

Permutation Equivariance of Transformers and Its Applications

Supplementary Material

7. Structure of Transformer

Transformer-based models are the state-of-the-art deep neural networks and have attracted great attention in both areas of computer vision and natural language processing. Models including transformer encoder blocks as their backbone, such as Bert [6], ViT [8], T2T-ViT [35], ViTGAN [12], BEiT [2] and CoCa [34], have been achieving exceeding performance in a great many tasks.

Transformer encoder blocks, as shown in Fig. 9, mainly contain two critical components: Multi-head Scaled-dot-product self-attention and a feed-forward network (MLP). Inputs are fed in the form of patches, which are usually embedding vectors for words in Bert, or for fractions of images in ViT. The relative position of patches are learned by position embeddings [32], which are injected into the model. Fig. 9 shows the main operators in a Transformer where the shortcut and the linear projection in the Attention block are left out for simplicity.

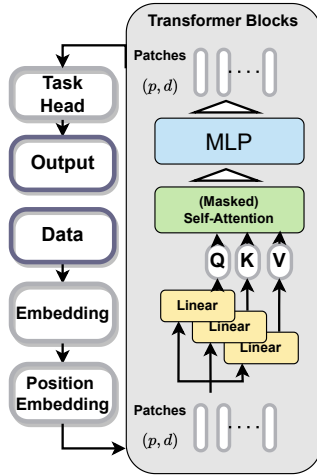


Figure 9. Transformer Encoder Block

The Transformer encoder block is denoted as Enc and the loss is ℓ . The patch embedding of a single input \mathbf{X} is expressed as \mathbf{Z} of shape (p, d) . The first layer in the self-attention contains three parallel linear layers projecting \mathbf{Z} to $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ as

$$\mathbf{Q} = \mathbf{Z}\mathbf{W}_Q^\top, \quad (14)$$

$$\mathbf{K} = \mathbf{Z}\mathbf{W}_K^\top, \quad (15)$$

$$\mathbf{V} = \mathbf{Z}\mathbf{W}_V^\top. \quad (16)$$

$\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are fed to the following attention operation

$$\mathbf{S} = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right), \quad (17)$$

$$\mathbf{A} = \mathbf{S}\mathbf{V}, \quad (18)$$

where \mathbf{S} and \mathbf{A} are the softmax output, and the attention output, respectively.

We neglect the attention projection and the residual connection for simplicity. The part following the attention layer is the MLP layer:

$$\mathbf{A}_1 = \mathbf{A}\mathbf{W}_1^\top, \quad (19)$$

$$\mathbf{H} = a(\mathbf{A}_1), \quad (20)$$

$$\mathbf{A}_2 = \mathbf{H}\mathbf{W}_2^\top \quad (21)$$

where $\mathbf{A}_1, \mathbf{A}_2$ are the outputs of the linear layers with weights $\mathbf{W}_1, \mathbf{W}_2$, respectively, and \mathbf{H} is the output of the element-wise activation function a which can be ReLu, Tanh, etc.

The backward propagation of Transformer encoder block is as following, we calculate the all the gradients from the final layer back to the first. Gradients are expressed as

$$\begin{aligned} dl &= \text{tr}\left(\frac{\partial l}{\partial \mathbf{A}_2}^\top d\mathbf{A}_2\right) \\ &= \text{tr}\left(\frac{\partial l}{\partial \mathbf{A}_2}^\top (d\mathbf{H})\mathbf{W}_2^\top\right) + \text{tr}\left(\frac{\partial l}{\partial \mathbf{A}_2}^\top \mathbf{H}d(\mathbf{W}_2^\top)\right). \end{aligned}$$

The two additive terms are inspected in the following. Let's study \mathbf{H} first:

$$\begin{aligned} dl_1 &\triangleq \text{tr}\left(\frac{\partial l}{\partial \mathbf{A}_2}^\top (d\mathbf{H})\mathbf{W}_2^\top\right) \\ &= \text{tr}\left(\mathbf{W}_2^\top \frac{\partial l}{\partial \mathbf{A}_2}^\top d\mathbf{H}\right) \\ &= \text{tr}\left(\left(\frac{\partial l}{\partial \mathbf{A}_2}^\top \mathbf{W}_2\right)^\top d\mathbf{H}\right), \end{aligned}$$

indicating

$$\frac{\partial l}{\partial \mathbf{H}} = \frac{\partial l}{\partial \mathbf{A}_2}^\top \mathbf{W}_2. \quad (22)$$

For \mathbf{W}_2 ,

$$\begin{aligned} dl_2 &\triangleq \text{tr}\left(\frac{\partial l}{\partial \mathbf{A}_2}^\top \mathbf{H}d(\mathbf{W}_2^\top)\right) \\ &= \text{tr}\left(d\mathbf{W}_2\mathbf{H}^\top \frac{\partial l}{\partial \mathbf{A}_2}^\top\right) \\ &= \text{tr}\left(\left(\frac{\partial l}{\partial \mathbf{A}_2}^\top \mathbf{H}\right)^\top d\mathbf{W}_2\right), \end{aligned}$$

and

$$\frac{\partial l}{\partial \mathbf{W}_2} = \frac{\partial l}{\partial \mathbf{A}_2}^\top \mathbf{H}. \quad (23)$$

For \mathbf{A}_1 :

$$\begin{aligned} dl_1 &= \text{tr}\left(\left(\frac{\partial l}{\partial \mathbf{H}}\right)^\top d\mathbf{H}\right) \\ &= \text{tr}\left(\frac{\partial l}{\partial \mathbf{H}}\right)^\top d(a(\mathbf{A}_1)) \\ &= \text{tr}\left(\frac{\partial l}{\partial \mathbf{H}}\right)^\top a'(\mathbf{A}_1) \odot d\mathbf{A}_1 \\ &= \text{tr}\left(\left(\frac{\partial l}{\partial \mathbf{H}}\right) \odot a'(\mathbf{A}_1)\right)^\top d\mathbf{A}_1, \end{aligned}$$

by Eq. 22, we have

$$\frac{\partial l}{\partial \mathbf{A}_1} = \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{W}_2 \odot a'(\mathbf{A}_1). \quad (24)$$

Similarly, we calculate the gradients of \mathbf{A} and \mathbf{W}_1 :

$$\frac{\partial l}{\partial \mathbf{A}} = \frac{\partial l}{\partial \mathbf{A}_1} \mathbf{W}_1, \quad (25)$$

$$\frac{\partial l}{\partial \mathbf{W}_1} = \frac{\partial l}{\partial \mathbf{A}_1}^\top \mathbf{A}. \quad (26)$$

In the attention operation:

$$\begin{aligned} dl_3 &\triangleq \text{tr}\left(\frac{\partial l}{\partial \mathbf{A}}\right)^\top d\mathbf{A} \\ &= \text{tr}\left(\frac{\partial l}{\partial \mathbf{A}}\right)^\top (d\mathbf{S})\mathbf{V} + \text{tr}\left(\frac{\partial l}{\partial \mathbf{A}}\right)^\top \mathbf{S}d\mathbf{V} \\ &= \text{tr}\left(\left(\frac{\partial l}{\partial \mathbf{A}}\mathbf{V}^\top\right)^\top d\mathbf{S}\right) + \text{tr}\left(\left(\mathbf{S}^\top \frac{\partial l}{\partial \mathbf{A}}\right)^\top d\mathbf{V}\right), \end{aligned}$$

and

$$\frac{\partial l}{\partial \mathbf{S}} = \frac{\partial l}{\partial \mathbf{A}} \mathbf{V}^\top, \quad (27)$$

$$\frac{\partial l}{\partial \mathbf{V}} = \mathbf{S}^\top \frac{\partial l}{\partial \mathbf{A}}. \quad (28)$$

First, for $\mathbf{V} = \mathbf{Z}\mathbf{W}_V^\top$:

$$\begin{aligned} dl_4 &\triangleq \text{tr}\left(\frac{\partial l}{\partial \mathbf{V}}\right)^\top d\mathbf{V} \\ &= \text{tr}\left(\frac{\partial l}{\partial \mathbf{V}}\right)^\top (d\mathbf{Z})\mathbf{W}_V^\top + \text{tr}\left(\frac{\partial l}{\partial \mathbf{V}}\right)^\top \mathbf{Z}d\mathbf{W}_V^\top. \end{aligned}$$

Similarly, the gradients of \mathbf{Z} and \mathbf{W}_V are:

$$\frac{\partial l}{\partial \mathbf{Z}} = \frac{\partial l}{\partial \mathbf{V}} \mathbf{W}_V, \quad (29)$$

$$\frac{\partial l}{\partial \mathbf{W}_V} = \frac{\partial l}{\partial \mathbf{V}}^\top \mathbf{Z}. \quad (30)$$

Now we focus on $\mathbf{S} = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right)$:

$$\begin{aligned} dl_5 &\triangleq \text{tr}\left(\frac{\partial l}{\partial \mathbf{S}}\right)^\top d\mathbf{S} \\ &= \text{tr}\left(\frac{\partial l}{\partial \mathbf{S}}\right)^\top (\text{diag}(\mathbf{S}) - \mathbf{S}^\top \mathbf{S})d\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \\ &= \text{tr}\left(\left((\text{diag}(\mathbf{S}) - \mathbf{S}^\top \mathbf{S})^\top \frac{\partial l}{\partial \mathbf{S}}\right)^\top d\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right)\right), \end{aligned}$$

and thus

$$\frac{\partial l}{\partial \mathbf{Q}} = \frac{1}{\sqrt{d}} \left((\text{diag}(\mathbf{S}) - \mathbf{S}^\top \mathbf{S})^\top \frac{\partial l}{\partial \mathbf{S}} \right) \mathbf{K}, \quad (31)$$

$$\frac{\partial l}{\partial \mathbf{K}} = \frac{1}{\sqrt{d}} \left((\text{diag}(\mathbf{S}) - \mathbf{S}^\top \mathbf{S})^\top \frac{\partial l}{\partial \mathbf{S}} \right)^\top \mathbf{Q}. \quad (32)$$

And similarly the gradients of \mathbf{W}_Q and \mathbf{W}_K are:

$$\frac{\partial l}{\partial \mathbf{W}_Q} = \frac{\partial l}{\partial \mathbf{Q}}^\top \mathbf{Z}, \quad (33)$$

$$\frac{\partial l}{\partial \mathbf{W}_K} = \frac{\partial l}{\partial \mathbf{K}}^\top \mathbf{Z}. \quad (34)$$

8. Alg. on Permuted Training

Our permuted training is described by pseudo code in Alg. 1. It should be noted that the permutation takes place not on the dimension of ‘batches’ but on the rest two dimensions. Taking ViT for example, each image is transformed into a (p, d) matrix representing p patches, and each patch denotes a fraction of the image. Each fraction is embedded into a d -dimensional vector.

We further provide a toy example. Let \mathbf{Z} of shape $(3, 4)$ and the row shuffle matrix \mathbf{P}_R be

$$\mathbf{Z} = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \end{pmatrix} \quad \mathbf{P}_R = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}.$$

The row permuted feature is

$$\mathbf{P}_R \mathbf{Z} = \begin{pmatrix} 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 1 & 2 & 3 & 4 \end{pmatrix}.$$

Column permutation is performed in a similar way. Note that permutation matrices are orthogonal, i.e., $\mathbf{P}_R^{-1} = \mathbf{P}_R^\top$ and $\mathbf{P}_C^{-1} = \mathbf{P}_C^\top$.

Algorithm 1 Permuted Training

- 1: Initialization: Initialize the model. Load permutation matrices P_R, P_C .
- 2: Start training
- 3: **repeat**
- 4: Start a new epoch
- 5: **repeat**
- 6: Get a batch of data \mathbf{X} from data loader
- 7: Get embedding \mathbf{Z} of size $(batch_size, p, d)$.
- 8: **if** using row permutation **then**
- 9: $\mathbf{Z} = \text{matmul}(P_R, \mathbf{Z})$
- 10: **end if**
- 11: **if** using column permutation **then**
- 12: $\mathbf{Z} = \text{matmul}(\mathbf{Z}, P_C)$
- 13: **end if**
- 14: Send \mathbf{Z} to the Transformer Backbone and retrieve the output $\hat{\mathbf{Y}}$
- 15: **if** using row permutation **then**
- 16: $\hat{\mathbf{Y}} = \text{matmul}(P_R^{-1}, \hat{\mathbf{Y}})$
- 17: **end if**
- 18: **if** using column permutaton **then**
- 19: $\hat{\mathbf{Y}} = \text{matmul}(\hat{\mathbf{Y}}, P_C^{-1})$
- 20: **end if**
- 21: Perform backward propagation
- 22: **until** done all batches
- 23: **until** done all epochs

9. Permutation-Equivariant Operators

As far as we will show, the following operators are permutation-equivariant:

- Element-wise operators,
- Softmax,
- Linear layer,
- MLP,
- LayerNorm and BatchNorm,
- Attention.

In the following, we will prove the permutation-equivariance of each operator.

Element-wise operators including shortcut, Hadamard product, matrix addition/subtraction and other element-wise functions. We have

Lemma 9.1. *Element-wise operators are permutation-equivariant that*

$$(P_R \mathbf{A} P_C) \odot (P_R \mathbf{B} P_C) = P_R (\mathbf{A} \odot \mathbf{B}) P_C. \quad (35)$$

where \odot denotes the element-wise operation.

On the left hand-side of the equation, a_{ij} in \mathbf{A} and b_{ij} in \mathbf{B} are permuted to the same position before being performed the operation. On the right hand-side, a_{ij} and b_{ij}

are performed the operation of which the results are permuted. The two are obviously equivariant. Lemma 9.1 also holds for matrix addition and activation function:

$$a(P_R \mathbf{A} P_C) = P_R a(\mathbf{A}) P_C \quad (36)$$

where a is an element-wise activation function, or other element-wise functions like scalar multiplication, division, etc.

Lemma 9.2. *Softmax is permutation-equivariant:*

$$\text{Softmax}(P_R \mathbf{A} P_C) = P_R \text{Softmax}(\mathbf{A}) P_C. \quad (37)$$

This is because an element is always normalized with the same group of elements, which are not changed in permutations. Thus Softmax is permutation-equivariant.

Lemma 9.3. *Linear layer is permutation-equivariant:*

$$f_{(P)}(P_R \mathbf{X} P_C) = P_R f(\mathbf{X}) P_C \quad (38)$$

where $f(\mathbf{X}) = \mathbf{X} \mathbf{W}^\top + b$ and $f_{(P)}(\mathbf{X}) = \mathbf{X} \mathbf{W}_{(P)}^\top + b_{(P)}$ and:

$$\begin{aligned} \mathbf{W}_{(P)} &= P_C^\top \mathbf{W} P_C, \\ b_{(P)} &= b P_C. \end{aligned}$$

Proof.

$$\begin{aligned} f_{(P)}(P_R \mathbf{X} P_C) &= P_R \mathbf{X} P_C \mathbf{W}_{(P)}^\top + b_{(P)} \\ &= P_R \mathbf{X} P_C \cdot P_C^\top \mathbf{W}^\top P_C + b P_C \\ &= P_R \mathbf{X} \mathbf{W}^\top P_C + b P_C \\ &= P_R f(\mathbf{X}) P_C, \end{aligned}$$

where the bias b is broadcast to each row. Note that if $P_C \neq \mathbf{I}$, the identity matrix, this lemma is limited to linear layer with square weight matrix. If P_C is not included, i.e. only row shuffle is used, then all linear layers are row-permutation equivariant. \square

Lemma 9.4. *MLP is permutation-equivariant:*

$$f_{(P)}(P_R \mathbf{X} P_C) = P_R f(\mathbf{X}) P_C \quad (39)$$

where

$$\begin{aligned} f(\mathbf{X}) &= \sigma(\mathbf{X} \mathbf{W}_1^\top + b_1) \mathbf{W}_2^\top + b_2, \\ f_{(P)}(\mathbf{X}) &= \sigma(\mathbf{X} \mathbf{W}_{1(P)}^\top + b_{1(P)}) \mathbf{W}_{2(P)}^\top + b_{2(P)}, \end{aligned}$$

and:

$$\begin{aligned} \mathbf{W}_{1(P)} &= \mathbf{W}_1 P_C, \mathbf{W}_{2(P)} = P_C^\top \mathbf{W}_2, \\ b_{1(P)} &= b_1, b_{2(P)} = b_2 P_C, \end{aligned}$$

where σ is the activation function, $\mathbf{W}_1 \in \mathbb{R}^{t \times d}$, $\mathbf{W}_2 \in \mathbb{R}^{d \times t}$, $b_1 \in \mathbb{R}^t$, $b_2 \in \mathbb{R}^d$, and t is the hidden dimension of MLP.

Proof.

$$\begin{aligned}
f_{(P)}(\mathbf{P}_R \mathbf{X} \mathbf{P}_C) &= \sigma(\mathbf{P}_R \mathbf{X} \mathbf{P}_C \mathbf{W}_{1(P)}^\top + b_{1(P)}) \mathbf{W}_{2(P)}^\top + b_{2(P)} \\
&= \sigma(\mathbf{P}_R \mathbf{X} \mathbf{W}_1^\top + b_1) \mathbf{W}_2^\top \mathbf{P}_C + b_2 \mathbf{P}_C \\
&= \mathbf{P}_R (\sigma(\mathbf{X} \mathbf{W}_1^\top + b_1) \mathbf{W}_2^\top + b_2) \mathbf{P}_C \\
&= \mathbf{P}_R f(\mathbf{X}) \mathbf{P}_C
\end{aligned}$$

where the third equation holds due to Lem. 9.1 and the broadcast of bias. \square

Lemma 9.5. *Normalization (LayerNorm for example, LN for short) is permutation-equivariant:*

$$\text{LN}_{(P)}(\mathbf{P}_R \mathbf{X} \mathbf{P}_C) = \text{LN}(\mathbf{X}) \quad (40)$$

where $\text{LN}(\mathbf{X}) = \frac{\mathbf{X} - \mathbb{E}(\mathbf{X})}{\sqrt{\text{Var}(\mathbf{X}) - \epsilon}} * \gamma + b$, and:

$$\gamma_{(P)} = \gamma \mathbf{P}_C, \quad b = b \mathbf{P}_C.$$

Since the same $\mathbb{E}(\mathbf{X})$ and $\text{Var}(\mathbf{X})$ work on each element, permutation dose not affect the normalization operation. And the affine operation is ‘column-wise’, weight γ and bias b are broadcast to each row.

Lemma 9.6. *Attention ($\mathbf{A} = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}})\mathbf{V}$) is permutation-equivariant:*

$$\begin{aligned}
\text{Attention}(\mathbf{P}_R \mathbf{Q} \mathbf{P}_C, \mathbf{P}_R \mathbf{K} \mathbf{P}_C, \mathbf{P}_R \mathbf{V} \mathbf{P}_C) \\
= \mathbf{P}_R \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \mathbf{P}_C.
\end{aligned} \quad (41)$$

Proof.

$$\begin{aligned}
&\text{Attention}(\mathbf{P}_R \mathbf{Q} \mathbf{P}_C, \mathbf{P}_R \mathbf{K} \mathbf{P}_C, \mathbf{P}_R \mathbf{V} \mathbf{P}_C) \\
&= \text{Softmax}\left(\frac{\mathbf{P}_R \mathbf{Q} \mathbf{P}_C \cdot \mathbf{P}_C^\top \mathbf{K}^\top \mathbf{P}_R^\top}{\sqrt{d}}\right) \mathbf{P}_R \mathbf{V} \mathbf{P}_C \\
&= \mathbf{P}_R \text{Softmax}\left(\frac{\mathbf{Q} \mathbf{K}^\top}{\sqrt{d}}\right) \mathbf{P}_R^\top \cdot \mathbf{P}_R \mathbf{V} \mathbf{P}_C \\
&= \mathbf{P}_R \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \mathbf{P}_C
\end{aligned}$$

where the second equality holds because of the permutation equivariance of Softmax. \square

Multihead attention is a special case. The validity of Thm. 4.4 and Thm. 4.5 is contingent on constraining the permutation \mathbf{P}_C to operate within a single head. It means permutation equivariance holds if permutation is performed within each head, or on different heads, but not across heads. In that case, the feasible permutation space shrinks but the space is still considerable. Take a base Transformer for example, the possible permutations is reduced from $768!$ to $12! \times 64!$.

We later provide detailed proofs of Thm. 4.1, Thm. 4.2, Thm. 4.4, Thm. 4.5 specifically for Transformer architecture.

10. Proofs on Transformer Encoder Blocks

We show the detailed proof on the Transformer encoder blocks. The notations are shown in Appendix 7, and Transformer Encoder Block is denoted as Enc for short.

10.1. Enc is Forward Permutation-Equivariant

Enc is forward permutation-equivariant. As proven above, all the basic operators in Enc are permutation-equivariant. The following section shows how the combination of the operators still holds in detail. The proofs of Thm. 4.1 and Thm. 4.4 are organized into one where the weight matrices are permuted by \mathbf{P}_R and \mathbf{P}_C at the same time. The row permutation equivariance can be seen as a special case where $\mathbf{P}_C = \mathbf{I}$, and the column permutation equivariance is a special case of $\mathbf{P}_R = \mathbf{I}$.

Proof. First and foremost, we ‘encrypt’ all the weight matrices by Eq. 6:

$$\mathbf{W}_{i(P)} = \mathbf{P}_C^\top \mathbf{W}_i \mathbf{P}_C,$$

where \mathbf{P}_C is the column permutation matrix, \mathbf{W}_i is the weight of a normal Enc , and $i \in \{Q, K, V\}$. Weights in MLP are ‘encrypted’ by $\mathbf{W}_{1(C)} = \mathbf{W}_1 \mathbf{P}_C$, $\mathbf{W}_{2(C)} = \mathbf{P}_C^\top \mathbf{W}_2$. We denote the Transformer encoder block with such ‘encryption’ as $\text{Enc}_{(P)}$.

For \mathbf{Q} :

$$\mathbf{Q}_{(P)} = \mathbf{Z}_{(P)} \mathbf{W}_{\mathbf{Q}(P)}^\top \quad (42)$$

$$= \mathbf{P}_R \mathbf{Z} \mathbf{P}_C \cdot \mathbf{P}_C^\top \mathbf{W}_{\mathbf{Q}}^\top \mathbf{P}_C \quad (43)$$

$$= \mathbf{P}_R \mathbf{Z} \mathbf{W}_{\mathbf{Q}}^\top \mathbf{P}_C \quad (44)$$

$$= \mathbf{P}_R \mathbf{Q} \mathbf{P}_C. \quad (45)$$

Similarly for \mathbf{K}, \mathbf{V} :

$$\mathbf{K}_{(P)} = \mathbf{P}_R \mathbf{K} \mathbf{P}_C, \quad (46)$$

$$\mathbf{V}_{(P)} = \mathbf{P}_R \mathbf{V} \mathbf{P}_C. \quad (47)$$

For $\mathbf{S} = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}})$:

$$\mathbf{S}_{(P)} = \text{Softmax}\left(\frac{\mathbf{Q}_{(P)} \mathbf{K}_{(P)}^\top}{\sqrt{d}}\right) \quad (48)$$

$$= \text{Softmax}\left(\frac{\mathbf{P}_R \mathbf{Q} \mathbf{P}_C \cdot \mathbf{P}_C^\top \mathbf{K}^\top \mathbf{P}_R^\top}{\sqrt{d}}\right) \quad (49)$$

$$= \text{Softmax}\left(\frac{\mathbf{P}_R \mathbf{Q} \mathbf{K}^\top \mathbf{P}_R^\top}{\sqrt{d}}\right) \quad (50)$$

$$= \mathbf{P}_R \text{Softmax}\left(\frac{\mathbf{Q} \mathbf{K}^\top}{\sqrt{d}}\right) \mathbf{P}_R^\top \quad (51)$$

$$= \mathbf{P}_R \mathbf{S} \mathbf{P}_R^\top. \quad (52)$$

So for \mathbf{A} :

$$\mathbf{A}_{(P)} = \mathbf{S}_{(P)} \mathbf{V}_{(P)} \quad (53)$$

$$= \mathbf{P}_R \mathbf{S} \mathbf{P}_R^\top \cdot \mathbf{P}_R \mathbf{V} \mathbf{P}_C \quad (54)$$

$$= \mathbf{P}_R \mathbf{S} \mathbf{V} \mathbf{P}_C \quad (55)$$

$$= \mathbf{P}_R \mathbf{A} \mathbf{P}_C. \quad (56)$$

Following the attention layer, \mathbf{A} is fed to the MLP layer:

$$\mathbf{A}_{1(P)} = \mathbf{A}_{(P)} \mathbf{W}_{1(P)}^\top \quad (57)$$

$$= \mathbf{P}_R \mathbf{A} \mathbf{P}_C \cdot \mathbf{P}_C^\top \mathbf{W}_1 \quad (58)$$

$$= \mathbf{P}_R \mathbf{A} \mathbf{W}_1 \quad (59)$$

$$= \mathbf{P}_R \mathbf{A}_1. \quad (60)$$

Similarly for \mathbf{A}_2 ,

$$\mathbf{A}_{2(P)} = \mathbf{P}_R \mathbf{A}_2 \mathbf{P}_C. \quad (61)$$

As for the activation in the middle, the element-wise activation function is permutation-equivariant:

$$\mathbf{H}_{(P)} = \mathbf{P}_R \mathbf{H}. \quad (62)$$

Overall, we have proved Enc satisfies permutation forward equivariance. \square

10.2. Enc is Backward Permutation-Invariant

According to Thm. 4.7, since all the operators in Enc are forward permutation-equivariant, the feature of Enc is backward permutation-equivariant and the weight in Enc is permutation-invariant. The following section shows how the combination of the operators still holds in detail. Similar to the proof of forward permutation equivariance, we prove Thm. 4.2 and Thm. 4.5 altogether in one proof where the weight matrices are permuted by \mathbf{P}_R and \mathbf{P}_C at the same time. The row permutation equivariance can be seen as a special case where $\mathbf{P}_C = \mathbf{I}$, and the column permutation equivariance is a special case of $\mathbf{P}_R = \mathbf{I}$.

Proof. Due to the shuffling and unshuffling procedures of Alg. 1, we have the forward and backward propagation outside of the backbone no different from the normal ones. Hence we only focus on the propagation of the Transformer encoder blocks.

We denote $\mathbf{A}_{3(P)}$ as the reversed intermediate feature that the down-stream head receives:

$$\mathbf{A}_{3(P)} = \mathbf{P}_R^\top \mathbf{A}_{2(P)} \mathbf{P}_C^\top. \quad (63)$$

Since the feature is unshuffled, we have

$$\mathbf{A}_{3(P)} = \mathbf{A}_3 = \mathbf{A}_2. \quad (64)$$

First, we focus on the MLP layer:

$$\begin{aligned} dl &= \text{tr} \left(\frac{\partial l}{\partial \mathbf{A}_{3(P)}} \mathbf{P}_R^\top d(\mathbf{A}_{2(P)}) \mathbf{P}_C^\top \right) \\ &= \text{tr} \left(\mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{A}_{3(P)}} \mathbf{P}_R^\top d\mathbf{A}_{2(P)} \right) \\ &= \text{tr} \left(\left(\mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}_{3(P)}} \mathbf{P}_C \right)^\top d\mathbf{A}_{2(P)} \right), \end{aligned}$$

that is:

$$\frac{\partial l}{\partial \mathbf{A}_{2(P)}} = \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}_{3(P)}} \mathbf{P}_C = \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{P}_C \quad (65)$$

by Eq. 64.

With $\mathbf{H}_{(P)} = \mathbf{P}_R \mathbf{H} \mathbf{P}_C^\top$ and Eq. 23, the gradient:

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{W}_{2(P)}} &= \frac{\partial l}{\partial \mathbf{A}_{2(P)}} \mathbf{H}_{(P)} \\ &= \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{P}_R^\top \cdot \mathbf{P}_R \mathbf{H} \\ &= \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{H} \\ &= \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{W}_2}, \end{aligned}$$

that is:

$$\frac{\partial l}{\partial \mathbf{W}_{2(P)}} = \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{W}_2}. \quad (66)$$

By Eq. 24 and Eq. 66, we have

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{A}_{1(P)}} &= \frac{\partial l}{\partial \mathbf{A}_{2(P)}} \mathbf{W}_{2(P)} \odot a'(\mathbf{A}_{1(P)}) \\ &= \left[\mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{P}_C \cdot \mathbf{P}_C^\top \mathbf{W}_2 \right] \odot \left[\mathbf{P}_R a'(\mathbf{A}_1) \right] \\ &= \left[\mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{W}_2 \right] \odot \left[\mathbf{P}_R a'(\mathbf{A}_1) \right] \\ &= \mathbf{P}_R \left[\frac{\partial l}{\partial \mathbf{A}_2} \mathbf{W}_2 \odot a'(\mathbf{A}_1) \right] \\ &= \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}_1}, \end{aligned}$$

that is:

$$\frac{\partial l}{\partial \mathbf{A}_{1(P)}} = \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}_1}. \quad (67)$$

The weight $\mathbf{W}_{1(P)}$ in the MLP has the following gradient by Eq. 26:

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{W}_{1(P)}} &= \frac{\partial l}{\partial \mathbf{A}_{1(P)}} \mathbf{A}_{(P)} \\ &= \frac{\partial l}{\partial \mathbf{A}_1} \mathbf{P}_R^\top \cdot \mathbf{P}_R \mathbf{A} \mathbf{P}_C \\ &= \frac{\partial l}{\partial \mathbf{W}_1} \mathbf{P}_C, \end{aligned}$$

that is:

$$\frac{\partial l}{\partial \mathbf{W}_{1(P)}} = \frac{\partial l}{\partial \mathbf{W}_1} \mathbf{P}_C. \quad (68)$$

And we come to the attention operation, from Eq. 25, we have

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{A}_{(P)}} &= \frac{\partial l}{\partial \mathbf{A}_{1(P)}} \mathbf{W}_{1(P)} \\ &= \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}_1} \mathbf{W}_1 \mathbf{P}_C \\ &= \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}} \mathbf{P}_C, \end{aligned}$$

that is:

$$\frac{\partial l}{\partial \mathbf{A}_{(P)}} = \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}} \mathbf{P}_C. \quad (69)$$

Hence we observe the permutation rules for the gradients of the intermediate-layer outputs vary from the gradients of the weights. As for the gradients of the softmax-layer output, we have

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{S}_{(P)}} &= \frac{\partial l}{\partial \mathbf{A}_{(P)}} \mathbf{V}_{(P)}^\top \\ &= \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}} \mathbf{P}_C \cdot \mathbf{P}_C^\top \mathbf{V}^\top \mathbf{P}_R^\top \\ &= \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}} \mathbf{V}^\top \mathbf{P}_R^\top \\ &= \mathbf{P}_R \frac{\partial l}{\partial \mathbf{S}} \mathbf{P}_R^\top, \end{aligned}$$

that is:

$$\frac{\partial l}{\partial \mathbf{S}_{(P)}} = \mathbf{P}_R \frac{\partial l}{\partial \mathbf{S}} \mathbf{P}_R^\top. \quad (70)$$

Since $\mathbf{S}_{(P)}$ follows Eq. 52, we have the gradients for

$\mathbf{Q}_{(P)}$ combining with Eq. 70:

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{Q}_{(P)}} &= \frac{1}{\sqrt{d}} [(\text{diag}(\mathbf{S}_{(P)}) - \mathbf{S}_{(P)}^\top \mathbf{S}_{(P)}) \frac{\partial l}{\partial \mathbf{S}_{(P)}}] \mathbf{K}_{(P)} \\ &= \frac{1}{\sqrt{d}} [(\mathbf{P}_R \text{diag}(\mathbf{S}) \mathbf{P}_R^\top - \mathbf{P}_R \mathbf{S}^\top \mathbf{P}_R^\top \cdot \mathbf{P}_R \mathbf{S} \mathbf{P}_R^\top) \\ &\quad \cdot \mathbf{P}_R \frac{\partial l}{\partial \mathbf{S}} \mathbf{P}_R^\top] \mathbf{P}_R \mathbf{K} \mathbf{P}_C^\top \\ &= \frac{1}{\sqrt{d}} [(\mathbf{P}_R \text{diag}(\mathbf{S}) \mathbf{P}_R^\top - \mathbf{P}_R \mathbf{S}^\top \mathbf{S} \mathbf{P}_R^\top) \\ &\quad \cdot \mathbf{P}_R \frac{\partial l}{\partial \mathbf{S}} \mathbf{P}_R^\top] \mathbf{P}_R \mathbf{K} \mathbf{P}_C^\top \\ &= \frac{1}{\sqrt{d}} [\mathbf{P}_R (\text{diag}(\mathbf{S}) - \mathbf{S}^\top \mathbf{S}) \\ &\quad \cdot \mathbf{P}_R^\top \cdot \mathbf{P}_R \frac{\partial l}{\partial \mathbf{S}} \mathbf{P}_R^\top] \mathbf{P}_R \mathbf{K} \mathbf{P}_C \\ &= \frac{1}{\sqrt{d}} [\mathbf{P}_R (\text{diag}(\mathbf{S}) - \mathbf{S}^\top \mathbf{S}) \frac{\partial l}{\partial \mathbf{S}} \mathbf{P}_R^\top] \mathbf{P}_R \mathbf{K} \mathbf{P}_C \\ &= \frac{1}{\sqrt{d}} \mathbf{P}_R (\text{diag}(\mathbf{S}) - \mathbf{S}^\top \mathbf{S}) \frac{\partial l}{\partial \mathbf{S}} \mathbf{P}_R^\top \cdot \mathbf{P}_R \mathbf{K} \mathbf{P}_C \\ &= \mathbf{P}_R \frac{1}{\sqrt{d}} (\text{diag}(\mathbf{S}) - \mathbf{S}^\top \mathbf{S}) \frac{\partial l}{\partial \mathbf{S}} \mathbf{K} \mathbf{P}_C \\ &= \mathbf{P}_R \frac{\partial l}{\partial \mathbf{Q}} \mathbf{P}_C. \end{aligned}$$

By a similar derivation on \mathbf{K} we obtain:

$$\frac{\partial l}{\partial \mathbf{K}_{(P)}} = \mathbf{P}_R \frac{\partial l}{\partial \mathbf{K}} \mathbf{P}_C. \quad (71)$$

Following a similar proof to the gradients of $\mathbf{W}_{1(P)}$ or $\mathbf{W}_{2(P)}$, we could easily derive:

$$\frac{\partial l}{\partial \mathbf{W}_{Q(P)}} = \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{W}_Q} \mathbf{P}_C, \quad (72)$$

$$\frac{\partial l}{\partial \mathbf{W}_{K(P)}} = \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{W}_K} \mathbf{P}_C. \quad (73)$$

By Eq. 28, the gradient of $\mathbf{V}_{(P)}$ is

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{V}_{(P)}} &= \mathbf{S}_{(P)}^\top \frac{\partial l}{\partial \mathbf{A}_{(P)}} \\ &= \mathbf{P}_R \mathbf{S} \mathbf{P}_R^\top \cdot \mathbf{P}_R \frac{\partial l}{\partial \mathbf{A}} \mathbf{P}_C \\ &= \mathbf{P}_R \frac{\partial l}{\partial \mathbf{V}} \mathbf{P}_C, \end{aligned}$$

and thus we have

$$\frac{\partial l}{\partial \mathbf{V}_{(P)}} = \mathbf{P}_R \frac{\partial l}{\partial \mathbf{V}} \mathbf{P}_C, \quad (74)$$

$$\frac{\partial l}{\partial \mathbf{W}_{V(P)}} = \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{W}_V} \mathbf{P}_C. \quad (75)$$

So far, we have proved the rule for the gradient of weight matrices:

$$\frac{\partial l}{\partial \mathbf{W}_{i(P)}} = \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{W}_i} \mathbf{P}_C, \quad i \in \{Q, K, V\}. \quad (76)$$

$$\frac{\partial l}{\partial \mathbf{W}_{1(P)}} = \frac{\partial l}{\partial \mathbf{W}_1} \mathbf{P}_C, \quad \frac{\partial l}{\partial \mathbf{W}_{2(P)}} = \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{W}_2}. \quad (77)$$

$\mathbf{W}_{i(P)}$ are the weights of $\text{Enc}_{(P)}$ while \mathbf{W}_i are the weights of Enc . By induction, we can reach the conclusion that if a Transformer encoder block is randomly initialized and trained with $\mathbf{Z}_{(P)}$, it would eventually learn to become $\text{Enc}_{(P)}$, the weights of which are associated with Enc by Eq. 76 and Eq. 77. The proof of backward equivariance on the linear projection in the attention is omitted as its proof is similar and the conclusion is the same with Eq. 76. Hence we have proved backward permutation in/equi-variance. \square

10.3. Proofs on Embeddings

We show in this section that the parameters of the embedding layer F_1 , including the position embeddings, are the same despite Alg. 1 is applied or not.

Theorem 10.1. *The parameters of F_1 trained with or without permutation are the same.*

Proof. We denote the output of F_1 as \mathbf{Z}_0 , and the input of the Transformer backbone as \mathbf{Z} . In normal setting (Eq. 1), \mathbf{Z} equals to \mathbf{Z}_0 , and so do their gradients. In the permuted setting where we use subscript (P) to denote all the variables, the input to the backbone is the permuted output of $F_{1(P)}$:

$$\mathbf{Z}_{(P)} = \mathbf{P}_R \mathbf{Z}_{0(P)} \mathbf{P}_C. \quad (78)$$

To prove the weights of $F_{1(P)}$ is equivalent to those of F_1 , we need to prove:

$$\frac{\partial l}{\partial \mathbf{Z}_{0(P)}} = \frac{\partial l}{\partial \mathbf{Z}_0}. \quad (79)$$

It is clear that

$$\begin{aligned} dl &= \text{tr} \left(\frac{\partial l}{\partial \mathbf{Z}_{(P)}}^\top d\mathbf{Z}_{(P)} \right) \\ &= \text{tr} \left(\mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{Z}}^\top \mathbf{P}_R^\top \mathbf{P}_R d\mathbf{Z}_{0(P)} \mathbf{P}_C \right) \\ &= \text{tr} \left(\mathbf{P}_C \mathbf{P}_C^\top \frac{\partial l}{\partial \mathbf{Z}}^\top d\mathbf{Z}_{0(P)} \right) \\ &= \text{tr} \left(\frac{\partial l}{\partial \mathbf{Z}}^\top d\mathbf{Z}_{0(P)} \right), \end{aligned}$$

where the second equality holds by Lemma 4.8. Hence,

$$\frac{\partial l}{\partial \mathbf{Z}_{0(P)}} = \frac{\partial l}{\partial \mathbf{Z}} = \frac{\partial l}{\partial \mathbf{Z}_0}. \quad (80)$$

The invariance of the weights of F_1 can be derived from Eq. 80. \square

It is worth noting that the position embeddings are added before the permutation, so Transformer gets the right position information, since Transformer modeling the position by the position embeddings instead of the order of the input.

10.4. Proofs on Masked Attention

Masked attention is an essential component in the Transformer Decoder Block which is the backbone of the generative language model [3, 27]. The token permutation can hardly pass through the nonlinear effect of the masked attention, but the column permutation cancels out before the mask takes effect. Thus for the masked attention, we apply column permutations only.

Theorem 10.2. *The results of Masked Softmax with or without column permutation are the same:*

$$MS_{(P)} = MS. \quad (81)$$

Proof. For $MS = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}})$:

$$MS_{(P)} = \text{Masked} - \text{Softmax} \left(\frac{\mathbf{Q}_{(P)} \mathbf{K}_{(P)}^\top}{\sqrt{d}} \right) \quad (82)$$

$$= \text{Masked} - \text{Softmax} \left(\frac{\mathbf{Q} \mathbf{P}_C \cdot \mathbf{P}_C^\top \mathbf{K}^\top}{\sqrt{d}} \right) \quad (83)$$

$$= \text{Masked} - \text{Softmax} \left(\frac{\mathbf{Q} \mathbf{K}^\top}{\sqrt{d}} \right) \quad (84)$$

$$= \text{Masked} - \text{Softmax} \left(\frac{\mathbf{Q} \mathbf{K}^\top}{\sqrt{d}} \right) \quad (85)$$

$$= MS. \quad (86)$$

\square

11. Detailed Experimental Setup

11.1. Training Setup

In fine-tuning ViT-Base on Cifar10, an Adam optimizer with a fixed learning rate 10^{-4} is used. The model is trained for 5 epochs with a cross-entropy loss. In fine-tuning Bert and GPT2 on IMDB classification, an Adam optimizer with a fixed learning rate 10^{-5} is used. The models are trained for 2 epochs.

For the ‘unauthorized’ and train-from-scratch setting on Cifar10 of ViT, the default training setting of ViT is used. Random data augmentation and cosine scheduler are added and the model is trained for 200 epochs till convergence.

The text generation in the properties validation experiments is zero-shot learning on huggingface pre-trained GPT2. The model is fine-tuned with Adam optimizer and learning rate 10^{-5} for 1 epoch.

For the CelebA attribute classification task in the ‘privacy-preserving split learning’ experiments, we adopt `timm` pre-trained model `vit_base_patch16_224` (ViT-Base) to transfer to a 40-binary-attribute classification task. A SGD optimizer is used with a cosine scheduler, for which the (initial, final) learning rate are set to $(0.05, 2 \times 10^{-4})$ and $(5 \times 10^{-4}, 2 \times 10^{-6})$ for the classification head and the encoder blocks, respectively.

11.2. Threat Model of Privacy-Preserving Split Learning

The threat model consists with [15, 33]. We assume the edge holds a private training data set $\mathbb{D}_{train} = \{X, Y\}$, where X are the private data and Y are the private labels. The edge aims at training a model with the assistance of the cloud, yet without exposing the private data. The cloud possesses powerful computing power and is honest-but-curious, meaning it obeys the protocol and performs the learning task accordingly, but is curious about the private data. The edge selects a model and splits it into three parts: F_1, Enc, F_2 . F_1, F_2 are parts close to the input layer and the output layer, respectively, which are sufficiently lightweight to deploy on the edge, whereas Enc is the major part run in the cloud.

Referring to the loss function as L_{task} and the local privacy-preserving method as M , the ultimate goal is to jointly train F_1, Enc, F_2 to

$$\underset{F_1, \text{Enc}, F_2}{\text{minimize}} L_{task}(F_2(\text{Enc}(F_1(X))), Y), \quad (87)$$

without the edge revealing X, Y to the cloud. Accessing $F_1(X)$ should not permit the cloud to infer about X . The cloud could launch black-box inversion attacks:

Black-box attackers collect the auxiliary data set X_{aux} and the corresponding features under protection mechanism M as $M(F_1(X_{aux}))$, maybe over multiple training rounds. The attacker trains an inversion model G over $(X_{aux}, M(F_1(X_{aux})))$ to invert the raw input from features [7, 10, 24] by

$$\underset{G}{\text{minimize}} L_{atk}(G(M(F_1(X_{aux}))), X_{aux}). \quad (88)$$

The loss L_{atk} can be the mean square error (MSE) between the reconstructed input \tilde{X}_{aux} and X_{aux} . At convergence, G works as a decoder to invert features into inputs. We adopt the MAE decoder G with base-size Transformer backbone and a Tanh activation layer, pre-trained on ImageNet, as the model inversion model. We train G with an AdamW optimizer with a learning rate of 10^{-4} for 30 epochs.