

RNN Introduction

講者：Isaac

Outline

- ▶ Word2vect
- ▶ RNN Introduction
- ▶ LSTM/GRU Cell



Word2vect

Can machine understand the meaning of words?

“Dog” is close to “Cat”

“Love” is close to “Hate”

“Korea” is close to “America”

Word representation

What happen if

we let machine read a lot of articles and use one hot encoding on each word

"a"	"abbreviations"		"zoology"	"zoom"
1	0		0	0
0	1		0	1
0	0		0	0
.
.	.		.	.
.	.		.	.
0	0		0	0
0	0		1	0
0	0		0	1

Word representation

"a"	"abbreviations"		"zoology"	"zoom"
1	0		0	0
0	1		0	1
0	0		0	0
.
.	.		.	.
.	.		.	.
0	0		0	0
0	0		1	0
0	0		0	1

Drawbacks

1. Dog is not close to Cat
 - machine can not understand the meaning of word
2. Waste a lot of entity
 - most of entity are zero
 - dimension of vector = vocabulary size (very large)

Word representation

This is why we need better representation of words

Word representation

- ▶ We will introduce two methods in this course
 - ▶ Skip-Gram
 - ▶ CBOW

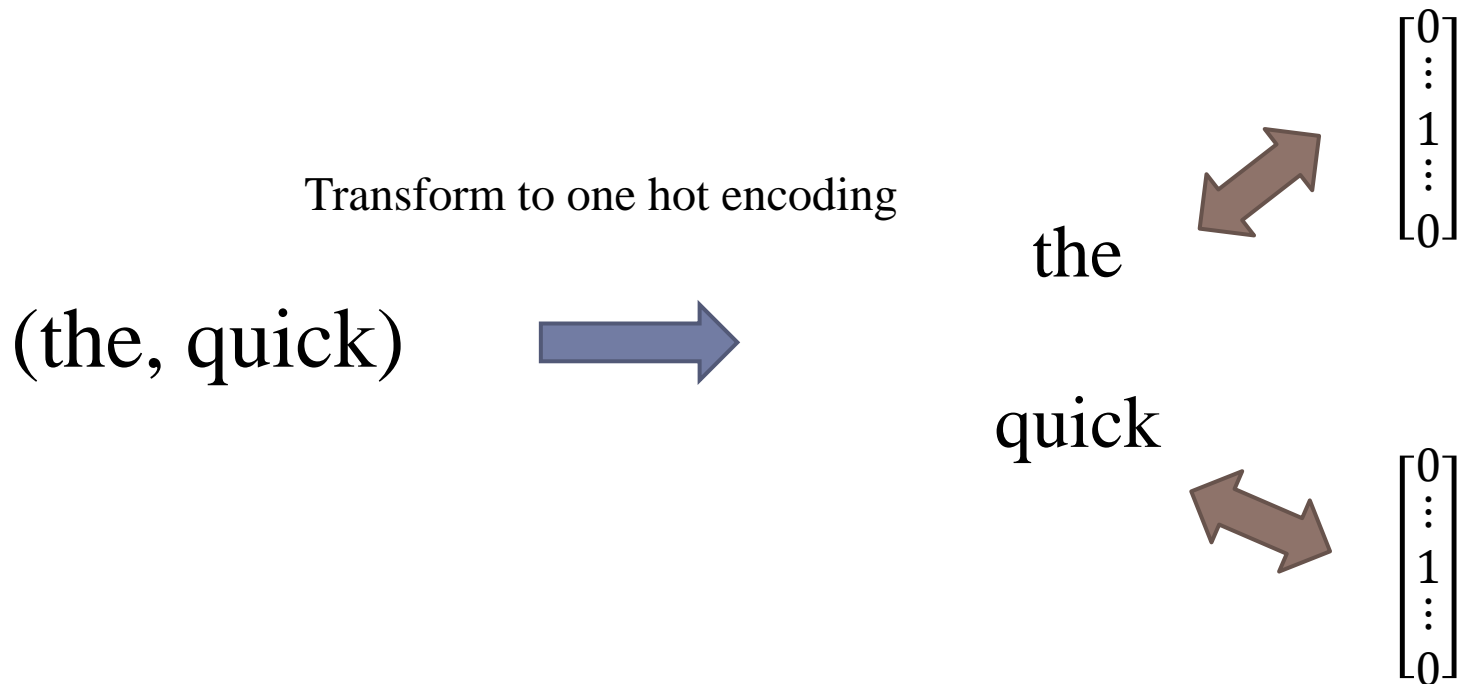
Skip-Gram

Transform many sentences into training pairs

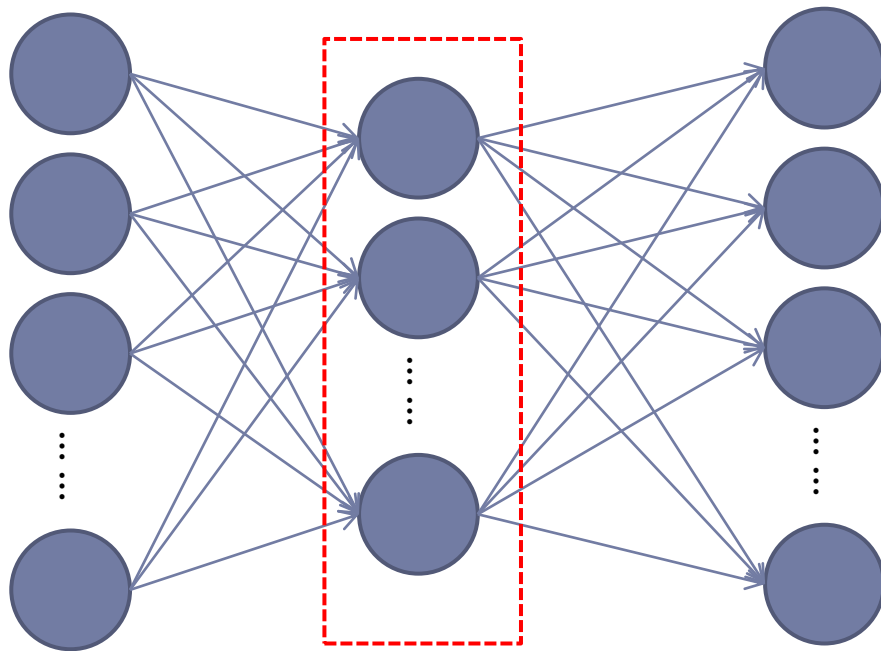
Source Text	Training Samples						
<table><tr><td>The</td><td>quick</td><td>brown</td></tr></table> fox jumps over the lazy dog. ➡	The	quick	brown	(the, quick) (the, brown)			
The	quick	brown					
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td></tr></table> jumps over the lazy dog. ➡	The	quick	brown	fox	(quick, the) (quick, brown) (quick, fox)		
The	quick	brown	fox				
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td></tr></table> over the lazy dog. ➡	The	quick	brown	fox	jumps	(brown, the) (brown, quick) (brown, fox) (brown, jumps)	
The	quick	brown	fox	jumps			
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td><td>over</td></tr></table> the lazy dog. ➡	The	quick	brown	fox	jumps	over	(fox, quick) (fox, brown) (fox, jumps) (fox, over)
The	quick	brown	fox	jumps	over		

Window size = 2

Skip-Gram

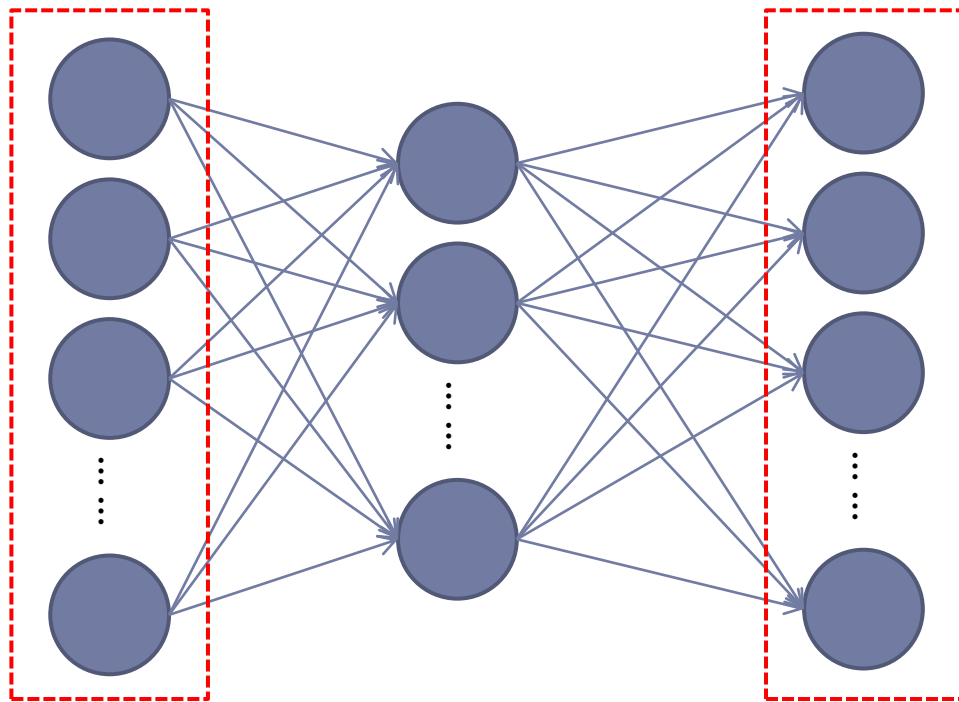


Skip-Gram



of hidden node is less than input node

Skip-Gram



of input/output node = # of vocabulary

Skip-Gram

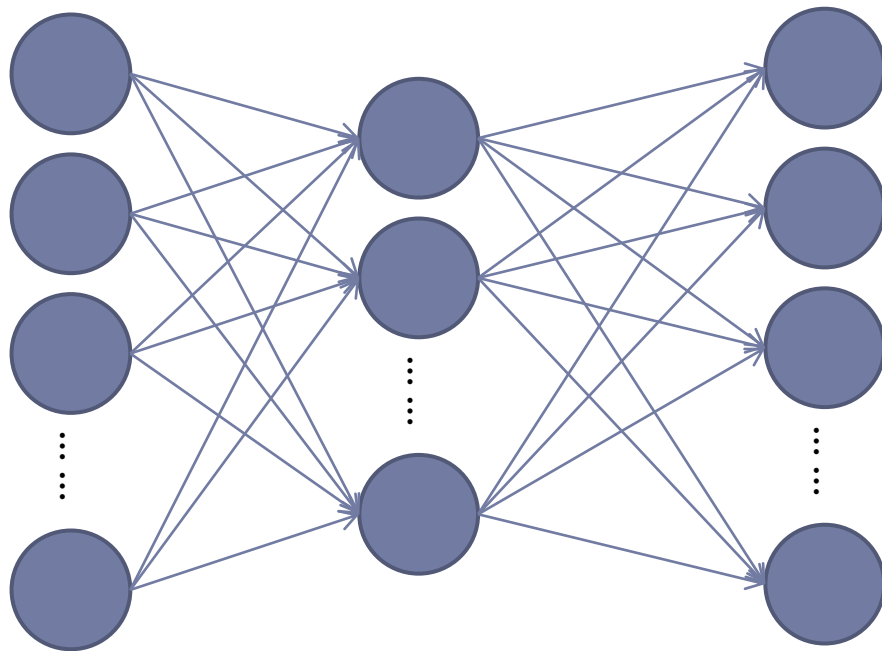
Training pairs

(the, quick)
(the, brown)

(quick, the)
(quick, brown)
(quick, fox)

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)
⋮

(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)



Softmax + cross entropy

Note that there is no activation function on hidden layer

Skip-Gram

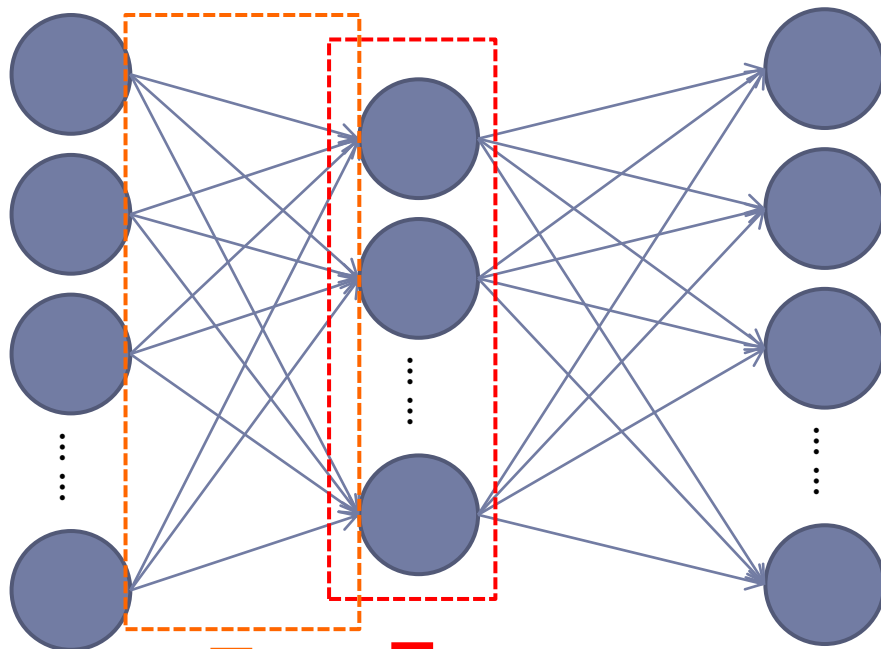
Training pairs

(the, quick)
(the, brown)

(quick, the)
(quick, brown)
(quick, fox)

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

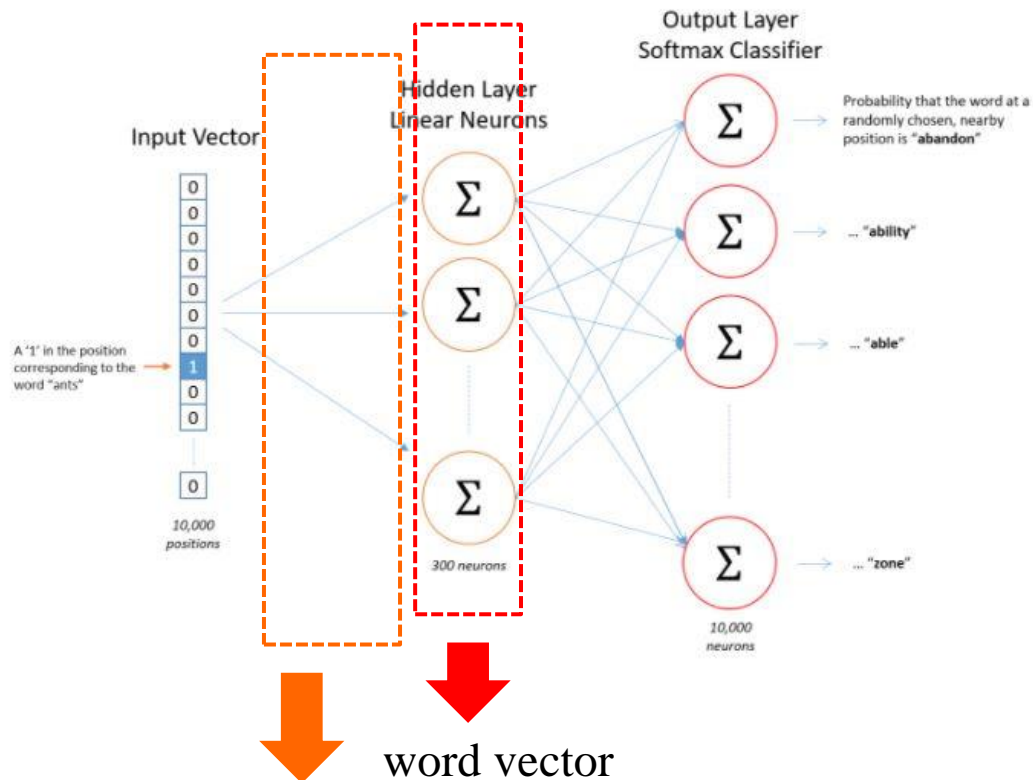
(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)



word vector

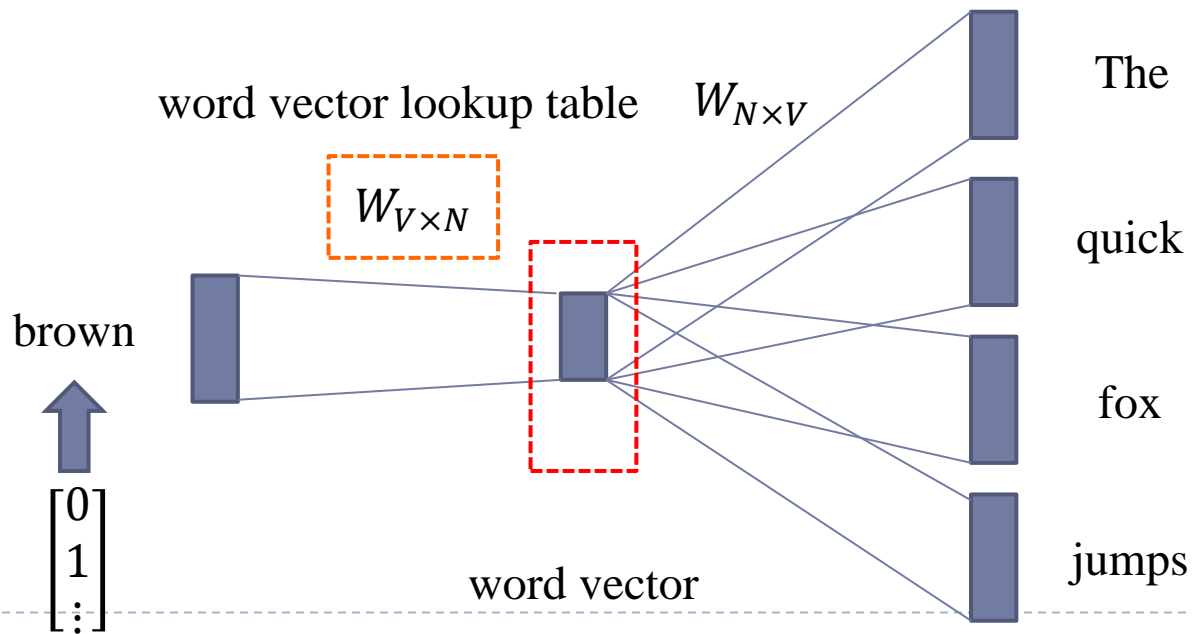
word vector lookup table

Skip-Gram

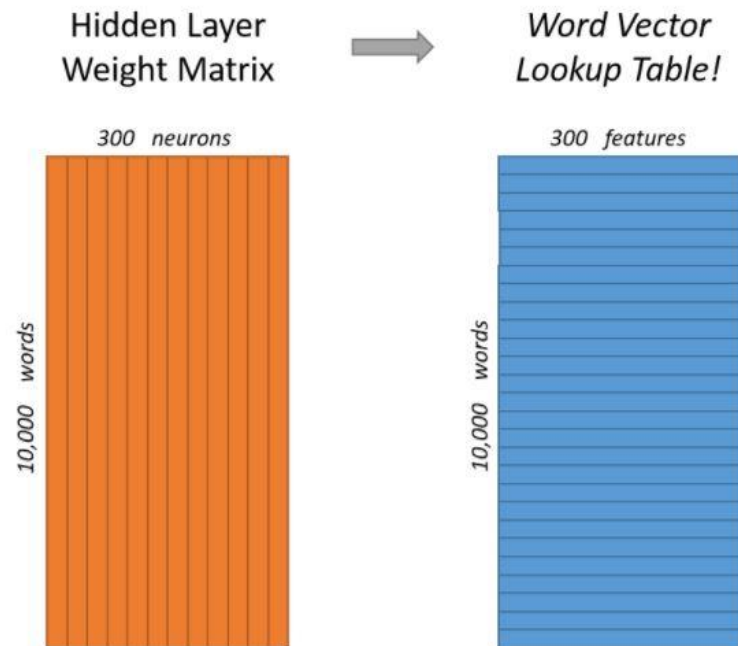


Skip-Gram

The quick brown fox jumps over the lazy dog



Skip-Gram

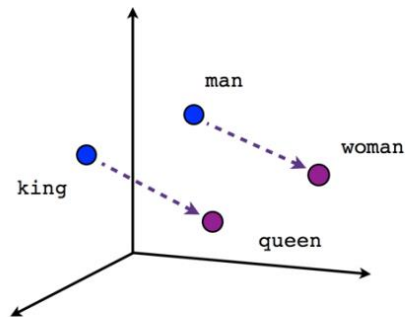


$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$

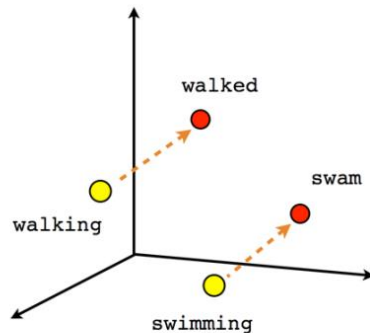
Assume vocabulary = 10000

Small example

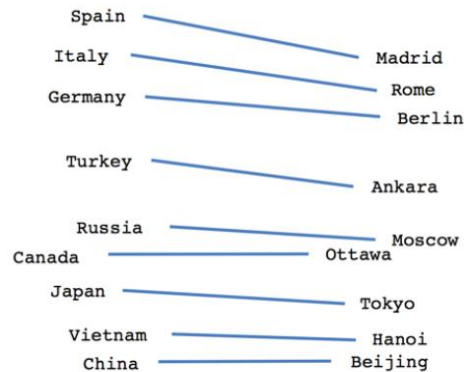
Skip-Gram



Male-Female



Verb tense



Country-Capital

$$\text{vector}(\text{king}) - \text{vector}(\text{queen}) = \text{vector}(\text{man}) - \text{vector}(\text{woman})$$

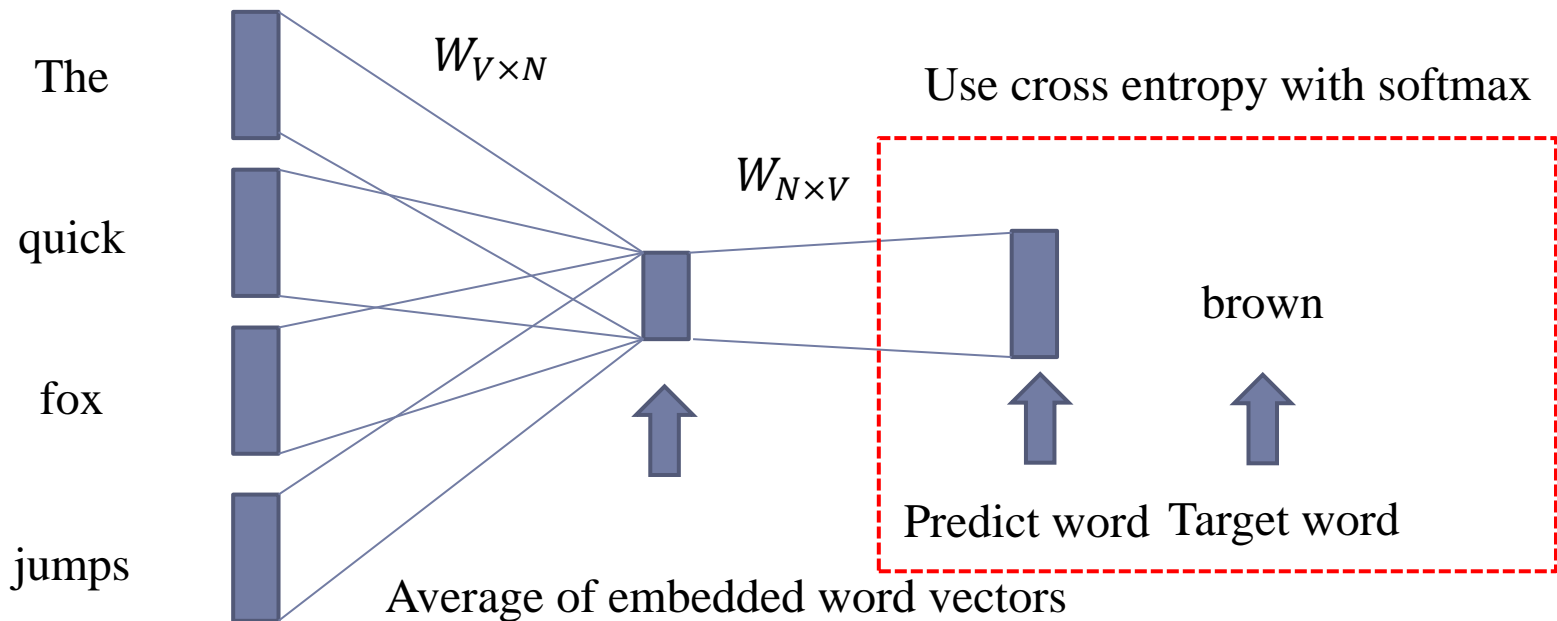
$$\text{vector}(\text{walking}) - \text{vector}(\text{walked}) = \text{vector}(\text{swimming}) - \text{vector}(\text{swam})$$

$$\text{vector}(\text{Spain}) - \text{vector}(\text{Italy}) = \text{vector}(\text{Madrid}) - \text{vector}(\text{Rome})$$

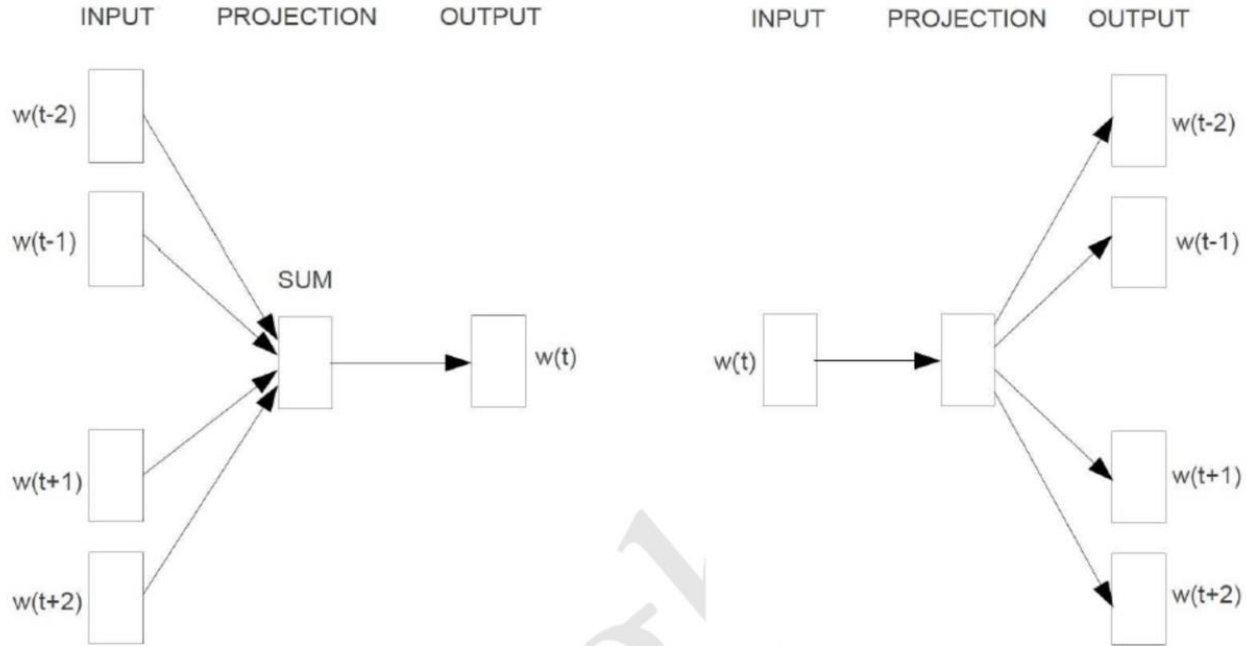


CBOW

The quick brown fox jumps over the lazy dog



Skip-Gram V.S. CBOW



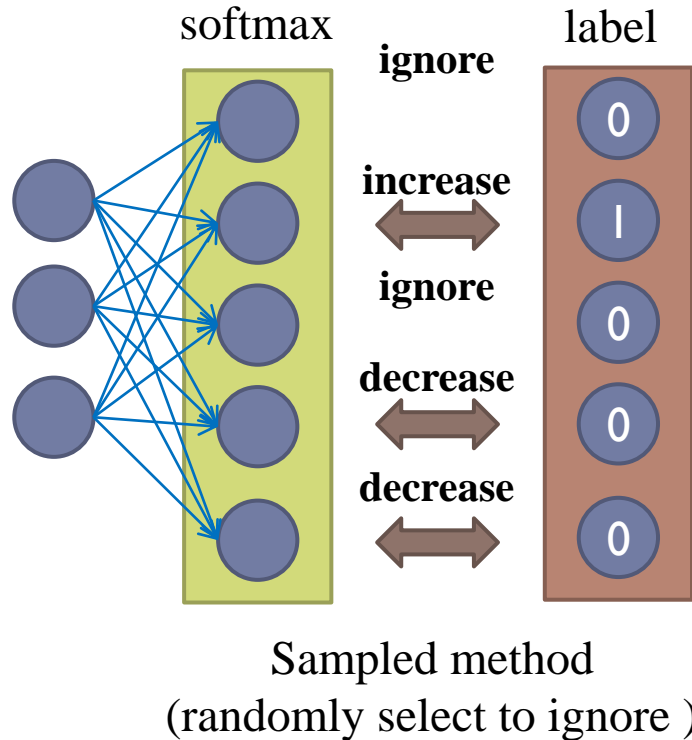
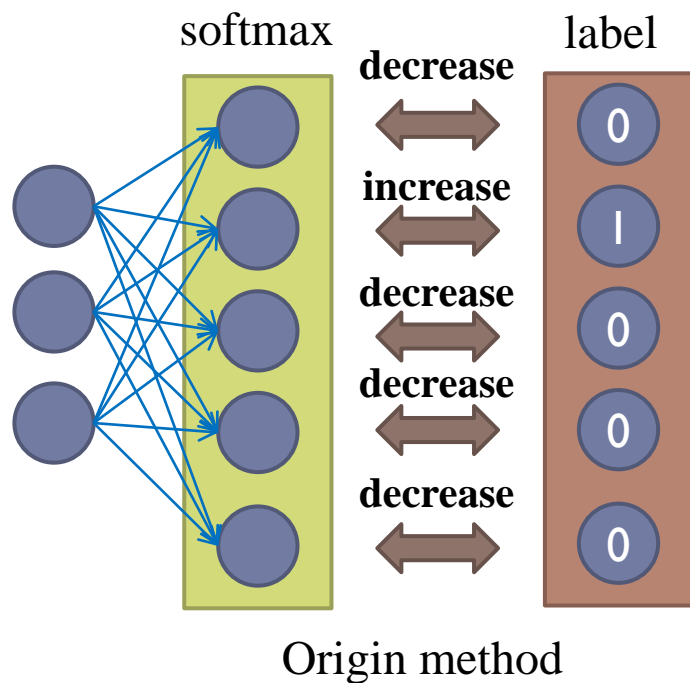
Cost function in Word2vect

- ▶ **Softmax function**
 - ▶ Computation expensive
 - ▶ Use sampled version of softmax function

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

K = number of vocabulary (very large)

Cost function in Word2vect



Perplexity

- ▶ A measurement of how well a probability distribution predict a sample

$$2^{H(p)} = 2^{-\sum_x p(x) \log_2 p(x)}$$

RNN Introduction

Why RNN?

- ▶ Position of words is important
- ▶ Slightly change a word may change meaning of whole sentence

I am handsome



Am I handsome

I will **leave** Taiwan in January



I will **arrive** Taiwan in January

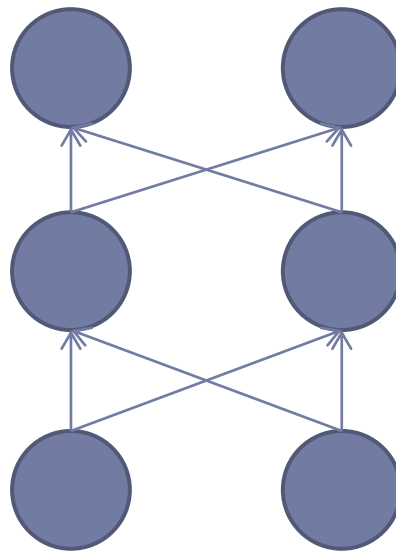


Why RNN?

Lack of sequence concept in ordinary NN

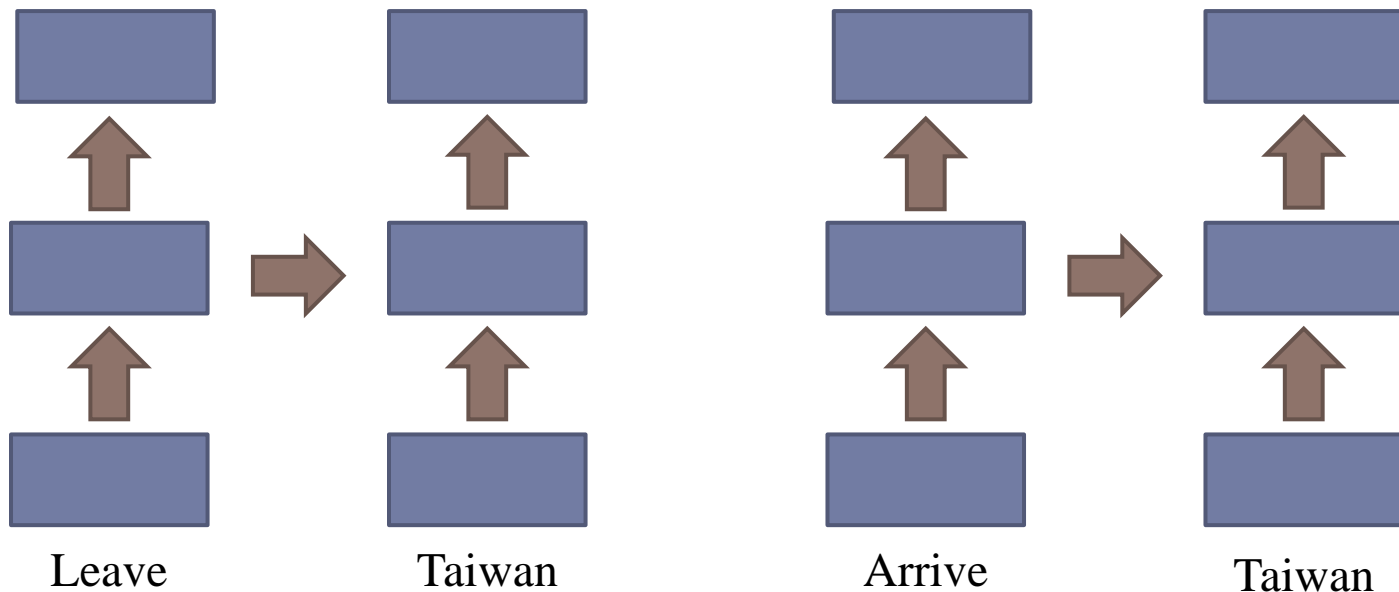
I will **leave** Taiwan on January

I will **arrive** Taiwan on January



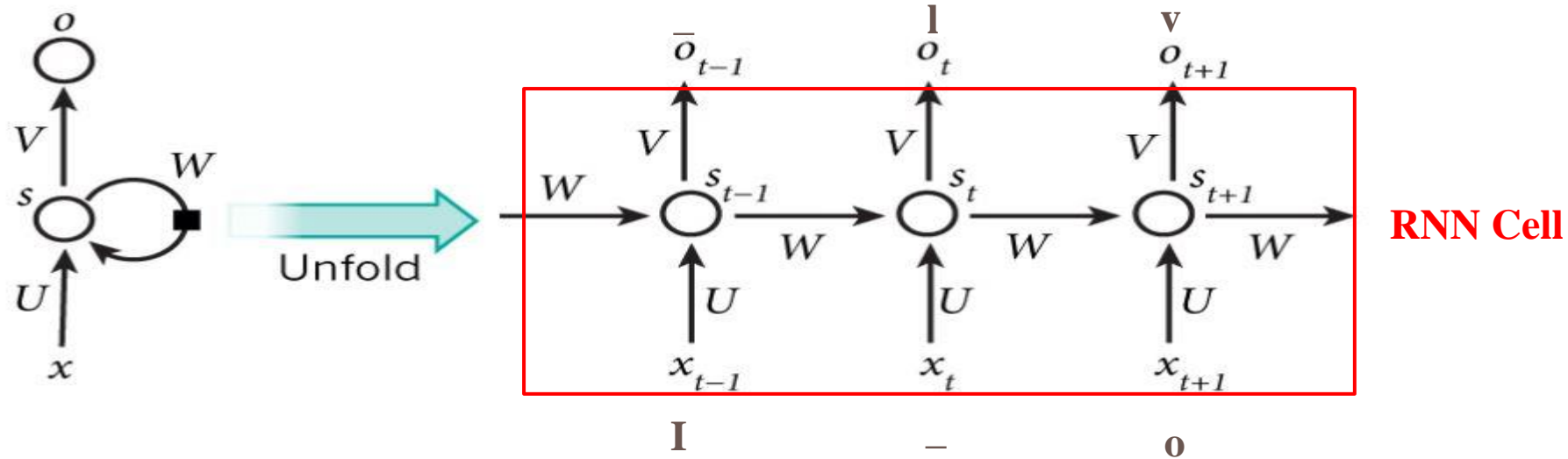
Taiwan

Why RNN?



This type of structure contain sequence concept

Memory in RNN



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

$$o_t = \text{softmax}(Vs_t).$$

Input X_t can be as following:

I love you

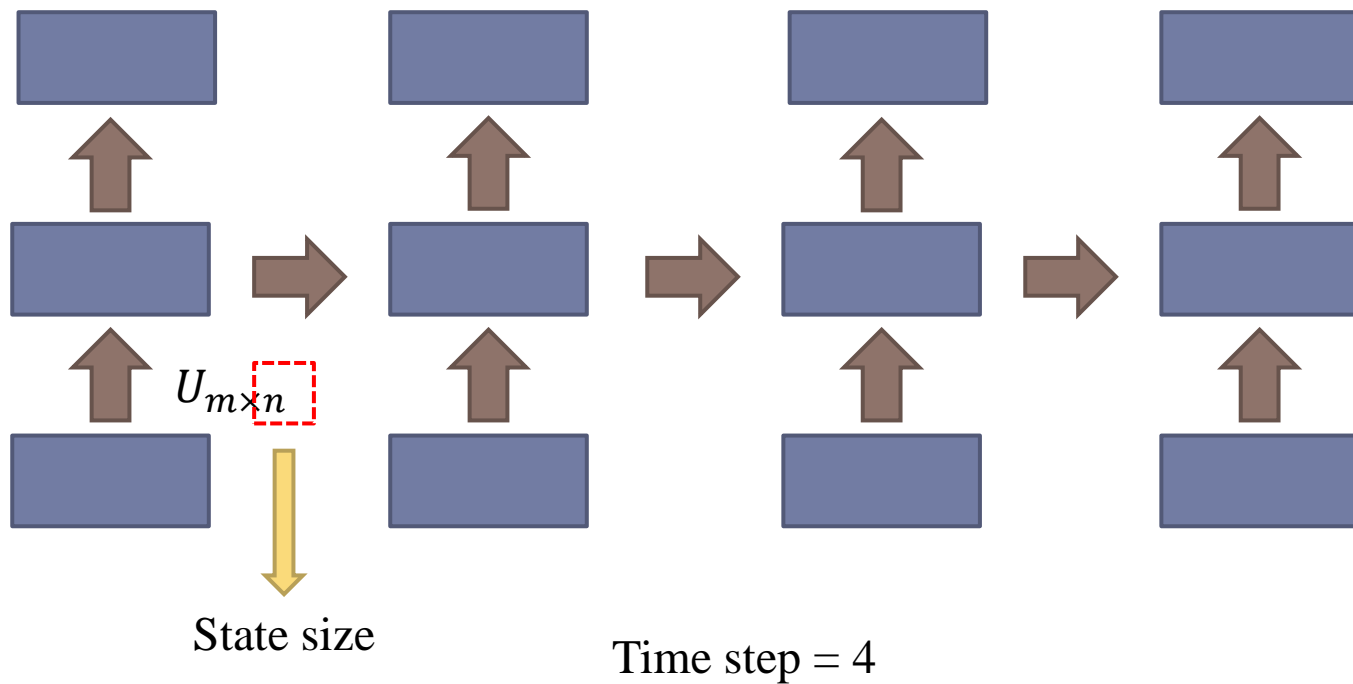
I love you

I love you

I love you (word based)

Etc

RNN



Example

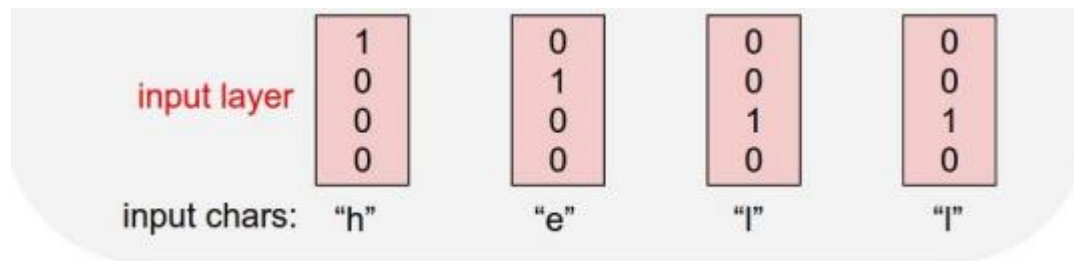
- ▶ We would like our computer to read a novel
 - ▶ One common solution is feed data into RNN character by character
 - ▶ Also called character-level language model

hello, I am Isaac



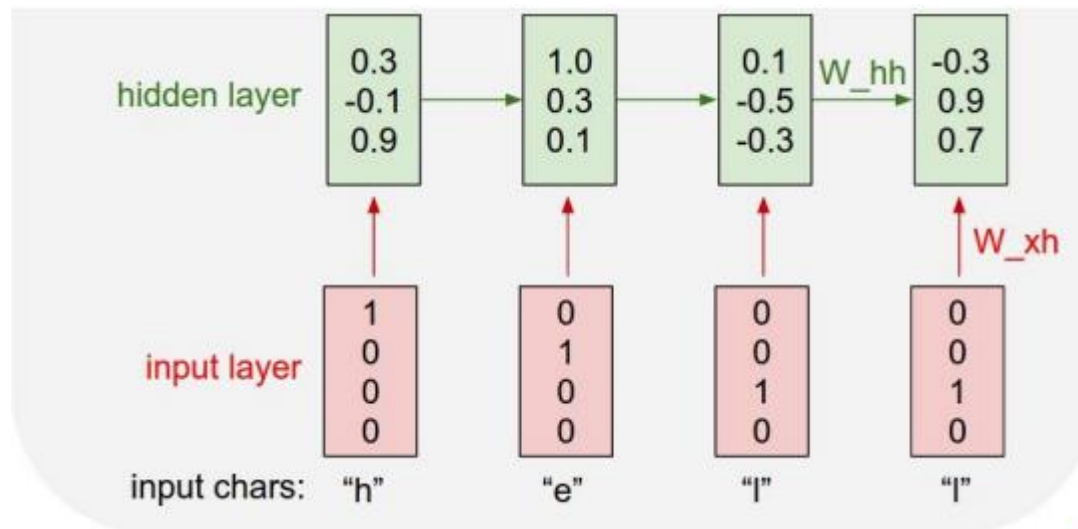
'h' 'e' 'l' 'l' 'o' ',' ' ' 'I' ' ' 'a' 'm' ' ' 'I' 's' 'a' 'a' 'c'

Example



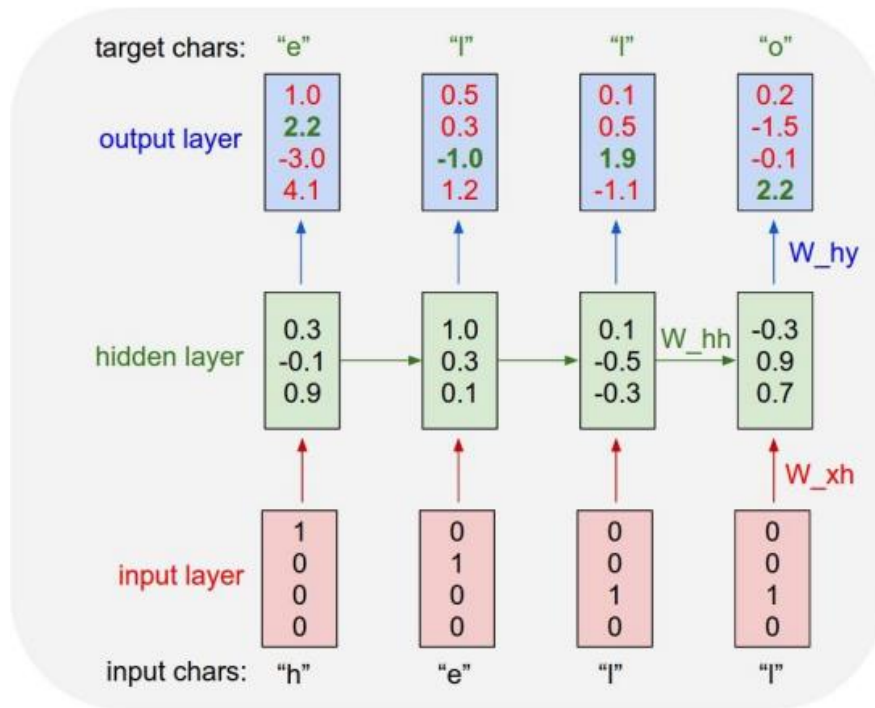
Character-level Language Model

Example

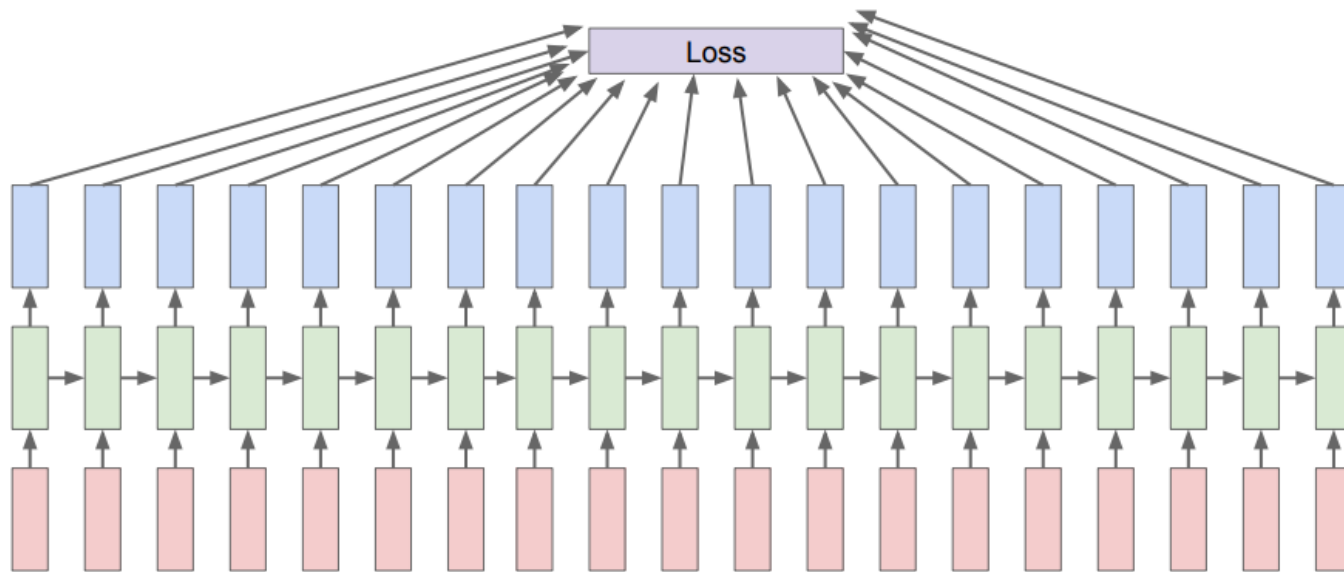


Character-level Language Model

Example



Example



Character-level Language Model

Example

- ▶ After training, RNN can learn relationship of characters and even generate it

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.

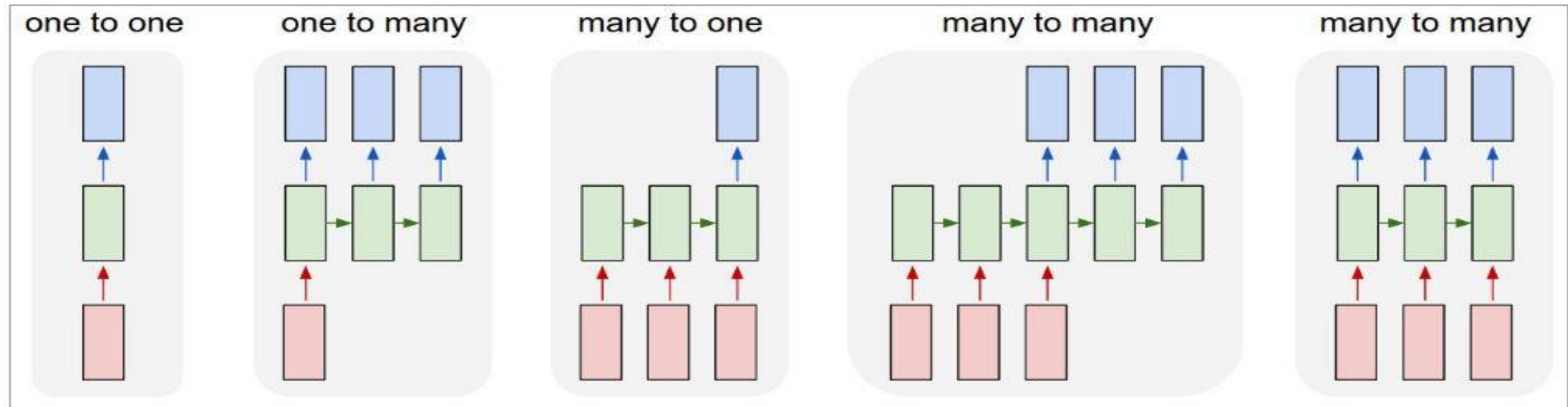
DUKE VINCENTIO:
Well, your wit is in the care of side and that.

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.

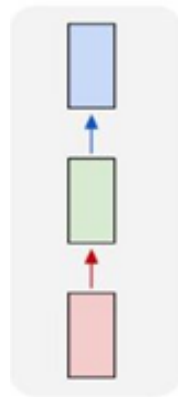
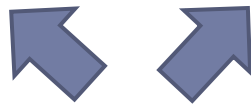
VIOLA:
I'll drink it.

Different Type of RNN Application



Example

Cat Dog



Classification task

One to one (DNN)



Example

This is a cat

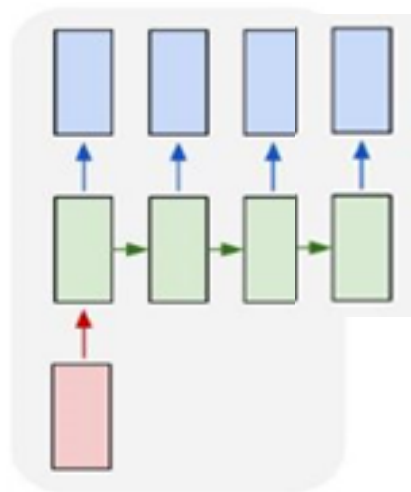


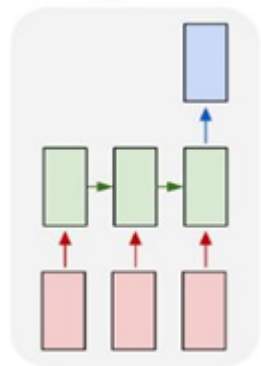
Image Caption

One to many



Example

Good/Bad ?

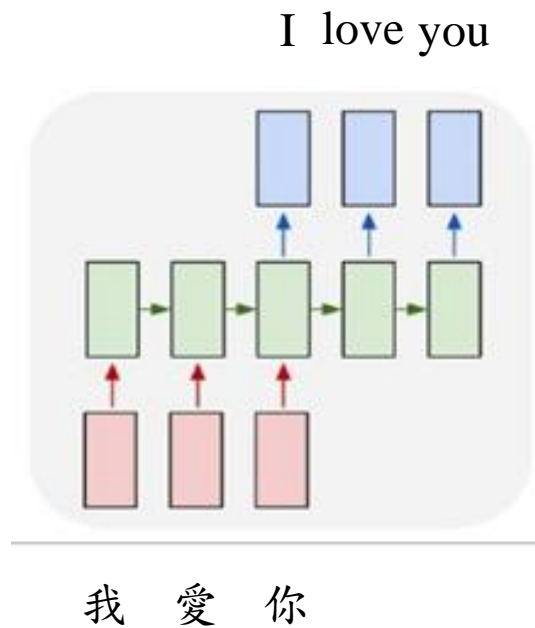


Movie Review
(positive/negative)

Many to one

It is good

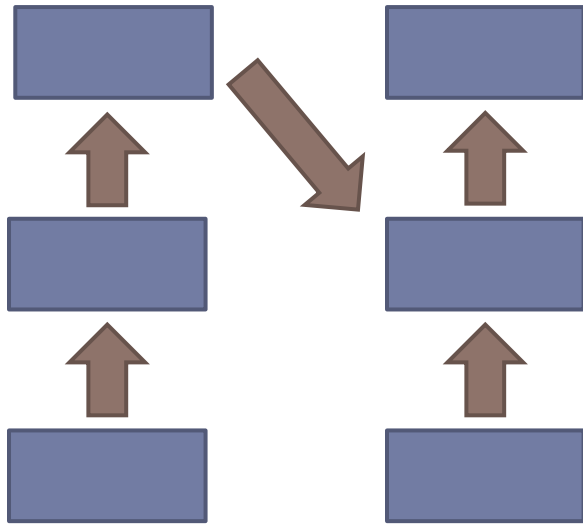
Example



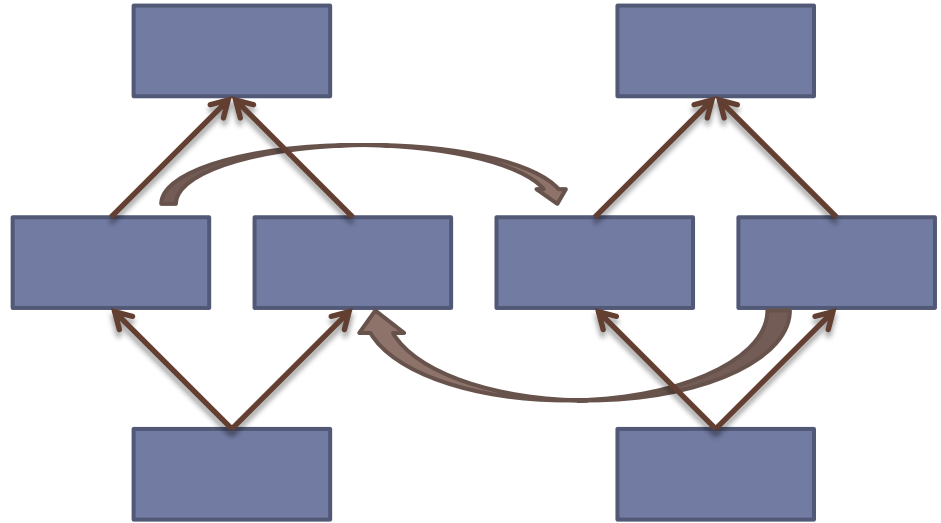
Language translation

Many to many

Different Type of RNN Structure

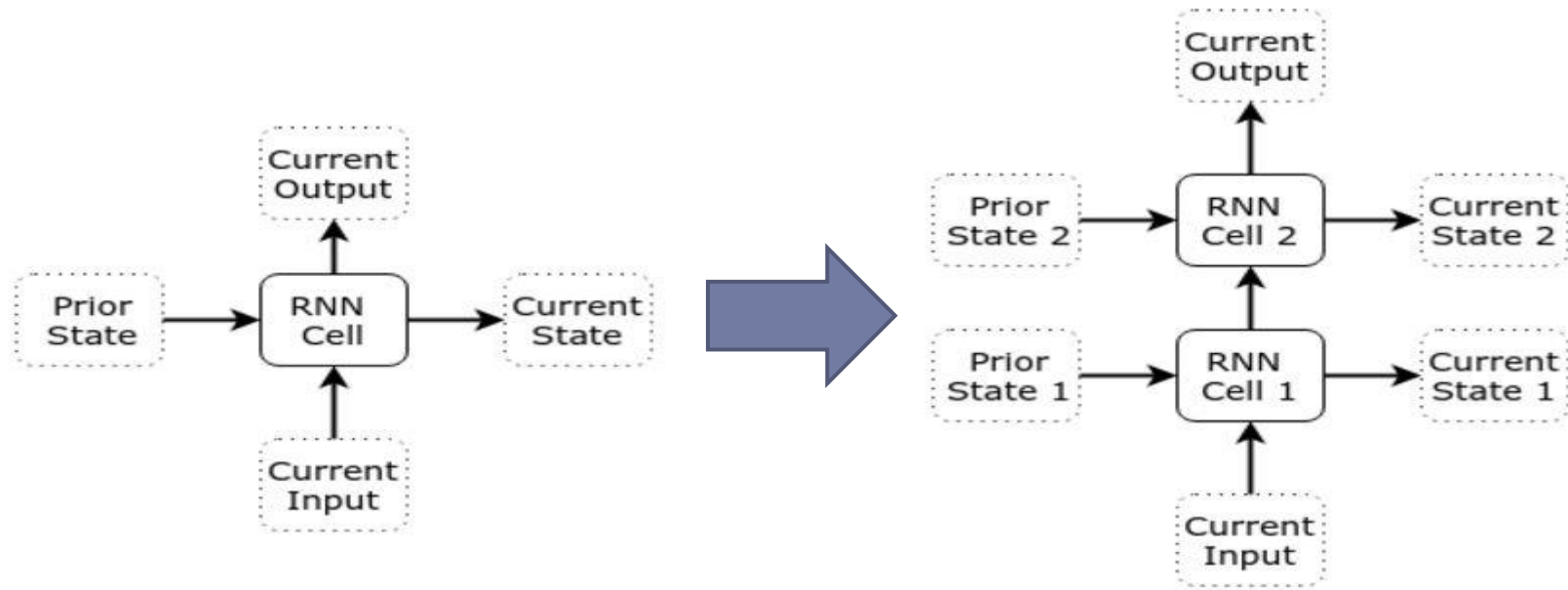


Jordan Network

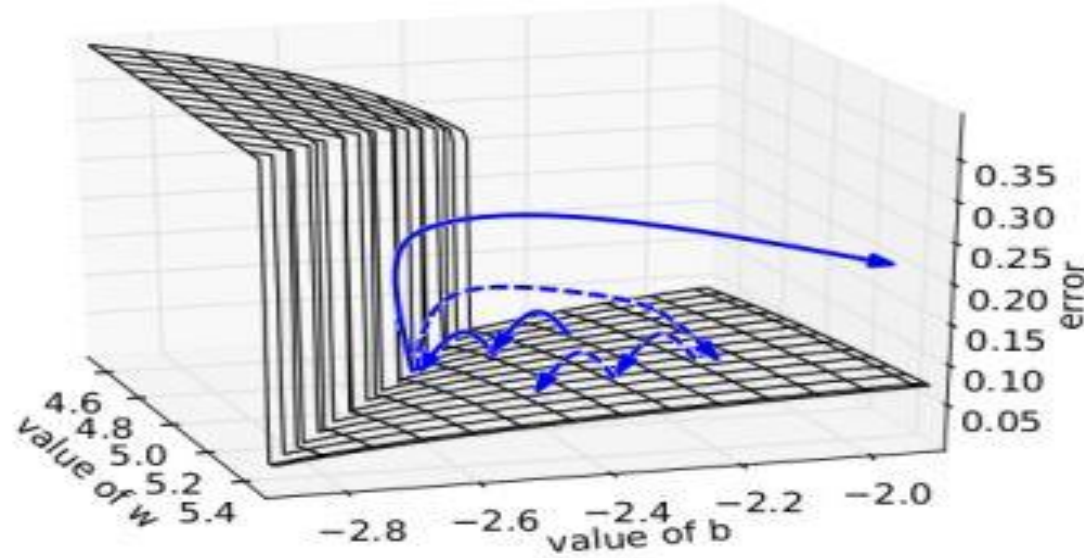


Bidirectional RNN

Go Deep in RNN

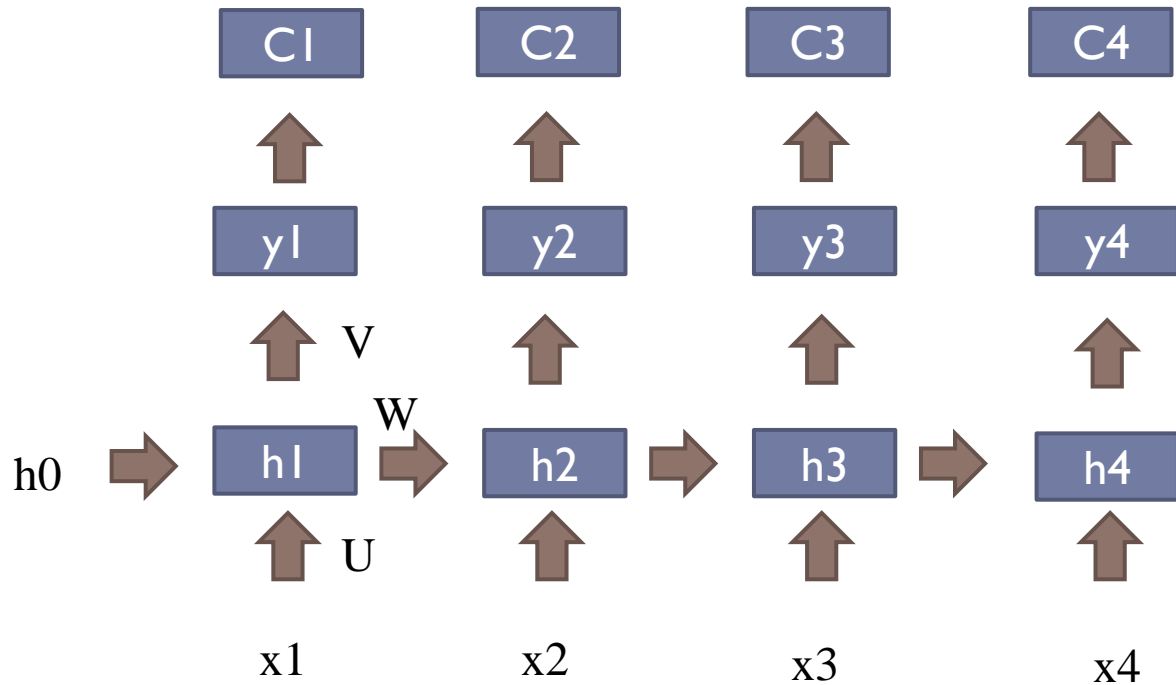


Error surface in RNN



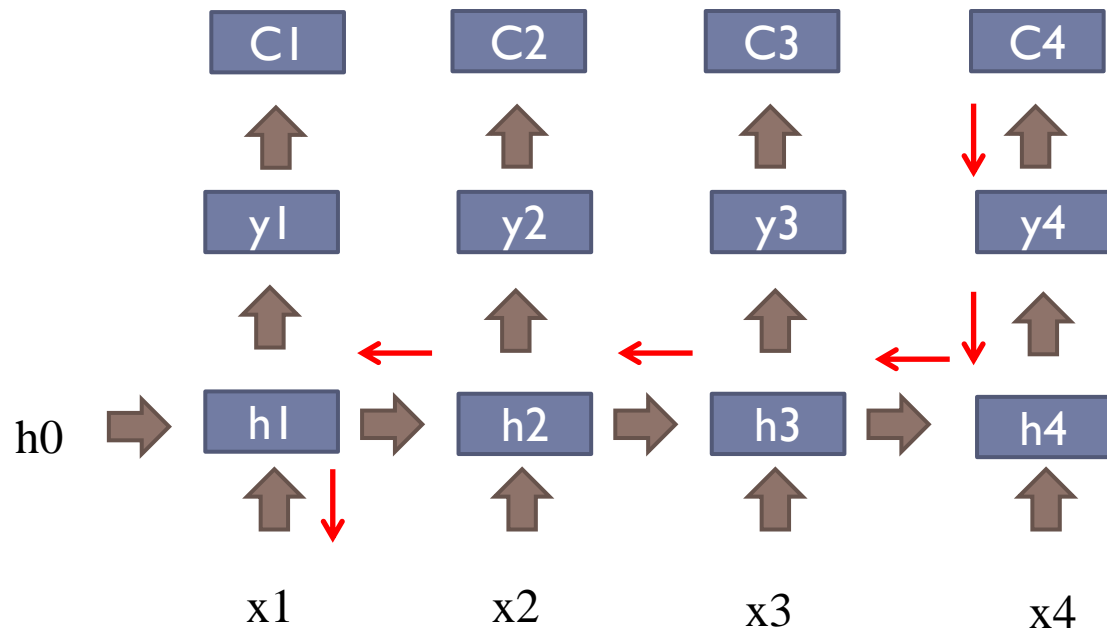
Very flat or very steep on error surface in practice

Learning RNN



BPTT

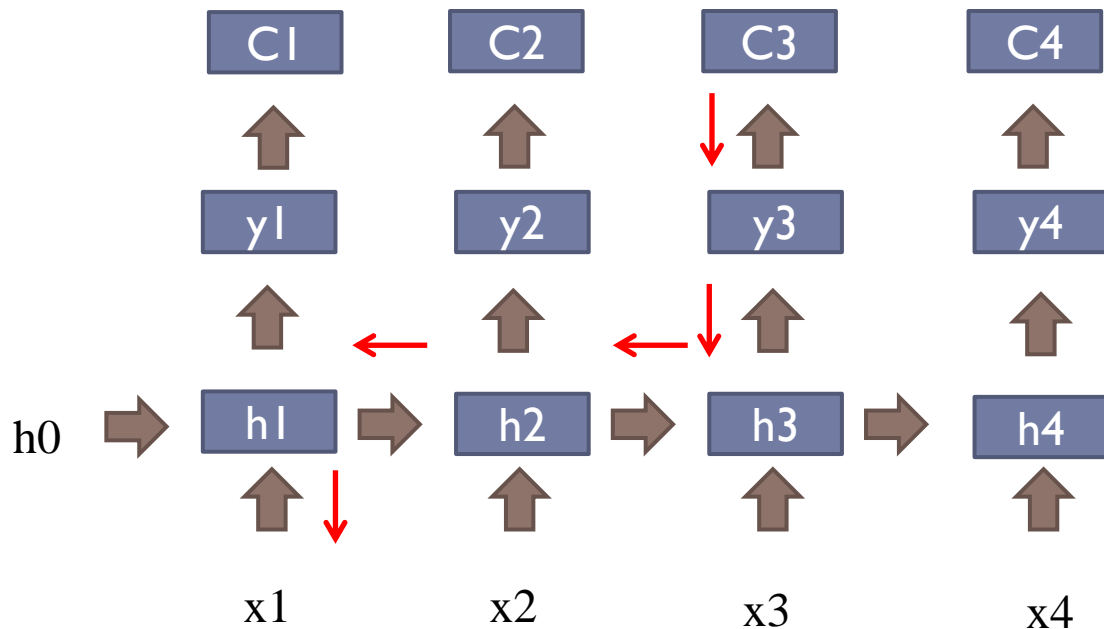
backpropagation through time (BPTT)



$$\frac{\partial C_4}{\partial w} = \frac{\partial C_4}{\partial y_4} * \frac{\partial y_4}{\partial h_4} * \frac{\partial h_4}{\partial h_3} * \frac{\partial h_3}{\partial h_2} * \frac{\partial h_2}{\partial h_1} * \frac{\partial h_1}{\partial w}$$

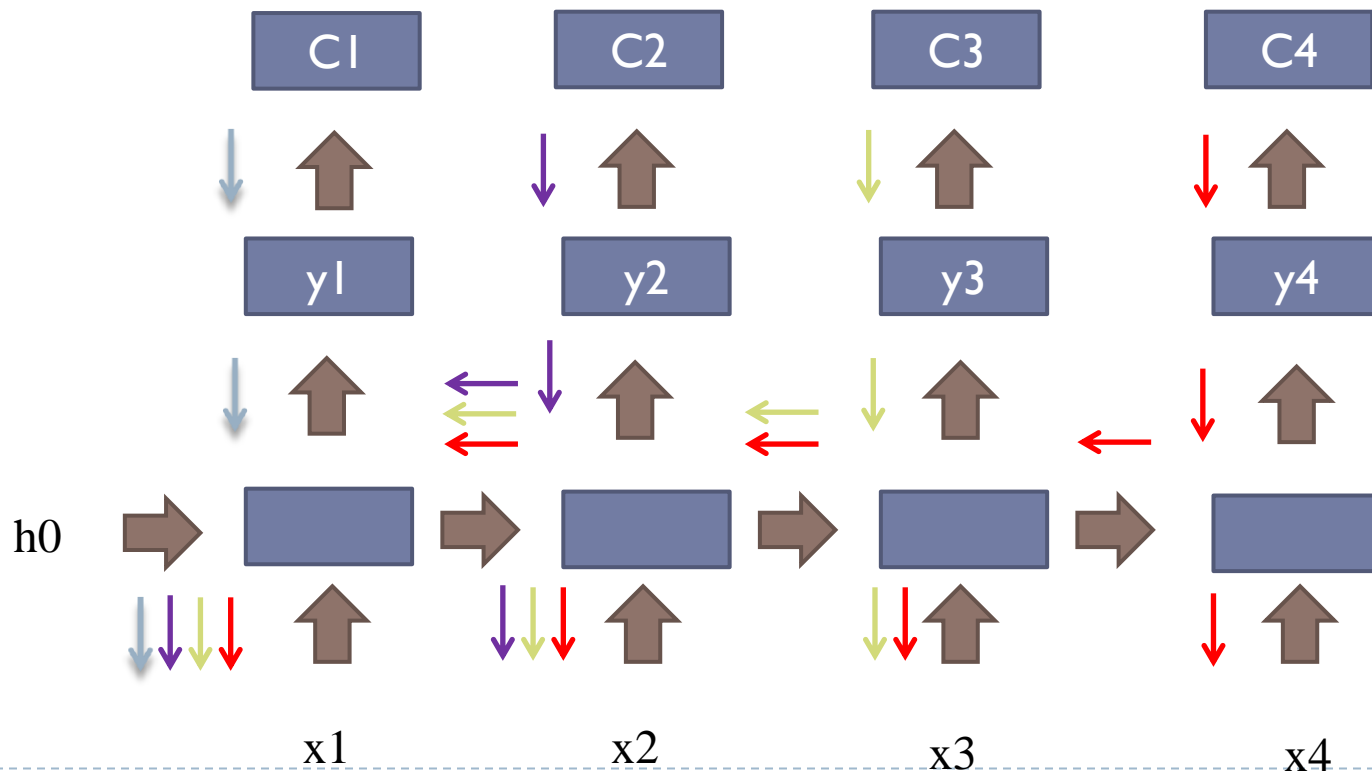
BPTT

backpropagation through time (BPTT)



$$\frac{\partial C_3}{\partial h_1} = \frac{\partial C_3}{\partial y_3} * \frac{\partial y_3}{\partial h_3} * \frac{\partial h_3}{\partial h_2} * \frac{\partial h_2}{\partial h_1} * \frac{\partial h_1}{\partial w}$$

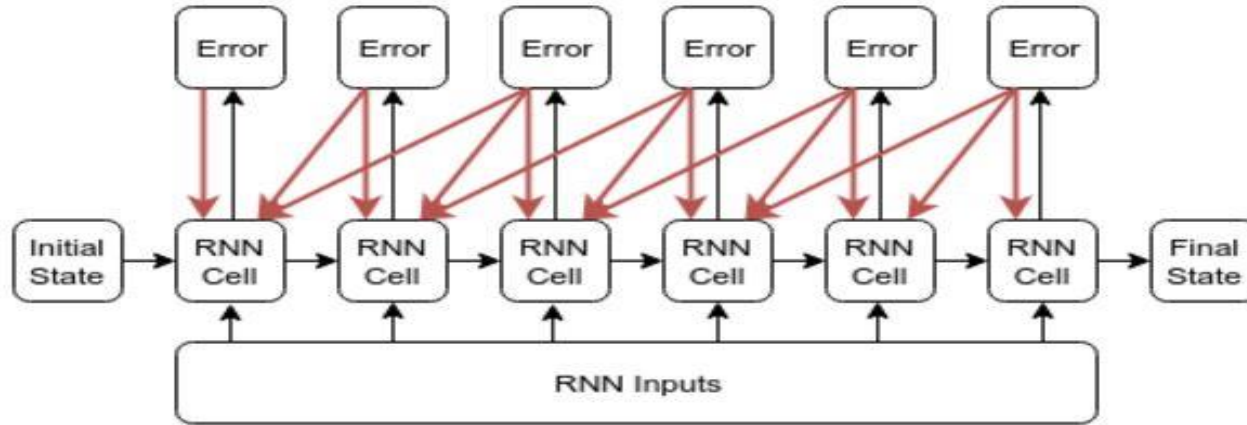
BPTT



Some problem in BPTT

- ▶ Inefficiency on propagation under long time step RNN structure
 - ▶ Truncated BPTT
- ▶ Gradient vanish/exploding problem
 - ▶ Main reason why RNN fail in early day
 - ▶ New kind of RNN cell

Truncated BPTT



Truncated backpropagation

Gradient Vanish/Exploding problem

$$\frac{\partial C_t}{\partial h_1} = \frac{\partial C_t}{\partial y_t} * \frac{\partial y_t}{\partial h_t} * \boxed{\frac{\partial h_t}{\partial h_{t-1}} * \frac{\partial h_{t-1}}{\partial h_{t-2}} * \dots * \frac{\partial h_3}{\partial h_2} * \frac{\partial h_2}{\partial h_1}}$$



- If very big (gradient exploding)
 - Clip the value

Gradient Vanish/Exploding problem

$$\frac{\partial C_t}{\partial h_1} = \frac{\partial C_t}{\partial y_t} * \frac{\partial y_t}{\partial h_t} * \boxed{\frac{\partial h_t}{\partial h_{t-1}} * \frac{\partial h_{t-1}}{\partial h_{t-2}} * \dots * \frac{\partial h_3}{\partial h_2} * \frac{\partial h_2}{\partial h_1}}$$



- If very small (gradient vanish)
 - need $\frac{\partial h_n}{\partial h_{n-1}}$ to be constant

Gradient Vanish/Exploding problem

How to avoid $\frac{\partial h_n}{\partial h_{n-1}}$ too small ?



If two hidden state is recursive

$$h_n = h_{n-1} + \dots \dots$$

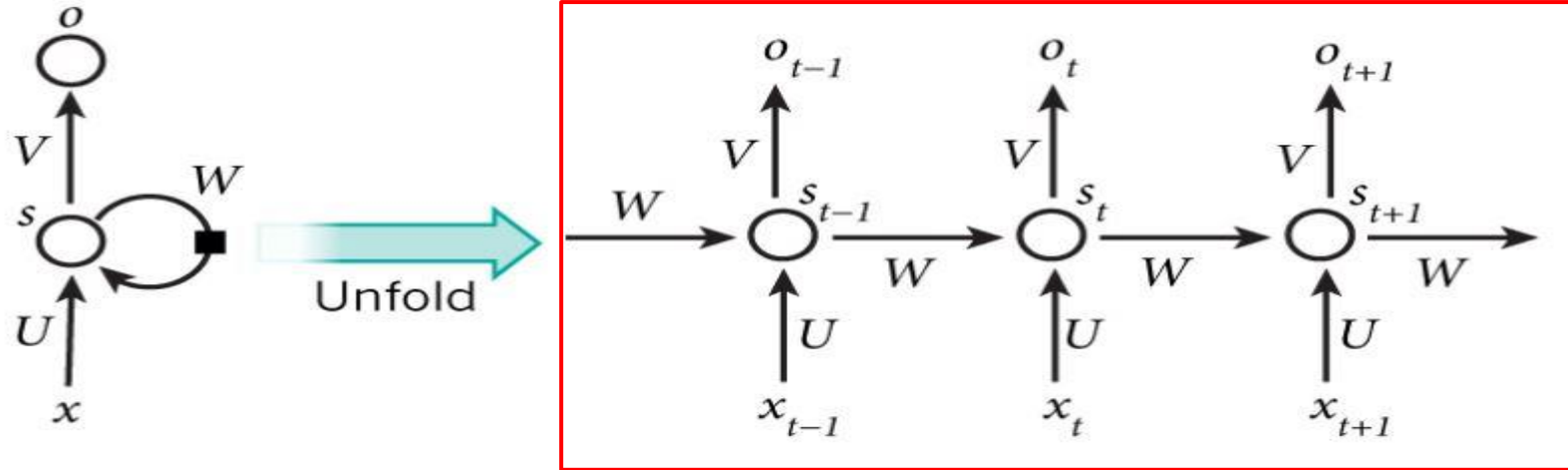


$$\frac{\partial C_t}{\partial h_1} = \frac{\partial C_t}{\partial y_t} * \frac{\partial y_t}{\partial h_t} * \boxed{\frac{\partial h_t}{\partial h_{t-1}} * \frac{\partial h_{t-1}}{\partial h_{t-2}} * \dots * \frac{\partial h_3}{\partial h_2} * \frac{\partial h_2}{\partial h_1}}$$



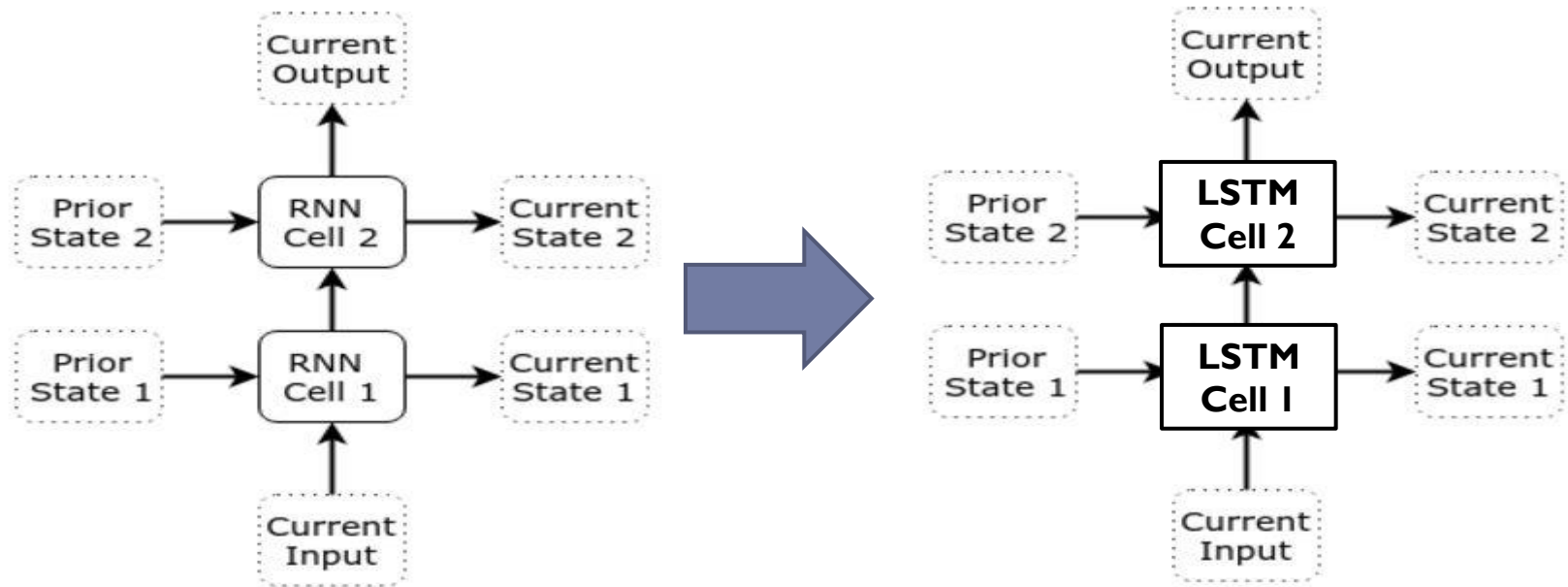
LSTM/GRU Cell

Recall RNN Cell

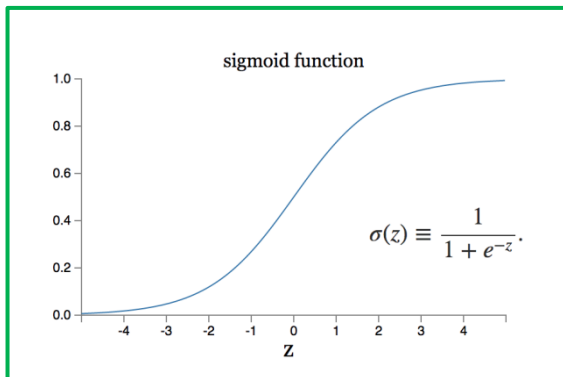
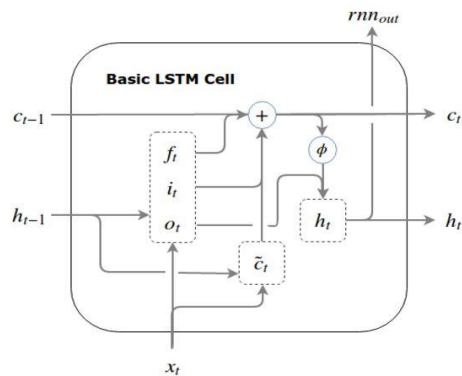


RNN Cell

LSTM Cell



LSTM



$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

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

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

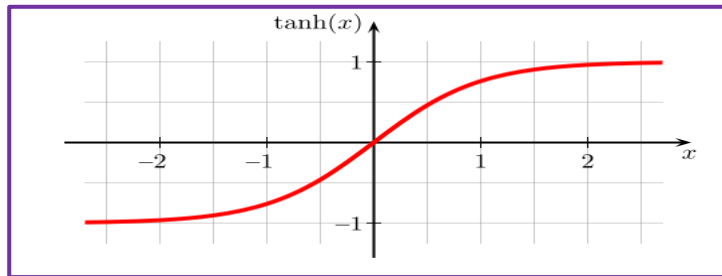
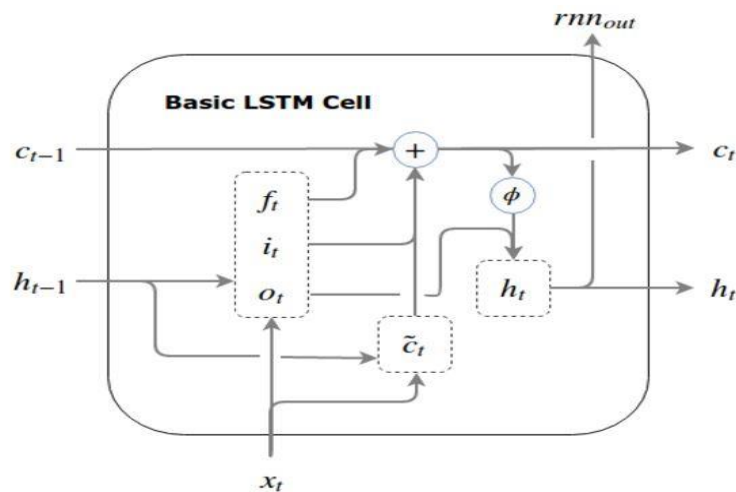
$$\tilde{c}_t = \phi(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot \phi(c_t)$$

$$\text{rnn}_{out} = h_t$$

LSTM



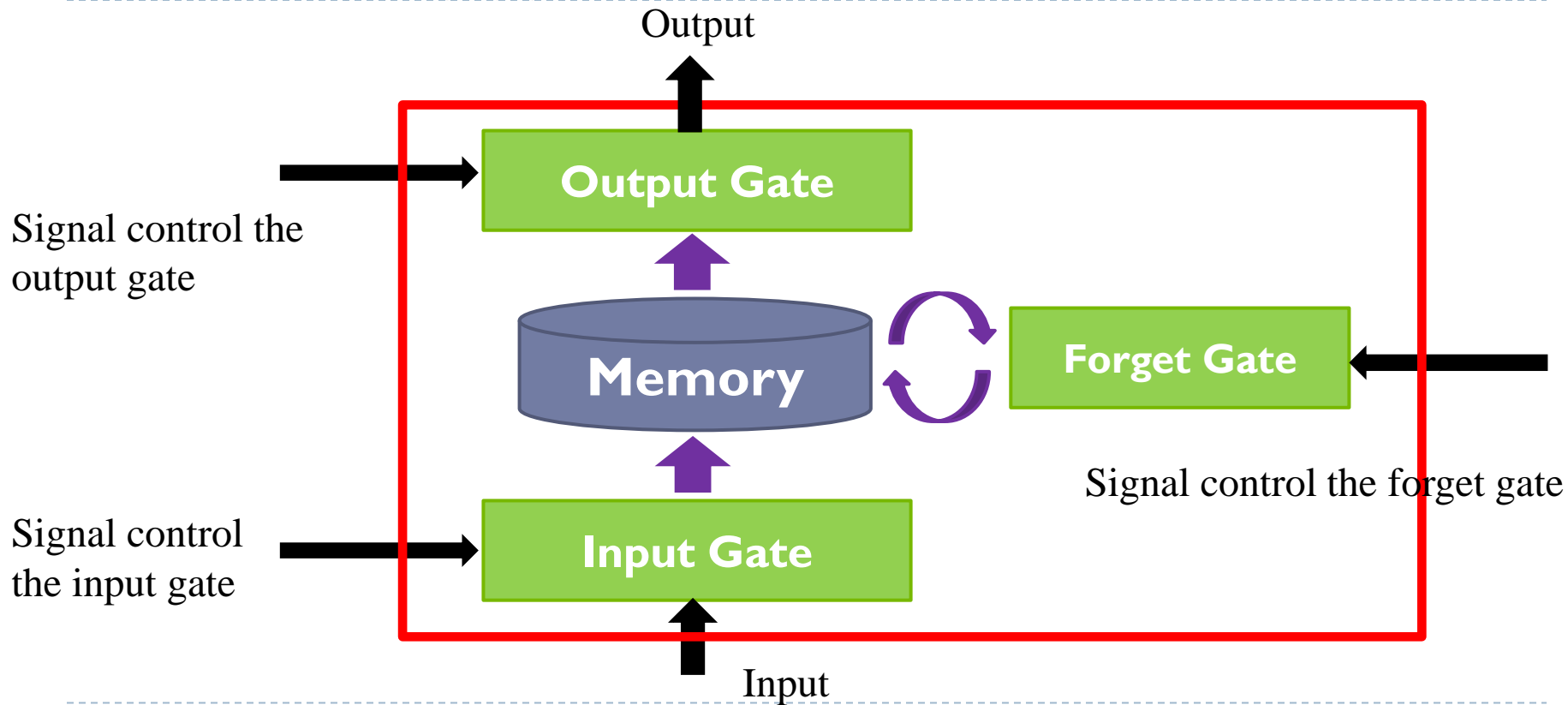
$$\begin{aligned}i_t &= \sigma(W_i h_{t-1} + U_i x_t + b_i) \\o_t &= \sigma(W_o h_{t-1} + U_o x_t + b_o) \\f_t &= \sigma(W_f h_{t-1} + U_f x_t + b_f)\end{aligned}$$

$$\begin{aligned}\tilde{c}_t &= \phi(W h_{t-1} + U x_t + b) \\c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t\end{aligned}$$

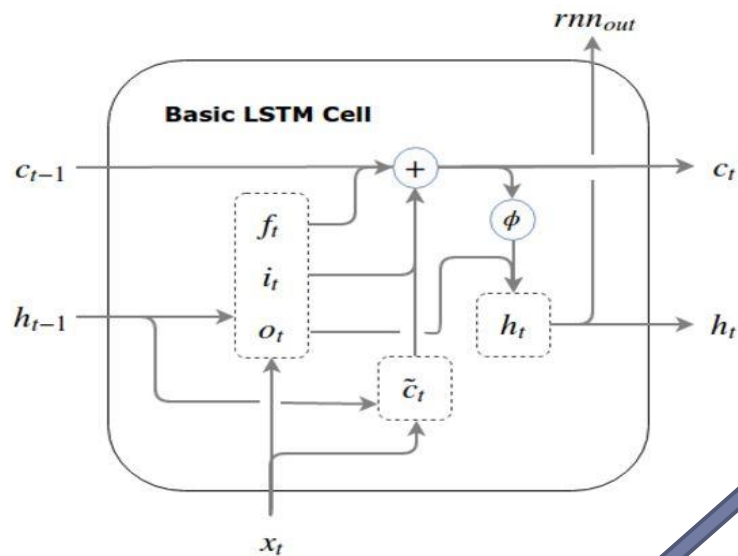
$$h_t = o_t \odot \phi(c_t)$$

$$rnn_{out} = h_t$$

LSTM Structure



LSTM



$$\begin{aligned}i_t &= \sigma(W_i h_{t-1} + U_i x_t + b_i) \\o_t &= \sigma(W_o h_{t-1} + U_o x_t + b_o) \\f_t &= \sigma(W_f h_{t-1} + U_f x_t + b_f)\end{aligned}$$

$$\tilde{c}_t = \phi(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot \phi(c_t)$$

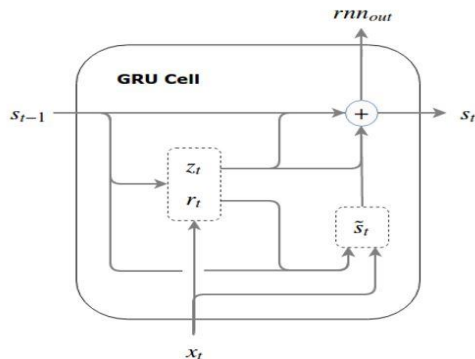
$$rnn_{out} = h_t$$

Key point of LSTM

Varient LSTM

► GRU

- input gate, forget gate, and output gate are replaced by update gate and reset gate



$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

$$\tilde{s}_t = \phi(W(r_t \odot s_{t-1}) + Ux_t + b)$$

$$s_t = z_t \odot s_{t-1} + (1 - z_t) \odot \tilde{s}_t$$

Variant LSTM

- ▶ There are more.....
- ▶ What kind of LSTM is the best?



Example



GRU/LSTM cell

This is a cat

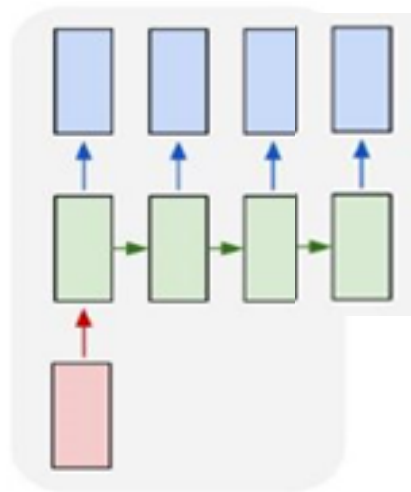


Image Caption

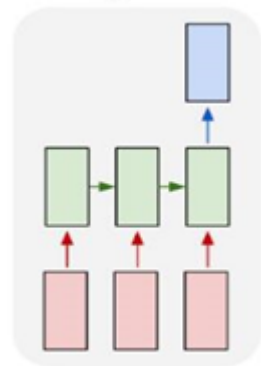
One to many



Example



GRU/LSTM cell



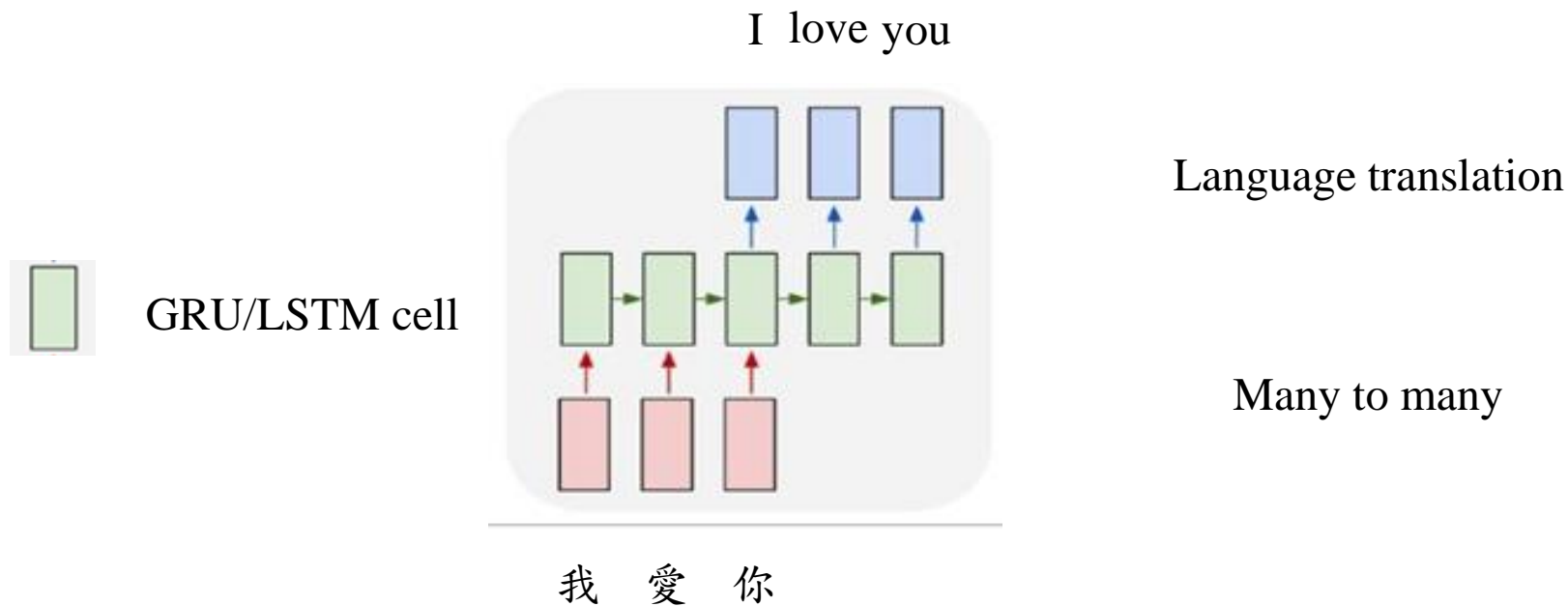
Good/Bad ?

Movie Review
(positive/negative)

Many to one

It is good

Example



Summary

- ▶ **Build your RNN model**
 - ▶ Observe your application which kind of RNN is suitable
 - ▶ Many to one, many to many,
- ▶ **Use LSTM/GRU cell now**
 - ▶ Perform better than traditional RNN