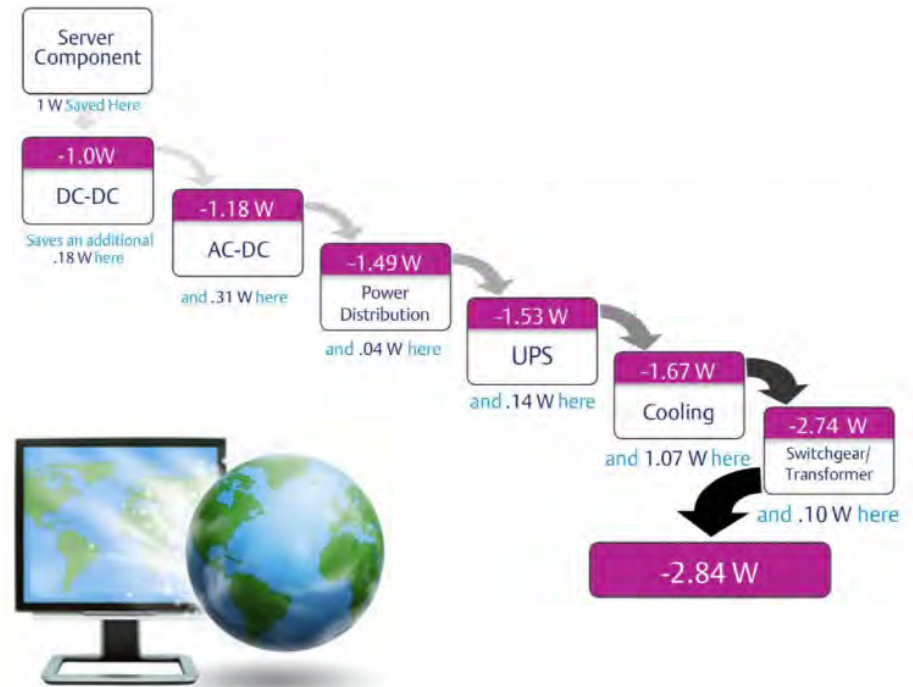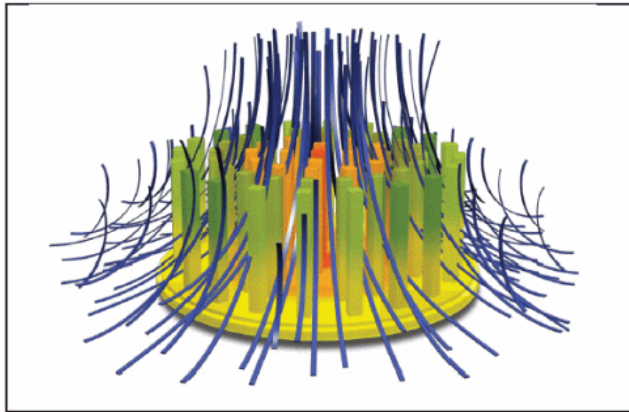# Spart energieeffziente Hardware Strom?

Markus Herber— Senior Technology Consultant, HP
1. Dezember 2009
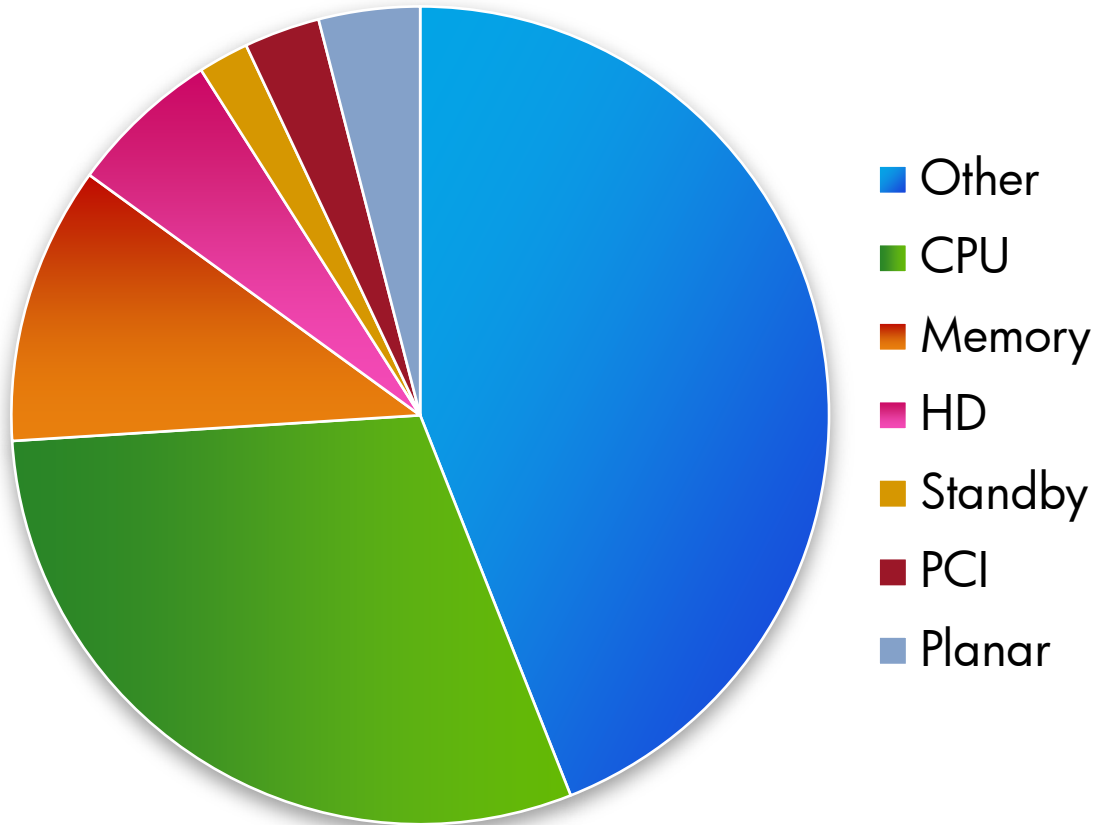
# Was bedeutet "energieeffiziente Hardware"?

– Je weniger Energie für die gleiche Leistung benötigt wird, desto höher ist die Energieeffizienz. Die Energieeffizienz lässt sich in erster Linie durch effizientere Technik steigern.
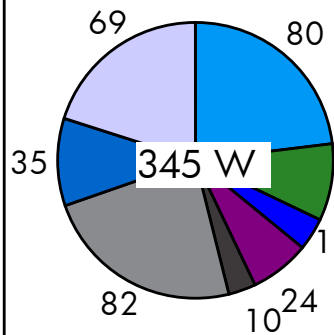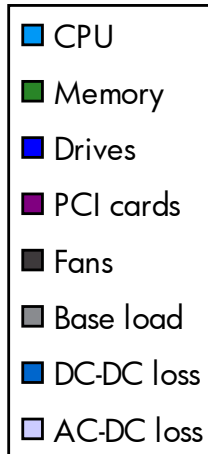
# Wo bestehen Einflussmöglichkeiten?



Legend:
- Other
- CPU
- Memory
- HD
- Standby
- PCI
- Planar

OTHER?
- AC to DC Transitions (25%)
- DC to DC Deliveries (10%)
- Fans and air movement (9%)

# Es gibt beliebige viele Verteilungen



Legend:
- CPU
- Memory
- Drives
- PCI cards
- Fans
- Base load
- DC-DC loss
- AC-DC loss

DL380 G6
2-socket rack*

DL580 G6
4-socket rack*

BL460c/490c
G6 blade
2-socket*†

Typical config
Rank:
1. CPU
2. Base
3. AC-DC
4. DC-DC

Max memory
Rank:
1. Mem
2. CPU
3. AC-DC
4. Base

* Predictions based on early analysis
† Assumes 2 OA's, 2 switches, 10 fans @22 °C amortized across 16 blades

# Power Management Elemente



Processor

Chipset

Platform

Operating System (not user settable policy)

Operating System (user settable policy)

Application

# Low Power Prozessoren – Einsparung?

– Lohnt sich die LP Variante?

– Gibt es Vergleiche?

# Low Power Prozessoren – Einsparung?

Beispiel auf Basis Intel E/L 5520 QuadCore

- Intel® Xeon® Processor E5520 (2.26 GHz, 8MB L3 Cache, **80W**, DDR3-1066, HT, Turbo 1/1/2/2)

- Intel® Xeon® Processor L5520 (2.26 GHz, 8MB L3 Cache, **60W**, DDR3-1066, HT, Turbo 1/1/2/2)

Gleiche Performance aber 2* 20 Watt Unterschied in einem 2 Sockel System

- Preisdifferenz ~200$, bei 2 Sockeln ~400$

- Preisdifferenz Energiekosten nach 3 Jahren: ~263 $

**Welche Performance ist überhaupt notwenig?**

*Berechnungen auf Basis Listpreis und Power Advisor

# Low Power Prozessoren – reicht das?

- Effektives Design (heute 45nm, morgen 32 nm)
- Ausnutzung von Prozessor Performance States
  - Intel's Demand Based Switching
  - AMD's PowerNow!
- AMD's Dual Dynamic Power Management
  - Getrennte Ansteuerung von CPU und Memory Controller
- Intel Nehalem Turbo Boost: mehr Flexibilität im Rahmen der TDP
  - Übertacktung solange die TDP nicht überschritten wird

# P-State

– Processor performance states (P-states) are a predefined set of frequency and voltage combinations at which a given processor can operate correctly. The processor P-state is the capability of running the processor at different voltage and/or frequency levels. Generally, P0 is the highest state resulting in maximum performance, while P1, P2, and so on, will save power but at some penalty to CPU performance.

| Power State | Voltage | Frequency | Heat Generated |
|---|---|---|---|
| P0 | 1.4v | 3.6 GHz | 103 Watts |
| P1 | 1.35v | 3.4 GHz | 94 Watts |
| P2 | 1.3v | 3.2 GHz | 85 Watts |
| P3 | 1.25v | 3.0 GHz | 76 Watts |
| P4 | 1.2v | 2.8 GHz | 68 Watts |

AMD POWERNOW™ TECHNOLOGY WITH OPTIMIZED POWER MANAGEMENT (OPM)

**P-State**

| P0 | 2600MHz 1.40V ~95watts | HIGH |
| P1 | 2400MHz 1.35V ~90watts | |
| P2 | 2200MHz 1.30V ~76watts | PROCESSOR UTILIZATION |
| P3 | 2000MHz 1.25V ~65watts | |
| P4 | 1800MHz 1.20V ~55watts | |
| P5 | 1000MHz 1.10V ~32watts | LOW |

# Demand Based Switching



**Demand Based Switching from P0 to P4 state**

Is DBS enabled? → Terminate

Yes

Is the workload (processor utilization) reduced? → No → Wait and check again

Yes

Use DBS to downshift voltage for the processor

Reduction in voltage in-turn reduces clock-speed

Reduction in voltage and clock speed reduces the heat generated

Cooling needs are reduced enabling cost reduction in data centers

1.4v    3.6 GHz    103 Watts

1.2v

2.8 GHz

68 Watts

**Example power-states**

| Power State | Voltage | Frequency | Heat Generated |
|---|---|---|---|
| P0 | 1.4v | 3.6 GHz | 103 Watts |
| P1 | 1.35v | 3.4 GHz | 94 Watts |
| P2 | 1.3v | 3.2 GHz | 85 Watts |
| P3 | 1.25v | 3.0 GHz | 76 Watts |
| P4 | 1.2v | 2.8 GHz | 68 Watts |

# C-States

| Power State | Execution | Wake-Up Time | CPU Power | Platform | Core Voltage | Cache Shrink | Loss Of Context |
|---|---|---|---|---|---|---|---|
| C0 | Yes | 0ns | large | normal | normal | no | no |
| C1 | No | 10ns | 30% | normal | normal | no | no |
| C2 | No | 100ns | 30% | no I/O buffer | normal | no | no |
| C3 | No | 50,000ns | 30% | I/O + no snoop | normal | no | no |
| C4 | No | 160,000ns | 2% | I/O + no snoop | C4_VID | yes | no |
| C5 | No | 200,000ns | N/A | N/A | C4_VID | L2 = 0KB | no |
| C6 | No | N/A | N/A | N/A | C6_VID | L2 = 0KB | yes |

– With the exception of C0, where the CPU is active and busy doing something, a C-state is an idle state.

– It provides a power savings tradeoff which depends on the length of time the CPU sleeps. The deeper the sleep, the longer it takes for the CPU to wake up, but the more power you save. The operating system selects which state you enter, based on when you anticipate the CPU will be waking up.

– Intel® processors based on *Nehalem* microarchitecture support core C0, C1, C3, and C6. C0 and C1 are always supported; the availability of the remaining C-states may vary by processor number. Any core within the processor can go into any C-state independent of the state of the other cores.
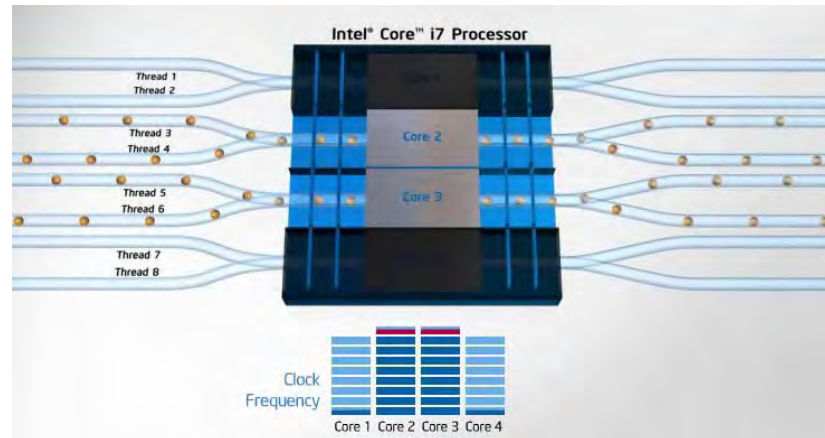
# T-State

– "T"-states (throttling states) which will further throttle down a CPU (but not the actual clock rate) by inserting STPCLK (stop clock) signals and thus omitting duty cycles.

– The T state is one of the three execution states that CPUs execute code in. Much like the emergency brake of a car, the T state is used to forcefully reduce the CPU's execution speed.

# Intel® Turbo Boost Technology

– **Intel Turbo Boost Technology** allows a processor's cores to run faster than the base operating frequency if the package is operating below its power, current, and temperature specification limits. Intel Turbo Boost Technology is activated when the OS requests the highest processor performance state (P0). Maximum frequency depends on the number of active cores. The amount of time the processor spends in the Intel Turbo Boost Technology state depends on the workload and operating environment, providing the extra performance.
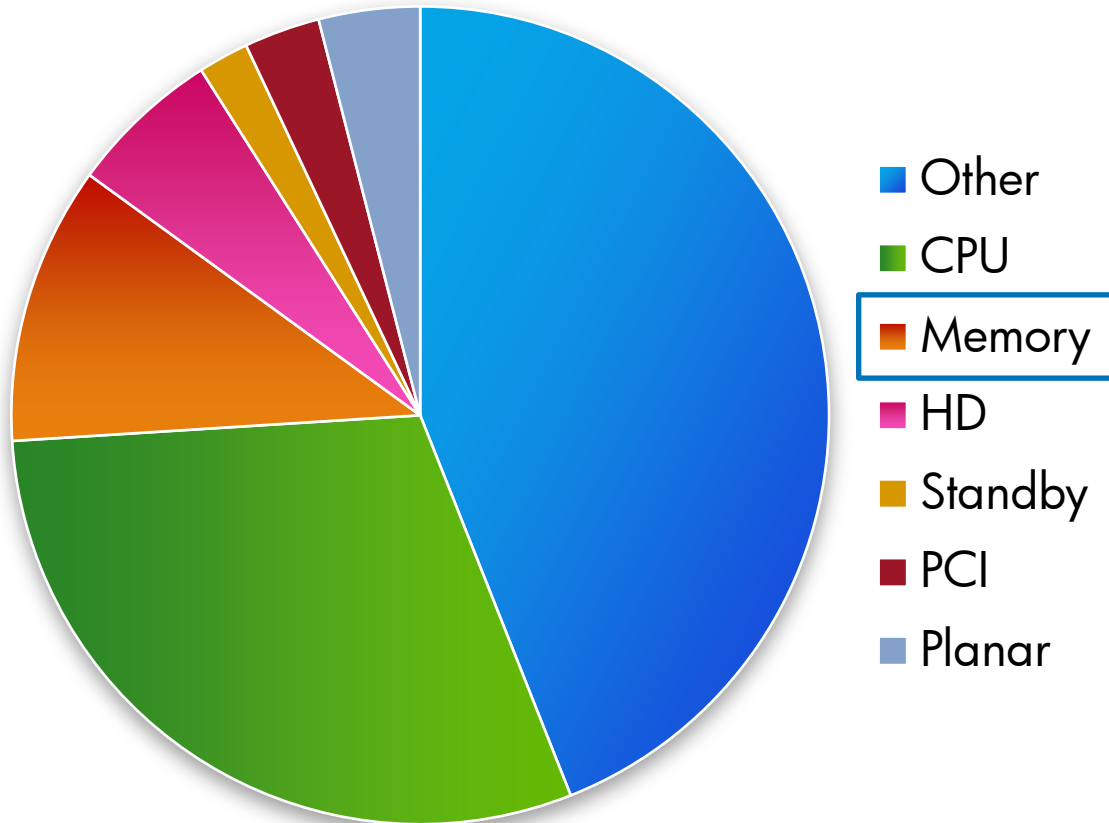
# AMD-P Technologies

AMD Opteron processors offer outstanding energy efficiency

- **Enhanced AMD PowerNow!™** Technology with Independent Dynamic Core Technology allows each core to vary its frequency, based on the specific needs of the application. This allows for more precise power management to reduce data center energy consumption and thereby reduce total cost of ownership (TCO).

- **Dual Dynamic Power Management™** allows each processor to maximize the power-saving benefits of Enhanced AMD PowerNow! technology without compromising performance. Dual Dynamic Power Management can reduce idle power consumption and allow for per-processor power management in multi-socket systems to decrease power consumption.

- **AMD CoolCore™ Technology** evaluates which parts of the die - the cores, the memory, or both - are needed to support currently running applications. It can cut power to unused transistor areas to reduce power consumption and lower heat generation.

- **AMD PowerCap Manager** gives an IT manager the ability to put a cap on the P-state level of the cores via the BIOS. This can help reduce processor power consumption of a system.

- **AMD Smart Fetch Technology** enables inactive cores to write contents of their L1 and L2 caches to the shared L3 cache. This can allow the inactive cores to enter a "halt" state and draw less power, reducing CPU power consumption
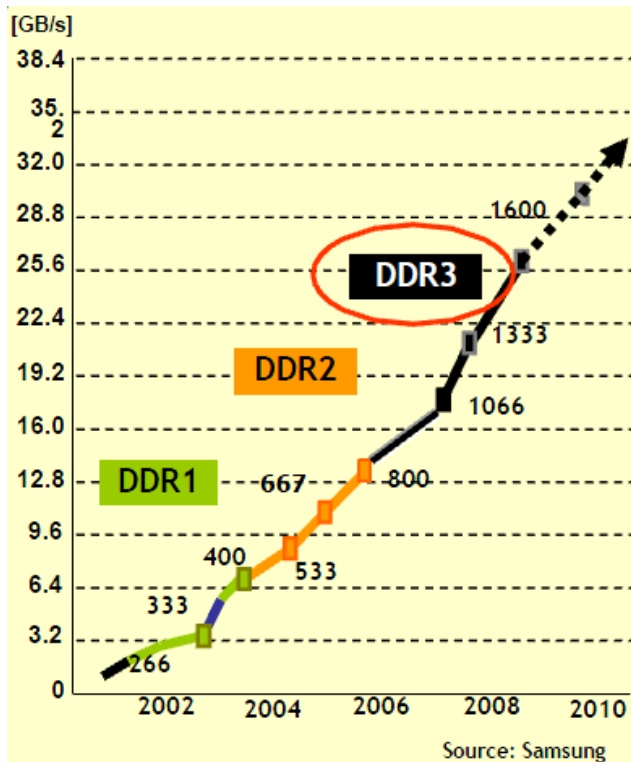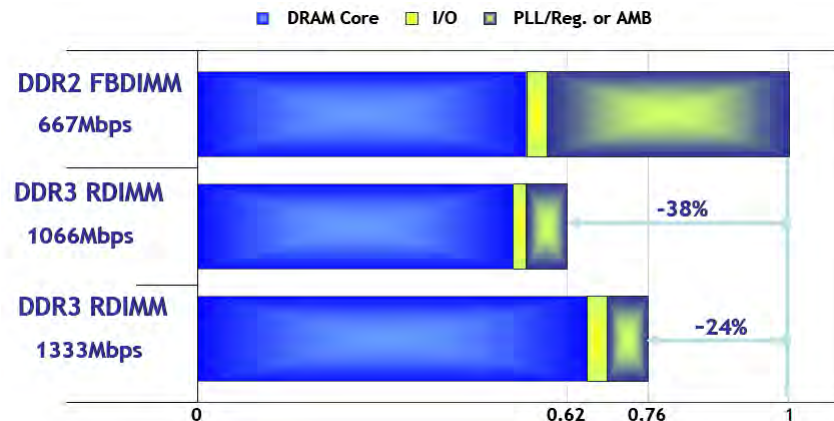
# Wo bestehen Einflussmöglichkeiten?



- Other
- CPU
- Memory
- HD
- Standby
- PCI
- Planar

OTHER?
- AC to DC Transitions (25%)
- DC to DC Deliveries (10%)
- Fans and air movement (9%)

# Speicherausbau/Entwicklung



Source: Samsung

- DDR3 has lower power architecture, due to lower core voltage

- >25% power savings over DDR2 (DDR2-800 vs. DDR3-800)
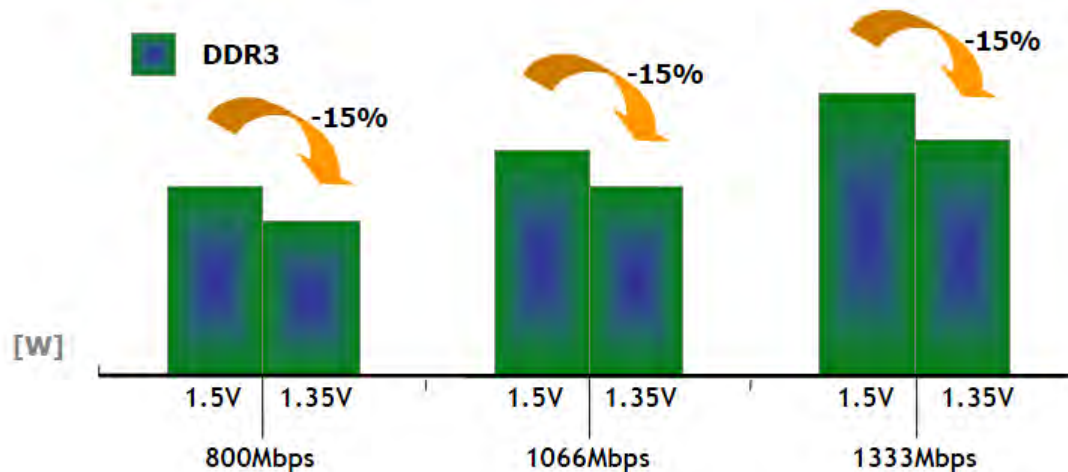
- DDR3-1066 consumes less power than DDR2-800.



Source: Samsung

# Speicherausbau – was bring die Zukunft?

– Over time, DDR3 memory will consist of three voltage ratings; Standard at announce (1.5V); future plans call for Low Voltage (1.35V) and Ultra Low Voltage (TBD; ~1.25V)

– Early data with 1.35V shows approx 15 % reduction in total power

– Reduce output-swing voltage for the driver reduces with 18 DIMMs the power usasge of the system by 22W
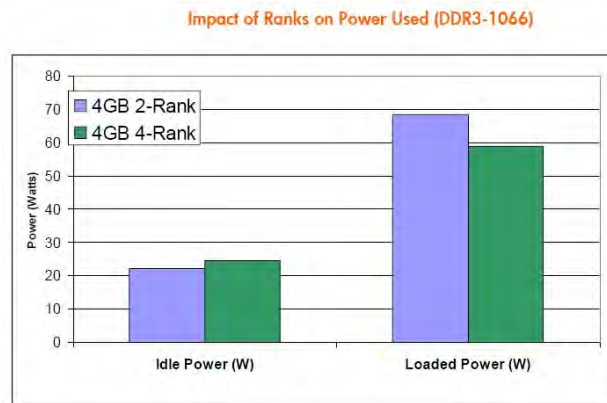


Source: Samsung

# Speicherausbau – Einsparung?

– Kleiner Module sind günstiger, Ausbaufähigkeit begrenzt, geringere Verbrauch

– Größere Module sind wesentlich teurer, besser Skalierbar, mehr Verbrauch

– Quad-Ranked DIMMS verbrauchen ca. 15% weniger als dual-ranked DIMMS

**Impact of Ranks on Power Used (DDR3-1066)**

- 4GB 2-Rank
- 4GB 4-Rank

Power (Watts): Idle Power (W), Loaded Power (W)

– Low Power DIMMs: In the dual-rank DIMM half of the DRAM chips are active at a time. In the quad-rank DIMM, a quarter of the DRAM chips are active at a time resulting in 15% less power consumption.

# Speicher Optimierung: Energieverbrauch

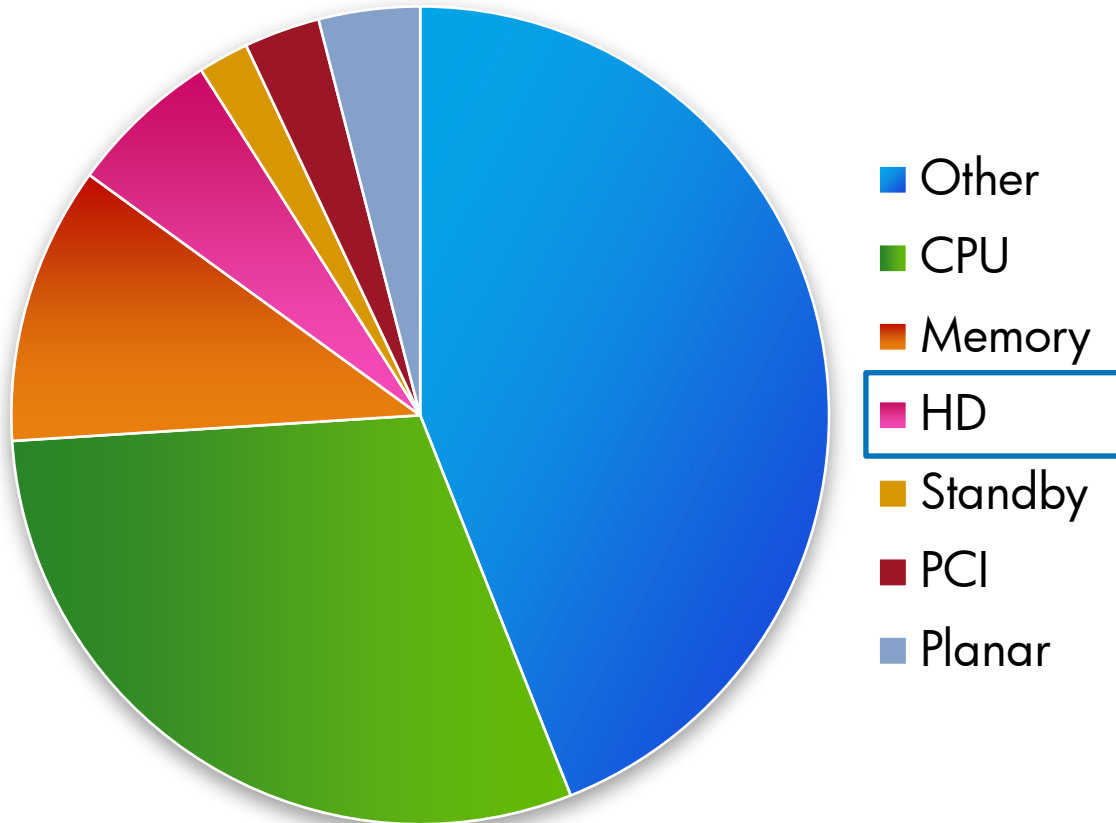**Table 14 : Power of various 24GB Configurations**

| Total Memory (GB) | Memory Config | Number of DIMMs | DIMM Size | DIMM Rank | UDIMM or RDIMM | DIMM Speed | DDR3 Loaded Power (W) |
|---|---|---|---|---|---|---|---|
| 24 | 6x4G4R_800_R | 6 | 4GB | 4 | RDIMM | 800 | 51.82 |
| 24 | 6x4G4R_1067_R | 6 | 4GB | 4 | RDIMM | 1066 | 58.96 |
| 24 | 6x4G2R_800_R | 6 | 4GB | 2 | RDIMM | 800 | 59.34 |
| 24 | 12x2G2R_800_U | 12 | 2GB | 2 | UDIMM | 800 | 59.88 |
| 24 | 12x2G2R_1067_U | 12 | 2GB | 2 | UDIMM | 1066 | 67.33 |
| 24 | 6x4G2R_1067_R | 6 | 4GB | 2 | RDIMM | 1066 | 68.45 |
| 24 | 12x2G2R_800_R | 12 | 2GB | 2 | RDIMM | 800 | 72.99 |
| 24 | 6x4G2R_1333_R | 6 | 4GB | 2 | RDIMM | 1333 | 75.24 |
| 24 | 12x2G2R_1067_R | 12 | 2GB | 2 | RDIMM | 1066 | 80.29 |
| 24 | 12x2G2R_1333_R | 12 | 2GB | 2 | RDIMM | 1333 | 87.04 |

**Table 15 : Power of various 48GB Configurations**

| Total Memory (GB) | Memory Config | Number of DIMMs | DIMM Size | DIMM Rank | UDIMM or RDIMM | DIMM Speed | DDR3 Loaded Power (W) |
|---|---|---|---|---|---|---|---|
| 48 | 6x8G2R_800_R | 6 | 8GB | 2 | RDIMM | 800 | 52.79 |
| 48 | 6x8G2R_1067_R | 6 | 8GB | 2 | RDIMM | 1066 | 58.92 |
| 48 | 6x8G2R_1333_R | 6 | 8GB | 2 | RDIMM | 1333 | 62.96 |
| 48 | 12x4G4R_800_R | 12 | 4GB | 4 | RDIMM | 800 | 90.23 |
| 48 | 12x4G2R_800_R | 12 | 4GB | 2 | RDIMM | 800 | 106.08 |
| 48 | 12x4G2R_1067_R | 12 | 4GB | 2 | RDIMM | 1066 | 119.97 |
| 48 | 12x4G2R_1333_R | 12 | 4GB | 2 | RDIMM | 1333 | 132.80 |

# Wo bestehen Einflussmöglichkeiten?



Legend:
- Other
- CPU
- Memory
- HD
- Standby
- PCI
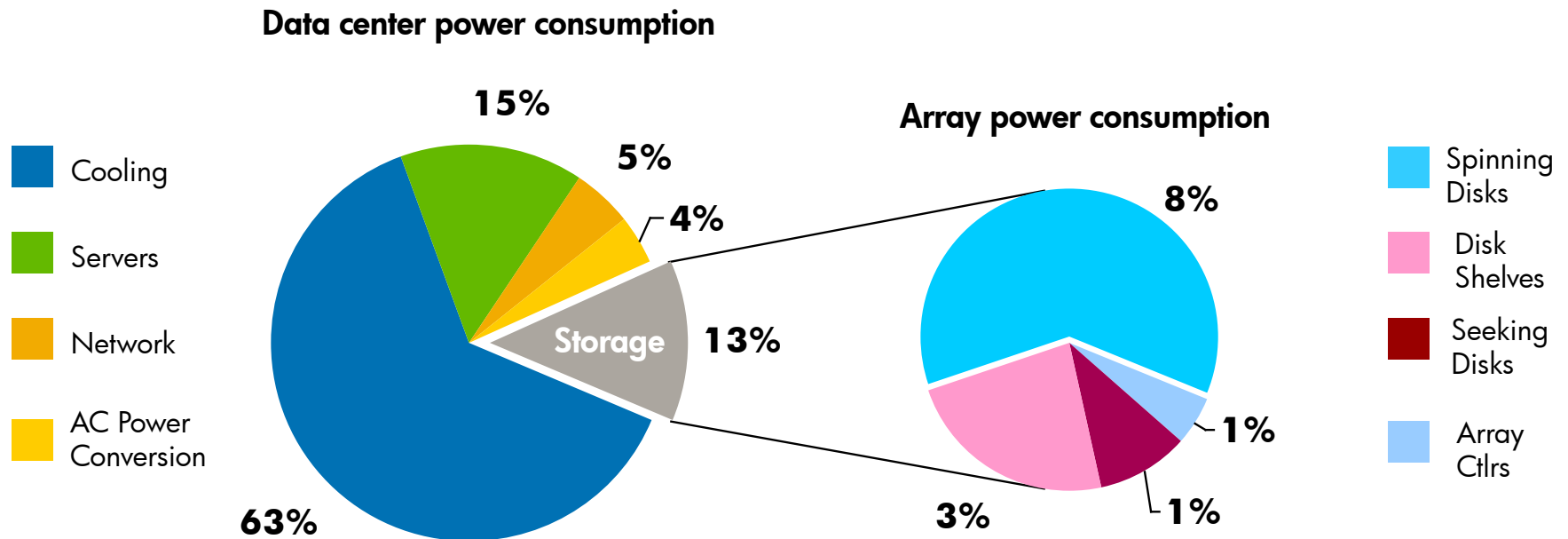- Planar

OTHER?
- AC to DC Transitions (25%)
- DC to DC Deliveries (10%)
- Fans and air movement (9%)

# Stromverbrauch bei Storagesystemen

**Data center power consumption**

**Array power consumption**

| | |
|---|---|
| ■ Cooling | |
| ■ Servers | |
| ■ Network | |
| ■ AC Power Conversion | |

15%
5%
4%
**Storage** 13%
63%

8%
1%
3%
1%

| | |
|---|---|
| ■ Spinning Disks | |
| ■ Disk Shelves | |
| ■ Seeking Disks | |
| ■ Array Ctlrs | |

## ~**60%** des Stromverbrauchs von Storage Systemen ist für den Betrieb (spinning) der Platten notwendig

# Solid State Drive (SSD)

– **Extreme Ruggedness**
  - Extended Operating Temperature (0° up to 70°C)
  - Shock and Vibration almost a non-issue

– **High Read Performance**
  - \> 50x SATA random read performance
  - \> 15x SAS random read performance
  - No seek time means high IOPS
  - Limited write performance (relative to 15k SAS)

– **Increased Reliability**
  - No moving parts
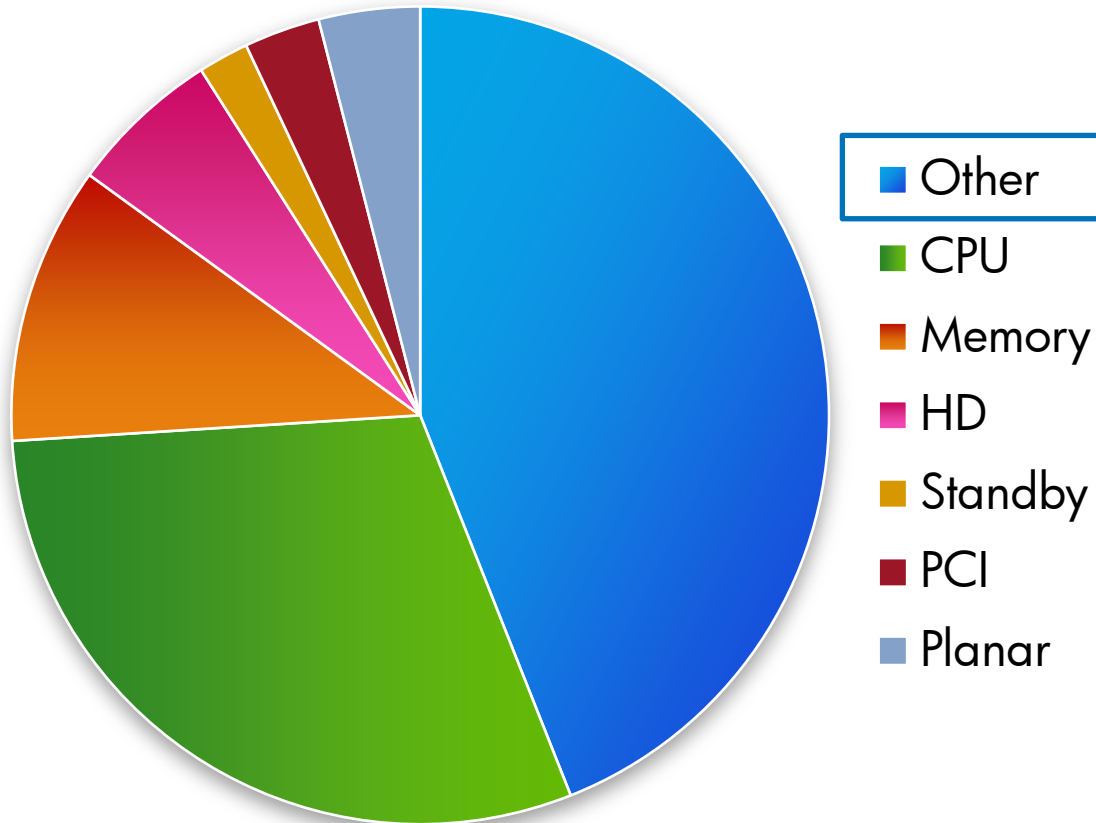  - Virtually eliminates the need to RAID SSDs

– **Up to 10x Reduced Power**
  - < 2 Watts, versus 9 Watt for 15k 2.5" SAS

– **Thermal, Size and Acoustic Advantages**
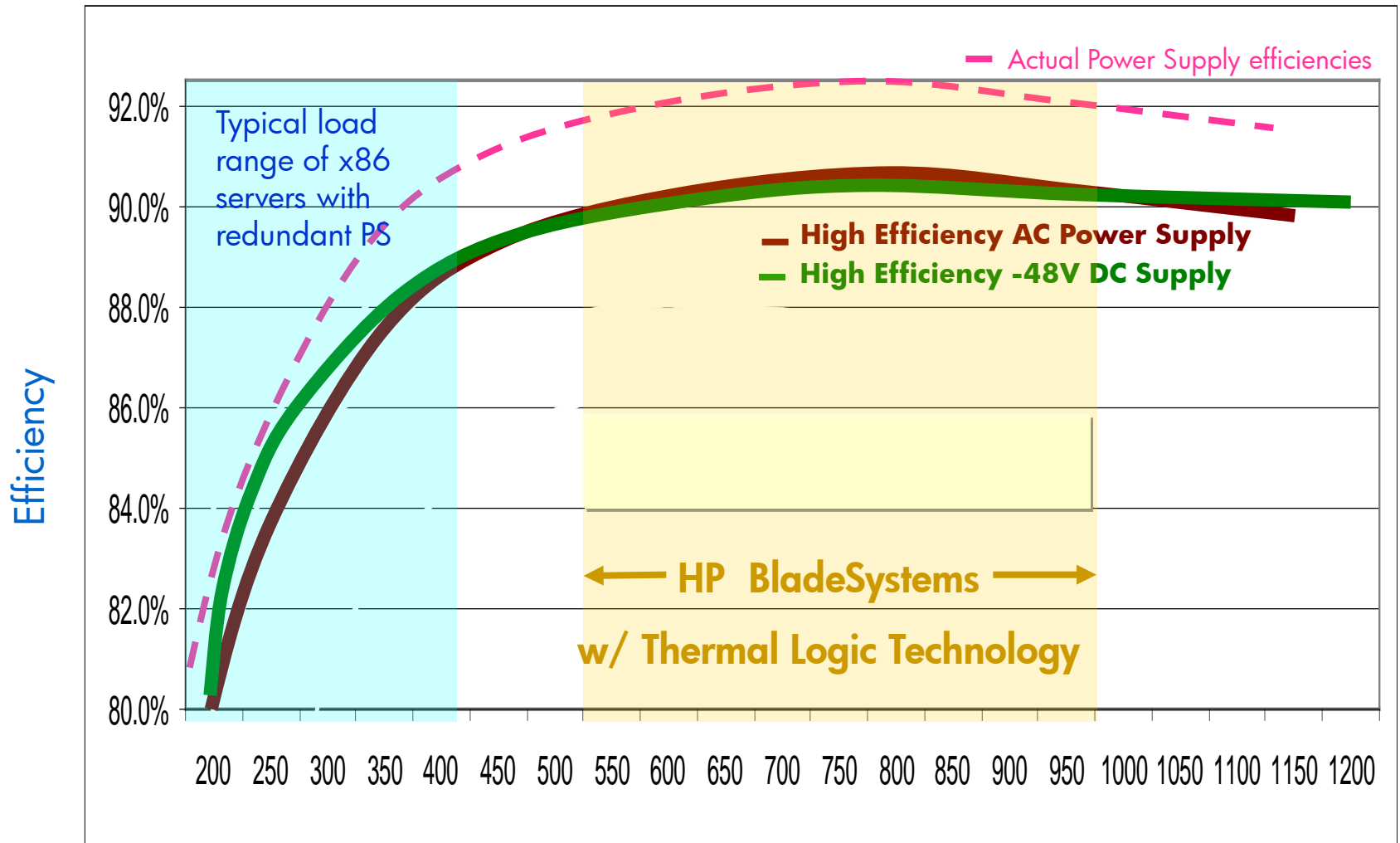  - No Noise, Low Heat
  - Small and light weight

# Wo bestehen Einflussmöglichkeiten?



Other
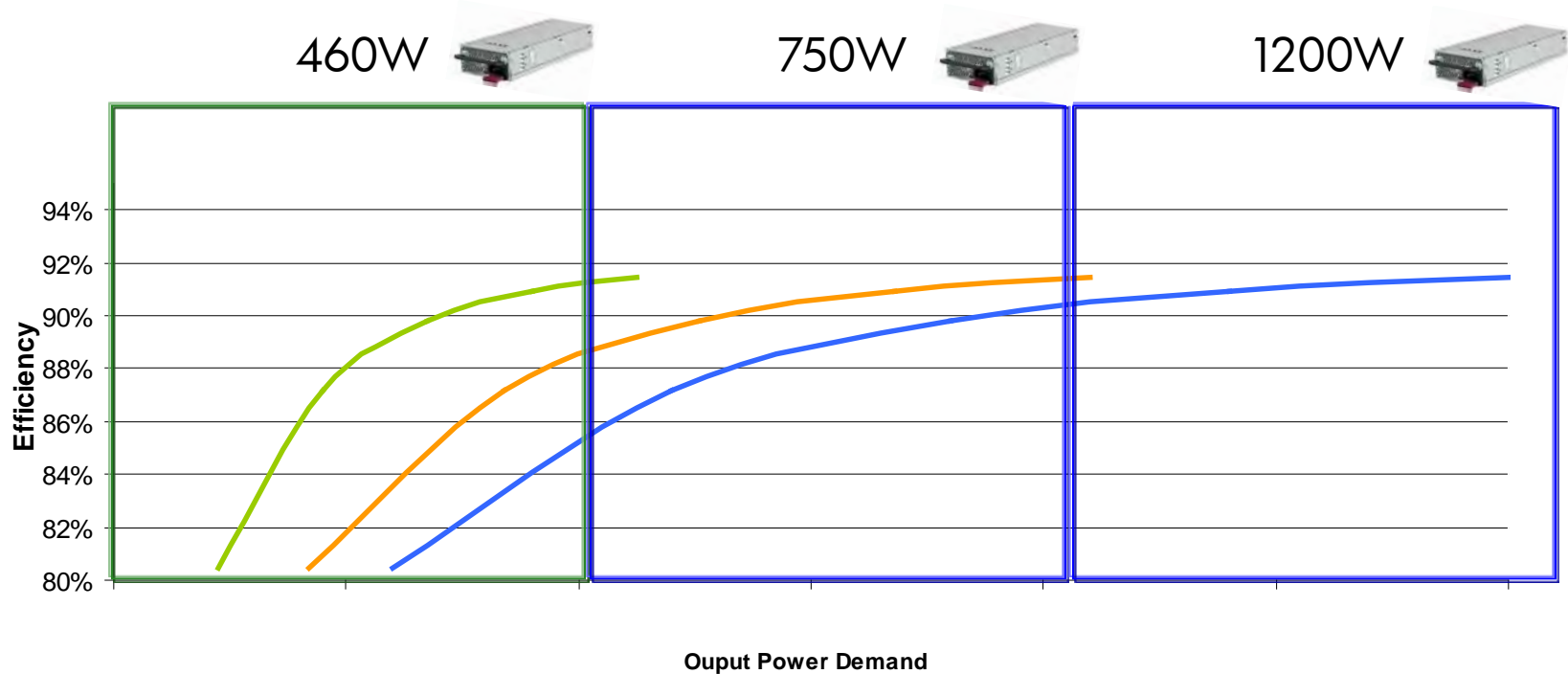CPU
Memory
HD
Standby
PCI
Planar

OTHER?
• AC to DC Transitions (25%)
• DC to DC Deliveries (10%)
• Fans and air movement  (9%)

# Power Supply Efficiency Curves

# Right-sizing your Power Supply

– Stay on the healthy part of the curve…

– choose the Power Supply to fit your application

# Netzteile

Ein Netzteil-Steckplatz für unterschiedliche Systeme



ProLiant DL160 G5
ProLiant DL165 G5
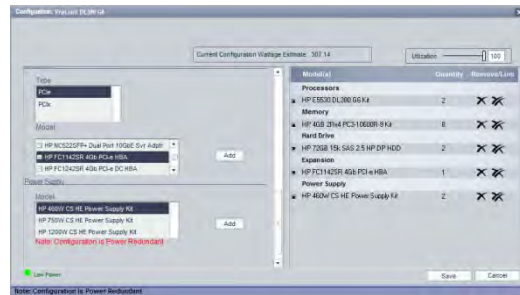ProLiant DL180 G5
ProLiant DL185 G5
ProLiant DL380 G6
ProLiant DL385 G6
ProLiant DL580 G5
ProLiant DL785 G5
ProLiant ML/DL350 G6 and ML/DL370 G6

**Common Slot**

**460W AC up to 92% efficiency**

**750W AC up to 92% efficiency**

**1200W AC up to 90% efficiency**

**48Vdc 1200W up to 90% efficiency**

HP Proliant Power Advisior

# Netzteile

307 W - 460 W Netzteil
352 W - 1200W Netzteil

~13% Stromersparnis



geringere
Anschaffungskosten

# Mehr Effektivität durch Dynamic Power Saver



| 1+1 up to 33% | 2+2 up to 66% | 3+3 up to 100% |

**Power Supply Efficiency** (vertical axis)

More than 90% efficiency maintained from 8% load

With DPS    without DPS    First Gen

**Power Output** (horizontal axis)
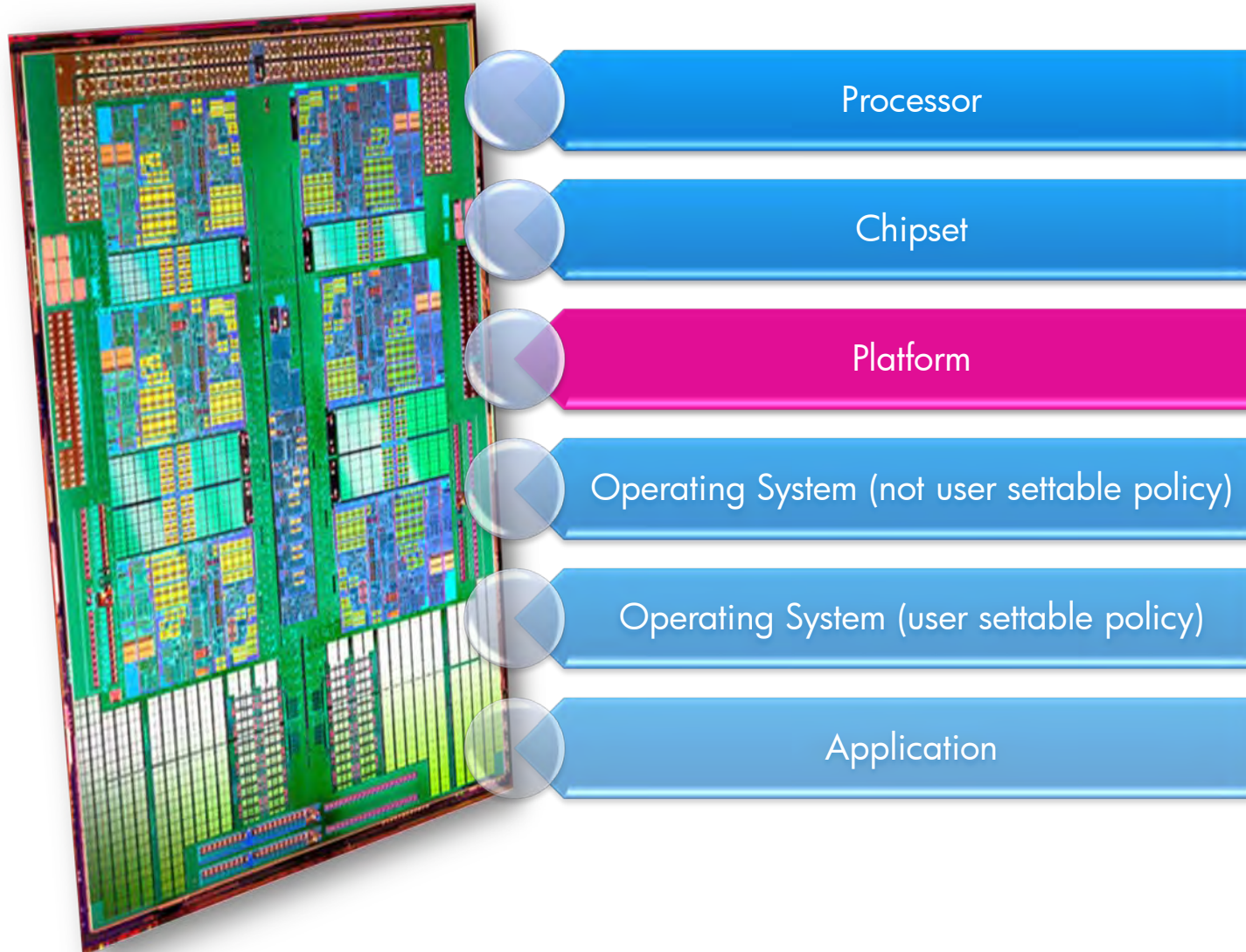
# Intelligente Lüfter Steuerung

## Power Efficient Fan Control

- 30+ temperature sensors in G6 servers
- DIMM temperature monitoring
- Hard Drive temperature monitoring
- CPU sensor(s), air inlet sensor, power supply temps
- Additional fan zones – individual fan control
- Intelligent iLO fan speeds using process control algorithms
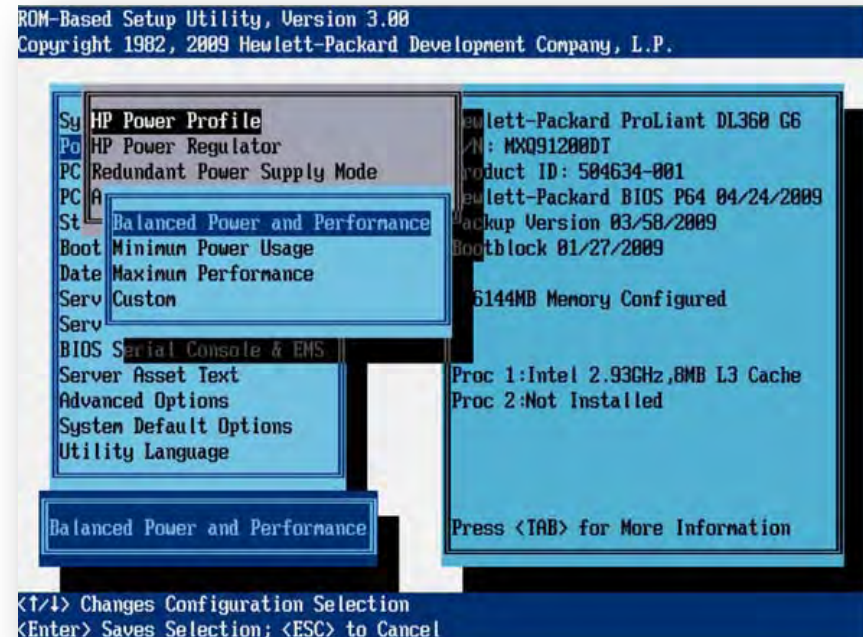- Lower fan RPMs result in power and noise savings

# Power Management Elemente

Processor

Chipset

Platform

Operating System (not user settable policy)

Operating System (user settable policy)

Application

# Power Profile

- **Balanced Power and Performance** provides the optimum settings to maximize power savings with minimal performance impact for most operating systems and applications.

- **Minimum Power Usage** enables power reduction mechanisms that may affect performance negatively. This mode guarantees a lower maximum power usage by the system.

- **Maximum Performance** disables all power management options that may affect performance negatively.

# Power Regulator

- **HP Dynamic Power Savings Mode**
  - Automatically varies processor speed and power usage based on processor use (adjust P-States)
  - Reduces overall power consumption with little or no impact to performance
  - Does not require OS support
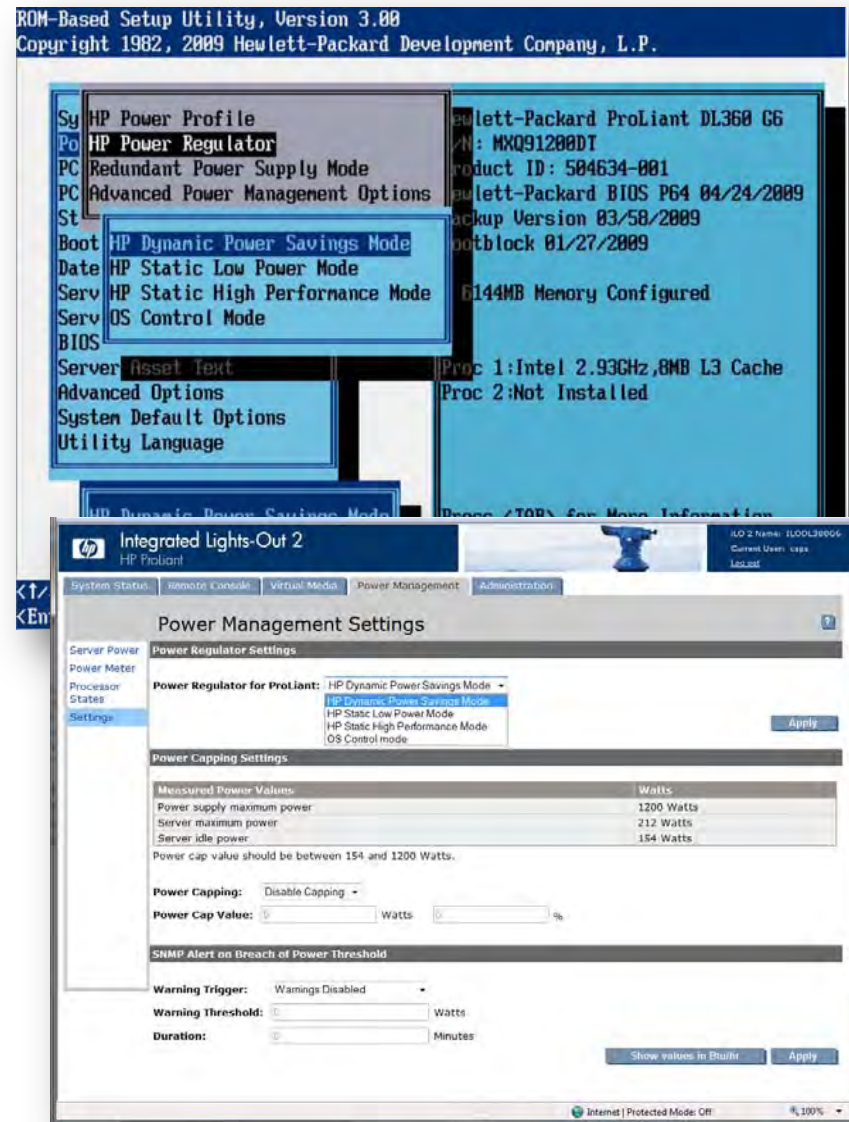
- **HP Static Low Power Mode**
  - Reduces processor speed and power usage
  - Guarantees a lower maximum power usage for the system
  - The impact on performance is greater for environments with higher processor utilization.

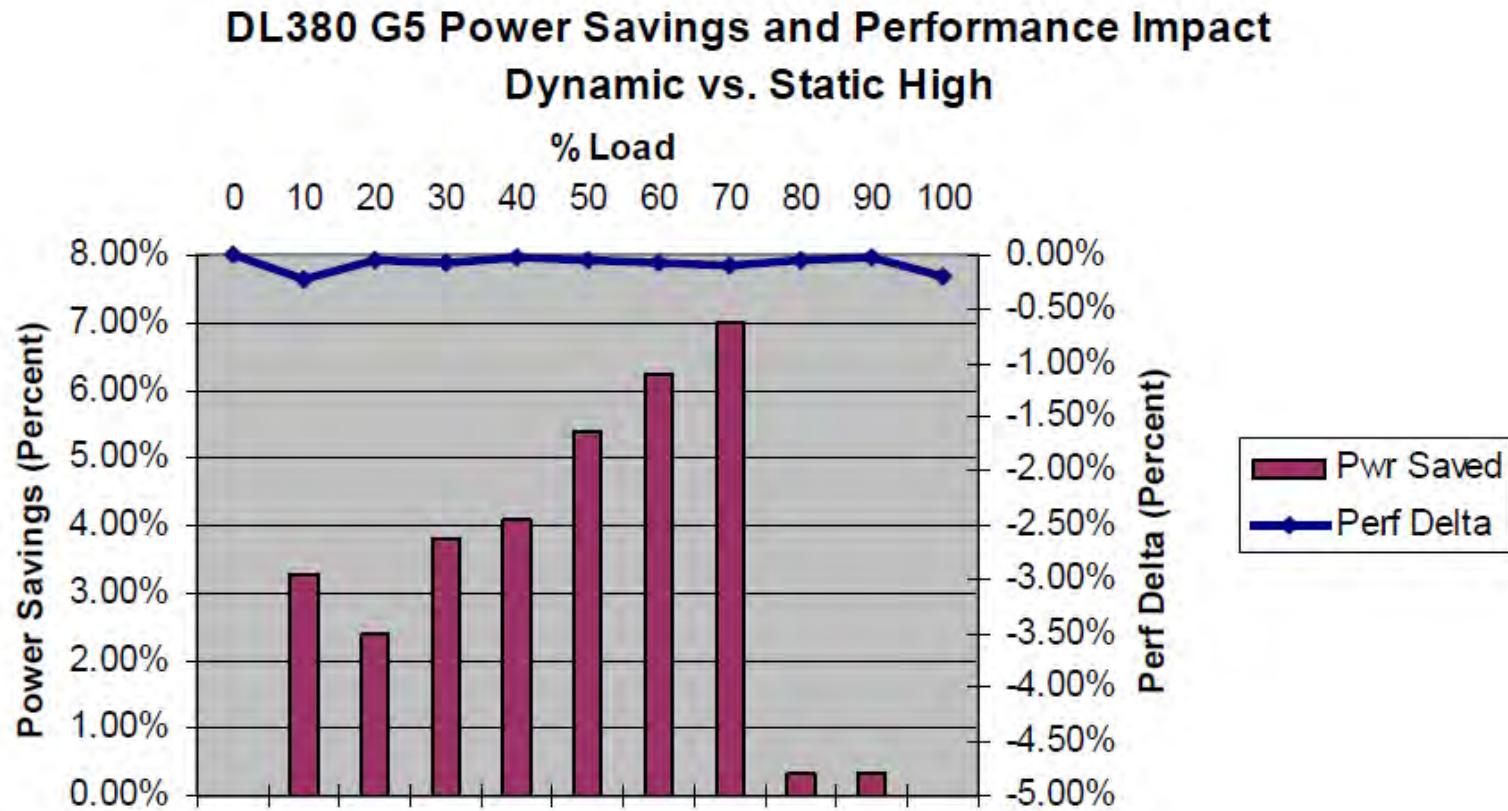- **HP Static High Performance Mode**
  - Processors run in the maximum power and performance state, regardless of the OS power management policy.

- **OS Control Mode**
  - Processors run in the maximum power and performance state, unless the OS enables a power management policy.

# Dynamic Power Savings Mode im Vergleich



DL380 G5 Power Savings and Performance Impact
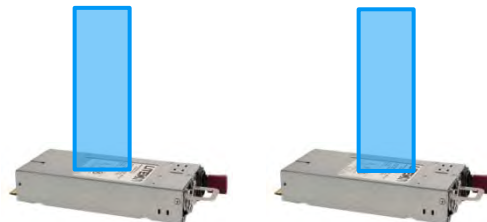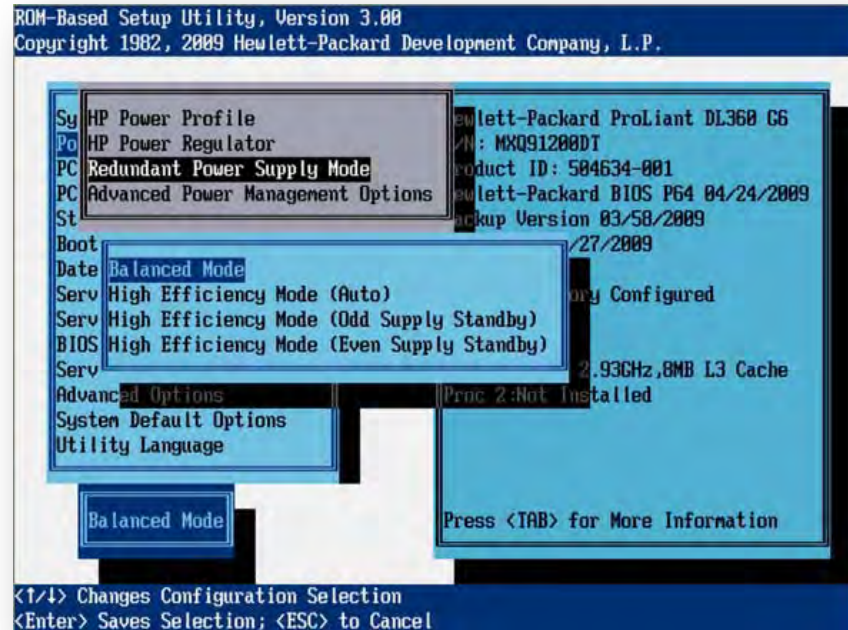Dynamic vs. Static High

# Redundant Power Supply Mode

– This feature enables the user to configure how the system handles redundant power supply configurations

– The High Efficiency Mode options allow the user to choose which power supply is placed in standby.
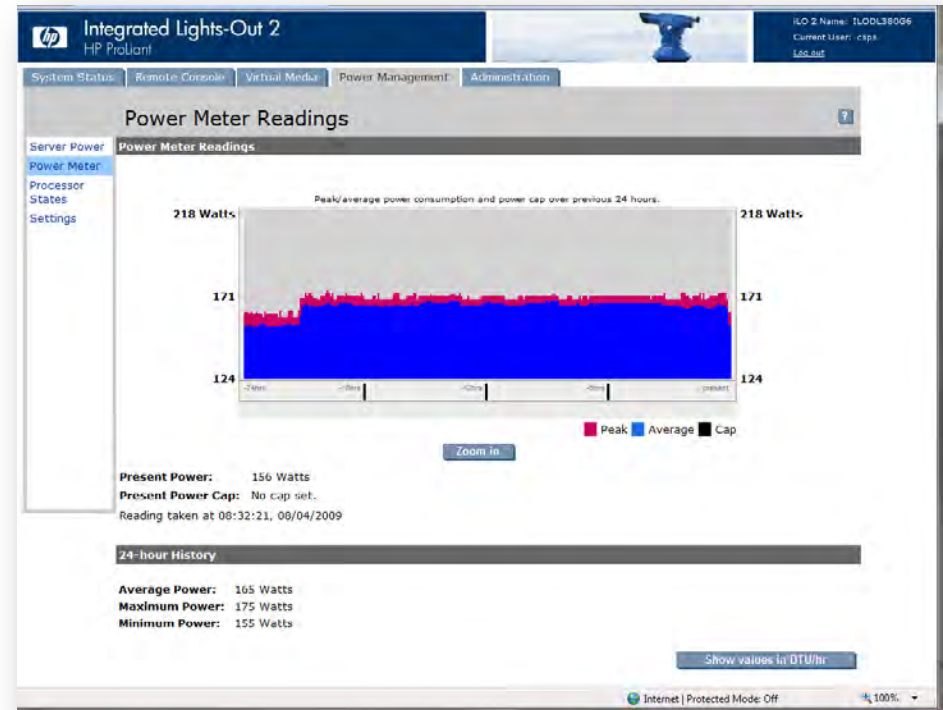


Balanced Mode

**Optimum power efficiency even under light loads**



```
ROM-Based Setup Utility, Version 3.00
Copyright 1982, 2009 Hewlett-Packard Development Company, L.P.

Sy  HP Power Profile                    ewlett-Packard ProLiant DL360 G6
Po  HP Power Regulator                  N: MXQ91200DT
PC  Redundant Power Supply Mode         oduct ID: 504634-001
PC  Advanced Power Management Options   ewlett-Packard BIOS P64 04/24/2009
St                                      ackup Version 03/58/2009
Boot  Balanced Mode                     /27/2009
Date  Balanced Mode
Serv  High Efficiency Mode (Auto)              ory Configured
Serv  High Efficiency Mode (Odd Supply Standby)
BIOS  High Efficiency Mode (Even Supply Standby)
Serv                                    2.93GHz,8MB L3 Cache
Advanced Options                        Proc 2:Not Installed
System Default Options
Utility Language

Balanced Mode                          Press <TAB> for More Information

<↑/↓> Changes Configuration Selection
<Enter> Saves Selection; <ESC> to Cancel
```

# Visualiserung: Power Meter

– The graph shows a 24-hour history with 5 minute samples.

– A 20-minute history with 10 second samples is available with the "Real Time" graph when a power cap is supported and configured.

– iLO periodically samples peak power, average power and power cap.

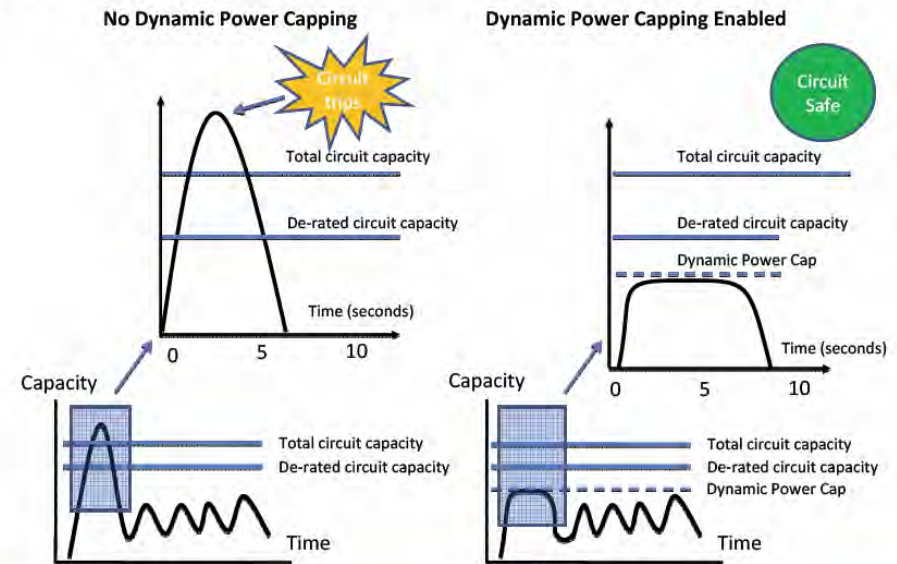– An iLO 2 license is required to view the data

# Differences between Dynamic Power Capping and Power Capping

– HP Dynamic Power Capping monitors power consumption and maintains a server's power cap much more rapidly than HP Power Capping.
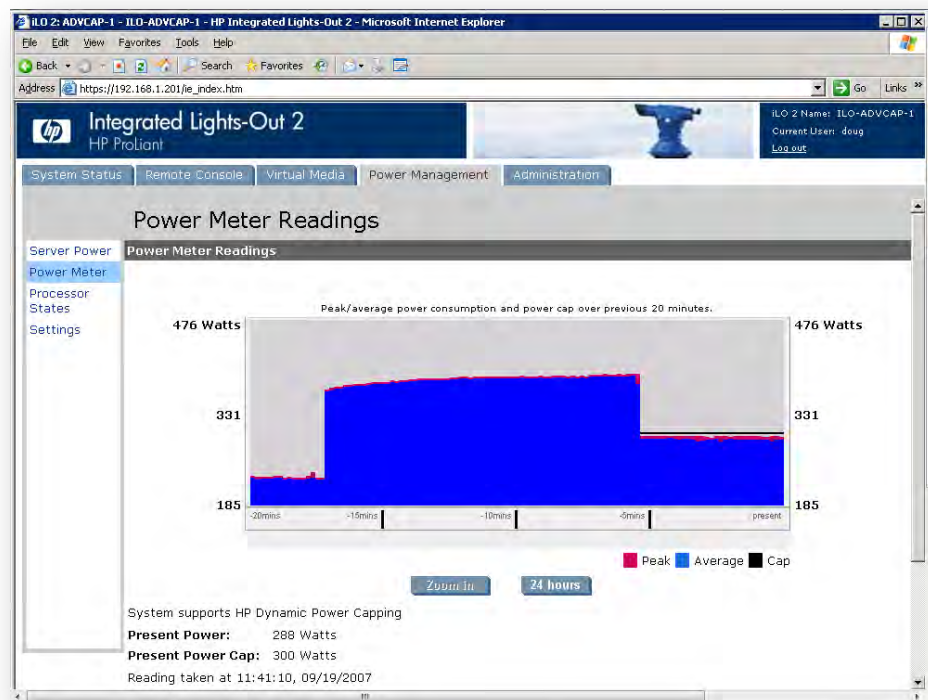
Figure 2. Rapid response of Dynamic Power Capping avoids circuit breaker trips

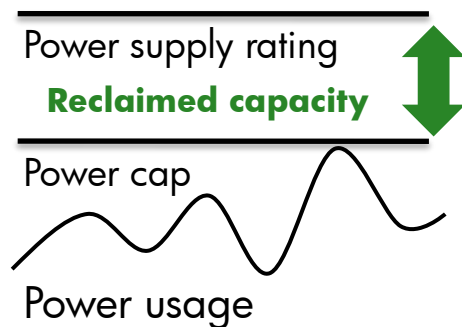| | Dynamic Power Capping | Basic Power Capping |
|---|---|---|
| Power capping executed from | Power management microcontroller | iLO and system ROM BIOS |
| Control of processor power | Direct hardware connection to processor to control P-state/clock throttling at the processor hardware level | Firmware control of P-state/clock throttling through processor registers |
| Power monitoring cycle | More than 5 times per second | Once every 5 seconds |
| Time to bring server power consumption back under its cap | Less than 0.5 seconds | 10 – 30 seconds |
| Intended application | Managing power and cooling provisioning | Managing cooling provisioning |

# Single Server Dynamic Power Cap

- In this graph, the server starts out idle.

- The workload and the system utilization and power consumption went up dramatically.  The peak power consumption was about 350 W.

- Setting a 300 W Dynamic Power Cap.  The Dynamic Power Cap is represented by the black line.
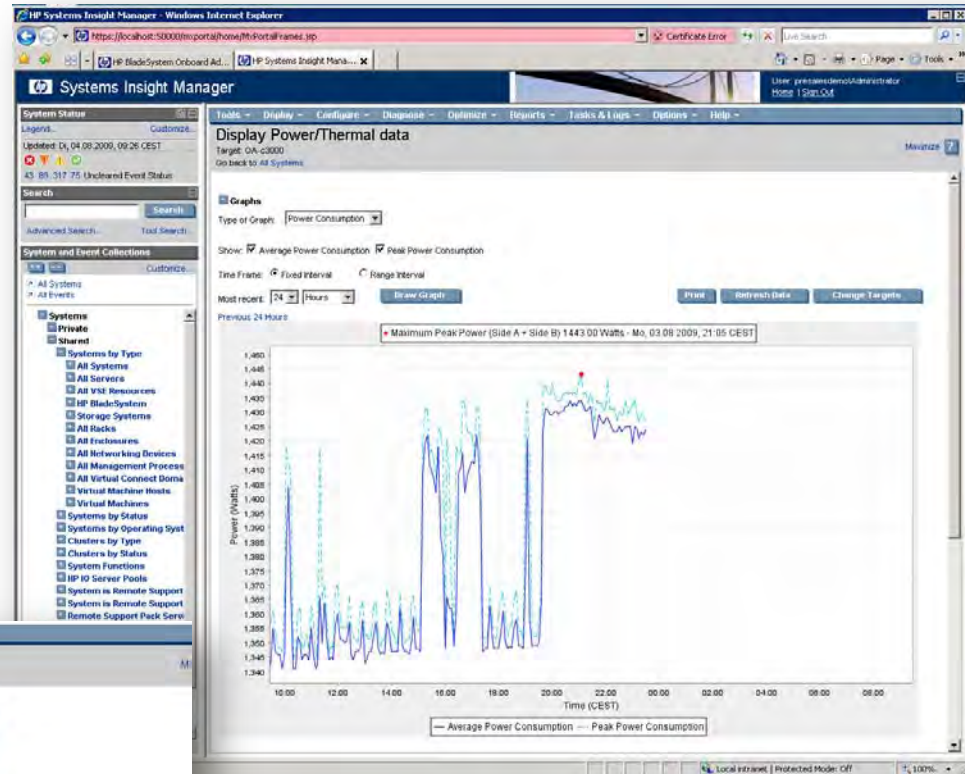
# Enclosure Dynamic Power Capping

- Measure server power consumption
- Cap servers based on measurements
- Workload changes with time
- Change the power caps
  - Busy servers get more power
  - Idle servers get less



Power supply rating

**Reclaimed capacity**

Power cap

Power usage

# Insight Power Manager

- Talks directly to iLO

- Use power history to determine actual power usage

- Set power cap to observed peak

- Monitor power usage to detect potential performance impact

# Power Management Elemente



Processor

Chipset

Platform

Operating System (not user settable policy)

Operating System (user settable policy)

Application

# OS basiertes Processor Power Management (PPM)

– Offers considerable power savings

– Negligible impact to server performance, responsiveness

– Capable processors are prevalent in the market today

– Mature, reliable technology

   • Significant deployments in mobile and desktop systems

Requirements

– Hardware must support PPM capabilities

– ACPI namespace must described capabilities and contain processor objects

– Windows processor driver required for specific CPU make/model

# Operating System: Vergleich

**The evaluation shows that Windows Server 2008 requires approximately 10% less**

– Multiprocessor support
The P-State and C-State controls of Windows Server 2003 support only single processors, whereas Windows Server 2008 supports multiple processors and can individually control processors and cores.

– Improved P-State and C-State Controls
The algorithms of P-State and C-State control have been refined, and power management in Windows Server 2008 improved over that of Windows Server 2003.

| Power options | | Settings |
|---|---|---|
| Windows Server 2003 | Windows Server 2008 | |
| Always On (Default) | High Performance | Always sets P-State to "P0" and demands high performance. |
| Server Balanced Processor Power and Performance | Balanced (Default) | Sets P-State appropriately and balances performance and power consumption. |
| - | Power Saver | Always sets P-State to "Pn". Power consumption is reduced but performance decreases. |

# Processor Power Management in Windows 2008

– Power policy will always use "Demand Based Switching" (DBS) between the range defined by max, min frequency

- Full range of available states, or
- A subset of available states
- Will not include linear clock throttle states

– Policy may be set to use only one performance state

- Max, min, or any intermediate state



| State | Freq | % | Type |
|---|---|---|---|
| 0 | 2800 | 100 | Performance |
| 1 | 2520 | 90 | Performance |
| 2 | 2380 | 85 | Performance |
| 3 | 2100 | 75 | Performance |
| 4 | 1680 | 60 | Performance |
| 5 | 1400 | 50 | Performance |
| 6 | 1400 | 50 | Throttle |
| 7 | 1120 | 40 | Throttle |
| 8 | 840 | 30 | Throttle |
| 9 | 560 | 20 | Throttle |

DBS Allowed

No DBS Allowed

# Die Entwicklung spricht für sich

– Low-power options: CPU, DIMMs, drives
– Scalable performance
– Efficient component selection and design



DL380 Idle and Max Power

# HP ENERGY STAR Servers



- – HP is *first* to publicly announce ENERGY STAR servers

- – Thermal Logic technology enables ENERGY STAR qualification

- – ENERGY STAR DL360 G6 and 380 G6 configurations are currently available

- – More ENERGY STAR servers will be announced soon

**"The EPA is glad to be working with HP, the leading server vendor in the world, to accelerate the reduction of energy consumption in the data center, help businesses reduce operating costs, and encourage companies to be more environmentally responsible."** *-Andrew Fanara, Director, EPA´s ENERGY STAR specifications team.*

©2009

# Links

– HP Power Advisor
[http://h18000.www1.hp.com/products/servers/power-advisor/index.html](http://h18000.www1.hp.com/products/servers/power-advisor/index.html)

– DDR3 Memory Configurations Recommendations
[http://h20195.www2.hp.com/v2/GetPDF.aspx/c01750914.pdf](http://h20195.www2.hp.com/v2/GetPDF.aspx/c01750914.pdf)

– DDR3 Memory Configuration Tool
[http://h18000.www1.hp.com/products/servers/options/tool/hp_memtool.html](http://h18000.www1.hp.com/products/servers/options/tool/hp_memtool.html)

–  Power Regulator for ProLiant servers
[http://h20000.www2.hp.com/bc/docs/support/SupportManual/c00300430/c00300430.pdf](http://h20000.www2.hp.com/bc/docs/support/SupportManual/c00300430/c00300430.pdf)