

ESE5320: System-on-a-Chip Architecture

Day 8: September 25, 2024
Spatial Computations

Preclass 3 pages



Penn ESE5320 Fall 2024 -- DeHon

1

Today

- Accelerator Pipelines (Part 1)
- FPGAs (Part 2)
- Computational Capacity (Part 3)
 - Zynq

Penn ESE5320 Fall 2024 -- DeHon

2

2

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

Penn ESE5320 Fall 2024 -- DeHon

3

3

Accelerator Datapaths

Penn ESE5320 Fall 2024 -- DeHon

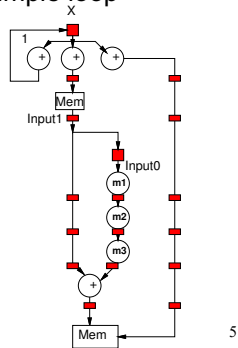
4

4

Pipeline Graph

- Last time: pipelined simple loop

```
unsigned char Input0, Input1;
Input1=Input[Y * INPUT_WIDTH + X + 0];
for (int X = 0; X < OUTPUT_WIDTH; X++)
{
  unsigned int Sum;
  Input0=Input1;
  Input1=Input[Y * INPUT_WIDTH + X + 1];
  Sum = Coefficients_0 * Input0 + Input1;
  Output[Y * OUTPUT_WIDTH + X] = Sum;
}
```

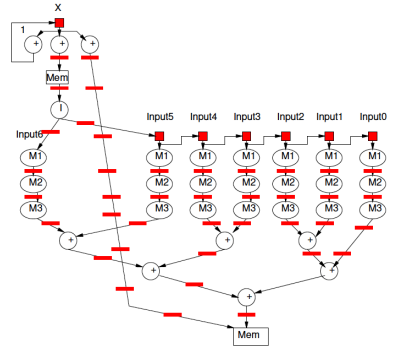


Penn ESE5320 Fall 2024 -- DeHon

5

5

Pipeline for Unrolled Loop



Penn ESE5320 Fall 2024 -- DeHon

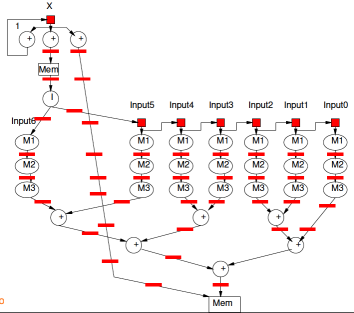
6

6

Preclass 1

- For fully unrolled loop shown, how many instructions per pipeline cycle?

- Add
- Mpy
- Load
- Store



Penn ESE5320 Fall 2024 -- DeHon

7

7

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 ESE5320 Fall 2024 -- DeHon

8

8

Operators

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

Penn ESE5320 Fall 2024 -- DeHon

9

9

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 ESE5320 Fall 2024 -- DeHon

10

10

Streaming Dataflow

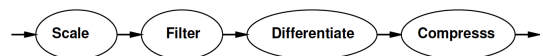
- Replace operator with custom accelerator
- Stream data to/from it

Penn ESE5320 Fall 2024 -- DeHon

11

11

Streaming Dataflow Example



Penn ESE5320 Fall 2024 -- DeHon

12

12

Application-Specific SoCs

- For dedicated applications may build custom hardware for accelerators
 - Layout VLSI, fab unique chips
 - ESE3700, 5700
- Video-encoder – include custom DCT, motion-estimation engines

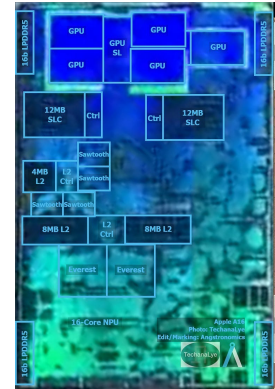
Penn ESE5320 Fall 2024 -- DeHon

13

13

Apple A16 Bionic

- ? 110+mm², 4nm
- 16 Billion Tr.
- iPhone 14
- 6 ARM cores
 - 2 fast (3.5GHz)
 - 4 low energy (2GHz)
- 5 custom GPUs (1.4GHz)
- 16 Neural Engines
 - 17 Trillion ops/s?



Penn ESE5320 Fall 2024 -- DeHon

https://en.wikipedia.org/wiki/Apple_A16

14

Customizable Accelerators

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

Penn ESE5320 Fall 2024 -- DeHon

15

15

Field-Programmable Gate Arrays

FPGAs
Part 2

Penn ESE5320 Fall 2024 -- DeHon

16

16

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 ESE5320 Fall 2024 -- DeHon

17

17

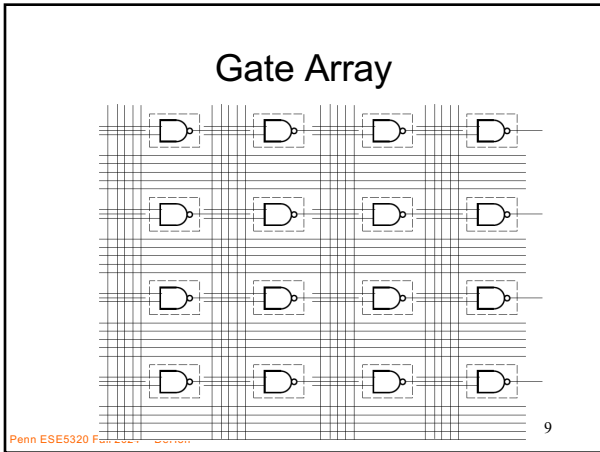
Gate Array

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

Penn ESE5320 Fall 2024 -- DeHon

18

18



19

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
 - Configuration bits define gates, wiring

20

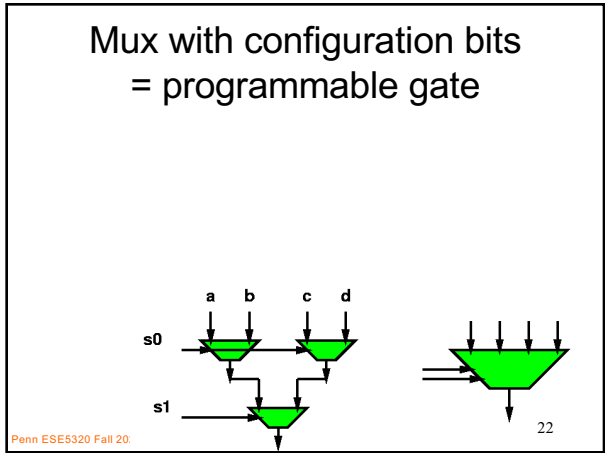
Multiplexer Gate

- MUX
 - When S=0, output=i0
 - When S=1, output=i1

$$\text{Out} = \bar{s} \cdot i_0 + s \cdot i_1$$

21

21



22

Preclass 4a

- How do we program to behave as and2?

23

23

Preclass 4b

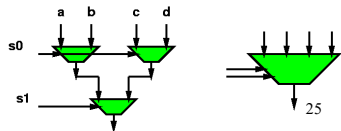
- How do we program to behave as xor2?

24

24

Mux as Logic

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

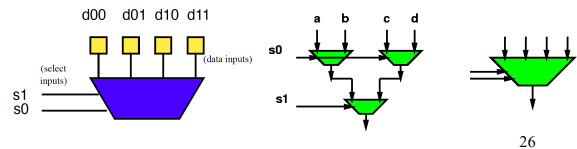


Penn ESE5320 Fall 2024 -- DeHon

25

LUT – LookUp Table

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

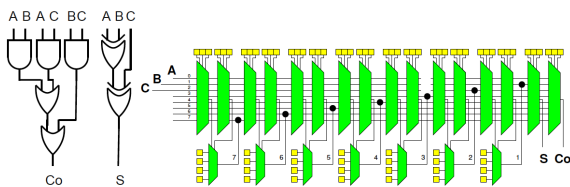


Penn ESE5320 Fall 2024 -- DeHon

26

Preclass 6

- How do we program full adder?



Penn ESE5320 Fall 2024 -- DeHon

[Switch Google Doc]

27

27

FPGA

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

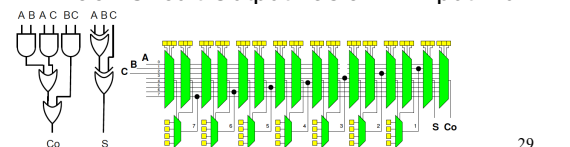
Penn ESE5320 Fall 2024 -- DeHon

28

28

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



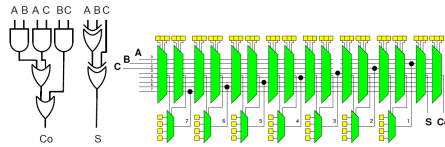
Penn ESE5320 Fall 2024 -- DeHon

29

29

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 ESE5320 Fall 2024 -- DeHon

30

30

Preclass 3

How big is an m-input mux?

- In terms of 2-input muxes?
 - Warmup: how many for 4-input (Preclass 2)
 - Warmup: how many for 8-input (below)

m inputs – what we are selecting from

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

Penn ESE5320 Fall 2024 -- DeHon 31

31

Math: Series Sums

- $A_0(1+r+r^2+r^3+r^4+\dots)$
- $A_0(1+r+r^2+r^3+r^4+\dots)*(1-r)$

$$= A_0 + A_0 r + A_0 r^2 + A_0 r^3 + A_0 r^4 + \dots$$

$$- A_0 r - A_0 r^2 - A_0 r^3 - A_0 r^4 - \dots$$

$$= A_0 \quad (\text{when } r < 1)$$
- $A_0(1+r+r^2+r^3+r^4+\dots)*(1-r) = A_0$
- $A_0(1+r+r^2+r^3+r^4+\dots) = A_0/(1-r)$

Penn ESE5320 Fall 2024 -- DeHon 32

32

Receding Sum

Penn ESE5320 Fall 2024 -- DeHon 33

33

Preclass 3

How big is an m-input mux?

- m-input mux will require m-1 2-input muxes

m inputs – what we are selecting from

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

Penn ESE5320 Fall 2024 -- DeHon 34

34

Simplistic FPGA

(illustrate possibility...and expense)

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

Penn ESI 35

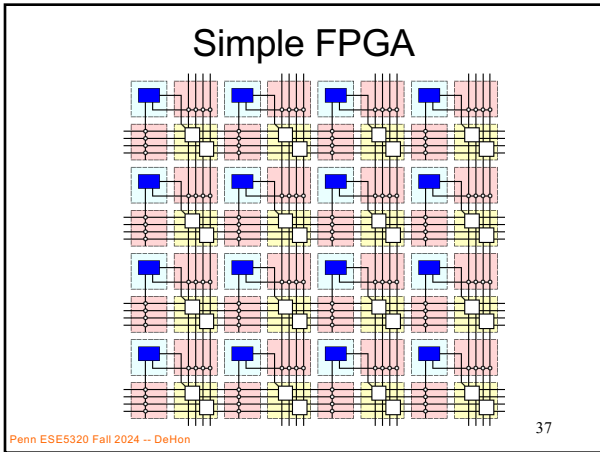
35

Interconnect

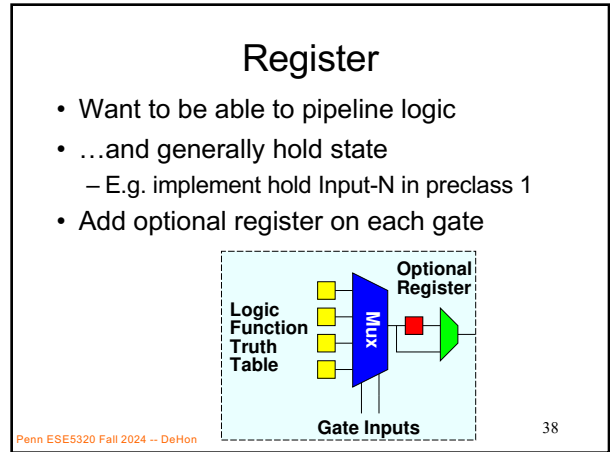
- Fully connected mux input is too expensive, growing as N^2
 - ...and not necessary
- Want
 - To be able to wire up gates
 - Economical with wires and muxes
 - ...and configuration bits
 - Exploit locality (keep wires short)

Penn ESE5320 Fall 2024 -- DeHon 36

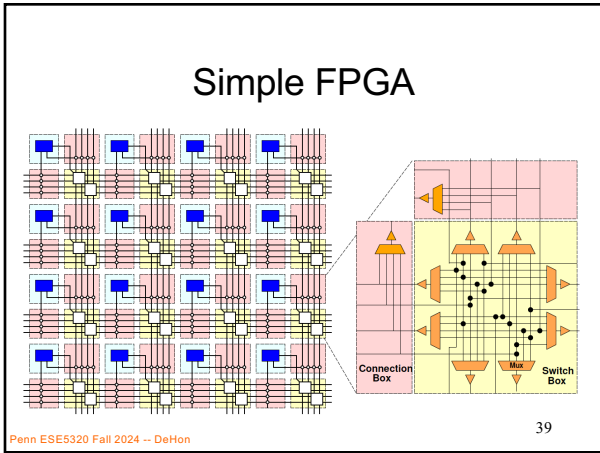
36



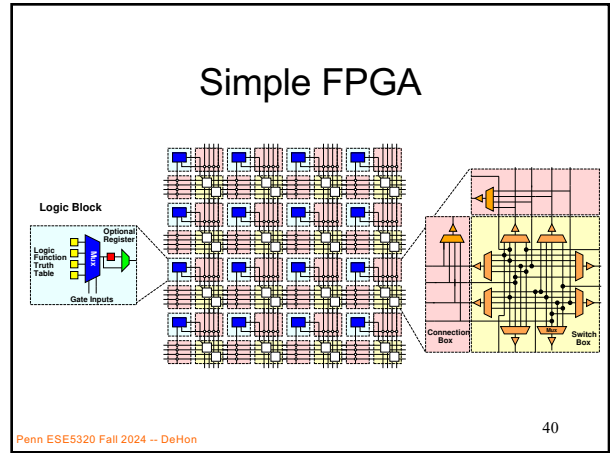
37



38



39



40

FPGA Design

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

Penn ESE5320 Fall 2024 -- DeHon

41

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 ESE5320 Fall 2024 -- DeHon

42

More than LUTs

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

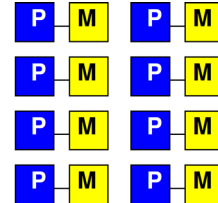
Penn ESE5320 Fall 2024 -- DeHon

43

43

Embedded Memory

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



Penn ESE5320 Fall 2024 -- DeHon

44

44

Embed Memory in Array

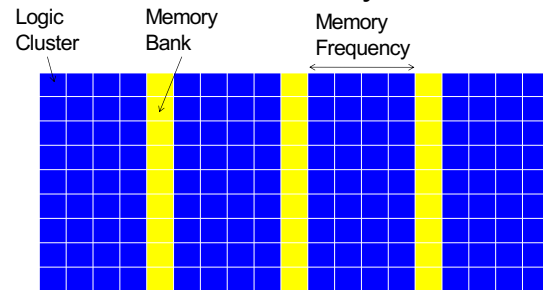
- Replace logic clusters
- Convenient to replace columns
 - Since area of memory may not match area of logic cluster

Penn ESE5320 Fall 2024 -- DeHon

45

45

Embedded Memory in FPGA



Memory banks on AMD/Xilinx called BRAMs (Block RAMs)

Penn ESE5320 Fall 2024 -- DeHon

46

46

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

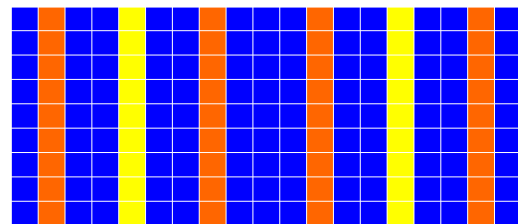
Penn ESE5320 Fall 2024 -- DeHon

47

47

Multiplier Integration

- Integrate like memories
 - Replace columns



Penn ESE5320 Fall 2024 -- DeHon

48

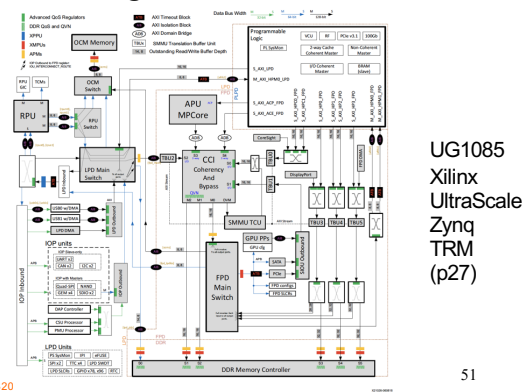
More FPGA Architecture Design Questions

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

Zynq MPSoC

Part 3

Programmable SoC



UG1085
Xilinx
UltraScale
Zynq
TRM
(p27)

ZU3EG (Ultra96)

- 6-LUTs: 70,560
- DSP Blocks: 360
 - 18x27 multiply, 48b accumulate
- Block RAMs (BRAMs): 216
 - 36Kb
 - Dual port
 - Up to 72b wide (512x72)

DSP48

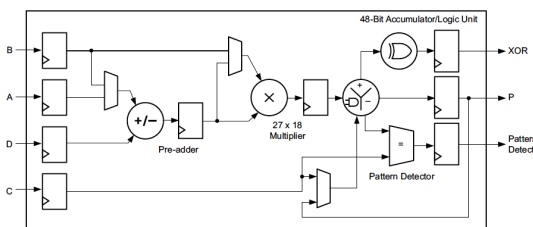


Figure 1-1: Basic DSP48E2 Functionality

Xilinx UG579 UltraScale DSP Slice User's Guide 53

Preclass 5

Approximating	Resources	Cycle	Per second
Zynq LUTs	70,000 adder bits	0.5 GHz	
4x ARM Scalar	4x2x64 adder bits	1.2 GHz	
4x ARM Neon	4x1x64 adder bits	1.2 GHz	
Zynq DSP	360 multiply-accumulates	0.5 GHz	
4x ARM Scalar	4x(1 mpy+1add)	1.2 GHz	
4x ARM Neon	4x1x4 multiply-accumulates	1.2 GHz	

- How compare between ARM scalar, 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 (Day 6)
 - How wide SIMD can you exploit?

Penn ESE5320 Fall 2024 -- DeHon

55

55

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 ESE5320 Fall 2024 -- DeHon

56

56

90/10 Rule

- 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 ESE5320 Fall 2024 -- DeHon

57

57

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 ESE5320 Fall 2024 -- DeHon

58

58

Admin

- Reading for Day 9 on canvas
- HW4 due on Friday
- HW5 out soon
 - Heavier – start early...
 - Vivado HLS synthesis slow (plan for it)

Penn ESE5320 Fall 2024 -- DeHon

59

59