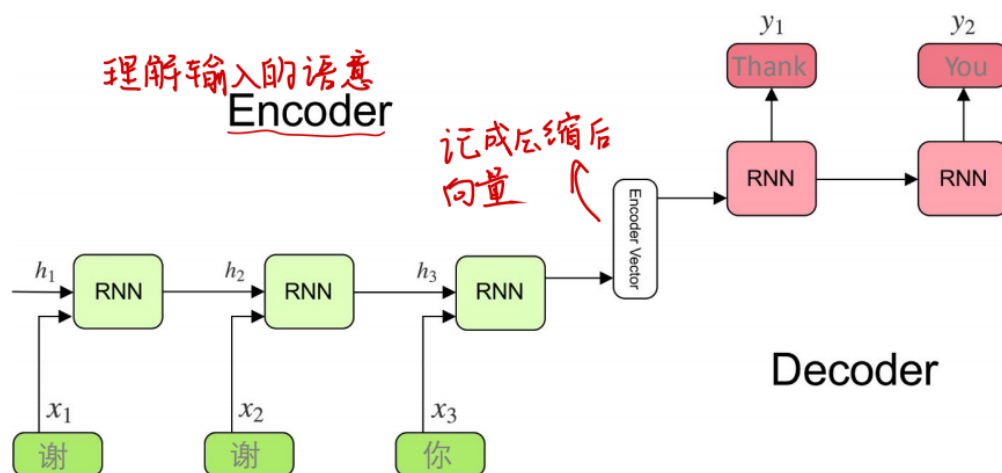# ch12: Natural Language Processing

## RNN Encoder-Decoder

### overview(P6)

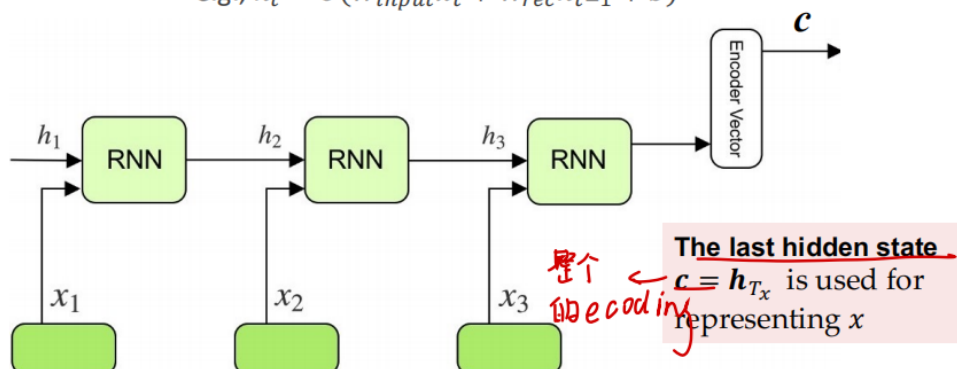- given $x = (x_1, \cdots, x_{T_x})$, generate $y = (y_1, \cdots, y_{T_y})$



### Encoder(P8)

- An RNN which learns a dense representation of a sequence.
- **Compresses** a sequence of tokens into a context vector c.

压缩

$$h_t = \text{RNN}_{in}(x_t, h_{t-1})$$ 信息压缩

e.g., $h_t = \sigma(W_{input}x_t + W_{rec}h_{t-1} + b)$



The last hidden state
整个 ← $c = h_{T_x}$ is used for
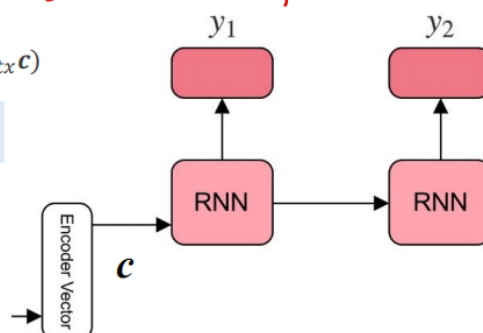做decoding representing $x$

## Decoder(P9)

- An RNN which generates an sequence **conditioned on** the intermediate representation.
逐个生成下一个token y
- Sequentially predicts the next token $y_i$ given the context vector $c$ and the hidden state of past-generated sequence.

也要看 encoder 的输入

$$s_i = \text{RNN}_{out}(y_{i-1}, s_{i-1}, c)$$

e.g., $s_i = \sigma(W'_{input}y_{i-1} + W'_{rec}s_{i-1} + W_{ctx}c)$

$$p(y_i | y_{<i}, c) = \text{softmax}(g(s_i))$$

e.g., $g = V^T(W_1 y_{t-1} + W_2 s_t + W_3 c)$
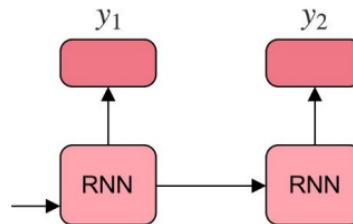


## Training(P10)

**Data:** $\{(x^{(1)}, y^{(1)}), \ldots, (x^{(N)}, y^{(N)})\}$,

where $x^{(\ell)} = (x_1, \ldots x_{Tx})$ and $y^{(\ell)} = (y_1, \ldots y_{Ty})$.

**Loss Function** – minimize the **cross-entropy** loss:

$$L(\theta) = -\frac{1}{N} \sum_{\ell=1}^{N} \sum_{t=1}^{Ty} \log p_\theta(y_t^{(\ell)} | x^{(\ell)})$$

**Optimization** – gradient descend

## Applications(P12-14)

- Translation
- Conversation
- Image Captioning
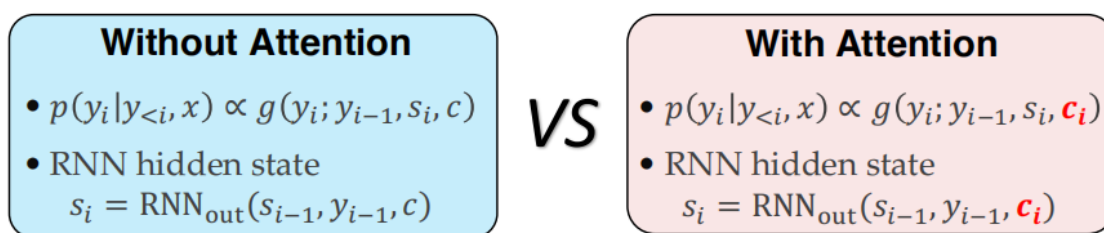- **Example: Chatbot(Hierarchical RNN Encoder-Decoder)P14**

# Attention

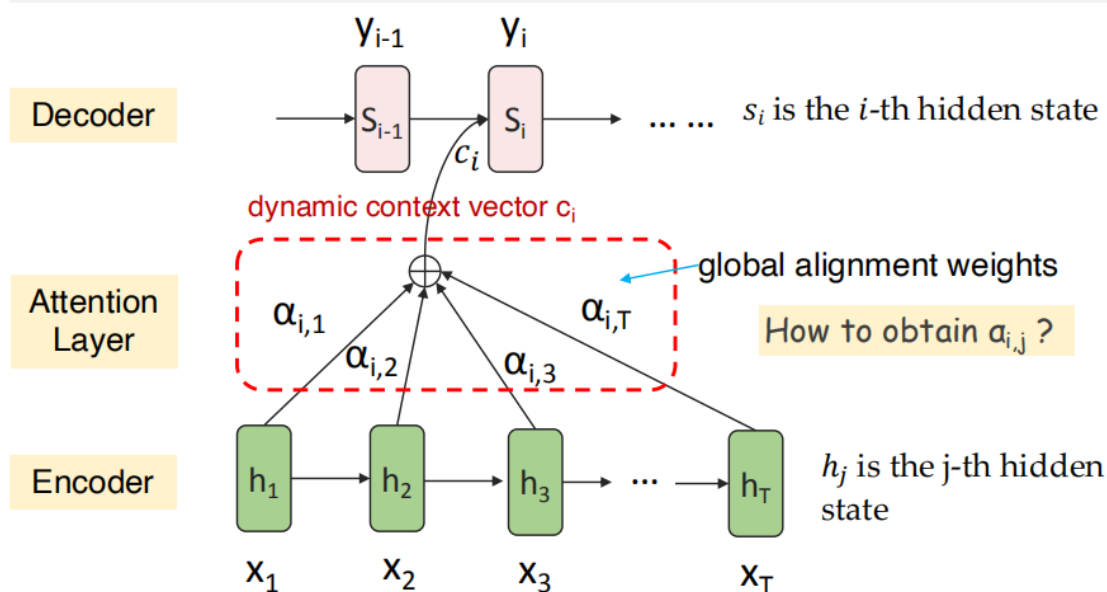## Sequence-to-Sequence with Attention(P18-19)

- When decoding each $y_i$ in $y = (y_1, \cdots, y_{T_y})$, we use a **dynamic context vector** $c_i$ which corresponds to a linear combination of different positions in $x$.

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

编码之后输入的信息

| Without Attention | VS | With Attention |
|---|---|---|
| • $p(y_i \vert y_{<i}, x) \propto g(y_i; y_{i-1}, s_i, c)$ <br> • RNN hidden state <br> $\quad s_i = \text{RNN}_{\text{out}}(s_{i-1}, y_{i-1}, c)$ | | • $p(y_i \vert y_{<i}, x) \propto g(y_i; y_{i-1}, s_i, \boldsymbol{c_i})$ <br> • RNN hidden state <br> $\quad s_i = \text{RNN}_{\text{out}}(s_{i-1}, y_{i-1}, \boldsymbol{c_i})$ |

The decoder **dynamically** pays attention to different tokens in the source sentences during decoding.



$s_i$ is the $i$-th hidden state

global alignment weights

How to obtain $a_{i,j}$ ?

$h_j$ is the j-th hidden state

**AEention-based Model(P19)**

$$\alpha_{ij} = \frac{\exp\left(a(s_{i-1}, h_j)\right)}{\sum_{k=1}^{Tx} \exp\left(a(s_{i-1}, h_k)\right)}$$

- scores how well the inputs around position $j$ and the output at position $i$ match.
- where $a(\cdot)$ denotes a **neural network**:
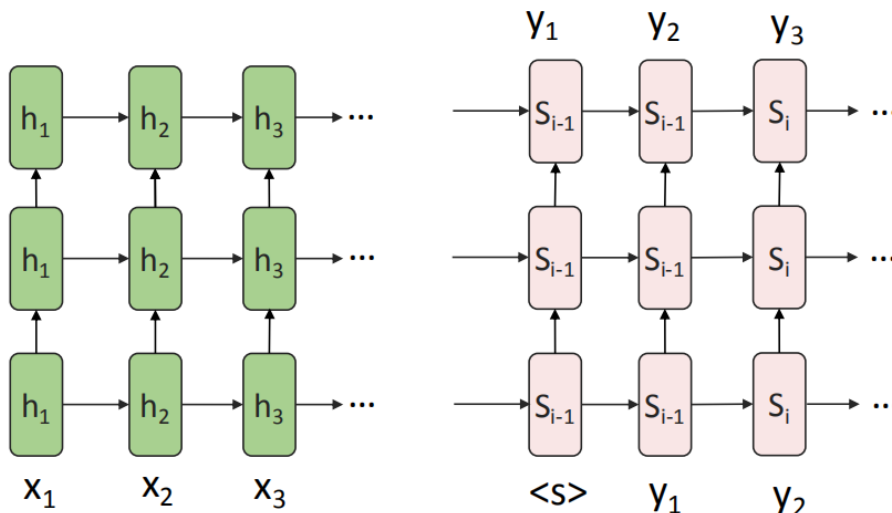  e.g., $a(s_{i-1}, h_j) = v_a^T \tanh\left(W_a s_{i-1} + U_a h_j\right)$



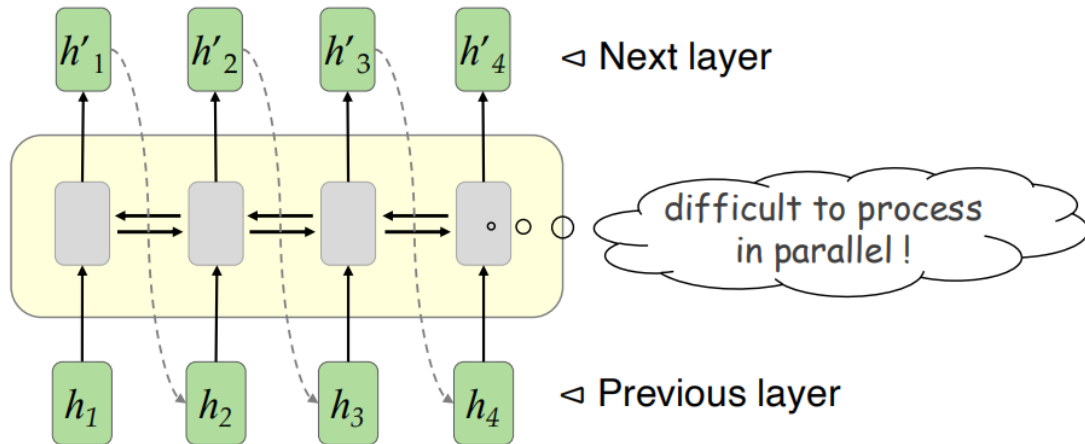## Transformer: Seq2seq model with "Self-attention"
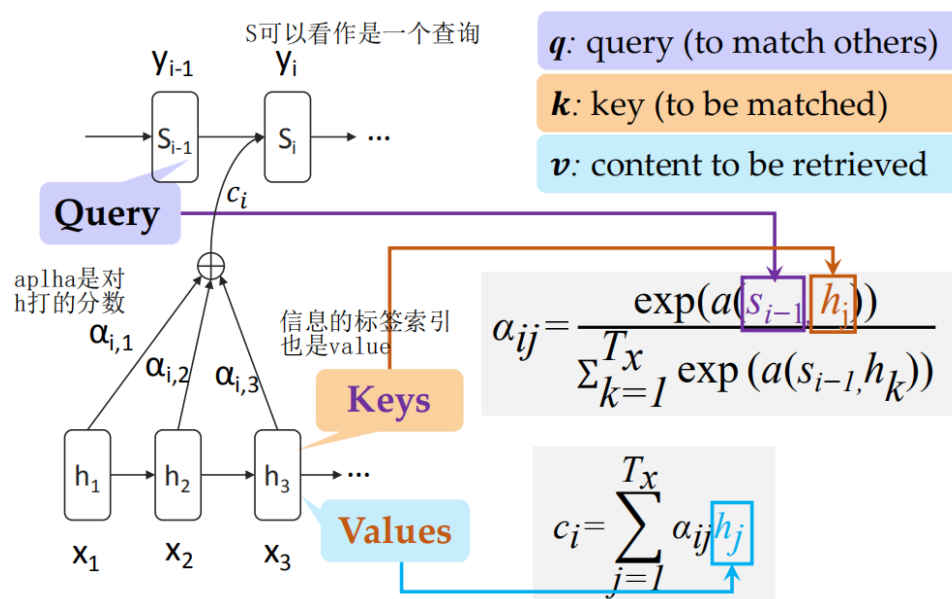
### Limitations of RNN（P22-23）

- Multilayer RNNs    问题：无法并行化，全部串行

- **A problem for the RNN encoder/decoder:** the hidden state for one position is dependent on the computation of the preceding position.



## Attention Revisited（P26）



$q$: query (to match others)

$k$: key (to be matched)

$v$: content to be retrieved

$$\alpha_{ij} = \frac{\exp(a(s_{i-1}, h_j))}{\sum_{k=1}^{T_x} \exp(a(s_{i-1}, h_k))}$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

## Self-Attention: The Idea(P27-30)

> **self-attention: 在理解任何单词的语义时，同时查看其他所有的单词，并对语义进行总结**

- **Let each word pay attention to all other words.**

- Multiplying the query vector by each key vector produces a score for each value (technically: dot product followed by softmax)
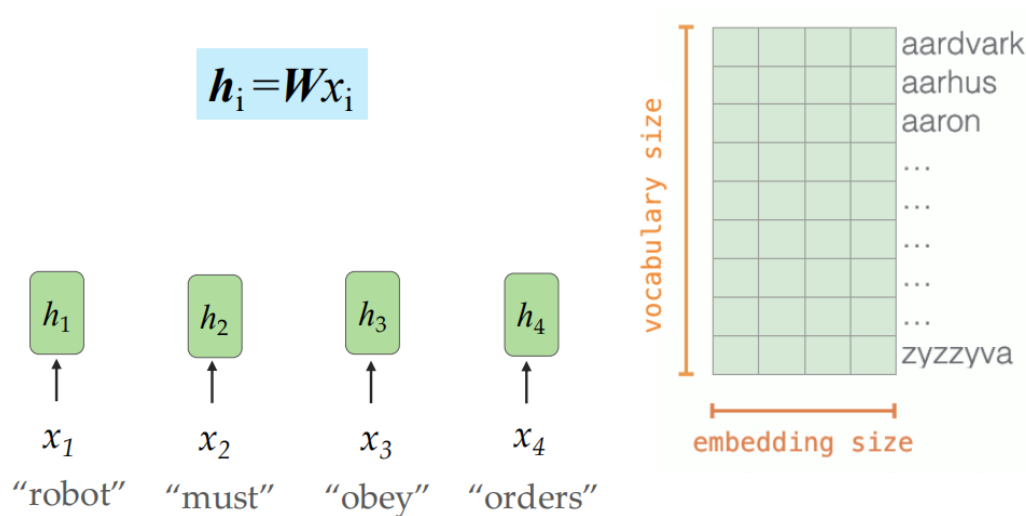
- We multiply each value by its score and sum them up.

| Word | Value vector | Score | Value X Score |
|---|---|---|---|
| \<s\> | | 0.001 | |
| a | | 0.3 | |
| robot | | 0.5 | |
| must | | 0.002 | |
| obey | | 0.001 | |
| the | | 0.0003 | |
| orders | | 0.005 | |
| given | | 0.002 | |
| it | | 0.19 | |
| | | | |
| | | Sum: | |

做完内积后每个词都获得一个分数，表示和 it 的相关程度

- The outcome vector represents the new state (refreshed memory) for the query word.
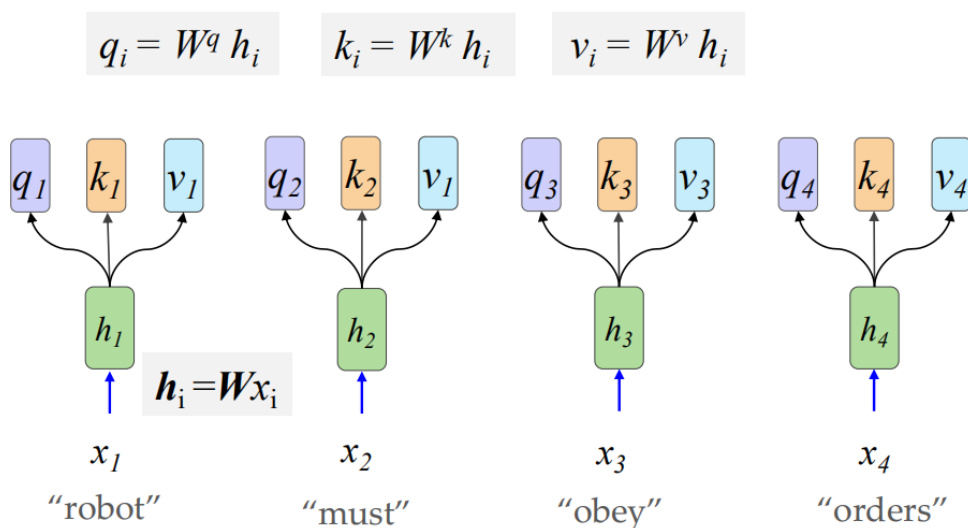
## Step 1: Token Embedding(P31)

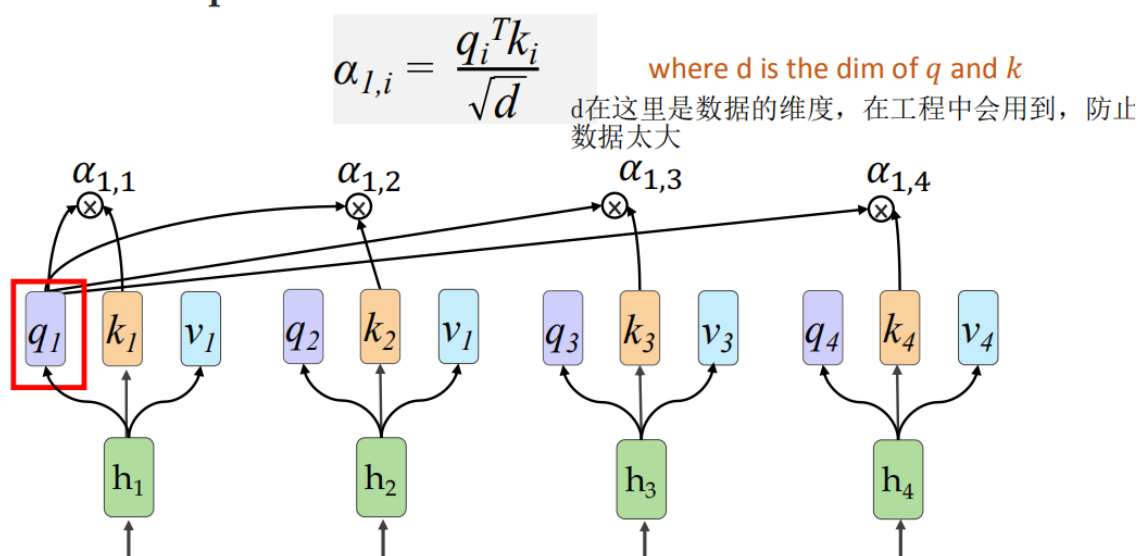- Embedding tokens (integer id) into vectors (hidden states).

$$h_i = W x_i$$



## Step 2 : Q, K, V vectors(P32)

- Transform each hidden state into **query/key/value** vectors:

$$q_i = W^q h_i \qquad k_i = W^k h_i \qquad v_i = W^v h_i$$

$q_1$ $k_1$ $v_1$ $\qquad$ $q_2$ $k_2$ $v_1$ $\qquad$ $q_3$ $k_3$ $v_3$ $\qquad$ $q_4$ $k_4$ $v_4$

$h_1$ $\qquad$ $h_2$ $\qquad$ $h_3$ $\qquad$ $h_4$

$h_i = Wx_i$

$x_1$ $\qquad$ $x_2$ $\qquad$ $x_3$ $\qquad$ $x_4$

"robot" $\qquad$ "must" $\qquad$ "obey" $\qquad$ "orders"

## Step 3: Calculate Attention Scores(P33)

- Calculate an attention score for each <query, key> pair using **scaled dot-product**.

$$\alpha_{1,i} = \frac{q_i^T k_i}{\sqrt{d}}$$

where d is the dim of $q$ and $k$

d在这里是数据的维度，在工程中会用到，防止数据太大

$\alpha_{1,1}$ $\qquad$ $\alpha_{1,2}$ $\qquad$ $\alpha_{1,3}$ $\qquad$ $\alpha_{1,4}$

$q_1$ $k_1$ $v_1$ $\qquad$ $q_2$ $k_2$ $v_1$ $\qquad$ $q_3$ $k_3$ $v_3$ $\qquad$ $q_4$ $k_4$ $v_4$

$h_1$ $\qquad$ $h_2$ $\qquad$ $h_3$ $\qquad$ $h_4$

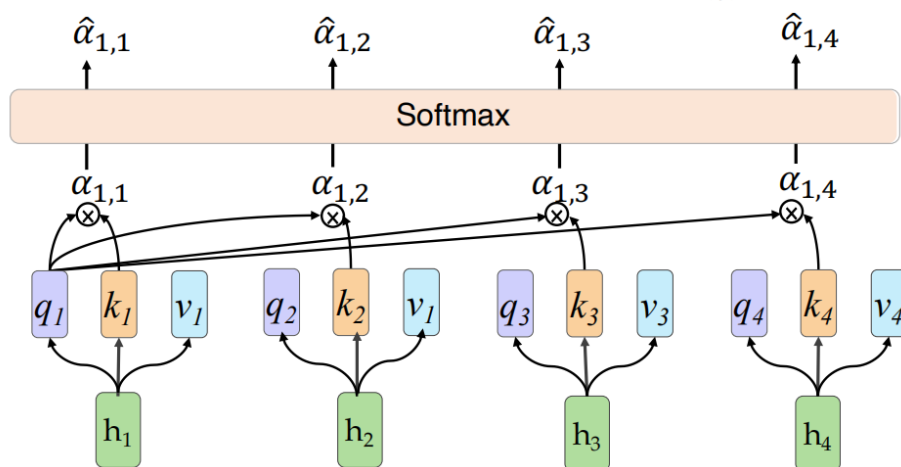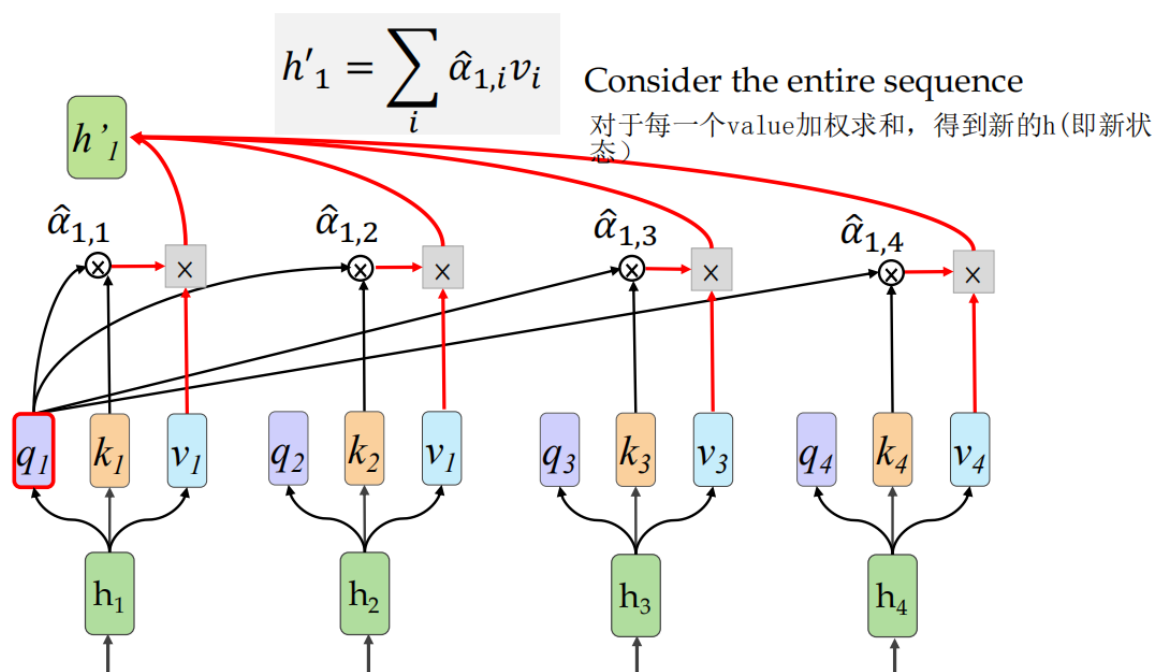## Step 4: Normalize Attention Scores(P34)

- Normalize attention scores by softmax to obtain attention weights

$$\hat{\alpha}_{1,i}=\exp(\alpha_{1,i})/\sum_j \exp(\alpha_{1,j})$$
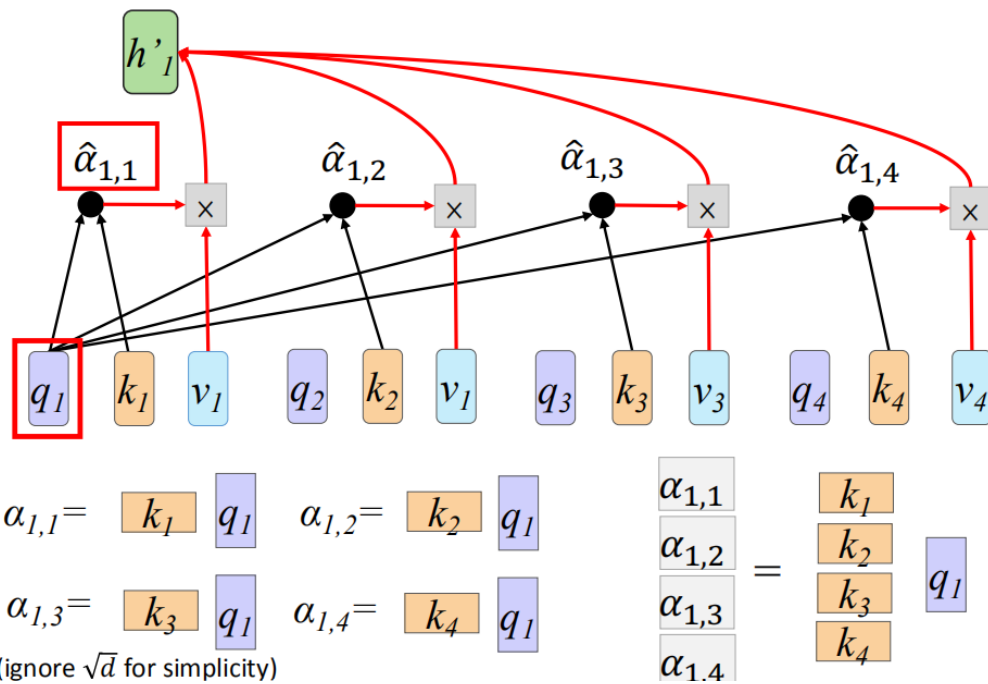
分数可能很大，用softmax归一化，得到0-1之间的分数



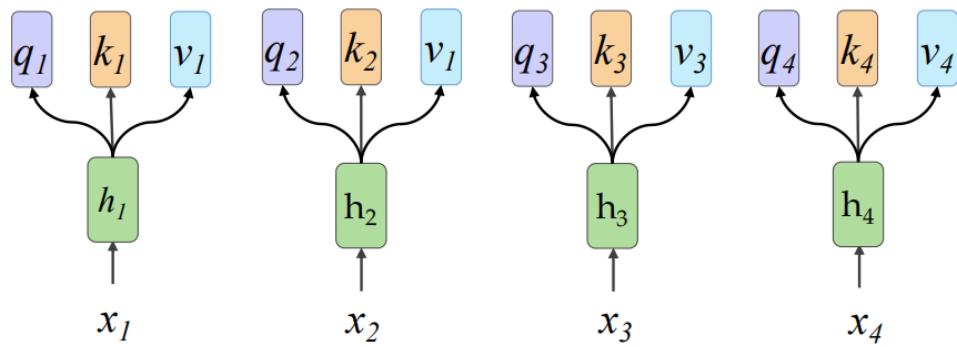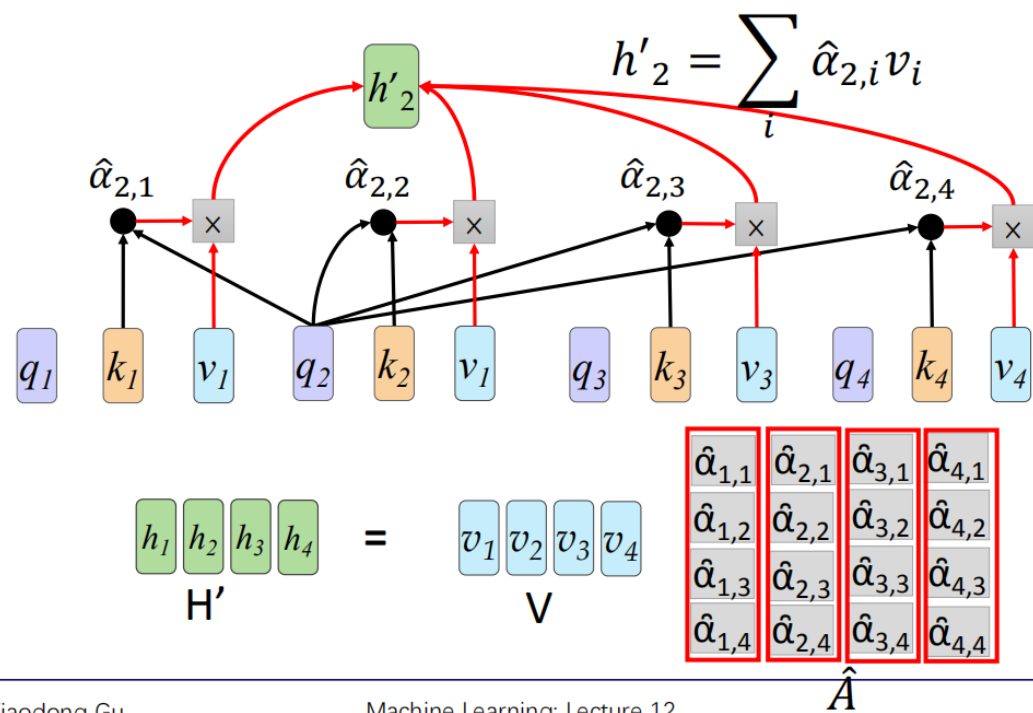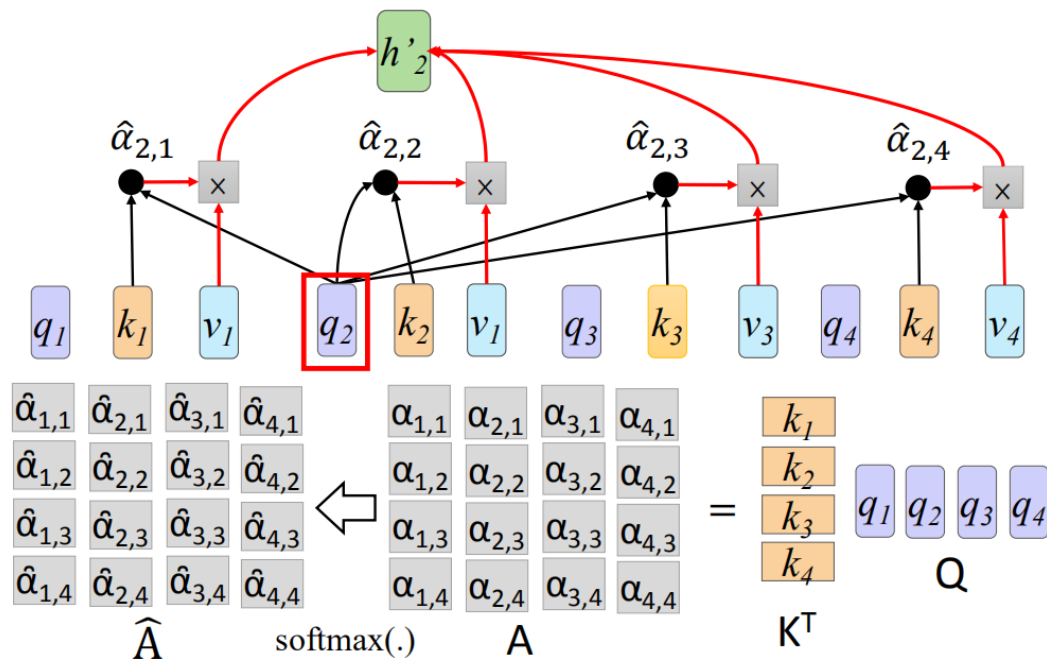## Step 5: Aggregate Values Based on Attention Weights(P35)

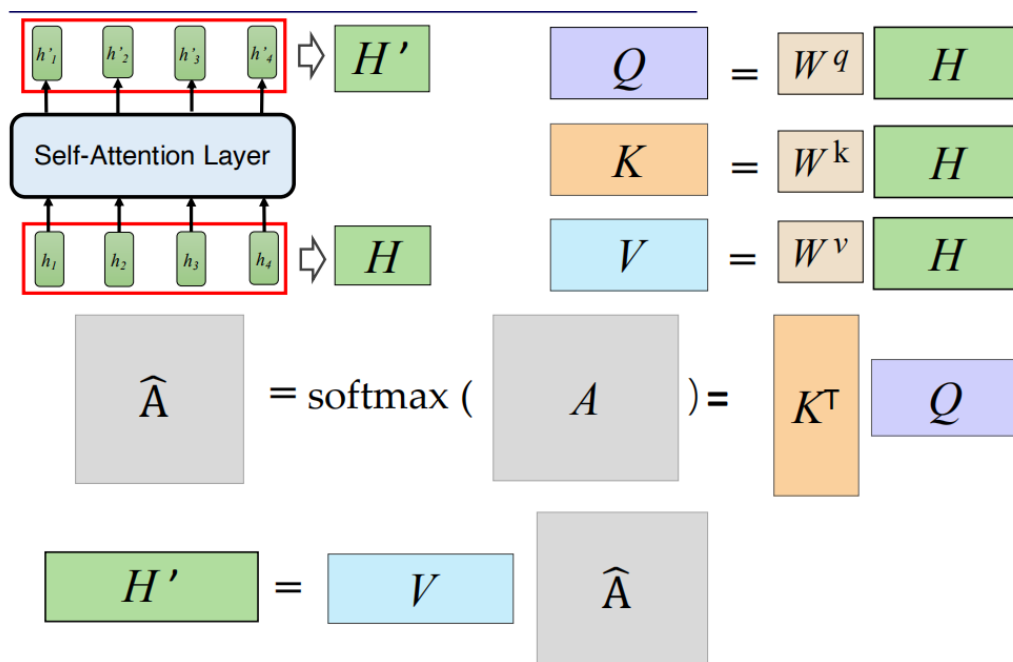$$h'_1 = \sum_i \hat{\alpha}_{1,i} v_i$$

Consider the entire sequence

对于每一个value加权求和，得到新的h（即新状态）



- h1,h2的更新其实是同步进行的，彼此之间互不影响

## Self-Attention as Matrix Multiplication(P40-45)

$$q_i = W^q h_i$$

$$Q \quad \boxed{q_1}\,\boxed{q_2}\,\boxed{q_3}\,\boxed{q_4} \;=\; \boxed{W^q} \quad \boxed{h_1}\,\boxed{h_2}\,\boxed{h_3}\,\boxed{h_4} \quad H$$

$$k_i = W^k h_i \quad \Rightarrow \quad K \quad \boxed{k_1}\,\boxed{k_2}\,\boxed{k_3}\,\boxed{k_4} \;=\; \boxed{W^k} \quad \boxed{h_1}\,\boxed{h_2}\,\boxed{h_3}\,\boxed{h_4} \quad H$$

$$v_i = W^v h_i \quad\quad V \quad \boxed{v_1}\,\boxed{v_2}\,\boxed{v_3}\,\boxed{v_4} \;=\; \boxed{W^v} \quad \boxed{h_1}\,\boxed{h_2}\,\boxed{h_3}\,\boxed{h_4} \quad H$$



$$\alpha_{1,1} = \boxed{k_1}\,\boxed{q_1} \quad\quad \alpha_{1,2} = \boxed{k_2}\,\boxed{q_1}$$

$$\alpha_{1,3} = \boxed{k_3}\,\boxed{q_1} \quad\quad \alpha_{1,4} = \boxed{k_4}\,\boxed{q_1}$$

(ignore $\sqrt{d}$ for simplicity)

$$\begin{array}{c} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{array} = \begin{array}{c} k_1 \\ k_2 \\ k_3 \\ k_4 \end{array} \; q_1$$

$$h'_2 = \sum_i \hat{\alpha}_{2,i} v_i$$
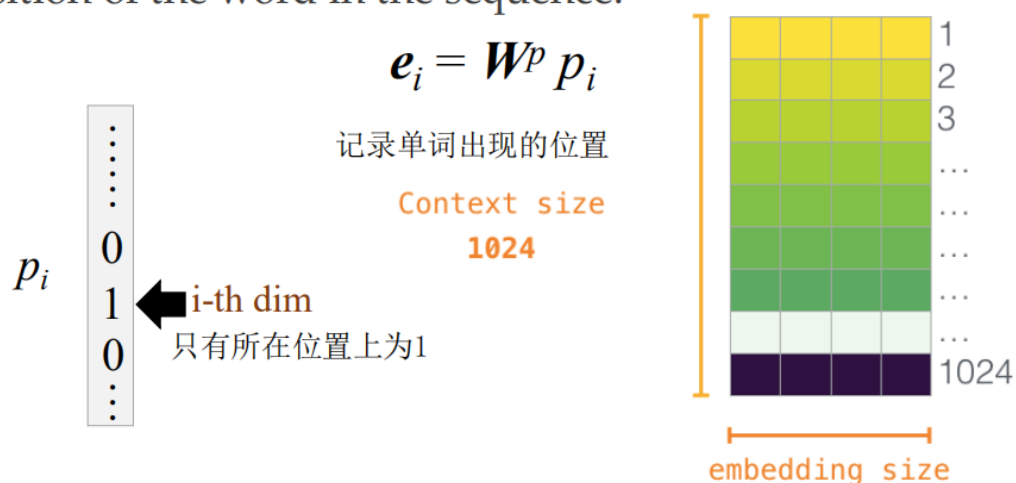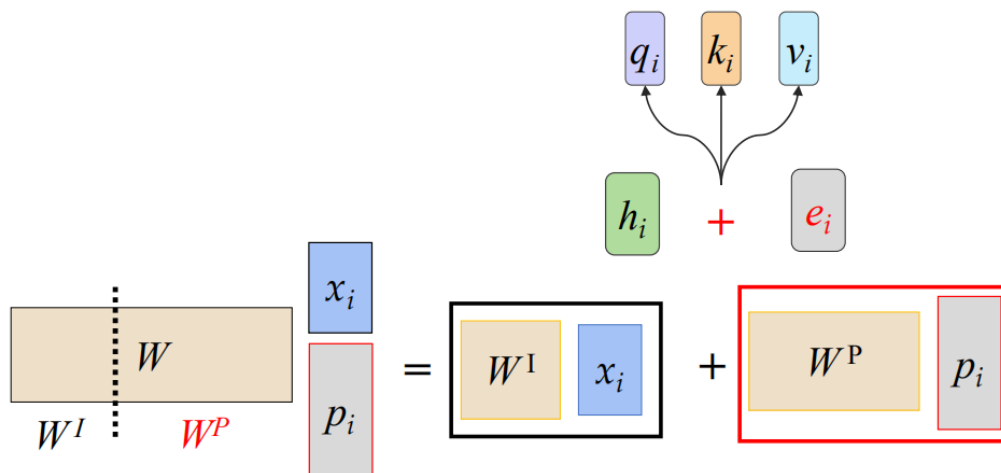
Xiaodong Gu Machine Learning: Lecture 12

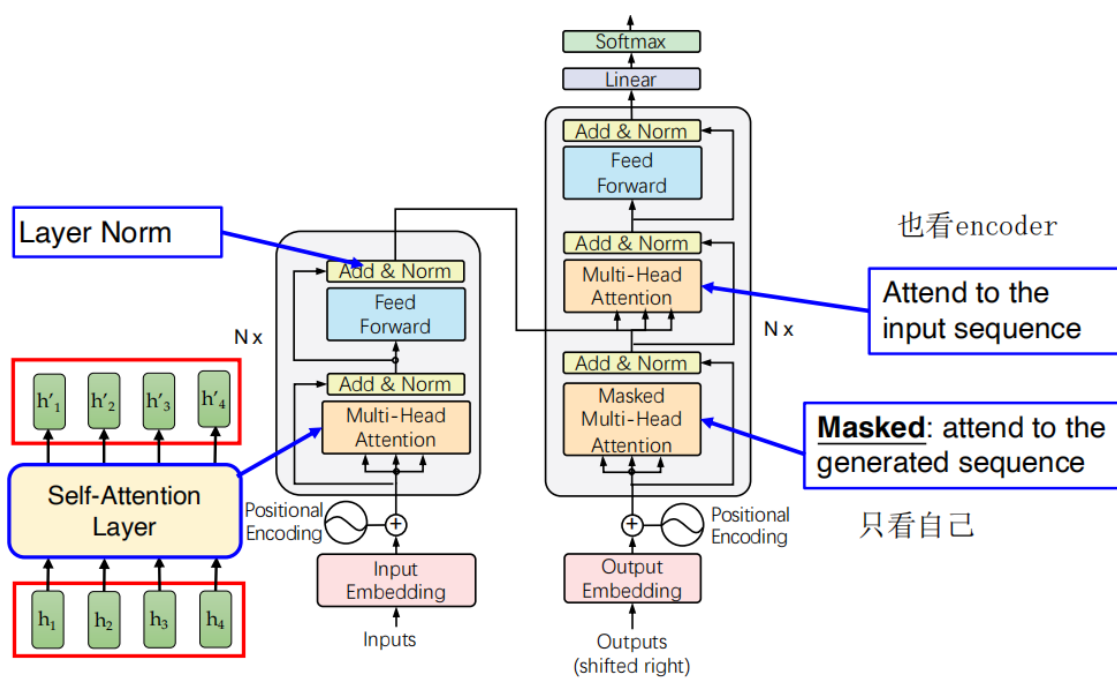## Self-Attention with Position Encoding(P48-49)

- Position Encoding: each position has a unique positional vector $e_i$ (typically learned from data)

- For each $x_i$ we append a one-hot vector $p_i$ indicating the position of the word in the sequence.

$$e_i = W^p \, p_i$$

记录单词出现的位置

Context size
1024



$p_i$ 
0
1 ◄ i-th dim
0
只有所在位置上为1

- The position encoding $e_i$, combined with the word embedding $h_i$, is feed into the self-attention layer.

$$q_i \quad k_i \quad v_i$$

$$h_i \quad + \quad e_i$$

$$\begin{bmatrix} W^I & \vdots & W^P \end{bmatrix} \begin{bmatrix} x_i \\ p_i \end{bmatrix} = \begin{bmatrix} W^I & x_i \end{bmatrix} + \begin{bmatrix} W^P & p_i \end{bmatrix}$$

## Transformer Architecture(P52)



Layer Norm

Self-Attention Layer

$h'_1 \quad h'_2 \quad h'_3 \quad h'_4$

$h_1 \quad h_2 \quad h_3 \quad h_4$

也看 encoder

Attend to the input sequence

**Masked**: attend to the generated sequence

只看自己

## Training(P64)

## The Same as RNN Encoder-Decoder

**Data:** $\{(x^{(1)}, y^{(1)}), \ldots, (x^{(N)}, y^{(N)})\}$,

where $x^{(\ell)} = (x_1, \ldots x_{Tx})$ and $y^{(\ell)} = (y_1, \ldots y_{Ty})$.

**Loss Function** – minimize the cross-entropy loss:

$$L(\theta) = -\frac{1}{N} \sum_{\ell=1}^{N} \sum_{t=1}^{Ty} \log p_\theta(y_t^{(\ell)} | x^{(\ell)})$$

**Optimization** – gradient descend

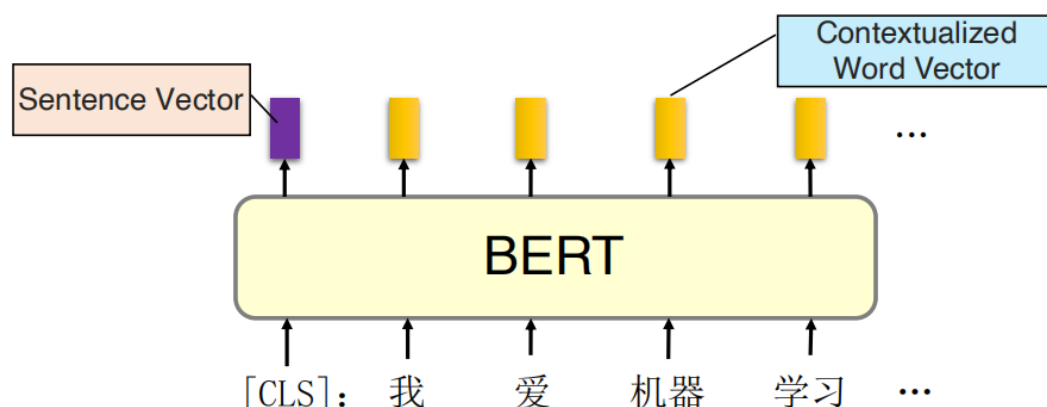# Pretrained Language Model

💡 **Pretrained Language Models def(P69):** 预训练模型，解决数据不够的问题，在公共的数据上先训练出一个模型（在大规模数据上面作无监督学习）

## Bidirectional Encoder Representations from Transformers (BERT) overview(P72)

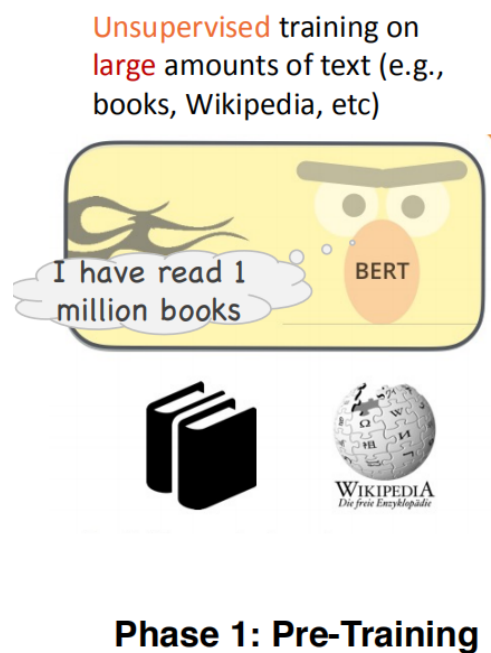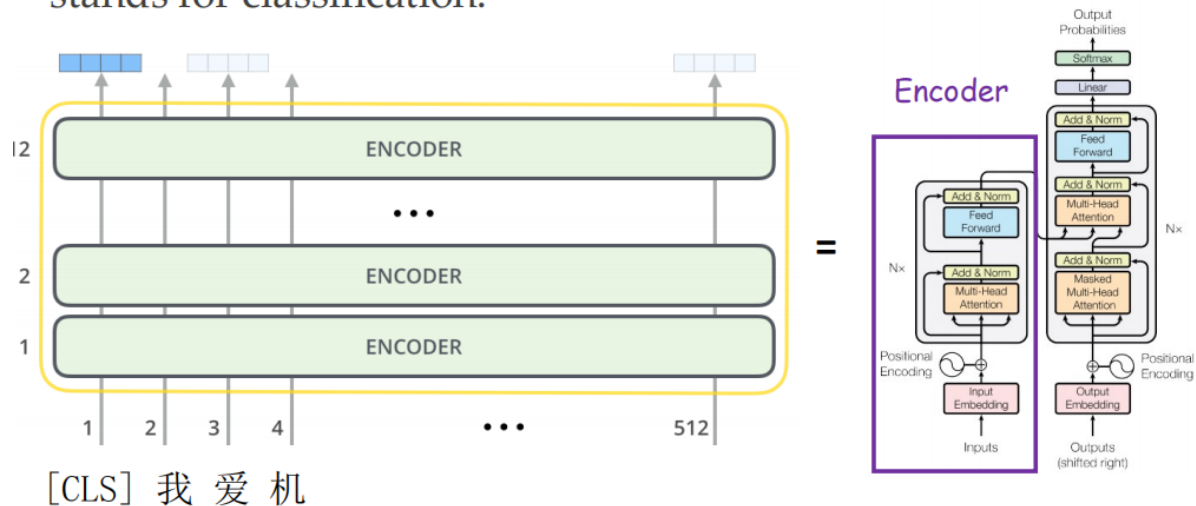### A Transformer Encoder that

- allows for learning representations of words and sentences.
- pre-trained on large-scale text corpora and then（无监督）;
- fine-tuned on small task-specific datasets (e.g., classification, QA.)

1. **Model Architecture**

- Just like Transformer encoder, BERT takes a sequence of words as input which keep flowing up the stack.
- The first input token is always a special [CLS] token which stands for classification.



Unsupervised training on large amounts of text (e.g., books, Wikipedia, etc)

**Phase 1: Pre-Training**

Supervised training on a specific task with a labeled dataset. (e.g., spam detection)

**Phase 2: Fine-Tuning**

- **Pre-Training(P75-76)**

## Task 1: Masked Language Model

| | |
|---|---|
| 0.1% | 我 |
| ... | ... |
| 10% | 学习 |
| ... | ... |
| 0% | 好 |

learn to predict the masked token

mask:挖掉一个词，推导这个被mask的词

**Classifier**
(MLP+Softmax)

**Transformer Encoder**

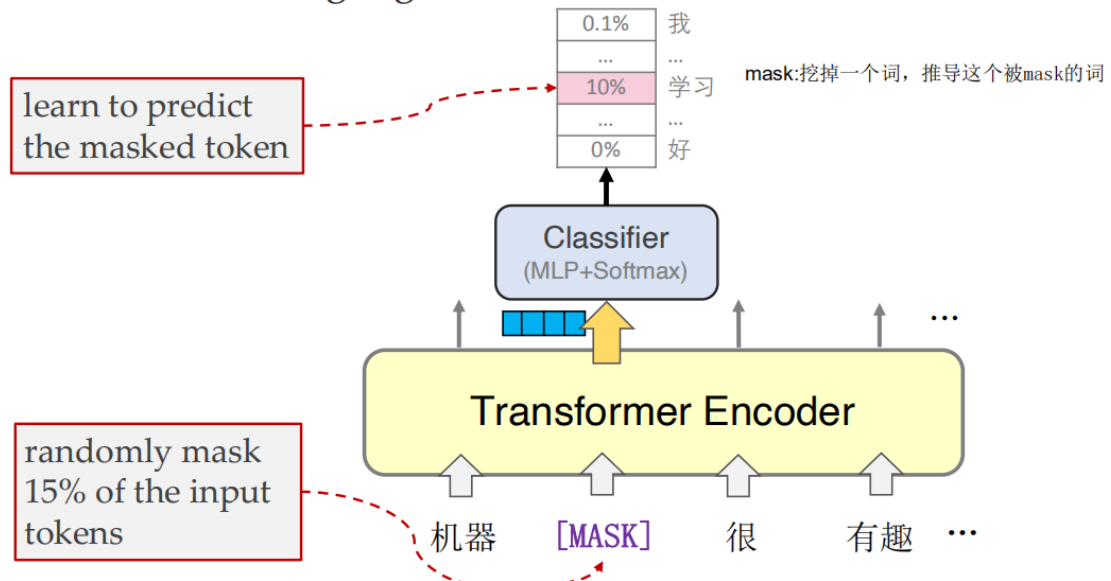randomly mask 15% of the input tokens

机器　　[MASK]　　很　　有趣　…

## Task 2: Next Sentence Prediction

| 85% | Is Next |
|---|---|
| 15% | Not Next |

predict the likelihood that sentence B comes after sentence A

切分，用前面的预测后面的

**Classifier**
(MLP+Softmax)

... 512

**Transformer Encoder**

... 512

[CLS]　机　　器　[SEP]　学　　习　　很　[SEP]

the position that outputs classification results

A

B

the boundary of two sentences

Machine Learning: Lecture 12

- **Finetuning(P77-79)**
  - **Sentence Classification**

- **input**: a single sentence, **output**: class

class

Classifier (MLP+Softmax)

注重classfier的训练

trained from scratch

Fine-tuned

BERT

[CLS] $w_1$ $w_2$ $w_3$

sentence

Examples:
- sentiment analysis
- document classification

○ **Word Tagging**

Input: a single sentence    Output: class of each word

class    class    class

Classifier    Classifier    Classifier

BERT

[CLS] $w_1$ $w_2$ $w_3$

sentence

Example: slot filling

Arrive Shanghai on Nov 2nd

other dest other time time

○ **Classifying Sentence Pairs**

**input**: two sentences    **output**: class

自然语言推理给定一个"前提"，确定一个"假设"是否为T/F/未知

Example: natural language inference given a "premise", determining whether a "hypothesis" is T/F/unknown.

class

Classifier
(MLP+Softmax)

classifier:关注具体的任务

BERT    BERT：关注模型整体

[CLS]    $w_1$    $w_2$    [SEP]    $w_3$    $w_4$    $w_5$

sentence1                    sentence 2

○ **QA (Reading Comprehension)**

s = 2    e = 3

The answer is "$d_2$ $d_3$".

$p(e|d)$    0.1    0.2    **0.7**

softmax

Learned from scratch

BERT

[CLS]    $q_1$    $q_2$    [SEP]    $d_1$    $d_2$    $d_3$

question                    document

# BART (Denoising Seq-to-Seq Pretraining)（P86）

- An **encoder-decoder** architecture 既用encoder也用decoder，来预训练

- Pre-training by reconstructing inputs that are **corrupted** by **5** methods: (token masking, token deletion, text infilling, sentence permutation, document rotation)

- More efficient for sequence-to-sequence tasks (e.g., generation, translation, comprehension)

机器　学习　很　有趣

⇧　⇧　⇧　⇧

| Transformer Encoder | → | Transformer Decoder |

⇧　⇧　⇧　⇧

机器　[MASK]　有趣　很