



厦门大学信息学院 本科选修课

2021-2022 第二学期

模式识别

Pattern Recognition

主讲：王程



第九章

循环神经网络 (RNN)

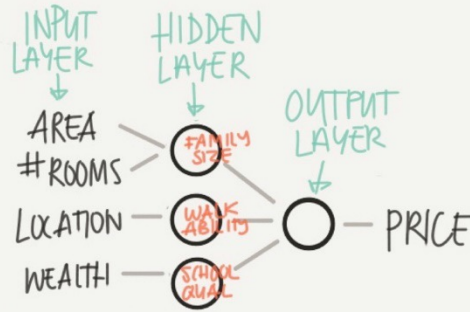
陈龙彪

longbiaochen@xmu.edu.cn

厦门大学信息学院

2022年4月

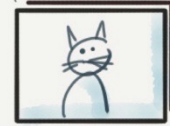
INTRO TO DEEP LEARNING



NNs CAN DEAL WITH BOTH STRUCTURED & UNSTRUCTURED DATA



STRUCTURED



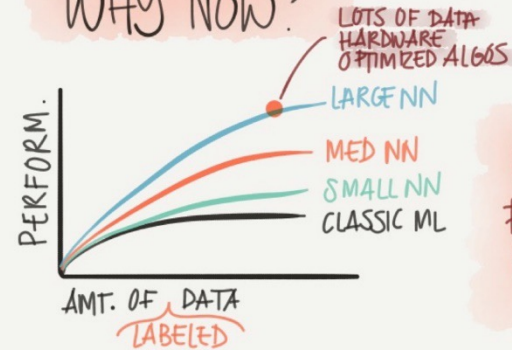
"THE QUICK BROWN FOX"
UNSTRUCTURED

HUMANS ARE GOOD AT THIS

SUPERVISED LEARNING

INPUT: X	OUTPUT: Y	NN TYPE
HOME FEATURES AD+ USER INFO	PRICE WILL CLICK ON AD (0/1)	STANDARD NN
IMAGE	OBJECT (1...1000)	CONV. NN (CNN)
AUDIO ENGLISH	TEXT TRANSCRIPT CHINESE	RECURRENT NN (RNN)
IMAGE/RADAR	POS OF OTHER CARS	CUSTOM/HYBRID

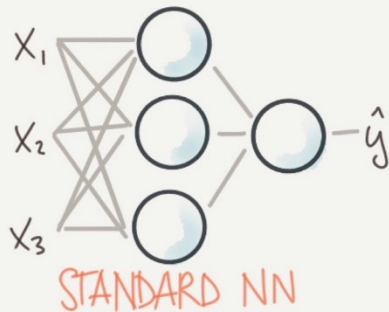
WHY NOW?



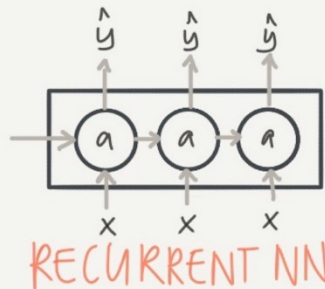
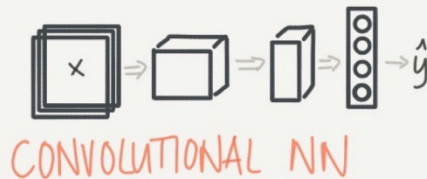
ONE OF THE BIG BREAKTHROUGHS HAS BEEN MOVING FROM SIGMOID TO RELU FOR FASTER GRADIENT DESCENT



FASTER COMPUTATION IS IMPORTANT TO SPEED UP THE ITERATIVE PROCESS



NETWORK ARCHITECTURES



机器学习 \approx 寻找一个函数

- 语音识别

$$f(\text{语音波形}) = \text{"How are you"}$$

- 图片识别

$$f(\text{小猫图片}) = \text{"Cat"}$$

- 下围棋

$$f(\text{围棋棋盘}) = \text{"5-5"} \quad (\text{下一步棋的结果})$$

- 对话系统

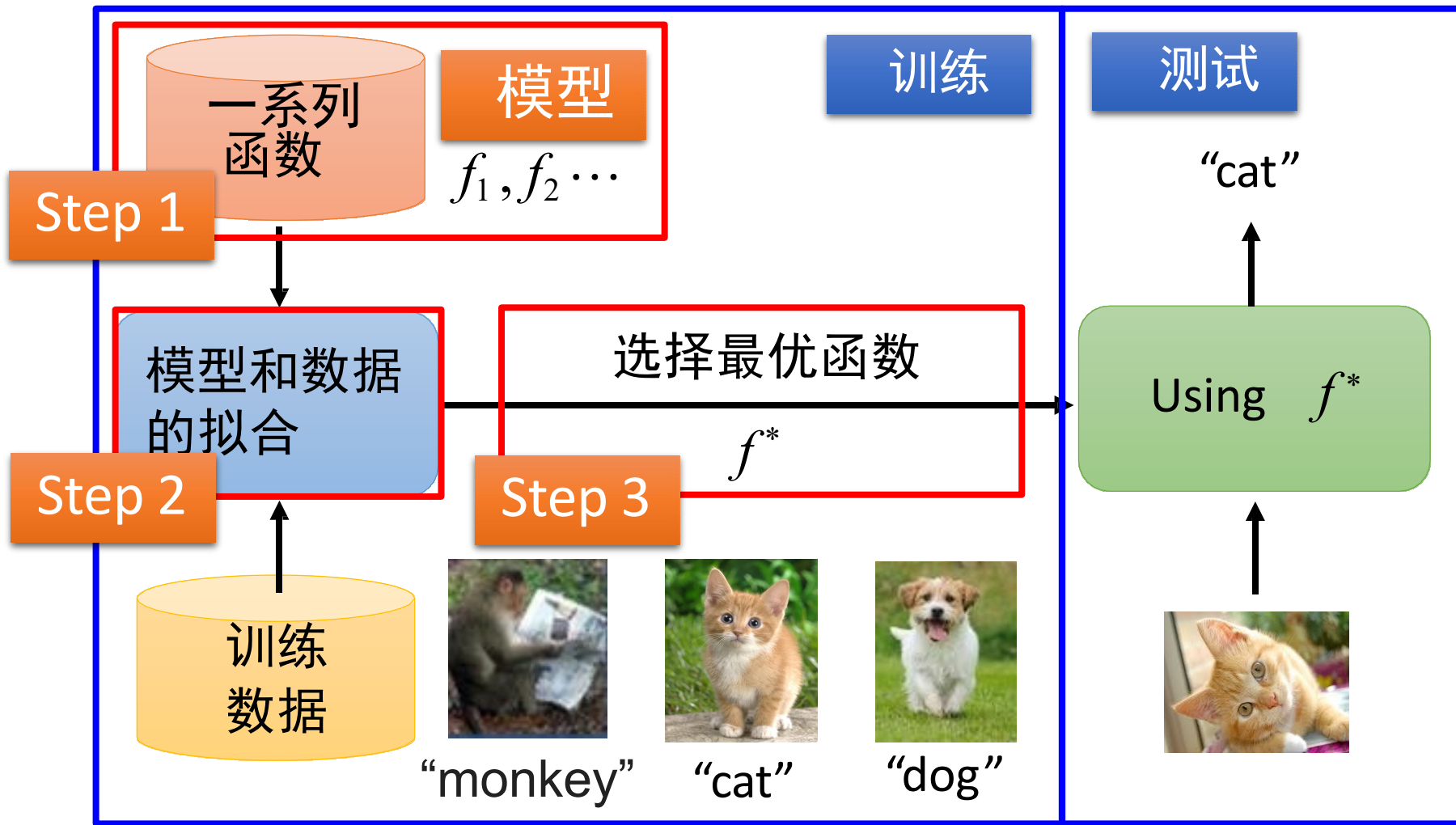
$$f(\text{"导航回家"}) = \text{"你说啥?"}$$

(你所说的内容) (系统的回答)

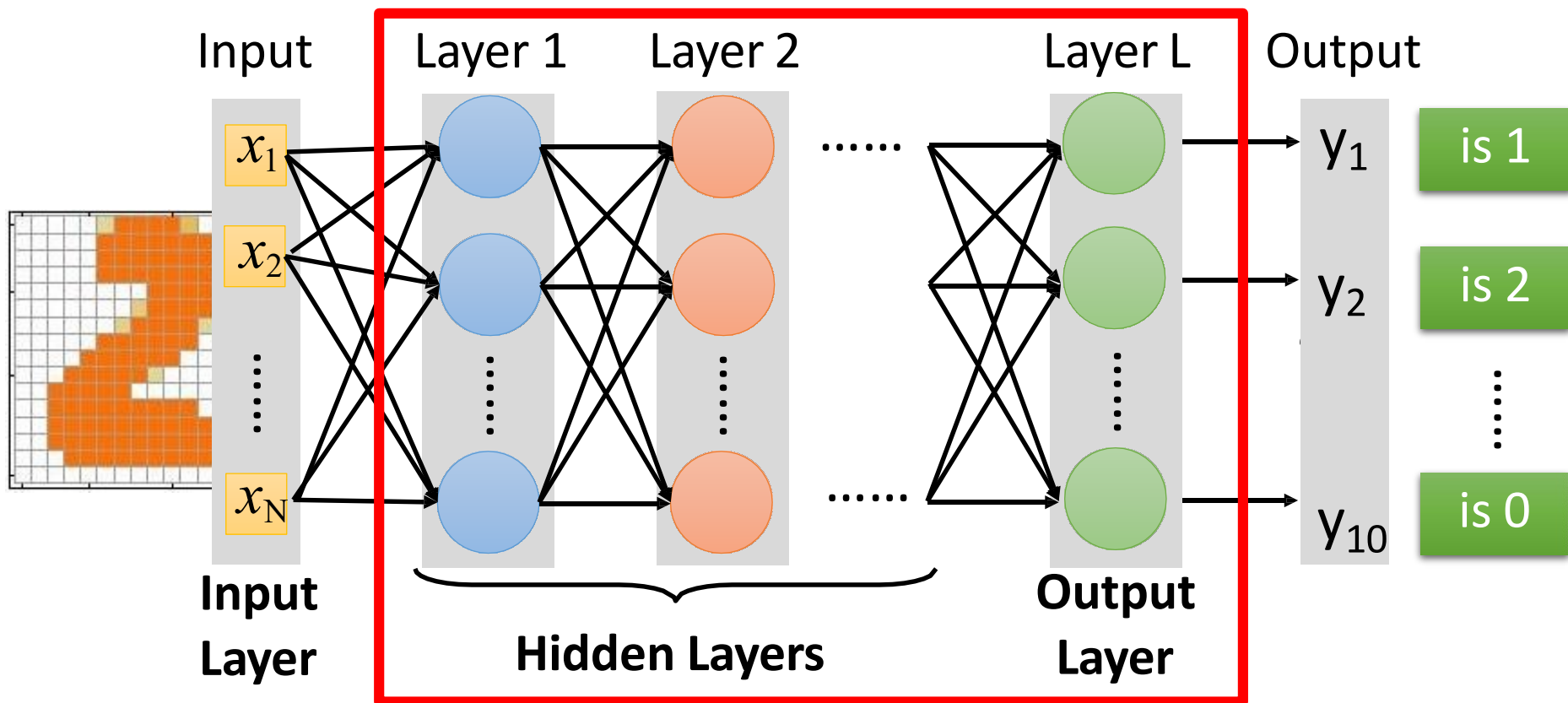
机器学习框架

图像识别:

$$f\left(\text{img}\right) = \text{"cat"}$$



深度神经网络结构



你需要决定网络结构，让你的函数有好的功能

神经网络的变体

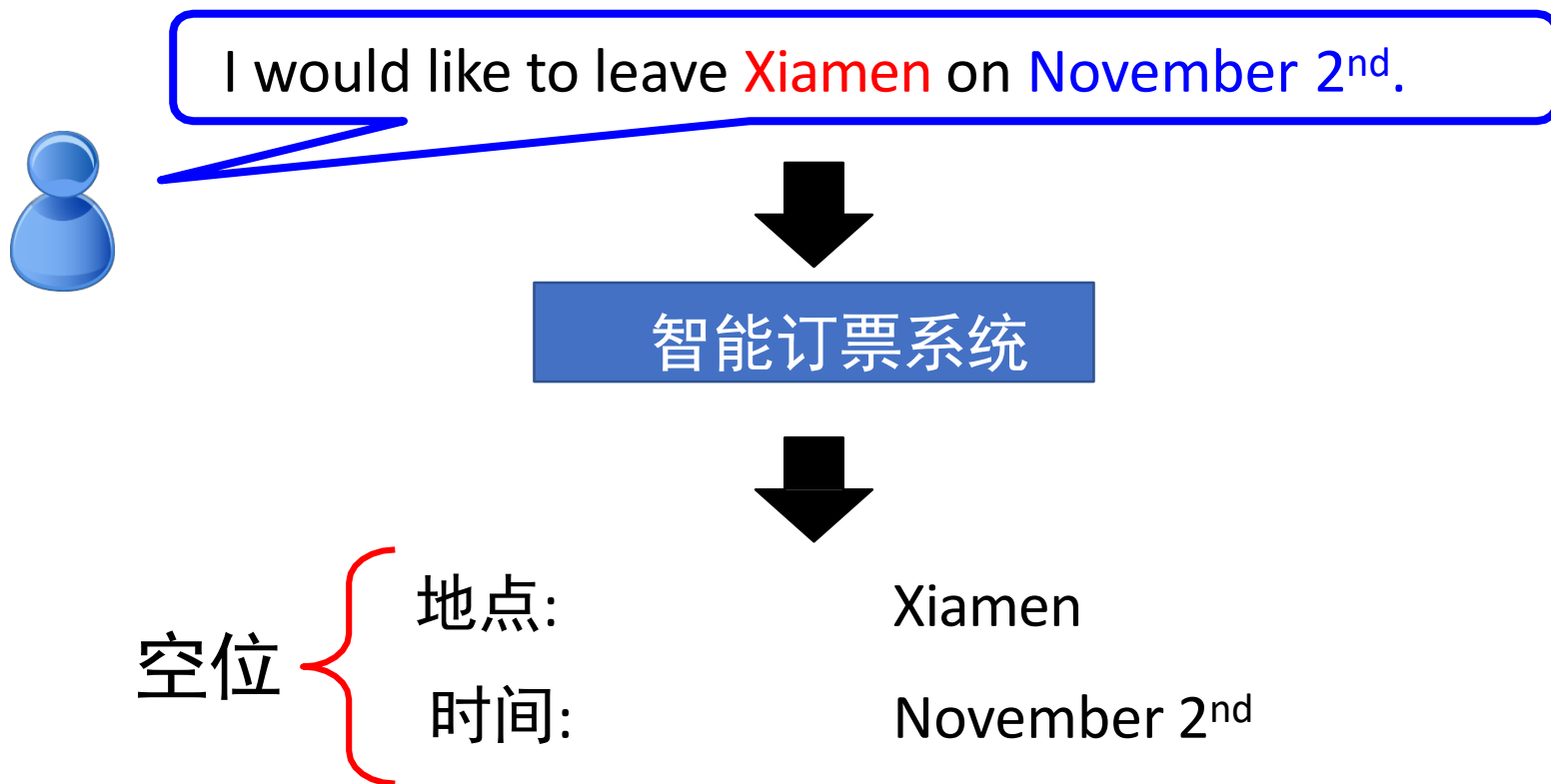
卷积神经网络
(CNN)

循环神经网络 (RNN)

具有记忆的神经网络

示例应用程序

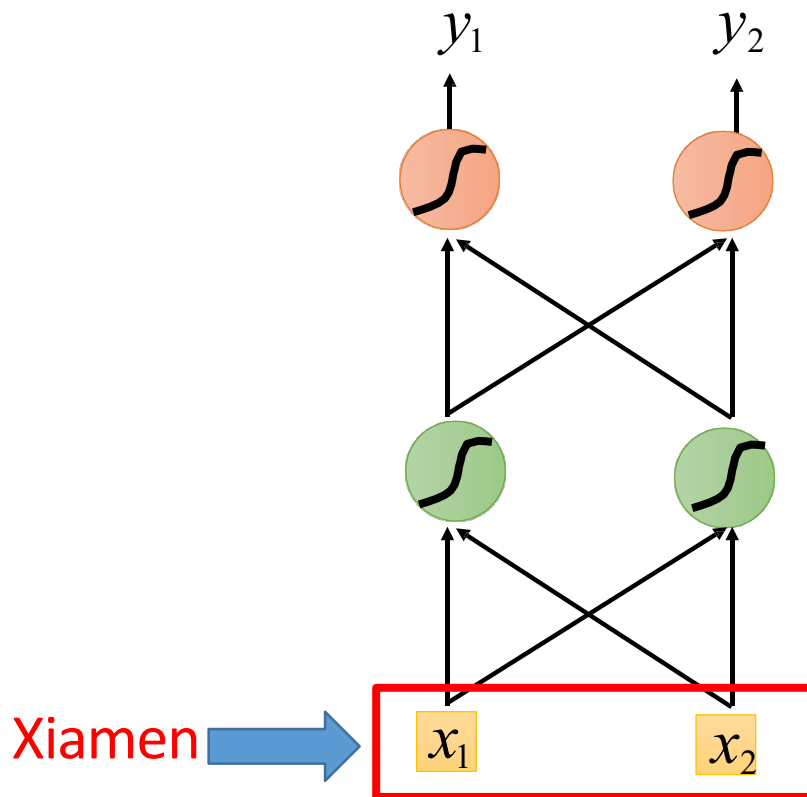
- 填空问题



示例应用程序

□ 通过前馈神经网络解决填空问题？

- 输入: a word (每个单词都被表示成一个向量)



1-of-N encoding

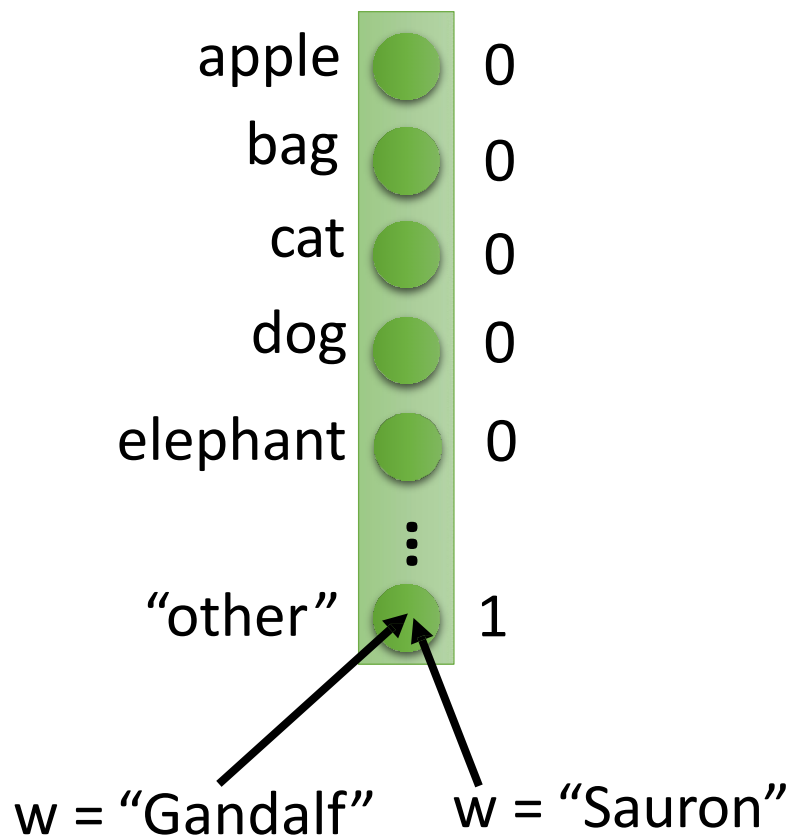
□ 如何将每个单词表示成一个词向量?

1-of-N Encoding 词典 = {apple, bag, cat, dog, elephant}

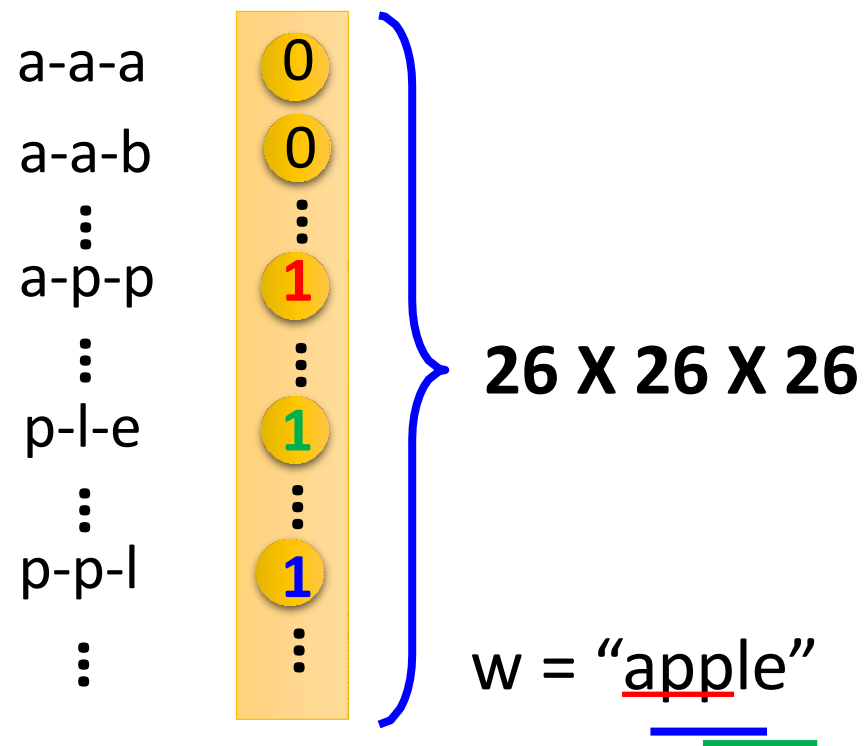
- 向量长度是词典大小
 - 每个维度对应于词典中的一个单词
 - 这个词的维度是1，其他是0
- | | | |
|----------|---|---------------|
| apple | = | [1 0 0 0 0] |
| bag | = | [0 1 0 0 0] |
| cat | = | [0 0 1 0 0] |
| dog | = | [0 0 0 1 0] |
| elephant | = | [0 0 0 0 1] |

Beyond 1-of-N encoding

Dimension for "Other"



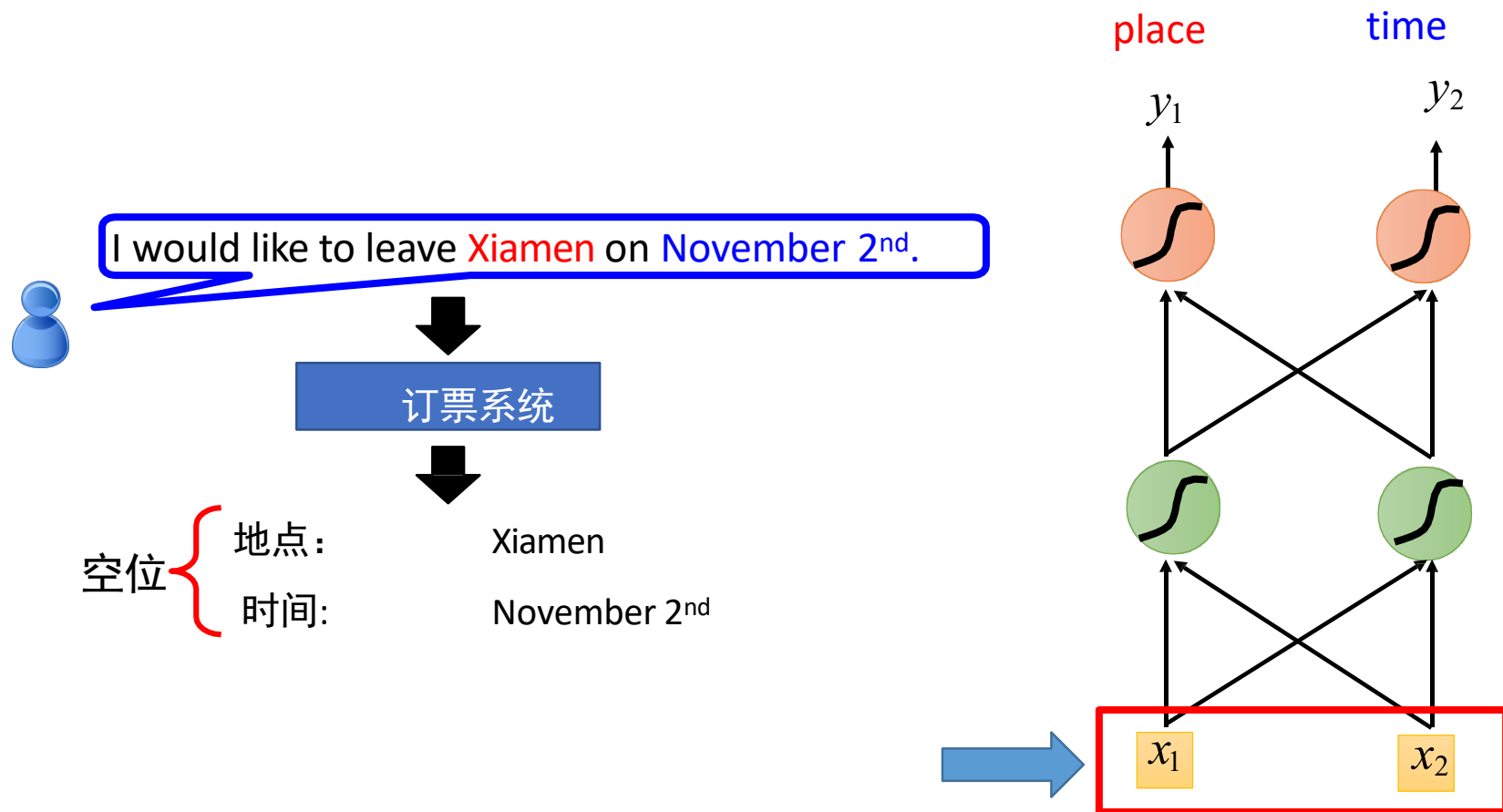
Word hashing



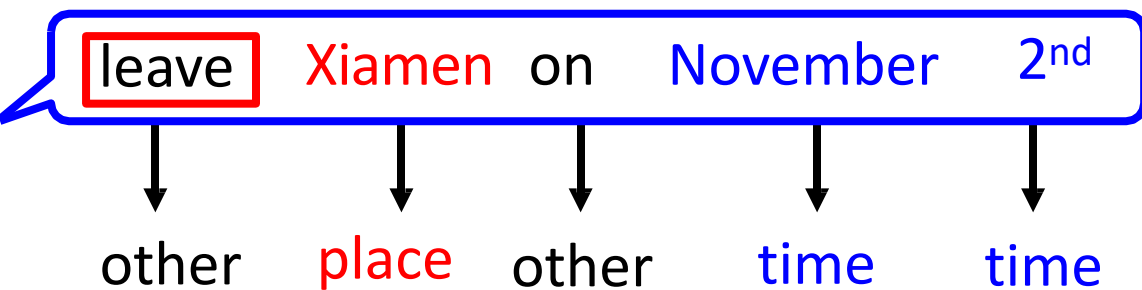
示例应用

□ 通过前馈神经网络解决填空问题？

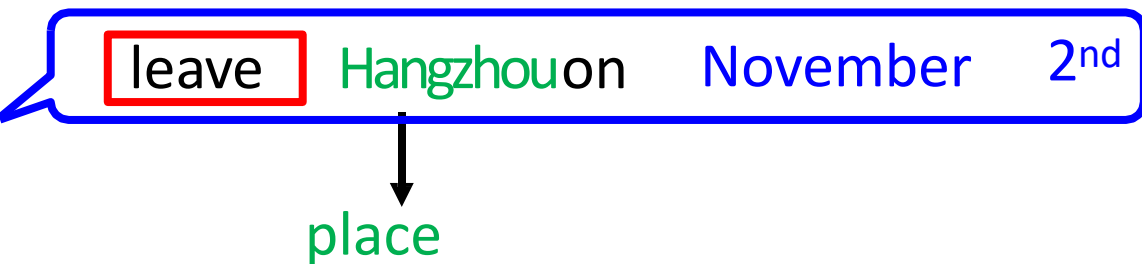
- 输入：a word (每个单词都被表示成一个向量)
- 输出：输入单词属于空位的概率分布



示例应用程序

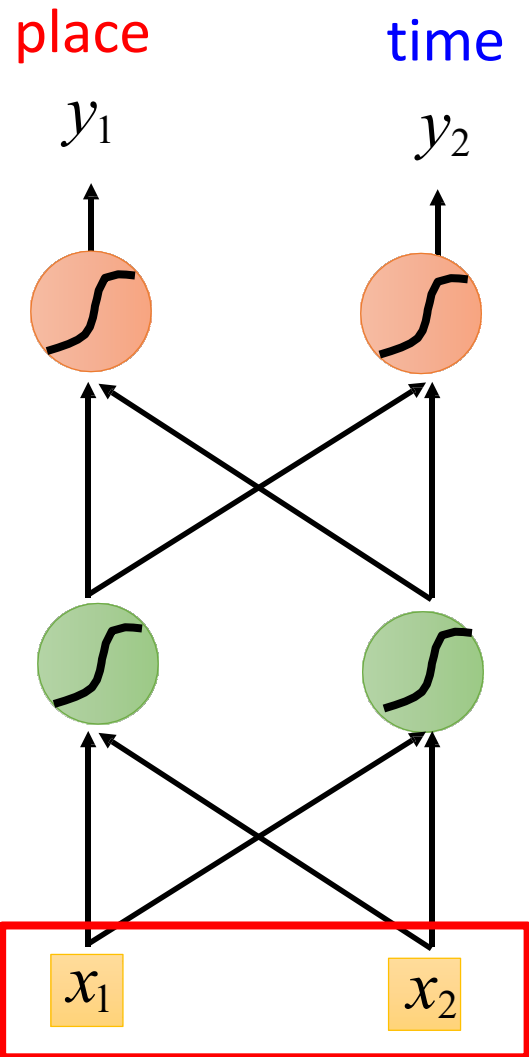


Problem?



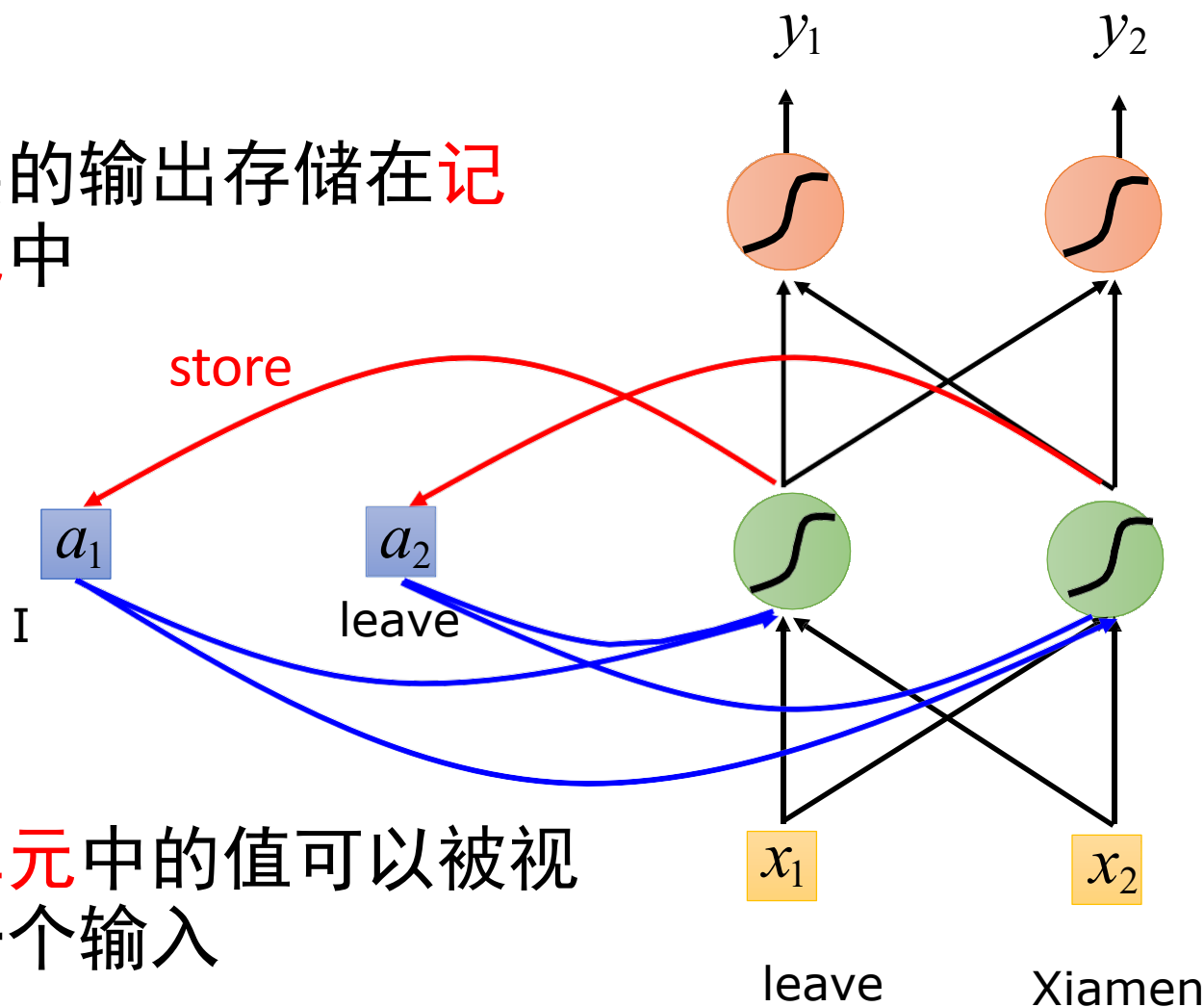
神经网络需要
记忆

Xiamen



循环神经网络 (RNN)

隐藏层的输出存储在**记忆单元**中



记忆单元中的值可以被视为另一个输入

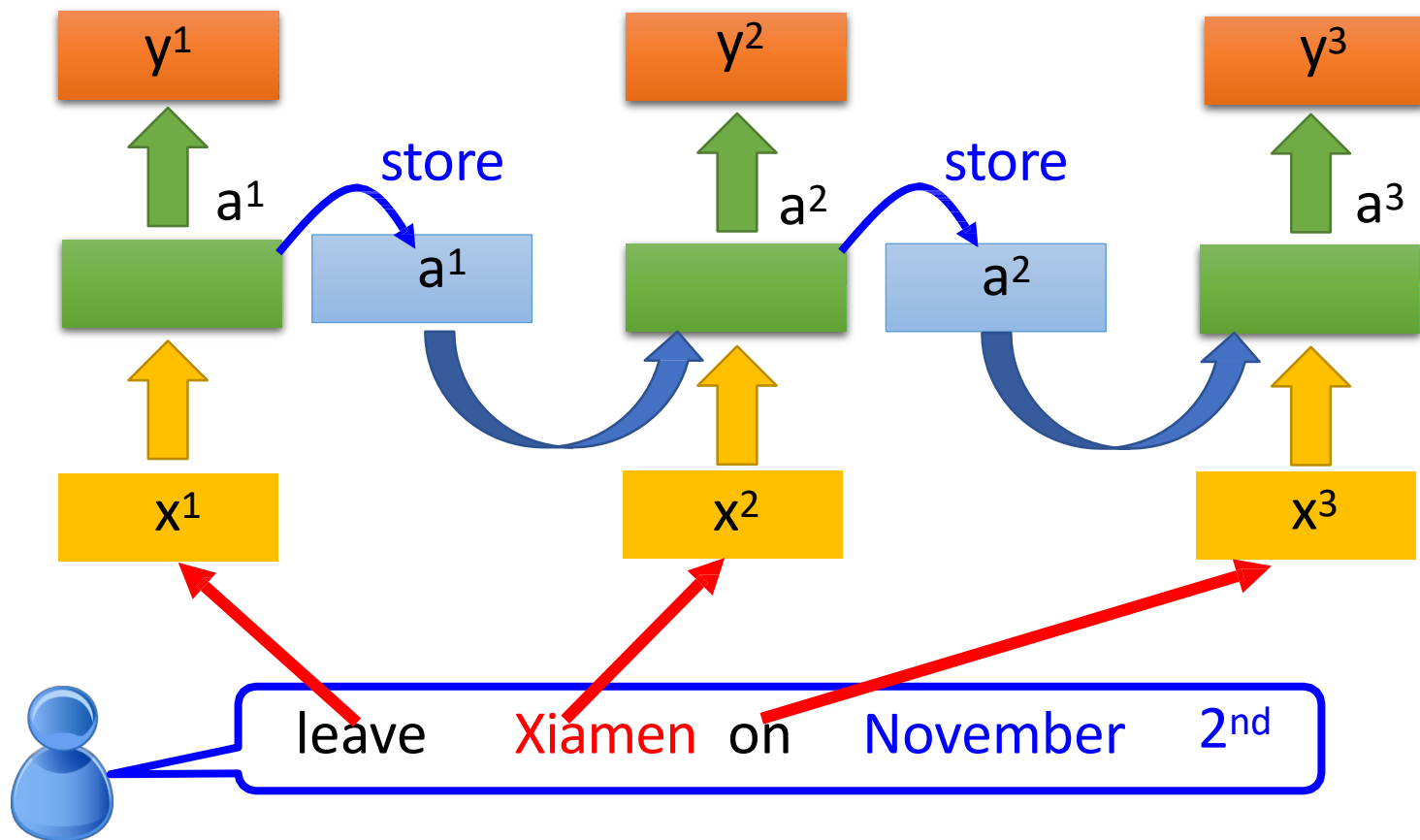
循环神经网络 (RNN)

同样的网络多次使用

Probability of
“arrive” in each slot

Probability of
“Xiamen” in each slot

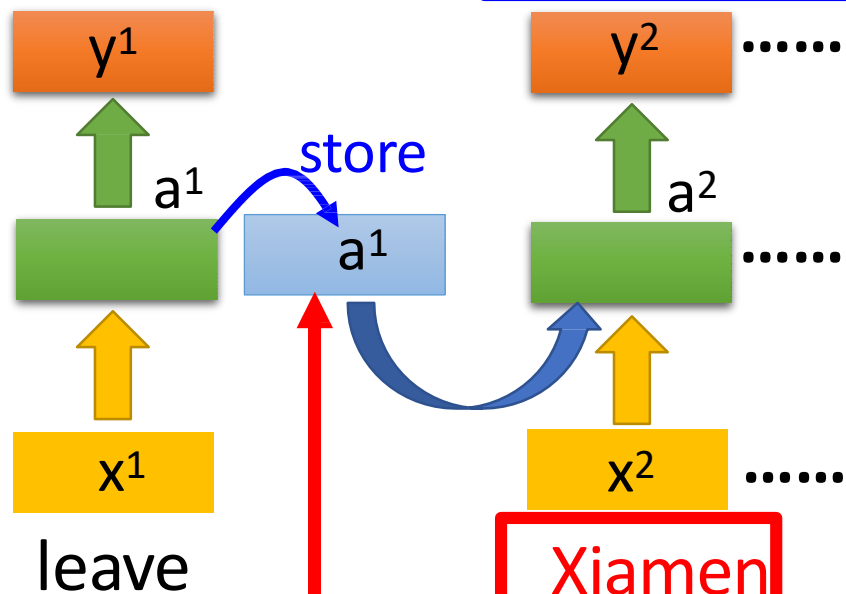
Probability of
“on” in each slot



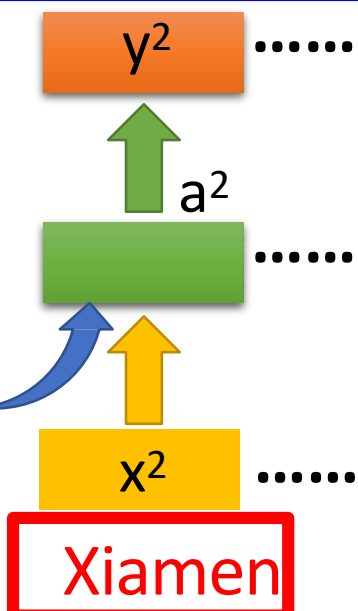
循环神经网络 (RNN)

Different

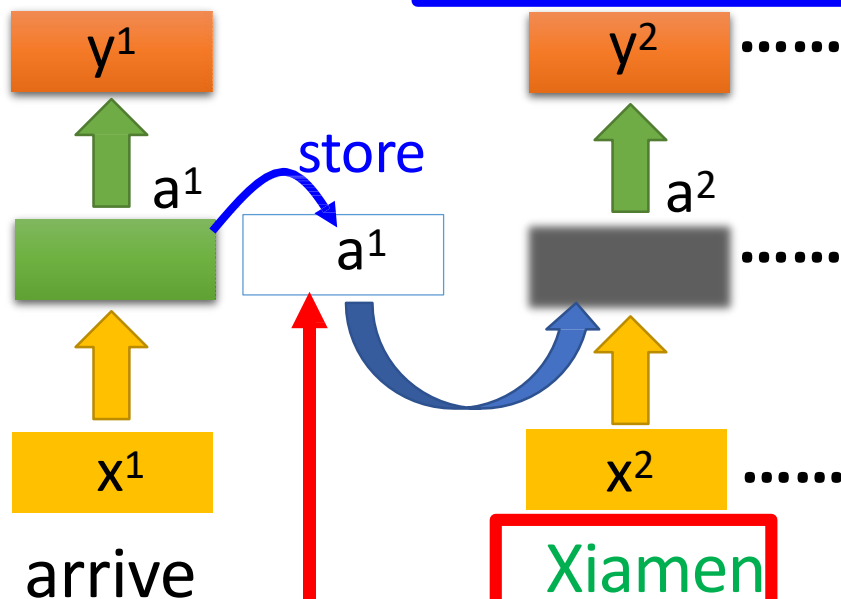
Prob of "leave"
in each slot



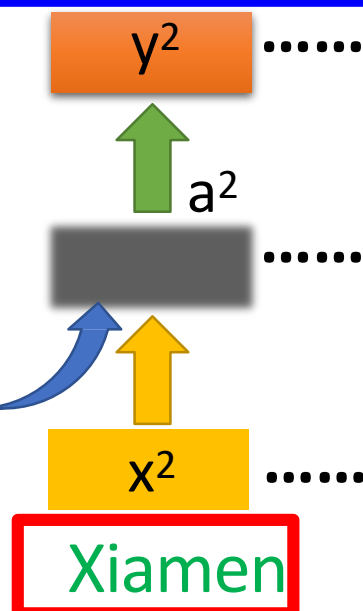
Prob of "Xiamen"
in each slot



Prob of "arrive"
in each slot

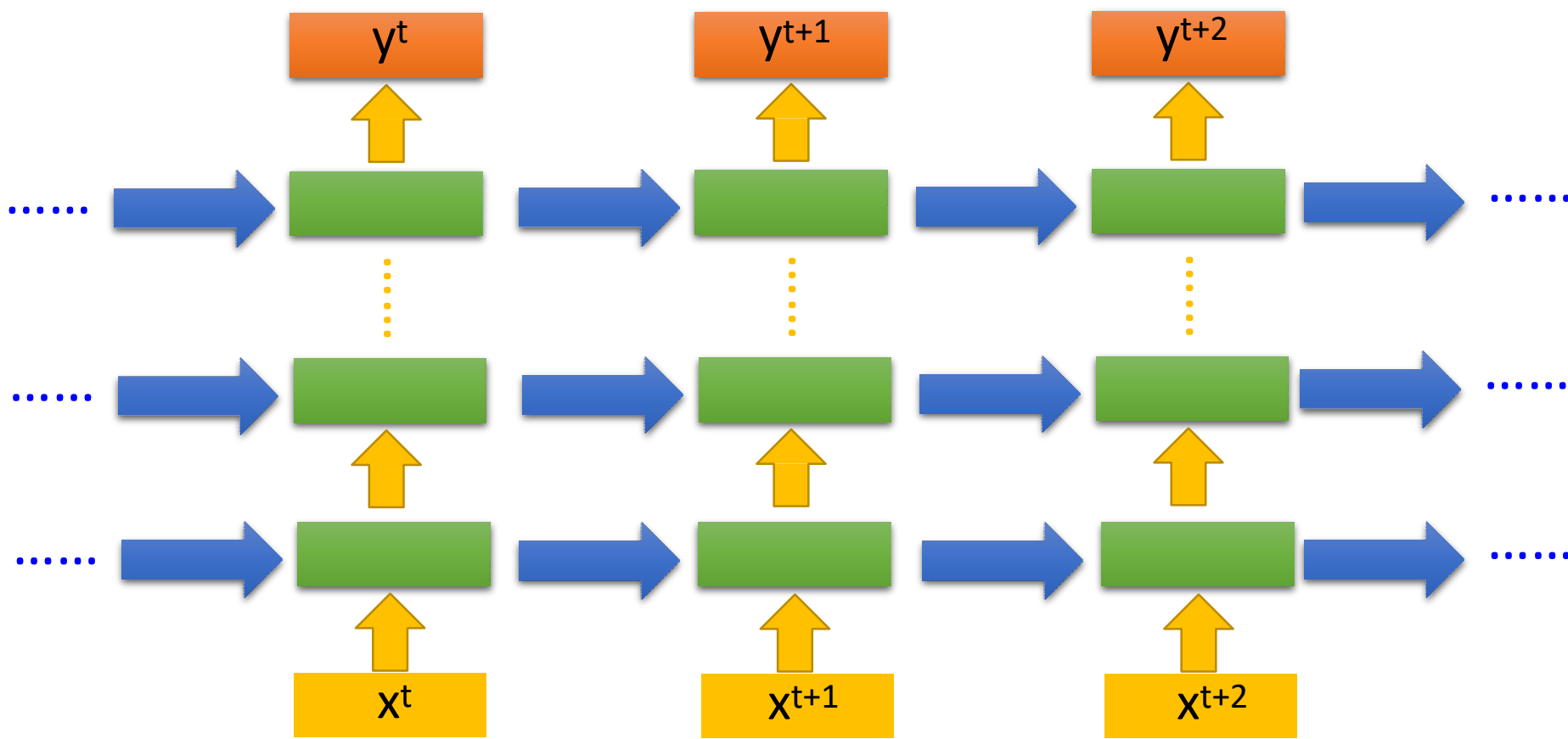


Prob of "Xiamen"
in each slot

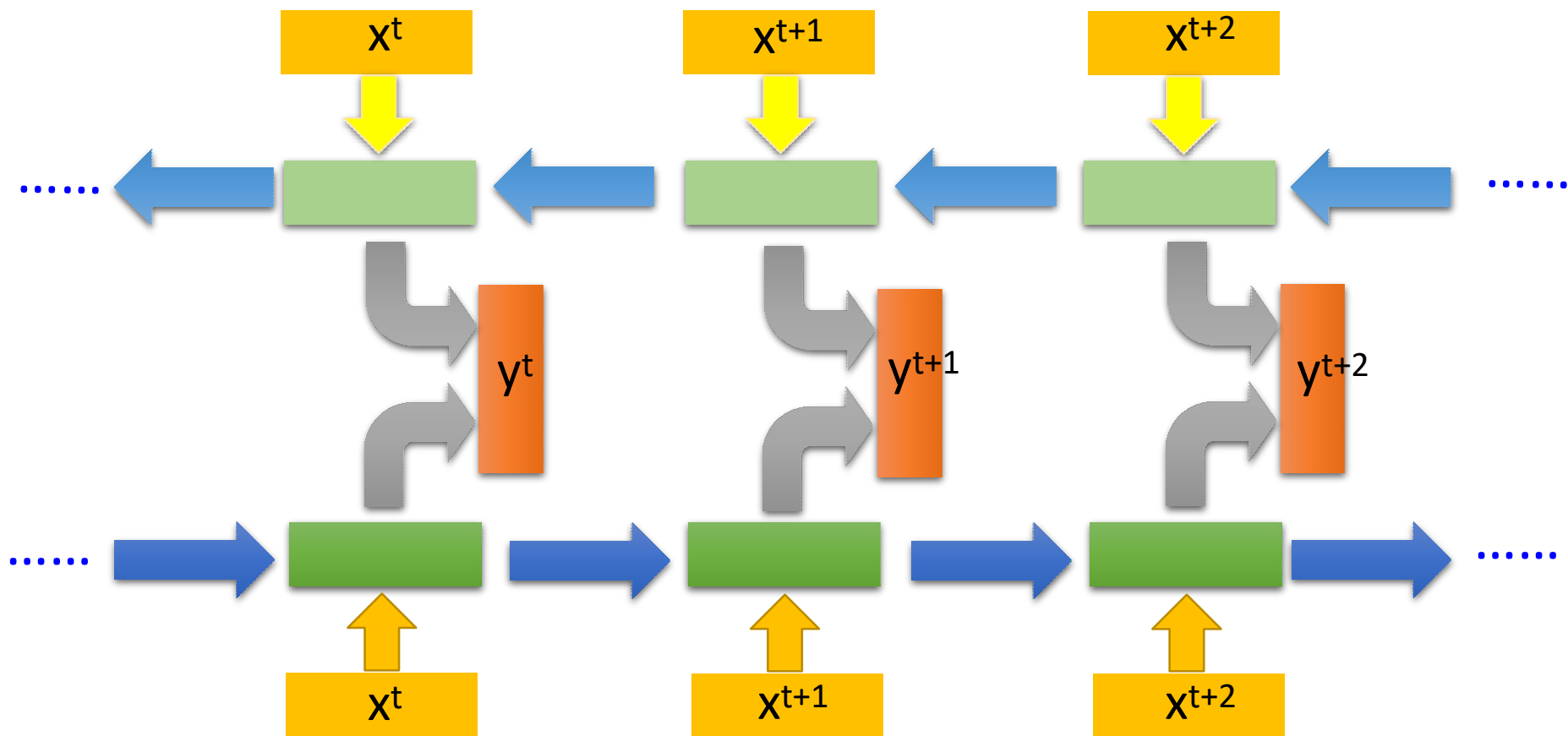


存在“记忆”里的值是不同的

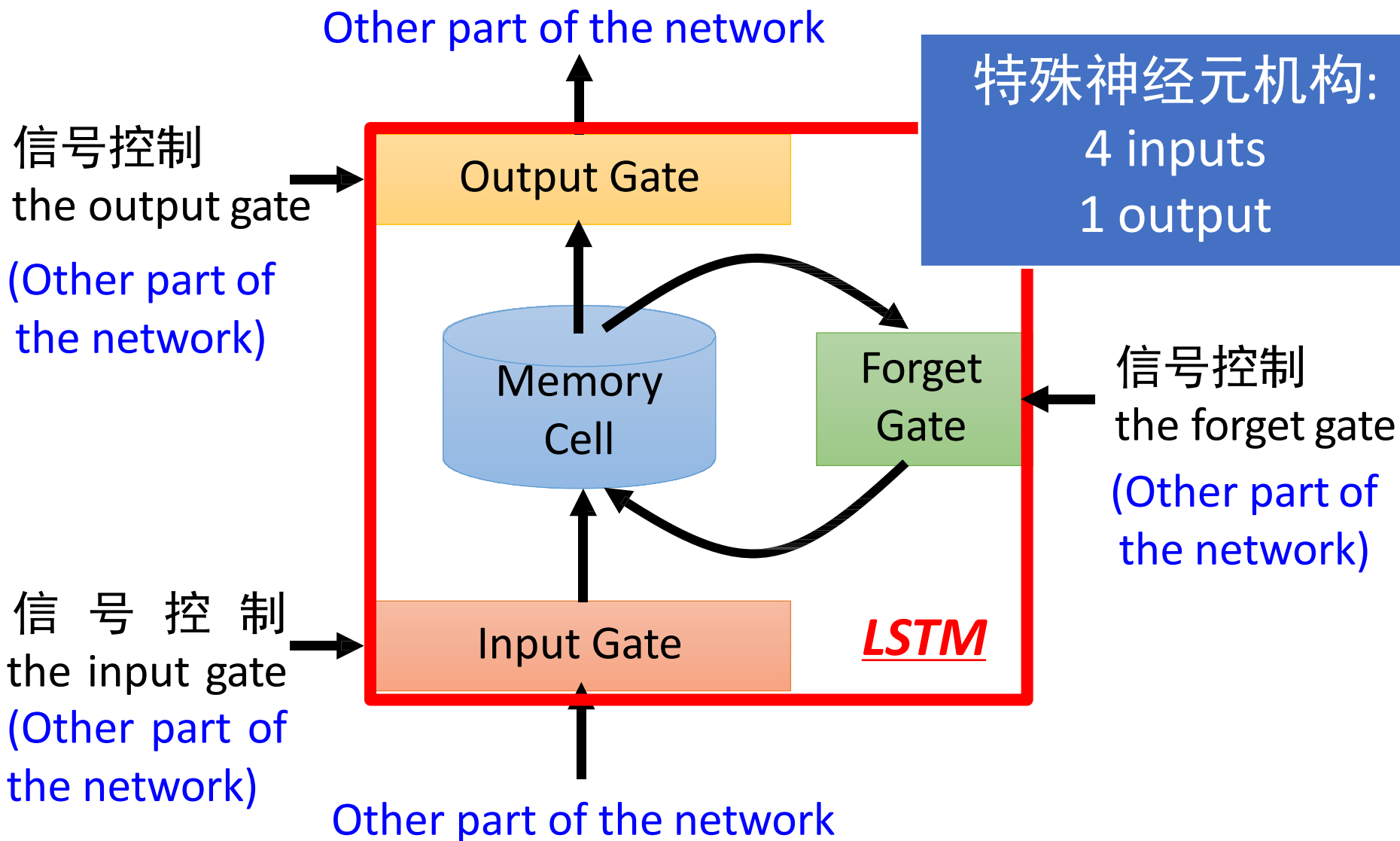
当然，它可以是深层网络 ...



双向 RNN



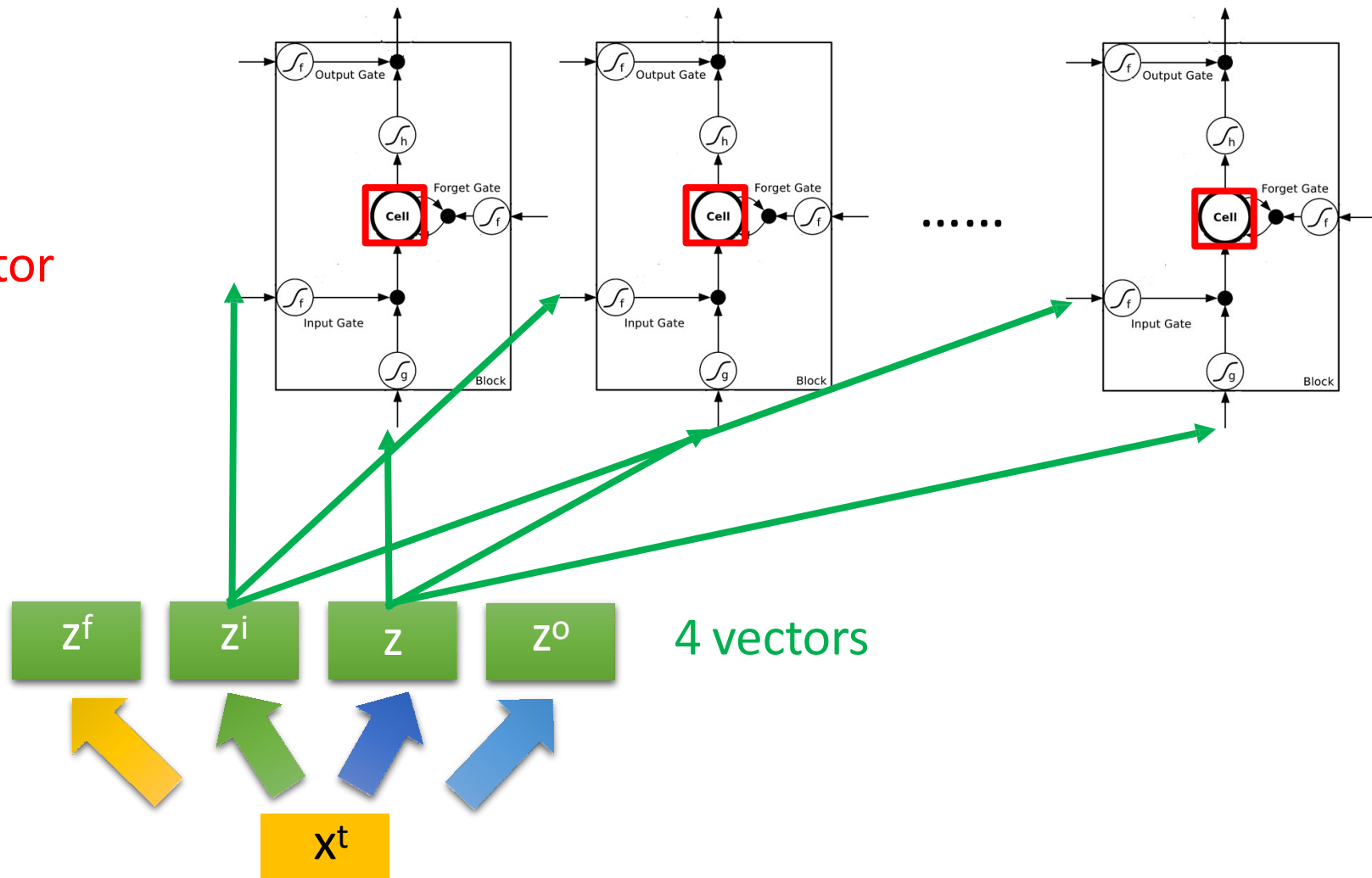
Long Short-term Memory (LSTM)



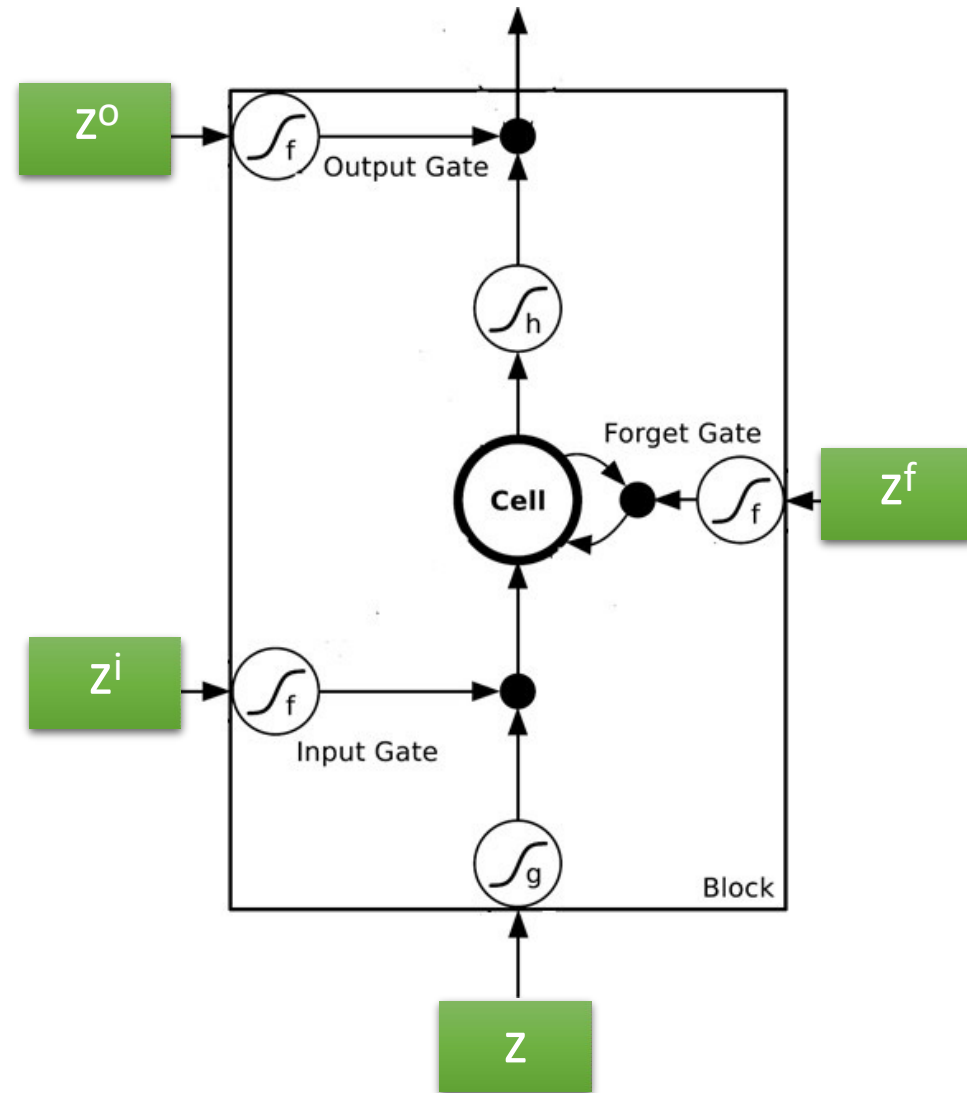
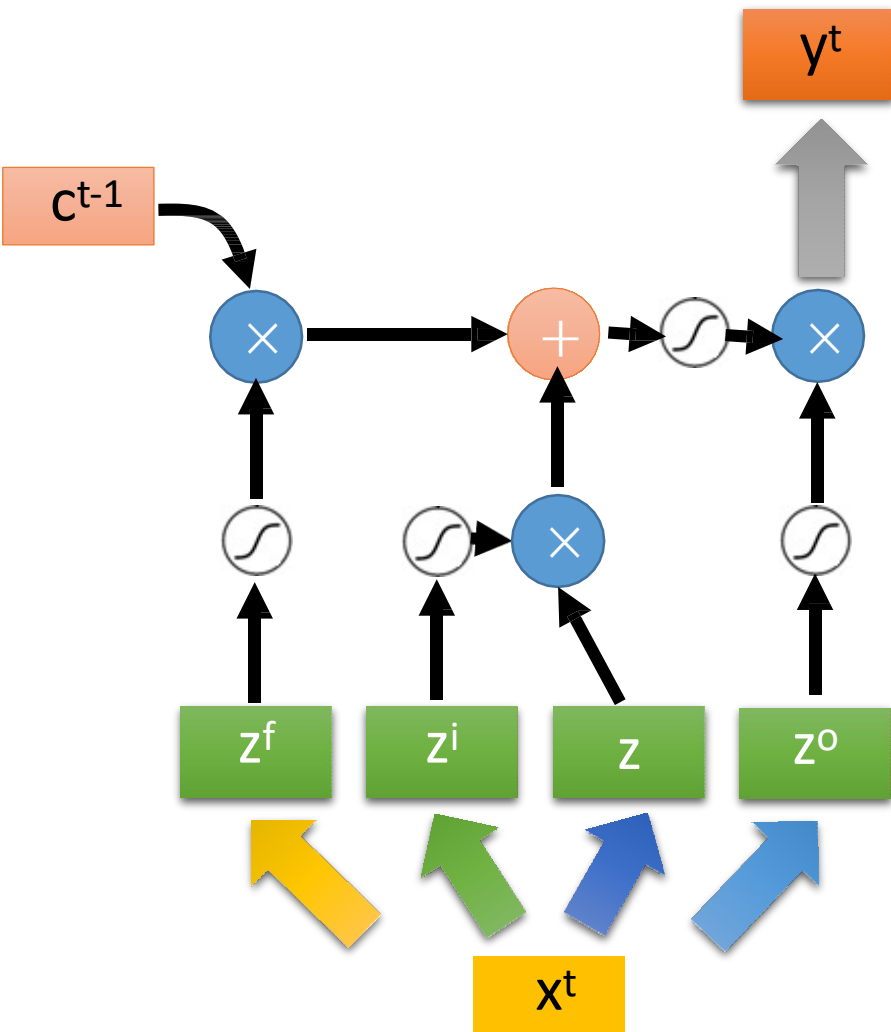
LSTM

 c^{t-1}

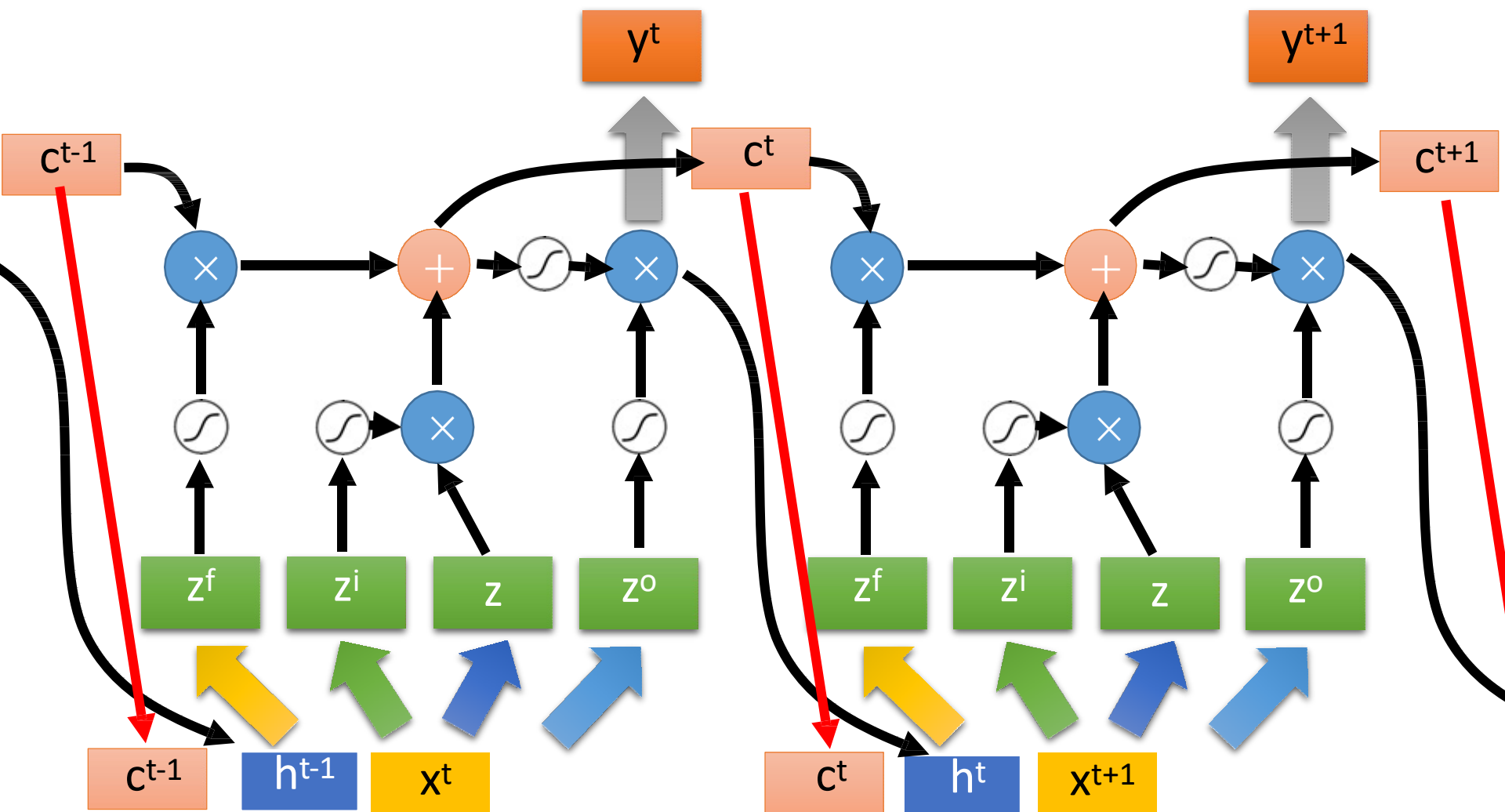
vector



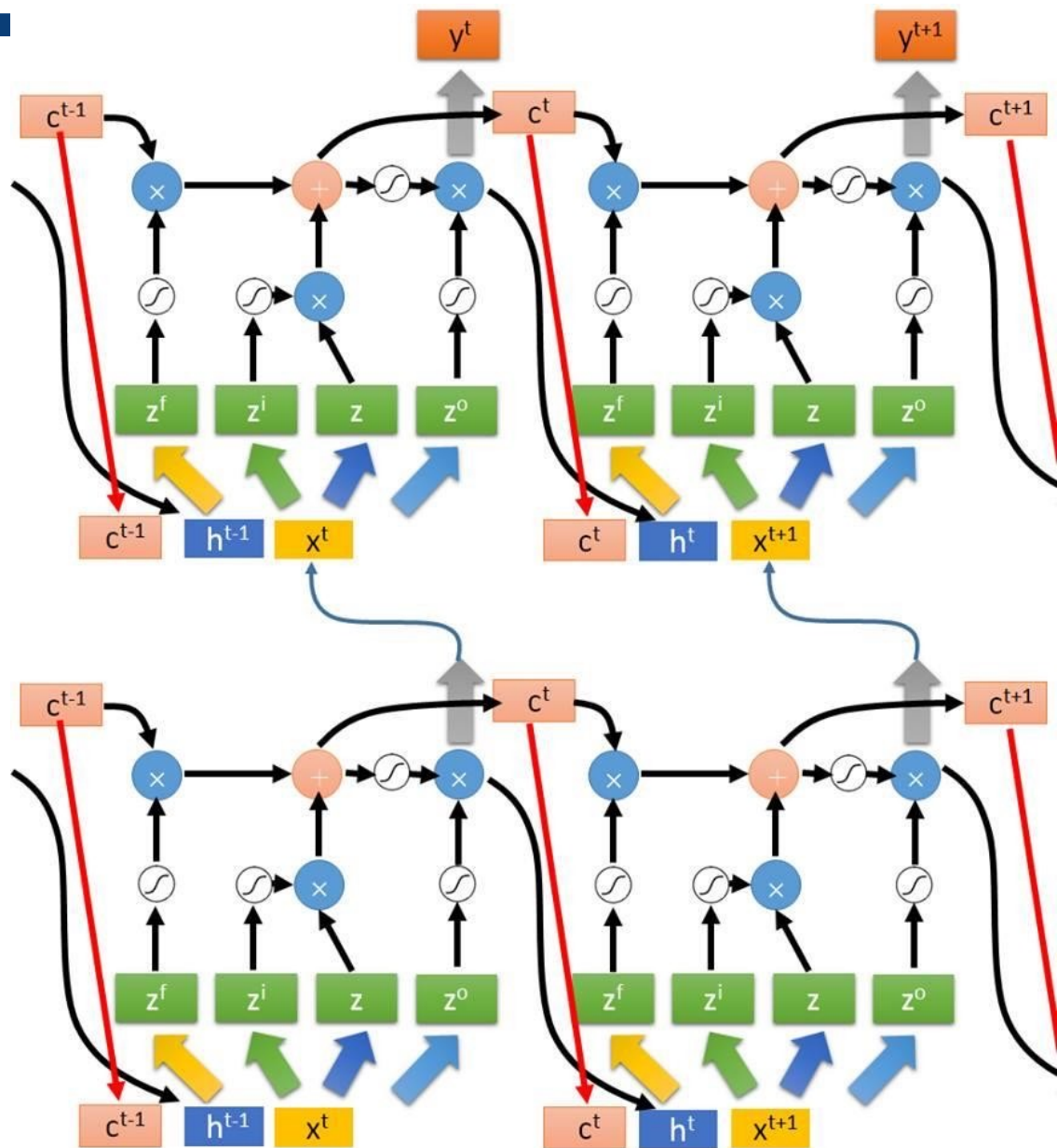
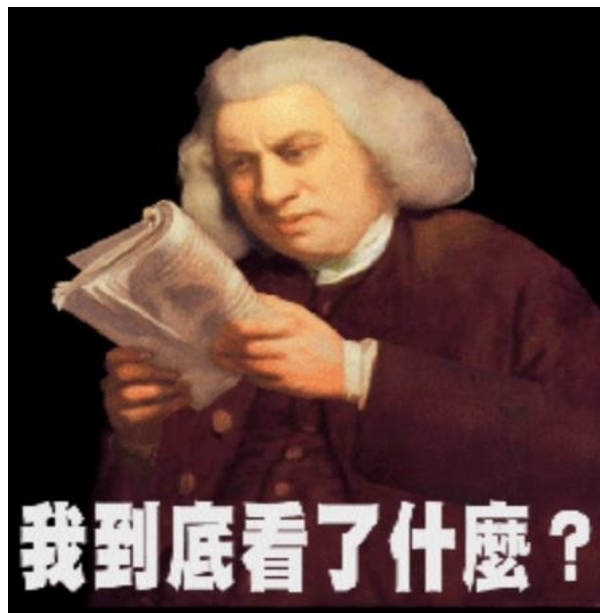
LSTM



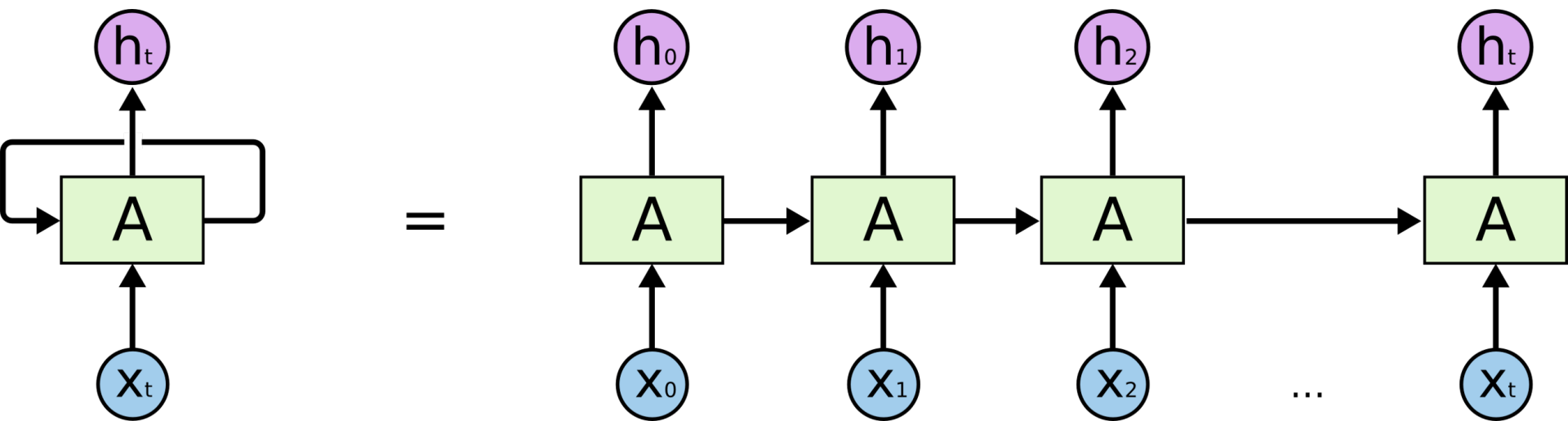
LSTM



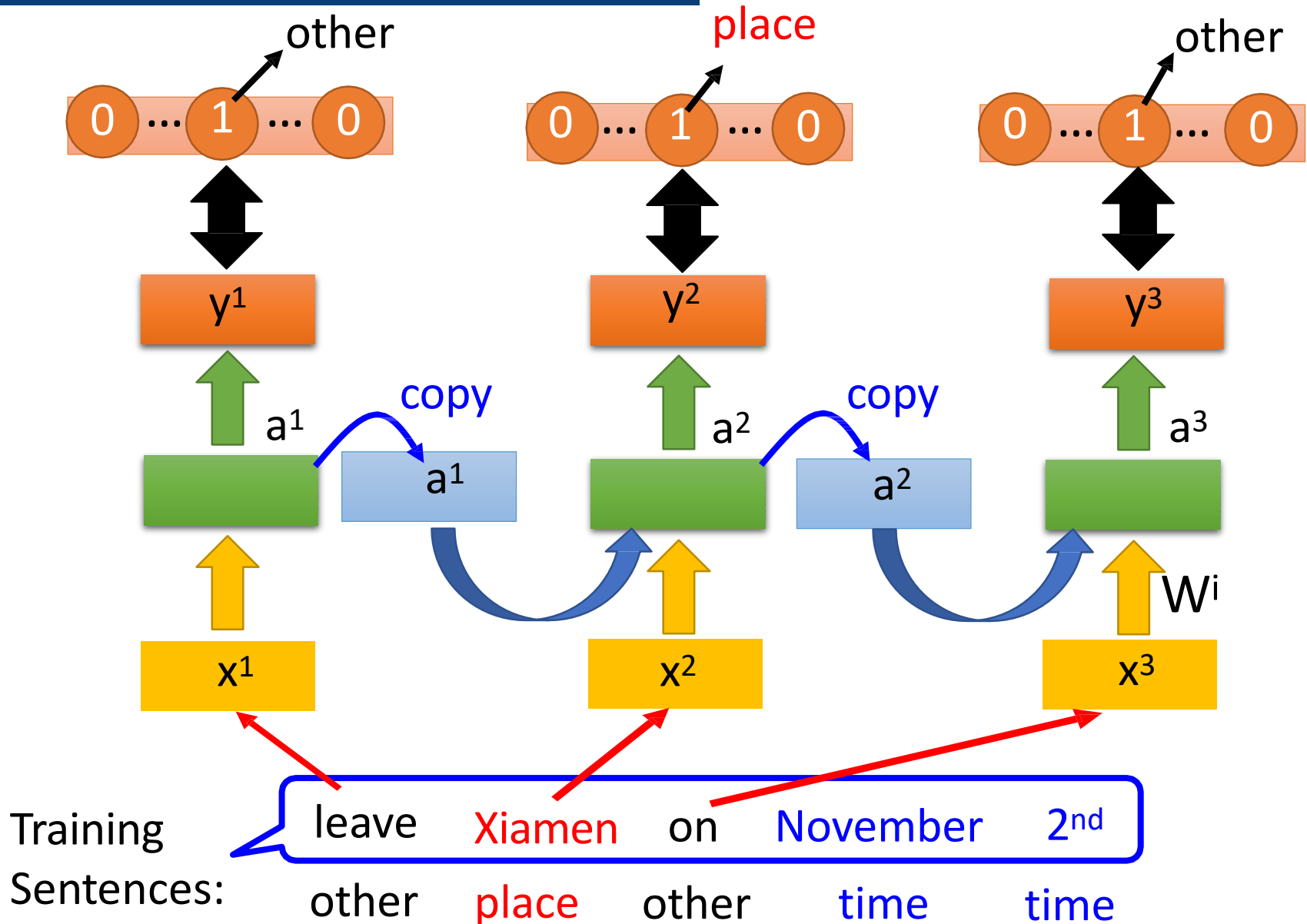
Multiple-layer LSTM



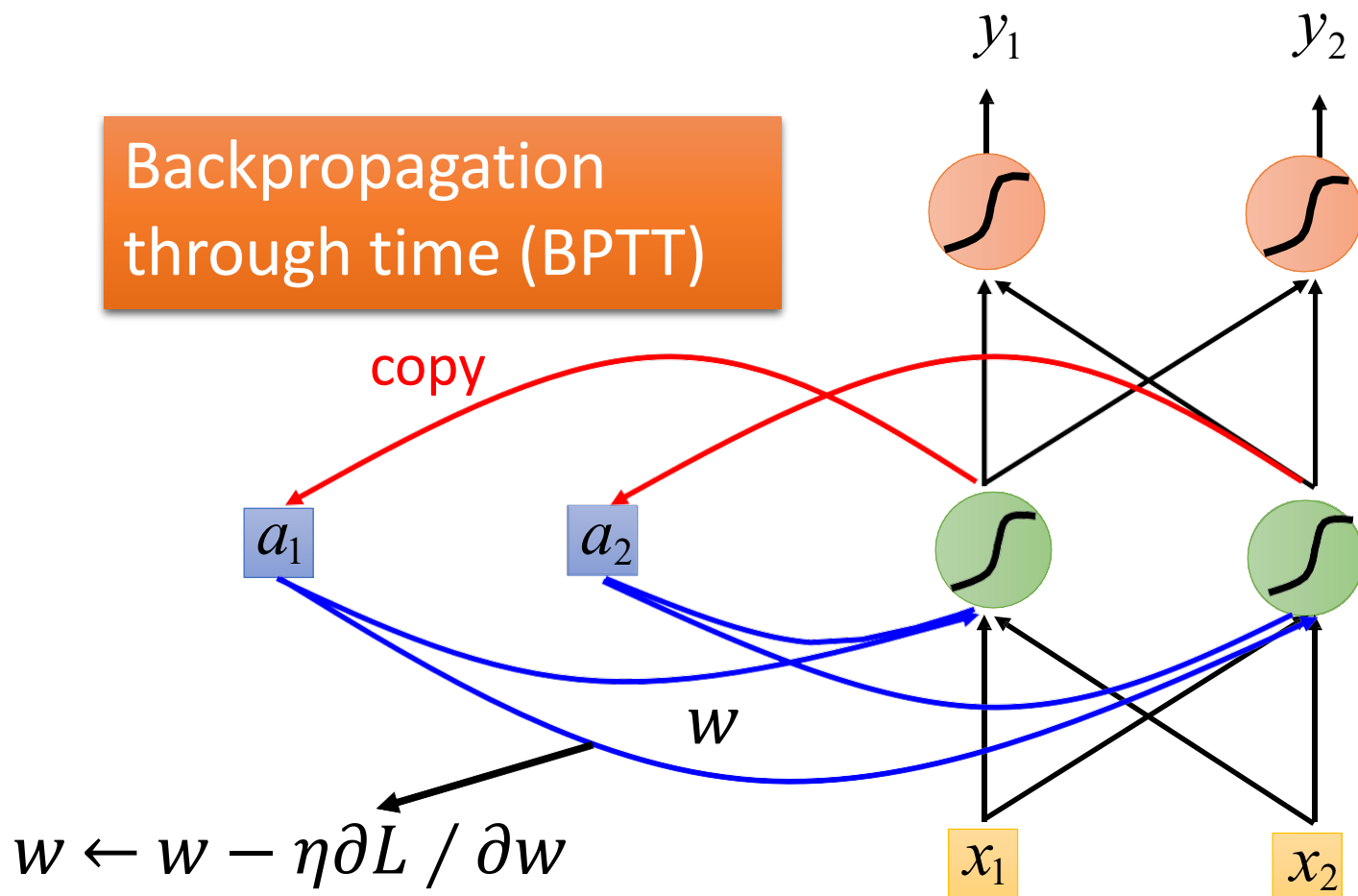
RNN结构



Learning Target



Learning



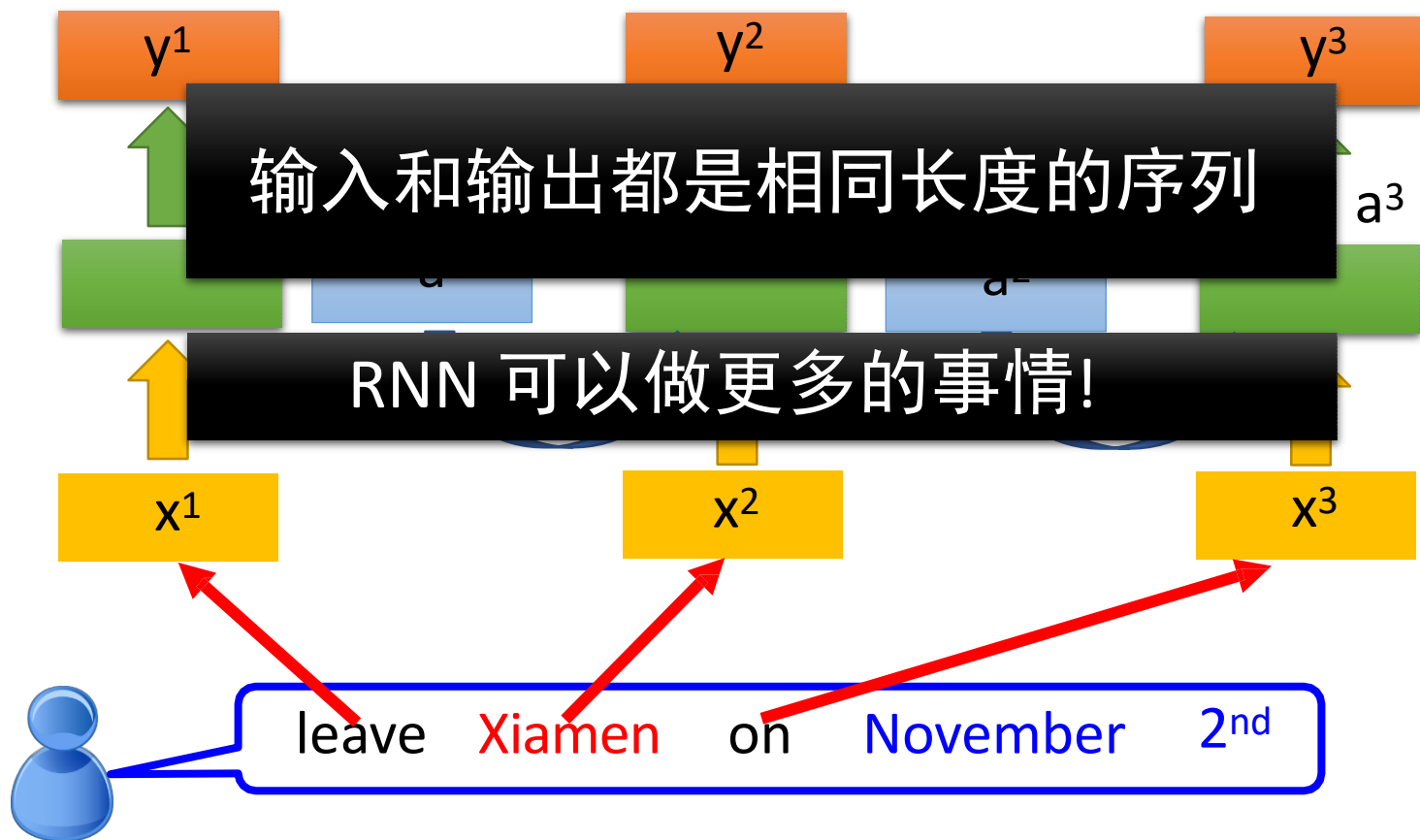
RNN的学习在实践中是非常困难的

更多应用

Probability of
“leave” in each slot

Probability of
“Xiamen” in each slot

Probability of
“on” in each slot



Many to one

- 输入是一个向量序列，但输出只是一个向量

情感分析

看了这部电影觉得很高兴

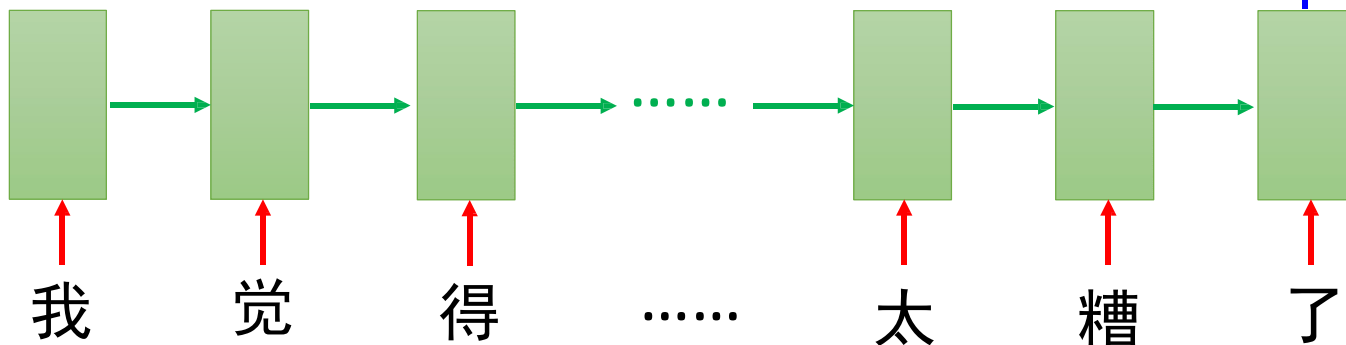
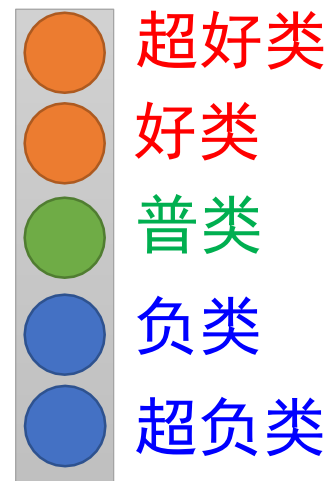
Positive (正类)

这部电影太糟了

Negative (负类)

这部电影真的很棒

Positive (正类)



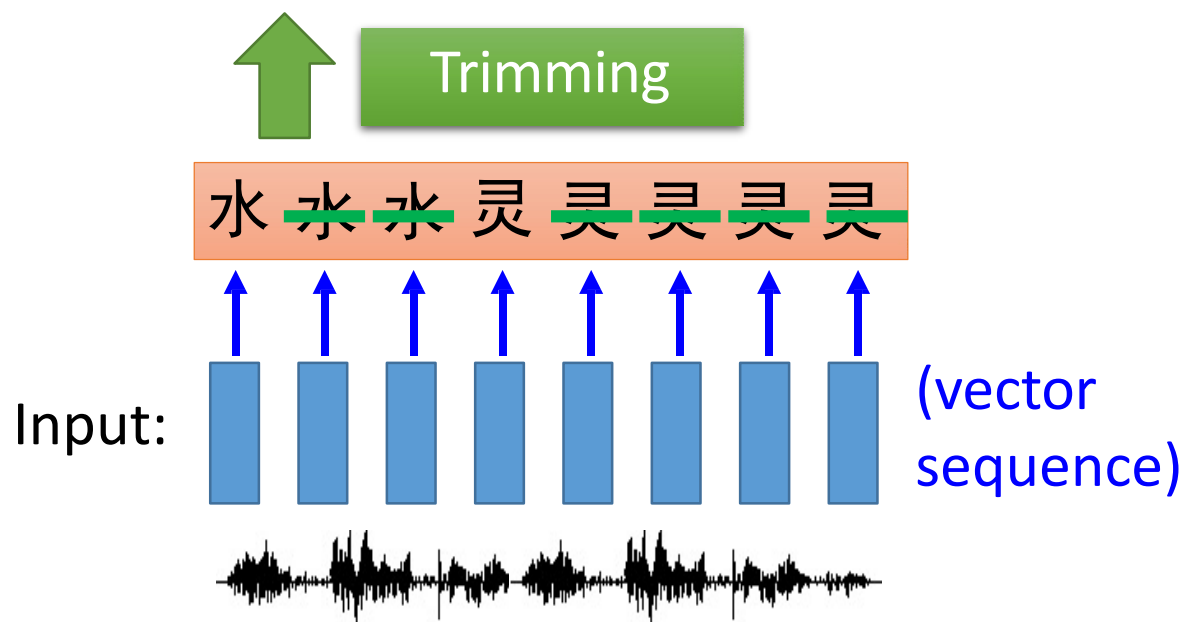
Many to Many (输出较短)

- 输入和输出都是序列, 但输出比较短
 - E.g. 语音识别

Output: “水灵” (character sequence)

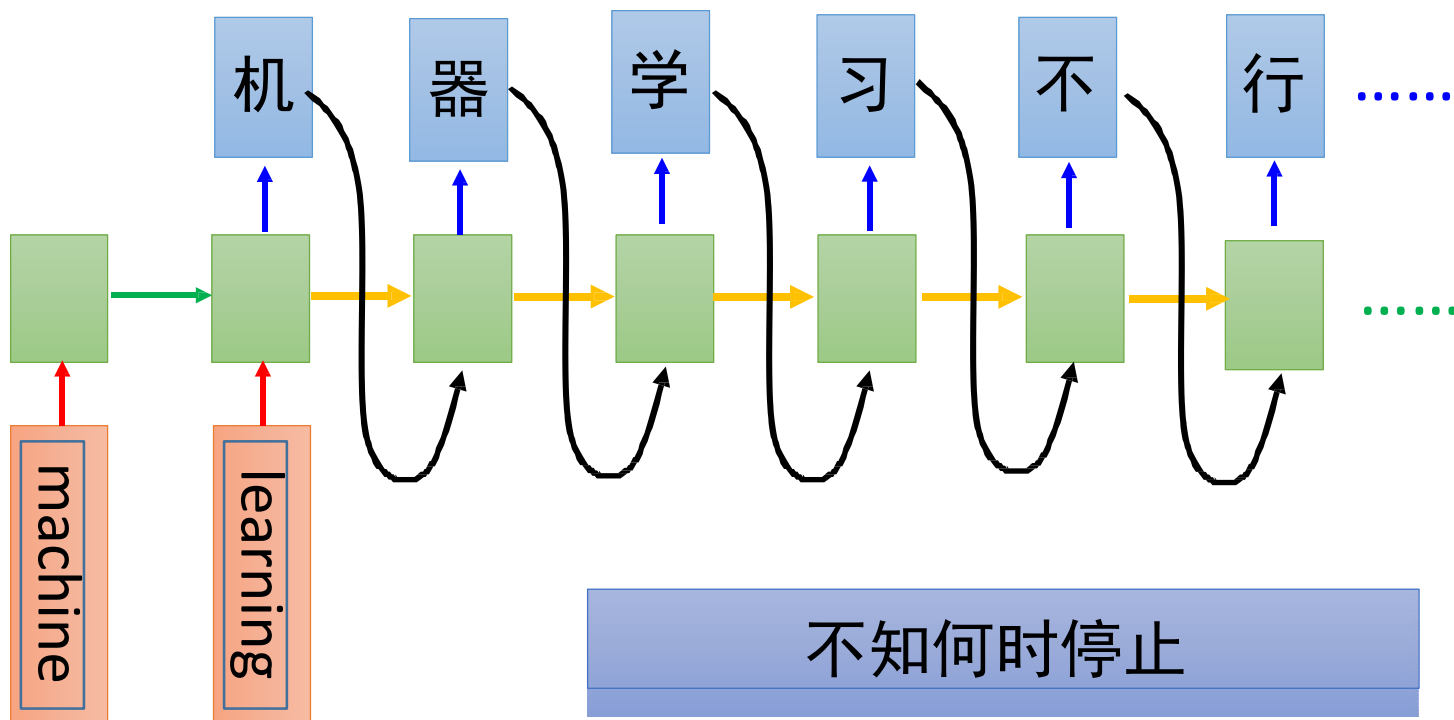
Problem?

Why can't it be
“水灵灵”



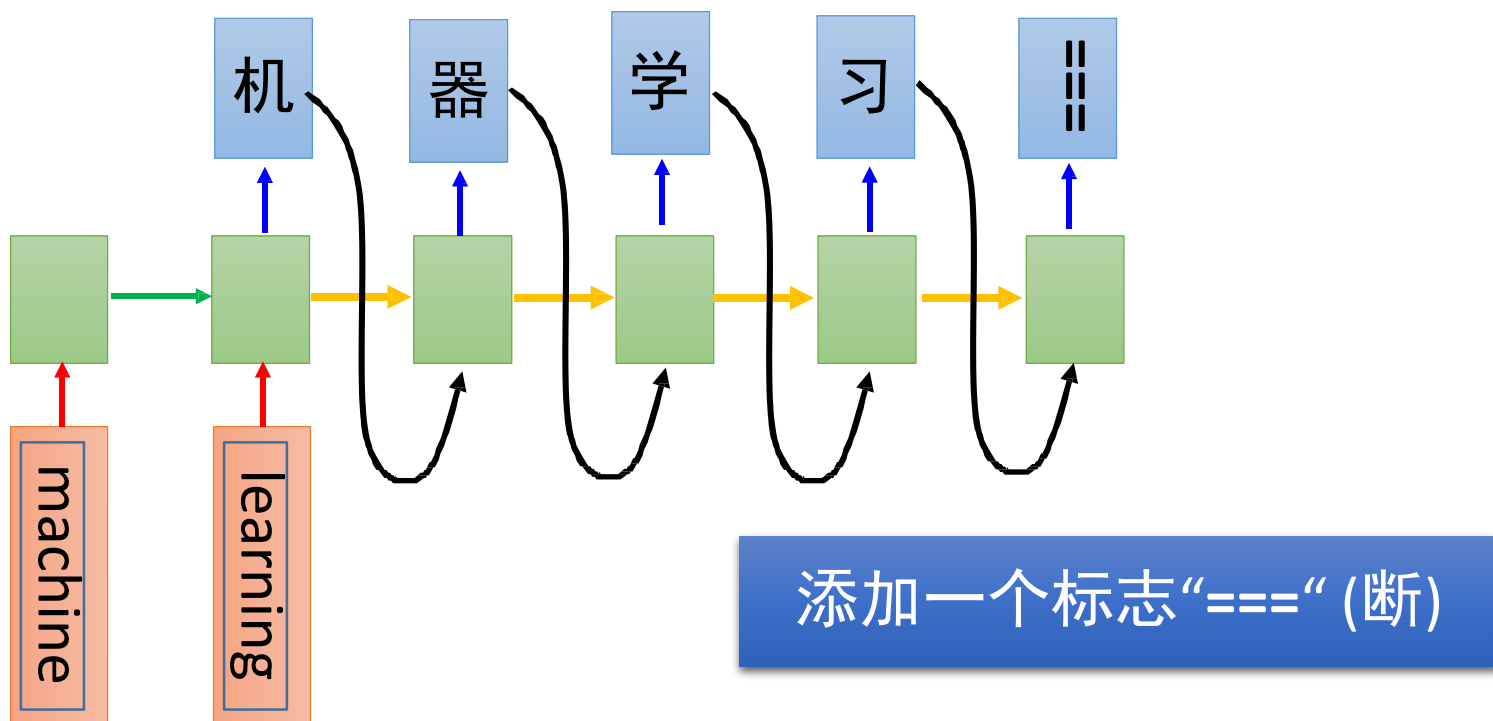
Many to Many (No Limitation)

- 输入和输出都是不同长度的序列 → Sequence to sequence learning
 - E.g. 机器翻译 (machine learning → 机器学习)



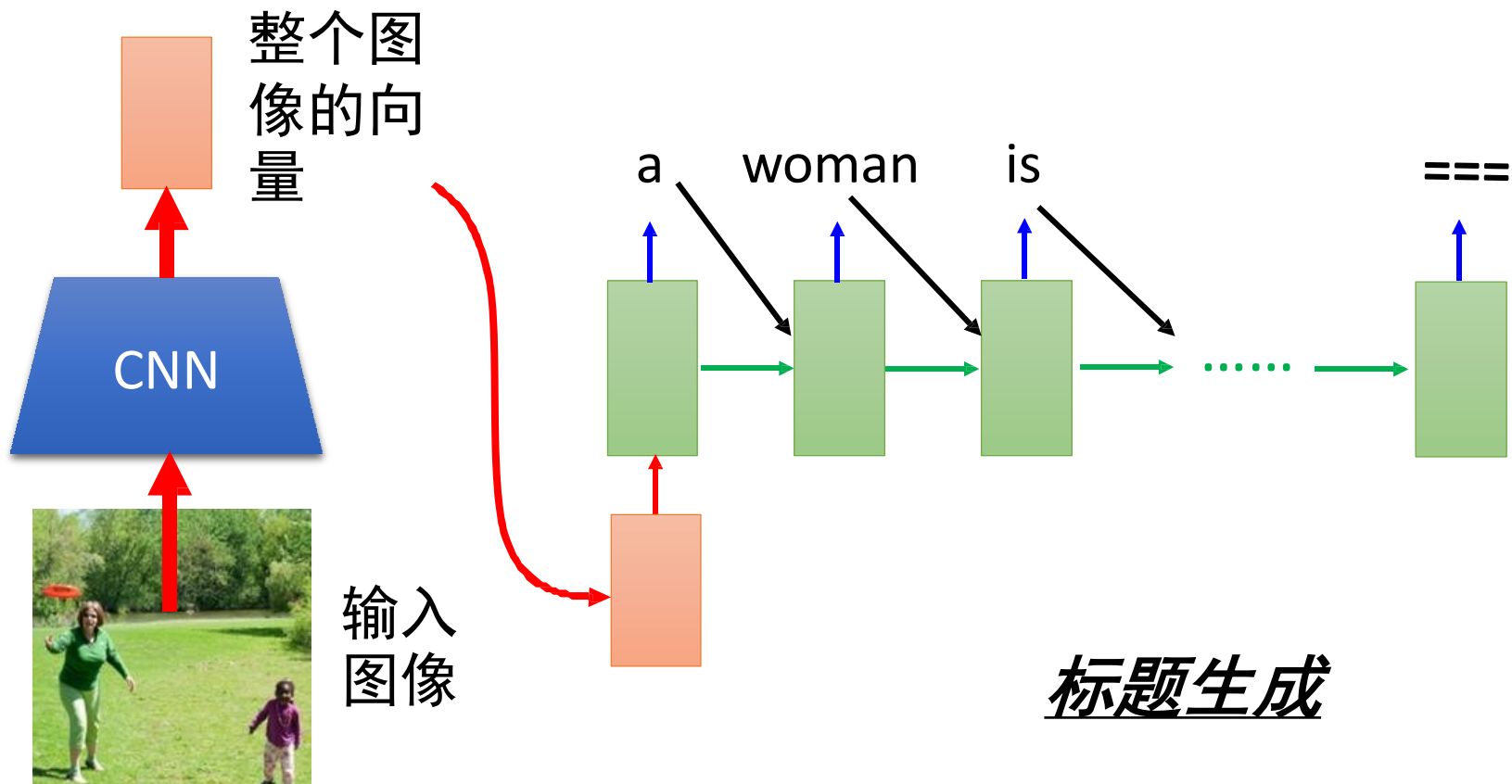
Many to Many (没有限制)

- 输入和输出都是不同长度的序列 → Sequence to sequence learning
 - E.g. 机器翻译 (machine learning → 机器学习)



One to Many

- 输入图像，但输出一系列单词



Demo

□ Image Captioning

■ <https://www.captionbot.ai>

□ [Prof. Cheng WANG](#)

□ [皮卡丘](#)

□ [皮卡](#)

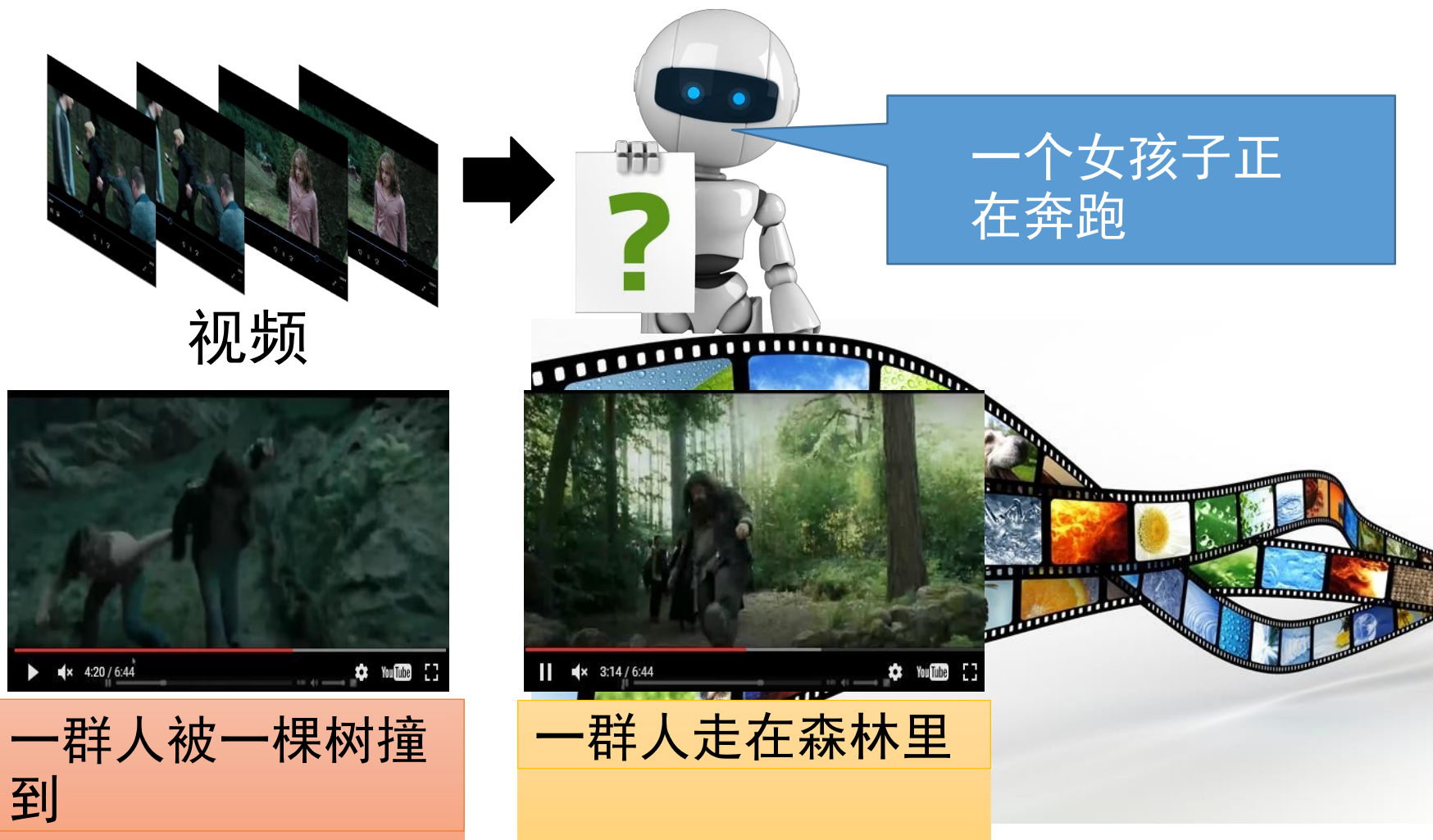
□ [Selfie Lisa !](#)



I am not really confident, but I think it's Mona Lisa et al. taking a selfie.



应用：视频字幕生成



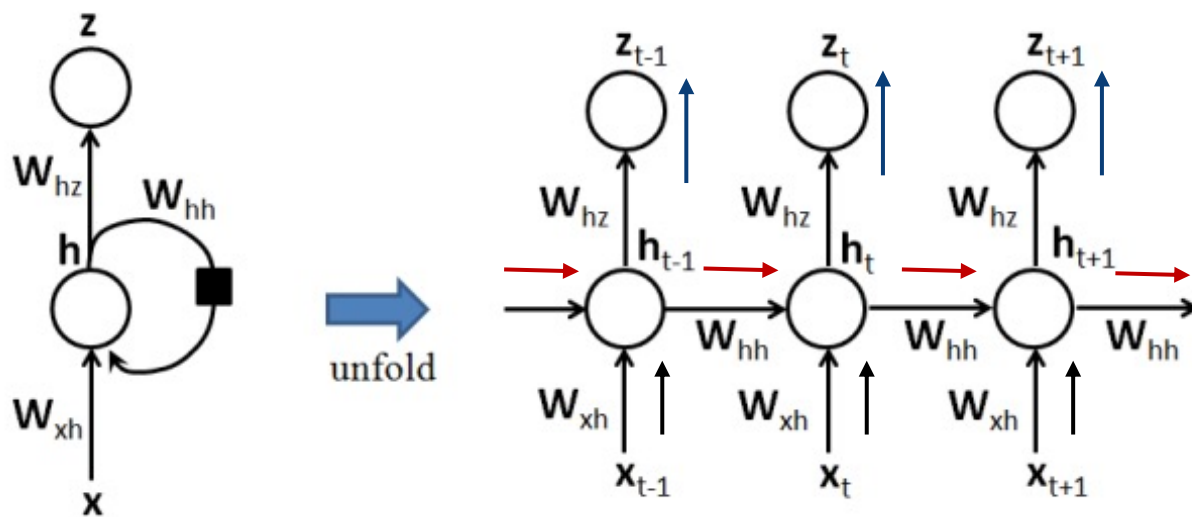
Demo

□ Sketch-RNN

- https://magenta.tensorflow.org/assets/sketch_rnn_demo/index.html
- Cat and Fish

循环神经网络模型

□ 时序扩展



$$\mathbf{h}_t = \tanh(\mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h)$$

$$\mathbf{z}_t = \text{softmax}(\mathbf{W}_{hz}\mathbf{h}_t + \mathbf{b}_z)$$

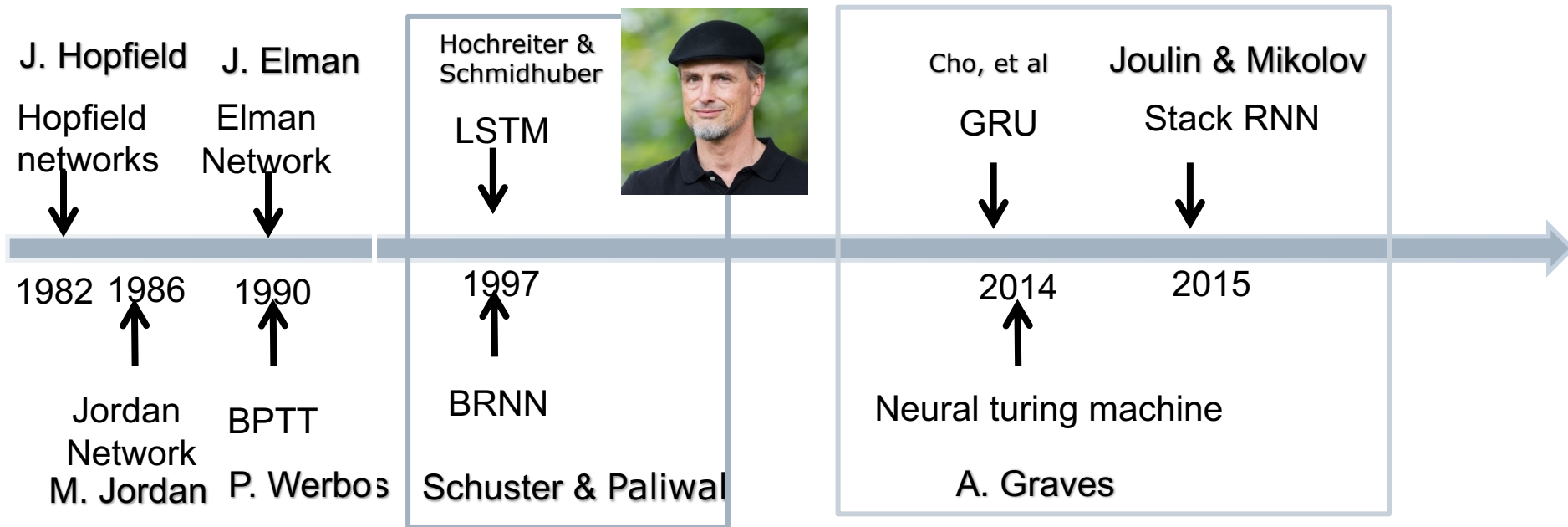
Recurrent Neural Network

RNN是一类扩展的人工神经网络，它是为了对序列数据进行建模而产生的。

- 针对对象：序列数据。例如文本，是字母和词汇的序列；语音，是音节的序列；视频，是图像的序列；气象观测数据，股票交易数据等等，也都是序列数据。
- 核心思想：样本间存在顺序关系，每个样本和它之前的样本存在关联。通过神经网络在时序上的展开，我们能够找到样本之间的序列相关性。

$$h_t = \mathcal{H}(W_{ih}x_t + W_{hh}h_{t-1} + b_h)$$

RNN发展历史



早期（80、90年代）
主要思想：重新使用
参数和计算

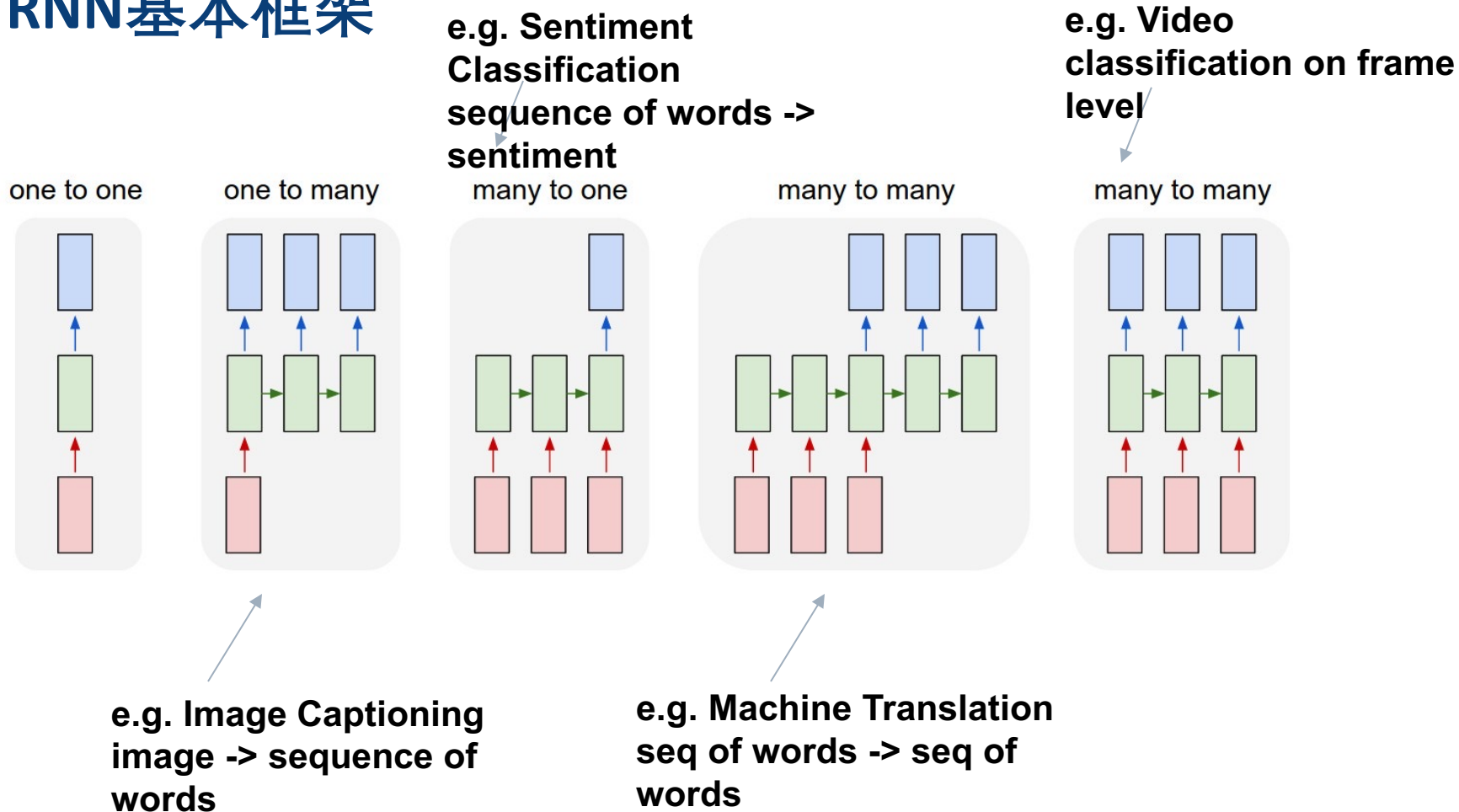
中期（90-2010）
除LSTM以外，RNN
基本从主流研究中消
失了。

当前（2010 - ）应用
广泛：
自然语言应用
视频建模，手写识别
，用户意图预测

开源工具包：
Theano
Torch
PyBrain
TensorFlow
///

Recurrent Neural Network

□ RNN基本框架

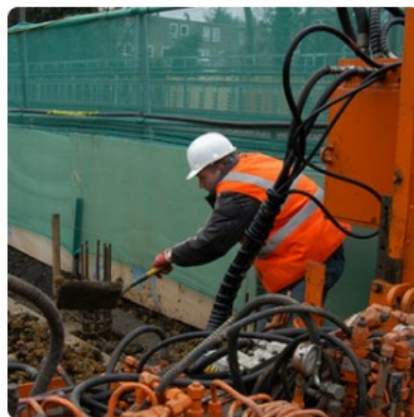


Recurrent Neural Network

□ 典型应用：图像标注



"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."

Recurrent Neural Network

□ 典型应用：语言生成

PANDARUS:

Alas, I think he shall be come approached and the day
When little strain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and
my fair nudes begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

```

/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clearl(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECOND
    return segtable;
}

```

Recurrent Neural Network

□ 典型应用：音乐作曲

Lea Oxlee

The musical score is written in G major (one sharp) and 4/4 time. It consists of four staves of music. The first staff has a treble clef and a key signature of one sharp (F#). The melody starts with a quarter note G, followed by quarter notes A and B, then a quarter rest, and finally a quarter note G. The second staff continues the melody with quarter notes A, B, C, and D, followed by quarter notes E, F#, and G. The third staff has a bass clef and continues the melody with quarter notes G, F#, E, and D, followed by quarter notes C, B, and A. The fourth staff continues the melody with quarter notes G, F#, E, and D, followed by quarter notes C, B, and A. The score includes various chords and a double bar line with repeat dots at the end of the second and third staves.

Chords: A, D, G, A, D, Em, A7, D, Em, F#, E7, A7, A, E7, A, B, A, D, A, D, A, D#m, A7, D



循环神经网络模型

□ softmax

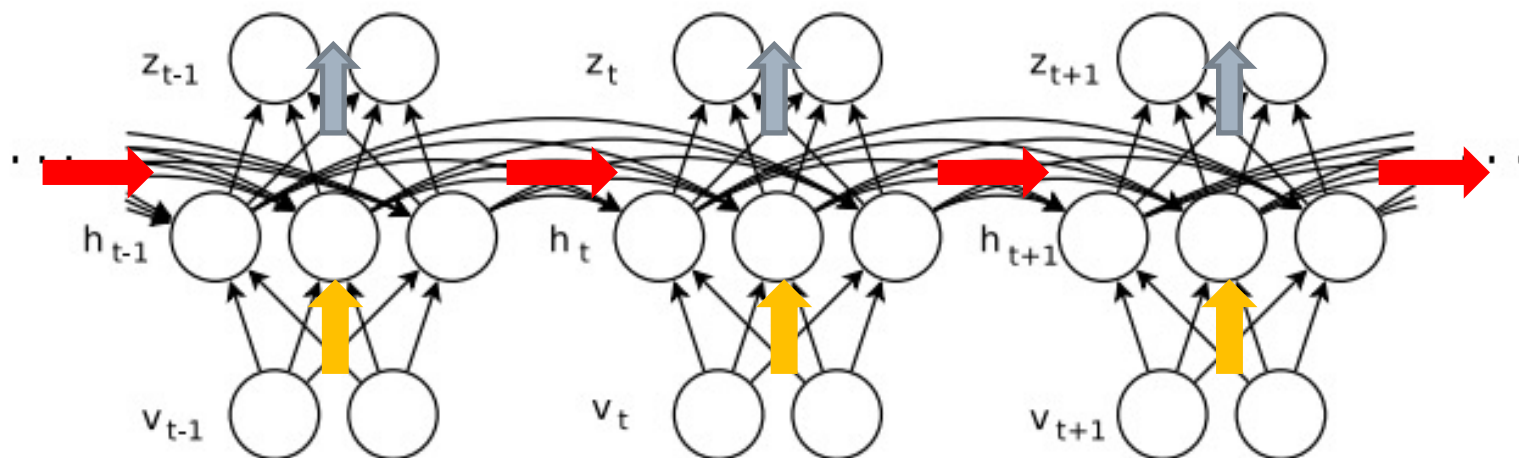
- Softmax函数是sigmoid函数的一个变种，通常我们将其用在多分类任务的输出层，将输入转化成标签的概率。

$$h_{\theta}(x^{(i)}) = \begin{bmatrix} p(y^{(i)} = 1|x^{(i)}; \theta) \\ p(y^{(i)} = 2|x^{(i)}; \theta) \\ \vdots \\ p(y^{(i)} = k|x^{(i)}; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \begin{bmatrix} e^{\theta_1^T x^{(i)}} \\ e^{\theta_2^T x^{(i)}} \\ \vdots \\ e^{\theta_k^T x^{(i)}} \end{bmatrix}$$

本质就是将一个K维的任意实数向量压缩（映射）成另一个K维的实数向量，其中向量中的每个元素取值都介于（0，1）之间。

循环神经网络模型

□ 简单循环网络SRN



神经元之间的连接权重在时域上不变。

循环神经网络模型

□ 随时间反向传播算法BPTT

BP回顾：定义损失函数 E 来表示输出 \hat{y} 和真实标签 y 的误差，通过链式法则自顶向下求得 E 对网络权重的偏导。沿梯度的反方向更新权重的值，直到 E 收敛。

BPTT的本质其实和BP很像，就是加上了时序演化。

定义权重 U, V, W 。

$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$\hat{y}_t = \text{softmax}(Vs_t)$$

定义损失函数：

$$E_t(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

$$E(y, \hat{y}) = \sum_t E_t(y_t, \hat{y}_t)$$

$$= - \sum_t y_t \log \hat{y}_t$$

我们将整个序列作为一次训练，所以需要每个时刻的误差进行求和。

循环神经网络模型

□ 随时间反向传播算法BPTT

目前的任务是求 E 对于 U, V, W 的梯度。

定义 E 对于 W 的梯度 (U, V 同理):

$$\frac{\partial E}{\partial W} = \sum_t \frac{\partial E_t}{\partial W}$$

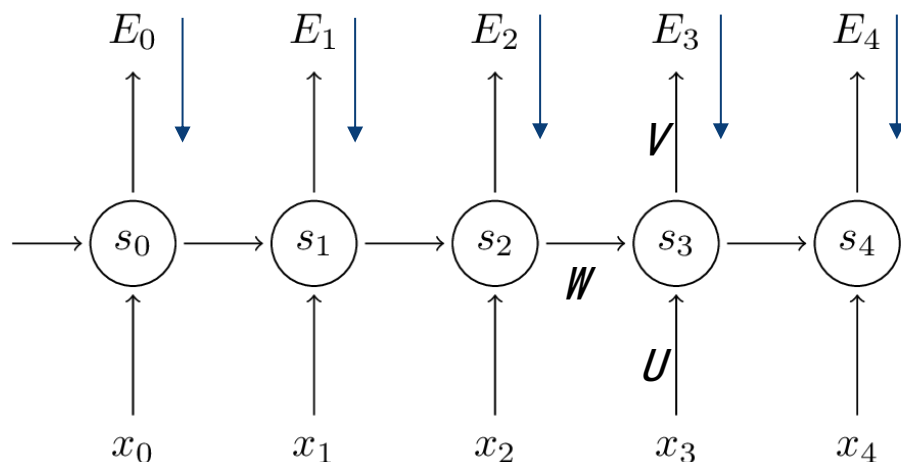
(1) 求 E 对于 V 的梯度。

先求 E_3 对于 V 的梯度:

$$\begin{aligned} \frac{\partial E_3}{\partial V} &= \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial V} \\ &= \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial z_3} \frac{\partial z_3}{\partial V} \end{aligned}$$

其中:

$$z_3 = V s_3$$



$\frac{\partial E}{\partial V}$ 求和可得。

循环神经网络模型

□ 随时间反向传播算法BPTT

(2) 求 E 对于 W 的梯度。注意，现在情况开始变得复杂起来。

先求 E_3 对于 W 的梯度：

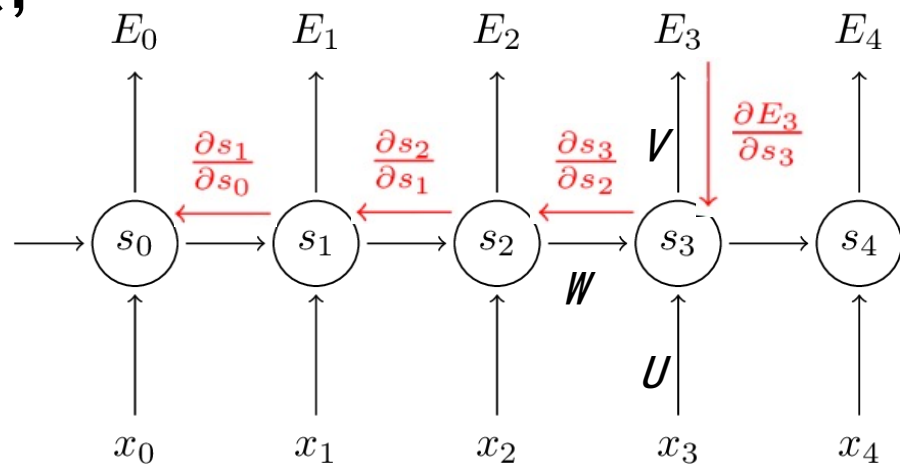
$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial W}$$

当我们求 s_3 对于 W 的偏导时。

注意到： $s_3 = \tanh(Ux_t + Ws_2)$

其中： s_3 依赖于 s_2 ，而 s_2 又依赖于 s_1 和 W ，依赖关系一直传递到 $t = 0$ 的时刻。**因此，当我们计算对于 W 的偏导数时，不能把 s_2 看作是常数项！**

$\frac{\partial E}{\partial W}$ 求和可得。



$$\frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial W}$$

循环神经网络模型

□ 随时间反向传播算法BPTT

(3) 求 E 对于 U 的梯度。

情况与 W 类似。

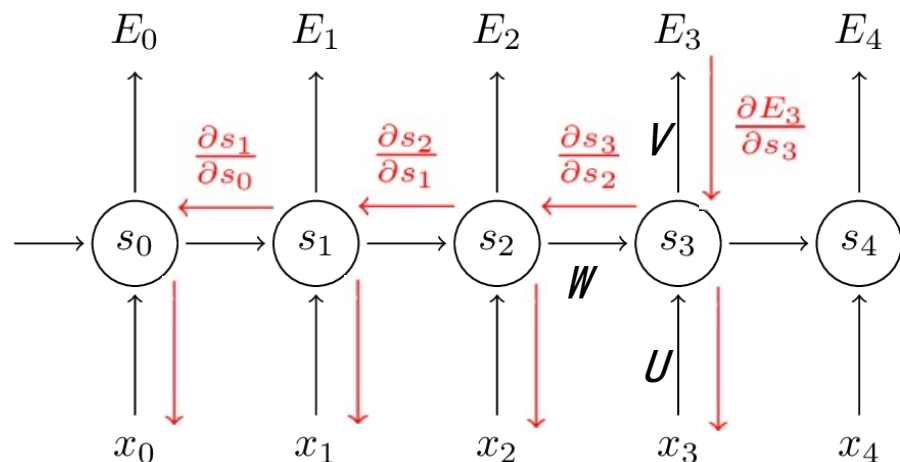
先求 E_3 对于 U 的梯度：

$$\frac{\partial E_3}{\partial U} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial U}$$

当我们求 s_3 对于 W 的偏导时。

注意到： $s_3 = \tanh(Ux_t + Ws_2)$

$\frac{\partial E}{\partial U}$ 求和可得。



$$\frac{\partial E_3}{\partial U} = \sum_{k=1}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial U}$$

同样： s_3 依赖于 s_2 ，而 s_2 又依赖于 s_1 和 U 。

类似求 W ，当我们计算对于 U 的偏导数时，也不能把 s_2 看作是常数项！

循环神经网络模型

□ 随时间反向传播算法BPTT

```

1: for  $t$  from  $T$  downto 1 do
2:    $do_t \leftarrow g'(o_t) \cdot dL(z_t; y_t)/dz_t$ 
3:    $db_o \leftarrow db_o + do_t$ 
4:    $dW_{oh} \leftarrow dW_{oh} + do_t h_t^\top$ 
5:    $dh_t \leftarrow dh_t + W_{oh}^\top do_t$ 
6:    $dz_t \leftarrow e'(z_t) \cdot dh_t$ 
7:    $dW_{hv} \leftarrow dW_{hv} + dz_t v_t^\top$ 
8:    $db_h \leftarrow db_h + dz_t$ 
9:    $dW_{hh} \leftarrow dW_{hh} + dz_t h_{t-1}^\top$ 
10:   $dh_{t-1} \leftarrow W_{hh}^\top dz_t$ 
11: end for
12: Return  $d\theta = [dW_{hv}, dW_{hh}, dW_{oh}, db_h, db_o, dh_0]$ .

```

参数意义:

W_{hv} : 输入层到隐含层的权重参数,

W_{hh} : 隐含层到隐含层的权重参数,

W_{oh} : 隐含层到输出层的权重参数,

b_h : 隐含层的偏移量,

b_o : 输出层的偏移量,

h_0 : 起始状态的隐含层的输出,
一般初始为0。

递归神经网络模型

□ 随时间反向传播算法BPTT

面临的问题：

- 梯度消失问题
- 梯度爆炸问题

解决方案：

- 选择其他的激活函数。例如ReLU。
- 引入改进网络结构的机制，例如LSTM，GRU。
- 现在在自然语言处理上应用十分广的就是LSTM。



End

