



**Samueli**  
Computer Science



# CS145 Discussion: Week 6 (Add-On) A Story of Computing

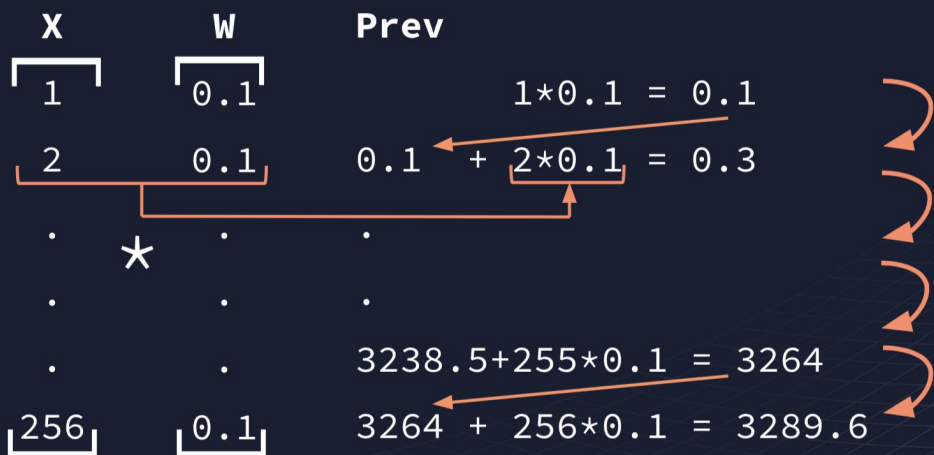
Junheng Hao  
Friday, 11/13/2020

## Single-threaded Execution

$X = [1.0, 2.0, \dots, 256.0]$  # Let's say we have 256 input values

$W = [0.1, 0.1, \dots, 0.1]$  # Then we need to have 256 weight values

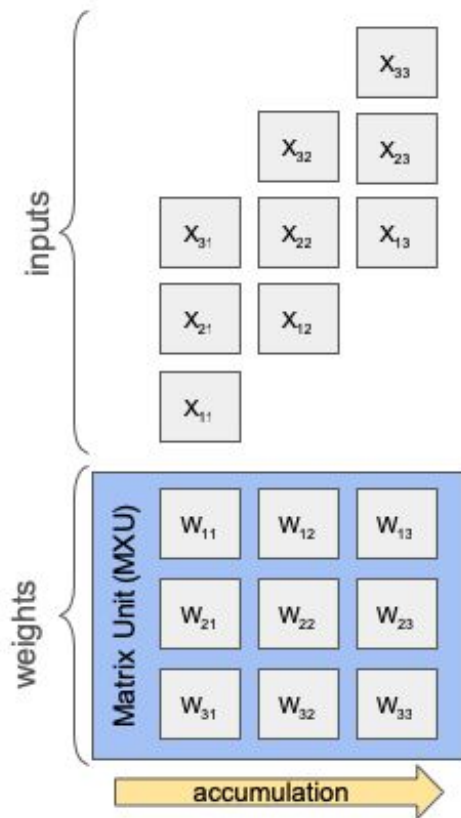
$h_{0,0} = X * W$  #  $[1*0.1 + 2*0.1 + \dots + 256*0.1] == 32389.6$



Single-threaded  
Execution

$256 * \blacktriangle t$

# UCLA Neural Networks: Computation Example



## Matrix Unit Systolic Array

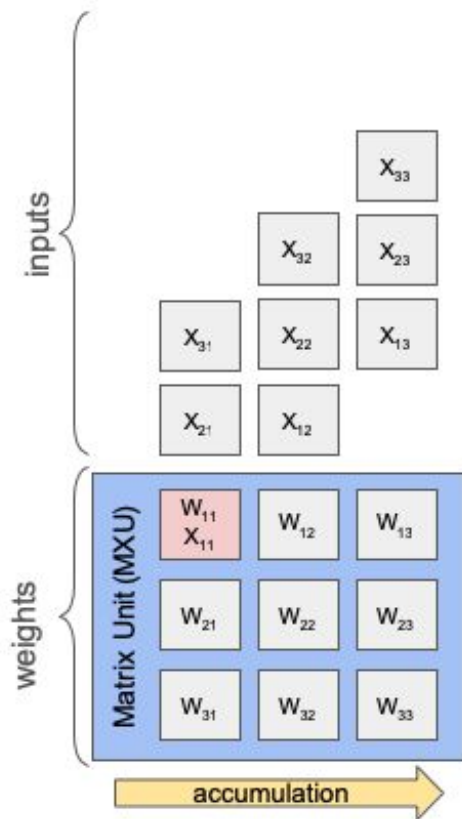
Computing  $y = Wx$

3x3 systolic array

$W = 3 \times 3$  matrix

Batch-size( $x$ ) = 3

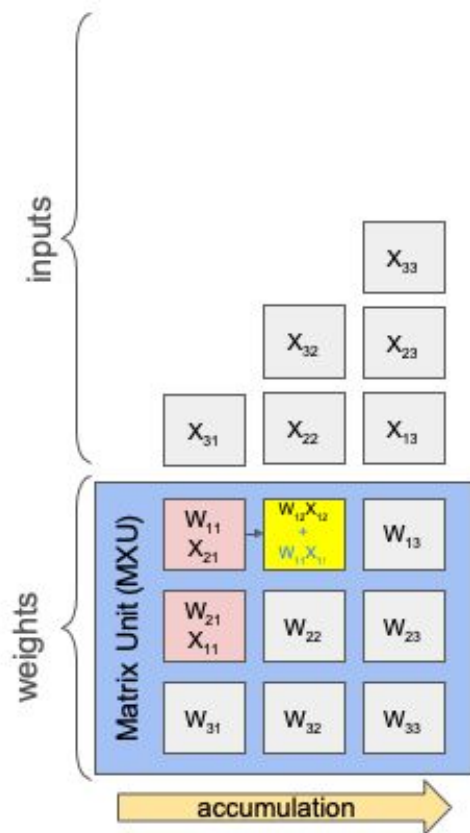
# UCLA Neural Networks: Computation Example



## Matrix Unit Systolic Array

Computing  $y = Wx$   
with  $W = 3 \times 3$ , batch-size( $x$ ) = 3

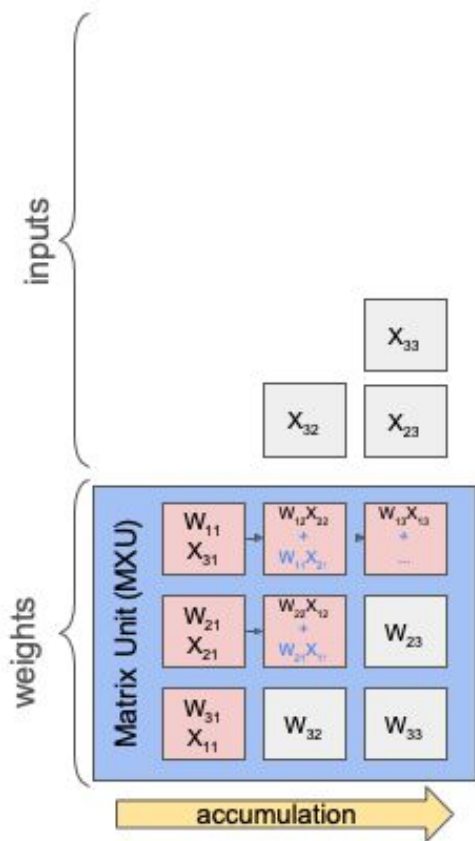
# UCLA Neural Networks: Computation Example



## Matrix Unit Systolic Array

Computing  $y = Wx$   
with  $W = 3 \times 3$ ,  $\text{batch-size}(x) = 3$

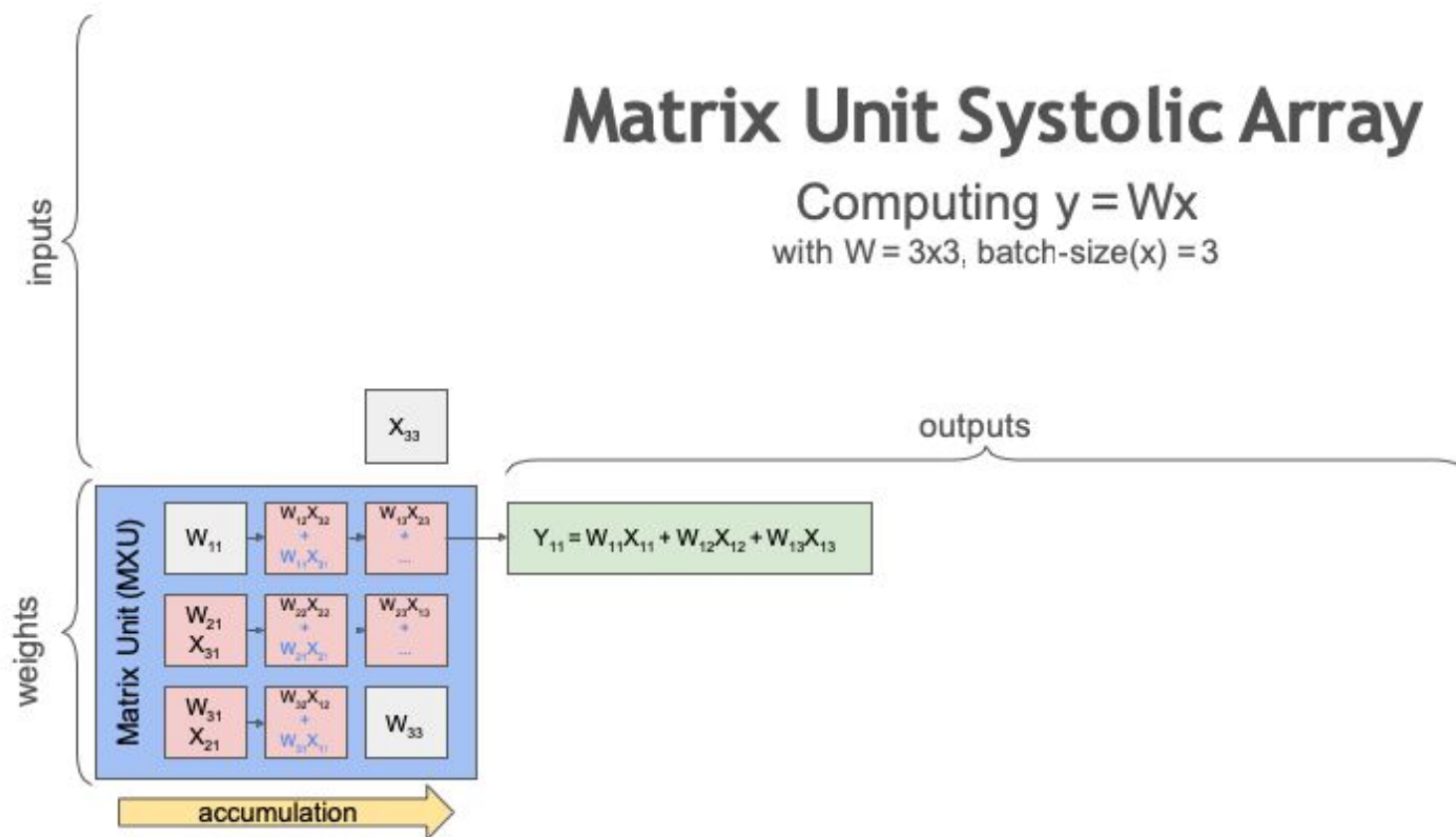
# UCLA Neural Networks: Computation Example



## Matrix Unit Systolic Array

Computing  $y = Wx$   
with  $W = 3 \times 3$ , batch-size( $x$ ) = 3

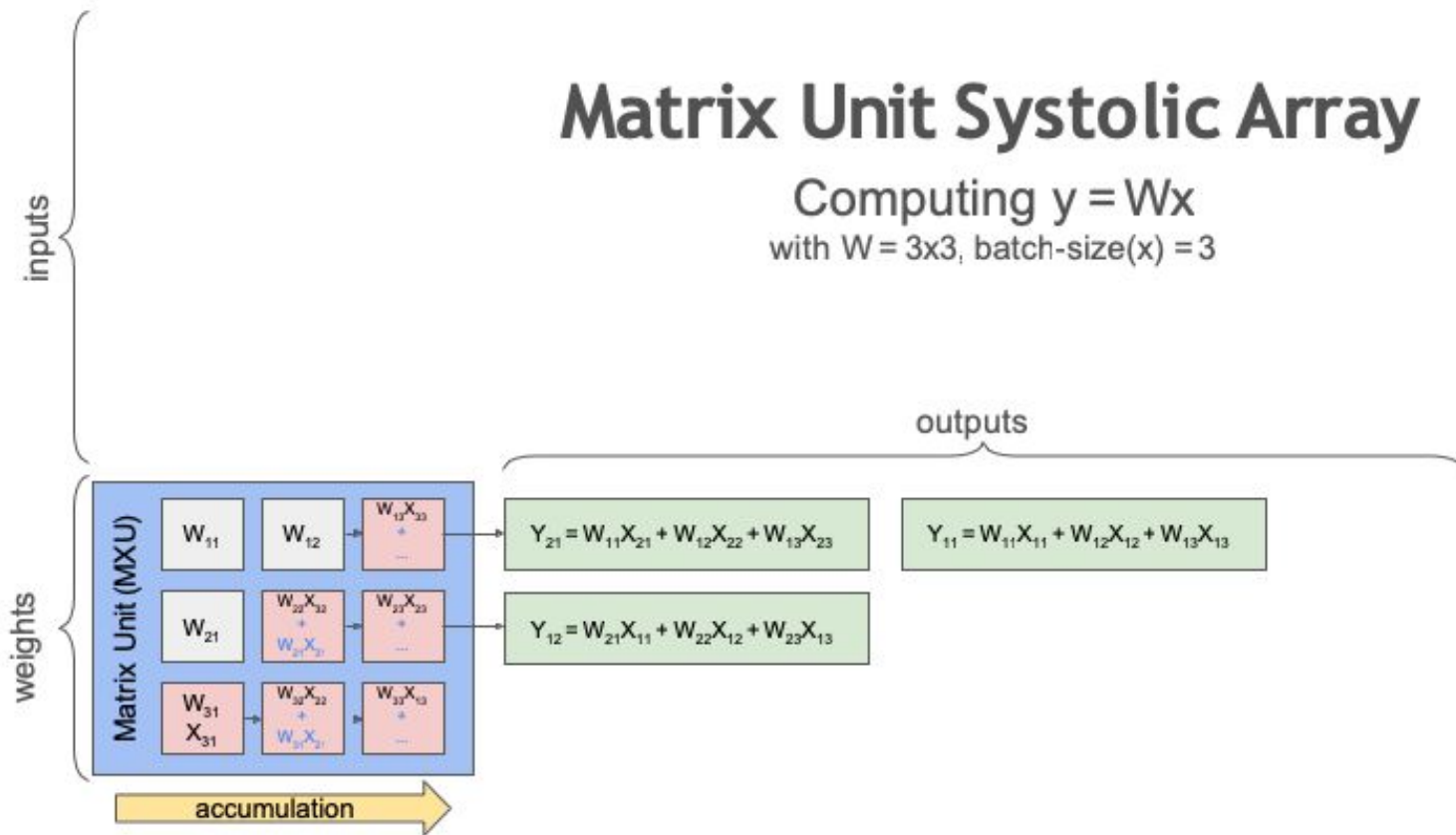
# UCLA Neural Networks: Computation Example



## Matrix Unit Systolic Array

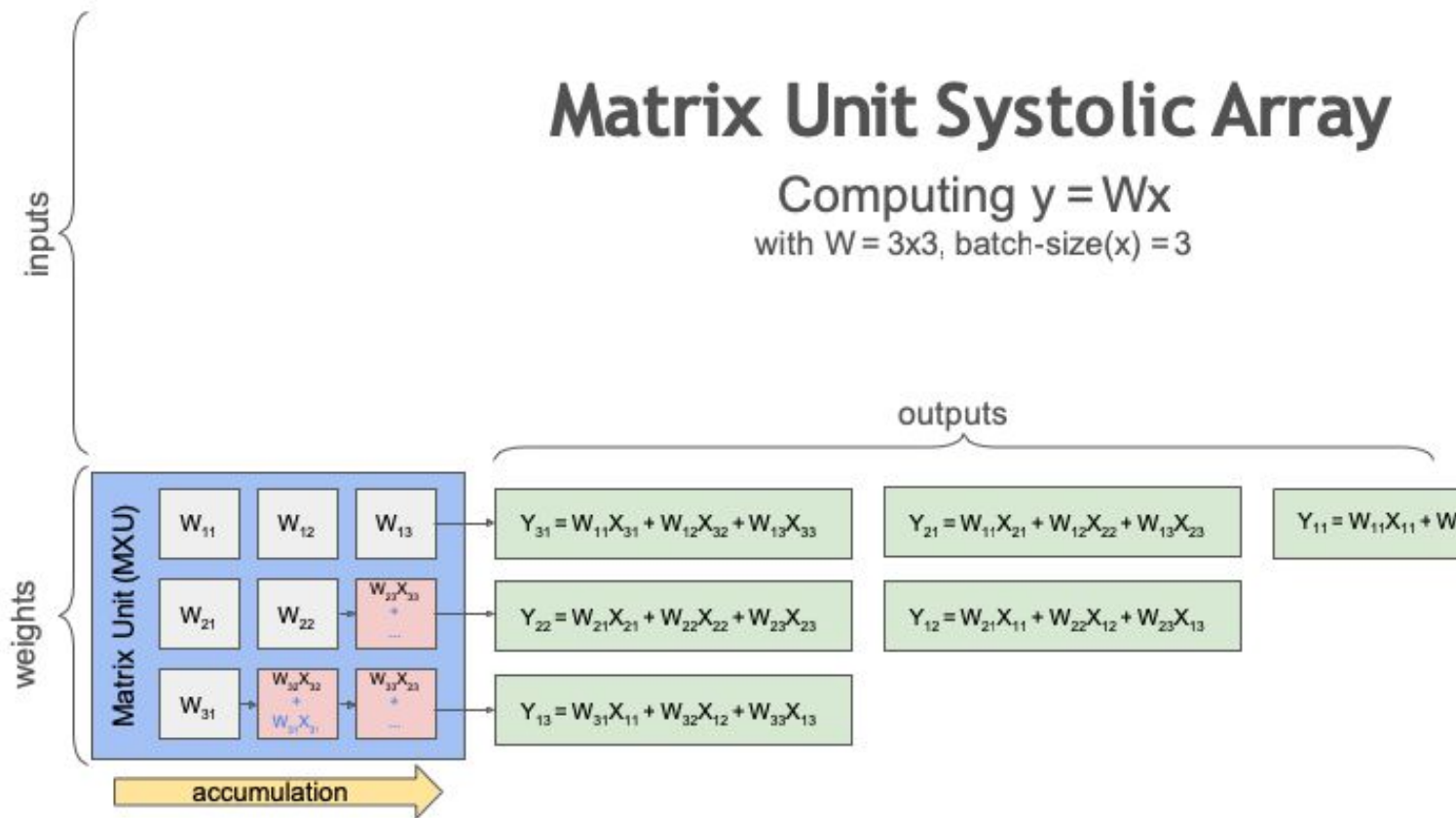
## Computing $y = Wx$

with  $W = 3 \times 3$ ,  $\text{batch-size}(x) = 3$

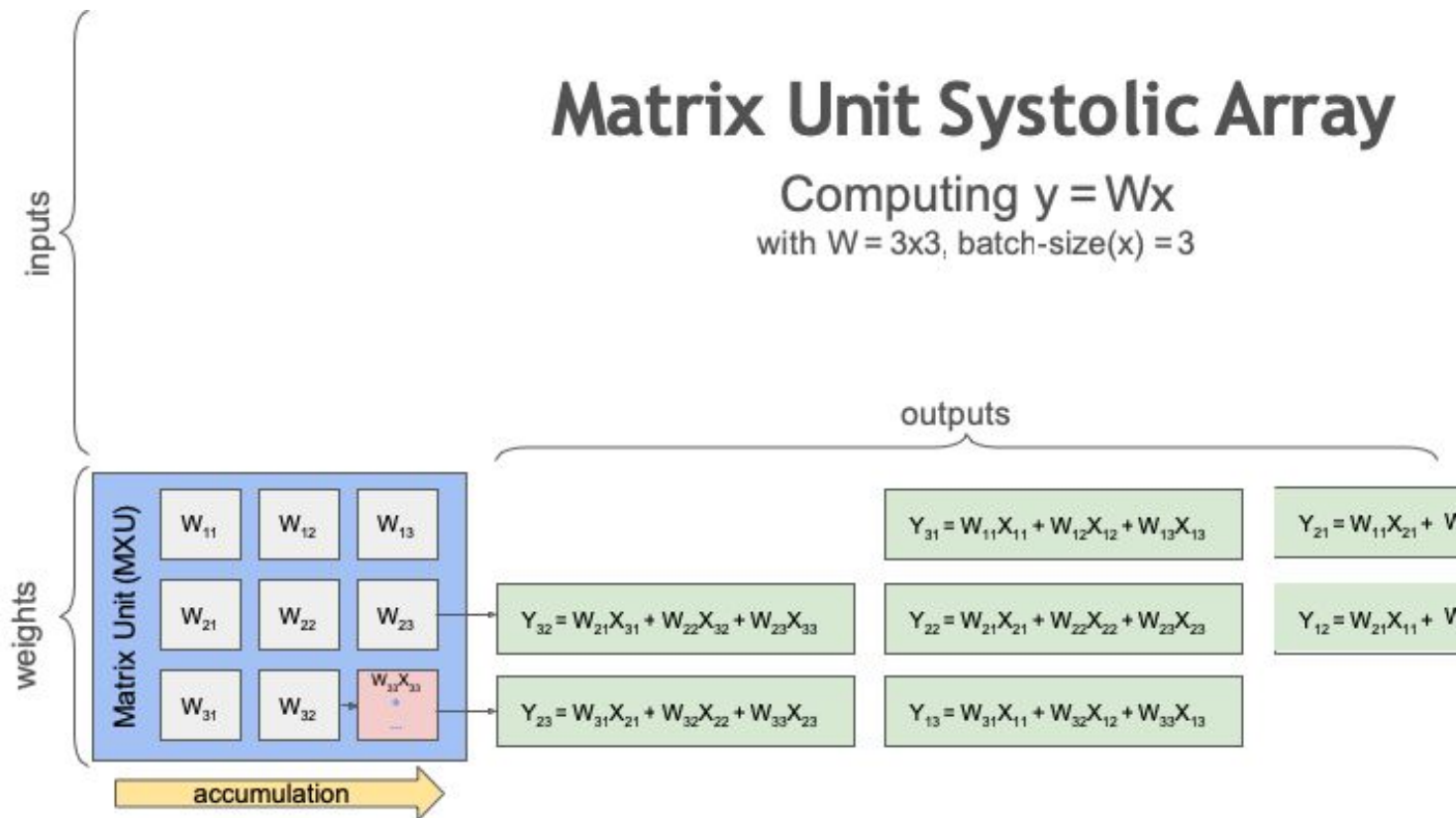




# UCLA Neural Networks: Computation Example



# UCLA Neural Networks: Computation Example



# UCLA Neural Networks: Computation Example

