

# Rope (Rotary pos Encod) : Linearly token pos assign e.g.  $t_1, t_2, t_3, \dots, t_n \rightarrow n$  sequence

REPO (Content Reposition) : Dynamically token pos assign

Two Components:

## ① Pos Representation

$$r_i = \text{Swish}(h_i W_g) \odot (h_i W_c) \quad * R^{dp} = d/8 \text{ for efficiency}$$

$r_i \xrightarrow{R^{dp}}$

$h_i \rightarrow$  hidden Emb  
 $W_g \rightarrow$  weight matrix for gating mechanism [ $\in \mathbb{R}^{d \times dp}$ ]  
 $W_c \rightarrow$  Content pathway [ $\in \mathbb{R}^{d \times dp}$ ]  
 $\odot \rightarrow$  hadamard prod / element-wise multi

$W_g$  # when a particular token gets reorganized in Content, how that token will influence other tokens, that is learned by gating weight matrix ( $W_g$ ). It will score that token based on this, if score is high it will make it attend forward in pos, vice-versa.

$$\text{gate}_i[j] = (h_i W_g)[j] \cdot \sigma(h_i W_g)[j]$$

$\xrightarrow{\text{Each dim } j \in [1, d_p]}$  (sigmoid) values stay in  $[0, 1]$  when inputs are neg.

$W_c$  # After the scoring by the gate, the Content matrix gives pos signal. If the words semantically similar they will group/cluster together even if seq distant. It might learn event causality rather than chronological order like repositioning the "reason" before the "consequence".

$$h_i W_c : \mathbb{R}^d \rightarrow \mathbb{R}^{dp}$$

Semantic  $\rightarrow$  Pos state

## ② Pos Assignment

$$z_i = r_i W_z \quad W_z \in \mathbb{R}^{dp \times 1}$$

$r_i$  gives a diff pos hypothesis for each  $dp$ :

Dim 1: How imp is this token?  
 Dim 2: Is this a sub or obj?  
 and so on ...

$W_z$  learns to combine these hypotheses into a single pos val:

$$z_i = \sum_{j=1}^{dp} r_i[j] \cdot W_z[j]$$

Component ① + ② :

$$f_\phi(h_i) = \left[ \text{Swish}(h_i W_g) \odot (h_i W_c) \right] W_z$$

## REPO Attn:

$$A_{ij}^{\text{REPO}} = q_i^T g_\phi(f_\phi(h_j) - f_\phi(h_i)) k_j$$

$$= \left[ q_i^T g_\phi(z_j - z_i) k_j \right] \checkmark$$