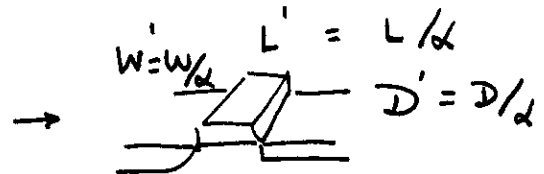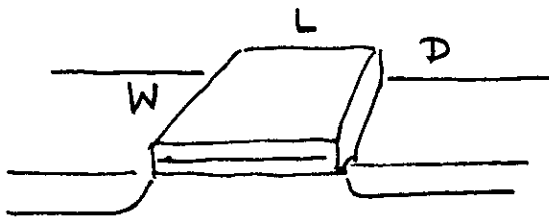# LECTURE # 15        NOVEMBER 9

- Collect Proj. Assign. # 2

- <u>LAB /PROJ</u>: Use only Boxes at right #'s. Later software won't support more.
  Also: We have priority in Lab > BPM.
  POINT OUT # COSTS: ONCE DES RIGHT, STAMP OUT LIKE COOKIE AT ~ #1 or so /piece.

- <u>TODAY</u>: CONTINUE SCALING.
  DISCUSS LIMITING FACTORS.

- <u>SUMMARIZE SCALING SO FAR</u>:



All Voltages  $V' = V/\alpha$

Using this simple scaling by $\alpha$ of all dimensions including vertical, and all voltages (VDD, $V_{TH}$, $V_{Dep}$) we found:

$$\boxed{\tau' = \tau/\alpha} \quad ; \quad \boxed{I_{ds}' = I_{ds}/\alpha} \quad ; \quad \text{PROBLEM ENCOUNTERED!}$$

[ex.: EQ 1:  $\tau = l^3/\mu V_{ds} \therefore 5$]

But # of devices per unit area goes up by $\alpha^2$
So mean current density : <u>CURRENT INTO AREA OF CHIP</u>
<u>GOES UP by $\alpha$</u>.  This is real problem, since in our scaling wires would get thinner. Current limited wires would either have to have aspect ratios increase by $\alpha^2$ (which cant go on for long) or get proportionally wider.

$$\boxed{C_g' = C_g/\alpha} \quad ; \quad \text{But } \underline{\text{capacitances}} \text{ in general scale up by } \alpha \text{ if given in terms of } C/\mu m^2 \text{ since vertical dimensions shrinking (oxide getting thinner).}$$

<u>Resistance/□</u> <u>scale up by</u> $\alpha$ since lines getting thinner.
(Except note that we can't really do that with curr. lim. metl.)

However, if calculate it out, find that R/□ of FETS will stay about the same.

<u>So another problem</u>: R/□ of poly, diff getting proportionally larger while R/□ FET staying same. This is aggravated by crystal clumping in POLY --- makes effect worse as scale Down.

<u>DC Power Dissipation</u>: Per Device $P_{dc} = I \cdot V$

$$P'_{dc} = \frac{P_{dc}}{\alpha^2}$$

; # Dev. / Area goes up as $\alpha^2$

So (whew!) Power dissipation/unit area stays approx constant.

<u>We discussed power dissipation limits</u>: Is much more diff. to pin down to single constraint as in current density in wires. <u>Dependent on next level context</u>

<u>Tabulated</u>:
$< 1$ w/cm$^2$ no prob.
$2$ w/cm$^2$ } begin to need
$4$ w/cm$^2$ } reasonable heat sinking
$> 8$ w/cm$^2$ need way to remove heat: forced cooling

<u>We examined ≈ worst case in our methodology</u>:

Large array of shift registers are in Fig 86 CH 4.

Found power/unit area ≈ 10 W/cm$^2$.

BUT NOTED THAT WE USE PASS-T logic BETWEEN THESE AND USUALLY DON'T USE SUCH SHORT PULSES, etc.

So NORMALLY WE DON'T NEED TO WORRY. BUT SHOULD CALCULATE IF IN DOUBT. TRY TO STAY $< 2$ W/cm$^2$, IF CAN. (over less area to alleviate)

Scaling of
<u>SWITCHING POWER</u>: The drivers which operate pass gates, charging
& discharging capacitances dissipate switching power.

<u>The power is dissipated at the drivers</u>, but we calculate

the amount based on the capacitances & voltages &
clock period:

$P_{sw}$ = energy stored on capacitance divided by the
clock period or time between successive
chargings / discharging.

But $T' \propto T$. Thus, $P_{sw} \propto \dfrac{C V^2}{T}$ ; $T \propto \dfrac{L^2}{V}$

$\therefore P_{sw}' \propto \dfrac{WL}{D} \cdot V^2 \cdot \dfrac{V}{L^2} = \dfrac{W V^3}{D L}$  $\therefore$ $\boxed{P_{sw}' = \dfrac{P_{sw}}{\alpha^2}}$

<u>So</u>: Since $P_{sw}$ per device goes down by $\alpha^2$, and
# dev./area goes up by $\alpha^2$, the switching
power also stays constant per unit area
as we scale things down.

<u>NOTE</u>: AVERAGE DC POWER DISS. IN MOST SYSTEMS CAN
BE APPROXIMATED BY ADDING TOTAL $P_{sw}$
TO 1/2 DC POWER RESULTING IF ALL
LEVEL RESTORING LOGIC WERE TURNED ON.

# SWITCHING ENERGY

We've noted before that there are various ways to trade off power vs delays. We can often use less power if we can tolerate longer delays, and vice-versa. This can be done by binding it into the design, or sometimes can be controlled dynamically. This reflects an important metric of device performance: <u>SWITCHING ENERGY per DEVICE</u>

$E_{sw}$ = power consumed by device at max clock Freq. multiplied by the delay: i.e., it is a "power × delay" product.
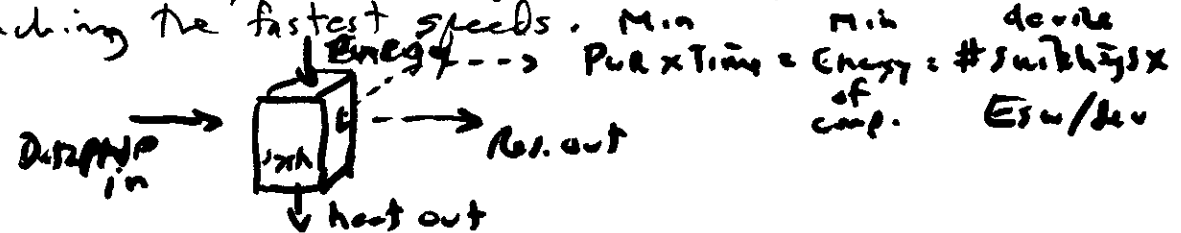
<u>Rationale:</u> In a sense, to do any computation, we must <u>at minimum</u> switch a large collection of switches, switching (Think of toggle switches) them in a particular order, and some number of times.

Switches have some switching energy which measures the work done to throw the switch. We often have the option (by design or control) to choose to put the energy into a system slowly, (and thus less power) taking longer for a calculation. Or- put it in faster, flipping the switches faster. (SEE CHAP 9)

However, there are usually constraints on both speed and power, and these limit ultimate performance:

No matter how much power we put in, we can't reduce delays below the minimum value of $\gamma$.

Also, if we put in too much power, we may reach a power density limit in a particular design even before hitting the fastest speeds. Min

Min Energy ---> Pwr × Time = Energy = #Switchings × Esw/dev

Min device of comp.

## HOW DOES SWITCHING ENERGY SCALE?

Our basic FET switches have $E_{sw} \propto CV^2$

and $\therefore$

$$E'_{sw} = \frac{E_{sw}}{\alpha^3}$$

So this crucial metric of device performance scales incredibly favorably. This is why scaling down sizes is so important.

## Summary So Far: Suppose we

Scale down an entire system by $\underline{\alpha = 10}$.

> Resulting system will have <u>100X as many devices/unit area.</u>
  (or take only 1/100 th as many chips)
> Power Density remains constant.

> All voltages reduced by <u>factor of 10</u>
> Current /Area increased by <u>factor of 10</u>
       (on chip)
> Time delay /stage <u>decreased</u> by factor of 10
> Power-Delay Product decreased by " " 1000
  of Devices

This is very attractive scaling except for the current density problem. The delivery of the required average dc current presents an important obstacle to scaling. Even in today's systems, many wires are operated at near there current limit. So, wires must become relatively wider, or have much higher aspect ratios, or both.

$\left(\text{& Don't forget the problems of Poly, Diff res/}\square \text{ rel to FET}\right)$

CONSIDER

- Delays to Outside World : (Read SPACE vs TIME in CH 1)

What is effect of scaling on output driver design ? delays?
We can't just scale them down : the outside world stays
big. Remember the result in chap 1 :

Min Delay when use factor of $e$, and $N = \ln Y = \ln \dfrac{C_L}{C_{g_{min}}}$

In this case: Min Tot Delay $\approx \tau e \ln\left[\dfrac{C_L}{C_g}\right]$

Now, scale everything down by $\alpha$, including Voltages.
(This we do scale even in the external world)

$$\tau' = \tau/\alpha \; ; \; C_g' = C_g/\alpha \; ; \quad \therefore \quad Y' = \alpha Y$$

DERIVE
↓ ⌐→ $* \quad \boxed{\;t'_{min} = t_{min} \cdot \dfrac{1}{\alpha}\left[1 + \dfrac{\ln\alpha}{\ln y}\right]\;}$

below for
deriv.
if necessary

So, as inverters get smaller, more stages are required
to obtain minimum offchip delay.

The relative delay to outside world increases,
But the absolute delay decreases!

ALSO: At Least Driver Designs must change ;
     can't be just scaled down like        (cont.)
     rest of system

---

Derive * :  $\quad t_{min} = \tau e \ln Y$

$t'_{min} = \tau' e \ln Y' = \dfrac{\tau}{\alpha} e \ln(\alpha \cdot Y) = \dfrac{\tau}{\alpha} e[\ln\alpha + \ln Y]$

$t'_{min} = \dfrac{\tau e \ln Y}{\alpha}\left[1 + \dfrac{\ln\alpha}{\ln y}\right] = \dfrac{t_{min}}{\alpha}\left[1 + \dfrac{\ln\alpha}{\ln y}\right]$

SO, SCALING PRODUCES SOME GREAT EFFECTS. ⟶ ↑

↓ CONT.

BUT WE SHOULD ASK: HOW SMALL CAN WE MAKE NMOS THESE DEVICES AND STILL HAVE THEM WORK?
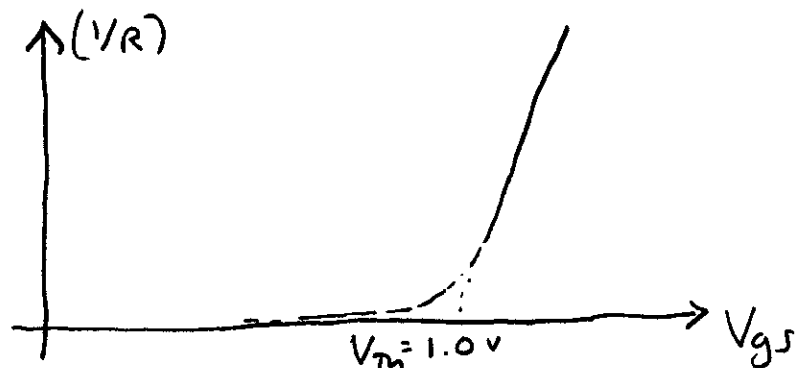
WE MUST SUSPECT THAT THERMAL, STATISTICAL, QUANTUM EFFECTS ARE ULTIMATELY GOING TO MESS THINGS UP!

QUESTION: IF PATT'S, FAB WERE NOT LIMITING US, HOW SMALL COULD WE MAKE FETi AND STILL HAVE THEM WORK?

[
MANY FACTORS TO CONSIDER. I'LL DISCUSS SEVERAL OF THE MAJOR ONES, ONE IN SOME DETAIL. IF YOU'RE INTERESTED IN THIS: I SUGGEST READING THE SURVEY PAPER BY KEYES, AND ALSO BROWSING IN CHAPTER 9
]

- **SUBTHRESHOLD CONDUCTANCE:**

IF WE REALLY PLOT DETAILS OF COND. VS $V_{gs}$, IT IS NOT SIMPLY A STRAIGHT LINE RUNNING DOWN TO $V_{TH}$, BUT HAS AN EXPONENTIAL TAIL:



$V_{Th} = 1.0 V$

Below $V_{Th}$, the conductance $1/R$ is not zero but depends on $V_{gs}$ and temperature:
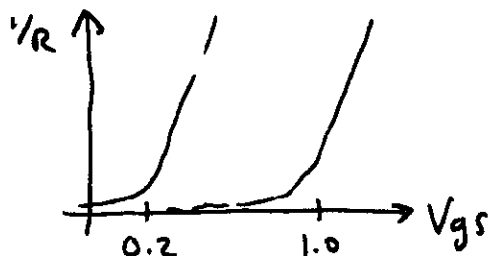
$$\frac{1}{R} \propto e^{(V_{gs} - V_{Th})/(kT/q)}$$

$T$ = absolute temp
$k$ = Boltzmann's constant
$q$ = charge on electron

<u>At room temperature</u>, $\frac{kT}{q} \sim 0.025$ volts

Thus, at <u>present</u> threshold voltages, an off device, below threshold by perhaps 0.5 volts, is below threshold by $20 kT/q$. Thus its conductance is decreased by a factor of $\approx 10^7$ over that when on near threshold. Said another way: if used as a pass-T, Q taking T to pass thru on device will take $10^7 T$ to pass thru off device.

<u>BUT</u> suppose scale down by factor of 5 :



$\left[ \begin{array}{l} \text{Curves keep } \sim \text{some form.} \\ \text{But shifted left.} \end{array} \right]$

NOW the OFF FET is DOWN ONLY BY $4kT/q$, ∴ may have as much as $\frac{1}{100}$ conductance when off as when on.

<u>Use of dynamic</u> storage especially in memories, where stored for many $T$, will be increasingly harder, espec. below $1 \mu m$. OF course trying to do things statically causes us power dissipation problems. Ah --- you can now see how we're going to get boxed in.

<u>We could</u> scale without continuing to scale voltages say when VDD reaches about 1V. But this also causes us power dissipation problems.
EVEN NOW VOLTAGE SCALING SHOULD BE DONE TO AV. MOS PERF.
<u>Could reduce prob by</u> <u>Reducing T</u>. See interesting paper by Gaensslen, Rideout, Walker CMS where very small mosfets were op. at liquid N temps and measurements confirm improvements.

BUT
THE
TR
Cap
HAD
at
5V

( Ah, related to low temp operation: Side point:
─ So→ Reducing temperature also reduces $E_{sw}$. Figures
    quoted were at room temperature. But be careful!

CH9 shows interesting comparison: FET's vs Jos.Jcns:

I won't go into full detail, read if you are interested:

Although switching energy at device is lowered,
you must put in energy into the refrigerator to
keep it at the low temp, at least as much
as difference resulting from lower $T_{\text{operation}}$.
       NOW, A COMP DIFF TECH: USE FLUX NOT Q TO ST. INFO
IBM CS's quote the low $E_{sw}$ of Josephson junctions
at the low temperature environment. It turns out
that if you scale FET's down to ∼ ½ μm,
gaining the $\alpha^3$ improvement in $E_{sw}$,
and operate them at the same temp as J-J's —
They will have the same $E_{sw}$!

But to calculate the energy required for a computation,
must use T at the heat sink temperature. Refrigerating
devices to reduce energy of computation is the logical
equivalent of constructing a perpetual motion machine.

So: Viewed as a system: FET system and JJ system will
have similar energy switching energy requirements.
FET will be simple (operate ≈ room temp.). J-J has
advantage of trading the power-delay trade off to
lower values of delay (∼ 1/50 best FET's). J-J's
at quantum limits at ∼1μm sizes. can't be scaled down smaller. ---

So the factor of 100,000 quoted by IBM CS's
of switching energy is wiped out by   X1000
(possibly) improvement in FET's by scaling, and X100
    due to the perpetual motion machine error.

MAYBE DON'T BOTHER WITH THIS

JUST MENTION THESE BRIEFLY:

OTHER LIMITING FACTORS: (SEE REF BY KEYES, SEE CH 9)

- STATISTICAL VARIATIONS IN THRESHOLD VOLTAGE:

  AS WE SCALE DOWN, WE'LL FIND THAT $\dfrac{\Delta V_{Th}}{V_{Th}}$ is proportional to Scaling factor $\alpha$.

  Results from granularity ; statistical distribution of substrate impurity charges which determine the threshold voltages.

  At same time, # devices increasing. If pullup threshold goes one way, and pulldown another, may end up with inverter which doesn't work. As shown in ch. 9, this may also limit how small supply voltages can be made. In VLSI system contain $10^7$ inverters, if we require probability (that all FETs being within threshold limits) = 0.9, may require VDD $\approx 0.7$ V.

- <u>Quantum Effects</u>. Gate oxide is already only 1000 Å = 0.1 μm thick. Positional uncertainty for electrons is related to uncertainty in momentum by

$$\Delta p \, \Delta x \approx \hbar$$

  For energy barrier of ~1 ev, calculating corresponding $\Delta p$, we find $\Delta x$ is about .001 μm. Gate oxides and junction depletion layers must be many times this or electrons will "tunnel" through. Thus we are near a fundamental size limitation due to quantum phenomena.

SUMMARIZING HOW THINGS MAY GO:

| | 1978 | MID-80's | 19XX |
|---|---|---|---|
| MIN FEAT. SIZE (2λ): | 6μm | 1μm | 0.3μm |
| T : | 0.3 to 1.0 ns | ~0.05 to 0.15 ns | ~0.02 ns to 0.04 ns |
| $E_{SW}$ : | ~$10^{-12}$ J | ~$5 \times 10^{-15}$ J | ~$2 \times 10^{-16}$ J |
| LOCAL SYNCH SYS : CLK. PERIOD (~100T) | ~30 to 100 ns | ~5 to 15 ns | ~2 to 4 ns |

> The mid 80's column we'll probably reach without major hassles. Voltage will be scaled to 1½ or 1$^V$, and power density won't be too much of a problem.

> Subthreshold current will be emerging as a problem, but not within our digital processing structures where "refresh" occurs every 50T or so.

> Current density will be a rapidly emerging problem, but will be handled with more area devoted to power lines, and higher-aspect ratio wires.

• Getting the last order of magnitude out of the technology before fundamental physical limits are finally hit will, however, be a major hassle. It will require close collaboration of researchers spanning the range from CS, to Arch., to E.E., to Device phys., to materials, in order to provide the small catalyst to help narrow down, select the alternatives to explore.

ON ORDER OF OF COURSE, WE ARE STILL LEFT WITH THE PROBLEM:
WAVELENGTH .4 to .7μm     MV ≈ .3μm
OF LIGHT!   HOW DO WE MAKE SYSTEMS THIS SMALL?