

Gigascale Integration— Design Challenges & Opportunities

Shekhar Borkar

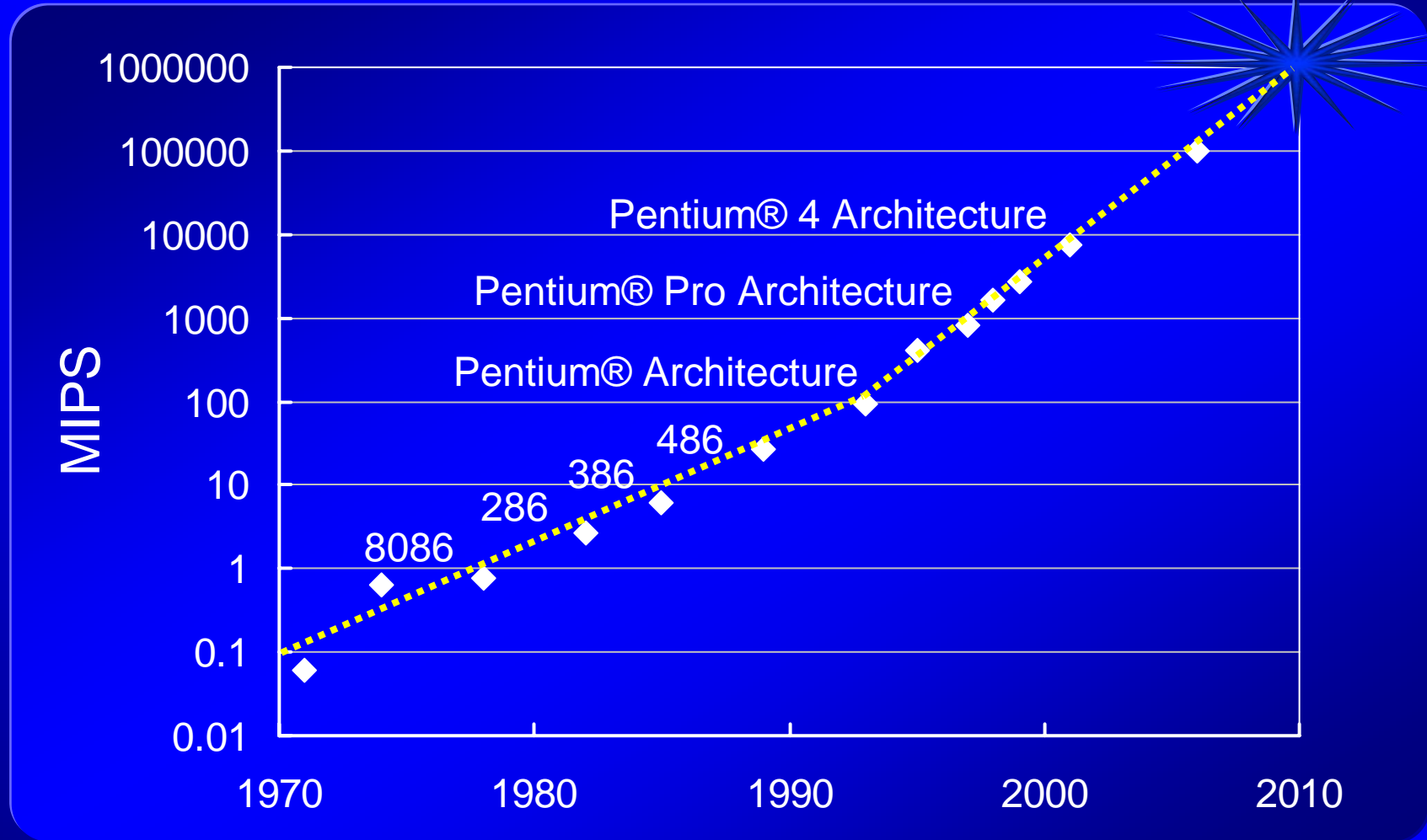
Circuit Research, Intel Labs

October 24, 2004

Outline

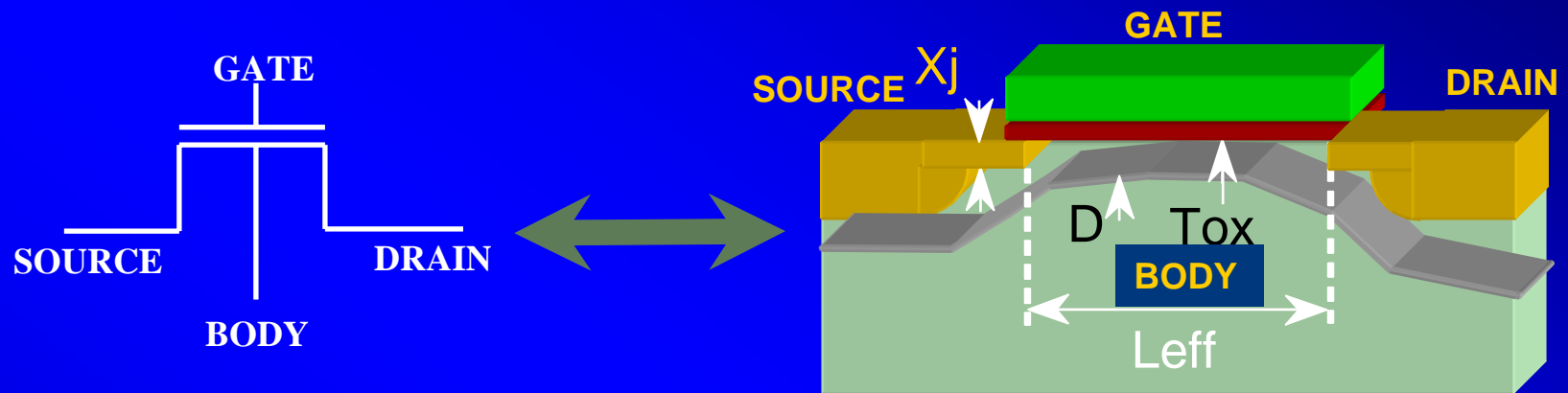
- **CMOS technology challenges**
- **Technology, circuit and μ Architecture solutions**
- **Integration opportunities**
- **Summary**

Goal: 1TIPS by 2010



How do you get there?

CMOS Technology Scaling



**Dimensions scale
down by 30%**

**Oxide thickness
scales down**

Vdd & Vt scaling

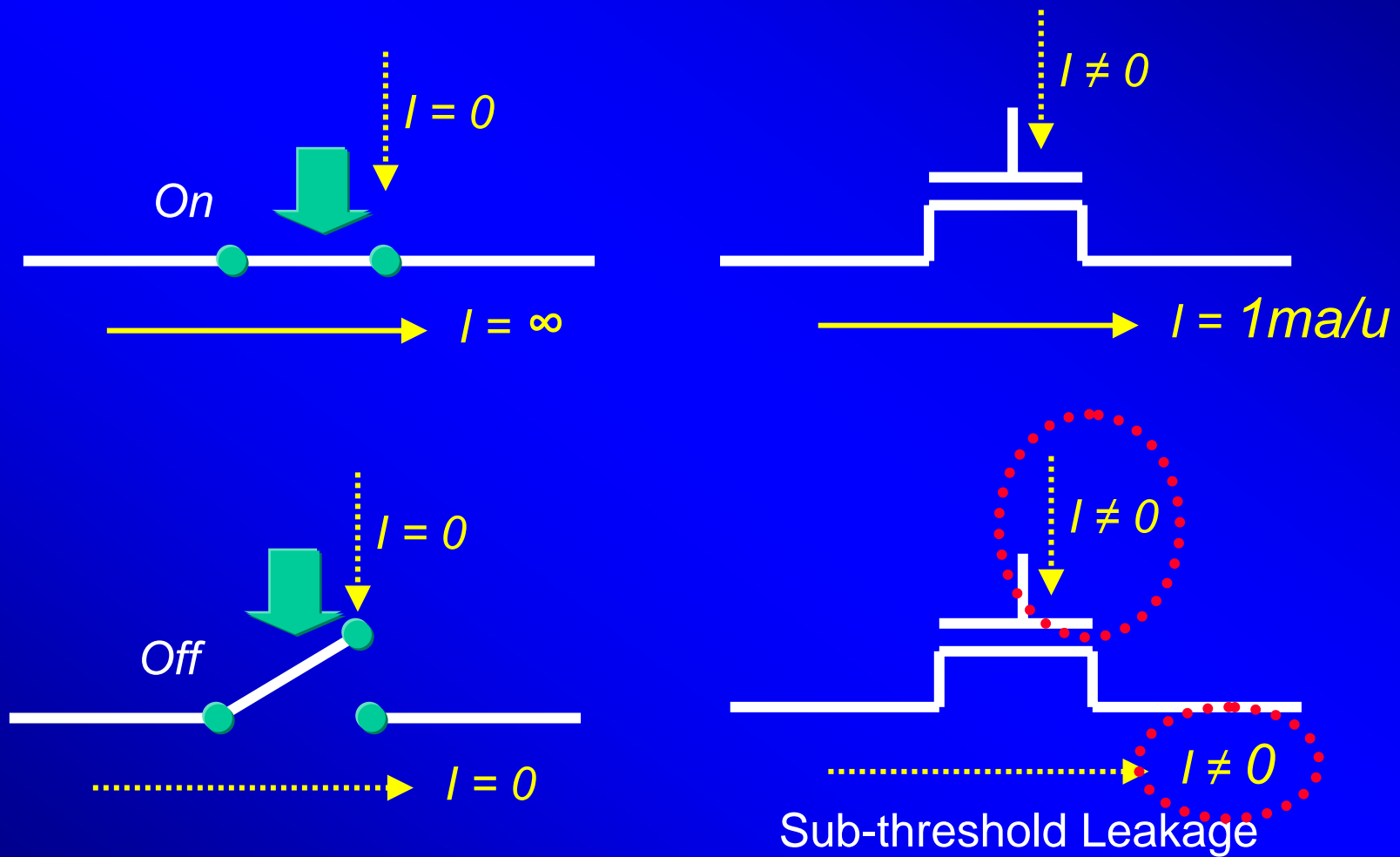
**Doubles transistor
density**

**Faster transistor,
higher performance**

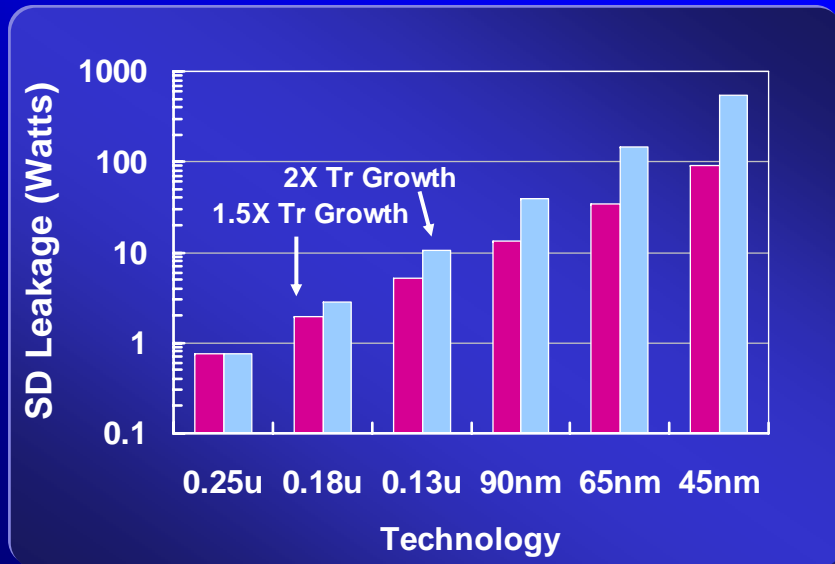
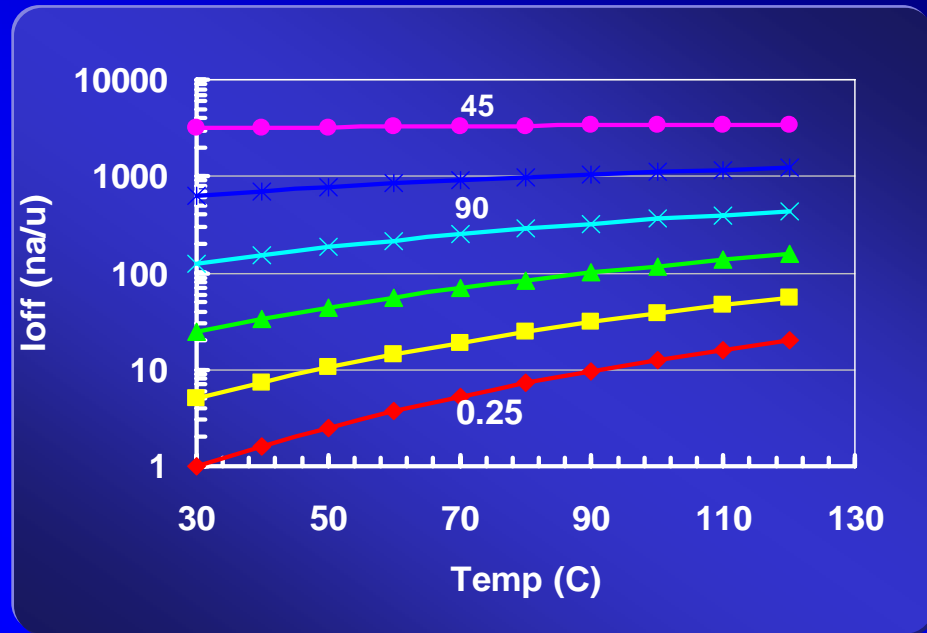
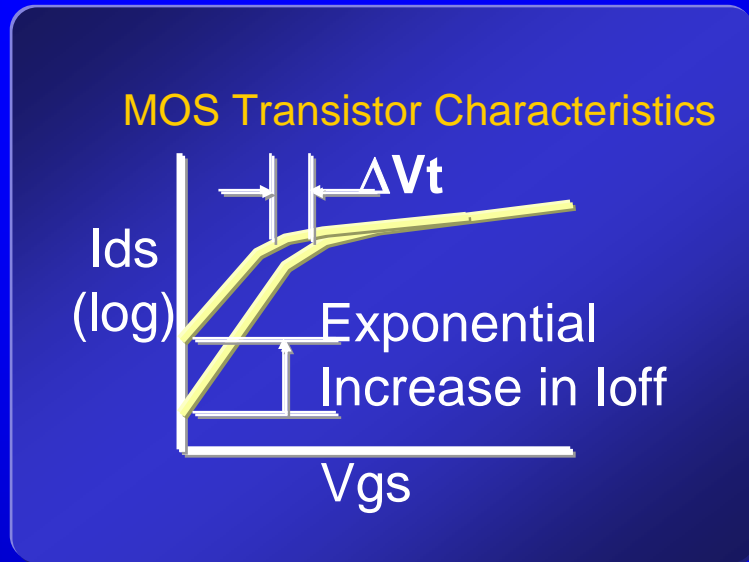
Lower active power

Technology has scaled well, will it in the future?

Is Transistor a Good Switch?

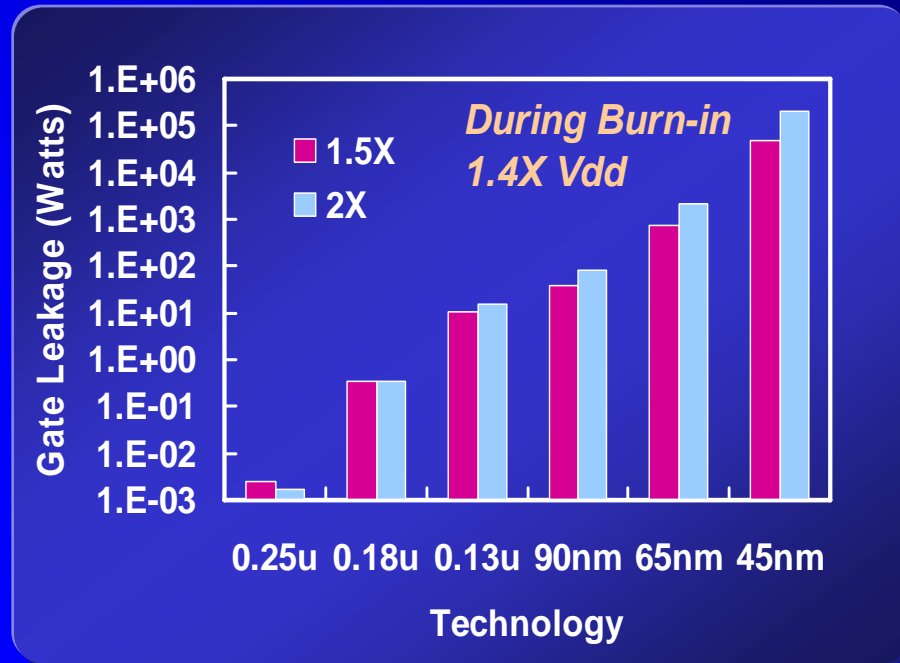
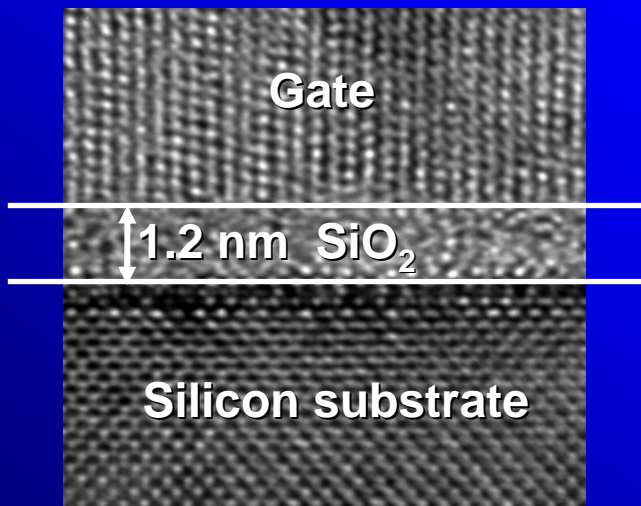
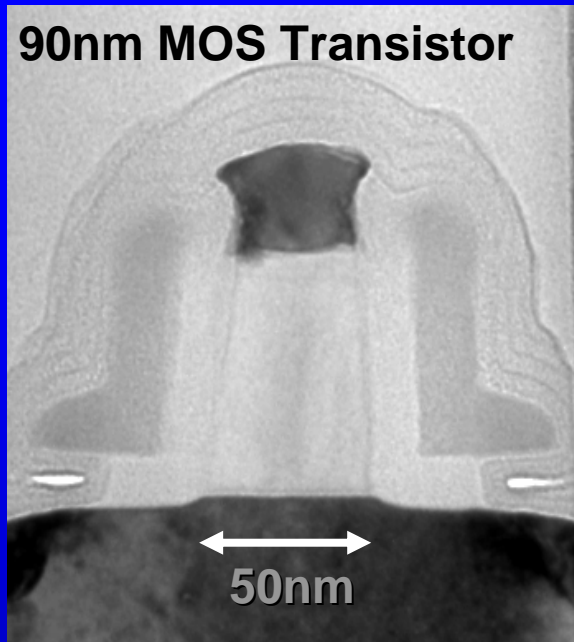


Sub-threshold Leakage



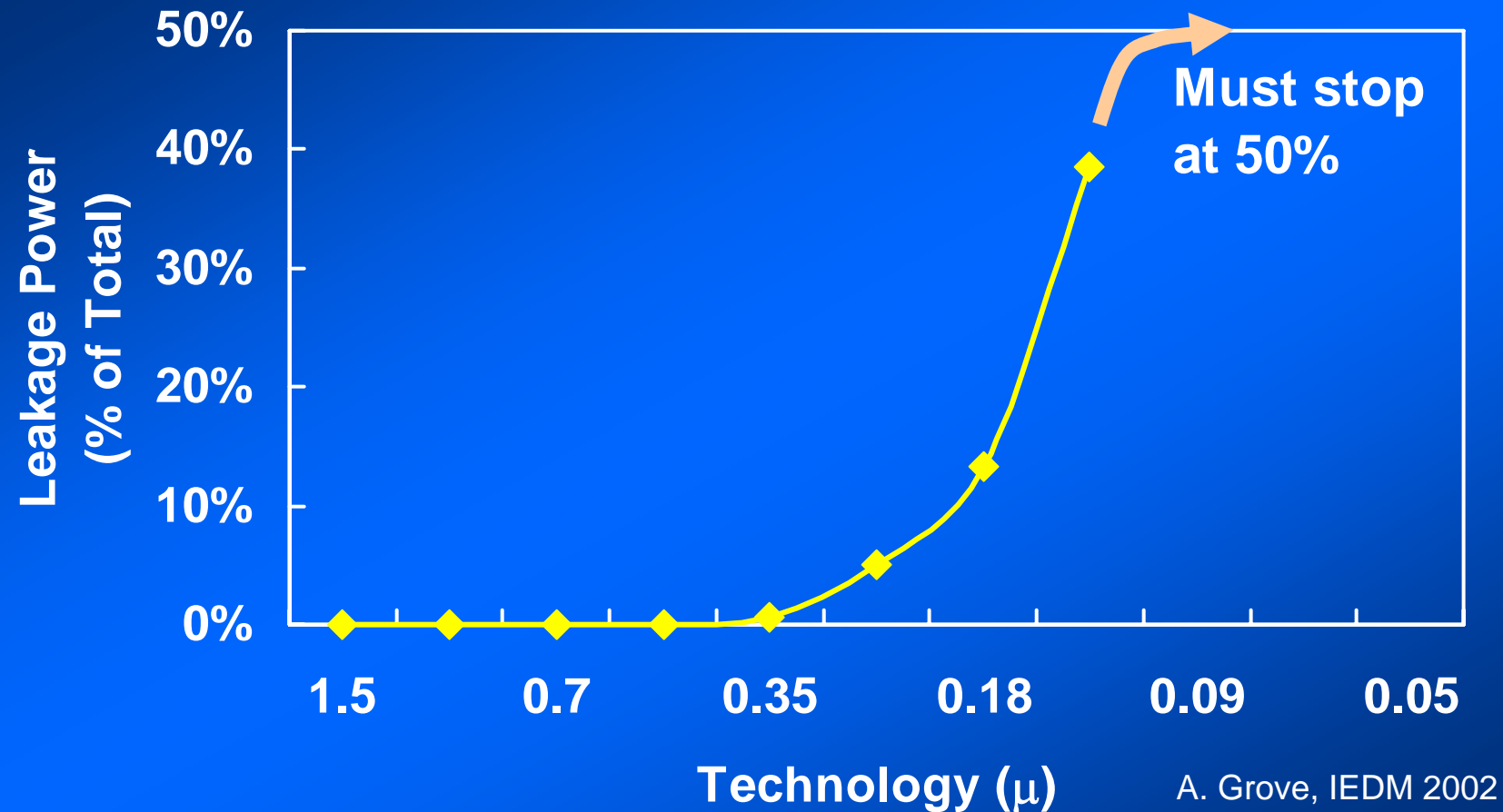
Transistors will not be switches, but dimmers...
 ...Dennis Buss, TI

Gate Oxide is Near Limit



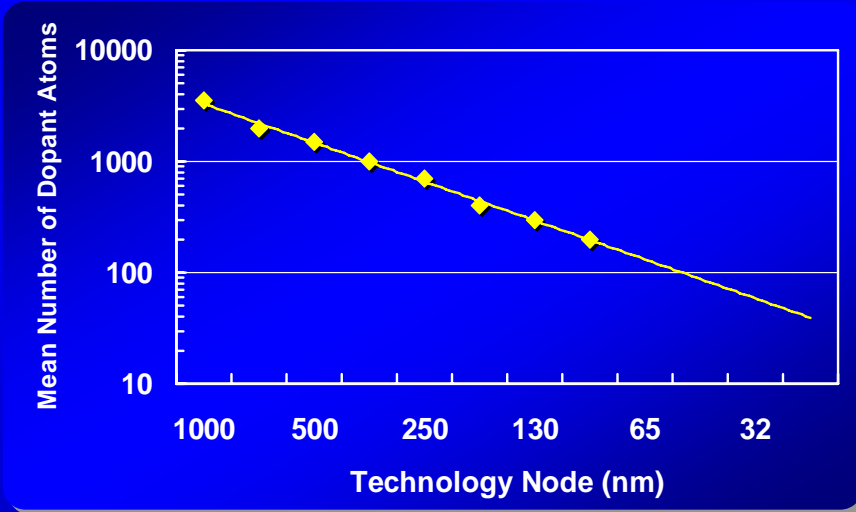
**High K & Metal Gate
crucial for the industry**

Leakage Power

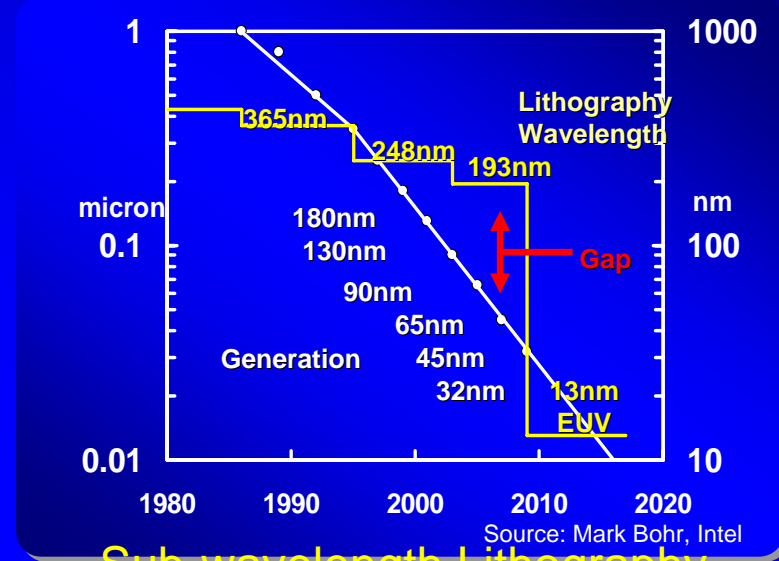


Leakage power limits Vt scaling

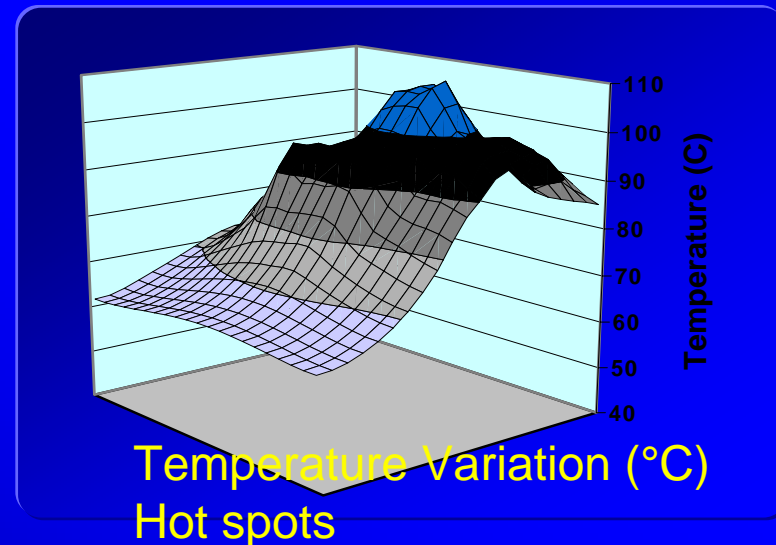
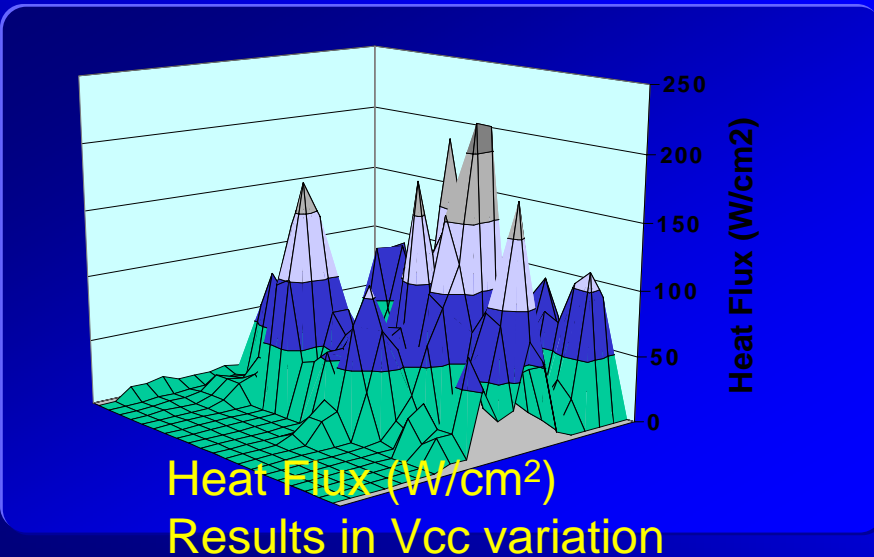
Sources of Variations



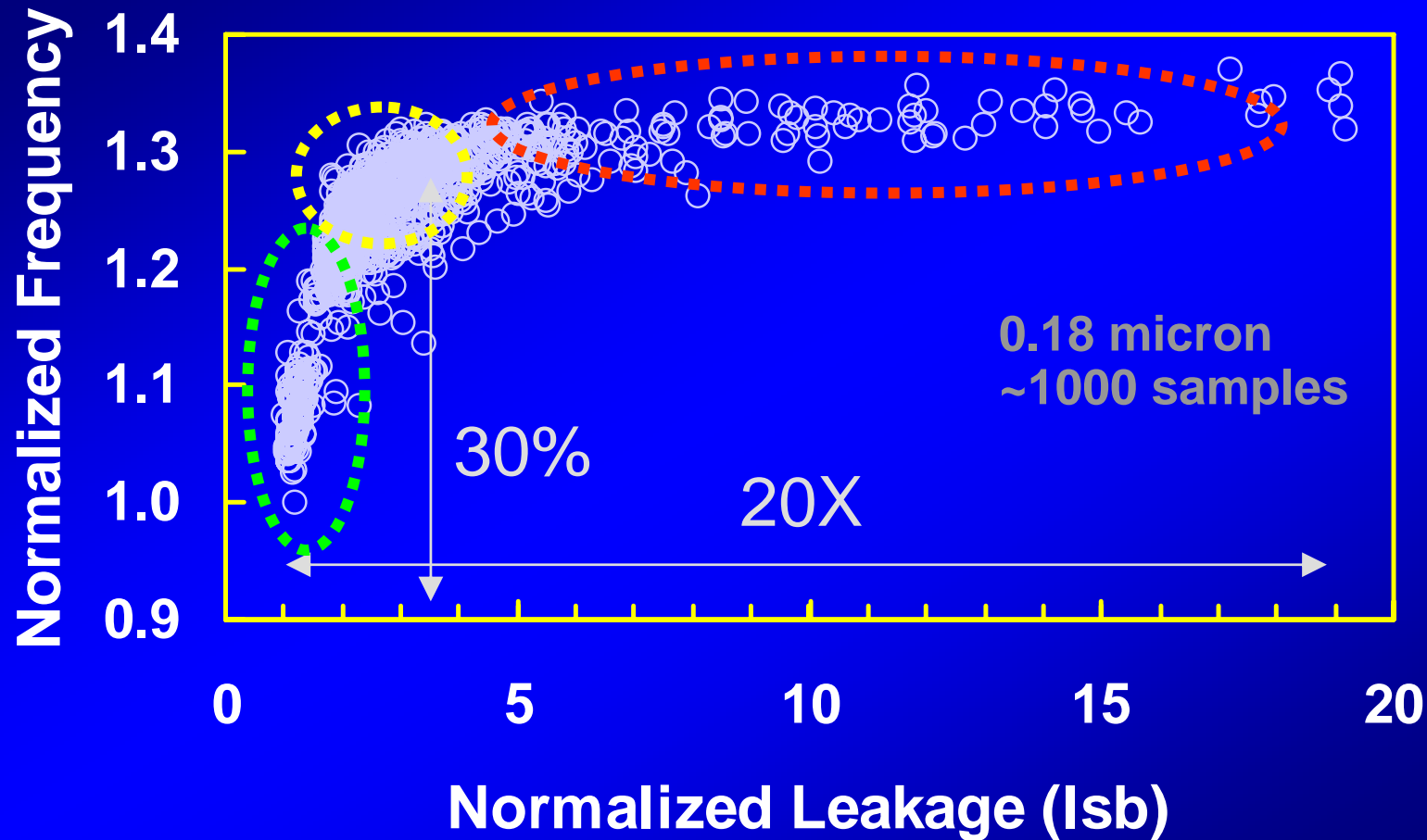
Random Dopant Fluctuations



Sub-wavelength Lithography



Impact of Variations

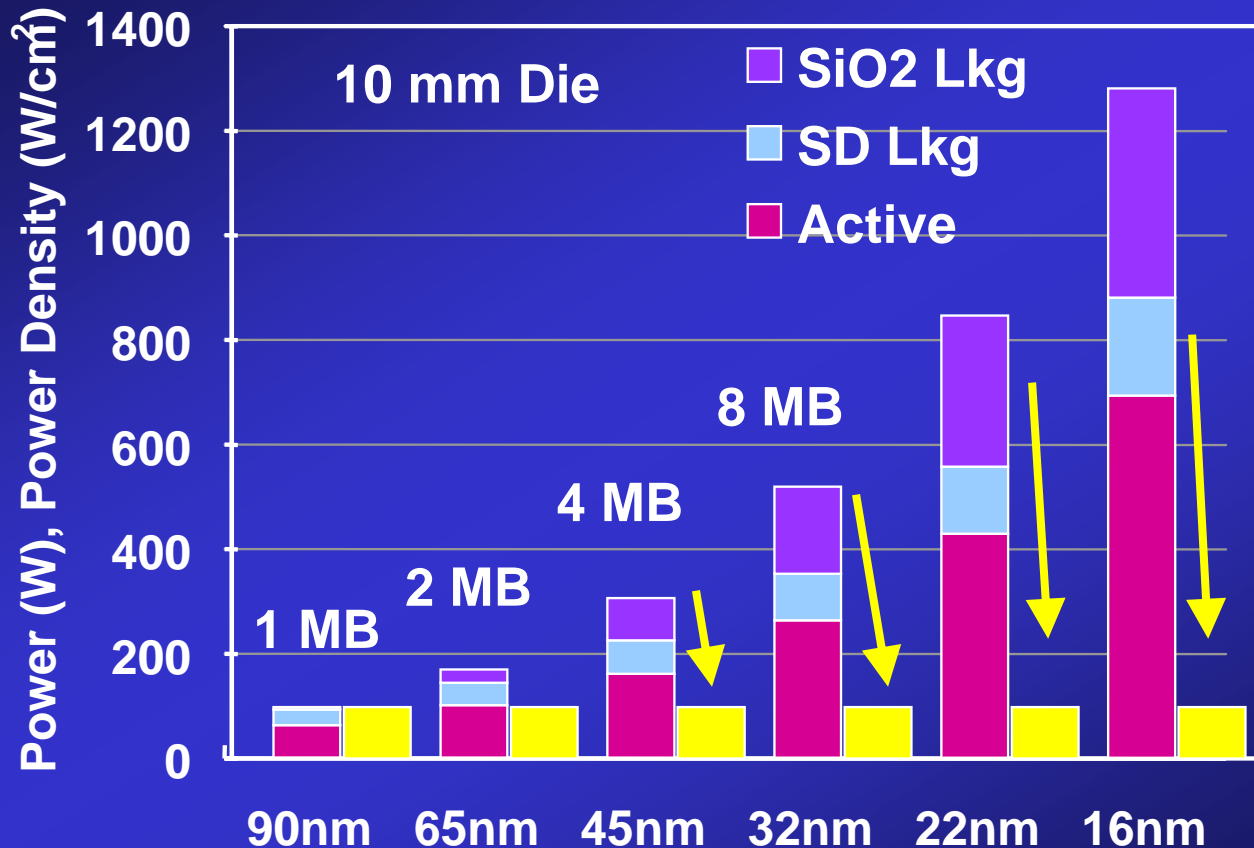


Low Freq
Low I_{sb}

High Freq
Medium I_{sb}

High Freq
High I_{sb}

The Power Envelope...



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

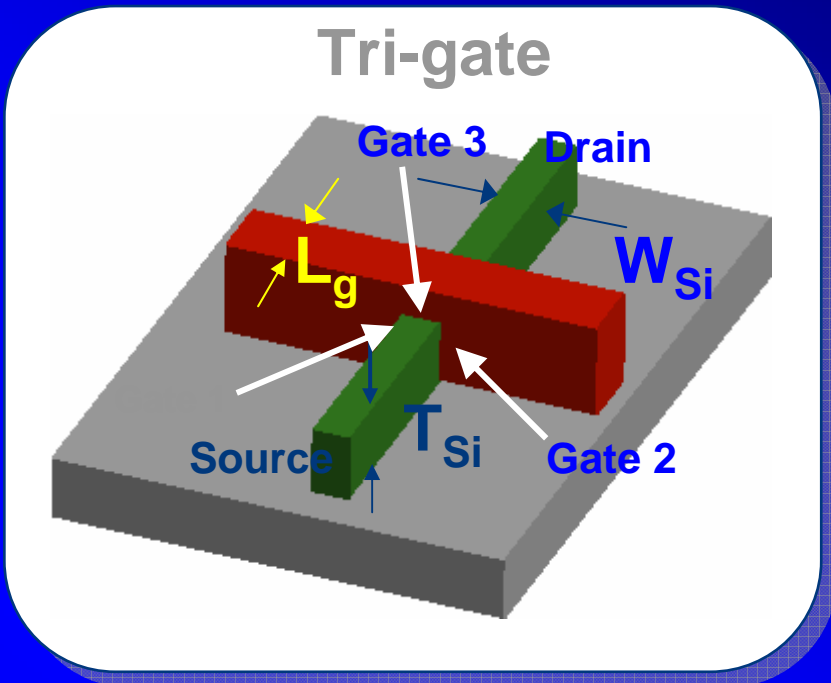
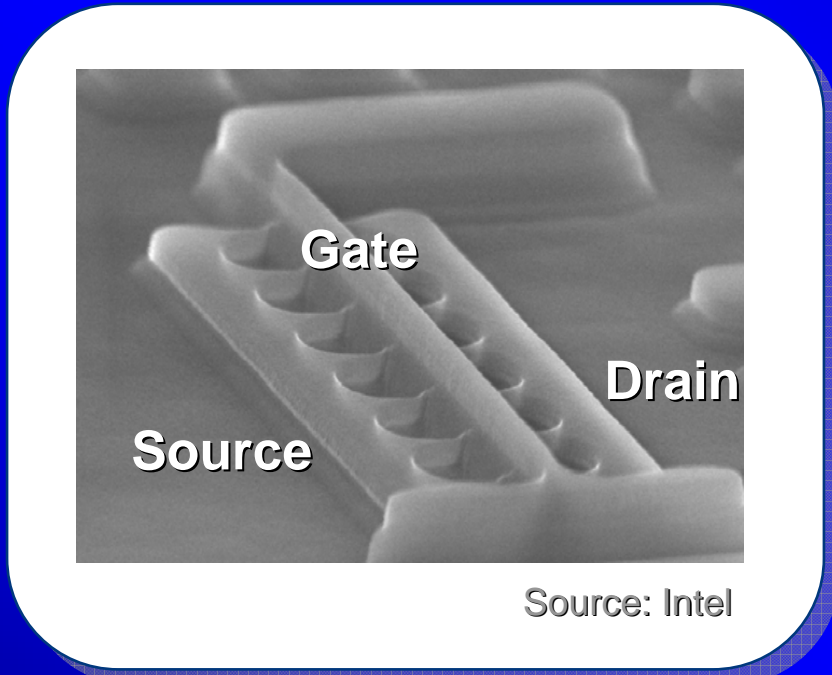
The Gigascale Dilemma

- Huge transistor integration capacity
- But unusable due to power
- Logic T growth will have to slow down
- Transistor performance will be limited

Solutions

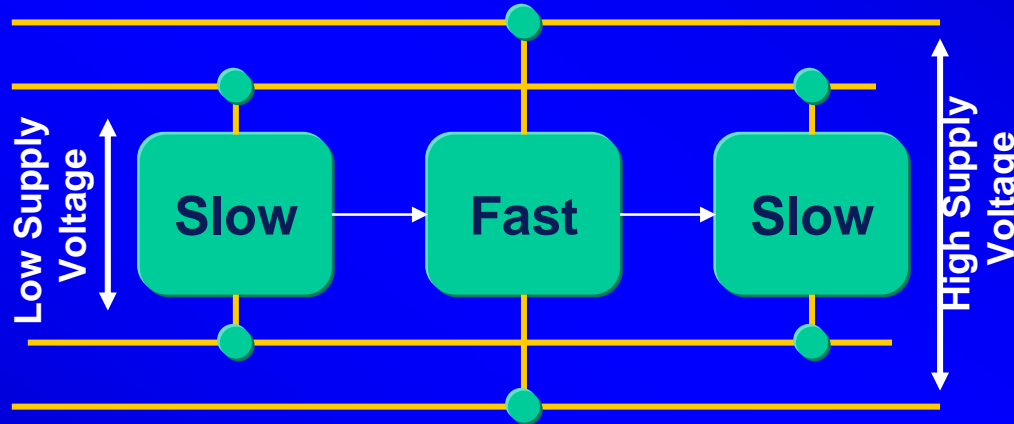
- Low power design techniques
- Improve design efficiency—Multi everywhere
- Valued performance by even higher integration (of potentially slower transistors)

New Transistors: Tri-Gate...



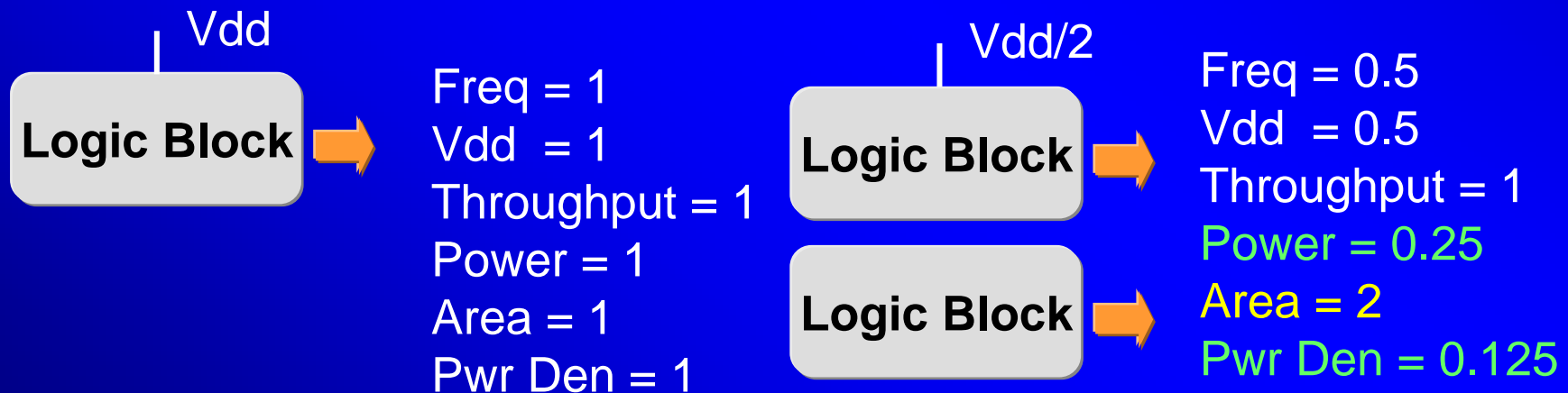
Improved short-channel effects
Higher ON current for lower SD Leakage
Manufacturing control: research underway

Active Power Reduction



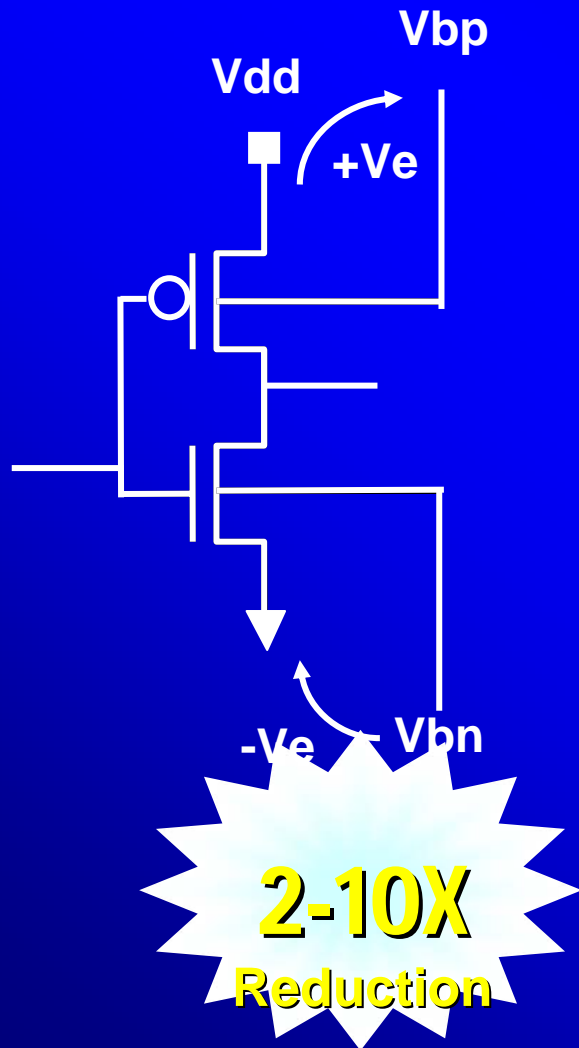
Multiple Supply Voltages

Replicated Designs

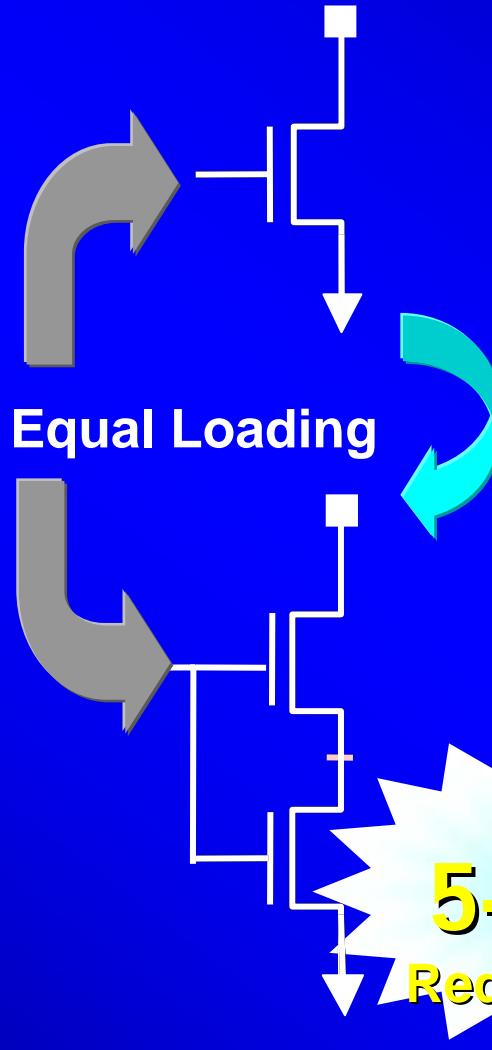


Leakage Control

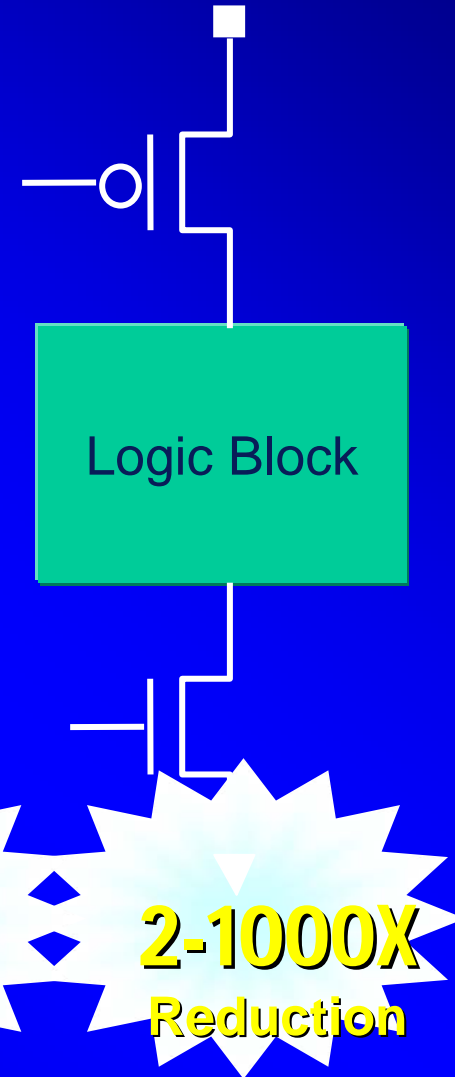
Body Bias



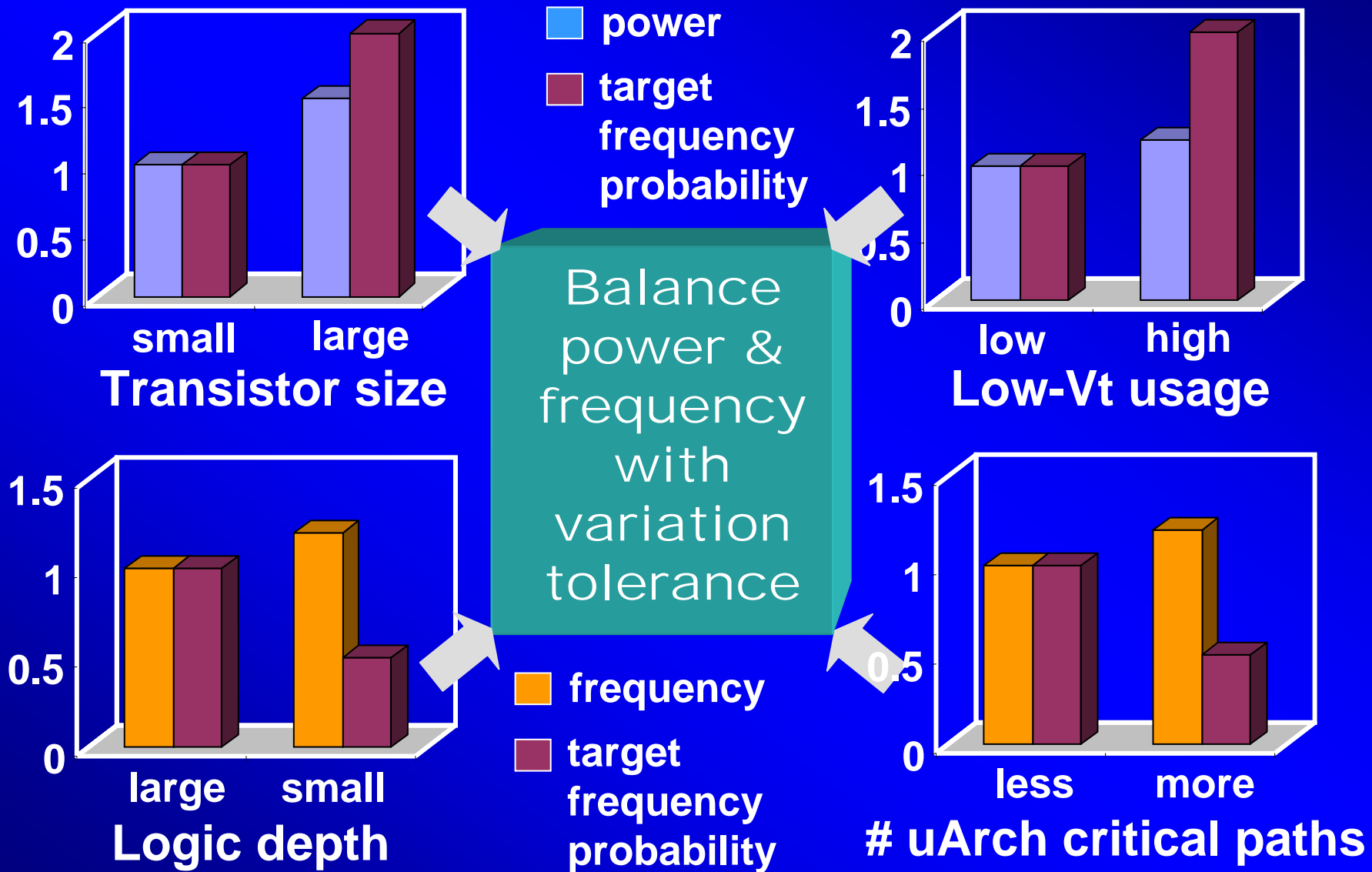
Stack Effect



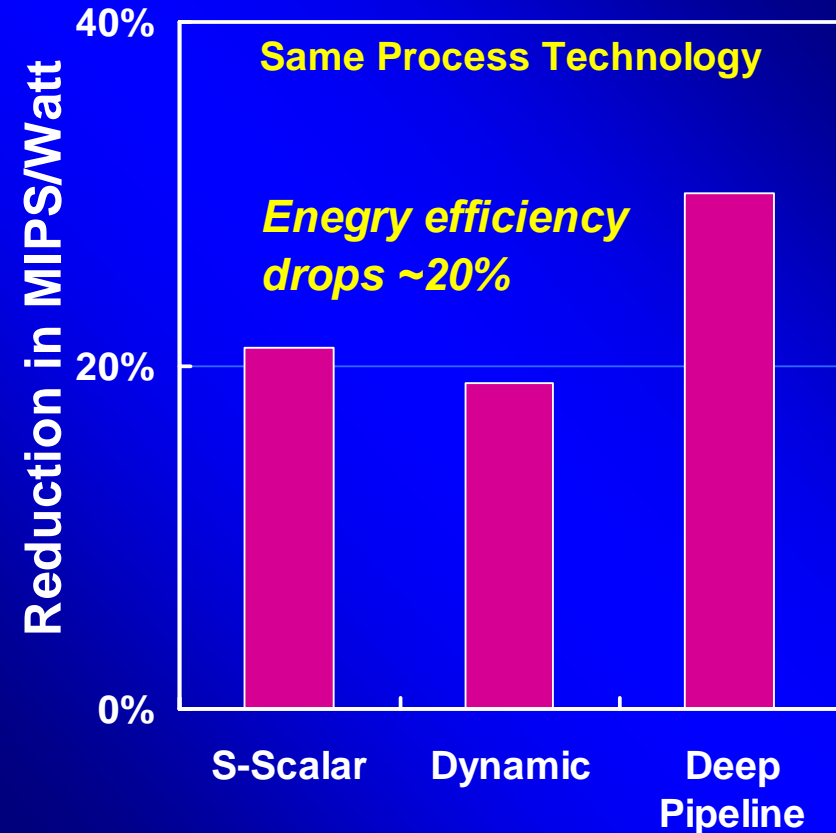
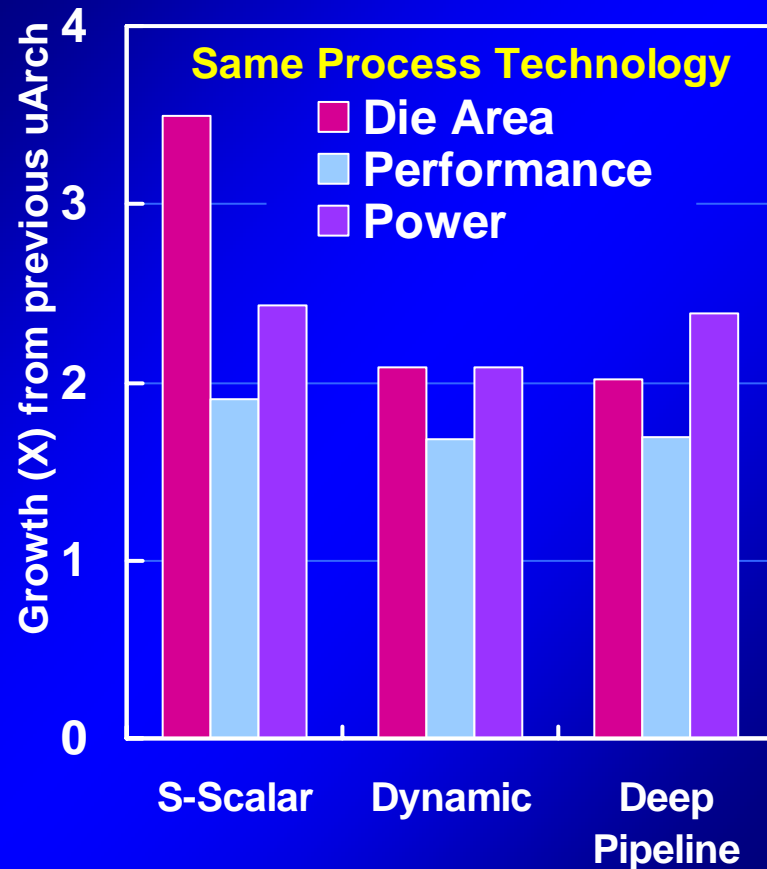
Sleep Transistor



Variation-tolerant Design

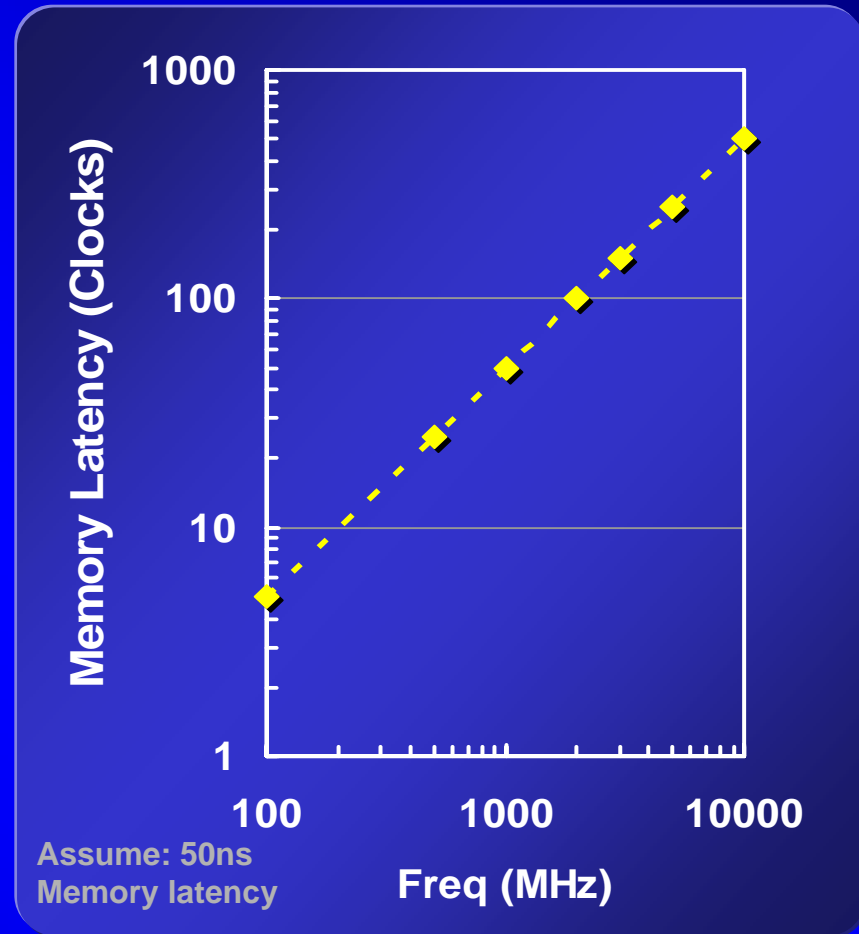
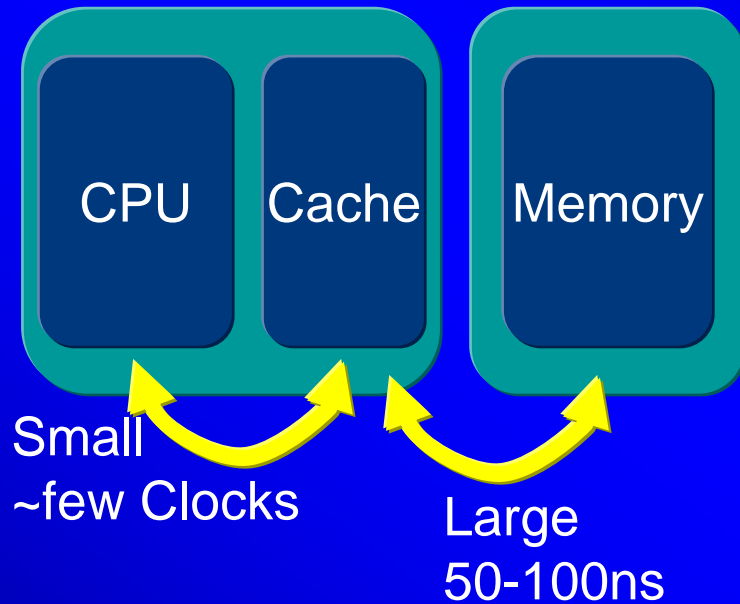


Design & μ Arch Efficiency



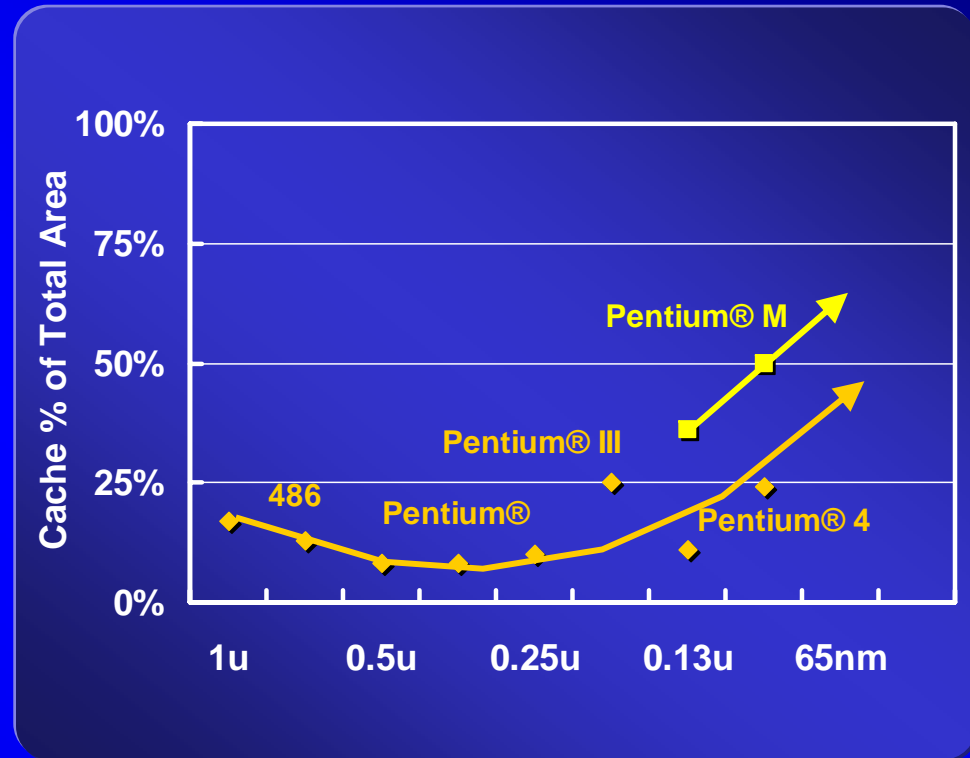
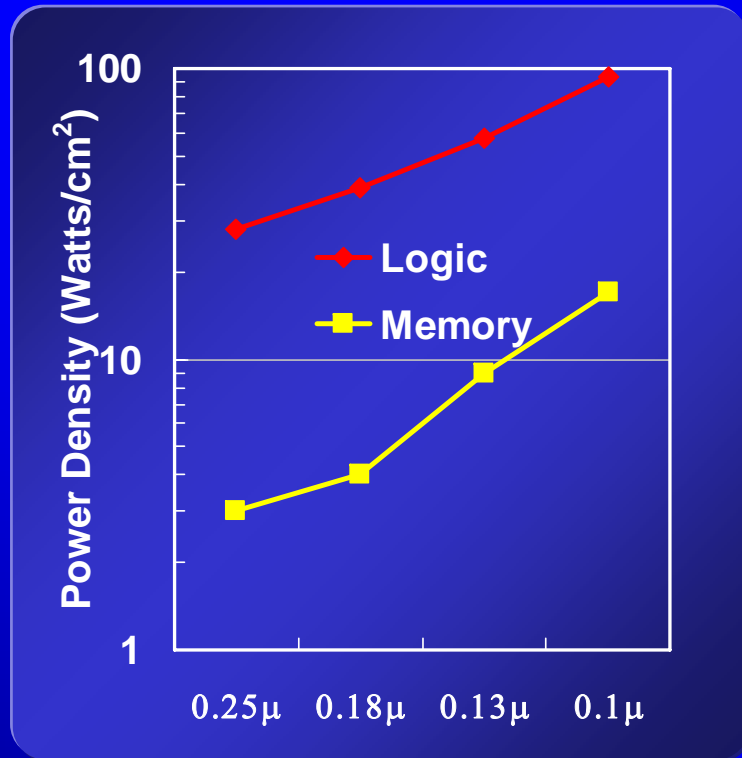
Employ efficient design & μ Architectures

Memory Latency



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

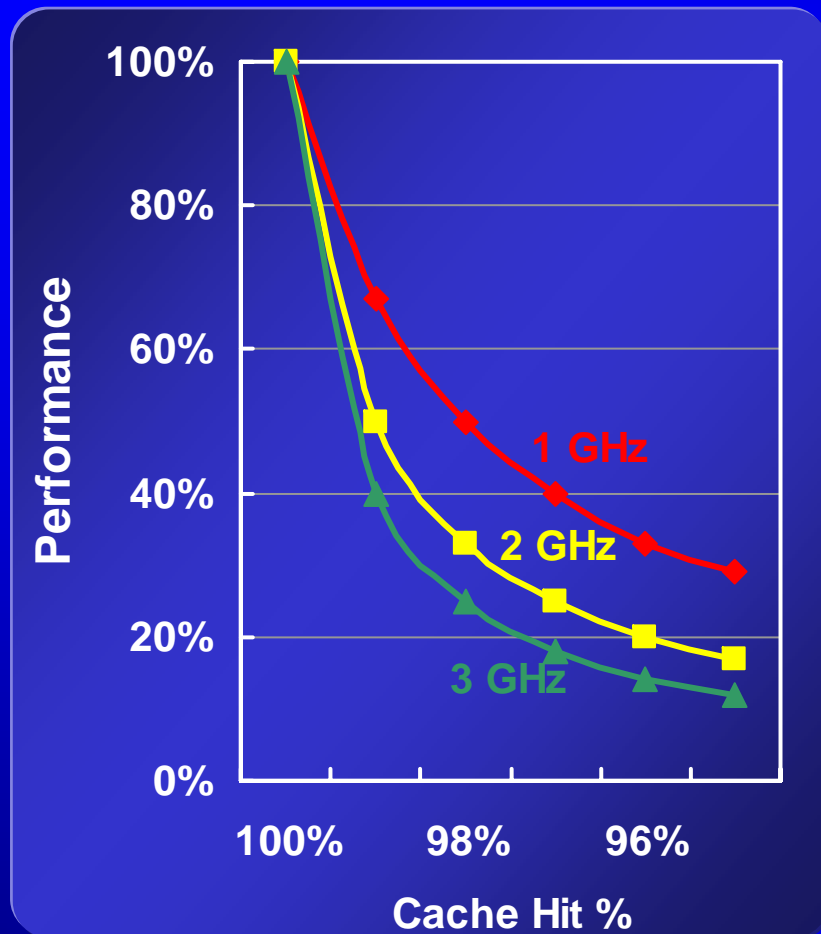
Increase on-die Memory



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

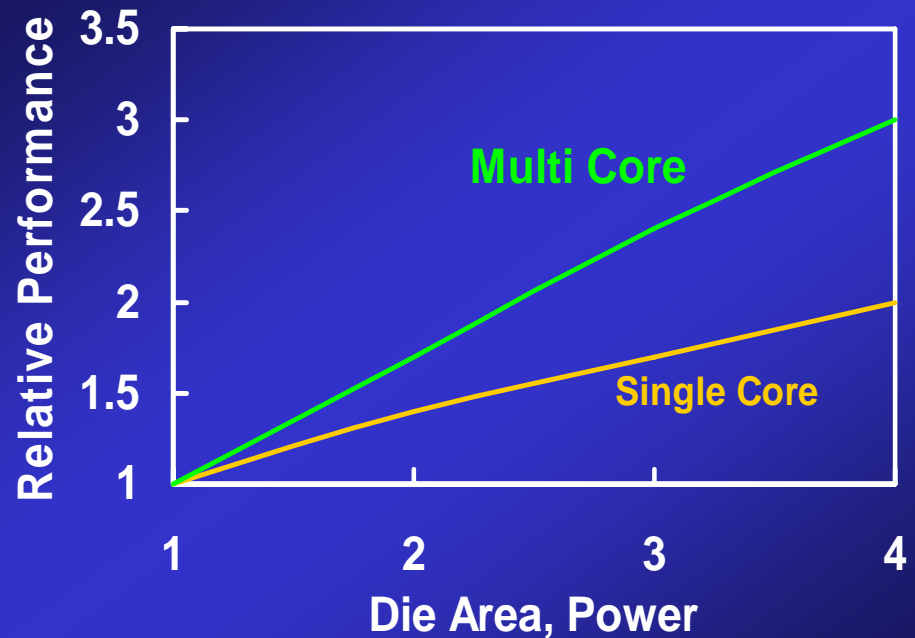
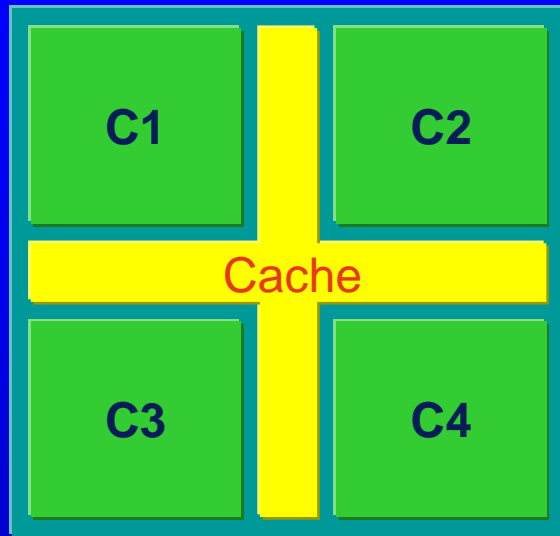
MT2

Wait

MT3

Multi-threading improves performance without impacting thermals & power delivery

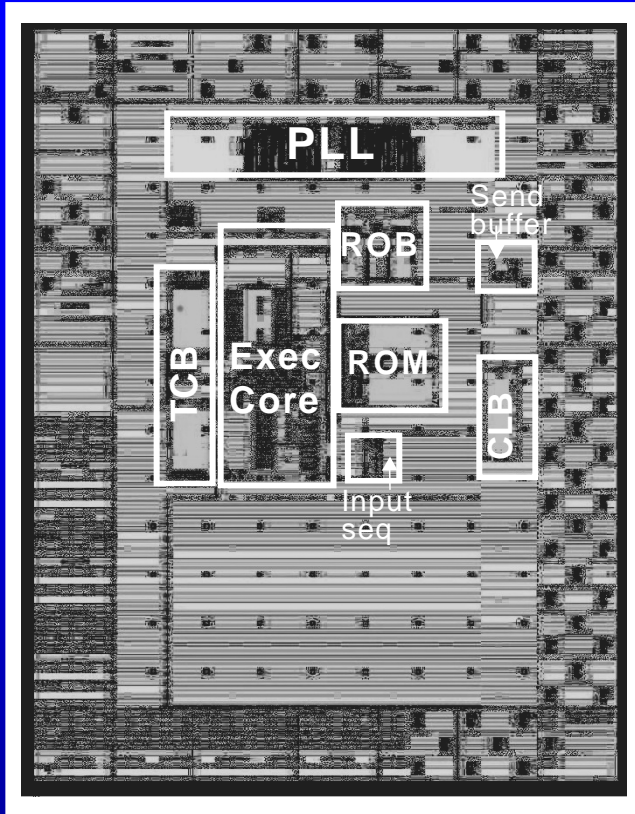
Chip Multi-Processing



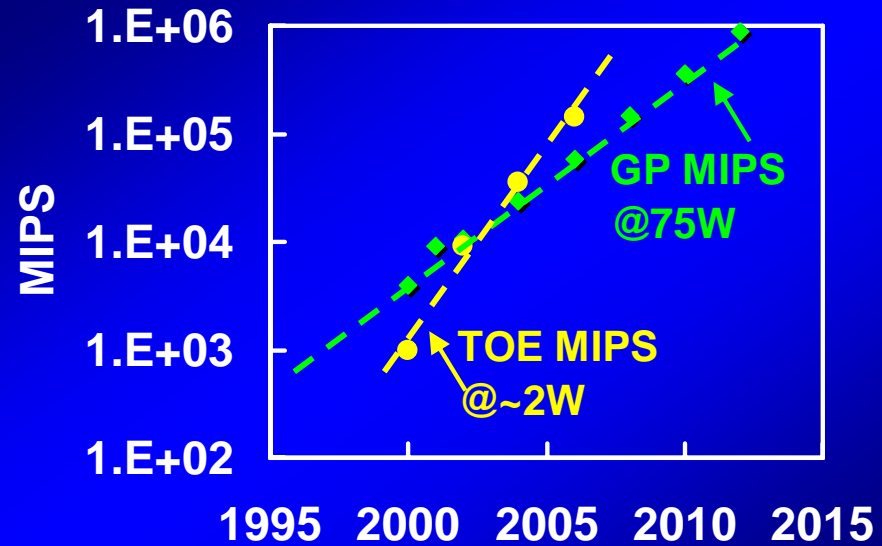
- Multi-core, each core Multi-threaded
- Shared cache and front side bus
- Each core has different Vdd & Freq
- Core hopping to spread hot spots
- Lower junction temperature

Special Purpose Hardware

TCP Offload Engine



2.23 mm X 3.54 mm, 260K transistors



Opportunities:

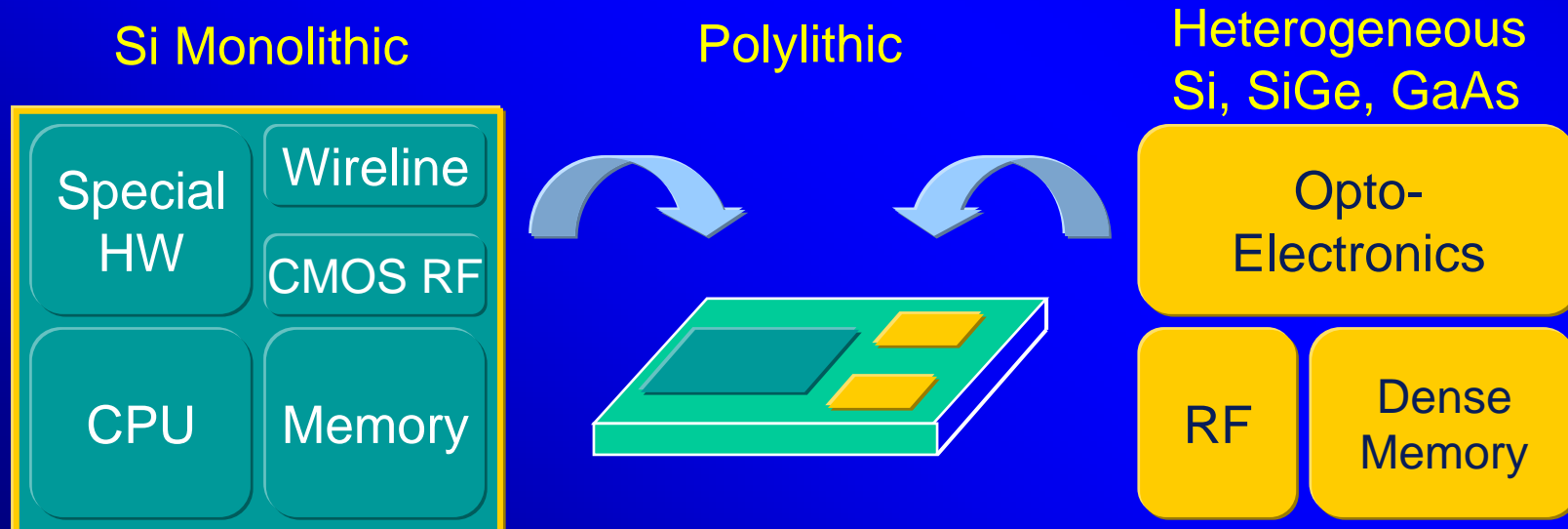
- Network processing engines
- MPEG Encode/Decode engines
- Speech engines

Special purpose HW—Best Mips/Watt

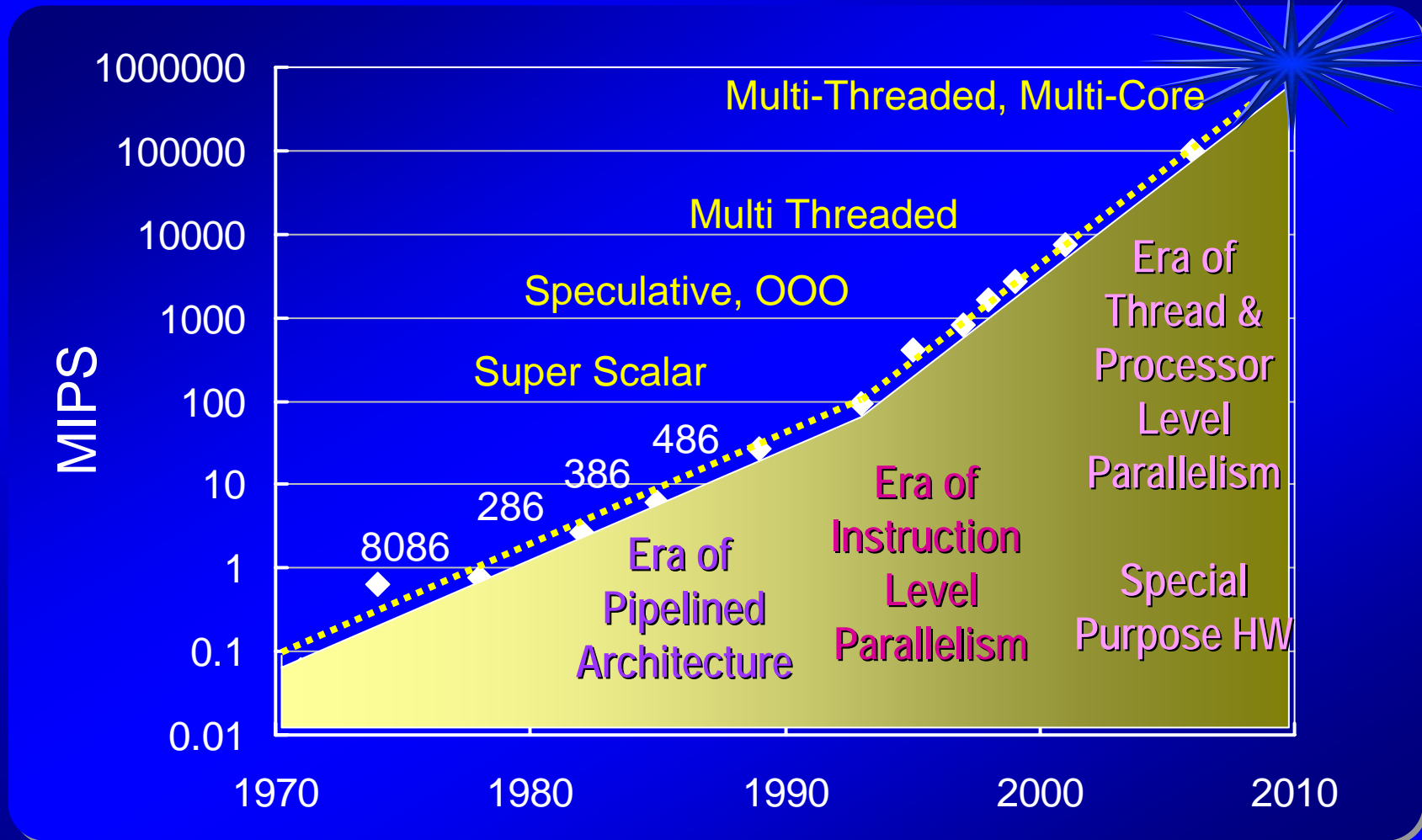
Valued Performance: SOC (System on a Chip)

- Special-purpose hardware → more MIPS/mm²
- 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



Roadmap to TIPS...



Multi-everywhere: MT, CMP

Summary

- ***Business as usual* is not an option**
 - Performance at any cost is history
 - Move away from frequency alone to deliver performance
- **Future μ Architectures and designs**
 - More memory (larger caches)
 - Multi-threading
 - Multi-processing
 - Special purpose hardware
 - Valued performance with higher integration