

## Stronger Normalization-Free Transformers

Mingzhi Chen<sup>1</sup> Taiming Lu<sup>1</sup> Jiachen Zhu<sup>2</sup> Mingjie Sun<sup>3</sup> Zhuang Liu<sup>1</sup>

<sup>1</sup>Princeton University <sup>2</sup>New York University <sup>3</sup>Carnegie Mellon University

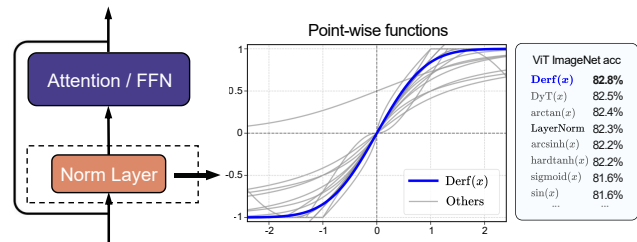
### Abstract

Although normalization layers have long been viewed as indispensable components of deep learning architectures, the recent introduction of Dynamic Tanh (DyT) [67] has demonstrated that alternatives are possible. The point-wise function DyT constrains extreme values for stable convergence and reaches normalization-level performance; this work seeks further for function designs that can surpass it. We first study how the intrinsic properties of point-wise functions influence training and performance. Building on these findings, we conduct a large-scale search for a more effective function design. Through this exploration, we introduce  $\text{Derf}(x) = \text{erf}(\alpha x + s)$ , where  $\text{erf}(x)$  is the rescaled Gaussian cumulative distribution function, and identify it as the most performant design. *Derf* outperforms *LayerNorm*, *RMSNorm*, and *DyT* across a wide range of domains, including visual recognition and generation, speech representation, and DNA sequence modeling. Our analysis also suggests that the performance gains of *Derf* largely stem from its improved generalization rather than stronger fitting capacity. Its simplicity and stronger performance make *Derf* a practical choice for normalization-free Transformer architectures. Our code is available at this [link](#).

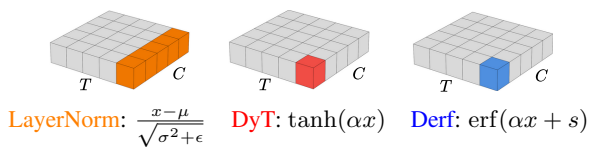
### 1. Introduction

Normalization layers have become a critical component in modern deep neural networks. Since the invention of Batch Normalization [25], more and more variants have been developed to adapt normalization to various architectures and model modalities [2, 43, 54, 56, 63]. By regulating the distribution of intermediate activations, normalization layers have long demonstrated their strong capability in stabilizing training and accelerating model convergence [7, 44].

Due to the inherent formulation of normalization layers, they heavily rely on activation statistics during training. This introduces additional memory access and synchronization overhead [11, 61, 63]. Moreover, some normalization methods are highly sensitive to batch size, and inappropriate batch settings can lead to unstable training [29, 48, 56]. These issues motivate recent efforts to de-



a) We search for different shapes of norm layer replacement.



b) Formulation of LayerNorm (LN), DyT, and Derf (ours).

method	ViT acc (↑)	DiT FID (↓)	DNA acc (↑)
LN	82.3%	45.91	86.9%
DyT	82.5%	45.66	86.9%
Derf	<b>82.8%</b>	<b>43.94</b>	<b>87.3%</b>

c) Performance across domains.

Figure 1. We introduce **Dynamic erf (Derf)**, a **point-wise function to replace normalization layers, which outperforms other normalization-based and -free designs across domains**. (a) Evaluating all feasible normalization replacements, we identify and introduce *Derf* as the strongest choice. (b) *LayerNorm*, *DyT* [67], and *Derf* operate in fundamentally different dimensions. (c) Across Imagenet-1K classification and generation, and DNA modeling, *Derf* consistently outperforms *LayerNorm* and *DyT*.

velop normalization-free methods. Among these attempts, Dynamic Tanh [67], an S-shaped point-wise function, has emerged as a simple yet effective drop-in replacement for normalization layers. This work has established the foundation for point-wise functions that match the performance of normalization layers, yet functions that can surpass them remain unexplored. In this work, we aim to discover point-wise functions that outperform normalization layers to push toward stronger Transformer architectures [16, 55].

We first systematically study how the intrinsic properties of point-wise functions affect the training dynam-

ics and final performance. Specifically, we focus on four fundamental and representative function properties: *zero-centeredness*, *boundedness*, *center sensitivity*, and *monotonicity*. Each property is independently examined through controlled experiments on a diverse set of point-wise functions to assess its impact on the training result. This analysis isolates a subset of point-wise functions as effective normalization replacements and yields a concrete design principle for normalization-free Transformers.

Guided by these principles, we identify a set of promising functions that have the potential to surpass normalization layers. Within this set, we empirically search for the optimal designs, among which Dynamic erf (Derf) emerges as a simple yet the most performant function (Fig. 1a). Derf augments  $\text{erf}(x)$  with learnable parameters, where the error function  $\text{erf}(x)$  is an S-shaped, rescaled cumulative distribution of a standard Gaussian around zero.

We evaluate Derf spanning multiple modalities (vision, language, speech, and DNA sequences); covering various tasks (classification, generation, and sequence modeling), under different training paradigms (supervised and self-supervised). Across all these settings, Derf consistently surpasses LayerNorm, and Dynamic Tanh (Fig. 1b and Fig. 1c). To pinpoint the source of these gains, we measure the training loss in evaluation mode after optimization. Derf exhibits higher training loss than normalization-based models, indicating that its superior performance stems from stronger generalization rather than enhanced fitting capacity. Overall, our work demonstrates that well-designed point-wise functions can outperform normalization layers.

## 2. Background

**Normalization layers.** Normalization layers have become pivotal components of modern neural networks. Among the various normalization techniques, Batch Normalization (BN) [25], Layer Normalization (LN) [2], and Root Mean Square Normalization (RMSNorm) [63] are the three most widely used in deep learning models.

$$y = \gamma * \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta \quad (1)$$

All normalization methods adhere to a unified paradigm, formalized in Eq. (1), where activations within each group are centered and scaled by their mean  $\mu$  and standard deviation  $\sigma$  (with  $\epsilon$  for numerical stability) to maintain consistent scale and stable gradient flow. The main distinction among different normalization methods lies in how the activations are grouped when computing  $\mu$  and  $\sigma$ . For example, LN computes the statistics along the channel dimension for each token independently. Given a token representation  $x \in \mathbb{R}^C$ , the mean and variance are computed as Eq. (2), where  $C$  denotes the number of hidden features (channels). Due to its per-token normalization, LN is par-

ticularly well-suited for Transformer architectures, where activations across tokens exhibit diverse statistics.

$$\mu = \frac{1}{C} \sum_{k=1}^C x_k, \quad \sigma^2 = \frac{1}{C} \sum_{k=1}^C (x_k - \mu)^2, \quad (2)$$

**Point-wise functions.** The strong reliance of normalization layers on activation statistics has motivated further exploration of statistics-free methods [21, 22, 26, 67]. Among these approaches, point-wise functions [67] have emerged as simple yet effective alternatives to traditional normalization methods. Unlike normalization, a point-wise function applies the same parametric mapping  $f(x; \theta)$  to each activation independently. The parameters  $\theta$  are fixed or learned, rather than being computed from batch-, token-, or channel-level statistics. A recent study [67] introduces the Dynamic Tanh (DyT) function (Eq. (3)), where  $\alpha$  is a learnable parameter. This design is motivated by the observation that Layer Normalization often produces an S-shaped input-output mapping in practice. The saturating nature of the tanh function squashes extreme activations, thereby fulfilling a role analogous to the re-centering and re-scaling effects of normalization layers.

$$\text{DyT}(x) = \gamma * \tanh(\alpha x) + \beta \quad (3)$$

While DyT has shown similar performance to normalization layers across various Transformer-based models, a comprehensive analysis of the design space for these statistics-free operators remains missing. In this work, we target the optimal form of the point-wise function as normalization replacement. We identify the function properties crucial for convergence and performance, and we introduce Derf, a point-wise function outperforming normalization layers rather than merely matching their performance.

## 3. Functional Property Analysis

Training Transformers without normalization requires understanding the factors that make a point-wise function stable and effective as a replacement. In this section, we examine four essential properties: *zero-centeredness*, *boundedness*, *center sensitivity*, and *monotonicity* (see Fig. 2). These properties collectively characterize the fundamental shape of point-wise functions and their behavior on activations. By isolating the impact of each property, we explore its influence on optimization and final performance.

To investigate these properties, we replace each normalization layer with a point-wise function of the form:

$$y = \gamma \cdot f(\alpha x) + \beta, \quad (4)$$

where  $f(\cdot)$  denotes the chosen base function with learnable  $\alpha$  rescaling the input.  $\gamma$  and  $\beta$  are affine parameters, similar to those in normalization layers. We begin with three base functions:  $\tanh(x)$ ,  $\text{erf}(x)$ , and  $\arctan(x)$ . In

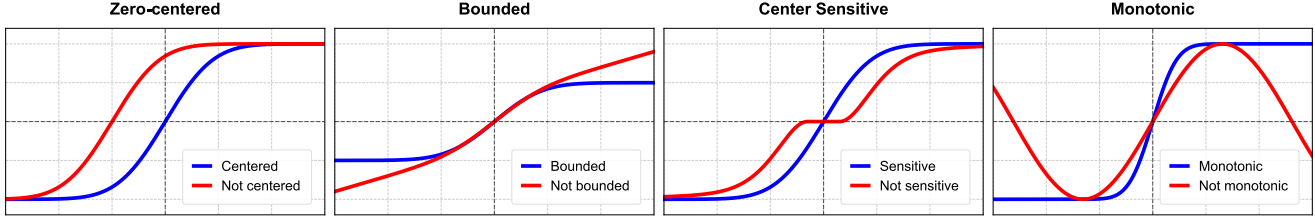


Figure 2. **Key properties of point-wise function.** The four properties: *zero-centeredness*, *boundedness*, *center sensitivity*, and *monotonicity* collectively characterize functional behavior on activations and influence training dynamics. Blue curves represent functions that satisfy each property, while red curves violate them.

subsequent experiments, we modify these functions with controlled transformations to examine the impact of each property. All experiments are conducted with ViT-Base [16], and top-1 accuracy on ImageNet-1K [15] is reported. In App. A, we provide more detailed training results.

### 3.1. Zero-centeredness

*Zero-centeredness* means that the function’s outputs are balanced around zero, with positive and negative values of similar magnitude and symmetry. Because normalization layers inherently recenter activations to the origin for stable gradients, keeping this property could reduce internal covariate shifts and promote smoother gradient flow during training.

**Setup.** Under the ViT setup, we manipulate the centering of the functions. For each base function, we consider two types of shifts: horizontal and vertical, defined in Eq. (5). In this form,  $\lambda_{\text{horiz}}$  and  $\lambda_{\text{vert}}$  respectively denote the magnitudes of horizontal and vertical shifts. For both types of shifts, we vary  $\lambda$  over  $\{\pm\frac{1}{2}, \pm 1, \pm 2\}$  to examine how increasing deviation from zero-centeredness affects the function’s behavior. All other training settings remain unchanged.

$$f_{\text{horiz}}(x) = f(x + \lambda_{\text{horiz}}), \quad f_{\text{vert}}(x) = f(x) + \lambda_{\text{vert}}, \quad (5)$$

**Results.** As shown in Tab. 1, the results are consistent across different base functions: for horizontal shifts, performance remains largely comparable to the zero-centered base function when  $|\lambda_{\text{horiz}}| \leq 0.5$ . However, as  $|\lambda_{\text{horiz}}|$  increases, performance gradually degrades, and training diverges when  $|\lambda_{\text{horiz}}| \geq 2$ . Similarly, vertical shifts consis-

function	shift type	-2	-1	-0.5	-0.1	$\lambda = 0$	+0.1	+0.5	+1	+2
$\text{erf}(x)$	horizontal	×	82.0%	82.5%	82.6%	82.6%	<b>82.7%</b>	82.5%	82.1%	×
	vertical	×	81.8%	82.3%	82.4%	<b>82.6%</b>	82.5%	82.3%	81.6%	×
$\tanh(x)$	horizontal	×	82.1%	82.5%	<b>82.6%</b>	82.5%	<b>82.6%</b>	82.4%	82.2%	×
	vertical	×	81.5%	81.9%	82.4%	<b>82.5%</b>	82.3%	81.9%	81.4%	×
$\arctan(x)$	horizontal	×	81.9%	82.3%	82.3%	82.3%	<b>82.4%</b>	82.2%	82.0%	×
	vertical	×	81.4%	81.9%	82.2%	<b>82.3%</b>	<b>82.3%</b>	82.0%	81.2%	×

Table 1. **Results of zero-centeredness on ViT-Base.** Horizontal shift corresponds to modifying the input as  $f(\alpha x \pm \lambda)$ , while vertical shift adds or subtracts a constant to the output as  $f(\alpha x) \pm \lambda$ . “×” indicates training failure.

tently lead to a decline in performance as  $|\lambda_{\text{vert}}|$  grows with training failure once  $|\lambda_{\text{vert}}| \geq 2$ . These results show that *zero-centeredness* is a requirement for effective training.

### 3.2. Boundedness

*Boundedness* refers to the property of a function whose output is constrained within a finite range. Formally, a function  $f(\cdot)$  is bounded if there exist constants  $a, b \in \mathbb{R}$  such that  $a \leq f(x) \leq b$  for all  $x$  in its domain. This ensures that activations remain finite and do not accumulate variance across layers. Unbounded functions, in contrast, may induce signal explosion and gradient instability.

**Setup.** Under the same ViT setup, we study the role of *boundedness* with two methods. Firstly, we select three inherently unbounded S-shaped functions (e.g.,  $\text{arcsinh}(x)$ ) and compare them with their clamped versions shown in Eq. (6), where  $f_u(x)$  denotes the unbounded function, and  $\lambda$  is a chosen value specifying the clipping range.

$$y = \text{clip}(f_u(x), -\lambda_u, \lambda_u), \quad (6)$$

Secondly, we gradually transition bounded functions (e.g.,  $\text{erf}(x)$ ) toward unbounded linear form:

$$y = (1 - \lambda)f_b(x) + \lambda b x, \quad \lambda_b \in (0, 1). \quad (7)$$

where  $f_b$  denotes a bounded point-wise function, and  $\lambda$  controls how quickly the function becomes unbounded. We vary  $\lambda_u$  over  $\{0.5, 0.8, 1.0, 2.0, 3.0, 5.0\}$  in the first method and  $\lambda_b$  over  $\{0.01, 0.1, 0.5\}$  for the second. The original unmodified function is also included as a baseline.

$\lambda_u$	$\text{arcsinh}(x)$	$\text{logsign}(x)$	$\text{linear}(x)$
—	82.2%	82.2%	×
0.5	82.3%	<b>82.4%</b>	82.1%
0.8	82.3%	<b>82.4%</b>	<b>82.2%</b>
1.0	<b>82.4%</b>	<b>82.4%</b>	<b>82.2%</b>
2.0	<b>82.4%</b>	<b>82.4%</b>	82.1%
3.0	<b>82.4%</b>	82.3%	82.1%

Table 2. **Results of clamping for boundedness on ViT-Base.** “—” denotes the original unmodified function. “×” indicates training failure.

**Results.** For the first method, among the three unbounded functions in Tab. 2, only  $\text{arcsinh}(x)$  and  $\text{logsign}(x)$  converge effectively, while  $\text{linear}(x)$  does not. For the convergent functions, their clipped versions consistently outperform the unbounded baselines across all tested  $\lambda$  values. These results indicate that incorporating *boundedness* can improve optimization and result in better performance. For the second, as shown in Tab. 3, the results are consistent with clipping the intrinsic unbounded functions: the unbounded variant yields slightly lower accuracy than the bounded baseline.

$\lambda_b$	$\text{erf}(x)$	$\text{tanh}(x)$	$\text{arctan}(x)$	$\text{isru}(x)$
—	<b>82.6%</b>	<b>82.5%</b>	<b>82.3%</b>	<b>82.2%</b>
0.01	82.4%	82.4%	82.1%	<b>82.2%</b>
0.1	82.3%	82.3%	82.1%	82.1%
0.5	×	×	×	×

Table 3. **Results of removing boundedness on ViT-Base.** Performance drops with reduced boundedness. “—” denotes the original function without modification and “×” denotes training failure.

**Limitation of growth rate.** From Tab. 2 and Tab. 3, we observe that there is an upper limit on their acceptable growth rate. Large growth rates often lead to training failure. To determine this limit, we evaluate a family of inherently unbounded functions with varying growth rates, as illustrated in Fig. 3. Among them,  $\text{logquad}(x)$  exhibits the fastest growth that still allows training convergence (see Tab. 4). Functions with faster growth, such as  $\text{linear}(x)$  and  $\text{power23}(x)$ , tend to cause optimization divergence in the early stages of training. This failure occurs because rapidly growing functions fail to suppress variance effectively, leading to large gradient norms at the start of optimization.

$\text{logsign}(x)$	$\text{arcsinh}(x)$	$\text{logquad}(x)$	$\text{power23}(x)$	$\text{linear}(x)$
82.2%	82.2%	82.1%	×	×

Table 4. **Results of unbounded functions with different growth rates on ViT-Base.** Point-wise functions have a growth rate upper bound, with  $\text{logquad}(x)$  being the fastest function that still converges. “×” indicates training failure.

### 3.3. Center Sensitivity

We use *center sensitivity* to characterize how quickly a point-wise function becomes responsive to input variations around zero. Without *center sensitivity*, a function is locally flat around the origin, returning zero or near-zero over a finite interval. The region around zero is particularly important, as most activations tend to concentrate near the origin during training. Consequently, the responsiveness of a function in this area directly influences how effectively small signals can propagate through the network.

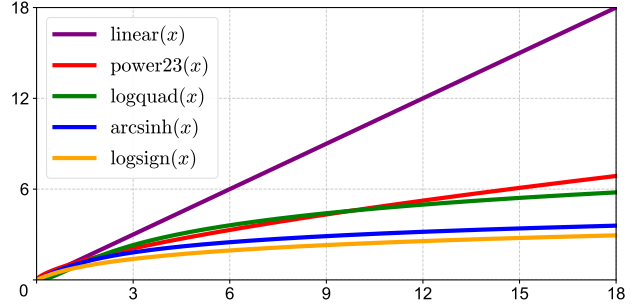


Figure 3. **Visualization of several unbounded point-wise functions with different growth rates** on the positive half-axis.  $\text{arcsinh}(x)$  refers to its standard analytical form. The remaining functions are defined as  $\text{linear}(x) = x$ ,  $\text{power23}(x) = \text{sign}(x) \cdot x^{\frac{23}{3}}$ ,  $\text{logsign}(x) = \text{sign}(x) \ln(|x| + 1)$ ,  $\text{smoothsign}(x) = \frac{x}{1+|x|}$ , and  $\text{logquad}(x) = \text{sign}(x) \ln(x^2 + 1)$ . Among them,  $\text{logquad}(x)$  shows the fastest growth that still ensures stable convergence.

**Setup.** Since *center sensitivity* is difficult to isolate independently, we approximate it using a controllable near-zero inactive region. Under the same ViT setup, we modify each base function to incorporate a symmetric flat region around the origin with a sensitivity scale  $\lambda > 0$  to control the extent of this region. Specifically, for inputs in the range  $x \in [-\lambda, \lambda]$ , we enforce  $f(x) = 0$  and smoothly shift the positive and negative parts outward for  $|x| > \lambda$  to ensure continuity at the boundaries. A smaller  $\lambda$  results in a narrower flat region and higher sensitivity near zero, while a larger  $\lambda$  leads to lower sensitivity. We vary  $\lambda$  over  $\{0.1, 0.5, 1.0, 2.0, 3.0\}$  across three base functions.

**Results.** As shown in Tab. 5, the best performance is achieved at  $\lambda = 0$ . As  $\lambda$  increases, the performance consistently degrades. This trend is not very clear when  $\lambda \leq 0.5$ , but once  $\lambda$  exceeds 1.0, the degradation becomes much more obvious. Finally, when  $\lambda \geq 3.0$ , the training process diverges at an early stage.

function	$\lambda = 0$	0.1	0.5	1.0	2.0	3.0
$\text{erf}(x)$	<b>82.6%</b>	82.5%	82.5%	82.1%	81.3%	×
$\text{tanh}(x)$	<b>82.5%</b>	<b>82.5%</b>	82.4%	82.1%	81.1%	×
$\text{arctan}(x)$	<b>82.3%</b>	<b>82.3%</b>	82.1%	81.8%	80.9%	×

Table 5. **Results of center sensitivity ( $\lambda$ ) on ViT-Base.** “×” indicates training failure. The best performance is achieved when no flat region is given, showing the importance of center sensitivity.

### 3.4. Monotonicity

*Monotonicity* ensures a function’s output consistently increases (or decreases) as the input increases, preserving the relative order of inputs throughout the transformation. Non-monotonic functions may disrupt the relative ordering of activations. Furthermore, since a non-monotonic function necessarily has regions where its derivative changes sign, it may also produce flipped gradient signals during training.

**Setup.** Each base function selected can serve as the monotonically increasing case, while its negated counterpart is defined as  $f_{\text{neg}}(x) = -f(x)$ , representing the monotonically decreasing variant. As non-monotonic comparisons, we include hump-shaped functions and oscillatory functions (see Fig. 5) to examine how violations of *monotonicity* influence the training performance. To control potential confounding factors, we rescale each function so that its output range matches that of the monotonic functions. After rescaling, all functions are aligned in terms of *zero-centeredness*, *boundedness*, and *center sensitivity*.

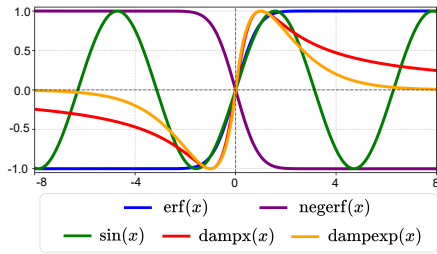


Figure 5. **Visualization of point-wise functions with different monotonicity behaviors.**  $\text{erf}(x)$  and  $\sin(x)$  refer to their standard form. The remaining functions are defined as  $\text{ngerf}(x) = -\text{erf}(x)$ ,  $\text{dampx}(x) = \frac{2x}{1+x^2}$ ,  $\text{dampexp}(x) = 2.72x \cdot e^{-|x|}$ .

**Results.** As shown in Tab. 6, both increasing and decreasing monotonic functions train stably and achieve high accuracy. In contrast, non-monotonic functions, whether hump-shaped or oscillatory, consistently perform worse than monotonic functions and lead to a clear drop in final accuracy. These results highlight *monotonicity* as a key property for point-wise functions to ensure effective learning.

function	$f(x)$	$f_{\text{neg}}(x)$	function	$f(x)$
$\text{erf}(x)$	82.6%	82.5%	$\sin(x)$	81.6%
$\tanh(x)$	82.5%	82.5%	$\text{dampx}(x)$	80.7%
$\arctan(x)$	82.3%	82.2%	$\text{dampexp}(x)$	81.2%

(a) Monotonic functions

(b) Non-monotonic functions

Table 6. **Results of monotonicity** on ViT-Base. Top-1 accuracy on ImageNet-1K for both monotonic and non-monotonic functions.

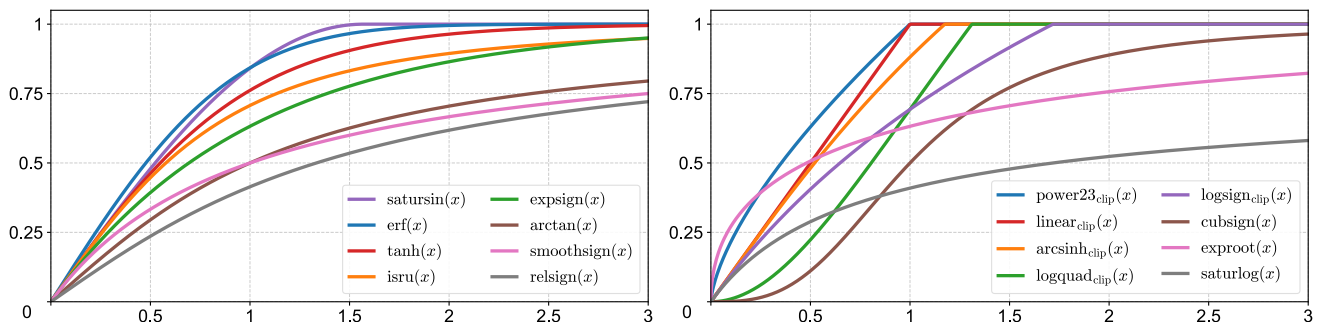


Figure 4. **Visualization of candidate point-wise functions on the positive half-axis.** All are self-symmetric with respect to the origin.

## 4. Function Search

From the previous section, we observe that functions that are near *zero-centered*, *bounded*, *center-sensitive* (responsive to input variations around zero), and *monotonic* tend to yield better optimization performance. Building upon these insights, we start to construct our function set from widely used scalar functions and cumulative distribution functions (CDFs), including polynomial, rational, exponential, logarithmic, and trigonometric forms. We then generate variants via simple transformations such as scaling, mirroring, rotation, and clipping. Functions that satisfy our four function properties after these transformations are retained as the candidate subset used in the search. For example, we transform the unbounded function  $\text{arcsinh}(x)$  by clipping it to the range  $[-1, 1]$ , limiting it to a finite range and conforming to all four principles. In App. B, we provide further details about how we obtain these candidate functions. Within this set, we evaluate their performance, and Derf emerges as the most effective function.

**Setup.** We conduct an empirical search on two representative vision architectures: Vision Transformer [16] and Diffusion Transformer [41]. Models are trained on ImageNet-1K [15] under their default training settings. For ViT, model performance is measured using top-1 accuracy on the ImageNet-1K validation set. For DiT, we follow the standard ImageNet reference batch evaluation and report the Fréchet Inception Distance (FID) as the metric.

**Formulation.** We quantitatively evaluate a set of functions under the constraint of our function properties, as illustrated in Fig. 4. Each function is instantiated in a unified form in Eq. (8), where  $f(\cdot)$  denotes a candidate point-wise function, with learnable parameter  $s$  and  $\alpha$  recentering and rescaling the input. The parameters  $\gamma$  and  $\beta$  follow the same role as in standard normalization layers. We introduce a learnable shift parameter  $s$ , as it improves the final performance to varying degrees across different functions. Detailed ablation results on the effect of  $s$  are provided in Sec. 7.1.

$$y = \gamma * f(\alpha x + s) + \beta, \quad (8)$$

function	top-1 acc $\uparrow$	FID $\downarrow$	
	ViT-Base	DiT-B/4	DiT-L/4
LayerNorm	82.3%	64.93	45.91
erf( $x$ )	<b>82.8%</b>	<b>63.23</b>	<b>43.94</b>
tanh( $x$ )	82.6%	63.71	45.48
satursin( $x$ )	82.6%	63.90	44.83
arcsinh <sub>clip</sub> ( $x$ )	82.5%	64.72	45.48
arctan( $x$ )	82.4%	67.07	46.62
smoothsign( $x$ )	82.4%	68.84	47.29
exproot( $x$ )	82.4%	65.20	46.91
logsign <sub>clip</sub> ( $x$ )	82.4%	65.59	46.34
reلسign( $x$ )	82.3%	68.42	48.33
isru( $x$ )	82.3%	65.72	45.93
linear <sub>clip</sub> ( $x$ )	82.3%	66.08	45.49
expsign( $x$ )	82.2%	64.85	45.82
logquad <sub>clip</sub> ( $x$ )	82.2%	65.92	47.12
power23 <sub>clip</sub> ( $x$ )	82.1%	66.11	46.47
saturllog( $x$ )	81.8%	68.23	47.44
cubsign( $x$ )	81.4%	70.22	49.16

Table 7. **Top-1 accuracy on ViT and FID DiT.** Different functions show noticeable differences in performance. Among all the point-wise functions and LayerNorm, erf( $x$ ) shows the best performance in both top-1 accuracy and FID. Visualization and analytical form of each function is included in App. B.

**Quantitative evaluation.** As shown in Tab. 7, even though these S-shaped functions appear highly similar in form, their empirical training results show noticeable differences in final performance. Among all the point-wise functions, erf( $x$ ) with the introduced transformations stands out as the best-performing function, consistently surpassing all other candidates and the baseline normalization layers.

## 5. Dynamic erf (Derf)

From the search, we identify erf( $x$ ) as the most performant point-wise function. The error function erf( $\cdot$ ) is closely related to the CDF of a standard Gaussian distribution. Specifically, erf( $x$ ) can be defined by Eq. (9). In our setup, erf( $x$ ) is in the form augmented with learnable parameters, which we introduce as Derf, **Dynamic erf**. Given an input tensor  $x$ , a Derf layer is defined in Eq. (10), where both the shift  $s$  and the scale  $\alpha$  are learnable scalars.  $\gamma$  and  $\beta$  are learnable per-channel vectors. To integrate Derf into a transformer-based architecture, we replace each normalization layer with a corresponding Derf layer. In particular, the pre-attention, the pre-FFN, and the final normalization layers are all substituted in a one-to-one manner, ensuring consistent incorporation of Derf across the entire model.

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt \quad (9)$$

$$\text{Derf}(x) = \gamma * \text{erf}(\alpha x + s) + \beta \quad (10)$$

**Parameter initialization.** We initialize  $\gamma$  to an all-one vector and  $\beta$  to an all-zero vector following the same strategy as in standard normalization layers. For the additional scalar parameters introduced by Derf, the scaling parameter  $\alpha$  is initialized to 0.5, while the shift parameter  $s$  is initialized to 0. Unless otherwise specified, these initialization settings are adopted throughout all experiments.

## 6. Experiments

We evaluate the effectiveness of Derf across various transformer-based and a few other modern architectures. For each model, we replace the normalization layers with DyT and Derf, following the standard training and evaluation protocols, as detailed in App. C. Across all tested architectures, Derf consistently achieves stronger performance over the baseline normalization methods and DyT. Besides each model’s default normalization, we also report results with other common normalization methods in App. D.

**Vision Transformers.** We train ViT-Base and ViT-Large models [16] on ImageNet-1K [15] using LayerNorm (LN), DyT, and Derf for comparison. Tab. 8 reports the top-1 classification accuracy. Compared to LN and DyT, Derf achieves clearly higher top-1 accuracy.

model	LN	DyT	Derf	$\Delta_{\text{LN}}$	$\Delta_{\text{DyT}}$
ViT-B	82.3%	82.5%	<b>82.8%</b>	$\uparrow 0.5\%$	$\uparrow 0.3\%$
ViT-L	83.1%	83.6%	<b>83.8%</b>	$\uparrow 0.7\%$	$\uparrow 0.2\%$

Table 8. **Top-1 accuracy on ImageNet-1K.** Derf achieves higher top-1 accuracy than both LN and DyT on different model sizes.

**Diffusion Transformers.** We train three Diffusion Transformer (DiT) [41] models on ImageNet-1K [15]. Consistent with the original DiT setup, the affine parameters in the normalization layers are retained for class conditioning across LN, DyT, and Derf. After training, we evaluate the FID scores using the standard ImageNet “reference batch” to measure image generation quality. In Tab. 9, Derf achieves a clear improvement in FID compared to both LN and DyT.

model	LN	DyT	Derf	$\Delta_{\text{LN}}$	$\Delta_{\text{DyT}}$
DiT-B/4	64.93	63.94	<b>63.23</b>	$\downarrow 1.70$	$\downarrow 0.71$
DiT-L/4	45.91	45.66	<b>43.94</b>	$\downarrow 1.97$	$\downarrow 1.72$
DiT-XL/2	19.94	20.83	<b>18.92</b>	$\downarrow 1.02$	$\downarrow 1.91$

Table 9. **Image generation quality (FID) on ImageNet.** Lower FID indicates better image generation quality. Derf achieves lower FID scores than both LN and DyT across all DiT model sizes.

**Speech models.** We train two wav2vec 2.0 Transformer models [4] on the LibriSpeech dataset [40] for speech representation learning. We report the final validation loss in Tab. 10. Compared to LayerNorm and DyT, Derf yields lower validation loss on different model sizes.

model	LN	DyT	Derf	$\Delta_{LN}$	$\Delta_{DyT}$
wav2vec 2.0 Base	1.95	1.95	<b>1.93</b>	$\downarrow 0.02$	$\downarrow 0.02$
wav2vec 2.0 Large	1.92	1.91	<b>1.90</b>	$\downarrow 0.02$	$\downarrow 0.01$

Table 10. **Speech pretraining validation loss on LibriSpeech.** Derf achieves lower validation loss than both LN and DyT across two wav2vec 2.0 models.

**DNA models.** For the long-range DNA sequence modeling task, we pretrain the HyenaDNA model [38] and the Caduceus model [45] using the human reference genome from [18]. Model evaluation is conducted on the GenomicBenchmarks dataset [19]. We report the averaged accuracy over all subtasks. As shown in Tab. 11, Derf surpasses both normalization layers and DyT in performance, demonstrating its robustness in genomic sequence modeling.

model	Norm	DyT	Derf	$\Delta_{Norm}$	$\Delta_{DyT}$
Hyena	85.2%	85.2%	<b>85.7%</b>	$\uparrow 0.5\%$	$\uparrow 0.5\%$
Caduceus	86.9%	86.9%	<b>87.3%</b>	$\uparrow 0.4\%$	$\uparrow 0.4\%$

Table 11. **DNA classification accuracy on GenomicBenchmarks,** averaged over each dataset. Each model is evaluated with its native normalization (LN for Heyna, RMSNorm for Caduceus); Derf consistently outperforms both normalization layers and DyT.

**Language models.** We pretrain a GPT-2 (124M) model on the OpenWebText dataset and report the validation loss in Tab. 12. For DyT and Derf, we additionally finetune the initialization of the learnable parameter  $\alpha$ . We observe that Derf achieves comparable performance to LN, while clearly outperforming DyT.

model	LN	DyT	Derf	$\Delta_{LN}$	$\Delta_{DyT}$
GPT-2	2.94	2.97	2.94	0.00	$\downarrow 0.03$

Table 12. **GPT2 validation loss on OpenWebText.** Derf matches LN’s performance while achieving lower validation loss than DyT.

### 6.1. Stronger Generalization or Better Fitting?

Given Derf’s superior performance, we aim to determine whether the gains arise from improved fitting capacity or stronger generalization. To this end, we compare the training loss of models respectively trained with normalization layers, DyT, and Derf. Since lower training loss indicates stronger fitting ability, this comparison helps us assess whether Derf improves training or enhances generalization.

**Setup.** We compute training losses across diverse architectures and scales. To measure fitting capacity fairly, we do not use the loss during optimization, which is confounded by stochastic regularization (e.g., stochastic depth [24]) and train-time augmentations. Instead, after training, we switch to evaluation mode, disable stochastic depth, adopt the test-time preprocessing pipeline, and compute the loss on the

training set. This yields a fair estimate of each model’s fitting capacity. In App. E, we provide the detailed procedure for computing the evaluation-mode training loss.

**Results.** Across all architectures and scales, both Derf and DyT result in higher training loss than normalization-based models, with Derf generally yielding slightly lower training loss than DyT, as shown in Tab. 13. This consistent pattern indicates that neither Derf nor DyT improves fitting capacity over normalization layers.

model	Norm	Derf	DyT
ViT-B	<b>0.2623</b>	<u>0.2681</u>	0.2714
ViT-L	<b>0.2034</b>	<u>0.2066</u>	0.2083
DiT-B	<b>0.1531</b>	<u>0.1533</u>	0.1535
DiT-L	<b>0.1501</b>	<u>0.1510</u>	0.1518
DiT-XL	<b>0.1432</b>	<u>0.1436</u>	0.1440
wav2vec 2.0 B	<b>1.8509</b>	<u>1.8821</u>	1.8946
wav2vec 2.0 L	<b>1.8241</b>	<u>1.8563</u>	1.8641
Hyena	<b>1.1297</b>	<u>1.1526</u>	1.1631
Caduceus	<b>0.8917</b>	<u>0.9129</u>	0.9203
GPT-2	<b>2.9478</b>	<u>2.9702</u>	2.9822

Table 13. **Evaluation-mode training loss of normalization layers (Norm), Derf, and DyT after optimization.** Bolded indicates the lowest loss, and underlined means the second-lowest loss. Across all model architectures, the training loss follows the relation: Norm < Derf < DyT. Both DyT and Derf exhibit higher training loss than normalization layers, while Derf achieves slightly lower loss than DyT.

**Discussion.** Despite the reduced fitting capacity, Derf achieves better performance across domains. We assume that these gains arise from both better generalization than normalization layers and stronger fitting capacity than DyT.

Firstly, point-wise functions promote stronger generalization. Although Derf yields higher training loss, it achieves superior downstream performance, indicating that its benefits stem not from improved fitting but from enhanced generalization. This difference likely originates from the contrasting operational principles between normalization layers and point-wise functions. Normalization layers adapt their transformation based on training statistics, allowing them to dynamically fit activation distributions throughout training. In contrast, point-wise functions are controlled by only a small set of learnable scalar parameters (e.g.,  $\alpha$  for DyT and  $\alpha, s$  for Derf) that do not adapt to activation statistics after training. They apply the same transformation regardless of activation distribution. This limited adaptability constrains overfitting and effectively serves as an implicit regularizer, leading to improved generalization.

Secondly, Derf exhibits stronger fitting power than DyT. It achieves lower training loss while retaining the implicit regularization of point-wise functions, combining higher fitting capacity with strong generalization to outperform both DyT and normalization-based models.

## 7. Analysis

### 7.1. Effect of $s$

**Removing  $s$ .** We investigate the effect of the learnable scalar parameter  $s$  by removing it from the point-wise function. As shown in Tab. 14, introducing this learnable shift consistently improves the overall training performance, and the degree of improvement varies across different functions. The stronger results of  $\text{erf}(x)$  over  $\text{tanh}(x)$  indicate that Derf surpasses DyT not only because of the shift  $s$ .

function	top-1 acc $\uparrow$		FID $\downarrow$	
	without $s$	with $s$	without $s$	with $s$
$\text{erf}(x)$	82.6%	<b>82.8%</b>	63.39	<b>63.23</b>
$\text{tanh}(x)$	82.5%	<b>82.6%</b>	63.94	<b>63.71</b>
$\text{satur}\sin(x)$	82.4%	<b>82.6%</b>	65.28	<b>63.90</b>
$\text{isru}(x)$	82.2%	<b>82.3%</b>	66.14	<b>65.72</b>
$\text{arctan}(x)$	82.3%	<b>82.4%</b>	67.41	<b>67.07</b>
$\text{arcsinh}_{\text{clip}}(x)$	82.4%	<b>82.5%</b>	65.19	<b>64.72</b>

Table 14. **Ablation study of  $s$ .** Top-1 accuracy on ViT-Base and FID score on DiT-B/4, comparing models with and without  $s$ .

**Scalar vs. vector  $s$ .** We further examine whether using a per-channel vector parameter instead of a scalar  $s$  leads to any performance improvement. As shown in Tab. 15, across all three point-wise functions, the choice between a scalar and a per-channel vector shows no significant impact on the final performance. Therefore, we adopt the scalar form of  $s$  for efficiency and simplicity during training.

function	vector		function	base	
	vector	scalar		base	large
$\text{erf}(x)$	<b>82.8%</b>	<b>82.8%</b>	$\text{tanh}(x)$	82.6%	83.6%
$\text{arctan}(x)$	<b>82.5%</b>	82.4%	$\text{tanh}(\varepsilon x)$	82.7%	83.7%
$\text{arcsinh}_{\text{clip}}(x)$	<b>82.5%</b>	<b>82.5%</b>	$\text{erf}(x)$	<b>82.8%</b>	<b>83.8%</b>

Table 15. **Top-1 accuracy of scalar vs. vector  $s$  on ViT-Base.** Table 16. **Top-1 accuracy of  $\tanh(\varepsilon x)$  on ViT.**

### 7.2. Approximating Derf

Given the superior performance of  $\text{erf}(x)$  over  $\text{tanh}(x)$ , we approximate  $\text{erf}(x)$  by scaling  $\text{tanh}(x)$  and examine whether this modification can lead to performance improvement. We introduce a fixed coefficient  $\varepsilon$  and use  $\text{tanh}(\varepsilon x)$ , where  $\varepsilon$  is obtained by minimizing the following objective:

$$\min_{\varepsilon} \int_{-\infty}^{+\infty} |\text{tanh}(\varepsilon x) - \text{erf}(x)| dx. \quad (11)$$

The optimal value is found to be  $\varepsilon \approx 1.205$ . In Tab. 16,  $\text{tanh}(\varepsilon x)$  achieves a comparable or slightly improved performance over the original  $\text{tanh}(x)$ , while still performing worse than  $\text{erf}(x)$ . This indicates that simply scaling  $\text{tanh}(x)$  is insufficient to match the behavior or performance of  $\text{erf}(x)$ .

## 8. Related Work

**Normalization layers.** Since the introduction of Batch Normalization (BN) [25], various normalization methods have been proposed to better stabilize training. To address BN’s limitations with small batches, several alternatives [43, 47, 48, 56, 59] have been explored. In parallel, LayerNorm [2, 39, 57, 58] and RMSNorm [63] were designed for RNN [23] and Transformer architectures [55]. Task-specific variants [47, 54, 56] further adapt normalization to applications such as object detection and style transfer.

**Mechanisms of normalization.** A series of studies has investigated how normalization layers contribute to model convergence. From an optimization perspective, normalization stabilizes gradient flow [6, 13, 35], reduces sensitivity to initialization [14, 46, 65], and implicitly tunes learning rates [1, 51]. It has also been shown to smooth the loss landscape [7, 27, 44] and reduce sharpness [12, 36, 37], promoting more stable optimization dynamics. Understanding these underlying functionalities provides valuable guidance for designing normalization-free training methods.

**Normalization-free methods.** Building on this understanding of normalization, recent work explores how to achieve stable convergence without normalization. One line of work operates at the parameter and optimization level, using tailored initialization schemes [3, 14, 65], self-normalizing activations [28], weight normalization [8, 43], or adaptive gradient clipping [9] to maintain stable gradient propagation. Another line of work modifies the architecture through structural simplifications [21], Softmax-only formulations [26], and bounded convolutional operators [31, 32]. More recently, point-wise functions such as Dynamic Tanh [67] have been proposed, with theoretical analyses revealing their similarity to normalization operations [49]. Unlike previous methods that aim to match the performance of normalization layers, Derf consistently delivers stronger performance across diverse models.

## 9. Conclusion

In this work, we demonstrate that well-designed point-wise functions do not merely match the performance of normalization layers, but can surpass them. By revisiting the design space of point-wise functions, we identify zero-centeredness, boundedness, center sensitivity, and monotonicity as four key properties that enable strong performance in Transformer-based models. Among the functions satisfying these properties, Derf stands out as the most effective design: it outperforms normalization-based methods and another notable point-wise function, DyT, across a wide range of modalities and tasks. Its simplicity and strong empirical performance make Derf a compelling replacement for normalization layers in many Transformer architectures.

## Acknowledgments

We gratefully acknowledge the use of the Neuronic GPU computing cluster maintained by the Department of Computer Science at Princeton University. This work was substantially performed using Princeton Research Computing resources, a consortium led by the Princeton Institute for Computational Science and Engineering (PICSciE) and Research Computing at Princeton University. This work is also supported by the computational resources generously provided by Google’s TPU Research Cloud program.

## References

- [1] Sanjeev Arora, Zhiyuan Li, and Kaifeng Lyu. Theoretical analysis of auto rate-tuning by batch normalization. *ICLR*, 2019. 8
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 1, 2, 8
- [3] Thomas Bachlechner, Bodhisattwa Prasad Majumder, Henry Mao, Gary Cottrell, and Julian McAuley. Rezero is all you need: Fast convergence at large depth. In *UAI*, 2021. 8
- [4] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *NeurIPS*, 2020. 6, 14, 16
- [5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 15
- [6] David Balduzzi, Marcus Frean, Lennox Leary, JP Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question? In *ICML*, 2017. 8
- [7] Nils Björck, Carla P Gomes, Bart Selman, and Kilian Q Weinberger. Understanding batch normalization. In *NeurIPS*, 2018. 1, 8
- [8] Andrew Brock, Soham De, and Samuel L Smith. Characterizing signal propagation to close the performance gap in unnormalized resnets. *ICLR*, 2021. 8
- [9] Andrew Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-performance large-scale image recognition without normalization. In *ICML*, 2021. 8
- [10] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 15
- [11] Zhaodong Chen, Lei Deng, Guoqi Li, Jiawei Sun, Xing Hu, Ling Liang, Yufei Ding, and Yuan Xie. Effective and efficient batch normalization using a few uncorrelated data for statistics estimation. *IEEE Transactions on Neural Networks and Learning Systems*, 2020. 1
- [12] Yan Dai, Kwangjun Ahn, and Suvrit Sra. The crucial role of normalization in sharpness-aware minimization. In *NeurIPS*, 2023. 8
- [13] Hadi Daneshmand, Jonas Kohler, Francis Bach, Thomas Hofmann, and Aurelien Lucchi. Batch normalization provably avoids ranks collapse for randomly initialised deep networks. In *NeurIPS*, 2020. 8
- [14] Soham De and Sam Smith. Batch normalization biases residual blocks towards the identity function in deep networks. In *NeurIPS*, 2020. 8
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 3, 5, 6
- [16] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. 1, 3, 5, 6, 15
- [17] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 15
- [18] Ensembl GRCh38. p13 (genome reference consortium human build 38), insdc assembly, 2013. 7
- [19] Katarína Grešová, Vlastimil Martinek, David Čechák, Petr Šimeček, and Panagiotis Alexiou. Genomic benchmarks: a collection of datasets for genomic sequence classification. *BMC Genomic Data*, 2023. 7
- [20] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 15
- [21] Bobby He and Thomas Hofmann. Simplifying transformer blocks. *ICLR*, 2024. 2, 8
- [22] Stefan Heimersheim. You can remove gpt2’s layernorm by fine-tuning. *arXiv preprint arXiv:2409.13710*, 2024. 2
- [23] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997. 8
- [24] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016. 7, 15
- [25] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 1, 2, 8
- [26] Nandan Kumar Jha and Brandon Reagen. Aero: Softmax-only llms for efficient private inference. *arXiv preprint arXiv:2410.13060*, 2024. 2, 8
- [27] Ryo Karakida, Shotaro Akaho, and Shun-ichi Amari. The normalization method for alleviating pathological sharpness in wide neural networks. In *NeurIPS*, 2019. 8
- [28] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks. In *NeurIPS*, 2017. 8
- [29] Xiangru Lian and Ji Liu. Revisit batch normalization: New understanding and refinement via composition optimization. In *AISTATS*, 2019. 1
- [30] Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*, 2024. 15

- [31] Weiyang Liu, Yan-Ming Zhang, Xingguo Li, Zhiding Yu, Bo Dai, Tuo Zhao, and Le Song. Deep hyperspherical learning. *NeurIPS*, 2017. 8
- [32] Weiyang Liu, Zhen Liu, Zhiding Yu, Bo Dai, Rongmei Lin, Yisen Wang, James M Rehg, and Le Song. Decoupled networks. In *CVPR*, 2018. 8
- [33] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 15
- [34] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, 2022. 15
- [35] Ekdeep S Lubana, Robert Dick, and Hidenori Tanaka. Beyond batchnorm: Towards a unified understanding of normalization in deep learning. In *NeurIPS*, 2021. 8
- [36] Kaifeng Lyu, Zhiyuan Li, and Sanjeev Arora. Understanding the generalization benefit of normalization layers: Sharpness reduction. In *NeurIPS*, 2022. 8
- [37] Maximilian Mueller, Tiffany Vlaar, David Rolnick, and Matthias Hein. Normalization layers are all that sharpness-aware minimization needs. In *NeurIPS*, 2023. 8
- [38] Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Michael Wornow, Callum Birch-Sykes, Stefano Massaroli, Aman Patel, Clayton Rabideau, Yoshua Bengio, et al. Hye-nadna: Long-range genomic sequence modeling at single nucleotide resolution. In *NeurIPS*, 2023. 7, 15, 16
- [39] Toan Q Nguyen and Julian Salazar. Transformers without tears: Improving the normalization of self-attention. *IWSLT*, 2019. 8
- [40] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus based on public domain audio books. In *ICASSP*, 2015. 6
- [41] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, 2023. 5, 6, 13, 16
- [42] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 2020. 15
- [43] Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *NeurIPS*, 2016. 1, 8
- [44] Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization? In *NeurIPS*, 2018. 1, 8
- [45] Yair Schiff, Chia-Hsiang Kao, Aaron Gokaslan, Tri Dao, Albert Gu, and Volodymyr Kuleshov. Caduceus: Bi-directional equivariant long-range dna sequence modeling. In *ICML*, 2024. 7, 15, 16
- [46] Jie Shao, Kai Hu, Changhu Wang, Xiangyang Xue, and Bhiksha Raj. Is normalization indispensable for training deep neural network? In *NeurIPS*, 2020. 8
- [47] Sheng Shen, Zhewei Yao, Amir Gholami, Michael Mahoney, and Kurt Keutzer. Powernorm: Rethinking batch normalization in transformers. In *ICML*, 2020. 8
- [48] Saurabh Singh and Shankar Krishnan. Filter response normalization layer: Eliminating batch dependence in the training of deep neural networks. In *CVPR*, 2020. 1, 8
- [49] Felix Stollenwerk. The mathematical relationship between layer normalization and dynamic activation functions. *arXiv preprint arXiv:2503.21708*, 2025. 8
- [50] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 15
- [51] Hidenori Tanaka and Daniel Kunin. Noether’s learning dynamics: Role of symmetry breaking in neural networks. In *NeurIPS*, 2021. 8
- [52] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 15
- [53] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 15
- [54] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. 1, 8
- [55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 1, 8
- [56] Yuxin Wu and Kaiming He. Group normalization. In *ECCV*, 2018. 1, 8, 16
- [57] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In *ICML*, 2020. 8
- [58] Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and Junyang Lin. Understanding and improving layer normalization. In *NeurIPS*, 2019. 8
- [59] Junjie Yan, Ruosi Wan, Xiangyu Zhang, Wei Zhang, Yichen Wei, and Jian Sun. Towards stabilizing batch statistics in backward propagation of batch normalization. *ICLR*, 2020. 8
- [60] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. 15
- [61] Qiming Yang, Kai Zhang, Chaoxiang Lan, Zhi Yang, Zheyang Li, Wenming Tan, Jun Xiao, and Shiliang Pu. Unified normalization for accelerating and stabilizing transformers. In *ACM MM*, 2022. 1
- [62] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019. 15
- [63] Biao Zhang and Rico Sennrich. Root mean square layer normalization. In *NeurIPS*, 2019. 1, 2, 8, 15, 16
- [64] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. 2018. 15

- [65] Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without normalization. *ICLR*, 2019. [8](#)
- [66] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, 2020. [15](#)
- [67] Jiachen Zhu, Xinlei Chen, Kaiming He, Yann LeCun, and Zhuang Liu. Transformers without normalization. In *CVPR*, 2025. [1](#), [2](#), [8](#)