# ESE532:
## System-on-a-Chip Architecture

Day 6: September 20, 2017
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 die

---

## 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

7

## Thread Exploit DP

- How exploit threads for data-parallel text search?

8

## SPMD

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

9

## Pipeline Exploit

- How exploit hardware pipeline for text search?

10

## Pipeline Text Search

**Text Document**

**Match Targets**

**and**

11

## Common Examples

- What are common examples of DLP?
  - Simulation
  - Numerical Linear Algebra
  - Graphics
  - Signal Processing
  - Image Processing
  - Optimization
  - Other?
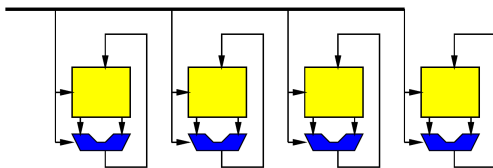
12

2

# Hardware Architectures

13

# Idea

- If we're going to perform the same operations on different data, exploit that to reduce area, energy

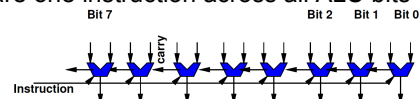- Reduced area means can have more computation on a fixed-size die.

14

# SIMD

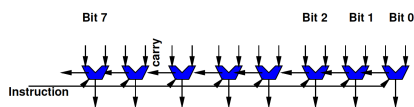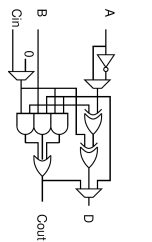- Single Instruction Multiple Data

**Shared Instruction**

15

# W-bit ALU as SIMD

- Familiar idea
- A W-bit ALU (W=8, 16, 32, 64, …) 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)
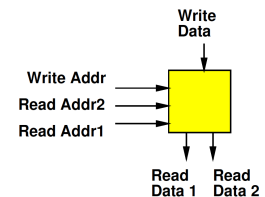- Share one instruction across all ALU bits

Bit 7    carry    Bit 2    Bit 1    Bit 0

Instruction

16

# ALU Bit Slice

Cin    B    A

Cout    D

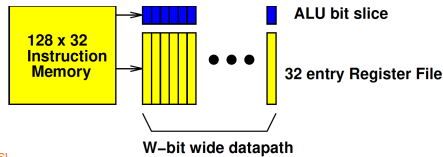Bit 7    carry    Bit 2    Bit 1    Bit 0

Instruction

17

# 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

Write Data

Write Addr

Read Addr2

Read Addr1

Read Data 1    Read Data 2

18

## Preclass 2

- Area W=16?
- Area W=128?

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

**128 x 32 Instruction Memory**

**ALU bit slice**

**32 entry Register File**

**W–bit wide datapath**

19

## Preclass 2

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

- Compare W=128 perfect pack ratio?

**128 x 32 Instruction Memory**

**ALU bit slice**
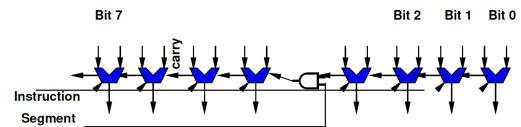
**32 entry Register File**

**W–bit wide datapath**

20

## ALU vs. SIMD ?

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

21

## 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

Bit 7          Bit 2   Bit 1   Bit 0

carry

Instruction

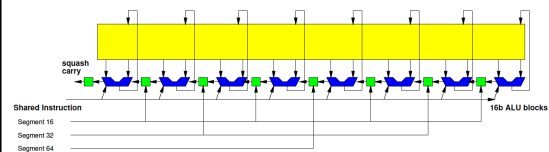Segment

22

## 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

23

## Segmented 128b Datapath

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

squash carry

Shared Instruction                                   16b ALU blocks

Segment 16
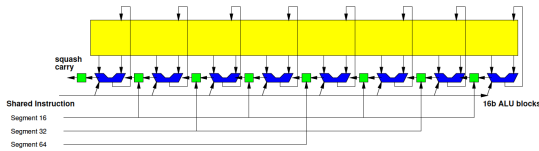Segment 32
Segment 64

24

## Terminology: Vector Lane

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

25

## 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?

26

## Vector Computation

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

27

## Concepts

28

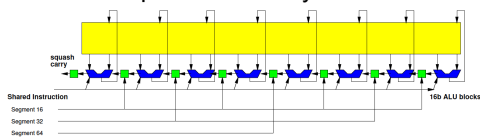## Vector Register File

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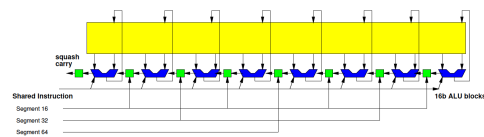- Wide RF to supply
- With wide path to memory

## Point-wise Vector Operations

- Easy – just like wide-word operations (now with segmentation)
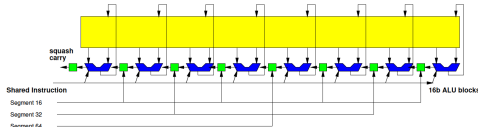
30

5

## Point-wise Vector Operations

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

31

## Ideal

- for (i=0;i<64;i=i++)
  - c[i]=a[i]+b[i]

- No data dependencies
- Access every element
- Number of operations is a multiple of number of vector lanes

32

## Vector Length

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

33

## Skipping Elements?

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

34

## 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

35

## 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

36

# Dot Product

- What happens when need a dot product?
- res=0;
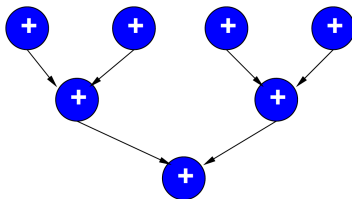- for (i=0;i<N;i++)
  – res+=a[i]*b[i]

# Reduction

- Common operations where want to perform a combining operation to reduce a vector to a scalar
  – Sum values in vector
  – AND, OR
- Reduce Operation
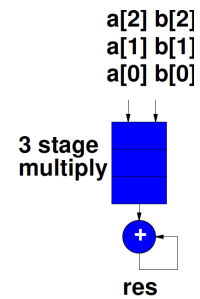
# Reduce Tree

- Efficiently handled with reduce tree

# Reduce in Pipeline

- Comes almost for free in pipeline

**a[2] b[2]**
**a[1] b[1]**
**a[0] b[0]**

**3 stage multiply**

**res**

# Vector Reduce Instruction

- Usually include support for vector reduce operation
  – Doesn't need to add much to delay
  – Maybe even faster than performing larger operation
    - 8 16x16 multiplies with sum reduce less complex than one 128x128 multiply
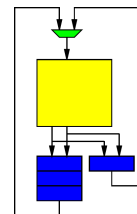    - …can exploit datapath of larger operation

# Dot Product Revisited

- for (i=0;i<N;i++)
  – res+=a[i]*b[i]
- With 3 cycle pipelined multiply
- What happens if try to implement dot product as:
  – MPY R0, R4, R14
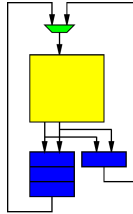  – ADD R14, R15, R15
  – MPY R1, R5, R14
  – ADD R14, R15, R15
  – …
- a in R0—R4
- b in R4—R7

## Dot Product Revisited

- How should order (reformulate) instructions exploiting data-level parallelism?

- for (i=0;i<N;i++)
  - res+=a[i]*b[i]
- a in R0—R4
- b in R4—R7
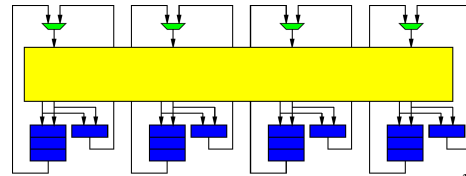
43

## Pipelined Vector Units

- Will get both pipelining and parallel vector lanes
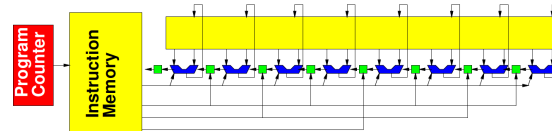- Exploit data-level parallelism for both

44

## Conditionals?

- What happens if want to do something different?
- For (i=0;i<8;i++)
  - if (a[i]<b[i])
    - d[i]=a[i]+c[i]
  - else
    - d[i]=b[i]+c[i]

45

## Conditionals

- Only have one Program Counter
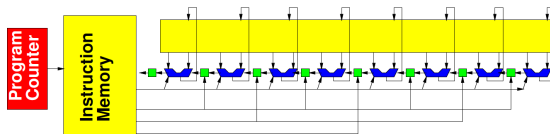  - Cannot implement conditional via branching

**Program Counter**

**Instruction Memory**

46

## Conditionals

- Only have one instruction
  - Cannot perform separate operations on each ALU in datapath

**Program Counter**

**Instruction Memory**

47

## Conditionals

- Only have one Program Counter
  - Cannot implement conditional via branching
- Only have one instruction
  - Cannot perform separate operations on each ALU in datapath
- Must perform an invariant operation sequence
- Simple answer: prevent using SIMD unit
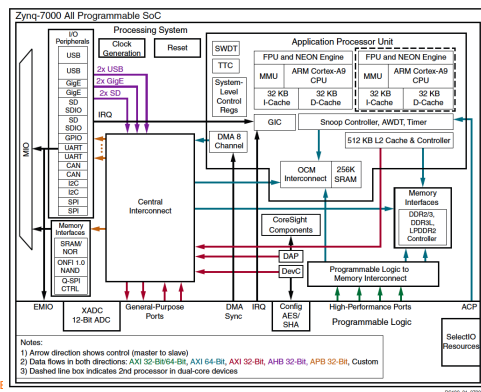- Better: predicated execution

48

8

## Predicated Operation

- Many architectures will provide a predicated operation
- Only perform operation when predicate matches instruction

- p[i]=a[i]<b[i]
- p[i]:   d[i]=c[i] + a[i]
- ~p[i]: d[i]=c[i] + b[i]

49

## Predicated Operation

- What does this do to instructions must be issued?
- What does this do to efficiency?
  - Useful operations performed per cycle

- p[i]=a[i]<b[i]
- p[i]:   d[i]=c[i] + a[i]
- ~p[i]: d[i]=c[i] + b[i]

50

## Neon

ARM Vector Accelerator on Zynq

51

## Neon Vector

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

53

## Sample Instructions

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

54

9

## Neon Notes

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

## 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

## Admin

- No reading for day 7
- HW3 due Friday
- HW4 out