# Big Data ...
## and the Next Wave of *InfraStress*

## John  R. Mashey
## Chief Scientist, SGI

**Technology Waves:**
**NOT technology for technology's sake**   OK!
**IT'S WHAT YOU DO WITH IT**
**But if you don't understand the trends**
**IT'S WHAT IT WILL DO TO YOU**

Uh–oh!

# Big Data
## And The Next Wave of InfraStress

sgi

1. *Big data:* storage growing bigger faster
   DRAM: 1.6X/year (4X/3 years) continues
   Disk density:
   1.3X/year CAGR: historical trendline
   1.6X/year since ~1990
   2.0X/year leap ~1998/1999

2. *Net* continues raising user expectations
   More data (image, graphics, models)
   (Some) more difficult data (audio, video)
   Pressure on net, especially last mile

=> Explosion of WIDELY–accessible data
Create, understand, store, move ... or else ...
Drown in Wave of Infrastructure Stress

# InfraStress
# = <u>Infra</u>structure <u>Stress</u>

**sgi**

**in·fra·stress**. *n.*
1. Bad effects of faster change in computer
 *subsystems & usage*:
 CPUs, memory, disks, demand ...
 than in underlying *infrastructure:*
 bandwidths, addressability & naming,
 scalability of interconnect,
 operating systems, file systems, backup ...
Symptoms: bottlenecks, odd limits, workarounds,
 instability, unpredictability, nonlinear surprise,
 over−frequent releases, multiple versions,
 hardware obsolete before depreciated

2. In organizations that grow quickly, stress on
 management and support infrastructure.

# Environment:
# 4*X Data Problems



#1 Have data, cannot find & understand it     insight <– data

#2 Cannot create data from outside     creativity –> data

#3 Cannot have/process data, system limits     (data)

    Server always needs (30%?) headroom     power

#4 Have the data, but in wrong place/form     data <–> data

    Internal interconnect; network; firewalls     unleash

#X Rapid change, surprise amplify all 4 DATA problems

    Data distribution more troublesome than CPU distribution

# http://www.botham.co.uk

**sgi**

**Hidden flag**

## ELIZABETH BOTHAM & SONS

*Over 130 Years of Baking Excellence*

**Finest Quality North Yorkshire**

**Biscuits, Cakes and Plum Bread**

**World–wide delivery service**

*Est. 1865*

### ABOUT THE COMPANY

Since 1865, Elizabeth Botham & Sons has been a family run craft bakery in the ancient port of Whitby on the North Yorkshire coast. Following our original recipes, the finest ingredients are skilfully combined to produce Biscuits, Plum Bread and Cakes of the highest possible standard.

### CAKES

Luxury iced greetings cakes of the finest quality to be found. These rich fruit cakes are hand–crafted from premium ingredients, generously topped with our own handmade almond paste, iced and individually inscribed with either our standard message "Merry Christmas", "Happy Birthday", "Happy Mother's Day", "Happy Anniversary", or any five words of your choice (at a small surcharge), gift boxed with a personalised message of your choice and delivered to the door.

**Family bakery in Yorkshire + Website => suddenly begin selling outside UK.**

**Predict this?**

**No ... just predict change & surprise. But, some technology predictions easier...**

# 1. CPUs
## CMOS Microprocessors

sgi

Infra–
Stress

16–bit
micros
OK

**Change
minis–>
micros,
16 –> 32**

32–bit
micros
OK

**Change
32 –>
64/32**

64–bit
micros
OK

100%

% 32–bit
systems
shipped
(vs 16–bit)

% 64–bit
systems
shipped
(vs 32–bit)

1980    1983    1986    1989    1992    1995    1998    2001    2004    2007

**64**

1st 64–bit micro (MIPS R4000)

# 2. Big Memory & Micros

**sgi**

Infra–
Stress

16–bit
micros
OK

Change
minis–>
micros,
16 –> 32

32–bit micros
OK

Change
32 –>
64/32

64–bit micros
OK

**PCs:
640K
painful limit
1MB hack...**

**large servers:
<4GB limit painful**

**large servers
> 4GB useful**

1980    1983    1986    1989    1992    1995    1998    2001    2004    2007

# 3. Big Net

**Infra–Stress**

**Everybody knows this one!**

**Note: does not mean effects stop, just that most organizations will have Web–ized operations by 2002.**

**Networks Organizations Procedures**

**BIG NET: The Net, WWW**

1980  1983  1986  1989  1992  1995  1998  2001  2004  2007

# 4. Bigger (Disk) Data

**sgi**

Infra–
Stress

1.3X   1.6X   2X

**Disk file systems
Backups
I/O systems**

**Many must
rewrite
critical
software**

BIGGER DATA
3.5" disk density

1980   1983   1986   1989   1992   1995   1998   2001   2004   2007

http://www.quantum.com/src/history, http://www.disktrend.com
http://www.ibm.com/storage/microdrive: 340MB Microdrive, 1999. 1.7"x1.4"x.19"

# 5. HUGE Data (Maybe)
## Storage Hierarchy

**sgi**

1) Tapes, near–line storage

2) Laser–enhanced magnetics
 for removables, maybe fixed disks
10X: TeraStor
NFR: "Near–Field Recording"
   5.25", removable, 2400 RPM, 18ms
   2Q99: 10GB, 6 MB/sec, <$800
   4Q99: 20GB, 11 MB/sec, <$1200
   ??   : 40GB, 2–sided
3–5X: Quinta (Seagate), demo 11/98
OAW: Optically assisted Winchester

**Like bigger, but worse**

| 1980 | 1983 | 1986 | 1989 | 1992 | 1995 | 1998 | 2001 | 2004 | 2007 |
|------|------|------|------|------|------|------|------|------|------|

~1999: Laser=enhanced magnetic disks (removable)
http://www.quinta.com, http://www.terastor.com

# InfraStress Addup

**sgi**

Infra–
Stress

5. HUGE DATA:
   Storage hierarchy

4. BIGGER DATA:
   1.3X –> 1.6X –> 2X

3. BIG NET:
   The Net, WWW

2. BIG MEMORY:
   DRAM vs 32–bit

1. CPUS: Microprocessors
   32 –> 64

1980   1983   1986   1989   1992   1995   1998   2001   2004   2007

# Technology Change Rates
## Example: Large Server*

sgi

| | Years | Large Server | # Revisions in 6 years | |
|---|---|---|---|---|
| **H/W chassis** | 4..6 | | 0 | |
| **Interconnects** | | | | |
| I/O bus (PCI...) | 4–6+ | | 0–(1) | |
| CPU==mem | 3–5 | | 0–(1) | |
| Backplane | 3–5 | | 0 | |
| Network | varies | | 1–2 | |
| **Subsystems** | | | | |
| **CPU MHz** | .75–1.5 | | **4–8** | |
| **4X DRAM** | 3 | | **2–(3)** | |
| **Disks** | 1 | | **6** | |
| **Graphics** | 1.5–2.5 | | | |
| **Software** | | | | |
| File system | 8–10 | | 0–1 | |
| OS release | 1–2 | | **2–6** | |
| App release | 1–2 | | **2–6** | |
| **Data** | forever | | | |
| Media | not long | | | |

0     3     6
**Years**

*Desktops &
other access devices
cycle faster, maybe

# Technology Trends

**sgi**

**Capacities – Great News**

**Latencies – Not–so–great News**

**Bandwidths – InfraStress**

**Interactions – Surprises**

**Tradeoffs – keep changing**

# 1"x 3.5" Disk Capacity

sgi



**Capacity**

| | | |
|---|---|---|
| 1.3X | 1.6X | 2X |

>4X / 3 years

"Fear is **not** an option ..."

90 GB
80 GB
70 GB
60 GB
50 GB
40 GB
30 GB
20 GB
10 GB
0 GB

Traditional disk density growth

These are 1" (LP) drives only.

1.6" (HH) drives have higher capacity, (36–50GB available 1Q99).

72

1.6X

36

16.8*

18

9

4.5

.5   1

1.3X

1980  1983  1986  1989  1992  1995  1998  2001  2004  2007

**"Disks are binary devices ... new and full"**

*IBM Desktap 16GP, Giant Magnetoresistive heads (GMR), 4Q97.

# Log–scale charts ahead

sgi

**Linear scale**

**Logarithmic scale
Huge differences do
not *look* so big at top**

==>

**Parallel = same ratio
Inflection points clear**

# DRAM Capacity: 1.6X CAGR
## 4X / 3 years

sgi

**Capacity**

- 1 TB
- 100 GB
- 10 GB
- **4GB**
- 1 GB
- 100 MB
- 10 MB
- 1 MB
- 100 KB
- 10 KB
- 1 KB

**Supers**

**Big T3E ~220GB**
**Multi–rack Origin2000 128GB**
Origin2000 (1 Rack) 32GB
Power Challenge 16GB

Challenge 2GB

Power Series 256MB

**Total DRAM:**
**actually sold,**
**1–rack system**

MIPS M/500
32MB

"4Gb"??

"1Gb"

"256Mb"

"64Mb"

"16Mb"

**1 DRAM:**
**Bytes/chip**

1Q92: 1st 64–bit micro
4Q94: technical use

1980   1983   1986   1989   1992   1995   1998   2001   2004   2007

**64**   **T64**

**See: John R. Mashey, "64–bit Computing", BYTE, September 1991, 135–141.**

# Disk Capacity:
## 1.3X –> 1.6X –> 2X

sgi

Capacity

1 TB

100 GB

10 GB

1 GB

100 MB

10 MB

1 MB

100 KB

10 KB

1 KB

**1 Disk ~= 300–500 DRAMs**

**1"X 3.5" Disk Bytes/disk**

144?

72

36

18

9

4.5

1

.5

**Historical trend 1.3X**

**DRAM Bytes/chip**

4Gb??

1Gb

256Mb

64Mb

16Mb

1980   1983   1986   1989   1992   1995   1998   2001   2004   2007

See: John R. Mashey, Darryl Ramm, "Databases on RISC: still The Future",
UNIX Review, September 1996, 47–54.

# 3.5" Disk Review

Height (1" or 1.6") X (4" X 5.75")
Capacity (1MB = 1,000,000 B)
Seek Times (msecs)
 Track–to–track (Read/Write)
 Average (Read/Write)
  Typical < Average (OS & controllers)
 Maximum (Read/Write)
Rotational latency (msecs)
 Average Latency = .5 * rev = 30000/RPM
Bandwidths (MB/sec)
 Internal Formatted Transfer
 ZBR range
 External Rate (Bus)
Density (Gbit/sq inch)

Controller

# 3.5" Disk Review

sgi

- – **Capacity/drive ~ # platters (varies)**
- – **Capacity/platter ~ areal density**
- – **Bandwidth ~ RPM * Linear density**
- – **Seek time ... improves slowly**
- – **Combine several drives onto one: take care, may lose seeks/second**
- – **IOPS vs MB/s applications**

**System (OS)**
**I/O Bus (~PCI)**
**Peripheral Connect (~SCSI)**
**Embedded Disk Controller**
**Disk  Seek**
**        Rotate**
**        Read**

**Time –>**

# Common Disk Types

1. By capacity
   A. Large (1.6" x 3.5", HH) ~8–10 platters
   B. Medium (1" X 3.5", LP), ~4–5 platters
   C. "Depopulated", 1 platter
   D. Smaller platters ...
   E. "Microdrive", 1 small platter

2. By target
   – High–performance (B: high RPM)
   – High–capacity (A)
   – By IOPs (multiples of C & D)
   – By cost [ATA, IDE versions of A, B, C]
   – By physical size (mobile, consumer)Bad

Huge disks => long backup times
   Good for archive–like applications

# Storage Densities

**sgi**

**10,000,000 Billion Atoms/in2**

Density/in²

- 10,000 Tb
- 1,000 Tb
- 100 Tb
- 10 Tb
- 1 Tb
- 100 Gb
- 10 Gb
- 1 Gb
- 100 Mb

"IBM and other vendors, universities, and the government are working on a holographic storage system they say will achieve 100Gb per square inch and data transfer rates of 30Mb per second by November 1998. Future targets are 100Gb per square inch and 100Mb per second data rates by January 1999, and 100Gb per square inch and 1Gb per second transfer by April 1999.

OptiTek, in Mountain View, Calif., is developing holography products, promising 5.25 disk capacities of 100GB with cartridges backward–compatible to current automated libraries. The company will release evaluation models in the second half of 1999, and plans to release "write–once" products for use in archiving applications by early 2000."

InfoWorld Electric, "When Data Explodes", http://www.idg.net

**1 TB/in3 Tape density**

**300 Gb/in2 Atomic Force microscope(?)**

**45 Gb/in2 AF demo**

**40–70 Gb/in2**

**Near–field recording**

GMR: 2.4–2.6 (1997)
10 (2001), 40 (2004)

**2.0–2.8 Gb/in2**

**super– paramagnetic limit**

**1.0–1.5 Gb/in2**

**.660–.981 Gb/in2**

**.129 Gb/in2: Tape: DDS–3**

**1980  1983  1986  1989  1992  1995  1998  2001  2004  2007**

**See: Merrit E. Jones, The MITRE Corp, "The Limits That Await Us", THIC Meeting April 23, 1997, Falls Church, Va.**
**See http://www.terastor.com on Near–field recording.**

# Disk Issues
## Workloads Converge

"IOPS" – Transaction / seeks/second
Classic OLTP, small blocks

"MB/s" – Bandwidth (& backup!)
Classic technical, larger blocks

Some commercial now more like technical

**Classic Technical**

**Classic Commercial**

Gflops   **Big Data**   tpms          other

Silicon Graphics

# Disk Issues – Implications

1. Huge capacity leap breaks old filesystems
   Hard limits (2GB, 8GB, etc) OR
    Algorithmic performance, scaling issues

2. More memory, more bandwidth, everywhere
   Small disk blocks even less efficient
   => 64–bit addressing more useful
   => Big pages, map more pages, MMUs
   => More memory => more bandwidth
   => More interconnect bandwidth

3. BACKUP ...
   Must run many tapes, full–speed, parallel
   Sometimes use HSM, RAID, mirror
   New cartridge disks may be useful

# Disk Rotational Latencies
## High−performance − 1/2 Rotation

Faster rotation ~ 2−3 years
Average Latency = .5 * (60/RPM)

Clock

1 GHz — 1 ns

100 MHz — 10 ns

10 MHz — 100 ns

1 MHz — 1 mics

100 Khz — 10 mics

10 KHz — 100 mics

1 KHz — 1 msec

100 Hz — 10 msec

10 Hz

**Platters shrink**

1.5

2.0

3.0

4.17

5.55

8.3 msec

20000

15000

3600    5400    7200   10000              RPM

1980    1983    1986    1989    1992    1995    1998    2001    2004    2007

**Money can buy bandwidth, but latency is forever.**

# Disk Average Seek
## High–performance disks

sgi

**Faster rotation ~ 2–3 years**
**Average Latency = .5 * (60/RPM)**

**1/2 Rotation faster than average seek ...**
**But of course, short seeks are faster**

**Short random blocks dominated by seek**
**Large blocks dominated by transfer time**

Clock

1 GHz — 1 ns

100 MHz — 10 ns

10 MHz — 100 ns

1 MHz — 1 mics

100 Khz — 10 mics

10 KHz — 100 mics    16 msec 15 14  12  9              8    6  5    4          Avg Seek

1 KHz — 1 msec                                                    3.0   2.0    1.5

100 Hz — 10 msec         8.3 msec    5.55    4.17                              1/2 Rotation

                                                                              Avg Seek

10 Hz —

1980   1983   1986   1989   1992   1995   1998   2001   2004   2007

# Disk Total Latencies
## 1/2 Rotation + Average Seek

sgi

Faster rotation ~ 2–3 years
Average Latency = .5 * (60/RPM)

1/2 Rotation faster than average seek ...
    But of course, short seeks are faster

Short random blocks dominated by seeks
Large blocks dominated by transfer time

Clock

| | |
|---|---|
| 1 GHz | 1 ns |
| 100 MHz | 10 ns |
| 10 MHz | 100 ns |
| 1 MHz | 1 mics |
| 100 Khz | 10 mics |
| 10 KHz | 100 mics |
| 1 KHz | 1 msec |
| 100 Hz | 10 msec |
| 10 Hz | |

16 msec 15 14  12  9          8     6  5     4      Avg Seek

                                          1.5
                              3.0   2.0
                        4.17
                5.55                               1/2 Rotation

8.3 msec                                           Avg Seek

                                                   Latency
24 msec 23 20 18 15   13  12    11 9  7    5.5     1.1X CAGR

1980   1983   1986   1989   1992   1995   1998   2001   2004   2007

4/25/98    page 2 6

# CPU Latency, Performance

sgi

Effective instruction latency =
DRAM ... CPU cycle/peak issue

Clock

10 GHz — .1 ns

1 GHz — 1 ns

100 MHz — 10 ns

10 MHz — 100 ns

1 MHz — 1 mics

100 Khz — 10 mics

10 KHz — 100 mics

1 KHz — 1 msec

100 Hz — 10 msec

10 Hz —

1 ns

4 ns

10ns

40ns

125ns

120ns

100ns

80ns

60ns

40ns

Upper edge = raw DRAM access time
Lower edge = lean memory system,
including overhead, for acual load
2000: 40ns nominal  –> 150ns+

1980   1983   1986   1989   1992   1995   1998   2001   2004   2007

**CPU perform
1.4X–1.6X**

**CPU cycle
1.4X CAGR**

**Raw DRAM
1.1X CAGR**

**CPU:DRAM:
  40X (cycle)
  100X (real)
  400X  (instrs)
Soon:
1000X (instrs)**

# Latency & Performance

**sgi**

**Clock**

Effective instruction latency =
DRAM ... CPU cycle/peak issue

| | |
|---|---|
| 10 GHz | .1 ns |
| 1 GHz | 1 ns |
| 100 MHz | 10 ns |
| 10 MHz | 100 ns |
| 1 MHz | 1 mics |
| 100 Khz | 10 mics |
| 10 KHz | 100 mics |
| 1 KHz | 1 msec |
| 100 Hz | 10 msec |
| 10 Hz | |

1 ns

4 ns

10ns

40ns

125ns

40ns

100ns    80ns

60ns

120ns

Lower edge = memory system

CPU:Disk:1986
200K instrs

24 msec  23  20  18 15    13  12    11  9   7    5.5

**1980   1983   1986   1989   1992   1995   1998   2001   2004   2007**

CPU perform
1.4X–1.6X

CPU cycle
1.4X CAGR

Raw DRAM
1.1X CAGR

CPU:DRAM
1000X (insts)

CPU:Disk
>5M instrs now
>30M soon

Disk latency
1.1X CAGR

Humans
1X/ ...

# Latencies – Implications

**sgi**

**1. CPU <–> DRAM <–> disk
   Latency ratios already bad, getting worse.**
   "Money can buy bandwidth, but latency is forever."

**==> More latency tolerance in CPUs
==> Trade (bandwidth, memory, CPU,
     PROGRAMMING) for latency
==> Already worth 1M instruction
     to avoid a disk I/O**

**2. RDBMS huge buffer areas for indices,
   small tables, to avoid latency**

**3. Networks: be alert for latency issues**

# Input/Output: A Sad History

sgi

"I/O certainly has been lagging in the last decade."
     Seymour Cray
     Public Lecture (1976)

"Also, I/O needs a lot of work."
     David Kuck
     Keynote Address, 15th Annual Symposium
     on Computer Architecture (1988)

"Input/output has been the orphan of
  computer architecture ... I/O's revenge is at hand"
     David A. Patterson, John. L. Hennessy
     Computer Architecture: A Quantitative Approach,
     2nd Ed (1996), Morgan Kaufmann.

# I/O Single-Channel Bandwidth



**I/O Busses falling behind 4X/3 growth, need faster I/O**

**4X/3**

XIO (4Q96) [1.2 GB/s (2X .64)]

GigaRing ●

Indigo2, Indy GIO64 [.2]

Indigo, GIO32 [.1]

PCI64-66 [.4]

PCI64 [.2]

PCI32 [.1]

Sun SBUS64 [.1]

EISA (.033 p)

ISA (.007 p)

Y-axis (GB/sec): 1000, 100, 10, 1, 0.1 / 100 MBs, 0.01 / 10 MBs, 0.001 / 1MBs

X-axis: 1980, 1983, 1986, 1989, 1992, 1995, 1998, 2001, 2004, 2007

# Bus–Based SMP
## Bandwidth Wall

sgi

**SMP Busses falling behind 4X/3 growth, need change**

**4X/3**

**Data gap, big, growing**

**Laws of physics ... are laws ...**

Sun UE X000
2Q96 (2.5)

DEC 8400
2Q95 (1.6)

SGI Challenge
1Q93 (1.22)

Sun SC2000
2Q93, (.5)

Intel SHV
2Q96, (.534p)

SGI Power
Series 4Q88 (.064)

Sequent Highly
Scalable Bus
1994, (.107, [.240 p])

Sequent Bus
4Q87 (.053)

−2.5 GB/s
2X / 3 growth,
slowing

SMP Bus,
Memory, Total I/O

**Y-axis:** 1000 GB/sec, 100, 10, 1, 0.1 / 100 MBs, 0.01 / 10 MBs, 0.001 / 1MBs

**X-axis:** 1980, 1983, 1986, 1989, 1992, 1995, 1998, 2001, 2004, 2007

# Bandwidths (ccNUMA, XBAR)

sgi



**Why ccNUMA?**
**A: Central XBAR $$.**

**4X/3**

**128p**

128p Origin, Onyx2: up to
**80GB/s I/O**
40 GB/s memory,
**20 GB/s Bisection**

**SMP Bus Bandwidth**

**1 XIO, 1.28 GB/s**

**1p**

**Origin200 PCI64,**
**.2 GB/s**

**I/O Bus Bandwidth**

**Start small**
**Buy incrementally**
**Scale big**

Y-axis: 1000 GB/sec, 100, 10, 1, 0.1 / 100 MBs, 0.01 / 10 MBs, 0.001 / 1MBs

X-axis: 1980, 1983, 1986, 1989, 1992, 1995, 1998, 2001, 2004, 2007

# LAN, Interconnect Bandwidths

**sgi**

**Networks improving faster than SMP Bus & I/O Busses**

**4X/3**

**Networks must improve to stay ahead of disks**

Origin ccNUMA I/O

High–end SMP bus bandwidth

Gigabyte System Network (GSN)

HIPPI 800

ATM OC12

Ethernet 1000BT

**1000BT coming faster**

ATM OC3

Ethernet 100BT

Ethernet 10BT

| Y-axis | |
|---|---|
| 1000 GB/sec | |
| 100 | |
| 10 | |
| 1 | |
| 0.1 / 100 MBs | |
| 0.01 / 10 MBs | |
| 0.001 / 1MBs | |

X-axis: 1980 1983 1986 1989 1992 1995 1998 2001 2004 2007

# Beyond the LAN
## (Different Scale!)

sgi



Gigabyte System Network (GSN)

HIPPI 800

Ethernet 1000BT

ATM OC12

DS–4, 274 Mbs Mbs

ATM OC3

Ethernet 100BT

T3, 43.2 Mbs, 5.4 MBs

Ethernet 10BT

*DSL, 2 Mbs – 7 Mbs

3Mbs Cable Modem (375 KBs)

T1, 1.544 Mbs

ISDN (128Kb, 16 KBs)

56Kbs Modem (7 KBs)

28.8Kbs Modem (3.6 KBs)

All these are theoretical peaks, reality = less

Y-axis:
1 / 1 GBs
0.1
0.01 / 10 MBs
0.001 / 1MBs
0.0001 / 100 KBs
0.00001 / 10 KBs
0.000001 / 1 KBs

X-axis: 1980  1983  1986  1989  1992  1995  1998  2001  2004  2007

# Disk Bandwidths (Highest)

**sgi**

1"X 3.5" Disk
Bytes/disk

Striped Bandwidth/
4 disks
3 disks
2 disks
Bandwidth/1 disk

4
3
2
1

2001: Guess 40 MB/s

1999 – 18GB,
10000 RPM, 28 MB/s

1998 – 9GB, 7200 RPM, 13 MB/s
10000 RPM , 15 MB/s

Y-axis: 1000 GB/sec, 100, 10, 1, 0.1 / 100 MBs, 0.01 / 10 MBs, 0.001 / 1MBs

X-axis: 1980  1983  1986  1989  1992  1995  1998  2001  2004  2007

# Fast Disk Bandwidth
## vs Peripheral Connections

**Disk bandwidth growth overpowers peripheral connection growth!**

| # 10MB/s Disks | FW SCSI 20 MB/s | F20W 40 MB/s | FC100 100 MB/s |
|---|---|---|---|
| 1 | 10 | 10 | 10 |
| 2 | 18* | 20 | 20 |
| 3 | * | 30 | 30 |
| 4 | * | 32* | 40 |
| ... | ... | ... | ... |
| 10 | * | * | 95* |

**\* Already saturated on bandwidth tasks, like backup or striped–disk I/O.**

**Peripheral Connections MB/s**

**x = 4 disks exhaust bus in bandwidth apps**

**4 disks**
**3 disks**
**2 disks**
**Bandwidth/1 disk**

200 FC200
160 SCSI
100 FC100
80 SCSI LV
40 SCSI F20W
20 FW SCSI
10 F SCSI

40 MB/s
28 MB/s
10 MB/s

Y-axis: 1000 GB/sec, 100, 10, 1, 0.1 / 100 MBs, 0.01 / 10 MBs, 0.001 / 1MBs

X-axis: 1980 1983 1986 1989 1992 1995 1998 2001 2004 2007

# Fast Disk Bandwidth
## vs Networks & Peripheral Connections

sgi

| | | |
|---|---|---|
| 10BaseT | = .1 | 1997 fast disks (bottleneck |
| 100BaseT | = 1 | 1997 fast disk |
| 1000BaseT | = 2 | 2001 fast disks (2 X 40 MBs) |
| | = 1 | 2001 dual–head fast disk (80 MBs) |
| GSN | = many disks, still not enough for all! |

1000 GB/sec

100

10

**Theoretical ... reality much less**

1

GSN

0.1
100 MBs

Ethernet 1000BaseT

**Peripheral Connections MB/s**

100 FC100
80 SCSI LV
40 SCSI F20W
20 FW SCSI
10 F SCSI

40 MB/s

15 MB/s

0.01
10 MBs

Ethernet 100BaseT

X

10 MB/s

4
3
2
1

0.001
1MBs

Ethernet 10BaseT

**Bandwidth/1 disk**

1980  1983  1986  1989  1992  1995  1998  2001  2004  2007

# Bandwidths – Summary

sgi

**Networks and disks pressure on I/O Bus and SMP Bus**

**4X/3**

Disks
  InfraStress on networks

Disks + networks
  InfraStress on I/O bus

Disks + nets + memory
  InfraStress on SMP bus

Origin
ccNUMA
I/O

High–end
SMP bus
bandwidth

Network
bandwidth

Disk
bandwidth

1 I/O bus
bandwidth

4
3
2
1

| | 1000 GB/sec |
| 100 |
| 10 |
| 1 |
| 0.1 / 100 MBs |
| 0.01 / 10 MBs |
| 0.001 / 1MBs |

1980  1983  1986  1989  1992  1995  1998  2001  2004  2007

# Bandwidths – Implications

**sgi**

1. SMP Busses not growing with 4X/3
   Interconnect and memory bandwidth limits
   ==> Crossbars
      Centralized (mainframe)
      Distributed (ccNUMA)

2. Some I/O busses, peripheral connects,
   and especially networks under pressure
   to keep up with disk bandwidth

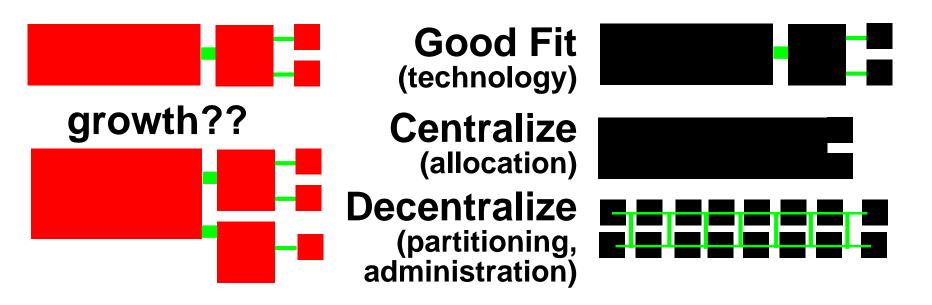3. Disks are faster than tapes ... backup?

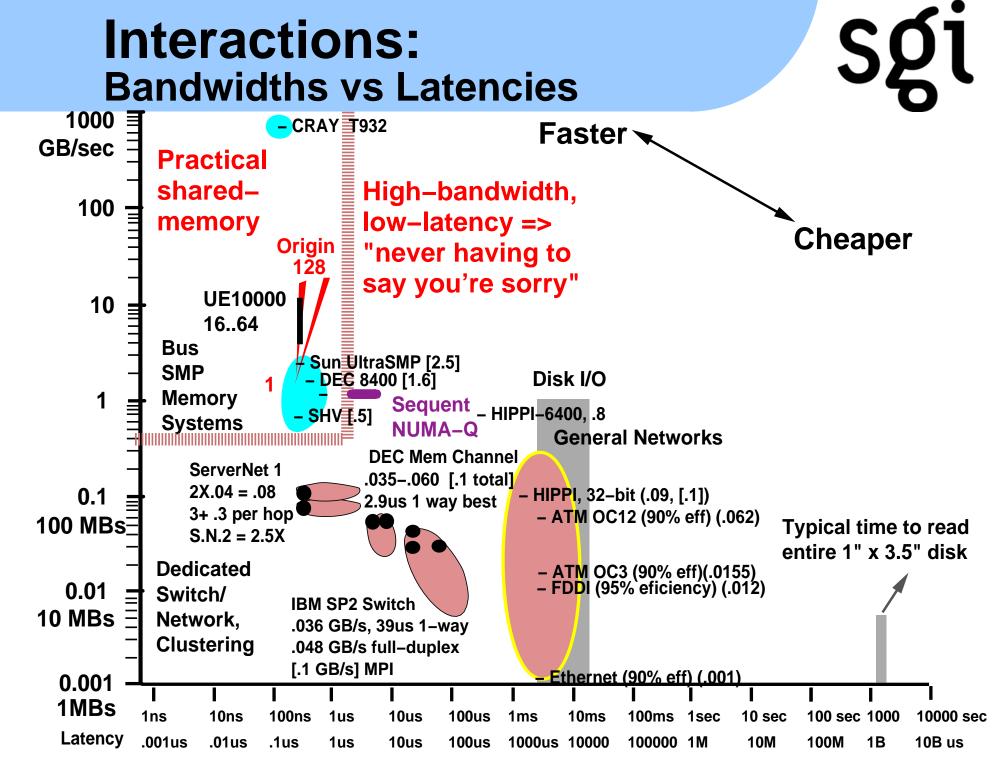4. SANs for  bandwidth and latency

# Interactions: Distributed Data

**Shape of solution driven by shape of hardware?**
**"Natural" distribution of work: cost–effective**
**"Unnatural" data distribution: very painful**
*High bandwidth, low latency, or else...*

**Better: make hardware match shape of problem**

**Problem Shape**                    **Solution Shape?**

growth??

**Good Fit**
(technology)

**Centralize**
(allocation)

**Decentralize**
(partitioning,
administration)

# Interactions:
## Bandwidths vs Latencies

sgi

**Practical shared–memory**

**High–bandwidth, low–latency => "never having to say you're sorry"**

Faster

Cheaper

1000 GB/sec
— CRAY T932

100

10

Origin 128

UE10000 16..64

**Bus SMP Memory Systems**

1

— Sun UltraSMP [2.5]
— DEC 8400 [1.6]
—
— SHV [.5]

**Sequent NUMA–Q**

— HIPPI–6400, .8

**Disk I/O**

**General Networks**

ServerNet 1
2X.04 = .08
3+ .3 per hop
S.N.2 = 2.5X

DEC Mem Channel
.035–.060 [.1 total]
2.9us 1 way best

0.1
100 MBs

— HIPPI, 32–bit (.09, [.1])
— ATM OC12 (90% eff) (.062)

**Typical time to read entire 1" x 3.5" disk**

**Dedicated Switch/ Network, Clustering**

0.01
10 MBs

IBM SP2 Switch
.036 GB/s, 39us 1–way
.048 GB/s full–duplex
[.1 GB/s] MPI

— ATM OC3 (90% eff)(.0155)
— FDDI (95% eficiency) (.012)

— Ethernet (90% eff) (.001)

0.001
1MBs

| Latency | 1ns | 10ns | 100ns | 1us | 10us | 100us | 1ms | 10ms | 100ms | 1sec | 10 sec | 100 sec | 1000 | 10000 sec |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | .001us | .01us | .1us | 1us | 10us | 100us | 1000us | 10000 | 100000 | 1M | 10M | 100M | 1B | 10B us |

# Interactions:
## Disk Technology Trends

**sgi**

**Capacities**
> Grow very fast

**Latencies**
> Barely improve for small blocks
> Improve moderately for large blocks

**Bandwidths**
> Improve, but not so fast as capacity
> Capacity/bandwidth ratios get worse
> Pressure –> more smaller disks

**Interactions**
> 100BaseT, PCI32, F+W SCSI overrun
> Backup rethinking
>> Desktop & 2 half–empty disks?
>> Backup servers?

# Technology Summary

| | Good | Bad | Ugly |
|---|---|---|---|
| CPU | Mhz | Parallelism | Latency |
| SRAM | On–chip | | Latency |
| RAM | Capacity | | Latency |
| Disk | Capacity | | Latency |
| Tape | Capacity | Bandwidth | Latency |
| Network | Bandwidth | | Latency |
| Software | | | Work! |
| Sysadmin Technology | | | Exciting |

# Conclusion: InfraStress
## Wishlist for Overcoming It

1. Find/understand: insight
   Tools: Navigate, organize, visualize
2. Input: creativity
   Tools: create content from ideas

3. Store and process the data: power
   Big addressing, modern file system
   Big I/O (number and individual speed)
   Big compute (HPC or commercial)
4. Move it: unleash
   Scalable interconnect
   High−performance networking

5. Change: survive!
   Incremental scalability, headroom
   Infrastructure already upgraded

# References

1. http://www.storage.ibm.com/hardsoft/diskdrdl/library/technolo.htm
   IBM storage web page