

A TWO-CHARACTER CHANGE IN TRANSFORMER ARCHITECTURE PROMOTES IDEAL TOKEN GEOMETRY

Anonymous authors

Paper under double-blind review

ABSTRACT

We hypothesize that in the optimal geometric configuration of token embeddings for transformer classifiers, tokens should collapse to single points according to their classes, and these points themselves should exhibit Neural Collapse. We study whether current transformers achieve this configuration through principal component projections, cosine similarity measurements, analysis of variance on token embeddings, and Neural Collapse measurements and find that they fall far short of the conjectured ideal. To address this, we introduce a simple modification to attention that brings token embeddings markedly closer to the conjectured configuration and yields consistent performance improvements across benchmarks.

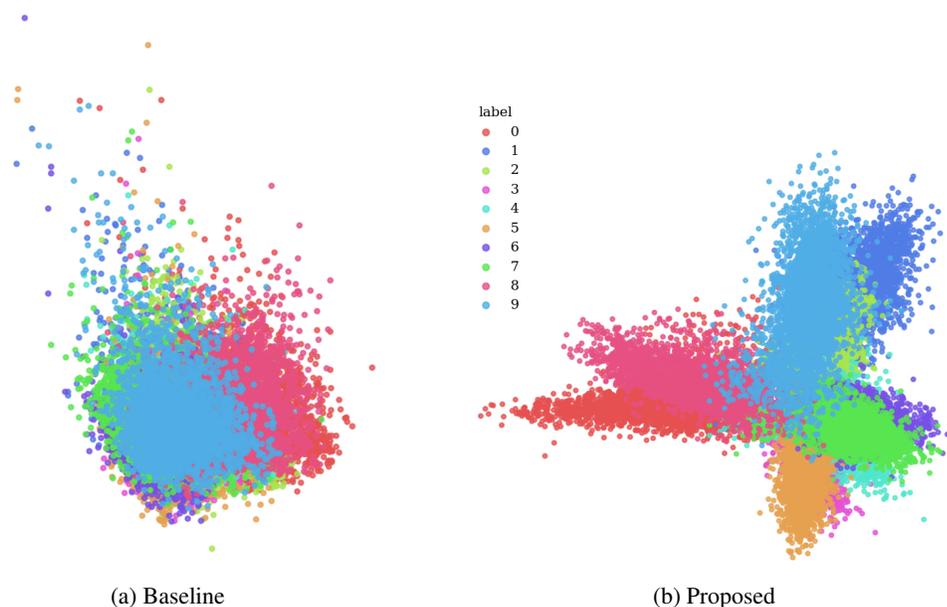


Figure 1: Proposed modification cleanly separates tokens into clusters according to class.

1 INTRODUCTION

1.1 MOTIVATION AND PROBLEM STATEMENT

The transformer architecture has achieved remarkable success across a wide range of tasks since its introduction by Vaswani et al. (2017). From image classification to language modeling, many of these tasks reduce to classification problems in which the transformer learns input-token embeddings and then applies a final-layer classifier to predict the class of the input sequence.

In the context of deep networks, Pappayan et al. (2020) showed that as training approaches zero error, the final-layer representations of standard classifiers exhibit a distinctive geometry—termed **Neural Collapse (NC)**—with four defining properties:

- 054 **NC1** Final-layer representations collapse to their class means.
 055 **NC2** The class means converge to a simplex equiangular tight frame (ETF).
 056 **NC3** The linear classifier weights align with the corresponding class means.
 057 **NC4** The classifier predicts each class via nearest-class-mean decision boundaries.
 058
 059

060 While these phenomena have been extensively studied in convolutional networks for image classification (Papayan et al., 2020; Han et al., 2022), transformers present a fundamentally different setting: inputs are tokenized sequences, and the learned representations consist of token embeddings rather than a single feature vector. This naturally raises the following questions:

- 064 **Q1:** *What is the ideal geometric configuration of token embeddings for classification?*
 065 **Q2:** *Do current transformers achieve this configuration?*
 066 **Q3:** *If not, can architectural modifications induce it?*
 067

068 At its core, the transformer consists of a sequence of attention layers (Vaswani et al., 2017). Several theoretical works show that attention leads to **rank collapse**, where token embeddings converge to a single vector (Dong et al., 2023; Noci et al., 2022; Geshkovski et al., 2024). These results indicate an inherent inductive bias of attention toward collapsing tokens.
 072

073 1.2 CONTRIBUTIONS

074 To address Question Q1, we draw on insights from NC and prior work on rank collapse to introduce Neural Token Collapse (NTC), our hypothesized ideal geometry for token embeddings in classification tasks. NTC is defined formally in Section 3 and characterized by three core properties:

- 075 **NTC0** Tokens within the same sequence collapse to a single vector.
 076 **NTC1** All sequences of the same class collapse to the class mean.
 077 **NTC2-4** The class means and classifier weights satisfy NC2, NC3 and NC4.
 078

079 To address Question Q2, in Section 4, we assess to what extent the token embeddings of current transformers conform to NTC by applying principal component analysis, cosine similarity, analysis of variance, and NC measurements. Our findings show that standard transformers fail to achieve NTC.
 087

088 Finally, to address Question Q3, we introduce in Section 2 a simple architectural modification designed to overcome this limitation. It is motivated by the observation that standard transformers uses attention to compute a contextual mean and adds it to each token, increasing the tokens’ mean relative to their variance. Layer normalization then projects the tokens back onto a sphere, causing them to cluster more closely together. This process is inefficient because it introduces norm-changing radial movements that are subsequently removed by layer normalization (Figure 2). Our modification moves the tokens directly along their variance directions, so the tokens can collapse without relying on the intermediate mean shift and projection (Figure 2). As shown in Section 4, this approach promotes NTC more effectively and consistently improves performance across multiple datasets.
 097

098 2 PROPOSED ARCHITECTURE MODIFICATION

099 In this section, we propose our modification to the transformer architecture, which we motivate by analyzing how attention affects the geometry of token embeddings.
 102

103 2.1 STANDARD ATTENTION MECHANISM

104 Consider a sequence of T token embeddings (x_1, \dots, x_T) with each $x_i \in \mathbb{R}^d$, collected into a matrix $X \in \mathbb{R}^{T \times d}$. Attention (**Attn**) projects X into queries, keys, and values using trainable matrices $W_Q, W_K, W_V \in \mathbb{R}^{d \times d_k}$:

$$105 Q = XW_Q, \quad 106 K = XW_K, \quad 107 V = XW_V.$$

108 It then computes scaled dot-product weights and outputs a weighted sum of the values:

$$109 P = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right), \quad \text{Attn}(X) = PV. \quad (1)$$

110 Building on this definition, multi-head attention (MHA) computes h attention heads in parallel, concatenates their outputs, then applies a trainable projection matrix $W_o \in \mathbb{R}^{hd_k \times d}$. The result is added back to the input via a residual connection (He et al., 2015):

$$111 X' = X + \text{MHA}(X). \quad (2)$$

112 After this residual update, the token embeddings are passed through LayerNorm (Ba et al., 2016), which projects the tokens back onto a sphere with radius \sqrt{d} . This process repeats at each layer.

113 **Interpretation.** Equation 2 increases the tokens' mean relative to their variance, so the subsequent renormalization draws them closer together on the sphere. Figure 2 illustrates this process: token embeddings \vec{x}_1, \vec{x}_2 (black arrows) are shifted by their attention-weighted means $(PX)_1, (PX)_2$ (blue arrows) along the orange arrows, and after projection they end up closer together on the sphere (green points).

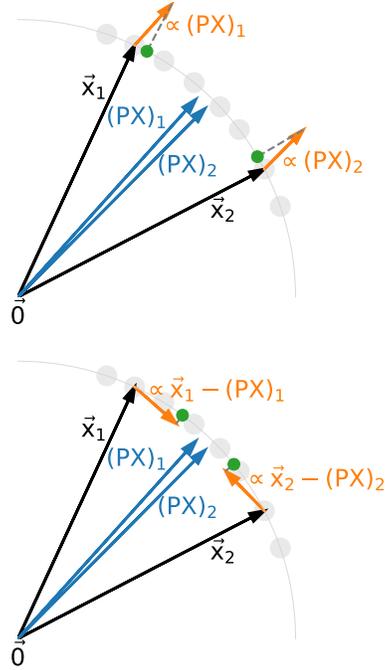


Figure 2: Baseline (top) and proposed (bottom)

114 2.2 THE TWO-CHARACTER CHANGE

115 We modify the standard attention mechanism into a *Laplacian mechanism*, denoted by \mathcal{L} :

$$116 P = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right), \quad \mathcal{L}(X) = V - PV. \quad (3)$$

117 This modification amounts to adding two extra characters, “-” and “V”, hence the title. The term “Laplacian” will be justified in Section 5.

118 **Remark.** *It is important to emphasize that this modification differs from the standard skip connection applied around Attn (or MHA), as the residual stream carries an additional XW_v term.*

119 **Interpretation.** While **Attn** computes the mean of tokens, \mathcal{L} computes the difference between tokens and the mean. If every attention head is replaced with \mathcal{L} , then Equation 2 allows each token to move directly toward the mean. This eliminates the redundant radial movement later canceled by projection. Figure 2 shows this effect: token embeddings \vec{x}_1, \vec{x}_2 (black arrows) move directly towards the means $(PV)_1, (PV)_2$ (blue arrows) along the variance direction (orange arrows). The green points are the projected tokens, which end up much closer together with the projection step playing only a minor role. [We provide empirical evidence for this geometric interpretation in Appendix D.](#)

120 2.3 MIXING BOTH MECHANISMS

121 If one is willing to change more than two characters, then the mechanism can be extended further by mixing **Attn** and \mathcal{L} in transformers. Specifically, we propose using \mathcal{L} for some heads and **Attn** for others. There are many possible ways to do this and we find the following two simple strategies to be effective:

- 122 1. Fix $m < h$. In each layer, assign **Attn** to m heads and \mathcal{L} to the remaining $h - m$ heads.
- 123 2. Assign **Attn** to all heads in the first half of the transformer blocks and \mathcal{L} to all heads in the second half.

124 Both strategies are straightforward to implement and introduce no additional trainable parameters. We describe more mixing strategies we have tried in Appendix E.

Interpretation. This setup allows some heads to move tokens toward their mean while others can move tokens along their variance direction, giving the model greater flexibility to steer tokens toward the ideal geometry.

3 IDEAL TOKEN GEOMETRY

Before formally defining our conjectured ideal token geometry, we need to define several quantities that describe the distribution of token embeddings. We consider a dataset D with N data points and C classes, where class c contains N_c data points so that $N = \sum_{c=1}^C N_c$. Each data point is tokenized into a sequence of T tokens. For indices $1 \leq c \leq C$, $1 \leq i \leq N_c$, and $1 \leq t \leq T$, let $X_{t,i,c}$ denote the embedding of the t -th token of the i -th data point in class c .

3.1 TOKEN MEANS

We first define several means. For each data point we compute the *sequence token mean*, which is the average of its T token embeddings:

$$\mu_{i,c} = \mathbf{Ave}_t X_{t,i,c}.$$

For each class, we compute the *class token mean*, obtained by averaging the sequence means over all N_c data points in that class:

$$\mu_c = \mathbf{Ave}_i \mu_{i,c} = \mathbf{Ave}_{t,i} X_{t,i,c}.$$

Finally, we compute the *global token mean* by averaging the class means over all C classes:

$$\mu_G = \mathbf{Ave}_c \mu_c = \mathbf{Ave}_{t,i,c} X_{t,i,c}.$$

3.2 TOKEN VARIANCES

Using $\mu_{i,c}$, μ_c , and μ_G , we define several measures of variance that quantify different aspects of token distribution.

The *within-sequence variance* measures how far individual token embeddings of a single data point deviate from the token mean of that data point itself on average:

$$\text{WithinSeqVar} = \mathbf{Ave}_{t,i,c} \|X_{t,i,c} - \mu_{i,c}\|^2.$$

The *within-class variance* measures how far each sequence’s token mean deviates from the corresponding class token mean on average:

$$\text{WithinClassVar} = \mathbf{Ave}_{i,c} \|\mu_{i,c} - \mu_c\|^2.$$

The *between-class variance* measures how far each class token mean deviates from the global token mean:

$$\text{BetweenClassVar} = \mathbf{Ave}_c \|\mu_c - \mu_G\|^2.$$

Finally, the *total variance* quantifies the overall variance of all token embeddings around the global token mean:

$$\text{TotalVar} = \mathbf{Ave}_{t,i,c} \|X_{t,i,c} - \mu_G\|^2.$$

3.3 HIERARCHICAL VARIANCE DECOMPOSITION

By expanding the squared norm, the total variance decomposes additively into the three components above:

$$\text{TotalVar} = \text{BetweenClassVar} + \text{WithinClassVar} + \text{WithinSeqVar}.$$

This decomposition mirrors an ANOVA-style breakdown of the total energy into within-sequence, within-class, and between-class contributions.

3.4 NEURAL TOKEN COLLAPSE

Using the means and variances defined in the previous section, we formally define *Neural Token Collapse* (NTC), our conjectured ideal token geometry for classification. This regime is characterized by the following properties:

NTC0 All tokens within a sequence coincide with that sequence’s mean,

$$X_{t,i,c} = \mu_{i,c} \quad \forall 1 \leq c \leq C, 1 \leq t \leq T, 1 \leq i \leq N_c,$$

which implies

$$\text{WithinSeqVar} = 0.$$

NTC1 All sequence means within a class coincide with the class mean,

$$\mu_{i,c} = \mu_c \quad \forall 1 \leq i \leq N_c,$$

which implies

$$\text{WithinClassVar} = 0.$$

Combined with the token-level collapse above, this further implies

$$\text{TotalVar} = \text{BetweenClassVar}.$$

NTC2-4 The class means

$$\{\mu_c : 1 \leq c \leq C\},$$

together with the classifier weights, satisfy NC2, NC3, and NC4 as defined in Papyan et al. (2020).

4 EXPERIMENTS

4.1 SETUP

Our experiments are based on the DeiT-3 vision transformer (Touvron et al., 2022), which is widely considered a strong baseline for image classification. By default, we applied the modification described in Section 2.3 to the ViT-B model. **For strategy 1, we focus mainly on three options: $m = 0, 1, 3$. Note that $m = 0$ is equivalent to using the Laplacian mechanism for all heads, and we discuss the impact of different values of m in Section 4.6. Throughout the paper, we denote these three options as “ViT-B-0P”, “ViT-B-1P” and “ViT-B-3P”, respectively.** We refer to the baseline vision transformer as “ViT-B” and to the approach described in strategy 2 as “ViT-B-Mix-Depth.”

We trained these models on CIFAR-10, CIFAR-100 (Krizhevsky, 2009), and ImageNet-1k (Deng et al., 2009). For each dataset, all models shared the same training recipe. Full details of the training setup is provided in Appendix B.

In the following subsections, we use a range of metrics to analyze the token geometry of the models, assessing whether the proposed modifications promote NTC more relative to the ViT-B baseline.

4.2 PRINCIPAL COMPONENT ANALYSIS OF TOKEN EMBEDDINGS

Let X denote a batch of B sequences of token embeddings, represented as a tensor of shape (B, T, d) , where T is the sequence length and d is the embedding dimension. We apply principal component analysis (PCA) to X and project the tokens onto \mathbb{R}^2 using the top-two principal components, as formally defined in Algorithm 1 in Appendix A. For visualization, tokens belonging to the same class are plotted in the same color.

Figure 1a illustrates the resulting two-dimensional PCA projection of the last-layer token embeddings for ViT-B and ViT-B-1P trained on CIFAR-10. For ViT-B, tokens from different classes overlap and exhibit no clear geometric structure. In contrast, the token embeddings of ViT-B-1P form well-separated clusters. The projections for ViT-B-0P, ViT-B-3P, and ViT-B-Mix-Depth display a pattern similar to that of ViT-B-1P (see Appendix C).

4.3 ANALYSIS OF VARIANCE (ANOVA) OF TOKEN EMBEDDINGS

Following the definitions in Section 3.2, we compute the total variance (TotalVar) of the token embeddings and its three components—between-class variance (BetweenClassVar), within-class variance (WithinClassVar), and within-sequence variance (WithinSeqVar)—each expressed as a fraction of TotalVar so that they sum to one.

Figure 3 presents these fractions for ViT-B and ViT-B-1P trained on CIFAR10. For ViT-B, WithinSeqVar constitutes the majority of the total variance, indicating a limited amount of token-level collapse. In contrast, ViT-B-1P substantially reduces WithinSeqVar while leaving WithinClassVar roughly unchanged, implying a much higher degree of token-level collapse. Moreover, variance energy shifts from WithinSeqVar to BetweenClassVar, suggesting that ViT-B-1P pushes the class means farther apart than ViT-B. These findings reinforce the patterns observed in the PCA visualization.

4.4 LAYERWISE NEURAL TOKEN COLLAPSE

We measure the average cosine similarity between all pairs of tokens within the same sequence, defined as

$$\frac{1}{B} \sum_{b=1}^B \frac{1}{T(T-1)} \sum_{\substack{i,j=1 \\ i \neq j}}^T \frac{\langle X_{b,i}, X_{b,j} \rangle}{\|X_{b,i}\| \|X_{b,j}\|},$$

where $X_{b,i} \in \mathbb{R}^d$ denotes the i -th token in the b -th sequence. This quantity, denoted CosSim, lies between -1 and 1 . A larger value indicates that, on average, tokens within the same sequence are more strongly aligned.

Figure 4 shows CosSim for the output tokens of each transformer block (after the final normalization layer) in models trained on ImageNet. For ViT-B, CosSim remains small across all layers, indicating limited token-level collapse. In contrast, all four proposed models exhibit a much steeper increase in CosSim, reaching substantially higher values in the deeper layers. This behavior suggests that the proposed models induce markedly stronger alignment among tokens than the ViT-B baseline.

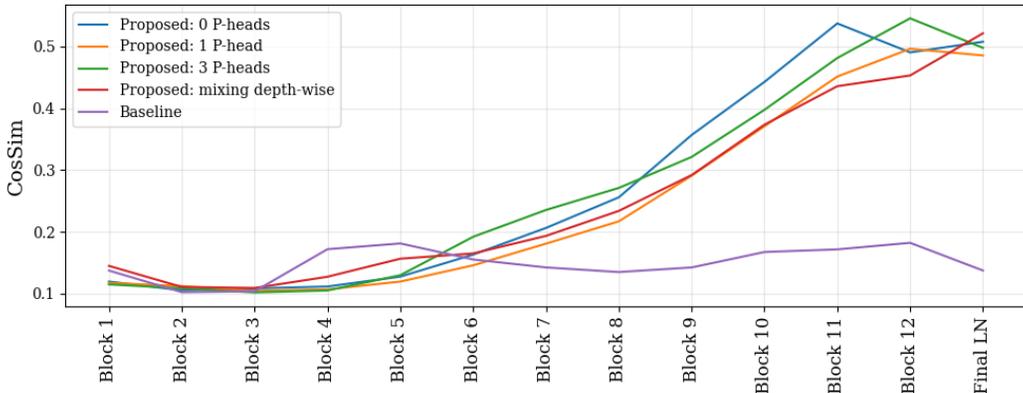


Figure 4: Proposed method induces a significantly steeper increase in CosSim across depth. Measured on ImageNet.

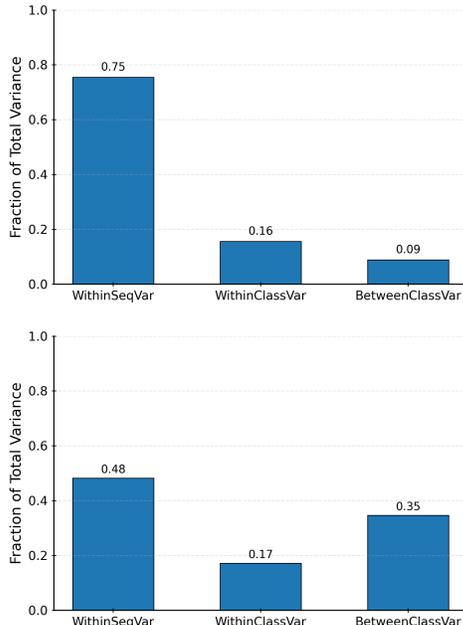


Figure 3: ViT-B-1P (bottom) moves variance from within-image to between-class compared to ViT-B (top). The total variance decomposes additively as TotalVar = BetweenClassVar + WithinClassVar + WithinSeqVar.

4.5 NEURAL COLLAPSE METRICS AND VISUALIZATION

We assess NC2–NC4 using the methodology of Han et al. (2022) and complement these metrics with a visualization technique introduced by Fisher et al. (2024) (Algorithm 2). Full details are provided in Appendix F.

Figure 5 compares the four NC metrics across the different models trained on CIFAR10. Figure 5a shows that the class means are more equinorm for all three proposed models than for the baseline. Figure 5b further indicates that all four proposed models produces class means that are more equiangular. Overall, these results suggest that our modifications better promote a simplex ETF structure in the class means. We observe no significant differences in self-duality or in the NC2 metrics for the classifier weights across models (Figures 5b, 5c). Finally, Figure 5d shows that the proposed models move significantly closer to an NCC classifier.

Figure 5e visualizes the token embeddings projected onto the classifier. While the baseline displays diffuse, overlapping clouds of tokens from different classes, the proposed modification produces well-separated clusters with a clear simplex-like arrangement. This pattern reinforces the NC2–NC4 metrics by illustrating directly that the modified models achieve stronger token-level collapse and more distinct class separation than the ViT-B baseline.

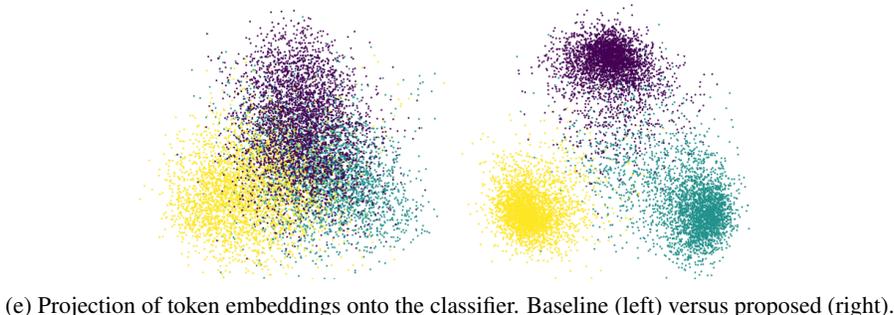
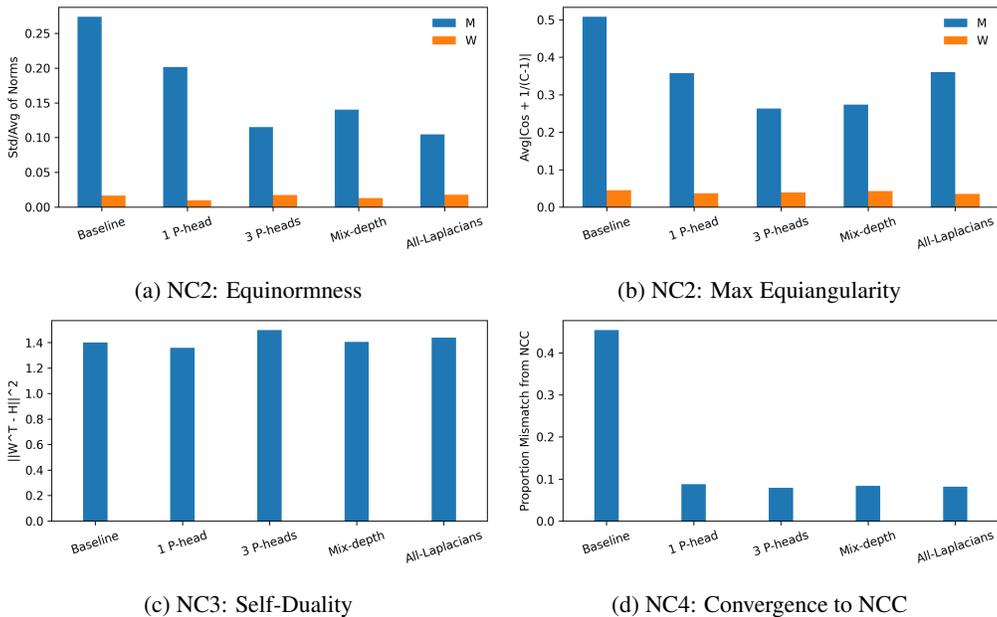


Figure 5: Neural-collapse metrics on CIFAR-10 (top) and projection of token embeddings onto the classifier (bottom). Here M denotes the matrix of class means—averaged over all tokens and instances within each class—while W denotes the classifier weight matrix.

4.6 PROPOSED METHOD IMPROVES PERFORMANCE

4.6.1 IMAGE CLASSIFICATION

Table 1 compares the top-1 test accuracies on multiple datasets between our proposal and the baseline. It clearly illustrates that the proposed models consistently produce meaningful improvements upon the baseline across all datasets. Most notably, the improvements on ImageNet-1k provide evidence that the proposed models work well on large-scale datasets. The significant improvements produced by ViT-B-1P and ViT-B-3P on CIFAR100 (4-5%) also suggest that the proposed models are more data-efficient as CIFAR100 has limited data points per class.

Table 1: Top-1 test accuracy (%) of models across datasets (mean \pm standard deviation). The proposed models yield consistent performance improvement.

| Model | CIFAR-10 | CIFAR-100 | ImageNet-1k |
|----------------------------|----------------------------------|----------------------------------|--------------|
| Baseline (ViT-B) | 90.41 \pm 0.15 | 61.41 \pm 0.36 | 81.2 |
| Proposed (ViT-B-0P) | 91.74 \pm 0.10 | 65.39 \pm 0.07 | 82.02 |
| Proposed (ViT-B-1P) | 91.83\pm0.08 | 66.05\pm0.22 | 82.18 |
| Proposed (ViT-B-3P) | 91.83 \pm 0.10 | 65.44 \pm 0.36 | 82.17 |
| Proposed (ViT-B-Mix-Depth) | 91.79 \pm 0.18 | 61.55 \pm 0.39 | 82.07 |

Next, we compare the top-1 accuracies of models with different numbers m of standard attention heads (which we refer to as P heads). Since the baseline model ViT-B has 12 P heads in each layer, we trained models with $m \in \{0, 1, 3, 6, 9, 12\}$. Figure 6 illustrates that while using the Laplacian heads only (0P) already induces noticeable improvement upon the baseline (12P), incorporating a small number of P heads ($m = 1, 3$) yields further improvements on ImageNet-1k. Since this pattern is consistent across datasets (see Appendix E), we recommend using a small number of P heads as the default strategy.

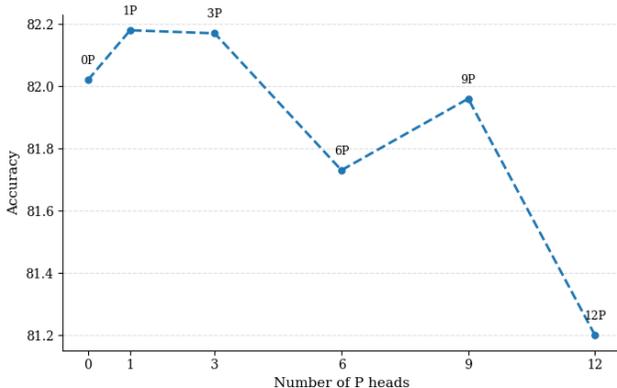


Figure 6: ImageNet accuracy as a function of the number of P heads.

4.6.2 AUTOREGRESSIVE LANGUAGE MODELLING

We trained decoder-only transformer models for autoregressive next-token prediction to investigate the impact of the Laplacian mechanism on language modelling tasks. Specifically, we took a GPT-2 style transformer model with 836 million parameters as the baseline. Following Strategy 1, we assigned the Laplacian mechanism to 15 heads and standard attention to 5 heads in each layer. We trained the two models on 20 billion tokens from the FineWebEdu dataset (Lozhkov et al., 2024) and evaluated them on a variety of zero-shot downstream tasks. Full details of the experiments can be found in Appendix B.

Table 2: Zero-shot results of 836M GPT2-style decoder-only transformers on downstream datasets. “ACC” means accuracy (higher is better).

| | ARC-Easy (ACC) | ARC-Challenge (ACC) | HellaSwag (ACC) | PIQA (ACC) | RACE (ACC) | OpenBookQA (ACC) | WinoGrande (ACC) | SciQ (ACC) |
|----------------------------|-------------------|------------------------|--------------------|---------------|---------------|---------------------|---------------------|---------------|
| Baseline (836M) | 63.8 | 31.48 | 36.22 | 68.55 | 30.62 | 33.2 | 53.99 | 84.1 |
| Proposed: 5 P-heads (836M) | 64.56 | 32.08 | 36.4 | 68.71 | 32.15 | 34.6 | 54.78 | 85.2 |

Our results (Table 2) show that incorporating the Laplacian heads leads to better performance on various downstream datasets. This results corroborates several previous works on the token collapse phenomenon in autoregressive transformers, a connection we discuss in more detail in Section 6.

5 RELATION TO DIFFUSION OVER GRAPHS

Our Laplacian mechanism can be understood through the lens of diffusion on graphs, in particular the discrete heat equation.

Connection to graphs Consider a sequence of T tokens. Treat each token as a vertex in a graph $G = (V, E)$ with $|V| = T$. The attention weights from equation 1,

$$P = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right) \in \mathbb{R}^{T \times T},$$

define directed edge weights: P_{ij} encodes the normalized influence of token j on token i . Because each row sums to one, P is a row-stochastic adjacency matrix.

The graph Laplacian. Given P , the random-walk normalized graph Laplacian is

$$\mathcal{L} = I - P,$$

where $I \in \mathbb{R}^{T \times T}$ is the identity matrix. Since P is row-stochastic, applying \mathcal{L} to a signal $x \in \mathbb{R}^T$ measures its deviation from the local row-normalized average over neighbors.

Diffusion as the discrete heat equation. Diffusion on a graph evolves $x(t)$ according to

$$\frac{d}{dt}x(t) = -\mathcal{L}x(t),$$

whose solution progressively smooths x along the edges. A single explicit Euler step with step size $\Delta t = 1$ gives

$$x_{t+1} = x_t - \mathcal{L}x_t = x_t - (I - P)x_t.$$

Connecting diffusion to our mechanism. This update matches exactly what our Laplacian head together with the residual connection in equation 2 performs:

$$X' = X - \mathcal{L}(V) = X - (I - P)V.$$

The residual connection in equation 2 is written with the opposite sign convention, i.e. $X' = X + \mathcal{L}(V)$. This is purely a notational choice: the sign can always be absorbed into the learnable matrices (for example, replacing W_V by $-W_V$ leaves the mechanism mathematically equivalent).

Consequences for token geometry. Repeated application of the normalized graph Laplacian derived from attention scores drives token representations toward their locally averaged state, reducing within-sequence variance and moving them closer to the collapses described in Subsection 3.

6 RELATED WORKS

Neural Collapse NC is a set of phenomena that describe the training dynamics of deep classifier networks in the terminal phase (Papayan et al., 2020). It reveals the emergence of a distinctive

486 geometry in the final-layer representations. NC has been widely studied in various settings (Han
487 et al., 2022; Zhu et al., 2021; Zhou et al., 2022a;b; Rangamani et al., 2023; Jacot et al., 2024; Wu
488 & Pappan, 2024; Fisher et al., 2024; Yan et al., 2024). Recent work by S ukenik et al. (2025) shows
489 that NC is optimal for transformers and ResNet, and several works focus on inducing NC-related
490 properties in the last-layer features and weights (Markou et al., 2024; Chen et al., 2024). However,
491 these works do not consider the collapse of tokens within the same sequence. Many other works
492 leverage the properties of NC to perform specific tasks (Ammar et al., 2024; Pham et al., 2025).

493
494 **Rank Collapse** Our work is closely related to the widely studied phenomenon of *rank collapse*,
495 where the dimension of token embeddings progressively decreases as they pass through transformer
496 blocks (Dong et al., 2023; Noci et al., 2022; Saada et al., 2025; Geshkovski et al., 2025; Kirsanov
497 et al., 2025; Zhou et al., 2025; Bruno et al., 2025). Rank collapse is commonly regarded as a degen-
498 erate behaviour to be mitigated because it causes training difficulties (Noci et al., 2022) and limits
499 the model’s expressivity (Dong et al., 2023; Barbero et al., 2024). Our work shows that tokens could
500 collapse into a specific geometry that is beneficial for classification, so that collapse is not univer-
501 sally bad. Relatedly, Geshkovski et al. (2024) prove that self-attention collapses tokens into various
502 geometries depending on the initial conditions. However, their work does not investigate how those
503 geometries impact performance. Viswanathan et al. (2025) take one step further to measure the ge-
504 ometric properties of token embeddings in large language models. They show that input prompts
505 with a higher loss produce token embeddings with higher intrinsic dimension. This finding agrees
506 with our hypothesis that NTC0 benefits classification.

507 **Token Collapse in Autoregressive Language Models** A recent work by Zhang et al. (2025) ob-
508 serves a low-rank structure in the token embeddings of transformers trained for next-token predic-
509 tion. More explicitly, Zhao et al. (2025) show that next-token prediction implicitly favors learning
510 logits with a sparse-plus-low-rank structure, where the low-rank component becomes dominant
511 during training and depends only on the support pattern of the context–token co-occurrence matrix.
512 Consequently, when projected onto an appropriate subspace, contexts that share similar next-token
513 supports collapse toward shared low-dimensional directions—a phenomenon they term *subspace*
514 *collapse*. Our results in Section 4.6.2 empirically support their theories: the Laplacian heads could
515 induce subspace collapse more efficiently, leading to improved downstream performance.

516 **Oversmoothing in Graph Neural Networks** Since transformers can be viewed as graph neural
517 networks (GNNs) (Joshi, 2025), our work is related to the oversmoothing phenomenon in GNNs
518 (Roth & Liebig, 2024; Rusch et al., 2023;  lvvaro Arroyo et al., 2025). Like rank collapse, over-
519 smoothing is regarded as an issue to be mitigated (Roth et al., 2024; Rusch et al., 2022; Nguyen
520 et al., 2023; Bodnar et al., 2023; Wang et al., 2022). In particular, Rusch et al. (2022) addresses
521 oversmoothing by modifying the underlying GNN dynamics. From this perspective, our modifica-
522 tion can be understood as discrete heat diffusion on graphs, a connection that we discuss in Section 5.

524 7 CONCLUSION

525
526 Neural Token Collapse represents an ideal state where token embeddings collapse to a simplex-
527 ETF geometry that maximizes class separability. Standard transformers diverge from this geometry
528 because attention moves tokens by inflating their means rather than directly along their variance
529 directions. To address this, we proposed a simple Laplacian-based modification of attention that
530 enables tokens to move along their variance directions, promoting a more efficient convergence
531 toward the desired structure. The Laplacian head can be interpreted as performing diffusion—via a
532 heat equation—of token embeddings over the graph implicitly learned by the attention mechanism.
533 Experiments across datasets show that this modification moves token embeddings markedly closer
534 to NTC and yields consistent gains in performance. A promising direction for future work is to
535 explore how these mechanisms can be modified further to fully realize the envisioned NTC ideal.

536 It is important to note that while NTC is ideal for classification, the optimal representation geometry
537 for other tasks might be different. By allowing tokens to move more freely along their variance
538 directions, the Laplacian mechanism can, in principle, encourage tokens to deviate from their means.
539 Therefore, another important direction for future work is to characterize the optimal geometries for
other tasks and study whether the Laplacian can steer representations towards those ideals.

REFERENCES

- 540
541
542 Mouin Ben Ammar, Nacim Belkhir, Sebastian Popescu, Antoine Manzanera, and Gianni Franchi.
543 Neco: Neural collapse based out-of-distribution detection, 2024. URL <https://arxiv.org/abs/2310.06823>.
544
- 545 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL
546 <https://arxiv.org/abs/1607.06450>.
547
- 548 Federico Barbero, Andrea Banino, Steven Kapturowski, Dharshan Kumaran, João G. M. Araújo,
549 Alex Vitvitskiy, Razvan Pascanu, and Petar Veličković. Transformers need glasses! information
550 over-squashing in language tasks, 2024. URL <https://arxiv.org/abs/2406.04267>.
- 551 Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about
552 physical commonsense in natural language, 2019. URL <https://arxiv.org/abs/1911.11641>.
553
554
- 555 Cristian Bodnar, Francesco Di Giovanni, Benjamin Paul Chamberlain, Pietro Liò, and Michael M.
556 Bronstein. Neural sheaf diffusion: A topological perspective on heterophily and oversmoothing
557 in gnns, 2023. URL <https://arxiv.org/abs/2202.04579>.
- 558 Giuseppe Bruno, Federico Pasqualotto, and Andrea Agazzi. Emergence of meta-stable clustering in
559 mean-field transformer models, 2025. URL <https://arxiv.org/abs/2410.23228>.
560
- 561 Zhikang Chen, Min Zhang, Sen Cui, Haoxuan Li, Gang Niu, Mingming Gong, Changshui Zhang,
562 and Kun Zhang. Neural collapse inspired feature alignment for out-of-distribution generalization.
563 In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL
564 <https://openreview.net/forum?id=wQpNG9JnPK>.
- 565 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
566 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge,
567 2018. URL <https://arxiv.org/abs/1803.05457>.
568
- 569 Ekin D. Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V. Le. Randaugment: Practical automated
570 data augmentation with a reduced search space, 2019. URL <https://arxiv.org/abs/1909.13719>.
571
572
- 573 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-
574 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,
575 pp. 248–255. IEEE, 2009.
- 576 Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure
577 attention loses rank doubly exponentially with depth, 2023. URL <https://arxiv.org/abs/2103.03404>.
578
579
- 580 Quinn Fisher, Haoming Meng, and Vardan Papyan. Pushing boundaries: Mixup’s influence on
581 neural collapse, 2024. URL <https://arxiv.org/abs/2402.06171>.
- 582 Borjan Geshkovski, Cyril Letrouit, Yury Polyanskiy, and Philippe Rigollet. The emergence of clus-
583 ters in self-attention dynamics, 2024. URL <https://arxiv.org/abs/2305.05465>.
584
- 585 Borjan Geshkovski, Cyril Letrouit, Yury Polyanskiy, and Philippe Rigollet. A mathematical per-
586 spective on transformers, 2025. URL <https://arxiv.org/abs/2312.10794>.
- 587 X. Y. Han, Vardan Papyan, and David L. Donoho. Neural collapse under mse loss: Proximity to and
588 dynamics on the central path, 2022. URL <https://arxiv.org/abs/2106.02073>.
589
- 590 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
591 nition, 2015. URL <https://arxiv.org/abs/1512.03385>.
592
593
- Alex Henry, Prudhvi Raj Dachapally, Shubham Pawar, and Yuxuan Chen. Query-key normalization
for transformers, 2020. URL <https://arxiv.org/abs/2010.04245>.

- 594 Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang,
595 Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zheng Leng Thai, Kaihuo Zhang, Chongyi Wang,
596 Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang
597 Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm: Unveiling the potential of small
598 language models with scalable training strategies, 2024. URL [https://arxiv.org/abs/
599 2404.06395](https://arxiv.org/abs/2404.06395).
- 600 Arthur Jacot, Peter Sukenk, Zihan Wang, and Marco Mondelli. Wide neural networks trained
601 with weight decay provably exhibit neural collapse, 2024. URL [https://arxiv.org/abs/
602 2410.04887](https://arxiv.org/abs/2410.04887).
- 603
604 Matt Gardner Johannes Welbl, Nelson F. Liu. Crowdsourcing multiple choice science questions.
605 2017.
- 606
607 Keller Jordan, Jeremy Bernstein, Brendan Rappazzo, @fernbear.bsky.social, Boza Vlado,
608 You Jiacheng, Franz Cesista, Braden Koszarsky, and @Grad62304977. modded-nanogpt:
609 Speedrunning the nanogpt baseline, 2024. URL [https://github.com/KellerJordan/
610 modded-nanogpt](https://github.com/KellerJordan/modded-nanogpt).
- 611 Chaitanya K. Joshi. Transformers are graph neural networks, 2025. URL [https://arxiv.org/
612 abs/2506.22084](https://arxiv.org/abs/2506.22084).
- 613
614 Andrej Karpathy. nanoGPT: The simplest, fastest repository for training/finetuning medium-sized
615 gpts. <https://github.com/karpathy/nanoGPT>, 2022. GitHub repository, accessed:
616 YYYY-MM-DD.
- 617
618 Andrej Karpathy. nanochat: The best chatgpt that \$100 can buy. [https://github.com/
619 karpathy/nanochat](https://github.com/karpathy/nanochat), 2025.
- 620 Artem Kirsanov, Chi-Ning Chou, Kyunghyun Cho, and SueYeon Chung. The geometry of prompt-
621 ing: Unveiling distinct mechanisms of task adaptation in language models, 2025. URL [https :
622 //arxiv.org/abs/2502.08009](https://arxiv.org/abs/2502.08009).
- 623
624 Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University
625 of Toronto, 2009. Describes the CIFAR-10 dataset.
- 626
627 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading
628 comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017.
- 629 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL [https :
630 //arxiv.org/abs/1711.05101](https://arxiv.org/abs/1711.05101).
- 631
632 Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. Fineweb-edu: the finest
633 collection of educational content, 2024. URL [https://huggingface.co/datasets/
634 HuggingFaceFW/fineweb-edu](https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu).
- 635
636 Evan Markou, Thalaiyasingam Ajanthan, and Stephen Gould. Guiding neural collapse: Optimising
637 towards the nearest simplex equiangular tight frame. In *The Thirty-eighth Annual Conference on
638 Neural Information Processing Systems*, 2024. URL [https://openreview.net/forum?
639 id=z4FaPUslma](https://openreview.net/forum?id=z4FaPUslma).
- 640
641 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct
642 electricity? a new dataset for open book question answering, 2018. URL [https://arxiv .
643 org/abs/1809.02789](https://arxiv.org/abs/1809.02789).
- 644
645 Tam Nguyen, Tan M. Nguyen, and Richard G. Baraniuk. Mitigating over-smoothing in transformers
646 via regularized nonlocal functionals, 2023. URL <https://arxiv.org/abs/2312.00751>.
- 647
648 Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurelien
649 Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse,
650 2022. URL <https://arxiv.org/abs/2206.03126>.

- 648 Nvidia, :, Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H. Anh, Pallab Bhattacharya, Annika
649 Brundyn, Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, Sirshak Das, Ayush
650 Dattagupta, Olivier Delalleau, Leon Derczynski, Yi Dong, Daniel Egert, Ellie Evans, Alek-
651 sander Ficek, Denys Fridman, Shaona Ghosh, Boris Ginsburg, Igor Gitman, Tomasz Grze-
652 gorzek, Robert Hero, Jining Huang, Vibhu Jawa, Joseph Jennings, Aastha Jhunjhunwala, John
653 Kamalu, Sadaf Khan, Oleksii Kuchaiev, Patrick LeGresley, Hui Li, Jiwei Liu, Zihan Liu, Eileen
654 Long, Ameya Sunil Mahabaleshwarkar, Somshubra Majumdar, James Maki, Miguel Martinez,
655 Maer Rodrigues de Melo, Ivan Moshkov, Deepak Narayanan, Sean Narenthiran, Jesus Navarro,
656 Phong Nguyen, Osvald Nitski, Vahid Noroozi, Guruprasad Nutheti, Christopher Parisien, Jupinder
657 Parmar, Mostofa Patwary, Krzysztof Pawelec, Wei Ping, Shrimai Prabhumoye, Rajarshi Roy,
658 Trisha Saar, Vasanth Rao Naik Sabavat, Sanjeev Satheesh, Jane Polak Scowcroft, Jason Se-
659 wall, Pavel Shamis, Gerald Shen, Mohammad Shoeybi, Dave Sizer, Misha Smelyanskiy, Felipe
660 Soares, Makesh Narsimhan Sreedhar, Dan Su, Sandeep Subramanian, Shengyang Sun, Shub-
661 ham Toshniwal, Hao Wang, Zhilin Wang, Jiaxuan You, Jiaqi Zeng, Jimmy Zhang, Jing Zhang,
662 Vivienne Zhang, Yian Zhang, and Chen Zhu. NemoTron-4 340b technical report, 2024. URL
663 <https://arxiv.org/abs/2406.11704>.
- 664 Vardan Papyan, X. Y. Han, and David L. Donoho. Prevalence of neural collapse during the terminal
665 phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):
666 24652–24663, 2020. doi: 10.1073/pnas.2015509117. URL [https://www.pnas.org/doi/](https://www.pnas.org/doi/abs/10.1073/pnas.2015509117)
667 [abs/10.1073/pnas.2015509117](https://www.pnas.org/doi/abs/10.1073/pnas.2015509117).
- 668 Duc Thanh Pham, Hong Dang Nguyen, Nhat Minh Nguyen Quoc, Linh Ngo Van, Sang Dinh Viet,
669 and Duc Anh Nguyen. Hierarchical neural collapse detection transformer for class incremental
670 object detection, 2025. URL <https://arxiv.org/abs/2506.08562>.
- 671 Akshay Rangamani, Marius Lindegaard, Tomer Galanti, and Tomaso A Poggio. Feature learn-
672 ing in deep classifiers through intermediate neural collapse. In Andreas Krause, Emma Brun-
673 skill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Pro-*
674 *ceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proce-*
675 *edings of Machine Learning Research*, pp. 28729–28745. PMLR, 23–29 Jul 2023. URL [https://](https://proceedings.mlr.press/v202/rangamani23a.html)
676 proceedings.mlr.press/v202/rangamani23a.html.
- 677 Andreas Roth and Thomas Liebig. Rank collapse causes over-smoothing and over-correlation in
678 graph neural networks, 2024. URL <https://arxiv.org/abs/2308.16800>.
- 680 Andreas Roth, Franka Bause, Nils M. Kriege, and Thomas Liebig. Preventing representational rank
681 collapse in mpnns by splitting the computational graph, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2409.11504)
682 [abs/2409.11504](https://arxiv.org/abs/2409.11504).
- 683 T. Konstantin Rusch, Benjamin P. Chamberlain, James Rowbottom, Siddhartha Mishra, and
684 Michael M. Bronstein. Graph-coupled oscillator networks, 2022. URL [https://arxiv.](https://arxiv.org/abs/2202.02296)
685 [org/abs/2202.02296](https://arxiv.org/abs/2202.02296).
- 687 T. Konstantin Rusch, Michael M. Bronstein, and Siddhartha Mishra. A survey on oversmoothing in
688 graph neural networks, 2023. URL <https://arxiv.org/abs/2303.10993>.
- 689 Thiziri Nait Saada, Alireza Naderi, and Jared Tanner. Mind the gap: a spectral analysis of rank
690 collapse and signal propagation in attention layers, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2410.07799)
691 [2410.07799](https://arxiv.org/abs/2410.07799).
- 692 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-
693 sarial winograd schema challenge at scale, 2019. URL [https://arxiv.org/abs/1907.](https://arxiv.org/abs/1907.10641)
694 [10641](https://arxiv.org/abs/1907.10641).
- 696 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: En-
697 hanced transformer with rotary position embedding, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2104.09864)
698 [2104.09864](https://arxiv.org/abs/2104.09864).
- 700 Peter Súkeník, Christoph H. Lampert, and Marco Mondelli. Neural collapse is globally optimal in
701 deep regularized resnets and transformers, 2025. URL [https://arxiv.org/abs/2505.](https://arxiv.org/abs/2505.15239)
[15239](https://arxiv.org/abs/2505.15239).

- 702 Hugo Touvron, Matthieu Cord, and Herve Jegou. Deit iii: Revenge of the vit. *arXiv preprint*
703 *arXiv:2204.07118*, 2022.
704
- 705 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
706 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
707 *tion processing systems*, 30, 2017.
708
- 709 Karthik Viswanathan, Yuri Gardinazzi, Giada Panerai, Alberto Cazzaniga, and Matteo Biagetti. The
710 geometry of tokens in internal representations of large language models, 2025. URL <https://arxiv.org/abs/2501.10573>.
711
- 712 Peihao Wang, Wenqing Zheng, Tianlong Chen, and Zhangyang Wang. Anti-oversmoothing in deep
713 vision transformers via the fourier domain analysis: From theory to practice, 2022. URL <https://arxiv.org/abs/2203.05962>.
714
- 715 Robert Wu and Vardan Papyan. Linguistic collapse: Neural collapse in (large) language models,
716 2024. URL <https://arxiv.org/abs/2405.17767>.
717
- 718 Hongren Yan, Yuhua Qian, Furong Peng, Jiachen Luo, Zheqing Zhu, and Feijiang Li. Neural col-
719 lapse to multiple centers for imbalanced data. In *Advances in Neural Information Processing*
720 *Systems*, 2024. URL <https://openreview.net/forum?id=unknown>. NeurIPS 2024
721 poster.
722
- 723 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
724 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
725 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
726 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
727 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
728 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
729 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
730 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
731 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- 732 Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan
733 Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learn-
734 ing: Training bert in 76 minutes, 2020. URL <https://arxiv.org/abs/1904.00962>.
- 735 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-
736 chine really finish your sentence?, 2019. URL <https://arxiv.org/abs/1905.07830>.
737
- 738 Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empiri-
739 cal risk minimization, 2018. URL <https://arxiv.org/abs/1710.09412>.
740
- 741 Stephen Zhang, Mustafa Khan, and Vardan Papyan. Attention sinks: A ‘catch, tag, release’ mecha-
742 nism for embeddings, 2025. URL <https://arxiv.org/abs/2502.00919>.
- 743 Yize Zhao, Tina Behnia, Vala Vakilian, and Christos Thrampoulidis. Implicit geometry of next-
744 token prediction: From language sparsity patterns to model representations, 2025. URL <https://arxiv.org/abs/2408.15417>.
745
- 746 Jinxin Zhou, Xiao Li, Tianyu Ding, Chong You, Qing Qu, and Zhihui Zhu. On the optimization land-
747 scape of neural collapse under mse loss: Global optimality with unconstrained features, 2022a.
748 URL <https://arxiv.org/abs/2203.01238>.
749
- 750 Jinxin Zhou, Chong You, Xiao Li, Kangning Liu, Sheng Liu, Qing Qu, and Zhihui Zhu. Are all
751 losses created equal: A neural collapse perspective, 2022b. URL <https://arxiv.org/abs/2210.02192>.
752
- 753 Yuqi Zhou, Sunhao Dai, Zhanshuo Cao, Xiao Zhang, and Jun Xu. Length-induced embedding
754 collapse in transformer-based models, 2025. URL [https://openreview.net/forum?](https://openreview.net/forum?id=jgISClwdYy)
755 [id=jgISClwdYy](https://openreview.net/forum?id=jgISClwdYy).

756 Zihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu.
757 A geometric analysis of neural collapse with unconstrained features, 2021. URL <https://arxiv.org/abs/2105.02375>.
758
759
760 Álvaro Arroyo, Alessio Gravina, Benjamin Gutteridge, Federico Barbero, Claudio Gallicchio, Xi-
761 aowen Dong, Michael Bronstein, and Pierre Vandergheynst. On vanishing gradients, over-
762 smoothing, and over-squashing in gnns: Bridging recurrent and graph learning, 2025. URL
763 <https://arxiv.org/abs/2502.10818>.
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A PCA VISUALIZATION METHOD

Algorithm 1 PCA Projection of Tokens to \mathbb{R}^2

Require: $X \in \mathbb{R}^{B \times T \times d}$

- 1: $X \leftarrow \text{reshape}(X, [B \cdot T, d])$
- 2: $\mu \leftarrow \text{mean}(X, \text{axis} = 0)$
- 3: $U, \Sigma, V^\top \leftarrow \text{SVD}(X - \mu)$
- 4: $V_2 \leftarrow V[:, 0 : 2]$
- 5: **return** $(X - \mu) \cdot V_2$

B EXPERIMENT DETAILS

B.1 IMAGE CLASSIFICATION

To produce results in Table 1, we trained the ViT-B model from Touvron et al. (2022). The model consists of 12 blocks, with 12 attention heads in each block, and an embedding dimension of 768. The number of trainable parameters is around 86.6 million. For each dataset, we sweep the peak learning rate over the set $\{4e - 5, 3e - 4, 5e - 4, 3e - 3, 4e - 3\}$ and weight decay over the set $\{0.01, 0.02, 0.05\}$ and select whichever combination that works the best. Many other hyperparameters (such as the mixup (Zhang et al., 2018) α) were selected following the training recipe detailed in Touvron et al. (2022). Depending on the dataset, we use RandAugment (Cubuk et al., 2019) or 3-Aug Touvron et al. (2022) for data augmentation and AdamW (Loshchilov & Hutter, 2019) or LAMB (You et al., 2020) as the optimizer. All models in Table 1 were trained for 300 epochs using 3 random seeds. Full details of our training setup, including the hyperparameters that were eventually selected for the experiments, are provided in Table 3.

| | CIFAR-10 | CIFAR-100 | ImageNet |
|------------------------|------------------|------------------|----------------------|
| Loss | Cross Entropy | Cross Entropy | Binary Cross Entropy |
| Optimizer | AdamW | AdamW | LAMB |
| AdamW β_1 | 0.9 | 0.9 | 0.9 |
| AdamW β_2 | 0.99 | 0.99 | 0.999 |
| Starting Learning Rate | 3e-6 | 3e-6 | 1e-3 |
| Peak Learning Rate | 3e-4 | 3e-4 | 3e-3 |
| Minimum Learning Rate | 0 | 0 | 1e-6 |
| Weight Decay | 0.02 | 0.02 | 0.02 |
| Drop Path Rate | 0.1 | 0.1 | 0.3 |
| Batch Size | 512 | 512 | 2048 |
| Gradient Clipping | 1.0 | 1.0 | 1.0 |
| LR Scheduler | Cosine Annealing | Cosine Annealing | Cosine Annealing |
| Warmup Epochs | 5 | 5 | 5 |
| Data Augmentation | RandAugment | RandAugment | 3-Aug |
| Mixup α | 0.8 | 0.8 | 0.8 |
| Mixup Probability | 1.0 | 1.0 | 1.0 |
| Input Size | 32×32 | 32×32 | 224×224 |
| Patch Size | 4×4 | 4×4 | 16×16 |
| Precision | float32 | float32 | bfloat16 |

Table 3: Training setup for CIFAR-10, CIFAR-100, and ImageNet with hyper-parameter selection informed by Touvron et al. (2022).

B.2 AUTOREGRESSIVE LANGUAGE MODELLING

The model architecture we used was based on (Yang et al., 2025; Karpathy, 2025; 2022; Jordan et al., 2024). The model is a decoder-only transformer similar to GPT-2, with the following architectural modification:

- Query-key normalization (Henry et al., 2020; Yang et al., 2025).
- Rotary positional encodings (Su et al., 2023).
- Squared ReLU activation (Nvidia et al., 2024).
- Untied weights for the first-layer embeddings and the last linear layer.

We chose this architecture as the baseline since we found it to have more stable training and converge to a lower validation loss with our limited compute resources. We trained two transformer models with 836M parameters. Each model was trained on roughly 20B FineWeb-Edu tokens. We followed a WSD learning rate schedule (Hu et al., 2024): the first 5% of the total training duration was used for linear learning rate warmup, and the last 20% of the training duration was used for linear decay. We used the same training setup for both models (Table 4).

| | |
|------------------------|----------|
| Optimizer | AdamW |
| Peak Learning Rate | 2.5e-4 |
| Minimum Learning Rate | 0.0 |
| Learning Rate Schedule | WSD |
| Batch Size | 0.49M |
| Precision | bfloat16 |

Table 4: Training setup for language modelling on FineWeb-Edu.

After training, we evaluated the models on the following zero-shot tasks: HellaSwag(Zellers et al., 2019), PIQA(Bisk et al., 2019), ARC-Easy(Clark et al., 2018), ARC-Challenge(Clark et al., 2018), OpenBookQA(Mihaylov et al., 2018), WinoGrande(Sakaguchi et al., 2019), RACE(Lai et al., 2017), and SciQ(Johannes Welbl, 2017).

C ADDITIONAL EXPERIMENT RESULTS

This section contains additional measurements and visualizations for the models in Section 4.1.

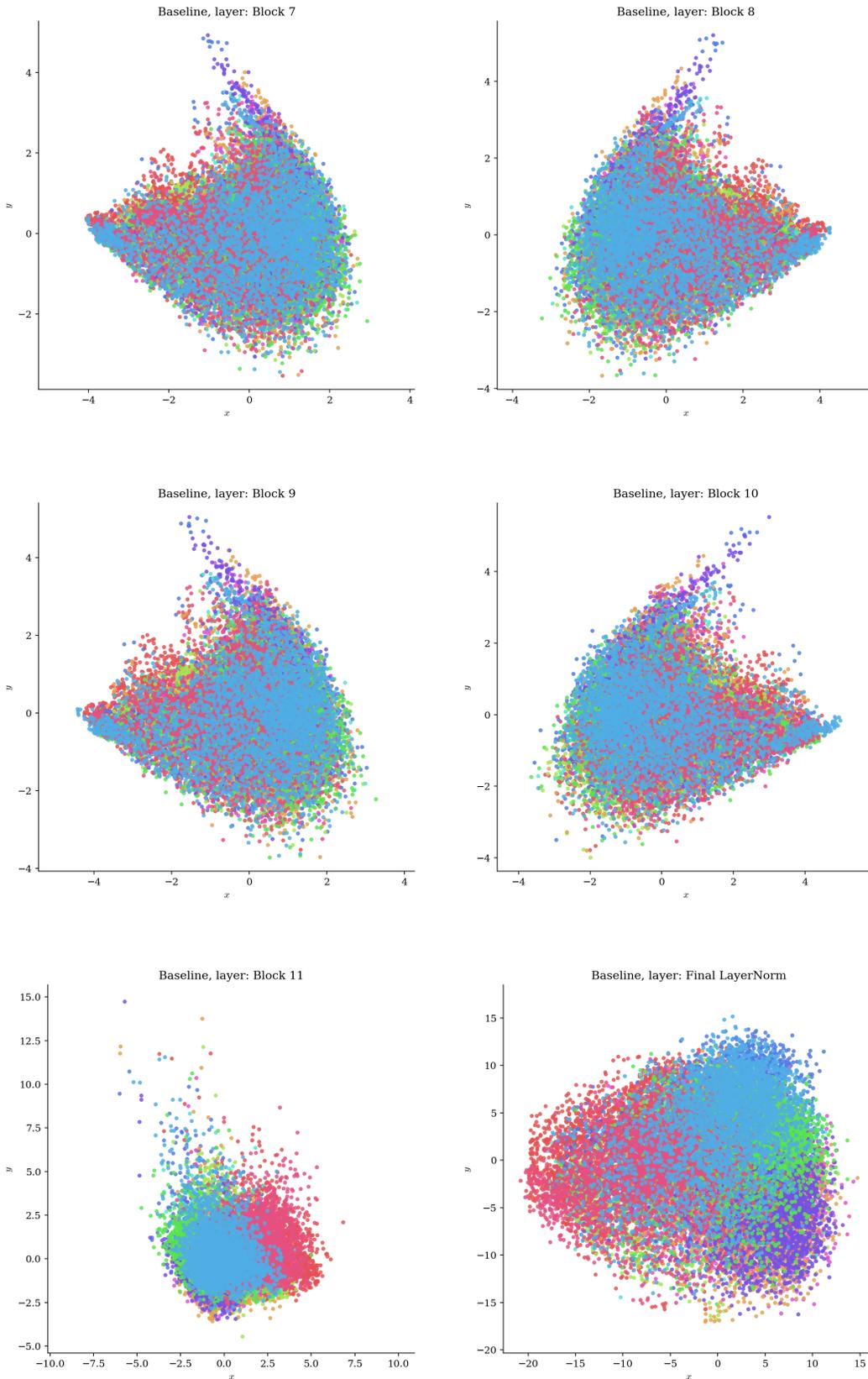
C.1 PCA VISUALIZATIONS

For CIFAR100 and ImageNet-1k, we first randomly sample 10 and 3 classes, respectively. Then we sample a batch of 512 images that contain roughly an equal number of images across the selected classes. All PCA visualizations in this section are performed using the output of some layer (usually a block from the later half, or the final normalization layer), as indicated in the plot.

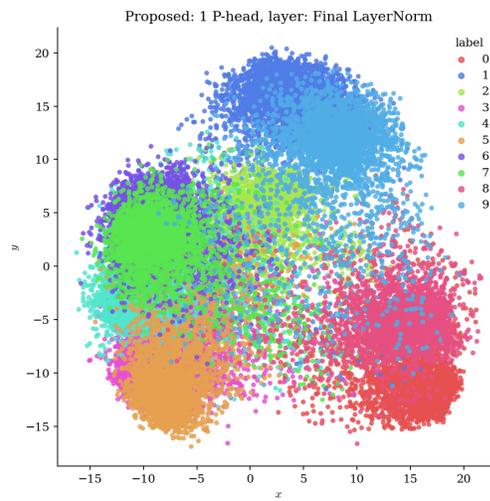
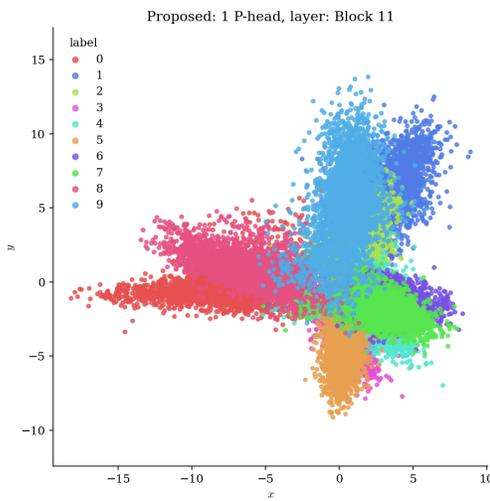
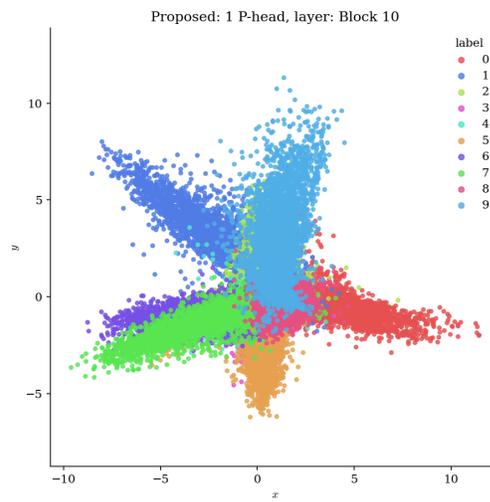
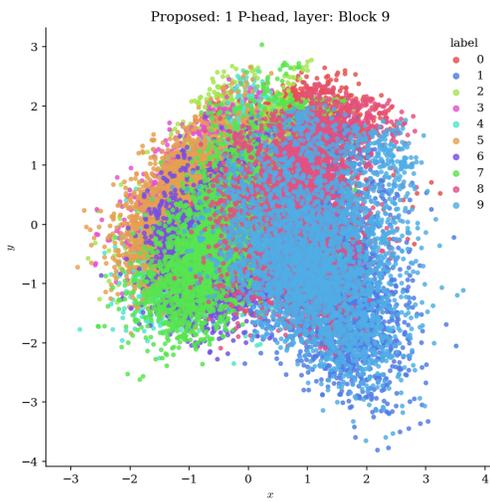
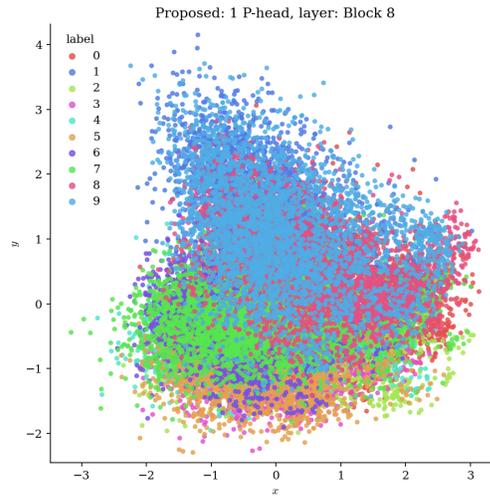
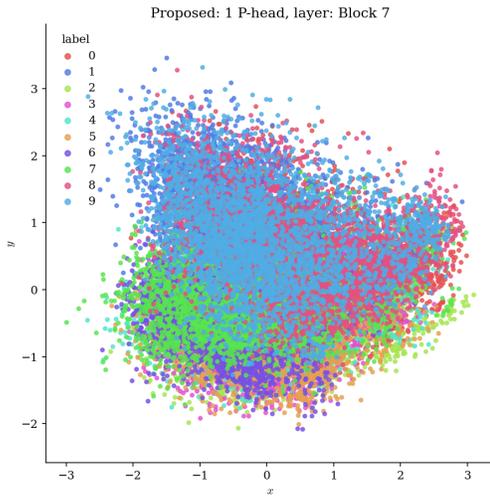
C.1.1 CIFAR10

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

ViT-B (Baseline):



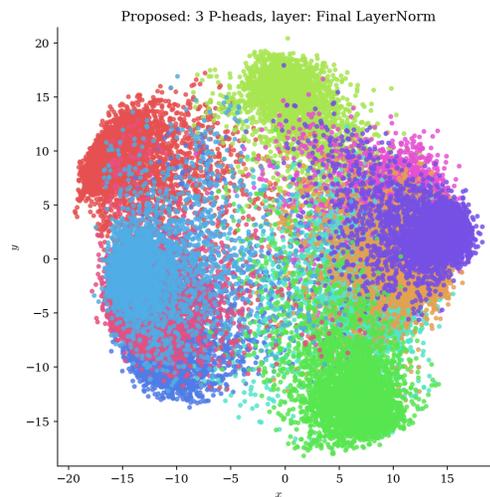
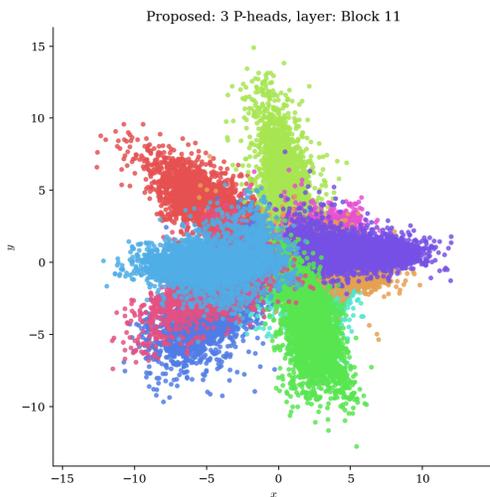
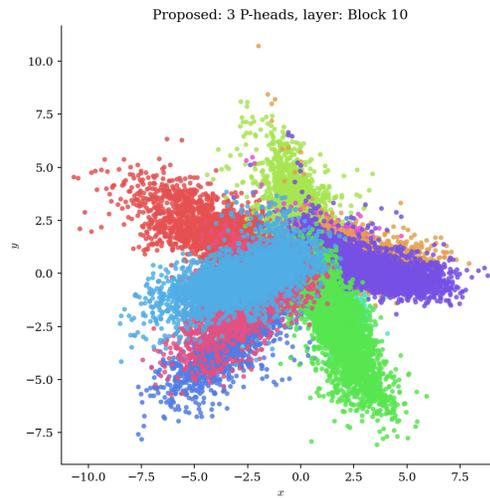
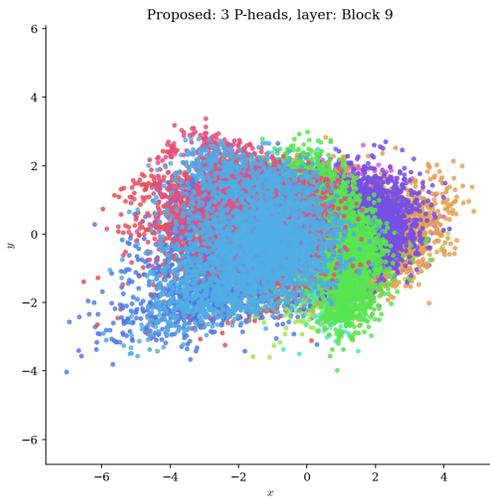
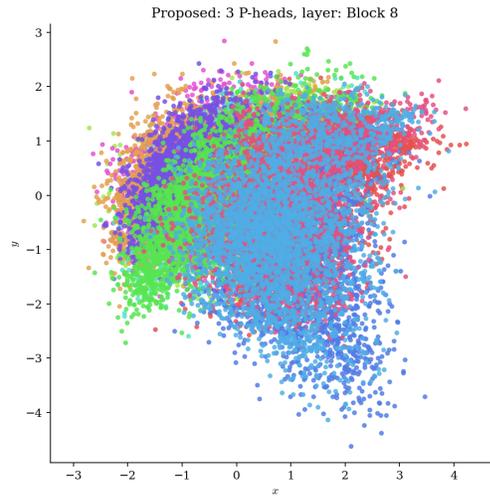
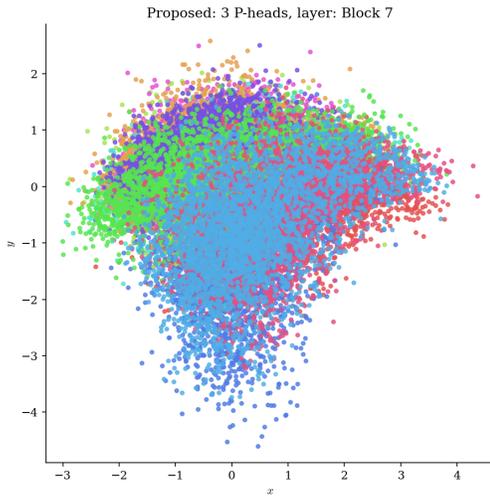
ViT-1P (Proposed):



972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

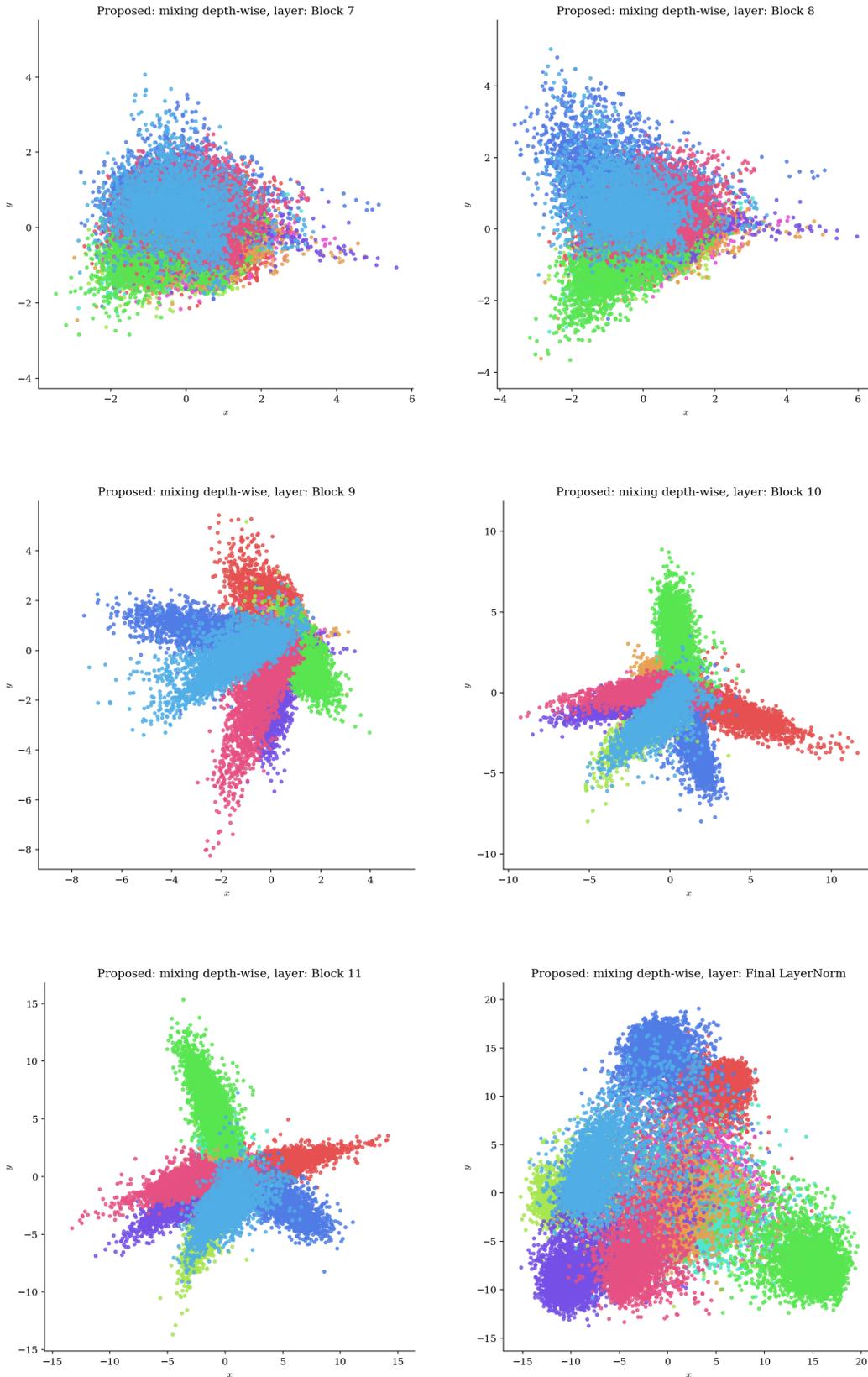
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

ViT-3P (Proposed):



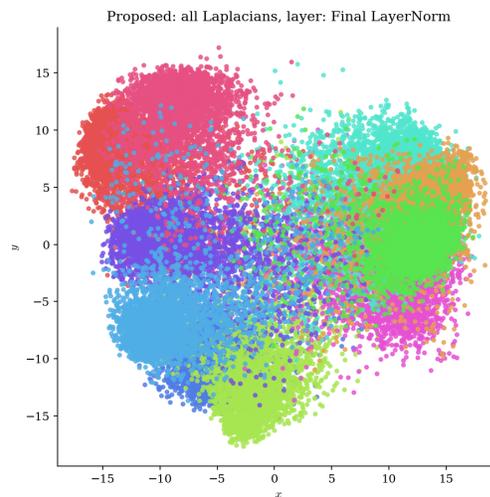
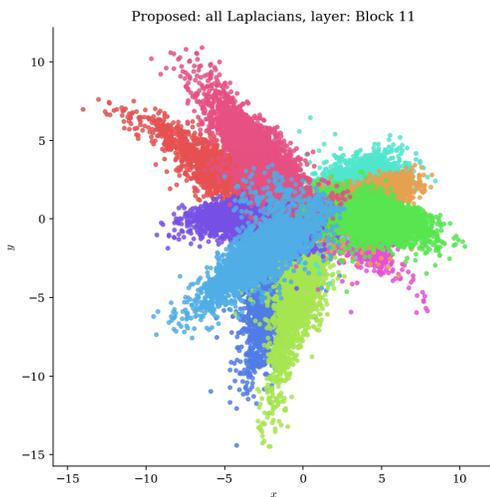
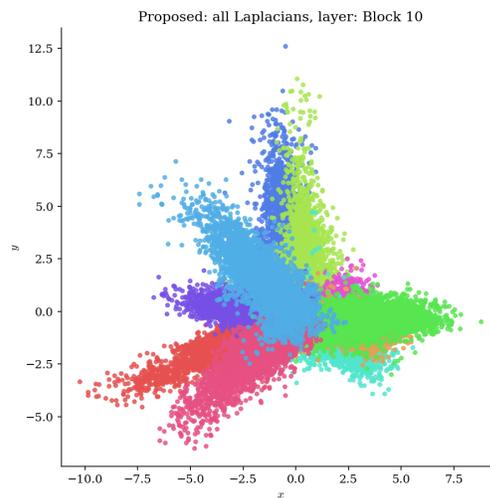
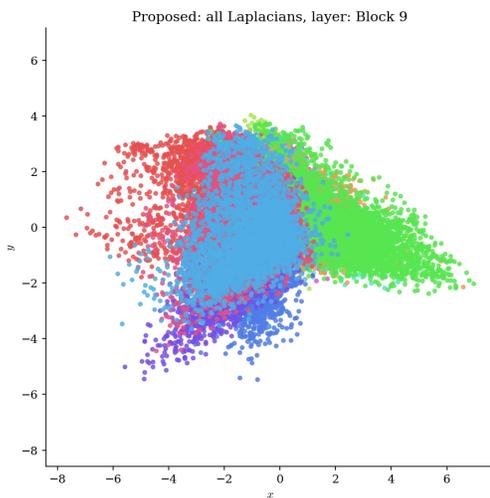
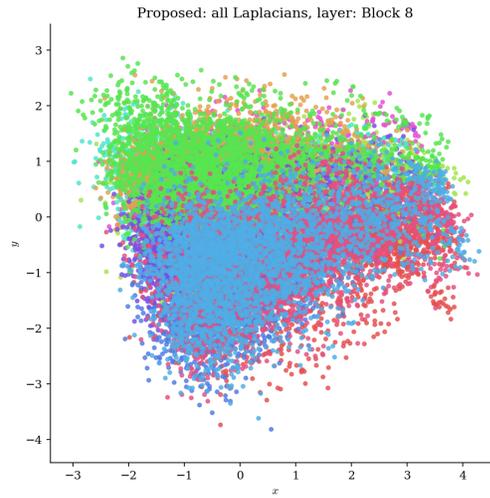
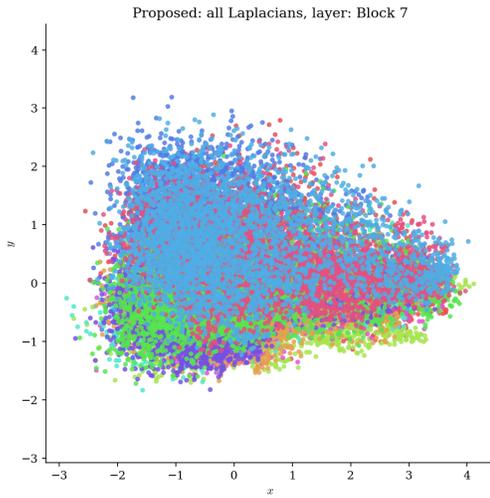
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

ViT-mix-depth (Proposed):



1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

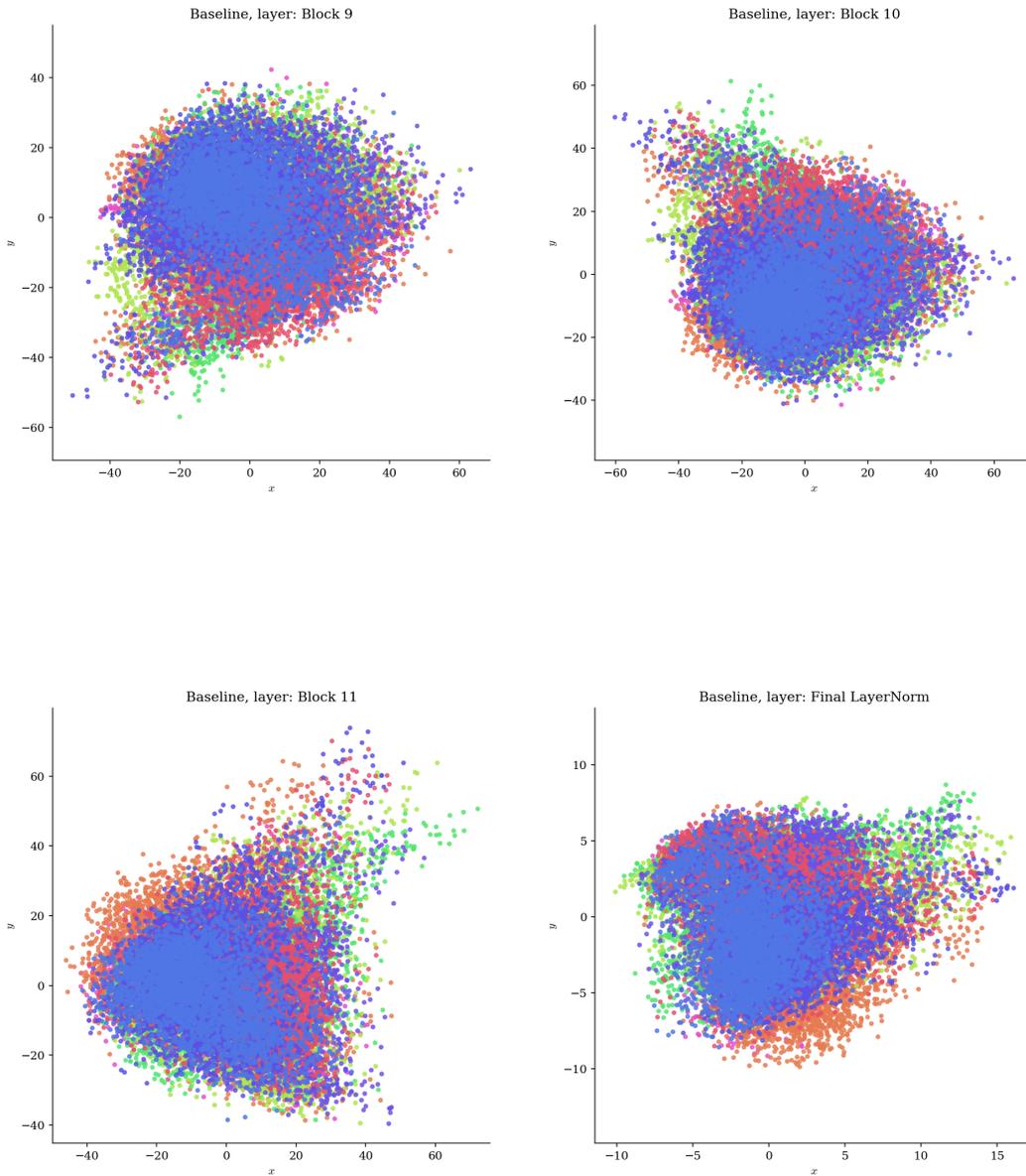
ViT-0P (Proposed):



1188 Notice the emergence of the distinct geometry in the token embeddings for the proposed models
1189 (block 9 to block 10 for ViT-0P, ViT-1P, and ViT-3P, and block 8 to block 9 for ViT-mix-depth). We
1190 do not observe this phenomenon in the baseline model.
1191
1192
1193

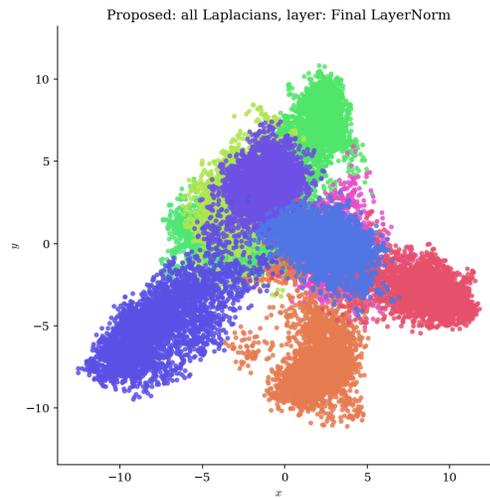
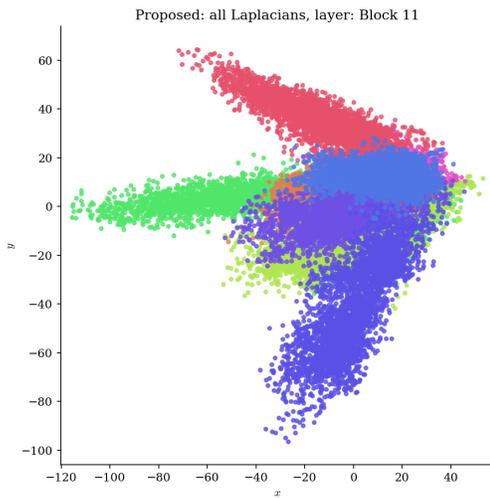
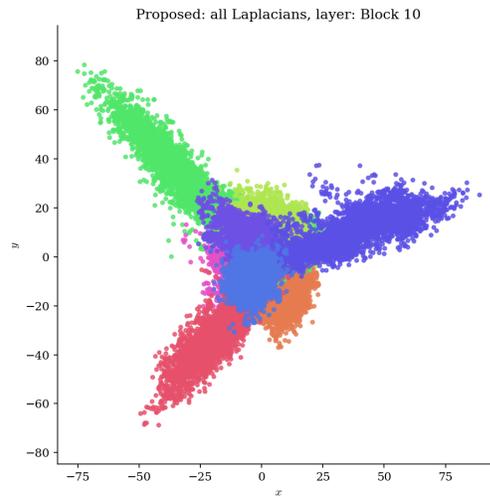
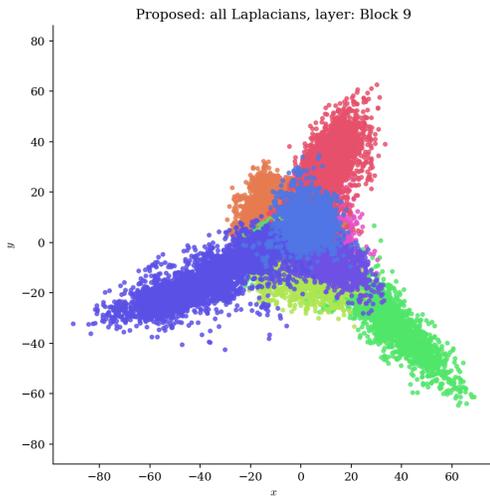
1194 C.1.2 CIFAR100

1195 ViT-B (Baseline):

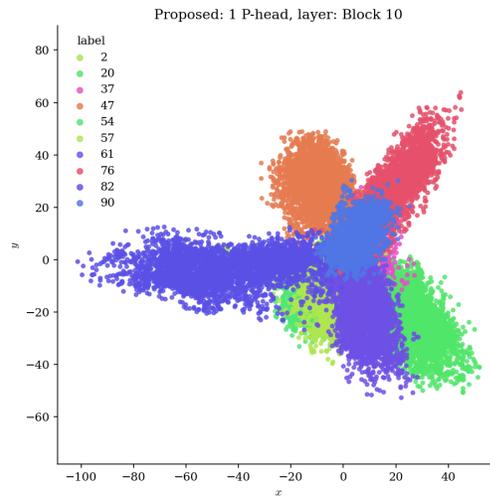
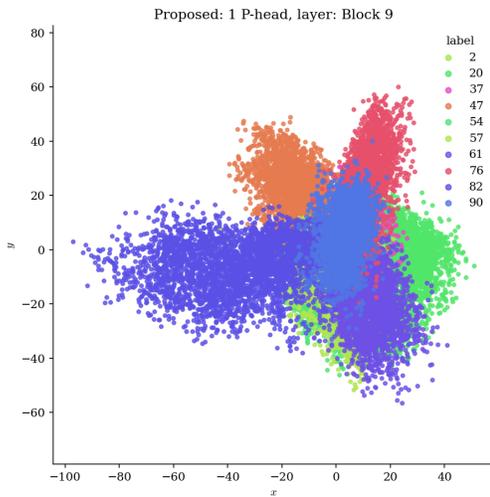


ViT-0P (Proposed):

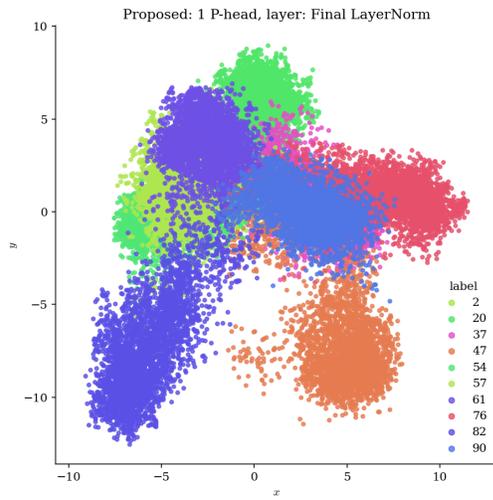
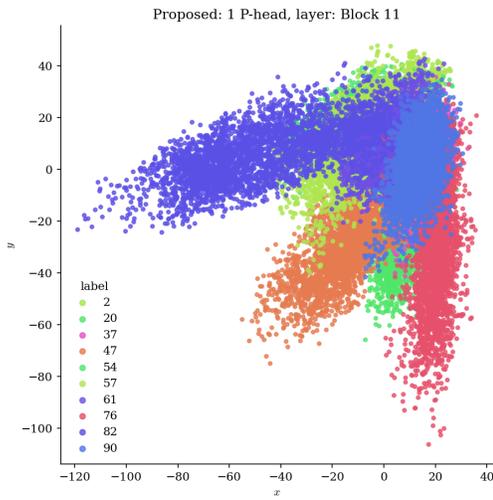
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295



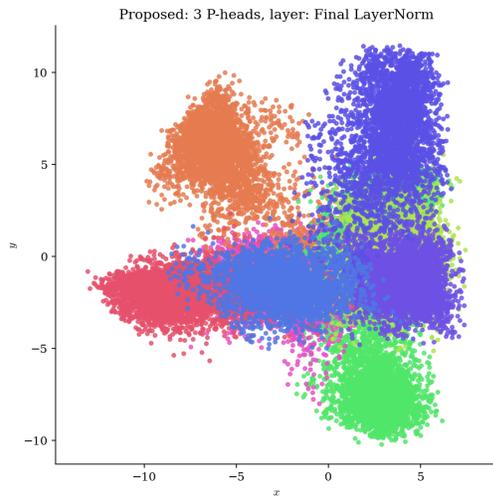
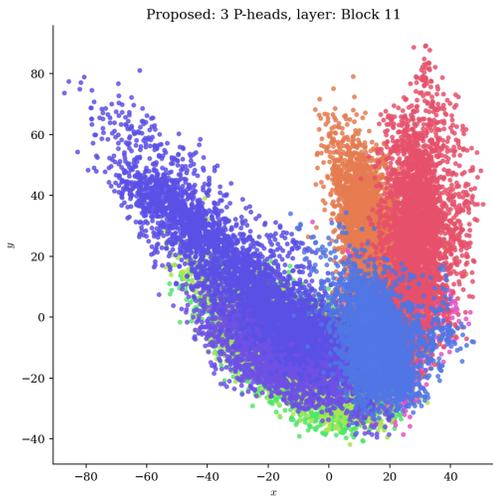
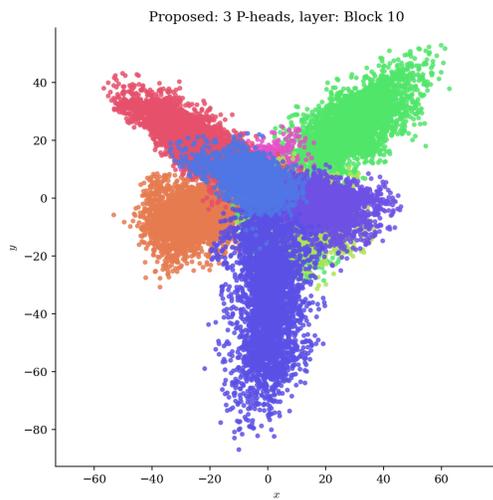
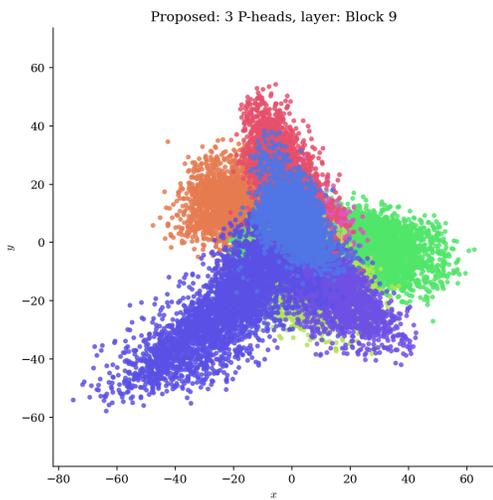
ViT-1P (Proposed):



1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

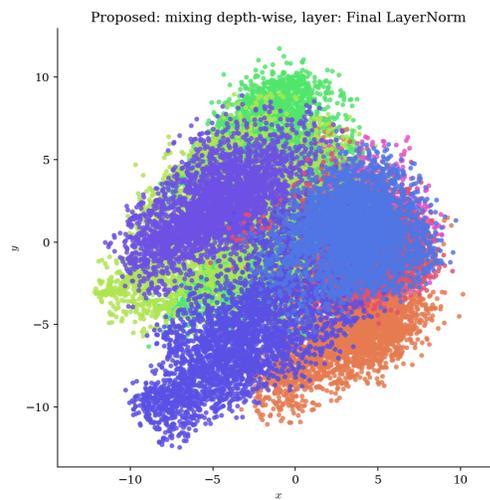
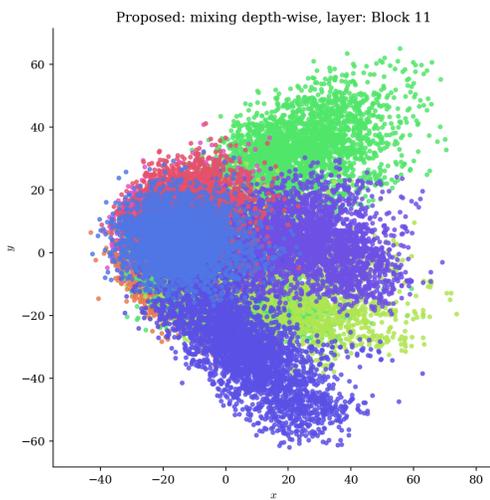
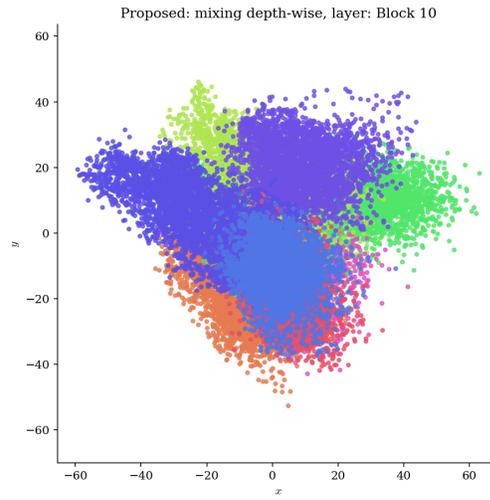
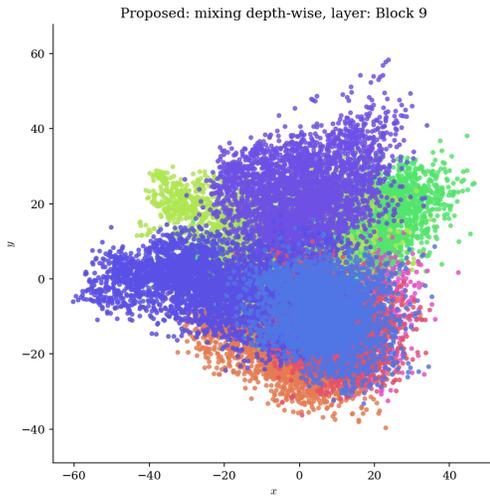


ViT-3P (Proposed):



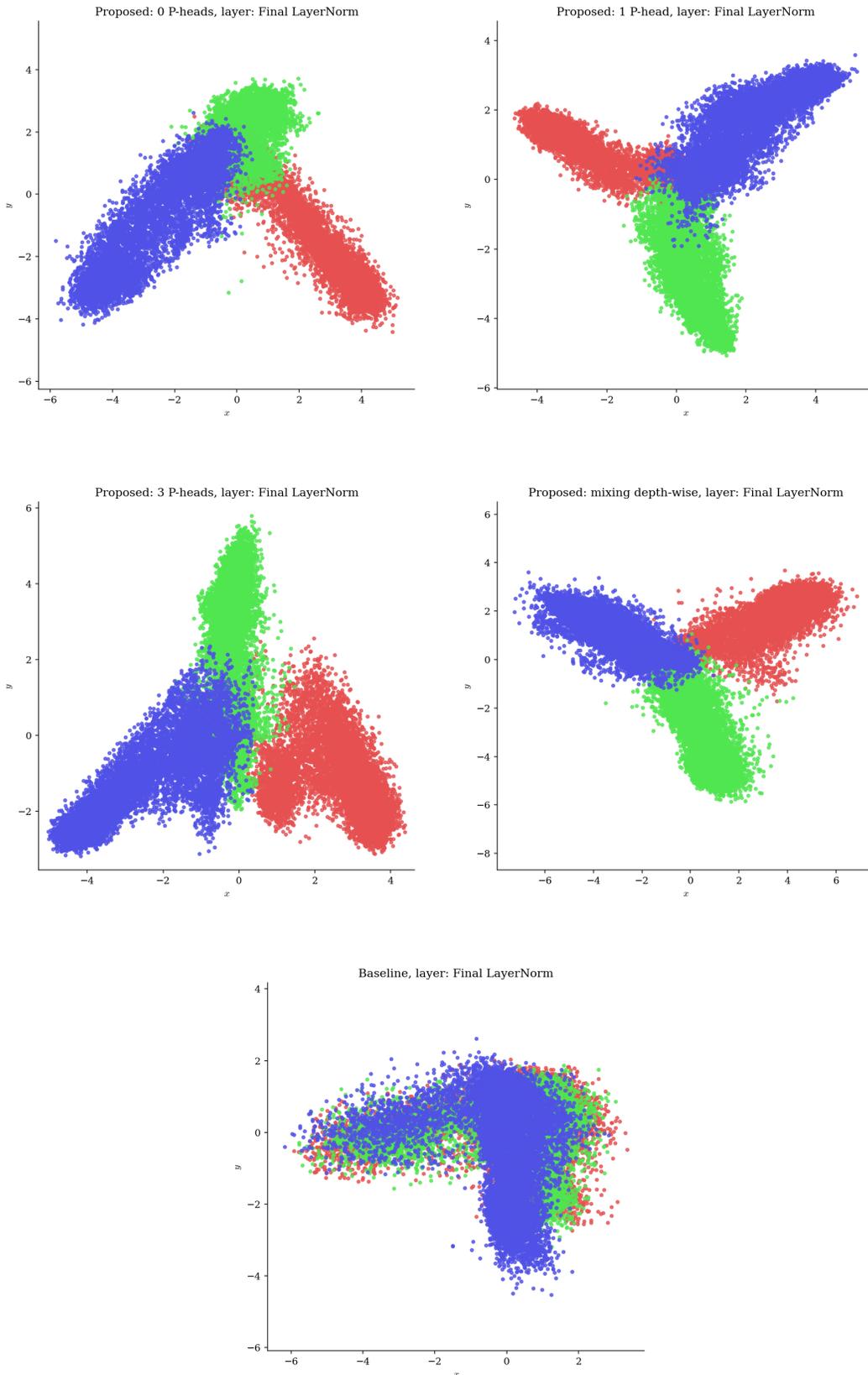
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

ViT-mix-depth (Proposed):



1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

C.1.3 IMAGENET-1K



1458 C.2 ANOVA DECOMPOSITION

1459

1460

1461

1462

1463 C.2.1 CIFAR10

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

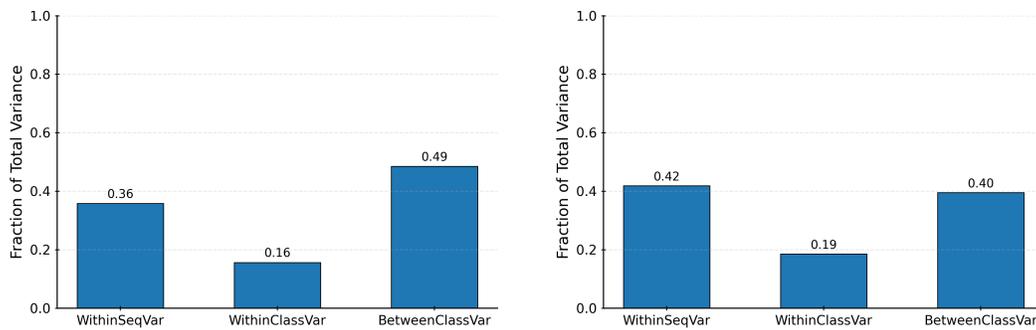
1479

1480

1481

1482

1483



1484 Figure 7: ANOVA decomposition of normalized variance for different model variants on CIFAR-10.
 1485 Left: ViT-B-3P; Right: ViT-B-Mix-Depth.

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511



Figure 8: ANOVA decomposition of normalized variance for ViT-B-0P on CIFAR-10.

C.2.2 CIFAR100

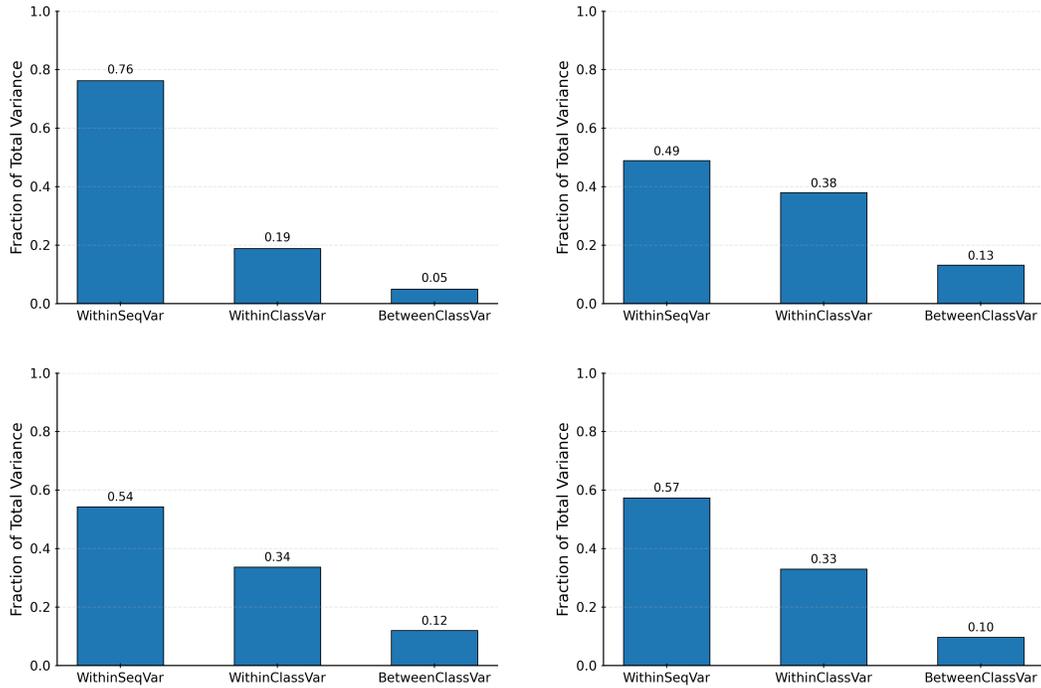


Figure 9: ANOVA comparison of normalized variance for different model variants on CIFAR-100. Top row: ViT-B (left) and ViT-B-1P (right). Bottom row: ViT-B-3P (left) and ViT-B-Mix-Depth (right).

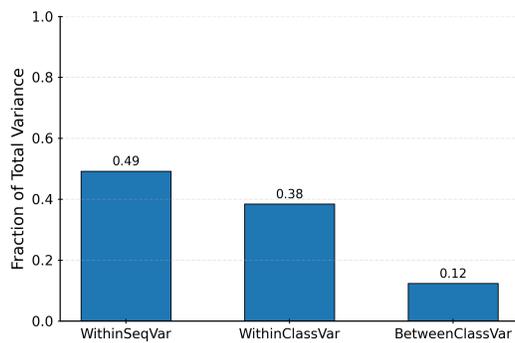


Figure 10: ANOVA decomposition of normalized variance for ViT-B-0P on CIFAR-100.

C.2.3 IMAGENET-1K

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

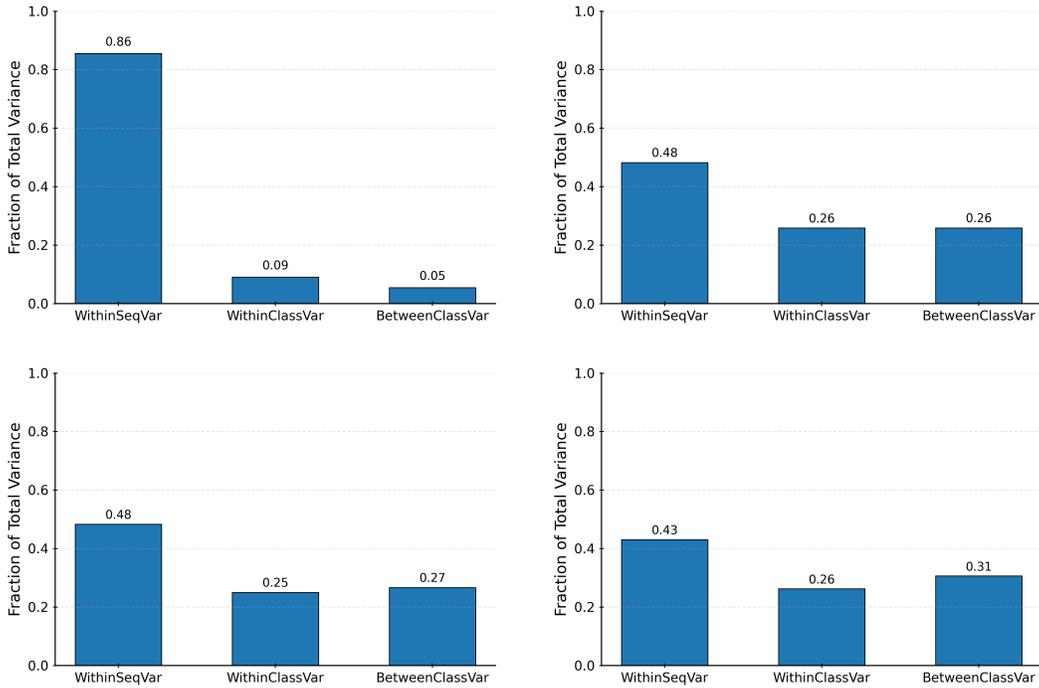


Figure 11: ANOVA comparison of normalized variance for different model variants on ImageNet-1k. Top row: ViT-B (left) and ViT-B-IP (right). Bottom row: ViT-B-3P (left) and ViT-B-Mix-Depth (right).

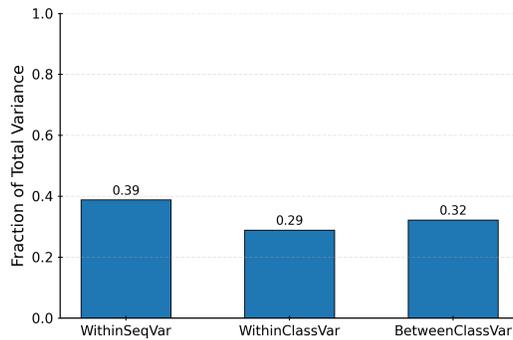


Figure 12: ANOVA decomposition of normalized variance for ViT-B-0P on ImageNet-1k.

C.3 CosSim METRIC

C.3.1 CIFAR10

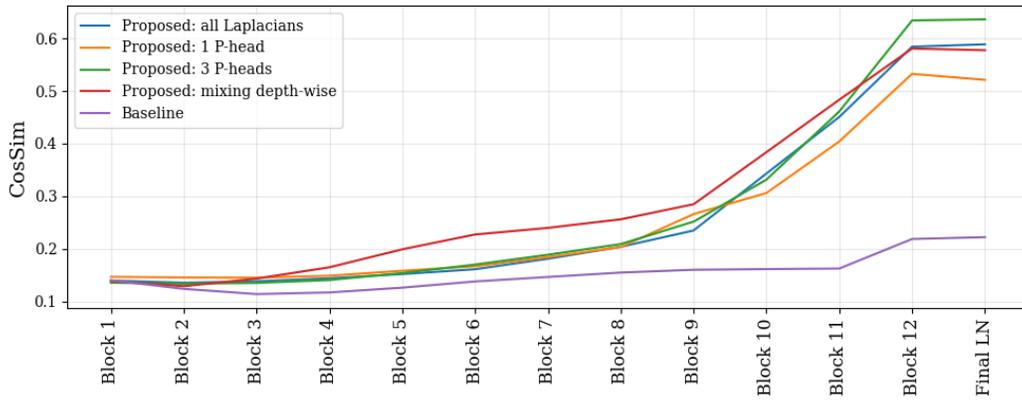


Figure 13: CosSim across depth on CIFAR10.

C.3.2 CIFAR100

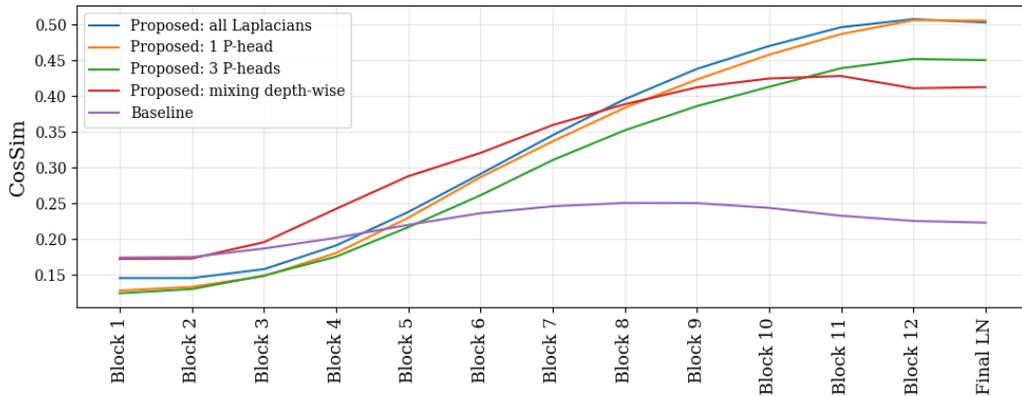


Figure 14: CosSim across depth on CIFAR100.

C.4 NC METRICS AND VISUALIZATION

C.4.1 CIFAR100

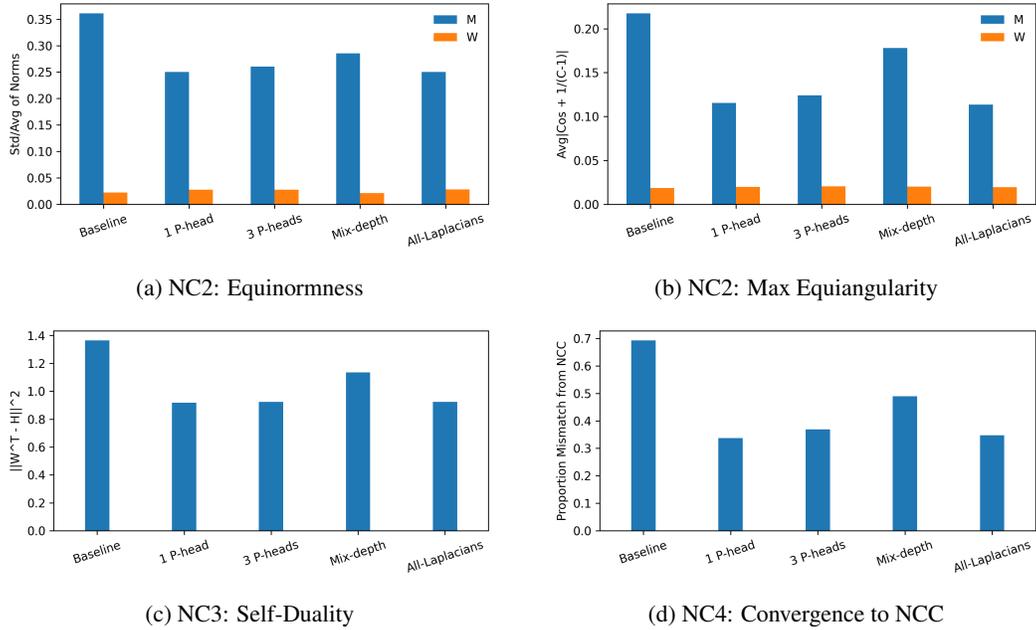


Figure 15: Neural Collapse metrics on CIFAR100.

C.4.2 IMAGENET-1K

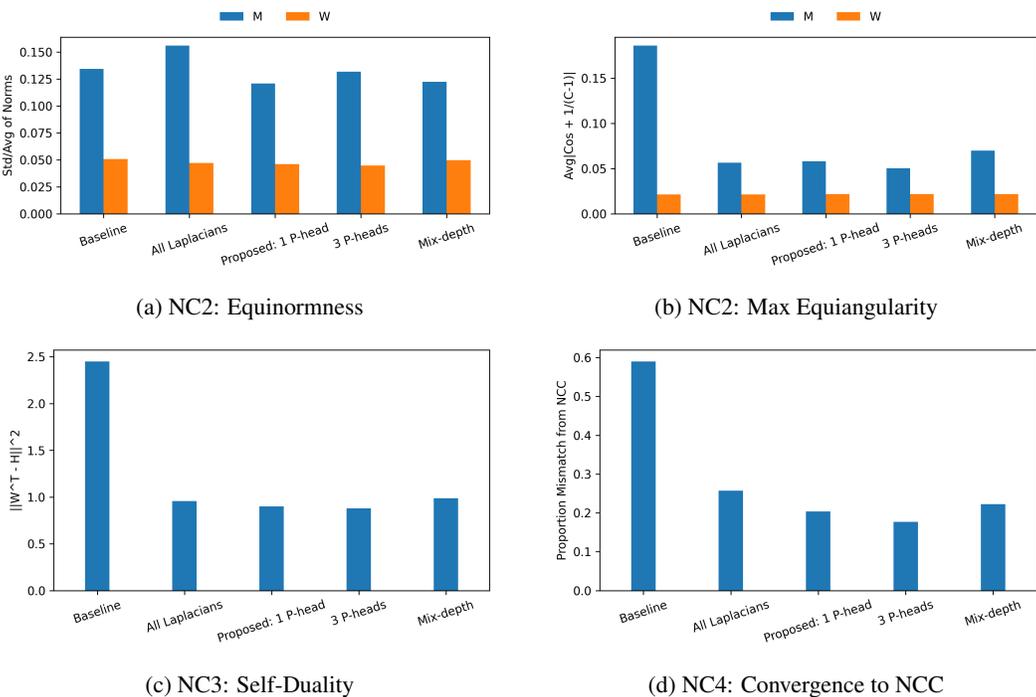


Figure 16: Neural Collapse metrics on ImageNet-1k.

C.5 VISUALIZATIONS OF PROJECTIONS ONTO A SIMPLEX

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

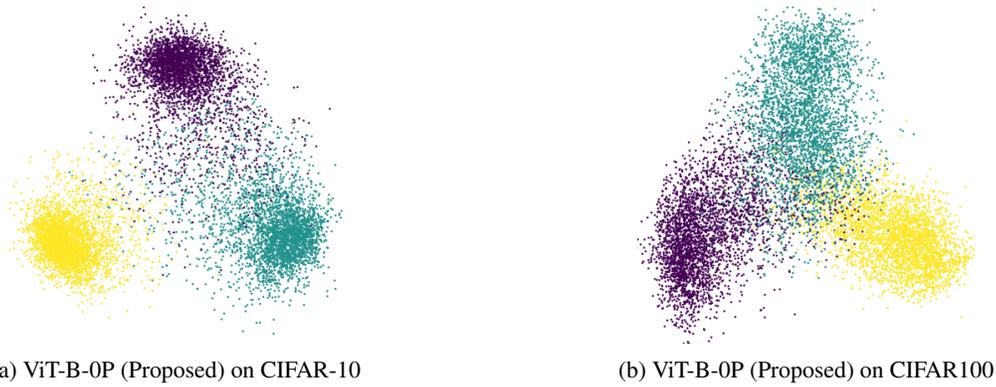


Figure 17: Visualization of projections onto a simplex for ViT-B-0P.

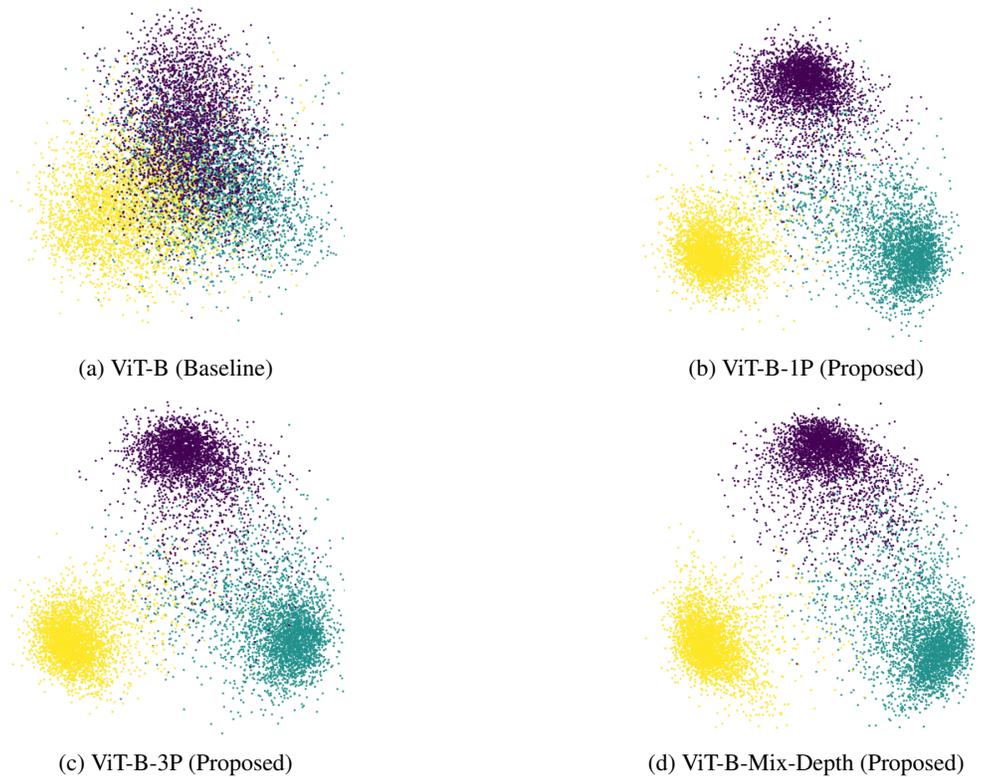


Figure 18: Visualization of projections onto a simplex on CIFAR10.

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

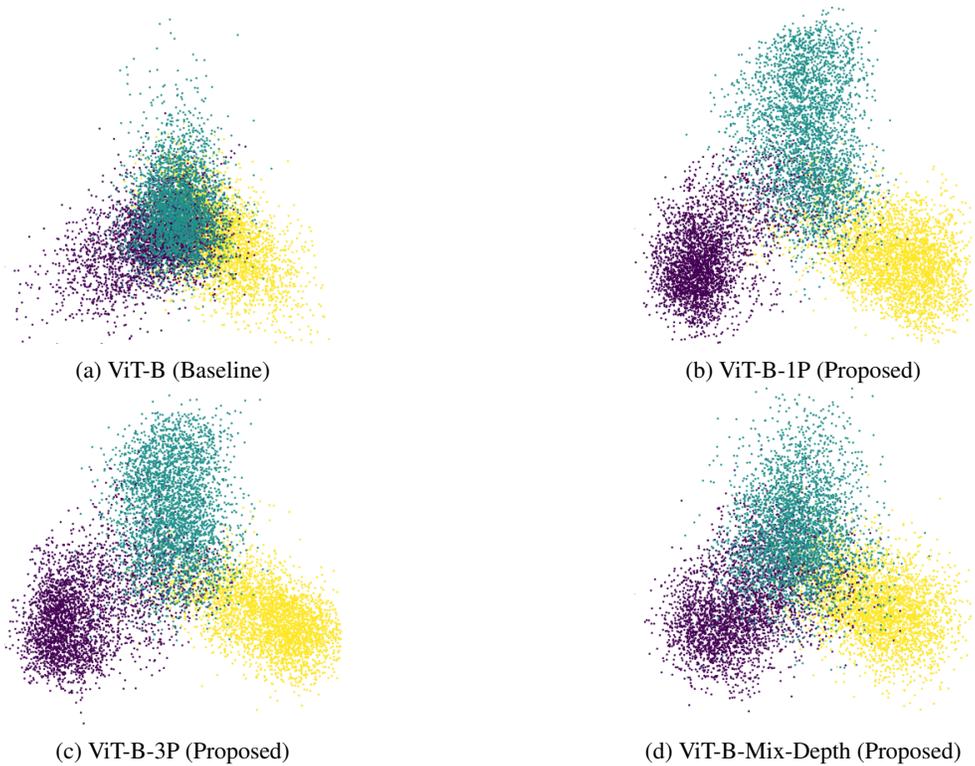


Figure 19: Visualization of projections onto a simplex on CIFAR100.

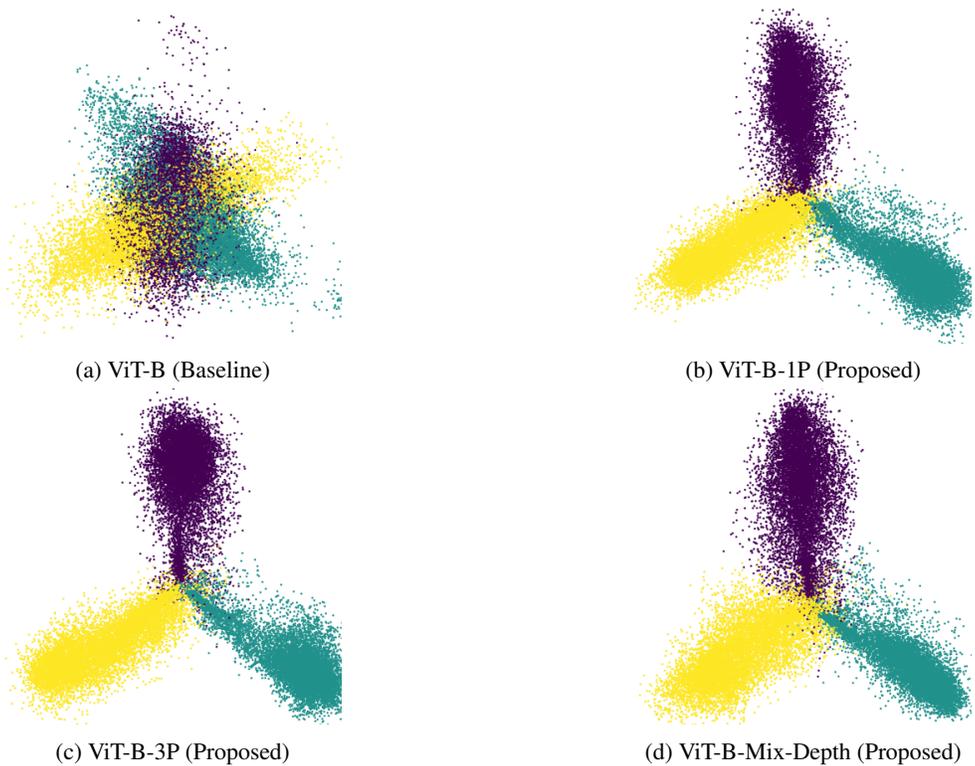


Figure 20: Visualization of projections onto a simplex on ImageNet-1k.

D EMPIRICAL EVIDENCE FOR THE GEOMETRIC ROLE OF THE LAPLACIAN

Given a batch of sequences of token embeddings $X \in \mathbb{R}^{B \times T \times d}$, we define the average signal-to-noise ratio (SNR) of X as

$$\text{SNR}(X) = \frac{1}{B} \sum_{b=1}^B \frac{\|\text{Mean}(X_b)\|_2}{\text{Std}(X_b)},$$

where

$$\text{Mean}(X_b) = \frac{1}{T} \sum_{i=1}^T X_{b,i} \quad \text{and} \quad \text{Std}(X_b) = \sqrt{\frac{1}{T} \sum_{i=1}^T \|X_{b,i} - \text{Mean}(X_b)\|_2^2}$$

Here, $X_b \in \mathbb{R}^{T \times d}$ denote the b th sequence in the batch and $X_{b,i} \in \mathbb{R}^d$ denote the i th token embedding within the sequence.

The SNR directly measures how large (in l_2 norm) the mean of a sequence of tokens is relative to their variance/standard deviation. A larger SNR implies that the sequence is more collapsed since the mean is larger relative to the variance. To validate our geometric interpretation in Sections 2.1 and 2.2 (see Figure 2), we directly measure the SNR of the output of the Layer Normalization module right before the MLP layer. In other words, we measure

$$\text{LayerNorm}(X + \text{MHA}(X))$$

for every transformer block. Figure 21 plots the SNR of the ImageNet token embeddings for the baseline and the proposed models as a function of depth. It clearly illustrates that for all models that use the Laplacian heads, the output of the Pre-MLP LayerNorm has significantly higher SNR than the baseline. Moreover, as depth increases, the SNRs for the proposed models grow more drastically. This measurement empirically supports our interpretation of the mechanism by which transformers collapse tokens, and it directly confirms our geometric intuition that the Laplacian induces more efficient collapse of token embeddings.

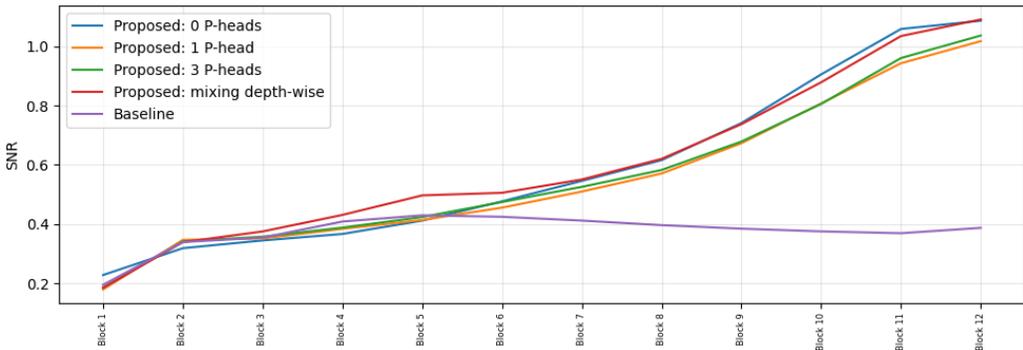


Figure 21: The Laplacian mechanism collapses tokens more effectively.

E RESULTS FOR ALTERNATIVE WAYS OF MIXING Attn AND \mathcal{L}

For a transformer with n blocks and h heads, there are in total $(h + 1)^n$ possibilities of assigning **Attn** and \mathcal{L} to different heads. Although the search space is huge, we hypothesize that the number of options that exhibit meaningful differences is much smaller. Here, we describe a limited subset of options that we experimented with.

The most obvious extension of strategy 1 is to vary the number m of heads that use **Attn** and keep this number the same across depth. The ViT-B model has 12 heads and 12 blocks, and we experimented with $m \in \{0, 1, 3, 6, 9, 12\}$. As discussed in Section 4.6.1 We observed that while $m = 0$ already produces noticeable improvements upon the baseline, incorporating a small number

of standard attention heads sometimes induces further improvements. This agrees with our intuition that it is beneficial to allow movement of tokens in both the mean and variance directions. In our experiments, $m = 1, 3$ consistently produced the best results across all datasets, and as m increased, the model’s performance converged to that of the standard transformer (equivalent to $m = 12$). We report the performance for different m values in Table 5.

Table 5: Top-1 test accuracy (%) of models with different numbers of standard attention heads.

| m | CIFAR-10 | CIFAR-100 | ImageNet-1k |
|-----|----------|-----------|-------------|
| 0 | 91.74 | 65.39 | 82.02 |
| 1 | 91.83 | 66.05 | 82.18 |
| 3 | 91.83 | 65.44 | 82.17 |
| 6 | 91.27 | 64.33 | 81.73 |
| 9 | 90.93 | 62.37 | 81.96 |
| 12 | 90.41 | 61.41 | 81.2 |

Another obvious strategy inverts strategy 2 by using \mathcal{L} for all heads in the first half of blocks and use **Attn** for all heads in the second half. However, this strategy consistently produced worse results than strategy 2. We also tried interleaving blocks that only used **Attn** and blocks that only used \mathcal{L} . Interestingly, the order in which the two types of blocks is interleaved appeared to impact performance significantly, where the order **Attn** \rightarrow \mathcal{L} \rightarrow **Attn** \rightarrow \mathcal{L} consistently performed better. More investigation is needed to understand these phenomena.

F NEURAL COLLAPSE

F.1 METRICS

Let $M \in \mathbb{R}^{d \times C}$ be the matrix whose columns are the class means $\{\mu_i : 1 \leq i \leq C\}$ and $W \in \mathbb{R}^{C \times d}$ be the weight matrix of the final-layer classifier. We quantify NC2 - NC4 following Han et al. (2022):

- **NC2 (Equinorm and Maximal Equiangularity):**

- *Equinorm*: Measures how uniform the vector norms are within the class means or weights, using the coefficient of variation (CoV):

$$\frac{\text{std}(\|\mu_c\|)}{\text{mean}(\|\mu_c\|)} \quad \text{and} \quad \frac{\text{std}(\|w_c\|)}{\text{mean}(\|w_c\|)},$$

where w_c is the classifier weight vector corresponding to class c .

- *Maximal Equiangularity*: Measures how close the vectors are to forming a maximally equiangular tight frame (ETF):

$$\frac{1}{C(C-1)} \sum_{i \neq j} \left| \langle \hat{v}_i, \hat{v}_j \rangle + \frac{1}{C-1} \right|,$$

where \hat{v}_i and \hat{v}_j are ℓ_2 -normalized class means or weight vectors. A lower value indicates greater conformity to an ETF structure.

- **NC3 (Self-Duality)**: Measures the alignment between the classifier weights and the centered class means:

$$\left\| \frac{W^T}{\|W^T\|_F} - \frac{M'}{\|M'\|_F} \right\|_F^2,$$

where $M' = M - \mu_G \mathbf{1}^T$ is the matrix of class means centered by their global mean μ_G .

- **NC4 (Convergence to NCC)**: Measures how close the learned classifier is to a Nearest Class Center (NCC) classifier:

$$1 - \frac{1}{N} \sum_{i=1}^N \mathbb{1} \left[\arg \max f(x_i) = \arg \min_c \|h_i - \mu_c\| \right],$$

where $f(x_i)$ are the logits, h_i is the feature of sample x_i , and μ_c is the mean feature for class c .

1944 F.2 VISUALIZATION OF PROJECTION ONTO SIMPLEX ETF
1945

1946 Each token embedding is first projected onto the classifier W for a random subset of three classes,
1947 then the result is projected again onto a two-dimensional representation of a three-dimensional sim-
1948 plex ETF. The result is visualized with each point colored according to its ground truth class. This
1949 visualization aims to illustrate the conformity of token embeddings to a simplex ETF.

1950 **Algorithm 2** Projection of Tokens to a simplex ETF

1951

1952 **Require:** $X \in \mathbb{R}^{B \times T \times d}, W \in \mathbb{R}^{C \times d}$
1953 1: $X \leftarrow \text{reshape}(X, [B \cdot T, d]), W' \in \mathbb{R}^{3 \times d} \leftarrow \text{random sample}(W)$
1954 2: $U, S, V^T = \text{SVD}(\text{normalize}(W'))$
1955 3: $A \leftarrow \sqrt{2} \cdot \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} & 0 \\ 0 & 0 & \frac{\sqrt{3}}{2} \end{bmatrix} \cdot (I_3 - \frac{1}{3}\mathbf{1}\mathbf{1}^T)$
1956
1957 4: output $AUV^T X^T$

1958

1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997