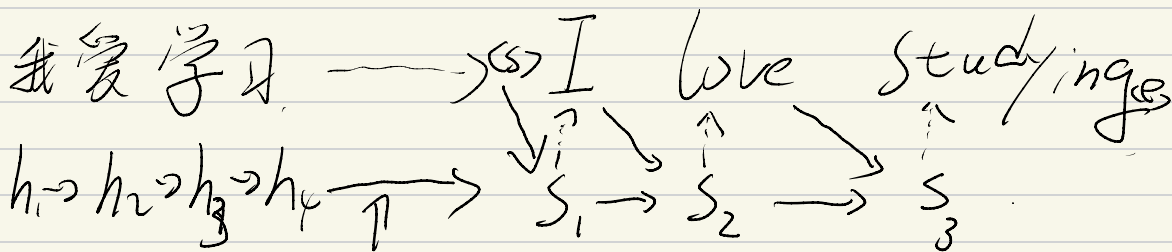


25.02.11

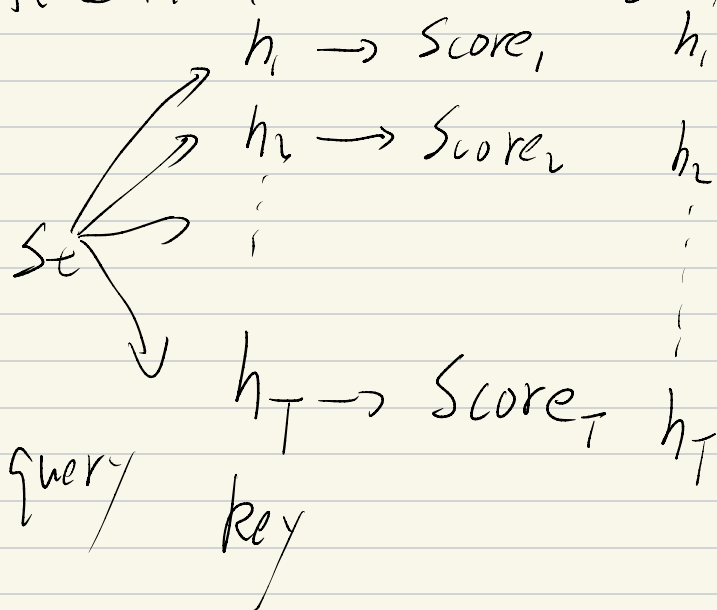
seq2seq



③ O(seq. length) Information Bottleneck

② Recurrent structure is not parallelizable

Attention Mechanism



$$\alpha_t = \frac{1}{\sum_{i=1}^T \text{Score}_i} \text{Score}_i h_i$$

Attention is All U Need

我 爱 学 习

$$x_1 \ x_2 \ x_3 \ x_4 \in \mathbb{R}^d$$

$$q_i = W_q \cdot x_i \quad ; \quad k_i = W_k \cdot x_i \quad , \quad v_i = W_v \cdot x_i$$

query key value

$$W_q, W_k, W_v \in \mathbb{R}^{d \times d}$$

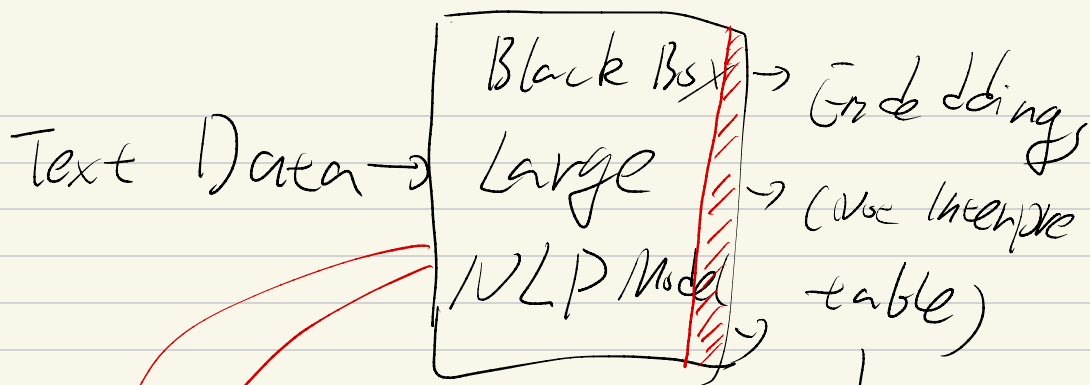
$$\tilde{w}_{ij} = \frac{q_i^T k_j}{\sqrt{d}} \quad , \quad i, j \in \{1, 2, \dots, n\}$$

↓ softmax

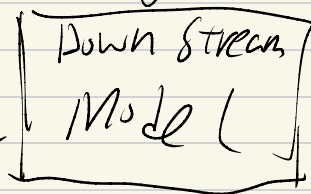
$$w_{ij} = \frac{\exp(\tilde{w}_{ij})}{\sum_{l=1}^n \exp(\tilde{w}_{il})}$$

$$y_i = \sum_{j=1}^n w_{ij} v_j$$

- ① Parallelization
v
- ② Information Bottleneck
- ③ $O(\text{seq. length})$



Causal



- ① Off-the-Shelf
- ② Fine-tuning
- ③ pre-training