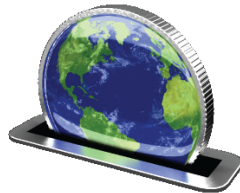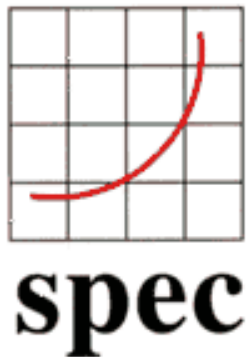# The Powers That Be (in HPC)

Kirk W. Cameron

Computer Science

Virginia Tech

# ENA-HPC Street Credit

- Over $6M related federal funding (since '04) (NSF, DOE, SBIR, IBM, Intel, and others)
- EPA Energy Star for servers (since '05)
- SPECPower Founding Member (since '05)
- Co-founder Green500 (since '06)
- Green IT Columnist (*IEEE Computer*)
- CEO and Founder, MiserWare Inc. (since '07)
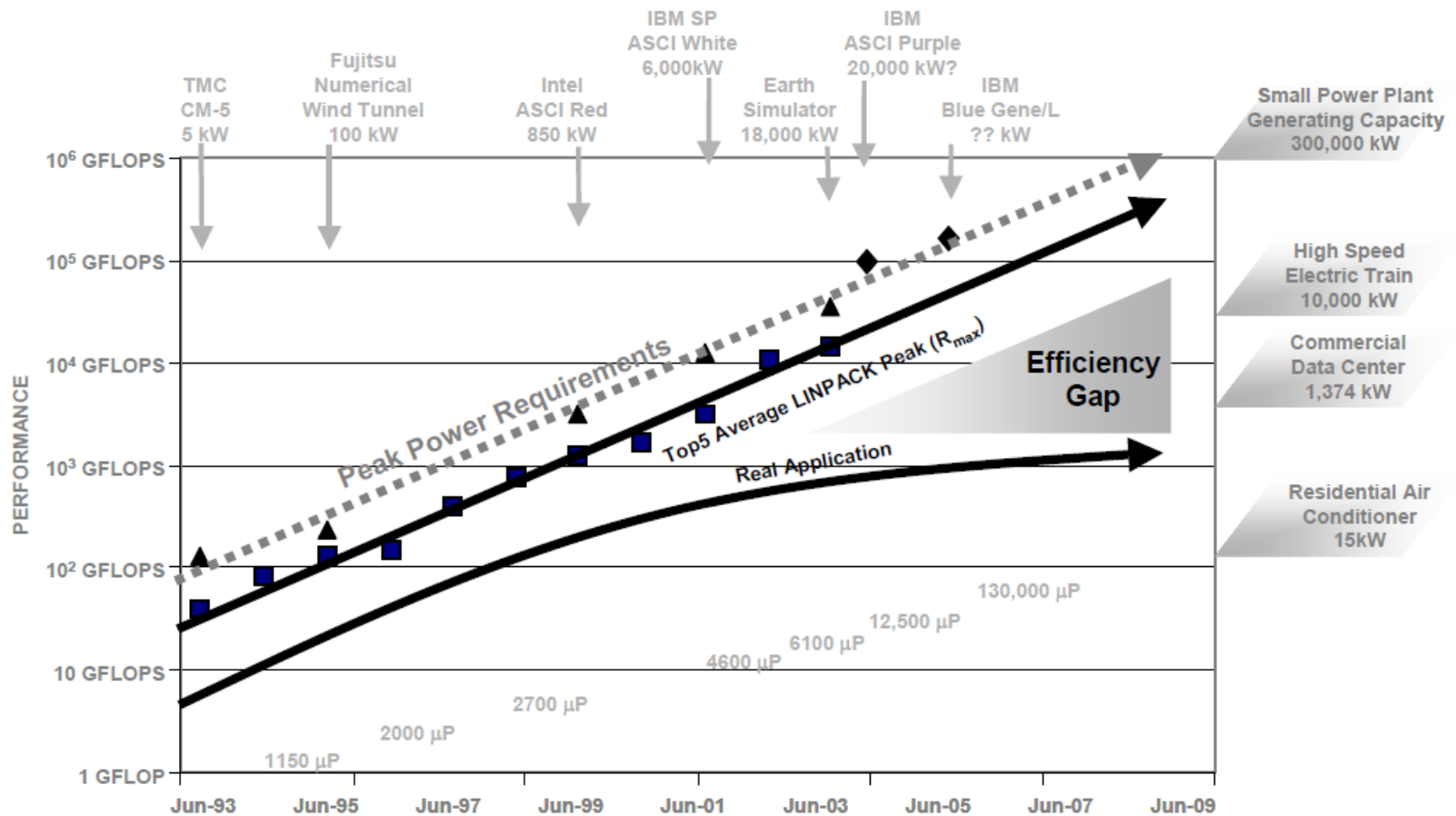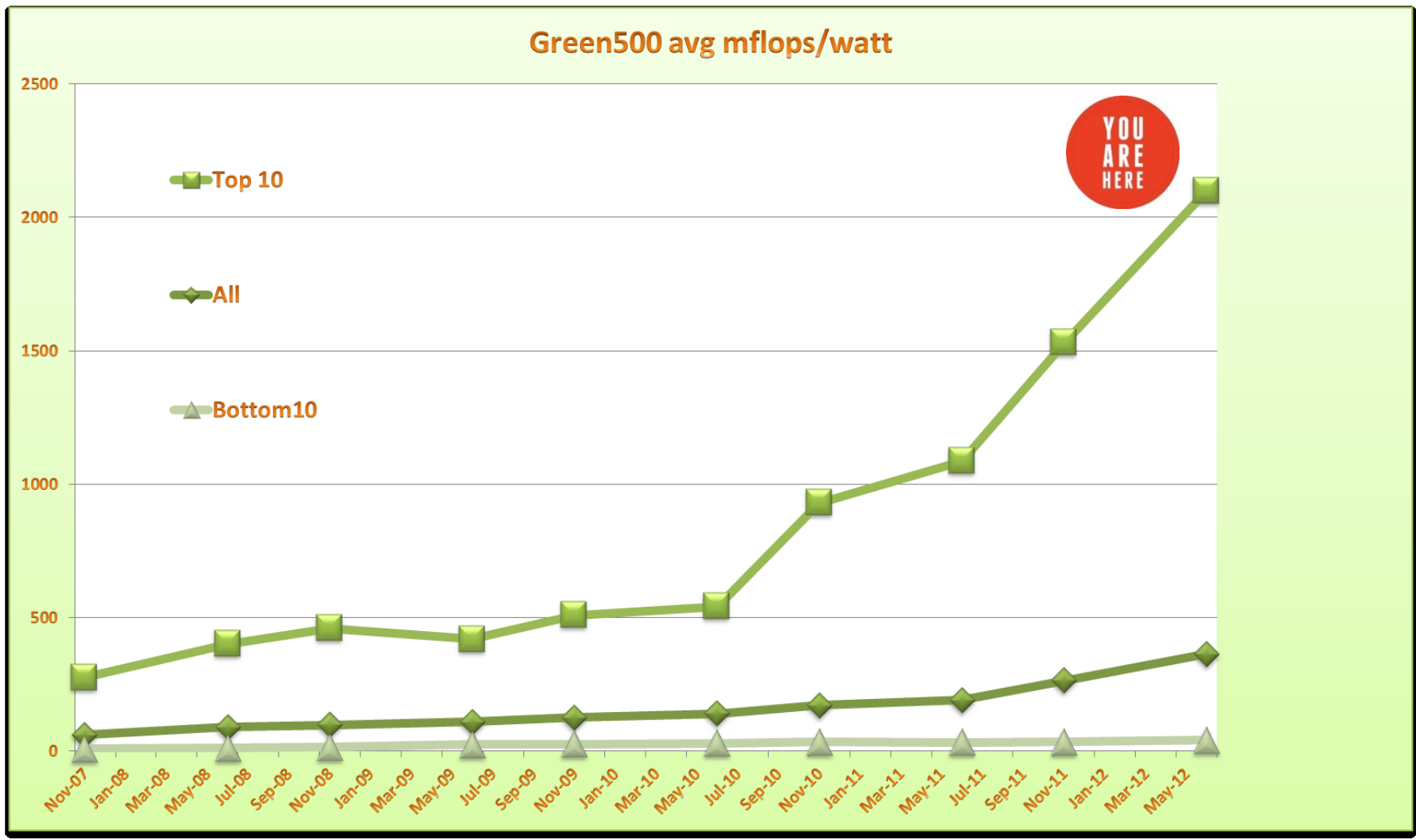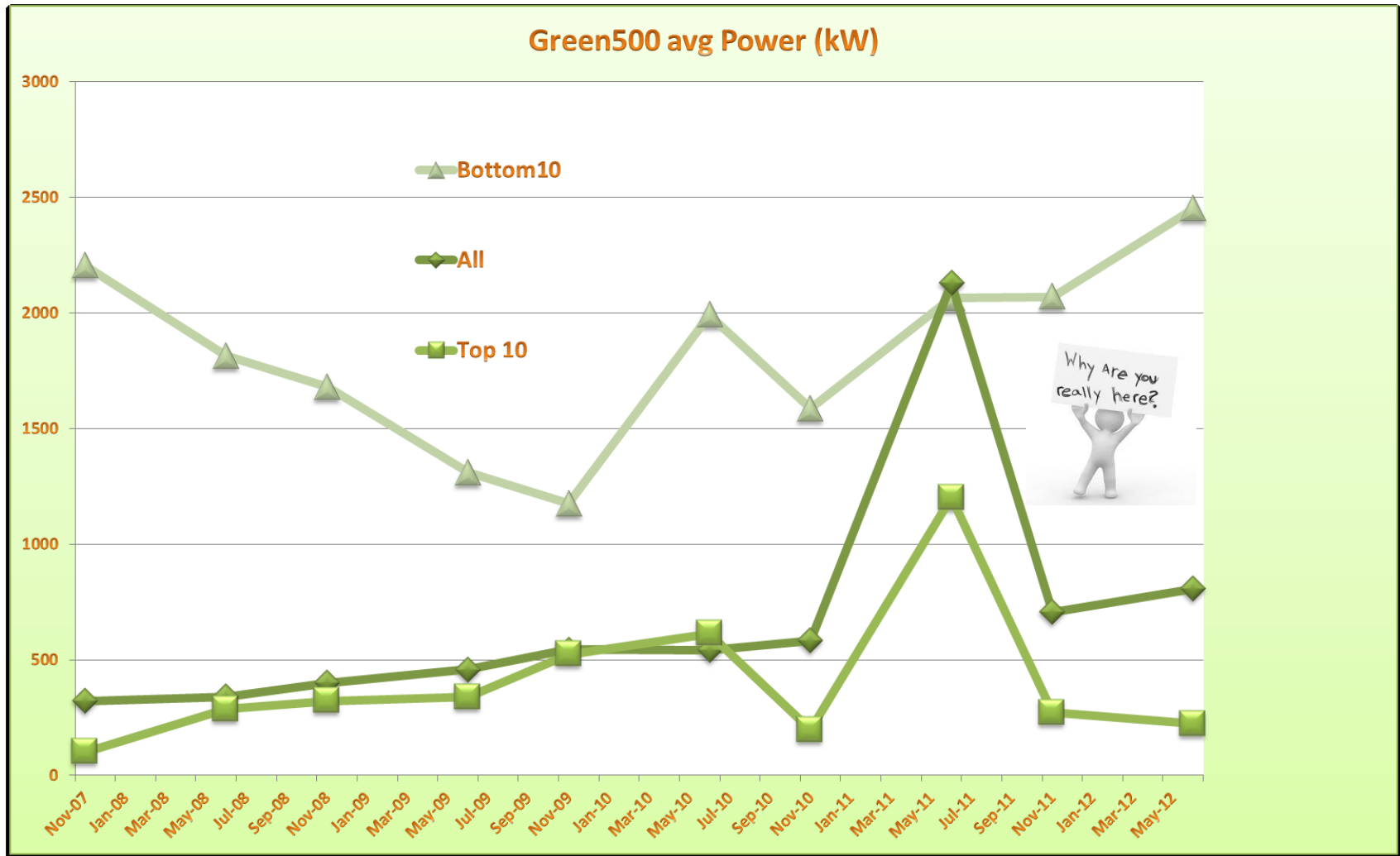
# The way we were (circa 2003)



Fig. 1 Power-performance trends in the supercomputer industry. The computational demands of scientific applications have led to exponential increases in peak system performance (shown as average of peak LINPACK measurements), system power consumption (shown for several supercomputers), and

# You are here (September 2012)

# Or are you really here?

# Getting there...

From 2007-2012...

[6x ↑ Flops/watt]

[~2.5x ↑ power consumption]

[Commodity systems catch efficiency of top 10 in 18 mo.]

Projections for 2012-2019...

[2100 to ~15,000 MFlops/Watt]
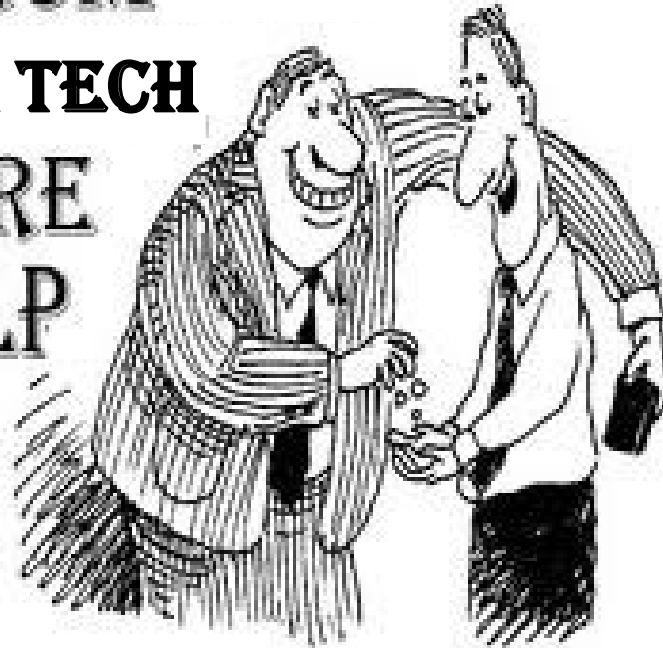
[66 kW for 1 Petaflop System]

[66 MW for 1 Exaflop System}

[Need 50,000 Mflops/Watt for 1 Exaflop @ 20 MW by 2019!!!]

# Conclusion: We need help.

# How can we...help you...help us...

# What do we need...?



**<u>Insight</u>**

**Where does energy go?**



**<u>Understanding</u>**

**Why does energy go?**



**<u>Action</u>**

**What can we do?**

# SCAPE Research (circa 2002)

- ## My observations
  - ### Power will become disruptive to HPC
  - ### Laptops outselling PC's
  - ### Commercial power-aware not appropriate for HPC

**$800,000 per year per megawatt!**

$4,000/yr

$12,000/yr

$680,000/yr

$8 million/yr

$9.6 million/yr

TM CM-5
.005 Megawatts

Residential A/C
.015 Megawatts

Conventional Power Plant
300 Megawatts

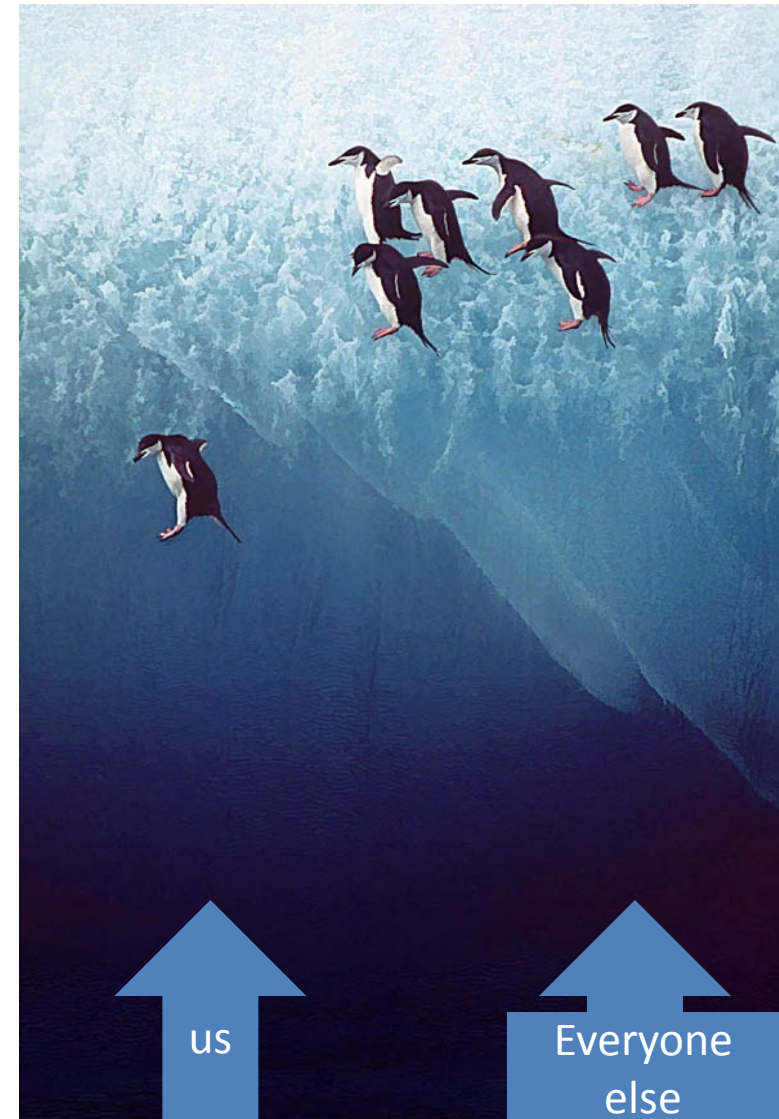Intel ASCI Red
850 Megawatts

High-speed train
10 Megawatts

K Supercomputer
12 Megawatts

# SCAPE Launches HPPAC

- **High-performance, Power-aware Computing**
  - **Maintain Performance**
  - **Reduce energy waste**
- **Measurement tools**
- **No funding initially**



us

Everyone else

12

# We were right! Whew.

## IT confronts the datacenter power crisis

As energy costs escalate, conserving resources tops the list of challenges for today's IT managers

By Dan Goodin
October 06, 2006

E-mail    Printer Friendly    Reprints    Slashdot It!

When David Young told his colocation provider late last year that his online applications startup, Joyent, planned to add 10 servers to its 150-system datacenter, he received a rude awakening. The local power utility in Southern California wouldn't be able to provide the additional electricity needed. Joyent's upgrade would have to wait.

## In the Data Center, the Heat Is On

Halamka John    Today's Top Stories ▸    or  Other Servers Stories ▸

October 23, 2006 (Computerworld) -- I recently began a project to consolidate two dat

## Data Center Budgets Face Radical Changes

Consortium head says facilities costs are surpassing the price of hardware

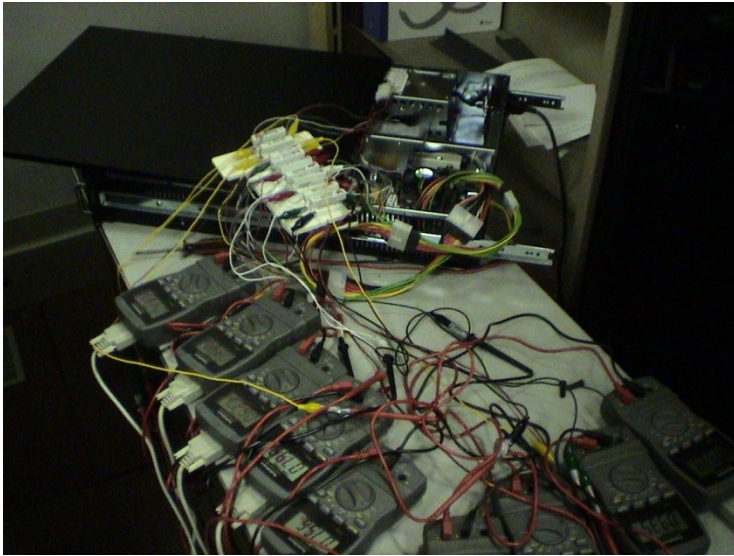Patrick Thibodeau and Patrick Thibodeau    Today's Top Stories ▸    or  Other IT Management Stories ▸

October 30. 2006 (Computerworld) -- The business value arising from Moore's Law, which says the number of

"You can only manage what you can measure."

*Peter Drucker, writer*
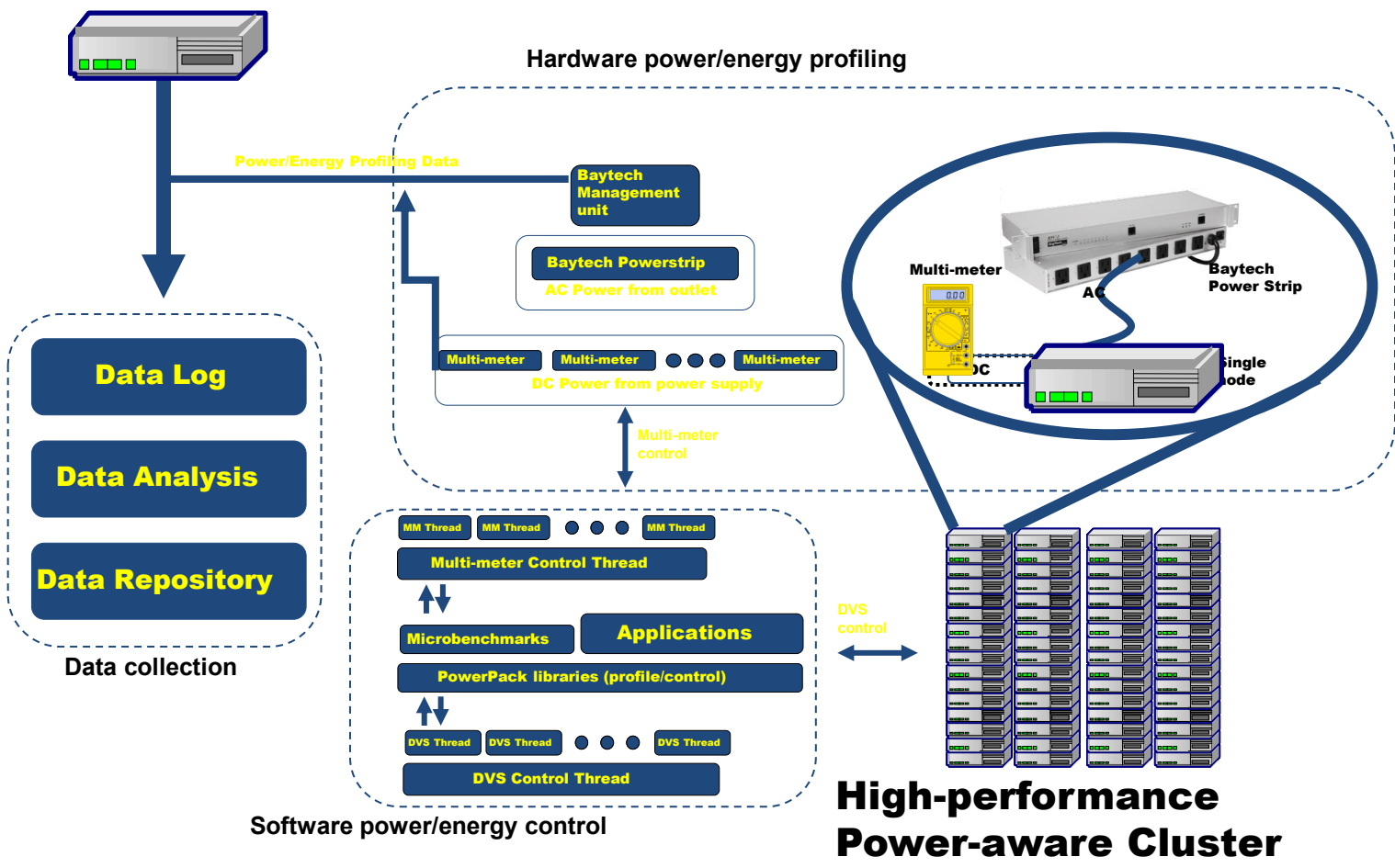
# Measuring power is "tough"

# HPPAC Tools

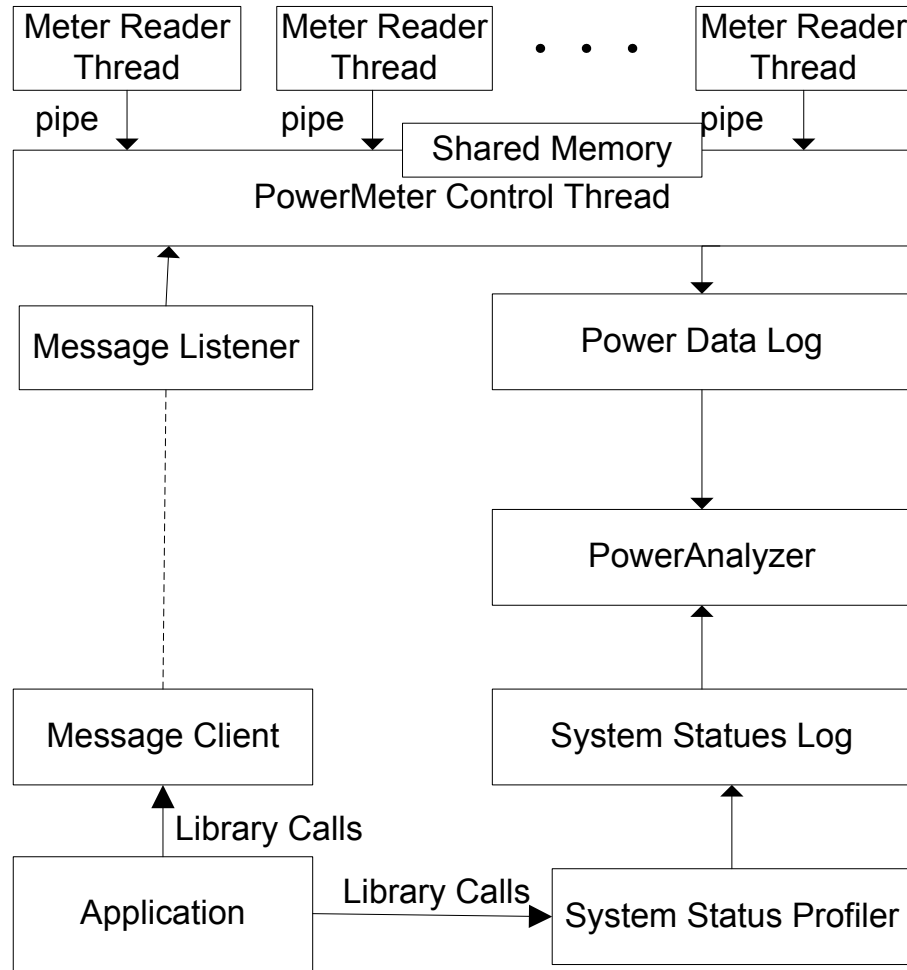- ## PowerPack
  - Modularized software + HW sensors
  - Extended analytics for applicability
  - Extended to support thermals

- ## SysteMISER (evolves to MiserWare/Granola)
  - Improved analytics to weigh tradeoffs at runtime
  - Automated cluster-wide, DVS scheduling
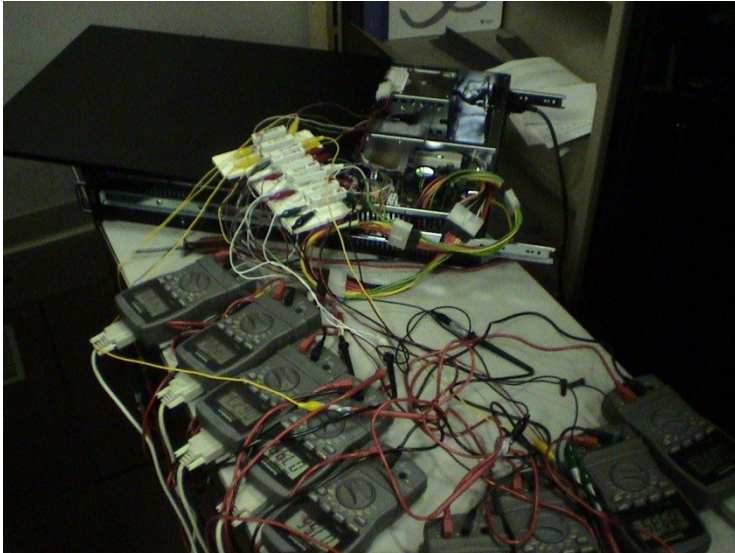  - Support for automated power-aware memory

# PowerPack

**Scalable, synchronized, and accurate.**



Hardware power/energy profiling

Power/Energy Profiling Data

Baytech Management unit

Baytech Powerstrip

AC Power from outlet

Multi-meter   Multi-meter   Multi-meter

DC Power from power supply

Multi-meter control

Multi-meter

Baytech Power Strip

AC

DC

Single node

Data Log

Data Analysis

Data Repository

**Data collection**

MM Thread   MM Thread   MM Thread

Multi-meter Control Thread

Microbenchmarks   Applications

PowerPack libraries (profile/control)

DVS Thread   DVS Thread   DVS Thread

DVS Control Thread

DVS control

**Software power/energy control**

## High-performance Power-aware Cluster

# PowerPack

# DC Power Profiling



```
If node .eq. root then
        call pmeter_init (xmhost,xmport)
        call pmeter_log (pmlog,NEW_LOG)
endif

<CODE SEGMENT>

If node .eq. root then
        call pmeter_start_session(pm_label)
endif

<CODE SEGMENT>

If node .eq. root then
        call pmeter_pause()
        call pmeter_log(pmlog,CLOSE_LOG)
        call pmeter_finalize()
endif
```

**Multi-meters + 32-node Beowulf**

# Power Profiles – Single Node



(a) Power distribution for **system idle**: system power 152.5 Watts

(b) Power distribution for **164.gzip**: system power 206.5 Watts

(c) Power distribution for **171.swim**: system power 209.2 Watts

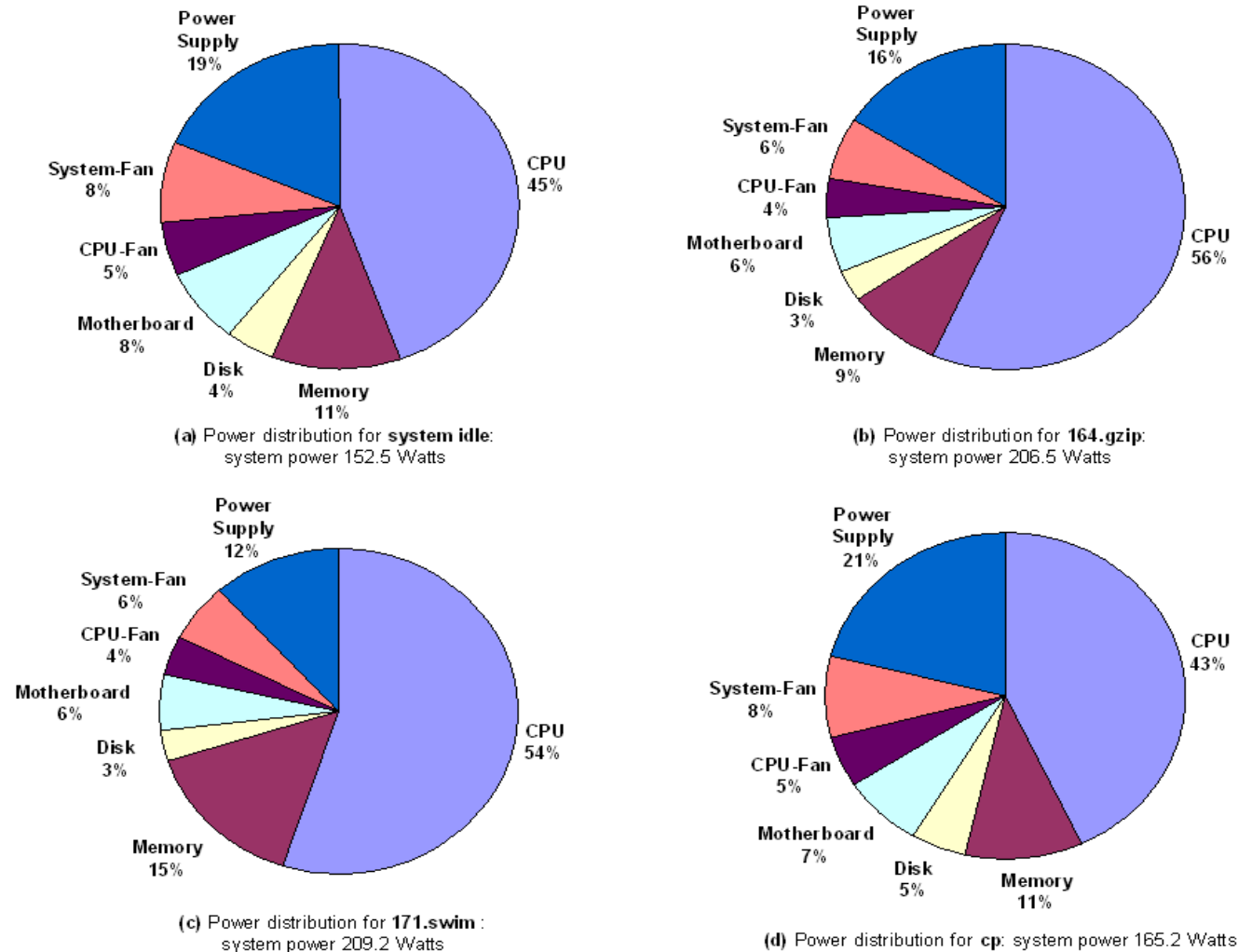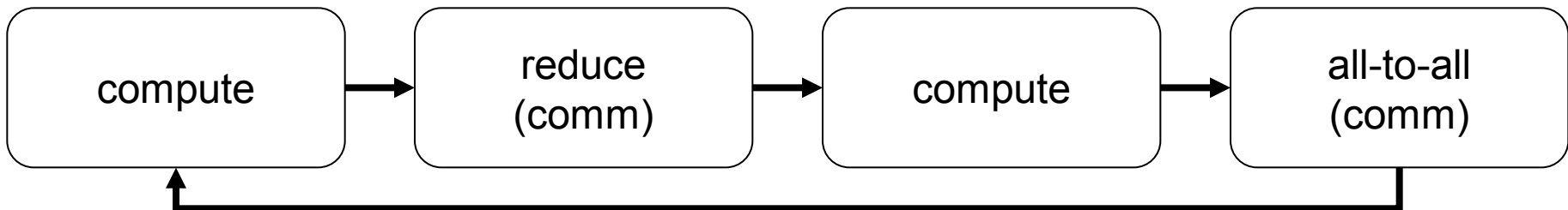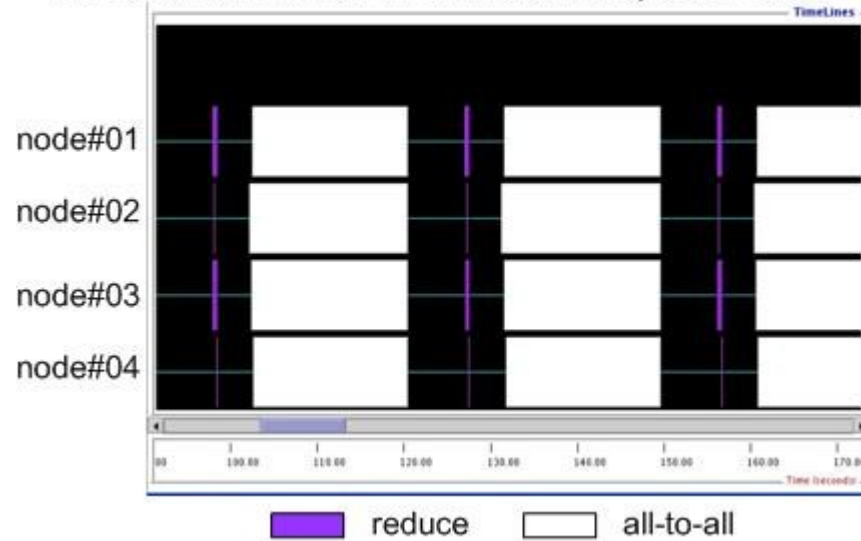(d) Power distribution for **cp**: system power 165.2 Watts

Fig. 5. Power distribution for a single node under different workloads: (a) zero workload (system is in idle state); (b) CPU bounded workload; (c) memory bounded workload; (d) disk bounded workload.

# NAS PB FT – Performance Profiling

Part of the timeline of FT.B.4 visualized by JUMPSHOT

node#01
node#02
node#03
node#04

130.00  110.00  120.00  130.00  140.00  150.00  160.00  170.0

reduce    all-to-all

compute → reduce (comm) → compute → all-to-all (comm)

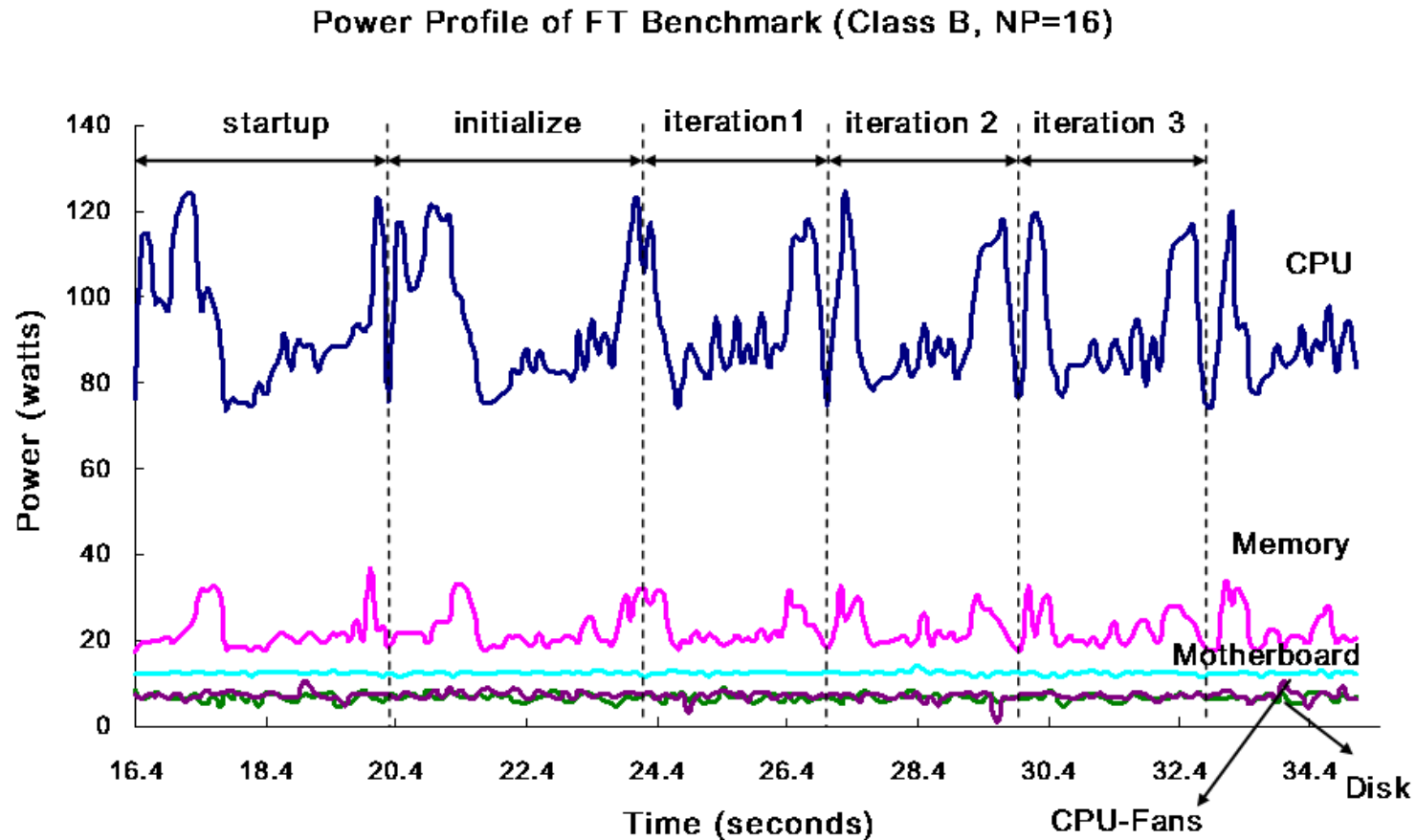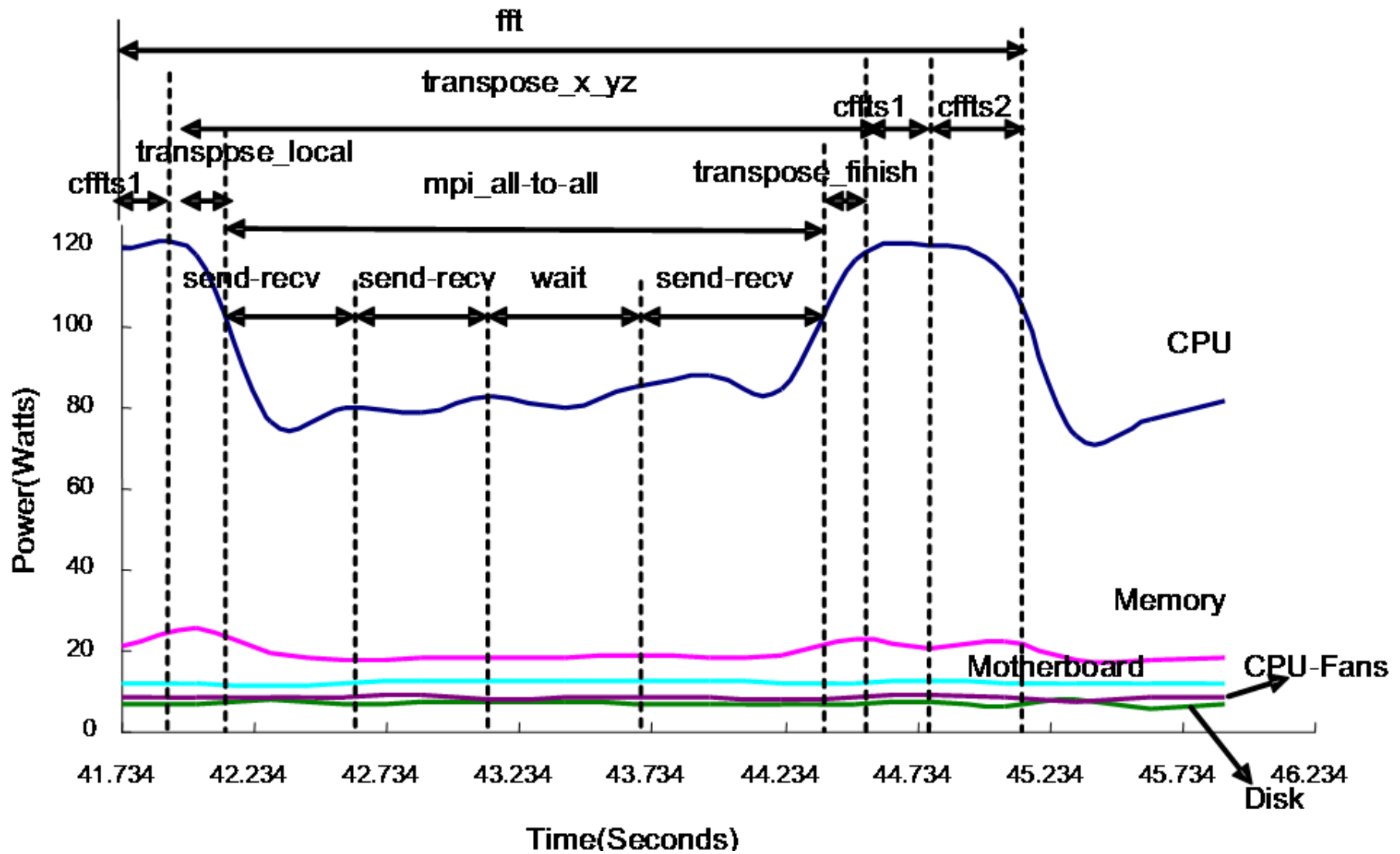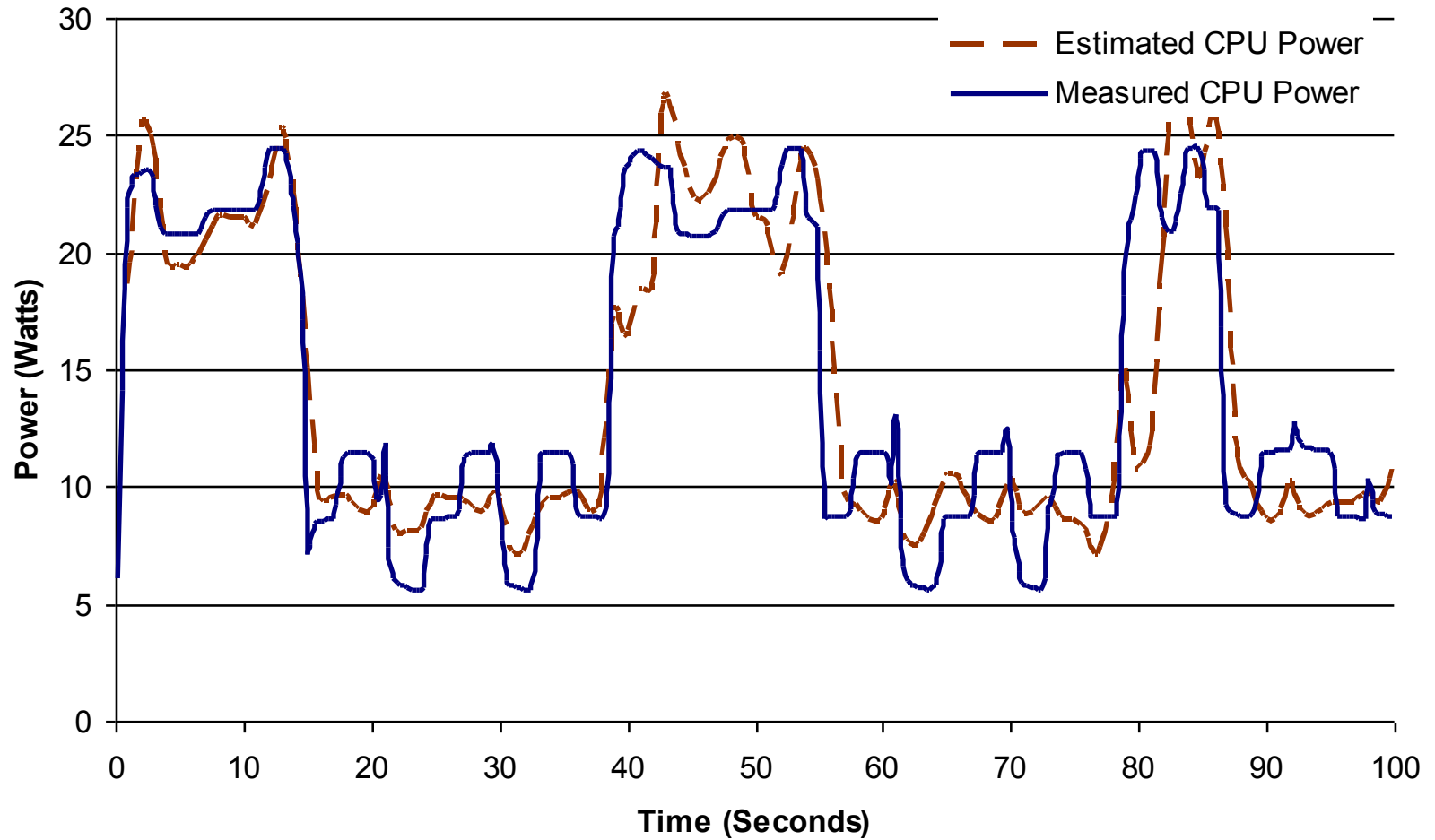# Power Profiles – Single Node



Power Profile of FT Benchmark (Class B, NP=16)

Fig. 6. shows the power use on one node of four for the FT benchmark, class B workload. Note: x-axis is overlaid for ease of presentation.

# PowerPack
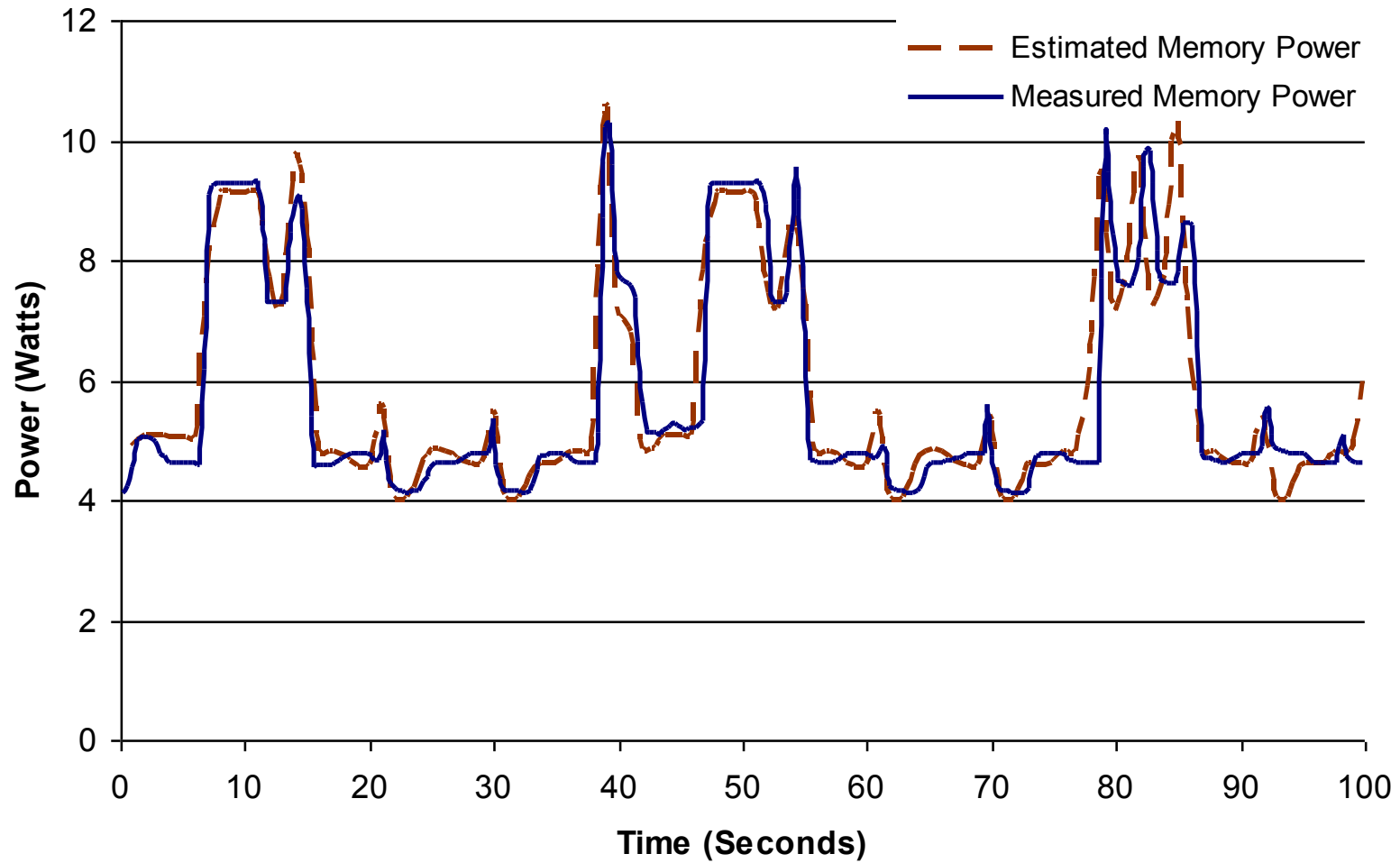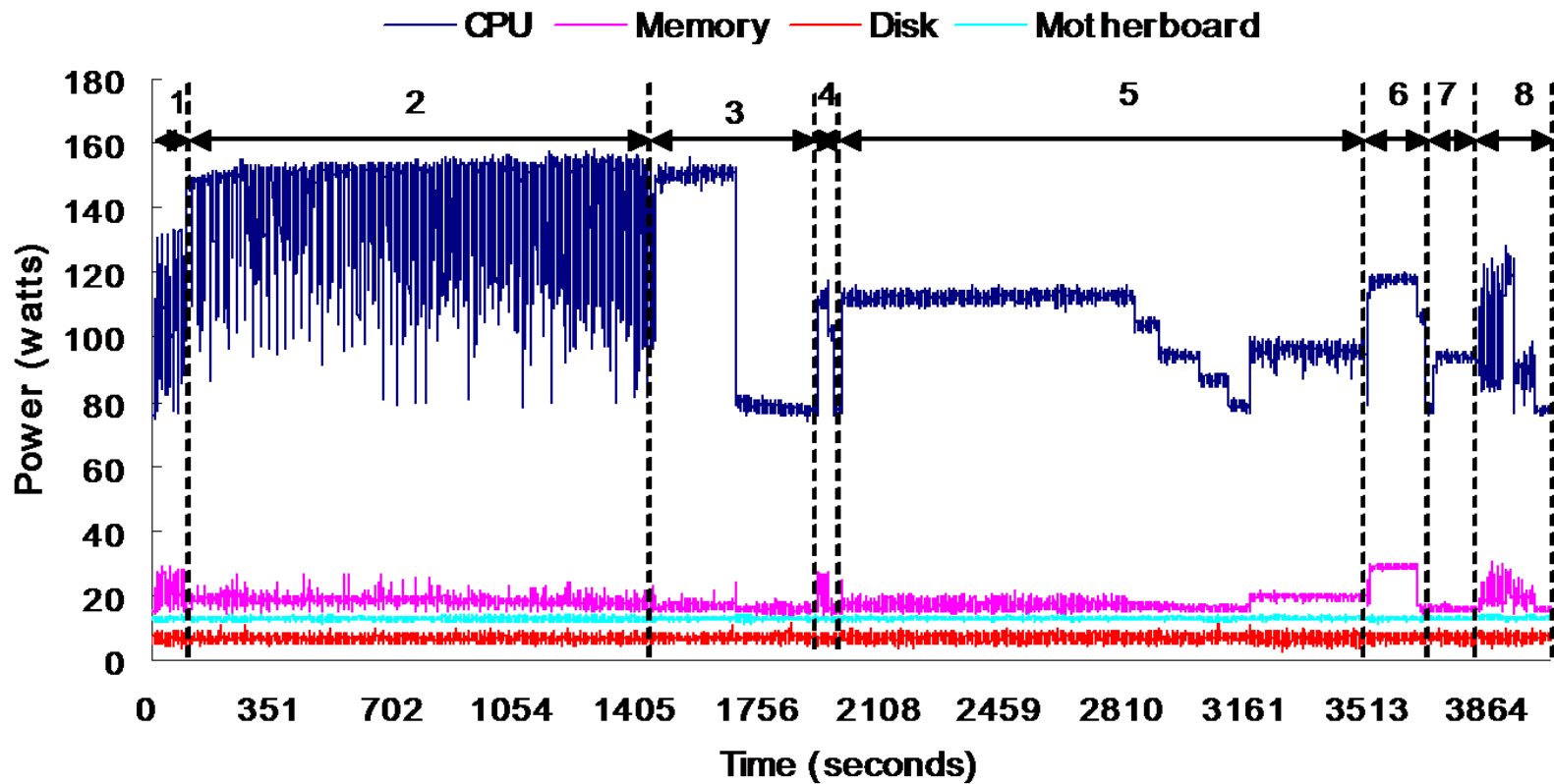
# Predicting CPU Power

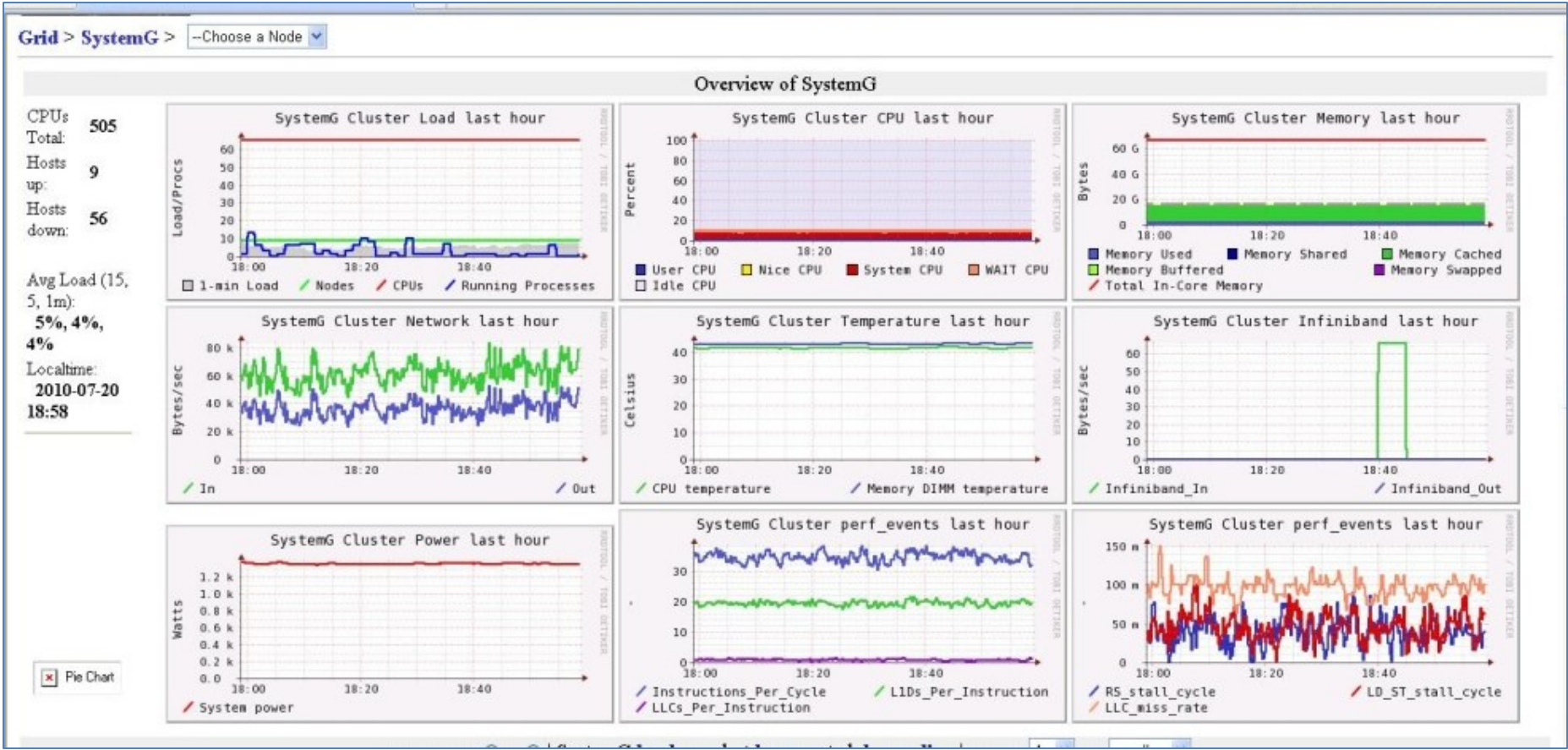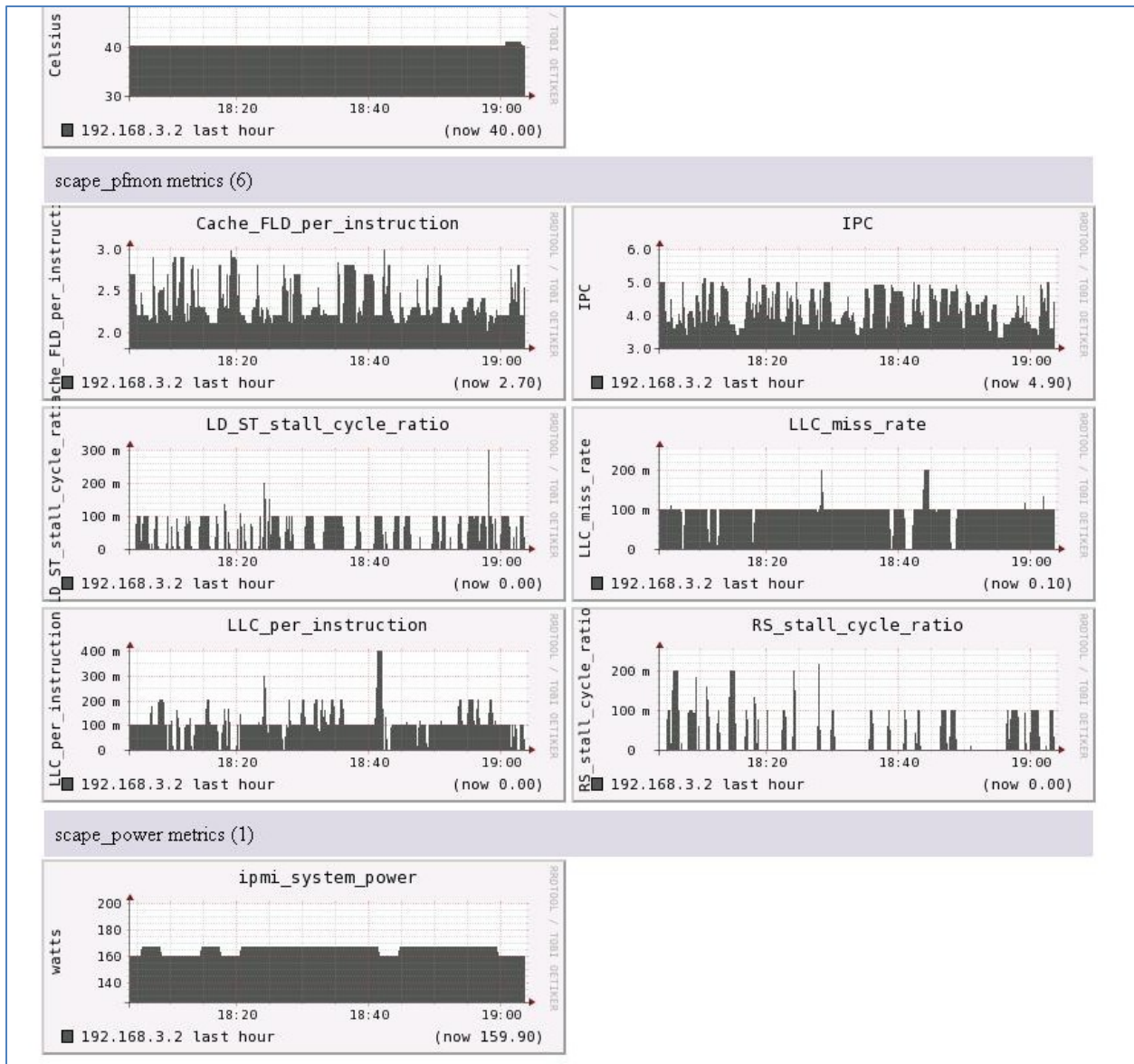# Predicting Memory Power

# SystemG Supercomputer

# PowerPack



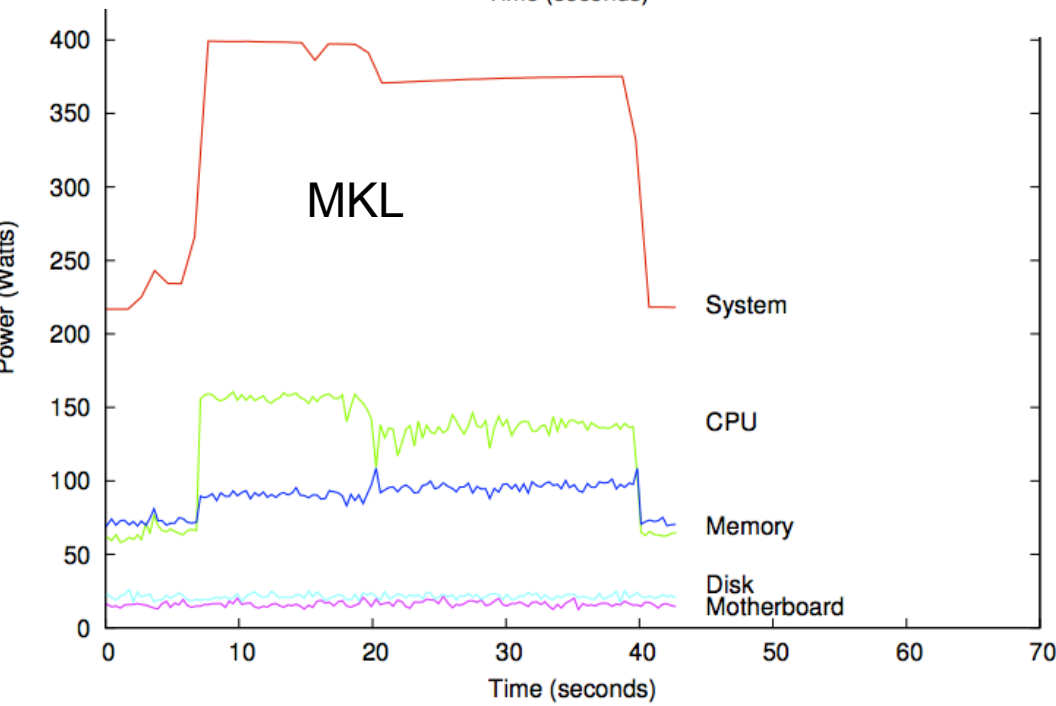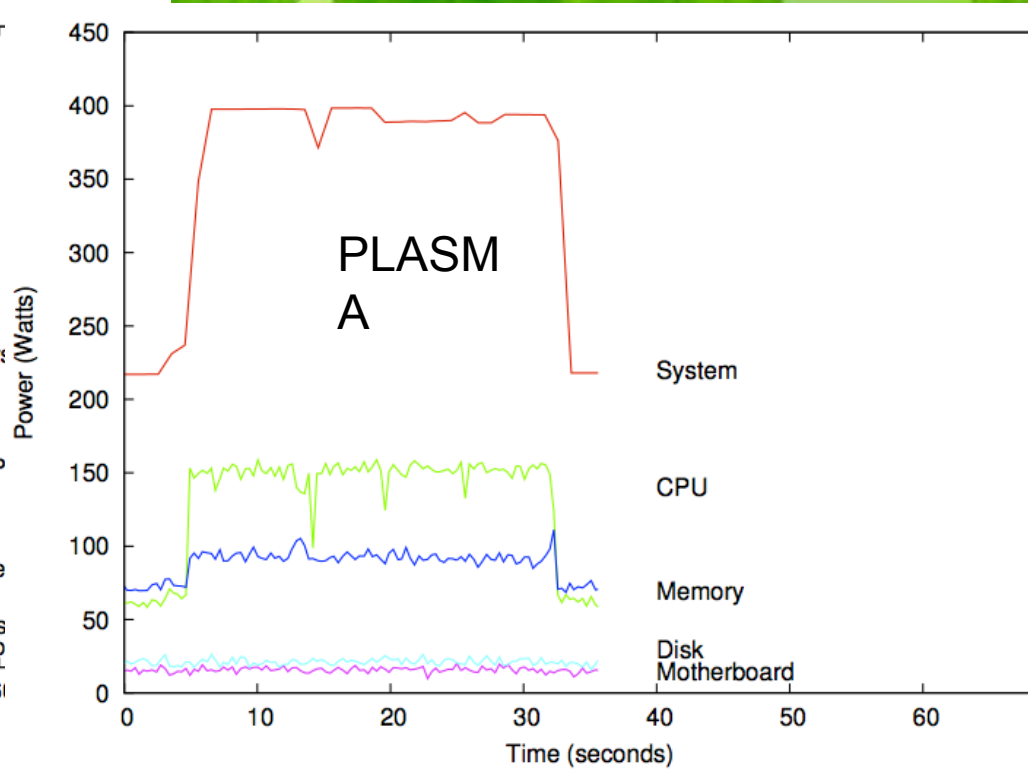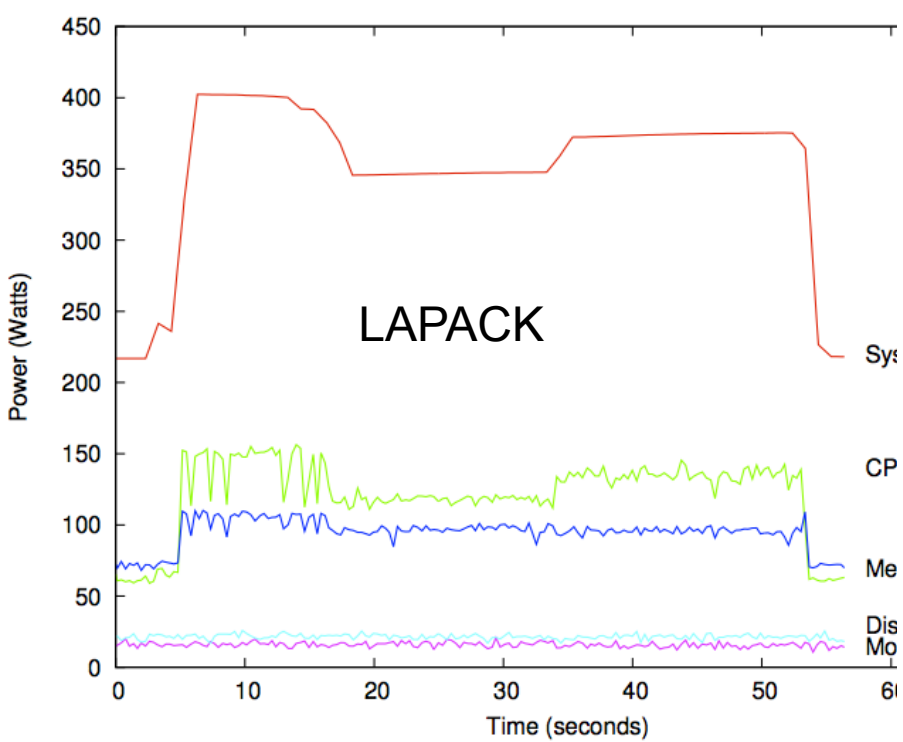Power Profile for HPCC benchmarks running on 8 cores of 2 nodes

# PowerPack 3.0

# PowerPack 3.0

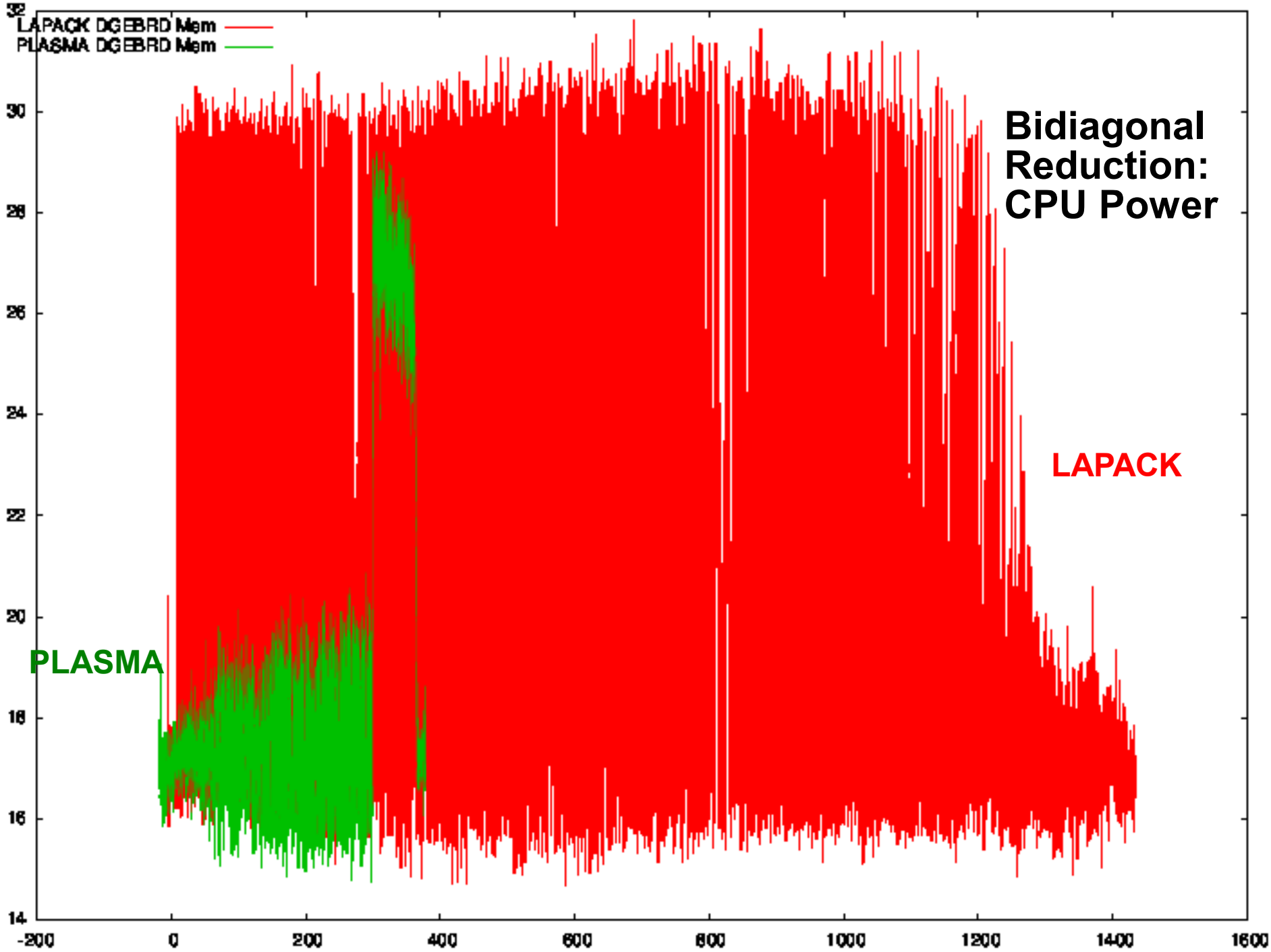# Who uses PowerPack? SystemG?

- Texas A&M (Taylor et al)
- UTenn-Knoxville (Moore, Dongarra, et al)
- Oxford University
- Lawrence Livermore National Lab
- Pacific Northwest National Lab
- Oak Ridge National Lab
- University of Florida
- KAUST (Saudi Arabia)
- University of Madrid (Spain)
  ...and many others

**Power consumption over time**

**Matrix inverse**

Sources:
Piotr Luszczek    Hatem Ltaief

Bidiagonal Reduction: CPU Power

# "To know is to understand."

## *Aristotle*

# Power-Performance Efficiency



Model & Optimize Performance

Model Effects of Power

Improve Power-Performance Efficiency

Profile & Evaluate Power

Optimize Effects of Power
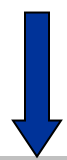
# First power-aware "HPC" cluster
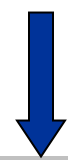
# How DVFS affects HPC efficiency



Communication bound

Memory bound

CPU bound

Lower *f* for energy savings with minimal perf. loss

Higher *f* for better perf. and less energy

# Understanding power-performance

Early system level approaches focus on power mode
**predictor and controller** design:  This is great for *reacting* to change.



**Focus of previous work**

user → policy → Measured error → controller → System input → system → System output

Measured output

Prediction data

sensor → Measured output → predictor

What's missing?

➡ What are the bounds on efficiency? In HPC?
How does power-performance quantitatively affect efficiency?
How do we create policies to guarantee power-performance?

Strong need to improve <u>understanding</u> of power-performance.

# Amdahl's Law

- **Classical speedup**
  - Amdahl's law for 1 enhancement (parallelism)

$$S_N(w) = \frac{T_1(w)}{T_N(w)} = \left[ (1 - FE) + \frac{FE}{SE} \right]^{-1}$$

Energy

Time

Degree of Parallelism

## Time ~ energy. Right?

So we only get energy savings by reducing time. Right?

Then why does PM (e.g. DVFS) save energy? And sometimes without affecting time?

## Amdahl = no overhead

But, overhead is the key to savings energy without loss!

# Power-Aware Speedup

- **Definition**
  - **Speedup**

$$S_N(w, f) = \frac{T_1(w, f_0)}{T_N(w, f) + O(w, f)}$$

  - *w:* workload
  - *N:* number of nodes
  - *f:* the clock frequency and $f_0$ is the base value
  - $T_1(w, f_0)$: sequential execution time at base frequency $f_0$
  - $T_N(w, f)$: parallel execution time at $N$ processors at frequency $f$

# Bounding Efficiency at Scale



EDP values for LU

- # **Optimal system configuration**
  - – # processors: 256
  - – CPU frequency: 1200MHz

# Understanding power-performance

Early system level approaches focus on power mode
**predictor and controller** design:  This is great for *reacting* to change.



What's missing?
  What are the bounds on efficiency? In HPC?
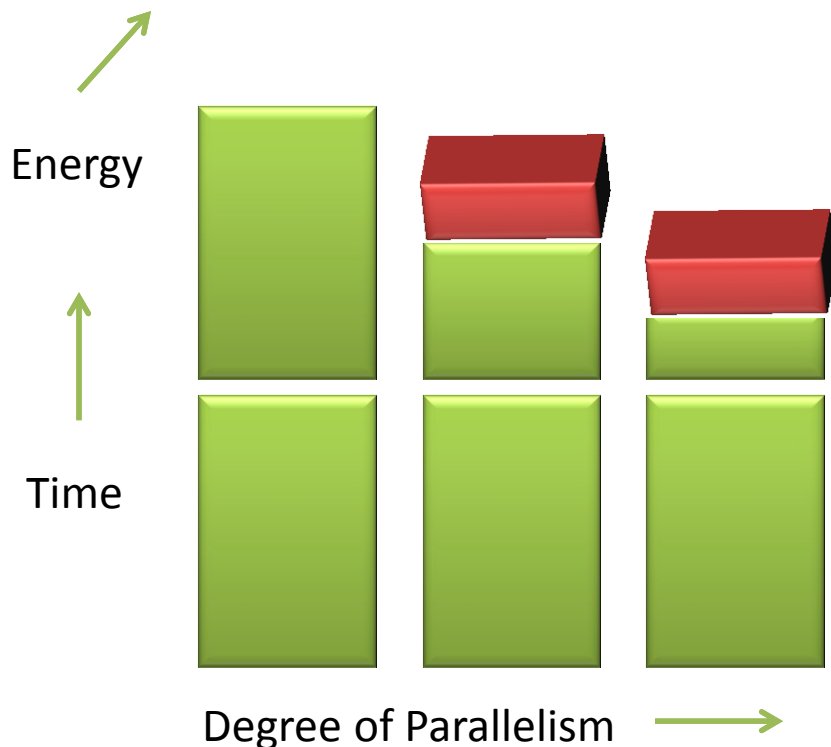  How does power-performance quantitatively affect efficiency?
  How do we create policies to guarantee power-performance?

Strong need to improve <u>understanding</u> of power-performance.

# Iso-energy-efficiency

Grama et al: performance efficiency can be held constant if we increase both number of processors and problem size simultaneously.

Algorithm + Scale → fixed performance

**Iso-energy-efficiency**

Algorithm + Scale + Power Modes → (power, performance)
- Requires accurate performance model
- Requires accurate power model
- Must be accurate, useful, usable

# Iso-energy-efficiency Derivation

General form of our Iso-energy-efficiency model:

$$EE = \frac{E_1}{E_p} = \frac{E_1}{E_1 + E_o} = \frac{1}{1 + E_o / E_1}$$

$EE$ : *system-wide energy efficiency*

$E_1$ *(baseline):* total energy consumption of sequential execution on one processor
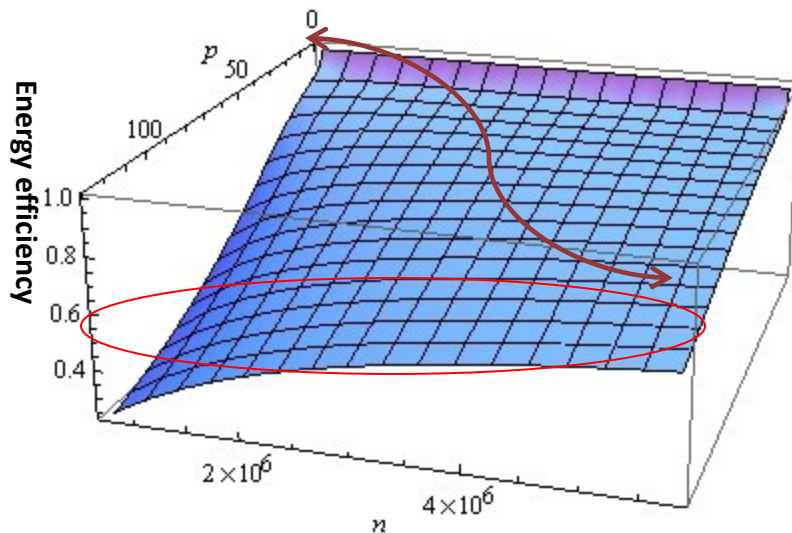
$E_p$ : the total energy consumption of parallel execution for a given application on *p* parallel processors

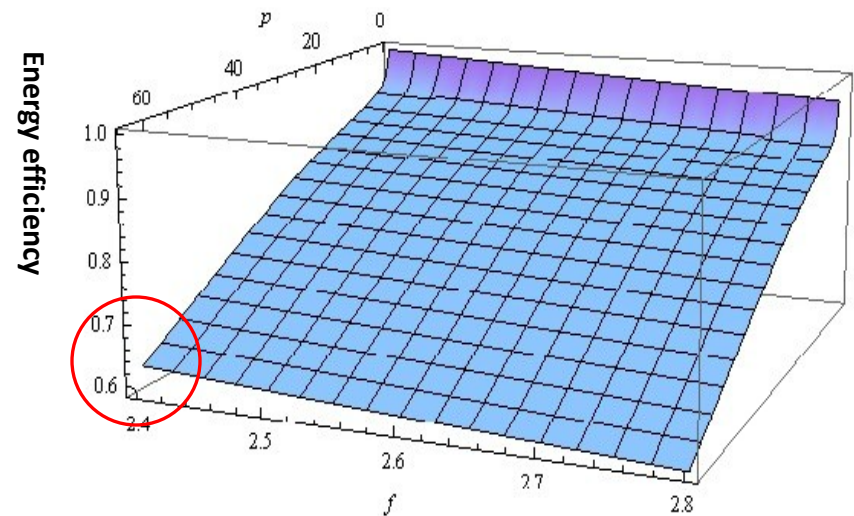$E_o$ : the additional energy overhead required for parallel execution and running extra system components

43

# Maintaining Efficiency in 3-D FFT

$$EE_{FFT} = \cfrac{1}{1 + \cfrac{6.87 \log_2 p - 1.75 f \log_2 p + p(p-1) f \left(\dfrac{11500}{n} + \dfrac{0.376}{4^{\log_2 p - 2}}\right)}{163 + 22.7f}}$$

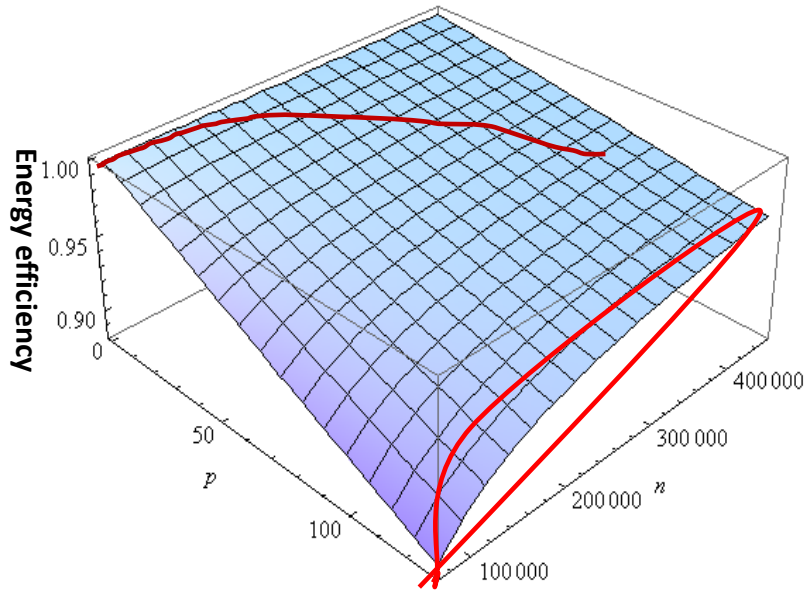FT's system-wide energy efficiency with p and n as variables

FT's system-wide energy efficiency with p and f as variables
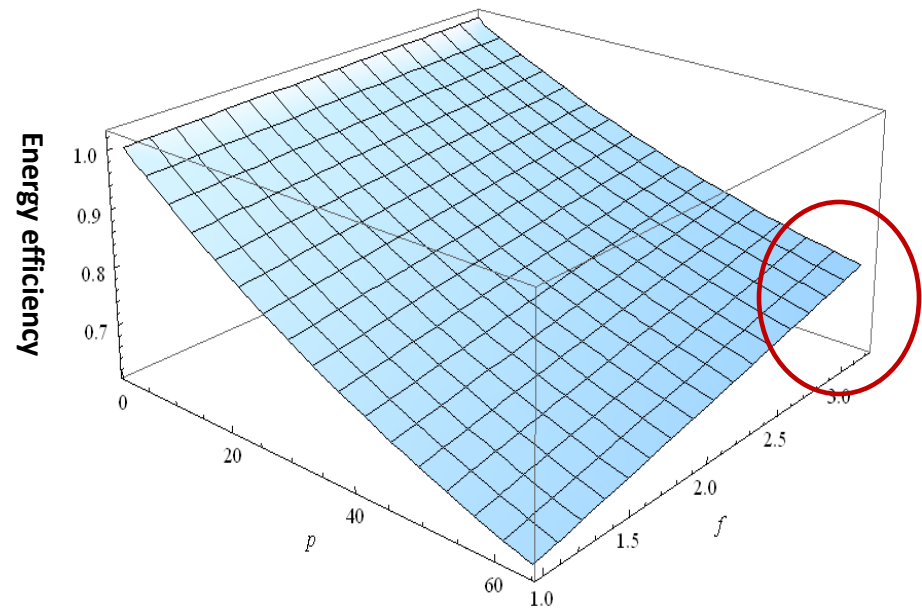


- ➢ *Problem size scaling effective in maintaining overall system energy*
- ➢ *CPU frequency scaling: only slightly improves EE*
- ➢ *But, the effects of CPU clock frequency on on-chip workload diminish while scaling up system size.*

44

# Maintaining Efficiency in CG

CG's system-wide energy efficiency with p and n as variables       CG's system-wide energy efficiency with p and f as variables



- ➤ *Overall EE decreases with system size*
- ➤ *EE can be maintained or improved by scaling up problem size N.*
- ➤ *Applying higher frequency will improve system-wide EE while system size scales up.*
- ➤ *In contrast to FT, effects of frequency on on-chip workload diminish at a slower rate.*

45

"Those that can, do.
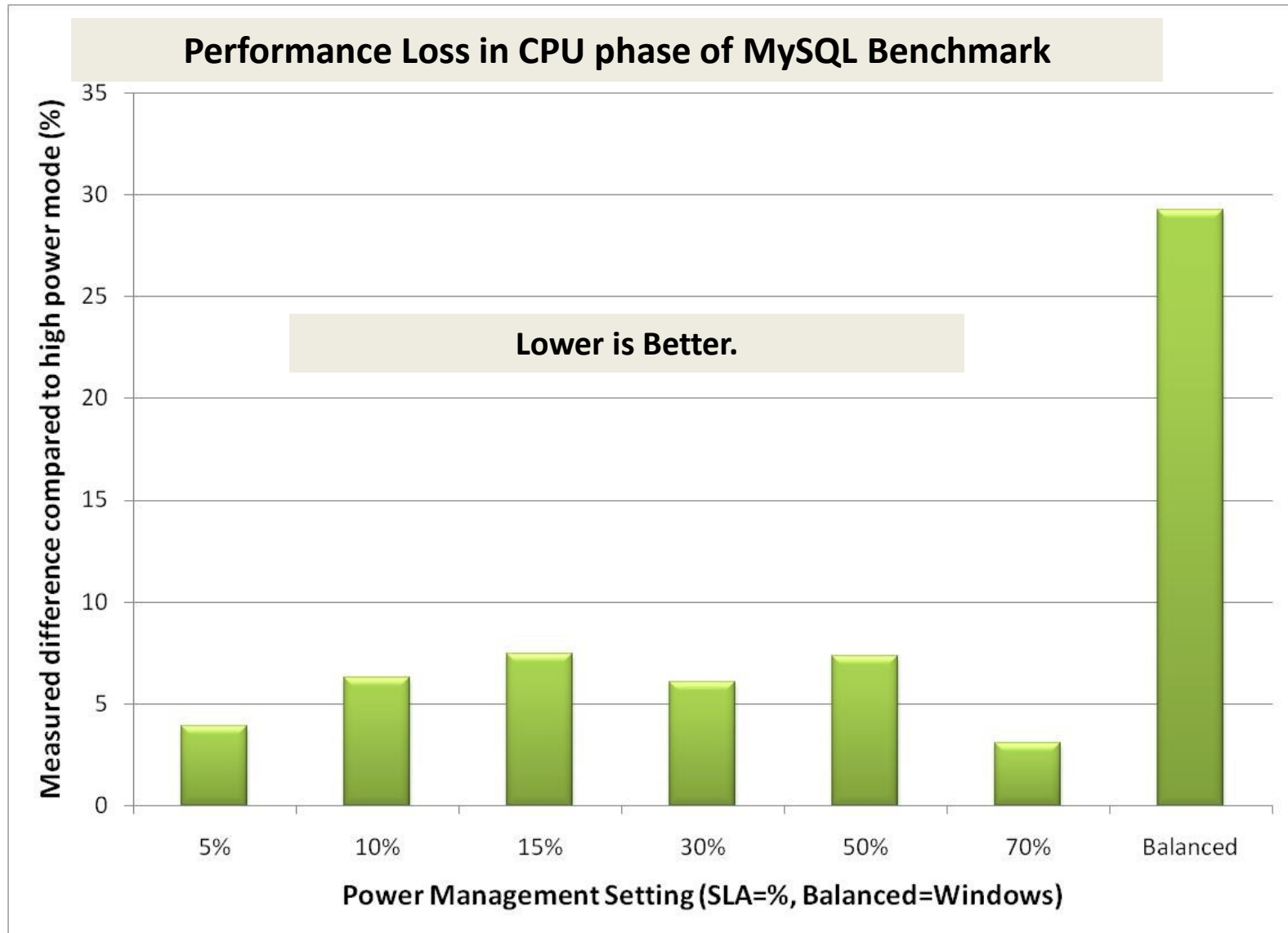Those that can't, complain."

*Linus Torvalds*

# State of the art PM

Amount and cost of power continues to increase.



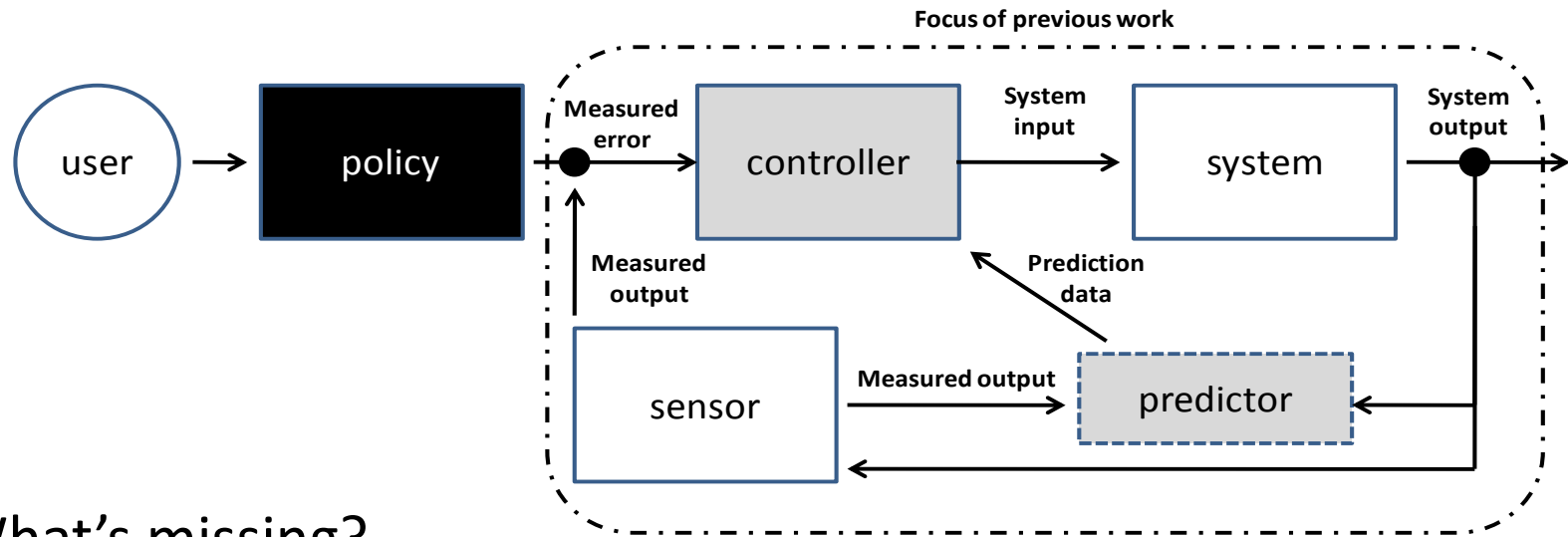Power management features disabled by default.

# Why is PM turned off?



**Performance Loss in CPU phase of MySQL Benchmark**

Lower is Better.

Measured difference compared to high power mode (%)

Power Management Setting (SLA=%, Balanced=Windows)

# Understanding power-performance

Early system level approaches focus on power mode
**predictor and controller** design:  This is great for *reacting* to change.



What's missing?

What are the bounds on efficiency? In HPC?
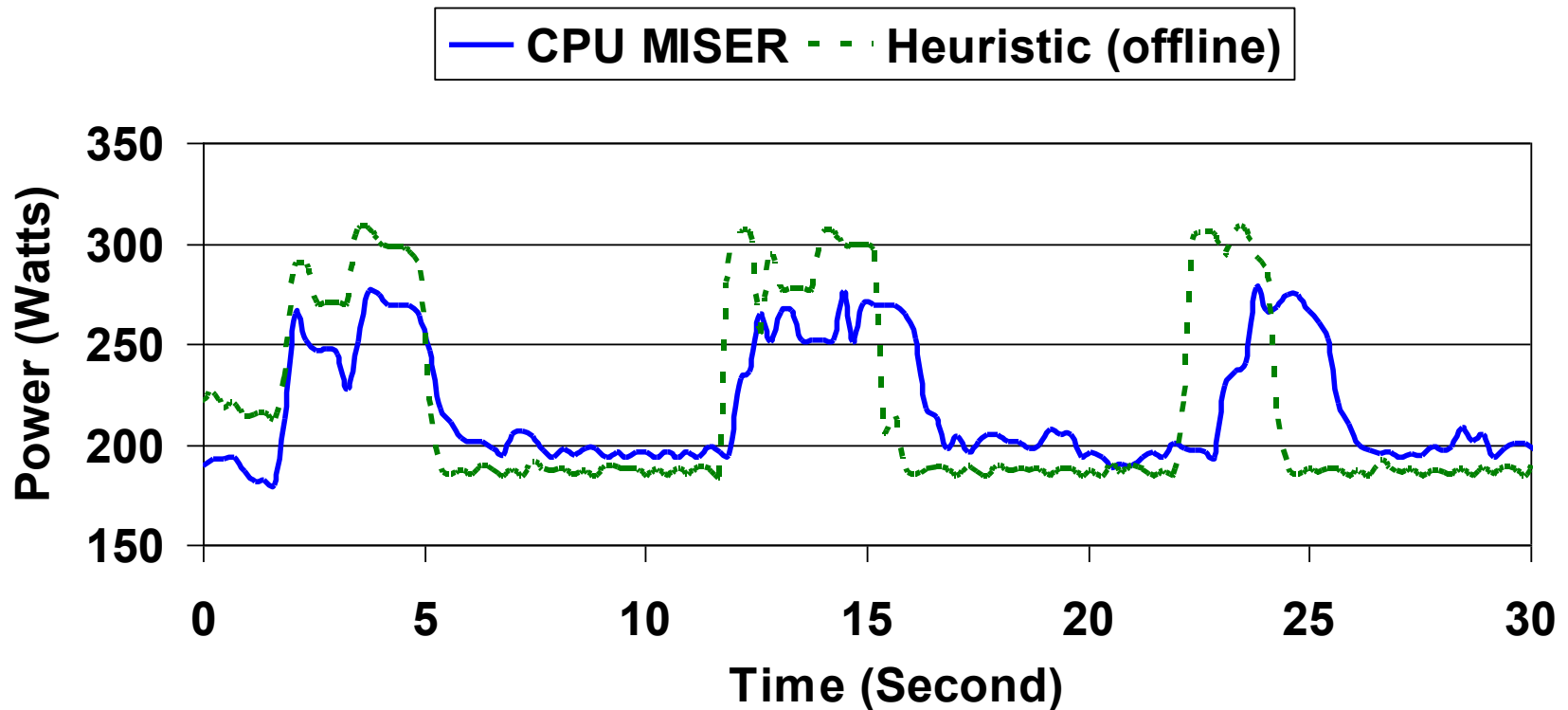How does power-performance quantitatively affect efficiency?
How do we create policies to guarantee power-performance?

Strong need to improve <u>understanding</u> of power-performance.
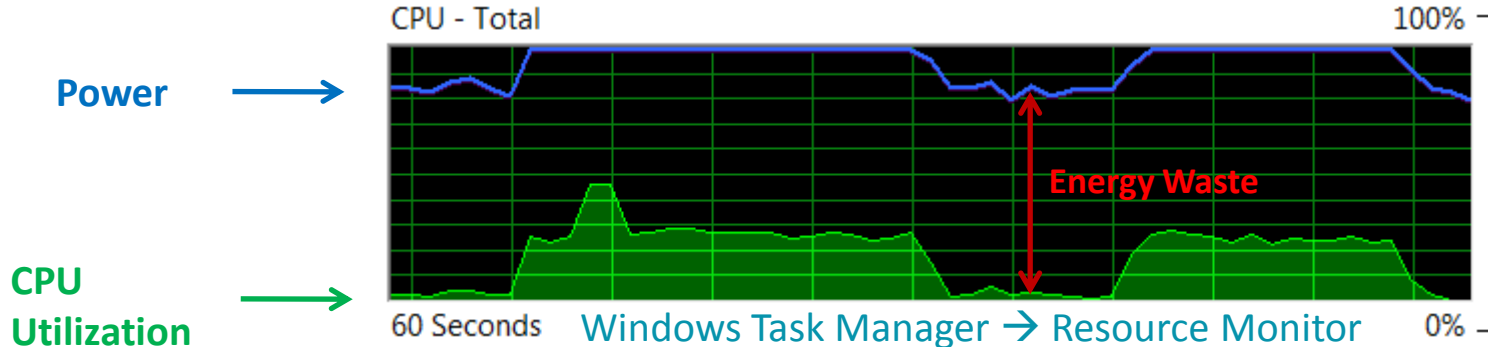
# Model-directed Scheduling



**System Power Traces for FT**

Legend: —— CPU MISER  - - - Heuristic (offline)

- Automatically and transparently schedule CPU frequency to reduce power

# Better SW for the masses...

## Before Granola:

Power →

CPU Utilization →

Energy Waste

CPU - Total    100%
60 Seconds    Windows Task Manager → Resource Monitor    0%

## After Granola:

Power →

CPU Utilization →

Reduced Energy Waste with no noticeable performance loss

CPU - Total    100%
60 Seconds    Windows Task Manager → Resource Monitor    0%
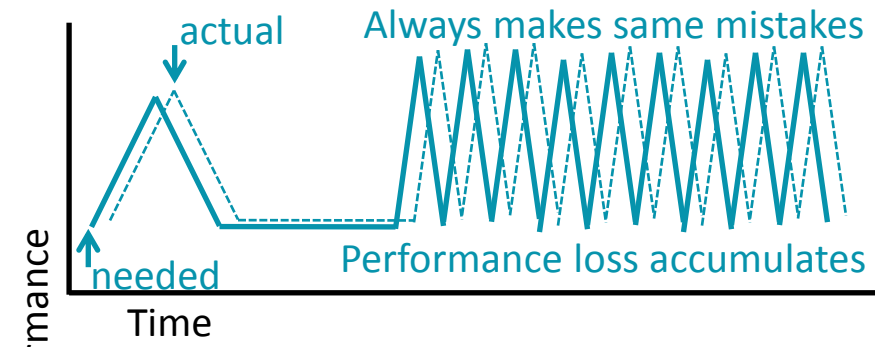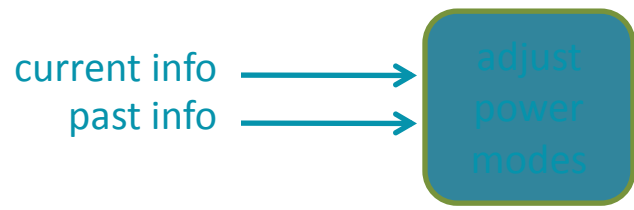
# Reducing IT Costs *Without* Performance Loss

EVERYONE ELSE[1]
Intel, HP,
Windows,
VMWare, …

current info
past info

adjust power modes

actual
Always makes same mistakes

Performance loss accumulates

needed

Performance

Time

GRANOLA
POWER TUNING

Performance Guarantee
Technology
current info
past info
user-defined SLA
current level of service

adjust power modes

algorithm metric feedback

actual
Save power always within SLA

needed
Starts conservative, then adapts

Performance

Time

[1]Note: Verdiem, 1E, and others *only* turn systems off when not in use. We offer that too as needed.

# Commercial grade measurement...



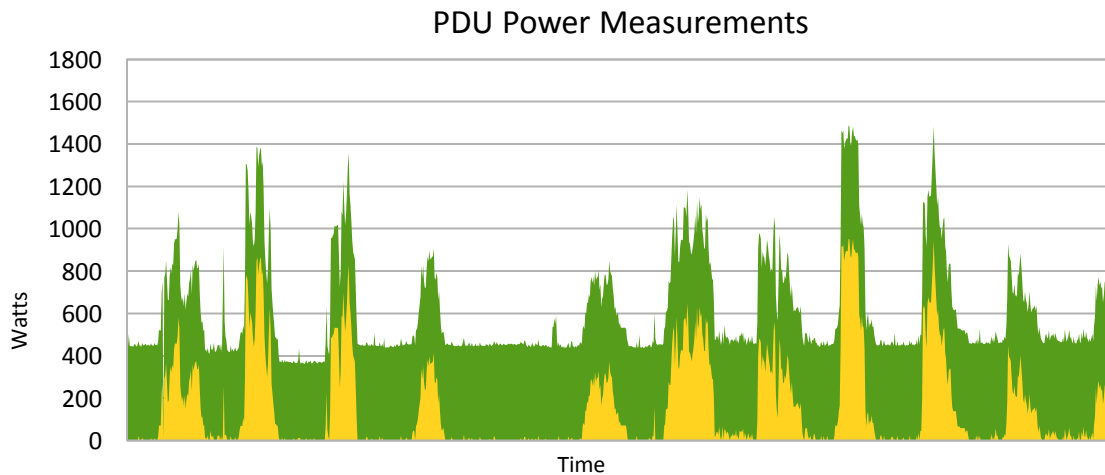Granola Enterprise Power Estimates

Granola software gives more detail...

■ CPU
■ System
■ Monitor

PDU Power Measurements

...same accuracy as expensive hardware

■ System +CPU
■ Monitor

# Granola (http://grano.la)

- **Launched Earth Day 2010**
- **Free home version**
- **300K+ Downloads so far...**
- **160+ Countries**
- **Uses: laptops, PCs, servers**
- *Performance Guarantees*



granola

You'll save 236.3 kWh yearly

Enough to power 31 electric furnaces, an air conditioner, and 6 refrigerators for an hour

You'll save 28.36 USD yearly

Enough for a monkey wrench to throw in the gears, 3 shirts from the thrift store, and a political bumper sticker

You'll save 321.4 lbs CO2 yearly

As much as a 500-mile flight, a tree, and 18 miles in a compact car

You've saved 45.4% CPU energy

...and you didn't even notice!

*224,404 trees*

You and the Granola community will offset 224,404 trees worth of CO2 this year.

# The hard truth about the future

### Measurement



DEFINITIONS
- Experts needed
- Easy to get a wrong answer/conclusion
- Scalability questionable

### Analysis



CHAOTIC
- Power-performance relationship not well understood
- How can we help?
- Who are we helping?

### Optimization



CONTENTIOUS
- Many point solutions
- Reactive
- Making something no one wants

# Where do we go from here?



We need lots of help.

Disruptive vs. Incremental.

Silver bullet is unlikely.

Commodity matters.

Markets matter.

Tools matter.

Wanted: Major catastrophe.

Custom system is likely the only answer by 2019. Energy wall?

"Victory" is inevitable when you change the game.

# Thank you.

# Fine-Grain Parameterization

- ## Assumptions
  - Workload perfectly parallelizable: $T_s^{on}=T_s^{off}=0$
- ## Methodology
  - Measure system prior to application execution
    - CPI/f for on-chip workload for all frequencies
    - $t^{off}$ for off-chip workload
    - Empirically estimate $T_{PO}$
  - Profile workload at base frequency
    - Accesses for on-chip workload
    - Accesses for off-chip workload
  - Predict perf of node and frequency combinations