# Sequence to sequence model

$x^{<1>}$  $x^{<2>}$     $x^{<3>}$     $x^{<4>}$   $x^{<5>}$
Jane  visite  l'Afrique  en  septembre

$\longrightarrow$  Jane  is  visiting  Africa  in  September.
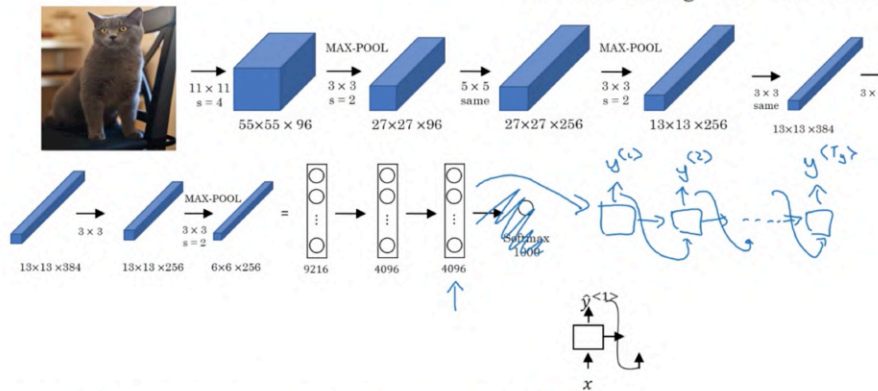  $y^{<1>}$  $y^{<2>}$  $y^{<3>}$     $y^{<4>}$   $y^{<5>}$     $y^{<6>}$

[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

# Image captioning

$y^{<1>}$ $y^{<2>}$     $y^{<3>}$     $y^{<4>}$   $y^{<5>}$   $y^{<6>}$
A   cat   sitting   on    a   chair

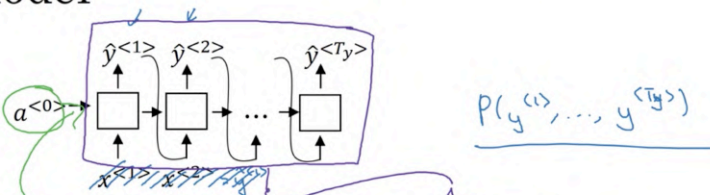[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]
[Vinyals et. al., 2014. Show and tell: Neural image caption generator]
[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]
Andrew Ng
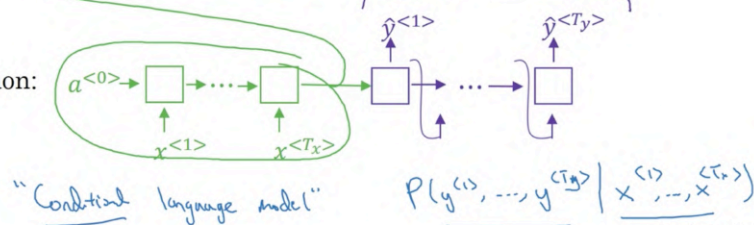
# Machine translation as building a conditional language model

Language model:

Machine translation:

$$P(y^{<1>}, \dots, y^{<T_y>})$$

"Conditional language model"

$$P(y^{<1>}, \dots, y^{<T_y>} \mid x^{<1>}, \dots, x^{<T_x>})$$



Andrew Ng

# Finding the most likely translation

Jane visite l'Afrique en septembre.

$P(y^{<1>}, \dots, y^{<T_y>} | x)$ · English · French

→ Jane is visiting Africa in September.

→ Jane is going to be visiting Africa in September.

→ In September, Jane will visit Africa.

→ Her African friend welcomed Jane in September.

$$\arg\max_{y^{<1>}, \dots, y^{<T_y>}} P(y^{<1>}, \dots, y^{<T_y>} | x)$$

Andrew Ng

---

# Why not a greedy search?

$P(\hat{y}^{<1>} | x)$

Considers Only 1 word at a time



$\arg\max_y P(\hat{y}^{<1>}, \hat{y}^{<2>}, \dots, \hat{y}^{<T_y>} | x)$

10,000
10
$10,000^{10}$

$P(y|x)$

→ Jane is visiting Africa in September.

→ Jane is going to be visiting Africa in September.

$P(\text{Jane is goiy} | x) > P(\text{Jane is visity} | x)$

Andrew Ng

---

# Beam search algorithm   $B = 3$  (beam width)

Step 1

$\rightarrow P(y^{<1>} | x)$



$$10000 \begin{bmatrix} a \\ \vdots \\ \text{in} \\ \vdots \\ \text{jane} \\ \vdots \\ \text{september} \\ \vdots \\ \text{zulu} \end{bmatrix}$$

instantiate B copies of the network

Beam width / Considers B words!

Andrew Ng

# Beam search algorithm $(B = 3)$

Step 1     Step 2



$P(y^{<1>}, y^{<2>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})$

$P(y^{<2>} | x, "in")$

$P(y^{<2>} | x, "jane")$

---

# Beam search $(B = 3)$

$B = 1 \implies$ greedy search



in september

jane is

jane visits

$P(y^{<1>}, y^{<2>} | x)$

jane visits africa in september. <EOS>

$P(y^{<3>} | x, "in \, \hat{y}^{<3>} \, (september")$

---

# Length normalization

$P(y^{<1>} \ldots y^{<T_y>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>}) \cdots$

$P(y^{<T_y>} | x, y^{<1>} \ldots, y^{<T_y - 1>})$

$$\arg\max_{y} \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \ldots, y^{<t-1>})$$

$\log$

$\log P(y | x) \leftarrow$

$P(y | x) \leftarrow$

$$\arg\max_{y} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \ldots, y^{<t-1>}) \leftarrow$$

$T_y = 1, 2, 3, \ldots, 30.$

$$\frac{1}{T_y^{\alpha}} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \ldots, y^{<t-1>})$$

$\alpha = 0.7$    $\alpha = 1$

$\alpha = 0$

# Beam search discussion

Beam width B?

$1 \to 3 \to 10, \quad 100, \quad 1000 \to 3000$

large B: better result, slower
Small B: worse result, faster

Unlike exact search algorithms like BFS (Breadth First Search) or
DFS (Depth First Search), Beam Search runs faster but is not
guaranteed to find exact maximum for $\arg\max\limits_{y} P(y|x)$.

# Length normalization

$P(y^{<1>} \dots y^{<T_y>} | x) = P(y^{<1>} | x) \, P(y^{<2>} | x, y^{<1>}) \dots$
$P(y^{<T_y>} | x, y^{<1>} \dots, y^{<T_y - 1>})$

$$\arg\max_{y} \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

Numbers less than 1

log

$\log P(y|x) \leftarrow$

$P(y|x) \leftarrow$

more numerically stable with logs

$$\arg\max_{y} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>}) \leftarrow$$

$T_y = 1, 2, 3, \dots, 30.$

Average of log probabilities

$$\frac{1}{T_y^{\alpha}} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

$\alpha = 0.7$

$\alpha = 1$

$\alpha = 0$

# Beam search discussion

# Example

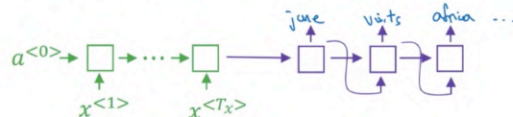Jane visite l'Afrique en septembre.

→ RNN
→ Beam Search    B↑

Human: Jane visits Africa in September. $(y^*)$

Algorithm: Jane visited Africa last September. $(\hat{y})$ ←

RNN computes $P(y^*|x) \gtrless P(\hat{y}|x)$

---

# Error analysis on beam search

Human: Jane visits Africa in September. $(y^*)$

Algorithm: Jane visited Africa last September. $(\hat{y})$

$P(y^*|x)$

$P(\hat{y}|x)$

Case 1: $P(y^*|x) > P(\hat{y}|x)$ ←    $\underset{y}{\arg\max} P(y|x)$

Beam search chose $\hat{y}$. But $y^*$ attains higher $\boxed{P(y|x)}$.

Conclusion: Beam search is at fault.

Case 2: $P(y^*|x) \leq P(\hat{y}|x)$ ←

$y^*$ is a better translation than $\hat{y}$. But RNN predicted $\boxed{P(y^*|x)} < P(\hat{y}|x)$.

Conclusion: RNN model is at fault.

---

# Error analysis process

| Human | Algorithm | $P(y^*|x)$ | $P(\hat{y}|x)$ | At fault? |
|---|---|---|---|---|
| Jane visits Africa in September. | Jane visited Africa last September. | $2 \times 10^{-10}$ | $1 \times 10^{-10}$ | B |
| | | | | R |
| | | | | B |
| | | | | R |
| | | | | R |
| | | | | ⋮ |

Figures out what faction of errors are "due to" beam search vs. RNN model

# Evaluating machine translation

French: Le chat est sur le tapis.

→ Reference 1: The cat is on the mat.

→ Reference 2: There is a cat on the mat.

→ MT output: the the the the the the the.

Precision: $\dfrac{7}{7}$  Modified precision: $\dfrac{2}{7}$

*Bleu — bilingual evaluation understudy*

*Count$_{clip}$("the")* ← 2

*Count("the")* ← 7

*2 appearance*

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]  Andrew Ng

---

# Bleu score on bigrams

Example: Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

MT output: The cat the cat on the mat.

| | Count | Count$_{clip}$ |
|---|---|---|
| the cat | 2 | 1 |
| cat the | 1 | 0 |
| cat on | 1 | 1 |
| on the | 1 | 1 |
| the mat | 1 | 1 |

$\dfrac{4}{6}$

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]  Andrew Ng

---

# Bleu score on unigrams

Example: Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

→ MT output: The cat the cat on the mat. $(\hat{y})$

$P_1 \cdot P_2 = 1.0$

$$P_1 = \frac{\sum\limits_{\text{unigrams} \in \hat{y}} \left( Count_{clip}(\text{unigram}) \right)}{\sum\limits_{\text{unigram} \in \hat{y}} Count(\text{unigram})}$$

*unigram*

$$P_n = \frac{\sum\limits_{n\text{-grams} \in \hat{y}} Count_{clip}(n\text{-gram})}{\sum\limits_{n\text{-grams} \in \hat{y}} Count(n\text{-gram})}$$

*n-gram*

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]  Andrew Ng

# Bleu details

$p_n$ = Bleu score on n-grams only        $P_1, P_2, P_3, P_4$

Combined Bleu score:        $BP \; exp\left(\frac{1}{4}\sum_{n=1}^{4} P_n\right)$

BP = brevity penalty

$$BP = \begin{cases} 1 & \text{if MT\_output\_length} > \text{reference\_output\_length} \\ \exp(1 - \text{MT\_output\_length}/\text{reference\_output\_length}) & \text{otherwise} \end{cases}$$

# Bleu details

$p_n$ = Bleu score on n-grams only        $P_1, P_2, P_3, P_4$

Combined Bleu score:        $BP \; exp\left(\frac{1}{4}\sum_{n=1}^{4} P_n\right)$

BP = brevity penalty

$$BP = \begin{cases} 1 & \text{if MT\_output\_length} > \text{reference\_output\_length} \\ \exp(1 - \text{MT\_output\_length}/\text{reference\_output\_length}) & \text{otherwise} \end{cases}$$

exp(1 − reference_output_length/MT_output_length)

# The problem of long sequences



$a^{<0>}$        $\hat{y}^{<1>}$        $\hat{y}^{<T_y>}$        $x^{<1>}$        $x^{<T_x>}$

Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de
gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente
d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people;
she came back raving about how wonderful her trip was, and is tempting me to go too.

Bleu
score        10   20   30   40   50        Sentence length

# Attention model intuition

Jane visits Africa <EOS>

$\alpha^{<t,t'>}$

$S^{<0>}$ $S^{<1>}$ $S^{<2>}$ $S^{<3>}$ ...

$S^{(2)}$

$\vec{a}^{(t)}$, $\overleftarrow{a}^{(t)}$

How much attention should we pay?

$c$ $\alpha^{<2,1>}$ $c$ $\alpha^{<2,1>}$

$\alpha^{<1,1>}$ $\alpha^{<1,2>}$ $\alpha^{<1,3>}$ $c$ $\alpha^{<3,t>}$

$a^{<0>} \rightarrow$

$x^{<1>}$  $x^{<2>}$  $x^{<3>}$  $x^{<4>}$  $x^{<5>}$

jane   visite   l'Afrique   en   septembre

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

---

# Attention model

Generate New context

$\alpha^{<t,t'>}$ amount of "attention" $y^{<t>}$ should pay to $a^{<t'>}$

$y^{<1>}$  $y^{<2>}$

$S^{<0>} \rightarrow S^{<1>}$  $S^{<2>} \rightarrow$ ...
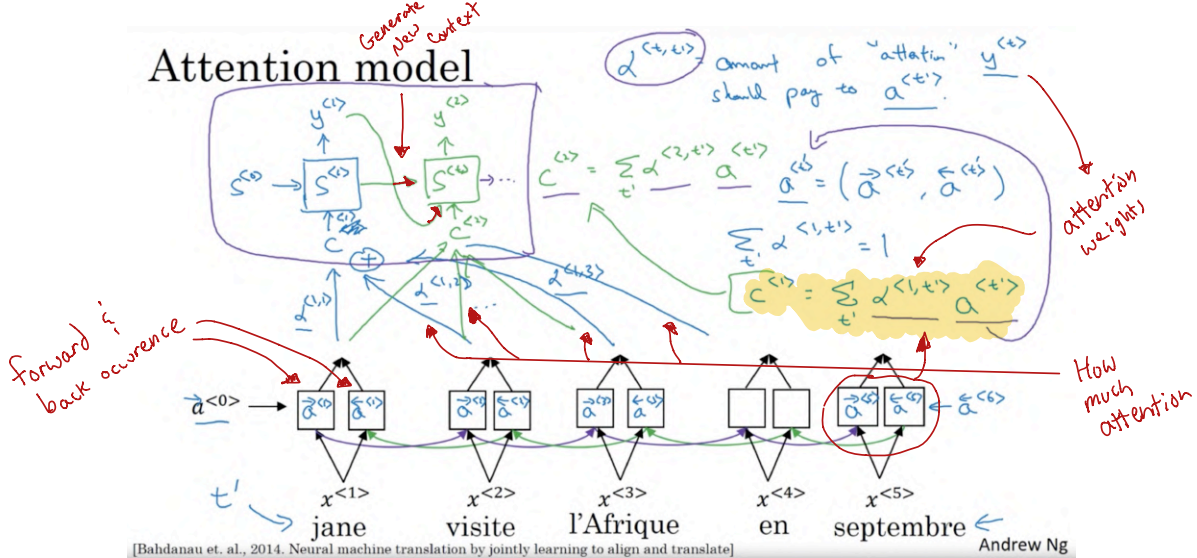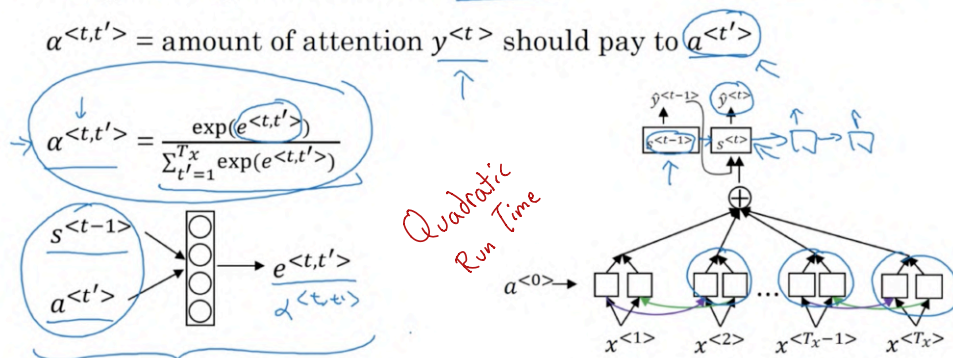
$C^{<2>} = \sum_{t'} \alpha^{<2,t'>} a^{<t'>}$

$a^{<t'>} = (\vec{a}^{<t'>}, \overleftarrow{a}^{<t'>})$

attention weights

$\sum_{t'} \alpha^{<1,t'>} = 1$

$C^{<1>}$ $C^{<2>}$

$c^{<1>} = \sum_{t'} \alpha^{<1,t'>} a^{<t'>}$

$\alpha^{<1,1>}$ $\alpha^{<1,2>}$ $\alpha^{<1,3>}$

forward & back ocurrence

How much attention

$\vec{a}^{<0>} \rightarrow$

$t'$

$x^{<1>}$  $x^{<2>}$  $x^{<3>}$  $x^{<4>}$  $x^{<5>}$

jane   visite   l'Afrique   en   septembre

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

---

# Computing attention $\alpha^{<t,t'>}$

$T_x$    $T_y$

$\alpha^{<t,t'>}$ = amount of attention $y^{<t>}$ should pay to $a^{<t'>}$

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})}$$

$\hat{y}^{<t-1>}$  $\hat{y}^{<t>}$

$s^{<t-1>}$  $s^{<t>}$

$s^{<t-1>}$

$a^{<t'>}$ $\rightarrow e^{<t,t'>}$

$\alpha^{<t,t'>}$

Quadratic Run Time

$a^{<0>} \rightarrow$

$x^{<1>}$  $x^{<2>}$  ...  $x^{<T_x-1>}$  $x^{<T_x>}$

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

[Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention]

Andrew Ng

# Attention examples

$$July\ 20th\ 1969 \longrightarrow 1969-07-20$$

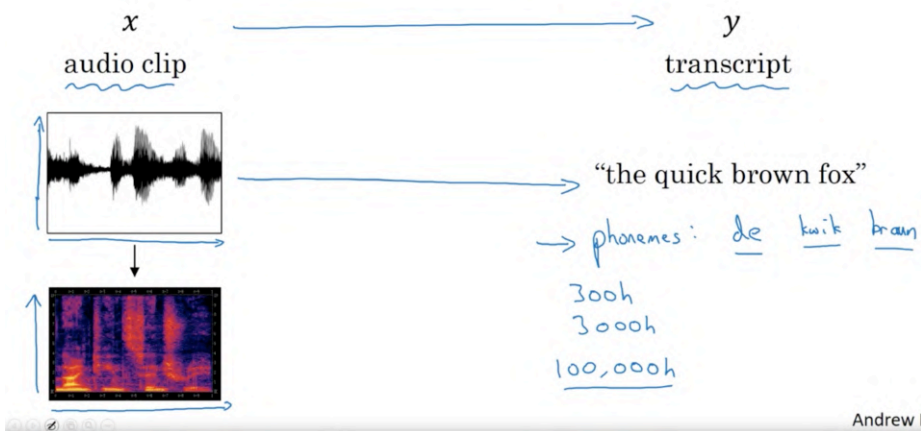$$23\ April,\ 1564 \longrightarrow 1564-04-23$$

Visualization of $\alpha^{<t,t'>}$:



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]
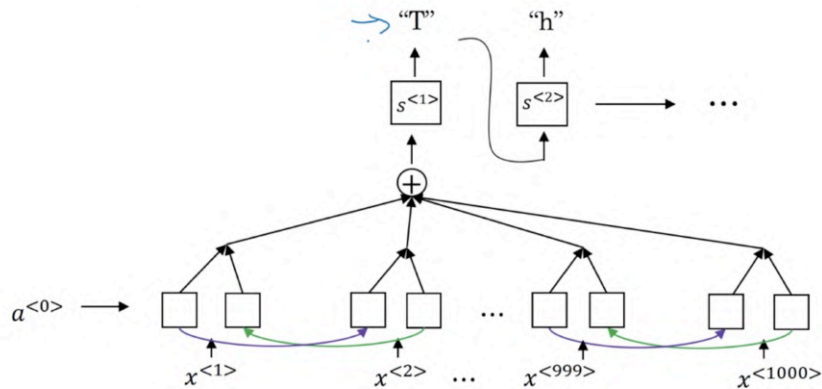
Andrew Ng

# Speech recognition problem

$x$ audio clip $\longrightarrow$ $y$ transcript



"the quick brown fox"

$\longrightarrow$ phonemes: de kwik bran

300h
3000h
100,000h

Andrew Ng

# Attention model for speech recognition



Andrew Ng

# CTC cost for speech recognition

(Connectionist temporal classification)



"the quick brown fox" — 19 characters

$a^{<0>} \rightarrow \square \rightarrow \square \rightarrow \cdots \rightarrow \square$

$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \hat{y}^{<1000>}$

$x^{<1>} \quad x^{<2>} \quad x^{<1000>}$

ttt_h_eee - - - - ⊔ - - - - qqq - - the q

"space"    "blank"

Basic rule: collapse repeated characters not separated by "blank"

[Graves et al., 2006. Connectionist Temporal Classification: Labeling unsegmented sequence data with recurrent neural networks]   Andrew Ng

# What is trigger word detection?



| Amazon Echo (Alexa) | Baidu DuerOS (xiaodunihao) | Apple Siri (Hey Siri) | Google Home (Okay Google) |

Andrew Ng

# Trigger word detection algorithm



$a^{<0>} \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow$

$x^{<1>} \quad x^{<2>} \quad x^{<3>}$

Andrew Ng