

- Moore's Law - 2x transistors/2 years, ending
- Dennard Scaling - how scale this w/ transistor size
- Geopony Flat, power flat
- Pareto Optimal Frontier - best can do at any point
- performance vs cost vs speed tradeoff
- abstraction, modular design
- top down vs bottom up
- Hardware/Software - smart/human
- latency - time
- energy, power

NRE cost (neglecting) vs Recurring cost vs wafer cost

$$\text{Die (chip) cost} = \frac{\text{recurring cost}}{\text{die area} \cdot \text{die yield}}$$

$$\text{die yield} = \frac{\pi \cdot (\text{die area})^2}{\text{die area}} = \frac{\pi \cdot \text{die area}}{\sqrt{2 \cdot \text{die area}}}$$

$$\text{die yield} = \left(\frac{\text{die area}}{\sqrt{2 \cdot \text{die area}}} \right)^{-\alpha}$$

$\alpha \approx 3$

- abstract/resolution - "clean" signals
- combinational logic (CL) - just inputs
- sequential logic - fine of output given state for
- usually with tick
- state elements - storage
- RTL - abstract to CL + state elements
- Implementation
- Full custom
- Standard-cell (usually ASIC)
- Gate array
- FPGA
- Microprocessor
- Domain specific processor
- "birds home" - where do you finish with?
- Timing closure

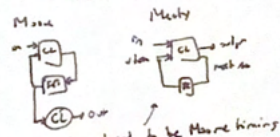
- Verilog
- on HDL
- structural - hierarchical, one-to-one
- behavioral - high level constructs (variables)
- module names (symbols); can be used
- end module
- always combinatorial (combinational); (combinational)
- reduction op - on all bits
- unsigned by default
- add signals if want
- continuous assignment - assign a = b;
- non-continuous assign
- always @(*) begin
- reg type (not recorded = 0)
- end
- always case, if else
- default assign!!!

- generator - like for loop, starts with
- generator is
- generate
- for(i=0; i<N; i++) begin: bit
- end
- assign
- use @(*) assign, always @(*) assign
- no register inference
- use @(*) assign!!!

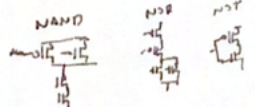
- FPGA
- MUX, latch, LUT
- has configurable logic blocks (CLB) "SRAM" based
- between program latches
- in LUT - 2ⁿ outputs
- Partitioning - draws circle how inputs, 1 output
- usually has dedicated blocks for bus, memory, expansion
- combine LUTs -> use MUX!

- Boolean Logic
- bool expr - bool logic - gate up
- dual - swap OR/AND
- SOP - sum of products, minterm
- POS - product of sums, maxterm
- K-map
- group odd on sides, 2 or ones
- circle first group, power of 2
- can overlap
- each is product term
- POS version - sum of products
- sum initial of product
- cost delay tradeoff
- Demorgan's Law!!!

- FSM
- diagram table
- Moore's output depends on state
- Mealy: output depends on input
- faster, but clock dependent
- extra output
- feedback
- issue: if input fluctuates, but

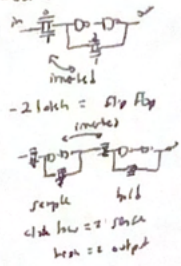


- State encoding
- least FFS -> 2 bits = 2 BFF
- easier to see one bit state encoding
- 1 bit per state - good for FPGA
- variables can be defined values
- but use is used
- CMOS
- MOSFET NMOS
- increase in Vgs ~
- 1st order for manufacturing
- accurate for imperfection
- make UV hardened
- HF drain p⁺ except mask
- make to implant n⁺
- driver possible gate
- self aligned:
- gate oxide strip
- etch tabs
- oxidation
- etch again
- FZMFGT
- Gate all around
- "through" or $\frac{W}{L}$
- some is one more rail
- metal is
- effective
- $\tau_{in} \propto V_r \cdot \tau_{eff}$
- use approx -> 1st order ODE



- Complex Gate
- Pullup/down (PUN, PDN)
- duals of each other!
- Pull -> PUN -> dual for PDN
- at node

- Transmission Gate
- $A \frac{a^1}{1} B$
- more real interactions
- Tri-state buffer - 0, 1, Z
- Latch



- Timing
- CL time
- setup
- hold
- setup
- hold
- usually setup CL with a
- FO4 delay
- reflections at input
- at this small
- Retiming
- avoid setup first!
- push thru gate
- push around latches
- constraints!
- in a loop - try not to
- combine
- push thru gate? and check

- Delay Model
- $C_0 = T C_0$ (usually)
- $t_p = 0.69 (3\tau) R_{eq} C_0$
- no dependencies!
- $f = \frac{C_L}{C_{in}}$
- inverts $t_p (1 + \frac{f}{2})$
- $C_{in} = 3W C_0$
- wire $0.38 C_{wire} L^2$
- $0.69 R_{eq} (C_{in} + C_{wire}) = 0.69 (R_{eq} C_{in} + R_{eq} C_{wire})$
- $\frac{R_{eq}}{C_{in}} \frac{C_{wire}}{L} = \frac{R_{eq}}{C_{in}} \frac{0.38 C_{wire} L^2}{L} = 0.38 R_{eq} C_{wire} L$

- Large load
- $C_{in} \frac{1}{2} L^2 \dots N \frac{1}{2} C_L$
- $f = \sqrt{C_L / C_{in}} = \sqrt{f}$
- $f = e^{\frac{(1+\tau)}{f}}$

- RISC-V
 - Pipelines
 - hazard - read instr dep on past
 - data hazard - read instr written
 - control hazard - branch/jump
 - forward
 - stall
 - Power
 - performance, area, unit cost, heat
 - servers - Total Cost of Ownership (TCO)
 - Switching Energy CMOS
 - 1/2 CV² loss, independent of tech
 - reduce N static transistors
 - reduce V_{dd} - limit clock
 - Rearrange circuit
 - reduce C per node - scale process
 - Dynamic Power
 - $P = \frac{1}{2} \alpha C V_{dd}^3 F$
 - where $\alpha = 0.4$
 - Short circuit current to enter both on \downarrow
 - leakage current - continue, but lower I_{on}
 - raise I_{on} → ↑ V_{dd} or ↓ V_t & doping
 - Parallel = P_{branch} + P_{FC} + P_{clock}
 - Technique for Low Power
 - Parallelism Pipelining - split block up
 - $P \propto F \cdot V_{dd}^2$
 - ↓ V_{dd} & F but some performance
 - F_{eff} roughly
 - more C_{int} also works
 - Scale V and F (DVFS)
 - better heat and/or energy
 - Power Down Mode
 - sleep transistors - better leak
 - slow clock rate
 - Slow down "static gates"
 - wait for clock to fall
 - use multiple V_{dd} for min critical path
 - or multi V_t (low V_t → slower)
 - Clock Gating
 - or clock enable only if state change
 - Thermal Management
 - keep cool
- Memory
 - volatile - reads power
 - static - refresh $\frac{DTR}{T_{refresh}}$
 - dynamic - charge $\frac{DTR}{T_{refresh}}$
 - generic architecture
 - word line - select row
 - bit line - data to edge
 - core access rates - keep cool
 - address bit dir. bit to row/col
 - row decoder - select row
 - col decoder - select col
 - SRAM
 - CMOS $\frac{1}{2} CV_{dd}^2$ + requires current sizing
 - read
 - precharge BL to V_{DD}
 - write H26H
 - cell pull 1 down
 - measure and differential
 - write
 - diff drive bit line
 - sense overpower pull-up
 - Decoder
 - predecoder - 2 bit = 4 inputs
 - few steps, work noisy
 - col decodes to log MUX
 - row transistors
 - DRAM
 - special process for low C
 - 3T - row drive both read, write pass good
 - 1T - row & write after read, need refresh
 - 1 bit read from array
 - int to EQ equal
 - then sense BL
 - then overwrite

- Multi-Port Memory
 - add port
 - add on the side of peripheral logic
- Combining Memory Blocks
 - tile (smaller blocks) leaf
 - leaf size limit by RC delay
 - low power by activate banks
 - delay + energy mostly from I/O interconnect
 - larger memory? (higher words)
 - IP just split the data
 - large memory (high speed)
 - use MUX w/ 1-bit bit
 - add read port?
 - write to both
 - add write port?
 - MUX w/ "predecoder"
 - write to both
 - write to both - parallel ports
 - FIFO - queue
 - 2 ports
 - full, empty, interrupt, multi
 - Cache
 - temporal, spatial locality
 - Tag Index, Dirty/Valid
 - dirty assoc. - no index
 - set assoc. - no dirty
 - direct
 - replacement policy
 - Parallelism
 - arithmetic by log a burst/word
 - tree
 - broken hardware out
 - Time Multiplex
 - ↓ cost ↑ time
 - Loop
 - better throughput, hardware
 - planning - limited by PE speed
 - Pipeline limit
 - clock skew
 - unequal stages
 - FF dominates cost
 - clock distribution
 - hardware - parallel "loop carry dependency"
 - reorder ops - check hazard
 - C-Flow
 - run n-indep stream in parallel
 - register each reg w/ n reg
 - loop F per stream, bit n of them
 - read throughput
 - MUX streams
 - multi threaded parallel
 - SIMD
 - RT Language
 - 1 reg cycle / 1 sum cycle
 - FSM follow directly
 - Design Pattern
 - Dataflow graph
 - "Algorithm"
 - Mem. calculation
 - Cost, Perf, Power, Area
 - Optimization, Variation
 - Detailed Design
 - Resource Utilization Chart
 - show power impact
 - Modulo Scheduling
 - pick characteristic subset
 - schedule over time wrap, +1
 - if not fit, increase size

- Adder
 - Carry Ripple - 10cm, cost 10n
 - subtraction just XOR
 - 2's complement, Cin = 1
 - Carry Select - do both, select 1
 - multiple ripple stages
 - fix size, fix bits, fixed
 - T_{prop}, C_{in} = 2ⁿ ripple + time
 - N = 3 better
 - if unequal, can fix to better match delay
 - Carry Look Ahead (CLA) O(log n)
 - propagate - P_i = a_i ⊕ b_i; dx out
 - generate - G_i = a_i b_i
 - group output - Cout = G + P Cin
 - S_i = P_i ⊙ C_i
 - P = P₁ P₂
 - G = G₂ + G₁ P₂
 - C_{i+1} = G_i + P_i C_i
 - arithmetic!
 - parallel prefix adder - carry group
 - Karnaugh
 - Ladner-Fischer
 - Brent-Kung
 - Han-Carlson
 - mix w/ ripple, use CLA to make carry
 - bit serial adder
 - O(n) hardware, O(n) cycle
 - Multiplier
 - shift and add
 - add partial products
 - sign ext
 - adders, 2's complement
 - Array multiplier
 - layout, n² cost
 - Carry save
 - 3 bit in, 2 out
 - need CLA at each full adder
 - final Wallace Tree multiplier
 - Booth recoding
 - sign of 2 bit, half delay 1
 - Bit cost
 - O(n²) time, cost 1
 - Booth multiplier - signed multi
 - add/subtract sign bit, no sign extension
 - all inputs of cell
 - scan through cell a universal
 - Constant coefficient (KCM)
 - shift and add table
 - use shift register
 - Canonical Signed Digit (CSD)
 - 15 → 102 T 2ⁿ⁻¹ binary
 - 11 → 01 2ⁿ⁻¹ = 2ⁿ⁻¹
 - do factorization to simplify further
 - exhaustive search
 - Shift Register
 - O(log n) - A power of 2 shift
 - barrel shifter - n² cost
 - can be useful
 - any permutation, sorted, etc
 - 10000000
 - Clock non-idealities
 - skew - time diff between edges +/-
 - jitter - random variation
 - setup time - t_{su}
 - hold time - t_{hd}
 - setup - margin
 - max skew - offset min cycle time
 - you skew - hold time
 - Clock Distribution
 - H - has eq in both & buffer
 - Chip Packaging
 - bond wires on pad - 1mm
 - flip-chip - double pad
 - micro ball on pad - fine pad
 - ESD - diodes + guard ring

- Power Distribution
 - IR drop, L, and resistance
 - max power, IR drop
 - metal resistance
 - current sum multi
 - $\frac{I^2}{A}$ drop, IR drop
 - write 5-10% P_{su} variation
 - flip chip better
 - decoupling cap
- Faults
 - Design
 - Manufacturing
 - Runtime
- Modeling
 - ATPG - test vector
 - BIST - self test
 - scan chain
 - extra hardware sacrificed
 - laser fix
- Routing
 - placement, transit, interconnect
- Physical Full
 - NBT, PBT
 - HCT
 - TADB - show
 - electro migration
- ECC
 - parity
 - Hamming code
 - SEC - one error correct
 - P_{su} → DED - data error detection
 - undetectable - bits of parity are error
 - probability even parity
 - error falls when it is!
 - C=0, P=0 no error
 - C=0, P=1 SEC
 - C=0, P=0 DED
 - C=0, P=1 error in P

