

# ESE532: System-on-a-Chip Architecture

Day 8: September 26, 2018  
Spatial Computations



## Today

- Graph Cycles (from Day 7)
- Accelerator Pipelines
- FPGAs
- Zynq Computational Capacity

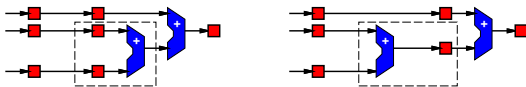
## Message

- Custom accelerators efficient for large computations
  - Exploit Instruction-level parallelism
  - Run many low-level operations in parallel
- Field Programmable Gate Arrays (FPGAs)
  - Allow post-fabrication configuration of custom accelerator pipelines
  - Can offer high computational capacity

## Graph Cycles

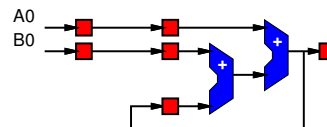
## Reminder

- Able to pipeline and retiming to reduce cycle time on acyclic dataflow graphs



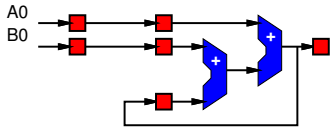
## Preclass 1

- Can we retiming to reduce cycle



## Retiming Limits?

- What prevents us from retiming?



Penn ESE532 Fall 2018 -- DeHon

7

## Cycle Observation

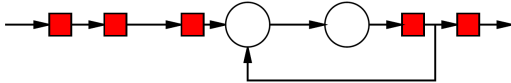
- Retiming does not allow us to change the *number of registers inside a cycle*.
- Limit to cycle time
  - Max delay in cycle / Registers in cycle
- Pipelining doesn't help inside cycle
  - Cannot push registers into cycle

Penn ESE532 Fall 2018 -- DeHon

8

## Simple Cycle

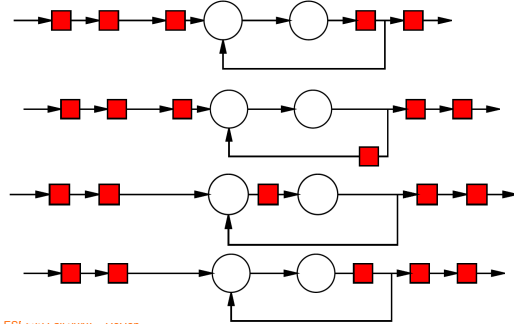
- Delay of cycle?
- Registers in cycle?
- What happens to cycle if try to apply lead/lag?



Penn ESE532 Fall 2018 -- DeHon

9

## Retiming



Penn ESE532 Fall 2018 -- DeHon

## Loop

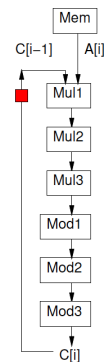
- What does graph look like for this loop body?
  - [multiply and mod each take 3 cycles]
- For (i=0;i<N;i++)  
 $C[i] = (C[i-1] * A[i]) \% N;$

Penn ESE532 Fall 2018 -- DeHon

11

## Loop

- For (i=0;i<N;i++)  
 $C[i] = (C[i-1] * A[i]) \% N;$

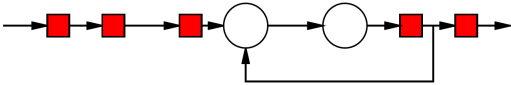


Penn ESE532 Fall 2018 -- DeHon

12

## Initiation Interval (II)

- Cyclic dependencies can limit throughput
- Due to dependent cycles,
  - May not be able to initiate a new computation on every cycle
- II – cycles (delay) before can initiate
- Throughput =  $1/II$

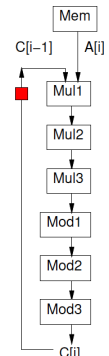


Penn ESE532 Fall 2018 -- DeHon

13

## Loop

- For ( $i=0; i<N; i++$ )  
 $C[i] = (C[i-1] * A[i]) \% N;$
- Initiation Interval?



Penn ESE532 Fall 2018 -- DeHon

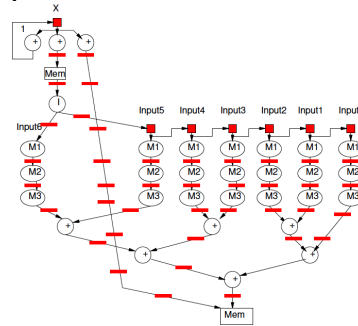
14

## Accelerator Datapaths

Penn ESE532 Fall 2018 -- DeHon

15

## Pipeline for Unrolled Loop

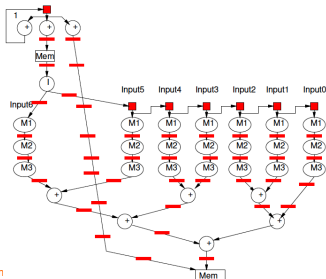


Penn ESE532 Fall 2018 -- DeHon

16

## Preclass 1

- For fully unrolled loop shown, how many instructions per pipeline cycle?
  - Add
  - Mpy
  - Load
  - Store



Penn ESE532 Fall 2018 -- DeHon

17

## Spatial Pipeline

- Can compute equivalent of tens of "instructions" in a cycle
- Wire up primitive operators
  - No indirection through register file, memory
- Pipeline for operator latencies
- Any dataflow graph of computational operations

Penn ESE532 Fall 2018 -- DeHon

18

## Operators

- Can assemble any custom operators
  - Ones may not have in generic processor
- Processor
  - Add, bitwise-xor/and/or, multiply
  - Maybe: floating-point add, multiply
- Less likely
  - Square-root, exponent, cosine, encryption (AES) step, polynomial evaluate, log-number-system

Penn ESE532 Fall 2018 -- DeHon

19

## Accelerators

- Compression/decompression
- Encryption/decryption
- Encoding (ECC, Checksum)
- Discrete Cosine Transform (DCT)
- Sorter
- Taylor Series Approximation of function
- Transistor evaluator
- Tensor or Neural Network evaluator

Penn ESE532 Fall 2018 -- DeHon

20

## Streaming Dataflow

- Replace operator with custom accelerator
- Stream data to/from it

Penn ESE532 Fall 2018 -- DeHon

21

## Streaming Dataflow Example



Penn ESE532 Fall 2018 -- DeHon

22

## Application-Specific SoCs

- For dedicated applications may build custom hardware for accelerators
  - Layout VLSI, fab unique chips
  - ESE370, 570
- Video-encoder – include custom DCT, motion-estimation engines

Penn ESE532 Fall 2018 -- DeHon

23

## Customizable Accelerators

- With post-fabrication configurability can exploit without unique fabrication
- Need programmable substrate that allows us to wire-up computations

Penn ESE532 Fall 2018 -- DeHon

24

## Field-Programmable Gate Arrays

FPGAs

Penn ESE532 Fall 2018 -- DeHon

25

## FPGA

- Idea: Can wire up programmable gates in the “field”
  - After fabrication
  - At your desk
  - When part “boots”
- Like a “Gate Array”
  - But not hardwired

Penn ESE532 Fall 2018 -- DeHon

26

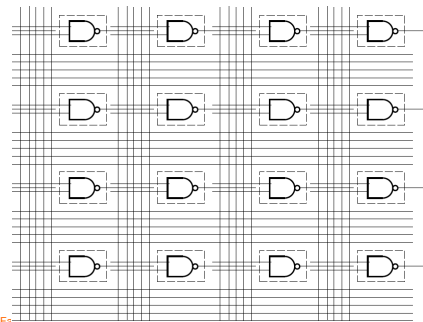
## Gate Array

- Idea: Provide a collection of uncommitted gates
- Create your “custom” logic by wiring together the gates
- Less layout, fewer masks than full custom
  - Since only wiring together pre-fab gates
  - lower cost (fewer masks)
  - lower manufacturing delay

Penn ESE532 Fall 2018 -- DeHon

27

## Gate Array



Penn ESE532 Fall 2018 -- DeHon

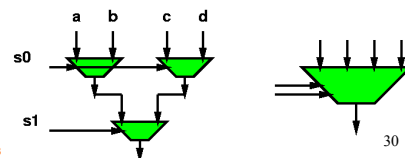
## GA → FPGA

- Remove the need to even fabricate the wiring mask
- Make “customization” soft
- Key trick:
  - Use reprogrammable configuration bits
  - Typically: static-RAM bits
    - Like SRAM cells or latches in memory
    - Hold a configuration value

Penn ESE532 Fall 2018 -- DeHon

29

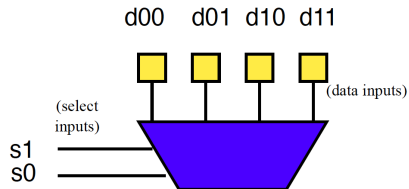
## Mux with configuration bits = programmable gate



Penn ESE532 Fall 2018

## Preclass 2a

- How do we program to behave as and2?

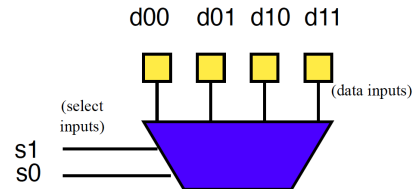


Penn ESE532 Fall 2018 -- DeHon

31

## Preclass 2b

- How do we program to behave as xor2?

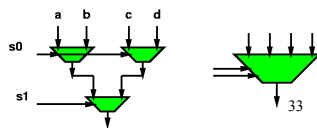


Penn ESE532 Fall 2018 -- DeHon

32

## Mux as Logic

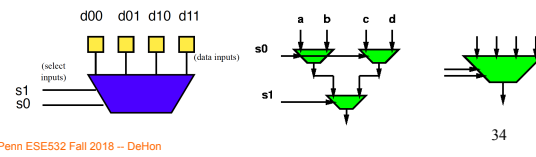
- Just by “configuring” data into this mux4,
  - Can select **any** two input function



Penn ESE532 Fall 2018 -- DeHon

## LUT – LookUp Table

- When use a mux as programmable gate
  - Call it a **LookUp Table (LUT)**
  - Implementing the Truth Table for small # of inputs
    - # of inputs = k (need mux-2<sup>k</sup>)
  - Just lookup the output result in the table

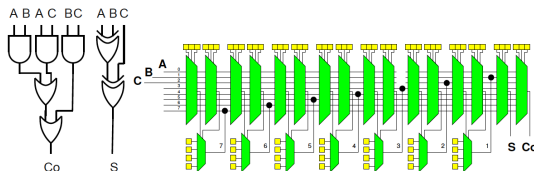


Penn ESE532 Fall 2018 -- DeHon

34

## Preclass 3

- How do we program full adder?



Penn ESE532 Fall 2018 -- DeHon

35

## FPGA

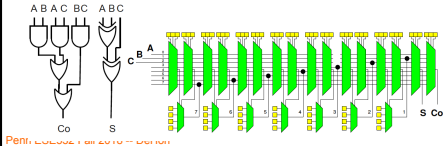
- Programmable gates + wiring
  - (both built from muxes w/ config. bits)
- Can wire up any collection of gates
  - Like a gate array

Penn ESE532 Fall 2018 -- DeHon

36

## Simplistic FPGA (illustrate possibility)

- Every LUT input has a mux
- Every such mux has  $m=(N+1)$  inputs
  - An input for each LUT output (N 2-LUTs)
  - An input for each Circuit Input (I Circuit inputs)
- Each Circuit Output has an m-input mux

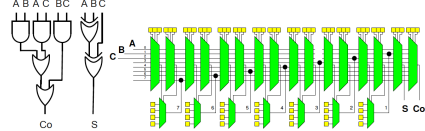


Penn ESE532 Fall 2018 -- DeHon

37

## Simplistic FPGA (illustrate possibility)

- N 2-LUTs, I Circuit Inputs, O Circuit Outputs
- $2N+O$  muxes to connect
- Can build **any** combinational logic circuit that doesn't need more than N 2-input gates, I inputs, O outputs



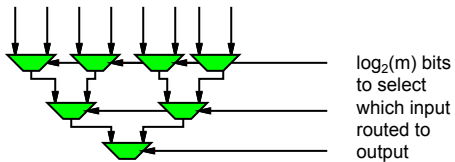
Penn ESE532 Fall 2018 -- DeHon

38

## How big is an m-input mux?

- In terms of 2-input muxes?

m inputs – what we are selecting from



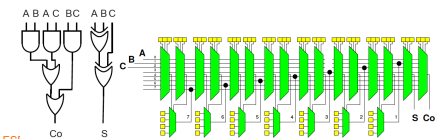
$\log_2(m)$  bits to select which input routed to output

Penn ESE532 Fall 2018 -- DeHon

39

## Simplistic FPGA (illustrate possibility...and expense)

- $2N+O$  m-input muxes;  $m=N+I$
- Each m-input mux is  $m-1$  2-input muxes
- Requires:  $(2N+O)*(N+I-1)$  2-input muxes
- Mux area grows as  $\sim N^2$ 
  - when gate (LUT) area grows as N



Penn ESI

40

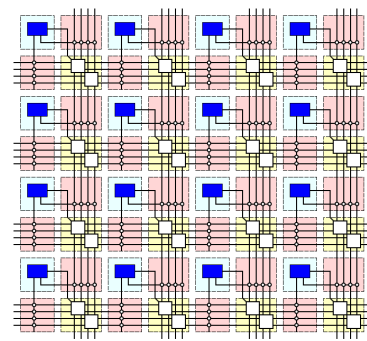
## Interconnect

- Fully connected mux input too expensive
  - ...and not necessary
- Want
  - To be able to wire up gates
  - Economical with wires and muxes
    - ...and configuration bits
  - Exploit locality (keep wires short)

Penn ESE532 Fall 2018 -- DeHon

41

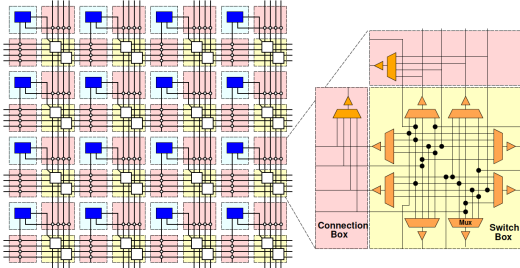
## Simple FPGA



Penn ESE532 Fall 2018 -- DeHon

42

## Simple FPGA

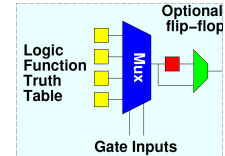


Penn ESE532 Fall 2018 -- DeHon

43

## Flip-Flops

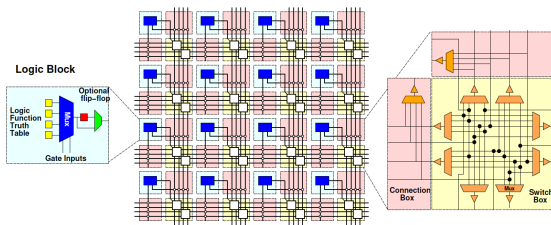
- Want to be able to pipeline logic
- ...and generally hold state
  - E.g. implement hold Input-N in preclass 2
- Add optional flip-flop on each gate



Penn ESE532 Fall 2018 -- DeHon

44

## Simple FPGA



Penn ESE532 Fall 2018 -- DeHon

45

## FPGA Design

- Raises many architectural design questions
  - How big (many inputs) should the gates have?
    - Are LUTs really the right thing...
  - How rich is the interconnect?
    - Wires/channel
    - Wire length
    - Switching options

Penn ESE532 Fall 2018 -- DeHon

46

## Modern FPGAs

- Logic Blocks
  - hardwired fast-carry logic
    - Can implement adder bit in single "LUT"
  - Speed optimized: 6-LUTs
  - Energy, Cost optimization: 4-LUTs
  - Clusters many LUTs into a tile
- Interconnect
  - Mesh, segments of length 4 and longer

Penn ESE532 Fall 2018 -- DeHon

47

## More than LUTs

- Should there be more than LUTs in the "array" fabric?
  - What else might we want?

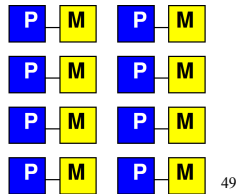
Penn ESE532 Fall 2018 -- DeHon

48



## Embedded Memory

- One flip-flop per LUT doesn't store state densely
- Want memory close to logic



Penn ESE532 Fall 2018 -- DeHon

49

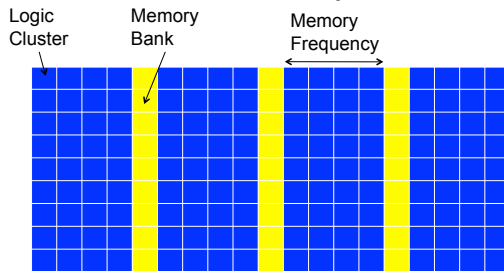
## Embed Memory in Array

- Replace logic clusters
- Convenient to replace columns
  - Since area of memory may not match area of logic cluster

Penn ESE532 Fall 2018 -- DeHon

50

## Embedded Memory in FPGA



Memory banks on Xilinx called BRAMs (Block RAMs)

Penn ESE532 Fall 2018 -- DeHon

51

## Hardwired Multipliers

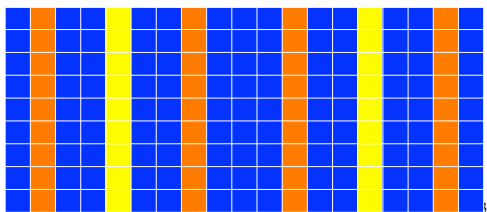
- Can build multipliers out of LUTs
  - Just as can implement multiplies on processor out of adds
- But, custom multiplier is smaller than LUT-configured multiplier
  - ...and multipliers common in signal processing, scientific/engineering compute

Penn ESE532 Fall 2018 -- DeHon

52

## Multiplier Integration

- Integrate like memories
  - Replace columns



Penn ESE532 Fall 2018 -- DeHon

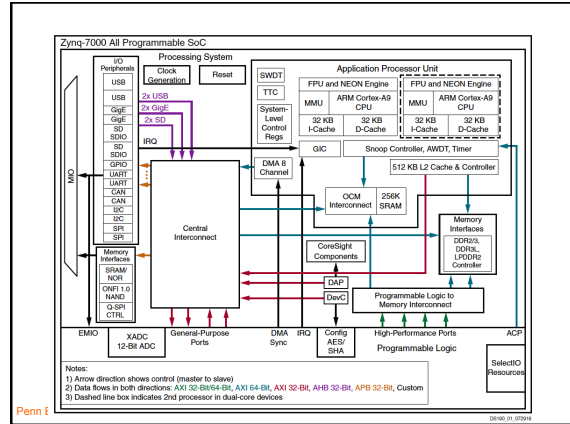
## More FPGA Architecture Design Questions

- Size of Memories? Multipliers?
- Mix of LUTs, Memories, Multipliers?
- Add processors? Floating-point?
- Other hardwired blocks?
- How manage configuration?

Penn ESE532 Fall 2018 -- DeHon

54

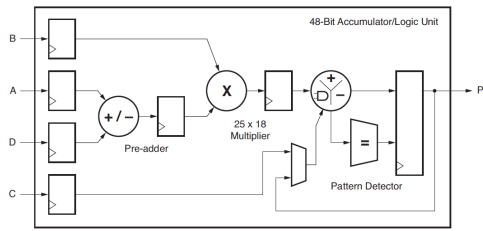
# Zynq



# XC7Z020

- 6-LUTs: 53,200
- DSP Blocks: 220
  - 18x25 multiply, 48b accumulate
- Block RAMs (BRAMs): 140
  - 36Kb
  - Dual port
  - Up to 72b wide (512x72)

# DSP48



# Preclass 4

Resoruces	Cycle	per second
50,000 adder bits	0.2 GHz	
128x2 adder bits	1.0 GHz	
200 multiply-accumulates	0.2 GHz	
8x2 multiply-accumulates	1.0 GHz	

# Compute Capacity

- How compare between ARM/NEON and FPGA array?
  - Adder-bits/second?
  - Multiply-accumulators/second?

## Capacity → Density

- Says Zynq has high computational capacity in FPGA
- More broadly
  - FPGA can have more compute/area than processor
    - E.g., more adder bits in some fixed area
  - SIMD can have more compute/area than processor
    - How wide SIMD can you exploit?

Penn ESE532 Fall 2018 -- DeHon

61

## FPGA Potential

- FPGA Array has high raw capacity
- Exploitable when computation has high regularity
  - Uses the same computation over-and-over
  - High throughput on a computation
  - Build customized accelerator pipeline to match the computation
- Low-hanging fruit
  - Operator/function takes most of the compute time

Penn ESE532 Fall 2018 -- DeHon

62

## 90/10 Rule (of Thumb)

- Observation that code is not used uniformly
- 90% of the time is spent in 10% of the code
- Knuth: 50% of the time in 2% of the code
- Opportunity
  - Build custom datapath in FPGA (hardware) for that 10% (or 2%) of the code

Penn ESE532 Fall 2018 -- DeHon

63

## Big Ideas

- Custom accelerators efficient for large computations
  - Exploit Instruction-level parallelism
  - Run many low-level operations in parallel
- Field Programmable Gate Arrays (FPGAs)
  - Allow post-fabrication configuration of custom accelerator pipelines
  - Can offer high computational capacity

Penn ESE532 Fall 2018 -- DeHon

64

## Admin

- Reading for Day 9 on canvas
- HW4 due Friday
- No homework due 10/5 (Fall Break)
- HW5 out
  - Due 10/12
  - SDSoC synthesis slow (plan for it)
  - May need to run Vivado HLS on linux

Penn ESE532 Fall 2018 -- DeHon

65