

SIT330-770: Natural Language Processing

Week 6 - Neural Networks for NLP

Dr. Mohamed Reda Bouadjenek

School of Information Technology, Faculty of
Sci Eng & Built Env


reda.bouadjenek@deakin.edu.au



1

Andrew Ng

Neural Networks and Deep Learning
(Optional)




2

SIT330-770: Natural Language Processing

Week 6.11 - Applying feedforward networks to NLP tasks

Dr. Mohamed Reda Bouadjenek


School of Information Technology,
Faculty of Sci Eng & Built Env



3

Use cases for feedforward networks

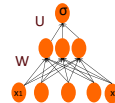
- Let's consider 2 (simplified) sample tasks:
 - Text classification
 - Language modeling
- State-of-the-art systems use more powerful neural architectures, but simple models are useful to consider!



4

Classification: Sentiment Analysis


- We could do exactly what we did with logistic regression
- Input layer are binary features as before
- Output layer is 0 or 1



5

Sentiment Features

Var	Definition
x_1	count(positive lexicon) \in doc
x_2	count(negative lexicon) \in doc
x_3	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_4	count(1st and 2nd pronouns \in doc)
x_5	$\begin{cases} 1 & \text{if "I"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
x_6	$\log(\text{word count of doc})$



6

Feedforward nets for simple classification

- Just adding a hidden layer to logistic regression
 - allows the network to use non-linear interactions between features
 - which may (or may not) improve performance.

7

Even better: representation learning

- The real power of deep learning comes from the ability to **learn** features from the data
- Instead of using hand-built human-engineered features for classification
- Use learned representations like embeddings!

8

Neural Net Classification with embeddings as input features!

9

Issue: texts come in different sizes

- This assumes a fixed size length (3)!
 - Kind of unrealistic.
- Some simple solutions (more sophisticated solutions later)
 - If shorter then pad with zero embeddings
 - Truncate if you get longer reviews at test time
- 1. Make the input the length of the longest review
 - Take the mean of all the word embeddings
 - Take the element-wise max of all the word embeddings
 - For each dimension, pick the max value from all words
- 2. Create a single "sentence embedding" (the same dimensionality as a word) to represent all the words

10

Reminder: Multiclass Outputs

- What if you have more than two output classes?
 - Add more output units (one for each class)
 - And use a "softmax layer"

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^D e^{z_j}} \quad 1 \leq i \leq D$$

11

Neural Language Models (LMs)

- Language Modeling:** Calculating the probability of the next word in a sequence given some history.
 - We've seen N-gram based LMs
 - But neural network LMs far outperform n-gram language models
- State-of-the-art neural LMs are based on more powerful neural network technology like Transformers
- But **simple feedforward LMs** can do almost as well!

12

Simple feedforward Neural Language Models

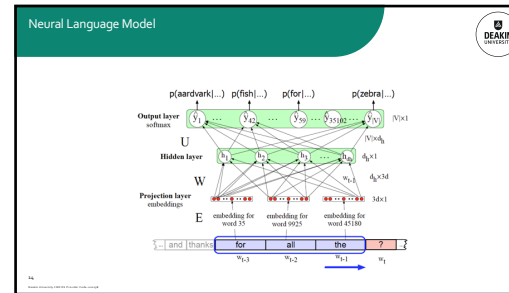
Task: predict next word w_t
given prior words $w_{t-1}, w_{t-2}, w_{t-3}, \dots$

Problem: Now we're dealing with sequences of arbitrary length.

Solution: Sliding windows (of fixed length)

$$P(w_t | w_{t-1}^{t-1}) \approx P(w_t | w_{t-N+1}^{t-1})$$

13



14

Why Neural LMs work better than N-gram LMs

- **Training data:**
 - We've seen: *I have to make sure that the cat gets fed.*
 - Never seen: *dog gets fed*
- **Test data:**
 - *I forgot to make sure that the dog gets ____*
- N-gram LM can't predict "fed"!
- Neural LM can use similarity of "cat" and "dog" embeddings to generalize and predict "fed" after dog

15

SIT330-770: Natural Language Processing

Week 6 - Neural Networks and Neural LMs

Dr. Mohamed Reda Bouadjene

School of Information Technology, Faculty of Sci Eng & Built Env


reda.bouadjene@deakin.edu.au

16

Andrew Ng

Neural Networks and Deep Learning

17




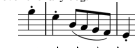
Recurrent Neural Networks

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 → AGCCCTGTGAGGAAGTAG → AGCCCTGTGAGGAAGTAG


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, 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.

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

y: 1 1 0 1 0 0 0 0 0

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

$T_x = 9$

$T_y = 9$

$x^{(i)} < t >$

$T_x^{(i)} = 9$

$T_y^{(i)} = 15$

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
invented	4
a	5
new	6
spell	7
with	8
Harry	9
Potter	10
and	11
Hermione	12
Granger	13
invented	14
a	15
new	16
spell	17

One-hot

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

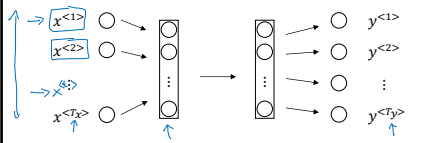


Recurrent Neural Networks

Recurrent Neural Network Model

7

Why not a standard network?



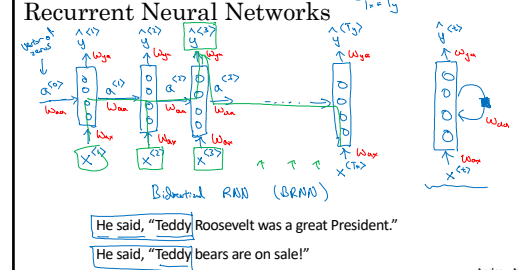
Problems:

- 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

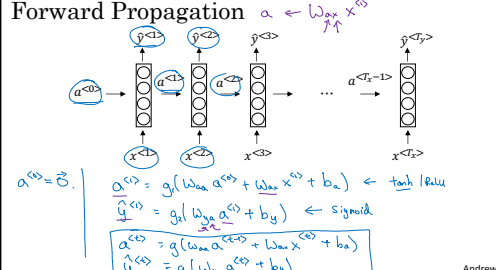


He said, "Teddy" Roosevelt was a great President.
He said, "Teddy" bears are on sale!

Andrew Ng

9

Forward Propagation

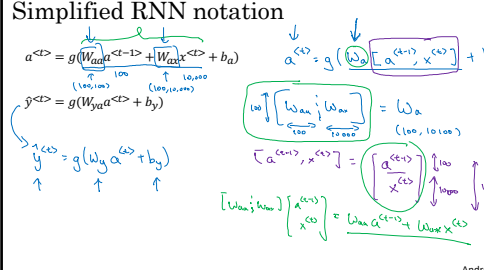


$a^{(0)} = 0$
 $a^{(t)} = g(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a)$
 $y^{(t)} = g(W_{ya} a^{(t)} + b_y)$
 $\hat{y}^{(t)} = a(W_{yx} x^{(t)} + b_y)$

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: x (audio waveform) \rightarrow "The quick brown fox jumped over the lazy dog." y (text)

Music generation: x (empty box) \rightarrow (musical notes)

Sentiment classification: "There is nothing to like in this movie." \rightarrow ★★★★★

DNA sequence analysis: AGCCCTGTGAGGAAGTAG \rightarrow AGCCCTGTGAGGAAGTAG

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

Video activity recognition: (video frames) \rightarrow Running

Name entity recognition: Yesterday, Harry Potter met Hermione Granger. \rightarrow Yesterday, Harry Potter met Hermione Granger.

Andrew Ng

13

Examples of RNN architectures

$T_x = T_y$

Many-to-many: $x^{(1)} \rightarrow y^{(1)}, x^{(2)} \rightarrow y^{(2)}, \dots, x^{(T)} \rightarrow y^{(T)}$

Sentiment classification: $x = \text{text}$, $y = 0/1$ (1-5)

Many-to-one: $x^{(1)}, x^{(2)}, \dots, x^{(T)} \rightarrow y$

One-to-one: $x \rightarrow y$

Andrew Ng

14

Examples of RNN architectures

Music generation: $x \rightarrow y^{(1)}, y^{(2)}, \dots, y^{(T)}$

Machine translation: $x^{(1)}, x^{(2)}, \dots, x^{(T)} \rightarrow y^{(1)}, y^{(2)}, \dots, y^{(T)}$

One-to-many: $x \rightarrow y^{(1)}, y^{(2)}, \dots, y^{(T)}$

Many-to-many: $x^{(1)}, x^{(2)}, \dots, x^{(T)} \rightarrow y^{(1)}, y^{(2)}, \dots, y^{(T)}$

Andrew Ng

15

Summary of RNN types

One to one: $x^{(1)} \rightarrow y^{(1)}$

One to many: $x^{(1)} \rightarrow y^{(1)}, y^{(2)}, \dots, y^{(T)}$

Many to one: $x^{(1)}, x^{(2)}, \dots, x^{(T)} \rightarrow y$

Many to many: $x^{(1)}, x^{(2)}, \dots, x^{(T)} \rightarrow y^{(1)}, y^{(2)}, \dots, y^{(T)}$

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. \rightarrow 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

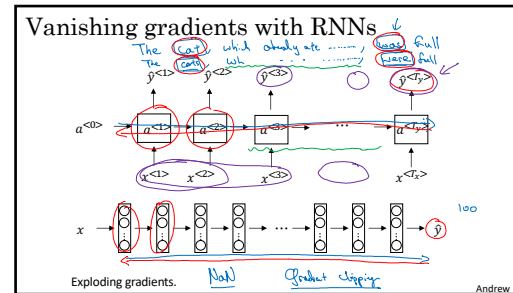
18




Recurrent Neural Networks

Vanishing gradients with RNNs

25



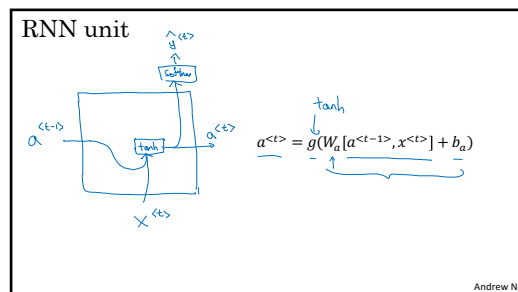
26



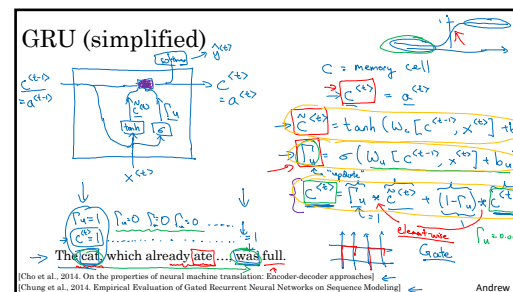
Recurrent Neural Networks

Gated Recurrent Unit (GRU)

27



28



29


Full GRU

$$\begin{aligned} \tilde{c}^{<t>} &= \tanh(W_c[\tilde{a}^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_r &= \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r) \\ c^{<t>} &= \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>} \end{aligned}$$

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

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

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

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

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

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

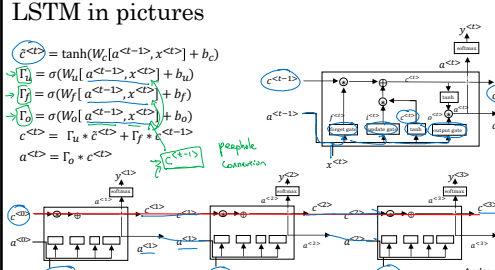
$$a^{<t>} = \Gamma_o * c^{<t>}$$

(Handwritten notes: "forget", "update", "output gate")

Hochreiter & Schmidhuber 1997, Long short-term memory

32

LSTM in pictures



Andrew Ng

34



Recurrent Neural Networks

Bidirectional RNN

35

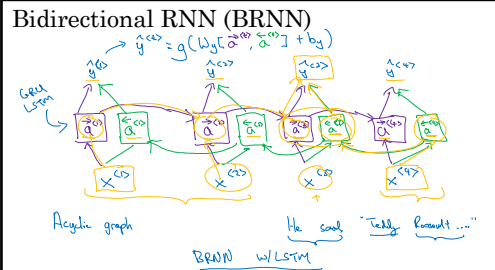
Getting information from the future

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

Andrew Ng


36

Bidirectional RNN (BRNN)



Andrew Ng

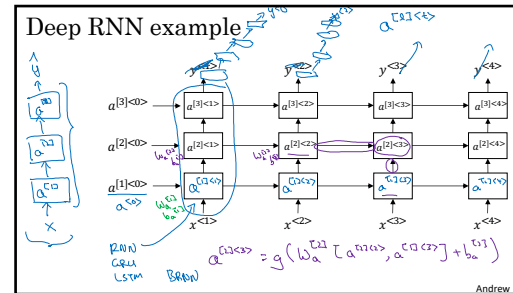
37




Recurrent Neural
Networks

Deep RNNs

38



39



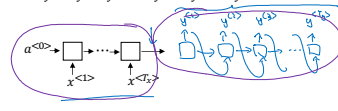
Sequence to
sequence models

Basic models

40

Sequence to sequence model

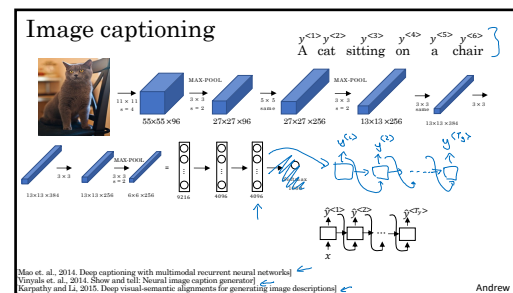
$x^{<1>} x^{<2>} x^{<3>} x^{<4>} x^{<5>}$
 Jane visite l'Afrique en septembre
 \rightarrow Jane is visiting Africa in September.
 $y^{<1>} y^{<2>} y^{<3>} y^{<4>} y^{<5>}$




[Sutskever 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



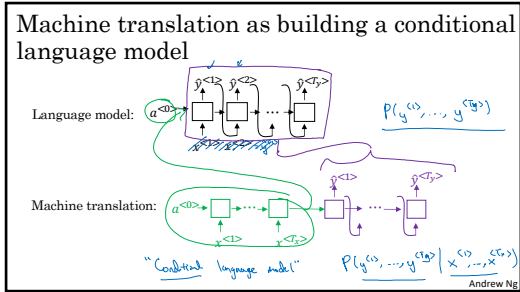
42



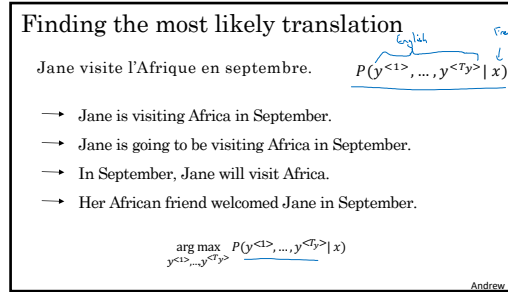
Sequence to
sequence models

Picking the most
likely sentence

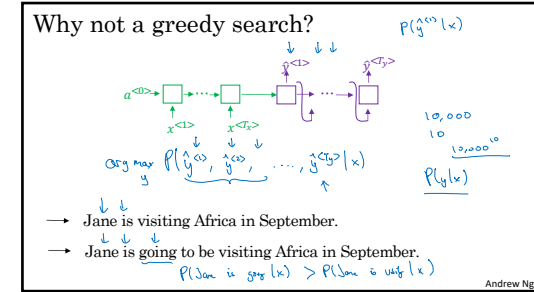
43



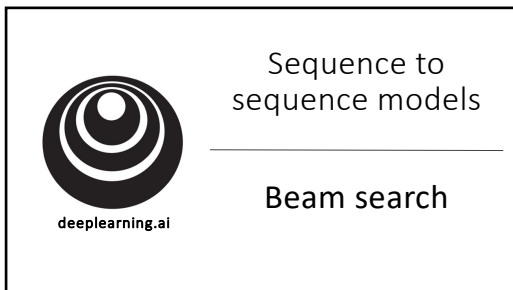
44



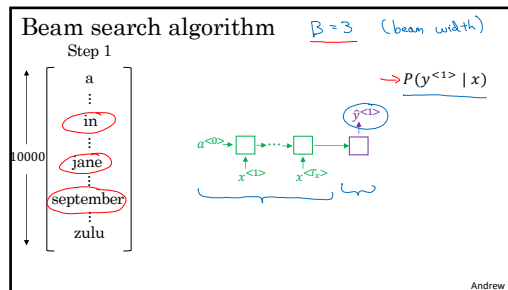
45



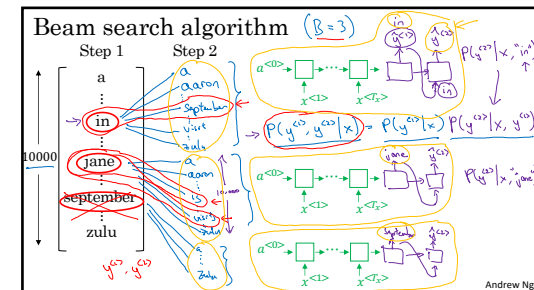
46



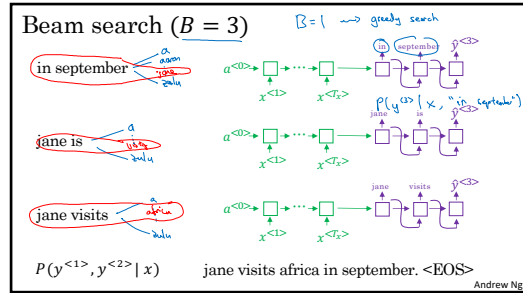
47



48



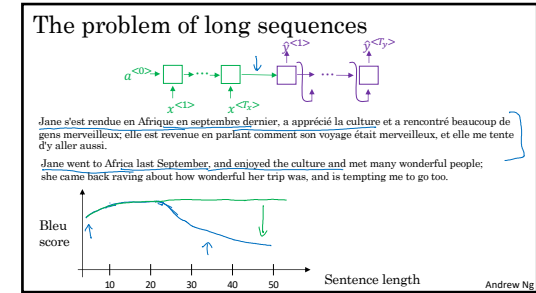
49



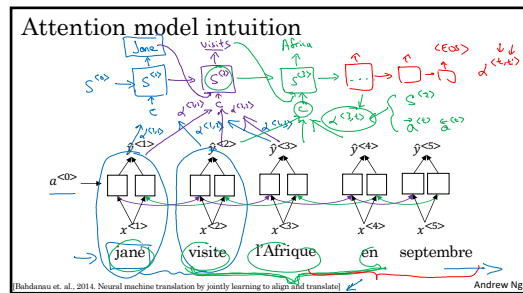
50



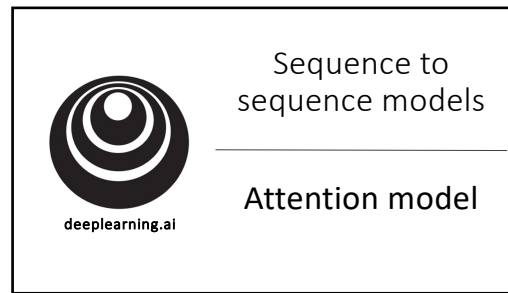
51



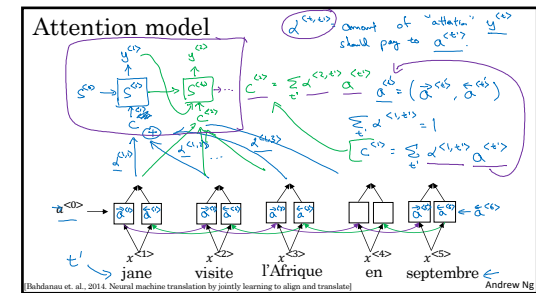
52



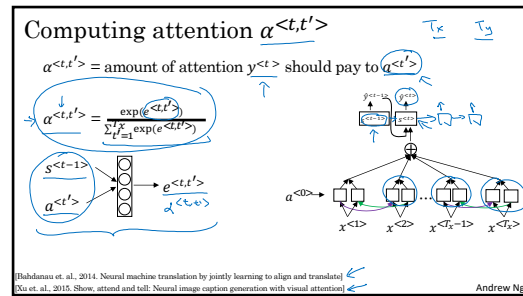
53



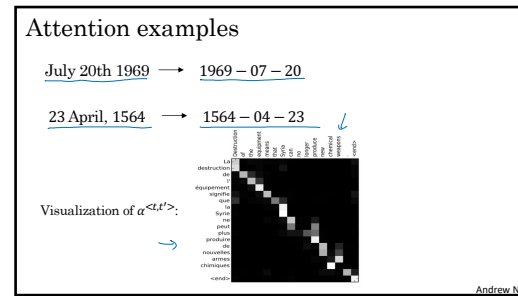
54



55



56



57