


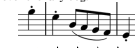
# Recurrent Neural Networks

## Why sequence models?

1

### Examples of sequence data


Speech recognition →  → "The quick brown fox jumped over the lazy dog."

Music generation →  → "There is nothing to like in this movie." ★☆☆☆☆

Sentiment classification → AGCCCCCTGTGAGGAAGTAG → AGCCCCCTGTGAGGAAGTAG


DNA sequence analysis → AGCCCCCTGTGAGGAAGTAG → AGCCCCCTGTGAGGAAGTAG

Machine translation → Voulez-vous chanter avec moi? → Do you want to sing with me?

Video activity recognition →  → Running

Name entity recognition → Yesterday, Harry Potter met Hermione Granger. → Yesterday, Harry Potter met Hermione Granger. Andrew Ng

2



# Recurrent Neural Networks

## Notation

3

### Motivating example

NLP

x: Harry Potter and Hermione Granger invented a new spell.

→  $x^{(1)}$   $x^{(2)}$   $x^{(3)}$  ...  $x^{(9)}$

$T_x = 9$

y: 1 1 0 1 0 0 0 0 0

$y^{(1)}$   $y^{(2)}$   $y^{(3)}$  ...  $y^{(9)}$

$T_y = 9$

$x^{(i)} \in \{1, \dots, 10,000\}$

$T_x = 9$  15

$y^{(i)} \in \{0, 1\}$

Andrew Ng

4

### Representing words

x: Harry Potter and Hermione Granger invented a new spell.

$x^{(1)}$   $x^{(2)}$   $x^{(3)}$  ...  $x^{(9)}$

Vocabulary:

a	1
and	2
Harry	3
potter	4
and	5
hermione	6
granger	7
invented	8
a	9
new	10
spell	11

One-hot

$x^{(i)} \in \{1, \dots, 10,000\}$

$T_x = 9$

$y^{(i)} \in \{0, 1\}$

Andrew Ng

5

### Representing words


x: Harry Potter and Hermione Granger invented a new spell.

$x^{(1)}$   $x^{(2)}$   $x^{(3)}$  ...  $x^{(9)}$

And = 367  
Invented = 4700  
A = 1  
New = 5976  
Spell = 8376  
Harry = 4075  
Potter = 6830  
Hermione = 4200  
Granger = 4000

Andrew Ng

6



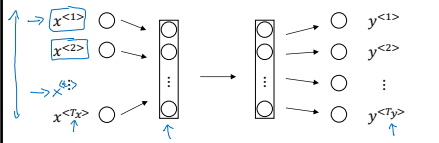
# Recurrent Neural Networks

---

## Recurrent Neural Network Model

7

### Why not a standard network?



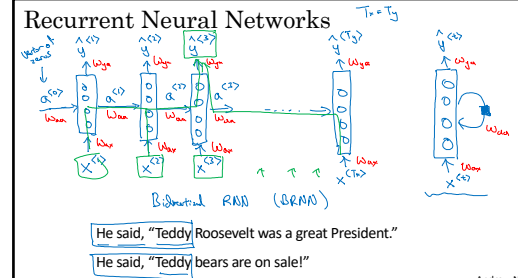
Problems:

- Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.

Andrew Ng

8

### Recurrent Neural Networks



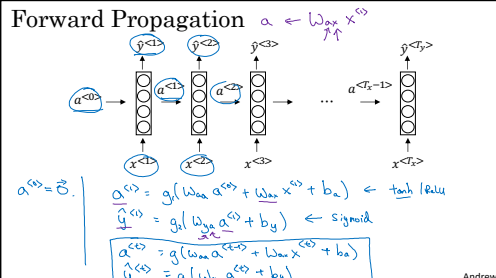
Bi-directional RNN (BRNN)

He said, "Teddy" Roosevelt was a great President.  
He said, "Teddy" bears are on sale!

Andrew Ng

9

### Forward Propagation

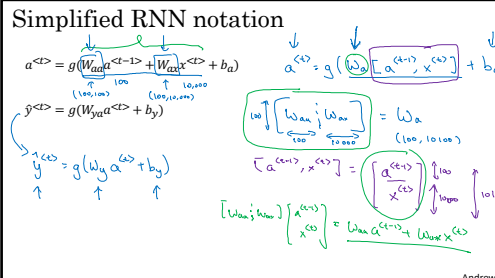


$a^{(0)} = 0$   
 $a^{(t)} = g(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a)$   
 $y^{(t)} = g(W_{ya} a^{(t)} + b_y)$   
 $\hat{y}^{(t)} = a(W_{yx} x^{(t)} + b_y)$

Andrew Ng

10


### Simplified RNN notation



$a^{(t)} = g(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a)$   
 $y^{(t)} = g(W_{ya} a^{(t)} + b_y)$   
 $\hat{y}^{(t)} = a(W_{yx} x^{(t)} + b_y)$

Andrew Ng

11




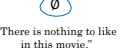

# Recurrent Neural Networks

---

## Different types of RNNs

12

### Examples of sequence data

Speech recognition		→ "The quick brown fox jumped over the lazy dog."
Music generation		→
Sentiment classification	"There is nothing to like in this movie."	→ ★★★★★
DNA sequence analysis	AGCCCTGTGAGGAAGTAG	→ AGCCCTGTGAGGAAGTAG
Machine translation	Voulez-vous chanter avec moi?	→ Do you want to sing with me?
Video activity recognition		→ Running
Name entity recognition	Yesterday, Harry Potter met Hermione Granger.	→ Yesterday, <u>Harry Potter</u> met <u>Hermione Granger</u> .

Andrew Ng

13

### Examples of RNN architectures

$T_x = T_y$

Sentiment classification  
 $x = \text{text}$   
 $y = 0/1 \quad 1-5$

Many-to-many

Many-to-one

One-to-one

Andrew Ng

14

### Examples of RNN architectures

Music generation

Machine translation

One-to-many

Many-to-many

Andrew Ng

15

### Summary of RNN types

One to one

One to many

Many to one

Many to many

Andrew Ng

16

### Recurrent Neural Networks

Language model and sequence generation

deeplearning.ai

17

### What is language modelling?

Speech recognition

The apple and pair salad.

→ The apple and pear salad.

$P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$

$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$

$P(\text{sentence}) = ?$

$P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$

Andrew Ng

18

### Language modelling with an RNN

Training set: large corpus of english text.

Tokense

Cats average 15 hours of sleep a day. <EOS>

$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(n)}$

$x^{(1)} = y^{(1-1)}$

The Egyptian Mau is a bread of cat. <EOS>

$10,000$  <UNK>

Andrew Ng

19

### RNN model

$P(y^{(1)} | x^{(1)}) \dots P(y^{(n)} | x^{(n)})$

$P(y^{(1)} | x^{(1)}) = P(y^{(1)} | a^{(0)})$

$P(y^{(2)} | x^{(2)}) = P(y^{(2)} | a^{(1)})$

$P(y^{(3)} | x^{(3)}) = P(y^{(3)} | a^{(2)})$

$P(y^{(n)} | x^{(n)}) = P(y^{(n)} | a^{(n-1)})$

$a^{(0)} = \vec{0}$

$x^{(1)} = y^{(0)}$

$x^{(2)} = y^{(1)}$

$x^{(3)} = y^{(2)}$

$x^{(n)} = y^{(n-1)}$

$P(\text{cats} | \text{cats})$

$P(\text{average} | \text{cats average})$

$P(\text{EOS} | \dots)$

$\rightarrow$  Cats average 15 hours of sleep a day. <EOS>

$L(y^{(1)}, y^{(2)}, y^{(3)}, \dots) = -\sum_i y_i^{(i)} \log \hat{y}_i^{(i)}$

$L = \sum_i L^{(i)}(y^{(i)}, \hat{y}^{(i)})$

$P(y^{(1)}, y^{(2)}, y^{(3)}, \dots) = P(y^{(1)}) P(y^{(2)} | y^{(1)}) P(y^{(3)} | y^{(1)}, y^{(2)})$

Andrew Ng

20

### Recurrent Neural Networks

Sampling novel sequences

deeplearning.ai

21

### Sampling a sequence from a trained RNN

$P(y^{(1)}, \dots, y^{(n)})$

Training:  $a^{(0)} \rightarrow a^{(1)} \rightarrow a^{(2)} \rightarrow \dots \rightarrow a^{(n)}$

$x^{(1)} = y^{(0)}$

$x^{(2)} = y^{(1)}$

$x^{(3)} = y^{(2)}$

$x^{(n)} = y^{(n-1)}$

$\rightarrow P(a^{(1)} | a^{(0)}) \dots P(a^{(n)} | a^{(n-1)})$  a.p. random choice

$P(\text{the})$

Andrew Ng

22

### Character-level language model

$\rightarrow$  Vocabulary = [a, aaron, ..., zulu, <UNK>]

$\rightarrow$  Vocabulary = [a, b, c, ..., z, ., , , , , 0, ..., 9, A, -, , 2]

$\rightarrow$   $y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad y^{(4)}$

$\rightarrow$   $\hat{y}^{(1)} \quad \hat{y}^{(2)} \quad \hat{y}^{(3)} \quad \hat{y}^{(4)}$

$\rightarrow$   $\hat{y}^{(1)} = \text{Cost}$

$\rightarrow$   $\hat{y}^{(2)} = \text{average}$

$\rightarrow$   $\hat{y}^{(3)} = \text{Mau}$

$\rightarrow$   $\hat{y}^{(4)} = \text{EOS}$

$\rightarrow$   $\hat{y}^{(5)} = \text{UNK}$

Andrew Ng

23

### Sequence generation

News

President Enrique Peña Nieto, announced sench's sulk former coming football langston paring.

"I was not at all surprised," said hich langston.

"Concussion epidemic", to be examined.

The gray football the told some and this has on the uefa icon, should money as.

Shakespeare

The mortal moon hath her eclipse in love.


And subject of this thou art another this fold.

When better be my love to me see sabl's.

For whose are ruse of mine eyes heaves.

Andrew Ng

24

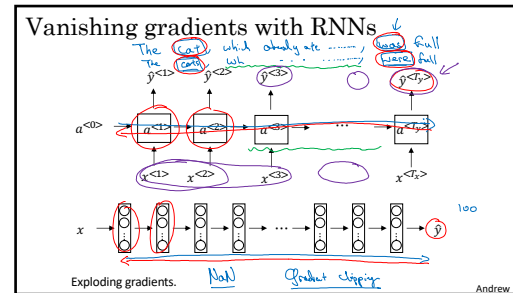


Recurrent Neural Networks


---

Vanishing gradients with RNNs

25



26

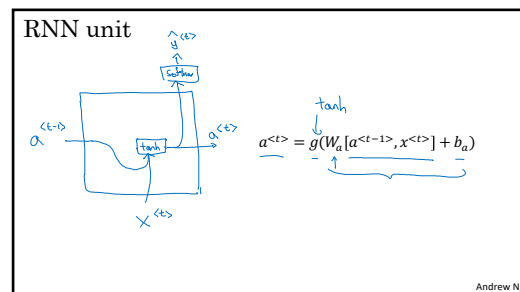


Recurrent Neural Networks

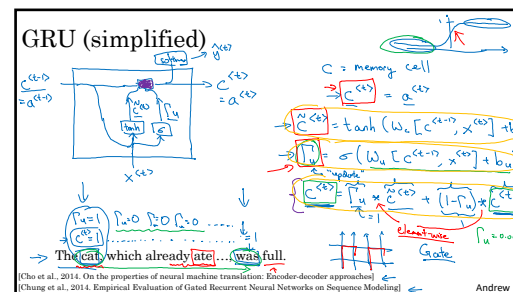
---

Gated Recurrent Unit (GRU)

27



28



29

Full GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r^{<t>} * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$


$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

The cat, which ate already, was full.

Andrew Ng

30



Recurrent Neural Networks

---

LSTM (long short term memory) unit

31

### GRU and LSTM

#### GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

#### LSTM

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

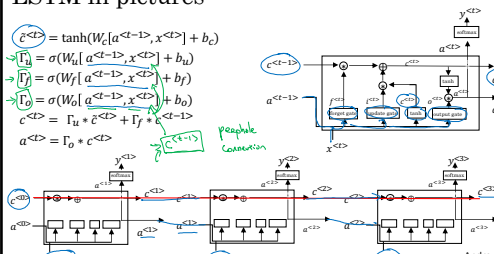
$$a^{<t>} = \Gamma_o * c^{<t>}$$

*(Handwritten notes: "forgetful" for  $\Gamma_f$ , "update" for  $\Gamma_u$ )*

Hochreiter & Schmidhuber 1997, Long short-term memory

32


### LSTM in pictures



*(Handwritten notes: "forgetful" for  $\Gamma_f$ , "update" for  $\Gamma_u$ )*

Andrew Ng

34



Recurrent Neural Networks

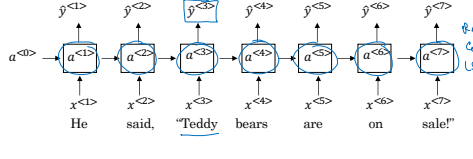
---

Bidirectional RNN

35

### Getting information from the future

He said, "Teddy bears are on sale!"  
He said, "Teddy Roosevelt was a great President!"

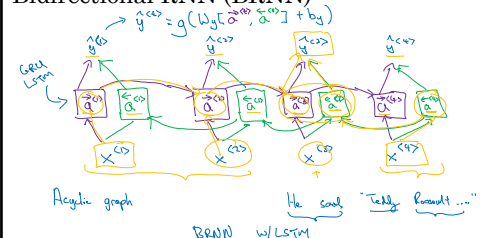


*(Handwritten notes: "BRNN GRU LSTM")*

Andrew Ng

36


### Bidirectional RNN (BRNN)



*(Handwritten notes: "peephole connection", "BRNN w/ LSTM")*

Andrew Ng

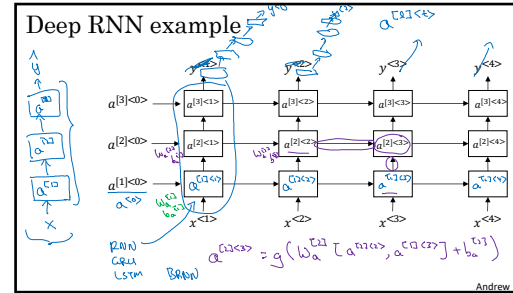
37




# Recurrent Neural Networks

## Deep RNNs

38



39



# Sequence to sequence models

## Basic models

40

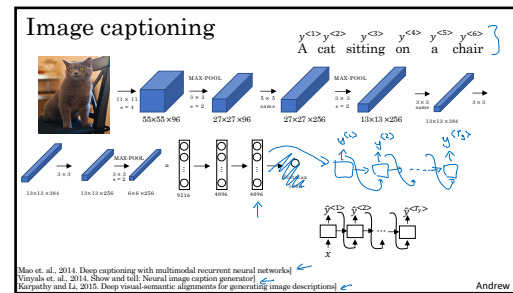
### Sequence to sequence model

$x^{<1>}$   $x^{<2>}$   $x^{<3>}$   $x^{<4>}$   $x^{<5>}$   
 Jane visite l'Afrique en septembre  
 → Jane is visiting Africa in September.  
 $y^{<1>}$   $y^{<2>}$   $y^{<3>}$   $y^{<4>}$   $y^{<5>}$


[Bastien et al., 2014. Sequence to sequence learning with neural networks](#)  
[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation](#)

Andrew Ng

41



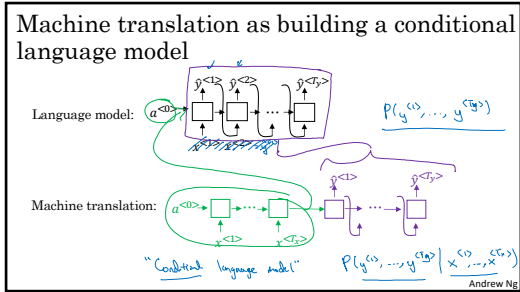
42



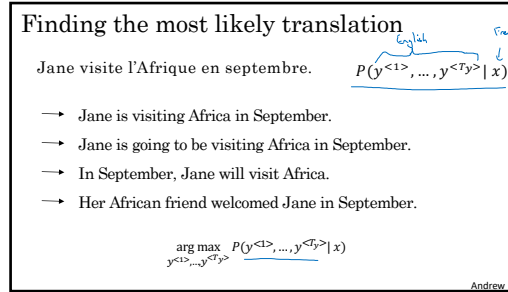
# Sequence to sequence models

## Picking the most likely sentence

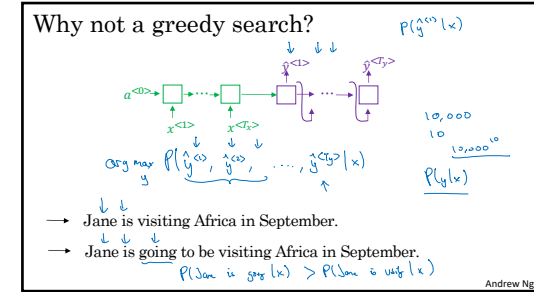
43



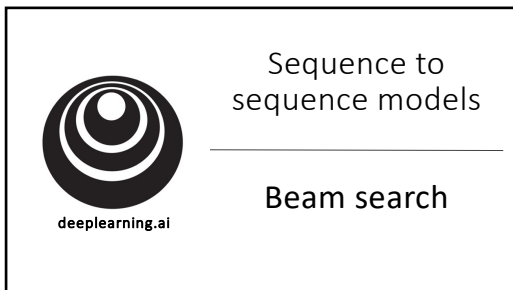
44



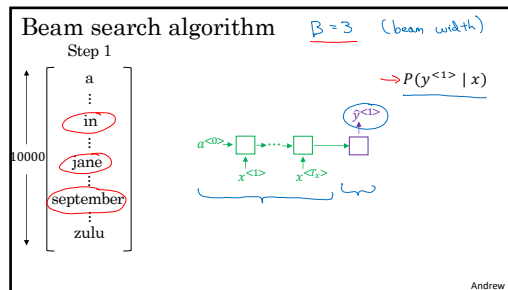
45



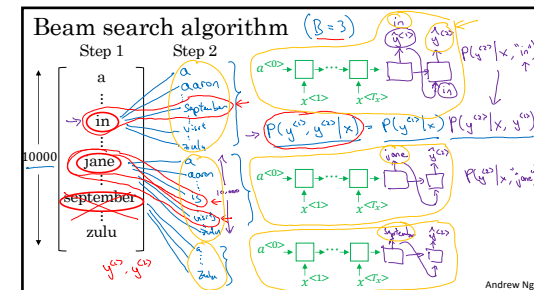
46



47

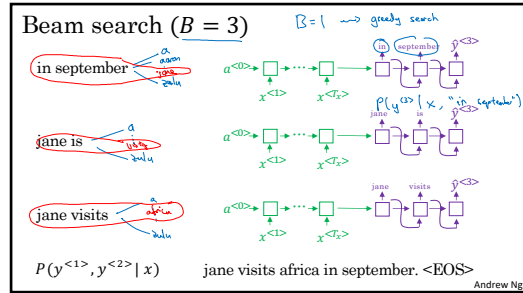


48

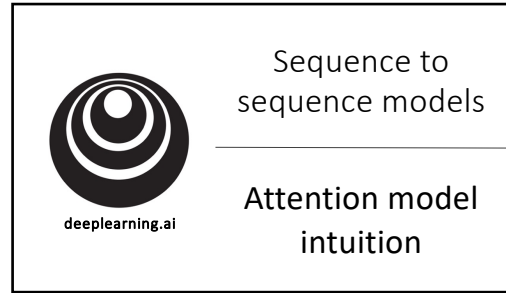


49

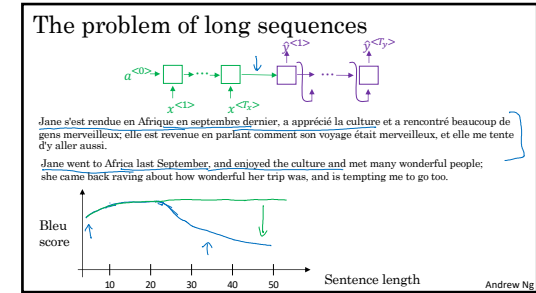




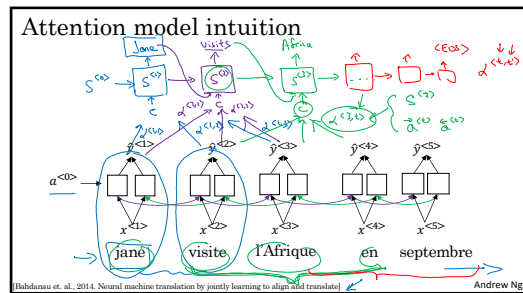
50



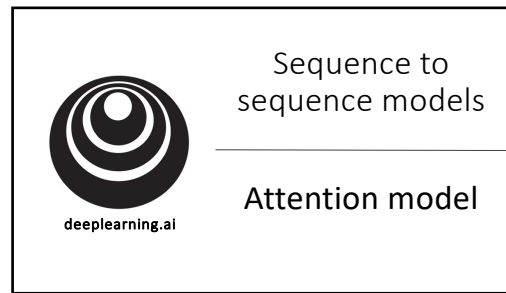
51



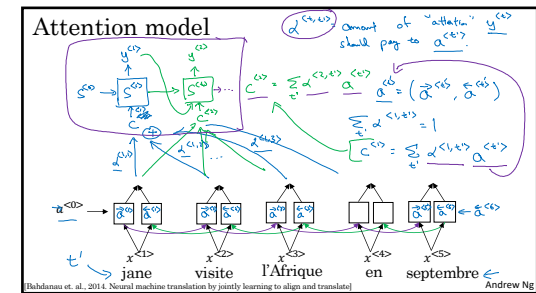
52



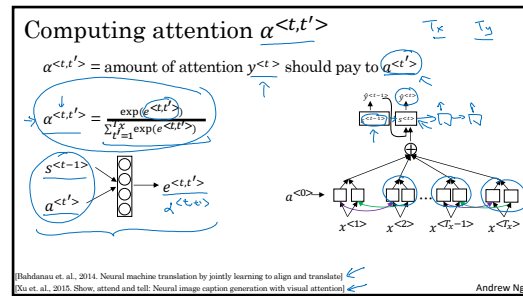
53



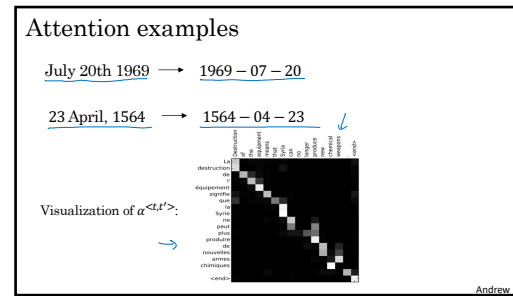
54



55



56



57