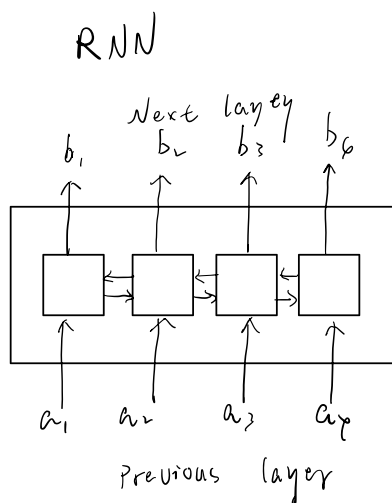
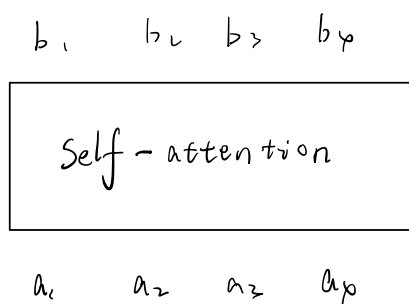


Transformer $\xrightarrow{\text{unsupervised}}$ BERT

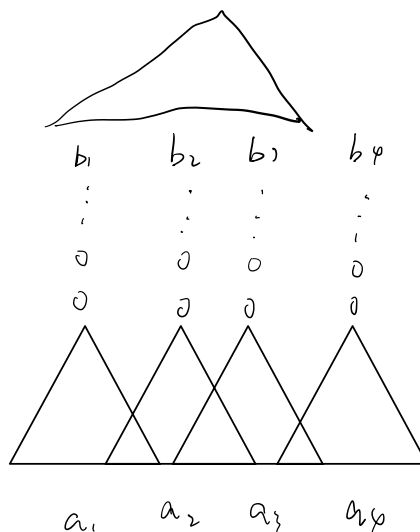
Sequence to sequence



Hard to parallel

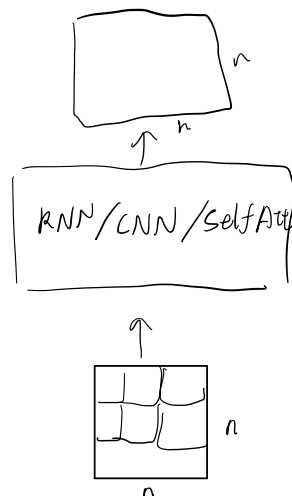


b_i is obtained based on the whole input sequence
can be parallelly computed



GNN replace RNN

视野受限, 多层才能看全



parallel computed

include whole info

b_2

b_3

b_6

b_i

$$b_i = \sum_j d_{ij} v_j$$

$\hat{a}_{i1} \rightarrow \otimes$

$\hat{a}_{i2} \rightarrow \otimes$

$\hat{a}_{i3} \rightarrow \otimes$

$\hat{a}_{i6} \rightarrow \otimes$

soft-max

\hat{a}_{i1}

\hat{a}_{i2}

\hat{a}_{i3}

\hat{a}_{i6}

$$a_{ij} = q^i \cdot k^j / \sqrt{d}$$

q^1, k^1, v^1

q^2, k^2, v^2

q^3, k^3, v^3

q^6, k^6, v^6

$$\begin{aligned} q^i &= W^q a^i \\ k^i &= W^k a^i \\ v^i &= W^v a^i \end{aligned}$$

a_1

a_2

a_3

a_6

$$a^i = W x^i$$

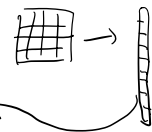
x^1

x^2

x^3

x^6

How to parallelly computed?



$$\begin{bmatrix} q^1 & q^2 & q^3 & q^p \end{bmatrix} = W^q \begin{bmatrix} a^1 & a^2 & a^3 & a^p \end{bmatrix}$$

$\underbrace{\quad}_{256 \times 1} \quad \underbrace{\quad}_{256 \times 1} \quad \underbrace{\quad}_{256 \times 1} \quad \underbrace{\quad}_{256 \times 1} \quad \underbrace{\quad}_{256 \times 256} \quad \underbrace{\quad}_{256 \times 1} \quad \underbrace{\quad}_{256 \times 1} \quad \underbrace{\quad}_{256 \times 1} \quad \underbrace{\quad}_{256 \times 1}$

$$a = W^q I$$

$$K = W^k I$$

$$v = W^v I$$

$$a_{i1} = k^1 q^1$$

$$a_{i2} = k^2 q^1$$

$$a_{i3} = k^3 q^1$$

$$a_{ip} = k^p q^1$$

$$\begin{bmatrix} a_{i1} \\ a_{i2} \\ a_{i3} \\ a_{ip} \end{bmatrix} = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \\ k^p \end{bmatrix} q^1$$

$$a^{ij} \in \mathbb{R}^{1 \times 1}$$

$$k^i \in \mathbb{R}^{16 \times 1}$$

$$q^j \in \mathbb{R}^{16 \times 1}$$

$$\hat{a}_{11} \quad \hat{a}_{12} \quad \hat{a}_{13} \quad \hat{a}_{1p}$$

$$\hat{a}_{21} \quad \hat{a}_{22} \quad \hat{a}_{23} \quad \hat{a}_{2p}$$

$$\hat{a}_{31} \quad \hat{a}_{32} \quad \hat{a}_{33} \quad \hat{a}_{3p}$$

$$\hat{a}_{p1} \quad \hat{a}_{p2} \quad \hat{a}_{p3} \quad \hat{a}_{pp}$$

$$\hat{A}$$

softmax

$$a_{11} \quad a_{12} \quad a_{13} \quad a_{1p}$$

$$a_{21} \quad a_{22} \quad a_{23} \quad a_{2p}$$

$$a_{31} \quad a_{32} \quad a_{33} \quad a_{3p}$$

$$a_{p1} \quad a_{p2} \quad a_{p3} \quad a_{pp}$$

$$A$$

$$k_1$$

$$k_2$$

$$k_3$$

$$k_6$$

$$K^T$$

$$Q$$

$$\bullet q^1 \quad q^2 \quad q^3 \quad q^p$$

$$b_1 \quad b_2 \quad b_3 \quad b_p = v_1 \quad v_2 \quad v_3 \quad v^p$$

$$Q$$

$$v$$

$$\hat{a}_{11} \quad \hat{a}_{12} \quad \hat{a}_{13} \quad \hat{a}_{1p}$$

$$\hat{a}_{21} \quad \hat{a}_{22} \quad \hat{a}_{23} \quad \hat{a}_{2p}$$

$$\hat{a}_{31} \quad \hat{a}_{32} \quad \hat{a}_{33} \quad \hat{a}_{3p}$$

$$\hat{a}_{p1} \quad \hat{a}_{p2} \quad \hat{a}_{p3} \quad \hat{a}_{pp}$$

$$\hat{A}$$

$$b_1 \quad b_2 \quad b_3 \quad b_p \Rightarrow 0$$

Self-attention

$$a_1 \quad a_2 \quad a_3 \quad a_p \Rightarrow I$$

$$\begin{aligned} Q &= W^q I \\ K &= W^k I \\ V &= W^v I \end{aligned}$$

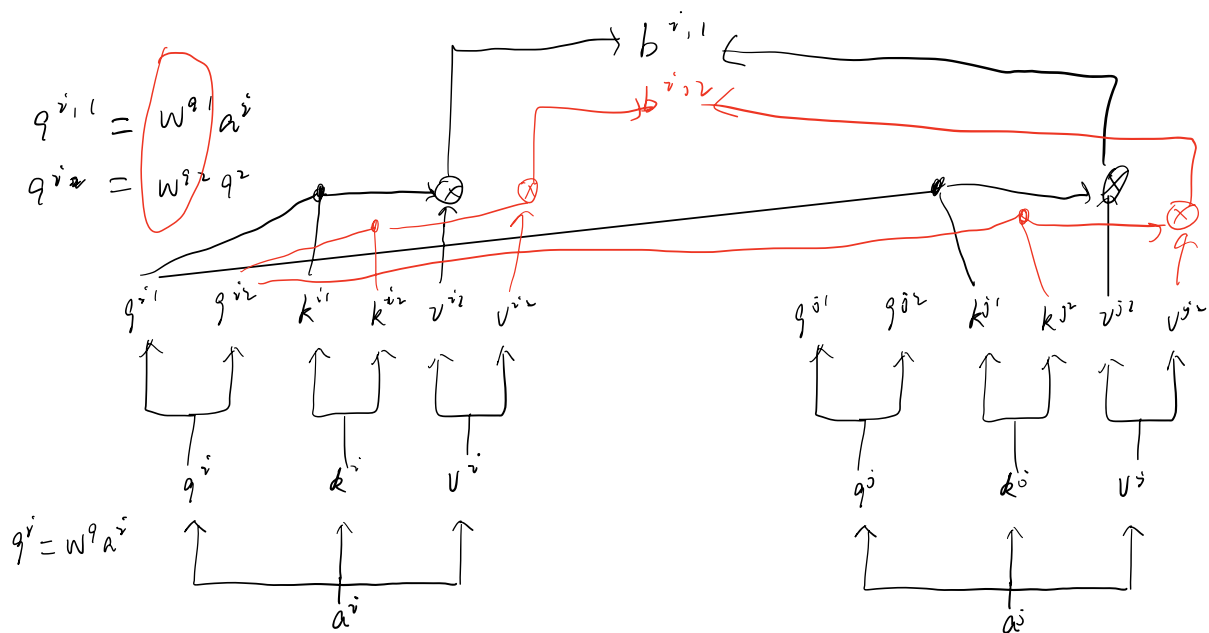
$$\hat{A} \Leftarrow A = K^T Q$$

$$O = V \hat{A}$$

Multi-head self-attention

不同 head

不同关注点 (local vs long time)



$$b^i = W^o \begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix} \quad \text{降维}$$

Position Encoding

1. no position info in self-attention
2. each position has a unique position vector e^i not learned from data

