

19: Parallelism II

ENGR 315: Hardware/Software CoDesign

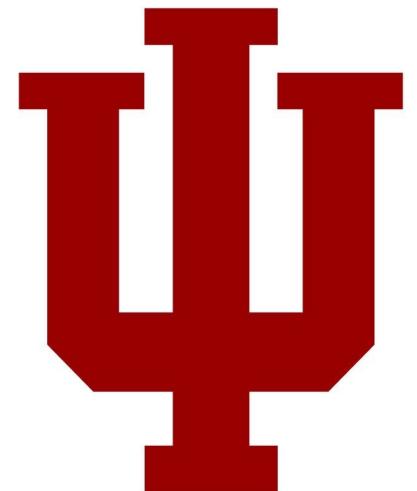
Andrew Lukefahr

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/>



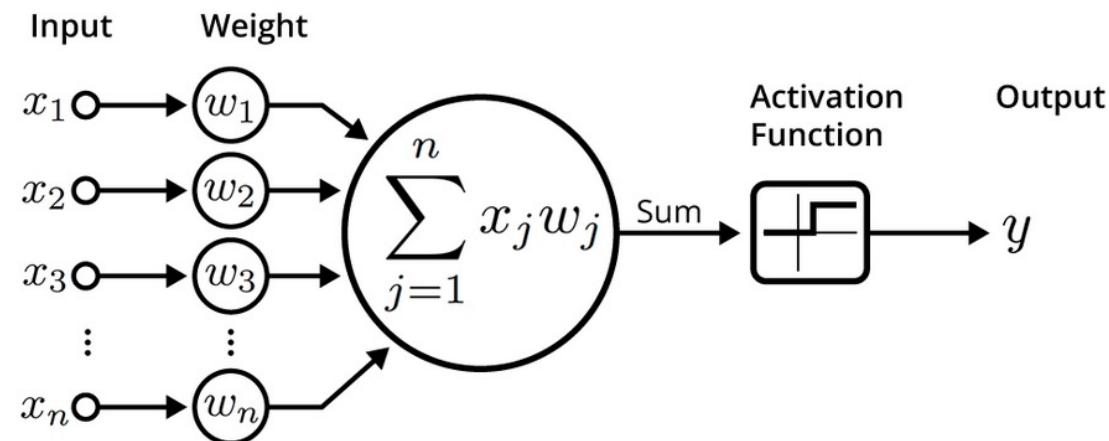
Announcements

- P8: Pipelining
- P9: Parallelism

P8+ are part of the “Design Challenge”

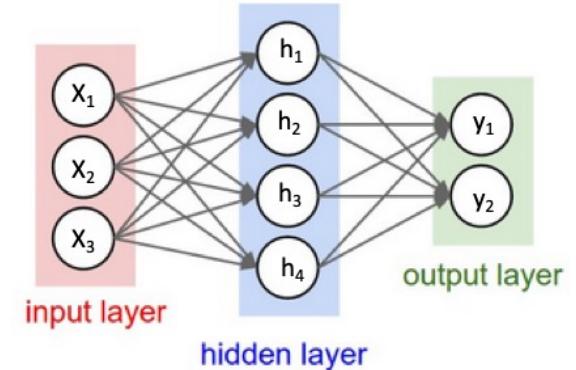
- Goal: Accelerate reference neural network
- Harder, more open-ended projects

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)

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)

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

Multiply-Accumulate Dot Computations

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

round 1

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

round 2

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

Multiply-Accumulate Dot Computations

Step 0:

$$\begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.2 & 1.5 \end{bmatrix}$$
$$= \begin{bmatrix} 0^{\sim} & 0^{\sim} & 0^{\sim} \end{bmatrix}$$
$$\begin{array}{c} 0.1 \cdot 1 + 0 \\ 0.1 \cdot 2 + 0 \\ 0.1 \cdot 3 + 0 \end{array} \rightarrow \boxed{[0.1]} \quad \begin{array}{c} 0.2 \cdot 4 + 0 \\ 0.2 \cdot 5 + 0.2 \\ 0.2 \cdot 6 + 0.3 \end{array} \rightarrow \boxed{[0.9 \quad 1.2 \quad 1.5]}$$

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

[0.1 0.2 0.3 0.4]

$$0.1 \cdot 1 + 0$$

$$0.1 \cdot 2 + 0$$

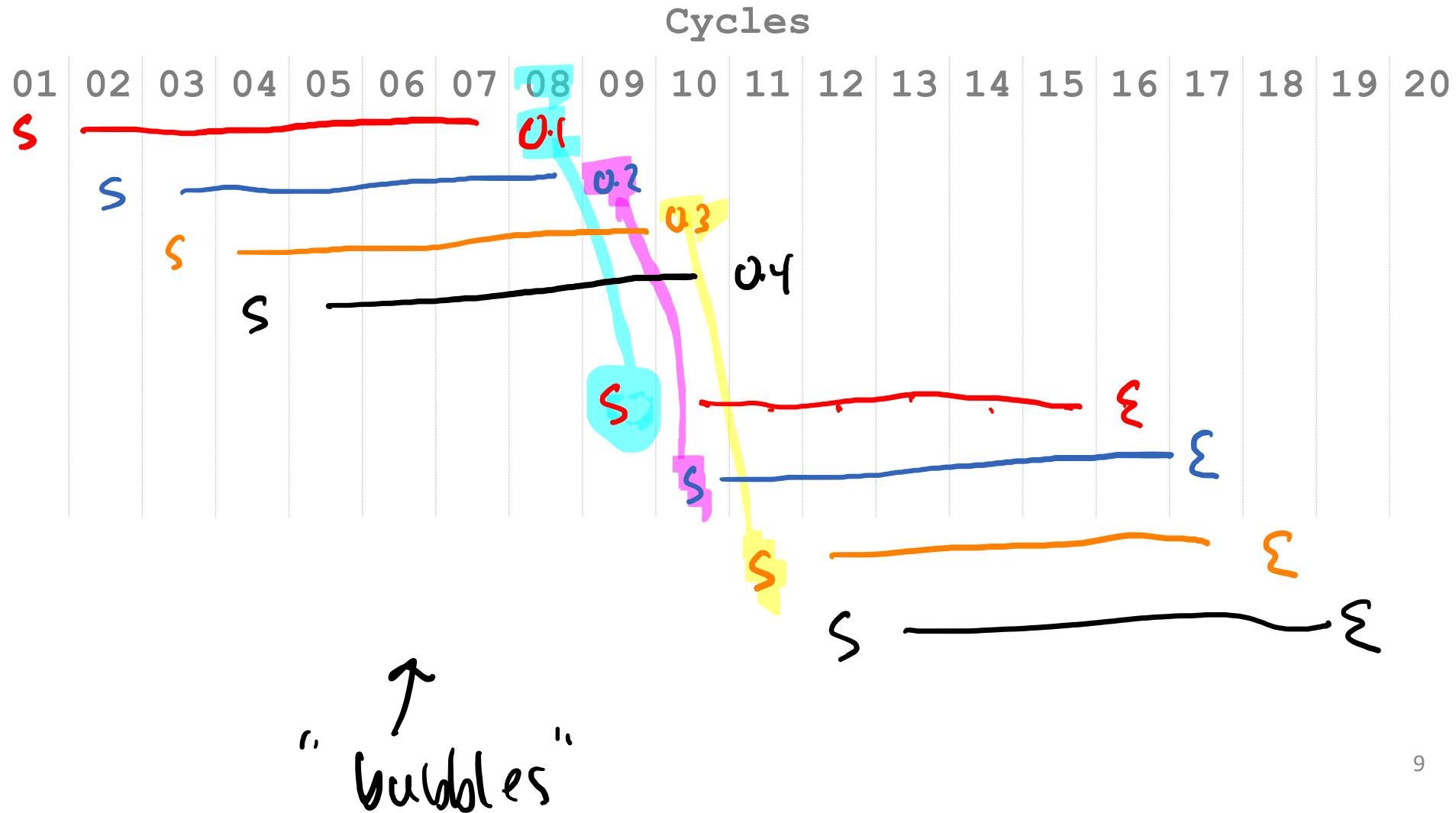
$$0.1 \cdot 3 + 0$$

$$0.1 \cdot 4 + 0$$

$$0.2 \cdot 5 + 0.1$$

$$0.2 \cdot 6 + 0.2$$

$$0.2 \cdot 7 + 0.3$$



Latency on Pipelined FMAC

- Solution: Stall at the end of a row.
 - “Drain” the pipeline.
-
- Restart new row
 - “Refill” the pipeline

P8: Dot Organization

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,
    input [BITS-1:0] b,
    output [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,  
    input [BITS-1:0] b,  
    output [BITS-1:0] c) );  
    assign c = a + b;  
endmodule
```

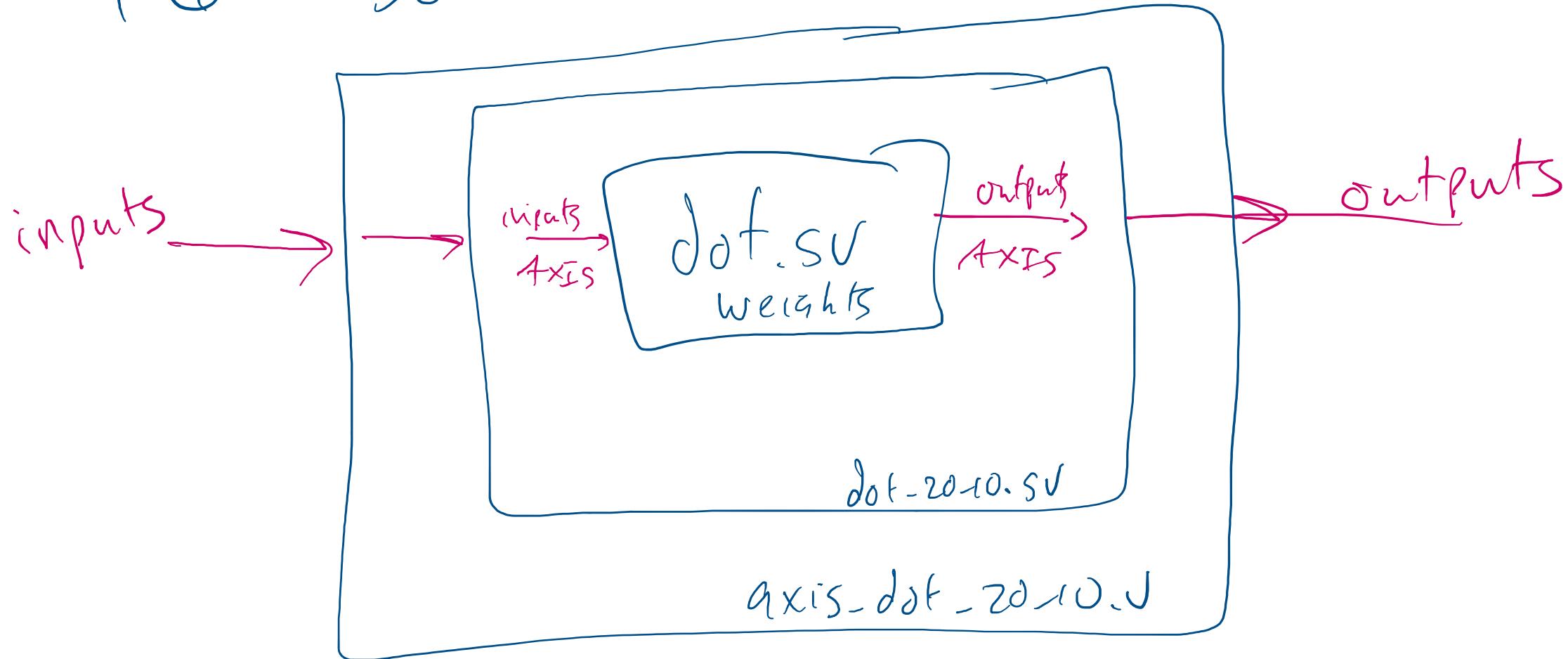
Verilog Parameters in Dot_20_10.sv

```
// This is autogenerated. see python/dot_20_10.py for details.  
localparam ROWS = 20;  
localparam COLS = 10;  
  
localparam [31:0] weights [0:ROWS-1] [0:COLS-1] = '{  
    '{...}... };  
  
dot #(  
    .ROWS(ROWS),  
    .COLS(COLS)  
) dot0(  
    // AXI4-Stream Interface  
    .clk(clk),  
    .rst(rst),  
  
    .weights(weights),
```

P6

Dot

Data flow



Data flow graph

- 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.

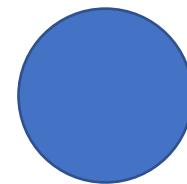
[\[ref\]](#)

Data Flow Graph

- Illustrates the dependencies between inputs and operations to produce an output



Inputs



Operations

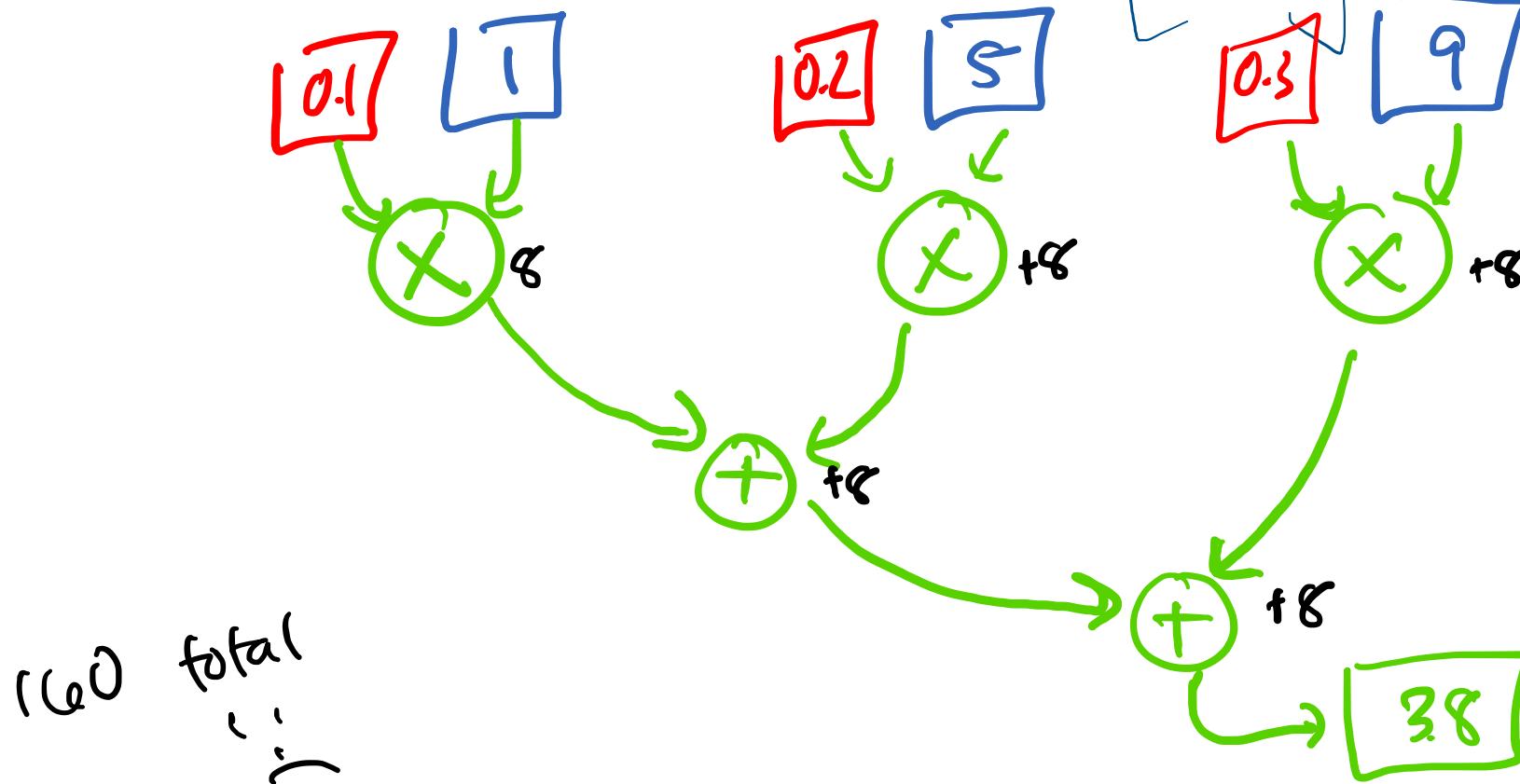


Dependencies

Data Flow Graph

+ NC
+ NJ
+ MAC for Now
Pipeline

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} = \begin{bmatrix} 3.8 \end{bmatrix}$$



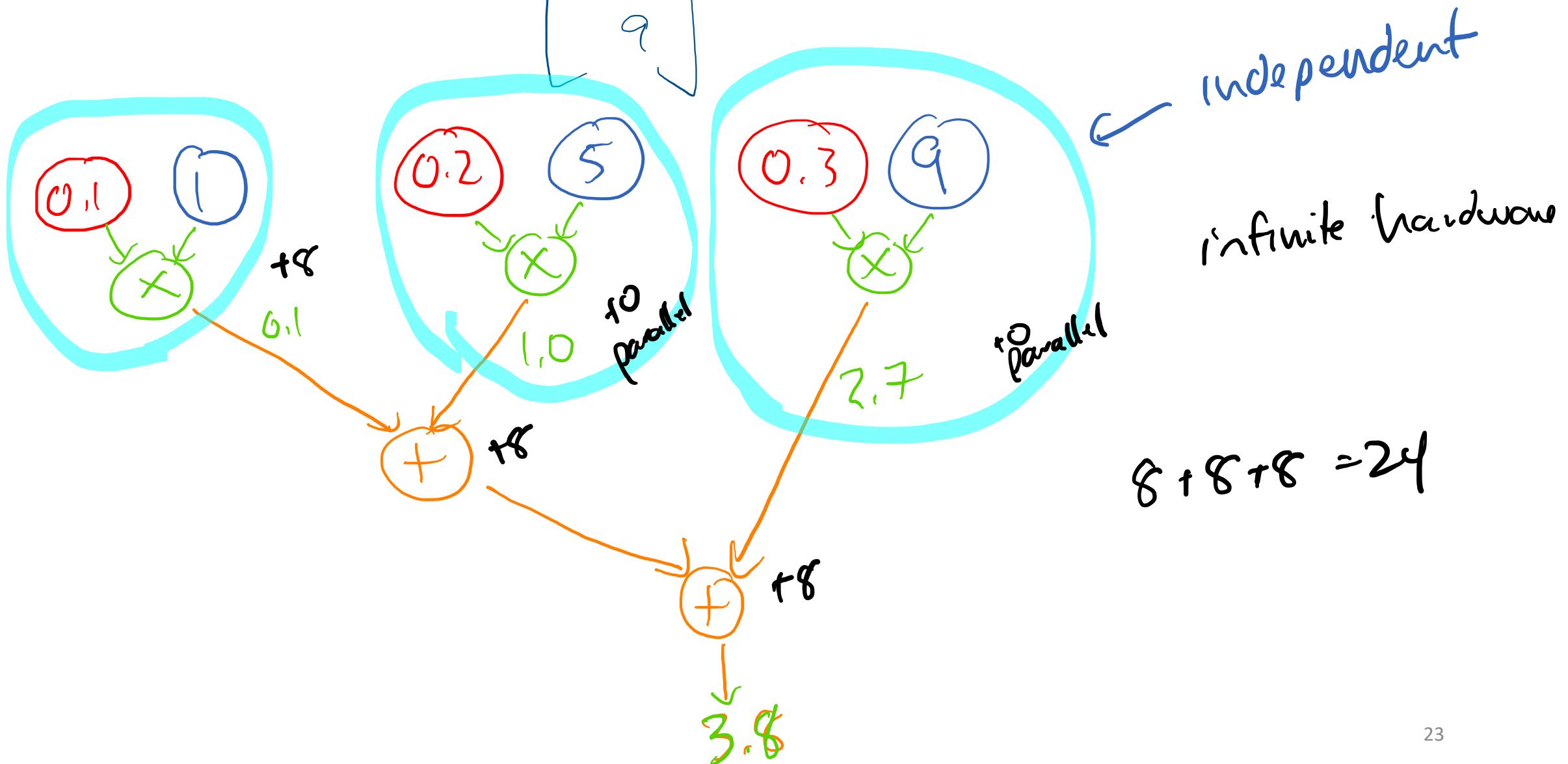
MULT: 8 cycles
ADD: 8 cycles

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

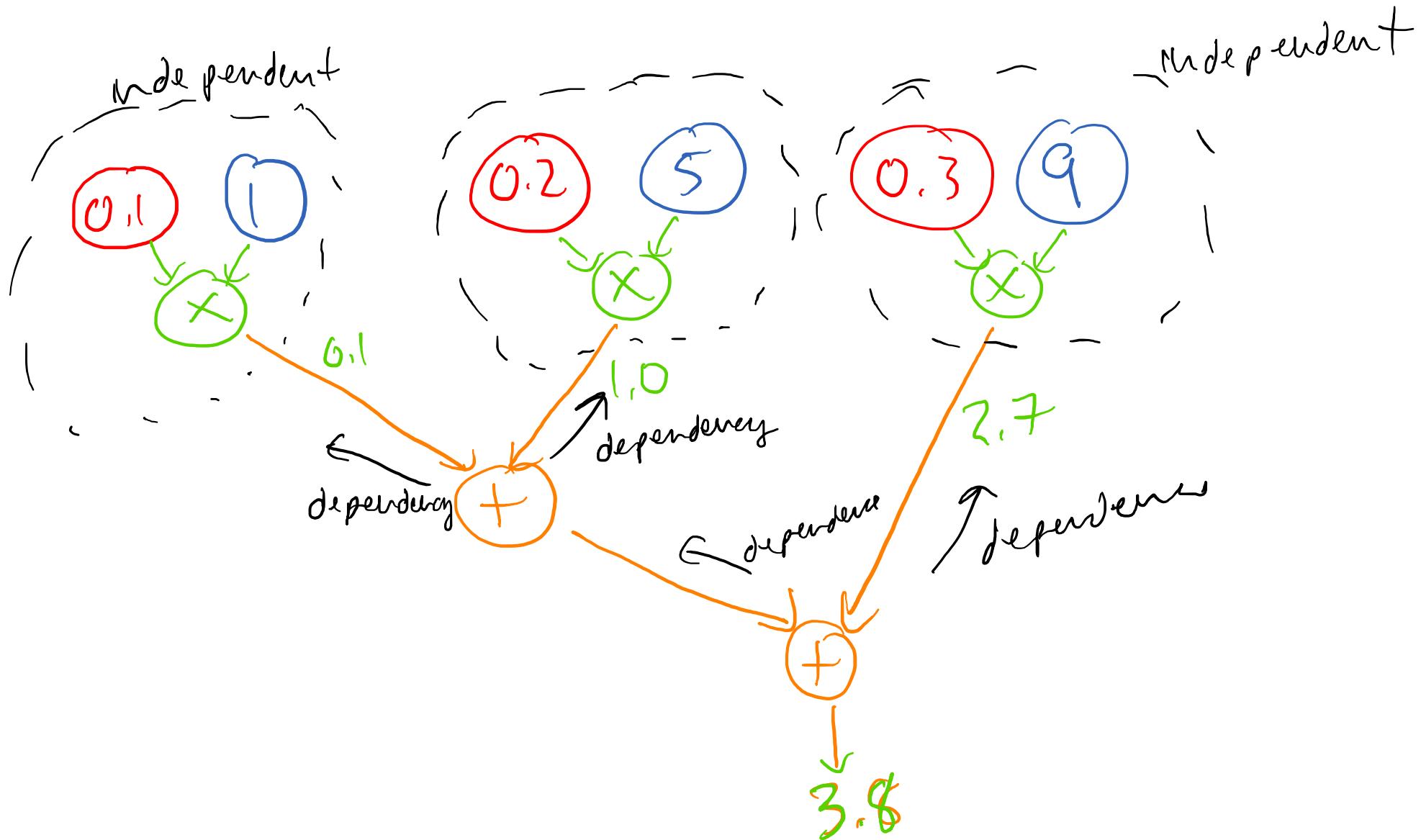
Finding Parallelism

- Find some computation that doesn't depend on other computation's results.
- No arrows between blocks.
- Shared Inputs are OK.

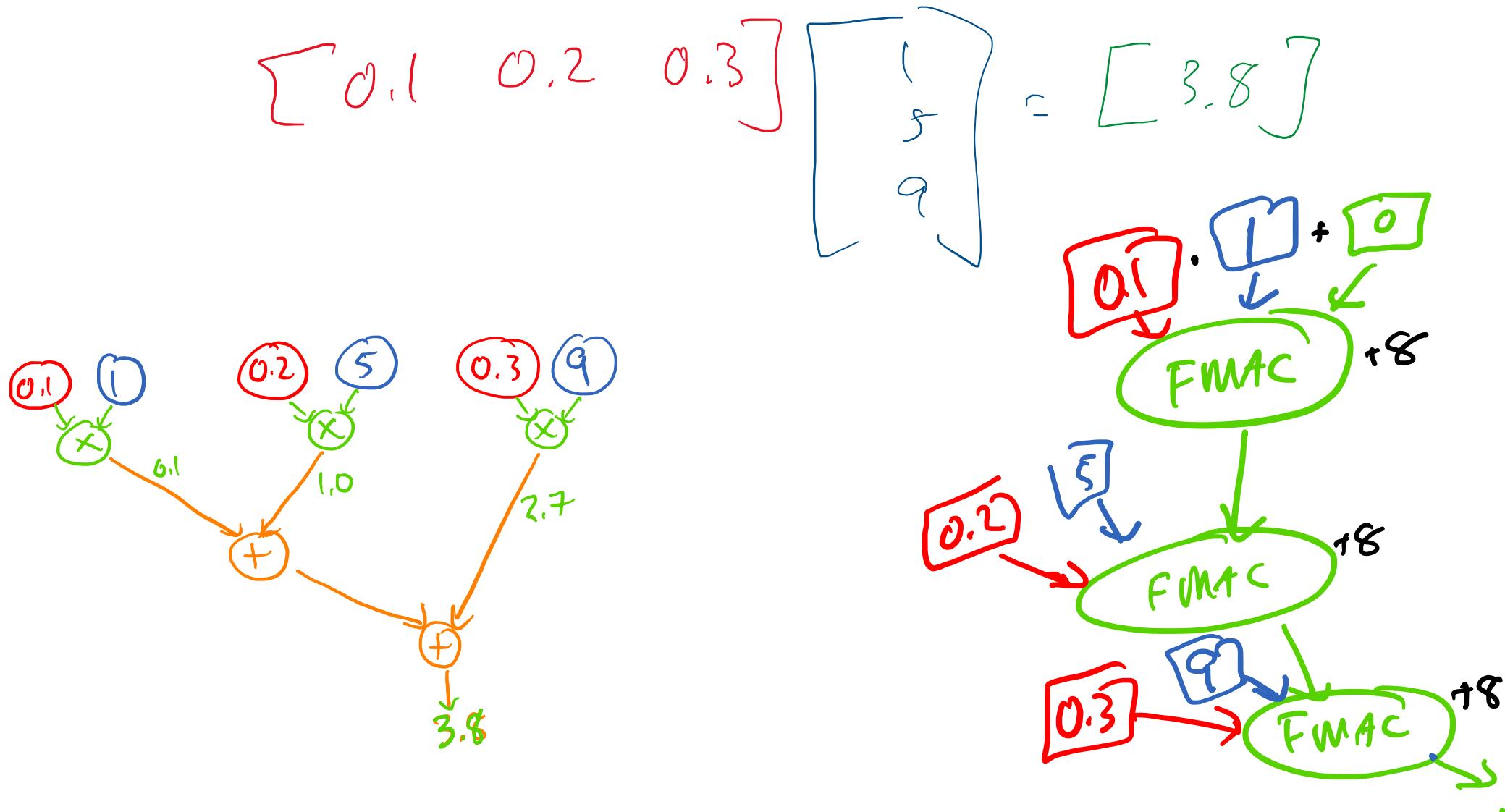
$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 5 \\ 9 \end{bmatrix} = \begin{bmatrix} 3.8 \end{bmatrix}$$



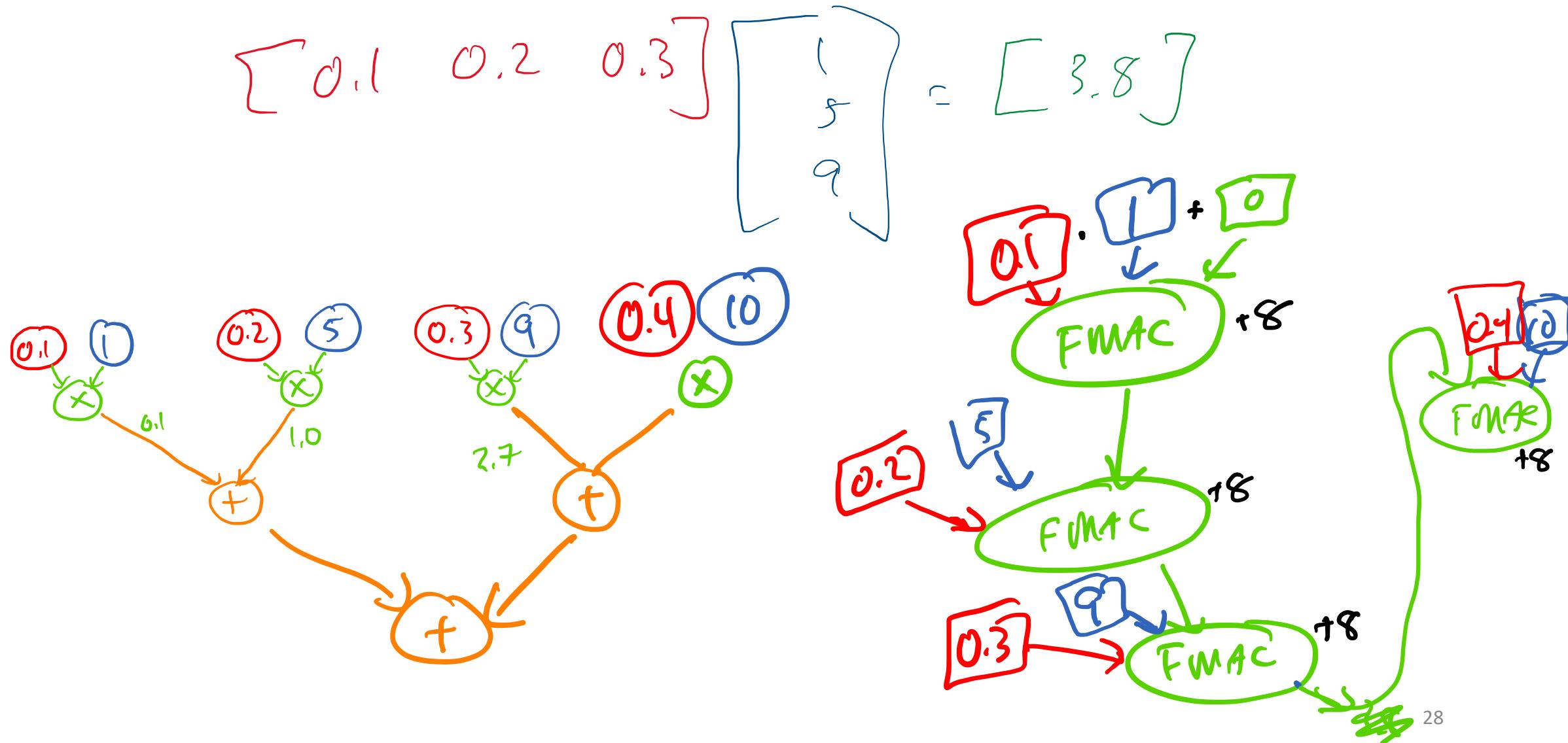
How many cycles should this take?



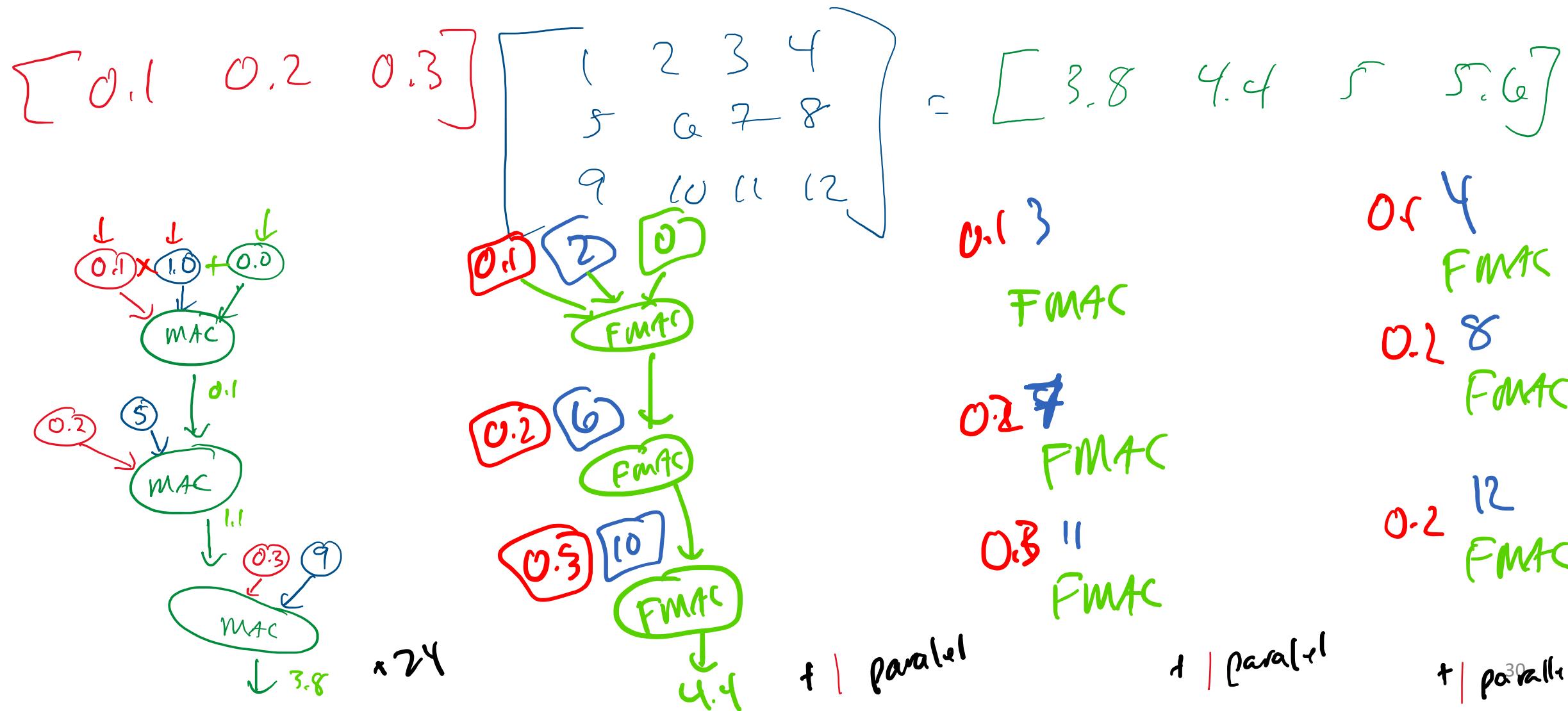
How many cycles should this take?

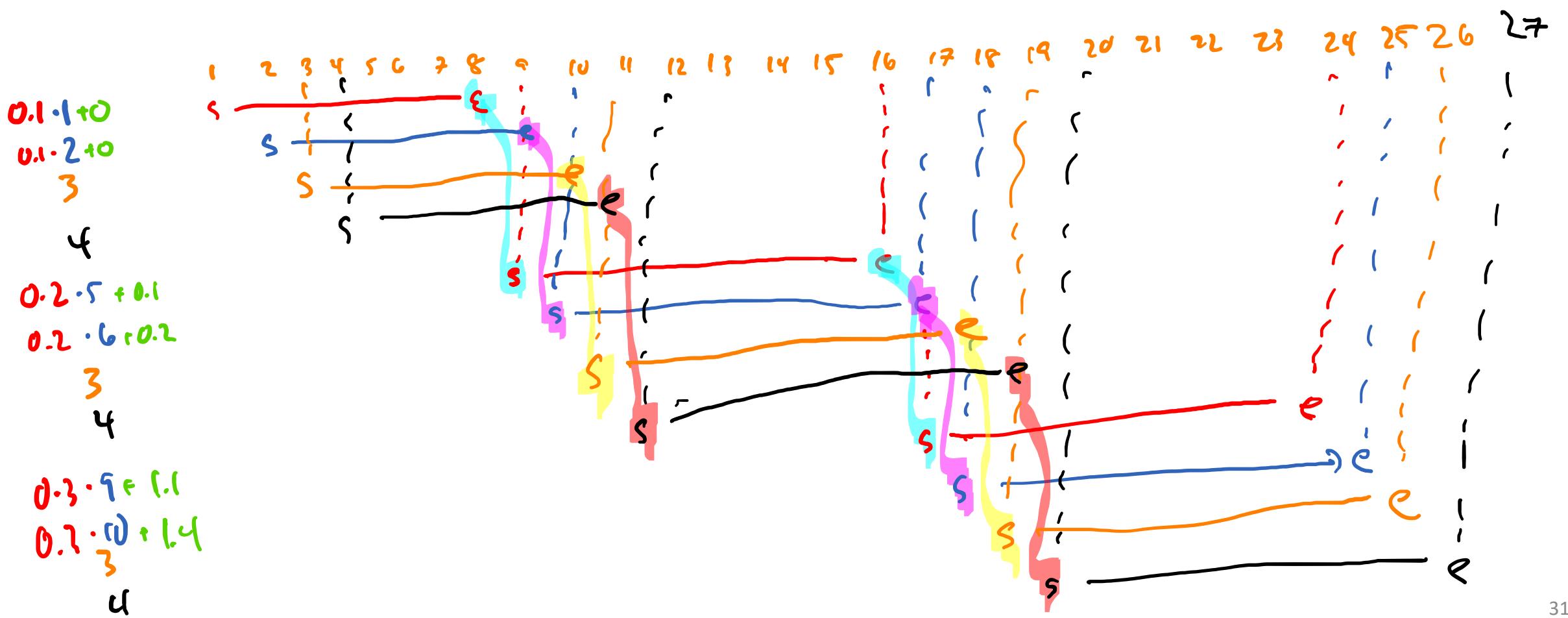


How many cycles should this take?



More Data Flow





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

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

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

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

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

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

One Output

Multiply & Add separate:

24

24

24

+1 extra column

24

Pipelined MAC

Maximum Parallel:

All 4 Outputs

Serial

96

96

Parallel

24

24

$$24 + 1 + 1 + 1 = 27$$

24

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

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

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

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

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

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

Next Tim: Working out parallelism