

人工智能之NLP

Seq2Seq

主讲人：GerryLiu

课程要求

- 课上课下“九字”真言
 - 认真听，善摘录，勤思考
 - 多温故，乐实践，再发散
- 四不原则
 - 不懒散惰性，不迟到早退
 - 不请假旷课，不拖延作业
- 一点注意事项
 - 违反“四不原则”，不推荐就业

课程内容

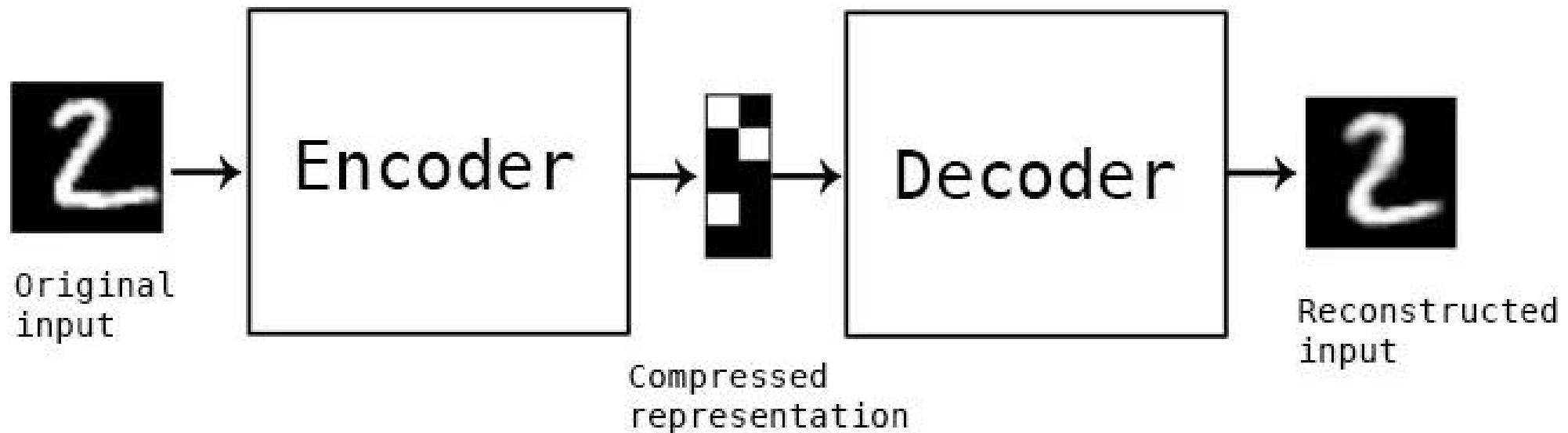
- 自编码神经网络回顾
- RNN、LSTM神经网络回顾
- Seq2Seq网络结构讲解
- Attention结构讲解
- Seq2Seq+Attention项目

自编码器回顾

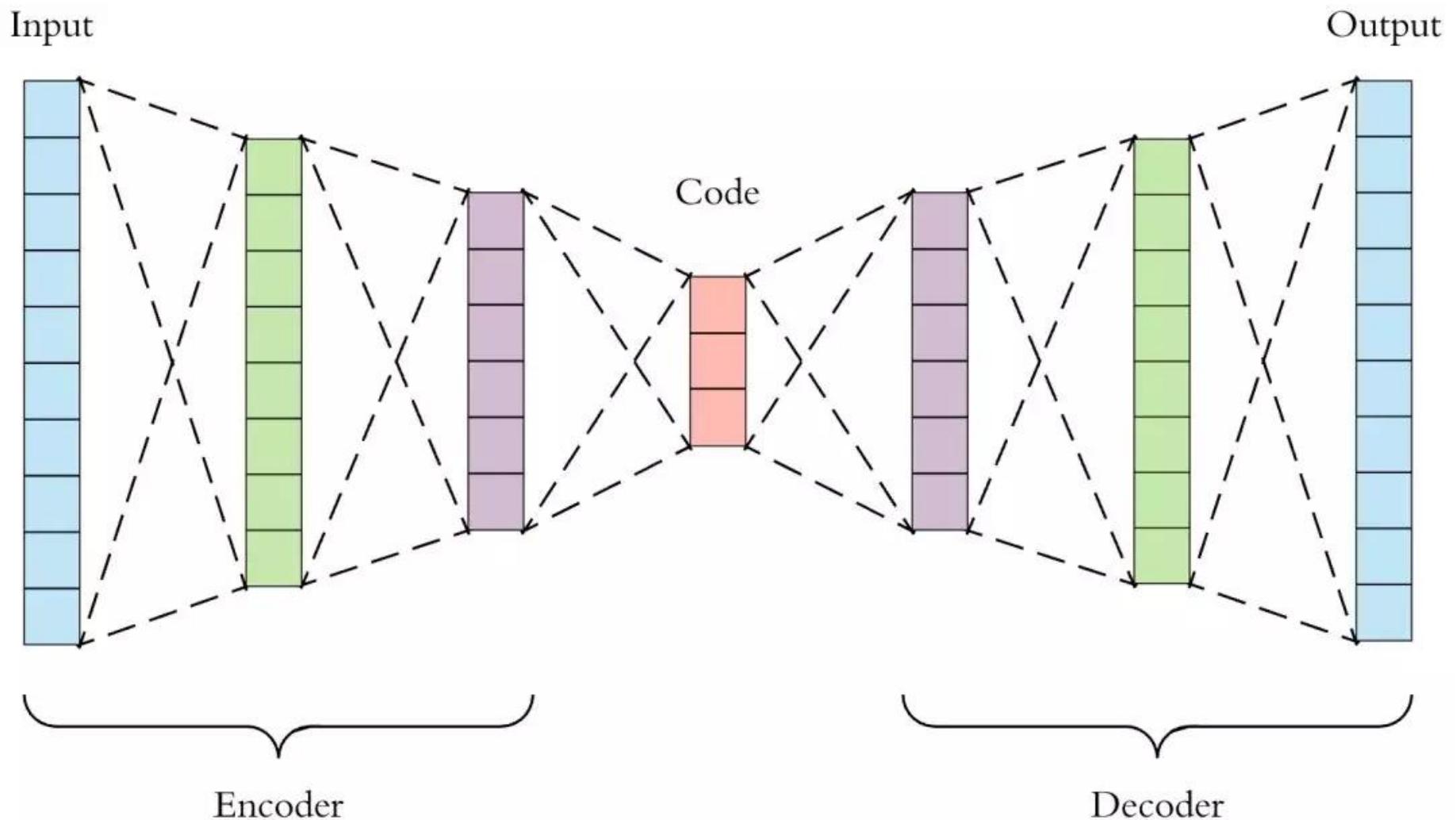
- 自编码器作为一种深度学习领域无监督的算法，本质上是一种数据压缩算法，和生成对抗网络一样，属于生成算法的一种。
- 自编码器(AutoEncoder, AE)就是一种利用反向传播使得输出值等于输入值的神经网络，它将输入压缩成潜在特征/高阶特征，然后将这种表征重构输出。主要包含以下三个特征：
 - 数据相关性。
 - 数据有损性。
 - 自动学习性。

自编码器回顾

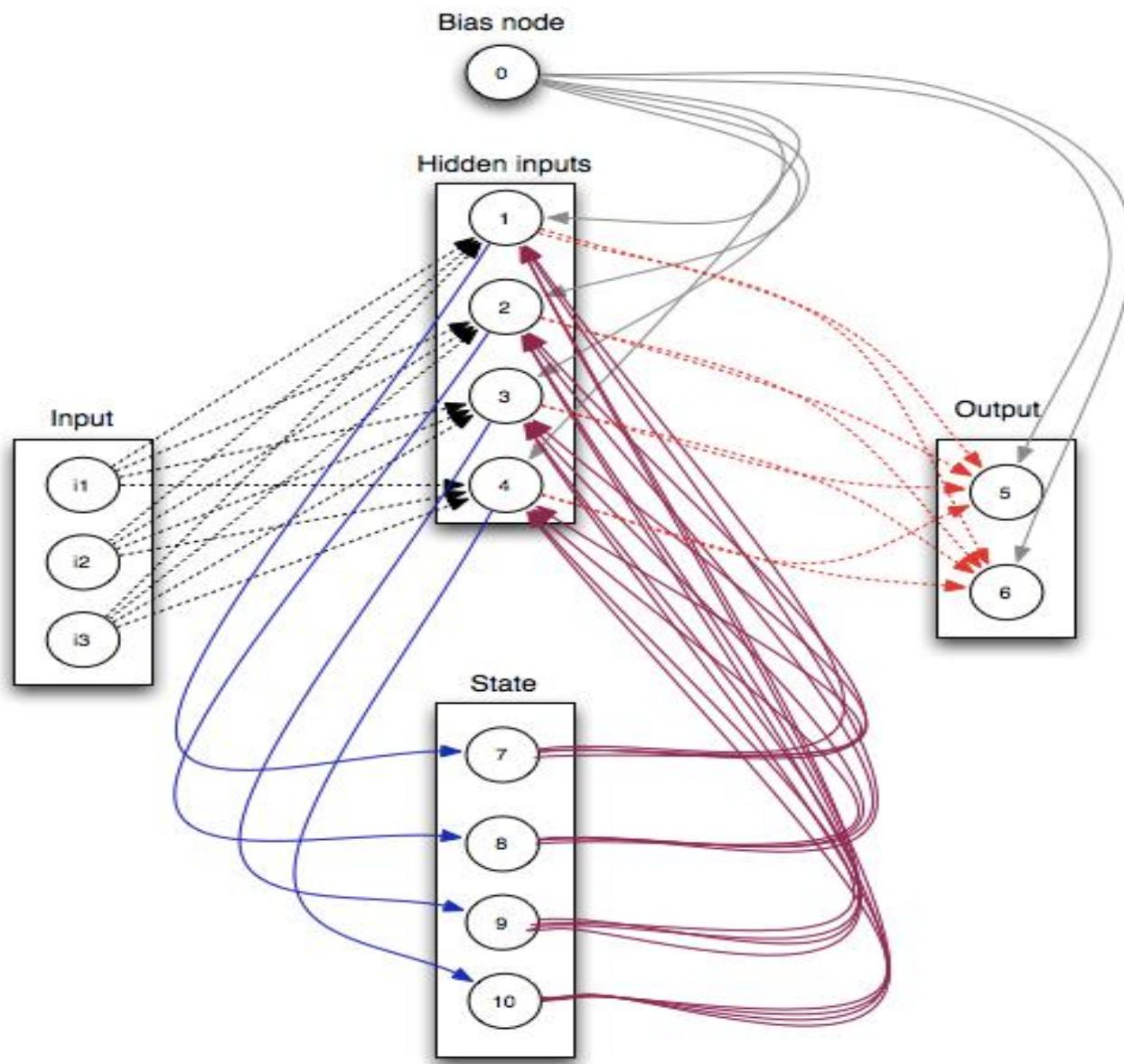
- 构建一个自编码器主要包括两部分：编码器(Encoder)和解码器(Decoder)。编码器将输入压缩为潜在空间特征，解码器将潜在空间特征重构输出。
- 自编码的核心价值是在于提取潜在的高阶空间特征信息。主要应用是两个方面：数据去燥以及进行可视化降维。



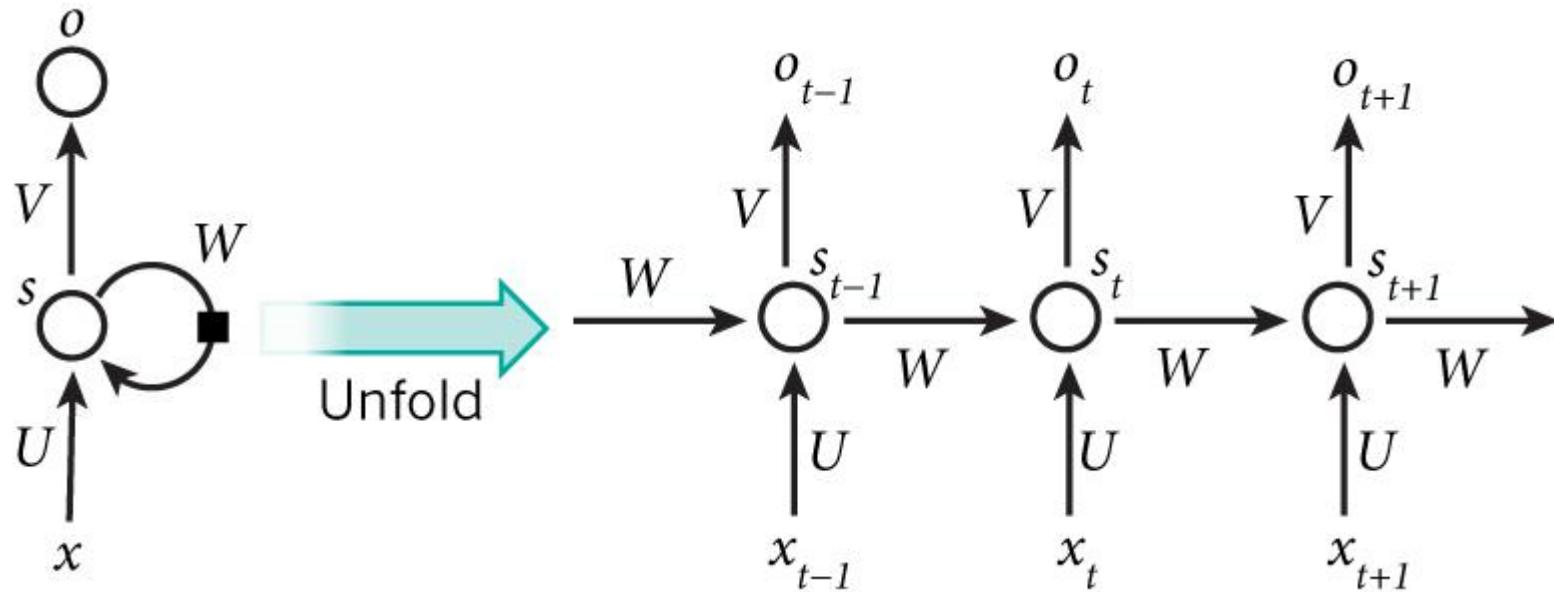
自编码器回顾



RNN回顾



RNN回顾

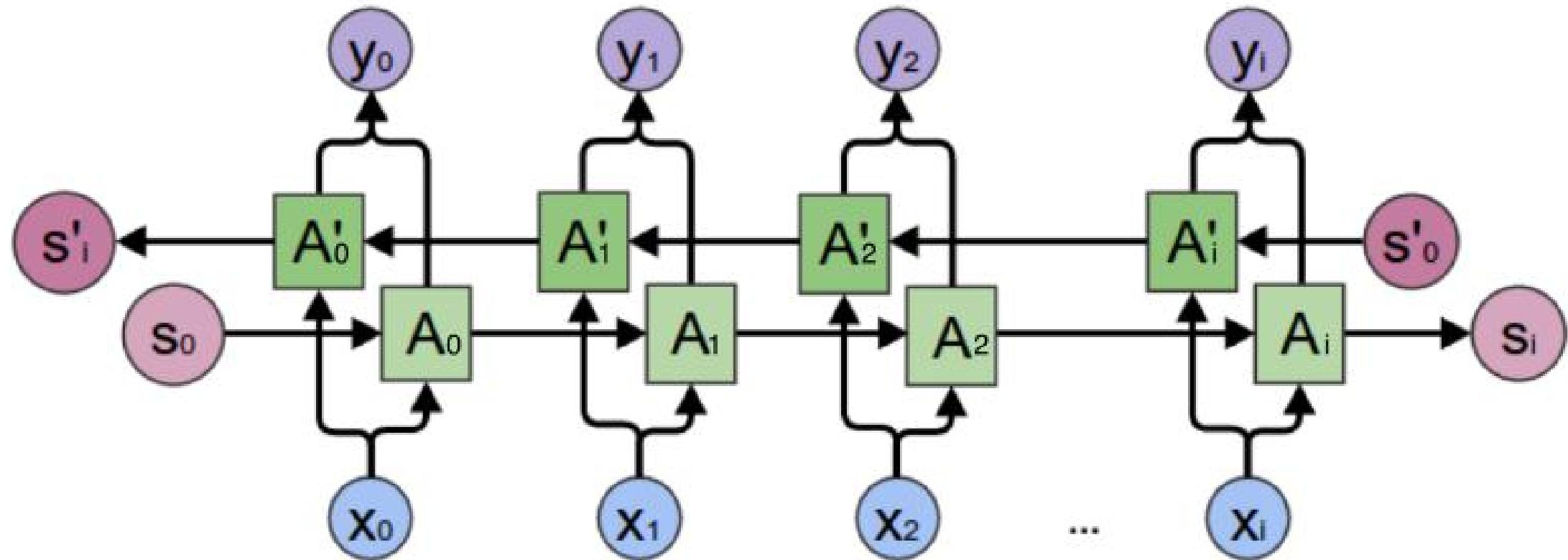
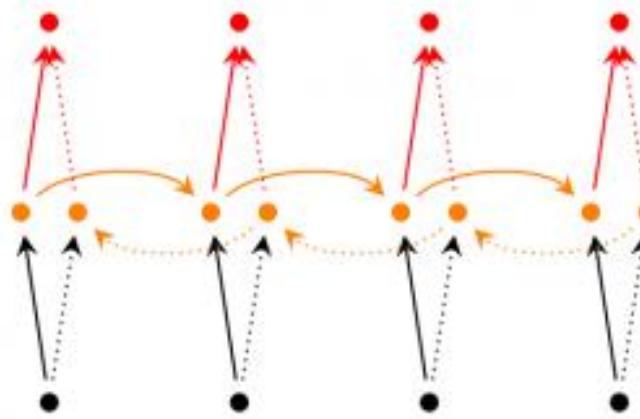


$$s_t = Ux_t + Wh_{t-1}$$

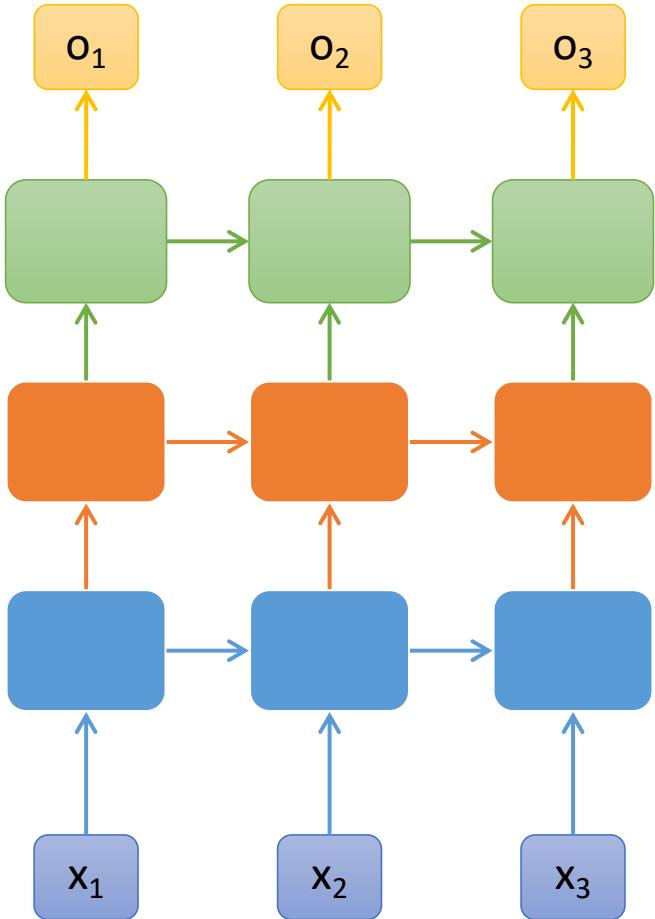
$$h_t = f(Ux_t + Wh_{t-1})$$

$$o_t = g(Vh_t)$$

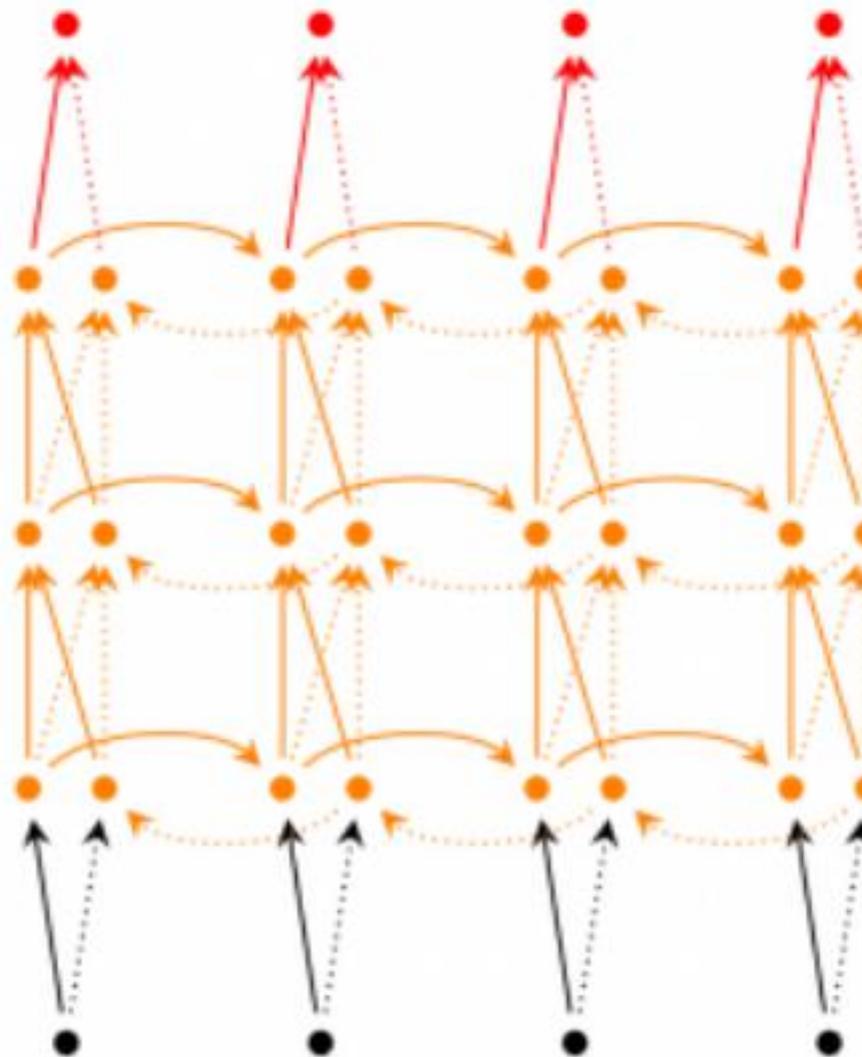
Bidirectional RNN回顾



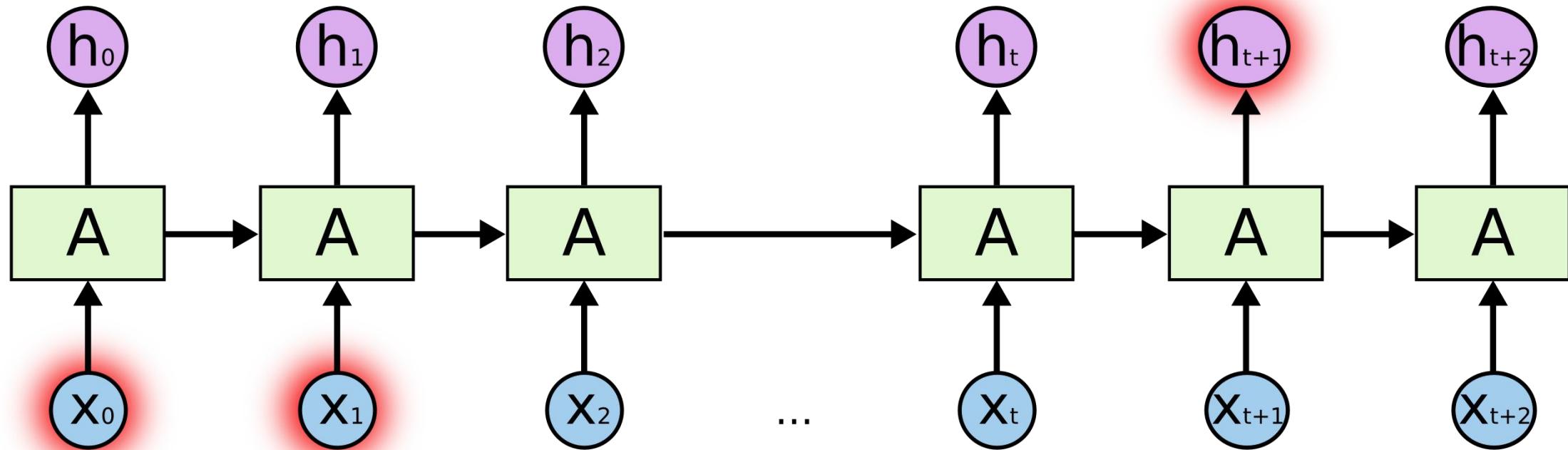
Deep RNN回顾



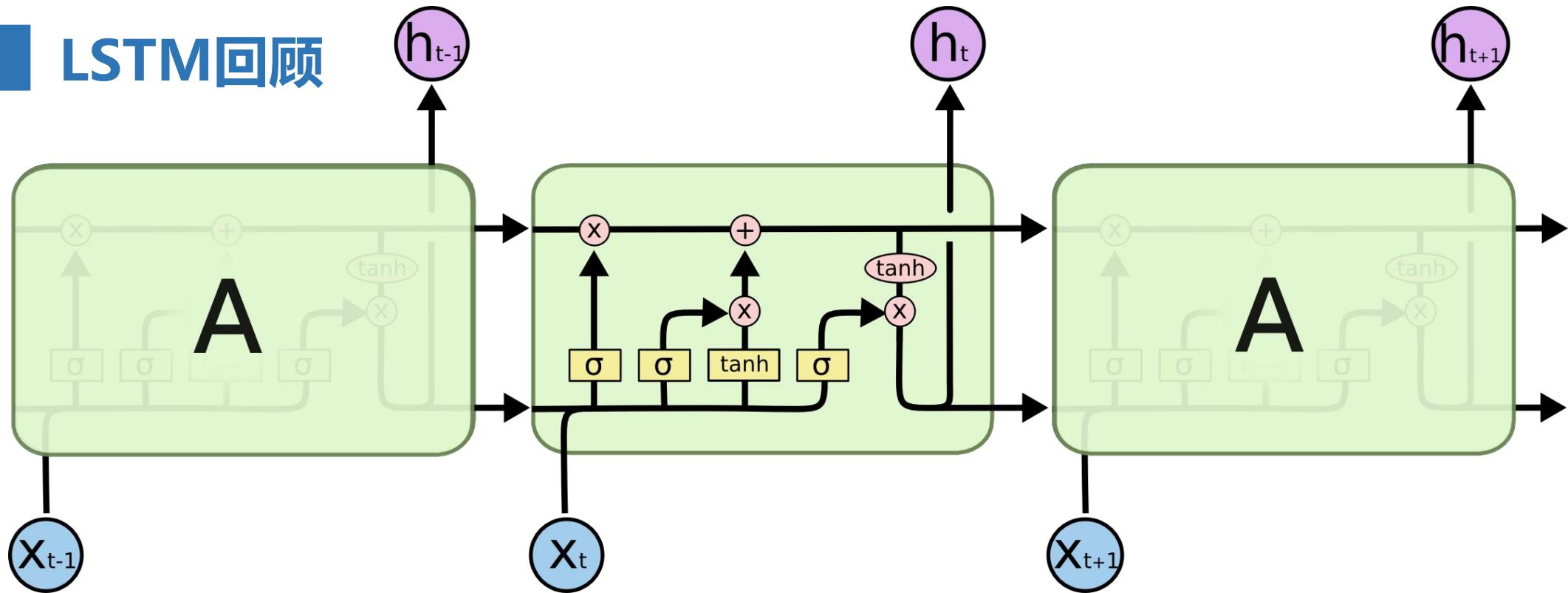
Deep Bidirectional RNN回顾



LSTM回顾



LSTM回顾



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

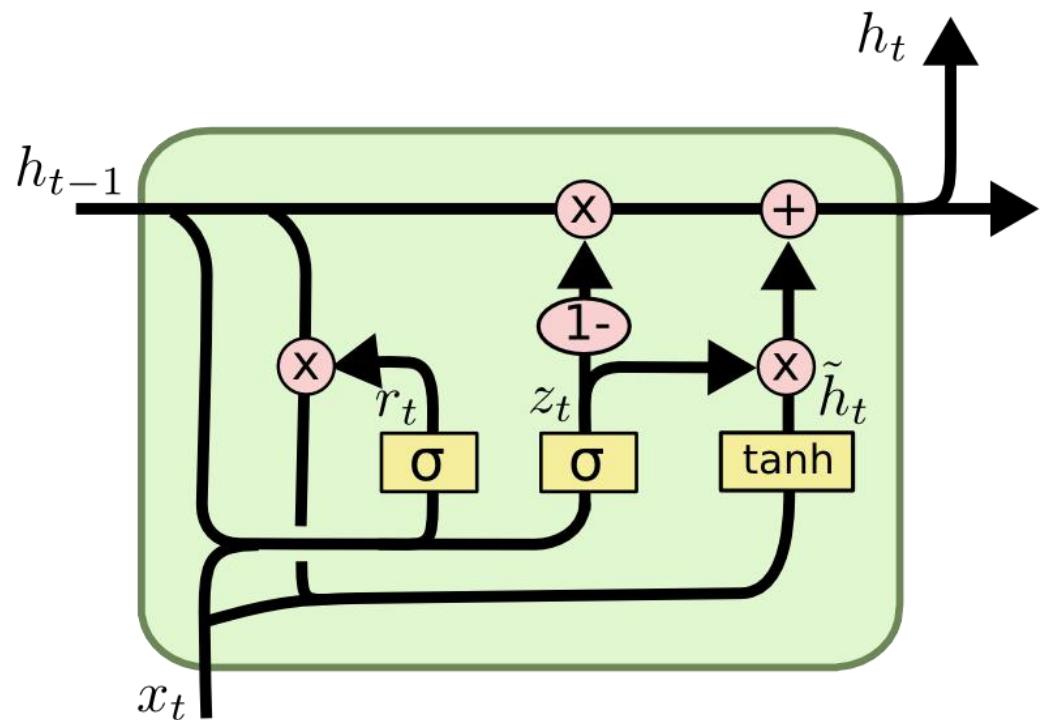
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

GRU回顾



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

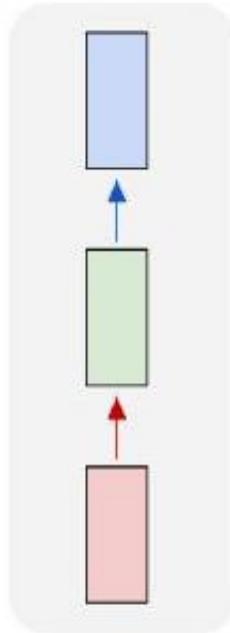
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

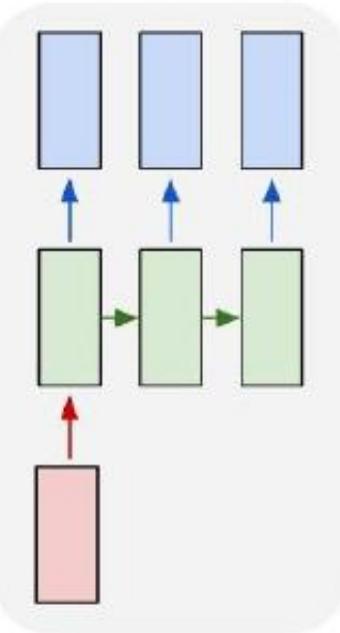
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

RNN结构回顾

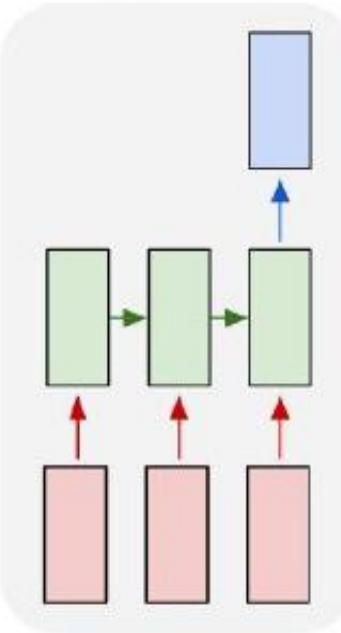
one to one



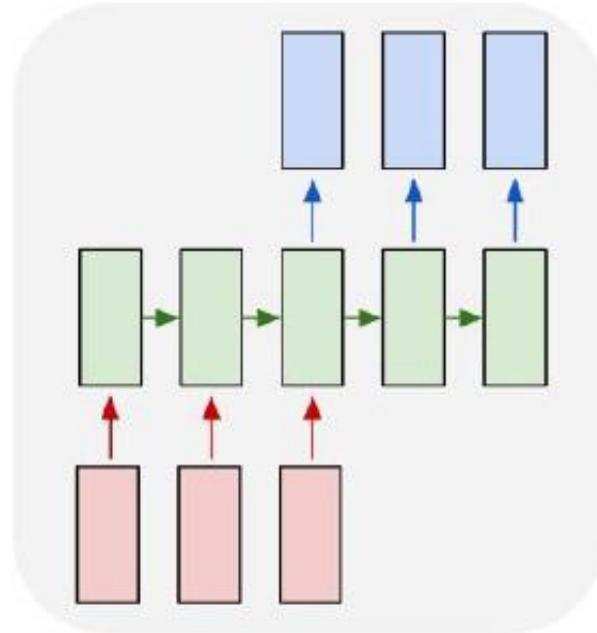
one to many



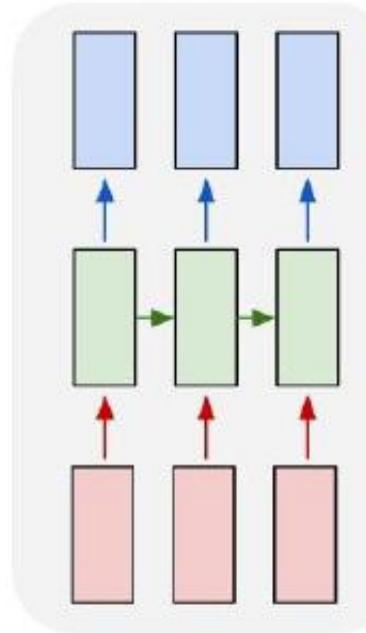
many to one



many to many

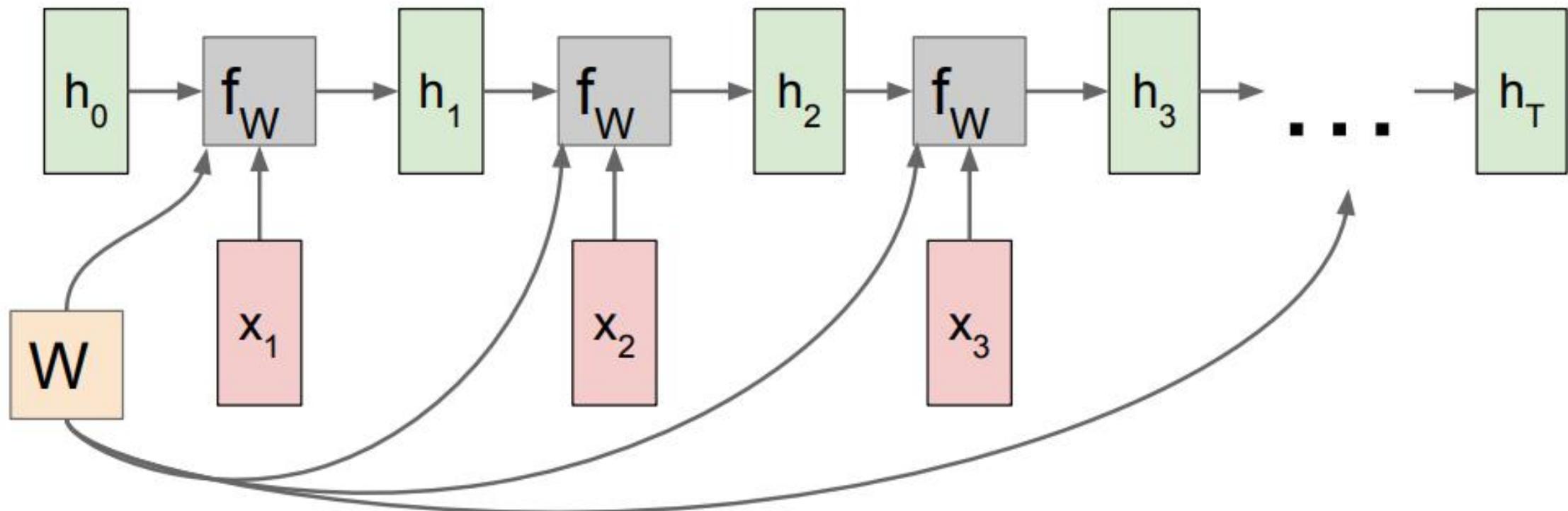


many to many



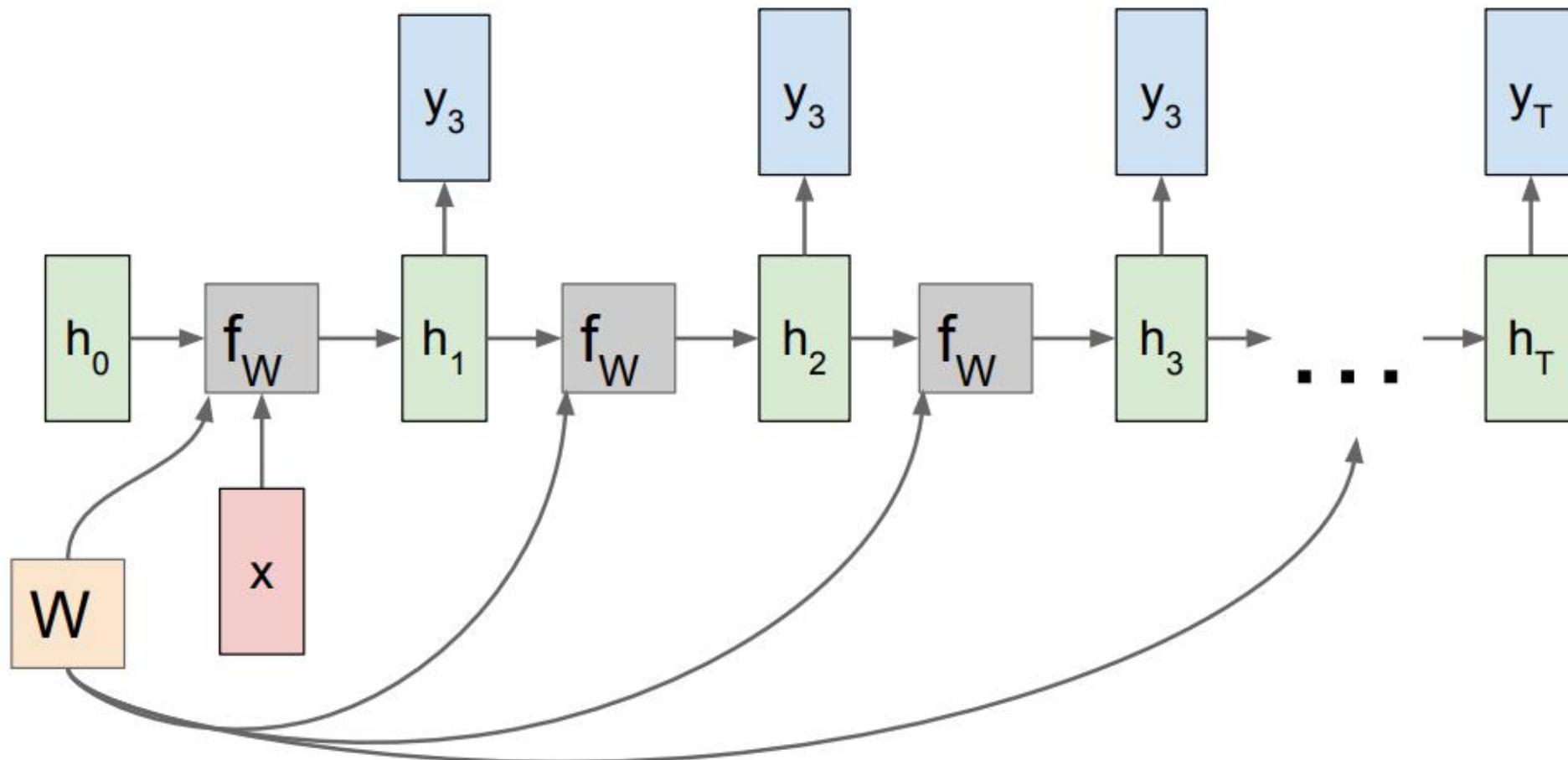
RNN结构回顾

Re-use the same weight matrix at every time-step



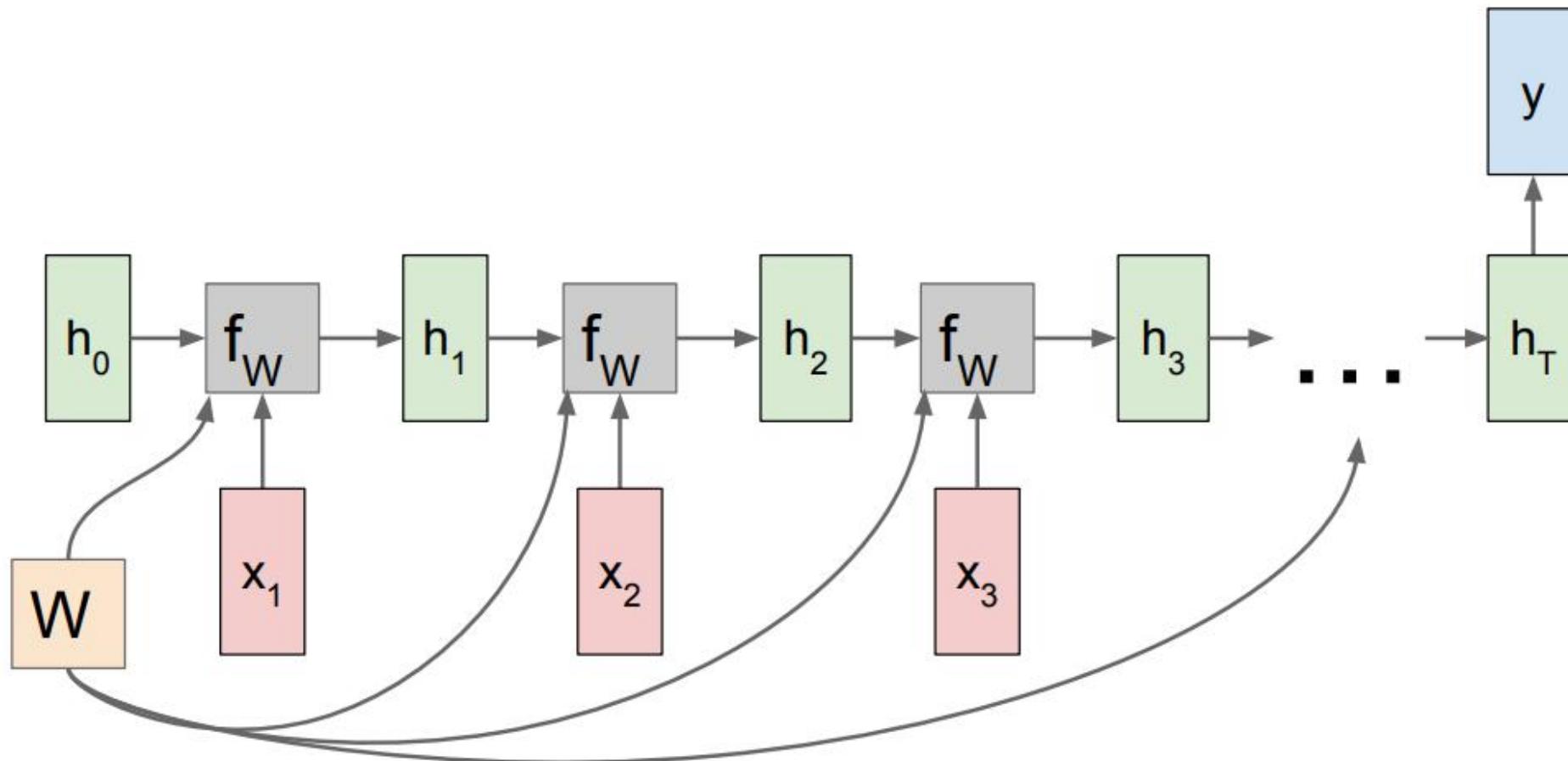
RNN结构回顾

RNN: Computational Graph: One to Many



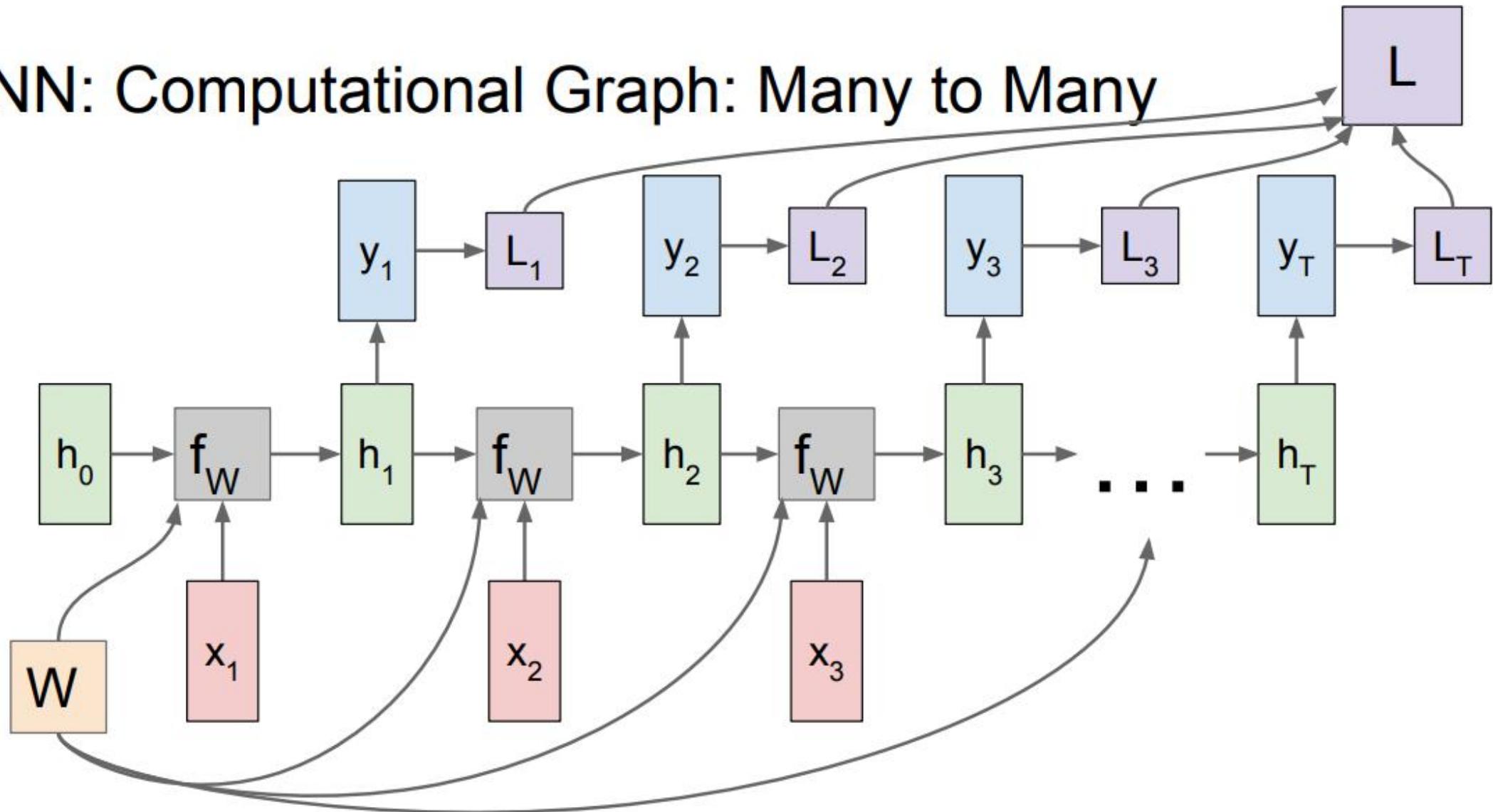
RNN结构回顾

RNN: Computational Graph: Many to One

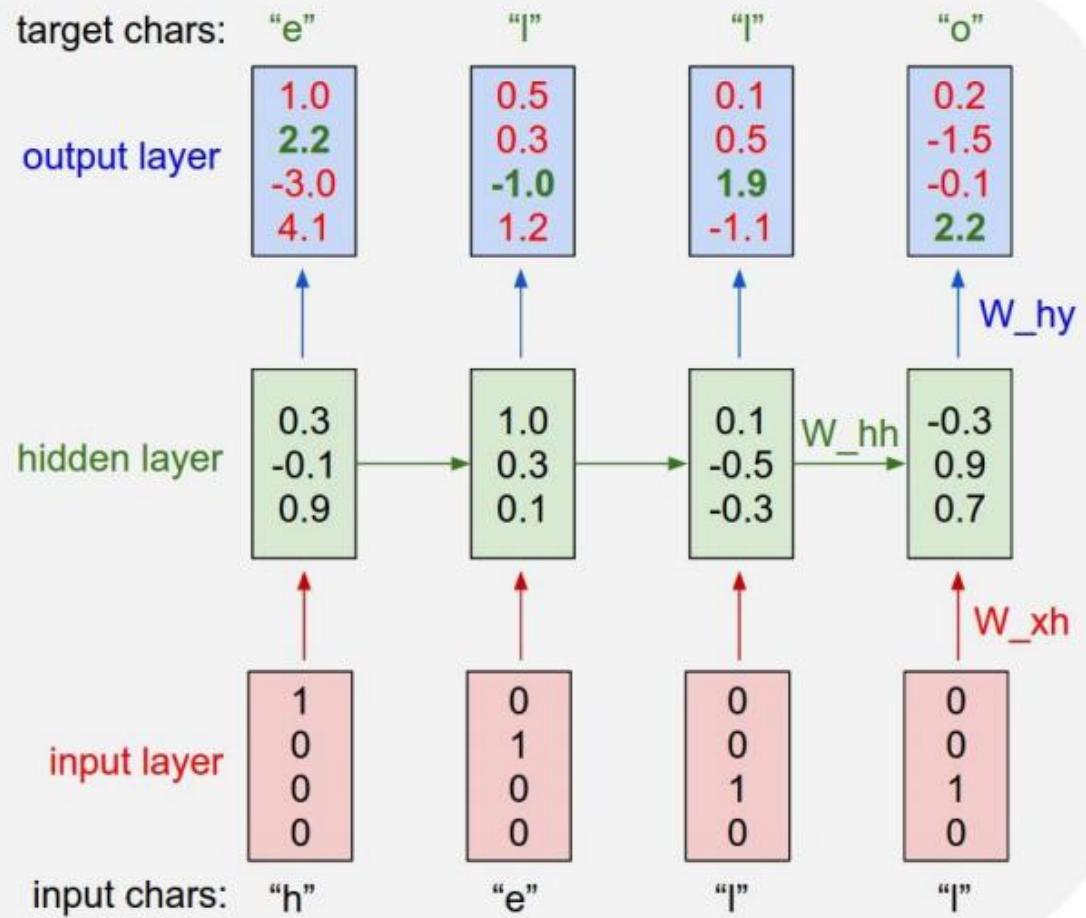


RNN结构回顾

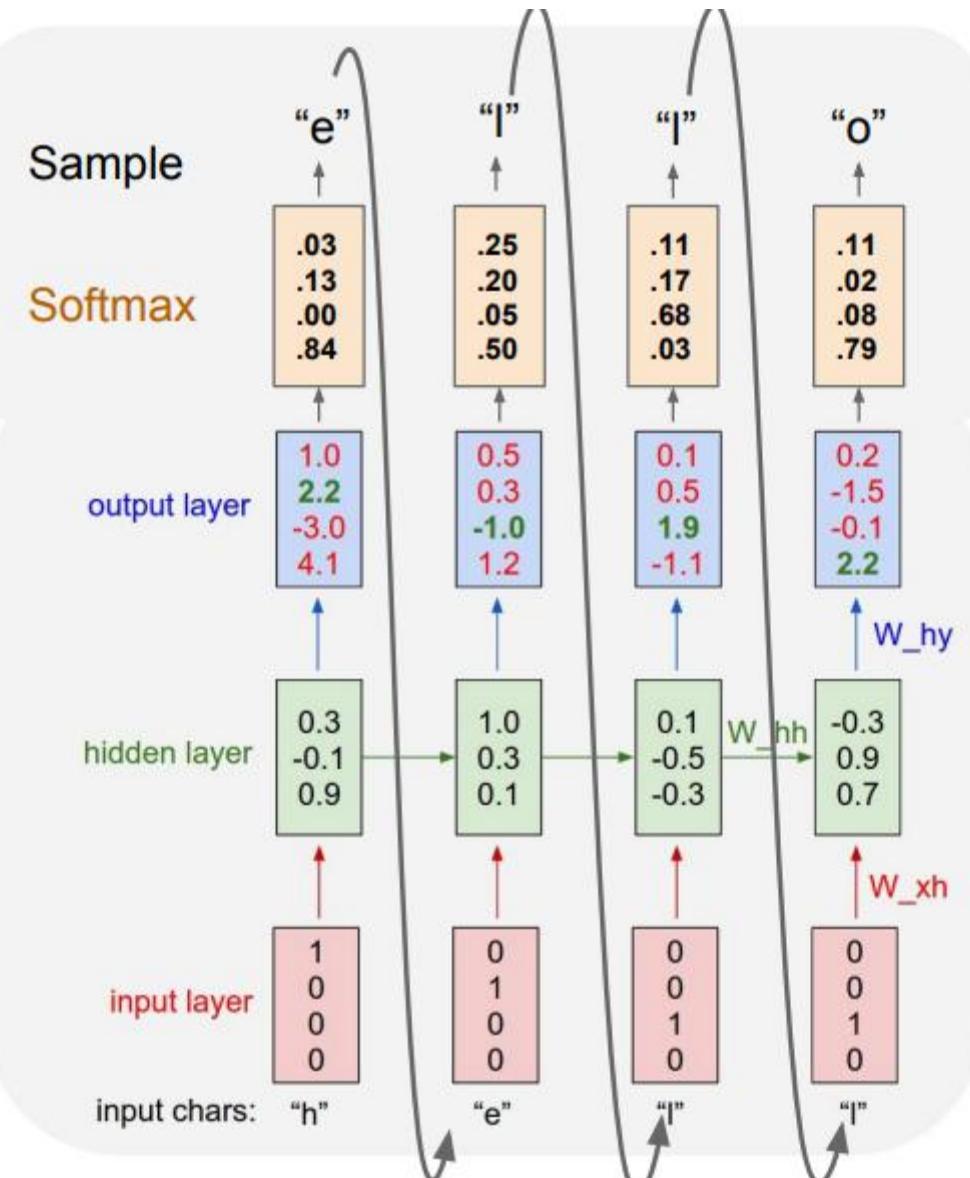
RNN: Computational Graph: Many to Many



RNN结构回顾



训练阶段



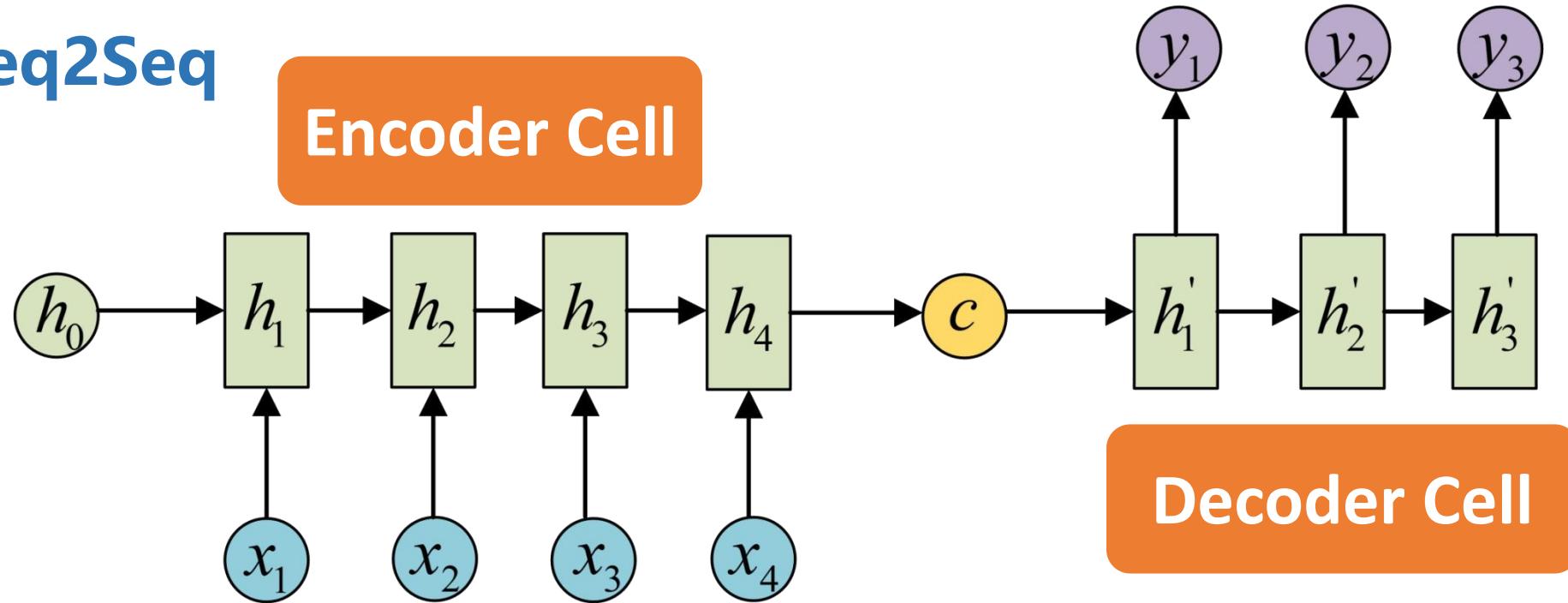
预测阶段

Seq2Seq

- Seq2Seq(Sequence to Sequence)，它被提出于2014年，最早由两篇文章独立地阐述了它主要思想，分别是Google Brain团队的《Sequence to Sequence Learning with Neural Networks》和Yoshua Bengio团队的《Learning Phrase Representation using RNN Encoder-Decoder for Statistical Machine Translation》。
- Seq2Seq属于一种Encoder-Decoder结构。



Seq2Seq

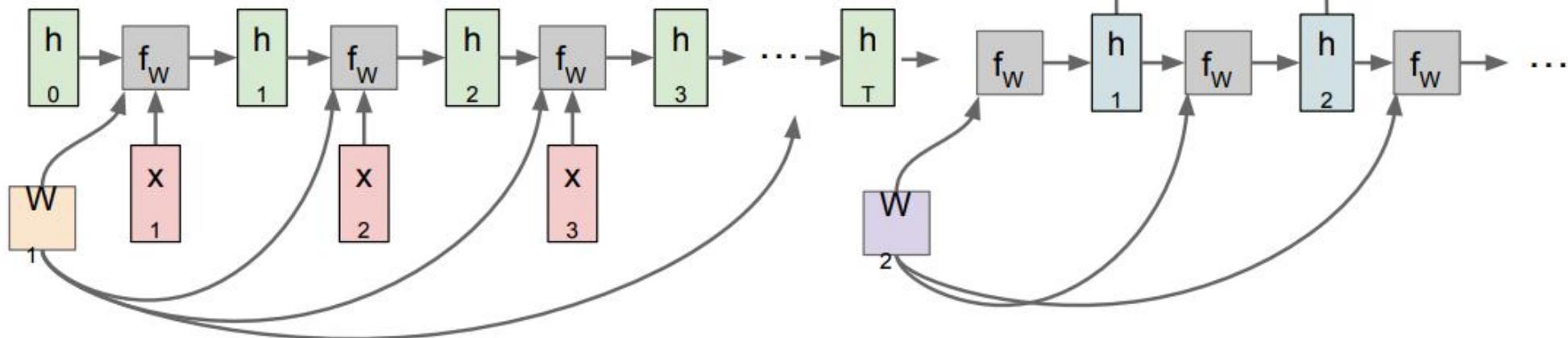


- Encoder-Decoder 这种结构的，其中 Encoder 是一个 RNNCell (RNN, GRU, LSTM 等) 结构。每个 timestep，我们向 Encoder 中输入一个字/词（一般是表示这个字/词的一个实数向量），直到我们输入这个句子的最后一个字/词 x_T ，然后输出整个句子的语义向量 c (一般情况下， $c = h_T = F([x_T; h_{T-1}]W)$ ， x_T 是最后一个 timestep 输入)。因为 RNN 的特点就是把前面每一步的输入信息都考虑进来了，所以理论上这个 c 就能够把整个句子的信息都包含了，我们可以把这个 c 当成这个句子的一个语义表示，也就是一个句向量。在 Decoder 中，我们根据 Encoder 得到的句向量 c ，一步一步地把蕴含在其中的信息分析出来。

Seq2Seq

Sequence to Sequence: Many-to-one + one-to-many

Many to one: Encode input sequence in a single vector



One to many: Produce output sequence from single input vector

Seq2Seq

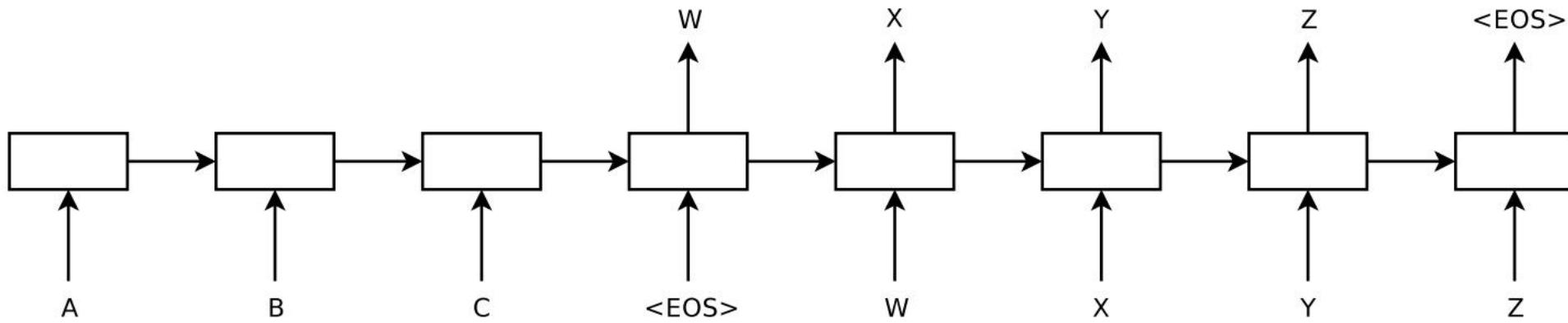
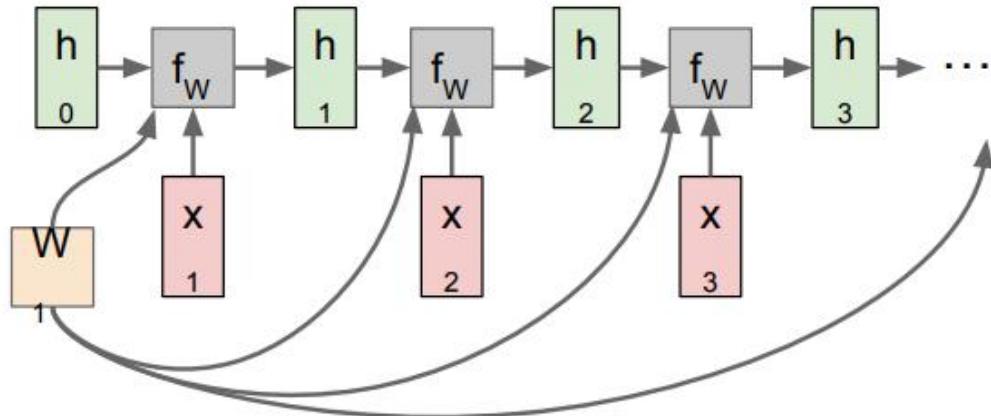


Figure 1: Our model reads an input sentence “ABC” and produces “WXYZ” as the output sentence. The model stops making predictions after outputting the end-of-sentence token. Note that the LSTM reads the input sentence in reverse, because doing so introduces many short term dependencies in the data that make the optimization problem much easier.

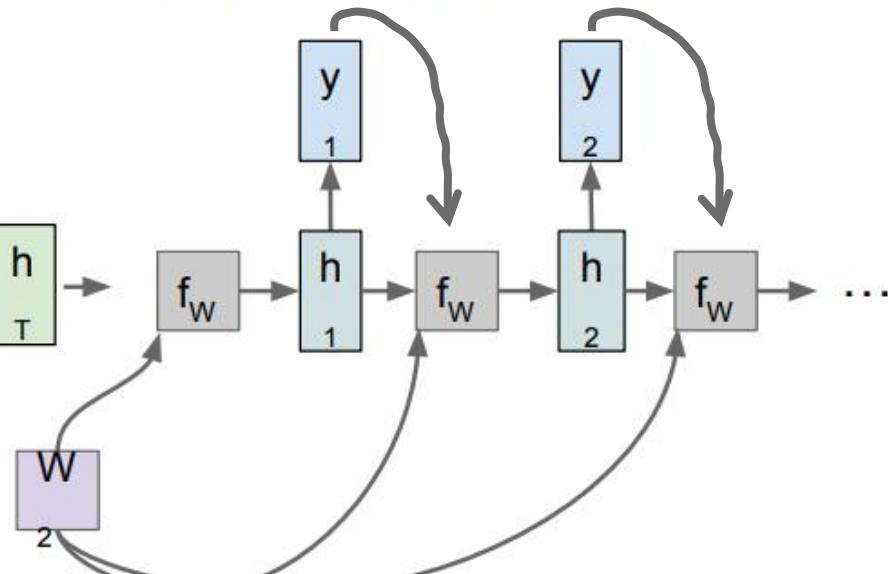
Seq2Seq

Sequence to Sequence: Many-to-one + one-to-many

Many to one: Encode input sequence in a single vector

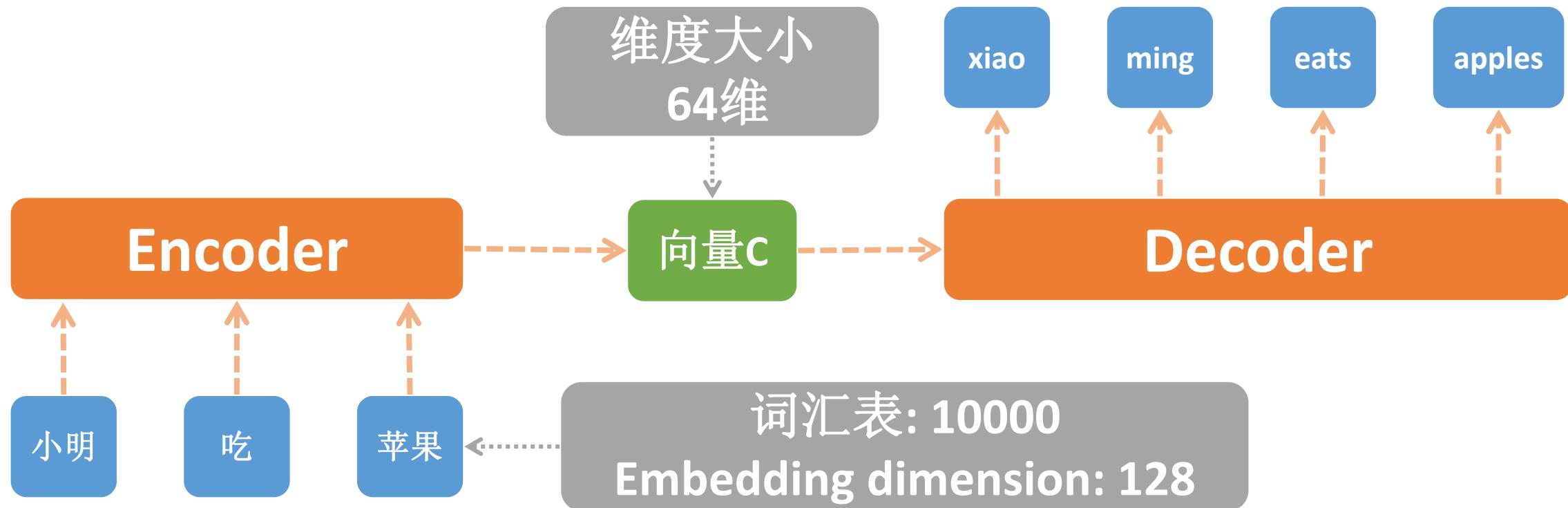


One to many: Produce output sequence from single input vector

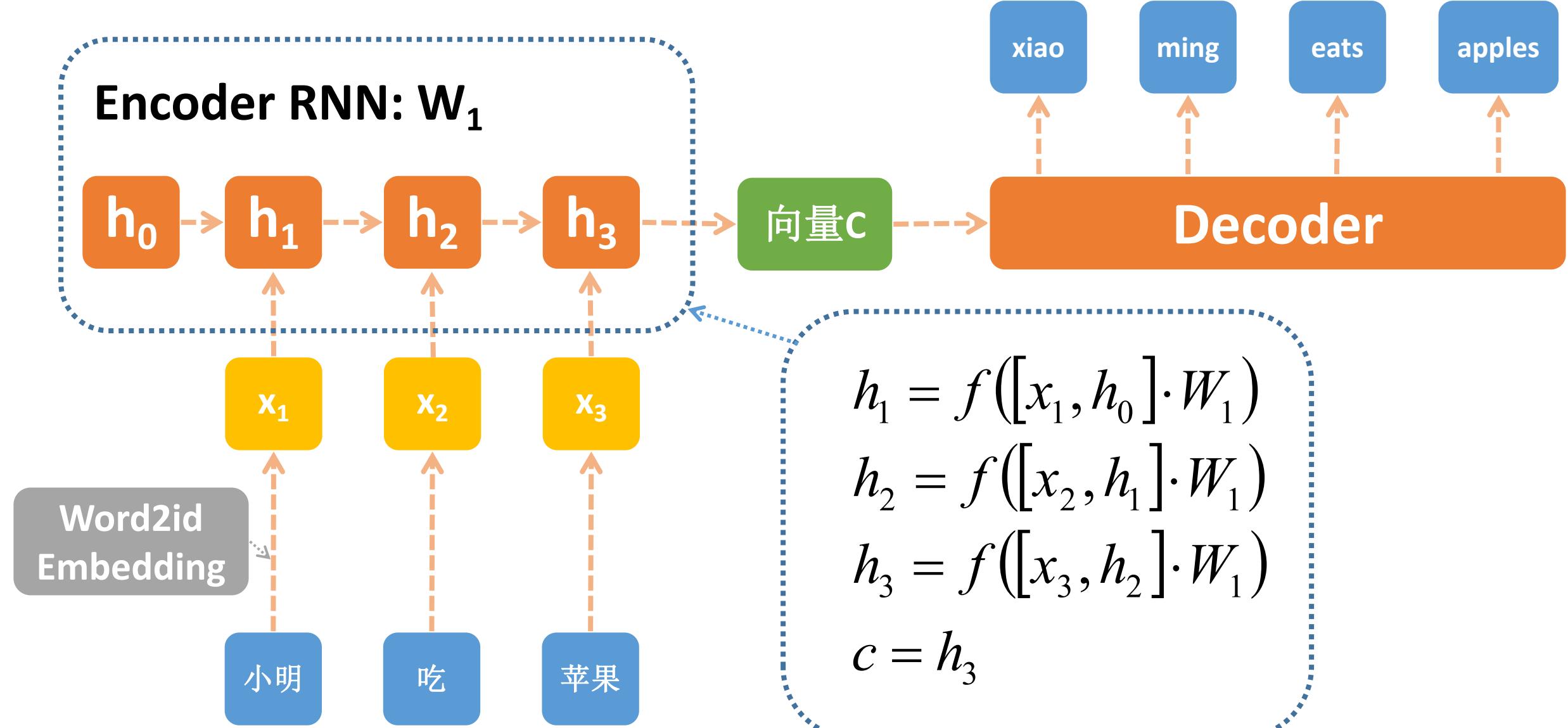


Seq2Seq理解

- 输入: 小明 吃 苹果
- 希望输出: xiao ming eats apples



Seq2Seq理解



Seq2Seq理解

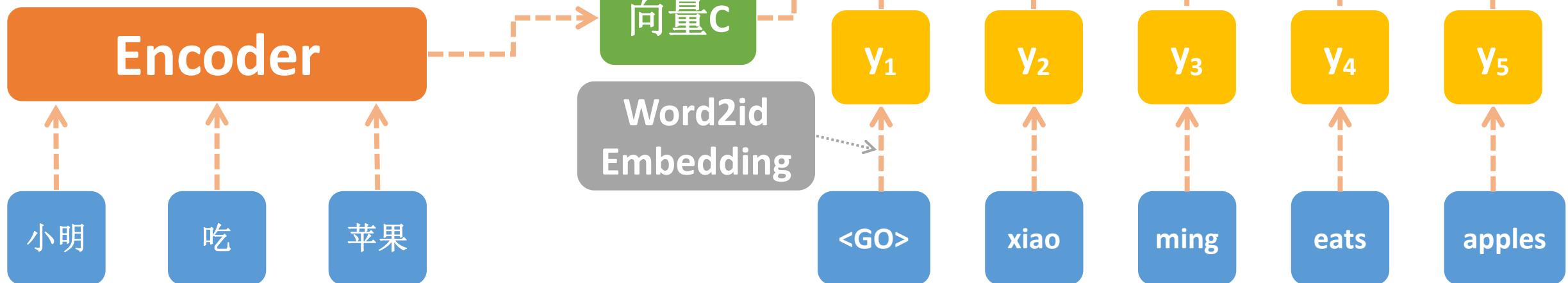
$$h_1 = f([y_1, c] \cdot W_2)$$

$$h_2 = f([y_2, h_1] \cdot W_2)$$

$$h_3 = f([y_3, h_2] \cdot W_2)$$

$$h_4 = f([y_4, h_3] \cdot W_2)$$

$$h_5 = f([y_5, h_4] \cdot W_2)$$



Seq2Seq理解

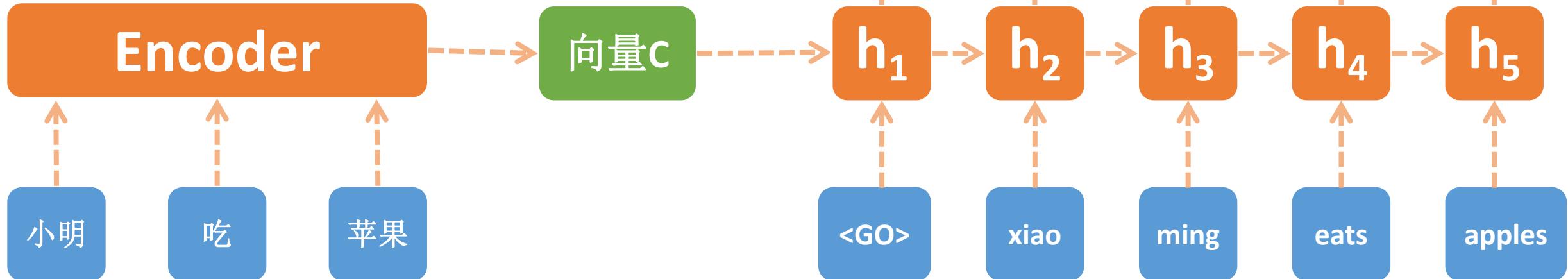
$$o_1 = g(h_1 \cdot U)$$

$$o_2 = g(h_2 \cdot U)$$

$$o_3 = g(h_3 \cdot U)$$

$$o_4 = g(h_4 \cdot U)$$

$$o_5 = g(h_5 \cdot U)$$



Seq2Seq理解

希望输出

xiao

ming

eats

apples

<EOS>

id2word

Y_1

Y_2

Y_3

Y_4

Y_5

argmax

$$Y_i = \arg \max(O_i, 1), i = 1, 2, 3, 4, 5$$

O_1

O_2

O_3

O_4

O_5

Encoder

向量C

Decoder

小明

吃

苹果

<GO>

xiao

ming

eats

apples

Seq2Seq理解

希望输出

xiao

ming

eats

apples

<EOS>

id2word

$$Y_i = \arg \max(O_i, 1), i = 1, 2, 3, 4, 5$$

Y_1

Y_2

Y_3

Y_4

Y_5

argmax

O_1

O_2

O_3

O_4

O_5

$$k_i = word2id(word_i), i = 1, 2, 3, 4, 5$$

$$p_i = O_i$$

$$E = \sum_{i=1}^5 E_i$$

$$E_i = -\ln p_{i, k_i}$$

LOSS

Seq2Seq理解预测

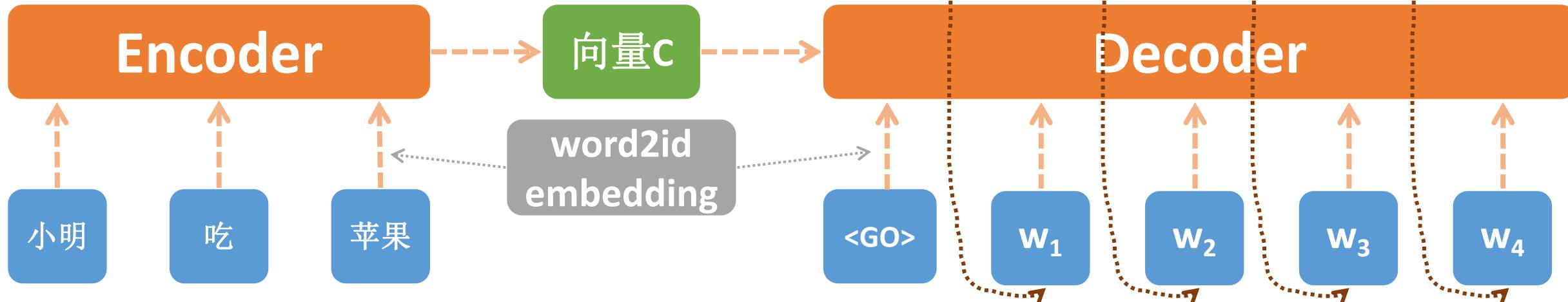
$$c = f(x, W_1)$$

$$h_i = f(x, c, w_1, w_2, \dots, w_{i-1}, W_2)$$

$$o_i = g(h_i, U)$$

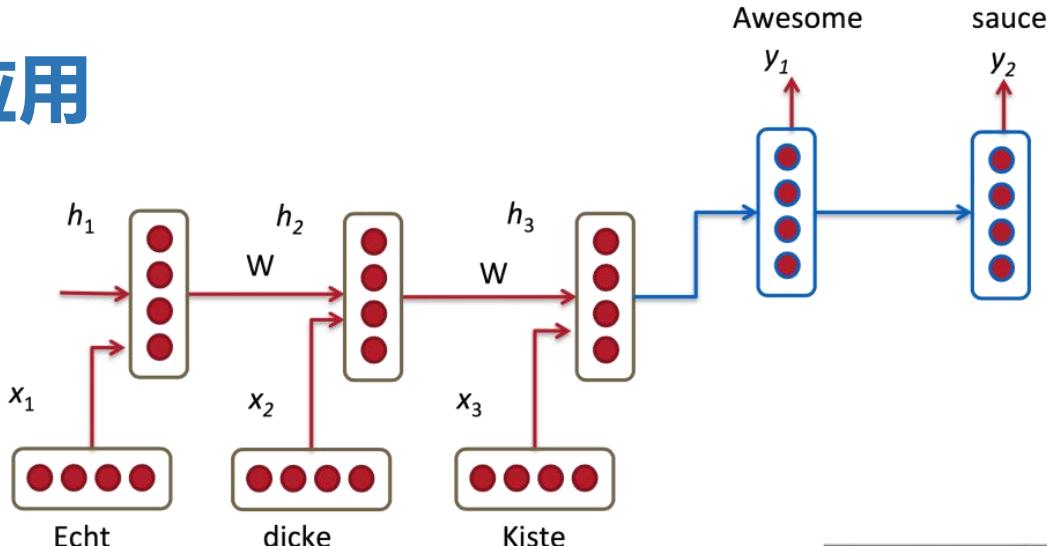
$$Y_i = \text{argmax}(o_i, 1)$$

$$w_i = id2word(Y_i)$$



Seq2Seq应用

- 机器翻译



- 自动对话机器人

Post	Model	Response(Original)	Response(Translated)
比情人节更庸俗的事就是，对情人节的接待规格待遇提出要求。 What is more vulgar than Valentine's day is to make specific demands on the reception of the Valentine's day.	Seq2Seq ECM	这才是真正的情人节。	This is a real Valentine's day.
		Like 情人节快乐！	Happy Valentine's Day!
		Happiness 哈哈，这也太浪漫了吧！	Aha, this is too romantic!
		Sadness 我也想要这样的情人节。	I also want this kind of Valentine's day, sigh.
		Disgust 这就是所谓的情人节！	This is the so-called Valentine's day!
		Anger 这是赤裸裸的炫耀啊！	This is shameless show-off!
知道真相的我眼泪笑下来。 Knowing the truth, I laughed myself to tears.	Seq2Seq ECM	知道真相的我眼泪掉下来。	After knowing the truth, I burst into tears.
		Like 太可爱了，太可爱了！	Too lovely, it's too lovely.
		Happiness 笑得我眼泪都出来了。	I can't help laugh myself to tears.
		Sadness 知道真相的我眼泪掉下来。	After knowing the truth, I burst into tears.
		Disgust 这是真的还是假的？	Is this true?
		Anger 什么意思？没看懂。	What do you mean? I don't understand.

ID	Post	NRM-glo	NRM-loc	NRM-hyb	Rtr-based
1	每年元旦都来高烧。 High fever attacks me every New Year's Day	多休息，增强免疫力！ Rest more to boost your immune system.	每年元旦都去医院看病 Go to see a doctor every New Year's Day	祝早日康复，身体健康！ Get well soon and stay healthy!	高烧... High fever...
2	又老了一岁，感谢我们组，我很高兴 I gain one more year. Grateful to my group, so happy	生日快乐！永远支持你！ Happy birthday! Will stand by you forever!	生日快乐！永远年轻漂亮！ Happy birthday! Wish you always stay young and pretty!	老了老了，岁月不饶人啊！ Getting old now. Time has no mercy.	我只是老了一岁 I just gain only one more year.
3	吃素第一天，坚持住，崔朵拉。 First day of being a vegetarian. Hold on, CuiDuoLa	必须来支持！加油！ I come to offer my support. Keep it up!	坚持住，坚持下去。 Hold on, keep it up.	每天都是素食主义者 I have been being a vegetarian everyday	坚持住！ Hold on!
4	这种天气真应该带个帐篷出去踏青赏花 We should go out with some cute guys to enjoy a great outing in such a nice weather.	这是哪儿啊？真美！ Where is it? It is so beautiful!!	阳光明媚，心情舒畅。 Such a nice sunny day!	这种天气真不是盖的。 I am in a great mood.	文山啊出去踏青寻找灵感哈哈 WenShan, let's go out to get some inspiration. Ha! Ha!

Seq2Seq应用

- 文档摘要自动生成

Good quality summary output	
S: a man charged with the murder last year of a british backpacker confessed to the slaying on the night he was charged with her killing , according to police evidence presented at a court hearing tuesday . ian douglas previte , ## , is charged with murdering caroline stuttle , ## , of yorkshire , england	
T: man charged with british backpacker 's death confessed to crime police officer claims	
O: man charged with murdering british backpacker confessed to murder	

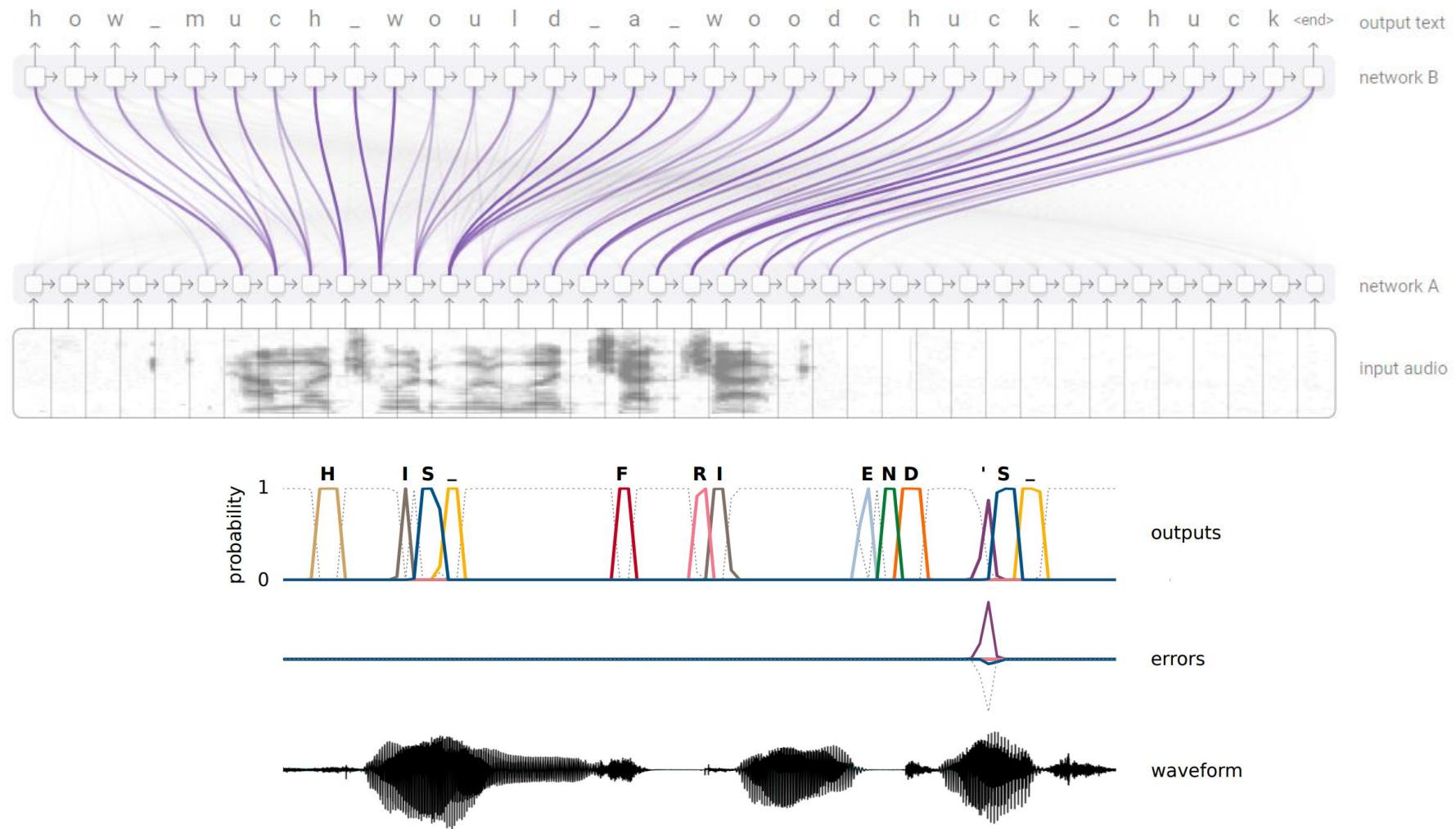
- 文本生成

<p>秋夕湖上 By a Lake at Autumn Sunset 一夜秋凉雨湿衣， A cold autumn rain wetted my clothes last night, 西窗独坐对夕晖。 And I sit alone by the window and enjoy the sunset. 湖波荡漾千山色。 With mountain scenery mirrored on the rippling lake, 山鸟徘徊万籁微。 A silence prevails over all except the hovering birds.</p>	<p>秋夕湖上 By a Lake at Autumn Sunset 荻花风里桂花浮。 The wind blows reeds with osmanthus flying, 根竹生云翠欲流。 And the bamboos under clouds are so green as if to flow down. 谁拂半湖新镜面。 The misty rain ripples the smooth surface of lake, 飞来烟雨暮天愁。 And I feel blue at sunset .</p>
---	---

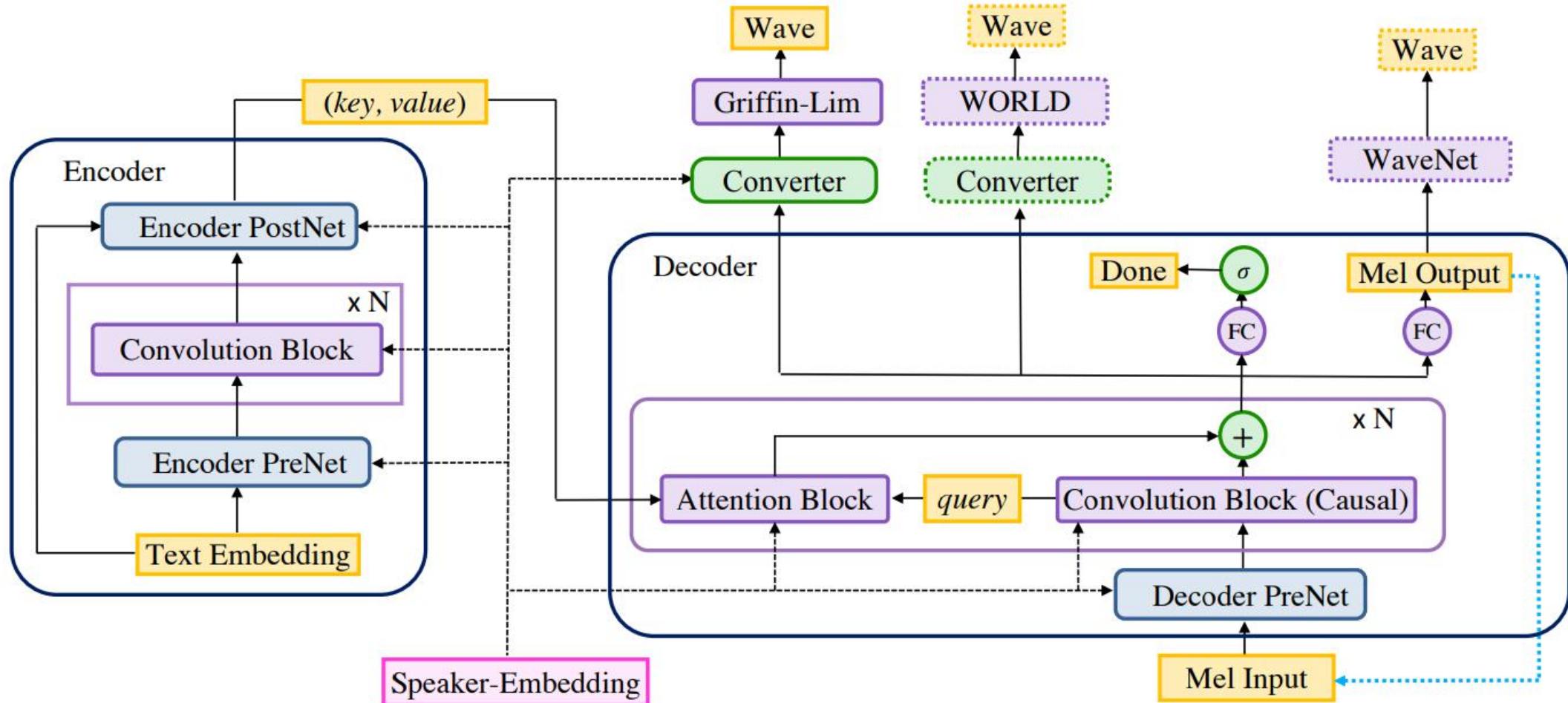
▲ 图5：左边是机器生成的诗词，右边是一首宋代诗词

Seq2Seq应用

- 语音识别/合成/语音-文本转换



Seq2Seq应用



百度Deep Voice v3

Seq2Seq应用

- 图片描述自动生成



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."

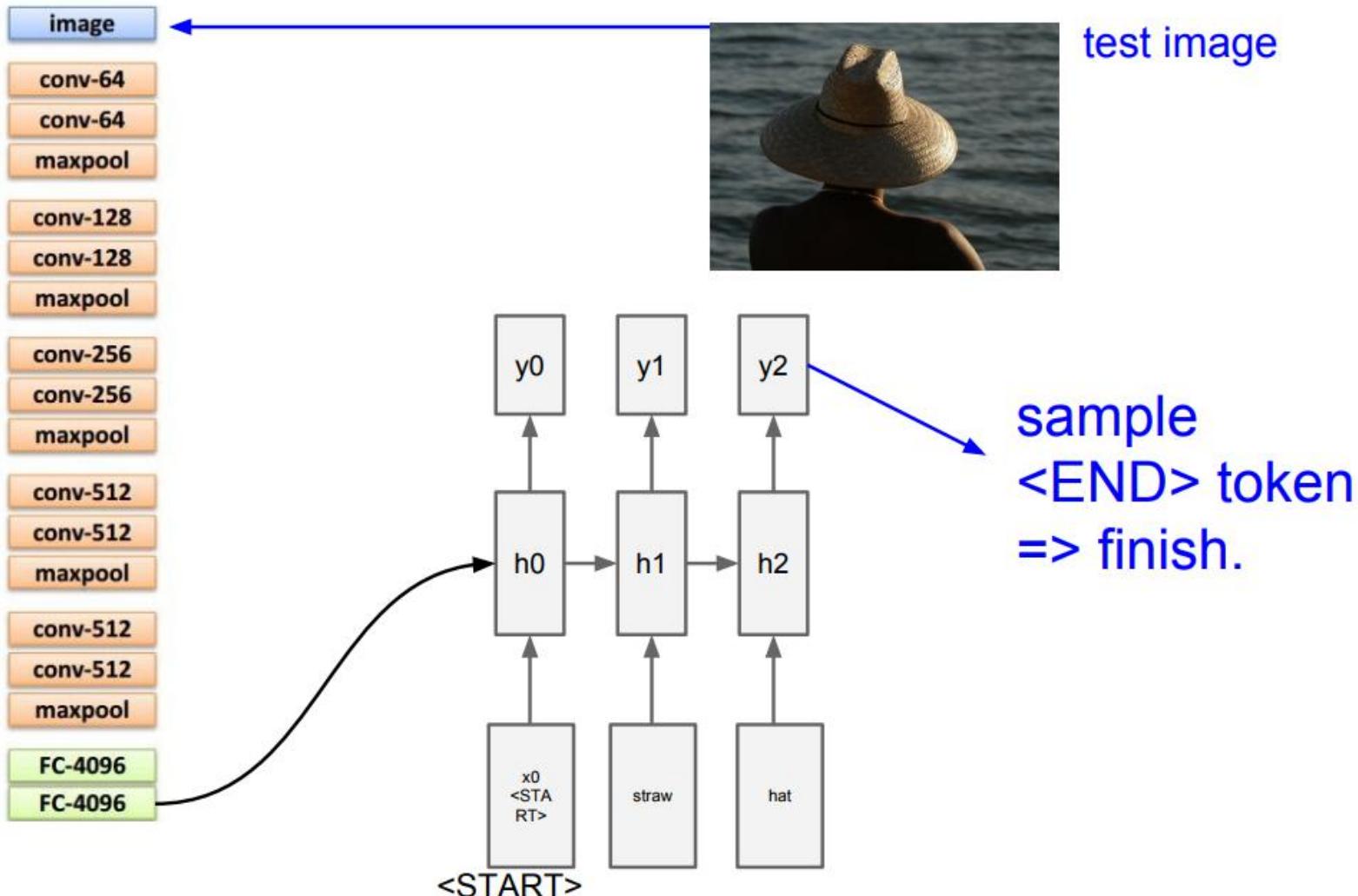


"young girl in pink shirt is swinging on swing."



"man in blue wetsuit is surfing on wave."

Seq2Seq应用



Seq2Seq应用

- Visual Question Answering(VQA, 视觉问答系统)

- <https://visualqa.org/index.html>

A photograph of a young man with dark hair and a beard, wearing a red t-shirt with a white logo, blue pants, and a green knit cap, performing a trick on a longboard on a wooden ramp. He is leaning forward, pushing off with one foot. The ramp is made of light-colored wood and is set on a paved surface. In the background, there is a colorful carnival booth with various games and prizes visible through the glass windows. Other people are standing near the booth, and a person in a red shirt is walking away from the camera. The overall atmosphere is casual and suggests a fun day at the fair.

GT Question 戴帽子的男孩在干什么?
What is the boy in green cap doing?

GT Answer 他在玩滑板。
He is playing skateboard.



GT Question 房间里的沙发是什么质地的?
What is the texture of the sofa in
the room?

GT Answer 布艺。
Cloth.



图片中有人么?
Is there any person
in the image?

有。
Yes.



这个人挑菜么?
Is the man trying to
buy vegetables?
是的。
Yes.

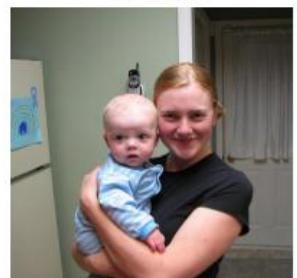
Who is wearing glasses?
man woman



Is the umbrella upside down?
yes no



Where is the child sitting?
fridge arms



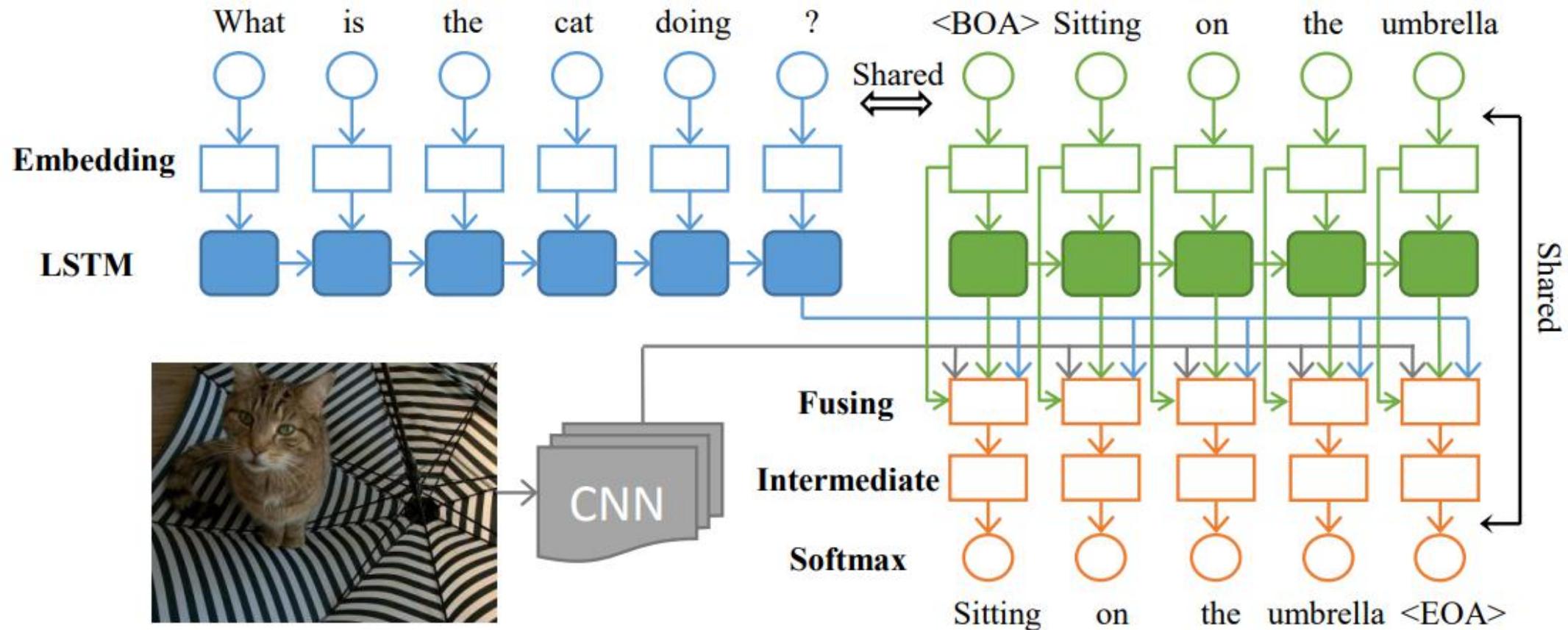
How many children are in the bed?
2 1

2

1



Seq2Seq应用

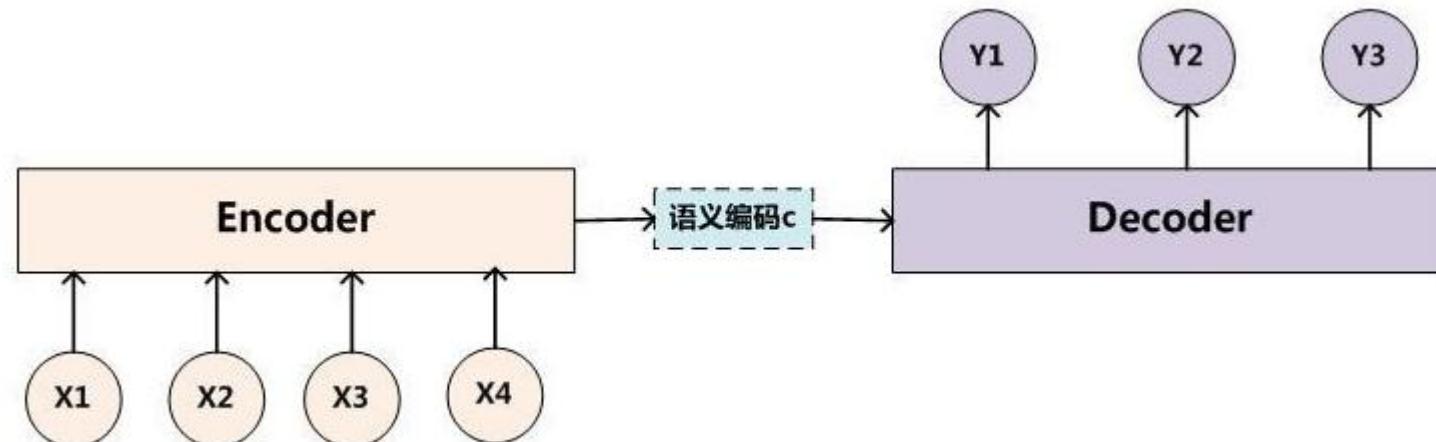


Seq2Seq应用总结

- 总而言之，Seq2Seq应用场景，包括了经典的机器翻译、文本摘要和对话生成等，也包括了一些非常有趣的应用，比如：根据公式图片生成 latex 代码，生成 commit message 等。自然语言生成（NLG）是一个非常有意思，也非常有前途的研究领域，简单地说，就是解决一个条件概率 $p(\text{output} | \text{context})$ 的建模问题，即根据 context 来生成 output，这里的 context 可以非常灵活多样，大家都是利用深度学习模型对这个条件概率进行建模，同时加上大量的训练数据和丰富的想象力，可以实现很多有趣的工作。Seq2Seq 是一个简单易用的框架，开源的实现也非常多，但并不意味着直接生搬硬套就可以了，需要具体问题具体分析。此外，对于生成内容的控制，即 decoding 部分的研究也是一个非常有意思的方向，比如：如何控制生成文本的长度，控制生成文本的多样性，控制生成文本的信息量大小，控制生成文本的情感等等。

Seq2Seq原理

- 最基础的Seq2Seq模型包含了三个部分，即Encoder、Decoder以及连接两者的中间状态向量，Encoder通过学习输入，将其编码成一个固定大小的状态向量c，继而将c传给Decoder，Decoder再通过对状态向量c的学习来进行输出。下图中，图中每一个box代表了一个**RNN Cell**单元，通常是**LSTM**或者**GRU**。



Seq2Seq

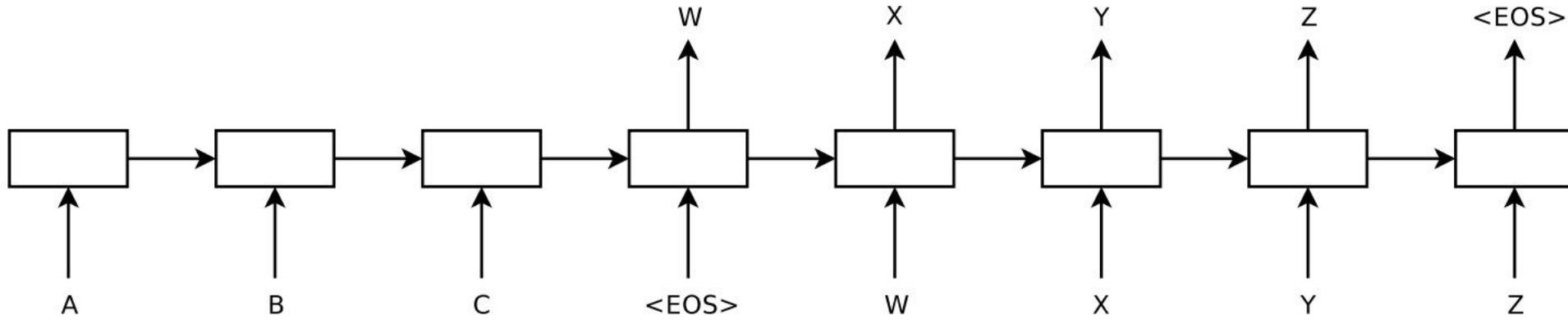


Figure 1: Our model reads an input sentence “ABC” and produces “WXYZ” as the output sentence. The model stops making predictions after outputting the end-of-sentence token. Note that the LSTM reads the input sentence in reverse, because doing so introduces many short term dependencies in the data that make the optimization problem much easier.

Seq2Seq

- Encoder-Decoder框架可以这么直观地去理解：可以把它看作适合处理由一个句子（或篇章）生成另外一个句子（或篇章）的通用处理模型。对于句子对 $\langle X, Y \rangle$ ，我们的目标是给定输入句子 X ，期待通过Encoder-Decoder框架来生成目标句子 Y 。 X 和 Y 可以是同一种语言，也可以是两种不同的语言。而 X 和 Y 分别由各自的单词序列构成：

$$X = (x_1, x_2, \dots, x_m)$$

$$Y = (y_1, y_2, \dots, y_n)$$

Seq2Seq

- Encoder顾名思义就是对输入句子 x 进行编码，将输入句子通过非线性变换转化为中间语义表示 C :

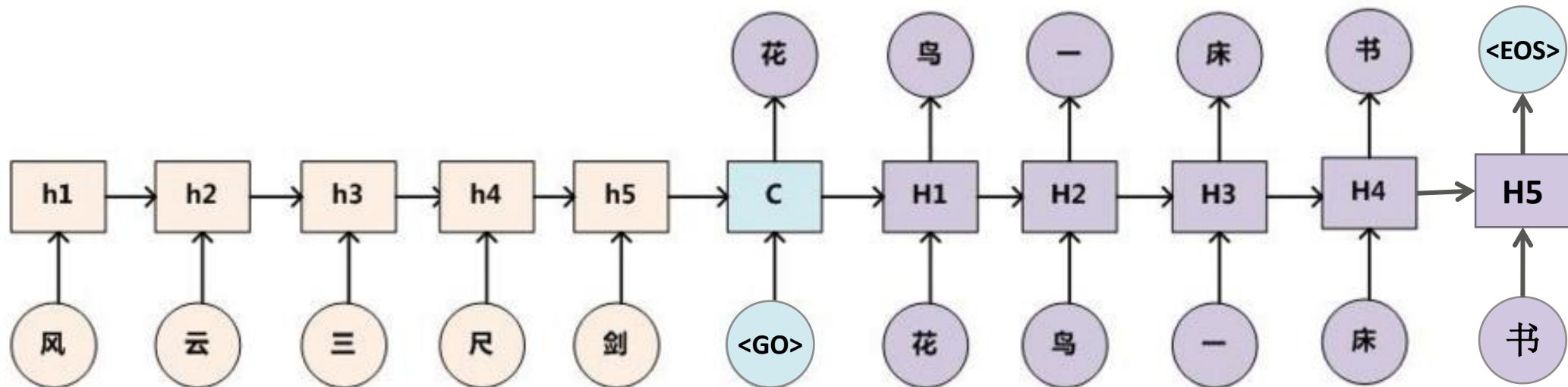
$$C = F(x_1, x_2, \dots, x_m)$$

- 对于解码器Decoder来说，其任务是根据句子 x 的中间语义表示 C 和之前已经生成的历史信息 y_1, y_2, \dots, y_{i-1} 来生成 i 时刻要生成的单词 y_i :每个 y_i 都依次这么产生，那么看起来就是整个系统根据输入句子 x 生成了目标句子 y 。

$$y_i = G(C, y_1, y_2, \dots, y_{i-1})$$

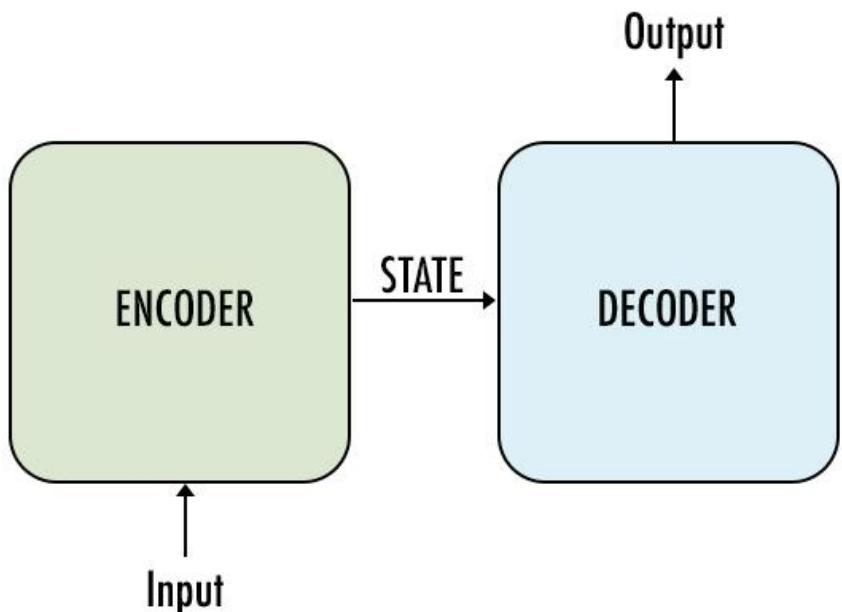
Seq2Seq案例

- 只需要找到大量的对联数据对这个模型进行训练，那么即可利用这个模型，输入上联，机器自动产生下联了。



Seq2Seq案例

Inputs	Target
How are you?	I am good
Can you fly that thing?	Not yet



编码器: [`tf.nn.dynamic_rnn`](#)
解码器: [`tf.contrib.seq2seq.dynamic_rnn_decoder`](#)

举例: `tf.nn.dynamic_rnn(cell, inputs, sequence_length=None, initial_state=None, dtype=None, parallel_iterations=None, swap_memory=False, time_major=False, scope=None)`

Cell为前面构建的RNN

`cell(tf.contrib.rnn.BasicLSTMCell);`

Inputs,为输入的文本数据, 通常是嵌入层的输出。
以及initial_state

Seq2Seq

- <PAD> 在训练中，我们将数据按批次输入。但同一批次中必须有相同的Sequence Length(序列长度 /time_steps)。所以我们会用<PAD>填充较短的输入。
- <EOS> 它能告诉解码器句子在哪里结束，并且它允许解码器在其输出中表明句子结束的位置
- <UNK> 忽视词汇表中出现频率不够高而不足以考虑在内的文字,将这些单词替换为 <UNK>
- <GO> 解码器的第一个时间步骤的输入，以使解码器知道何时开始产生输出

Seq2Seq案例

0	<PAD>	11	can
1	<EOS>	12	you
2	<UNK>	13	fly
3	<GO>	14	that
4	how	15	thing
5	are	16	not
6	you	17	yet
7	?		
8	i		
9	am		
10	good		

Seq2Seq案例

Inputs

How are you?

Can you fly that thing?

Target

I am good

Not yet



how
can

are
you

you
fly

?
that

<PAD>
thing

<PAD>
?

Seq2Seq案例

Inputs

How are you?

Can you fly that thing?



Target

I am good

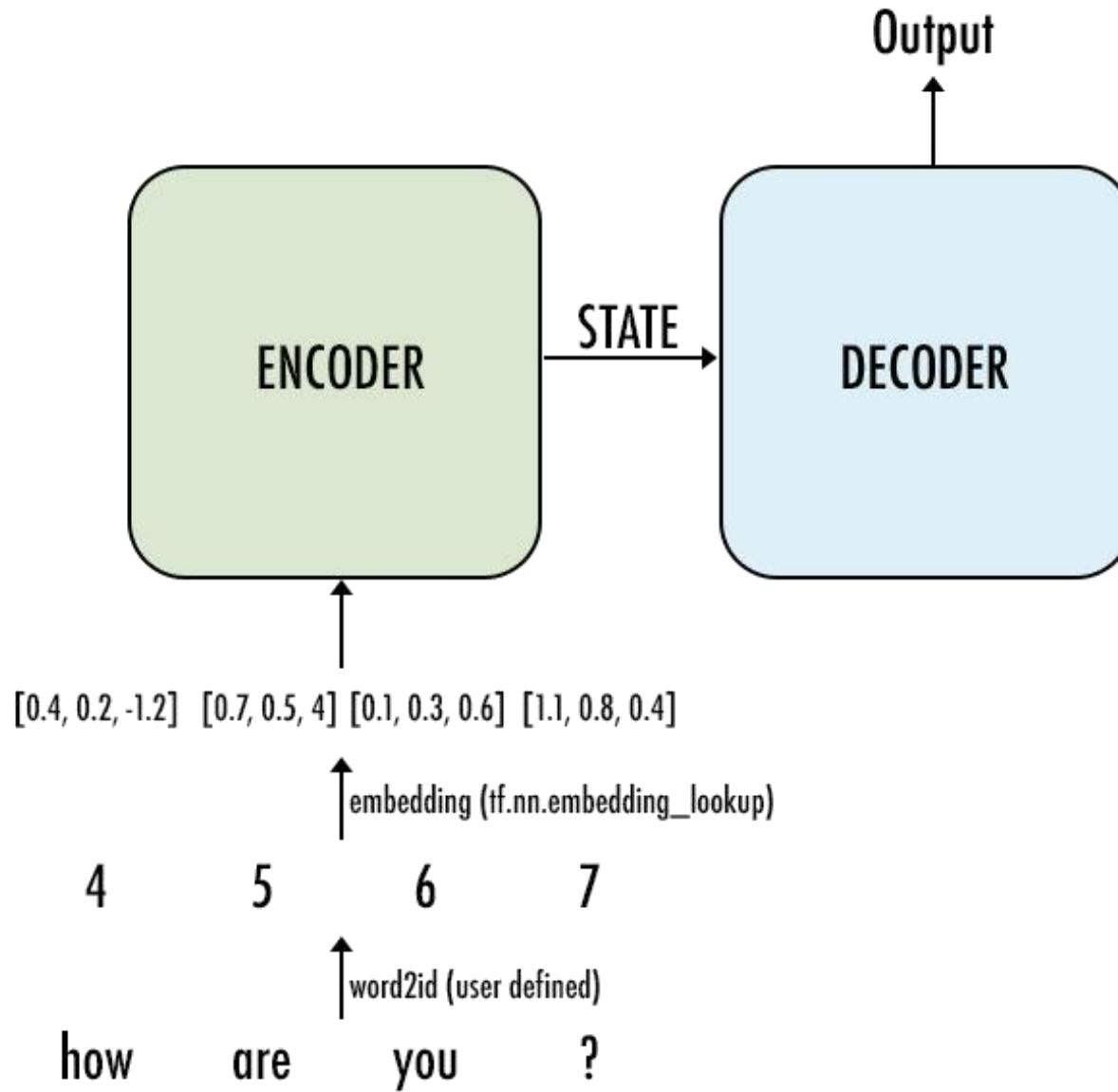
Not yet

how	are	you	?	<PAD>	<PAD>
can	you	fly	that	thing	?



4	5	6	7	0	0
11	12	13	14	15	7

Seq2Seq案例



Seq2Seq案例

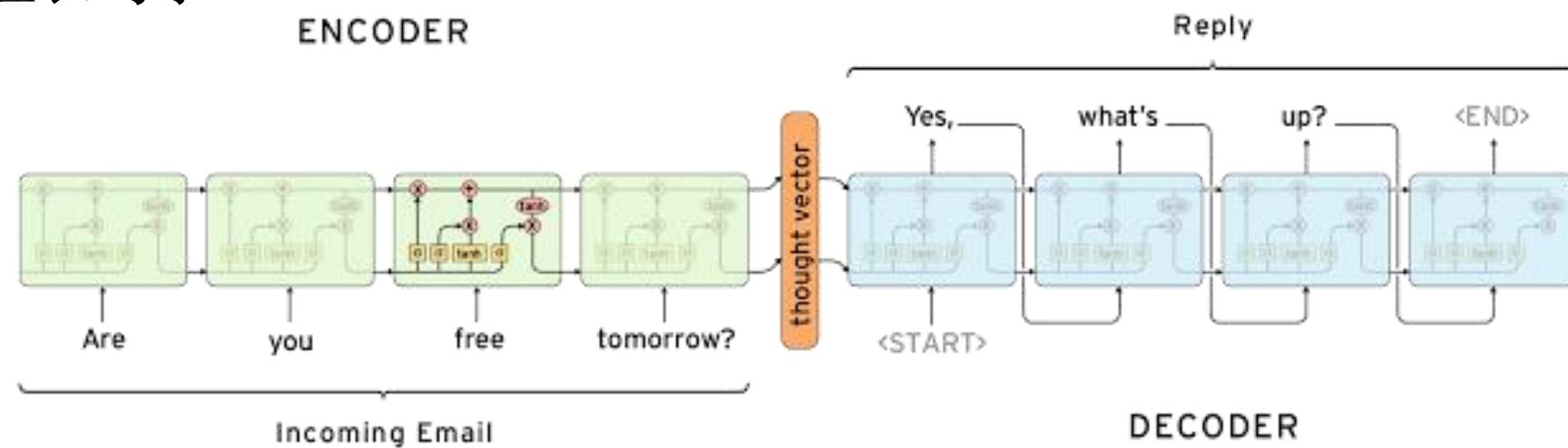


Seq2Seq案例

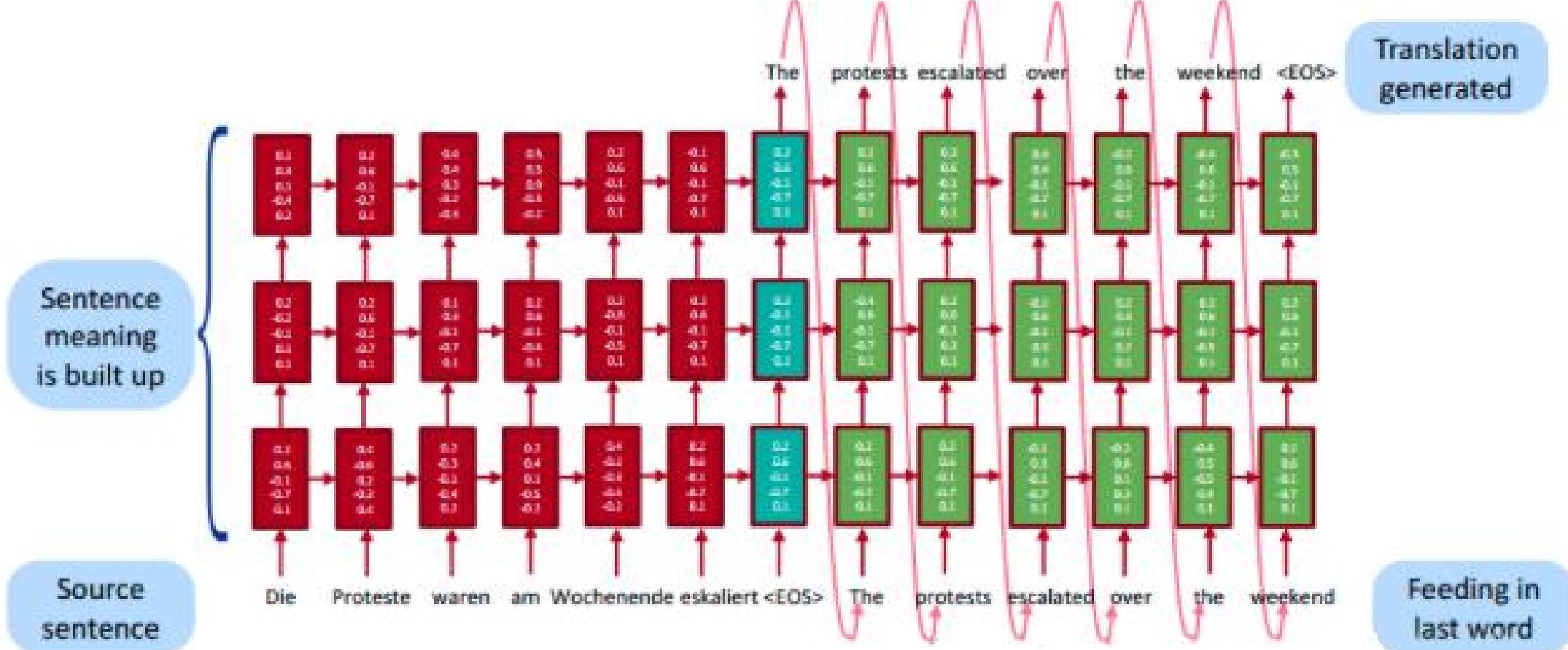


Seq2Seq

- 将RNN模块换成LSTM，则效果如下图。Encoder 和 Decoder 都是 4 个时间步长的 LSTM(但是只有两个RNN Cell)。小技巧：将源句子顺序颠倒后再输入 Encoder 中，比如源句子为“ABC”，那么输入 Encoder 的顺序为“CBA”，经过这样的处理后，取得了很大的提升，而且这样的处理使得模型能够很好地处理长句子。

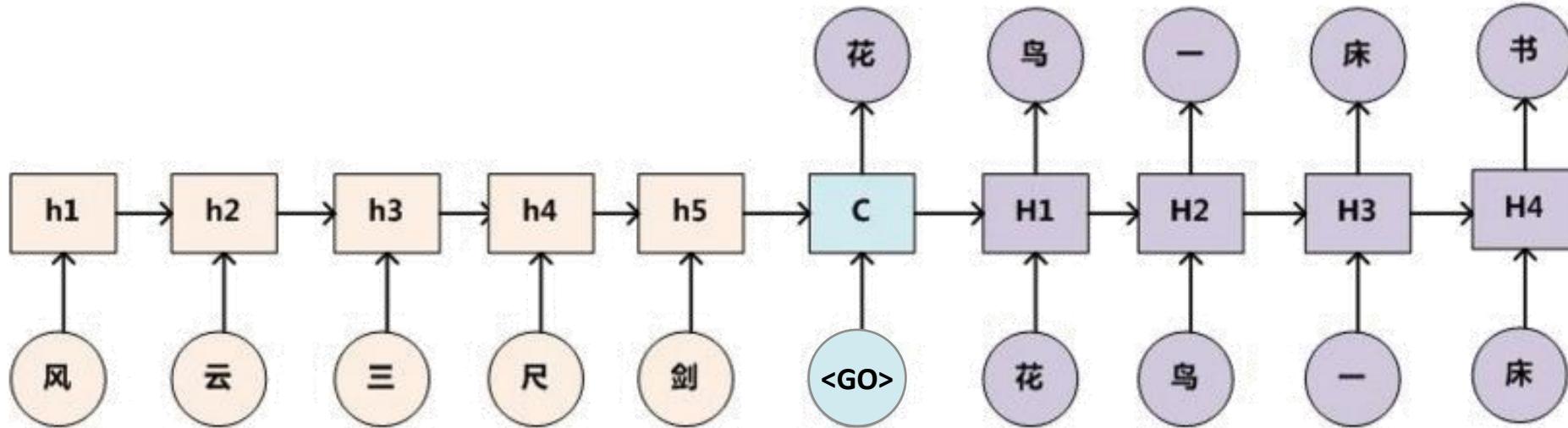


Seq2Seq



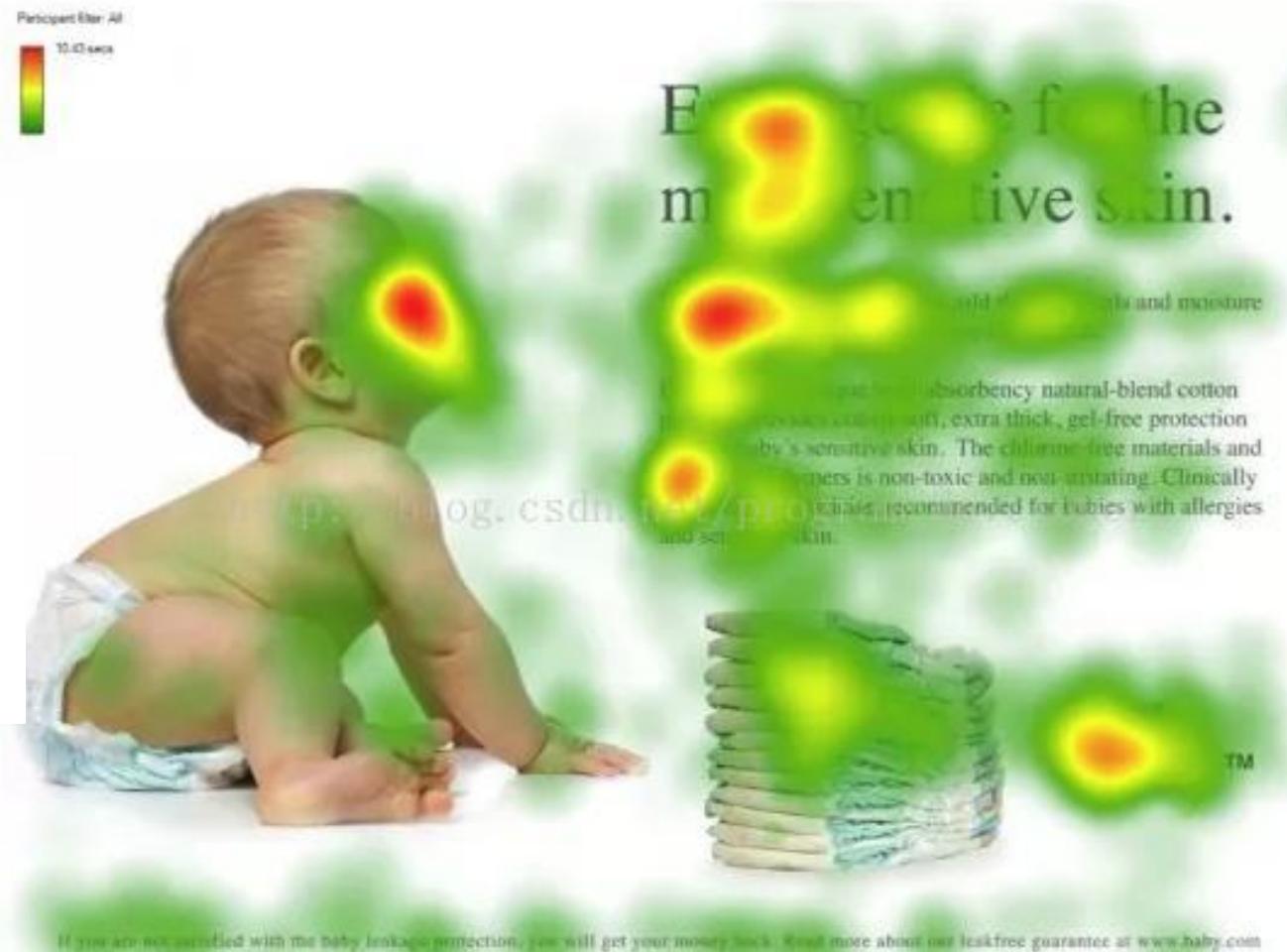
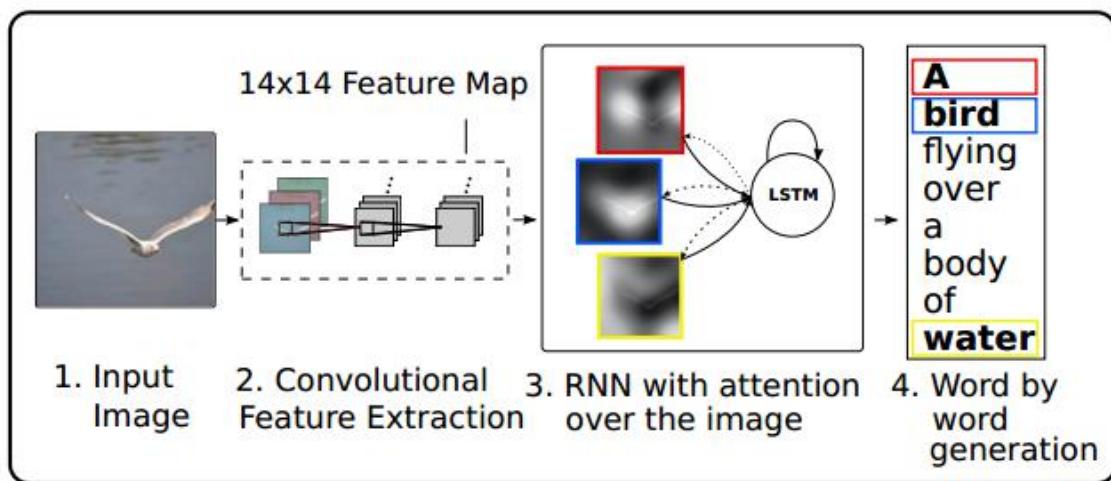
Seq2Seq问题

- 问题描述：“风”对应的特征对于下联的影响是最弱的。



1. 字句对等；
2. 词性对品；
3. 结构对应；
4. 节律对拍；
5. 平仄对立；
6. 形对意联；

Attention



Seq2Seq Attention

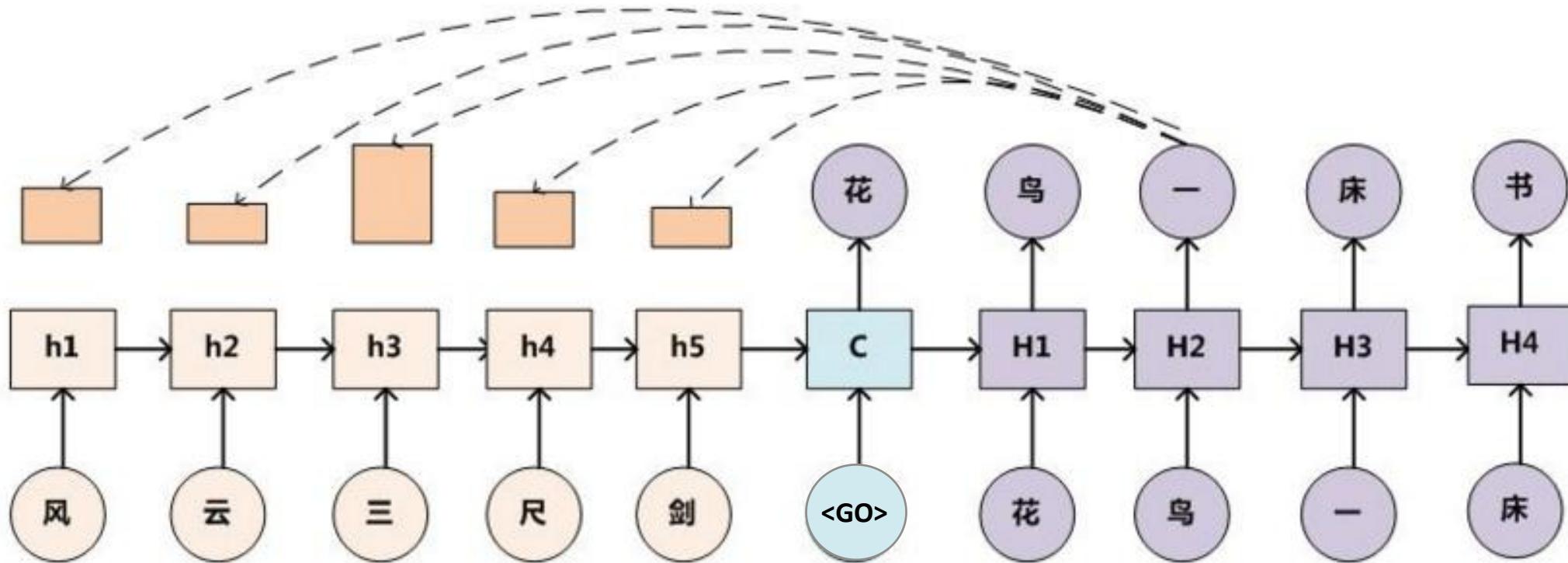
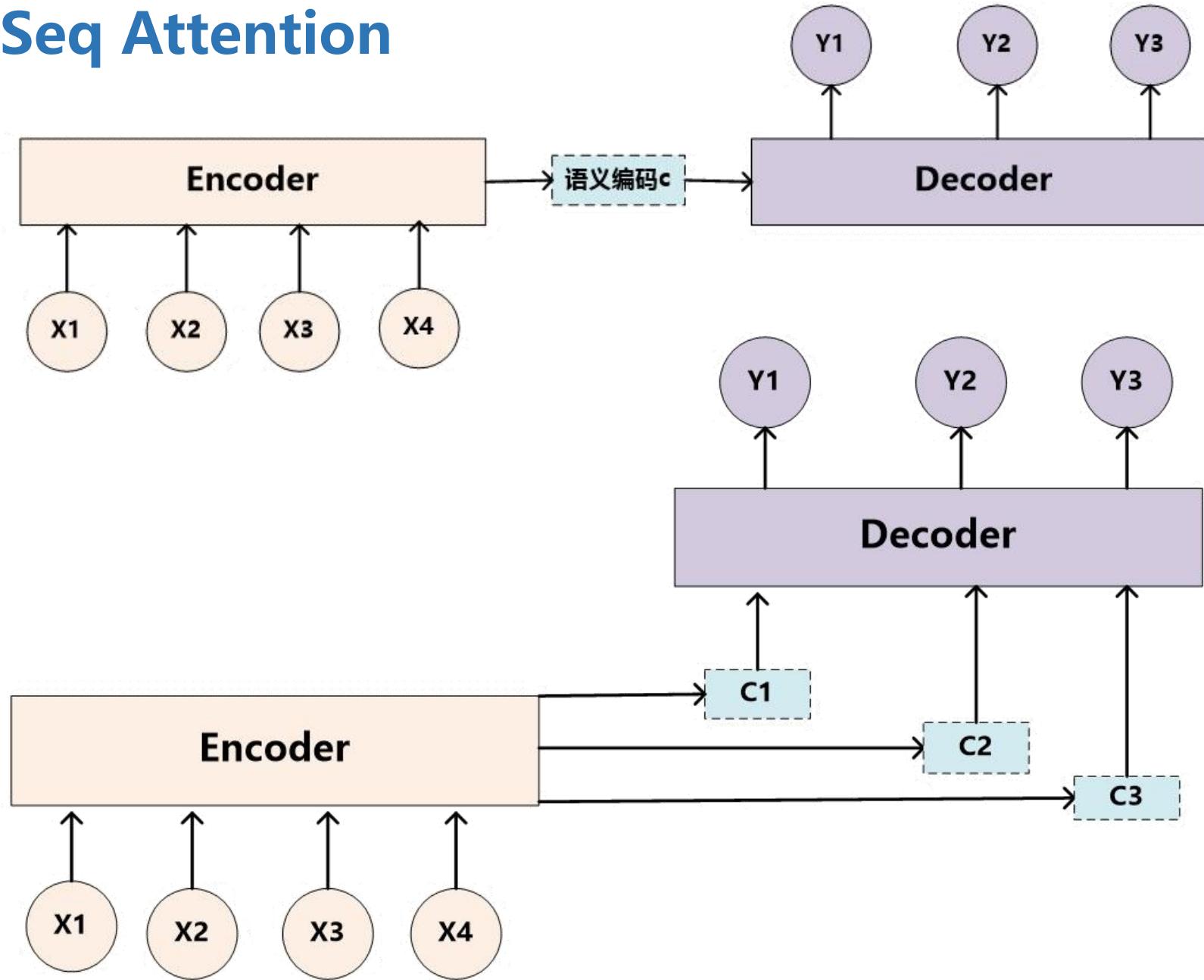
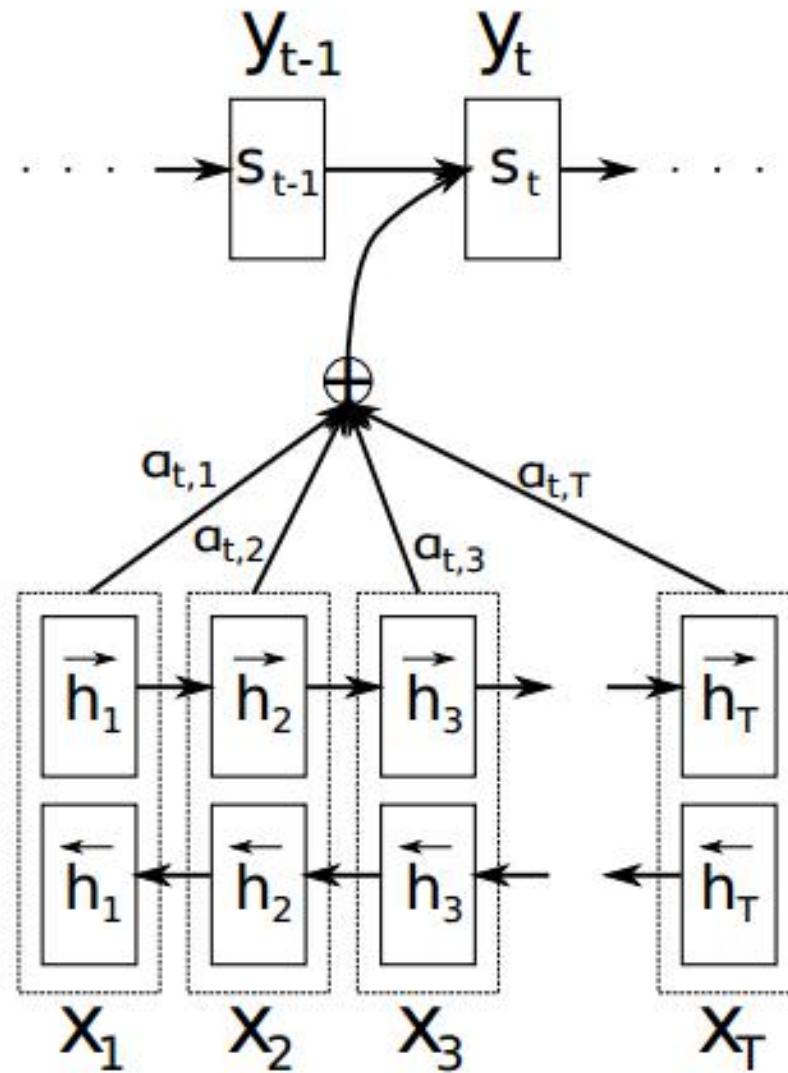


图3. Attention模型

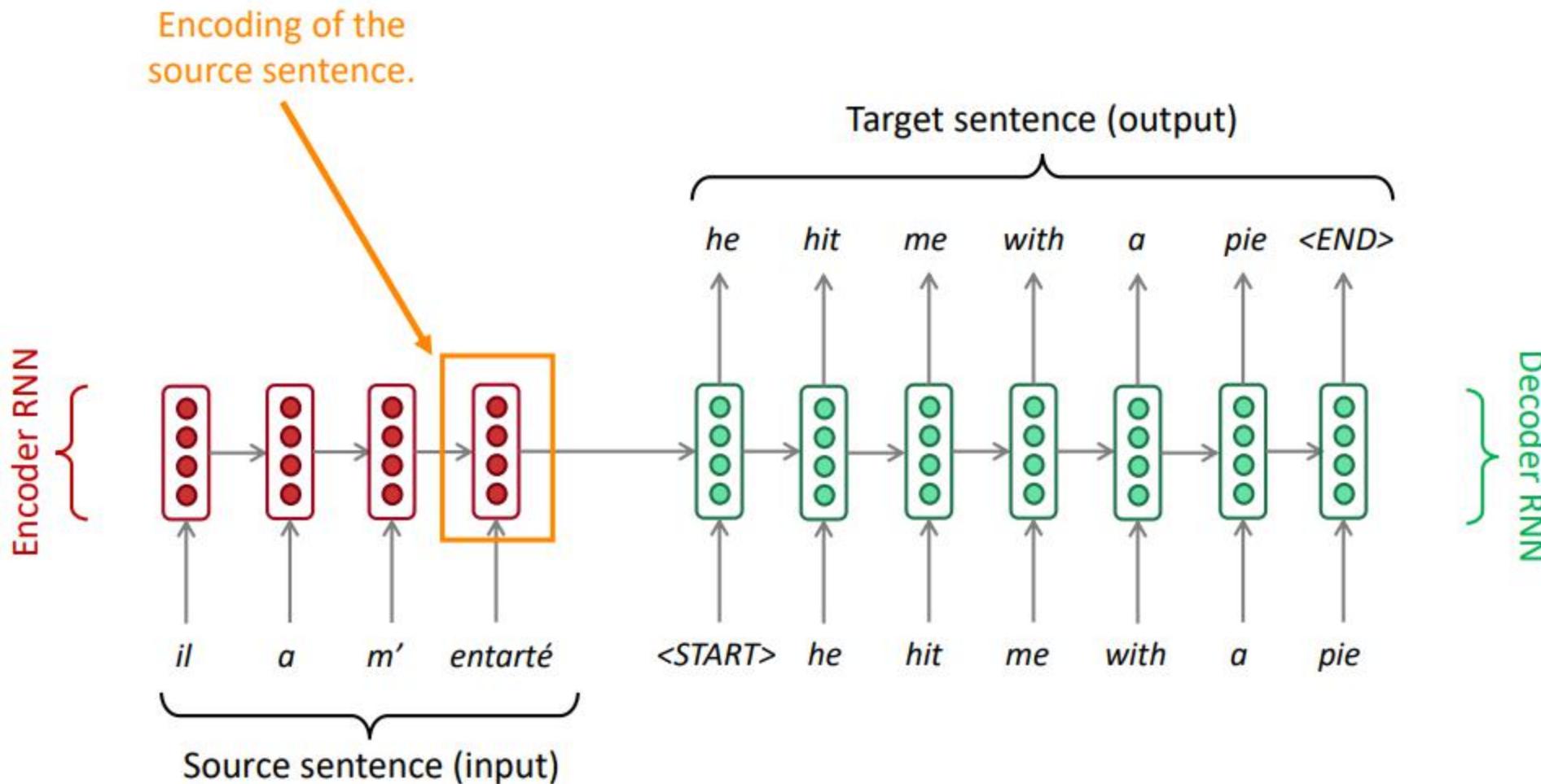
Seq2Seq Attention



Seq2Seq Attention

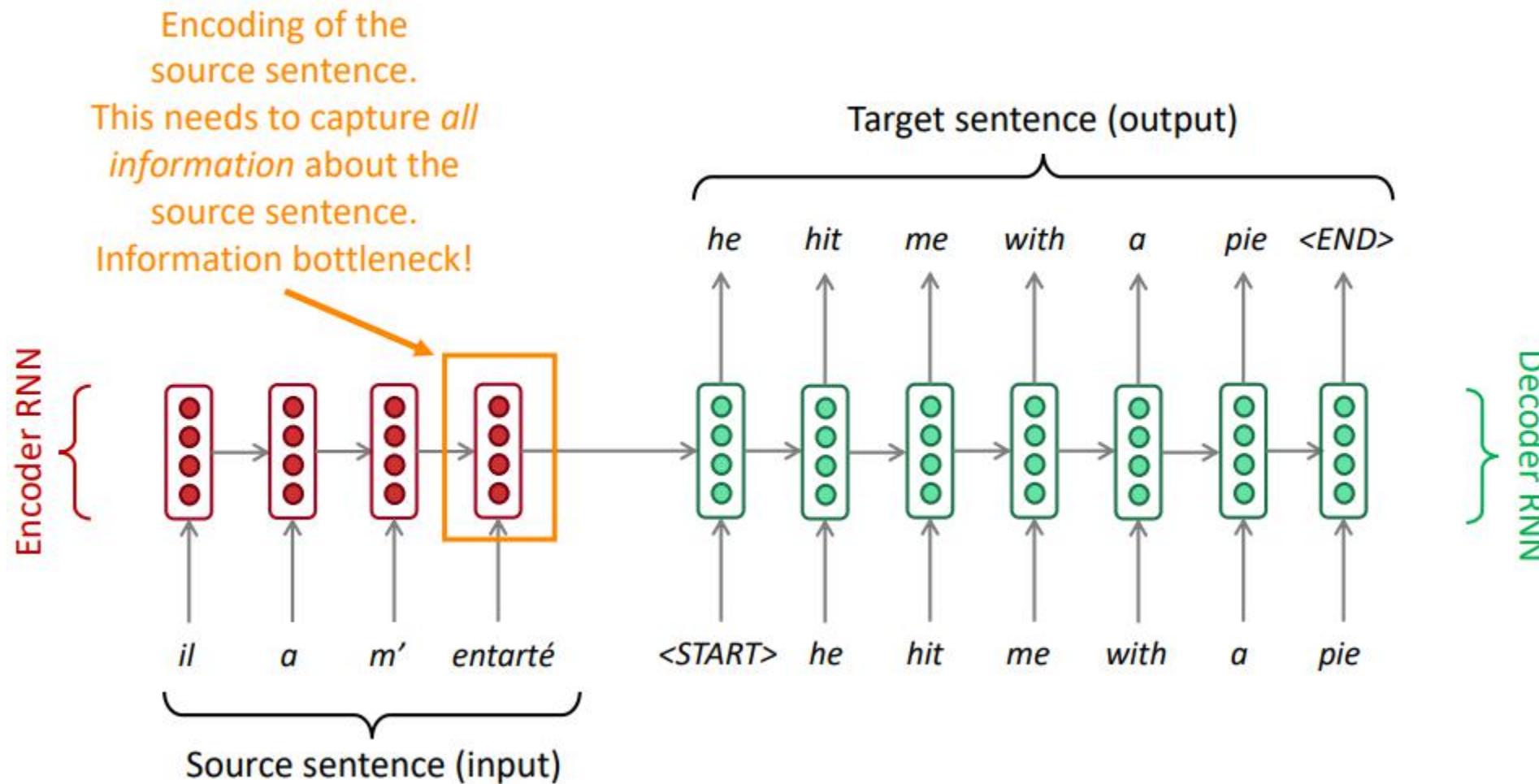


Seq2Seq Attention计算过程

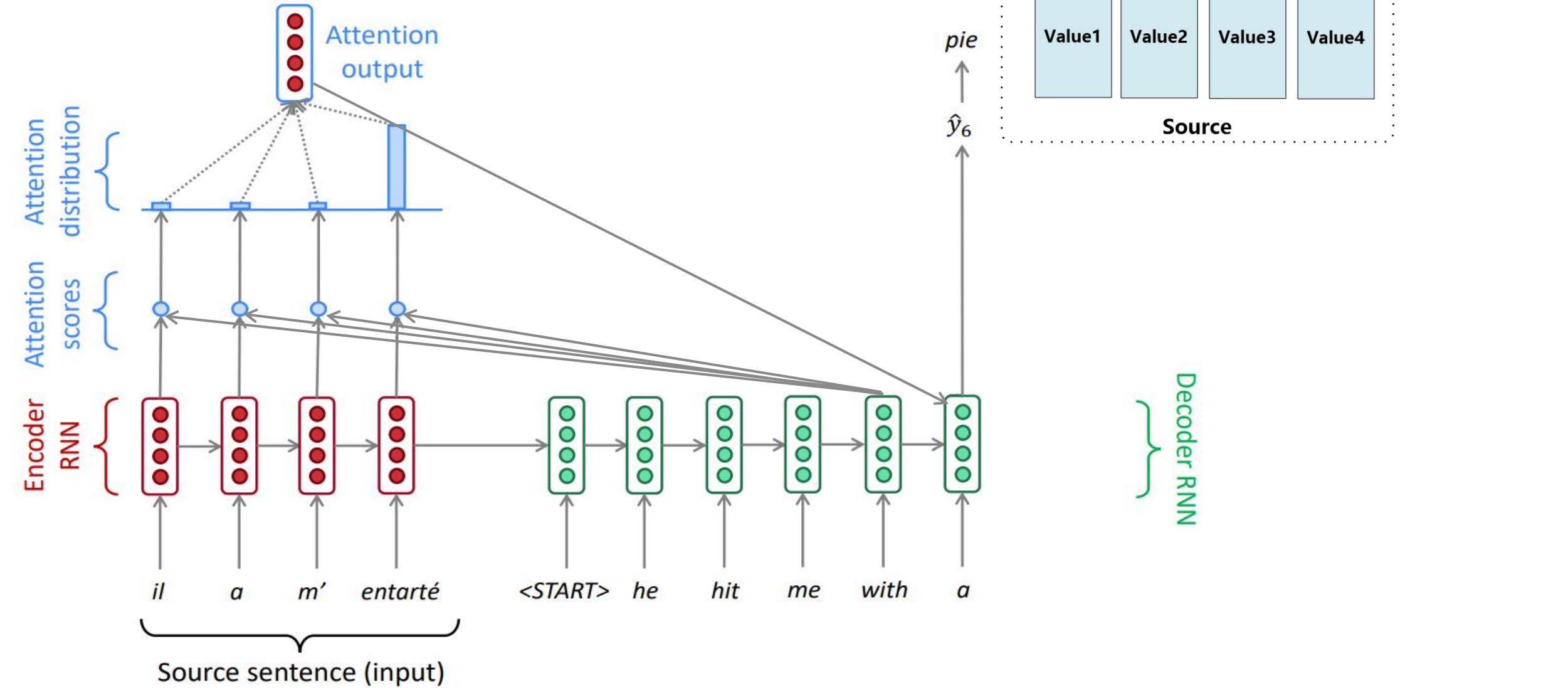


Problems with this architecture?

Seq2Seq Attention计算过程



Seq2Seq Attention计算过程



Seq2Seq Attention计算过程

- Encoder hidden states/output values: h_i ;
- 时刻t, Decoder hidden states: s_t ;
- 基于每一个Encoder状态以及上一个Decoder状态来构建Attention Scores:
$$e_{t,i} = F(h_i, s_{t-1}) \quad e_t = (e_{t,1}, e_{t,2}, \dots, e_{t,n})$$
- 对 e 进行softmax转换, 得到概率分布: $\alpha_t = \text{softmax}(e_t)$
- 基于概率分布以及所有Encoder的状态计算出Attention值: $a_t = \sum_{i=1}^N \alpha_{t,i} h_i$
- 将Decoder当前时刻的输入和Attention值结合, 然后进行普通的RNN操作。

$$y'_t = [y_t; a_t]$$

Seq2Seq Attention计算过程

- Attention Scores的计算函数F在不同论文中有很多形式，主要方式

如下：

$$e_{t,i} = s_{t-1}^T h_i$$

- 乘法Attention: $e_{t,i} = s_{t-1}^T h_i / \sqrt{d}$

$$e_{t,i} = s_{t-1}^T W h_i$$

- 加法Attention:

$$e_{t,i} = u^T \tanh(W_1 h_i + W_2 s_{t-1})$$

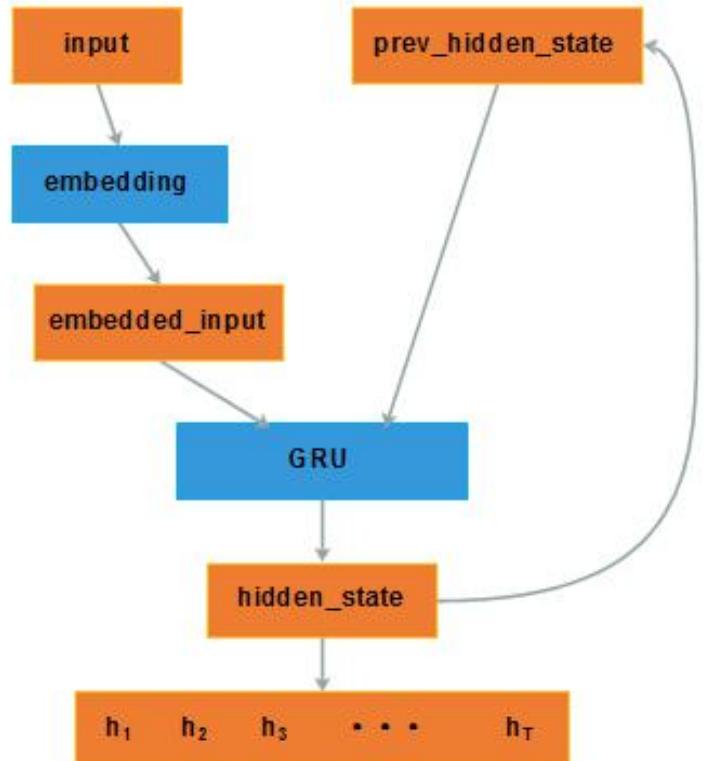
$$e_{t,i} = W_1 h_i + W_2 s_{t-1}$$

$$e_{t,i} = W h_i$$

TensorFlow默认

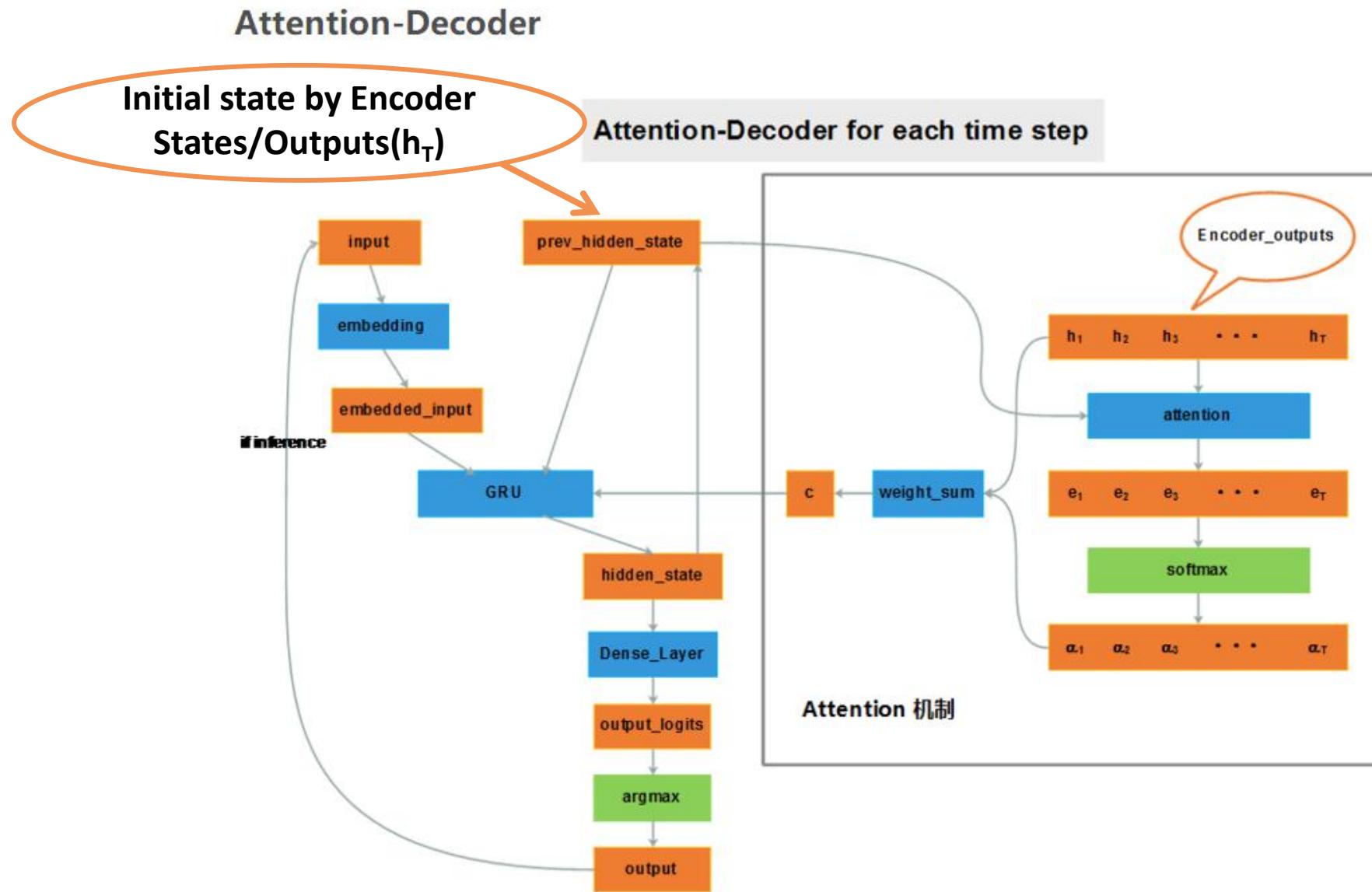
Seq2Seq Attention计算过程

Bi-RNN Encoder



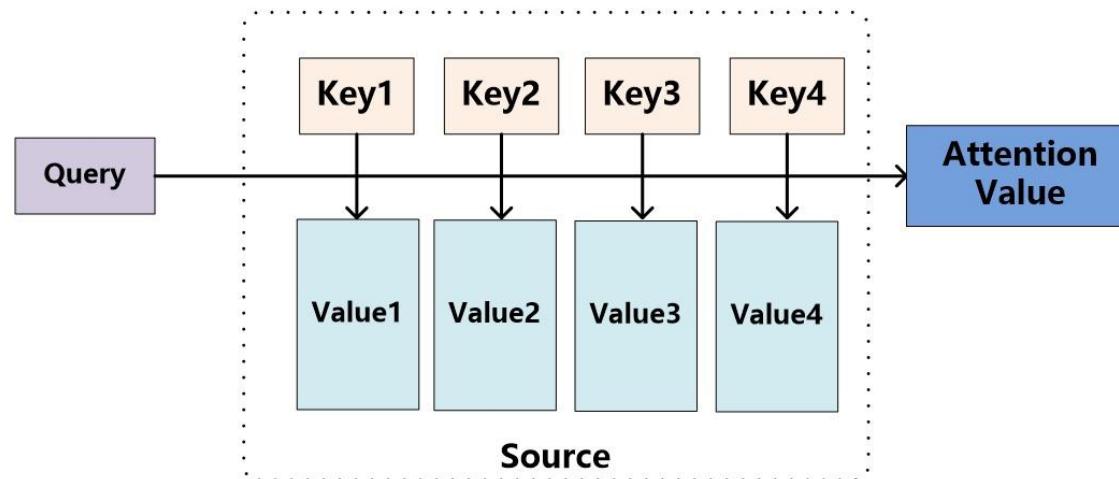
Encoder的流程如上图所示，最终的输出结果是每个时刻的hidden_state $h_1, h_2, h_3, \dots, h_T$ 。

Seq2Seq Attention计算过程



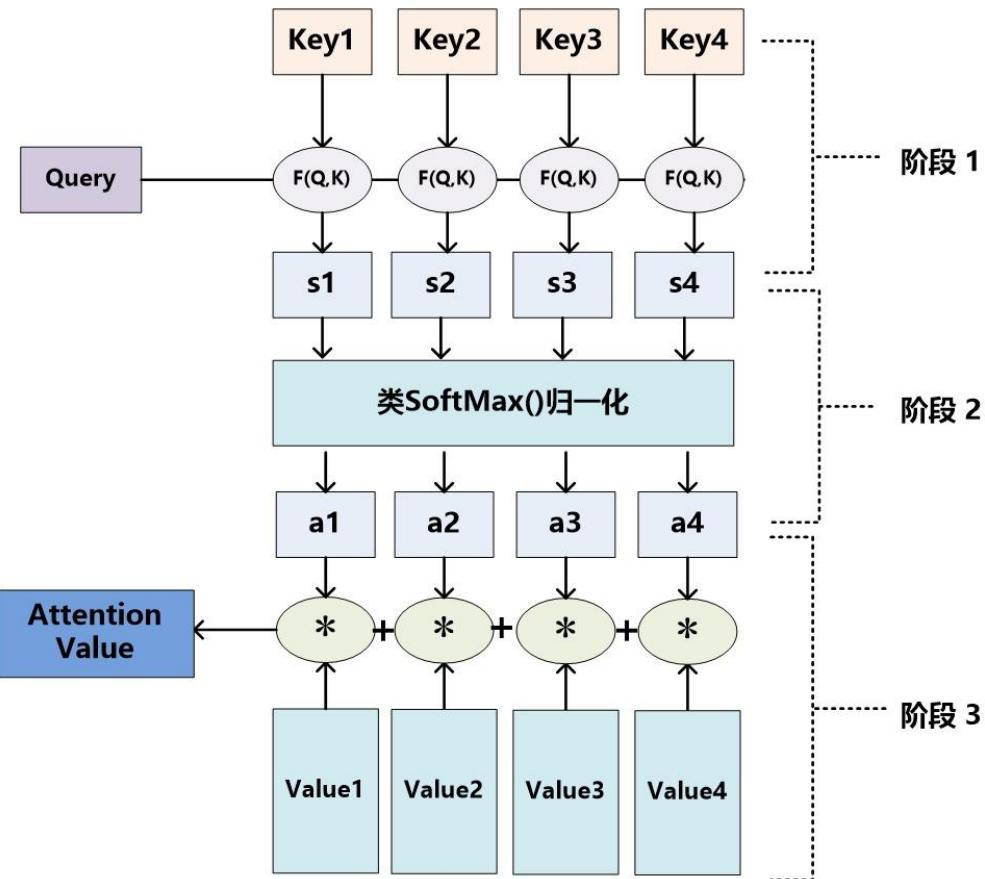
Seq2Seq Attention计算过程(另一种理解方式)

- 此时给定Target中的某个元素Query，通过计算Query和各个Key的相似性或者相关性，得到每个Key对应Value的权重系数，然后对Value进行加权求和，即得到了最终的Attention数值。所以本质上Attention机制是对Source中元素的Value值进行加权求和，而Query和Key用来计算对应Value的权重系数。



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

Seq2Seq Attention计算过程(另一种理解方式)

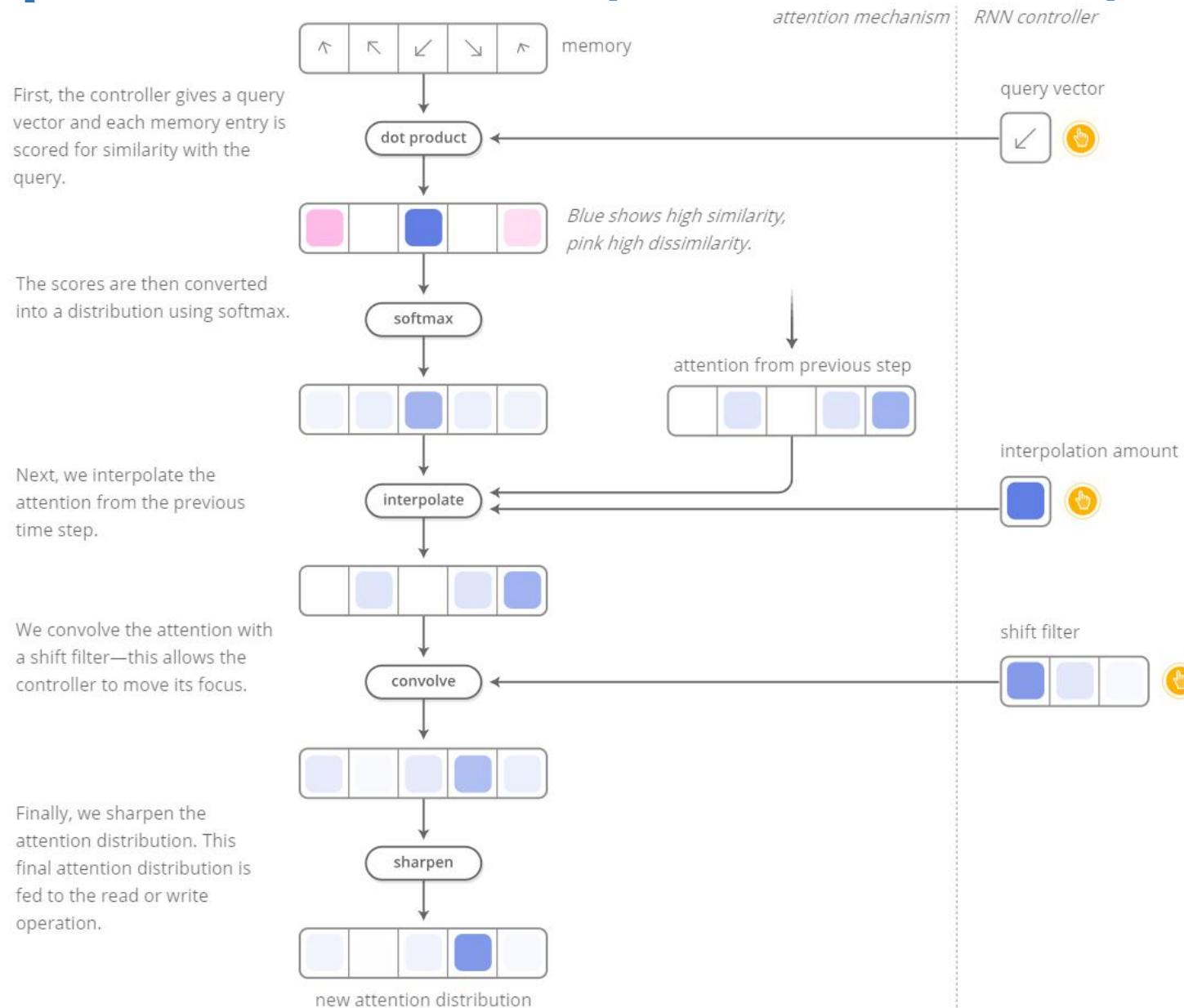


点积: $\text{Similarity}(\text{Query}, \text{Key}_i) = \text{Query} \cdot \text{Key}_i$

Cosine 相似性: $\text{Similarity}(\text{Query}, \text{Key}_i) = \frac{\text{Query} \cdot \text{Key}_i}{\|\text{Query}\| \cdot \|\text{Key}_i\|}$

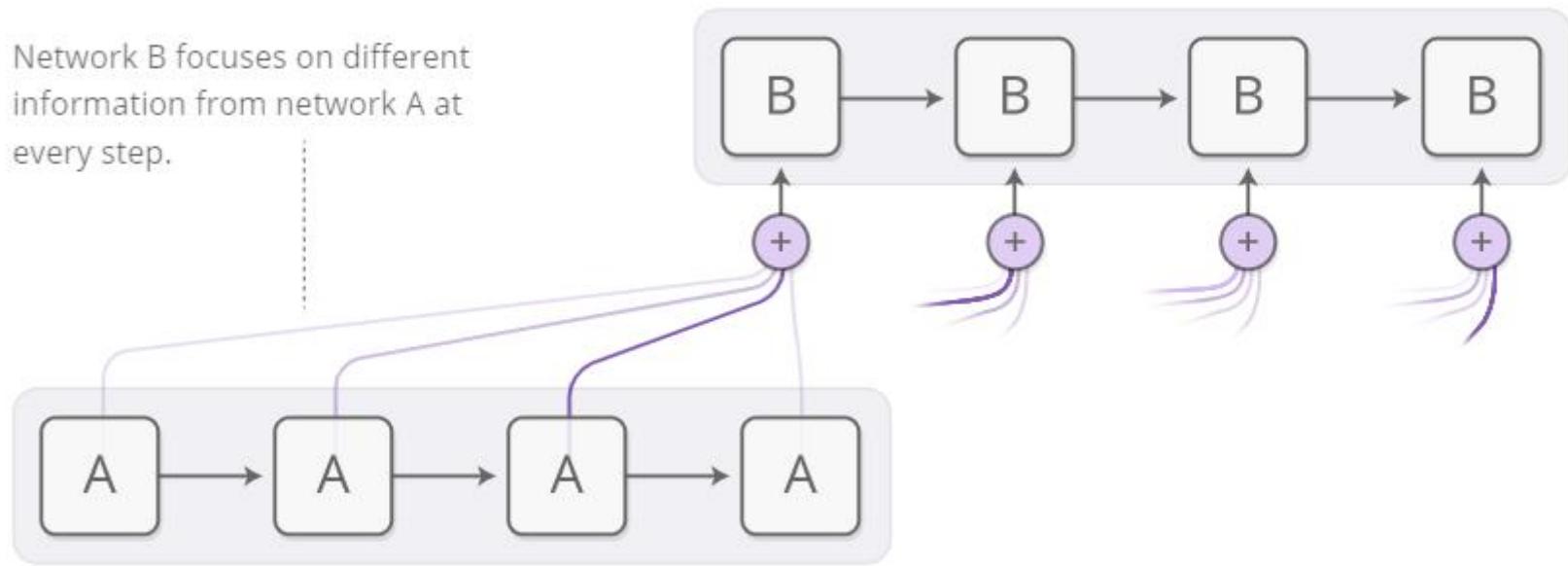
MLP 网络: $\text{Similarity}(\text{Query}, \text{Key}_i) = \text{MLP}(\text{Query}, \text{Key}_i)$

Seq2Seq Attention计算过程(另另一种理解方式)

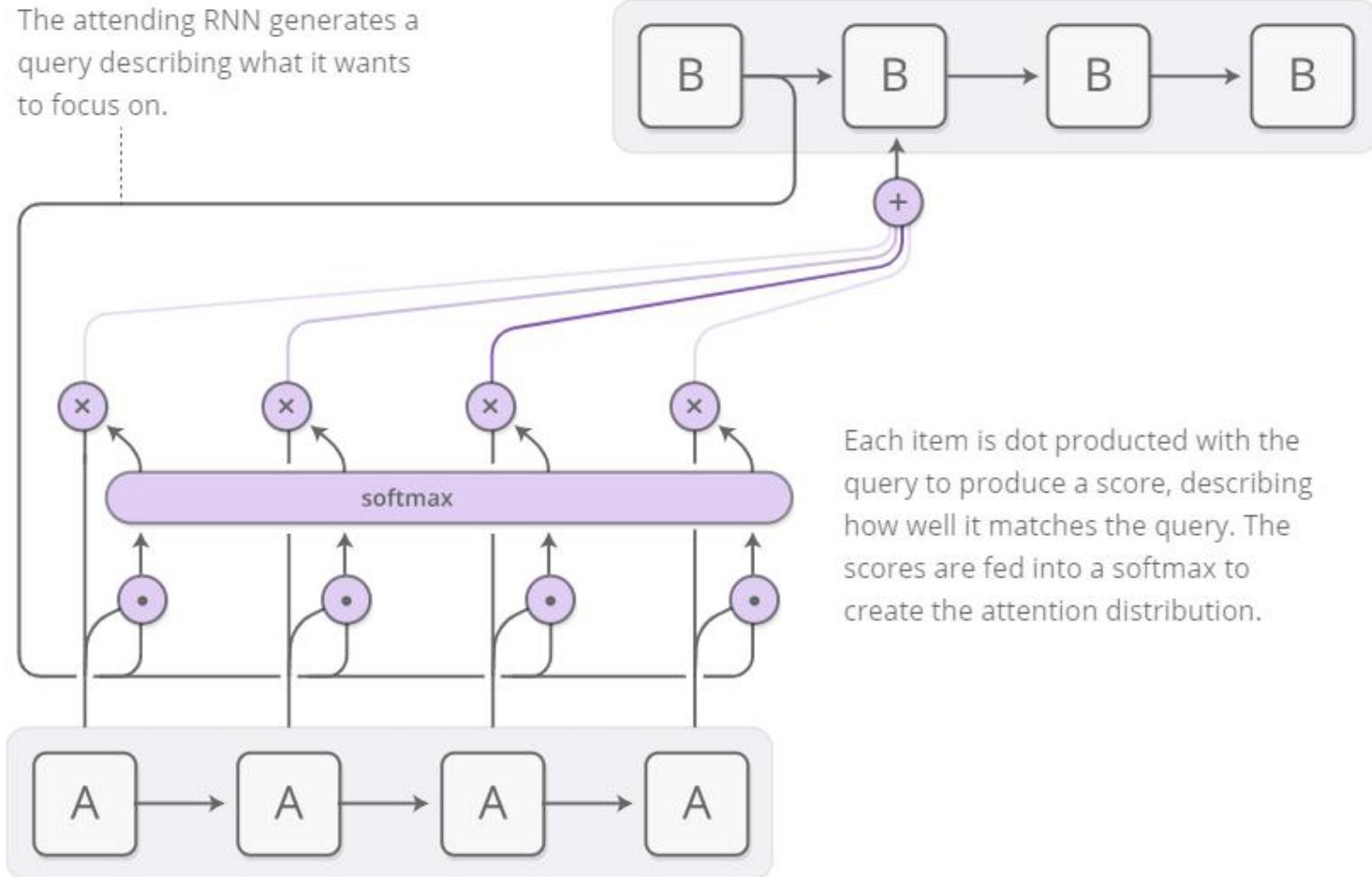


Seq2Seq Attention计算过程(另另一种理解方式)

Network B focuses on different information from network A at every step.



Seq2Seq Attention计算过程(另另一种理解方式)



Seq2Seq Attention效果

L'accord sur la zone économique européenne a été signé en août 1992.

The agreement on the European Economic Area was signed in August 1992.

(a)

Il convient de noter que l'environnement marin est le moins connu de l'environnement.

It should be noted that the marine environment is the least known of environments.

(b)

La destruction de l'équipement signifie que la Syrie ne peut plus produire de nouvelles armes chimiques.

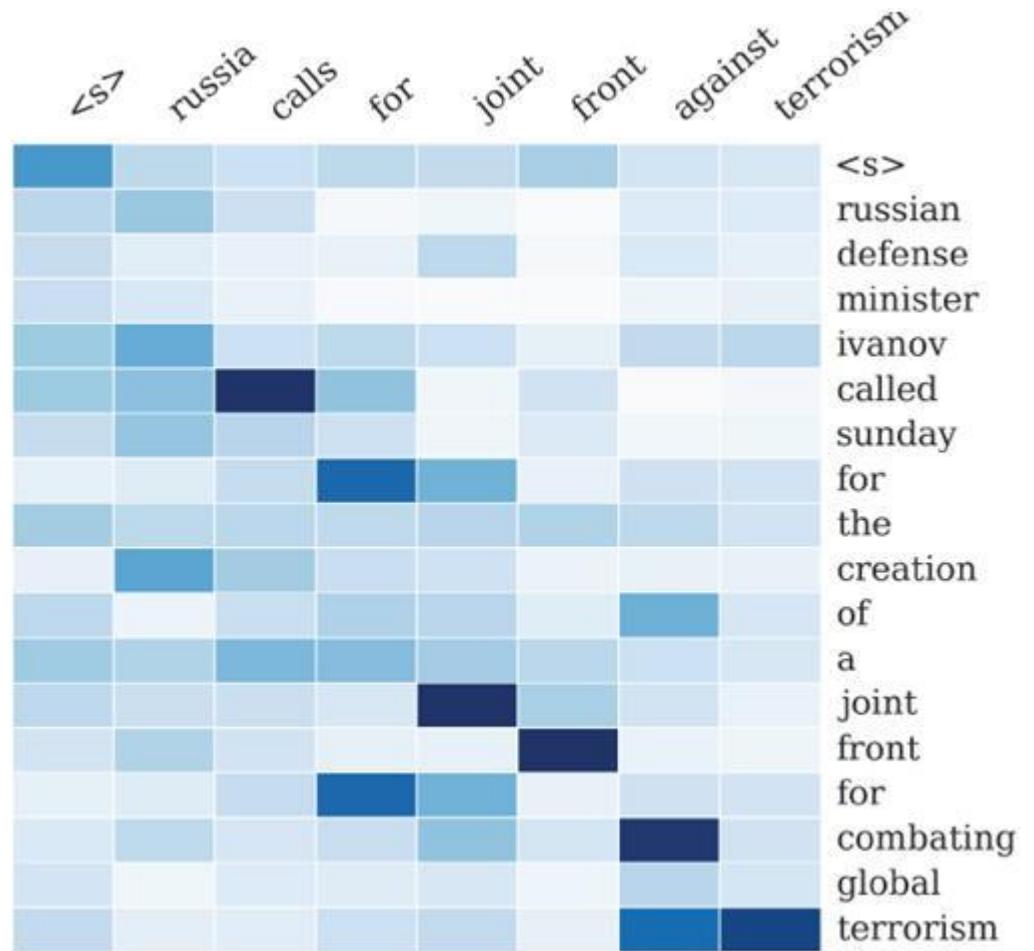
Destruction of the equipment means that Syria can no longer produce new chemical weapons.

(c)

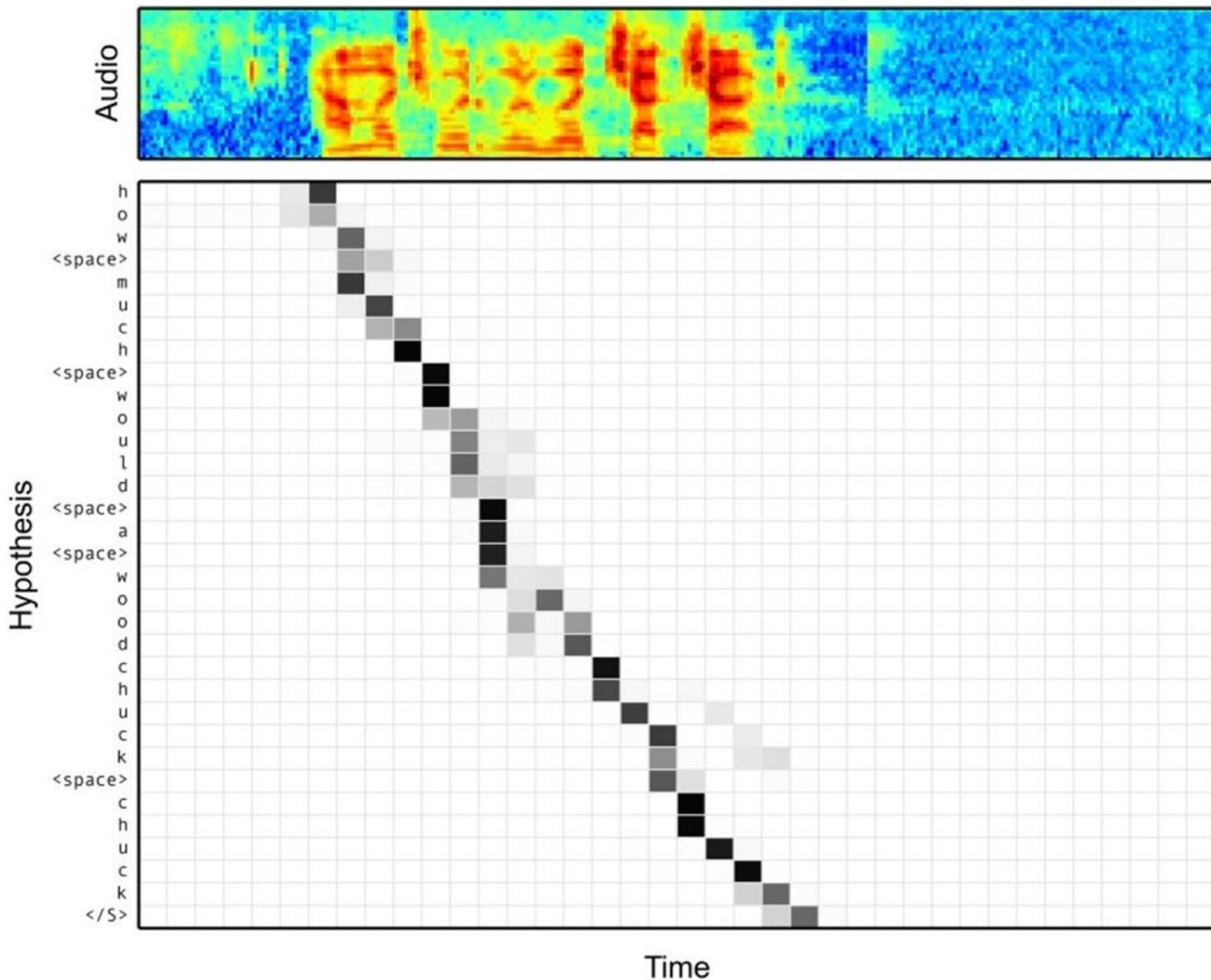
Cela va changer mon avenir avec ma famille. " , a dit l'homme.

"This will change my future with my family . " the man said .

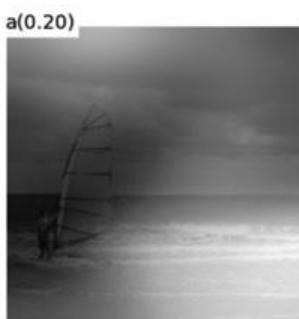
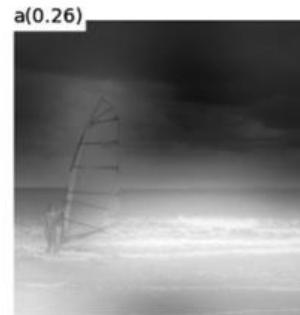
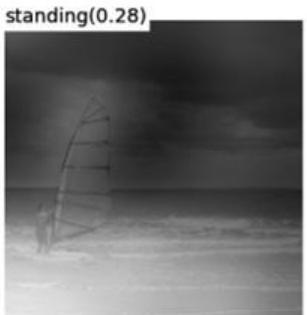
(d)



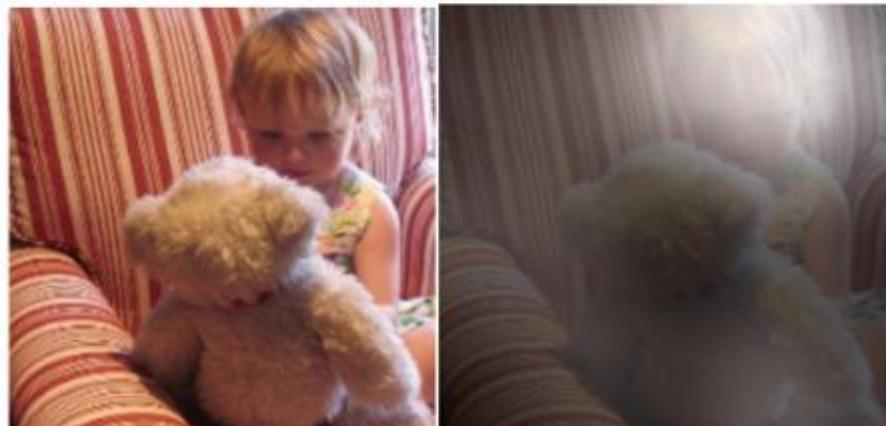
Seq2Seq Attention效果



Seq2Seq Attention效果

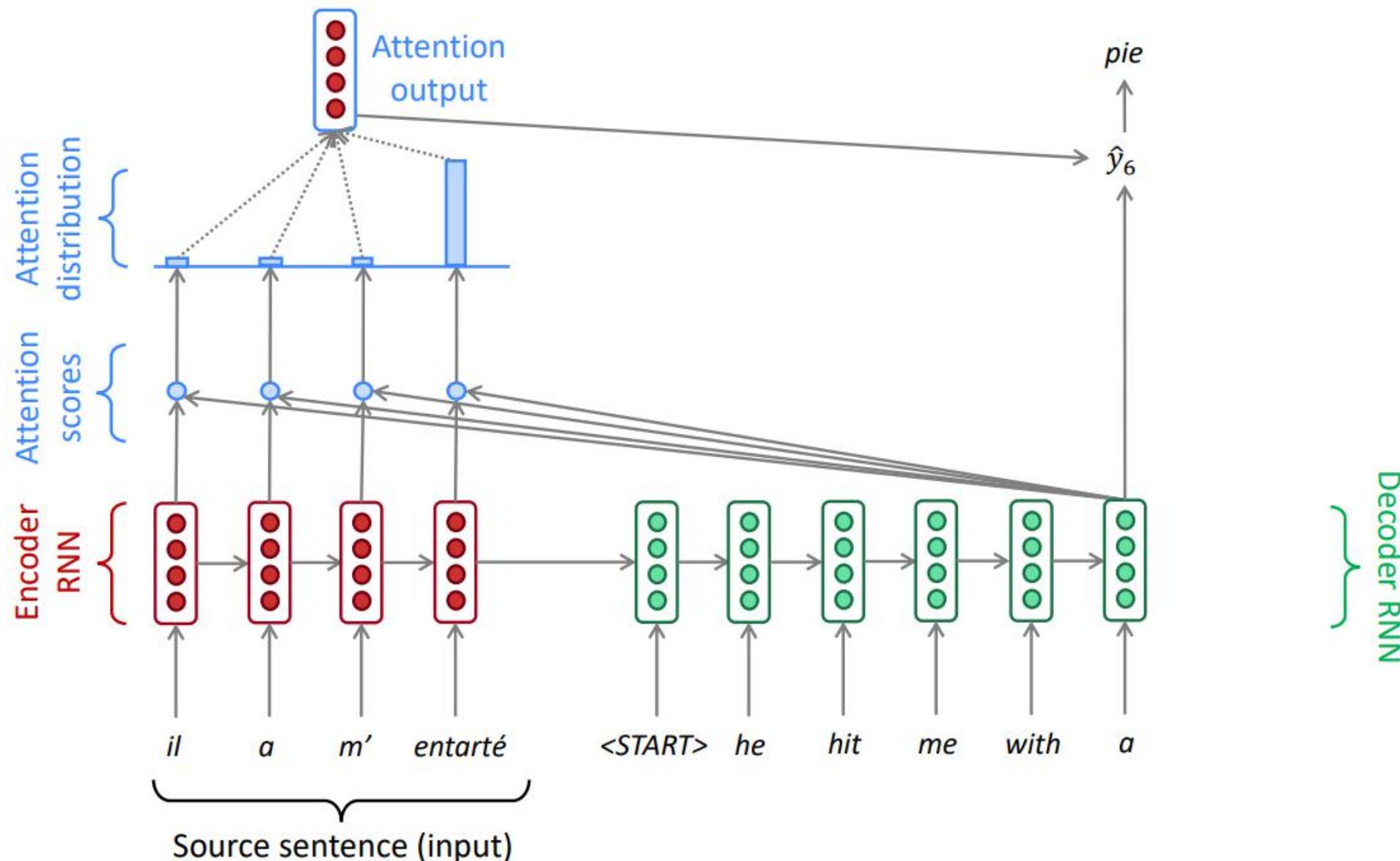


A woman is throwing a frisbee in a park.



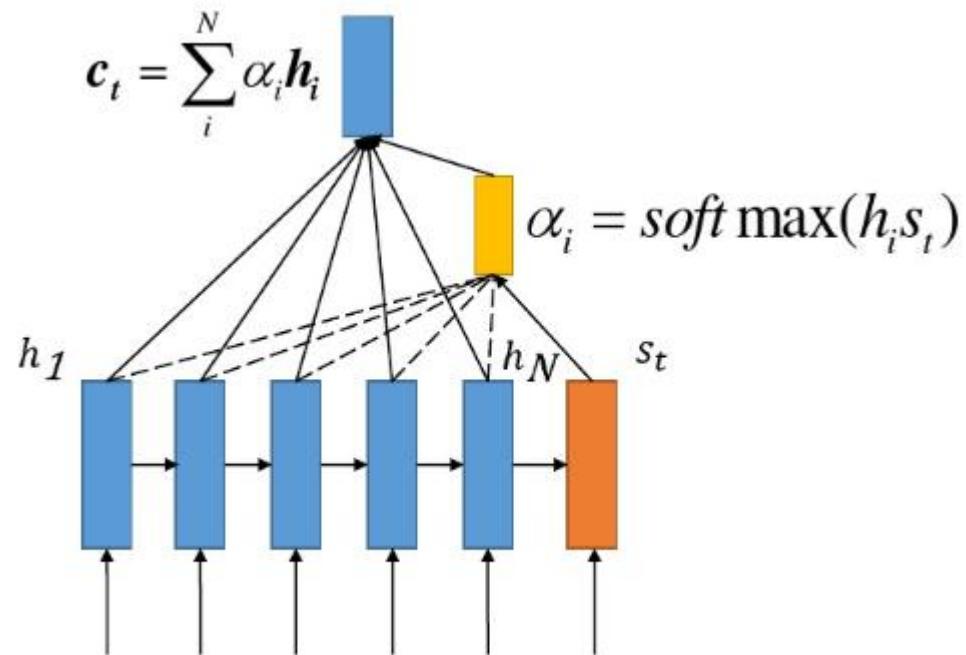
A little girl sitting on a bed with a teddy bear.

Seq2Seq Attention常规形状



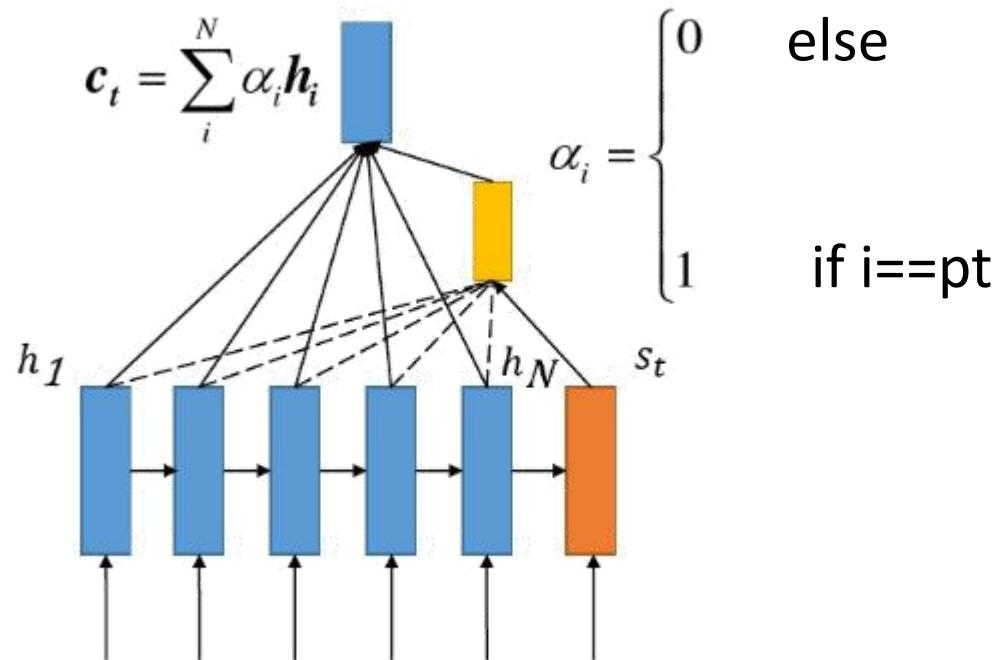
Seq2Seq Attention Soft Attention

- 15年被提出于《Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, Kelvin Xu》



Seq2Seq Attention Hard Attention

- 和Soft Attention在15年同时在同一篇论文中被提出；
- Soft Attention中是对于每个Encoder的Hidden State会match一个概率值，而在Hard Attention会直接找一个特定的单词概率为1，而其它对应概率为0.



Seq2Seq Attention Global Attention

- 在15年被提出于《Effective Approaches to Attention-based Neural Machine Translation, Minh-Thang Luong》，和Soft Attention类似。

$$\mathbf{h}_j = f(\mathbf{h}_{j-1}, \mathbf{s})$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t])$$

$$p(y_t | y_{<t}, x) = \text{softmax}(\mathbf{W}_s \tilde{\mathbf{h}}_t)$$

$$\begin{aligned} \mathbf{a}_t(s) &= \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \\ &= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \end{aligned}$$

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & dot \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & general \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s]) & concat \end{cases}$$

$$\mathbf{a}_t = \text{softmax}(\mathbf{W}_a \mathbf{h}_t) \quad location$$

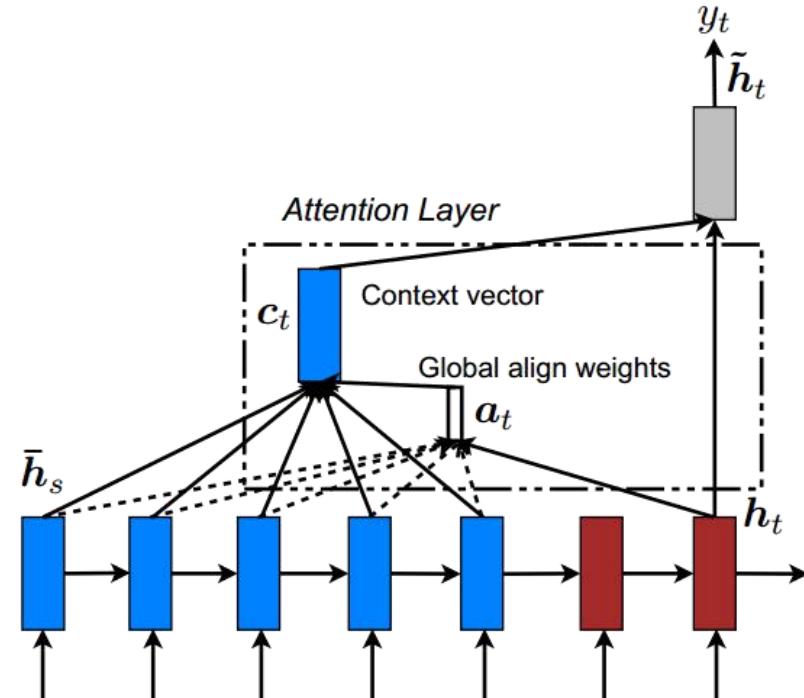
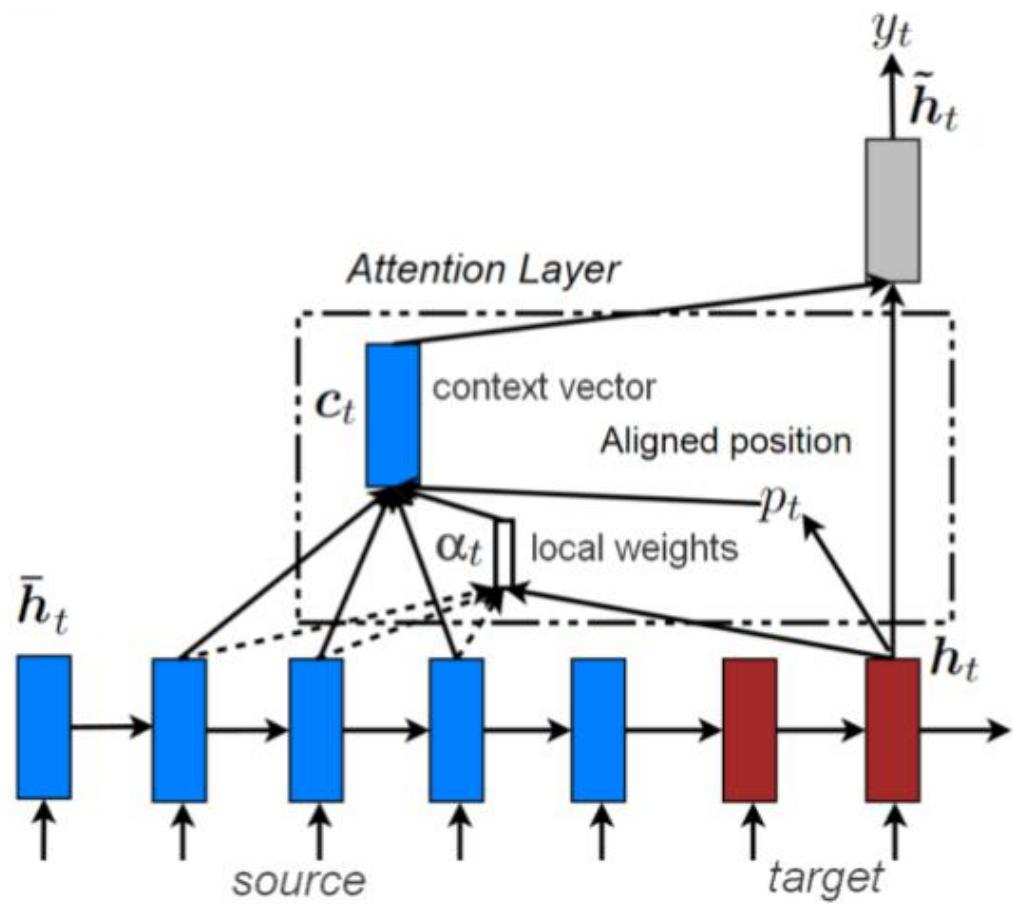


Figure 2: **Global attentional model** – at each time step t , the model infers a *variable-length* alignment weight vector \mathbf{a}_t based on the current target state \mathbf{h}_t and all source states $\bar{\mathbf{h}}_s$. A global context vector \mathbf{c}_t is then computed as the weighted average, according to \mathbf{a}_t , over all the source states.

Seq2Seq Attention Local Attention

- 和Global Attention在同一篇论文中被提出；相当于Soft Attention和Hard Attention中间状态(半硬半软Attention)
- 对于时刻t的词汇，模型首先产生一个对齐位置 pt (aligned position)， context vector(c)由编码器中的隐状态计算得到，编码器的隐状态不是所有的隐状态，而是在区间 $[pt-D, pt+D]$ 中， D 的大小由经验给定。

Seq2Seq Attention Local Attention



$$P[y_t | \{y_1, \dots, y_{t-1}\}, c_t] = \text{softmax}(W_s \tilde{h}_t)$$

attentional hidden state

$$\tilde{h}_t = \tanh(W_c [c_t; h_t])$$

decoder hidden state

context vector

$$p_t = T_x \cdot \sigma(v_p^\top \tanh(W_p \tilde{h}_t))$$

i^{th} encoder hidden state

$$c_t = \sum_{i=p_t-D}^{p_t+D} \alpha_{t,i} \bar{h}_i$$

alignment vector

$$\alpha_{t,i} = \frac{\exp(\text{score}(h_t, \bar{h}_i))}{\sum_{i'=p_t-D}^{p_t+D} \exp(\text{score}(h_t, \bar{h}_{i'}))} \exp\left(-\frac{(i - p_t)^2}{2(D/2)^2}\right)$$

p_t

$-D/2$ $D/2$

$$\text{score}(h_t, \bar{h}_i) = h_t^\top W_\alpha \bar{h}_i$$

Seq2Seq Attention Self Attention

- 在17年被提出于《Attention Is All You Need, Ashish Vaswani》, 也称为Transformer结构; 内部包含Multi-Head Attention以及Residual残差结构。
- Transformer是Bert网络结构的基础。

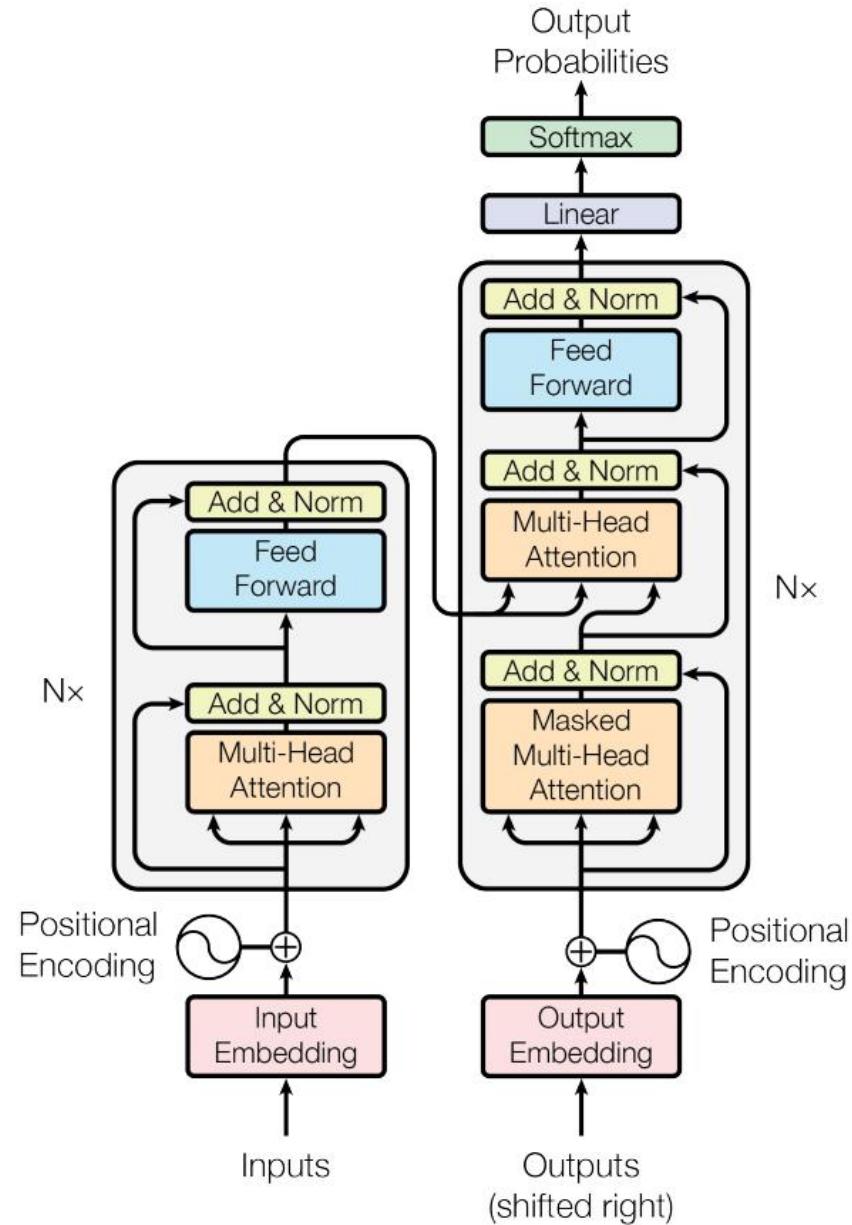
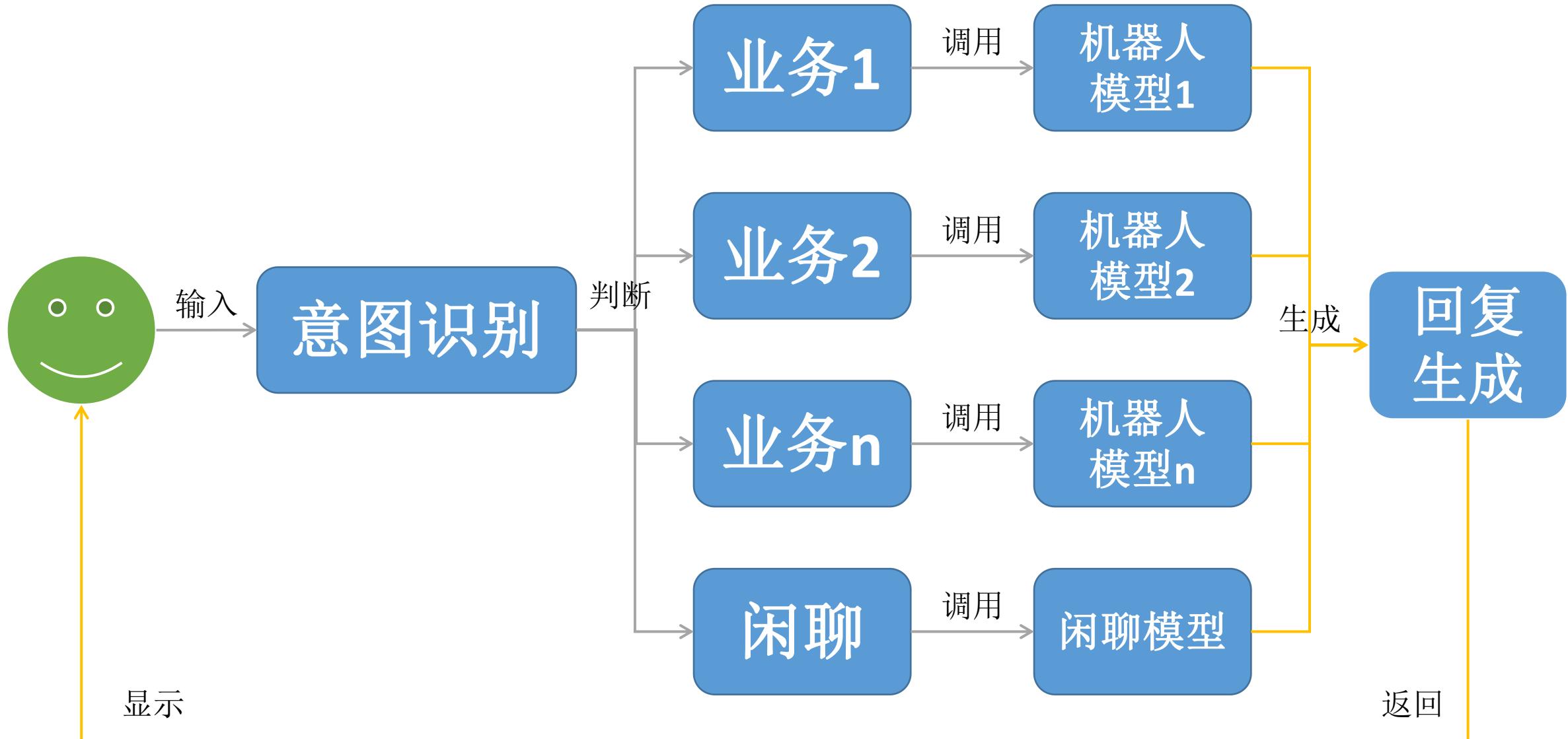


Figure 1: The Transformer - model architecture.

Seq2Seq+Attention项目_聊天机器人



Seq2Seq Attention TensorFlow实现

