

Tanh Works Better with Asymmetry

Dongjin Kim^{1,3}, Woojeong Kim², Suhyun Kim³

¹Korea University, ²Cornell University, ³Korea Institute of Science and Technology



KOREA
UNIVERSITY



CORNELL
UNIVERSITY

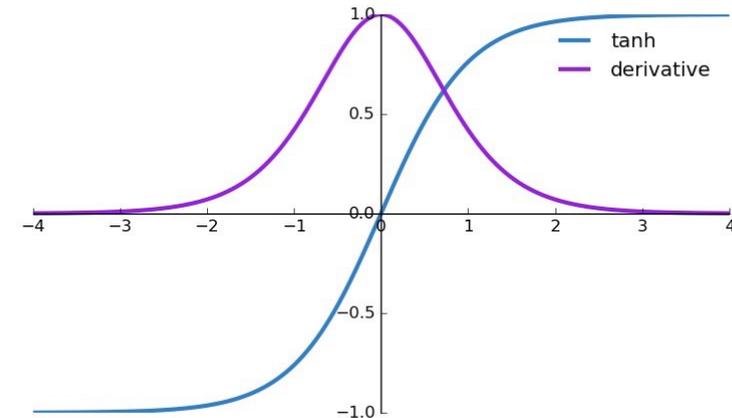
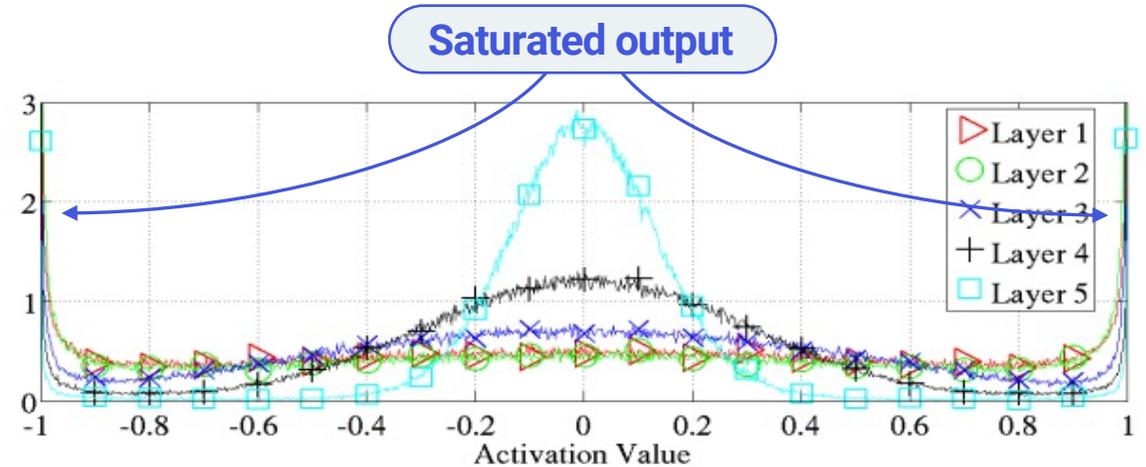


An activation function with two boundaries

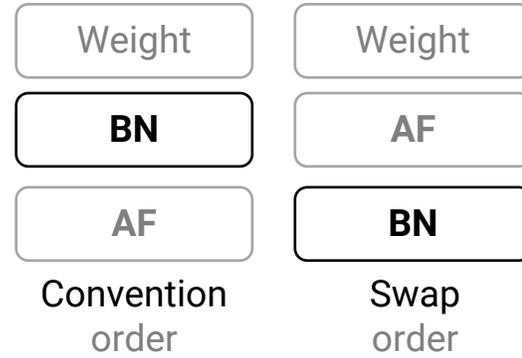
- It has a **saturation state**.
The activation is close to the asymptotic value.

- The saturated output suffered from the vanishing gradients problem

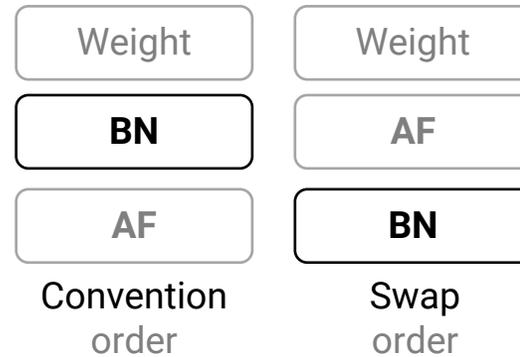
- Due to poor performance, Tanh becomes forgotten.



Batch Normalization(BN)
places between the weight
and activation function(AF) .



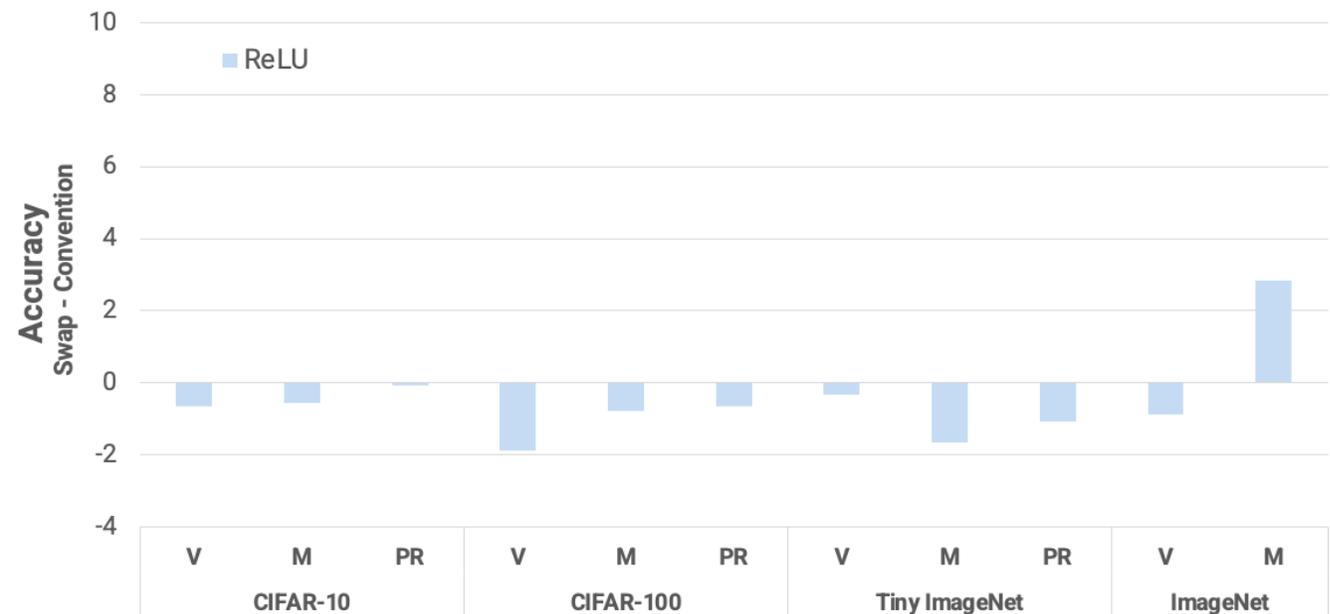
Batch Normalization(BN)
places between the weight
and activation function(AF) .



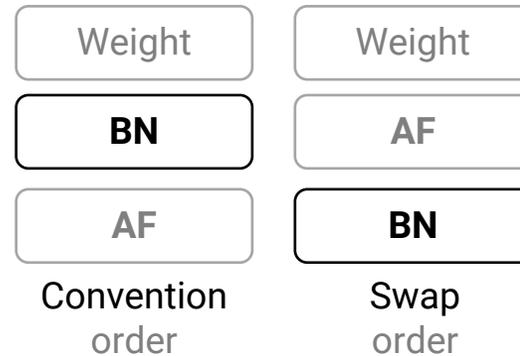
ReLU

The layer order does not
dramatically change accuracy.

Datasets	Models	ReLU	
		Convention	Swap
CIFAR-10	VGG16	93.69	93.04
	MobileNet	92.48	91.93
	PreAct-ResNet18	94.94	94.86
CIFAR-100	VGG16	73.68	71.79
	MobileNet	70.27	69.49
	PreAct-ResNet18	78.06	77.39
Tiny ImageNet	VGG16	59.37	59.05
	MobileNet	51.90	50.25
	PreAct-ResNet34	67.28	66.21
ImageNet	VGG16	73.83	72.95
	MobileNet	68.27	71.1



Batch Normalization(BN)
places between the weight
and activation function(AF) .



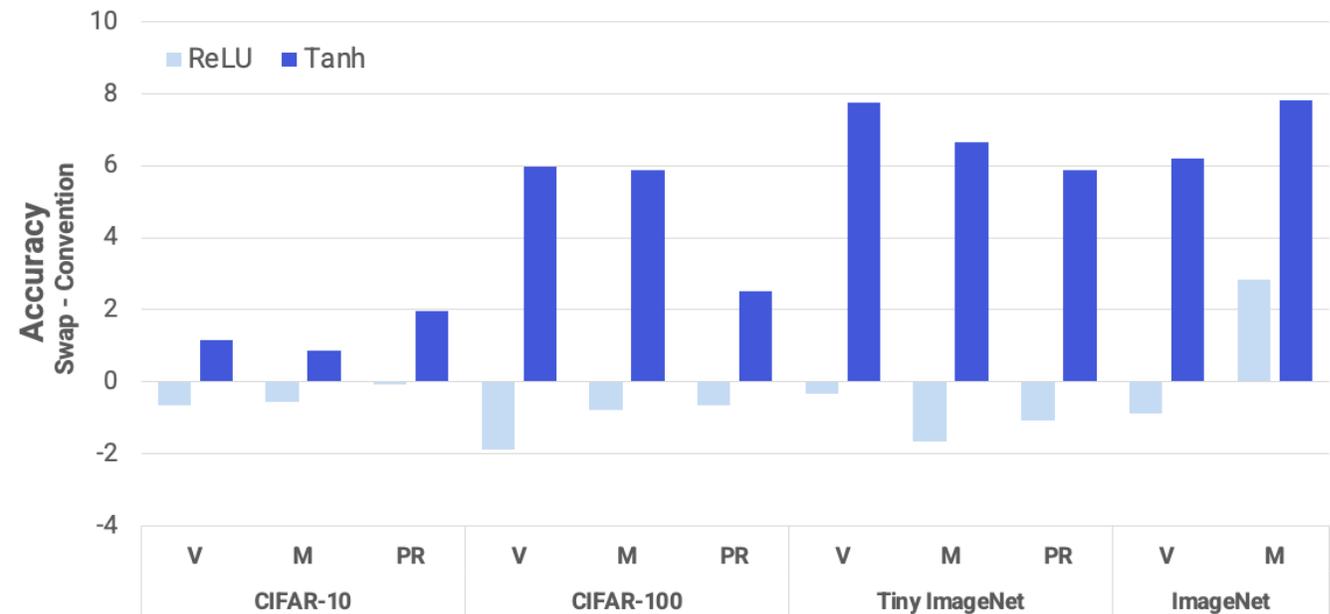
ReLU

The layer order does not
dramatically change accuracy.

Tanh

The Swap order significantly
outperforms the Convention order.

Datasets	Models	ReLU		Tanh	
		Convention	Swap	Convention	Swap
CIFAR-10	VGG16	93.69	93.04	91.75	92.90
	MobileNet	92.48	91.93	91.66	92.53
	PreAct-ResNet18	94.94	94.86	92.46	94.41
CIFAR-100	VGG16	73.68	71.79	64.95	70.93
	MobileNet	70.27	69.49	64.50	70.39
	PreAct-ResNet18	78.06	77.39	73.26	75.76
Tiny ImageNet	VGG16	59.37	59.05	49.29	57.05
	MobileNet	51.90	50.25	45.38	52.05
	PreAct-ResNet34	67.28	66.21	59.06	64.94
ImageNet	VGG16	73.83	72.95	60.85	67.04
	MobileNet	68.27	71.1	64.26	72.07



Batch Normalization(BN)
places between the weight
and activation function(AF).



Datasets	Models	ReLU		Tanh	
		Convention	Swap	Convention	Swap
CIFAR-10	VGG16	93.69	93.04	91.75	92.90
	MobileNet	92.48	91.93	91.66	92.53
	PreAct-ResNet18	94.94	94.86	92.46	94.41
CIFAR-100	VGG16	73.68	71.79	64.95	70.93
	MobileNet	70.27	69.49	64.50	70.39

Goal

1. Reveal a hidden property

Why is the Swap order effective on Tanh?

2. Modified Activation Function

How can we redesign the order-agnostic Tanh with improved accuracy?

Tanh

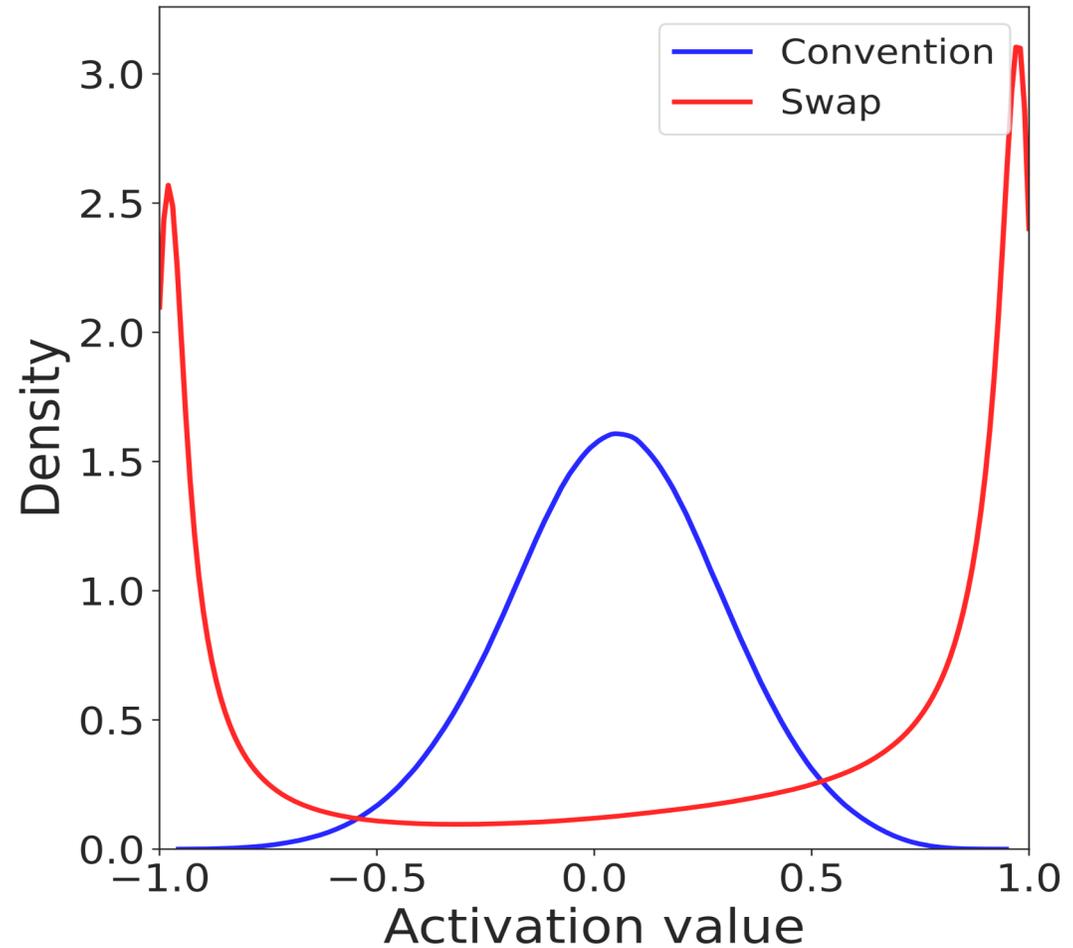
The Swap order significantly
outperforms the Convention order.

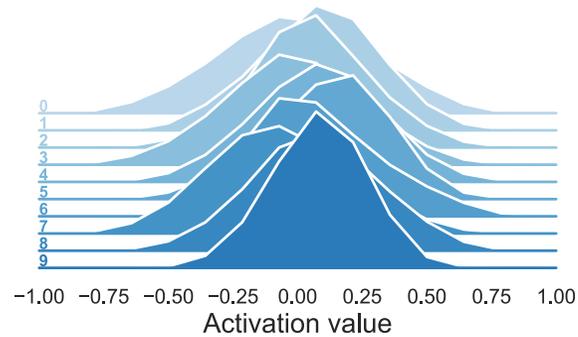
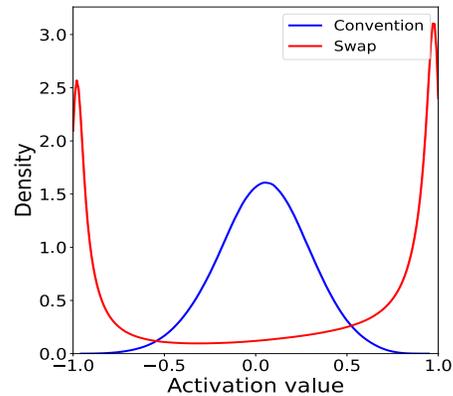


Layer-wise Activations of Tanh

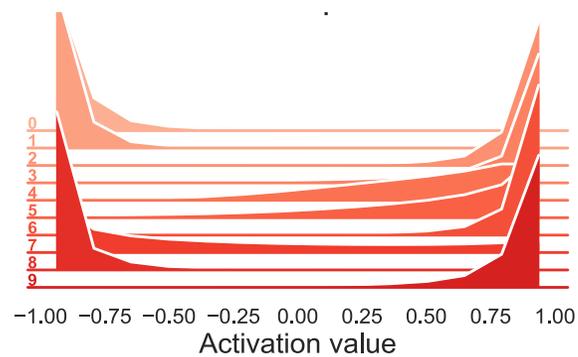


**Symmetrically
Distributed**





Convention model
Symmetric activation

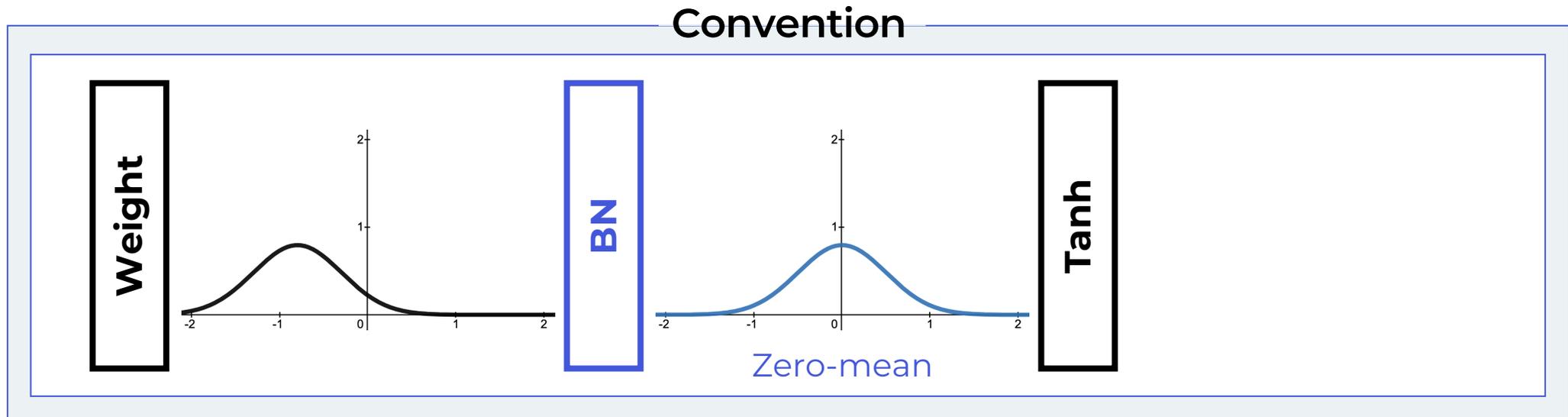


Swap model
Asymmetric activation

Channel-wise activation of the Swap order is **asymmetrically distributed**.

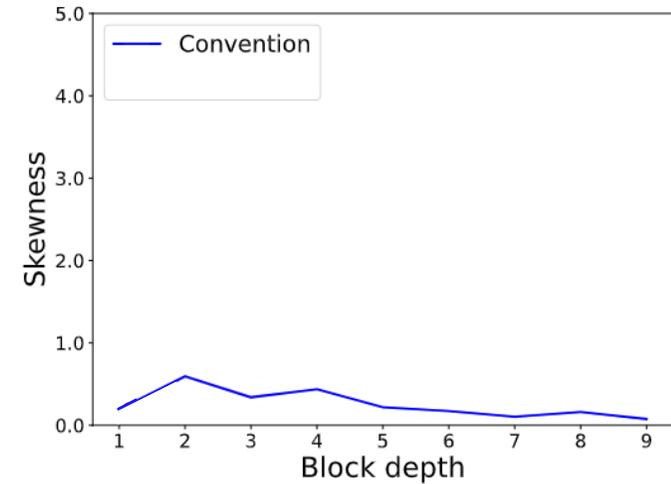
The Reason for Low Asymmetry in the Convention Order

- Batch Normalization shifts the biased weighted sum outputs to **zero**.
- In the Convention order, the zero mean distribution generates **symmetric activation on Tanh**.

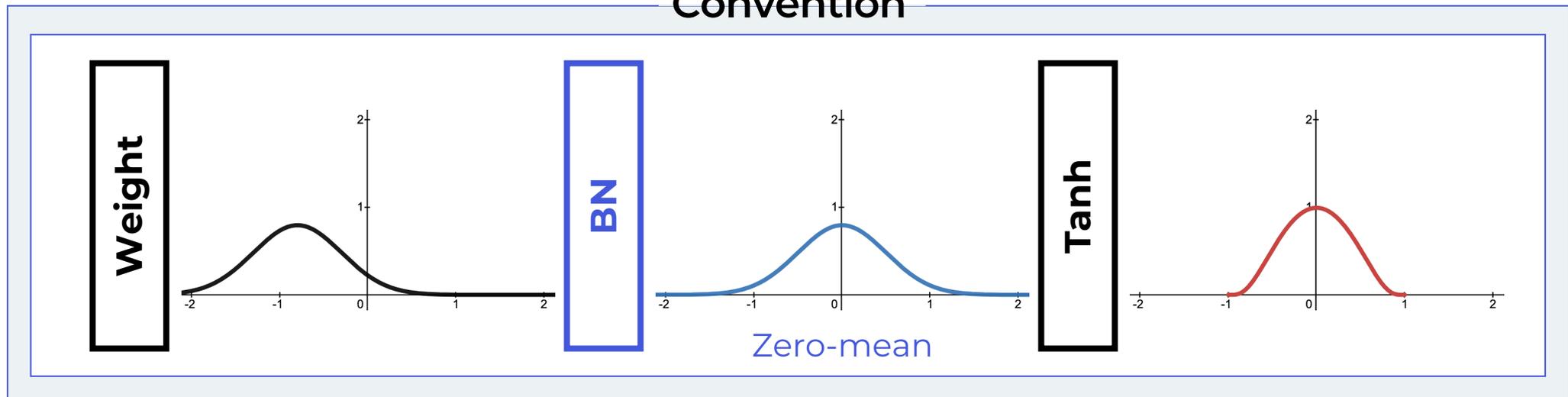


The Reason for Low Asymmetry in the Convention Order

- Batch Normalization shifts the biased weighted sum outputs to **zero**.
- In the Convention order, the zero mean distribution generates **symmetric activation on Tanh**.



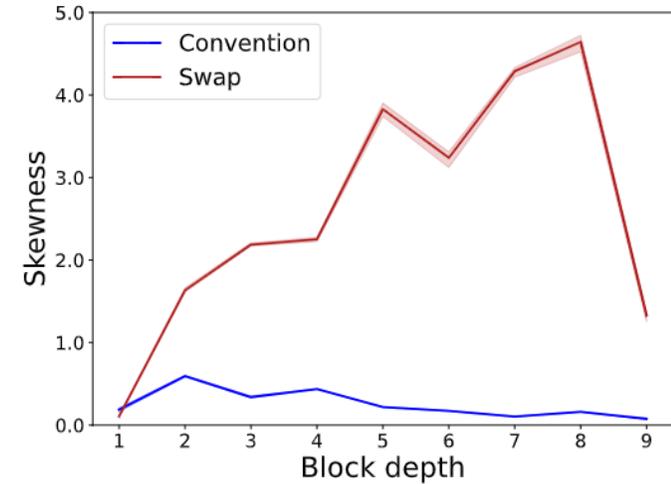
Convention



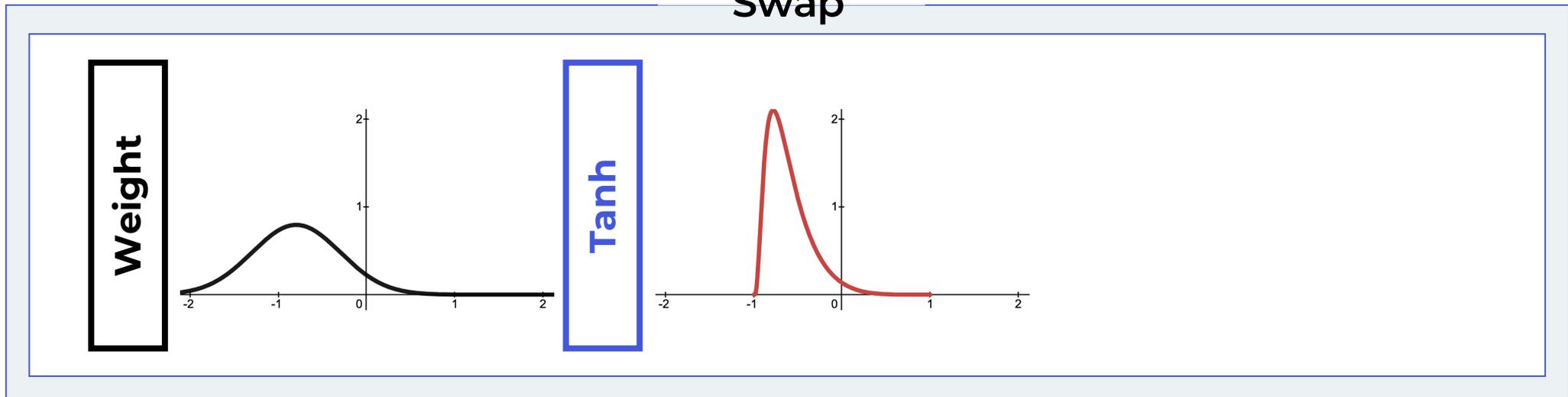
What Brings Asymmetry in the Swap Order?

The elimination of the preceding
Batch Normalization

The biased distribution to Tanh encourages
asymmetric saturation in the Swap order.



Swap



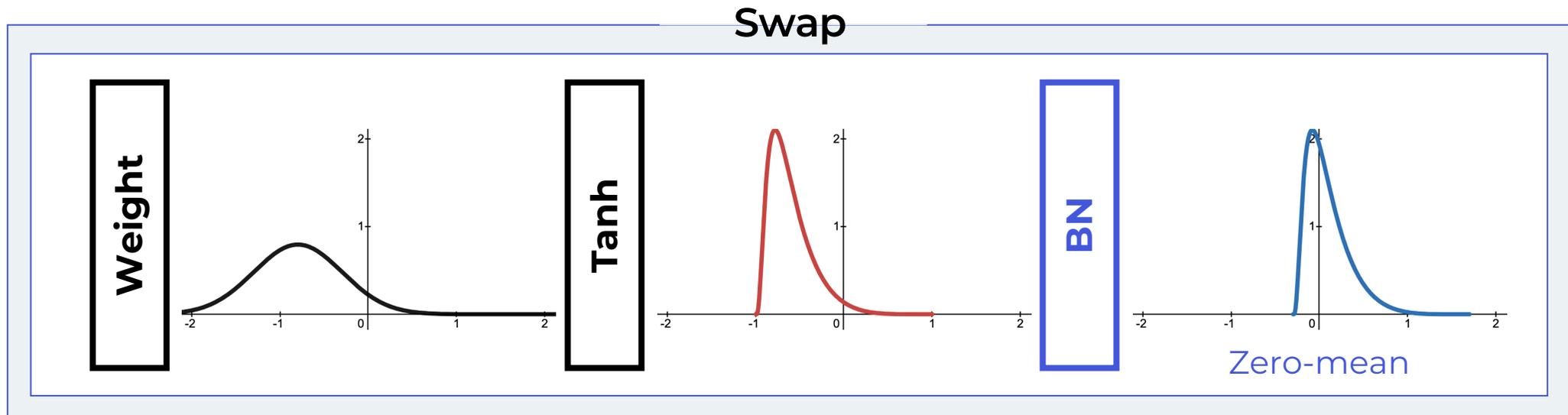
Additional Improvement incurred by Asymmetry

The elimination of Batch Normalization before the Tanh

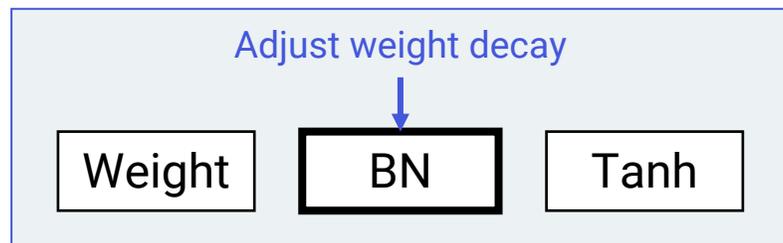
The biased distribution to Tanh encourages asymmetric saturation in the Swap order.

Asymmetric saturation in the Swap order incurs sparsity

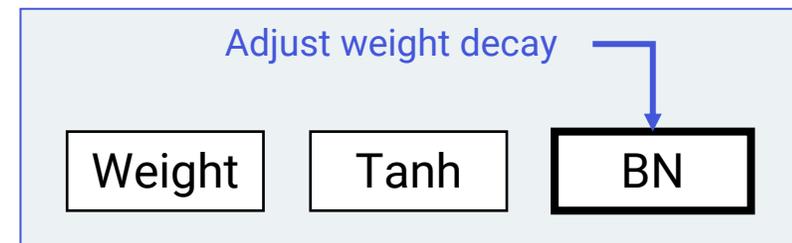
Asymmetric saturation incurs sparsity by a zero mean shifting in normalization.



The Effect of Asymmetry and Sparsity on Accuracy

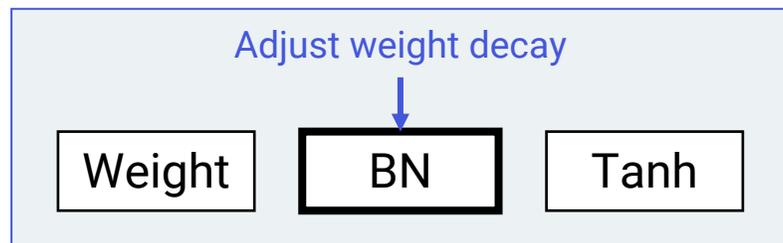
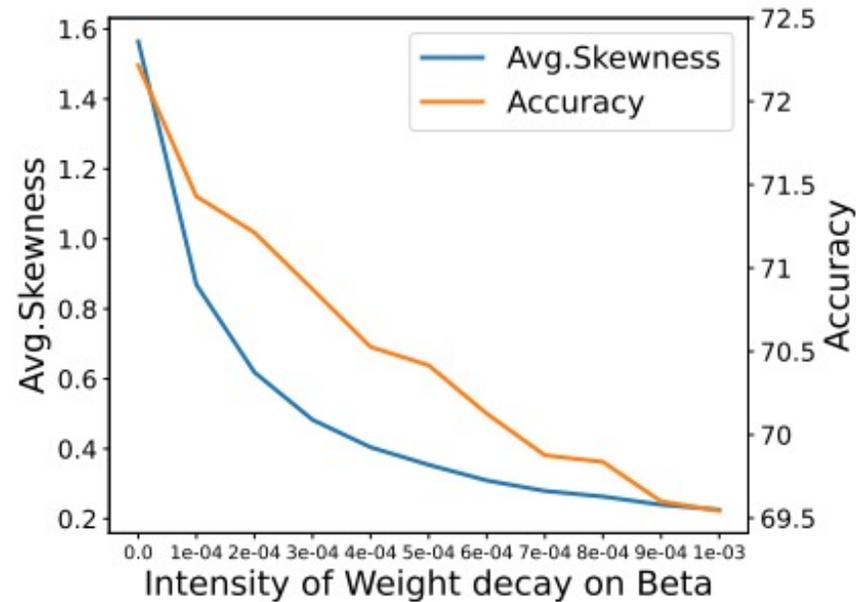


Convention order

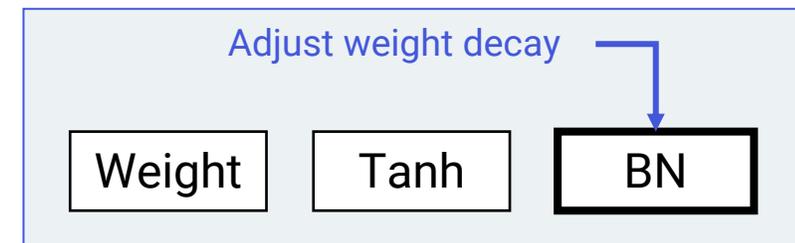
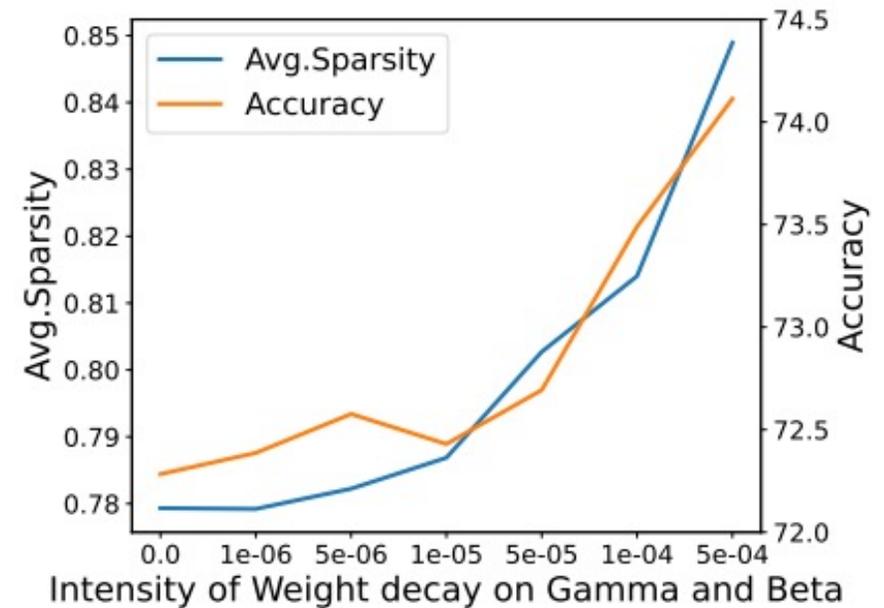


Swap order

The Effect of Asymmetry and Sparsity on Accuracy



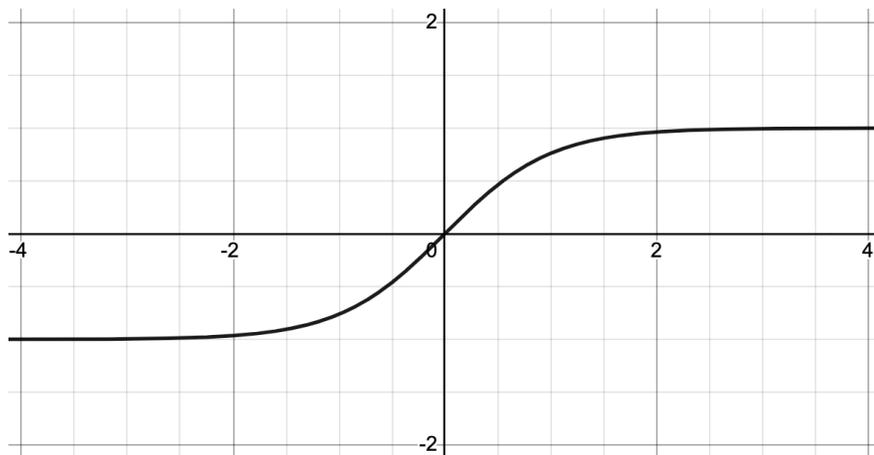
Convention order



Swap order

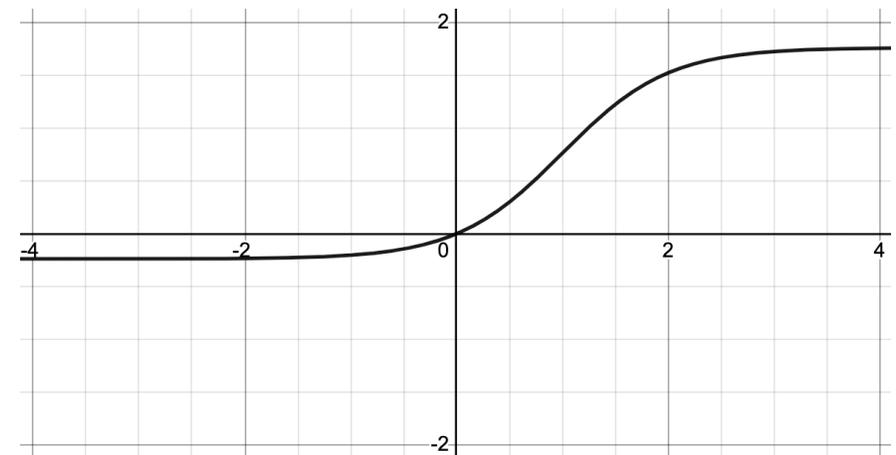
The shifted Tanh introduce asymmetric and sparse activation easily.

$$\tanh(x + \tau) - \tanh(\tau)$$



$$\tau = 0$$

Original Tanh

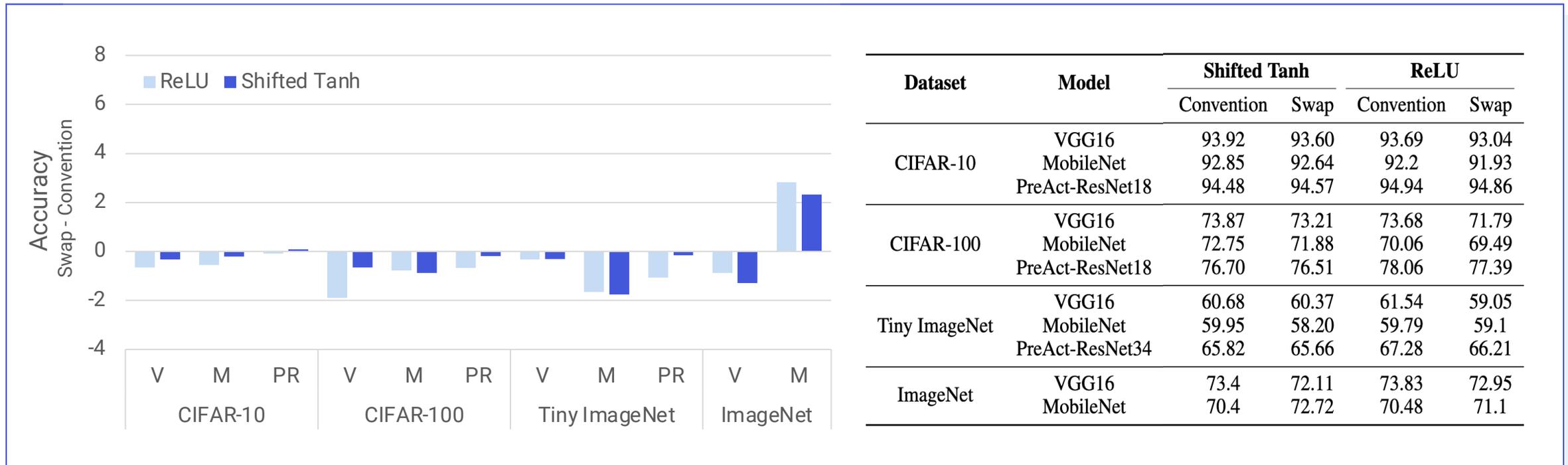


$$\tau = -1$$

Shifted Tanh

The properties of the shifted Tanh

- It shows improved accuracy comparable with the ReLU model.
- The accuracy discrepancy between orders decreased.



Other Bounded Activation Functions

- The Swap model with other bounded functions outperforms the Convention model.
-
- Softsign, which is a slower approach to its asymptotes than Tanh, underperforms Tanh on the Swap order, even if it performs better on the Convention order.

Activation functions	Order		Δ avg. Skewness (Swap - Convention)
	Convention	Swap	
Tanh	69.5	74.11	2.38
Softsign	70.01	73.65	1.28
LeCun Tanh	67.82	74.46	1.90

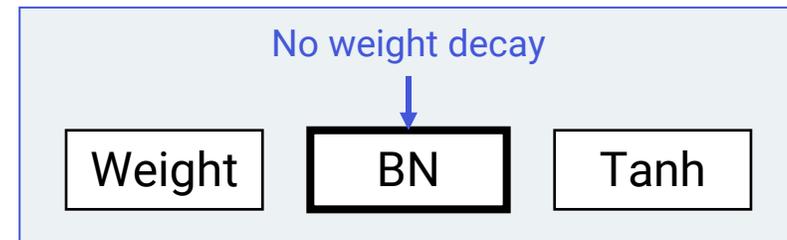
Other Bounded Activation Functions

- The Swap model with other bounded functions outperforms the Convention model.
-
- Softsign, which is a slower approach to its asymptotes than Tanh, underperforms Tanh on the Swap order, even if it performs better on the Convention order.

Activation functions	Order		Δ avg. Skewness (Swap - Convention)
	Convention	Swap	
Tanh	69.5	74.11	2.38
Softsign	70.01	73.65	1.28
LeCun Tanh	67.82	74.46	1.90

Dominance Between Asymmetry and Sparsity

- The NWDBN model with encouraged asymmetry outperforms the Convention model even if the sparsity is decreased.



NWDBN

Configurations	Accuracy	Avg. Skewness	Avg. Sparsity
Convention	69.5	0.254	0.718
NWDBN	72.22	0.718	0.288

Thank You!

npclinic3@gmail.com