

More code PS

17: Pipelining II

ENGR 315: Hardware/Software CoDesign

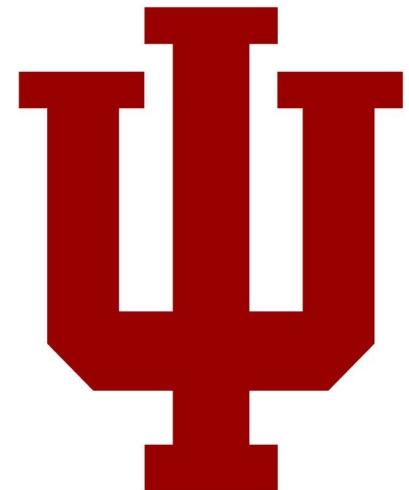
Andrew Lukefahr

Indiana University

Some material taken from:

https://github.com/trekhleb/homemade-machine-learning/tree/master/homemade/neural_network

<http://cs231n.github.io/neural-networks-1/>



Announcements

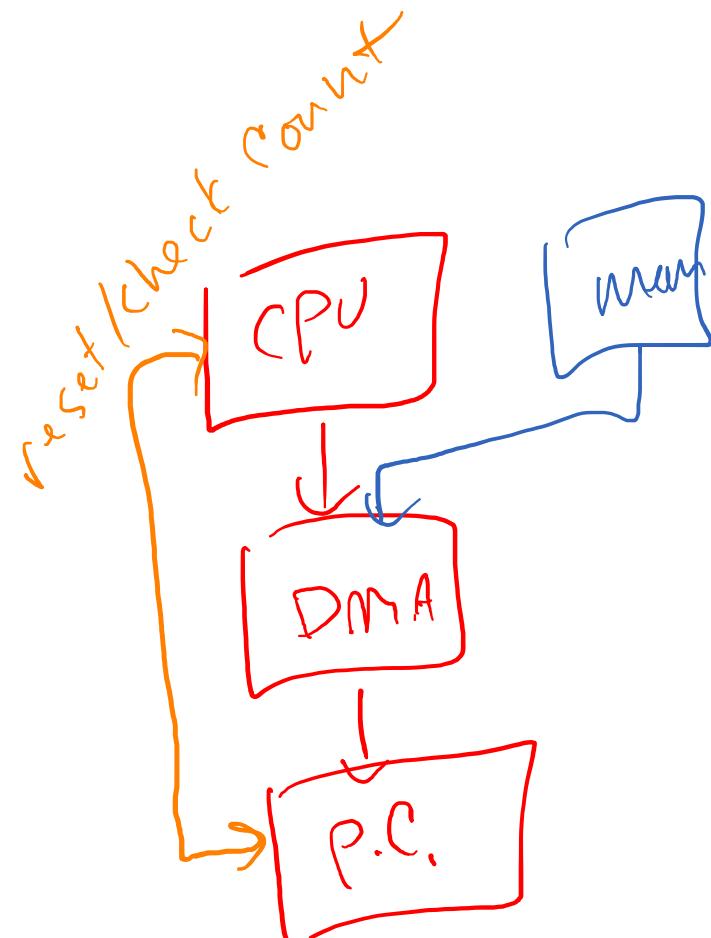
- P6 out. Due Wednesday.

- No Class on Wednesday (Oct 26th)

Review - Prior Wed Nov. 1st
TA session

- Exam Nov 6th *- Mon*

P6: Adds DMA + AXI-Stream to Popcount



- DMA
 - Add DMA engine to move data via AXI4-Full to AXI-Stream interface
- Popcount.sv:
 - Add AXI-Stream Interface
 - Keep AXI4-Lite Interface to read result

P7 – DMA from C

A screenshot of a search results page. The search bar at the top contains the query "xilinx dma ip". Below the search bar are navigation links: All (highlighted), News, Shopping, Images, Videos, More, and Tools. The text "About 294,000 results (0.50 seconds)" is displayed. The first result is a link to the "AXI DMA v7.1 LogiCORE IP Product Guide - Xilinx" document, which is a PDF file from Jun 14, 2019.

xilinx dma ip

All News Shopping Images Videos More Tools

About 294,000 results (0.50 seconds)

[https://www.xilinx.com › support › pg021_axi_dma](https://www.xilinx.com/support/documentation/ip_documentation/axi_dma/v7_1/pg021_axi_dma.pdf) PDF

AXI DMA v7.1 LogiCORE IP Product Guide - Xilinx

Jun 14, 2019 — The AXI Direct Memory Access (AXI DMA) IP core provides high-bandwidth direct memory access between the AXI4 memory mapped and AXI4-Stream IP ...
97 pages

https://www.xilinx.com/support/documentation/ip_documentation/axi_dma/v7_1/pg021_axi_dma.pdf

P7 – DMA from C

Programming Sequence

Direct Register Mode (Simple DMA)

P7 – DMA from C

1. Start the MM2S channel running by setting the run/stop bit to 1 (MM2S_DMACR.RS = 1). The halted bit (DMASR.Halted) should deassert indicating the MM2S channel is running.
2. Skip
3. Write a valid source address to the MM2S_SA register.
4. Write the number of bytes to transfer in the MM2S_LENGTH register.
The MM2S_LENGTH register must be written last.
5. Wait until MM2S_DMASR.Idle==1 for completion

P8+ Accelerate Machine Learning

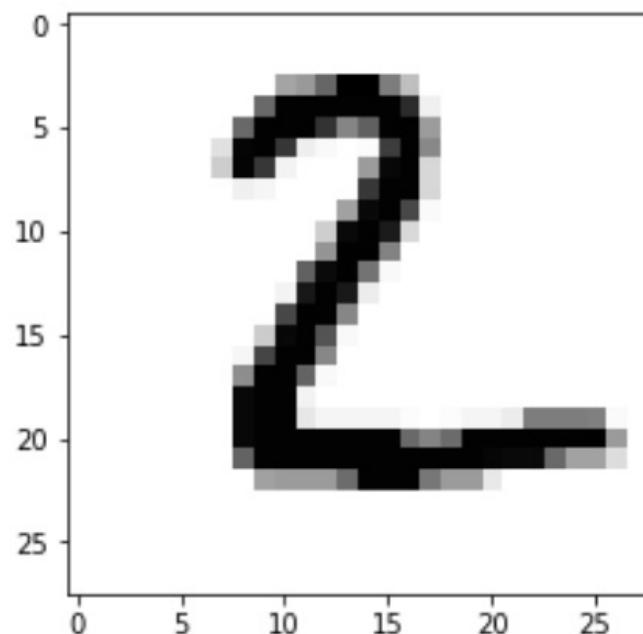
- Goal: Accelerate reference neural network
- Harder, more open-ended projects

Simple Neural Network

=====

Index: 0

Image:



- Takes in image of number
- Returns integer value
- How? artificial neural network

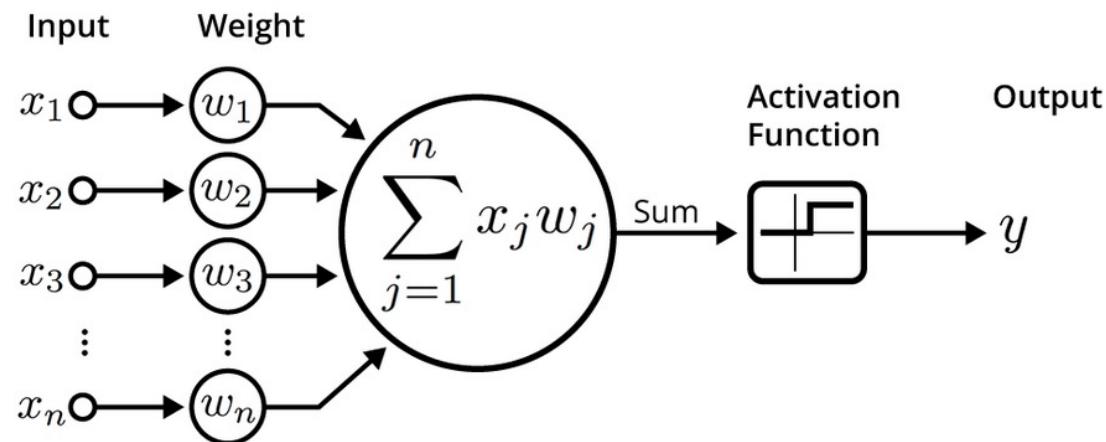
→ ML Classification Result: 2

Real Value: 2

Correct Result: True

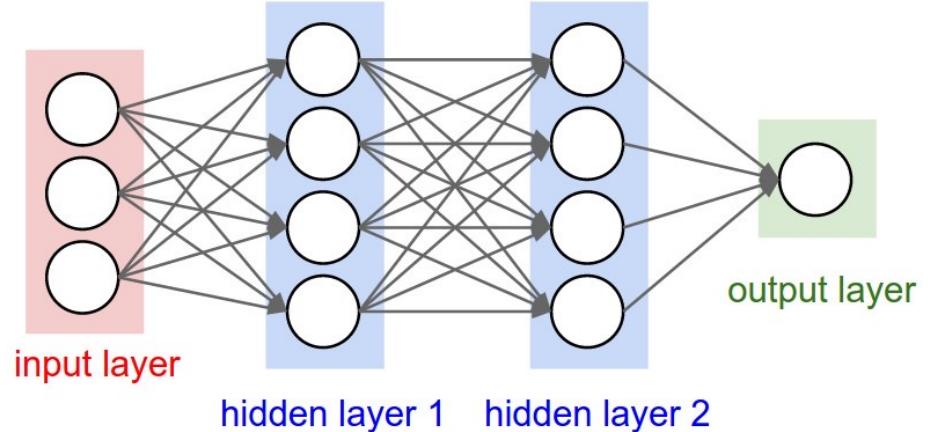
=====

Python Neuron

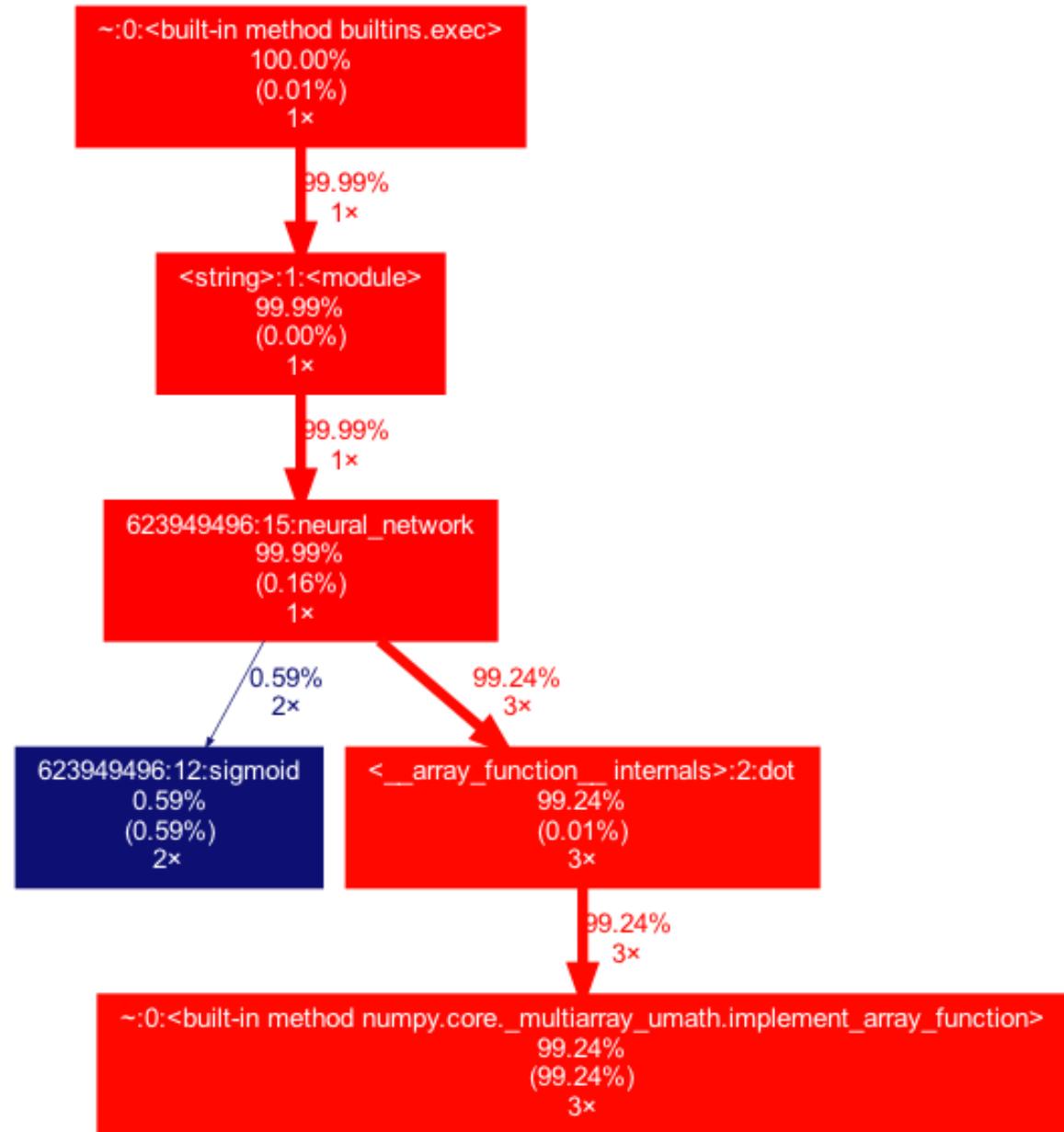


```
class Neuron(object):
    ...
    def forward(self, inputs):
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """
        cell_body_sum = np.sum(inputs * self.weights) + self.bias
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation function
        return firing_rate
```

Why focus on Dot Product?



```
# forward-pass of a 3-layer neural network:  
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)  
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)  
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)  
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)  
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```



Matrix Multiplication (Dot Product)

$$\begin{bmatrix} i_0 & i_1 \end{bmatrix} \times \begin{bmatrix} w_{00} & w_{10} & w_{20} \\ w_{01} & w_{11} & w_{21} \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} o_0 & o_1 & o_2 \end{bmatrix}$$

$$o_0 = i_0 \cdot w_{00} + i_1 \cdot w_{01}$$

$$o_1 = i_0 \cdot w_{10} + i_1 \cdot w_{11}$$

$$o_2 = i_0 \cdot w_{20} + i_1 \cdot w_{21}$$

Matrix Multiplication (Dot Product)

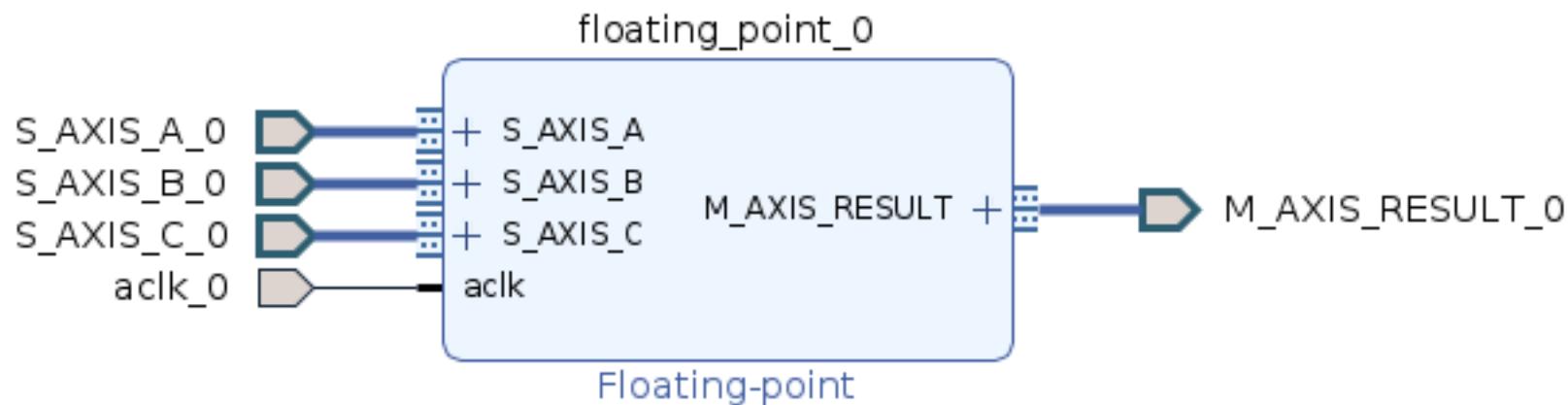
inputs

$$\begin{bmatrix} 0.1 \\ \underline{0.2} \end{bmatrix} \times \begin{bmatrix} \text{weights} \\ 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} =$$

Floating-Point Multiply-Accumulate (FMAC)

- Math: $a * b + c$

Floating-Point Multiply in Hardware

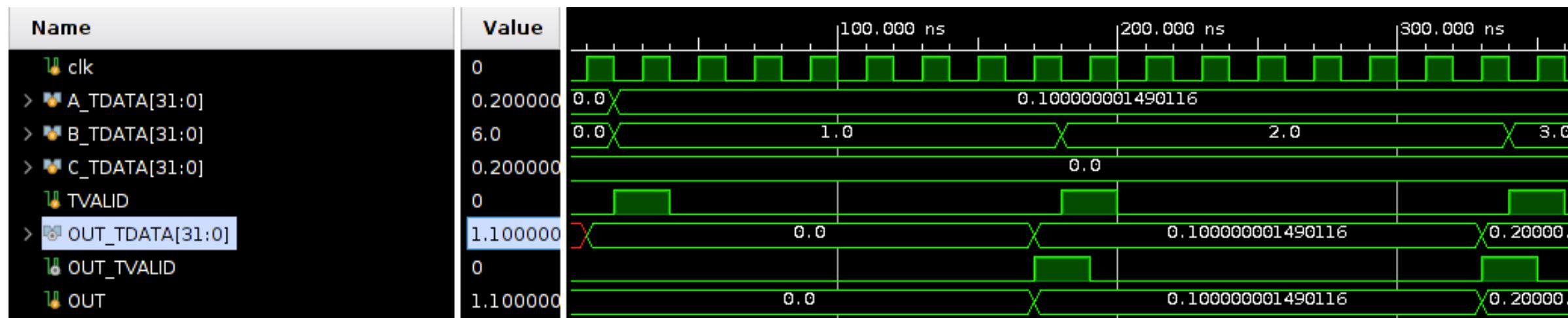


- $\text{result} = a * b + c$

Floating-Point math takes 8 cycles.

- Floating-Point is complicated.
- 8 cycles of complicated.

Demo Time



Floating-Point math takes 8 cycles.

- Floating-Point is complicated.
- 8 cycles of complicated.
- How do we work around an 8 cycle latency?
- Pipelining!

How does Pipelining work in Math?

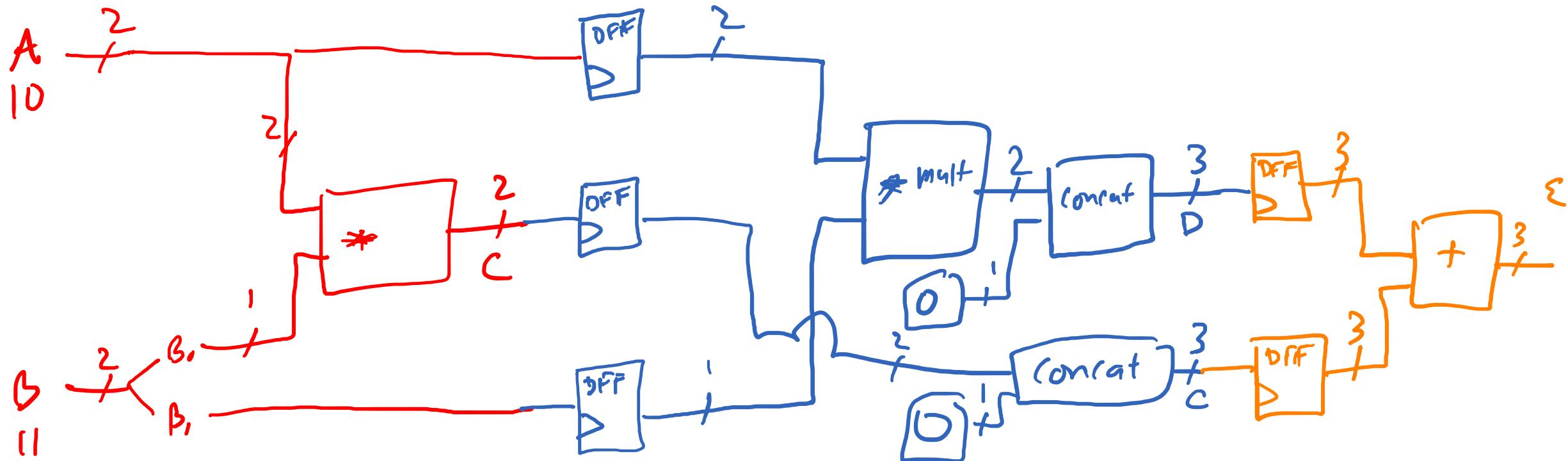
$$\begin{array}{r} 2 \\ \times 3 \\ \hline 6 \end{array} \Rightarrow \begin{array}{r} 10 \\ \times 11 \\ \hline 110 \end{array} \Rightarrow \begin{array}{r} 10 \\ \times 1 \\ \hline 10 \end{array} + \begin{array}{r} 10 \\ \times 10 \\ \hline 100 \end{array} = 110$$

$$\begin{array}{r} 20 \\ \times 31 \\ \hline 20 \\ + 600 \\ \hline \end{array}$$

How does pipelining work in circuits?

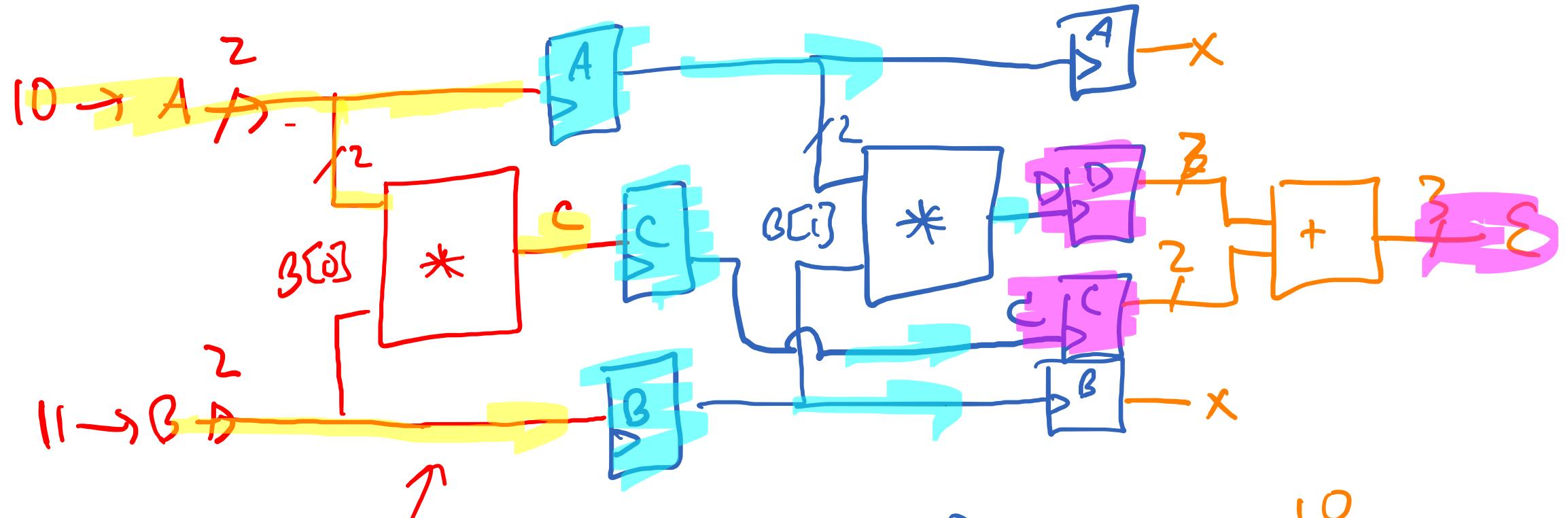
$$\begin{array}{r} 2 \\ \times 3 \\ \hline 6 \end{array} \quad \begin{array}{r} 10 \\ \times 1 \\ \hline 10 \end{array} \quad \begin{array}{r} 10 \\ + 10 \\ \hline 110 \end{array}$$

Handwritten annotations: A red circle highlights the top row of the first multiplication step. A red circle highlights the result of the first multiplication (110). A red circle highlights the carry bit from the first multiplication (10). A blue circle highlights the result of the second multiplication (100). An orange circle highlights the final sum (110).



Pipelining in hardware

$$\begin{array}{r}
 \begin{array}{r}
 A \\
 b \\
 \times 3 \\
 \hline
 6
 \end{array}
 \Rightarrow
 \begin{array}{r}
 10 \\
 \times 11 \\
 \hline
 110
 \end{array}
 \Rightarrow
 \begin{array}{r}
 10 \\
 \times 1 \\
 \hline
 c_{10}
 \end{array}
 + D = \underline{\underline{100}} = 110
 \end{array}$$

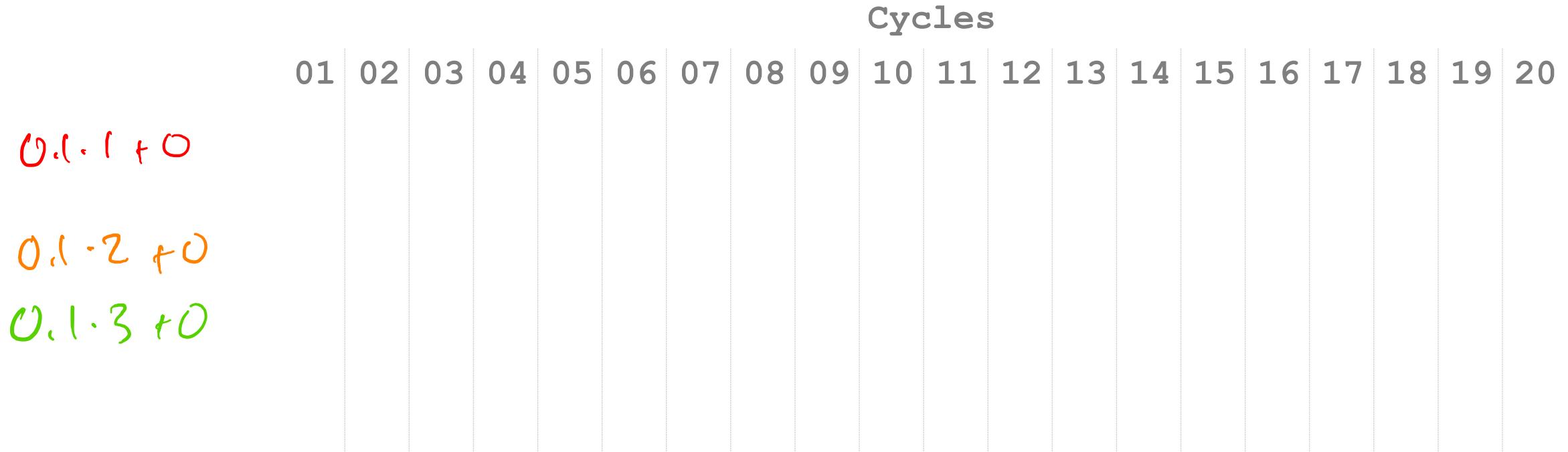
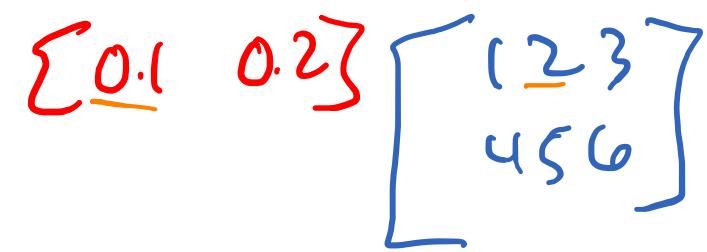


$$\begin{array}{r}
 (3) \quad \begin{array}{r}
 10 \\
 \times 1 \\
 \hline
 c = 10
 \end{array}
 \end{array}$$

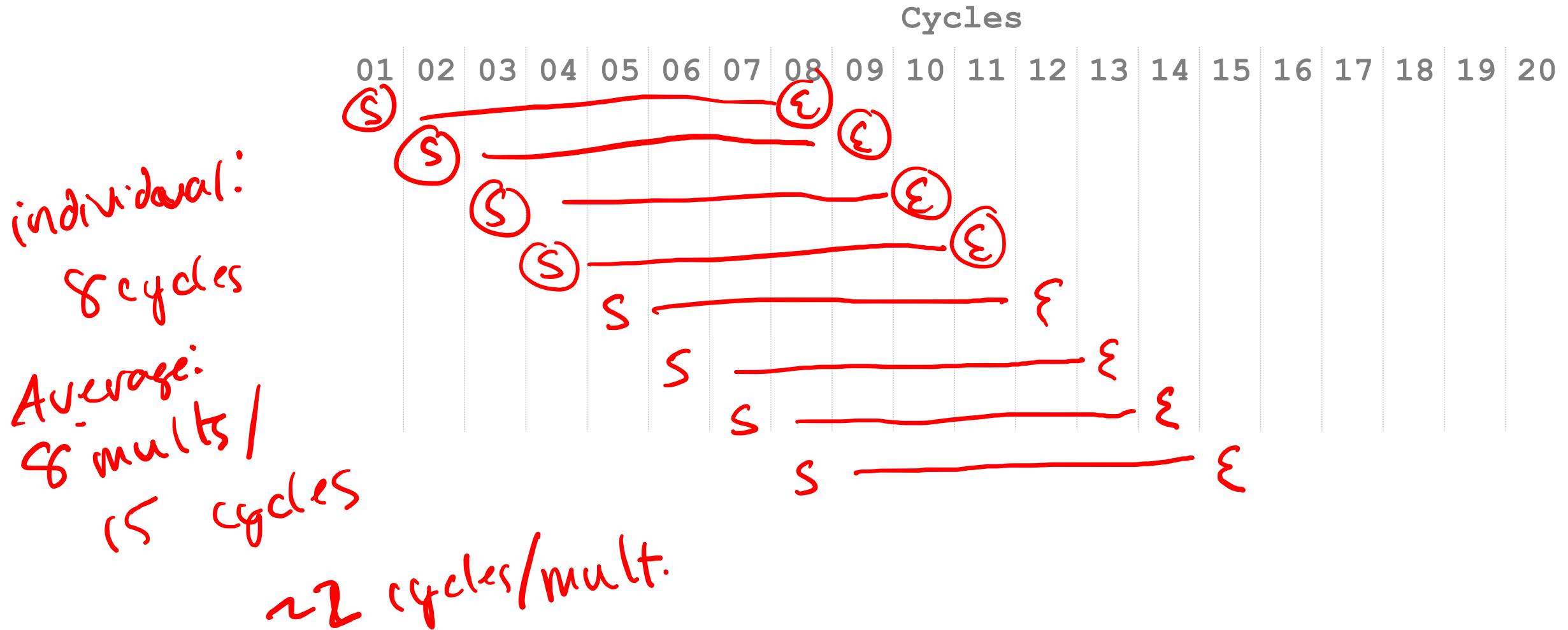
$$D = \frac{10}{100}$$

$$\begin{array}{r}
 10 \\
 + 100 \\
 \hline
 e = 110
 \end{array}$$

FMAC Pipelining



FMAC Pipelining



Latency vs. Throughput

- **Latency:** How long does an individual operation take to complete?

individual mult: 8 cycles

- **Throughput:** How many operations can you complete per second (or per cycle)?

Average output per cycles:

~ 1 mult / cycle

Pipelining

- FMAC takes 8 cycles for 1 value
 - But can accept a new value every cycle.
-
- What is Latency: 8
 - What is Throughput: ~1 / cycle

Recall: Matrix Multiplication (Dot Product)

inputs

$$\begin{bmatrix} 0.1 \\ \underline{0.2} \end{bmatrix} \times \begin{bmatrix} 1 & \overset{\text{weights}}{2} & 3 \\ 4 & 5 & 6 \end{bmatrix} =$$
$$= \begin{bmatrix} (0.1 \times 1 + 0.2 \times 4) & (0.1 \times 2 + 0.2 \times 5) & (0.1 \times 3 + 0.2 \times 6) \end{bmatrix}$$

(Answer)

$$= \begin{bmatrix} \underline{0.9} \\ 1.2 \\ \underline{1.5} \end{bmatrix}$$

Alternative Dot Computations

$$\begin{bmatrix} 0.1 & 0.2 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.2 & 1.5 \end{bmatrix}$$

$$\begin{bmatrix} 0.1 \cdot 1 & 0.1 \cdot 2 & 0.1 \cdot 3 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \cdot 4 & 0.2 \cdot 5 & 0.2 \cdot 6 \end{bmatrix} \stackrel{+}{=} \begin{bmatrix} 0.8 & 1.0 & 1.2 \end{bmatrix}$$
$$0.9 \quad 1.2 \quad 1.5$$

Multiply-Accumulate Dot Computations

$$\begin{bmatrix} 0.1 & 0.2 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.2 & 1.5 \\ 0 & 0 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 0.1 \cdot 1 + 0, 0.1 \cdot 2 + 0, 0.1 \cdot 3 + 0 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix}$$
$$\begin{bmatrix} 0.2 \cdot 4 + 0.1, 0.2 \cdot 5 + 0.2, 0.2 \cdot 6 + 0.3 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.2 & 1.5 \end{bmatrix}$$

Python Time

inputs

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \times \begin{array}{c} \text{Weights} \\ \left[\begin{array}{cccc} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{array} \right] \end{array} =$$

```
weights = np.array( [[1,2,3,4],[5,6,7,8],[9,10,11,12]], dtype=np.float32)
inputs = np.array([[0.1,0.2,0.3]], dtype=np.float32)
outputs = np.dot(inputs, weights)
```

Input	Weights	Output
[0.1 0.2 0.3]	[1. 2. 3. 4.] [5. 6. 7. 8.] [9. 10. 11. 12.]	= [3.8000002 4.4 5. 5.6000004]

Input [[0.1 0.2 0.3]] .	Weights [1. 2. 3. 4.] [5. 6. 7. 8.] [9. 10. 11. 12.]	= [3.8000002 4.4	5.	5.6000004]
----------------------------	--	------------------	----	------------

Mult-Accum Dot

```
# how its done in dot.sv
def pydot(inputs,weights):
    inputs = inputs[0] # remove outer nesting
    outs = np.zeros(weights.shape[1], dtype=np.float32)
    for i in range(weights.shape[0]): # input length
        for j in range(weights.shape[1]): # output length
            outs[j] = outs[j] + weights[i][j] * inputs[i]
    return outs
```

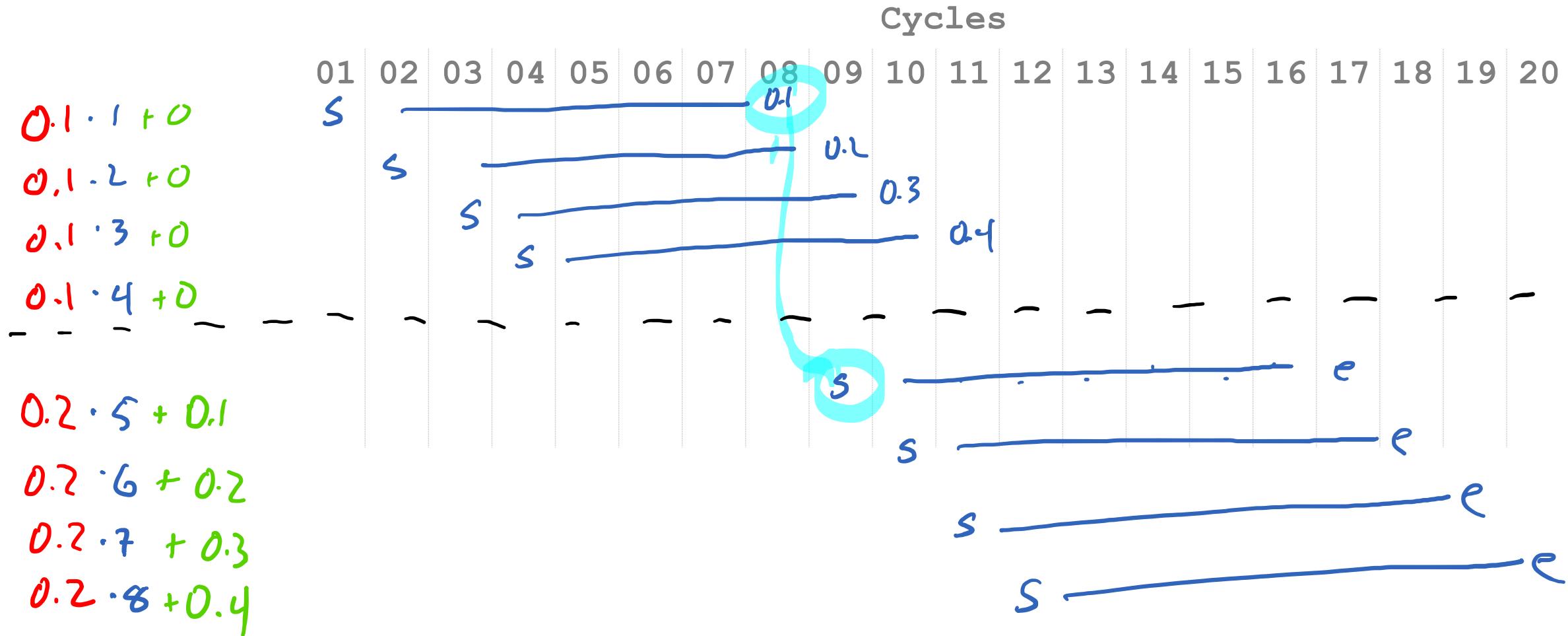
Inputs (Shape):
(1, 3)
Output (Shape):
(1, 4)
Weights (Shape):
(3, 4)

Input
[[0.1 0.2 0.3]]

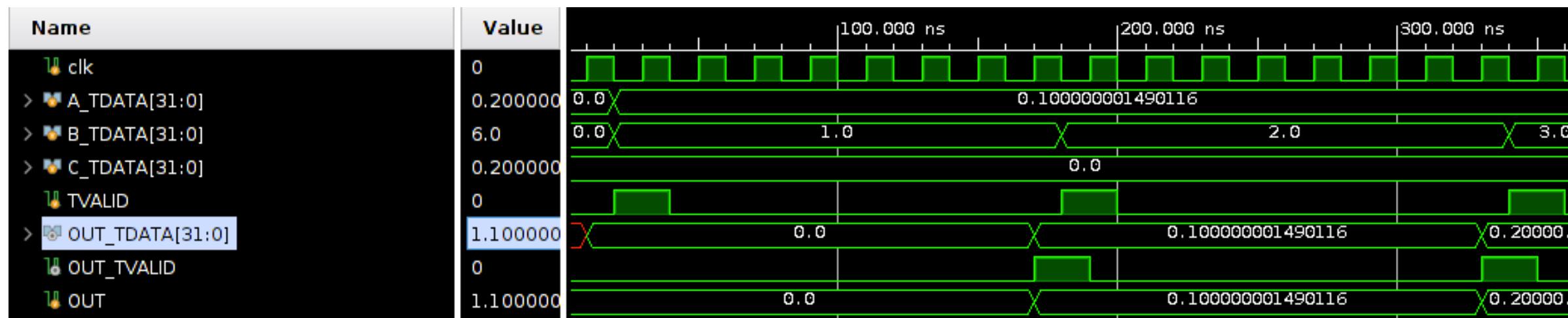
Weights
[1. 2. 3. 4.]
[5. 6. 7. 8.]
[9. 10. 11. 12.]

Output
[3.8000002 4.4
5. 5.6000004]

Dependencies



Demo Time



Next Time: More on
Dependencies

Latency on Pipelined FMAC

- Solution: Stall at the end of a row.
- Drain the pipeline.

Hardware Parallelism

- CPU: 1 Floating-Point Unit
- FPGA? *(10 Floating-Point Units?
20?
100?)*

Finding Parallelism

- Some computation that doesn't depend on other computation's results
- Shared Inputs are OK.

Next Time: Can we use 2+ FMACs?

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix} = \begin{bmatrix} 3.8 & 4.4 & 5 & 5.6 \end{bmatrix}$$

Option 1

Stopped here!

	<u>Cycle</u>	<u>fMAC</u>	<u>comp</u>
0	0	$0.1 \cdot 1 + 0$	
1	1	$0.1 \cdot 2 + 0$	
2		$0.1 \cdot 3 + 0$	
3		$0.1 \cdot 4 + 0$	

	<u>Cycle</u>	<u>fMAC</u>
0	0	$0.1 \cdot 1 + 0$
1	1	$0.1 \cdot 3 + 0$

fMAC

2

$0.1 \cdot 2 + 0$

$0.1 \cdot 3 + 0$

Option 2

$\rightarrow 0.1 * 1 + 0$
 $0.1 * 2 + 0$
 $0.1 * 3 + 0$
 $0.1 * 4 + 0$

$0.2 * 5 + 0$

$0.2 * 6 + 0$

$0.2 * 7 + 0$

$0.2 * 8 + 0$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix} = \begin{bmatrix} 3.8 & 4.4 & 5 & 5.6 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix}$$

$$0.1 \cdot \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \end{bmatrix}$$

$$0.2 \cdot \begin{bmatrix} 5 & 6 & 7 & 8 \end{bmatrix} + \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \end{bmatrix} = \begin{bmatrix} 1.1 & 1.4 & 1.7 & 2.0 \end{bmatrix}$$

$$0.3 \cdot \begin{bmatrix} 9 & 10 & 11 & 12 \end{bmatrix} + \begin{bmatrix} 1.1 & 1.4 & 1.7 & 2.0 \end{bmatrix} = \underbrace{\begin{bmatrix} 3.8 & 4.4 & 5 & 5.6 \end{bmatrix}}$$

Parallelize Alternative Dot Computations?

$$\begin{bmatrix} 0.1 & 0.2 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.2 & 1.5 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \leftarrow \text{result}$$

$$0.1 \cdot \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} =$$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \leftarrow \text{temp result}$$

$$0.2 \cdot \begin{bmatrix} 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} =$$

$$\begin{bmatrix} 0.8 & 1.0 & 1.2 \end{bmatrix} + \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.2 & 1.5 \end{bmatrix}$$

Can we parallelize Dot?

```
# how its done in dot.sv
def pydot(inputs,weights):
    inputs = inputs[0] # remove outer nesting
    outs = np.zeros(weights.shape[1], dtype=np.float32)
    for i in range(weights.shape[0]): # input length
        for j in range(weights.shape[1]): # output length
            outs[j] = outs[j] + weights[i][j] * inputs[i]
    return outs
```

Can we parallelize Dot?

```
# how its done in dot.sv
def pydot(inputs,weights):
    inputs = inputs[0] # remove outer nesting
    outs = np.zeros(weights.shape[1], dtype=np.float32)
    for i in range(weights.shape[0]): # input length
        for j in range(weights.shape[1]): # output length
            outs[j] = outs[j] + weights[i][j] * inputs[i]
    return outs
```

```
def par_pydot(inputs, weights):
    par_inputs = [inputs[:,::2], inputs[:,1::2]]
    par_weights = [weights[::2,:,:], weights[1::2,:,:]]

    par_outputs = [pydot(par_inputs[0], par_weights[0]),
                  pydot(par_inputs[1], par_weights[1])]

    outputs = par_outputs[0] + par_outputs[1]
    return outputs
```

Can we parallelize Dot?

```
# how its done in dot.sv
def pydot(inputs,weights):
    inputs = inputs[0] # remove outer nesting
    outs = np.zeros(weights.shape[1], dtype=np.float32)
    for i in range(weights.shape[0]): # input length
        for j in range(weights.shape[1]): # output length
            outs[j] = outs[j] + weights[i][j] * inputs[i]
    return outs
```

```
def par_pydot(inputs, weights):
    par_inputs = [inputs[:,::2], inputs[:,1::2]]
    par_weights = [weights[::2,:,:], weights[1::2,:,:]]

    → par_outputs = [pydot(par_inputs[0], par_weights[0]),
                     pydot(par_inputs[1], par_weights[1])]

    outputs = par_outputs[0] + par_outputs[1]
    return outputs
```

19: Hardware Acceleration III

Engr 315: Hardware / Software Codesign

Andrew Lukefahr

Indiana University

