

- Introduction
 - Transistor switch model
- Logic Circuits
 - Regeneration of the level
 - Capacitance
 - * Effective fan-out
 - * Signal in reality
 - Pass gate logic
 - * Level restoration
- Circuit Timing / Dynamic Logic
 - Latch/Register Implementation
 - Sequencing, pipelining revisited
 - Dynamic logic

Introduction

One of the main thing about this class is finding how to optimize the ratio of $Op/s/W = Op/J$. We want to enhance this ratio. Less and less foundry have smaller and smaller technology. We want to reduce it both for *plugged* and *battery* device. Either we don't have the requirements to pull out or the space to store the energy.

This class is all about Gate and transistor. Low level is king.

Transistor switch model

We can model a switch (MOSFET) with an ideal switch and a resistor :

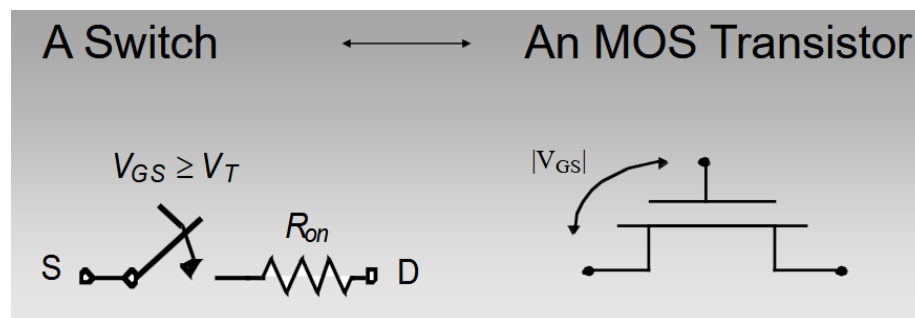


Figure 1: Switch

By the transistor scaling, short channel are behaving differently due to the **velocity saturation**. The relation becomes linear and not quadratic like it used to be.

- case 1 : $V_{\min} = V_{DS} \rightarrow$ Linear region

$$I_{DSAT} = \mu C_{ox} \frac{W}{L} \left((V_{GS} - V_T) V_{DS} - \frac{V_{DS}^2}{2} \right)$$

- case 2 : $V_{\min} = V_{GS} - V_T \rightarrow$ (Channel) Saturation
- $$I_{DSAT} = \frac{1}{2} \mu C_{ox} \frac{W}{L} (V_{GS} - V_T)^2 (1 + \lambda V_{DS})$$

- case 3 : $V_{\min} = V_{DSAT} \rightarrow$ Velocity Saturation

$$I_{DSAT} = \mu C_{ox} \frac{W}{L} \left((V_{GS} - V_T) V_{DSAT} - \frac{V_{DSAT}^2}{2} \right) (1 + \lambda V_{DS})$$

Figure 2: Equations

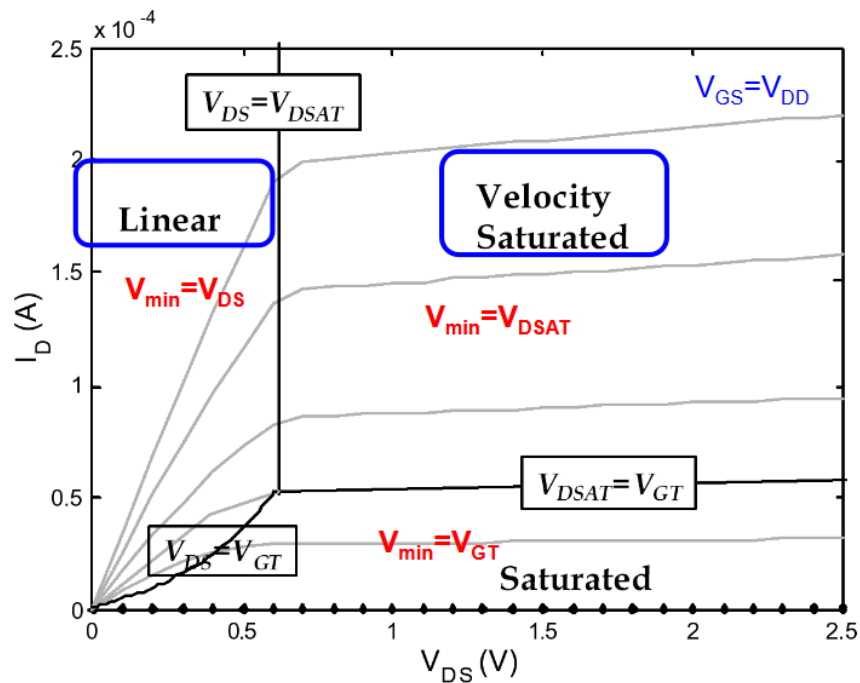


Figure 3: Regions

Logic Circuits

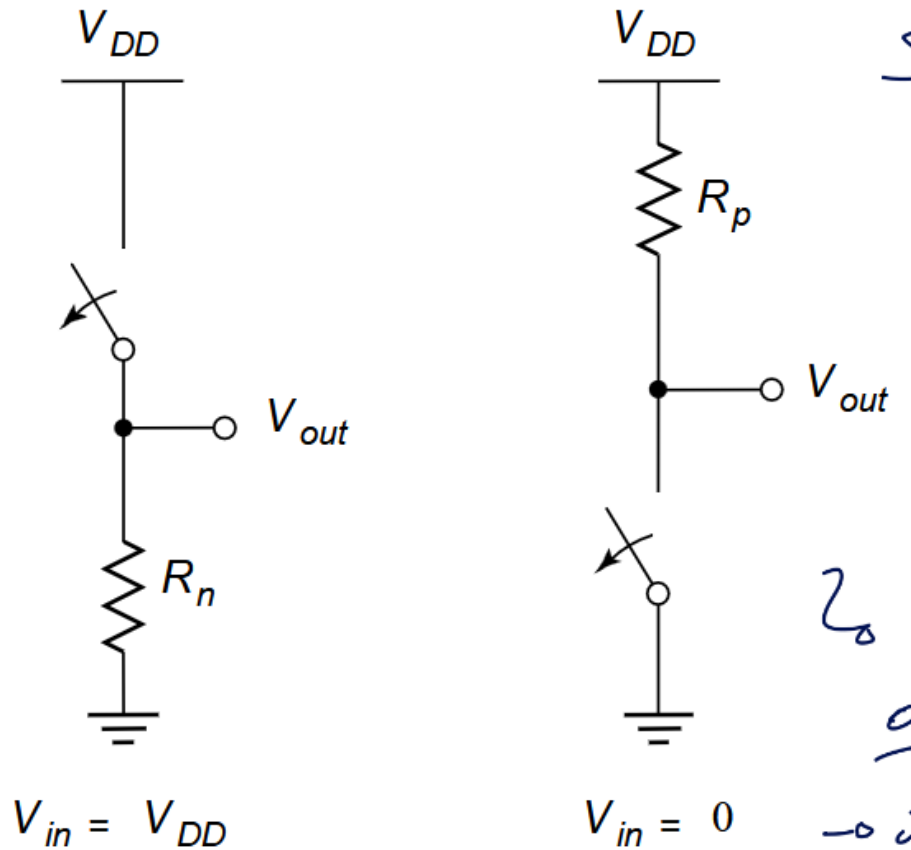


Figure 4: The static model

The swing here is equal to V_{dd} so we have a high noise margin. It is not a **ratioed** logic so we can't use tricks to minimize mismatch by taking advantage of ratios. We only have 1 resistor on so low output impedance but the input is the gate of MOS so we have a high input impedance. There is no static power consumption since no direct path from Vdd to ground. That's nice :)

In the dynamic model we need to add an output cap C_L . The load cap is simply the sum of all capacitance at the output node. Transition time is determined by the charging of this cap by a resistor. The sizing impacts the dynamic behavior of the gate.

Ideally we want V_M to be at the middle of the other nominal voltage. We call the region in between the *undefined region*.

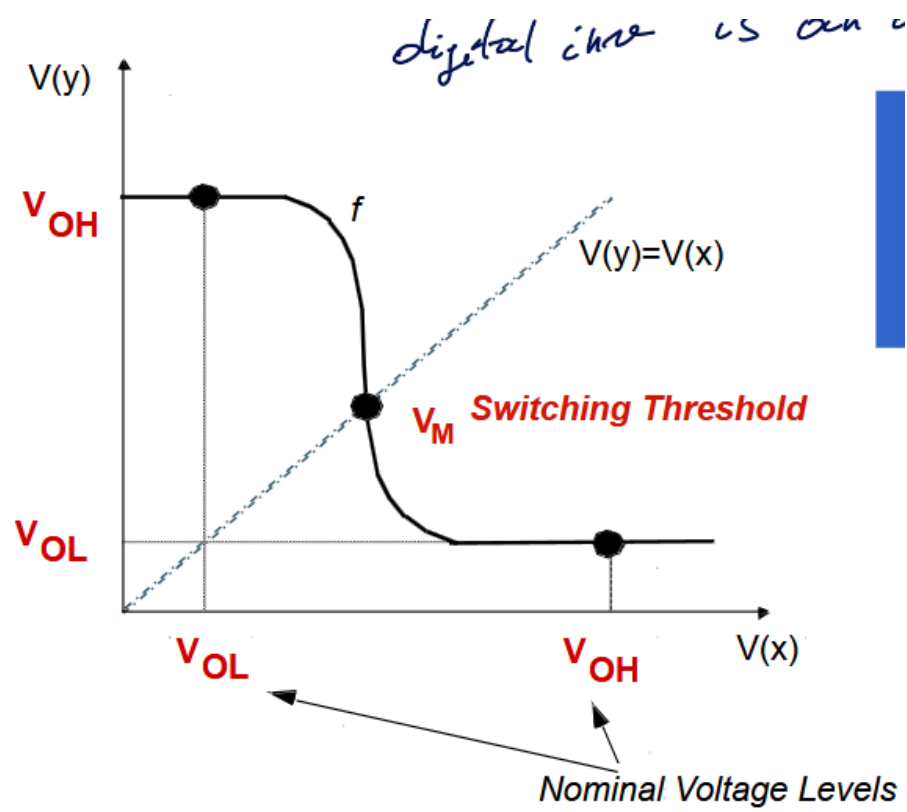


Figure 5: Switching threshold

Regeneration of the level

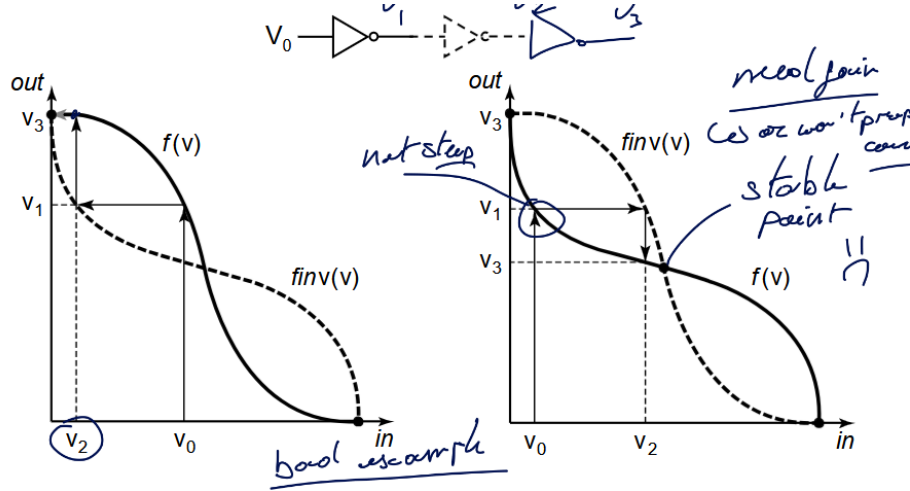


Figure 6: NAND gate regeneration

With using this type of gate we have some regeneration level, it will amplify the signal and so we won't have undefined level and we will have the signal that will reach one defined state. If we have no regeneration, we will reach meta-stability. We have to meet some conditions :

- The transient or undefined region in the VTC should have a gain $|dV_{out}/dV_{in}|$ larger than 1
- In the "legal" or defined regions the VTC gain should be smaller than 1 in absolute value
- The boundary between the defined and undefined regions are V_{IH} and V_{IL} where the gain = -1

We need gain or the signal will be lost.

We have 3 different types of noise :

1. Inductive coupling
2. Capacitive coupling
3. Power and ground noise

The noise margin in CMOS is rather high which is a good thing seeing the low output impedance.

We see that the ratio of PMOS and NMOS determine the V_M voltages.

Capacitance

We know that the delay of a switch is $t_{phl} = f(R_{on}C_L) = \ln(1/2)R_{on}C_L = \ln(1/2)(R_{eqn} + R_{eqp})/2 \cdot C_L$. We are still observing some glitches when we switch

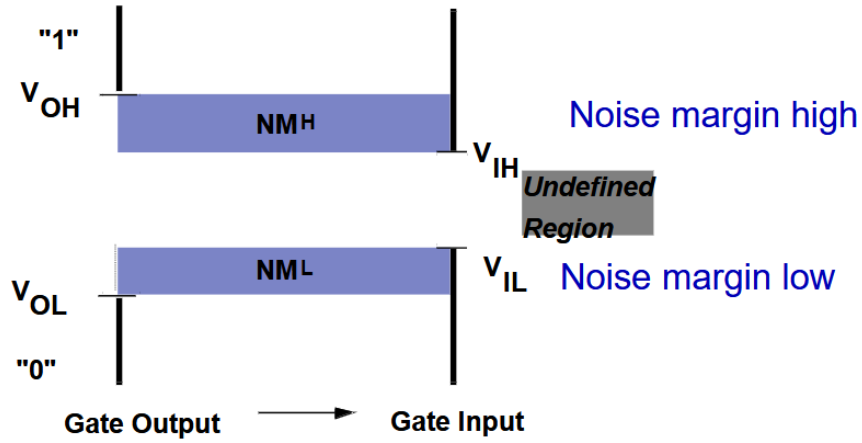


Figure 7: Noise margin

$$I_n(V_{GS} = V_M) + I_p(V_{GS} = V_M - V_{DD}) = 0$$

$$k_n V_{DSATn} (V_M - V_{Tn} - \frac{V_{DSATn}}{2}) + k_p V_{DSATp} (V_M - V_{DD} - V_{Tp} - \frac{V_{DSATp}}{2}) = 0$$

Solving for V_M yields

$$V_M = \frac{(V_{Tn} + \frac{V_{DSATn}}{2}) + r(V_{DD} + V_{Tp} + \frac{V_{DSATp}}{2})}{1 + r} \text{ with } r = \frac{k_p V_{DSATp}}{k_n V_{DSATn}}$$

The ratio $(W/L)_p / (W/L)_n$ to set a certain V_M is given by

$$\frac{(W/L)_p}{(W/L)_n} = \frac{k_n' V_{DSATn} (V_M - V_{Tn} - V_{DSATn}/2)}{k_p' V_{DSATp} (V_{DD} - V_M + V_{Tp} - V_{DSATp}/2)} \text{ (both TOR in velocity saturation!)}$$

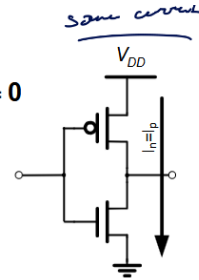


Figure 8: Switching threshold for INV

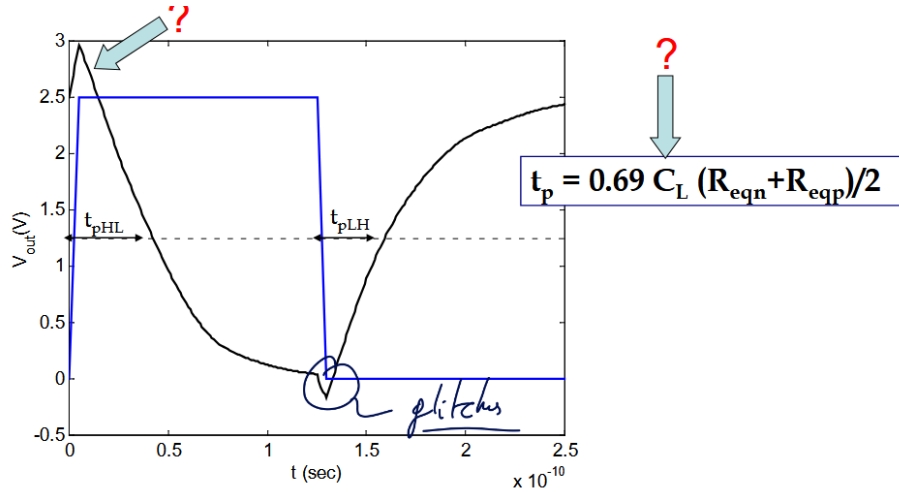


Figure 9: Glitches and Delay

on and off. This **isn't** due to the miller effect. This **overshoot** is due to the gate drain capacitor.

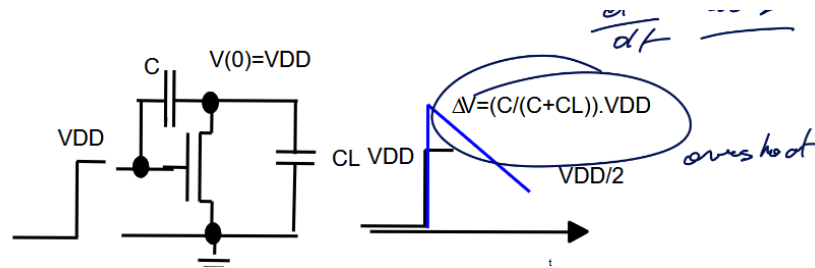


Figure 10: Explanation of the overshoot

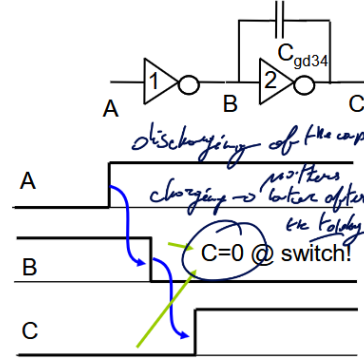
This is due to charges and sudden and “infinite” steep step at the input which will create an extra unwanted voltage. Thankfully the input isn't as steep in reality and so the effect is less severe but still noticeable. We call it the *digital miller effect*:

$$t_d = \frac{C_{total} \cdot (\Delta V + V_{DD}/2)}{I_{max}} = \frac{(C + C_L)}{I_{max}} \left(\frac{C}{C + C_L} V_{DD} + \frac{V_{DD}}{2} \right) = \frac{(3C + C_L) \cdot V_{DD}}{2I_{max}}$$

So it is like the cap becomes 3 times larger (similar to the miller effect) but in reality we are closer to 2 since we never have a perfect step at the input.

We can move those C_{gd} to the inside and see it as an impact on C_L . Again

- We analyse the influence of C_{gd34} on the delay of Inverter I_1
- C_{gd34} is loading node B, but
 - C_{gd34} is first charged $\frac{+}{B} \parallel \frac{-}{C}$
 - when node B is switching node C is at ground,
- Charge required to bring node at 0 and discharge C_{gd34} is $C_{gd34} \cdot V_{DD}$
- When node C finally switches
 - C_{gd34} is charged $\frac{-}{B} \parallel \frac{+}{C}$
 - requiring another charge $C_{gd34} \cdot V_{DD}$



Conclusion: “Recharging” C_{gd34} requires a charge of $2 \cdot C_{gd34} V_{DD}$, but only half of it is exchanged during the switching of I_1
 $\Rightarrow C_{gd34}$ is seen only “once” in the delay of I_1

Figure 11: Loading issue

we can reuse the theory of DDP with the intrinsic and extrinsic load where $C_L = C_{int} + C_{ext}$:

$$t_p = 0.69R_{eq}(C_{int} + C_{ext}) = 0.69R_{eq}C_{int} \left(1 + \frac{C_{ext}}{C_{int}}\right) = t_{p0} \left(1 + \frac{C_{ext}}{C_{int}}\right)$$

So the sizing can help up to a certain point where we have an *irreducible* delay.

Effective fan-out

The input C_g and intrinsic cap are always proportional to the sizing : $C_{int} = \gamma C_g$ where γ is a technological constant. Same goes for the extrinsic load C where it is the input C of the next inverter proportional to the sizing $C_{ext} = f C_g$. So we can summarize the delay t_p by :

$$t_p = t_{p0} \left(1 + \frac{f}{\gamma}\right)$$

The delay depends on the ratio between its external load capacitance and its input capacitance, the ratio is called the *effective fan-out*.

The newly introduced γ is not valid for dynamic logic or more exotic technologies. If this $\gamma = 1$ then it means that $C_{int} = C_{in}$. It is an acceptable approximation in standard CMOS logic.

$$t_p = kR_w C_{int} (1 + C_L / C_{int}) = t_{p0} (1 + f / \gamma)$$

$$C_{int} = \gamma C_{gin} \text{ with } \gamma \approx 1$$

$$f = C_L / C_{gin} - \text{effective fanout}$$

$$R = R_{unit} / W ; C_{int} = W C_{unit}$$

$$t_{p0} = 0.69 R_{unit} C_{unit}$$

$$R_{unit} \sim 1 / V_{DD}$$

Figure 12: The golden formula

For the ring oscillator where we assume equal size $f = 1$ is independent of the size. In real technology we see only a weak dependency of timing on sizing.

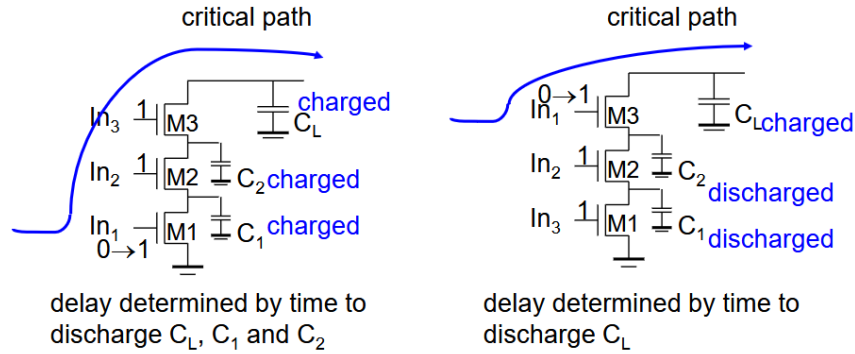
Signal in reality

We know that we can't have infinitely steep signal and even worse, a too slow rise and fall time could lead to metastability issues ! We could also have some actual short circuit for a brief amount of time leading to waste of energy. So in most of design software we will leave some headroom to avoid possible short-circuit and we will flag it with *max transition violation*.

Note on sizing When we see the a or b next to a transistor it is its *relative sizing* compared to a classic $\frac{W_{min}}{L_{min}}$. For complex gate and due to the various sizing we can have, we will transform a little bit our formula :

$$t_p = t_{p0} \left(p + \frac{gf}{\gamma} \right) = t_{p0} d$$

- f : electrical effort
- g : logical effort
 - due to the fact a logical gate is always slower than an inverter with equal current drive. Ratio of input cap to the cap of an inverter that delivers equal current.
- p : ratio of intrinsic delay of the gate to the intrinsic delay of an inverter
- d : gate delay
 - relative to the intrinsic delay of the reference inverter



Critical path should be connected to the input TOR closest to the output

Figure 13: Critical path & Charging

Pass gate logic

Another approach than PUN and PDN as seen previously is the pass gate logic where we will not only use the gate but also the source of a transistor to create some gate.

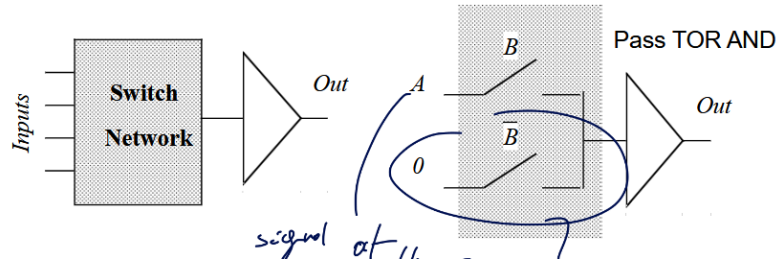


Figure 14: Pass gate logic

This technique uses less transistors. But the NMOS isn't a really good pull up transistor. It won't give a nice and crisp high voltage but rather a voltage with a $\Delta V = V_{T_{loss}}$. So to keep this signal crisp and nice we are obligated to add an INV to output a valid signal. This goes in the same idea as the *regeneration of the signal* logic.

So cascading of passes tor logic isn't a smart choice since the signal will rapidly deprecate.

Level restoration

So we need to do what we call *level restoration*.

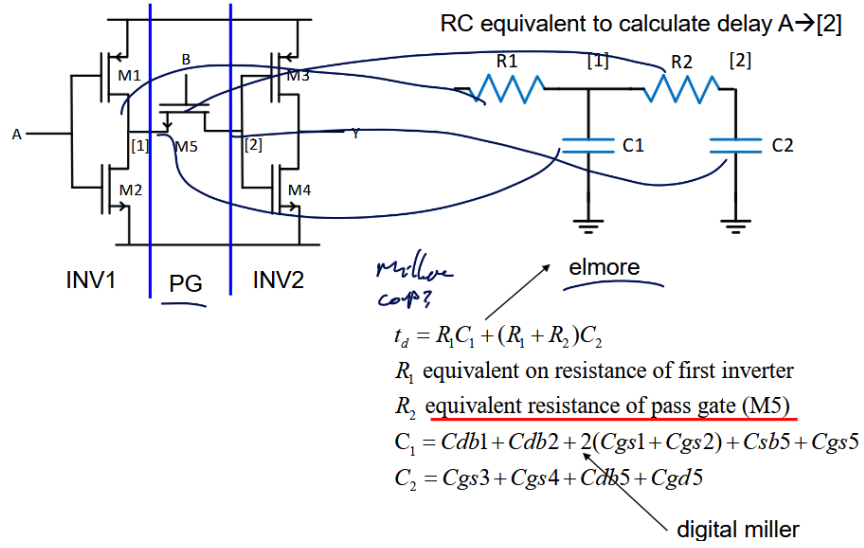


Figure 15: Static analysis

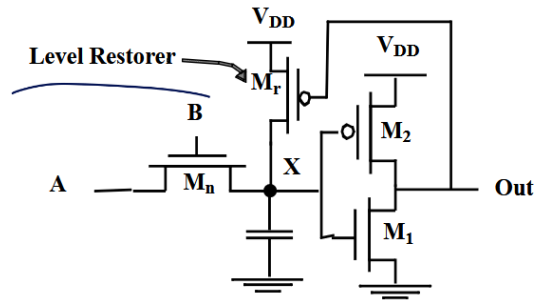


Figure 16: Level restorer

It is compatible with the full swing and the static power consumption is gone. But we need a bleeder which will increase capacitance at node X and takes away the pull down current. Leads to speed degradation and needs proper sizing.

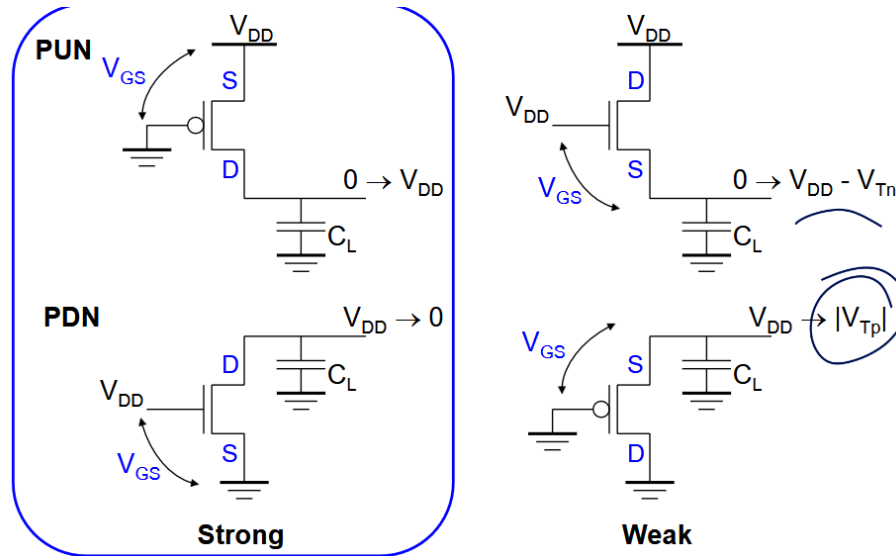


Figure 17: Threshold drops generalized

Transmission gate We can add a PMOS to the switch to create a **transmission gate**. This requires another transistor but also a complementary signal is needed so this will result in extra cost.

Circuit Timing / Dynamic Logic

Latch/Register Implementation

We know that having 2 INV back to back could lead to meta-stability so how can we reliably store data ? We want to remember the data *no triggering* but also sample and so stop looping the data *triggering*. In the second approach we simply *overpower the feedback loop* due to the asymmetry of the INV and so the input D will be more important than the output of the small INV (**David-Goliath latch**). Or we can cut the loop like in the second example.

Specifications	Positive Latch	Register
Storage ?	store when clock is low	Stores on rising edge
Transparent ?	Yes when clock is high	No

- overpower the feedback loop

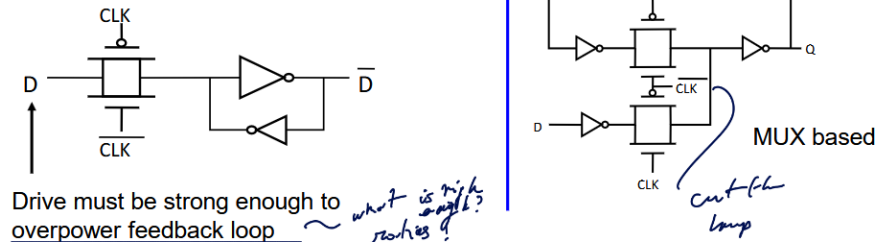


Figure 18: Latches

For latches we can either go for positive or negative latches and a register is simply two latch back to back. We can do latches using MUX for example.

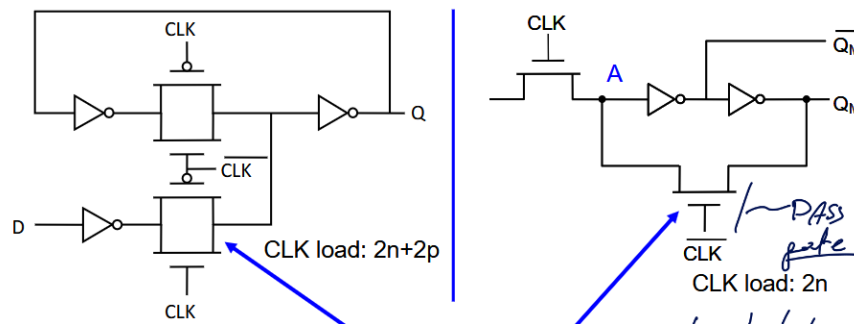


Figure 19: Mux-Based latch : Transmission Gate vs Switch

The second option is more preferable because it will have less load on the CLK which reduce the energy waste. But we need to remember we have to *regenerate* the signal since we will have a V_{Tloss} .

The latch is pretty similar to the idea we have seen with pass gate. Here we can see again the bleeder on top and NMOS on the bottom which forms a basic INV.

NEED TO READ THIS PART CAREFULLY

Sequencing, pipelining revisited

Dynamic logic

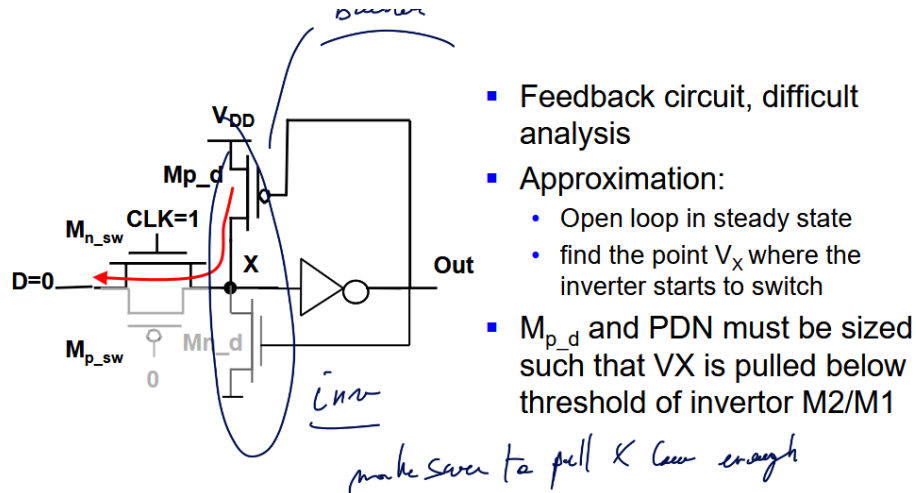
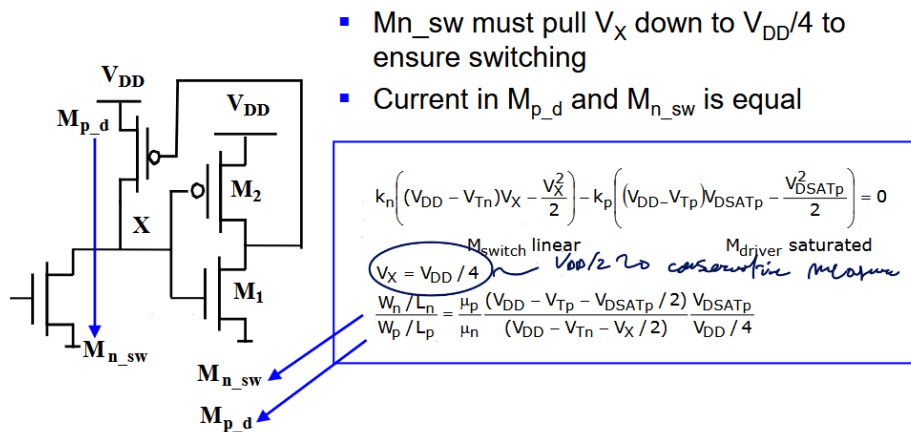


Figure 20: David Goliath



This is a worst case approach. In practice the bleeder will already be "weakened" because at $V_X = V_{DD}/4$, the inverter has already begun to switch off.

Figure 21: Sizing of the latch