

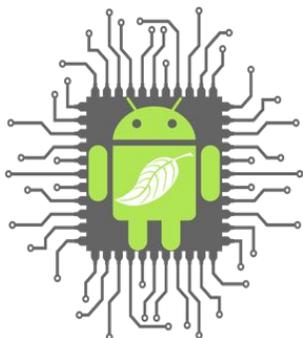
# Is Dark Silicon Useful?

## *Harnessing the Four Horsemen of the Coming Dark Silicon Apocalypse*

---

***Michael B. Taylor***

*Associate Professor (July 2012)  
University of California, San Diego*



*Presented at DAC 2012 and DaSi 2012*



# Is Dark Silicon Useful?

Dark Silicon

*Harnessing the Four Horsemen  
of the Coming Dark Silicon Apocalypse*



Prof. Michael B. Taylor

UC San Diego

# **This Talk**

*The Dark Silicon Apocalypse*



*Explaining the Source of Dark Silicon*

*The Four Horsemen*

# **ISCA 2002 Session I: We Had It All Figured Out**

- The Optimum Pipeline Depth for a Microprocessor  
IBM (22-36 pipeline stages)
- The Optimal Logic Depth Per Pipeline Stage is 6 to 8 FO4 Inverter Delays (~40 pipeline stages)  
Dec/Compaq/HP
- Increasing Processor Performance by Implementing Deeper Pipelines (~50-60 stages)  
Intel

**Universal Conclusion: Frequency-Boosted Microarch == Future**

# 2004: *Santa Clara, we have a problem!*

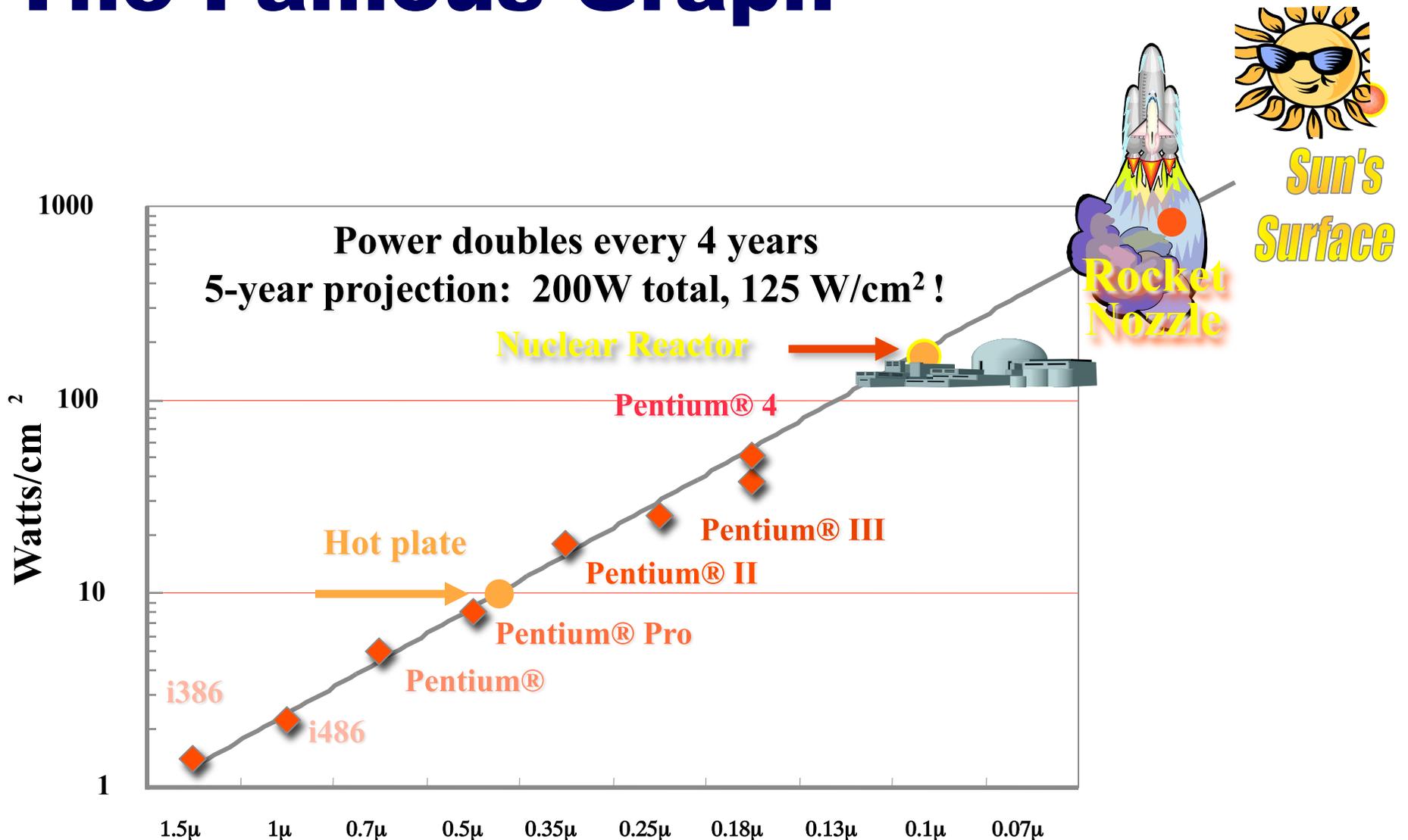
More pipeline stages,  
less efficient, more power.

Just can't remove  
> 100 watts  
without great expense on  
a desktop.

All computing is now  
Low Power Computing!

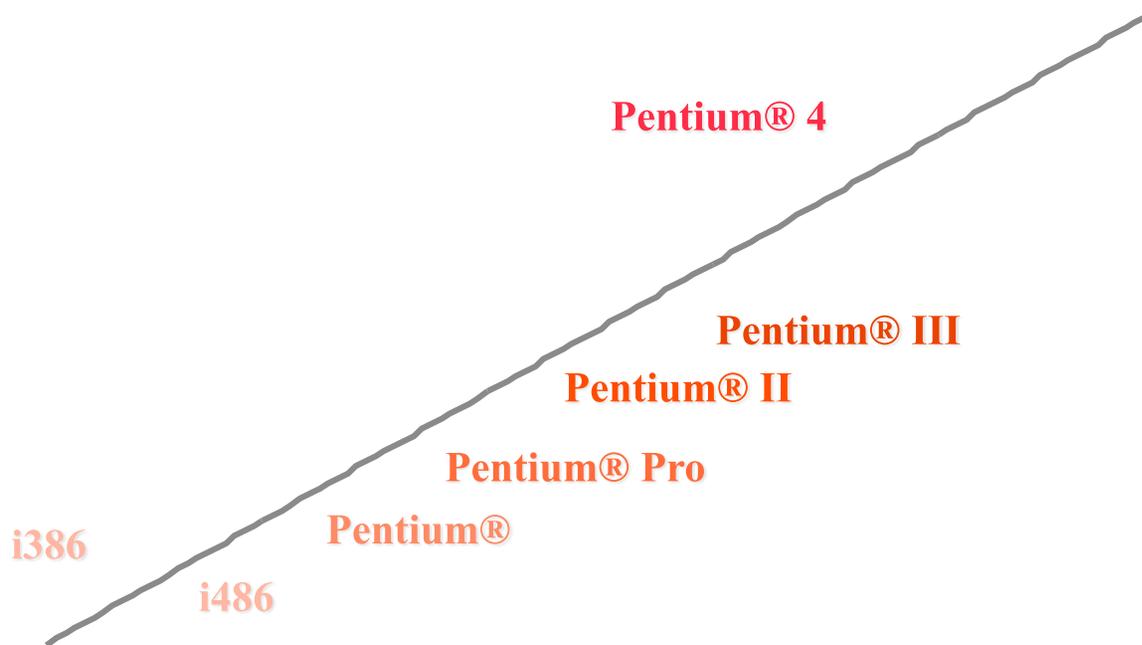


# The Famous Graph



From "New Microarchitecture Challenges in the Coming Generations of CMOS Process Technologies"  
– Fred Pollack, Intel Corp. Micro32 conference key note - 1999.

# **Widespread Assumption: *Microarchitecture was the cause of the power problem***



From “New Microarchitecture Challenges in the Coming Generations of CMOS Process Technologies” – Fred Pollack, Intel Corp. Micro32 conference key note - 1999.

# Back to the future ...

PPro/P3:

**12 stages**

~~P4 (h4 paper):~~

~~**20 stages**~~

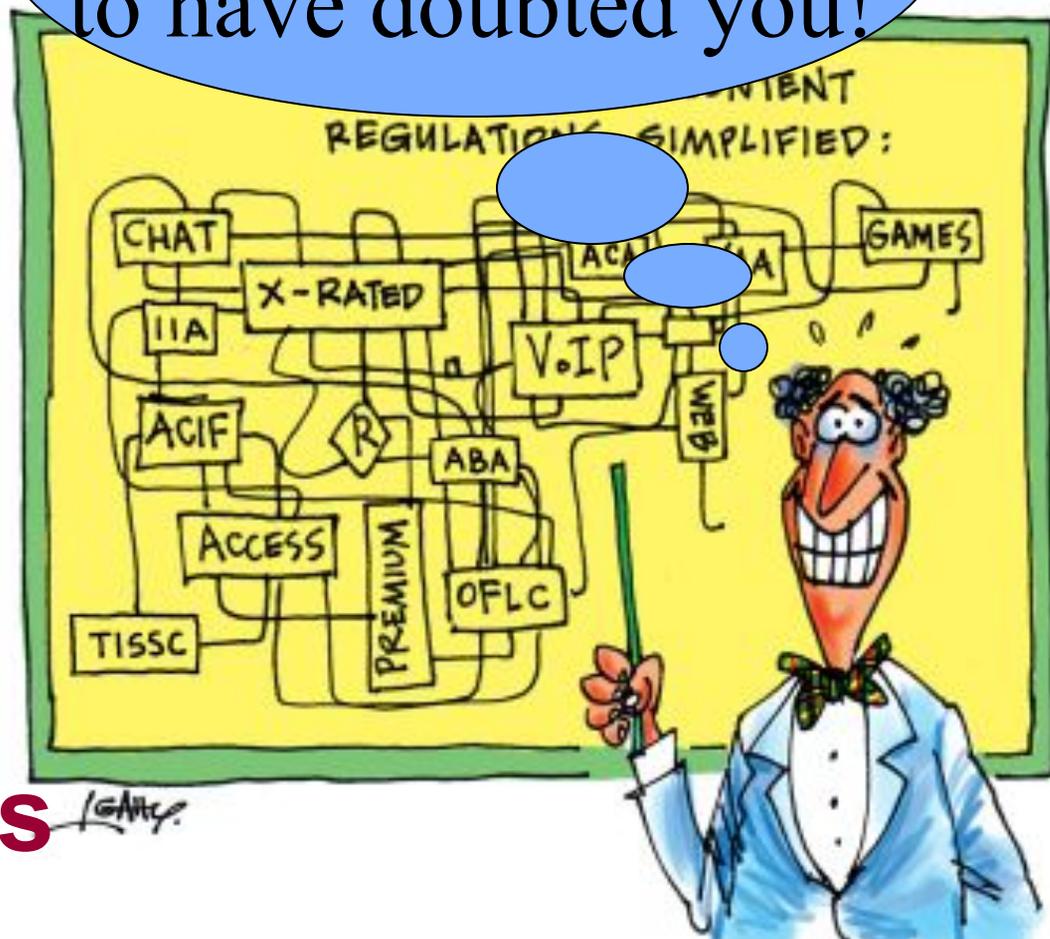
~~P4/prescott:~~

~~**31 stages**~~

~~P5/Tejas:~~

~~**>> 31 stages**~~

Oh P Pro, I'm sorry to have doubted you!



# And forward to multicore...

PPro/P3:

**12 stages**

~~P4 (b4 paper):~~

~~**20 stages**~~

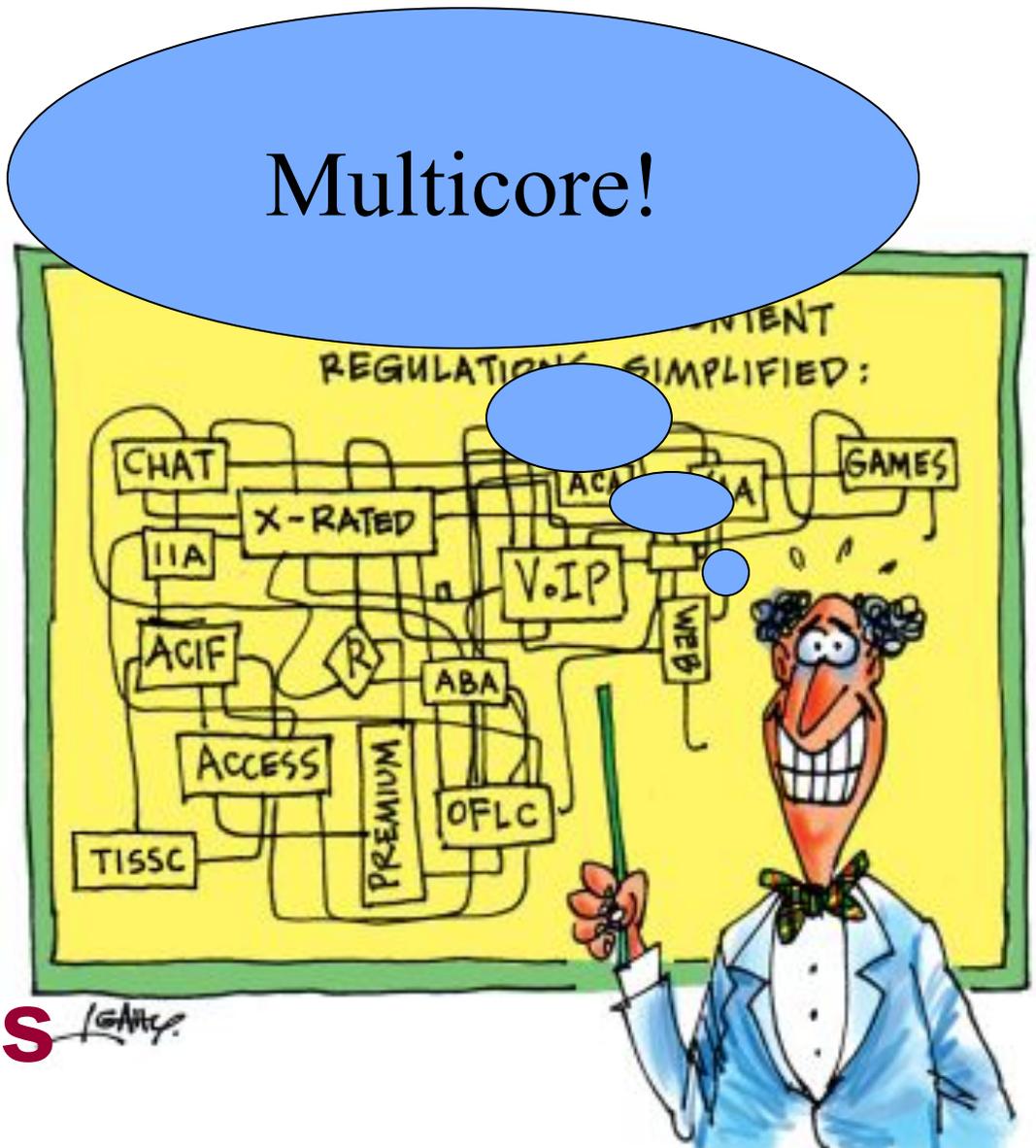
~~P4/prescott:~~

~~**31 stages**~~

~~P5/Tejas:~~

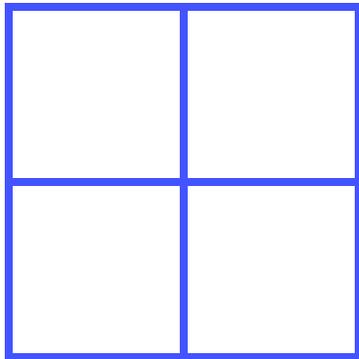
~~**>> 31 stages**~~

Multicore!



# The Scaling Promise of Multicore

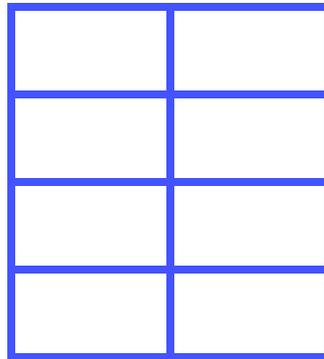
4 cores  
1.8 GHz



65 nm



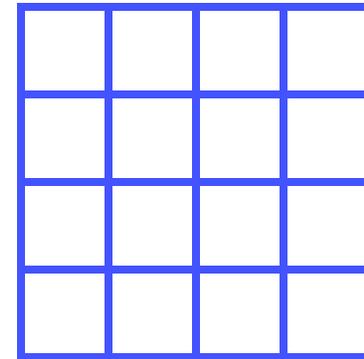
8 cores  
 $\geq 1.8$  GHz



45 nm



16 cores  
 $\geq 1.8$  GHz

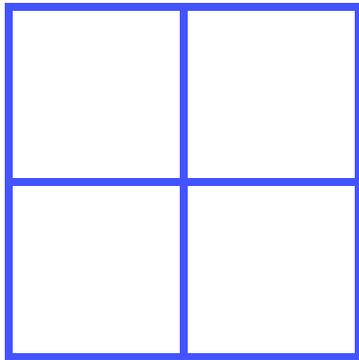


32 nm

2x cores per generation,  
flat or slightly growing frequency

# But actually, that's not what's happening

4 cores  
1.8 GHz

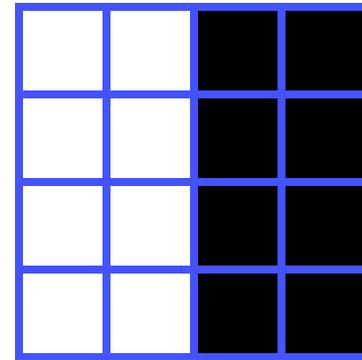


65 nm



45 nm

8 cores  
 $\geq 1.8$  GHz

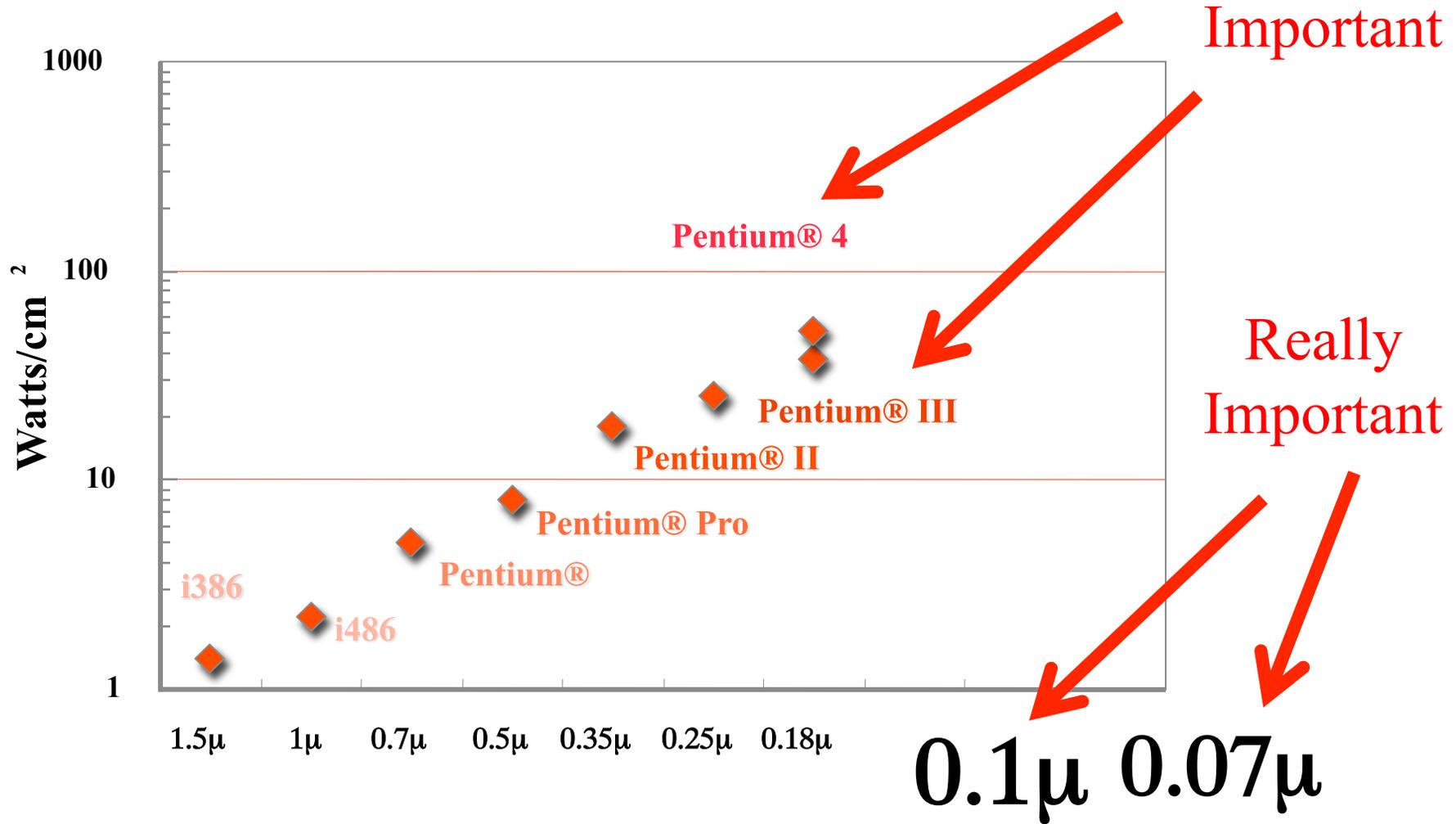


32 nm

1.4x cores per generation,  
flat or slightly growing frequency

Dark or Dim  
Silicon (“uncore”)

# ***Energy Scaling of Process Technology is the Bigger Problem – microarch/multicore just gave us some breathing room.***



# a•poc•a•lypse

noun

(Greek: ἀποκάλυψις apokálypsis; lifting of the veil or revelation)

A disclosure of something hidden from the majority of mankind in an era dominated by misconception ...

# **a•poc•a•lypse**

noun

(Greek: ἀποκάλυψις apokálypsis; lifting of the veil or revelation)

A disclosure of something hidden from the majority of mankind in an era dominated by misconception ...

# **dark sil•i•con a•poc•a•lypse**

noun

Us figuring out what the heck we should do in this new dark silicon design regime.

# **This Talk**

*The Dark Silicon Apocalypse*

*Explaining the Source of Dark Silicon*



*The Four Horsemen*

# **Where does dark silicon come from? And how dark is it going to be?**

The Utilization Wall:

With each successive process generation, the percentage of a chip that can switch at full frequency drops exponentially due to power constraints.

[Venkatesh, ASPLOS '10]

# Scaling 101: Moore's Law

90 65 45 32 22 16 11 8 nm



$$S = \frac{22}{16} = \sim 1.4x$$

# Scaling 101:

*Transistors scale as  $S^2$*

180 nm

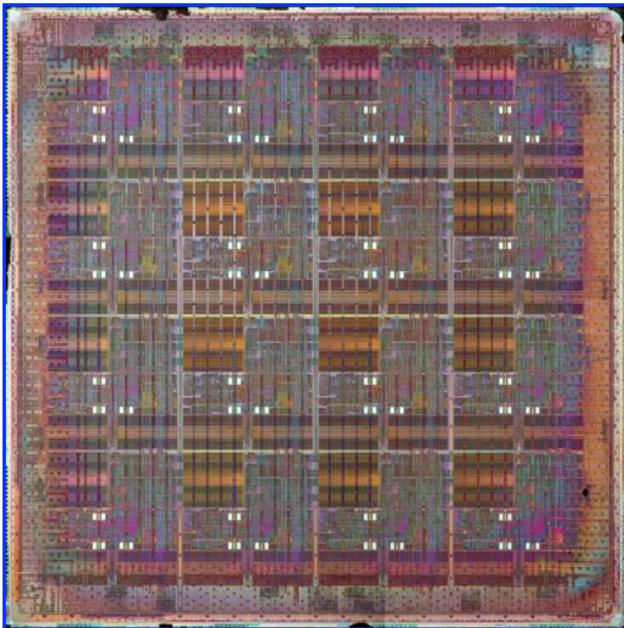
16 cores

$S = 2x$

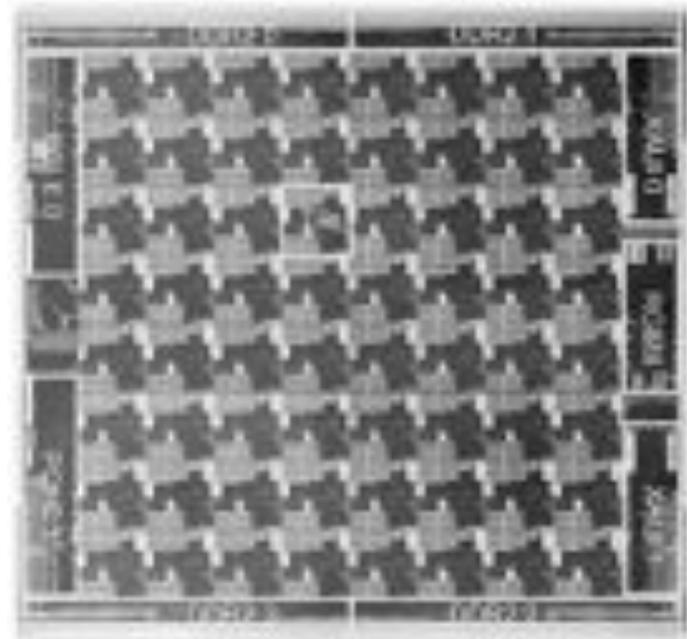
Transistors = 4x

90 nm

64 cores



MIT Raw



Tiler TILE64

# Advanced Scaling:

***Dennard: “Computing Capabilities***

If  $S=1.4x \dots$

***Scale by  $S^3 = 2.8x$ ”***



Design of Ion-Implanted MOSFETs with Very Small Dimensions

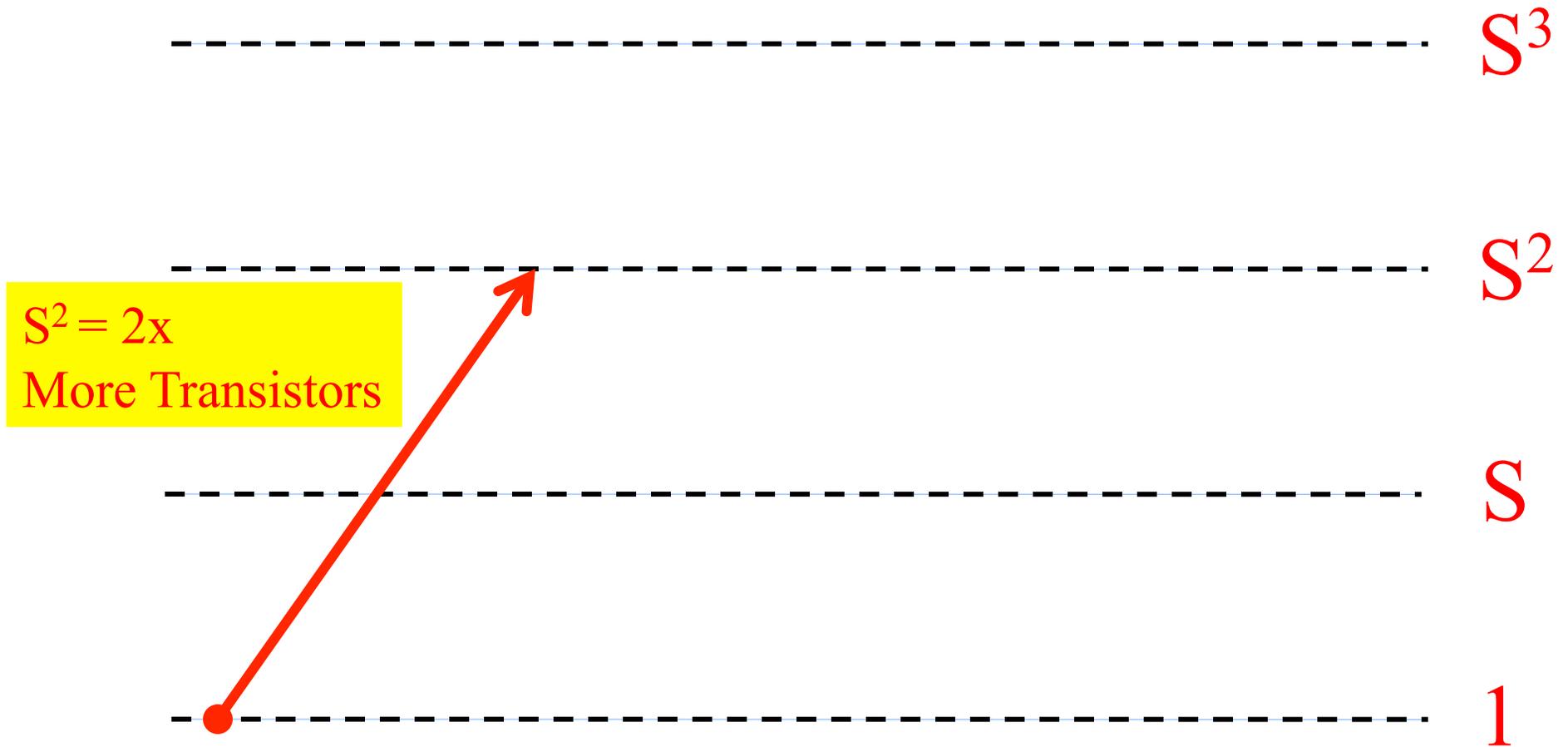
Dennard et al, 1974

# Advanced Scaling:

**Dennard: “Computing Capabilities**

If  $S=1.4x$  ...

**Scale by  $S^3 = 2.8x$ ”**

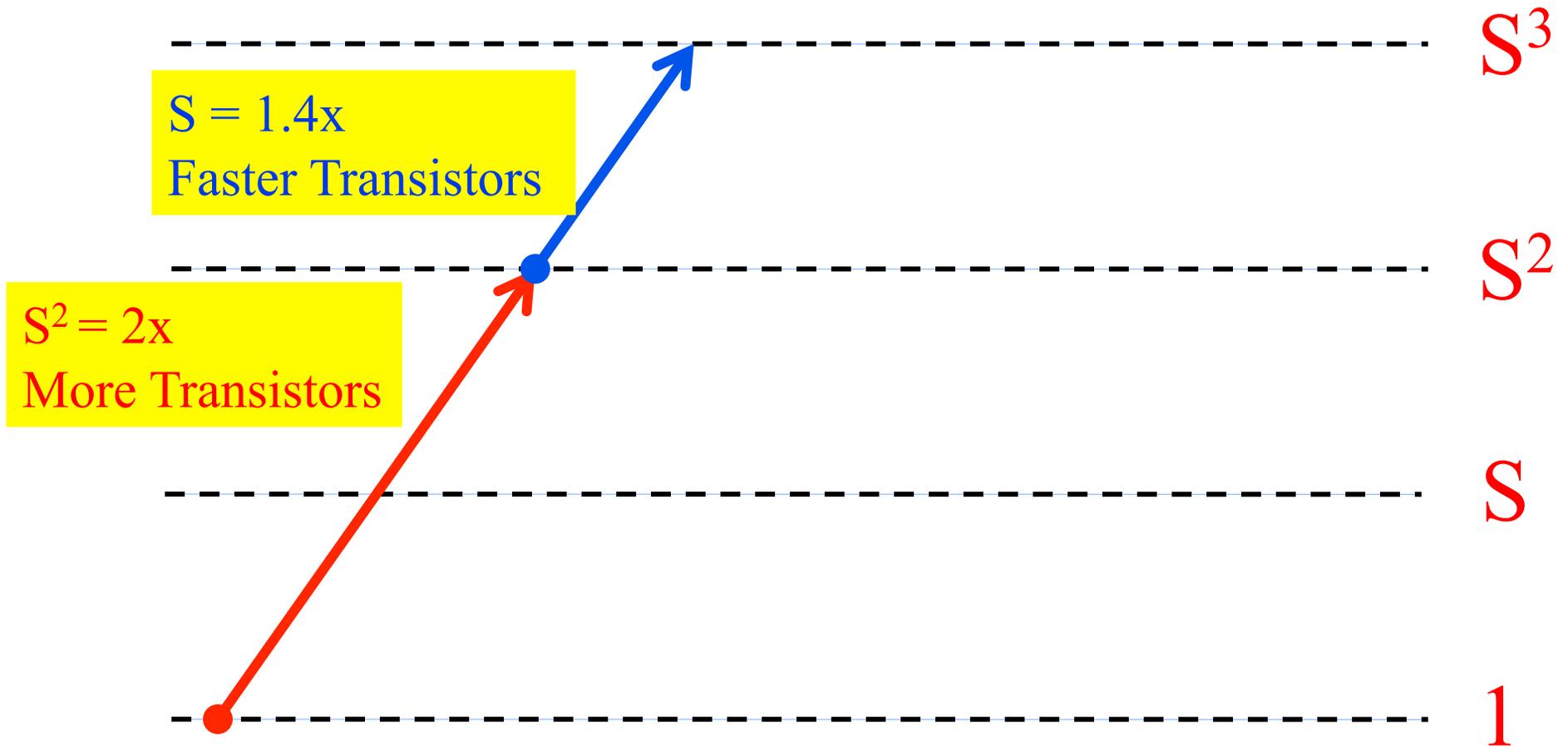


# Advanced Scaling:

**Dennard: “Computing Capabilities**

If  $S=1.4x$  ...

**Scale by  $S^3 = 2.8x$ ”**

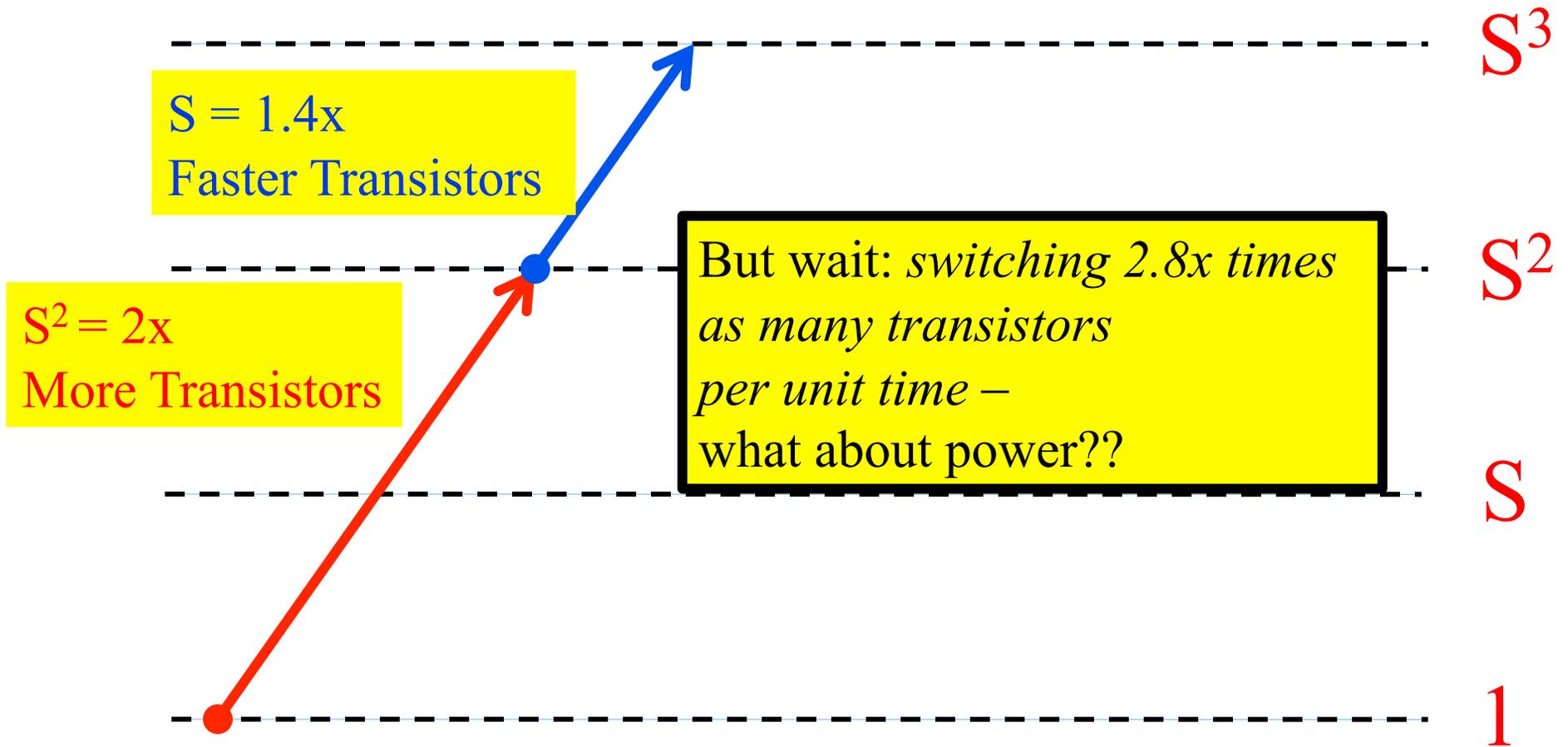


# Advanced Scaling:

**Dennard: “Computing Capabilities**

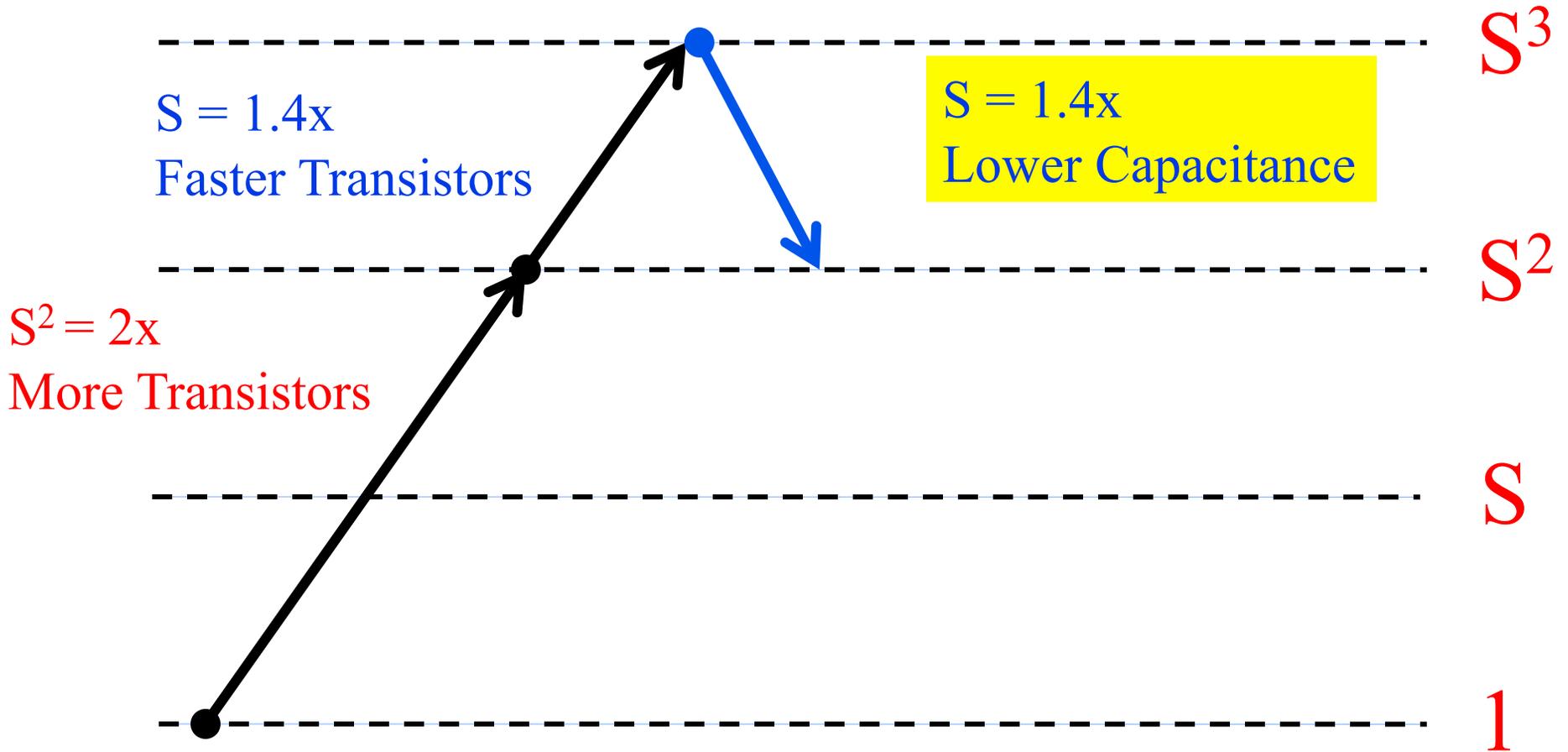
If  $S=1.4x$  ...

**Scale by  $S^3 = 2.8x$ ”**



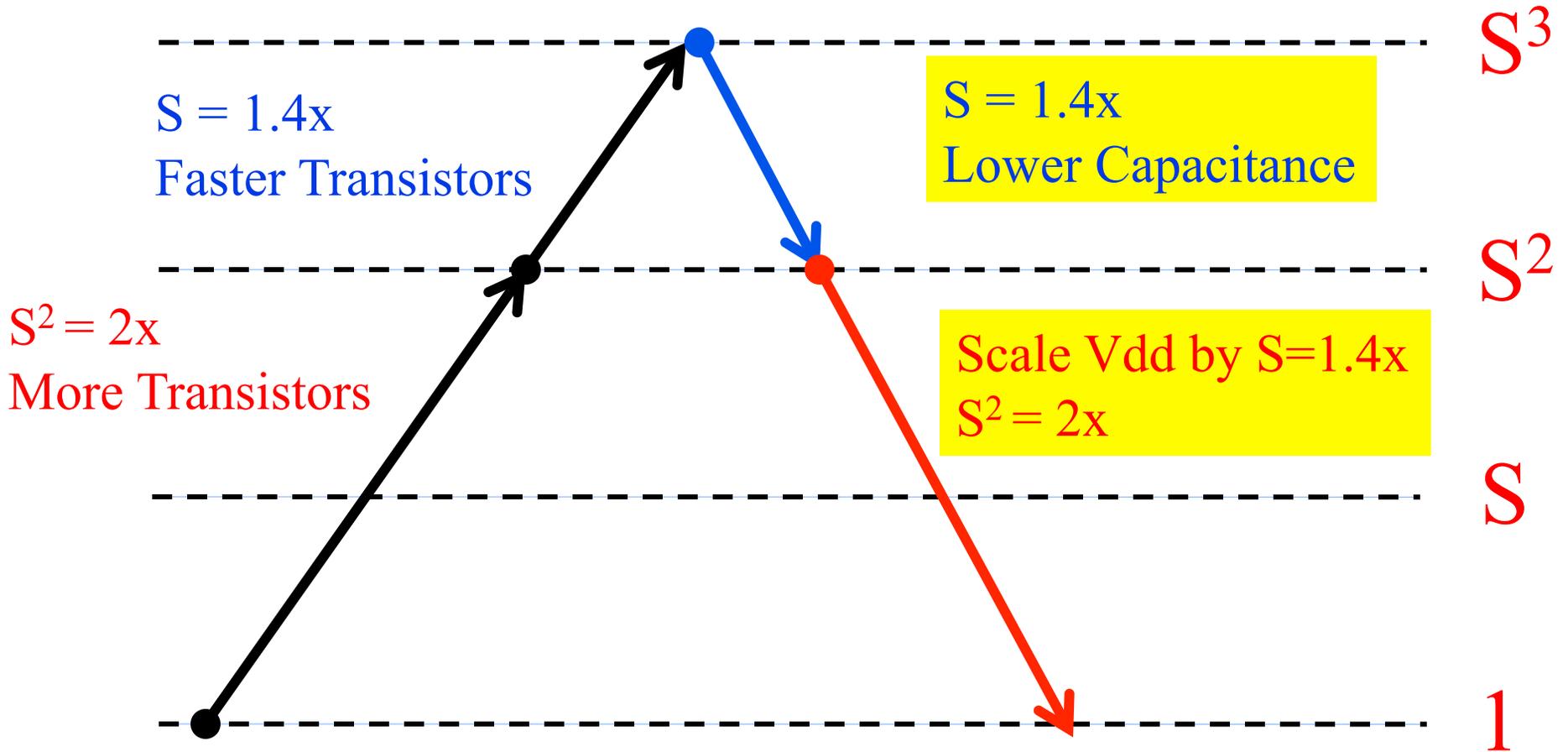
# Dennard:

***“We can keep power consumption constant”***



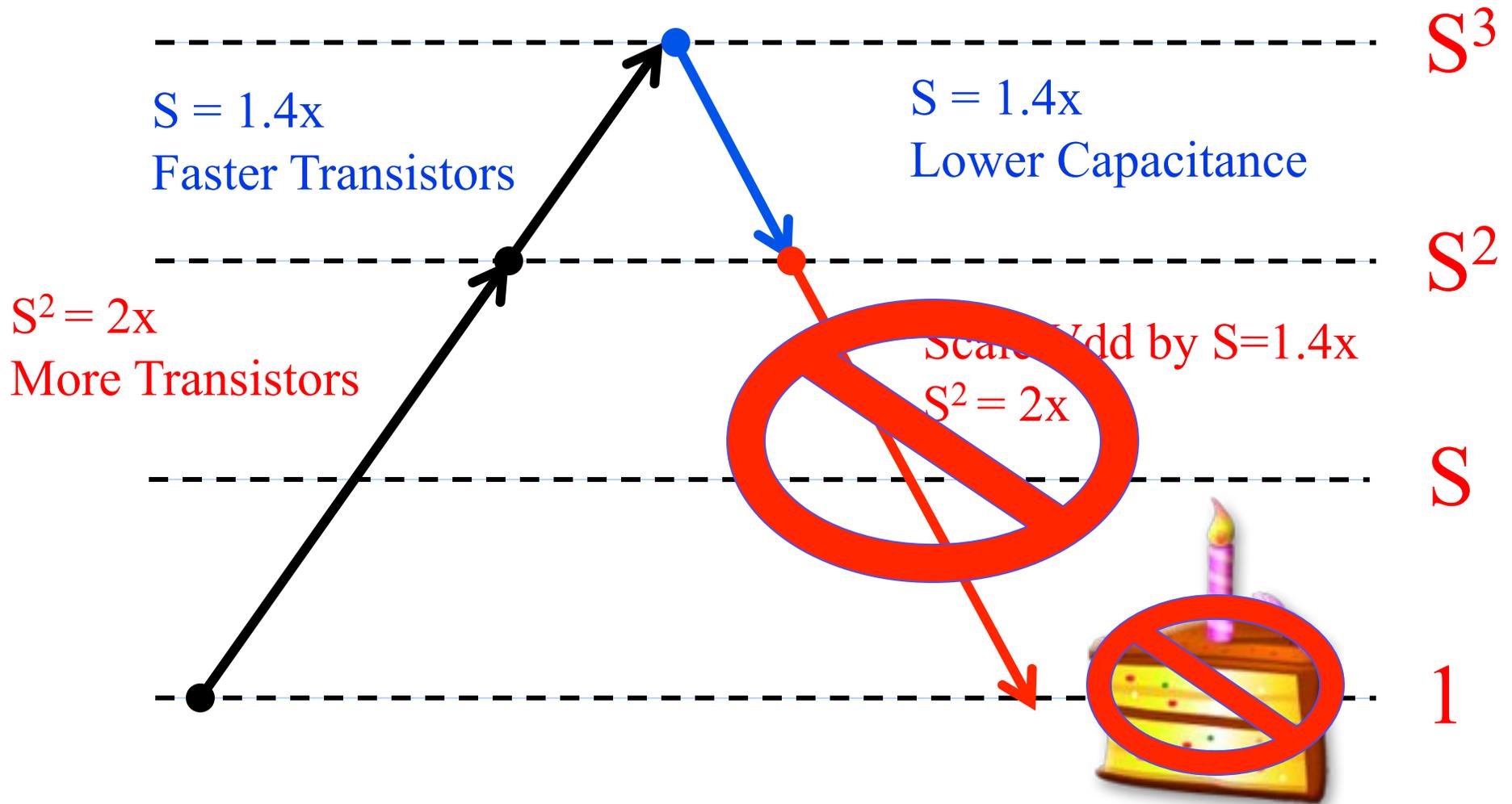
# Dennard:

***“We can keep power consumption constant”***

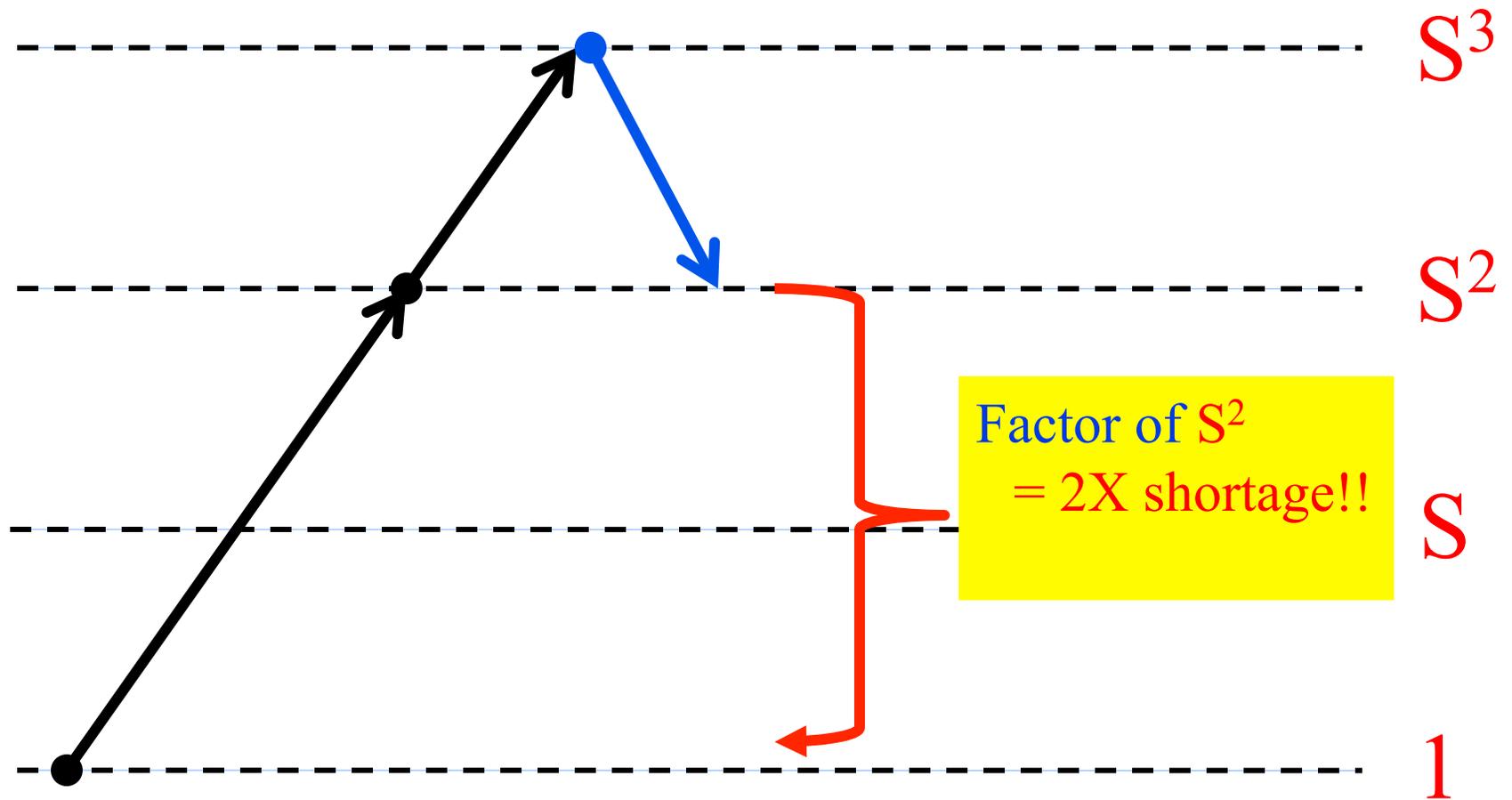


# Fast forward to 2005:

## *Threshold Scaling Problems due to Leakage Prevents Us From Scaling Voltage*



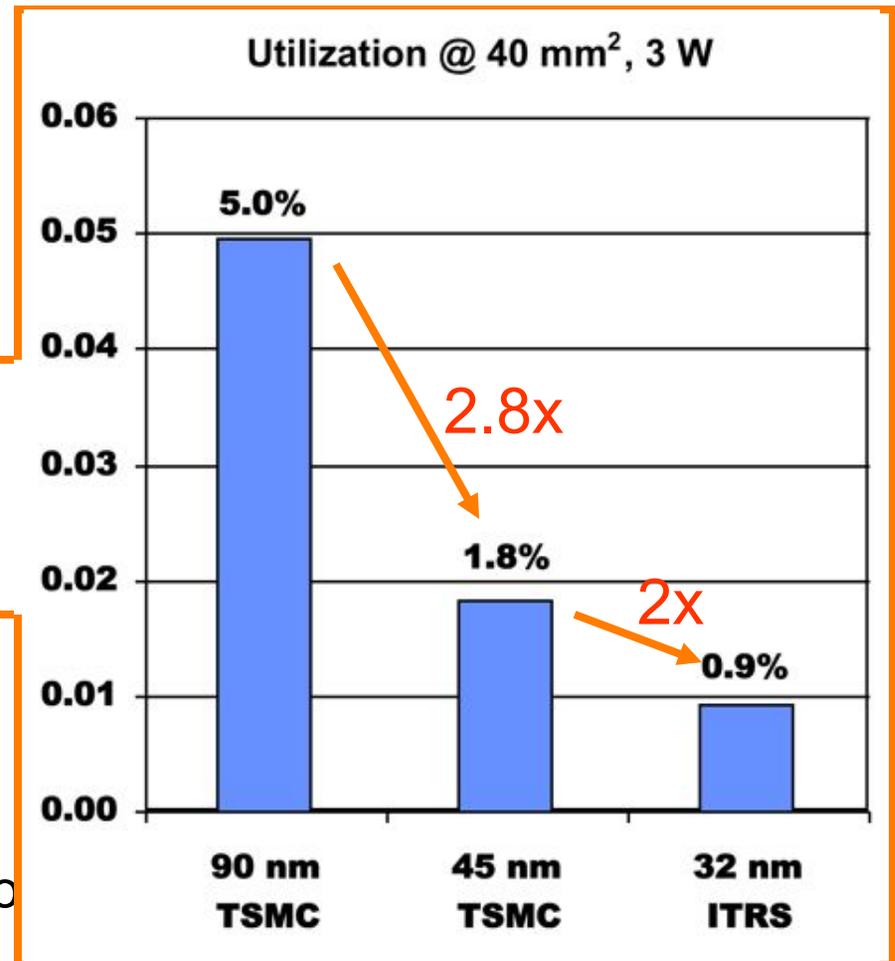
**Full Chip, Full Frequency Power Dissipation  
Is increasing exponentially by 2x with  
every process generation**



# We've Hit The Utilization Wall

*Utilization Wall: With each successive process generation, the percentage of a chip that can actively switch drops exponentially due to power constraints.*

- **Scaling theory**
  - Transistor and power budgets are no longer balanced
  - Exponentially increasing problem!
- **Experimental results**
  - Replicated a small datapath
  - More "dark silicon" than active
- **Observations in the wild**
  - Flat frequency curve
  - "Turbo Mode"
  - Increasing cache/processor ratio



[Venkatesh, ASPLOS '10]

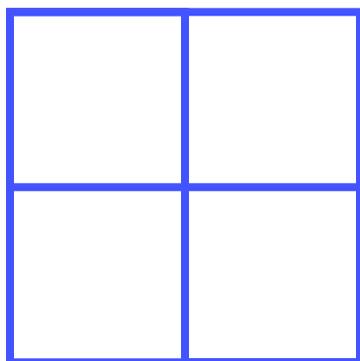
# Multicore has hit the Utilization Wall

Spectrum of tradeoffs  
between # of cores and  
frequency

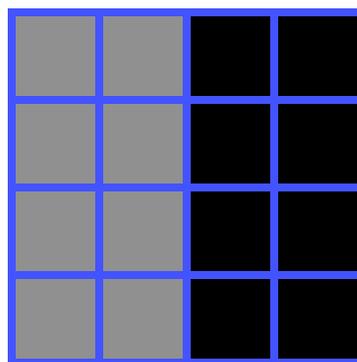
Example:

65 nm  $\rightarrow$  32 nm ( $S = 2$ )

4 cores @ 1.8 GHz



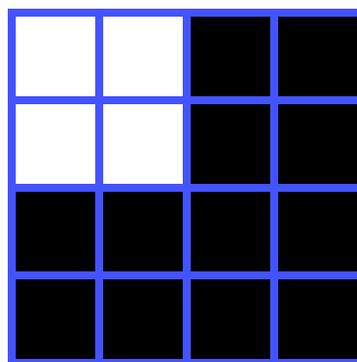
65 nm



4x4 cores @ .9 GHz  
(*GPUs of future?*)

2x4 cores @ 1.8 GHz  
(8 cores dark, 8 dim)

(*Intel/x86 Choice,*  
*next slide*)



4 cores @ 2x1.8 GHz  
(12 cores dark)

32 nm

[Goulding, Hotchips 2010,  
IEEE Micro 2011]

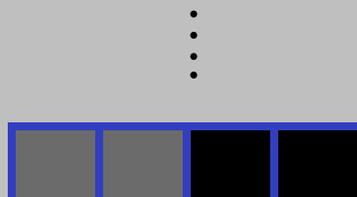
[Esmailzadeh ISCA 2011]

[Skadron IEEE Micro 2011]

[Hardavellas, IEEE Micro 2011]

# Multicore has hit the Utilization Wall

Spectrum of tradeoffs  
between # of cores and  
frequency

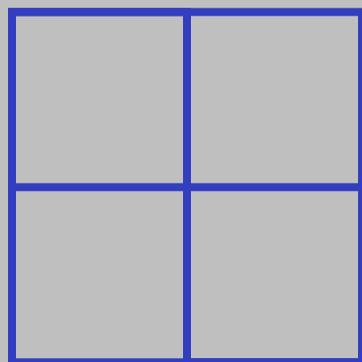


2x4 cores @ 1.8 GHz  
(m)

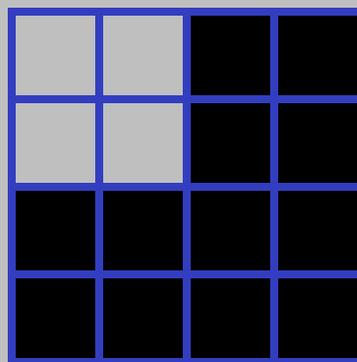
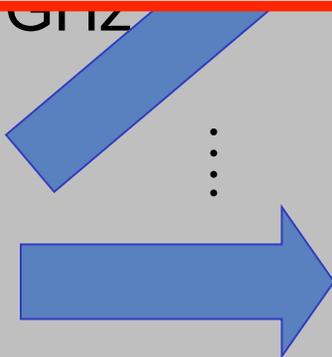
Exa  
65 n

The utilization wall will change the way  
everyone builds chips.

4 cores @ 1.8 GHz



65 nm



4 cores @ 2x1.8 GHz  
(12 cores dark)

32 nm

# **This Talk**

*The Dark Silicon Apocalypse*

*Explaining the Source of Dark Silicon*

*The Four Horsemen*



# The Four Horsemen

What do we do with this dark silicon?

*Four top contenders, each of which seemed like an unlikely candidate from the beginning, carrying unwelcome burdens in design, manufacturing and programming. None is ideal, but each has its benefit and the optimal solution probably incorporates all four of them...*



I



II



III

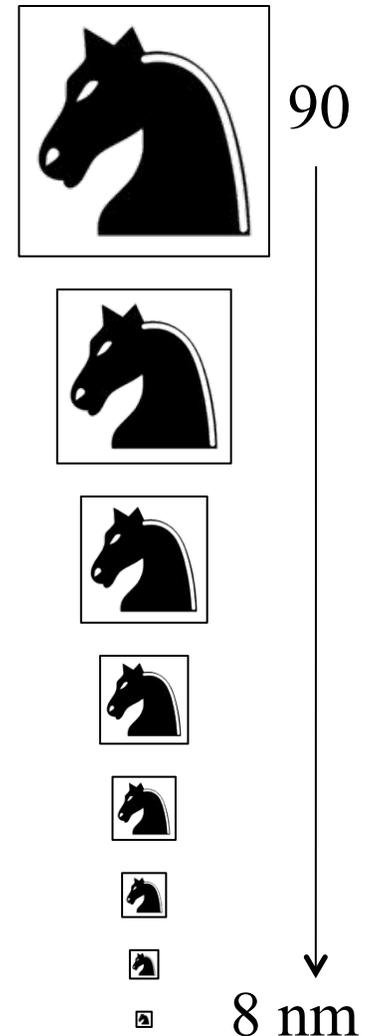


IV

# The Shrinking Horseman (#1)

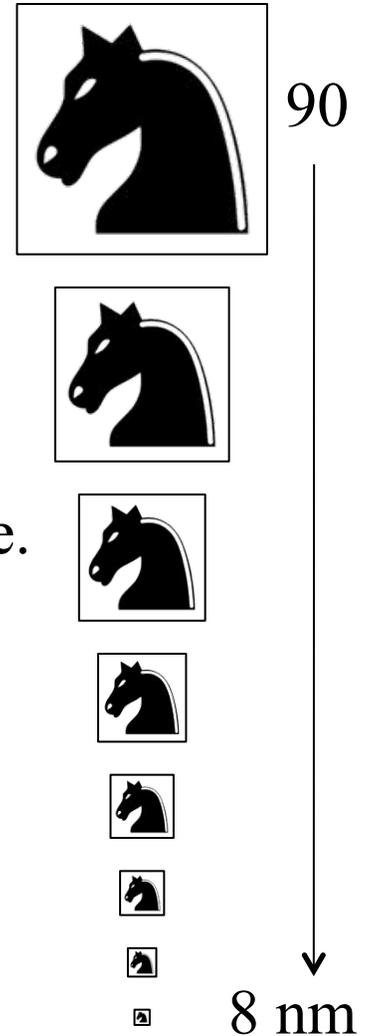
*“Area is expensive. Chip designers will just build smaller chips instead of having dark silicon in their designs!”*

*(if you work on Dark Silicon research, you will hear this a lot...)*



# The Shrinking Horseman (#1)

*“Area is expensive. Chip designers will just build smaller chips instead of having dark silicon in their designs!”*



First, dark silicon doesn't mean *useless silicon*, it just means it's under-clocked or not used all of the time.

There's lots of dark silicon in current chips:

- On-chip GPU on AMD Fusion or Intel Sandybridge for GCC
- L3 cache is very dark for applications with small working sets
- SSE units for integer apps
- Many of the resources in FPGAs not used by many designs (DSP blocks, PCI-E, Gig-E etc)

# The Shrinking Horseman (#1)

*“Just build smaller chips!”*

Possibly – but why didn't we shrink all of our chips before the dark silicon days? This too would be cheaper!

- **Competition and Margins**

- *If there is an advantage to be had from using dark silicon, you have to use it too, to keep up with the Jones.*

- **Diminished Returns**

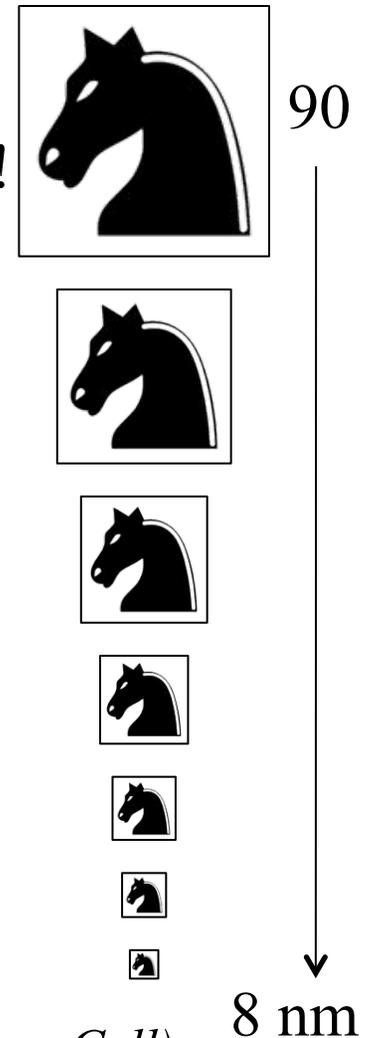
- **e.g., \$10 silicon selling for \$200 today**
  - *Savings Exponentially Diminishing: \$5, \$2.5, \$1.25, 63c*
  - *Overheads: packaging, test, marketing, etc.*
  - *Chip structures like I/O Pad Area do not scale*

- **Exponential increase in Power Density →**

*Exponential Rise in Temperature [Skadron]*

- **But, some chips will shrink**

- *Nasty low margin, high competition chips; or a monopoly (Sony Cell)*



# The Four Horsemen

*The Dark Silicon Apocalypse*

*Explaining the Source of Dark Silicon*

*The Four Horsemen*



I



II



III

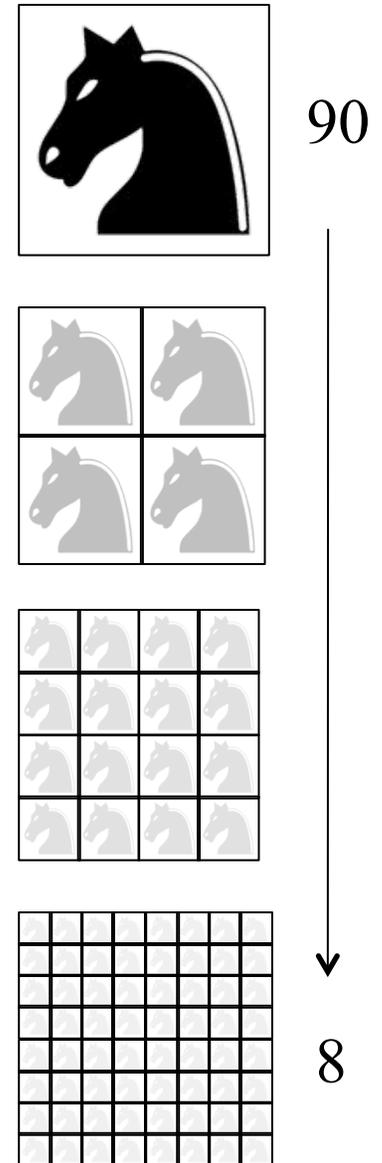


IV

# The Dim Horseman (#2)

*“We will fill the chip with homogeneous cores that would exceed the power budget but we will underclock them (spatial dimming), or use them all only in bursts (temporal dimming)*

*... “dim silicon”.*

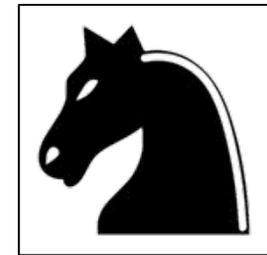


# The Dim Horseman (#2)

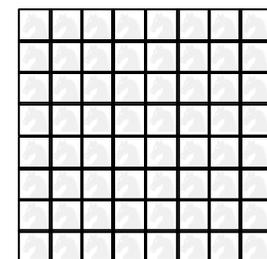
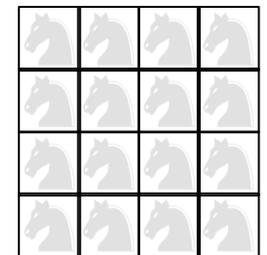
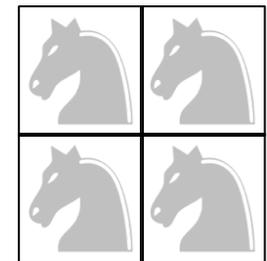
## Spatial Dimming

Gen 1 & 2 Multicores (higher core count → lower freqs)  
Near Threshold Voltage (NTV) Operation

- Delay Loss > Energy Gain
  - But, make it up with lots of dim cores
  - Watch for Non-Ideal Speedups / Amdahl's Law
- Manycore (e.g., [Michigan's Centipede \[ISSCC 2012\]](#))
- SIMD (e.g., [Synctium \[CAL 2010\]](#))
  - Attack issues with Variability and synchronization
- x86 [[Intel, ISSCC 2012](#)]
  - "Solar Powered x86"



90



8

# The Dim Horseman (#2)

## Temporal Dimming

### - Thermally Limited Systems

Turbo Boost 2.0 [Intel, Rotem et al., HOTCHIPS 2011]

- Leverage Thermal Cap for DVFS – “overspend” if cold

Computational Sprinting, [Raghavan HPCA 2012]

- Phase Change, use surplus to power dark silicon instead of DVFS

ARM A15 Core in mobile phone [DAC 2012]

- A15 power usage way above sustainable for phone  
→ 10 second bursts at most ->big.LITTLE

### - Battery Limited Systems

Quad-core mobile application processors



# The Four Horsemen

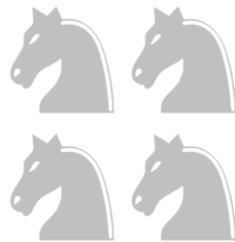
*The Dark Silicon Apocalypse*

*Explaining the Source of Dark Silicon*

*The Four Horsemen*



I



II



III

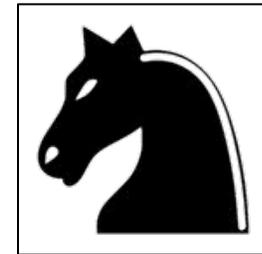


IV

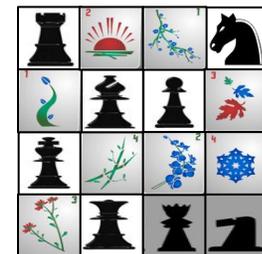
# The Specialized Horseman (#3)

*“We will use all of that dark silicon area to build specialized cores, each of them tuned for the task at hand (10-100x more energy efficient), and only turn on the ones we need...”*

[e.g., Venkatesh et al., ASPLOS 2010,  
Lyons et al., CAL 2010,  
Goulding et al., Hotchips 2010,  
Hardavellas et al. IEEE Micro 2011]



90



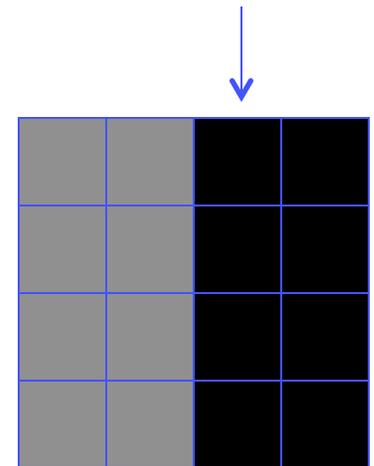
8

# The Specialized Horseman (#3)

## Ex: Conservation Cores (w/ Steven Swanson)

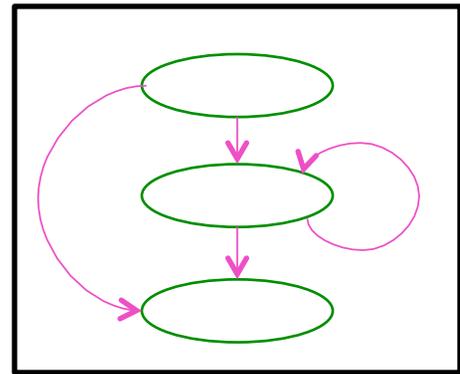
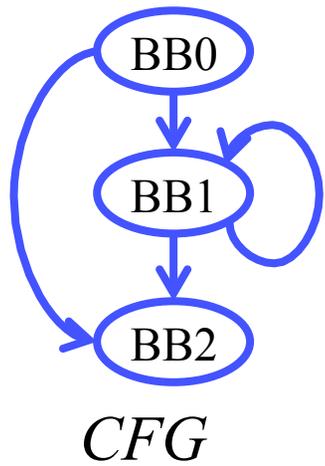
- Idea: Leverage dark silicon to “fight” the utilization wall
- Insights:
  - Power is now more expensive than area
  - Specialized logic can improve energy efficiency by 10-1000x
- C-cores Approach:
  - Fill dark silicon with *Conservation Cores*, or c-cores, which are automatically-generated, specialized energy-saving coprocessors that save energy on common apps
  - Execution jumps among c-cores (hot code) and a host CPU (cold code)
  - Power-gate HW that is not currently in use
  - Coherent Memory & Patching Support for C-cores

Dark Silicon

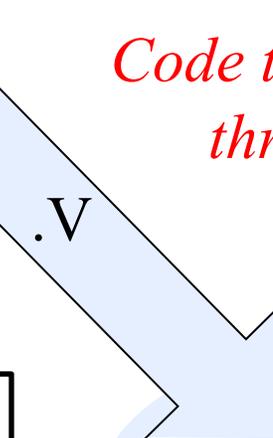
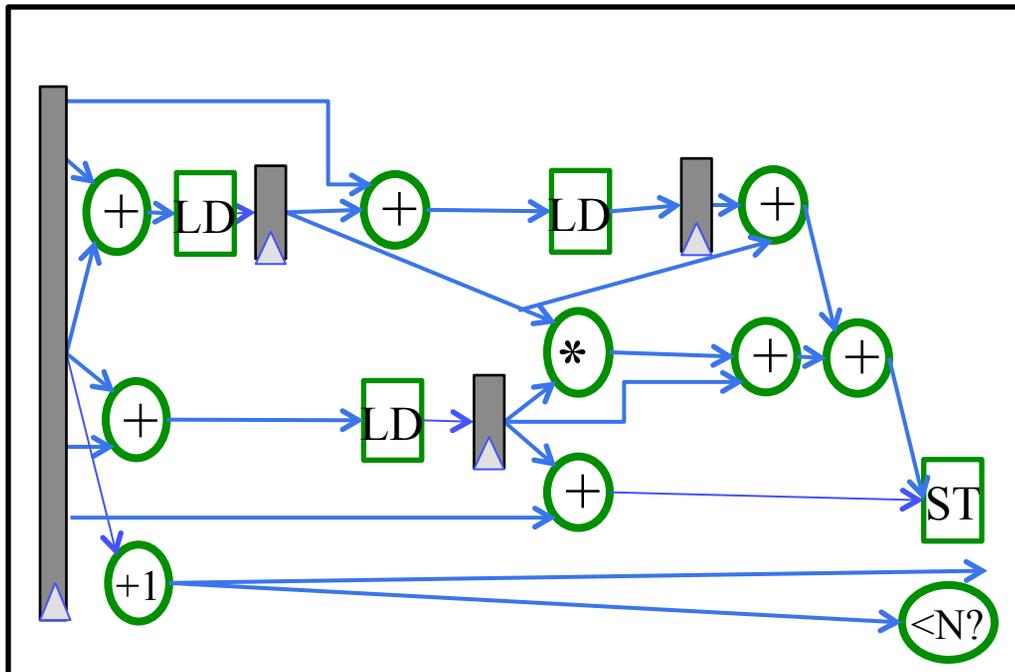


# C-core Generation

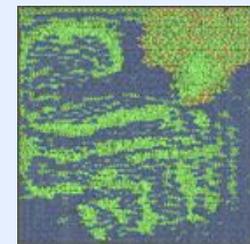
*Code to Stylized Verilog and  
through a CAD flow.*



*Datapath*

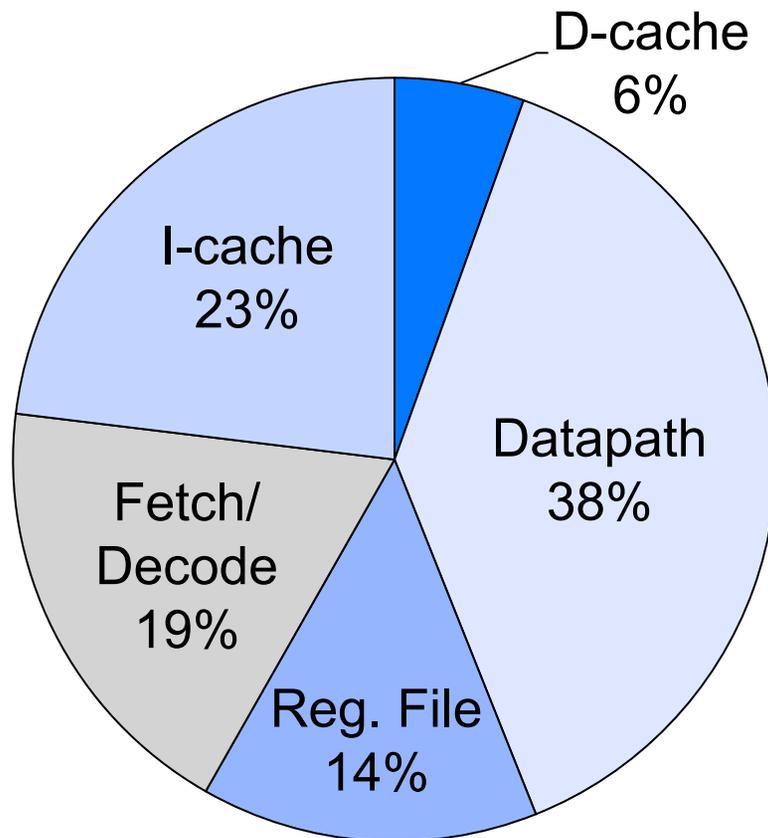
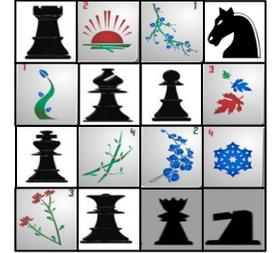


Synopsys  
IC Compiler,  
P&R, CTS



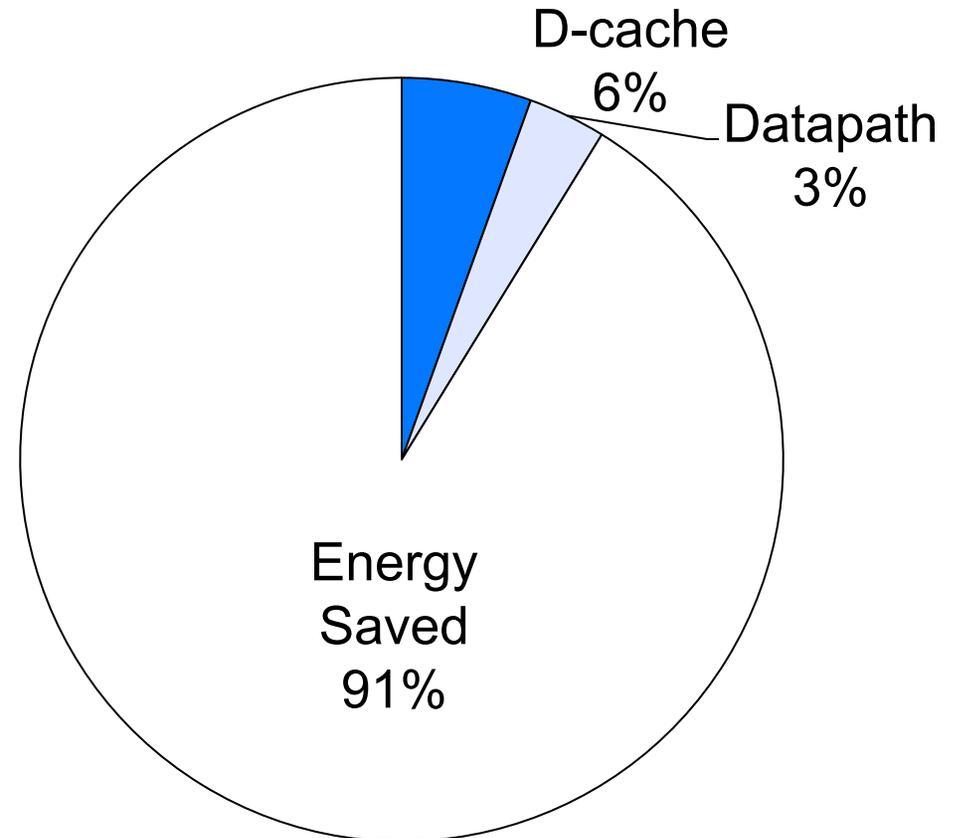
*0.01 mm<sup>2</sup> in 45 nm TSMC  
runs at 1.4 GHz*

# Typical Energy Savings



**RISC baseline**  
91 pJ/instr.

← ~11x →



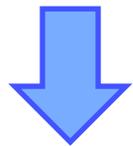
**C-cores**  
8 pJ/instr.

# “GreenDroid: A Mobile Application Processor for a Future of Dark Silicon”

*HOTCHIPS AUG 2010*

*IEEE Micro Mar 2011*

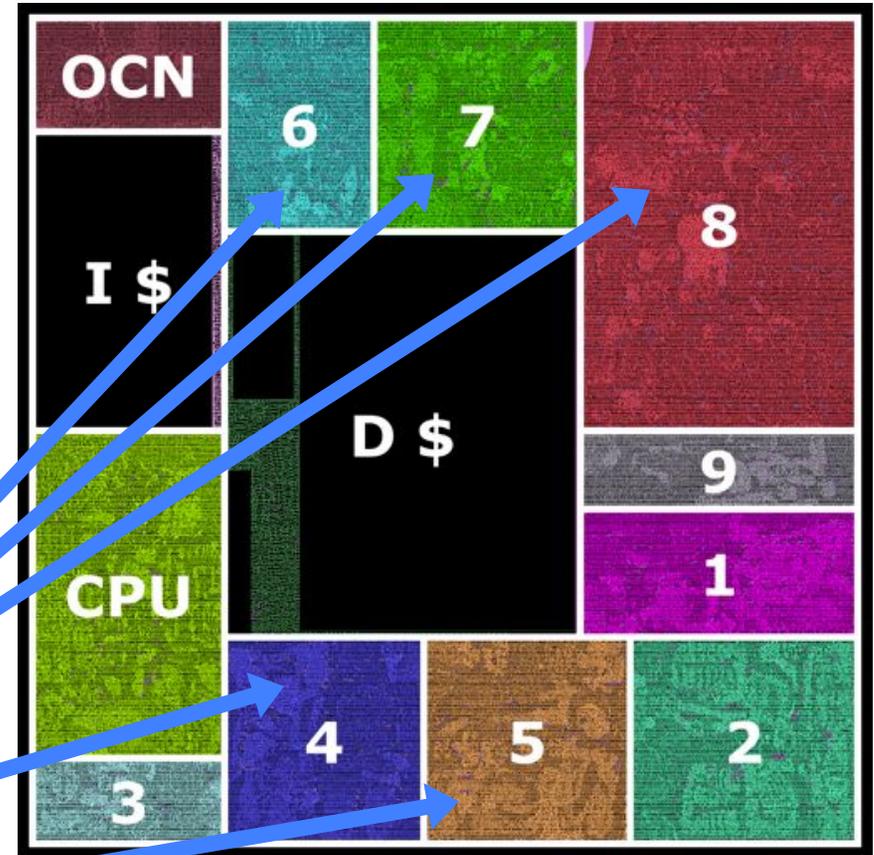
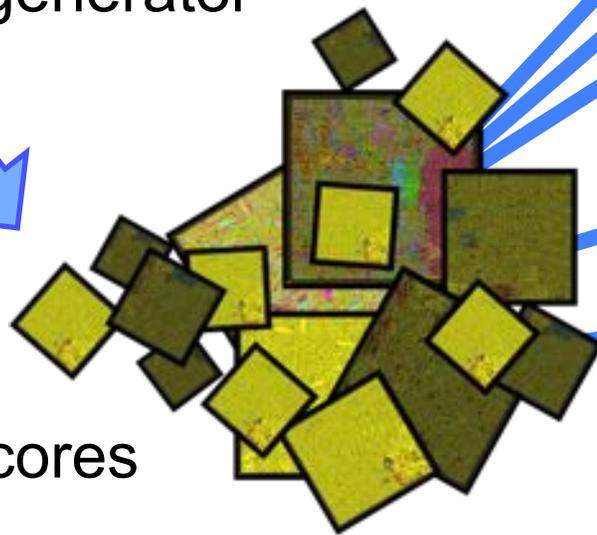
*ASPAC 2012*



Automatic  
c-core  
generator

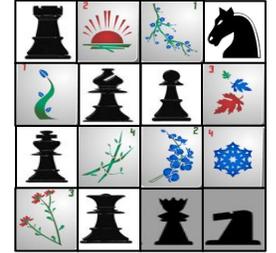


C-cores

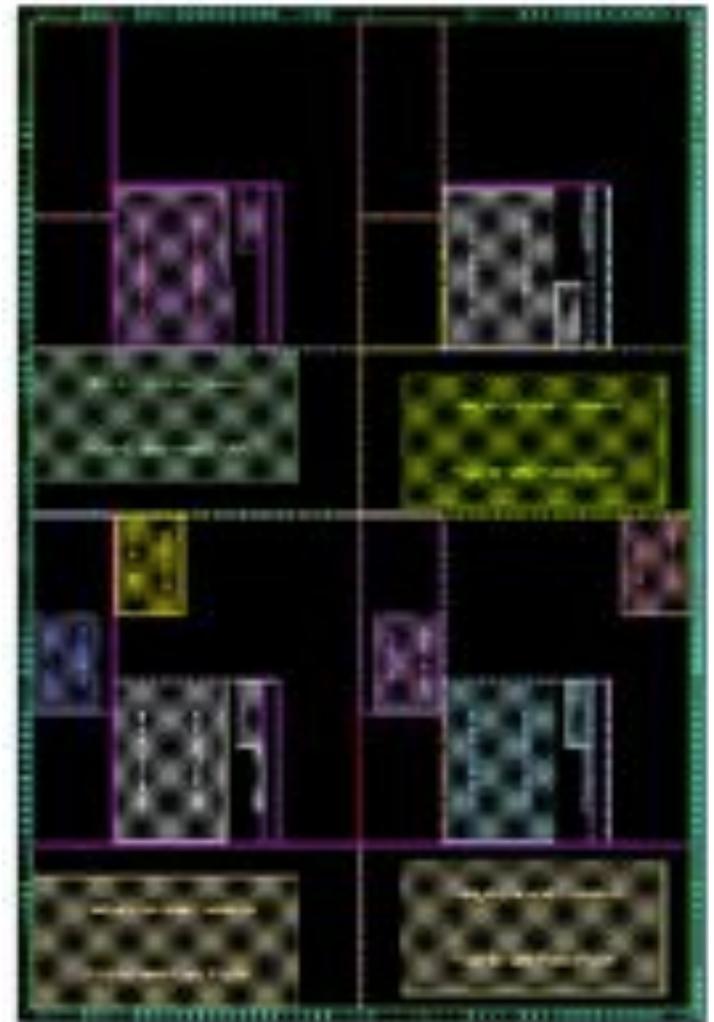


Placed-and-routed chip  
with 9 Android c-cores

# Quad-Core UCSD GreenDroid Prototype



- Four heterogeneous tiles with ~40 C-cores.
- Synopsys IC Compiler
- 28-nm Global Foundries
- ~1.5 GHz
- 2 mm<sup>2</sup>
- In backend/verification stages
- Multiproject Tapeout w/ UCSC  
November 2012



# The Four Horsemen

*The Dark Silicon Apocalypse*

*Explaining the Source of Dark Silicon*

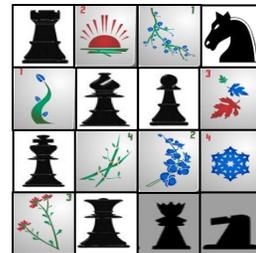
*The Four Horsemen*



I



II



III



IV

# The Deus Ex Machina Horseman

Latin        [/*dayus ex makeena*/]

American   [/*duece ex mashina*/]

**deus ex machina** /*dayus ex makeena*/

A plot device whereby a seemingly unsolvable problem is suddenly and abruptly solved with the unexpected intervention of some new event, character, ability or object.



# The Deus Ex Machina Horseman

*“MOSFETs are the fundamental problem.”*

**We can switch to FinFets, Trigate, High-K, nanotubes, 3D, for one-time improvements, but none are sustainable solutions across process generations.**

**Device physics (“thermionic emission of carriers across a potential well”) limit MOSFETS to 60 mV/decade subthreshold slope, which means the leakage problem is always there..”**



# The Deus Ex Machina Horseman

## Possible “Beyond CMOS” Device Directions (none are there yet, imho)

- Nano-electrical Mechanical Relays

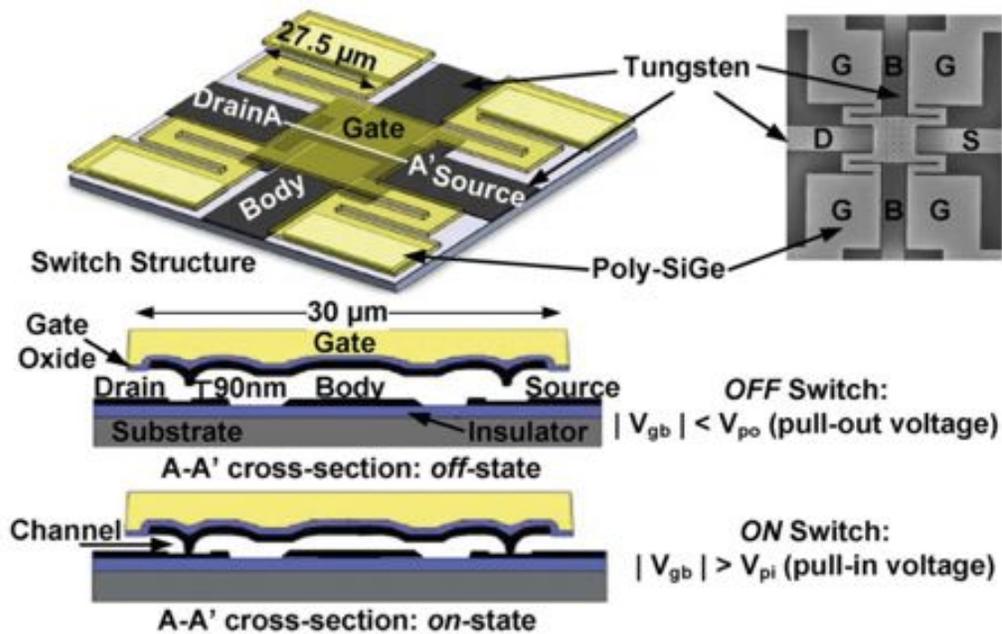


Fig. 1. SEM, diagram, and operating states of the MEM relay device.

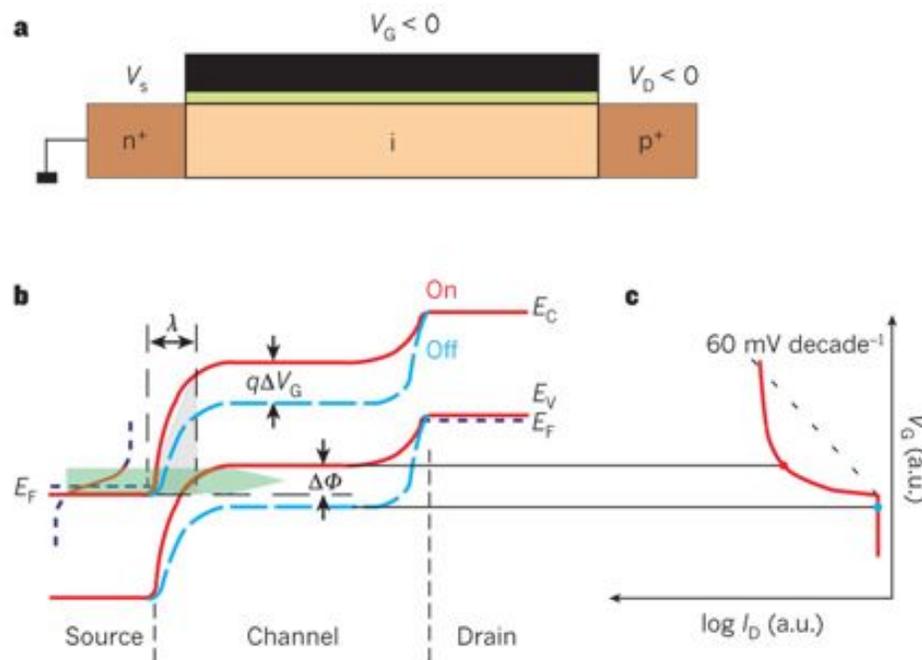
[e.g, Spencer et al JSSC 2011]



# The Deus Ex Machina Horseman

## “Beyond CMOS” Device Directions

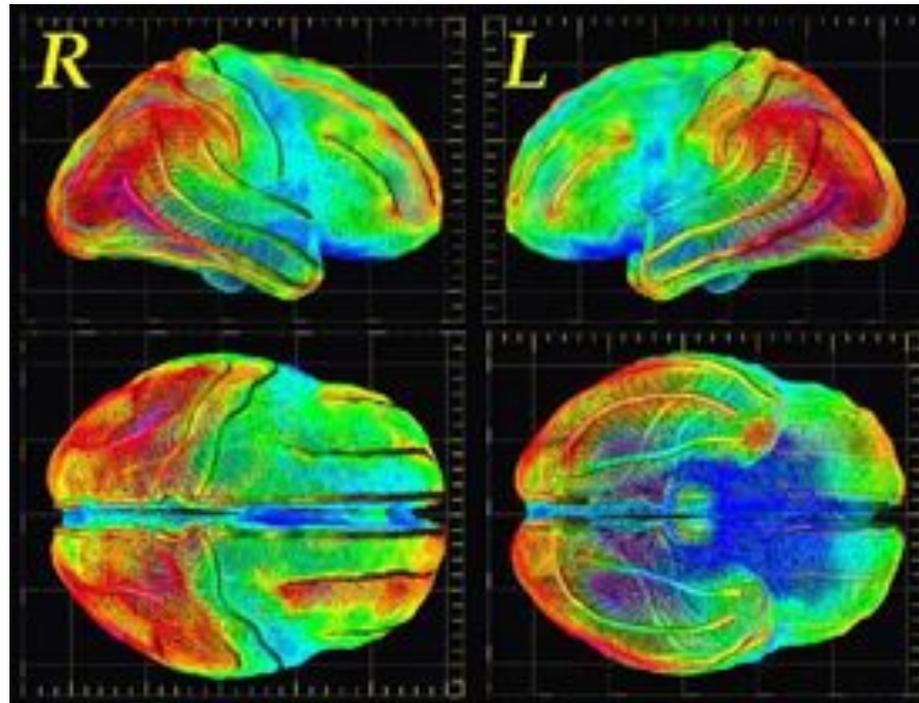
- Tunnel Field Effect Transistors (TFETs) [e.g., Ionescu et al, Nature 2011]
- Use Tunneling Effects to overcome MOSFET Limits



# The Deus Ex Machina Horseman (“*Before CMOS*” Directions)

- Human Brain

- 100 trillion synapses @ 20 W!
- Very “dark” circuits



# The Four Horsemen

*The Dark Silicon Apocalypse*

*Explaining the Source of Dark Silicon*

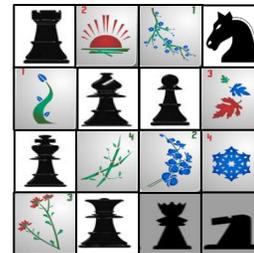
*The Four Horsemen*



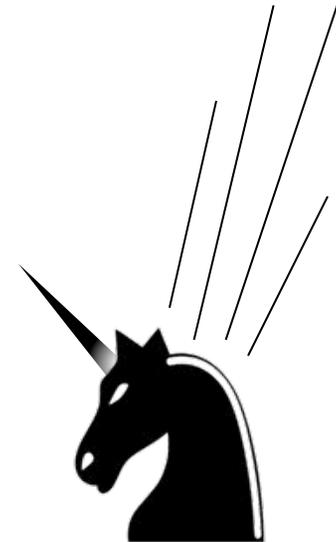
I



II



III



IV

# Conclusion

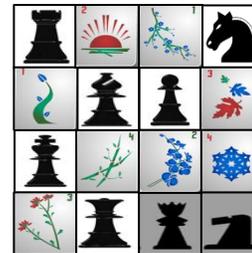
- Dark Silicon is opening up a whole new class of exciting new architectural directions which many folks are starting to move into – which I have termed the “four horsemen”.
- Probably the final answers will be a heterogeneous combination of all of these.
- Excited to see even more new ideas today!



I



II



III



IV

**darksilicon.org/horsemen**  
**for more details (also, 2012 DAC)**

*You are already attending the*  
Dark Silicon Workshop (DaSI) at ISCA 2012

*So, submit to the*  
IEEE Micro Special Issue on Dark Silicon!

[mbtaylor@ucsd.edu](mailto:mbtaylor@ucsd.edu)