

APR

self attention: $A^{(1)}, A^{(2)}, \dots, A^{(n)}$

[MHA]: for loop (parallelization over SA)

$A(q, k, v)$ = attention-based vector representation of a word
calculate for each word $\rightarrow A^{(1)}, \dots, A^{(n)}$

Jane visit Africa on September

RNN attention

Transformer Attention

$$\alpha(t, t') = \frac{\exp(e^{q \cdot k})}{\sum_{t=1}^T \exp(e^{q \cdot k})}$$

$$A(q, k, v) = \sum_i \frac{\exp(q \cdot k^{(i)})}{\sum_j \exp(q \cdot k^{(j)})} v^{(i)}$$

Q

K

V

Query

Key reason value

$q^{(1)}$ $k^{(1)}$ action $v^{(1)}$

$q^{(2)}$ $k^{(2)}$ $v^{(2)}$

$q^{(3)}$ $k^{(3)}$ $v^{(3)}$

$q^{(4)}$ $k^{(4)}$ $v^{(4)}$

$q^{(5)}$ $k^{(5)}$ $v^{(5)}$

$q^{(2)}k^{(2)}v^{(2)}$ $q^{(4)}k^{(4)}v^{(4)}$

$q^{(1)}k^{(1)}v^{(1)}$ $q^{(3)}k^{(3)}v^{(3)}$ $q^{(5)}k^{(5)}v^{(5)}$

Jane $x^{(1)}$

visit $x^{(2)}$

Africa $x^{(3)}$

on September $x^{(4)}$

$x^{(5)}$

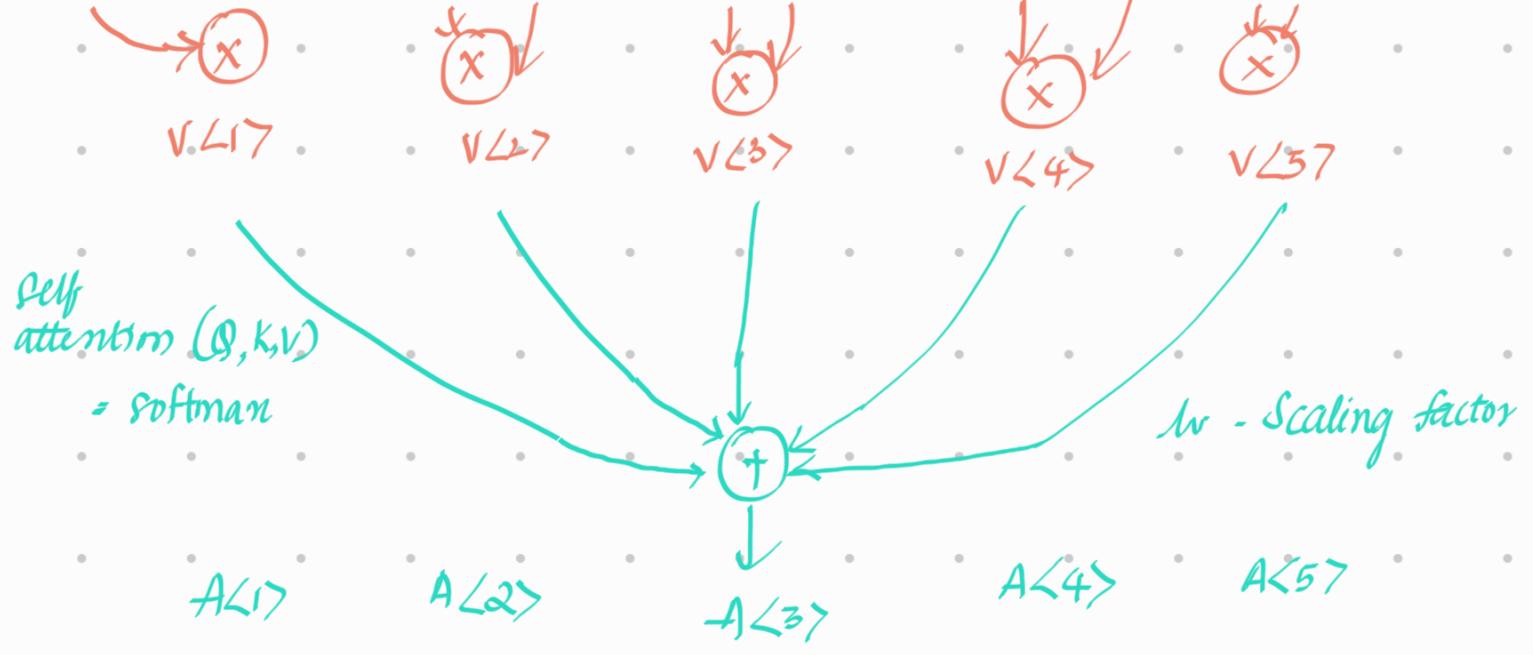
$q^{(3)}k^{(1)}$

$q^{(3)}k^{(2)}$

$q^{(3)}k^{(3)}$

$q^{(3)}k^{(4)}$

$q^{(3)}k^{(5)}$



$$A(Q, K, V) = \sum_i \frac{\exp(Q \cdot K^T \langle i \rangle)}{\sum_j \exp(Q \cdot K^T \langle j \rangle)} V \langle i \rangle$$

self attention (or) scaled dot product attention

MHA : multihead attention

Q, K, V

