit
 - » access disjoint registers
 - » access registers in same order as previous instructions
 - » access contiguous memory words or known pattern
 - » hides memory latency (and any other latency)
 - Compiler technology already developed, for sale!

Operation & Instruction Count: RISC v. “VSIW” Processor

(from F. Quintana, U. Barcelona.)

Spec92fp Program	Operations (M)			Instructions (M)		
	RISC	VSIW	R / V	RISC	VSIW	R / V
swim256	115	95	1.1x	115	0.8	142x
hydro2d	58	40	1.4x	58	0.8	71x
nasa7	69	41	1.7x	69	2.2	31x
su2cor	51	35	1.4x	51	1.8	29x
tomcatv	15	10	1.4x	15	1.3	11x
wave5	27	25	1.1x	27	7.2	4x
mdljdp2	32	52	0.6x	32	15.8	2x

VSIW reduces ops by 1.2X, instructions by 20X!

Revive Vector (= VSIW) Architecture!

- Cost: \$1M each?
- Low latency, high BW memory system?
- Code density?
- Compilers?
- Vector Performance?
- Power/Energy?
- Scalar performance?
- Real-time?
- Limited to scientific applications?
- Single-chip CMOS MPU/IRAM
- ? (new media?)
- Much smaller than VLIW/EPIC
- For sale, mature (>20 years)
- Easy scale speed with technology
- Parallel to save energy, keep perf
- Include modern, modest CPU
OK scalar (MIPS 5K v. 10k)
- No caches, no speculation
repeatable speed as vary input
- Multimedia apps vectorizable too:
N*64b, 2N*32b, 4N*16b

Vector Surprise

- Use vectors for inner loop parallelism (no surprise)
 - One dimension of array: $A[0, \underline{0}]$, $A[0, \underline{1}]$, $A[0, \underline{2}]$, ...
 - think of machine as 32 vector regs each with 64 elements
 - 1 instruction updates 64 elements of 1 vector register
- and for outer loop parallelism!
 - 1 element from each column: $A[\underline{0}, 0]$, $A[\underline{1}, 0]$, $A[\underline{2}, 0]$, ...
 - think of machine as 64 “virtual processors” (VPs) each with 32 scalar registers! (multithreaded processor)
 - 1 instruction updates 1 scalar register in 64 VPs
- Hardware identical, just 2 compiler perspectives

Vector Multiply with dependency

```
/* Multiply a[m][k] * b[k][n] to get
   c[m][n] */
for (i=1; i<m; i++)
{
    for (j=1; j<n; j++)
    {
        sum = 0;
        for (t=1; t<k; t++)
        {
            sum += a[i][t] * b[t][j];
        }
        c[i][j] = sum;
    }
}
```

Novel Matrix Multiply Solution

- You don't need to do reductions for matrix multiply
- You can calculate multiple independent sums within one vector register
- You can vectorize the outer (j) loop to perform 32 dot-products at the same time
- Or you can think of each 32 Virtual Processors doing one of the dot products
 - (Assume Maximum Vector Length is 32)
- Show it in C source code, but can imagine the assembly vector instructions from it

Optimized Vector Example

```
/* Multiply a[m][k] * b[k][n] to get c[m][n] */
for (i=1; i<m; i++)
{
  for (j=1; j<n; j+=32)//* Step j 32 at a time. */
  {
    sum[0:31] = 0; /* Initialize a vector
                    register to zeros. */
    for (t=1; t<k; t++)
    {
      a_scalar = a[i][t]; /* Get scalar from
                          a matrix. */
      b_vector[0:31] = b[t][j:j+31];
                          /* Get vector from
                          b matrix. */
      prod[0:31] = b_vector[0:31]*a_scalar;
      /* Do a vector-scalar multiply. */
    }
  }
}
```

Optimized Vector Example cont'd

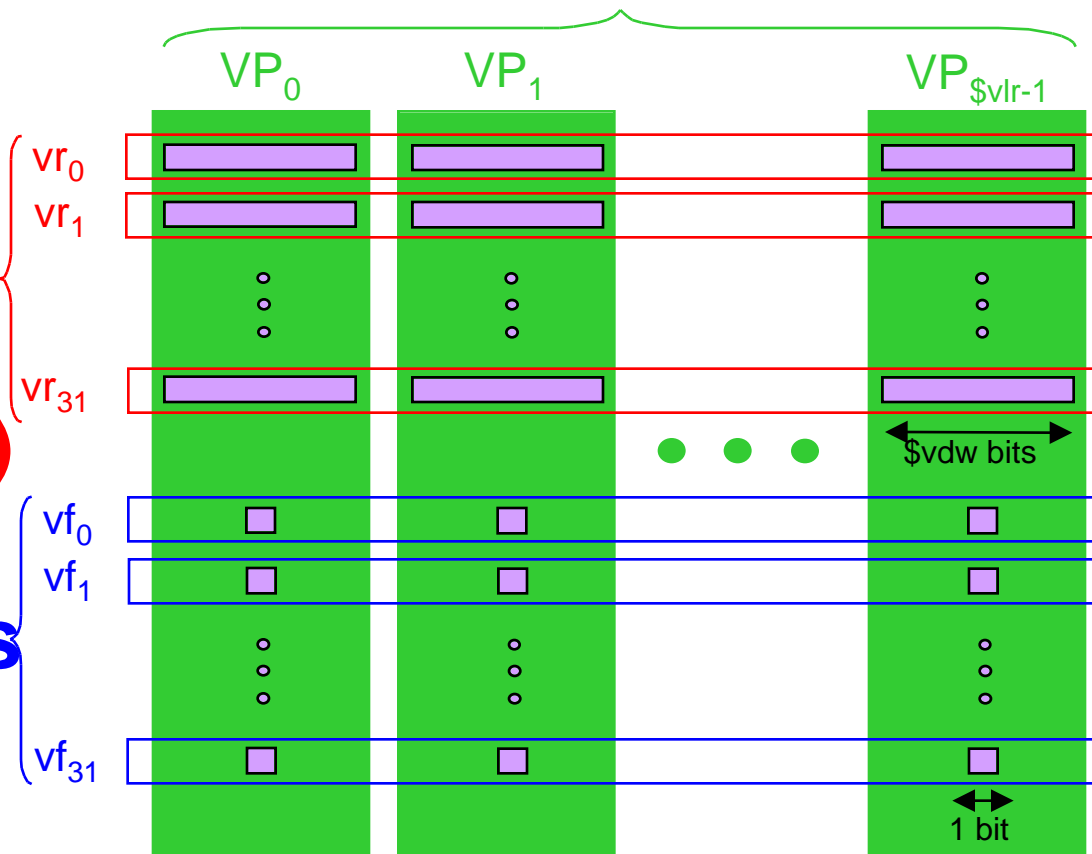
```
        /* Vector-vector add into results. */  
        sum[0:31] += prod[0:31];  
    }  
  
    /* Unit-stride store of vector of  
       results. */  
    c[i][j:j+31] = sum[0:31];  
} }  
}
```

Vector Multimedia Architectural State

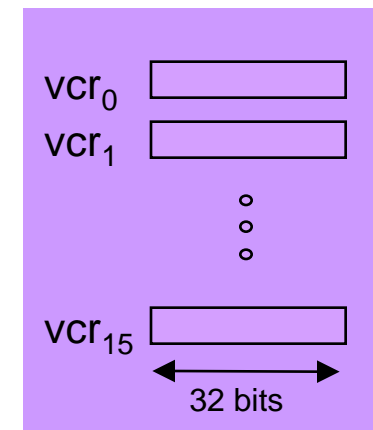
General Purpose Registers
(32 x 32/64/128 x 64/32/16)

Flag Registers
(32 x 128 x 1)

Virtual Processors (\$vlr)



Control Registers



Vector Multimedia Instruction Set

Scalar Standard scalar instruction set (e.g., ARM, MIPS)

Vector ALU

$\left\{ \begin{array}{c} + \\ - \\ \times \\ \div \\ \& \\ \\ \text{shl} \\ \text{shr} \end{array} \right\}$	$\left\{ \begin{array}{c} \text{s.int} \\ \text{u.int} \\ \text{s.fp} \\ \text{d.fp} \end{array} \right\}$	$\left\{ \begin{array}{c} 8 \\ 16 \\ 32 \\ 64 \end{array} \right\}$	$\left\{ \begin{array}{c} \text{.VV} \\ \text{.VS} \\ \text{.SV} \end{array} \right\}$	$\left\{ \begin{array}{c} \text{saturate} \\ \text{overflow} \end{array} \right\}$	$\left\{ \begin{array}{c} \text{masked} \\ \text{unmasked} \end{array} \right\}$
-----------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------	----------------------------------------------------------------------------------------	------------------------------------------------------------------------------------	----------------------------------------------------------------------------------

Vector Memory

$\left\{ \begin{array}{c} \text{load} \\ \text{store} \end{array} \right\}$	$\left\{ \begin{array}{c} \text{s.int} \\ \text{u.int} \end{array} \right\}$	$\left\{ \begin{array}{c} 8 \\ 16 \\ 32 \\ 64 \end{array} \right\}$	$\left\{ \begin{array}{c} 8 \\ 16 \\ 32 \\ 64 \end{array} \right\}$	$\left\{ \begin{array}{c} \text{unit} \\ \text{constant} \\ \text{indexed} \end{array} \right\}$	$\left\{ \begin{array}{c} \text{masked} \\ \text{unmasked} \end{array} \right\}$
-----------------------------------------------------------------------------	------------------------------------------------------------------------------	---------------------------------------------------------------------	---------------------------------------------------------------------	--------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------

Vector Registers 32 x 32 x 64b (or 32 x 64 x 32b or 32 x 128 x 16b)
 + 32 x 128 x 1b flag

Plus: **flag**, **convert**, **DSP**, and **transfer** operations

Software Technology Trends Affecting New Direction?

- any CPU + vector coprocessor/memory
 - scalar/vector interactions are limited, simple
 - Example architecture based on ARM 9, MIPS
- Vectorizing compilers built for 25 years
 - can buy one for new machine from The Portland Group
- Microsoft “Win CE”/ Java OS for non-x86 platforms
- Library solutions (e.g., MMX); retarget packages
- Software distribution model is evolving?
 - New Model: Java byte codes over network?
 - + Just-In-Time compiler to tailor program to machine?

Outline

- Desktop/Server Microprocessor State of the Art
- Mobile Multimedia Computing as New Direction
- A New Architecture for Mobile Multimedia Computing
- A New Technology for Mobile Multimedia Computing
- Berkeley's Mobile Multimedia Microprocessor
- Radical Bonus Application
- Challenges & Potential Industrial Impact

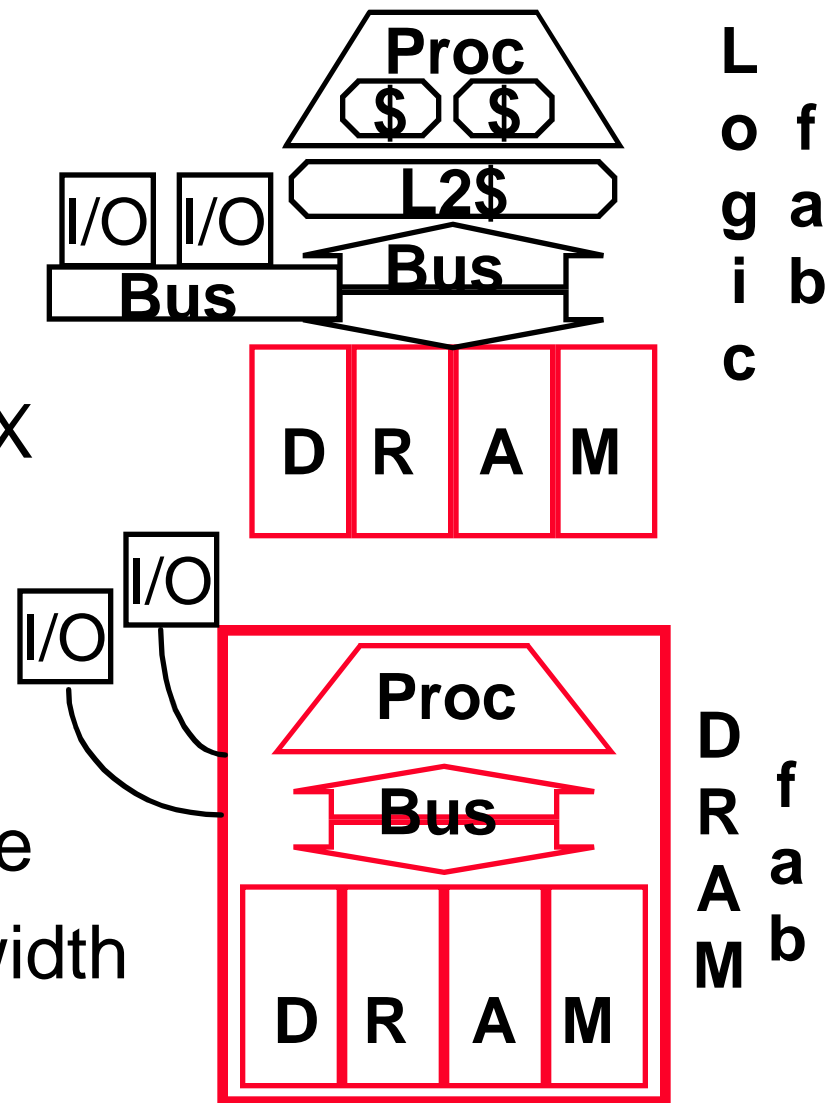
A Better Media for Mobile Multimedia MPUs: Logic+DRAM

- Crash of DRAM market inspires new use of wafers
- Faster logic in DRAM process
 - DRAM vendors offer faster transistors + same number metal layers as good logic process?
@ 20% higher cost per wafer?
 - As die cost $f(\text{die area}^4)$, 4% die shrink equal cost
- Called **Intelligent RAM** (“**IRAM**”) since most of transistors will be DRAM

IRAM Vision Statement

Microprocessor & DRAM on a single chip:

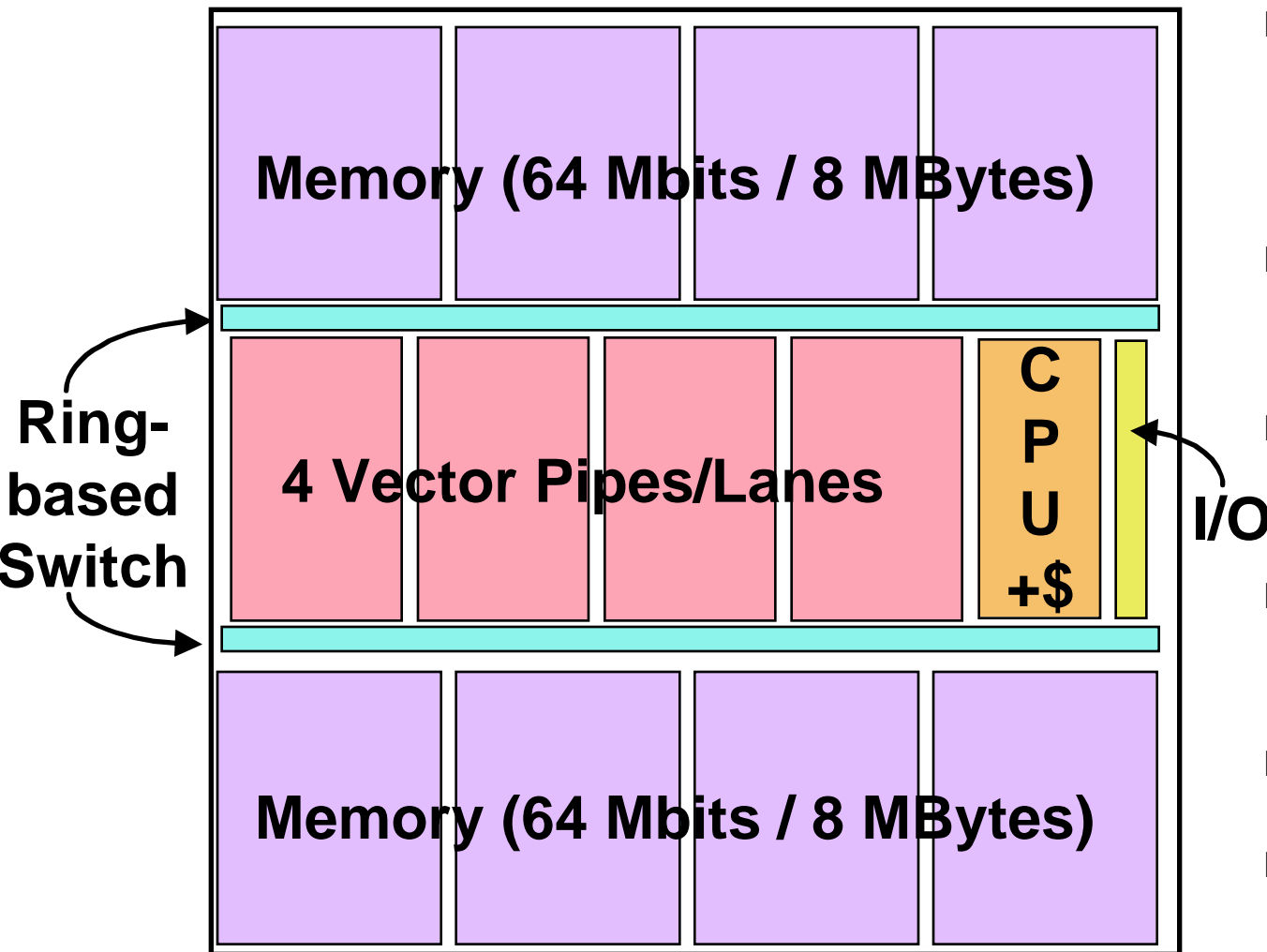
- on-chip memory latency 5-10X, bandwidth 50-100X
- improve energy efficiency 2X-4X (no off-chip bus)
- serial I/O 5-10X v. buses
- smaller board area/volume
- adjustable memory size/width



Outline

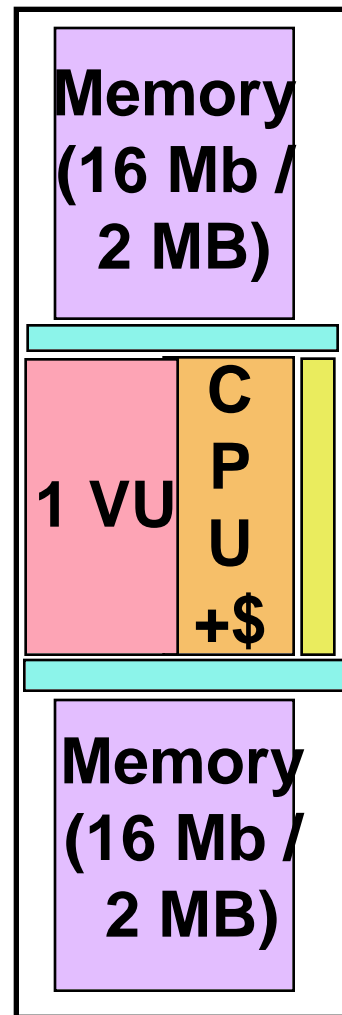
- Desktop/Server Microprocessor State of the Art
- Mobile Multimedia Computing as New Direction
- A New Architecture for Mobile Multimedia Computing
- A New Technology for Mobile Multimedia Computing
- Berkeley's Mobile Multimedia Microprocessor
- Radical Bonus Application
- Challenges & Potential Industrial Impact

Tentative VIRAM-1 Floorplan



- 0.25 μm DRAM
16 MB in 8 banks x 256b, 64 subbanks
- 0.25 μm ,
5 Metal Logic
- 200 MHz MIPS IV,
16K I\$, 16K D\$
- 4 200 MHz
FP/int. vector units
- die: 20x20 mm
- xtors: 130M
- power: 2 Watts

Tentative VIRAM-"0.25" Floorplan



- Demonstrate scalability via 2nd layout (automatic from 1st)
- 4 MB in 2 banks x 256b, 32 subbanks
- 200 MHz CPU, 8K I\$, 8K D\$
- 1 200 MHz FP/int. vector units
- die: 5 x 20 mm
- xtors: 35M
- power: 0.5 Watts ³³

VIRAM-1 Specs/Goals

Technology	0.18-0.25 micron, 5-6 metal layers, fast xtor
Memory	16-32 MB
Die size	250-400 mm²
Vector pipes/lanes	4 64-bit (or 8 32-bit or 16 16-bit)
Serial I/O	4 lines @ 1 Gbit/s
Power _{university}	2 w @ 1-1.5 volt logic
Clock _{university}	200scalar/200vector MHz
Perf _{university}	1.6 GFLOPS₆₄ – 6 GOPS₁₆
<hr/>	
Power _{industry}	1 w @ 1-1.5 volt logic
Clock _{industry}	400scalar/400vector MHz
Perf _{industry}	3.2 GFLOPS₆₄ – 12 GOPS₁₆

2X



V-IRAM-1 Tentative Plan

- Phase I: Feasibility stage (H2'98)
 - Test chip, CAD agreement, architecture defined
- Phase 2: Design & Layout Stage (H1'99)
 - Test chip, Simulated design and layout
- Phase 3: Verification (H2'99)
 - Tape-out
- Phase 4: Fabrication, Testing, and Demonstration (H1'00)
 - Functional integrated circuit
- 100M transistor microprocessor before Intel?

Grading VIRAM

Stationary Metrics

Mobile Multimedia Metrics

	VIRAM		VIRAM
SPEC Int	-	Energy/power	+
SPEC FP	+	Code Size	+
TPC (DataBse)	-	Real-time response	+
SW Effort	=	Continous Data-types	+
Design Scal.	+	Memory BW	+
Physical	=	Fine-grain Parallelism	+
Design Complexity		Coarse-gr. Parallelism	=

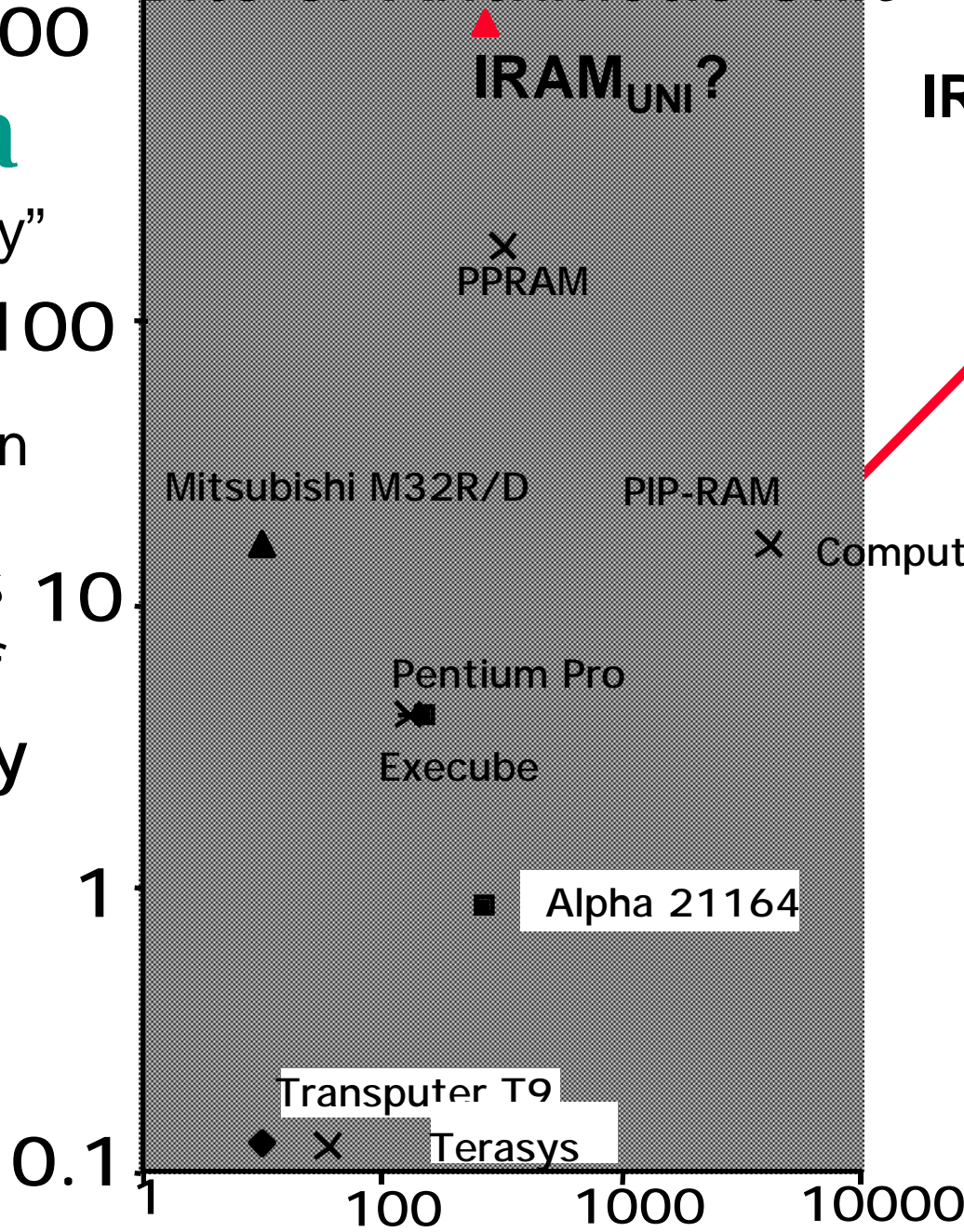
IRAM

not a new idea

- Stone, '70 "Logic-in memory"
- Barron, '78 "Transputer"
- Dally, '90 "J-machine"
- Patterson, '90 panel session
- Kogge, '94 "Execube"

Mbits of Memory

Bits of Arithmetic Unit



- × SIMD on chip (DRAM)
- Uniprocessor (SRAM)
- × MIMD on chip (DRAM)
- ▲ Uniprocessor (DRAM)
- ◆ MIMD component (SRAM)

Why IRAM now?

Lower risk than before

- Faster Logic + DRAM available now/soon
- DRAM manufacturers now willing to listen
 - Before not interested, so early IRAM = SRAM
- Past efforts memory limited multiple chips
 - 1st solve the unsolved (parallel processing)
- Gigabit DRAM 100 MB; OK for many apps?
- Systems headed to 2 chips: CPU + memory
- Embedded apps leverage energy efficiency, adjustable mem. capacity, smaller board area
 - OK market v. desktop (55M 32b RISC '96)

IRAM Challenges

■ Chip

- Good performance and reasonable power?
- Speed, area, power, yield, cost in embedded DRAM process? (time delay vs. state-of-the-art logic, DRAM)
- Testing time of IRAM vs DRAM vs microprocessor?

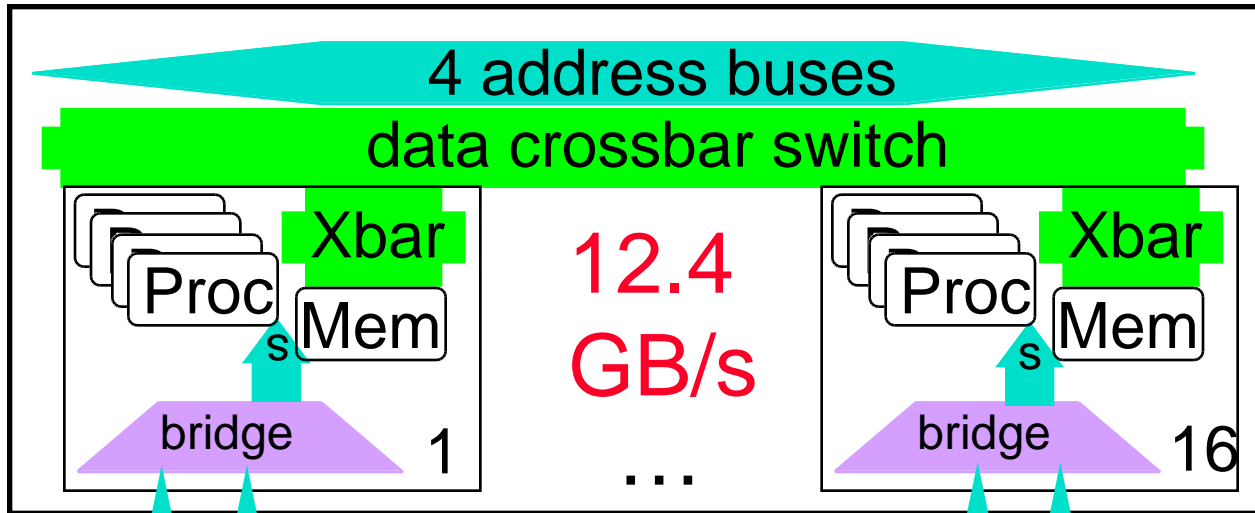
■ Architecture

- How to turn high memory bandwidth into performance for real applications?
- Extensible IRAM: Large program/data solution? (e.g., external DRAM, clusters, CC-NUMA, IDISK ...)

Outline

- Desktop/Server Microprocessor State of the Art
- Mobile Multimedia Computing as New Direction
- A New Architecture for Mobile Multimedia Computing
- A New Technology for Mobile Multimedia Computing
- Berkeley's Mobile Multimedia Microprocessor
- Radical Bonus Application
- Challenges & Potential Industrial Impact

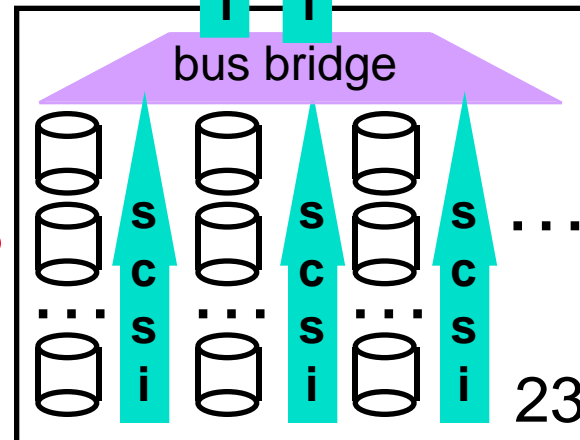
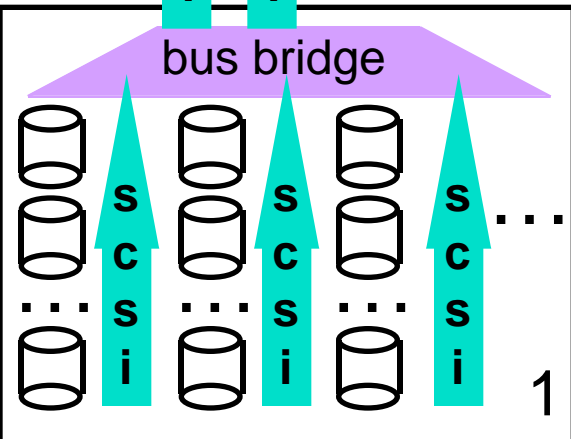
Revolutionary App: Decision Support?



12.4
GB/s

2.6
GB/s

6.0
GB/s



Sun 10000 (Oracle 8):

- TPC-D (1TB) leader
- SMP 64 CPUs, 64GB dram, 603 disks

Disks,encl. \$2,348k

DRAM \$2,328k

Boards,encl. \$983k

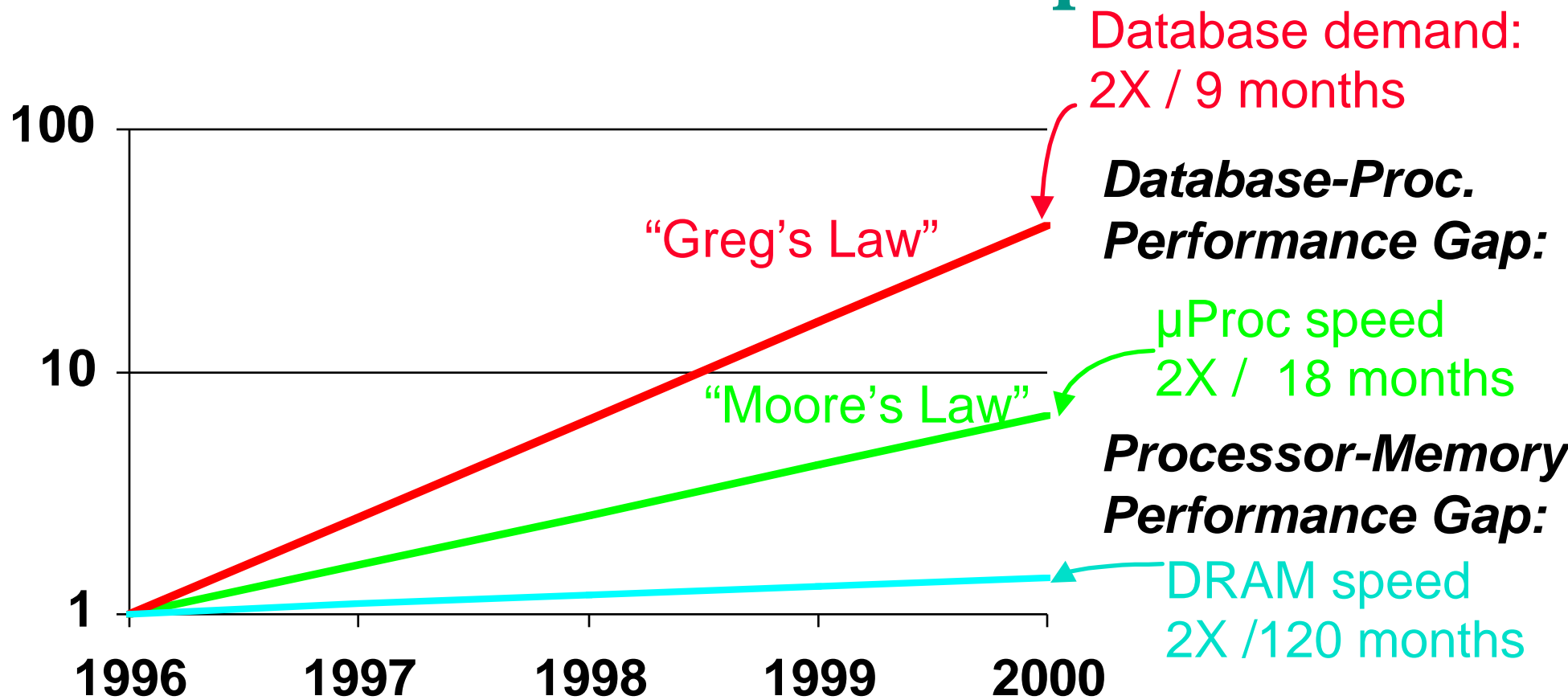
CPUs \$912k

Cables,I/O \$139k

Misc. \$65k

HW total \$6,775k

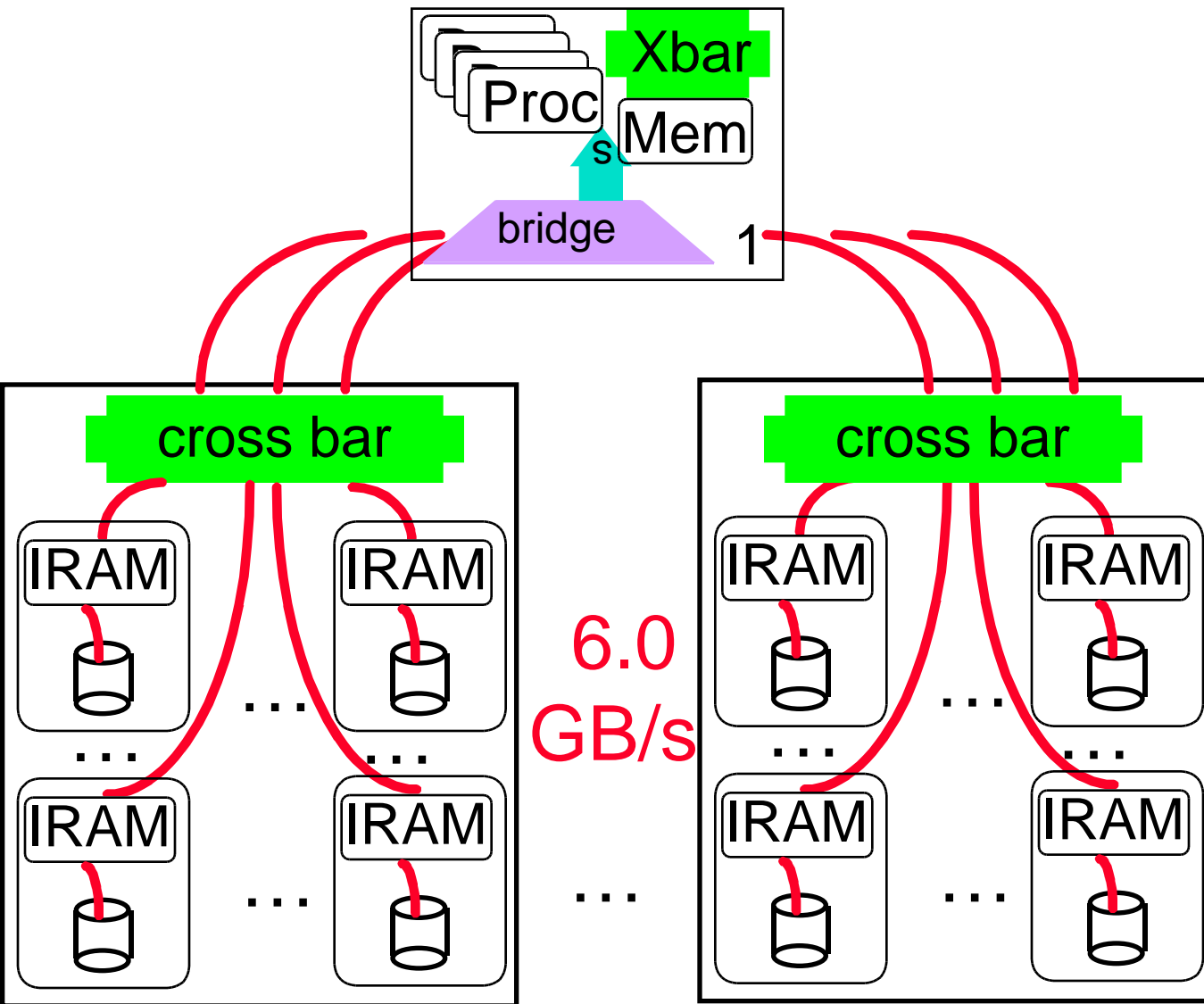
IRAM Application Inspiration: Database Demand vs. Processor/DRAM speed



IRAM Application Inspiration: Cost of Ownership

- Annual system administration cost:
3X - 8X cost of disk (!)
- Current computer generation
emphasizes cost-performance,
neglects cost of use, ease of use

App #2: “Intelligent Storage” (ISTORE): Scaleable Decision Support?



- 1 IRAM/disk + xbar + fast serial link v. conventional SMP
- Network latency = f(SW overhead), not link distance
- Move function to data v. data to CPU (scan, sort, join,...)
- Cheaper, more scalable (1/3 \$, 3X perf)

Mobile Multimedia Conclusion

- 10000X cost-performance increase in “stationary” computers, consolidation of industry
=> time for architecture/OS/compiler researchers declare victory, search for new horizons?
- Mobile Multimedia offer many new challenges: energy efficiency, size, real time performance, ...
- VIRAM-1 one example, hope others will follow
- Apps/metrics of future to design computer of future!
 - Suppose PDA replaces desktop as primary computer?
 - Work on FPPP on PC vs. Speech on PDA?

Infrastructure for Next Generation

- Applications of ISTORE systems
 - Database-powered information appliances providing data-intensive services over WWW
 - » decision support, data mining, rent-a-server, ...
- Lego-like model of system design gives advantages in administration, scalability
 - HW+SW for self-maintenance, self-tuning
 - Configured to match resource needs of workload
 - Easily adapted/scaled to changes in workload

IRAM Conclusion

- IRAM potential in mem/IO BW, energy, board area; challenges in power/performance, testing, yield
- 10X-100X improvements based on technology shipping for 20 years (not JJ, photons, MEMS, ...)
- Suppose IRAM is successful
- **Revolution in computer implementation v. Instr Set**
 - Potential Impact #1: turn server industry inside-out?
- **Potential #2: shift semiconductor balance of power?**
 - Who ships the most memory? Most microprocessors?

Interested in Participating?

- Looking for ideas of VIRAM enabled apps
- Contact us if you're interested:
email: patterson@cs.berkeley.edu
<http://iram.cs.berkeley.edu/>
 - iram.cs.berkeley.edu/papers/direction/paper.html
- Thanks for advice/support: DARPA, California MICRO, Hitachi, IBM, Intel, LG Semicon, Microsoft, Neomagic, Sandcraft, SGI/Cray, Sun Microsystems, TI, TSMC

IRAM Project Team

Jim Beck, Aaron Brown,
Ben Gribstad, Richard Fromm,
Joe Gebis, Jason Golbus, Kimberly Keeton,
Christoforos Kozyrakis, John Kubiatoicz,
David Martin, Morley Mao, David Oppenheimer,
David Patterson, Steve Pope, Randi Thomas,
Noah Treuhaft, and Katherine Yelick

Backup Slides

(The following slides are used to help answer questions)

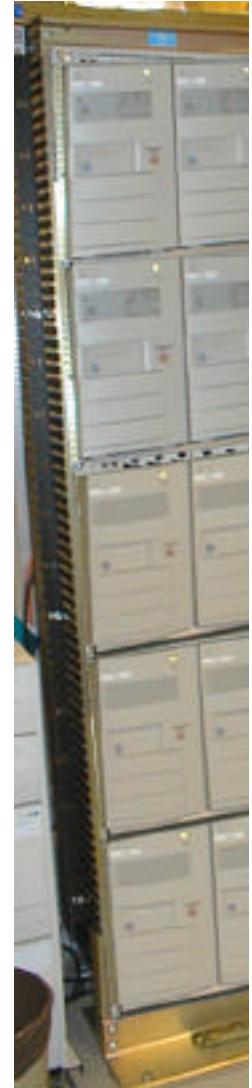
ISTORE Cluster?

- 8 disks / enclosure
- 15 enclosures /rack = 120 disks/rack



Cluster of PCs?

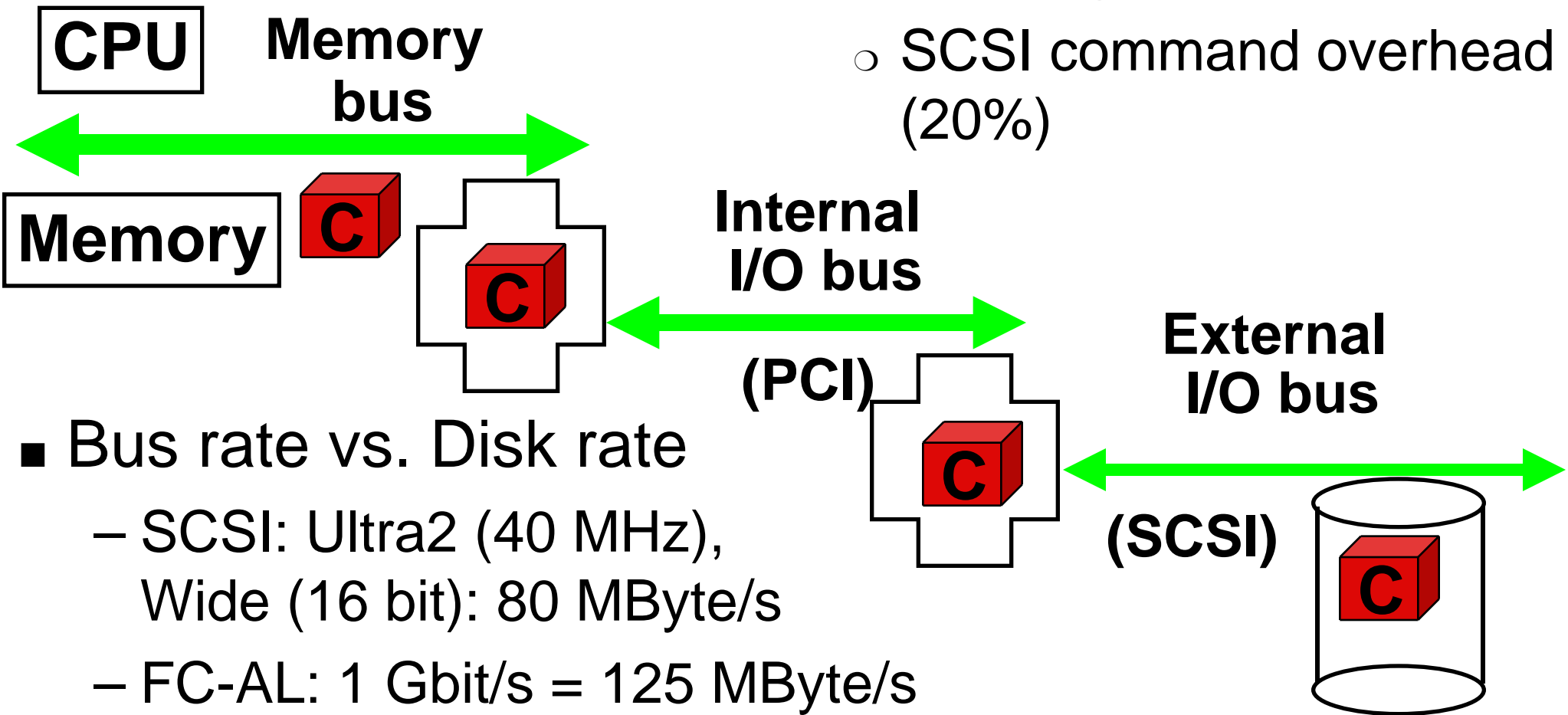
- 2 disks / PC
- 10 PCs /rack = 20 disks/rack
- Quality of Equipment?
- Ease of Repair?
- System Admin.?



Disk Limit: I/O Buses

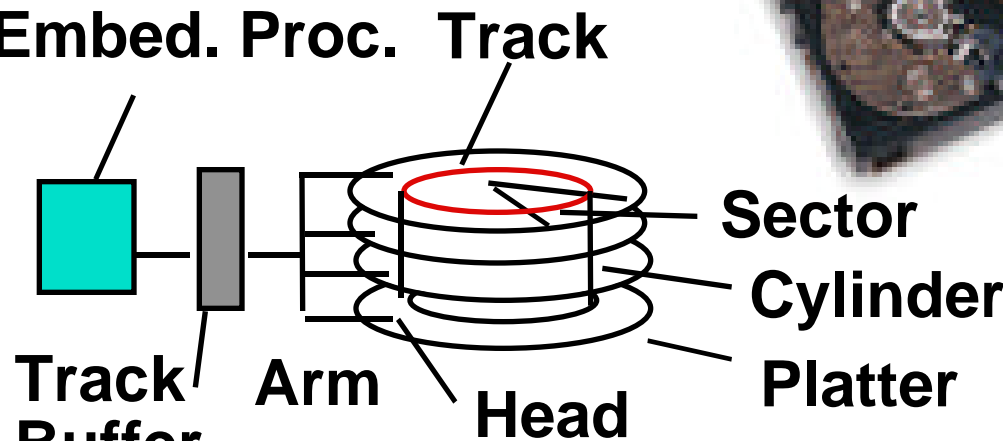
- Multiple copies of data
- Cannot use 100% of bus

- Queuing Theory (< 70%)
- SCSI command overhead (20%)



- Bus rate vs. Disk rate
 - SCSI: Ultra2 (40 MHz), Wide (16 bit): 80 MByte/s
 - FC-AL: 1 Gbit/s = 125 MByte/s (single disk in 2002)

State of the Art: Seagate Cheetah 18



Latency =

per access { **Queuing Time +**
Controller time +
Seek Time +

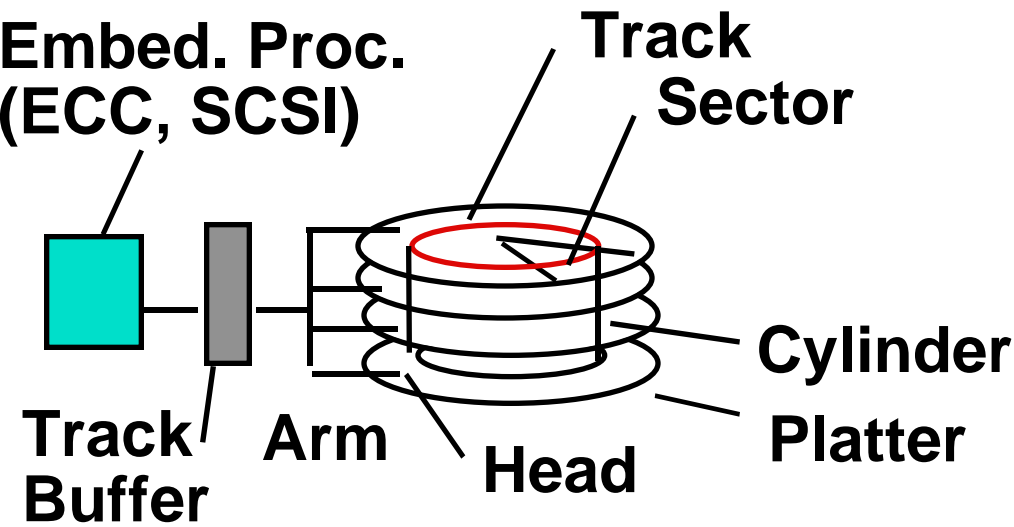
+

per byte { **Rotation Time +**
Size / Bandwidth

- 18.2 GB, 3.5 inch disk
- \$1647 or 11MB/\$ (9¢/MB)
- 1MB track buffer (+ 4MB optional expansion)
- 6962 cylinders, 12 platters
- 19 watts
- 0.15 ms controller time
- 6 ms avg. seek (seek 1 track => 1 ms)
- 3 ms = 1/2 rotation
- 21 to 15 MB/s media (x 75% => 16 to 11 MB/s)

source: www.seagate.com;
www.pricewatch.com; 5/21/98

Description/Trends



Latency =

per access { Queuing Time +
Controller time +
Seek Time +
per byte { Rotation Time +
Size / Bandwidth

- Capacity
 - + 60%/year (2X / 1.5 yrs)
- MB/\$
 - > 60%/year (2X / <1.5 yrs)
 - Fewer chips + areal density
- Rotation + Seek time
 - – 8%/ year (1/2 in 10 yrs)
- Transfer rate (BW)
 - + 40%/year (2X / 2.0 yrs)
 - deliver 75% of quoted rate (ECC, gaps, servo...)

source: Ed Grochowski, 1996,
 "IBM leadership in disk drive technology";
www.storage.ibm.com/storage/technolo/grochows/grocho01.htm,

Vectors Lower Power

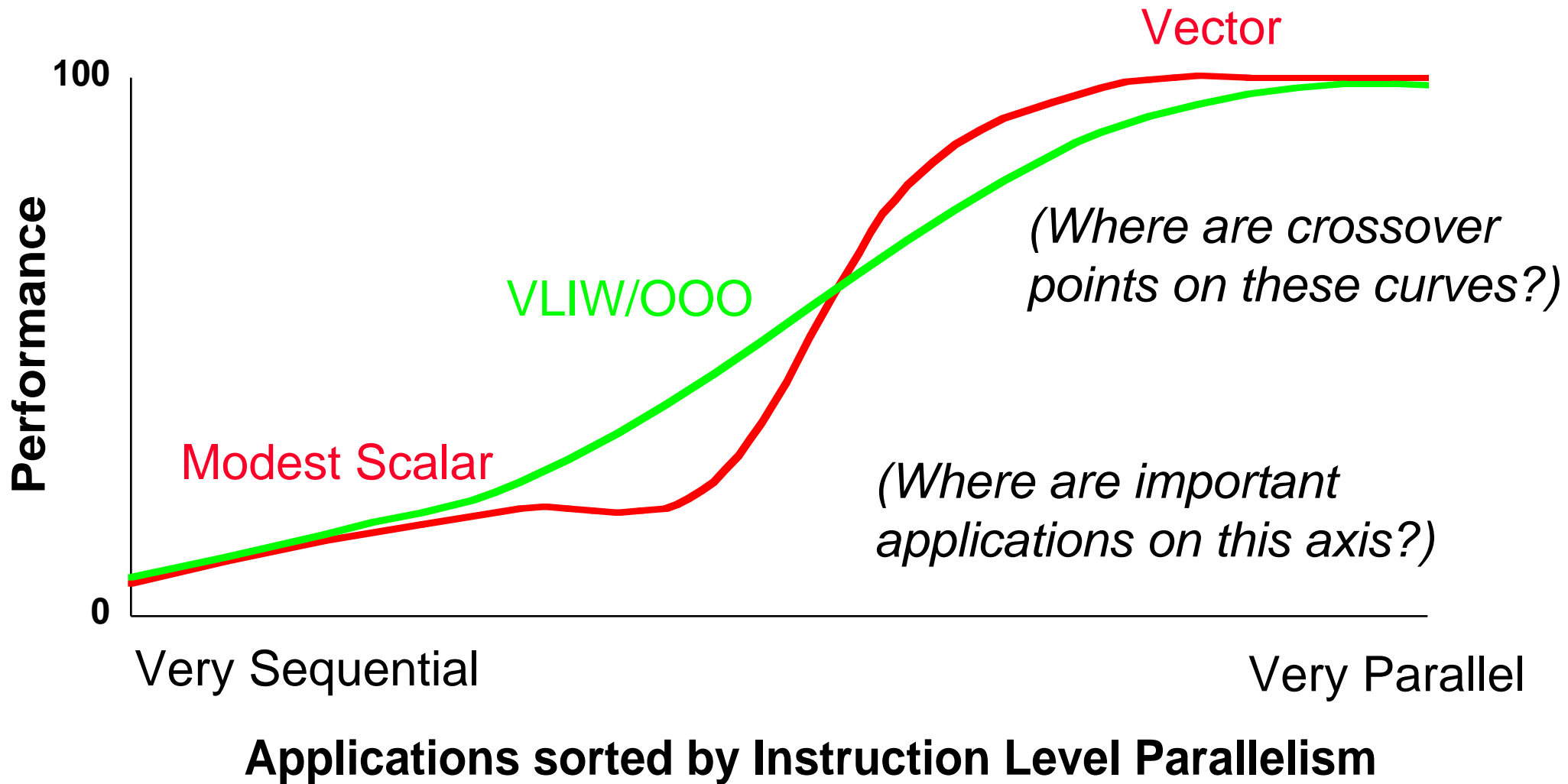
Single-issue Scalar

- One instruction fetch, decode, dispatch per operation
- Arbitrary register accesses, adds area and power
- Loop unrolling and software pipelining for high performance increases instruction cache footprint
- All data passes through cache; waste power if no temporal locality
- One TLB lookup per load or store
- Off-chip access in whole cache lines

Vector

- One instruction fetch, decode, dispatch per vector
- Structured register accesses
- Smaller code for high performance, less power in instruction cache misses
- Bypass cache
- One TLB lookup per group of loads or stores
- Move only necessary data across chip boundary

VLIW/Out-of-Order vs. Modest Scalar+Vector



Potential IRAM Latency: 5 - 10X

- No parallel DRAMs, memory controller, bus to turn around, SIMM module, pins...
- New focus: Latency oriented DRAM?
 - Dominant delay = RC of the word lines
 - keep wire length short & block sizes small?
- 10-30 ns for 64b-256b IRAM “RAS/CAS”?
- AlphaSta. 600: 180 ns=128b, 270 ns= 512b
Next generation (21264): 180 ns for 512b?

Potential IRAM Bandwidth: 100X

- 1024 1Mbit modules(1Gb), each 256b wide
 - 20% @ 20 ns RAS/CAS = 320 GBytes/sec
- If cross bar switch delivers 1/3 to 2/3 of BW of 20% of modules
 - 100 - 200 GBytes/sec
- FYI: AlphaServer 8400 = 1.2 GBytes/sec
 - 75 MHz, 256-bit memory bus, 4 banks

Potential Energy Efficiency: 2X-4X

- Case study of StrongARM memory hierarchy vs. IRAM memory hierarchy
 - cell size advantages much larger cache
fewer off-chip references
up to 2X-4X energy efficiency for memory
 - less energy per bit access for DRAM
- Memory cell area ratio/process: P6, '164, SArm
cache/logic : SRAM/SRAM : DRAM/DRAM
20-50 : 8-11 : 1

Potential Innovation in Standard DRAM Interfaces

- Optimizations when chip is a system vs. chip is a memory component
 - Lower power via on-demand memory module activation?
 - “Map out” bad memory modules to improve yield?
 - Improve yield with variable refresh rate?
 - Reduce test cases/testing time during manufacturing?
- IRAM advantages even greater if innovate inside DRAM memory interface?

Mediaprocesing Functions (Dubey)

Kernel

Vector length

- Matrix transpose/multiply # vertices at once
- DCT (video, comm.) image width
- FFT (audio) 256-1024
- Motion estimation (video) image width, i.w./16
- Gamma correction (video) image width
- Haar transform (media mining) image width
- Median filter (image process.) image width
- Separable convolution (“”) image width

(from <http://www.research.ibm.com/people/p/pradeep/tutor.html>)

“Architectural Issues for the 1990s”

(From Microprocessor Forum 10-10-90):

● Given:

Superscalar, superpipelined RISCs and
Amdahl's Law will not be repealed

=> High performance in 1990s is not limited by CPU

● Predictions for 1990s:

"Either/Or" CPU/Memory will disappear (*“nonblocking cache”*)

Multipronged attack on memory bottleneck
cache conscious compilers
lockup free caches / prefetching

All programs will become I/O bound; design accordingly

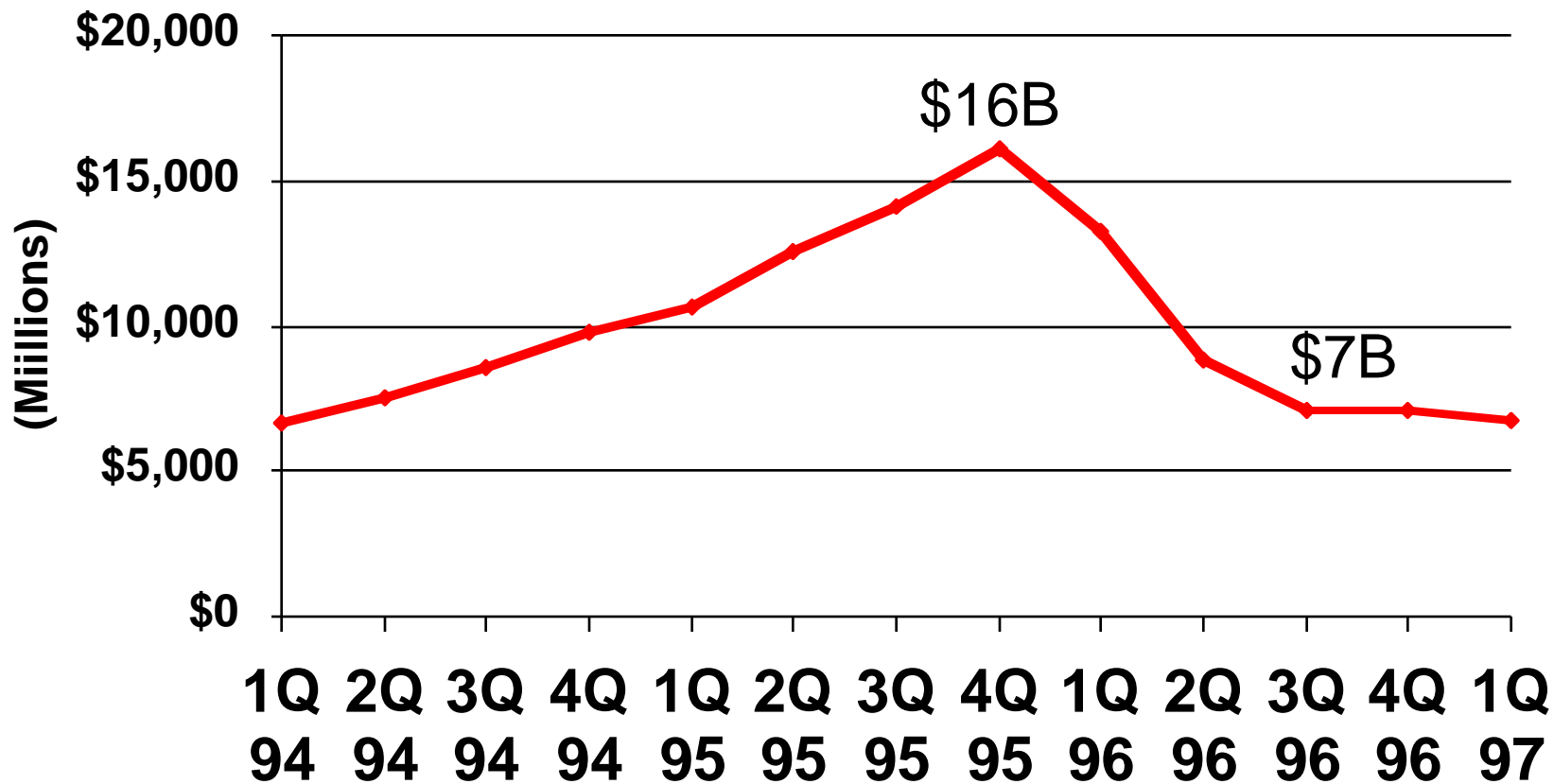
Most important CPU of 1990s is in DRAM: "IRAM"
(Intelligent RAM: 64Mb + 0.3M transistor CPU = 100.5%)
=> CPUs are genuinely free with IRAM

“Vanilla” Approach to IRAM

- Estimate performance IRAM version of Alpha (same caches, benchmarks, standard DRAM)
 - Used optimistic and pessimistic factors for logic (1.3-2.0 slower), SRAM (1.1-1.3 slower), DRAM speed (5X-10X faster) for standard DRAM
 - SPEC92 benchmark 1.2 to 1.8 times slower
 - Database 1.1 times slower to 1.1 times faster
 - Sparse matrix 1.2 to 1.8 times faster

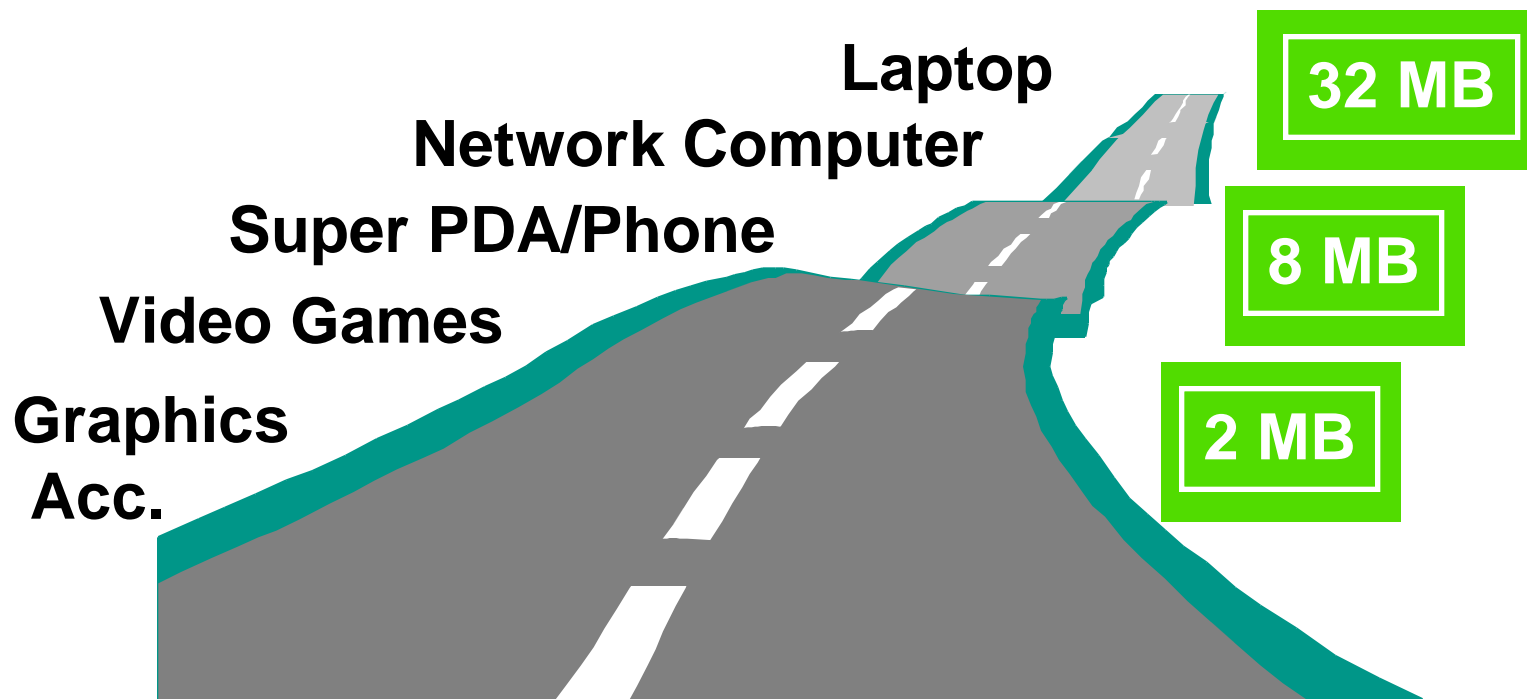
Today's Situation: DRAM

DRAM Revenue per Quarter



- Intel: 30%/year since 1987; 1/3 income profit

Commercial IRAM highway is governed by memory per IRAM?



Near-term IRAM Applications

- “Intelligent” Set-top
 - 2.6M Nintendo 64 (\$150) sold in 1st year
 - 4-chip Nintendo 1-chip: 3D graphics, sound, fun!
- “Intelligent” Personal Digital Assistant
 - 0.6M PalmPilots (\$300) sold in 1st 6 months
 - Handwriting + learn new alphabet (= K, \neg = T, L = 4)
v. Speech input

Vector Memory Operations

- Load/store operations move groups of data between registers and memory
- Three types of addressing
 - Unit stride
 - » Fastest
 - Non-unit (constant) stride
 - Indexed (gather-scatter)
 - » Vector equivalent of register indirect
 - » Good for sparse arrays of data
 - » Increases number of programs that vectorize

Variable Data Width

- Programmer thinks in terms of vectors of data of some width (16, 32, or 64 bits)
- Good for multimedia
 - More elegant than MMX-style extensions
- Shouldn't have to worry about how it is stored in memory
 - No need for explicit pack/unpack operations

Vectors Are Inexpensive

Scalar

- N ops per cycle
(N^2) circuitry
- HP PA-8000
 - 4-way issue
 - reorder buffer:
850K transistors
 - incl. 6,720 5-bit register
number comparators

Vector

- N ops per cycle
($N + N^2$) circuitry
- T0 vector micro*
 - 24 ops per cycle
 - 730K transistors total
 - only 23 5-bit register
number comparators
 - No floating point

*See <http://www.icsi.berkeley.edu/real/spert/t0-intro.html>

What about I/O?

- Current system architectures have limitations
- I/O bus performance lags other components
- Parallel I/O bus performance scaled by increasing clock speed and/or bus width
 - Eg. 32-bit PCI: ~50 pins; 64-bit PCI: ~90 pins
 - Greater number of pins greater packaging costs
- Are there alternatives to parallel I/O buses for IRAM?

Serial I/O and IRAM

- Communication advances: fast (Gbps) serial I/O lines [YankHorowitz96], [DallyPoulton96]
 - Serial lines require 1-2 pins per unidirectional link
 - Access to standardized I/O devices
 - » Fiber Channel-Arbitrated Loop (FC-AL) disks
 - » Gbps Ethernet networks
- Serial I/O lines a natural match for IRAM
- Benefits
 - Serial lines provide high I/O bandwidth for I/O-intensive applications
 - I/O bandwidth incrementally scalable by adding more lines
 - » Number of pins required still lower than parallel bus
- How to overcome limited memory capacity of single IRAM?
 - SmartSIMM: collection of IRAMs (and optionally external DRAMs)
 - Can leverage high-bandwidth I/O to compensate for limited memory

ISIMM/IDISK Example: Sort

- Berkeley NOW cluster has world record sort: 8.6GB disk-to-disk using 95 processors in 1 minute
- Balanced system ratios for processor:memory:I/O
 - Processor: N MIPS
 - Large memory: N Mbit/s disk I/O & $2N$ Mb/s Network
 - Small memory: $2N$ Mbit/s disk I/O & $2N$ Mb/s Network
- Serial I/O at 2-4 GHz today (v. 0.1 GHz bus)
- IRAM: $2-4$ GIPS + 2 2-4Gb/s I/O + 2 2-4Gb/s Net
- ISIMM: 16 IRAMs+net switch+ FC-AL links (+disks)
- 1 IRAM sorts 9 GB, Smart SIMM sorts 100 GB

How to get Low Power, High Clock rate IRAM?

- Digital Strong ARM 110 (1996): 2.1M Xtors
 - 160 MHz @ 1.5 v = 184 “MIPS” < 0.5 W
 - 215 MHz @ 2.0 v = 245 “MIPS” < 1.0 W
- Start with Alpha 21064 @ 3.5v, 26 W
 - Vdd reduction 5.3X 4.9 W
 - Reduce functions 3.0X 1.6 W
 - Scale process 2.0X 0.8 W
 - Clock load 1.3X 0.6 W
 - Clock rate 1.2X 0.5 W
- 12/97: 233 MHz, 268 MIPS, 0.36W typ., \$49

DRAM v. Desktop Microprocessors

Standards	pinout, package, refresh rate, capacity, ...	binary compatibility, IEEE 754, I/O bus
Sources	Multiple	Single
Figures of Merit	1) capacity, 1a) \$/bit 2) BW, 3) latency	1) SPEC speed 2) cost
Improve Rate/year	1) 60%, 1a) 25%, 2) 20%, 3) 7%	1) 60%, 2) little change

Testing in DRAM

- Importance of testing over time
 - Testing time affects time to qualification of new DRAM, time to First Customer Ship
 - Goal is to get 10% of market by being one of the first companies to FCS with good yield
 - Testing 10% to 15% of cost of early DRAM
- Built In Self Test of memory:
 - BIST v. External tester?
 - Vector Processor 10X v. Scalar Processor?
- System v. component may reduce testing cost

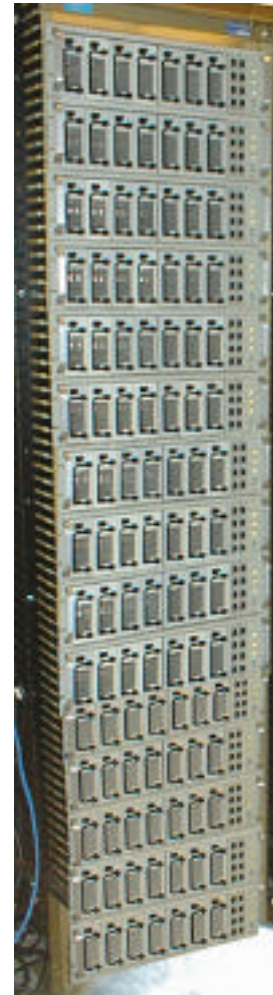
Words to Remember

“...a strategic inflection point is a time in the life of a business when its fundamentals are about to change. ... Let's not mince words: A strategic inflection point can be deadly when unattended to. Companies that begin a decline as a result of its changes rarely recover their previous greatness.”

– *Only the Paranoid Survive*, Andrew S. Grove, 1996

IDISK Cluster

- 8 disks / enclosure
- 15 enclosures /rack
= 120 disks/rack
- $1312 \text{ disks} / 120 = 11 \text{ racks}$
- $1312 / 8 = 164 \text{ 1.5 Gbit links}$
- $164 / 16 = 12 \text{ 32x32 switch}$
- $12 \text{ racks} / 4 = 3 \text{ UPS}$
- Floor space: wider
 $12 / 8 \times 200 = 300 \text{ sq. ft.}$



- HW,
assembly
cost:
\$1.5 M
- Quality,
Repair
good
- System
Admin.
better?