# A 2D Semantic-Aware Position Encoding for Vision Transformers

**Xi Chen**[1][*] **Shiyang Zhou**[1][*] **Muqi Huang**[2], **Jiaxu Feng**[3], **Yun Xiong**[1][†] **Kun Zhou**[2],
**Biao Yang**[1], **Yuhui Zhang**[1], **Huishuai Bao**[1], **Sijia Peng**[1], **Chuan Li**[2], **Feng Shi**[2]

**[1]Shanghai Key Laboratory of Data Science, School of Computer Science,
Fudan University, Shanghai, China
[2]Alibaba Group, Shanghai, China
[3]Department of Physics, Fudan University, Shanghai, China**

{x_chen21, 22210240079}@m.fudan.edu.cn, huangmugi.hmg@alibaba-inc.com, feng.jiaxu@outlook.com,

yunx@fudan.edu.cn, kun.zhouk@alibaba-inc.com, {22210240343, 23212010049, 22212010001}@m.fudan.edu.cn,

pengsijia1216@gmail.com, {lc357677, sam.sf}@alibaba-inc.com

## Abstract

Vision transformers have demonstrated significant advantages in computer vision tasks due to their ability to capture long-range dependencies and contextual relationships through self-attention. However, existing position encoding techniques, which are largely borrowed from natural language processing, fail to effectively capture semantic-aware positional relationships between image patches. Traditional approaches like absolute position encoding and relative position encoding primarily focus on 1D linear position relationship, often neglecting the semantic similarity between distant yet contextually related patches. These limitations hinder model generalization, translation equivariance, and the ability to effectively handle repetitive or structured patterns in images. In this paper, we propose 2-Dimensional Semantic-Aware Position Encoding ($SaPE^2$), a novel position encoding method with semantic awareness that dynamically adapts position representations by leveraging local content instead of fixed linear position relationship or spatial coordinates. Our method enhances the model's ability to generalize across varying image resolutions and scales, improves translation equivariance, and better aggregates features for visually similar but spatially distant patches. By integrating $SaPE^2$ into vision transformers, we bridge the gap between position encoding and perceptual similarity, thereby improving performance on computer vision tasks.

## 1 Introduction

Transformer has been widely adopted in both natural language processing (NLP) [1] and computer vision (CV) [2, 3] due to its superior ability to capture long-range dependencies, efficient parallel computation, and flexible contextual modeling. Compared to traditional convolutional neural networks (CNNs) [4, 5], transformer-based architectures exhibit greater flexibility, enabling seamless adaptation to complex tasks such as cross-modal learning (e.g., joint image-text processing), image generation, and high-level semantic reasoning. Their self-attention mechanism provides strong expressive power, making them highly competitive in visual understanding.

Unlike CNNs [4], which operate on spatially structured grid-based feature maps, vision transformers process images in a sequence-based manner [2]. Typically, an input image is partitioned into a set

---

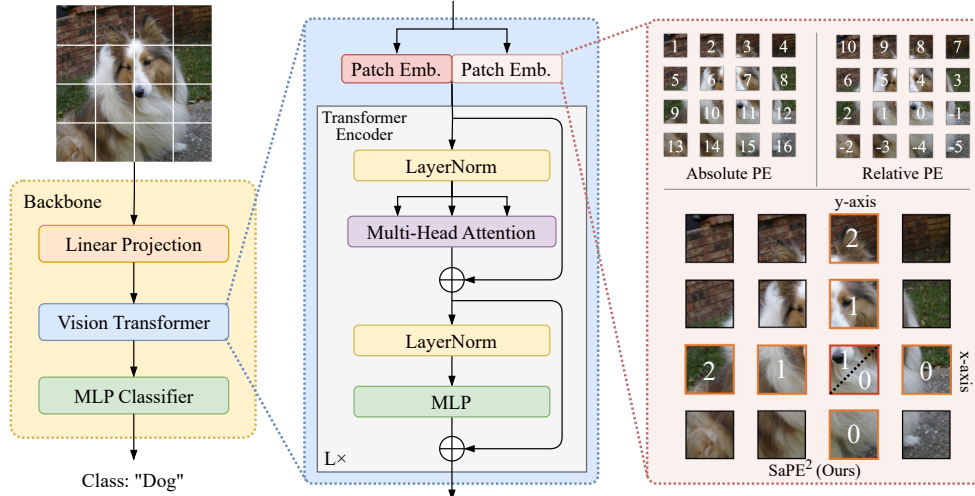[*]Equal Contribution.

[†]Corresponding Author.

Figure 1: The overall pipeline of the vision model and a comparison between the absolute PE, relative PE, and SaPE$^2$. The patch with the red border is the target patch for relative PE and our proposed method (SaPE$^2$). The number in each patch indicates the position information of the patch. This is just a schematic diagram; please refer to Figure 2 for details on SaPE$^2$.

of non-overlapping patches, which are then linearly embedded into tokens before being fed into a transformer model [2, 3, 6]. These tokens are treated equally and interact through self-attention, where pairwise attention scores determine their relationships. However, since the self-attention mechanism itself is position-agnostic (permutation-invariance of self-attention operation), explicit position encoding (PE) is required to introduce location information into the model. To this end, position encodings are integrated into the model's attention mechanism, ensuring that the computed attention scores incorporate spatial relationships between tokens.

While existing position encoding techniques in computer vision effectively incorporate spatial awareness into vision transformers, they rely exclusively on explicit spatial coordinates while overlooking semantic-aware positional relationships between patches. Specifically, absolute position encoding (APE) [1, 7] is coordinate-dependent, making it sensitive to image resolution and scale variations, thereby limiting its adaptability to different input sizes. Relative position encoding (RPE) [8, 9] (including Rotary Position Embedding [10]) focuses only on spatial displacement, without incorporating the semantic similarity between different regions of the image. None of these methods consider the semantic relationships among patches, meaning they fail to capture perceptual similarity beyond local adjacency.

These limitations are particularly problematic in structured and repetitive patterns, where visually similar regions are treated as completely independent simply due to their spatial separation. For instance, two patches representing the same object category (e.g., two instances of a car in an image) might receive distinct encodings based on their absolute coordinates, even though they share strong semantic similarity. This discrepancy prevents the model from effectively recognizing and leveraging their relatedness. Similarly, in images with repetitive textures or grid-like structures, conventional position encodings may fail to exploit the redundancy in attention computation, as semantically correlated tokens are treated as independent, leading to suboptimal information aggregation.

Moreover, position encoding methods that rely solely on coordinate-based distances hinder translation equivariance. When an object shifts within an image, traditional position encoding schemes alter its spatial representation, leading to inconsistencies in feature extraction. This can be particularly detrimental to model robustness in image classification, where maintaining spatially consistent representations is crucial for recognizing objects regardless of their position in the image.

These challenges highlight the need for a more flexible, semantic-aware position encoding approach—one that adapts dynamically to local content rather than relying solely on fixed spatial coordinates.

To address the aforementioned limitations, we propose 2-Dimensional Semantic-Aware Position Encoding (SaPE$^2$), which encodes positional relationships based on semantic similarity rather than

predefined spatial distances. By dynamically adjusting position representations according to local context, our method enables vision transformers to (1) generalize more effectively across different resolutions and scales, (2) improve translation equivariance, and (3) enhance feature aggregation for visually similar yet spatially distant regions.

Our proposed semantic-aware position encoding allows the model to assign similar positional representations to semantically related regions, enhancing its understanding of global context. For instance, in the given image in Figure 1 (a dog resting near a corner wall), traditional position encoding methods would treat each patch independently based solely on spatial coordinates. In contrast, our method recognizes that patches corresponding to the dog belong to a coherent foreground entity, while patches representing the brick wall and grass form the background. By assigning similar position encodings to semantically related patches, our approach helps the model better differentiate between the foreground and background, leading to more consistent feature representations. By integrating SaPE[2] into vision transformers, our method bridges the gap between previous position encoding and perceptual similarity, leading to improved generalization in downstream tasks.

We summarize our contributions as follows:

- We propose a semantic-aware approach to position encoding, integrating it with the inherent 2D spatial structure of image inputs. This enables vision transformers to capture semantic relationships between patches rather than relying solely on fixed spatial coordinates. Moreover, it ensures a cohesive representation that enhances both spatial and semantic understanding in transformer-based vision models.
- We conduct extensive experiments, benchmarking our method against traditional position encoding techniques and thoroughly analyzing its effectiveness and impact on model performance.
- We identify the limitations of commonly used position encoding methods in transformer-based vision models and empirically demonstrate that our proposed method offers a promising solution.

## 2 Preliminaries

**Absolute Position Encoding (APE)**    To incorporate the sequential order information into ViTs, APE assigns a unique positional representation to each token. Given the token embeddings $\mathbf{x}_i \in \mathbb{R}^{d_x}$, the absolute position encoding $\mathbf{p}_i \in \mathbb{R}^{d_x}$ is added directly to the embeddings, yielding the final input representation:

$$\mathbf{x}'_i = \mathbf{x}_i + \mathbf{p}_i. \tag{1}$$

The position encoding $\mathbf{p}_i$ can be either fixed or learnable. Learnable encodings associate each position $i$ with a unique trainable parameter vector $\mathbf{e}[i]$ [2], while fixed encodings often leverage sinusoidal functions with varying frequencies [1]. The sinusoidal position encoding defines each dimension of the position vector as:

$$\mathbf{p}_i^{(2j)} = \sin\left(\frac{i}{10000^{\frac{2j}{d_x}}}\right), \mathbf{p}_i^{(2j+1)} = \cos\left(\frac{i}{10000^{\frac{2j}{d_x}}}\right), \tag{2}$$

where $i$ denotes the position index, $j$ is the dimension index, and $d_x$ represents the dimension of the token embeddings. The frequency of the sinusoidal functions is determined by the denominator, which scales the position index $i$ by an exponential function of the dimension index. This design enables the model to generalize to longer sequences.

**Relative Position Encoding (RPE)**    Unlike APE, which assigns absolute positional information to each token, RPE models the pairwise distance between tokens, making the encoding invariant to translation. For a pair of tokens at positions $i$ and $j$, the relative position encoding is represented as $r_{i,j}$, and the final attention score is modified as:

$$a_{ij} = \frac{\mathbf{q}_i^\top \mathbf{k}_j + r_{i,j}}{\sqrt{d_k}}, \tag{3}$$

where $\mathbf{q}_i, \mathbf{k}_j \in \mathbb{R}^{d_k}$ are the query and key representations, and $d_k$ is the dimensionality of the key vectors. The relative position encoding $r_{i,j}$ can be implemented as a learnable parameter [8] or a function of the relative distance $i - j$.
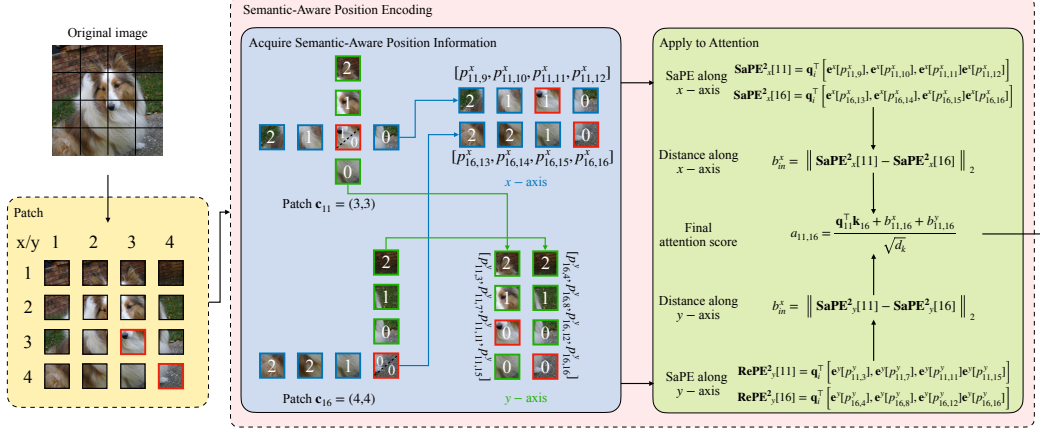
Figure 2: The pipeline of calculating the relative positional relationship between two patches which in the different x-axes and y-axes.

**Rotary Position Encoding (RoPE)**   RoPE is a special form of RPE. It encodes positional information by applying a rotation matrix to the query and key embeddings. This method preserves the relative positional relationship between tokens in the inner product space:

$$\mathbf{q}'_n = \mathbf{q}_n e^{in\theta}, \mathbf{k}'_m = \mathbf{k}_m e^{im\theta}, \tag{4}$$

$$\mathbf{A}'_{(n,m)} = \mathrm{Re}\left[\mathbf{q}'_n \mathbf{k}'^*_m\right] = \mathrm{Re}\left[\mathbf{q}_n \mathbf{k}^*_m e^{i(n-m)\theta}\right], \tag{5}$$

where $n$ and $m$ denote the positions of the corresponding tokens, and $\theta$ represents the angular frequency that varies across different feature dimensions. The varying frequencies $\theta_t$ are defined as:

$$\theta_t = 10000^{-t/(d_{\text{head}}/2)}, \text{ where } t \in \{0, 1, \ldots, d_{\text{head}}/2\}, \tag{6}$$

where $d_{\text{head}}$ is the dimensionality of each attention head. This design assigns higher frequencies to lower-dimensional features and lower frequencies to higher-dimensional features. The exponential scaling ensures that different dimensions encode positional information at different granularities, with lower dimensions capturing finer positional distinctions and higher dimensions capturing coarser positional patterns.

**Contextual Position Encoding (CoPE)**   CoPE [11] is initially proposed in NLP tasks and it measures positions in a context-dependent way, rather than simple token counts. First, a gate value is computed to determine whether a key token should be included in the position measurement:

$$g_{ij} = \sigma\left(\mathbf{q}_i^\top \mathbf{k}_j\right), \tag{7}$$

where $j < i$ and $\sigma$ denotes the sigmoid function. The gate value ranges from 0 to 1, indicating whether the token is counted or ignored. The position value is then obtained by summing the gate values between the current and target tokens:

$$p_{ij} = \sum_{k=j}^{i} g_{ik}. \tag{8}$$

If all gates are 1, $p_{ij}$ corresponds to the token-based relative position, making CoPE a generalization of relative position encoding. Generally, $p_{ij}$ can be the count of specific words or word types like nouns or numbers, the number of sentences, or other useful concepts. The final PE can be calculated by interpolating trainable position embeddings $\mathbf{e}[p]$.

## 3   Method

In this section, we present the detailed implementation of our proposed SaPE$^2$ method. The overall architecture is illustrated in Figure 1 and 2.

Given that images are naturally represented as 2D grids, the core idea is to decompose the 2D position encoding into two independent 1D position encodings along the horizontal and vertical axes. This decomposition enables the model to capture spatial information in a separable manner, enhancing both local and global spatial dependencies.

Let $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$ denote an input image, where $H$ and $W$ represent the number of patches along the vertical ($y$) and horizontal ($x$) axes, respectively, and $C$ is the number of channels. The coordinate of the $i$-th patch is denoted by $\mathbf{c}_i = (x_i, y_i)$, where $x_i \in \{1, 2, \ldots, W\}$ and $y_i \in \{1, 2, \ldots, H\}$.

SaPE$^2$ encodes spatial information by applying independent 1D position encodings along each axis. For the horizontal axis, let $\mathbf{q}_i$, $\mathbf{k}_i$, and $\mathbf{v}_i$ denote the query, key, and value representations of the $i$-th patch. The gate value for each patch pair sharing the same $y$ coordinate is computed as:

$$g_{ij}^x = \sigma\left(\mathbf{q}_i^\top \mathbf{k}_j\right), \quad \text{where } y_i = y_j. \tag{9}$$

The gate value $g_{ij}^x$ encodes the relative position between the $i$-th and $j$-th patches, enabling the model to perceive spatial relationships such as foreground and background information. The position value along the $x$-axis is obtained by summing the gate values:

$$p_{im}^x = \sum_{j \in \mathcal{I}_{im}} g_{ij}, \quad \mathcal{I}_{im} = \{j \mid y_i = y_j = y_m,\ x_m \le x_j\}. \tag{10}$$

Here, $p_{im}^x$ answers the question: what is the relative position of the $m$-th patch from the $i$-th patch's perspective along the $x$-axis? In general, it captures transitions in color, object boundaries, or other latent patterns that the transformer deems useful for understanding the image content.

Since $p_{im}^x$ is not restricted to integer values, we cannot directly convert a position value to a position vector in the same way as the relative PE. Inspired by CoPE [11], we use interpolation between integer embeddings to obtain continuous positional representations. Let $\mathbf{e}[p]$ denote the learnable embedding for integer position $p \in [0, M]$. The interpolated embedding is computed as:

$$\mathbf{e}^x\left[p_{im}^x\right] = \left(p_{im}^x - \lfloor p_{im}^x \rfloor\right) \mathbf{e}^x\left[\lceil p_{im}^x \rceil\right] + \left(1 - p_{im}^x + \lfloor p_{im}^x \rfloor\right) \mathbf{e}^x\left[\lfloor p_{im}^x \rfloor\right], \tag{11}$$

where $\mathbf{e}^x\left[p_{im}^x\right]$ represents the position embedding of the $m$-th patch relative to the $i$-th patch along the $x$-axis. This embedding encodes the spatial relationship between the two patches, answering the question: what does it mean for the $m$-th patch to be at this particular position from the $i$-th patch's perspective?

Finally, the comprehensive position embedding that characterizes the location of the $i$-th patch along the $x$-axis is obtained through the following equation:

$$\mathbf{SaPE^2}_x[i] = \mathbf{q}_i^\top \left[\mathbf{e}^x\left[p_{im_1}^x\right], \mathbf{e}^x\left[p_{im_2}^x\right] \ldots \mathbf{e}^x\left[p_{im_W}^x\right]\right] = [z_{i1}, z_{i2}, \ldots, z_{iW}],$$
$$\text{where } x_{m_1} = 1, x_{m_2} = 2, \ldots, x_{m_W} = W. \tag{12}$$

Here $z_{ir} = \mathbf{q}_i^\top \mathbf{e}^x\left[p_{im_r}\right]$ answers the question: what is the contribution of the $r$-th patch's position embedding to the representation of the $i$-th patch along the $x$-axis? Except for the query mode, denoted as SaPE$^2$(Q), we can also switch to the key mode $\tilde{z}_{ir} = \mathbf{k}_i^\top \mathbf{e}^x\left[p_{im_r}\right]$, denoted as SaPE$^2$(K). The implementation details are illustrated in Figure 3. The term $z_{ir}$ serves as the attention bias added to $\mathbf{q}_i^\top \mathbf{k}_r$ in the 1D setting commonly used in NLP tasks. However, since we need to model the interactions between patches across both rows and columns, the neighboring content surrounding each patch must be incorporated. Therefore, the vector $\mathbf{SaPE^2}_x$ is constructed to provide a richer contextual representation by aggregating the relative position information of neighboring patches.

In practice, directly computing and storing the vectors $\mathbf{q}_i^\top \mathbf{e}^x\left[p_{im_r}\right]$ incurs additional computational and memory overhead. A more efficient approach is to first compute the multiplications $\mathbf{q}_i^\top \mathbf{e}^x[p]$ for all integer positions $p$, and then obtain the desired values through interpolation. The $z_i\left[p_{im_r}^x\right]$ acquired from Equation (14) is equivalent to $z_{ir}$ in Equation (12):

$$z_i[p] = \mathbf{q}_i^\top \mathbf{e}^x[p] \quad \text{for } p \in [0, 1, \ldots, M] \tag{13}$$

$$z_i\left[p_{im}^x\right] = \left(p_{im}^x - \lfloor p_{im}^x \rfloor\right) z_i\left[\lceil p_{im}^x \rceil\right] + \left(1 - p_{im}^x + \lfloor p_{im}^x \rfloor\right) z_i\left[\lfloor p_{im}^x \rfloor\right]. \tag{14}$$

Using the $\mathbf{SaPE^2}_x$ vector, the relative positional bias between the $i$-th and the $n$-th patch along the $x$-axis can be computed based on their Euclidean distance. This formulation generalizes the
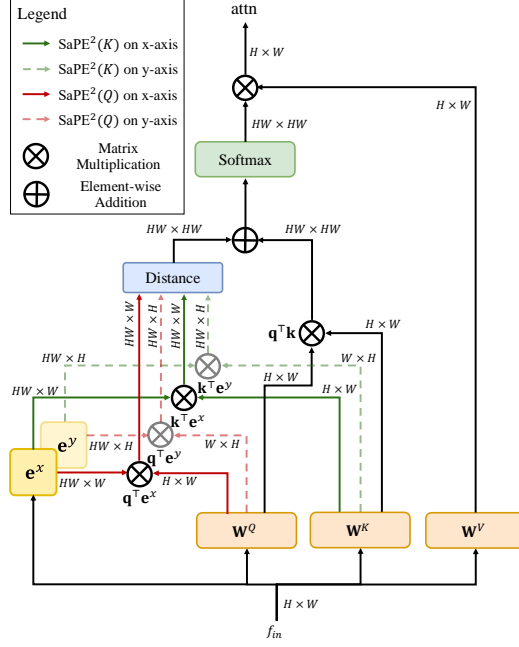
Figure 3: Self-attention mechanism with SaPE$^2$, the position encoding can be applied to either key or query. For simplicity, we ignore the channel dimension here. The dimensions and symbols outside the parentheses indicate operations on the x-axis, while those inside indicate operations on the y-axis.

bias calculation to arbitrary pairs of patches. The $n$-th patch is not constrained to share the same $y$ coordinate with the $i$-th patch, enabling the model to capture cross-row spatial interactions and enhancing its capacity to model complex spatial structures across the entire image:

$$b_{in}^x = \left\| \mathbf{SaPE^2}_x[i] - \mathbf{SaPE^2}_x[n] \right\|_2. \tag{15}$$

By calculating the bias for the $y$-axis in the same way, we can obtain the final attention weights with SaPE$^2$:

$$a_{im} = \frac{\mathbf{q}_i^\top \mathbf{k}_m + b_{in}^x + b_{in}^y}{\sqrt{d_k}}. \tag{16}$$

This decomposition-based approach improves the model's ability to capture both local and global spatial dependencies. By encoding each axis separately, the model gains greater flexibility in adapting to the inherent structure of image data.

## 4 Complexity Analysis

### 4.1 Space Complexity Analysis

The space complexity of SaPE$^2$ is mainly attributed to storing intermediate attention gates, positional values, and bias terms. For each axis, the attention gates require $\mathcal{O}(N \cdot W)$ storage, as they capture all pairwise interactions along a row or column, where $N = H * W$, representing the number of patches. The interpolated relative positions and their embeddings further introduce a cost of $\mathcal{O}(N \cdot W)$ and $\mathcal{O}(M \cdot d)$, respectively, where $M$ is the maximum integer position index. The interpolated position vectors $\mathbf{SaPE^2}$ occupy $\mathcal{O}(N \cdot W)$ memory, and the pairwise positional biases require $\mathcal{O}(N^2)$ space per axis. Consequently, the total space complexity is dominated by the storage of pairwise bias matrices, which scales quadratically with the number of patches.

### 4.2 Time Complexity Analysis

The time complexity of the proposed SaPE$^2$ method arises primarily from computing gated inter-actions and interpolated positional embeddings. For each row (or column), the gate values $g_{ij}$ are

Table 1: Experimental results in terms of Top 1 Accuracy (%). Higher accuracy shows better performance.

| Model | CIFAR10 | CIFAR100 |
|---|---|---|
| ViT | 87.41 | 66.54 |
| 2D RoPE | 89.92 | 70.71 |
| 2D RoPE + APE | 89.98 | 71.32 |
| CoPE + APE | 92.54 | 72.10 |
| SaPE$^2$ | 91.06 | 71.29 |
| SaPE$^2$ + APE | **93.98** | **72.23** |

computed by evaluating inner products between all patch pairs along the same axis, resulting in a time complexity of $\mathcal{O}(N \cdot W \cdot d)$ for horizontal attention and similarly for vertical. The relative position values $p_{im}$ are accumulated through masked summations over gate values, requiring $\mathcal{O}(N \cdot W)$ operations. Interpolated embeddings are obtained by bilinear interpolation between integer positional embeddings, which introduces $\mathcal{O}(N \cdot W \cdot d)$ cost. To enhance efficiency, we adopt an optimized computation strategy that first computes dot products for all integer positions and then performs interpolation, reducing redundancy. Finally, pairwise Euclidean distances between learned positional vectors are computed using $\mathcal{O}(N^2 \cdot W)$ operations per axis, followed by attention bias calculation at $\mathcal{O}(N^2 \cdot d)$. Overall, the computational bottleneck lies in the $\mathcal{O}(N^2)$ pairwise operations.

## 5 Experiments

To verify the effectiveness of our proposed position encoding, we apply SaPE$^2$ to the ViT model and modify the original position encoding module for training on an image classification task. We begin by evaluating the performance of SaPE$^2$ against commonly used position encoding methods, such as APE and 2D RoPE, on image classification tasks. Next, we examine the effect of position encoding on different components of the attention mechanism (e.g., Q and K, as illustrated in Figure 3) and its impact on downstream task performance. Finally, we empirically demonstrate that our proposed method addresses the limitations of traditional position encoding in capturing semantic information by optimizing position encoding and enhancing its capacity for semantic awareness.

### 5.1 Implementation Details

All experiments are conducted on a server equipped with four Nvidia H20 GPUs. All models are trained for 400 epochs, with a patch size of 4 and an input size of 32. The hidden dimension is set to 384, and the batch size is 128. We use the Adam optimizer to train the model. Model performance is evaluated using Top-1 Accuracy, where a higher score indicates better performance.

### 5.2 Comparison between Different PE Methods

To demonstrate the effectiveness of our approach, we compare SaPE$^2$ against commonly used position encoding methods from prior work by integrating different position encoding methods into a vanilla ViT model. The vanilla ViT model employs learnable absolute position encoding (APE). We replace the APE module with 2D RoPE, CoPE, and SaPE$^2$ for comparison. The details of the compared methods are as follows:

- **Vanilla ViT** [2]: We use ViT-Small as the base model, which adopts APE as its position encoding. All other comparison methods replace the position encoding in this model. We denote each method directly by the name of its position encoding module.

- **2D RoPE** [12]: Instead of applying the original RoPE, we directly adopt 2D RoPE as a baseline, as it introduces optimizations tailored for image data and vision models.

- **CoPE** [11]: We directly apply the CoPE method from NLP to image data and vision models. CoPE serves as a generalized form of the standard relative position encoding (RPE).

7

Table 2: Experimental results of image classification on CIFAR10 with APE and different settings of SaPE$^2$.

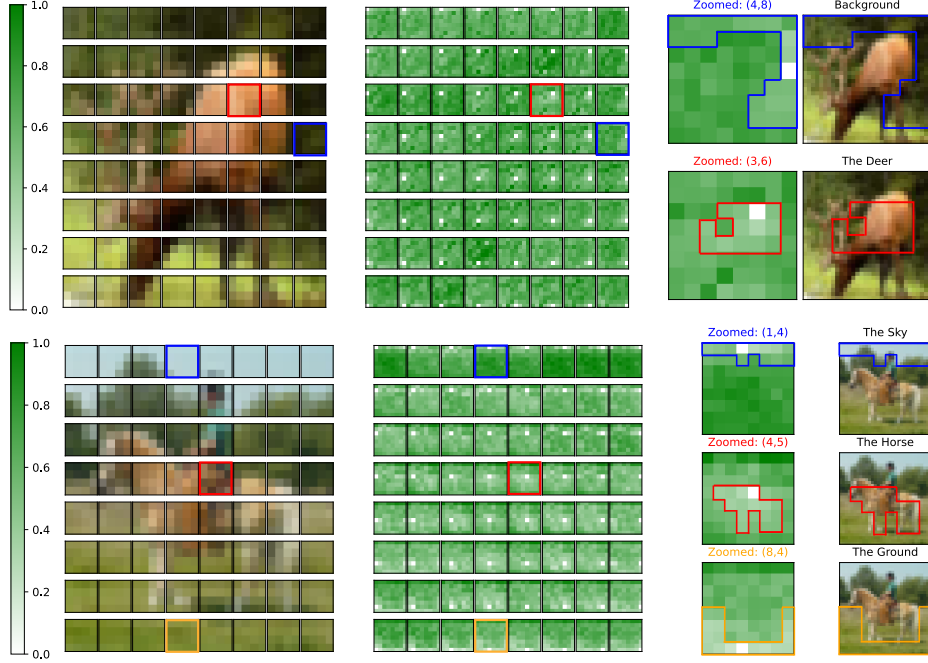| | Top 1 Acc.(%) | Top 5 Acc.(%) |
|---|---|---|
| APE | 87.41 | 99.47 |
| SaPE$^2$(K) + APE | **93.98** | **99.71** |
| SaPE$^2$(Q) + APE | 90.17 | 99.46 |



Figure 4: Visualization of the attention bias calculated based on SaPE$^2$. The patch size is set to 4, resulting in 8×8 patches for 32×32 images. The color transition from white to green indicates further semantic distance. **Left**: the original image; **Middle**: the SaPE$^2$ value between different patches; **Right**: visualization of SaPE$^2$ on patch from the background and the subject.The results indicate that our position encoding method effectively helps the model perceive contextual information.

- **\*+APE**: This approach retains the original APE in ViT while incorporating additional position encoding methods. Here, APE refers to the learnable version used in the vanilla ViT.

Since image classification tasks require model retraining, and our goal is to efficiently evaluate the effectiveness of the position encoding module, we conduct experiments on two public datasets, CIFAR10 [3] and CIFAR1100 [4].

The experimental results are presented in Table 1.

As observed, SaPE$^2$ achieves the best performance. We further provide the following analysis:

(1) Incorporating semantic information of patches into position encoding enhances model performance. Compared to APE and 2D RoPE, strategies that incorporate semantic information, such as CoPE and SaPE$^2$, enable the model to better differentiate instances of similar objects, thereby achieving superior performance.

(2) The 2D version of position encodings enhance the model's ability to capture spatial structures, thereby improving prediction accuracy. The comparison of RoPE and 2D RoPE, CoPE+APE and

---

[3]https://www.cs.toronto.edu/ kriz/cifar-10-python.tar.gz
[4]https://www.cs.toronto.edu/ kriz/cifar-100-python.tar.gz

SaPE$^2$+APE demonstrates the effectiveness of integrating 2D spatial information, enabling the model to better capture spatial relationships within the input image.

(3) Integrating both absolute and relative position information provides additional positional cues to the model, thereby improving its performance. This suggests that absolute and relative encodings provide complementary information: while relative encoding captures positional relationships, absolute encoding offers explicit positional differentiation, enhancing the model's ability to understand sequence structure.

### 5.3 Analysis of Integrating PE with Attention

As discussed in the methodology section, our proposed SaPE$^2$ can be integrated into the query (Q) and the key (K) within the self-attention module. To assess the effectiveness of applying SaPE$^2$ to different components of the attention module in vision transformers, we design three experimental settings: applying SaPE$^2$ to Q or K, or both Q and K simultaneously, denoted as SaPE$^2$(Q), SaPE$^2$(K). All other experimental settings follow the description in Section 5.2.

The results are presented in Table 2. We analyze the results as follows:

(1) Compared to solely using APE (i.e., the default PE setting in ViT), SaPE$^2$ significantly enhances model performance (Table 2 Line 1 vs. Lines 2 and 3). This demonstrates that semantic-aware position encoding enables the model to more effectively capture both semantic and spatial relationships.

(2) Integrating SaPE$^2$ with K yields the best performance, indicating that improving the query's semantic and spatial awareness within the attention module is the most effective approach. This is because, during attention computation, enriching K with semantic-aware PE enables the query to more effectively capture spatial relationships within the context. As a result, the model exhibits an enhanced ability to capture spatial dependencies among patches, leading to improved performance.

### 5.4 Case Study

To further validate the effectiveness of our proposed PE strategy, we visualize the attention bias induced by SaPE$^2$ in ViT-Small on the CIFAR-10 validation set. Figure 4 presents two examples. The images are divided into 8×8 patches, and each patch is represented as a small block in the attention map, illustrating its attention bias with respect to other patches.

Effectively recognizing and distinguishing the subject and background plays a crucial role in computer vision tasks, contributing to improved prediction performance. Take the "horse" image (Figure 4 bottom) as an example. The image can generally be divided into three parts: the sky above, the human and horse subjects in the center, and the ground below. Correspondingly, the PE maps on the right reveal that patches in the top two rows exhibit smaller distance from patches in the same region (i.e., the sky) while distant from the human, the horse, and the ground below. Similar patterns are observed for the subject and the ground. This indicates that patches tend to focus more on others with similar semantic contexts.

We attribute this effective attention mechanism to our proposed semantic-aware position encoding, SaPE$^2$, which enables ViT to better capture positional relationships between patches based on semantic similarity, thereby enhancing its performance in vision tasks.

## 6 Conclusion

In this paper, we propose a novel 2-dimensional semantic-aware position encoding method, SaPE$^2$, to address the limitations of existing position encoding techniques in transformer-based vision models, which primarily focus on absolute or relative spatial coordinates while overlooking relationships between semantically similar yet spatially distant regions. To this end, SaPE$^2$ dynamically adapts position representations based on contextual information, effectively capturing semantic relationships between different patches. Experimental results demonstrate that by bridging the gap between position encoding and perceptual similarity, SaPE$^2$ significantly enhances the performance of transformer-based vision models, presenting a promising direction for future research. The code of our proposed method will be open-sourced upon publication.

# References

[1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, NIPS'17, page 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc.

[2] 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.

[3] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.

[4] Keiron O'shea and Ryan Nash. An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458*, 2015.

[5] Zewen Li, Fan Liu, Wenjie Yang, Shouheng Peng, and Jun Zhou. A survey of convolutional neural networks: analysis, applications, and prospects. *IEEE transactions on neural networks and learning systems*, 33(12):6999–7019, 2021.

[6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.

[7] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional sequence to sequence learning. In *International conference on machine learning*, pages 1243–1252. PMLR, 2017.

[8] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.

[9] 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.

[10] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.

[11] Olga Golovneva, Tianlu Wang, Jason Weston, and Sainbayar Sukhbaatar. Contextual position encoding: Learning to count what's important. *arXiv preprint arXiv:2405.18719*, 2024.

[12] Byeongho Heo, Song Park, Dongyoon Han, and Sangdoo Yun. Rotary position embedding for vision transformer. In *ECCV*, pages 289–305. Springer, 2024.

[13] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 558–567, 2021.

[14] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021.

[15] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 568–578, 2021.

[16] Byeongho Heo, Sangdoo Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon Oh. Rethinking spatial dimensions of vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11936–11945, 2021.

[17] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 22–31, 2021.

[18] Kun Yuan, Shaopeng Guo, Ziwei Liu, Aojun Zhou, Fengwei Yu, and Wei Wu. Incorporating convolution designs into visual transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 579–588, 2021.

[19] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020.

[20] Minghang Zheng, Peng Gao, Renrui Zhang, Kunchang Li, Xiaogang Wang, Hongsheng Li, and Hao Dong. End-to-end object detection with adaptive clustering transformer. *arXiv preprint arXiv:2011.09315*, 2020.

[21] Peng Gao, Minghang Zheng, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Fast convergence of detr with spatially modulated co-attention. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3621–3630, 2021.

[22] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *European conference on computer vision*, pages 280–296. Springer, 2022.

[23] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. In *The Eleventh International Conference on Learning Representations*, 2022.

[24] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021.

[25] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7262–7272, 2021.

[26] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34:12077–12090, 2021.

[27] 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*, pages 10033–10041, 2021.

[28] 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*, pages 5886–5896, 2023.

# A  Related Work

## A.1  Position Encodings in NLP

Position encoding was first introduced in NLP transformers, where absolute position encoding (APE) [1, 7] and relative position encoding (RPE) [8, 9] became standard approaches. APE assigns fixed or learned embeddings to tokens based on absolute positions, often using sinusoidal functions or trainable vectors. However, APE does not capture relative positional information, making it sensitive to sequence length variations.

To address this, relative position encoding (RPE) modifies the attention mechanism to incorporate relative position biases, improving generalization to different sequence lengths. Rotary Position Embedding (RoPE) [10] extends this by encoding relative angular displacement through rotation matrices, which enhances the model's ability to handle longer sequences effectively.

Contextual position encoding (CoPE) [11] extends the idea of position encoding by dynamically adjusting positional representations based on contextual information rather than relying solely on predefined spatial coordinates. Unlike traditional absolute or relative position encodings, CoPE modulates position embeddings based on local content, allowing the model to better capture hierarchical and structural dependencies within a sequence. This adaptability enables improved generalization in scenarios where positional importance varies dynamically, such as in structured sequences or tasks requiring position-dependent reasoning.

## A.2  Vision Transformers

The impact of Transformer models has expanded beyond natural language processing to various domains, including computer vision. Vision Transformer (ViT) [2] is the first model in the field of computer vision to adopt a Transformer-based architecture as its backbone. ViT partitions an input image into multiple fixed-size patches, which are treated as tokens in a sequence and processed through a Transformer encoder. This approach has demonstrated superior performance compared to traditional CNN-based models in image classification tasks. Subsequent models have introduced various modifications to further enhance performance based on the ViT framework. T2T-ViT [13] enhances the tokenization process by progressively aggregating neighboring tokens, thereby capturing local structure information more effectively. DeiT [14] introduces a distillation token during training, enabling the model to learn from a teacher network and improving data efficiency. Pyramid ViT [15] and PiT [16] employ hierarchical architectures, facilitating multi-scale feature representation and reducing computational complexity. CvT [17] and CeiT [18] integrate convolutional layers into the Transformer framework, providing inductive biases that enhance the modeling of local features. Swin Transformer [3] utilizes a shifted window-based attention mechanism, which limits self-attention computation to non-overlapping local windows and achieves efficient modeling of long-range dependencies.

Transformer were later extended to downstream tasks in computer vision, including object detection and semantic segmentation. DETR [6] introduced an end-to-end object detection framework using a CNN backbone with a Transformer encoder-decoder, inspiring further optimizations [19, 20, 21]. Several models have employed vanilla ViT as a backbone, such as ViTDet [22], which incorporates upsampling and downsampling; and ViT-Adapter [23], which adds inductive biases for detection and segmentation tasks. SETR [24] first adopted ViT with a CNN decoder for semantic segmentation, while Segmenter [25] introduced a Transformer decoder. SegFormer [26] employed a hierarchical encoder to improve segmentation performance.

## A.3  Position Encodings in CV

Previous transformer-based vision models often adopted absolute position encoding, relative position encoding, or rotary position encoding, similar to methods used in natural language processing. Given that vision models inherently process 2D inputs, some subsequent models have adapted these NLP position encodings to two-dimensional formats.

A straightforward and simple idea is to extend relative position encoding to a 2D formulation [27]. This involves computing relative position encodings separately for the horizontal (x-axis) and vertical (y-axis) directions. These directional encodings are then combined to effectively represent the 2D relative positional relationships between pixels, enhancing the model's ability to understand spatial hierarchies and patterns in visual data. Furthermore, 2D RoPE [12] generalizes RoPE [10] to a two-dimensional setting. It achieve this by implementing axial and mixed learnable frequency approaches. The axial frequency method applies RoPE separately to the x-axis and y-axis, while the mixed learnable frequency approach introduces learnable parameters to handle diagonal directions, aiming to capture more complex spatial relationships in images.

Following works further improve the position encoding for CV in different aspects. A representative work is LaPE [28]. It identify that PEs are added to patch embeddings before entering the Transformer encoders, with both undergoing the same layer normalization. This shared normalization can restrict the expressiveness of PEs, potentially hindering model performance. LaPE involves assigning separate layer normalizations to

token embeddings and PEs within each Transformer layer, allowing each to maintain its unique characteristics. Additionally, PEs are progressively adapted and delivered across layers, providing layer-specific positional information and enabling the model to adjust to varying levels of abstraction throughout the network.

# B   Datasets

We conduct experiments on public datasets including CIFAR10 [5] and CIFAR100 [6]. Detailed statistics are shown in Tab. 3.

Table 3: Detailed statistics, including image size, categories, and data splits.

| Dataset | size | categories | training size | eval size |
|---------|------|------------|---------------|-----------|
| CIFAR10 | 32 | 10 | 50,000 | 10,000 |
| CIFAR100 | 32 | 100 | 50,000 | 10,000 |

# C   Limatations and Future Work

While our method demonstrates strong performance, a current limitation lies in its limited generalizability across different backbone architectures. Specifically, when the vision transformer backbone is changed, it may be necessary to re-pretrain the model with our proposed $SaPE^2$ as the positional embedding, which introduces additional computational cost. Nevertheless, this also opens up a promising research direction. In future work, we plan to investigate fine-tuning strategies to improve the adaptability of $SaPE^2$ and reduce the need for full retraining across diverse architectures.

# D   Codes Implementation

Below is the core implementation code of $SaPE^2$, which can be equipped within the transformer's attention calculation module. To use it, you can integrate $SaPE^2$ into the transformer's attention layer. Specifically, this can be achieved by adding the output of SaPE to the attention computation.

```python
class SaPE_unit(nn.Module):
    def __init__(self, npos_max, head_dim, scale):
        super(SaPE_unit, self).__init__()
        self.npos_max = npos_max
        self.pos_emb = nn.Parameter(torch.zeros(1, head_dim, npos_max))
        self.scale = scale
        trunc_normal_(self.pos_emb, std=.02)

    def forward(self, q, k, v):
        # query: (batch_size, heads, seq_len, head_dim)
        # attn_logits: (batch_size, heads, seq_len, seq_len)
        # compute positions

        attn_logits = (q * self.scale) @ k.transpose(-2, -1)
        gates = torch.sigmoid(attn_logits)

        pos = gates.flip(-1).cumsum(dim=-1).flip(-1)
        pos = pos.clamp(max=self.npos_max - 1)
        # interpolate from integer positions
        pos_ceil = pos.ceil()
        pos_floor = pos.floor()
        logits_int = torch.matmul(q, self.pos_emb)
        logits_ceil = logits_int.gather(-1, pos_ceil)
        logits_floor = logits_int.gather(-1, pos_floor)
        w = pos - pos_floor
        return logits_ceil * w + logits_floor * (1 - w)
```

---

[5]https://www.cs.toronto.edu/ kriz/cifar-10-python.tar.gz
[6]https://www.cs.toronto.edu/ kriz/cifar-100-python.tar.gz

```python
class SaPE(nn.Module):
    def __init__(self, npos_max, head_dim, scale, num_heads, num_patches, img_size):
        super(SaPE, self).__init__()
        self.npos_max = npos_max
        self.sape_x = SaPE_unit(npos_max, head_dim, scale)
        self.sape_y = SaPE_unit(npos_max, head_dim, scale)

    def forward(self, query, key, value):
        # query: (batch_size, heads, seq_len, head_dim)

        B, heads, seq_len, head_dim = query.shape
        W = H = int(math.sqrt(seq_len))
        query = query.reshape(B, heads, H, W, head_dim)
        key = key.reshape(B, heads, H, W, head_dim)
        value = value.reshape(B, heads, H, W, head_dim)
        # [batch, heads, patch_h, patch_w, hidden], [batch, h/2, w/2]

        # split q, k, v
        q_x = query.permute(0, 2, 1, 3, 4).reshape(B*H, heads, W, head_dim)
        q_y = query.permute(0, 3, 1, 2, 4).reshape(B*W, heads, H, head_dim)

        k_x = key.permute(0, 2, 1, 3, 4).reshape(B*H, heads, W, head_dim)
        k_y = key.permute(0, 3, 1, 2, 4).reshape(B*W, heads, H, head_dim)

        v_x = value.permute(0, 2, 1, 3, 4).reshape(B*H, heads, W, head_dim)
        v_y = value.permute(0, 3, 1, 2, 4).reshape(B*W, heads, H, head_dim)

        w_x = self.sape_x(q_x, k_x, v_x).reshape(B, H, heads, W, H)
        w_y = self.sape_y(q_y, k_y, v_y).reshape(B, W, heads, H, W)

        w_x = w_x.permute(0, 2, 1, 3, 4).reshape(B * heads, H * W, H)
        w_y = w_y.permute(0, 2, 3, 1, 4).reshape(B * heads, H * W, H)

        # Calculate the 2-dimentional relative distance
        pe_res_x = torch.cdist(w_x, w_x, p=2).reshape(B, heads, H*W, H*W)
        pe_res_y = torch.cdist(w_y, w_y, p=2).reshape(B, heads, H*W, H*W)

        return pe_res_x + pe_res_y
```