

ESE532: System-on-a-Chip Architecture

Day 9: February 13, 2017
Spatial Computations



Today

- 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

Preclass 1

- Spatial pipeline?

```
for (i=0;i<MAX;i++) {  
    t=0;  
    for (j=0;j<5;j++)  
        t+=x[i+j]*w[j];  
    y[i]=t;  
}
```

Preclass 1

- Adders/cycle?
- Multipliers/cycle?
- Pipeline Depth?

Spatial Pipeline

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

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 Spring 2017 -- DeHon

7

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 Spring 2017 -- DeHon

8

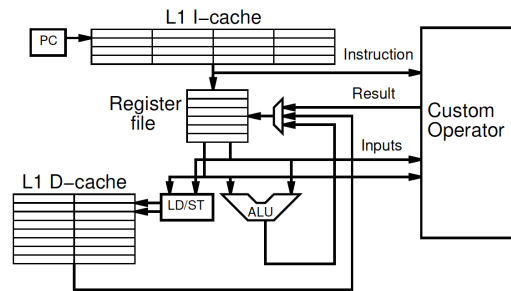
Integration System Architectures

- Instruction Augmentation
- Streaming Dataflow Operators

Penn ESE532 Spring 2017 -- DeHon

9

Instruction Augmentation

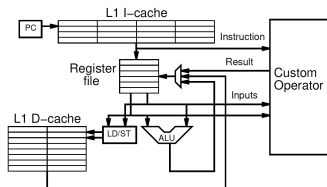


Penn ESE532 Spring 2017 -- DeHon

10

Instruction Augmentation

- Looks like additional instructions
- Issue from processor
- Examples
 - Floating-point
 - SIMD



Penn ESE532 Spring 2017 -- DeHon

11

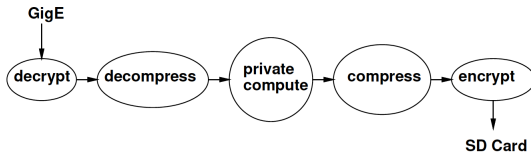
Streaming Dataflow

- Replace operator with custom accelerator
- Stream data to/from it
 - Maybe dedicated path
 - Operates from memory
 - Maybe DMA

Penn ESE532 Spring 2017 -- DeHon

12

Streaming Dataflow Example



Penn ESE532 Spring 2017 -- DeHon

13

Application-Specific SoCs

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

Penn ESE532 Spring 2017 -- DeHon

14

Customizable Accelerators

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

Penn ESE532 Spring 2017 -- DeHon

15

Field-Programmable Gate Arrays

FPGAs

Penn ESE532 Spring 2017 -- DeHon

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 ESE532 Spring 2017 -- DeHon

17

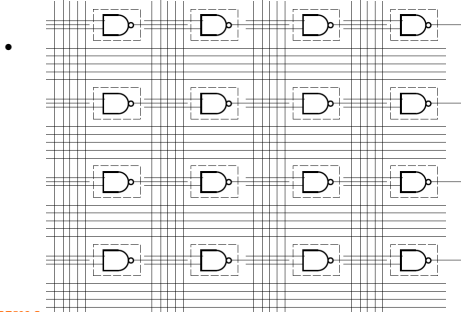
Gate Array

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

Penn ESE532 Spring 2017 -- DeHon

18

Gate Array



Penn ESE532 Spring 2017 -- DeHon

9

GA → FPGA

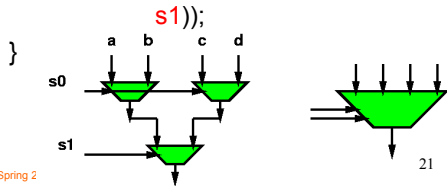
- Remove the need to even fabricate the wiring mask
- Make “customization” soft
- Key trick:
 - Use reprogrammable configuration bits
 - Typically: static-RAM bits
 - Like SRAM cells or latches
 - Hold a configuration value

Penn ESE532 Spring 2017 -- DeHon

20

Mux with configuration bits = programmable gate

- `bool mux4(bool a, b, c, d, s0, s1) {
return(mux2(mux2(a,b,s0),
mux2(c,d,s0),
s1));
}`

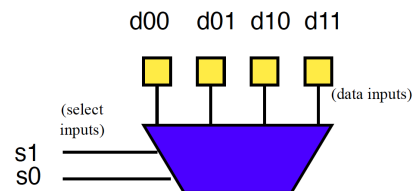


Penn ESE532 Spring 2

21

Preclass 2a

- How do we program to behave as and2?

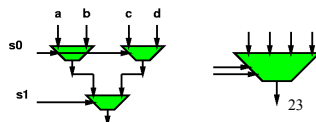


Penn ESE532 Spring 2017 -- DeHon

22

Mux as Logic

- `bool and2(bool x, y)
{return (mux4(false,false,false,true,x,y));}`

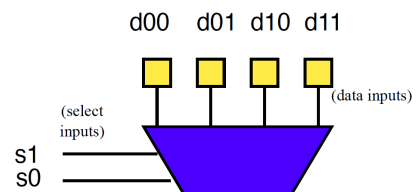


Penn ESE532 Spring 2017 -- DeHon

23

Preclass 2b

- How do we program to behave as xor2?



Penn ESE532 Spring 2017 -- DeHon

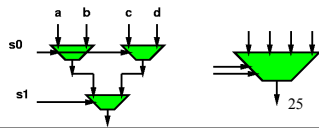
24

Mux as Logic

- bool and2(bool x, y)


```
{return (mux4(false,false,false,true,x,y));}
```
- bool xor2(bool x, y)

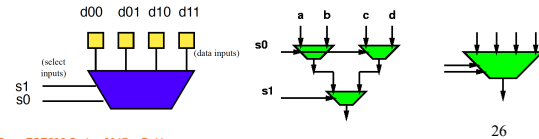

```
{return (mux4(false,true,true,false,x,y));}
```
- Just by “configuring” data into this mux4,
 - Can select **any** two input function



Penn ESE532 Spring 2017 -- DeHon

LUT – LookUp Table

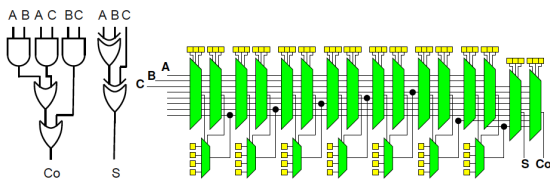
- When use a mux as programmable gate
 - Call it a **LookUp Table (LUT)**
 - Implementing the Truth Table for small # of inputs
 - # of inputs =k (need mux-2^k)
 - Just lookup the output result in the table



Penn ESE532 Spring 2017 -- DeHon

Preclass 3

- How do we program full adder?



27

Penn ESE532 Spring 2017 -- DeHon

FPGA

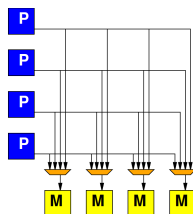
- Programmable gates + wiring
 - (both built from muxes w/ config. bits)
- Can wire up any collection of gates
 - Like a gate array

28

Penn ESE532 Spring 2017 -- DeHon

Crossbar Interconnect

- How did we say the crossbar interconnect scaled?



29

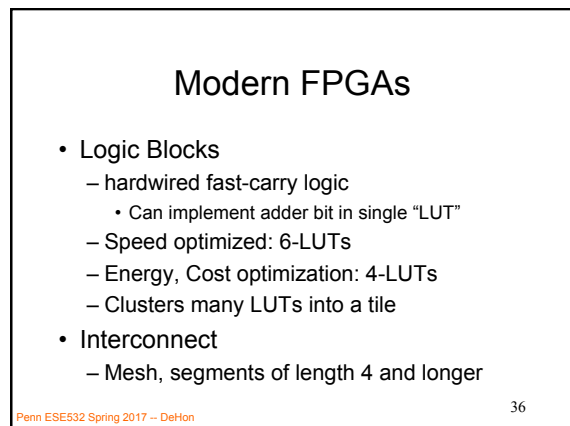
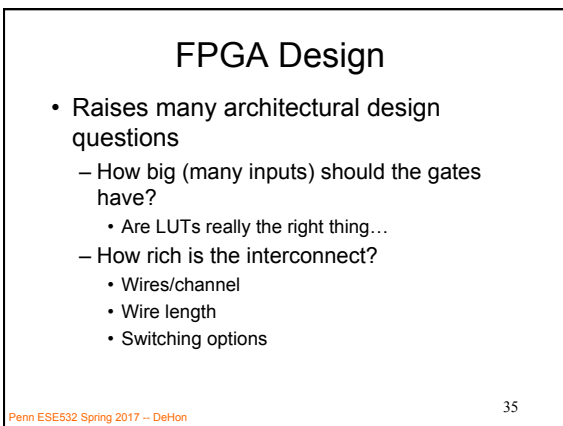
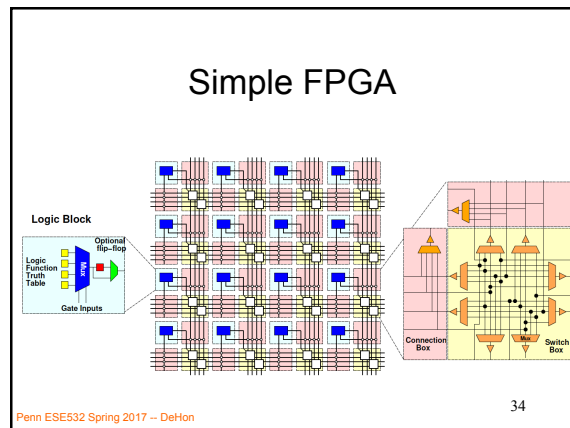
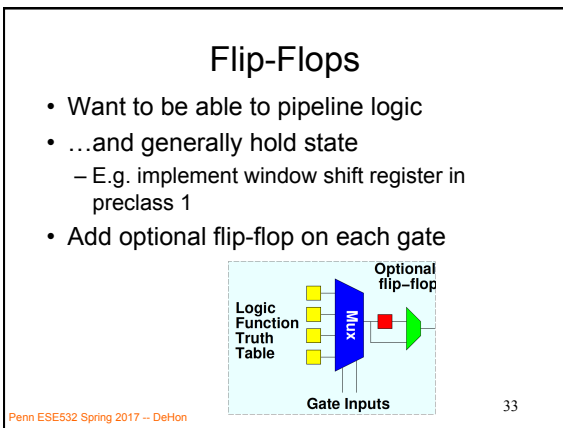
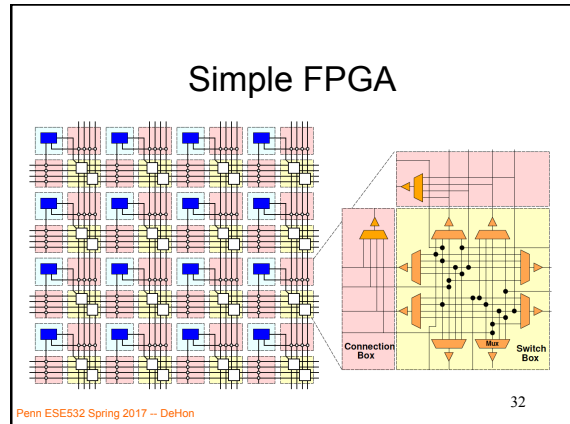
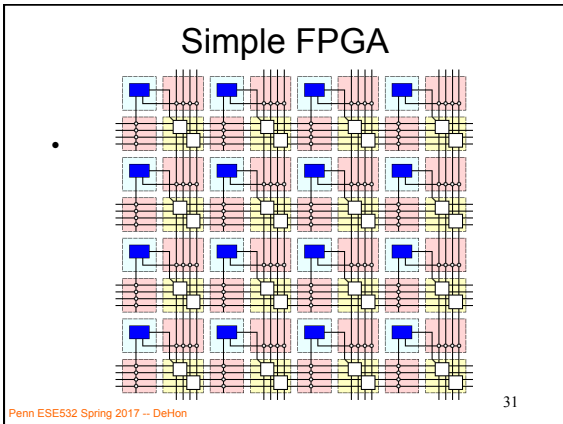
Penn ESE532 Spring 2017 -- DeHon

Crossbar Interconnect

- Crossbar interconnect is 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)

30

Penn ESE532 Spring 2017 -- DeHon



Instruction View

- FPGA: collection of bit-processing units
 - Must do the same thing on every cycle
 - Each BPU unit can do its own thing
 - Has own (configured) instruction
- SIMD: collection of bit-processing units
 - Can do a different thing on every cycle
 - All BPU must do the same thing on a cycle
 - Single broadcast instruction to all BPUs

Penn ESE532 Spring 2017 -- DeHon

37

FPGA vs. SIMD Instruction Architecture

	Each cycle	Each BPU
FPGA	same	unique
SIMD	unique	same

FPGAs: instructions do not change on cycle-by-cycle basis
SIMD: instructions shared across all units on a cycle

Penn ESE532 Spring 2017 -- DeHon

38

More than LUTs

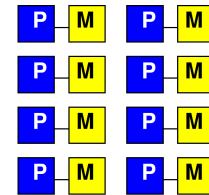
- Should there be more than LUTs in the “array” fabric?
- What else might we want?

Penn ESE532 Spring 2017 -- DeHon

39

Embedded Memory

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



Penn ESE532 Spring 2017 -- DeHon

40

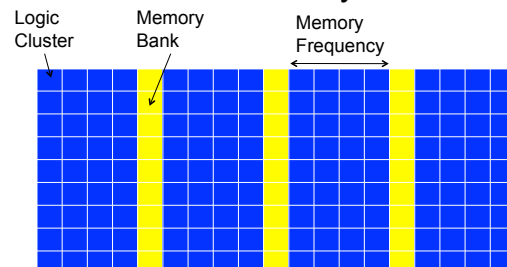
Embed Memory in Array

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

Penn ESE532 Spring 2017 -- DeHon

41

Embedded Memory in FPGA



Penn ESE532 Spring 2017 -- DeHon

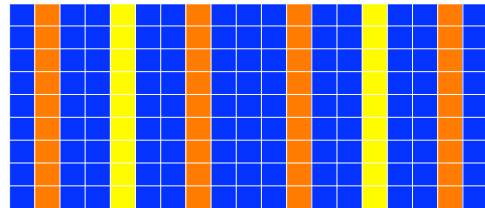
42

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

Multiplier Integration

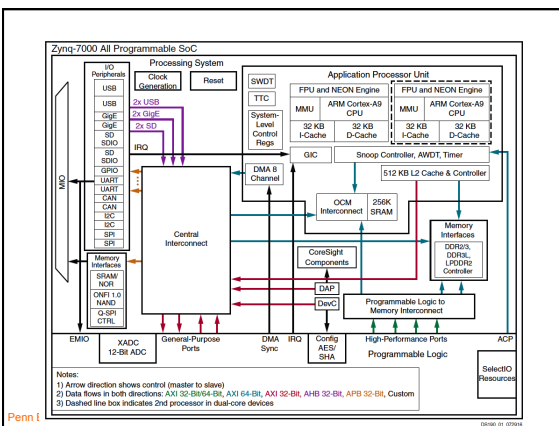
- Integrate like memories
 - Replace columns



More FPGA Architecture Design Questions

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

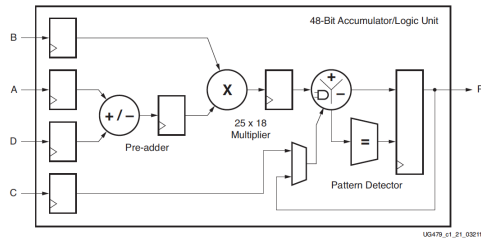
Zynq



XC7Z020

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

DSP48



Xilinx UG479 DSP48E1 User's Guide 49

Penn ESE532 Spring 2017 -- DeHon

Preclass 4

Resoruces	Cycle	per second
50,000 adder bits	0.5 GHz	
128 × 2 adder bits	1.0 GHz	
200 multiply-accumulates	0.5 GHz	
8 × 2 multiply-accumulates	1.0 GHz	

Penn ESE532 Spring 2017 -- DeHon

50

Compute Capacity

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

51

Penn ESE532 Spring 2017 -- DeHon

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?

52

Penn ESE532 Spring 2017 -- DeHon

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

53

Penn ESE532 Spring 2017 -- DeHon

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

54

Penn ESE532 Spring 2017 -- DeHon

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

Admin

- Reading for Day 10 on canvas
- HW5 due Friday