

Jacob Abraham

Department of Electrical and Computer Engineering The University of Texas at Austin

> VLSI Design Fall 2020

November 10, 2020

Lecture 21. Skews, Scaling

Jacob Abraham, November 10, 2020 1 / 48

Jacob Abraham, November 10, 2020 1 / 48

Clock Distribution

ECE Department, University of Texas at Austin

ECE Department, University of Texas at Austin



- On practical chips, the RC delay of the wire resistance and gate load is very long
 - Variations in this delay cause clock to get to different elements at different times
 - This is called clock skew
- Most chips use repeaters to buffer the clock and equalize the delay

Lecture 21. Skews, Scaling

Reduces but doesn't eliminate skew



Review: Skew Impact





ECE Department, University of Texas at Austi



ob Abraham, November 10, 2020 4 / 48

Dynamic Circuit Review

- Static circuits are slow because fat pMOS load input
- Dynamic gates use precharge to remove pMOS transistors from the inputs
 - Precharge: $\phi = 0$, output forced high
 - Evaluate: $\phi = 1$, output may pull low



Domino Circuits

ECE Department, University of Texas at Austin

Dynamic inputs must monotonically rise during evaluation

Jacob Abraham, November 10, 2020 6 / 48

- Place inverting stage between each dynamic gate
- Dynamic/static pair called domino gate
- Domino gates can be safely cascaded





Traditional Domino Circuits

ECE Department, University of Texas at Aust



Jacob Abraham, November 10, 2020 8 / 48

Clock Skew

- Skew increases sequencing overhead
 - Traditional domino has hard edges
 - Evaluate at latest rising edge
 - Setup at latch by earliest falling edge



 $t_{pd} = T_c - 2t_{pdq} - 2t_{skew}$

Time Borrowing

ECE Department, University of Texas at Austin

- Logic may not exactly fit half-cycle
 - No flexibility to borrow time to balance logic between half cycles
 - Traditional domino sequencing overhead is about 25% of cycle time in fast systems!



m, November 10, 2020 10 / 48



Skew-Tolerant Domino

ECE Department, University of Texas at Austin



7

Jacob Abraham, November 10, 2020 12 / 48

Full Keeper

ECE Department, University of Texas at a

- After second phase evaluates, first phase precharges
- Input to second phase falls
 - Violates monotonicity?
- But we no longer need the value
- Now the second gate has a floating output
 - Need full keeper to hold it either high or low





Abraham, November 10, 2020 14 / 48

Multiple Phases

ECE Department, University of Texas at Aus

With more clock phases, each phase overlaps more
Permits more skew tolerance and time borrowing





Department of Electrical and Computer Engineering, The University of Texas at Austin J. A. Abraham, November 10, 2020

b Abraham, November 10, 2020 16 / 48

Opportunistic Time Borrowing

U. S. Patent no. 5517136 (Harris et al., May 14, 1996, assigned to Intel Corporation)

Pipelined domino logic allowing a slow stage to "borrow" from the time normally allocated to a faster stage



Clocking of Time-Borrowing Pipeline

ECE Department, University of Texas at A

- Delayed falling edges on clocks allow evaluation to continue into subsequent half cycle
 - Time delay t_d should be greater than of equal to the hold time of the domino logic gate plus any global clock skew

Jacob Abraham, November 10, 2020 18 / 48

• Can generate the clocks by a local reference driven by the chip's global reference clock signal





CLE

DCLK

D1

LATCH

ANYTHING

Half-cycles 1 and 3 evaluate when CLK is high, half-cycle 2 when CLK is low

ECE Department, University of Texas at





















Chip Densities increase with Scaling

- In 1965, Gordon Moore predicted the exponential growth of the number of transistors on an IC (Moore's Law)
- Transistor count doubled every year since invention
- Predicted > 65,000 transistors by 1975!
- Growth limited by power





Scaling

- The only constant in VLSI is constant change
- Feature size shrinks by 30% every 2-3 years
 - Transistors become cheaper, and faster
 - Wires do not improve (and may get worse)
- Scale factor S (typical technology nodes) $S = \sqrt{2}$





ECE Department, University of Texas at Austin

Scaling			
Table 4.15 Influence of scaling o	n MOS device	characteris	tics
Parameter	Sensitivity	Constant Field	Lateral
Scaling	g Parameters		
Length: L		1/S	1/S
Width: W		1/S	1
Gate oxide thickness: t_{ox}		1/S	1
Supply voltage: V _{DD}		1/S	1
Threshold voltage: V_{tn} , V_{tp}		1/S	1
Substrate doping: N_A		S	1
Device C	Characteristics		
β	$\frac{W}{L}\frac{1}{t_{\text{ox}}}$	S	S
Current: I_{di}	$\beta \left(V_{DD} - V_t \right)^2$	1/S	S
Resistance: R	$\frac{V_{DD}}{I_{ds}}$	1	1/8
Gate capacitance: <i>C</i>	$\frac{WL}{t_{\text{ox}}}$	1/S	1/8
Gate delay: τ	RC	1/S	$1/S^{2}$
Clock frequency: f	1/τ	S	S^2
Dynamic power dissipation (per gate): P	CV^2f	$1/S^{2}$	S
Chip area: A		$1/S^{2}$	1
Power density	P/A	1	S
Current density	I_{d}/A	S	S

Lecture 21. Skews, Scaling Jacob Abraham, November 10, 2020 32 / 48

Observations

- Gate capacitance per micron is nearly independent of process
- But ON resistance × micron improves with process
- Gates get faster with scaling (good)
- Dynamic power goes down with scaling (good)
- Current density goes up with scaling (bad)
- Velocity saturation makes lateral scaling unsustainable

Solution

ECE Department, University of Texas at Austin

- Gate capacitance is typically about 2 fF/ μ m
- The FO4 inverter delay in the TT corner for a process of feature size f (in nm) is about 0.5f ps
- Estimate the ON resistance of a unit $(4/2 \lambda)$ transistor

Lecture 21. Skews, Scal

- FO4 = 5 τ = 15 RC
- RC = (0.5f)/15 = (f/30) ps/nm
- If W = 2f, R = 8.33 k Ω

Unit resistance is roughly independent of f



r 10, 2020 34 / 48

Interconnect Scaling

Parameter	Sensitivity	Reduced Thickness	Constant Thickness	
Scaling Pa	rameters			
Width: w		1/S		
Spacing: s		1/S		
Thickness: t		1/S	1	
Interlayer oxide height: h		1/S		
Characteristics Per Unit Length				
Wire resistance per unit length: R_w	$\frac{1}{wt}$	S^2	S	
Fringing capacitance per unit length: $C_{\!w\!f}$	$\frac{t}{s}$	1	S	
Parallel plate capacitance per unit length: $C_{\rm tep}$	$\frac{w}{b}$	1	1	
Total wire capacitance per unit length: $C_{\!w}$	C_{wf} + C_{wp}	1	between 1, S	
Unrepeated RC constant per unit length: t _{ww}	$R_w C_w$	S^2	between S, S ²	
Repeated wire RC delay per unit length: t_{wr} (assuming constant field scaling of gates in Table 4.15)	$\sqrt{RCR_wC_w}$	\sqrt{s}	between 1, \sqrt{S}	
Crosstalk noise	$\frac{t}{s}$	1	S	

Interconnect Delay

ECE Department, University of Texas at Austin

ECE Department, University of Texas at Austin

Parameter	Sensitivity	Reduced Thickness	Constant Thickness	
S	caling Parameters			
Width: w		1/S		
Spacing: s		1/S		
Thickness: t		1/S	1	
Interlayer oxide height: h			1/S	
Local/Scaled Interconnect Characteri	stics			
Length: /		1/S		
Unrepeated wire RC delay	$l^2 t_{wu}$	1 betweer 1/ <i>S</i> , 1		
Repeated wire delay	lt _{wr}	$\sqrt{1/S}$ betwee $1/S$, $\sqrt{1/S}$		
Global Interconnect Characteristics				
Length: /			D_c	
Unrepeated wire RC delay	l ² t _{wu}	$S^2 D_c^2$	between SD_c^2 , $S^2D_c^2$	
Repeated wire delay	lt _{wr}	$D_c \sqrt{S}$	between D_c	

Jacob Abraham, November 10, 2020 36 / 48

Jacob Abraham, November 10, 2020 37 / 48

Lecture 21. Sk



Previously: International Technology Roadmap for Semiconductors (ITRS)

ECE Department, University of Texas at Austin

ECE Department, University of Texas at Austin

Table 4.17 Predictions from the 2002 ITRS								
Year	2001	2004	2007	2010	2013	2016		
Feature size (nm)	130	90	65	45	32	22		
$V_{DD}(\mathbf{V})$	1.1-1.2	1-1.2	0.7–1.1	0.6-1.0	0.5–0.9	0.4–0.9		
Millions of transistors/die	193	385	773	1564	3092	6184		
Wiring levels	8-10	9-13	10-14	10-14	11-15	11-15		
Intermediate wire pitch (nm)	450	275	195	135	95	65		
Interconnect dielectric	3-3.6	2.6-3.1	2.3-2.7	2.1	1.9	1.8		
constant								
I/O signals	1024	1024	1024	1280	1408	1472		
Clock rate (MHz)	1684	3990	6739	11511	19348	28751		
FO4 delays/cycle	13.7	8.4	6.8	5.8	4.8	4.7		
Maximum power (W)	130	160	190	218	251	288		
DRAM capacity (Gbits)	0.5	1	4	8	32	64		

Jacob Abraham, November 10, 2020 38 / 48

Jacob Abraham, November 10, 2020 39 / 48



Technology Roadmaps

- Many steps needed to produce an IC
- Each step requires specialized (and expensive) equipment produced by different vendors
- Roadmaps give equipment manufacturers an idea what equipment would be used, and when the capability would be needed (example, scaling factor of sqrt(2))

ITRS established in 2013

ECE Department. University of Texas at Austir

ECE Department, University of Texas at Austin

- Scaling projections to 2028
- With Moore's law coming to an end, the final roadmap was issued in 2016
- Through IEEE's Rebooting Computing initiative, the IRDS was started

Some of the Focus Team Topics in the International Roadmap for Devices and Systems

Application Benchmarking Systems and Architectures More Moore Beyond CMOS Packaging Integration Outside System Connectivity Factory Integration Lithography Metrology Emerging Research Materials Environment, Safety, Health, and Sustainability Yield Enhancement Cryogenic Electronics and Quantum Information Processing (added in 2018)

am, November 10, 2020 41 / 48



Reachable Radius

ECE Department, University of Texas at Austin

• We can't send a signal across a large fast chip in one cycle anymore

Jacob Abraham, November 10, 2020 42 / 48

- But the microarchitect can plan around this
 - Just as off-chip memory latencies were tolerated









Solution

- Digital designers
 - \$70k salary
 - \$30k overhead
 - \$10k computer
 - \$10k CAD tools
 - Total:
 - $120k \times 7 = 840k$
- Analog designers
 - \$100k salary
 - \$30k overhead
 - \$10k computer
 - \$100k CAD tools
 - Total:

ECE Department, University of Texas at Austin

 $240k \times 3 = 720k$

- Support staff
 - \$45k salary
 - \$20k overhead
 - \$5k computer
 - Total:
 - $70k \times 5 = 350k$
- Fabrication
 - Back-end tools: \$1M
 - Masks: \$1M
 - Total: \$2M/year
- Summary:
 - 2 years at \$3.91M/year
 - \$8M design and prototype

Jacob Abraham, November 10, 2020 47 / 48

