

# Sequence Modeling Methods and Applications

**Prof. Mingkui Tan**

SCUT Machine Intelligence Laboratory (SMIL)



SMIL内部资料 请勿外泄

# Outline

1 Backgrounds

2 Recurrent Neural Networks (RNNs)

3 Attention Mechanisms

4 Attention is All You Need: Transformer

5 Long Short-Term Memory (LSTM) Networks

6 BERT(Bidirectional Encoder Representation from Transformers)

7 Conclusions

# Outline

1 Backgrounds

2 Recurrent Neural Networks (RNNs)

3 Attention Mechanisms

4 Attention is All You Need: Transformer

5 Long Short-Term Memory (LSTM) Networks

6 BERT(Bidirectional Encoder Representation from Transformers)

7 Conclusions

SMIL内部资料 请勿外泄

# Sequence Modeling

## Definition:

- **Sequence Modeling** is the task to model, interpret, make predictions about or generate any type of **sequential data**, such as audio, text, etc.



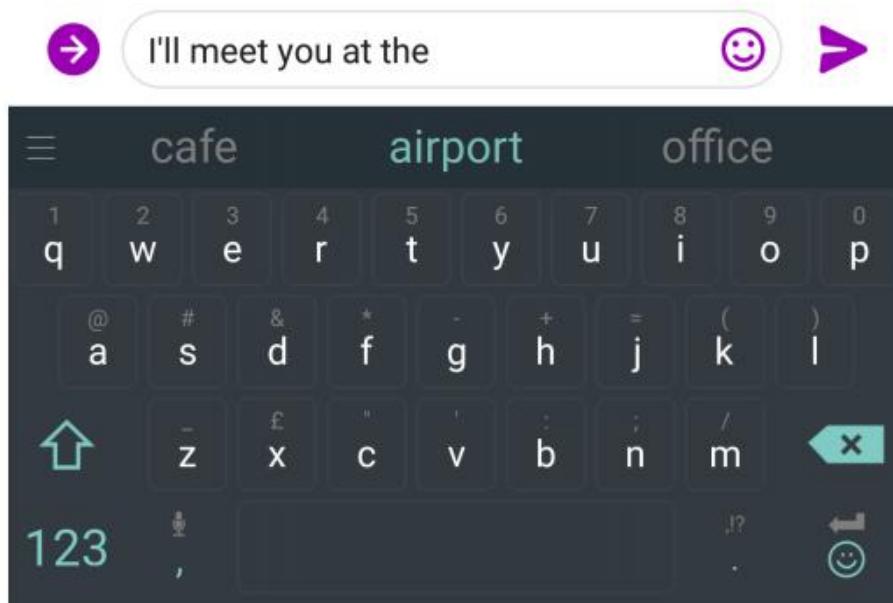
1 Second



# General Application

## Auto-Completion:

- Auto-Completion is the task of predicting what word/letter comes next.

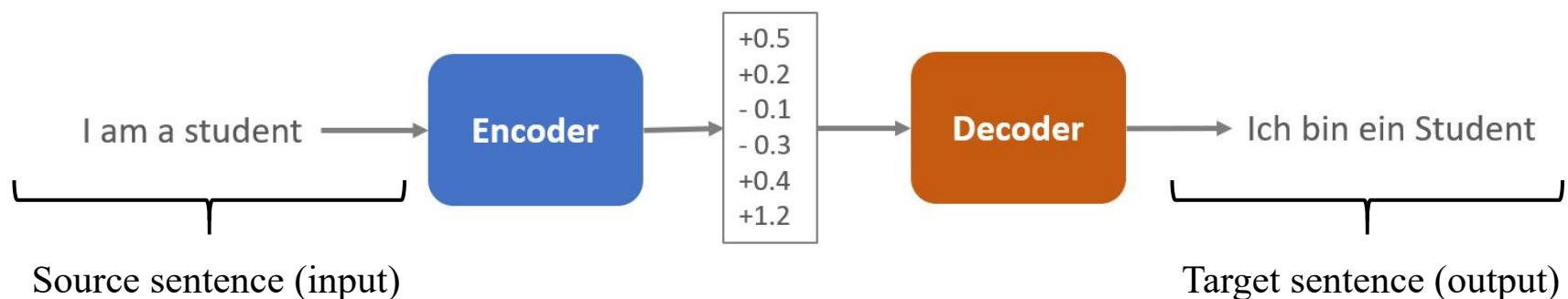
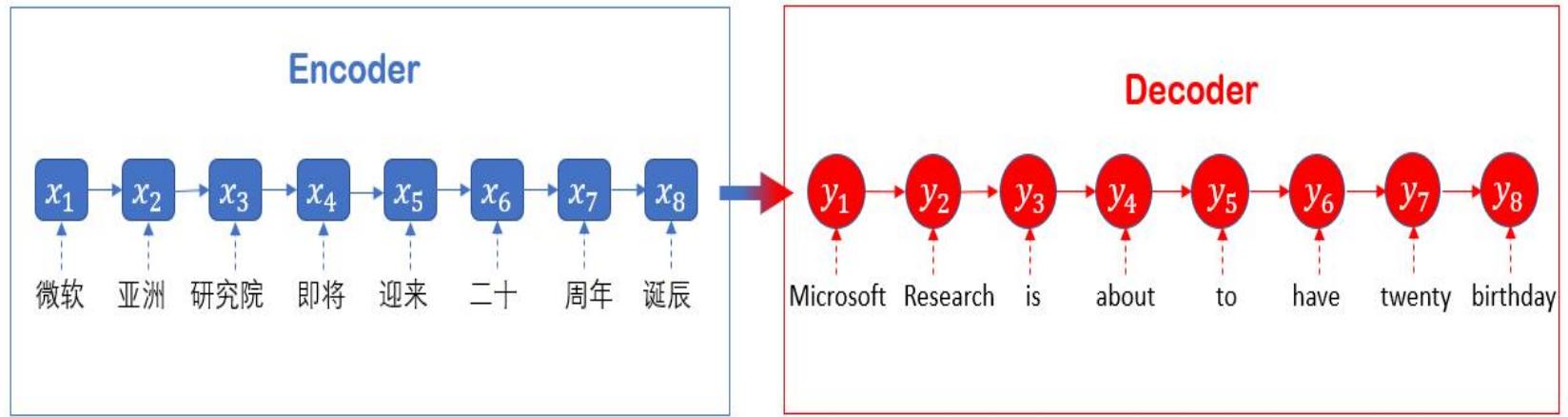


## Challenges:

- The inputs and outputs have a **sequential dependence**.
- The current output is dependent on the previous input.

# Neural Machine Translation (NMT)

- Definition: an approach to **machine translation** that uses an **artificial neural network** to predict the likelihood of a sequence of words.



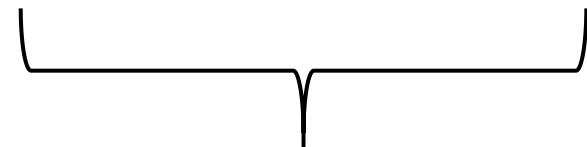
SMIL内部资料 请勿外泄

# Neural Machine Translation (NMT)

How to predict the **likelihood** of a sequence of words?

- NMT directly calculates  $P(\mathbf{y}|\mathbf{x})$ :

$$P(\mathbf{y}|\mathbf{x}) = P(\mathbf{y}_1|\mathbf{x})P(\mathbf{y}_2|\mathbf{y}_1, \mathbf{x})P(\mathbf{y}_3|\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) \dots P(\mathbf{y}_T|\mathbf{y}_1, \dots, \mathbf{y}_{T-1}, \mathbf{x})$$

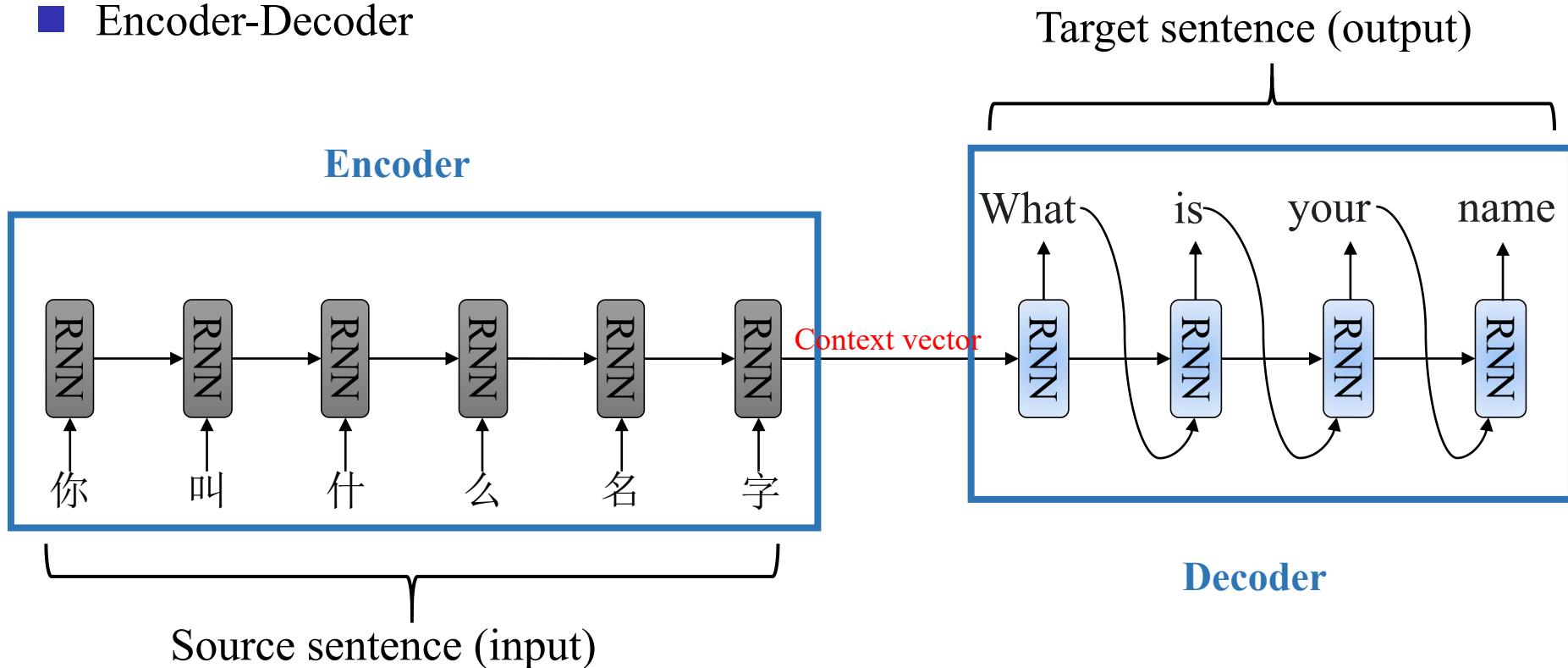


Probability of next target word, given target words so far and source sentence  $\mathbf{x}$

- Input: encode of the source sentence  $\mathbf{x}_1, \dots, \mathbf{x}_T$ .
- Output: generate target sentence  $\mathbf{y}_1, \dots, \mathbf{y}_T$ .

# Neural Machine Translation (NMT)

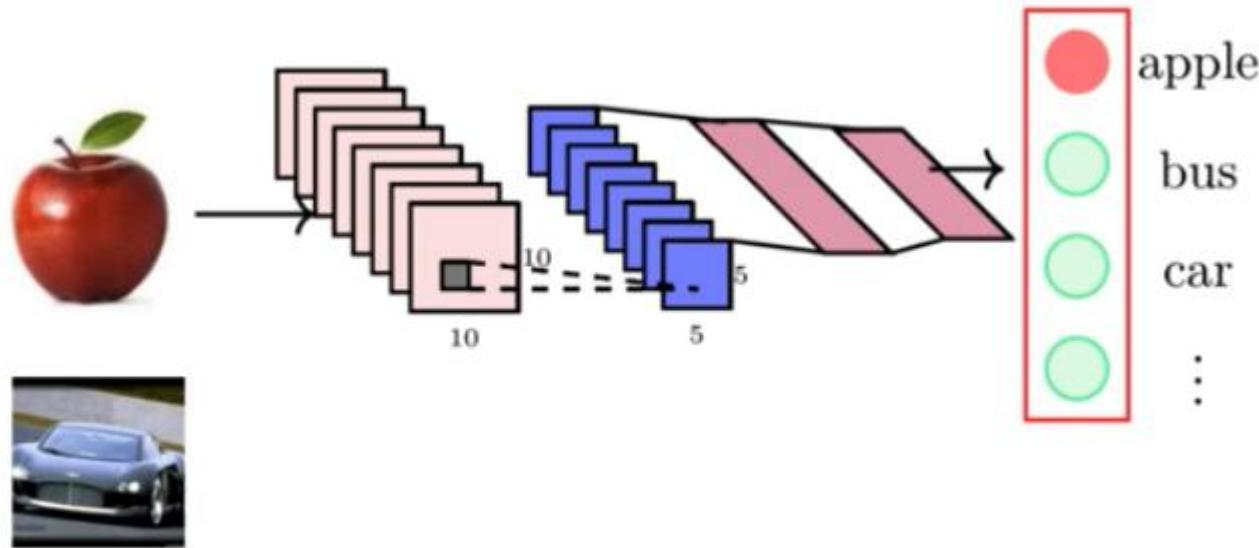
## ■ Encoder-Decoder



- **Encoder:** produce an encoding of the source sentence into a context vector.
- **Decoder:** a language Model that generates target sentence from the context vector iteratively.

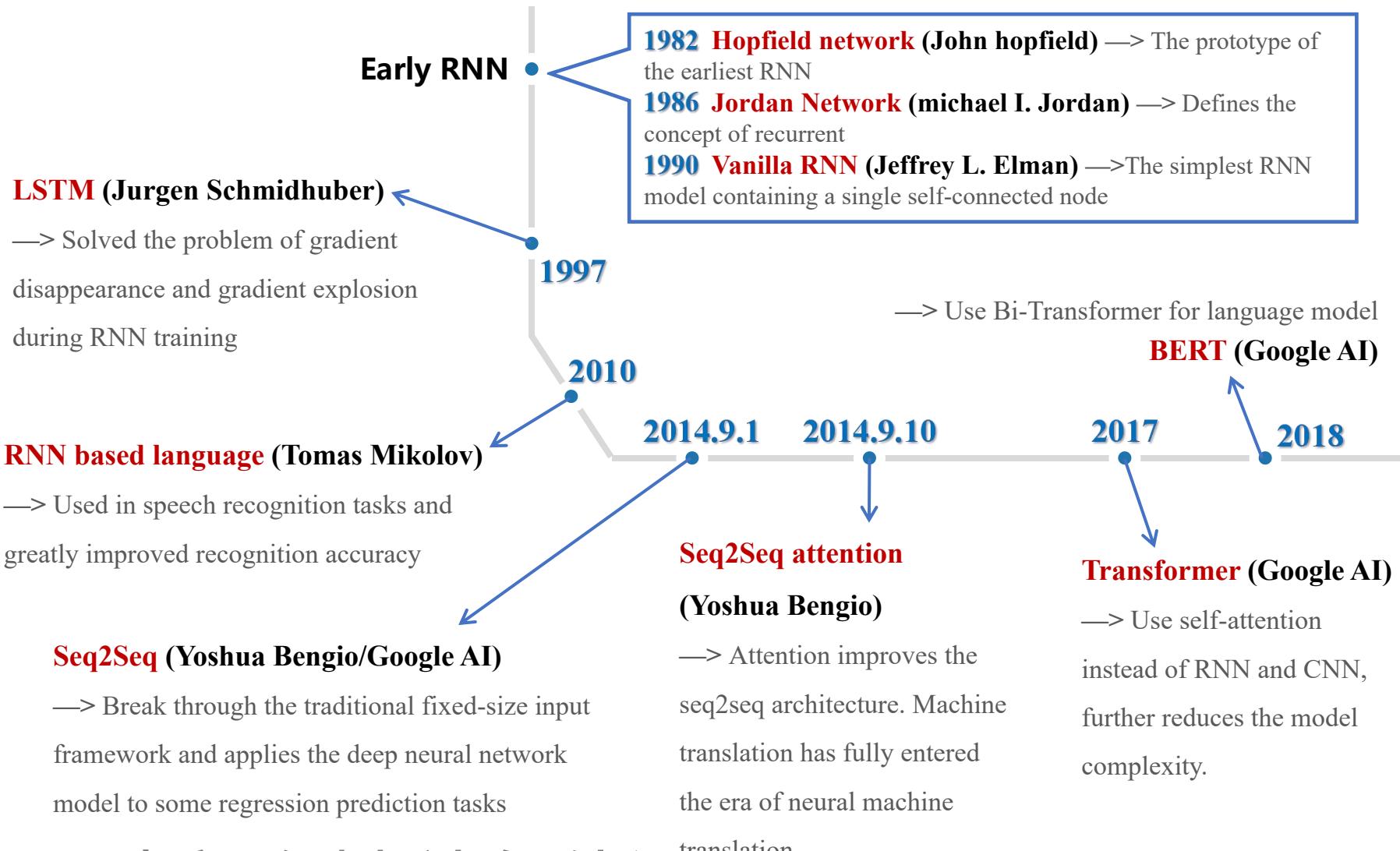
SMIL内部资料 请勿外泄

# Why Don't We Use CNN



- Traditional Convolution Neural Networks (CNNs) only consider the current input and all inputs are **independent** of each other.
- CNN is more suitable for image processing.

# Evolution of Sequence Modeling methods

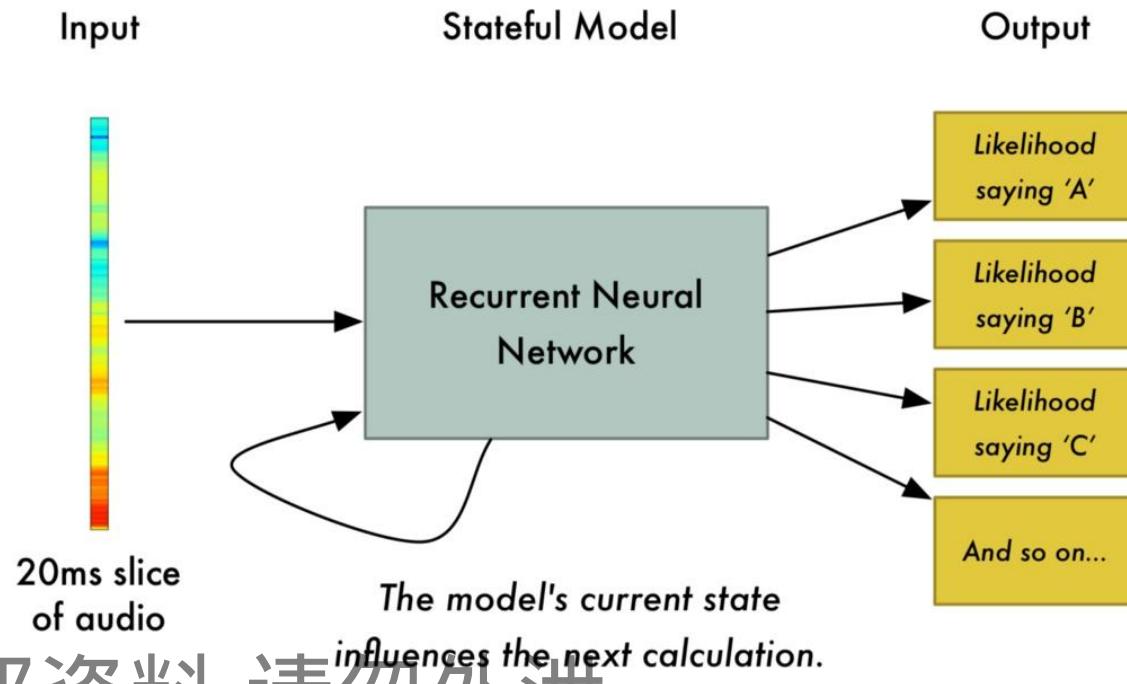


SMIL内部资料 请勿外泄

# Sequence to Sequence Model Applications

## Speech Recognition

- **Definition:** Speech recognition aims to convert spoken language into readable text.
- **Difficulty:** Vocabulary is hard to recognize if the speech contains confusing words.

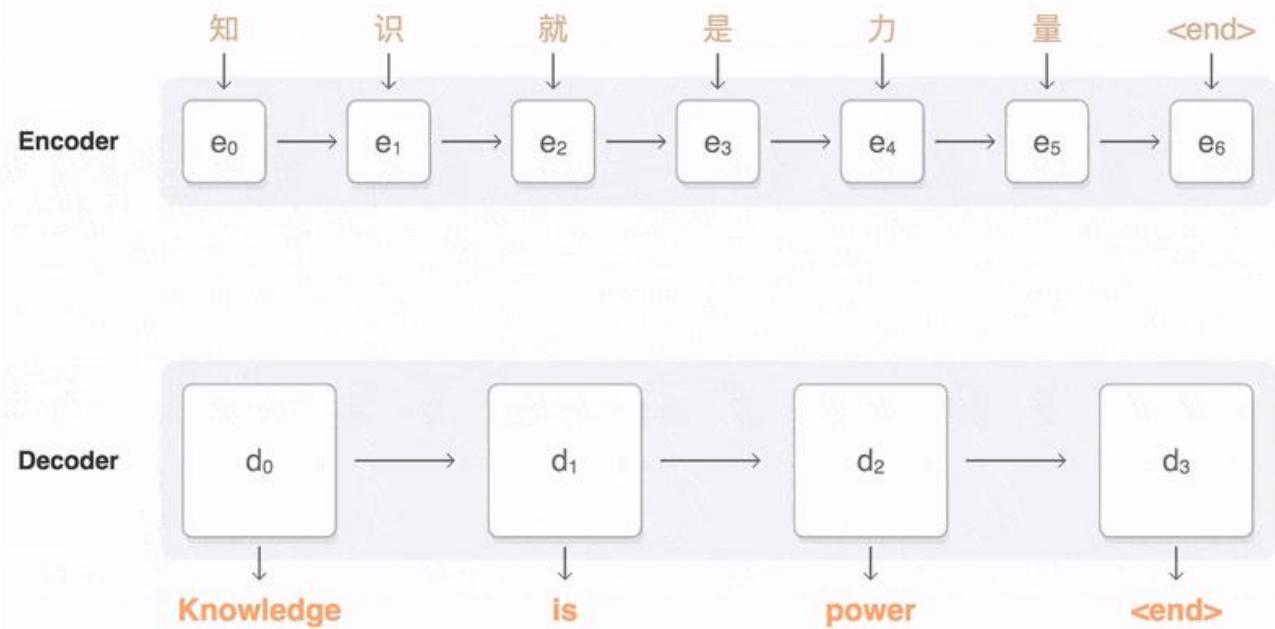


SMIL内部资料 请勿外泄

# Sequence to Sequence Model Applications

## Machine Translation

- **Definition:** Machine translation aims to translate text from one language to another.
- **Difficulties:** Word sense disambiguation and Syntactic disambiguation



# Sequence to Sequence Model Applications

## Named Entity Recognition

- **Definition:** Named-entity recognition seeks to locate and classify named entities mentioned in unstructured text into pre-defined categories.
- **Challenge:** Devise models to deal with linguistically complex contexts.

Sundar Pichai is CEO of Google having headquarter in Mountain View, CA

Person

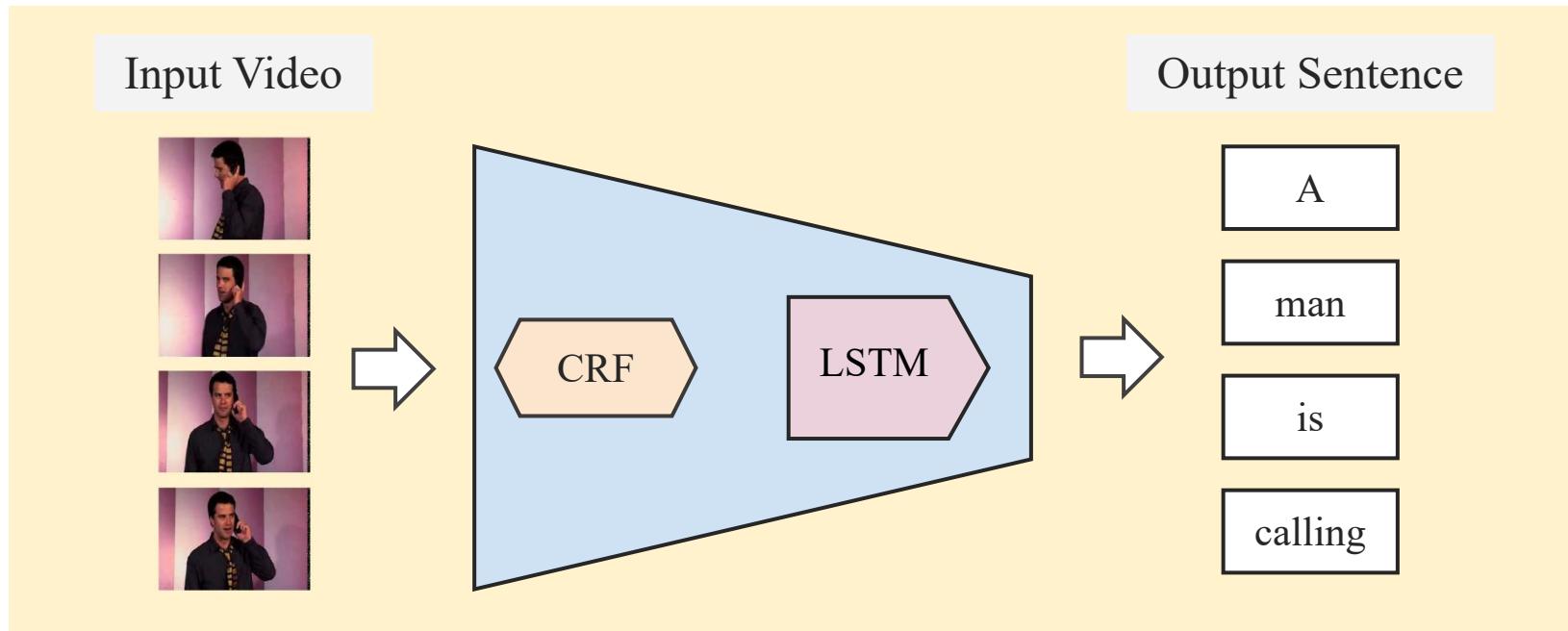
Organization

Location

# Sequence to Sequence Model Applications

## Video Captioning

- **Definition:** Video Captioning aims to generate text descriptions for videos.
- **Challenge:** Model variable length videos with a fixed number of parameters.



SMIL内部资料 请勿外泄

# Outline

1 Backgrounds

2 Recurrent Neural Networks (RNNs)

3 Attention Mechanisms

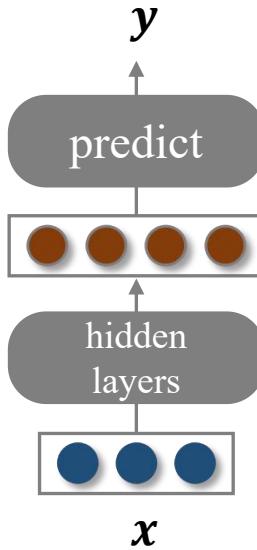
4 Attention is All You Need: Transformer

5 Long Short-Term Memory (LSTM) Networks

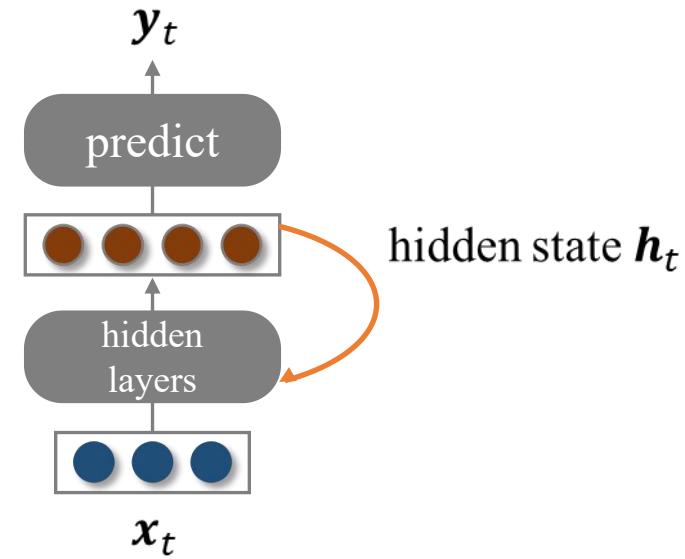
6 BERT(Bidirectional Encoder Representation from Transformers)

7 Conclusions

RNNs are capable of remembering information in time dimension.



Feedforward Networks



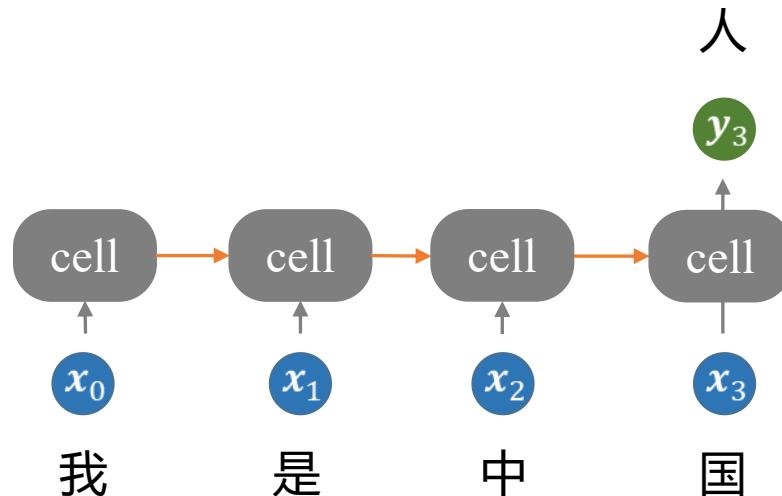
Recurrent Networks

- Input: sequence of words  $x_1, \dots, x_t$
- Output: sequence of words  $y_1, \dots, y_t$

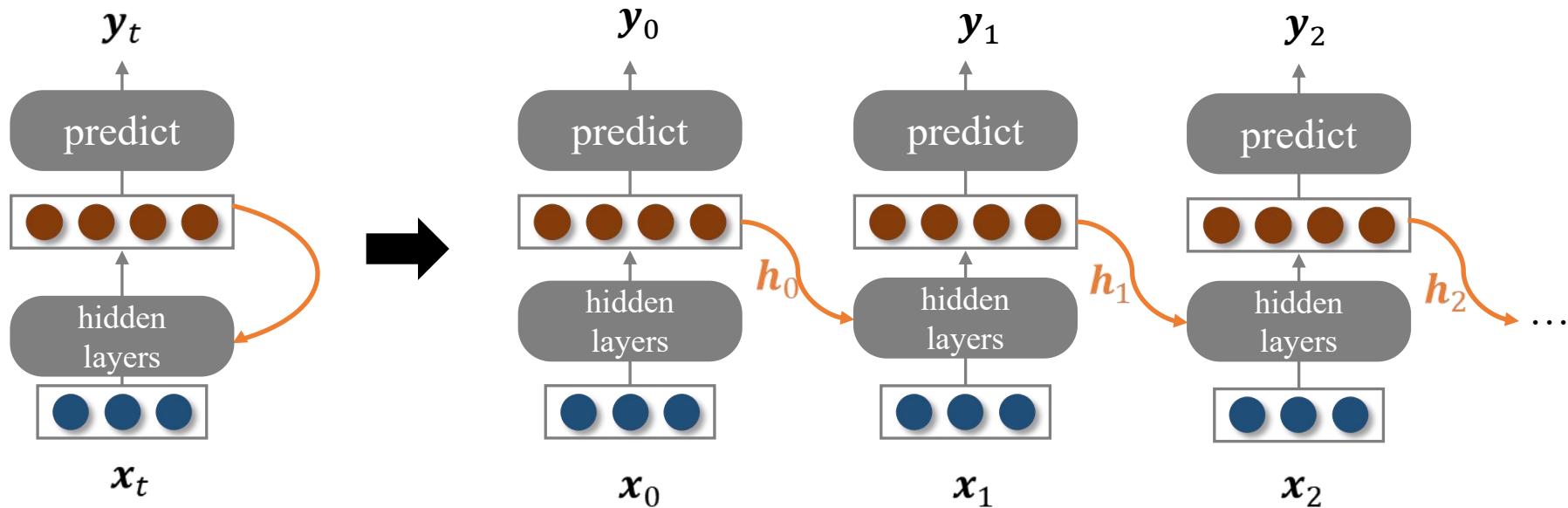
- $x_t \leftarrow$  input: sequence of words
- $y_t \leftarrow$  output: sequence of words
- $h_t \leftarrow$  hidden node
- $U, V, W \leftarrow$  coefficients
- $E \leftarrow$  embedding matrix

- Predict a single output:

Maximize the likely-hood of  $P(y_3|x_0, x_1, x_2, x_3)$



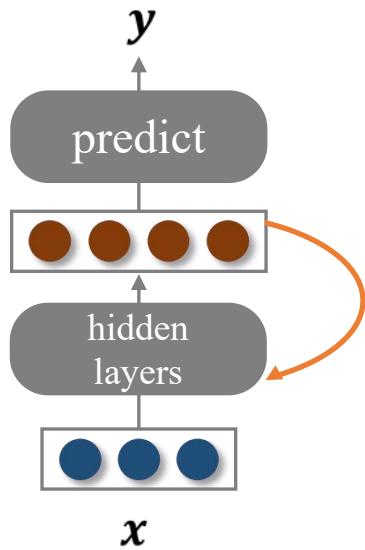
Unroll the recurrent loop:



A RNN can be thought of as multiple copies of the same network, each passing a message to a successor.

SMIL内部资料 请勿外泄

We can model the recurrence formula by a sequence of vectors  $x$  at every time step:



$$h_t = g(h_{t-1}, x_t)$$

Annotations for the equation:

- $h_t$ : new hidden node
- $g$ : some activation function
- $h_{t-1}$ : old hidden node
- $x_t$ : input vector at  $t$  time step

The same function and the same set of parameters are used at every time step.

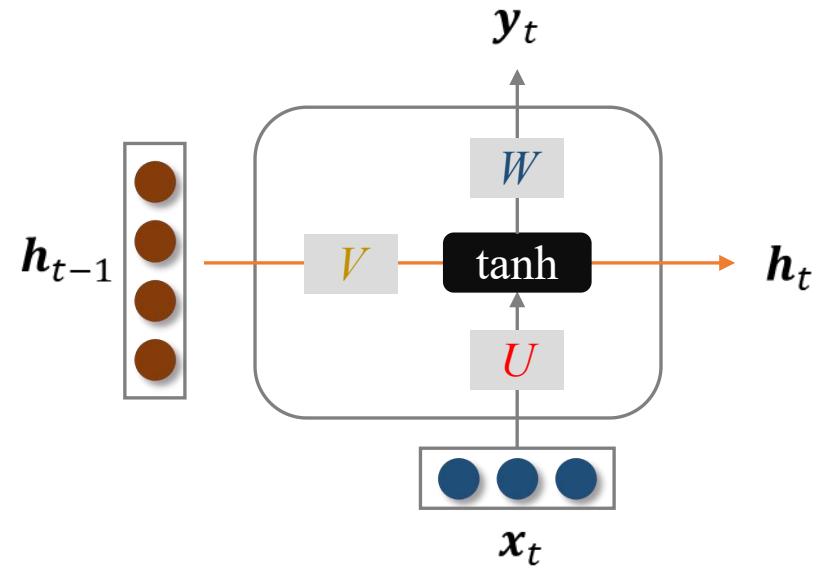
A typical RNN cell and time step  $t$ :

$$\mathbf{h}_t = g(\mathbf{h}_{t-1}, \mathbf{x}_t)$$



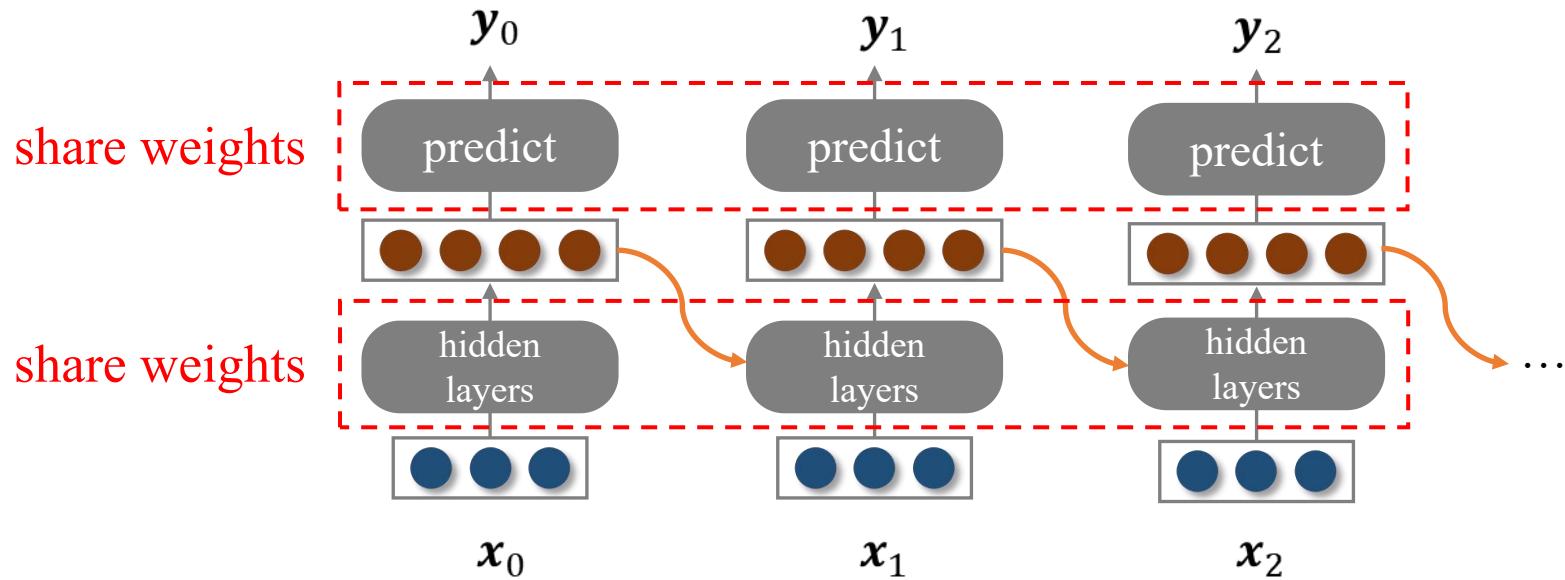
$$\mathbf{h}_t = \tanh(\mathbf{U}^T \mathbf{x}_t + \mathbf{V}^T \mathbf{h}_{t-1} + b)$$

$$\mathbf{y}_t = \mathbf{W}^T \mathbf{h}_t$$



- The **tanh** function is considered as the activation function.
- We can replace the **tanh** function with **sigmoid** or **ReLU** function.

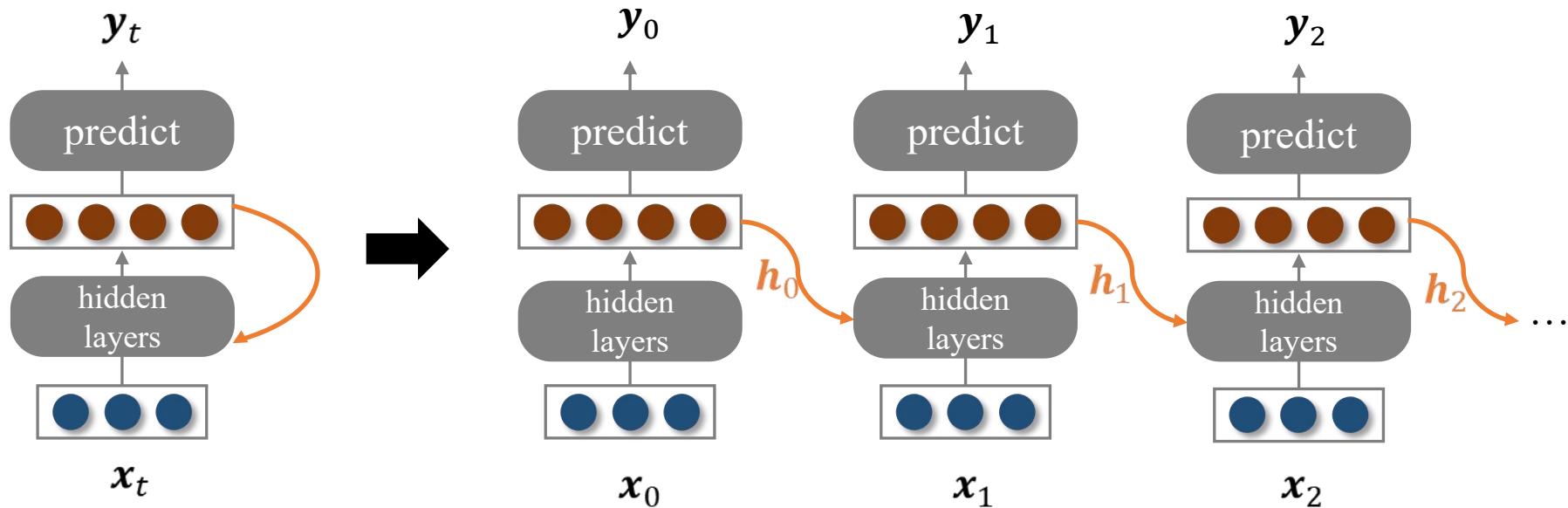
Unroll the recurrent loop:



A RNN can be thought of as multiple copies of the same network, each passing a message to a successor.

SMIL内部资料 请勿外泄

Unroll the recurrent loop:

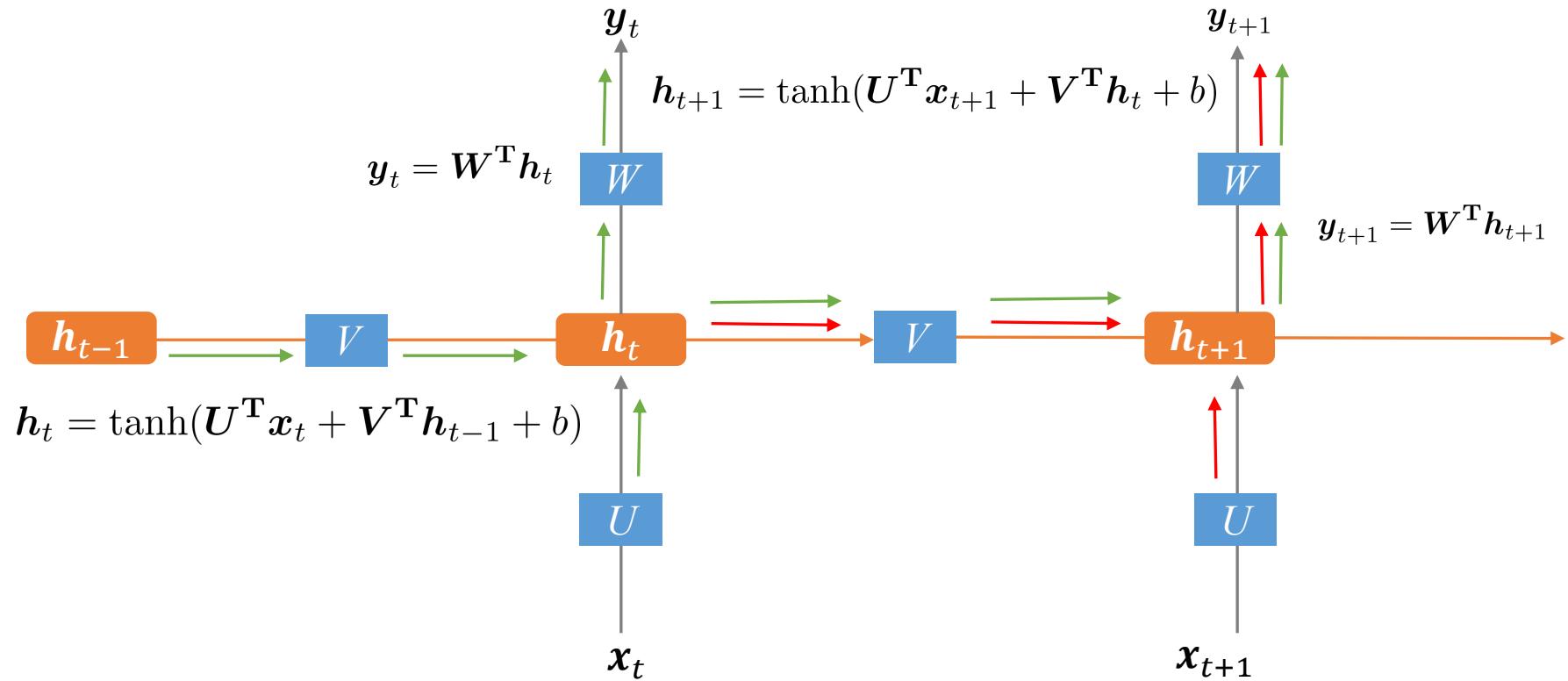


A RNN can be thought of as multiple copies of the same network, each passing a message to a successor.

SMIL内部资料 请勿外泄

# Forward Propagation: Prediction for Networks

Forward Propagation Through Time (FPTT):



# Forward Propagation: Prediction for Networks

- Given input a sequence of words  $x_1, \dots, x_T$ , the forward propagation of a RNN cell is to predict the output  $y_1$  of  $x_1$ .

- Detailed prediction process:

for  $t$  from 1 to  $T$  do

$$\mathbf{u}_t = \mathbf{U}^T \mathbf{x}_t + \mathbf{V}^T \mathbf{h}_{t-1} + b$$

$$\mathbf{h}_t = \sigma(\mathbf{u}_t)$$

$$\mathbf{y}_t = \mathbf{W}^T \mathbf{h}_t$$

end for



Example:

$$\sigma = \tanh(\cdot)$$

$\sigma$ : some activation function

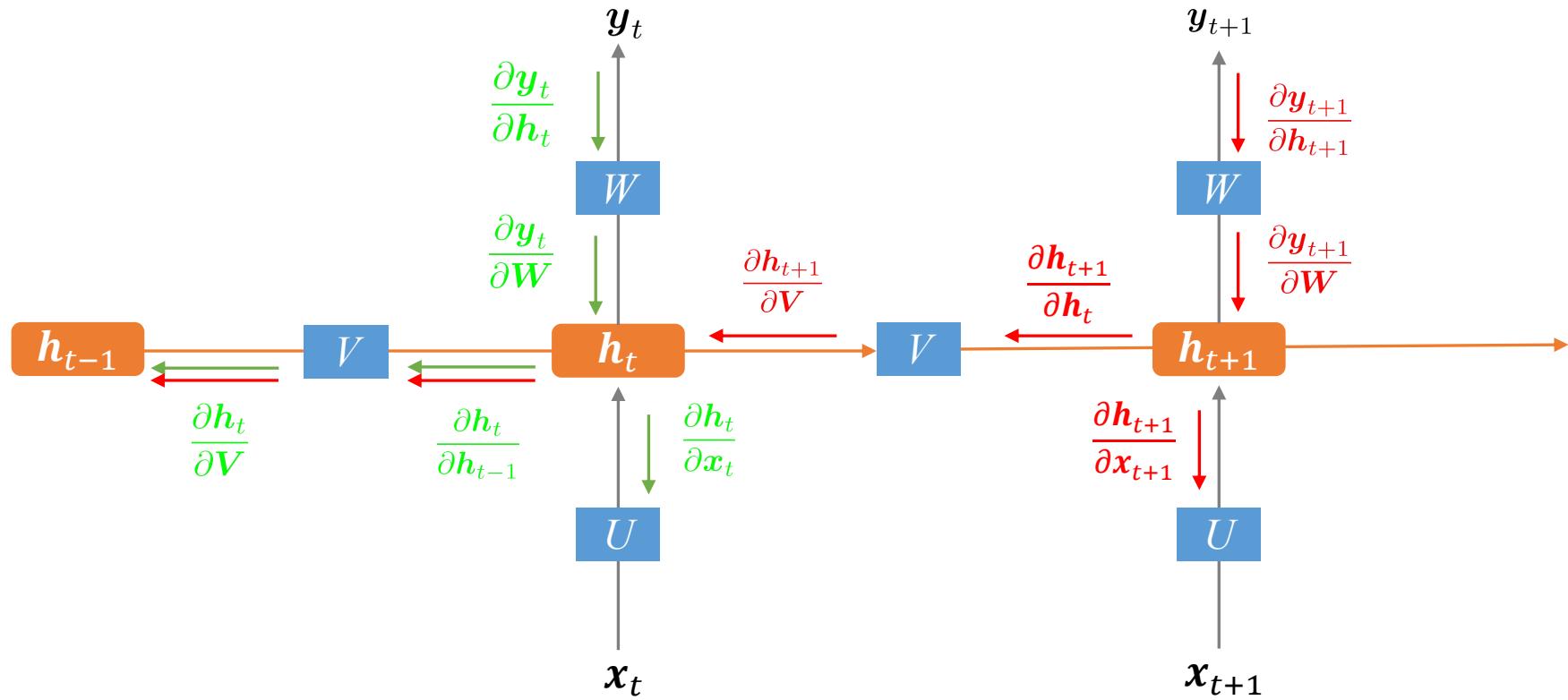
- Total loss:  $L(y, \hat{y}) = \sum_{i=1}^m \sum_{t=1}^T L_i(y_t, \hat{y}_t)$

$\hat{y}_t$  : predicted label

$y_t$  : ground-truth label

# Backpropagation of Neural Networks

Backward Propagation Through Time (BPTT):



SMIL内部资料 请勿外泄

# Backpropagation of Neural Networks

What's the derivative of  $L(\mathbf{y}, \hat{\mathbf{y}})$ , the repeated weight matrix  $U, V, W$ ?

- Let us resort to Chain Rule to compute the gradient.

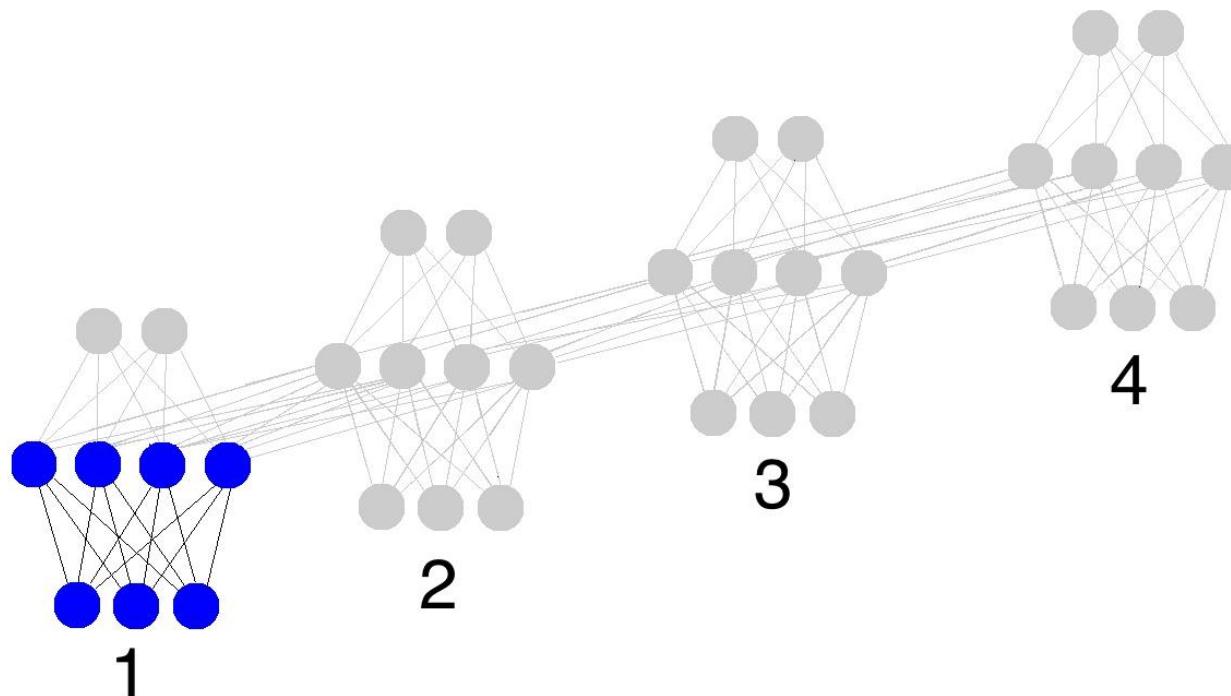
$$\frac{\partial L_3(\mathbf{y}, \hat{\mathbf{y}})}{\partial W} = \frac{\partial L_3(\mathbf{y}, \hat{\mathbf{y}})}{\partial y_3} \frac{\partial y_3}{\partial W}$$

Use  $L_3(\mathbf{y}, \hat{\mathbf{y}})$  as an example

$$\frac{\partial L_3(\mathbf{y}, \hat{\mathbf{y}})}{\partial U} = \sum_{k=0}^3 \frac{\partial L_3(\mathbf{y}, \hat{\mathbf{y}})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial h_k} \frac{\partial h_k}{\partial U}$$

$$\frac{\partial L_3(\mathbf{y}, \hat{\mathbf{y}})}{\partial V} = \sum_{k=0}^3 \frac{\partial L_3(\mathbf{y}, \hat{\mathbf{y}})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial h_k} \frac{\partial h_k}{\partial V}$$

## Backward Propagation Through Time (BPTT):

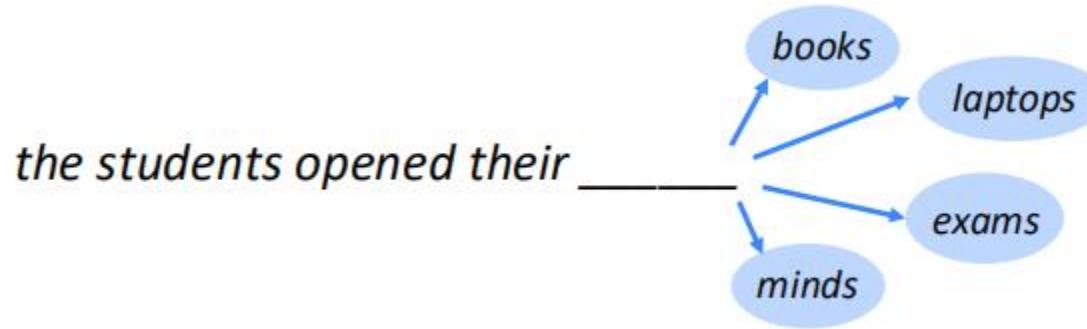


<https://blog.csdn.net/Heartbougan>  
MakeAGIF.com

Black represents the predicted value, yellow represents the prediction error, and dark yellow represents gradient back propagation.

# Application: Language Modeling (LM)

- Language Modeling: predict the coming word basis on a piece of words.



- Given a sequence of words  $x_1, \dots, x_t$ , compute the probability distribution of the next word  $x_{t+1}$ :

$$P(x_{t+1} | x_1, \dots, x_t)$$

where  $x_{t+1}$  can be any word in the vocabulary  $V = \{w_1, \dots, w_{|V|}\}$ .

- We need to consider lots of words at a time if we want to model language well.

# Training RNN Base on a Language Model

$$\hat{y}_4 = P(x_5 | \text{the students opened their})$$

- output distribution

$$\hat{y}_t = \text{softmax}(\mathbf{U}h_t + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

- hidden states

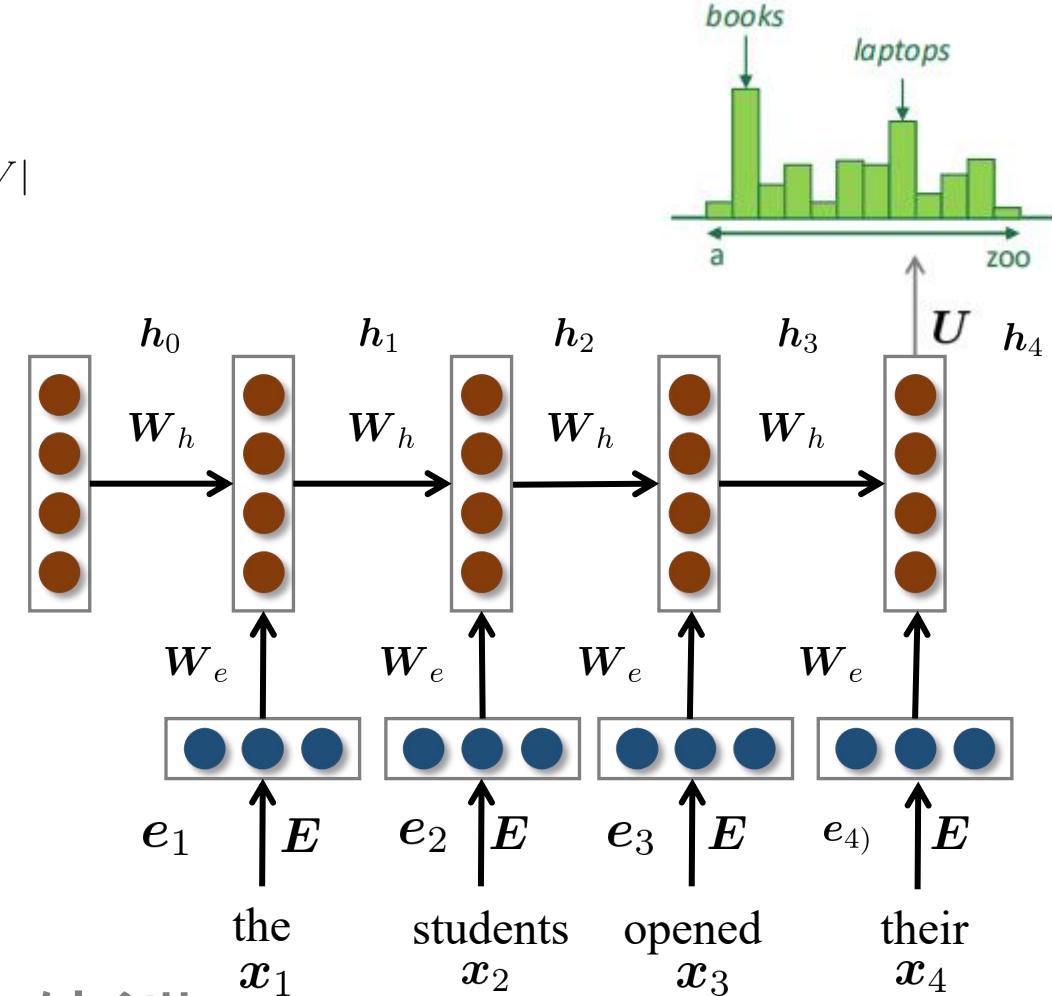
$$h_t = \sigma(\mathbf{W}_h h_{t-1} + \mathbf{W}_e e_t + \mathbf{b}_1)$$

- word embedding

$$e_t = \mathbf{E}x_t$$

- words/one-hot vectors

$$x_t \in \mathbb{R}^{|V|}$$



SMIL内部资料 请勿外泄

# Training RNN Base on a Language Model

- Get a **big corpus of text** which is a sequence of words  $x_1, \dots, x_T$ .
- Feed into RNN-LM, compute output distribution  $\hat{y}_t$  **for every step  $t$** .
- Apply the **cross-entropy function** between predicted probability distribution  $\hat{y}_t$ , and the true next word  $y_t$ :

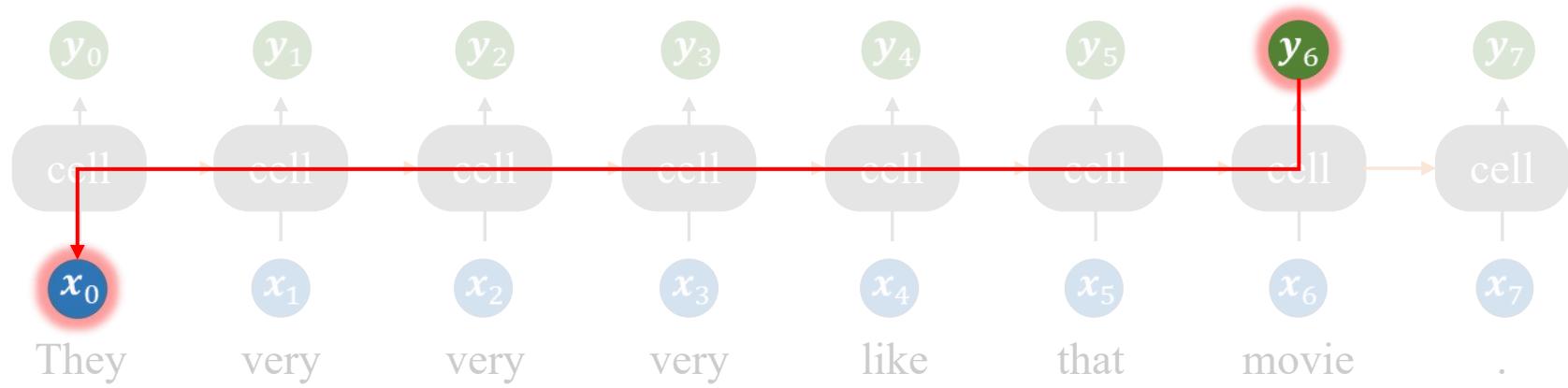
$$J_t(\theta) = CE(y_t, \hat{y}_t) = - \sum y_t \log \hat{y}_t.$$

- **Total loss** for the entire training set:

$$J(\theta) = \frac{1}{T} \sum_{i=1}^m \sum_{t=1}^T J_{i,t}(\theta)$$

# Problems with RNN

The problem of long-term dependencies:



$$y_6 = W_3 \cdot \tanh(W_1 \cdot x_6 + \underbrace{W_2 \cdot \tanh(\cdots \underbrace{W_2 \cdot \tanh(W_1 \cdot x_0 + \underbrace{W_2 \cdot h_0 + b} \cdots)}_{6 \text{ times!}} + b)}$$

# Problems with RNN

The problem of long-term dependencies:

The gradient of  $\varepsilon_t$  on  $x_k$  ( $k < t$ ) :

$$\frac{\partial \varepsilon_t}{\partial x_k} = W_3 \cdot \left( \prod_{i=k}^t \frac{\partial h_i}{\partial h_{i-1}} \right) \cdot \frac{\partial h_k}{\partial x_k} = W_3 \cdot \left( \prod_{i=k}^t W_2 \cdot \text{diag}(\tanh'(h_i)) \right) \cdot W_1$$

can vanish or explode due to the Matrix-chain multiplication.

- Lots of operations are hard to capture in the hidden state path.

# Outline

1 Backgrounds

2 Recurrent Neural Networks (RNNs)

3 Attention Mechanisms

4 Attention is All You Need: Transformer

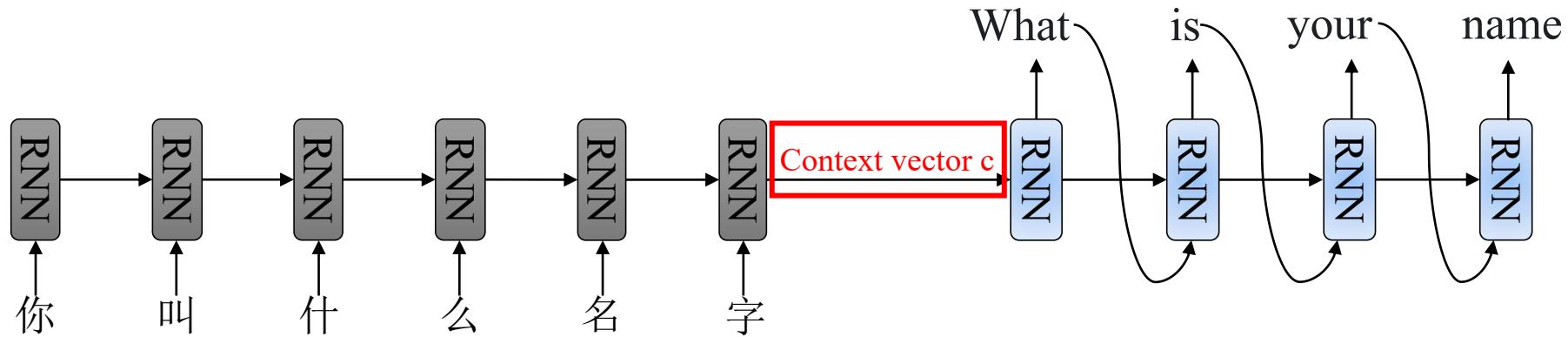
5 Long Short-Term Memory (LSTM) Networks

6 BERT(Bidirectional Encoder Representation from Transformers)

7 Conclusions

# Motivation

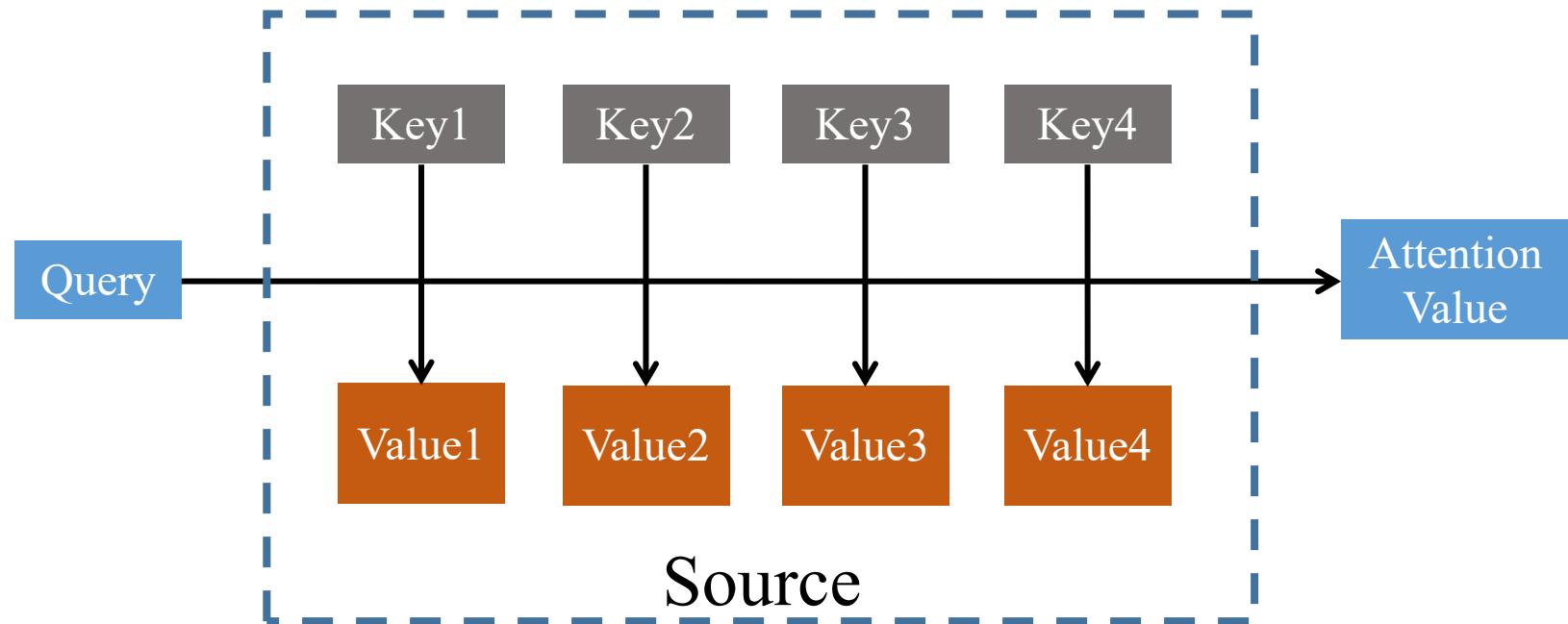
Compress all the necessary information of a source sentence into a static context vector  $c$ .



- **The problem of long-term dependencies:** the translation effect will deteriorate rapidly as the length of text increases.

# Motivation

On each step of the decoder, use **direct connection to the encoder** to **focus on a particular part** of the source sequence.



$$\text{Attention}(\text{Query}, \text{Source}) = \sum_{i=1}^t \text{Similarity}(\text{Query}, \text{Key}_i) * \text{Value}_i$$

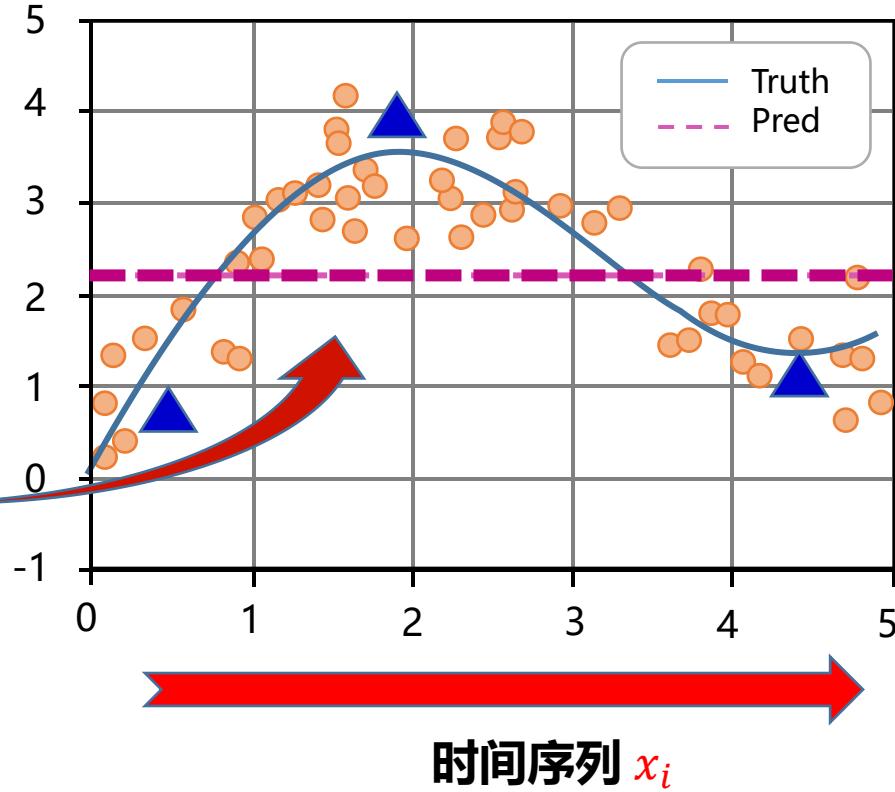
# Motivation of Attention

Considering a function fitting task: **Ground Truth:**  $y_i = 2\sin(x_i) + x_i^{0.8}$

Linear model  $y = w'x$  no longer works

Simple Prediction for a new point  $x$ :

$$f(x) = \sum_{i=1}^n y_i$$



Simply using the **average of  $y_i$**  as prediction will lead to **very poor performance!**

# Motivation of Attention

Considering a function fitting task: **Ground Truth:**  $y_i = 2\sin(x_i) + x_i^{0.8}$

**Weighted** Prediction for a new point  $x$ :

$$f(x) = \sum_{i=1}^n \alpha_i y_i$$

How to compute  $\alpha_i$ ?

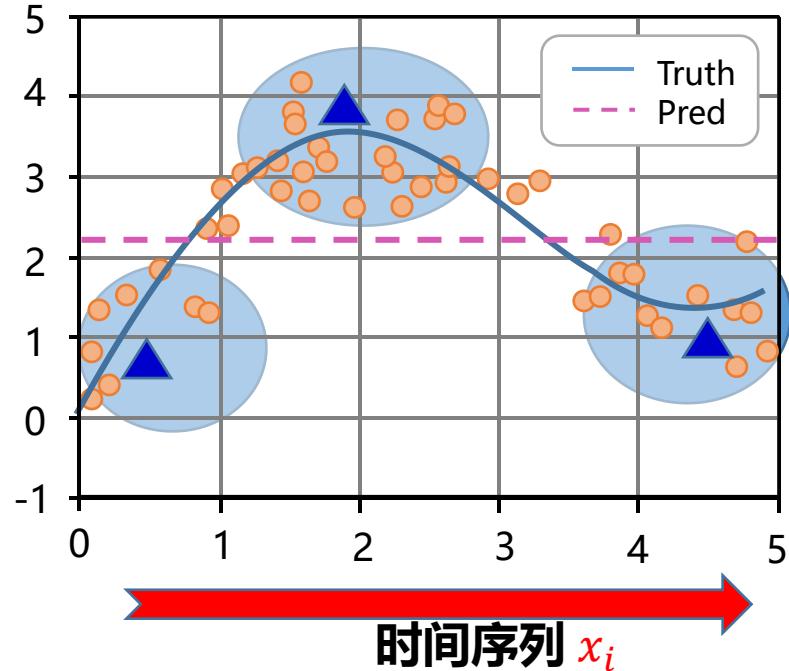
A solution is to get information from **the position of samples**, which can determine where is important:

$$f(x) = \sum_{i=1}^n \frac{K(x - x_i)}{\sum_{j=1}^n K(x - x_j)} y_i$$

where  $K$  is a kernel function. Then, we have

$$f(x) = \sum_{i=1}^n \alpha(x, x_i) y_i$$

where  $\alpha(x, x_i)$  denotes the weights of  $y_i$ .



Relation with **Kernel SVM**:

$$f(\mathbf{x}) = \sum_i^N \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b$$

<https://www.robots.ox.ac.uk/~az/lectures/ml/lect3.pdf>

# Motivation of Attention

## ■ Nonparametric attention pooling

$$f(x) = \sum_{i=1}^n \alpha(x, x_i) y_i$$

$$\alpha_i = \frac{K(x - x_i)}{\sum_{j=1}^n K(x - x_j)}$$

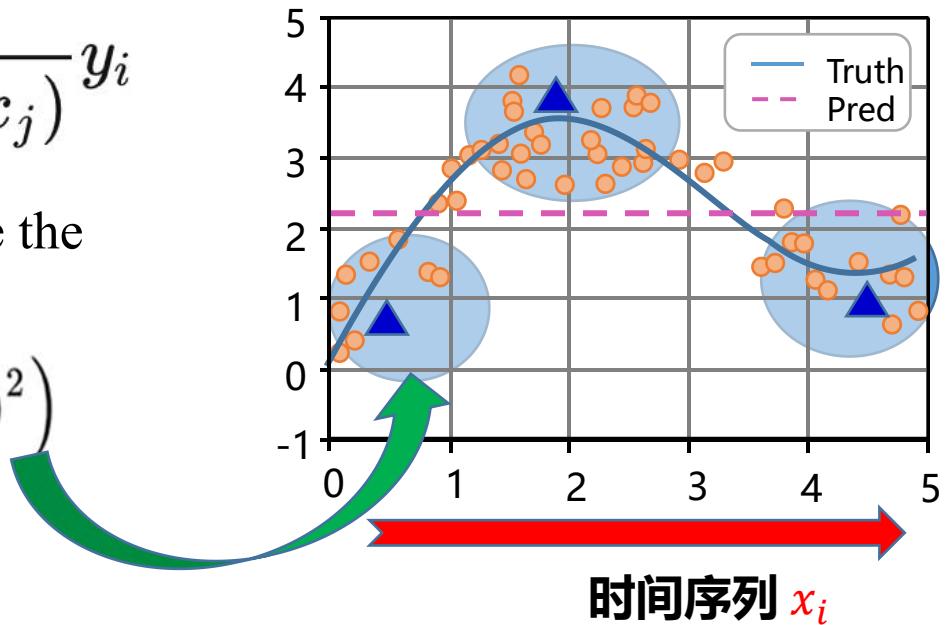
Normalization term

where  $\alpha(x, x_i)$  denotes the weights of  $y_i$ .

$$f(x) = \sum_{i=1}^n \frac{K(x - x_i)}{\sum_{j=1}^n K(x - x_j)} y_i$$

We can use **Gaussian kernel** to compute the distance of  $x$  to any point  $x_i$ :

$$K(x - x_i) = \exp\left(-\frac{1}{2}(x - x_i)^2\right)$$

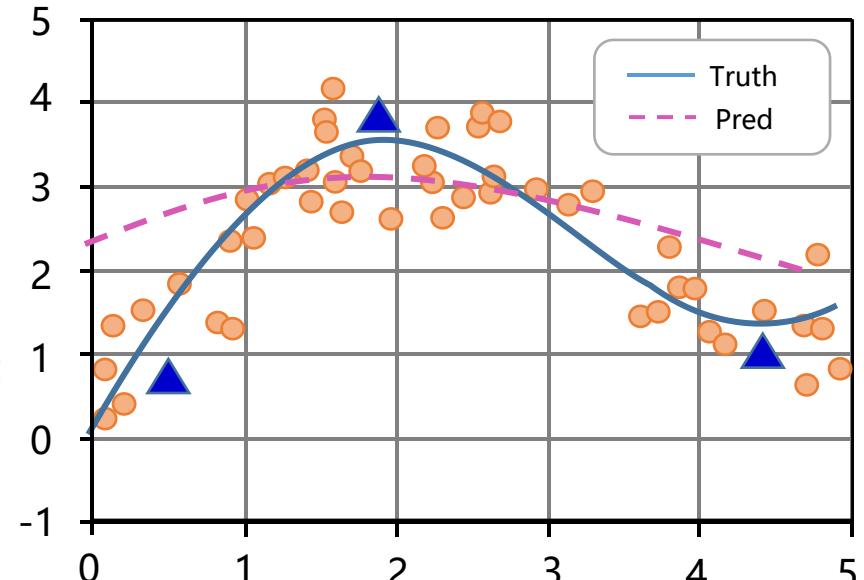


# Motivation of Attention

## ■ Nonparametric attention pooling

Then, the prediction becomes

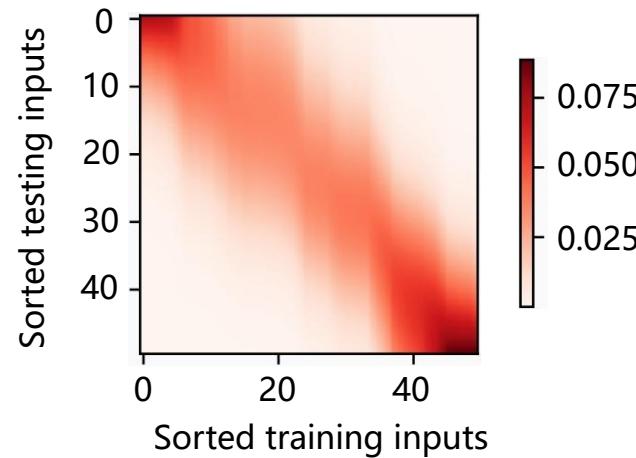
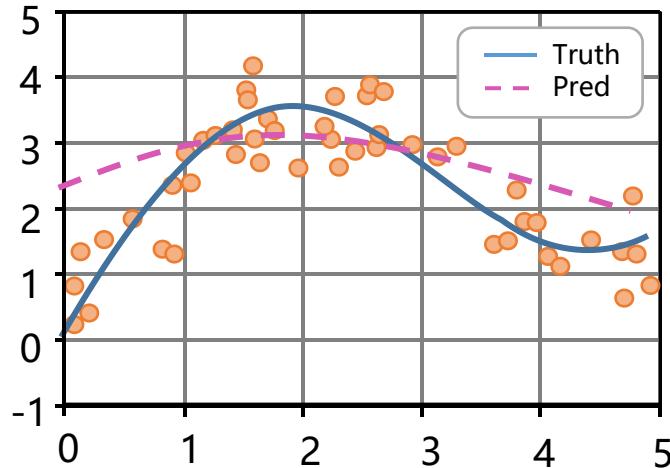
$$\begin{aligned}f(x) &= \sum_{i=1}^n \alpha(x, x_i) y_i \\&= \sum_{i=1}^n \frac{\exp\left(-\frac{1}{2}(x - x_i)^2\right)}{\sum_{j=1}^n \exp\left(-\frac{1}{2}(x - x_j)^2\right)} y_i \\&= \sum_{i=1}^n \text{softmax}\left(-\frac{1}{2}(x - x_i)^2\right) y_i.\end{aligned}$$



- The performance is improved with **nonparametric** attention pooling.

# Core Idea of Attention

## ■ Nonparametric attention pooling



$$f(x) = \sum_{i=1}^n \text{softmax}\left(-\frac{1}{2}(x - x_i)^2\right)y_i.$$

- When  $x_i$  is more **closer** to  $x$ , the weight of corresponding  $y_i$  is larger, also means **more important**.
- $x, x_i, y_i$  correspond to **Query, Key, Value** in Attention Mechanism,

# Core Idea of Attention

## ■ Parametric attention pooling

A variant of attention pooling: Project  $(x - x_i)$  by trainable matrix  $w$ .

Nonparametric  
attention pooling

$$f(x) = \sum_{i=1}^n \text{softmax}\left(-\frac{1}{2}(x - x_i)^2\right) y_i.$$

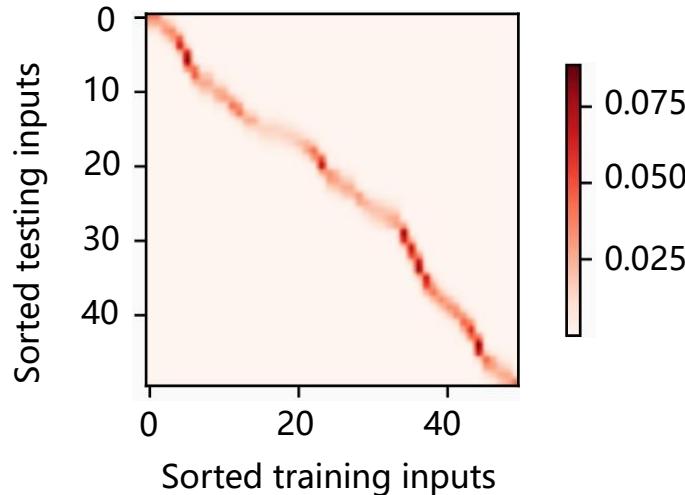
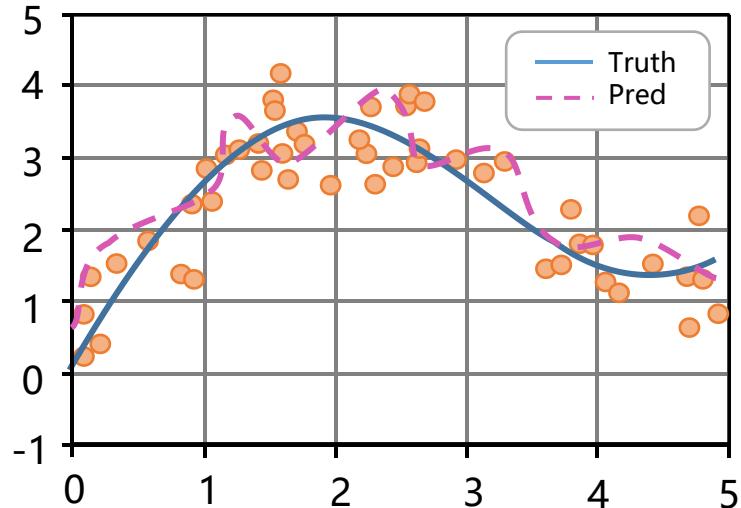


Parametric  
attention pooling

$$\begin{aligned} f(x) &= \sum_{i=1}^n \alpha(x, x_i) y_i \\ &= \sum_{i=1}^n \frac{\exp\left(-\frac{1}{2}((x - x_i) w)^2\right)}{\sum_{j=1}^n \exp\left(-\frac{1}{2}((x - x_j) w)^2\right)} y_i \\ &= \sum_{i=1}^n \text{softmax}\left(-\frac{1}{2}((x - x_i) w)^2\right) y_i. \end{aligned}$$

# Core Idea of Attention

## ■ Parametric attention pooling



$$f(x) := \sum_{i=1}^n \text{softmax}\left(-\frac{1}{2}((x - x_i) w)^2\right) y_i.$$

- **Parametric** attention pooling can make attention **more concentration**, but probability **over-fitting**.

## Summary

- Attention pooling can be parametric or nonparametric.
- In attention pooling, each  $Value_i$  corresponds to a weight  $\alpha(Query, Key_i)$ .
- Formulation of Attention Function

$$f(\mathbf{q}, (\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_n, \mathbf{v}_n)) = \sum_{i=1}^n \alpha(\mathbf{q}, \mathbf{k}_i) \mathbf{v}_i$$

where  $\mathbf{q}, \mathbf{k}, \mathbf{v}$  denote **Query**, **Key** and **Value**, respectively.

$\alpha(\mathbf{q}, \mathbf{k}_i)$  is scoring function.

# Core Idea of Attention

## ■ An example of scoring function

$$f(\mathbf{q}, (\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_n, \mathbf{v}_n)) = \sum_{i=1}^n \alpha(\mathbf{q}, \mathbf{k}_i) \mathbf{v}_i$$

$\alpha(\mathbf{q}, \mathbf{k}_i)$  is usually *softmax* function:

$$\alpha(\mathbf{q}, \mathbf{k}_i) = softmax(\alpha(\mathbf{q}, \mathbf{k}_i)) = \frac{\exp(\alpha(\mathbf{q}, \mathbf{k}_i))}{\sum_{j=1}^m \exp(\alpha(\mathbf{q}, \mathbf{k}_j))} \in \mathbb{R}$$

## ■ More forms of scoring function $\alpha(\mathbf{q}, \mathbf{k}_i)$

➤ Additive Attention:

$$\alpha(\mathbf{q}, \mathbf{k}) = \mathbf{w}_v^\top \tanh(\mathbf{W}_q \mathbf{q} + \mathbf{W}_k \mathbf{k}) \in \mathbb{R}$$

➤ Scaled Dot-Product Attention:

SMIL 内部资料 请勿外泄  $\alpha(\mathbf{q}, \mathbf{k}) = \mathbf{q}^\top \mathbf{k} / d$

# Implementation of Attention Layer

## ■ Additive Attention:

$$\alpha(\mathbf{q}, \mathbf{k}) = \mathbf{w}_v^\top \tanh(\mathbf{W}_q \mathbf{q} + \mathbf{W}_k \mathbf{k}) \in \mathbb{R}$$

```
class AdditiveAttention(nn.Module):
    def __init__(self, key_size, query_size, num_hiddens, dropout, **kwargs):
        super(AdditiveAttention, self).__init__(**kwargs)
        self.W_k = nn.Linear(key_size, num_hiddens, bias=False)
        self.W_q = nn.Linear(query_size, num_hiddens, bias=False)
        self.w_v = nn.Linear(num_hiddens, 1, bias=False)
        self.dropout = nn.Dropout(dropout)

    def forward(self, queries, keys, values, valid_lens):
        # 下一步将Q和K分别对独自的W进行点乘
        queries, keys = self.W_q(queries), self.W_k(keys)
        # 输出的: queries[batch_size, num_query, num_hidden]
        # 输出的: keys[batch_size, num_keys, num_hidden]
        features = queries.unsqueeze(2) + keys.unsqueeze(1)
        # broadcasting 相加, features[batch_size, num_query, num_keys, num_hidden]
        features = torch.tanh(features)
        scores = self.w_v(features).squeeze(-1)
        # scores[batch_size, num_query, num_keys]
        self.attention_weights = masked_softmax(scores, valid_lens)
        return torch.bmm(self.dropout(self.attention_weights), values)
        # values 的长度和keys 的长度是一样的。
```

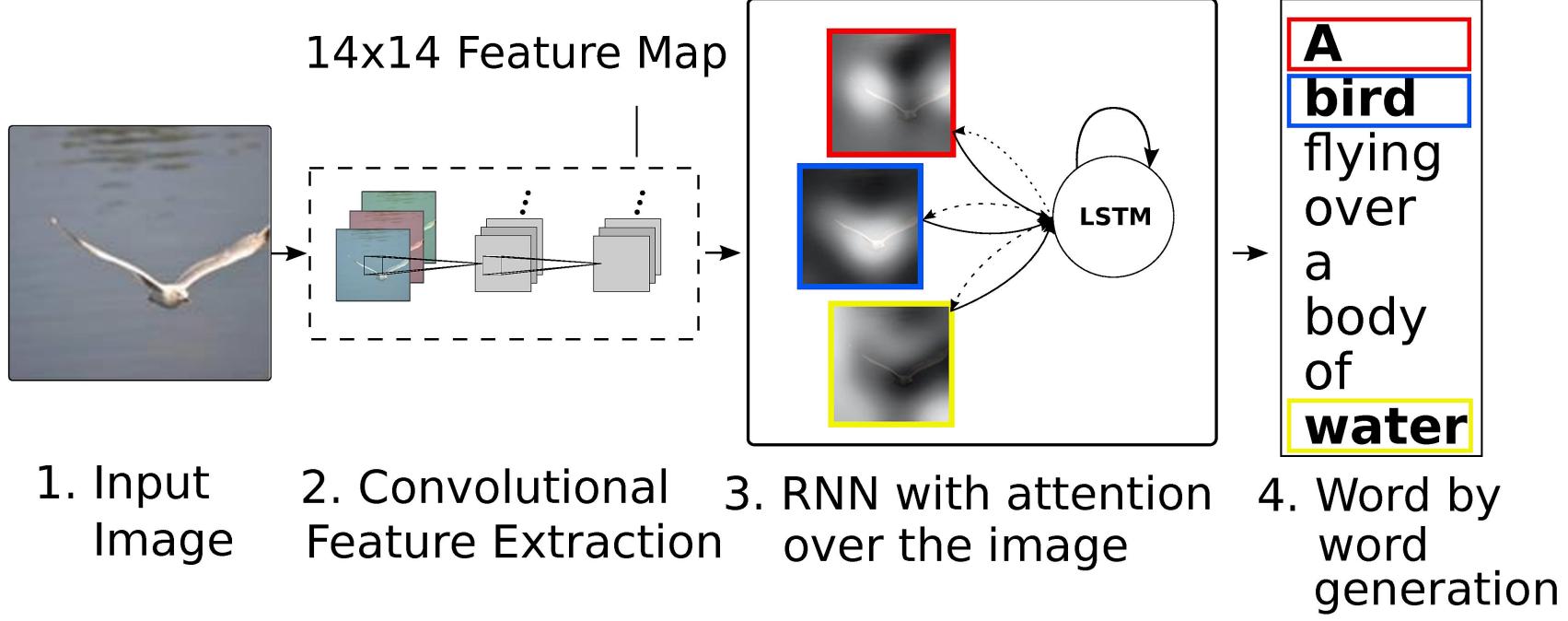
# Implementation of Attention Layer

## ■ Scaled Dot-Product Attention:

$$\alpha(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k} / d$$

```
class DotProductAttention(nn.Module):
    def __init__(self, dropout, **kwargs):
        super(DotProductAttention, self).__init__(**kwargs)
        self.dropout = nn.Dropout(dropout)
    def forward(self, queries, keys, values, valid_lens=None):
        d = queries.shape[-1]
        scores = torch.bmm(queries, keys.transpose(1, 2)) / math.sqrt(d)
        self.attention_weights = masked_softmax(scores, valid_lens)
        return torch.bmm(self.dropout(self.attention_weights), values)
```

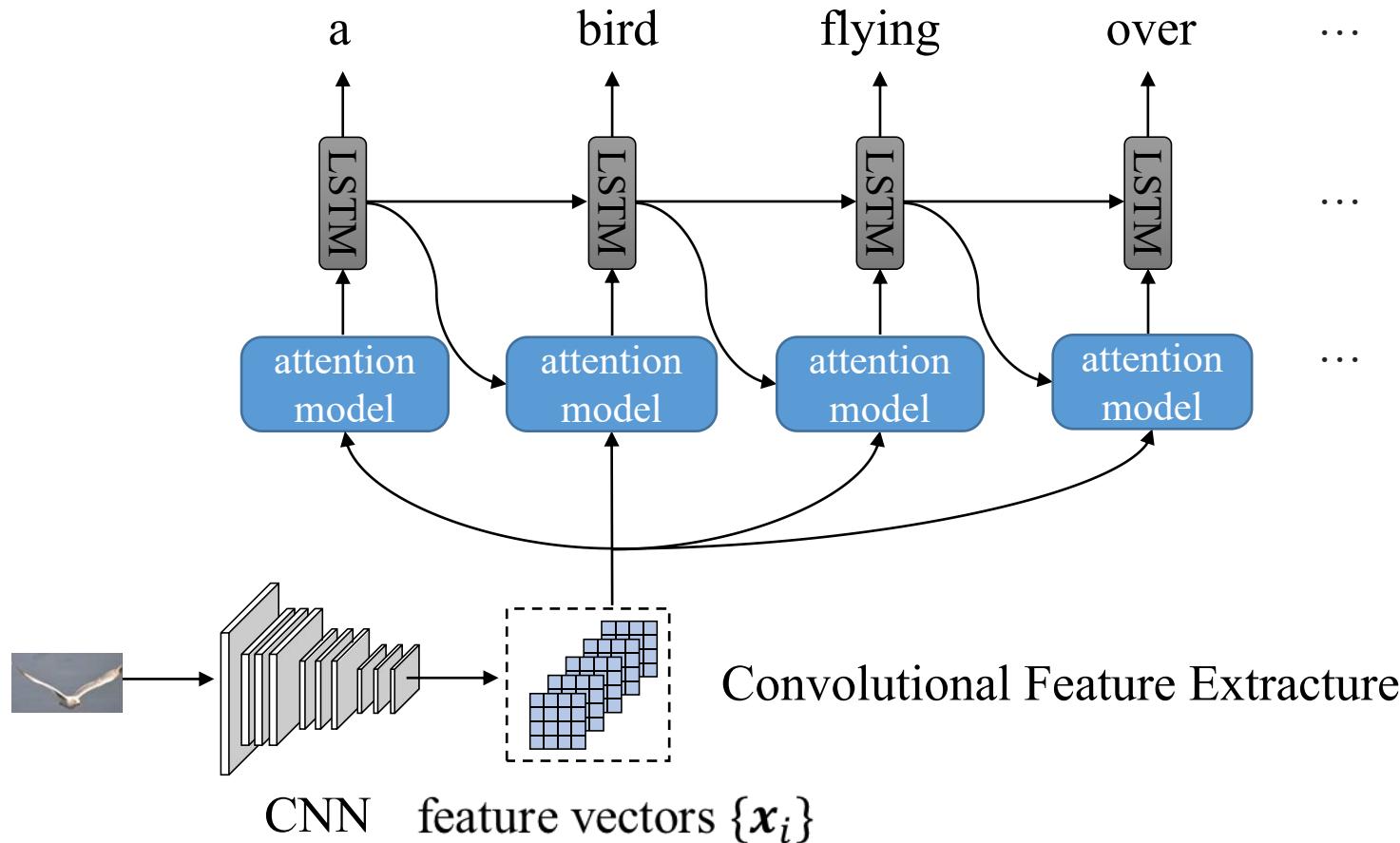
# Attention in Image Captioning



SMIL内部资料 请勿外泄

# Attention in Image Caption

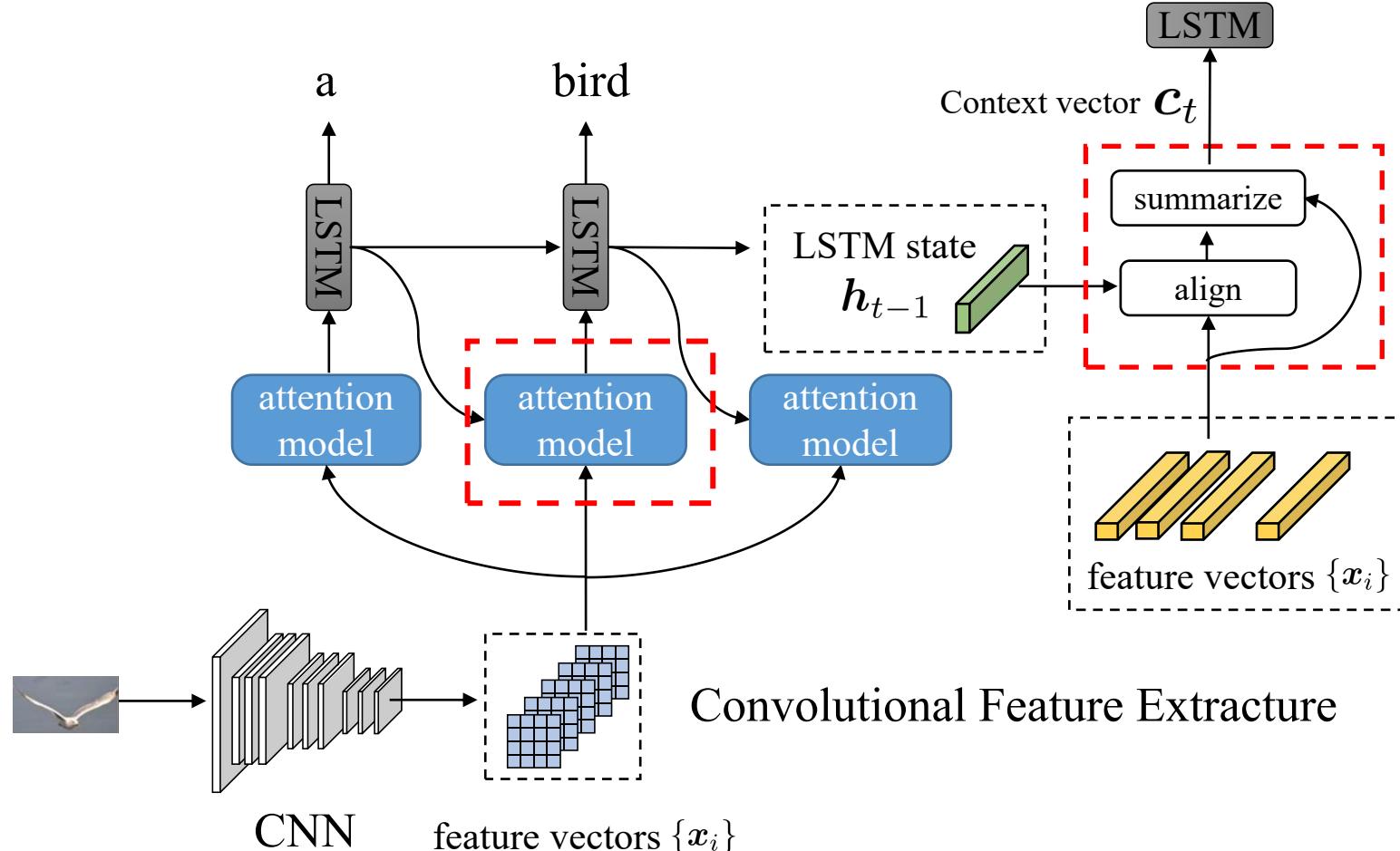
LSTM with attention for image captioning.



SMIL 内部资料 请勿外泄

# Attention in Image Caption

## The attention module



SMIL内部资料 请勿外泄

# Attention in Image Caption

align decodes which encoder state  $x_i$  should pay attention to according the decoder state  $h_{t-1}$ .

- Compute the alignment score between  $x_i$  and  $h_{t-1}$ :

$$e_{t,i} = \mathbf{W}^T \tanh(\mathbf{W}_x^T \mathbf{x}_i + \mathbf{W}_h^T \mathbf{h}_{t-1}).$$

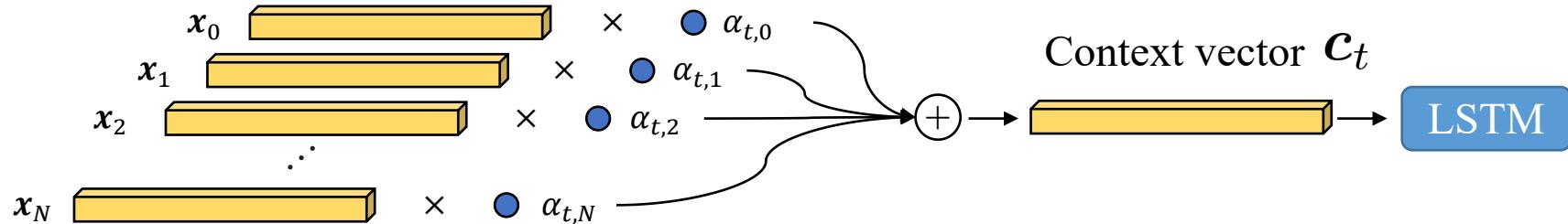
- Compute the attention weight for  $x_i$  :

$$\alpha_{t,i} = \frac{\exp(e_{t,i})}{\sum_j \exp(e_{t,j})}$$

$\mathbf{W}, \mathbf{W}_h, \mathbf{W}_s$  are jointly trained within the whole architecture.

# Attention in Image Caption

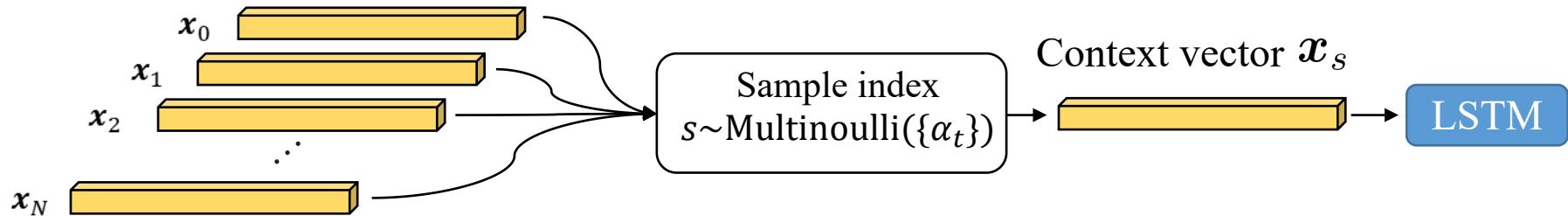
summarize



- **Soft attention:** summarize all encoder states into a context vector with the attention weights, and feed it into the decoder LSTM.

# Attention in Image Caption

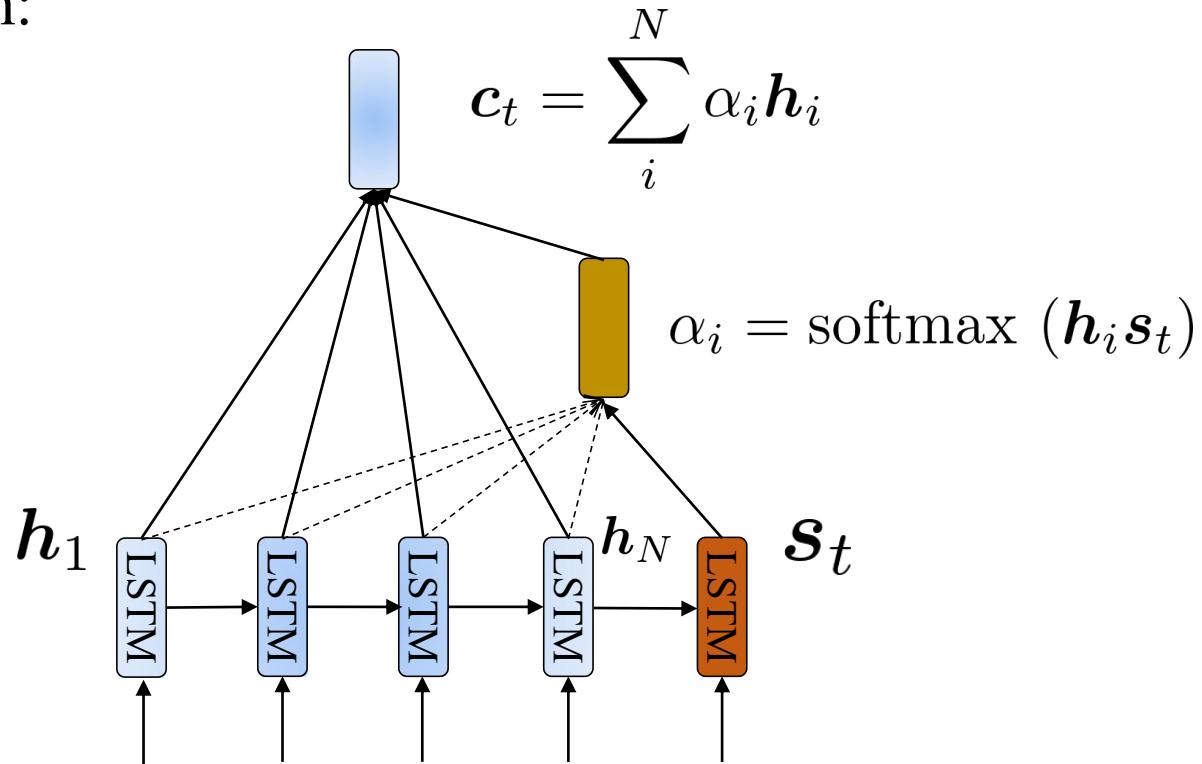
summarize



- **Hard attention:** parameterize a multinoulli distribution wth the attention weights, and sample an feature vector  $\mathbf{x}_s$  from  $\{\mathbf{x}_i\}$ according to it.

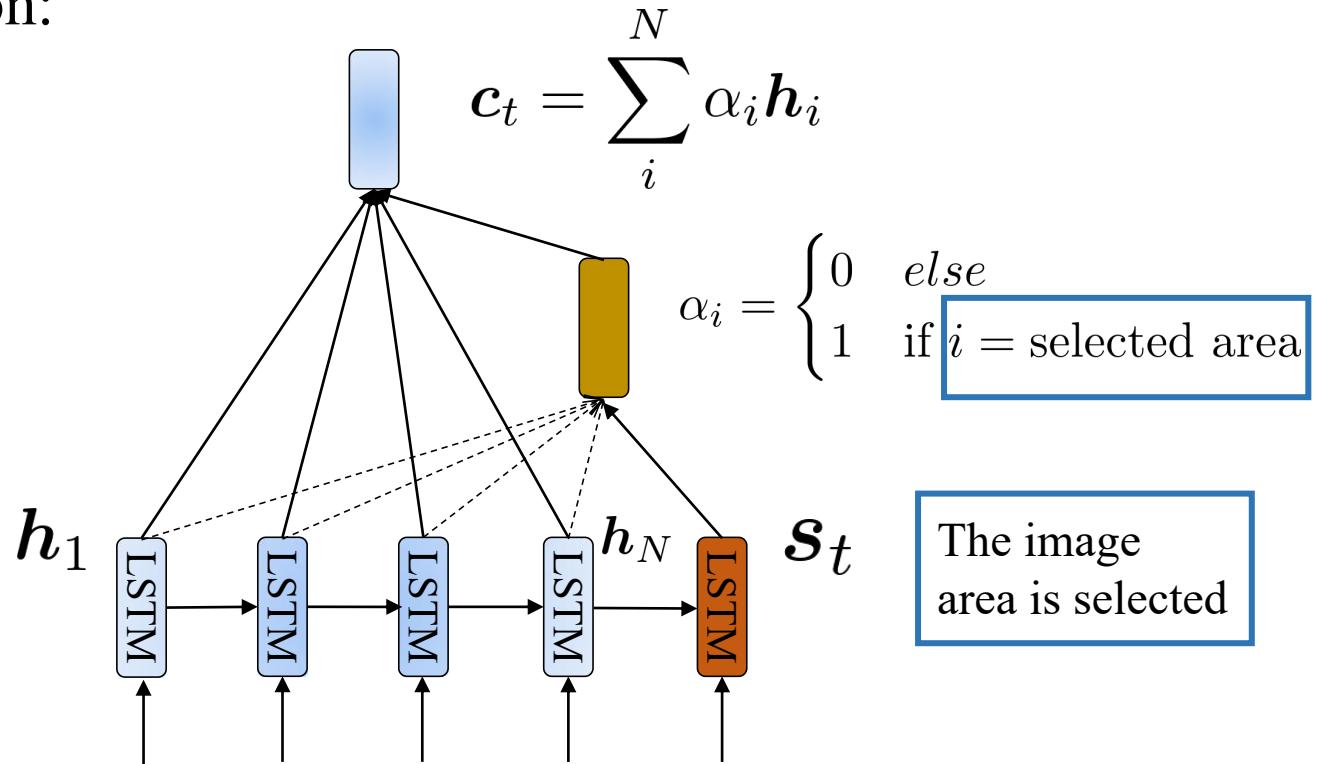
# Variants of Attention

Soft attention:



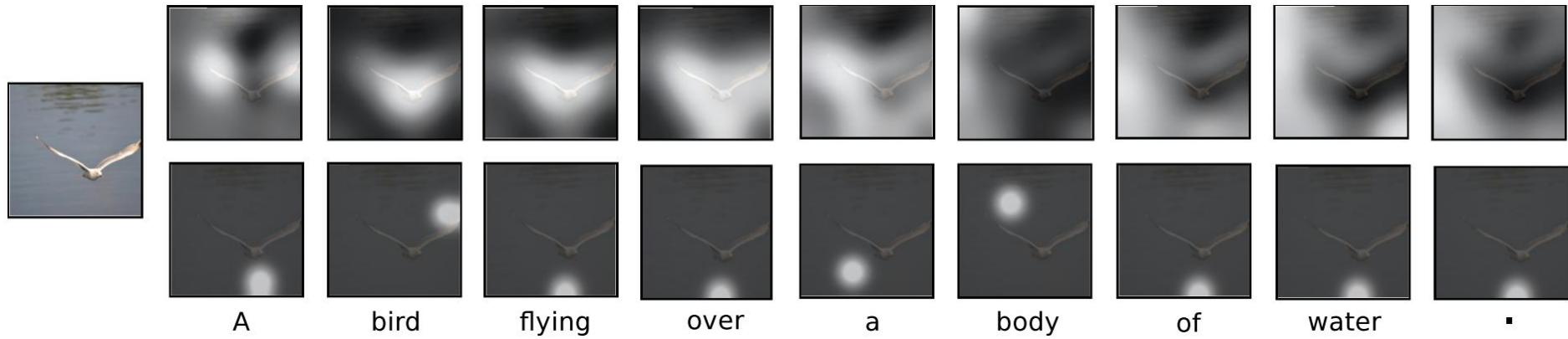
# Variants of Attention

Hard attention:

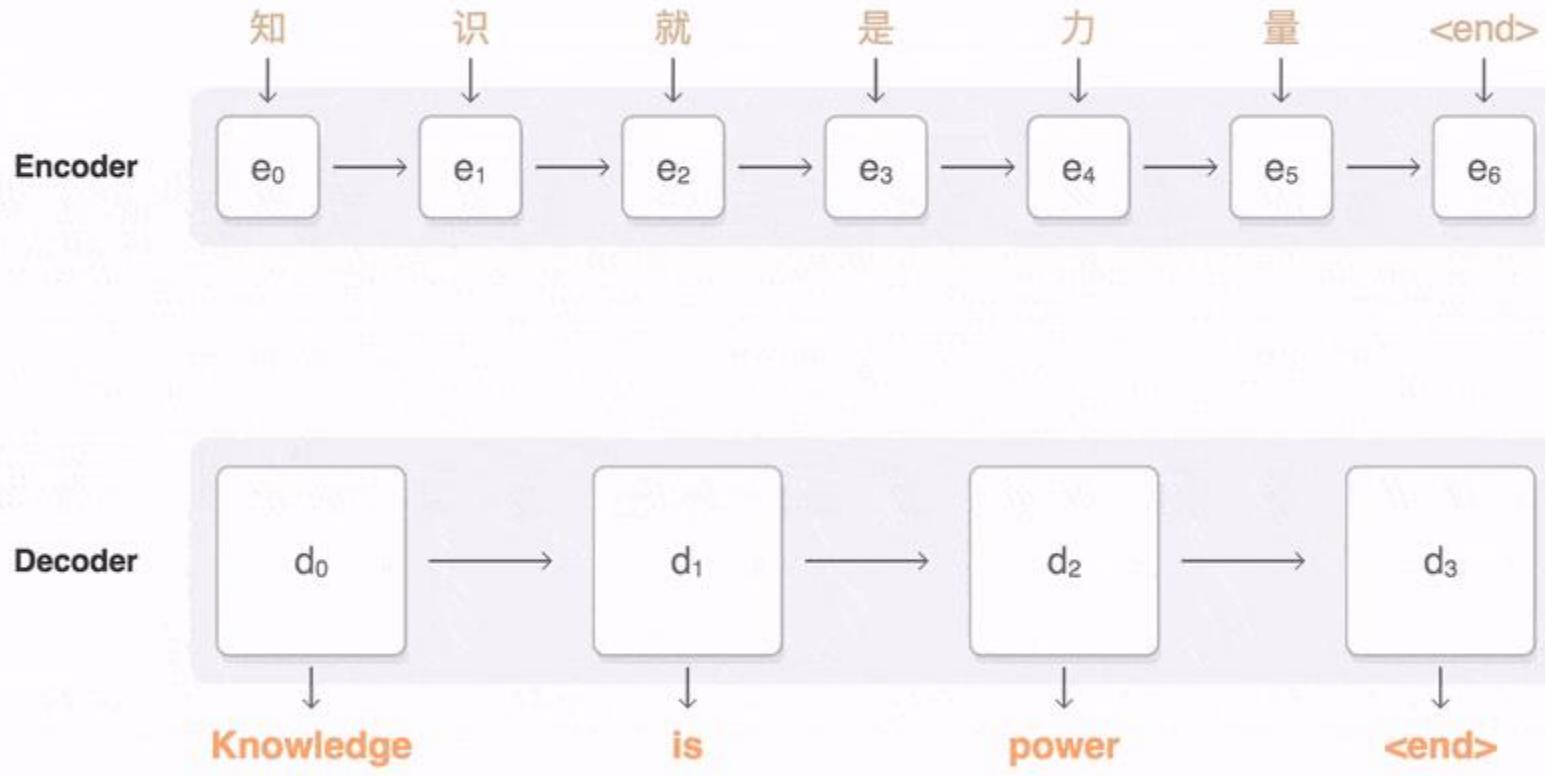


# Attention in Image Caption

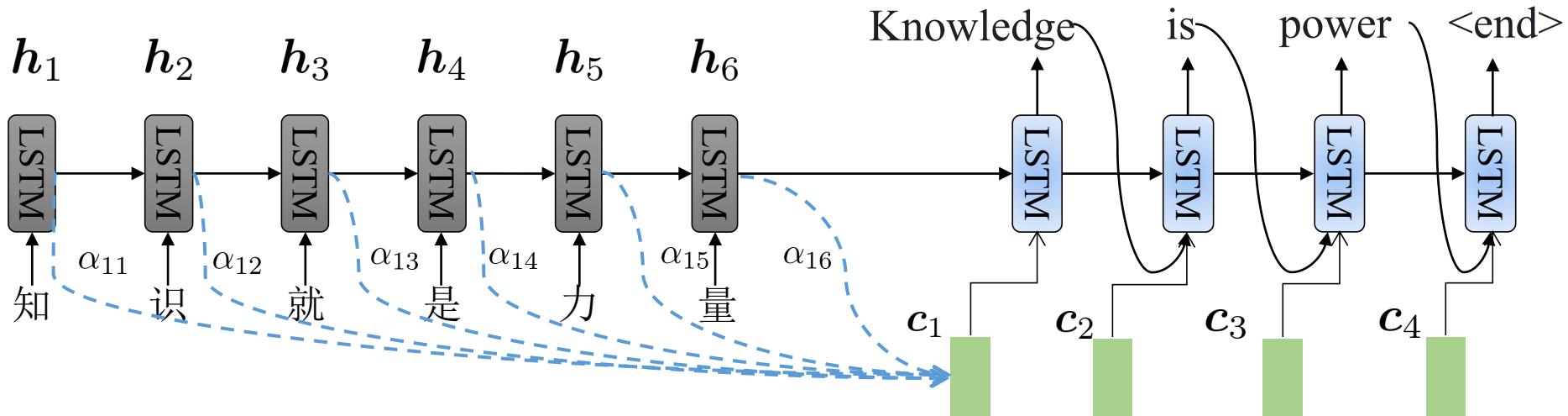
Soft attention vs. Hard attention



# Attention in Machine Translation



# Attention in Machine Translation



知              识              就              是              力              量              Context vector

$$h_1 \times \alpha_{11} + h_2 \times \alpha_{12} + h_3 \times \alpha_{13} + h_4 \times \alpha_{14} + h_5 \times \alpha_{15} + h_6 \times \alpha_{16} = c_1 \rightarrow \text{Knowledge}$$

$$h_1 \times \alpha_{41} + h_2 \times \alpha_{42} + h_3 \times \alpha_{43} + h_4 \times \alpha_{44} + h_5 \times \alpha_{45} + h_6 \times \alpha_{46} = c_4 \rightarrow \text{is}$$

# Attention in Machine Translation

How to compute the attention weights  $\alpha$ ?

- Compute the alignment score between  $h_i$  and  $s_{t-1}$  :

$$e_{t,i} = \mathbf{W}^T \tanh(\mathbf{W}_h^T h_i + \mathbf{W}_s^T s_{t-1})$$

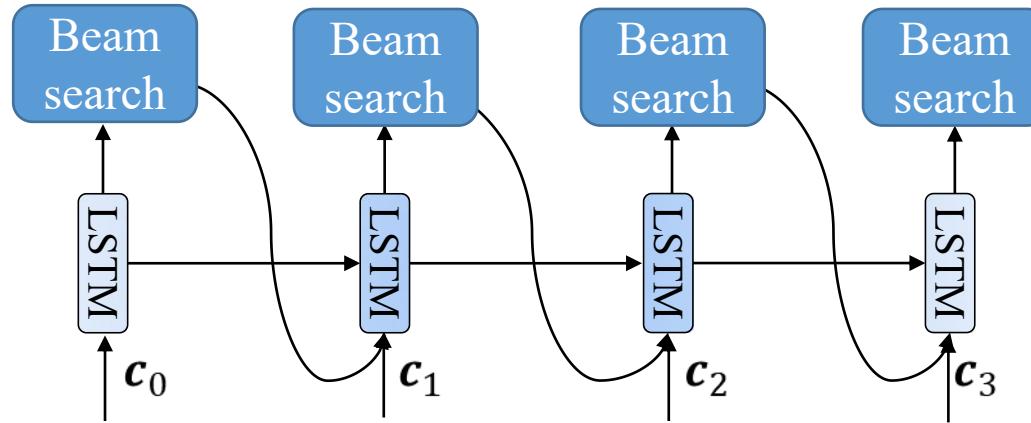
- Compute the attention weight for  $h_i$  :

$$\alpha_{t,i} = \frac{\exp(e_{t,i})}{\sum_j \exp(e_{t,j})}$$

$\mathbf{W}, \mathbf{W}_h, \mathbf{W}_s$  are jointly trained within the whole architecture.

# Attention in Machine Translation

Inference:



- Beam Search: at every step, select the sequences with top- $k$  probabilities, and feed its last word into the next time step.

# Outline

1 Backgrounds

2 Recurrent Neural Networks (RNNs)

3 Attention Mechanisms

4 Attention is All You Need: Transformer

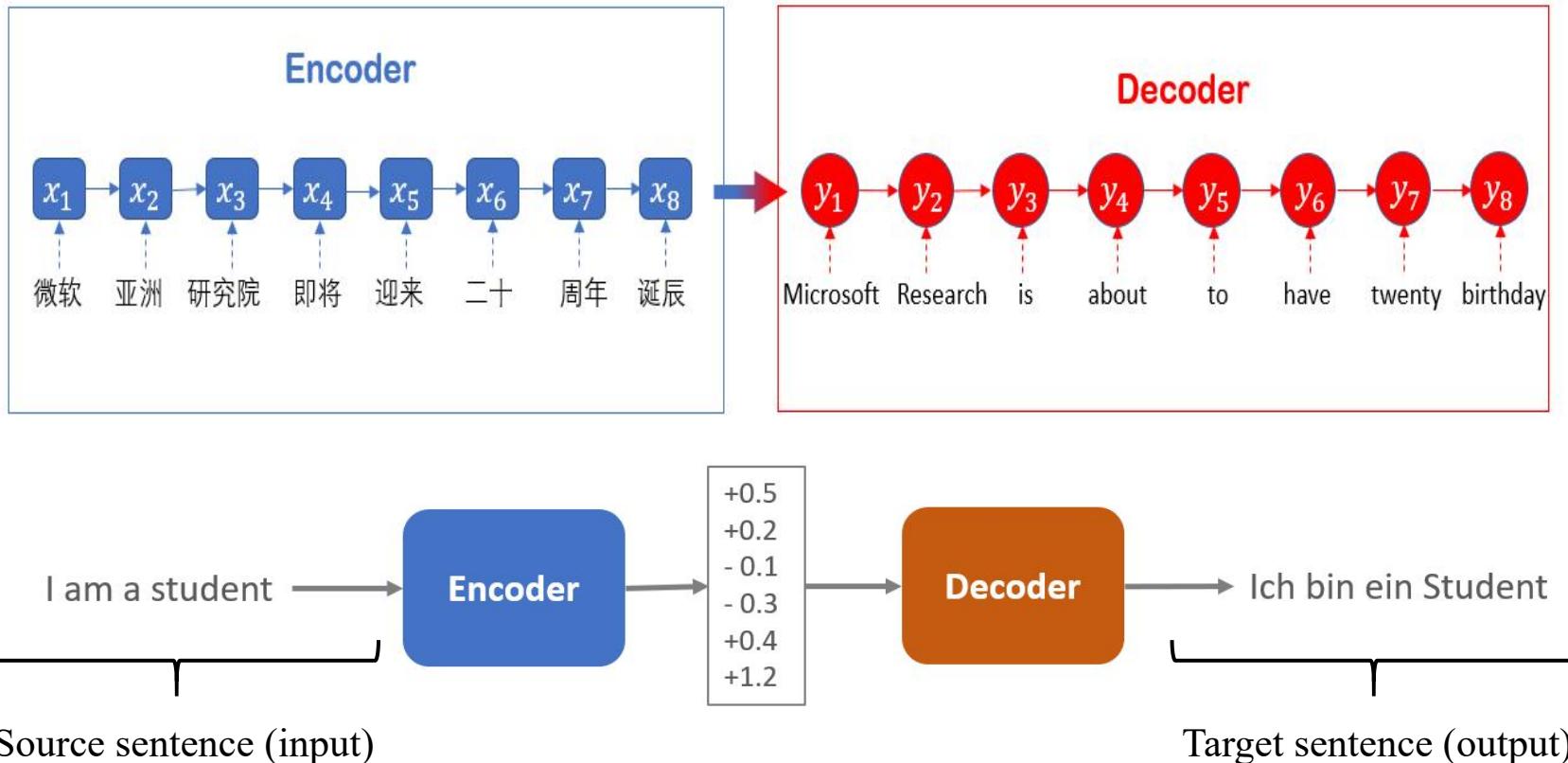
5 Long Short-Term Memory (LSTM) Networks

6 BERT(Bidirectional Encoder Representation from Transformers)

7 Conclusions

# Motivation: Neural Machine Translation (NMT)

- Definition: an approach to **machine translation** that uses an **artificial neural network** to predict the likelihood of a sequence of words.



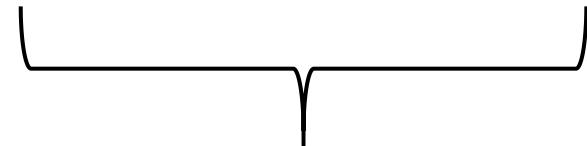
SMIL内部资料 请勿外泄

# Neural Machine Translation (NMT)

How to predict the **likelihood** of a sequence of words?

- NMT directly calculates  $P(\mathbf{y}|\mathbf{x})$ :

$$P(\mathbf{y}|\mathbf{x}) = P(\mathbf{y}_1|\mathbf{x})P(\mathbf{y}_2|\mathbf{y}_1, \mathbf{x})P(\mathbf{y}_3|\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) \dots P(\mathbf{y}_T|\mathbf{y}_1, \dots, \mathbf{y}_{T-1}, \mathbf{x})$$

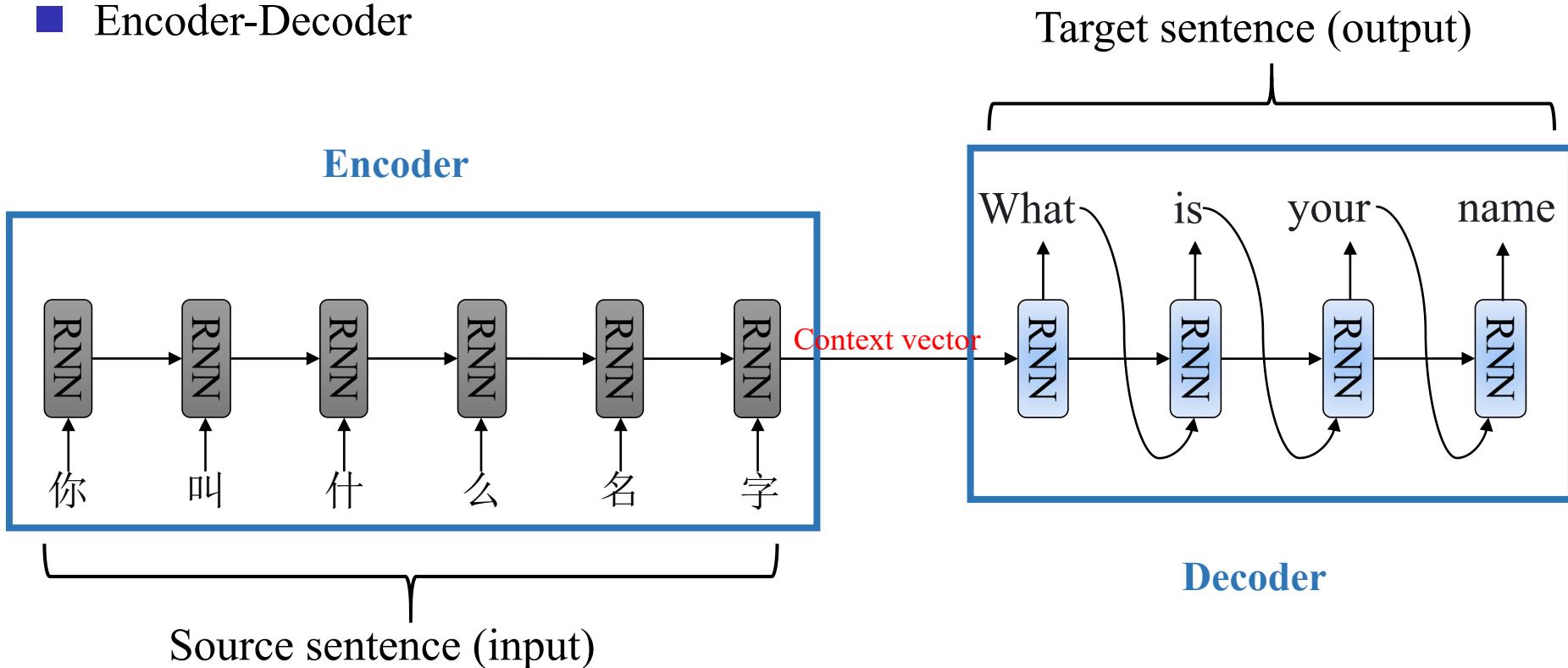


Probability of next target word, given target words so far and source sentence  $\mathbf{x}$

- Input: encode of the source sentence  $\mathbf{x}_1, \dots, \mathbf{x}_T$ .
- Output: generate target sentence  $\mathbf{y}_1, \dots, \mathbf{y}_T$ .

# Neural Machine Translation (NMT)

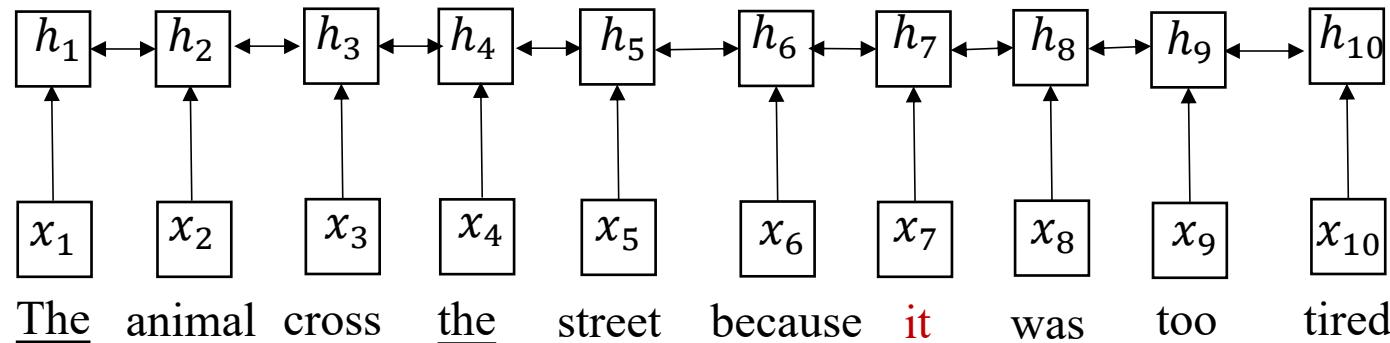
## ■ Encoder-Decoder



- Encoder RNN produces an encoding of the source sentence into a context vector.
- Decoder RNN is a language Model that generates target sentence from the context vector iteratively.

# Motivation

RNN: difficult to capture long-distance dependencies due to **vanishing gradient**

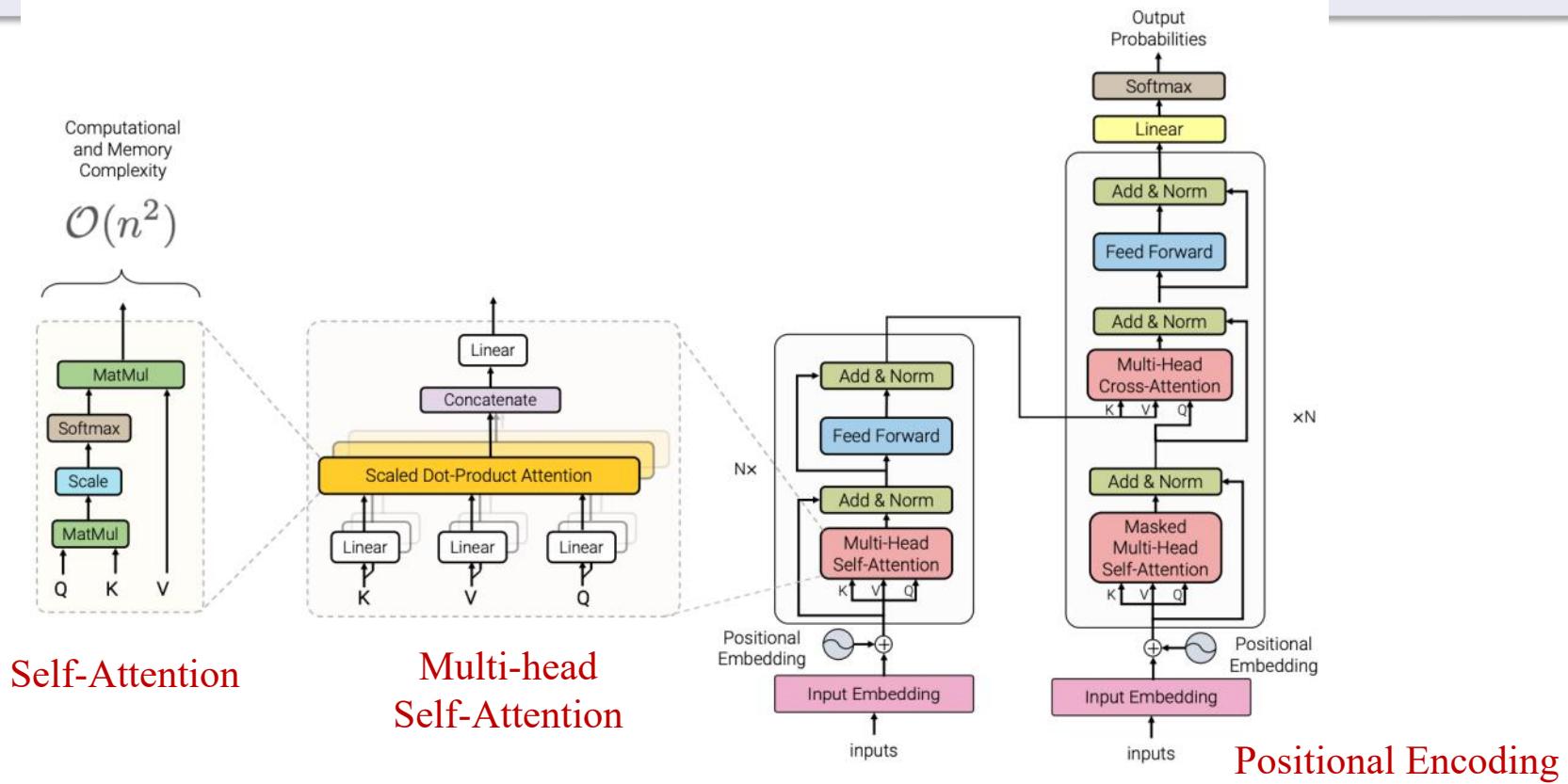


Bidirectional RNN modeling process

**Self-Attention in Transformer** can solve the long-distance dependence problem by calculating the direct dependence between words.

# Transformer

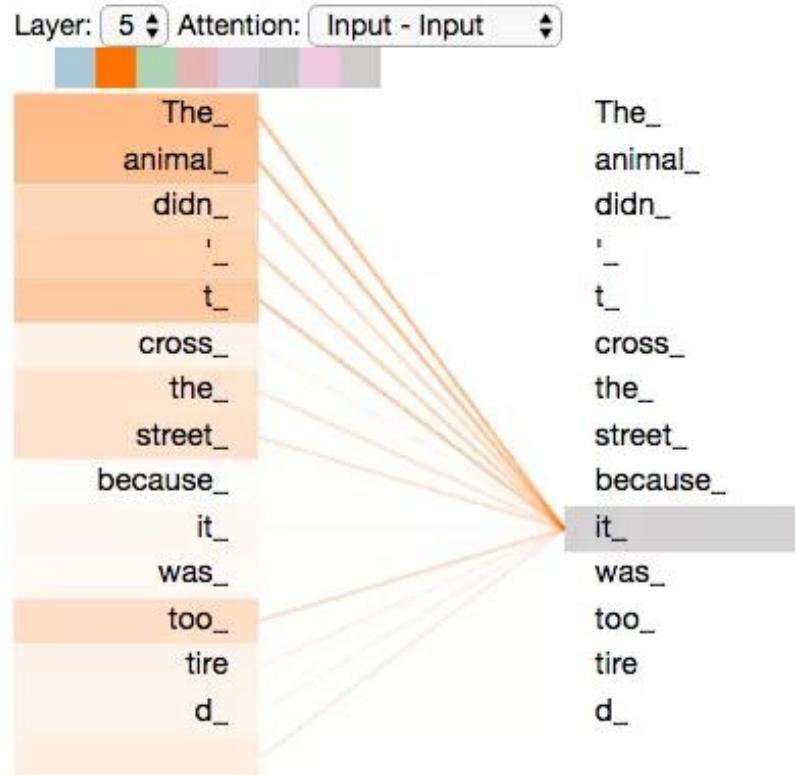
- Transformer has a encoder-decoder architecture similar to the previous RNN models.
- Both encoder and decoder use N stacked self-attention layers.



SMIL内部资料 请勿外泄

# Self-Attention in Transformer

**Self-Attention in Transformer** can bake the “understanding” of other relevant words into the one we’re currently processing.

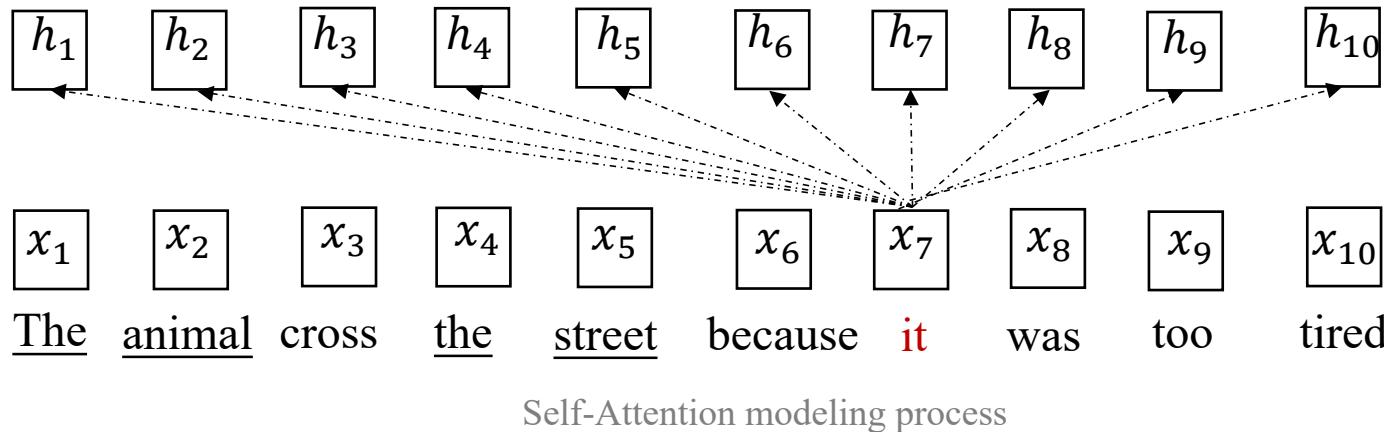


As we are encoding the word "it", part of the attention mechanism was focusing on "The Animal", and baked a part of its representation into the encoding of "it".

SMIL内部资料 请勿外泄

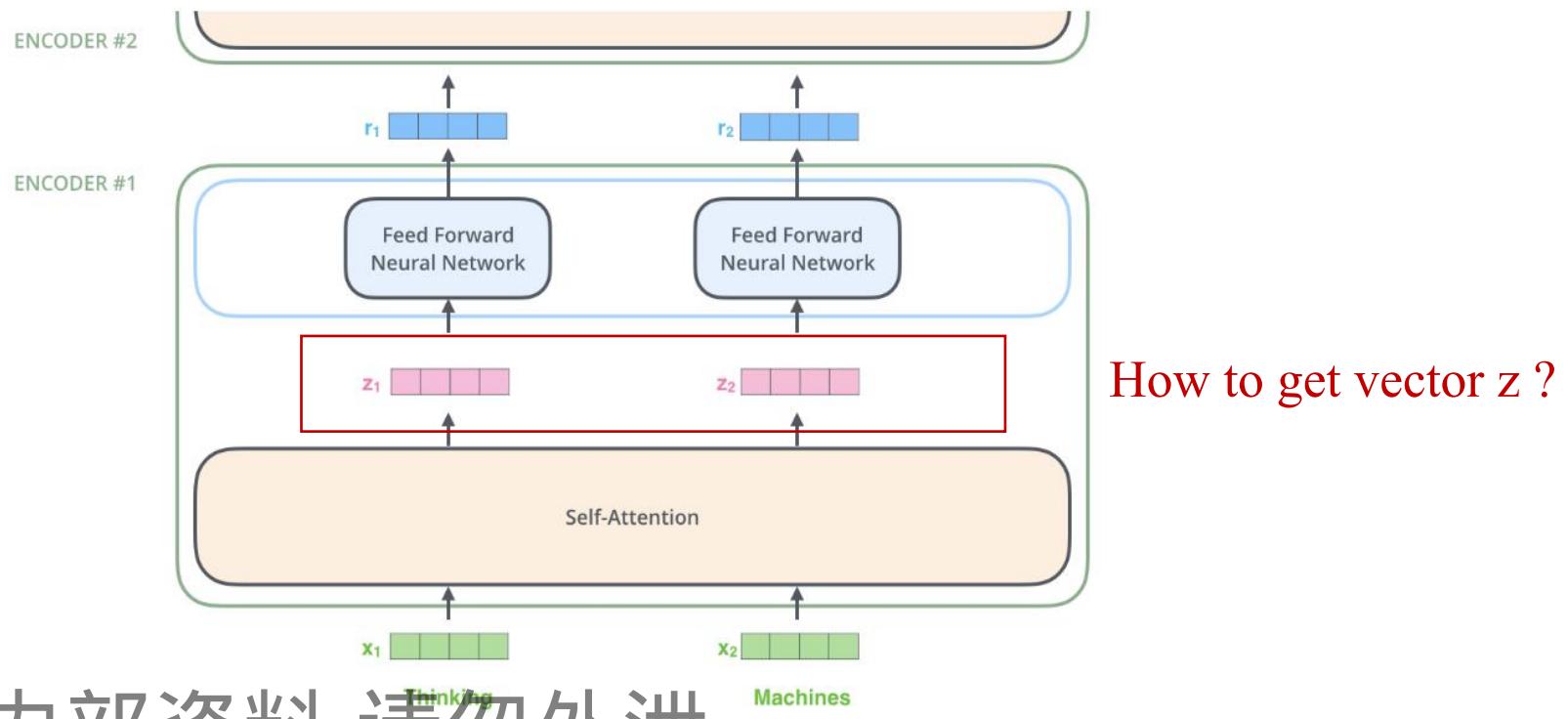
# Self-Attention in Transformer

- Self-Attention calculates attention between each word and all words.
- The maximum path length between each word is only 1.
- Self-Attention can **capture long-distance dependencies**.



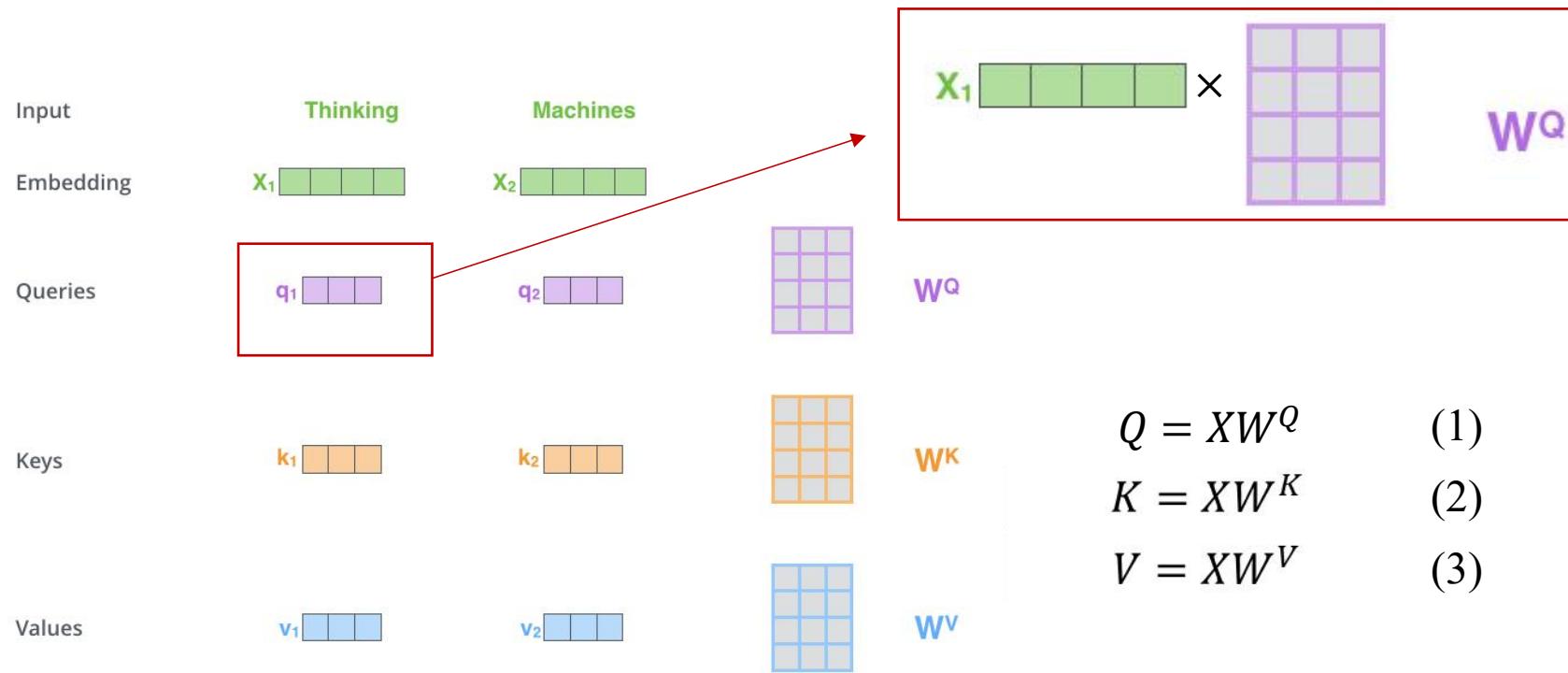
# Self-Attention in Transformer

- The encoder's inputs first flow through a **self-attention layer**.
- It looks at other words in the sentence when it encodes a specific word.



SMIL内部资料 请勿外泄

# Self-Attention in Transformer



Step 1. Create a **query vector**( $q_i$ ), a **key vector**( $k_i$ ), and a **value vector**( $v_1$ ) from each of the encoder's input vectors.

# Self-Attention in Transformer

Input	Thinking		Machines	
Embedding	$x_1$		$x_2$	
Queries	$q_1$		$q_2$	
Keys	$k_1$		$k_2$	
Values	$v_1$		$v_2$	
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	

Step 2. Calculate a score between each word.

The score determines how much focus to place on other parts of the input sentence as we encode a word at a certain position.

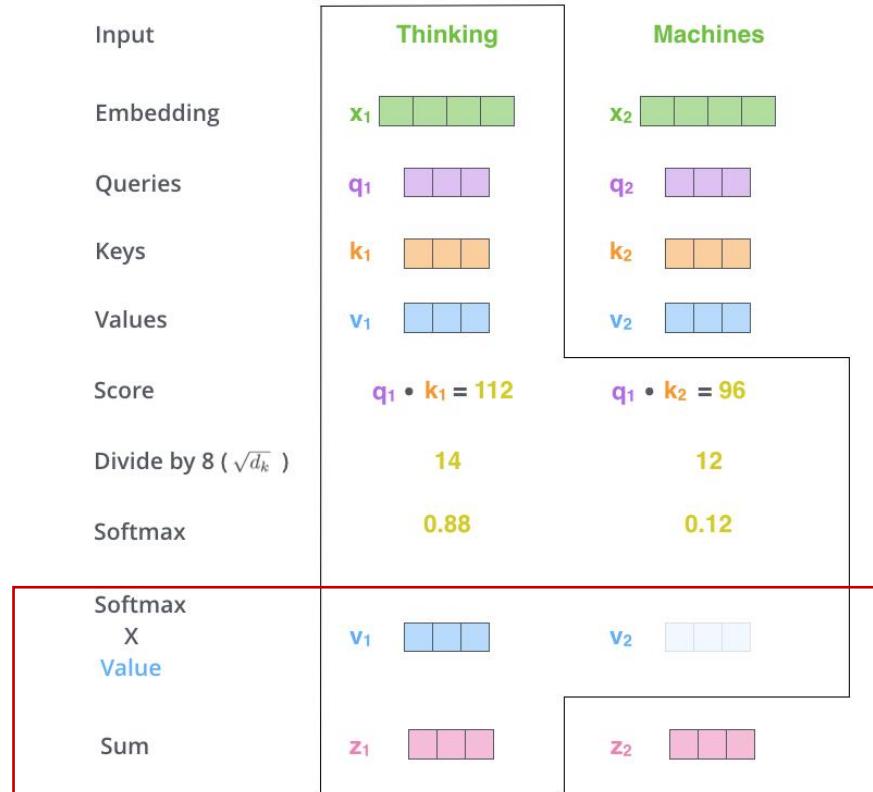
# Self-Attention in Transformer

Input	Thinking		Machines	
Embedding	$x_1$		$x_2$	
Queries	$q_1$		$q_2$	
Keys	$k_1$		$k_2$	
Values	$v_1$		$v_2$	
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	
Divide by $8 (\sqrt{d_k})$	14		12	
Softmax	0.88		0.12	

- Step 3. Divide the scores by the square root of the dimension of the key vectors.  
Step 4. SoftMax normalizes the scores so they're all positive and add up to 1.

SMIL内部资料 请勿外泄

# Self-Attention in Transformer



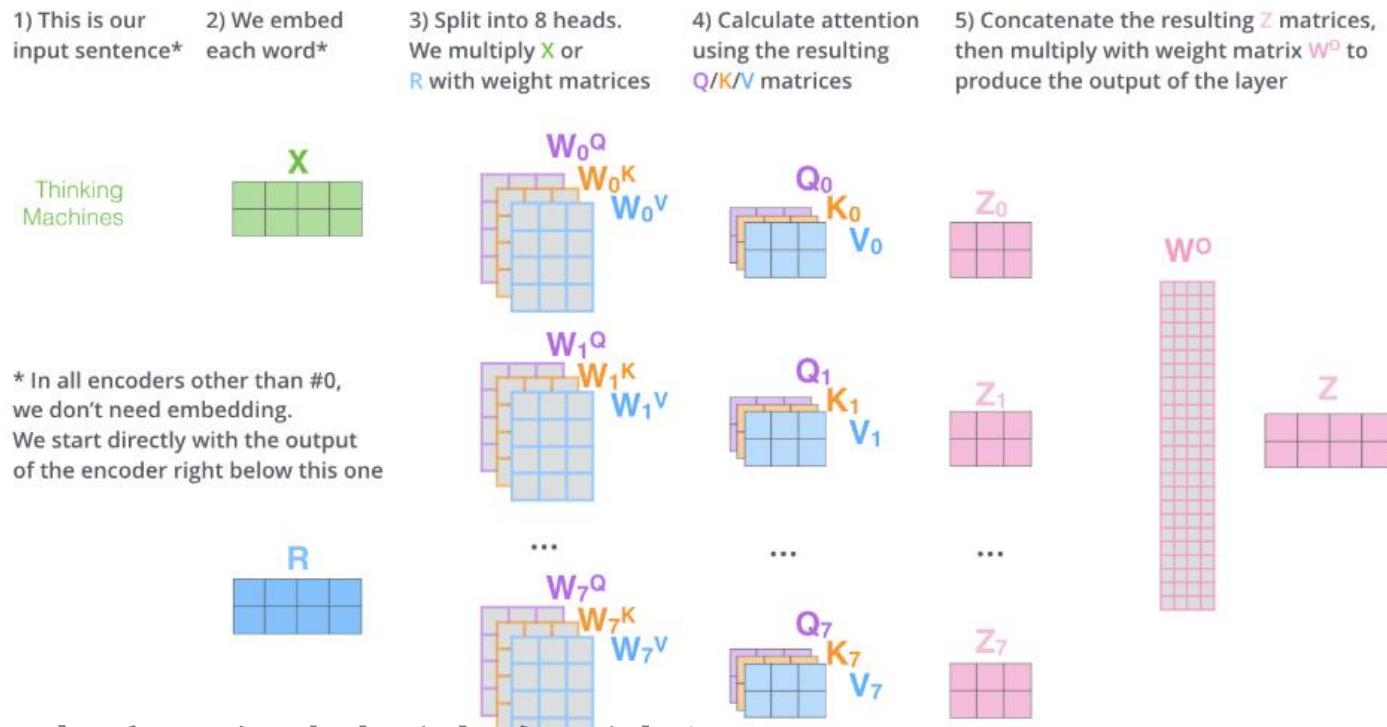
$$Z = \text{softmax} \left( \frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

Step 5. Multiply each value vector by the SoftMax score to keep intact the values of the word(s) we want to focus on, and drown-out irrelevant words

Step 6. Sum up the weighted value vectors to produce self-attention outputs( $Z_i$ )

# Multi-head Self-Attention

- Expands the model's ability to focus on **different positions**.
- Gives the attention layer multiple “representation subspaces”.



SMIL内部资料 请勿外泄

# Position Embedding in Transformer

## Why we need positional encoding in Transformer?

- A sentence of different word sequence convey a different meaning.
- All building blocks in Transformer such as self-attention is not aware of the word position information.

### Why position matters?

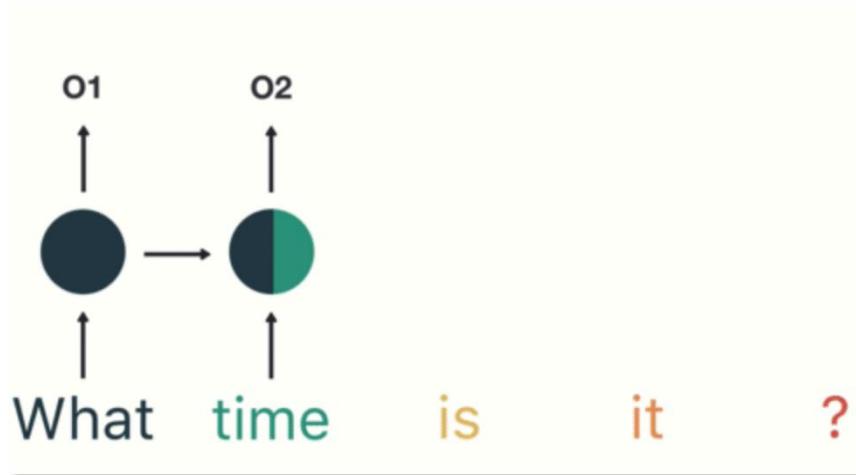
Even though she did **not** win the award, she was satisfied.

Even though she did win the award, she was **not** satisfied.

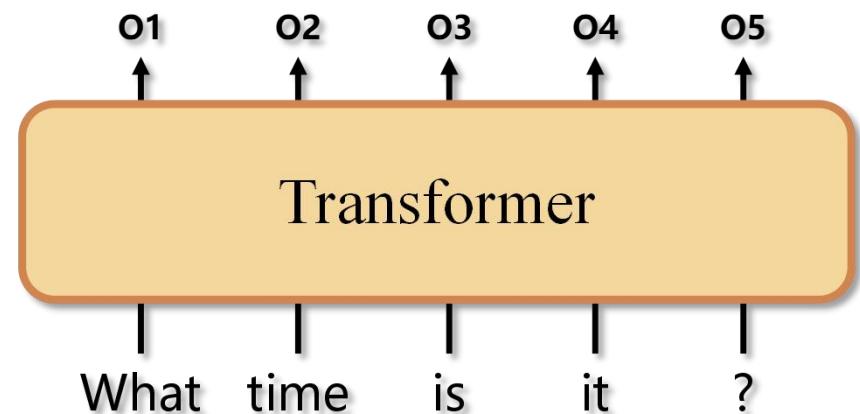
# Positional Encoding in Transformer

Think: Why does RNN has positional information?

- The words are inputted one-by-one into the RNN model.



In RNN, one word's representation is based on the previous ones.

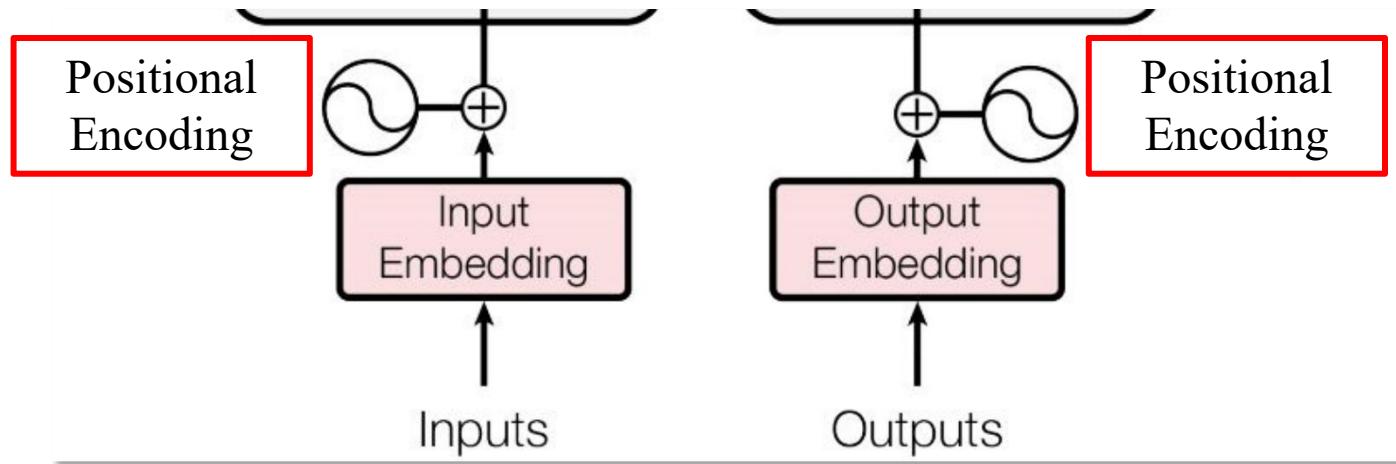


In Transformer, all words are inputted and processed simultaneously.

# Positional Encoding in Transformer

How to add the position information efficiently into Transformer?

- Add on the input level: **Positional Encoding Layer**.
- Why efficient? Does not modify the Transformer architecture.

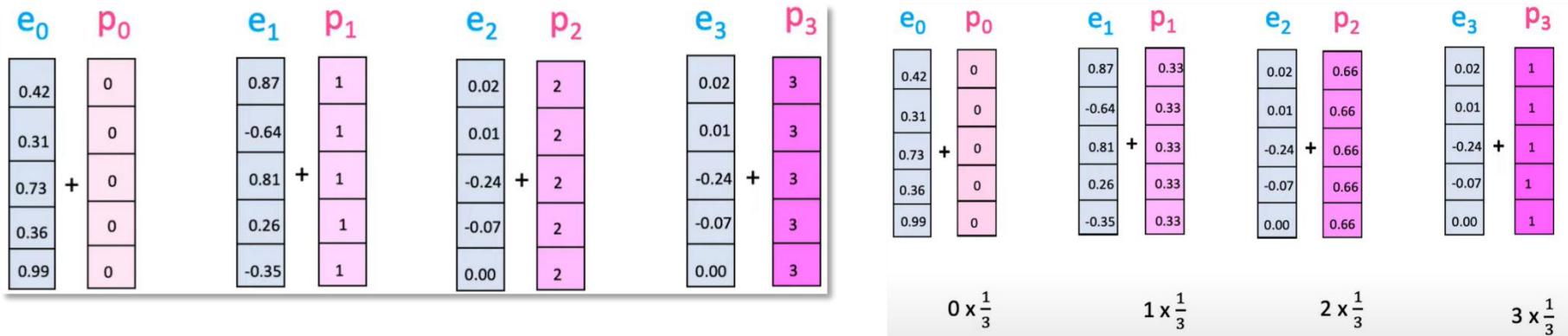


The positional encoding layer encodes the position information and generates position embeddings of inputs.

# Positional Encoding in Transformer

## Trivial Solutions

- Increase by word IDs: The words latter in the text have huge embeddings.
- Normalize by sentence length: The same word of the same position in sentences with different lengths does not have the same embedding



Ideally, the position embedding values

- (1) Should have **similar scale** in different position
- (2) Should **remain the same at a given position** irrespective of the text total length

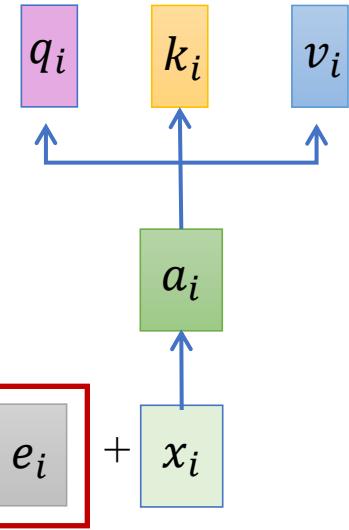
# Positional Encoding in Transformer

- To incorporate position information.
- Each position has a unique positional vector  $e_i$  (not learned from data)

$$a_i = x_i + e_i$$

$$e_{i,2j} = \sin(i/10000^{2j/d})$$

$$e_{i,2j+1} = \cos(i/10000^{2j/d})$$



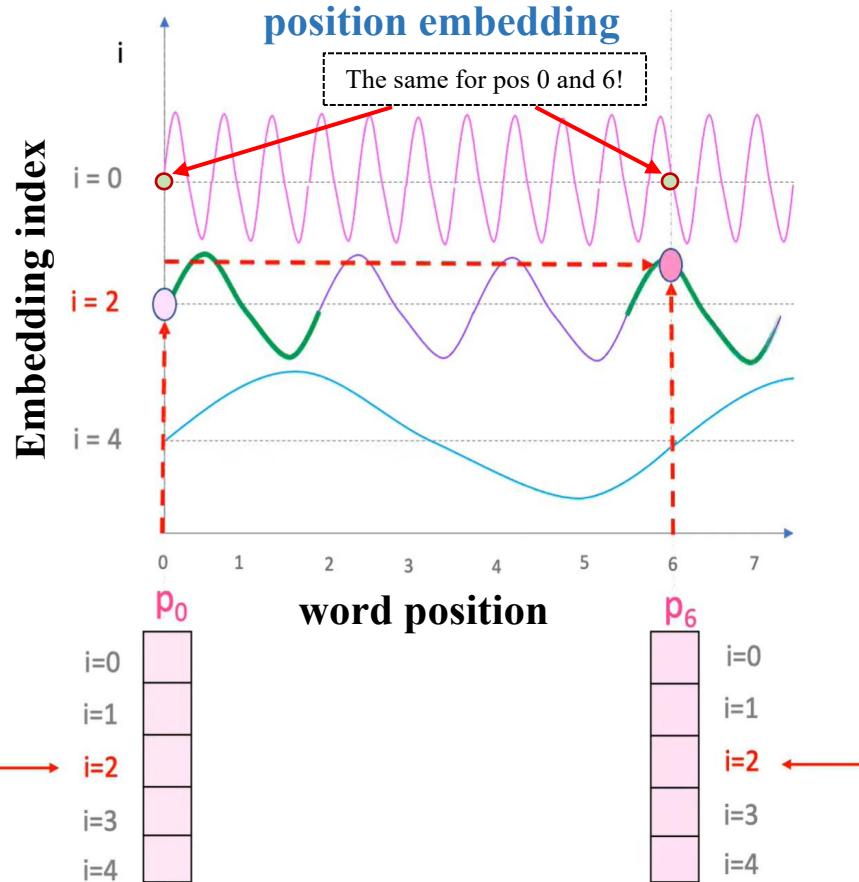
Contain two parts  
(1) Sin part for even position  
(2) Cos part for odd position

Where  $i$  is the position,  $j$  is the dimension.

The positional encodings  $e_i$  have the same dimension  $d$  as the word embedding  $x_i$ .

# Positional Encoding in Transformer

- A deep look of encoding function (sin part)



Position Embedding is a 2D matrix

- Axis  $i$ : embedding index ( $0, 1, \dots, 128$ )
- Axis  $pos$ : word position ( $0, 1, \dots, \infty$ )

Why two axes?

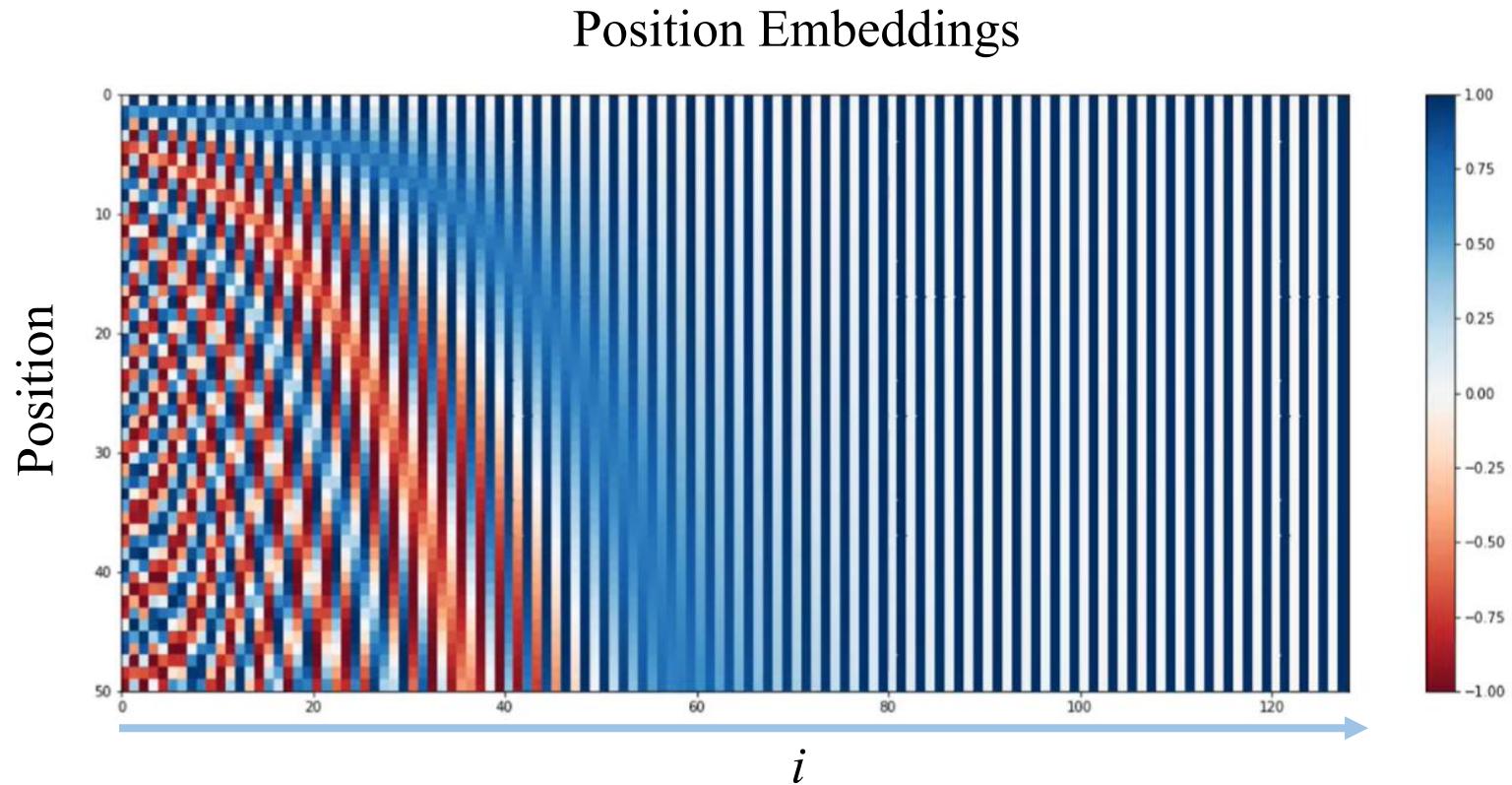
- (1) The axis  $pos$  is required to distinguish different word position.
- (2) Then, why we need the axis  $i$  ?

**Answer: Sin is periodic function!**

- At pos 0 and 6, position embedding values are the same at index 0.
- At pos 0 and 6, position embedding values are differ at index 2.

# Positional Encoding in Transformer

- Visualization of  $PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$



# Positional Encoding in Transformer

## ■ Implementation of positional encoding layer

```
class PositionalEncoding(nn.Module):
    def __init__(self, emb_size: int, dropout, maxlen: int = 5000):
        super(PositionalEncoding, self).__init__()
        den = torch.exp(- torch.arange(0, emb_size, 2) * math.log(10000) / emb_size)
        pos = torch.arange(0, maxlen).reshape(maxlen, 1)
        pos_embedding = torch.zeros((maxlen, emb_size))
        pos_embedding[:, 0::2] = torch.sin(pos * den)
        pos_embedding[:, 1::2] = torch.cos(pos * den)
        pos_embedding = pos_embedding.unsqueeze(-2)

        self.dropout = nn.Dropout(dropout)
        self.register_buffer('pos_embedding', pos_embedding)

    def forward(self, token_embedding: Tensor):
        return self.dropout(token_embedding +
                           self.pos_embedding[:token_embedding.size(0), :])
```

# Outline

1 Backgrounds

2 Recurrent Neural Networks (RNNs)

3 Attention Mechanisms

4 Attention is All You Need: Transformer

5 Long Short-Term Memory (LSTM) Networks

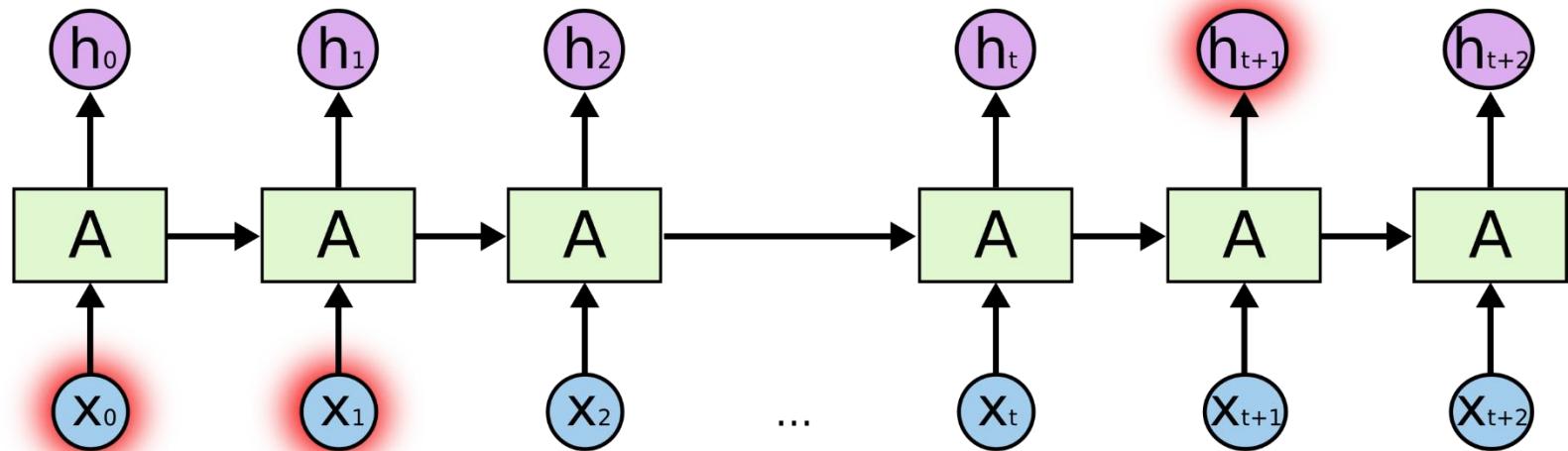
6 BERT(Bidirectional Encoder Representation from Transformers)

7 Conclusions

# Long-Term Dependency

- Long distance between previous context and present location.

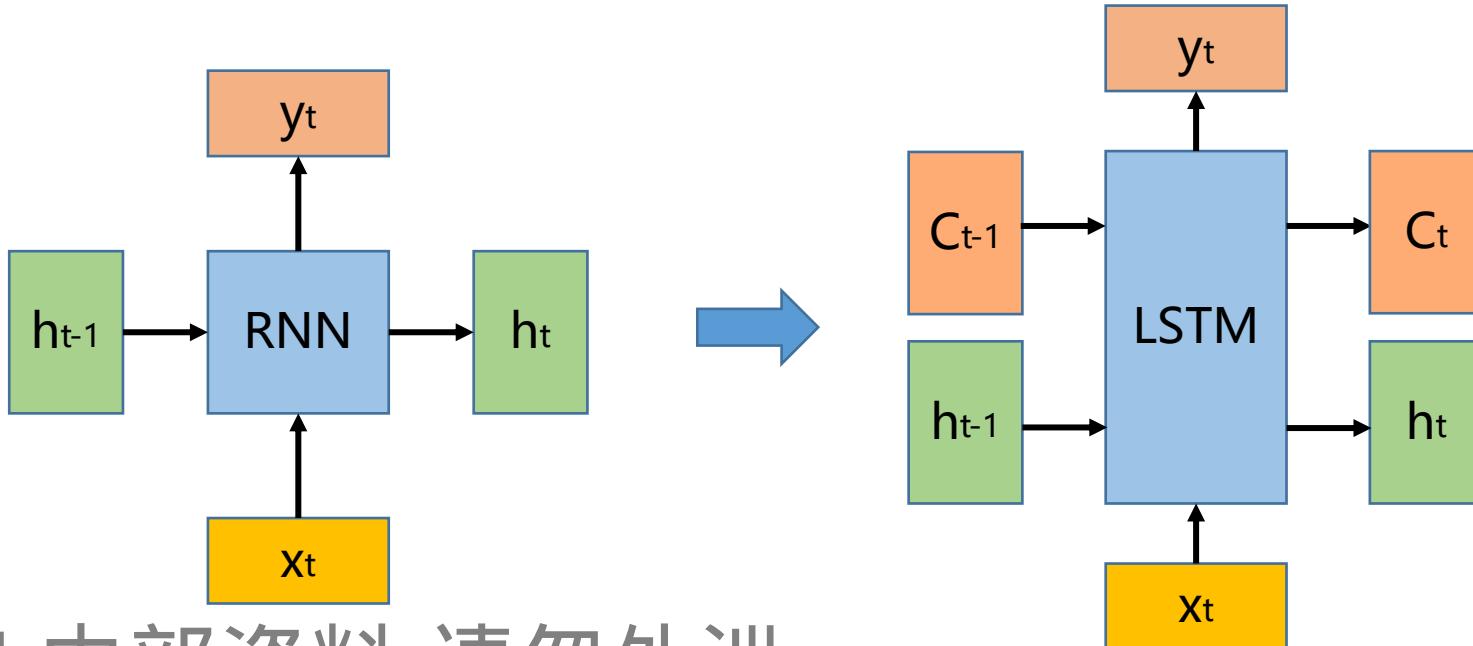
"I grew up in France... I speak fluent French"



- RNN is unable to connect to information too far ahead.

# Long Short-Term Memory Model

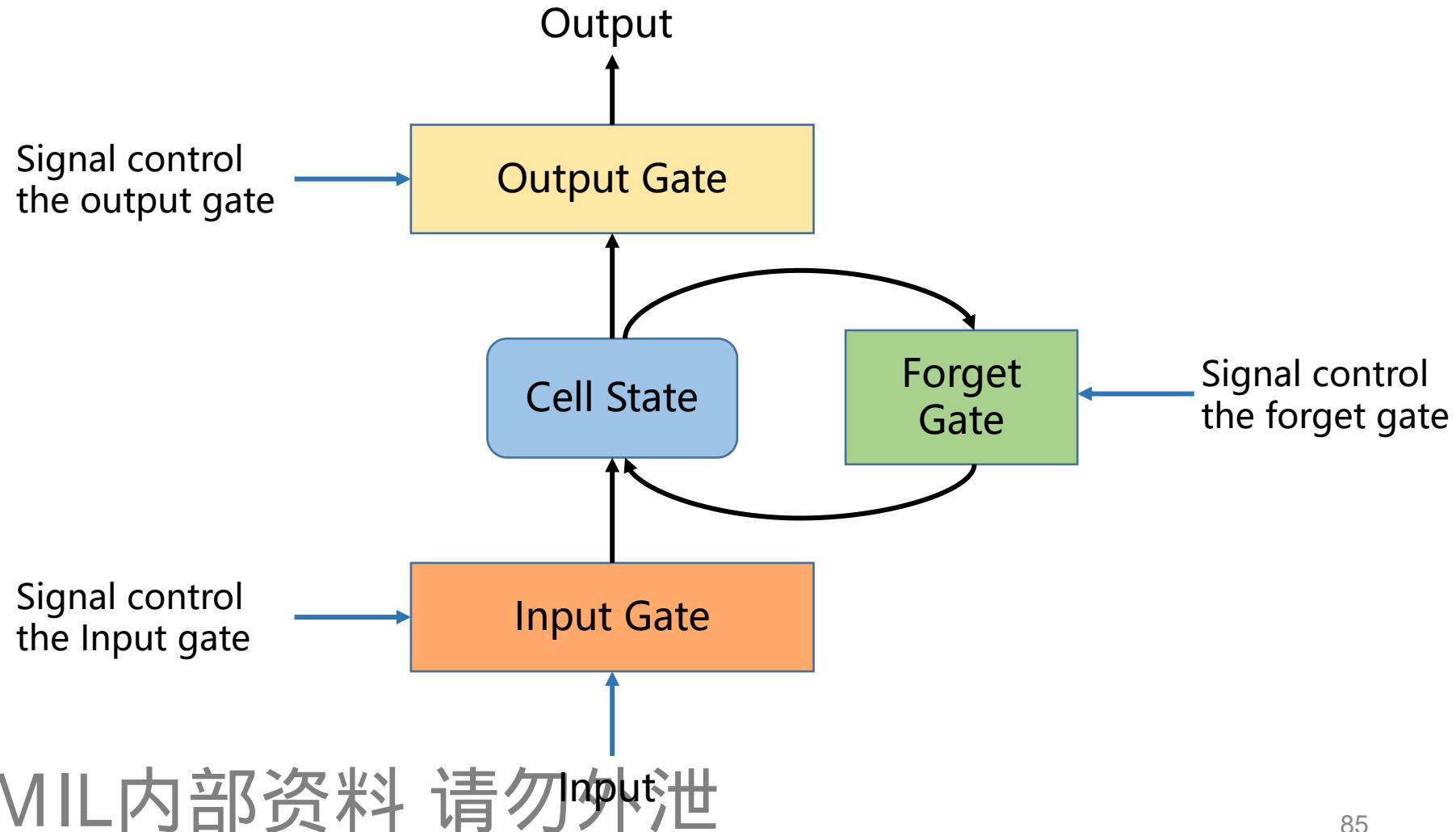
- Long Short-term Memory Models (LSTMs) make dependencies cross long distance easier to learn.
- LSTM has two states  $h_t$ (hidden state) and  $C_t$ (cell state) , while RNN only has one state  $h_t$ .



SMIL内部资料 请勿外泄

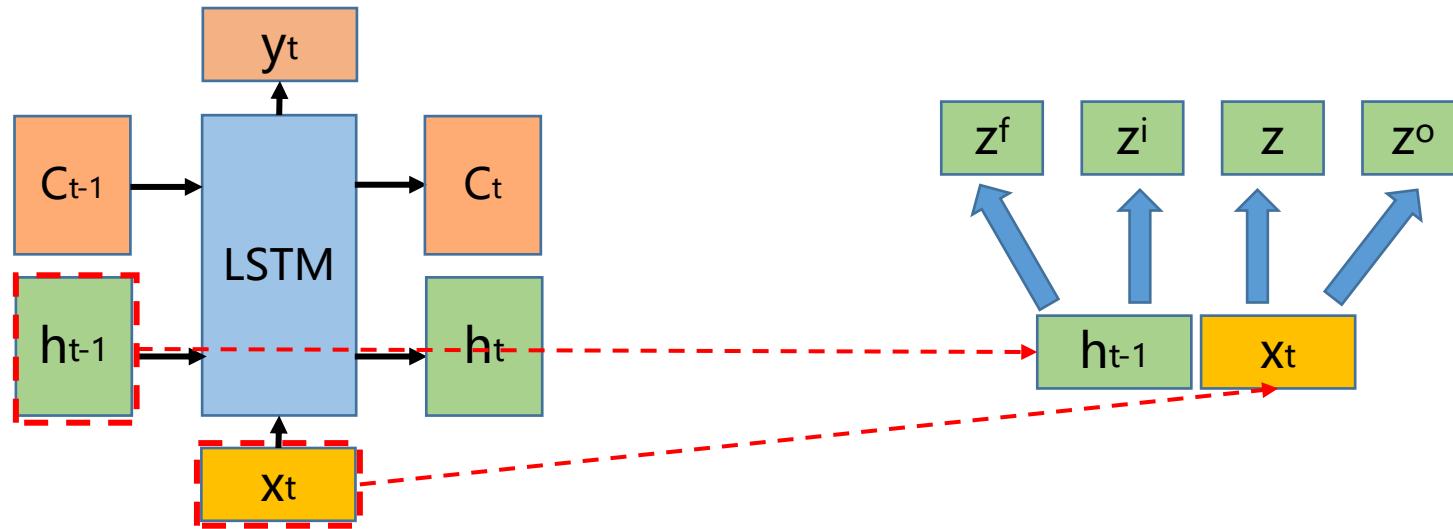
# Long Short-Term Memory Model

## ■ Basic structure of LSTMs



# Long Short-Term Memory Model

## ■ Three gates of LSTM



$$\text{Input gate: } z^i = \sigma(W_i [h_{t-1}, x_t] + b_i)$$

$$\text{Forget gate: } z^f = \sigma(W_f [h_{t-1}, x_t] + b_f)$$

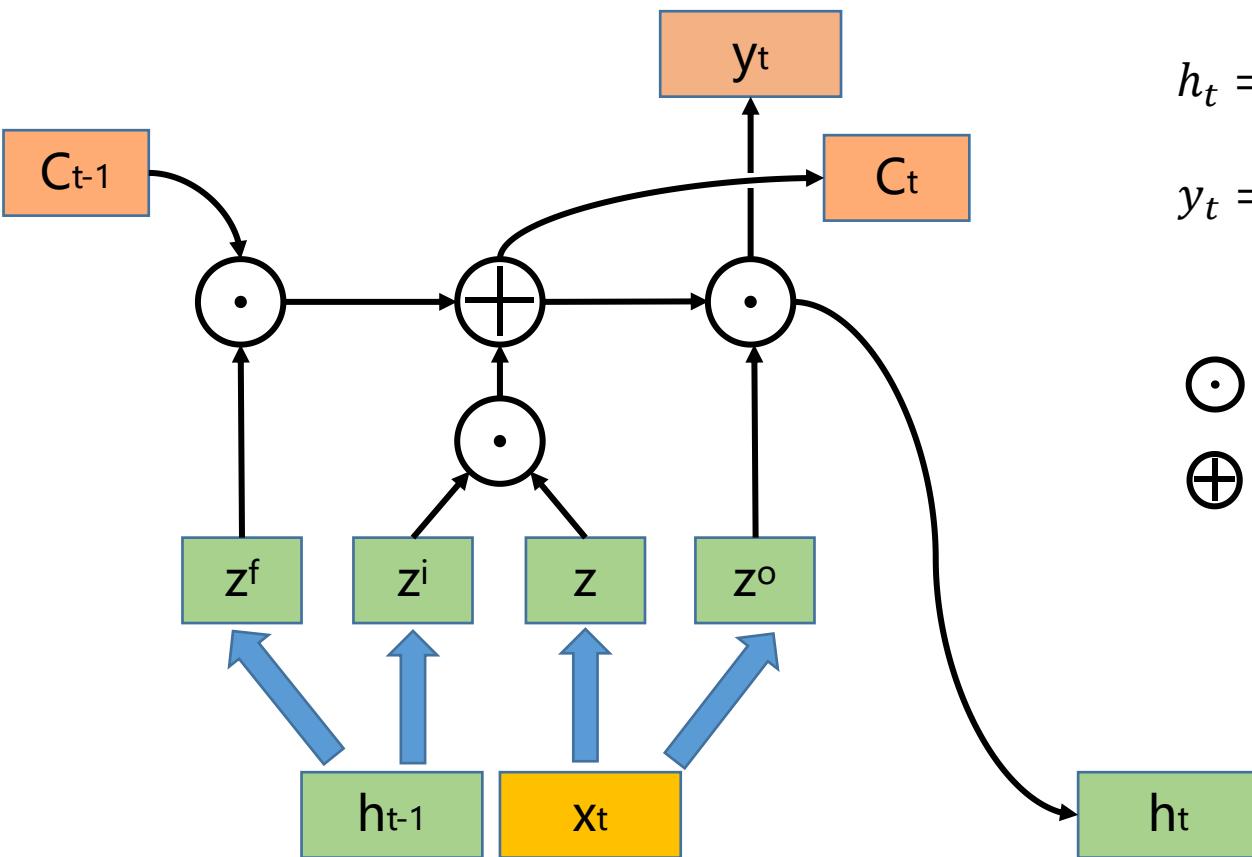
$$\text{Output gate: } z^o = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$\text{Input activation: } z = \tanh(W[h_{t-1}, x_t] + b)$$

- Apply sigmoid function  $\sigma$  and  $z^i, z^f, z^o \in (0,1)$ , simulate the gates' open and close.
- Apply tanh activation function to get  $z$ .

# Long Short-Term Memory Model

## ■ How the gates work



$$C_t = z^f \odot C_{t-1} + z^i \odot z$$

$$h_t = z^o \odot \tanh(C_t)$$

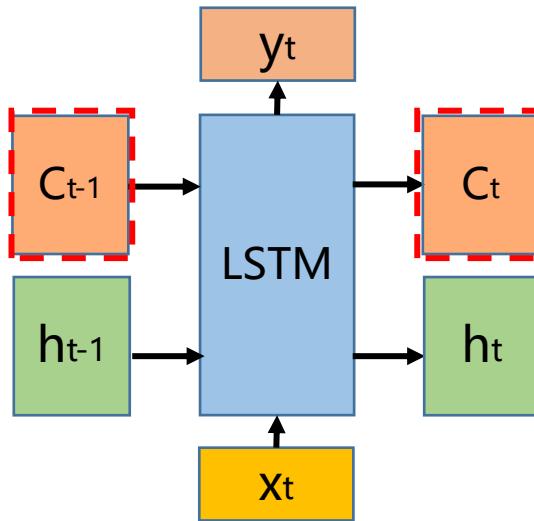
$$y_t = \sigma(W' h_t)$$

○ Hadamard Product

⊕ Matrix Addition

# Long Short-Term Memory Model

- How LSTMs handle long-term dependencies



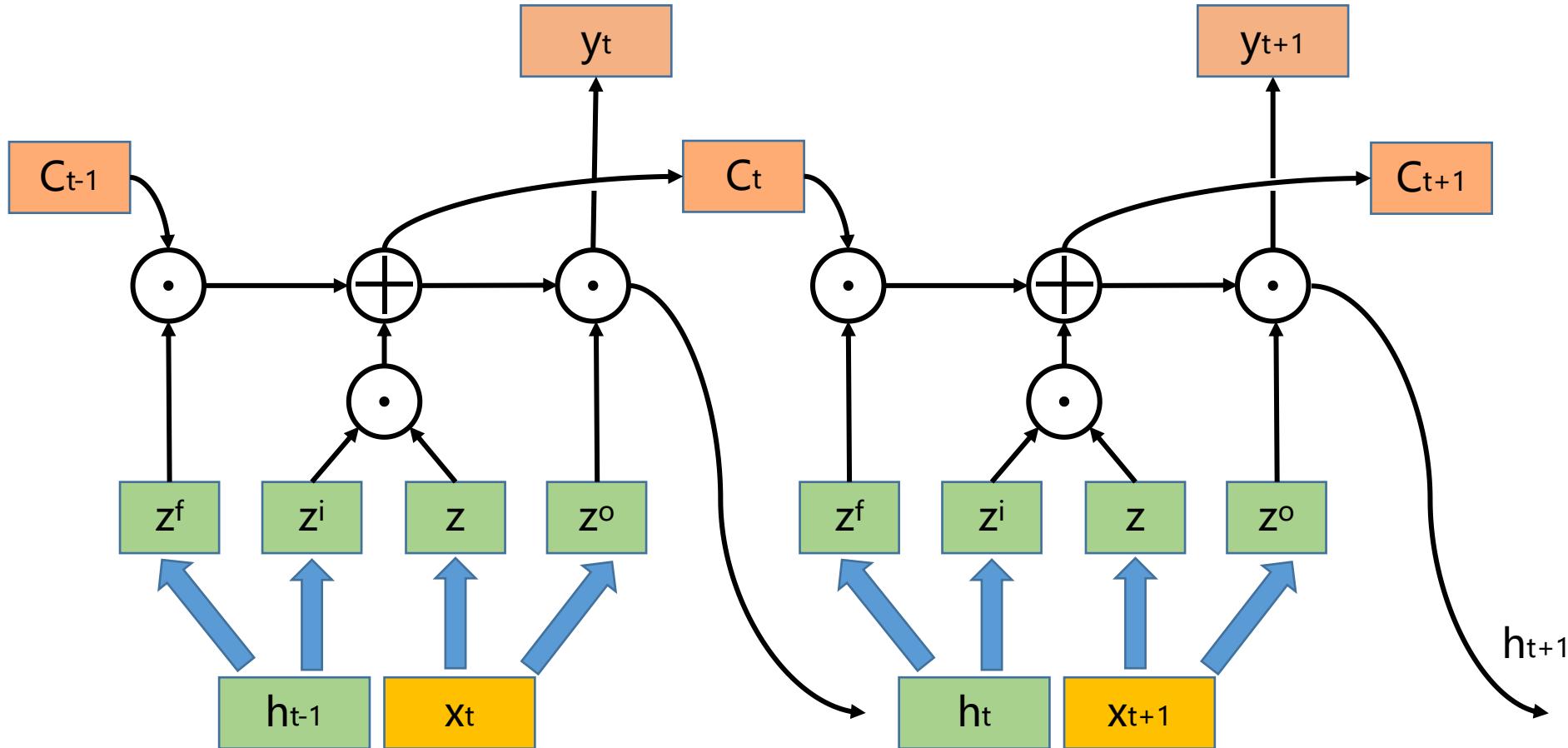
$$C_t = z^f \odot C_{t-1} + z^i \odot z$$

With the forget gate  $z^f$  and input gate  $z^i$ , cell state  $C_t$  can selectively store the information.

- By introducing gates and cell state  $C_t$ , LSTMs are able to keep previous information for long time and deal with long-term dependencies.

# Long Short-Term Memory Model

## ■ Forward Propagation Through Time (BPTT)



SMIL内部资料 请勿外泄

# Forward Propagation: Prediction for Networks

- Given input a sequence of words  $x_1, \dots, x_T$ , the forward propagation of a LSTM cell is to predict the output  $y_1$  of  $x_1$ .

- Detailed prediction process

for  $t$  from 1 to  $T$  do

$$\mathbf{z}^i = \sigma(W_i [h_{t-1}, x_t] + b_i),$$

$$\mathbf{z}^f = \sigma(W_f [h_{t-1}, x_t] + b_f),$$

$$\mathbf{z}^o = \sigma(W_o [h_{t-1}, x_t] + b_o),$$

$$\mathbf{z} = \tanh(W [h_{t-1}, x_t] + b)$$

$$C_t = z^f \odot C_{t-1} + z^i \odot \mathbf{z}$$

$$h_t = z^o \odot \tanh(C_t)$$

$$y_t = \sigma(W' h_t)$$

end for



$\sigma$ : sigmoid function  
 $\tanh$ : hyperbolic tangent function

$\hat{y}_t$  : ground-truth label

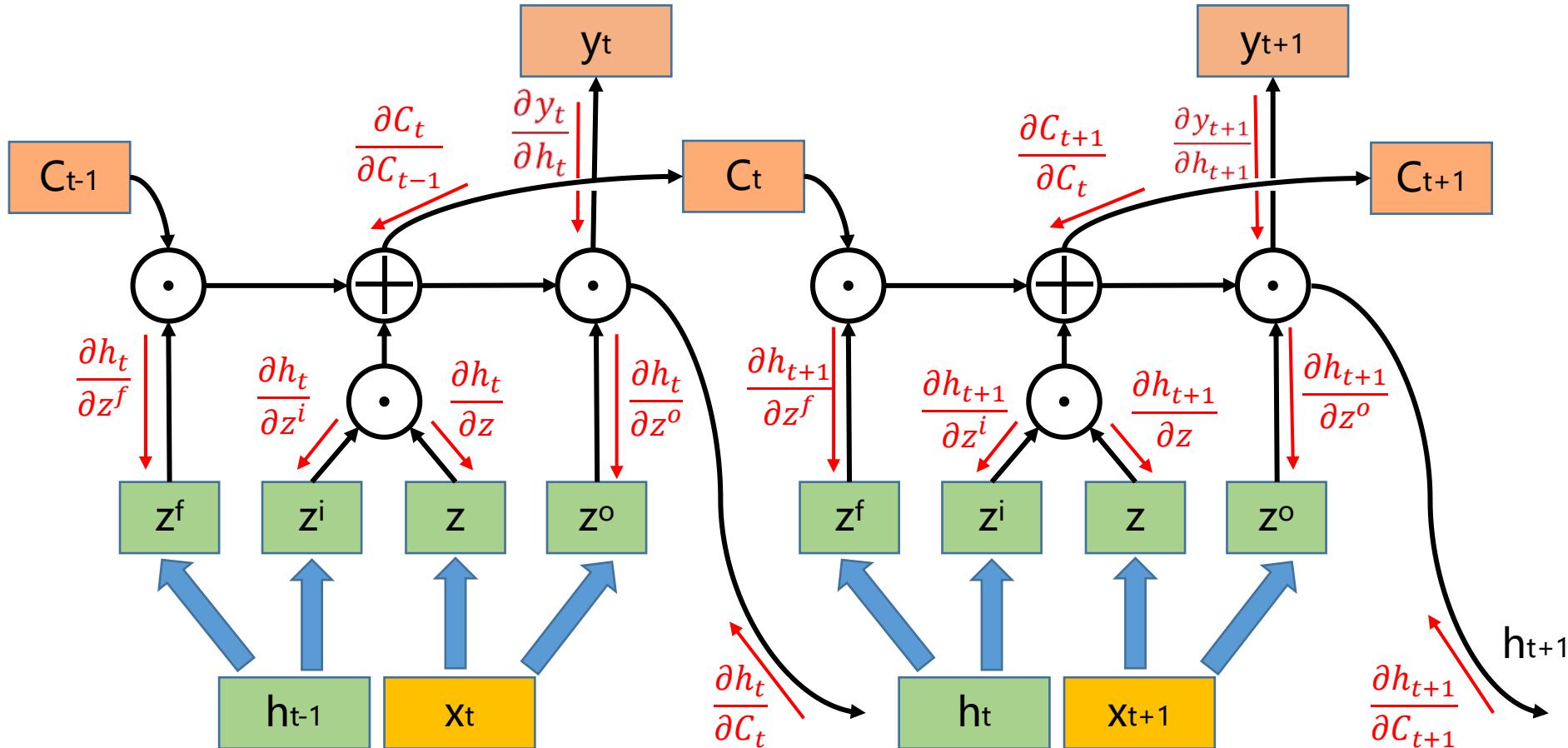
$y_t$  : predicted label

■ Total loss:  $L(y, \hat{y}) = \sum_{t=1}^T L(y_t, \hat{y}_t)$

SMIL 内部资料 请勿外泄

# Long Short-Term Memory Model

## ■ Backward Propagation Through Time (BPTT)



SMIL内部资料 请勿外泄

# Backpropagation of Neural Networks

- Let us resort to Chain Rule to compute the gradient.

$$\frac{\partial L_3(y, \hat{y})}{\partial W'} = \frac{\partial L_3(y, \hat{y})}{\partial y_3} \frac{\partial y_3}{\partial W'}$$

Use  $L_3(y, \hat{y})$  as an example

$$\frac{\partial L_3(y, \hat{y})}{\partial W} = \sum_{k=0}^3 \frac{\partial L_3(y, \hat{y})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial h_k} \frac{\partial h_k}{\partial C_k} \frac{\partial C_k}{\partial z} \frac{\partial z}{\partial W}$$

$$\frac{\partial L_3(y, \hat{y})}{\partial W_i} = \sum_{k=0}^3 \frac{\partial L_3(y, \hat{y})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial h_k} \frac{\partial h_k}{\partial C_k} \frac{\partial C_k}{\partial z^i} \frac{\partial z^i}{\partial W_i}$$

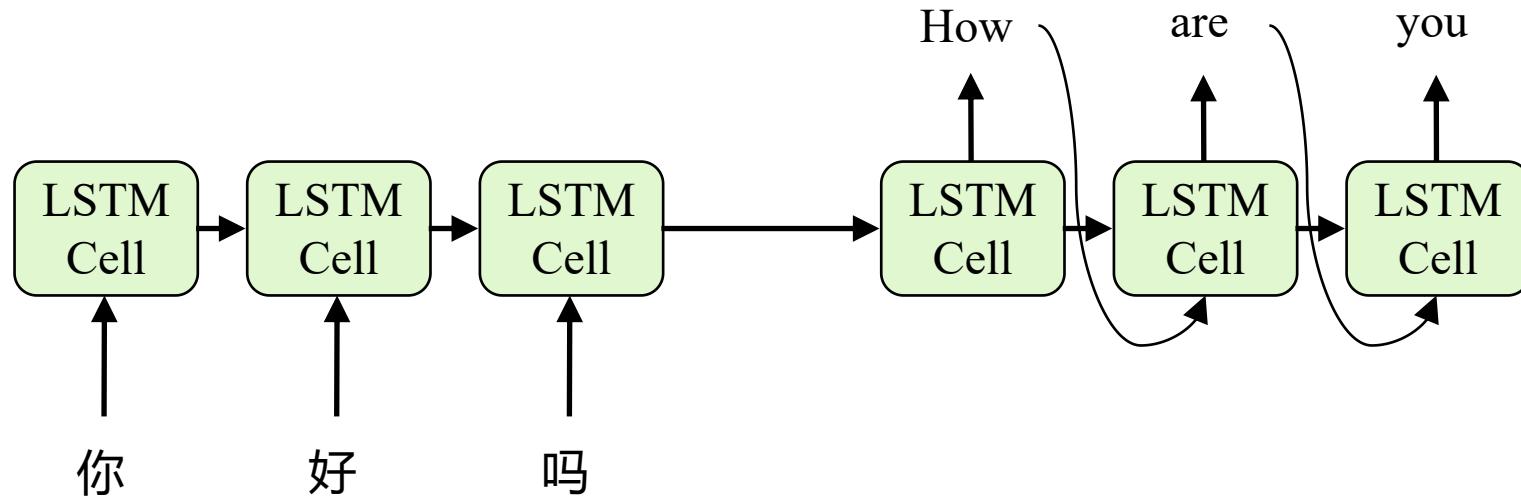
$$\frac{\partial L_3(y, \hat{y})}{\partial W_f} = \sum_{k=0}^3 \frac{\partial L_3(y, \hat{y})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial h_k} \frac{\partial h_k}{\partial C_k} \frac{\partial C_k}{\partial z^f} \frac{\partial z^f}{\partial W_f}$$

$$\frac{\partial L_3(y, \hat{y})}{\partial W_o} = \sum_{k=0}^3 \frac{\partial L_3(y, \hat{y})}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial h_k} \frac{\partial h_k}{\partial C_k} \frac{\partial C_k}{\partial z^o} \frac{\partial z^o}{\partial W_o}$$

# Application

What can LSTM do?

Predict a sequence:



- Neural machine translation.

- Encoder-decoder structure.

SMIL内部资料 请勿外泄

# Outline

1 Backgrounds

2 Recurrent Neural Networks (RNNs)

3 Attention Mechanisms

4 Attention is All You Need: Transformer

5 Long Short-Term Memory (LSTM) Networks

6 BERT(Bidirectional Encoder Representation from Transformers)

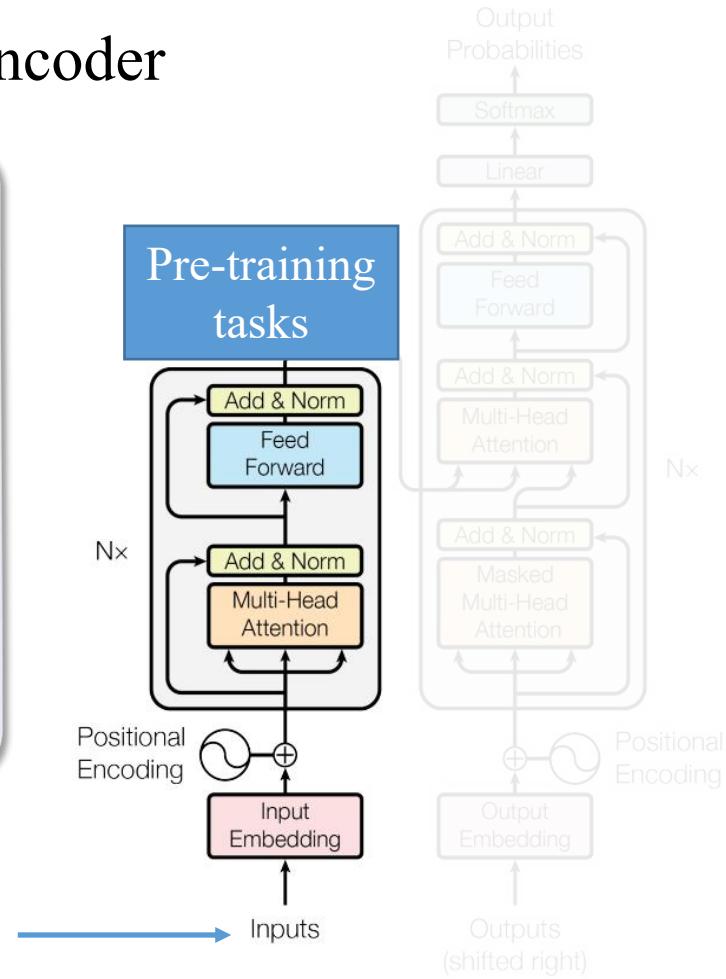
7 Conclusions

BERT = Pre-trained Transformer Encoder

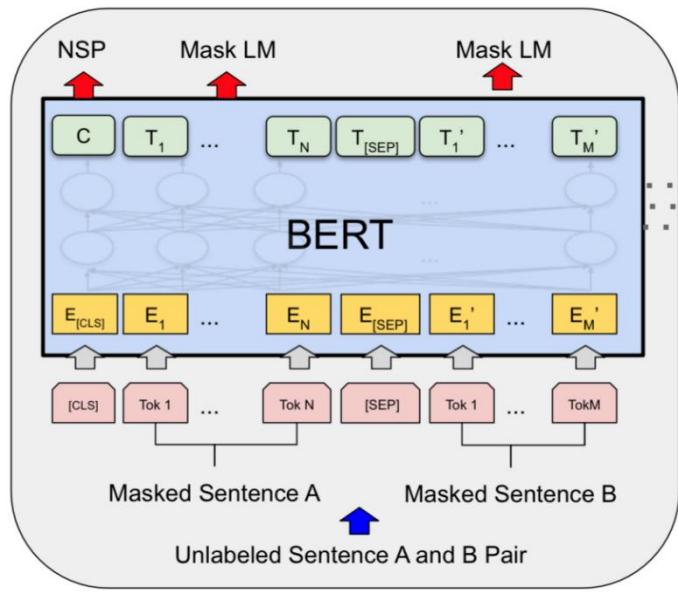
The same as Transformers:

- Multi-head self attention
- Feed-forward layers
- Layer normalization and residual connections
- Input Embedding

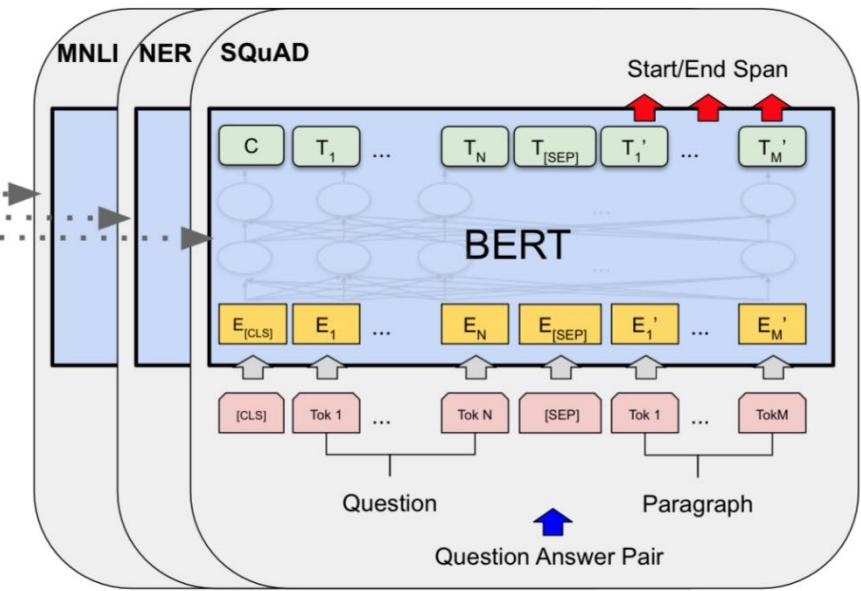
large amount of text without annotation



Pre-training, then fine-tuning on various downstream tasks.



Pre-training



Fine-Tuning

# Task #1: Masked Language Model (cloze task)

How to best utilize the **bi-directional** architecture of Transformers using **pre-training**?

- Mask out k% of the input words, then predict the masked words.
- The prediction is conditioned on context from both direction.

BERT always use k=15%

The dog [MASK] across the street because it was [MASK] tired.

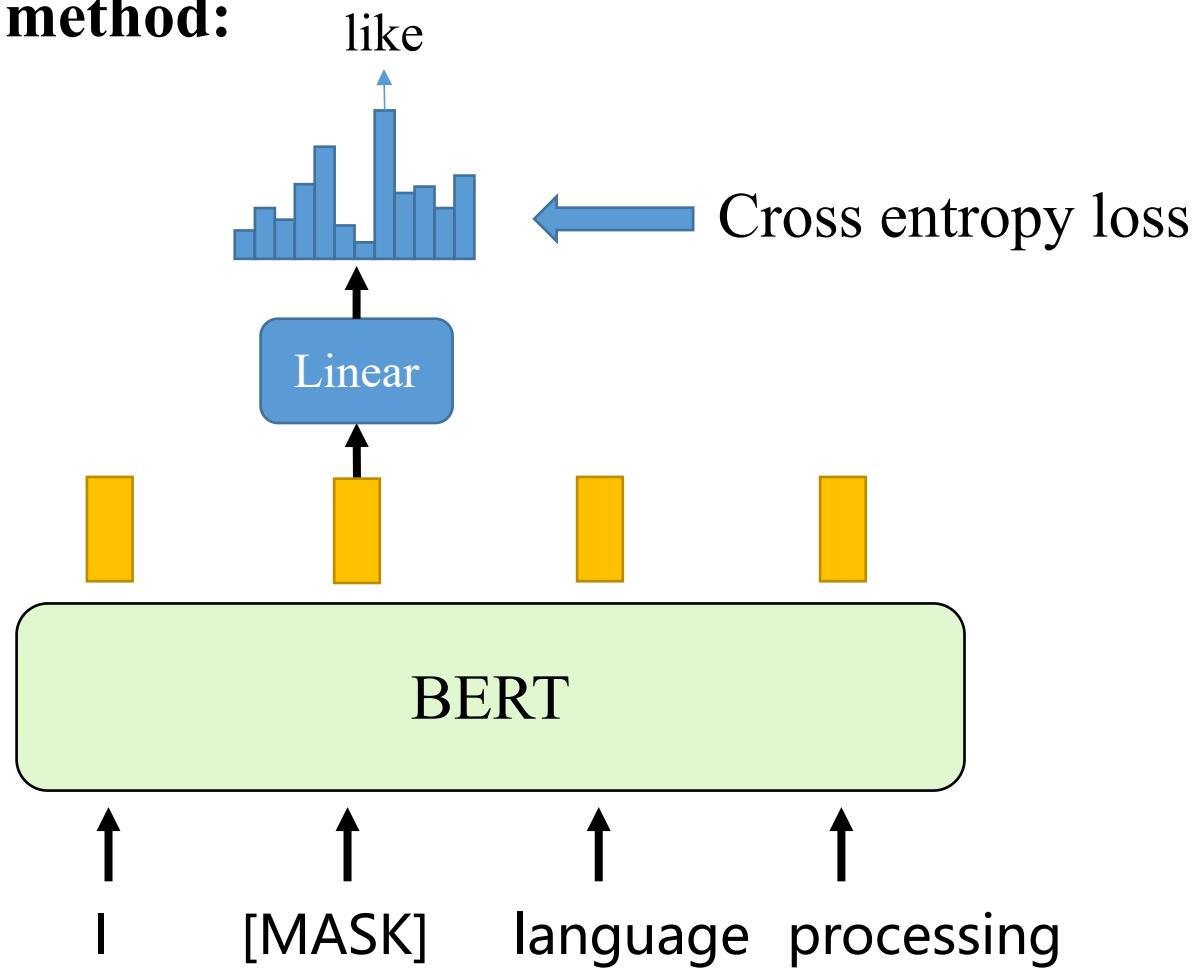


Too little masking: too expensive to train

Too much masking: not enough context

# Task #1: Masked Language Model (cloze task)

**Training method:**



# Task #1: Masked Language Model (cloze task)

- **Problem:** [MASK] does not appear during fine-tuning.
- **Solution:** The training data generator randomly chooses 15% of the token positions for prediction.

1. 80% of the time, replace with [MASK]

my dog is cute → my dog is [MASK]

2. 10% of the time, replace random word

my dog is cute → my dog is apple

3. 10% of the time, keep same

my dog is cute → my dog is cute

## Task #2: Next Sentence Prediction

**Motivation:** understanding the **relationship** between two sentences is important to some tasks, e.g.:

- Question Answering
- Natural Language Inference

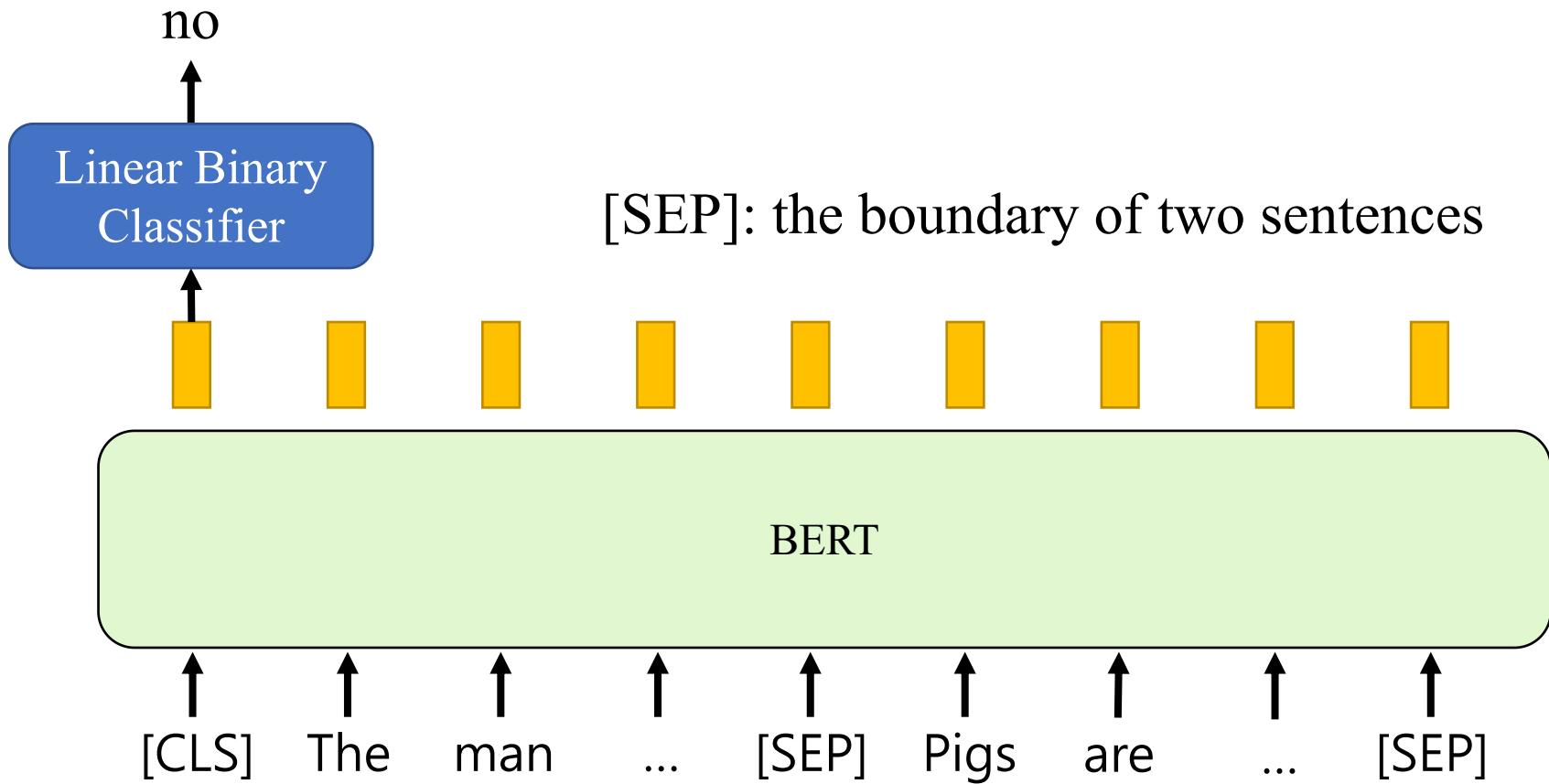
Choose two sentence A, B:

- 50% of the time, B is the actual next sentence that follows A.
- 50% of the time, B is a random sentence from the corpus.
- Let neural networks predict which is the case.

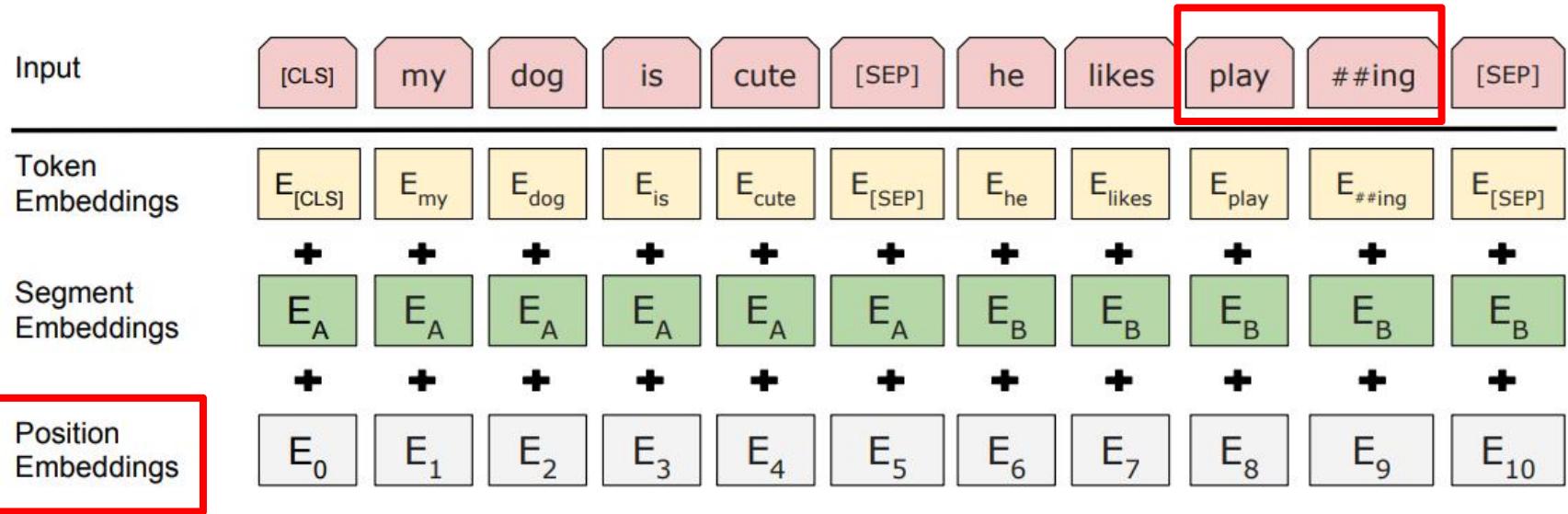
Sentence A = The man went to the store.      Sentence A = The man went to the store.  
Sentence B = He bought a glass of milk.      Sentence B = Pigs are flightless.  
Label = IsNextSentence                          Label = NotNextSentence

## Task #2: Next Sentence Prediction

Training method:

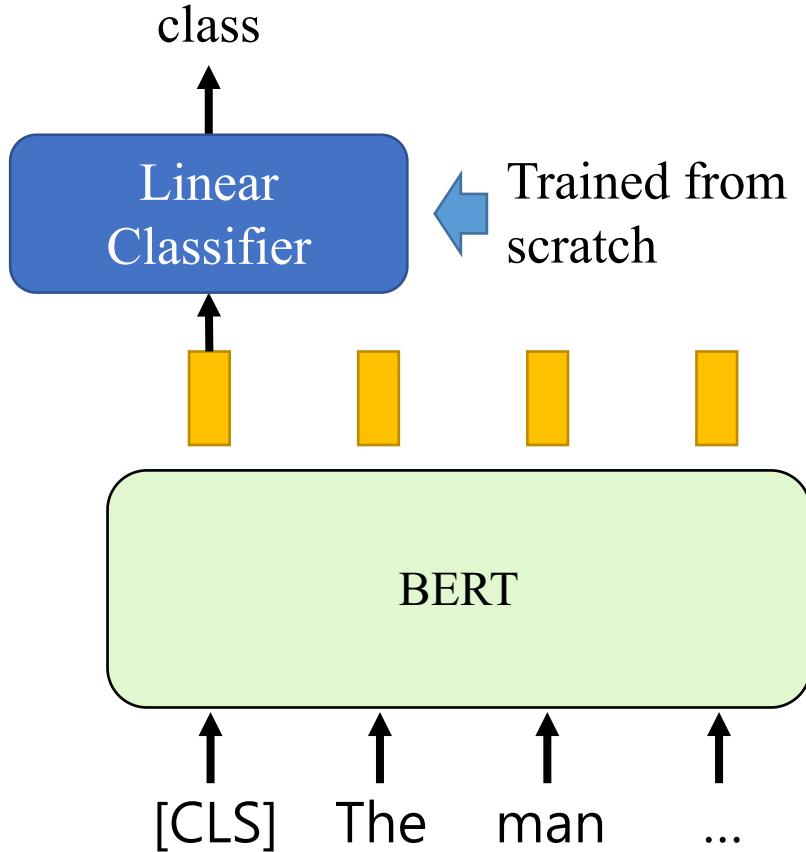


# BERT



- Use WordPiece embeddings with a 30k token vocabulary.
- Segment embeddings are used to differentiate sentences.
- Position embeddings are different from Transformer.

# Down-stream task #1: Sentence classification

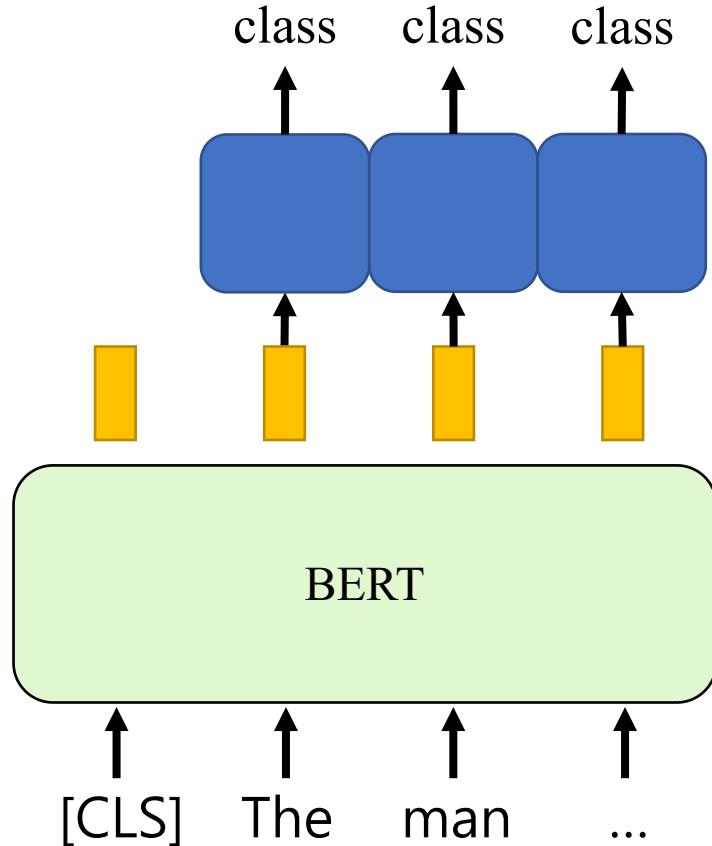


**Input:** A sentence or a pair of sentence  
**Output:** class

**Example:** MNLI (Multi-Genre Natural Language Inference)

- **Fiction:** The Old One always comforted Jack, except today.
- **Hypothesis:** Jack knew the Old One very well.
- **Label:** neutral (neutral / contradiction / entailment)

# Down-stream task #2: Sentence tagging



**Input:** A sentence

**Output:** class of each token

**Example:** NER (Named Entity Recognition)

U.N. official Ekeus heads for Baghdad .  
ORG - PER - - LOC -

# Outline

1 Backgrounds

2 Recurrent Neural Networks (RNNs)

3 Attention Mechanisms

4 Attention is All You Need: Transformer

5 Long Short-Term Memory (LSTM) Networks

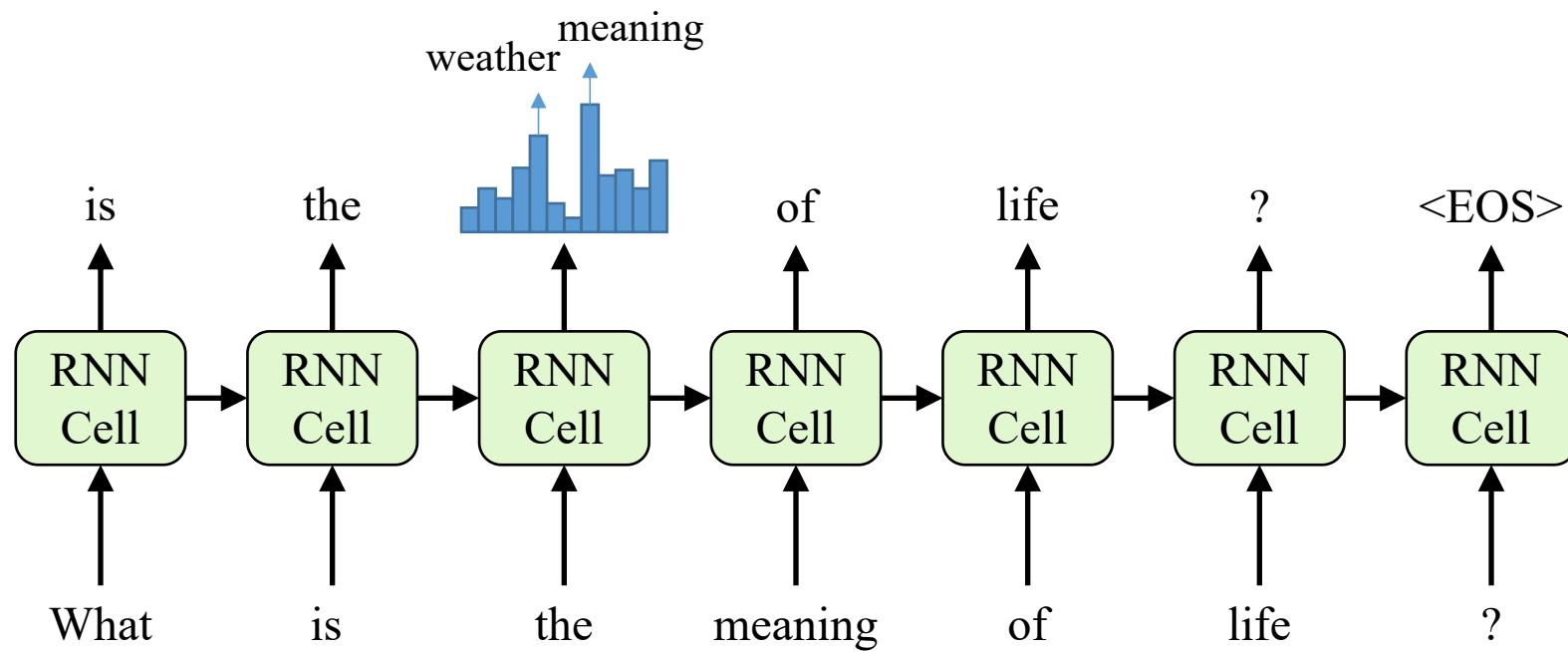
6 BERT(Bidirectional Encoder Representation from Transformers)

7 Conclusions

# Vanilla RNN

## Vanilla Recurrent Neural Network:

- Introduce **hidden state**, which remembers information about a sequence.
- Use **variable-length** sequences as input/output.
- Example application: Language Model (LM)

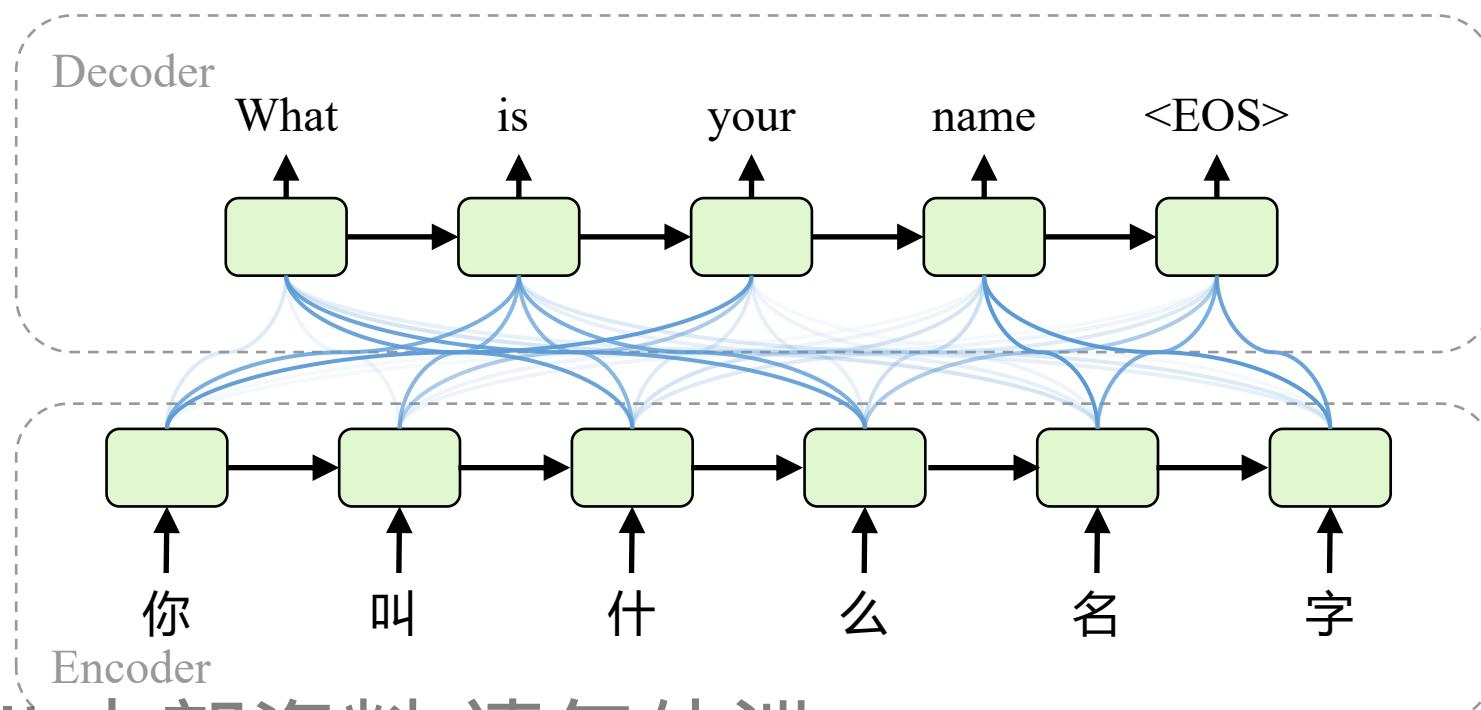


SMIL内部资料 请勿外泄

# Attention

## Attention in encoder-decoder

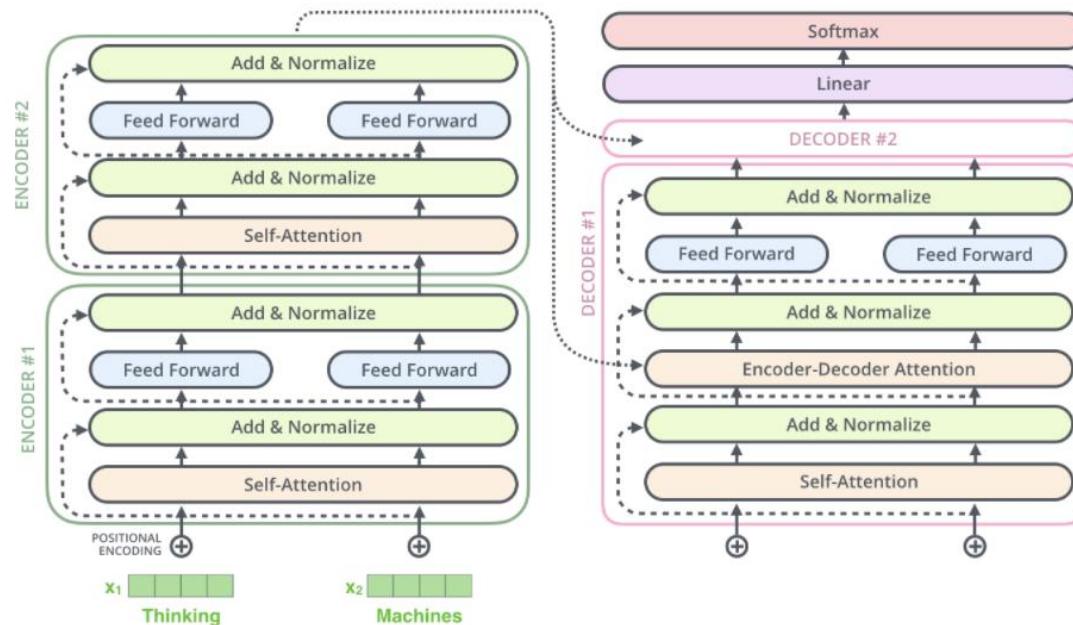
- Imitate human attention mechanism.
- Avoid encoding all information in a fixed-length vector.



# Transformer

## Attention is all you need: Transformer

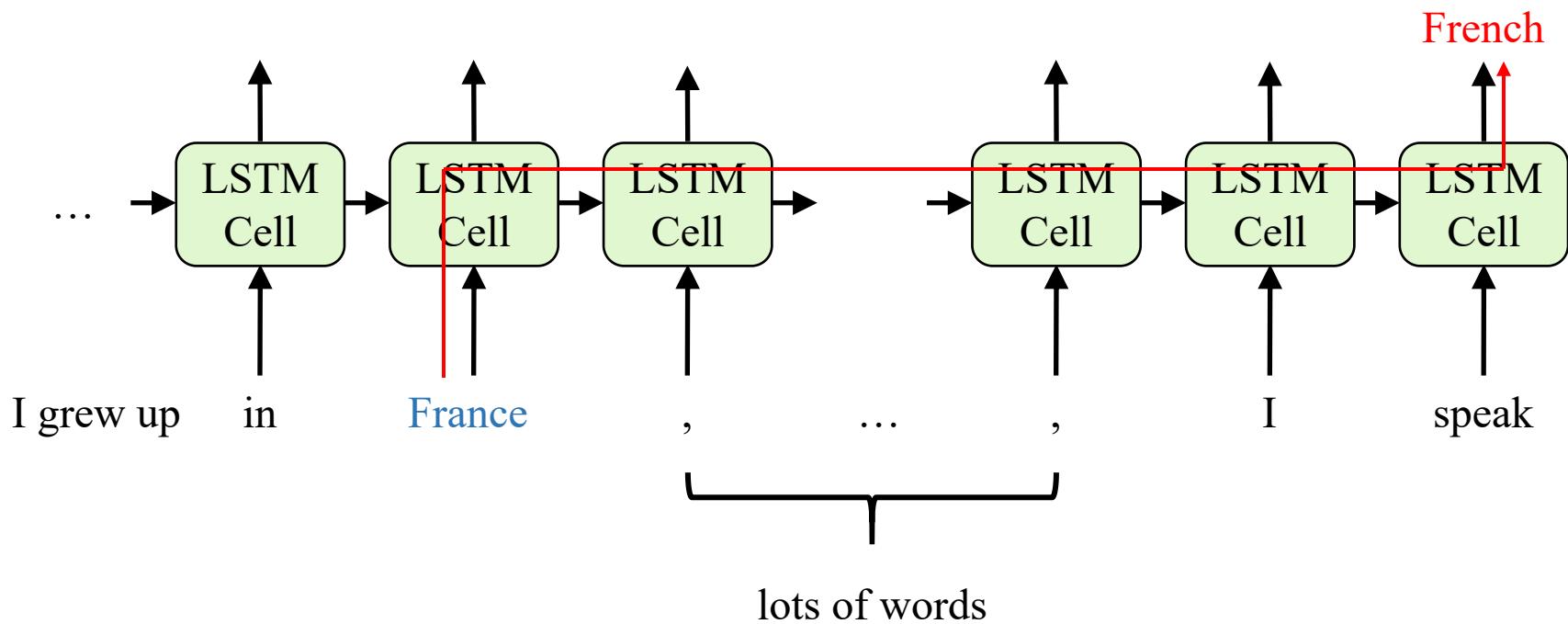
- A new simple network architecture based solely on attention mechanisms.
- The **self-attention** solves the problem of **long distance dependence**.



SMIL内部资料 请勿外泄

## Long Short-Term Memory

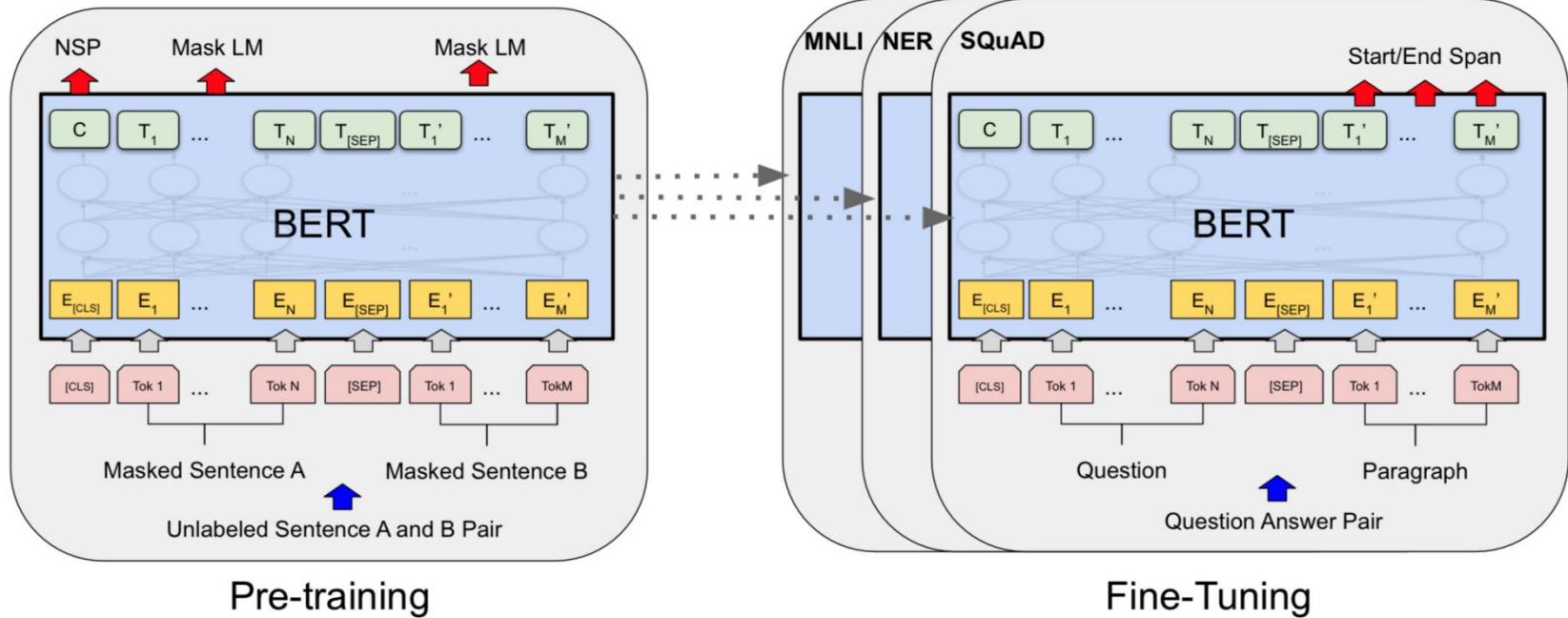
- Take good care of the long-term dependencies.



# BERT

## Bidirectional Encoder Representation from Transformers (BERT):

- Pre-trained Transformer encoders
- One pre-trained model, transferred to 11 tasks, got state-of-the-art result.



SMIL内部资料 请勿外泄