Digital Integrated Circuits (83-313)

Lecture 9: Interconnect

Semester B, 2016-17

Lecturer: Dr. Adam Teman

Itamar Levi, Robert Giterman

23 May 2017



Emerging Nanoscaled Integrated Circuits and Systems Labs

EnICS

TAs:

Disclaimer: This course was prepared, in its entirety, by Adam Teman. Many materials were copied from sources freely available on the internet. When possible, these sources have been cited; however, some references may have been cited incorrectly or overlooked. If you feel that a picture, graph, or code example has been copied from you and either needs to be cited or removed, please feel free to email <u>adam.teman@biu.ac.il</u> and I will address this as soon as possible.

Lecture Content





A First Glance at Interconnect









schematic view

physical realization

All-inclusive model

Capacitance-only



Impact of Interconnect Parasitics

- Interconnect parasitics affect all the metrics we care about
 - Reliability
 - Performance
 - Power Consumption
 - Cost

Classes of parasitics

- Capacitive
- Resistive
- Inductive





Modern Interconnect





Capacitance





Capacitance of Wire Interconnect





Capacitance: The Parallel Plate Model

• How can we reduce this capacitance?



Typical numbers:

- Wire cap ~0.2 fF/um
- Gate cap ~2 fF/um
- Diffusion cap ~2 fF/um

Electrical-field lines

 $\frac{\varepsilon_{di}}{t}WL$ C_{pp}



Permittivity

Material	ε _r
Free space	1
Aerogels	~1.5
Polyimides (organic)	3-4
Silicon dioxide	3.9
Glass-epoxy (PC board)	5
Silicon Nitride (Si ₃ N ₄)	7.5
Alumina (package)	9.5
Silicon	11.7



Fringing Capacitance



Η





W - H/2

Fringing versus Parallel Plate



A simple model for deriving wire cap

• Wiring capacitances in $0.25 \mu m$

$$C_{wire} = C_{parallel_plate} \cdot W \cdot L$$



Impact of Interwire Capacitance





Coupling Capacitance and Delay



 $C_{tot} = C_L$



Coupling Capacitance and Delay



 $C_{tot} = C_L + C_{C1} + C_{C2}$



Coupling Capacitance and Delay



 $C_{tot} = C_L + 2(C_{C1} + C_{C2})$



Example – Coupling Cap

- A pair of wires, each with a capacitance to ground of 5pF, have a 1pF coupling capacitance between them.
- A square pulse of 1.8V (relative to ground) is connected to one of the wires.
- How high will the noise pulse be on the other wire?



Example – Coupling Cap

• Draw an Equivalent Circuit:



$$V_{C2} = \frac{V_{in} \cdot C_{coupled}}{C_{coupled} + C_2} = \frac{1.8 \cdot 1p}{1p + 5p} = 0.3V$$



Coupling Waveforms

• Simulated coupling for $C_{agg} = C_{victim}$







Feedthrough Cap



Measuring Capacitance



Resistance





Wire Resistance



Metal	Bulk resistivity (μΩ*cm)
Silver (Ag)	1.6
Copper (Cu)	1.7
Gold (Au)	2.2
Aluminum (Al)	2.8
Tungsten (W)	5.3
Molybdenum (Mo)	5.3



Sheet Resistance

Typical sheet resistances for 180nm process

Layer	Sheet Resistance (Ω/\Box)
N-Well/P-Well	1000-1500
Diffusion (silicided)	3-10
Diffusion (no silicide)	50-200
Polysilicon (silicided)	3-10
Polysilicon (no silicide)	50-400
Metal1	0.08
Metal2	0.05
Metal3	0.05
Metal4	0.03
Metal5	0.02
Metal6	0.02



Silicides: WSi _{2,} TiSi ₂, PtSi ₂ and TaSi

Conductivity: 8-10 times better than Poly



Contact Resistance

Contact/Vias add extra resistance

- Similar to changing between roads on the way to a destination...
- Contact resistance is generally 2-20 Ω



Make contacts bigger

- BUT... current "crowds" around the perimeter of a contact.
- There are also problems in deposition...
- Contacts/Vias have a maximum practical size.
- Use multiple contacts
 - But does this add overlap capacitance?



Dealing with Resistance

- Selective Technology Scaling
 - Don't scale the H
- Use Better Interconnect Materials
 - reduce average wire-length
 - e.g. copper, silicides
- More Interconnect Layers
 - reduce average wire-length
- Minimize Contact Resistance
 - Use single layer routing
 - When changing layers, use lots of contacts.



90nm Process



Interconnect Modeling





The Ideal Model

- In schematics, a wire has no parasitics:
 - The wire is a single equipotential region.
 - No effect on circuit behavior.
 - Effective in first stages of design and for very short wires.



The Lumped Model





The Distributed RC-line

• But actually, our wire is a distributed entity.

• We can find its behavior by breaking it up into small RC segments.



Step-response of RC wire

• Step-response of RC wire as a function of time and space





Elmore Delay Approximation

- Solving the diffusion equation for a given network is complex.
- Elmore proposed a reasonably accurate method to achieve an approximation of the dominate pole.





Elmore Delay Approximation

For a complex network use the following method:

• Find all the resistors on the path from in to out.

• For every capacitor:

- Find all the resistors on the path from the input to the capacitor.
- Multiply the capacitance by the resistors that are also on the path to out.
- The dominant pole is approximately the sum of all these time constants.

$$R_{ik} = \sum R_j \Rightarrow (R_j \in [path(s \to i) \cap path(s \to k)]) \quad \tau_{Di} = \sum_{k=1}^{k} C_k R_{ik}$$

Simple Elmore Delay Example



$$\tau_{elmore} = R_1 C_1 + (R_1 + R_2) C_2 + (R_1) C_2$$



General Elmore Delay Example



 $\tau_{elmore} = R_1 C_1 + R_1 C_2 + (R_1 + R_3) C_3 + (R_1 + R_3) C_4 + (R_1 + R_3 + R_i) C_i$



Generalized Ladder Chain

- Lets apply the Elmore approximation for our original distributed wire.
 - Divide the wire into N equal segments of dx=L/N length with capacitance cdx and resistance rdx.

$$\begin{aligned} \tau_N &= c \left(\frac{L}{N}\right) \left(r\frac{L}{N} + 2r\frac{L}{N} + ... + Nr\frac{L}{N}\right) \\ &= \left(\frac{L}{N}\right)^2 \left(rc + 2rc + ... + Nrc\right) = rcL^2 \left(\frac{N(N+1)}{2N^2}\right) \\ &\lim_{N \to \infty} \tau_D = \frac{rcL^2}{2} = \frac{RC}{2} \end{aligned}$$



RC-Models

Voltage Range	Lumped RC- network	Distributed RC-network
$0 \rightarrow 50\%$ (t _p)	0.69 RC	0.38 RC
0→63% (7)	RC	0.5 RC
10% \rightarrow 90% (t _r)	2.2 RC	0.9 RC

Step Response of Lumped and Distributed RC Networks:

Points of Interest.



EnICS

Wire Delay Example

• Inverter driving a wire and a load cap.

$$\tau_{driver} = \left(C_d + \frac{C_w}{2} \right) R_{inv} + \left(C_{ext} + \frac{C_w}{2} \right) \left(R_{inv} + R_w \right)$$



A different look...

- Again we'll look at our driver with a distributed wire.
 - For the driver resistance, we can lump the output load as a capacitor.
 - For the wire resistance, we will use the distributed time constant.
 - For the load capacitance, we can lump the wire and driver resistance.

$$C_{w} \approx 0.2 \frac{\text{fF}}{\mu \text{m}}$$
$$R_{\Box} \approx 0.1 \frac{\Omega}{\Box}$$

$$\tau_{D} = 0.69R_{inv} \left(C_{d} + C_{W} \right) + 0.38R_{W}C_{W} + 0.69 \left(R_{inv} + R_{W} \right)C_{L}$$



Dealing with long wires

Repeater Insertion



Dealing with long wires

Buffer Tree Insertion



Wire Scaling





Wire Scaling

- We could try to scale interconnect at the same rate (S) as device dimensions.
 - This makes sense for *local wires* that connect smaller devices/gates.
 - But *global interconnections*, such as clock signals, buses, etc., won't scale in length.
- Length of global interconnect is proportional to die size or system complexity.
 - Die Size has increased by 6% per year (X2 @10 years)
 - Devices have scaled, but complexity has grown!



Nature of Interconnect







Local Wire Scaling

- Looking at local interconnect:
 - W, H, t, L all scale at 1/S
 - $C = LW/t \rightarrow 1/S$
 - $R = L/WH \rightarrow S$
 - *RC*=1



- Reminder Full Scaling of Transistors
 - $R_{\rm on} = V_{\rm DD} / I_{\rm on} \alpha 1$
 - $t_{\rm pd} = R_{\rm on} C_{\rm g} \alpha 1/S$

But the delay of local interconnect increases relative to transistors!

Local Wire Scaling – Full Scaling

• What about fringe cap?

48



Local Wire Scaling - Constant Thickness

• Wire thickness (height) wasn't scaled!



Local Wire Scaling – Interwire Capacitance

- Without scaling height, coupling gets much worse.
- Aspect ratio is limited and we eventually have to scale the height.
- Therefore, different metal layers have different heights.



Global Wire Scaling

- Looking at global interconnect:
 - W, H, t scale at 1/S
 - L doesn't scale!
 - $C = LW/t \rightarrow 1$
 - $R = L/WH \rightarrow S^2$
 - $RC = S^2 |||$

Long wire delay increases



• And if chip size grows, *L* actually increases!



Global Wire Scaling – Constant Thickness

- Leave thickness constant for global wires
- But wire delay still gets quadratically worse than gate delay...



Wire Scaling

- So whereas device speed increases with scaling:
 - Local interconnect speed stays constant.
 - Global interconnect delays increase quadratically.
- Therefore:
 - Interconnect delay is often the limiting factor for speed.

• What can we do?

- Keep the wire thickness (*H*) fixed.
- This would provide 1/S for local wire delays and S for constant length global wires.
- But fringing capacitance increases, so this is optimistic.



Wire Scaling

- What is done today?
 - Low resistance metals.
 - Low-*K* insulation.
 - Low metals (*M1*, *M2*) are used for local interconnect, so they are thin and dense.
 - Higher metals are used for global routing, so they are thicker, wider and spaced farther apart.



Modern Interconnect



1 µm

Cross Sectional View: For Height, Width and Spacing

[Moon08]



Further Reading

- J. Rabaey, "Digital Integrated Circuits" 2003, Chapter 4
- E. Alon, Berkeley *EE-141*, Lectures 15,16 (Fall 2009) http://bwrc.eecs.berkeley.edu/classes/icdesign/ee141_f09/
- B. Nicolic, Berkeley EE-241, Lecture 3 (Spring 2011) http://bwrc.eecs.berkeley.edu/classes/icdesign/ee241_s11
- Stanford EE311

