

Appendix

A. Proof of Permutation Equivariance of Group Shuffle Attention

Lemma 1 (Permutation matrix and permutation function). $\forall \Lambda \in \mathbb{R}^{N \times N}$, \forall permutation matrix P of size N , $\exists p : \{1, 2, \dots, N\} \rightarrow \{1, 2, \dots, N\}$ is a permutation function:

$$\Lambda_{ij} = (P \cdot \Lambda)_{p(i)p(j)} = (\Lambda \cdot P^T)_{ip(j)} = (P \cdot \Lambda \cdot P^T)_{p(i)p(j)}. \quad (15)$$

Lemma 2. Given $A \in \mathbb{R}^{N \times N}$, \forall permutation matrix P of size N ,

$$\text{softmax}(PAP^T) = P\text{softmax}(A)P^T. \quad (16)$$

Proof. $\text{softmax}(A)_{ij} = \frac{e^{A_{ij}}}{\sum_{n=1}^N e^{A_{in}}}$.

Consider a permutation function p for P , using Eq. 15, we get:

$$\begin{aligned} (P \cdot \text{softmax}(A) \cdot P^T)_{p(i)p(j)} &= \text{softmax}(A)_{ij} \\ &= \frac{e^{A_{ij}}}{\sum_{n=1}^N e^{A_{in}}} \\ &= \frac{e^{(PAP^T)_{p(i)p(j)}}}{\sum_{n=1}^N e^{(PAP^T)_{p(i)p(n)}}} \\ &= \text{softmax}(P \cdot A \cdot P^T)_{p(i)p(j)}. \end{aligned}$$

which implies $P \cdot \text{softmax}(A) \cdot P^T = \text{softmax}(P \cdot A \cdot P^T)$. \square

Lemma 3. Non-linear self-attention $\text{Attn}_\sigma(Q, X) = S(Q, X) \cdot \sigma(X)$, with $S(Q, X) = \text{softmax}(QX^T/\sqrt{c})$ is permutation-equivariant.

Proof. For a self-attention, $Q = X$.

σ is an element-wise function, thus $\sigma(P \cdot X) = P \cdot \sigma(X)$.

$$\begin{aligned} \text{Attn}_\sigma(PX, PX) &= S(PX, PX) \cdot \sigma(PX) \\ &= \text{softmax}(P \cdot XX^T \cdot P^T / \sqrt{c}) \cdot P \cdot \sigma(X) \\ &= P \cdot \text{softmax}(XX^T / \sqrt{c}) \cdot P^T \cdot P \cdot \sigma(X) \\ &= P \cdot \text{softmax}(XX^T / \sqrt{c}) \cdot \sigma(X) \\ &= P \cdot \text{Attn}_\sigma(X, X). \end{aligned}$$

which implies non-linear self-attention is permutation-equivariant. \square

Proposition 1. The Group Shuffle Attention operation is permutation-equivariant, i.e., given input $X \in \mathbb{R}^{N \times c}$, \forall permutation matrix P of size N ,

$$\text{GSA}(P \cdot X) = P \cdot \text{GSA}(X).$$

Proof.

$$\text{GSA}(X) = \mathcal{GN}(\psi(\text{GroupAttn}(X)) + X), \quad (11)$$

where,

$$\begin{aligned} \text{GroupAttn}(X) &= \\ \text{concat}\{\text{Attn}_\sigma(X_i, X_i) \mid X_i = X^{(i)} W_i\}_{i=1, \dots, g}. \end{aligned} \quad (8)$$

GSA only introduces element-wise operations, which does not change the permutation-equivariance of Attn_σ . \square

B. Proof of Permutation Invariance of Gumbel Subset Sampling

Proposition 1. The Gumbel Subset Sampling operation is permutation-invariant, i.e., given input $X \in \mathbb{R}^{N \times c}$, \forall permutation matrix P of size N ,

$$\text{GSS}(P \cdot X) = \text{GSS}(X).$$

Proof. According to Eq. 14,

$$\text{GSS}(X) = \text{gumbel_softmax}(WX^T) \cdot X, \quad W \in \mathbb{R}^{N \times c}.$$

Similar to Lemma 2,

$$\text{softmax}(AP^T) = \text{softmax}(A)P^T. \quad (17)$$

Since gumbel_softmax add element-wise operations on softmax , it does not change the permutation property. In this way,

$$\begin{aligned} \text{GSS}(P \cdot X) &= \text{gumbel_softmax}(WX^T P^T) \cdot PX \\ &= \text{gumbel_softmax}(WX^T) P^T \cdot PX \\ &= \text{gumbel_softmax}(WX^T) \cdot X \\ &= \text{GSS}(X). \end{aligned}$$

\square