







Text Processing using Machine Learning

Attention

Liling Tan
05 Dec 2019

5,500 GRADUATE ALUMNI

120 & LEADERSHIP PROGRAMMES

TRAINING OVER

120,000 DIGITAL LEADERS

PROFESSIONALS

Overview





<u>Lecture</u>

Attention

Hands-on

Seq2Seq with Attention





Attention

Recurrent Continuous Translation Models





Add the encoded hidden state at every decoder time step

(Kalchbrenner & Blunsom, 2013) is never **Predict Predict RNN** RNN RNN h_{x} language is never

Sentence Representation







"You can't cram the meaning of a whole !@#\$-ing sentence into a single %^&*-ing vector!"

- Raymond Mooney

Sentence Representation

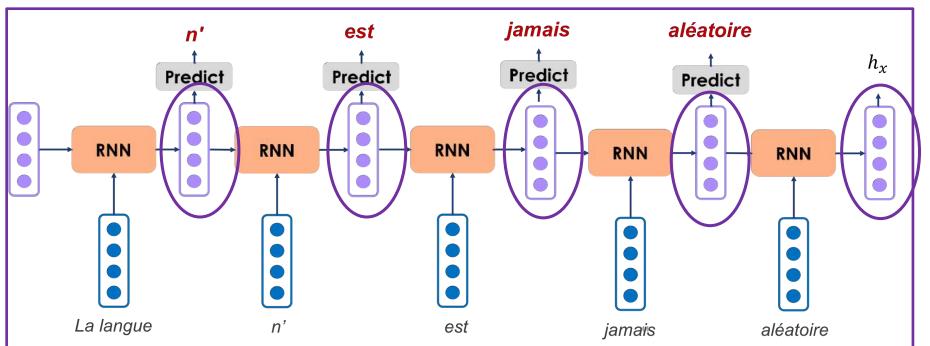






"You can't cram the meaning of a whole !@#\$-ing sentence into a single %^&*-ing vector!"

- Raymond Mooney

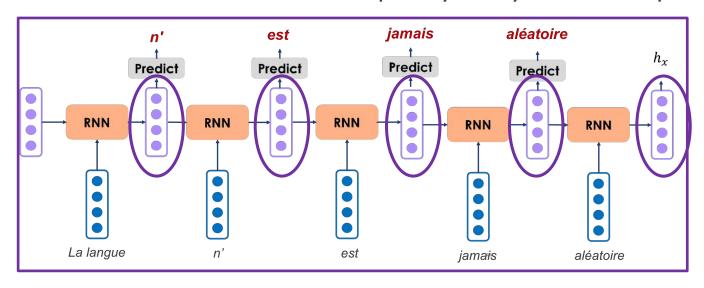


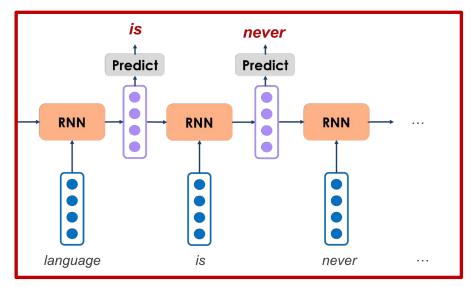
What if instead of taking just h_x , we take <u>all</u> the hidden states?





- Encode each word in input into a vector
- When decoding, linearly combine encoded vector vectors weighted by "attention"
- Bahdanau et al. (2015) proposes to learn a MLP mapping between the query-key vector pairs

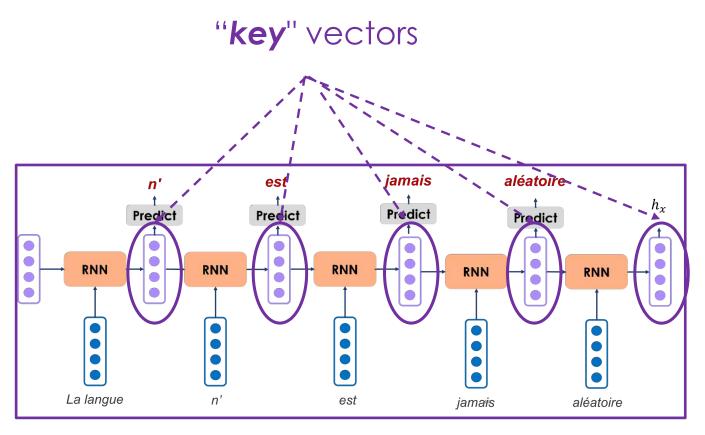


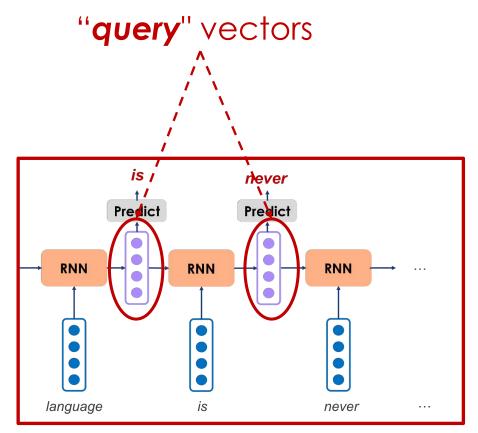






 <u>Bahdanau et al. (2015)</u> proposes to learn a MLP mapping between the query-key vector pairs



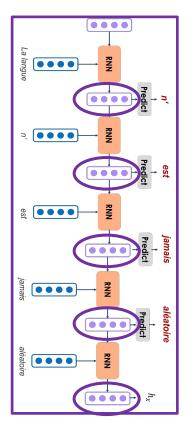


Encoder Model

Decoder Model







$$a(q_1, k_1) = 7.5$$

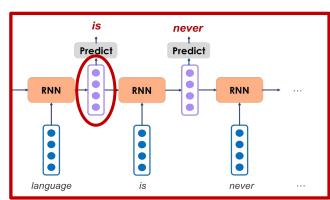
$$a(q_1, k_2) = 1.9$$

$$a(q_1, k_3) = 2.1$$

$$a(q_1, k_4) = -3.2$$

$$a(q_1, k_5) = -0.5$$

Encoder Model



Alignment score between "key" and "query"

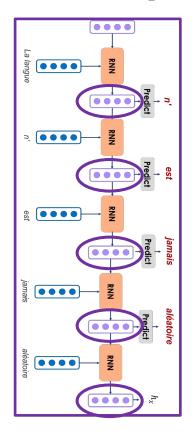
$$a(\mathbf{q}, k) = w_2^T(W_1[\mathbf{q}; k])$$

Can be learned by a w_2^T vector and W_1 matrix



 $a(q, k) = w_2^T(W_1[q; k])$





$$a(q_1, k_1) = 7.5$$

$$a(q_1, k_2) = 1.9$$

$$a(q_1, k_3) = 2.1$$

$$a(q_1, k_4) = -3.2$$

$$a(q_1, k_5) = -0.5$$

0.9

m

X

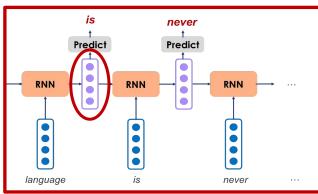
0.003

0.004

0.00005

0.0003

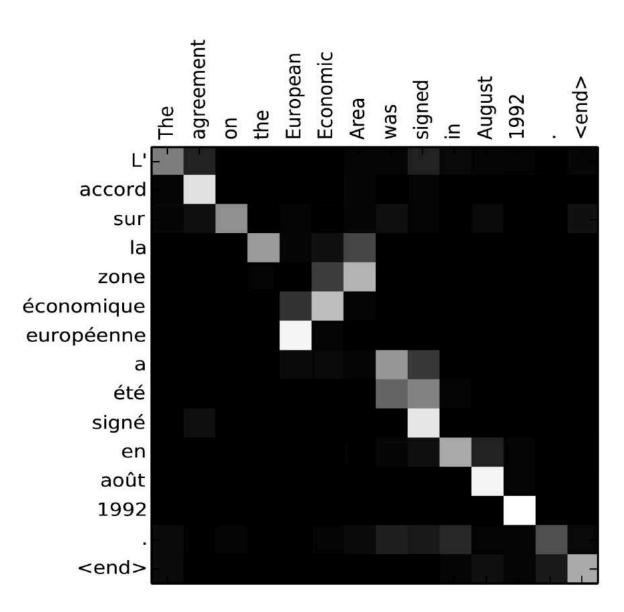
Encoder Model







- Bahdanau et al. (2015)
 proposes to learn a MLP mapping between the query-key vector pairs
- Each decoder state
 generates a probabilistic
 "attention" vector that
 aligns the predicted word
 at the time step to every
 input word in the encoder
- As a by-product, it creates a word-alignment matrix (e.g. figure on right)







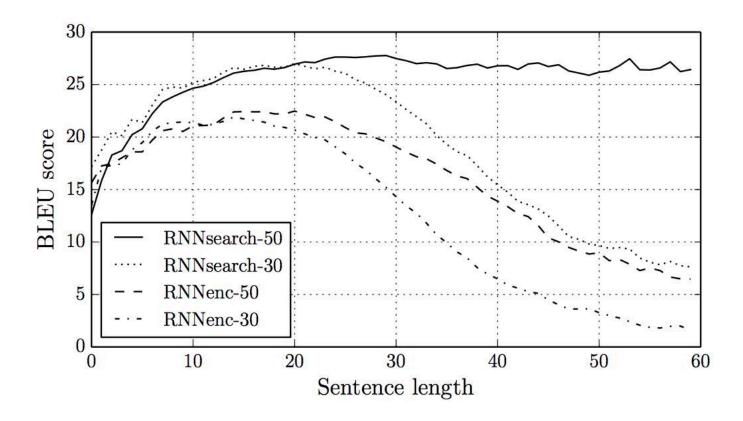


Figure 2: The BLEU scores of the generated translations on the test set with respect to the lengths of the sentences. The results are on the full test set which includes sentences having unknown words to the models.

Bahdanau et al. (2015)

Many Other Attentions





Name	Alignment score function	Citation	
Content-base attention	$score(s_t, \boldsymbol{h}_i) = cosine[s_t, \boldsymbol{h}_i]$	Graves2014	
Additive(*)	$\operatorname{score}(\boldsymbol{s}_t, \boldsymbol{h}_i) = \mathbf{v}_a^{\top} \tanh(\mathbf{W}_a[\boldsymbol{s}_t; \boldsymbol{h}_i])$	Bahdanau201	
Location-	$\alpha_{t,i} = \operatorname{softmax}(\mathbf{W}_a \mathbf{s}_t)$	Luong2015	
Base	Note: This simplifies the softmax alignment to only depend on the target position.		
General	$\operatorname{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^{\top} \mathbf{W}_a \mathbf{h}_i$ where \mathbf{W}_a is a trainable weight matrix in the attention layer.	Luong2015	
Dot-Product	$score(s_t, \boldsymbol{h}_i) = s_t^{\top} \boldsymbol{h}_i$	Luong2015	
Scaled Dot- Product(^)	$score(s_t, h_i) = \frac{s_t^{T} h_i}{\sqrt{n}}$	Vaswani2017	
	Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.		

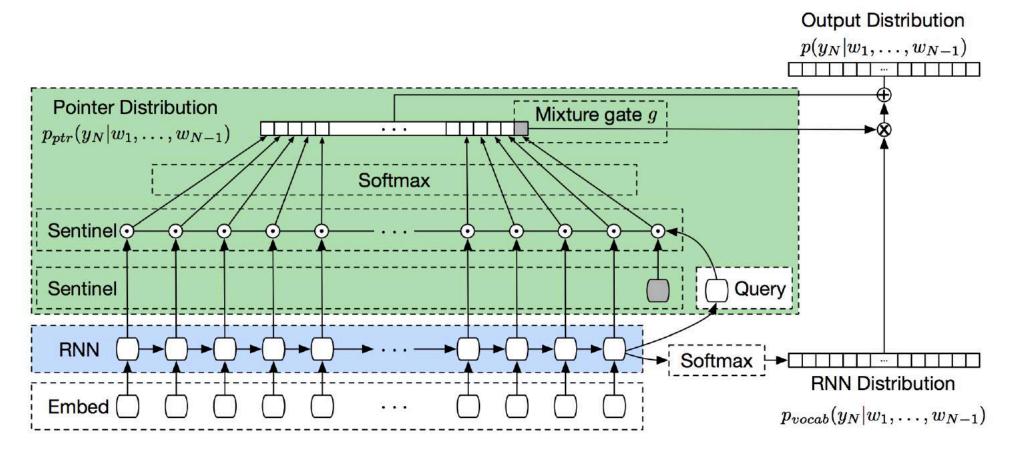
- (*) Referred to as "concat" in Luong, et al., 2015 and as "additive attention" in Vaswani, et al., 2017.
- (^) It adds a scaling factor $1/\sqrt{n}$, motivated by the concern when the input is large, the softmax function may have an extremely small gradient, hard for efficient learning.

Reproducing Previous Generated Word





 Merity et al. (2016) applied a mixture model of softmax and pointers.

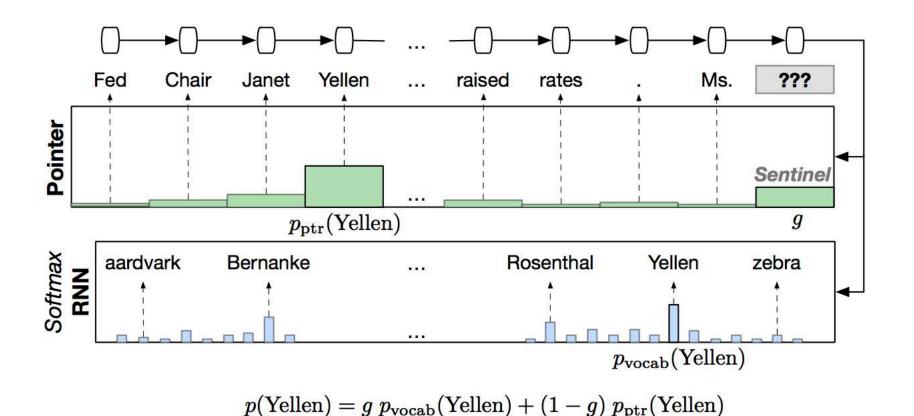


Reproducing Previous Generated Word





 Merity et al. (2016) applied a mixture model of softmax and pointers.



Show, Attend and Tell





<u>Xu et al. (2015)</u> generate captions from images and found that attention activates for words with respect to specific regions in the image























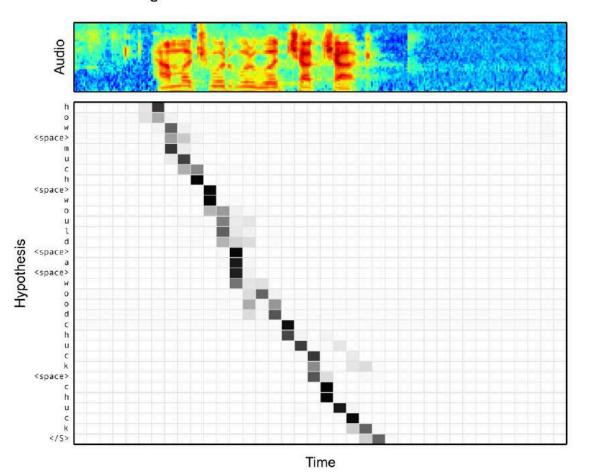
Listen, Attend and Spell





Chan et al. (2015) generate characters from audios and found that attention learns start and end of utterance

Alignment between the Characters and Audio



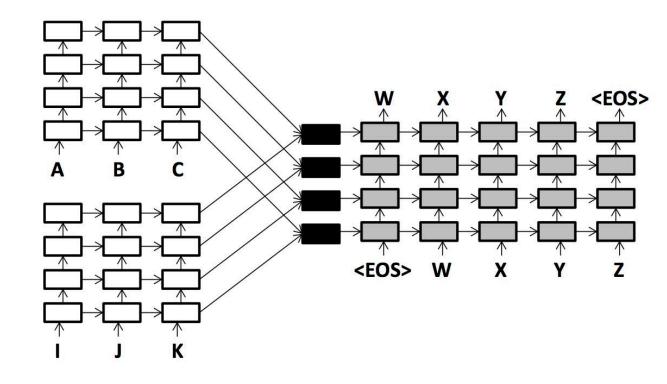
Multiple Source Attention





Zoph and Knight (2015) combines multiple encoders through the attention mechanism for NMT

Libovicky and Helcl (2017) suggested a couple of strategies for combined attention



Source 1: UNK Aspekte sind ebenfalls wichtig.

Target: UNK aspects are important, too.

Source 2: Les aspects UNK sont également importants.

Self-Attention





Cheng et al. (2016) proposed self-attention relating the current word with previous words

```
The FBI is chasing a criminal on the run.
The FBI is chasing a criminal on the run.
The FBI is chasing a criminal on the run.
     FBI is chasing a criminal on the run.
The
              chasing a criminal on the run.
The
     FBI is chasing a criminal on the run.
     FBI is
              chasing a criminal on the run.
The
              chasing a
                          criminal on the run.
The
              chasing a
                          criminal on
                                        the run.
              chasing a
                          criminal
                                    on
                                        the run.
```

Multi-headed Attention





Create multiple attention "heads" to focus on different parts of the sentence.

Allamanis et al. (2016) creates different heads for "copy" vs regular convolutional attention

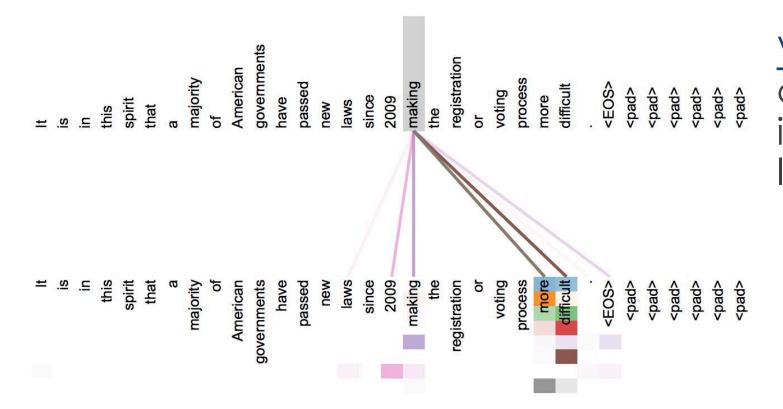
Target		Attention Vectors	
m_1	set	$lpha = \frac{\text{ss} \{ \text{ this . use } \underline{\text{Browser }} \underline{\text{Cache}} = \underline{\text{use }} \underline{\text{Browser }} \underline{\text{Cache}}; \} }{\kappa = \frac{\text{ss} \{ \text{ this . use }}{\text{Browser }} \underline{\text{Cache}} = \underline{\text{use }} \underline{\text{Browser }} \underline{\text{Cache}}; \} }$	0.012
m_2	use	$lpha=$ <s> { this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } </s> $\kappa=$ <s> { this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } </s>	0.974
m_3	browser	$lpha=$ <s>{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache; } </s> $\kappa=$ <s>{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache; } </s>	0.969
m_4	cache	$lpha = \frac{\sline Browser}{\sline Browser} \sline Browser} \sline Browser Cache = use Browser Cache; } \kappa = \frac{\sline Browser}{\sline Browser} \sline Browser} \sline Browser Cache; } $	0.583
m_5	END	$lpha = \frac{\mbox{{ this . use } Browser \ Cache = use Browser \ Cache ; } \kappa = \frac{\mbox{{ this . use } Browser \ Cache = use Browser \ Cache ; } }$	0.066

Multi-headed Attention





Create multiple attention "heads" to focus on different parts of the sentence.



Vaswani et al. (2017) created multiple independently learned "heads"

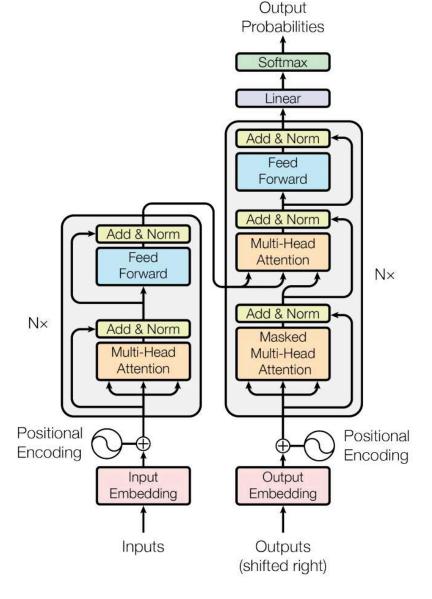
Attention is All You Need





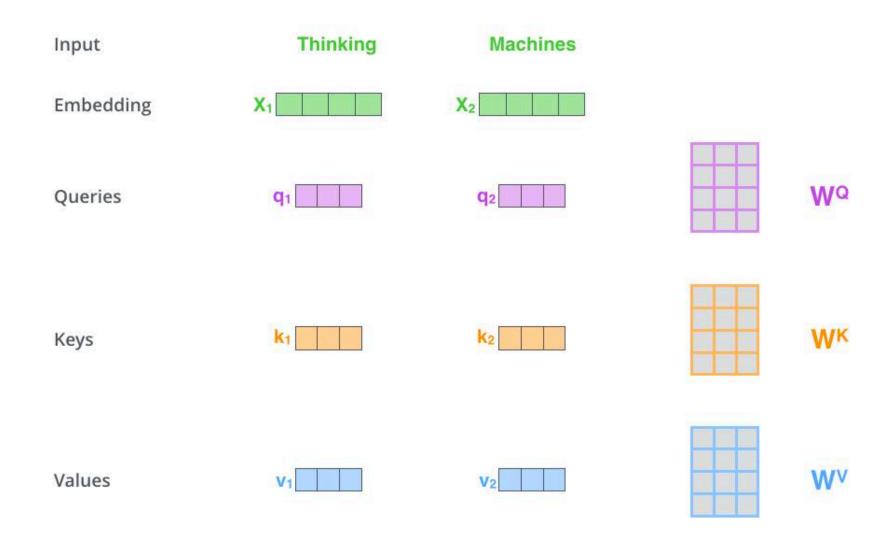
 Sequence-to-Sequence model entirely based on attention, no recurrence, no convolution

- Fast because only matrix multiplication operation
- It's everywhere...



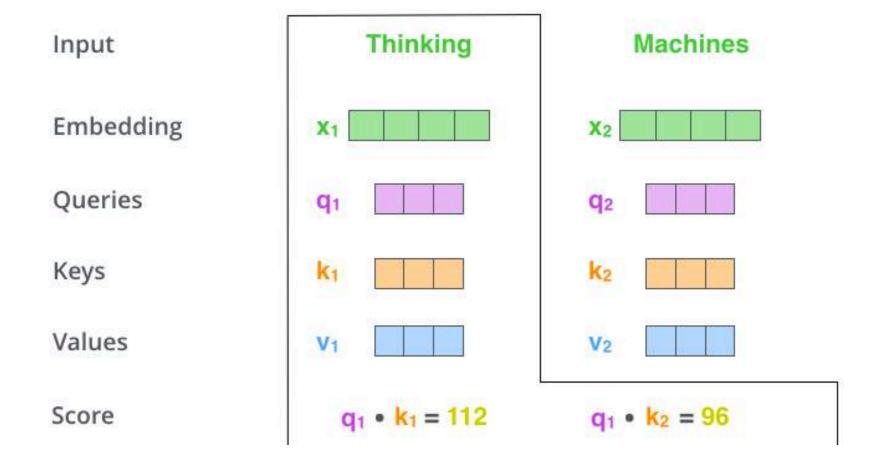
















Input	Thinking	Machines
Embedding	X1	X ₂
Queries	q ₁	q ₂
Keys	k ₁	k ₂
Values	V ₁	V ₂
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Divide by 8 ($\sqrt{d_k}$)	14	12
Softmax	0.88	0.12

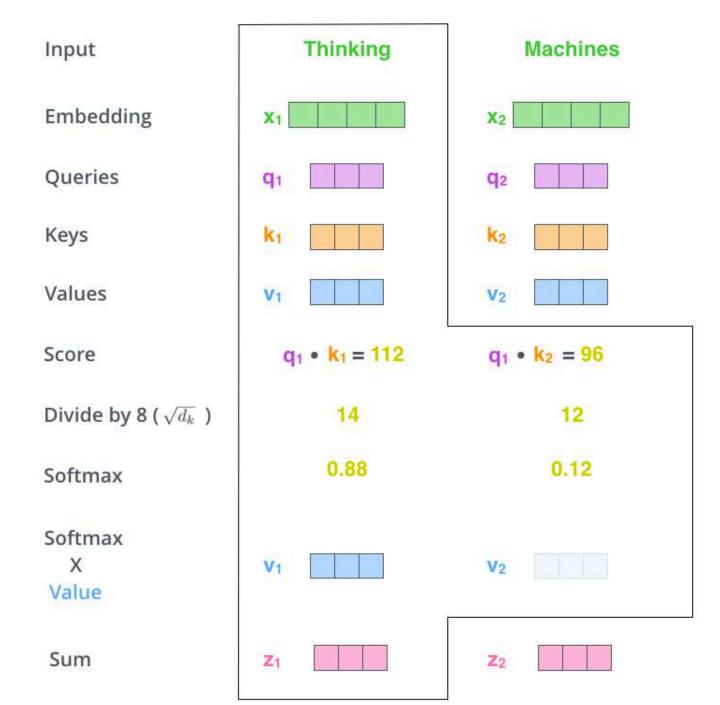


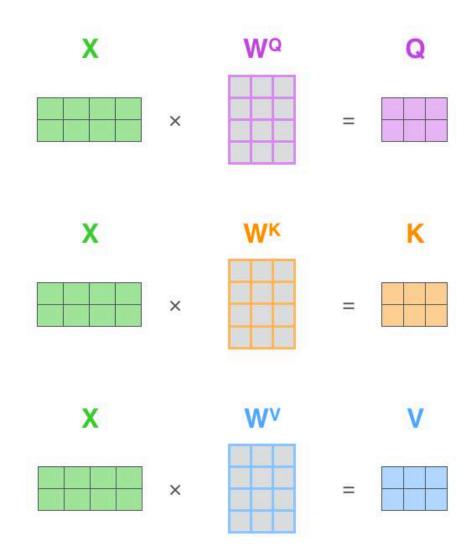


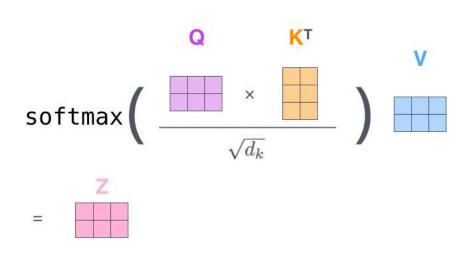


Image from http://jalammar.github.io/illustrate d-transformer/



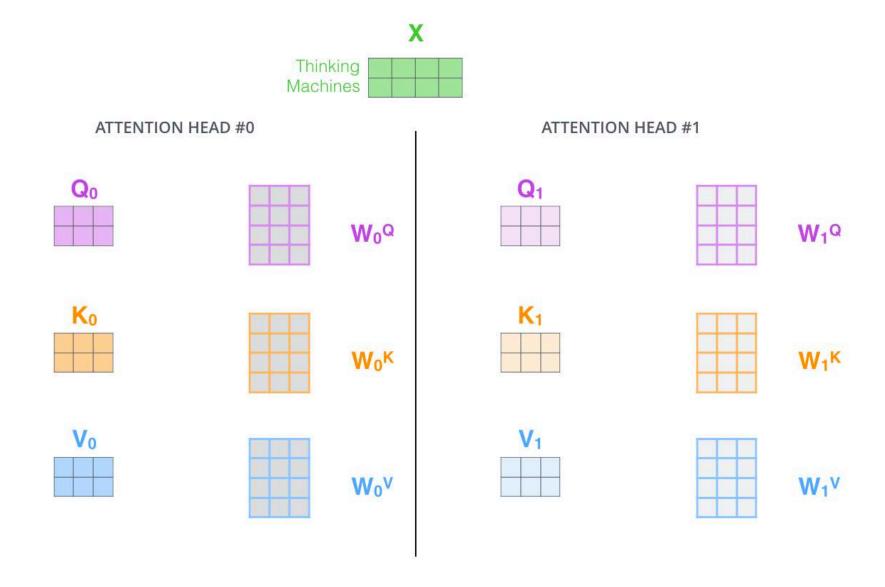






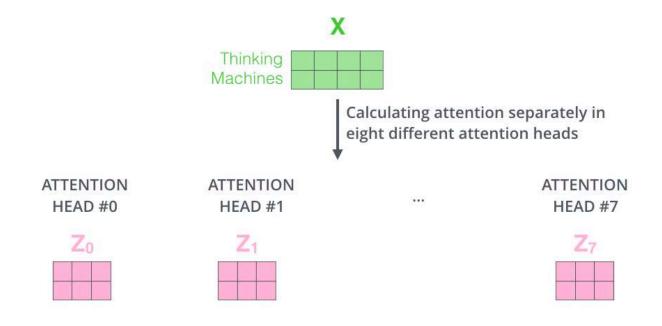














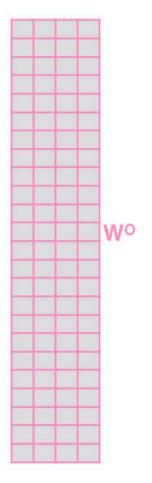


1) Concatenate all the attention heads



2) Multiply with a weight matrix W^o that was trained jointly with the model

X

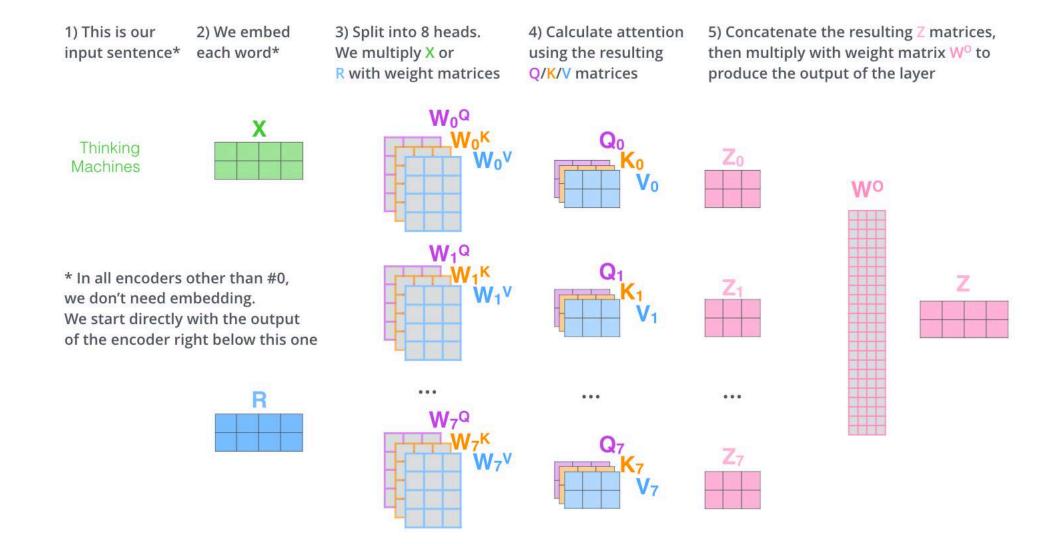


3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

= Z

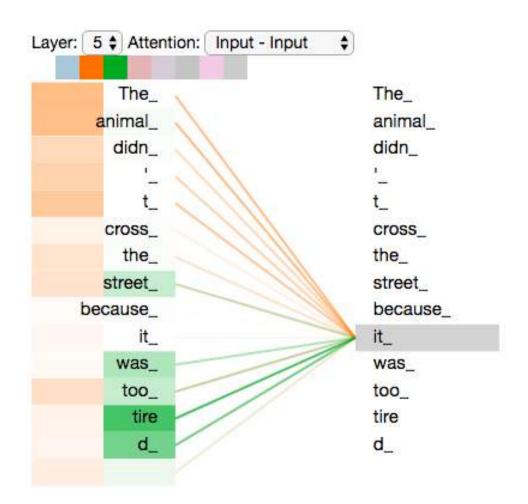


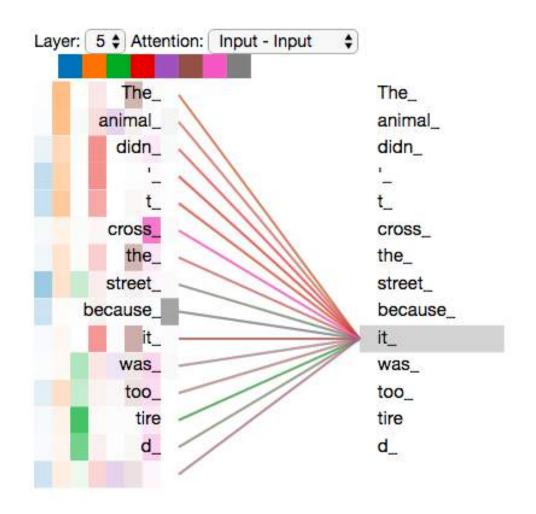












Attention Tricks





Self Attention: Each layer combines words with every other word in sentence

- Multi-headed Attention: Multiple attention heads learnt independently
- Normalized dot product attention: removes bias

 Positional Encodings: Can distinguish positive without stepwise recurrence

Training Tricks





- Layer Normalization: Ensure layer outputs remains in range
- Specialized Training Schedule: Warm-up and cool-up scheduler adjust defaults Adam learning rates
- Label Smoothing: Insert some uncertainty during training to prevent overfitting

The Annotated Transformer





- Lets walk through the code
- See http://bit.ly/ANLP-AnnTransformer

Fin