

1. MOTIVATION

Learning Sentence Representation

- Transformer translation systems rely on the position index of each word to learn a positional embedding to encode order information.
- Positional embeddings are then added to corresponding word embeddings to form a input representation.
- The input representation allows Transformer translation systems to parallel learn the sentence representation with order dependencies.

Several Observations

- The input representation only involves static order dependencies based on discrete numerical information. That is, any word in the vocabulary has the same positional embedding on the same position index.
- In addition, the order dependencies encoded by the existing positional embedding are independent of word contents, which may further hinder the improvement of translation capacity.

2. RECURRENT POSITIONAL EMBEDDING

- The sequence of word vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_J\}$ are split into $\{\mathbf{x}_1^p, \dots, \mathbf{x}_J^p\}$ and $\{\mathbf{x}_1^r, \dots, \mathbf{x}_J^r\}$ and their dimensions are d_p and d_r ($d_{model}=d_p+d_r$), respectively.
- An RNN with a non-linear projection layer is used to learn its recurrent positional embedding (RPE) \mathbf{r}_j for each word over $\{\mathbf{x}_1^r, \dots, \mathbf{x}_J^r\}$, and gains a sequence of RPE $\mathbf{R}=\{\mathbf{r}_1, \dots, \mathbf{r}_J\}$:

$$\mathbf{r}_j = \text{Tanh}(f_{\text{RNN}}(\mathbf{x}_j^r, \mathbf{r}_{j-1})\mathbf{W}^r + \mathbf{b}^r). \quad (1)$$

- The traditional positional encoding is used to learn a reduced dimension input representation $\mathbf{P}=\{\mathbf{p}_1, \dots, \mathbf{p}_J\}$ for other sub-sequence $\{\mathbf{x}_1^p, \dots, \mathbf{x}_J^p\}$:

$$\mathbf{pe}_{(j,2i)} = \sin(j/10000^{2i/d_p}), \quad (2)$$

$$\mathbf{pe}_{(j,2i+1)} = \cos(j/10000^{2i/d_p}).$$

$$\mathbf{p}_j = \mathbf{x}_j^p + \mathbf{pe}_j. \quad (3)$$

3. METHODS

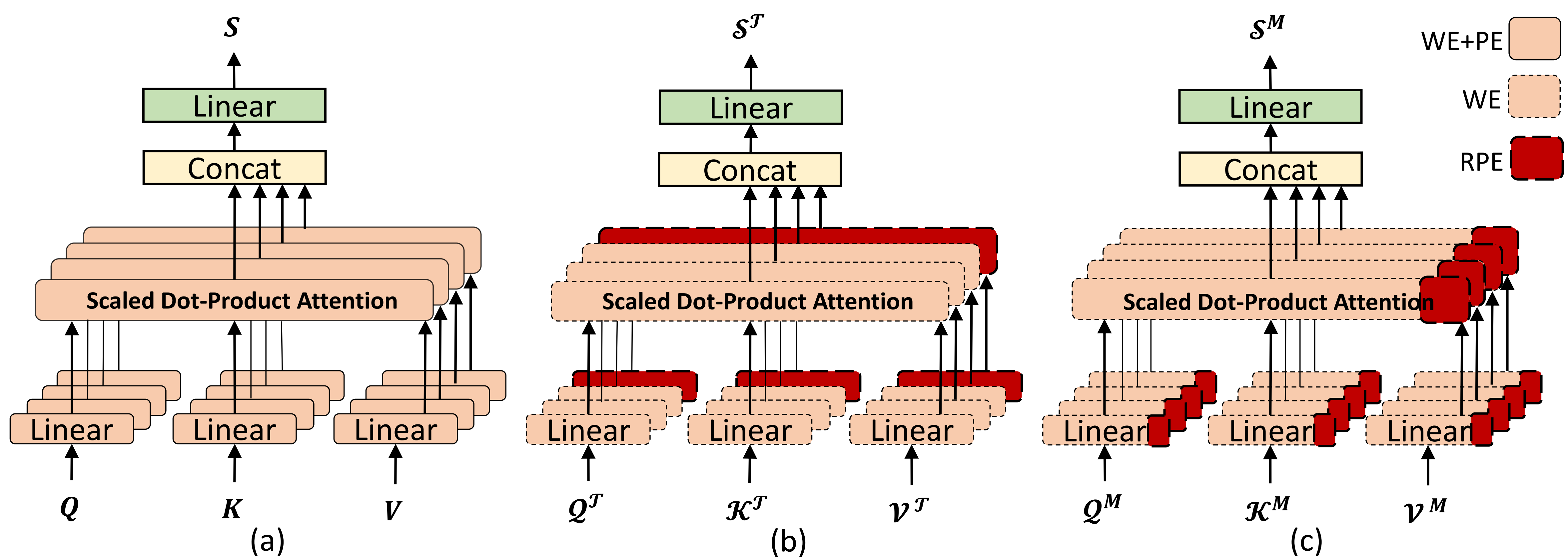


Fig 1: (a) Vanilla Multi-Head Self-Attention; (b) RPEHead Self-Attention; (c) MPRHead Self-Attention.

RPEHead Self-Attention

- These learned RPEs are integrated into multi-head self-attention (Fig 1(a)) as several independent heads to learn the sentence representation, as shown in Fig 1(b). To guarantee that there are two types of heads over the combined input representation \mathcal{T} , d_r and d_p are set to t^*d_{model}/H and $d_{model}-d_r$, respectively.

$$\begin{aligned} \mathcal{T} &= \{\mathbf{P} : \mathbf{R}\}, \\ \mathcal{O}_h^{\mathcal{T}} &= \text{Att}(\mathcal{Q}^{\mathcal{T}} \mathbf{W}_h^{\mathcal{Q}}, \mathcal{K}^{\mathcal{T}} \mathbf{W}_h^{\mathcal{K}}, \mathcal{V}^{\mathcal{T}} \mathbf{W}_h^{\mathcal{V}}), \\ \mathcal{S}^{\mathcal{T}} &= \text{Concat}(\mathcal{O}_1^{\mathcal{T}}, \dots, \mathcal{O}_{H-1}^{\mathcal{T}}, \mathcal{O}_H^{\mathcal{T}}) \mathbf{W}^{\mathcal{O}}, \end{aligned} \quad (4)$$

where $\{\mathcal{Q}^{\mathcal{T}}, \mathcal{K}^{\mathcal{T}}, \mathcal{V}^{\mathcal{T}}\}$ depends on \mathcal{T} .

MPRHead Self-Attention

- Compared with the RPEHead with independent RPE Head, the MPRHead mixed \mathbf{P} and \mathbf{R} into a combined input representation \mathbf{M} for learning the sentence representation \mathcal{S}^M , as shown in Fig 1(c).

$$\begin{aligned} \mathbf{m}_j &= \{\mathbf{p}_j^1 : \mathbf{r}_j^1, \dots, \mathbf{p}_j^H : \mathbf{r}_j^H\}, \\ \mathbf{M} &= \{\mathbf{m}_1, \dots, \mathbf{m}_J\}, \\ \mathcal{O}_h^M &= \text{Att}(\mathcal{Q}^M \mathbf{W}_h^{\mathcal{Q}}, \mathcal{K}^M \mathbf{W}_h^{\mathcal{K}}, \mathcal{V}^M \mathbf{W}_h^{\mathcal{V}}), \\ \mathcal{S}^M &= \text{Concat}(\mathcal{O}_1^M, \dots, \mathcal{O}_{H-1}^M, \mathcal{O}_H^M) \mathbf{W}^{\mathcal{O}}, \end{aligned} \quad (5)$$

where $\{\mathcal{Q}^M, \mathcal{K}^M, \mathcal{V}^M\}$ depends on \mathbf{M} .

EXPERIMENTS

Main Results

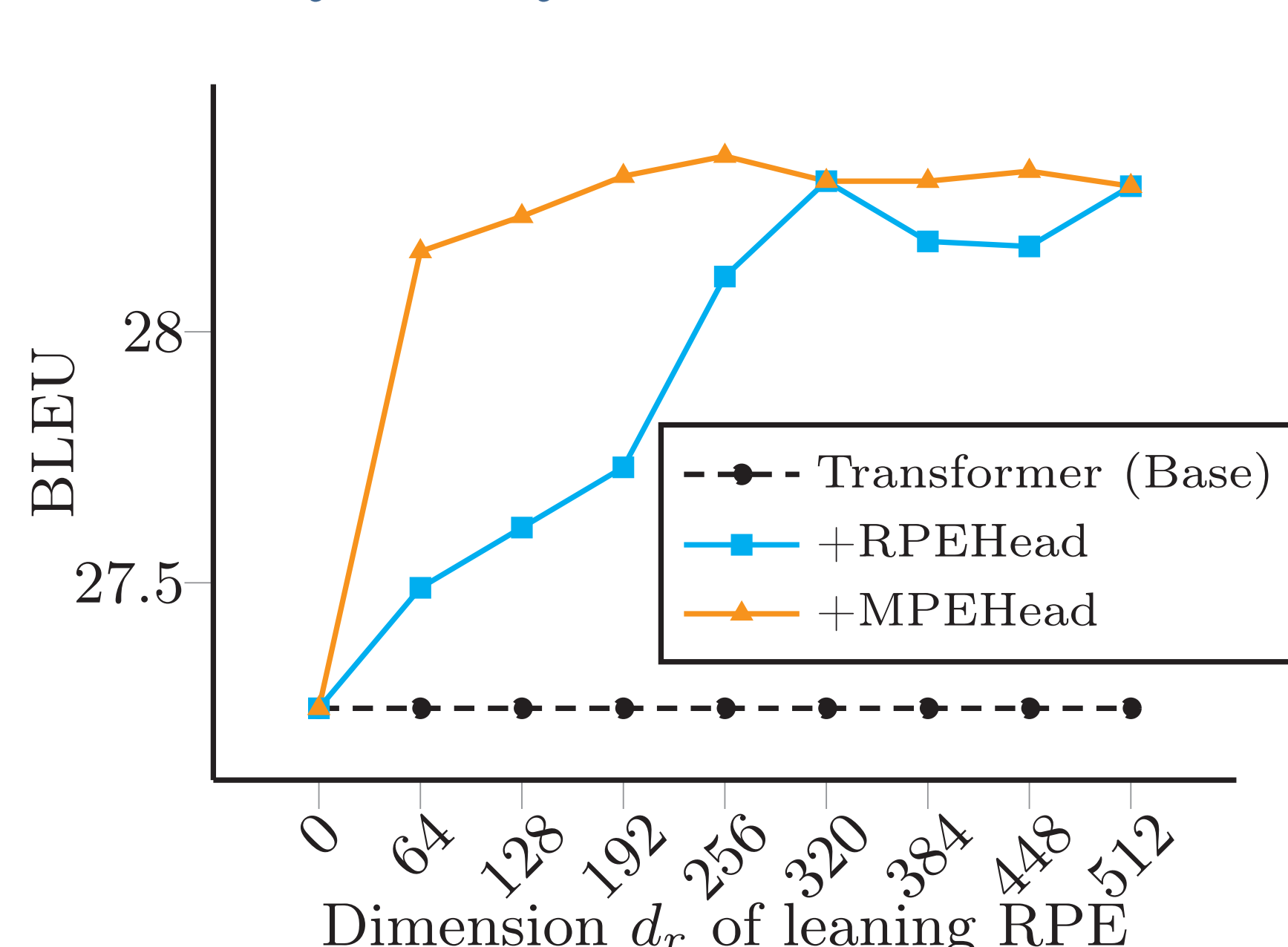
Model	NIST ZH-EN				WMT14 EN-DE	
	MT02	MT03	MT04	#Param	newstst14	#Param
Transformer(Base)	46.62	45.90	45.95	77.94	27.25	97.35M
+Relative PEs	46.91†	46.47†	45.96†	77.01M†	27.60†	97.42M
+RPEHead	47.06†	46.21†	46.32†	78.04M†	28.11†	97.84M
+MPRHead	47.69†	47.27†	47.31†	78.08†	28.35†	97.72M
Transformer(Big)	47.88	47.58	47.37	243.7M	28.22	272.6M
+MPRHead	48.73†	48.61†	48.21†	245.5M †	29.11†	289.1M

“†” indicates statistically better than Transformer(Base/Big)

($\rho < 0.01$).

“#Param” denotes the trainable parameter size of NMT model.

Efficiency Analysis for RPEs



Ablation Experiments

	newstest2014	#Speed
Transformer (Base)	27.25	11.7K
-PE of Dec	26.73	11.7K
-PE of Enc	10.71	11.7K
RPEHead (base)	28.11	10.5K
-RPE&PE of Dec	27.54	11.2K
-RPE&PE of Enc	11.16	10.7K
MPRHead (base)	28.35	9.9K
-RPE&PE of Dec	27.74	11.2K
-RPE/PE of Enc	10.89	10.7K
-RPE&PE of Enc&Dec	10.43	11.7K

“#Speed” denotes the training speed (tokens/second).

