

# ESE532: System-on-a-Chip Architecture

Day 13: October 17, 2022  
Vitis OpenCL  
Data Transfer Model



Penn ESE5320 Fall 2022 -- DeHon

1

## Previously

- Value of small, local memories
- Tune local memories in C
- DMA as data-movement threads

Penn ESE5320 Fall 2022 -- DeHon

2

2

## Today

Vitis/OpenCL Interface

- Execution Model (Part 1)
- Host-Side Programming (Part 2)
- FPGA-Side Programming (Part 3)

Penn ESE5320 Fall 2022 -- DeHon

3

3

## Message

- Vitis uses an explicit DMA transfer interface
- Can build our streaming model on top of primitives offered
- Will need to tune parameters and use carefully to get maximize performance
  - Good and bad → largely exposed

Penn ESE5320 Fall 2022 -- DeHon

4

4

## OpenCL

- OpenCL – Open Computing Language
  - Standard for programming accelerators
  - Open alternative to CUDA
  - Write things in one language
    - Compile to many compute platforms
- Xilinx(AMD) uses for host interface to accelerator
  - Won't be using for kernel computation on FPGA

Penn ESE5320 Fall 2022 -- DeHon

5

5

## Execution Model(s)

Part 1

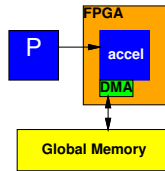
Penn ESE5320 Fall 2022 -- DeHon

6

6

## Model

- Host and Accelerator
- Host control
  - Set accelerator arguments
  - DMA bulk data to and from accelerator
  - Tell it to start
  - DMA, Accelerator runs as independent thread
  - Synchronize on data return



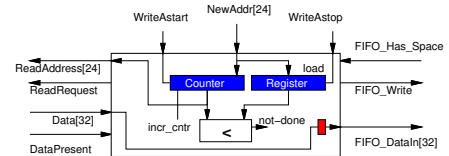
Penn ESE5320 Fall 2022 -- DeHon

7

7

## Remember: Preclass 4a

P1: for(i=0;i<MAX;i++) Astream.write(a[i]);  
Where do we set DMA arguments?



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

Penn ESE5320 Fall 2022 -- DeHon

8

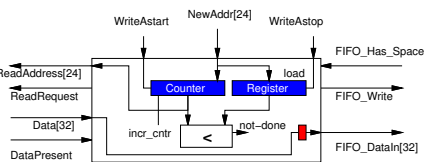
8

## Remember: Preclass 4a

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

Set arguments to accelerator

Accelerator Uses to grab inputs



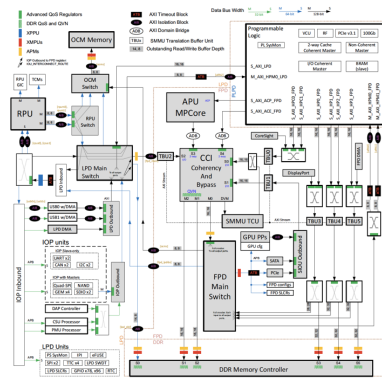
int \*p;  
P1: for(p=&(a[0]);p<&(a[MAX]);p++) Astream.write(\*p);

Penn ESE5320 Fall 2022 -- DeHon

9

9

## Programmable SoC



UG1085  
Xilinx  
UltraScale  
Zynq  
TRM  
(p27)

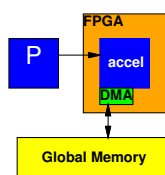
Penn ESE5320

10

10

## Function Call

- Host and Accelerator
- Host control
  - Set accelerator arguments
  - DMA bulk data to accelerator
  - Tell it to start
  - DMA result back
  - Host program wait on accelerator completion



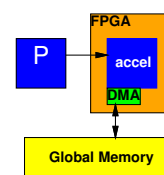
Penn ESE5320 Fall 2022 -- DeHon

11

11

## Separate Accelerator Thread

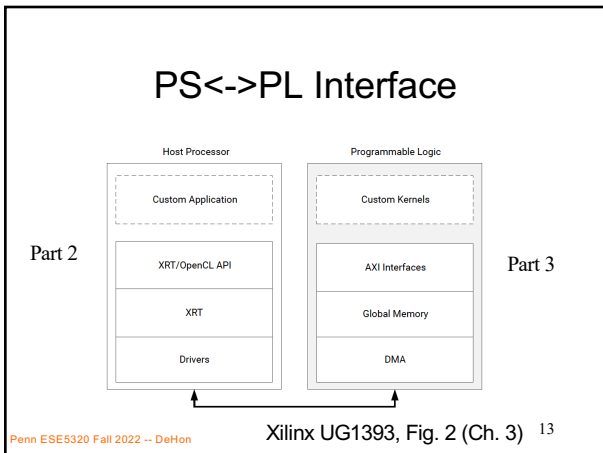
- Host and Accelerator
- Host control
  - Set accelerator arguments
  - Tell it to start
    - accelerator
      - DMA bulk data to accelerator
      - DMA result back
    - Host: execute something else
  - Host program synchronize before use accelerator result



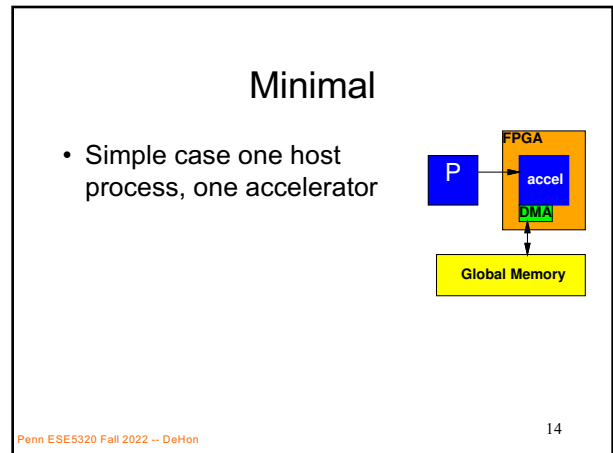
Penn ESE5320 Fall 2022 -- DeHon

12

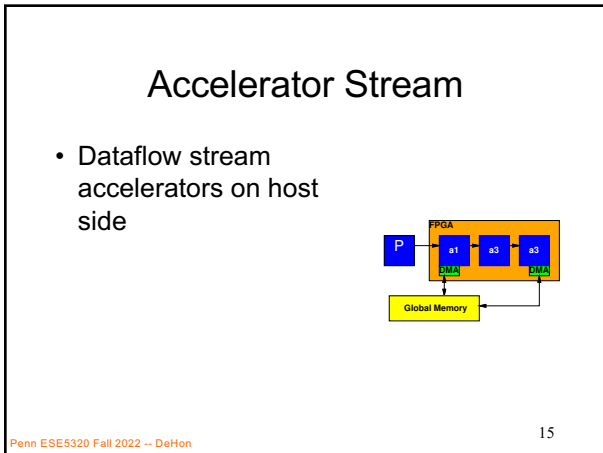
12



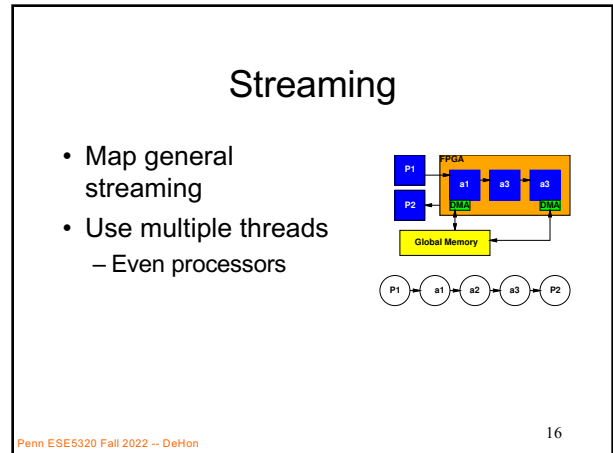
13



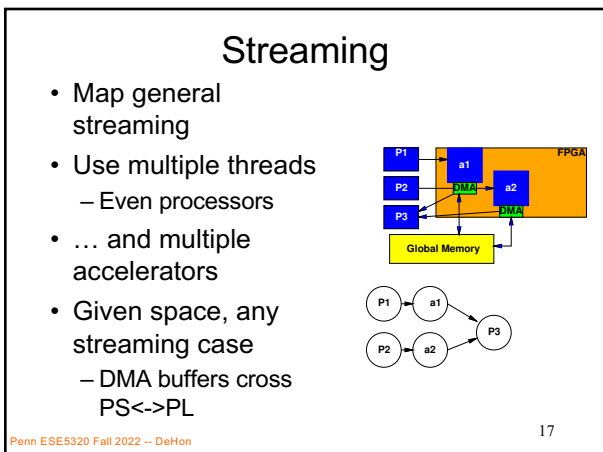
14



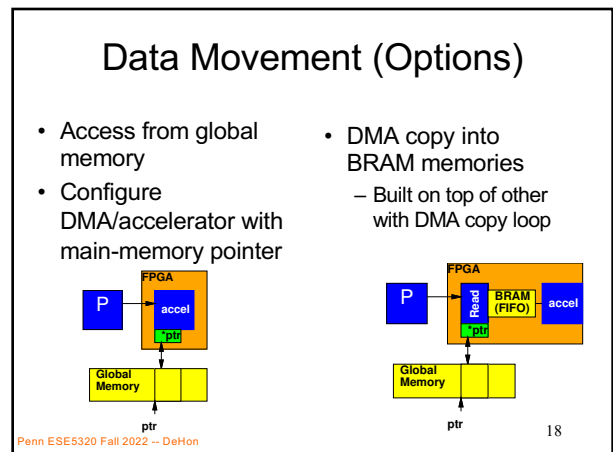
15



16



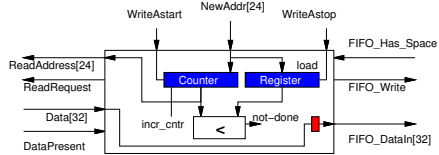
17



18

## Remember: Preclass 4a

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



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

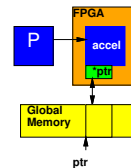
Penn ESE5320 Fall 2022 -- DeHon

19

19

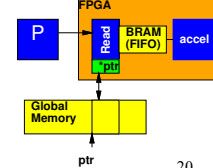
## Data Movement (Options)

- How copy loop improve performance?
  - Assume dataflow streaming between copy loop and accelerator computation?



Penn ESE5320 Fall 2022 -- DeHon

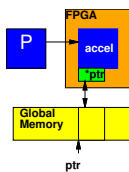
20



20

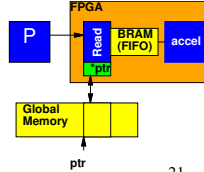
## Copy Loop Benefits

- Run concurrently – improve throughput
  - Overlap compute and communication
- Hide latency of read



Penn ESE5320 Fall 2022 -- DeHon

21



21

## Host-Side Programming

Part 2

Penn ESE5320 Fall 2022 -- DeHon

22

## Outline

- Platform
- Device
- Context
- Command Queue
- Load FPGA
- Buffers and Transfer
- Kernel Execution
- Synchronize on Result
- Overlap Data and Compute
- Running Multiple Kernels

Penn ESE5320 Fall 2022 -- DeHon

23

23

## Platform

```
cl_platform_id platform_id; // platform id
err = clGetPlatformIDs(10, platforms, &platform_count);
// Find Xilinx Platform
for (unsigned int iplat=0; iplat<platform_count; iplat++) {
    err = clGetPlatformInfo(platforms[iplat],
        CL_PLATFORM_VENDOR,
        1000,
        (void *)&cl_platform_vendor,
        NULL);
    if (strcmp(cl_platform_vendor, "Xilinx") == 0) {
        // Xilinx Platform found
        platform_id = platforms[iplat];
    }
}
```

- For general case where have multiple OpenCL platforms
- Silly for Ultra96 SoC
  - and, looks like we can partially hide

Penn ESE5320 Fall 2022 -- DeHon

24

24

## Device

```
cl_device_id devices[16]; // compute device id
char cl_device_name[100];

err = clGetDeviceIDs(platform_id, CL_DEVICE_TYPE_ACCELERATOR,
    16, devices, &num_devices);

printf("INFO: Found %d devices\n", num_devices);

//iterate all devices to select the target device.
for (uint i=0; i<num_devices; i++) {
    err = clGetDeviceInfo(devices[i], CL_DEVICE_NAME, 1024, cl_device_name,
    0);
    printf("CL_DEVICE_NAME %s\n", cl_device_name);
}
```

- For case of multiple FPGAs → which FPGA
- Also somewhat silly for Ultra96 with single FPGA (PL)
- Amazon F1.x16large has multiple FPGAs

Penn ESE5320 Fall 2022 -- DeHon

25

25

## Device

- Our Utilities.cpp,h – hide Platform  
`std::vector<cl::Device> devices = get_xilinx_devices();`  
`devices.resize(1);`  
`cl::Device device = devices[0];`
- Get away with simplification
  - Know want only Xilinx platform
  - Know only one device (Ultra96)

Penn ESE5320 Fall 2022 -- DeHon

26

26

## Context

```
context = clCreateContext(0, 1, &device_id, NULL, NULL, &err);
```

- Context for device
  - Not sure what flexibility/differentiation this is giving
    - maybe for DFX overlay shells
  - Supports some error callback

Penn ESE5320 Fall 2022 -- DeHon

27

27

## Command Queue

```
// Out-of-order Command queue
commands = clCreateCommandQueue(context, device_id,
    CL_QUEUE_OUT_OF_ORDER_EXEC_MODE_ENABLE, &err);

// In-order Command Queue
commands = clCreateCommandQueue(context, device_id, 0, &err);
```

- Setup queue to issue commands
- For “programming” interactions with accelerators
- Use to set them up to happen as separate thread
  - Separate when host requests and when accelerator actually run

Penn ESE5320 Fall 2022 -- DeHon

28

28

## Load Bitstream on FPGA

```
int size=load_file_to_memory(xclbin, (char **) &kernelbinary);
size_t size_var = size;

cl_program program = clCreateProgramWithBinary(context, 1, &device_id,
    &size_var, (const unsigned char **) &kernelbinary,
    &status, &err);
```

- Where we actually reconfigure FPGA
  - Not really for our current platform
    - But confirms bitstream using
    - ...and needs read some information from xclbin
  - For Partial Reconfiguration (DFX) platform
    - Option to choose which to load
    - Option to reconfigure during operation under host control

Penn ESE5320 Fall 2022 -- DeHon

29

29

## Kernels

```
kernel1 = clCreateKernel(program, "<kernel_name_1>", &err);
kernel2 = clCreateKernel(program, "<kernel_name_2>", &err);
```

- Setup the actual accelerators to call
  - Extract handle to reference them from the loaded bitstream (program)
  - Just naming them at this point

Penn ESE5320 Fall 2022 -- DeHon

30

30

## Preclass 1

- Kernel with
  - 10ms input transfer
  - 30ms compute
  - 5ms outputs
- Latency?
- Throughput if must serialize?
- Throughput if overlap data transfers and computation?

Penn ESE5320 Fall 2022 -- DeHon 31

31

## Preclass 1

Penn ESE5320 Fall 2022 -- DeHon 32

32

## Host and Device Memory View

- Host pointers not necessarily meaningful to device
  - Need to map
  - Different memories? (depend on platform)
    - Embedded incl. Ultra96 (same), data-center (may be different)
  - Virtual vs. Physical addresses?
  - Same memory, different address map?

Penn ESE5320 Fall 2022 -- DeHon 33

33

## Buffers and Transfer

```

// Two cl_mem buffer, for read and write by kernel
cl_mem dev_mem_read_ptr = clCreateBuffer(context,
    CL_MEM_READ_ONLY,
    sizeof(int) * number_of_words, NULL, NULL);

cl_mem dev_mem_write_ptr = clCreateBuffer(context,
    CL_MEM_WRITE_ONLY,
    sizeof(int) * number_of_words, NULL, NULL);

// Setting arguments
clSetKernelArg(kernel, 0, sizeof(cl_mem), &dev_mem_read_ptr);
clSetKernelArg(kernel, 1, sizeof(cl_mem), &dev_mem_write_ptr);

// Get host side pointer of the cl_mem buffer object
auto host_write_ptr =
    clEnqueueMapBuffer(queue, dev_mem_read_ptr, true, CL_MAP_WRITE, 0, nullptr,
        r, nullptr, &err);
auto host_read_ptr =
    clEnqueueMapBuffer(queue, dev_mem_write_ptr, true, CL_MAP_READ, 0, nullptr,
        r, nullptr, &err);

// Fill up the host_write_ptr to send the data to the FPGA
for(int i=0; i< MAX; i++) {
    host_write_ptr[i] = <... >
}

// Migrate
cl_mem mems[2] = {host_write_ptr, host_read_ptr};
clEnqueueMigrateMemObjects(queue, 2, mems, 0, 0, nullptr, &migrate_event);

// Schedule the kernel
clEnqueueTask(queue, kernel, 1, &migrate_event, &enqueue_event);

// Migrate data back to host
clEnqueueMigrateMemObjects(queue, 1, &dev_mem_write_ptr,
    &data_read_event);
    
```

Penn ESE5320 Fall 2022 -- DeHon 34

34

## Remember: Preclass 4a

```

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

```

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

Penn ESE5320 Fall 2022 -- DeHon 35

35

## Movement Options

- clEnqueueReadBuffer (WriteBuffer)
  - Not guarantee timing of transfer
- clEnqueueMigrateObjects
  - Supports overlap of compute and communicate
  - recommended

Penn ESE5320 Fall 2022 -- DeHon 36

36

## Remap Device Pointer

- Migration remap dev\_mem\_ptr
- (don't know if this is happening this way)

Penn ESE5320 Fall 2022 -- DeHon 37

37

## Kernel Invocation

```
err = clEnqueueTask(commands, kernel, 0, NULL, NULL);
```

- After
  - Kernel arguments set
  - Buffer transfers have been enqueued
- Causes to actual run on data transferred
- Commands – command queue
- Kernel – naming of kernel to run
- Last two arguments for synchronization <sup>38</sup>

Penn ESE5320 Fall 2022 -- DeHon

38

## Synchronize

```
q.enqueueMigrateMemObjects({in1_buf, in2_buf},
    0, NULL, &event_sp);
clWaitForEvents(1, (const cl_event *)&event_sp);
q.enqueueTask(krn1_mmult, NULL, &exec_event_sp);
clWaitForEvents(1,
    (const cl_event *)&exec_event_sp);
```

- Options to synchronize in calls
- Must synchronize to know finished

Penn ESE5320 Fall 2022 -- DeHon 39

39

## Synchronize

```
q.enqueueMigrateMemObjects({in1_buf, in2_buf},
    0, NULL, &event_sp);
clWaitForEvents(1, (const cl_event *)&event_sp);
q.enqueueTask(krn1_mmult, NULL, &exec_event_sp);
clWaitForEvents(1,
    (const cl_event *)&exec_event_sp);
```

- But, don't always want to synchronize after every operation

Penn ESE5320 Fall 2022 -- DeHon 40

40

## Preclass 2a

- Break kernel into three 10ms pieces and use dataflow between them
  - 10ms input transfer
  - 10ms compute x 3
  - 5ms output transfer
- Throughput?

Penn ESE5320 Fall 2022 -- DeHon 41

41

## Preclass 1 and 2

Penn ESE5320 Fall 2022 -- DeHon 42

42

## Dataflow: Host-to-kernel

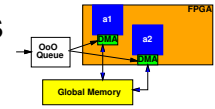
- Create a dataflow stream between host process and accelerator
  - Like DATAFLOW pragma
  - Control on FPGA/kernel side (coming up)
- Overlap next computation (setup) on host with FPGA kernel
- UG1393 says embedded platforms not support?
  - ...but maybe not need since single global memory

Penn ESE5320 Fall 2022 -- DeHon

43

43

## Running Multiple Kernel Instances



- Can have multiple instances of same kernel on FPGA
- Dispatch from Out-of-Order command queue
- Will distribute among identical kernels
- Out-of-Order queue
  - Helps with compute/communication overlap for single kernel. (see HW6)
  - Not force synchronization on return result.

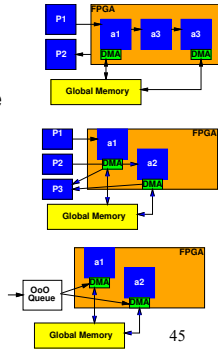
Penn ESE5320 Fall 2022 -- DeHon

44

44

## Parallelism

- Pipeline/Dataflow
  - With Host Dataflow
  - DMA transfers and compute
  - Within compute
- Thread Parallelism
  - Launch multiple accelerator
- Data Parallelism
  - Define multiple instances of one accelerator and enqueue work for



Penn ESE5320 Fall 2022 -- DeHon

45

45

## Cleanup

```

clReleaseCommandQueue(Command_Queue);
clReleaseContext(Context);
clReleaseDevice(Target_Device_ID);
clReleaseKernel(Kernel);
clReleaseProgram(Program);
free(Platform_IDs);
free(Device_IDs);
    
```

- Need to release/free resources setup

Penn ESE5320 Fall 2022 -- DeHon

46

46

## FPGA-Side Programming

### Part 3

Penn ESE5320 Fall 2022 -- DeHon

47

47

## Interface Pragmas

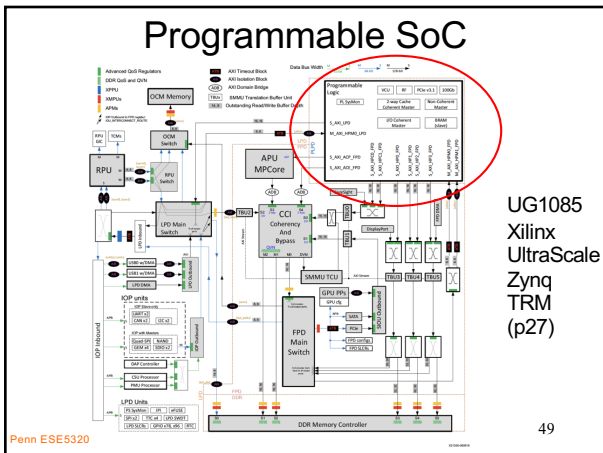
- Assign AXI
  - Bundles
- Host-to-Kernel
- Burst read
- Interface width

Penn ESE5320 Fall 2022 -- DeHon

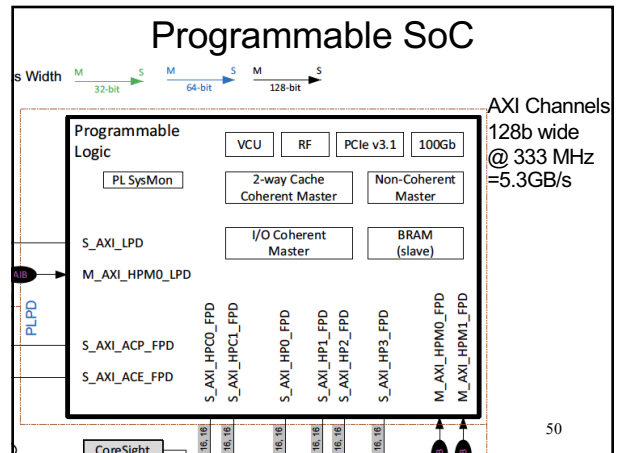
48

48





49



50

### AXI Channels

```
void cnn( int *pixel, // Input pixel
         int *weights, // Input Weight Matrix
         int *out, // Output pixel
         ... // Other input or Output ports

#pragma HLS INTERFACE m_axi port=pixel offset=slave bundle=gmem
#pragma HLS INTERFACE m_axi port=weights offset=slave bundle=gmem
#pragma HLS INTERFACE m_axi port=out offset=slave bundle=gmem
```

- Can tell it to use AXI Channels
- By default (as shown above) all share same channel

Penn ESE5320 Fall 2022 -- DeHon 51

51

### AXI Channels

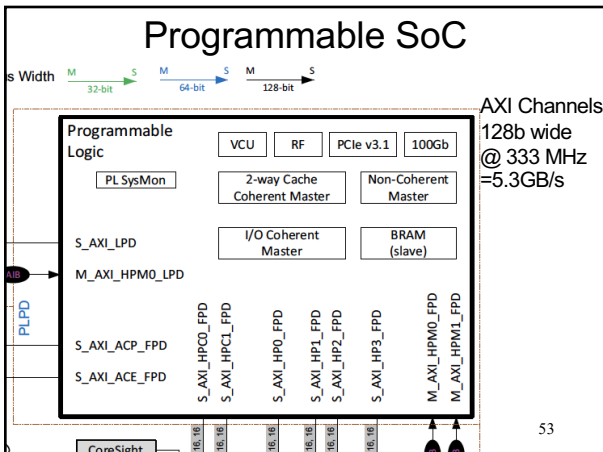
```
void cnn( int *pixel, // Input pixel
         int *weights, // Input Weight Matrix
         int *out, // Output pixel
         ... // Other input or Output ports

#pragma HLS INTERFACE m_axi port=pixel offset=slave bundle=gmem
#pragma HLS INTERFACE m_axi port=weights offset=slave bundle=gmem1
#pragma HLS INTERFACE m_axi port=out offset=slave bundle=gmem
```

- Separate names to use different
- Note, not naming to match ZCU3EG names  
– Portability across targets

Penn ESE5320 Fall 2022 -- DeHon 52

52



53

### Older and Defaults

```
• 2019.2:
#pragma HLS INTERFACE s_axilite port = in1 bundle = control
#pragma HLS INTERFACE s_axilite port = in2 bundle = control
#pragma HLS INTERFACE s_axilite port = out_r bundle = control
#pragma HLS INTERFACE s_axilite port = size bundle = control
#pragma HLS INTERFACE s_axilite port = return bundle = control
#pragma HLS INTERFACE m_axi port=in1 offset=slave bundle=gmem
#pragma HLS INTERFACE m_axi port=in2 offset=slave bundle=gmem
#pragma HLS INTERFACE m_axi port=out_r offset=slave bundle=gmem
• 2020.1:
#pragma HLS INTERFACE m_axi port=in1 bundle=gmem
#pragma HLS INTERFACE m_axi port=in2 bundle=gmem
#pragma HLS INTERFACE m_axi port=out_r bundle=gmem
```

Penn ESE5320 Fall 2022 -- DeHon 54

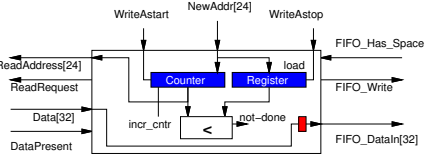
54

## Remember: Preclass 4a

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

### S-AXI Control – load pointers

M-AXI transfer  
Using pointers



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

Penn ESE5320 Fall 2022 -- DeHon

55

55

## Host-to-kernel Dataflow

```
void kernel_name( int *inputs,
... ) // Other input or Output ports
{
#pragma HLS INTERFACE ... // Other interface pragmas
#pragma HLS INTERFACE ap_ctrl_chain port=return bundle=control
```

- ap\_ctrl\_chain

Penn ESE5320 Fall 2022 -- DeHon

56

56

## Burst Transfer

```
hls::stream<datatype_t> str;
INPUT_READ: for(int i=0; i<INPUT_SIZE; i++) {
#pragma HLS PIPELINE
str.write(inp[i]); // Reading from Input interface
}
```

- Where have we seen this before?

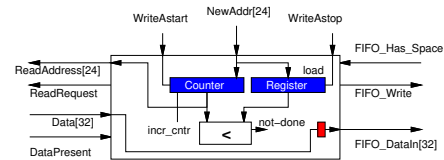
Penn ESE5320 Fall 2022 -- DeHon

57

57

## Remember: Preclass 4a

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



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

Penn ESE5320 Fall 2022 -- DeHon

58

58

## Burst Transfer

```
top_function(datatype_t * m_in, // Memory data Input
datatype_t * m_out, // Memory data Output
int inp1, // Other Input
int inp2) { // Other Input
#pragma HLS DATAFLOW
hls::stream<datatype_t> in_var1; // Internal stream to transfer
hls::stream<datatype_t> out_var1; // data through the dataflow region
read_function(m_in, inp1, in_var1); // Read function contains pipelined for
loop // to infer burst
execute_function(in_var1, out_var1, inp1, inp2); // Core compute function
write_function(out_var1, m_out); // Write function contains pipelined for
loop // to infer burst
}
```

- Use Dataflow to setup as separate threads on accelerator side
- May also be able to use C memcpy
  - Vitis UG1399 (2020.1, p. 212)

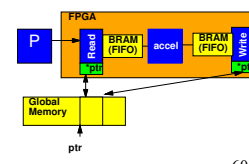
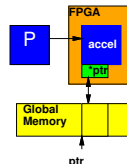
Penn ESE5320 Fall 2022 -- DeHon

59

59

## Data Movement (Options)

- Access from global memory
- Configure DMA/accelerator with main-memory pointer
- DMA burst copy into BRAM memories
  - Built on top of other with DMA copy loop



Penn ESE5320 Fall 2022 -- DeHon

60

60

## Preclass 3, 4

- 333MHz cycle time on bus
- Peak bandwidth if transfer per cycle
  - 32b
  - 64b
  - 128b
- Memory at 4 GB/s
  - Which widths make bus the bottleneck?
  - Memory?

Penn ESE5320 Fall 2022 -- DeHon

61

61

## Width

```
void cnn( ap_uint<512> *pixel, // Input pixel
         int *weights, // Input Weight Matrix
         ap_uint<512> *out, // Output pixel
         ... // Other input or output ports

#pragma HLS INTERFACE m_axi port=pixel offset=slave bundle=gmem
#pragma HLS INTERFACE m_axi port=weights offset=slave bundle=gmem
#pragma HLS INTERFACE m_axi port=out offset=slave bundle=gmem
```

- In some cases will auto size data transfer
  - See Vitis UG 1399 (2020.1) p. 214  
“Automatic Port Width Resizing”
- Others transfer at width of data type give
  - May need to pack and unpack data to exploit full width

Penn ESE5320 Fall 2022 -- DeHon

62

62

## Big Ideas

- Data between PS and PL managed as DMA
- Can exploit familiar parallelism
  - Thread, Data, Streaming Dataflow
  - Overlap compute and communicate
- Must manage some pieces pretty explicitly

Penn ESE5320 Fall 2022 -- DeHon

63

63

## Admin

- Feedback
- Reading for Wednesday on web
- HW6
  - Due on Friday

Penn ESE5320 Fall 2022 -- DeHon

64

64