

# ESE532: System-on-a-Chip Architecture

Day 6: September 18, 2019  
Data-Level Parallelism



## Today

- Data-level Parallelism
  - For Parallel Decomposition
  - Architectures
  - Concepts
  - NEON

## Message

- Data Parallelism easy basis for decomposition
- Data Parallel architectures can be compact
  - pack more computations onto a fixed-size IC die
  - OR perform computation in less area

## Preclass 1

- 400 news articles
- Count total occurrences of a string
- How can we exploit data-level parallelism on task?
- How much parallelism can we exploit?

## Parallel Decomposition

## Data Parallel

- Data-level parallelism can serve as an organizing principle for parallel task decomposition
- Run computation on independent data in parallel

## Exploit

- Can exploit with
  - Threads
  - Pipeline Parallelism
  - Instruction-level Parallelism
  - Fine-grained Data-Level Parallelism

## Performance Benefit

- Ideally linear in number of processors (resources)
- Resource Bound:
  - $T_{dp} = (T_{single} \times N_{data}) / P$
- $T_{single}$  = Latency on single data item
- $T_{dp} = \max(T_{single}, N_{data} / P)$

## SPMD

- Single Program Multiple Data
- Only need to write code once
  - Get to use many times

## Preclass 2 Common Examples

- What are common examples of DLP?
  - Simulation
  - Numerical Linear Algebra
  - Signal or Image Processing
  - Image Processing
  - Optimization

## Hardware Architectures

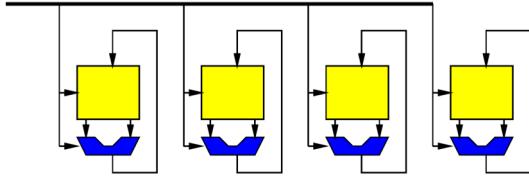
## Idea

- If we're going to perform the same operations on different data, exploit that to reduce area, energy
- Reduced area means can have more computation on a fixed-size die.

## SIMD

- Single Instruction Multiple Data

Shared Instruction

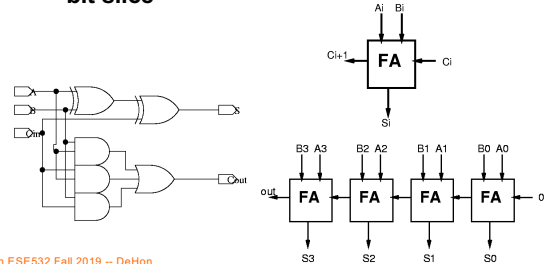


Penn ESE532 Fall 2019 -- DeHon

13

## Ripple Carry Addition

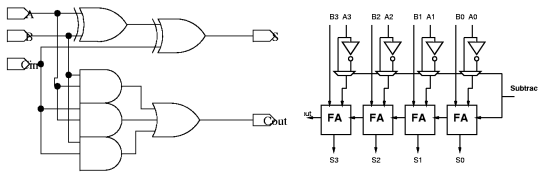
- Can define logic for each bit, then assemble:  
– bit slice



Penn ESE532 Fall 2019 -- DeHon

## Arithmetic Logic Unit (ALU)

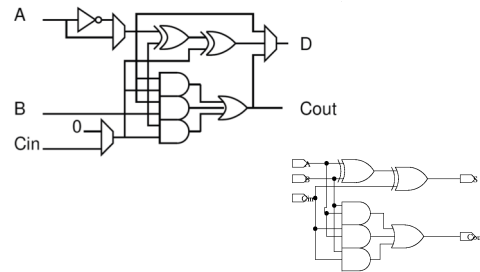
- Observe:  
– with small tweaks can get many functions with basic adder components



Penn ESE532 Fall 2019 -- DeHon

15

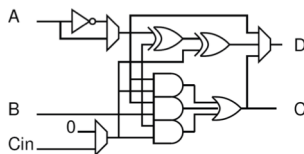
## Arithmetic and Logic Unit



Penn ESE532 Fall 2019 -- DeHon

16

## ALU Functions



- A+B w/ Carry
- B-A
- A xor B (squash carry)
- A\*B (squash carry)
- /A

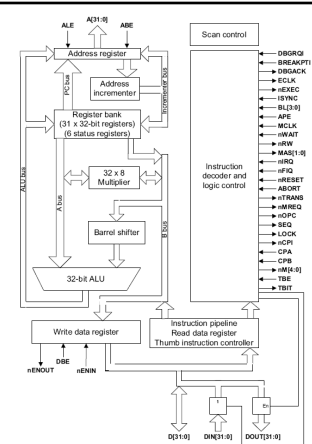
**Key observation:** every ALU bit does the same thing on different bits of the data word(s).

Penn ESE532 Fall 2019 -- DeHon

17

## ARM v7 Core

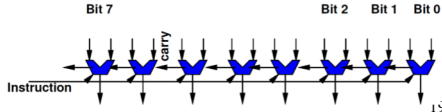
- ALU is Key compute operator in a processor



Penn ESE532 Fall 2019 -- DeHon

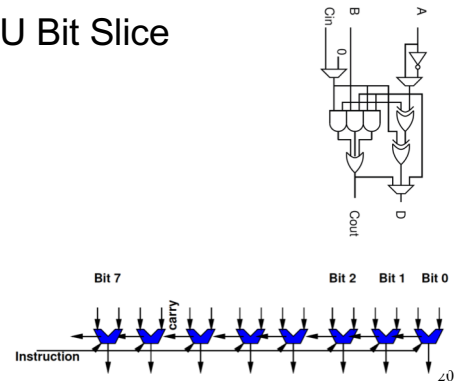
## W-bit ALU as SIMD

- Familiar idea
- A W-bit ALU ( $W=8, 16, 32, 64, \dots$ ) is SIMD
- Each bit of ALU works on separate bits
  - Performing the same operation on it
    - Trivial to see bitwise AND, OR, XOR
    - Also true for ADD (each bit performing Full Adder)
- Share one instruction across all ALU bits



Penn ESE532 Fall 2019 -- DeHon

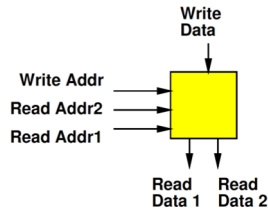
## ALU Bit Slice



Penn ESE532 Fall 2019 -- DeHon

## Register File

- Small Memory
- Usually with multiple ports
  - Ability to perform multiple reads and writes simultaneously
- Small
  - To make it fast (small memories fast)
  - Multiple ports are expensive

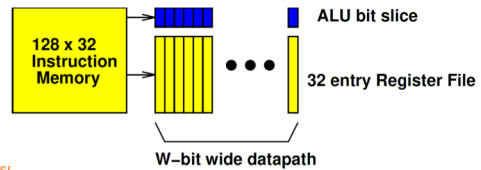


21

Penn ESE532 Fall 2019 -- DeHon

## Preclass 4

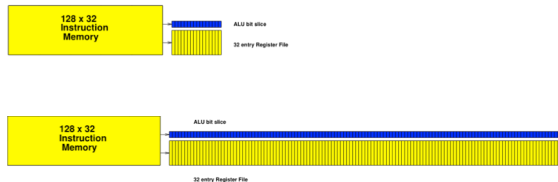
- Number in  $10^8$ 
  - $W=16$
  - $W=128$
- Area  $W=16$ ?
- Area  $W=128$ ?
- Perfect Pack Ratio?
- Why?



Penn ESL

22

## To Scale Comparison

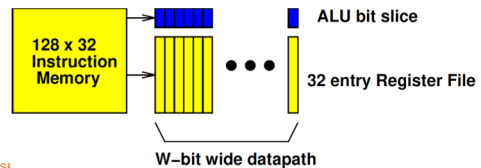


Penn ESE532 Fall 2019 -- DeHon

23

## Preclass 4

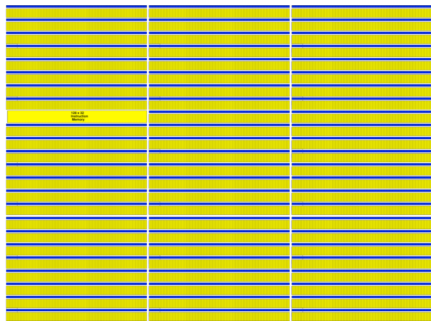
- W for single datapath in  $10^8$ ?
- Compare  $W=128$  perfect pack ratio?
- Perfect 16b pack ratio?



Penn ESI

24

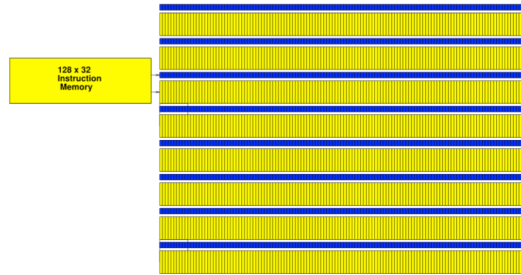
## To Scale W=9088



Penn ESE532 Fall 2019 -- DeHon

25

## To Scale W=1024



Penn ESE532 Fall 2019 -- DeHon

26

## Preclass 6

- What do we get when add 65280 to 257
  - 32b unsigned add?
  - 16b unsigned add?

Penn ESE532 Fall 2019 -- DeHon

27

## ALU vs. SIMD ?

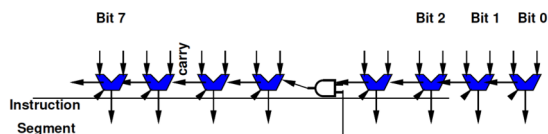
- What's different between
  - 128b wide ALU
  - SIMD datapath supporting eight 16b ALU operations

Penn ESE532 Fall 2019 -- DeHon

28

## Segmented Datapath

- Relatively easy (few additional gates) to convert a wide datapath into one supporting a set of smaller operations
  - Just need to squash the carry at points



Penn ESE532 Fall 2019 -- DeHon

29

## Segmented Datapath

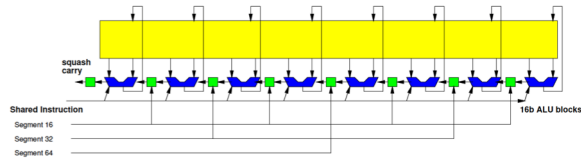
- Relatively easy (few additional gates) to convert a wide datapath into one supporting a set of smaller operations
  - Just need to squash the carry at points
- But need to keep instructions (description) small
  - So typically have limited, homogeneous widths supported

Penn ESE532 Fall 2019 -- DeHon

30

## Segmented 128b Datapath

- 1x128b, 2x64b, 4x32b, 8x16b

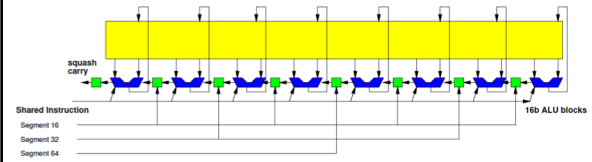


Penn ESE532 Fall 2019 -- DeHon

31

## Terminology: Vector Lane

- Each of the separate segments called a **Vector Lane**
- For 16b data, this provides 8 vector lanes

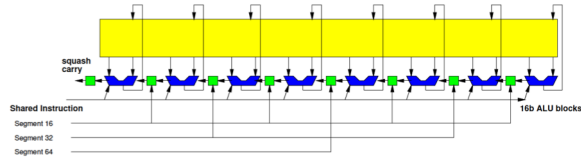


Penn ESE532 Fall 2019 -- DeHon

32

## Performance

- Ideally, pack into vector lanes
- Resource Bound:  $T_{\text{vector}} = N_{\text{op}}/VL$



Penn ESE532 Fall 2019 -- DeHon

33

## Preclass 5: Vector Length

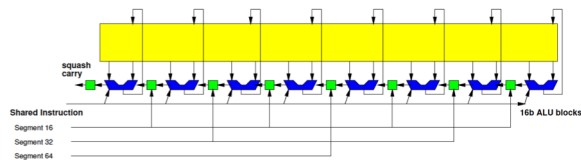
- May not match physical hardware length
- What happens when
  - Vector length > Vector Lanes?
  - Vector length < Vector Lanes?
  - Vector length % (Vector Lanes) != 0
    - E.g. vector length 20, for 8 hwd operators

Penn ESE532 Fall 2019 -- DeHon

34

## Performance

- Resource Bound
- $$T_{\text{vector}} = \text{ceil}(N_{\text{op}}/VL) * T_{\text{cycle}}$$



Penn ESE532 Fall 2019 -- DeHon

35

## Preclass 3: Opportunity

- Don't need 64b variables for lots of things
- Natural data sizes?
  - Audio samples?
  - Input from A/D?
  - Video Pixels?
  - X, Y coordinates for 4K x 4K image?

Penn ESE532 Fall 2019 -- DeHon

36

## Vector Computation

- Easy to map to SIMD flow if can express computation as operation on vectors
  - Vector Add
  - Vector Multiply
  - Dot Product

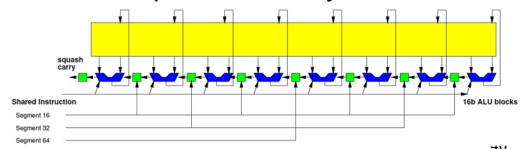
## Concepts

## Terminology: Scalar

- Simple: non-vector
- When we have a vector unit controlled by a normal (non-vector) processor core often need to distinguish:
  - Vector operations that are performed on the vector unit
  - Normal=non-vector=scalar operations performed on the base processor core

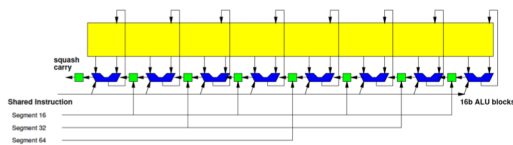
## Vector Register File

- Need to be able to feed the SIMD compute units
  - Not be bottlenecked on data movement to the SIMD ALU
- Wide RF to supply
- With wide path to memory



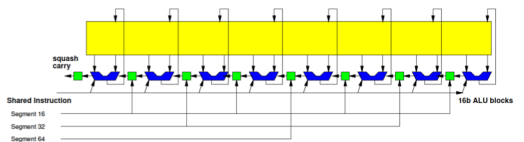
## Point-wise Vector Operations

- Easy – just like wide-word operations (now with segmentation)



## Point-wise Vector Operations

- ...but alignment matters.
- If not aligned, need to perform data movement operations to get aligned



## Ideal

- for (i=0;i<64;i++)
  - $c[i]=a[i]+b[i]$
- No data dependencies
- Access every element
- Number of operations is a multiple of number of vector lanes

## Skipping Elements?

- How does this work with datapath?
  - Assume loaded  $a[0], a[1], \dots, a[63]$  and  $b[0], b[1], \dots, b[63]$  into vector register file
- for (i=0;i<64;i=i+2)
  - $c[i/2]=a[i]+b[i]$

## Stride

- Stride: the distance between vector elements used
- for (i=0;i<64;i=i+2)
  - $c[i/2]=a[i]+b[i]$
- Accessing data with stride=2

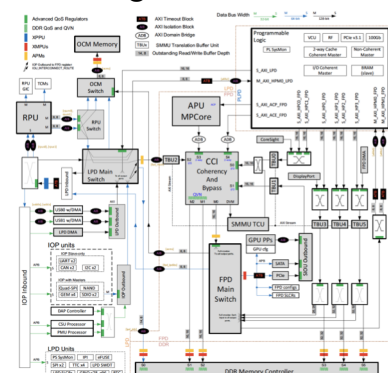
## Load/Store

- Strided load/stores
  - Some architectures will provide strided memory access that compact when read into register file
- Scatter/gather
  - Some architectures will provide memory operations to grab data from different places to construct a dense vector

## Neon

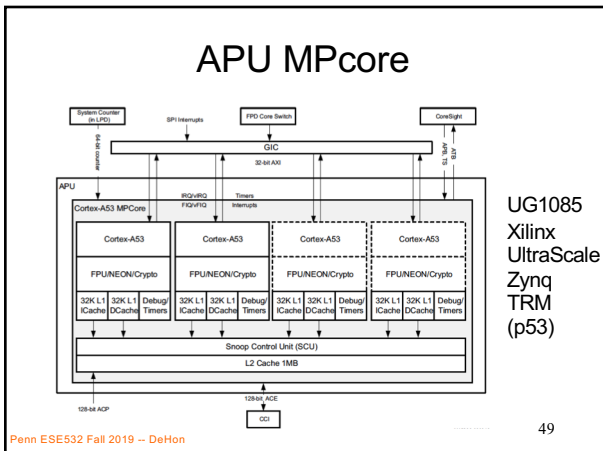
ARM Vector Accelerator on Zynq

## Programmable SoC



UG1085  
Xilinx  
UltraScale  
Zynq  
TRM  
(p27)





### Neon Vector

- 128b wide register file, 16 registers
- Support
  - 2x64b
  - 4x32b (also Single-Precision Float)
  - 8x16b
  - 16x8b

50

Penn ESE532 Fall 2019 -- DeHon

### Sample Instructions

- VADD – basic vector
- VCEQ – compare equal
  - Sets to all 0s or 1s, useful for masking
- VMIN – avoid using ifs
- VMLA – accumulating multiply
- VPADAL – maybe useful for reduce
  - Vector pair-wise add
- VEXT – for “shifting” vector alignment
- VLDn – deinterleaving load

51

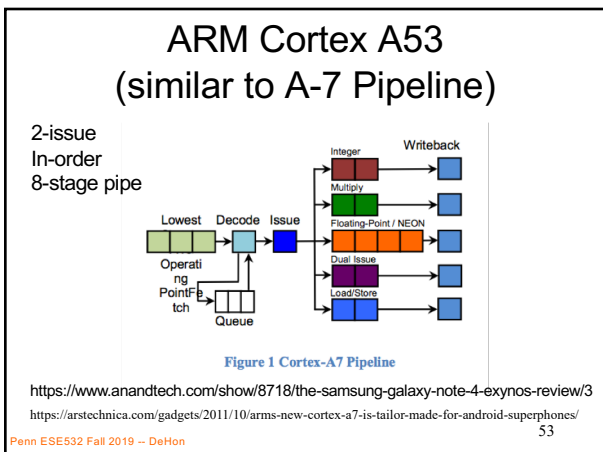
Penn ESE532 Fall 2019 -- DeHon

### Neon Notes

- Didn't see
  - Vector-wide reduce operation
- Do need to think about operations being pipelined within lanes

52

Penn ESE532 Fall 2019 -- DeHon



### Big Ideas

- Data Parallelism easy basis for decomposition
- Data Parallel architectures can be compact – pack more computations onto a chip
  - SIMD, Pipelined
  - Benefit by sharing (instructions)
  - Performance can be brittle
    - Drop from peak as mismatch

55

Penn ESE532 Fall 2019 -- DeHon

## Admin

- SDSoC available on Linux machines
  - See piazza
- Reading for Day 7 online
- HW3 due Friday
- HW4 out