

DNT: a Deeply Normalized Transformer

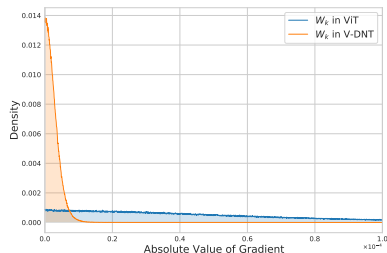
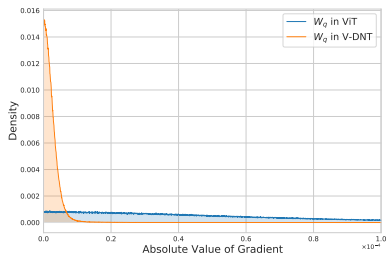
That Can Be Trained by Momentum SGD

Xianbiao Qi¹, Marco Chen², Wenjie Xiao³, Jiaquan Ye¹,
Yelin He¹, Chun-Guang Li^{*4}, Zhouchen Lin^{*5,6}
¹Intellifusion Inc., ²Tsinghua University, ³JHU,
⁴BUPT, ⁵Peking University

ICLR 2026

Motivation: Why Can't SGD Train Transformers?

- ▶ CNNs can be trained with vanilla momentum SGD (mSGD)
- ▶ Transformers **require** adaptive optimizers like AdamW
- ▶ **Root cause:** Gradients in Transformers exhibit **heavy-tailed distributions**



Standard ViT: gradients spread in $[0, 10^{-4}]$ vs V-DNT: concentrated around $[0, 10^{-5}]$

Root Cause: Jacobian Matrix Perspective

For a forward layer $\mathbf{x}^l = \mathbf{W}^l \mathbf{x}^{l-1}$, the gradient is:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^l} = \underbrace{\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}}}_{\text{upstream grad}} \cdot \underbrace{\frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l}}_{\text{Jacobian}} \cdot \mathbf{x}^{l-1 \top}$$

Key insight: Heavy-tail gradients \Leftarrow **diverse singular values** in the Jacobian matrix

Causes of diverse singular values

1. Weight matrices \mathbf{W} have very diverse singular values
2. Activations span widely \Rightarrow uneven Jacobian singular values
3. Large condition number \Rightarrow stretches input differently along directions

Solution: Constrain Jacobian singular values via **proper normalization**

Five Normalization Settings

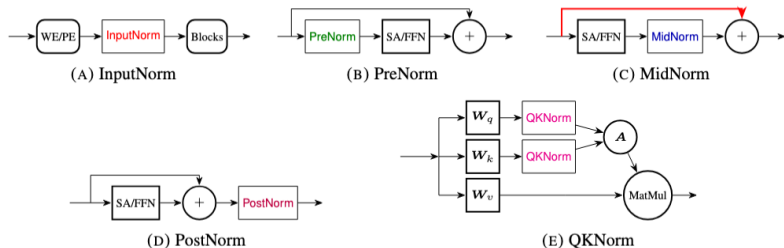


Figure: Five different normalization methods. The only difference between them is the position of normalization. In (A), “WE/PE” indicates Word Embedding (WE) and Patch Embedding (PE).

DNT Architecture

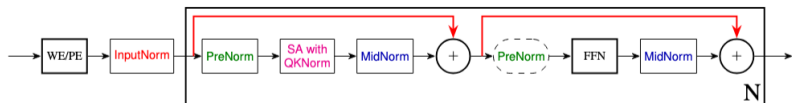


Figure: DNT architecture. The second PreNorm marked with dashed and rounded corners is optional. By default, we do not use the second PreNorm.

Why InputNorm?

For a residual network: $\mathbf{x}^{l+1} = \mathbf{x}^0 + f(\mathbf{x}^0) + f(\mathbf{x}^1) + \dots + f(\mathbf{x}^l)$

In high dimensions, random vectors are almost orthogonal, so:

$$\|\mathbf{x}^{l+1}\|_2 \approx \sqrt{\|\mathbf{x}^0\|_2^2 + \|f(\mathbf{x}^0)\|_2^2 + \dots + \|f(\mathbf{x}^l)\|_2^2}$$

The Jacobian of RMSNorm is:

$$\frac{\partial \text{RMSN}(\mathbf{x}^{l+1})}{\partial \mathbf{x}^{l+1}} = \frac{\sqrt{d}}{\sqrt{\|\mathbf{x}^{l+1}\|_2^2 + \epsilon}} \text{diag}(\gamma) \left(\mathbf{I} - \frac{\mathbf{x}^{l+1} \mathbf{x}^{l+1 \top}}{\|\mathbf{x}^{l+1}\|_2^2 + \epsilon} \right)$$

Why PreNorm?

Self-attention: $\mathbf{Y} = \mathbf{W}_v \mathbf{X} \mathbf{A}$, where $\mathbf{A} = \text{softmax}\left(\frac{\mathbf{x}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{x}}{\sqrt{d_q}}\right)$

The Jacobian matrix of self-attention w.r.t. \mathbf{X} :

$$\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})} = (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{x}) \frac{\mathbf{J}}{\sqrt{d_q}} \left((\mathbf{x}^\top \mathbf{W}_k^\top \mathbf{W}_q \otimes \mathbf{I}_n) \mathbf{C} + (\mathbf{I}_n \otimes \mathbf{x}^\top \mathbf{W}_q^\top \mathbf{W}_k) \right)$$

Why MidNorm?

For FFN: $\mathbf{z} = \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{x})$, followed by MidNorm
 $\mathbf{y} = \text{MidNorm}(\mathbf{z})$

The Jacobian of FFN + MidNorm:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \sqrt{d} \text{diag}(\gamma) \left(\mathbf{I} - \frac{\mathbf{z}\mathbf{z}^\top}{\|\mathbf{z}\|_2^2} \right) \cdot \underbrace{\frac{\mathbf{W}_2 \text{diag}(\mathbf{1}(\mathbf{W}_1 \mathbf{x} > \mathbf{0})) \mathbf{W}_1}{\|\mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{x})\|_2}}_M$$

Why **NOT** PostNorm?

PostNorm normalizes *after* the residual connection:

$$\mathbf{x}^{l+1} = \text{PostNorm}(\mathbf{z}^{l+1}), \quad \text{where} \quad \mathbf{z}^{l+1} = \mathbf{x}^l + f(\mathbf{x}^l; \mathbf{W}^{l+1})$$

The Jacobian:

$$\frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{z}^{l+1}} = \frac{\sqrt{d}}{\|\mathbf{z}^{l+1}\|_2} \text{diag}(\gamma) \left(\mathbf{I} - \frac{\mathbf{z}^{l+1} \mathbf{z}^{l+1 \top}}{\|\mathbf{z}^{l+1}\|_2^2} \right)$$

PostNorm is **unstable** — we do **not** use it in DNT.

Why QKNorm?

Normalize queries and keys:

$$\mathbf{q}'_i = \sqrt{d_h} \operatorname{diag}(\gamma_q) \frac{\mathbf{W}_q \mathbf{x}_i}{\|\mathbf{W}_q \mathbf{x}_i\|_2}, \quad \mathbf{k}'_j = \sqrt{d_h} \operatorname{diag}(\gamma_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\|\mathbf{W}_k \mathbf{x}_j\|_2}$$

The gradient of the logit $P'_{ij} = \mathbf{q}'_i{}^\top \mathbf{k}'_j$:

$$\frac{\partial P'_{ij}}{\partial \mathbf{x}} = \sqrt{d_h} \operatorname{diag}(\gamma_q) \mathbf{k}'_j{}^\top \left(\mathbf{I} - \frac{\mathbf{q}'_i \mathbf{q}'_i{}^\top}{\|\mathbf{q}'_i\|_2^2} \right) \frac{\mathbf{W}_q}{\|\mathbf{W}_q \mathbf{x}_i\|_2} + \sqrt{d_h} \operatorname{diag}(\gamma_k) \mathbf{q}'_i{}^\top \left(\mathbf{I} - \frac{\mathbf{k}'_j \mathbf{k}'_j{}^\top}{\|\mathbf{k}'_j\|_2^2} \right) \frac{\mathbf{W}_k}{\|\mathbf{W}_k \mathbf{x}_j\|_2}$$

Experimental Results

Optimizer	Model	ImageNet (Acc. \uparrow)		OpenWebText (Val Loss \downarrow)		
		307M	632M	124M	774M	1436M
AdamW	ViT/GPT2	81.7	80.8	2.867	2.492	2.435
AdamW	V-DNT/L-DNT	82.1	81.9	2.863	2.481	2.396
mSGDW	ViT/GPT2	78.2	73.5	2.906	2.544	2.472
mSGDW	V-DNT/L-DNT	81.5	81.2	2.849	2.503	2.408

Conclusion

1. **Theoretical analysis:** heavy-tail gradients stem from diverse Jacobian singular values
2. Each normalization has a **clear theoretical justification**
3. **DNT architecture:** four strategically placed normalizations to constrain Jacobians