Mitigating Over-smoothing in Transformers via Regularized Nonlocal Functionals

Tam Nguven

Department of Electrical & Computer Engineering Rice University Houston, USA mn72@rice.edu

Tan M. Nguyen

Department of Mathematics National University of Singapore Singapore tanmn@nus.edu.sg

Richard G. Baraniuk

Department of Electrical & Computer Engineering Rice University Houston, USA richb@rice.edu

Abstract

Transformers have achieved remarkable success in a wide range of natural language processing and computer vision applications. However, the representation capacity of a deep transformer model is degraded due to the over-smoothing issue in which the token representations become identical when the model's depth grows. In this work, we show that self-attention layers in transformers minimize a functional which promotes smoothness, thereby causing token uniformity. We then propose a novel regularizer that penalizes the norm of the difference between the smooth output tokens from self-attention and the input tokens to preserve the fidelity of the tokens. Minimizing the resulting regularized energy functional, we derive the Neural Transformer with a Regularized Nonlocal Functional (NeuTRENO), a novel class of transformer models that can mitigate the over-smoothing issue. We empirically demonstrate the advantages of NeuTRENO over the baseline transformers and state-of-the-art methods in reducing the over-smoothing of token representations on various practical tasks, including object classification, image segmentation, and language modeling.

1 Introduction

Transformer models [62] have achieved substantial success in natural language processing [16, 2, 13, 10, 47, 4, 6, 14], reinforcement learning [9, 32], computer vision [19, 40, 59, 49, 44, 3, 41, 71, 27], and other practical applications [50, 33, 70, 26, 66]. Transformers also excel at transferring knowledge from pre-trained models to new tasks, even when limited supervision is available [45, 46, 16, 69, 39]. At the heart of transformers lies the self-attention mechanism, which computes a weighted average of token representations within a sequence. These weights are determined based on the similarity scores between pairs of tokens, determining their relative importance in the sequence [11, 43, 38]. This flexibility in capturing diverse syntactic and semantic relationships has been identified as a crucial factor contributing to the success of transformers [57, 63, 12, 64, 31].

1.1 Background: Self-Attention

For a given input sequence $\mathbf{X} := [\boldsymbol{x}(1), \cdots, \boldsymbol{x}(N)]^{\top} \in \mathbb{R}^{N \times D_x}$ of N feature vectors, self-attention transforms \mathbf{X} into the output sequence \mathbf{H} in the following two steps:

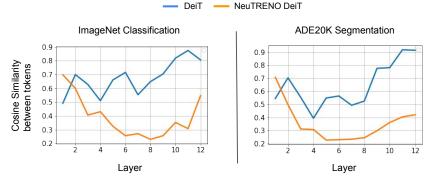


Figure 1: The cosine similarity between tokens representations across layers of NeuTRENO DeiT vs. the baseline DeiT models on the Imagenet classification and ADE20K image segmentation tasks. In both tasks, the DeiT baseline suffers from over-smoothing as tokens become similar to identical when the model gets deeper. In contrast, tokens in NeuTRENO models are significantly more diverse, suggesting a reduction in over-smoothing. Further details regarding this analysis can be found in Appendix E.

Step 1. The input sequence X is projected into the query matrix Q, the key matrix K, and the value matrix V via three linear transformations

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_{Q}^{\top}; \mathbf{K} = \mathbf{X}\mathbf{W}_{K}^{\top}; \mathbf{V} = \mathbf{X}\mathbf{W}_{V}^{\top}, \tag{1}$$

where $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{D_{qk} \times D_x}$, and $\mathbf{W}_V \in \mathbb{R}^{D \times D_x}$ are the weight matrices. We denote $\mathbf{Q} := [\mathbf{q}(1), \dots, \mathbf{q}(N)]^{\top}, \mathbf{K} := [\mathbf{k}(1), \dots, \mathbf{k}(N)]^{\top}$, and $\mathbf{V} := [\mathbf{v}(1), \dots, \mathbf{v}(N)]^{\top}$, where the vectors $\mathbf{q}(i), \mathbf{k}(i)$, and $\mathbf{v}(i)$, for $i = 1, \dots, N$ are the query, key, and value vectors, respectively.

Step 2. The output sequence $\mathbf{U} := [\boldsymbol{u}(1), \dots, \boldsymbol{u}(N)]^{\top} \in \mathbb{R}^{N \times D_{qk}}$ is then computed as follows

$$\mathbf{U} = \operatorname{softmax} \left(\mathbf{Q} \mathbf{K}^{\top} / \sqrt{D_{qk}} \right) \mathbf{V} := \mathbf{A} \mathbf{V}, \tag{2}$$

where the softmax function is applied to each row of the matrix $\mathbf{Q}\mathbf{K}^{\top}/\sqrt{D_{qk}}$. The matrix $\mathbf{A} := \operatorname{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{D_{qk}}}\right) \in \mathbb{R}^{N \times N}$ and its component a_{ij} for $i, j = 1, \cdots, N$ are called the attention matrix and attention scores, respectively. For each query vector $\mathbf{q}(i)$ for $i = 1, \cdots, N$, an equivalent form of Eqn. (2) to compute the output vector $\mathbf{u}(i)$ is given by

$$\boldsymbol{u}(i) = \sum_{j=1}^{N} \operatorname{softmax} \left(\boldsymbol{q}(i)^{\top} \boldsymbol{k}(j) / \sqrt{D_{qk}} \right) \boldsymbol{v}(j).$$
 (3)

The self-attention computed by Eqn. (2) and (3) is referred as softmax attention. In our work, we refer to a transformer that uses softmax attention as a softmax transformer.

1.2 Over-smoothing in Transformers

Despite their remarkable success, deep transformer-based models have been observed to suffer from the over-smoothing issue, in which all token representations become identical when more layers are added to the models [55, 65, 18]. This over-smoothing phenomenon, also known as the "token uniformity" problem, significantly limits the representation capacity of transformers. To illustrate this phenomenon, we examine the average cosine similarity between pairs of token representations across different layers in a softmax transformer trained for the Imagenet object classification and ADK20 image segmentation tasks [73]. As depicted in Fig. 1, in both tasks, this cosine similarity between tokens increases as the models become deeper. Particularly, in the last two layers, the cosine similarity scores are approximately 0.9, indicating a high degree of similarity among tokens.

1.3 Contribution

We develop a nonlocal variational denoising framework for self-attention, providing insights into the over-smoothing phenomenon in transformers. In particular, by viewing self-attention as a gradient descent step toward minimizing a nonlocal functional that penalizes high-frequency noise in the signal, we uncover the diffusive nature of self-attention, which explains the over-smoothing issue of transformers. Motivated by this understanding, we propose the Neural Transformer with a Regularized

Nonlocal Functional (NeuTRENO), a novel class of transformers designed to mitigate over-smoothing. NeuTRENO is derived by optimizing a regularized nonlocal functional, which includes an additional convex fidelity term. This fidelity term penalizes the norm of the difference between the smooth output tokens from self-attention and the input tokens, thereby reducing the over-smoothing effect. Our contribution is three-fold.

- 1. We develop a nonlocal variational denoising framework for self-attention and shed light on the over-smoothing issue that hampers the representation capacity of transformers.
- 2. We develop NeuTRENO, a novel class of transformers that are capable of alleviating the over-smoothing issue.
- 3. We theoretically prove that transformers with softmax self-attention are prone to oversmoothing while NeuTRENO can avoid this issue.

We empirically demonstrate the benefits of NeuTRENO on various large-scale applications, including the ImageNet object classification, ADE20K image segmentation, and WikiText-103 language modeling tasks.

Organization: We organize our paper as follows: in Section 2, we develop a nonlocal variational denoising framework for self-attention and provide an explanation for the over-smoothing issue in transformer-based models. In section 3, we propose NeuTRENO, and present a theoretical result that guarantees NeuTRENO's capability of mitigating over-smoothing. In Section 4, we empirically validate the benefits of NeuTRENO. We discuss the related work in Section 6. Finally, we conclude our main contributions and remarks. Further results, details, and proofs are provided in the Appendix.

2 A Nonlocal Variational Denoising Framework for Self-attention

We first consider the output matrix $\mathbf{U} := [\boldsymbol{u}(1), \cdots, \boldsymbol{u}(N)]^{\top} \in \mathbb{R}^{N \times D}$ in self-attention as given by Eqn. 2 in Section 1.1. Let $\Omega \subset \mathbb{R}, x \in \Omega$, and $\boldsymbol{u}(x) := [u_1(x), \dots, u_D(x)]^T$ be a real vector-valued function, $\boldsymbol{u} : \Omega \to \mathbb{R}^D$, $\boldsymbol{u} \in L^2(\Omega)$. The output matrix \mathbf{U} in self-attention discretizes the function $\boldsymbol{u}(x)$ on a 1-D grid. In the context of signal/image denoising, \mathbf{U} can be considered as the *desired clean signal*, and $\boldsymbol{u}(x)$ is its corresponding intensity function denoting the signal values at the position $x \in \Omega$. We further let the observed intensity function $\boldsymbol{f}(x)$ denote the values of the *observed noisy signal* at $x \in \Omega$, $\boldsymbol{f} : \Omega \to \mathbb{R}^D$, $\boldsymbol{f} \in L^2(\Omega)$. For example, $\boldsymbol{f}(x)$ can be given as

$$f(x) = u(x) + n(x), \tag{4}$$

where n is the additive noise. We wish to reconstruct u(x) from f(x). Following the variational denoising method proposed in [23] and [24], the denoised image u(x) can be obtained by minimizing the following regularized functional with respect to u:

$$E(\boldsymbol{u}, \boldsymbol{f}) = J(\boldsymbol{u}) + G(\boldsymbol{u}, \boldsymbol{f})$$

$$= \frac{1}{2} \int_{\Omega \times \Omega} \|\boldsymbol{u}(x) - \boldsymbol{u}(y)\|_{2}^{2} k(x, y) dx dy + \frac{\lambda}{2} \int_{\Omega} \|\boldsymbol{u}(x) - \boldsymbol{f}(x)\|_{2}^{2} dx.$$
(5)

Here, $J(\boldsymbol{u}) = \frac{1}{2} \int_{\Omega \times \Omega} \|\boldsymbol{u}(x) - \boldsymbol{u}(y)\|_2^2 k(x,y) dx dy$ is a nonlocal functional of weighted differences. The weights k(x,y) represent the affinity between signal values at positions x and y. For example, for images, k(x,y) captures the proximity between pixels x and y in the image. $J(\boldsymbol{u})$ works as a regularizer. Minimizing $J(\boldsymbol{u})$ promotes the smoothness of \boldsymbol{u} and penalizes high-frequency noise in the signal. Adding the convex fidelity term $G(\boldsymbol{u},\boldsymbol{f}) = \frac{\lambda}{2} \int_{\Omega} \|\boldsymbol{u}(x) - \boldsymbol{f}(x)\|_2^2 dx$ to the functional $J(\boldsymbol{u})$ allows the denoised signal $\boldsymbol{u}(x)$ to preserve relevant information in the observed noisy signal $\boldsymbol{f}(x)$. The regularized functional $E(\boldsymbol{u},\boldsymbol{f})$ can be considered as an energy functional.

2.1 Self-attention as a Gradient Descent Step to Minimize the Nonlocal Functional J

We show that self-attention is equivalent to taking a gradient descent step toward minimizing the functional J(u) in the energy functional E(u, f). We expand J(u) as follows

$$J(u) = \frac{1}{2} \int_{\Omega \times \Omega} \sum_{i=1}^{D} (u_j(x) - u_j(y))^2 k(x, y) dx dy$$
 (6)

The gradient of J with respect to u is then given by

$$\nabla_{\boldsymbol{u}}J(\boldsymbol{u}) = \left[\frac{\partial J}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_D}\right]^T. \tag{7}$$

The partial derivative $\partial J/\partial u_j$, $j=1,2,\ldots,D$, is defined through its dot product with an arbitrary function $h_j \in L^2(\Omega)$ as follows

$$\begin{split} \frac{\partial J}{\partial u_j} \cdot h_j(x) &= \frac{d}{d\tau} J(u_j + \tau h_j) \big|_{\tau=0} \\ &= \frac{1}{2} \left(\frac{d}{d\tau} \int_{\Omega \times \Omega} (u_j(x) - u_j(y) + \tau h_j(x) - \tau h_j(y))^2 k(x, y) dx dy \right) \Big|_{\tau=0} \\ &= \left(\int_{\Omega \times \Omega} (u_j(x) - u_j(y) + \tau h_j(x) - \tau h_j(y)) (h_j(x) - h_j(y)) k(x, y) dx dy \right) \Big|_{\tau=0} \\ &= \int_{\Omega \times \Omega} (u_j(x) - u_j(y)) (h_j(x) - h_j(y)) k(x, y) dx dy \\ &= \int_{\Omega \times \Omega} (u_j(x) - u_j(y)) h_j(x) k(x, y) dx dy - \int_{\Omega \times \Omega} (u_j(x) - u_j(y)) h_j(y) k(x, y) dx dy \end{split}$$

Applying a change of variables $(x,y) \to (y,x)$ to the second term of the above integral, we have

$$\frac{\partial J}{\partial u_j} \cdot h_j(x) = \int_{\Omega \times \Omega} (u_j(x) - u_j(y)) h_j(x) k(x, y) dx dy - \int_{\Omega \times \Omega} (u_j(y) - u_j(x)) h_j(x) k(y, x) dx dy$$
$$= \int_{\Omega \times \Omega} (u_j(x) - u_j(y)) (k(x, y) + k(y, x)) dy h_j(x) dx$$

Thus, the Frechet derivative of J with respect to u_i is given by

$$\frac{\partial J}{\partial u_j} = \int_{\Omega} (u_j(x) - u_j(y)(k(x,y) + k(y,x))dy.$$
 (8)

Substituting the formula for $\partial J/\partial u_j$ in Eqn. 8 into Eqn. 7 for $\nabla_{\boldsymbol{u}} J(\boldsymbol{u})(x)$, we obtain the following gradient flow

$$\frac{d\mathbf{u}(x,t)}{dt} = -\nabla_{\mathbf{u}}J(\mathbf{u}) = \int_{\Omega} (\mathbf{u}(y,t) - \mathbf{u}(x,t)) (k(x,y) + k(y,x)) dy, \tag{9}$$

where t is the time variable we introduce to capture the dynamics of \boldsymbol{u} when gradient descent is applied to minimize $J(\boldsymbol{u})$. Let $\boldsymbol{v}(x) := [v_1(x), \dots, v_D(x)]^T$ be a real vector-valued function, $\boldsymbol{v}: \Omega \to \mathbb{R}^D$, $\boldsymbol{v} \in L^2(\Omega)$. We discretize $\boldsymbol{v}(x)$ on a 1-D grid to attain the value vectors $\boldsymbol{v}(1), \dots, \boldsymbol{v}(N) \in \mathbb{R}^D$, which form the value matrix $\mathbf{V} := [\boldsymbol{v}(1), \dots, \boldsymbol{v}(N)]^\top \in \mathbb{R}^{N \times D}$ in self-attention as defined in Eqn. 2. We initialize \boldsymbol{u} at t=0 with $\boldsymbol{v}(x)$, i.e., $\boldsymbol{u}(x,0) = \boldsymbol{v}(x)$.

Self-attention is an Euler Discretization of the Gradient Flow Given in 9. We discretize the gradient flow in Eqn. 9 using the Euler method [21] with step size $\Delta t(x) = 1/\int_{\Omega} \left(k(x,y) + k(y,x)\right) dy$ and obtain the following update

$$\mathbf{u}(x,\Delta t(x)) = \mathbf{u}(x,0) + \Delta t(x) \int_{\Omega} (\mathbf{u}(y,0) - \mathbf{u}(x,0)) (k(x,y) + k(y,x)) dy$$

$$= \int_{\Omega} \frac{(k(x,y) + k(y,x)) \mathbf{u}(y,0)}{\int_{\Omega} (k(x,y') + k(y',x)) dy'} dy = \int_{\Omega} \frac{K(x,y) \mathbf{v}(y)}{\int_{\Omega} K(x,y') dy'} dy.$$
(10)

Here, K(x,y):=k(x,y)+k(y,x) is a symmetric kernel and $\boldsymbol{u}(y,0)=\boldsymbol{v}(y)$ since \boldsymbol{u} is initialized at t=0 with \boldsymbol{v} as aforementioned. Let $\boldsymbol{k}(x):=[k_1(x),\ldots,k_{D_{qk}}(x)]^T$ be a real vector-valued function, $\boldsymbol{k}:\Omega\to\mathbb{R}^{D_{qk}},\,\boldsymbol{k}\in L^2(\Omega)$. Similar to $\boldsymbol{u}(x)$ and $\boldsymbol{v}(x)$, we can discretize $\boldsymbol{k}(x)$ on a 1-D grid to attain the key vectors $\boldsymbol{k}(1),\ldots,\boldsymbol{k}(N)\in\mathbb{R}^{D_{qk}}$, which form the key matrix $\mathbf{K}:=[\boldsymbol{k}(1),\cdots,\boldsymbol{k}(N)]^\top\in\mathbb{R}^{N\times D_{qk}}$ in self-attention as defined in Eqn. 2. We choose $K(x,y)=\exp(\boldsymbol{k}(x)^T\boldsymbol{k}(y)/\sqrt{D_{qk}})$ and rewrite Eqn. 10 as follows

$$\boldsymbol{u}(x,\Delta t(x)) = \int_{\Omega} \frac{\exp(\boldsymbol{k}(x)^T \boldsymbol{k}(y) / \sqrt{D_{qk}})}{\int_{\Omega} \exp(\boldsymbol{k}(x)^T \boldsymbol{k}(y') / \sqrt{D_{qk}}) dy'} \boldsymbol{v}(y) dy.$$
(11)

Estimating the integrals in Eqn. 11 via Monte-Carlo approximation using the key vectors $k(1), \ldots, k(N) \in \mathbb{R}^{D_{qk}}$ and and value vectors $v(1), \ldots, v(N) \in \mathbb{R}^{D}$, we obtain

$$\boldsymbol{u}(x, \Delta t(x)) \approx \sum_{j=1}^{N} \frac{\exp(\boldsymbol{k}(x)^{T} \boldsymbol{k}(j) / \sqrt{D_{qk}})}{\sum_{j'=1}^{N} \exp(\boldsymbol{k}(x)^{T} \boldsymbol{k}(j') / \sqrt{D_{qk}})} \boldsymbol{v}(j).$$
(12)

Discretizing $u(x, \Delta t(x))$ on another 1-D grid, we attain

$$\boldsymbol{u}(i) \approx \sum_{j=1}^{N} \frac{\exp(\boldsymbol{k}(i)^{T} \boldsymbol{k}(j) / \sqrt{D_{qk}})}{\sum_{j'=1}^{N} \exp(\boldsymbol{k}(i)^{T} \boldsymbol{k}(j') / \sqrt{D_{qk}})} \boldsymbol{v}(j)$$

$$= \sum_{j=1}^{N} \operatorname{softmax} \left(\boldsymbol{k}(i)^{T} \boldsymbol{k}(j) / \sqrt{D_{qk}}\right) \boldsymbol{v}(j), \quad i = 1, \dots, N.$$
(13)

Comparing Eqn. 13 and Eqn. 3, we observe that Eqn. 13 implement a symmetric self-attention, in which the query matrix \mathbf{Q} and the key matrix \mathbf{K} are the same, i.e. $\mathbf{W}_Q = \mathbf{W}_K$ where \mathbf{W}_Q and \mathbf{W}_K are the linear projections that map the input sequence \mathbf{X} into \mathbf{Q} and \mathbf{K} as given in Eqn. 1. This symmetry of the attention scores is desirable in some image processing tasks due to the symmetric similarities between pixels, but can be relaxed for other tasks. To break the symmetry of attention scores in Eqn. 13, we replace the key vectors $\mathbf{k}(i)$ by the query vectors $\mathbf{q}(i)$, $i=1,\ldots,N$, to obtain the exact formula of self-attention given by Eqn. 3. The following theorem summarizes our results:

Theorem 1 (Self-attention as a Gradient Descent Step to Minimize a Nonlocal Functional). Given the nonlocal functional $J(\boldsymbol{u}) = \frac{1}{2} \int_{\Omega \times \Omega} \|\boldsymbol{u}(x) - \boldsymbol{u}(y)\|_2^2 k(x,y) dxdy$ of a vector-valued function $\boldsymbol{u}: \Omega \to \mathbb{R}^D$, $\boldsymbol{u} \in L^2(\Omega)$, and let $K(x,y) := k(x,y) + k(y,x) = \exp(\boldsymbol{k}(x)^T \boldsymbol{k}(y) / \sqrt{D_{qk}})$, where $\boldsymbol{k}: \Omega \to \mathbb{R}^{D_{qk}}$, $\boldsymbol{k} \in L^2(\Omega)$. Then, taking a gradient descent step on \boldsymbol{u} at time t=0, where $\boldsymbol{u}(x,0) = \boldsymbol{v}(x)$, with an adaptive step size $\Delta t(x) := \frac{1}{\int_{\Omega} (k(x,y) + k(y,x)) dy}$ to minimize J is equivalent to updating \boldsymbol{u} via a symmetric self-attention

$$\boldsymbol{u}(x, \Delta t(x)) = \sum_{j=1}^{N} \operatorname{softmax} \left(\boldsymbol{k}(x)^{\top} \boldsymbol{k}(j) / \sqrt{D_{qk}} \right) \boldsymbol{v}(j),$$

which results in

$$\boldsymbol{u}(i) = \sum_{j=1}^{N} \operatorname{softmax} \left(\boldsymbol{k}(i)^{\top} \boldsymbol{k}(j) / \sqrt{D_{qk}} \right) \boldsymbol{v}(j), \quad i = 1, \dots, N.$$
(14)

Here, $\mathbf{u}(n)$, $\mathbf{v}(n)$, and $\mathbf{u}(n)$, $n=1,\ldots,N$, are the key, value, and output vectors in self-attention, respectively. Breaking the symmetry of the attention scores by replacing $\mathbf{k}(i)$ with $\mathbf{q}(i)$, $i=1,\ldots,N$, in Eqn. 14, we obtain the exact formula of self-attention

$$\boldsymbol{u}(i) = \sum_{j=1}^{N} \operatorname{softmax} \left(\boldsymbol{q}(i)^{\top} \boldsymbol{k}(j) / \sqrt{D_{qk}} \right) \boldsymbol{v}(j), \ i = 1, \dots, N.$$

Remark 1. In Eqn. 9, the change in \mathbf{u} at position x is proportional to the sum of differences between $\mathbf{u}(x)$ and \mathbf{u} at other position in the domain Ω . In particular, when $\mathbf{u}(x)$ is smaller or larger than the values at other positions, it will increase or decrease, respectively. This is analogous to a diffusion process in which particles or substances move from high-concentration to low-concentration regions. It has been proved that a diffusion process converges to a saturating state in which the concentrations at all positions are the same. This suggests that $\mathbf{u}(x)$ tends to suffer from the over-smoothing issue.

2.2 Random Walk Analysis of Over-smoothing

The diffusion process and random walk are closely related concepts, as diffusion can be seen as a collective behavior of numerous random walks performed by individual particles or molecules. Inspired by the analogy between the dynamics of \boldsymbol{u} in Eqn 9 and a diffusion process, as well as the relationship between diffusion process and random walk, in this section, we show the connection

between the evolution of u and a random walk. By adopting a random walk perspective on graph neural network [58], we demonstrate that u(x) under the dynamics given in Eqn 9 suffers from over-smoothing.

Recall from the gradient flow in Eqn 9, by using Euler method discretization, after k update steps starting from the initial $\boldsymbol{u}(x,0) = \boldsymbol{v}(x)$, with adaptive stepsize $\Delta t = 1/\int_{\Omega} \left(k(x,y) + k(y,x)\right) dy$, we obtain the following

$$\mathbf{u}(x, k\Delta t(x)) = \int_{\Omega} \frac{K(x, y)\mathbf{u}(y, (k-1)\Delta t(x))}{\int_{\Omega} K(x, y')dy'} dy.$$
 (15)

Discretizing $u(x, k\Delta t(x))$ and using Monte-Carlo approximation for the integrals in 15, we attain

$$\mathbf{u}^{(k)}(i) = \sum_{j=1}^{N} \mathbf{A}_{ij} \mathbf{u}^{(k-1)}(j)$$
 (16)

where \mathbf{A}_{ij} is computed using the keys and queries as either $\operatorname{softmax}\left(\mathbf{k}(i)^{\top}\mathbf{k}(j)/\sqrt{D_{qk}}\right)$ or $\operatorname{softmax}\left(\mathbf{q}(i)^{\top}\mathbf{k}(j)/\sqrt{D_{qk}}\right)$. Let $\{\mathbf{B}^{(k)}(i)\}_{k\in K}$ be a random walk on $\{\mathbf{v}(i)\}_{i=1}^{N}$ as defined:

$$\mathbf{B}^{(0)}(i) = \boldsymbol{v}(i)$$

$$\mathbb{P}(\mathbf{B}^{(k+1)}(l) = \boldsymbol{v}(j)|\mathbf{B}^{(k)}(l) = \boldsymbol{v}(i)) = \mathbf{A}_{ij}$$
(17)

where $\mathbf{B}^{(k)}(n)$ is the random value of a k-step walk, starts at node n, and $\boldsymbol{v}(n)$ is the initial value at node n, respectively, for $n=1,2,\ldots,N$. The transition probability \mathbf{A} is defined as above. To investigate the connection between the update process of \boldsymbol{u} and the random walk defined in 17, we show that, for $i=1,2,\ldots,N$, after k update steps as in 16, with initial value $\boldsymbol{u}^{(0)}(i)=\boldsymbol{v}(i),\boldsymbol{u}(i)^{(k)}$ equals to the expected value of the k-step walk, starting at node i:

Lemma 1. Let $u^{(k)}(i)$ defined in 16 and $\{\mathbf{B}^{(k)}(i)\}_{k\in K}$ is the random walk defined by 17. Then

$$\boldsymbol{u}^{(k)}(i) = \mathbb{E}[\mathbf{B}^{(k)}(i)]. \tag{18}$$

We next present the Lemma 2 which is necessary to show the convergence of $u^{(k)}(i)$.

Lemma 2. The random walk $\mathbf{B}^{(k)}(i)$ in 17 with the transition matrix \mathbf{A} either be $\mathbf{A}_{ij} = \operatorname{softmax}\left(\mathbf{k}(i)^{\top}\mathbf{k}(j)/\sqrt{D_{qk}}\right)$ or $\mathbf{A}_{ij} = \operatorname{softmax}\left(\mathbf{q}(i)^{\top}\mathbf{k}(j)/\sqrt{D_{qk}}\right)$, has a unique stationary distribution $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_N]$ such that $\pi_i := P(\mathbf{B}^{(k)}(j) = \boldsymbol{v}(i))$, for $i, j = 1, 2, \dots, N, \sum_{i=1}^N \pi_i = 1$, and $\boldsymbol{\pi}^T = \boldsymbol{\pi}^T \mathbf{A}$.

If $\mathbf{A}_{ij} = \operatorname{softmax} \left(\mathbf{k}(i)^{\top} \mathbf{k}(j) / \sqrt{D_{qk}} \right)$, the stationary distribution is:

$$\boldsymbol{\pi} = \left(\frac{d_1}{\sum_{j=1}^{N} d_j}, \frac{d_2}{\sum_{j=1}^{N} d_j}, \dots, \frac{d_n}{\sum_{j=1}^{N} d_j}\right),\tag{19}$$

where $d_i = \sum_{j=1}^N \exp(\mathbf{k}(i)^\top \mathbf{k}(j)/\sqrt{D_{qk}})$, $\mathbf{k}(1), \mathbf{k}(2), \dots, \mathbf{k}(N)$ are the key vectos.

In general, π_i can be found by finding the left eigenvector of **A** corresponding to the dominant eigenvalue 1.

From the Lemma 1 and Lemma 2, we see that, for all i = 1, 2, ..., N,

$$\mathbf{u}^{(k)}(i) = \mathbb{E}[\mathbf{B}^{(k)}(i)] = \sum_{j=1}^{N} \mathbf{v}(j) \mathbb{P}(\mathbf{B}^{(k-1)}(i) = \mathbf{v}(j)) \to \sum_{j=1}^{N} \pi_{j} \mathbf{v}(j) =: \bar{\mathbf{v}}.$$
 (20)

as $k \to \infty$. This shows that when k increases, $u(i)^{(k)}$ converges to a constant vector, indicating that u(x), under the dynamic in 9, suffers from over-smoothing.

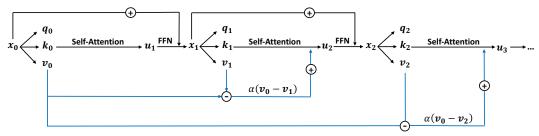


Figure 2: Our proposed NeuTRENO model adds a proportion of the difference between the values of the first and that of the current layer to the self-attention's output at each layer.

3 NeuTRENO: Mitigating the Over-smoothing in Transformers via Minimizing a Regularized Functional

In Section 2.1, we have shown that self-attention implicitly performs a gradient descent step to minimize the nonlocal functional J(u) in Eqn. 5, which results in the diffusive characteristics of u and causes the over-smoothing phenomenon in transformers, as proved in Section 2.2. Fortunately, our objective is not to minimize J(u) but the energy/regularized functional E(u, f) defined by Eqn. 5. This regularized functional consists of not only J(u) but also the convex fidelity term $G(u, f) = \frac{\lambda}{2} \int_{\Omega} \|u(x) - f(x)\|_2^2 dx$. This fidelity term aims to preserve the relevant information in the observed noisy signal f(x) by penalizing solution u(x) that deviates significantly from f(x), thereby mitigating the effects of over-smoothing caused by minimizing J(u).

In this section, we will derive our Neural Transformer with a Regularized Nonlocal Functional (NeuTRENO) by minimizing the regularized functional $E(\boldsymbol{u}, \boldsymbol{f})$. We then provide a theoretical result to prove that NeuTRENO does not suffer from over-smoothing. Recall from Eqn. 5 that $E(\boldsymbol{u}, \boldsymbol{f})$ is given by

$$E(\boldsymbol{u}, \boldsymbol{f}) = J(\boldsymbol{u}) + G(\boldsymbol{u}, \boldsymbol{f}) = J(\boldsymbol{u}) + \frac{\lambda}{2} \int_{\Omega} \sum_{j=1}^{D} (u_j(x) - f_j(x))^2 dx$$

Following a similar derivation as in Section 2.1 (see Appendix C for the detailed derivation), we obtain the following gradient flow when minimizing E(u, f) using gradient descent

$$\frac{d\mathbf{u}(x,t)}{dt} = -\nabla_{\mathbf{u}}E(\mathbf{u},\mathbf{f}) = -\nabla_{\mathbf{u}}J(\mathbf{u}) - \lambda(\mathbf{u}(x) - \mathbf{f}(x)), \tag{21}$$

NeuTRENO-attention is an Euler Discretization of the Gradient Flow Given in 21. Following the similar derivation in Section 2.1, we discretize the gradient flow in Eqn. 21 using the Euler method [21] with step size $\Delta t(x) = 1/\int_{\Omega} (k(x,y) + k(y,x)) dy$ and initializing \boldsymbol{u} at t=0 with $\boldsymbol{v}(x)$, i.e., $\boldsymbol{u}(x,0) = \boldsymbol{v}(x)$. Choosing $\lambda = \tilde{\lambda}/\Delta t(x)$, we obtain the following update

$$\mathbf{u}(x, \Delta t(x)) = \mathbf{u}(x, 0) - \Delta t(x) \nabla_{\mathbf{u}} J - \lambda \Delta t(x) (\mathbf{u}(x, 0) - \mathbf{f}(x))$$

$$= \int_{\Omega} \frac{K(x, y) \mathbf{v}(y)}{\int_{\Omega} K(x, y') dy'} dy + \tilde{\lambda} (\mathbf{f}(x) - \mathbf{v}(x)). \tag{22}$$

We choose the observed noisy signal $f(x) = v^0(x)$ where $v^0(x)$ is v(x) at the first layer in the transformer model. The update in Eqn. 22 becomes

$$\boldsymbol{u}(x,\Delta t(x)) = \int_{\Omega} \frac{K(x,y)\boldsymbol{v}(y)}{\int_{\Omega} K(x,y')dy'} dy + \tilde{\lambda} (\boldsymbol{v}^{0}(x) - \boldsymbol{v}(x)).$$
 (23)

Applying the Monte-Carlo method to approximate the integrals in Eqn. 23 and discretizing $u(x, \Delta t(x))$, v(x), and $v^0(x)$ on a 1-D grid, we attain the following new formula for calculating symmetric self-attention:

$$\boldsymbol{u}(i) = \sum_{j=1}^{N} \operatorname{softmax} \left(\boldsymbol{k}(i)^{\top} \boldsymbol{k}(j) / \sqrt{D_{qk}} \right) \boldsymbol{v}(j) + \tilde{\lambda} (\boldsymbol{v}^{0}(i) - \boldsymbol{v}(i)), \quad i = 1, \dots, N.$$
 (24)

Table 1: Top-1 and Top-5 accuracy (%) of NeuTRENO DeiT vs. DeiT on the ImageNet benchmark. We also present the performance of adapting NeuTRENO to the pre-trained DeiT baseline, NeuTRENO Adaptation. In addition, we compare NeuTRENO with FeatScale [65] and incorporate our method with FeatScale model.

Model/Metric	Top-1 Acc (%)	Top-5 Acc (%)
Softmax DeiT	72.17	91.02
NeuTRENO-DeiT	73.01	91.56
NeuTRENO Adaptation	72.63	91.38
DeiT + FeatScale	72.346	91.22
NeuTRENO DeiT + FeatScale	73.23	91.73

Its corresponding asymmetric self-attention is obtained by replacing the key vectors k(i) with the query vectors q(i), i = 1, ..., N, and given by

$$\boldsymbol{u}(i) = \sum_{j=1}^{N} \operatorname{softmax} \left(\boldsymbol{q}(i)^{\top} \boldsymbol{k}(j) / \sqrt{D_{qk}} \right) \boldsymbol{v}(j) + \tilde{\lambda} (\boldsymbol{v}^{0}(i) - \boldsymbol{v}(i)), \quad i = 1, \dots, N.$$
 (25)

Leveraging Eqn. 25, we define the Neural Transformer with a Regularized Nonlocal Functional (NeuTRENO) as follows.

Definition 1 (Neural Transformer with a Regularized Nonlocal Functional (NeuTRENO)). Given a set of key and value vectors $\{k^{\ell}(j), v^{\ell}(j)\}_{j=1}^{N}$ in each layer ℓ , $\ell=1,\ldots,L$, for each query vector $q^{\ell}(i)$, $i=1,\ldots,N$, in the same layer, the self-attention unit at layer ℓ in a Neural Transformer with a Regularized Nonlocal Functional (NeuTRENO) computes the corresponding output vector $\mathbf{u}^{\ell}(i)$ of the query $\mathbf{q}^{\ell}(i)$ by the following attention formula:

$$\boldsymbol{u}^{\ell}(i) = \sum_{j=1}^{N} \operatorname{softmax} \left(\boldsymbol{q}^{\ell}(i)^{\top} \boldsymbol{k}^{\ell}(j) / \sqrt{D_{qk}} \right) \boldsymbol{v}^{\ell}(j) + \tilde{\lambda}(\boldsymbol{v}^{0}(i) - \boldsymbol{v}^{\ell}(i)), \quad i = 1, \dots, N.$$
 (26)

where $v^0(1), \dots v^0(N) \in \mathbb{R}^D$ are the value vectors in the first layer of NeuTRENO.

Fig. 2 illustrates the architecture of NeuTRENO.

Proposition 1. The evolution of u(x) under the dynamic in 21 does not converge to a constant vector.

Proposition 1 indicates that our NeuTRENO mitigates the over-smoothing issue, suggesting the benefit of our method. The proof for Proposition 1 is given in Appendix B.3.

4 Experimental Results

In this section, we empirically demonstrate the advantages of our proposed NeuTRENO approach across various tasks, including ImageNet classification [15], ADE20K image segmentation [73], and language modeling on the WikiText-103 [42]. Our aim to show: (i) NeuTRENO significantly outperforms the transformer baseline with softmax-attention defined in 2 across various tasks; moreover, NeuTRENO surpass FeatScale, a vision transformer that addresses over-smoothing, combining NeuTRENO with FeatScale is beneficial; (ii) the advantages of incorporating our proposed method with pre-trained models. We also demonstrate the benefits of our NeuTRENO in the symmetry setting and we point to Appendix D for the results. Throughout our experiments, we compare the performance of our proposed models with baselines of the same configuration. For additional details regarding datasets, models, and training procedures, please refer to Appendix A.

Object classification on ImageNet. To demonstrate the advantage of our NeuTRENO method, we compare it with the DeiT baseline [59] on the ImageNet image classification task. Our NeuTRENO DeiT surpasses the DeiT baseline, as shown in Table 1. Notably, our NeuTRENO DeiT achieves significantly higher performance in terms of both Top-1 Accuracy and Top-5 Accuracy. We also compare our method with FeatScale [65], a vision transformer model addressing over-smoothing (see Table 1). Our NeuTRENO significantly outperforms FeatScale, and combining NeuTRENO with FeatScale leads to substantial improvements. These results confirm the benefits of our model.

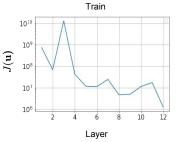
Image Segmentation on ADE20K dataset. To further validate the advantages of our proposed methods, we compare the performance of the Segmenter models [56] using the NeuTRENO DeiT

Table 2: Single-scale (SS) MIoU and multi-scale MIoU (MS) of the NeuTRENO DeiT vs. the DeiT on the ADE20K image segmentation.

Model/Metric	SS MIoU	MS MIoU (%)
Softmax DeiT	35.72	36.68
NeuTRENO DeiT	3 7.24	38.06

Table 3: Test and valid perplexity (Test PPL and Valid PPL) on WikiText-103 of NeuTRENO compared to the softmax transformer. Our proposed method achieves a significantly better performance PPL than the baseline.

Method/Metric	Valid PPL	Test PPL
Softmax Transformer	33.15	34.29
NeuTRENO	32.60	33.70



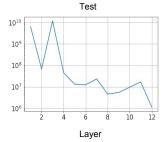


Figure 3: The average value of functional J(u) over 1000 training (Left) samples and test (Right) samples. When softmax attention is applied, the functional decreases as the depth of the trained DeiT increases. and DeiT backbones the on ADE20K image segmentation task [72], as shown in Table 2. The results demonstrate the substantial performance improvements achieved by utilizing the NeuTRENO DeiT backbone over the DeiT backbone, in terms of both single-scale (SS) MIoU and multi-scale (MS) MIoU metrics. These results strongly emphasize the effectiveness of our NeuTRENO approach in enhancing image segmentation performance.

Language Model on WikiText-103. In addition to computer vision tasks, we also evaluate the effectiveness of our model on a large-scale natural language processing application, specifically language modeling on WikiText-103. Our NeuTRENO language model demonstrates better performance in terms of both test perplexity and valid perplexity when compared to the softmax transformer language model [68]. These findings, combined with the results obtained across various tasks, empirically confirm the significant benefits of our NeuTRENO models.

Combine with pre-trained models. Furthermore, our proposed method is also beneficial to combine with pre-trained models. To empirically demonstrate that we incorporate NeuTRENO with pre-trained DeiT and fine-tune on the ImageNet dataset with one-third number of epochs that are used in training. The result is presented in Table 1, showing that combined with our method improves both the Top-1 and Top-5 accuracies of the pre-trained models.

5 Empirical Analysis

Applying Softmax-Attention Reduces the functional J(u). We present evidence supporting that the employment of softmax attention minimizes the functional J(u). Initially, we observe that the average cosine similarity between the numerical approximation of $\nabla_{\boldsymbol{u}}J(\boldsymbol{u})$ using symmetric or asymmetric kernel K(x,y) for both the trained Sym-DeiT (using symmetric self-attention 14) and DeiT models, closed 1, as shown in Table 4. This suggests that reversing the direction of the asymmetric approximation effectively decreases $J(\boldsymbol{u})$. Considering that softmax attention takes steps in this reversed direction numerically, its application leads to a reduction in $J(\boldsymbol{u})$. This is further substantiated by Fig. 3, which demonstrates a decrease in $J(\boldsymbol{u})$ as the depth of the trained DeiT increases when softmax attention is employed. More details of this analysis are in Appendix E

Over-smoothing Analysis. We empirically illustrate the effectiveness of NeuTRENOs in mitigating the over-smoothing problem in transformers. Fig. 1 compares the cosine similarity between token representations across layers for both NeuTRENO and softmax baseline models, specifically focusing on the Imagenet classification task (Left) and ADE20K image segmentation (Right). The token

Table 4: The average cosine similarity between the numerical approximation of $\nabla J(u)(x)$ using symmetric or asymmetric kernel K(x,y), for the trained Sym-DeiT and softmax DeiT models. The metric is evaluated on 1000 training and 1000 test data samples. The average score close to 1 shows a strong alignment between symmetric and asymmetric gradient approximations, suggesting that reversing the direction of the asymmetric approximation effectively reduces the functional J(u).

Model	Training data	Test data
Sym-DeiT	0.982	0.976
Softmax DeiT	0.973	0.964

features extracted by NeuTRENOs exhibit significantly lower similarity, particularly in the final layers. This finding highlights the ability of NeuTRENOs to address the over-smoothing issue and improve the diversity of token representations. We provide more details of this analysis in Appendix E.

6 Related Work

Over-smoothing in Transformers. Over-smoothing in deep transformers has been observed in various domains and applications from natural language processing [55] to computer vision [65, 18]. In vision tasks, [74] observes that the performance of the vision transformer (ViT [20]) quickly saturates as more layers are added to the model. Moreover, experiments in [74] show that the 32-layer ViT underperforms the 24-layer ViT, indicating the difficulty of ViTs in gaining benefits from deeper architectures. The authors point out that over-smoothing results in this phenomenon by causing the token representations to become identical when the model grows deeper. Based on this observation, the authors propose a cross-head communication method that helps enhance the diversity of both token representations and attention matrices. Furthermore, it has been shown in [60] that the training of ViT models encounters instability with greater depths. [25] proposes that this instability arises from the over-smoothing, where token representation for patches within an image becomes progressively alike as the model's depth increases. To explain this issue, [65] finds out that self-attention acts as a low-pass filter and smoothens the token representations in ViTs. This leads to the proposal of the FeatScale method [65], which regulates feature frequencies, whether low or high, to counteract the consequences of over-smoothing.

In addition, [55] observes the phenomenon in BERT [16], a deep language model, and explores over-smoothing through the graph perspective. The work utilizes hierarchical fusion strategies by preserving the output of self-attention through all layers, which is memory-costly. On the other hand, [65, 18] investigate over-smoothing in the image domain through the lens of Fourier spectrum, showing that self-attentions are low-pass filters, retaining only low-frequency, causing over-smoothed outputs. Our work is an orthogonal explanation of the previous work. We focus on developing a variational denoising framework to understand the self-attention of transformers as a gradient descent approximation of a functional. Our new finding explains the over-smoothing issue of transformers due to self-attention minimizing a functional and inspires us to derive the novel NeuTRENO method to overcome over-smoothing.

Nonlocal Functionals for Image Processing. Total variation [51] is well-known as an image-denoising technique. It denoises a noisy image by solving a constraint optimization problem. The method is also related to PDE-flow-based image-denoising techniques [24], namely isotropic and anisotropic diffusion [67] models. The method is edge preserving, meaning to avoid over-blurring edges' information [7]. Nonlocal functionals [35, 24] is considered as an extension of total variation to a nonlocal scale. Nonlocal functional and edge preservation properties are the motivation of our work to explain and overcome over-smoothing in transformers.

7 Concluding Remarks

In this paper, we establish a nonlocal variational denoising framework for self-attention. From this variational perspective, we explain over-smoothing in self-attention, which hinders the representation capacity of transformer models. We also derive the novel Neural Transformer with a Regularized Nonlocal Functional (NeuTRENO) to alleviate the over-smoothing. We empirically verify the benefits of NeuTRENO with a wide range of large-scale applications including ImageNet object classification, ADE20K object segmentation, and WikiText-103 language modeling. A limitation of our paper is that the privacy-preserving of NeuTRENO has not been addressed. It is interesting to explore if regularized nonlocal functional can also help improve the privacy-preserving of transformer models. We leave this exciting research idea as future work.

Acknowledgments and Disclosure of Funding

RGB acknowledges support from the NSF grants CCF-1911094, IIS-1838177, and IIS-1730574; ONR grants N00014-18-12571, N00014-20-1-2534, and MURI N00014-20-1-2787; AFOSR grant FA9550-22-1-0060; and a Vannevar Bush Faculty Fellowship, ONR grant N00014-18-1-2047. TMN acknowledges support from his start-up grant at the National University of Singapore (Grant Number: A-0009807-00-00).

References

- [1] Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. SemEval-2012 task 6: A pilot on semantic textual similarity. In *SEM 2012: The First Joint Conference on Lexical and Computational Semantics Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012), pages 385–393, Montréal, Canada, 7-8 June 2012. Association for Computational Linguistics.
- [2] Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy Guo, and Llion Jones. Character-level language modeling with deeper self-attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3159–3166, 2019.
- [3] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pages 6816–6826, 2021.
- [4] Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling. In *International Conference on Learning Representations*, 2019.
- [5] Alfonso S. Bandeira, Amit Singer, and Thomas Strohmer. Mathematics of Data Science. 2020.
- [6] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [7] A. Buades, B. Coll, and J. M. Morel. A review of image denoising algorithms, with a new one. *Multiscale Modeling & Simulation*, 4(2):490–530, 2005.
- [8] Kung-Ching Chang, Kelly Pearson, and Tan Zhang. Perron-frobenius theorem for nonnegative tensors. *Communications in Mathematical Sciences*, 6(2):507–520, 2008.
- [9] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. Advances in neural information processing systems, 34:15084–15097, 2021.
- [10] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv* preprint arXiv:1904.10509, 2019.
- [11] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar, October 2014. Association for Computational Linguistics.
- [12] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does BERT look at? an analysis of BERT's attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy, August 2019. Association for Computational Linguistics.
- [13] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, Florence, Italy, July 2019. Association for Computational Linguistics.

- [14] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Universal transformers. *arXiv preprint arXiv:1807.03819*, 2018.
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [18] Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *International Conference on Machine Learning*, pages 2793–2803. PMLR, 2021.
- [19] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021.
- [21] Leonhard Euler. *Institutiones calculi integralis*, volume 1. impensis Academiae imperialis scientiarum, 1792.
- [22] Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [23] Guy Gilboa and S. Osher. Nonlocal linear image regularization and supervised segmentation. *Multiscale Model. Simul.*, 6:595–630, 2007.
- [24] Guy Gilboa and S. Osher. Nonlocal operators with applications to image processing. *Multiscale Model. Simul.*, 7:1005–1028, 2008.
- [25] Chengyue Gong, Dilin Wang, Meng Li, Vikas Chandra, and Qiang Liu. Vision transformers with patch diversification. *arXiv* preprint arXiv:2104.12753, 2021.
- [26] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*, 2020.
- [27] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, 2021.
- [28] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021.
- [29] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.

- [30] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 15262–15271, 2021.
- [31] John Hewitt and Percy Liang. Designing and interpreting probes with control tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [32] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence modeling problem. *Advances in neural information processing systems*, 34:1273–1286, 2021.
- [33] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [34] Patrick Kahardipraja, Brielen Madureira, and David Schlangen. Towards incremental transformers: An empirical analysis of transformer models for incremental NLU. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1178–1189, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [35] Stefan Kindermann, S. Osher, and Peter W. Jones. Deblurring and denoising of images by nonlocal functionals. *Multiscale Model. Simul.*, 4:1091–1115, 2005.
- [36] Dimitrios Kotzias, Misha Denil, Nando de Freitas, and Padhraic Smyth. From group to individual labels using deep features. *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015.
- [37] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [38] Zhouhan Lin, Minwei Feng, Cícero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. CoRR, abs/1703.03130, 2017.
- [39] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [40] 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 *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pages 10012–10022, 2021.
- [41] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [42] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net, 2017.
- [43] Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2249–2255, Austin, Texas, November 2016. Association for Computational Linguistics.
- [44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.

- [45] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. *OpenAI report*, 2018.
- [46] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [47] 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. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [48] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics.
- [49] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [50] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15), 2021.
- [51] Leonid I. Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1):259–268, 1992.
- [52] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [53] Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight programmers. In *International Conference on Machine Learning*, pages 9355–9366. PMLR, 2021.
- [54] Matthias Seeger. Gaussian processes for machine learning. *International journal of neural systems*, 14(02):69–106, 2004.
- [55] Han Shi, JIAHUI GAO, Hang Xu, Xiaodan Liang, Zhenguo Li, Lingpeng Kong, Stephen M. S. Lee, and James Kwok. Revisiting over-smoothing in BERT from the perspective of graph. In *International Conference on Learning Representations*, 2022.
- [56] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7262–7272, 2021.
- [57] Ian Tenney, Dipanjan Das, and Ellie Pavlick. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy, July 2019. Association for Computational Linguistics.
- [58] Matthew Thorpe, Tan Minh Nguyen, Hedi Xia, Thomas Strohmer, A. Bertozzi, Stanley J. Osher, and Bao Wang. Grand++: Graph neural diffusion with a source term. In *International Conference on Learning Representations*, 2022.
- [59] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers distillation through attention. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 10347–10357. PMLR, 18–24 Jul 2021.
- [60] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 32–42, October 2021.

- [61] Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhutdinov. Transformer dissection: An unified understanding for transformer's attention via the lens of kernel. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 4344–4353, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [62] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information* processing systems, pages 5998–6008, 2017.
- [63] Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language model. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpret*ing Neural Networks for NLP, pages 63–76, Florence, Italy, August 2019. Association for Computational Linguistics.
- [64] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5797–5808, Florence, Italy, July 2019. Association for Computational Linguistics.
- [65] Peihao Wang, Wenqing Zheng, Tianlong Chen, and Zhangyang Wang. Anti-oversmoothing in deep vision transformers via the fourier domain analysis: From theory to practice. In *International Conference on Learning Representations*, 2022.
- [66] Zifeng Wang and Jimeng Sun. TransTab: Learning Transferable Tabular Transformers Across Tables. In *Advances in Neural Information Processing Systems (NeurIPS 2022)*, 2022.
- [67] Joachim Weickert, Wissenschaftlicher Werdegang, Steven Zucker, Allan Dobbins, Lee Iverson, Benjamin Kimia, and Allen Tannenbaum. Anisotropic diffusion in image processing. 01 1996.
- [68] Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and Vikas Singh. Nyströmformer: A Nyström-based Algorithm for Approximating Self-Attention. 2021.
- [69] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. arXiv preprint arXiv:1906.08237, 2019.
- [70] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. Deep learning based recommender system: A survey and new perspectives. ACM Computing Surveys (CSUR), 52(1):1–38, 2019.
- [71] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 16259–16268, 2021.
- [72] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017.
- [73] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset, 2018.
- [74] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer, 2021.

Supplement to "Mitigating Over-smoothing in Transformers via Regularized Nonlocal Functionals"

Table	of	Contents

A	Add	itional Details on the Experiments in Section 4	16
	A.1	Image Classification on Imagenet	17
	A.2	Image Segmentation on ADK20 dataset	17
	A.3	Language Modeling on WikiText-103	17
В	Tech	nnical Proofs	17
	B .1	Proof of Lemma 1	17
	B.2	Proof of Lemma 2	18
	B.3	Proof of Proposition 1	18
C	Deri	vation of Gradient of E as Given in Eqn. 21	19
D	Resu	ults of Symmetric Setting	19
E	Add	itional Details on the Empirical Analysis in Section 5	20
	E .1	Average Cosine Similarity between Gradient Approximations	20
	E.2	Average Value of Function	20
	E.3	Over-smoothing Analysis	20
F	Add	itional Experimental Results	21
	F.1	Object classification on Imagenet with DeiT-small baseline	21
	F.2	Beyond Softmax-Attention	21
G	Add	itional Empirical Analysis Results	21
	G .1	Visualizing Attention Matrices	21
	G.2	Head Redundancy between Layers	21
	G.3	NeuTRENO Inherently Mitigates Over-smoothing, even without Training the Models	21
	G .4	Efficiency Analysis	21
	G.5	Stability and Significance of NeuTRENO	22
	G.6	Robustness of NeuTRENO	23
	G.7	NeuTRENO in Incremental Learning	23
	G.8	Ablatation study on the choice of $\tilde{\lambda}$	23
	G .9	Scalability of NeuTRENO	23

A Additional Details on the Experiments in Section 4

This section provides datasets, models, and training details for experiments in Section 4. The code to reproduce our experimental results is included in our Supplementary Material submission.

A.1 Image Classification on Imagenet

Datasets and Metrics. The ImageNet dataset [15, 52] comprises 1.28 million training images and 50,000 validation images, encompassing the classification of 1000 categories. The evaluation metrics used for performance assessment are the top-1 and top-5 accuracies.

Models and Baselines. Our baseline model is the DeiT-tiny model [59], which consists of 12 transformer layers, 3 attention heads per layer, and a model dimension of 192. For model setting and setting and configuration, we follow [59]. Their implementation is available at https://github.com/facebookresearch/deit. The $\tilde{\lambda}$ used for our NeuTRENO method is 0.6.

A.2 Image Segmentation on ADK20 dataset

Datasets and Metrics. The ADE20K dataset is recognized for its inclusion of challenging scenes with fine-grained labels, making it one of the most demanding semantic segmentation datasets. The training set consists of 20,210 images encompassing 150 semantic classes. Additionally, there are 2,000 images in the validation set and 3,352 images in the test set. This in task the Single-scale mean Intersection over Union (SS mIoU) and the Multi-scale (MS mIoU).

Models and baselines. The training configuration and setting for our models are followed by [56]. The baseline model is finetuned with the pretrained DeiT-tiny backbone while our segmenter model used the pretrained NeuTRENO DeiT-tiny, with $\tilde{\lambda}=0.6$.

A.3 Language Modeling on WikiText-103

Datasets and Metrics. The WikiText-103 dataset consists of articles extracted from Wikipedia and is specifically designed to capture long contextual dependencies. The training set comprises approximately 28,000 articles, totaling 103 million running words. Each article contains text blocks consisting of approximately 3,600 words. The validation and test sets contain 218,000 and 246,000 running words, respectively, with each set consisting of 60 articles and approximately 268,000 words. Our experiment follows the standard setting [42,53], which involves dividing the training data into independent long segments of L words. For evaluation, we employ a batch size of 1 and process the text sequence using a sliding window of size L. When computing perplexity (PPL), we consider only the last position, except for the first segment where all positions are evaluated, following the approach in [2,53].

Models and baselines. For our language modeling implementation, we rely on the publicly available code https://github.com/IDSIA/Imtool-fwp developed by [53]. In our experiments, we set the dimensions of keys, values, and queries to 128, while the training and evaluation context length is set to 256. In this experiment, $\tilde{\lambda}=0.4$ yields the best performance of NeuTRENO language model.

B Technical Proofs

B.1 Proof of Lemma 1

For all i = 1, ..., N, we have $\mathbb{E}[\mathbf{B}^{(0)}(i)] = \mathbf{v}(i)$. Assume that $\mathbb{E}[\mathbf{B}^{(k)}(i)] = \mathbf{u}^{(k)}(i)$, then

$$\begin{split} \mathbb{E}[\mathbf{B}^{(k+1)}(i)] &= \sum_{j=1}^{N} \boldsymbol{v}(j) \mathbb{P}(\mathbf{B}^{(k+1)}(i) = \boldsymbol{v}(j)) \\ &= \sum_{j=1}^{N} \boldsymbol{v}(j) \sum_{l=1}^{N} \mathbb{P}(\mathbf{B}^{(k+1)}(i) = \boldsymbol{v}(j) | \mathbf{B}^{(1)}(i) = \boldsymbol{v}(l)) \mathbb{P}(\mathbf{B}^{(1)}(i) = \boldsymbol{v}(l)) \\ &= \sum_{j=1}^{N} \boldsymbol{v}_{j} \sum_{l=1}^{N} \mathbb{P}(\mathbf{B}^{(k)}(l) = \boldsymbol{v}(j)) \mathbb{P}(\mathbf{B}^{(1)}(i) = \boldsymbol{v}(l) | \mathbf{B}^{(0)}(i) = \boldsymbol{v}(i)) \\ &= \sum_{j=1}^{N} \boldsymbol{v}(j) \sum_{l=1}^{N} \mathbf{A}_{il} \mathbb{P}(\mathbf{B}^{(k)}(l) = \boldsymbol{v}(j)) \\ &= \sum_{l=1}^{N} \mathbf{A}_{il} \mathbb{E}[\mathbf{B}^{(k)}(l)] = \sum_{l=1}^{N} \mathbf{A}_{il} \boldsymbol{u}^{(k)}(l) \\ &= \boldsymbol{u}^{(k+1)}(i). \end{split}$$

Thus, by induction, we obtain the conclusion of the lemma.

B.2 Proof of Lemma 2

Since the transition matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ is right-stochastic, its largest eigenvalue is 1 (see Theorem 4.1 in [5]). Also, \mathbf{A} is a regular positive matrix since its elements are positive. Thus, the Perron-Frebenius theorem [8] implies the existence of a unique probability distribution $\boldsymbol{\pi}$, which is a positive left eigenvector of the transition matrix \mathbf{A} associated with its largest eigenvalue 1. In particular, in the case of symmetricity constraint, $\boldsymbol{\pi}$ can be chosen as follows

$$\pi = \left(\frac{d_1}{\sum_{j=1}^N d_j}, \frac{d_2}{\sum_{j=1}^N d_j}, \dots, \frac{d_n}{\sum_{j=1}^N d_j}\right),$$

where $d_i = \sum_{j=1}^N \exp\Bigl({m k}(i)^{ op} {m k}(j)/\sqrt{D_{qk}}\Bigr)$. It is easy to see that

$$\sum_{i=1}^{N} \pi_i \mathbf{A}_{ij} = \sum_{i=1}^{N} \frac{d_i}{\sum_{l=1}^{N} d_l} \frac{\exp\left(\mathbf{k}(i)^{\top} \mathbf{k}(j) / \sqrt{D_{qk}}\right)}{d_i}$$
$$= \frac{\sum_{i=1}^{N} \left(\exp\left(\mathbf{k}(i)^{\top} \mathbf{k}(j) / \sqrt{D_{qk}}\right)\right)}{\sum_{l=1}^{N} d_l}$$
$$= \frac{d_j}{\sum_{l=1}^{N} d_l} = \pi_j.$$

As a consequence, π must be the unique stationary distribution of the random walk $\{\mathbf{B}^{(k)}(i)\}_{k\in K}$. This concludes the proof.

B.3 Proof of Proposition 1

Recall from the gradient flow in Eqn 21, by using the method of Euler discretization, after k update steps starting from the initial u(x,0) = v(x) with adaptive stepsize $\Delta t = 1/\int_{\Omega} (k(x,y) + k(y,x)) dy$ and by choosing $\lambda = \tilde{\lambda}/\Delta t(x)$, we obtain the following

$$\mathbf{u}(x, k\Delta t(x)) = \mathbf{u}(x, (k-1)\Delta t(x)) - \Delta t(x)\nabla_{\mathbf{u}}J - \lambda \Delta t(x)\left(\mathbf{u}(x, (k-1)\Delta t(x)) - \mathbf{f}(x)\right)$$

$$= \int_{\Omega} \frac{K(x, y)\mathbf{u}(y, (k-1)\Delta t(x))}{\int_{\Omega} K(x, y')dy'}dy + \tilde{\lambda}\left(\mathbf{f}(x) - \mathbf{u}(x, (k-1)\Delta t(x))\right). \tag{27}$$

Discretizing $u(x, k\Delta t(x))$ and using Monte-Carlo approximation for the integrals in 27, we obtain

$$\mathbf{u}^{(k)}(i) = \sum_{j=1}^{N} \mathbf{A}_{ij} \mathbf{u}^{(k-1)}(j) + \tilde{\lambda} (\mathbf{f}(i) - \mathbf{u}^{(k-1)}(i)),$$
(28)

where \mathbf{A}_{ij} is computed using the keys and queries as either $\operatorname{softmax}\left(\mathbf{k}(i)^{\top}\mathbf{k}(j)/\sqrt{D_{qk}}\right)$ or $\operatorname{softmax}\left(\mathbf{q}(i)^{\top}\mathbf{k}(j)/\sqrt{D_{qk}}\right)$.

Suppose that $u^{(k)}(i)$, defined as Eqn. 28, converges to a constant vector \bar{u} as $k \to \infty$. We have

$$\mathbf{u}^{(k+1)}(i) - \mathbf{u}^{(k+1)}(j)
= \sum_{l=1}^{N} \mathbf{A}_{il} \mathbf{u}^{(k)}(l) - \sum_{l=1}^{N} \mathbf{A}_{jl} \mathbf{u}^{(k)}(l) + \tilde{\lambda}(\mathbf{u}^{(k)}(j) - \mathbf{u}^{(k)}(i)) + \tilde{\lambda}(\mathbf{f}(i) - \mathbf{f}(j))
= (\sum_{l=1}^{N} \mathbf{A}_{il} \mathbf{u}^{(k)}(l) - \mathbf{u}^{(k)}(i) \sum_{l=1}^{N} \mathbf{A}_{il}) - (\sum_{l=1}^{N} \mathbf{A}_{jl} \mathbf{u}^{(k)}(l) - \mathbf{u}^{(k)}(j) \sum_{l=1}^{N} \mathbf{A}_{jl})
+ (\tilde{\lambda} - 1)(\mathbf{u}^{(k)}(j) - \mathbf{u}^{(k)}(i)) + \tilde{\lambda}(\mathbf{f}(i) - \mathbf{f}(j))
= \sum_{l=1}^{N} \mathbf{A}_{il}(\mathbf{u}^{(k)}(l) - \mathbf{u}^{(k)}(i)) - \sum_{l=1}^{N} \mathbf{A}_{jl}(\mathbf{u}^{(k)}(l) - \mathbf{u}^{(k)}(j)) + (\tilde{\lambda} - 1)(\mathbf{u}^{(k)}(j) - \mathbf{u}^{(k)}(i))
+ \tilde{\lambda}(\mathbf{f}(i) - \mathbf{f}(j))$$
(29)

$$\text{Since } \boldsymbol{u}^{(k)}(i) \to \bar{\boldsymbol{u}} \text{, for } i = 1, 2, \dots, N \text{, as } k \to \infty \text{, we have} \begin{cases} (\boldsymbol{u}^{(k+1)}(i) - \boldsymbol{u}^{(k+1)}(j)) \to \boldsymbol{0} \\ (\boldsymbol{u}^{(k)}(l) - \boldsymbol{u}^{(k)}(i)) \to \boldsymbol{0} \\ (\boldsymbol{u}^{(k)}(l) - \boldsymbol{u}^{(k)}(j)) \to \boldsymbol{0} \\ (\boldsymbol{u}^{(k)}(j) - \boldsymbol{u}^{(k)}(i)) \to \boldsymbol{0} \end{cases}$$

as $k \to \infty$. This is a contradiction since while the LHS of 29 approaches 0, its RHS approaches $\tilde{\lambda}(f(i) - f(j))$, which is not 0 in general. Thus, we obtain the conclusion of Proposition 1.

C Derivation of Gradient of E as Given in Eqn. 21

Taking the gradient of E(u, f) with respect to u, we obtain

$$\nabla_{\boldsymbol{u}}E = \nabla_{\boldsymbol{u}}J + \left[\frac{\partial G}{\partial u_1}, \frac{\partial G}{\partial u_2}, \dots, \frac{\partial G}{\partial u_D}\right]^T.$$
(30)

The partial derivative $\partial G/\partial u_j$, $j=1,2,\ldots,D$, is defined through its dot product with an arbitrary function $h_j \in L^2(\Omega)$ as follows

$$\begin{split} \frac{\partial G}{\partial u_j} \cdot h_j(x) &= \frac{d}{d\tau} G(u_j + \tau h_j) \big|_{\tau=0} \\ &= \frac{\lambda}{2} \left(\frac{d}{d\tau} \int_{\Omega} (u_j(x) - f_j(x) + \tau h_j(x))^2 dx \right) \Big|_{\tau=0} \\ &= \lambda \int_{\Omega} (u_j(x) - f_j(x)) h_j(x) dx. \end{split}$$

Thus, the Frechet derivative of F with respect to u_i is given by

$$\frac{\partial G}{\partial u_j} = \lambda(u_j(x) - f_j(x)) \tag{31}$$

Substituting the formula for $\partial G/\partial u_j$ in Eqn. 31 into Eqn. 30 for $\nabla_{\boldsymbol{u}} E(\boldsymbol{u}, \boldsymbol{f})$, we obtain the following gradient flow

$$\frac{d\boldsymbol{u}(x,t)}{dt} = -\nabla_{\boldsymbol{v}}E(\boldsymbol{u},\boldsymbol{f}) = -\nabla_{\boldsymbol{u}}J(\boldsymbol{u})(x) + \lambda(\boldsymbol{f}(x) - \boldsymbol{u}(x)), \tag{32}$$

where t is a dummy time variable and $-\nabla_{\boldsymbol{u}}J(\boldsymbol{u})$ is defined as in 9.

D Results of Symmetric Setting

In this section, we show that NeuTRENO significantly improves the performance of a symmetric transformer baseline, which utilizes symmetric self-attention. We refer to the DeiT with symmetric attention, defined in 14, as Sym-DeiT and the Sym-DeiT combined with our NeuTRENO method as Sym-NeuTRENO DeiT.

Object classification on Imagenet To further illustrate the advantage of our NeuTRENO method, we compare Sym-NeuTRENO DeiT with the Sym-DeiT baseline on the ImageNet image

Table 5: Top-1 and Top-5 accuracy (%) of Sym-NeuTRENO DeiT vs. Sym-DeiT on the ImageNet classification task. The Sym-NeuTRENO DeiT models significantly outperform the Sym-DeiT in terms of accuracy, indicating the benefit of NeuTRENO method.

Model/Metric	Top-1 Acc (%)	Top-5 Acc (%)
Sym-DeiT	71.14	90.54
Sym-NeuTRENO DeiT	72.07	91.22

Table 6: Single-scale (SS) MIoU and multi-scale (MS) MIoU of the Sym-NeuTRENO DeiT vs. Sym-DeiT. The Sym-NeuTRENO DeiT model is beneficial since they significantly outperform the Sym-DeiT.

Model/Metric	SS MIoU	MS MIoU (%)
Sym-DeiT	35.18	36.00
Sym-NeuTRENO DeiT	35.68	36.39

classification task. Our Sym-NeuTRENO DeiT outperforms the Sym-DeiT baseline, as shown in Table 5. Notably, the Sym-NeuTRENO DeiT achieves higher performance in terms of both top-1 accuracy and top-5 accuracy than Sym-DeiT baseline. These results further confirm the benefits of our proposed NeuTRENO model.

Image Segmentation on ADE20K dataset We also compare the performance of the Segmenter models [56] using the Sym-NeuTRENO DeiT backbone with models using the Sym-DeiT backbone on ADE20K image segmentation [72], as shown in Table 6. The results demonstrate the substantial performance improvements achieved by utilizing the Sym-NeuTRENO DeiT backbone compared to the Sym-DeiT backbone in terms of both single-scale (SS) MIoU and multi-scale (MS) MIoU metrics. This result further validates the advantages of our NeuTRENO models in enhancing image segmentation performance in the symmetric setting.

E Additional Details on the Empirical Analysis in Section 5

In this section, we provide the details for the empirical analysis in Section 5.

E.1 Average Cosine Similarity between Gradient Approximations

To produce the results in Table 4, we derive the approximation for the gradient $\nabla_{\boldsymbol{u}} J(\boldsymbol{u})$, from Eqn 9, at time t=0:

$$\nabla_{\boldsymbol{u}}J(\boldsymbol{u}) = \int_{\Omega} (\boldsymbol{u}(x,0) - \boldsymbol{u}(y,0))K(x,y)dy = \int_{\Omega} (\boldsymbol{v}(x) - \boldsymbol{v}(y))K(x,y)dy,$$

where K(x,y) := k(x,y) + k(y,x). Using Monte-Carlo approximation for the integral and choosing $K(x,y) = \exp(\mathbf{k}(x)^T\mathbf{k}(y)/\sqrt{D_{qk}})$, the symmetric approximation of the gradient is derived as $\sum_{j=1}^N (\mathbf{v}(i) - \mathbf{v}(j)) \exp(\mathbf{k}(i)^T\mathbf{k}(j)/\sqrt{D_{qk}})$. Otherwise, by choosing $K(x,y) = \exp(\mathbf{q}(x)^T\mathbf{k}(y)/\sqrt{D_{qk}})$, the assymmetric approximation of the gradient is derived as $\sum_{j=1}^N (\mathbf{v}(i) - \mathbf{v}(j)) \exp(\mathbf{q}(i)^T\mathbf{k}(j)/\sqrt{D_{qk}})$. In this analysis, we take the dot product between the symmetric and asymmetric approximation of the gradient $\nabla_{\mathbf{u}}J(\mathbf{u})$ and average these dot products over positions. We finally report the average cosine similarity over 1000 training data and 1000 test data, as shown in Table 4.

E.2 Average Value of Function

In order to report the average value of function J(u) in Fig. 3, we follow the process of computing J(u) for 1000 data points for each transformer block. Subsequently, the average value is reported for each layer. This procedure is carried out for both the training and test datasets.

E.3 Over-smoothing Analysis

The average cosine similarity between all pairs of token's representations (x_i, x_j) in a sequence is computed as

$$\frac{1}{N(N-1)} \sum_{i \neq j} \frac{\boldsymbol{x}_i^T \boldsymbol{x}_j}{\|\boldsymbol{x}_i\|_2 \|\boldsymbol{x}_j\|_2}.$$

The result is then averaged over 1000 randomly chosen test data in ImageNet and ADE20K. The result is then reported for each layer, as in Fig. 1.

Table 7: Top-1 and Top-5 accuracy (%) of NeuTRENO DeiT-small vs. DeiT-small on the ImageNet benchmark. The NeuTRENO DeiT-small significantly outperform the DeiT-small in terms of accuracy. We also compare NeuTRENO DeiT-small with DeiT plus FeatScale, a vision transformer model that addresses over-smoothing, showing the advantage of NeuTRENO. The accuracies reported in [59] for DeiT-small and [65] for DeiT-small plus FeatScale, respectively, are in parentheses.

Model/Metric	Top-1 Acc (%)	Top-5 Acc (%)
DeiT-small DeiT-small + FeatScale NeuTRENO DeiT-small	79.97 (79.9) 79.96 (80.9) 80.68	95.05 (95.0) 95.06 95.30

Table 8: Accuracy of NeuTRENO vs.Kernel Transformerr on the CIFAR-10 dataset [37]. The NeuTRENO model significantly outperforms the in terms of accuracy.

Model/Metric	Accuracy (%)
Kernel Transformer	75.89
NeuTRENO	76.75

F Additional Experimental Results

F.1 Object classification on Imagenet with DeiT-small baseline

In this section, we show the advantages of our method when we further scale up the model by doubling the model dimension and the number of heads compared to that of the DeiT-tiny. In particular, the NeuTRENO DeiT-small achieves better results in both Top-1 Accuracy and Top-5 Accuracy, as shown in Table 7. Our method also outperforms DeiT plus FeatScale. Here, we did our best to reproduce the results of DeiT-small plus FeatScale [65]. In Table 7, we include our reproduced results and the results reported in [59] for DeiT-small and [65] for DeiT-small plus FeatScale, respectively.

F.2 Beyond Softmax-Attention

We show that NeuTRENO can be combined with other baseline attention mechanisms other than softmax attention. In particular, our NeuTRENO significantly improves transformer-based models with kernel attention [54, 61], on the CIFAR-10 image classification task [37], as shown in Table 8. This further confirms the benefits of our model. Here, both models share the same configuration regarding training, the model's size, and the model's depth (12 layers).

G Additional Empirical Analysis Results

This section provides extra empirical analysis to further demonstrate the benefits of NeuTRENO models in mitigating over-smoothing.

G.1 Visualizing Attention Matrices

Fig. 4 displays the 3-head attention matrices obtained from layer [1,6,12] of both the pre-trained NeuTRENO DeiT-tiny and the DeiT-tiny baseline models, using a random sample from the ImageNet dataset.

G.2 Head Redundancy between Layers

NeuTRENO mitigates head redundancy between layers, particularly in the final transformer layers where over-smoothing is most pronounced. Fig. 5 shows the average cosine similarity of attention matrices between two successive layers, over 1000 randomly sampled data. The trained NeuTRENO DeiT obtains lower cosine similarity than that of the trained DeiT as the model depth increases.

G.3 NeuTRENO Inherently Mitigates Over-smoothing, even without Training the Models

Randomly-initialized NeuTRENO DeiT-tiny significantly reduces the average cosine similarity between token representations of 12-layer randomly-initialized DeiT-tiny model, as shown in Fig. 6, on the Imagenet classification task. This observation highlights the ability of our NeuTRENO models in mitigating over-smoothing.

G.4 Efficiency Analysis

We report the ratios of the floating-point operations per second (FLOPs), the inference memory, and the inference real-time running of NeuTRENO DeiT vs. DeiT per sample on the ImageNet dataset, which are 1.00005, 1.000002, 1.00013, respectively. This indicates that the significant gain in the performance of NeuTRENO does not come with the cost of efficiency.

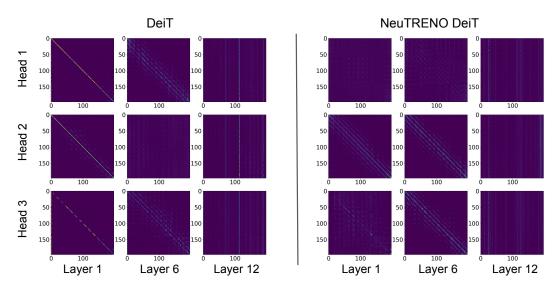


Figure 4: Plot of attention matrices attained from layer [1, 6, 12] of both the pretrained DeiT-tiny baseline (Left) and the NeuTRENO DeiT-tiny (Right) models, for each head, using a random sample from the Imagenet dataset.

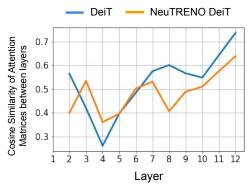


Figure 5: The average cosine similarity of attention matrices between two successive layers, over 1000 randomly sampled data, of the trained NeuTRENO DeiT and trained DeiT models on the Imagenet classification task.

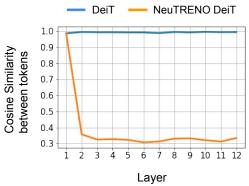


Figure 6: The average cosine similarity between token representations of 12-layer randomly-initialized NeuTRENO DeiT and DeiT models, on the Imagenet classification task. Here, 1000 data are randomly sampled for the analysis.

G.5 Stability and Significance of NeuTRENO

To further confirm the stability and significance of NeuTRENO's performance, we provide the standard deviations from five runs for both the NeuTRENO and baseline models for each experiment (in the main text) in Tables 9, 10, 11.

Table 9: Means and standard deviations over five runs with different random seeds of models trained on the Imagenet Classification task.

Model/Metric	Top-1 Acc (%)	Top-5 Acc (%)
Softmax DeiT-Tiny	72.17 ± 0.07	91.02 ± 0.04
NeuTRENO DeiT-Tiny	73.01 ± 0.09	91.56 ± 0.05
NeuTRENO Adaptation	72.63 ± 0.07	91.38 ± 0.03
DeiT-Tiny + FeatScale	72.346 ± 0.06	91.22 ± 0.04
NeuTRENO DeiT-Tiny + FeatScale	$\textbf{73.23} \pm \textbf{0.08}$	$\textbf{91.73} \pm \textbf{0.05}$

Table 10: Means and standard deviations over five runs with different random seeds of models trained on the ADE20K image segmentation task

Metric/Model	Pretrained Softmax Deit-Tiny	Pretrained NeuTRENO DeiT-Tiny
SS MIoU MS MIoU	$35.72 \pm 0.57 \\ 36.68 \pm 0.42$	$37.24 \pm 0.62 \ 38.06 \pm 0.54$

Table 11: Means and standard deviations over five runs with different random seeds of models trained on the WikiText-103 language model task.

Metric/Model	Softmax Transformer	NeuTRENO
Valid PPL Test PPL	$ \begin{vmatrix} 33.15 \pm 0.07 \\ 34.29 \pm 0.09 \end{vmatrix} $	$\begin{array}{c} \textbf{32.60} \pm \textbf{0.08} \\ \textbf{33.70} \pm \textbf{0.07} \end{array}$

Table 12: Evaluation of NeuTRENO DeiT-Tiny vs. Softmax DeiT-Tiny on the ImageNet-C (mean corruption error mCE), Imagenet-A (Accuracy), and Imagenet-R (Accuracy) datasets.

Model/Dataset (Metric)	Imagenet-C (mCE)	Imagenet-A (Accuracy %)	Imagenet-R (Accuracy %)		
Softmax DeiT-Tiny	71.6	6.9	32.83		
NeuTRENO-DeiT-Tiny	70.1	8.2	33.82		

Table 13: Ablation study of different values hyperparameter $\tilde{\lambda}$ of NeuTRENO DeiT-Tiny on the ADE20K Image Segmentation task.

Metric/Model	Baseline	$\tilde{\lambda} = 0.1$	$\tilde{\lambda} = 0.2$	$\tilde{\lambda} = 0.4$	$\tilde{\lambda} = 0.5$	$\tilde{\lambda} = 0.6$	$\tilde{\lambda} = 0.8$	$\tilde{\lambda} = 1.0$	$\tilde{\lambda} = 2.0$
	35.72 36.68								

G.6 Robustness of NeuTRENO

In addition to the standard metrics, we evaluate the robustness of our NeuTRENO model compared to the baseline transformer model, particularly under adversarial examples and for out-of-distribution generalization. Table 12 demonstrates that NeuTRENO DeiT-Tiny is consistently more robust than the DeiT-Tiny baseline on the Imagenet-C (common data corruption and perturbations, such as adding noise and blurring the images) [29], Imagenet-A (adversarial examples) [30], and Imagenet-R (out of distribution generalization) [28] datasets, which are widely used to test the model's robustness.

G.7 NeuTRENO in Incremental Learning

In an incremental learning setting [34], our 8-layer NeuTRENO achieves 1.97% higher accuracy on the sentiment classification task [36] than the 8-layer baseline transformer.

G.8 Ablatation study on the choice of λ

We also conduct an ablation study on the impact of the hyperparameter. In particular, on the ADE20K image segmentation task, we train NeuTRENO with different values. We summarize our results in Table 13. Our findings reveal that within the range of [0.2,1], NeuTRENO consistently outperforms the softmax baseline. However, when values become small or big (below 0.2 or above 1, respectively), NeuTRENO's performance declines.

G.9 Scalability of NeuTRENO

To demonstrate the scalability of our proposed model, we conduct additional experiments to show that our NeuTRENO method can effectively mitigate the oversmoothing issue in the BERT-base model. In particular, in Figure 7 (Left), we plot the cosine similarity between token representations across layers of a pre-trained BERT-base model [17] on the SQuAD v1.1 question answering task [48] and observe

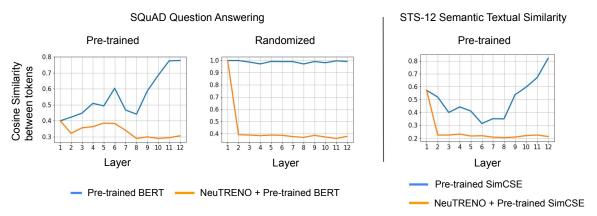


Figure 7: The average cosine similarity between token representations of 12-layer trained (Left) and randomly-initialized (Middle) BERT-base and NeuTRENO BERT-base on the SQuAD question answering task. We also plot the same cosine similarity scores for the trained SimCSE and NeuTRENO SimCSE models (Right) on the STS-12 semantic textual similarity task. Here, 1000 and 500 data are randomly sampled for the analysis on the SQuAD and STS-12 datasets, respectively.

the presence of the oversmoothing issue as the model gets deeper, causing tokens to become identical. We then apply NeuTRENO on the same pre-trained BERT model, and without any fine-tuning, we observe a significant reduction in the cosine similarity between token embeddings in each layer (see Figure 7 (Left)), indicating that NeuTRENO effectively mitigates the oversmoothing problem in BERT. Additionally, our NeuTRENO BERT finetuned on the task yields better accuracy than the finetuned BERT (81.39 exact match score and 88.62 F1-score vs. 80.77 exact match score and 88.12 F1-score). Moreover, we have conducted the same analysis for a randomized BERT-base model and a randomized NeuTRENO BERT-base model and obtained the same encouraging results (see Figure 7 (Middle)). These results further suggest that NeuTRENO helps alleviate the over-smoothing issue in large-scale transformer models.

We also obtain additional results and show that our NeuTRENO SimCSE, after fine-tuned on the STS-12 semantic textual similarity task [1], gains a significant improvement over the baseline SimCSE [22], which is also fine-tuned on the same task (77.32% vs. 75.29% Spearman's correlation. Here, the higher correlation, the better). This additional result further verifies that decreasing the cosine dissimilarity between tokens within trained transformer-based models leads to improved empirical performance.