

LOOPE: Learnable Optimal Patch Order for Positional Encoders in Vision Transformers

Anonymous CVPR submission

Abstract

Transformers are inherently permutation invariant, making positional encoding (PE) fundamental to their success. While prior work has focused on designing absolute or relative positional encoders, nearly all approaches implicitly assume that the 1D ordering of image patches is fixed. Yet for higher-dimensional data such as images, mapping a 2D grid into a 1D sequence is itself a geometric operation—one that can distort locality, disrupt neighborhood structure, and ultimately limit how effectively PEs capture spatial relations. Surprisingly, the impact of the order of patches on positional representation has received little attention. We revisit this overlooked design dimension and show that patch order is not merely a preprocessing choice but a learnable degree of freedom that fundamentally shapes positional embedding quality. We propose LOOPE, a lightweight framework that learns an image-dependent ordering of patches. Starting from a locality-preserving space-filling curve and applying small context-aware refinements, LOOPE yields an interpretable and content-adaptive coordinate system on which standard sinusoidal encodings operate robustly. Across multiple ViT and ViT-based architectures, LOOPE improves performance and stabilizes PEs under challenging settings such as non-optimal frequency schedules. To complement standard benchmarks, we introduce a simple diagnostic task revealing that effective PEs can preserve far more positional information than previously reported. Our findings highlight learnable patch ordering as a powerful and largely untapped tool for improving positional encoding in vision Transformers.

1. Introduction

Transformers have become the backbone of modern computer vision tasks—from image classification to image segmentation—yet they inherently lack explicit spatial ordering. Positional encoding (PE) is therefore essential for capturing spatial relationships. While numerous PE vari-

ants exist, effectively modeling 2D spatial structure remains challenging.

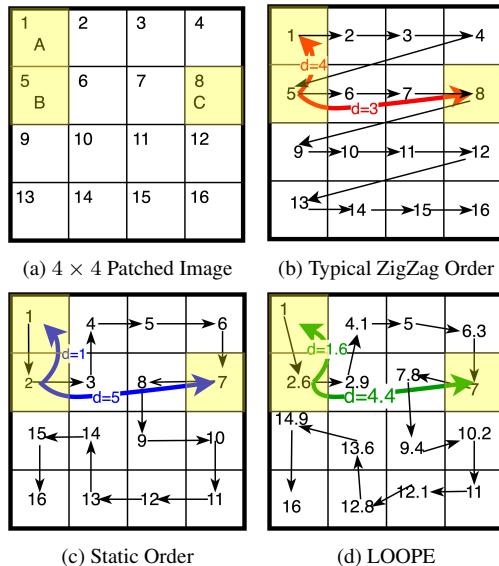


Figure 1. Fig. 1(b) shows that under Generic ZigZag Order, the distances between patches are $d(B, A) = 4$ and $d(B, C) = 3$, even though B is spatially and contextually closer to A. Static Order (Fig. 1(c)) improves this with $d(B, A) = 1$ and $d(B, C) = 5$, but it remains static. LOOPE (Fig. 1(d)) dynamically reorders based on context, yielding $d(B, A) = 1.6$ and $d(B, C) = 4.4$.

Traditional PEs—such as sinusoidal encodings [1], learned absolute encodings [1], and relative positional encodings (RPEs) [25]—either impose rigid spatial biases or introduce substantial parameter overhead. Conditional encoders like CPE [5] incorporate local information but discard global coordinates, while frequency-based approaches [22, 31, 37] and Learnable Fourier Features (LFF) emphasize spectral design. However, a common assumption underlies all these methods: the 1D ordering of image patches is fixed (e.g., raster or zigzag) and independent of image content.

Most positional encoders can be represented by Eq. 1,

050 where PE construction depends on a frequency set, a sequence of patch indices, and an optional phase:
 051

$$052 \quad \text{PE} = \text{Policy}\left(\begin{bmatrix} \omega_1 \\ \vdots \\ \omega_j \\ \vdots \\ \omega_L \end{bmatrix}, [\underline{\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N}], \frac{\pi}{2}\delta\right) \quad (1)$$

053 Yet none of these approaches examine **how the patch**
 054 **sequence itself should be constructed.** This is notable
 055 because $2D \rightarrow 1D$ flattening determines which patches be-
 056 come neighbors in the sequence, directly influencing local-
 057 ity preservation, similarity patterns, and how frequencies in-
 058 teract with spatial structure.

059 We observe that periodic PEs depend on three coupled
 060 factors: **context**, **order**, and **frequencies**. Contextual en-
 061 coders [15] and frequency-based methods, including LFF,
 062 explore portions of this space, but the joint effect of con-
 063 text and patch order is almost entirely unexplored. Addi-
 064 tionally, although RPEs often perform well on downstream
 065 tasks, they encode positional information inside attention
 066 weights and do not preserve a reusable absolute coordi-
 067 nate system. As shown later in the Three-Cell Experiment
 068 (Sec. 4.4), when labels depend purely on geometry, well-
 069 structured absolute PEs can outperform strong RPEs, high-
 070 lighting the continuing importance of absolute encodings.
 071 In this work, we address this gap by treating patch ordering
 072 as a learnable, context-aware variable and by studying its
 073 interaction with sinusoidal PEs. To summarize,

- 074 1. Our primary contribution is the introduction of **LOOPE**,
 075 a lightweight patch-ordering framework that combines
 076 a topology-preserving static Gilbert curve X_G with
 077 a context-aware refinement X_C , yielding an image-
 078 dependent coordinate system that minimizes struc-
 079 tural distortion and enhances robustness to frequency
 080 choices. To s
- 081 2. To prove the locality preserving nature of LOOPE,
 082 we additionally introduce lightweight structural probes
 083 (PESI) for analyzing monotonicity and symmetry in
 084 learned positional representations.
- 085 3. Finally, we propose a simple diagnostic task (Three-Cell
 086 Experiment) revealing that carefully structured absolute
 087 PEs can retain substantially richer spatial information
 088 than task-oriented RPEs.

089 Extensive experiments across ViT and ViT-based archi-
 090 tectures demonstrate that LOOPE not only improves per-
 091 formance but also provides a clearer geometric under-
 092 standing of how positional information is encoded within trans-
 093 former architectures.

094 2. Related Works

095 Positional encoding (PE) is crucial for ViTs as self-attention
 096 lacks spatial awareness. Absolute positional encodings

(APEs) [1], first introduced in NLP, were adapted for vision
 097 tasks where images are tokenized into patches. Models like
 098 ViT [10] and DeiT [32] use fixed sinusoidal or learnable en-
 099 coders. APEs provide positional awareness but suffer from
 100 fixed-size constraints, limiting adaptability to varying reso-
 101 lutions in models such as CrossViT [2] and Swin [19], and
 102 fail to capture relative spatial relationships.
 103

104 Relative positional encodings (RPEs) [26, 36] encode
 105 pairwise relationships, generalizing across sequence lengths
 106 for NLP [8] and vision tasks, e.g., axial attention [34] and
 107 2D-aware encodings [23]. However, they discard absolute
 108 position information, limiting fine-grained spatial represen-
 109 tation, particularly in object localization [27].

110 Hybrid PE strategies balance absolute and relative cues.
 111 CP-RPE [36] encodes horizontal and vertical distances sep-
 112 arately but struggles with optimal mapping. CPE [5] gener-
 113 ates encodings from local context, improving translation
 114 invariance but is hyperparameter-sensitive and lacks struc-
 115 tured positional order. RoPE [12, 29] preserves relative po-
 116 sitioning via rotational invariance but misses 2D spatial hi-
 117 erarchies, whereas AS2DRoPE [6] introduces scalable 2D
 118 priors. Learned Fourier Features (LFF) [14] further refine
 119 spatial representation by learning a spectral basis over fixed
 120 coordinates; however, LFF assumes a predetermined patch
 121 ordering and does not provide a mechanism for constructing
 122 or adapting the coordinate system itself.

123 ViTs are widely used for segmentation. SETR [38], Seg-
 124 menter [28], and DPT [24] use ViT backbones with learned
 125 APEs. Mask2Former [4] uses Swin or ResNet [11] back-
 126 bones, inheriting RPE-based design.

127 3. Methodology

128 3.1. Motivation

129 3.1.1. A Geometric Perspective on the Problem

130 Let an image be partitioned into an $h \times w$ grid of patches,
 131 giving $N = hw$ tokens. A Vision Transformer must process
 132 these tokens as a 1D sequence, which requires a bijection

$$\pi : \{1, \dots, h\} \times \{1, \dots, w\} \longrightarrow \{1, \dots, N\} \quad (2)$$

133 mapping every 2D patch coordinate to a unique sequence
 134 index. The quality of this mapping depends on how well
 135 it preserves spatial relationships: ideally, patches that are
 136 close in the grid remain close in the sequence. In geom-
 137 etric terms, the goal is to minimize distortion of Euclidean
 138 distances under the mapping π , i.e.
 139

$$140 \| (i, j) - (i', j') \|_2 \approx |\pi(i, j) - \pi(i', j')| \quad (3)$$

141 so that neighborhood structure is not destroyed when flat-
 142 tening the image.

143 This is precisely the principle behind space-filling curves
 144 such as Hilbert [13] or Peano [21] traversals, which are

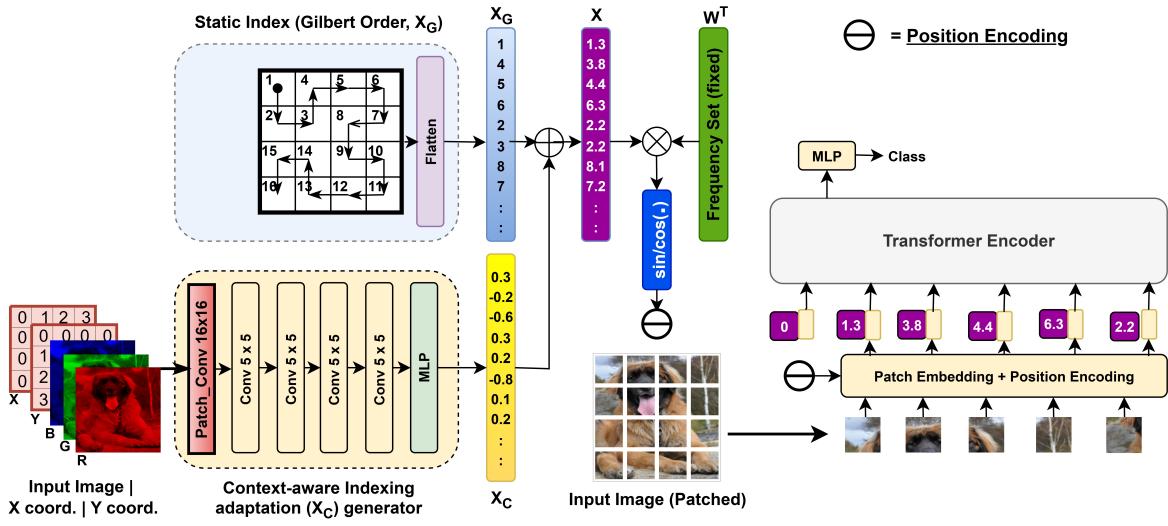


Figure 2. The complete architecture of LOOPE. For architectural details refer to *Supp. A.2*

designed to unfold two-dimensional domains into one-dimensional paths with strong locality preservation.

For ViTs, this ordering problem is not incidental but fundamental. Since self-attention is permutation-invariant, positional encodings are the only source of spatial information. If the initial $2D \rightarrow 1D$ mapping π introduces large distortions, the positional encoder must spend capacity correcting for these artifacts rather than representing meaningful structure—and even sophisticated APEs or RPEs cannot fully recover geometric relations once the coordinate system is distorted. Frequency-based methods such as LFF learn powerful positional functions but still assume a fixed ordering; our perspective instead optimizes the coordinate system itself.

In this light, patch ordering is not a cosmetic design choice but a structural constraint that determines how much geometric information a positional encoding can ultimately preserve.

3.1.2. Limitations of Existing Dynamic SFC Approaches

Prior attempts to learn adaptive space-filling curves (SFCs) [3, 35] highlight the importance of patch ordering, but they face two key obstacles. Most rely on discrete procedures such as Minimum Spanning Trees (MSTs) to build Hamiltonian traversals, which cannot be optimized end-to-end because they lack differentiable loss functions. GNN-based variants [35] estimate edge weights before applying MST heuristics, but this introduces substantial computational overhead. Moreover, these methods optimize curves independently of image content and therefore do not yield sample-specific orderings. While suitable for small graphs, such pipelines are far too heavy for positional encodings in ViTs, which must remain lightweight and effi-

cient. In short, existing contextual SFC generators are either non-differentiable or computationally impractical—leaving open the need for a method that is both trainable and efficient in the transformer setting.

3.2. Proposed Method

To address these limitations, we propose **LOOPE**, a learnable patch-ordering framework that unifies two complementary components: a static space-filling order $X_G \in \mathbb{Z}_+^{1 \times N}$ and a context-aware refinement $X_C \in [-1, 1]^{1 \times N}$, where N is the number of patches. **The static part ensures stability and locality preservation, while the dynamic part adapts to image content**, yielding an indexing scheme that is both efficient and trainable. In the following, we first describe the static order and then detail the context-aware refinement.

Static Patch Index (Gilbert Order, X_G): The Hilbert curve maps a $2^n \times 2^n$ grid, $n \in \mathbb{Z}_+$, into a 1D sequence while preserving locality, but it is restricted to square grids of power-of-two size. To generalize, we adopt the **Gilbert order** [40], which recursively partitions arbitrary rectangular grids while maintaining strong spatial coherence (see Fig. 1c). This improves over typical zigzag order by reducing locality distortion, but its effectiveness still depends on the chosen frequency set, as analyzed in Table 6.

Context-Aware Index Adaptation (X_C): The static Gilbert order X_G offers stability and locality preservation, but it remains tied to the hand-crafted choice of frequency set (Table 6). To address this, we introduce a learnable *context-aware bias* X_C that adapts the patch order to the input image itself. As shown in Fig. 2, the generator G takes as input the concatenated tensor $[Image = I_0 \in \mathbb{R}^{3 \times H \times W}, coordinates = x, y \in \mathbb{R}^{1 \times H \times W}] \in \mathbb{R}^{5 \times H \times W}$

and outputs continuous index offsets. These offsets are added to X_G , yielding a refined ordering like in Fig. 3, 1d that is both content-adaptive and differentiable.

This refinement introduces two decisive properties. First, it **mitigates frequency sensitivity**: unlike methods tied to a fixed frequency design, X_C adaptively corrects the encoding, keeping it stable across different frequency sets (Table 6). Second, it enables **fractional indexing**, where patch indices are no longer restricted to integers—allowing patches to be pulled closer, pushed apart, or locally reordered while preserving differentiability. This flexibility gives X_C the freedom to offset poor frequency choices and to adjust neighborhood geometry in a lightweight manner (see $\partial \text{PE}/\partial W$ analysis in Supp.A.3). Thus, while X_C introduces spatial awareness, its central role is to extend X_G into a geometry-preserving, *frequency-robust* framework—far more than a simple spatial inductive bias.

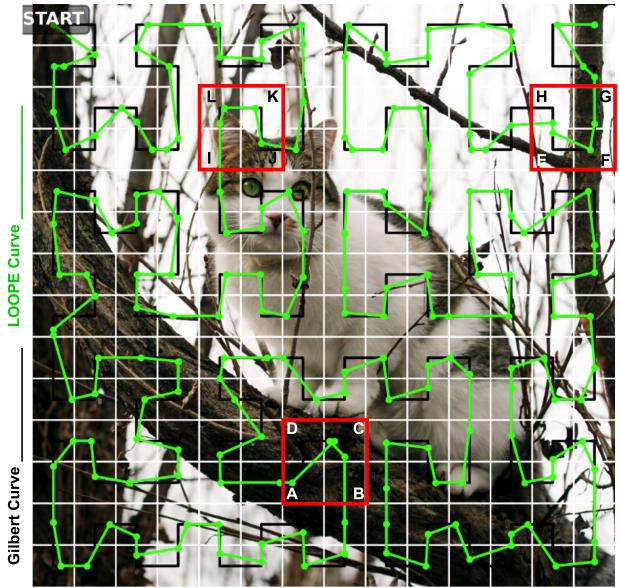


Figure 3. Here is an example of the complete generated curve. i) $d_{IL} \downarrow$ to capture ear. ii) $d_{EH} \downarrow$ (sky to sky) and $d_{EF} \uparrow$ (sky to trunk). iii) $d_{CD} \downarrow$ to swap C(Branch), D(Branch+Twig). Seemingly, A(Branch) is more similar to C than D, making $ADCB \rightarrow ACDB$. Contextually similar patches are drawn closer. Here, \downarrow = distance decreased, \uparrow = increased.

Final Formulation: Combining the static Gilbert order X_G with the context-aware refinement X_C , LOOPE defines the final positional encoding as

$$\text{PE}(X) = \sin\left(W(X_G + X_C) + \frac{\pi}{2} \delta\right), \quad (4)$$

where $W = [\omega_1, \dots, \omega_L]^T$ denotes the chosen frequency set (a design choice) and $\delta \in \mathbb{Z}_2^{L \times N}$ is a fixed phase matrix with entries $\delta_{i,j} = (i+j) \bmod 2$.

To understand why X_C improves robustness to frequency choices, consider the j -th sinusoidal channel of Eq. 4: $\text{PE}_j = \sin(\omega_j(X_G + X_C) + \frac{\pi}{2} \delta_j)$. A simple sensitivity analysis gives

$$\frac{\partial \text{PE}_j}{\partial \omega_j} = (X_G + X_C + \omega_j \frac{\partial X_C}{\partial \omega_j}) \cos(\omega_j(X_G + X_C) + \frac{\pi}{2} \delta_j). \quad (237)$$

Here X_C does not take ω_j as an explicit input; its dependence on ω_j is implicit through training, since different frequency sets lead to different learned refinements. This post-training sensitivity analysis (see Supplementary Sec. A.3 for the full derivation) highlights two terms: a static geometric component ($X_G + X_C$) and an adaptive correction $\omega_j \partial X_C / \partial \omega_j$. The latter enables LOOPE to reshape its effective positional geometry under different frequency sets, explaining its empirical frequency robustness (Table 6).

Equation 4 summarizes the essence of LOOPE: a **stable static order** enriched by a **learnable, content-adaptive bias**, producing a lightweight yet flexible encoding that preserves geometry while reducing design sensitivity.

3.3. Positional Embeddings Structural Integrity (PESI) Metrics

PESI is designed as a set of structural probes that assess how well a positional embedding maintains radial and directional monotonicity, and how much radial symmetry it sacrifices to achieve higher precision.

Given a positional embedding tensor $P \in \mathbb{R}^{h \times w \times D}$, define the *cosine similarity matrix* centered at (x, y) :

$$E_{(x,y)}(i, j) = \frac{P_{(x,y)} \cdot P_{(i,j)}}{\|P_{(x,y)}\| \|P_{(i,j)}\|}. \quad (5)$$

Undirected monotonicity. The *radial average similarity function* is:

$$\mu_{(x,y)}(r) = \frac{1}{|B_r|} \sum_{(i,j) \in B_r} E_{(x,y)}(i, j) \quad (6)$$

where B_r is the set of positions at radius r . We compute *Spearman's rank correlation*:

$$\rho_{(x,y)} = 1 - \frac{6 \sum d_r^2}{|R|(|R|^2 - 1)} \quad (7)$$

where d_r is the rank difference between $\text{level}(r)$ and $\mu_{(x,y)}(r)$, and $|R|$ is the total number of radial levels. The undirected monotonicity score is then

$$M_u = \frac{1}{hw} \sum_{x=0}^{h-1} \sum_{y=0}^{w-1} (1 - \rho_{(x,y)}), \quad (8)$$

where a higher M_u indicates stronger undirected monotonicity across the grid. Ideally, it should approach 2.

272 **Directed monotonicity.** We quantize 2π into $N = \frac{2\pi}{\delta}$
 273 directional buckets (with quantization angle δ). For each
 274 cell (i, j) relative to (x, y) , compute:

$$275 \quad \theta_{(x,y)}(i, j) = \text{atan2}(j - y, i - x) \quad (9)$$

276 where $\text{atan2}(y, x)$ returns the angle in $(-\pi, \pi]$. Assign the
 277 cell to bucket

$$278 \quad k = \left\lfloor \frac{\theta_{(x,y)}(i, j)}{\delta} \right\rfloor \mod N. \quad (10)$$

279 Within each bucket k , order the cells by radial distance and
 280 compute Spearman’s rank correlation:

$$281 \quad \rho_{(x,y)}^k = 1 - \frac{6 \sum_r d_r^2}{|R_k|(|R_k|^2 - 1)}, \quad (11)$$

282 where d_r is the rank difference between the radius r and the
 283 corresponding similarity values in bucket k , and $|R_k|$ is the
 284 number of elements in bucket k . The mean correlation per
 285 cell is then

$$286 \quad \bar{\rho}_{(x,y)} = \frac{1}{N} \sum_{k=0}^{N-1} \rho_{(x,y)}^k \quad (12)$$

287 and the global directed monotonicity measure is defined as

$$288 \quad M_D = \frac{1}{hw} \sum_{x=0}^{h-1} \sum_{y=0}^{w-1} (1 - \bar{\rho}_{(x,y)}). \quad (13)$$

289 With varying N and δ , we can investigate how precisely the
 290 encoder maintains monotonicity along different directions.
 291 As $N \rightarrow 1$, M_D reduces to the undirected measure M_U .
 292 A higher M_D indicates stronger directional monotonicity;
 293 ideally it approaches 2.

294 **Undirected asymmetry.** For each center (x, y) , let B_r
 295 be the set of cells at radius r from (x, y) . We define the
 296 standard deviation of the cosine similarity values as

$$297 \quad \sigma_{(x,y)}(r) = \sqrt{\frac{1}{|B_r|} \sum_{(i,j) \in B_r} (E_{(x,y)}(i, j) - \mu_{(x,y)}(r))^2}. \quad (14)$$

298 The coefficient of variation at radius r is then

$$299 \quad \text{CV}_{(x,y)}(r) = \frac{\sigma_{(x,y)}(r)}{\mu_{(x,y)}(r)}. \quad (15)$$

300 Averaging over all radial distances $r \in R$ yields the undi-
 301 rected asymmetry measure at (x, y) :

$$302 \quad A'_{SU}(x, y) = \frac{1}{|R|} \sum_{r \in R} \text{CV}_{(x,y)}(r), \quad (16)$$

303 and the global undirected asymmetry is defined as

$$304 \quad A_{SU} = \frac{1}{HW} \sum_{x=1}^H \sum_{y=1}^W A'_{SU}(x, y). \quad (17)$$

For complete symmetry, $|A_{SU}| \rightarrow 0$. In practice, there is no single ideal value of undirected asymmetry, since many encoders increase A_{SU} to achieve sharper directional monotonicity M_D ; if $A_{SU} = 0$, there is no directional information in the embedding. Detailed algorithms are provided in *Supp. A.5.1, A.5.2, A.5.3*.

4. Experiments

4.1. Experimental Setup

All experiments were conducted on a single NVIDIA RTX 5090 GPU (32 GB VRAM). Unless specified otherwise, we follow the standard training protocol used in prior ViT literature to ensure fair comparison with existing positional encoders.

We train all models for 150 epochs using the Adam optimizer with a cosine learning-rate schedule (max LR = 1×10^{-3} , min LR = 2.5×10^{-5}) and Cross-Entropy loss. Batch sizes are 96 for Oxford-IIIT [20] and 64 for CIFAR-100 [17], ImageNet-1k [9], and our Three-Cell probing dataset. For the resolution-scaling experiment, Oxford-IIIT images are trained at 384×384 with batch size 32. CrossViT models use 240×240 inputs with mixed patch sizes (12×12 and 16×16), following the original implementation.

To ensure a strictly controlled comparison across positional encodings, all models—including baseline APE/RPE variants—use the same ImageNet-1K pretrained backbones (DeiT-Base unless otherwise stated) while PEs are trained from scratch. This removes differences due to random initialization and isolates the contribution of the positional encoding itself, which is the primary quantity under investigation.

We adopt an 80-10-10 train-validation-test split for all datasets and apply standard augmentations (horizontal flips, rotations, brightness jitter, and elastic deformations) consistently across all PE variants.

ViT-Base contains **85.8M** parameters. LOOPE increases this to only **85.9M** (a **0.12%** overhead), confirming that the proposed content-aware patch-ordering mechanism adds negligible computational cost while providing significant gains in positional fidelity and downstream accuracy.

4.2. Experiments on Image Classification

4.2.1. Comparison against 1-D Positional Encoders

From Table 1, LOOPE achieves the highest accuracy across all models on Oxford-IIIT, most notably with ViT, indicating enhanced fine-grained feature learning. CIFAR-100 presents greater inter-class variability, still our PE outperforms all other PEs, especially with DeiT-Small. Finally, for ImageNet-1k, LOOPE again surpasses all PEs across all models. These results demonstrate that it effectively bal-

Dataset	Model	No PE	Learnable	Sinusoid	Static (X_G)	LOOPE ($X_G + X_c$)
Oxford-IIIT	ViT-Base [10]	83.6%	84.6%	85.3%	84.2%	88.1%
	DeiT-Base [32]	88.9%	89.4%	86.3%	89.0%	89.8%
	DeiT-Small [32]	83.8%	83.8%	83.7%	80.6%	84.5%
	CaiT [33]	87.4%	89.0%	90.0%	89.6%	90.5%
	Cross-ViT [2]	88.3%	90.9%	88.0%	89.3%	91.0%
CIFAR-100	ViT-Base [10]	79.8%	83.0%	85.2%	87.6%	88.3%
	DeiT-Base [32]	82.1%	86.3%	86.6%	86.9%	87.1%
	DeiT-Small [32]	68.6%	81.6%	71.9%	77.7%	82.0%
	CaiT [33]	77.3%	82.5%	82.3%	82.5%	83.1%
	Cross-ViT [2]	80.5%	84.6%	86.3%	85.3%	86.8%
ImageNet-1k	ViT-Base [10]	82.7%	83.6%	84.4%	84.5%	84.7%
	DeiT-Base [32]	81.7%	82.1%	82.0%	83.1%	83.5%
	DeiT-Small [32]	80.5%	82.2%	82.8%	82.8%	83.0%
	CaiT [33]	80.3%	81.6%	81.6%	82.8%	83.1%
	Cross-ViT [2]	82.4%	83.5%	83.1%	84.5%	84.8%

Table 1. Accuracy comparison of different ViT models with various PEs across Oxford-IIIT, CIFAR-100 and ImageNet-1k datasets.

ances structured spatial encoding with learnable adaptability, making it a robust solution.

4.2.2. More comparisons with Positional Encoders

PE	ImageNet-1k	Oxford-IIIT	CIFAR-100
CPE [5]	82.7%	83.9%	79.1%
RPE [36]	80.3%	80.5%	79.2%
LFF(Fourier)[18]	83.4%	90.5%	89.1%
2D Sinusoid [39]	83.2%	80.1%	86.3%
AS2DRoPE [6]	82.9%	88.7%	86.4%
Static (X_G) [13]	83.1%	89.0%	86.9%
LOOPE ($X_G + X_c$)	83.5%	89.8%	87.1%

Table 2. Accuracy comparison of LOOPE against other advanced PEs on DeiT-Base across Oxford-IIIT, CIFAR-100 and ImageNet-1k datasets

Table. 2 compares LOOPE against advanced PEs. Fourier PE achieves the highest accuracy for Oxford IIT and CIFAR-100, due to its rich frequency encoding properties. But, LOOPE outperforms Fourier in ImageNet-1k. These results validate the superiority of Fourier PE while demonstrating that LOOPE remains a strong alternative for vision tasks.

4.3. Experiments on Semantic Segmentation

To test the impact of LOOPE on downstream tasks, We used the Cityscapes [7] dataset for its strong sensitivity to positional information. SETR [38] and Segmenter [28] were chosen as base models for their state-of-the-art performance and exclusive use of ViT backbones, making them PE-sensitive. We excluded other SOTA models like Mask2Former [4] due to extra attention modules and custom loss, and SAM-2 [16] for its prompt encoder with cross-attention. All experiments were conducted with a

Model	PE	Acc	mIoU	mDice	Boundary F1	MCC
SETR [38]	No PE	80.84%	79.56%	84.53%	3.48%	88.75%
	1D Sinusoid [25]	80.96%	78.33%	82.51%	3.55%	91.99%
	Learnable [25]	81.35%	74.41%	79.26%	3.88%	88.87%
	CPE [5]	80.67%	78.42%	82.71%	3.31%	87.26%
	LFF(Fourier) [18]	82.69%	80.41%	85.20%	7.83%	85.64%
Segmenter [28]	2D Sinusoid [39]	80.74%	78.26%	82.46%	3.04%	91.88%
	LOOPE	83.63%	81.37%	86.17%	7.44%	89.56%
	No PE	79.55%	79.13%	84.48%	6.23%	83.34%
	1D Sinusoid [25]	81.82%	80.59%	85.79%	18.37%	84.74%
	Learnable [25]	81.48%	80.46%	85.73%	15.33%	84.68%
	CPE [5]	79.19%	78.71%	84.00%	8.80%	82.89%
	LFF(Fourier) [18]	83.20%	81.41%	86.46%	26.60%	85.51%
	2D Sinusoid [39]	80.93%	79.90%	85.13%	9.58%	84.08%
	LOOPE	84.05%	82.26%	87.31%	32.21%	86.37%

Table 3. Performance comparison of PEs with SETR and Segmenter models on the Cityscapes dataset.

batch size of 64 and trained for 100 epochs. We used a learning rate of 0.0001 and an 80-20 train-validation split.

From Table 3, we observe that LOOPE outperforms all PEs by learning patch order, with LFF (Fourier) closely behind — surpassing LOOPE only once in Boundary F1 for SETR — and 1D Sinusoid outperforming LOOPE in MCC.

4.4. Three-Cell Experiment: Positional Probing Task

Standard vision benchmarks often show only modest gains from positional encodings (PEs), because natural images exhibit strong correlations between neighboring patches and transformers can infer rough spatial structure from content alone. This makes it difficult to isolate how much positional information different encoders truly preserve, and in particular to compare absolute PEs (APEs) with task-oriented relative PEs (RPEs). To obtain a more controlled view, we design a synthetic probing task where labels depend purely on geometry and cannot be solved from appearance.

Setup. Each 224×224 image is divided into a 14×14 grid of 16×16 patches. Three patches are randomly selected and colored red, green, and blue at coordinates (x_r, y_r) , (x_g, y_g) , (x_b, y_b) , while all other patches contain color black. The model receives only these images and must answer geometric questions about the positions of the three colored cells. We consider four query types (Fig. 4) and cast them jointly as a 6-way classification problem. Full formulas and the sampling procedure are provided in *Supp. A.2*.

Geometric queries. The four query types are designed to probe complementary aspects of positional structure: (i) *distance comparison* between d_{RG} and d_{RB} , testing undirected monotonicity; (ii) *orientation* (clockwise vs. counterclockwise) of the RGB triangle, testing signed relational structure; (iii) *shadow area comparison* under RG vs. RB,

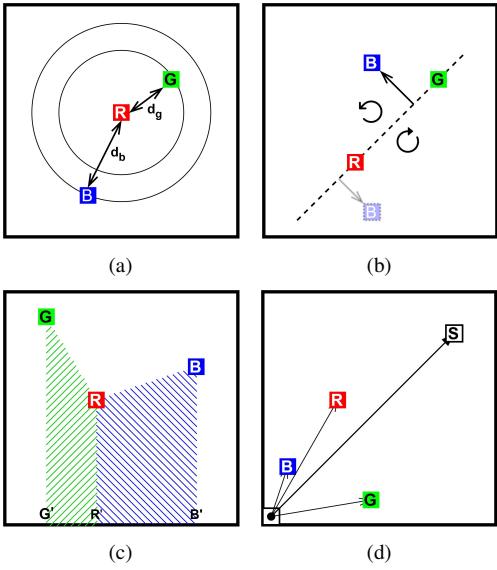


Figure 4. Three-Cell probing task. Given three colored patches (R,G,B), the model must: (a) decide whether d_{RG} or d_{RB} is larger; (b) predict the orientation of triangle RGB (clockwise vs. counter-clockwise); (c) compare “shadow areas” under lines RG and RB; (d) decide whether the sum vector $OR + OG + OB$ lies outside the grid. Labels depend purely on the coordinates of R, G, and B.

which requires combining offsets via addition; and (iv) *vector sum*, predicting whether $\mathbf{p}_r + \mathbf{p}_g + \mathbf{p}_b$ leaves the 14×14 grid, explicitly testing additive composition of coordinates. While the first two can in principle be derived from relative differences, the latter two are fundamentally difficult to reconstruct from purely RPE-based formulations.

Models	Relative Spatial Reasoning		Absolute Spatial Reasoning		
	Distance	Orientation	Area	Vector Sum	Average
ResNet-50[11]	90.8%	85.6%	89.7%	92.9%	89.8%
Inception-V3[30]	92.8%	96.1%	93.7%	95.3%	94.5%
No PE	61.9%	53.1%	60.1%	62.6%	59.4%
Learnable [1]	85.6%	89.3%	84.0%	92.2%	87.8%
1D Sinusoid [1]	90.9%	94.9%	91.3%	95.3%	93.1%
CPE [5]	72.8%	63.2%	72.6%	75.0%	70.9%
RPE [36]	73.7%	62.6%	71.0%	73.9%	70.3%
LFF(Fourier) [18]	93.1%	96.7%	92.4%	93.9%	94.0%
2D Sinusoid [18]	83.6%	60.0%	72.6%	92.3%	77.1%
Static (X_G) [13]	88.8%	95.0%	89.9%	93.8%	91.9%
LOOPE ($X_G + X_c$)	93.4%	95.8%	93.3%	94.6%	94.3%

Table 4. Accuracy comparison of baseline CNN models and PEs with DeiT-Base model on the Three-Cell Experiment dataset.

Findings. Table 4 shows three consistent trends. First, RPEs clearly outperform the no-PE baseline but remain close to chance on the area and vector-sum tasks, confirming that they struggle to encode additive coordinate structure. Second, APEs (sinusoidal, LFF) achieve much higher accuracy across all queries, indicating that absolute coordinates remain crucial even in modern ViTs. **Third, LOOPE**

yields the highest geometric fidelity, suggesting that learning a content-adaptive patch ordering further improves how absolute PEs expose positional information to the transformer.

Overall, this probing task supports our central claim: rather than relying solely on task-oriented RPEs, developing stronger APEs—and in particular learnable patch orderings such as LOOPE—remains a promising direction for positional encoding research.

4.5. Positional Embedding Structural Integrity (PESI) Metrics

PE	Undirected Monotonicity $M_U \uparrow$	Directed Monotonicity $M_D \uparrow$	Undirected Asymmetry A_{SU}
Learnable [1]	1.7493	1.2003	-0.7272
1D Sinusoid [1]	1.9567	1.4905	0.1243
LFF(Fourier) [18]	1.9623	1.5230	0.2683
Static (X_G) [13]	1.9670	1.2897	0.0945
LOOPE ($X_G + X_c$)	1.9674	1.2900	0.0939

Table 5. Comparison of PEs in terms of Undirected Monotonicity, Directed Monotonicity, and Undirected Asymmetry.

Table 5 presents the positional fidelity indices for various APEs. For calculating directed monotonicity, the number of buckets, N is set to 60 testing. So, the $\delta = 6^\circ$. The results indicate that LOOPE achieves the highest values in both **Undirected Monotonicity** and **Undirected Asymmetry**, demonstrating its robustness. Conversely, Learnable APE performs the worst across all three metrics, indicating that its embeddings are not highly monotone. A notable observation is the asymmetry value of Learnable, which is -0.72. This negative value arises because the average cosine similarity across all cells is predominantly negative, leading to an overall asymmetry value below zero. Meanwhile, Fourier exhibits strong **Directed Monotonicity** with stable results in the undirected setting. However, it compromises radial symmetry, meaning that values on a single radius show greater instability compared to other periodic APEs. In contrast, LOOPE demonstrates the most stable radial symmetry, reinforcing its reliability in positional encoding.

4.6. Ablation Studies

4.6.1. Robustness across varying Frequency Set

The results in Table 6 demonstrate that LOOPE consistently outperforms other PEs, even when using a non-optimal frequency set. Notably, there is no theoretically proven optimal frequency configuration. For instance, a simple geometric sequence achieves better accuracy than the original frequency set across all PEs. **This analysis clearly highlights the robustness of**

Freq.	PE	$\omega(i)$	1D Sinusoid	Static (X_G)	LOOPE ($X_G + X_c$)
	Original	0.978^i	85.3%	84.6%	88.1%
	Arithmetic	$1 - \frac{i(1-\lambda)}{L-1}$	73.6%	76.5%	86.6%
	Geometric	$r^i, r = 0.9$	89.8%	87.3%	89.9%
	Random	Uniform ($\lambda, 1$)	83.5%	87.5%	88.0%
	Sensitivity to Freq. ↓		0.027760	0.022516	0.005244

Table 6. Robustness of LOOPE across different frequency sequences ($\lambda = 0.0001$, $L = 768$, $i \in 1, \dots, L$) with ViT-Base as backbone on Oxford-IIIT Dataset. For Sensitivity computation, refer to *Supp.A.6*.

460 **LOOPE—particularly in the case of the arithmetic se-**
 461 **quence, where accuracy drops to 73.6%, yet LOOPE**
 462 **alone helps maintain it at 86.6%.** This experiment indicates
 463 that LOOPE effectively mitigates frequency selection bias.
 464

465 4.6.2. Impact of LOOPE on Variable Resolution

Resolution	Models	1D Sinusoid	LOOPE (Ours)
224 × 224	ViT-Base [10]	85.3%	88.1% (+2.8%)
	ViT-Small [10]	81.6%	83.8% (+2.2%)
	DeiT-Base [32]	86.3%	89.8% (+3.5%)
384 × 384	ViT-Base [10]	89.1%	92.2% (+3.1%)
	ViT-Small [10]	83.0%	86.1% (+3.1%)
	DeiT-Base [32]	88.5%	92.4% (+3.9%)

Table 7. Accuracy Comparison of different ViT models with Sinusoid and LOOPE PEs for different Image Resolutions on Oxford-IIIT. For both cases, patch shape is 16×16 .

466 Table 7 presents the performance of different ViT mod-
 467 els with Sinusoid and LOOPE PEs across two different res-
 468 olutions. Our method shows higher improvement in accu-
 469 racy with bigger resolution.

470 4.6.3. Visualization of Positional Encodings

471 Figure 5 clearly shows that LOOPE is generating more
 472 robust cosine similarity map for all positions. Due to tra-
 473 ditional zigzag order in sinusoidal APE, the boundary and
 474 corner cells have inconsistent similarity pattern, as the or-
 475 der propagates from right to all way back to left boundary.
 476 For 1D learnable APE, it loosely generates map with a lots
 477 of anomaly in near to far distance. Also, other than, central
 478 cells, edge cells are having non-monotone similarity map.
 479 More visualization can be found in (*Supp. B*)

480 4.6.4. Trends in PESI metrics

481 Figure 6, shows the interesting facts about directed mono-
 482 tonicity, with this tool one can investigate how precisely the
 483 positional encodings can maintain monotonicity. Clearly,
 484 increasing precision all positional encoders struggles to pro-
 485 vide a monotone trend in cosine similarity. As $N \rightarrow 1$, $M_D \rightarrow M_U$ which is exactly what we expected. At

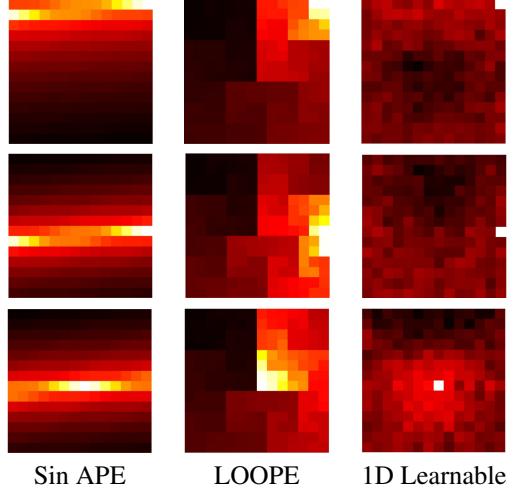


Figure 5. Cosine Similarity Maps for Three APEs: Top-Right Corner, Right-Boundary, and Middle Cell (Top to Bottom)

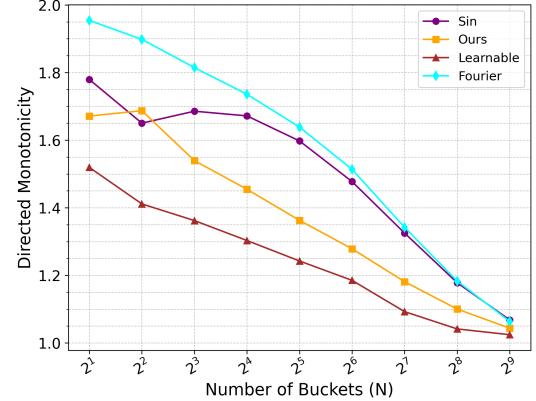


Figure 6. Trend of directed Monotonicity, M_D with increasing angle precision, $\delta = 2\pi/N$

$N = 4, \delta = \pi/2$, we can see that, LOOPE outperforms zigzag and learnable as it highly depends on hilbert order which propagates in square pattern.

5. Conclusion

We introduced **LOOPE**, a learnable framework for patch ordering that unifies geometric structure and context-aware adaptation, establishing reordering as a key design dimension with preserved geometry and frequency robustness. In addition, we proposed the **Three-Cell experiment** and **PESI metrics** as principled tools for evaluating PEs beyond downstream accuracy. This perspective opens new directions for rethinking positional information in transformers.

499 References

- 500 [1] Vaswani Ashish. Attention is all you need. *Advances in
501 neural information processing systems*, 30:I, 2017. 1, 2, 7
- 502 [2] Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda.
503 Crossvit: Cross-attention multi-scale vision transformer for
504 image classification. In *Proceedings of the IEEE/CVF international
505 conference on computer vision*, pages 357–366,
506 2021. 2, 6
- 507 [3] Wanli Chen, Xufeng Yao, Xinyun Zhang, and Bei Yu. Efficient
508 deep space filling curve. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 17525–17534, 2023. 3
- 509 [4] Bowen Cheng, Alexander Schwing, and Alexander Kirillov.
510 Masked-attention mask transformer for universal image seg-
511 mentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
512 1290–1299, 2022. 2, 6
- 513 [5] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xi-
514 aolin Wei, Huaxia Xia, and Chunhua Shen. Conditional pos-
515 itional encodings for vision transformers. In *International Conference on Learning Representations (ICLR)*, 2023. 1, 2,
516 6, 7
- 517 [6] Xiangxiang Chu, Jianlin Su, Bo Zhang, and Chunhua Shen.
518 Visionllama: A unified LLaMA backbone for vision tasks.
519 In *European Conference on Computer Vision (ECCV)*, 2024.
520 2, 6
- 521 [7] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo
522 Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe
523 Franke, Stefan Roth, and Bernt Schiele. The cityscapes
524 dataset for semantic urban scene understanding. In *Proc.
525 of the IEEE Conference on Computer Vision and Pattern
526 Recognition (CVPR)*, 2016. 6
- 527 [8] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell,
528 Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl:
529 Attentive language models beyond a fixed-length context.
530 *arXiv preprint arXiv:1901.02860*, 2019. 2
- 531 [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li,
532 and Li Fei-Fei. Imagenet: A large-scale hierarchical image
533 database. In *Proceedings of the IEEE Conference on Com-
534 puter Vision and Pattern Recognition (CVPR)*, pages 248–
535 255, 2009. 5
- 536 [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,
537 Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
538 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl-
539 vain Gelly, et al. An image is worth 16x16 words: Trans-
540 formers for image recognition at scale. *arXiv preprint
541 arXiv:2010.11929*, 2020. 2, 6, 8
- 542 [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
543 Deep residual learning for image recognition, 2015. 2, 7
- 544 [12] Byeongho Heo, Song Park, Dongyoon Han, and Sangdoo
545 Yun. Rotary position embedding for vision transformer,
546 2024. 2
- 547 [13] David Hilbert. Über die stetige abbildung einer linie auf ein
548 flächenstück. *Mathematische Annalen*, 38:459–460, 1891. 2,
549 6, 7
- 550 [14] Yun-Ning Hung, Ting-Yun Chang, Hsin-Ying Lee, Min Sun,
551 and Sylvain Paris. Learnable fourier features for multi-
552 dimensional spatial positional encoding, 2023. 2
- 553 [15] G Ke, D He, and TY Liu. Rethinking positional en-
554 coding in language pre-training. arxiv. *arXiv preprint
555 arXiv:2006.15595*, 2021. 2
- 556 [16] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi
557 Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer
558 Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár,
559 and Ross Girshick. Segment anything. *arXiv preprint
560 arXiv:2304.02643*, 2023. 6
- 561 [17] Alex Krizhevsky. Learning multiple layers of features from
562 tiny images. Technical report, University of Toronto, 2009.
563 5
- 564 [18] Yang Li, Si Si, Gang Li, Cho-Jui Hsieh, and Samy Bengio.
565 Learnable fourier features for multi-dimensional spatial po-
566 sitional encoding. *Advances in Neural Information Process-
567 ing Systems*, 34:15816–15829, 2021. 6, 7
- 568 [19] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie,
569 Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al.
570 Swin transformer v2: Scaling up capacity and resolution. In
571 *Proceedings of the IEEE/CVF conference on computer vi-
572 sion and pattern recognition*, pages 12009–12019, 2022. 2
- 573 [20] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and
574 C. V Jawahar. Cats and dogs. *Proceedings of the IEEE
575 Conference on Computer Vision and Pattern Recognition
576 (CVPR)*, 2012. 5
- 577 [21] Giuseppe Peano. Sur une courbe, qui remplit toute une aire
578 plane. *Mathematische Annalen*, 36(1):157–160, 1890. 2
- 579 [22] Ali Rahimi and Benjamin Recht. Random features for large-
580 scale kernel machines. *Advances in neural information pro-
581 cessing systems*, 20, 2007. 1
- 582 [23] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan
583 Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-
584 attention in vision models. *Advances in neural information
585 processing systems*, 32, 2019. 2
- 586 [24] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vi-
587 sion transformers for dense prediction. In *Proceedings of
588 the IEEE/CVF International Conference on Computer Vision
589 (ICCV)*, pages 12179–12188, 2021. 2
- 590 [25] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-
591 attention with relative position representations. *arXiv
592 preprint arXiv:1803.02155*, 2018. 1, 6
- 593 [26] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-
594 attention with relative position representations, 2018. 2
- 595 [27] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon
596 Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck
597 transformers for visual recognition. In *Proceedings of
598 the IEEE/CVF conference on computer vision and pattern
599 recognition*, pages 16519–16529, 2021. 2
- 600 [28] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia
601 Schmid. Segmenter: Transformer for semantic segmenta-
602 tion, 2021. 2, 6
- 603 [29] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen
604 Bo, and Yunfeng Liu. Roformer: Enhanced transformer with
605 rotary position embedding. *Neurocomputing*, 568:127063,
606 2024. 2

- 611 [30] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe,
612 Jonathon Shlens, and Zbigniew Wojna. Rethinking the in-
613ception architecture for computer vision, 2015. 7
- 614 [31] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara
615 Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ra-
616 mamoorthi, Jonathan Barron, and Ren Ng. Fourier features
617 let networks learn high frequency functions in low dimen-
618 sional domains. *Advances in neural information processing
619 systems*, 33:7537–7547, 2020. 1
- 620 [32] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco
621 Massa, Alexandre Sablayrolles, and Hervé Jégou. Training
622 data-efficient image transformers & distillation through at-
623 tention. In *International conference on machine learning*,
624 pages 10347–10357. PMLR, 2021. 2, 6, 8
- 625 [33] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles,
626 Gabriel Synnaeve, and Hervé Jégou. Going deeper with im-
627 age transformers. In *Proceedings of the IEEE/CVF interna-*
628 *tional conference on computer vision*, pages 32–42, 2021. 6
- 629 [34] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam,
630 Alan Yuille, and Liang-Chieh Chen. Axial-deeplab: Stand-
631 alone axial-attention for panoptic segmentation. In *European
632 conference on computer vision*, pages 108–126. Springer,
633 2020. 2
- 634 [35] Hanyu Wang, Kamal Gupta, Larry Davis, and Abhinav Shri-
635 vastava. Neural space-filling curves, 2022. 3
- 636 [36] Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and
637 Hongyang Chao. Rethinking and improving relative posi-
638 tion encoding for vision transformer. *CoRR*, abs/2107.14222,
639 2021. 2, 6, 7
- 640 [37] Rui Xu, Xintao Wang, Kai Chen, Bolei Zhou, and
641 Chen Change Loy. Positional encoding as spatial inductive
642 bias in gans. In *Proceedings of the IEEE/CVF Conference
643 on Computer Vision and Pattern Recognition*, pages 13569–
644 13578, 2021. 1
- 645 [38] Sixiao Zheng, Jiacheng Li, Yuandong Tian, Hengyu Cai,
646 Ting Liu, Zhenguo Li, and Nanning Zheng. Rethinking se-
647 mantic segmentation from a sequence-to-sequence perspec-
648 tive with transformers, 2020. 2, 6
- 649 [39] Shengcao Zhou, Ziwei Liu, Xiaohang Zhan, Chen Change
650 Loy, and Ping Luo. When vision transformers outperform
651 resnets without pretraining or strong data augmentations,
652 2022. 6
- 653 [40] Jakub Červený. gilbert: Space-filling curve for rectangu-
654 lar domains of arbitrary size. [https://github.com/
655 jakubcerveny/gilbert](https://github.com/jakubcerveny/gilbert). 3