# ESE532:
## System-on-a-Chip Architecture
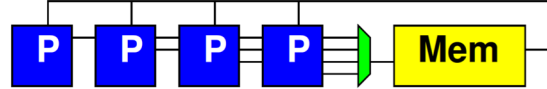
Day 12:  October 10, 2018
Data Movement
(Interconnect, DMA)

Penn

---

## Preclass 1

- N processors
- Each: 1 read, 10 cycle compute, 1 write
- Memory: 1 read or write per cycle
- How many processors can support before saturate memory capacity?
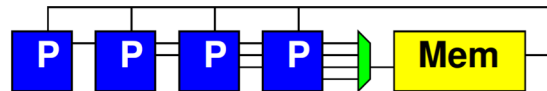
2

---

## Schedule Memory Port

| 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| P1.1 write | P1.2 read | P2.1 write | P2.2 read | P3.1 write | P3.2 read | P4.1 write | P4.2 read | P5.1 write | P5.2 read | P6.1 write | P6.2 read | P1.2 write | P1.3 read |

P1 compute f on 2nd iteration

P2 compute f on 2nd iteration

3

---

## Bottleneck
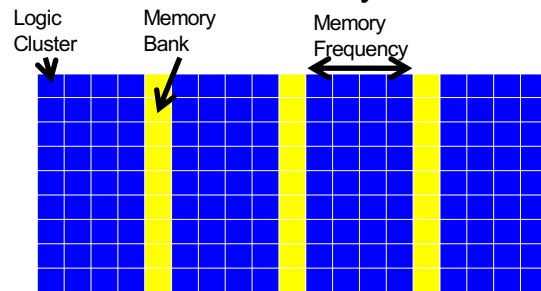
- Sequential access to a common memory can become the bottleneck

4

---

## Previously

- Want data in small memories
  - Low latency, high bandwidth
- FPGA has many memories all over fabric

5

---

## Embedded Memory in FPGA

Logic Cluster     Memory Bank     Memory Frequency

XC7Z020 (Zed Board) has 140 36Kb BRAMs

6

---

1

## Previously

- Want data in small memories
  - Low latency, high bandwidth
- FPGA has many memories all over fabric
- Want C arrays in small memories
  - Partitioned so can perform enough reads (writes) in a cycle to avoid memory bottleneck

## Today

- Interconnect Infrastructure
- Data Movement Threads
- Peripherals
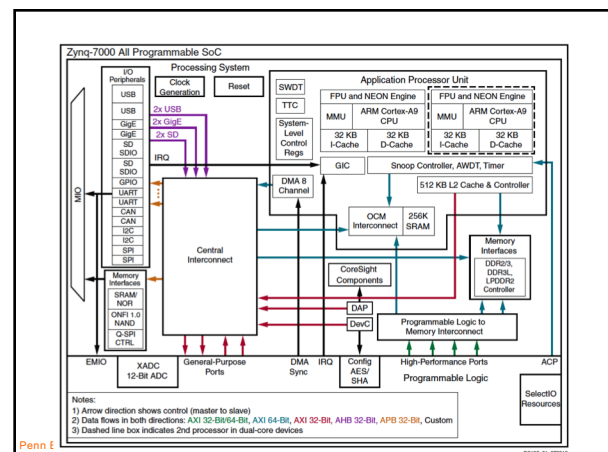- DMA -- Direct Memory Access

## Message

- Need to move data
- Shared interconnect to make physical connections
- Useful to move data as separate thread of control
  - Dedicating a processor is inefficient
  - Useful to have dedicated data-movement hardware: Direct Memory Access (DMA)

## Memory and I/O Organization

- Architecture contains
  - Large memories
    - For density, necessary sharing
  - Small memories local to compute
    - For high bandwidth, low latency, low energy
  - Peripherals for I/O
- Need to move data
  - Among memories and I/O
    - Large to small and back
    - Among small
    - From Inputs, To Outputs

## Term: Peripheral

- "On the edge (or perhiphery) of something"
- Peripheral device – device used to put information onto or get information off of a computer
  - E.g.
    - Keyboard, mouse, modem, USB flash drive, …

## Memory and I/O Organization

- Architecture contains
  - Large memories
    - For density, necessary sharing
  - Small memories local to compute
    - For high bandwidth, low latency, low energy
  - **Peripherals** for I/O
- Need to move data
  - Among memories and I/O
    - Large to small and back
    - Among small
    - From Inputs, To Outputs

## How move data?

- Abstractly, using stream links.
- Connect stream between producer and consumer.

- Ideally: dedicated wires

## Dedicated Wires?

- What might prevent us from having dedicated wires between all communicating units?
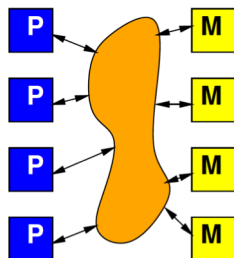
## Making Connections

- Cannot always be dedicated wires
  - Programmable
  - Wires take up area
  - Don't always have enough traffic to consume the bandwidth of point-to-point wire
  - May need to serialize use of resource
    - E.g. one memory read per cycle
  - Source or destination may be sequentialized on hardware
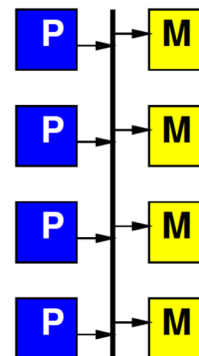
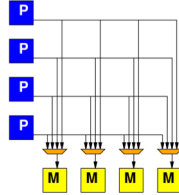## Model

- Programmable, possibly shared interconnect

## Simple Realization

Shared Bus
- Write to bus with address of destination
- When address match, take value off bus
- Pros?
- Cons?

## Alternate: Crossbar

- Provide programmable connection between all sources and destinations
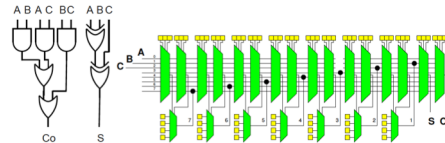- Any destination can be connected to any single source

19

## Simplistic FPGA (illustrate possibility)

- Every LUT input has a mux
- Every such mux has m=(N+I) 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

20

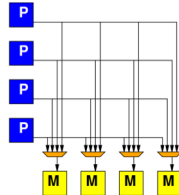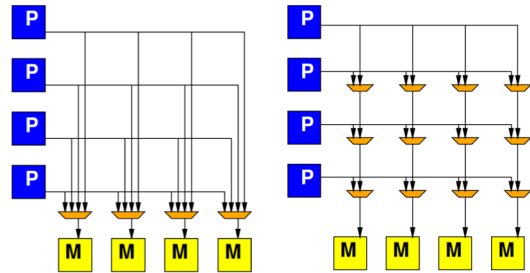## Alternate: Crossbar

- Provide programmable connection between all sources and destinations
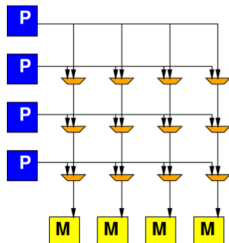- Any destination can be connected to any single source

21

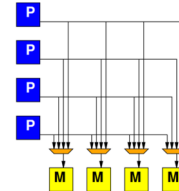## Crossbar

22

## Preclass 2

- K-input, O-output Crossbar
- How many 2-input muxes?

23

## Crossbar

- Provides high bandwidth
  – Minimal blocking
- Costs large amounts of area
  – Grows fast with inputs, outputs
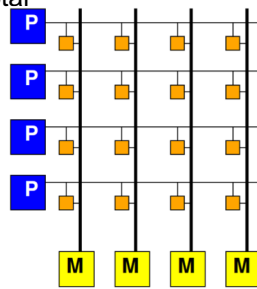
24

## General Interconnect

- Generally, want to be able to parameterize designs
- Here: tune area-bandwidth
  - Control how much bandwidth provide

25

## Interconnect

- How might get design points between bus and crossbar?
- How could reduce number
  - Inputs to crossbar?
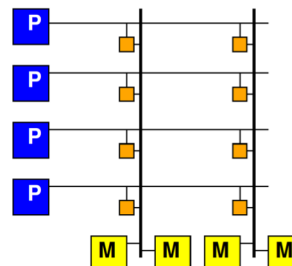  - Outputs from crossbar?

26

## Multiple Busses

- Think of crossbar as one bus per output
- Simple bus is one bus total
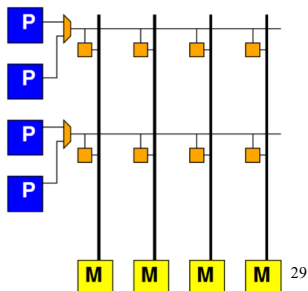- In between,
  - How many simultaneous busses support?

## Share Crossbar Outputs

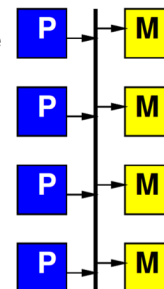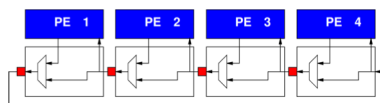- Group set of outputs together on a bus

28

## Share Crossbar Inputs

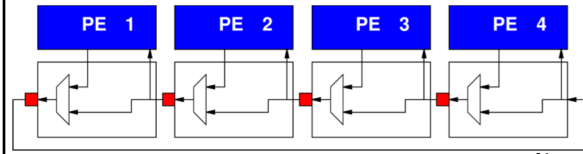- Group number of inputs together on an input port to crossbar

29

## Delay

- Delay proportional to distance
- Pipeline bus to keep cycle time down
  - Take many cycles to get travel long distance
  - …but fewer cycles when distance small

PE 1  PE 2  PE 3  PE 4

30

5

## Local Interconnect

- How many cycles from:
  - PE3 to PE2
  - PE3 to PE1
  - PE3 to PE4

31

## Mesh

32

## Hierarchical Busses

33

## Interconnect

- Will need an infrastructure for programmable connections
- Rich design space to tune area-bandwidth-locality
  - Will explore more later in course

34

## Long Latency Memory Operations

35

## Day 3

- Large memories are slow
  - Latency increases with memory size
- Distant memories are high latency
  - Multiple clock-cycles to cross chip
  - Off-chip memories even higher latency

36

## Day 3, Preclass 2

- 10 cycle latency to memory
- If must wait for data return, latency can degrade throughput
- 10 cycle latency + 10 op + (assorted)
  - More than 20 cycles / result

```
for(i=0;i<MAX;i++) {
    in=a[i]; // memory read
    out=f(in); // 10 cycle compute
    b[i]=out;
}
```

7

## Preclass 3

- Throughput using 3 threads?

```
P1:  for(i=0;i<MAX;i++) Astream.write(a[i]);
P2:  while(1) {Astream.read(aval); Bstream.write(f(aval));}
P3:  for(i=0;i<MAX;i++) Bstream.read(b[i]);
```

38

## Fetch (Write) Threads

- Potentially useful to move data in separate thread
- Especially when
  - Long (potentially variable) latency to data source (memory)
- Useful to split request/response

39

## Peripherals

40

## Input and Output

- Typical SoC has I/O with external world
  - Sensors
  - Actuators
  - Keyboard/mouse, display
  - Communications
- Also accessible from interconnect

41

## Masters and Slaves

- Two kinds of entities on interconnect
- Master – can initiate requests
  - E.g. processor that can perform a read or write
- Slaves – can only respond to requests
  - E.g. memory that can return the read data from a read request

42

## Simple Peripheral Model

- Peripherals are slave devices
  - Masters can read input data
  - Masters can write output data
  - To move data, master (e.g. processor) initiates

P → M → usb

P → M → ethernet

P → M → A/D

P → M → HDMI

## Simple Model Implications

- What implication to processor grabbing/moving each input (output) value?

44

P → M → usb

P → M → ethernet

P → M → A/D

P → M → HDMI

## Timing Demands

- Must read each input before overwritten
- Must write each output within real-time window
- Must guarantee processor scheduled to service each I/O at appropriate frequency
- How many cycles between 32b input words for 1Gb/s network and 32b, 1GHz processor?

45

## Refine Model

- Give each peripheral local FIFO
- Processor must still move data
- How does this change requirements and impact?

P → M → usb

P → M → ethernet

P → M → A/D

P → M → HDMI

46

## DMA

Direct Memory Access

47

## Preclass 4a

P1:   for(i=0;i<MAX;i++) Astream.write(a[i]);

WriteAstart   NewAddr[24]   WriteAstop

FIFO_Has_Space

ReadAddress[24]                  load
ReadRequest     Counter   Register    FIFO_Write

Data[32]                           FIFO_DataIn[32]
            incr_cntr    <   not–done
DataPresent

```
int *p;
P1:   for(p=&(a[0]);p<&(a[MAX]);i++) Astream.write(*p);
```
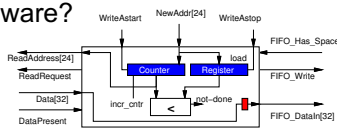
48

8

## Preclass 4

- How much hardware?
  - Counter bits?
  - Registers?
  - Comparators?
  - Control Logic gates? (4cd)



- Compare to MicroBlaze
  - small RISC Processor optimized for Xilinx
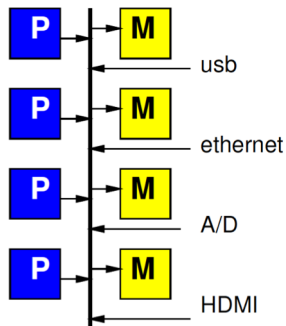  - minimum config 630 6-LUTs

## Observe

- Modest hardware can serve as data movement thread
  - Much less hardware than a processor
  - Offload work from processors

- Small hardware allow peripherals to be **master** devices on interconnect

## DMA

- Direct Memory Access (DMA)
- "Direct" – inputs (and outputs) don't have to be indirectly handled by the processor between memory and I/O
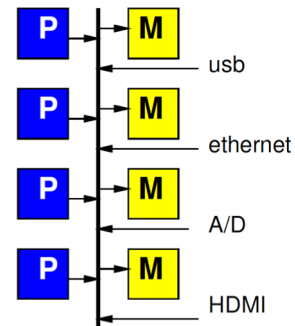- I/O unit directly reads/write memory

## DMA

- Direct Memory Access (DMA)
- Peripheral as Master
  - Can write **directly** into (read from) memory
  - Saves processor from copying
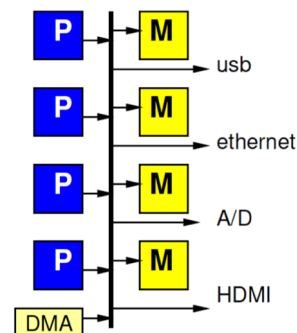  - Reduces demand to schedule processor to service

## DMA Engine

- Data Movement Thread
  - Specialized Processor that moves data
- Act independently
- Implement data movement
- Can build to move data between memories (Slave devices)
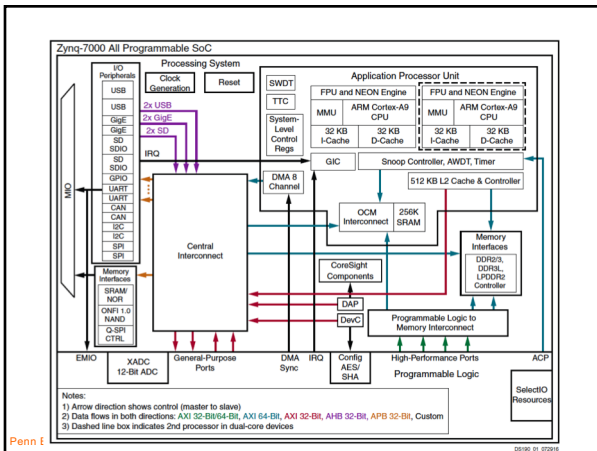- E.g., Implement P1, P3 in Preclass 3

## DMA Engine

9

## Programmable DMA Engine

- What copy from?
- How much?
- Where copy to?
- Stride?
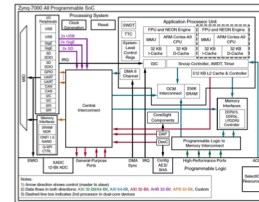- What size data?
- Loop?
- Transfer Rate?

55

## Multithreaded DMA Engine

- One copy task not necessarily saturate bandwidth of DMA Engine
- Share engine performing many transfers (channels)
- Separate transfer state for each
  – Hence thread (or channel)
- Swap among threads
  – Simplest: round-robin:
    - 1, 2, 3, .. K, 1, 2, 3, .. K, 1, ….

56



## Hardwired and Programmable

- Zynq has hardwired DMA engine
- Can also add data movement engines (Data Movers) in FPGA fabric

## Example

- Networking Application



- Header on processor
- Payload (encrypt, checksum) on FPGA
- DMA from ethernet→main memory
- DMA main memory→BRAM
- Stream between payload components
- DMA from chksum to ethernet out

59

## Automation

- SDSoC will automatically take care of DMA of memory to and from accelerators in FPGA fabric
  – Inserting logic, programming DMA
- Mostly need to be aware is happening
- Have some options to control with pragmas
- May not handle sophisticated communication with I/O devices…

60

10

## Big Ideas

- Need to move data
- Shared Interconnect to make physical connections – can tune area/bw/locality
- Useful to
  - move data as separate thread of control
  - Have dedicated data-movement hardware: DMA

61

## Admin

- Day 13
  - Chapter nine of *Parallel Programming for FPGAs* (available on web)
  - DRAM reading if not read on Day 3
- HW5 due Friday
- HW6 out

- Clear room for recitation at noon
- Turn in feedback sheets

62