

# Exponential Challenges, Exponential Rewards— The Future of Moore's Law

**Shekhar Borkar**

**Intel Fellow**

**Circuit Research, Intel Labs**

**Fall, 2004**

ISSCC 2003—  
Gordon Moore said...

**“No exponential is forever...**

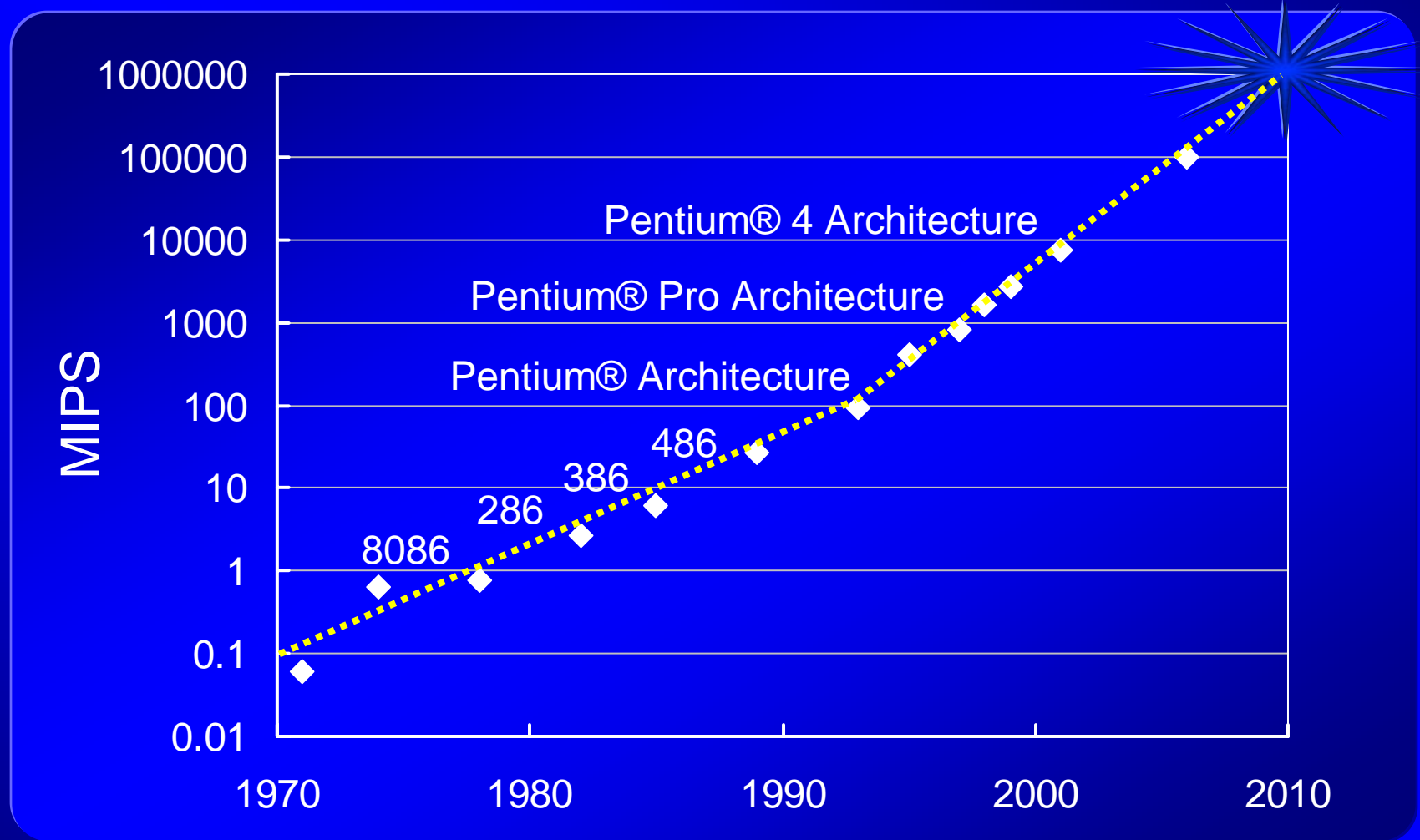
**But**

**We can delay *Forever*”**

# Outline

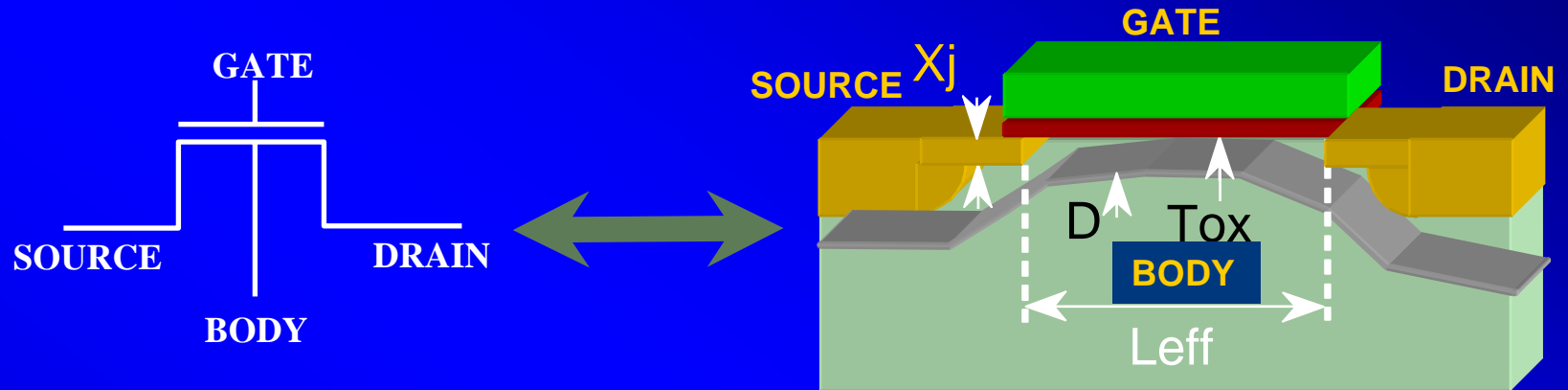
- The *exponential* challenges
- Circuit and  $\mu$ Arch solutions
- Major paradigm shifts in design
- Integration & SOC
- The *exponential* reward
- Summary

# Goal: 1TIPS by 2010



**How do you get there?**

# Technology Scaling



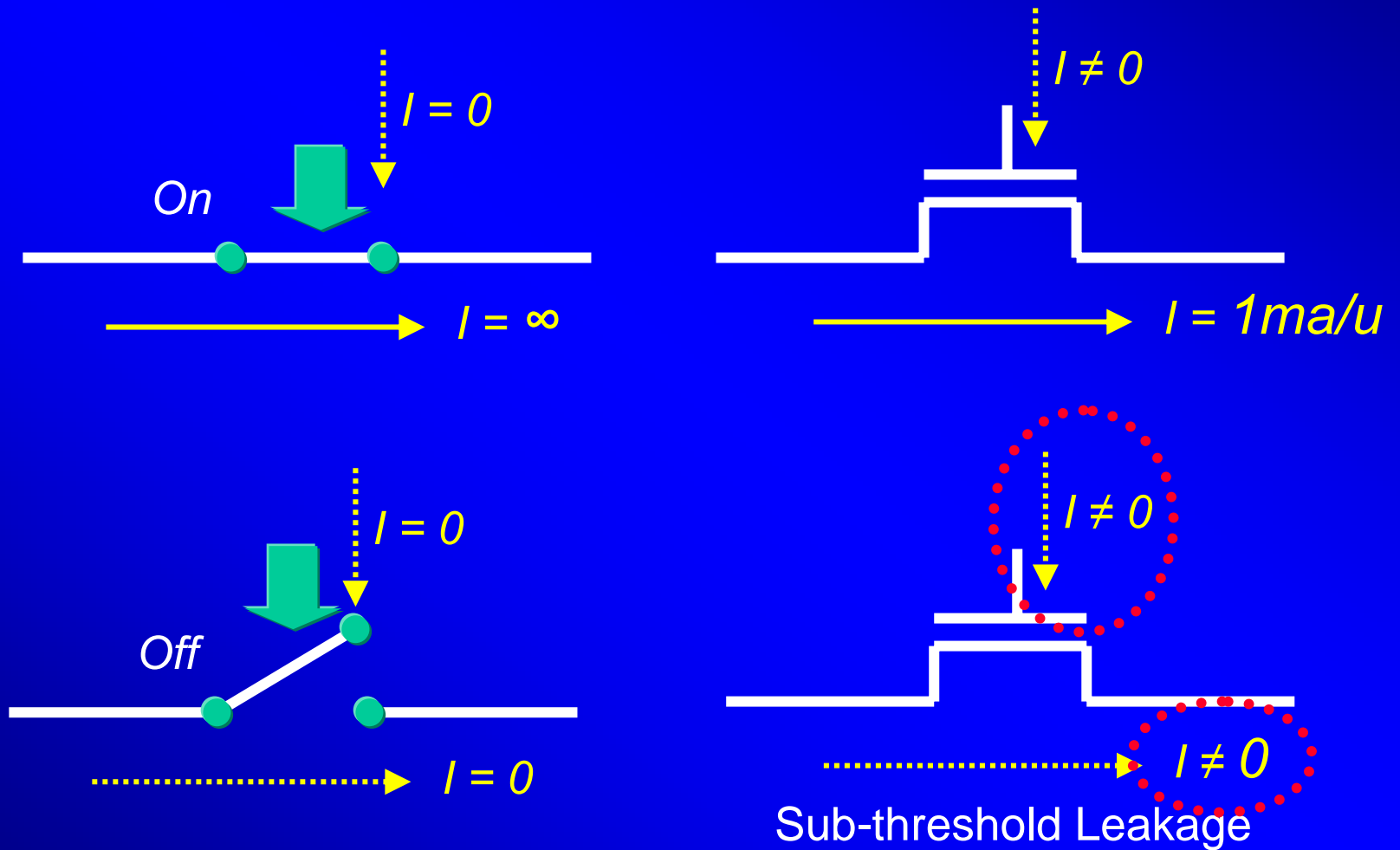
<b>Dimensions scale down by 30%</b>	<b>Doubles transistor density</b>
<b>Oxide thickness scales down</b>	<b>Faster transistor, higher performance</b>
<b>Vdd &amp; Vt scaling</b>	<b>Lower active power</b>

**Technology has scaled well, will it in the future?**

# Technology Outlook

High Volume Manufacturing	2004	2006	2008	2010	2012	2014	2016	2018
Technology Node (nm)	90	65	45	32	22	16	11	8
Integration Capacity (BT)	2	4	8	16	32	64	128	256
Delay = CV/I scaling	0.7	~0.7	>0.7	Delay scaling will slow down				
Energy/Logic Op scaling	>0.35	>0.5	>0.5	Energy scaling will slow down				
Bulk Planar CMOS	High Probability				Low Probability			
Alternate, 3G etc	Low Probability				High Probability			
Variability	Medium			High		Very High		
ILD (K)	~3	<3	Reduce slowly towards 2-2.5					
RC Delay	1	1	1	1	1	1	1	1
Metal Layers	6-7	7-8	8-9	0.5 to 1 layer per generation				

# Is Transistor a Good Switch?

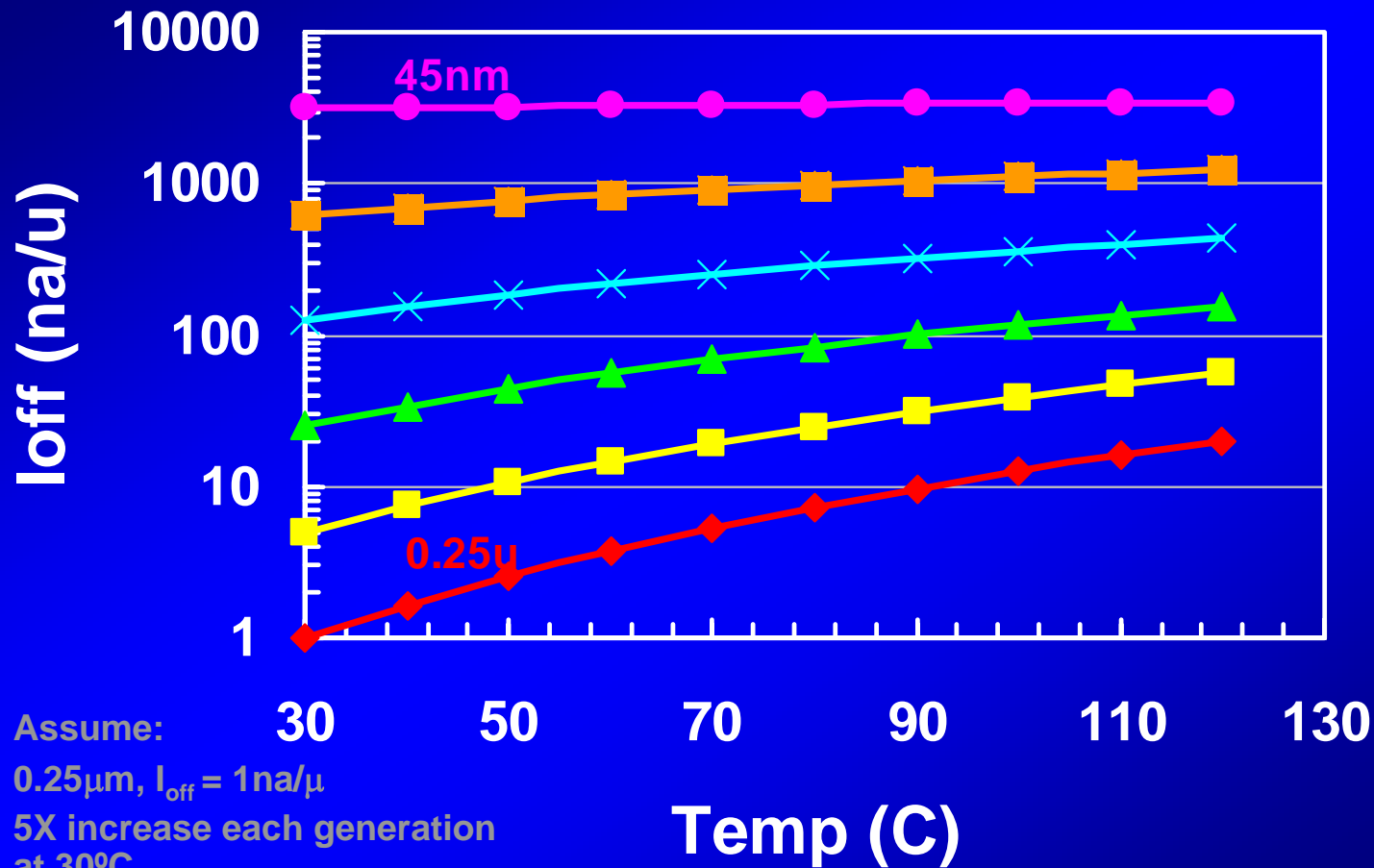


# Exponential Challenge #1

**Subthreshold  
Leakage**

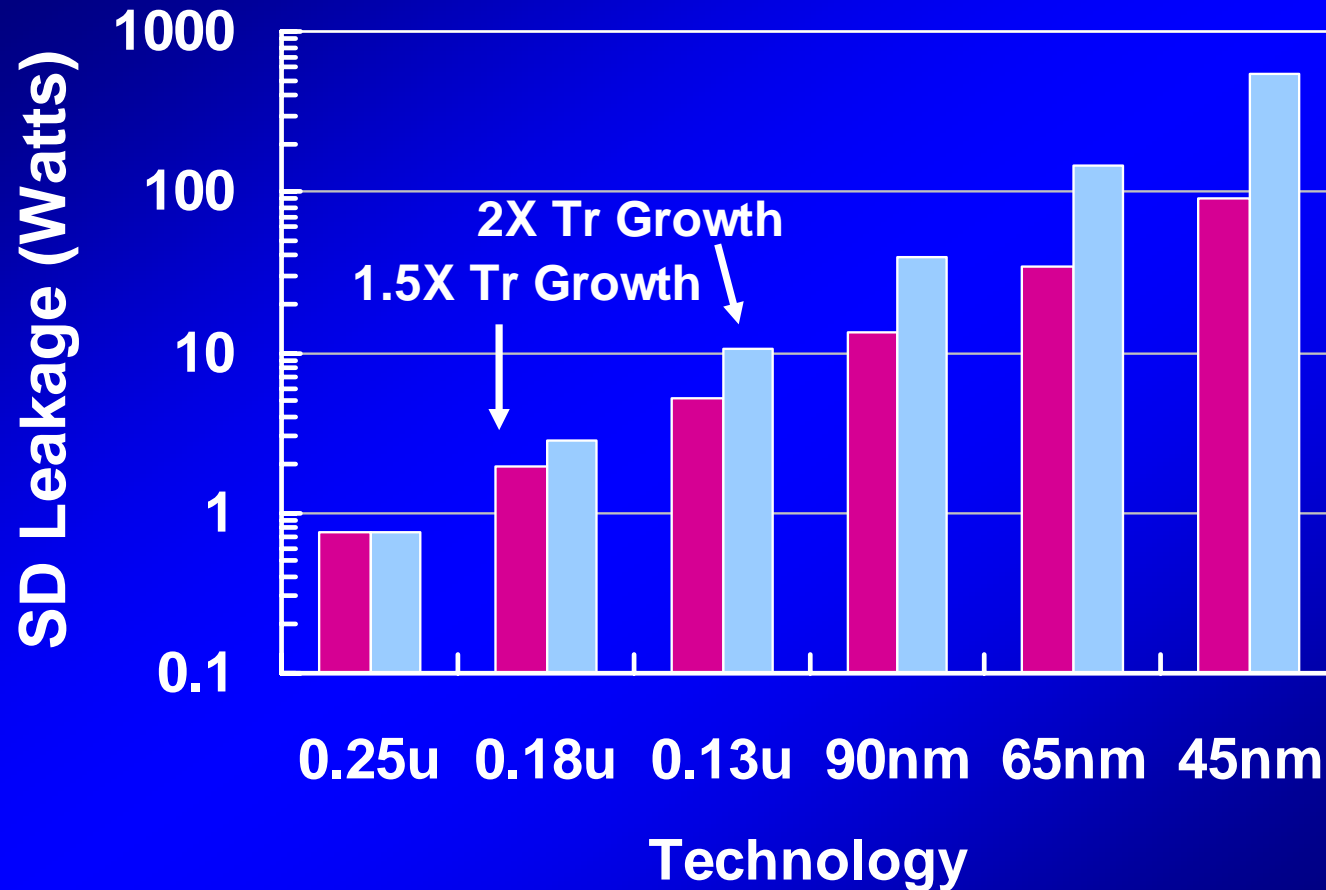


# Sub-threshold Leakage



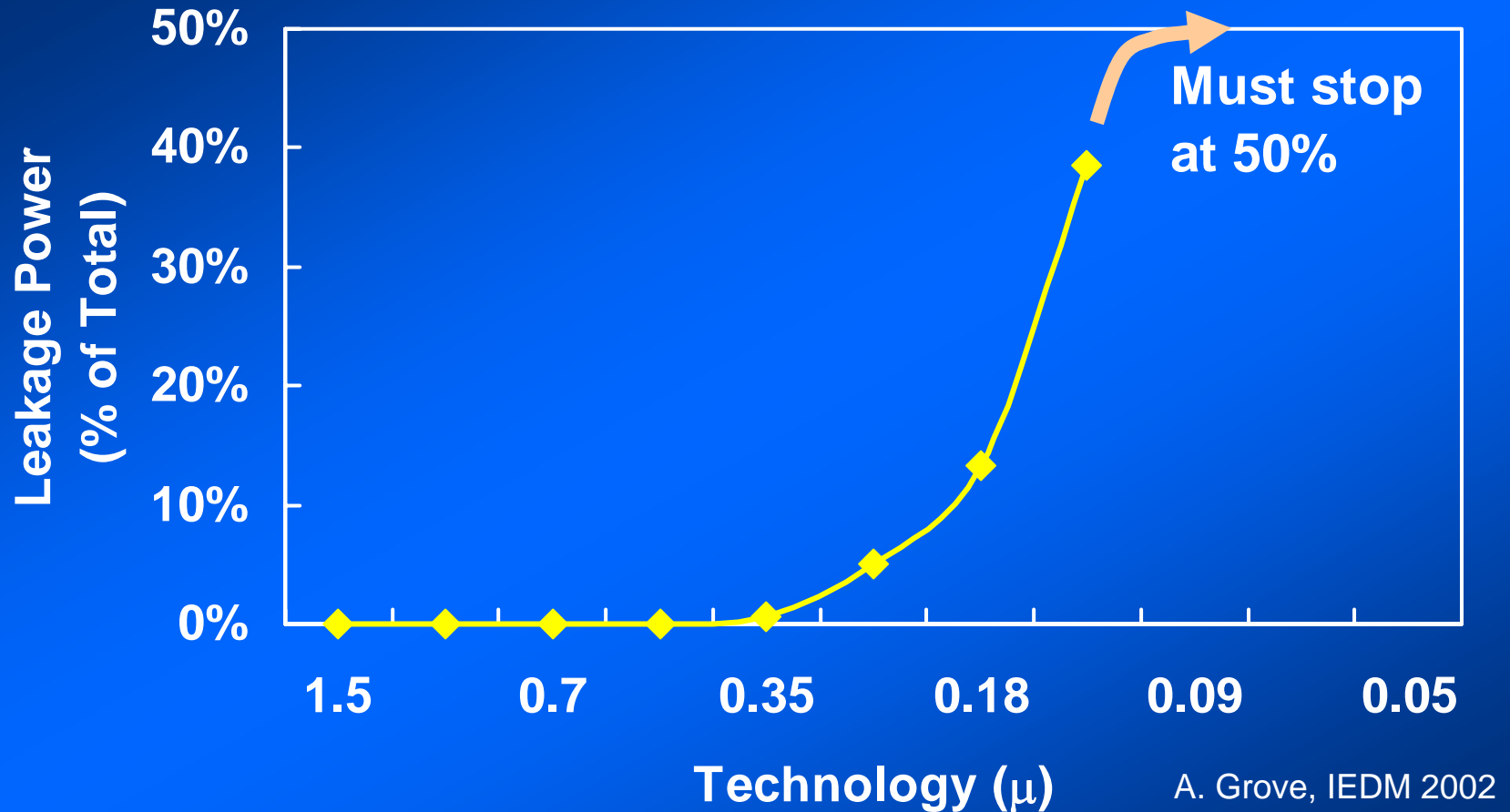
**Sub-threshold leakage increases exponentially**

# SD Leakage Power



**SD leakage power becomes prohibitive**

# Leakage Power

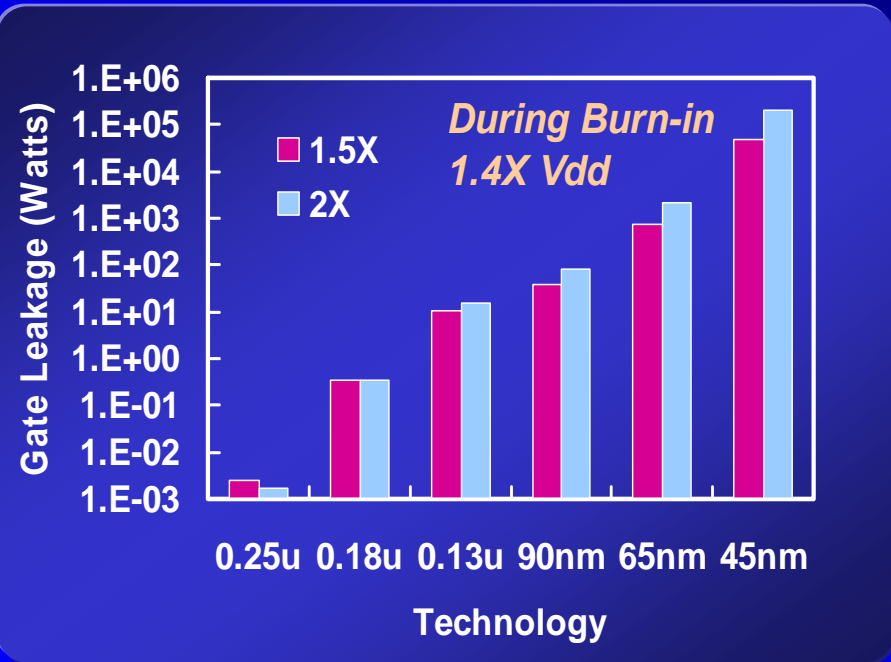
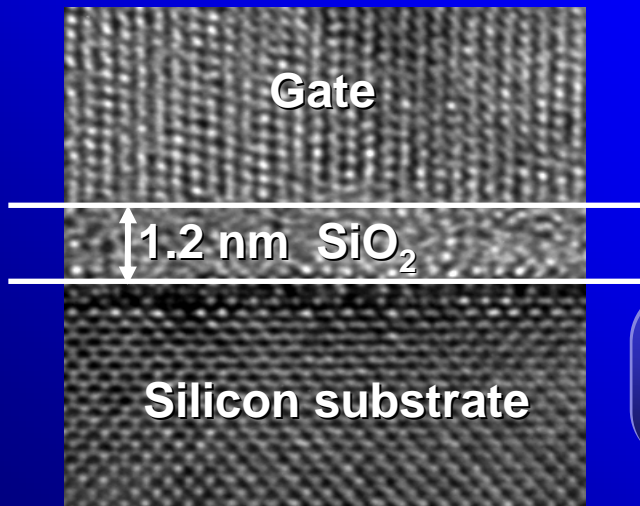
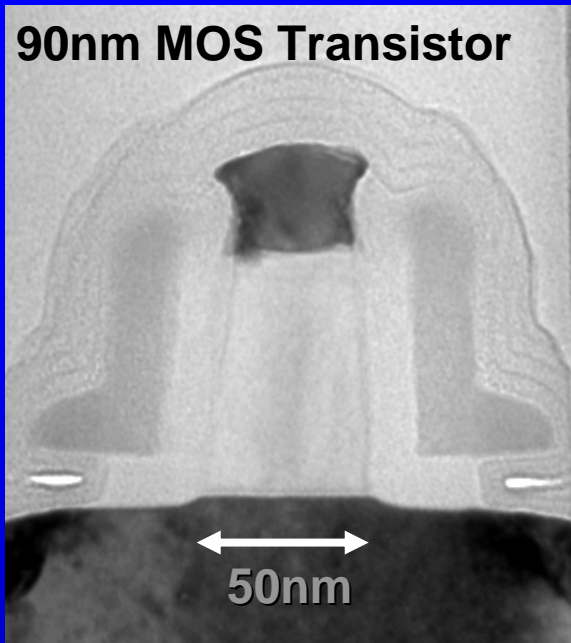


**Leakage power limits  $V_t$  scaling**

# Exponential Challenge #2

**Gate  
Leakage**

# Gate Oxide is Near Limit



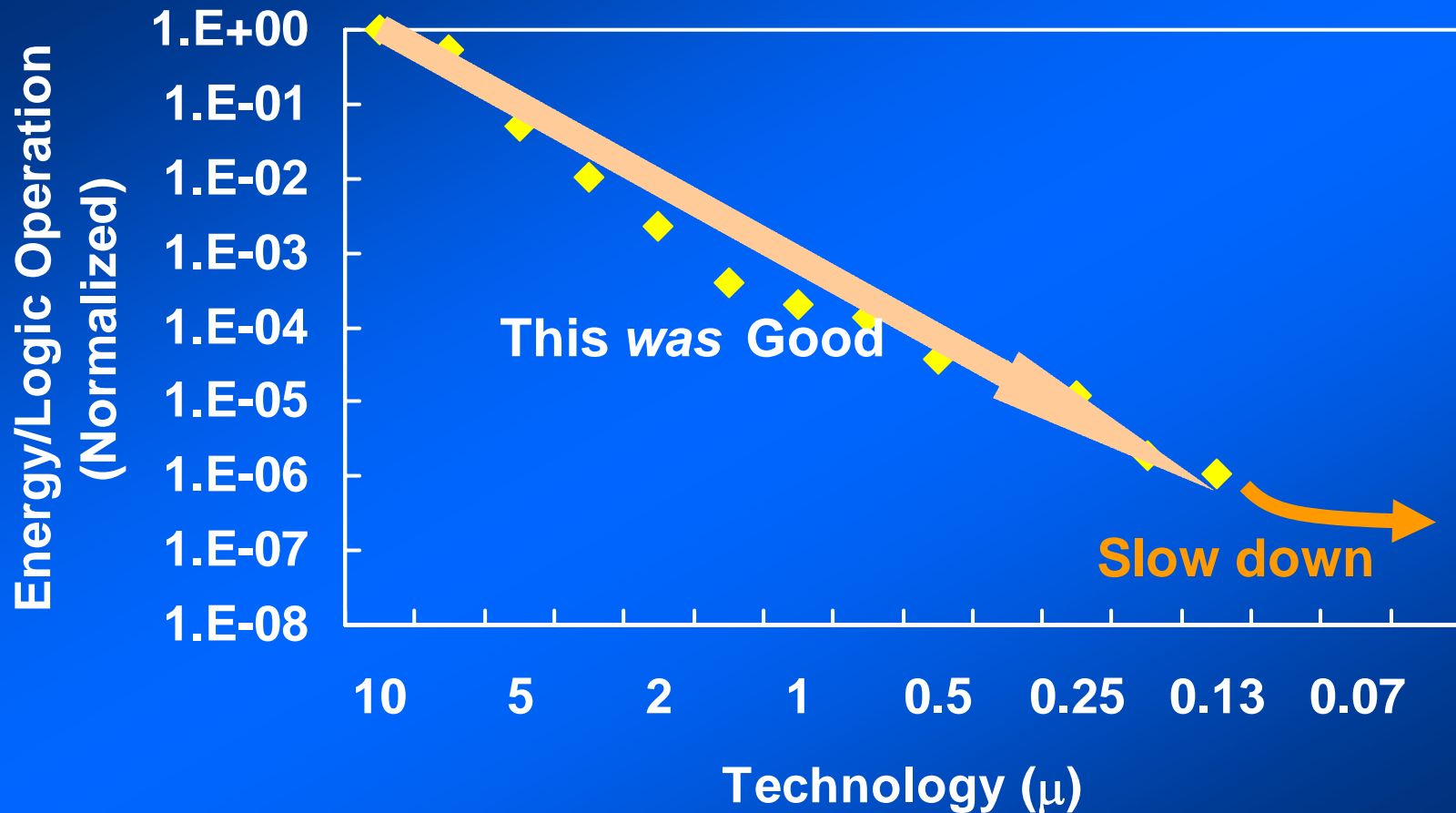
If  $T_{ox}$  scaling slows down, then Vdd scaling will have to slow down

High-K dielectric is crucial

# Exponential Challenge #3

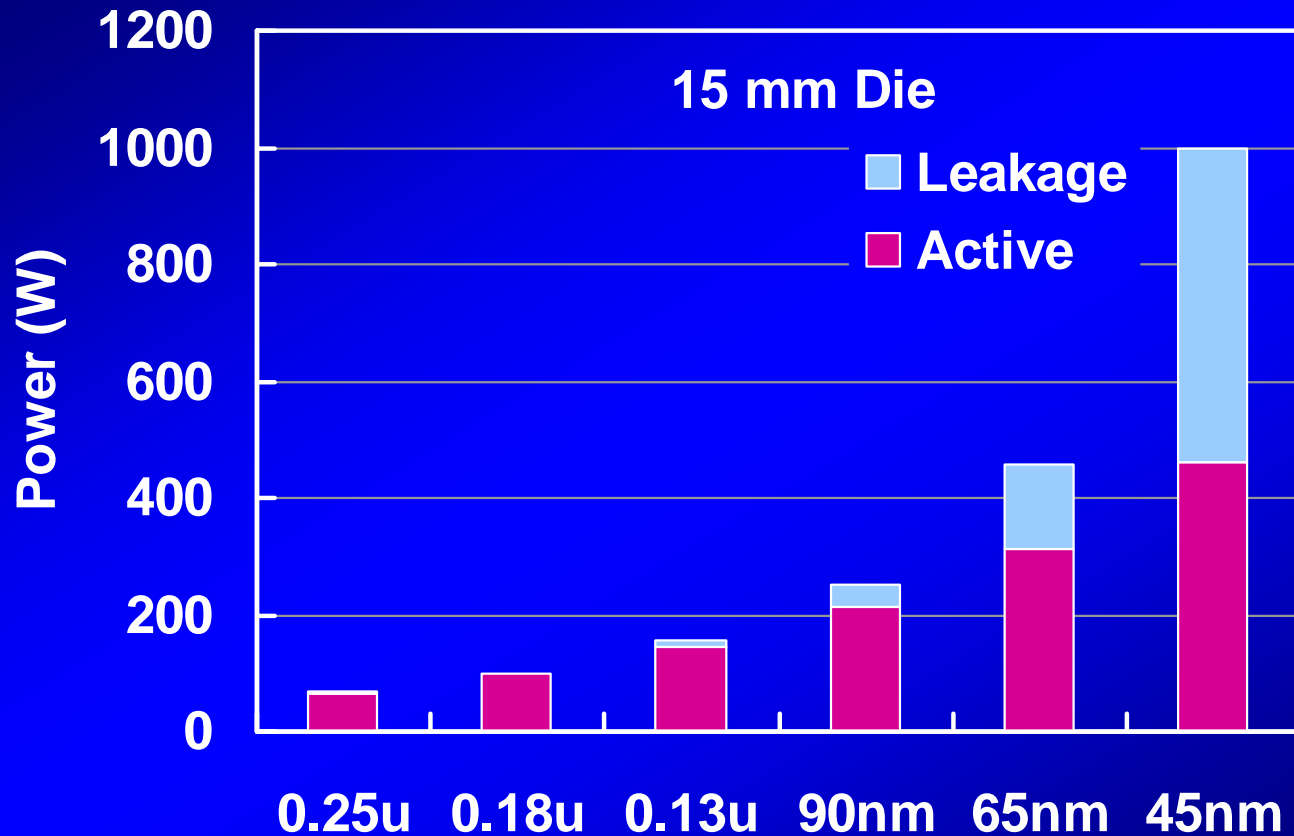
**ENERGY  
&  
POWER**

# Energy per Logic Operation



**Energy per logic operation scaling will slow down**

# The Power Crisis



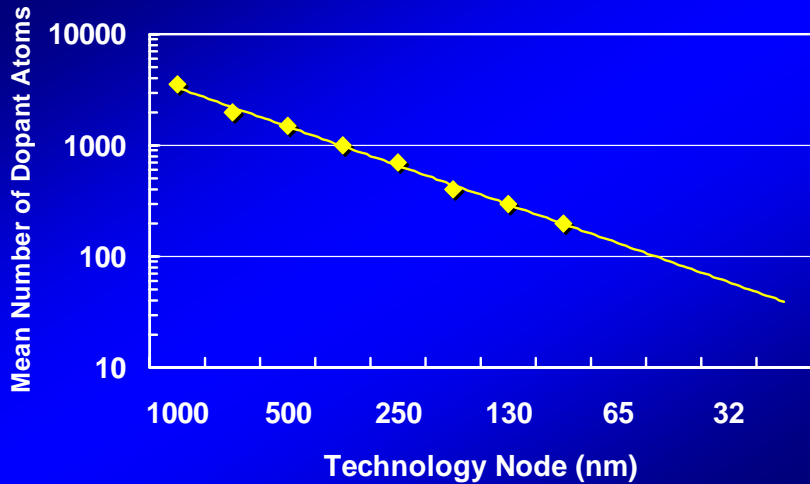
***Business as usual is not an option***



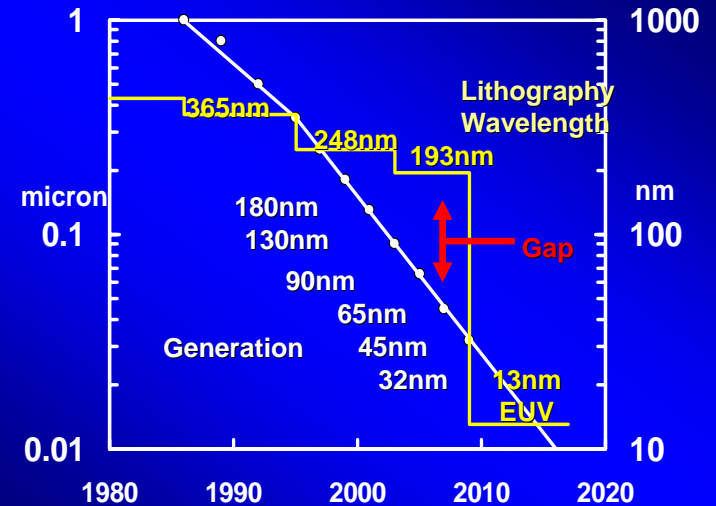
# Exponential Challenge #4

# Variations

# Sources of Variations

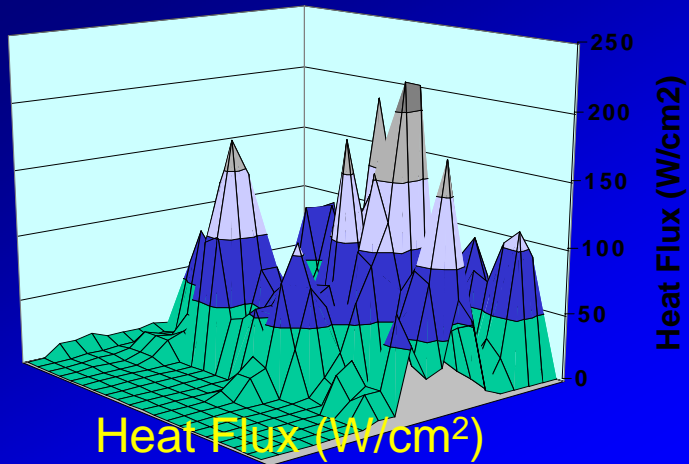


Random Dopant Fluctuations

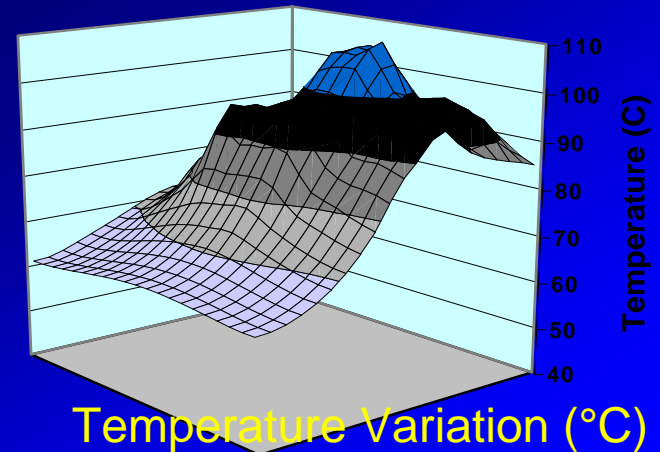


Sub-wavelength Lithography

Source: Mark Bohr, Intel

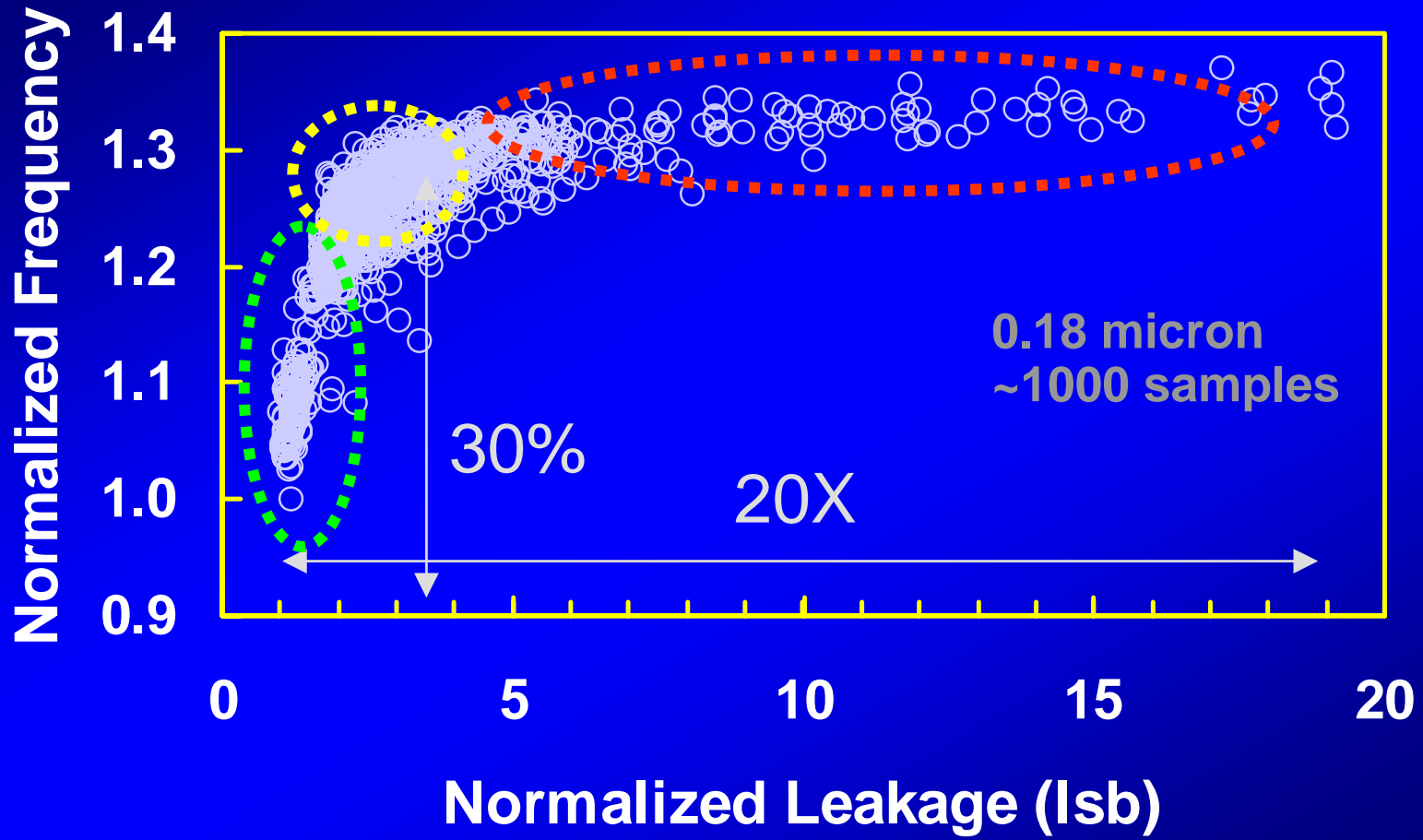


Heat Flux ( $W/cm^2$ )  
Results in Vcc variation



Temperature Variation ( $^{\circ}C$ )  
Hot spots

# Frequency & SD Leakage

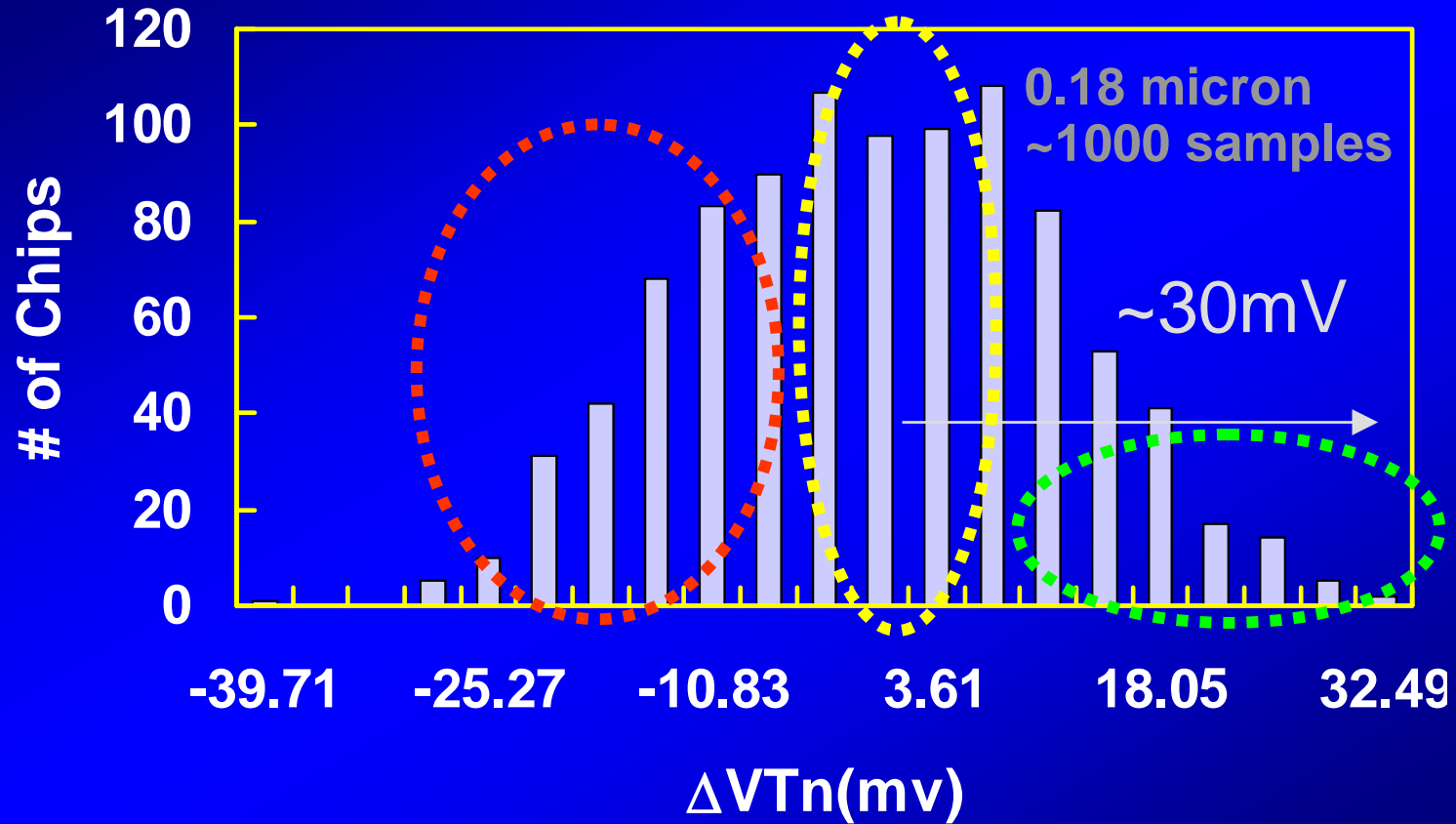


Low Freq  
Low Isb

High Freq  
Medium Isb

High Freq  
High Isb

# Vt Distribution



High Freq  
High Isb

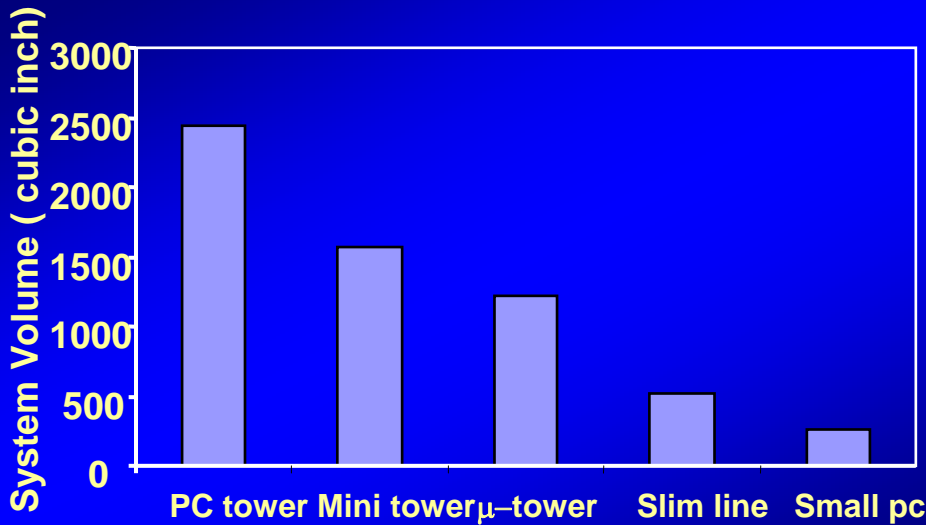
High Freq  
Medium Isb

Low Freq  
Low Isb

# Exponential Challenge #5

**Platform  
&  
System**

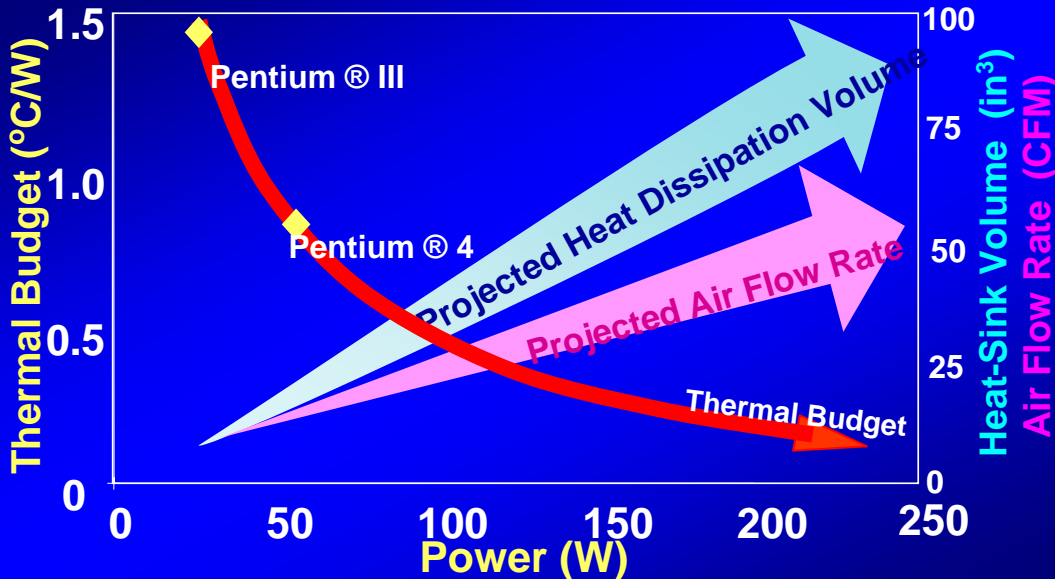
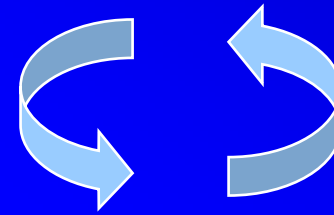
# Platform Requirements



Shrinking volume

Quieter

Yet, High Performance



Thermal budget decreasing

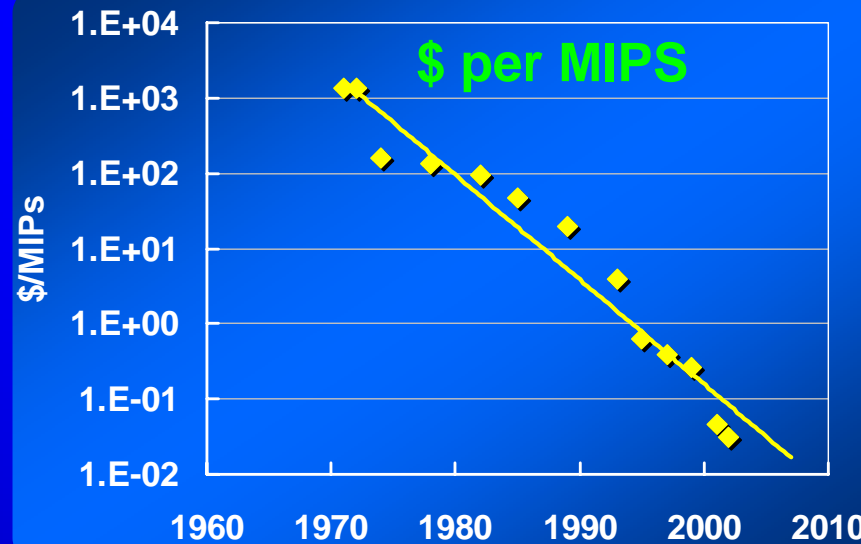
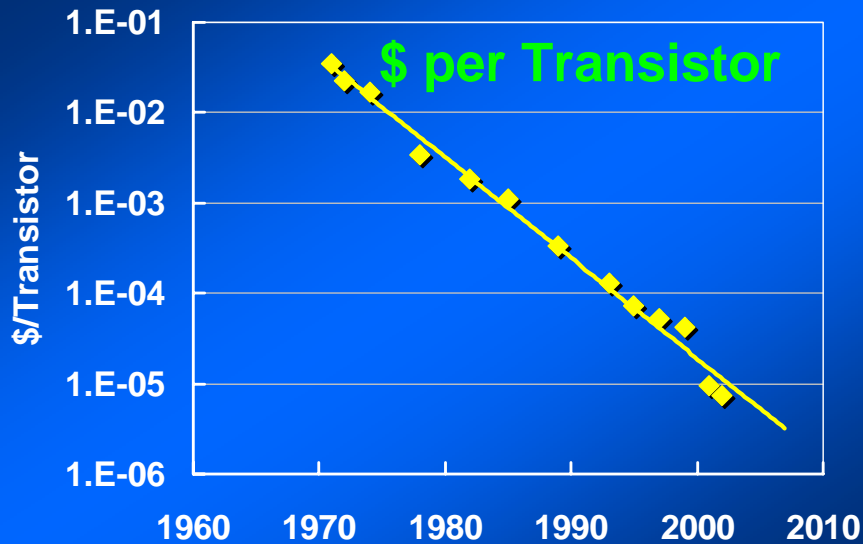
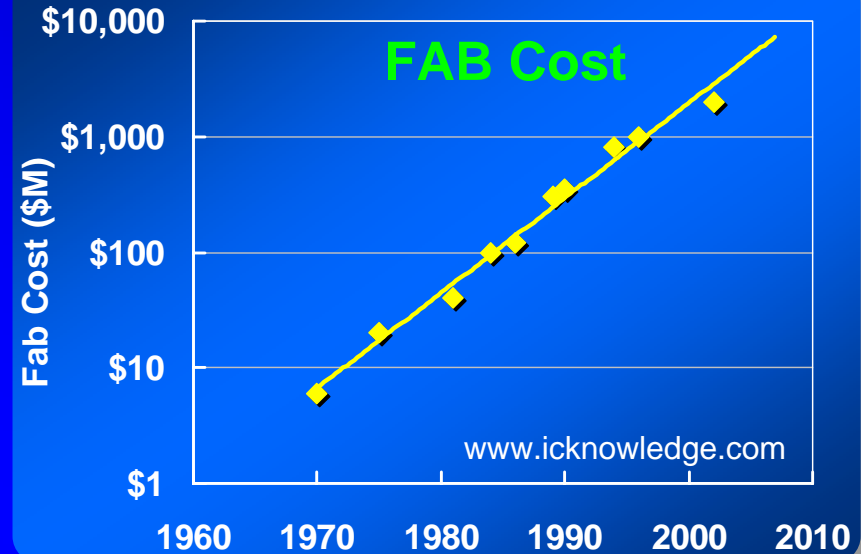
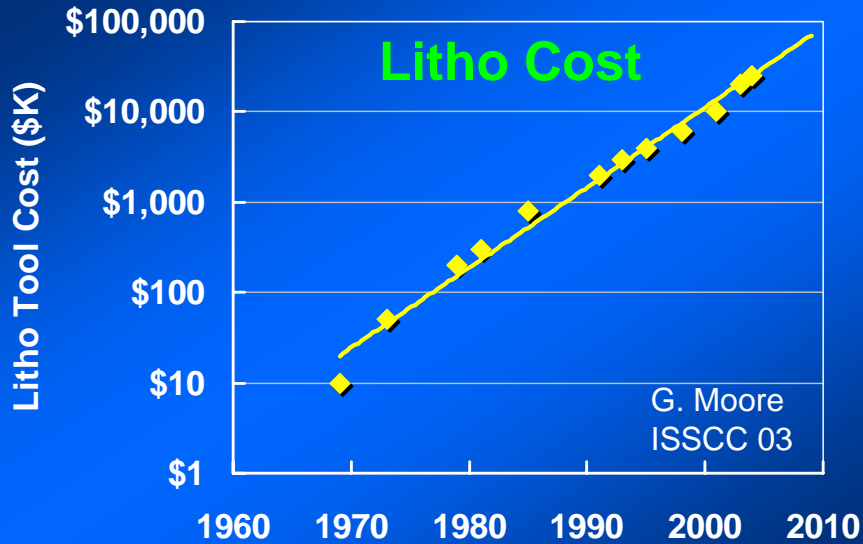
Higher heat sink volume

Higher air flow rate

# Exponential Challenge #6

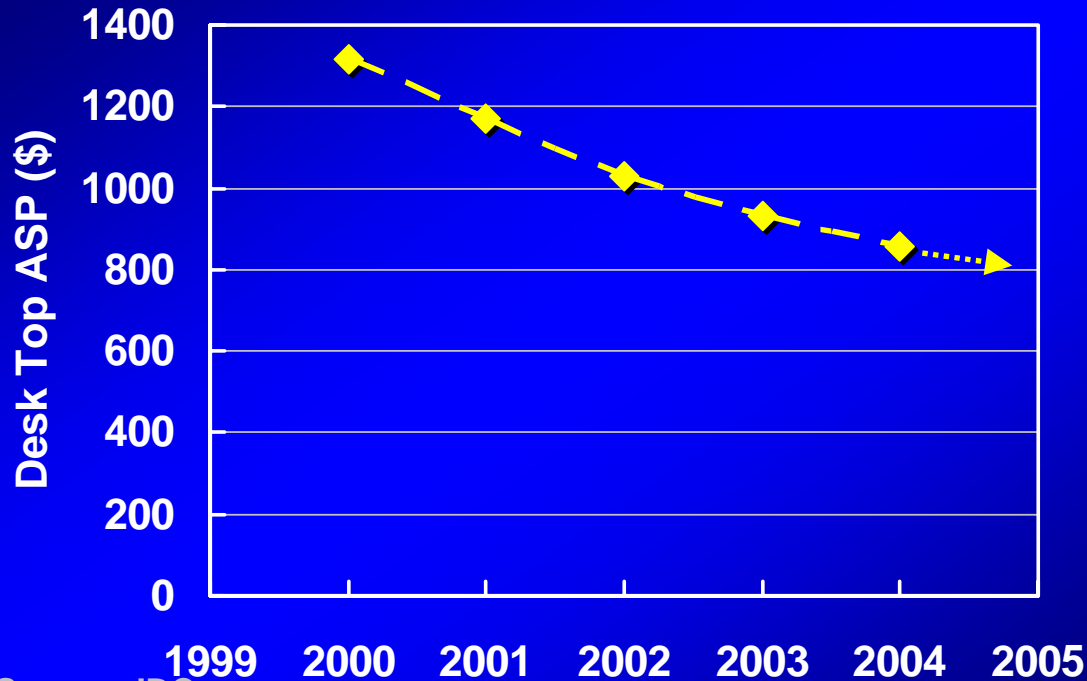
**Economics**

# Exponential Costs

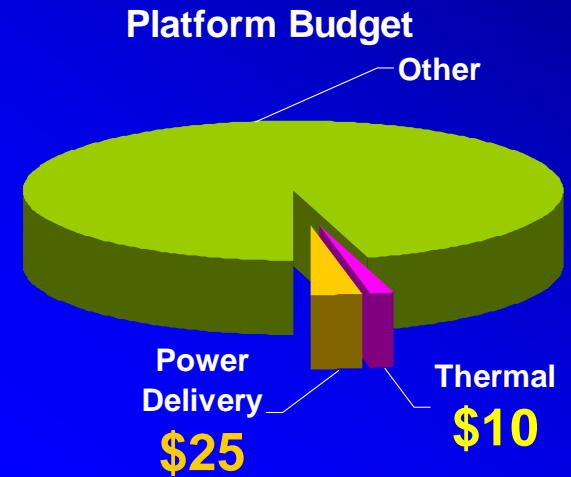




# Product Cost Pressure

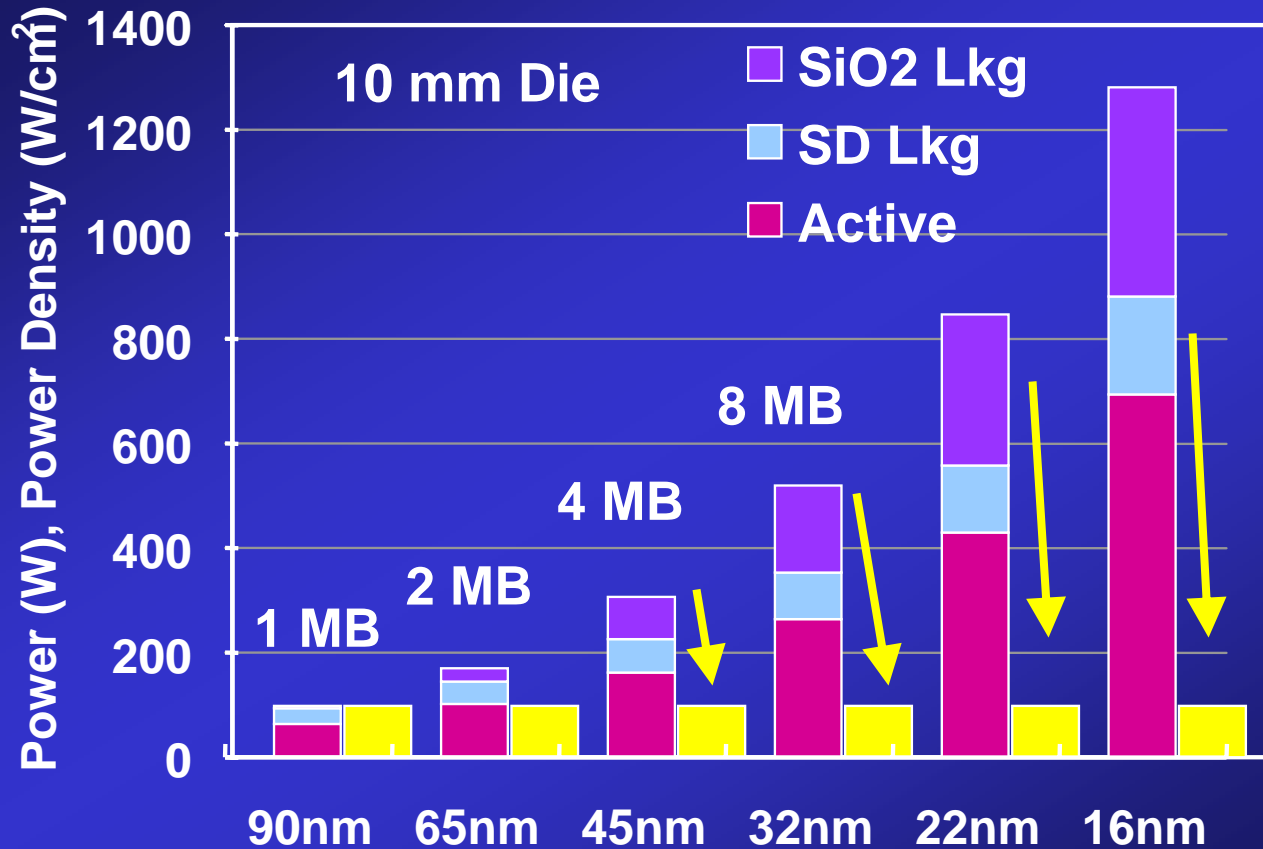


Source: IDC



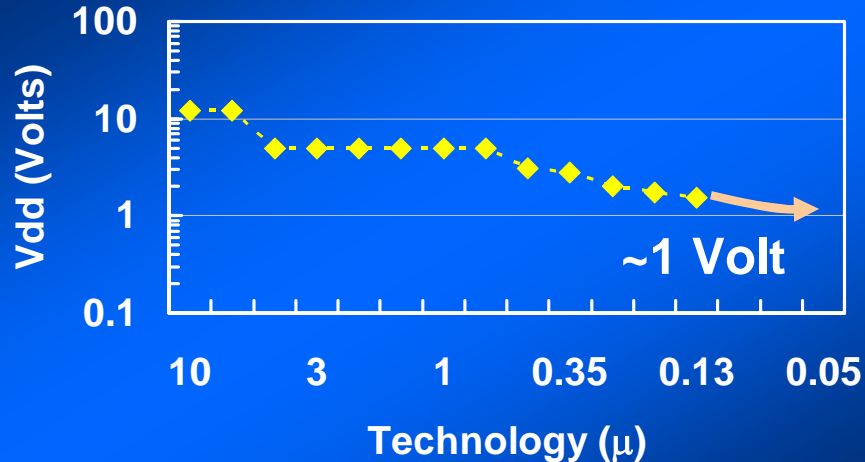
**Shrinking ASP, and shrinking \$ budget for power**

# Must Fit in Power Envelope

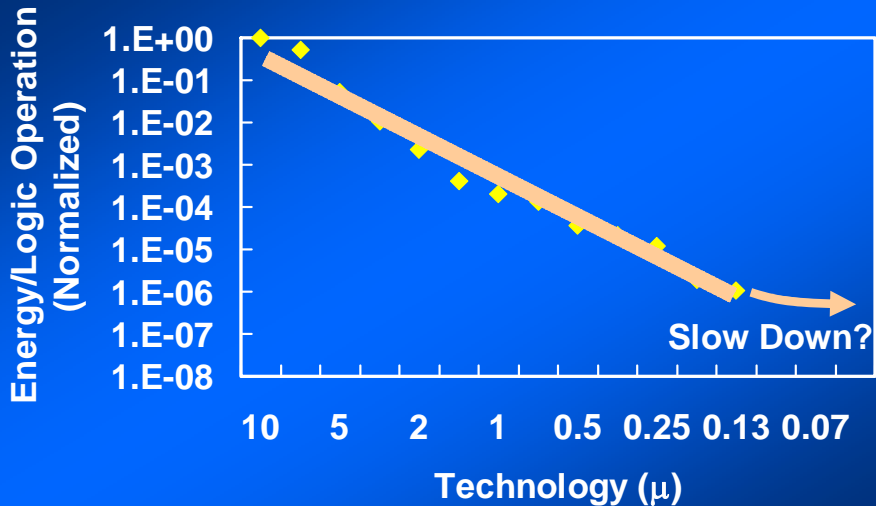


**Technology, Circuits, and  
Architecture to constrain the power**

# Some Implications



- $T_{ox}$  scaling will slow down—may stop?
- $V_{dd}$  scaling will slow down—may stop?
- $V_t$  scaling will slow down—may stop?
- Approaching constant  $V_{dd}$  scaling
- Energy/logic op will not scale



# The Gigascale Dilemma

- 1B T integration capacity will be available
- But could be unusable due to power
- Logic T growth will slow down
- Transistor performance will be limited

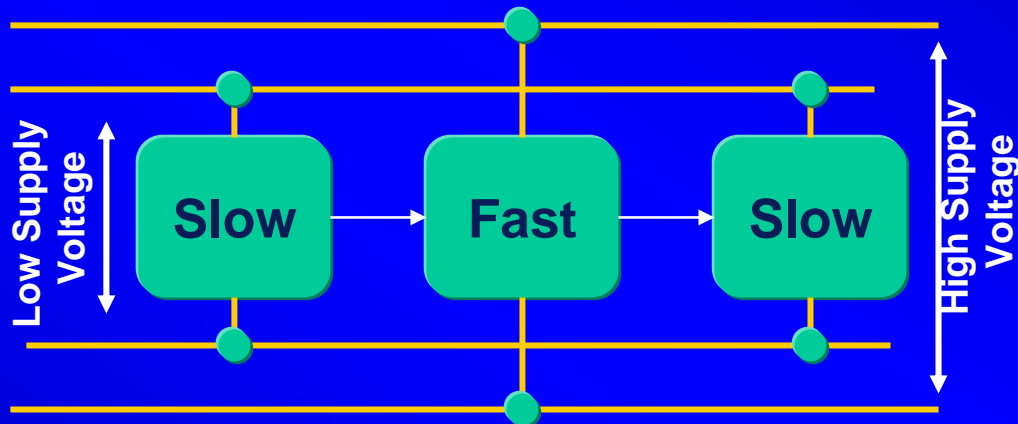
## *Solutions*

- Low power design techniques
- Improve design efficiency—Multi everywhere
- Valued performance by even higher integration (of potentially slower transistors)

# Solutions

**Power—active and leakage  
Variations  
Microarchitecture**

# Active Power Reduction



**Multiple Supply Voltages**

## Replicated Designs



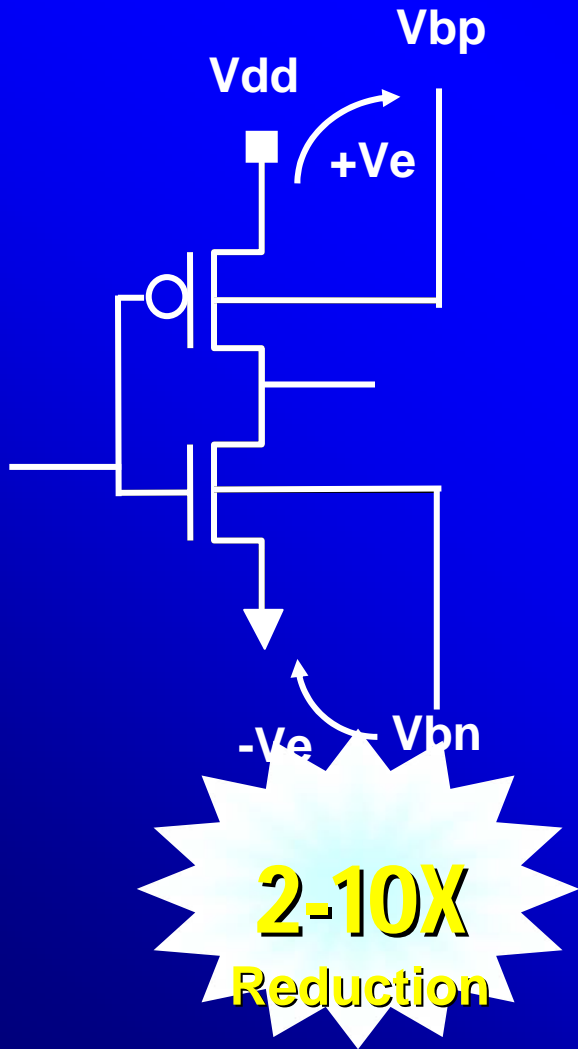
Freq = 1  
Vdd = 1  
Throughput = 1  
Power = 1  
Area = 1  
Pwr Den = 1



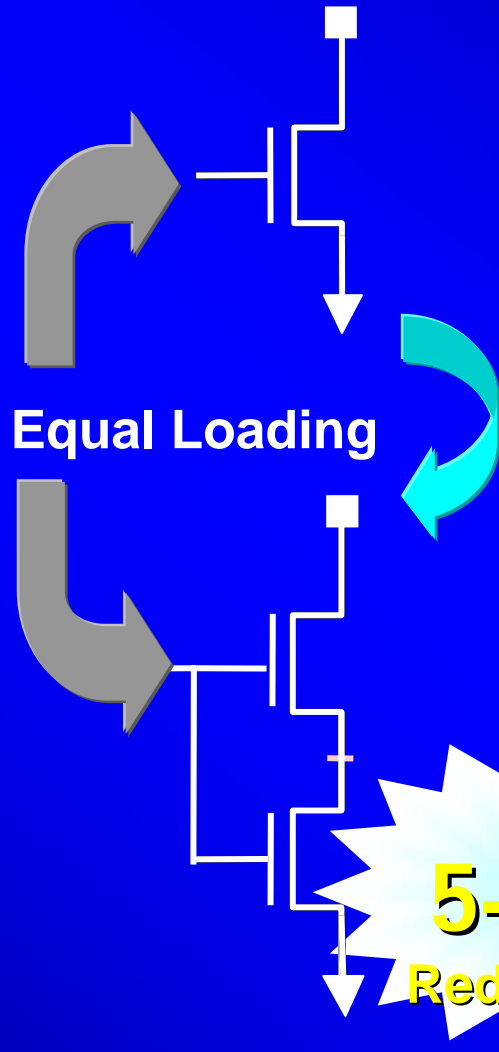
Freq = 0.5  
Vdd = 0.5  
Throughput = 1  
Power = 0.25  
Area = 2  
Pwr Den = 0.125

# Leakage Control

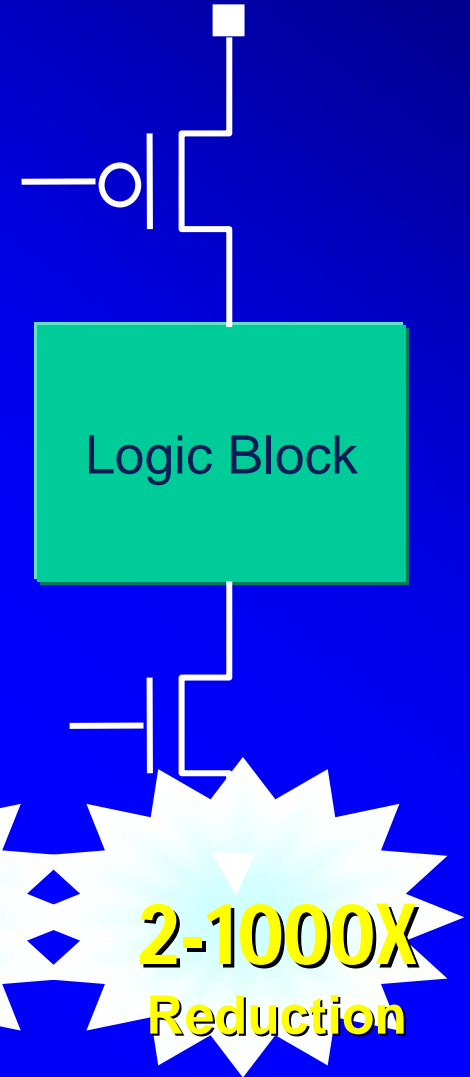
Body Bias



Stack Effect



Sleep Transistor



# Circuit Design Tradeoffs



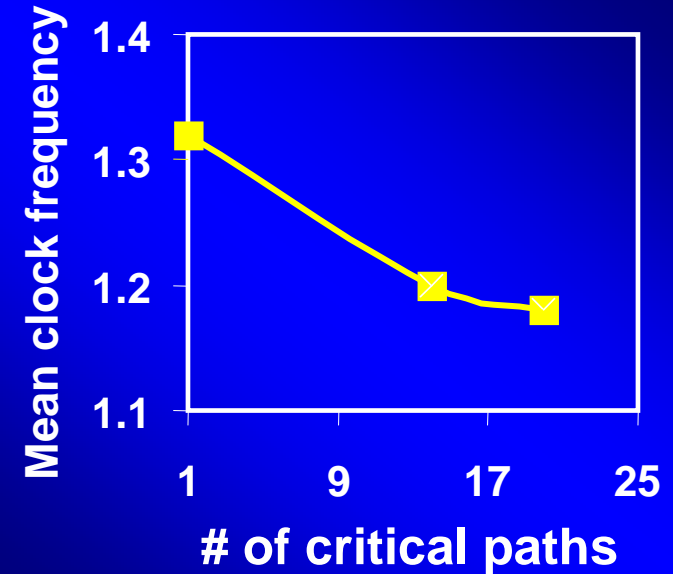
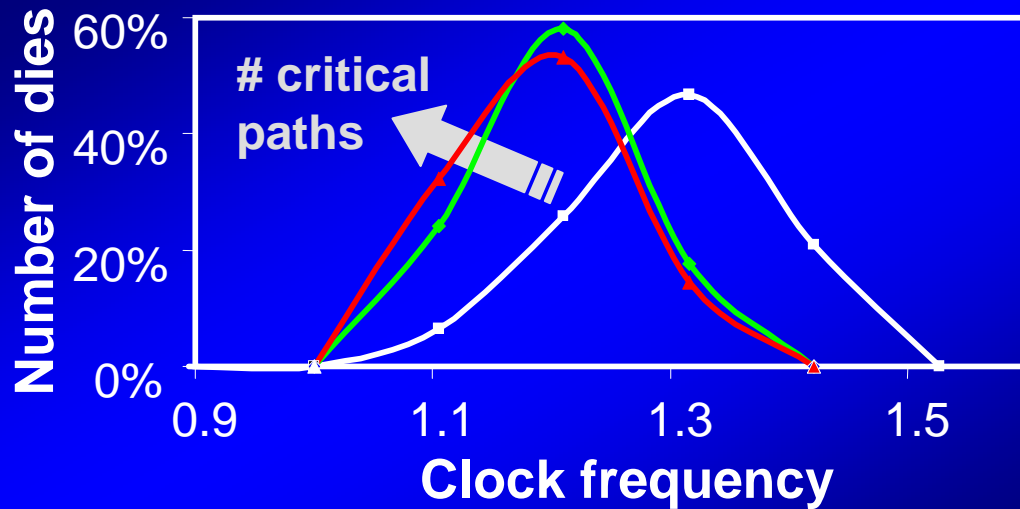
Higher probability of target frequency with:

1. Larger transistor sizes
2. Higher Low-Vt usage

But with power penalty

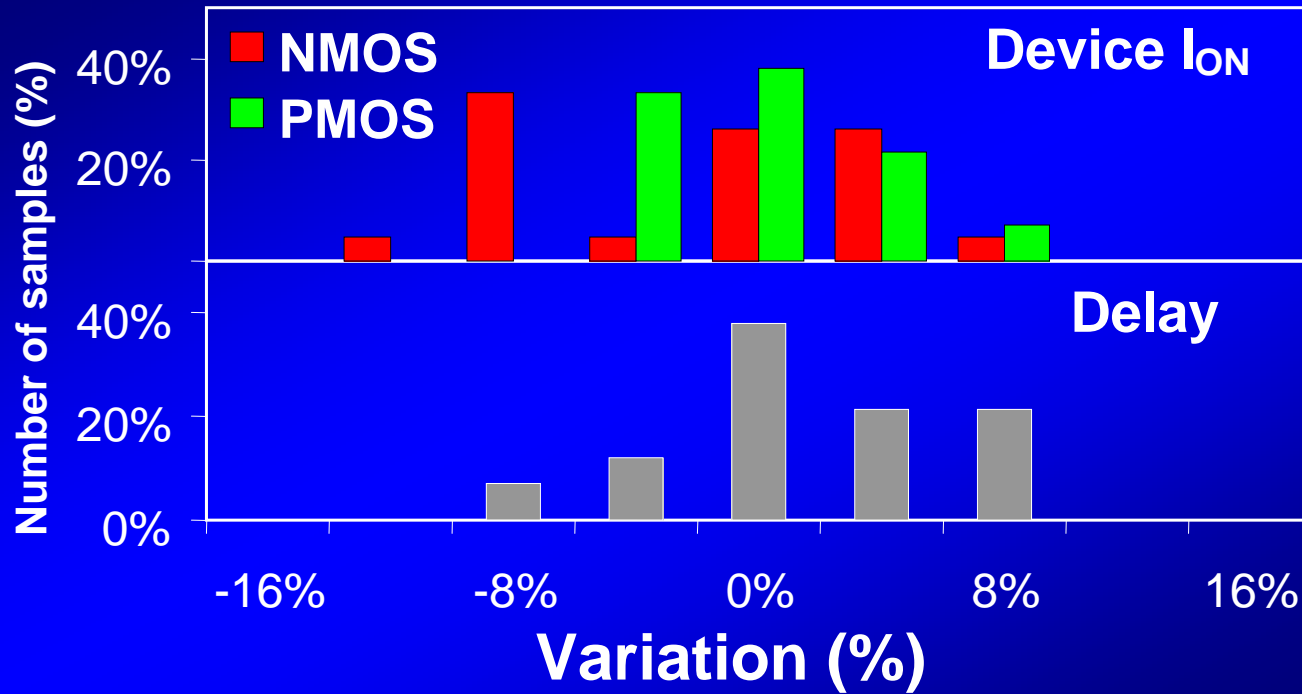


# Impact of Critical Paths

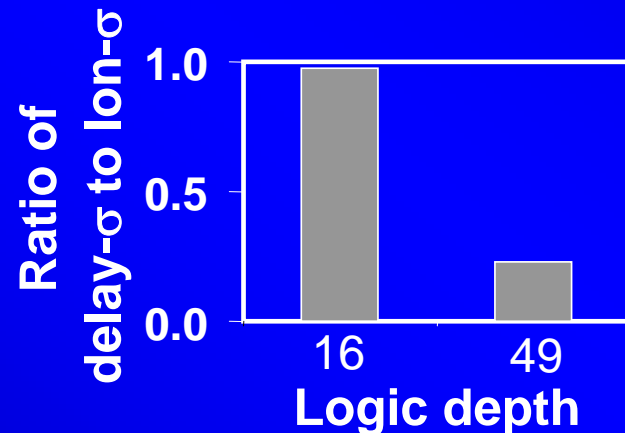


- **With increasing # of critical paths**
  - Both  $\sigma$  and  $\mu$  become smaller
  - Lower mean frequency

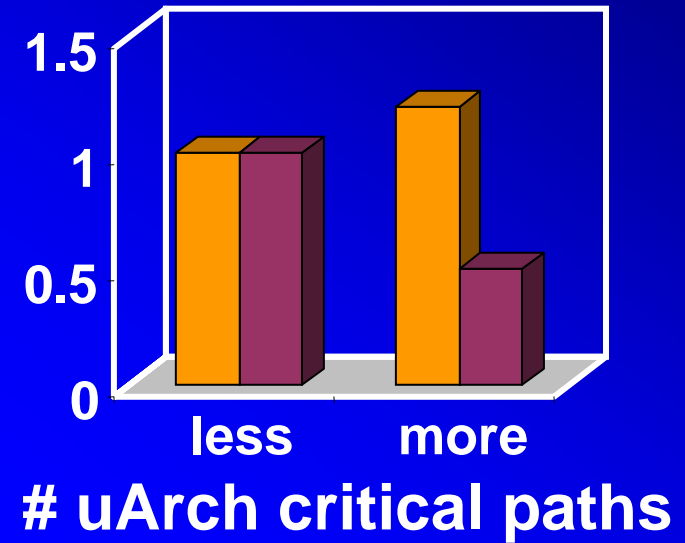
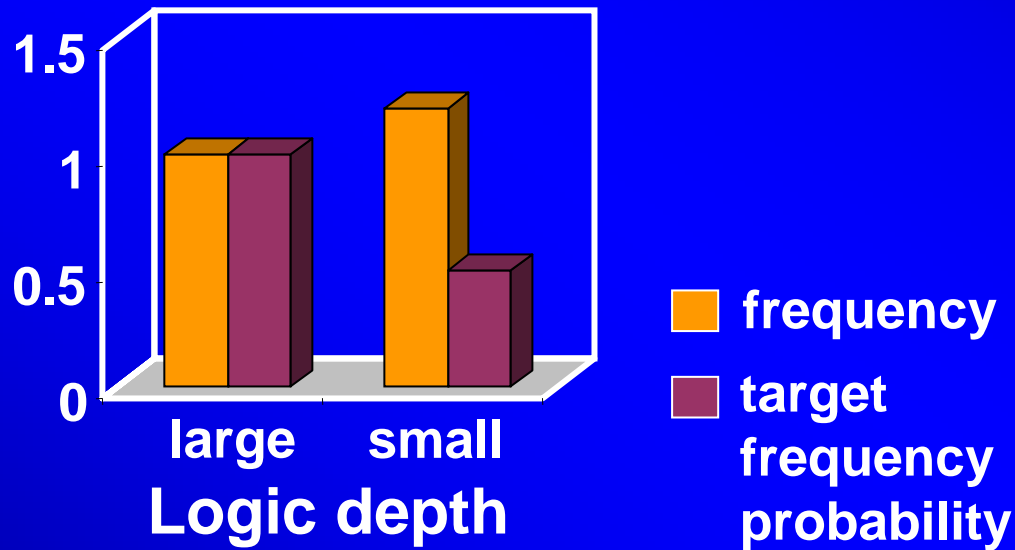
# Impact of Logic Depth



Logic depth: 16		
NMOS Ion	PMOS Ion	Delay
$\sigma/\mu$	$\sigma/\mu$	$\sigma/\mu$
5.6%	3.0%	4.2%



# $\mu$ Architecture Tradeoffs

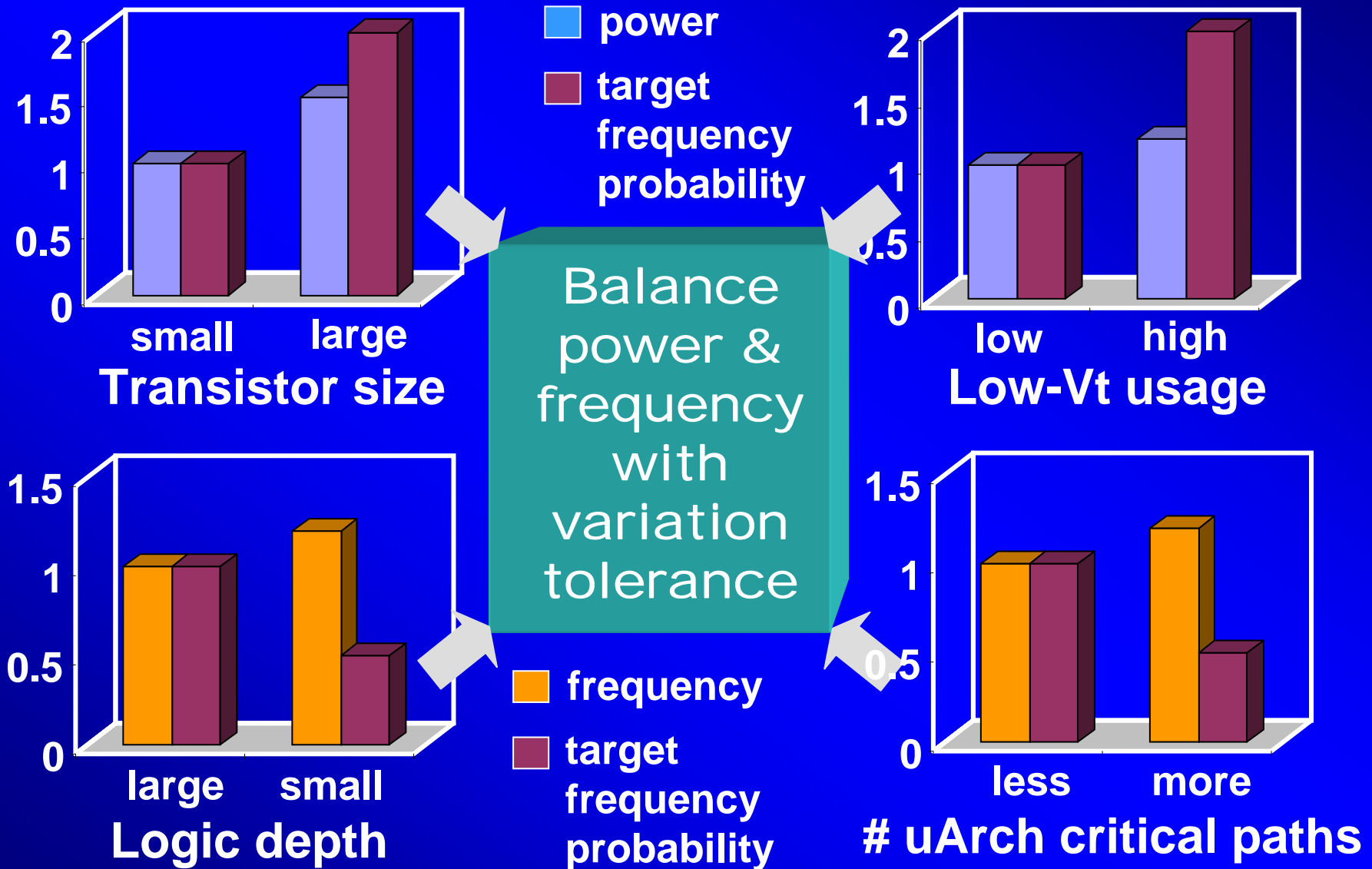


Higher target frequency with:

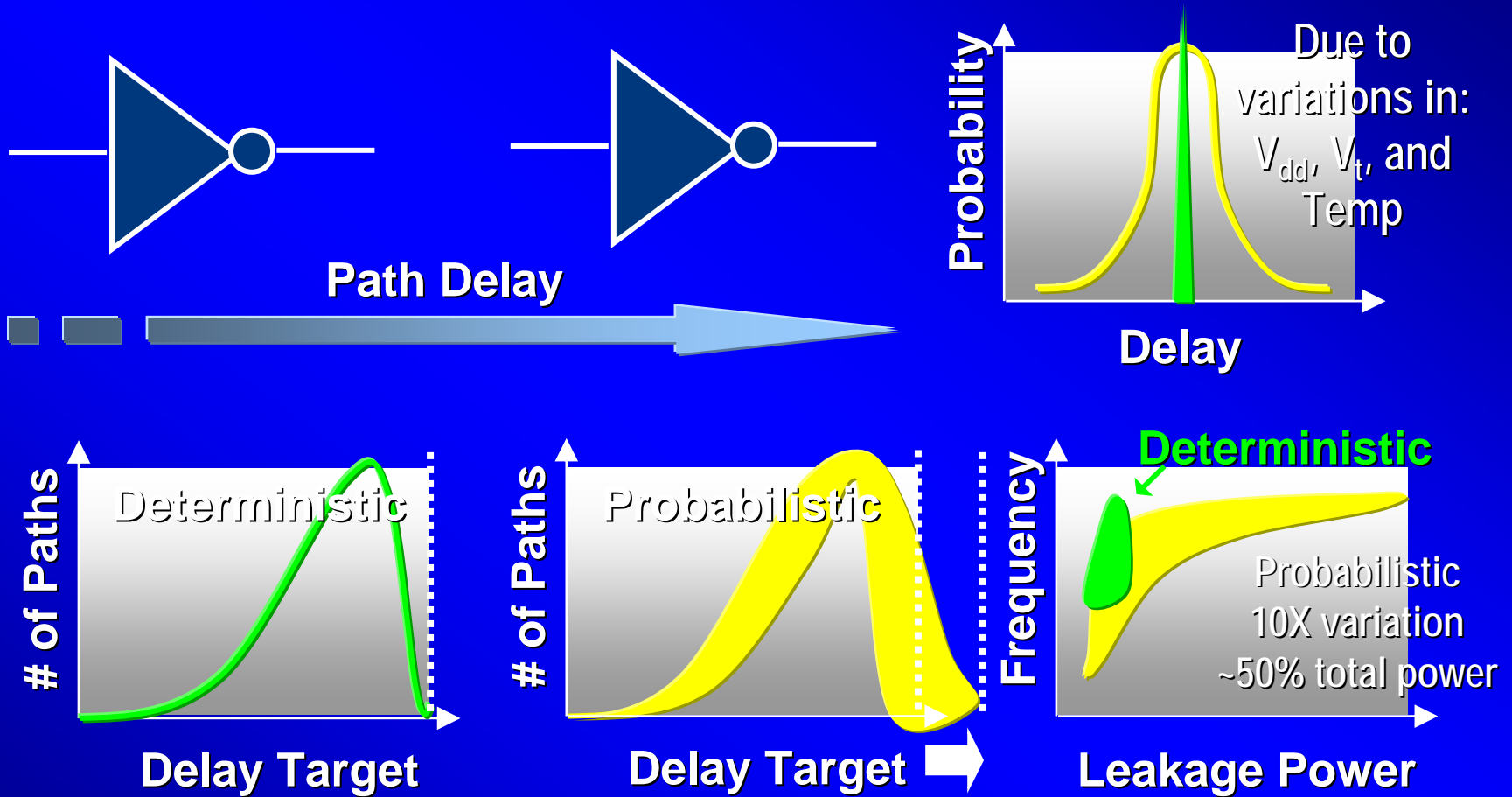
1. Shallow logic depth
2. Larger number of critical paths

But with lower probability

# Variation-tolerant Design



# Probabilistic Design



**Deterministic design techniques inadequate in the future**

# Shift in Design Paradigm

- **Multi-variable design optimization for:**
  - Yield and bin splits
  - Parameter variations
  - Active and leakage power
  - Performance

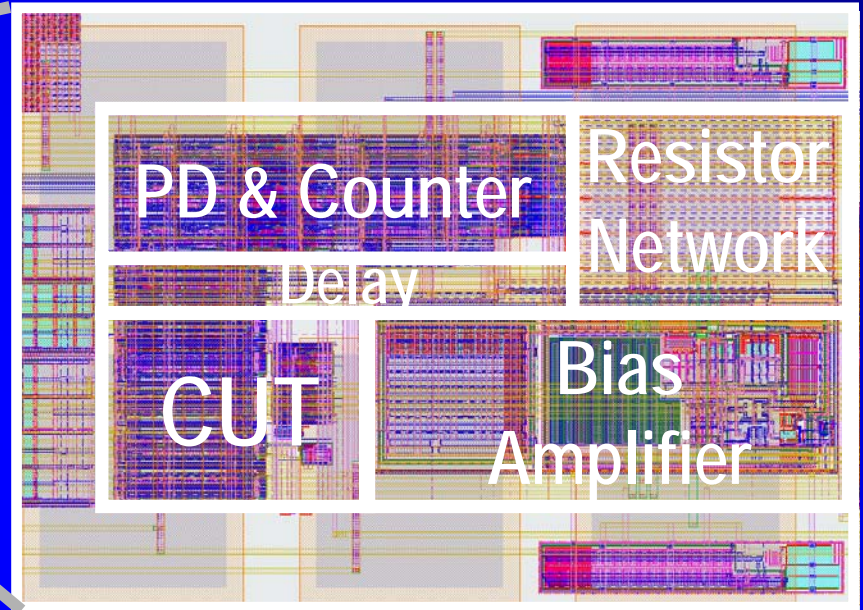
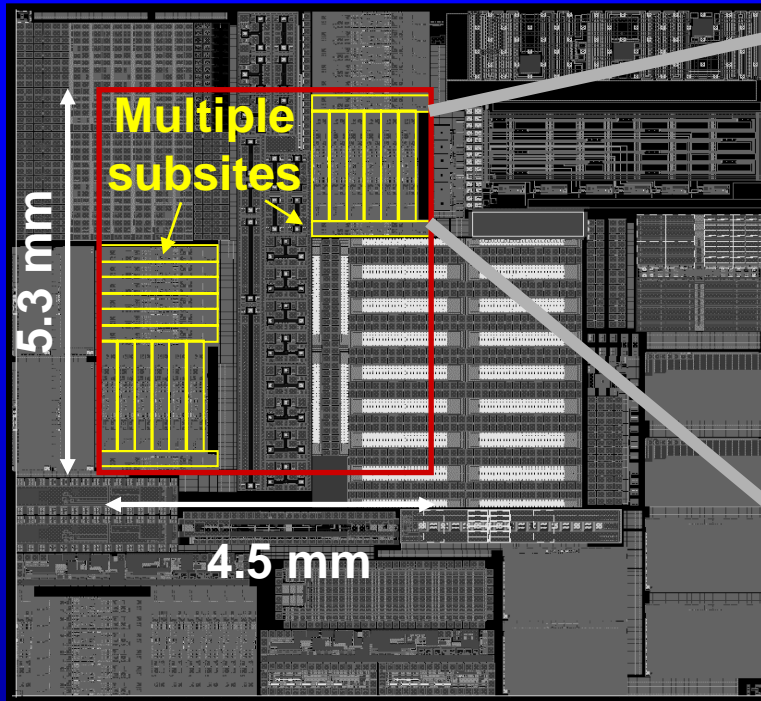
**Today:**

Local Optimization  
Single Variable

**Tomorrow:**

Global Optimization  
Multi-variate

# Adaptive Body Bias--Experiment



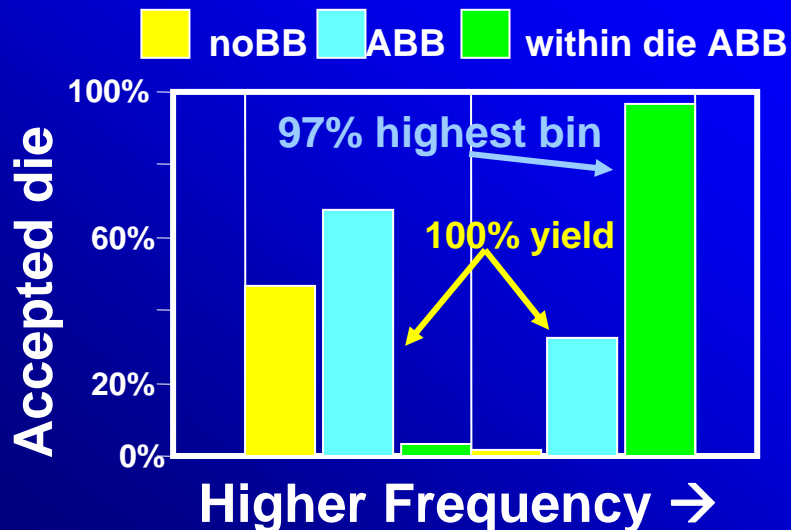
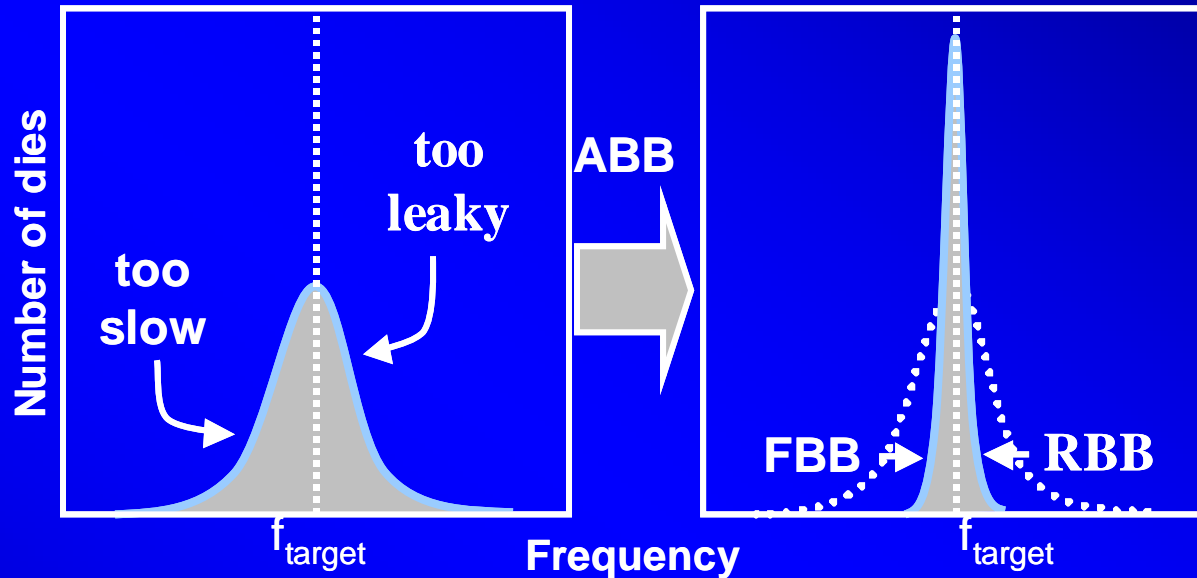
1.6 X 0.24 mm, 21 sites per die  
150nm CMOS

Technology	150nm CMOS
Number of subsites per die	21
Body bias range	0.5V FBB to 0.5V RBB
Bias resolution	32 mV

Die frequency:  $\text{Min}(F_1..F_{21})$

Die power:  $\text{Sum}(P_1..P_{21})$

# Adaptive Body Bias--Results

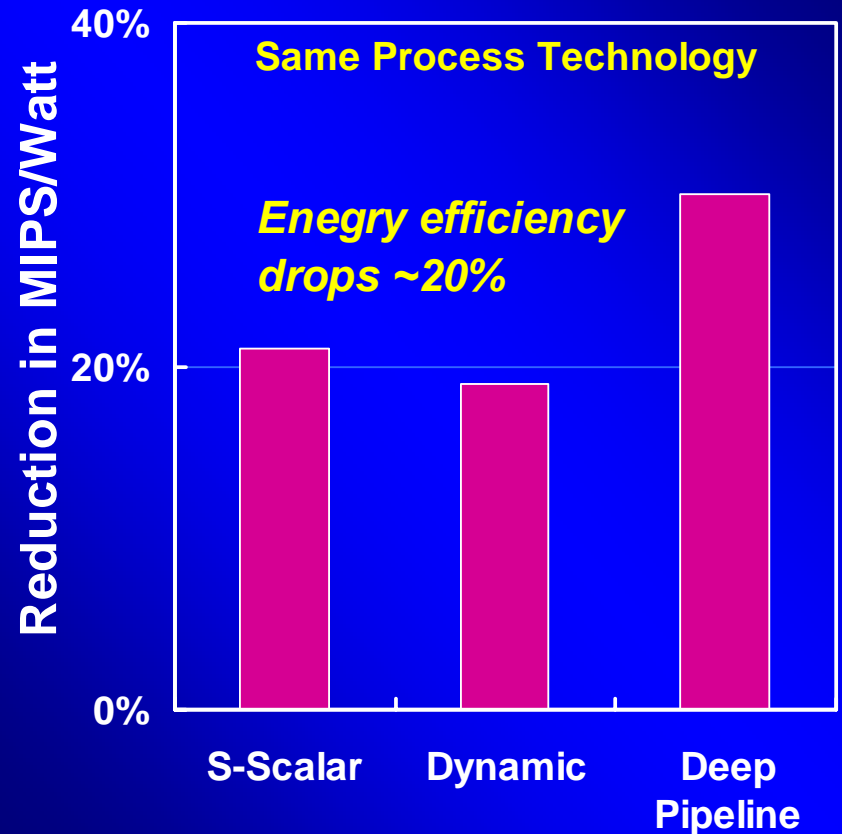
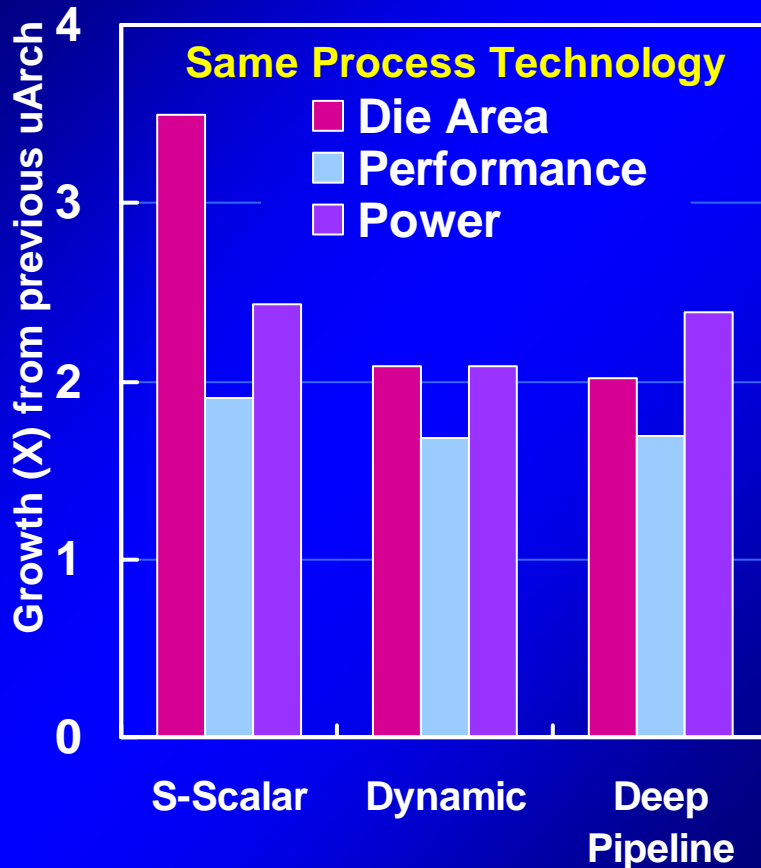


For given Freq and Power density

- 100% yield with ABB
- 97% highest freq bin with ABB for within die variability

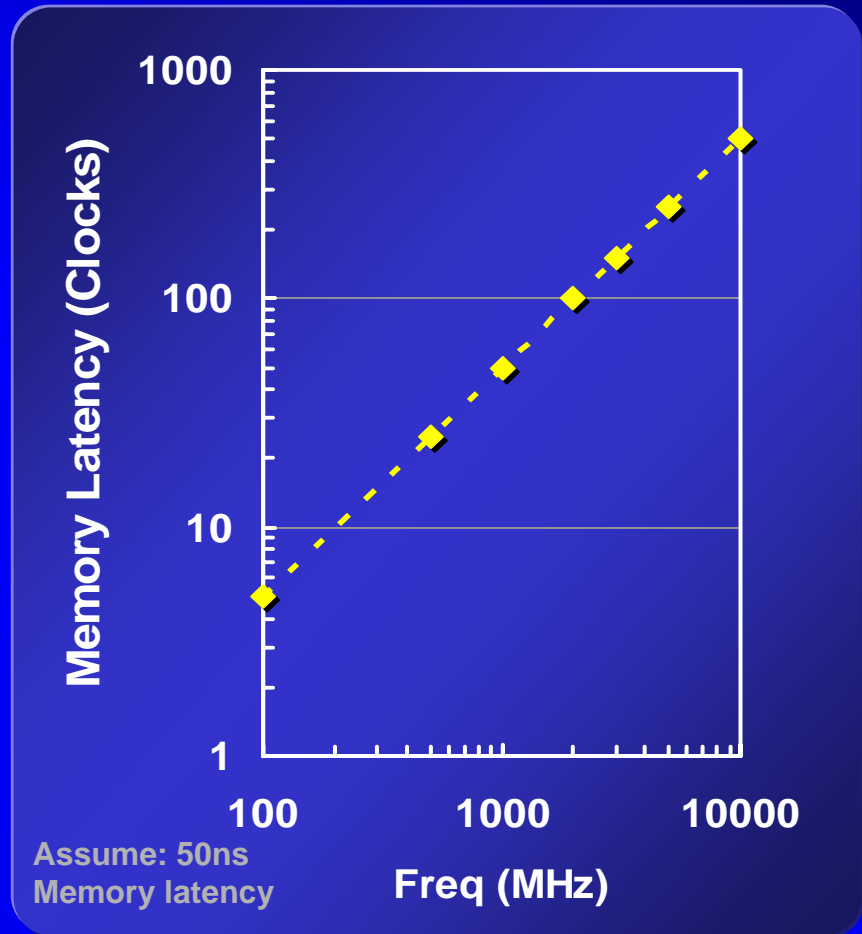
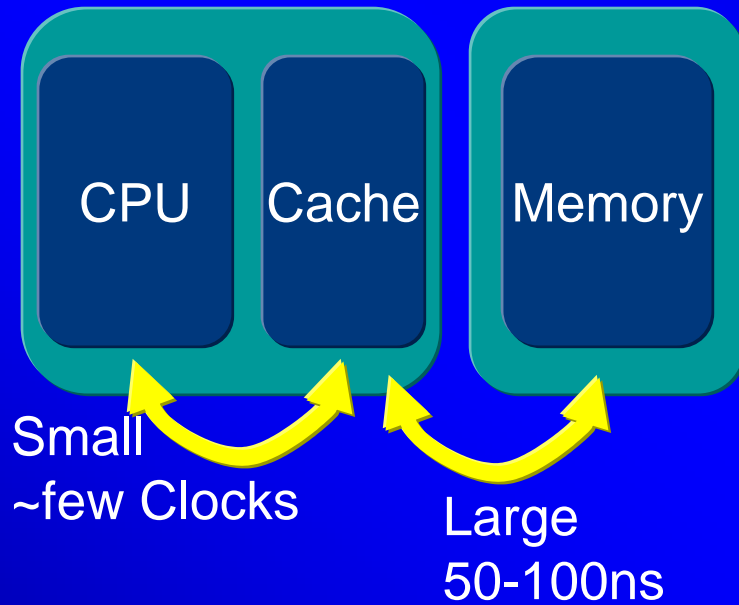


# Design & $\mu$ Arch Efficiency



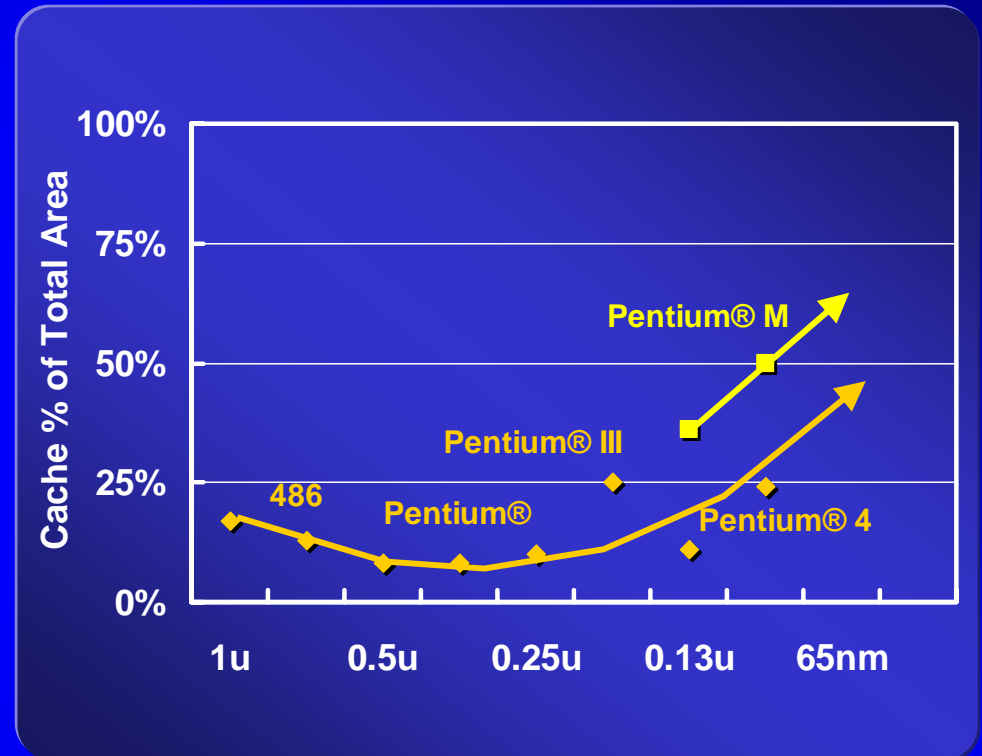
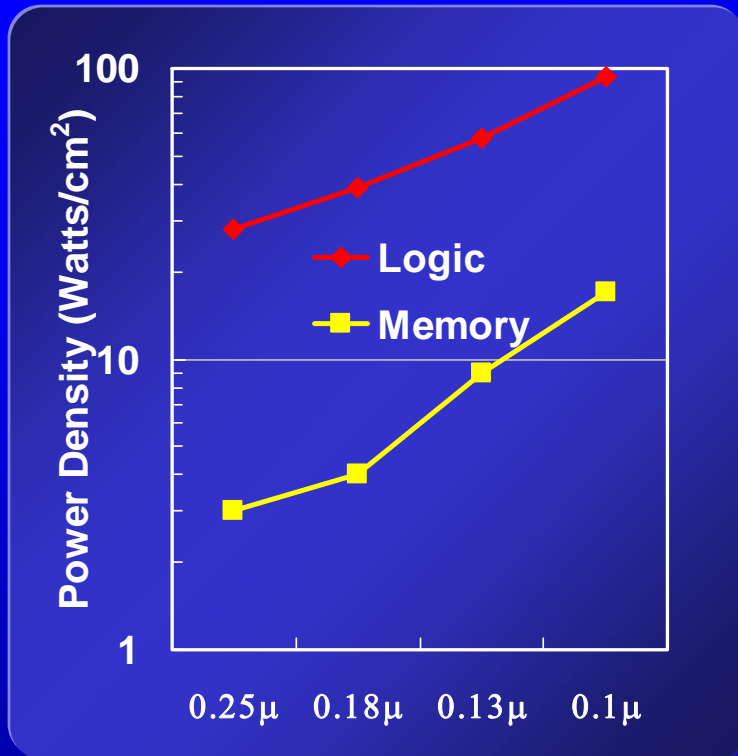
**Employ efficient design &  $\mu$ Architectures**

# Memory Latency



**Cache miss hurts performance  
Worse at higher frequency**

# Increase on-die Memory

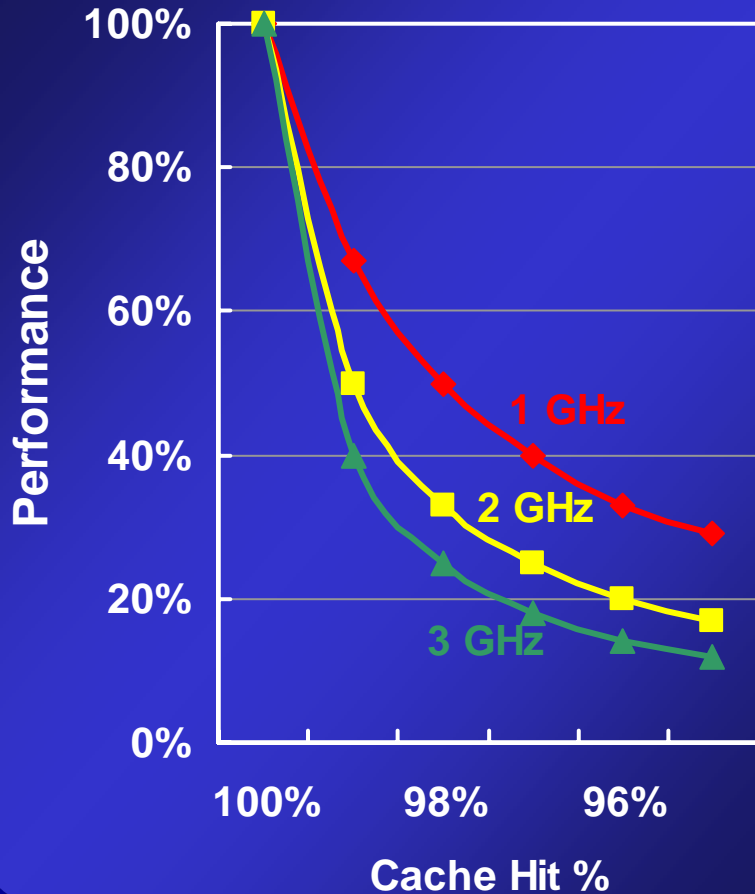


Large on die memory provides:

1. Increased Data Bandwidth & Reduced Latency
2. Hence, higher performance for much lower power

# Multi-threading

*Thermals & Power Delivery designed for full HW utilization*



## Single Thread

Full HW Utilization

ST

Wait for Mem

## Multi-Threading

MT1

Wait for Mem

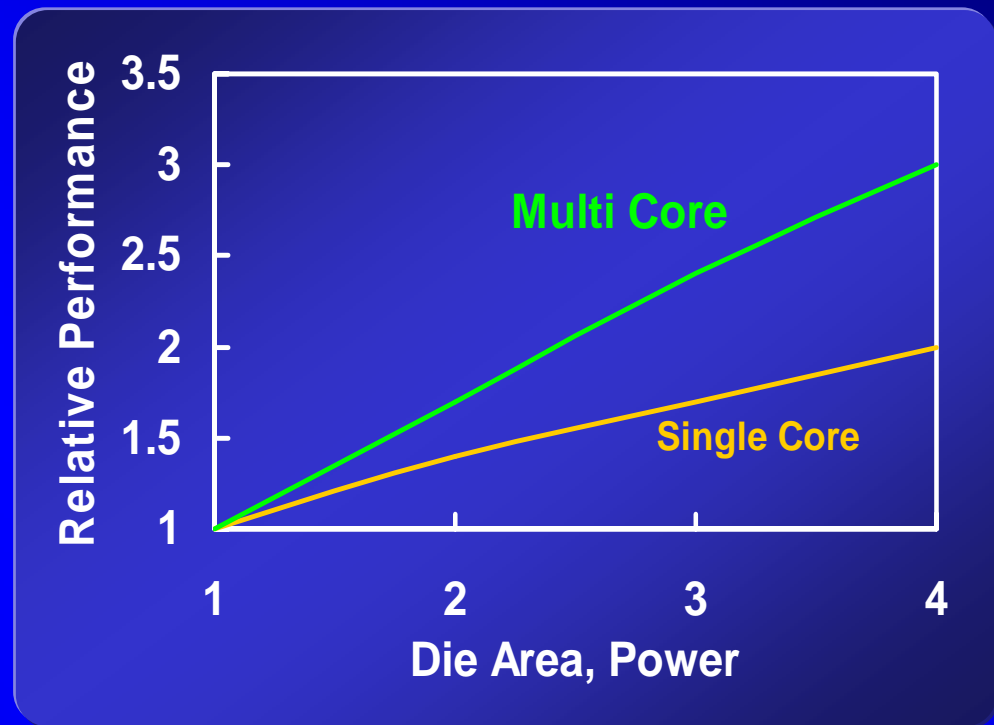
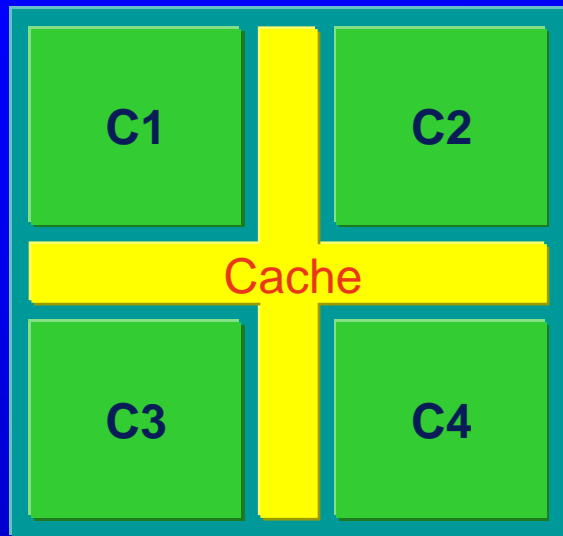
MT2

Wait

MT3

**Multi-threading improves performance without impacting thermals & power delivery**

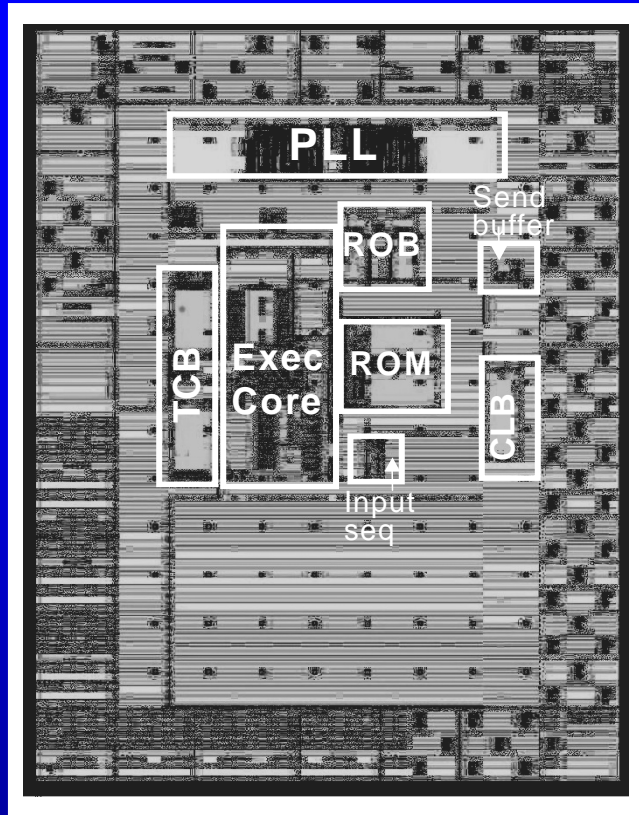
# Chip Multi-Processing



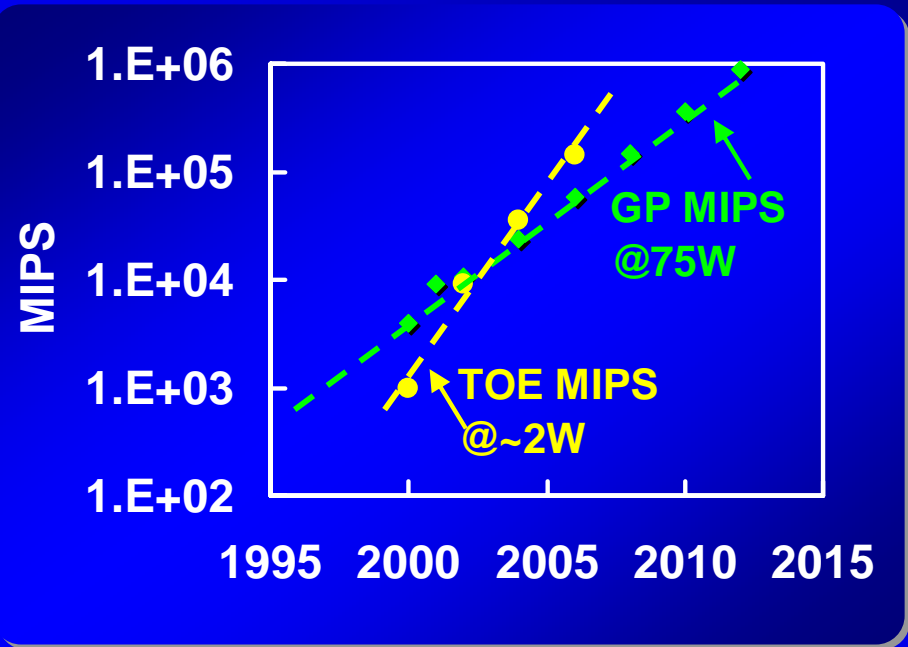
- Multi-core, each core Multi-threaded
- Shared cache and front side bus
- Each core has different Vdd & Freq
- Core hopping to spread hot spots
- Lower junction temperature

# Special Purpose Hardware

## TCP Offload Engine



2.23 mm X 3.54 mm, 260K transistors



Opportunities:  
Network processing engines  
MPEG Encode/Decode engines  
Speech engines

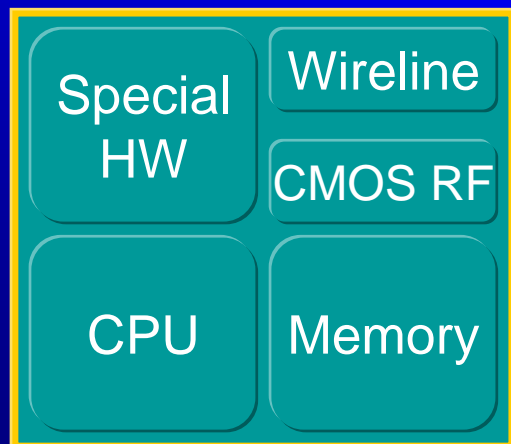
**Special purpose HW—Best Mips/Watt**

# Valued Performance: SOC (System on a Chip)

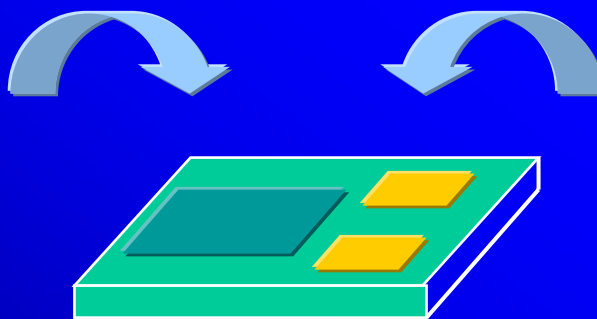
- Special-purpose hardware → more MIPS/mm<sup>2</sup>
- SIMD integer and FP instructions in several ISAs

	Die Area	Power	Performance
General Purpose	2X	2X	~1.4X
Multimedia Kernels	<10%	<10%	1.5 - 4X

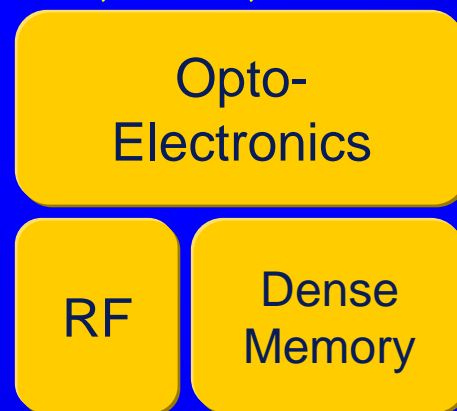
Si Monolithic



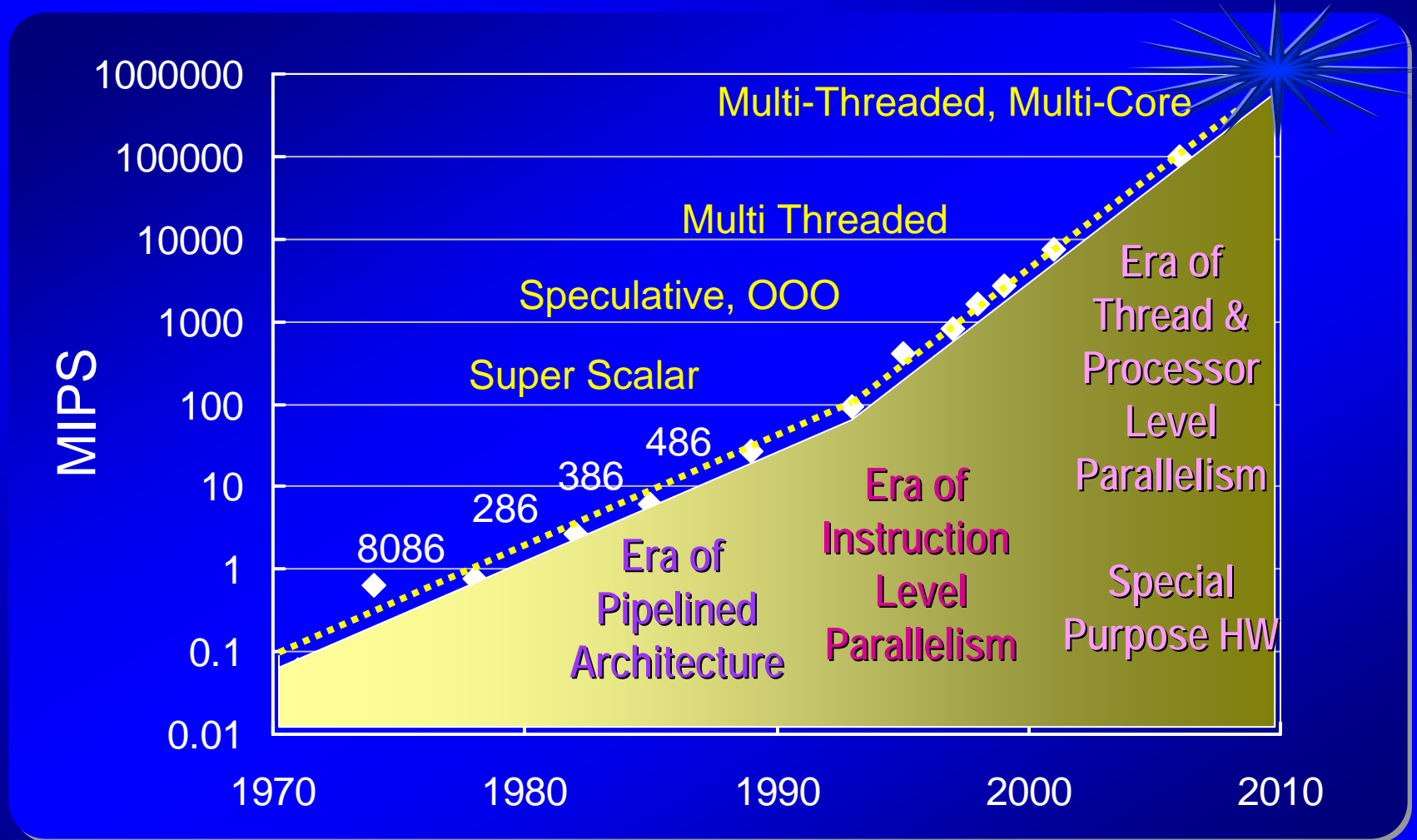
Polyolithic



Heterogeneous Si, SiGe, GaAs



# The Exponential Reward



**Multi-everywhere: MT, CMP**



# Summary—Delaying *Forever*

- **Gigascale transistor integration capacity will be available—Power and Energy are the barriers**
- **Variations will be even more prominent—shift from Deterministic to Probabilistic design**
- **Improve design efficiency**
- **Multi—everywhere, & SOC  $\Rightarrow$ valued performance**
- **Exploit integration capacity to deliver performance in power/cost envelope**