# The Future Is Low Power

## T. W. Williams, Ph.D.
## Synopsys Fellow

the future **is** low power
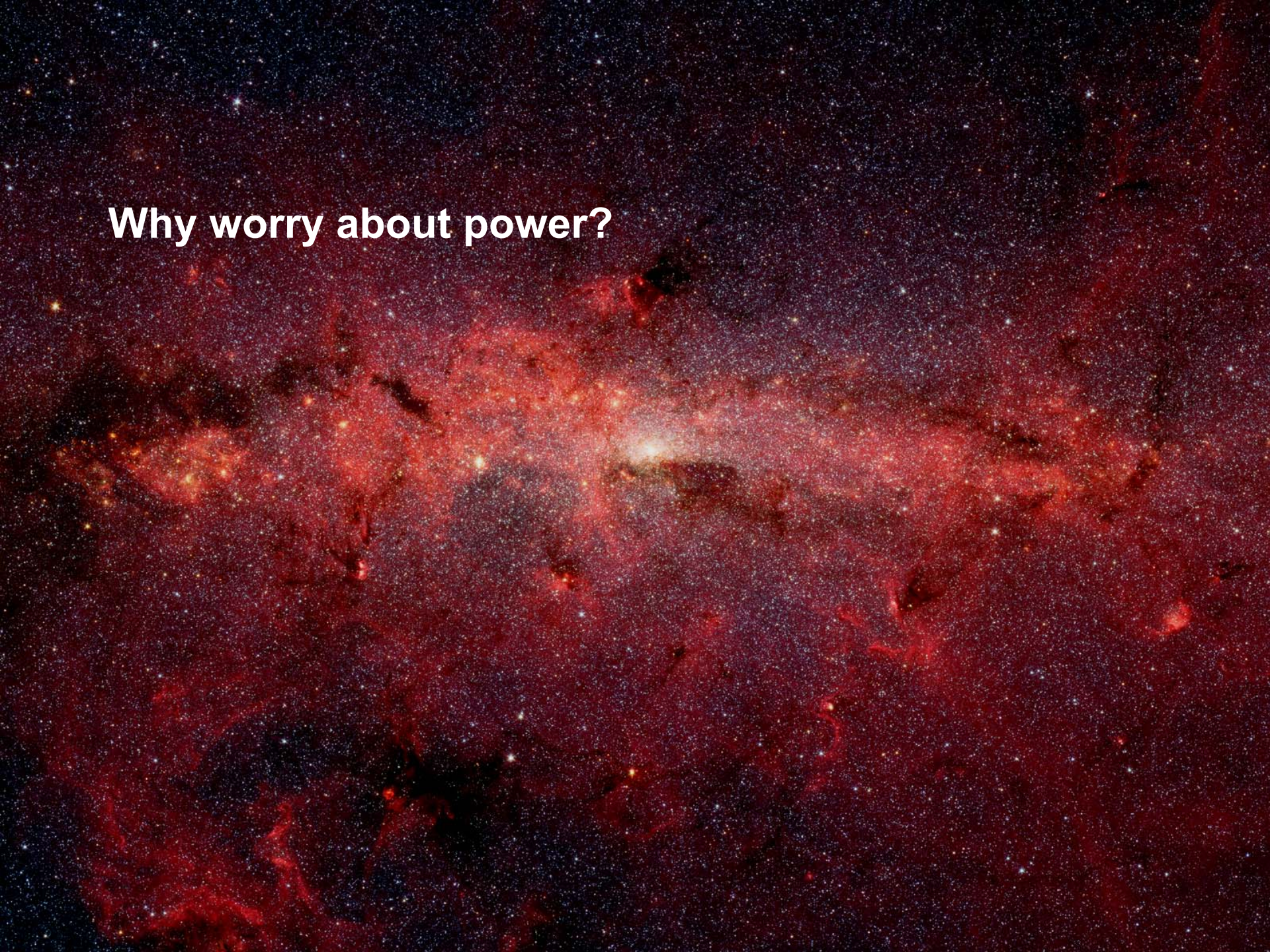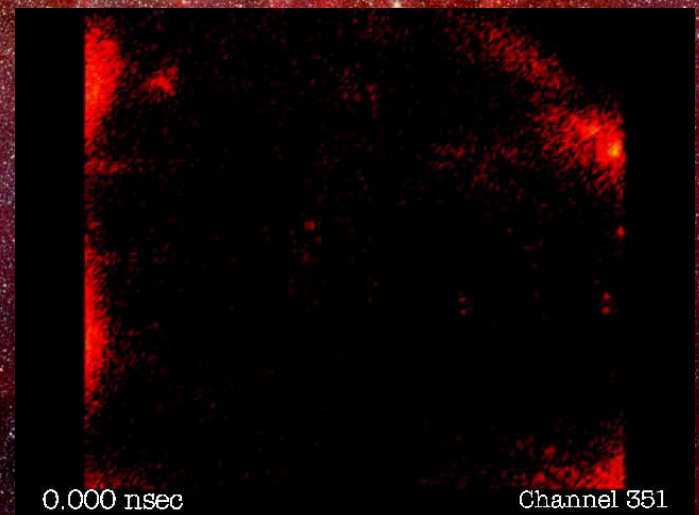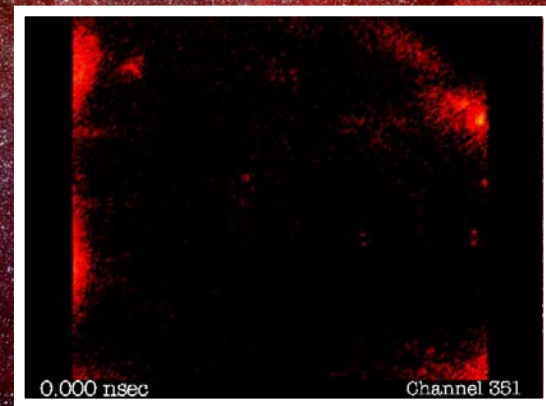
Why worry about power?

- Well, the total energy of the Milky Way galaxy is $\sim 10^{59}$ **Joules**
- The minimum energy required to switch a digital gate (1 electron @ 0.1 Volt) is $\sim 10^{-19}$ **Joules**
- This sets the hard upper bound of possible digital operations to $\sim 10^{78}$ **Joules**

Source: J. Rabaey, UCB 2005

- Now, assuming 1 billion Core 2 Duo Extreme at work all over the world, this is ~$10^{26}$ digital operations per year

- "Moore's Law is alive and well", isn't it? i.e. computing power doubles every 2 years

- **All galaxy's energy will be exhausted in ~343 years from now**

0.000 nsec                    Channel 351

# OK, Let's worry about power!

# The Era Of The Consumer
## *Households, Semiconductor Industry's #1 Customer*



Source: J.-H. Huang, SIA 2004

# Nanometer Market
## *4.5 Billions Customers Already Within Reach*



**Mobile Phone Subscriptions**

100Million

3.2 Billion

1.1B
$142B

| 1997 | 1999 | 2001 | 2004 | 2007E | 2010E |
|------|------|------|------|-------|-------|
| 250nm | 180nm | 130nm | 90nm | 65nm | 45nm |

Source: IC Insights 2007

**SYNOPSYS®**
**Predictable Success**

# Power Consumption & Supply Trends
## *The Battery Capacity Gap*



Power Consumption @ Max. Transmitter Output Power [W]
Battery Capacity Index [%]

Approximate power limit of a handheld device

Multimedia Call Space

Battery Capacity Gap or Reduced Operation Time

Nokia 6000 Series Introduced

An Enabler of Size Shrink

GPRS

Battery Capacity

GSM Voice Call

WCDMA Voice Call or Stand-Alone Application Space

Source: M. Ryynänen, Nokia 2005

# Where Do We Stand? Well…

**"Houston, We Have A Problem"**
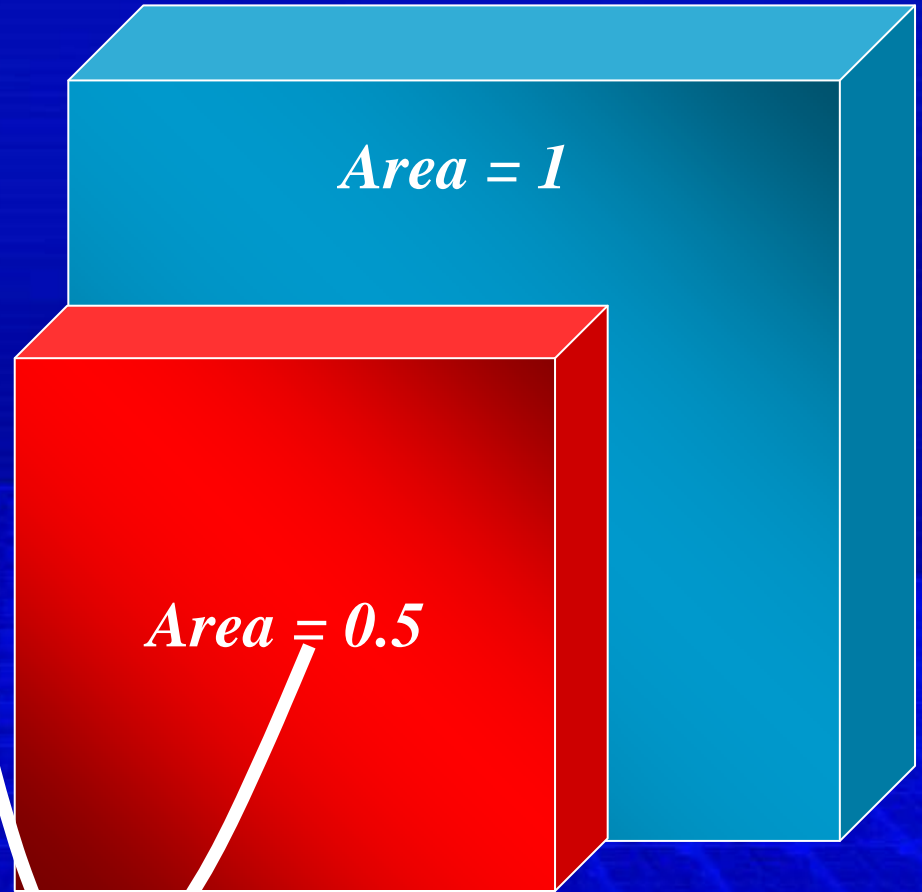
Source: NASA 1970

# Dr. Gordon E. Moore's Law
## *Integration's Capacity Doubles Every Two Years*
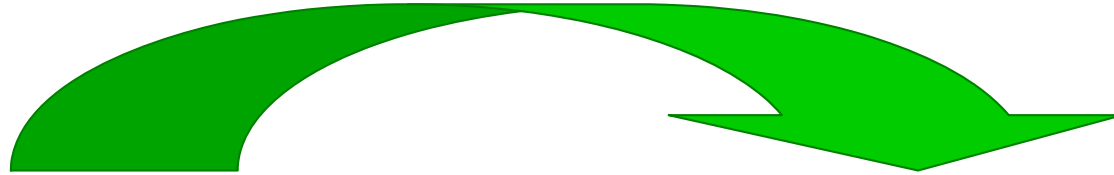
$$\sqrt{0.5} = \sim 0.7$$

## *The Scaling Factor*

*"The complexity for minimum component costs has increased at a rate of roughly a factor of two per year ... Certainly over the short term this rate can be expected to continue, if not to increase. Over the longer term, the rate of increase is a bit more uncertain, although there is no reason to believe it will not remain nearly constant for at least 10 years." Gordon E. Moore, Electronic Magazine, April 19th, 1965*

*Area = 1*

*Area = 0.5*

# Semiconductor Industry Cycle
## *Technology Advances, Market Grows*

⬆ **Density 2X**

⬆ **Performance 1.4X**

➡ **Power 1X**

⬇ **Cost 0.5X**

⬆ **Applications**

⬆ **Units**

⬆ **Users**

⬆ **Revenue**

**SYNOPSYS®**
**Predictable Success**

# The Signs of Crisis Are Visible Already
## *Voltage Has Broken the Rules of Scaling*

| Source: ITRS 2005/2006 | 90nm | 65nm | 45nm |
|---|---|---|---|
| Device Length (nm) ⇩ | 1X | 0.7X | 0.5X |
| Delay (ps) ⇩ | 1X | 0.7X | 0.5X |
| Frequency (GHz) ⇧ | 1X | **1.2X** | **1.45X** |
| Integration Capacity (BTx) ⇧ | 1X | 2X | 4X |
| Capacitance (fF) ⇩ | 1X | 0.7X | 0.5X |
| **Die Size (mm²) ⇨** | 1X | **1X** | **1X** |
| **Voltage (V) ⇲** | 1X | **0.85X** | **0.75X** |
| **Dynamic Power (W) ⇲** | 1X | **> 0.7X** | **> 0.5X** |
| Manufacturing (microcents/Tx) ⇩ | 1X | 0.35X | 0.12X |

# The [not so] Hidden Costs of Scaling
## *New Materials [and Devices] Are Badly Needed*

| Source: ITRS 2005 | 90nm | 65nm | 45nm |
|---|---|---|---|
| $V_{TH}$ (V) ⭦ | 1X | 0.85X | 0.75X |
| $I_{OFF}$ (nA/um) ⇧⇧ | 1X | ~3X | ~9X |
| Dynamic Power Density (W/cm²) ⇧ | 1X | 1.43X | 2X |
| Leakage Power Density (W/cm²) ⇧⇧ | 1X | ~2.5X | ~6.5X |
| Power Density (W/cm²) ⇧ | 1X | ~2X | ~4X |
| Cu Resistance (Ω) ⇧ | 1X | 2X | 4X |
| Interconnect RC Delay (ps) ⇧ | 1X | ~2X | ~5X |
| Packaging (cents/pin) ⭦ | 1X | 0.86X | 0.73X |
| Test (nanocents/Tx) ⇨ | 1X | 1X | 1X |

**SYNOPSYS®**
**Predictable Success**

# Semiconductor Industry Stalemate
## *Technology (CMOS) Shrinks, Doesn't Advance*



**Density 2X**

**Performance 1X**

**Power 2X**

**Cost 0.5X**

**Applications**

**Units**

**Users**

**Revenue**

**SYNOPSYS®**
**Predictable Success**

# Quasi-Atomic-Level Interconnect
## *Main Contributor to both Timing and Dynamic Power*

| Source: ITRS 2005/2006 | 90nm | 65nm | 45nm |
|---|---|---|---|
| ASIC (Gates/mm$^2$) | 500K | 1M | 2M |
| Total Interconnect (Mx + 5 My) (m/cm$^2$) | 1km | 1.4km | 2.2km |
| Interconnect RC Delay (My) (ps/mm) | 355ps | 682ps | 1.8ns |
| $\tau$ = RC Delay (My) (um) | 60um | 38um | 20um |
| On-Chip Frequency (Hertz) | 4.2GHz | 9.3GHz | 15.1GHz |

Source: R. Chau, Intel 2003

SYNOPSYS®
Predictable Success

# Quasi-Atomic-Level Interconnect
## *Main Contributor to both Timing and Dynamic Power*



| Source: ITRS 2005/2006 | 90nm | 65nm | 45nm |
|---|---|---|---|
| ASIC (Gates/mm$^2$) | 500K | 1M | 2M |
| Total Interconnect (Mx + 5 My) (m/cm$^2$) | 1km | 1.4km | 2.2km |
| Interconnect RC Delay (My) (ps/mm) | 355ps | 682ps | 1.8ns |
| $\tau$ = RC Delay (My) (um) | 60um | 38um | 20um |

Source: R. Chau, Intel 2003

SYNOPSYS®
Predictable Success

# Cu Resistivity Increase



Free Mean Path 400-600 A for cupper, Al is > 1000 A

# Atomic-Level Lithography
## *Main Contributor to both Leakage Power & Variability*



| Source: ITRS 2005/2006 | 90nm | 65nm | 45nm |
|---|---|---|---|
| ASIC (Gates/mm$^2$) | 500K | 1M | 2M |
| $T_{ox}$ (# of $SiO_2$ Molecules) | 4 - 5 | 3 - 4 | 2 - 3 |
| $V_{DD}$ Variability (%) | 10% | 10% | 10% |
| $V_{TH}$ Variability (%) | 20% | 30% | 40% |
| Power Variability (%) | 50% | 55% | 60% |

Source: R. Chau, Intel 2003

SYNOPSYS®
Predictable Success

# Molecules/Atoms Make The Difference
*50% Power Variability Already*

© Chandu Visweswariah, IBM 2004, ITRS 2005/2006

**SYNOPSYS®**
**Predictable Success**

# SiO$_2$ Dielectric/Polysilicon Gate
## *High Performance vs. Low Power Dilemma*



**Subthreshold Leakage**

$V_{TH}$ ⇩        $I_{OFF}$ ⇧

Junction Leakage

$I_{JUNC}$

**Gate Leakage**

$T_{OX}$ ⇩        $I_{GATE}$ ⇧

Source: R. Chau, Intel 2003

# 65 Nanometers, Ultra-Low Power…
## It's Actually 130nm ($T_{OX}$) 90nm ($L_{GATE}$) High Performance

Source: P. Bai, Intel, IEDM 2004; C.H. Jan, Intel, IEDM 2005

**SYNOPSYS**
**Predictable Success**

# High-k Dielectric & Metal Gate
## @ 45 Nanometers



Gate

Gate

1.2nm SiO$_2$

3.0nm High-k

Silicon substrate

Silicon substrate

Source: P. Gargini, Intel 2004

# Is High-k Dielectric *The* Solution?
## Benefits of Intel's 45nm Technology vs. 65nm One

> 20% improvement in transistor switching speed

OR

5X $I_{OFF}$ reduction

10X $I_{GATE}$ reduction

> 30% dynamic power reduction

# Power And Power Density

*Increase Exponentially! and*
*Leakage Does Increase Faster than Dynamic*

- $k = 0.7$, $n = 1, 2, 3, \ldots$
- **Dynamic Power Density $\approx O(k^{-n})$**
- **Leakage Power Density**
- $\approx O(1.5n(1+k)^n)$

250nm    180nm    130nm    90nm    65nm    45nm

Source: ITRS 2005/2006

**SYNOPSYS**
**Predictable Success**

# Probability Of Need For Design Re-Spin
## @ 45nm 50% of Design Re-Spins Due to Leakage



Legend: 130nm, 90nm, 65nm, 45mm, <= 32nm

Probability of Need for Design Re-Spin

Source: IBS 2007

**SYNOPSYS®**
**Predictable Success**

# From Exponential To Asymptotic
*Actual Number of Gates/mm² Increase, But Slower!*



Available Gates per mm2

Source: IBS 2007

SYNOPSYS®
Predictable Success

# It's All About [In]Efficiency
## *Gates Utilization as % of Available Gates/mm² Declines*



Gates Utilization as % of Available Gates/mm2

Source: IBS 2007

SYNOPSYS®
Predictable Success

# The Nanometer Application
## *The Most Obvious Solutions Aren't Necessarily the Right Ones*

# Design For Low Power
## *All the Right Ingredients Are Available*

| | Constant Throughput/Latency | | Variable Throughput/Latency |
|---|---|---|---|
| | Design Time | Non-Active Modules | Run Time |
| Dynamic & Short Circuit | Logic Re-Structuring, Logic Sizing Reduced $V_{DD}$ Multi-$V_{DD}$ **2.5X** | Clock Gating **2X** | Dynamic or Adaptive Frequency & Voltage Scaling **2.5X** |
| Leakage | Stack Effect + Multi-$V_{TH}$ **2X-10X** | Sleep Transistors Multi-$V_{DD}$ Variable $V_{TH}$ **10X-1000X** | Variable $V_{TH}$ **2X-10X** |

Source: J. Rabaey, UCB 2005

# CMOS Energy & Power Equations
## *Dynamic, Short Circuit & Leakage*

$$\text{Energy} = C_L V_{DD}^2 P_{0 \to 1} + t_{sc} V_{DD} I_{sc} P_{0 \to 1} + V_{DD} I_{leakage}$$

$$f_{0 \to 1} = P_{0 \to 1} \text{ frequency}$$

$$\text{Power} = C_L V_{DD}^2 f_{0 \to 1} + t_{sc} V_{DD} I_{sc} f_{0 \to 1} + V_{DD} I_{leakage}$$

**Dynamic Power**
50% @ 65nm
decreasing <u>relatively</u>

**Short Circuit Power**
1-2% @ 65nm,
decreasing <u>absolutely</u>

**Leakage Power**
50% @ 65nm,
increasing

# Dynamic Power Consumption
## *Depends on Output Load, $V_{DD}$ & Switching Activity*



Energy = $C_L * V_{DD}^2 * P_{0 \to 1}$          $f_{0 \to 1} = P_{0 \to 1} *$ frequency

Power = $C_L * V_{DD}^2 * f_{0 \to 1}$

**Transition Probability & Frequency Dominate**

# Dynamic Power Reduction
## *Voltage & Frequency Scaling*

**Operating System Load Monitor (SW)**

**Voltage/ Frequency Tables**

CPU Voltage

1.3V

1.2V

1.1V

100%

85%

62%

28% energy saving

55% energy saving

# Dynamic Power Reduction
## *Multi-Supply, Multi-Voltage, Dynamic & Adaptive Voltage Scaling*



**(1)**

Voltage Island — **A** — 1.2 V, 350 MHz
Voltage Island — **B** — 1.0 V, 250 MHz
Voltage Island — **C** — 1.5 V, 500 MHz

**(2)**

Programmable → Voltage Island **A**, Voltage Island **B**, Mode Control → Voltage Regulators, Voltage Island **C**

**(3)**

Voltage Island — Monitor — **A**
Voltage Island — Monitor — **B**
Mode Control → Voltage Regulators
Voltage Island — Monitor — **C**

**Adaptive Voltage Scaling (AVS)**

- Voltage areas with variable $V_{DD}$
- Software controlled

**Dynamic Voltage Scaling (DVS)**

- Voltage areas with multiple, but fixed voltages
- Software controlled

**Multiple Supply Multi-Voltage (MV) Islands**

- Voltage areas with fixed, single voltages

**SYNOPSYS®**
Predictable Success

# Short-Circuit Power Consumption
## *Depends on Input Slope, Output Load & $V_{DD}$*



Energy = $\mathbf{t_{sc}}$ * $\mathbf{V_{DD}}$ * $\mathbf{I_{sc}}$ * $\mathbf{P_{0\to1}}$    $f_{0\to1}$ = $P_{0\to1}$ * frequency

Power = $\mathbf{t_{sc}}$ * $\mathbf{V_{DD}}$ * $\mathbf{I_{sc}}$ * $\mathbf{f_{0\to1}}$

## Today, Interconnects Capacitance Dominates

# Impact Of $C_L$



**Large** capacitive load

**Small** capacitive load

Output fall time significantly larger than input rise time.

Output fall time substantially smaller than input rise time.

**Negligible When $V_{DD} < V_{THn} + |V_{THp}|$, Slope Engineering**

# Leakage Power Consumption
## *Depends on Sub-Threshold Current & $V_{DD}$*



Drain junction leakage

Gate leakage

Sub-threshold current

$V_{OUT}$

Energy = Power = **$V_{DD} * I_{leakage}$**

**Today Sub-Threshold Current, But Gate Leakage Coming**

**SYNOPSYS**®
**Predictable Success**

# Leakage Power Reduction
## *Body Bias (VTCMOS) and (SRPG)/MTCMOS*



Operating System Load Monitor (SW)

Active clock ratio table

CPU Voltage

1.0V

100%

85%

62%

15% energy saving

38% energy saving

**SYNOPSYS®**
**Predictable Success**

# Leakage Power Reduction
## *Body Bias (VTCMOS) vs. (SRPG)/MTCMOS*



Save    SRPG    Restore

Short Stop

Long Stop

Body Bias (VTCMOS)

MTCMOS

Reduced Leakage

Zero Leakage

# Body Bias (VTCMOS) Vs. MTCMOS



| Body Bias (VTCMOS) | MTCMOS |
|---|---|
| **$V_{TH}$ Control With Substrate Bias** | **On/Off Control of Internal $V_{DD}/V_{SS}$** |
| Needs Circuit Development | Conceptually Simpler |
| Compensate $\Delta V_{TH}$ Fluctuation | Compensate $\Delta V_{TH}$ Fluctuation |
| IDDQ Test | IDDQ Test |
| No Serial MOSFET | Large Serial MOSFET |
| Conventional EDA Tools | Conventional EDA Tools |
| Re-Use of Existing Design | Retention Registers |
| Triple Well Desirable | Conventional Well structure |

Source: J. Rabaey, UCB 2005

# Adaptive Voltage Scaling Vs. Adaptive Body Bias



**Adaptive Voltage** **Adaptive Body Bias**

**Adaptive Body Bias** **Adaptive Voltage & Body Bias**

Source: S. Borkar, Intel 2004

SYNOPSYS®
Predictable Success

# Energy Profiles
## *Power Gating (MTCMOS) with State Loss*

# Energy Profiles
## *Power Gating (MTCMOS) with State Retention*
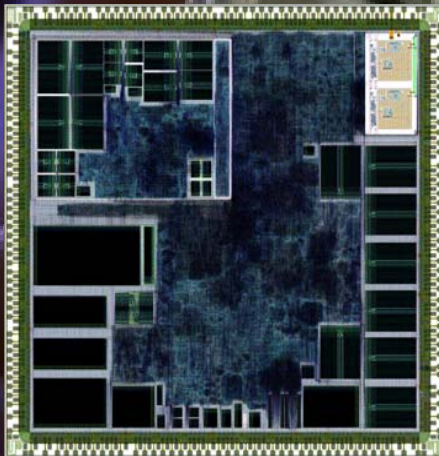
# Design For Low Power Helps!
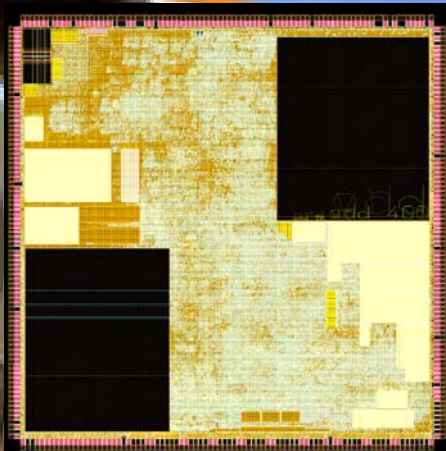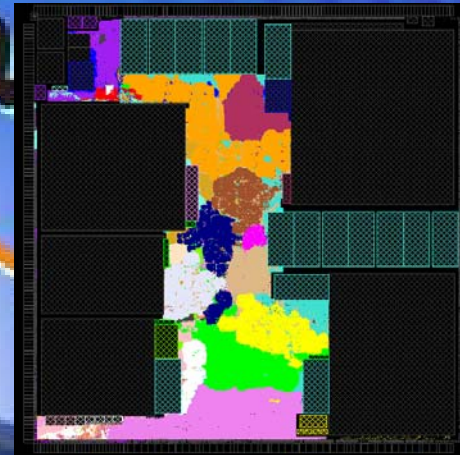## *Generations Of Low Power,*
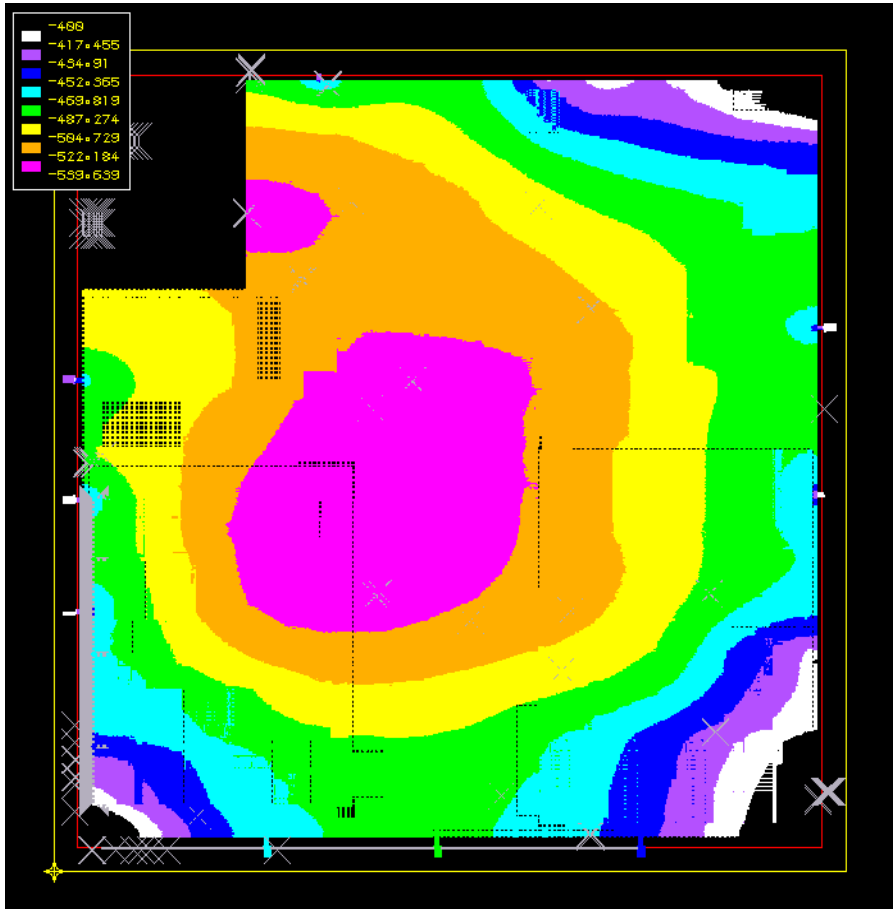## *First Time Silicon Success*

**2007**

**2006**

**2005**

**2003**

- 65 nanometers
- Clock Gating
- Multi-$V_{TH}$ & Nested Multi-Supply

- 90 nanometers
- Clock Gating
- Multi-$V_{TH}$ & Multi-Supply

- 130 nanometers
- Clock Gating
- Multi-$V_{TH}$ & Multi-Supply

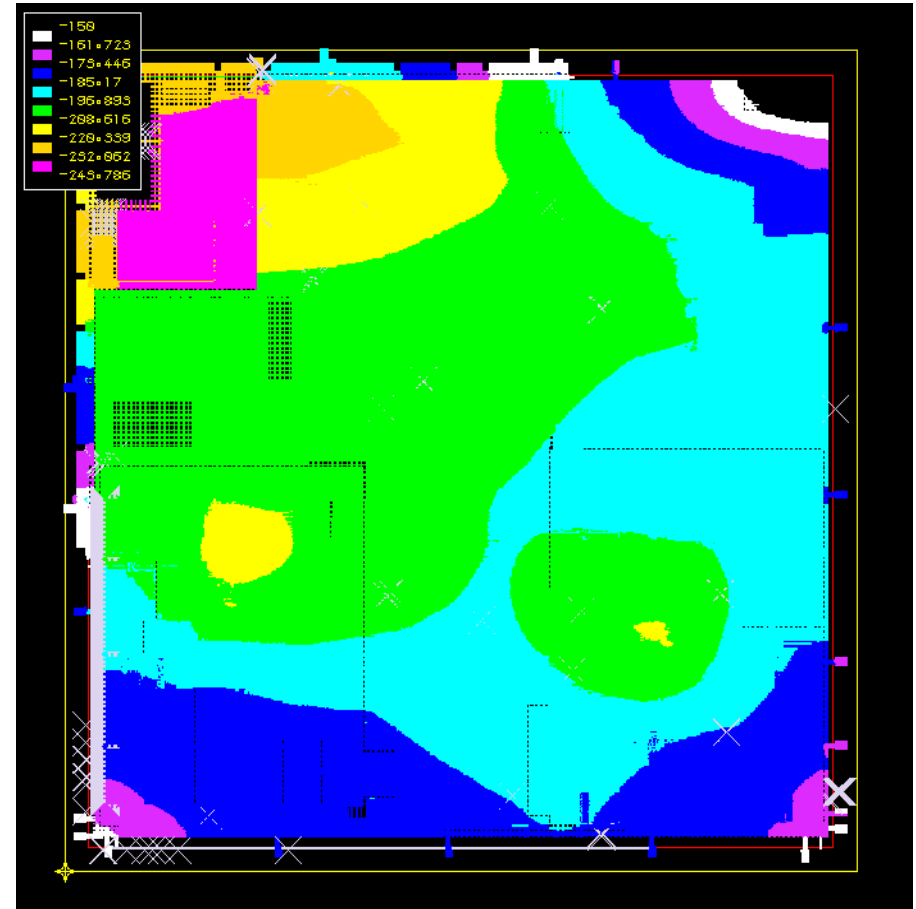- 130 nanometers
- Clock Gating
- Multi-$V_{TH}$

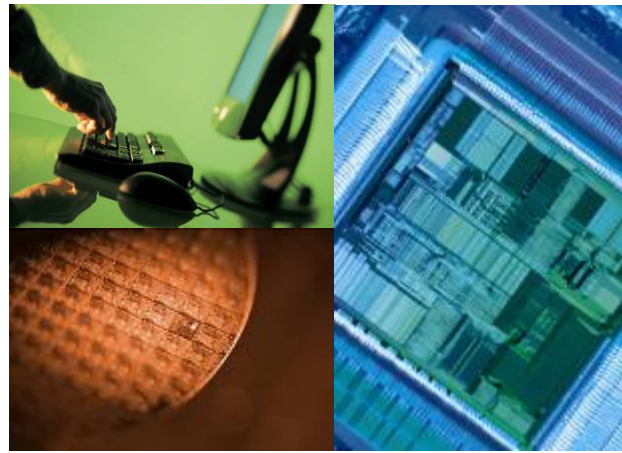# Design For Low Power Helps!
## *Voltage Drop in Implementation & Sign-Off*



*Without packaging parasitics*
*- Range (-400mV, -540mV)*

*With packaging parasitics + $C_{EXT}$*
*(330nf) - Range (-150mV, -243mV)*

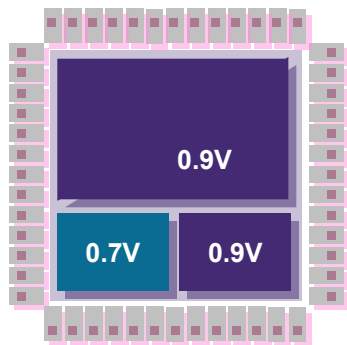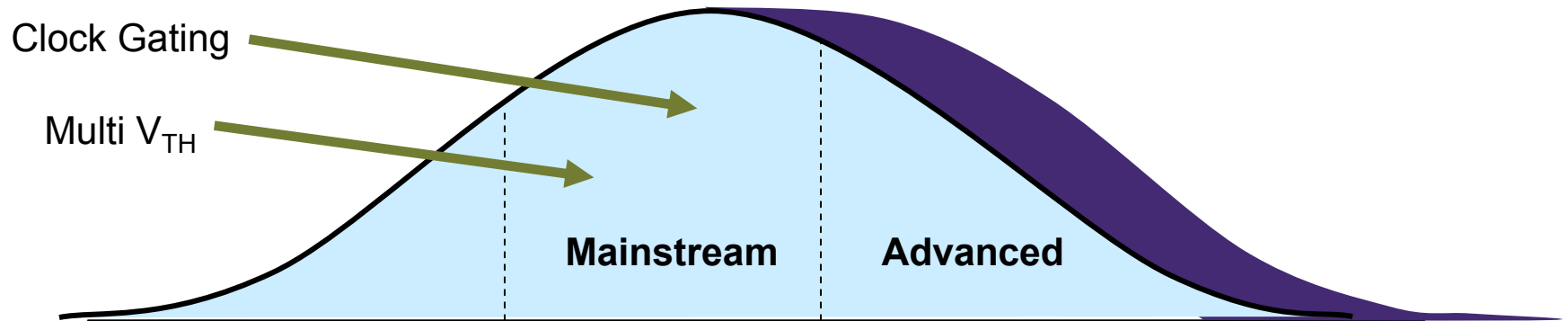Source: European IDM, Multimedia Application, 130nm

**SYNOPSYS®**
**Predictable Success**

# The Future Is Low Power

## T. W. Williams, Ph.D.
## Synopsys Fellow

**Predictable Success**

# Power Management Techniques



Clock Gating

Multi $V_{TH}$

Mainstream

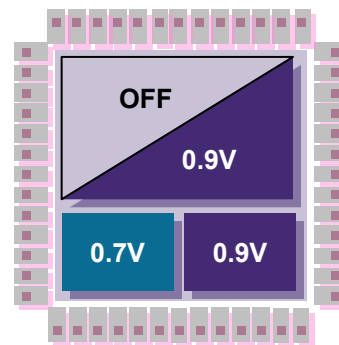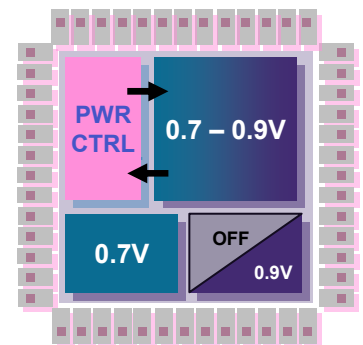Advanced

| Multi-Voltage | Power Gating (MTCMOS) | Multi-Voltage & Power Gating | Dynamic or Adaptive Voltage Scaling and Power Gating |

0.9V · 0.7V · 0.9V

OFF · 0.9V · 0.9V · 0.9V

OFF · 0.9V · 0.7V · 0.9V

PWR CTRL · 0.7 – 0.9V · 0.7V · OFF · 0.9V

**SYNOPSYS®**
**Predictable Success**