# LG8: (E) Engineering and Physical Considerations

Topics: Power consumption, scaling, size, logical effort and performance limits.

LG8.1 - E - 90 Nanometer Gate Length.
LG8.2 - E - Power Consumption
LG8.3 - E - Dynamic Power Gating
LG8.4 - E - Dynamic Frequency Scaling
LG8.5 - E - Dynamic Voltage Scaling
LG8.6 - E - Logical Effort
LG8.7 - E - Information Flux

What is the *Silicon End Point* ?

# LG8.1 - E - 90 Nanometer Gate Length.

The mainstream VLSI technology in the period 2004-2008.

Parameters from a 90 nanometer standard cell library:

| Parameter | Value | Unit |
|---|---|---|
| Drawn Gate Length | 0.08 | $\mu$m |
| Metal Layers | 6 to 9 | layers |
| Max Gate Density | 400K | gates/mm$^2$ |
| Track Width | 0.25 | $\mu$m |
| Track Spacing | 0.25 | $\mu$m |
| Tracking Capacitance | 1 | fF/mm |
| Core Supply Voltage | 0.9 to 1.4 | V |
| FO4 Delay | 51 | ps |

Typical processor core: 200k gates + 4 RAMs: one square millimeter.

Typical SoC chip area is 50-100 mm$^2$ ⇝ 20-40 million gates.

Actual gate and transistor count higher owing to custom blocks (RAMs mainly).

Now the industry is moving to 45 nanometer.

2007: Dual-core Intel Itanium2: 2 billion transistors.

http://en.wikipedia.org/wiki/Moore's_law

## LG8.2- E - Power Consumption

$P = V \times I = E \times f$

I (current) = Static Current + Dynamic Current.

Early CMOS (VCC= 5 volts): negligible static current. Today it's 30 % of dynamic.

Dynamic current = Short circuit current + Charge current.

Charge current:

- All energy in a net/gates is wasted each time it toggles.

- The energy in a capacitor is $E = CV^2/2$.

- Dominant capacitance is proportional to net length.

- Gate input and output capacitance also contribute to $C$.

Activity ratio, $a$: percentage of clock cycles that see a transition.

The net toggle rate = Operating frequency of the chip $(f) \times a$.

Short circuit power: see cell library.

`http://cis.poly.edu/cs2214rvs/powers03.htm`.

## LG8.2a – E – Power Consumption Example

Example: core area 64 mm$^2$; average net length 0.1 mm; 400K gates/mm$^2$, $a = 0.25$.

Net capacitance $= 0.1$ mm $\times$ 1 fF/mm $\times$ 400K $\times$ 64 mm$^2$ $= 2.5$ nF.

Compare 1.35V to 1.8V: twice the power and twice the clock frequency?

| Vcc Volts | Freq MHz | Static Power mW | Dynamic Power mW | Total Power mW |
|:---:|:---:|:---:|:---:|:---:|
| 0.8 | 100 | 40 | 24 | 64 |
| 1.35 | 100 | 67 | 68 | 135 |
| 1.35 | 200 | 67 | 136 | 204 |
| 1.8 | 100 | 90 | 121 | 211 |
| 1.8 | 200 | 90 | 243 | 333 |
| 1.8 | 400 | 90 | 486 | 576 |

In the past we were often *core-bound* or *pad-bound*.

Today's VLSI designs are commonly *power-bound*.

- 1 W/cm$^2$ can be dissipated from a plastic package.

- 2-4 W/cm$^2$ required a heat sink.

- More than 8 W/cm$^2$ required forced cooling.

Workstation microprocessors dissipate tens of Watts: hence cooling fans.

# LG8.3 - E - Dynamic Power Gating

Previously looked at dynamic clock gating: can also turn off power (coarser grain).

Use power gating cells in series with supply rails.

Use signal isolation and retention cells (t-latches) on data I/O nets.

No register and RAM data retention in block while off.

Requires some sequencing: several clock cycles to power up/down a region.

Generally, power of/on controller by software or top-level input pads.

Sometimes power off a whole chip except for a one or two RAMs and register files.

Can also retain contents using a lower supply (CMOS RAM data holding voltage).

# LG8.4 – E – Dynamic Frequency Scaling

Let's adjust the clock frquency (while keeping VCC constant for now).

Does frequency scaling help ?

Frequency scaling is software controlled: update PLL division ratio. PLL has inertia: e.g. 1 millisecond.

Let's compare with dynamic clock gating:

|  | Clock Gating. | Frequency Adjustment. |
|---|---|---|
| Control: | automatic, | manual. |
| Granularity: | register / FSM, | macroscopic. |
| Clock Tree: | mostly free runs, | slows down. |
| Response time: | instant, | acceptable. |
| Can vary voltage: | no, | yes. |

To compute quickly and halt we need a higher frequency clock but consume the same number of active cycles.

Work rate product, $af$ unchanged: so no power difference ?

Actually un-stopped regions consume power proportional to $f$.

Zeno: Tortoise and Achilles ?

Tortoise is best: keep going steadily and end just in time. (He becomes even righter when we vary the voltage.)

Clock gating still good for: bursty, localised activity.

# LG8.5 – E – Dynamic Voltage Scaling

Logic with higher-speed capabilities is smaller which means it consumes greater leakage current which is wasted even while we are halted.

CMOS delay is inversely proportional to supply voltage.

Voltage to a region may be varied dynamically. A higher supply voltage uses more power (square law) but allows a higher $f$.

Operating region of the frequency/voltage curve is roughly linear.

Overall: power has cubic dependence on $f$.

Let's only raise VCC when we ramp up $f$.

Still obtain peak performance under heavy loads: avoid cubic overhead when idle.

Method:

1. Adjust $f$ for just-in-time completion (e.g. in time to decode the next frame of a real-time video),

2. then adjust VCC so logic just works.

But zeno applies still: always aim for $a$ close to unity and a low work rate.

Large combinational blocks: can dip power supply to reduce static current when block is completely idle (detect with XORs): need to retain register state. NB: combinational logic cannot be clock gated (e.g. PAL and PLA).

# LG8.6 - E - Logical Effort

When sending a signal a long distance over a chip:

- can use powerful drivers and long nets,

- or normal drivers and repeaters made of asynchronous buffers,

- or normal drivers and registered repeaters (network-on-chip).

When we compute a function needed far away, should we spread out the gates?

When we compute a function used in several places, when is it cheaper to re-compute locally ? NB: most of the power use is driving nets.

`http://en.wikipedia.org/wiki/Logical_effort`

`http://www.cl.cam.ac.uk/teaching/0708/VLSI/logicaleffort.html`

It can be shown that in multistage logic networks, the minimum possible delay along a particular path can be achieved by designing the circuit such that the stage logical efforts are equal. For a given combination of gates and a known load, B, G, and H are all fixed causing F to be fixed; hence the individual gates should be sized such that the individual stage efforts are

$$f = F^{1/N}$$

where N is the number of stages in the circuit.

# LG8.7 – E – Information Flux

It is interesting to compare maximum bit capacity per square millimeter of per second per Joule for optics and electronics. What are the theoretical limits? Where is current technology ? There's also a thermodynamic limit on the energy cost of deleting a bit of information.... Interested: apply for our Masters or PhD....

END OF DOCUMENT.