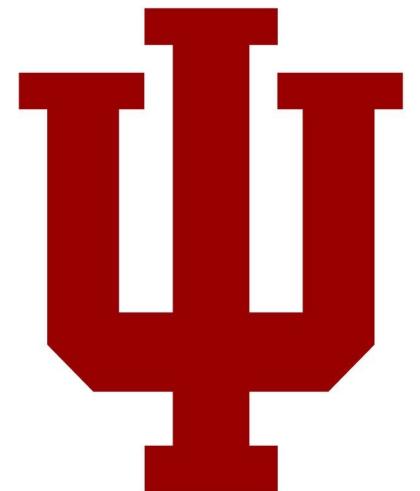


17: Hardware Acceleration IV

Engr 315: Hardware / Software Codesign

Andrew Lukefahr

Indiana University



Announcements

- Exams almost done grading.

- P7 is out → Dot product

≤ 300 cycles

- P8 is coming. → Dot product

how fast can
you go

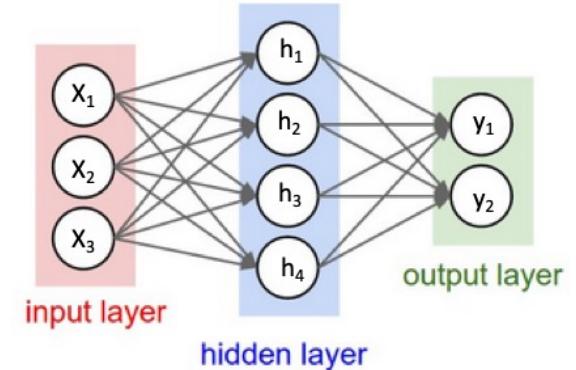
Project 7: Accelerate Dot Product

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

Project 8: Accelerate Dot Product More

- Your Task: Accelerate P7
- How? Pipelining + Parallelism
- Same groups of 2 allowed.
- Additional documentation requirements

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)
```

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

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

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

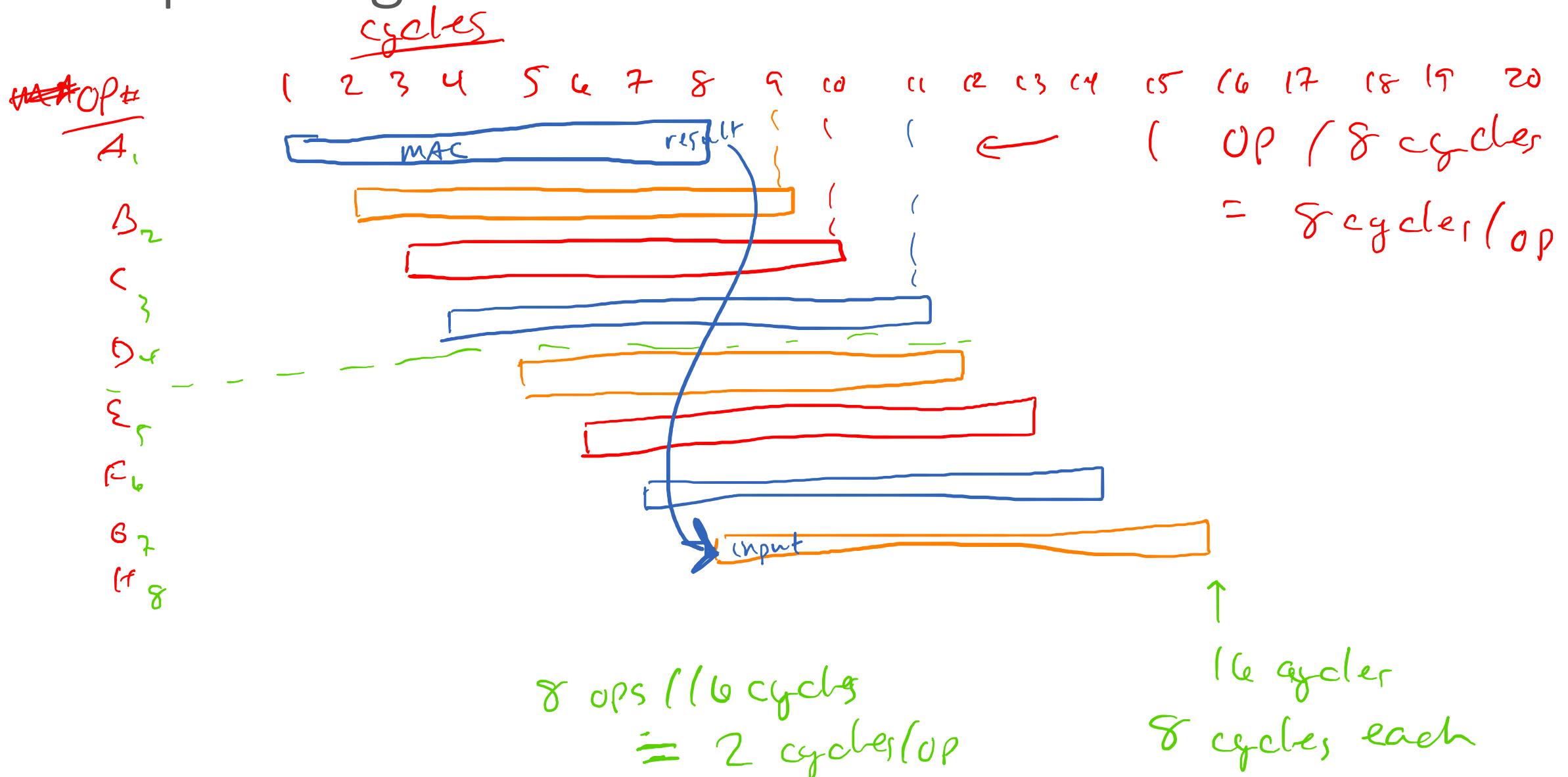
$$0.1 \cdot \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} =$$
$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \leftarrow^{\text{temp result}}$$

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

Pipelining

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

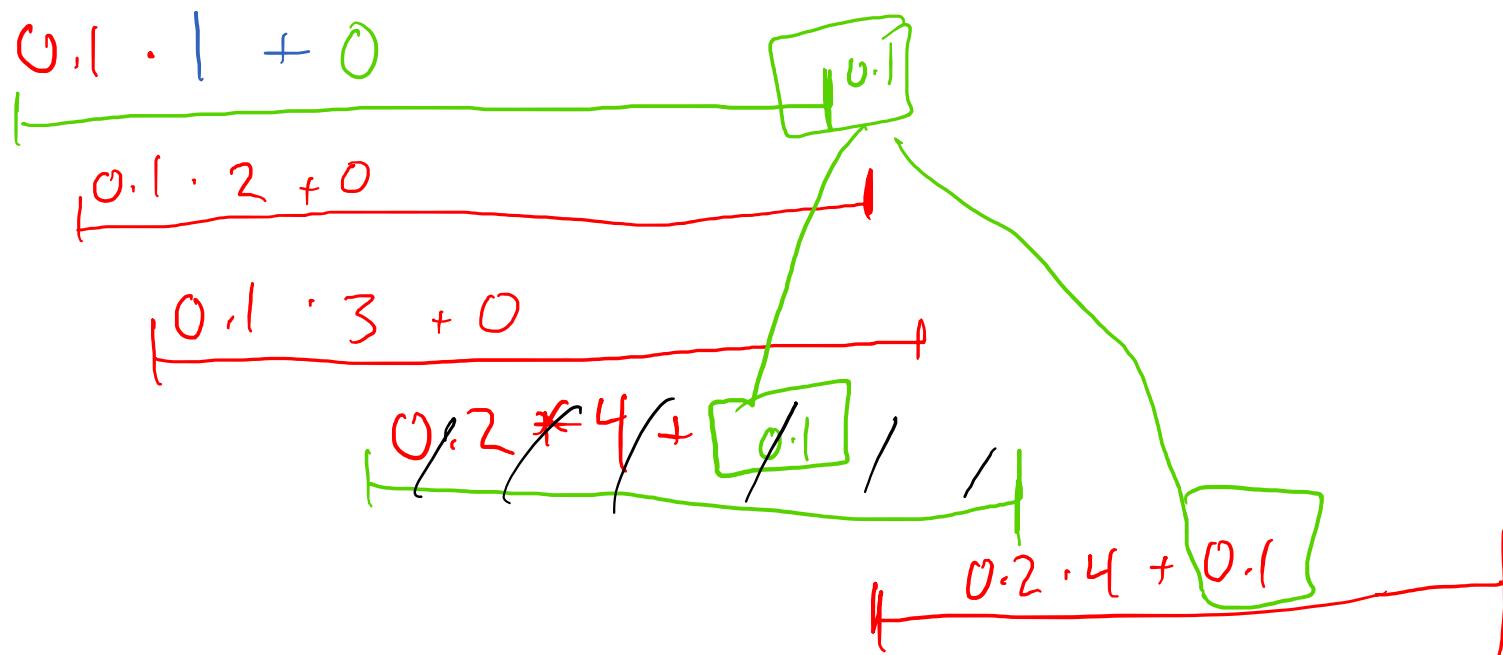
Pipelining



Pipelined Dot Computations

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

time →



Latency on Pipelined FMAC

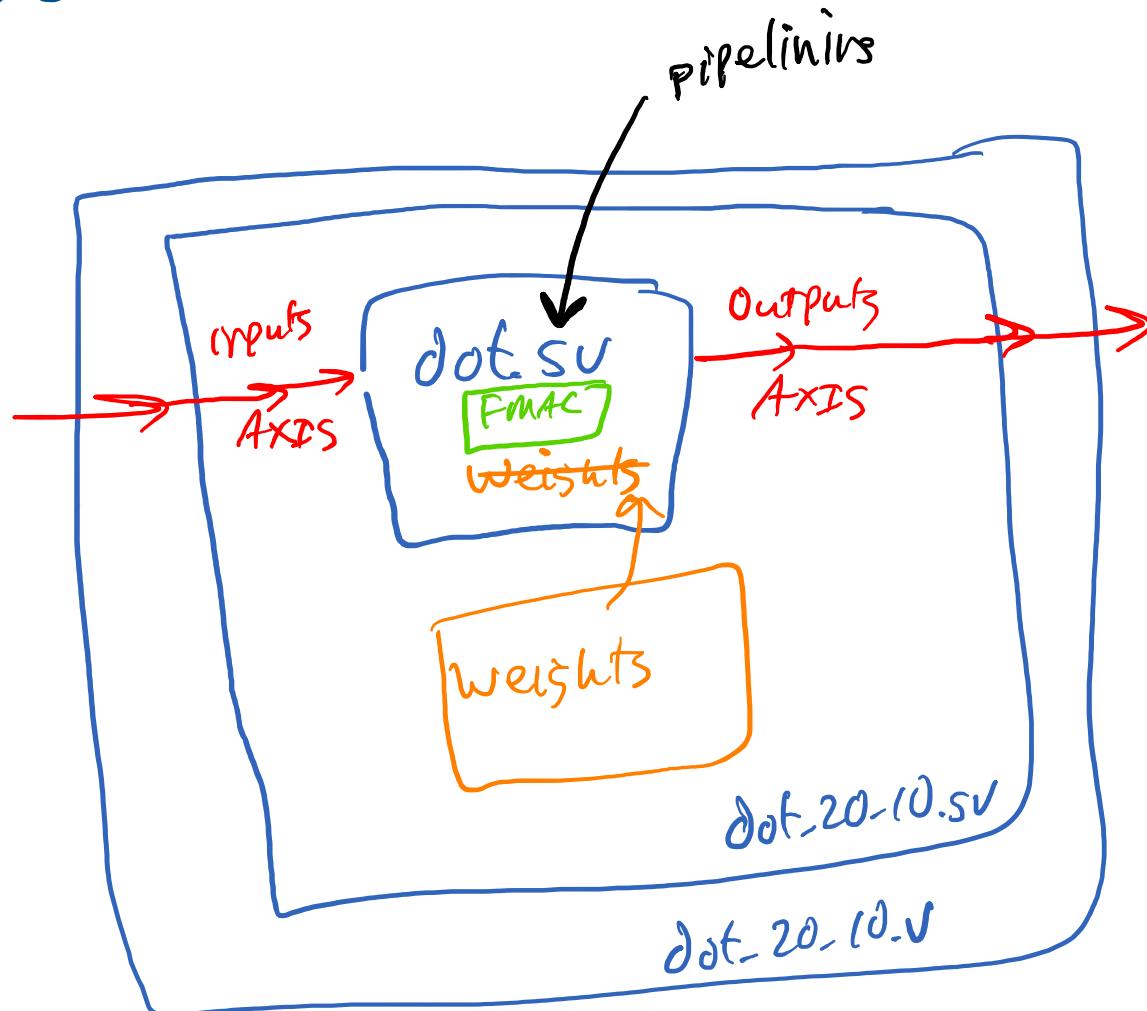
- **Pro Tip:** Stall at the end of a row.
- Drain the FMAC pipeline.

DT
PLo

Dot

Data flow

[



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

Verilog Parameters

- Parameters are Verilog constructs that allow a module to be **reused with a different specification.**

```
module adder #(parameter BITS = 2) (
    input [BITS-1:0] a,
    inputs [BITS-1:0] b,
    outputs [BITS-1:0] c);
    assign c = a + b;
endmodule
```

Verilog Parameters

```
// 2-bit adder  
adder add2 (a2, b2, c2);  
//another 2-bit adder  
adder #(BITS=2) add2b (a2b,b2b,c2b);  
//a 3-bit adder  
adder #(BITS=3) add3 (a3, b3, c3);  
  
Adder #(BITS=32) add32 (a32, b32, c32);
```

```
module adder #(parameter BITS = 2) (  
    input [BITS-1:0] a,  
    inputs [BITS-1:0] b,  
    outputs [BITS-1:0] c) );  
    assign c = a + b;  
endmodule
```

Verilog Parameters in Dot.sv

```
module dot #(  
    parameter ROWS = 3,  
    parameter COLS = 4,  
    parameter [31:0] weights [0:ROWS-1] [0:COLS-1] = '{  
        '{$shortrealtobits(1.0), $shortrealtobits(2.0), $shortrealtobits(3.0), $shortrealtobits(4.0)},  
        '{$shortrealtobits(5.0), $shortrealtobits(6.0), $shortrealtobits(7.0), $shortrealtobits(8.0)},  
        '{$shortrealtobits(9.0), $shortrealtobits(10.0), $shortrealtobits(11.0), $shortrealtobits(12.0)}  
    }  
)(  
    input clk, rst,  
    // Incomming Matrix AXI4-Stream  
    input [31:0] INPUT_AXIS_TDATA,  
    input INPUT_AXIS_TLAST,  
    input INPUT_AXIS_TVALID,  
    output reg INPUT_AXIS_TREADY,  
    // Outgoing Vector AXI4-Stream  
    output reg [31:0] OUTPUT_AXIS_TDATA,  
    output reg OUTPUT_AXIS_TLAST,  
    output reg OUTPUT_AXIS_TVALID,  
    input OUTPUT_AXIS_TREADY );
```

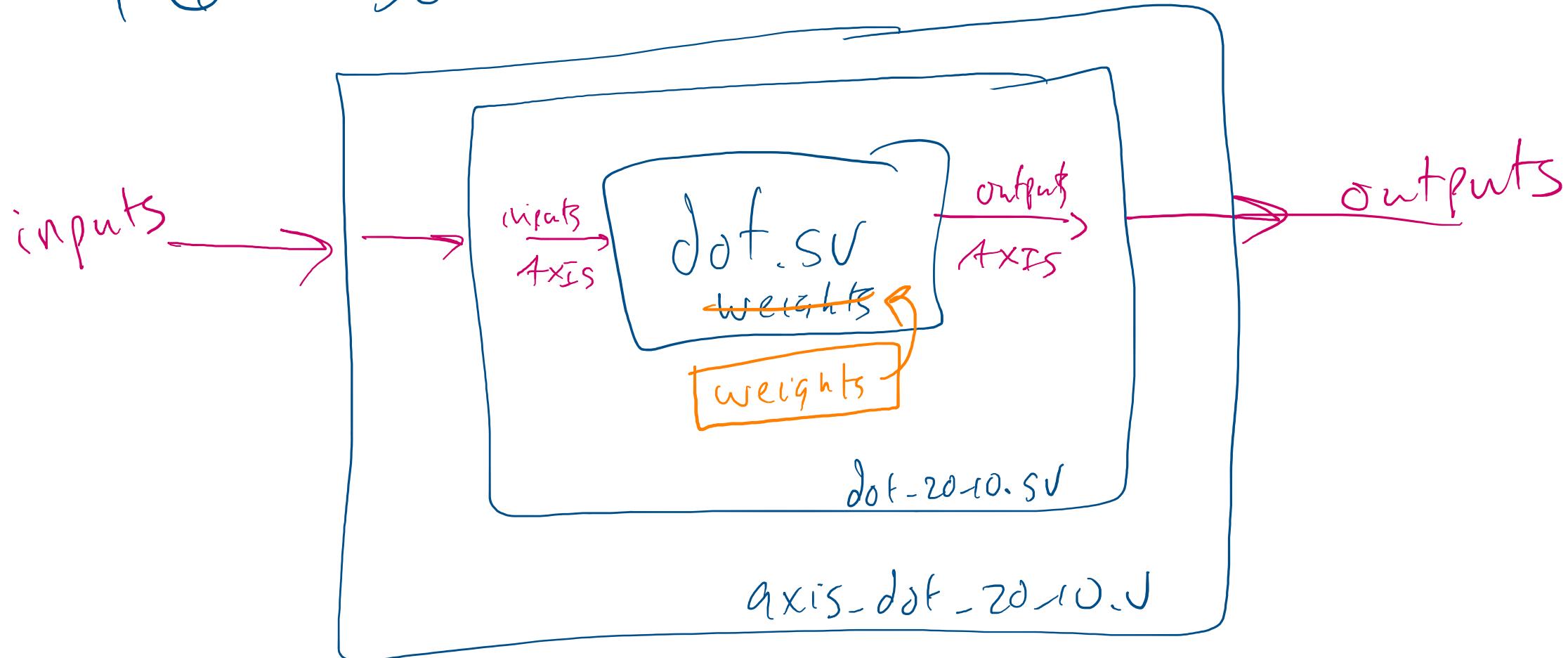
Verilog Parameters in dot_20_10.sv

```
dot #(  
    .ROWS(ROWS),  
    .COLS(COLS),  
    .weights(weights)  
) dot0 (  
    // AXI4-Stream Interface  
    .clk(clk),  
    .rst(rst),  
    .INPUT_AXIS_TDATA(INPUT_AXIS_TDATA),  
    .INPUT_AXIS_TLAST(INPUT_AXIS_TLAST),  
    .INPUT_AXIS_TVALID(INPUT_AXIS_TVALID),  
    .INPUT_AXIS_TREADY(INPUT_AXIS_TREADY),  
    .OUTPUT_AXIS_TDATA(OUTPUT_AXIS_TDATA),  
    .OUTPUT_AXIS_TLAST(OUTPUT_AXIS_TLAST),  
    .OUTPUT_AXIS_TVALID(OUTPUT_AXIS_TVALID),  
    .OUTPUT_AXIS_TREADY(OUTPUT_AXIS_TREADY)  
) ;
```

P6

Dot

Data flow



Data flow graph

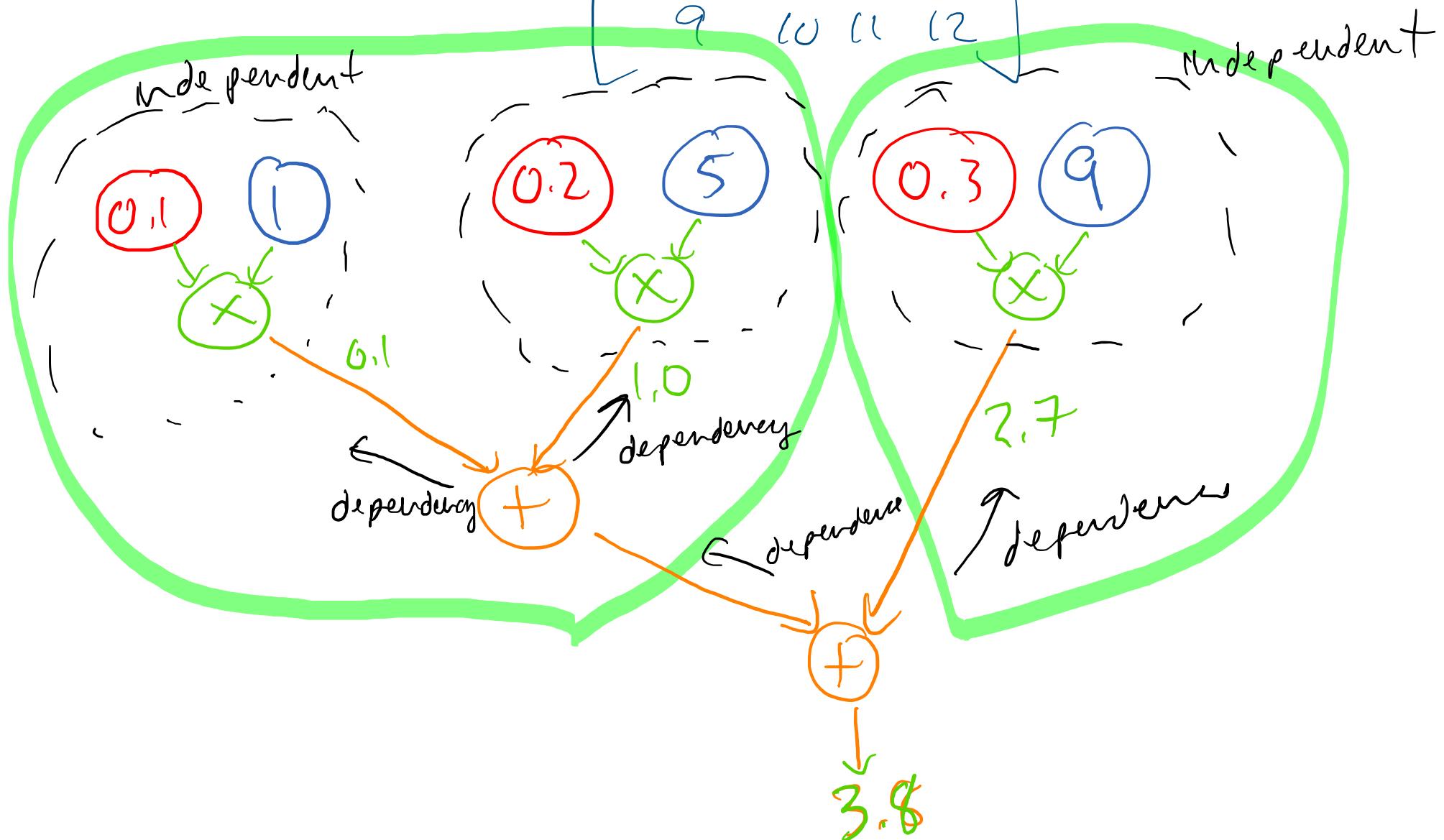
Fixme: add 0.4

{0.1 0.2 0.3 0.4}

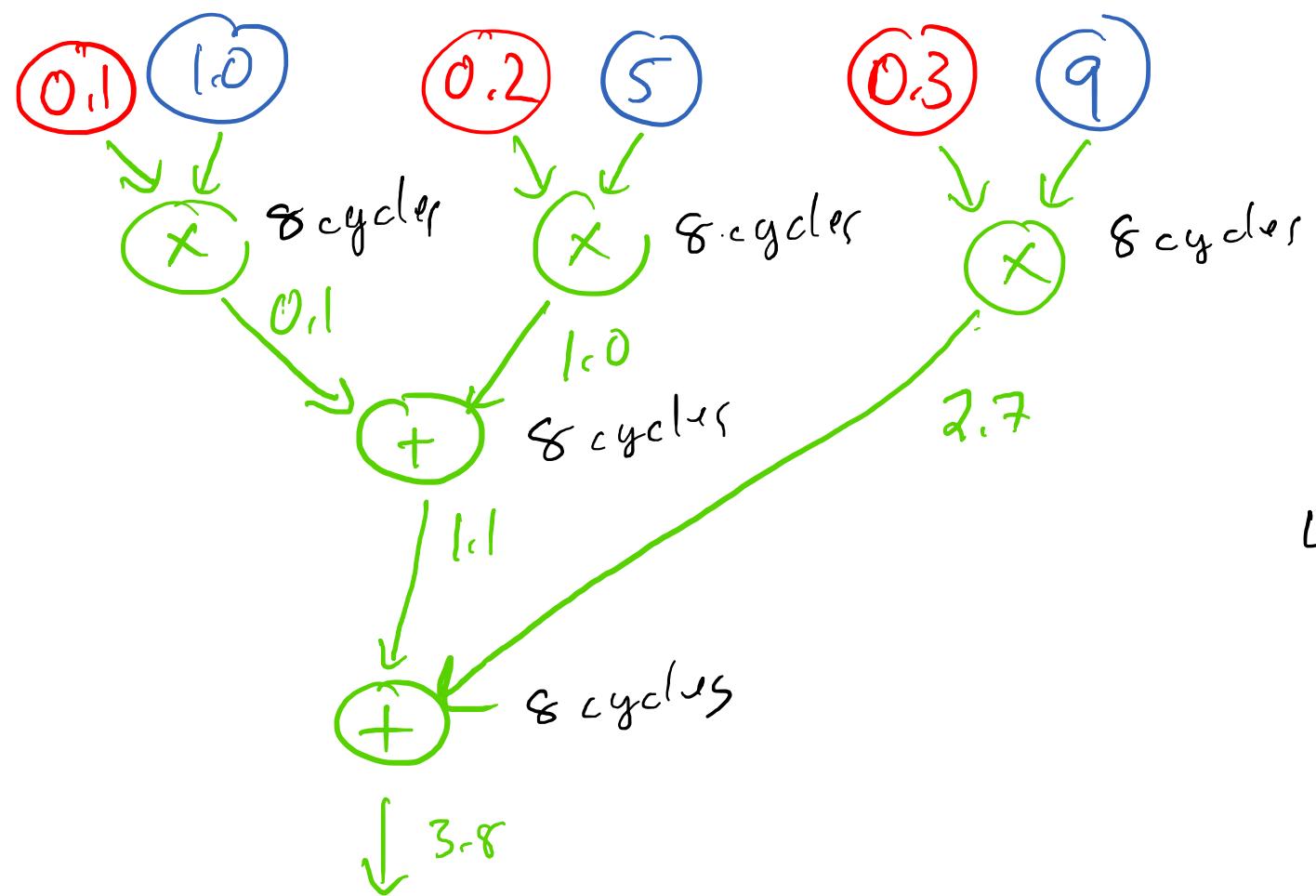
- A data-flow graph is a **collection of arcs and nodes** in which the nodes are either **places where variables are assigned** or used, and the arcs show the relationship between the places where a variable is assigned and where the assigned value is subsequently used.

[\[ref\]](#)

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



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

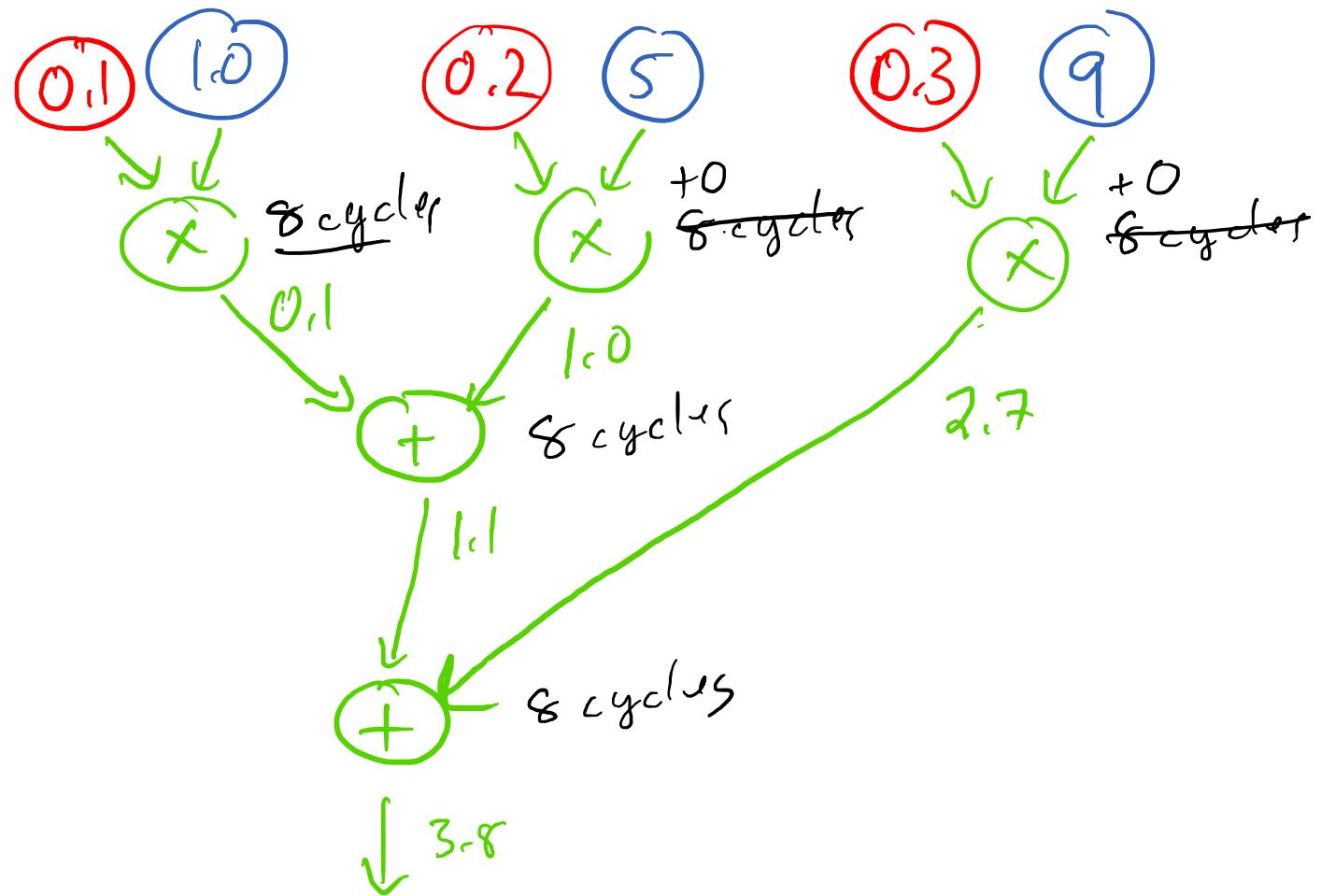


$$8 \cdot 5 = 40 \text{ cycles}$$

$$40 \text{ cycles} \cdot 4 \#^{'s} =$$

$$160 \text{ cycles}$$

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



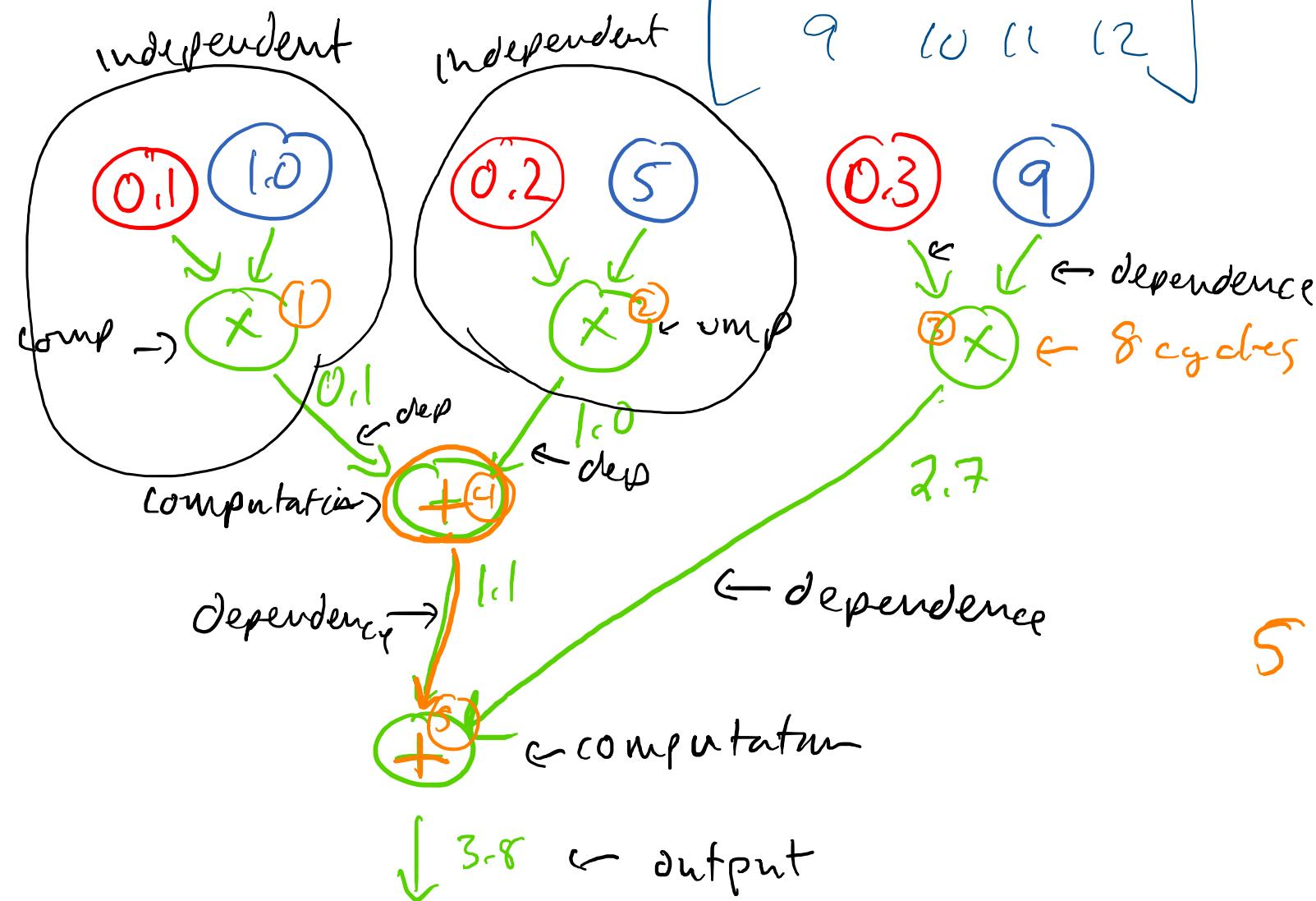
$3 \times$ multipliers
(or Adders)

$$\begin{array}{r}
 8 + 0 + 0 \\
 + 8 \\
 \hline
 + 8
 \end{array}$$

24 cycles / #

$$24 \cdot 4 = 96 \text{ cycles}$$

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



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

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

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

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

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

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

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

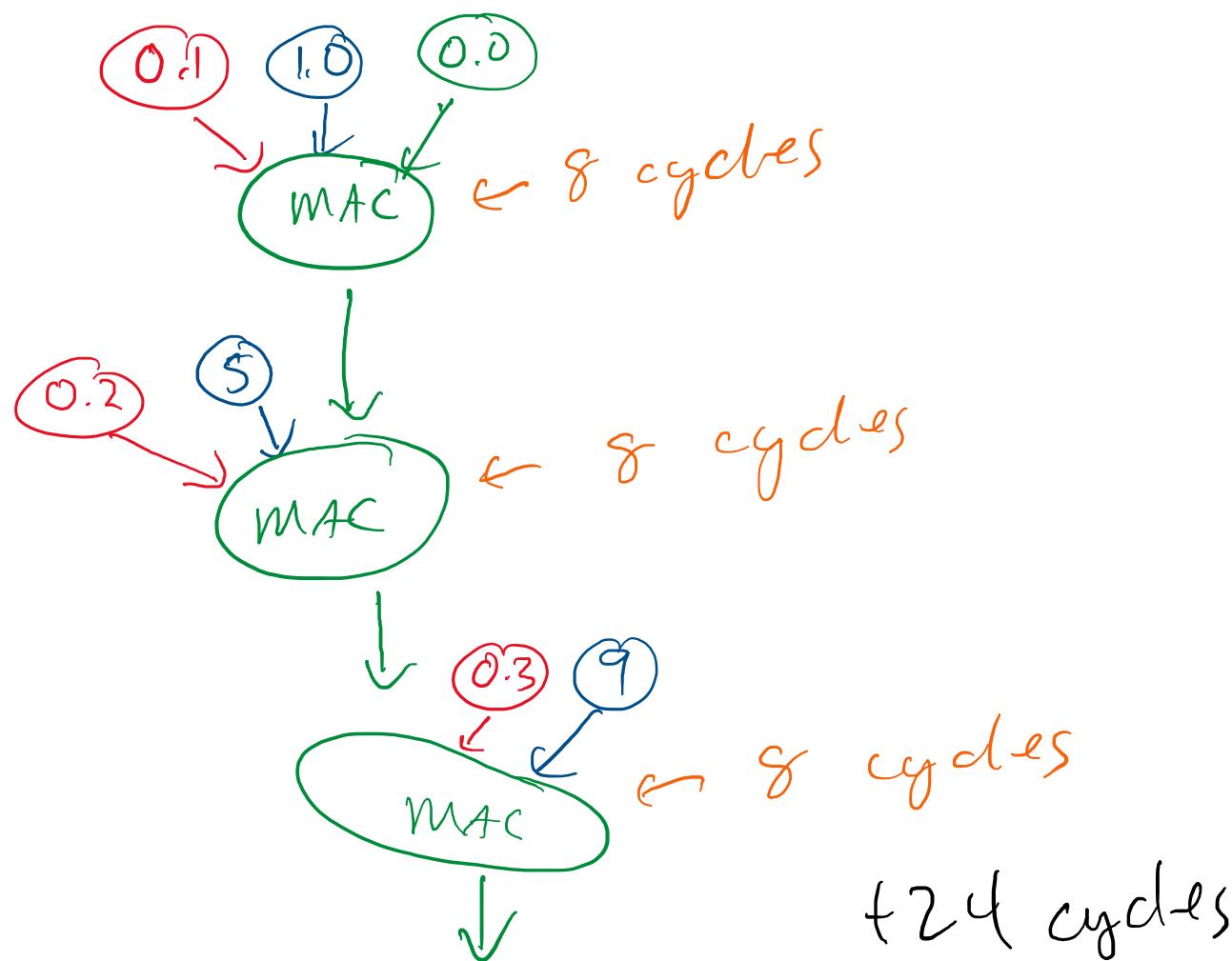
$$4.4 \leftarrow 24 \text{ cycles}$$

$$5 \leftarrow 24$$

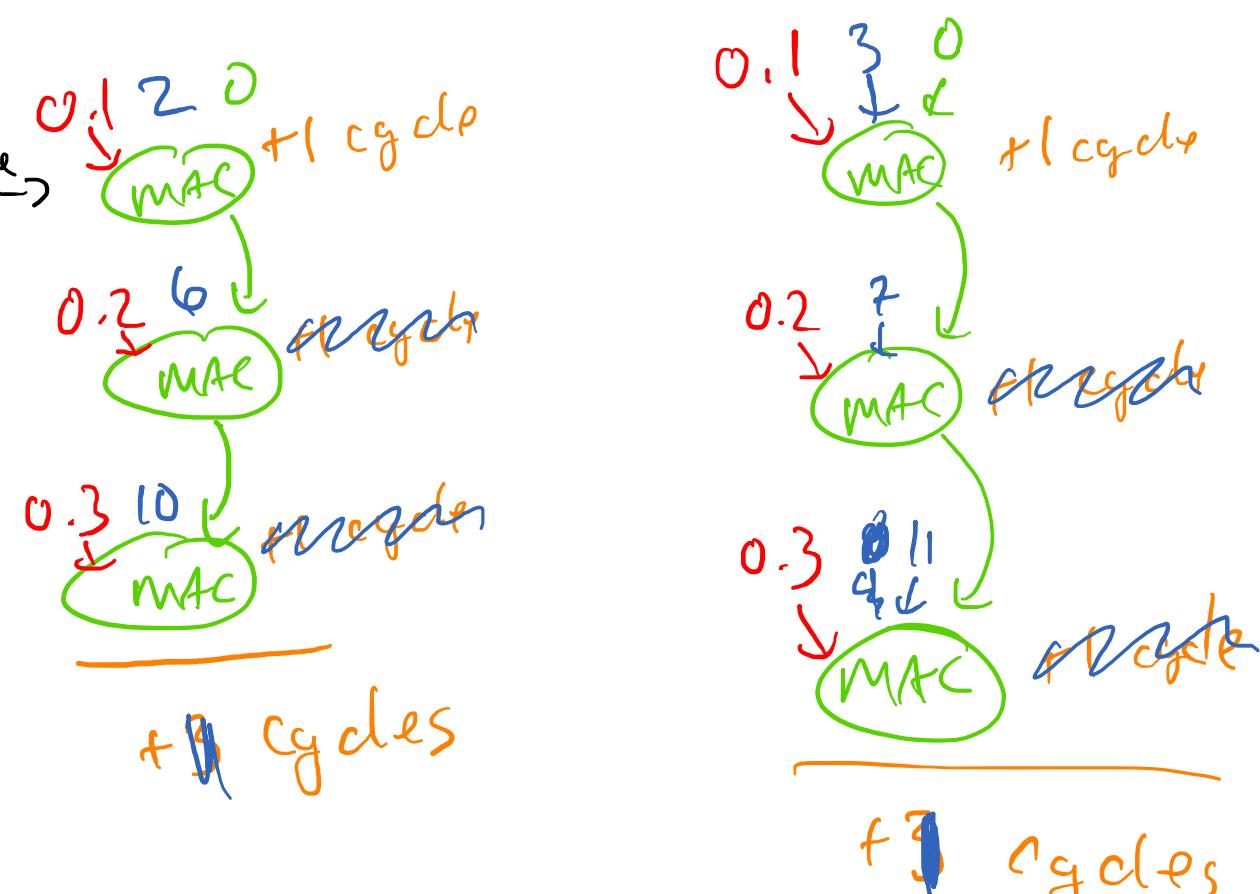
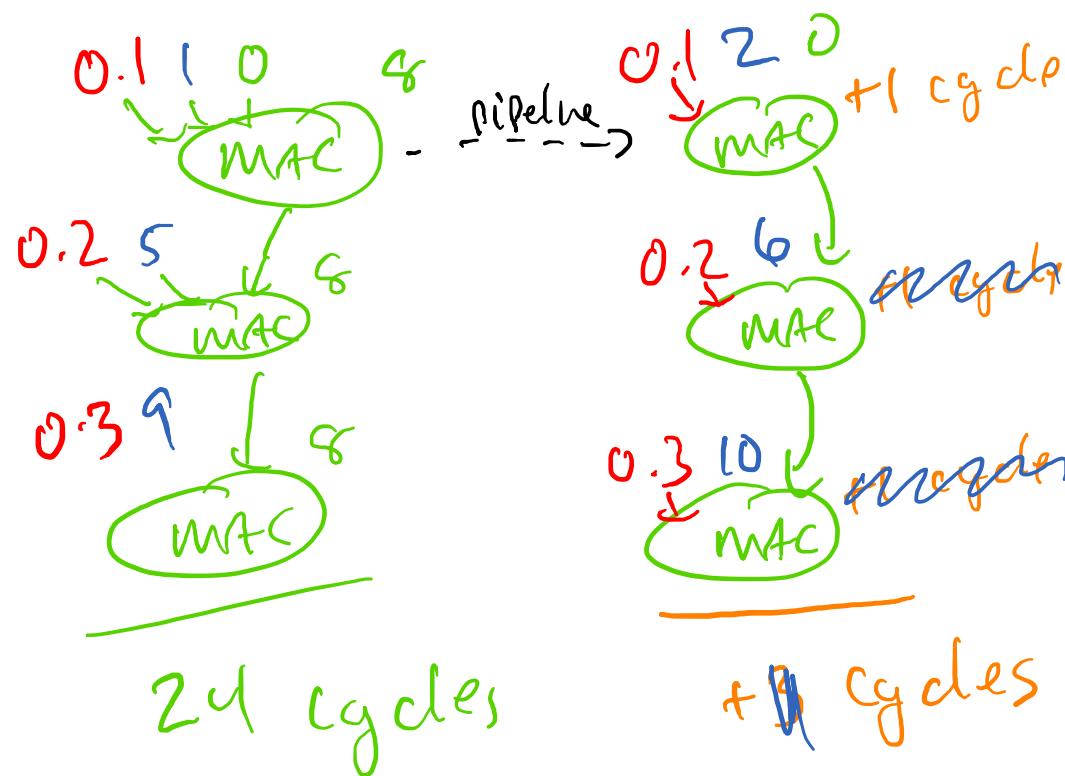
$$5.6 \leftarrow 24$$

96 cycles

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

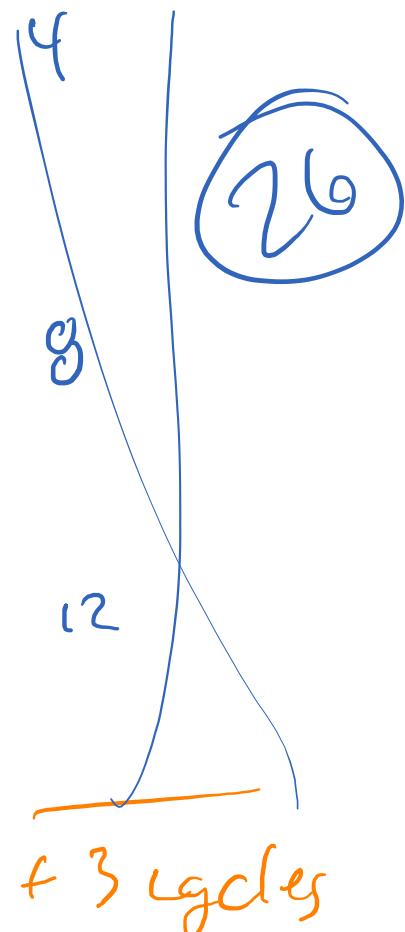


$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{bmatrix} = \begin{bmatrix} 3.8 & 4.4 & 5 & 5.6 \end{bmatrix}$$



$$24 + 3 + 3 = 30 \text{ cycles}$$

~~MAC = 27 cycles~~



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

time

$$\begin{array}{c}
 0.1 \cdot 1 + 0 = 0.1 \\
 \hline
 0.2 \cdot 5 + 0.1 = 1.1 \\
 \hline
 0.3 \cdot 9 + 1.1 = 2.4
 \end{array}$$

$$\begin{array}{c}
 0.1 \cdot 2 + 0 = 0.2 \\
 \hline
 0.2 \cdot 6 + 0.2 = 1.4 \\
 \hline
 0.3 \cdot 10 + 1.4 = 2.5
 \end{array}$$

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

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

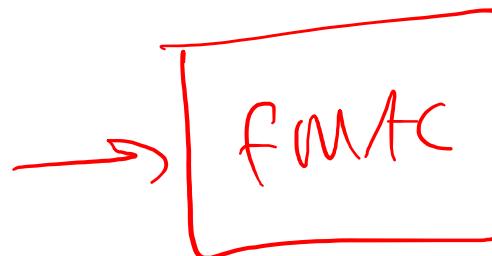
- CPU: 1 Floating-Point Unit
- FPGA? *(10 Floating-Point Units?
20?
100?)*

Finding Parallelism

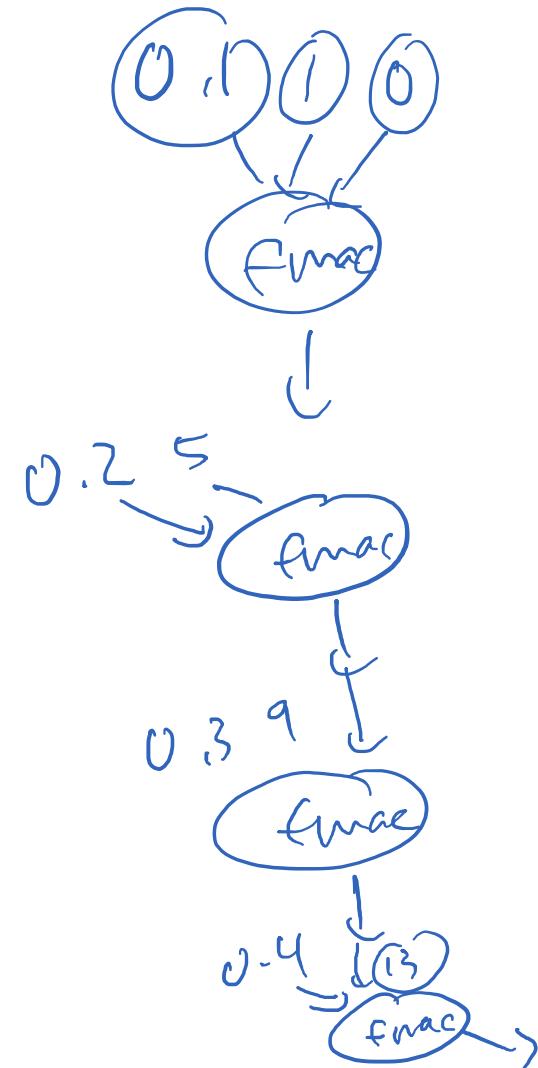
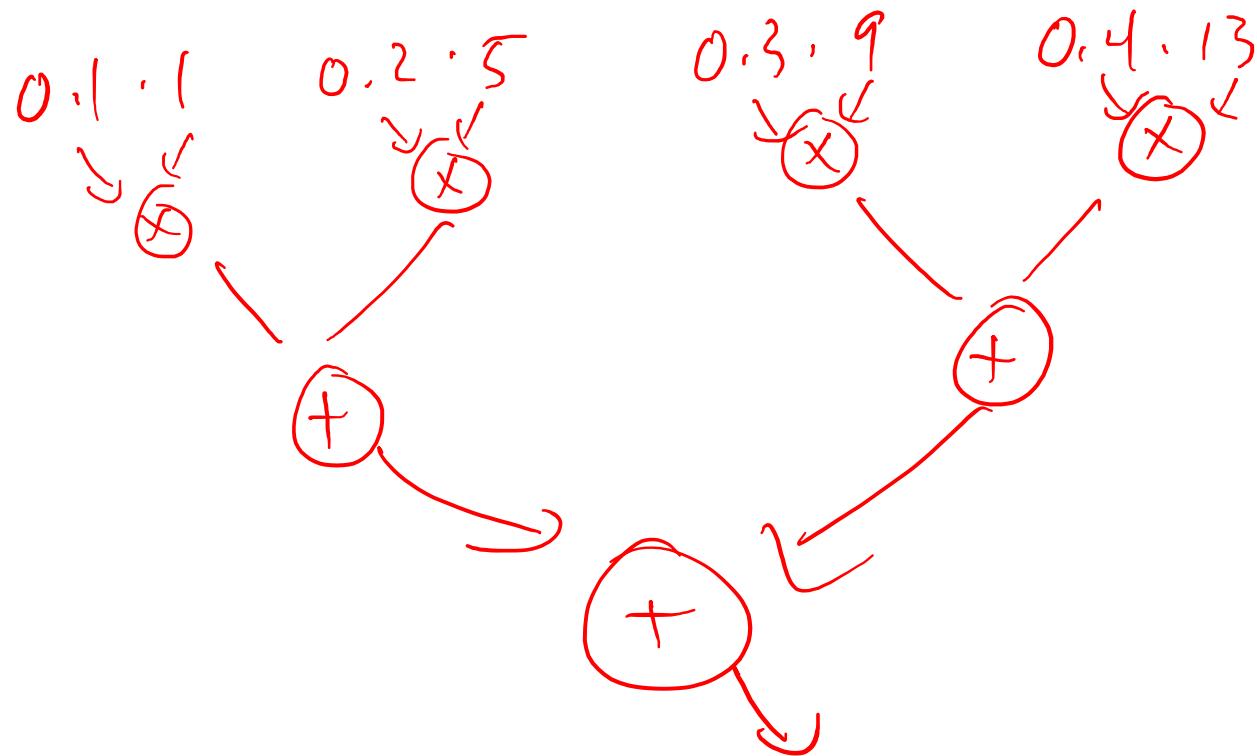
- Some computation that doesn't depend on other computation's results
- Shared Inputs are OK.

What can we do with 2+ FMACs?

dot



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



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

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

$$\begin{array}{r}
 1.1 \\
 + 7.9 \\
 \hline
 9
 \end{array}
 \quad
 \begin{array}{r}
 1.4 \\
 + 8.6 \\
 \hline
 10
 \end{array}
 \quad
 \begin{array}{r}
 1.7 \\
 + 9.3 \\
 \hline
 11
 \end{array}
 \quad
 \begin{array}{r}
 2.0 \\
 + 10 \\
 \hline
 12
 \end{array}$$

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

$$\begin{bmatrix} 0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix} \rightarrow \boxed{\text{dot1}} + \boxed{\text{dot2}}$$

$$\begin{bmatrix} 0.3 & 0.4 \end{bmatrix} \begin{bmatrix} 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix} \rightarrow \boxed{\text{dot1}}$$

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

$$[0.1 \ 0.3] \begin{bmatrix} 1 & 2 & 3 & 4 \\ 9 & 10 & 11 & 12 \end{bmatrix} = [2.8 \ 3.2 \ 3.6 \ 4.0]$$

$$[0.2 \ 0.4] \begin{bmatrix} 5 & 6 & 7 & 8 \\ 13 & 14 & 15 & 16 \end{bmatrix} = \begin{bmatrix} 6.2 & 6.8 & 7.4 & 8 \end{bmatrix}$$

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

+

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
```

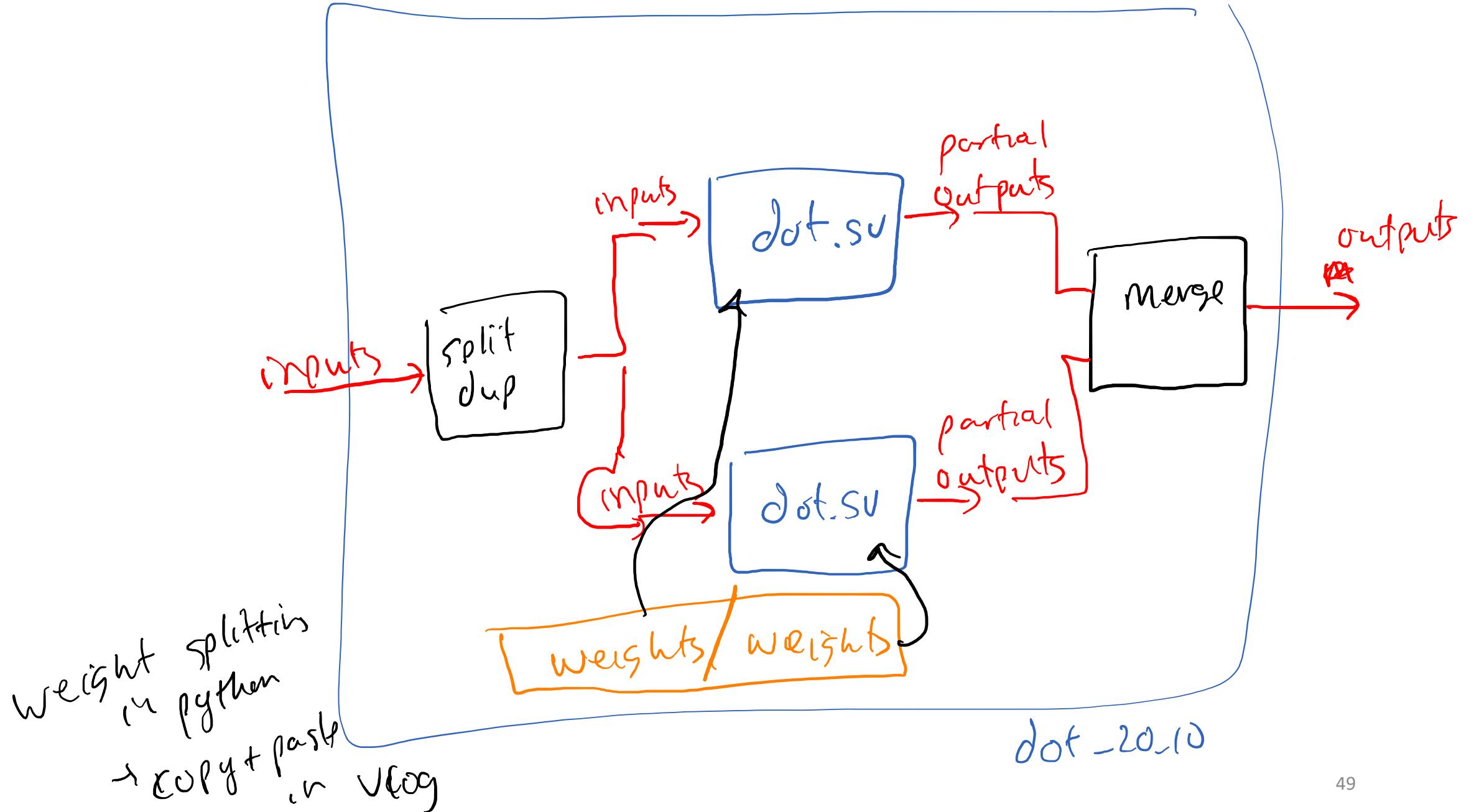
Can we parallelize Dot?

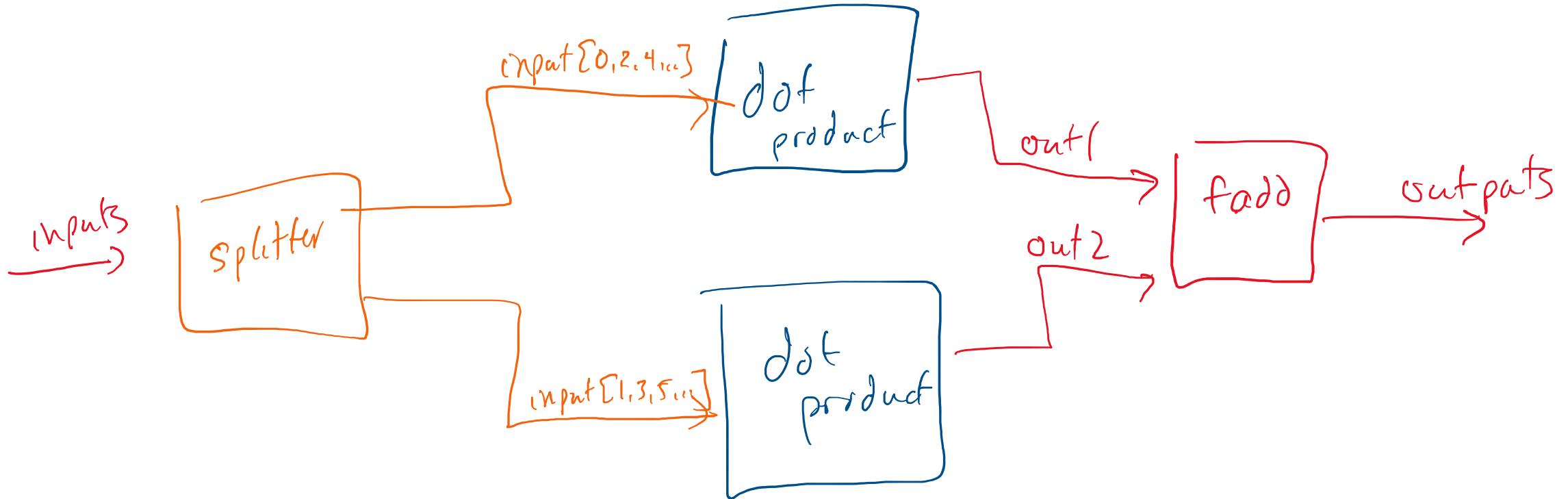
```
# how its done in dot.sv
def pydot(inputs,weights):
    inputs = inputs[0] # remove outer nesting
    outs = np.zeros(weights.shape[1], dtype=np.float32)
    for i in range(weights.shape[0]): # input length
        for j in range(weights.shape[1]): # output length
            outs[j] = outs[j] +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





17: Hardware Acceleration IV

Engr 315: Hardware / Software Codesign

Andrew Lukefahr

Indiana University

