

89688: Statistical Machine Translation **Neural Machine Translation (contd.)**

Roee Aharoni Computer Science Department Bar Ilan University

Based in part on slides by Kevin Duh and Hermann Ney from the <u>DL4MT winter school</u>

May 2020

• A language model $p(w_1^N)$ measures how likely is the sentence:

 $w_1^N = x_1, x_2, \dots, x_N$

- A language model $p(w_1^N)$ measures how likely is the sentence:
- Usually modelled as a product of conditional probabilities:

$$w_1^N = x_1, x_2, ..., x_N$$

$$p(x_1, x_2, \dots, x_N) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_N)$$



- A language model $p(w_1^N)$ measures how likely is the sentence:
- Usually modelled as a product of conditional probabilities:
- The conventional approach assume a Markov chain of order n:

es $w_1^N = x_1, x_2, ..., x_N$

$$p(x_1, x_2, \dots, x_N) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_t)$$

 $p(x_1, x_2, \dots, x_N) = \prod_{t=1}^T p(x_t \mid x_{t-n}, \dots, x_{t-1})$



- A language model $p(w_1^N)$ measures how likely is the sentence:
- Usually modelled as a product of conditional probabilities:
- The conventional approach assume a Markov chain of order r
- and count: $p(x_t \mid x_{t-n}, \ldots, x_{t-n})$

es $w_1^N = x_1, x_2, ..., x_N$

$$p(x_1, x_2, \dots, x_N) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_t)$$

$$p(x_1, x_2, \dots, x_N) = \prod_{t=1}^T p(x_t \mid x_{t-n}, \dots, x_N)$$
n:

$$_{-1}) = \frac{count(x_{t-n}, \dots, x_{t-1}, x_t)}{count(x_{t-n}, \dots, x_{t-1})}$$



word	unigram	bigram	trigram	4-8
i	6.684	3.197	3.197	3
would	8.342	2.884	2.791	2
like	9.129	2.026	1.031	1
to	5.081	0.402	0.144	0
commend	15.487	12.335	8.794	8
\mathbf{the}	3.885	1.402	1.084	0
rapporteur	10.840	7.319	2.763	2
on	6.765	4.140	4.150	1
his	10.678	7.316	2.367	1
work	9.993	4.816	3.498	2
•	4.896	3.020	1.785	1
	4.828	0.005	0.000	0
average	8.051	4.072	2.634	2
perplexity	265.136	16.817	6.206	4



• Let's compute:

p(i, would, like, to, ..., < /s >)

word	unigram	bigram	trigram	4-8
i	6.684	3.197	3.197	3
would	8.342	2.884	2.791	2
like	9.129	2.026	1.031	1
to	5.081	0.402	0.144	0
commend	15.487	12.335	8.794	8
\mathbf{the}	3.885	1.402	1.084	0
rapporteur	10.840	7.319	2.763	2
on	6.765	4.140	4.150	1
his	10.678	7.316	2.367	1
work	9.993	4.816	3.498	2
	4.896	3.020	1.785	1
	4.828	0.005	0.000	0
average	8.051	4.072	2.634	2
perplexity	265.136	16.817	6.206	4



• Let's compute:

p(i, would, like, to, ..., < /s >)

• Unigram LM:

p(i)p(would)p(like)...p(</s>)

word	unigram	bigram	trigram	4-g
i	6.684	3.197	3.197	3
would	8.342	2.884	2.791	2
like	9.129	2.026	1.031	1
to	5.081	0.402	0.144	0
commend	15.487	12.335	8.794	8
\mathbf{the}	3.885	1.402	1.084	0
rapporteur	10.840	7.319	2.763	2
on	6.765	4.140	4.150	1
his	10.678	7.316	2.367	1
work	9.993	4.816	3.498	2
	4.896	3.020	1.785	1
	4.828	0.005	0.000	0
average	8.051	4.072	2.634	2
perplexity	265.136	16.817	6.206	4



• Let's compute:

p(i, would, like, to, ..., < /s >)

• Unigram LM:

p(i)p(would)p(like)...p(</s>)

• Bi-gram LM:

 $p(i)p(would \mid i)p(like \mid would)...p(</s>|.)$

word	unigram	bigram	trigram	4-g
i	6.684	3.197	3.197	3
would	8.342	2.884	2.791	2
like	9.129	2.026	1.031	1
\mathbf{to}	5.081	0.402	0.144	0
commend	15.487	12.335	8.794	8
\mathbf{the}	3.885	1.402	1.084	0
rapporteur	10.840	7.319	2.763	2
on	6.765	4.140	4.150	1
his	10.678	7.316	2.367	1
work	9.993	4.816	3.498	2
	4.896	3.020	1.785	1
	4.828	0.005	0.000	0
average	8.051	4.072	2.634	2
perplexity	265.136	16.817	6.206	4

Perplexity - The lower, the better



• Let's compute:

p(i, would, like, to, ..., < /s >)

• Unigram LM:

p(i)p(would)p(like)...p(</s>)

• Bi-gram LM:

 $p(i)p(would \mid i)p(like \mid would)...p(</s>|.)$

• Tri-gram LM:

 $p(i)p(would \mid i)p(like \mid i, would)...p(</s > | work, .)$

word	unigram	bigram	trigram	4-g
i	6.684	3.197	3.197	3
would	8.342	2.884	2.791	2
like	9.129	2.026	1.031	1
\mathbf{to}	5.081	0.402	0.144	0
commend	15.487	12.335	8.794	8
\mathbf{the}	3.885	1.402	1.084	0
rapporteur	10.840	7.319	2.763	2
on	6.765	4.140	4.150	1
his	10.678	7.316	2.367	1
work	9.993	4.816	3.498	2
	4.896	3.020	1.785	1
	4.828	0.005	0.000	0
average	8.051	4.072	2.634	2
perplexity	265.136	16.817	6.206	4



 Perplexity is the inverse probability of the unseen test set, normalised by the number of words:

$PP(W) = P(w_1 w_2 ... w_N)^{-\frac{1}{N}}$

- Perplexity is the inverse probability of the unseen test set, normalised by the number of words:
- Can be seen as an exponentiation of the entropy:

 $PP(W) = P(w_1 w_2 ... w_N)^{-\frac{1}{N}}$

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

- Perplexity is the inverse probability of the unseen test set, normalised by the number of words:
- Can be seen as an exponentiation of the entropy:



- Perplexity is the inverse probability of the unseen test set, normalised by the number of words:
- Can be seen as an exponentiation of the entropy:
- Lower perplexity means lower entropy which means less uncertainty



- Perplexity is the inverse probability of the unseen test set, normalised by the number of words:
- Can be seen as an exponentiation of the entropy:
- Lower perplexity means lower entropy which means less uncertainty
 - So lower = better



 Data sparsity - many n-grams do not appear in the training data

- Data sparsity many n-grams do not appear in the training data
 - Can be handled by smoothing, back-off

- Data sparsity many n-grams do not appear in the training data
 - Can be handled by smoothing, back-off
- Lack of generalization

- Data sparsity many n-grams do not appear in the training data
 - Can be handled by smoothing, back-off
- Lack of generalization
 - "chases a cat", "chases a dog"...

- Data sparsity many n-grams do not appear in the training data
 - Can be handled by smoothing, back-off
- Lack of generalization
 - "chases a cat", "chases a dog"...
 - "chases an ostrich"?

- Data sparsity many n-grams do not appear in the training data
 - Can be handled by smoothing, back-off
- Lack of generalization
 - "chases a cat", "chases a dog"...
 - "chases an ostrich"?





 Start with one-hot encoding of each word



- Start with one-hot encoding of each word
- Learn word representations in a continuous space



- Start with one-hot encoding of each word
- Learn word representations in a continuous space
- Hidden layer using a non-linear activation



- Start with one-hot encoding of each word
- Learn word representations in a continuous space
- Hidden layer using a non-linear activation
- Output probabilities using softmax





• Experiment details (Sundermeyer, Ney & Schülter, 2015):



- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words



- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words
 - training text: 50M words



- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words
 - training text: 50M words
 - development corpus: 39k words



- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words
 - training text: 50M words
 - development corpus: 39k words
 - evaluation corpus: 35k words



- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words
 - training text: 50M words
 - development corpus: 39k words
 - evaluation corpus: 35k words
- Network structure:


- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words
 - training text: 50M words
 - development corpus: 39k words
 - evaluation corpus: 35k words
- Network structure:
 - projection layer: 300 nodes (per word)



- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words
 - training text: 50M words
 - development corpus: 39k words
 - evaluation corpus: 35k words
- Network structure:
 - projection layer: 300 nodes (per word)
 - hidden layer: 600 nodes



- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words
 - training text: 50M words
 - development corpus: 39k words
 - evaluation corpus: 35k words
- Network structure:
 - projection layer: 300 nodes (per word)
 - hidden layer: 600 nodes
 - total number of params: $128k \cdot 300 + 600 \cdot 128k = 115M$



- Experiment details (Sundermeyer, Ney & Schülter, 2015):
 - vocabulary size: 128k words
 - training text: 50M words
 - development corpus: 39k words
 - evaluation corpus: 35k words
- Network structure:
 - projection layer: 300 nodes (per word)
 - hidden layer: 600 nodes
 - total number of params: $128k \cdot 300 + 600 \cdot 128k = 115M$



Approach	PPL
4-gram count model	163.7
10-gram MLP	136.5
10-gram MLP with 2 layers	130.9



• MLP LM's are still limited in history (use n-gram assumption)



- MLP LM's are still limited in history (use) n-gram assumption)
- We would like to use RNN's to model the entire sentence "at once"



- MLP LM's are still limited in history (use) n-gram assumption)
- We would like to use RNN's to model the entire sentence "at once"
- Every input is a 1-hot vector, every output is the LM probabilities as softmax



• RNN's provide significant improvements over previous models

Approach	PPL
count model	163.7
10-gram MLP	136.5
RNN	125.2
LSTM-RNN	107.8
10-gram MLP with 2 layers	130.9
LSTM-RNN with 2 layers	100.5

- RNN's provide significant improvements over previous models
- A price to pay: longer training time

Approach	PPL
count model	163.7
10-gram MLP	136.5
RNN	125.2
LSTM-RNN	107.8
10-gram MLP with 2 layers	130.9
LSTM-RNN with 2 layers	100.5

Models	PPL	CPU Time (Order)
Count model	163.7	30 min
MLP	136.5	1 week
LSTM-RNN	107.8	3 weeks

 Continuous word representations are learned for each word during network training (a.k.a "word embeddings")

- Continuous word representations are learned for each word during network training (a.k.a "word embeddings")
- These representations can be useful for various tasks like word similarity and word analogies:



 word2vec (mikolov et al. 2003) introduced two similar models:

 word2vec (mikolov et al. 2003) introduced two similar models:



CBOW works and how we must learn the transfer matrices

- word2vec (mikolov et al. 2003) introduced two similar models:
- CBOW (left) and skip-gram (right), both can be seen as MLP's



- word2vec (mikolov et al. 2003) introduced two similar models:
- CBOW (left) and skip-gram (right), both can be seen as MLP's
- Have been shown to approximate the PMI matrix (Levy & Goldberg, 2015)



Reminder - Statistical MT



- Reminder Statistical MT
- Idea use an MLP to train a translation model



- Reminder Statistical MT
- Idea use an MLP to train a translation model
- Inputs are 1-hot encodings of the words in the aligned source language window





- Reminder Statistical MT
- Idea use an MLP to train a translation model
- Inputs are 1-hot encodings of the words in the aligned source language window
- Combine this model with the rest while decoding or rescoring



Another idea - Bilingual LM



- Another idea Bilingual LM
- Inputs are 1-hot encodings of:



- Another idea Bilingual LM
- Inputs are 1-hot encodings of:
 - The words in the aligned source language window



- Another idea Bilingual LM
- Inputs are 1-hot encodings of:
 - The words in the aligned source language window
 - Previous words in the translation hypothesis



- Another idea Bilingual LM
- Inputs are 1-hot encodings of:
 - The words in the aligned source language window
 - Previous words in the translation hypothesis



- Another idea Bilingual LM
- Inputs are 1-hot encodings of:
 - The words in the aligned source language window
 - Previous words in the translation hypothesis
- ACL 2014 Best Paper (Devlin et al, 2014)



Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014

• First (modern) models for end-to-end Neural Machine Translation presented by

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling

• First (modern) models for end-to-end Neural Machine Translation presented by

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the
- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)





• First (modern) models for end-to-end Neural Machine Translation presented by

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by



- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the



- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the

- Kalchbrenner et. al. 2013, Sutskever et al., 2014, Cho et al., 2014
- Inspired by RNN language modelling
- encoder-decoder architecture, sequence-to-sequence learning, seq2seq)



Encoder

• First (modern) models for end-to-end Neural Machine Translation presented by

• The Idea - 2 RNN's, one for reading the input and one for writing the output (a.k.a the

• More formally - model p(y|x) using a single neural network:

• More formally - model p(y|x) using a single neural network:

$y = y_1 \dots y_N$

• More formally - model p(y|x) using a single neural network: $y = y_1 \dots y_N$

 $p(y|x) = p(y_1|x)p(y_2|y_1, x)p(y_3|y_1, y_2, x)\dots p(y_N|y_1\dots y_{N-1}, x)$

• More formally - model p(y|x) using a single neural network: $y = y_1 \dots y_N$

$$p(y|x) = p(y_1|x)p(y_2|y_1, x)p$$

- $p(y_3|y_1, y_2, x) \dots p(y_N|y_1 \dots y_{N-1}, x)$
- $p(y_i = word_k | y_{\leq i}, x) = softmax_k(NN_{\Theta}(y_{\leq i}, x))$



"chat"

"le"



"chat"

input vocabulary size

"le"



"chat"



embedding size

input vocabulary size



"chat"

embedding size

input vocabulary size



"chat"



"chat"

output vocabulary size

hidden state size

embedding size

input vocabulary size







output vocabulary

output vocabulary

The Problem with Vanilla seq2seq

"You can't cram the meaning of a whole %&!\$# sentence into a single \$&!#* vector!" Ray Mooney



"attend" at each step to the relevant parts of the input

• Instead of using a single vector as a fixed representation of the input sequence,

- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state

- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state



- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state



- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state



- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state


- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state



- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state



- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state



- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state



- Instead of using a single vector as a fixed representation of the input sequence, "attend" at each step to the relevant parts of the input
- The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state



- "attend" at each step to the relevant parts of the input



• Instead of using a single vector as a fixed representation of the input sequence,

- "attend" at each step to the relevant parts of the input



• Instead of using a single vector as a fixed representation of the input sequence,

- "attend" at each step to the relevant parts of the input



• Instead of using a single vector as a fixed representation of the input sequence,

- "attend" at each step to the relevant parts of the input



• Instead of using a single vector as a fixed representation of the input sequence,

- "attend" at each step to the relevant parts of the input



• Instead of using a single vector as a fixed representation of the input sequence,

• The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state

Attention-based Decoder

- "attend" at each step to the relevant parts of the input



• Instead of using a single vector as a fixed representation of the input sequence,

• The "importance" of each input element to the current prediction is computed via a feed-forward network that gets the input element and the current decoder state

Attention-based Decoder

- "attend" at each step to the relevant parts of the input



• Instead of using a single vector as a fixed representation of the input sequence,

Attention-based Decoder

- "attend" at each step to the relevant parts of the input



• Instead of using a single vector as a fixed representation of the input sequence,

• And a bit more formally - in each decoder step:

- And a bit more formally in each decoder step: •
 - Compute attention scores for each input element: $\operatorname{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s) = \operatorname{tanh}(\boldsymbol{W}_{\boldsymbol{a}}[\boldsymbol{h}_t; \bar{\boldsymbol{h}}_s])$

- And a bit more formally in each decoder step: •
 - Compute attention scores for each input element: $\operatorname{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s) = \tanh(\boldsymbol{W}_{\boldsymbol{a}}[\boldsymbol{h}_t; \bar{\boldsymbol{h}}_s])$
 - Normalize the attention scores so they sum up to 1: $oldsymbol{a}_t(s) = ext{align}(oldsymbol{h}_t, oldsymbol{ar{h}}_s) = rac{ ext{exp}\left(ext{score}(oldsymbol{h}_t, oldsymbol{ar{h}}_s)
 ight)}{\sum_{s'} ext{exp}\left(ext{score}(oldsymbol{h}_t, oldsymbol{ar{h}}_{s'})
 ight)}$

- And a bit more formally in each decoder step: •
 - Compute attention scores for each input element: $\operatorname{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s) = \tanh(\boldsymbol{W}_{\boldsymbol{a}}[\boldsymbol{h}_t; \bar{\boldsymbol{h}}_s])$
 - Normalize the attention scores so they sum up to 1: $\boldsymbol{a}_{t}(s) = \operatorname{align}(\boldsymbol{h}_{t}, \bar{\boldsymbol{h}}_{s}) = \frac{\exp\left(\operatorname{score}(\boldsymbol{h}_{t}, \bar{\boldsymbol{h}}_{s})\right)}{\sum_{s'} \exp\left(\operatorname{score}(\boldsymbol{h}_{t}, \bar{\boldsymbol{h}}_{s'})\right)}$ • Compute ct: $c_{t} = \sum_{j=1}^{T_{x}} a_{j} \bar{h}_{j}$

- And a bit more formally in each decoder step:
 - Compute attention scores for each input element: $\operatorname{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s) = \tanh(\boldsymbol{W}_{\boldsymbol{a}}[\boldsymbol{h}_t; \bar{\boldsymbol{h}}_s])$
 - Normalize the attention scores so they sum up to 1: $a_t(s) = \operatorname{align}(h_t, \bar{h}_s) = \frac{\exp\left(\operatorname{score}(h_t, \bar{h}_s)\right)}{\sum_{s'} \exp\left(\operatorname{score}(h_t, \bar{h}_{s'})\right)}$ • Compute ct: $c_t = \sum_{j}^{T_x} a_j \bar{h}_j$
 - Compute attention output state:

$$ilde{m{h}}_t = anh(m{W}_{m{c}}[m{c}_t;m{h}_t])$$

- And a bit more formally in each decoder step:
 - Compute attention scores for each input element: $\operatorname{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s) = \tanh(\boldsymbol{W}_{\boldsymbol{a}}[\boldsymbol{h}_t; \bar{\boldsymbol{h}}_s])$
 - Normalize the attention scores so they sum up to 1: $a_t(s) = \operatorname{align}(h_t, \bar{h}_s) = \frac{\exp\left(\operatorname{score}(h_t, \bar{h}_s)\right)}{\sum_{s'} \exp\left(\operatorname{score}(h_t, \bar{h}_{s'})\right)}$ • Compute ct: $c_t = \sum_{j}^{T_x} a_j \bar{h}_j$
 - Compute attention output state:

$$\tilde{\boldsymbol{h}}_t = anh(\boldsymbol{W_c}[\boldsymbol{c}_t; \boldsymbol{h}_t])$$

• Compute output probability distribution: $p(y_t|y_{< t}, x) = \operatorname{softmax}(W_s h_t)$

which are updated as-you-go

Instead of keeping one best option on each time step, keep k best options

 Instead of keeping one best option which are updated as-you-go



Greedy Search

Instead of keeping one best option on each time step, keep k best options



Beam Search (k=2)

- which are updated as-you-go

Greedy Search

Instead of keeping one best option on each time step, keep k best options

• Requires to maintain different RNN states in-memory for each hypothesis



Beam Search (k=2)

- Instead of keeping one best option on each time step, keep k best options which are updated as-you-go
- Requires to maintain different RNN states in-memory for each hypothesis
- Usually a small beam size is enough (5-12)



Greedy Search



Beam Search (k=2)



Effect of Sentence Length

Effect of Sentence Length

 Results deteriorate for longer sentences as they are "compressed" to a fixed length vector



Bahdanau, Cho and Bengio (2014)

Effect of Sentence Length

- Results deteriorate for longer sentences as they are "compressed" to a fixed length vector
- The attention mechanism "opens the bottleneck"





Bahdanau, Cho and Bengio (2014)

• The model learns to "align" source and target representations





- The model learns to "align" source and target representations
- No explicit alignment supervision was given!





- The model learns to "align" source and target representations
- No explicit alignment supervision was given!
- But this isn't always perfect...



Figure 9: Mismatch between attention states and desired word alignments (German–English).

Koehn and Knowles, 2017

Self Attention and The Transformer Architecture

Self Attention and The Transformer Architecture

• Vaswani et al. (2017)

Self Attention and The Transformer Architecture

- Vaswani et al. (2017)
 - "Attention is All You Need"
- Vaswani et al. (2017)
 - "Attention is All You Need"
- Main idea: replace RNNs with self attention layers

- Vaswani et al. (2017)
 - "Attention is All You Need"
- Main idea: replace RNNs with self attention layers



- Vaswani et al. (2017)
 - "Attention is All You Need"
- Main idea: replace RNNs with self attention layers



- Vaswani et al. (2017)
 - "Attention is All You Need"
- Main idea: replace RNNs with self attention layers



- Vaswani et al. (2017)
 - "Attention is All You Need"
- Main idea: replace RNNs with self attention layers
- Can be parallelized at the sequence level - faster training of large networks





Positional encodings \bullet



- Positional encodings •
- Multi-head attention lacksquare



- Positional encodings •
- Multi-head attention ullet
- Layer normalization ullet



- Positional encodings
- Multi-head attention ullet
- Layer normalization \bullet
- Decoder masked self attention ●



- Positional encodings
- Multi-head attention ullet
- Layer normalization
- Decoder masked self attention ullet
- Unlike LSTM based models-



- Positional encodings
- Multi-head attention ullet
- Layer normalization
- Decoder masked self attention ullet
- Unlike LSTM based models-
 - encoder-decoder-attention in each layer!



- Positional encodings lacksquare
- Multi-head attention
- Layer normalization
- Decoder masked self attention
- Unlike LSTM based models-
 - encoder-decoder-attention in each layer!
 - Less interpretable



- Positional encodings
- Multi-head attention lacksquare
- Layer normalization
- Decoder masked self attention
- Unlike LSTM based models-
 - encoder-decoder-attention in each layer!
 - Less interpretable
- Learning rate schedule more sensitive to hyperparams than LSTM-based models





• Neural networks are very useful for language modeling



- Neural networks are very useful for language modeling
- Better generalization, more context



- Neural networks are very useful for language modeling
- Better generalization, more context
- Neural language models: from MLPs to RNNs



- Neural networks are very useful for language modeling
- Better generalization, more context
- Neural language models: from MLPs to RNNs
- Sequence-to-Sequence Learning



- Neural networks are very useful for language modeling
- Better generalization, more context
- Neural language models: from MLPs to RNNs
- Sequence-to-Sequence Learning
- The Attention Mechanism



- Neural networks are very useful for language modeling
- Better generalization, more context
- Neural language models: from MLPs to RNNs
- Sequence-to-Sequence Learning
- The Attention Mechanism
- The Transformer Architecture



