# Statistical Profiling-based Techniques for Effective Power Provisioning in Data Centers

**Sriram Govindan,**

Jeonghwan Choi, Bhuvan Urgaonkar, Anand Sivasubramaniam, Andrea Baldini

Penn State, KAIST, Tata Consultancy Services, Cisco Systems
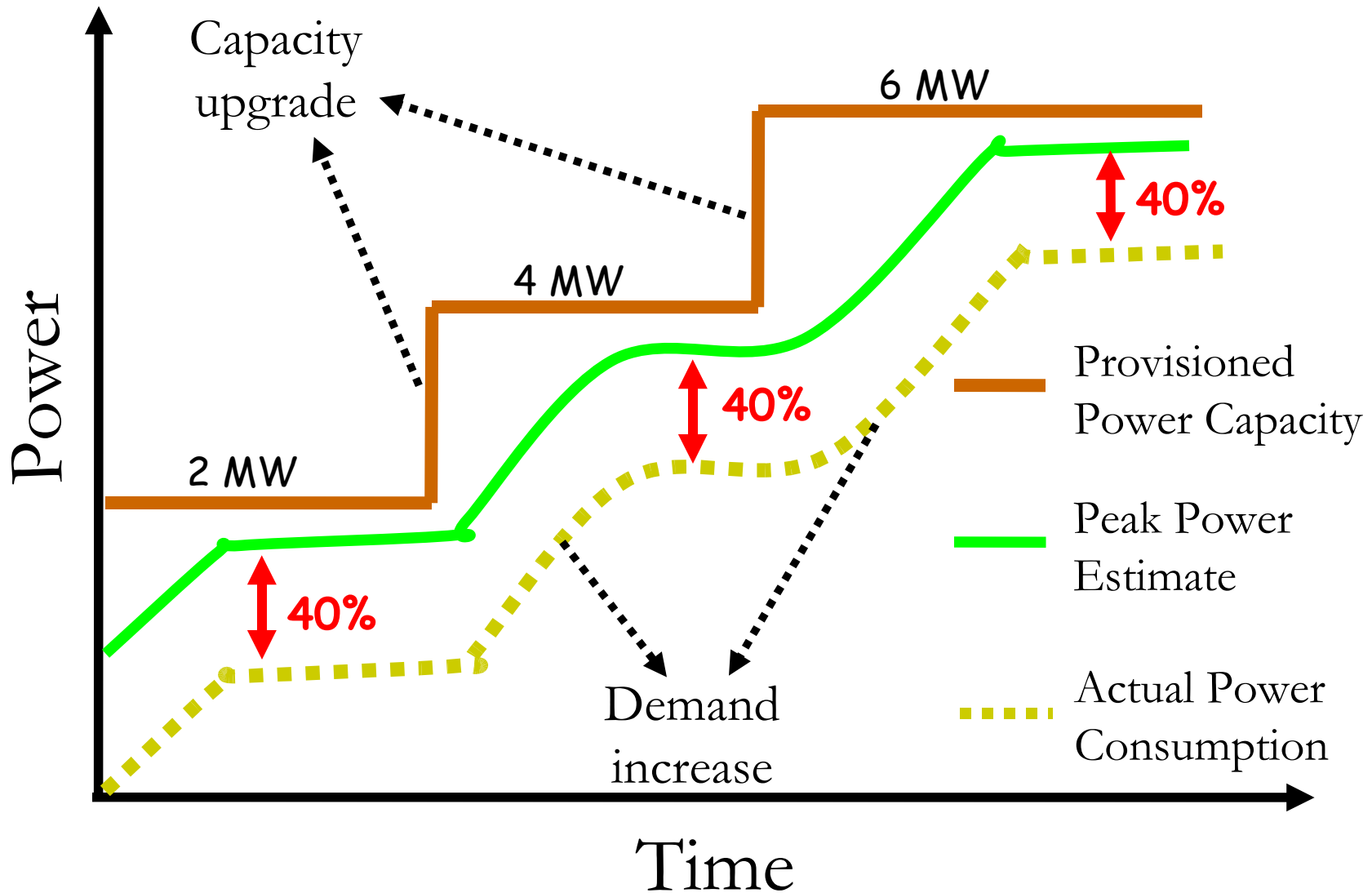
Eurosys 2009 , March 31st – April 3rd 2009

1

# Growing Energy Demands

- In 2006, U.S data centers

  - Spent $4.5 billion just for powering their infrastructure

  - 1.5% of the total electricity consumed in the U.S

  - Has more than doubled since 2000 - further expected to double by 2011

- Massive growth of installed hardware resources

  - By 2010, servers expected to triple from 2000

  - Average utilization of servers between 5% and 15%

# Data Center Energy Management

- Tackle server sprawl
  - *Server virtualization:* Consolidates workload on to fewer number of servers and switch off remaining idle servers
  - Growth in number of data centers – provisioning power infrastructure of a data center
    - *Provisioned power capacity:* Maximum power available to the data center as negotiated with the electricity provider
    - *Provisioning:* How many IT equipments (servers, disk arrays, etc.) can be hosted within a data center ?

# Data Center Power Provisioning



Capacity upgrade

6 MW

4 MW

2 MW

40%

40%

40%

Power

Time

Demand increase

Provisioned Power Capacity

Peak Power Estimate

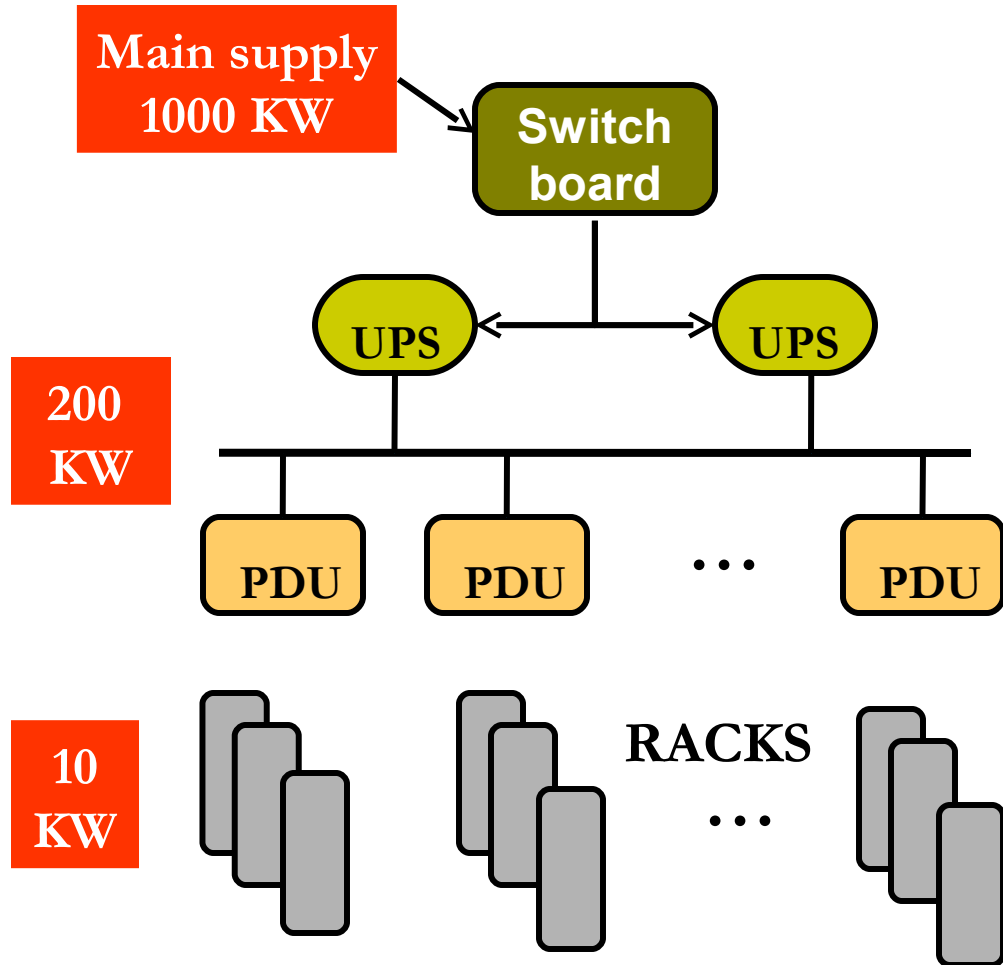Actual Power Consumption

- Hand drawn figure

# Over-provisioned Data Centers

- Current provisioning practices render data centers' power infrastructure highly under-utilized

  - **Reliability concerns**

- Over-provisioning hurts profitability of data centers due to

  - Unnecessary proliferation of data centers

    - Increase in management and installation costs

  - Electrical and cooling inefficiency

    - Efficiency is worse at lower loads

- **Goal: Improve utilization of the power infrastructure in data centers while adhering to reliability constraints**

# Talk Outline

- Data Center Power Hierarchy

  - Hardware reliability constraints

- Application Power Profiles

- Improved Power Provisioning

  - Threshold-based power budget enforcer
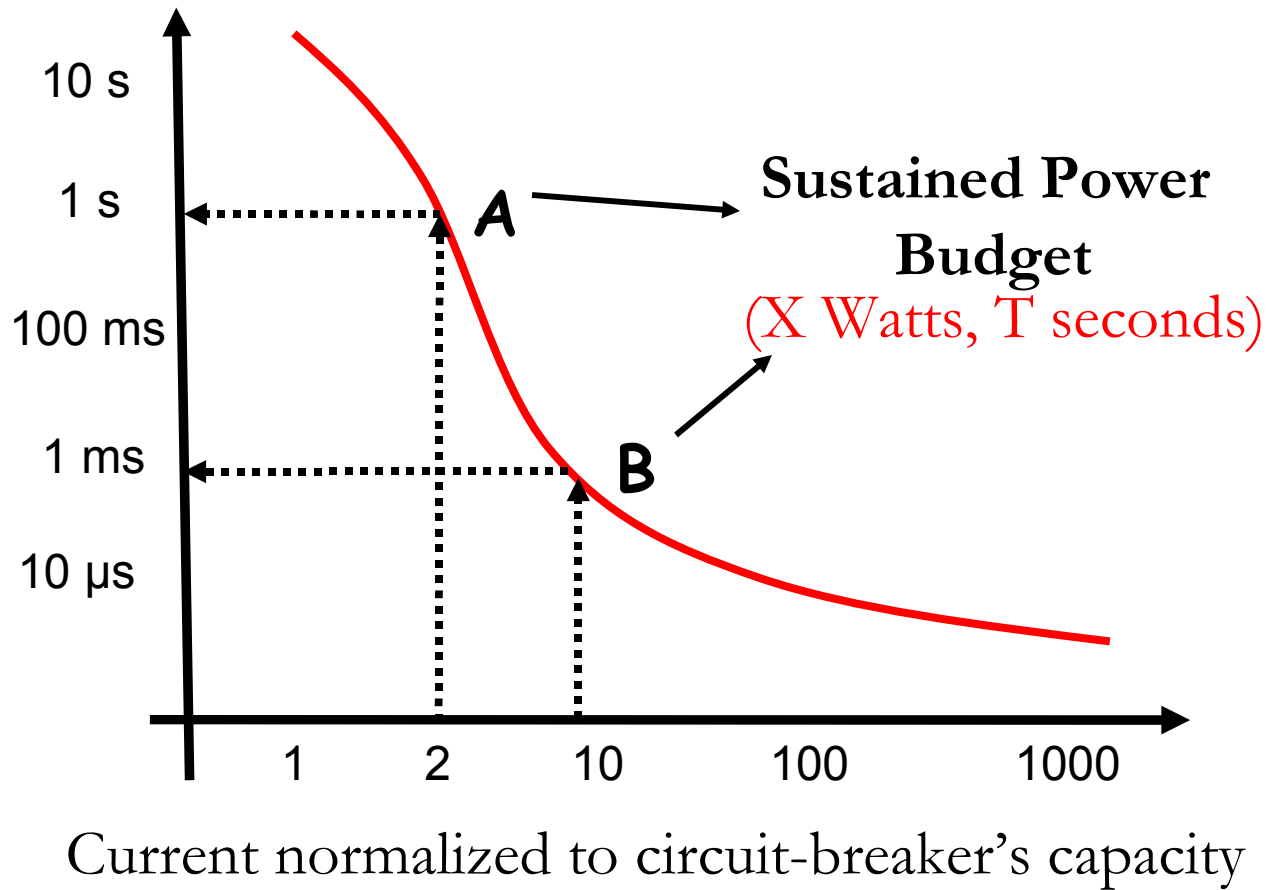
- Evaluation

# Data center Power Supply Hierarchy

**Main supply 1000 KW**

**Switch board**

**UPS** ⟷ **UPS**

**200 KW**

**PDU**  **PDU**  . . .  **PDU**

**10 KW**

**RACKS**
. . .

❑ *Circuit breakers* placed at each element of a data center power hierarchy to protect the underlying circuit from current overdraw or short-circuit situations

# Time-current characteristics Curve of a typical Circuit-breaker



Time for which current should be sustained before tripping the circuit breaker

10 s
1 s
100 ms
1 ms
10 µs

A — **Sustained Power Budget**
(X Watts, T seconds)

B

1    2    10    100    1000

Current normalized to circuit-breaker's capacity

- Hand drawn figure

8

# Profiling Application Power Consumption

**Application**

**Virtual Machine**

**Xen VMM**

Accuracy:
**1 µA**
Granularity:
**1 ms**

**Signametrics Multimeter (SM2040)**

**PDF**

Idle power ~ 160 W
Max power ~ 300 W

Probability

Power (W)

160

300

1

0

# Power Profiles - 2 ms Granularity



TPC-W
(60 sessions)

99th percentile
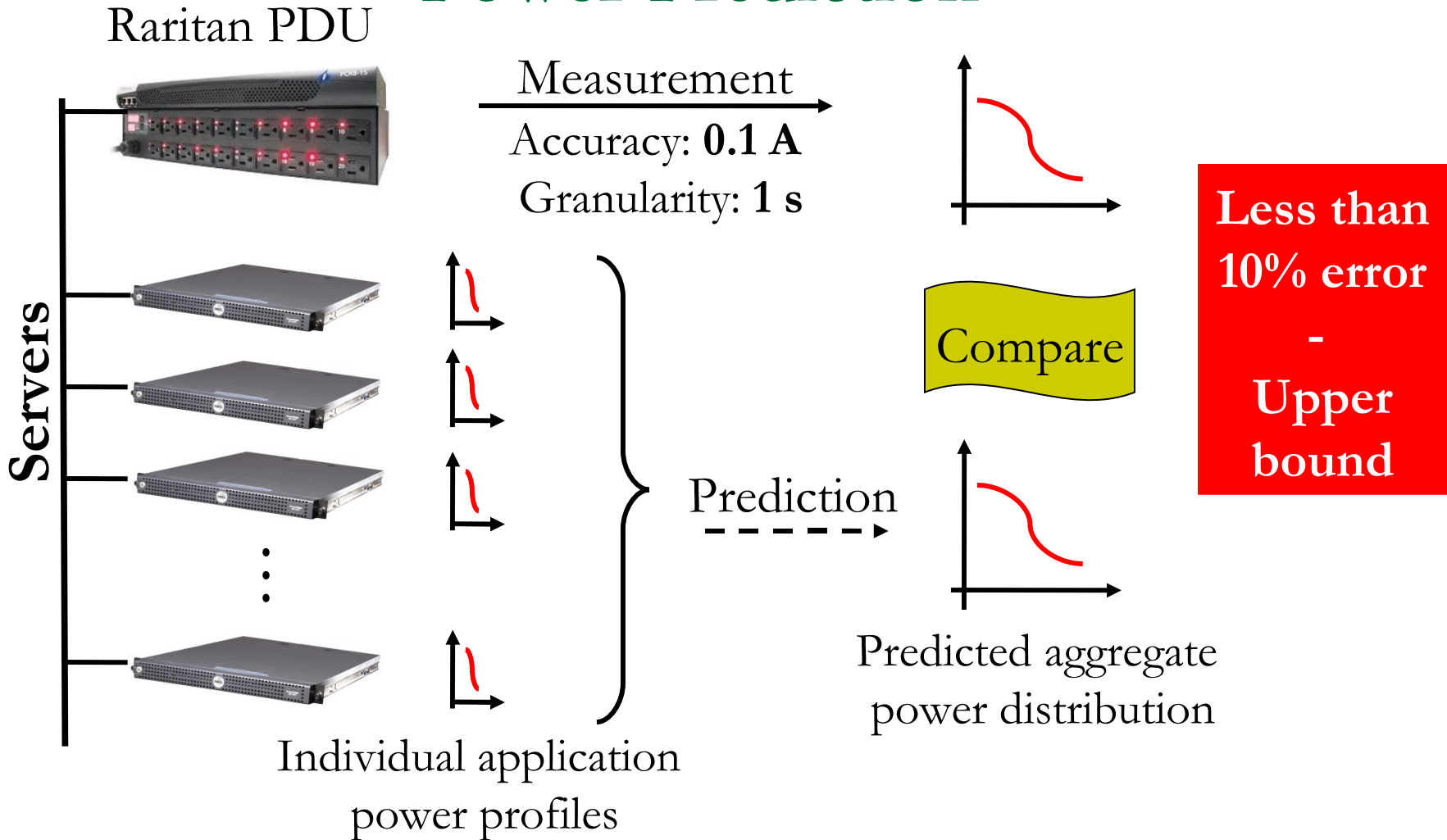
Peak

- **TPC-W**
  - Emulates a two-tiered implementation of an e-commerce book-store with front-end jboss web server and back-end mysql database.

# Statistical Multiplexing Based Sustained Power Prediction

Raritan PDU



**Servers**

Measurement

Accuracy: **0.1 A**

Granularity: **1 s**

Compare

Prediction

Individual application power profiles

Predicted aggregate power distribution

Less than 10% error - Upper bound

**Reference:** Profiling, prediction and capping of power-consumption for Consolidated Data-center environment, Choi et al., MASCOTS 2008
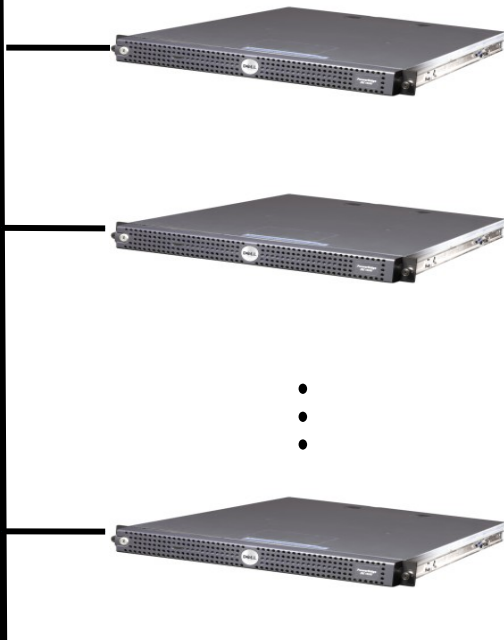
# Existing Power Provisioning Techniques

- Face-plate rating/Name-plate rating
  - Assumes all components are populated in the server
    - Eg: All processor sockets, DIMM slots, HDDs etc.,
  - Assumes all components consume peak power at the same time

- Vendor power calculators
  - Dell, IBM, HP etc.
  - Tuned for current server's configuration and coarse-level application load information.
  - Less conservative than Face-plate Rating
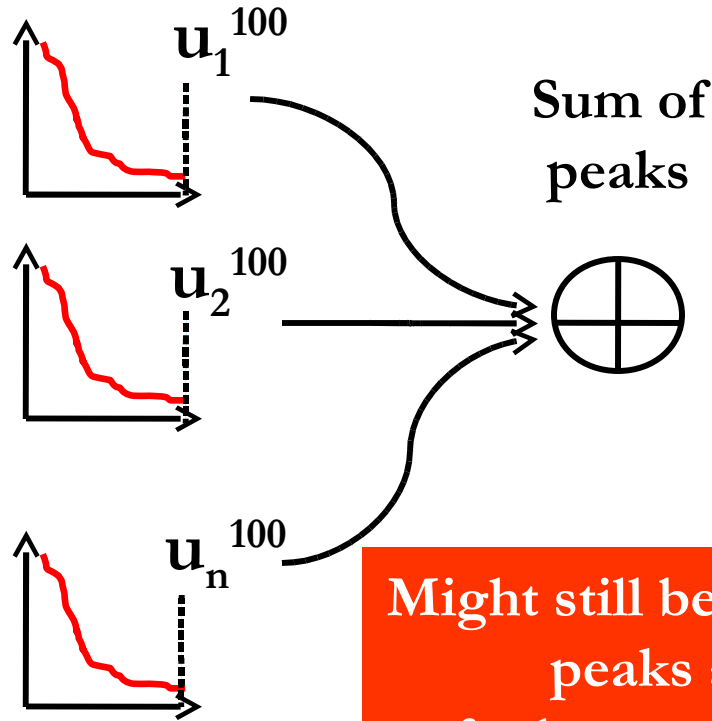
# Provisioning for Peak Power Needs
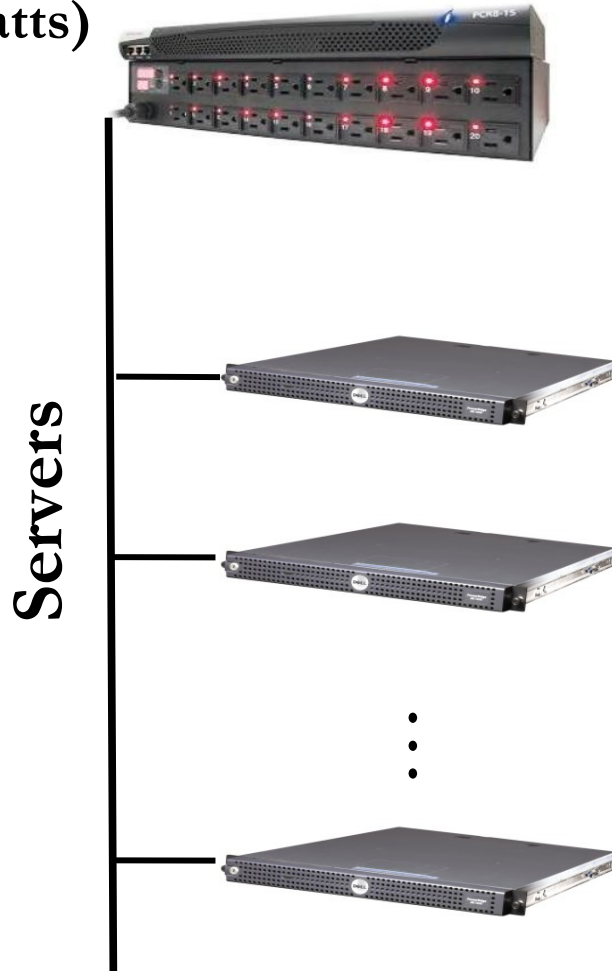
**PDU**
**(B Watts)**

**Servers**

$$\sum_{i=1}^{n} u_i^{100} \leq B$$

$u_1^{100}$

$u_2^{100}$

$u_n^{100}$

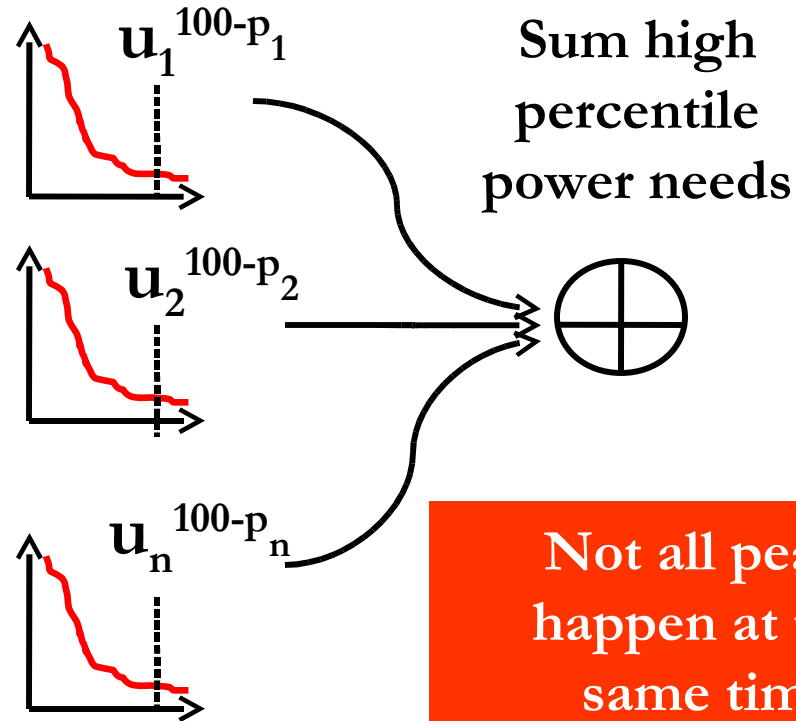**Sum of peaks**

$\oplus$

**Might still be conservative – peaks are rare for bursty applications**

# Under-provisioning Based on Power Profile Tails
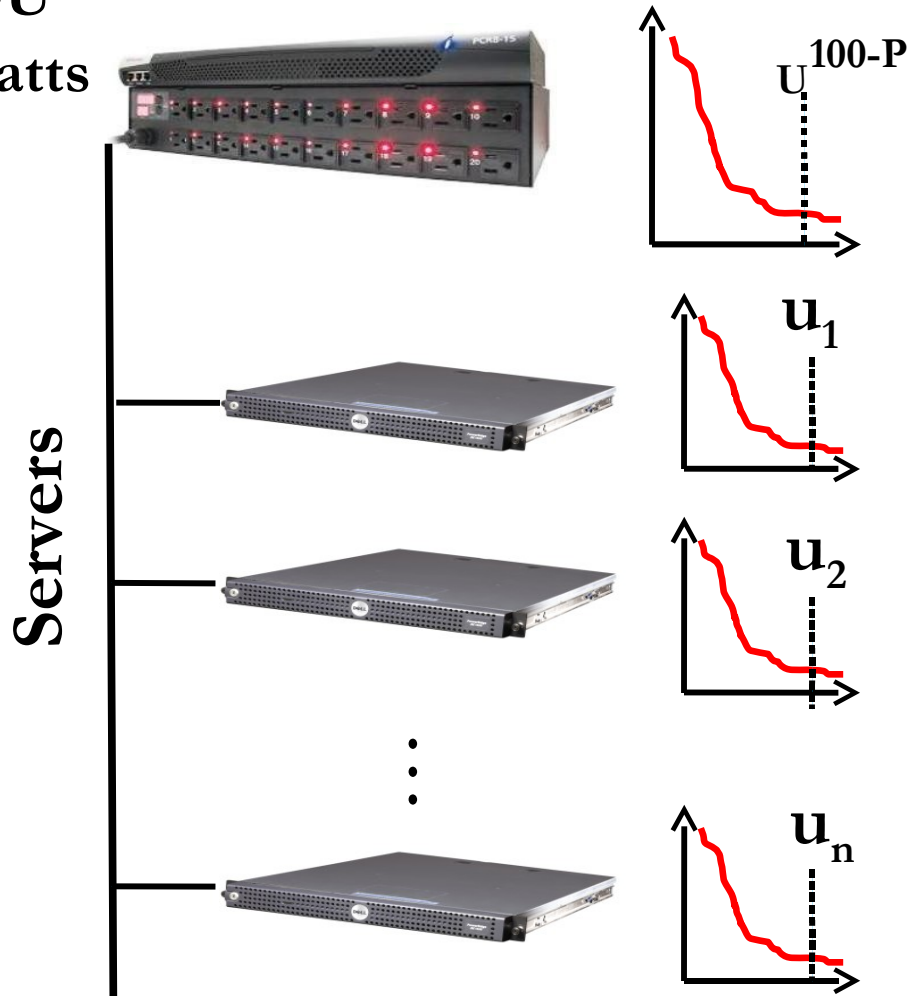
**PDU**
**(B Watts)**

$$\sum_{i=1}^{n} u_i^{100-p_i} \leq B$$

**Servers**

$u_1^{100\text{-}p_1}$

$u_2^{100\text{-}p_2}$

$u_n^{100\text{-}p_n}$

**Sum high percentile power needs**

**Not all peaks happen at the same time**

# Statistical-multiplexing Based Provisioning
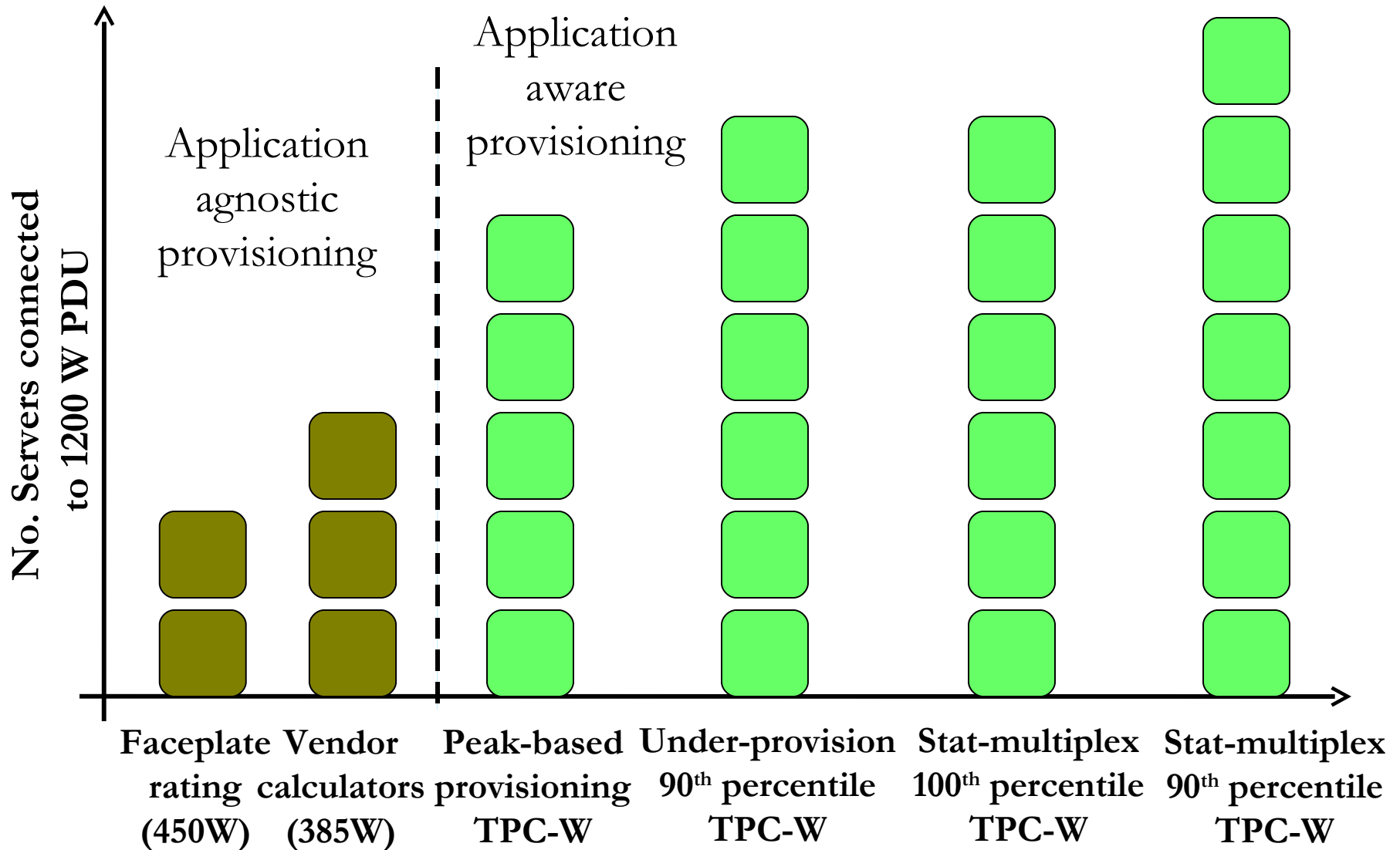


**PDU (B Watts**

**Servers**

$U^{100-P}$

$u_1$

$u_2$

$u_n$

$$U^{100-P} \leq B$$

Provision for the aggregated power profile of the PDU, 'U' as predicted by our sustained power prediction technique

# Provisioning Techniques -Evaluation



**No. Servers connected to 1200 W PDU** (y-axis)

Application agnostic provisioning

Application aware provisioning

| Faceplate rating (450W) | Vendor calculators (385W) | Peak-based provisioning TPC-W | Under-provision 90th percentile TPC-W | Stat-multiplex 100th percentile TPC-W | Stat-multiplex 90th percentile TPC-W |

# Threshold-based Soft-fuse Enforcement

PDU
(1200 W, 5 s)

Periodic power measurement (1s)

Threshold-based Enforcer

Soft fuse
(1200 W, 3 s)

No throttling

1200

Power (W)

Runtime power consumption of the PDU

1  2  3  4  5  6  7  8  9  10 11 12 13 14 15 16  Time (s)

- Hand drawn figure

# Threshold-based Soft-fuse Enforcement



PDU
(1200 W, 5 s)

Periodic power
measurement (1s)

Threshold-based
Enforcer

Soft fuse
(1200 W, 3 s)

Throttling initiated

**Guarantee ??**

1200

Power (W)

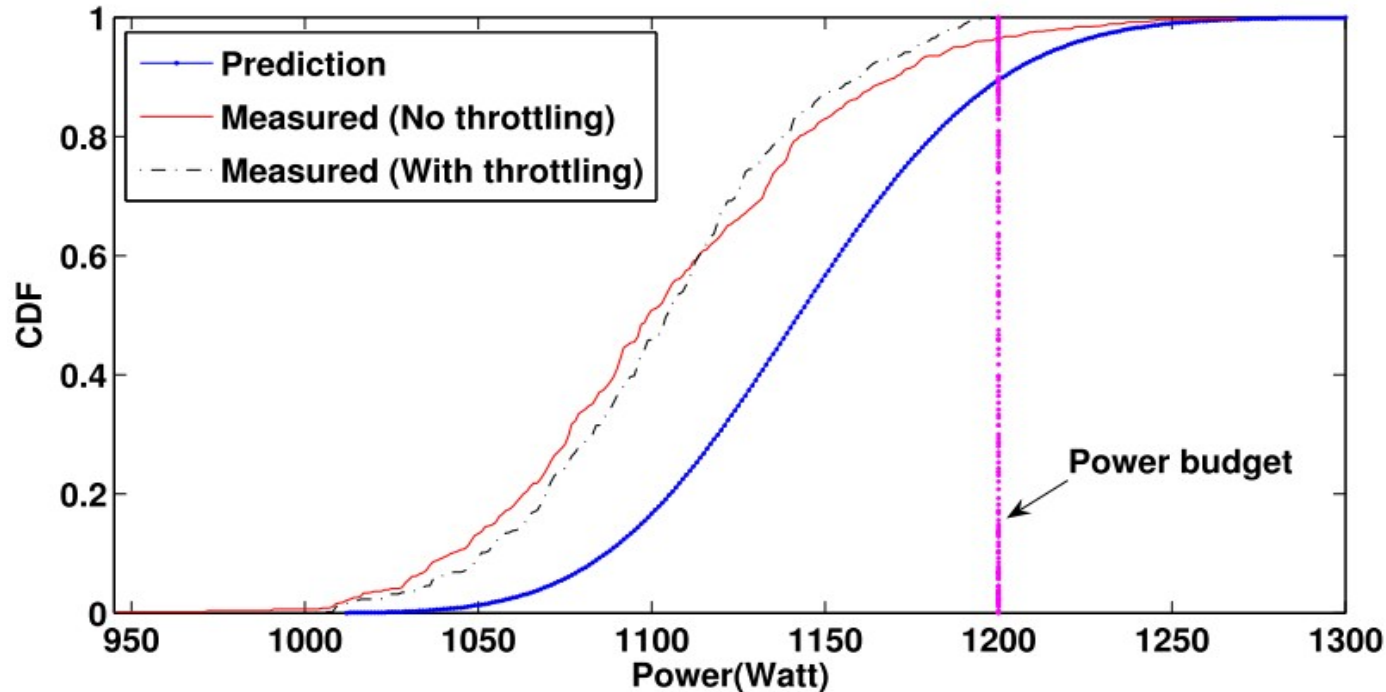1  2  3  4  5  6  7  8  9  10 11 12 13 14 15 16  Time (s)

- Hand drawn figure

# Threshold-based Soft-fuse Enforcement

| Sustained power consumption (100th percentile) of a PDU connected to servers hosting TPC-W | | | | |
|---|---|---|---|---|
| **Power State** | **6 Servers** | **7 Servers** | **8 Servers** | **9 Servers** |
| 3.4 Ghz | **1191.0 W** | 1300.0 W | 1481.0 W | 1672.0 W |
| 2.8 Ghz | 976.6 W | **1138.6 W** | 1308.2 W | 1478.2 W |
| 1.4 Ghz | 861.7 W | 1011.7 W | **1162.7 W** | 1313.6 W |

■ Choose appropriate throttling state that satisfies reliability constraint (1200W, 5s) as highlighted in the table

# Threshold-based Soft-fuse Enforcement



- **Provisioning for the 90th percentile power needs:** Threshold based enforcer is successfully able to enforce soft fuse of the PDU connected to 7 TPC-W servers

# Gains vs Performance Degradation

- **Experiment:** 7 TPC-W servers connected to 1200 W PDU

- **Gains:** Computation per Provisioned Watt

  - Increase in number of servers (computation cycles) hosted in the data center

  - Decrease in number of computation cycles due to throttling

  - CPW increased by 120% from vendor-based provisioning

- **Performance Degradation:**

  - Average response time of TPC-W not affected

  - 95[th] percentile response time of TPC-W increased from 1.59 s to 1.78 s (12% degradation)

# Concluding Remarks

- Power provisioning in data centers
  - Characterize hardware reliability constraints
  - Profile application power consumption
  - Improve provisioning of data center power infrastructure

- Future work
  - Correlated power peaks across servers
  - Handle dynamically varying workload phases

- Software URL: http://csl.cse.psu.edu/hotmap
  - Sustained power prediction scripts
  - Threshold-based soft-fuse enforcer
  - Xen kernel patch for enabling MSR writes