

# When Dataflows Converge: Reconfigurable and Approximate Computing for Emerging Neural Networks

Di Wu and Joshua San Miguel

[di.wu@ece.wisc.edu](mailto:di.wu@ece.wisc.edu) and [jsanmiguel@wisc.edu](mailto:jsanmiguel@wisc.edu)



# Outline

- ☐ Motivation
- ☐ Background
- ☐ Architecture
- ☐ Evaluation
- ☐ Conclusion

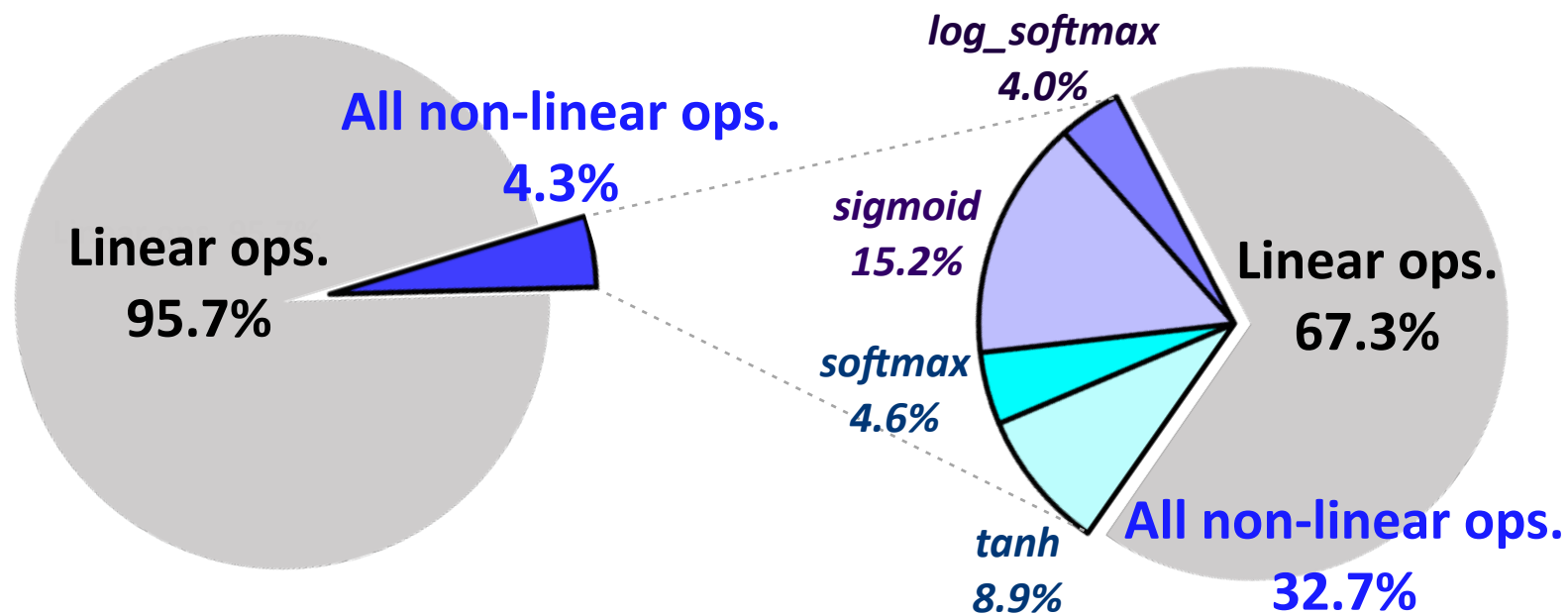
# Nonlinear operations are underrated

- Emerging deep neural networks have nonlinear operations with
  - High diversity: coexistence of multiple operations

Emerging NN	Nonlinear Operations
Capsule Neural Network (CapsNet)	<i>div, exp, log, sigmoid, softmax</i>
Graph Neural Network (GNN)	<i>div, exp, log, softmax</i>
Neural Machine Translation (NMT)	<i>log, sigmoid, softmax, tanh</i>

# Nonlinear operations are underrated

- Emerging deep neural networks have nonlinear operations with
  - High diversity
  - High complexity: beyond simple ReLU and max pooling



Operation count

Operation runtime

## Profiling NMT model on CPU

“UNO: Virtualizing and Unifying Nonlinear Operations for Emerging Neural Networks”, ISLPED’21



# Nonlinear operations are underrated

- Emerging deep neural networks have nonlinear operations with
  - High diversity
  - High complexity
  - High cost: costly naive processing element (PE)

Function	Performance	Naive PE	MAC PE	Overhead
<i>MAC, div, exp, log</i>	Area ( $\mu\text{m}^2$ )	9,323	1,559	6.0X
	Power (mW)	2.35	0.56	4.2X

# Nonlinear operations are underrated

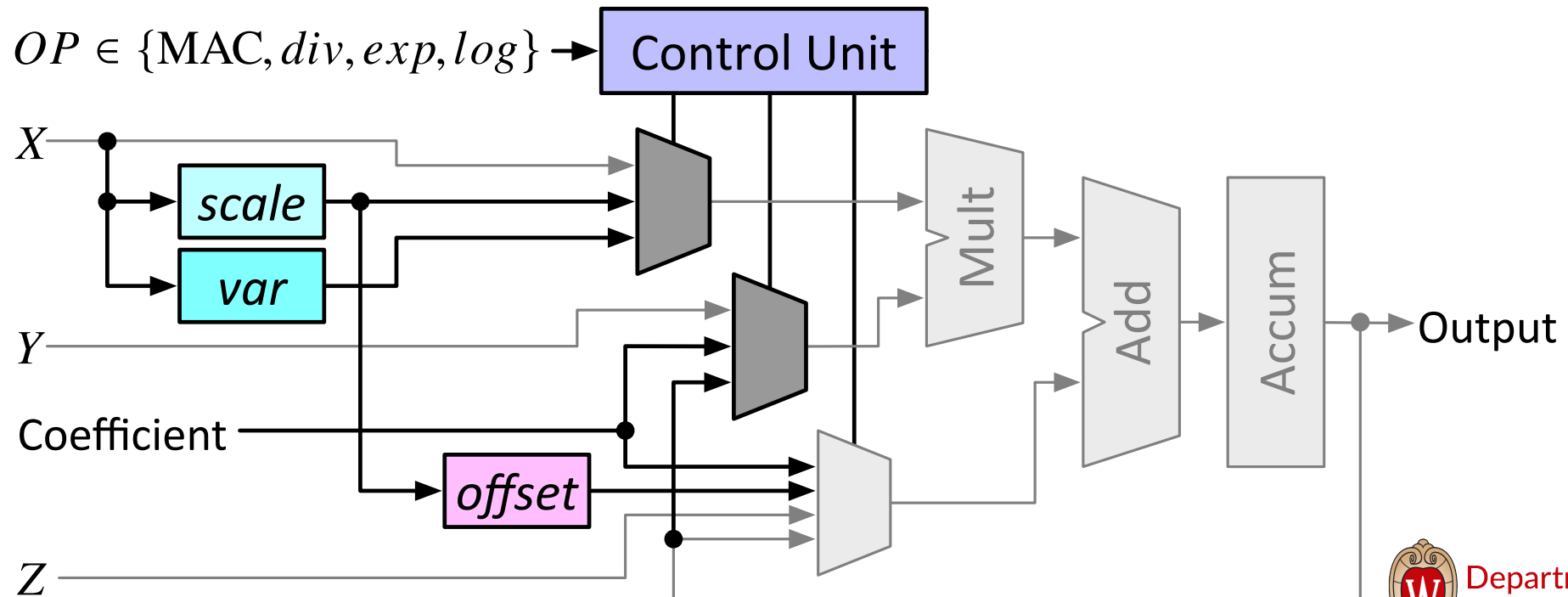
- Emerging deep neural networks have nonlinear operations with
  - High diversity
  - High complexity
  - High cost: costly naive processing element (PE)

Function	Performance	Naive PE	Optimized PE	Overhead
<i>MAC, div, exp, log</i>	Area ( $\mu\text{m}^2$ )	9,323	4,221	2.2X
	Power (mW)	2.35	0.93	2.5X

# Unify nonlinear operation (UNO)

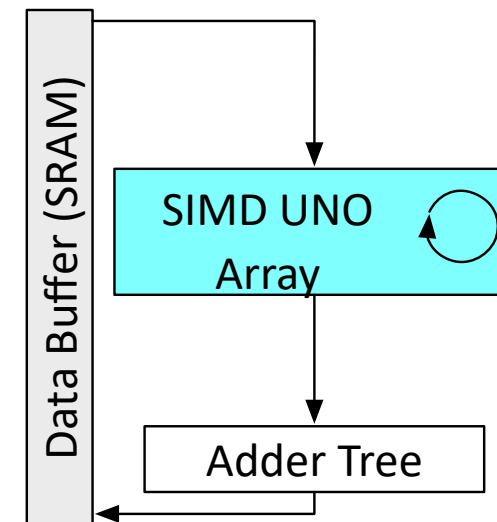
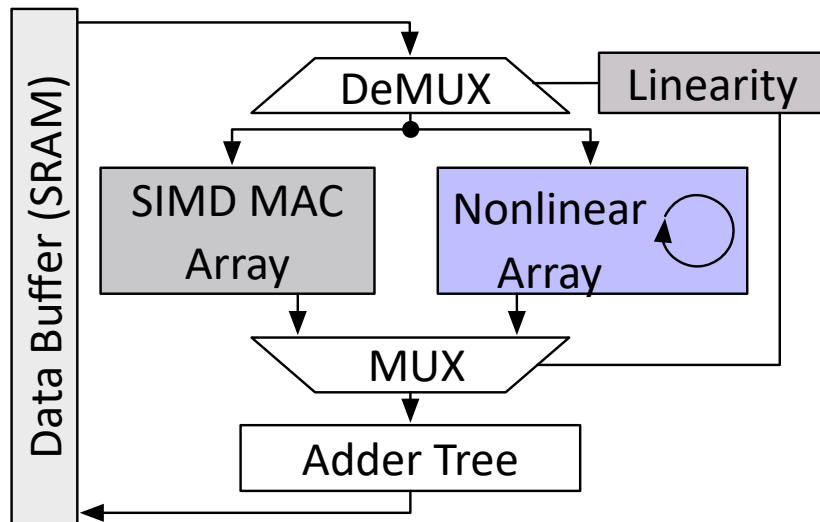
## ➤ Minimized overhead

- Gray blocks: original components for MAC
- Color blocks: extended components for nonlinear operations



# UNO SIMD array

- UNO SIMD array-based accelerators
  - Naive design (left): SIMD MAC array and SIMD nonlinear array





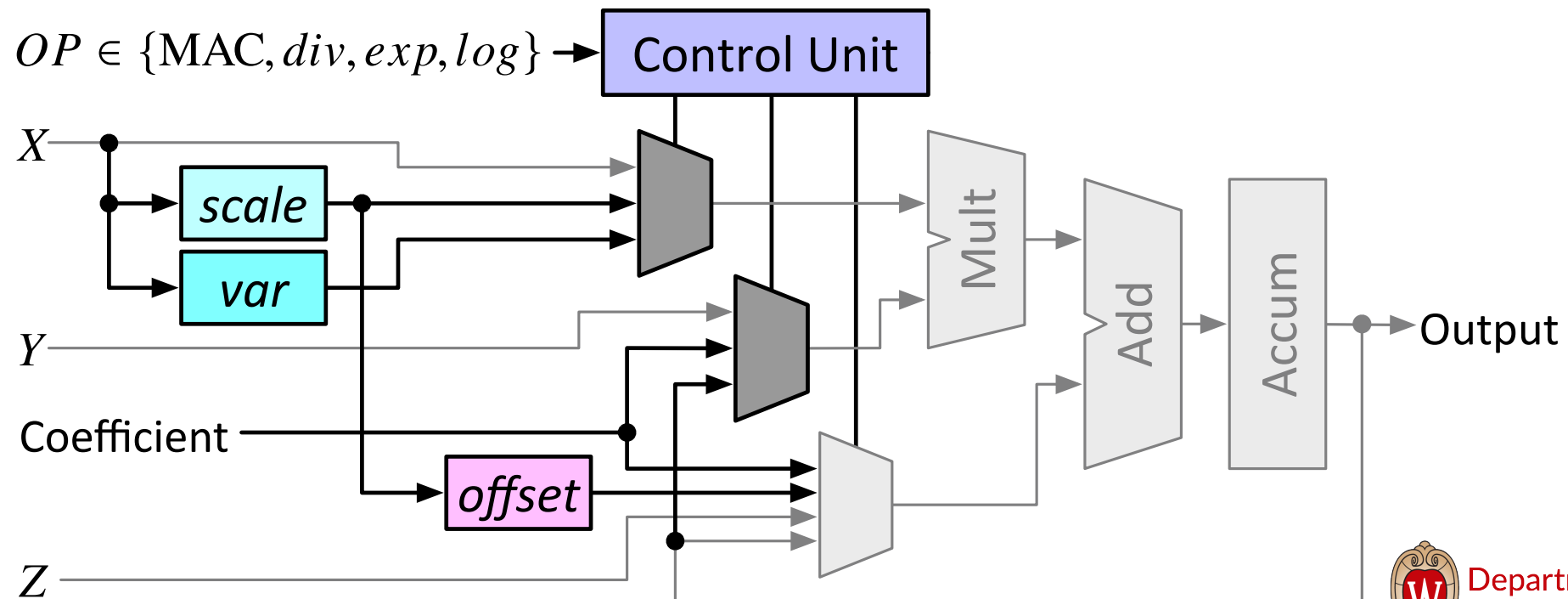
# Evaluate emerging neural networks

- UNO SIMD array-based accelerators
  - At least triples the energy efficiency

	Model	Naive design	UNO SIMD	Increase (%)
Area (mm <sup>2</sup> )	-	0.659	0.283	<b>-57.0</b>
Power (mW)	-	205.5	66.5	<b>-67.6</b>
Throughput (samples/s)	CapsNet	470.3	567.1	<b>+20.6</b>
	GNN	6.0	6.0	<b>+0.0</b>
	NMT	132.5	160.6	<b>+21.2</b>
Energy Efficiency (samples/s)	CapsNet	2288	8527	<b>+272.7</b>
	GNN	29	91	<b>+209.1</b>
	NMT	645	2415	<b>+274.5</b>

# Evaluate emerging neural networks

- UNO SIMD array-based accelerators
  - Ancillary logic exists in every PE
  - SIMD array suffers large overall overhead



# UNO systolic array as a salvage

- Two dataflows converge inside one architecture
  - Inherited GEMM dataflow and extended UNO dataflow
- PE overheads disperse in PE array for extra saving
  - Separation of ancillary logic into small parts
- Conventional approximation techniques are compatible
  - A heuristic microarchitecture approximation

# Outline

- Motivation
- Background
- Architecture
- Evaluation
- Conclusion

# Unify nonlinear operation (UNO)

## ➤ Leverage Horner's rule

- Taylor approximation: unify the math expression

$$f_n(x) = \sum_{i=0}^n \frac{f^{(i)}(a)}{i!} \cdot (x - a)^i = \sum_{i=0}^n c_i \cdot (x - a)^i$$

- MAC-compatible schedule: virtualize nonlinear operations using MAC from high to low degree

$$f_n(x) = offset + mac_n \cdot scale$$

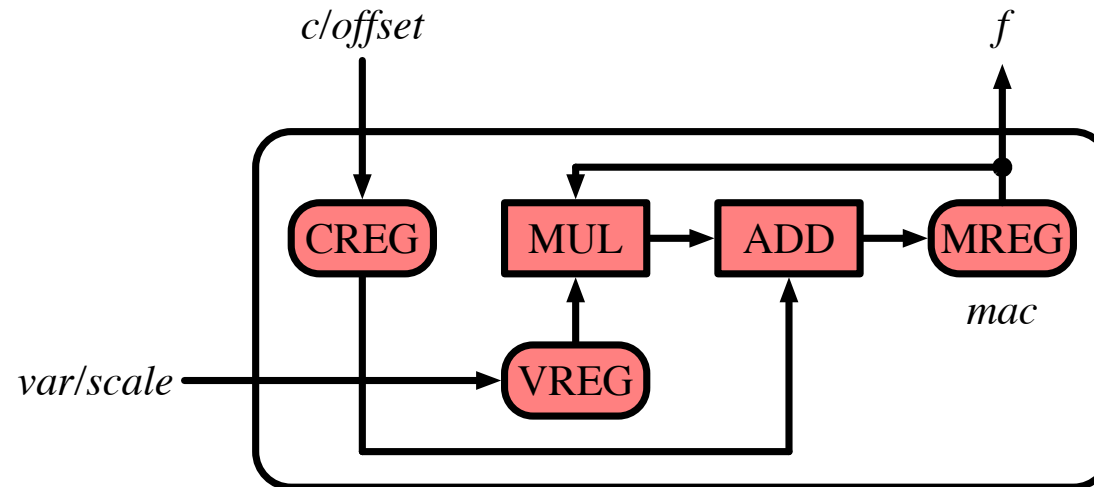
$$mac_i = \begin{cases} |c_{n-i}| + mac_{i-1} \cdot var & \text{if } 1 < i \leq n, \\ |c_{n-1}| + |c_n| \cdot var & \text{if } i = 1. \end{cases}$$

- Universal support: *div, exp, log, tanh, sigmoid, softmax, etc.*

# UNO PE

## ➤ Weight stationary

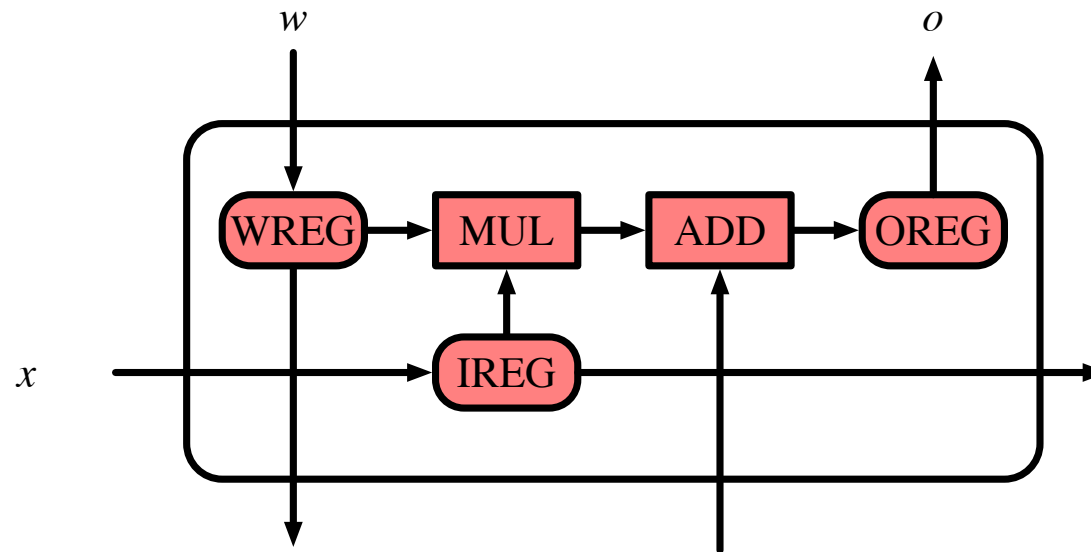
- Coefficient is stationary, i.e., statically stored, in CREG
- Var and scale in VREG are calculated from input
- Output is accumulated in MREG



# Systolic array PE for GEMM

## ➤ Weight stationary

- Weight is stationary (preloaded from top to bottom)
- Input is streamed from left to right
- Output is streamed from bottom to top



# Outline

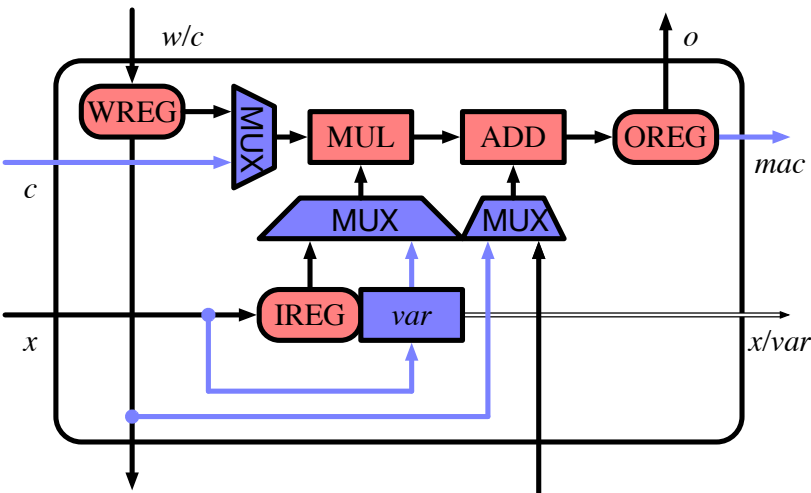
- Motivation
- Background
- Architecture
- Evaluation
- Conclusion



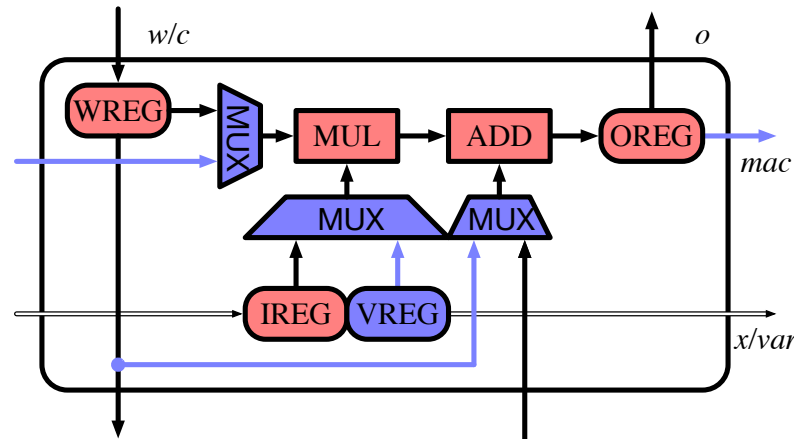
# UNO systolic array PE

## ➤ Heterogeneous PE

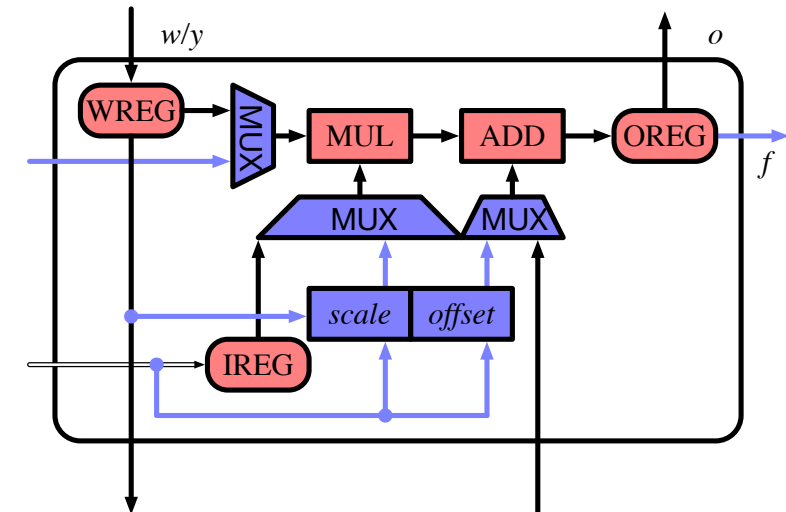
- Red block: original for GEMM dataflow
- Blue block: extended for UNO dataflow



Leftmost column



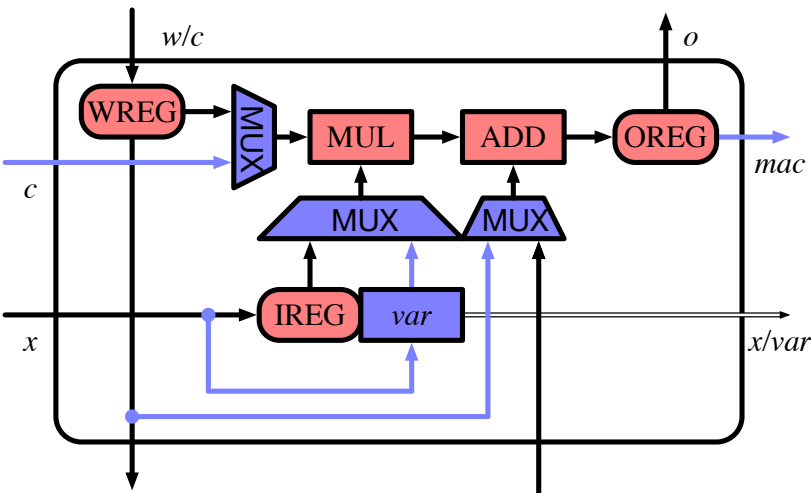
Middle column



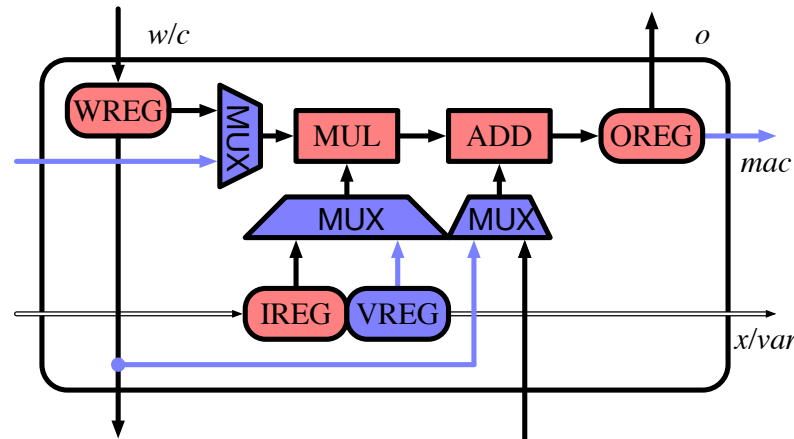
Rightmost column

# UNO systolic array PE

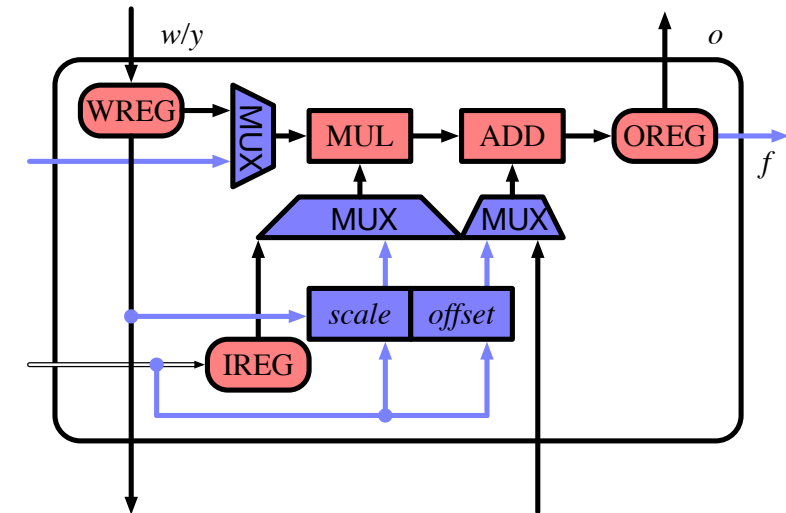
- Heterogeneous datapath
  - Black line: original for GEMM dataflow
  - Blue line: extended for UNO dataflow



Leftmost column



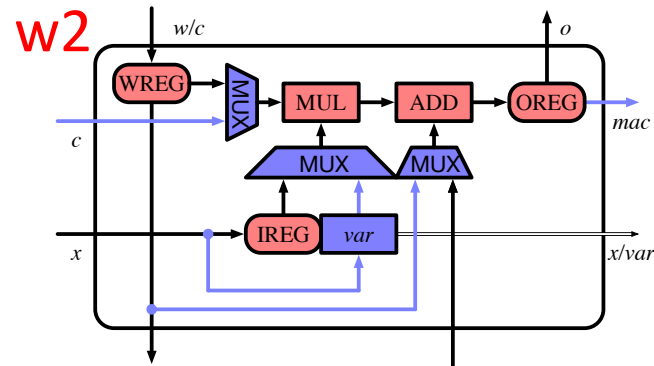
Middle column



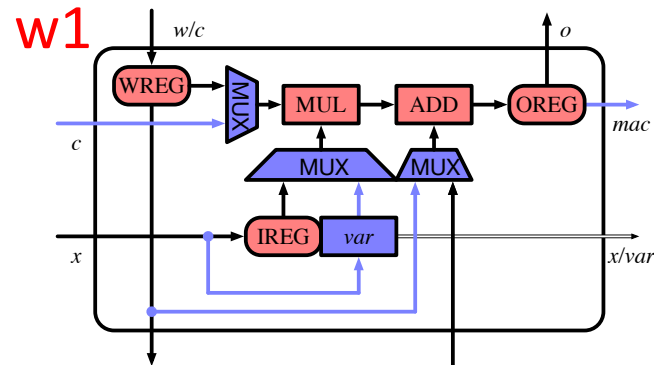
Rightmost column

# Walkthrough example-GEMM dataflow

- $o = x_1 * w_1 + x_2 * w_2$ 
  - Each vertical column performs identical operations

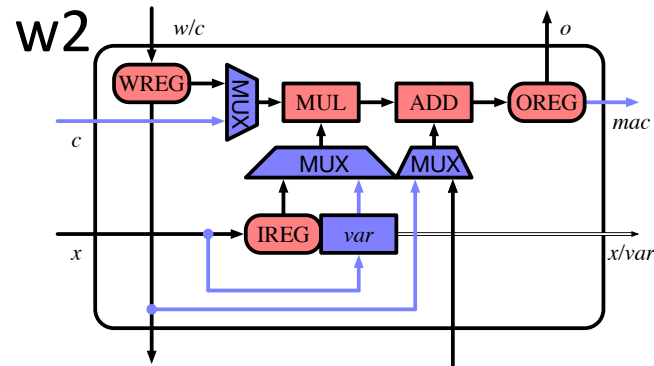


Cycle 0

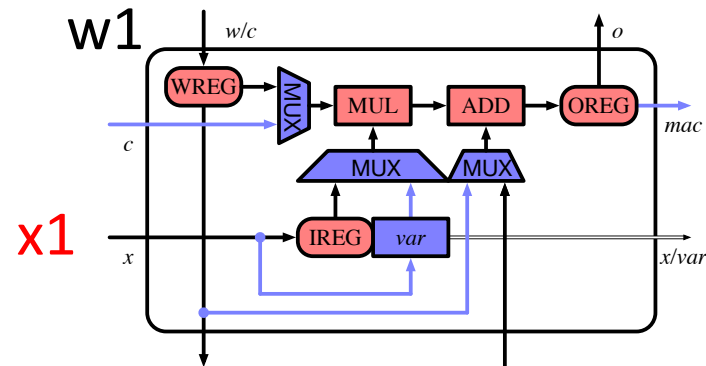


# Walkthrough example-GEMM dataflow

- $o = x_1 * w_1 + x_2 * w_2$ 
  - Each vertical column performs identical operations



Cycle 1

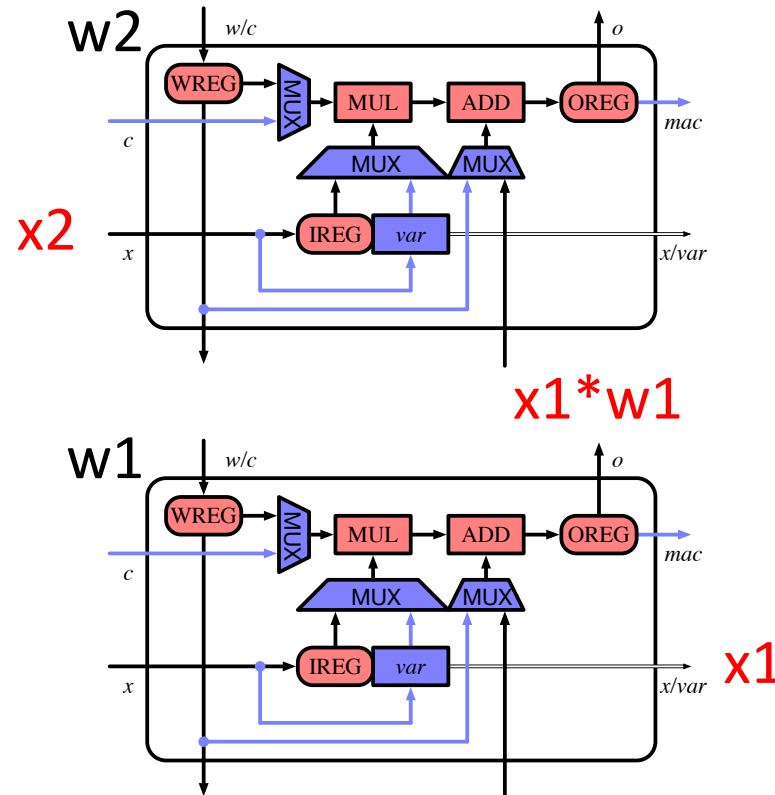


$x_1$

# Walkthrough example-GEMM dataflow

- $o = x_1 * w_1 + x_2 * w_2$ 
  - Each vertical column performs identical operations

Cycle 2

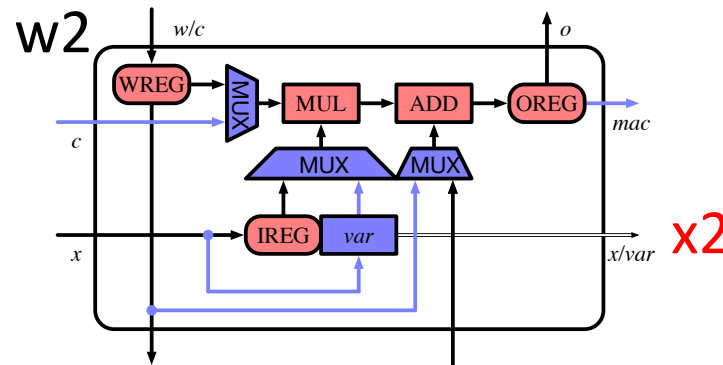


# Walkthrough example-GEMM dataflow

➤  $o = x_1 * w_1 + x_2 * w_2$

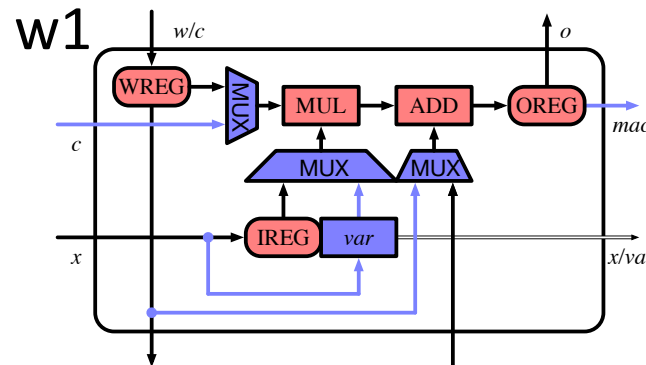
- Each vertical column performs identical operations

$o = x_1 * w_1 + x_2 * w_2$



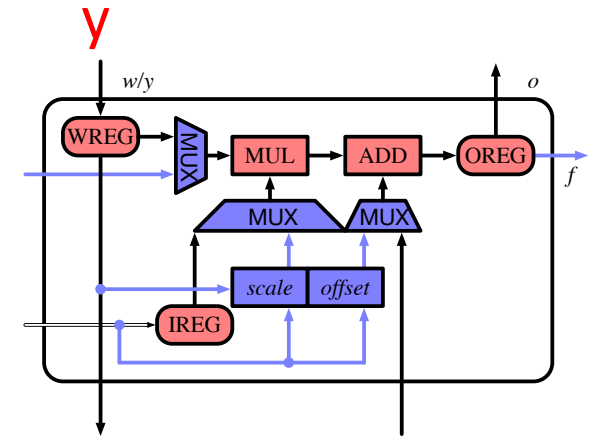
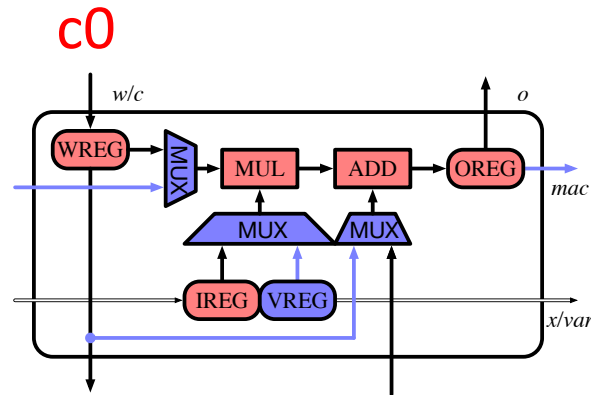
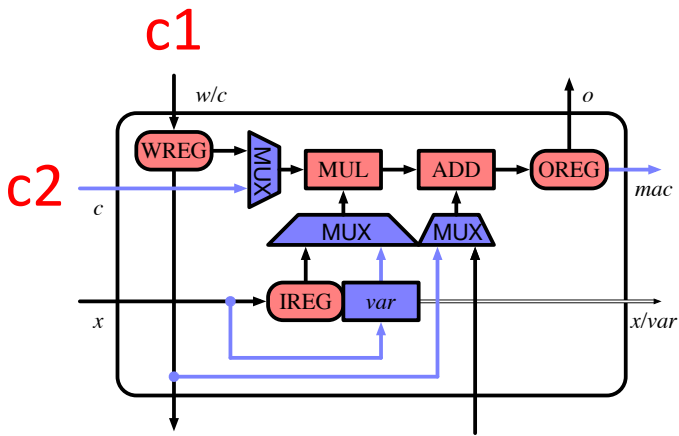
The previous MAC result is accumulated

Cycle 3



# Walkthrough example-UNO dataflow

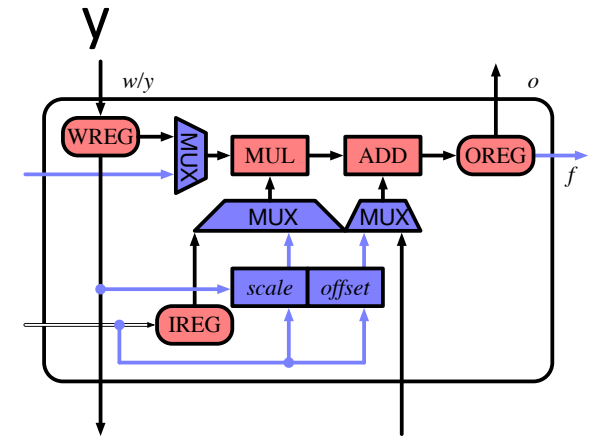
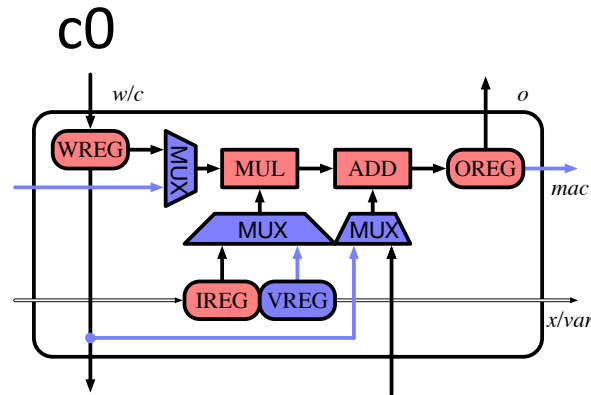
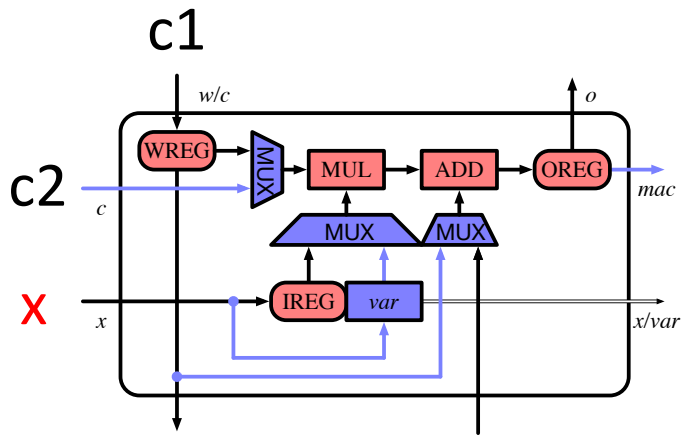
- $f = (c2 * var^2 + c1 * var + c0) * scale + offset$ 
  - Each horizontal row performs identical operations



Cycle 0

# Walkthrough example-UNO dataflow

- $f = (c2 * var^2 + c1 * var + c0) * scale + offset$ 
  - Each horizontal row performs identical operations

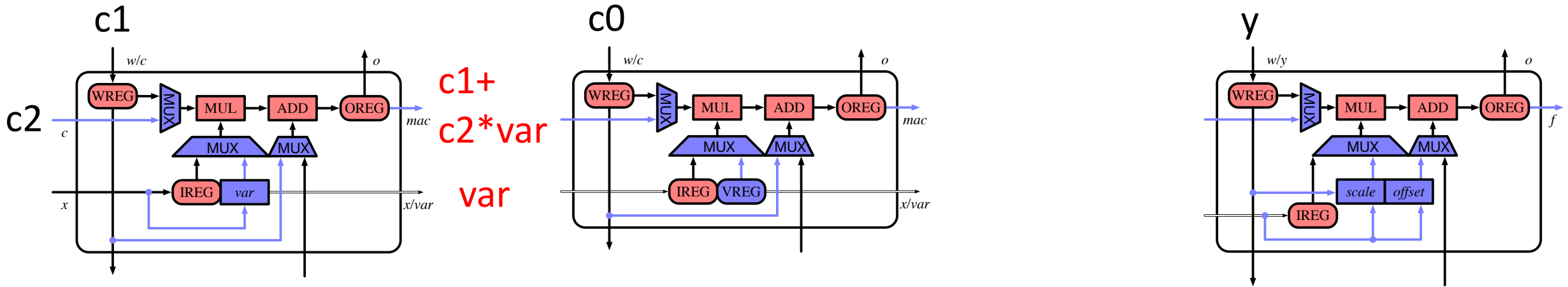


Cycle 1



# Walkthrough example-UNO dataflow

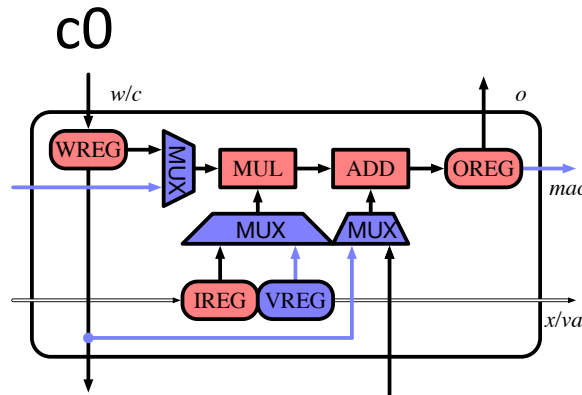
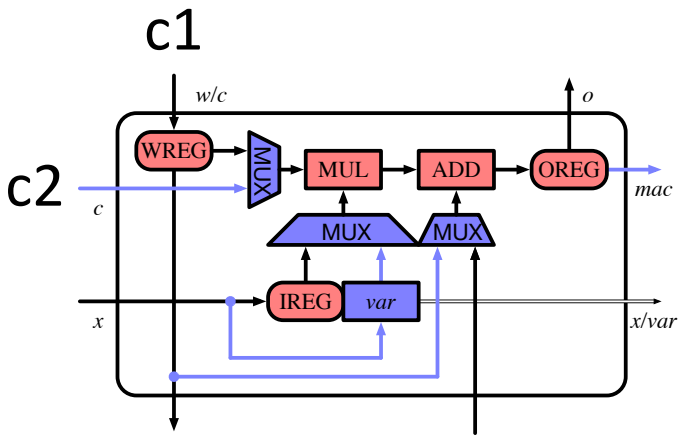
- $f = (c2 * var^2 + c1 * var + c0) * scale + offset$ 
  - Each horizontal row performs identical operations



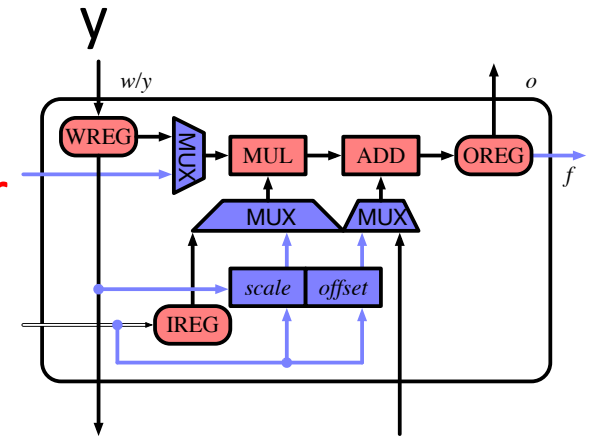
Cycle 2

# Walkthrough example-UNO dataflow

- $f = (c2 * var^2 + c1 * var + c0) * scale + offset$ 
  - Each horizontal row performs identical operations



$c0 +$   
 $(c1 + c2 * var) * var$   
 $var$

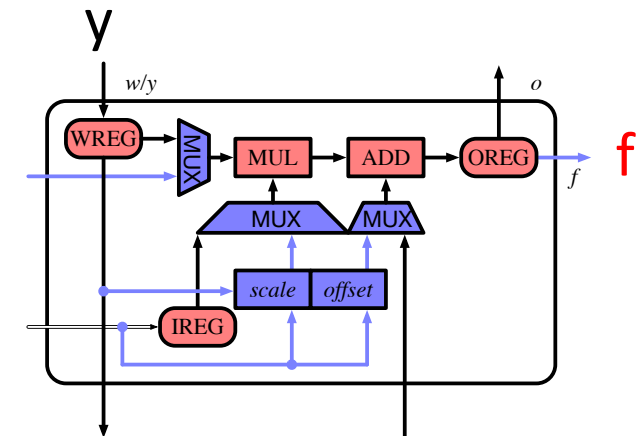
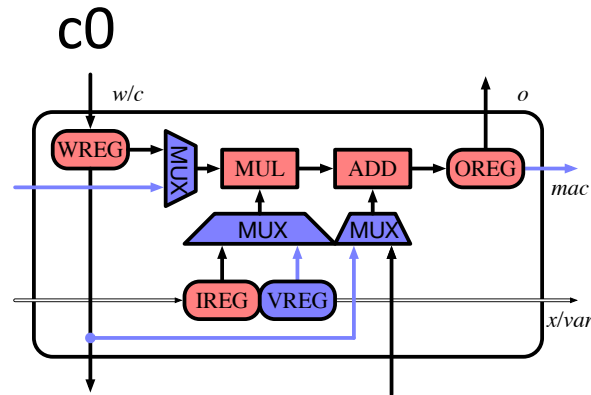
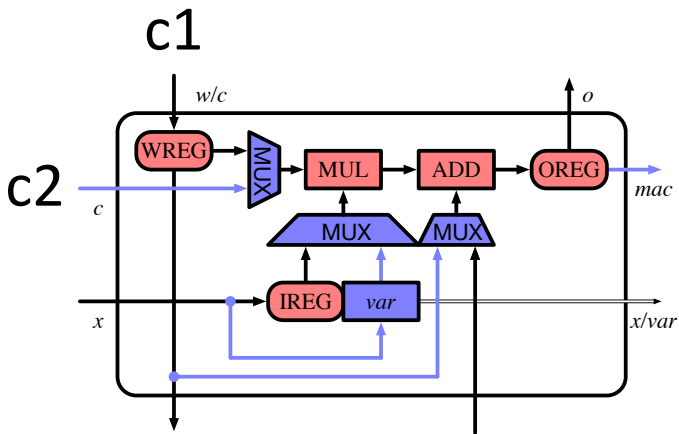


**The previous MAC result is multiplied**

Cycle 3

# Walkthrough example-UNO dataflow

- $f = (c2 * var^2 + c1 * var + c0) * scale + offset$ 
  - Each horizontal row performs identical operations



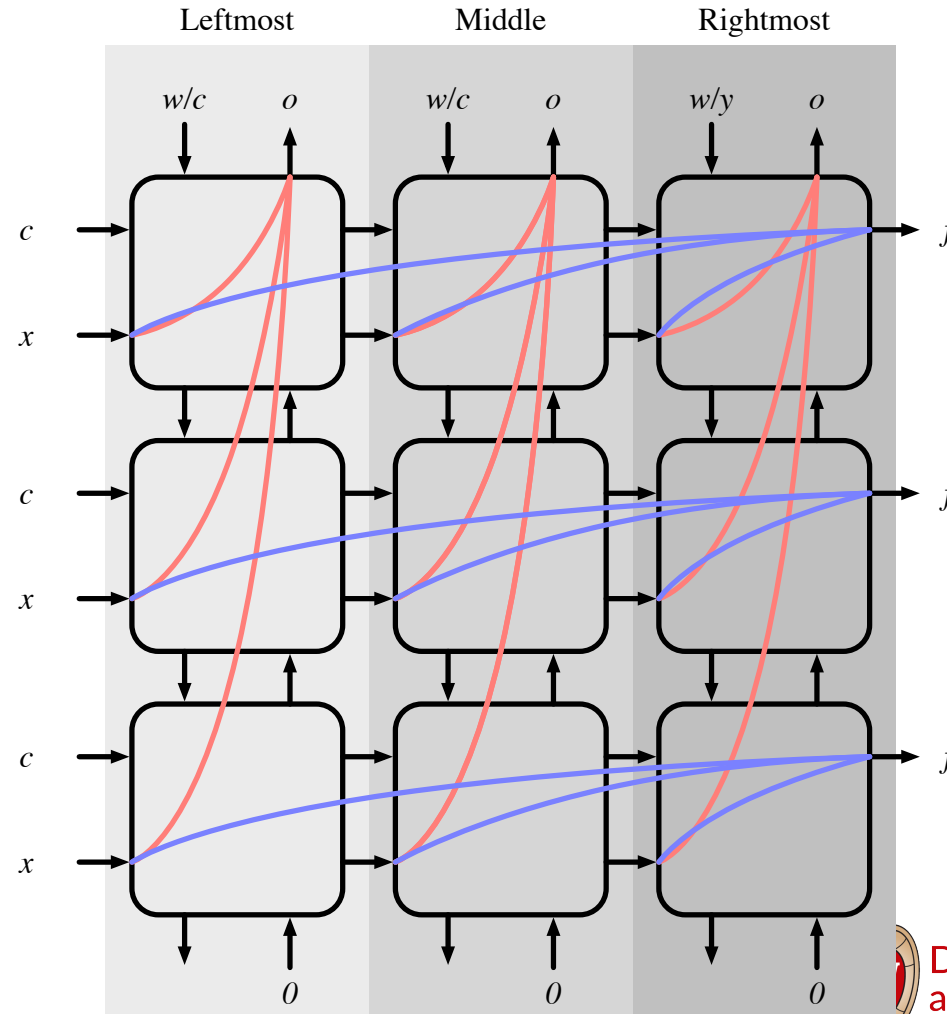
Cycle 4

# UNO systolic array

## ➤ Heterogeneous dataflows

Red line: GEMM dataflow in vertical column

Blue line: UNO dataflow in horizontal row



# Cross-level approximation

## ➤ Heuristic approximation

- UNO algorithm proceeds from high to low degree
  - Low degree term is multiplied less, and can be more eagerly approximated
- UNO dataflow proceeds from left to right
  - Right PE multiplies less, and can have one less bit for MAC than current PE
- N-bit UNO systolic array of size R-by-C
  - N bits for the leftmost PE
  - $(N-C+1)$  bits for the rightmost PE

# Outline

- Motivation
- Background
- Architecture
- Evaluation
- Conclusion

# Experimental setup

- DNN
  - CapsNet: *div, exp, log, sigmoid* and *softmax*
- Accuracy
  - Customized PyTorch UNO operator
  - Post-training quantization (no quantization-aware training)
- Hardware
  - Systolic array size  $R=C=4$  or  $8$  (SIMD array size  $R*C=16$  or  $64$ )
  - Bitwidth  $N=16$
  - Synthesized with TSMC 32nm technology @ 400MHz

# Accuracy

## ➤ Operation-level

- 16 bits: 1-bit sign, 5-bit integer and 10-bit fraction

Operation	UNO systolic (R,C,N)		UNO SIMD (R,C,N)	
	4,4,16	8,8,16	4,4,16	8,8,16
<i>div</i>	0.022	0.039	0.021	0.001
<i>exp</i>	0.080	0.032	0.079	0.001
<i>log</i>	0.005	0.026	0.005	0.000



# Accuracy

## ➤ Application-level

- FP32 accuracy = 98.78%

Sign-Int-Frac (bit)	UNO systolic (R,C,N)		UNO SIMD (R,C,N)	
	4,4,16	8,8,16	4,4,16	8,8,16
1-7-8	98.70	86.27	98.69	98.68
1-6-9	98.87	97.75	98.87	98.85
1-5-10	98.75	98.36	98.76	98.78
1-4-11	94.05	93.55	93.98	93.78

# Hardware performance

Performance	UNO systolic (R,C,N)		UNO SIMD (R,C,N)		Improve (%)	
	4,4,16	8,8,16	4,4,16	8,8,16	4,4,16	8,8,16
Area (mm <sup>2</sup> )	0.13	0.41	0.18	0.73	<b>29.4</b>	<b>43.3</b>
Power (mW)	8.33	28.52	11.45	46.12	<b>27.2</b>	<b>38.2</b>
Throughput (ksamples/s)	0.15	0.52	0.15	0.52	<b>0.0</b>	<b>0.0</b>
Energy efficiency (ksamples/J)	18.58	18.34	13.52	11.34	<b>37.4</b>	<b>61.7</b>

# Nonlinear operations are underrated?

- Cost of varying PEs over a MAC PE
  - Naively slabbing multiple PEs together

Function	Performance	Naive PE	MAC PE	Overhead
<i>MAC, div, exp, log</i>	Area ( $\mu\text{m}^2$ )	9,323	1,559	6.0X
	Power (mW)	2.35	0.56	4.2X

# Nonlinear operations are underrated?

- Cost of varying PEs over a MAC PE
  - UNO SIMD array PE

Function	Performance	SIMD PE	MAC PE	Overhead
<i>MAC, div, exp, log</i>	Area ( $\mu\text{m}^2$ )	4,221	1,559	2.7X
	Power (mW)	0.93	0.56	1.7X

# Nonlinear operations are not underrated

- Cost of varying PEs over a MAC PE
  - UNO systolic array PE (This work)

Function	Performance	Systolic PE	MAC PE	Overhead
<i>MAC, div, exp, log</i>	Area ( $\mu\text{m}^2$ )	1,832	1,559	1.2X
	Power (mW)	0.36	0.56	0.6X

# Outline

- Motivation
- Background
- Architecture
- Evaluation
- Conclusion

# Conclusion

- We identify the bottleneck of UNO SIMD array for emerging DNNs
  - Overhead repetition in every PE
- We propose a reconfigurable systolic array based on UNO
  - Flexibly execution of both linear and nonlinear operations
  - Cross-level approximation with heuristics
- We evaluate the proposed design on an emerging DNN
  - Tremendous performance gain without accuracy loss

Thank you!  
Q & A

Di Wu and Joshua San Miguel  
[di.wu@ece.wisc.edu](mailto:di.wu@ece.wisc.edu) and [jsanmiguel@wisc.edu](mailto:jsanmiguel@wisc.edu)

