

1 A Limitation

2 While this paper contains an in-depth analysis of Tanh, it does not propose a function that is superior
 3 to existing activation functions. Generally, it is known that Tanh does not perform better than ReLU,
 4 but it's hard to find a study that deeply analyzes the reasons for this. To our knowledge, this paper is
 5 the most comprehensive analysis of the performance degradation when Tanh is used in a conventional
 6 manner. Although it suggests strategies to utilize Tanh more effectively, achieving performance
 7 nearly on par with ReLU, it fundamentally does not assert that Tanh surpasses ReLU in terms of
 8 effectiveness.

9 B Accuracy of the Shifted Tanh According to τ

10 In this section, we focus on the hyperparameter τ in the shifted Tanh function. We carry out
 11 experiments using VGG16_11, MobileNet, and PreAct-ResNet models trained on CIFAR-10 and
 12 CIFAR-100 datasets, testing different values of the parameter τ (-1.5, -1.2, -1.0, -0.8, and -0.5). The
 13 performance results corresponding to each of these τ values are presented in Table B.1. In these
 14 experiments, all hyperparameters, except for τ , are maintained at their best-averaged accuracy settings
 15 based on each model trained on the CIFAR dataset. Our aim is to identify the effective τ value
 16 without pushing all inputs into the saturation zone due to excessive asymmetry. The best-performance
 17 τ is identified by averaging the accuracy across different models for each τ and selecting the one
 18 with the highest performance. In our case, -1 is the best.

Table B.1: The results of shifted Tanh models on various τ values.

τ	-1.5	-1.2	-1.0	-0.8	-0.5
Accuracy	80.56	83.86	84.15	84.01	83.58

19 C Metric

20 C.1 Saturation and Skewness on Different Distribution Patterns

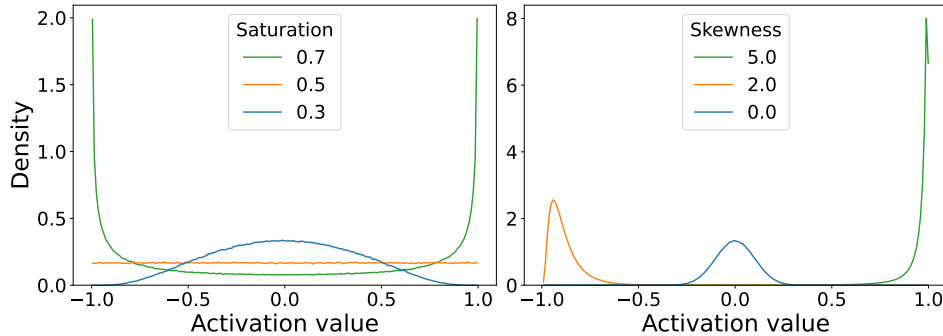


Figure C.1: The Saturation (left) and Skewness (right) values on different distributions.

21 Our saturation metric is designed to quantify the concentration of elements to the maximum absolute
 22 output values. The metric has bounds from 0 to 1, and it equates to 0 when all the elements are 0 in
 23 the Tanh case. Conversely, it grows larger as elements tend towards the maximum limit of the output
 24 range. For example, the saturation metric is evaluated at 0.5 in a uniform distribution. For saturation
 25 values of the Tanh output on the different Gaussian distributions, refer to Figure C.1 (left).

26 Skewness is a measure that quantifies the asymmetry within a distribution. Symmetric distribution
 27 will yield the Skewness of 0, and the Skewness increases with the rise in asymmetry. Notably, we
 28 calculate the absolute skewness value in our asymmetry metric, ensuring that the metric remains

29 unaffected by the direction of skewness. For a better understanding of how this measure applies to
 30 various distributions, Figure C.1 (right) provides Skewness on each distribution.

31 C.2 Impacts of Mean and Standard Deviation Variations on Tanh Asymmetry

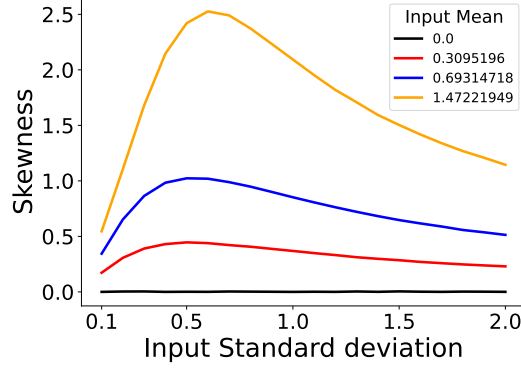


Figure C.2: Skewness value on the different means and standard deviations of the Gaussian distribution for Tanh input.

32 The mean and variance of input distribution on Tanh affect the asymmetry of Tanh output. The
 33 Skewness of Tanh output on various mean and standard deviation can be found in Figure C.2. In the
 34 Skewness of varied mean distribution, the Skewness is increased on the mean increase. However, the
 35 maximum Skewness does not align with the standard deviation increases. The Skewness decreases
 36 not only the small input standard deviation but also the large input standard deviation. Additionally,
 37 in the same mean condition, a decrease in standard deviation from the maximum skewness point
 38 more rapidly decreases the Skewness than an increase in standard deviation.

39 C.3 The Proof of Demonstrating the Sparsity Metric Criteria

40 This section examines the adequacy of our sparsity metric against the heuristic criteria for sparsity
 41 measures defined by Hurley & Rickard [1]. Our sparsity metric meets five of these criteria: Scaling,
 42 Rising Tide, Cloning, Bill Gates, and Babies. Without loss of generality, we consider the absolute
 43 operation as pre-applied to the metric input, according to Hurley & Rickard [1], and set the number
 44 of channels as one for simplicity. Thus, our sparsity metric is represented as $f : \mathbf{x} \mapsto s$, where
 45 $\mathbf{x} \in \mathbb{R}^{m+}$ is the vector of absolute values of one channel in BN output, i.e., $[|B_{i,0}|, |B_{i,1}|, \dots, |B_{i,m}|]$,
 46 and $x^{\max+}$ signifies the maximum value in \mathbf{x} , i.e., $x^{\max+} = \max_i x_i$. Below are the five criteria our
 47 sparsity metric meets, along with proof for each.

48 **Proof of Scaling.** Sparsity is scale invariant. That is, for any scalar $\alpha \in \mathbb{R}$, where $\alpha > 0$, the function
 49 f satisfies:

$$f(\alpha\mathbf{x}) = f(\mathbf{x}). \quad (1)$$

50 The invariance of scale means that sparsity considers relative differences rather than absolute magni-
 51 tudes.

52 *Proof.* If we scale the vector \mathbf{x} by α , we also scale the maximum positive value $x^{\max+}$ by α .
 53 Therefore, we have:

$$\begin{aligned} f(\alpha\mathbf{x}) &= 1 - \frac{\sum_{i=1}^m \alpha x_i}{m \alpha x^{\max+}} \\ &= 1 - \frac{\sum_{i=1}^m x_i}{m x^{\max+}} \\ &= f(\mathbf{x}). \end{aligned}$$

54 This demonstrates that the sparsity of the vector \mathbf{x} remains the same when the vector is scaled by any
 55 positive scalar α . \square

56 **Proof of Rising Tide.** Adding a constant to each element decreases sparsity. Formally, for any scalar
 57 $\alpha \in \mathbb{R}$, where $\alpha > 0$, the function f satisfies:

$$f(\alpha + \mathbf{x}) < f(\mathbf{x}).$$

58 We exclude the case, as mentioned in [1], where all elements of \mathbf{x} are the same.

59 This property also indicates that sparsity increases as more values approach zero.

60 *Proof.* We begin by noting the sparsity of the vector $\mathbf{x} + \alpha$:

$$f(\mathbf{x} + \alpha) = 1 - \frac{\sum_{i=1}^m x_i + m\alpha}{mx^{\max} + m\alpha}.$$

61 Assume for contradiction that $f(\mathbf{x} + \alpha) \geq f(\mathbf{x})$. This leads to the following inequalities:

$$mx^{\max} + > \sum_{i=1}^m x_i \implies \frac{\sum_{i=1}^m x_i + m\alpha}{mx^{\max} + m\alpha} > \frac{\sum_{i=1}^m x_i}{mx^{\max} +}, \quad (2)$$

$$1 - \frac{\sum_{i=1}^m x_i + m\alpha}{mx^{\max} + m\alpha} < 1 - \frac{\sum_{i=1}^m x_i}{mx^{\max} +}, \quad (3)$$

$$(4)$$

62 which contradicts the initial assumption, thus implying that

$$f(\alpha + \mathbf{x}) < f(\mathbf{x}),$$

63 as desired. □

64 **Proof of Cloning.** Sparsity is invariant under cloning.

65 That is, for a vector \mathbf{x} , the function f satisfies:

$$f(\mathbf{x}) = f(\mathbf{x}||\mathbf{x}) = f(\mathbf{x}||\mathbf{x}||\mathbf{x}) = \dots = f(\mathbf{x}||\mathbf{x}||\dots||\mathbf{x}),$$

66 where $||$ denotes concatenation, such that $\mathbf{x}||\mathbf{x} = [x_1, x_2, \dots, x_m, x_1, x_2, \dots, x_m]$.

67 The principle of cloning implies that sparsity remains the same even if a set of values is replicated.

68 *Proof.* Define a function $g : \mathbf{x}, \delta \mapsto \mathbf{y}$, where $\delta \in \mathbb{Z}^+$ and $\mathbf{y} \in \mathbb{R}^{m\delta}$ is the result of concatenating \mathbf{x}
 69 δ times. We then have:

$$\begin{aligned} f(g(\mathbf{x}, \delta)) &= 1 - \frac{\delta \sum_{i=1}^m x_i}{\delta mx^{\max} +} \\ &= 1 - \frac{\sum_{i=1}^m x_i}{mx^{\max} +} \\ &= f(\mathbf{x}). \end{aligned}$$

70 This shows that the sparsity of \mathbf{x} remains the same even when it is concatenated with itself δ times. □

71 **Proof of Bill Gates.** As one individual element becomes infinitely large, the sparsity increases.
 72 Formally, for every i , there exists $\rho > 0$, such that for all $\alpha > 0$:

$$f([x_1, \dots, x_i + \rho + \alpha, \dots, x_m]) > f([x_1, \dots, x_i + \rho, \dots, x_m])$$

73 The implication is that if a single value grows without bounds, the sparsity also increases indefinitely.

74 *Proof.* We first choose a sufficiently large ρ such that $\forall i, x_i + \rho > x^{\max} +$.

75 Assume for contradiction that $S([x_1, \dots, x_i + \rho + \alpha, \dots, x_m]) \leq S([x_1, \dots, x_i + \rho, \dots, x_m])$.

76 This implies:

$$\begin{aligned}
1 - \frac{\sum_{k=1}^m x_k + \rho + \alpha}{m(x_i + \rho + \alpha)} &\leq 1 - \frac{\sum_{k=1}^m x_k + \rho}{m(x_i + \rho)} \\
\Leftrightarrow \frac{\sum_{k=1}^m x_k + \rho + \alpha}{m(x_i + \rho + \alpha)} &\geq \frac{\sum_{k=1}^m x_k + \rho}{m(x_i + \rho)} \\
\Leftrightarrow \frac{\sum_{k \neq i} x_k + x_i + \rho + \alpha}{x_i + \rho + \alpha} &\geq \frac{\sum_{k \neq i} x_k + x_i + \rho}{x_i + \rho} \\
\Leftrightarrow \frac{\sum_{k \neq i} x_k}{x_i + \rho + \alpha} &\geq \frac{\sum_{k \neq i} x_k}{x_i + \rho} \\
\Leftrightarrow \frac{1}{x_i + \rho + \alpha} &\not\geq \frac{1}{x_i + \rho},
\end{aligned}$$

77 Leading to a contradiction. Therefore,

$$f([x_1, \dots, x_i + \rho + \alpha, \dots, x_m]) > f([x_1, \dots, x_i + \rho, \dots, x_m]).$$

78

□

79 **Proof of Babies.** Adding a new element of zero increases sparsity. Formally, for a vector \mathbf{x} , the
80 function f satisfies:

$$f(\mathbf{x}||0) > f(\mathbf{x}) \quad (5)$$

81 Concatenating zero elements to the existing values increases the relative difference to the other values,
82 which increases sparsity.

Proof.

$$1 - \frac{\sum_{k=1}^m x_k}{m+1} > 1 - \frac{\sum_{k=1}^m x_k}{m}$$

83 We start by subtracting one from both sides of the inequality:

$$-\frac{\sum_{k=1}^m x_k}{m+1} > -\frac{\sum_{k=1}^m x_k}{m}.$$

84 The negative sign can be removed by reversing the inequality:

$$\frac{\sum_{k=1}^m x_k}{m+1} < \frac{\sum_{k=1}^m x_k}{m}.$$

85 Multiplying both sides of the inequality by $m(m+1)$ (since m is a positive integer, $m+1$ is also
86 positive, and the inequality sign will not change), we get:

$$m \sum_{k=1}^m x_k < (m+1) \sum_{k=1}^m x_k.$$

87 Subtracting $m \sum_{k=1}^m x_k$ from both sides of the inequality, we obtain:

$$0 < \sum_{k=1}^m x_k.$$

88 Since x_k is not a negative value, the sum is also non-negative, which verifies the inequality. Therefore,
89 the original inequality is proven. □

90 D Deciding Model Configurations for Investigation

91 D.1 Investigation of Depth Impact on Accuracy in VGG16 Models

Table E.2: Experimental results of shortened VGG16 models with the Swap order for CIFAR-100. The number of removed convolution layers in the VGG16_n model is the difference between 16 and n.

	VGG16	VGG16_15	VGG16_14	VGG16_13	VGG16_12	VGG16_11	VGG16_10	VGG16_9	VGG16_8
Accuracy	72.17	73.02	73.48	73.85	73.76	73.92	72.57	70.91	70.69

Table E.3: Experimental comparison between the official VGG11 and our VGG16_11 models trained by CIFAR-100.

Models	Order	
	Convention	Swap
VGG11	64.55	69.94
VGG16_11	69.5	74.11

92 To identify a model for focused analysis using the CIFAR dataset, we examined various VGG16
 93 variants. This examination progressively removes convolution layers from the end towards the front
 94 of VGG16. The result can be seen in Table E.2. We observed an increase in accuracy until peaking at
 95 the VGG16_11 model, followed by a decline. Although a VGG11 model has already been proposed
 96 in Simonyan & Zisserman [2], the validation accuracy of VGG16_11 is significantly higher than
 97 VGG11 on CIFAR-100. The results can be seen in Table E.3

98 E Training Hyperparameters

99 We sweep the learning rate and weight decay hyperparameter. The learning rate was 0.1 and 0.01.
 100 For CIFAR and Tiny-ImageNet datasets, we trained models with a batch size of 128, and the learning
 101 rate was reduced by one-tenth at 100 and 150 of the total 200 epochs, and we swept four weight
 102 decay of 0.005, 0.001, 0.0005, and 0.0001. For ImageNet datasets, we trained models with a batch
 103 size of 256, and the learning rate was reduced by one-tenth at 30 and 60 of the total 100 epochs, and
 104 we swept three weight decay of 0.001, 0.0005, and 0.0001. We chose the best averaged-accuracy
 105 model for the three random seeds and averaged the values of these three models for all measurements
 106 for analysis. Because of the computation issue, we only use one seed for the ImageNet dataset with
 107 early stopping. The hyperparameters of VGG, MobileNet, and PreAct-ResNet are shown in E.1, E.2,
 108 E.3, respectively.

Table F.1: Hyperparameters of VGG11 with Tanh

	Convention	Swap	NWDBN
Training Epochs	200	200	200
Learning Rate	0.1	0.1	0.1
Learning Rate Drop	100, 150	100, 150	100, 150
Weight Decay	0.0005	0.0005	0.0005
Batch Size	128	128	128

Table F.2: Hyperparameters of VGG16 with Tanh

	Convention				Swap			
	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet
Training Epochs	200	200	200	90	200	200	200	90
Learning Rate	0.1	0.01	0.01	0.01	0.01	0.1	0.01	0.01
Learning Rate Drop	100, 150	100, 150	100, 150	30, 60	100, 150	100, 150	100, 150	30, 60
Weight Decay	0.0001	0.0005	0.001	0.0001	0.001	0.0005	0.001	0.001
Batch Size	128	128	128	256	128	128	128	256

Table F.3: Hyperparameters of VGG16 with the shifted Tanh

	Convention				Swap			
	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet
Training Epochs	200	200	200	90	200	200	200	90
Learning Rate	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
Learning Rate Drop	100, 150	100, 150	100, 150	30, 60	100, 150	100, 150	100, 150	30, 60
Weight Decay	0.0001	0.0005	0.0001	0.0001	0.0005	0.0005	0.0005	0.0001
Batch Size	128	128	128	256	128	128	128	256

Table F.4: Hyperparameters of VGG16 with the ReLU

	Convention				Swap			
	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet
Training Epochs	200	200	200	90	200	200	200	90
Learning Rate	0.01	0.01	0.1	0.1	0.01	0.01	0.01	0.01
Learning Rate Drop	100, 150	100, 150	100, 150	30, 60	100, 150	100, 150	100, 150	30, 60
Weight Decay	0.001	0.005	0.0001	0.0001	0.001	0.005	0.001	0.0005
Batch Size	128	128	128	256	128	128	128	256

Table F.5: Hyperparameters for MobileNet with Tanh

	Convention				Swap			
	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet
Training Epochs	200	200	200	90	200	200	200	90
Learning Rate	0.1	0.1	0.01	0.1	0.1	0.1	0.01	0.1
Learning Rate Drop	100, 150	100, 150	100, 150	30, 60	100, 150	100, 150	100, 150	30, 60
Weight Decay	0.0001	0.0005	0.005	0.0001	0.0001	0.0005	0.005	0.0001
Batch Size	128	128	128	256	128	128	128	256

Table F.6: Hyperparameters for MobileNet with the Shifted Tanh

	Convention				Swap			
	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet
Training Epochs	200	200	200	90	200	200	200	90
Learning Rate	0.1	0.1	0.01	0.1	0.1	0.1	0.01	0.01
Learning Rate Drop	100, 150	100, 150	100, 150	30, 60	100, 150	100, 150	100, 150	30, 60
Weight Decay	0.0001	0.0005	0.005	0.0001	0.0005	0.0005	0.005	0.0005
Batch Size	128	128	128	256	128	128	128	256

Table F.7: Hyperparameters for MobileNet with ReLU

	Convention				Swap			
	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet	ImageNet
Training Epochs	200	200	200	90	200	200	200	90
Learning Rate	0.01	0.01	0.1	0.01	0.01	0.01	0.1	0.1
Learning Rate Drop	100, 150	100, 150	100, 150	30, 60	100, 150	100, 150	100, 150	30, 60
Weight Decay	0.001	0.005	0.0001	0.0001	0.001	0.005	0.001	0.0001
Batch Size	128	128	128	256	128	128	128	256

Table F.8: Hyperparameters for PreAct-ResNet-18/34 with Tanh. PreAct-ResNet-18 is for the CIFAR dataset, and PreAct-ResNet-34 is for the Tiny ImageNet dataset.

	Convention			Swap		
	CIFAR-10	CIFAR-100	Tiny ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet
Training Epochs	200	200	200	200	200	200
Learning Rate	0.1	0.1	0.1	0.1	0.1	0.1
Learning Rate Drop	100, 150	100, 150	100, 150	100, 150	100, 150	100, 150
Weight Decay	0.0001	0.0001	0.0001	0.0005	0.0005	0.0005
Batch Size	128	128	128	128	128	128

Table F.9: Hyperparameters for PreAct-ResNet-18/34 with the Shifted Tanh. PreAct-ResNet-18 is for the CIFAR dataset, and PreAct-ResNet-34 is for the Tiny ImageNet dataset.

	Convention			Swap		
	CIFAR-10	CIFAR-100	Tiny ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet
Training Epochs	200	200	200	200	200	200
Learning Rate	0.1	0.1	0.01	0.1	0.1	0.1
Learning Rate Drop	100, 150	100, 150	100, 150	100, 150	100, 150	100, 150
Weight Decay	0.0001	0.0005	0.001	0.0005	0.0005	0.0005
Batch Size	128	128	128	128	128	128

Table F.10: Hyperparameters for PreAct-ResNet-18/34 with ReLU. PreAct-ResNet-18 is for the CIFAR dataset, and PreAct-ResNet-34 is for the Tiny ImageNet dataset.

	Convention			Swap		
	CIFAR-10	CIFAR-100	Tiny ImageNet	CIFAR-10	CIFAR-100	Tiny ImageNet
Training Epochs	200	200	200	200	200	200
Learning Rate	0.01	0.01	0.1	0.01	0.01	0.1
Learning Rate Drop	100, 150	100, 150	100, 150	100, 150	100, 150	100, 150
Weight Decay	0.005	0.005	0.0005	0.005	0.005	0.0005
Batch Size	128	128	128	128	128	128

112 **References**

- 113 [1] Hurley, N. and Rickard, S. Comparing measures of sparsity. *IEEE Transactions on Information*
114 *Theory*, 55(10):4723–4741, 2009.
- 115 [2] Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image
116 recognition. *arXiv preprint arXiv:1409.1556*, 2014.