

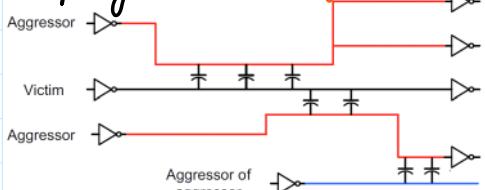
L3. Digital Noise

2023年1月3日 19:38

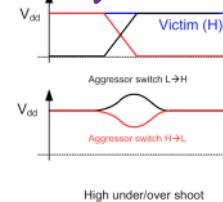
Noise Type:

- ① Digital: Deterministic, Repeatable (usually large)
- ② Analog: Random, Small.

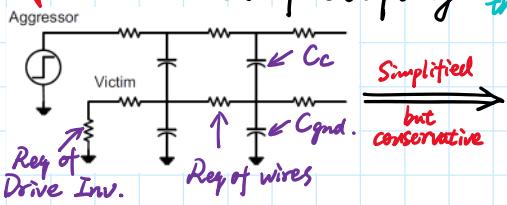
Coupling: Victim & Aggressor



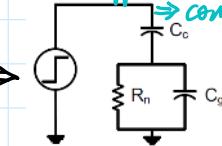
Four types of Coupling:



RL Model of Coupling:



the R of Aggressor is set as 0



Aggressor



$$tf = f(\text{Victim RC})$$

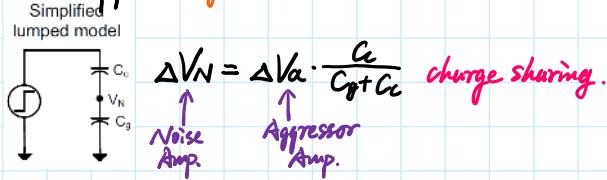
$$R \text{ discharge}$$

$$T = Rn Cg$$

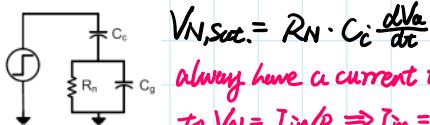
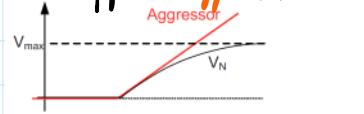
only about I_m in 1st order.

$$tr = f(\text{Aggressor slope})$$

1st Approx: Ignore $R_n \Rightarrow$ conservative



2nd Approx: Aggressor Never Sat.



$$V_{N,\text{sat}} = R_n \cdot C_g \frac{dV_a}{dt}$$

always have a current until C_g is charged to $V_N = I_m/R \Rightarrow I_m = I_R$, No more ΔV at C_g .

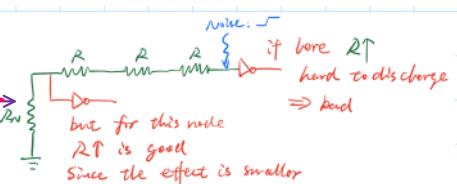
Coupling Parameter:

	V_N	t_r	t_f
Victim driver ↑	↓	same	↓
Aggressor ↓	↓	↑	same
C_g ↑	↓	same	↑
R_w ↓	Depends on topology		
C_c ↓	↓	same	↓

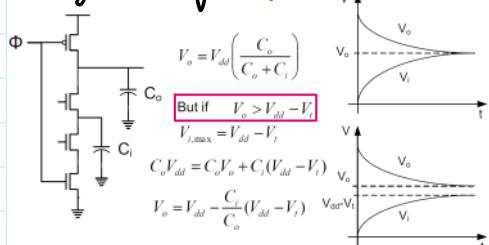
Best Solution

close to Noise → small; Far → large.

- Victim driver stronger makes it a strong aggressor
- Rise time depends on aggressor slope
- Fall time depends on victim's RC



Charge Sharing: C/Ctotal or related to Vdd. IR-Drop:



$V_o = V_{dd} \left(\frac{C_o}{C_o + C_i} \right)$

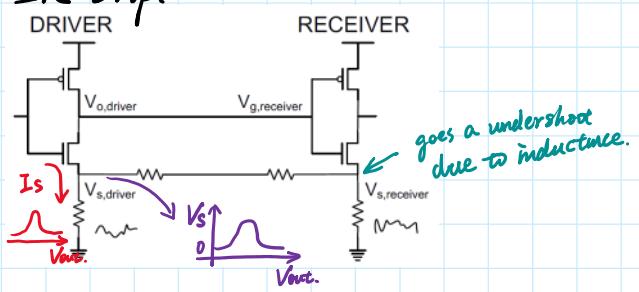
But if $V_o > V_{dd} - V_t$

$V_{o,\text{max}} = V_{dd} - V_t$

$C_o V_{dd} = C_o V_o + C_i (V_{dd} - V_t)$

$V_o = V_{dd} - \frac{C_i}{C_o} (V_{dd} - V_t)$

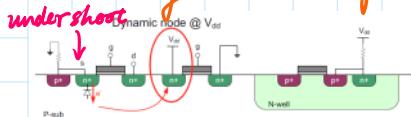
$V_{dd} - V_t$





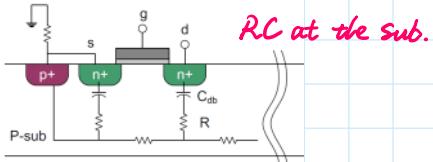
★ Substrate Coupling:

① Minority Carrier injection



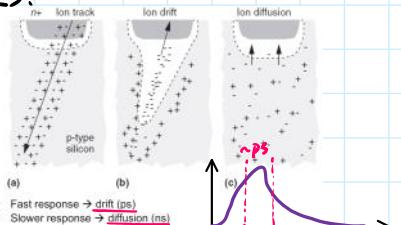
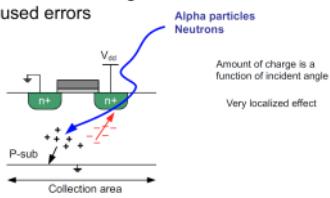
Carriers are injected by pn junction
Many hit the dynamic net.
⇒ Protect by guard ring.

② Ground Bounce:



★ Soft Error Rate (SER):

Non-deterministic generated radiation caused errors



★ Critical Parameter

1. Critical charge: Min charge stored on node
2. Collection Area: Area sensitive to SER

★ FIT: Number of Failures in 10^9 hours

★ Particuls:

Alpha Particles:

- Lead in chips and old packaging, materials:
 - Flux: 0.01 events/cm²h
 - Fix with cleaner processes (Clean lead, new materials)
 - Flux: 20x reduction with new materials
 - Background radiation:
 - Flux: 0.005 events/cm²h
 - Can be shielded (coating of the package with special materials)
 - Charge = 10-20 fC

Neutrons:

- Background cosmic radiation:
 - Flux: 0.005 events/cm²h
 - Charge: 100-200 fC
 - Cannot be shielded
- larger charge.*

More critical.

★ SER for Different Circuit.

- With scaling less charge is stored in memory cells → More sensitive to SER
- More cells/unit area → SER collection area is decreasing
- DRAM scaling → 2 effects cancel
- SRAM scaling → holding steady
- But: more multi-failure → interleaving

A small tech node:

Less Qcrit, but smaller Area.

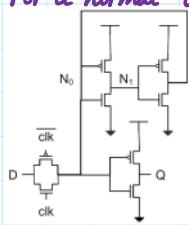
SRAM has error control code, usually detect 2-bit, fix 1-bit.

To avoid 2-bit flip, the bits of same word are placed separately.

e.g. W₁B₀ W₂B₀ W₃B₀ W₁B₁

★ Hardened Latch: Avoid SER

For a normal latch, the previous state before Soft error is unknown.



N ₀	N ₁
L	H
H	L
H	H

Normal
Alpha P.
Normal
Alpha P.

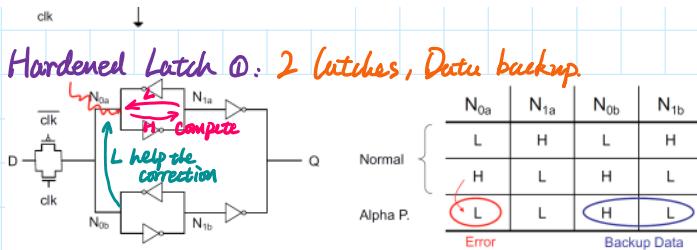
both High for N₀, N₁ — Soft error.

but 2 possible previous states

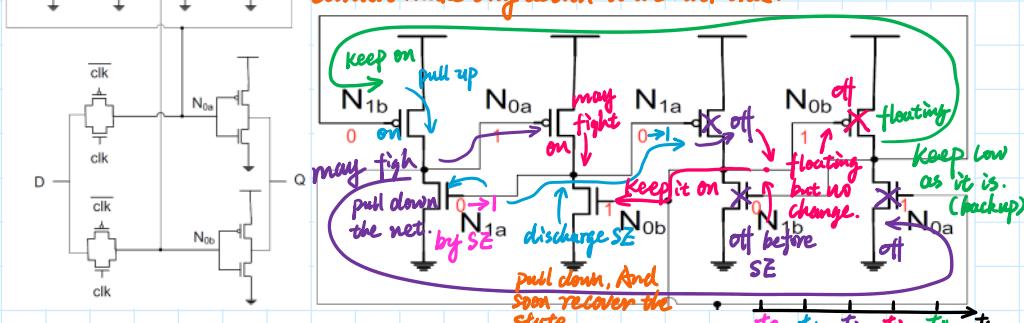
⇒ Cannot Recover.

Hardened Latch ①: 2 latches, Data backup





Hardened Latch ②: Cause floating nets for a single change if the other side keeps the same, it will be soonly corrected. Since the floating net cannot make any action to the other side.



For small changes, the SE will be recovered very fast and only a small fight. But leakage may dominate if charge is too much — may fail.

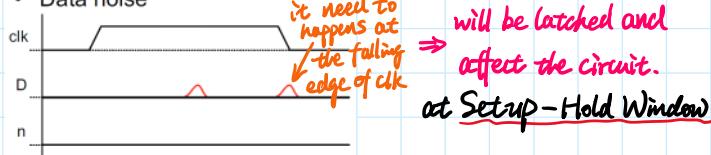
★ Noise Analysis - Global View

★ Latch Noise: Data Noise

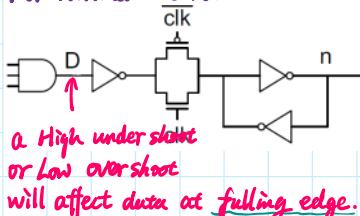
The clock noise is not common: Strong driver

Data noise is more common:

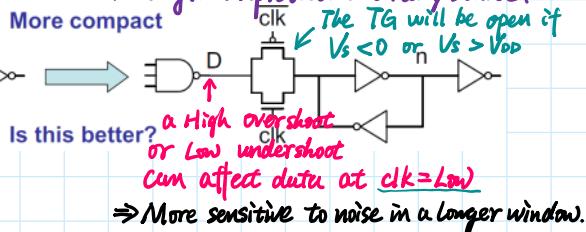
- Data noise



For normal latch:

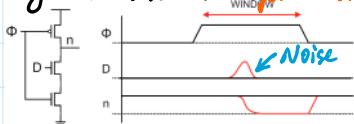


But for unprotected drain/source:

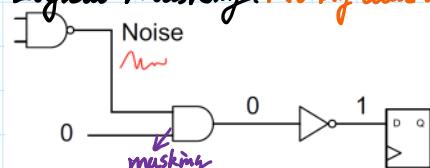


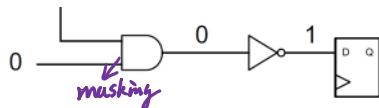
★ Latches' input must be protected.

★ Dynamic Node: Larger Sensitive Window

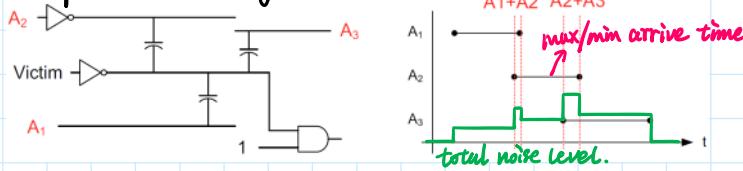


★ Logical Masking: No big deal if it can be masked at clk edge.





★ Temporal Masking:

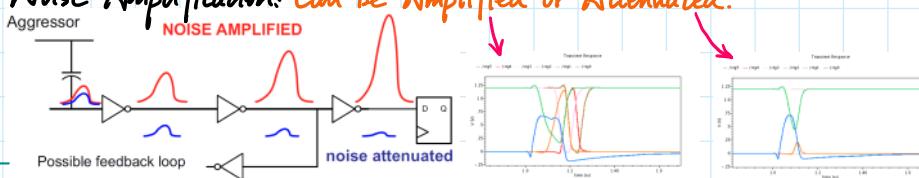


The total noise of one net is the sum of aggressors — may not be the same time
If the next stage is a FF, only need to worry about the noise at rising edge.

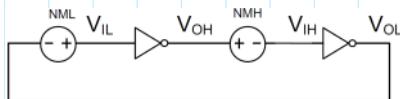
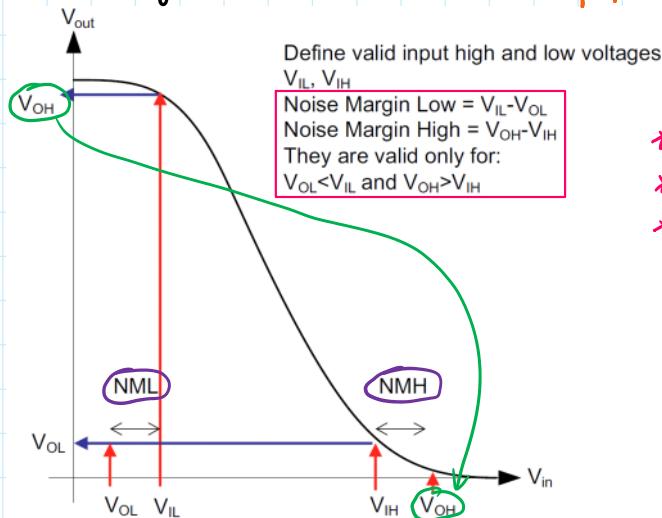
(fulling of Master)

★ Noise Margin - Local Analysis

★ Noise Amplification: Can be Amplified or Attenuated.



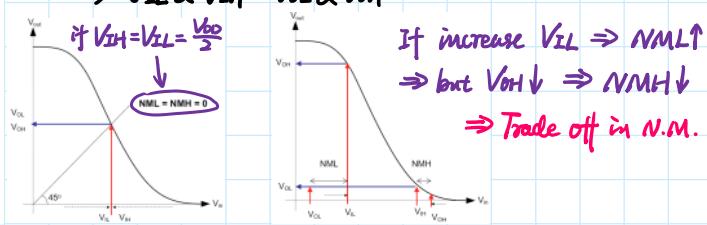
★ Noise Margin: A local constraint to avoid amplification



* V_{IL} & V_{IH} : Worst Condition allowed for input (set by human)

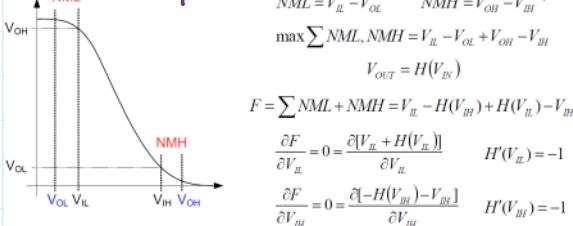
* V_{OL} & V_{OH} : Worst Condition of output (due to V_{IL} & V_{IH})

* Noise Margin: Margin left for noise under Worst case.
 $\Rightarrow V_{IL} \& V_{IH} - V_{OL} \& V_{OH}$

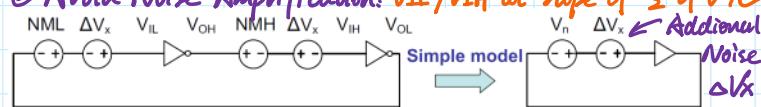


★ Optimal Noise Margin

① Max Sum of NML + NMH: V_{IL}/V_{IH} at Slope of -1 of VTC



② Avoid Noise Amplification: V_{IL}/V_{IH} at Slope of $\frac{1}{k}$ of VTC



k is the gain of the inverter

$$\Delta V_o = k \Delta V_i \quad \text{for } k < 1, \Delta V_o = k \cdot \left(\frac{1-k^n}{1-k} \right) \quad \text{Not safe.}$$

$$\Delta V_i = \Delta V_o + \Delta V_n \quad [\Delta V_x + k \Delta V_o] k = \Delta V_o \cdot k + \Delta V_x \dots \left[\frac{\partial V_o}{\partial V} = \infty \right]$$

• Unity gain $k=1$ oscillate

k is the gain of the inverter

$$\Delta V_o = k \Delta V_I$$

$$\text{for } k < 1, \Delta V_o = k \cdot \left(\frac{1-k^n}{1-k} \right)$$

$$\Delta V_I = \Delta V_x + \Delta V_o$$

$$[(\Delta V_x + k \Delta V_o)k + \Delta V_o]k + \Delta V_o \dots$$

$$\Delta V_o = k(\Delta V_x + \Delta V_o) \Rightarrow \Delta V_o (1+k+k^2+k^3 \dots) \cdot k$$

$$\Delta V_o (1-k) = k \Delta V_x$$

$$\Delta V_o = \Delta V_x \frac{k}{1-k}$$

$$\frac{\partial V_o}{\partial V_x} = \frac{k}{1-k}$$

- Unity gain $k=1$ oscillate
Not safe.
- For $\frac{\partial V_o}{\partial V_x} = 1 \quad k=1/2$
unity noise gain

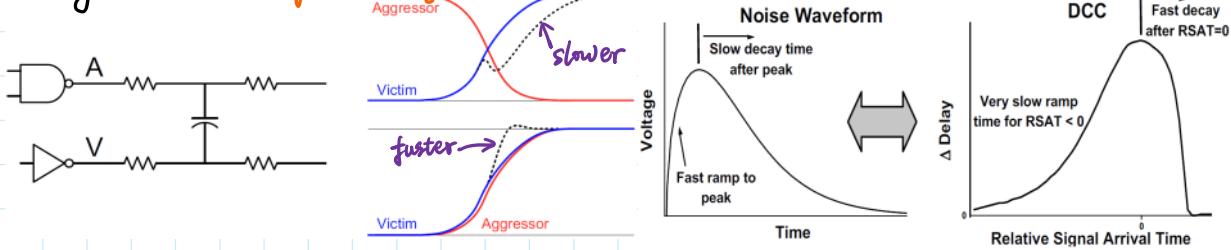
★ For safety use $k=-\frac{1}{2}$ point (actually close to $k=-1$ since transition is sharp)

★ Estimation in N.M. Analysis:

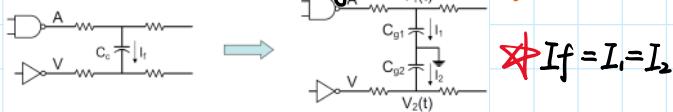
* Optimistic: 1. only VTC of INV (ignore other gates like NAND/NOR)
2. Assume nominal gates (No PVT)

* Pessimistic: 1. It is DC Analysis, actual noise has limited width
2. No masking
3. Not every gate is in the limit ($V_{IN} = V_{IL}/V_{IH}$) — N.M. is the worst case

★ Delay Noise: Change the delay of net, not Data.

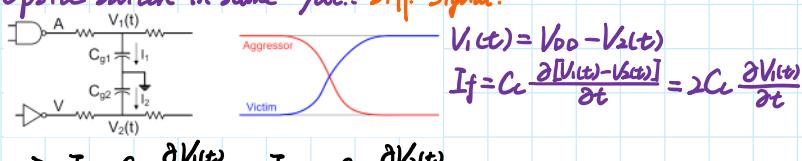


★ Effective Caps of Coupling: $C_c \Rightarrow C_{g1} \& C_{g2}$.



★ If $I_f = I_1 = I_2$

① Oppsite Switch in same dV/dt: Diff. Signal.



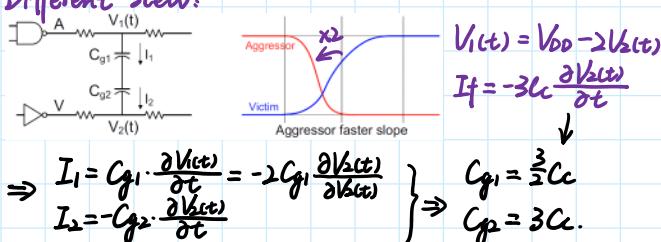
$$\Rightarrow I_1 = C_{g1} \frac{\partial V_1(t)}{\partial t} = -I_2 = -C_{g2} \frac{\partial V_2(t)}{\partial t}$$

$$\Rightarrow C_1 = C_2 = 2Cc$$

② Same switch with same slew:



③ Different slew:



★ Noise Avoidance:

- Size up victim driver :

- Makes it a stronger aggressor
- Only effective if wire is not too long

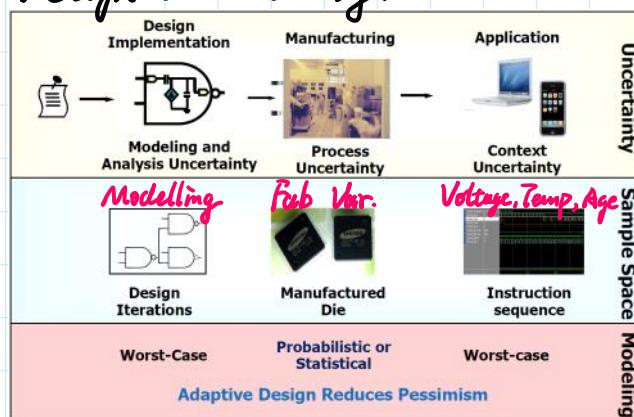
- Spacing

- Coupling capacitance is reduced
- Total capacitance is reduced

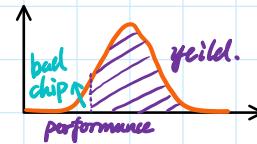
- Size up victim driver :
 - Makes it a stronger aggressor
 - Only effective if wire is not too long
 - Delay impact (increases input load)
- Widening wires:
 - Reduces resistance → Good for long wire with strong driver
 - Increases grounded capacitance
 - Can be good (or bad) for delay
- High level metal → Wider and more spacing
- Bus encoding
 - Minimize the number of transitions in a bus
 - Address/Data buses have different behavior
- Minimize capacitive coupling:
 - Switch the shields with the victim
 - Speed up delay
- *sacrificial wire.*

- Spacing
 - Coupling capacitance is reduced
 - Total capacitance is reduced
 - Good for delay
 - Area penalty
- Shielding (*always for Clk*)
 - Coupling capacitance is eliminated
 - Total capacitance is increased
- Insert buffers (Intel: every ~300um)
- Staggered buffers
 - Injects both positive/negative noise → cancel

★ Design Uncertainty:



It is a statistical process
Usually in random distribution



★ Modeling Errors: Model may not be accurate for some case (e.g. sub-Vth)

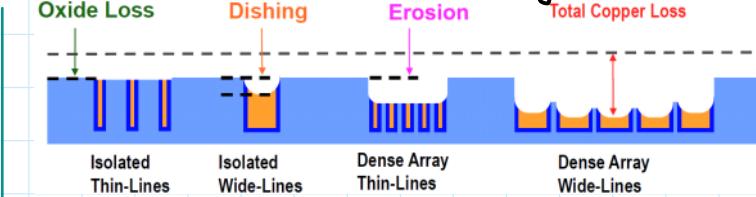
★ Manufacturing Variation

e.g. $T_{ox} : \frac{\delta}{\mu} = 3\%$; $L_{eff} : \frac{\delta}{\mu} = 5\sim7\%$; $V_{th} : \frac{\delta}{\mu} = 5\%$

⇒ Different Process has different effect to variation.

One physical para. can cause multi device para. variation.

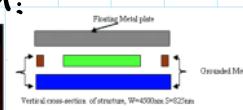
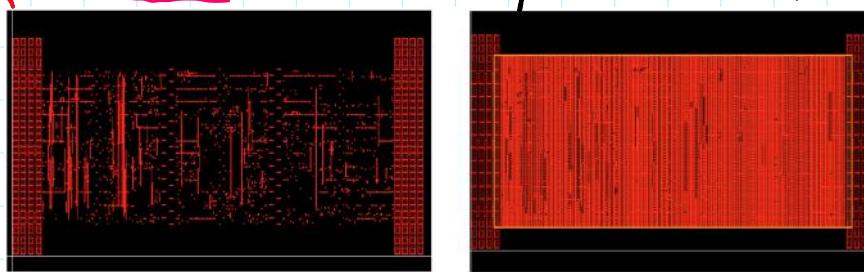
★ Chemical Mechanical Polishing: Metal thickness



Different Metal Density has different Polishing depth (SiO_2 is harder than metal).

$T_{metal} \downarrow : C \downarrow$ but $R \uparrow$

★ Use Metal Fill to have uniform metal distribution:

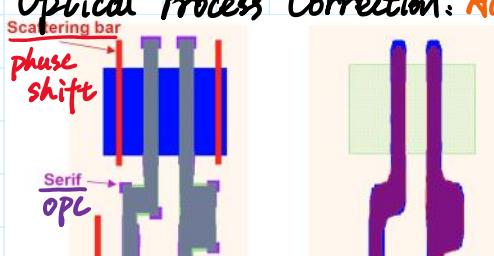


- * If M_{fill} is grounded: Large C_g but small C_c .
- * If M_{fill} is floating: Large C_c

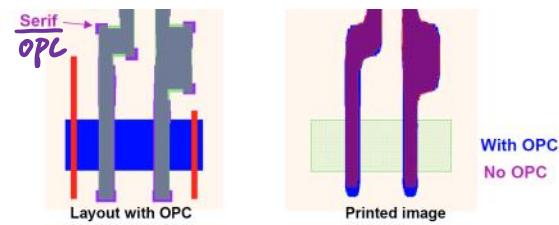
★ Lithography Related Variation: With resolution enhanced tech (RET)

Min Feature Size: $k_1 \cdot \lambda / NA$

★ Optical Process Correction: Add pattern to compensate

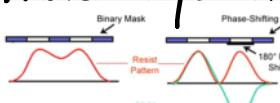


Since litho may not good at the edges of layout.



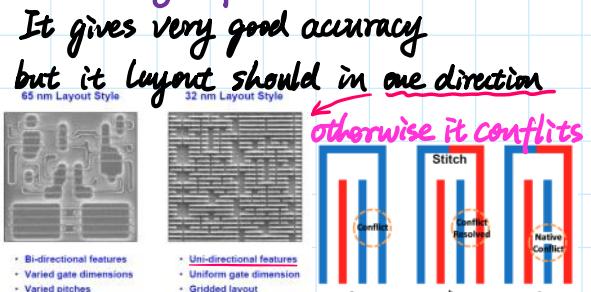
more more more given at the edges of layout.

★ Phase Shifted Masks Enhance Resolution

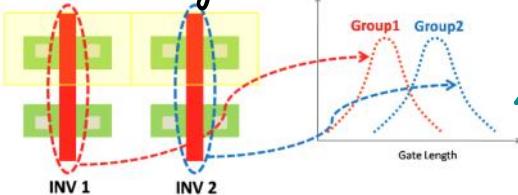


★ Double Patterning Lithography: closed gate may have different variation

↳ Two ways of Double Patterning:

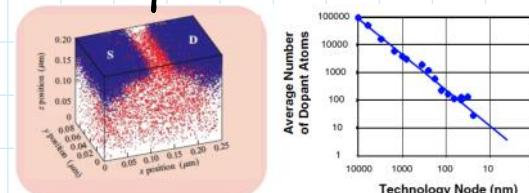


It is more easy to cause mismatch - Especially for SRAM



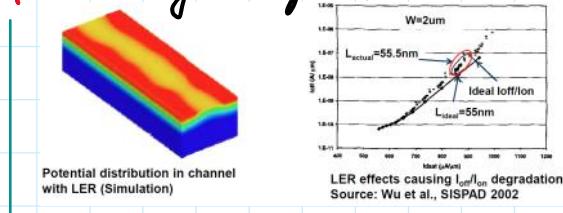
Different Process
Different Distribution

★ Ion Implantation - Random Dopants.



Number of Dopants is small in advanced node
Maybe only ~10 dopants
⇒ More easy to have variation

★ Line Edge Roughness: Random Distribution of Photon



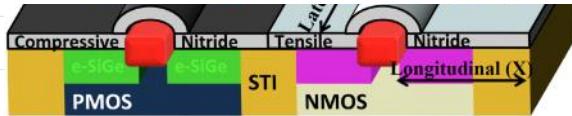
⇒ changed Ion/Ioff and on-off ratio if Lef changed

★ Mobility Enhancing - Stress

	NMOS	PMOS
X	Tensile	Compressive
Y	Tensile	Tensile
Z	Compressive	Tensile



X	Tensile	Compressive
Y	Tensile	Tensile
Z	Compressive	Tensile



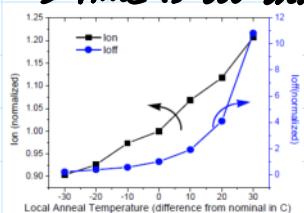
The stress changes with the layout of drain & drain contact

* RTA Variation:

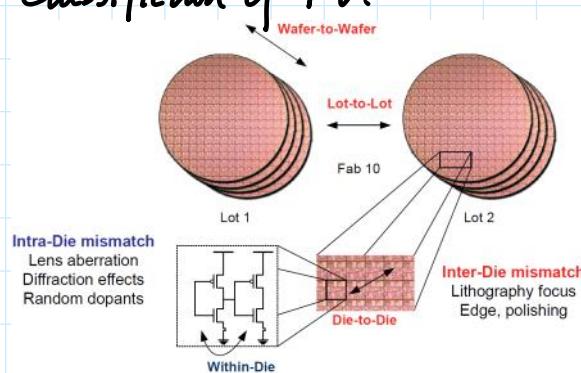
To avoid the gate metal damage, use rapid thermal annealing (RTA)

⇒ Time is too short, there may be Temp. Variation specifically.

⇒ changed V_{th} & R_{ext} and I_{on}/I_{off} .



* Classification of PV:



1. Lot-to-Lot
2. Wafer-to-Wafer
3. Die-to-Die
4. Within-Die

Process Variation

Systematic Error

- RTA
- Stress
- OPC
- (repeatable)

Random Variation

WID error
also DTD

Die-to-Die

Within Die

Special Correlated

- Gate Length
- Dose of Light
- Etch Rate

Independent

- RDF
- LZR

* There will be a batch calibration that test same chips and calibrate the whole batch ⇒ Less Die-to-Die variation is easy for calibration

* Simulation of PV:

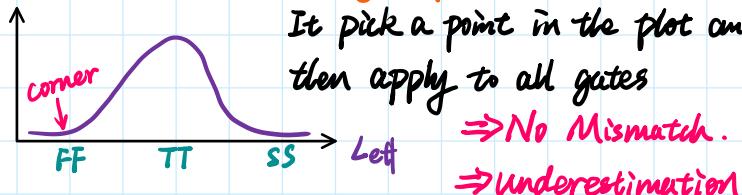
E.g. for gate length.

$$L_{eff} = \underbrace{L_{nom} + \Delta L_{system}}_{\text{Shared}} + \underbrace{\Delta L_{Die-to-Die}}_{\text{Die-to-Die}} + \underbrace{\Delta L_{random}}_{\text{Mismatch}}$$

* Corner Simulation: only captured Die-to-Die variation

Shared Mismatch

★ Corner Simulation: only captured Die-to-Die Variation



It pick a point in the plot and then apply to all gates

⇒ No Mismatch.

⇒ Underestimation

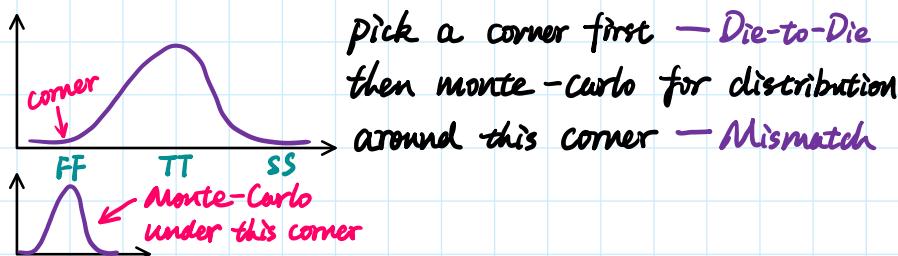
★ Global Monte-Carlo:

Randomly pick points on the distribution

⇒ It includes Die-to-Die Variation into Local Simulation

⇒ Overestimation

★ Local Monte-Carlo: Most Accurate



pick a corner first — Die-to-Die
then monte-carlo for distribution
around this corner — Mismatch

★ PV and Circuit Delay:

★ Gates in Series: t_d is the Sum.

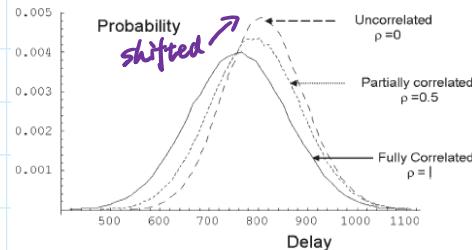
e.g. $\rightarrow \square \rightarrow \square \rightarrow \dots$

If uncorrelated paths: Mean adds, $Std = \sqrt{G_1^2 + G_2^2 + \dots} = \sqrt{n}G$

$$\star \left(\frac{G}{n}\right)_{path} = \frac{1}{\sqrt{n}} \cdot \left(\frac{G}{n}\right)_{gate}$$

If correlated:

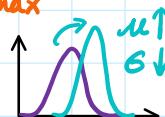
$$\star \left(\frac{G}{n}\right)_{path} = \sqrt{\frac{1+\rho(n-1)}{n}} \left(\frac{G}{n}\right)_{gate}$$



★ Gates in Parallel: $t_d = \lceil t_{path} \rceil_{max}$

e.g. $\overbrace{\square \square}^B - t_d$.

If uncorrelated paths: Mean ↑; Std ↓



If correlated paths and $\rho=1$: Mean & Std keeps the same

★ Operating Context Variation

1. Supply Voltage: IR Drop, L di/dt Drop

~~XP~~ Operating Context Variation

1. Supply Voltage: IR Drop, $L di/dt$ Drop
2. Temp. Hot Spot
3. Reliability: NBTI/PBTI, TDDDB, HCI (hot carrier degradation)
4. Electromigration

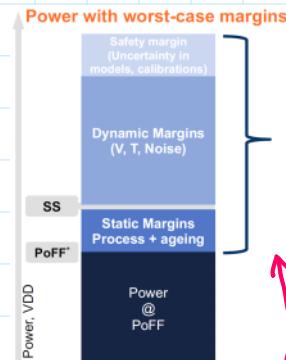
L5. Adaptive Design

2023年2月18日 23:40

Motivation:

Variations in Circuits: Time & Spacial

STATIC			DYNAMIC		
EXTREMELY SLOW	SLOW-CHANGING	FAST-CHANGING	EXTREMELY SLOW	SLOW-CHANGING	FAST-CHANGING
Inter-die process variations Life-time degradation (NBTI, TDB)	Package/Die VDD fluctuations Ambient temperature variations	PLL jitter IR Drop	Intra-die process variations Temperature hot-spots	IR Drop Coupling noise (capacitive and Ldi/dt) Local Clock-jitter (IR drop in clock-tree)	IR Drop Coupling noise (capacitive and Ldi/dt) Local Clock-jitter (IR drop in clock-tree)
Slow	TEMPORAL RATE OF CHANGE	Fast			



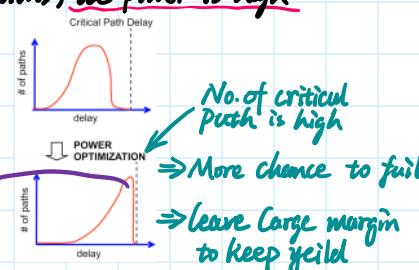
* To provide margin to all of the variations, the power is high

And for circuits that optimized well:

⇒ A lot of circuits has similar delay

⇒ the No. of critical path is high

* More sensitive to Variations



Adaptive Method:

Adaptive design is used to reduce the margin need by variation

It has several methods:

"Always correct" approaches ⇒ with Margin

- Predict point of first failure and add small margin to this point
- Look-up table based techniques
- "Canary" circuits
- In-situ delay detection

"let fail and correct" approaches ⇒ No margin

- Self-calibrating techniques - Circuits tune until point of failure
- Require recovery mechanisms
- Eliminates margins due to global and local variations

Always Correct Approaches:

Look-up Table Based DVFS:

Two ways: Design Based & Post Silicon (use tester)

Set some voltage & Frequency Combinations

Rise voltage first ⇒ rise freq.

Drop freq. first ⇒ Drop voltage } ⇒ keep circuit op. safe.

May use a glitchless Mux for PLL re-clocking

Voltage	Frequency	Power
1.65V	800MHz	900mW
1.3V	600MHz	450mW
0.75V	200MHz	50mW

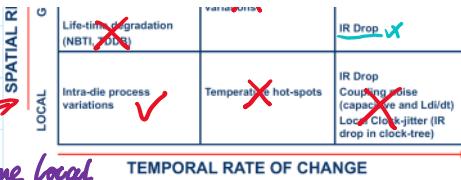
* For Design Based LUT, solves no variations ⇒ Not Adaptive.

STATIC			DYNAMIC		
EXTREMELY SLOW	SLOW-CHANGING	FAST-CHANGING	EXTREMELY SLOW	SLOW-CHANGING	FAST-CHANGING
Inter-die process variations ✓ Life-time degradation (NBTI, TDB) ✗	Package/Die VDD fluctuations ✗ Ambient temperature variations ✗	PLL jitter ✗ IR Drop ✗	Intra-die process ✗ Temperature hot-spots ✗	IR Drop ✗ Coupling noise ✗	IR Drop ✗ Coupling noise ✗
AL					

Variations \Rightarrow Not Adaptive.

* For Post-Silicon Based LUT:

Mostly static & global, but still some local



TEMPORAL RATE OF CHANGE

* Advantages

- Very easy to design and deploy
- Exploits low CPU utilization epochs through DVFS for better energy efficiency

* Disadvantages

- Design time tables do not eliminate margins – only gives frequency / energy trade-off (i.e., dynamic but not adaptive)
- Post Silicon tables: only addresses static variations (global/local process)
- Requires worst-case safety margins for all other variations
- Tester time calibration introduces significant inefficiency because of the different frequency points

* Canary Circuit: Mirror the critical path

Test the delay of canary circuit \Rightarrow Adjust VDD to fit.

* Canary Blocks:

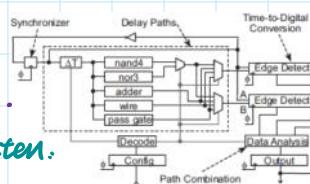
- (1) Use invertors: Ring Osc. Easy but not very accurate
 \Rightarrow May have different critical path – Different Delay Pattern

- (2) Critical Path Monitor:

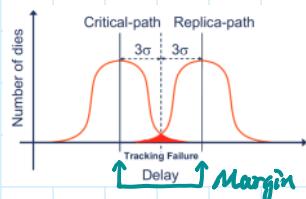
Combination of gates \Rightarrow More accurate.

Dif. Vdd can lead to different delay pattern.

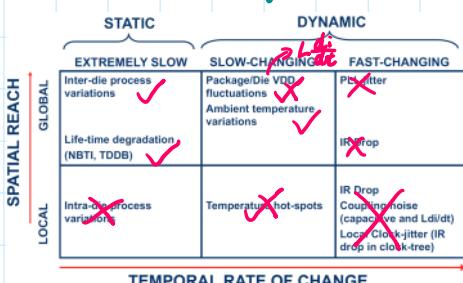
E.g. $pulse_1^{td}$ \xrightarrow{DVFS} with margin $\xrightarrow{\text{A small mistrack gives small margin}}$
 $pulse_2^{td}$ \xrightarrow{Vdd} $\xrightarrow{\text{Higher Efficiency}}$



* Local Variations:



The delay caused by local variation can not be tracked.
Even if the Canary circuit has margin, it still chance to fail.
 \Rightarrow May need localized \Rightarrow Put close critical path



Mainly track global & slow variation



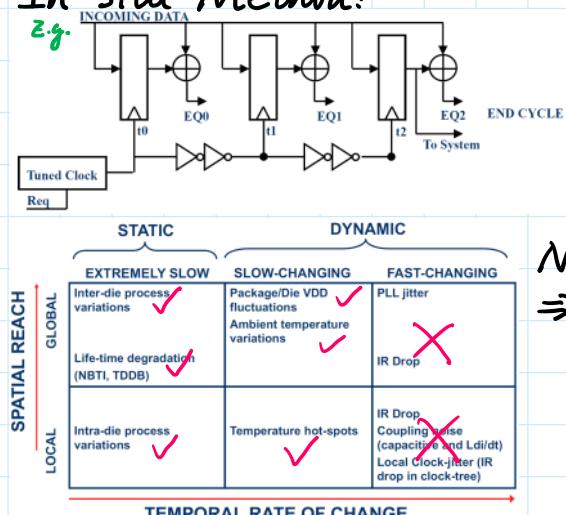
*Advantages

- Relatively easier to design and deploy
- No changes to processor design
- No need for tester time

*Disadvantages

- Adapts to more variations: slow and static, global variations
- Does not address: fast variations, local variations
- In some cases canary circuits can be tuned to reduce margins further but with more testing cost

*In-situ Method:



Detect Data by different FF.
⇒ Can see it is fast/slow/just in time.

Need at least 1 clk to response
⇒ Not work for very fast change.

*Advantages

- Captures Local and Global variations – monitors the actual circuit delay and not a copy
- Eliminates margin for all slow moving variations

*Disadvantages

- More invasive in circuit design – need to change the circuits themselves internally.
- Some delay overhead from monitors
- Requires periodic halting of processor to run worst-case vectors
→ only works if it is worst-case vector

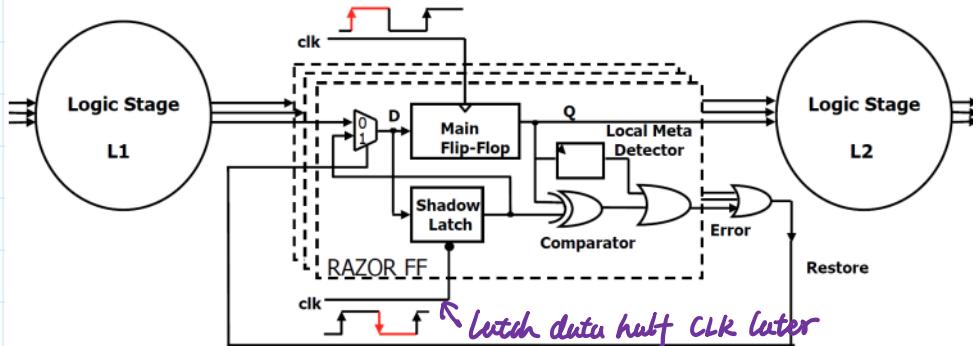
*Let Fail & Correct.

The processor will fail to execute a certain instruction.

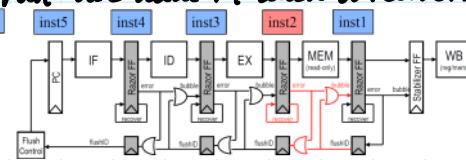
★ Let Fail & Correct.

The circuit will fail & be corrected with a certain error rate
 ⇒ It reduce the margin to zero

★ Razor I.



Compare Main FF with latch: if fail use data in latch to recover.

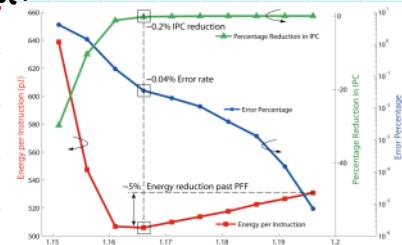
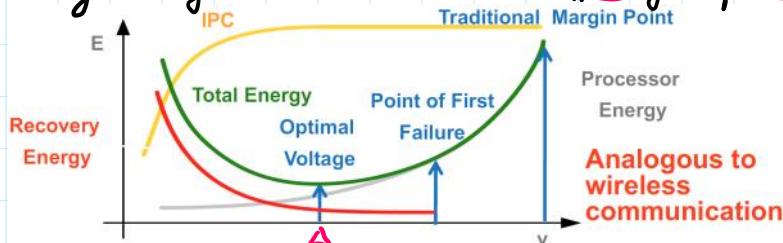


* Problem:

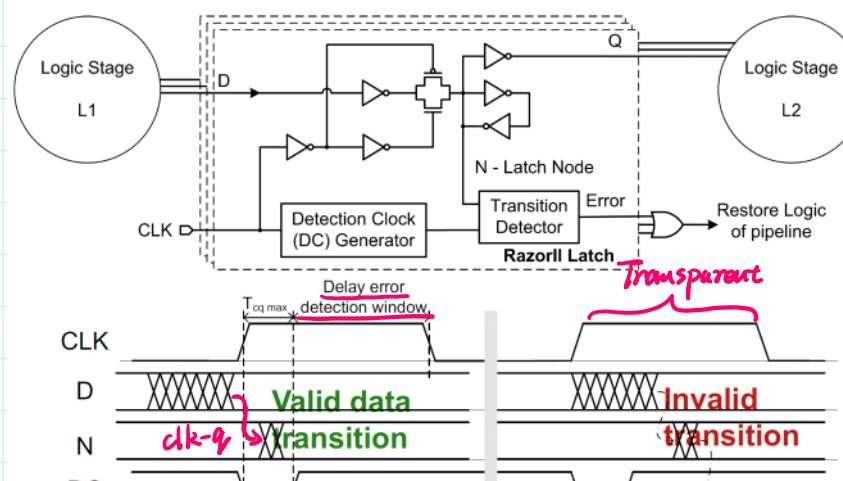
- (1) Hold & Max Detection Time. Data need to be held half CLK it is also the max detection time. otherwise latch will hold the next data.

- (2) Very Large OR gate for error detection: slow & large.

★ By choosing a suitable error rate, the efficiency is optimized.



★ Razor II:



- (1) Give a pulse at clk edge Block the transition detector

- (2) Detector transition in rest of the times ⇒ Yes, means fail (latch transparent)

- (3) Flag and Refetch Instruction



(3) Play and Refresh Instruction

* Conclusion :

		STATIC	DYNAMIC	
		EXTREMELY SLOW	SLOW-CHANGING	FAST-CHANGING
SPATIAL REACH	GLOBAL	Inter-die process variations Life-time degradation (NBBI, TDDB)	Package/Die VDD fluctuations Ambient temperature variations	PLL jitter IR Drop
	LOCAL	Intra-die process variations	Temperature hot-spots	IR Drop Coupling noise (capacitive and Ldi/dt) Local Clock-jitter (IR drop in clock-tree)

TEMPORAL RATE OF CHANGE

if delay within detection range.
⇒ Track all variations.

* Advantages

- Eliminates all margins
- Also captures instruction dependent delay variations

* Disadvantages

- Requires architectural changes in the processor design for error correction mechanism
- Hold time constraint increases overhead of approach with larger timing speculation window

★ Review of Dynamic Power:

See ZZCS427.

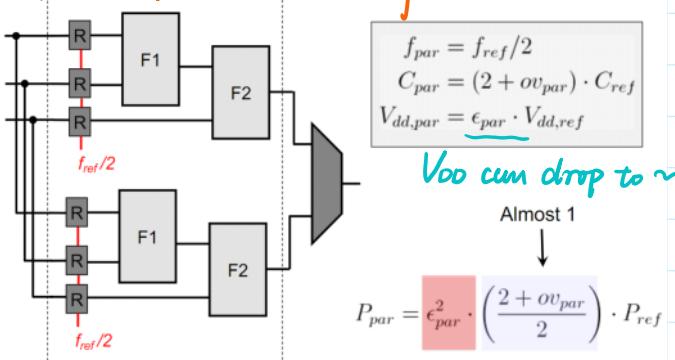
- Reduce C:
 - Sizing, P/N ratio, shorter routes
 - New circuit structures
- Reduce α :
 - Clock gating
 - Bus encoding → already discussed for a diff. sig. $C_g = 2C_L$.
- Scaling V_{dd} : Very effective
 - Parallel/Pipeline tradeoff
 - Low swing signaling/clocking
 - Dynamic Voltage Scaling (DVS)

★ Multiple V_{dd} :

Since Power: $P = \alpha C V_{dd}^2 f$, Reduce V_{dd} can reduce $V_{dd}^2 f$ — 3rd order.

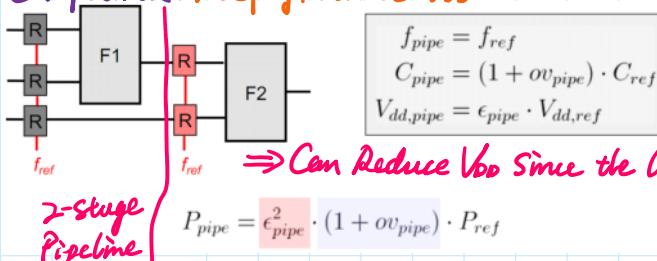
★ Parallel & Pipelined Datapath:

① Parallel: Reduce clk to half



⊖ 1. Large Area
2. Large delay / Data hazard.

② Pipelined: keep f, reduce Vdd



⊖ Large delay / Data hazard.

Area:

Parallel

Pipelined

⇒ Pipeline is a better choice

Delay:

↑

↑

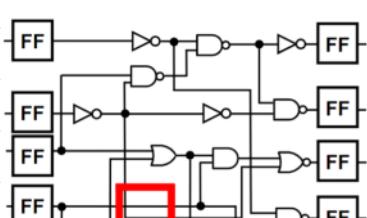
Leakage:

↑

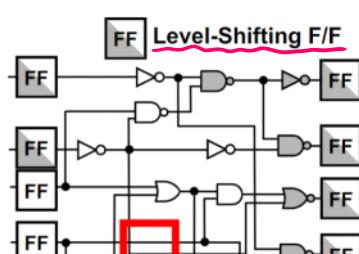
—

★ Multiple Supply in one Block: Low Vdd in non-critical path.

Conventional Design



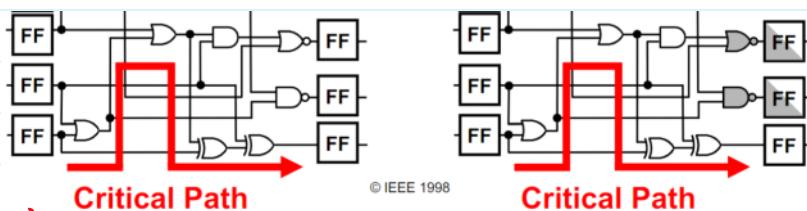
CVS Structure



① Level Shifted by FF

② Only go from $V_{ddH} \rightarrow V_{ddL}$ inside the block (No LS)

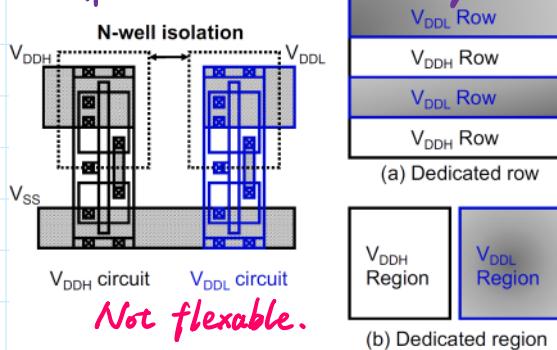
③ Keep Critical Path in V_{ddH}



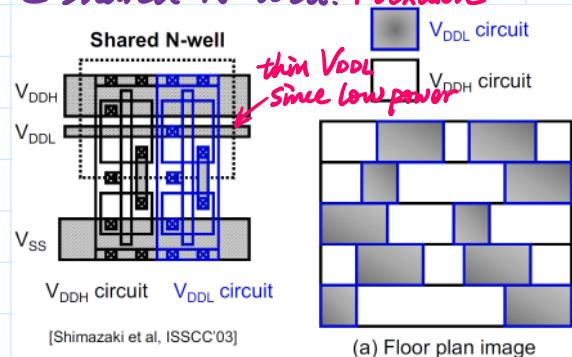
③ Keep Critical Path in V_{DDH}

* Standard Cell in multi-Vdd:

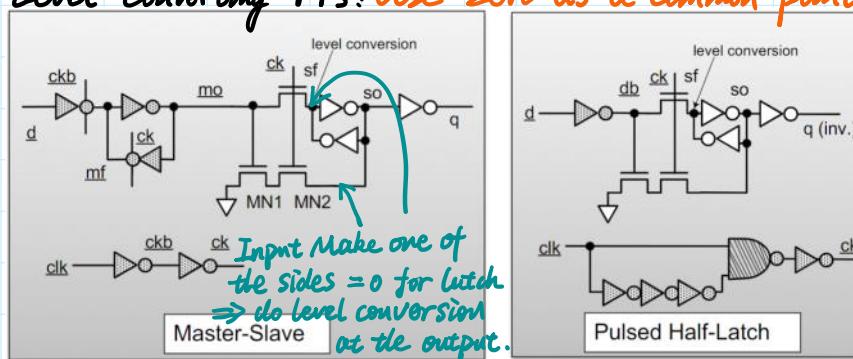
① Different N-well level and single Rail:



② Shared N-well. Flexible



* Level Converting FFs: Use zero as a common point to convert levels.



* Dynamic Power Reduction: DVFS

Reduce the power \propto duty³ by DVFS.

* Find the min power point for scaling: Bo Zhu, Blaauw, DAC, 2009

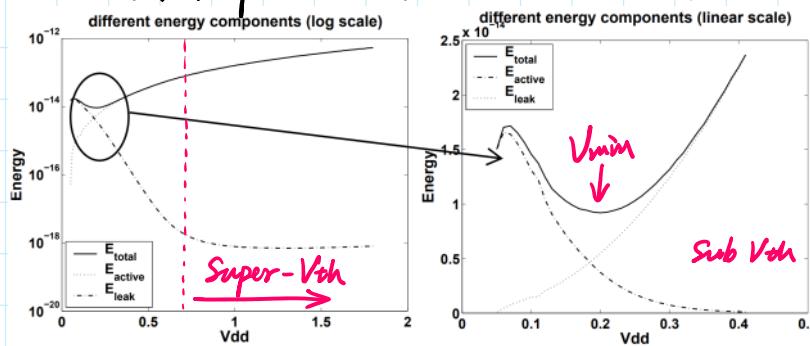
① V_{dd} > V_{th}: super-threshold region: Energy can be saved by smaller V_{dd}

1. Active Energy drops $\propto V_{dd}^2 \downarrow Z = \alpha C V_{dd}^2$

2. Delay rises $\propto V_{dd} \downarrow$

3. I_{leak} keeps the same.

} Static Energy keeps the same
 $Z = V_{dd} \cdot I_{leak} \cdot \text{Time}$



② Sub-V_{th} region: Has a V_{min} for min Energy

1. Active Energy drops $\propto V_{dd}^2 \downarrow$

② Sub-V_{th} region: Has a V_{min} for min Energy

1. Active Energy drops $\propto V_{dd}^2 \downarrow$
 2. Delay rises $\propto \exp(V_{dd}) \uparrow$
 3. I_{leak} keeps the same.
- } Static Energy rises $\propto \exp(V_{dd}) \downarrow$

e.g. For 0.18µm tech: $V_{min} = [1.587 \ln(\eta \cdot n_{eff}) - 2.355] \cdot m \cdot \frac{kT}{q}$ Not related to V_{th}

* There is a Neff = $\frac{n}{\alpha}$; n: stage Number, α : Active Factor

It means the ratio between leaking and active gates.

Neff < 10, No V_{min} — too many switching gates, the Energy can be reduced by keep reducing V_{dd}.

* DVFS works well for small Neff:

- less stages & large α

- Factors affecting V_{min} : $\alpha, n, T, S_S \downarrow$
when $\alpha \uparrow, n \downarrow, T \downarrow, S_S \uparrow (m \downarrow)$, $V_{min} \downarrow$

$$V_{min} = [1.587 \ln(\eta \cdot n_{eff}) - 2.355] \cdot m \cdot \frac{kT}{q}$$

- V_{min} dependency on α and n :

$$E_{leak} = (n * P_{leak, scaled}) * (n * t_{p, scaled})$$

$$E_{act} = n C_S V_{scaled}^2$$

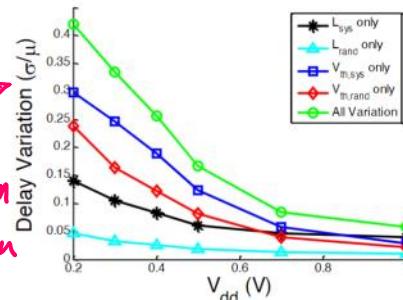
less gates are leaking
less time to leak due to path delay

* Problems in Low Voltage Design:

- Sensitivities \rightarrow exponential
 - Supply voltage
 - Threshold voltage
 - Temperature
- Complicates SRAM design
 - Alternative cell design again helpful

- Makes timing difficult
 - Short paths, long paths
 - Increased clock skew
 - Margins grow

Delay Variability in a 65nm Inverter Chain



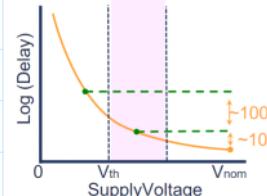
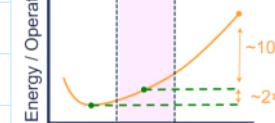
Low V_{dd}
 \Rightarrow Large Variation
Hard for Circuit Design

\Rightarrow Can use Near Threshold Computing (NTC)

* Save a large amount of energy

* But only add a small delay.

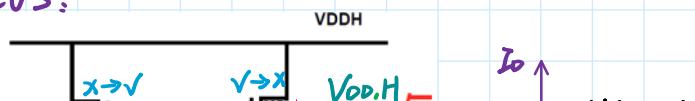
Sub-V_{th} NTC Super-V_{th}



* Level Converters:

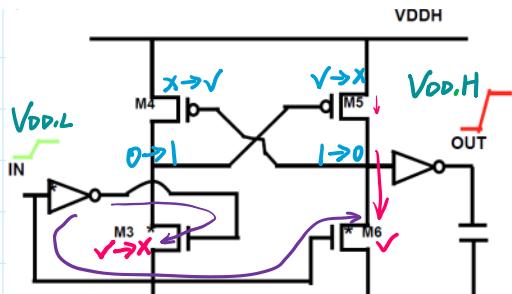
* DCVS & Pass Gate: Use Different Current level.

① DCVS.



② PG: Similar to DCVS

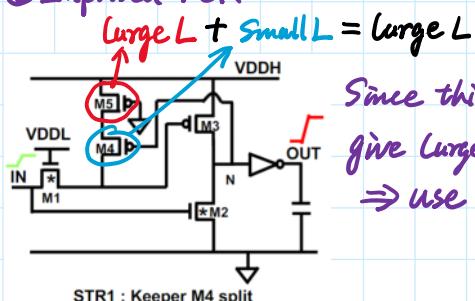




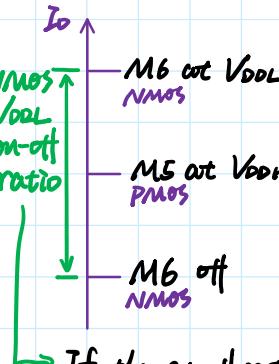
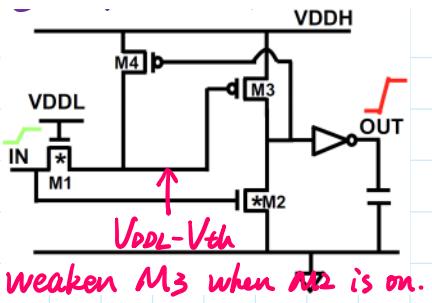
Even M_6 may not fully on at V_{DDL} by have a weak PMOS and feedback it will finally stabilized.

High power but great connection

③ Improved PG:

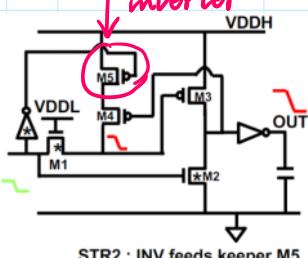


Since this PMOS should be weak \Rightarrow Large L give large load to output.
 \Rightarrow use M_4 & M_5 in series to reduce load.



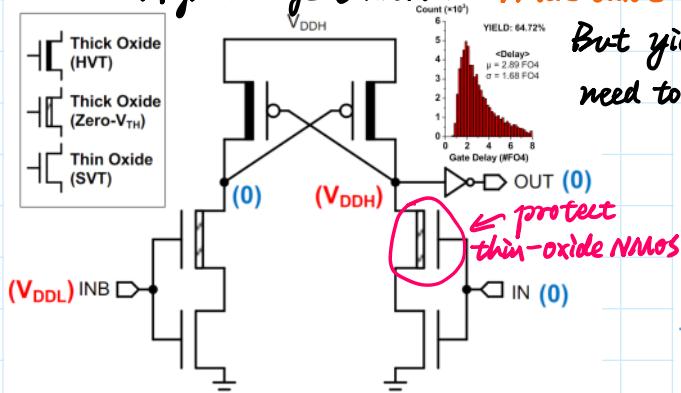
If the on-off ratio is not large (i.e. V_{DDL} is too small) The design margin for PVT variation is also small.

Similar but use inverter

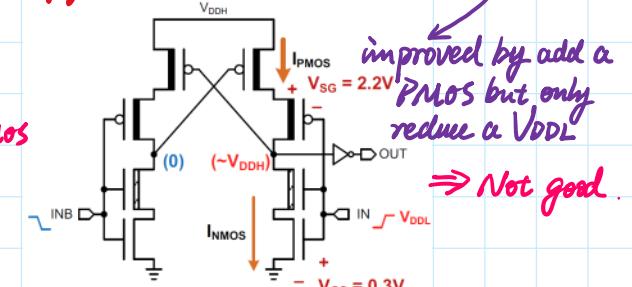


★ High Voltage Convertor

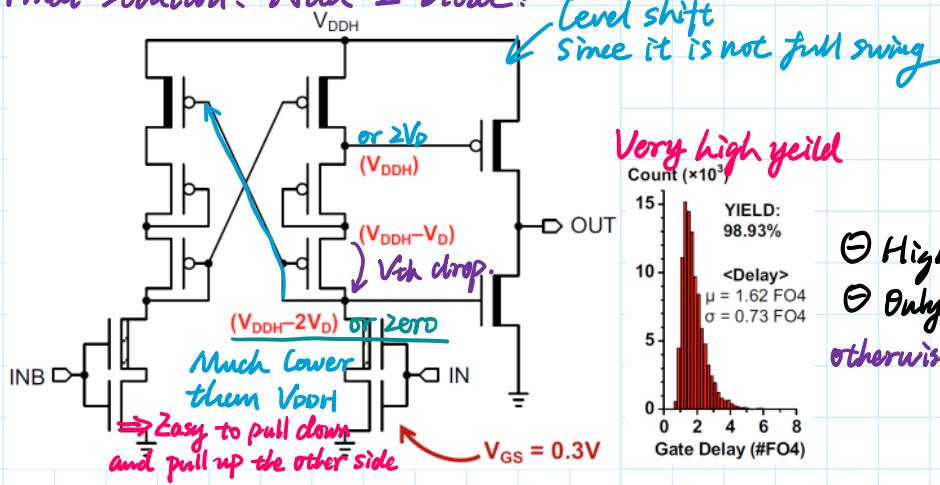
Since for a very low V_{DDL} (say 300mV) the previous converter may not be able to work
 \Rightarrow Use High Voltage Convertor: Thick oxide device is for high voltage.



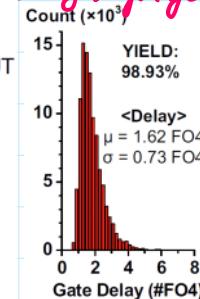
But yield is low since NMOS driven by V_{DDL} need to fight with PMOS driven V_{DDH}



Final Solution: Add 2 Diode:



Very high yield



① High leakage
② Only work for large V_{DDH}
otherwise $V_{DDH} - 2V_d$ too small

L

L7. Energy Recovery (Adiabatic)

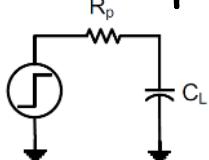
2023年2月18日 23:40

★ Energy in Digital Circuit:

Once there is $I & V$ acrossing resistor (Transistor)

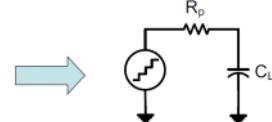
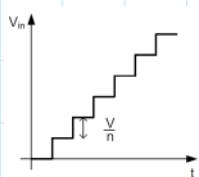
\Rightarrow There is power consumption.

For a step to RC : $E = \frac{1}{2}CV_{\text{step}}^2$



So, if make it to very small steps:

$$\star E = \left[\frac{1}{2}C \left(\frac{V}{n} \right)^2 \right] \cdot n = \frac{1}{2}C \frac{V^2}{n}$$

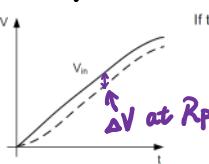


For $n \rightarrow \infty$ (Ramp)
 $P \rightarrow 0$

Since the ΔV between $R_p \rightarrow 0$

In practice: Use slow slew

$$E = I^2 RT$$



If there is a small change in voltage over time:

$$\begin{aligned} I &= \frac{CV}{T} \\ E &= \left(\frac{CV}{T} \right)^2 RT \\ &= \left(\frac{RC}{T} \right) CV_{dd}^2 \\ T \rightarrow \infty & \quad E = 0 \end{aligned}$$

If $T \gg RC$:

$$\Rightarrow E \rightarrow 0$$

E.g. $RC \sim 100ps \Rightarrow T = 1ns$

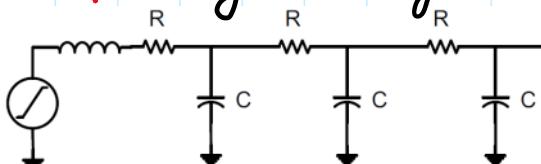
$f_{\text{clk}} = 1 \text{ GHz}$ - Not too slow

★ Zero Power Clock:

After making RC loss to be small, to achieve a zero clock power

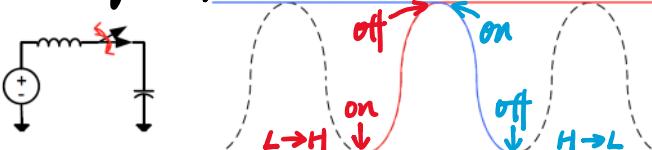
\Rightarrow Need to reuse the energy stored in Caps

★ Slowly Oscillating LC Networks.



if $T \gg RC$,
 $E \rightarrow 0$

For Logic operation: use switch to sample/hold the clock



ideally no power loss for LC.

only power loss at R_p .

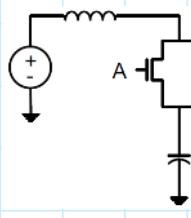
Traditional Clock loss all energy in cap.

★ Adiabatic Logic:

★ Pass-active Logic.

* Adiabatic Logic:

* Pass-gate Logic:



⊕ 1. Low Energy; 2. Auto Clock Gating

⊖ 1. Must switch at right time

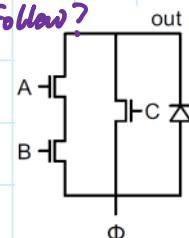
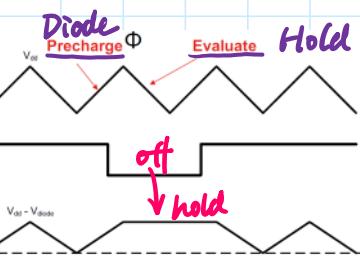
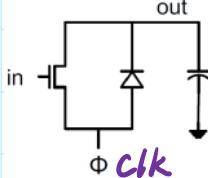
2. $t_d = 1$ clock phase

3. Swing (V_{th} drop)

4. Cap is changing (fc changing)

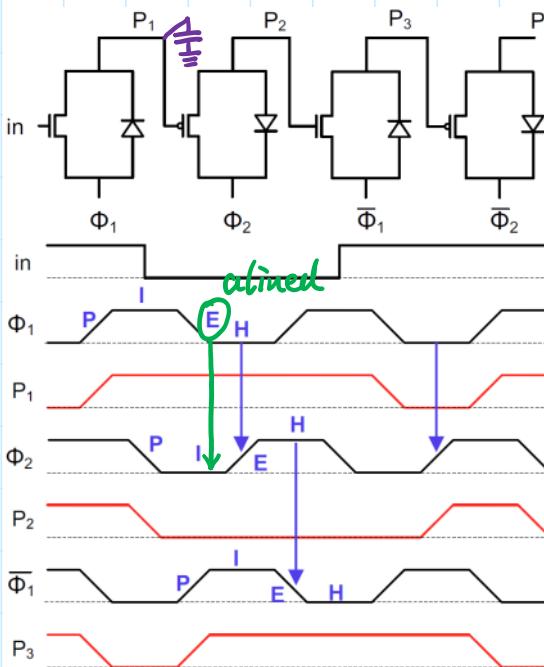
5. Need high Q for low loss

* Dynamic Logic:

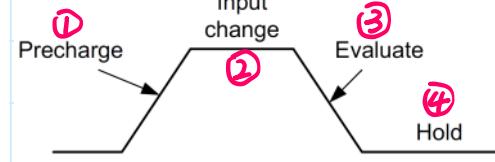


$$Out = \overline{A \cdot B} + C$$

* In cascade:



Need 4-phase Clock: $\Delta\phi = 90^\circ$



* Problems:

1. Diode give a V_D drop

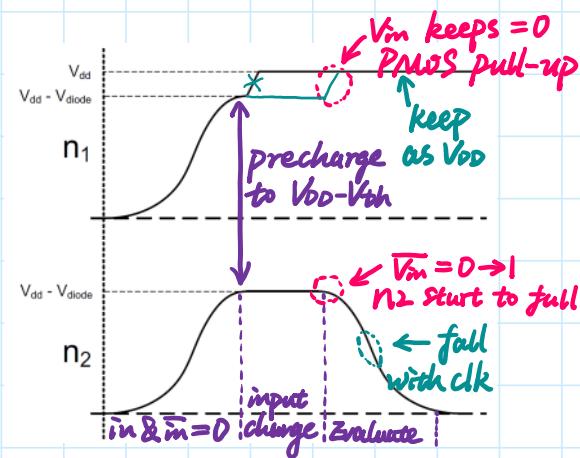
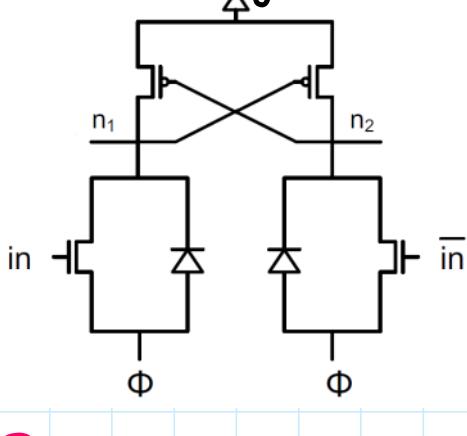
2. "Free Pipelining"

\Rightarrow but no arbitrary connection.

3. Cap is changing (fc changing)

\Rightarrow More loss if $f_{clk} \neq f_c$.

* Dual Rail Dynamic:



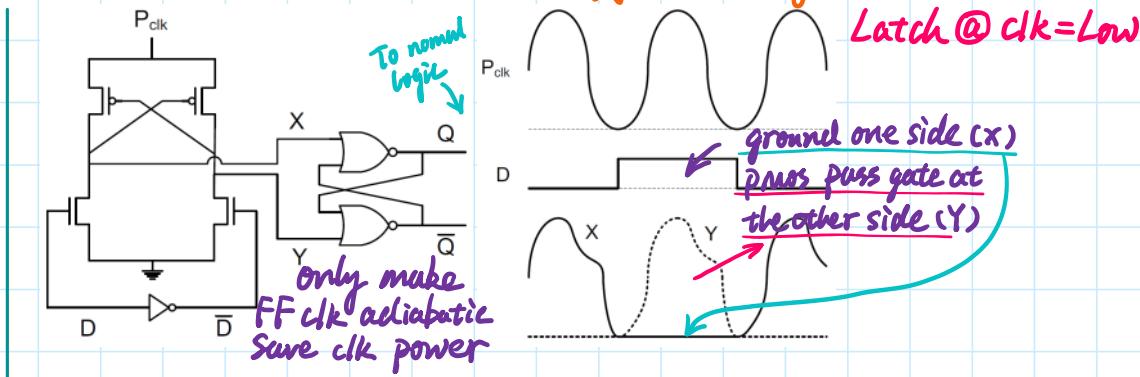
ϕ

ϕ

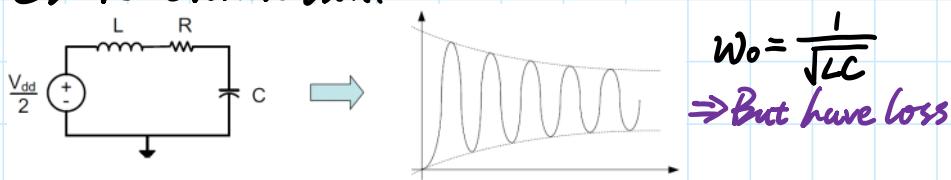
$i_{in} \& i_m = 0$ input charge evaluate

- ⊕ 1. Equal load Cap (No f_C change): see half of C every switch
- 2. Full-rail realized by cross-coupled PMOS
- ⊖ Too many transistors (6T for a Inv)

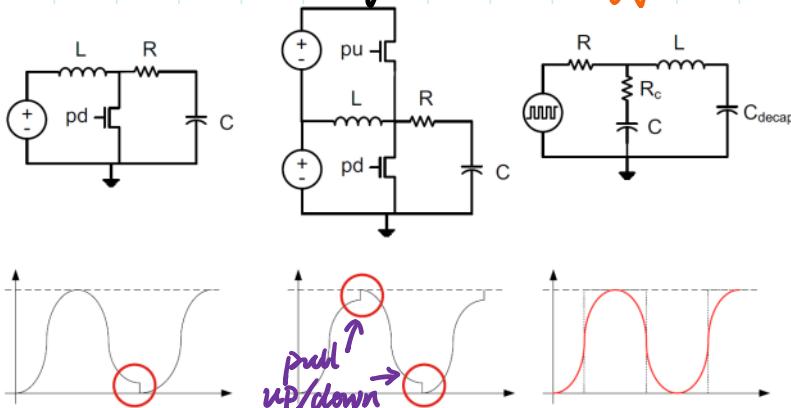
* Adiabatic FF: PMOS Energy Recovering FF (pTZRF)



* Clock Generation:



⇒ Injection Recovery: Add Energy to LC at clock = Low & High



L8. Leakage

2023年3月21日 0:24

★ Leakage Mechanisms.

★ Subthreshold Leakage. Covered in EECS427, L9.

I_{sub} formulae & DIBL:

$$I_{DS} = 2n\mu C_{ox} \frac{W}{L} \left(\frac{kT}{q} \right)^2 e^{\frac{V_{GS}-V_{TH}}{nkT/q}} \left(1 - e^{-\frac{V_{DS}}{kT/q}} \right) = I_s e^{\frac{V_{GS}-V_{TH}}{nkT/q}} \left(1 - e^{-\frac{V_{DS}}{kT/q}} \right) - w/o DIBL$$

$$P_{leak} = (I_0 \frac{W}{W_0}) 10^{\frac{-V_{TH}}{S}} (V_{DD} 10^{\frac{-V_{DD}}{S}}) - w/ DIBL$$

⇒ Subthreshold power is a strong function of V_{DD} .

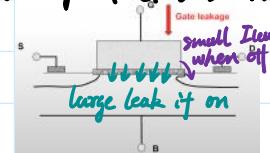
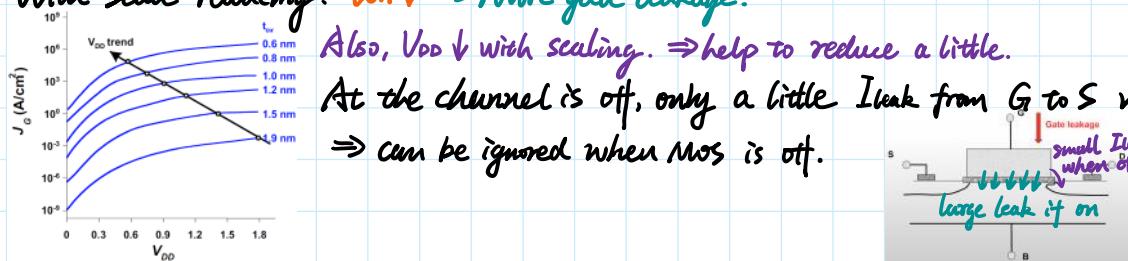
★ Gate-Induced Drain Leakage: GIDL

* PMOS has small gate leakage — hole tunneling is less.

With scale reducing: $t_{ox} \downarrow \Rightarrow$ More gate leakage.

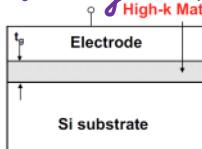
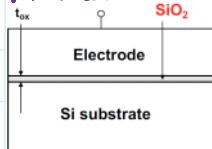
Also, $V_{DD} \downarrow$ with scaling. ⇒ help to reduce a little.

At the channel is off, only a little leak from G to S via overlap.
⇒ can be ignored when MOS is off.



★ High-k Gate Oxide:

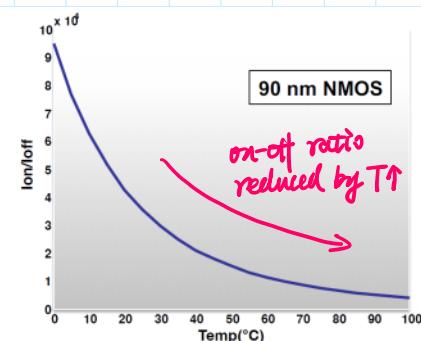
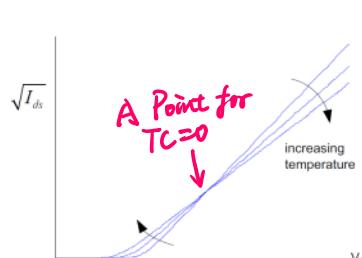
For advanced node, t_{ox} only \sim atoms ⇒ large leakage ⇒ use thick oxide but high-k



	High-k vs. SiO ₂	Benefits
Gate capacitance	60% greater	Faster transistors
Gate dielectric leakage	> 100% reduction	Lower power

★ Temperature Effect:

- Increasing temperature
 - Reduces mobility
 - Reduces V_{TH}
- I_{ON} decreases with temperature $T \downarrow$
- I_{OFF} increases with temperature $T \uparrow$



Since different region; different effects give different affects.

★ Other Leakage: Diode



- Electron-hole pair generation in depletion region of reverse-biased diodes
- Diffusion of minority carriers through junction
- For sub-50nm technologies with highly-doped pn junctions, tunneling through narrow depletion region becomes an issue

band-to-band Tunneling (BTBT)

Strong function of Temp.
 $HT \Rightarrow$ large leakage

- For sub-50nm technologies with highly-doped pn junctions, tunneling through narrow depletion region becomes an issue

*Strong function of temp.
HT \Rightarrow large leakage*

★ Leakage Reduction:

★ Reduce V_{th} : $\times 10$ leakage $\Rightarrow 20\% \sim 30\%$ delay

	Low-Vt; 0.9V	High-Vt; 0.9V	Low-Vt; 1.8V	High-Vt; 1.8V		
Leakage (norm)	1	10x	0.06	1	10x	0.07
Delay (norm)	1	30%	1.30	1	20%	1.20

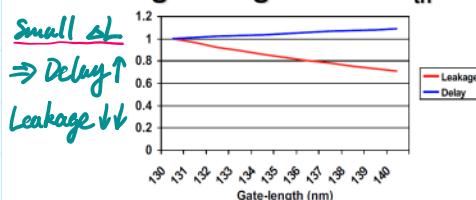
$$T_d \propto \frac{V_{th}}{(V_{DD} - V_{th})}$$

But for low V_{DD} , the advantage from HVT is less since delay increase more.

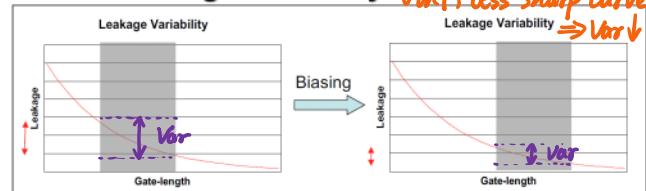
★ Lower Temperature: Cooling

★ Increase Gate Length: $L \uparrow, V_{th} \uparrow$

- Reducing leakage due to V_{th} roll-off



- Reduce leakage variability



★ Decrease V_{DD} : DIBL $I_{sub} \propto \exp(V_{DD})$

★ Low-Power Techniques:

★ Stacking Transistors:

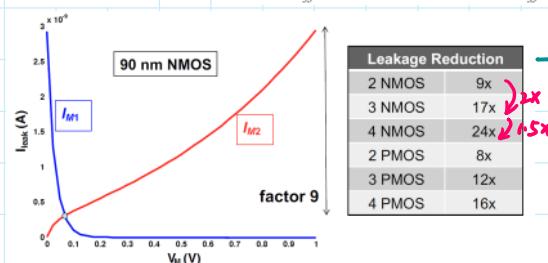
Same I_{sub} for N_1, N_2

$$I_{sub} = I_{off} 10^{\frac{\eta(V_x - V_{DD})}{S}} = I_{off} 10^{\frac{V_{gs} - V_{DD} + \eta((V_{DD} - V_x) - k_y V_x)}{S}}$$

$\star V_x = \frac{\eta V_{DD}}{1 + 2\eta + k_y} \approx 100 \text{ mV}$

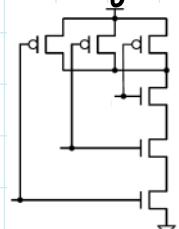
$$I_{sub} = I_{off} 10^{\frac{-\eta V_{DD}}{S(1 + 2\eta + k_y)}} \approx I_{off} 10^{\frac{-\eta V_{DD}}{S}}$$

\Rightarrow DIBL by V_x
DIBL $\frac{V_{gs} - V_{DD} + \eta((V_{DD} - V_x) - k_y V_x)}{S}$ Body effect.



\rightarrow but the advantages from stacking is not significant when $n > 2$.

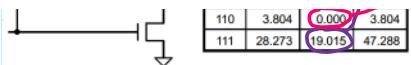
Stacking Order Dependence of Gate leakage:



State	I_{sub}	I_{gate}	I_{total}
000	0.382	0.000	0.382
001	0.709	6.339	7.048
010	0.709	1.275	1.275
011	5.626	12.677	18.303
100	0.676	0.000	0.676
101	3.804	6.339	10.141
110	3.804	0.000	3.804
111	28.273	19.015	47.288

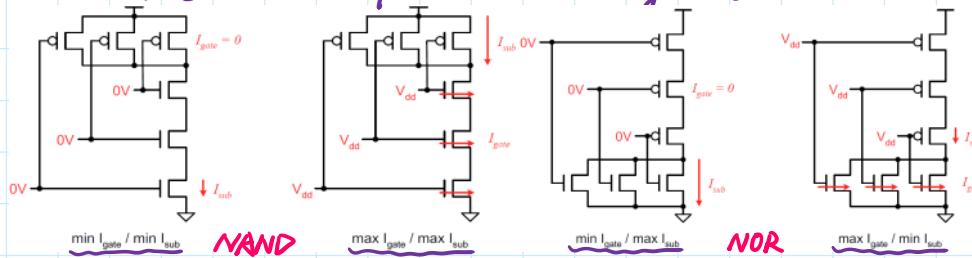
Gate leakage is significant when channel is on.
e.g. I_{gate} reduces with no. of on transistors.

$\rightarrow I_{gate} \approx 0$ when off or $V_D = V_g$.

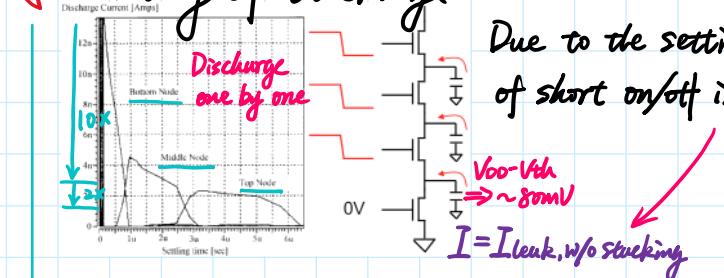


$I_{gate} \sim V_{DD} \ln(V_{DD}/V_t)$

This can be used to improve the state assignment.

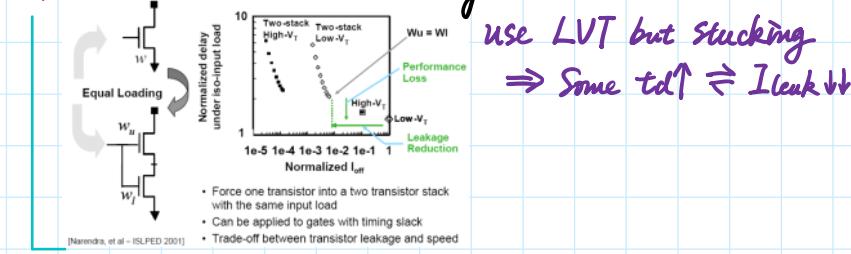


★ Setting of Stacking:

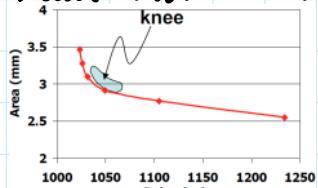


Due to the setting time, the advantage of short on/off is less.

★ Forced Transistor Stacking:



★ Dual Vth Circuit:



After the knee, by the same td reduction the Area & Ileak increase exponentially.

=> 1. Optimize to knee; 2. change some gate with HVT.

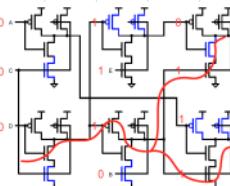
Not significantly rise td, but reduce leakage.

Mainly apply for non-critical path. And it is the only way reducing active leakage.

⊕ Simple to implement ⊖ only 2x reduction others only for sleep leakage.

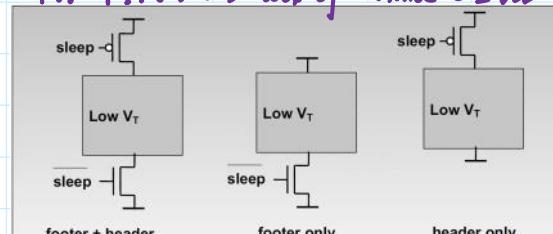
★ Combination of Vth & State Assignment:

=> Assign HVT for leaking state.

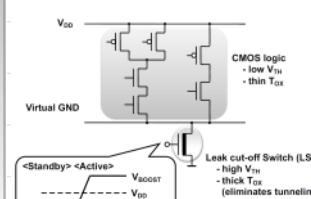


★ Power Gating: MTCMOS.

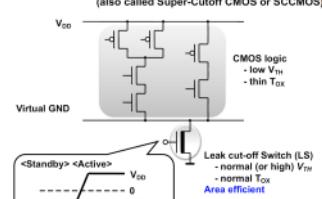
* For P/N MOS all of $V_{mid} \approx \frac{1}{2}V_{DD}$ (Not Stacking!)

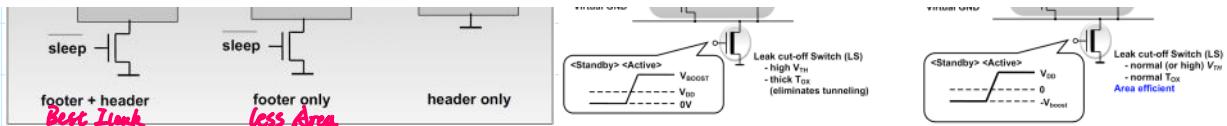


Other option: Boosted-Gate MOS (BGMOS)

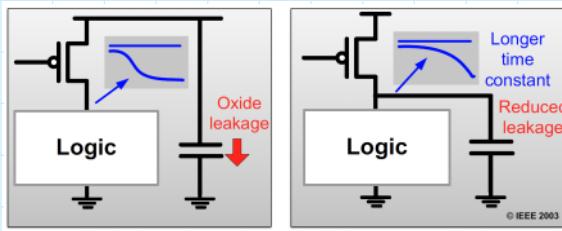
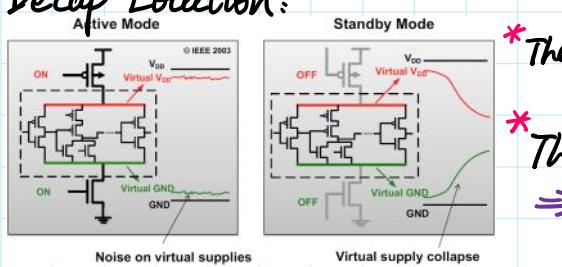


Other Option: Boosted-Sleep MOS
(also called Super-Cutoff CMOS or SCCMOS)





★ Decap Location:



* The power gating creates virtual supply

* The header/footer can be sized as ~10% total
⇒ Not all the gates switching in same time

on supply rails on virtual rails

Performance

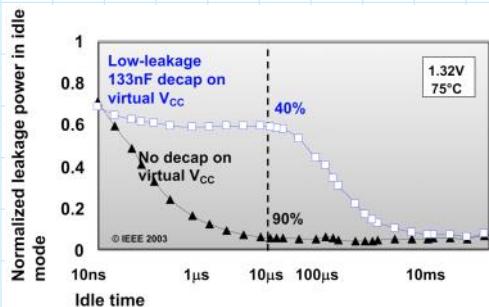


Convergence time



Oxide leakage savings

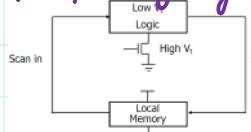
Also, for Cap on virtual rail:
the rush open may induce a supply noise
⇒ Affect other blocks. (turn-on gradually).



* The location of Decap depends on application
⇒ On virtual rail is not effective for fast on/off.

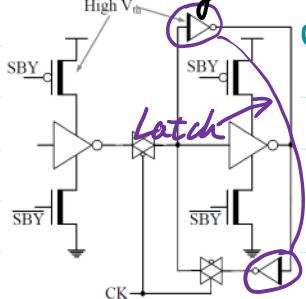
★ Preserving State:

For power gating, it requires retaining the state at sleep node.



One easy way is scan out all the state
⇒ But takes a long time.

* A usual way is state retaining FFs.



① When Sleep:

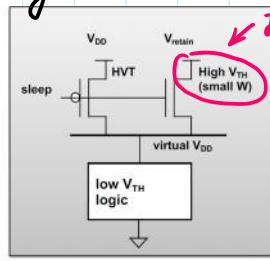
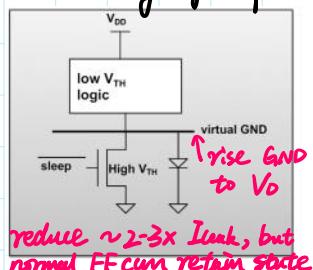
Use HVT Latch to reduce I_{sub}

② When work:

Just a normal latch.

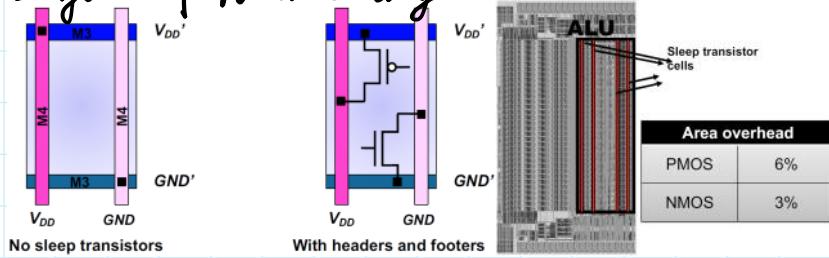
But it requires both footer & header
⇒ otherwise no leakage reduction

* other ways for preventing state:

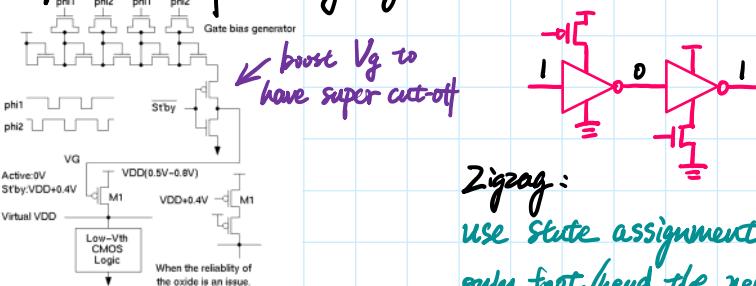




★ Layout of Power Gating



★ Super Cut-off & Zigzag MT-CMOS

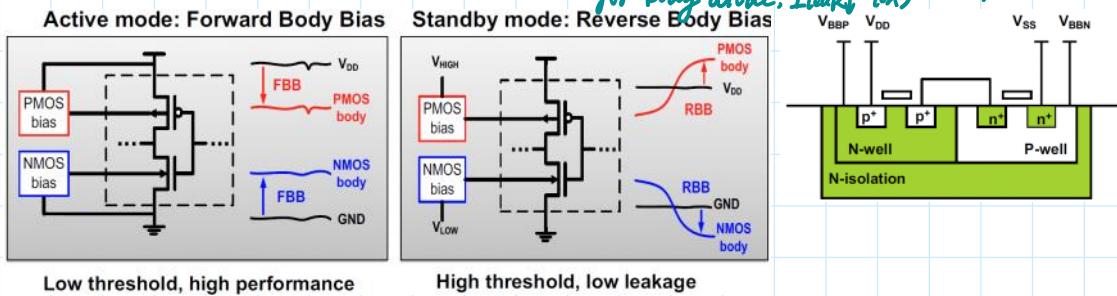


boost V_g to have super cut-off

Zigzag: use state assignment of logic only foot/head the needed one.

★ Dynamic Body Biasing:

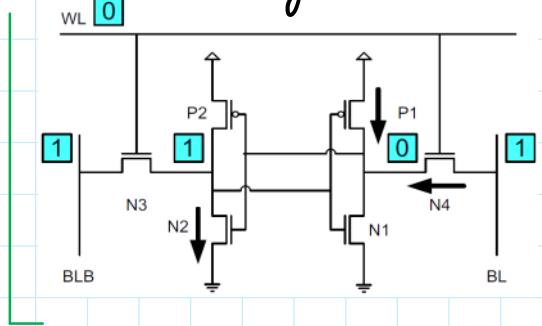
No delay reduction, but need triple-well, limited V_{th} adjustment, need charge body. ($\sim \alpha/V$, since $K=\alpha$ and $V_0=\alpha$) for body diode, leaky (~ αx) (power↑)



★ Conclusion:

	Transistor Stacking	Power Gating	Dynamic Body Biasing	Dual V_t
Pros	<ul style="list-style-type: none"> Conventional technology No performance impact 	<ul style="list-style-type: none"> Conventional technology Conceptually simple Most effective 	<ul style="list-style-type: none"> Reuse of standard designs No performance impact 	<ul style="list-style-type: none"> Works at run time No / minor performance loss No layout change extremely common
Cons	<ul style="list-style-type: none"> Limited impact Special registers 	<ul style="list-style-type: none"> Performance impact of serial transistor Changes in design flow 	<ul style="list-style-type: none"> Triple well Slow activation Does not fare well with technology scaling 	<ul style="list-style-type: none"> Limited leakage reduction More mask layer(s)
Potential Savings	• 1 - 2x <i>Very Fast</i>	• 2 - 1000x <i>Fast</i>	• 2 - 20x <i>slow</i>	• 2 - 5x <i>also works at run time</i>

★ SRAM Leakage:



1. Make V_{DD} Low / V_{SS} High for latch. (Drowsy SRAM)
2. Floating BL $\Rightarrow \frac{V_{DD}}{2}$ by leakage \Rightarrow Access Transistor $I_{leak} \downarrow$
3. Body Biasing
4. Underdrive WL. ($\sim -100mV$)

L9. Power Supply

2023年3月21日 0:24

★ Power Supply Model:

For power supply of the chip.

* Important to keep it low impedance for all frequency
⇒ low R & L ; more parallel Cap.

e.g. Power = P . Max ripple is $r \cdot V_{dd}$

$$\Rightarrow Z < r \cdot (V_{dd})^2 / P. \text{ For } P=100\text{mW}, r=0.1, Z < 1\text{m}\Omega$$

★ Voltage Drop: IR & $L \frac{dI}{dt}$

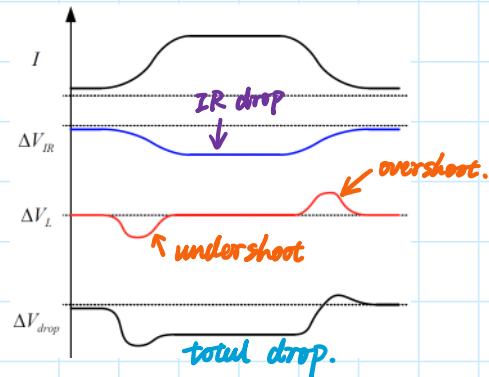
* IR drop:

$$\Delta V = IR; P = IV \Rightarrow \text{if } P \uparrow \Rightarrow V \downarrow \Rightarrow I \uparrow$$

⇒ Make ΔV larger

* $L \frac{dI}{dt}$ drop.

$$\text{if freq. } \uparrow; dI \uparrow, dt \downarrow \Rightarrow L \frac{dI}{dt} \uparrow \uparrow$$



★ Power Supply Noise:

① Sustained drop: $t_d \propto \frac{CV}{I} \propto \frac{1}{(V_{dd} - V_{ss} - V_t)}$

For ΔV at V_{dd} ⇒ Circuit delay $t_d \uparrow$

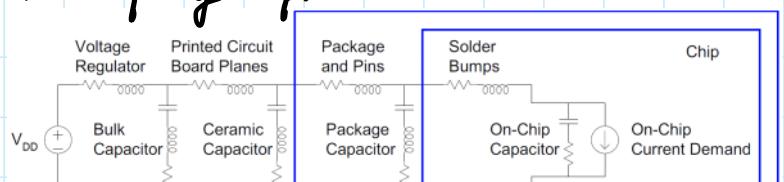
② AC Noise: Functional / Delay

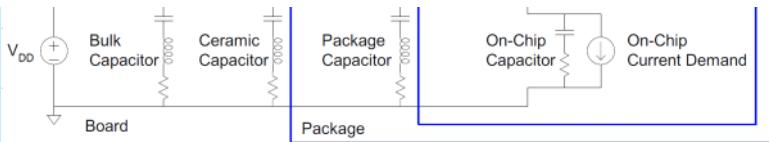
may have ringing.

③ Over/under shoot: Exacerbated TDDB

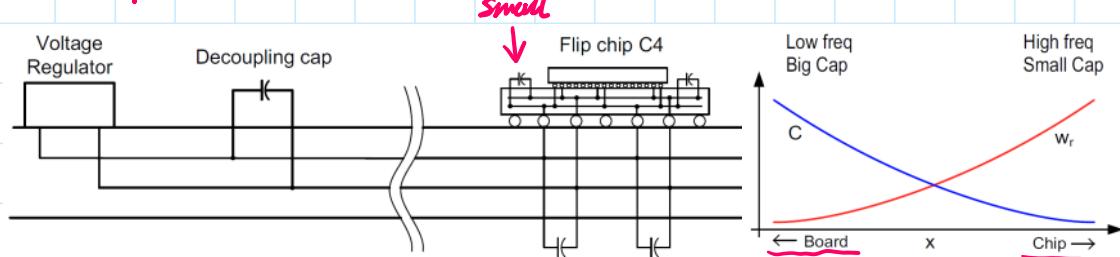
Also limited $V_{dd, max}$ to be small (margin for over/under shoot)

★ Decoupling Caps:





There are different decaps \Rightarrow Large C comes with large R_{ZSR} & L_{ZSL} .
Thus, they have different W_0 .

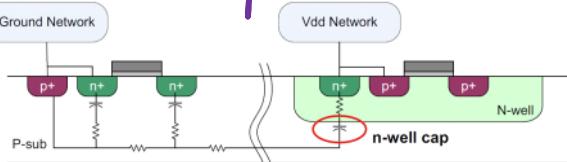


- * 1. usually the cap close to the chip is small but high frequency
- * 2. The decaps usually placed at where self-inductance is dominant.
i.e. $L_{decap} \gg L_{wire}$, otherwise it has no help to the circuit.

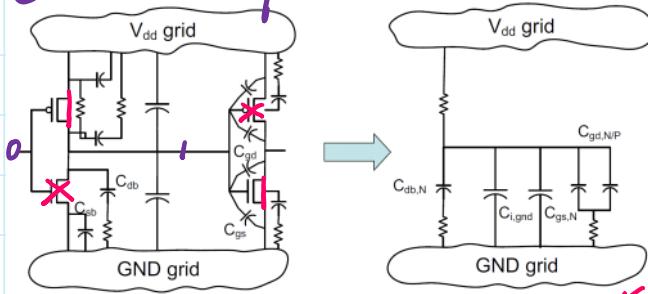
★ Decap in Circuit.

★ Implicit:

① N-well Cap:



② Parasitic Caps..



Some of the cap between V_{DD} & V_{SS}
can be considered as decap.

$$C_{eff} = \frac{1}{2}(C_{db,N} + C_{gs,N} + C_{db,P} + C_{gs,P}) + \frac{1}{2}(C_{i,gnd} + C_{i,Vdd}) + (C_{gd,N} + C_{gd,P})$$

half chance for $V_{in} = V_{SS}/V_{DD}$

always there.

For the cap in circuit: non-switching caps are decaps.

$$P = f C_{switch} \cdot V_{dd} \Rightarrow C_{switch} = P / f V_{dd}$$

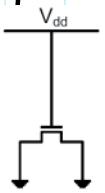
Due to the activity rate: $C_{switch} = S C_{total}$.

$$\text{I-S} \quad P$$

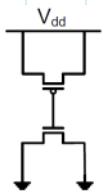
Due to the activity rate: $C_{switch} = sC_{total}$.

$$\star C_{decap} = (1-s)C_{total} = \frac{1-s}{s} \frac{P}{fV_{dd}}$$

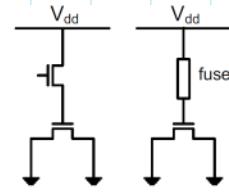
★ Explicit:



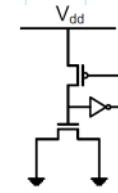
- Reliability: Thin oxide gate can short. Yield problem
- High gate leakage



- Short only if double failure
- Gate voltage is half, less sensitive to gate failure
- Less cap

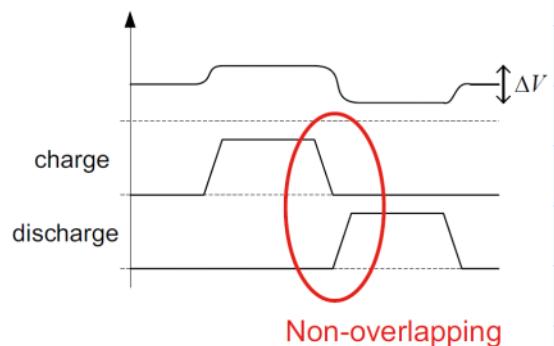
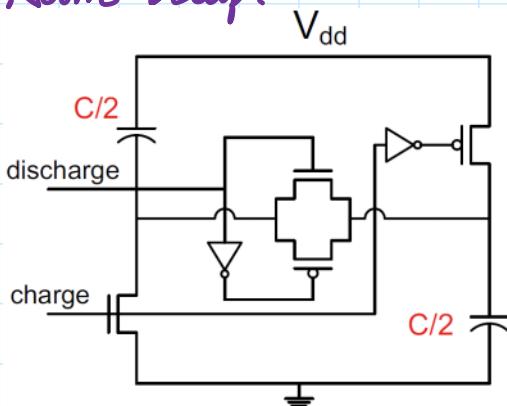


- Turns off the cap in case of oxide failure
- Less gate leakage
- Unusable if the transistor fails



- Variation of previous

Active Decap:



Charge in parallel; discharge in series: like charge pump.

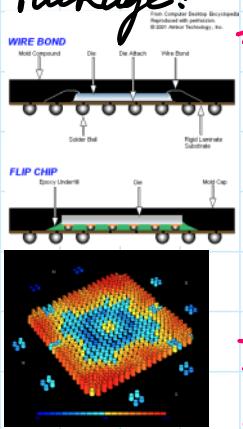
$(\frac{1}{2}V_{dd} \rightarrow V_{dd})$

$(2V_{dd} \rightarrow V_{dd})$

\Rightarrow Provide more charge ($Q = CV$)

Normal Cap: $Q = \sigma V C$ \Rightarrow Active: $Q = \frac{1}{2}CV_{dd} + \frac{C}{2}\Delta V$

★ Package:



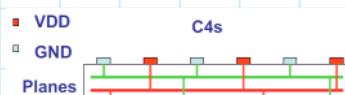
* Wire bond:

Cheap, but large L , small No. of pins
 $(\sim nH)$ (~ 200)

* Flip Chip:

Less L , expensive, more pins.
 $(\sim 0.1nH)$ (~ 1000)

There is package induced variation due to the pin location.



★ Heat Dissipation:

- ΔT : temperature rise on chip

heat sink usually.

★ Heat Dissipation:

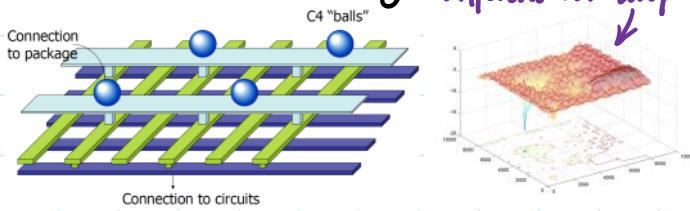
$$\star \Delta T = \theta_{ja} \cdot P$$

- ΔT : temperature rise on chip
- θ_{ja} : thermal resistance of chip junction to ambient *can be series/parallel*.
- P: power dissipation on chip
- Ohm's Law for heat

heat sink usually.

e.g. $\Delta T = 45^\circ\text{C}$, $\theta_{pack} = 1^\circ\text{C/W}$
 $\theta_{chip} = 0.5^\circ\text{C/W}$
 $\Rightarrow P = 45 / (4 + 0.5) = 30\text{W}$.

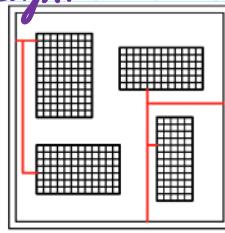
★ Chip Power Delivery:



Power grid design:

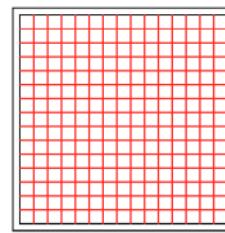
Tree

- Local grids
- Simple
- Metal utilization is efficient
- Single wire failure

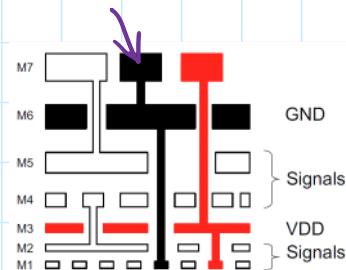


Grid

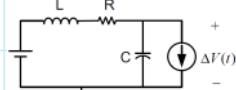
- High end processor
- Lots of metal
- Lower resistance
- Lines up with C4
- Less J, better electromigration
- Spatial variation is lower
- Highly redundant



Can use top layer metal for power



★ Power Grid Model:



$$\text{For the change of } I: \Delta V(t) = IR + I \sqrt{\frac{L}{C}} e^{\frac{-Rt}{2L}} \cos(\omega_r t - \theta)$$

\Rightarrow IR drop + resonate LC

$$\Delta V_L \propto I \sqrt{\frac{L}{C}} \Rightarrow \text{decap reduce} \propto \sqrt{C}$$

I_d reduce $\propto I_d$.
 R gives IR drop, but helps clamp the resonance.

Also, the LC may resonate.

For previous, $\omega_0 > \omega_{clk}$, the harmonic may induce resonance

Now, $\omega_{clk} > \omega_0$; harmonic is not a problem but the on-off of the blocks. *loop of on/off instruction.*

★ Solution:

Reduce R: Use more metal. Tree \rightarrow Grid \rightarrow Plane

$$\Delta V_{IR} \downarrow \quad \text{but less damping} \quad Q \uparrow$$

Reduce L: Thin package, bondwire, flip-chip. More pads

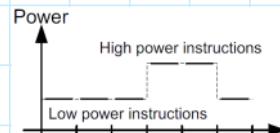
$$\Delta V_L \downarrow \quad \omega_r \uparrow \quad Q \downarrow \quad \text{difficult to control}$$

Increase C: Decoupling capacitance, but only \sqrt{C}

$$\Delta V_{IR} \downarrow \quad \Delta V_L \downarrow \quad Q \downarrow \quad \text{Area} \uparrow$$

Decrease I: Turn modules on slowly, larger reset latency

$$\Delta I \downarrow \quad \Delta V_{IR} \downarrow \quad \Delta V_L \downarrow$$

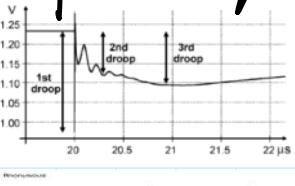


$$\omega_r = \frac{k \omega_{clk}}{l} \quad k \rightarrow \text{harmonic}$$

loop of on/off instruction.

$$\Delta I \downarrow \quad \Delta V_{IR} \downarrow \quad \Delta V_L \downarrow$$

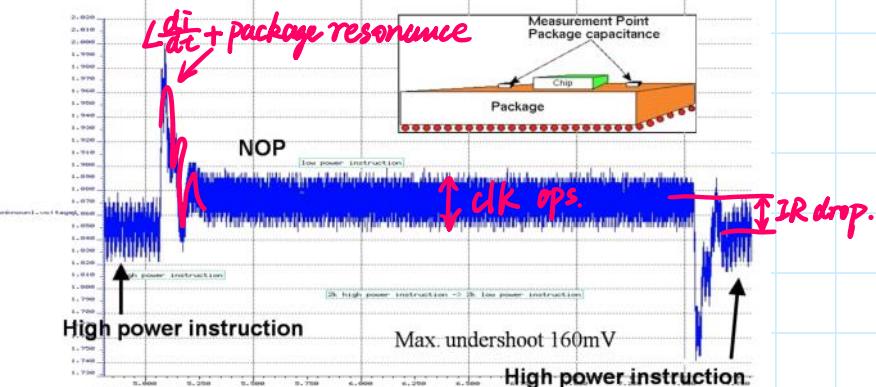
★ Different Frequency Drops:



1st droop: decoup (on-chip)

2nd droop: package cap.

3rd droop: PCB cap.

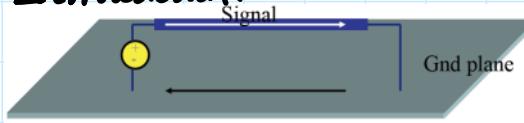


L10. Inductance

2023年3月27日 0:40



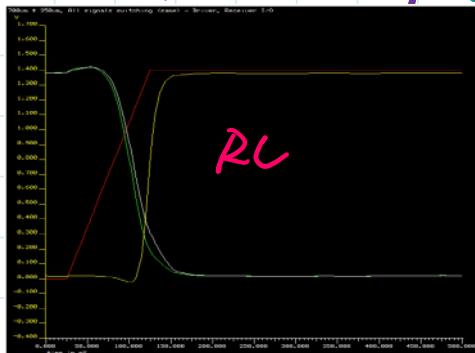
Introduction:



for inductance:

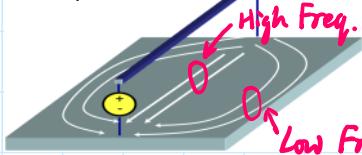
(1) Self: L & Mutual: M

(2) Must have return loops: large loop size \Rightarrow flux $\uparrow \Rightarrow L \uparrow$.



\Rightarrow Important for high speed IC chip.

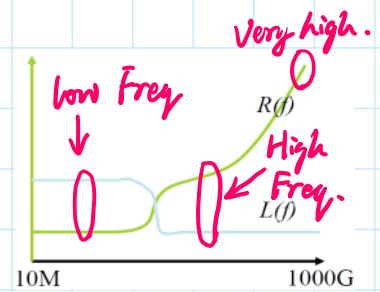
Frequency Dependent:



High Freq. For normal case: $Z = R + j\omega L$

but in chip with high freq.,

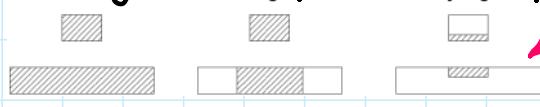
the L also change with f .



(1) Low Freq.: The current is wide spread: $L \uparrow, R \downarrow$

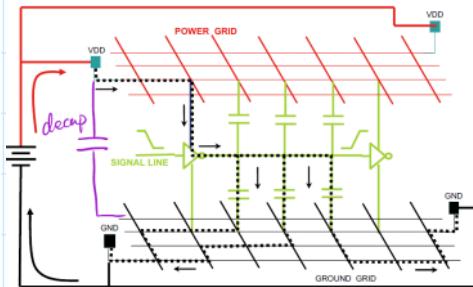
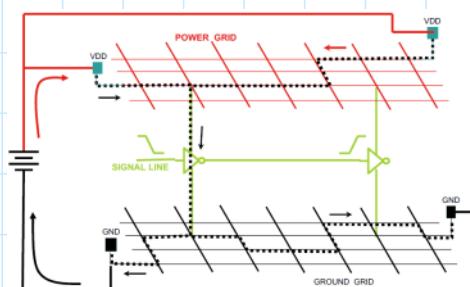
(2) High Freq.: The current is limited in a narrow path under the wire: $L \downarrow R \uparrow$

(3) Very High Freq.: The skin effect makes R even larger.



Current only flow at
at the surface.

Current in the Circuit:

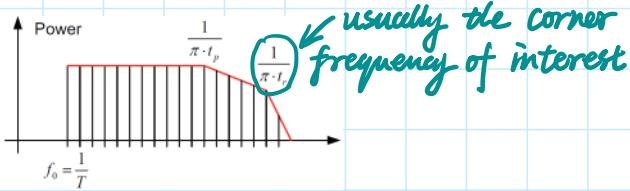
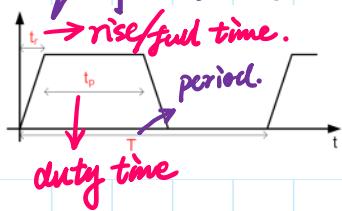


Flow through power grid, also with decaps.

\Rightarrow More decaps helps the current response.

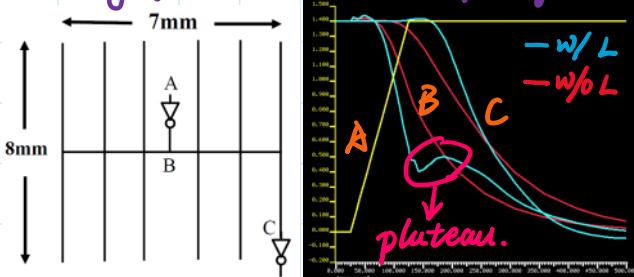
★ Inductance in Circuit:

Freq. of interest:



★ Switching with RLC:

usually for the clock net/large bus:



★ With L: at C

just switch but more delay

⊖ More noise

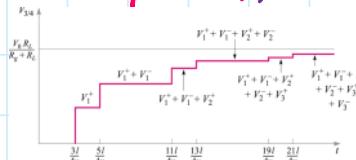
⊕ Less Ishort, fast switch.

For the circuit with L modelled: plateau.

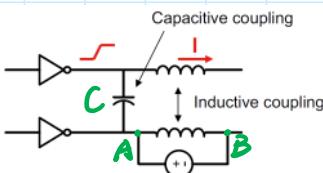
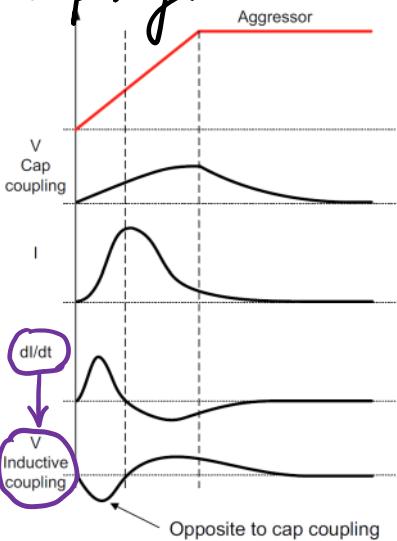
⇒ Current in L with step response will rise gradually.

like open first (fast switch to a plateau); then drop ($I \uparrow$)

Also can be explained by
DC response of a Tx line:



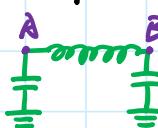
★ Coupling:



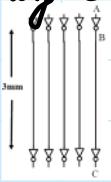
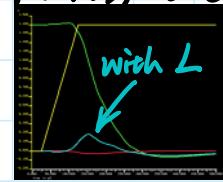
Two parallel wire may have inductive coupling:

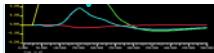
This coupling only provide ΔV ,

but the reference V_0 should be found.

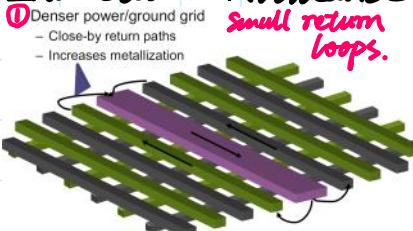
e.g.  $C_1 = C_2 \Rightarrow V_A = -\frac{1}{2}\Delta V; V_B = \frac{1}{2}\Delta V$
 $C_1 = C, C_2 = 0 \Rightarrow V_A = -\Delta V; V_B \text{ not change}$
only C_1 provide current.

For bus, there may be crosstalks





★ Inductance Avoidance.



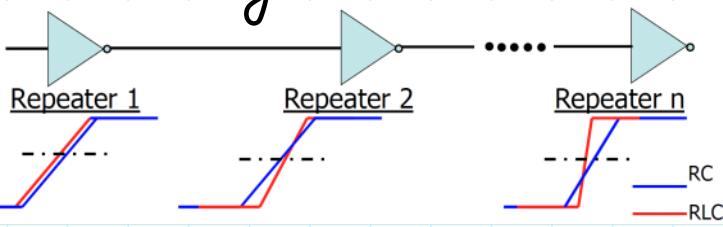
③ Reduction of inductance at low and high frequency
• More metallization

④ Split wide conductor into multiple wires with shields in between. [Massoud '98]

⑤ Reduce the overlap between adjacent wires
- Capacitive and inductive coupling reduced

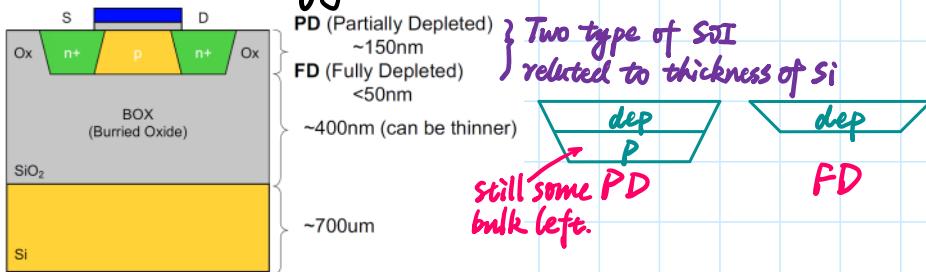
⑥ Nets routed to create opposite loops [Zhong '00]
- Magnetic flux from adjacent nets cancels out

★ Optimal Delay:



The inductance can be opt. to achieve a low delay:
by fast slow

★ SOI Technology:

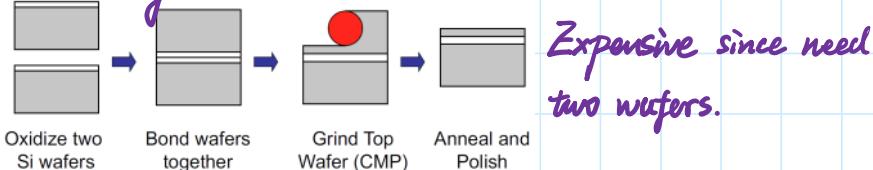


* Main Features:

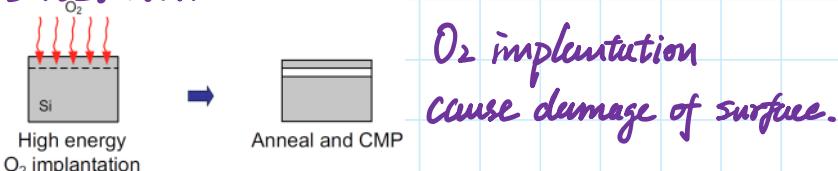
1. No junction leakage, less drain/source Cup - fast & Low power
2. Radiation hardening (CSR)
3. PD: Floating Body: Hysteresis, parasitic BJT, leaky.
FD: No Body: like a bulk device.
4. Thermal Isolation: SiO_2 has 2 order lower thermal conductivity.
 \Rightarrow Not significant globally, but can cause local hot point.

★ Fabrication:

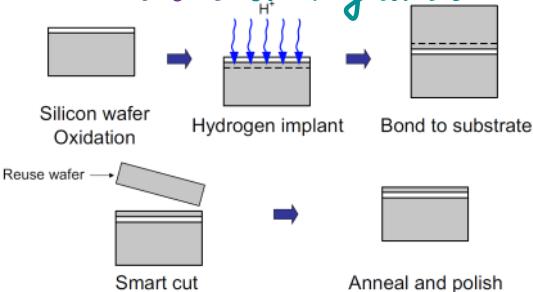
1. Bonding:



2. SIMOX:



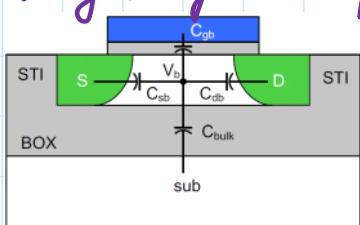
3. Smart Cut (mostly used)



★ PD SOI Device

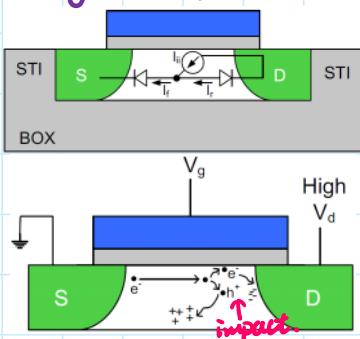
★ PD SOI Device

Body Charge - AC Cops: For fast changes (switching)



1. Change V_b
2. C_{gb} only there when no channel.
3. C_{sb} is voltage dependent.

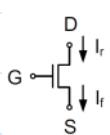
Body charge - DC: slow change (leakage)



1. SB forward biased current ($\sim nA$, in μs)
2. DB reverse biased current ($\sim pA$, in μs)
3. Tunneling current of gate ($\text{in } \mu s$)
4. Impact Ionization: Rise the body voltage.
5. Thermal generation/recombination

In actual case: AC + DC

★ Body Voltage in DC:



	G	S	D	V_b	V_t
①	0	0	0	0	V_{t0}
②	0	0	V_{dd}	0.35V	Reduced
③	0	V_{dd}	V_{dd}	V_{dd}	small

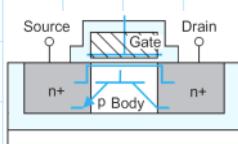
reduced by $V_b > 0$

very small V_{t0}

For ②: There is I_f & I_r . $I_f = I_r$, make $V_b \approx 0.35V$ if $V_g \neq 0$, there is impact ionization $\Rightarrow V_b \uparrow \Rightarrow V_{th} \downarrow$

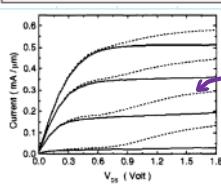


★ Parasitic BJT: at where V_b is high, but V_s drop to low



if there is a I_b (Body \rightarrow Source)

it will induce a large I_c (Drain \rightarrow Source)

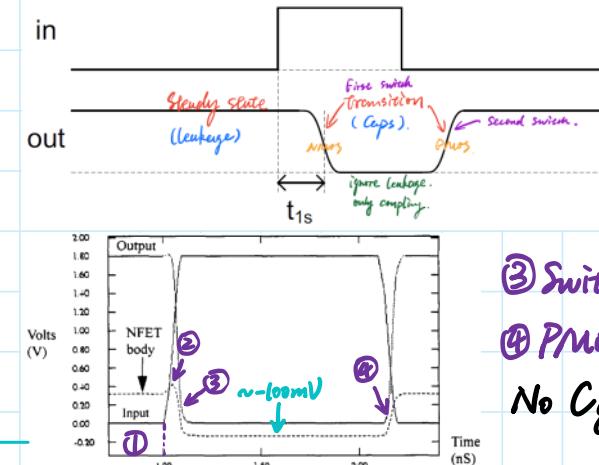


By having implant ionization & BJT
Current is higher

★ Switch Order: AC + DC

① NMOS First Switch:

① NMOS First Switch:



① Before switch:

leakage (DC) set the $V_b = 0.35V$

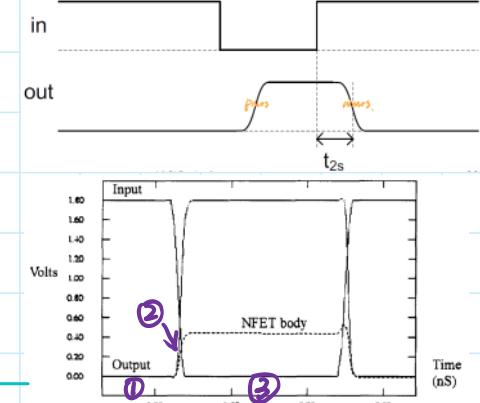
② Switching: C_{gb} coupling, $V_b \uparrow$

But C_{gb} disappear when NMOS is on.

③ Switching (V_{out} start drop): $V_b \uparrow$ (coupling)

④ PMOS switching: only C_{db} coupling ($\approx 450mV$)
No C_{gb} coupling since channel is on at beginning.

② NMOS Second Switch: PMOS Switch output from L → H first



① Before switch:

leakage (DC) set the $V_b = 0V$

② Switching (V_{out} start drop): $V_b \uparrow$ (coupling)

Small C_{gb} coupling since channel is on at beginning.

③ After 1st switch: V_b rises same amount

④ NMOS Switching: with C_{gb} & C_{db} Coupling

$V_b = 450mV \Rightarrow V_{th} \downarrow \Rightarrow$ fast switching.

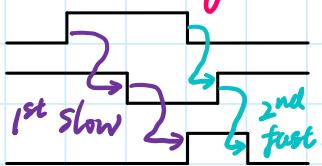
★ The first switching is slower than switch at second

It induces an effect called Pulse narrowing

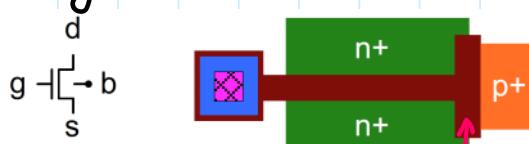
(Same for PMOS)



⇒ Pulse may disappear after a long propagation



★ Body Contact: Add a p⁺ to the well.

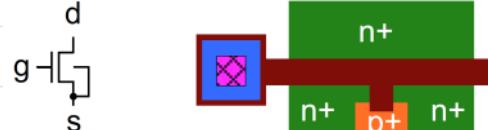


Tie body to GND

* Slower than bulk

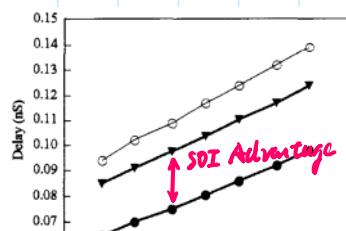
- But no delay variation

- No bipolar

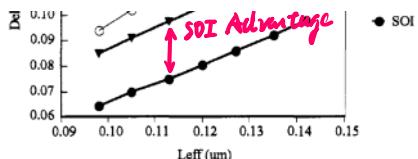


Source tie

- Increase gate cap
- Increase layout area
- Some delay variation
- No bipolar
- No body effect, $V_{sb}=0$

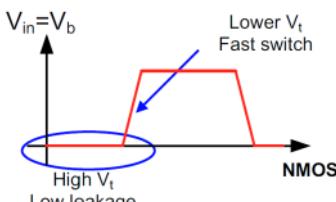
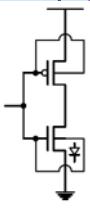


Slower than ↓ normal bulk device.



* Smart Body Contact:

① Tie to gate: $I_{leak} \downarrow$ & $t_{dv} \downarrow$



usually at the last buffer

↙ to drive heavy load.

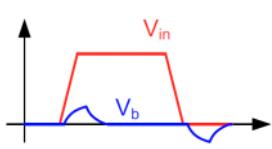
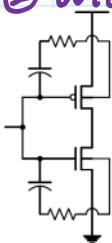
But: 1. Add gate cap.

2. Miller effect

3. R_L delay of back gate

4. limited $V_{dd} < 0.6V$ (forward of Diode)

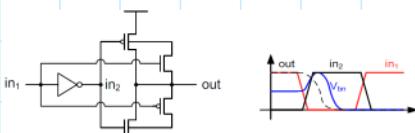
② With R_L:



No V_{dd} limitation

But hard to make

③ Smart Version:



No leakage reduction

Better as a driver

For Buffers.

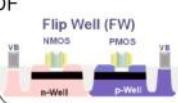


FD SOI: No body effect.

- Thin film thickness $< 50\text{nm}$
- Complete depletion when channel forms
- No bipolar
- Much less hysteresis
- Better subthreshold swing $\sim 70\text{mV/dec}$
- Lower V_{th} for same I_{off}
- But: S/D resistance
 - V_{th} sensitive to Si thickness
 - High V_{th} more difficult to fabricate

e.g. UTBB-SOI from ST.

- STMicroelectronics 28nm
- Ultra-thin body/box (7nm body, 25nm box)
- Advantages
 - strong back bias effect
 - undoped channel so low RDF
- Disadvantages
 - hard to manufacture
 - Ss penalized with thin box



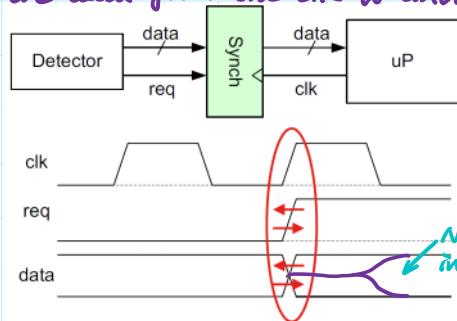
L12. Synchronization & Metastability

2023年4月6日 16:48

Synchronization:

In a system, there may be multiple clks.

The data from one clk to another requires synchronization.



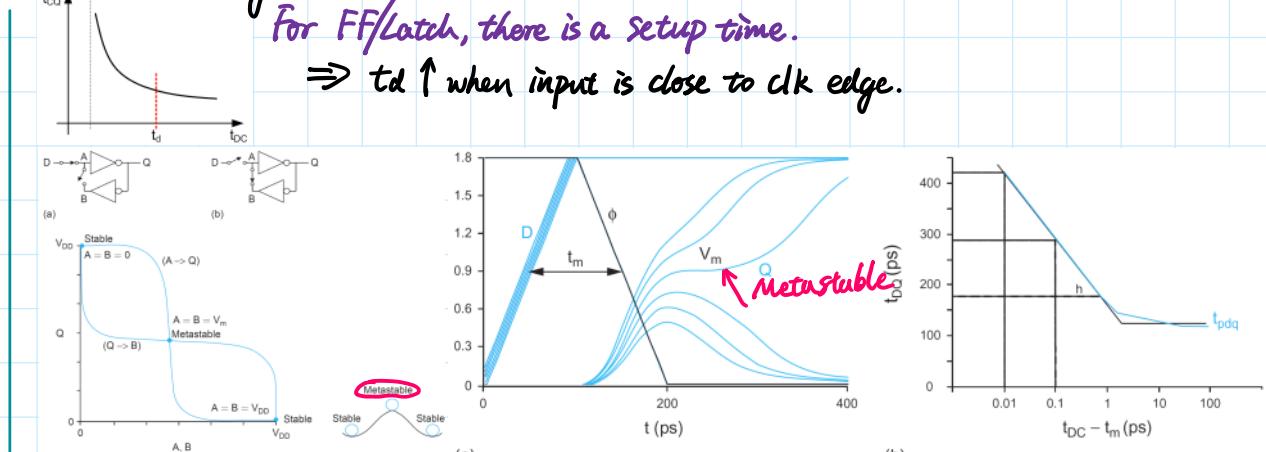
If the data/request change at the edge of clk
⇒ Induce an indeterminate state.

(Due to the metastability of the FF/latch)

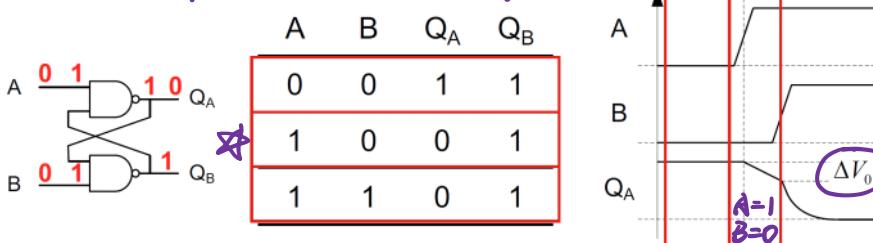
Metastability:

For FF/Latch, there is a setup time.

⇒ $t_d \uparrow$ when input is close to clk edge.



Use a crosscoupled nand as an example:



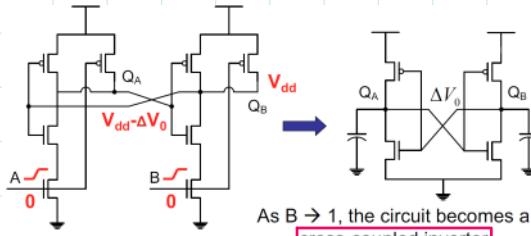
* At $A=1, B=0$. There is a ΔV drop by input: input data at edge.

Since $I = C \frac{dV}{dt} \Rightarrow \Delta V = \frac{I_0}{C_L} \cdot \Delta t$.

By set $t_a = I_0/C_L$ called Apperture Window.

$$\Rightarrow \star \Delta V_0 = \Delta t / t_a \quad - \text{the initial voltage difference.}$$

* At $A=1, B=1$, The circuit latches: evaluating the input data.



when the latch evaluating data:

$$\Delta I = g_m \Delta V \Rightarrow d\Delta V = \frac{\Delta I dt}{C_L} = \frac{g_m \Delta V dt}{C_L}$$

$$\Rightarrow \int \frac{d\Delta V}{\Delta V} = \int \frac{\Delta I}{C_L} \cdot dt \quad T = \frac{C_L}{g_m} \text{ (speed of circuit)}$$

$$\Rightarrow \star \Delta V(t) = \Delta V_0 e^{\left(\frac{g_m}{C_L} \cdot t\right)} = \Delta V_0 e^{\left(\frac{t}{T}\right)}$$



Equation to find the developed voltage in time t .

Define a resolving time τ_d for $\Delta V(\tau_d) = 1$ — the unit voltage that means fully developed.

$$\Rightarrow \tau_d = -T_r \cdot \ln(\Delta V_0) = -T_r \cdot \ln(\Delta t / t_d) \Rightarrow \text{fast circuit, less } \tau_d$$

If $\Delta V_0 = 0$ or $\Delta t = 0 \Rightarrow \tau_d \rightarrow \infty$: The data is defined by noise



★ Probability of failure:

Define $P(\Delta V(t_d) < 1)$ as failure rate.

for $\tau_d < t_d e^{-\frac{\tau_d}{\tau_r}}$, the sync fails.

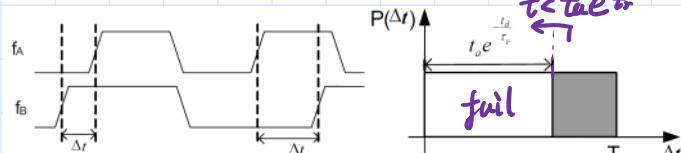
$$\star P(\Delta V(t_d) < 1) = P(\Delta t < t_d e^{-\frac{\tau_d}{\tau_r}}) = t_d e^{-\frac{\tau_d}{\tau_r}} / T_A = \text{failure} e^{-\frac{\tau_d}{\tau_r}}$$

The Δt is randomly given at the range of faster clk (clock A). The area of T_A that less than $t_d e^{-\frac{\tau_d}{\tau_r}}$ means fail.

The fail rate depends on the faster clock.

The sync event happens every clock of the slower clock (clock B)

$$\Rightarrow \star \text{Frequency of Failure} = f_A \cdot f_B \cdot t_d e^{-\frac{\tau_d}{\tau_r}}$$



★ Mean Time Between Failure (MTBF): $e^{\frac{\tau_d}{\tau_r}} / (f_A \cdot f_B \cdot t_d)$

e.g.

f_A = event frequency on line A = 50MHz

f_B = event frequency on line B = 300MHz

t_d = aperture of flop = 100ps

t_u = Wait time until decision = 10ns

τ_r = regeneration constant of flop = 200ps

Use higher frequency

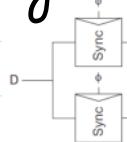
$$P_{failure} = t_d f_B e^{-\frac{t_d}{\tau_r}} = 100 \text{ ps} \cdot 300 \text{ MHz} \cdot e^{-\frac{10 \text{ ns}}{200 \text{ ps}}}$$

$$= 5.786 \times 10^{-24}$$

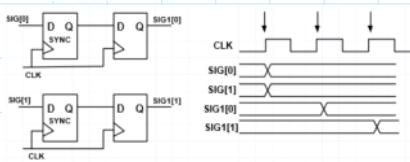
$$MTBF = 1 / P_{failure} = 1 / P_{failure} f_A = 1 / 5.786 \times 10^{-24} \cdot 50 \text{ MHz} = 3.45 \times 10^{15} \text{ sec}$$

$= 10^8 \text{ years} \Rightarrow \text{Not too good (x bus no. x chip no.)}$

★ Synchronization Pitfall:



1. Not sync one signal in two places.
⇒ May have different result.

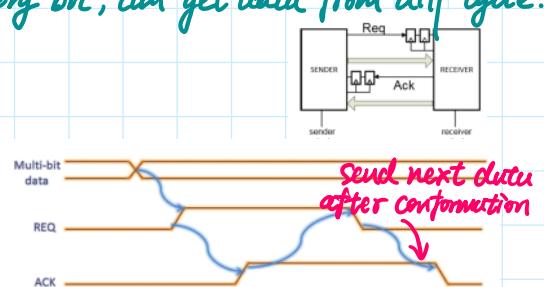


For bus, only use one synchronizer

⇒ if use for every bit, can get data from diff cycle.

★ Correct Way: Handshaking

- Sender outputs data and THEN asserts REQ
- Receiver latches data and THEN asserts ACK
- Sender deasserts REQ, will not reassert it until ACK deasserts
- Receiver sees REQ deasserted, deasserts ACK when ready to continue



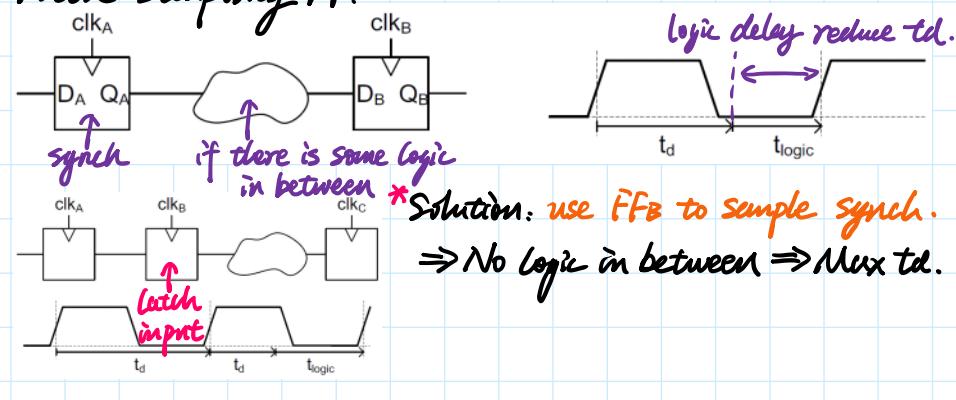
★ Fix Failure:

- Smaller cap $t_u = \frac{C_L}{I_d}$
 - Shield load cap at output of flops
- Wait longer
 - Application-specific, latency ok?
- Lower frequency

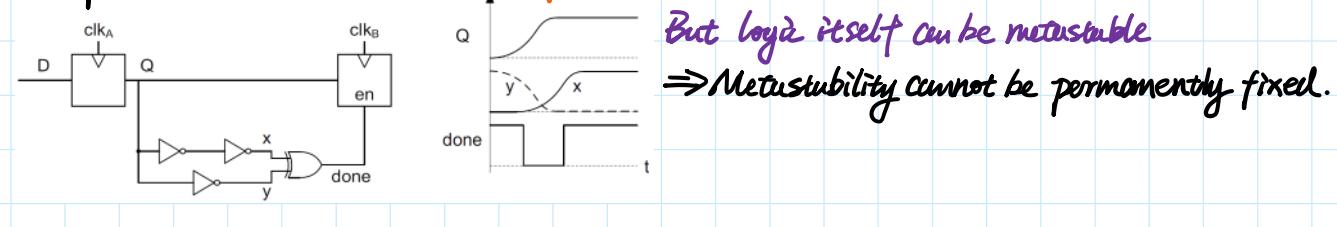
- Faster gates
 - Technology scaling helps but freq ↑ too
- Add completion logic
- Add more sampling FFs
- Multi-stage synchronizers

- Wait longer
 - Application-specific, latency ok?
- Lower frequency
 - Multi-stage synchronizers

More Sampling FF:



Completion Detector: Enable after transition



Jumb Latch: Small load Cap. fast.

