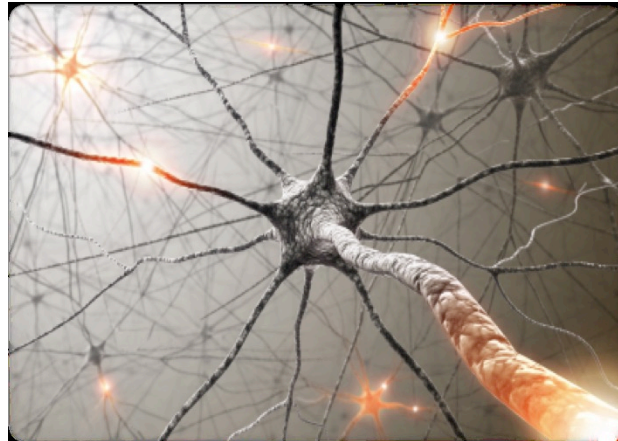


05

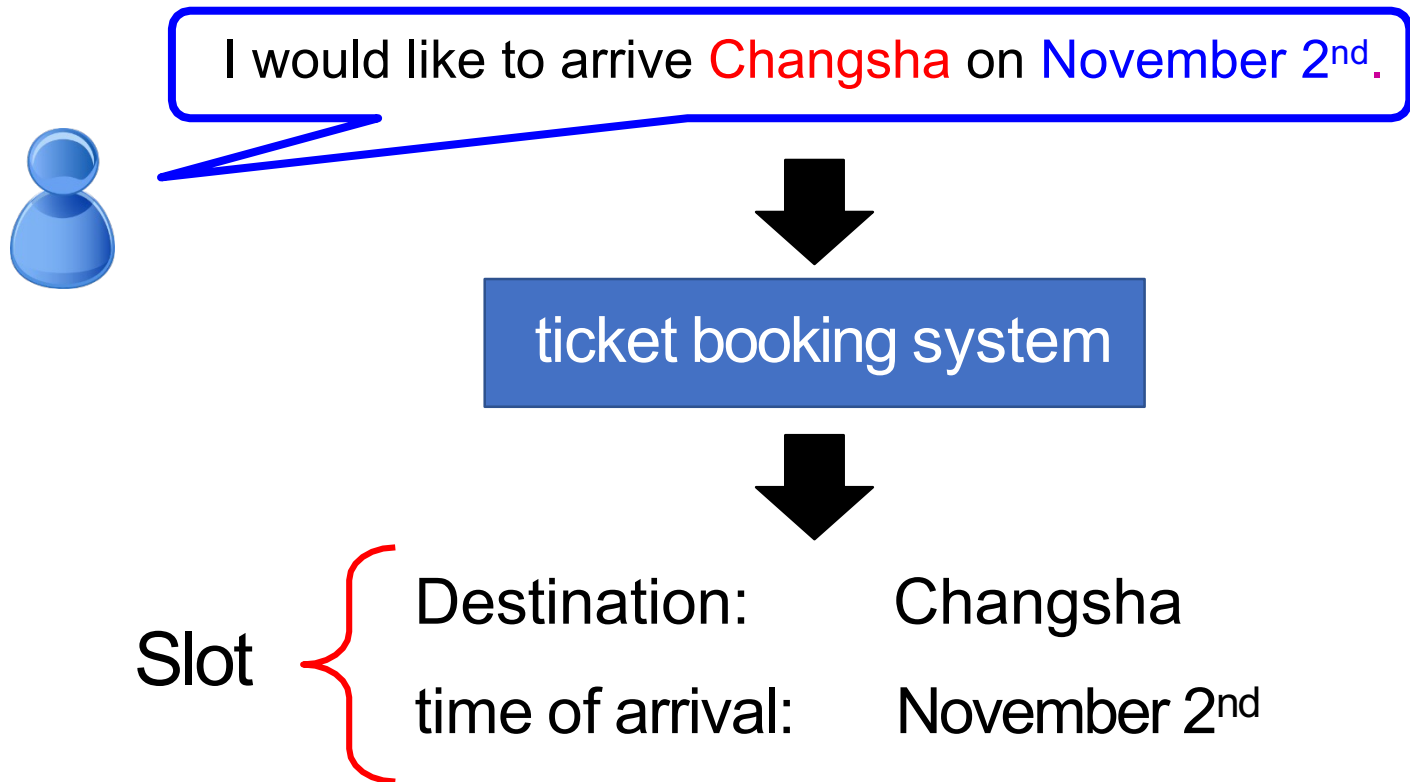
Neural Networks

Recurrent Neural Networks (RNNs)



Example Application

- Slot Filling



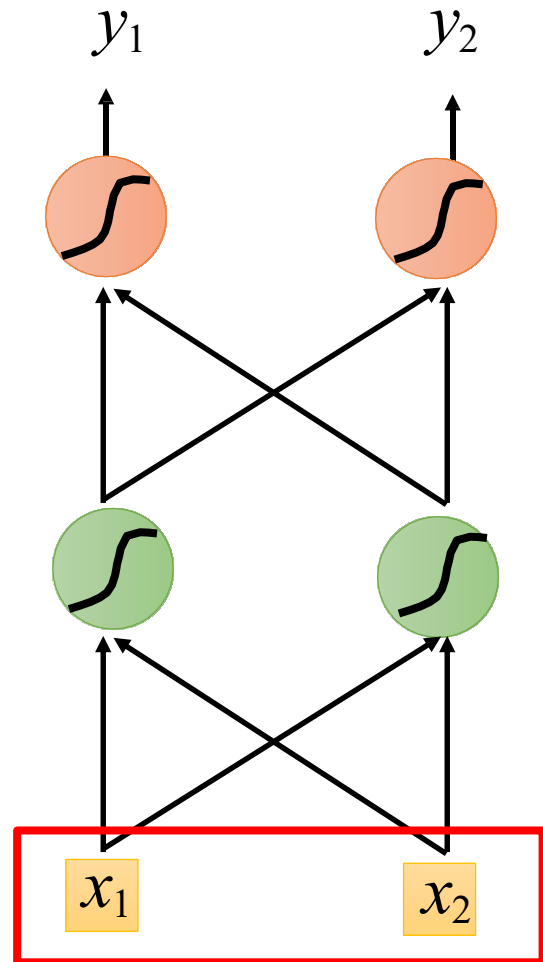
Example Application

Solving slot filling by
Feedforward network?

Input: a word

(Each word is represented as a
vector)

Changsha



1-of-N encoding

How to represent each word as a vector?

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

The vector is lexicon size.

Each dimension corresponds
to a word in the lexicon

The dimension for the word is
1, and others are 0

apple = [1 0 0 0 0]

bag = [0 1 0 0 0]

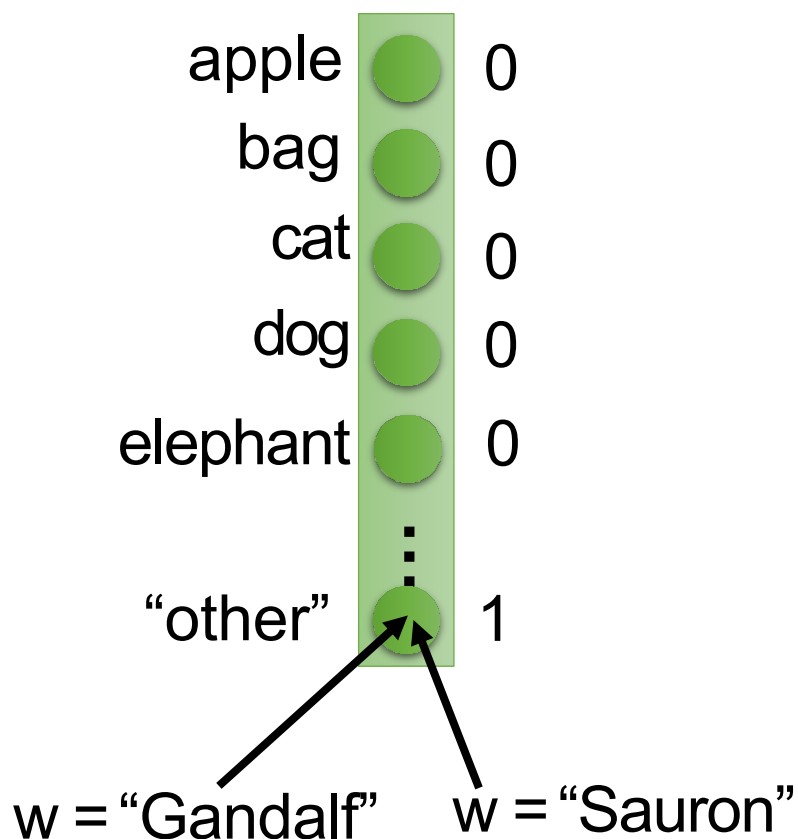
cat = [0 0 1 0 0]

dog = [0 0 0 1 0]

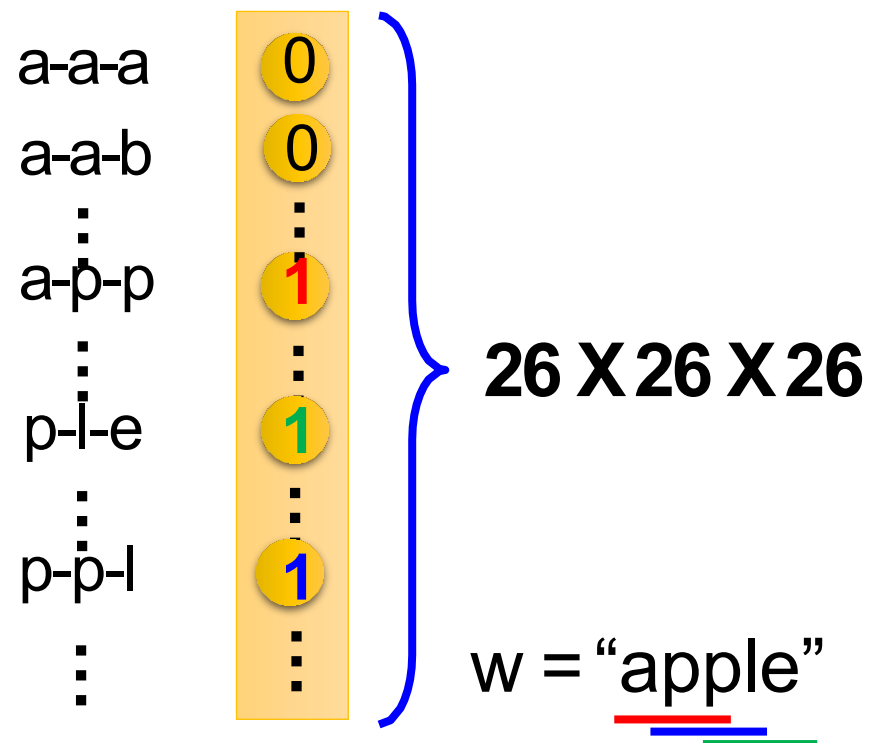
elephant = [0 0 0 0 1]

Beyond 1-of-N encoding

Dimension for “Other”



Word hashing



Example Application

Solving slot filling by
Feedforward network?

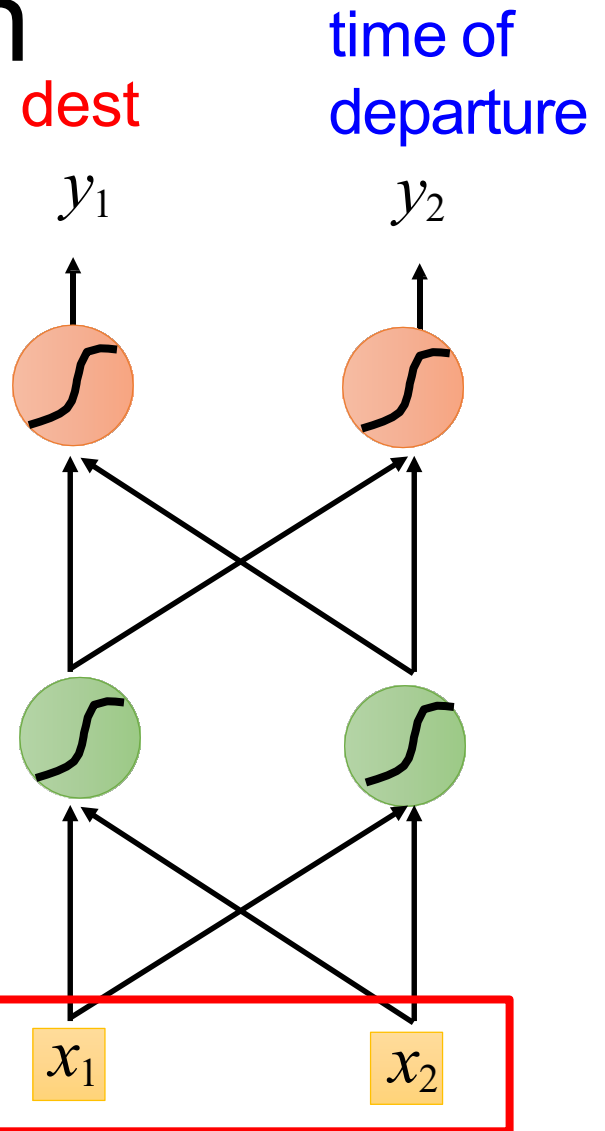
Input: a word

(Each word is represented
as a vector)

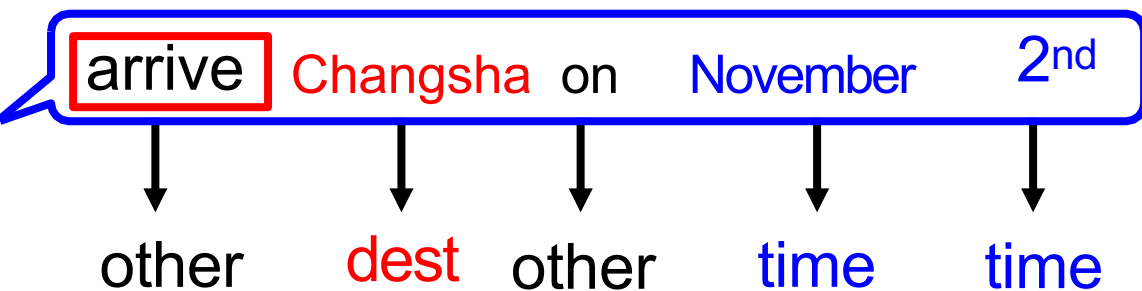
Output:

Probability distribution that the
input word belonging to the
slots

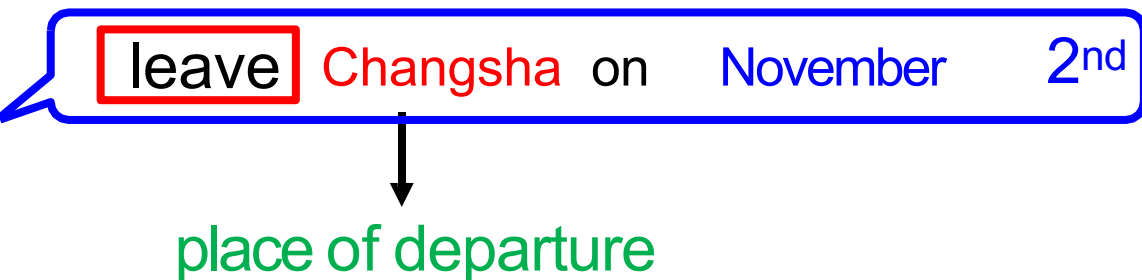
Changsha



Example Application

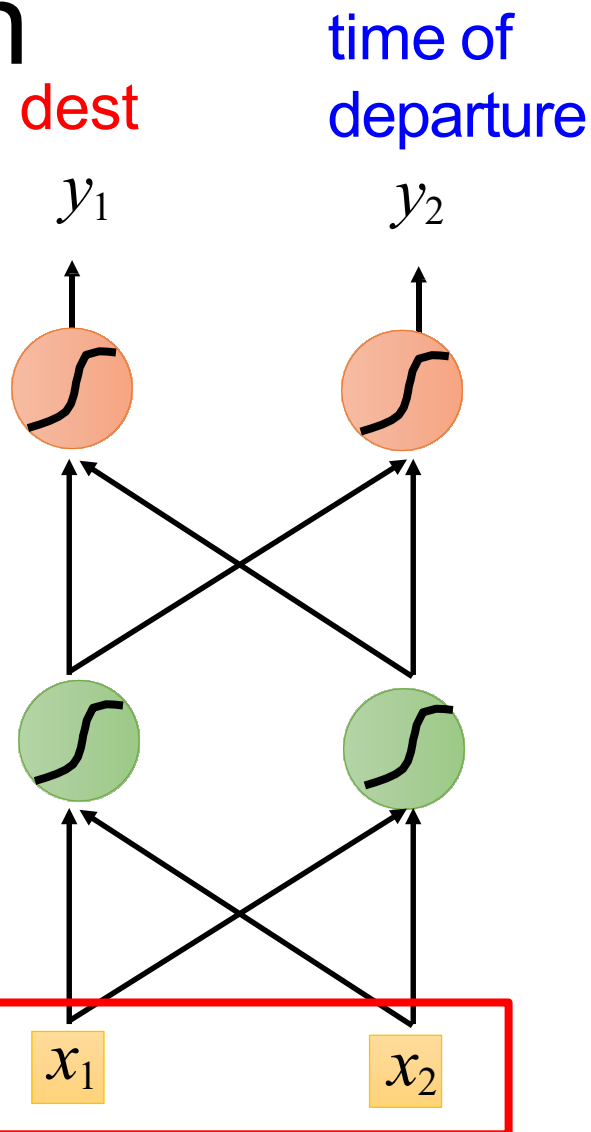


Problem?



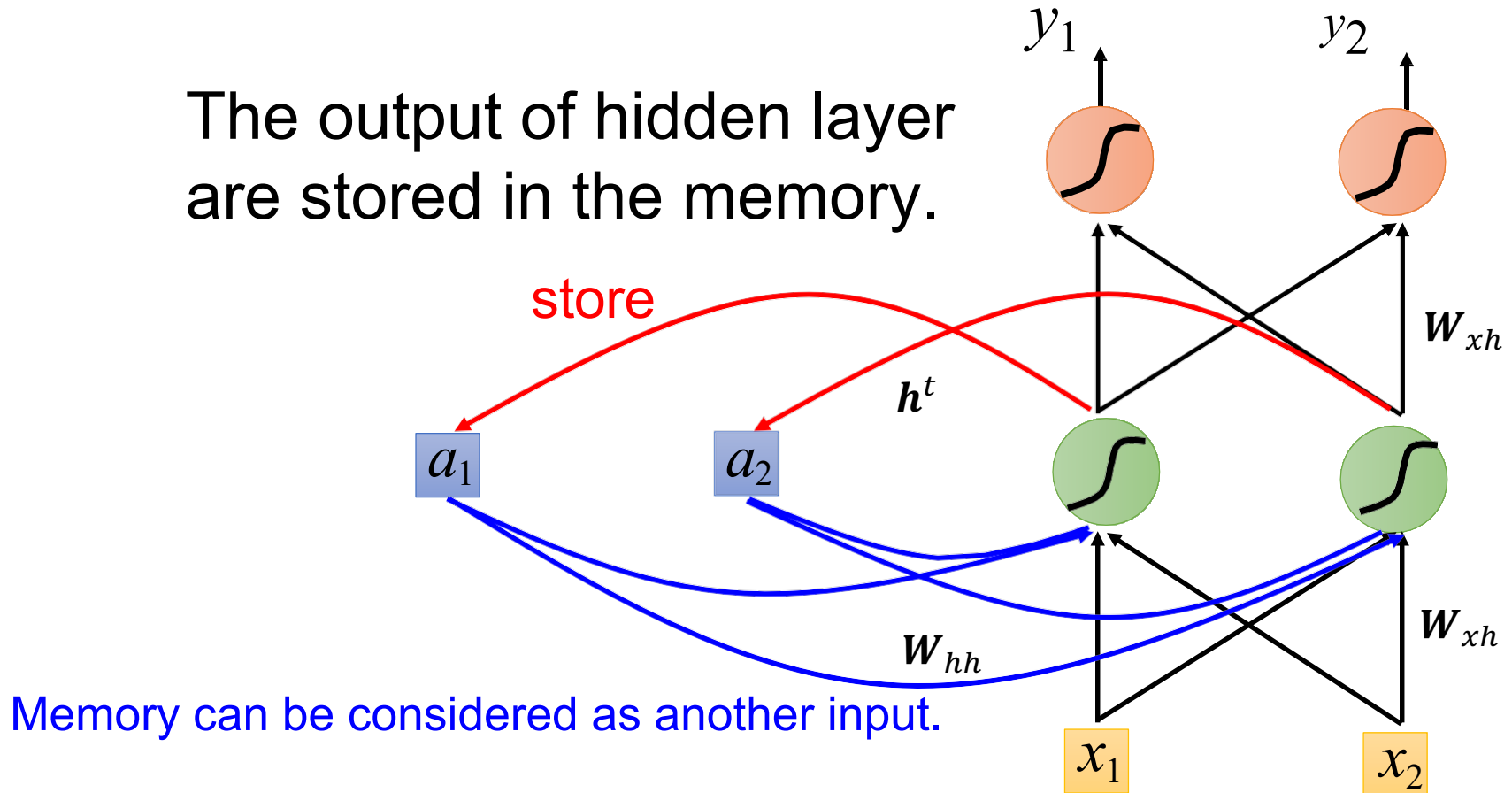
Neural network
needs memory!

Changsha



Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.

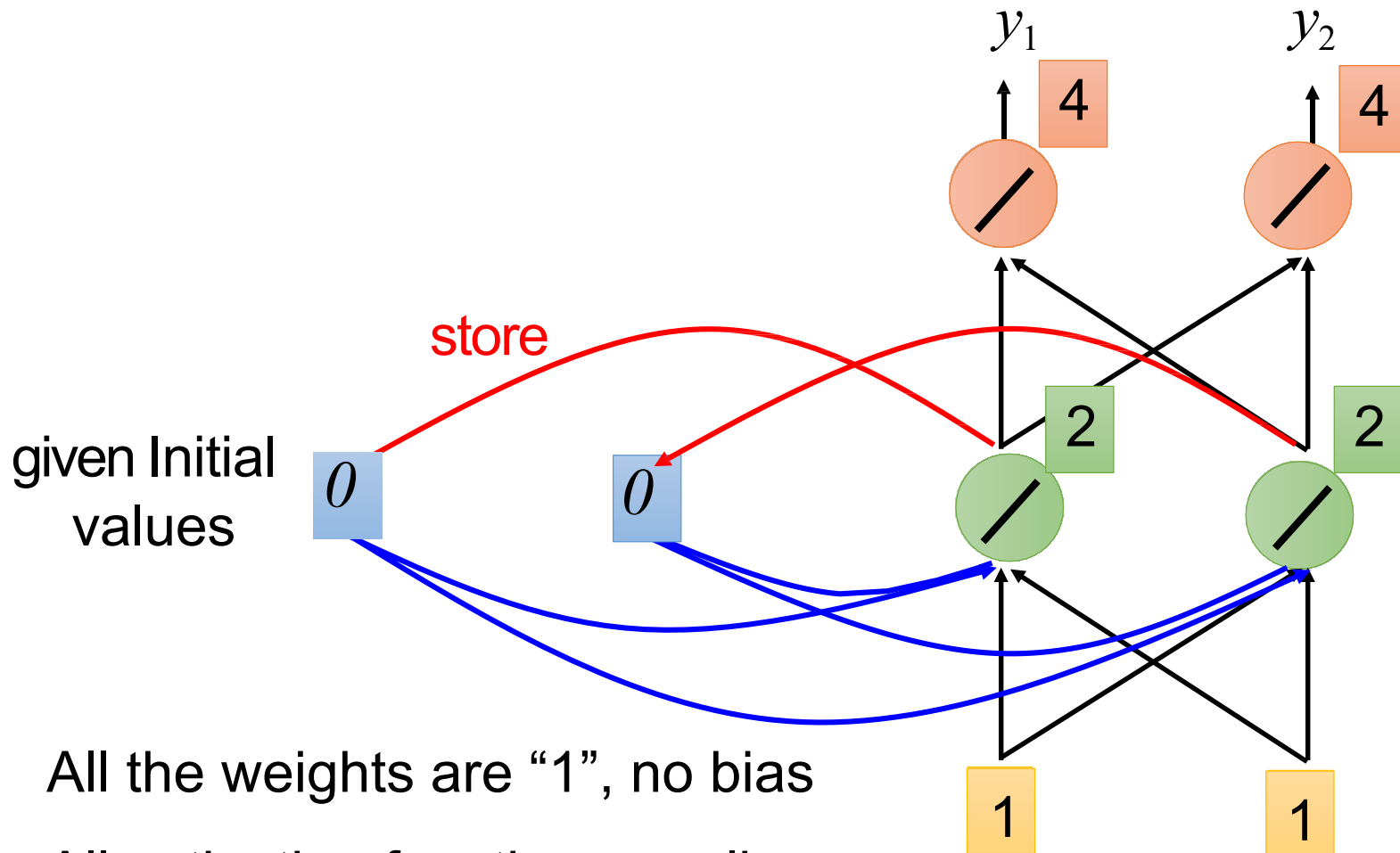


$$\begin{cases} \mathbf{h}^t = f(\mathbf{W}_{xh}\mathbf{x}^t + \mathbf{W}_{hh}\mathbf{h}^{t-1} + \mathbf{b}_h) \\ \mathbf{y}^t = f(\mathbf{W}_{hy}\mathbf{h}^t + \mathbf{b}_y) \end{cases}$$

Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$

Output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$



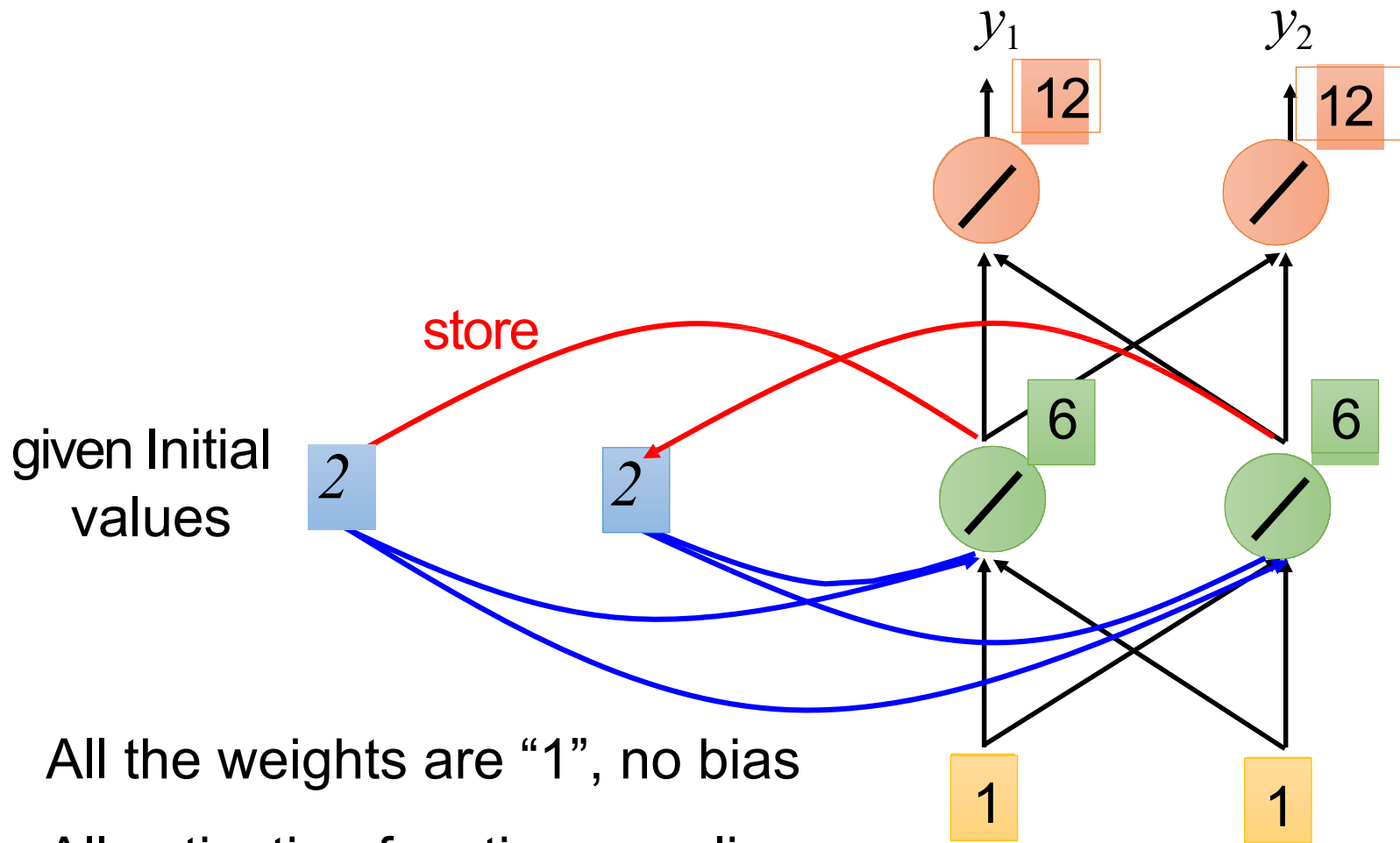
All the weights are “1”, no bias

All activation functions are linear

Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$

Output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$ $\begin{bmatrix} 12 \\ 12 \end{bmatrix}$



All the weights are "1", no bias

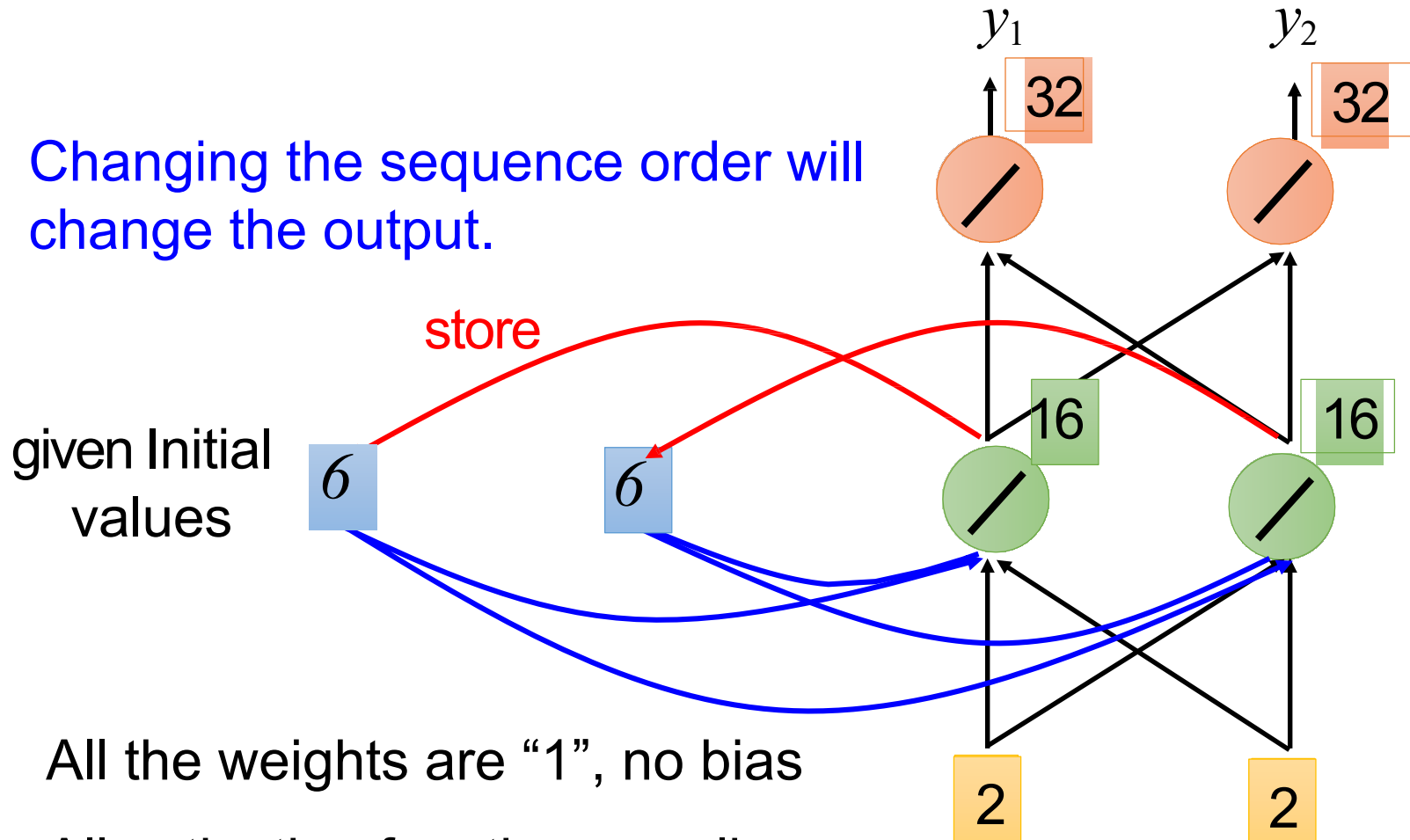
All activation functions are linear

Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$

Output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$ $\begin{bmatrix} 12 \\ 12 \end{bmatrix}$ $\begin{bmatrix} 32 \\ 32 \end{bmatrix}$

Changing the sequence order will change the output.



All the weights are "1", no bias

All activation functions are linear

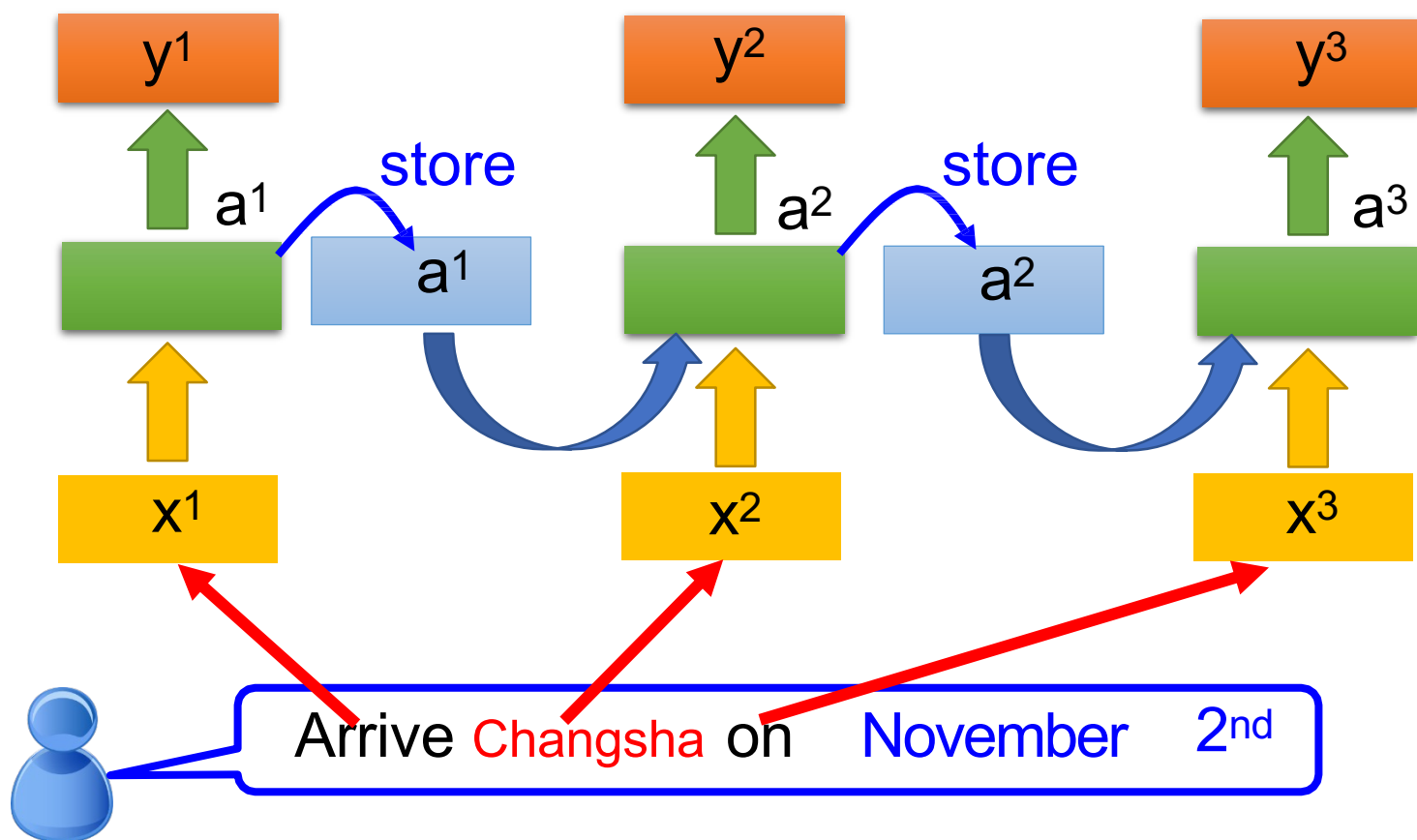
RNN

The same network is used again and again.

Probability of
“arrive” in each slot

Probability of
“**Changsha**” in each slot

Probability of
“on” in each slot



RNN

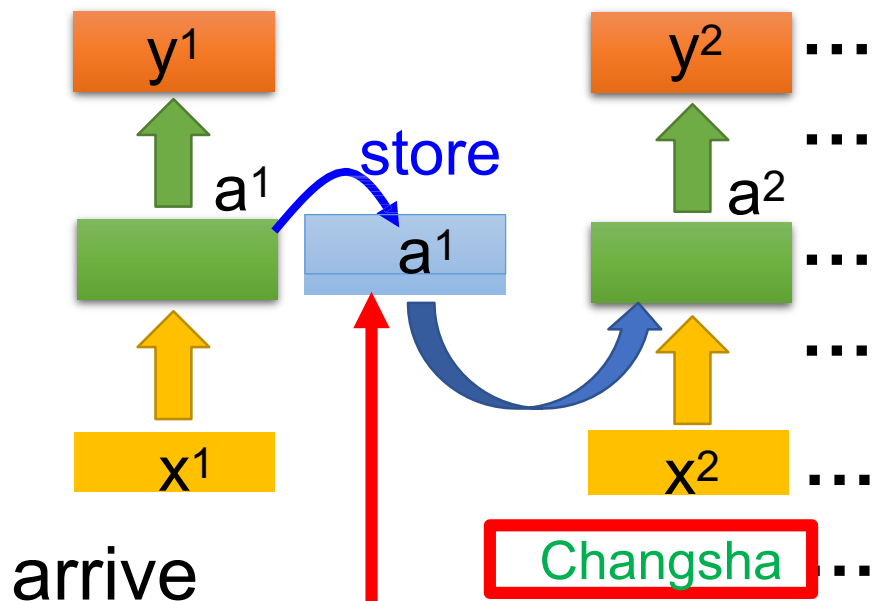
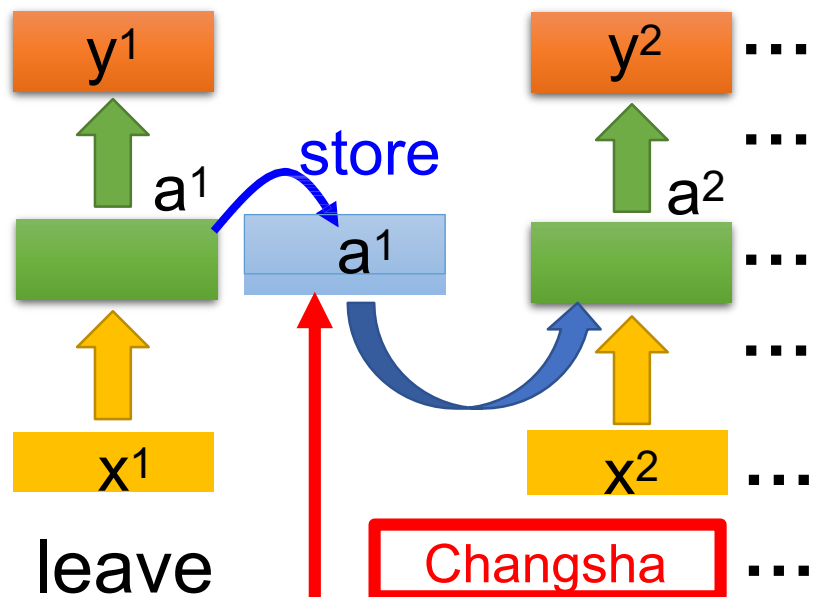
Different

Prob of “leave”
in each slot

Prob of “Changsha”
in each slot

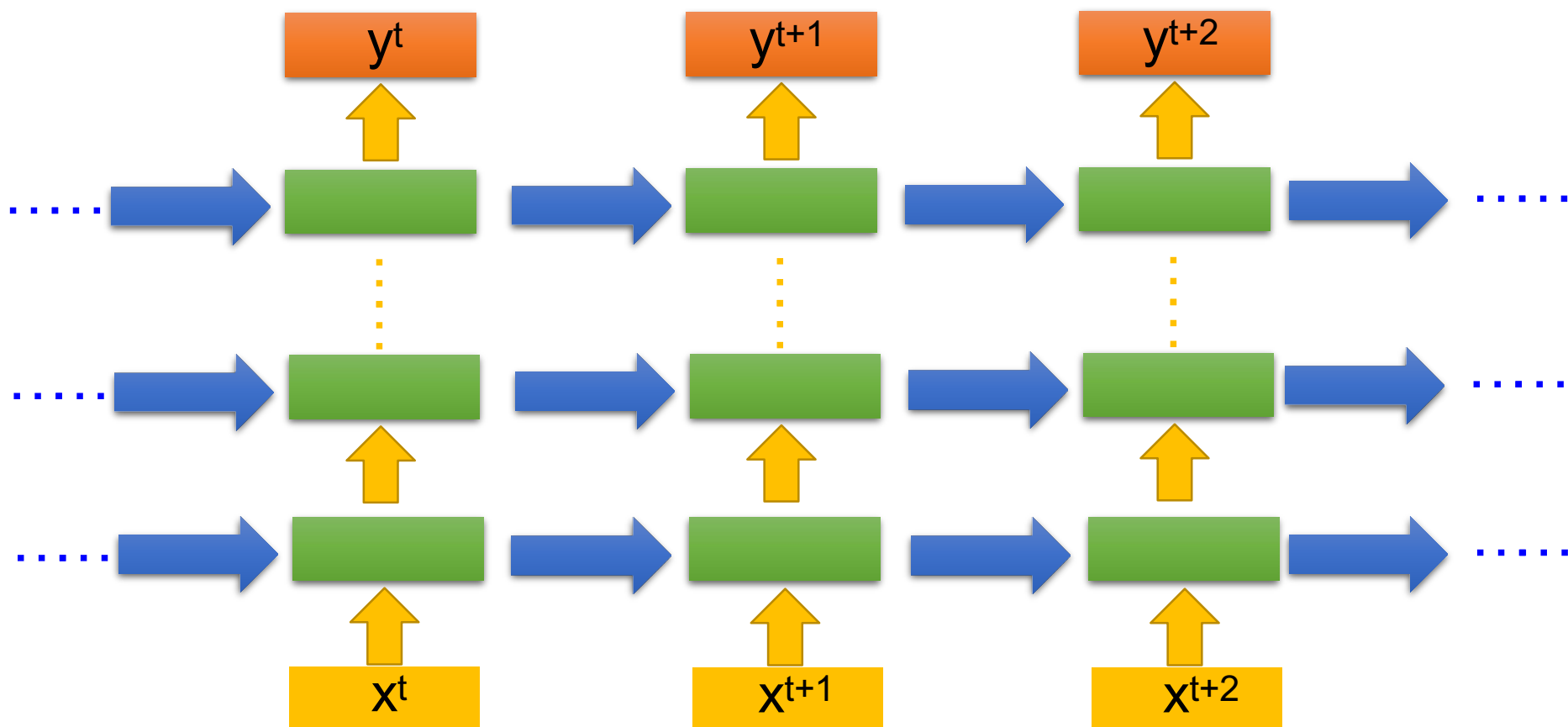
Prob of “arrive”
in each slot

Prob of “Changsha”
in each slot



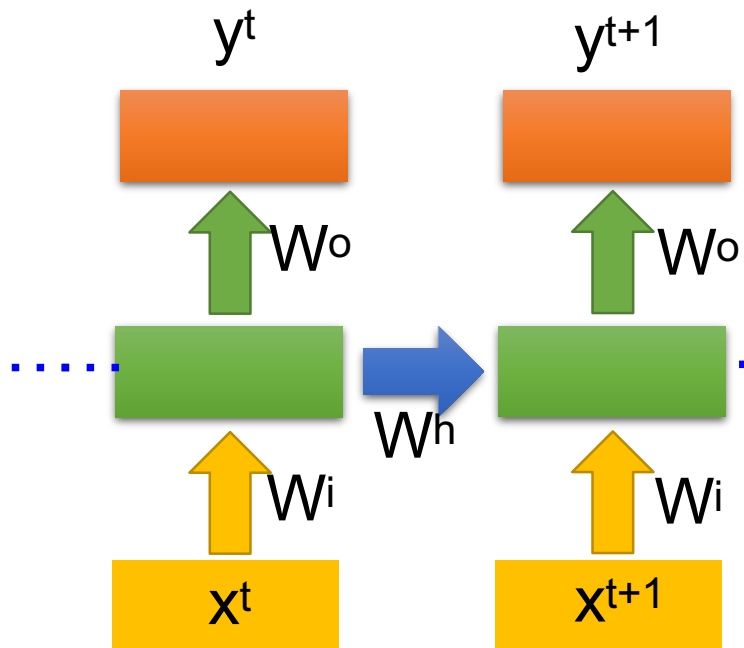
The values stored in the memory is different.

Of course it can be deep ...

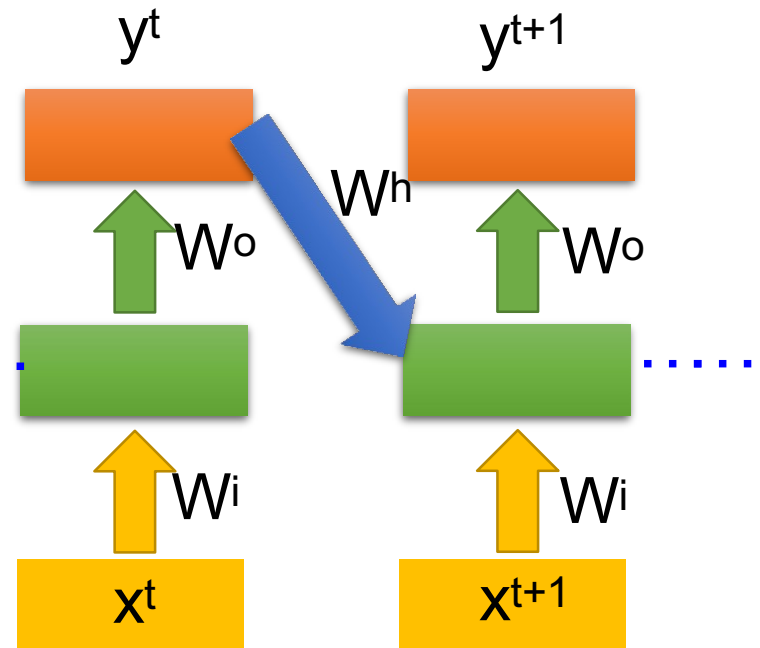


Elman Network & Jordan Network

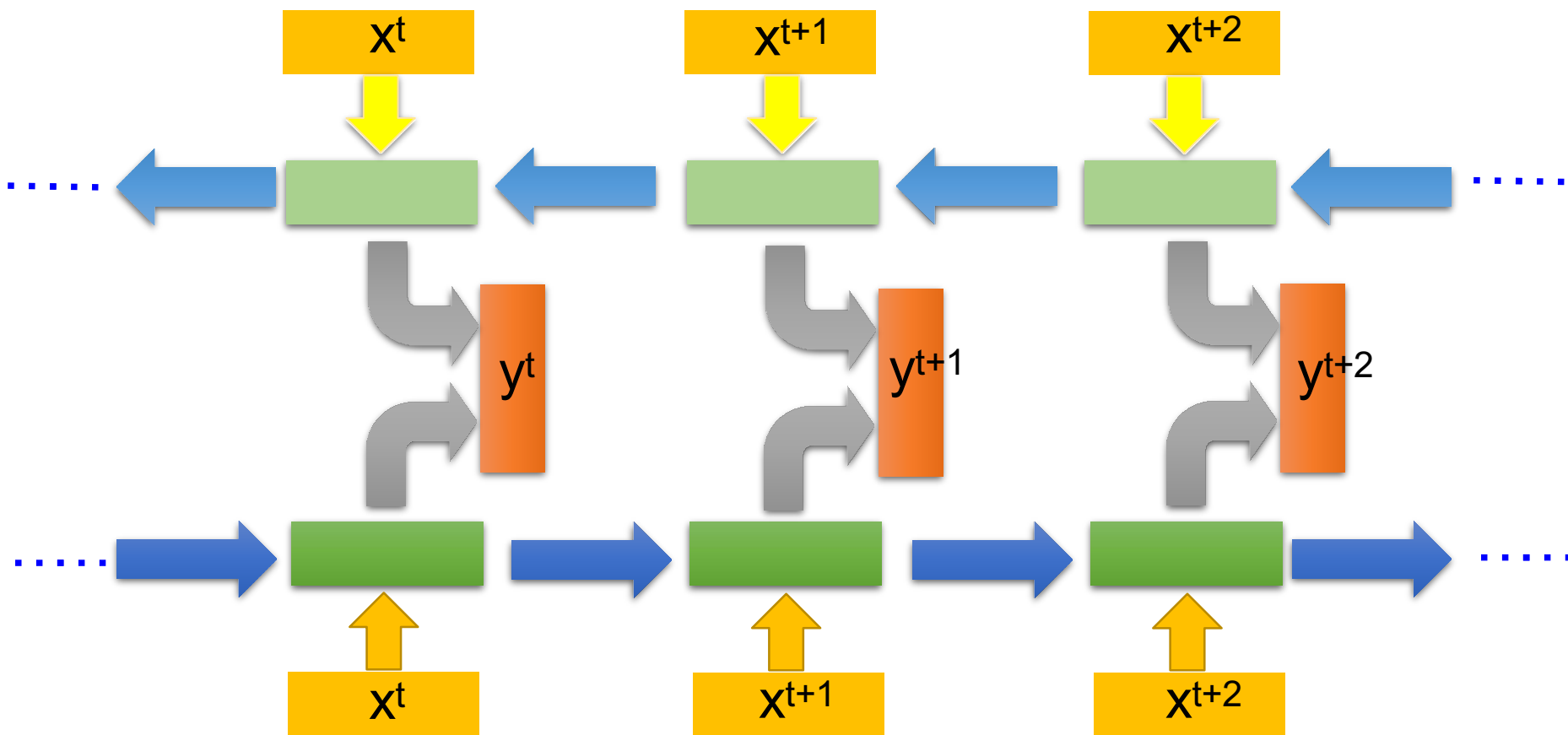
Elman Network



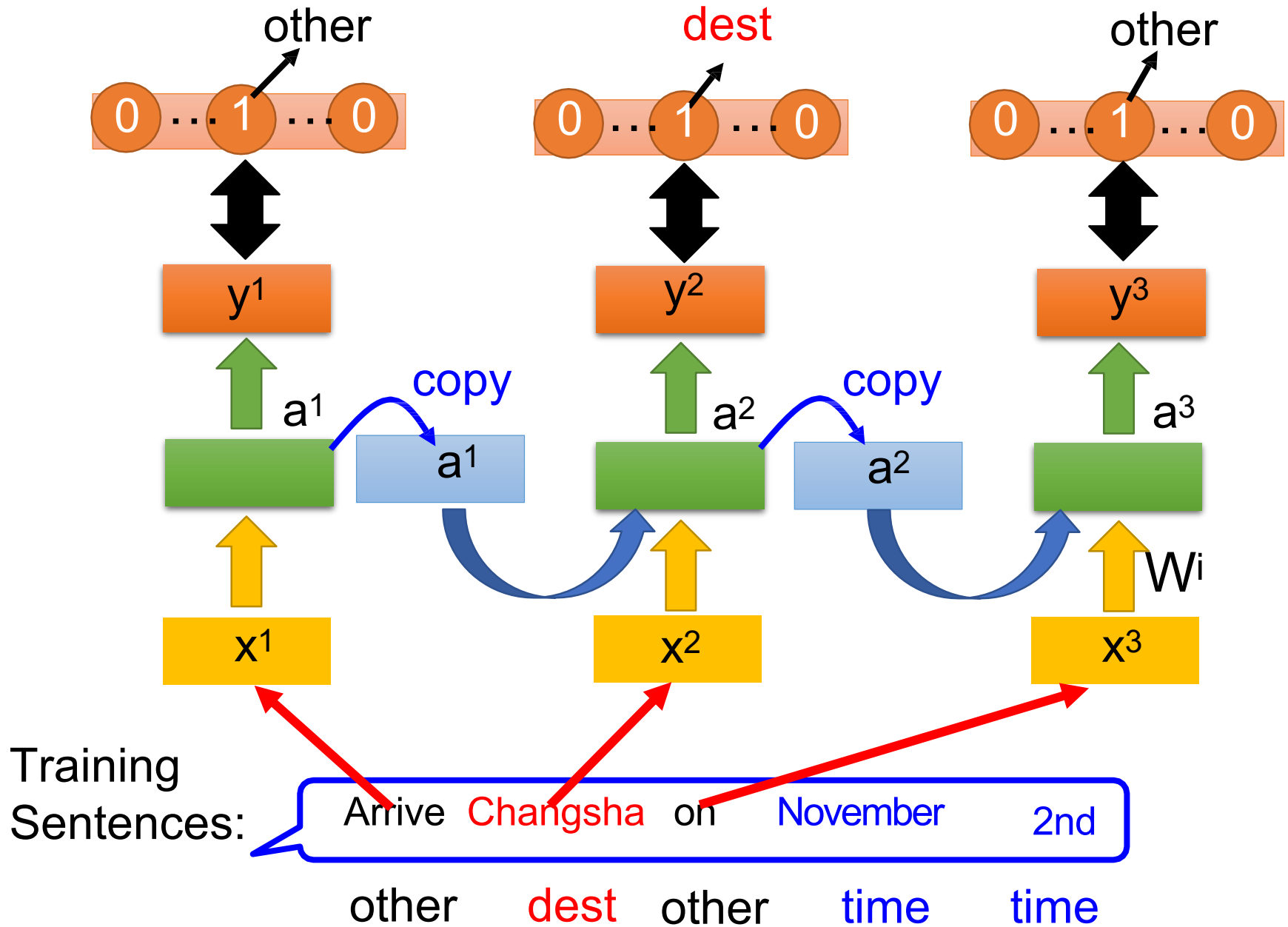
Jordan Network



Bidirectional RNN

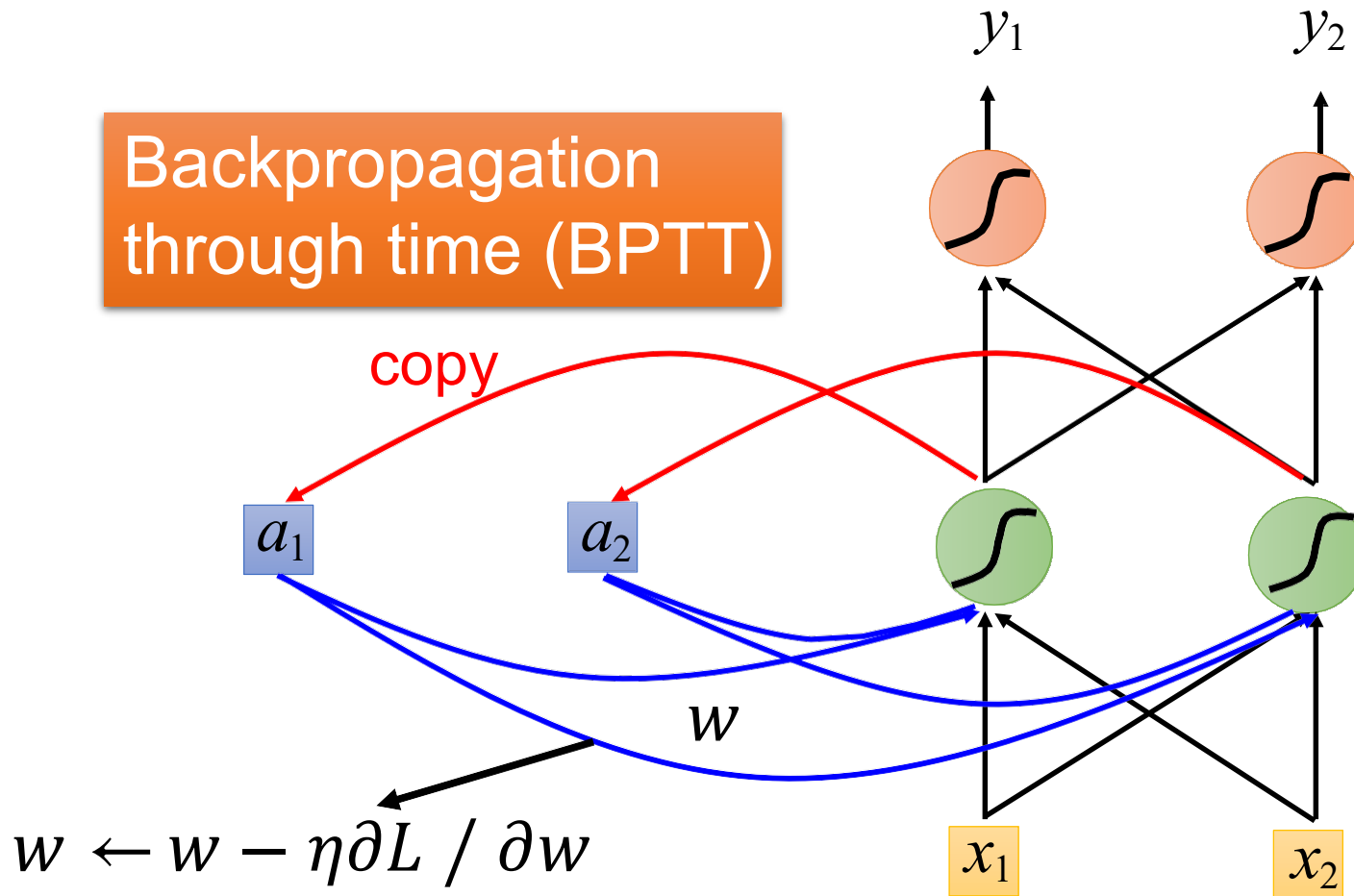


Learning Target



Learning

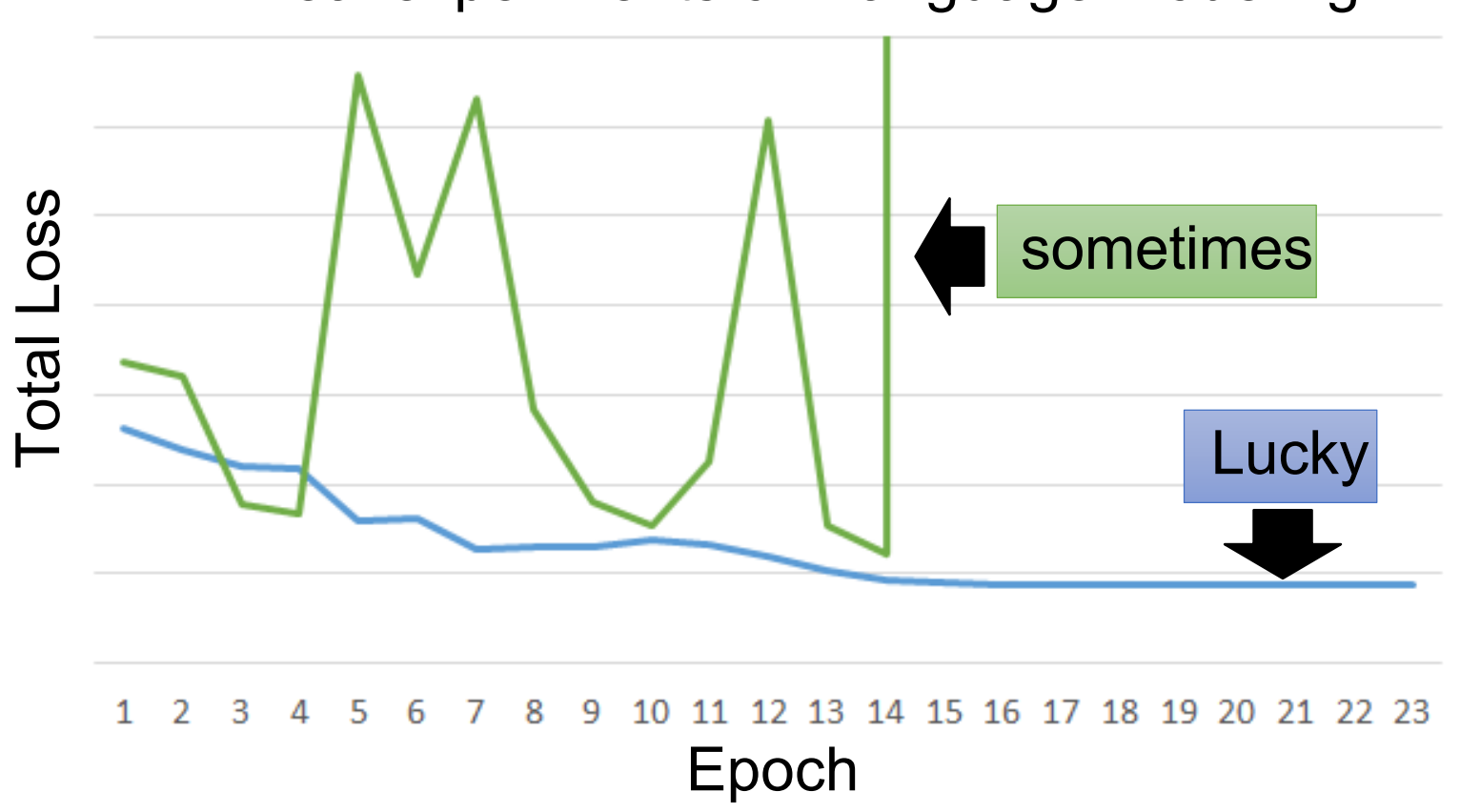
Backpropagation
through time (BPTT)



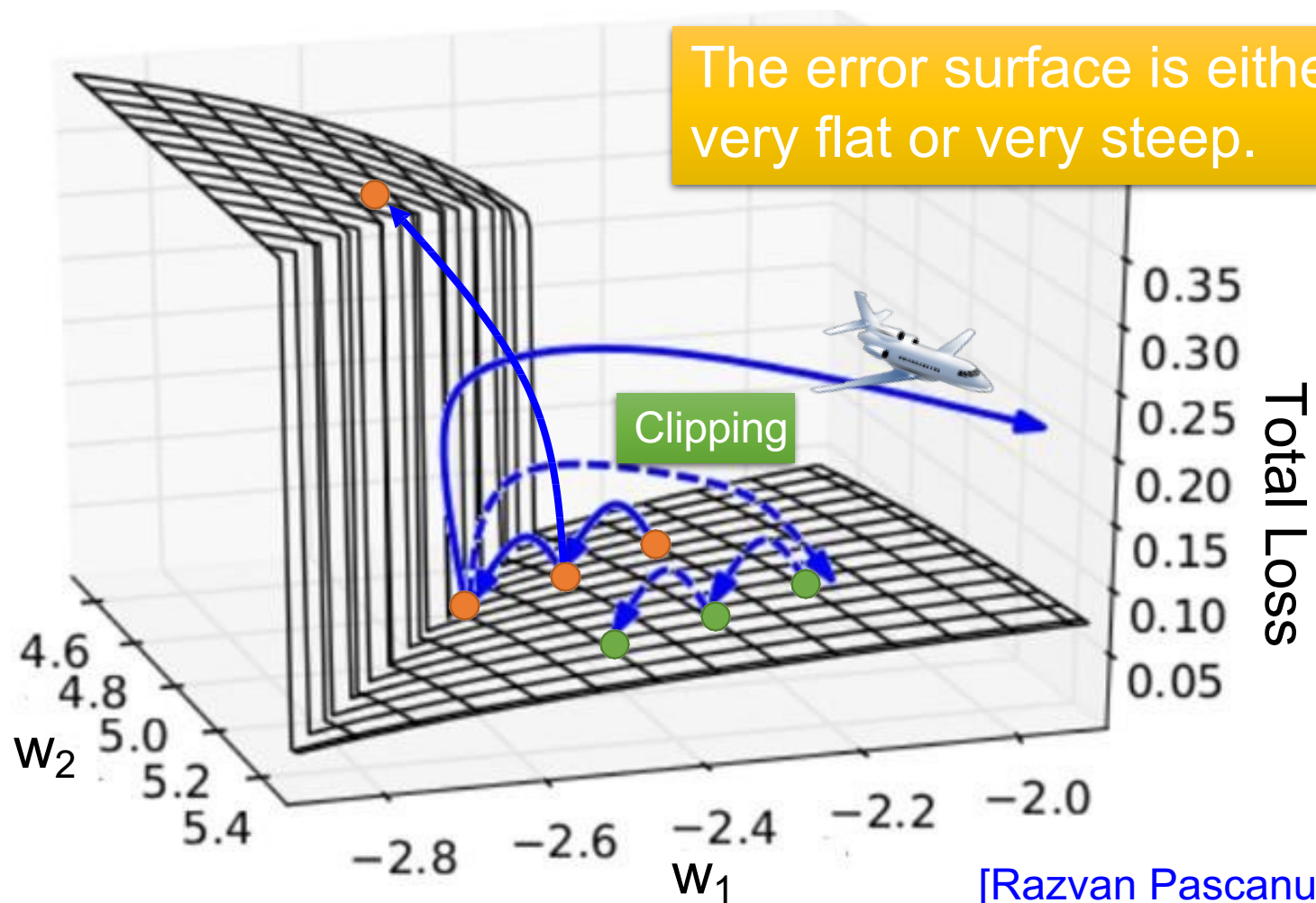
Unfortunately

- RNN-based network is not always easy to learn

Real experiments on Language modeling



The error surface is rough.



[Razvan Pascanu, ICML'13]

Why?

$$w = 1 \quad \longrightarrow \quad y^{1000} = 1$$

$$w = 1.01 \quad \longrightarrow \quad y^{1000} \approx 20000$$

$$w = 0.99 \quad \longrightarrow \quad y^{1000} \approx 0$$

$$w = 0.01 \quad \longrightarrow \quad y^{1000} \approx 0$$

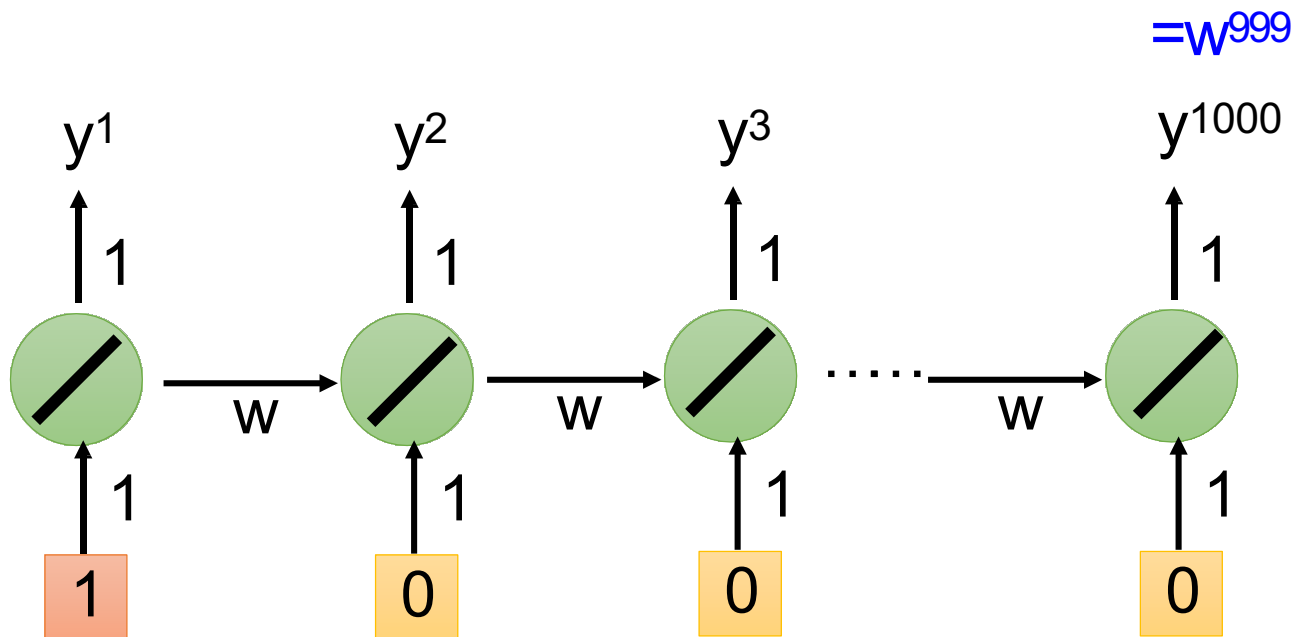
Large
 $\partial L / \partial w$

Small
Learning rate?

small
 $\partial L / \partial w$

Large
Learning rate?

Toy Example



Helpful Techniques

- Long Short-term Memory (LSTM)

- Can deal with gradient vanishing (not gradient explode)

- Memory and input are

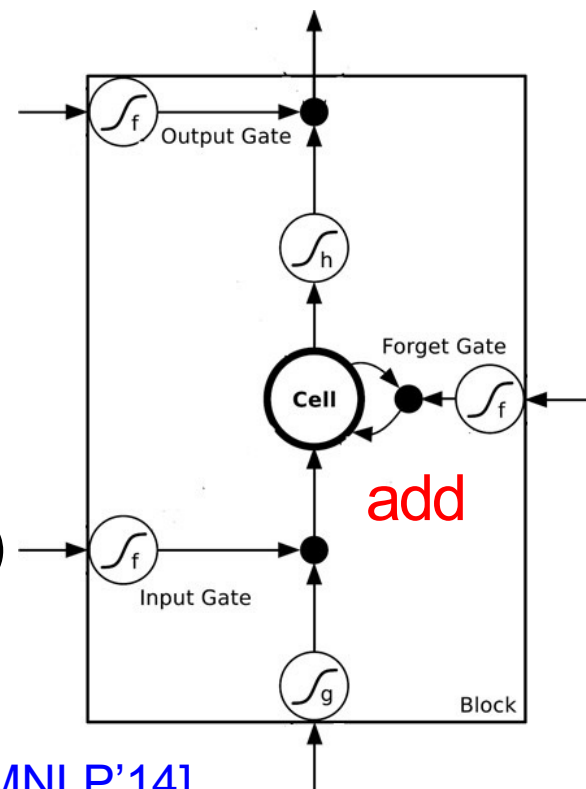
- added**

- The influence never disappears unless forget gate is closed



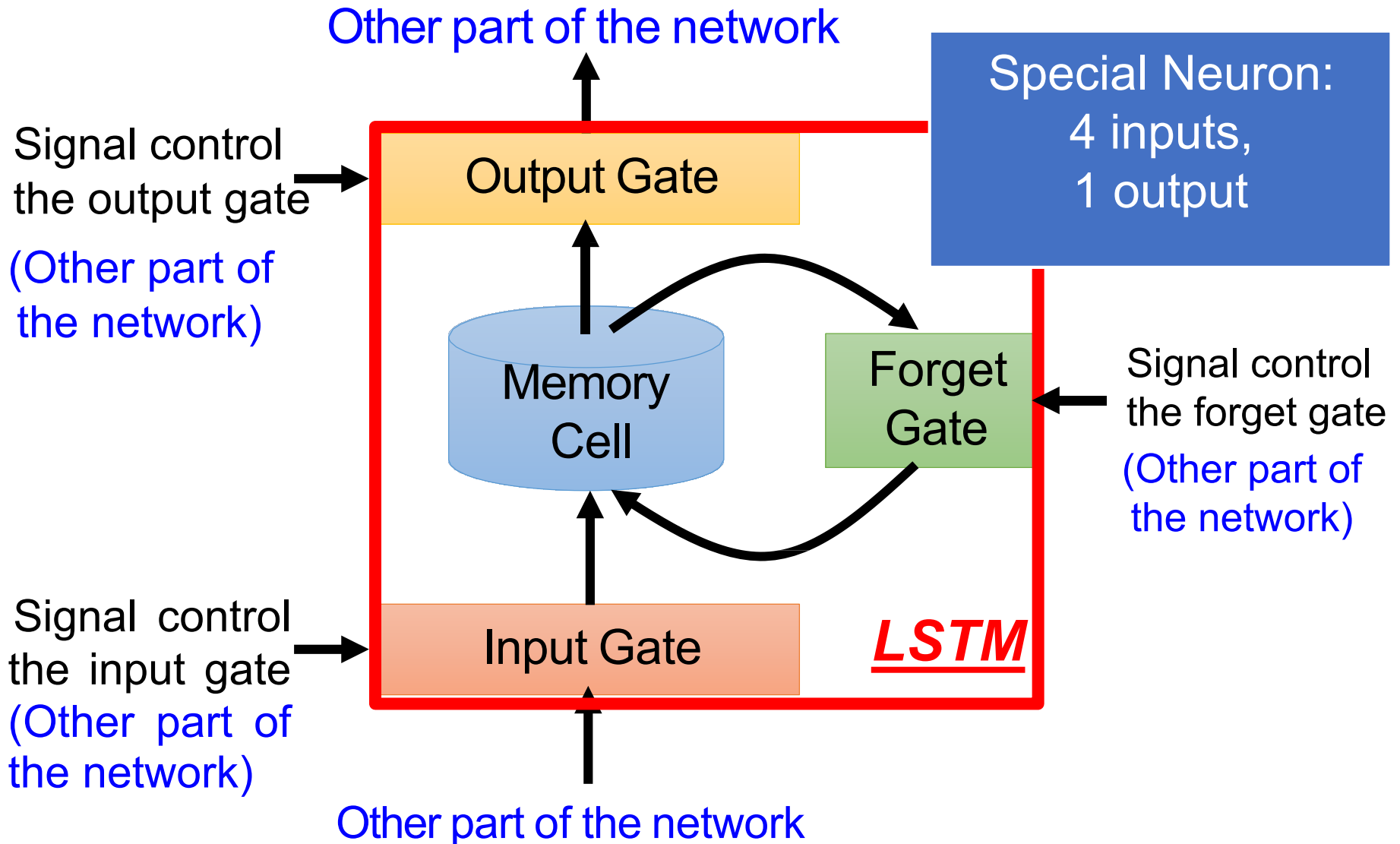
No Gradient vanishing
(If forget gate is opened.)

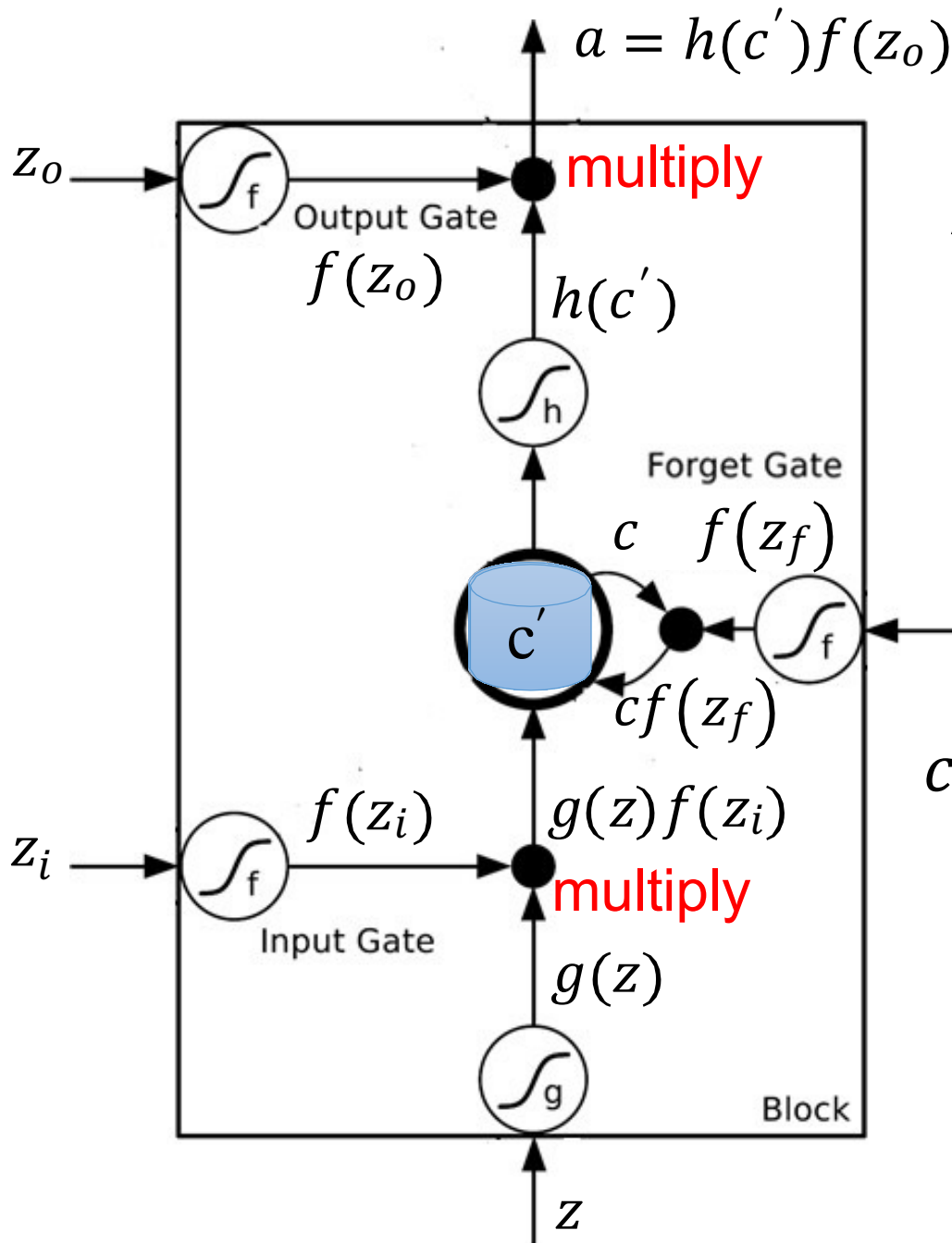
Gated Recurrent Unit (GRU):
simpler than LSTM



[Cho, EMNLP'14]

Long Short-Term Memory (LSTM)





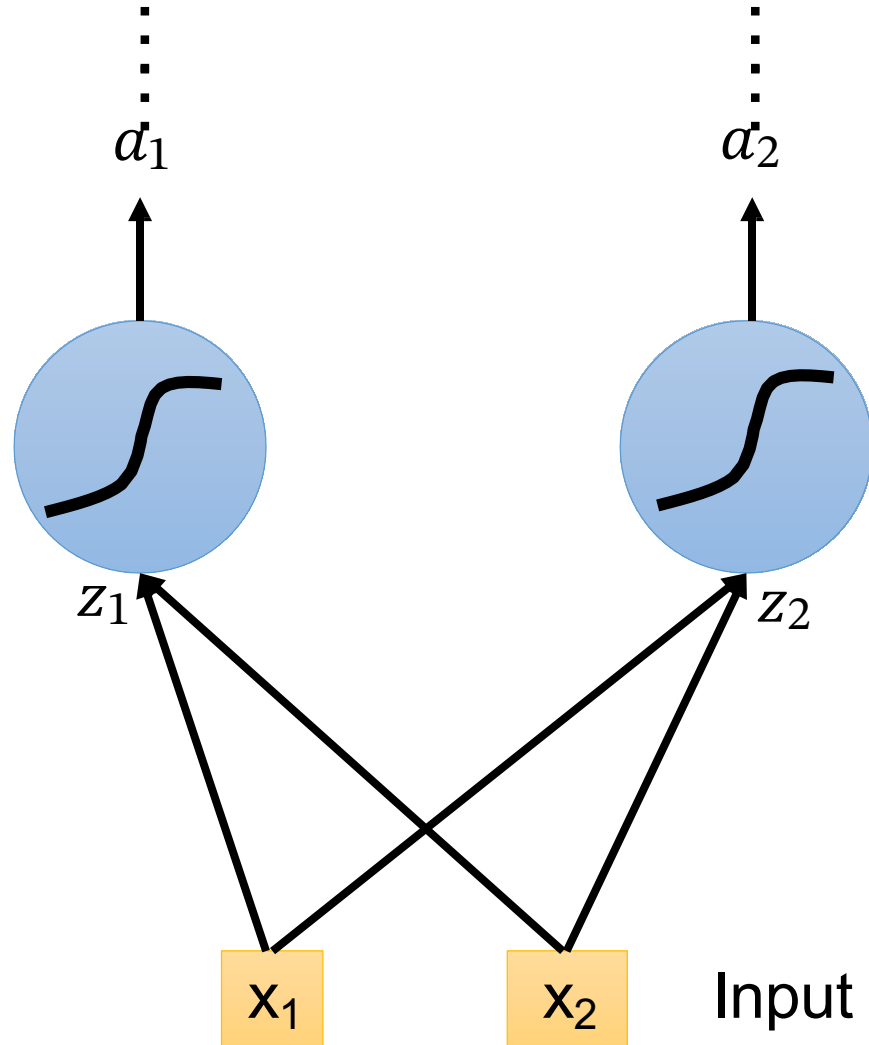
Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

Original Network:

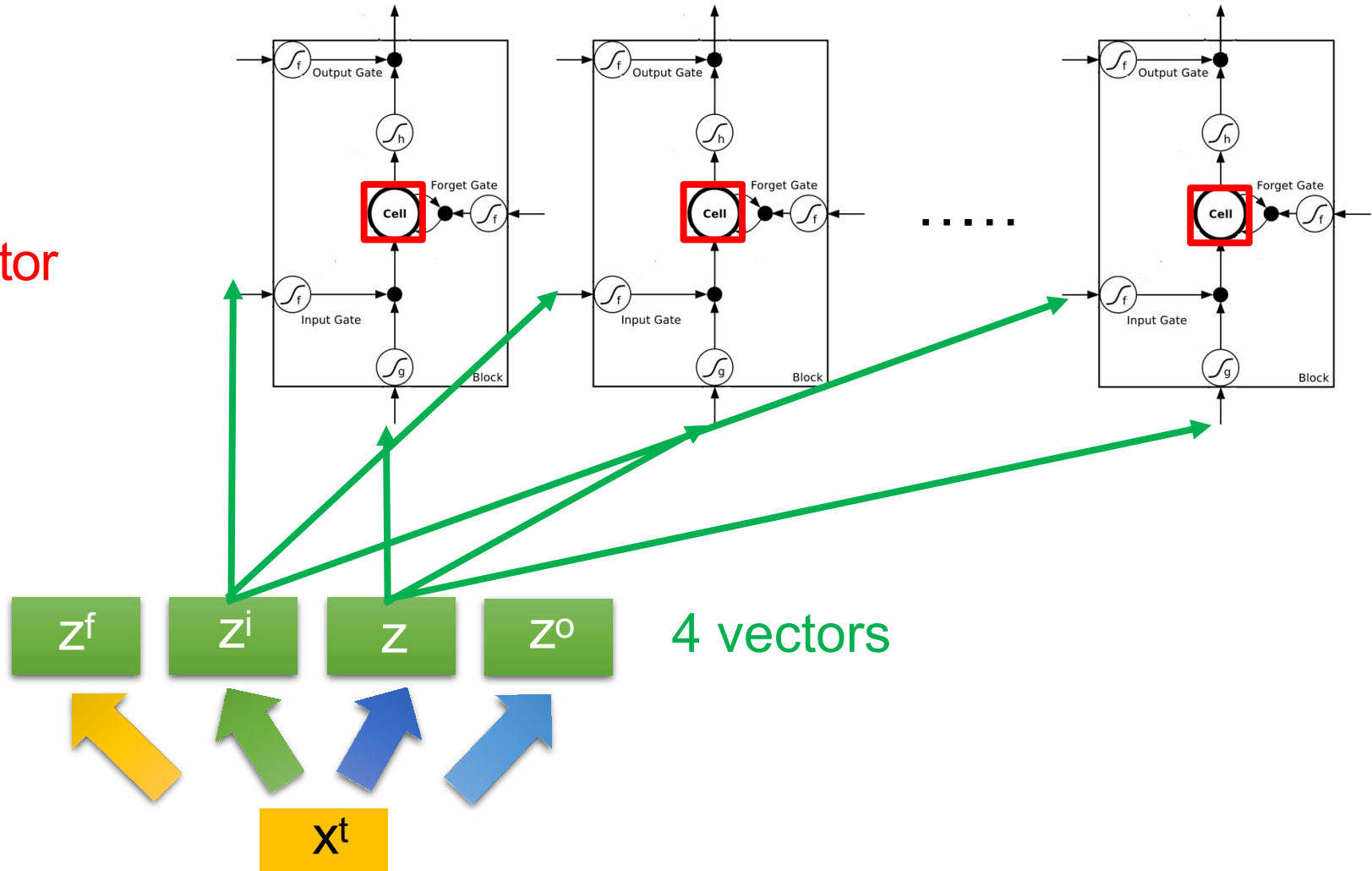
- Simply replace the neurons with LSTM



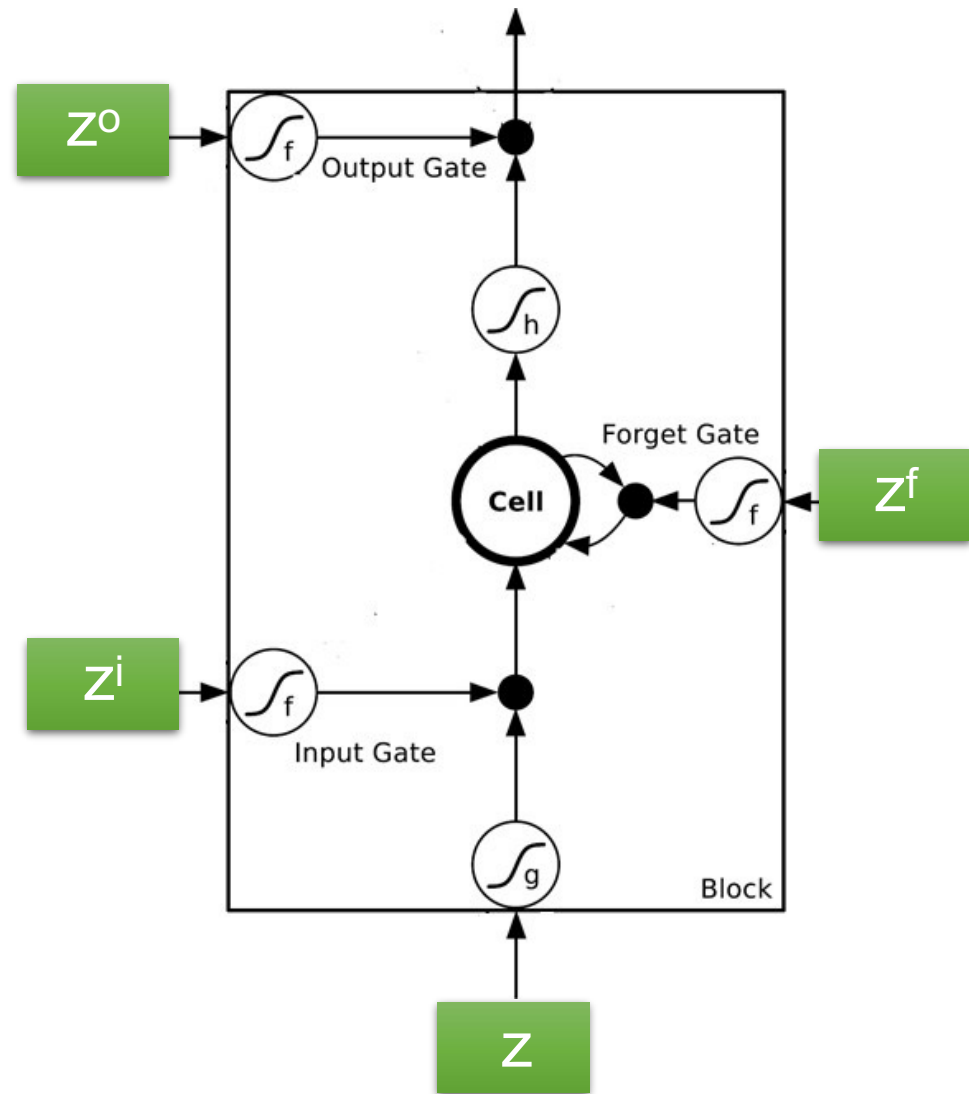
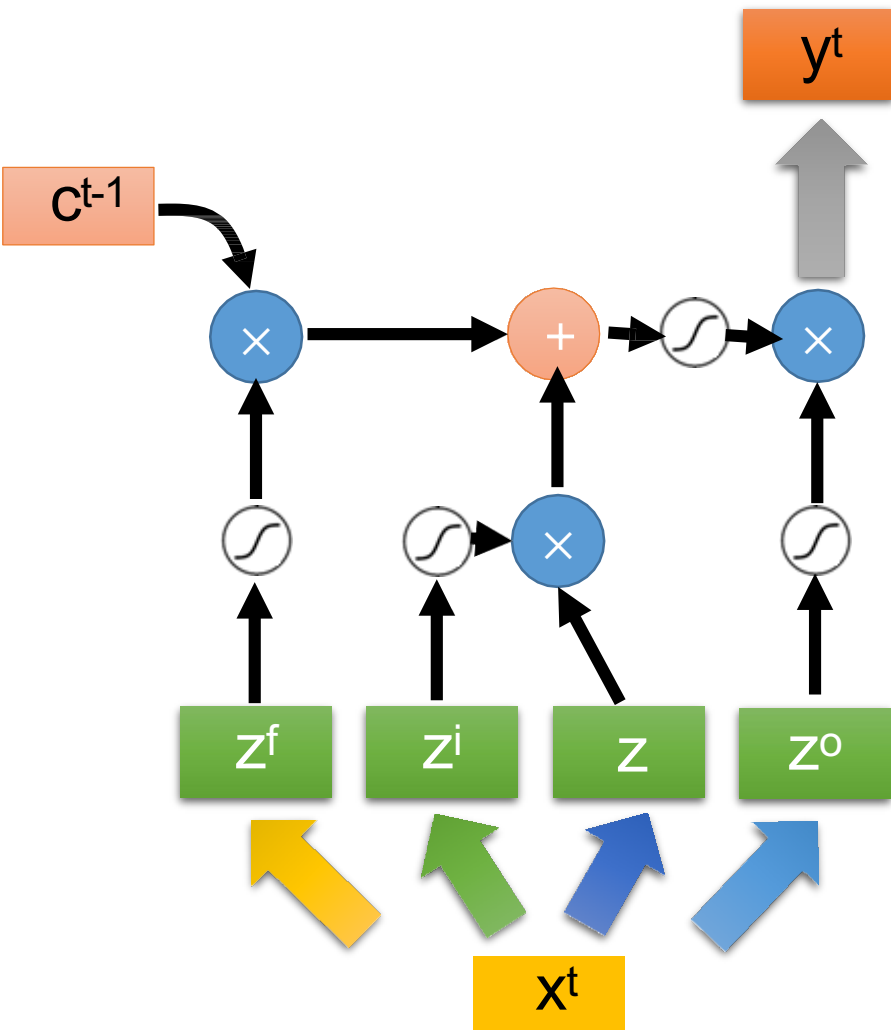
LSTM

 C^{t-1}

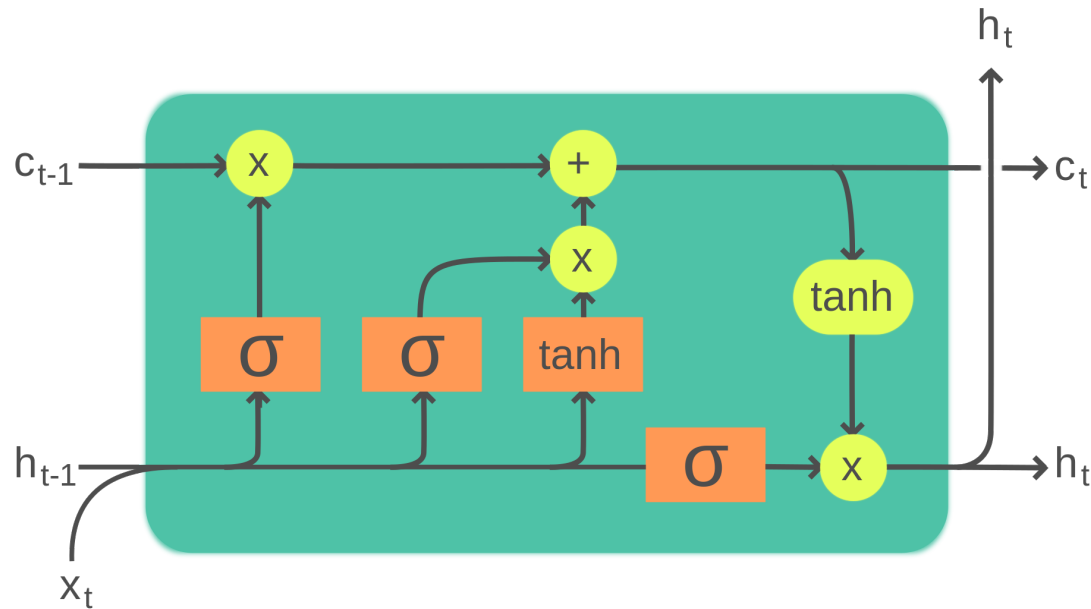
vector



LSTM



LSTM Detailed Structure



Legend:

Layer



Componentwise



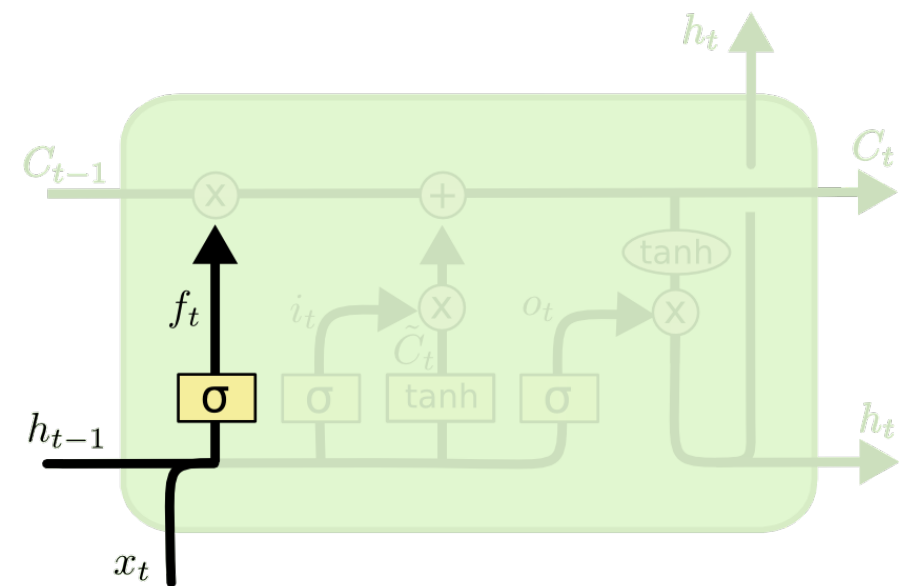
Copy



Concatenate



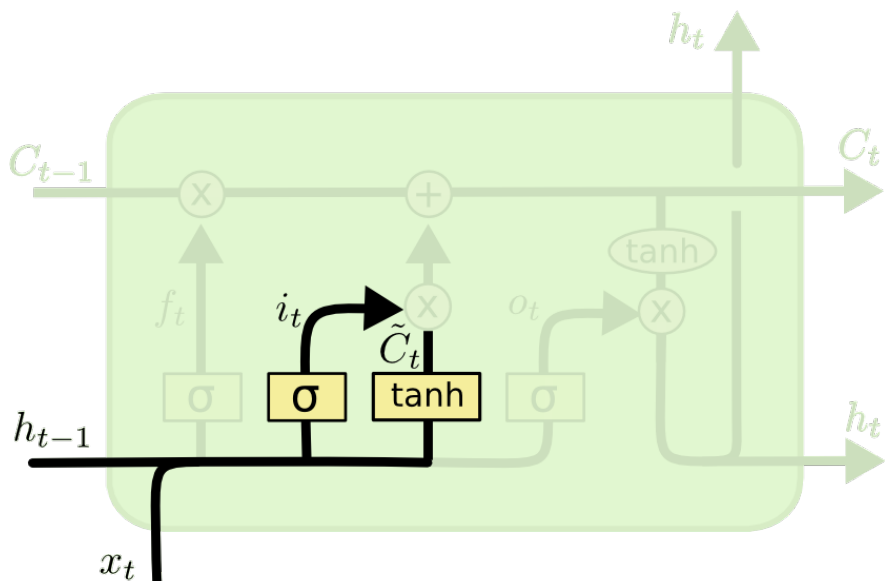
LSTM Detailed Structure



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

The first step in our LSTM is to decide what information we're going to throw away from the cell state. This decision is made by a sigmoid layer called the "forget gate layer." It looks at h_{t-1} and x_t , and outputs a number between 0 and 1 for each number in the cell state C_{t-1} . A 1 represents "completely keep this" while a 0 represents "completely get rid of this."

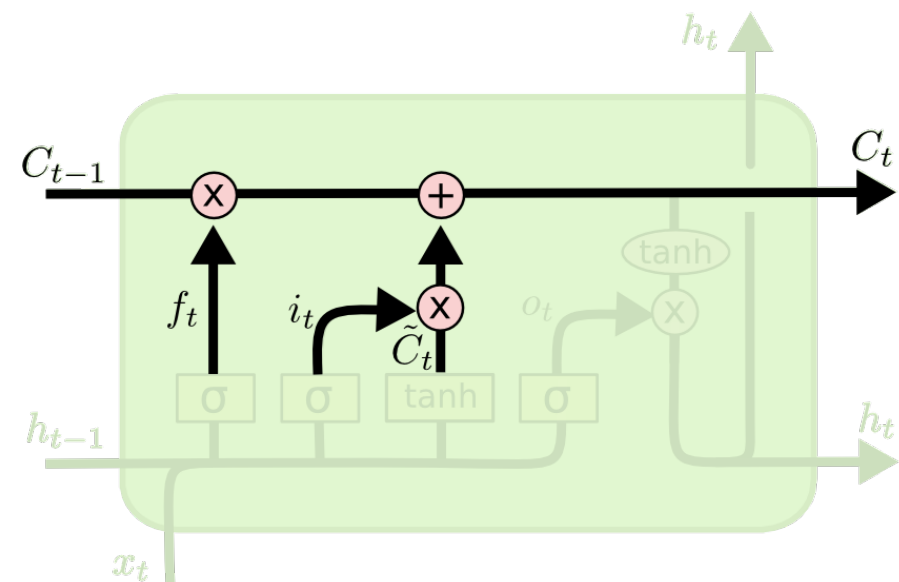
LSTM Detailed Structure



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

The next step is to decide what new information we're going to store in the cell state. This has two parts. First, a sigmoid layer called the "input gate layer" decides which values we'll update. Next, a tanh layer creates a vector of new candidate values, \tilde{C}_t , that could be added to the state. In the next step, we'll combine these two to create an update to the state.

LSTM Detailed Structure

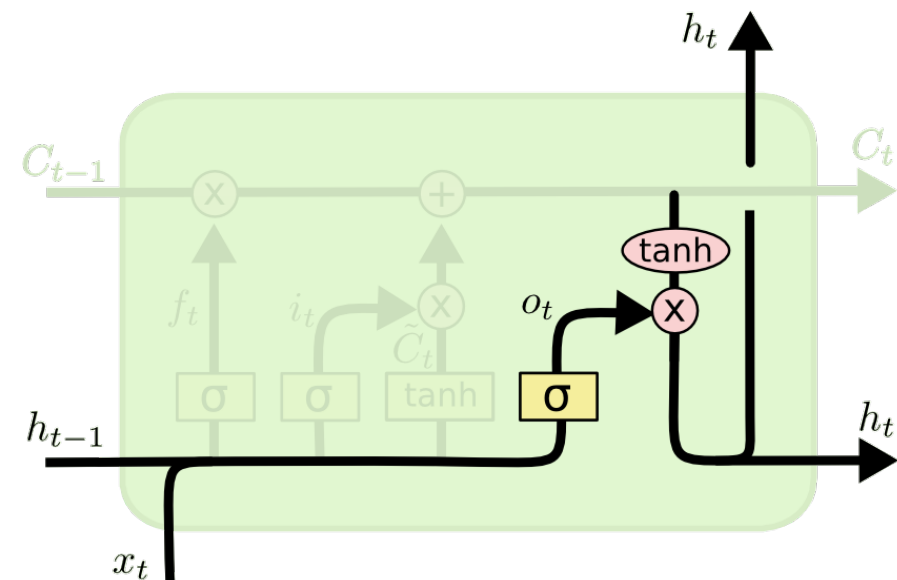


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

It's now time to update the old cell state, C_{t-1} , into the new cell state C_t . The previous steps already decided what to do, we just need to actually do it.

We multiply the old state by f_t , forgetting the things we decided to forget earlier. Then we add $i_t * \tilde{C}_t$. This is the new candidate values, scaled by how much we decided to update each state value.

LSTM Detailed Structure



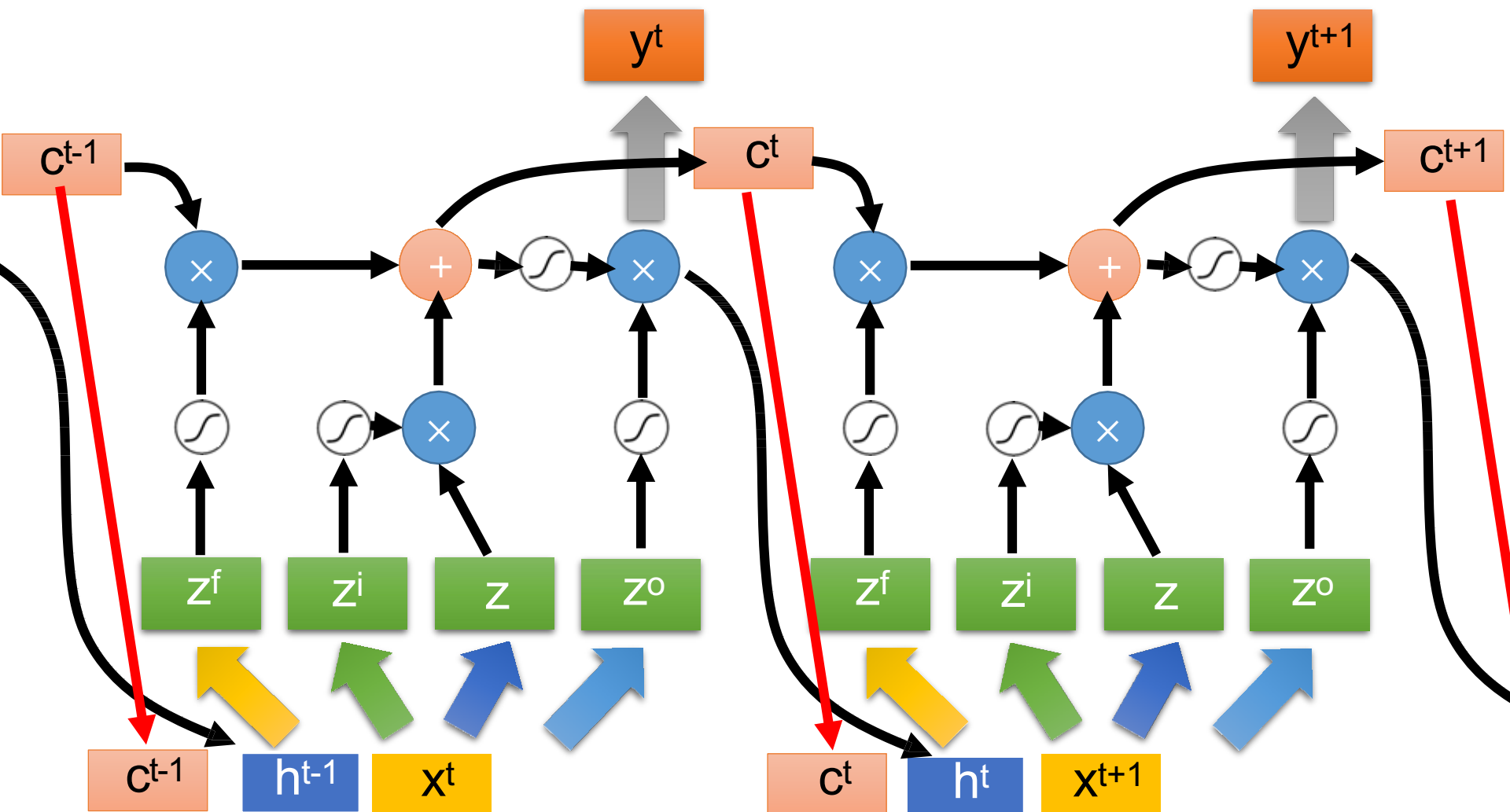
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

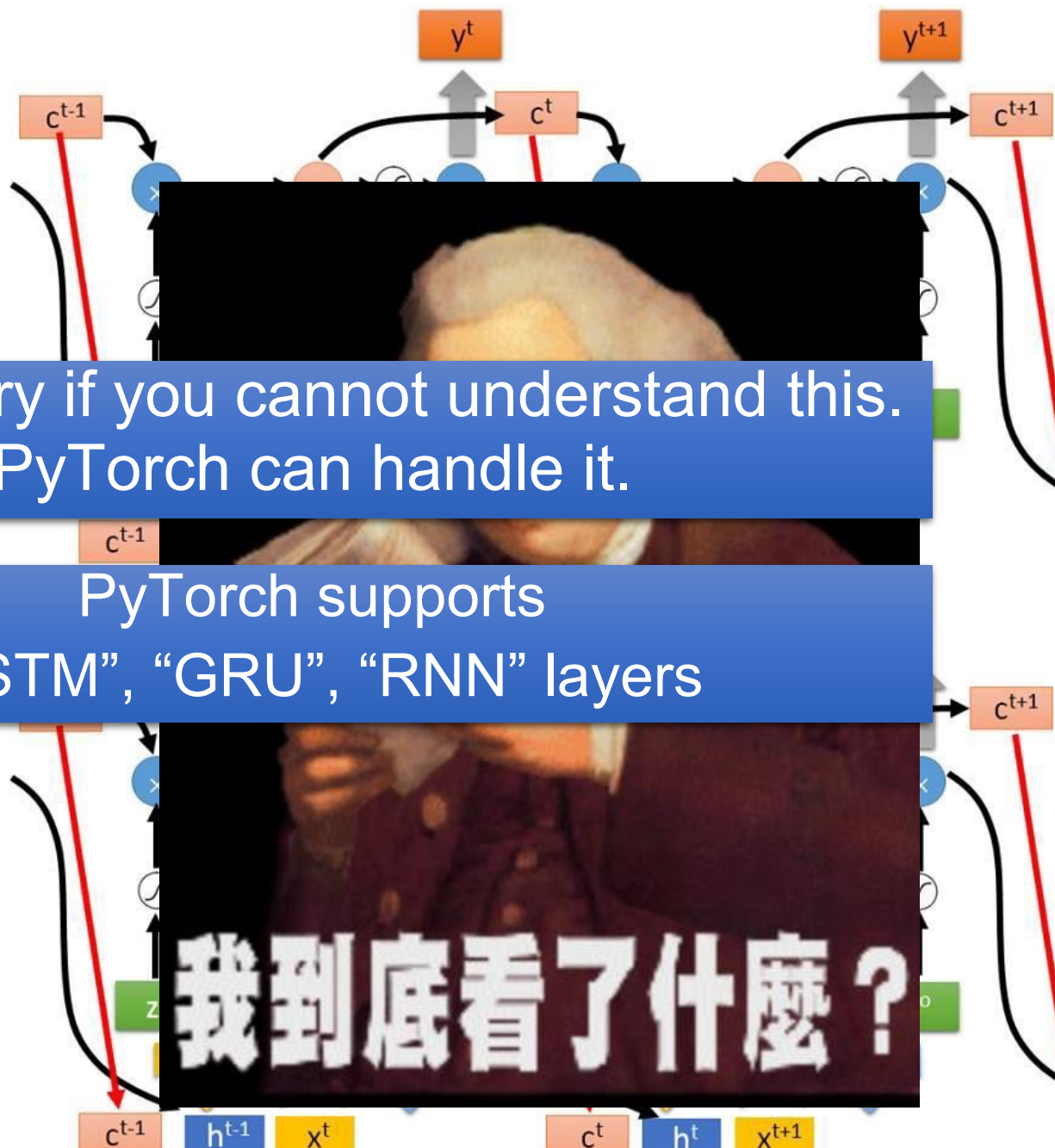
Finally, we need to decide what we're going to output. First, we run a sigmoid layer which decides what parts of the cell state we're going to output. Then, we put the cell state through tanh (to push the values to be between -1 and 1) and multiply it by the output of the sigmoid gate, so that we only output the parts we decided to.

LSTM

Extension: "peephole"



Multiple-layer LSTM



Don't worry if you cannot understand this.
PyTorch can handle it.

PyTorch supports
"LSTM", "GRU", "RNN" layers

This is quite
standard now.

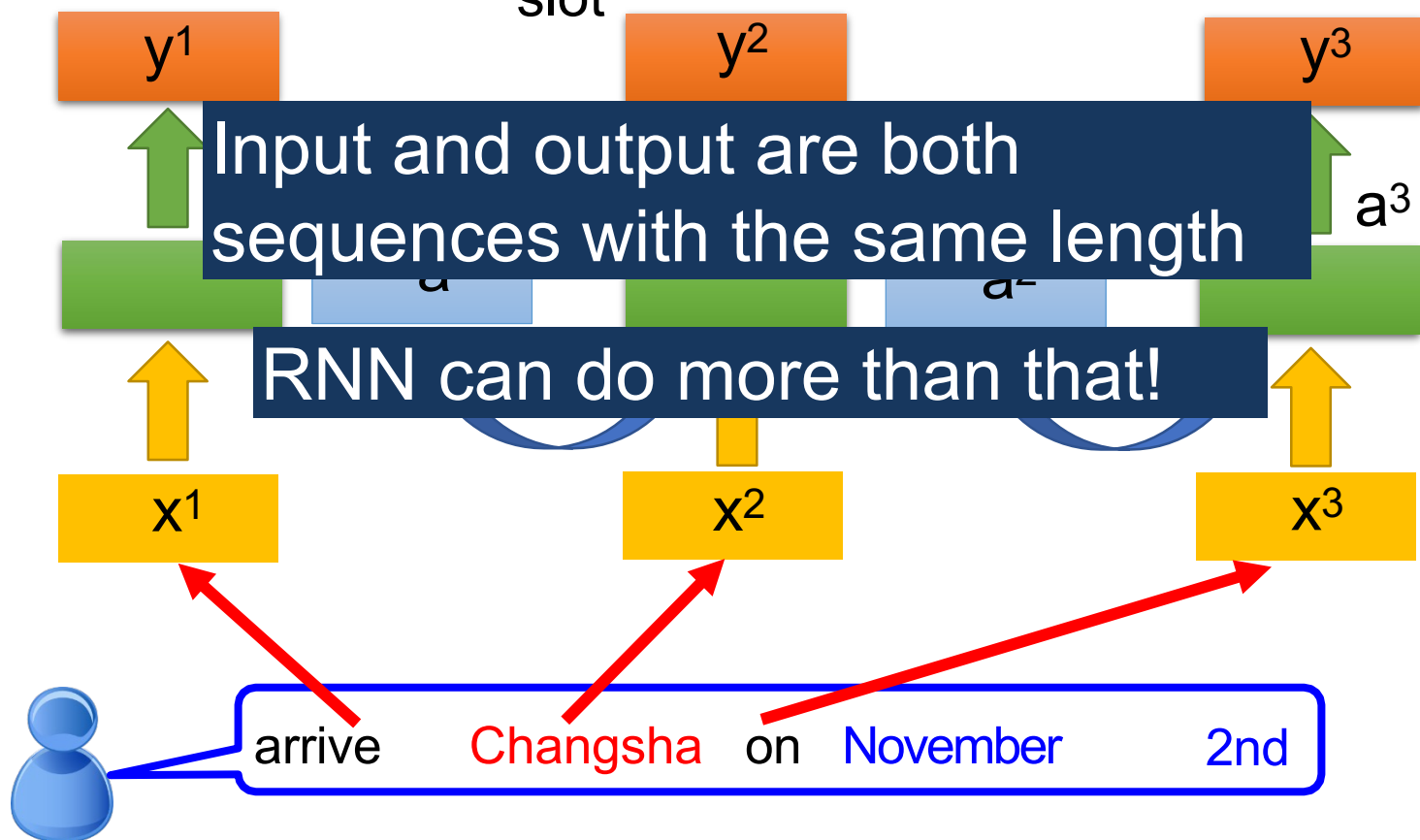
我到底看了什麼？

More Applications

Probability of
“arrive” in each slot

Probability of
“**Changsha**” in each
slot

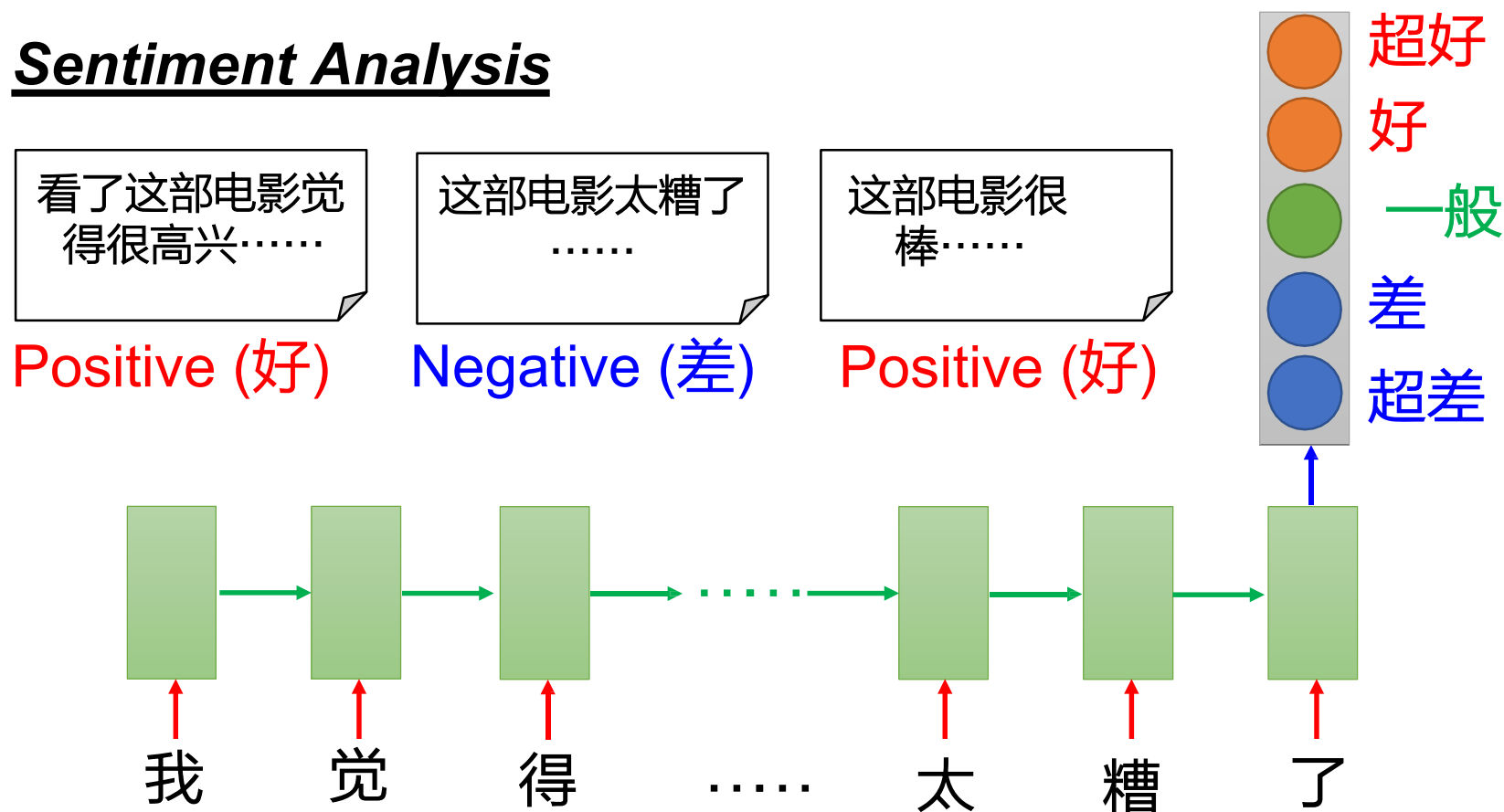
Probability of
“on” in each slot



Many to one

- Input is a vector sequence, but output is only one vector

Sentiment Analysis

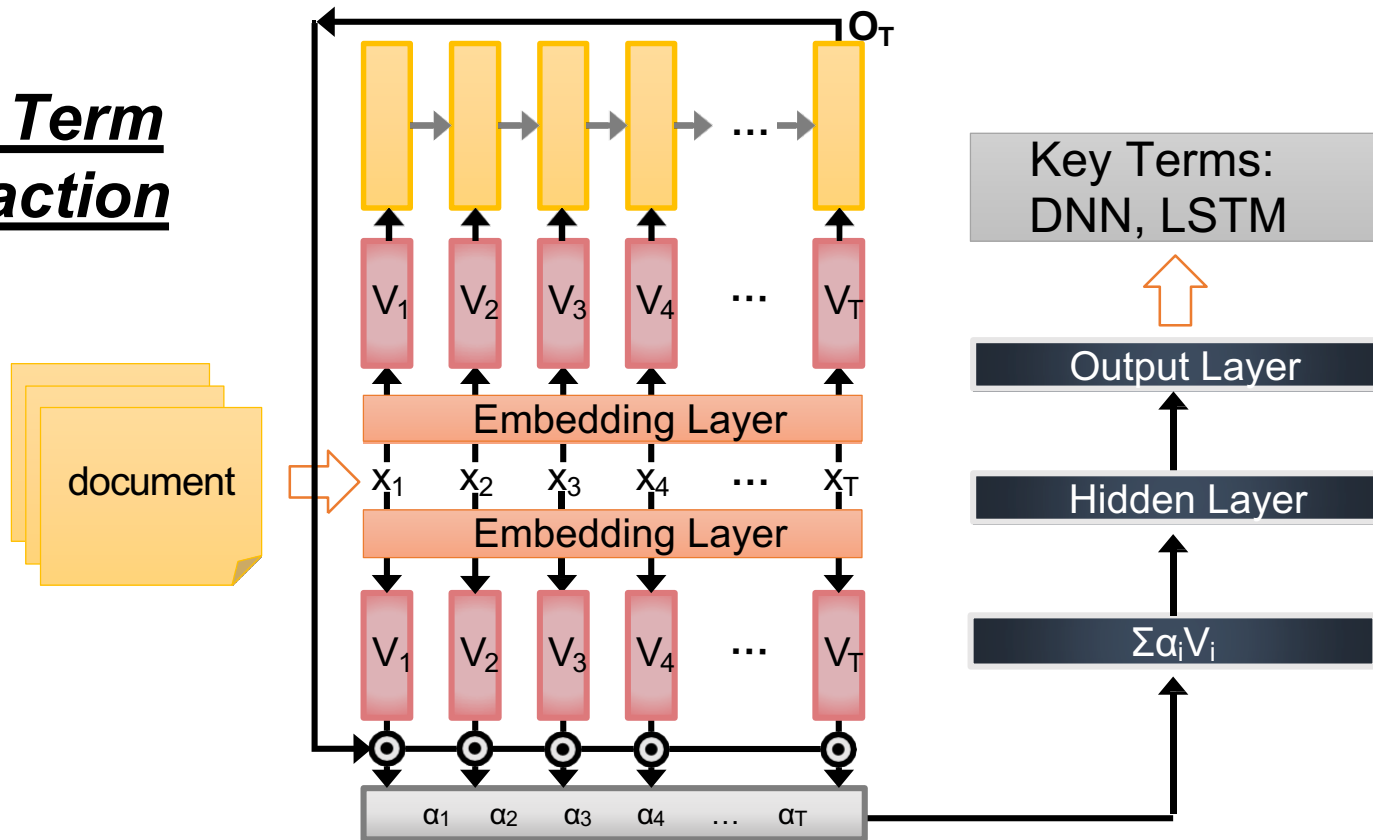


Many to one

[Shen & Lee, Interspeech 16]

- Input is a vector sequence, but output is only one vector

Key Term Extraction



Many to Many (Output is shorter)

- Both input and output are both sequences, **but the output is shorter.**
 - E.g. **Speech Recognition**

Problem?

Why can't it be
“好棒棒”

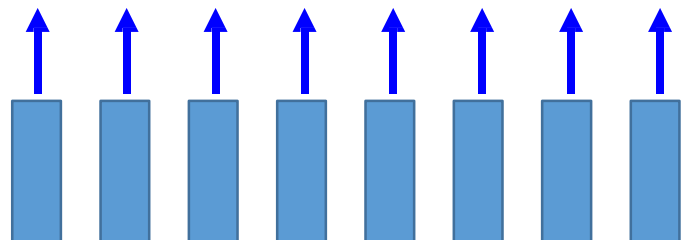
Output: “好棒” (character sequence)



Trimming

好 好 好 棒 棒 棒 棒 棒

Input:

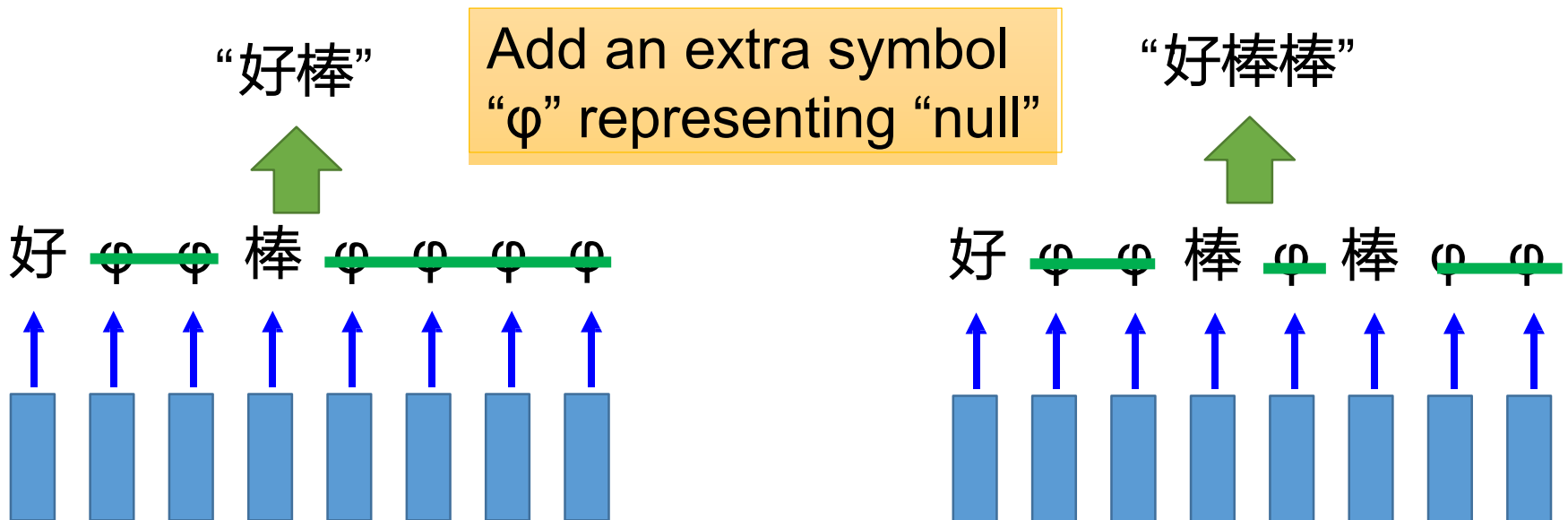


(vector sequence)



Many to Many (Output is shorter)

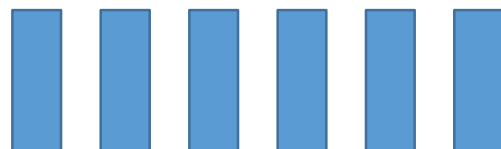
- Both input and output are both sequences, **but the output is shorter.**
- Connectionist Temporal Classification (CTC) [Alex Graves, ICML'06][Alex Graves, ICML'14][Hasim Sak, Interspeech'15][Jie Li, Interspeech'15][Andrew Senior, ASRU'15]



Many to Many (Output is shorter)







- CTC: Training







Acoustic
Features:









Label: 好 棒

All possible alignments are
considered as correct.

					
好	ϕ	棒	ϕ	ϕ	ϕ

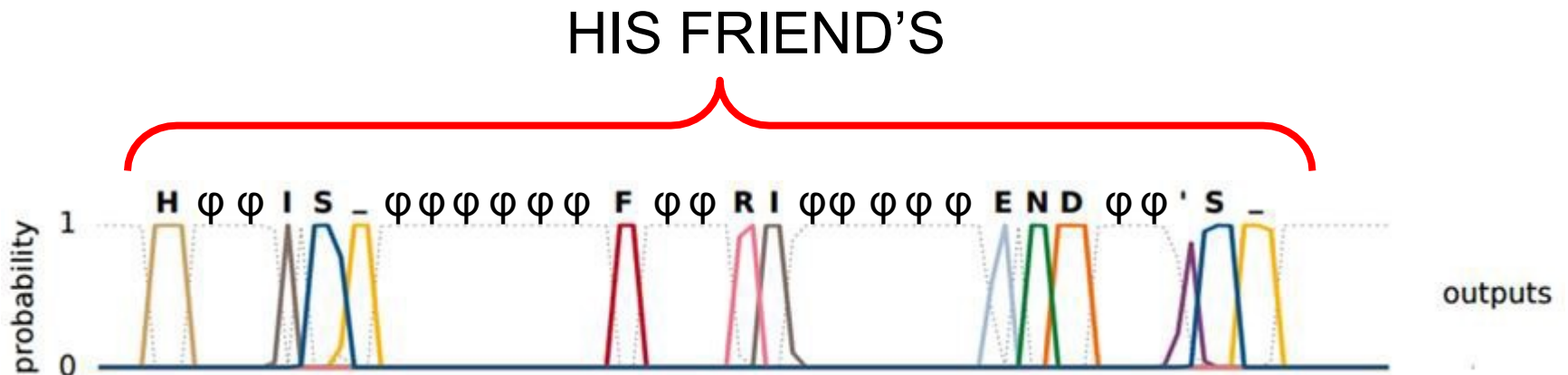
					
好	ϕ	ϕ	棒	ϕ	ϕ

					
好	ϕ	ϕ	ϕ	棒	ϕ

⋮

Many to Many (Output is shorter)

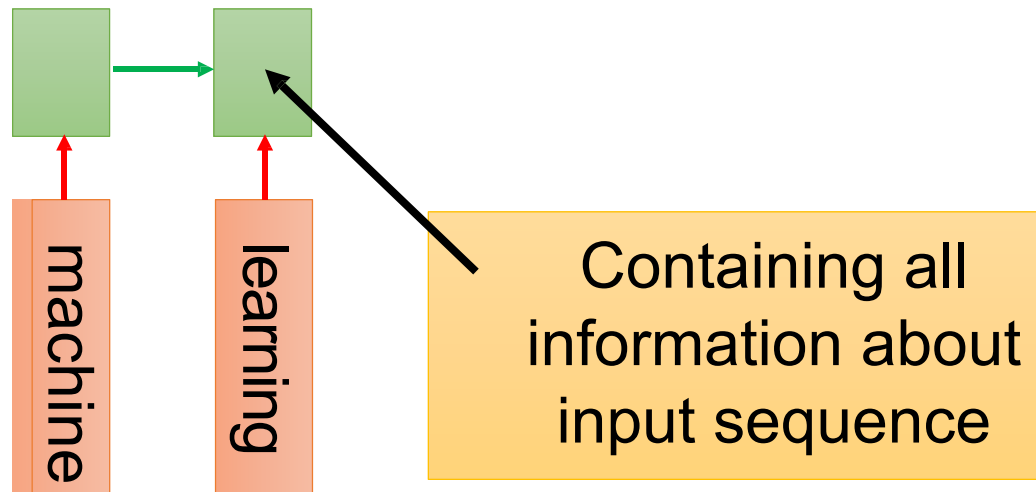
- CTC: example



Graves, Alex, and Navdeep Jaitly. "Towards end-to-end speech recognition with recurrent neural networks." *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*. 2014.

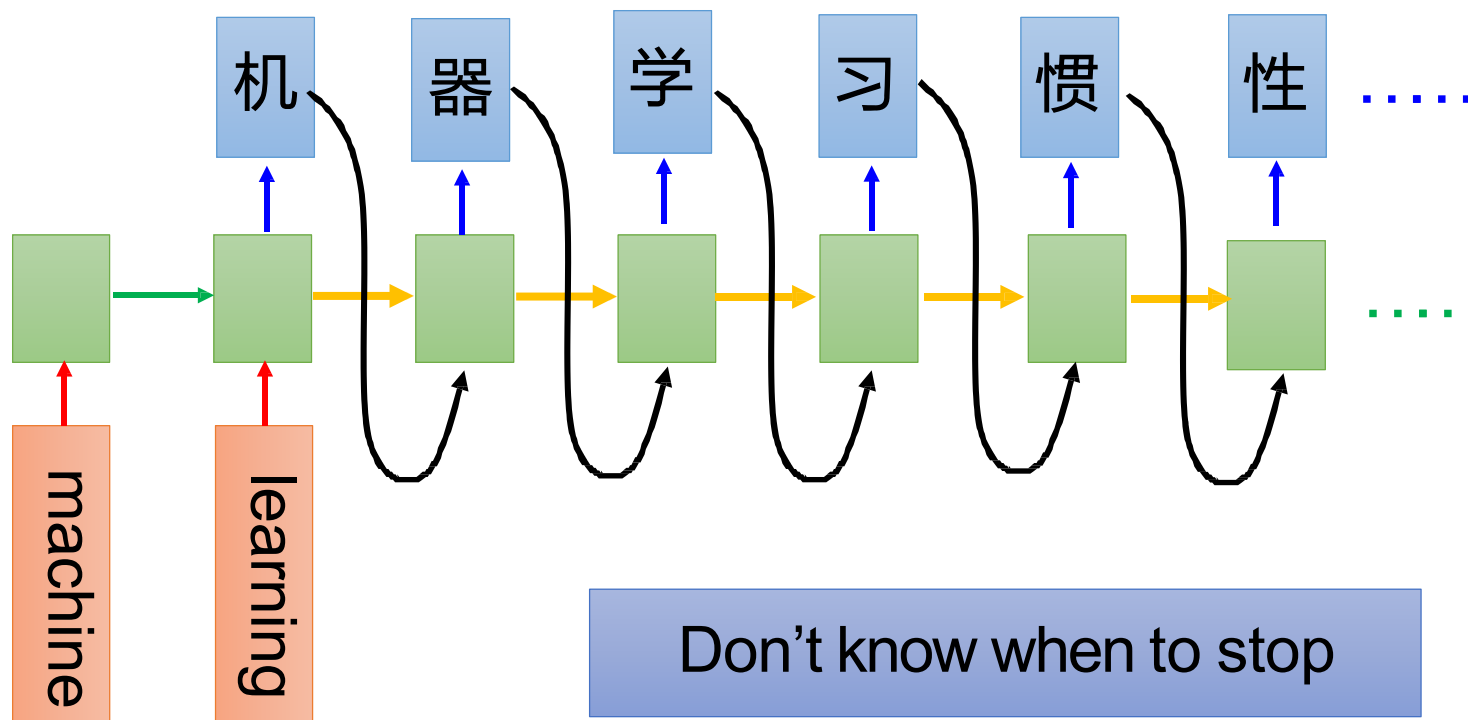
Many to Many (No Limitation)

- Both input and output are both sequences *with different lengths*. → *Sequence to sequence learning*
 - E.g. *Machine Translation* (machine learning → 机器学习)



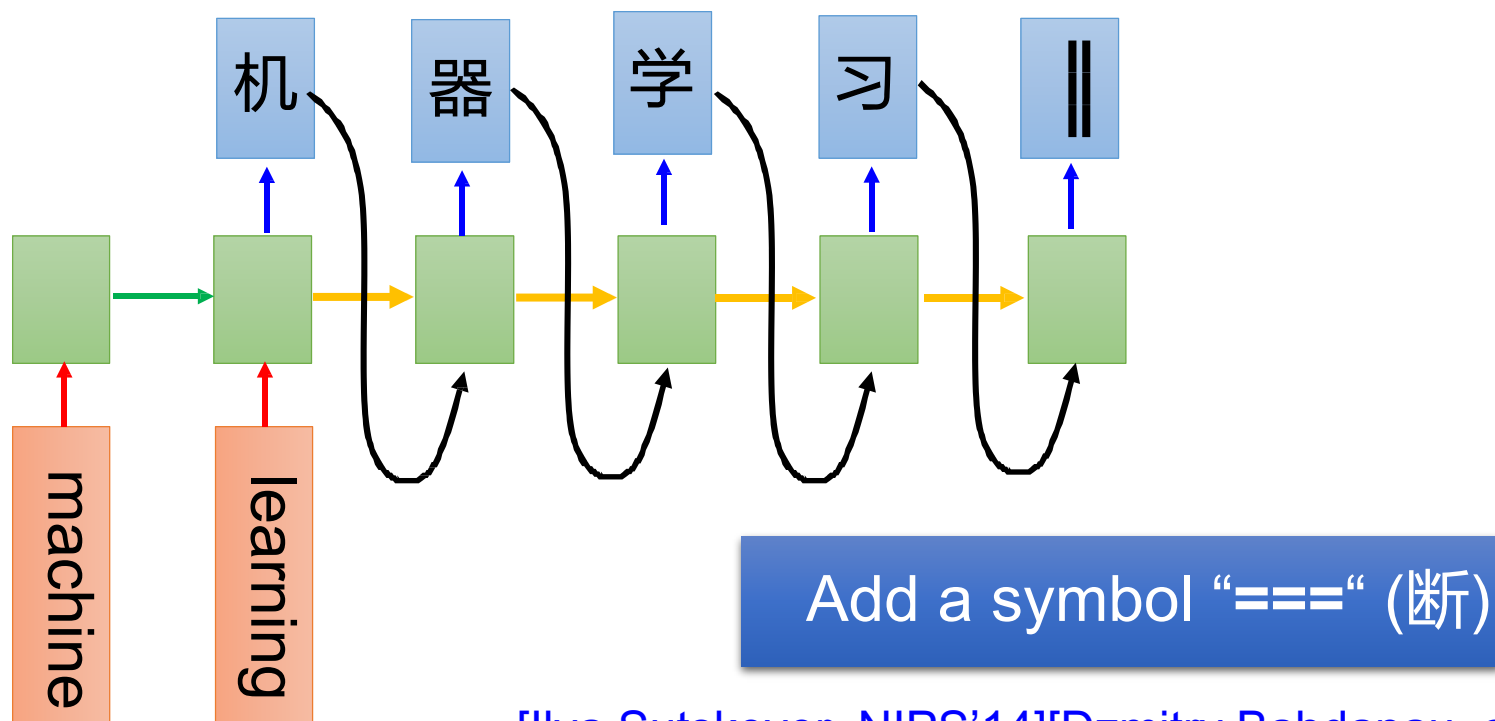
Many to Many (No Limitation)

- Both input and output are both sequences *with different lengths*. → *Sequence to sequence learning*
 - E.g. *Machine Translation* (machine learning → 机器学习)



Many to Many (No Limitation)

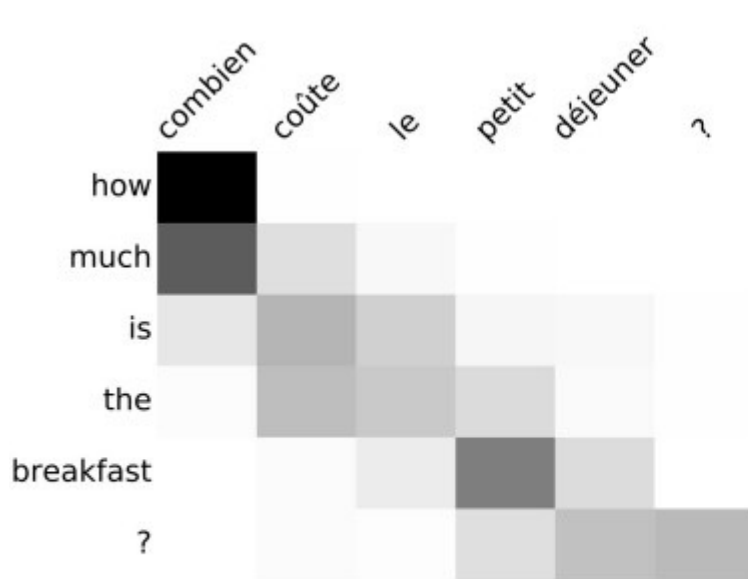
- Both input and output are both sequences *with different lengths*. → *Sequence to sequence learning*
 - E.g. *Machine Translation* (machine learning → 机器学习)



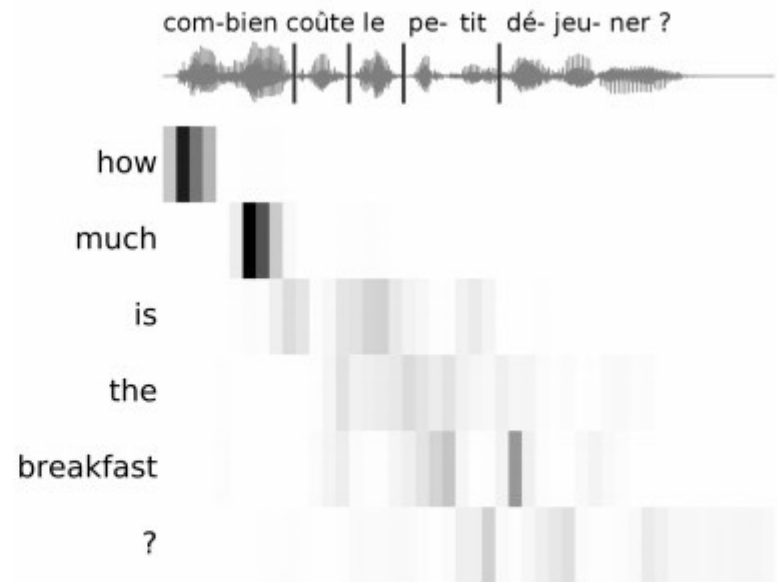
[Ilya Sutskever, NIPS'14][Dzmitry Bahdanau, arXiv'15]

Many to Many (No Limitation)

- Both input and output are both sequences **with different lengths**. → **Sequence to sequence learning**
 - E.g. **Machine Translation** (machine learning → 机器学习)



(a) Machine translation alignment

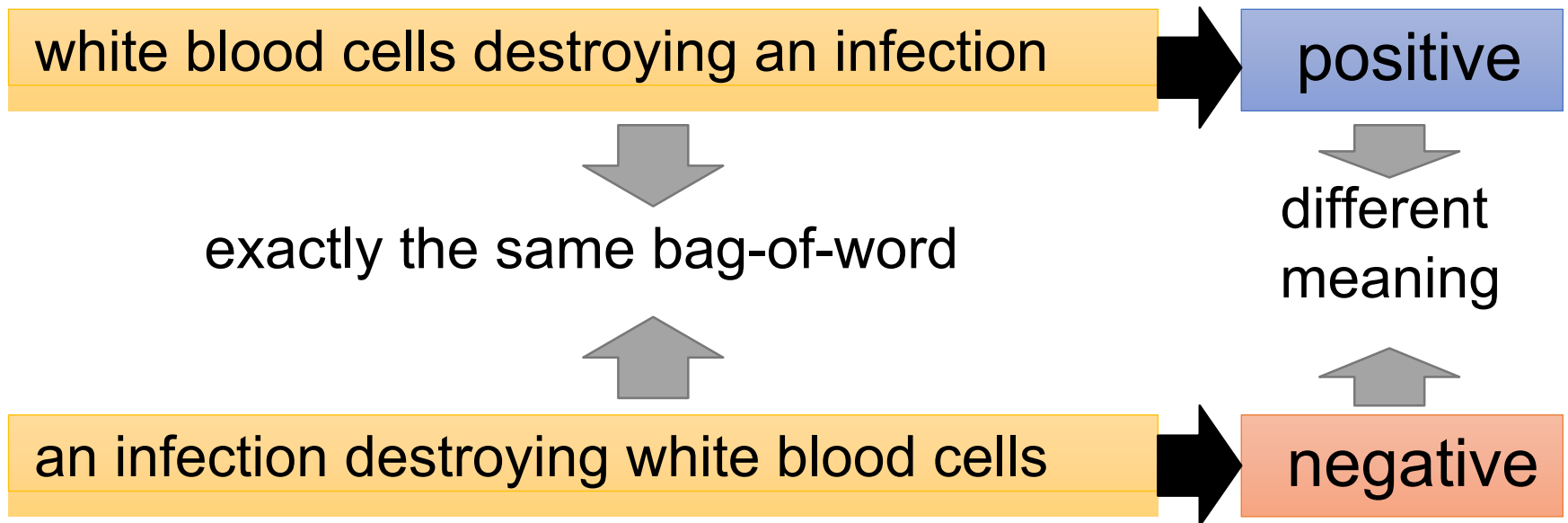


(b) Speech translation alignment

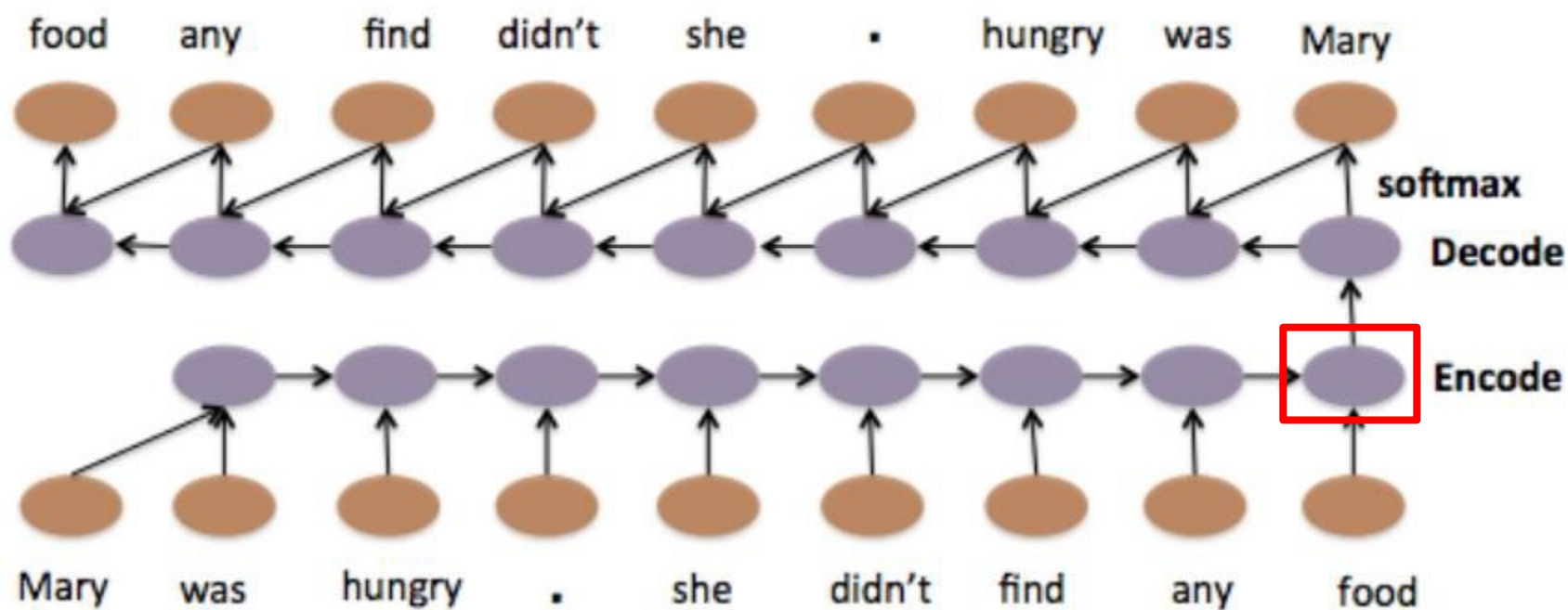
Figure 1: Alignments performed by the attention model during training

Sequence-to-sequence Auto-encoder -Text

- To understand the meaning of a word sequence, the order of the words can not be ignored.

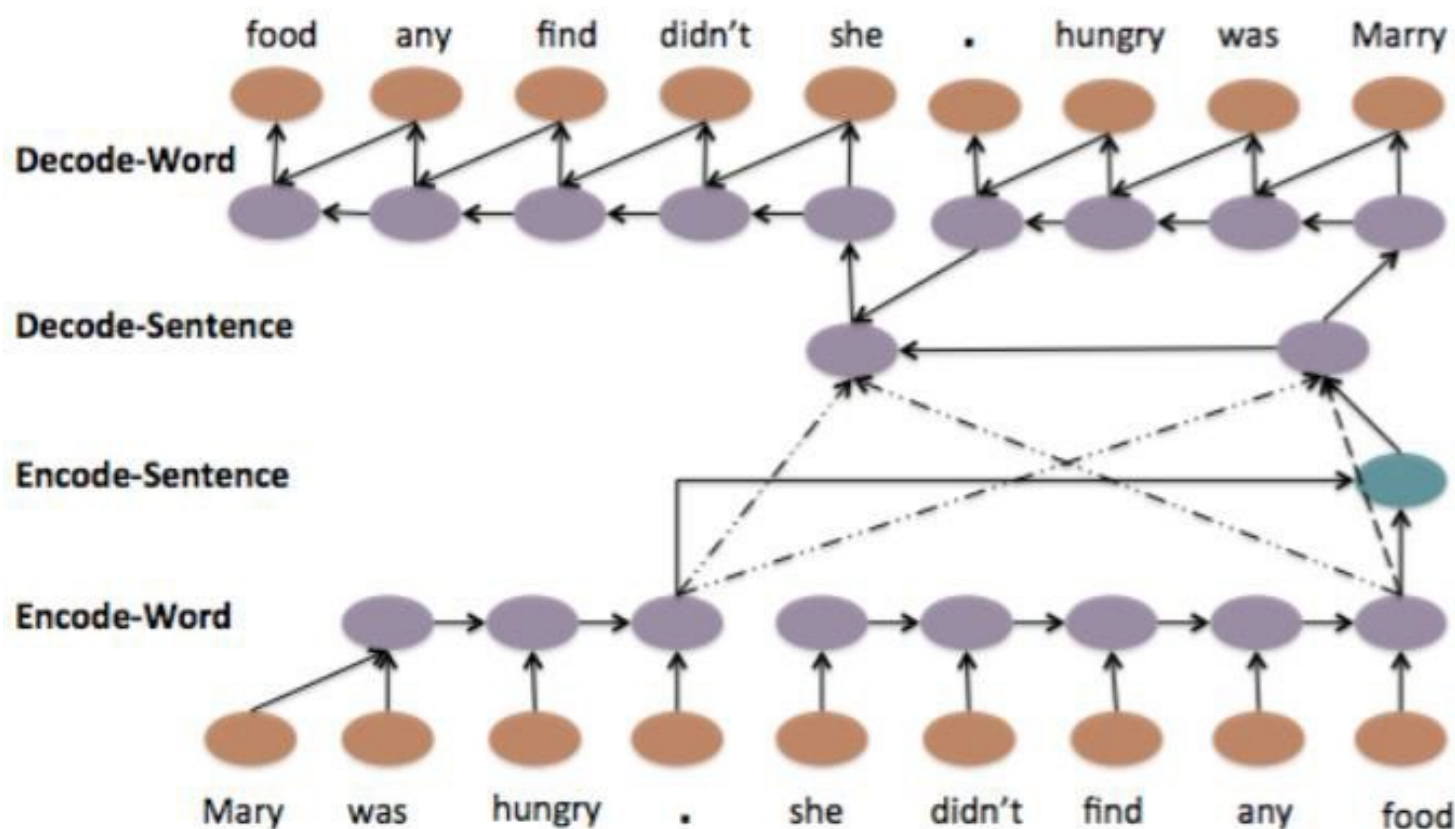


Sequence-to-sequence Auto-encoder -Text



Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." *arXiv preprint arXiv:1506.01057*(2015).

Sequence-to-sequence Auto-encoder -Text

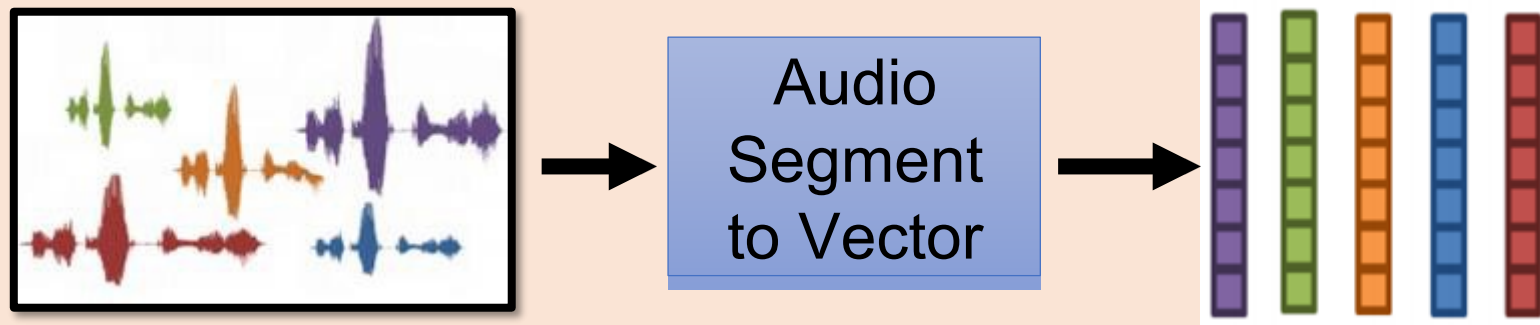


Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." *arXiv preprint arXiv:1506.01057*(2015).

Sequence-to-sequence Auto-encoder -Speech

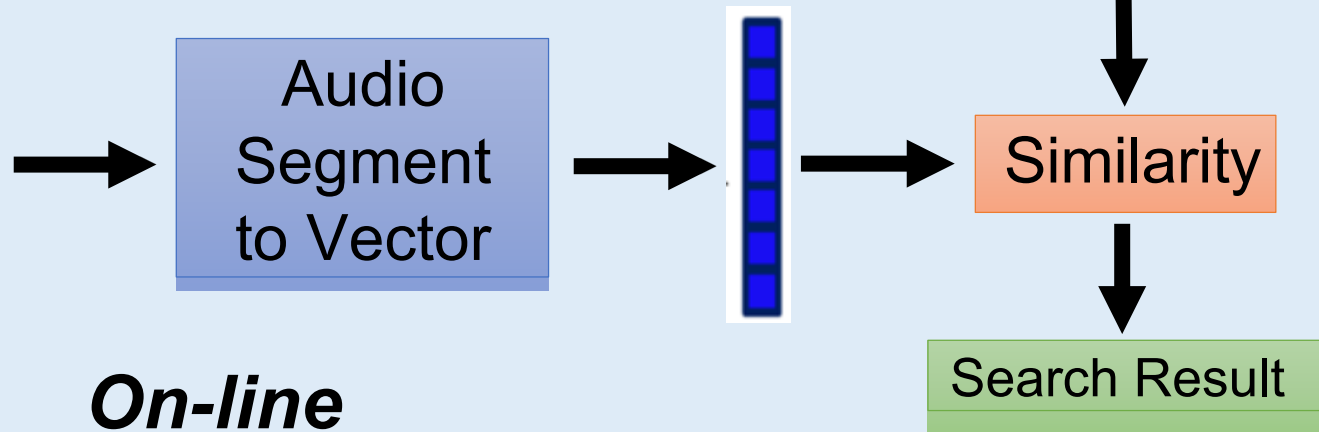
Audio archive divided into variable-length audio segments

Off-line

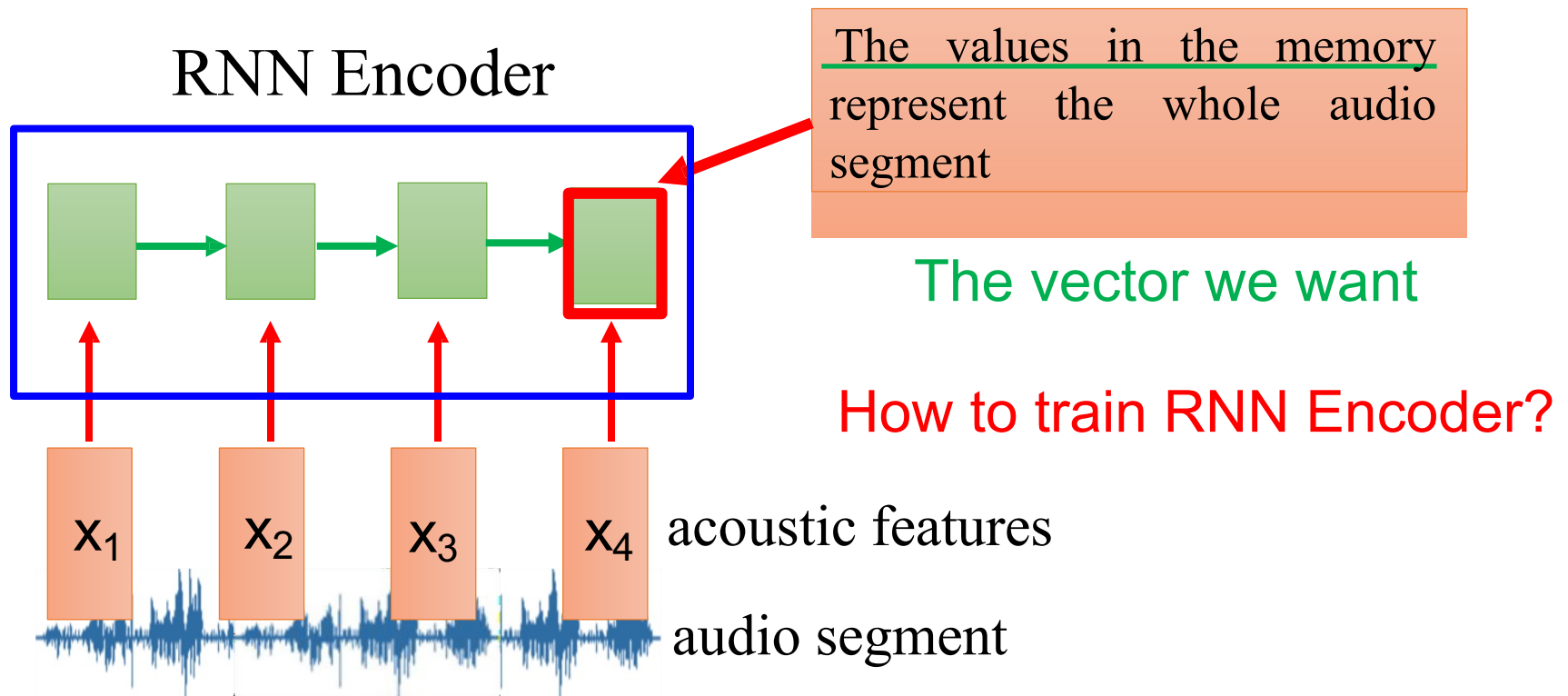


Spoken Query

On-line

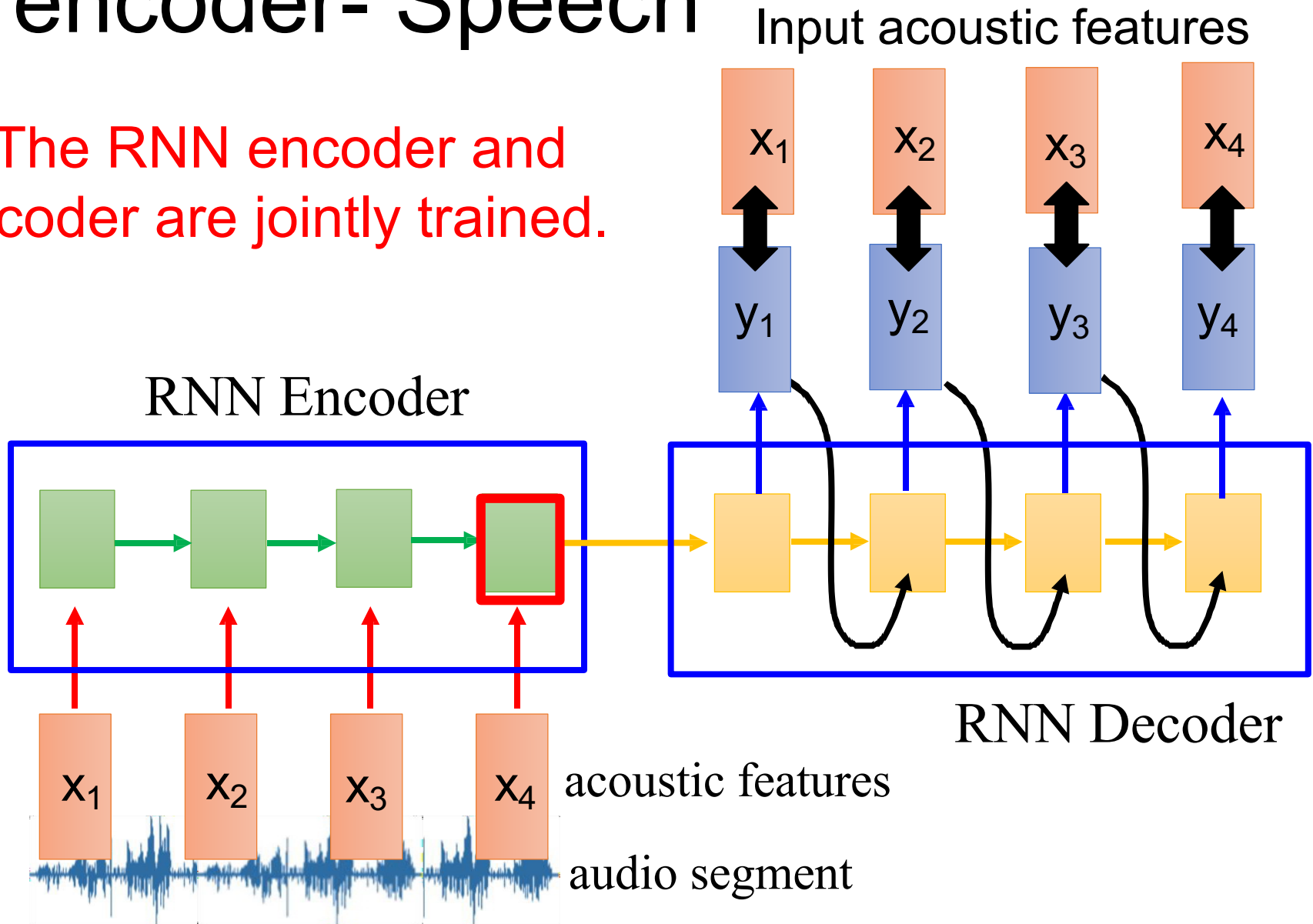


Sequence-to-sequence Auto-encoder -Speech



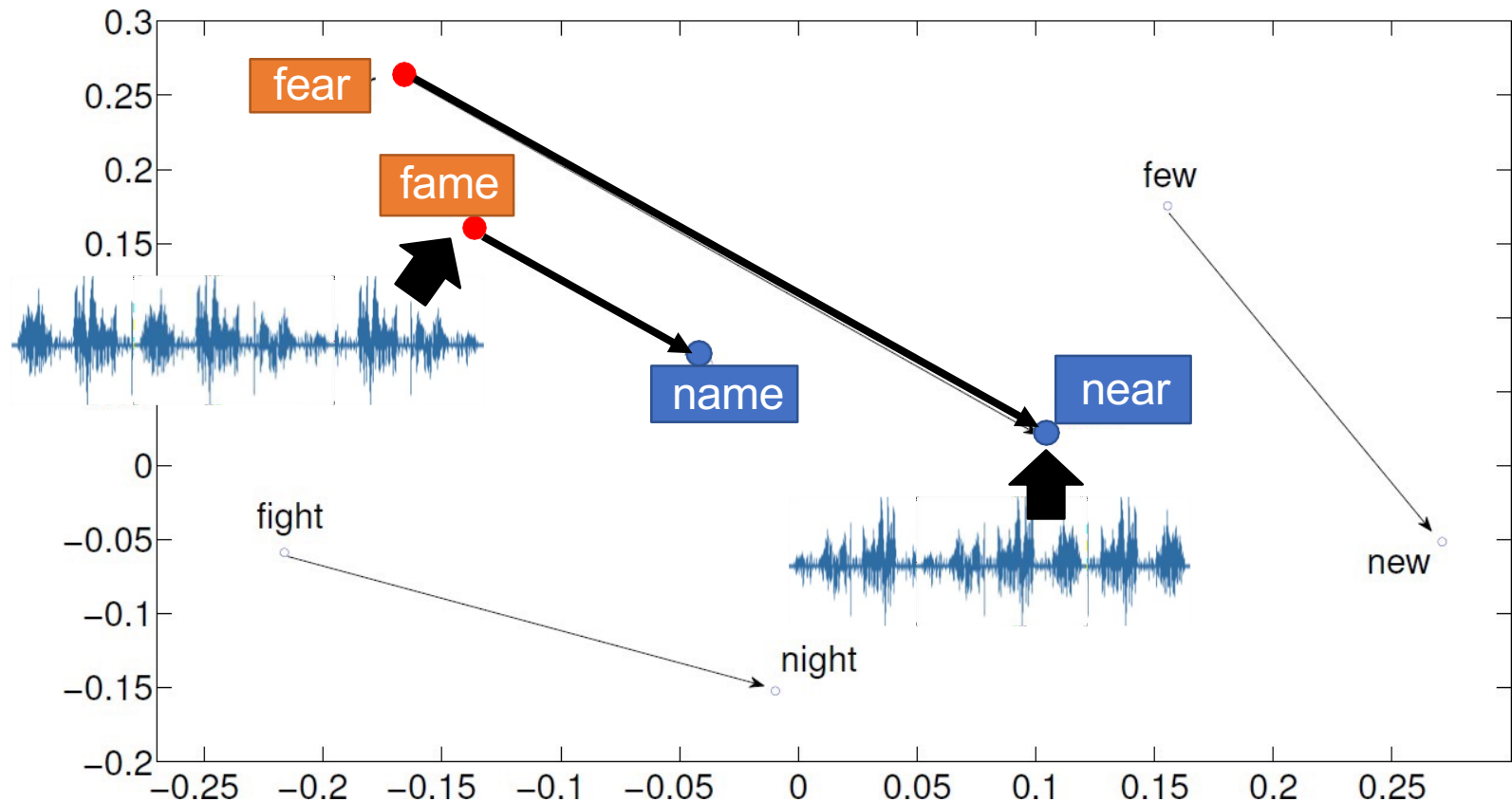
Sequence-to-sequence Auto-encoder- Speech

The RNN encoder and decoder are jointly trained.

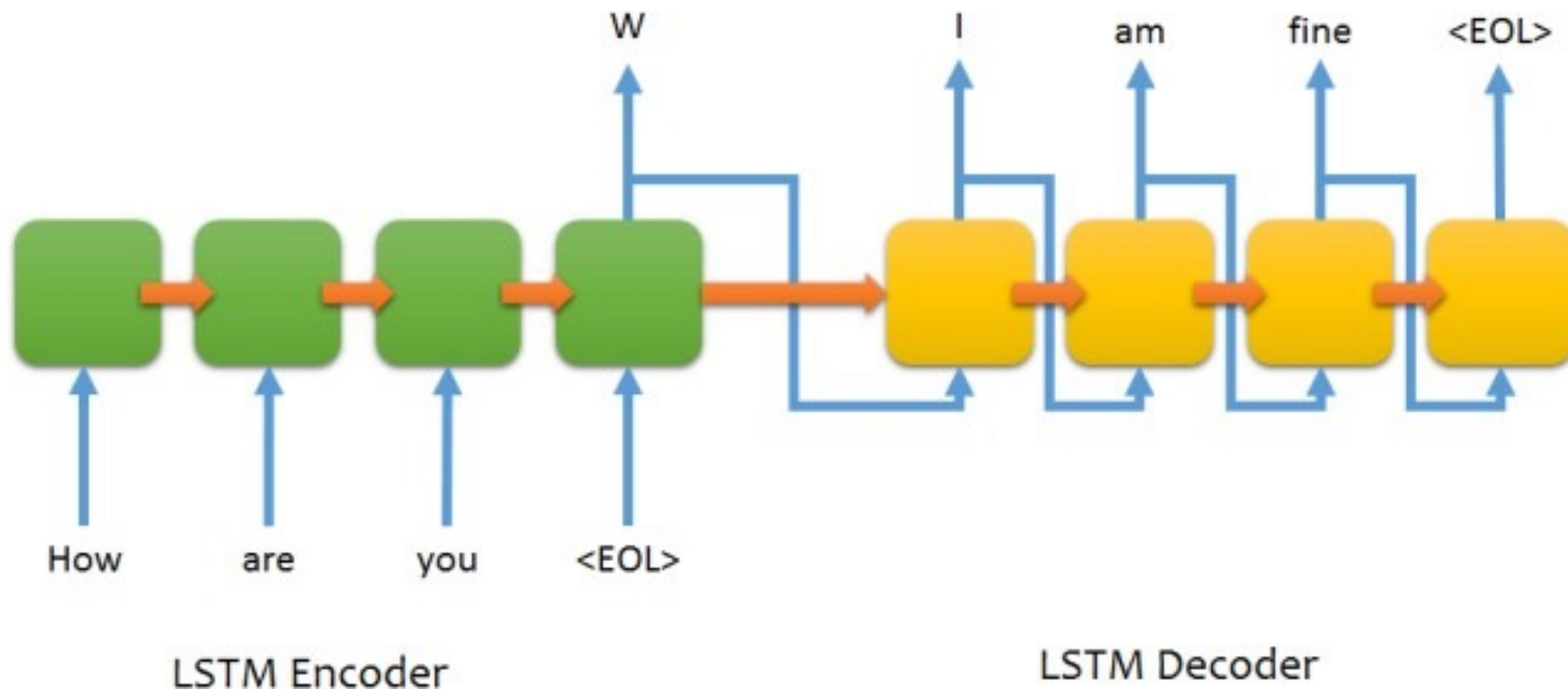


Sequence-to-sequence Auto-encoder -Speech

- Visualizing embedding vectors of the words

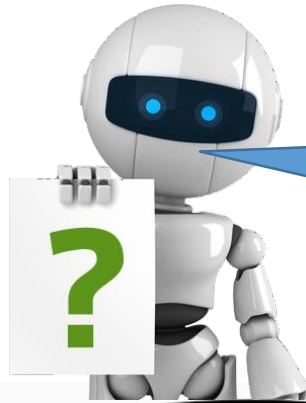
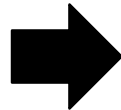


Example: Chat-bot



电视影集 (~40,000 sentences)、美国总统大选辩论

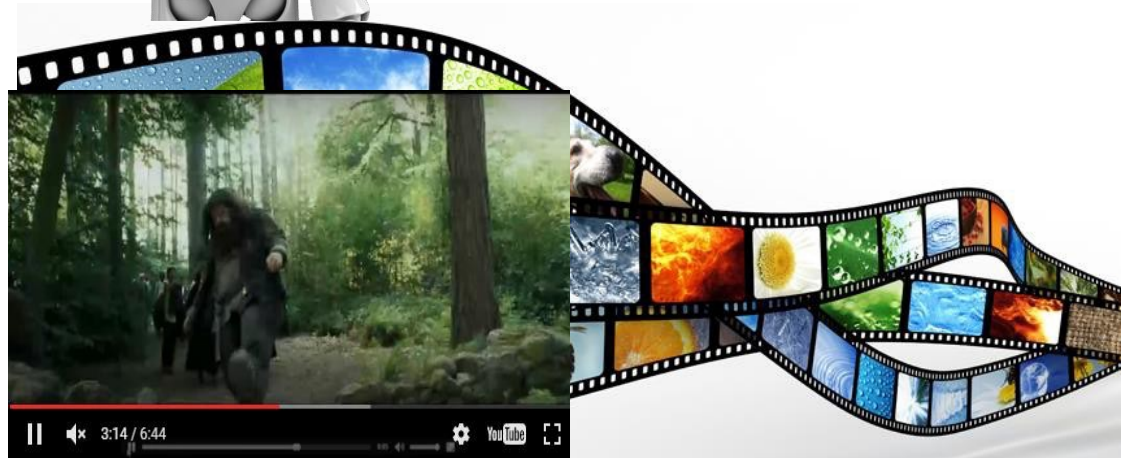
Example: Video Caption Generation



A girl is running.



A group of people is knocked by a tree.

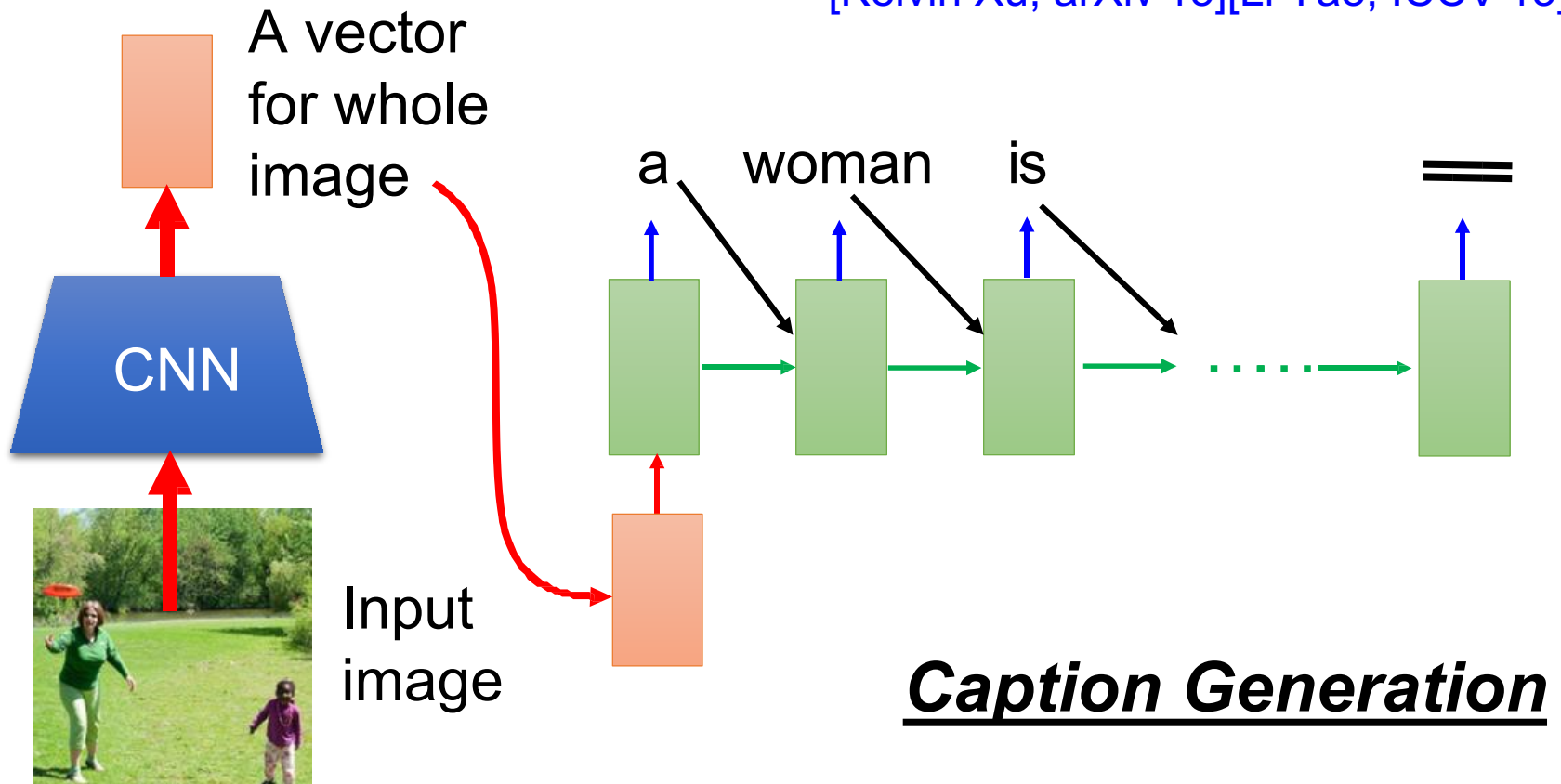


A group of people is walking in the forest.

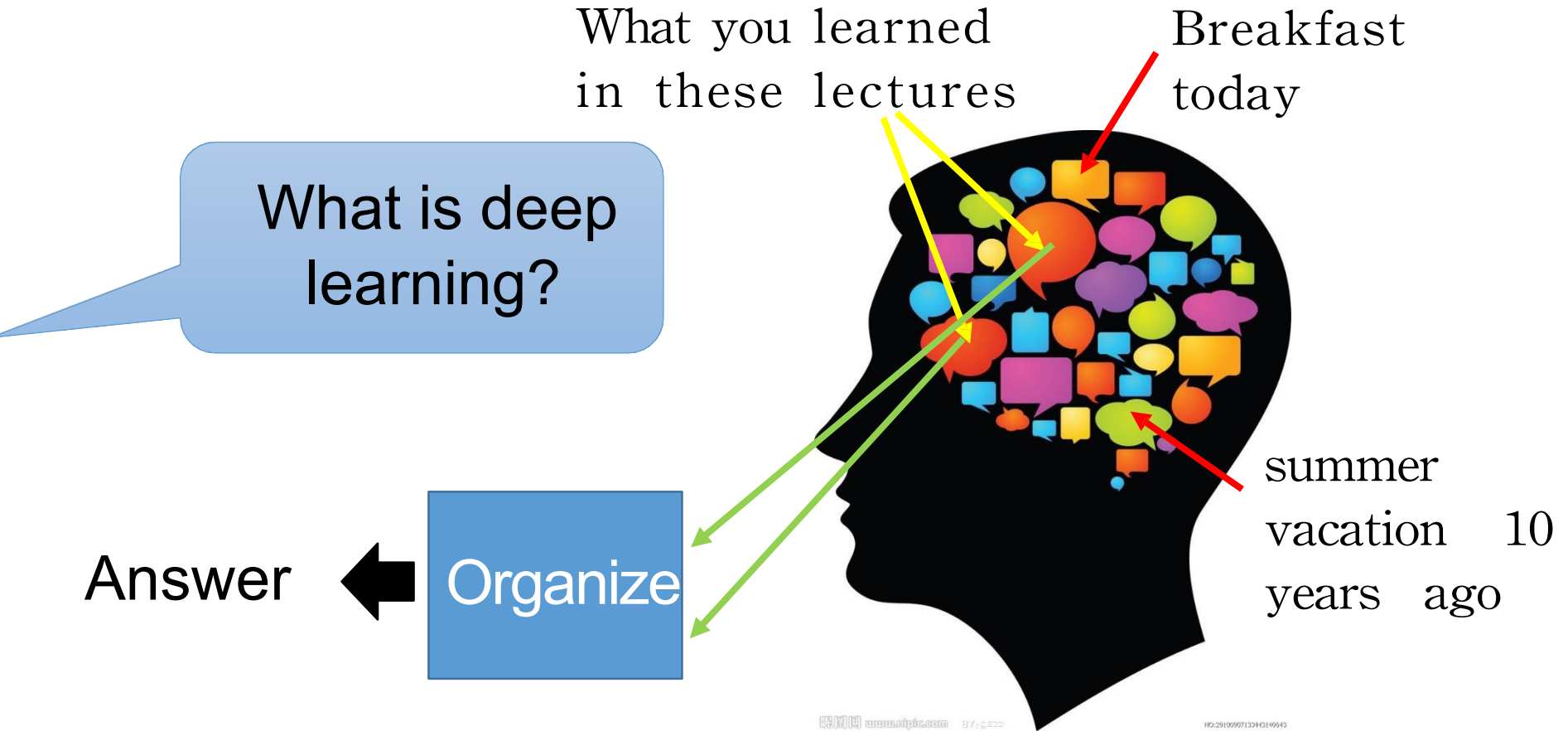
Example: Image Caption Generation

- Input an image, but output a sequence of words

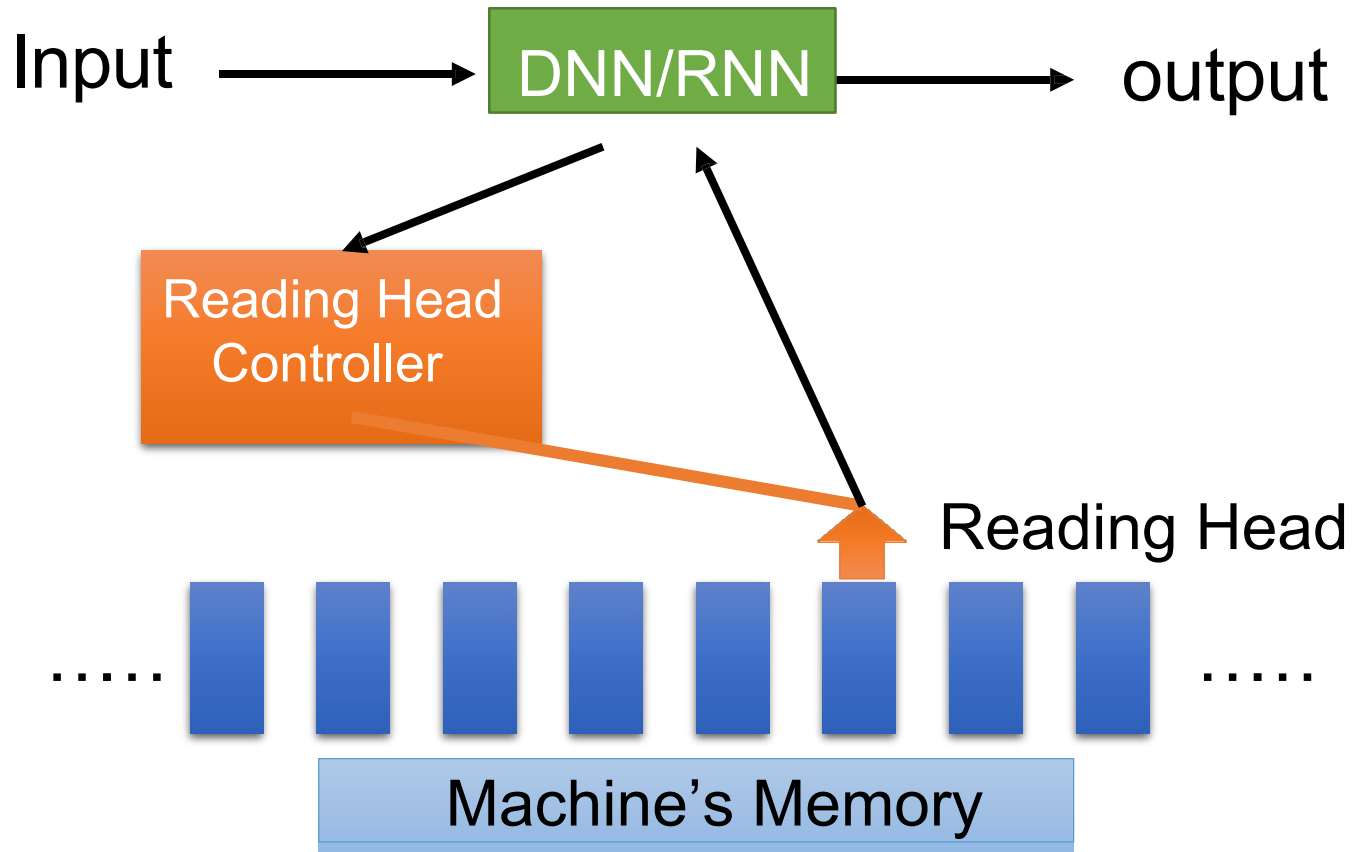
[Kelvin Xu, arXiv'15][Li Yao, ICCV'15]



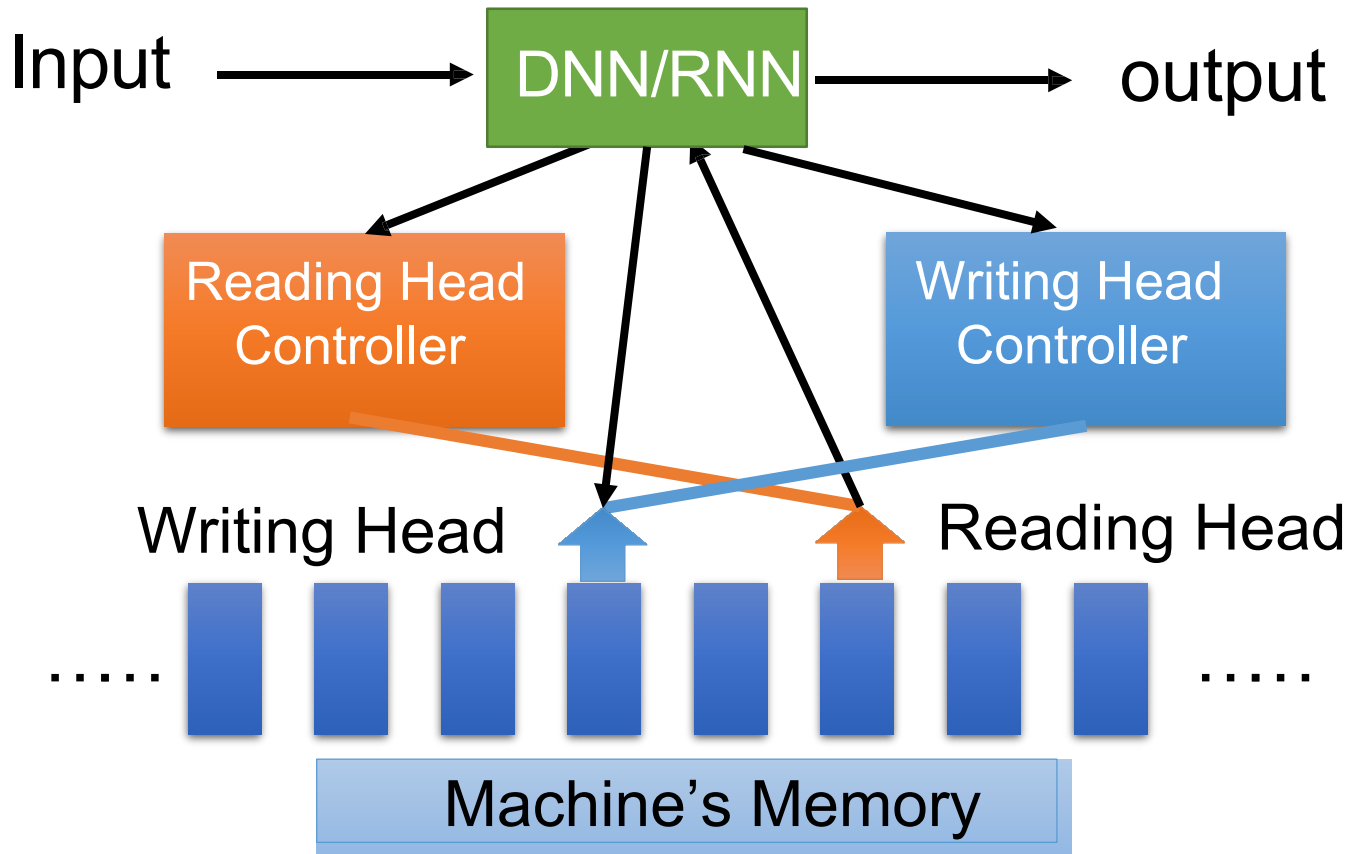
Attention-based Model



Attention-based Model

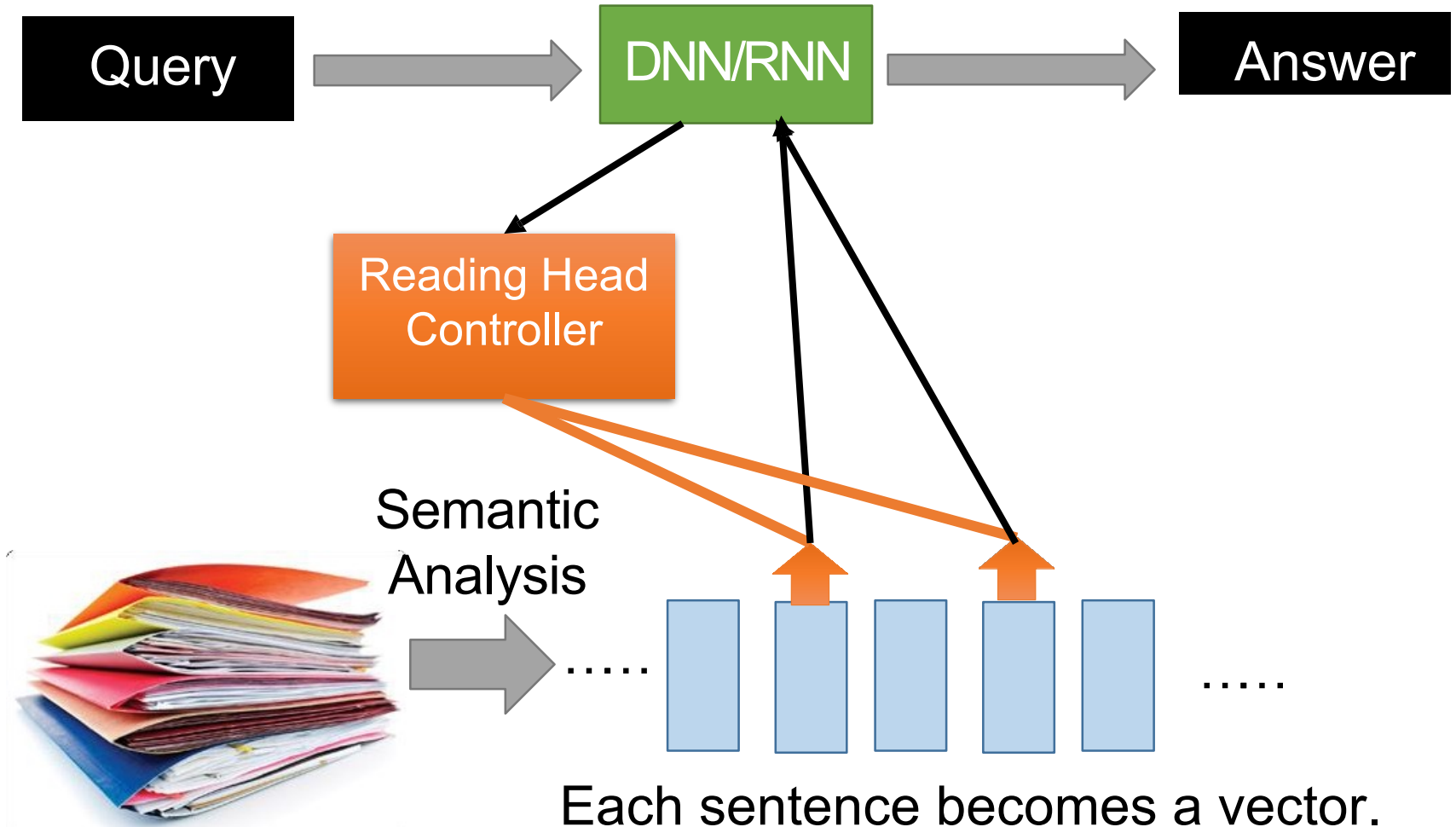


Attention-based Model v2



Neural Turing Machine

Reading Comprehension



Reading Comprehension

- End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. NIPS, 2015.

The position of reading head:

Story (16: basic induction)	Support	Hop 1	Hop 2	Hop 3
Brian is a frog.	yes	0.00	0.98	0.00
Lily is gray.		0.07	0.00	0.00
Brian is yellow.	yes	0.07	0.00	1.00
Julius is green.		0.06	0.00	0.00
Greg is a frog.	yes	0.76	0.02	0.00
What color is Greg? Answer: yellow Prediction: yellow				

Visual Question Answering



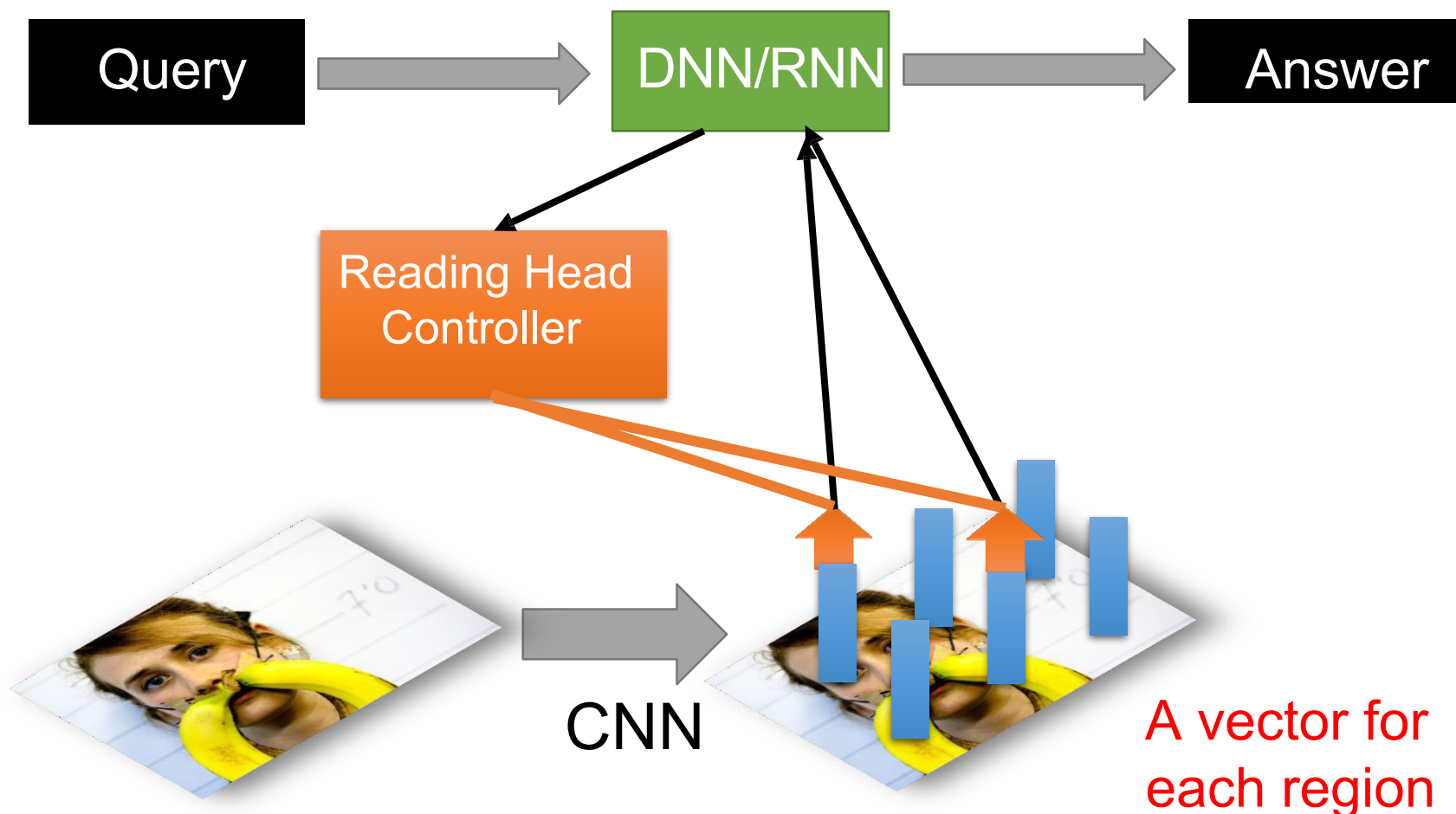
What is the mustache
made of?

AI System

bananas

source: <http://visualqa.org/>

Visual Question Answering



To Learn More

- The Unreasonable Effectiveness of Recurrent Neural Networks
 - <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- Understanding LSTM Networks
 - <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Acknowledgement

Reference and thanks to:

- **National Taiwan University ML2020 Course:**
Machine Learning

<https://speech.ee.ntu.edu.tw/~hylee/ml/2020-spring.php>

- **Understanding LSTM Networks**

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

- **LSTM Wikipedia**

https://en.wikipedia.org/wiki/Long_short-term_memory