

Sub-threshold Leakage Modeling and Reduction Techniques

James Kao^{1,3}

Siva Narendra^{2,3}

Anantha Chandrakasan³

¹Silicon Labs ²Intel Labs ³Massachusetts Institute of Technology

Outline

Technology scaling and motivation

Sub-threshold leakage modeling

Sub-threshold leakage reduction

Moore's Law on scaling

Cramming more components onto integrated circuits

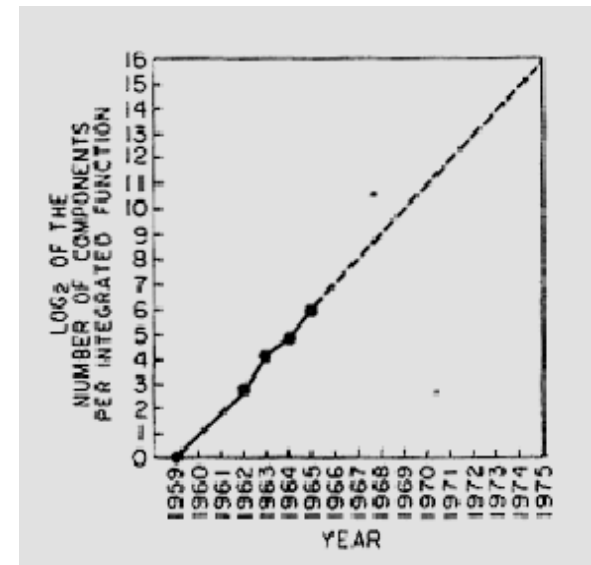
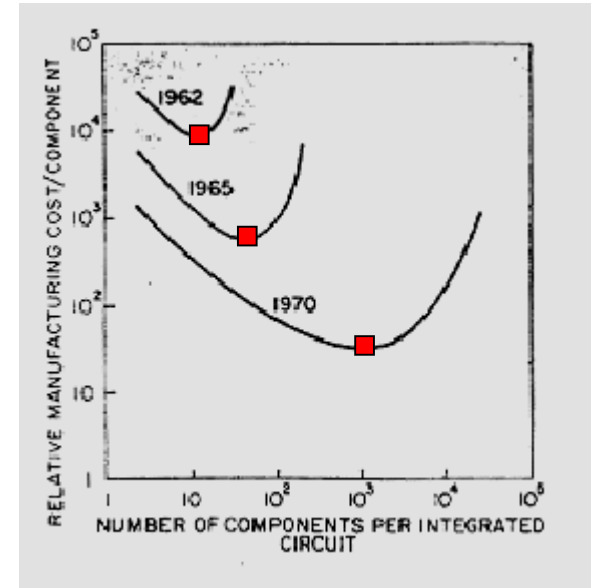
Electronics, Volume 38, Number 8, April 19, 1965

With unit cost falling as the number of components per circuit rises, by 1975 economics may dictate squeezing as many as 65,000 components on a single silicon chip

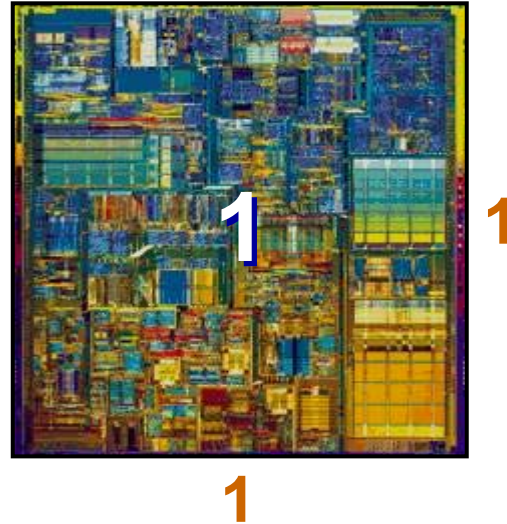
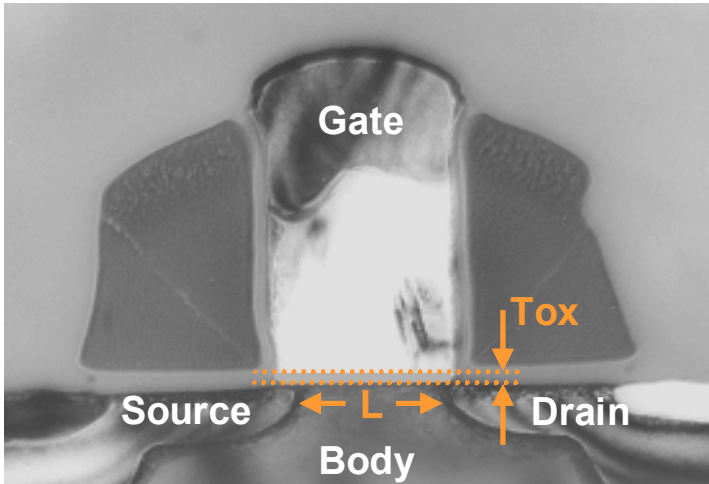
By Gordon E. Moore

Director, Research and Development Laboratories, Fairchild Semiconductor division of Fairchild Camera and Instrument Corp.

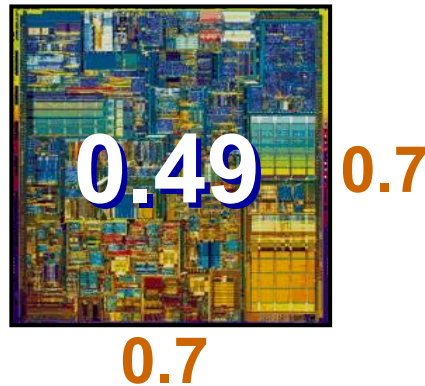
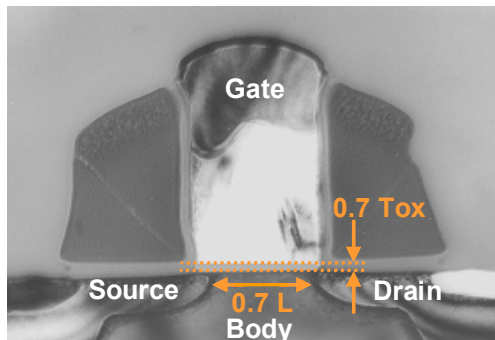
The future of integrated electronics is the future of electronics itself. The advantages of integration will bring about the proliferation of electronics, pushing this science



Scaling of dimensions

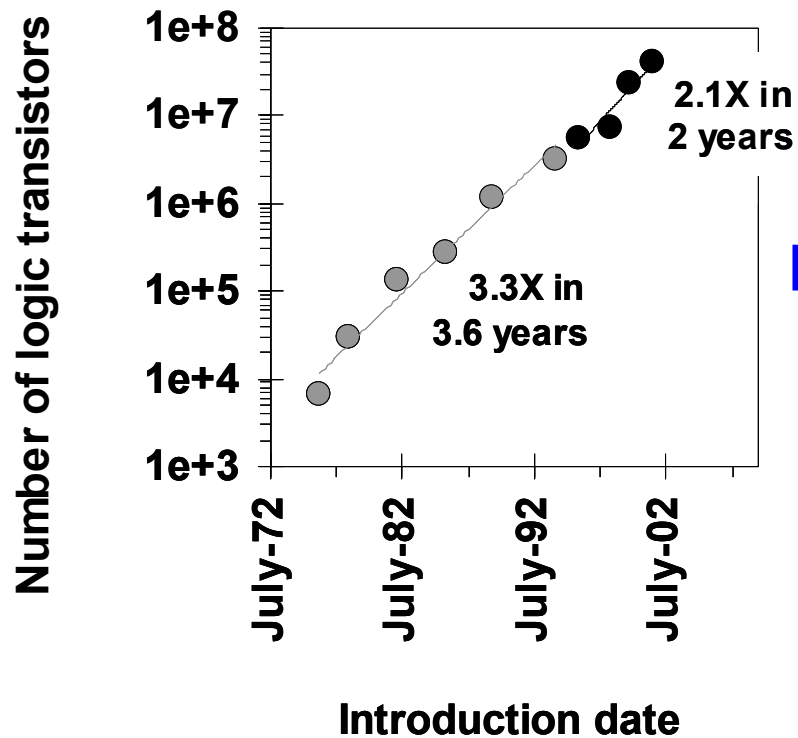
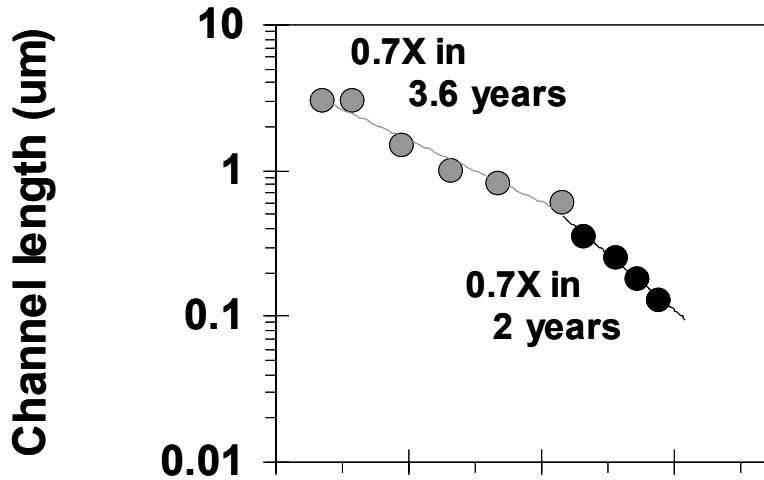


Delay = 1
Freq = 1

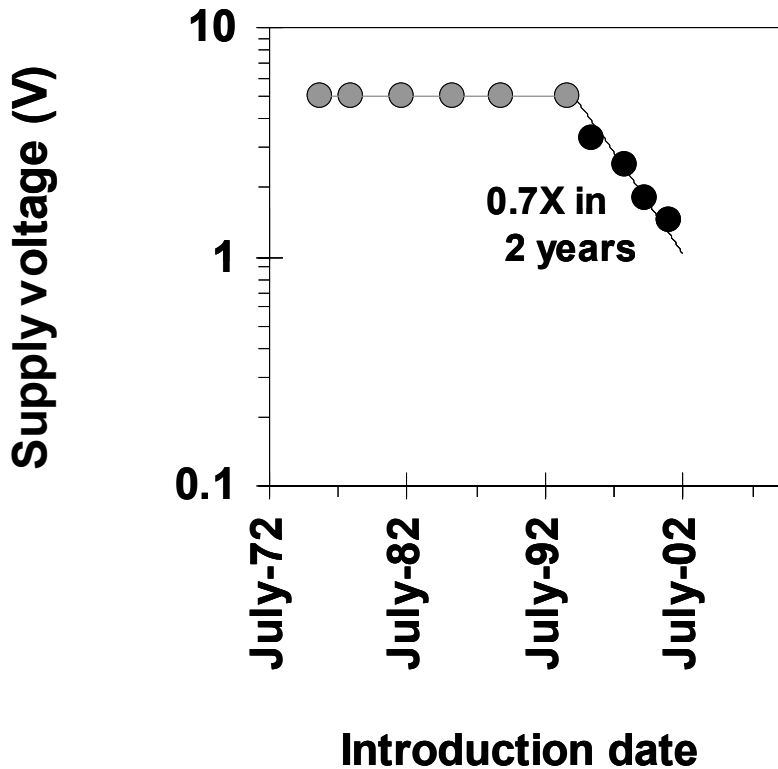
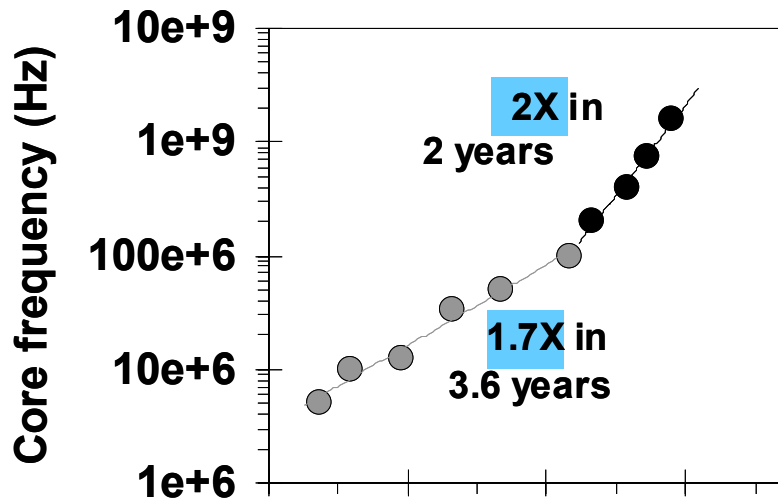


Delay ≈ 0.7

Freq $\approx \frac{1}{0.7} = 1.43$



**Requires die size growth
or same die size**



From early 90s to Present:

$$\text{Power} \propto \text{Area} \times \frac{\epsilon}{t} \times V_{DD}^2 \times F$$

$$\sim 1.0 \times \frac{1}{0.7} \times 1^2 \times 2 \approx 2.9$$

$$\sim 1.0 \times \frac{1}{0.7} \times 0.7^2 \times 2 \approx 1.4$$

$$\text{Delay} \propto \frac{1}{I_{ON}} = \frac{1}{(V_{DD} - V_T)^\alpha}$$

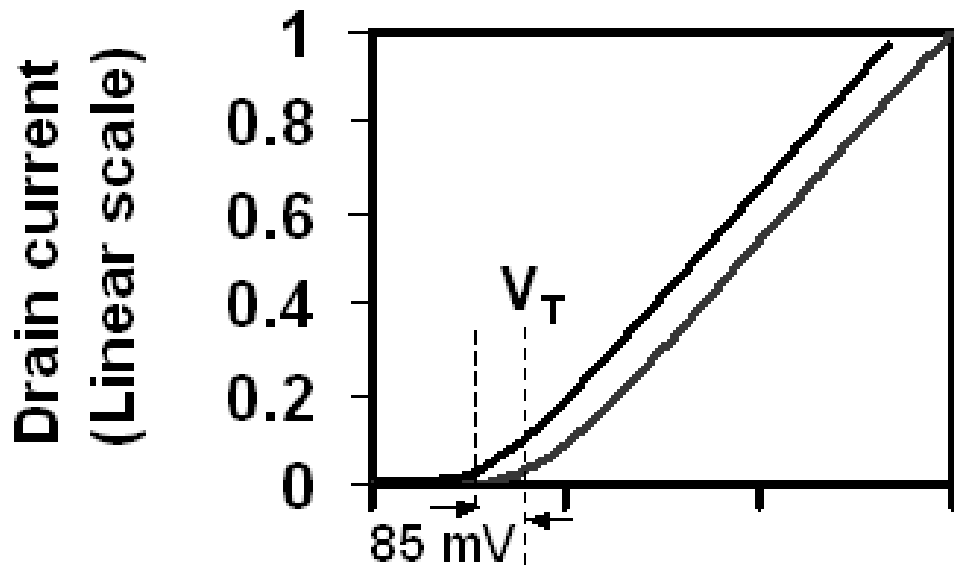
**V_{DD} scaling requires
 V_T scaling**

To continue scaling...

Dimensions including L and T_{ox} have to reduce

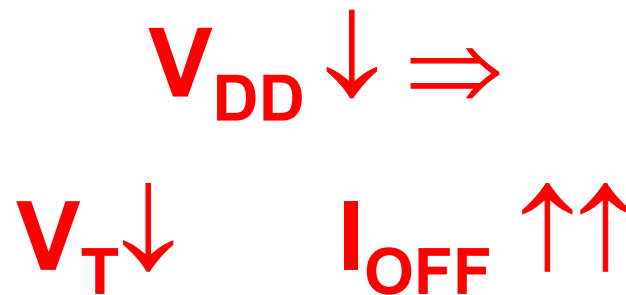
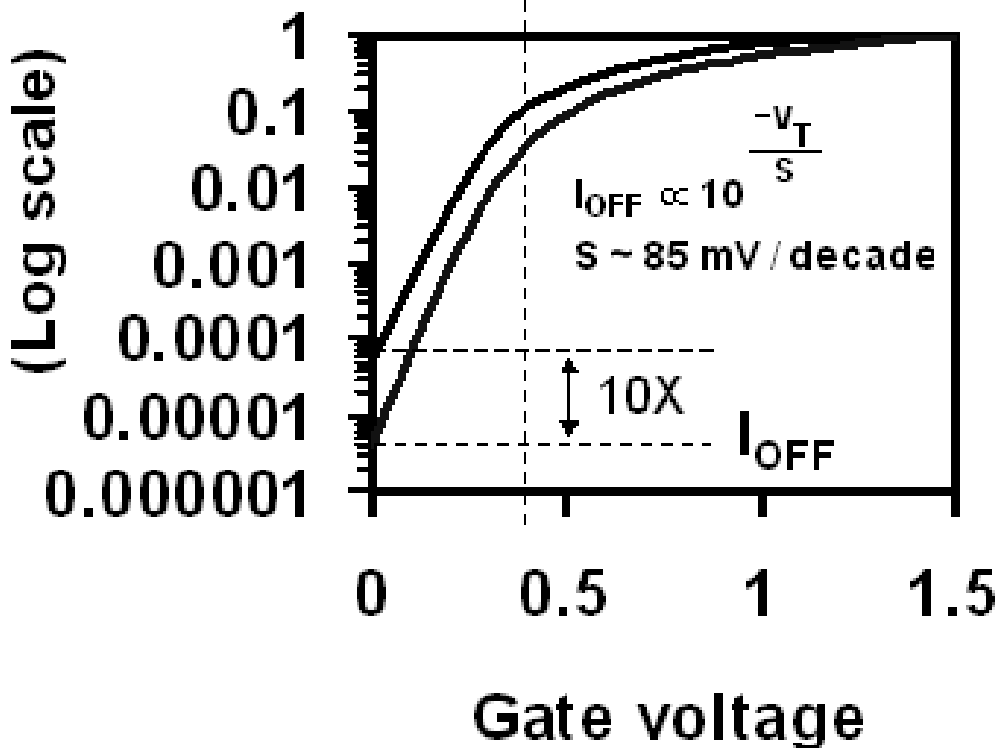
V_{DD} has to reduce to control switching power increase and oxide reliability

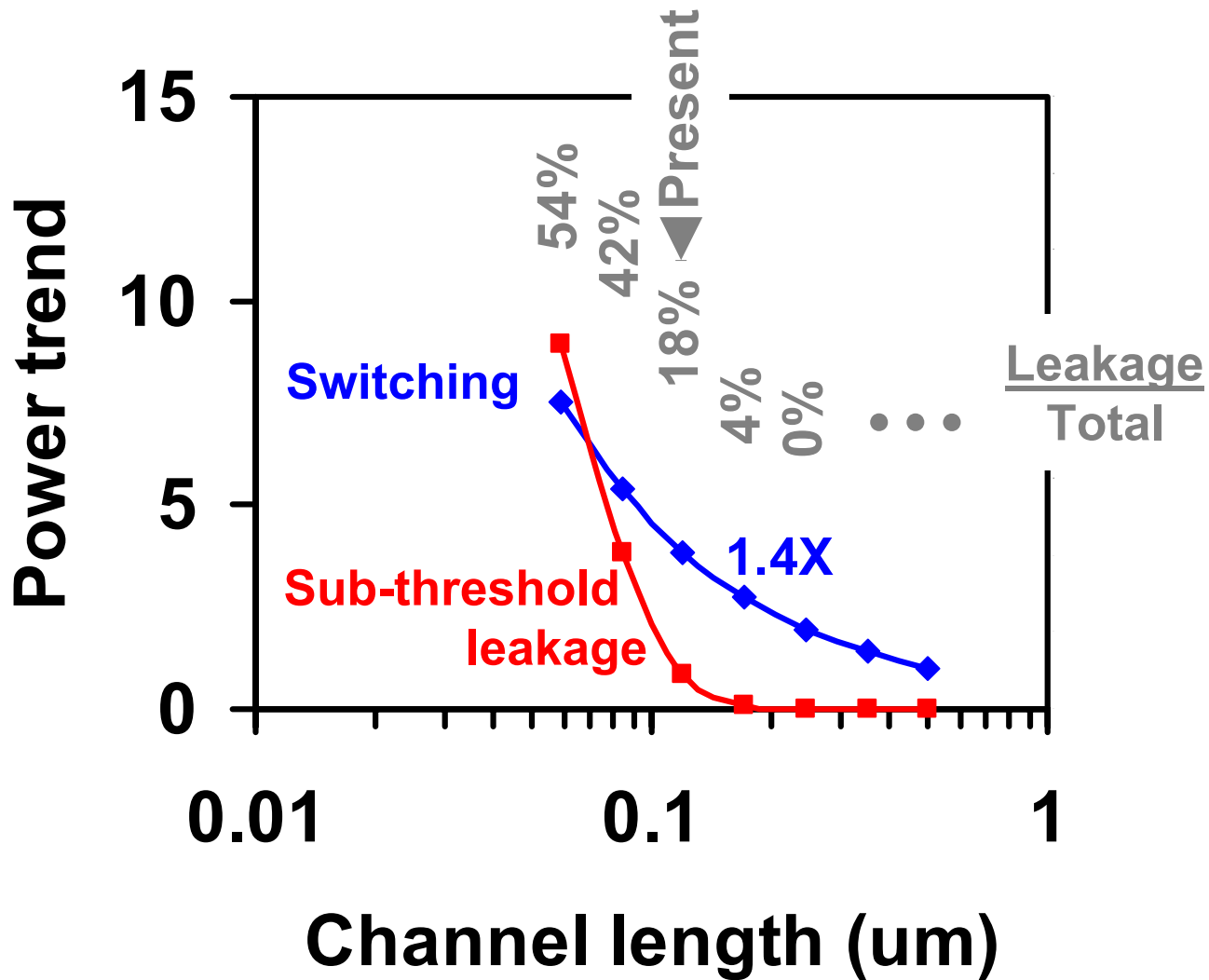
V_T has to reduce to maintain performance increase



$$I_{\text{OFF}} \propto 10^{\frac{-V_T}{S}}$$

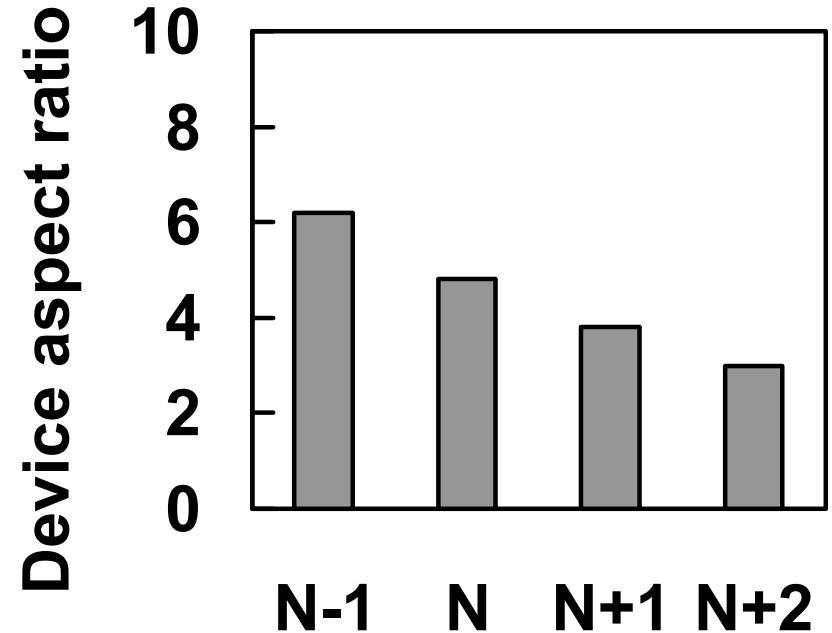
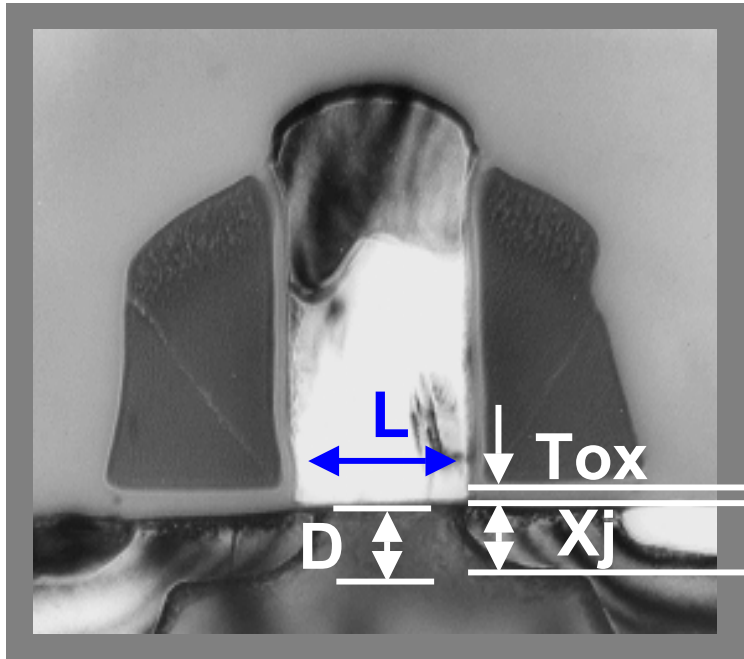
$S \sim 85 \text{ mV / decade}$





Necessary to estimate and reduce sub-threshold leakage power

Device scaling



Device aspect ratio $\approx \frac{L}{\sqrt[3]{T_{ox} \frac{\epsilon_{si}}{\epsilon_{ox}} X_j D}}$

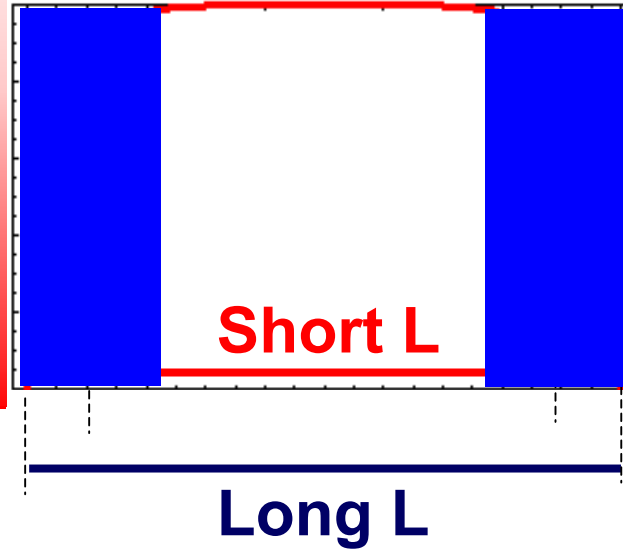
Technology generation

Short channel effects increase with scaling

Barrier Lowering (BL)

Increasing
electron
energy
(NMOS) ↑

Source (n^+)

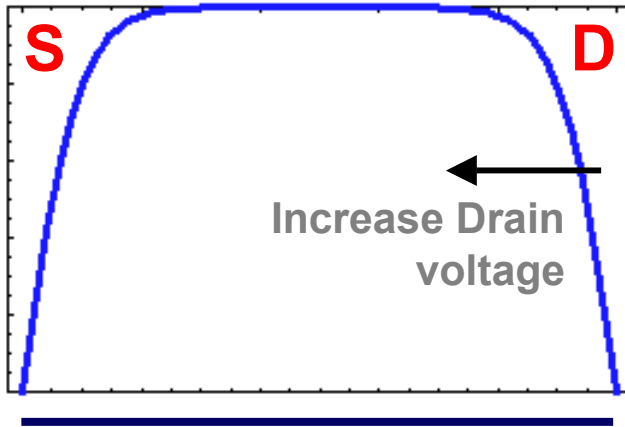


Drain (n^+)

$$L \downarrow \Rightarrow V_T \downarrow \quad I_{OFF} \uparrow \uparrow$$

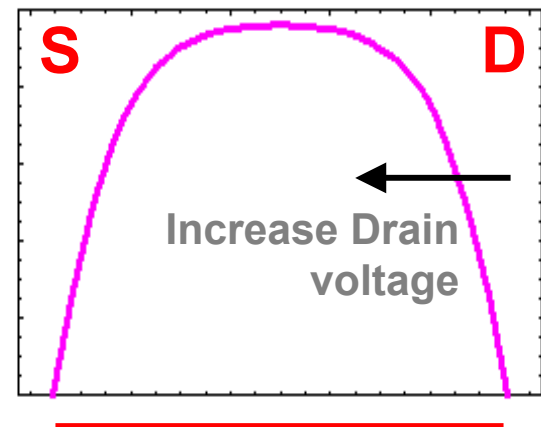
Drain Induced BL (DIBL)

Barrier height



Long L

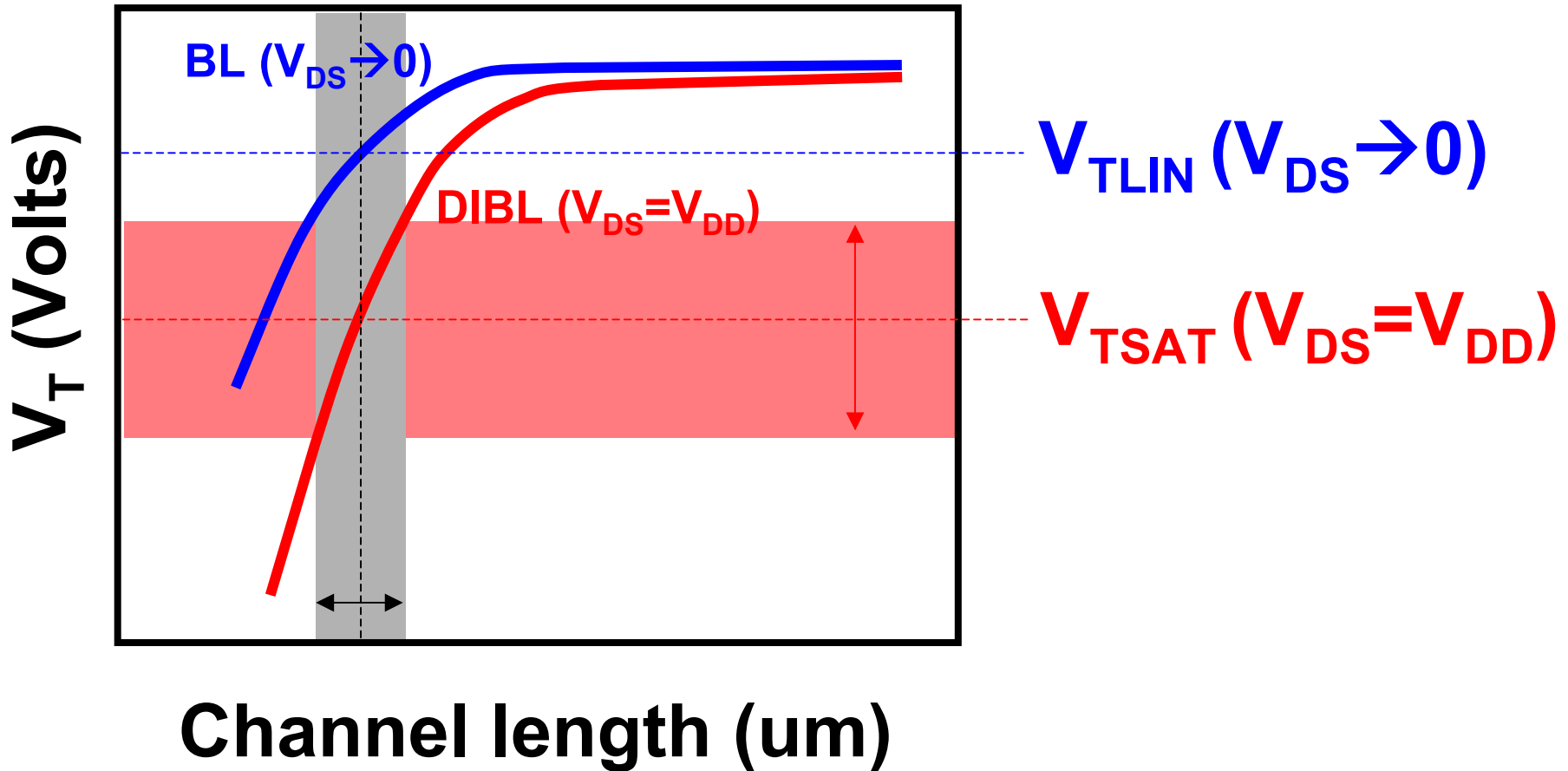
Barrier height



Short L

$$V_{DS} \uparrow \Rightarrow V_T \downarrow \quad I_{OFF} \uparrow \uparrow$$

Impact of variation in L



$$\Delta L \rightarrow \Delta V_T \rightarrow \Delta I_{ON} \text{ \& \ } \Delta I_{OFF}$$

Short Channel MOS V_t

$$V_t = V_{fb} + |2\phi_p| + \frac{\lambda_b}{C_{ox}} \sqrt{2qN\epsilon_s (|2\phi_p| + V_{sb})} - \lambda_d V_{ds}$$

Vt-roll-off factor

H. C. Poon et al., IEDM, pp. 156-159, 1973

$$\lambda_b = 1 - \left(\sqrt{1 + \frac{2W}{X_j}} - 1 \right) \frac{X_j}{L}$$

DIBL equation

K. K. Ng et al., IEEE TED, pp. 1895-1897, Oct. 1993

$$\lambda_d = \left[\frac{L}{2.2\mu m^{-2} (T_{ox} + 0.012\mu m) (W_{sd} + 0.15\mu m) (X_j + 2.9\mu m)} \right]^{-2.7}$$

Sub-threshold leakage

If the current scaling trend continues sub-threshold leakage power expected to be ~50% of the total power.

- Accurate prediction of chip leakage power**
- Techniques to reduce chip leakage power**

Estimation requirements

Require to include within-die variations in

Channel length

Supply voltage

Temperature

| Standby leakage

| Active & burn-in leakage

Leakage modeling

Prior techniques

Lower bound:

Assumes all devices in the die are nominal L

$$I_{leak-l} = \frac{W_p}{k_p} I_p^o + \frac{W_n}{k_n} I_n^o$$

Upper bound:

Assumes all devices in the die are minimum L

$$I_{leak-u} = \frac{W_p}{k_p} I_{off-p}^{3\sigma} + \frac{W_n}{k_n} I_{off-n}^{3\sigma}$$

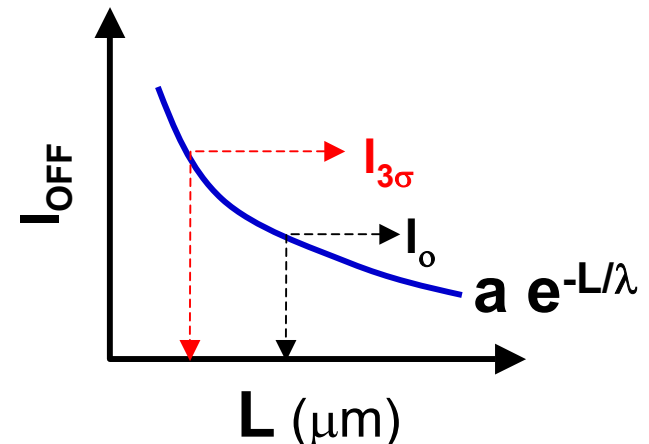
New model

Includes within-die variation

$$I_{leak} = \frac{I^o_w}{k} \frac{1}{\sigma \sqrt{2\pi}} \int_{l_{min}}^{l_{max}} e^{-\frac{(l-\mu)^2}{2\sigma^2}} e^{\frac{(\mu-l)}{\lambda}} dl$$

After simplification using error function properties,

$$\Rightarrow I_{leak} = \frac{I^o_w}{k} e^{\frac{\sigma^2}{2\lambda^2}}$$



Applications...

$$\sigma = \lambda \sqrt{2 \ln \left(\frac{k}{w} \frac{I_{leak}}{I^o} \right)}$$



A macroscopic standard deviation (σ) representing parameter variation in a chip

$$I_{leak-w} = \frac{I_p^o w_p}{k_p} e^{\frac{\sigma_p^2}{2\lambda_p^2}} + \frac{I_n^o w_n}{k_n} e^{\frac{\sigma_n^2}{2\lambda_n^2}}$$

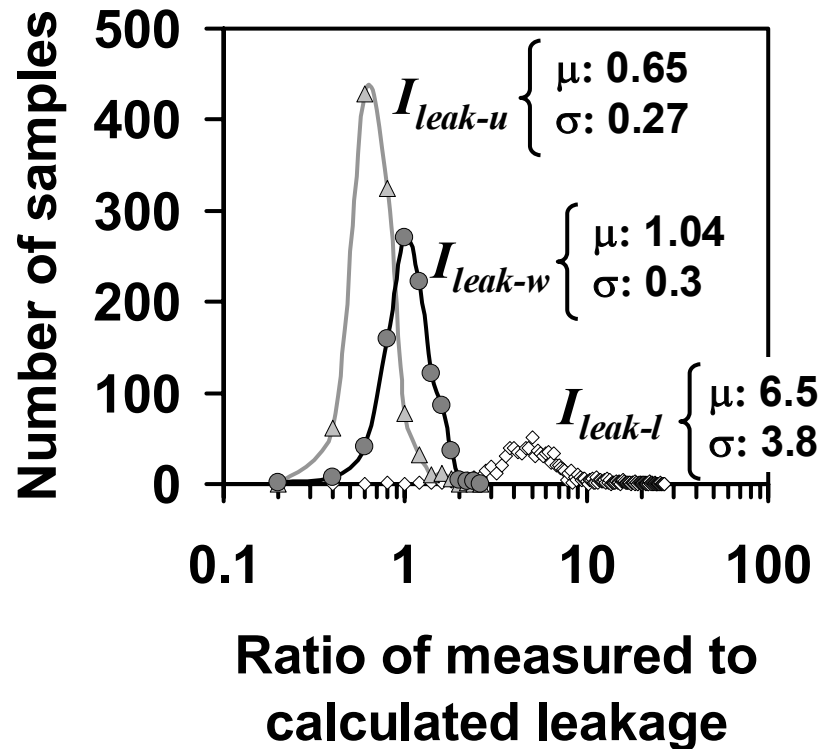


Leakage estimation

Depends on parameters that can be estimated

Measurement results

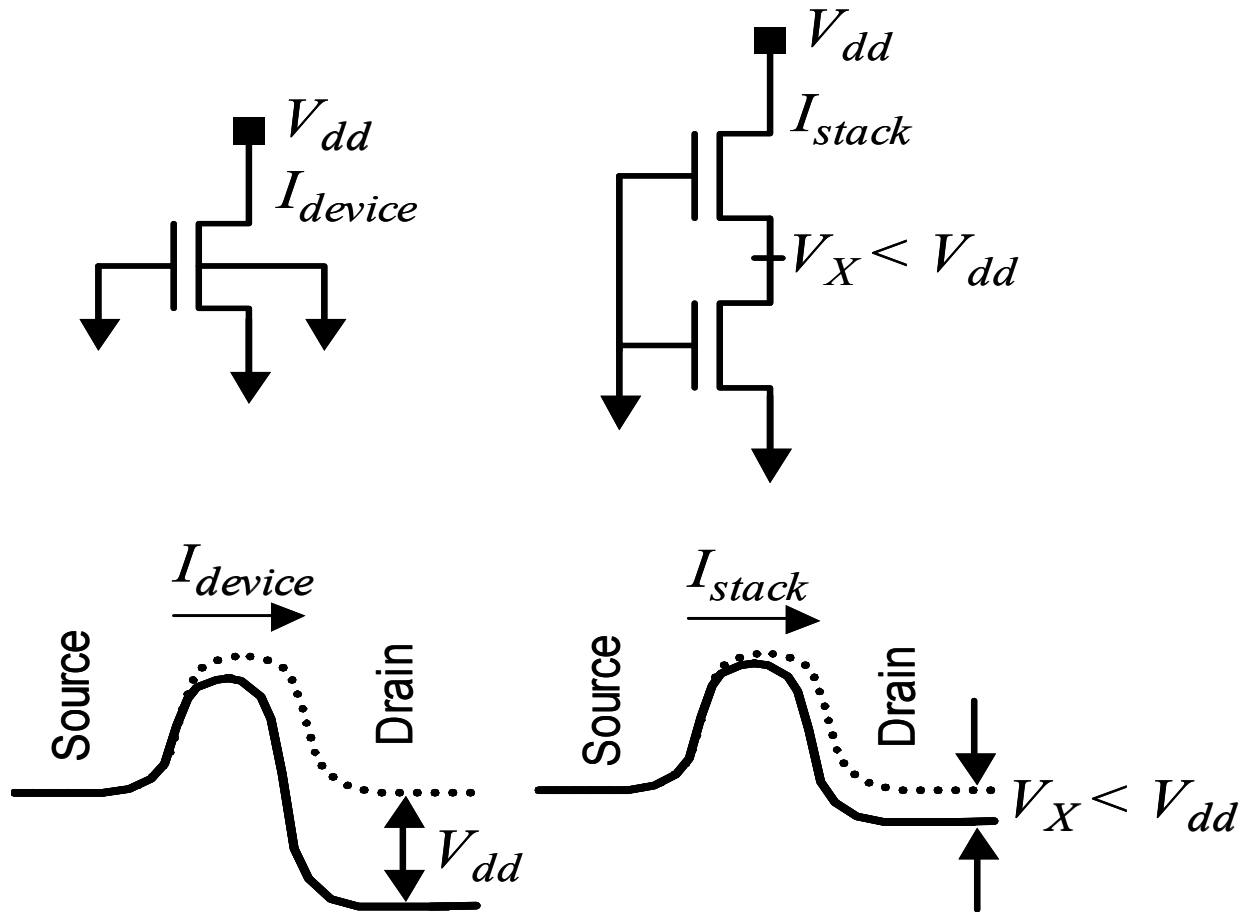
0.18 um 32-bit microprocessors (n=960)



50% of the samples within $\pm 20\%$ of the measured leakage
Compared 11% and 0.2% of the samples using other techniques

Stack effect

For a stack of two devices



Stack effect factor model

$$X = \frac{I_{\text{device}}}{I_{\text{stack}}}$$
$$= 10^{\frac{\lambda_d V_{dd}}{S} (1-\alpha)}$$

where $\alpha \approx \frac{\lambda_d}{1+2\lambda_d}$

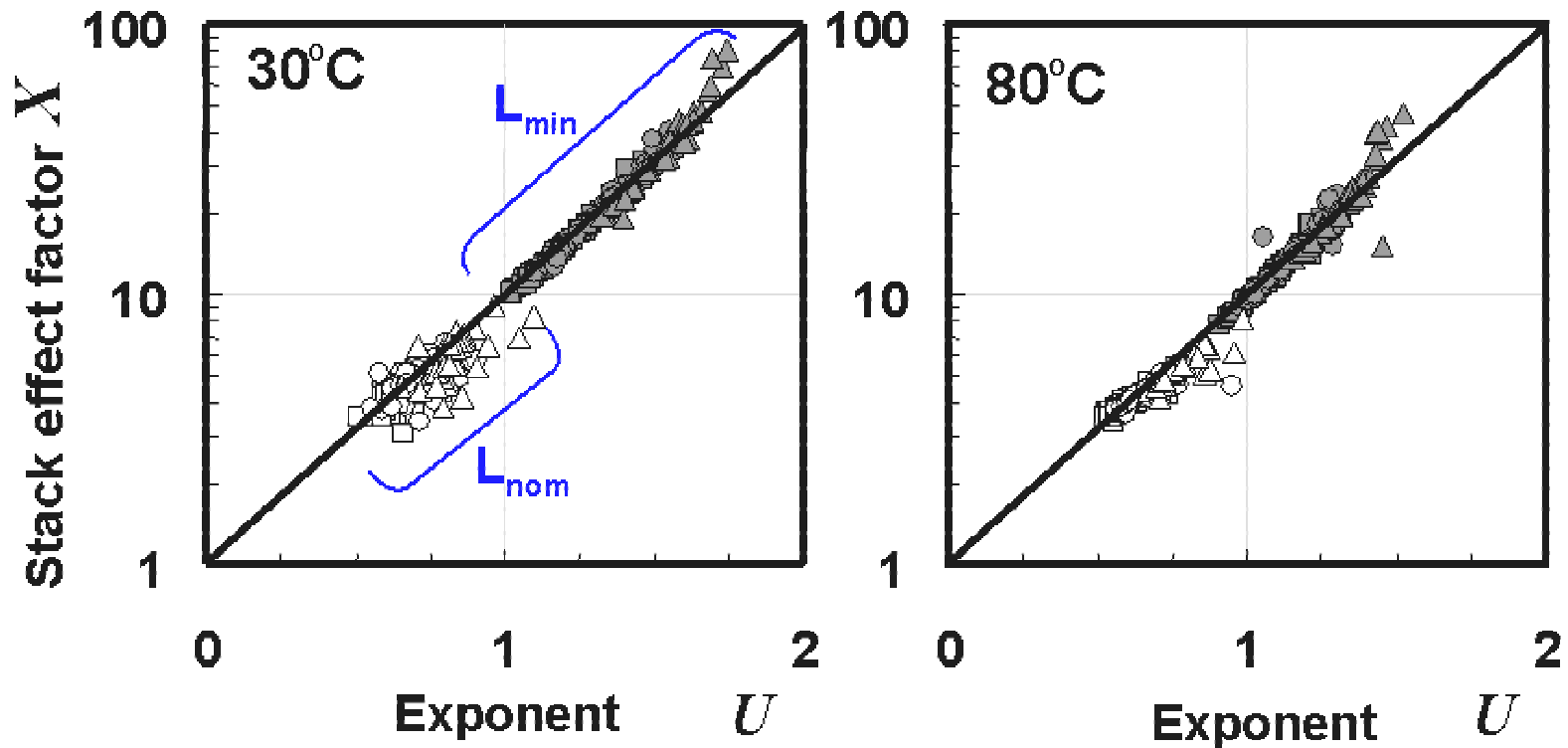
λ_d – DIBL S – Sub-threshold swing

**Stack effect factor can be predicted
based on fundamental device parameters**

Model verification

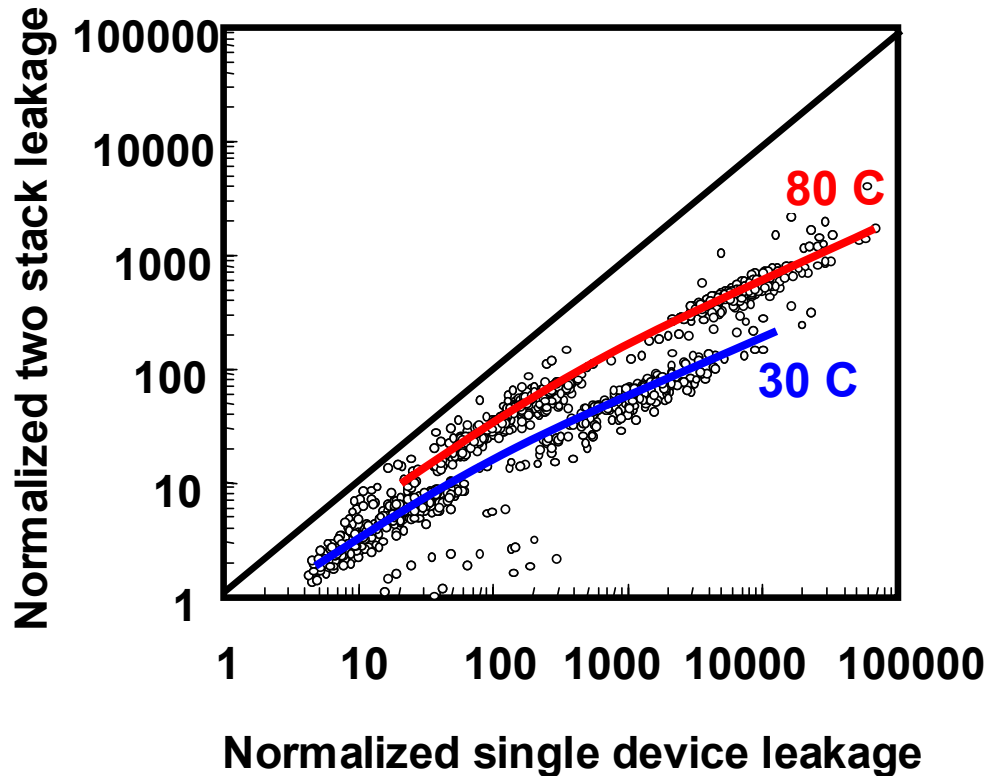
0.18 μm device measurements; zero body bias

$$X = 10^{\frac{\lambda_d V_{dd}}{S} \left(\frac{1 + \lambda_d}{1 + 2\lambda_d} \right)} = 10^U$$



I_{device} vs. I_{stack}

0.18 μm device measurements



**Leakage of a two stack increases
at a slower rate**

Summary

Scaling trend indicate sub-threshold leakage power expected to be ~50% of the total power

Sub-threshold leakage power models should include within-die parameter variation including L, Temperature, Vcc, and device connectivity

Reduction in impact of parameter variation and sub-threshold leakage are essential for scaling

Part II

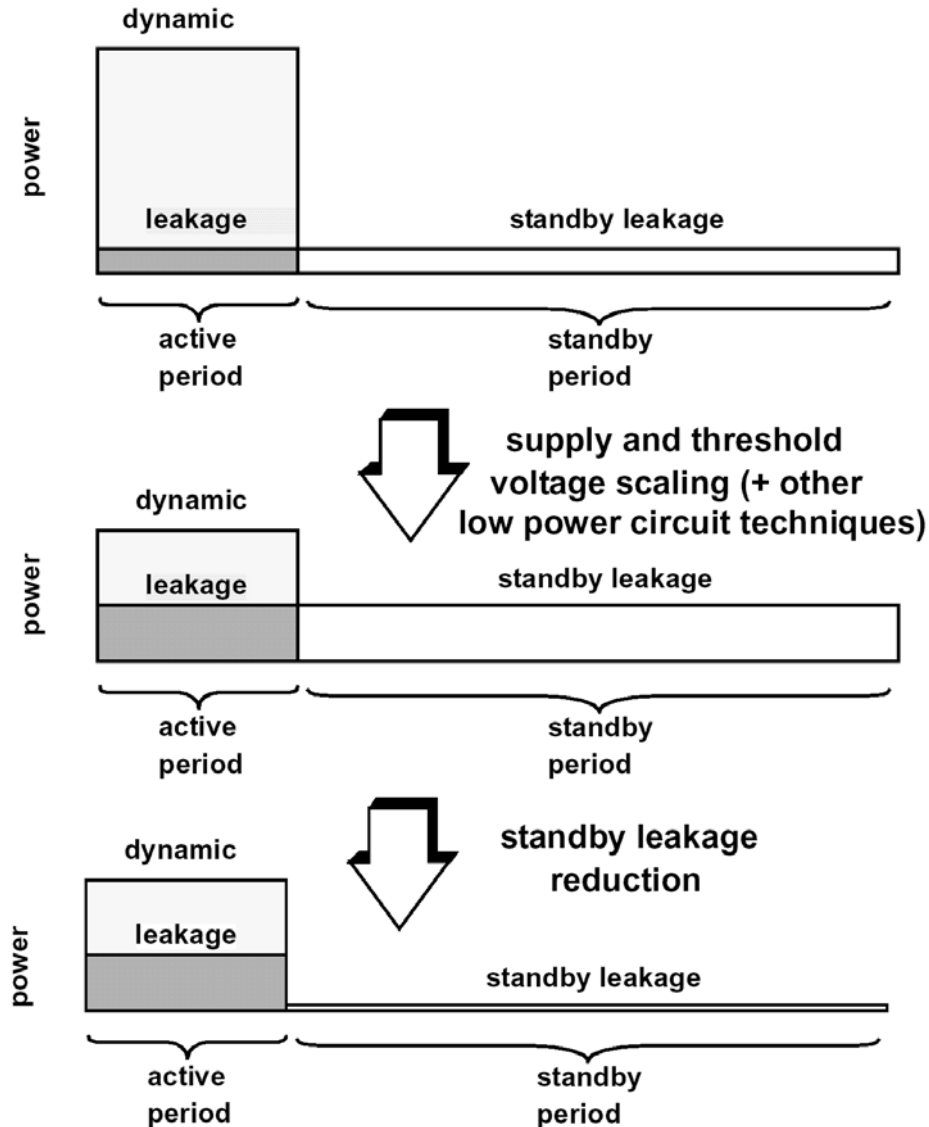
Leakage Reduction Techniques

Leakage Reduction Mechanisms

$$I_{leakage} = I_0 \exp \left(\frac{V_{gs} - V_{t0} - \gamma V_s + \eta V_{ds}}{nV_{th}} \right) * \left(1 - \exp \left(\frac{-V_{ds}}{V_{th}} \right) \right)$$

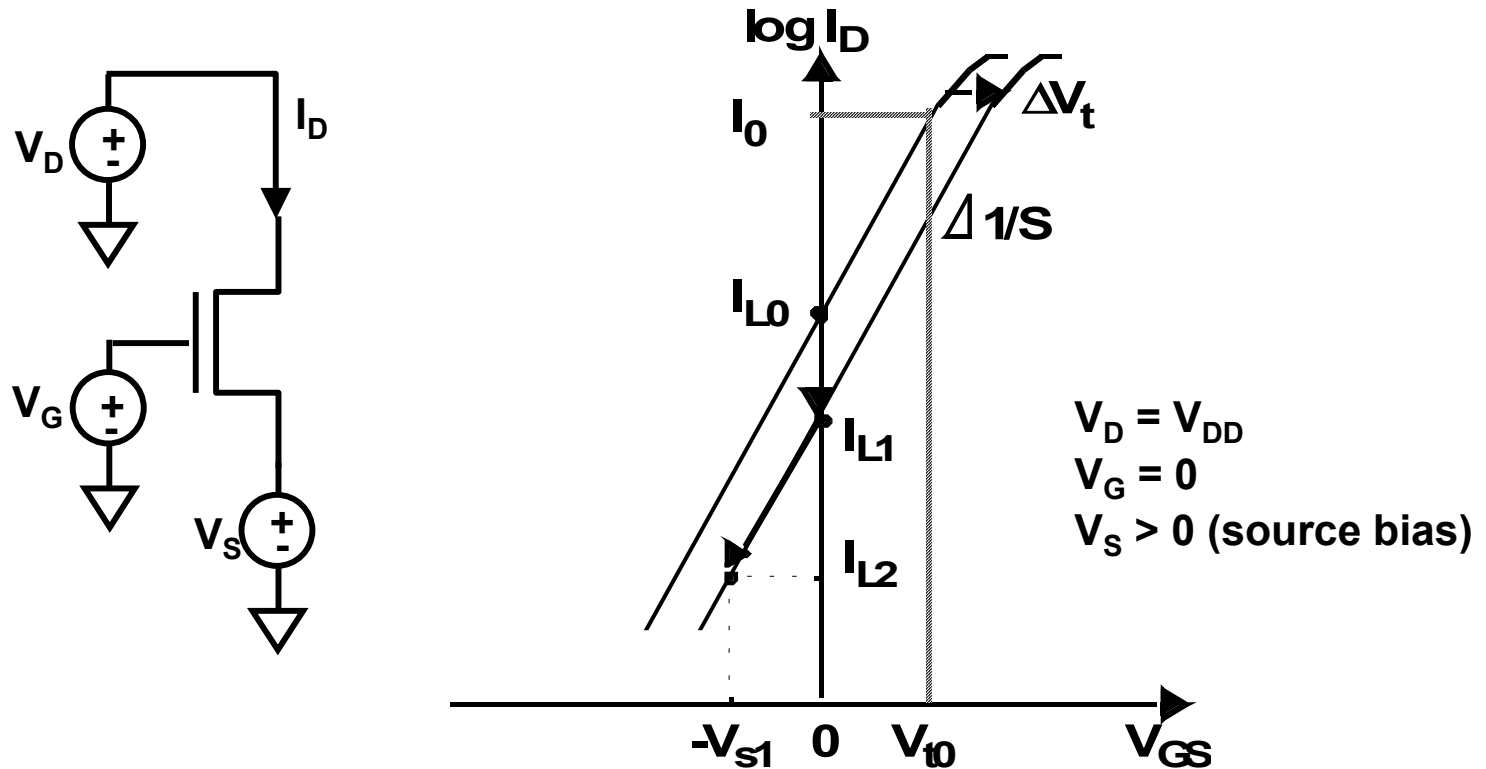
- **Increase V_t**
 - dual threshold Voltage/ MTCMOS/ VTCMOS
- **Increase V_s**
 - source biasing, self reverse biasing, stack effect
- **Decrease V_G**
 - Super cut-off CMOS
- **Decrease V_{DS}**
 - not practical (CMOS output full rail)

Standby and Active Leakage



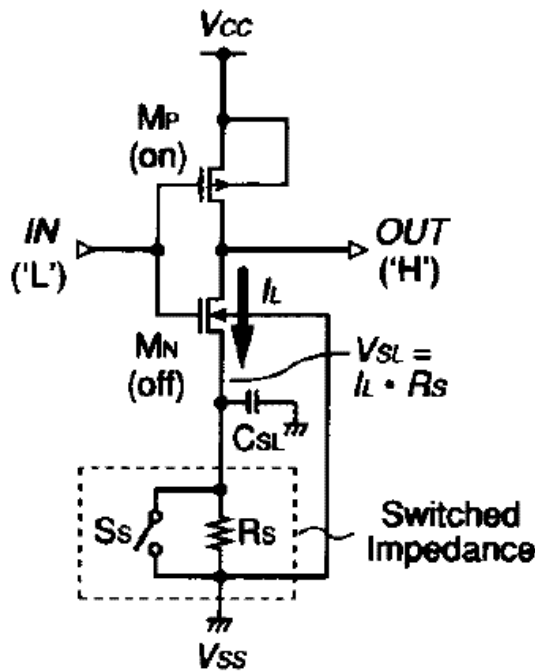
- V_t scaling causes exponential increase in leakage currents
- Dynamic power reduced with supply scaling
- Standby periods can be long (Burst Mode operation- cell phone, pager)
- Standby leakage problem more immediate
- Active leakage control can become important too

Source Biasing Principle

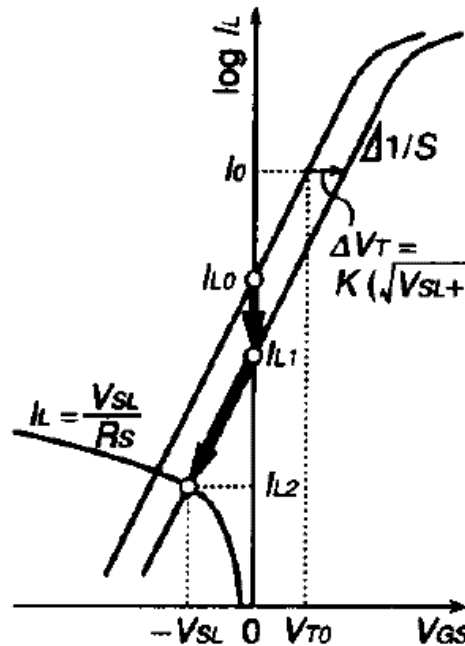


- V_t shift due to body effect γ
- V_{GS} becomes negative
- switched source impedance, self reverse biasing, stack effect

Switched Source Impedance



(a)



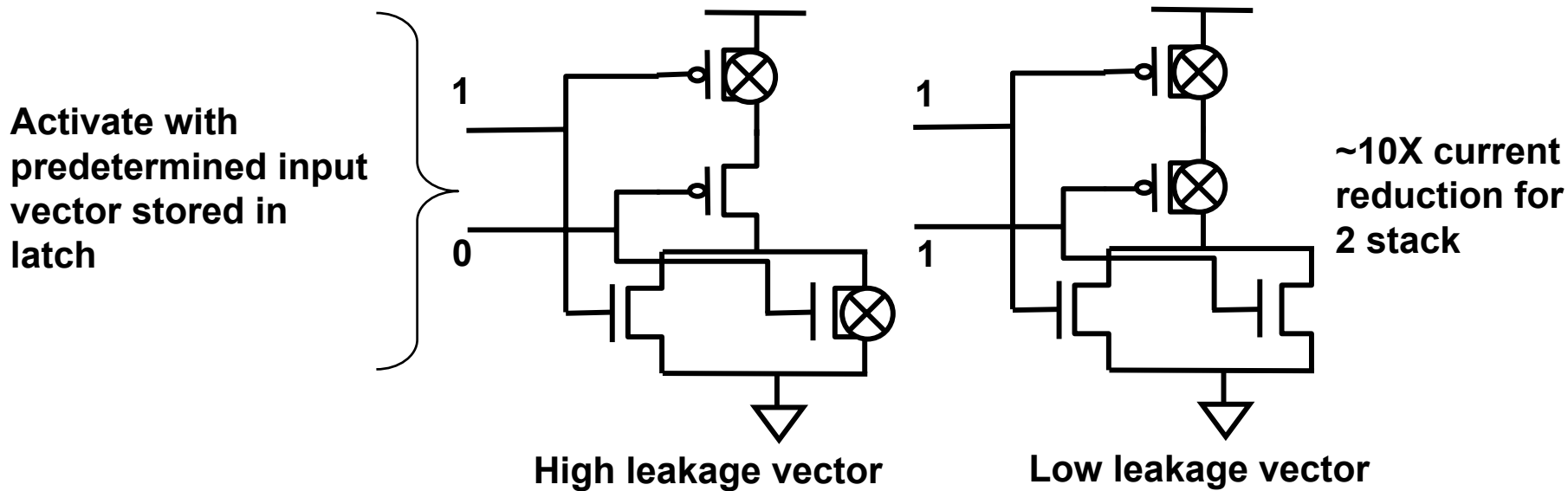
(b)

- R as source impedance
- Estimate >1000X reduction in standby leakage

M. Horiguchi, et al., "Switched-Source-Impedance CMOS Circuit for Low Standby Subthreshold Current Giga-Scale LSI's," JSSC November 1993.

Fig. 1. Principle of switched-source-impedance CMOS circuit: (a) schematic circuit diagram, (b) mechanism of subthreshold-current reduction.

Stack Effect By Vector Activation



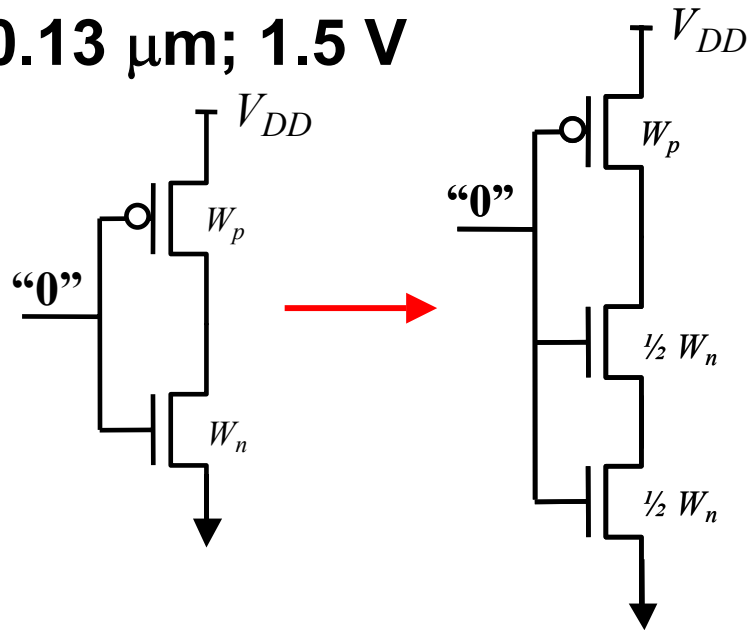
Eg. 32-bit static CMOS Kogg-Stone adder ~2X reduction in total leakage current

- limited by number of stacks available
- proper choice of activating vector (NP-hard algorithm -> use of heuristics)
- internal node settling time can be long
- single stacks are still HIGH leakage

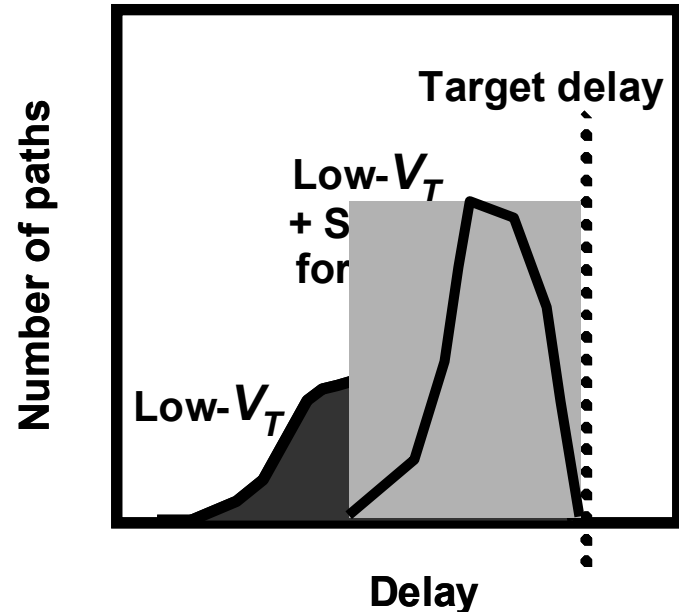
Y. Ye, "A New Technique for Standby Leakage Reduction in High-Performance Circuits" VLSI Symposium, 1998.
M. Johnson, "Models and algorithms for bounds on leakage in CMOS circuits," IEEE TCAD, June 1999.

Stack Forcing Principle

0.13 μm ; 1.5 V



10-30X leakage reduction
~100% higher delay



Force low- V_T stacks in non-critical paths
to reduce leakage

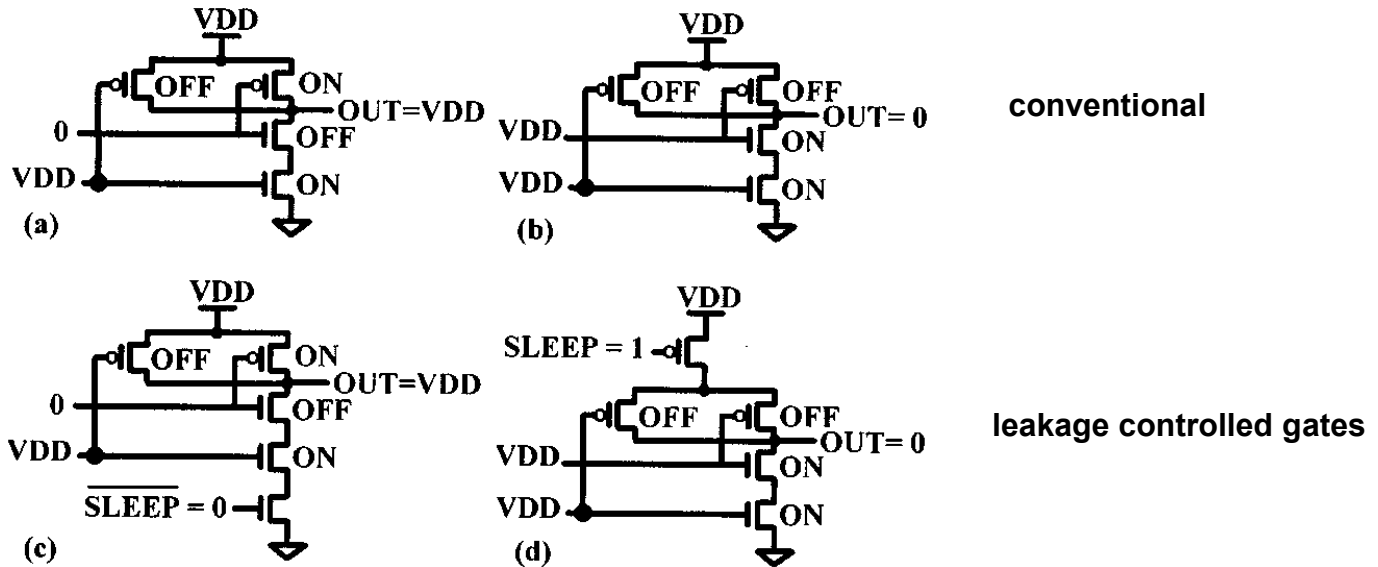
Stack Forcing Effectiveness

32-bit μ P instruction decode block 0.13 μ m; 1.5 V

Frequency of operation:	1.0 GHz		
Active power @ 10% activity:	45.9 mW		
All Low-Vt leakage:	39.1 mW		
Dual-Vt leakage:	9.0 mW	High-Vt usage:	94.2%
Forced stack in low-Vt:	13.2 mW	Forced stack usage:	70.2%

Leakage power reduction
4.3X with dual- V_t , 3X with stack forcing

Leakage Control Stack Devices



- Single V_t leakage reduction mechanism
- Insertion of extra stack devices (in addition to vector activation)
- Sleep devices can be shared among several gates
- Gives further 35% - 90% reduction compared to state dependence alone
- Boils down to single V_t version of MTCMOS (to be discussed)

Dual V_t CMOS

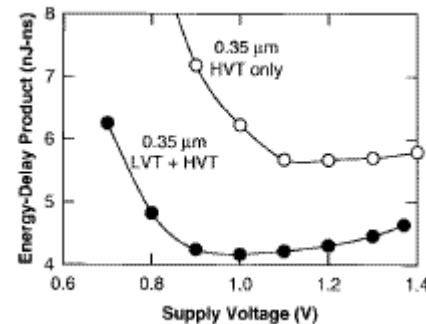
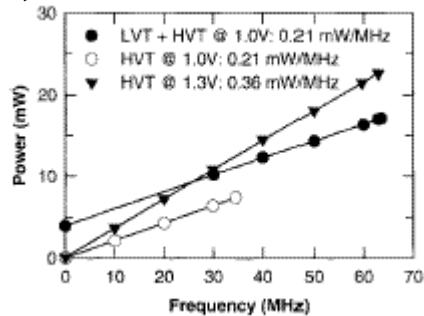
$$I_{leakage} = I_0 \exp \left(\frac{V_{gs} - V_{t0} - \gamma V_s + \eta V_{ds}}{nV_{th}} \right) * \left(1 - \exp \left(\frac{-V_{ds}}{V_{th}} \right) \right)$$

- **Dual V_t more effective at reducing leakage currents than source biasing**
- **Multiple threshold technologies more common**
- **For $S=85$ mV/Decade**
 - each 255mV shift = 3 orders of magnitude reduction
- **Low V_t device= fast, high leakage**
- **High V_t device= slow, low leakage**
- **Achieves both Active and Standby leakage reduction**

Dual V_t Gate Partitioning

A simple approach: Use Low V_t cells for time-critical paths to improve performance

- W. Lee, “A 1V Programmable DSP for Wireless Communications,” JSSC, Nov. 1997



- RN. Rohrer, “A 480 Mhz RISC uProcessor in a .12 μ m Leff CMOS Technology with Copper Interconnects”, JSSC Nov. 1998.
Use of LVT in 4% of standard cells yield 6.5% performance improvement
- T. Yamashita, “A 450 Mhz 64b RISC Processor Using Multiple threshold Voltage CMOS,” ISSSC Feb 2000
LVT + HVT improves performance by 12.5% (all LVT causes standby current to be so large as to cause thermal runaway)

Dual V_t Optimization

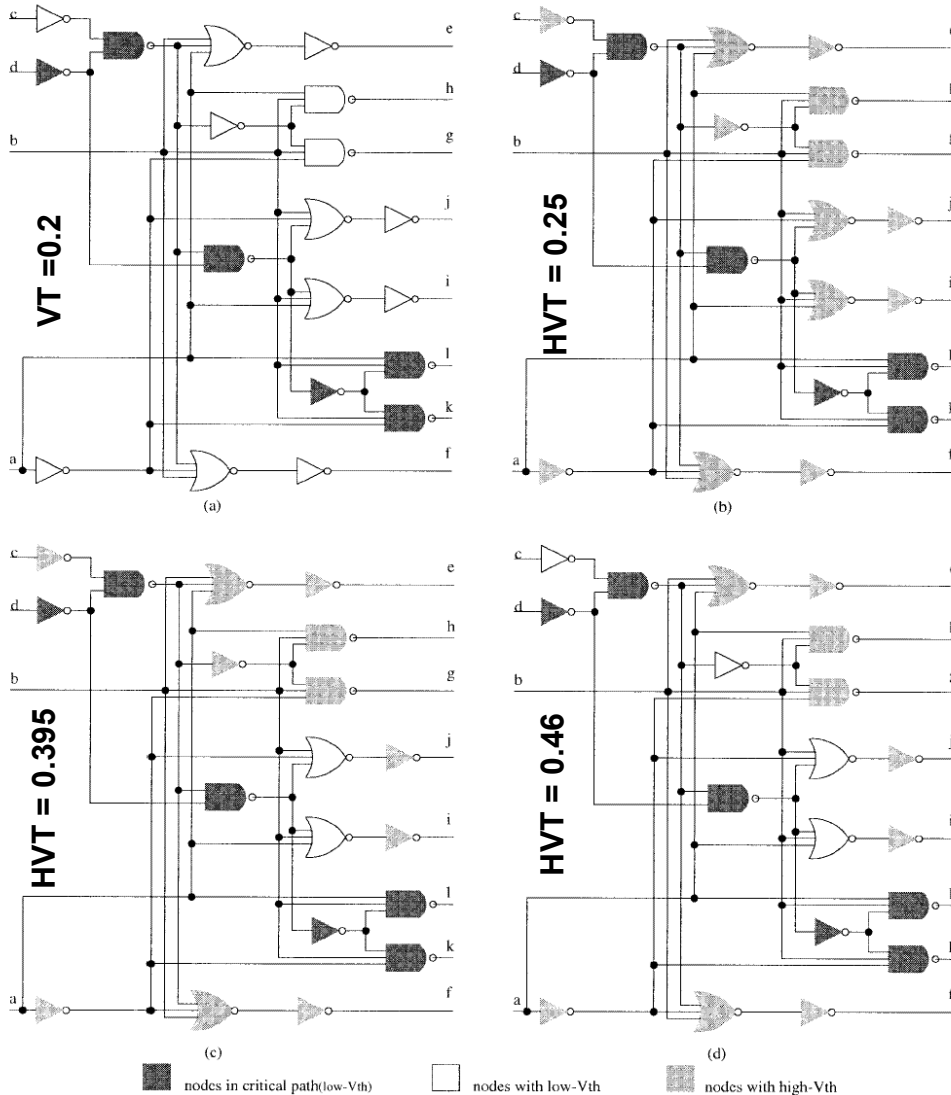


Fig. 4. An example circuit. (a) Original circuit $V_{dd} = 1$ V, $V_{th1} = 0.2$ V. (b) $V_{th2} = 0.25$ V. (c) $V_{th2} = 0.395$ V. (d) $V_{th2} = 0.46$ V.

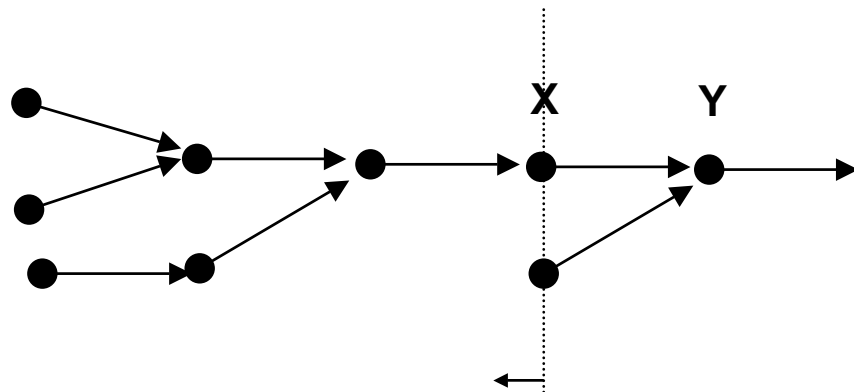
- Initially assume all LVT (for best performance)
- Some non-critical gates can be made HVT
- Choice of high V_t determines mixture
- Proposal for breadth first search algorithm to assign optimal high V_t value.

L. Wei, "Design and Optimization of Dual Threshold Circuits for Low-Voltage Low-Power Applications," TVLSI, March 1999.

A Dual V_t Partitioning Algorithm

L. Wei, "Design and Optimization of Dual-Threshold Circuits for Low-Voltage Low-Power Applications," TVLSI, March 1999

- Initially all low V_t
- For each node (gate) in graph calculate
 - Arrival time, Departure time, Propagation delay, graph level
- From output to input (back tracing level-by-level)
 - Determine slack availability for each node in level
 - Gates with enough slack set to high V_t
 - level = level -1
- Simulate and reiterate with other V_t choices



$$\Delta X (\text{slack}) = T_{\max\text{arrv}}(Y) - T_{\text{dep}}(X) + \Delta Y$$

If $\Delta X > 0$, X can be made high V_t

Update graph with new T_p , T_{dep} , ΔX
Move to next node/ level

Advanced Dual V_t Optimization

- Gate level dual V_t -> transistor level dual V_t
 - L. Wei, "Mixed- V_{th} (MVT) CMOS Circuit Design Meth. for Low Power Appl," DAC 1999
 - Improved leakage reduction
 - More involved partitioning algorithm (traverse transistors level by level)
- Combine dual V_t with transistor sizing:
 - S. Sirichotiyakul, D. Blaauw, "Stand-by Power Minimization through Simultaneous V_t Selection and Circuit Sizing," DAC 1999
 - High V_t to low V_t with same sizing can be too fast
 - Low V_t increases node capacitance seen by crossing paths
 - Use "Dominant Leakage State" + probability to estimate total leakage
 - Too complex to optimize: use heuristic approach
 - 1) Choose some V_t low for performance 2) resize circuit to win back area 3) repeat

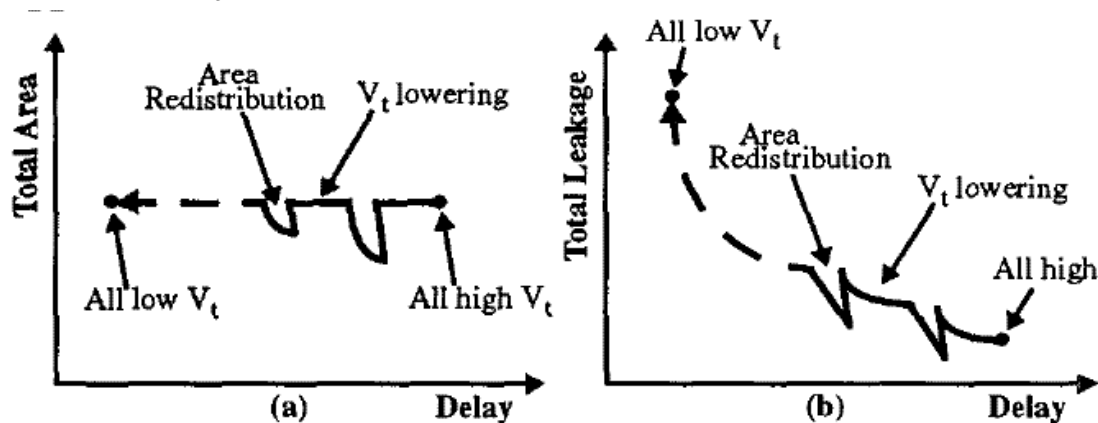
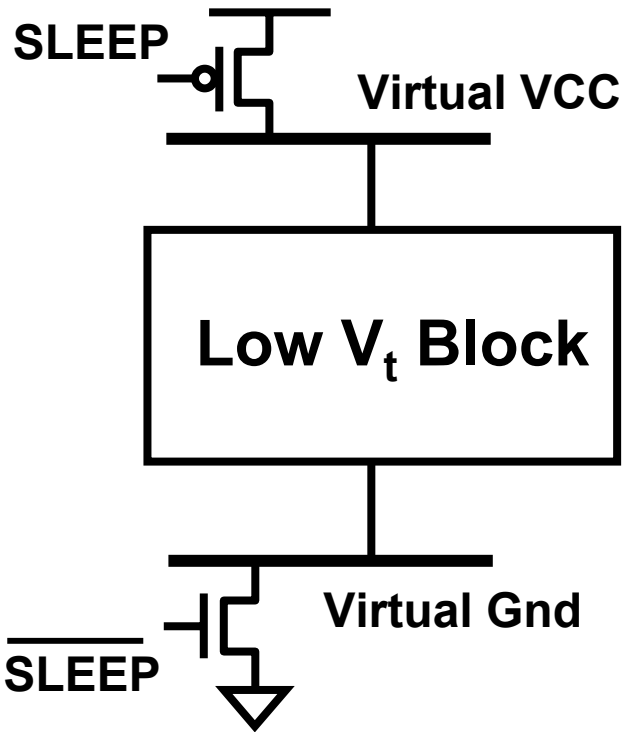


Figure 4. V_t selection and redistribution of area, two views

CAD for Dual V_t Optimization

- Leakage reduction principle simple
- Difficult to optimally choose parameters
 V_{th} , V_{tl} , V_{DD} , device selection, transistor sizing
- Need to develop fast, efficient CAD tools

MTCMOS Principle



Single polarity sleep device sufficient for combinational logic block

Active Mode

Low V_t circuit operation (or combined)

Standby Mode

Disconnect power supplies through High V_t devices

- For $S=85$ mV/Decade, $\Delta V_t = 225$ mV
~1000X reduction

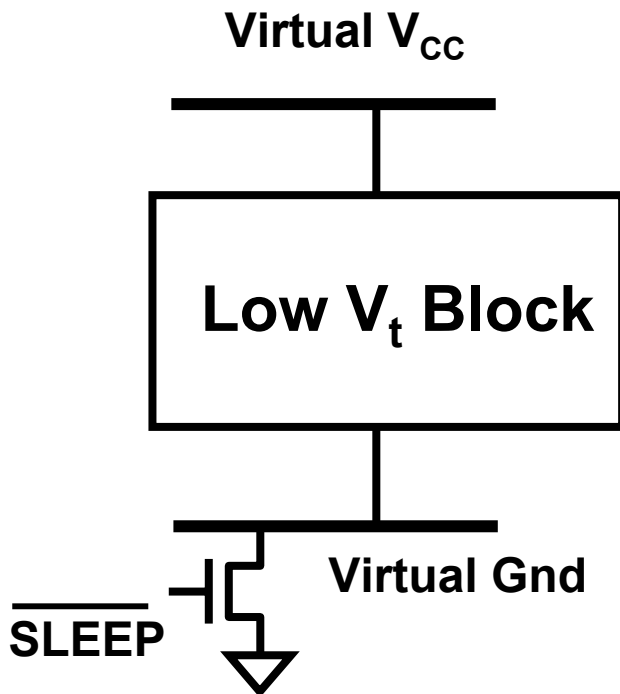
Use of LVT sleep with +/- gate (Super Cut-off/ Multi-Voltage CMOS, M. Stan, ISLPED 1998)

For fine grain sleep control

Sequential circuits must retain state

S. Mutoh, et. al., "1-V Power Supply High-Speed Digital Circuit Technology with Multithreshold-Voltage CMOS," JSSC August 1995.

MTCMOS Sleep Sizing



Virtual Ground Bounce

Gate drive decreases

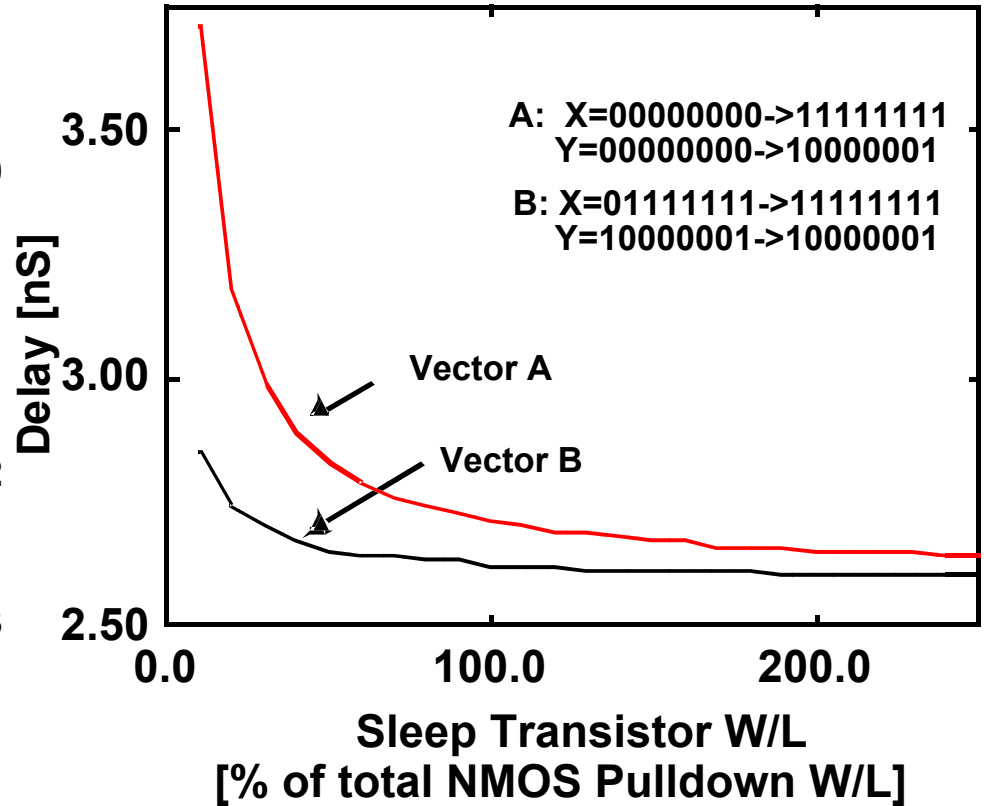
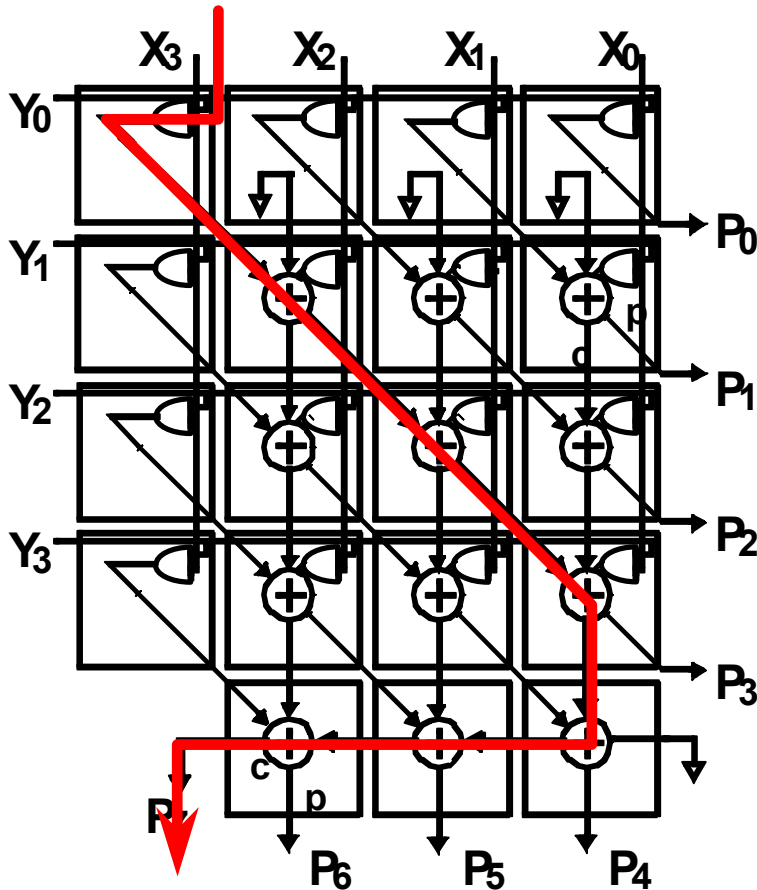
Body effect increases V_t

Reverse conduction noise concerns

Main Design Issue in MTCMOS

Properly size sleep transistor

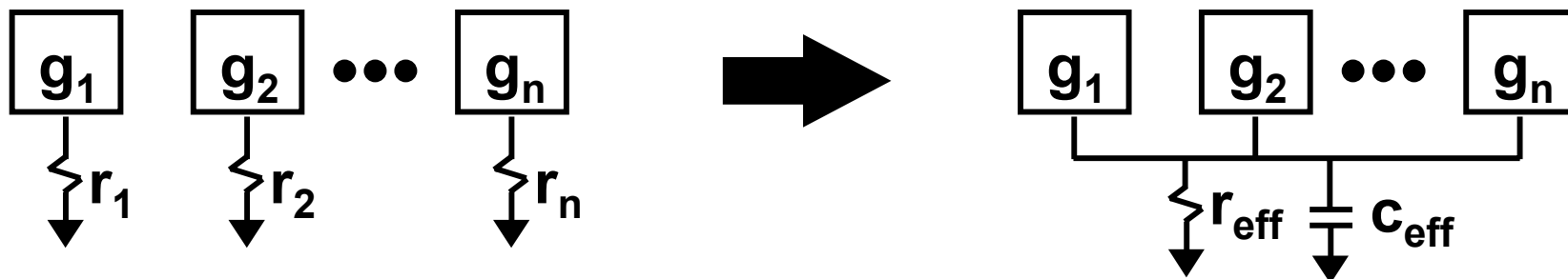
Input Vector Impact



Vector CMOS Delay % Degr (W/L=5.4%) % Degr (W/L=18%)

A	2.58 ns	15.4%	4.6%
B	2.59 ns	4.7%	1.6%

Hierarchical Sizing Approach



- **Compute effective sleep resistor for each gate**
 - Sets Maximum Gate Degradation
 - Overall delay is guaranteed
- **Mutual exclusive gates can share common sleep transistor**
- **Applied at multiple hierarchical levels**

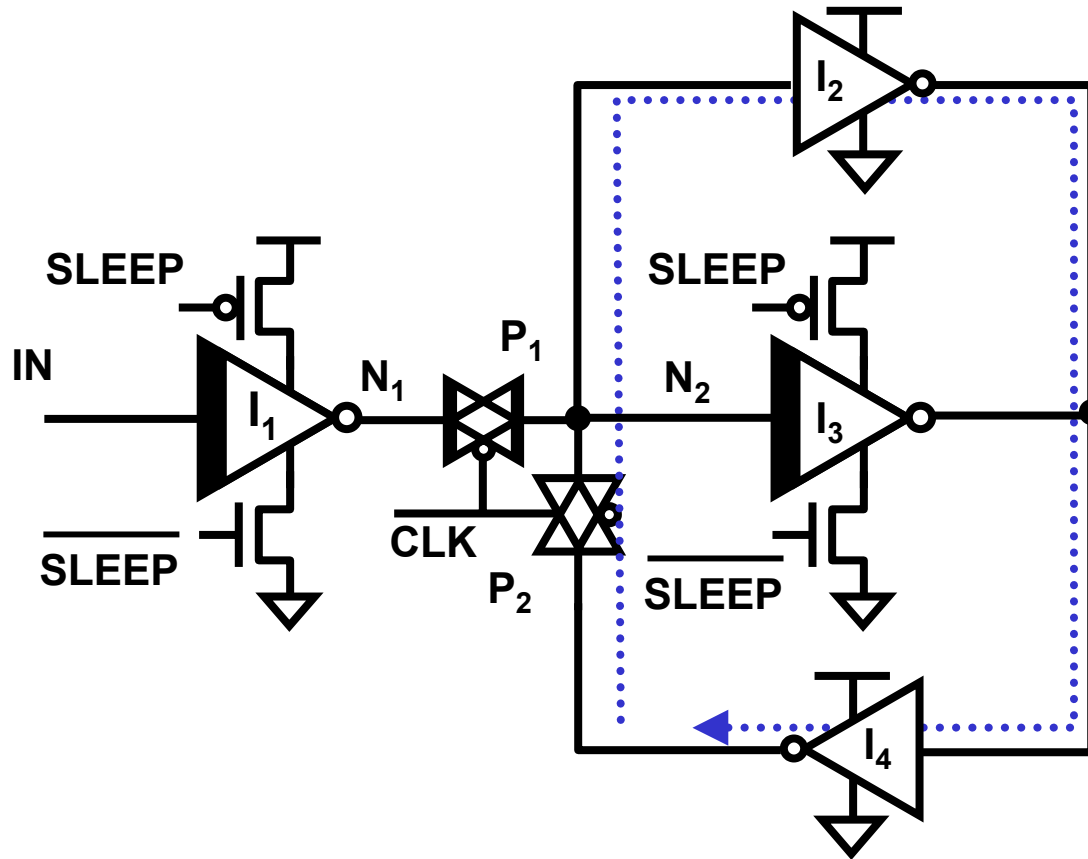
MTCMOS Sleep Sizing TBD

- **Need for improved sleep transistor sizing algorithms**
- **Static, functional timing techniques to better characterize MTCMOS discharge patterns**
- **Apply ideas from similar CAD research on power /gnd noise**
 - S. Bobba, I.N. Hajj, “Estimation of Maximum Current Envelope for Power Bus Analysis and Design,” Intl. Symp on Physical Design, April 1998
 - G. Bai, S. Booba, I.N. Hajj, “Static Timing Analysis Including Power Supply Noise Effect on Propagation Delay in VLSI Circuits,” DAC 2001
 - Y. Jiang, K. Cheng, “Dynamic Timing Analysis Considering Power Supply Noise Effects,” Int. Symp. on Quality of Electronic Design, March 2000
 - F. Najm, “Survey of Power Estimation Techniques in VLSI Circuits,” TVLSI Dec. 1994
 - H. Kriplani, F. Najm, I.N. Hajj, “Pattern Independent Maximum Current Estimation in Power and Ground Buses of CMOS VLSI circuits: Algorithms, Signal Correlations, and their Resolution,” TCAD, August 1995
 - S. Chowdhury, J. Barkatullah, “Estimation of Maximum Currents in MOS IC Logic Circuits,” TCAD, June 1990.
 - + many others ...

MTCMOS Sequential Circuits

- **MTCMOS Combinational Circuits**
 - Simple operation
 - Difficulty in sizing/ distributing sleep transistors
- **MTCMOS Sequential Circuits**
 - Virtual power/ gnd disconnected during sleep
 - Nodes will float
 - Techniques needed to maintain state
 - Need always powered circuits
 - Must avoid sneak leakage paths

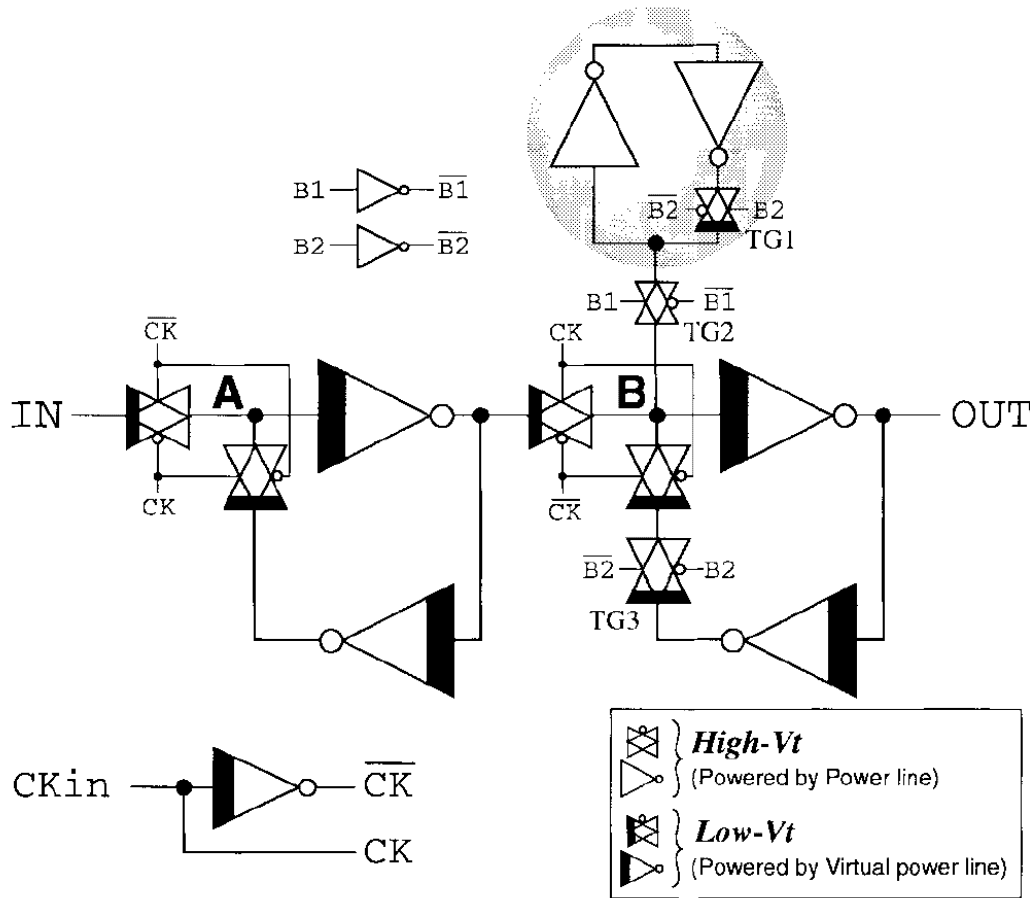
Basic MTCMOS Latch



- Use of always powered CMOS gates

S. Mutoh, et. al., "1-V Power Supply High-Speed Digital Circuit Technology with Multithreshold-Voltage CMOS," JSSC August 1995.

Balloon Flip Flop



- HVT storage “balloon” decoupled from LVT logic
- LVT blocks can share common virtual pwr/gnd
- Elimination of sneak leakage paths
- Complicated signalling

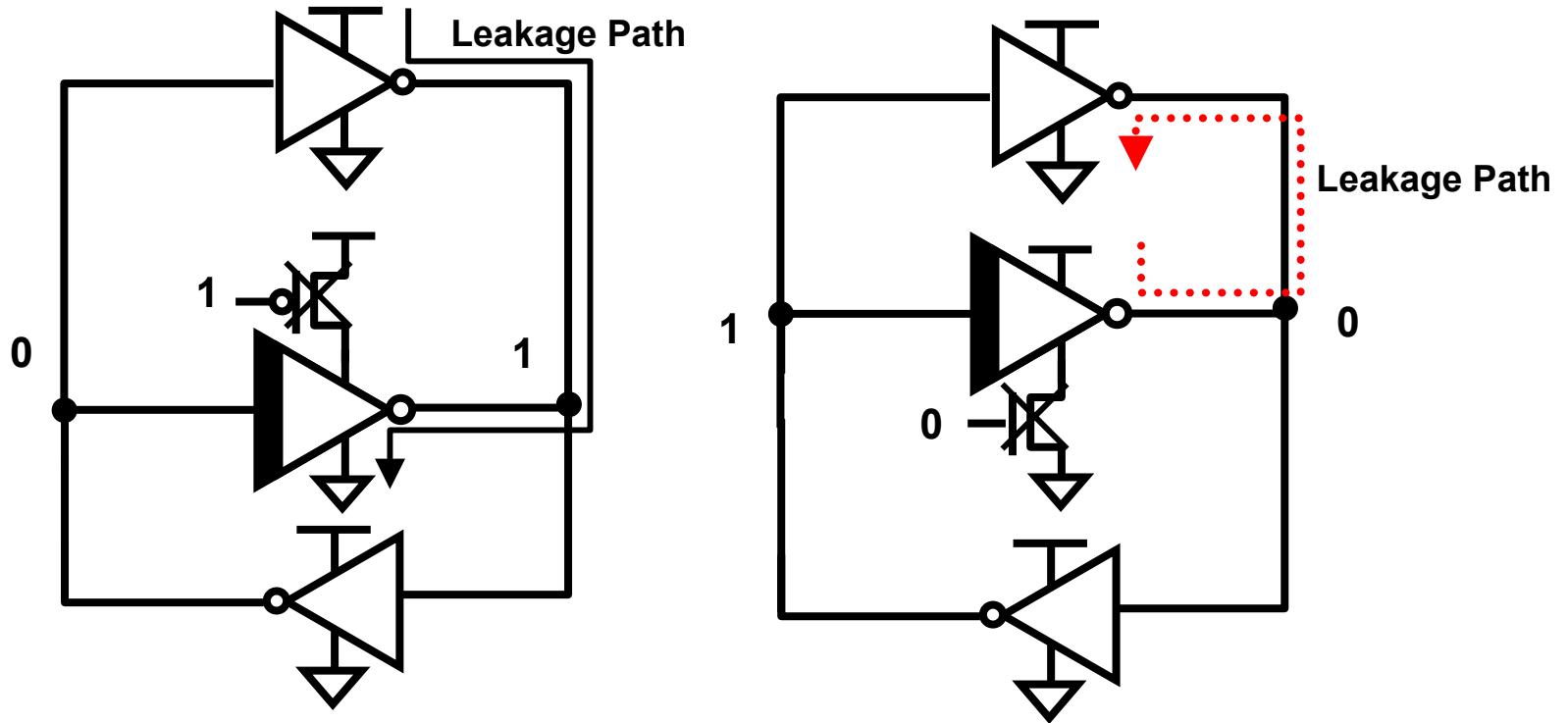
S. Shigematsu, et. al., “A 1-V High Speed MTCMOS Circuit Scheme for Power-DOWN Application Circuits,” JSSC June 1997

Fig. 9. A balloon circuit applied to a DFF circuit (clock-dependent type).

Sneak Leakage Paths

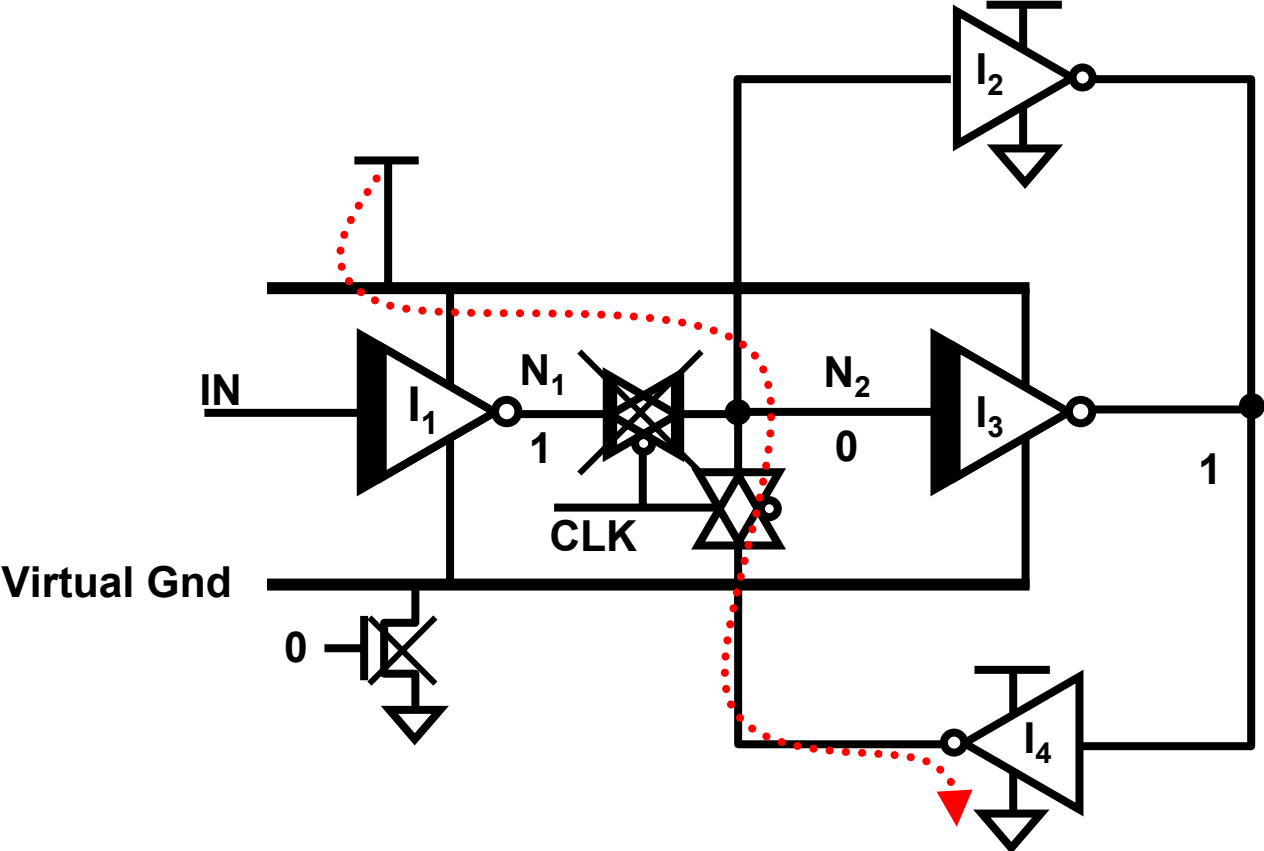
- **Sneak paths from MTCMOS/ CMOS interaction**
- **Leakage currents from V_{CC} to ground without passing through off high V_t device**
- **Need to utilize:**
 - both polarity sleep devices
 - local sleep (non shared sleep devices)
 - novel structures

Sneak Leakage Path From Parallel Combinations



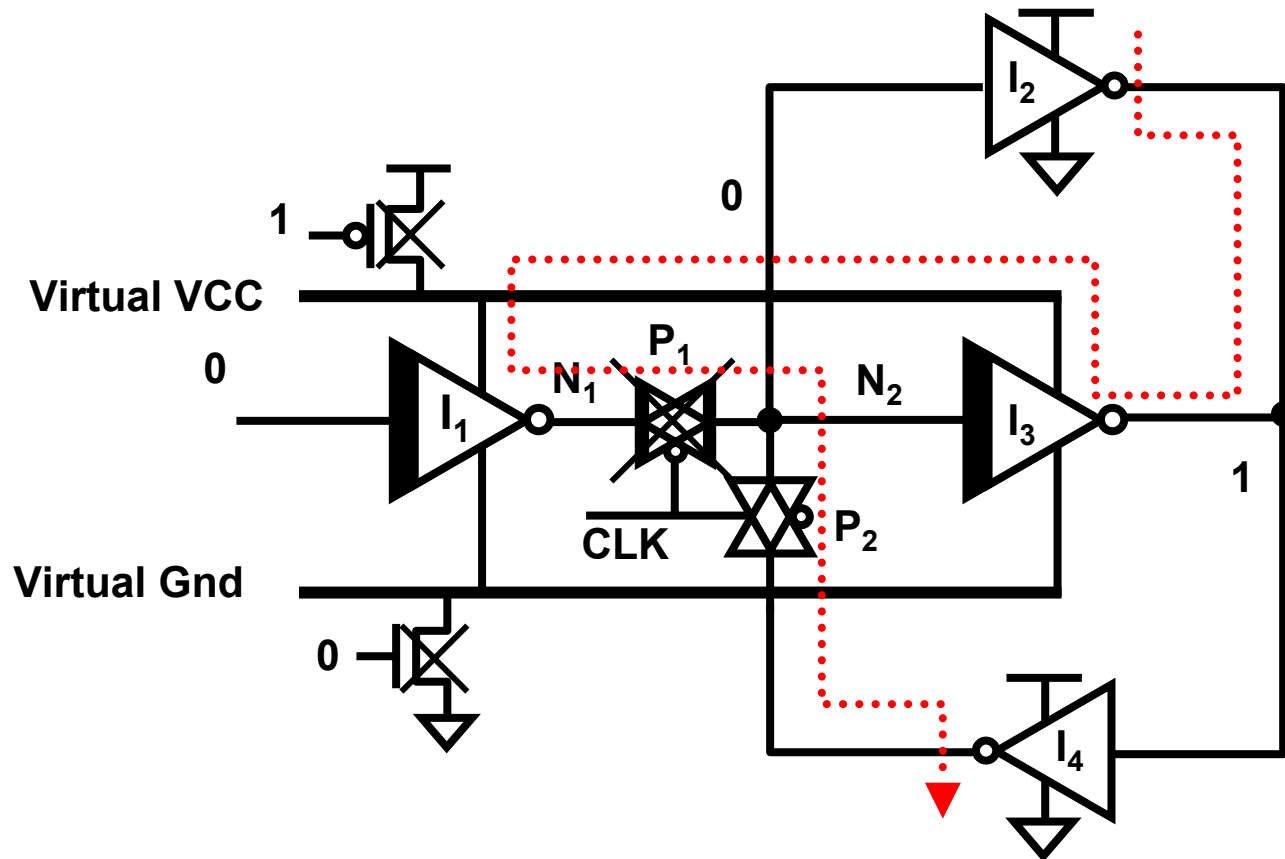
- Need for both polarity high V_t sleep devices

Sneak Leakage Path Through Low Threshold Passgate



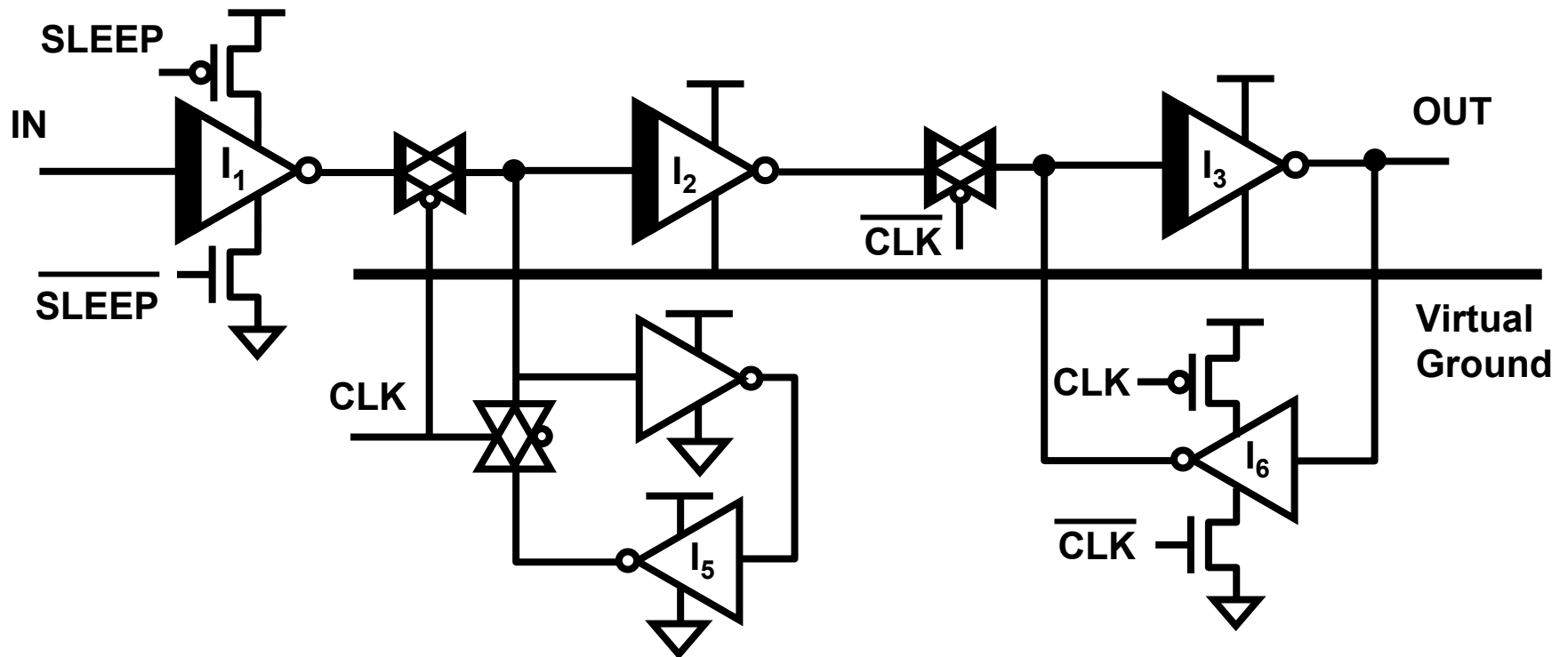
- Need for both polarity high V_t sleep devices

Sneak Leakage Path From Reverse Conduction Paths



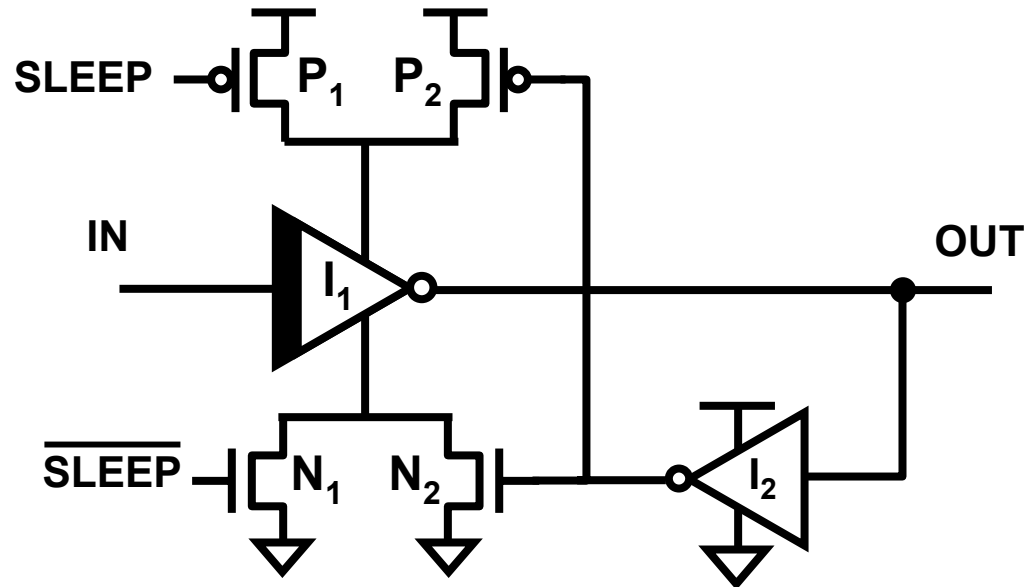
- Need for localized, non shared, high V_t sleep devices

Improved MTCMOS Flip Flop



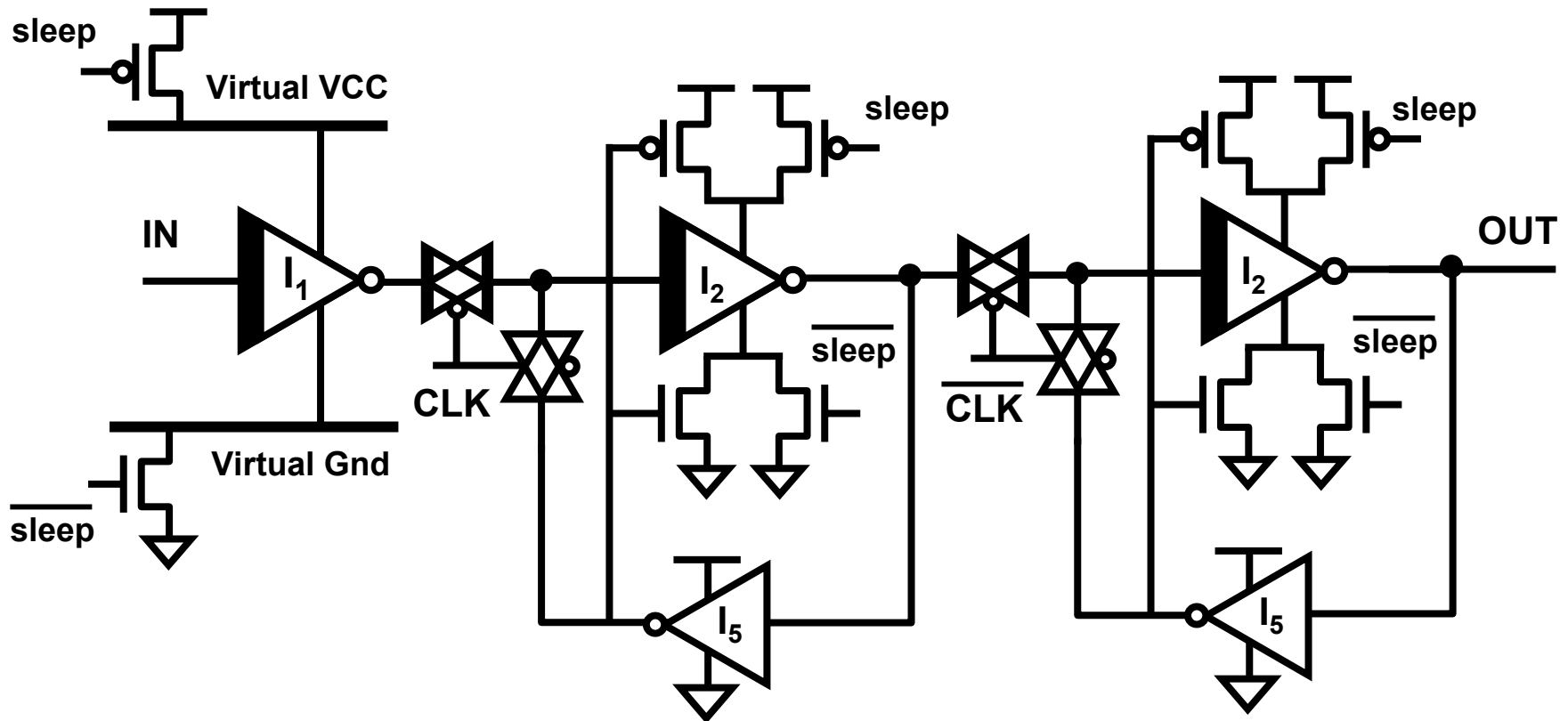
- Careful consideration of sneak leakage paths yields improved implementation

Leakage Feedback Gate



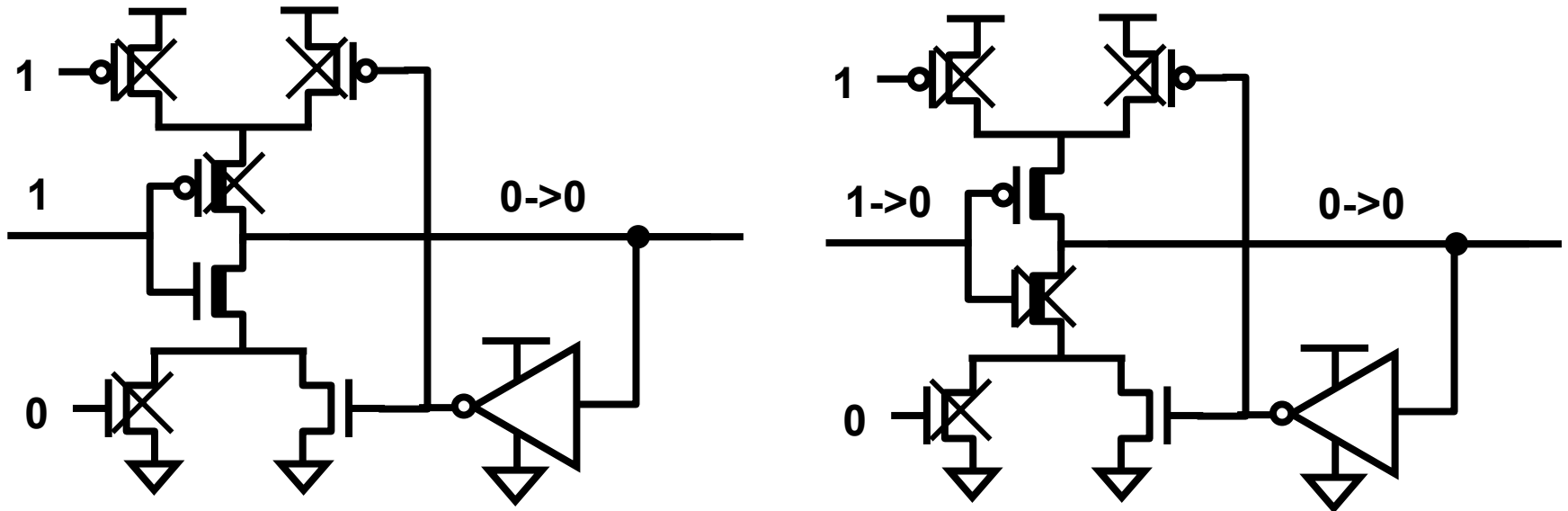
- Sufficient if either VCC or VSS path is cutoff
- Proper cutoff path yields actively driven output
- Low V_t operation + actively driven low leakage state
 - directly imported into CMOS structures

Leakage Feedback Flip Flop



- Virtually no extra loading -> performance is better than standard MTCMOS FF
- Same operation as a CMOS FF

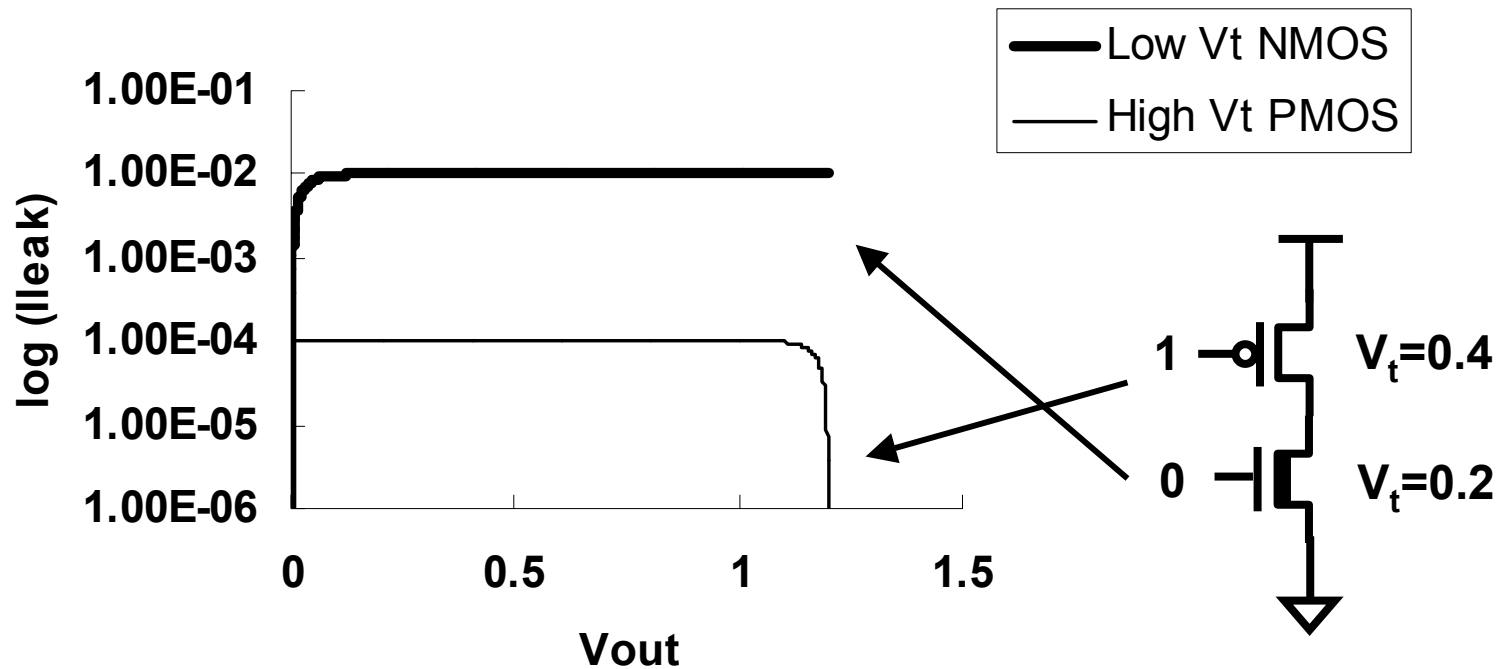
Leakage Feedback Effect



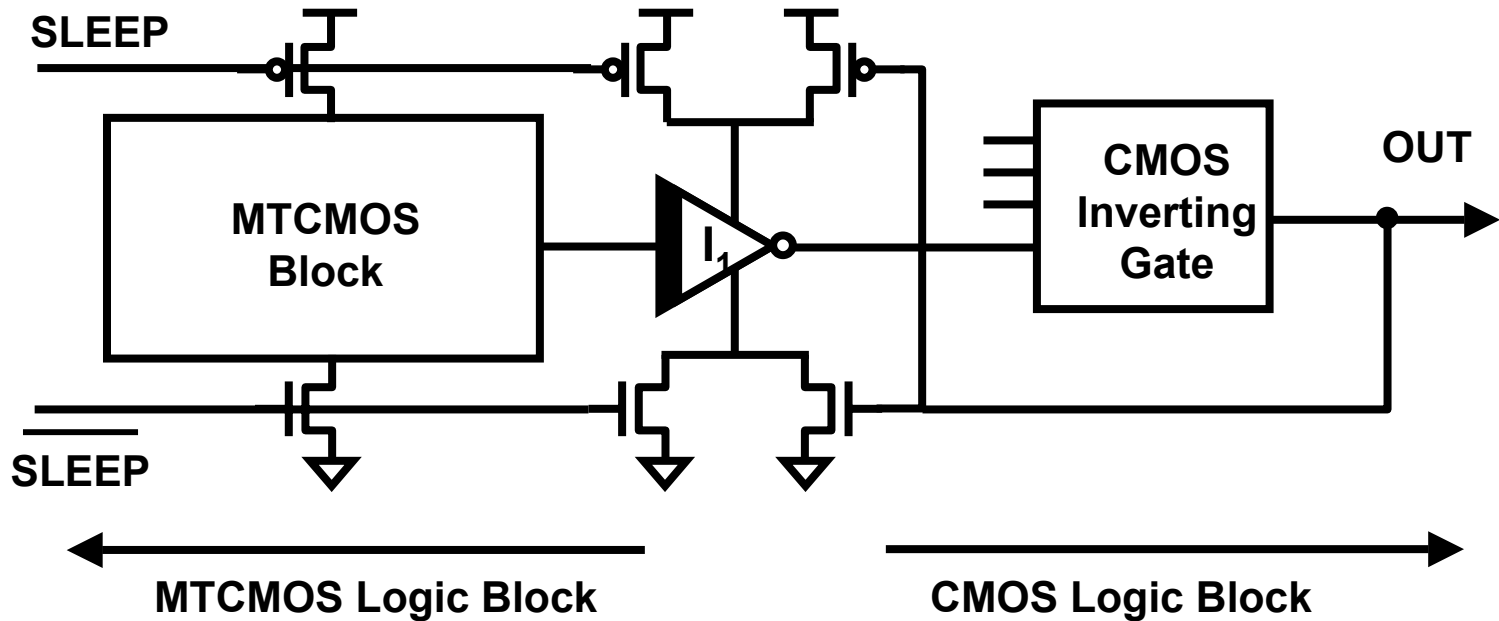
- Output data holds even if input floats
 - held by leakage mismatch
- Potential charge sharing if inputs change

Leakage Induced DC Operating Point

I-V Curves High Vt PMOS and Low Vt NMOS

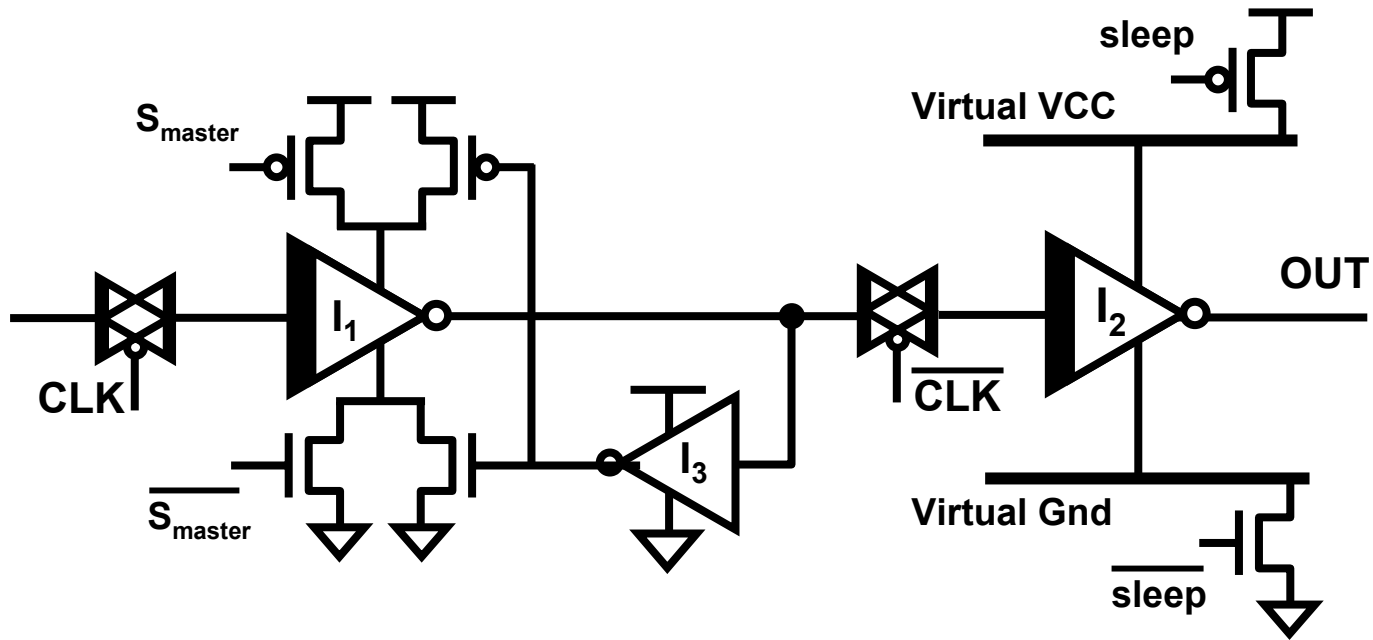


MTCMOS / CMOS Interface



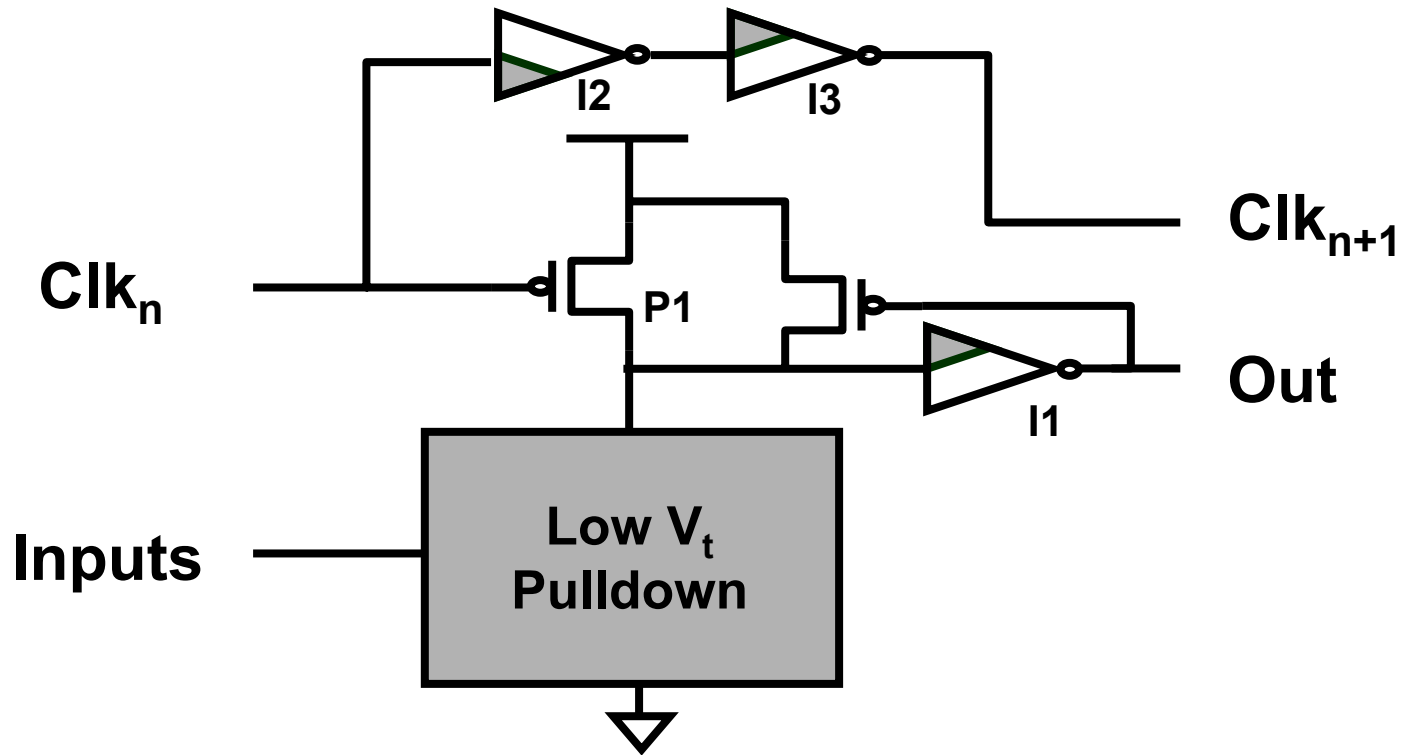
- Leakage feedback gate natural interface block between MTCMOS logic and CMOS logic

Dynamic Leakage Feedback FF



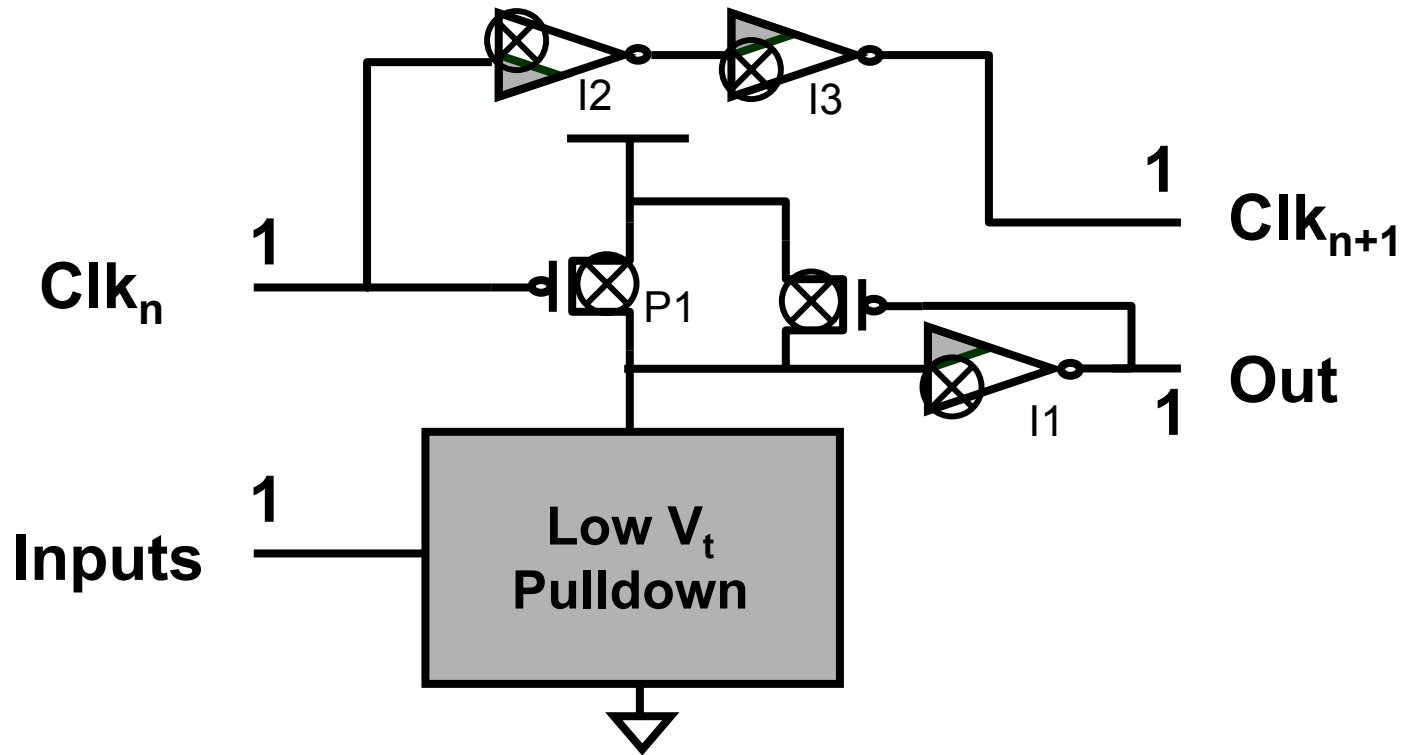
- Operates like standard dynamic FF during active mode
- Retains state during the standby mode (held by leakage)

Dual V_t Domino Gate



- Evaluate through LVT devices
- Precharge through HVT devices

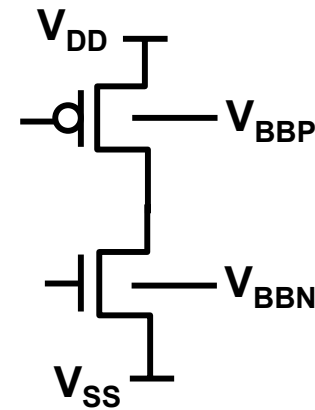
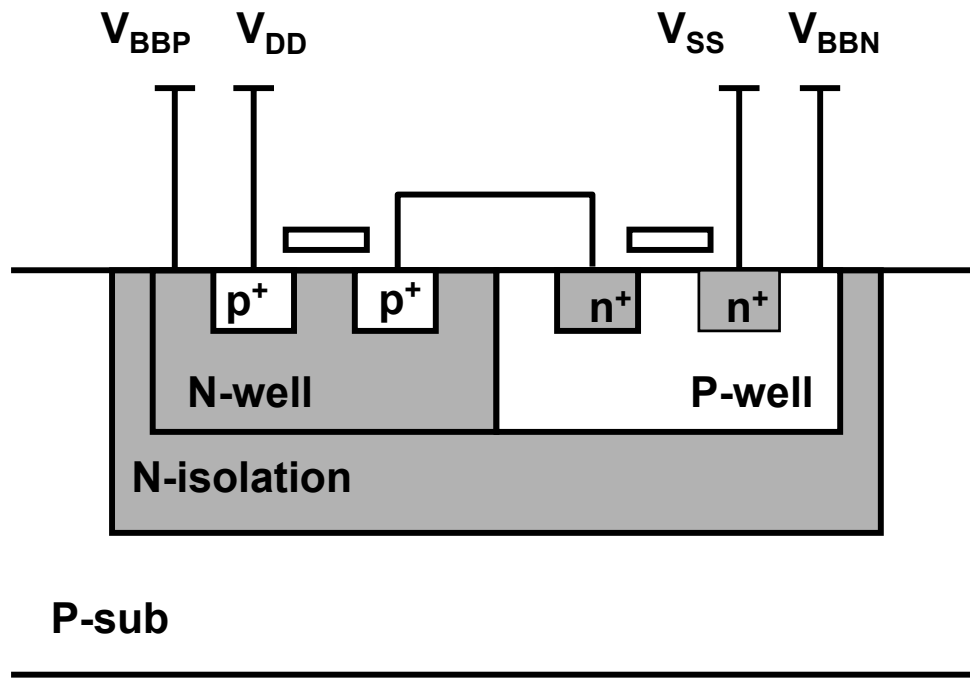
Leakage In Dual V_t Domino Gate



Sleep condition during evaluate mode

Variable Threshold CMOS (VTCMOS)

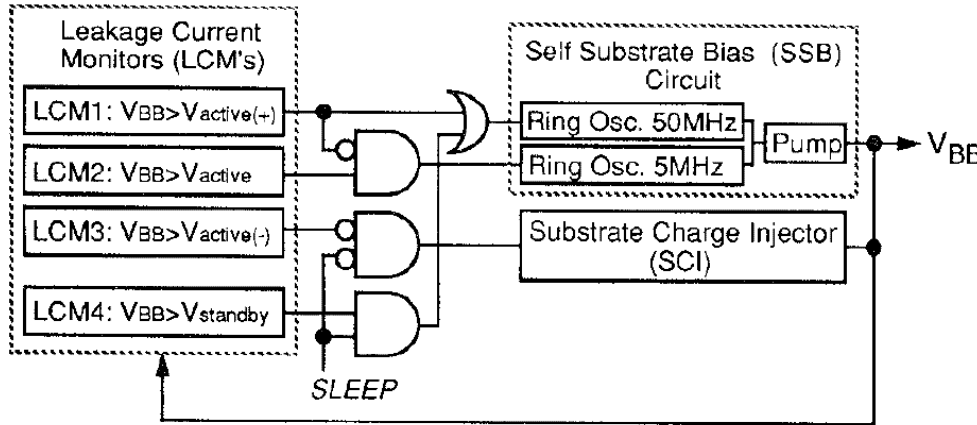
- Body effect to change device V_t
- Standby leakage reduction with maximum reverse bias
- Triple well structure



Body Effect:

$$V_t = V_{t0} + \gamma \left(\sqrt{2\phi_B - V_{BB}} - \sqrt{2\phi_B} \right)$$

VTCMOS Example



T. Kuroda, et al, "A 0.9V, 150MHz, 10mW, 4mm², 2-DCT Core Processor with Variable V_t Scheme," JSSC Nov. 1996

Fig. 3. VT block diagram.

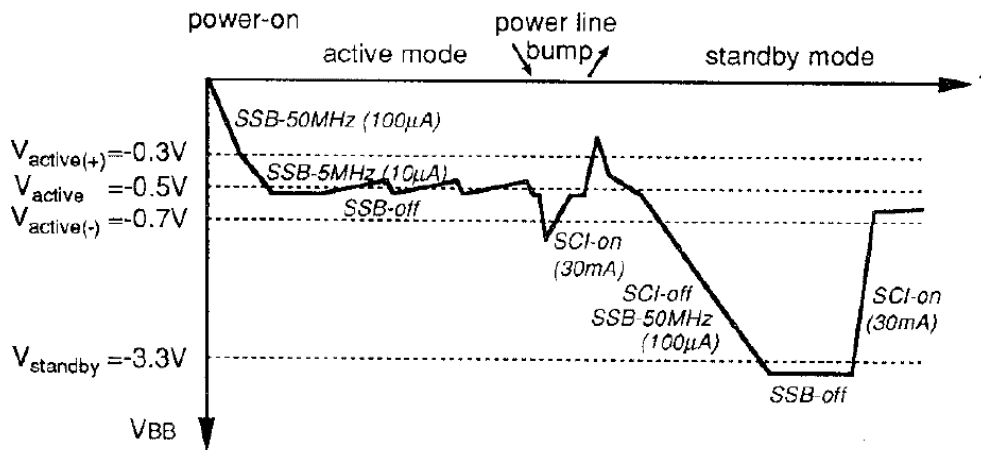


Fig. 4. Substrate-bias control in VT.

- VTCMOS principle applied to 4-mm² DCT core processor
- SSB increases V_t (more reverse bias)
- SCI decreases V_t (Standby \rightarrow Sleep)
- Leakage reduction
0.1mA active \rightarrow 10nA sleep (2.8v ΔV_{BB})
4 orders of magnitude
- Dynamically tunes V_t (by matching leakage current monitor) to minimize V_t variation

VTCMOS Pros/Cons

PROS:

- Significant standby leakage reduction
- Memory elements retain state
- No transistor sizing/ partitioning required
- Dynamically tunable V_t during runtime

CONS:

- Requires expensive triple well process
- Body factor decreases with scaling

Speed Adaptive V_t CMOS

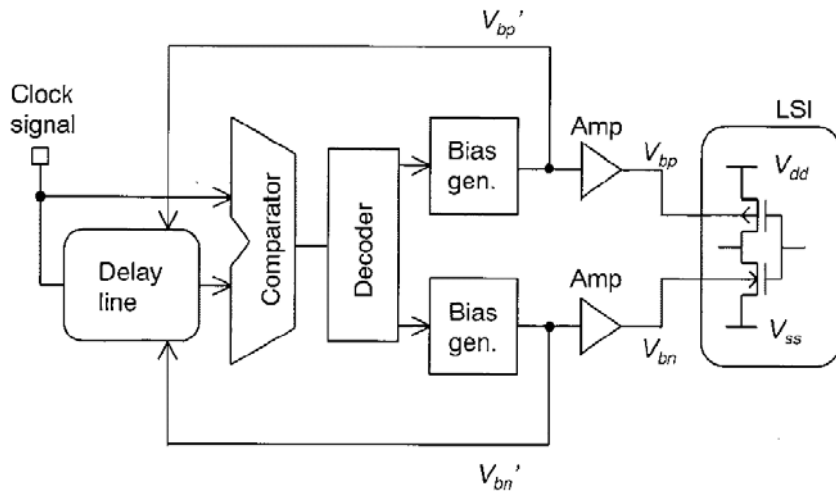
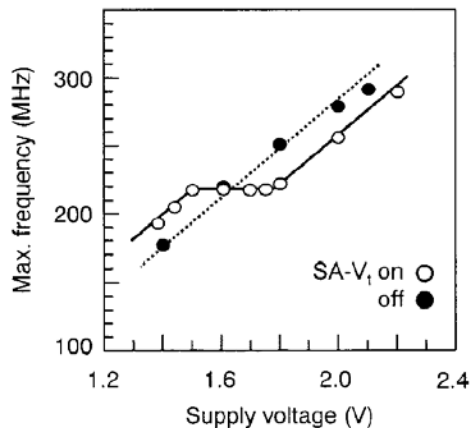


Fig. 2. Concept of SA- V_t CMOS scheme.

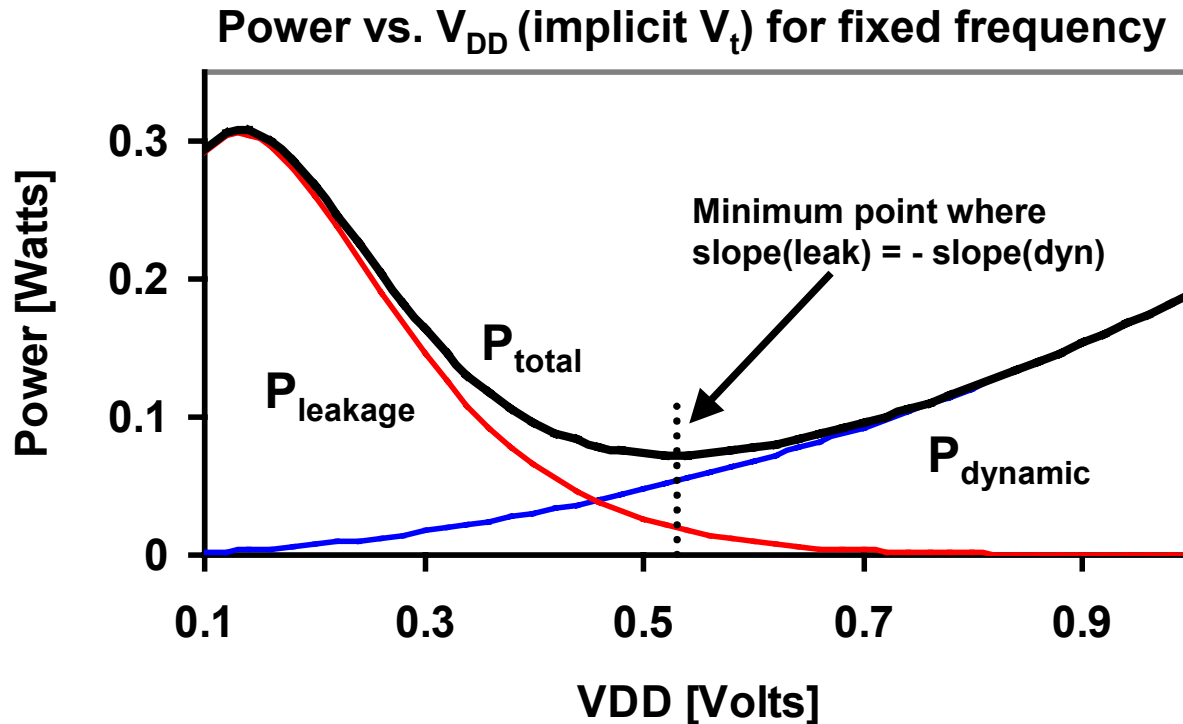


M. Miyazaki, et al, "A 1.2-GIPS/W uProc Using Speed-Adaptive V_t CMOS with Forward Bias," JSSC Feb 2002.

- Dynamically tune V_t so that critical path speed matched clock period
- Reduces chip-to-chip parameter variations
- Reverse bias:
 - Operate only as fast as necessary (reduces excess active leakage)
- Forward bias:
 - Speeds up slow chips
- Standby leakage with maximum reverse bias
- Also known as Adaptive Body Biasing (ABB)

Adaptive Supply & Body Bias (ASB)

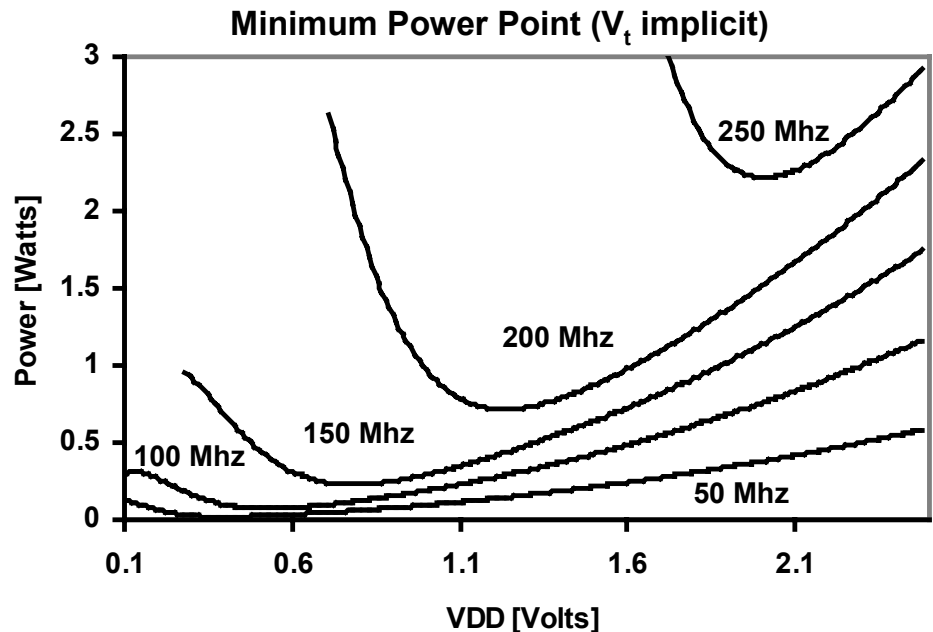
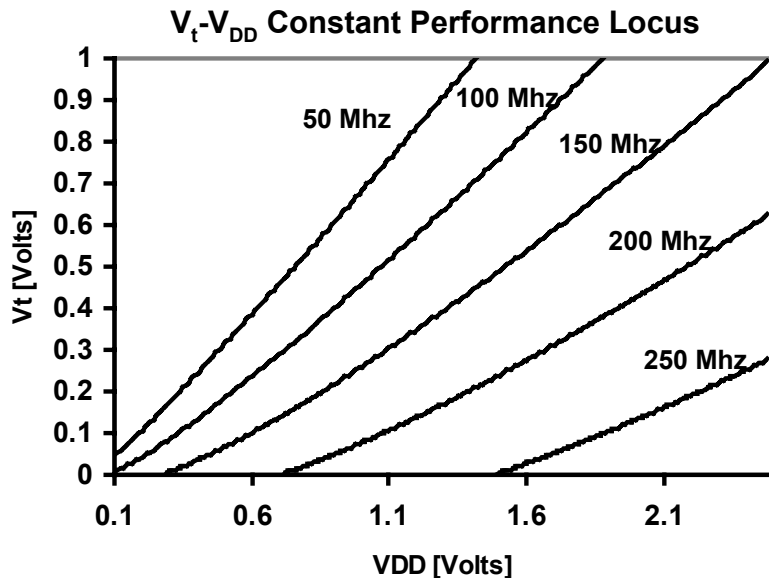
- Dynamically tune both V_{DD} & V_t as operating conditions change
- Trade-off between dynamic power (V_{DD} knob), leakage power (V_t)
- Minimize total ACTIVE power consumption
(higher active leakage current at expense of lowering dynamic power)



M. Miyazaki, et al, "A 175mV Multiply-Accumulate Unit using an Adaptive Supply Voltage and Body Bias (ASB) Architecture," ISSCC February 2002.

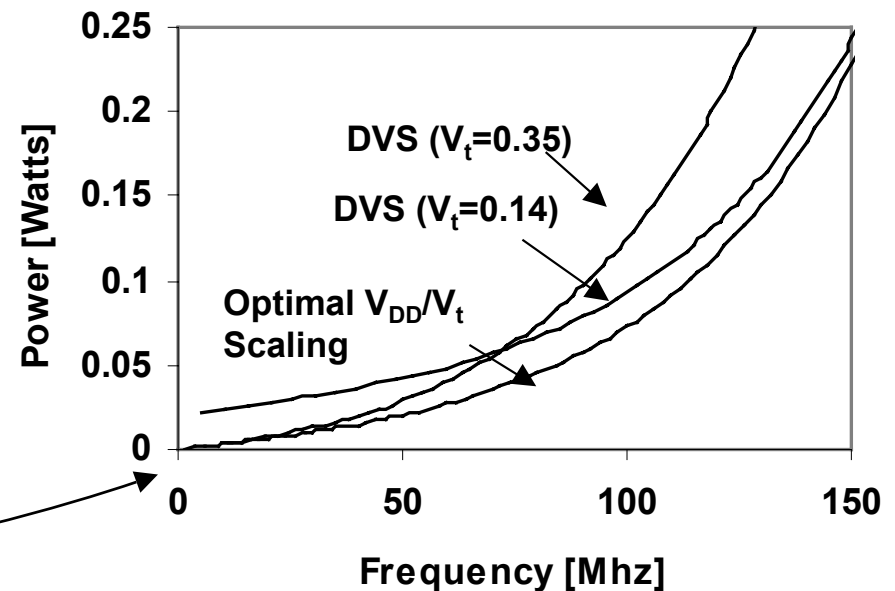
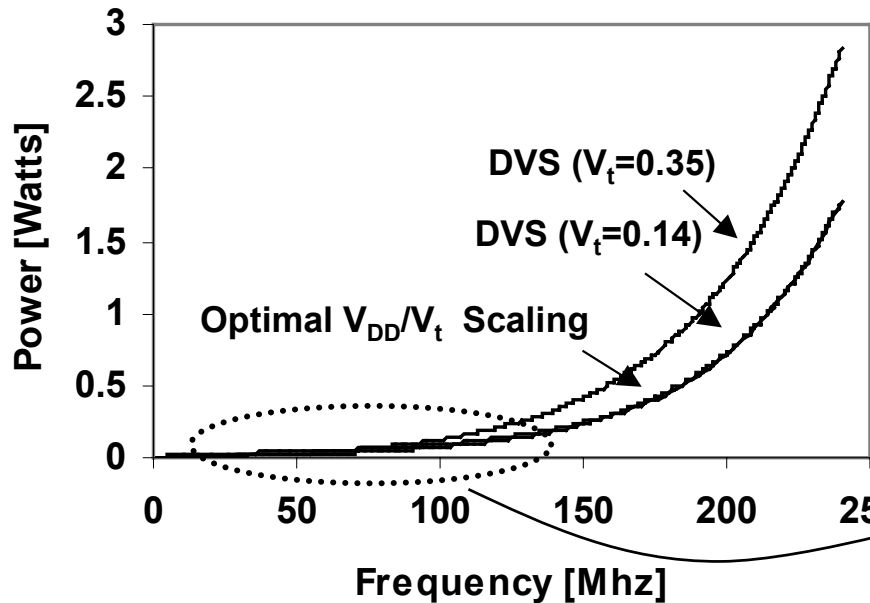
Optimal V_{DD}/V_T Selection

- Optimal V_{DD} & V_t target changes with operating conditions
 - e.g. Varying Workload
- Low frequencies high V_t more optimal
 - reduce leakage at expense of increased dynamic
- High Frequencies low V_{DD} more optimal
 - reduce dynamic at expense of increased leakage

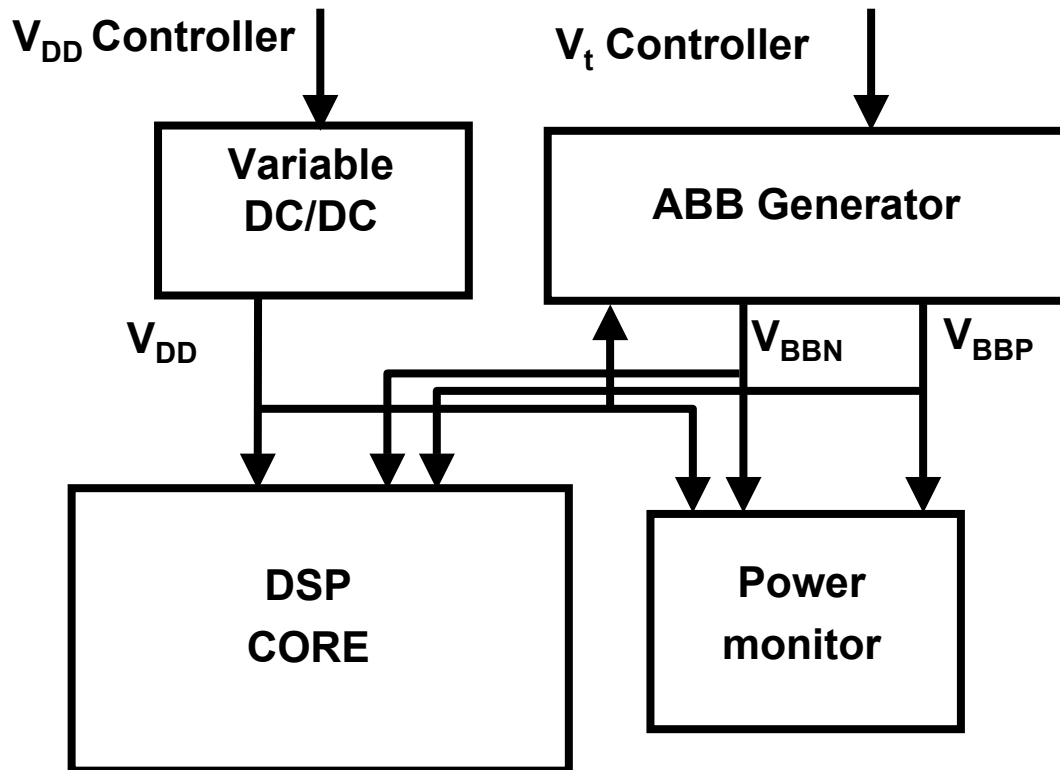


V_{DD}/V_T Optimization vs. DVS

- Dynamic voltage scaling ignores V_T influence
- DVS is sub-optimal over the frequency range



ASB Architecture



- Decouple V_{DD}/V_t tuning loops
- ABB (Auto Body Biasing) generator chooses V_t based on $V_{DD}/\text{Freq}/\text{etc.}$
- Simple V_{DD} sweep to search minimum active power point
- Architecture ensures minimum power for any operating condition

M. Miyazaki, J. Kao, A. Chandrakasan, "A 175mV Multiply-Accumulate Unit using an Adaptive Supply Voltage and Body Bias (ASB) Architecture," ISSCC February 2002.

Summary

- **Subthreshold leakage currents will grow exponentially**
- **Need to manage during STANDBY and ACTIVE**
- **Three main principles**
 - **Source Biasing**
 - **Multiple threshold voltage**
 - **Body biasing**
- **Need for CAD tools to model leakage currents**
- **Need for CAD tools to implement leakage reduction principles**