

## VLSI for Architects

Krste Asanović

Computer Architecture Group

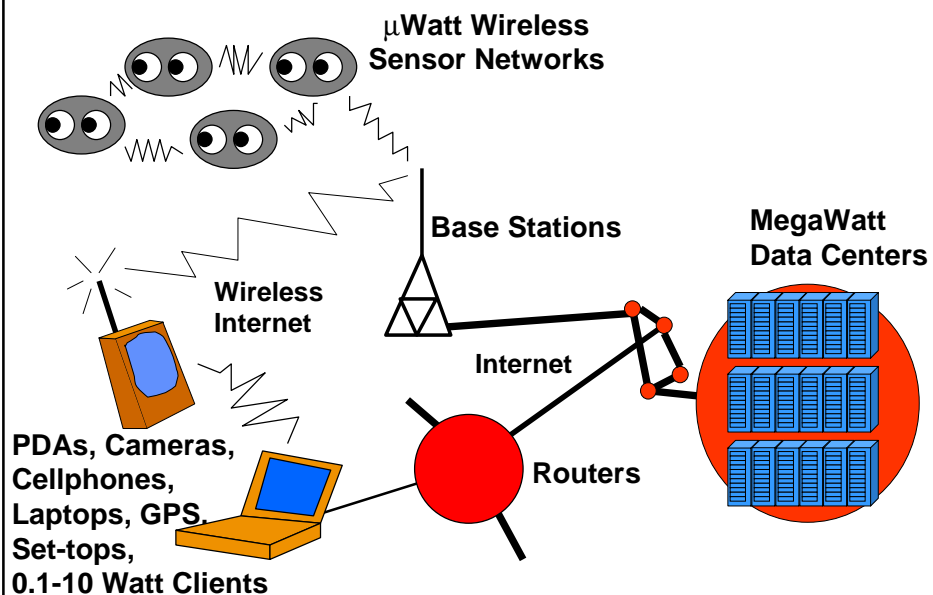
MIT Laboratory for Computer Science

krste@lcs.mit.edu

<http://www.cag.lcs.mit.edu/6.893-f2000/>

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 1. © Krste Asanović

## Future Computing Infrastructure



6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 2. © Krste Asanović

## Semiconductor Trends

- Non-Recurring Engineering (NRE) costs are increasing rapidly for new designs
  - >\$1M for masks to spin a new design
  - Engineers cost ~\$200K/year (salary+benefits+overhead)
  - Pentium Pro design verification took around 350 engineer years or ~\$70M

⇒ *Tremendous economies of scale*

(Can't sell <1,000,000 parts for <\$100 each)

- CMOS following Moore's Law until (at least) 2011-2014

- ITRS'99' roadmap 2011, 50nm technology
  - 64 Gb DRAMs (8 GB/chip)
  - 7 billion transistor CPUs
  - 10 GHz clocks (100 ps cycle time)

⇒ *Smallest viable chips have huge capacity*

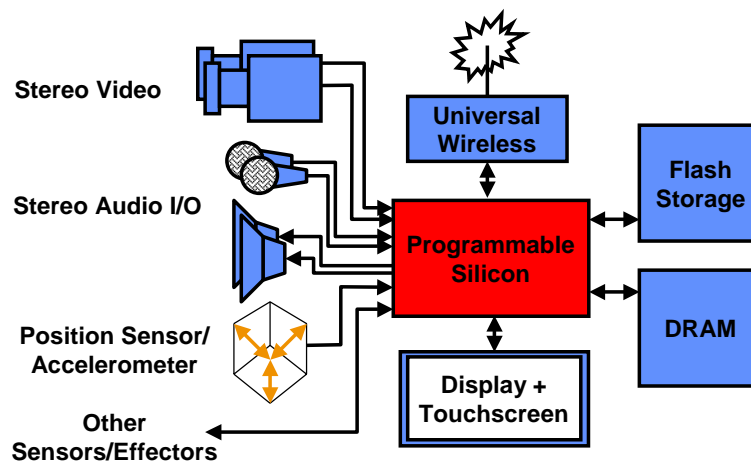
(~10 million transistors/mm<sup>2</sup>)

10 million transistors per person per day

[\*International Technology Roadmap for Semiconductors]

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 3. © Krste Asanović

## Programmable Silicon Replaces Custom Hardware



Programmable silicon replaces ASICs, or collections of DSPs, microprocessors and glue logic

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 4. © Krste Asanović

## Benchmarks & Metrics

- **Application space wider than desktop processors**
  - Benchmark as many applications as possible
  - Include apps done with special hardware now (graphics, audio, crypto)
  - Whole system measures
  - Real-time important
- **Primary metrics**
  - Cost (related to die area but also whole system cost)
  - Execution Time (latency and throughput, average and worst-case)
  - Energy (also peak power and peak switching current)
- **Compare against best possible solution for each application**
  - How much worse than application-specific circuitry?
  - Moore's law perhaps makes area the most forgiving dimension
    - try to keep energy and delay competitive, possibly at expense of area

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 5. © Krste Asanović

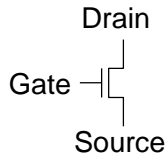
## VLSI for Architects

Two types of question architects ask:

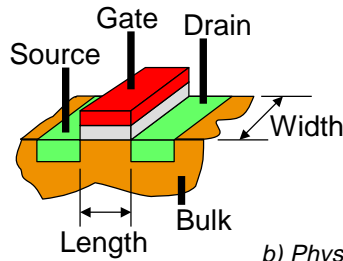
- How will this change affect area/delay/energy in current technology?
- How will this design scale to future technologies?
  
- For next 10-15 years, *the* technology is CMOS

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 6. © Krste Asanović

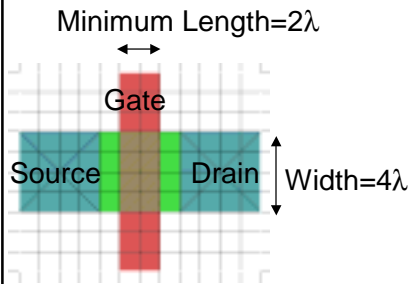
## Transistors



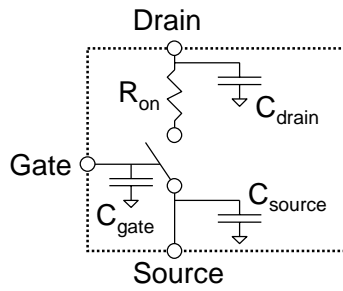
a) Circuit Symbol



b) Physical Realization

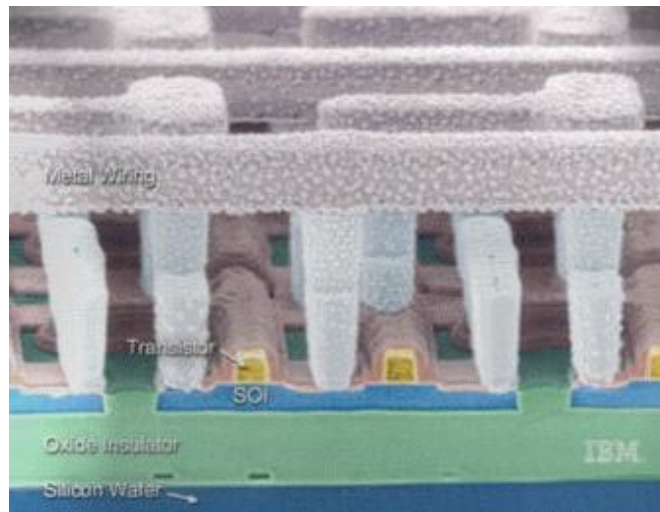


c) Layout View



d) Simple RC Model

## Transistors

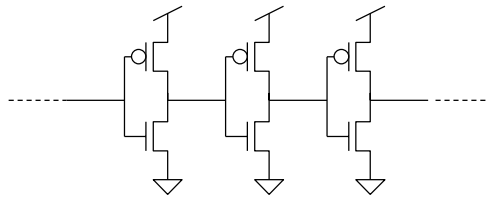


IBM SOI Technology

## Method of Logical Effort (Sutherland and Sproul)

- Easy way to estimate delays in CMOS process
- Indicates correct number of logic stages to use and transistor sizes
- Characterize process speed with single delay parameter:  $\tau$

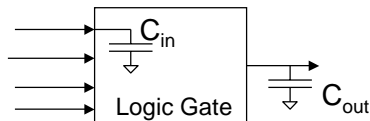
$\tau$ , delay of inverter driving same-sized inverter (no parasitics)



$\tau$  in range 10-15ps for 0.18 $\mu$ m processes

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 9. © Krste Asanović

## Gate Delay Components

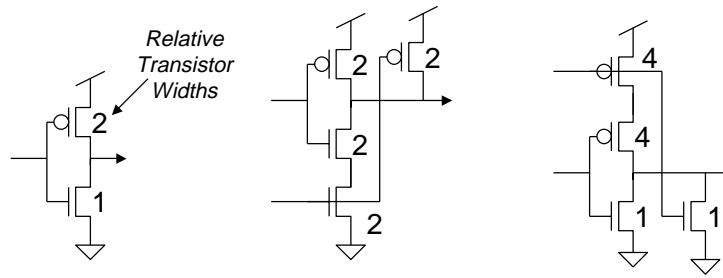


- Split delay of logic gate into three components  
Delay = Logical Effort x Electrical Effort + Parasitic Delay
- Logical Effort
  - Complexity of logic function (Invert, NAND, NOR, etc)
  - Define inverter has logical effort = 1
  - Depends only on topology not transistor sizing
- Electrical Effort
  - Ratio of output capacitance to input capacitance  $C_{out}/C_{in}$
- Parasitic Delay
  - Intrinsic self-loading of gate
  - Independent of transistor sizes and output load

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 10. © Krste Asanović

## Logical Effort for Simple Gates

- Define Logical Effort of Inverter = 1
- For other gates, size to give same current drive as inverter
- Logical Effort is ratio of logic gate's input cap. to inverter's input cap.



**Inverter**

**NAND**

**NOR**

Input Cap = 3 units

Input Cap = 4 units

Input Cap = 5 units

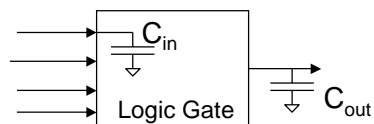
L.E.=1 (definition)

L.E.=4/3

L.E.=5/3

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 11. © Krste Asanović

## Electrical Effort



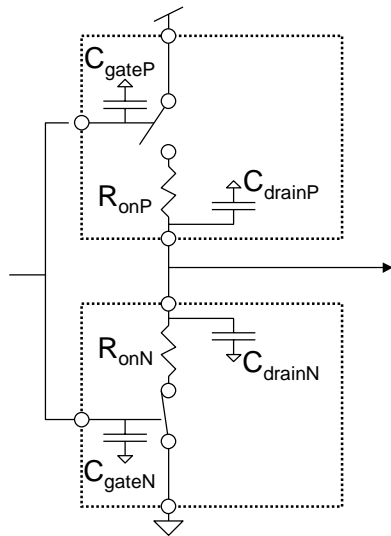
- Ratio of output load capacitance over input capacitance:

$$E.E. = C_{out}/C_{in}$$

- Usually, transistors have minimum length
- Input and output capacitances can be measured in units of transistor gate widths

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 12. © Krste Asanović

## Parasitic Delay



- Main cause is drain capacitances
- These scale with transistor width so P.D. independent of transistor sizes

- Useful approximation:

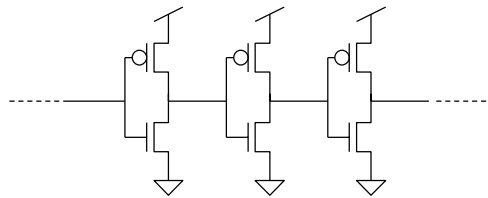
$$C_{gate} \approx C_{drain}$$

- For inverter:

$$\text{Parasitic Delay} \approx 1.0 \tau$$

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 13. © Krste Asanović

## Inverter Chain Delay



- For each stage:

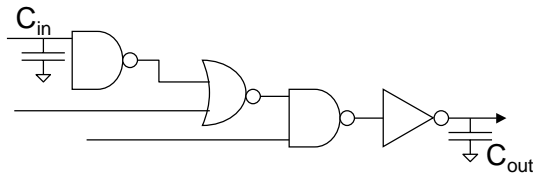
$$\text{Delay} = \text{Logical Effort} \times \text{Electrical Effort} + \text{Parasitic Delay}$$

$$= 1.0 \text{ (definition)} \times 1.0 \text{ (in = out)} + 1.0 \text{ (drain caps)}$$

$$= 2.0 \text{ units}$$

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 14. © Krste Asanović

## Optimizing Circuit Paths



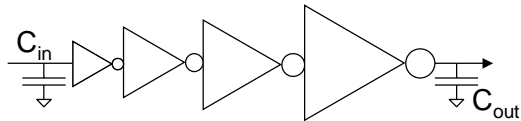
- Path logical effort,  $G = \prod g_i$  ( $g_i = \text{L.E. stage } i$ )
- Path electrical effort,  $H = C_{out}/C_{in}$  ( $h_i = \text{E.E. stage } i$ )
- Parasitic delay,  $P = \sum p_i$  ( $p_i = \text{P.D. stage } i$ )
- Path effort,  $F = GH$
- Minimum delay when each of  $N$  stages has equal effort

$$\text{Min. } D = NF^{1/N} + P$$

$$\text{i.e. } g_i h_i = F^{1/N}$$

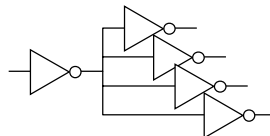
6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 15. © Krste Asanović

## Optimal Number of Stages



- Minimum delay when:
  - stage effort = logical effort x electrical effort  $\approx 3.4-3.8$
  - Some derivations have  $e = 2.718..$  as best stage effort – this ignores parasitics
  - Broad optimum, stage efforts of 2.4-6.0 within 15-20% of minimum
- Fan-out-of-four (FO4) is convenient design size ( $\sim 5\tau$ )

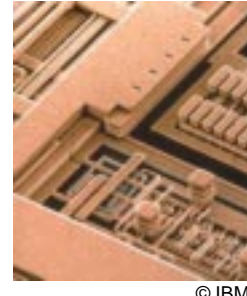
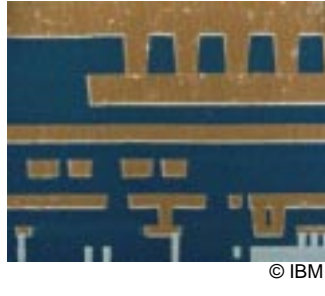
*FO4 delay: Delay of inverter driving four copies of itself*



6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 16. © Krste Asanović



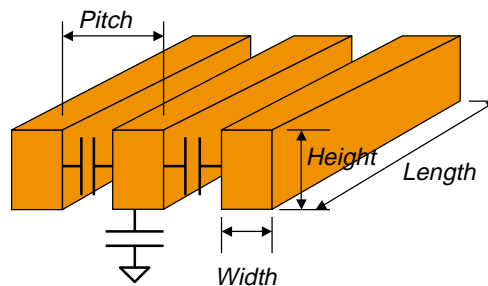
## Wires



IBM CMOS7  
process  
6 layers of  
copper wiring

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 17. © Krste Asanović

## Wires



- Resistance fixed by  $(\text{length} \cdot \text{resistivity}) / (\text{height} \cdot \text{width})$ 
  - bulk aluminum  $2.8 \mu\Omega\text{-cm}$ , bulk copper  $1.7 \mu\Omega\text{-cm}$
- Capacitance depends on geometry of surrounding wires and relative permittivity,  $\epsilon_r$ , of dielectric
  - silicon dioxide  $\epsilon_r = 3.9$ , new low-k dielectrics in range 1.2-3.1

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 18. © Krste Asanović

## Current Interconnect Densities

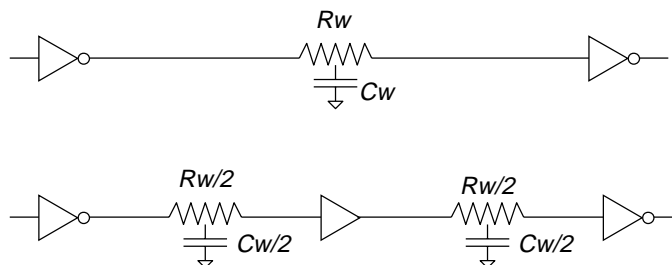
- Intel Pentium-III, 0.18 $\mu\text{m}$ , 6 aluminum layers, SiOF dielectric ( $\epsilon_r = 3.1$ )

Metal Layer	Pitch ( $\mu\text{m}$ )	Aspect Ratio (Height/Width)
<b>M1</b>	<b>0.50</b>	<b>1.9</b>
<b>M2</b>	<b>0.64</b>	<b>2.2</b>
<b>M3</b>	<b>0.64</b>	<b>2.2</b>
<b>M4</b>	<b>1.08</b>	<b>2.0</b>
<b>M5</b>	<b>1.60</b>	<b>2.0</b>
<b>M6</b>	<b>1.76</b>	<b>2.0</b>

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 19. © Krste Asanović

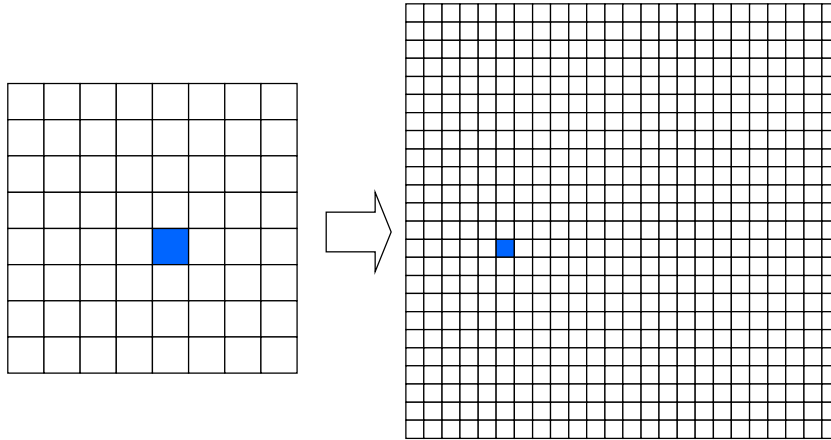
## Wire Delays

- Resistance,  $R$ , increases per unit length
  - in 0.25 $\mu\text{m}$  CMOS,  $\sim 1000\lambda$  thin M1 wire = minimum inverter resistance
- Capacitance,  $C$ , increases per unit length
  - in 0.25 $\mu\text{m}$  CMOS,  $\sim 1000\lambda$  thin M1 wire = minimum inverter capacitance
- Wire delay increases as  $RC$ , quadratic in length
  - in 0.25 $\mu\text{m}$ ,  $\sim 1000\lambda$  thin M1 wire = 30ps ( $\sim \mathcal{T}$ )
- Inserting repeaters makes delay linear with length



6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 20. © Krste Asanović

## Scaling



- Scale linear dimensions by factor  $S$  (around 0.7 / generation)
- Chip size also increases

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 21. © Krste Asanović

## Scaling Slides from Horowitz DAC'2000 Talk

[http://www.dac.com/37slides/05\\_2.ppt](http://www.dac.com/37slides/05_2.ppt)

(link on class web page)

6.893: Advanced VLSI Computer Architecture, September 12, 2000, Lecture 2, Slide 22. © Krste Asanović