# Is there a pathway to a Green Grid  ??
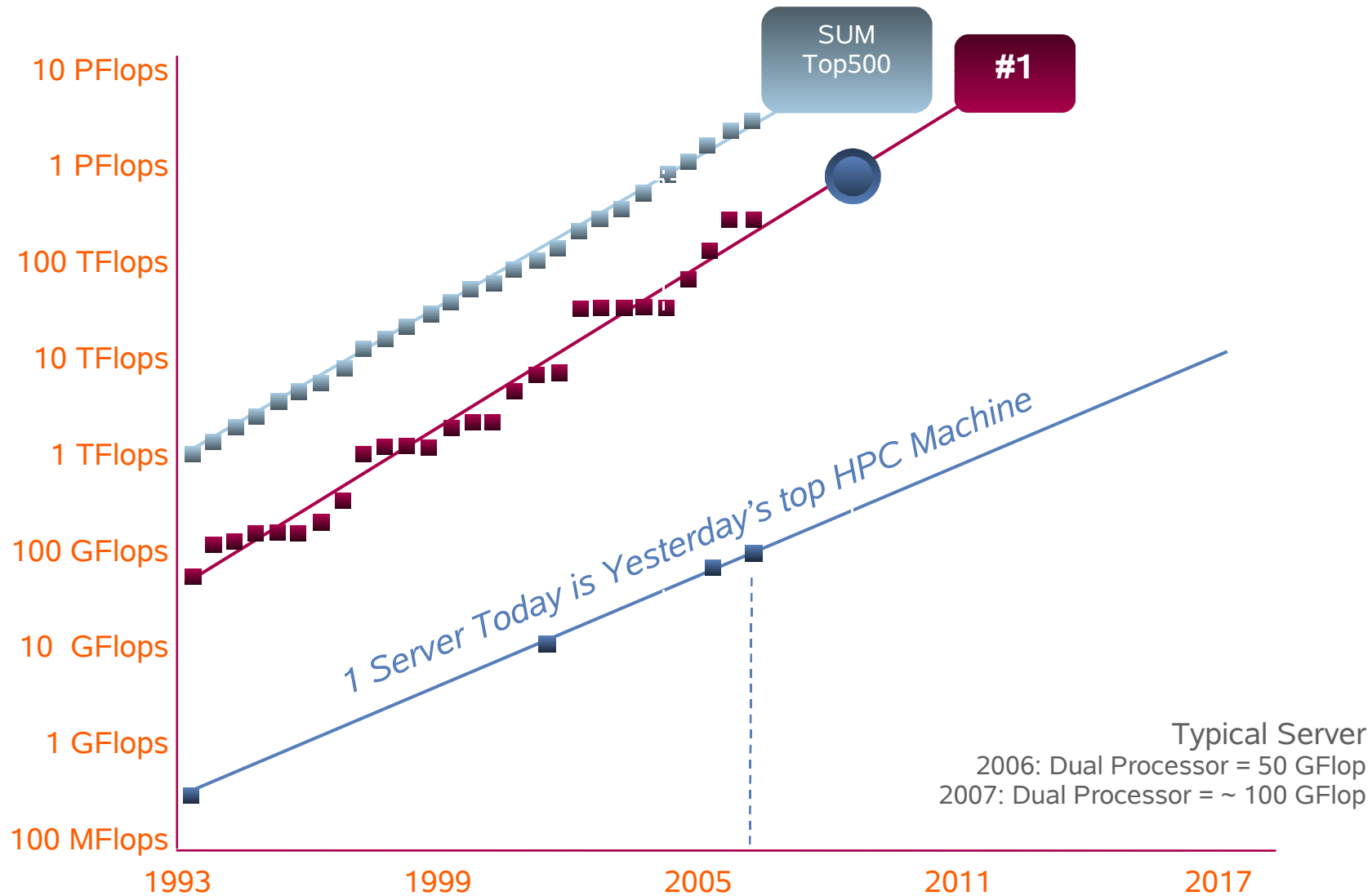
**Simon See, Ph.D**
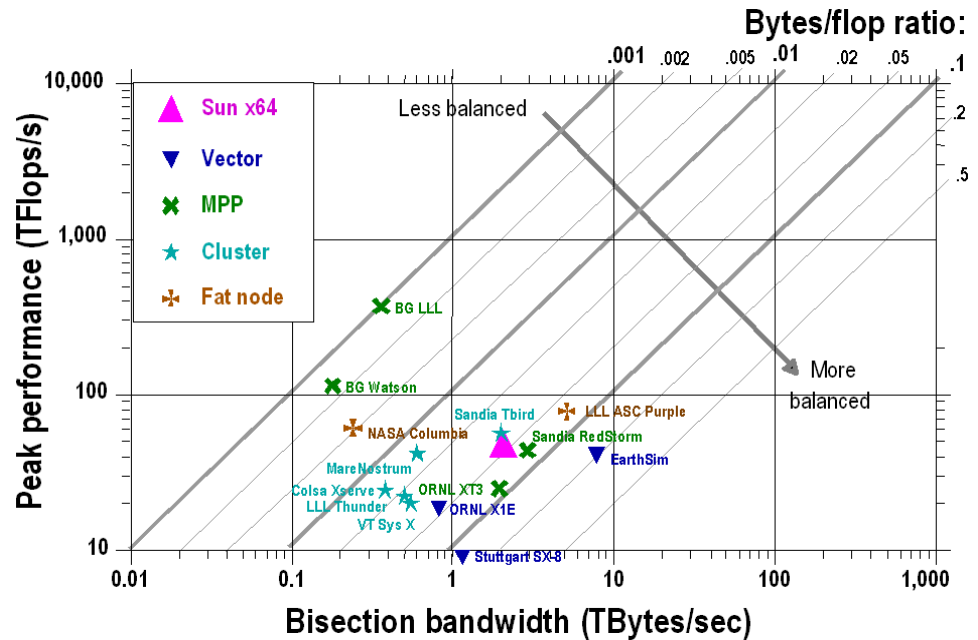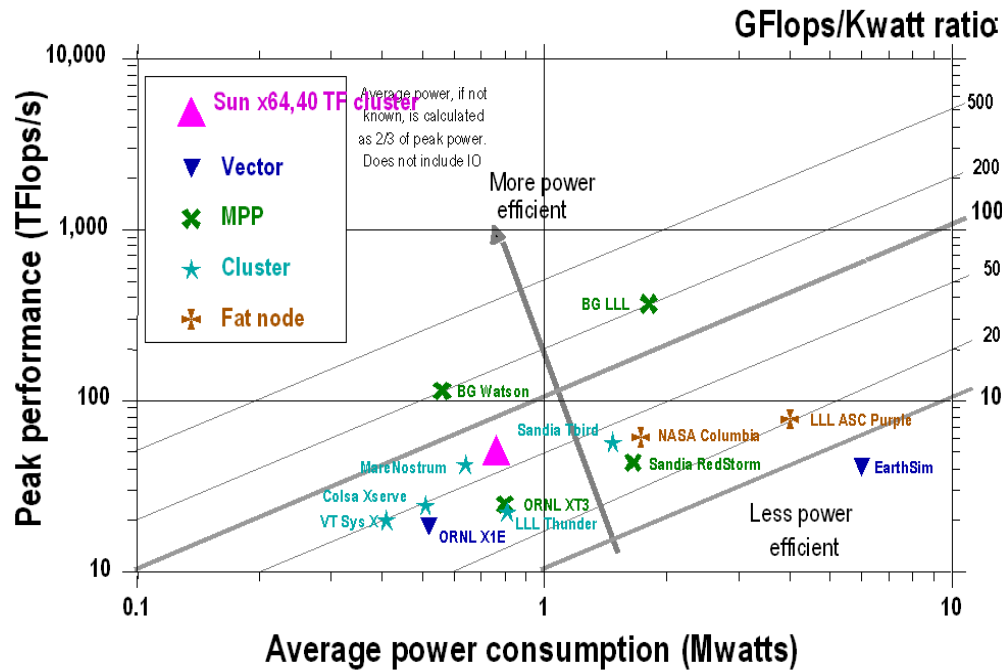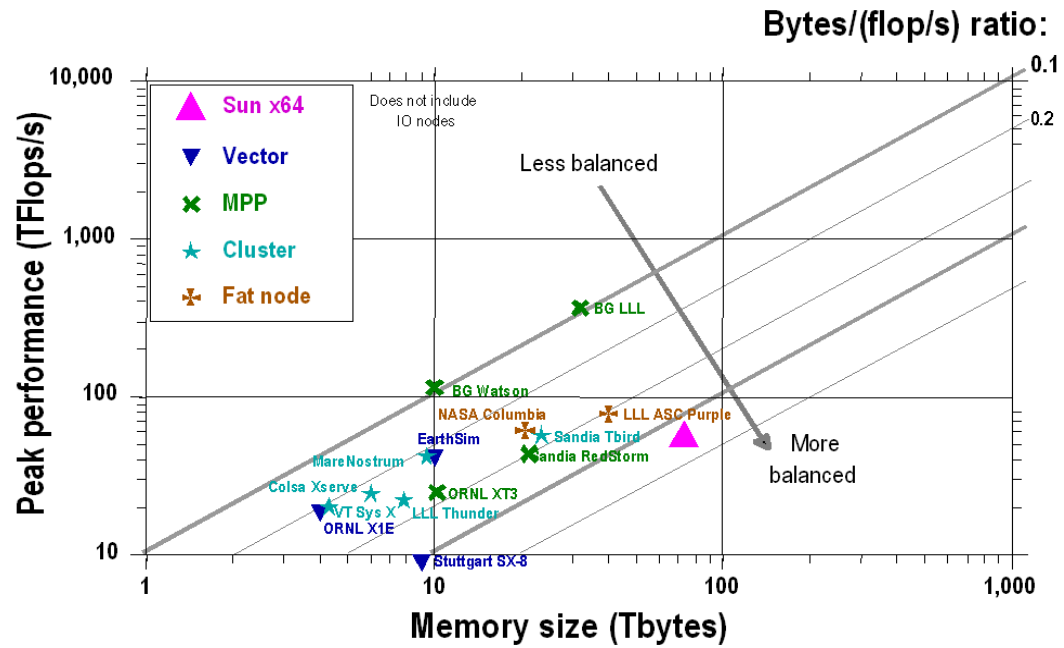
Director,  Global HPC Solution          Associate Prof.
Sun Microsystesms Inc.                  NTU and NUS

# HPC Top500: An Example of Moore's Law



SUM Top500

#1

10 PFlops

1 PFlops

100 TFlops

10 TFlops

1 TFlops

100 GFlops

10 GFlops

1 GFlops

100 MFlops

*1 Server Today is Yesterday's top HPC Machine*

Typical Server
2006: Dual Processor = 50 GFlop
2007: Dual Processor = ~ 100 GFlop

1993        1999        2005        2011        2017

# Number of Data Centers Growing

| | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 |
|---|---|---|---|---|---|---|---|---|---|---|
| Very Large | 100 | 106 | 120 | 160 | 175 | 185 | 210 | 230 | 250 | 270 |
| Large | 900 | 870 | 880 | 900 | 920 | 990 | 1,040 | 1,100 | 1,170 | 1,250 |
| Medium | 1,405 | 1,385 | 1,395 | 1,420 | 1,490 | 1,585 | 1,665 | 1,765 | 1,870 | 1,975 |
| Small | 2,180 | 2,100 | 2,110 | 2,190 | 2,230 | 2,290 | 2,360 | 2,430 | 2,500 | 2,570 |
| **Total** | **4,585** | **4,461** | **4,505** | **4,670** | **4,815** | **5,050** | **5,275** | **5,525** | **5,790** | **6,065** |
| Growth Rate | | -2.70% | 1.00% | 3.70% | 3.10% | 4.90% | 4.35% | 4.74% | 4.80% | 4.75% |

### Small Data Center
- Between 350 – 500 Servers Installed
- 15,000 Sq Feet Of Raised Floor
- Predominately Volume Server Architecture, with 1 -3 High End Server Systems

### Medium Data Center
- Between 1,500 – 1,700 Servers Installed
- 20,000 Sq Feet Of Raised Floor
- Four or Five High End Systems Form The Basis Of Enterprise Systems

### Large Data Center
- Between 2,000 – 2,500 Servers Installed
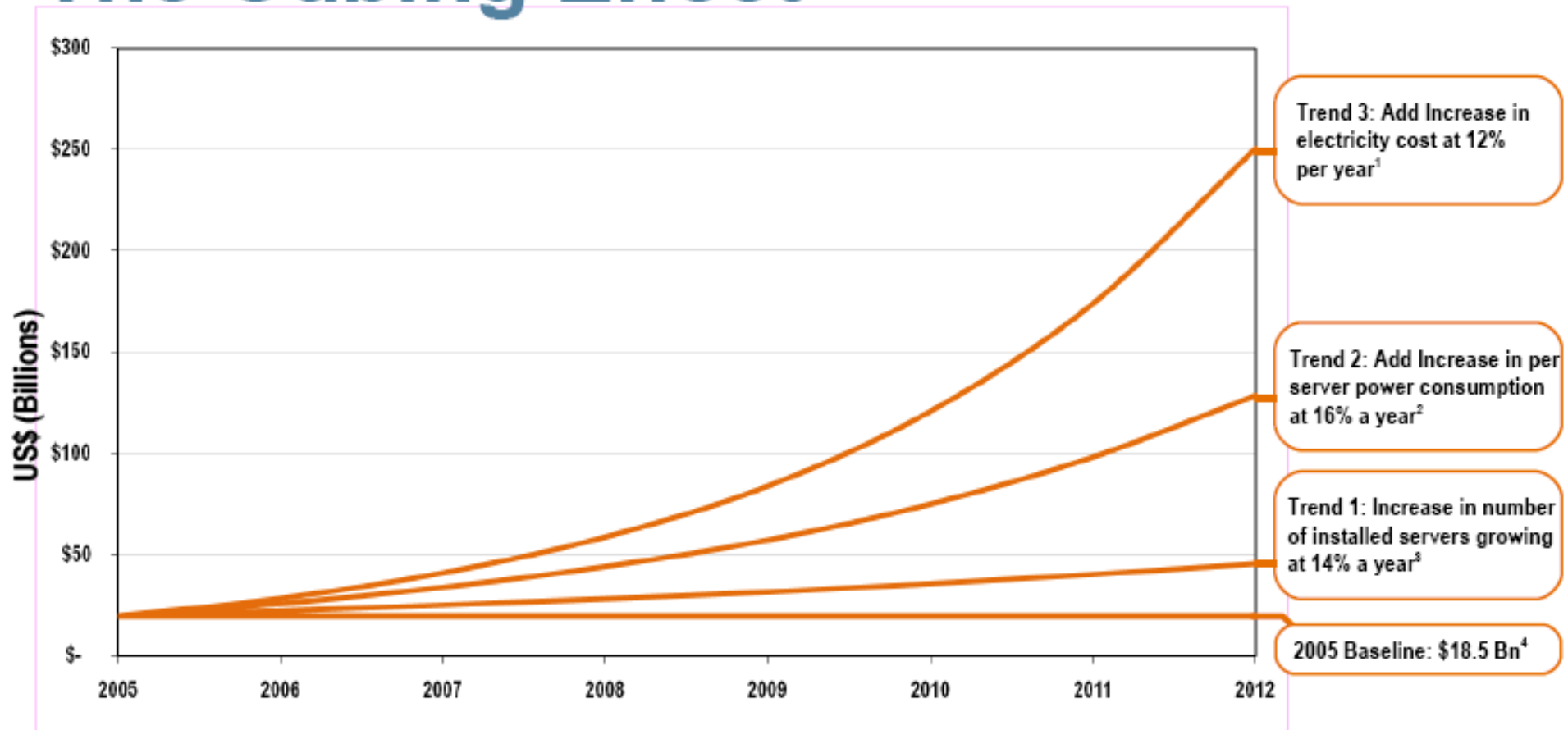- 35,000 Sq Feet Of Raised Floor
- House Up To 7 High End Systems

### Very Large Data Center
- Up to 25,000 Servers Installed
- 100,000+ Sq Feet Of Raised Floor
- Eight+ High End Systems

# Number of Servers in Data Centers Growing

| | 2008 | 2009 | 2010 | 2011 | 2012 |
|---|---|---|---|---|---|
| Very Large | 449 | 503 | 565 | 630 | 695 |
| Greenfield | 40 | 54 | 62 | 65 | 67 |
| Brownfield | 409 | 449 | 503 | 565 | 630 |
| Servers in Very Large | 11,227,357 | 12,572,577 | 14,115,316 | 15,754,862 | 17,509,290 |
| Large | 2,340 | 2,491 | 2,700 | 2,949 | 3,208 |
| Greenfield | 140 | 151 | 209 | 249 | 255 |
| Brownfield | 2,200 | 2,340 | 2,491 | 2,700 | 2,949 |
| Servers in Large | 5,265,040 | 5,603,777 | 6,075,723 | 6,635,948 | 7,295,537 |
| Medium | 3,746 | 3,987 | 4,333 | 4,714 | 5,095 |
| Greenfield | 220 | 241 | 346 | 381 | 385 |
| Brownfield | 3,526 | 3,746 | 3,987 | 4,333 | 4,714 |
| Servers in Medium | 5,994,227 | 6,379,685 | 6,932,461 | 7,542,168 | 8,196,942 |
| Small | 5,413 | 5,652 | 5,965 | 6,302 | 6,673 |
| Greenfield | 210 | 239 | 313 | 337 | 375 |
| Brownfield | 5,203 | 5,413 | 5,652 | 5,965 | 6,302 |
| Servers in Small | 2,300,426 | 2,401,961 | 2,535,233 | 2,678,327 | 2,833,263 |
| **Total Number of Servers** | **24,787,050** | **26,958,000** | **29,658,733** | **32,611,304** | **35,835,032** |

# The Cubing Effect



**US\$ (Billions)** — chart axis: \$300, \$250, \$200, \$150, \$100, \$50, \$-, years 2005–2012

**Trend 3:** Add Increase in electricity cost at 12% per year[1]

**Trend 2:** Add Increase in per server power consumption at 16% a year[2]

**Trend 1:** Increase in number of installed servers growing at 14% a year[3]

**2005 Baseline: \$18.5 Bn[4]**

By 2012 data center power consumption costs could grow to \$250B worldwide – demanding proactive energy management solutions

1. U.S. Energy Information Administration (www.eia.doe.gov)
2. Sun primary research
3. IDC#34867 U.S. and Worldwide Server Installed Base 2006-2009 Forecast (February 2006)
4. IDC Worldwide Server Power and Cooling Expense 2006-2010 Forecast

# Green Grid

**Industry and user organization focused on Energy Efficient Data Centers and Enterprise IT**

> Launched April 26$^{th}$ with 11 companies

> AMD, APC, Dell, HP, IBM, Intel, Microsoft, Rackable Systems, SprayCool, Sun Microsystems, and VMware

> Now at 40+ companies.

**Mission Statement:**

**A global consortium dedicated to advancing energy efficiency in data centers and business computing ecosystems.**

In furtherance of its mission, the Green Grid, in consultation with end-users, will:

-  Define meaningful, end-user-centric models and metrics;

-  Develop standards, measurement methods, processes and new technologies to improve performance against the defined metrics; and

- Promote the adoption of the energy efficient standards, processes, measurements and technologies.

# SPECpower

**1** Load varied in 10% decrements from 100% to 0%

**2** SSJ_ops calculated at each load point

**3** Power measured at wall socket with approved external meter at each load point

**4** SSJ_ops divided by Watts at each load point

## Benchmark Results Summary

| Performance | | | Power | Performance to Power Ratio |
|---|---|---|---|---|
| Target Load | Actual Load | ssj_ops | Average Power (W) | |
| 100% | 98.8% | 84,913 | 288 | 295 |
| 90% | 90.1% | 77,489 | 280 | 277 |
| 80% | 80.3% | 69,012 | 268 | 257 |
| 70% | 69.8% | 59,971 | 258 | 232 |
| 60% | 60.9% | 52,386 | 251 | 209 |
| 50% | 50.3% | 43,226 | 243 | 178 |
| 40% | 40.3% | 34,638 | 236 | 146 |
| 30% | 30.2% | 25,952 | 230 | 113 |
| 20% | 20.1% | 17,275 | 224 | 77.1 |
| 10% | 10.8% | 9,326 | 220 | 42.4 |
| Active Idle | | 0 | 215 | 0 |
| $\Sigma$ssj_ops / $\Sigma$power = | | | | 175 |

**Total SSJ_ops at all load points (including 0%) divided by total Watts at all load points**

**5**

**6** Final SPECpower number for system (higher is better)



Performance to Power Ratio

# Where has all the energy gone to?

# Where Does the Power in a Data Center Go?



Lighting, Security, Misc., 2%
Switch Gear / Generator, 1%
PDU Losses, 5%
UPS/Transformer Losses, 12%
CRACs/Air Movement, 10%
Humidifiers, 3%
Servers, Storage, Switches, 33%
Chillers, 34%

# Where Does Hardware Inefficiency Go?

99.9997% of **WATTS IN** becomes heat

Inefficiency (waste)

Useful Work

**Watts IN**

50%

50%

**Physical Layer**

40%

60%

**Computing Hardware**

99.999%

.001%

**Silicon**

**Business Apps**

.0003% of **WATTS IN** becomes computation

50%

60%

# Inefficiencies Create Consumption



Computing Inefficiencies

More Servers

More Power Consumption

Server Inefficiencies

More Power & Cooling

Power & Cooling Inefficiencies

**Inefficiencies drive both power consumption and material consumption**

# Cooling Basics

- All the heat generated by electronic equipment (server power) has to be removed out of the room.

- Traditional raised floor cooling can typically handle up to 5 kW per rack. This assumes:
    - > raised floor is high enough – higher than 24"
    - > no obstructions – cables, trays, etc…
    - > hot aisle / cold aisle equipment layout – servers front to front, back to back

# Data Center Mobile Hot Spots

# Power and Cooling Trends

- Raised Floor alone provide limited capabilities
- Rack Power Consumption > 10KW
- Blade Designs are increasing Density, Power, and Weight per Rack
- Close Coupling of Systems and Cooling will be required

# What's Wrong With This Data Center?

# Sun Blade Cooling Option with APC

- **Chilled Water or Refrigerant**
- **Variable Capacity Control**
- **kW Metering**
- **Front & Rear Serviceable**
- **Network Manageable**

# In-Row Chillers with Sun Blade 6000



Hot aisle air enters from rear preventing mixing

Cold air is supplied to the cold aisle

Heat captured and rejected to chilled water

InRow RC

Hot Aisle

Can Operate on hard floor or raised floor.

InfraStruXure InRow RC

© a company of Schneider Electric

APC | MGE
Critical Power and Cooling Services

# Liebert XDV



Liebert XDP/XDC

Liebert XDV
10 kW Cooling Capacity

Sensible Cooling –
No Condensation

Flexible Piping
With Quick Connects

Base Cooling

Cold Aisle

Double Stacked
Liebert XDV

Figure 4-11 Local cooling distribution using overhead cooling units mounted to the rack. ASHRAE, Datacom Equipment Power Trends and Cooling Applications, 2005. © American Society of Heating, Refrigerating and Air-Conditioning Engineers, Inc., www.ashrae.org.

Sun Blade Applications

Liebert.

EMERSON
Network Power

# Liebert XDH



**Flexible Piping With Quick Connects**

**Liebert XDP/XDC**

**Sensible Cooling – No Condensation**

**Liebert XDH 22 kW or 30 kW Cooling Capacity**

**Base Cooling**

**Cold Aisle**

**Liebert XDH**

Liebert

Sun Blade Applications

EMERSON
Network Power

# NSF TeriGrid - TACC Floorplan

*Size: approximately half a basketball court*



**2 Magnum switches**
16 line cards each
(2,304 4x IB ports each)

**82 blade compute racks**
(3,936 4S blades)

**12 IO racks**
(25 X4600 4 RU
 72 X4500 4 RU)

**112 APC Row coolers**

**1,312 12x cables**
(16 per rack)
***16 km total length***

**72 splitter cables**
6 per IO rack

12x cable lengths:    171 9m,  679 11m,  406 13m,  56 15m

Splitter cable lengths: 54 14m,    18 16m

# Extreme Datacenter – Hot Air Containment

**Rack Air Containment (RACS)**

**Hot Aisle Containment (HACS)**

# The World's First Containment Datacenter

## Module Data Center

# Cooling

- Air flows in circular path with fans and heat exchanger per rack
  - > Payload installed front to back

- Chiller size depends on the payload, 60-ton chiller for max 200kW load

# Steam and Ice

*"Today a data center just looks like a giant resistor in a multimegawatt circuit. It would be nice if it also was a capacitor."*

– Mark Bramfitt, PG&E

# Earth Pipes

# Disjointed Progress in the Data Center

**"Facilities is from Mars, IT is from Venus"**

## What Facilities is Doing

- Hot aisle containment
- Cold aisle containment
- Concrete slab floor
- Variable frequency drives
- Air side economizers

## What IT is Doing

- Server refresh
- Consolidation
- Virtualization
- Utilization Management

# Newthinking in Systems/Grid Design

# Introducing Niagara

- SPARC V9 implementation

- Up to eight 4-way multi-threaded cores for up to 32 simultaneous threads

- All cores connected through a 90GB/sec crossbar switch

- High-bandwidth 12-way associative 3MB Level-2 cache on chip

- 4 DDR2 channels (23GB/s)

- Power : < 70W !

- ~300M transistors

- 378 sq. mm die



1 of 8 Cores

DDR-2 SDRAM
DDR-2 SDRAM
DDR-2 SDRAM
DDR-2 SDRAM

L2$ L2$ L2$ L2$

FPU

Xbar

C1 C2 C3 C4 C5 C6 C7 C8

Sys I/F Buffer Switch Core

BUS

# The Power of CMT

Niagara Processor
Utilization: Up to 85%



**Chip Multi-threaded (CMT) Performance**

Thread 4
Thread 3
Thread 2
Thread 1

Time

■ Memory Latency    ■ Compute

# Niagara's Power Advantage

## "Cool Threads" Dramatically Reduce Power Consumption

Uses a Fraction of the Power/Thread vs. Xeon



Single-core Processor    *(Size Not to Scale)*    CMT Processor

# An integer linear programming (ILP) based static scheduling method that minimizes both thermal hot spots and temperature gradients to increase MPSoC reliability

### TABLE III
### ILP FORMULATION FOR MIN-TH&SP

Minimize $H + G$;

$H = max\{Q_p; p = 1...m,$ for a system of $m$ cores$\}$ where:

$$Q_p = \sum_{T_i \in T} \{x_{ip} \sum_{v_k} (y_{ik} q_{ik})\}$$

$$G = \sum_{p,r \in PU, p \neq r} \{n_{pr}\{ \sum_{i,j \in T, i \neq j} x_{ip} x_{jr} [p_{ij} d_{ij} (\tau_i - s_j) + p_{ji} d_{ji} (\tau_j - s_i)]\}\}$$

Subject to constraints:

| | |
|---|---|
| (a) $\forall T_i : \sum_p x_{ip} = 1$ | Each task is assigned to only one PU |
| (b) $\forall T_i : \sum_k y_{ik} = 1$ | Each task runs at only one V/f level |
| (c) $\tau_i = s_i + t_i$ | Execution finish time for $T_i$ |
| (d) $s_i \geq max_{E_{ji} \in E}\{\tau_j\}$ | Task precedence |
| (e) $\tau_i \leq D_i$ | Deadlines for all sink nodes |
| (f) $s_i \geq \tau_j; \ if \ p_{ji} = 1;$ | Precedence for tasks on the same core |
| (g) $p_{ij} + p_{ji} = 1;$ $\quad$ if $x_{ip} = x_{jp} = 1$ | If $T_i$ and $T_j$ are scheduled on the same core, either $T_i$ precedes $T_j$, or vice versa |



. Thermal maps: (a)DLB; (b)Adaptive-Random

Legend: >85°C, [75, 85] °C, <75°C

Source : Ayse Coskun, Tajana Simuni´c Rosing, Keith A. Whisnant, and Kenny C. Gross, Static and Dynamic Temperature-Aware Scheduling for Multiprocessor SoCs

# Desired Power/Performance Curve



100%

Power
Consumed

Near optimal linear performance

0%

25%   50%   75%   100%

Throughput

# The Reality

# Desired Behavior of Future 4 Socket Servers

# Getting to Greater Linearity

# Algorithm Profile

# Congugate Gradient Sparce Solver
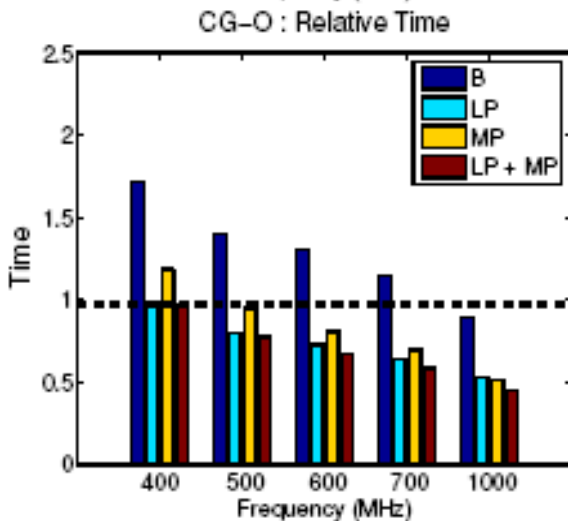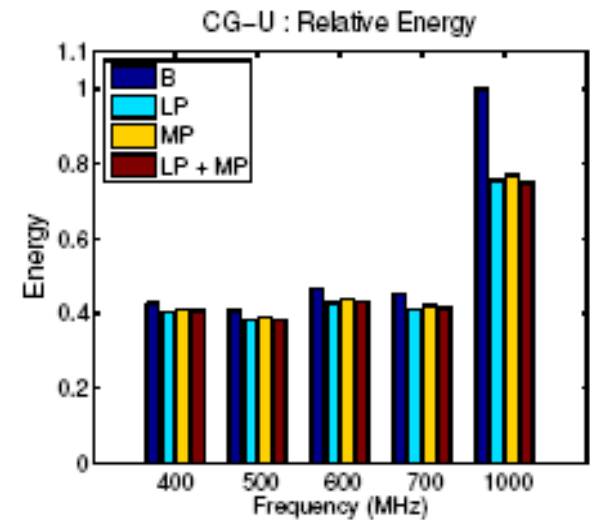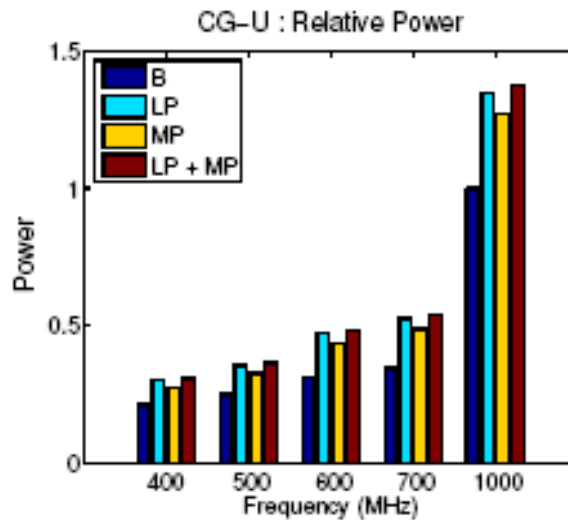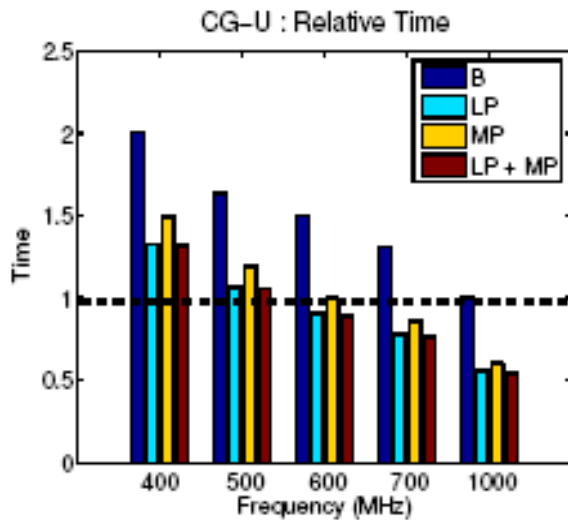
Source : Korad Malkowski, Ingyu Lee, Padma Raghavan, Mary Jane Irwin, Conjugate Gradient Sparse Solvers: Performance-Power Characteristics
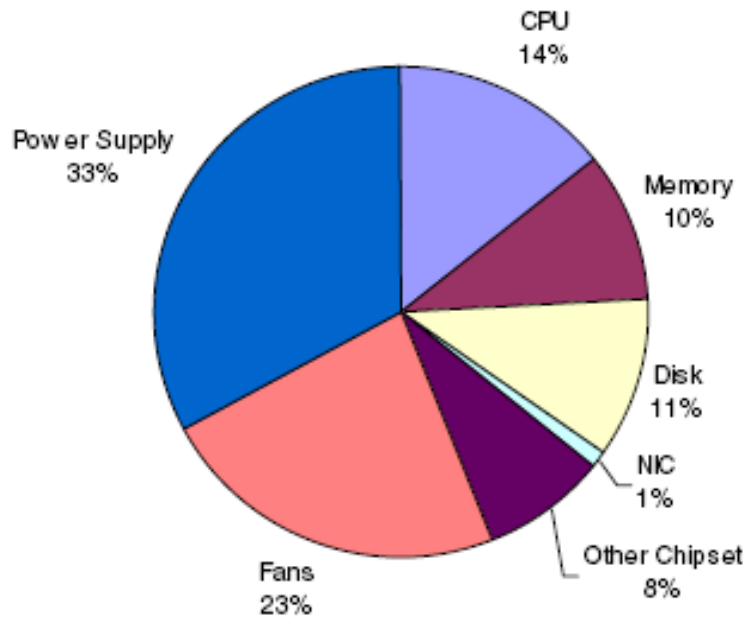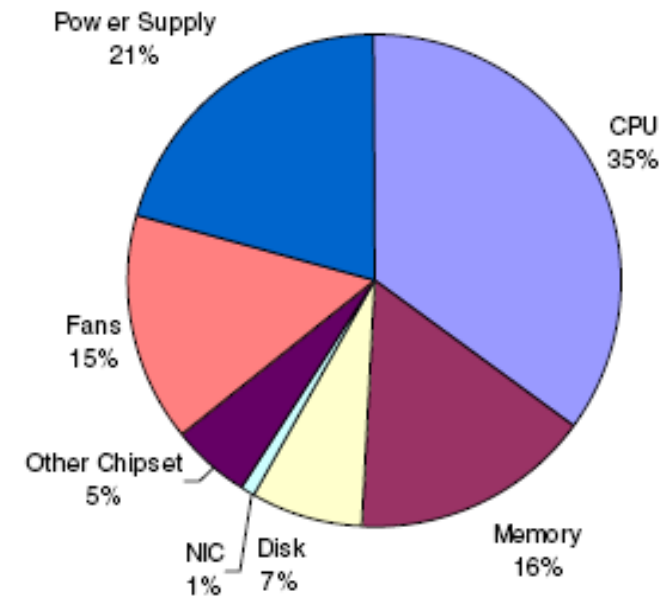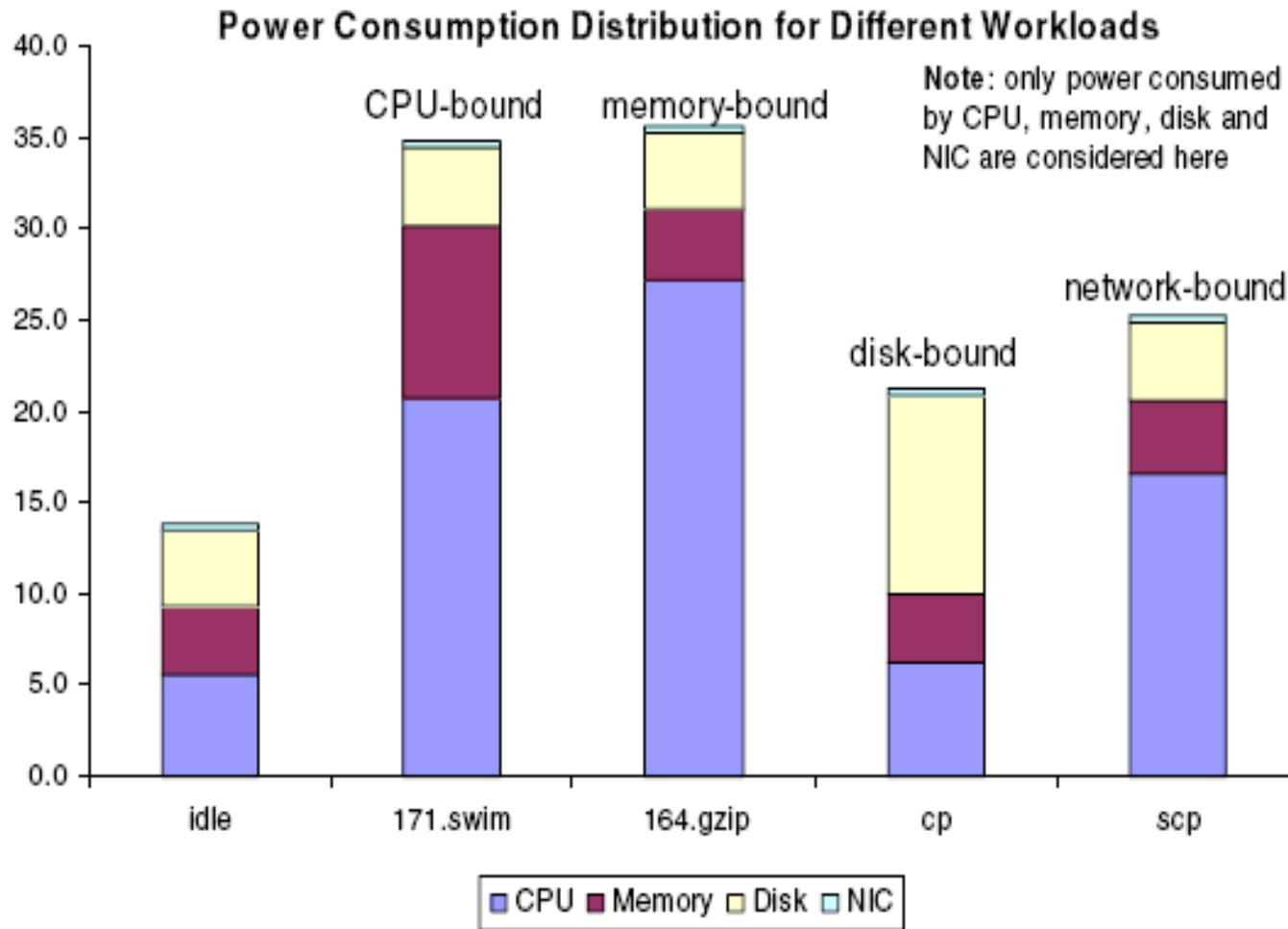
# Application Power Profiling

# Application Power Profile



Power consumption distribution for system idle
System Power: 39 Watt

CPU 14%
Memory 10%
Disk 11%
NIC 1%
Other Chipset 8%
Fans 23%
Power Supply 33%

Power consumption distribution for memory performance bound (171.swim)
System Power: 59 Watt

Power Supply 21%
CPU 35%
Memory 16%
Disk 7%
NIC 1%
Other Chipset 5%
Fans 15%

Source： Xizhou Feng, Rong Ge, Kirk W. Cameron, University of South Carolina, Columbia, SC 29208,
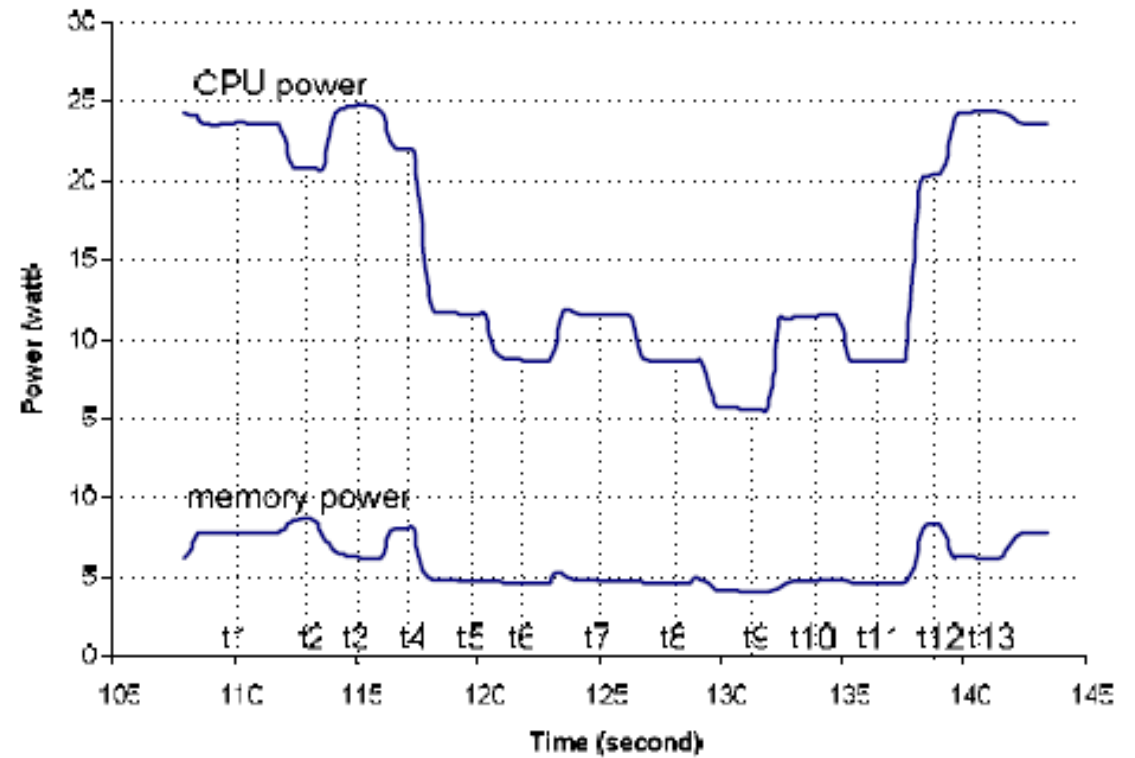Power and Energy Profiling of Scientific Applications on Distributed Systems

Source : Xizhou Feng, Rong Ge, Kirk W. Cameron, University of South Carolina, Columbia, SC 29208, Power and Energy Profiling of Scientific Applications on Distributed Systems
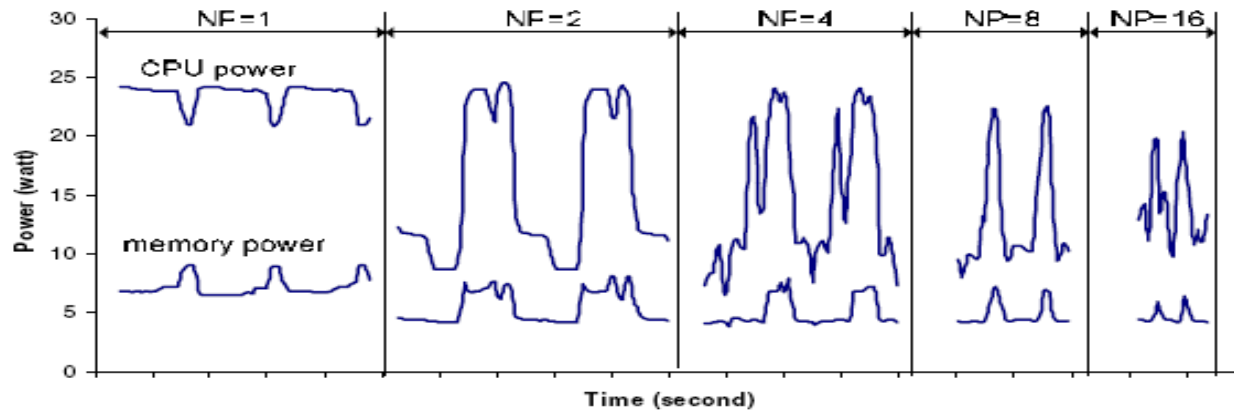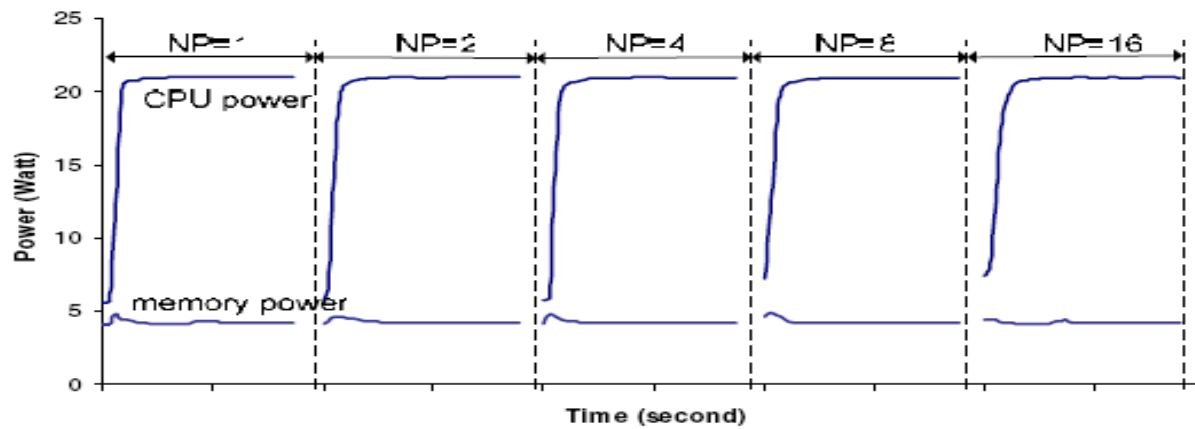
Power Profile of FT Benchmark (class B NP=4)

Expanded View of Power Profile of FT (class B NP=4)

Power Profile of FT Ber

Source : Xizhou Feng, Rong Ge, Kirk W. Cameron, University of South Carolina, Columbia, SC 29208,
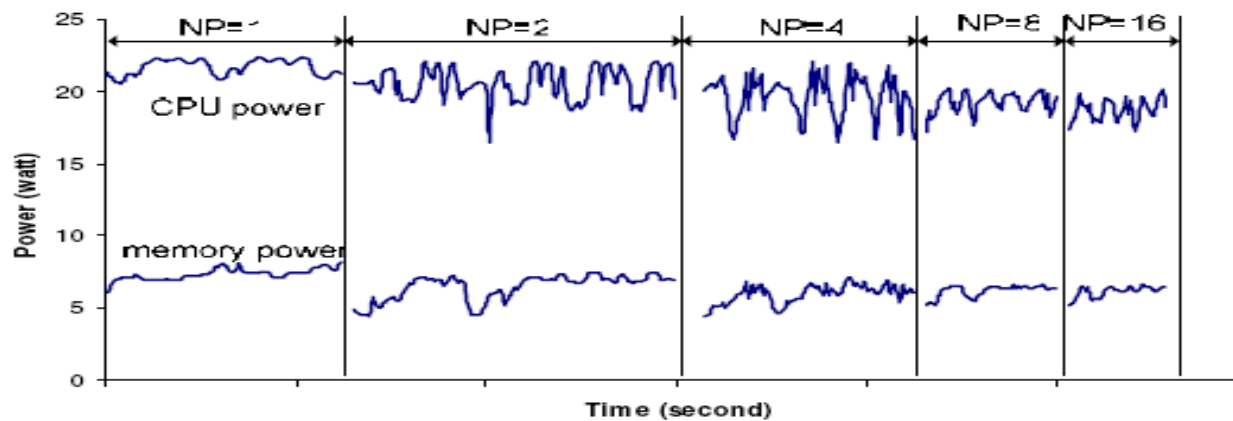Power and Energy Profiling of Scientific Applications on Distributed Systems

Power Profile of FT Benchmark (class A) with Different Number of Nodes

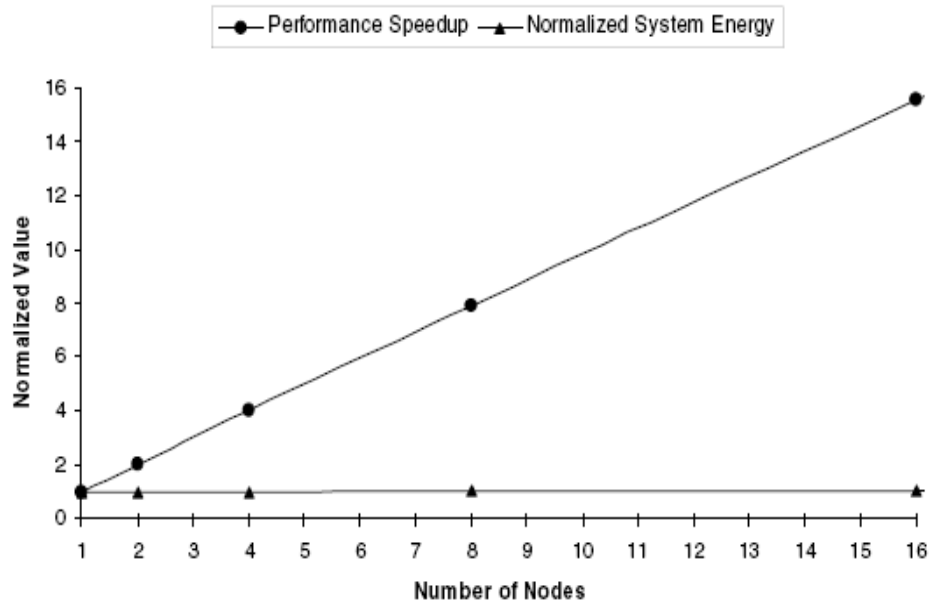Power Profile of EP (class A) with Differentnumber of Nodes

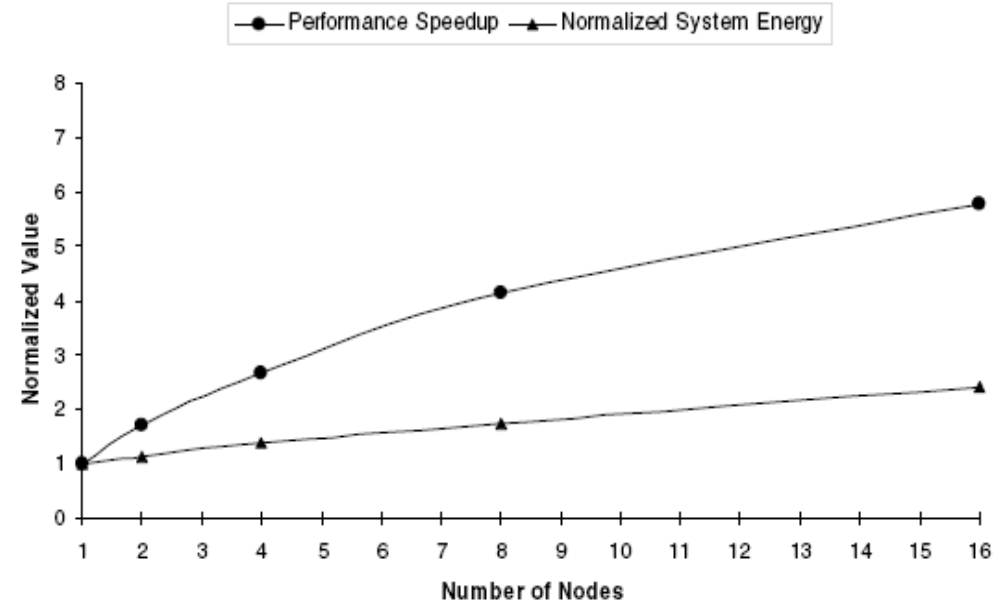Power Profile of MG (Class A) with Different Number of Nodes

# Observation of NPB

- CPU power consumption decreases as memory power consumption goes up;

-  Both CPU power and memory power decrease with message communication among different nodes;

-  For most parallel codes (except EP), the average power consumption goes down as the number of nodes increases;

- Communication distance and message size affects the power profile pattern (for example, LU has short and shallow power consumption in contrast with FT).
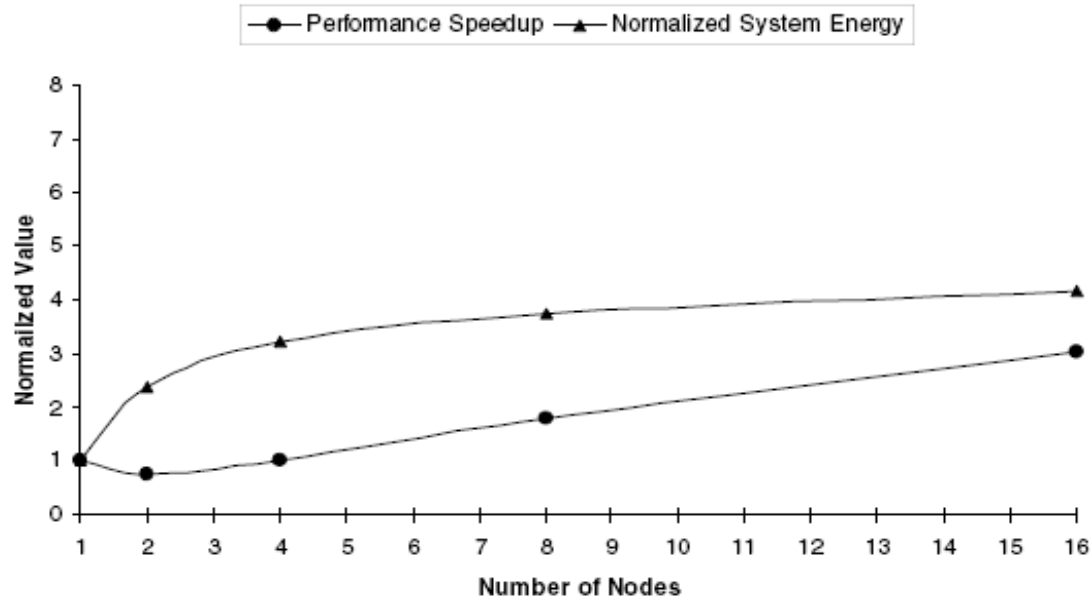
Performance and Energy Consumption for EP (class A) code

Performance and Energy Consumption for MG (class A) code

Performance and Energy Consumption for FT (class A) code

# AMD Opteron 2218 : DVFS

| Frequency (MHz) | Voltage (V) |
|---|---|
| 1000 | 1.10 |
| 1800 | 1.15 |
| 2000 | 1.15 |
| 2200 | 1.20 |
| 2400 | 1.25 |
| 2600 | 1.30 |

| Code | Frequency (MHz) | | | | | |
|---|---|---|---|---|---|---|
| | 1000 | 1800 | 2000 | 2200 | 2400 | 2600 |
| BT.C.16 | 1.66 | 1.17 | 1.08 | 1.07 | 1.05 | 1.00 |
| | 1.06 | 0.88 | 0.84 | 0.90 | 0.96 | 1.00 |
| CG.C.16 | 1.47 | 1.15 | 1.11 | 1.07 | 1.03 | 1.00 |
| | 0.98 | 0.88 | 0.88 | 0.91 | 0.94 | 1.00 |
| EP.C.16 | 2.57 | 1.45 | 1.30 | 1.18 | 1.08 | 1.00 |
| | 1.57 | 1.07 | 1.00 | 0.98 | 0.98 | 1.00 |
| FT.C.16 | 1.40 | 1.10 | 1.06 | 1.04 | 1.02 | 1.00 |
| | 0.92 | 0.84 | 0.83 | 0.88 | 0.94 | 1.00 |
| IS.C.16 | 1.52 | 1.07 | 0.99 | 1.01 | 1.01 | 1.00 |
| | 1.01 | 0.82 | 0.79 | 0.85 | 0.93 | 1.00 |
| LU.C.16 | 1.62 | 1.13 | 1.05 | 1.02 | 1.06 | 1.00 |
| | 1.03 | 0.86 | 0.83 | 0.86 | 0.96 | 1.00 |
| MG.C.16 | 1.41 | 1.11 | 1.03 | 1.05 | 0.99 | 1.00 |
| | 0.92 | 0.84 | 0.81 | 0.87 | 0.90 | 0.98 |
| SP.C.16 | 1.53 | 1.08 | 1.03 | 1.02 | 1.05 | 1.00 |
| | 1.00 | 0.84 | 0.81 | 0.87 | 0.96 | 1.00 |

# Intel 1.4 Ghz Pentium-M : FT Benchmark



Source :K. Cameron, Rong Ge, Xizhou Feng High-Performance, Power-Aware Distributed Computing for Scientific Applications

# CPU Miser

# Green Destiny



GREEN DESTINY: LOW-POWER SUPERCOMPUTER

Only Difference? The Processors

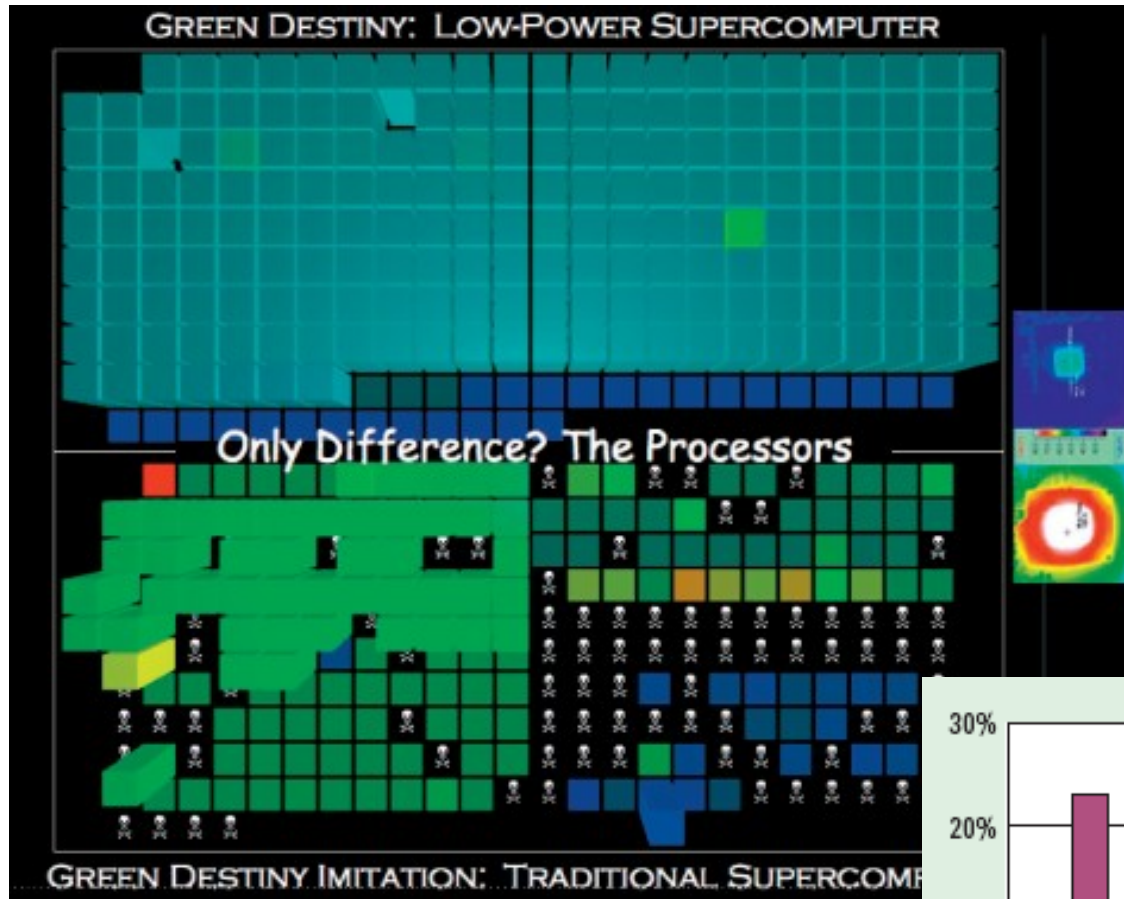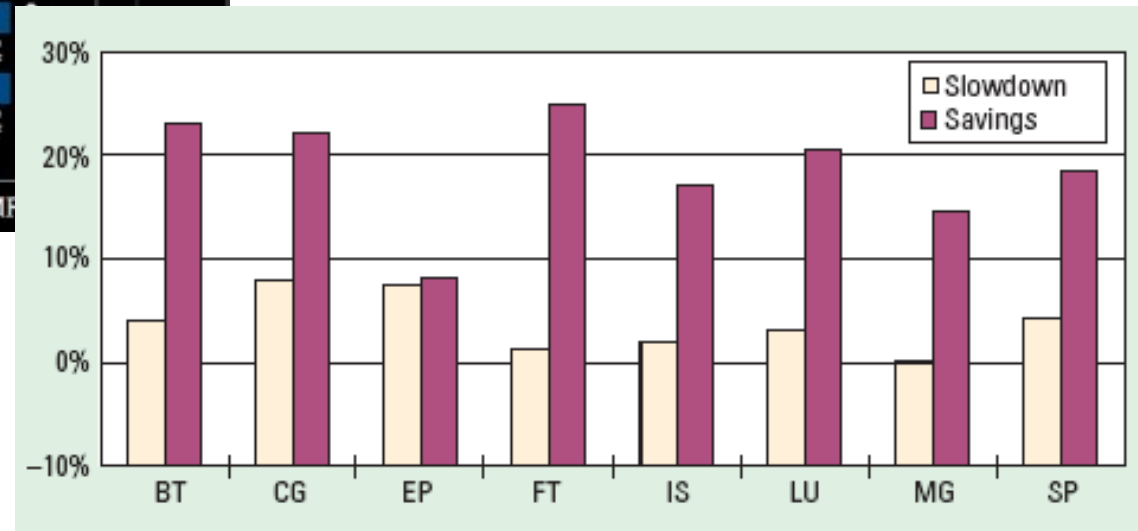GREEN DESTINY IMITATION: TRADITIONAL SUPERCOMP...

**Table 1. Comparison of supercomputing systems on the Linpack benchmark. ASCI White's top performance is shown in italics; Green Destiny's appears in bold.**
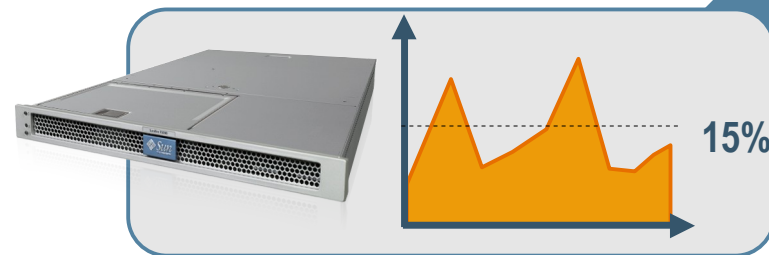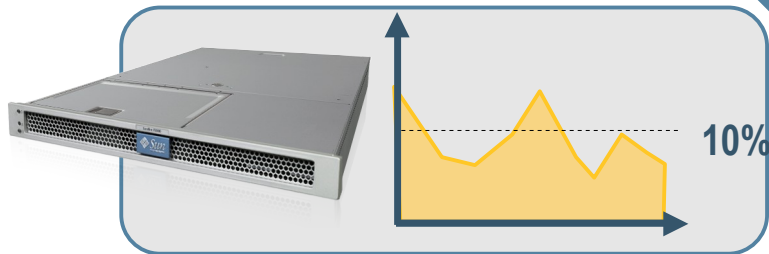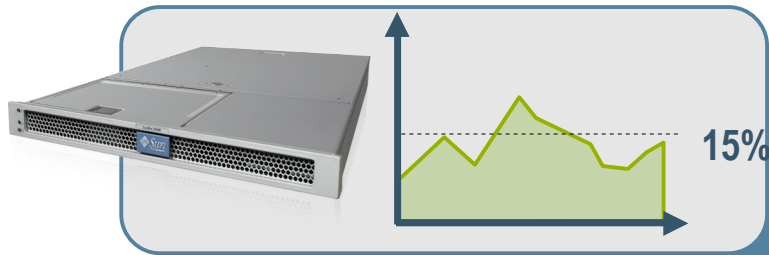
| Performance Metric | ASCI White | Green Destiny |
|---|---|---|
| Year | 2000 | 2002 |
| Number of processors | 8,192 | 240 |
| Performance (Gflops) | *7,226* | 101 |
| Space (ft$^2$) | *9,920* | **5** |
| Power (kW) | *2,000* | **5** |
| DRAM (Gbytes) | *6,200* | 150 |
| Disk (Tbytes) | *160.0* | **4.8** |
| DRAM density (Mbytes/ft$^2$) | 625 | **30,000** |
| Disk density (Gbytes/ft$^2$) | 16.1 | **960.0** |
| Perf/space (Gflops/ft$^2$) | 0.7 | **20.2** |
| Perf/space (Gflops/kW) | 4 | **20** |
| Reliability (hours) | 5.0 hours (2001), 40 hours (2003) | **No unscheduled downtime** |



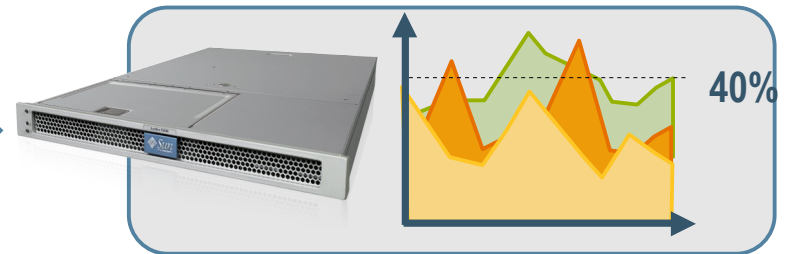Source: Wu-chun Feng, Xizhou Feng, and Rong Ge, Green Computing Comes of Age

# Virtualization & Grid

# Efficient Resource Utilization: Migration
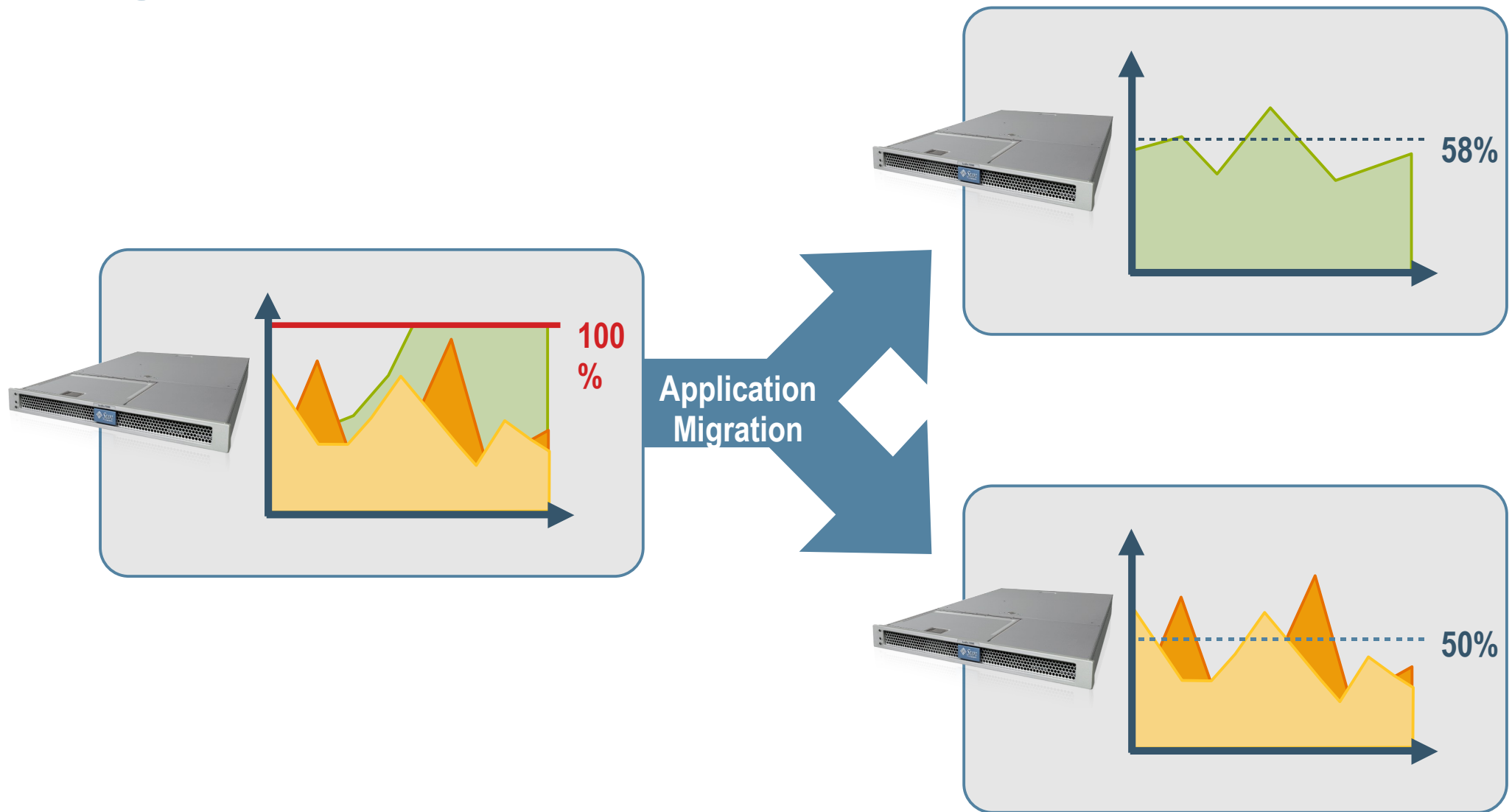
# Two Dimensions of Virtualization

**Make a data center with N servers...**

**Partitioning**

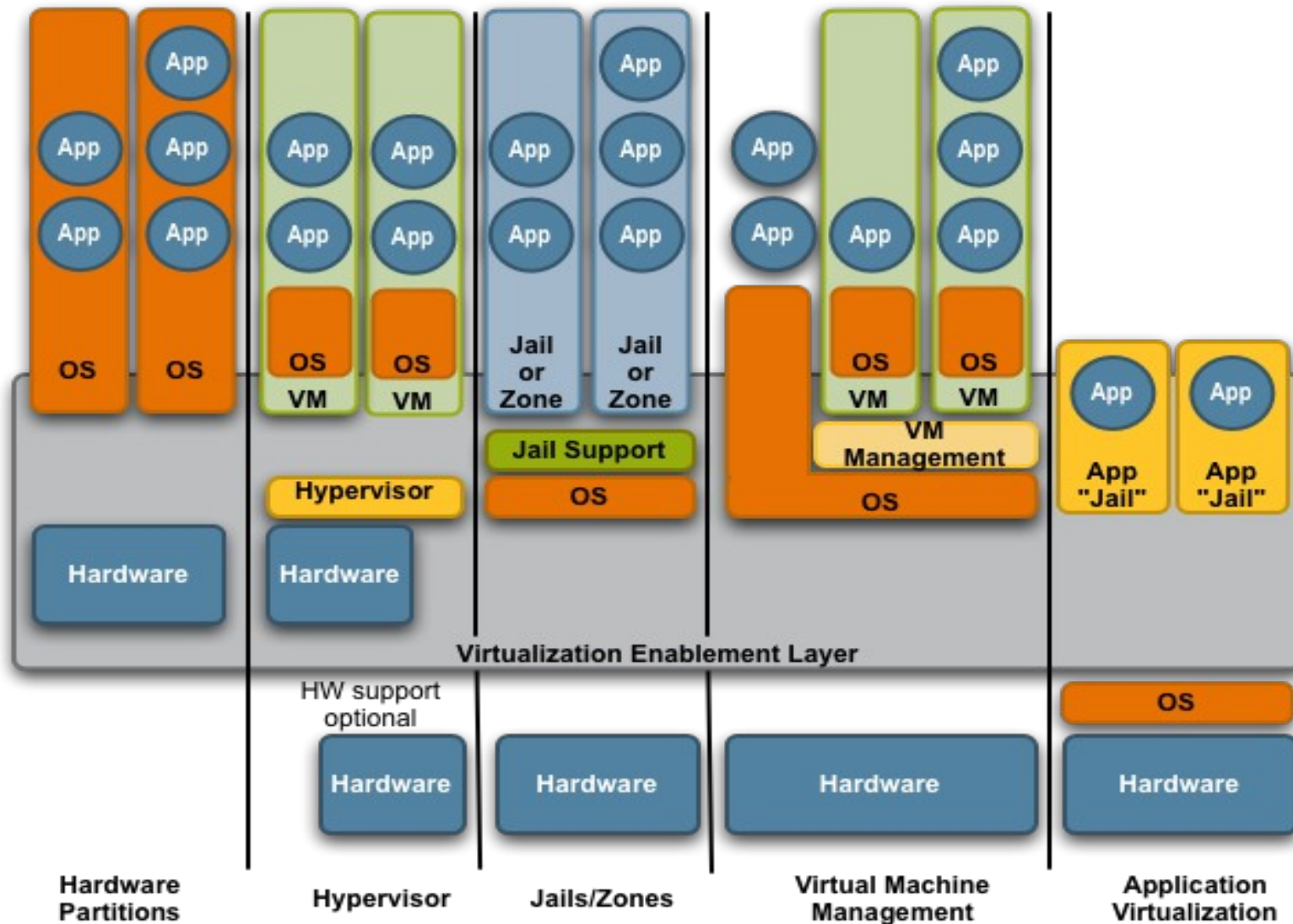**Aggregation**

**...look like it has >>N servers for utilization**

**...look like it has <<N servers for administration**

# Virtualization

# Virtualization Technology

**Hardware Partitions**

| Technology | Vendor |
| --- | --- |
| Dynamic Systems Domain | Sun |
| LPAR | IBM |
| VPAR NPAR | HP |

**Hypervisor**

| Technology | Vendor |
| --- | --- |
| XVM Server | Sun |
| Virtual Infr 3 (ESX) | Vmware |
| Xen | Xensource & Sun |
| Viridian | Microsoft |
| Logical Domains | Sun |
| KVM (Linux | Community |
| VM | IBM |

# Virtualization Technology

**OS Virtualization**

| Technology | Vendor |
|---|---|
| Solaris Containers/Zones | Sun |
| IBM Wpars | IBM |
| BSD Jails | HP |
| Virtuozzo | Swsoft |
| OpenVZ | Community |

**Application Virtualization**

| Technology | Vendor |
|---|---|
| Etude | Sun |
| Trigence | Trigence |
| Softgrid | Softricity |
| SVS – Software Virt Soln | Alteris |
| Logical Domains | Sun |
| Project Tarpon | Citrix |

# The Re-Entrant Grid



## The Holy Grail

Provision freeze-dried appliances in the most energy-efficient place in the data center

Software Appliance "Factory" Rack
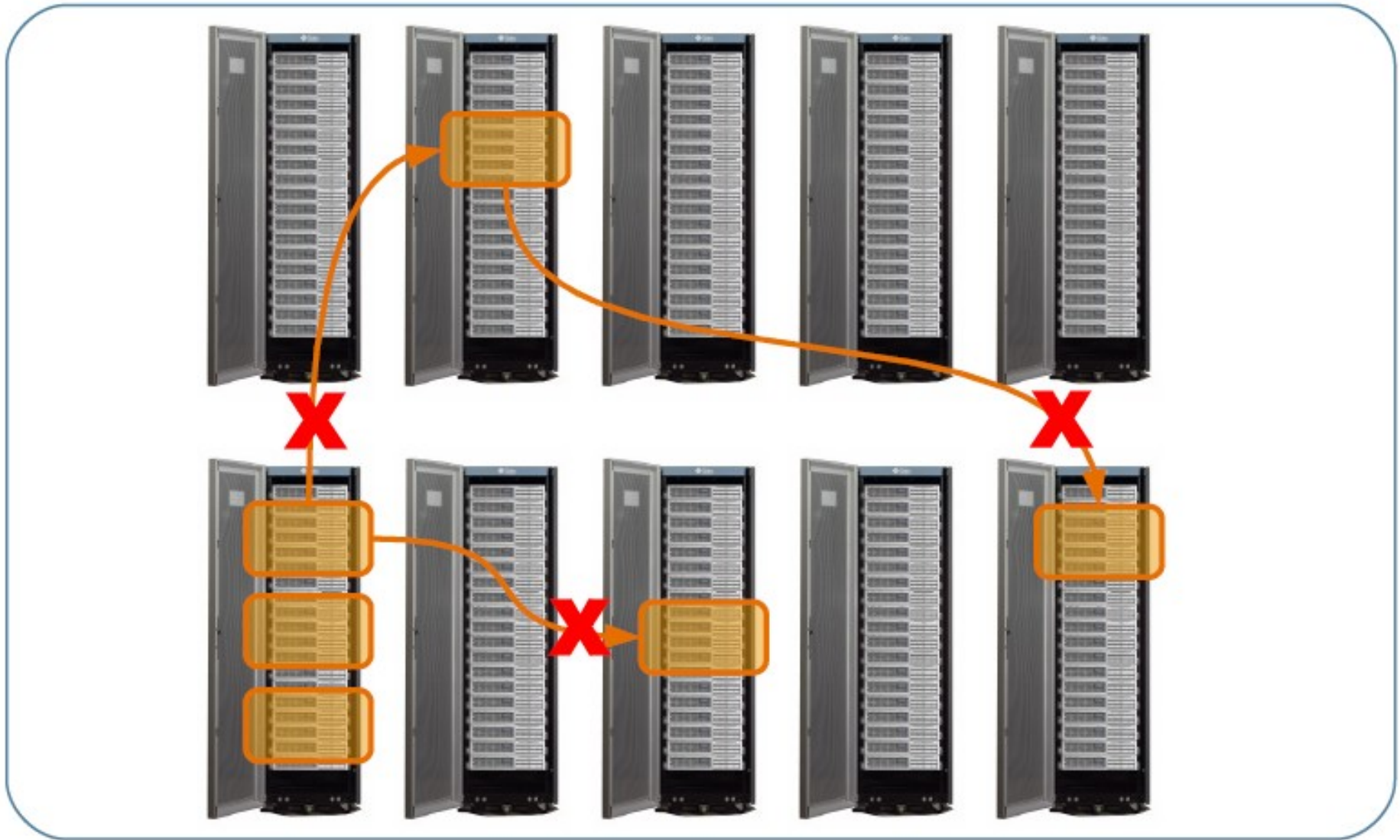
Migrate work-loads for energy efficiency reasons anywhere in the data center

# Why we cannot do this easily

**Because The Network Needs To Get Involved**

# Why we cannot do this easily

**Because The Network Needs To Get Involved**

**Corollary** : Energy efficiency isn't just a chip or hardware problem.
It is a Grid management problem,  a systems management problem,
an OS problem, a networking problem
a virtualisation problem,  a data grid problem (storage).

# What the community is up to?

- Spec power and performance
- Green Top 500
- Green Grid
- US Congress passed law 109-431
- EPA Report
- Others

# Ackwoledgement

- Subodh Bapat,  Sun Microsystems Inc.
- John Fragalla,  Sun Microsystems Inc.
- Dave Douglas, Sun Microsystems Inc.
- Many others