GEOPE: A Unified Geometric Positional Embedding for Structured Tensors

Anonymous authors

Paper under double-blind review

ABSTRACT

Rotary Positional Embedding (RoPE) excels at encoding relative positions in 1D sequences, but its generalization to higher-dimensional structured data like images and videos remains a challenge. Existing approaches often treat spatial axes independently or combine them heuristically, failing to capture their geometric coupling in a symmetric and consistent manner. To address this, we introduce Geometric Positional Embedding (GeoPE), a framework that extends rotations to 3D Euclidean space using quaternions. To overcome the non-commutativity of quaternion multiplication and ensure symmetry, GeoPE constructs a unified rotational operator by computing the geometric mean of rotations within the corresponding Lie algebra. We also propose a linear variant that preserves the strict relative positional encoding of 1D RoPE, offering superior extrapolation. Extensive experiments on image classification, object detection, and 3D semantic segmentation demonstrate that GeoPE consistently outperforms standard baselines and existing 2D RoPE variants, while retaining the strong extrapolation properties of its 1D predecessor.

1 Introduction

Transformer (Vaswani et al., 2017) has emerged as the backbone of large language models due to its capacity to capture global dependencies, enable parallel computation, and generalize across modalities. However, Transformer lacks an inherent mechanism for sequence order and thus relies on additional positional encodings (Devlin et al., 2019; Raffel et al., 2020; Shaw et al., 2018).

Conventional positional encodings include Absolute Positional Encodings (APE) (Devlin et al., 2019; Dosovitskiy et al., 2020; Chen et al., 2021) and Relative Positional Encodings (RPE) (Liu et al., 2021; Raffel et al., 2020; Park et al., 2022; Ke et al., 2020; Wu et al., 2021). APE injects absolute position indices into token embeddings but cannot capture relative distances, while RPE encodes relative positions at the cost of added complexity. Rotary Positional Encoding (RoPE) (Su et al., 2024) overcomes these limitations by rotating query and key vectors in a 2D plane according to their positions, providing attention with strong length generalization (Jiang et al., 2023; Roziere et al., 2023; Touvron et al., 2023; Yao, 2024).

With Transformer increasingly applied to vision tasks, researchers have explored extending RoPE to two dimensions (Fang et al., 2024; Lu et al., 2024a;b), such as Vision Transformer (ViT) (Dosovitskiy et al., 2020) and Swin Transformer (Liu et al., 2021). Existing methods often adopt axis-wise or mixed-frequency designs, processing horizontal and vertical encodings separately or in combination (Chu et al., 2024). For instance, Heo et al. (2024) partitions the embedding space to allow independent or mixed-frequency rotations per axis, enabling both axis-aligned and diagonal interactions. Nevertheless, these approaches remain essentially 1D ROPE, as axes are treated independently, and mixed-frequency schemes only partially capture diagonal dependencies, leaving the weak cross-axis coupling of high-dimensional RoPEs unresolved.

The challenge is amplified in multi-modal learning, where structured data such as images, videos, and audio demand flexible positional encodings (Dao et al., 2024; Yin et al., 2025; Shu et al., 2023). Some works extend RoPE to higher dimensions via Lie group/algebra frameworks. For example, Liu & Zhou (2025) formalizes RoPE using a maximal abelian subalgebra (MASA) and introduces cross-dimensional interactions through orthogonal basis changes, parameterized by Householder or Cayley methods to balance local and global effects. However, this can overly constrain or globally

entangle representations, incurring a high computational cost. Alternatively, Ostmeier et al. learn dense skew-symmetric matrices to build rotation operators, enabling block-wise interaction control, yet lacking theoretical guarantees and remaining computationally expensive.

We propose Geometric Positional Embedding (GeoPE), which extends RoPE's 2D complex-plane rotations to 3D Euclidean space using quaternions, enabling the modeling of coupled rotations in high-dimensional data in Section 3.1 and Section 3.3. To overcome the non-commutativity of quaternion multiplication, we perform a symmetric average in the logarithmic tangent space in Section 3.2 and also introduce a linear variant for direct relative encoding in Section 3.4. This method enriches self-attention with a geometrically meaningful understanding of space, thereby fostering superior spatial reasoning, as discussed in Section 4. Experiments in Section 5 show that GeoPE achieves significant performance gains in image classification, object detection, and 3D semantic segmentation, while also retaining the strong extrapolation properties of its 1D predecessor.

2 RELATED WORK

Position Encodings. Transformers lack inherent positional awareness and thus rely on encodings to capture the order of tokens. The original Transformer (Vaswani et al., 2017) employs sinusoidal absolute positional encodings (APE), which generalize poorly to long sequences. In contrast, learnable APE (Shaw et al., 2018) improves flexibility and representation for tasks such as sentence alignment and context modeling. Vision Transformers (ViT) (Dosovitskiy et al., 2020) similarly adopt learnable APE for image patches. Relative positional encodings (RPE) model pairwise token distances, supporting long sequences and cross-sequence dependencies (Liu et al., 2021; Shaw et al., 2018), though naive designs incur quadratic cost. Rotary Positional Encoding (RoPE) (Su et al., 2024) encodes relative positions via complex-plane rotations and is widely used in large language models; however, its performance degrades when extrapolated to much longer contexts. More recent approaches learn semanticized position structures. Contextual positional encodings (CoPE) (Golovneva et al., 2024) enhance reasoning and mathematical capabilities. Abacus embeddings (McLeish et al., 2024) capture numerical structures for arithmetic, while lightweight methods, such as LaPE (Yu et al., 2023), apply adaptive normalization to improve robustness across architectures.

RoPE in Visual Model. RoPE has demonstrated strong extrapolation capabilities in long-text modeling and dialogue, motivating its extension to vision and multimodal tasks (Lu et al., 2024b; Wang et al., 2024; Yao et al., 2024). A straightforward adaptation applies 1D RoPE to ViT variants, as in Hybrid X-former (Jeevan & Sethi, 2022). However, gains are modest and have been validated only on small datasets (e.g., CIFAR, Tiny ImageNet). To better handle 2D inputs, works such as EVA-02 (Fang et al., 2024) and Unified-IO 2 (Lu et al., 2024a) have incorporated axial 2D RoPE into multimodal and diffusion models; however, these fail to capture diagonal interactions. RoPE for ViT (Heo et al., 2024) further proposed RoPE-Mixed, which combines axial frequencies to enhance 2D encodings and downstream performance. However, this approach remains essentially frequency composition, offering only loose dimensional coupling and limited generality.

Shape Bias. Cognitive science has shown that humans rely primarily on global shape, rather than texture or color, for object recognition and lexical learning, whereas CNNs exhibit a different tendency. Hosseini et al. (2018) demonstrated that standard CNNs often lack shape bias, instead depending heavily on local texture or color cues. However, Ritter et al. (2017) reported that networks can develop shape preference under certain conditions. To examine this systematically, Geirhos et al. (2018) compared CNNs and humans using style-transferred images with conflicting shape and texture information. While humans consistently prioritize shape, CNNs tend to favor texture. To mitigate this bias, they introduced Stylized-ImageNet, which reduced texture reliance and induced stronger shape bias, yielding models with improved robustness and transferability. These findings suggest that enhancing shape bias can make models more human-like while also strengthening generalization.

3 METHODOLOGY

In this section, we detail the formulation and implementation of GeoPE. We first establish the geometric requirements for multi-axial rotation in Section 3.1, then construct a symmetric rotational

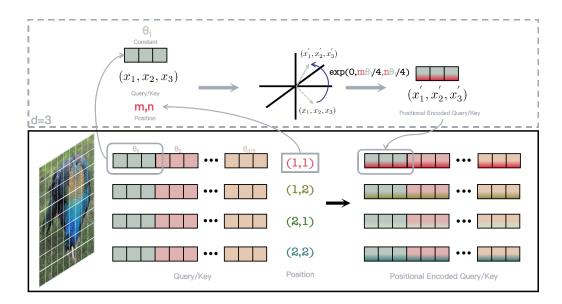


Figure 1: Geometric Transform of Geometric Position Embedding(GeoPE).

operator using Lie theory in Section 3.2. Finally, we demonstrate the framework's extension to 3D in Section 3.3 and propose a linear variant in Section 3.4.

3.1 GENERALIZING ROTATIONS TO 3D SPACE

While RoPE's complex-plane rotations are effective for 1D data, they are insufficient for higher-dimensional data like images, which require coupled multi-axis rotations. We therefore extend the rotational domain to 3D Euclidean space as illustrated in Figure 1, using quaternions to formulate these transformations robustly and avoid issues such as gimbal lock.

Mathematically, a feature vector $\mathbf{x} \in \mathbb{R}^d$ is first partitioned into d/3 sub-vectors, $\{\mathbf{v}_i\}_{i=1}^{d/3}$, where each $\mathbf{v}_i = (v_x, v_y, v_z) \in \mathbb{R}^3$. Each sub-vector \mathbf{v}_i is then "lifted" into the quaternion space \mathbb{H} as a pure quaternion (i.e., a quaternion with a zero scalar part):

$$\mathbf{p} = 0 + v_x \mathbf{i} + v_y \mathbf{j} + v_z \mathbf{k} \tag{1}$$

Given a unit quaternion \mathbf{r} that represents a desired rotation, the transformation of \mathbf{p} is given by the sandwich product:

$$\mathbf{p}' = \mathbf{r}\mathbf{p}\mathbf{r}^* \tag{2}$$

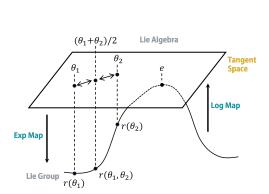
where \mathbf{r}^* is the conjugate of \mathbf{r} , which for a unit quaternion is equivalent to its inverse (\mathbf{r}^{-1}). A crucial property of this operation is that the result \mathbf{p}' remains a pure quaternion. Its vector part corresponds to the rotated vector \mathbf{v}'_i in \mathbb{R}^3 . This rotational operation is, by construction, an isometry for each 3D sub-vector, preserving its norm $\|\mathbf{v}_i\|$.

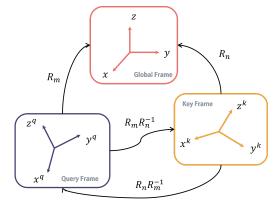
The rotational quaternion ${\bf r}$ is a function of positional indices, e.g., (h,w) for a 2D image, which encode phase information θ_h and θ_w . For a position (p_h,p_w) and a given sub-vector $i\in\{1,\ldots,d/3\}$, these are defined as $\theta_h=p_h\cdot\lambda^{2i/d}$ and $\theta_w=p_w\cdot\lambda^{2i/d}$, where λ is a chosen base which is set as $\lambda=100$ as usual(Heo et al., 2024).

3.2 Constructing a Symmetric Operator

For 2D data, positional information along the height and width dimensions can be encoded as rotations about distinct axes. A natural choice is to associate them with rotations about the y-axis (j) and z-axis (k), respectively. This yields two base quaternions:

$$\mathbf{r}_h(\theta_h) = \cos\left(\frac{\theta_h}{2}\right) + \sin\left(\frac{\theta_h}{2}\right)\mathbf{j}, \quad \mathbf{r}_w(\theta_w) = \cos\left(\frac{\theta_w}{2}\right) + \sin\left(\frac{\theta_w}{2}\right)\mathbf{k}$$





- (a) The Log-Exp average of Lie Algebra and Lie group, where e is identity element.
- (b) The transform of Global Frame and Relative frame.

Figure 2: Illustration of mathematical structure and coordinate transform.

A naive composition of these rotations via quaternion multiplication, such as $\mathbf{r}_{hw} = \mathbf{r}_h \mathbf{r}_w$, is ill-suited for our purpose. Quaternion multiplication is non-commutative ($\mathbf{r}_h \mathbf{r}_w \neq \mathbf{r}_w \mathbf{r}_h$), meaning the resulting rotation would be arbitrarily dependent on the chosen order of operations, creating an undesirable symmetric bias between the height and width encodings.

To construct an operator that treats each spatial dimension symmetrically, we turn to the tools of Lie theory. The core idea is to compute the geometric mean of the rotations. This is achieved by mapping the quaternions from the non-linear Lie group of 3D rotations, SO(3), to its corresponding linear Lie algebra, $\mathfrak{so}(3)$, via the logarithm map. In this tangent vector space, a simple averaging operation is well-defined and commutative. The result is then mapped back to the Lie group via the exponential map.

Accordingly, we define our symmetric rotational operator \mathbf{r} as:

$$\mathbf{r}(\theta_h, \theta_w) = \exp\left(\frac{1}{2}\left(\log(\mathbf{r}_h(\theta_h)) + \log(\mathbf{r}_w(\theta_w))\right)\right)$$
(3)

As derived in Appendix A, the intermediate averaged vector in the Lie algebra $\mathfrak{so}(3)$ is $(0,\theta_h/4,\theta_w/4)$. The exponential map yields an elegant closed-form solution for the resulting quaternion:

$$\mathbf{r} = \cos\left(\frac{\Theta}{2}\right) + \sin\left(\frac{\Theta}{2}\right) \frac{\theta_h}{2\Theta} \mathbf{j} + \sin\left(\frac{\Theta}{2}\right) \frac{\theta_w}{2\Theta} \mathbf{k} \tag{4}$$

where $\Theta = \frac{1}{2}\sqrt{\theta_h^2 + \theta_w^2}$. The coupled phase Θ is proportional to the Euclidean distance between (θ_h, θ_w) and the origin, while the vector components ensure that the influence of each positional phase remains monotonic. As illustrated in Figure 2a, this log-exp average provides a commutative and geometrically sound method for combining rotations.

The quaternion rotation in Equation 2 is equivalent to a matrix-vector product, $\mathbf{v}' = \mathbf{R}\mathbf{v}$, where $\mathbf{R} \in SO(3)$ is the rotation matrix corresponding to \mathbf{r} . The complete transformation on a d-dimensional query vector \mathbf{q} or key vector \mathbf{k} is thus a block-diagonal matrix:

$$\mathbf{R}_{\text{GeoPE}} = \begin{pmatrix} \mathbf{R}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{R}_{d/3} \end{pmatrix}, \mathbf{R}_i = \begin{pmatrix} \cos(\Theta) & -\frac{\theta_w \sin(\Theta)}{\sqrt{\theta_h^2 + \theta_w^2}} & \frac{\theta_h \sin(\Theta)}{\sqrt{\theta_h^2 + \theta_w^2}} \\ \frac{\theta_w \sin(\Theta)}{\sqrt{\theta_h^2 + \theta_w^2}} & 1 - \frac{\theta_w^2 (1 - \cos(\Theta))}{\theta_h^2 + \theta_w^2} & \frac{\theta_h \theta_w (1 - \cos(\Theta))}{\theta_h^2 + \theta_w^2} \\ -\frac{\theta_h \sin(\Theta)}{\sqrt{\theta_h^2 + \theta_w^2}} & \frac{\theta_h \theta_w (1 - \cos(\Theta))}{\theta_h^2 + \theta_w^2} & 1 - \frac{\theta_h^2 (1 - \cos(\Theta))}{\theta_h^2 + \theta_w^2} \end{pmatrix}$$

where each \mathbf{R}_i is a 3×3 rotation matrix derived from the quaternion \mathbf{r} computed with phases $(\theta_{h,i},\theta_{w,i})$ specific to that block. When the structured tensor is one-dimensional, GeoPE as disscussed in Appendix D gracefully degenerates to a 2D rotation equivalent to the original RoPE (Su et al., 2024). Meanwhile, GeoPE keep long distance decay as disscussed in Appendix C

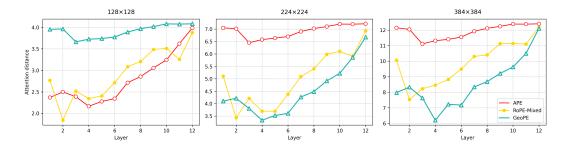


Figure 3: Mean attention distance as a function of layer depth across different input resolutions. While all methods exhibit an expanding receptive field in deeper layers, APE's consistently higher distance suggests an inefficient and unfocused global search. In contrast, GeoPE maintains a more moderate distance, indicating a more structured and efficient strategy for balancing local and global information gathering. These relative trends remain consistent across all tested resolutions.

3.3 EXTENSION TO THREE SPATIAL DIMENSIONS

The GeoPE framework extends naturally to three spatial dimensions (e.g., for video data or volumetric scans) with positions (d, h, w). We introduce a third base quaternion for depth, $\mathbf{r}_d(\theta_d) = \cos(\frac{\theta_d}{2}) + \sin(\frac{\theta_d}{2})\mathbf{i}$, and compute the symmetric average of the three rotations:

$$\mathbf{r}(\theta_d, \theta_h, \theta_w) = \exp\left(\frac{1}{3}\left(\log(\mathbf{r}_d) + \log(\mathbf{r}_h) + \log(\mathbf{r}_w)\right)\right)$$
 (5)

This yields the three-dimensional GeoPE operator by results in Appendix E:

$$\mathbf{r} = \cos\left(\frac{\Theta}{2}\right) + \sin\left(\frac{\Theta}{2}\right) \left(\frac{\theta_d}{3\Theta}\mathbf{i} + \frac{\theta_h}{3\Theta}\mathbf{j} + \frac{\theta_w}{3\Theta}\mathbf{k}\right)$$
(6)

where the composite phase is now $\Theta = \frac{1}{3}\sqrt{\theta_d^2 + \theta_h^2 + \theta_w^2}$. This demonstrates the flexibility and scalability of our proposed geometric approach.

3.4 Linear Formulation for Relative Position Encoding

A critical property of positional embeddings in Transformer architectures is the ability to encode relative position, as the attention mechanism is fundamentally relational. For a query \mathbf{q} at position m and a key \mathbf{k} at position n, the attention score is a function of $\langle \mathbf{R}_m \mathbf{q}, \mathbf{R}_n \mathbf{k} \rangle = \langle \mathbf{q}, \mathbf{R}_m^{\dagger} \mathbf{R}_n \mathbf{k} \rangle$. Ideally, the relative rotation matrix $\mathbf{R}_{m \to n} = \mathbf{R}_m^{\top} \mathbf{R}_n$ should depend only on the displacement n - m.

Our symmetric operator, while geometrically sound, does not inherently guarantee this linear relationship in the parameter space. That is, $\mathbf{r}(\theta_h, \theta_w) \neq \mathbf{r}(\phi_h, \phi_w)\mathbf{r}(\theta_h - \phi_h, \theta_w - \phi_w)$. To recover an inductive bias analogous to the simple phase subtraction in 1D RoPE, we propose a 'Linear GeoPE' formulation. The core insight is to enforce a linear relationship in the Lie algebra, where rotational composition is approximated by vector addition. By defining the relative rotation based on the difference of the Lie algebra vectors, i.e., $\mathbf{u}_{\text{rel}} = \mathbf{u}_k - \mathbf{u}_q$, we ensure the resulting rotation depends on the simple linear displacement of positional phases, mirroring the behavior of the original RoPE.

Let the Lie algebra vectors for a query at position (h_q, w_q) and a key at position (h_k, w_k) be $\mathbf{u}_q = (0, \theta_{h_q}/4, \theta_{w_q}/4)$ and $\mathbf{u}_k = (0, \theta_{h_k}/4, \theta_{w_k}/4)$, respectively. We define the relative Lie algebra vector as their difference:

$$\mathbf{u}_{\text{rel}} = \mathbf{u}_k - \mathbf{u}_q = \left(0, \frac{\theta_{h_k} - \theta_{h_q}}{4}, \frac{\theta_{w_k} - \theta_{w_q}}{4}\right) \tag{7}$$

The relative rotation is then obtained by mapping this difference back to the Lie group: $\mathbf{r}_{rel} = \exp(\mathbf{u}_{rel})$. This construction ensures that the transformation between any two positions depends solely on their relative displacement.

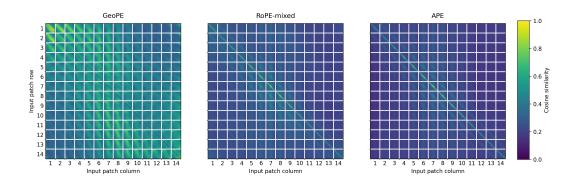


Figure 4: Comparison of attention patterns from the final layer of models using different positional embeddings. The heatmaps visualize the cosine similarity between patch representations, averaged across all attention heads. APE results in highly localized attention focused on the diagonal. RoPE-mixed shows a more distributed local pattern. In contrast, GeoPE facilitates complex, long-range attention, indicating a significantly more global receptive field.

This allows the attention score to be computed as $\langle \mathbf{q}, \mathbf{R}_{rel} \mathbf{k} \rangle$. However, unlike the 1D case where the relative rotation matrix is a simple 2D rotation, the 3×3 matrix \mathbf{R}_{rel} is generally dense. Applying this transformation explicitly is computationally more demanding than the standard GeoPE formulation, presenting a trade-off between enforcing a strict linear inductive bias and computational efficiency.

4 DISCUSSION

In this section, we further explore the properties of GeoPE to provide a deeper understanding of its mechanism and impact. We analyze the geometric interpretation of the attention score under 3D rotations and discuss how GeoPE influences the model's spatial reasoning capabilities.

4.1 Geometric Interpretation of the GeoPE

GeoPE enriches the self-attention mechanism by incorporating a geometrically meaningful understanding of space. The attention score between a query \mathbf{q} at position $m=(h_m,w_m)$ and a key \mathbf{k} at position $n=(h_n,w_n)$ is computed on their rotated counterparts:

$$AttnScore(\mathbf{q}_m, \mathbf{k}_n) = \langle \mathbf{R}_m \mathbf{q}, \mathbf{R}_n \mathbf{k} \rangle = \langle \mathbf{q}, \mathbf{R}_m^{\top} \mathbf{R}_n \mathbf{k} \rangle$$
 (8)

This formulation offers two powerful, complementary geometric interpretations as shown in Figure 2b.

Global Coordinate Frame. One perspective is that \mathbf{R}_m and \mathbf{R}_n transform the query and key vectors from their local, position-agnostic feature spaces into a shared global coordinate frame defined by their absolute positions. The inner product is then computed in this global frame, allowing for a direct, spatially-aware comparison.

Relative Coordinate Frame. Alternatively, and perhaps more intuitively for attention, the term $\mathbf{R}_{\text{rel}} = \mathbf{R}_m^{\top} \mathbf{R}_n$ can be interpreted as a relative rotation operator. It transforms the key vector \mathbf{k} from its own positional frame at n into the query's positional frame at m. The attention score is thus a measure of feature similarity after aligning the key to the query's geometric context.

Unlike the simple phase difference in Heo et al. (2024), this 3D relative rotation depends not only on the magnitude of the displacement $(h_n - h_m, w_n - w_m)$ but also on the direction of displacement. The attention score is governed by the inner product of a vector with its rotated version, which, according to Rodrigues' rotation formula, is a function of both the angle and the axis of this relative rotation. For a rotation of angle A about an axis \mathbf{n} , the inner product as discussed in Appendix B becomes:

$$\langle \mathbf{q}, \mathbf{R}_{\text{rel}} \mathbf{k} \rangle = \underbrace{\langle \mathbf{q}, \mathbf{k} \rangle \cos(A)}_{\text{Projected Similarity}} + \underbrace{(\mathbf{q} \cdot \mathbf{n})(\mathbf{k} \cdot \mathbf{n})(1 - \cos(A))}_{\text{Axial Alignment}} - \underbrace{(\mathbf{n} \times \mathbf{q}) \cdot \mathbf{k} \sin(A)}_{\text{Torsional Component}}$$
(9)

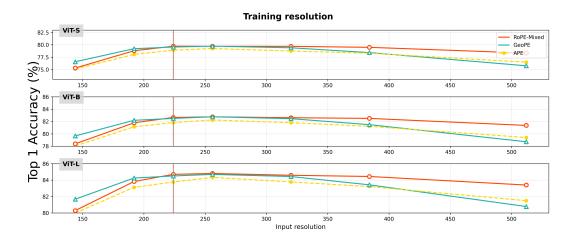


Figure 5: Generalization performance to unseen input resolutions for ViT-S, -B, and -L models. All models are trained at a fixed 224x224 resolution (marked by the vertical line) and evaluated on a range of different resolutions. Absolute Positional Embedding (APE) fails to generalize, with its accuracy collapsing at higher resolutions. In contrast, relative embeddings like RoPE-Mixed and GeoPE show strong robustness as their performance degrades gracefully, highlighting their suitability for real-world applications with variable input sizes.

This decomposition provides a clear geometric intuition. The Projected Similarity term generalizes RoPE by modulating similarity based on displacement magnitude (via angle A). Crucially, the Axial Alignment term—which has no 2D equivalent—adds sensitivity to the displacement direction (via axis n), while the Torsional Component captures the features' relative rotational orientation. Thus, the attention score moves beyond simple distance to depend on the full geometric relationship between tokens. For Linear GeoPE, the angle A_i and axis n_i for the i-th sub-vector are defined as:

$$A_{i} = \frac{1}{2} \sqrt{(\Delta \theta_{h,i})^{2} + (\Delta \theta_{w,i})^{2}}, \ \mathbf{n}_{i} = \frac{\frac{\Delta \theta_{h,i}}{4} \mathbf{j} + \frac{\Delta \theta_{w,i}}{4} \mathbf{k}}{\frac{1}{4} \sqrt{(\Delta \theta_{h,i})^{2} + (\Delta \theta_{w,i})^{2}}} = \frac{\Delta \theta_{h,i} \mathbf{j} + \Delta \theta_{w,i} \mathbf{k}}{\sqrt{(\Delta \theta_{h,i})^{2} + (\Delta \theta_{w,i})^{2}}}$$

where $\Delta\theta_{h,i} = \theta_{h_k,i} - \theta_{h_q,i} = (p_{h_k} - p_{h_q}) \cdot \lambda^{2i/d} = \Delta p_h \cdot \lambda^{2i/d}$ and $\Delta\theta_{w,i} = \theta_{w_k,i} - \theta_{w_q,i} = (p_{w_k} - p_{w_q}) \cdot \lambda^{2i/d} = \Delta p_w \cdot \lambda^{2i/d}$. This shows that the interaction is a complex blend of the original similarity $\langle \mathbf{q}, \mathbf{k} \rangle$ and terms modulated by the alignment of the vectors with the relative rotation axis, endowing the model with a richer, more expressive spatial bias.

4.2 IMPACT ON ATTENTION PATTERNS AND SPATIAL AWARENESS

We hypothesize that GeoPE's geometric inductive bias fosters more effective spatial reasoning by enabling more meaningful attention patterns. Our analyses support this: models equipped with GeoPE exhibit longer attention distances in Figure 3 and more global attention maps in Figure 4. This behavior allows the model to capture long-range dependencies and integrate information across the entire spatial domain, rather than focusing only on local texture. We posit that this enhanced global awareness directly contributes to the performance gains and improved shape-texture bias observed in our experiments.

5 EXPERIMENTS

We validate our methods, GeoPE and Linear GeoPE, through comprehensive experiments on image classification, object detection, and 3D semantic segmentation, benchmarking them against standard baselines and existing 2D rotational embeddings.

Table 1: ViTs were trained on ImageNet-1K(Deng et al., 2009) following the 400-epoch training protocol of DeiT-III(Touvron et al., 2022) using cross-entropy loss. In Swin Transformers, we substitute the RPB with GeoPE and train each model according to the original 300-epoch protocol(Liu et al., 2021).

Backbone	Train Resolution	PE	Top 1 Acc
ViT-small	192×192	GeoPE	78.5
	192×192	LinGeoPE	78.8
	224×224	APE	79.9
	224×224	CPE	80.7
	224×224	GeoPE	81.2
ViT-base	224×224	APE	81.3
	224×224	CPE	82.2
	224×224	GeoPE	82.5
ViT-large	224×224	APE	83.3
	224×224	CPE	83.6
	224×224	GeoPE	83.9
Swin-S	224×224	PRB	83.0
	224×224	Rope-Mixed	83.4
	224×224	GeoPE	83.5
Swin-B	224×224	PRB	83.5
	224×224	Rope-Mixed	83.8
	224×224	GeoPE	83.6

5.1 IMAGE CLASSIFICATION

We evaluate our methods on the ImageNet-1K classification task using Vision Transformer (ViT) (Dosovitskiy et al., 2020) and Swin Transformer (Liu et al., 2021) backbones, following the DeiT3 training protocol (Armeni et al., 2016b) with CE loss.

As shown in Table 1, GeoPE consistently improves Top-1 accuracy across all backbones. It outperforms standard baselines like APE and CPE (Chu et al., 2021) on ViT models and matches or exceeds the performance of PRB and Rope-Mixed (Heo et al., 2024) on Swin Transformers, demonstrating the broad applicability of its geometric prior. Furthermore, as depicted in Figure 5, Linear GeoPE exhibits exceptional zero-shot inference capabilities across multiple resolutions, confirming its superior extrapolation properties as a natural high-dimensional extension of RoPE (Su et al., 2024).

5.2 OBJECT DETECTION

To assess GeoPE's impact on tasks requiring fine-grained spatial awareness, we evaluate it on the MS-COCO (Lin et al., 2014) object detection benchmark. We integrate GeoPE into the DINO-ViTDet (Zhang et al., 2022) framework, a strong object detection pipeline.

Table 2 shows that GeoPE consistently improves mAP for both ViT-B and ViT-L backbones. Compared with APE and Rope-Mixed (Heo et al., 2024), GeoPE provides the largest relative gains, highlighting the importance of explicit geometric priors in capturing global spatial relationships critical for accurate object detection.

5.3 3D SEMANTIC SEGMENTATION

To verify the hypothesis that GeoPE is suitable for any structured tensor data where spatial relationships are paramount, we apply it to 3D point cloud segmentation on the S3DIS dataset (Armeni et al., 2016b). We incorporate GeoPE into the Point Transformer architecture.

As reported in Table 3, GeoPE improves all major metrics, including overall accuracy, mean class accuracy, and mean IoU, relative to the RPE baseline. These improvements confirm that explicitly

Table 2: This table reports MSCOCO(Lin et al., 2014) detection performance (box AP). DINO(Zhang et al., 2022) is trained under the 12-epoch DINO-ViTDet setting(Ren et al., 2023). For GeoPE, we apply it to the ViT backbone, which is pre-trained on ImageNet-1K using the 400-epoch DeiT-III recipe.

Backbone	PE	mAP	
ViT-base	APE Rope-Mixed GeoPE	49.4 51.2 51.3	
ViT-large	APE Rope-Mixed GeoPE	51.1 52.9 53.1	

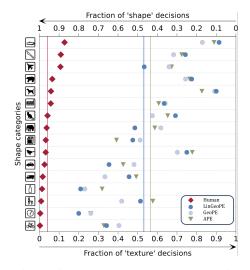


Table 3: Semantic segmentation performance on the S3DIS dataset(Armeni et al., 2016a), evaluated using 6-fold cross-validation.

Backbone	PE	OA	mAcc	mIoU
Point-	RPE	90.2	81.9	73.5
Transformer	GeoPE	90.5	82.1	74.4

Figure 6: Shape Bias Relation Analysis. The plot shows the fraction of 'shape' versus 'texture' decisions for different image categories. Models with GeoPE consistently show a higher shape bias.

encoding multi-axis spatial relationships allows the model to better capture 3D geometric structures, validating the general applicability of GeoPE beyond 2D vision tasks.

5.4 Shape-Texture Bias Analysis

To assess GeoPE's impact on spatial reasoning, we analyze its effect on the model's shape-texture bias. A strong shape bias, which prioritizes object structure over texture, is often correlated with better robustness and generalization. As shown in Figure 6, GeoPE significantly increases the model's shape bias compared to the baseline. This provides compelling evidence that our method encourages learning features more attuned to object geometry, fostering a more robust and holistic visual understanding by striking an effective balance between shape and texture.

6 Conclusion

We propose GeoPE, a method that extends 1D rotational embeddings to higher dimensions by using quaternions for 3D rotations in vector sub-spaces. To overcome the non-commutativity of quaternions, we introduce a symmetric averaging technique based on Lie theory and derive a Linear GeoPE variant that preserves relative position encoding. Extensive experiments show GeoPE significantly enhances spatial awareness in Transformers, boosting performance on 2D and 3D tasks while retaining the excellent extrapolation properties of its 1D predecessor. Our work offers a principled and effective path for position encoding in models for structured data.

REPRODUCIBILITY STATEMENT

We have made extensive efforts to ensure the reproducibility of our work. The main paper provides detailed descriptions of the proposed method, model architectures, and training procedures. Additional experimental details, ablation studies, and theoretical derivations are included in the Appendix. We also provide the complete data preprocessing steps and hyperparameter configurations

in the supplementary material. Furthermore, we submit the anonymized source code and training scripts as supplementary material to facilitate replication of all reported results.

ACKNOWLEDGEMENTS

We used a large language model to improve the grammar and clarity of the paper's text. All research ideas, experiments, and analyses are our own.

REFERENCES

- Iro Armeni, Ozan Sener, Amir R. Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016a.
- Iro Armeni, Ozan Sener, Amir R Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1534–1543, 2016b.
- Pu-Chin Chen, Henry Tsai, Srinadh Bhojanapalli, Hyung Won Chung, Yin-Wen Chang, and Chun-Sung Ferng. A simple and effective positional encoding for transformers. *arXiv preprint arXiv:2104.08698*, 2021.
- Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, and Chunhua Shen. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021.
- Xiangxiang Chu, Jianlin Su, Bo Zhang, and Chunhua Shen. Visionllama: A unified llama interface for vision tasks. *CoRR*, 2024.
- Hong N Dao, Tuyen Nguyen, Cherubin Mugisha, and Incheon Paik. A multimodal transfer learning approach using pubmedclip for medical image classification. *IEEE Access*, 12:75496–75507, 2024.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Yuxin Fang, Quan Sun, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva-02: A visual representation for neon genesis. *Image and Vision Computing*, 149:105171, 2024.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *International conference on learning representations*, 2018.
- Olga Golovneva, Tianlu Wang, Jason Weston, and Sainbayar Sukhbaatar. Contextual position encoding: Learning to count what's important. *arXiv preprint arXiv:2405.18719*, 2024.
- Byeongho Heo, Song Park, Dongyoon Han, and Sangdoo Yun. Rotary position embedding for vision transformer. In *European Conference on Computer Vision*, pp. 289–305. Springer, 2024.
- Hossein Hosseini, Baicen Xiao, Mayoore Jaiswal, and Radha Poovendran. Assessing shape bias property of convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1923–1931, 2018.

- Pranav Jeevan and Amit Sethi. Resource-efficient hybrid x-formers for vision. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 2982–2990, 2022.
 - Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.
 - Guolin Ke, Di He, and Tie-Yan Liu. Rethinking positional encoding in language pre-training. *arXiv* preprint arXiv:2006.15595, 2020.
 - Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pp. 740–755. Springer, 2014.
 - Haiping Liu and Hongpeng Zhou. Rethinking rope: A mathematical blueprint for n-dimensional positional encoding. *arXiv* preprint arXiv:2504.06308, 2025.
 - Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
 - Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26439–26455, 2024a.
 - Zeyu Lu, Zidong Wang, Di Huang, Chengyue Wu, Xihui Liu, Wanli Ouyang, and Lei Bai. Fit: Flexible vision transformer for diffusion model. *arXiv preprint arXiv:2402.12376*, 2024b.
 - Sean McLeish, Arpit Bansal, Alex Stein, Neel Jain, John Kirchenbauer, Brian Bartoldson, Bhavya Kailkhura, Abhinav Bhatele, Jonas Geiping, Avi Schwarzschild, et al. Transformers can do arithmetic with the right embeddings. *Advances in Neural Information Processing Systems*, 37: 108012–108041, 2024.
 - Sophie Ostmeier, Brian Axelrod, Maya Varma, Michael Moseley, Akshay S Chaudhari, and Curtis Langlotz. Liere: Lie rotational positional encodings. In *Forty-second International Conference on Machine Learning*.
 - Wonpyo Park, Woonggi Chang, Donggeon Lee, Juntae Kim, and Seung-won Hwang. Grpe: Relative positional encoding for graph transformer. *arXiv preprint arXiv:2201.12787*, 2022.
 - Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
 - Tianhe Ren, Shilong Liu, Feng Li, Hao Zhang, Ailing Zeng, Jie Yang, Xingyu Liao, Ding Jia, Hongyang Li, He Cao, Jianan Wang, Zhaoyang Zeng, Xianbiao Qi, Yuhui Yuan, Jianwei Yang, and Lei Zhang. detrex: Benchmarking detection transformers, 2023. URL https://arxiv.org/abs/2306.07265.
 - Samuel Ritter, David GT Barrett, Adam Santoro, and Matt M Botvinick. Cognitive psychology for deep neural networks: A shape bias case study. In *International conference on machine learning*, pp. 2940–2949. PMLR, 2017.
 - Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950, 2023.
 - Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.

- Changyong Shu, Jiajun Deng, Fisher Yu, and Yifan Liu. 3dppe: 3d point positional encoding for transformer-based multi-camera 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3580–3589, 2023.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. In *European conference on computer vision*, pp. 516–533. Springer, 2022.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and Hongyang Chao. Rethinking and improving relative position encoding for vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10033–10041, 2021.
- Yupu Yao. Generalize neural network through smooth hypothesis function. In *The Second Tiny Papers Track at ICLR 2024*, 2024. URL https://openreview.net/forum?id=RpiiKMGaK3.
- Yupu Yao, Shangqi Deng, Zihan Cao, Harry Zhang, and Liang-Jian Deng. Apla: Additional perturbation for latent noise with adversarial training enables consistency, 2024. URL https://arxiv.org/abs/2308.12605.
- Hao Yin, Guangzong Si, and Zilei Wang. Clearsight: Visual signal enhancement for object hallucination mitigation in multimodal large language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 14625–14634, 2025.
- Runyi Yu, Zhennan Wang, Yinhuai Wang, Kehan Li, Chang Liu, Haoyi Duan, Xiangyang Ji, and Jie Chen. Lape: Layer-adaptive position embedding for vision transformers with independent layer normalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5886–5896, 2023.
- Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. arXiv preprint arXiv:2203.03605, 2022.

A DERIVATION OF THE SYMMETRIC OPERATOR

- This section details the derivation of the closed-form solution for the symmetric rotational operator $\mathbf{r}(\theta_h, \theta_w)$ introduced in Section 3.2.
- Our goal is to compute the geometric mean of two base rotations, $\mathbf{r}_h(\theta_h)$ and $\mathbf{r}_w(\theta_w)$, using the log-exp map formalism:

$$\mathbf{r}(\theta_h, \theta_w) = \exp\left(\frac{1}{2}\left(\log(\mathbf{r}_h(\theta_h)) + \log(\mathbf{r}_w(\theta_w))\right)\right)$$
(10)

The logarithm map for a unit quaternion $\mathbf{r} = \cos(\alpha) + \sin(\alpha)\mathbf{n}$, where \mathbf{n} is a unit vector, is given by $\log(\mathbf{r}) = \alpha \mathbf{n}$. The vector $\alpha \mathbf{n}$ is an element of the Lie algebra $\mathfrak{so}(3)$.

The base quaternions are:

$$\mathbf{r}_h(\theta_h) = \cos\left(\frac{\theta_h}{2}\right) + \sin\left(\frac{\theta_h}{2}\right)\mathbf{j} \tag{11}$$

$$\mathbf{r}_w(\theta_w) = \cos\left(\frac{\theta_w}{2}\right) + \sin\left(\frac{\theta_w}{2}\right)\mathbf{k} \tag{12}$$

Applying the logarithm map to each, we get:

$$\log(\mathbf{r}_h(\theta_h)) = \frac{\theta_h}{2}\mathbf{j} \tag{13}$$

$$\log(\mathbf{r}_w(\theta_w)) = \frac{\theta_w}{2}\mathbf{k} \tag{14}$$

In the vector space $\mathfrak{so}(3) \cong \mathbb{R}^3$, these correspond to the vectors $(0, \theta_h/2, 0)$ and $(0, 0, \theta_w/2)$.

We compute the arithmetic mean of these vectors in the Lie algebra:

$$\mathbf{u} = \frac{1}{2} \left(\log(\mathbf{r}_h) + \log(\mathbf{r}_w) \right) = \frac{1}{2} \left(\frac{\theta_h}{2} \mathbf{j} + \frac{\theta_w}{2} \mathbf{k} \right) = \frac{\theta_h}{4} \mathbf{j} + \frac{\theta_w}{4} \mathbf{k}$$
 (15)

This corresponds to the vector $(0, \theta_h/4, \theta_w/4)$, as stated in the main text.

The exponential map for a Lie algebra vector \mathbf{u} is given by $\exp(\mathbf{u}) = \cos(\|\mathbf{u}\|) + \sin(\|\mathbf{u}\|) \frac{\mathbf{u}}{\|\mathbf{u}\|}$.

First, we compute the norm of our averaged vector u:

$$\|\mathbf{u}\| = \sqrt{\left(\frac{\theta_h}{4}\right)^2 + \left(\frac{\theta_w}{4}\right)^2} = \frac{1}{4}\sqrt{\theta_h^2 + \theta_w^2} \tag{16}$$

Let us define the coupled phase $\Theta = \frac{1}{2} \sqrt{\theta_h^2 + \theta_w^2}$. Then, $\|\mathbf{u}\| = \frac{\Theta}{2}$.

Next, we find the corresponding unit axis vector:

$$\frac{\mathbf{u}}{\|\mathbf{u}\|} = \frac{\frac{\theta_h}{4}\mathbf{j} + \frac{\theta_w}{4}\mathbf{k}}{\frac{1}{4}\sqrt{\theta_h^2 + \theta_w^2}} = \frac{\theta_h\mathbf{j} + \theta_w\mathbf{k}}{\sqrt{\theta_h^2 + \theta_w^2}} = \frac{\theta_h}{2\Theta}\mathbf{j} + \frac{\theta_w}{2\Theta}\mathbf{k}$$
(17)

Finally, applying the exponential map yields the desired symmetric operator:

$$\mathbf{r} = \exp(\mathbf{u}) = \cos\left(\frac{\Theta}{2}\right) + \sin\left(\frac{\Theta}{2}\right) \left(\frac{\theta_h}{2\Theta}\mathbf{j} + \frac{\theta_w}{2\Theta}\mathbf{k}\right)$$
(18)

This completes the derivation.

B INNER PRODUCT WITH ROTATED VECTORS

This section provides the derivation for the inner product of a vector \mathbf{q} with a rotated vector $\mathbf{R}\mathbf{k}$, as presented in the discussion on the geometric interpretation of attention.

A rotation of a vector $\mathbf{k} \in \mathbb{R}^3$ by an angle A around a unit axis vector $\mathbf{n} \in \mathbb{R}^3$ is given by Rodrigues' rotation formula:

$$\mathbf{R}\mathbf{k} = \mathbf{k}\cos(A) + (\mathbf{n} \times \mathbf{k})\sin(A) + \mathbf{n}(\mathbf{n} \cdot \mathbf{k})(1 - \cos(A))$$
(19)

To find the attention score, we compute the inner product of a query vector \mathbf{q} with this rotated key vector:

$$\langle \mathbf{q}, \mathbf{R} \mathbf{k} \rangle = \langle \mathbf{q}, \mathbf{k} \cos(A) + (\mathbf{n} \times \mathbf{k}) \sin(A) + \mathbf{n} (\mathbf{n} \cdot \mathbf{k}) (1 - \cos(A)) \rangle$$
 (20)

By the linearity of the inner product, we can distribute q across the terms:

$$\langle \mathbf{q}, \mathbf{R} \mathbf{k} \rangle = \langle \mathbf{q}, \mathbf{k} \rangle \cos(A) + \langle \mathbf{q}, (\mathbf{n} \times \mathbf{k}) \rangle \sin(A) + \langle \mathbf{q}, \mathbf{n} (\mathbf{n} \cdot \mathbf{k}) \rangle (1 - \cos(A))$$
 (21)

The last term can be simplified:

$$\langle \mathbf{q}, \mathbf{n}(\mathbf{n} \cdot \mathbf{k}) \rangle = (\mathbf{q} \cdot \mathbf{n})(\mathbf{n} \cdot \mathbf{k})$$
 (22)

The middle term involves the scalar triple product, which satisfies the identity $\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c}) = \mathbf{b} \cdot (\mathbf{c} \times \mathbf{a}) = \mathbf{c} \cdot (\mathbf{a} \times \mathbf{b})$. Let $\mathbf{a} = \mathbf{q}, \mathbf{b} = \mathbf{n}, \mathbf{c} = \mathbf{k}$. Then:

$$\langle \mathbf{q}, (\mathbf{n} \times \mathbf{k}) \rangle = \mathbf{k} \cdot (\mathbf{q} \times \mathbf{n}) = -\mathbf{k} \cdot (\mathbf{n} \times \mathbf{q})$$
 (23)

Substituting these back, we obtain the final expression for the inner product:

$$\langle \mathbf{q}, \mathbf{R} \mathbf{k} \rangle = \langle \mathbf{q}, \mathbf{k} \rangle \cos(A) + (\mathbf{q} \cdot \mathbf{n})(\mathbf{k} \cdot \mathbf{n})(1 - \cos(A)) - (\mathbf{n} \times \mathbf{q}) \cdot \mathbf{k} \sin(A)$$
 (24)

Note: The sign of the final term may vary depending on the convention used for the scalar triple product permutation, but the geometric intuition remains the same. The version in the main text is a common variant.

C LONG-DISTANCE DECAY PROPERTY

A key inductive bias of RoPE, which GeoPE is designed to generalize, is the decay of attention scores over large relative distances. We demonstrate that GeoPE preserves this behavior by analyzing the structure of the attention score's dominant term. The leading term in the total attention score is the sum:

$$S = \sum_{i=0}^{d/3-1} \langle \mathbf{q}_i, \mathbf{k}_i \rangle \cos(A_i)$$
 (25)

where the angle A_i for the *i*-th sub-vector is proportional to the relative distance $||\Delta p||$ and a frequency term $\lambda^{2i/d}$:

$$A_{i} = \frac{1}{2} \sqrt{(\Delta p_{h} \cdot \lambda^{2i/d})^{2} + (\Delta p_{w} \cdot \lambda^{2i/d})^{2}} = \frac{1}{2} ||\Delta p|| \lambda^{2i/d}$$
 (26)

Let $D = \frac{1}{2}||\Delta p||$ be the effective distance, $c_i = \langle \mathbf{q}_i, \mathbf{k}_i \rangle$ be the feature similarity, and $\phi_i = \lambda^{2i/d}$. The sum is $S = \sum_{i=0}^{N-1} c_i \cos(D\phi_i)$, where N = d/3.

To show this sum decays with D, we apply summation by parts (Abel transformation). Let $h_i = c_i$ and $g_i = \cos(D\phi_i)$. Let $G_k = \sum_{i=0}^k g_i$ be the partial sum of the cosine terms. The total sum is:

$$S = \sum_{i=0}^{N-1} h_i g_i = h_{N-1} G_{N-1} - \sum_{i=0}^{N-2} (h_{i+1} - h_i) G_i$$
 (27)

By the triangle inequality:

$$|S| \le |h_{N-1}||G_{N-1}| + \sum_{i=0}^{N-2} |h_{i+1} - h_i||G_i|$$
(28)

Assuming the feature similarities c_i are well-behaved (i.e., bounded and with small successive differences, a reasonable assumption for trained embeddings), the magnitude of S is primarily controlled by the magnitude of the partial sums $|G_k|$.

The partial sum $G_k = \sum_{i=0}^k \cos(D\phi_i)$ is a sum of cosines with geometrically increasing frequencies $(\phi_i = \lambda^{2i/d})$. For a large distance D, the arguments $D\phi_i$ grow rapidly, causing the cosine terms to oscillate with increasing frequency. Such sums are bounded due to destructive interference. While a simple closed-form bound is not available as in the arithmetic case (original RoPE), the geometric progression of frequencies ensures that the terms do not align constructively, keeping $|G_k|$ bounded for any k. As $D \to \infty$, the oscillations become more rapid, strengthening the cancellation effect. This implies that the average magnitude of the attention score decays with distance, preserving the crucial inductive bias for locality, analogous to the property shown in (Su et al., 2024).

D DEGENERATION OF GEOPE TO 1D ROPE

For 1D sequential data, GeoPE gracefully degenerates to a formulation equivalent to the original RoPE. In a 1D setting, we only have a single position index, say p, and its corresponding phase is $\theta = p \cdot \lambda^{2i/d}$. Since there is only one spatial dimension, the log-exp averaging process is unnecessary.

We can define a single base rotation directly. Following the original RoPE, this is a 2D rotation, which can be embedded in our 3D framework as a rotation around a single fixed axis (e.g., the y-axis, **j**).

The rotational quaternion for a phase θ is simply:

$$\mathbf{r}(\theta) = \cos\left(\frac{\theta}{2}\right) + \sin\left(\frac{\theta}{2}\right)\mathbf{j} \tag{29}$$

In GeoPE, feature vectors are partitioned into 3D sub-vectors $\mathbf{v}=(v_x,v_y,v_z)$. For a 1D application, we can effectively work with 2D sub-vectors by setting one component to zero, e.g., $\mathbf{v}=(v_x,0,v_z)$. This corresponds to a pure quaternion $\mathbf{p}=v_x\mathbf{i}+v_z\mathbf{k}$.

The rotation is applied via the sandwich product $\mathbf{p}' = \mathbf{r}(\theta)\mathbf{pr}(\theta)^*$. The rotation matrix corresponding to $\mathbf{r}(\theta)$ is a pure rotation around the y-axis:

$$\mathbf{R}(\theta) = \begin{pmatrix} \cos(\theta) & 0 & \sin(\theta) \\ 0 & 1 & 0 \\ -\sin(\theta) & 0 & \cos(\theta) \end{pmatrix}$$
(30)

Applying this rotation to our 2D-like sub-vector $\mathbf{v}' = \mathbf{R}(\theta)\mathbf{v}$:

$$\begin{pmatrix} v_x' \\ 0 \\ v_z' \end{pmatrix} = \begin{pmatrix} \cos(\theta) & 0 & \sin(\theta) \\ 0 & 1 & 0 \\ -\sin(\theta) & 0 & \cos(\theta) \end{pmatrix} \begin{pmatrix} v_x \\ 0 \\ v_z \end{pmatrix} = \begin{pmatrix} v_x \cos(\theta) + v_z \sin(\theta) \\ 0 \\ -v_x \sin(\theta) + v_z \cos(\theta) \end{pmatrix}$$
 (31)

This operation is a 2D rotation on the coordinates (v_x,v_z) . The original RoPE applies the matrix $\begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix}$ to a pair of features (f_1,f_2) . The resulting transformation is $\begin{pmatrix} v_x' \\ v_z' \end{pmatrix} = \begin{pmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{pmatrix} \begin{pmatrix} v_x \\ v_z \end{pmatrix}$. By identifying v_x with f_1 and v_z with f_2 , this is equivalent to the standard RoPE rotation matrix for an angle of $-\theta$. Thus, GeoPE contains RoPE as a special case, differing only by a sign convention on the rotation angle.

E THREE DIMENSION EXTENSION

This section details the derivation for the 3D symmetric rotational operator, extending the logic from Appendix A.

For 3D data with positions (d, h, w), we have three base quaternions corresponding to rotations about the **i**, **j**, and **k** axes:

$$\mathbf{r}_d(\theta_d) = \cos\left(\frac{\theta_d}{2}\right) + \sin\left(\frac{\theta_d}{2}\right)\mathbf{i} \tag{32}$$

$$\mathbf{r}_h(\theta_h) = \cos\left(\frac{\theta_h}{2}\right) + \sin\left(\frac{\theta_h}{2}\right)\mathbf{j} \tag{33}$$

$$\mathbf{r}_w(\theta_w) = \cos\left(\frac{\theta_w}{2}\right) + \sin\left(\frac{\theta_w}{2}\right)\mathbf{k} \tag{34}$$

We map these to the Lie algebra $\mathfrak{so}(3)$:

$$\log(\mathbf{r}_d) = \frac{\theta_d}{2}\mathbf{i} \tag{35}$$

$$\log(\mathbf{r}_h) = \frac{\theta_h}{2}\mathbf{j} \tag{36}$$

$$\log(\mathbf{r}_w) = \frac{\theta_w}{2}\mathbf{k} \tag{37}$$

We compute the arithmetic mean of these three vectors:

$$\mathbf{u} = \frac{1}{3} \left(\log(\mathbf{r}_d) + \log(\mathbf{r}_h) + \log(\mathbf{r}_w) \right) = \frac{\theta_d}{6} \mathbf{i} + \frac{\theta_h}{6} \mathbf{j} + \frac{\theta_w}{6} \mathbf{k}$$
(38)

Next, we compute the norm of this averaged vector **u**:

$$\|\mathbf{u}\| = \sqrt{\left(\frac{\theta_d}{6}\right)^2 + \left(\frac{\theta_h}{6}\right)^2 + \left(\frac{\theta_w}{6}\right)^2} = \frac{1}{6}\sqrt{\theta_d^2 + \theta_h^2 + \theta_w^2}$$
 (39)

Let's define the 3D composite phase $\Theta = \frac{1}{3}\sqrt{\theta_d^2 + \theta_h^2 + \theta_w^2}$. Then, $\|\mathbf{u}\| = \frac{\Theta}{2}$.

The unit axis vector is:

$$\frac{\mathbf{u}}{\|\mathbf{u}\|} = \frac{\frac{\theta_d}{6}\mathbf{i} + \frac{\theta_h}{6}\mathbf{j} + \frac{\theta_w}{6}\mathbf{k}}{\frac{1}{6}\sqrt{\theta_d^2 + \theta_h^2 + \theta_w^2}} = \frac{\theta_d\mathbf{i} + \theta_h\mathbf{j} + \theta_w\mathbf{k}}{\sqrt{\theta_d^2 + \theta_h^2 + \theta_w^2}} = \frac{\theta_d}{3\Theta}\mathbf{i} + \frac{\theta_h}{3\Theta}\mathbf{j} + \frac{\theta_w}{3\Theta}\mathbf{k}$$
(40)

Finally, applying the exponential map $\exp(\mathbf{u}) = \cos(\|\mathbf{u}\|) + \sin(\|\mathbf{u}\|) \frac{\mathbf{u}}{\|\mathbf{u}\|}$ yields the 3D GeoPE operator:

$$\mathbf{r} = \cos\left(\frac{\Theta}{2}\right) + \sin\left(\frac{\Theta}{2}\right) \left(\frac{\theta_d}{3\Theta}\mathbf{i} + \frac{\theta_h}{3\Theta}\mathbf{j} + \frac{\theta_w}{3\Theta}\mathbf{k}\right) \tag{41}$$

F LIMITATIONS

The Linear GeoPE variant enforces a strict linear inductive bias, but its current implementation incurs a significant memory footprint, with peak memory allocation increasing by over 200% compared to baselines, despite having similar FLOPs. In contrast, our standard GeoPE shows no such overhead; its memory usage is nearly identical to APE (less than 2% difference). This overhead is specific to Linear GeoPE's need to materialize relative rotation matrices and is an implementation-level challenge, not a fundamental limitation. We are confident it can be effectively mitigated with a custom CUDA kernel.