



18-600 Foundations of Computer Systems

Lecture 25:

"Performance, Power, & Energy of Computers"

John P. Shen

November 29, 2017



*The "Science" of
Computer Systems*

➤ Recommended References:

- "Energy per Instruction Trends in Intel® Microprocessors," by Ed Grochowski, Murali Annavaram, 2006.
- "Mitigating Amdahl's Law through EPI Throttling," by M. Annavaram, E. Grochowski, J. Shen. In 32nd ISCA 2005.



Electrical & Computer
ENGINEERING



18-600 Foundations of Computer Systems

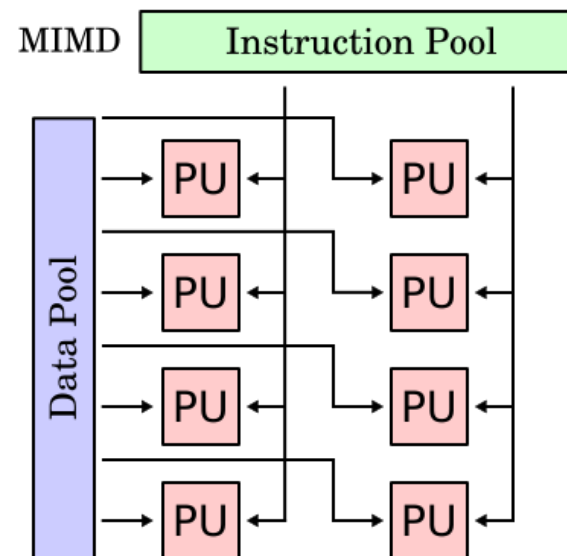
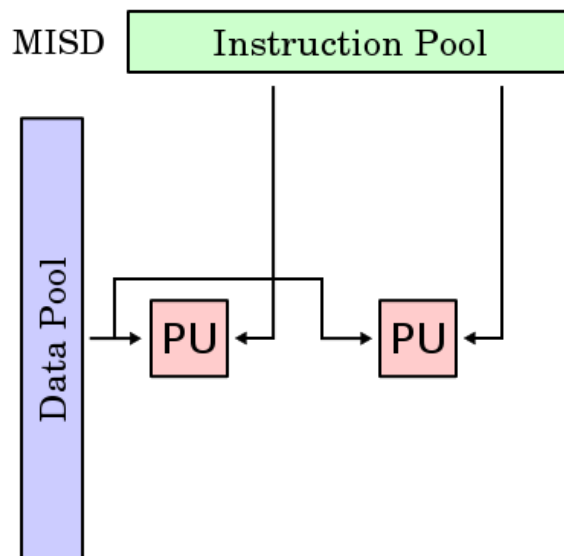
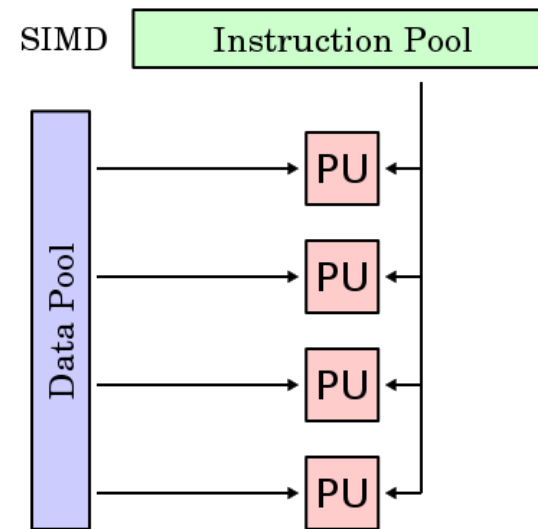
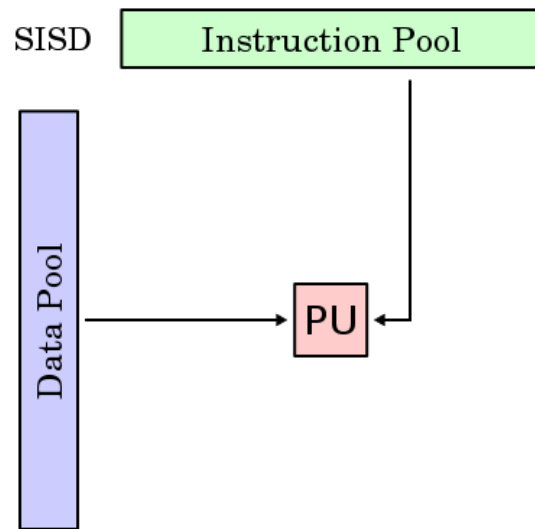
Lecture 25:

“Performance, Power, & Energy of Computers”

- A. Latency** (Single-Thread) Performance (Law #1)
- B. Throughput** (Multi-Thread) Performance (Law #2)
- C. Throughput Performance Scalability** (Law #3)
- D. Performance Scalability Impediments** (Law #4)
- E. Performance and Power** Scaling (Law #5)
- F. Power and Energy** Optimizations



Flynn's Taxonomy of Computer Systems [Mike Flynn, 1966]



SISD – Single Instruction Stream & Single Data Stream. (Sequential uni-processor.)

SIMD – Single Instruction Stream & Multiple Data Streams. (Vector processing, lockstep.)

MISD – Multiple Instruction Streams & Single Data Stream. (Stream data through multiple processing stages.)

MIMD – Multiple Instruction Streams & Multiple Data Streams. (Multi-threads & multi-processors, most general parallelism.)

- **SPMD** – Single Program, Multiple Data Streams. (GPUs, not lockstep.)
- **MPMD** – Multiple Programs, Multiple Data Streams.

Classification of Parallelism

- **SISD** – Traditional sequential program on single-core processor
 - Sequential code with sequential execution semantics (single PC)
 - Can support **concurrent** “multi-processing” through time sharing
 - Control-flow graph & data-flow graph embedded in sequential code
 - Achieve **Instruction Level Parallelism (ILP)** via aggressive control flow speculation and dataflow limit processing
- **MIMD** – Multi-threads & multi-cores, most general type of parallelism
 - Can support simultaneous **parallel** “multi-processing” (multiple PCs)
 - Simultaneous traversing of multiple control-flow graphs
 - Can support both “multi-processing” and “multi-threading”
 - Achieve **Thread Level Parallelism (TLP)** via program & machine parallelisms

A. Latency (Single-Thread) Performance (Law #1)

❖ Time to execute a program: T (latency)

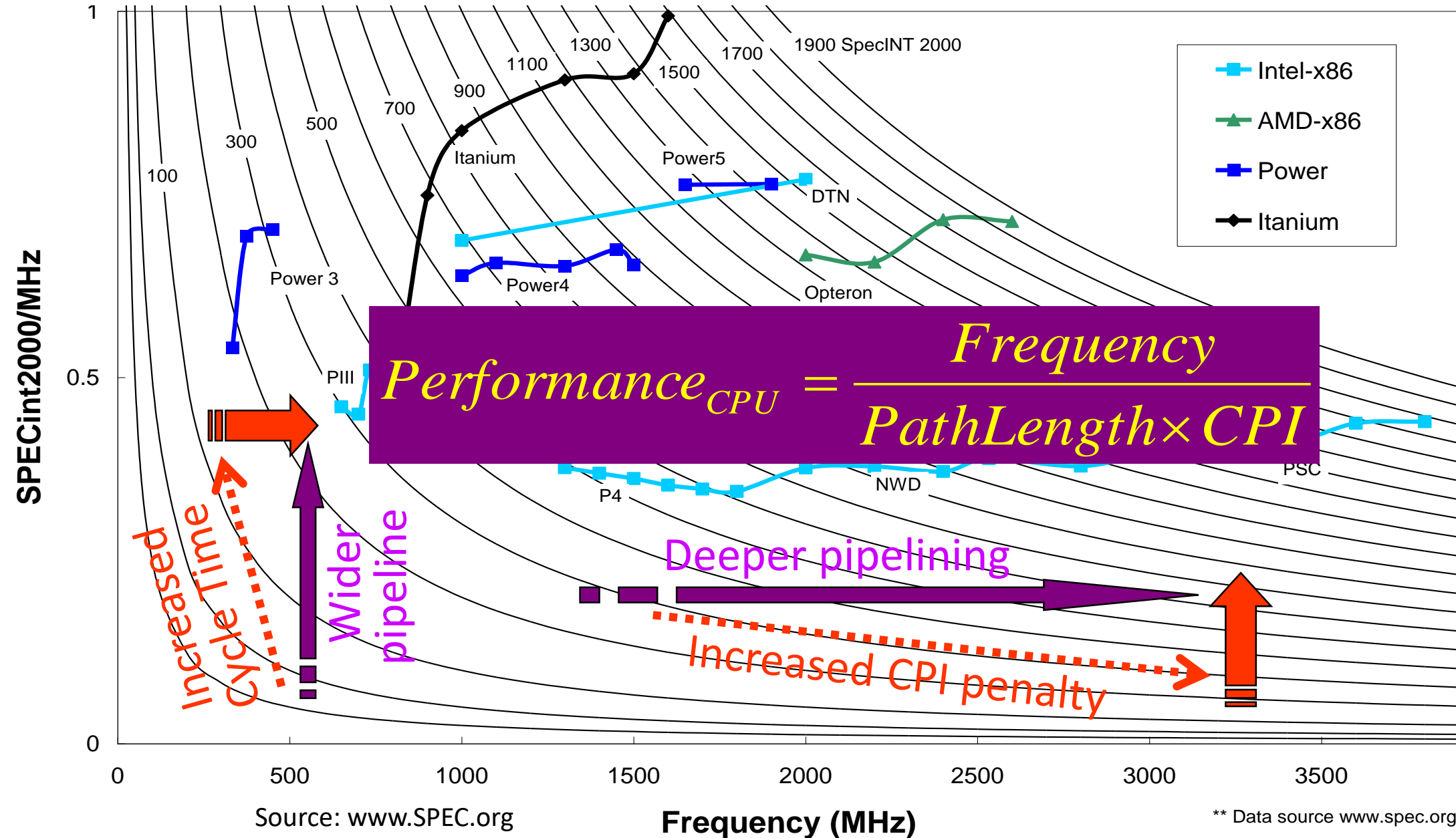
$$T = \frac{\text{instructions}}{\text{program}} \times \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{time}}{\text{cycle}}$$

$$T = \text{PathLength} \times \text{CPI} \times \text{CycleTime}$$

❖ Processor performance: $\text{Perf} = 1/T$

$$\text{Perf}_{\text{CPU}} = \frac{1}{\text{PathLength} \times \text{CPI} \times \text{CycleTime}} = \frac{\text{Frequency}}{\text{PathLength} \times \text{CPI}}$$

SPECint Latency Performance of Processors



Latency vs. Throughput Performance

❖ Reduce Latency of Application

- Uni-processor, Single Program
- Target Single-Thread (ST) Performance
- Examples: SPEC, PC and Workstations

❖ Increase Throughput of System

- Multi-cores/Multi-processors, Many Threads
- Target Multi-threaded/Multi-tasking Throughput
- Example: Database Transaction Processing

Ideal Throughput (Multi-Thread) Performance

- ❖ Time to process a thread: T (latency)

$$T = \frac{\text{instructions}}{\text{thread}} \times \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{time}}{\text{cycle}}$$

$$T = \text{PathLength} \times \text{CPI} \times \text{CycleTime}$$

- ❖ Multi-thread performance: $\text{Perf} = 1/T$

$$\text{Perf}_{MP} = \frac{n \times 1}{\text{PathLength} \times \text{CPI} \times \text{CycleTime}} = \frac{n \times \text{Frequency}}{\text{PathLength} \times \text{CPI}}$$

B. Throughput (Multi-Thread) Performance (Law #2)

❖ Multi-Core/Multi-Thread Performance:

$$Perf_{MC} = \frac{n \times Frequency}{PL(n) \times CPI(n)}$$

❖ Can Improve $Perf_{MC}$ by:

- Increasing: n (no. of CPUs or cores)
- Increasing: $Frequency$ (CPU clock frequency)
- Decreasing: PL (dynamic instruction count)
- Decreasing: CPI (cycles/instruction)

A Throughput Performance Scalability Model

- ❖ Multi-Core/Multi-Thread Speedup:

$$Speedup(n)_{MC} = \frac{\frac{n \times Frequency}{PL(n) \times CPI(n)}}{\frac{Frequency}{PL(1) \times CPI(1)}} = n \times \left(\frac{PL_1 \times CPI_1}{PL(n) \times CPI(n)} \right)$$

- ❖ A Rigorous Scalability Model (my proposal):

$$PL(n) \times CPI(n) \cong n^x PL_1 \times n^y CPI_1$$

C. Throughput Performance Scalability (Law #3)

- ❖ Multi-Core/Multi-Thread Speedup:

$$Speedup(n)_{MC} = n \times \left(\frac{PL_1 \times CPI_1}{n^x PL_1 \times n^y CPI_1} \right) = \left(\frac{n}{n^x \times n^y} \right)$$

- ❖ Scalability Impedance Functions:

$$\begin{cases} n^x \equiv PL(n) \text{ Impedance Function} \\ n^y \equiv CPI(n) \text{ Impedance Function} \end{cases}$$

$$\left(n^x \times n^y \right) = n^{(x+y)} \quad \left| \quad \begin{cases} (x+y) = 1.0 \Rightarrow \text{No Speedup} \\ (x+y) = 0.0 \Rightarrow \text{No Impedance} \end{cases}$$

Real World Example: Database Performance (OLTP)

❖ MIMD Database Performance: TPS

$$TPS = \frac{\text{Transactions}}{\text{Second}} = \frac{n \times \text{Frequency}}{IPX(n) \times CPI(n)} \quad [\text{Law \#2}]$$

❖ Can Improve TPS by:

- Increasing: n (no. of CPUs)
- Increasing: Frequency (CPU clock frequency)
- Decreasing: IPX (instructions/transaction) $\implies PL$
- Decreasing: CPI (cycles/instruction)

4-way Multiprocessor Experimental Setup

Component	Intel® Xeon™ System	Intel® Itanium® 2 System
Processors	4-way SMP, 1.6GHz	4-way SMP, 900MHz
Caches	256KB L2, 1MB L3	256KB L2, 3MB L3
Operating System	Red Hat® AS 2.1	Red Hat® AS 2.1
Disks	24 data + 2 log	32 data + 1 log
Main Memory	4GB	16GB
Database	Oracle® 9ir2	Oracle® 10g
OS Large Page Size	4MB	256MB
SGA	3GB	14GB

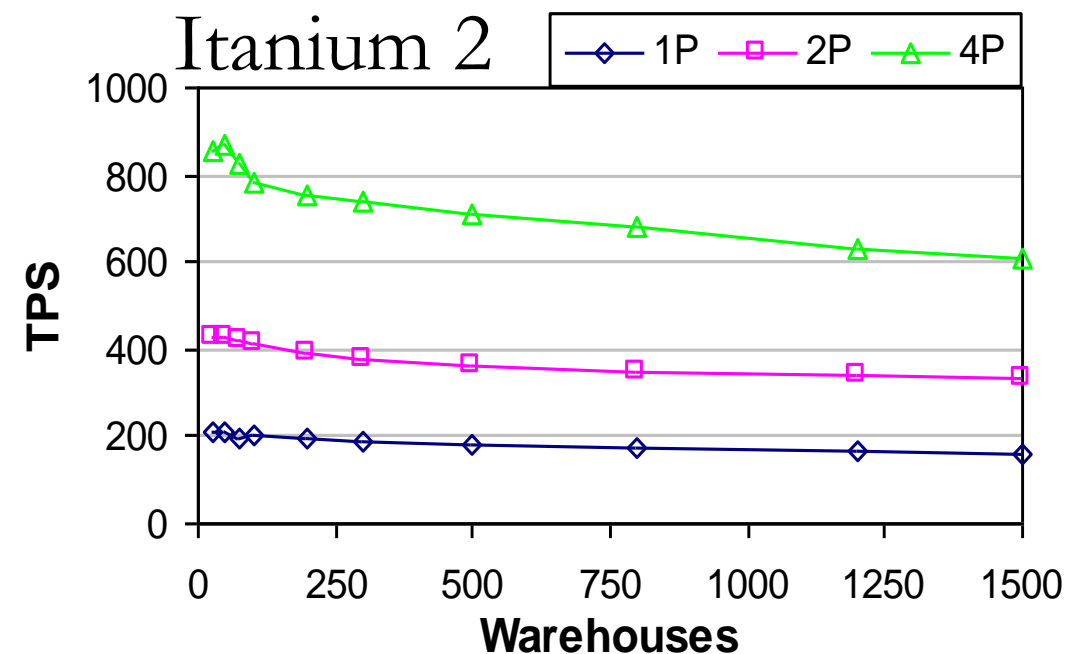
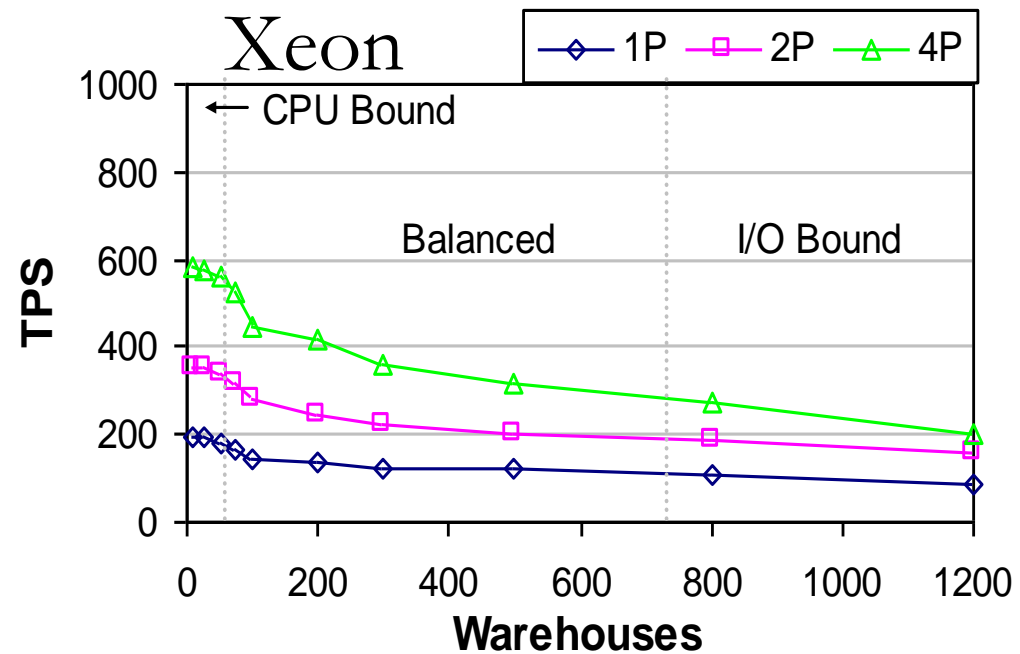
❖ Based on EMON Events

- Separate User and OS components for each event
- Use multiple long runs (20-min warm up, 10-min measurement)
- Strive for standard deviation <5% for TPS & CPU utilization > 90%
- Overall user execution time ~70-90%

$$TPS = \frac{n \times Frequency}{IPX(n) \times CPI(n)}$$

TPS: Throughput Scaling

$$TPS = \frac{n \times \text{Frequency}}{IPX(n) \times CPI(n)}$$

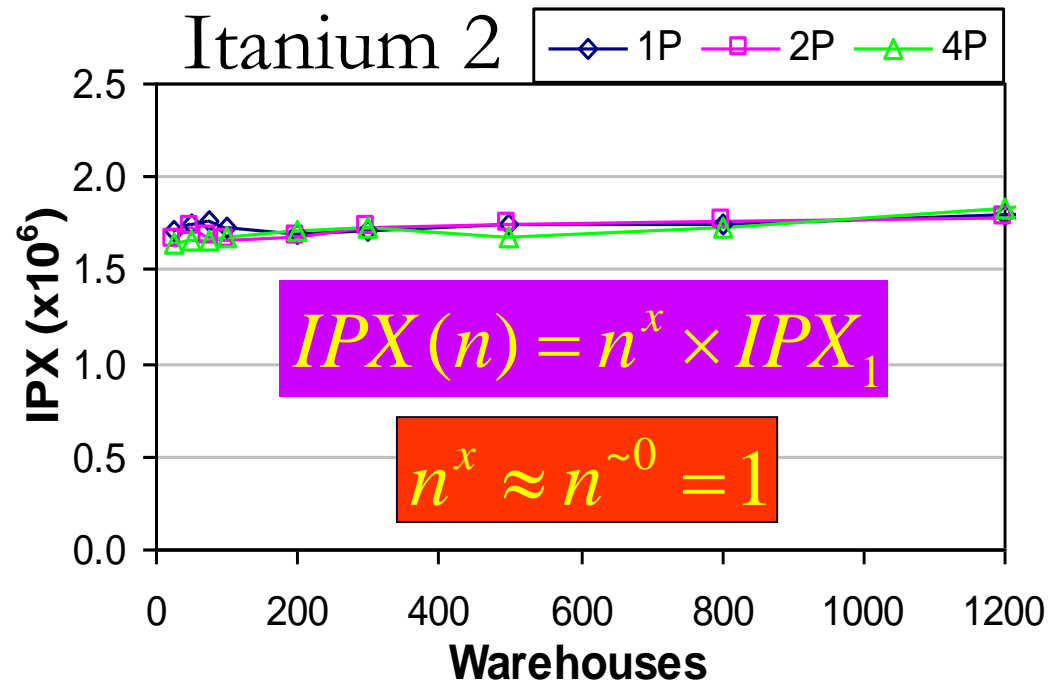
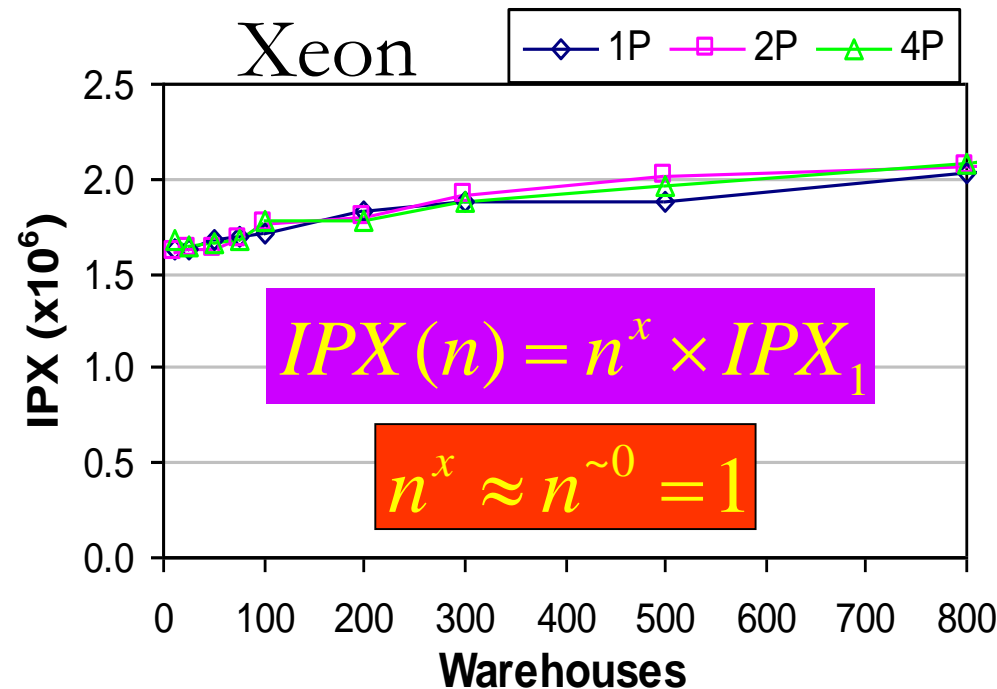


- ❖ **TPS scales more linearly on Itanium2**
 - Larger SGA implies slower I/O rate increase (Xeon I/O rate increases at 7KB/WH, & Itanium 2 at: 6KB/WH)
 - Bus utilization on Xeon higher than Itanium 2 (45% vs. 39%)
- ❖ **Increasing I/O rate → more processes & context switches**
 - Clients increase from 8 to 56 on Xeon & 4 to 64 on Itanium 2
 - OS time up to 20% on Xeon & only 10% on Itanium 2

- Performance degrades with increased data set size (due to I/O rate)
- Performance improves with increased n (on IPF almost linear increase)

IPX: Path-Length Scaling

$$TPS = \frac{n \times \text{Frequency}}{IPX(n) \times CPI(n)}$$

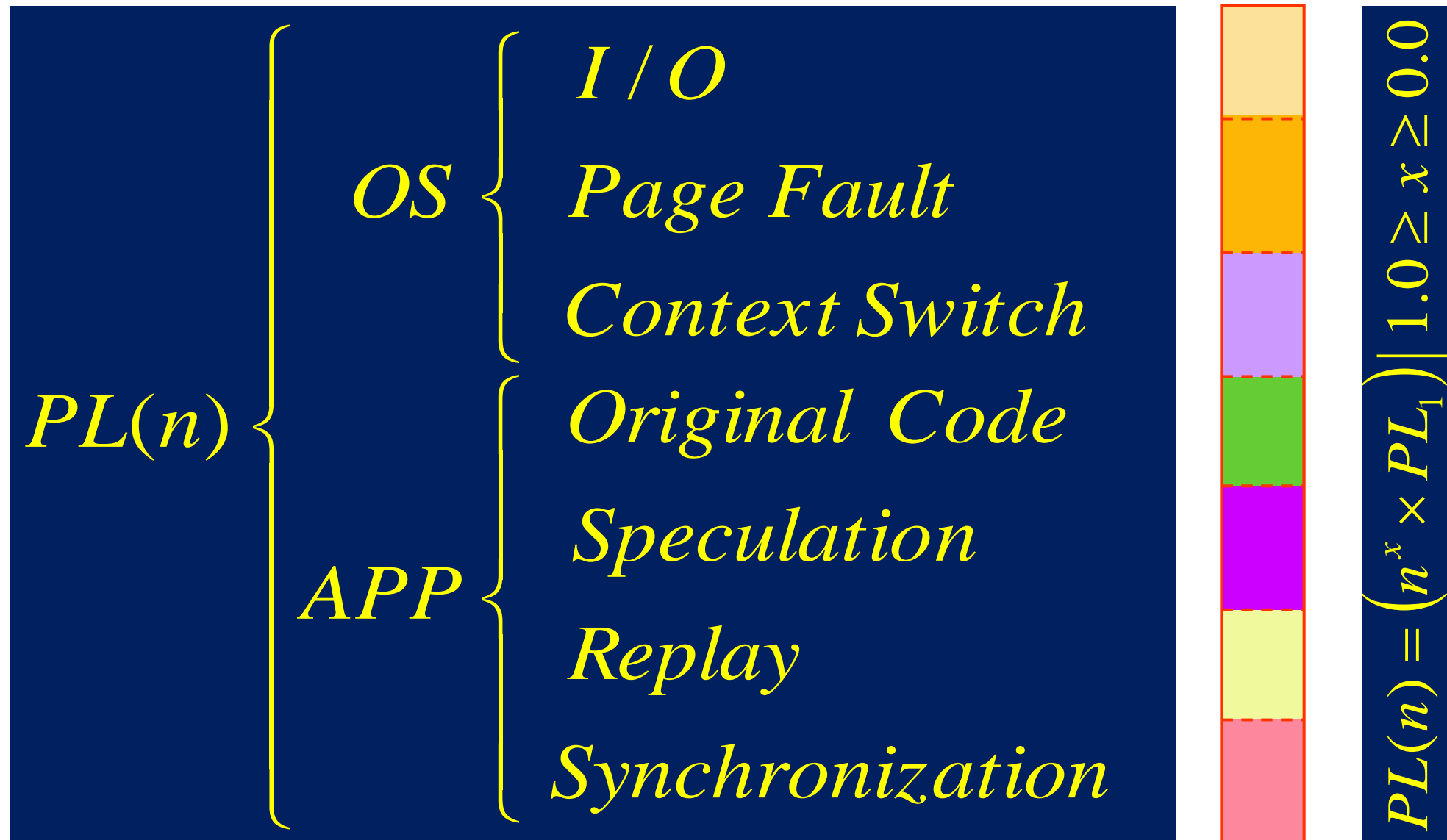


- ❖ Growth in IPX (quite linear) attributed to OS IPX increase
 - More I/O, more context switching
- ❖ User level IPX remains relatively constant in both systems
 - Code path through Oracle relatively constant
- ❖ Excluding NOPS, IPX at 25 WH similar for both systems!
- ❖ IPX growth less pronounced on Itanium 2

- IPX increases with increased data set size (Xeon)
- But IPX(n) doesn't increase much with increased n

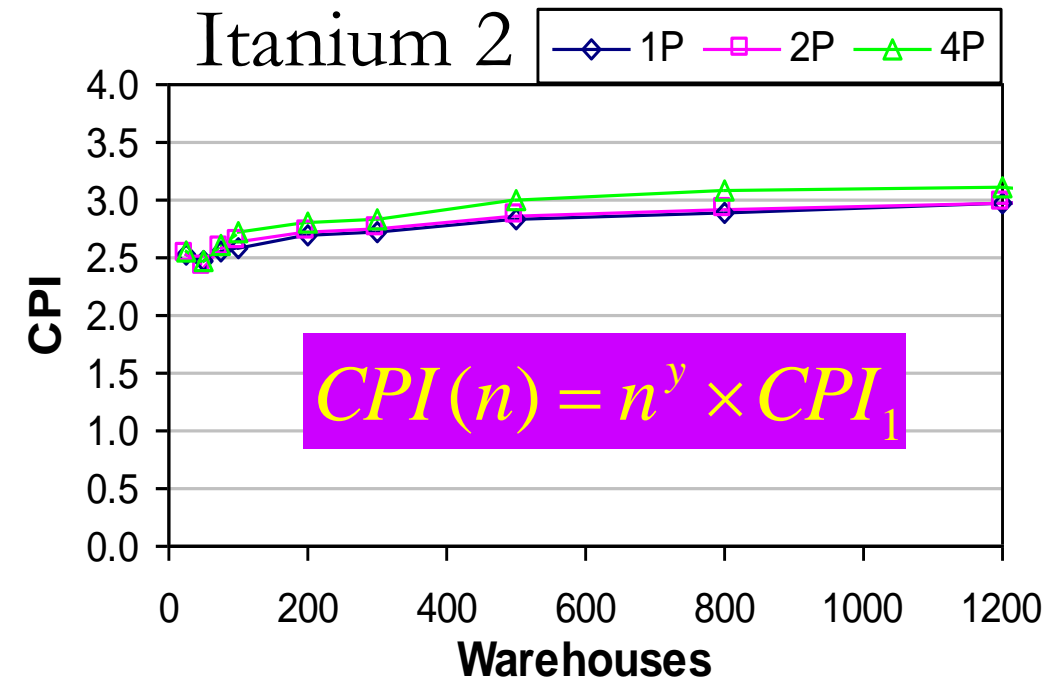
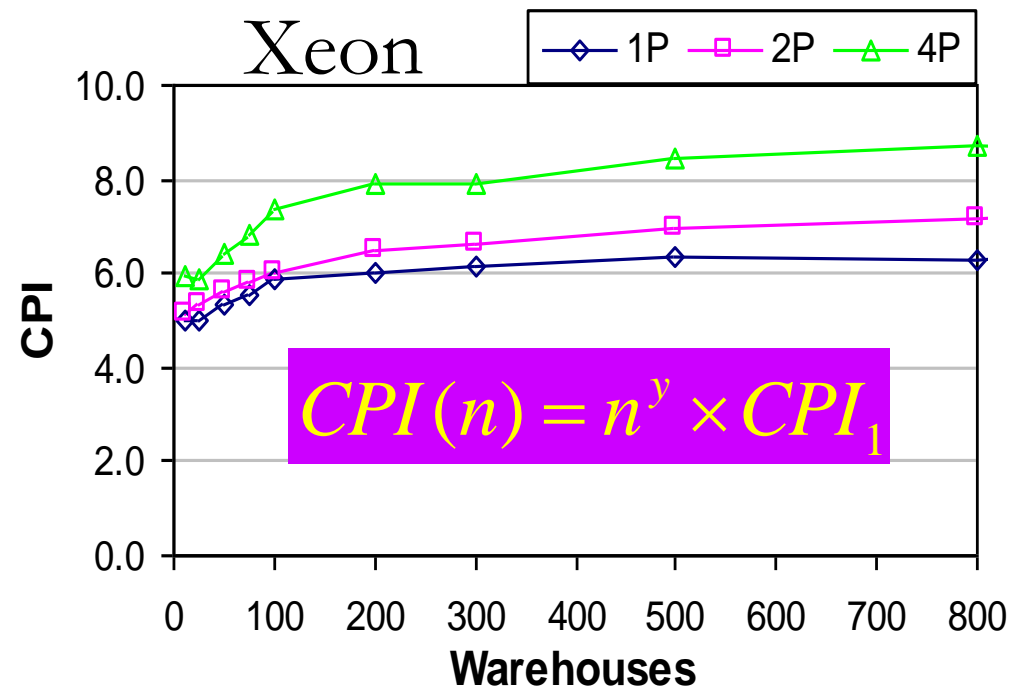
Path-Length Contributions

$$TPS = \frac{n \times \text{Frequency}}{IPX(n) \times CPI(n)}$$



CPI: Overhead Scaling

$$TPS = \frac{n \times \text{Frequency}}{IPX(n) \times CPI(n)}$$

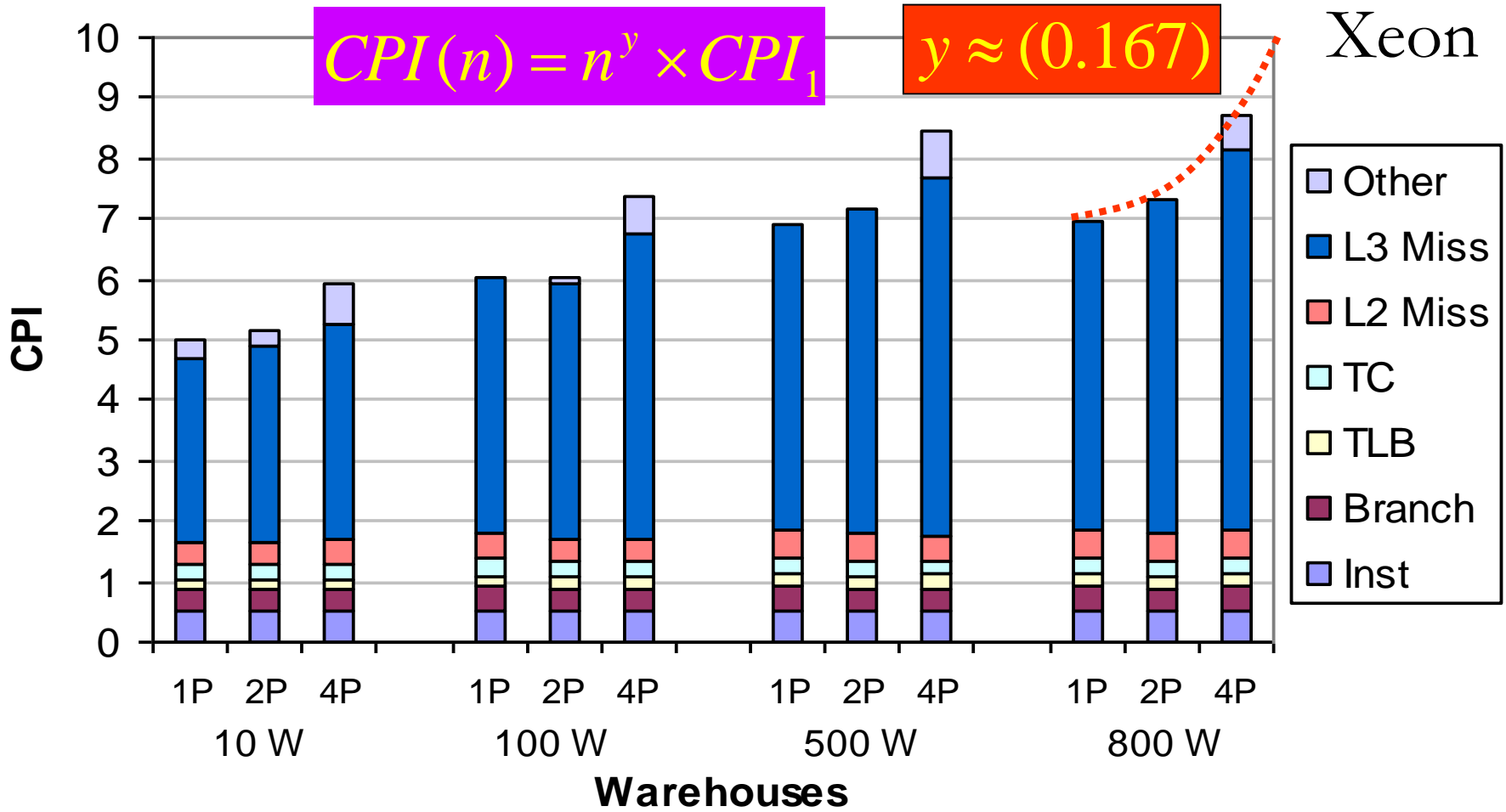


- ❖ CPI increases with a "knee" but less sharp on Itanium 2 !
- ❖ Overall CPI trend strongly determined by user CPI
 - User mode execution time more than 90% on IPF, 80% on Xeon
- ❖ Xeon CPI grows with P, Itanium 2 CPI does not
 - Growth attributed to higher bus utilization on Xeon

- CPI does increase with increased data set size
- CPI(n) also increases with increased n (esp. for Xeon)

CPI Breakdown: Xeon

$$TPS = \frac{n \times \text{Frequency}}{IPX(n) \times CPI(n)}$$



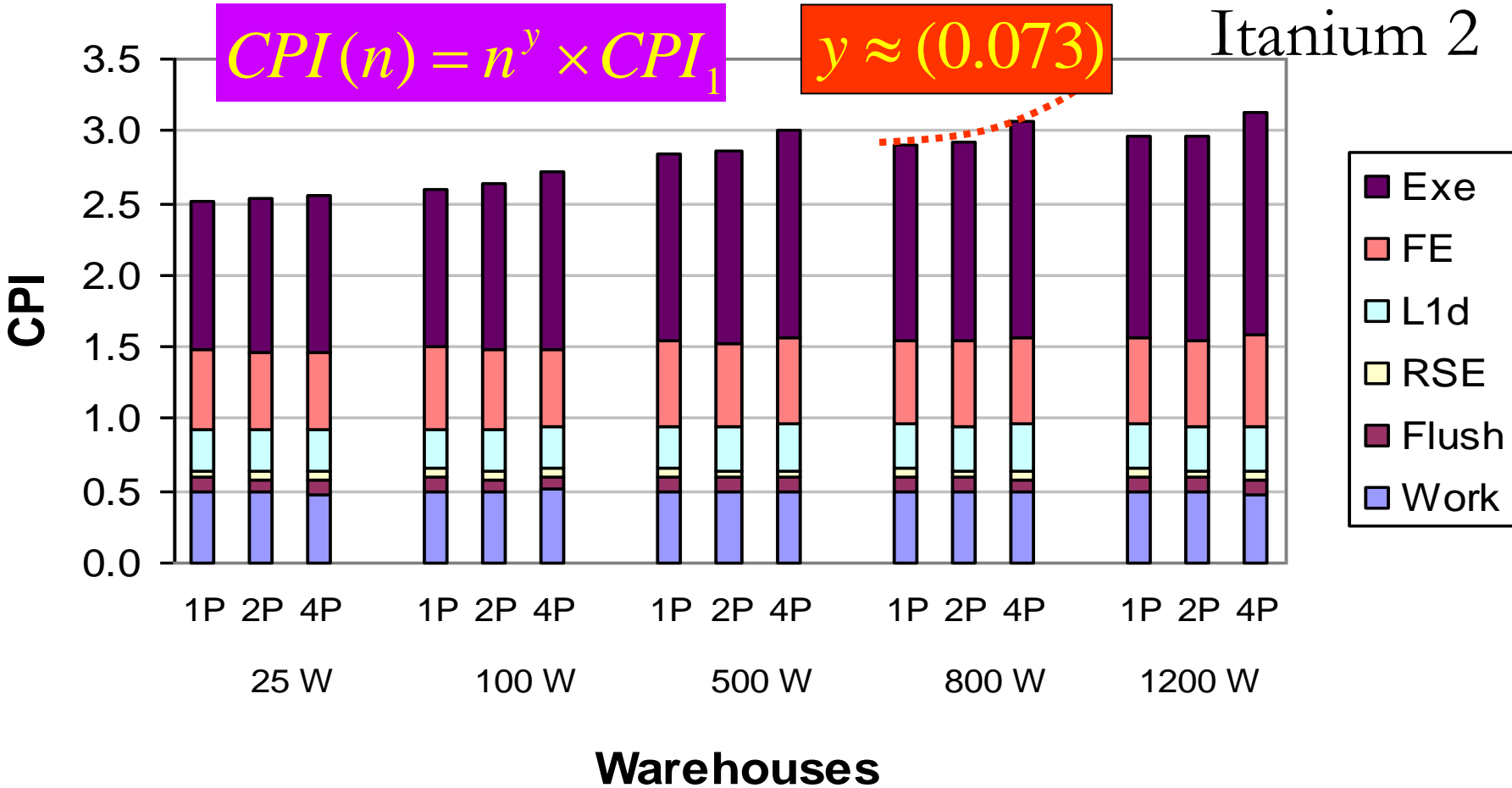
$$TPS = \frac{n \times \text{Frequency}}{(n^0 \times IPX_1) \times (n^{0.167} \times CPI_1)}$$

$$Speedup(n)_{MC} = \left(\frac{n}{n^x \times n^y} \right) = n^{0.833}$$

[Law #3]

CPI Breakdown: Itanium 2

$$TPS = \frac{n \times \text{Frequency}}{IPX(n) \times CPI(n)}$$



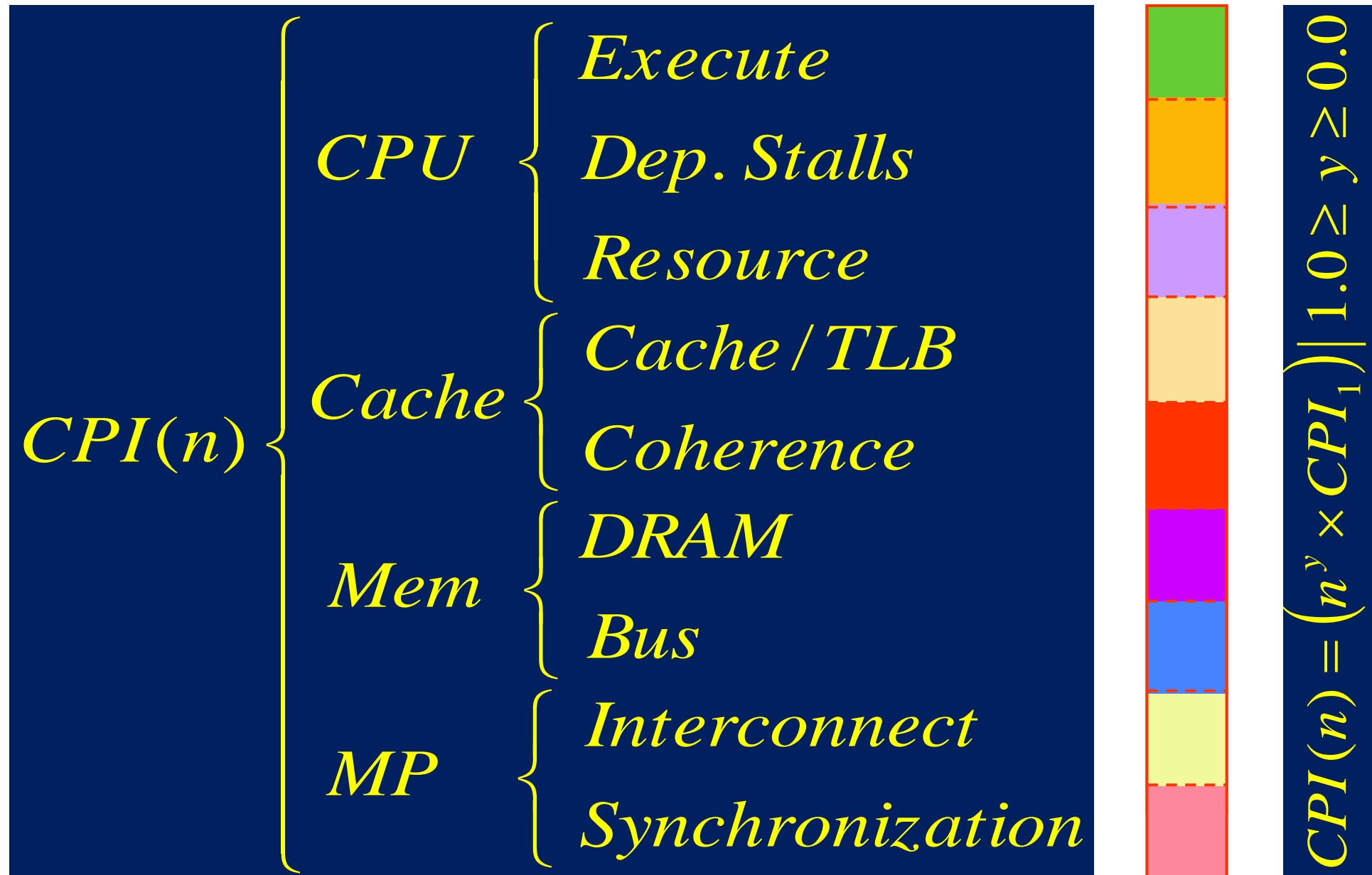
$$TPS = \frac{n \times \text{Frequency}}{(n^0 \times IPX_1) \times (n^{0.073} \times CPI_1)}$$

$$Speedup(n)_{MC} = \left(\frac{n}{n^x \times n^y} \right) = n^{0.927}$$

[Law #3]

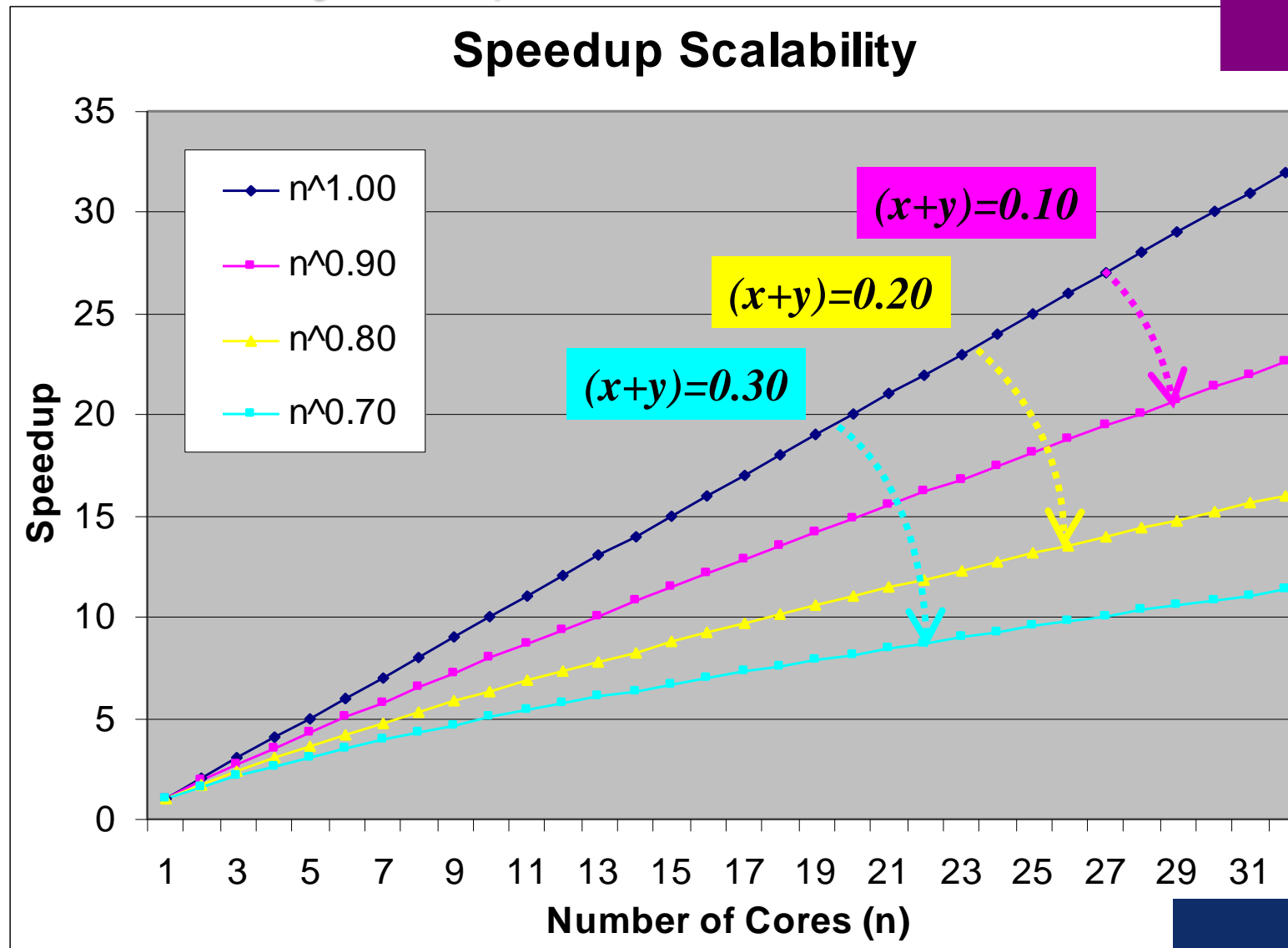
CPI Contributions

$$TPS = \frac{n \times \text{Frequency}}{IPX(n) \times CPI(n)}$$



Scalability Impediments

$$Perf_{MC} = \frac{n \times Frequency}{(n^x \times PL_1) \times (n^y \times CPI_1)}$$



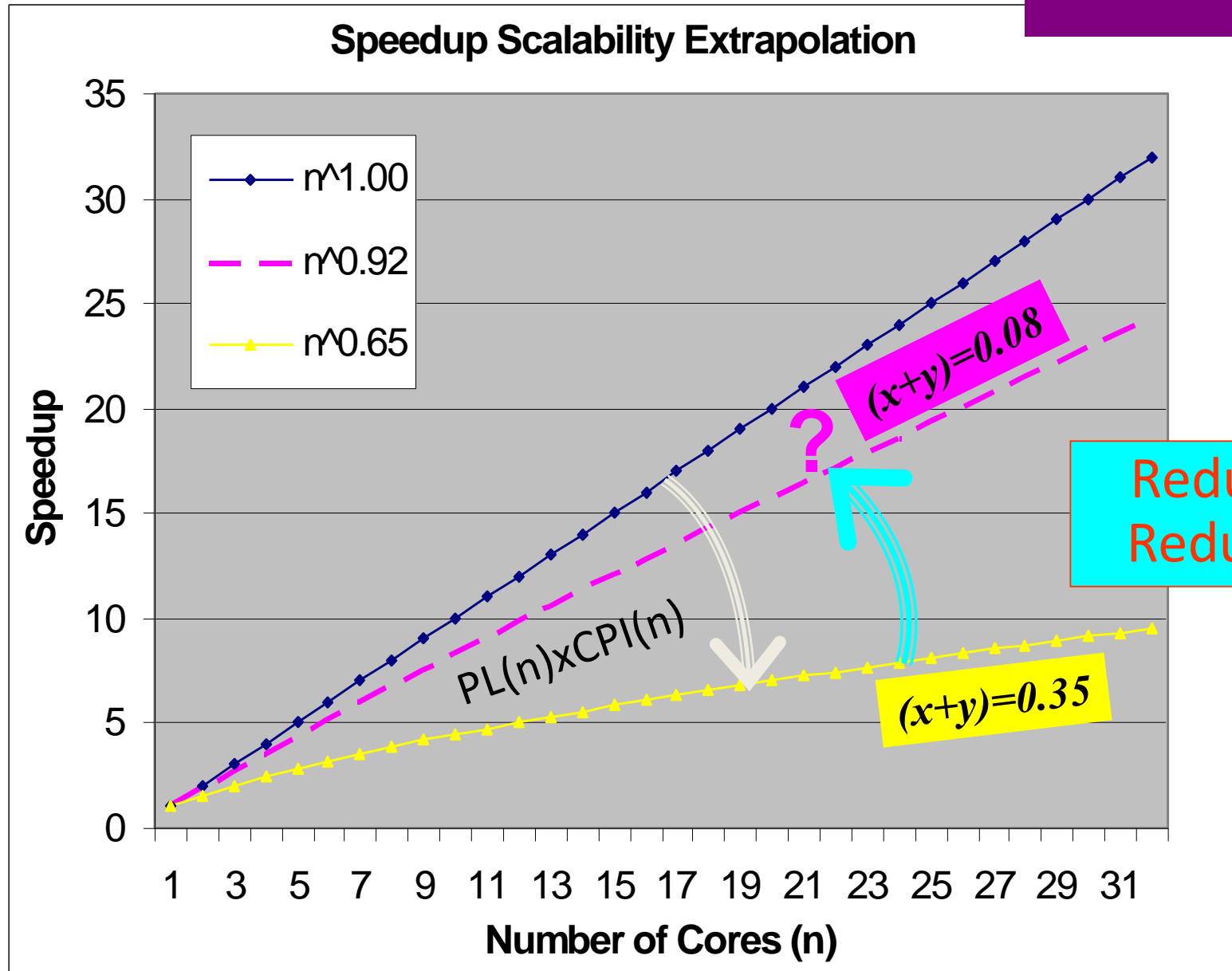
[Carole Dulong et al., 2005]

1p-16p scaling	$1-(x+y)$	R^2
SEMPHY	0.993	0.999
PLSA	0.963	0.999
Rsearch	0.931	0.997
SVM-RFE	0.786	0.970
SNPs	0.685	0.967
GeneNet	0.642	0.983

$$Speedup(n)_{MC} = \left(\frac{n}{n^x \times n^y} \right) = n^{1-(x+y)}$$

Scalability Headroom

$$Perf_{MC} = \frac{n \times Frequency}{(n^x \times PL_1) \times (n^y \times CPI_1)}$$



Reduce PL(n)
Reduce CPI(n)

Conspiring Forces Against Performance Scaling: Three Forms of Scalability Impediments

❖ **Algorithm**

- Limitation of Language and Algorithm
- Tyranny of Amdahl's Law (sequential bottleneck)

❖ **Architecture**

- Increase of Path-Length Undermines Scalability
- Increase of CPI also Undermines Scalability

❖ **Power/Thermal**

- Increased Complexity and Inefficiency of Design
- Super-linear Power Scaling Relative to Performance

D. Performance Scalability Impediments (Law #4)

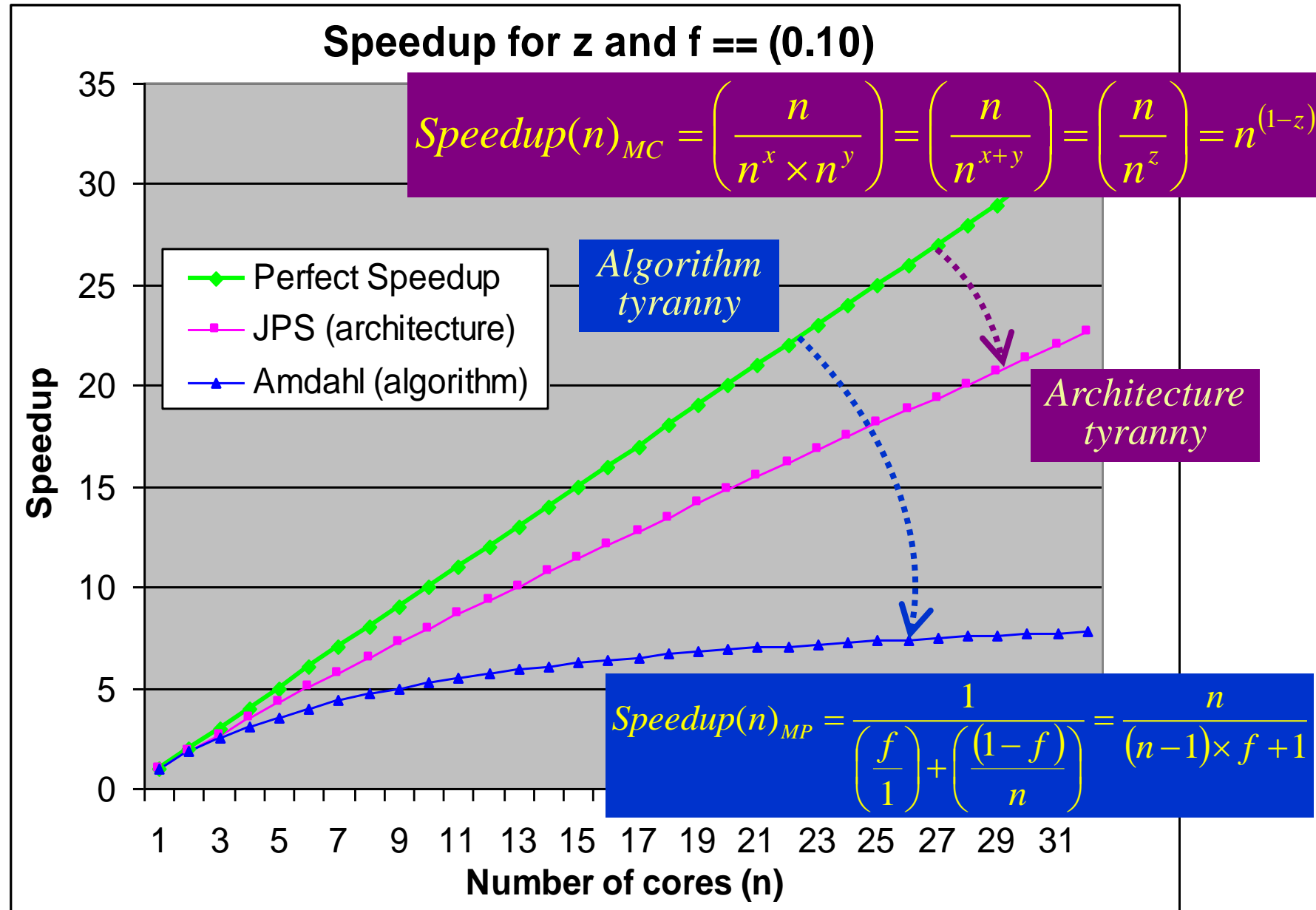
Algorithm Scalability: (Amdahl's Law)

$$Speedup(n)_{ALG} = \frac{1}{\left(\frac{f}{1}\right) + \left(\frac{(1-f)}{n}\right)} = \frac{n}{(n-1) \times f + 1}$$

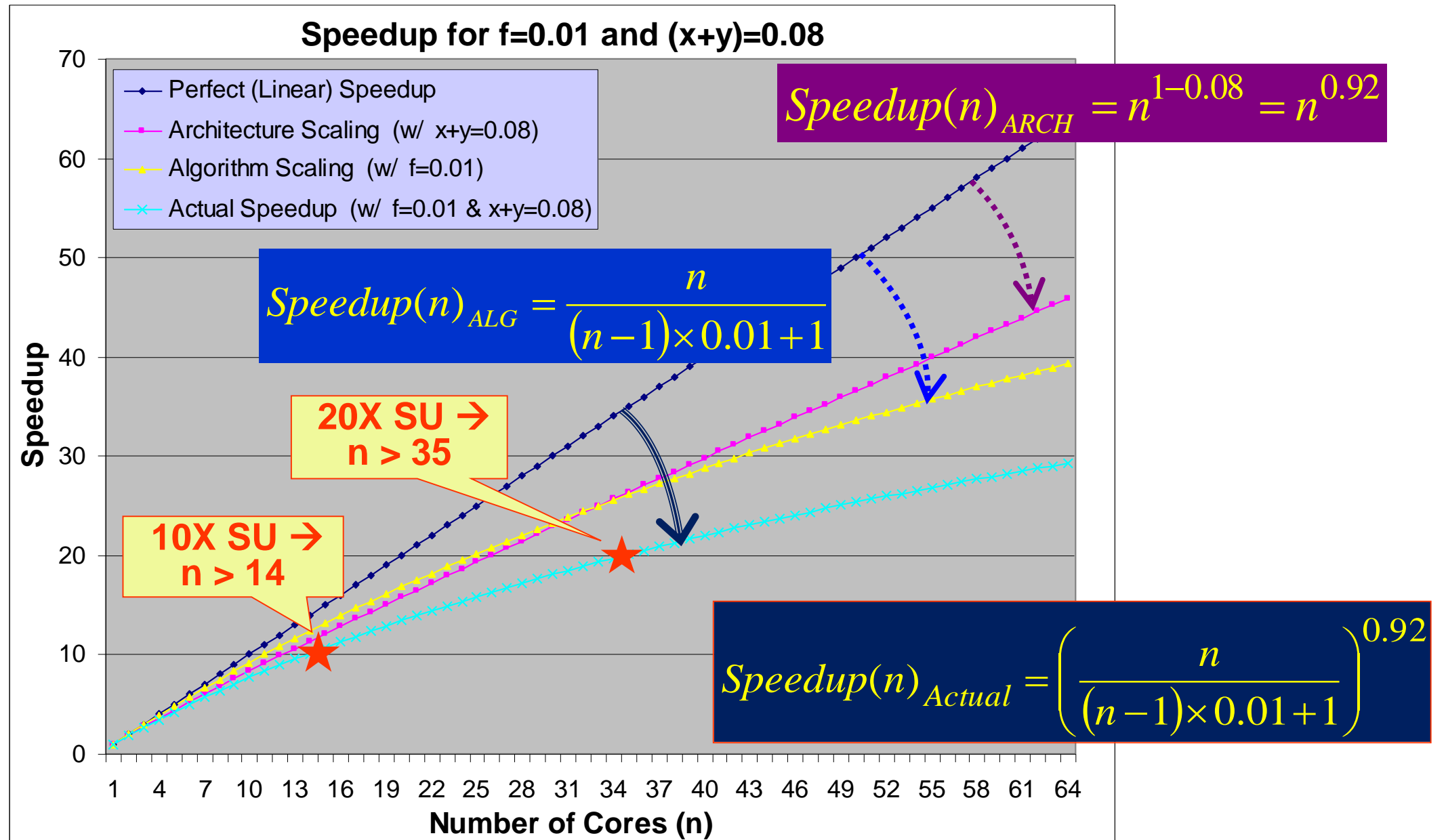
Architecture Scalability:

$$Speedup(n)_{ARCH} = \left(\frac{n}{n^x \times n^y}\right) = \left(\frac{n}{n^{x+y}}\right) = n^{1-(x+y)}$$

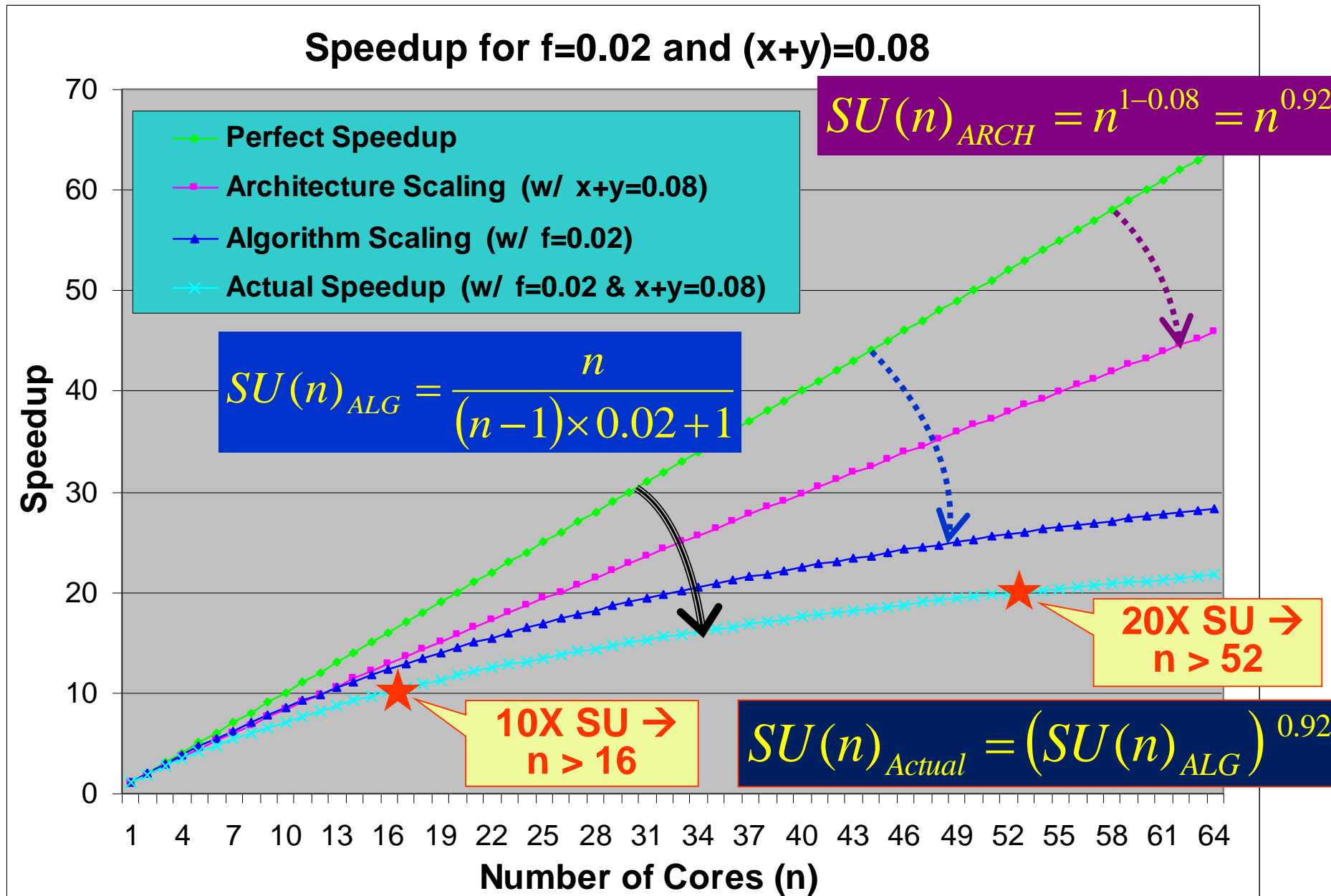
Relative (Additive) Degrees of Tyranny



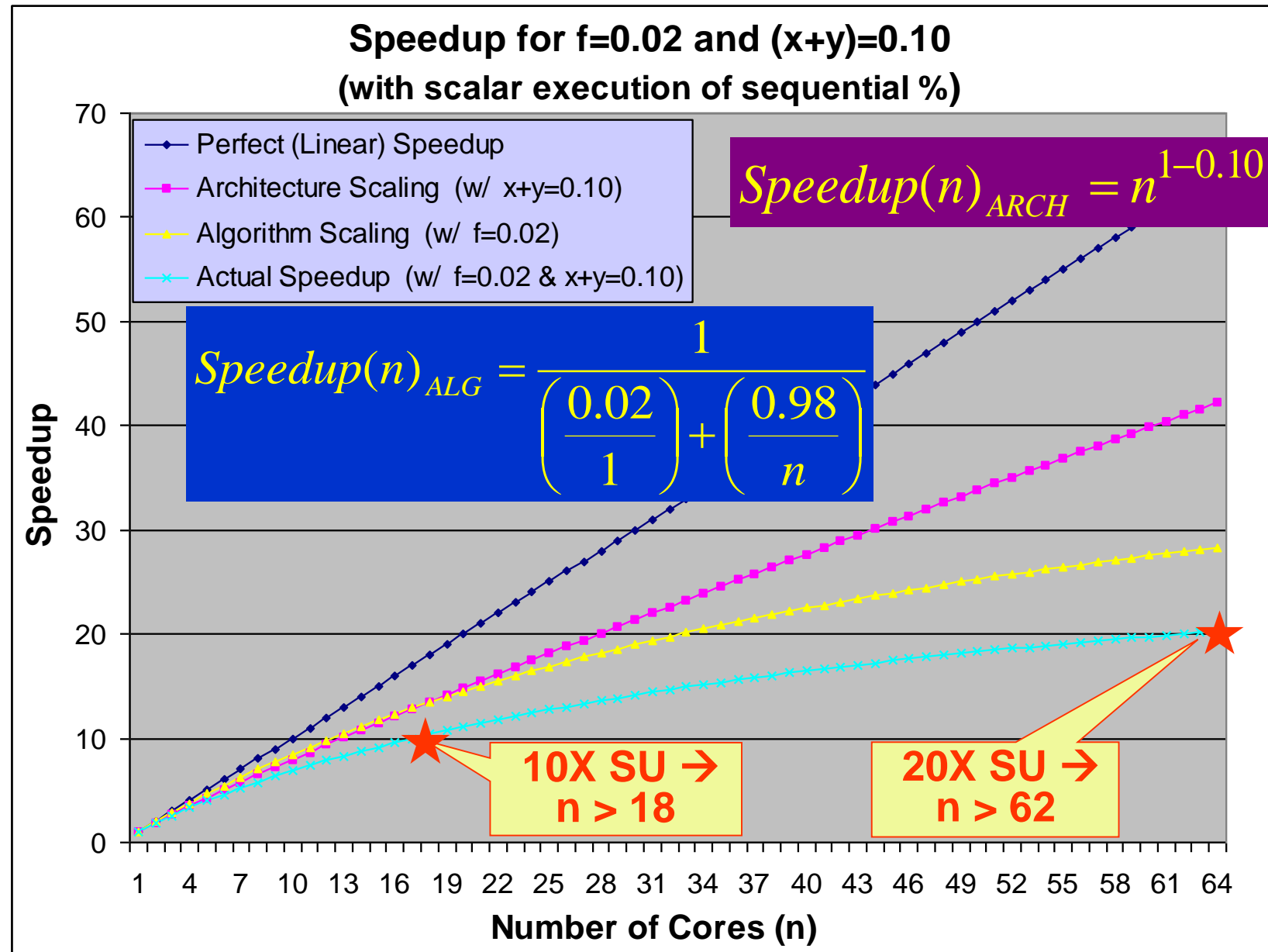
Combined Effect on Actual Speedup (f=0.01, x+y=0.08)



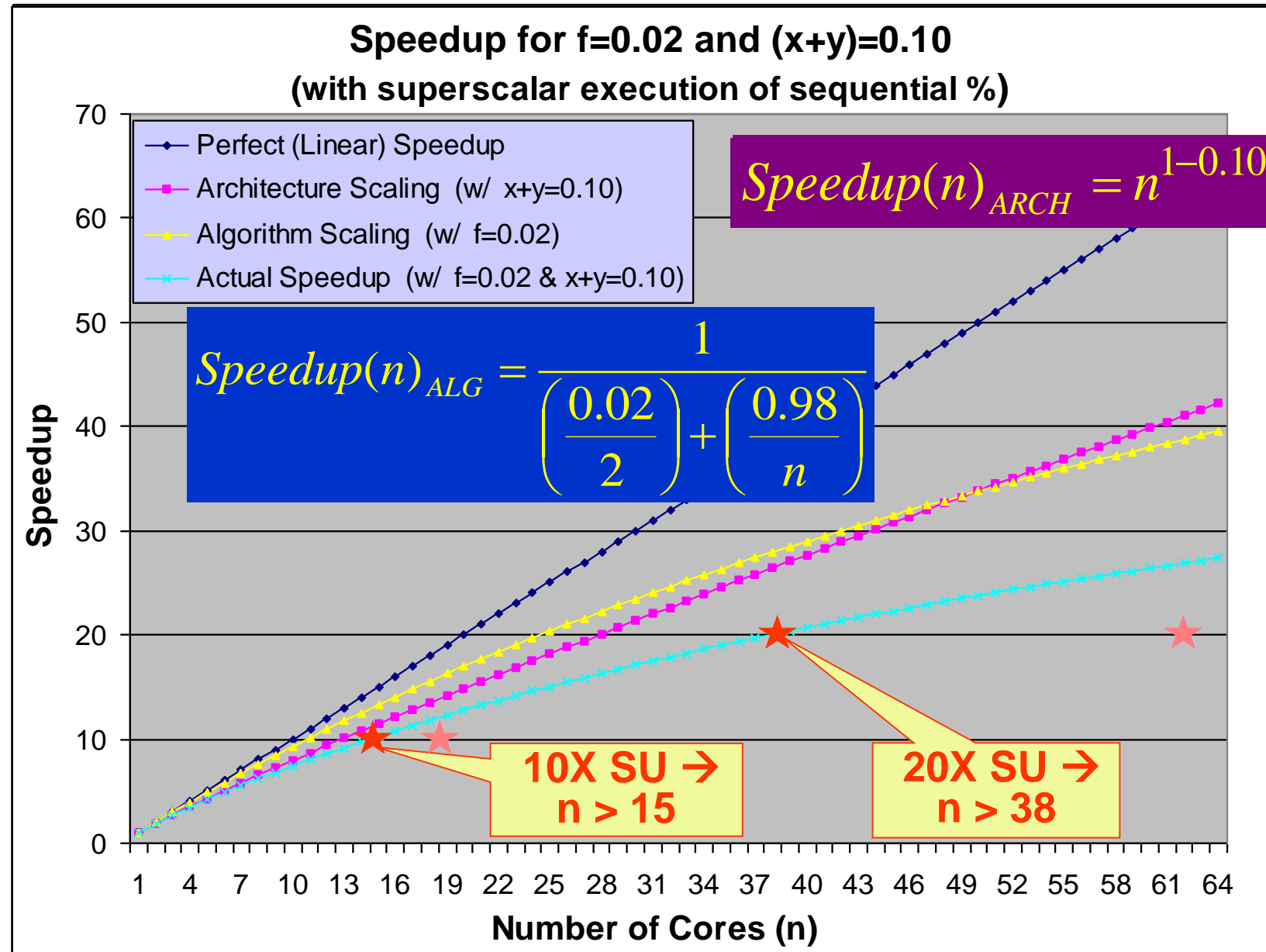
Combined Effect on Actual Speedup ($f=0.02, x+y=0.08$)



Combined Effect on Actual Speedup (f=0.02, x+y=0.10)



Impact of Latency Performance on MP Performance



E. Performance and Power Scaling (Law #5)

$$\text{Watt} = \frac{\text{Joule}}{\text{second}} = \frac{\text{Joule}}{\text{instruction}} \times \frac{\text{instruction}}{\text{cycle}} \times \frac{\text{cycle}}{\text{second}}$$

$$\text{Power} = \text{EPI} \times \text{IPC} \times \text{Frequency}$$

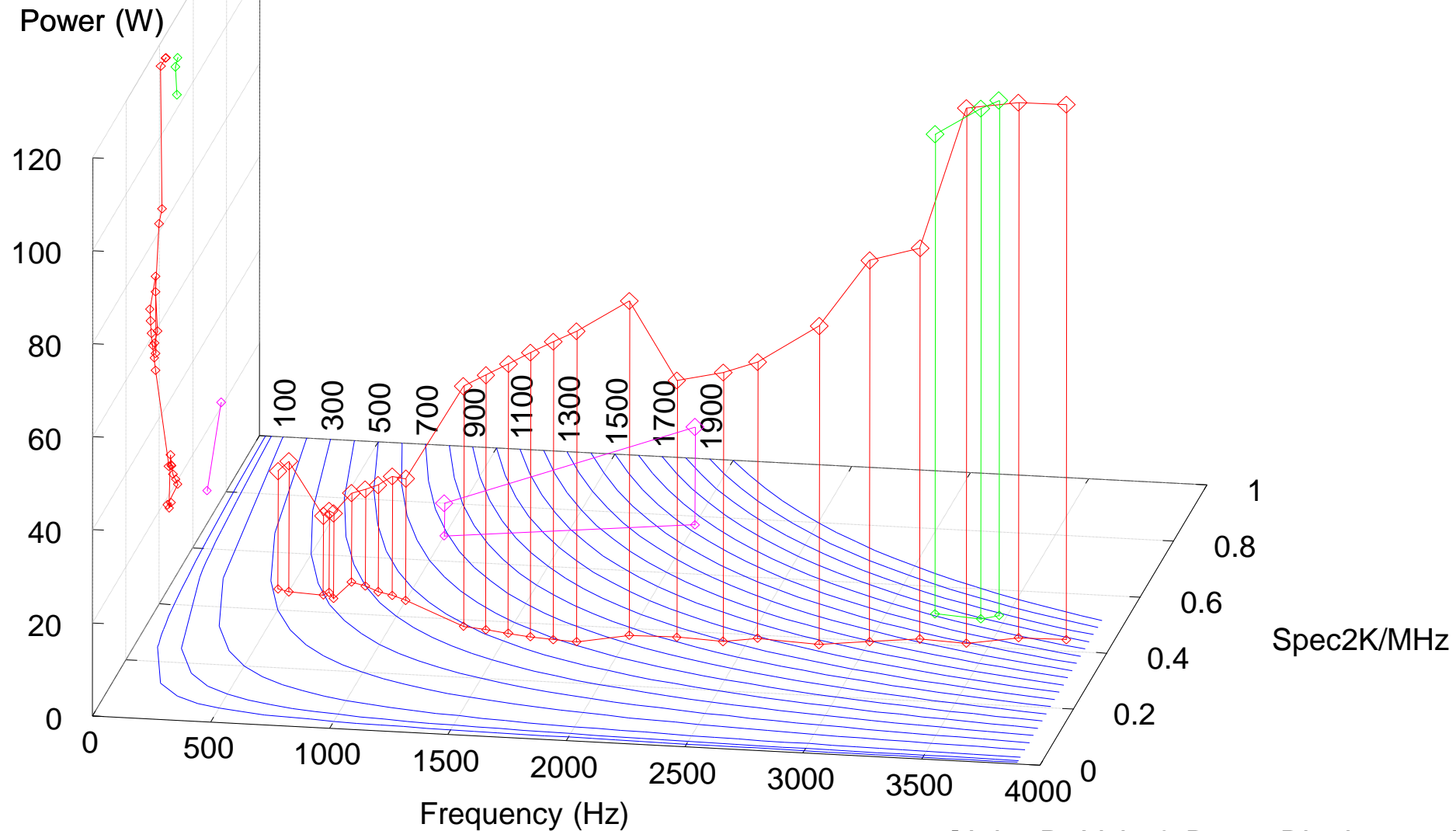
$$\text{Performance} = \frac{\text{Frequency}}{\text{PathLength} \times \text{CPI}} = \frac{\text{IPC} \times \text{Frequency}}{\text{PathLength}}$$

$$\text{Power} = \text{EPI} \times \text{Performance} \times \text{PathLength}$$

Power Scaling Relative to Performance Scaling

$$Power = EPI \times IPC \times Frequency$$

- Pentium —◇—
- Pentium EE —◇—
- Pentium M —◇—



[John DeVale & Bryan Black, 2005]

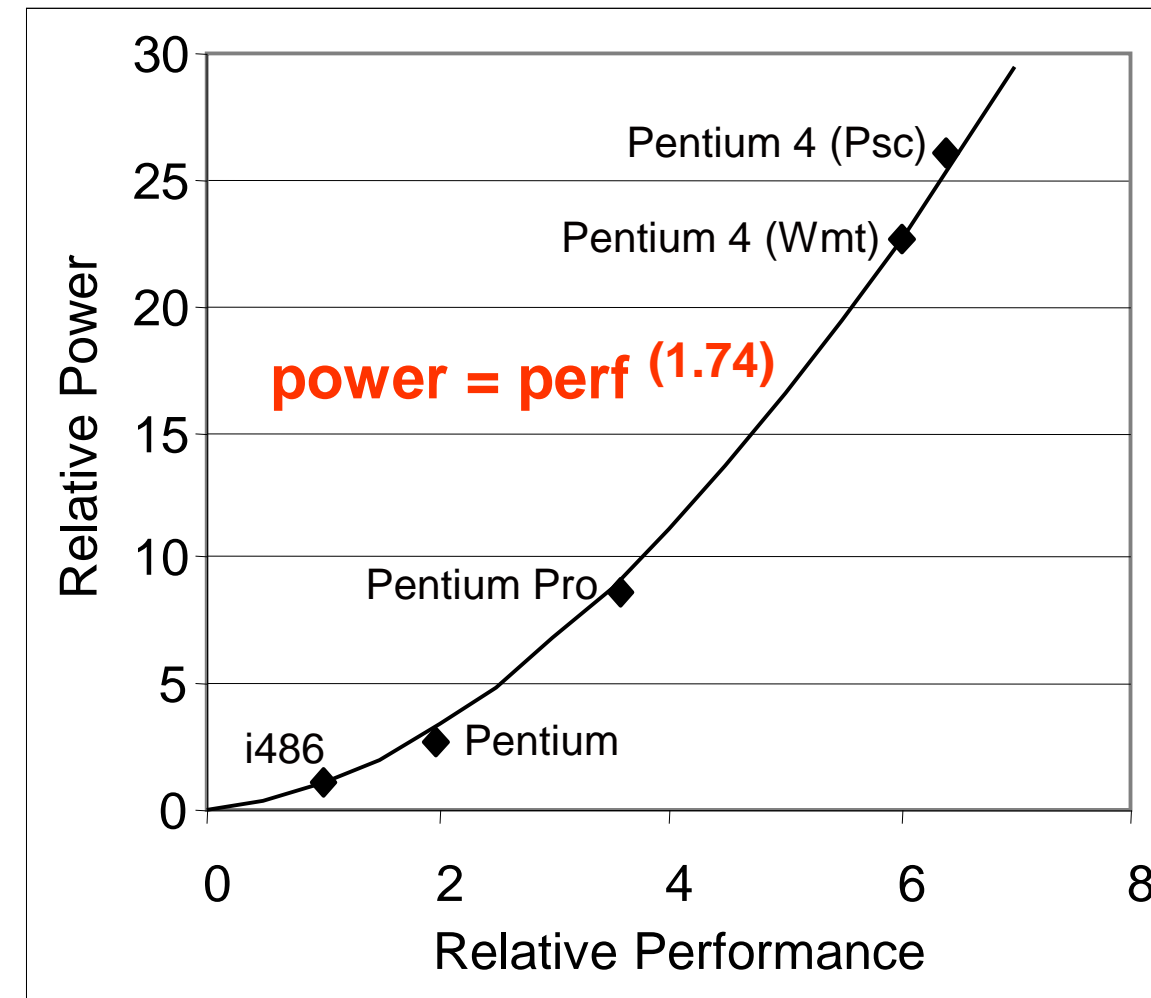
Power vs. Latency Performance

❖ For comparison

- Factor out contributions due to process technology
- Keep contributions due to microarchitecture design
- Normalize to i486™ processor

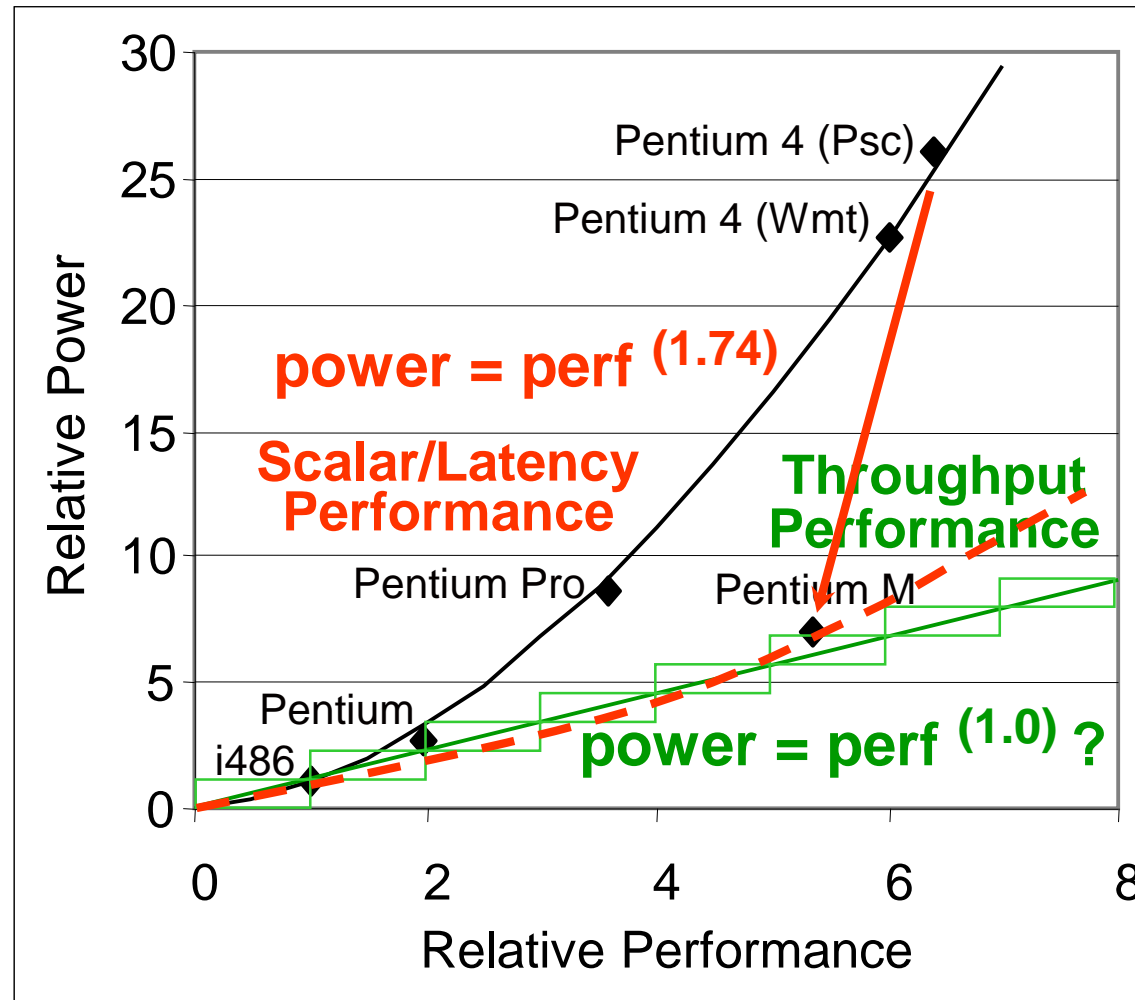
❖ Relative to i486™ Pentium® 4 (Wmt) processor is

- 6x faster (2X IPC at 3X frequency)
- 23x higher power
- Spending 4 units of power for every 1 unit of scalar performance



Power vs. Latency & Throughput Performances

[Ed Grochowski, 2005]



CPU	EPI
I486	7 nj
P5	10 nj
P6	17 nj
P4P-wmt	27 nj
P4P-psc	29 nj
Pentium M	9 nj

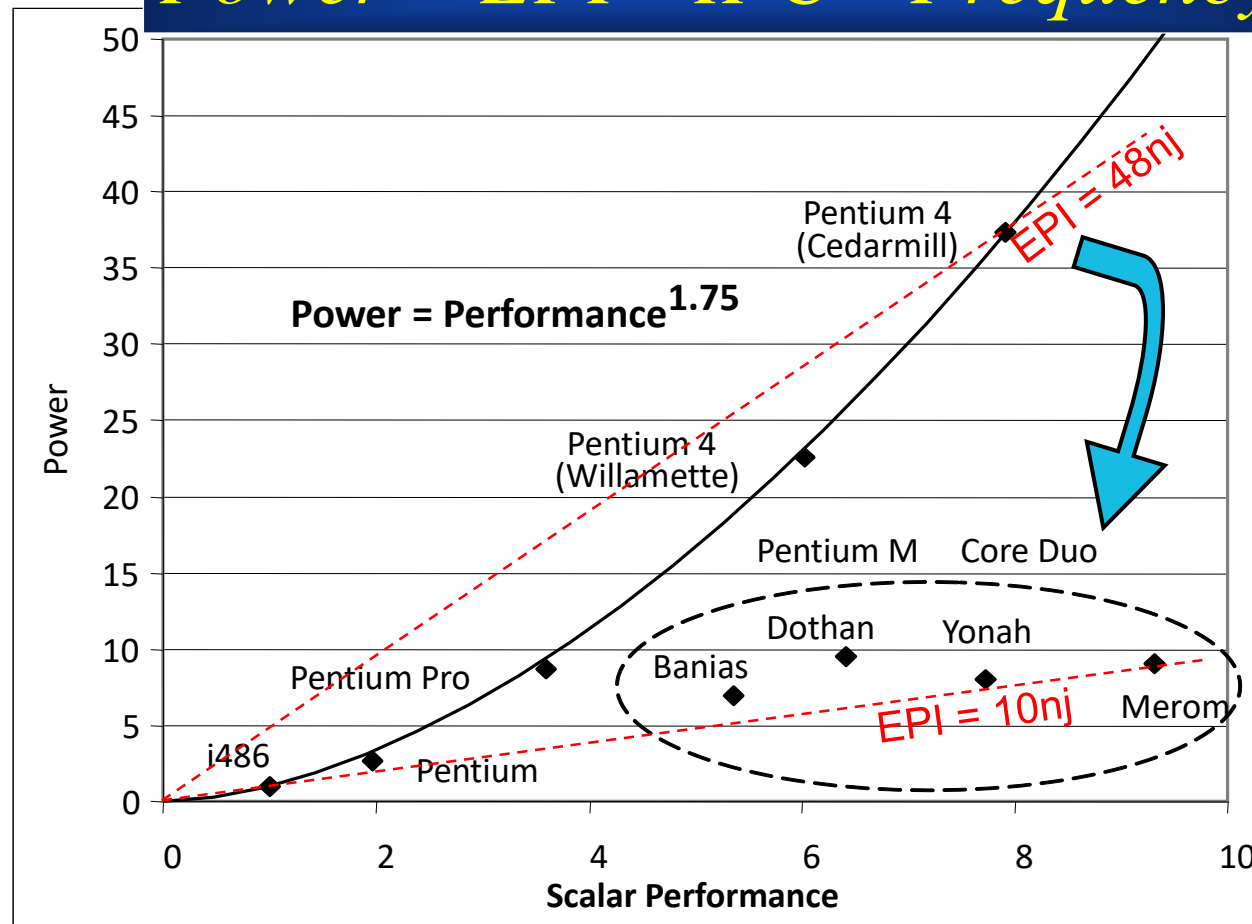
Low EPI

- ❖ Assume a large-scale CMP with potentially many cores
- ❖ Replication of cores results in (nearly) proportional increases to both throughput performance and power (hopefully).

Power/Performance (EPI) Evolution

$$Power = EPI \times IPC \times Frequency$$

[Ed Grochowski, 2004]

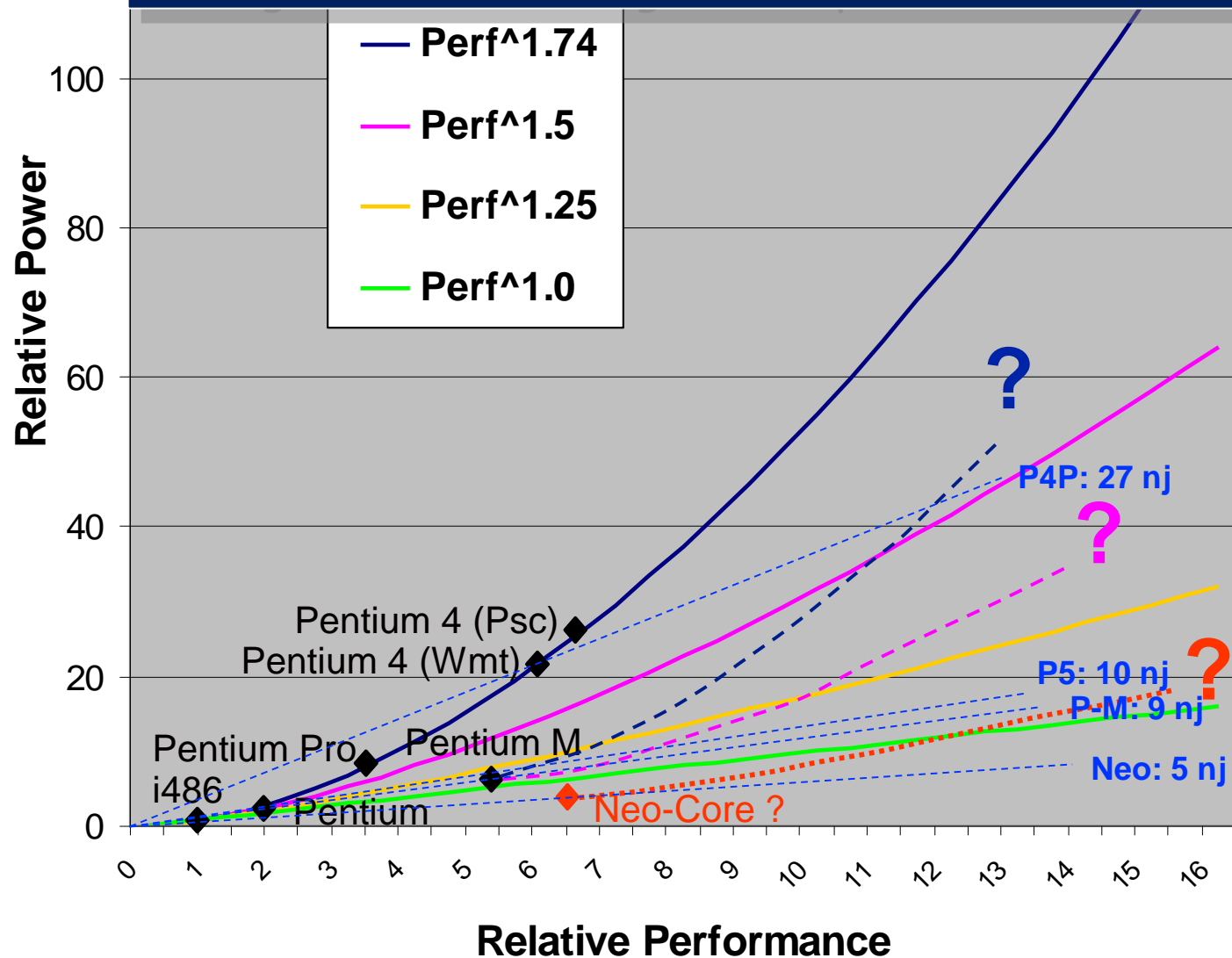


Intel Microprocessors	EPI (nj) 65nm at 1.33v
i486	10
Pentium	14
Pentium Pro	24
Pentium 4 (WMT)	38
Pentium 4 (CDM)	48
Pentium M (Baniyas)	13
Pentium M (Dothan)	15
Core Duo (Yonah)	11
Core Duo (Merom)	10

- ❖ Power: single core power (relative to i486 baseline)
- ❖ Performance: SPECint performance (relative to i486 baseline)
- ❖ EPI: average energy spent per instruction (in nano-joules)

Power/Performance Scaling

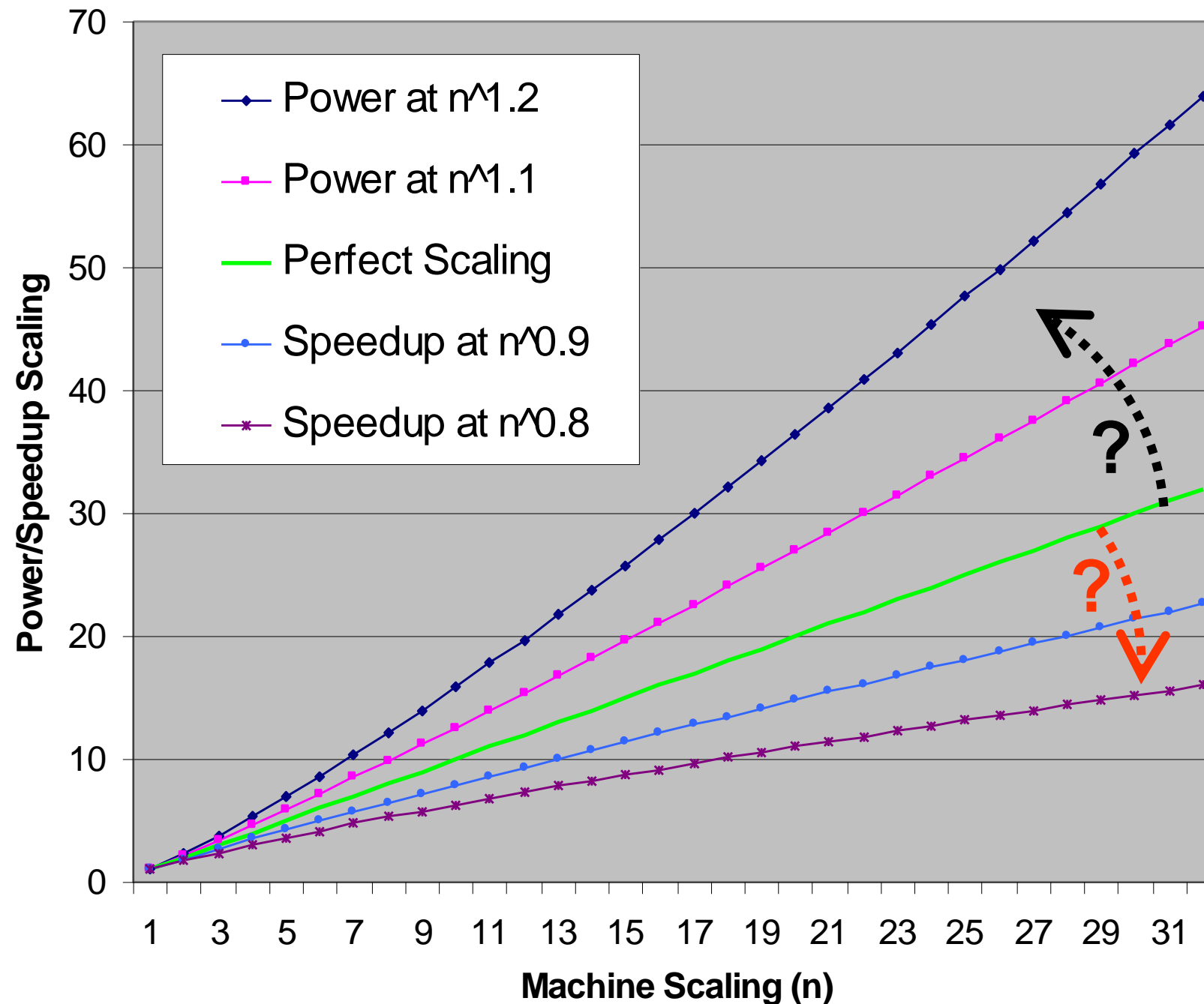
- The issue is not small vs. big cores, nor in-order vs. out-of-order cores. The key metric is EPI.
- The ideal core: ultra-low EPI with best possible single-thread or single-core performance.



CPU	EPI	SU
i486	7 nj	1
P5	10 nj	2
P6	17 nj	3.5
P4P-wmt	27 nj	6
P4P-psc	29 nj	6.5
Pentium M	9 nj	5.5
Neo-Core	5 nj	6.5

NP/DSP/GPU	EPI
IXP2800	1 nj
TMS320C6713	0.7 nj
GeF7800GTX	0.6 nj

Power and Speedup Scaling



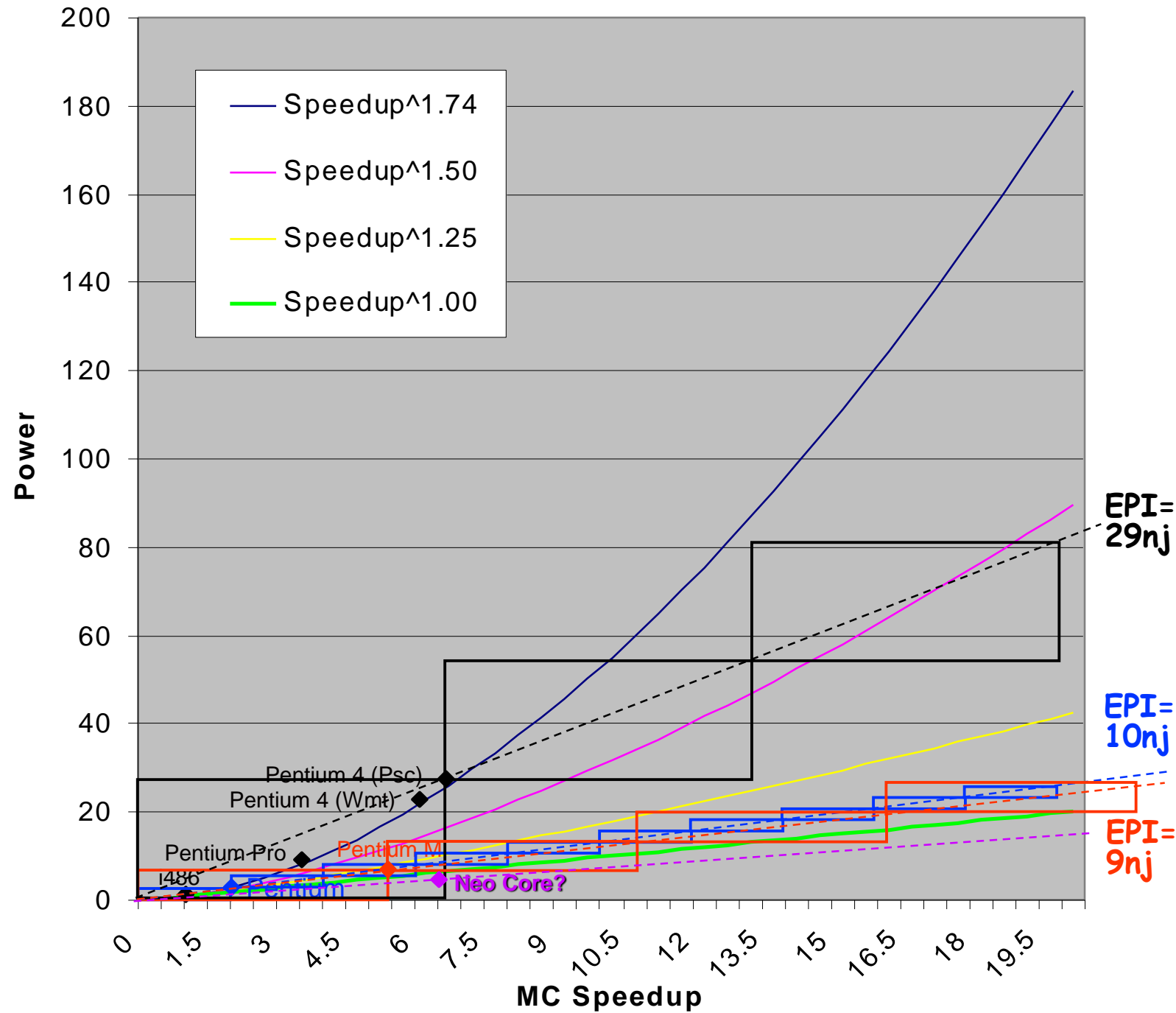
Power scaling challenges:

- Low EPI cores
- Un-core scaling

Speedup scaling challenges:

- Algorithm
 - Sequential %
- Architecture
 - PL scaling
 - CPI scaling

Power vs. MC Speedup Scaling



CPU	EPI	SU
i486	7 nj	1
P5	10 nj	2
P6	17 nj	3.5
P4P-wmt	27 nj	6
P4P-psc	29 nj	6.5
Pentium M	9 nj	5.5
Neo core?	< 5 nj	6.5

- The scaling goal is not just the number of cores but maximum throughput within fixed power envelope.
- The key issue is not the power scaling of replicated cores, but the un-core power scaling that may push total power scaling towards the square law again.

Future Directions

❖ Scalability Strategies:

- Algorithm – Languages & Specialized Parallelism
- Architecture – CPI and Path Length Reduction
- Power/Thermal – EPI Reduction & Scalable Un-core

❖ *Power/Energy is the new scalability wall*

❖ Research Challenges:

- Sequential % mitigation for compelling workloads
- Ultra-low EPI core with great latency performance
- Un-core fabric with near-linear power scaling

❖ *Un-core scaling is the new power goblin*

Useful Chart on Power, Voltage, Resistance, Current

Formula 1 – Electrical power equation:

Power $P = I \times V = R \times I^2 = V^2 / R$

where power P is in watts, voltage V is in volts and current I is in amperes (DC).

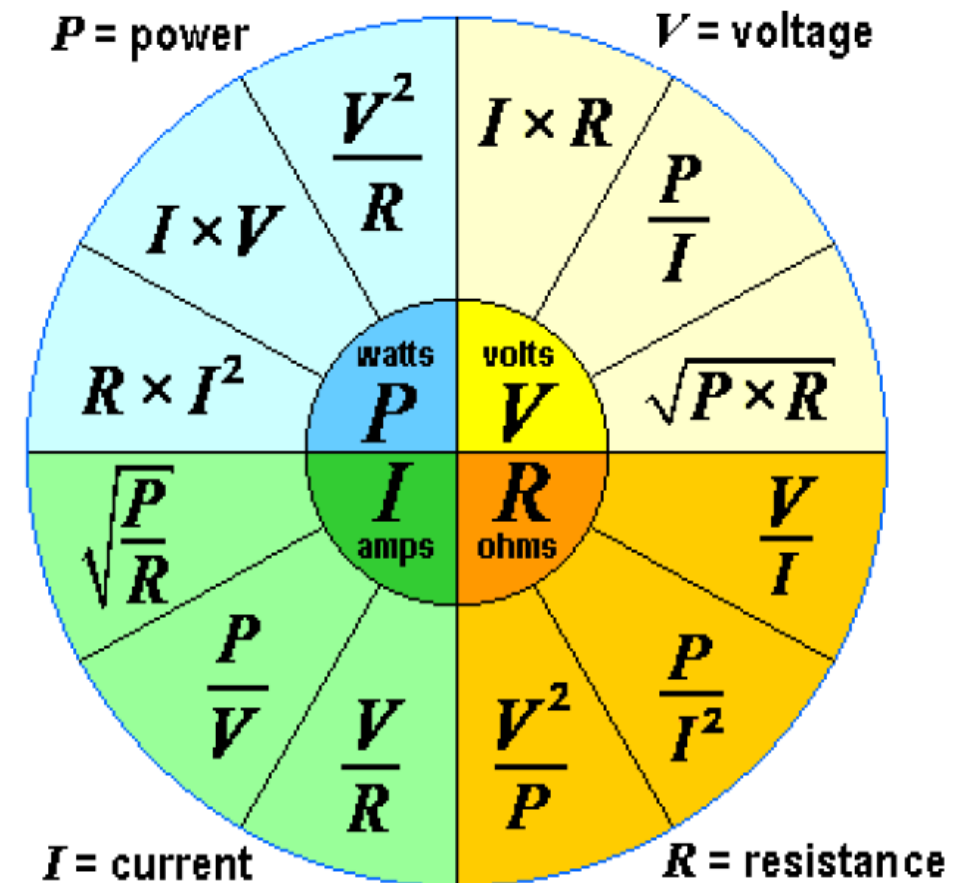
If there is AC, look also at the power factor $PF = \cos \varphi$ and φ = power factor angle (phase angle) between voltage and amperage.

Formula 2 – Mechanical power equation:

Power $P = E / t = W / t$

where power P is in watts, **Energy** E is in joules, and time t is in seconds. 1 W = 1 J/s.

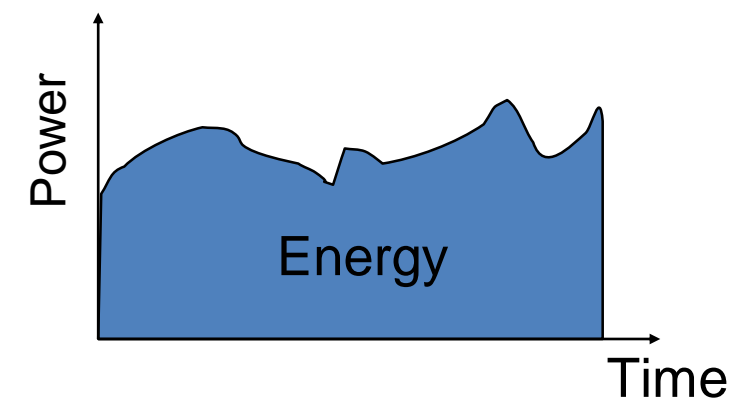
Electric **Energy** is $E = P \times t$ – measured in watt-hours, or also in kWh. 1J = 1N×m = 1W×s



<http://www.sengpielaudio.com/calculator-ohm.htm>

Power vs. Energy

- Energy: integral of power (area under the curve)
 - Energy and power driven by different design constraints
- Power issues:
 - Power delivery (supply current @ right voltage)
 - Thermal (don't fry the chip)
 - Reliability effects (chip lifetime)
- Energy issues:
 - Limited energy capacity (battery)
 - Efficiency (work per unit energy)
- Different usage models drive tradeoffs



F. Power and Energy Optimizations

- With constant time base, two are “equivalent”
 - 10% reduction in power => 10% reduction in energy
- Once time changes, must treat as separate metrics
 - E.g. reduce frequency to save power => reduce performance => increase time to completion => consume more energy (perhaps)
- Metric: energy-delay product per unit of work
 - Tries to capture both effects
 - Others advocate energy-delay²
 - Best to consider all
 - Plot performance (time), energy, ed, ed²

Performance/Power Efficiency Metrics

- Power is good metric for deciding on the thermal envelope of the processor
- Energy is good metric in battery constrained environments
 - Task executed at $\frac{1}{2}$ speed but $\frac{1}{4}$ power means $\frac{1}{2}$ the energy ($2T * \frac{1}{4} P = \frac{1}{2} E$)
 - 2X battery life!
- Energy*Delay metric gives higher weight to performance
 - Same example above, ED ($(2T)^2 * \frac{1}{4} P$) stays same
- Energy*Delay² gives even more weight to performance
 - Same example above shows that $\frac{1}{2}$ speed is 2X worse on ED² metric

<p>MIPS/watt</p>	<p>Common measure of power efficiency Equivalent to energy per instruction Independent of time Ideal metric for throughput performance</p> $\frac{\text{Mips}}{\text{Watt}} = \frac{\frac{\text{Instructions}}{\text{Second}}}{\frac{\text{Joules}}{\text{Second}}} = \frac{\text{Instructions}}{\text{Joule}}$	<table border="1"> <caption>mips/watt</caption> <thead> <tr> <th>Processor</th> <th>Normalized Metric</th> </tr> </thead> <tbody> <tr> <td>486</td> <td>1.0</td> </tr> <tr> <td>p5</td> <td>~0.75</td> </tr> <tr> <td>p6</td> <td>~0.42</td> </tr> <tr> <td>pentium 4</td> <td>~0.28</td> </tr> </tbody> </table>	Processor	Normalized Metric	486	1.0	p5	~0.75	p6	~0.42	pentium 4	~0.28
Processor	Normalized Metric											
486	1.0											
p5	~0.75											
p6	~0.42											
pentium 4	~0.28											
<p>MIPS^2/watt</p>	<p>Equivalent to energy • delay Common metric for comparing logic families</p>	<table border="1"> <caption>mips^2/watt</caption> <thead> <tr> <th>Processor</th> <th>Normalized Metric</th> </tr> </thead> <tbody> <tr> <td>486</td> <td>1.0</td> </tr> <tr> <td>p5</td> <td>~1.48</td> </tr> <tr> <td>p6</td> <td>~1.52</td> </tr> <tr> <td>pentium 4</td> <td>~1.62</td> </tr> </tbody> </table>	Processor	Normalized Metric	486	1.0	p5	~1.48	p6	~1.52	pentium 4	~1.62
Processor	Normalized Metric											
486	1.0											
p5	~1.48											
p6	~1.52											
pentium 4	~1.62											
<p>MIPS^3/watt</p>	<p>Equivalent to energy • delay^2 Assign increasing weight to time Appropriate metric for latency performance</p>	<table border="1"> <caption>mips^3/watt</caption> <thead> <tr> <th>Processor</th> <th>Normalized Metric</th> </tr> </thead> <tbody> <tr> <td>486</td> <td>1.0</td> </tr> <tr> <td>p5</td> <td>~2.9</td> </tr> <tr> <td>p6</td> <td>~5.5</td> </tr> <tr> <td>pentium 4</td> <td>~10.0</td> </tr> </tbody> </table>	Processor	Normalized Metric	486	1.0	p5	~2.9	p6	~5.5	pentium 4	~10.0
Processor	Normalized Metric											
486	1.0											
p5	~2.9											
p6	~5.5											
pentium 4	~10.0											

Energy Storage Challenge for Mobile Devices

How long between recharges?

iPad2

- 25Wh = 90kJ
- 10h use = avg 2.5W

Kindle3

- 6.5Wh = 23kJ
- “30 day use” = 30h
avg 220mW

Typical smartphone

- 32g, 13cc, 5.5Wh = 18kJ
- +5h charging @ max 1W
- 20mW static power = 10 days standby
- 150mW notifications = 1.3 day standby
- “typical usage” 5kJ active + 13kJ standby =
1 battery charge



Key Challenges...how far can we go?

[Per Ljung, 2012]

Key Areas

Description

Active Power Saving

- **Context Based Power Management:** Offloading of communication to available connectivity, and computation to companion devices or the cloud.
- **Workload Based Power Management:** Manage power consumption based on actual usage and workload scenarios by leveraging heterogeneous cores and smart parallelism to reduce overall power.

Standby Power Saving

- **Near-Zero-Power Standby Mode:** Low-power always-on transfective bistable (TF/BS) displays; eager hibernation with instant resume.
- **Ultra-Low-Power Always-On Device:** Device with minimal standby functionality and seamless quick switch over to companion devices.

Energy Harvesting

- **Casual Charging:** Wireless inductive charging; solar charging for large surface devices.
- **Anticipatory Preexecution:** Speculative cross-device or cloud-based preprocessing and content prefetching.

Active Power Saving

1.5h @ 900mW = 1.35Wh

avg app

goal <100mW

ideal 100%

5Wh battery: 30%

80%



20%



➤ Context-Based Power Management

- Offload Communication (Local ULP radio)
 - 1m radio instead of 1000m cellular
 - Dongle, access point, femtocell box
- Offload Computation (Cross device)
 - Cloud/companion pre-compute/pre-render

80% time in home/office
40mW vs 1W

25x

trade cheap communication
for expensive computation
40mW vs 1W

25x

➤ Workload-Based Power Management

- Power consumption based on usage scenario
 - Activity, location, history
- Approach zero energy waste
 - De-powering unused peripherals & power islands
- Heterogeneous cores & SP processor arrays (ULP)

0W vs 160mW (email, notifications)

10's mW vs 100's mW cores

Standby Power Saving

22.5h @ 160mW = 3.6Wh

default: email, skype

goal <10mW, 1kJ

ideal 0%

5Wh battery: 70%

➤ Near-Zero-Power Standby Mode

- Zero-power always-on displays
 - Transflective, bistable
- Zero-power OS idling
 - Android, Meego, WP7 use 15mW
 - iPhone uses 5mW

TF/BS, TF/LCD, TF/OLED
5mW vs OLED 500mW

100x

better firmware 2mW vs 15mW
with hibernation 1mW vs 15mW

15x

➤ Ultra-Low-Power Always-On Device

- Wearable accessory for notifications & voice
- Seamless quick switch over to companion
- Week-long battery life on small battery

4mW vs 200mW handset

50x

Energy Harvesting

Effectively Approach Unlimited Standby

➤ Casual Charging

- Wireless inductive chargers
- Instant charging for quick fix
- Solar for large surface areas
- Cross-device energy sharing

1W desk, nightstand, car

5h full charge

best case charge 3W
for 1W tablet

1h charge → 3h use

➤ Anticipatory Preexecution

- Speculative pre-processing and pre-fetching
- Cloud based or cross-device based
- Context and user behavior model driven

40mW vs 1W

25x

Refactoring the Mobile Form Factors?

Multiple Devices, One Seamless Experience

- Laptop (avg 10W → avg 2W)
 - Exploit large battery & connectivity
 - Runs offloaded apps from Phone
 - Gathers data for Wearable
 - + Display normally off
 - + 50 WH battery → 25 hours nonstop use
- Smartphone (avg 200mW → avg 40mW)
 - + Normally hibernating in standby
 - + Offload apps to Laptop
 - + Longer battery life: **5x battery life**
 - + 5 WH battery → 125 hours nonstop use
- Wearable (avg 3mW)
 - + Always-on display
 - + Always fresh data feeds
 - Respond using Laptop or Smartphone
 - + Reduces avg power of Laptop & Smartphone
 - + 1.2 WH battery → 400 hours nonstop use



18-600 Foundations of Computer Systems

Lecture 26: "Future of Computer Systems"

John P. Shen

December 4, 2017

Next Time ...

