

16: Hardware Acceleration III

Engr 315: Hardware / Software Codesign

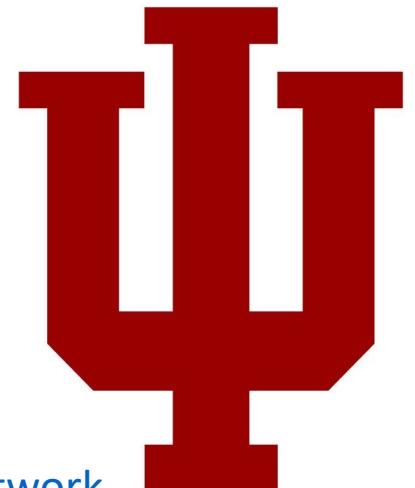
Andrew Lukefahr

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/>



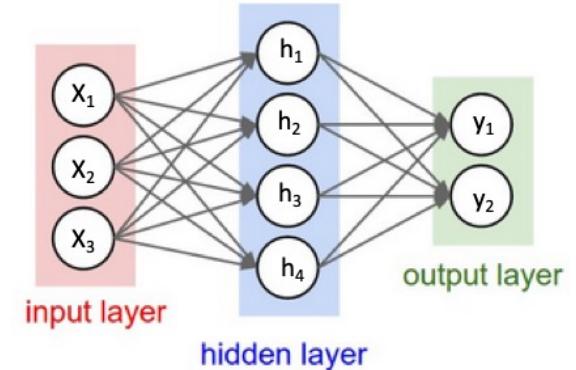
Announcements

- ~~P6~~ due Friday! DMA in C
- P7 is out
 - Maybe not AG yet?

Project 7: Accelerate Dot Product

- Given: Slow Dot-Product in Hardware
- Your Task: Accelerate it!
- How? Pipelining + Parallelism (later)
- Groups of 2 allowed.

Why Dot Product?



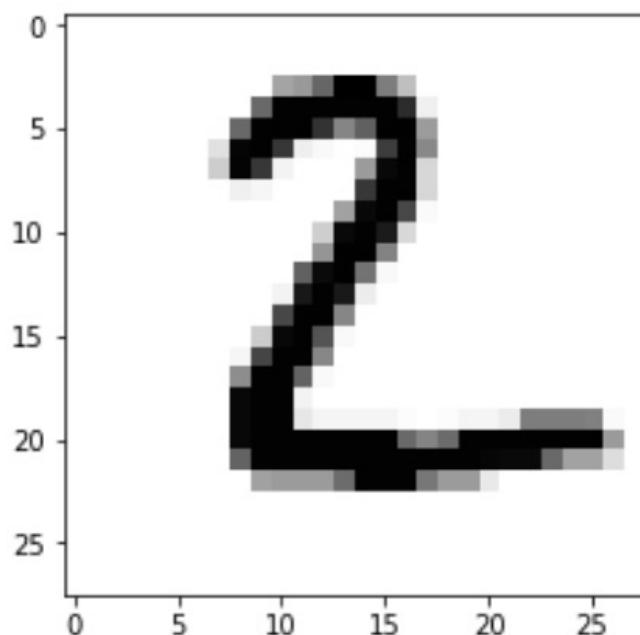
```
# forward-pass of a 3-layer neural network:  
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)  
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)  
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)  
h2 = f(np.dot(W2, h1) + b2)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

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

=====

Matrix Multiplication (Dot Product)

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

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

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

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

Alternative Dot Computations

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

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

$$0.1 \cdot \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} =$$
$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \leftarrow \text{temp result}$$

$$0.2 \cdot \begin{bmatrix} 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} =$$
$$\begin{bmatrix} 0.8 & 1.0 & 1.2 \end{bmatrix} + \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.2 & 1.5 \end{bmatrix}$$

Floating-Point math takes 8 cycles.

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

A · B + C
MAC

input:
 $(2) \begin{array}{r} 10 \\ 11 \end{array}$

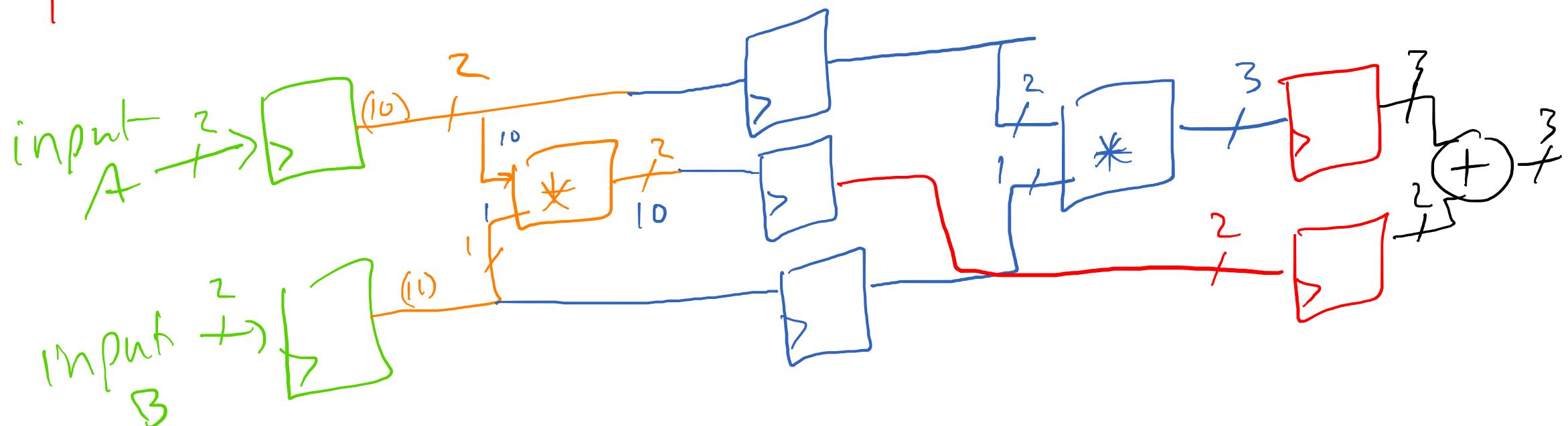
cycle 1
 $\begin{array}{r} 10 \\ * \quad 1 \\ \hline 10 \end{array}$

cycle 2
 $\begin{array}{r} 10 \\ * \quad 1 \quad 0 \\ \hline 10 \quad 0 \end{array}$

cycle 3

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

output
 $110 \quad (6)$



Pipelining

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

Latency vs. Throughput

- Latency: How long does an individual operation take?
- Throughput: How many operations can you complete in a given time?

Pipelining

Pipelined Dot Computations

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

Problems with Pipelined FMAC

Latency!

Latency on Pipelined FMAC

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

Parallelize Alternative Dot Computations?

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

$$\begin{bmatrix} 0.1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \quad | \text{ fMAC}$$

$$\begin{bmatrix} 0.2 \end{bmatrix} \begin{bmatrix} 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 0.8 & 1.0 & 1.2 \end{bmatrix} \quad | \text{ FMAC}$$

$$\begin{bmatrix} 0.9 & 1.2 & 1.5 \end{bmatrix} \quad | \text{ f, l ADD}$$

Parallelize Alternative Dot Computations?

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

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

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

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

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

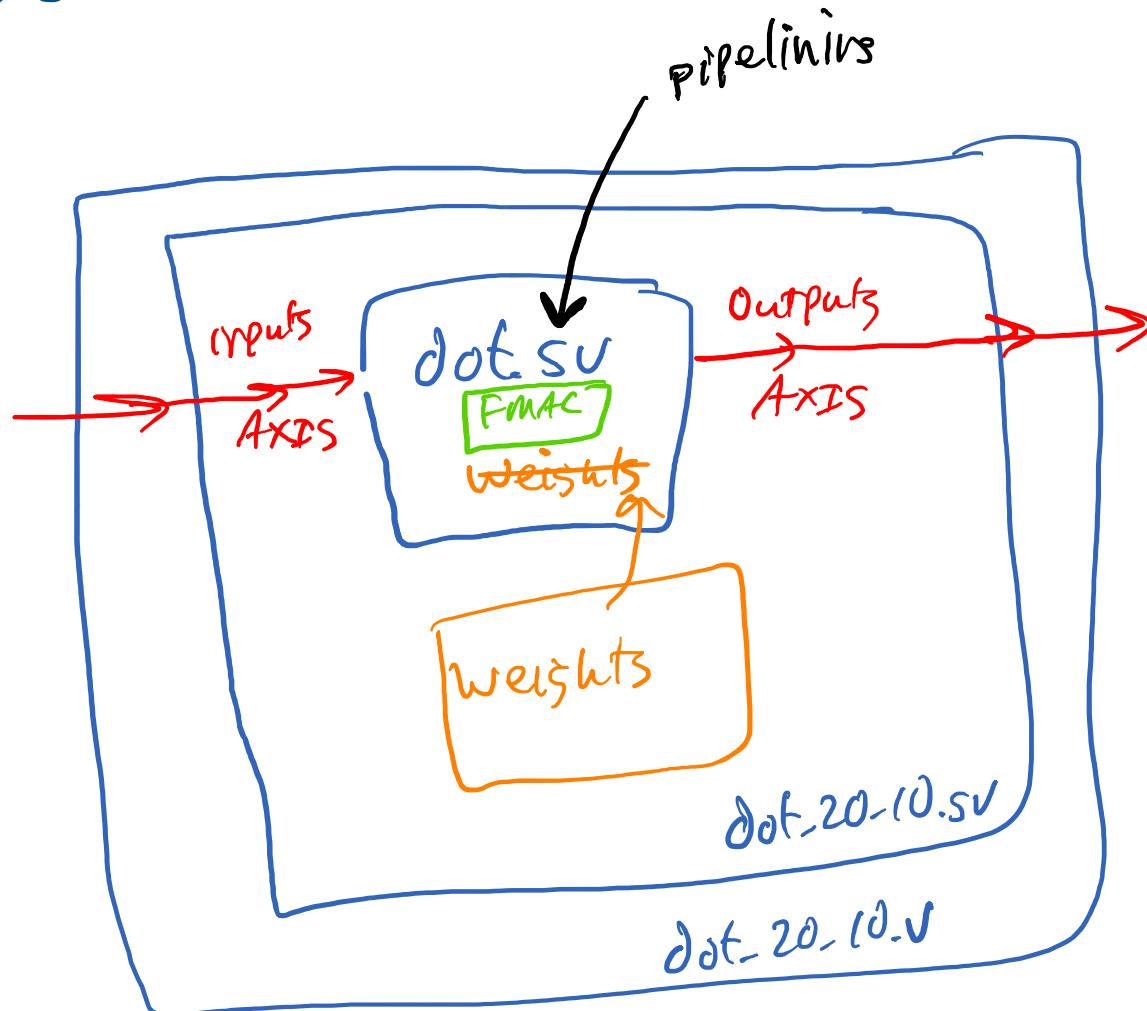
$$\begin{bmatrix} 0.8 & 1.0 & 1.2 \end{bmatrix} + \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.2 & 1.5 \end{bmatrix}$$

DT
PLo

Dot

Data flow

[

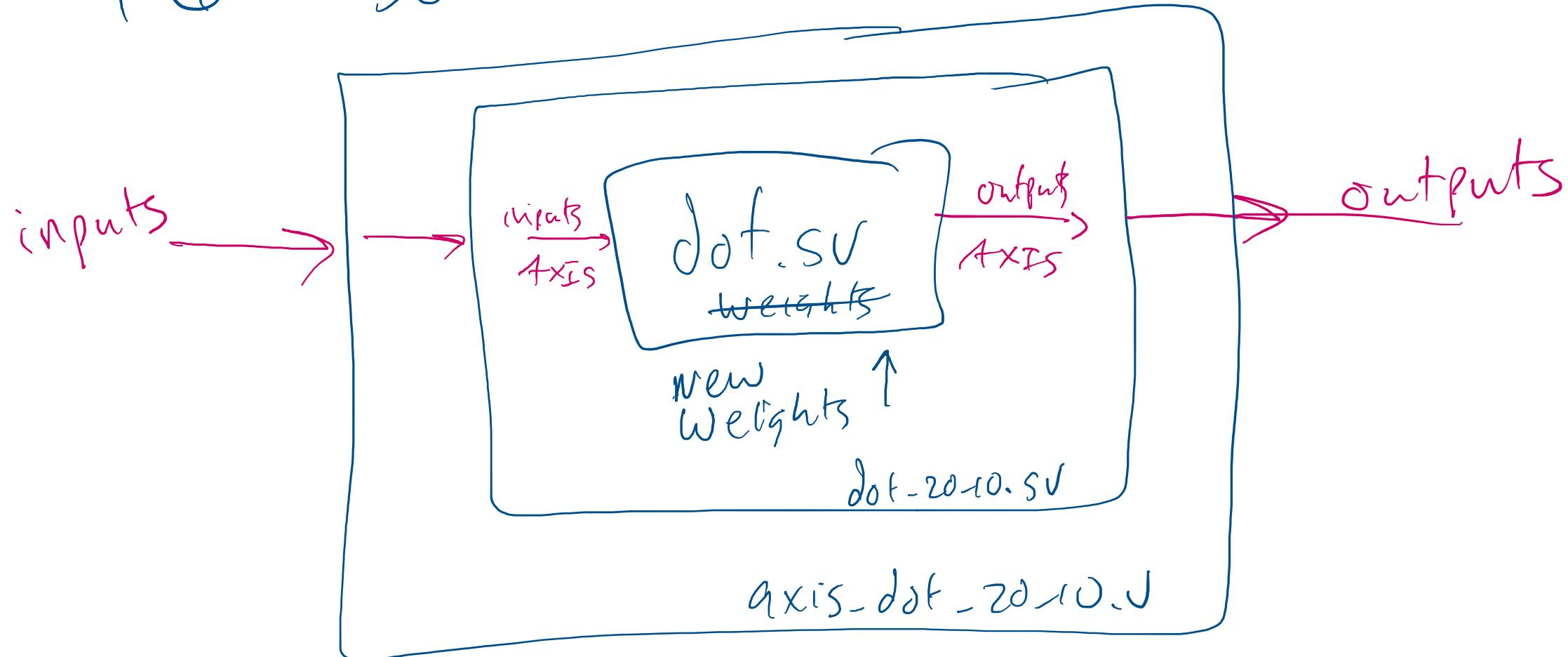


$$\begin{bmatrix} 0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 3 & 4 & 5 \\ 6 & 7 & 8 \end{bmatrix}$$

P6

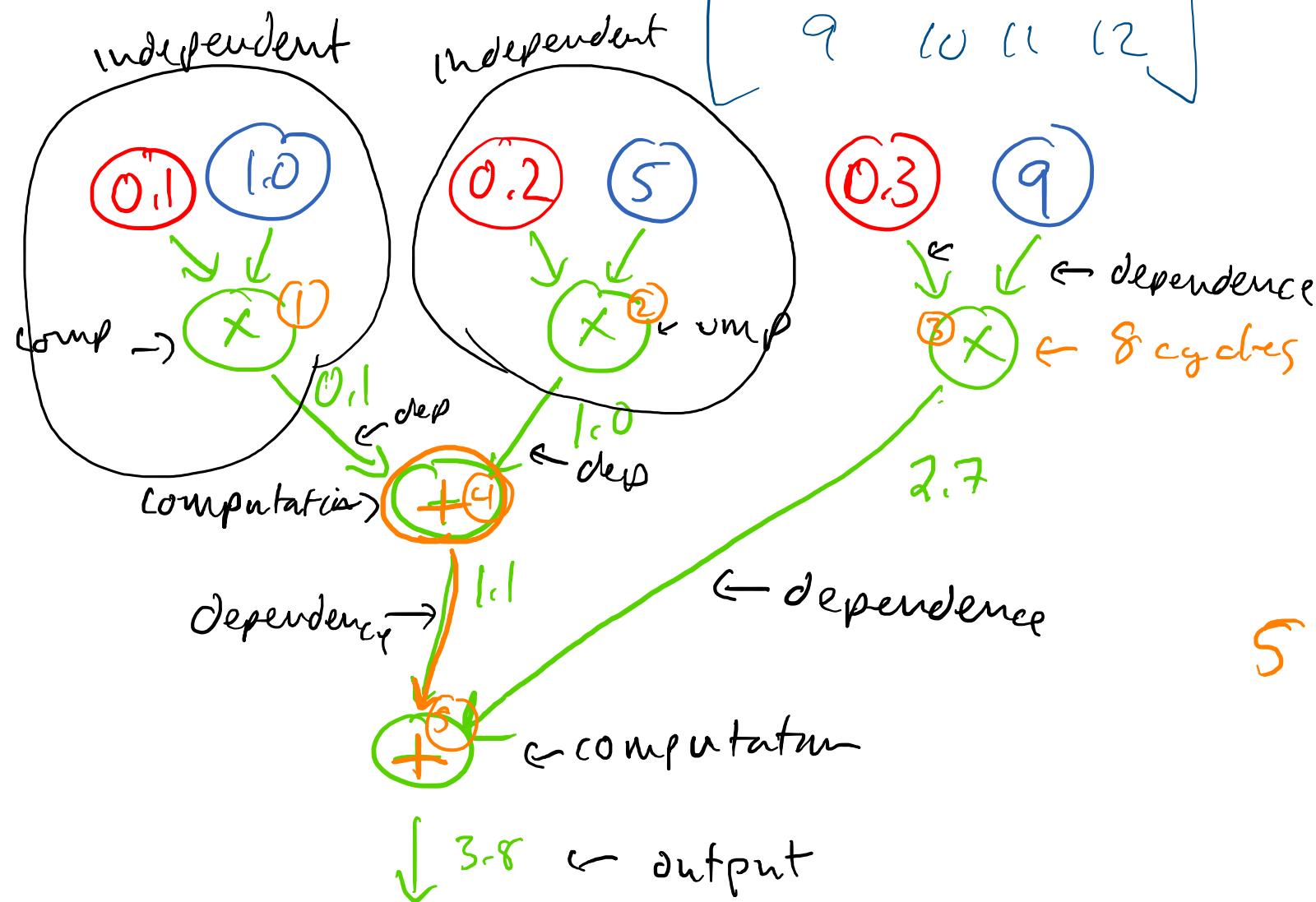
Dot

Data flow



Next Time: More parallelization

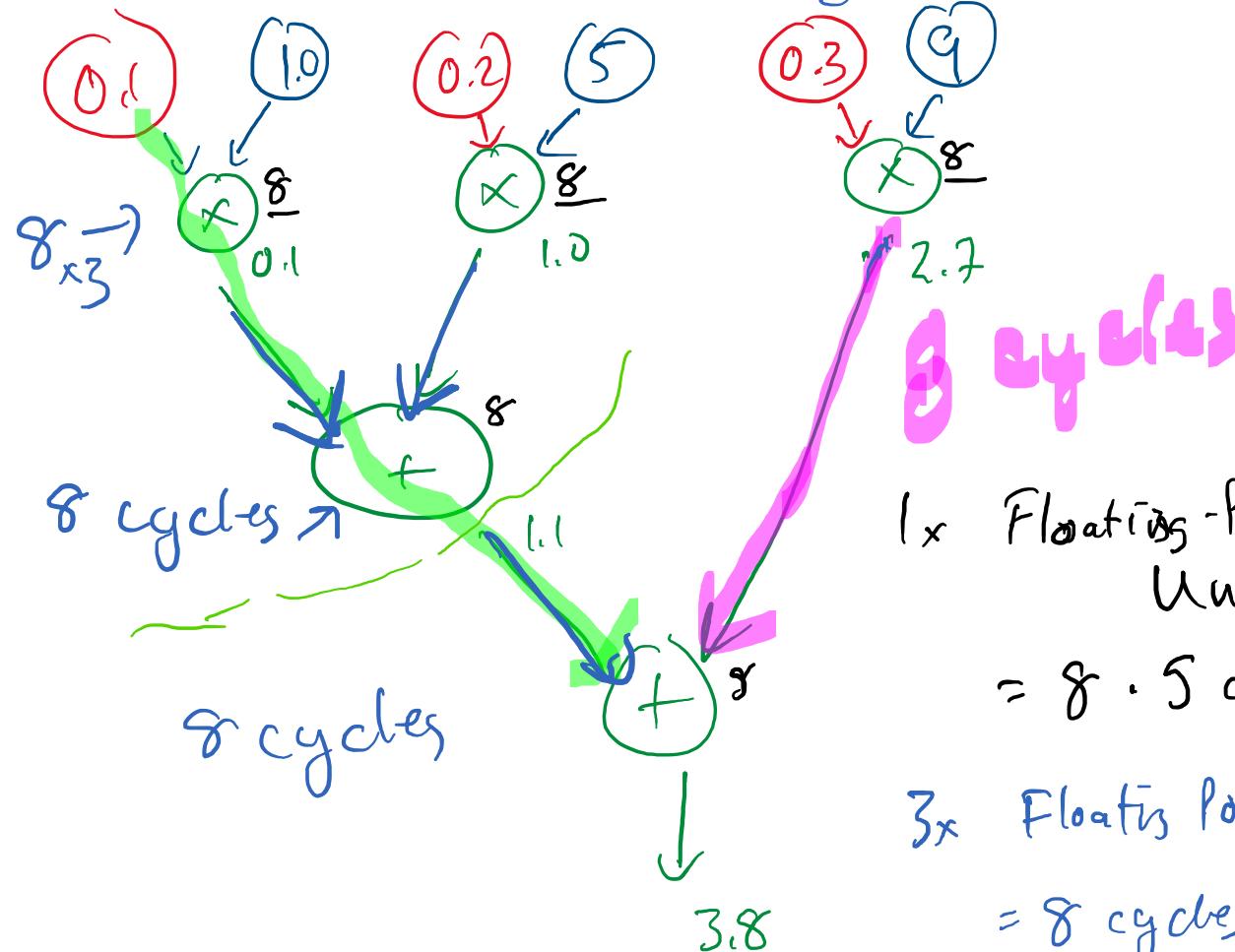
$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix} = \begin{bmatrix} 3.8 & 4.4 & 5 & 5.6 \end{bmatrix}$$



$$5 \times 8 = 40 \text{ cycles}$$

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

16 cycles

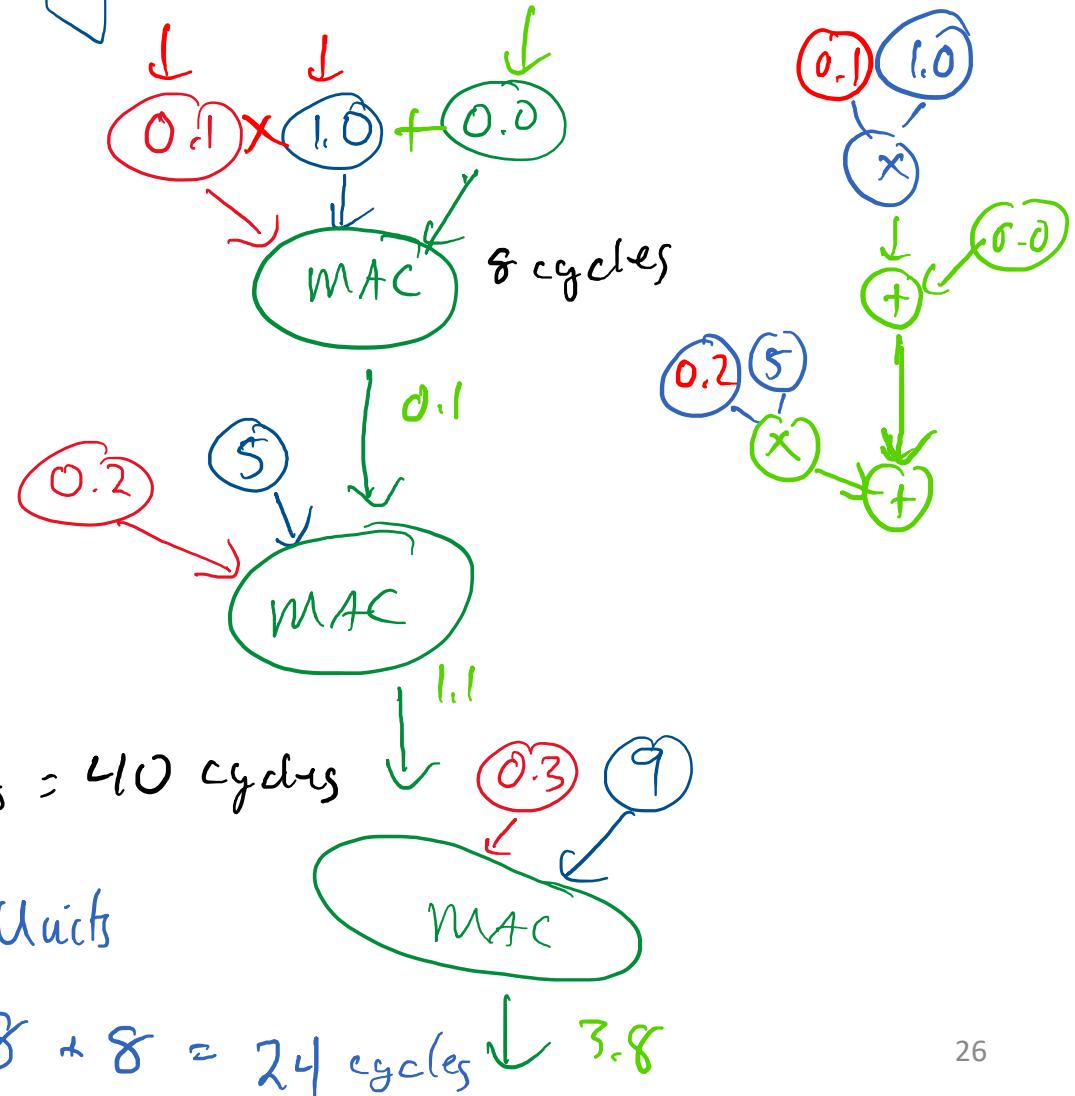


1x Floating-Point Unit

$$= 8 \cdot 5 \text{ cycles} = 40 \text{ cycles}$$

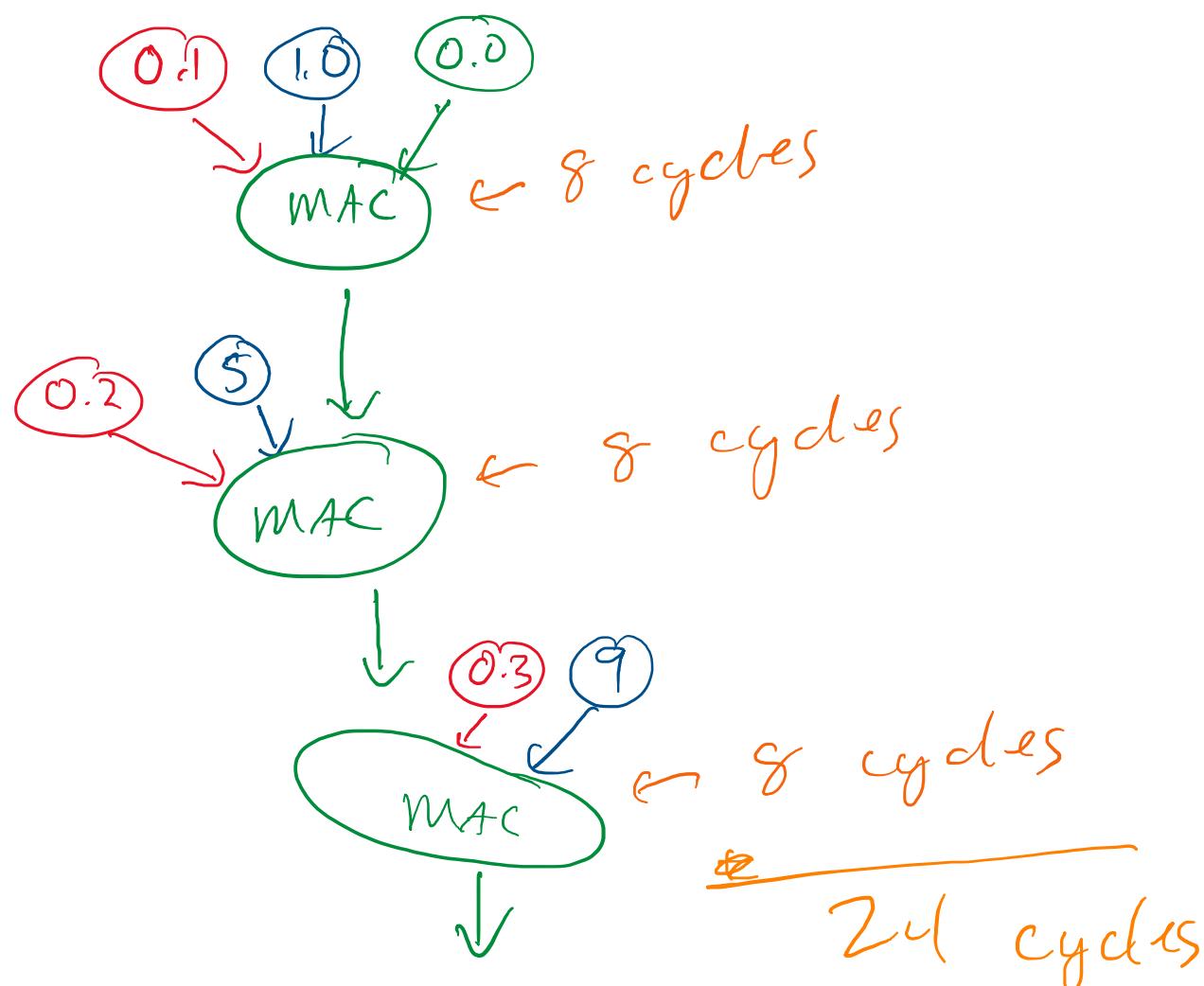
3x Floating Point Units

$$= 8 \text{ cycles} + 8 + 8 = 24 \text{ cycles} \downarrow 3.8$$



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

24 cycles 24 cycles 24 cycles 24 cycles



$$= 24 \cdot 4$$

$$= 96 \text{ cycles}$$

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

$$0.1 \cdot 1 + \underline{0} = \underline{0.1} \leftarrow +8 \text{ cycles}$$

$$\rightarrow 0.2 \cdot 5 + \underline{0.1} = \underline{1.1} \leftarrow +8 \text{ cycles}$$

$$0.3 \cdot 9 + \underline{1.1} = \underline{3.8} \leftarrow +8 \text{ cycles}$$

$$\rightarrow \text{result: } 3.8 \leftarrow 8 + 8 + 8 = 24 \text{ cycles}$$

$$4.4 \leftarrow 24 \text{ cycles}$$

$$5 \leftarrow 24$$

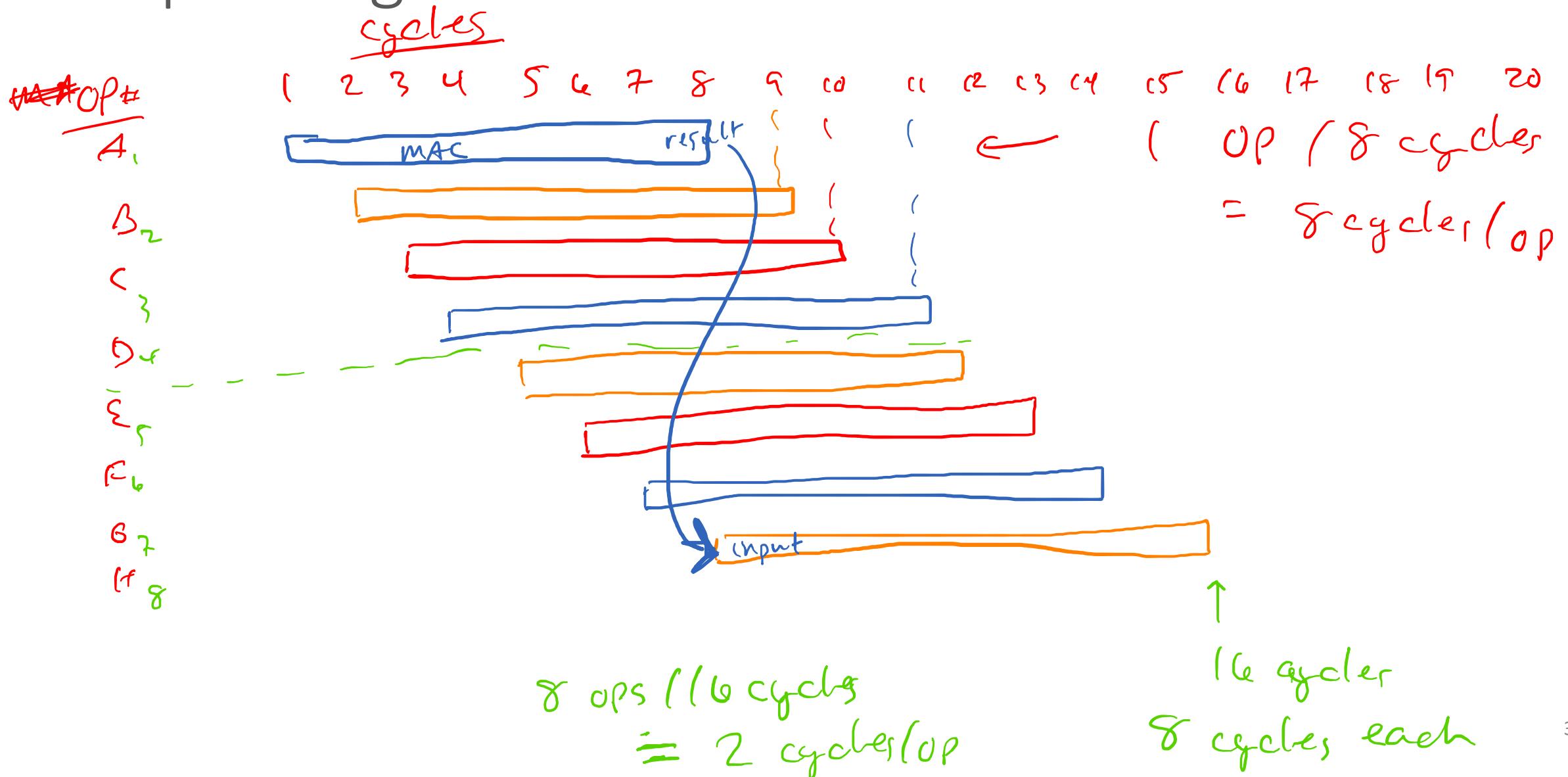
$$5.6 \leftarrow 24$$

96 cycles

Pipelining

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

Pipelining



$$\left[\begin{array}{ccc} \underline{0.1} & 0.2 & 0.3 \end{array} \right] \left[\begin{array}{cccc} 1 & 2 & 3 & 4 \\ \hline 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{array} \right] = \left[\begin{array}{cccc} 3.8 & 4.4 & 5 & 5.6 \end{array} \right]$$

$$\boxed{0.1 \cdot 1 + 0 = 0.1 \text{ (8 cycles)}}$$

$$0.2 \cdot 5 + 0.1 = 1.1 \text{ 8}$$

$$0.3 \cdot 9 + 1.1 = \underline{3.8} \text{ (8)}$$

$$0.1 \cdot 2 + 0 = 0.2$$

$$0.2 \cdot 6 + 0.2 = 1.4$$

$$0.3 \cdot 10 + 1.4 = \underline{4.4}$$

$$0.1 \cdot 1 + 0 = 0.1 \leftarrow$$

$$0.1 \cdot 2 + 0 = \underline{0.2} \leftarrow \text{4 pipelined MAC}$$

$$\rightarrow 0.1 \cdot 3 + 0 = 0.3 \leftarrow \text{in 11 cycles}$$

$$0.1 \cdot 4 + 0 = 0.4 \leftarrow$$

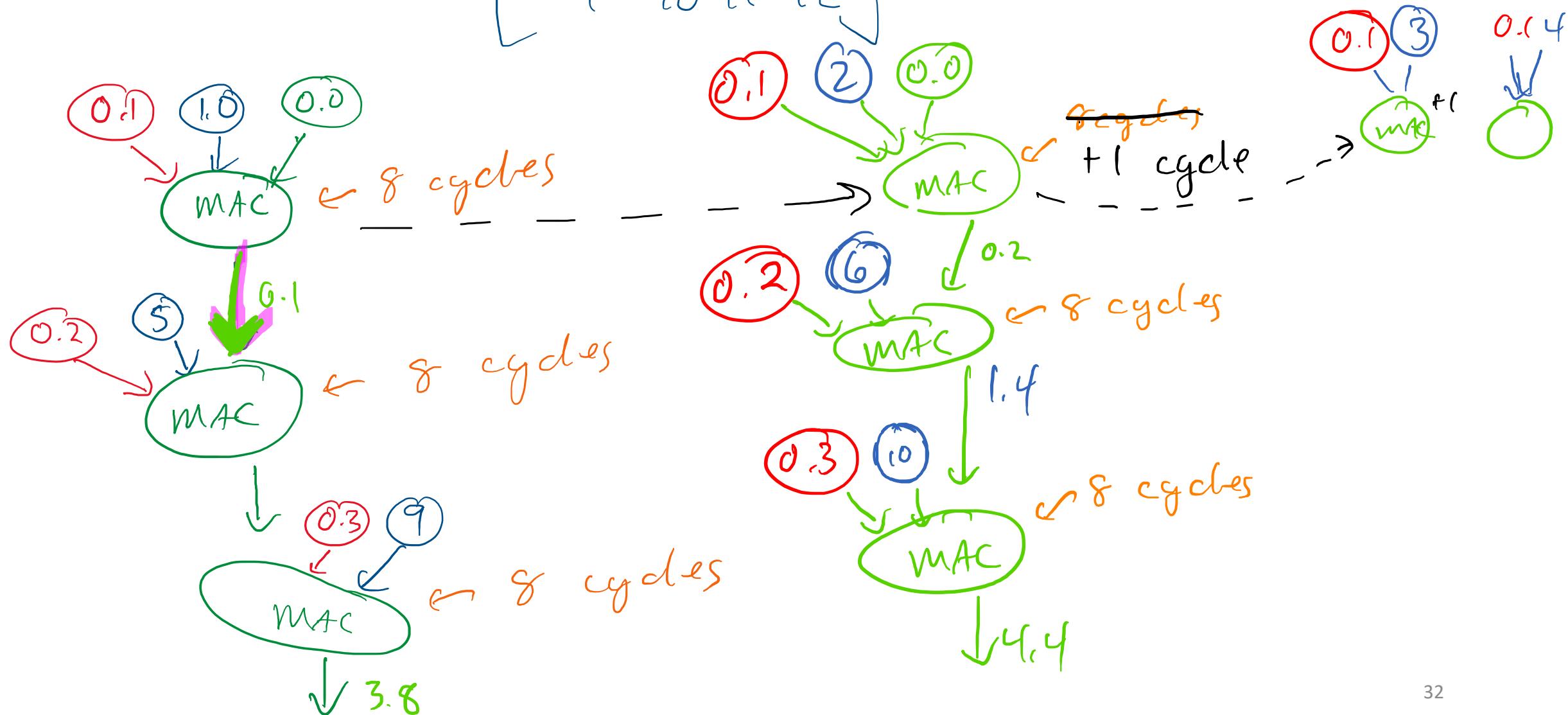
$$0.2 \cdot 5 + 0.1 =$$

$$0.2 \cdot 6 + \underline{0.2} =$$

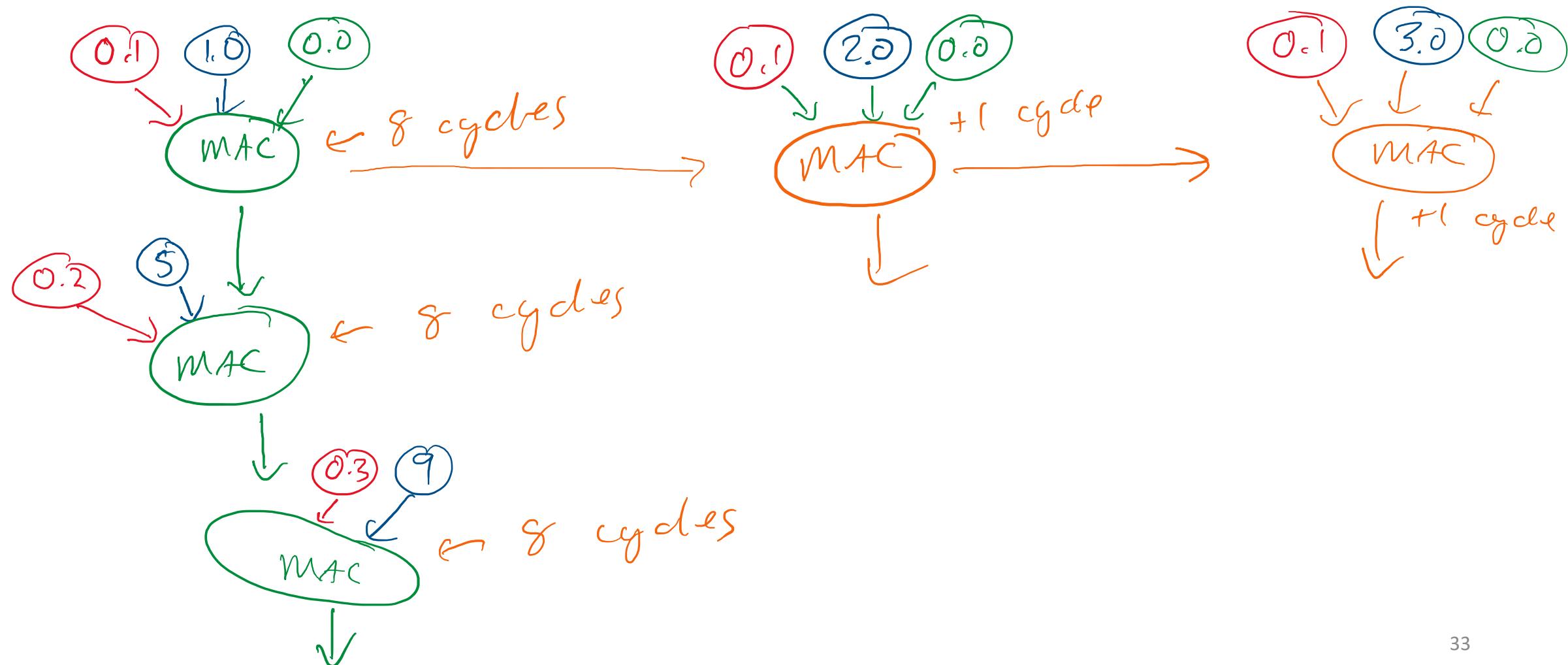
$$0.2 \cdot 7 + 0.3 =$$

$$0.2 \cdot 8 + 0.4 =$$

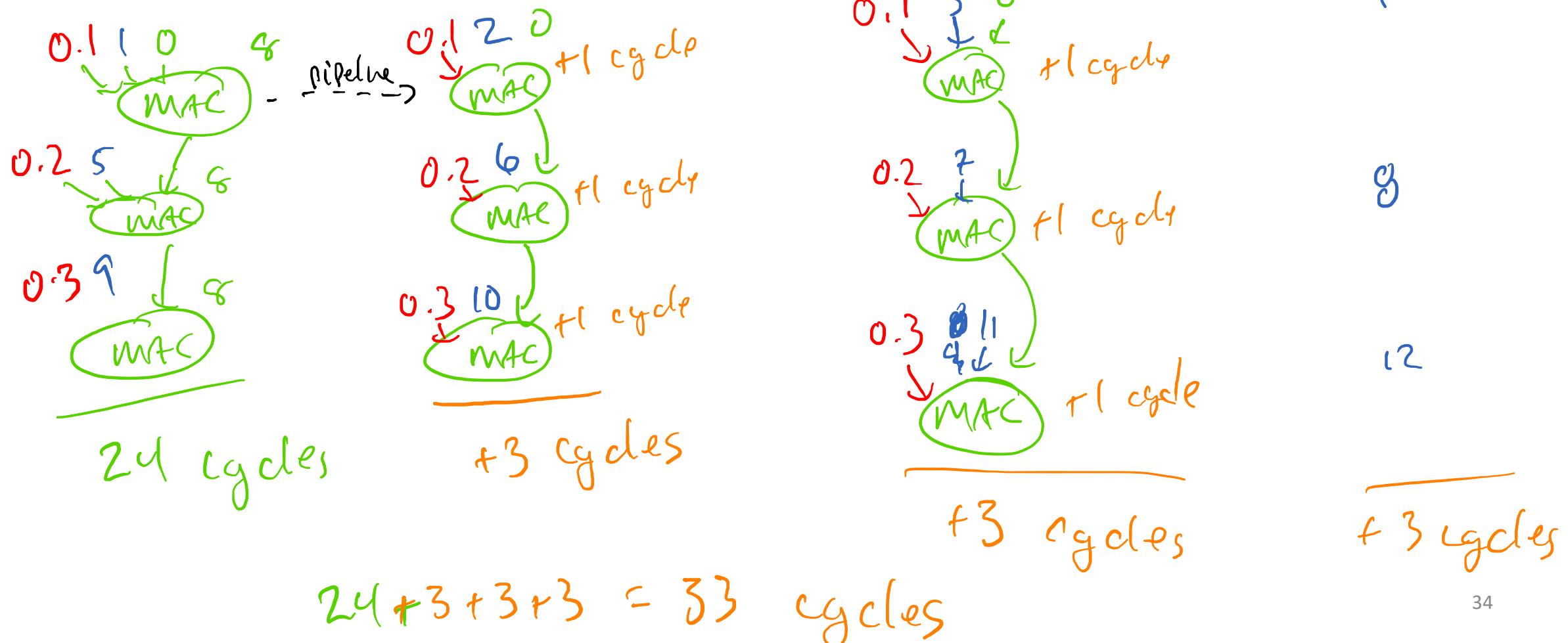
$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix} = \begin{bmatrix} 3.8 & 4.4 & 5 & 5.6 \end{bmatrix}$$



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



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



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

cycle 1:

$$\frac{10}{10}$$

$$\frac{2}{6} = \frac{10}{110}$$

$$\frac{10}{100} + \frac{1}{10} = 110$$

$$\frac{2}{4} + \frac{1}{2} = 6$$

cycle 2:

$$\frac{10}{100}$$

cycle 3:

$$10 + 100 = 110$$

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

Unpipelined:

8 cycles / MAC

X 4 MACs / input

X 3 inputs

96 cycles

pipelined:

8 cycles for 1st MAC

and +1 cycle for 2^{nd+} MAC

X 4 MACs / input

$8 + 1 + 1 + 1 = 11$ cycles

11 cycles / input

X 3 inputs

33 cycles

Hardware Parallelism

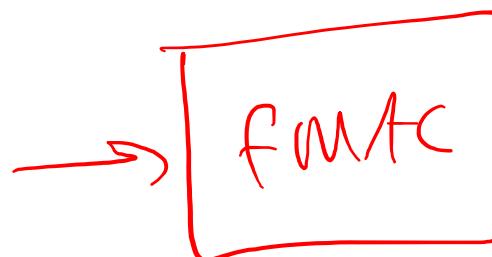
- CPU: 1 Floating-Point Unit
- FPGA? *(10 Floating-Point Units?
20?
100?)*

Finding Parallelism

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

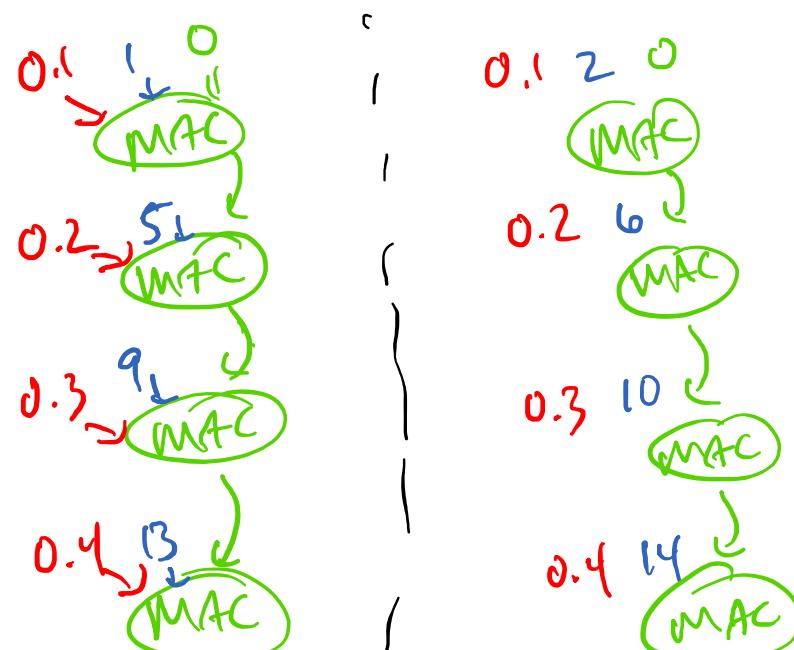
What can we do with 2+ FMACs?

dot



$$\begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} = \begin{bmatrix} 9 & 10 & 11 & 12 \end{bmatrix}$$

independence



\nwarrow \nearrow
totally independent

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 5 & 6 \\ 9 & 10 \\ 13 & 14 \end{bmatrix} = \begin{bmatrix} 9 & 10 \end{bmatrix}$$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \end{bmatrix} \begin{bmatrix} 3 & 4 \\ 7 & 8 \\ 11 & 12 \\ 15 & 16 \end{bmatrix} = \begin{bmatrix} 11 & 12 \end{bmatrix}$$

$$\begin{bmatrix} 9 & 10 & 11 & 12 \end{bmatrix}$$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} = \begin{bmatrix} 9 & 10 & 11 & 12 \end{bmatrix}$$

$$\begin{array}{r}
 \begin{array}{rr}
 0.1 & 0.2 \\
 \times 1 & \times 5 \\
 \hline
 0.1 & + 5.2 \\
 \hline
 5.2
 \end{array}
 \quad
 \begin{array}{rr}
 0.3 & 0.4 \\
 \times 9 & \times 13 \\
 \hline
 2.7 & + 5.2 \\
 \hline
 7.9
 \end{array}
 \quad
 \begin{array}{rr}
 0.1 & 0.2 \\
 \times 2 & \times 6 \\
 \hline
 0.2 & + 1.2 \\
 \hline
 1.4
 \end{array}
 \quad
 \begin{array}{rr}
 0.3 & 0.4 \\
 \times 10 & \times 14 \\
 \hline
 3.0 & + 5.6 \\
 \hline
 8.6
 \end{array}
 \quad
 \begin{array}{rr}
 0.1 & 0.2 \\
 \times 3 & \times 7 \\
 \hline
 0.3 & + 1.4 \\
 \hline
 1.7
 \end{array}
 \quad
 \begin{array}{rr}
 0.3 & 0.4 \\
 \times 11 & \times 15 \\
 \hline
 3.3 & + 6.0 \\
 \hline
 9.3
 \end{array}
 \quad
 \begin{array}{rr}
 0.1 & 0.2 \\
 \times 4 & \times 8 \\
 \hline
 0.4 & + 1.6 \\
 \hline
 2.0
 \end{array}
 \quad
 \begin{array}{rr}
 0.3 & 0.4 \\
 \times 12 & \times 16 \\
 \hline
 3.6 & + 6.4 \\
 \hline
 10
 \end{array}
 \end{array}
 \end{array}$$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} = \begin{bmatrix} 9 & 10 & 11 & 12 \end{bmatrix}$$

$$\begin{array}{r}
 \begin{array}{rr}
 0.1 & 0.2 \\
 \times 1 & \times 5 \\
 \hline
 0.1 & + 5.2 \\
 \hline
 5.2
 \end{array}
 \quad
 \begin{array}{rr}
 0.3 & 0.4 \\
 \times 9 & \times 13 \\
 \hline
 2.7 & + 5.2 \\
 \hline
 7.9
 \end{array}
 \quad
 \begin{array}{rr}
 0.1 & 0.2 \\
 \times 2 & \times 6 \\
 \hline
 0.2 & + 1.2 \\
 \hline
 1.4
 \end{array}
 \quad
 \begin{array}{rr}
 0.3 & 0.4 \\
 \times 10 & \times 14 \\
 \hline
 3.0 & + 5.6 \\
 \hline
 8.6
 \end{array}
 \quad
 \begin{array}{rr}
 0.1 & 0.2 \\
 \times 3 & \times 7 \\
 \hline
 0.3 & + 1.4 \\
 \hline
 1.7
 \end{array}
 \quad
 \begin{array}{rr}
 0.3 & 0.4 \\
 \times 11 & \times 15 \\
 \hline
 3.3 & + 6.0 \\
 \hline
 9.3
 \end{array}
 \quad
 \begin{array}{rr}
 0.1 & 0.2 \\
 \times 4 & \times 8 \\
 \hline
 0.4 & + 1.6 \\
 \hline
 2.0
 \end{array}
 \quad
 \begin{array}{rr}
 0.3 & 0.4 \\
 \times 12 & \times 16 \\
 \hline
 3.6 & + 6.4 \\
 \hline
 10
 \end{array}
 \end{array}$$

Independent Independent Independent Independent

$$[0.1 \ 0.2 \ 0.3 \ 0.4] \cdot \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} = [9 \ 10 \ 11 \ 12]$$

$$\begin{bmatrix} 0.1 & 0.3 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 5 & 7 \\ 9 & 11 \\ 13 & 15 \end{bmatrix} + \begin{bmatrix} 0.2 & 0.4 \end{bmatrix} \begin{bmatrix} 2 & 4 \\ 6 & 8 \\ 10 & 12 \\ 14 & 16 \end{bmatrix} = [9 \ 10 \ 11 \ 12]$$

WRONG!

$$[0.1 \ 0.2 \ 0.3 \ 0.4] \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} = [9 \ 10 \ 11 \ 12]$$

$$[0.1 \ 0.3] \begin{bmatrix} 1 & 2 & 3 & 4 \\ 9 & 10 & 11 & 12 \end{bmatrix} + [0.2 \ 0.4] \begin{bmatrix} 5 & 6 & 7 & 8 \\ 13 & 14 & 15 & 16 \end{bmatrix} =$$

$$[2.8 \ 3.2 \ 3.6 \ 4.0] + [6.2 \ 6.8 \ 7.4 \ 8]$$

$$= [9 \ 10 \ 11 \ 12]$$

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
```

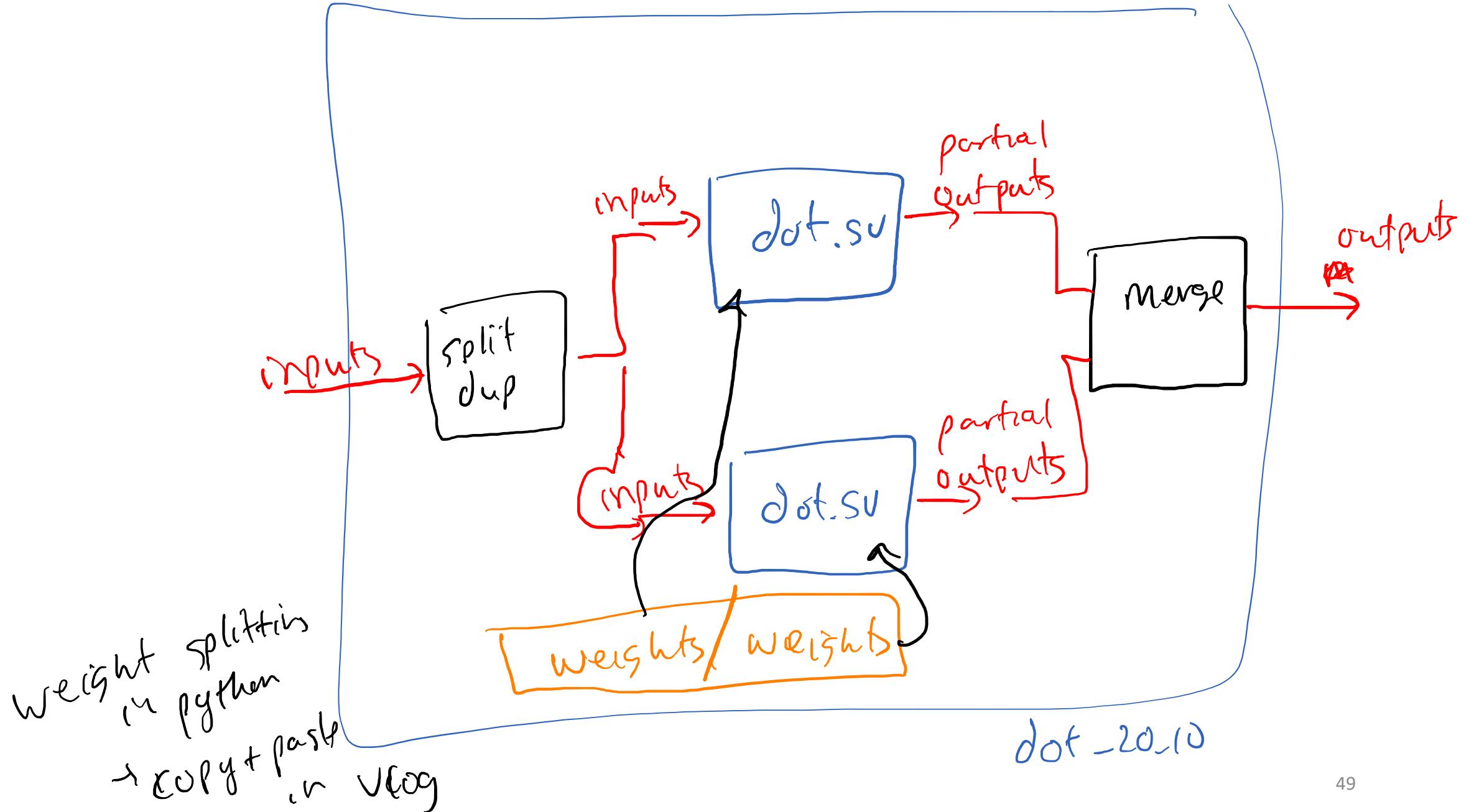
Can we parallelize Dot?

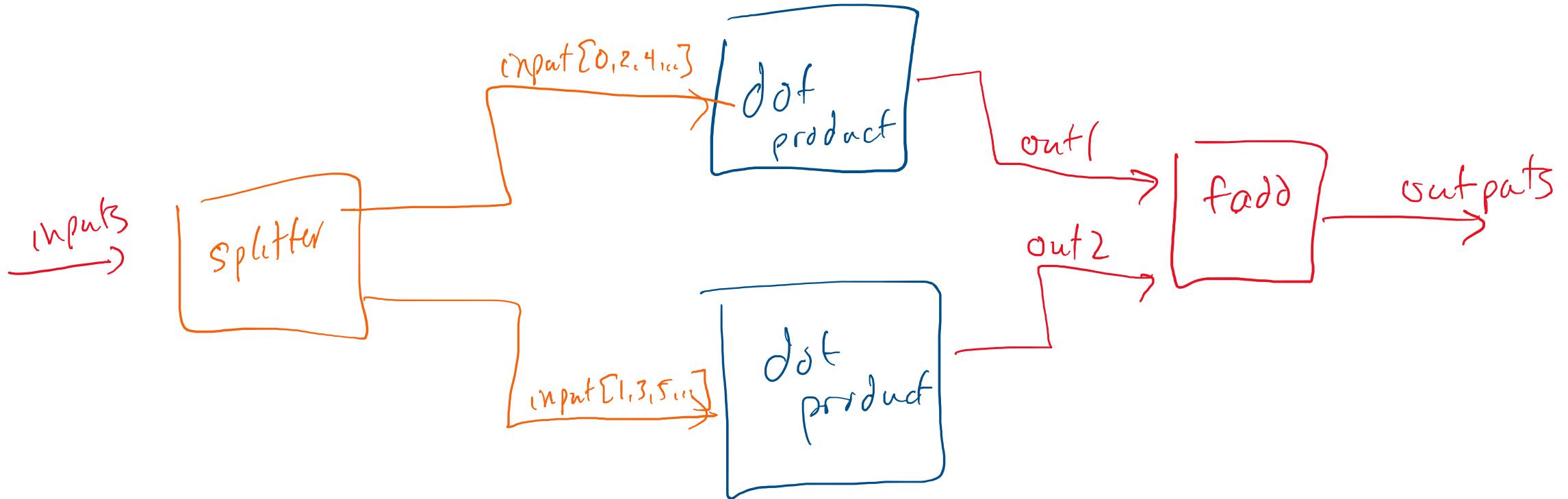
```
# how its done in dot.sv
def pydot(inputs,weights):
    inputs = inputs[0] # remove outer nesting
    outs = np.zeros(weights.shape[1], dtype=np.float32)
    for i in range(weights.shape[0]): # input length
        for j in range(weights.shape[1]): # output length
            outs[j] = outs[j] +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





16: Hardware Acceleration III

Engr 315: Hardware / Software Codesign

Andrew Lukefahr

Indiana University

