
What Positional Embeddings Really Do In Vision Transformers

Mahmoud Mannes¹

Abstract

Positional embeddings (PEs) in Vision Transformers (ViTs) are often viewed as simple spatial injectors that enable the model to encode absolute position. In this work, we show that even ViTs trained without PEs can recover meaningful spatial structure and distinguish relative position using patch content alone. This raises a central question: if spatial structure can be reconstructed without PEs, what functional role do PEs actually play? We find that PEs causally induce a dramatic increase in early-layer residual stream dimensionality and token heterogeneity, leading to richer image representations. PEs also introduce spatial structure prior to the encoder blocks and encourage representations that jointly rely on positional information and patch content. In contrast, models without PEs rely exclusively on content-based heuristics to infer spatial structure, a strategy that is significantly more fragile to distributional shifts.

1. Introduction

Vision Transformers (ViTs) have emerged as a powerful alternative to convolutional architectures for visual recognition by modeling images as sequences of patch tokens processed through self-attention mechanisms (Dosovitskiy et al., 2020). Unlike Convolutional Neural Networks (CNNs), however, ViTs lack strong built-in inductive biases toward locality and translation equivariance. As a result, most ViT architectures rely on Positional Embeddings (PEs) to inject explicit spatial information, enabling the model to distinguish between patches originating from different locations in the image.

Despite their widespread use, the precise role of positional embeddings in vision transformers remains incompletely understood. While PEs are generally assumed to be essential for encoding spatial structure, several empirical studies have shown that ViTs can retain non-trivial performance

even when explicit positional information is removed or degraded (Chu et al., 2021). These findings suggest that transformers may partially reconstruct spatial relationships from patch content alone, much like how CNNs learn implicit positional information from zero-padding (Islam et al., 2020), raising questions about what functional advantages positional embeddings actually provide beyond basic spatial identifiability.

Most prior work has investigated the role of positional embeddings primarily through downstream performance metrics or architectural variations. While informative, such analyses offer limited insight into how positional information influences the internal representations learned by the model. In particular, the impact of positional embeddings on the geometry, dimensionality, and stability of token representations throughout the transformer stack has received comparatively little attention.

In this work, we adopt a mechanistic perspective to study the role of positional embeddings in Vision Transformers. We analyze the evolution of token representations in the residual stream using tools from representational geometry (Raghu et al., 2021), with the goal of characterizing how spatial structure emerges and is maintained across layers. To support this analysis, we introduce the Spatial Structure–Diagonal Coefficient (SSDC), a metric designed to quantify the degree to which spatial relationships are reflected in token similarity patterns.

Using this framework, we conduct a comparative study of ViT models trained with and without positional embeddings. We examine how positional embeddings affect representational dimensionality, the nature of spatial reasoning employed by the model, and robustness to distributional shifts such as stylization and noise. Together, our results provide a more detailed picture of how explicit positional signals shape the internal organization of vision transformers, offering new insight into why positional embeddings play a critical role in stable and robust visual representation learning.

¹Independent Researcher. Correspondence to: Mahmoud Mannes <mannesmahmoud@gmail.com>.

2. Background and Setup

2.1. Vision Transformer Architecture

All models used in our experiments are vanilla Vision Transformers trained from scratch, with approximately 1.2M parameters. Images are divided into fixed-size patches, which are linearly projected into token embeddings and processed by a stack of self-attention and feedforward layers. When present, positional embeddings are learned parameters optimized jointly with the rest of the model. No architectural modifications or auxiliary inductive biases are introduced beyond standard ViT components.

2.2. Positional Embedding Ablation

To isolate the functional role of positional embeddings, we train a parallel set of models in which positional embeddings are entirely removed. Throughout the paper, we refer to models trained without positional embeddings as *ablated models*, and to models trained with positional embeddings as *intact models*. Apart from the presence or absence of positional embeddings, all architectural choices, optimization settings, and training procedures are held constant.

2.3. Datasets

We evaluate models on CIFAR-10 and a stylized variant derived from CIFAR-10 to assess robustness under distributional shift. The stylized dataset is generated using Adaptive Instance Normalization (AdaIN) with a mixing coefficient $\alpha = 0.1$, which significantly alters texture statistics while preserving coarse spatial structure. Models are trained on standard CIFAR-10 images and evaluated on both the original and stylized datasets.

3. Methods

3.1. Residual Stream Geometry

To analyze the evolution of internal representations across depth, we extract the residual stream at selected layers of the model (layers 0, 2, 4, and 9, where layer 9 corresponds to the final layer). At each layer, we represent the residual stream as a matrix $R \in \mathbb{R}^{T \times C}$, where T denotes the number of tokens and C the embedding dimension. Each row of R corresponds to the residual stream representation of a single token.

Given the singular values $\{\sigma_i\}$ of R , we compute the effective rank using the participation ratio:

$$\text{ER} = \frac{(\sum_i \sigma_i^2)^2}{\sum_i \sigma_i^4}.$$

Effective rank quantifies the number of dimensions that meaningfully contribute to the representation, with higher

values indicating more distributed and heterogeneous representations.

In addition, we compute pairwise cosine similarities between all token representations in R to form a token cosine similarity matrix, where the (i, j) -th entry corresponds to the cosine similarity between tokens i and j . This matrix is symmetric by construction. We average the token cosine similarity matrix across the batch dimension to obtain a layer-wise summary of inter-token relationships. This matrix serves both as a proxy for inter-token heterogeneity and as the basis for computing the Spatial Similarity Distance Correlation (SSDC).

3.2. Spatial Similarity Distance Correlation

To quantify the emergence of spatial structure, we introduce the Spatial Similarity Distance Correlation (SSDC). For a given layer, we compute the pairwise cosine similarity matrix between token representations and the corresponding matrix of pairwise spatial distances between token positions, where Spatial distance is defined as the Manhattan distance between patch coordinates on the image grid. SSDC is defined as the Spearman rank correlation between cosine similarity and the negative spatial distance, such that higher values indicate that tokens which are spatially closer tend to have more similar representations. Because SSDC measures the monotonic alignment between spatial proximity and representational similarity, it serves as a proxy for the presence of relative positional structure in the residual stream.

We use Spearman rank correlation to remain agnostic to the precise functional form relating spatial distance and representational similarity.

3.3. Fragility Score

To quantify a model’s sensitivity to distributional shifts, we define a simple *Fragility Score* (FS), which measures the relative drop in top-1 accuracy under distribution shift. It is defined as

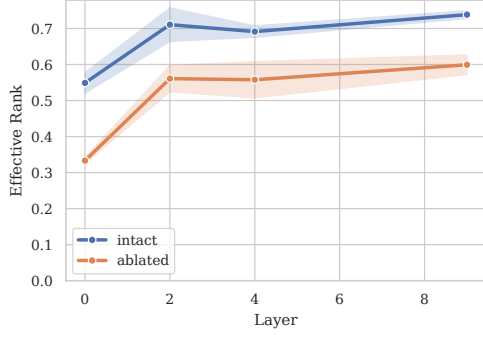
$$\text{FS} = 1 - \frac{A_{\text{shift}}}{A_{\text{normal}}},$$

where A_{normal} and A_{shift} denote top-1 accuracy on the normal and shifted datasets, respectively. Higher values indicate greater performance degradation.

4. Results

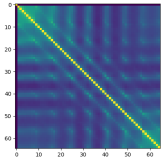
4.1. Evolution of Spatial Structure Across Depth

Experimental Setup: We evaluate SSDC at layers 0, 2, 4, and 9 on the CIFAR-10 dataset using the ablated and intact models. We also extract from each one of these layers the Token Cosine Similarity Matrix and plot them using the same color scale.

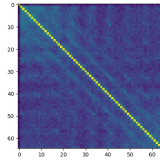


Raghu, Dosovitskiy, et al. Do vision transformers see like convolutional neural networks? 2021.

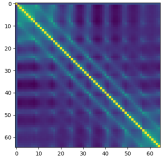
Figure 1. SSDC across layers for the intact and ablated models. Shaded regions indicate variability across runs.



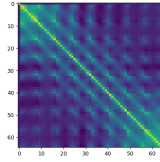
(a) Ablated model layer 0



(b) Intact model layer 0



(c) Ablated model layer 2



(d) Intact model layer 2

Figure 2. **Token relationships before and after going through encoder blocks** (a) Token Cosine Similarity on the ablated model at layer 0 (b) Token Cosine Similarity on the intact model at layer 0 (c) Token Cosine Similarity on the ablated model at layer 2 (d) Token Cosine Similarity on the intact model at layer 2

Results: Across all evaluated depths, the ablated model consistently exhibits a lower SSDC than the intact model, indicating a modestly reduced spatial organization. The difference is largest in layer 0 (i.e. tokens pre-encoder block) with an approximate value of 0.22, decreases at layer 2 to around 0.12, then remains roughly constant for the remaining layers. This trend is consistent across runs.

4.2.

References

- Chu, Zhi, et al. Conditional positional encodings for vision transformers. 2021.
- Dosovitskiy, Beyer, et al. An image is worth 16x16 words: Transformers for image recognition at scale. 2020.
- Islam, Jia, et al. How much positional information do convolutional neural networks encode? 2020.