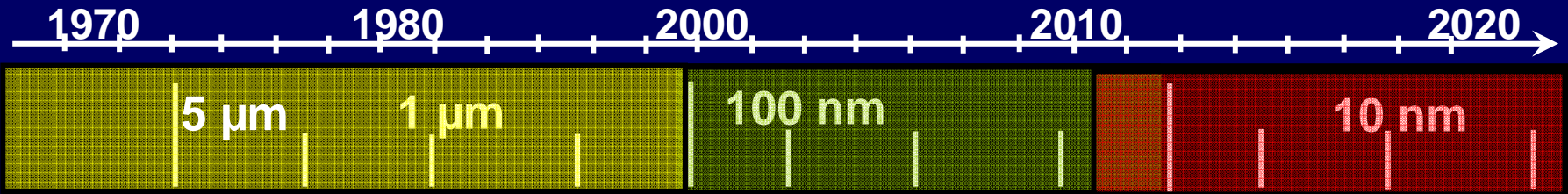# *Process-Tolerant Low-Power Design for the Nano-meter Regime*

## Kaushik Roy

### Electrical & Computer Engineering

### Purdue University

# Exponential Increase in Leakage

1970      1980      2000      2010      2020

| 5 µm    1 µm | 100 nm | 10 nm |

**Silicon Micro- electronics**

**Silicon Nano- electronics**

**Non-Silicon Technology**

$$\frac{I_{ON}}{I_{OFF}} = 10^6$$

$$\frac{I_{ON}}{I_{OFF}} = 10^3$$

$$\frac{I_{ON}}{I_{OFF}} \sim 10^{2\sim6}$$

**Subthreshold Leakage**

**Gate Leakage**

**Gate**

**Source**

**Drain**

n+          n+

**Junction leakage**

**Bulk**

**Must stop at 50%**

A. Grove, IEDM 2002

Leakage Power (% of Total)

50%
40%
30%
20%
10%
0%

1.5   0.7   0.35   0.18   0.09   0.05

**Technology (µ)**

# Technology Trend

2003

2009

2020

## Single gate device

### Bulk-CMOS

Source/Drain · S/D Extension · Gate · P · N+ · Retrograde Well · Halo

### PD/SOI

$V_G$ · Gate · $V_S$ · $V_D$ · Source · Floating Body · Drain · Buried Oxide (BOX) · Substrate

### FD/SOI

Fully-depleted body · $V_G$ · Gate · $V_S$ · $V_D$ · Source · Drain · Buried Oxide (BOX) · Substrate · $V_{back}$

## Multi-gate devices

### DGMOS

source · top gate · drain · $SiO_2$ · $SiO_2$ · SEG · Silicon (100)

### FinFET

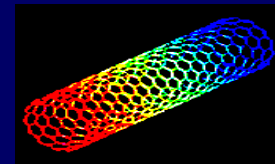Gate · Source · Drain · BOX · Si fin - Body!

### Trigate

Top Gate · Side Gate · Side Gate

## Nano devices

Carbon nanotube
III-V devices
nano-wires
Spintronics

Source · Drain · Nanowire · Oxide · Gate

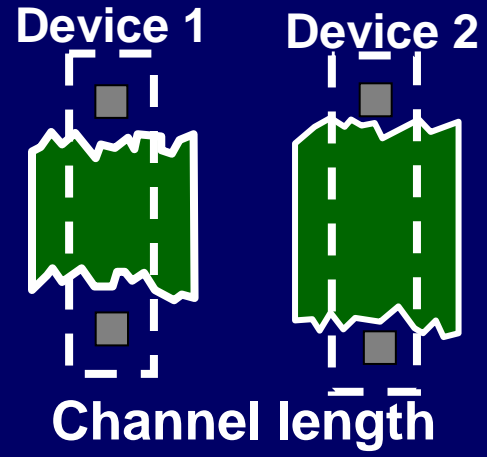## Design methods to exploit the advantages of technology innovations

# Variation in Process Parameters

**Device 1**   **Device 2**

**Channel length**

Delay and Leakage Spread

30%   130nm   Source: Intel   5 X

Normalized Frequency / Normalized Leakage (Isb)

Source: Intel

# dopant atoms / Technology Node (nm)

**Inter and Intra-die Variations**

**Random dopant fluctuation**

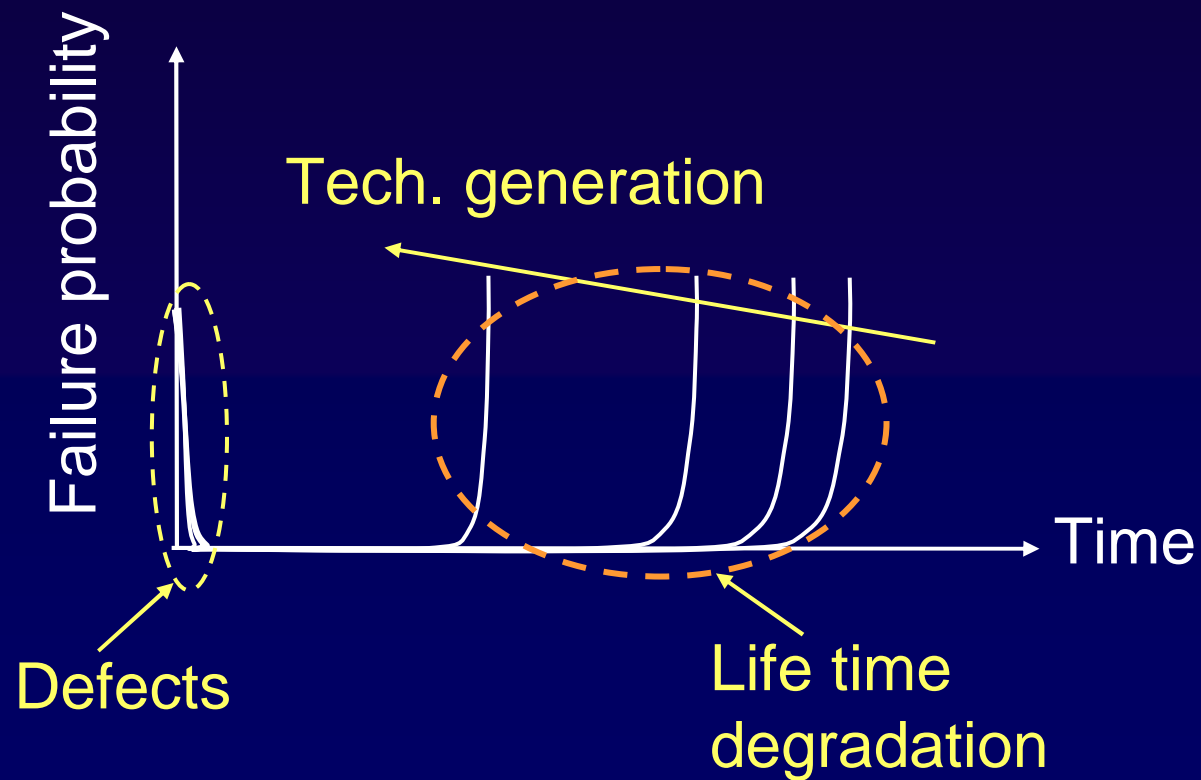**Device parameters are no longer deterministic**

# Reliability

Temporal degradation of performance -- NBTI

# Pessimistic Design Hurts Performance



- **Substantial variation in leakage across dies**
- **4X variation between nominal and worst-case leakage**
- **Performance determined at nominal leakage**
- **Robustness determined at worst-case leakage**

# Global and Local Variations

**Random Dopant Fluctuation**



$\delta V_{t-LOCAL}$

$\sigma_{LOCAL}$

**intra-die**

$\Delta V_{t-GLOBAL}$

$\sigma_{GLOBAL}$

**inter-die**

$$\delta V_t = \Delta V_{t-GLOBAL} + \delta V_{t-LOCAL}$$

# Process Tolerance: Memories

*S. Mukhopadhaya, Mahmoodi, Roy*
*VLSI Circuit Symposium 2006, JSSC 2006, TCAD*

# Parametric Failures: Read Failure



$$P_{RF} = P\left(V_{READ} > V_{TRIPRD}\right)$$

**Read failure => Flipping of Cell Data while Reading**

# Parametric Failures in SRAM

WL

PL   PR

'1'   '0'

AXL   AXR

NL   NR

BL   BR

High-Vt   Low-Vt

## Parametric failures

- – Read Failures
- – Write Failures
- – Access Failures
- – Hold Failures

Test & Repair using Redundancy

**Faulty chips**   **Working chips**

# Parametric failures can degrade SRAM yield

# Process Variations in On-chip SRAM

**Yield ≈ 33%**

■ Fault statistics

$\sigma_{Vt} \approx$ **30mv, using BPTM 45nm technology**

**Simulation of an 64KB Cache**

A. Agarwal, et. al, JSSC, 05

**Chip Count** (y-axis): 0, 50, 100, 150, 200, 250, 300, 350

**Number of faulty cells ($N_{Faulty\text{-}Cells}$)** (x-axis): 0, 52, 105, 157, 210, 262, 315, 367, 419, 472, 524, 577, 629, 682, 734, 786, 839, 890, 944, 996, 1049

# Parametric failures →Yield degradation

# Inter-die Variation & Cell Failures



## Low–Vt Corners
– **Read failure** ↑
– **Hold failure** ↑

## High–Vt Corners
– **Access failure** ↑
– **Write failure** ↑

$\sigma_{GLOBAL}$

inter-die Vt shift ($\Delta V_{\text{th-GLOBAL}}$)

# Inter-die Variation & Memory Failure



**Memory failure probabilities are high when inter-die shift in process is high**

# Post-Silicon Repair: Proposed Approach



$\sigma_{LOCAL}$   $\sigma_{LOCAL}$

**intra-die** **intra-die**

$\sigma_{GLOBAL}$

**inter-die**

**Apply correction to the global variation to reduce number of failures due to local variations**

# How to identify the inter-die Vt corner under a large intra-die variation ?



Monitor circuit parameters, e.g. leakage current

Effect of inter-die variation can be masked by intra-die variation

# Array Leakage Monitoring

$$Y = \sum_{i=1}^{N} X_i \Rightarrow \frac{\sigma_Y}{\mu_Y} = \frac{1}{\sqrt{N}} \frac{\sigma_X}{\mu_X}$$

- **Adding a large number of random variables reduces the effect of intra-die variation**

**Leakage of entire SRAM array is a reliable indicator of the inter-die Vt corner**

# Self-Repair using Leakage Monitoring

**Bypass Switch**

**Calibrate Signal**

$V_{DD}$

**On-chip Leakage Monitor**

$V_{out}$

$V_{REF1}$  $V_{REF2}$

**Comparator**

**SRAM Array**

**Body bias**

**Body-Bias selection**

FBB  ZBB  RBB

$V_{OUT}$

**SRAM ARRAY**

**Entire array leakage is monitored to detect inter-die corner and proper body-bias is selected**



BPTM 70nm

LVT

$V_{REF1}$

Nom. Vt

$V_{REF2}$

HVT  Nom. Vt  LVT

Sensor Output [V]

Memory Leakage [A]

# Yield Enhancement using Self-Repair



**Self-Repairing SRAM using body-bias can significantly improve design yield**

# Test-Chip of Self-Repairing SRAM

**VCO**

**16 KB block**

**64 KB LVT Array**

**Isolated cell**

**Sensor + Ref. gen.**

**BB gen**

simulations in 0.13μm CMOS — Conv. SRAM — Self-repair-SRAM

Total # of failures

Low-Vt   High-Vt

**Technology : IBM 0.13 μm**

**128KB SRAM**

**Dual-Vt Triple-well tech.**
**Number of Trans: ~ 7 million**
**Die size: 16mm2**
**VLSI CKT Symp. 2006, ITC 2005**

**Simulation results for 1MB array designed in IBM 0.13μm**

# Continuous vs Quantized Body Bias



**Quantized (3 Level: FBB, ZBB, RBB) body bias scheme is a cost effective solution with good yield enhancement possibility**

# Process Tolerance: Register Files

***Kim et. al. VLSI Circuit Symposium 2004***

# Process Compensating Dynamic Circuit Technology



Conventional Static Keeper

clk

RS0  RS1  RS7

D0  D1  D7

LBL0  N0  LBL1

- Keeper upsizing degrades average performance

# Process Compensating Dynamic Circuit Technology



C. Kim et al. , VLSI Circuits Symp. '03

● **Opportunistic speedup via keeper downsizing**

# On-Die Leakage Sensor For Measuring Process Variation



C. Kim et al. , VLSI Circuits Symp. '04

- **High leakage sensing gain – 90nm dual-Vt, Vdd=1.2V, 7 level resolution, 0.66 mW @80C°**

# Leakage Binning Results

Output codes from leakage sensor

# Self-Contained Process Compensation

**Fab**

**Wafer test**

*Process detection*

Leakage measurement

On-die leakage sensor

Program PCD using fuses

**Customer** **Package test** **Burn in** **Assembly**

# Self-Repair: Architecture Level

Agarawal, Roy TVLSI 2005

# Fault-Tolerant Cache Architecture



- **BIST detects the faulty blocks**
- **Config Storage stores the fault information**

**Idea is to resize the cache to avoid faulty blocks during regular operation**

# Mapping Issue

Address "one"  Address "two"  Location   TAG           DATA

"T R 00 Off"      "T R 01 Off"      "R 00" →  [    ]      [ FAULTY ]

STORE D "one"                         ↓ *Mapped by Controller*
LOAD "two" Register      "R 01" →  [ T ]      [ D ]

**Tag matches but wrong data**

**More than one INDEX are mapped to same block**

| TAG | INDEX | Off |
|-----|-------|-----|

**Column Address**

| | | Off |
|--|--|-----|

**New TAG**        **INDEX**

**Include column address bits into TAG bits**

Address "one"  Address "two"  Location   TAG           DATA

"T R 00 Off"      "T R 01 Off"      "R 00" →  [    ]      [ FAULTY ]

STORE D "one"                         ↓ *Mapped by Controller*
LOAD "two" Register      "R 01" →  [ T00 ]      [ D ]

**Tag does not match, cache miss**

**Resizing is transparent to processor → same memory address**

# Fault Tolerant Capability



- **Proposed architecture can handle more number of faulty cells than ECC, as high as 890 faulty cells**

- **Saves more number of chips than ECC for a given $N_{Faulty-Cells}$**

# CPU Performance Loss



For a 64K cache averaged over SPEC 2000 benchmarks

Y-axis: % CPU Performance Loss (0.0, 0.5, 1.0, 1.5, 2.0, 2.5)

X-axis: $N_{Faulty-Cells}$ (0, 105, 210, 315, 419, 524, 629, 734, 839)

- **Increase in miss rate due to downsizing of cache**

- **Average CPU performance loss over all SPEC 2000 benchmarks for a cache with 890 faulty cells is ~ 2%**

# Logic: Process Tolerance

# Logic: A New Paradigm for Low-Voltage, Variation Tolerant Circuit Synthesis Using Critical Path Isolation (CRISTA)

*Ghosh, Bhunia, Roy -- ICCAD 2006*

# Razor Approach

**Standard Latch**

D

CLK

Q

**Shadow Latch**

Delay

E

RAZOR: Dan Ernst et. al., MICRO 2003.

- *Post-Silicon* technique for *dynamic* supply scaling and *timing error* detection/correction
- Error correction overhead is 1% for a 10% error rate

# Vdd Scaling and Process Tolerance: Conventional Solutions

- Low power:

  - Reduce the supply voltage

    - Error rate increases

  - Dual-Vt/dual-VDD assignment

    - Number of critical paths increases

- Robustness:

  - Increase supply voltage

    - Power dissipation increases

  - Upsize the gates

    - Switching capacitance increases

**Low power and robustness: conflicting requirements**

# CRISTA: Basic Idea



- Important points:
  - Scale down the supply while making *delay failures predictable*
  - *Avoid* the failures by *adaptive clock stretching*
  - Ensure that critical paths are activated *rarely*

# Design Considerations for CRISTA



**Design A: conventional design**

**Design B: proposed design**

- Few predictable critical paths
- Low activation probability of critical paths
- Slack between critical and non-critical paths under variations

# Case Study: Adder



$C_{i,0} \rightarrow$ FA $\xrightarrow{C_{o,0}}$ FA $\xrightarrow{C_{o,1}}$ FA $\xrightarrow{C_{o,2}}$ FA $\rightarrow C_{o,3}$

($P_0$ $G_1$, $P_1$ $G_1$, $P_2$ $G_2$, $P_3$ $G_3$ inputs)

- Interesting features:
  - Single critical path (activated by $P_0 P_1 P_2 P_3 = 1$ & $C_{i,0} = 1$)
  - Low activation probability of critical path

VDD = 1V, TCLK = 260ps                    VDD = 0.8V, TCLK = 260ps

| | |
|---|---|
| Crit. path delay=260ps<br>longest non-crit. path delay=165ps<br>*P = 13uW (1-cycle)* | Crit. path delay=330ps<br>longest non-crit. path delay=260ps<br>*P = 7.4uW (rare 2-cycles, decoder)* |

**44% power saving by reducing voltage and, operating critical path at 2-cycle and other paths at 1-cycle**

Can we apply same technique to any random logic?

# Carry Select Adder



- ~20% power saving with ~6% area overhead

# Carry Save Multiplier



Vector Merging Adder

Critical Path

Longest off-critical path

LATENCY PREDICTOR BLOCK

HA    HA    HA

FA    FA    FA

FA    FA    FA

FA    FA    HA

- **25%** power saving with **~5%** area overhead

# Wallace Tree Multiplier

● Partial Products
■ Full Adders
■ Half Adder
■ Vector Merging Adder
  Critical path
  Longest off-critical path

Stage 1

Stage 2

Stage 3

Vector Merging Adder

Final Product

- **29%** power saving with **~4%** area overhead

# Simulation Results

# Random Logic: Shannon's Expansion

$$f(x_1,...,x_i,..., x_n) = x_i \bullet f(x_1,...,x_i=1,..., x_n) + x_i' \bullet f(x_1,...,x_i=0,..., x_n)$$

$$= x_i \bullet CF_1 + x_i' \bullet CF_2$$

$$CF_1 = f(x_1,...,x_i=1,..., x_n); \qquad CF_2 = f(x_1,...,x_i=0,..., x_n)$$



**Activation probability of cofactors can be reduced**
**How to choose *Control Variable* ?**

# Further Isolation and Slack Creation by Sizing



- Slack creation strategy
  - Lagrangian Relaxation based sizing (*B.C. Paul et. al., DAC 2004*) is used
  - Non-critical paths are selectively made faster
  - Critical paths are slightly slowed down

# Simulation Results

**MCNC benchmarks, 70nm Process**

% Imp. in power

- % imp in power with switching activity = 0.2
- % imp in power with switching activity = 0.5

cht  sct  pcle  mux  decod  cm150a  x2  alu2  count

**Power**

**MCNC benchmarks, 70nm Process**

Area (x10³)um^2

- Original design
- Proposed design

cht  sct  pcle  mux  decod  cm150a  x2  alu2  count

**Area**

- Average power saving = ~50%
- Average area overhead = 18%
- Avg performance penalty=5.9% (with 4 control variables) for signal prob=0.5

# Two-Stage Pipeline with Test Logic



~40% power saving with ~13% performance penalty

# VDD Scaling, Process Variation, and Quality Trade-off: DCT

*Banerjee, Karakonstantis, Roy*
*Design Automation and Test in Europe (DATE) 2007*

# Basic Idea

- All computations are "not equally important" for determining outputs

- Identify important and unimportant computations based on output "sensitivity"

- Compute important computations with "higher priority"

- Delay errors due to variations/ Vdd scaling "affect only" non-important computations

- "Gradual degradation" in output with voltage scaling and process variations

# DCT Based Image Compression Process

**8×8 blocks**

*Source image X*

**JPEG Encoder Block Diagram**



**512×512 image**

$$Round\left(\frac{T \cdot Z \cdot T\,'}{Q}\right)$$

FDCT → $Z$ → Quantizer → $V$ → Entropy Encoder → **Compressed Image Data**

$Z = T \cdot X \cdot T\,'$

$X$ → **1D DCT** → $W$ → **Transpose Memory** → $Y$ → **1D DCT** → $Z$

- **DCT is used in current international image/video coding standards**
  - **JPEG, MPEG, H.261, H.263**

# Energy Distribution of a 2D-DCT Output



■ **High energy components (important outputs 75% energy)**
■ **Low energy components (less important outputs)**

**Can important components be computed
with higher priority ?**

# Design Methodology



(a)  Input Block

$T.x^t$

1D-DCT

Faster Computation

Slower Computation

(b)  1D- intermediate DCT outputs

(c) Transpose Memory

1D-DCT

(d) Final DCT outputs

# Path Delays for 1D-DCT outputs



$(x_0 + x_7) \bullet d$
$(x_3 + x_4) \bullet d$
$w_0$
**(2 adders delay)**

$(x_2 + x_5) \bullet d$
$(x_1 + x_6) \bullet d$
$+$
$w_4$
$-$
**(3 adders delay)**

$(x_0 + x_7) \bullet e$
$(x_3 + x_4) \bullet e$
$(x_1 + x_6) \bullet f$
$(x_0 + x_7) \bullet f$
$(x_2 + x_5) \bullet f$
$(x_3 + x_4) \bullet f$
$(x_3 + x_4) \bullet f$
$(x_1 + x_6) \bullet f$
$(x_2 + x_5) \bullet e$
$(x_1 + x_6) \bullet e$
$w_2$
**(3 adders delay)**
$w_6$
**(4 adders delay)**

$(x_0 - x_7) \bullet a$
$(x_1 - x_6) \bullet a$
$(x_1 - x_6) \bullet e$
$(x_2 - x_5) \bullet e$
$(x_1 - x_6) \bullet f$
$(x_3 - x_4) \bullet f$
$(x_2 - x_5) \bullet e$
$(x_2 - x_5) \bullet a$
$(x_3 - x_4) \bullet a$
$(x_2 - x_5) \bullet f$
$(x_0 - x_7) \bullet f$
$w_1$
**(3 adders delay)**
$w_7$
**(4 adders delay)**

$(x_0 - x_7) \bullet a$
$(x_2 - x_5) \bullet a$
$(x_0 - x_7) \bullet e$
$(x_3 - x_4) \bullet e$
$(x_0 - x_7) \bullet f$
$(x_2 - x_5) \bullet f$
$(x_3 - x_4) \bullet e$
$(x_3 - x_4) \bullet a$
$(x_1 - x_6) \bullet a$
$(x_3 - x_4) \bullet f$
$(x_2 - x_5) \bullet f$
$w_3$
**(3 adders delay)**
$w_5$
**(4 adders delay)**

56

# Proposed DCT under Vdd scaling

**Proposed Design with high/low delay paths**

$W_0$
$W_1$
$W_2$
$W_3$
$W_4$     **Delay=D1**
$W_5$     **@ Vdd1**
$W_6$
$W_7$

**Important Computations**

**Longer Delays**

**Scaled Vdd: Longer paths under Vdd scaling**

$W_0$
$W_1$
$W_2$     **D1**
$W_3$     **@Vdd2**
$W_4$
$W_5$     **D2 >D1**
$W_6$     **@Vdd2**
$W_7$

**Paths Not Computed**

**Extreme Scaled Vdd: Shorter paths affected**

**Only DC component**

$W_0$ **D1 @Vdd3**
$W_1$
$W_2$     **D3 > D1**
$W_3$     **@Vdd3**
$W_4$
$W_5$     **D4 >D1**
$W_6$     **@Vdd3**
$W_7$

**Paths Not Computed**

$Vdd3 < Vdd2 < Vdd1(nominal)$

# 1D-DCT Path Delay Comparisons

58

# Effect of Vdd Scaling

**Different Architectures at Nominal Voltage**

| 1.0V | CSHM DCT (2 alphabets) | DCT with WTM | Proposed DCT |
|---|---|---|---|
| Power (mW) | 25.1 | 29.8 | 26 |
| Delay (ns) | 3.2 | 3.64 | 3.57 |
| Area (um²) | 80490 | 108738 | 90337 |
| PSNR (dB) | 21.97 | 33.23 | 33.22 |

**Proposed Architecture at Reduced Voltage**

| | Proposed DCT Vdd=0.9V | Proposed DCT Vdd=0.8V |
|---|---|---|
| Power (mW) | 17.53(41.2%) | 11.09(62.8%) |
| PSNR (dB) | 29 | 23.41 |

| | Convention al WTM DCT | CSHM DCT (2 alphabet) | Proposed DCT |
|---|---|---|---|
| 1.0 V |  |  |  |
| 0.9 V | FAILS | FAILS |  |
| 0.8 V | FAILS | FAILS |  |

- **Graceful degradation of proposed DCT architecture under Vdd scaling ( Vdd can be scaled to 0.75V)**

- **Conventional architectures fails**

59

# Temporal Degradation: NBTI

Kang, Roy, et. al. – TCAD, DAC-07

# Temporal Reliability Issues in CMOS Technology

- *HCI –* **Hot Carrier Injection**
- *NBTI –* **Negative Bias Temperature Instability**
  - ➤ **Increase in $V_T$ of PMOS with time**
  - ➤ **The dominant reliability factors in scaled tech.**
- *TDDB, etc.*

# NBTI: Negative Bias Temperature Instability



**Interface trap generation due to Si-H bond breaking**

- **Interface trap ($N_{IT}$) generation at the channel interface due to the Si-H bond breaking, when negative gate bias is applied**
- **With time, $V_T$ increases, subthreshold slope (S) increases, mobility degrades,**
- **Drive current ($I_{DS}$) reduces and affect the PMOS speed**
- **Overall reduces the lifetime of PMOS**

# NBTI: Experimental Data



- **PMOS $V_T$ degrades as a power of time due to NBTI**
- **Fixed exponent of 1/6 matches the simulation data\***

* V. Huard and M. Denais, IRPS 2004

# Power-law $V_T$ degradation Model



**Reaction rate:**
$$\frac{dN_{IT}}{dt} = k_F[N_0 - N_{IT}] - k_R N_{IT} N_H^{(0)} \approx 0$$

$$\frac{k_F}{k_R} \cdot [N_0 - N_{IT}] = N_{IT} \cdot N_H^{(0)}$$

**Conservation of hydrogen:**
$$N_{IT}(t) = \int_0^{\sqrt{D_H \cdot t}} N_H^{(0)}(y,t) \cdot dy$$

$$N_{IT} = \frac{1}{2} N_H^{(0)} \cdot \sqrt{D_H t}$$

$$\Delta V_T = \frac{q \cdot \Delta N_{IT}}{C_{OX}} \qquad \Longleftarrow \qquad N_{IT}(t) = \sqrt{\frac{k_F N_0}{2 k_R}} (D_H t)^{1/4}$$

## NBTI degrades in time of exponent 1/6

# Mobility degradation factor

- **Mobility degradation due to NBTI is expressed in an additional $V_T$ shift, noted as *m***

- **Overall temporal $V_T$ shift model is expressed as,**

$$\Delta V_T = (1+m)\frac{q\chi\sqrt{E_{ox}\,e^{\left(\frac{E_{ox}}{E_0}\right)}}\cdot t^{0.25}}{C_{OX}}$$

# Impact of NBTI on circuit performance

# Circuit Performance Degradation



- **Performance (delay) degradation also follows the power trend with same 0.17 exponent**
- **In CMOS logic, only the rising (L2H) delay's are affected**

# Circuit Performance degradation cont.



- **Delay degradation in ISCAS c432**
- **Activity factor (switching activity) does not affect much on the delay degradation**
  - ➢ **In reality, activity factor's are balanced in the normal operations**

# Design method considering the NBTI degradation

# NBTI-aware design method



- **Over-design is required to guarantee a lifetime stability of the circuit**
- **LR sizing is used to optimize the circuit**
  - Size the circuit considering the worst-case $V_T$ degradation over the lifetime

# LR Sizing considering NBTI

```
┌─────────────────────────────┐
│ 1. Delay Constraint (D_MAX) │ ◀──  Lifetime Constraint
│ 2. Required Lifetime (T_Life)│      A new design constraint
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│ Calibrate switching activity's│
│ (S_i) in each node          │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│ Compute V_T shift in each node│ ◀──  Power-law V_T model
│ considering S_i             │
└─────────────────────────────┘
              │
              ▼
┌─────────────────────────────┐
│ LR sizing with delay constraint│ ◀──  Optimal sizing from
│ D_MAX                       │      Lagrangian Relaxation*
└─────────────────────────────┘
              │
              ▼
      ⬡ NBTI-aware Design ⬡
              │
              ▼
```

**Guarantee a lifetime stability under NBTI degradations**

* C. Chen et. al., TCAD 1999

# Simulation results

**1. Delay degradation in ISCAS85 benchmark circuits after 10 years**

| Circuit | No. of Trans. | Nominal delay (ps) | % delay degrad. (10 yrs) | |
|---|---|---|---|---|
| | | | $S_i = 1$ | $S_i < 1$ |
| c432 | 590 | 525 | 8.90 | 7.32 |
| c499 | 1816 | 368 | 9.20 | 8.06 |
| c1908 | 1582 | 513.5 | 9.18 | 8.53 |
| c3540 | 3638 | 597.3 | 9.00 | 7.86 |
| c74181 | 372 | 194.6 | 9.89 | 8.68 |
| c74182 | 92 | 77.2 | 10.35 | 9.63 |
| c74283 | 188 | 131.9 | 7.90 | 6.83 |
| c74L85 | 148 | 115.1 | 9.50 | 7.60 |

**\* All benchmarks are synthesized in BPTM 70nm technology**

# Simulation results cont.

**2. Area overhead in NBTI-aware sizing**

| Circuit | Nominal delay (ps) | Nominal area (um) | % Area overhead | |
|---|---|---|---|---|
| | | | $S_i = 1$ | $S_i < 1$ |
| c432 | 385 | 196.7 | 14.8 | 13.6 |
| c499 | 340 | 581.47 | 7.82 | 6.71 |
| c1908 | 470 | 489.67 | 7.13 | 6.68 |
| c3540 | 500 | 1146.5 | 3.44 | 3.31 |
| c74181 | 180 | 111.1 | 9.45 | 9.0 |
| c74182 | 80 | 31.1 | 11.3 | 11.2 |
| c74283 | 125 | 66.71 | 10.0 | 10.0 |
| c74L85 | 120 | 42.59 | 5.85 | 5.8 |

**\* All benchmarks are synthesized in BPTM 70nm technology**

# Negative Bias Temperature Instability



V_GS < 0V

H₂ H₂ H₂ H₂ H₂

GATE OXIDE

$V_D = 0V$

$V_S = 0V$

H H H H H H H

X H X H

P+ Si Si Si Si P+ DRAIN

n-sub

$V_{BODY} = 0V$

After NBTI degradation

- PMOS specific *Aging Effect*
- Generation of (+) traps
- Reaction-Diffusion (RD) model*
- Time exponent ~ 1/6

$$N_{IT}(t) = \sqrt{\frac{k_F N_0}{2 k_R}} \left( D_H t \right)^{1/6}$$

$$\Delta V_T = \frac{q \cdot \Delta N_{IT}}{C_{OX}}$$

*M. A. Alam, IEDM'03

# NBTI in Digital Circuits



## *Logic Circuits*
- $f_{MAX}$ **decreases** ↓
- **Timing failure with time**

## *Memory Circuits*
- **Static Noise Margin (SNM)** ↓
- **Read & Write Stability**
- **Parametric Yield** ↓

Temporal $V_{Th}$ increase in PMOS affects critical performance factors of digital VLSI circuits

# NBTI: Random Logic Circuits

**Delay Degrad. STD cells**

| Logic Cell | fanin | Delay (ps) | | Δ (%) |
|---|---|---|---|---|
| | | t=0 | 3 years | |
| INV | 1 | 13.77 | 16.77 | 21.8 |
| NAND | 2 | 16.86 | 19.88 | 17.9 |
| NAND | 3 | 19.57 | 22.45 | 14.8 |
| NOR | 2 | 17.26 | 21.89 | 26.8 |
| NOR | 3 | 23.80 | 30.19 | 26.9 |



Chart legend:
- c2670
- c5315
- c1908
- c499
- c3540
- c74181

8% $f_{MAX}$ decrease

$n \sim 1/6$

PTM 65nm
ISCAS'85
Benchmarks

$f_{MAX}$ reduction (%)

Time (s)

3 years Lifetime

- ISCAS'85 Benchmark Circuits, PTM 65nm
- Gate delay: analytical delay model considering NBTI
- Circuit delay: NBTI-aware Static Timing Analysis (STA)
- Circuit $f_{MAX}$ → time exponent n ~ 1/6

# NBTI: 6T SRAM Cell

## Static Noise Margin (SNM)



PTM 60nm, 125°C
CR = 1.33
PR = 0.67

*1/6 trend line*

Degradation in SNM (%) vs Time (s)

## Distribution in Write Margin



- t=0
- t=$10^8$
- t=$10^7$
- t=$10^6$

*WM improves with time*

CDF vs Write Margin (V)

- ❑ **SNM degrades by more than 10% in 3 years**
- ❑ **% SNM Degradation → time exponent n ~ 1/6**
- ❑ **WM improves with time under NBTI**

# Design for Reliability under NBTI



**c1908**

*10ps*

*Area Saving From TR-based*

*Cell-based*

Legend:
- ▲ INIT
- ● TR-Based
- ■ Cell-Based

Y-axis: **Area (um)** — 550, 600, 650, 700, 750, 800

X-axis: **Delay (ps)** — 580, 600, 620, 640, 660, 680, 700, 720

❑ **Simulation Setup**
- ➢ **Synthesized in PTM 65nm**
- ➢ **1/6 $V_{Th}$ degradation model**
- ➢ **125°C Stress temperature**
- ➢ **50% Signal Probability at PI's**

❑ **Gate Sizing applied to guarantee lifetime functionality of design**

❑ **11.7% overhead for Cell-based sizing**

❑ **6.13% overhead for TR-based sizing**

- ➢ **45% improvement in area overhead**
- ➢ **Runtime complexity for TR-based sizing is identical to that of Cell-based sizing**

# $I_{DDQ}$ based NBTI Characterization

**Microphotograph**

**Layout**

Inverter Chain

$V_{DD}$

$V_{in}$

*1000 stages*

$I_{DDQ}$ **Measurement**

| Technology | CMOS 130nm |
|------------|------------|
| Die Size | 20 ($mm^2$) |
| I/O Pin | 209 |
| $T_{ox}$ | 1.6 (nm) |
| $V_{DD}$ | 1.2 (V) |

- Test Circuit Fabricated
- 1000 stage INV chain
- DC Stress signal @$V_{in}$
- $I_{DDQ}$ measurement @*GND*

# Correlation between $I_{DDQ}$ & $f_{MAX}$



- $D_M$ < 3ms, Temp=125°C, $V_{stress}$=1.7V
- $I_{DDQ}$ degradation → n~1/6 during
- Clear signature of NBTI
- *Correlation between $I_{DDQ}$ and $f_{MAX}$ can be used to predict circuit performance degradation under NBTI*

**% $I_{DDQ}$ decrease**
**% $f_{MAX}$ increase**  ← **n~1/6**

$$R_{leak} = \frac{\Delta I_{leak}(t)}{I_{leak}(0)} \propto t^{1/6} \propto \frac{\Delta f_{MAX}(t)}{f_{MAX}(0)} = R_{freq}$$

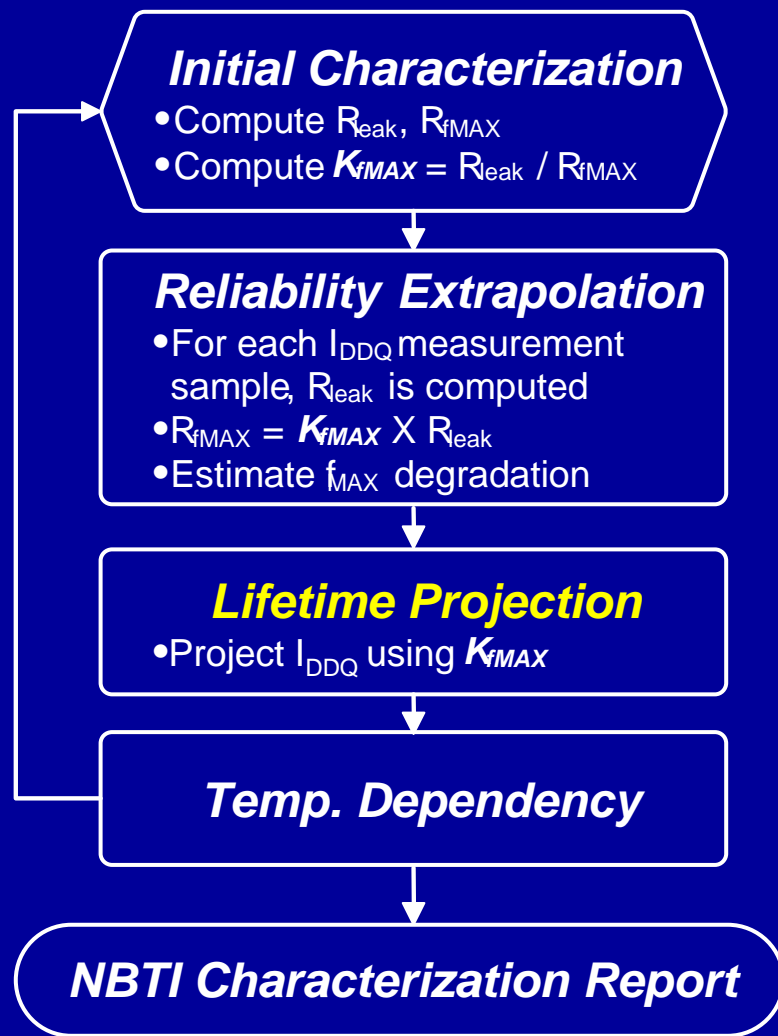$$R_{freq} = K \times R_{leak} \quad (K :\text{constant})$$

# $I_{DDQ}$ based Characterization Technique

**Initial Characterization**
- Compute $R_{eak}$, $R_{fMAX}$
- Compute $K_{fMAX} = R_{eak} / R_{fMAX}$

**Reliability Extrapolation**
- For each $I_{DDQ}$ measurement sample, $R_{eak}$ is computed
- $R_{fMAX} = K_{fMAX} \times R_{eak}$
- Estimate $f_{MAX}$ degradation

**Lifetime Projection**
- Project $I_{DDQ}$ using $K_{fMAX}$

**Temp. Dependency**

**NBTI Characterization Report**

*Design phase*

*Post-silicon phase*

- Circuit-level NBTI Reliability Characterization
- $I_{DDQ}$ test is used
- Expensive $f_{MAX}$ testing is avoided (or minimized)
- Accurate circuit level performance degradation can be predicted
- IC specific burn-in to qualify the target produce
- Efficient way of field monitoring: dynamic local signature of produce usage
- Possible usage in other reliability sources; HCI

# Conclusions

- **Process Variation and Process Tolerance is becoming important**

- **There is a need to optimize designs considering power/performance/yield**