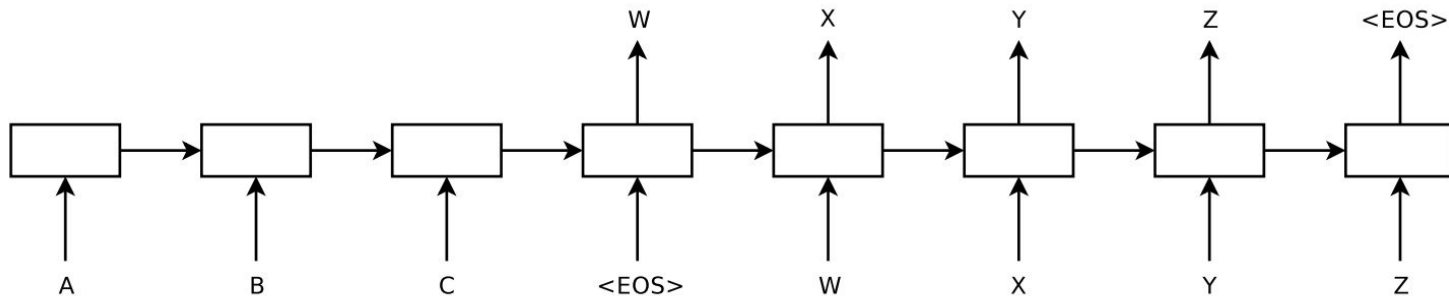


Effective Approaches to Attention-based Neural Machine Translation

Minh-Thang Luong, Hieu Pham, Christopher D. Manning

arxiv.org/abs/1508.04025

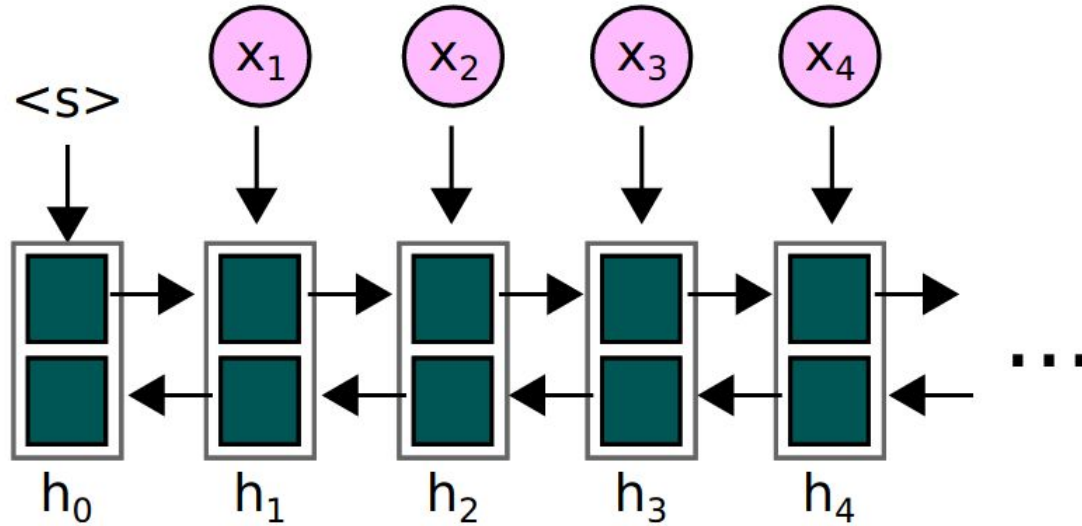
Vanilla Encoder-Decoder



Sentence probability: $\log p(y|x) = \sum_{j=1}^m \log p(y_j | y_{<j}, \mathbf{s})$

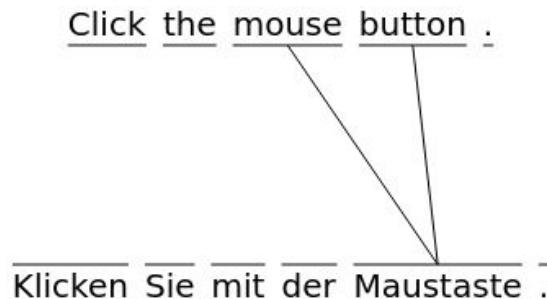
Training loss: $J_t = \sum_{(x,y) \in \mathbb{D}} -\log p(y|x)$

BiRNN Encoder-Decoder

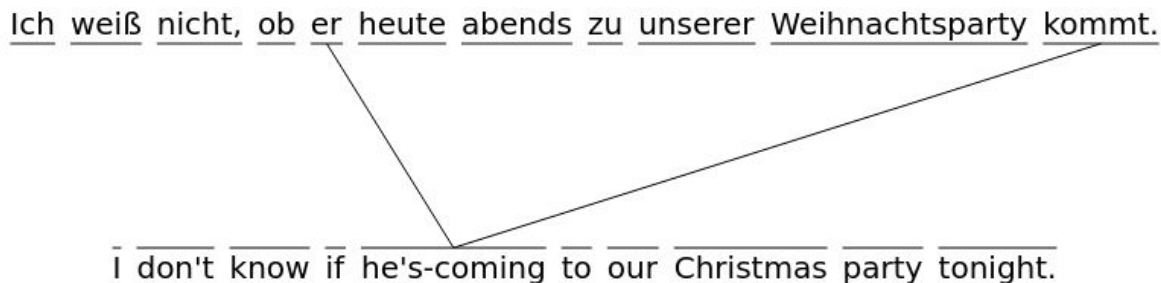


Alignment

Click the mouse button .
Klicken Sie mit der Maustaste .



Ich weiß nicht, ob er heute abends zu unserer Weihnachtsparty kommt.
I don't know if he's-coming to our Christmas party tonight.



Problem:

Vanishing gradients, long term dependencies.

Idea:

Tell the network where to look in the input.

Global Attention

Attention relevance
between the current
decoder state and
input state s

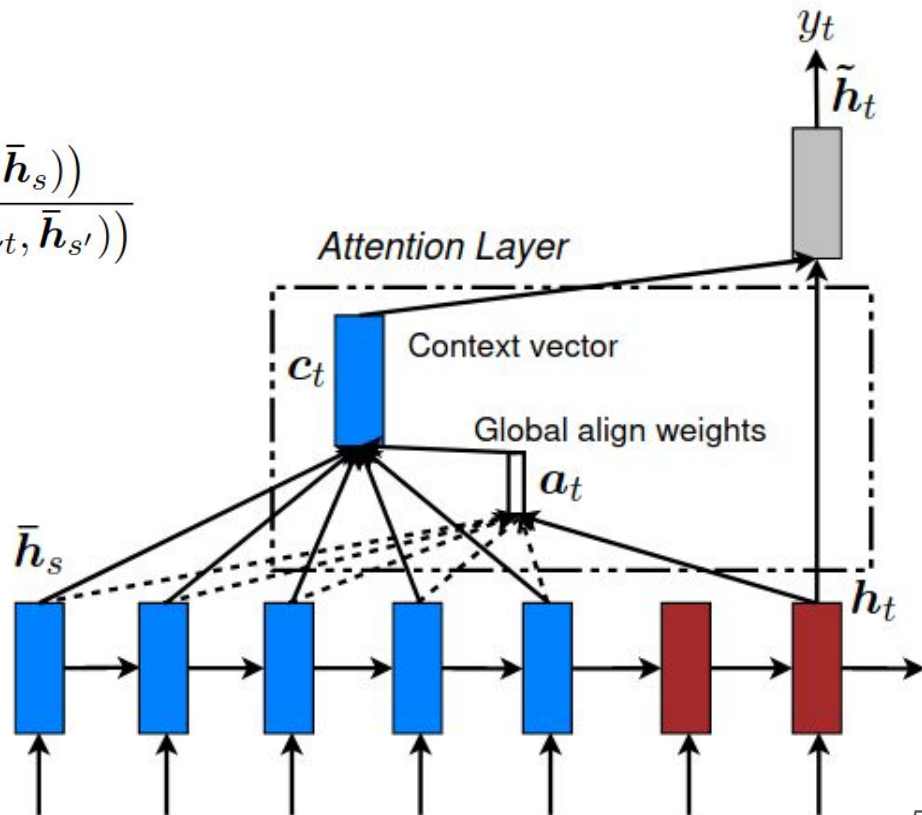
$$\begin{aligned} a_t(s) &= \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \\ &= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \end{aligned}$$

Context & state
concatenation

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t])$$

Word output

$$p(y_t | y_{<t}, x) = \text{softmax}(\mathbf{W}_s \tilde{\mathbf{h}}_t)$$



Computing Score (Globally)

$$\begin{aligned} \mathbf{a}_t(s) &= \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \\ &= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \end{aligned}$$

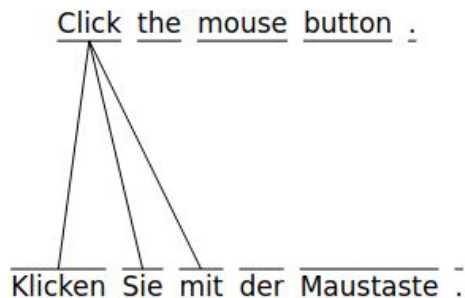
Attention score, not
attention!

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a [\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases}$$

Attention, not attention
score!

$$\mathbf{a}_t = \text{softmax}(\mathbf{W}_a \mathbf{h}_t)$$

Local Attention (local-m, local-p)



Diagonal alignment: $p_t = t$

Considered window of input tokens: $[p_t - D, p_t + D]$

Let the network predict the alignment itself:

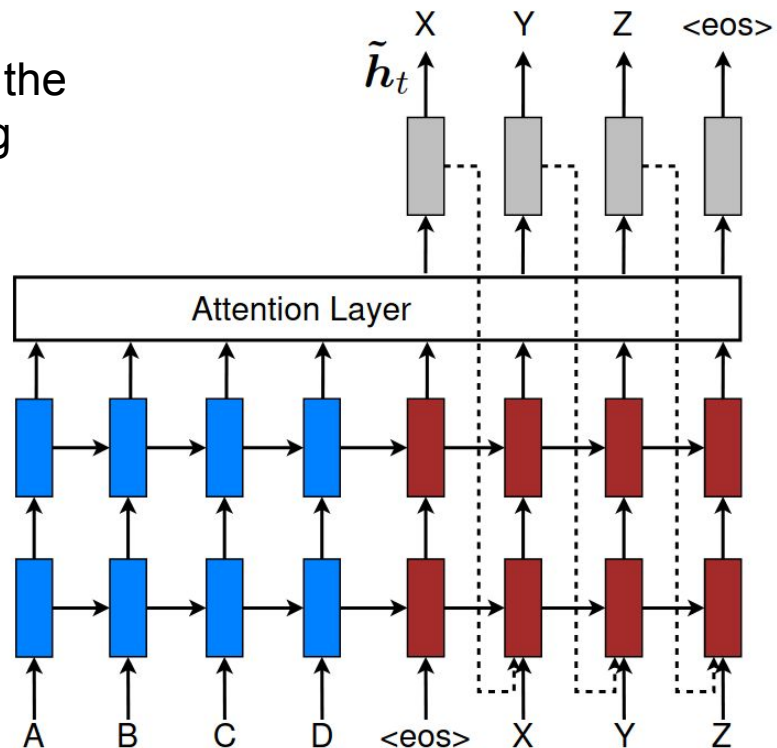
$$p_t = S \cdot \text{sigmoid}(\mathbf{v}_p^\top \tanh(\mathbf{W}_p \mathbf{h}_t)),$$

Force it to focus on tokens around the predicted location:

$$\mathbf{a}_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \exp\left(-\frac{(s - p_t)^2}{2\sigma^2}\right)$$

Input-feeding Approach

Also add the previous hidden state (not the generated word) to the current decoding computation.

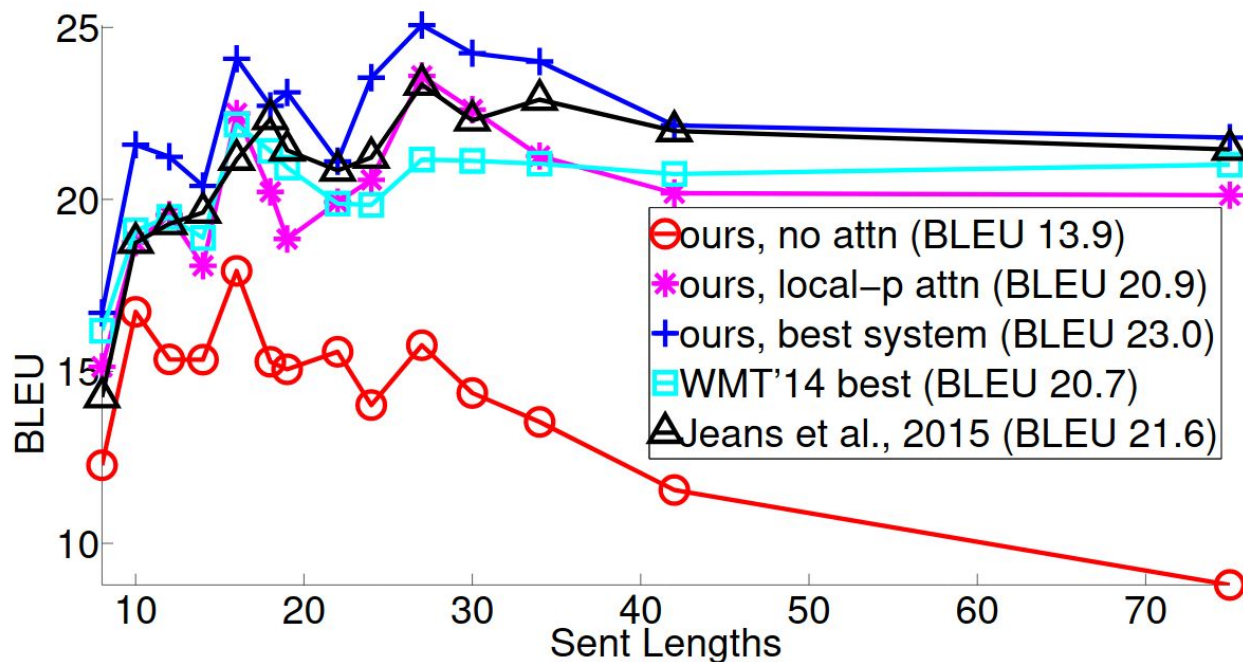


Experiments

System	Ppl	BLEU
Winning WMT'14 system – <i>phrase-based</i> + <i>large LM</i> (Buck et al., 2014)		20.7
<i>Existing NMT systems</i>		
RNNsearch (Jean et al., 2015)		16.5
RNNsearch + unk replace (Jean et al., 2015)		19.0
RNNsearch + unk replace + large vocab + <i>ensemble</i> 8 models (Jean et al., 2015)		21.6
<i>Our NMT systems</i>		
Base	10.6	11.3
Base + reverse	9.9	12.6 (+1.3)
Base + reverse + dropout	8.1	14.0 (+1.4)
Base + reverse + dropout + global attention (<i>location</i>)	7.3	16.8 (+2.8)
Base + reverse + dropout + global attention (<i>location</i>) + feed input	6.4	18.1 (+1.3)
Base + reverse + dropout + local-p attention (<i>general</i>) + feed input	5.9	19.0 (+0.9)
Base + reverse + dropout + local-p attention (<i>general</i>) + feed input + unk replace		20.9 (+1.9)
<i>Ensemble</i> 8 models + unk replace		23.0 (+2.1)

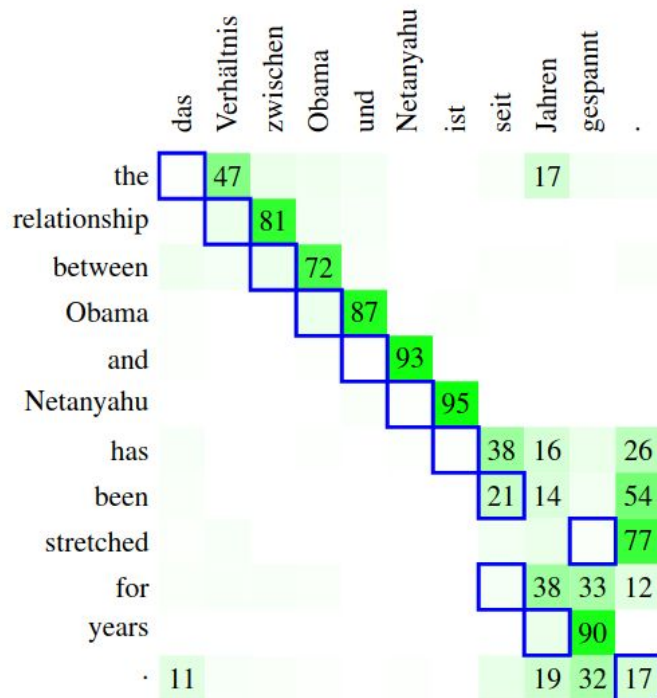
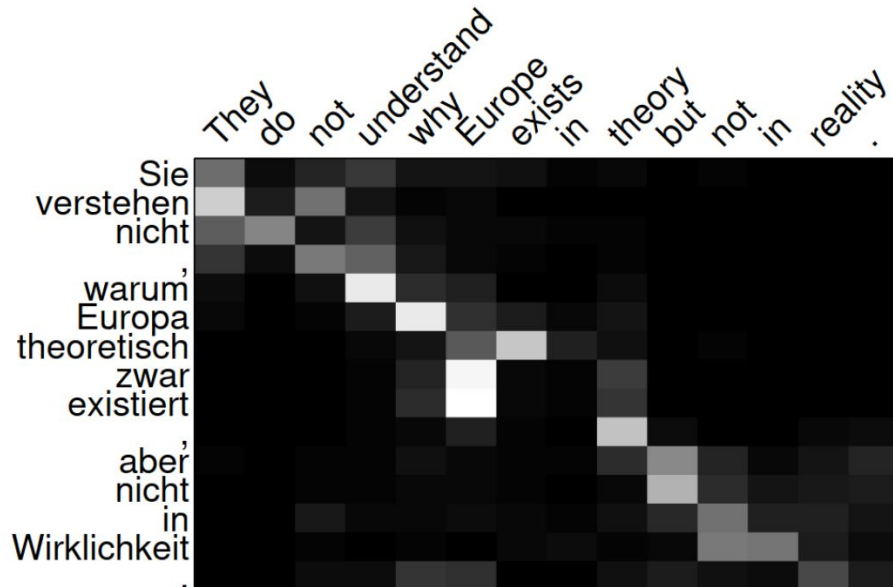
Sentence Length

Attention saves long sentences, especially compared to vanilla RNN.



Attention Note #1

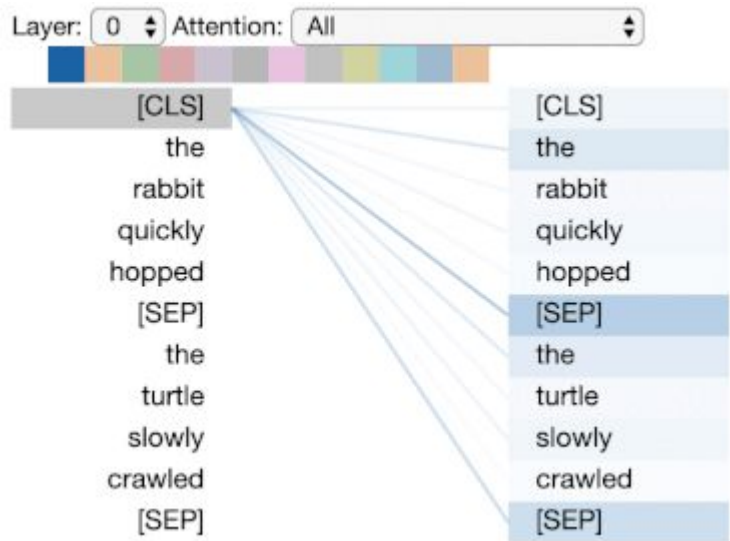
Why is alignment off by one?



Attention Note #2

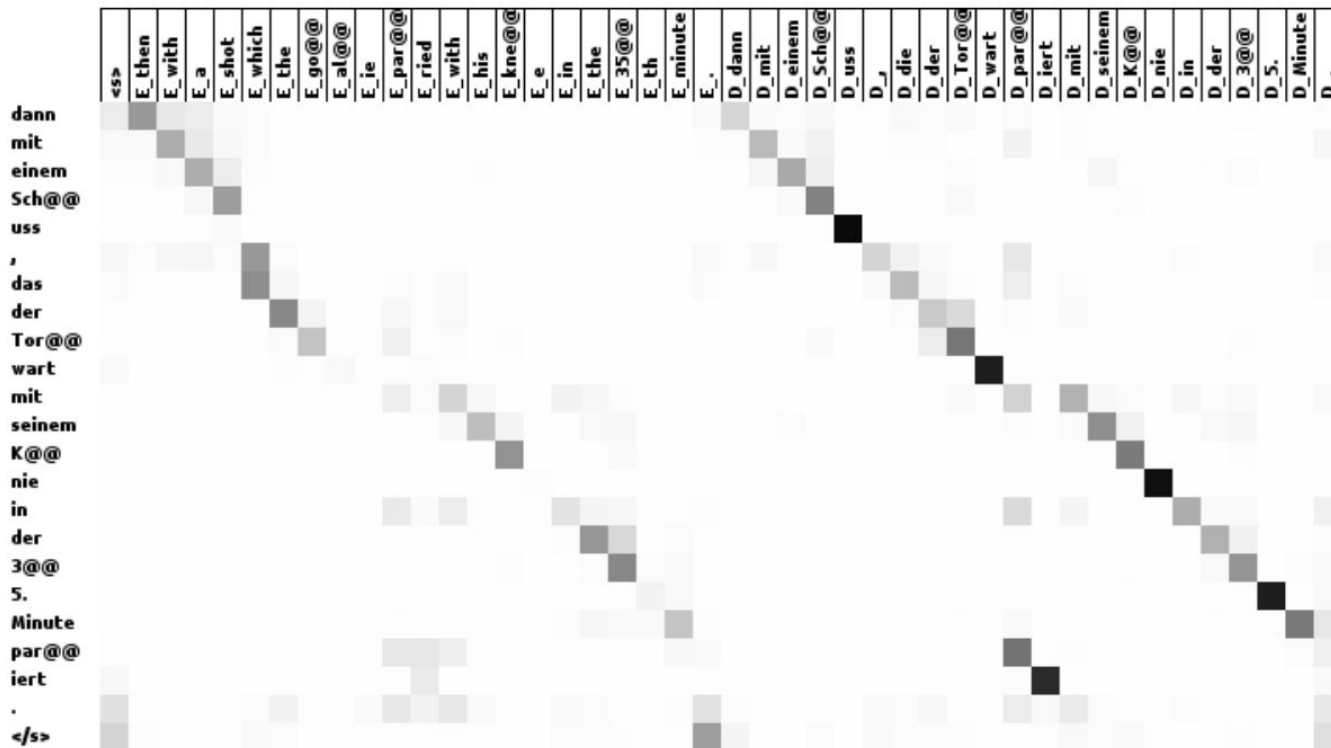
Attention heads encode linguistic structures.

Heads from Transformer self-attention resemble linguistic dependency.



Attention Note #3

Append less
performant MT
output as an
input to attention
based MT



Summary

- vanilla RNN suffers from vanishing gradients
- attention brings input closer to the output
- global attention - whole sentence
- local attention - window of tokens

Resources

- Effective Approaches to Attention-based Neural Machine Translation, Minh-Thang Luong, Hieu Pham, Christopher D. Manning
<https://arxiv.org/abs/1508.04025>
- Class on Statistical Machine Translation
Ondřej Bojar
<http://ufal.mff.cuni.cz/courses/npfl087>
- Enabling Outbound Machine Translation
My bachelor thesis
<https://dspace.cuni.cz/bitstream/handle/20.500.11956/119400/130284419.pdf>
- Six Challenges for Neural Machine Translation
Philipp Koehn, Rebecca Knowles
<https://arxiv.org/abs/1706.03872>
- A Multiscale Visualization of Attention in the Transformer Model
Jesse Vig
<https://github.com/jessevig/bertviz>
- Slow Align
<https://vilda.net/s/slowalign/>