

# UNO: Virtualizing and Unifying Nonlinear Operations for Emerging Neural Networks

**Di Wu, Jingjie Li, Setareh Behroozi,  
Younghyun Kim and Joshua San Miguel**  
University of Wisconsin–Madison



# Outline

- Motivation**
- Algorithm
- Microarchitecture
- Evaluation
- Conclusion

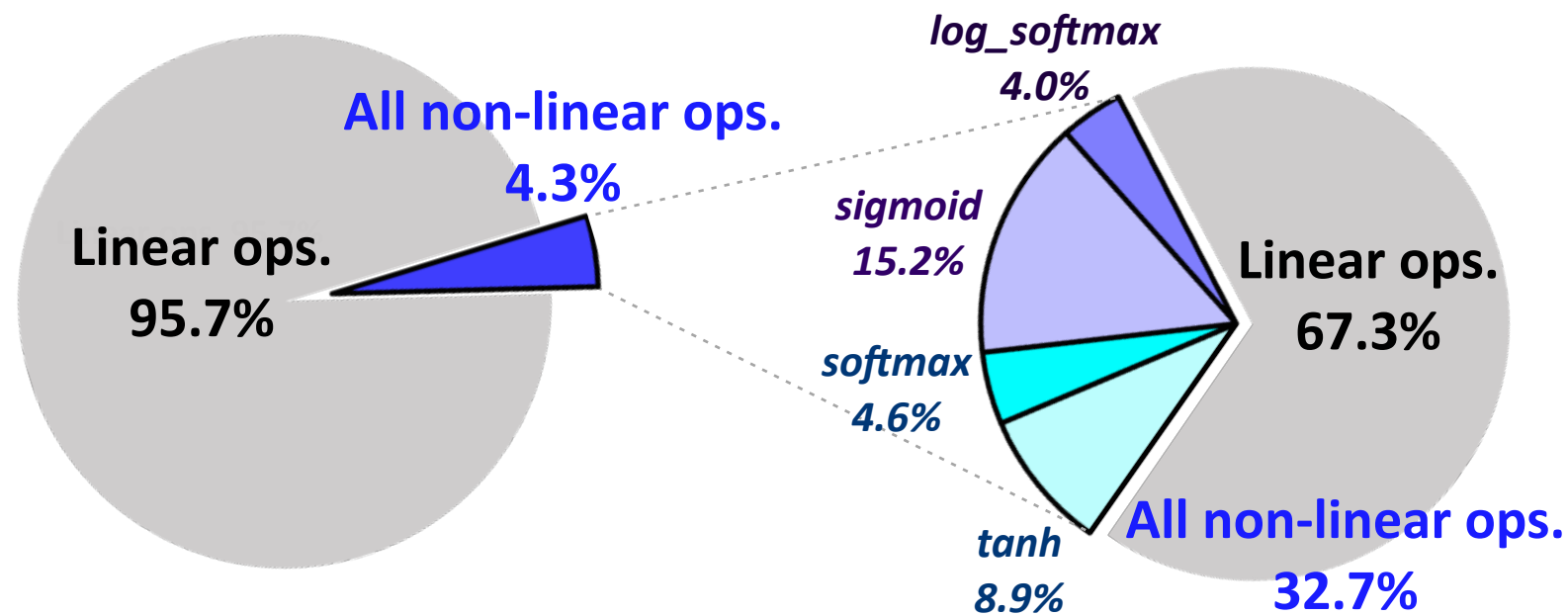
# Nonlinear operations are underrated

- 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



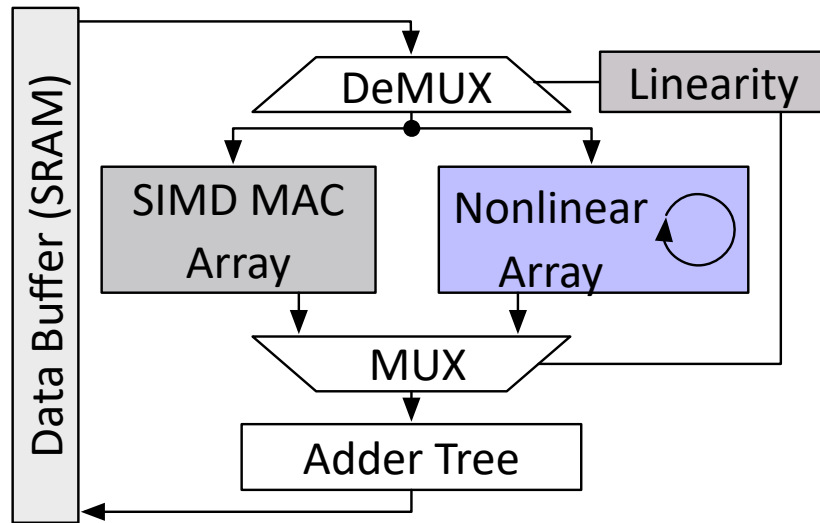
Operation count

Operation runtime

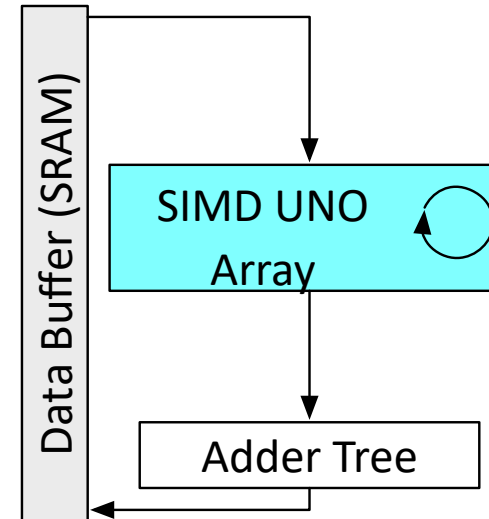
Profiling NMT model on CPU

# Nonlinear operations are underrated

- Emerging deep neural networks have nonlinear operations with
  - High diversity
  - High complexity
  - High cost: hardware overhead in accelerators



MAC array-based accelerators exhibit **low efficiency**



**UNO** array-based accelerators to **improve efficiency**

# Outline

- Motivation
- **Algorithm**
- Microarchitecture
- Evaluation
- Conclusion

# Unify nonlinear operations (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

$$f_n(x) = \text{offset} + \text{mac}_n \cdot \text{scale}$$

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

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

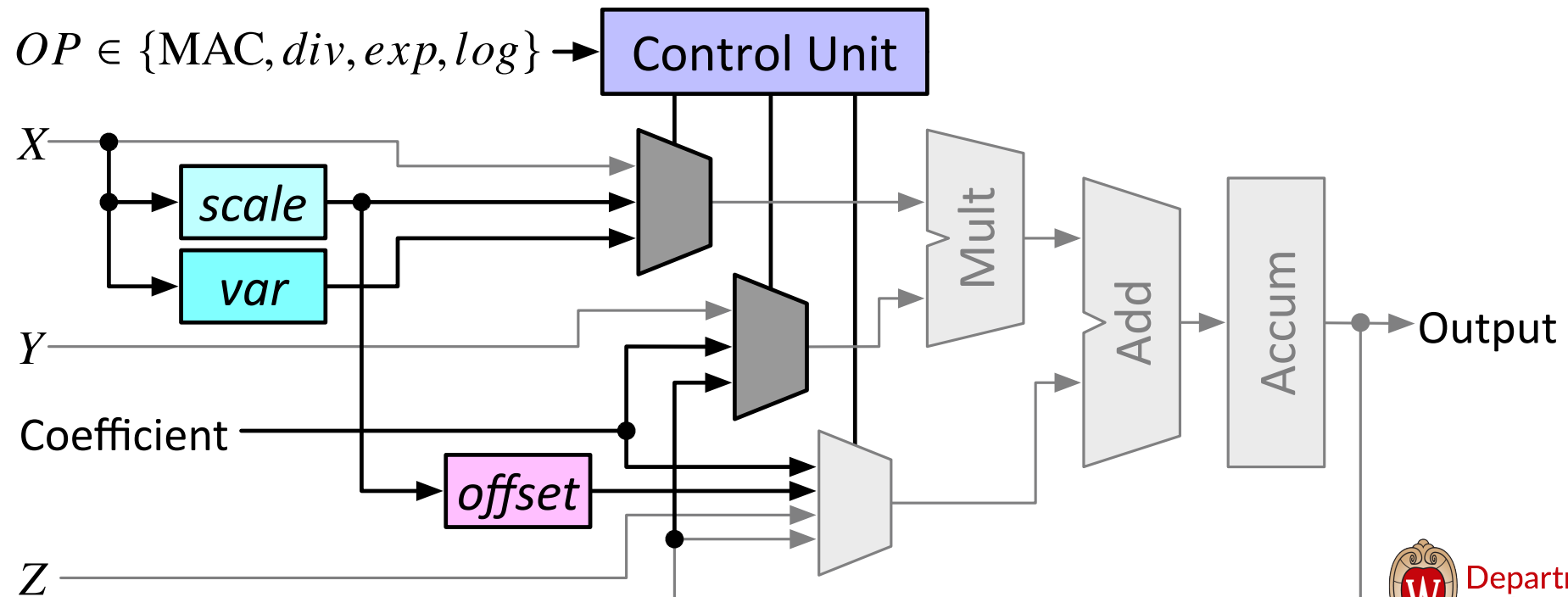
# Outline

- Motivation
- Algorithm
- Microarchitecture**
- Evaluation
- Conclusion



# Extend a standard MAC unit

- Minimize overhead to support nonlinear operations
  - Gray blocks: original MAC unit components
  - Color blocks: **extended** components



# Outline

- Motivation
- Algorithm
- Microarchitecture
- Evaluation**
- Conclusion

# Evaluate individual nonlinear operations

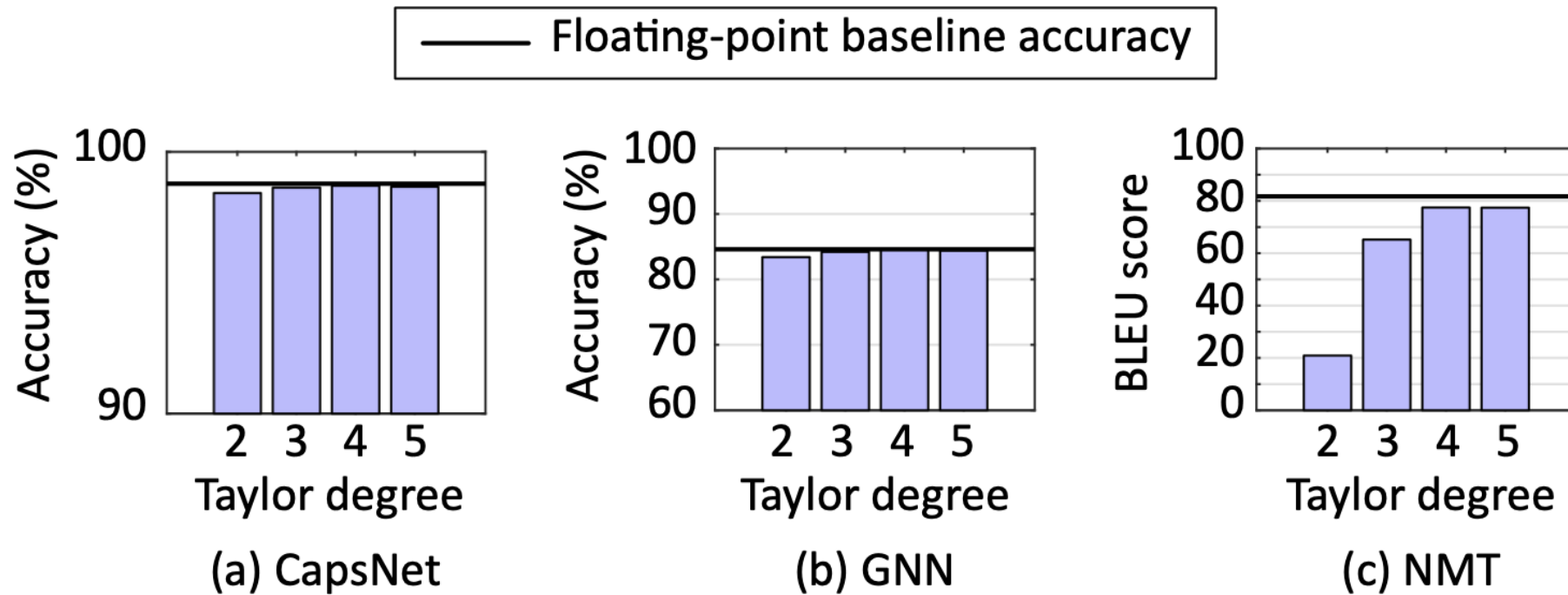
## ➤ Compare baseline and UNO processing elements

- Similar accuracy
- Lower area/power

div/exp/log refers to designs in [11,12,13] to compare against.

Design		Accuracy (%)	Area ( $\mu\text{m}^2$ )	Power (mW)
Baseline (MAC+ div+ exp+ log)	div	99.93	3,068	1.09
	exp	99.70	1,118	0.22
	log	98.79	3,578	0.48
	Total	-	<b>9,323</b>	<b>2.35</b>
UNO	div	99.92	-	-
	exp	99.91	-	-
	log	99.15	-	-
	Total	-	<b>4,221</b>	<b>0.93</b>

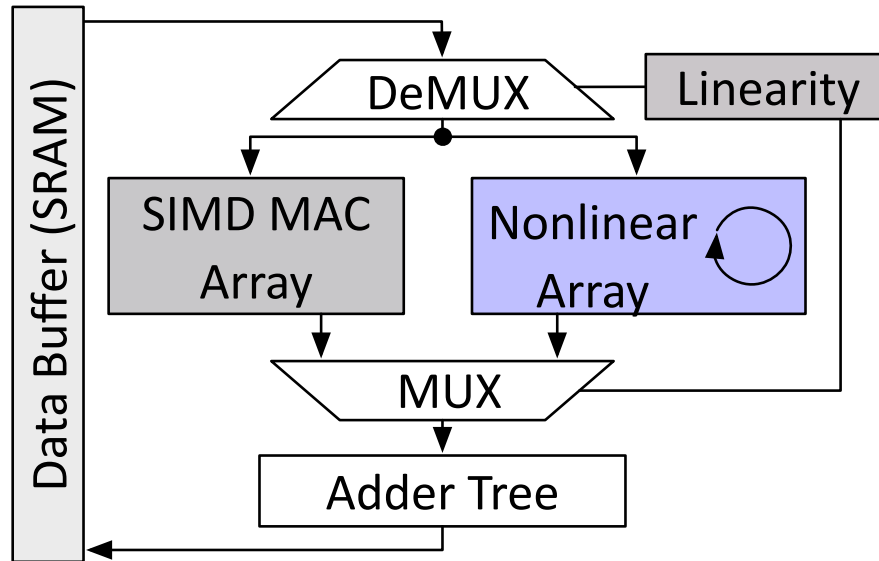
# Evaluate emerging neural networks



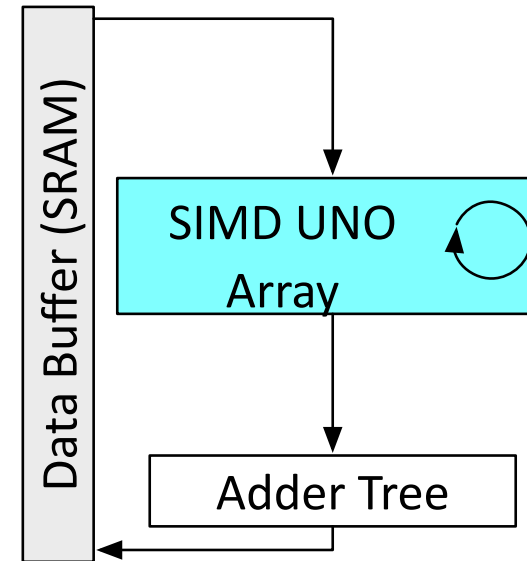
## ➤ Accuracy evaluation

- UNO attains **low accuracy loss** even with heavy approximation across different neural networks

# Evaluate emerging neural networks



MAC-based SIMD array



UNO-based SIMD array

## ➤ Efficiency evaluation (SIMD 64)

- a) Baseline: MAC array based
- b) Proposed: UNO array based

\*Nonlinear array in (a) refers to designs in [11,12,13].

# Evaluate emerging neural networks

	Model	Baseline	UNO	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>

## ➤ Efficiency evaluation (SIMD 64)

- UNO at least triples the energy efficiency for all evaluated models

# Evaluate emerging neural networks

	Model	Baseline	UNO	Increase (%)
Area (mm <sup>2</sup> )	-	0.659	0.283	-57.0
Power (mW)	-	205.5	66.5	-67.6
Throughput (samples/s)	CapsNet	470.3	567.1	+20.6
	GNN	6.0	6.0	+0.0
	NMT	132.5	160.6	+21.2
Energy Efficiency (samples/s)	CapsNet	4648	10953	+135.6
	GNN	68	116	+70.3
	NMT	1538	3036	+97.9

- Efficiency evaluation (SIMD 64)
  - UNO still doubles the energy efficiency **with power gating**

# Outline

- Motivation
- Algorithm
- Microarchitecture
- Evaluation
- Conclusion**



# Conclude this work

- UNO addresses the inefficiency of nonlinear operations in emerging neural networks.
- UNO unifies nonlinear operations with Taylor approximation.
- UNO virtualizes Taylor approximation on the MAC unit to minimize hardware overhead and boost efficiency.

**Thank you!**  
**Q & A**

**Di Wu, Jingjie Li, Setareh Behroozi,  
Younghyun Kim and Joshua San Miguel  
University of Wisconsin–Madison**

