# Seq2seq, NMT, Transformer

**Milan Straka**

📅 **Apr 11, 2021**

EUROPEAN UNION
European Structural and Investment Fund
Operational Programme Research,
Development and Education

**Sequence-to-Sequence Architecture**

Sequence-to-Sequence is a name for an architecture allowing to produce an arbitrary output sequence $y_1, \ldots, y_M$ from an input sequence $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N$.

Unlike CRF/CTC, no assumptions are necessary and we condition each output sequence element on all input sequence elements and all already generated output sequence elements:

$$P(y_i | \boldsymbol{x}_1, \ldots, \boldsymbol{x}_N, y_1, \ldots, y_{i-1}).$$

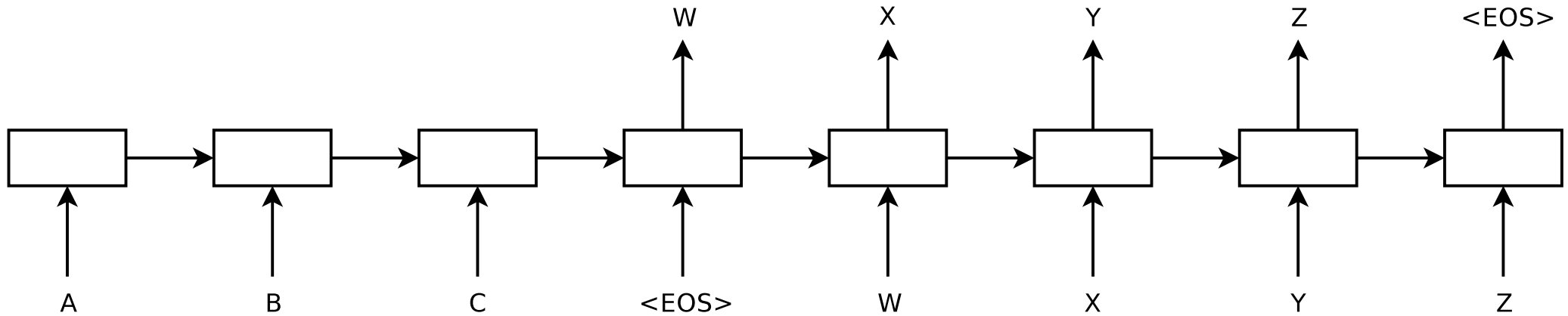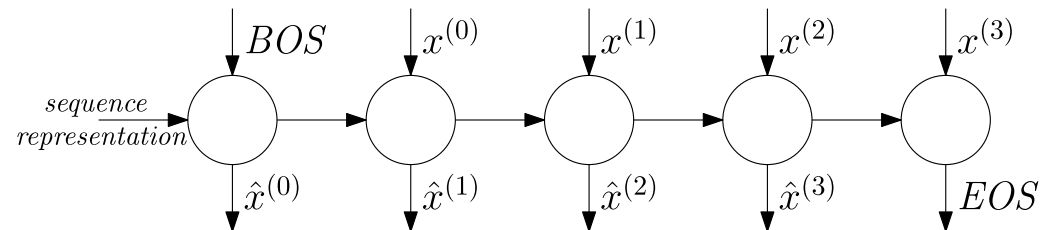*Figure 1 of "Sequence to Sequence Learning with Neural Networks", https://arxiv.org/abs/1409.0473*

Decoder



Figure 1 of "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation", https://arxiv.org/abs/1406.1078
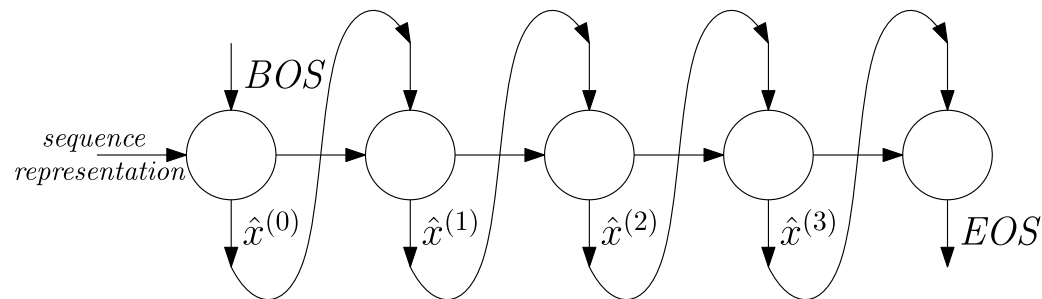
Encoder

## Training

The so-called **teacher forcing** is used during training – the gold outputs are used as inputs during training.



## Inference

During inference, the network processes its own predictions – such an approach is called **autoregressive decoding**.

Usually, the generated logits are processed by an $\arg\max$, the chosen word embedded and used as next input.
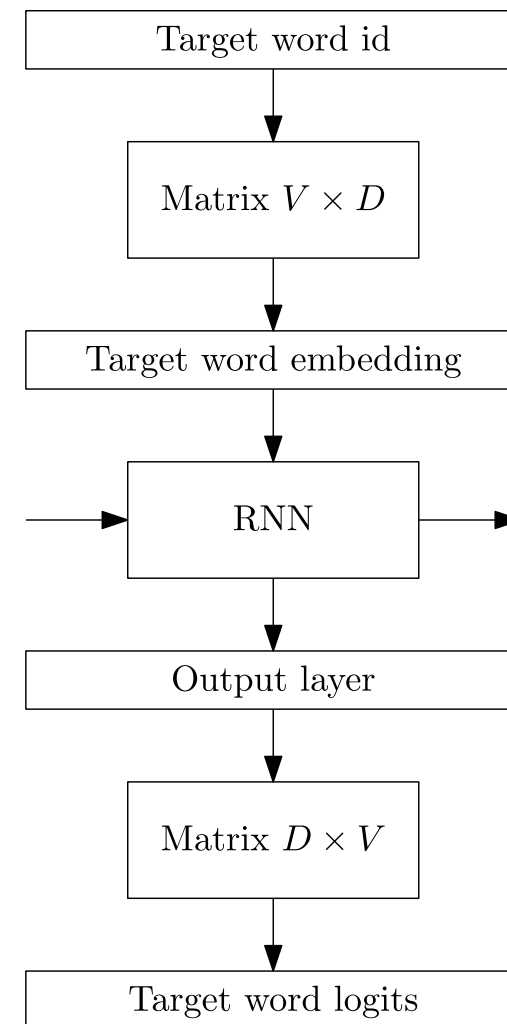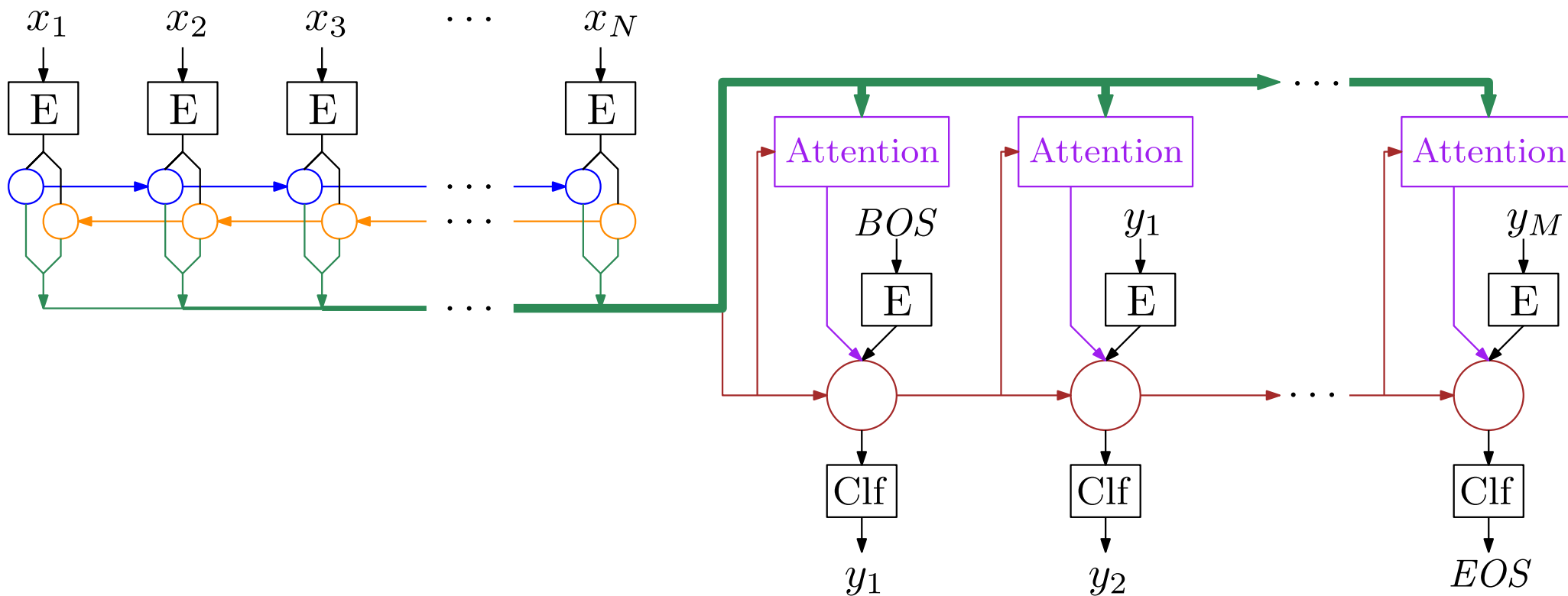
# Tying Word Embeddings

In the decoder, we both:

- embed the previous prediction, using a matrix of size $\mathbb{R}^{V \times D}$, where $V$ is the vocabulary size and $D$ is the embedding size;
- classify the hidden state into current prediction, using a matrix of size $\mathbb{R}^{D \times V}$.

Both these matrices have the same meaning – they represent the target-side words in the embedding space (the first explicitly represents the words by these embeddings, the second chooses the embedding in a sense "closest" to the produced hidden state).

Therefore, it makes sense to **tie** these matrices, i.e., to represent one of them as a transposition of the other.

```
Target word id
      ↓
Matrix V × D
      ↓
Target word embedding
      ↓
    RNN
      ↓
Output layer
      ↓
Matrix D × V
      ↓
Target word logits
```

# Attention

As another input during decoding, we add *context vector* $c_i$:

$$s_i = f(s_{i-1}, y_{i-1}, c_i).$$

We compute the context vector as a weighted combination of source sentence encoded outputs:

$$c_i = \sum_j \alpha_{ij} h_j$$

The weights $\alpha_{ij}$ are softmax of $e_{ij}$ over $j$,

$$\alpha_i = \mathrm{softmax}(e_i),$$

with $e_{ij}$ being

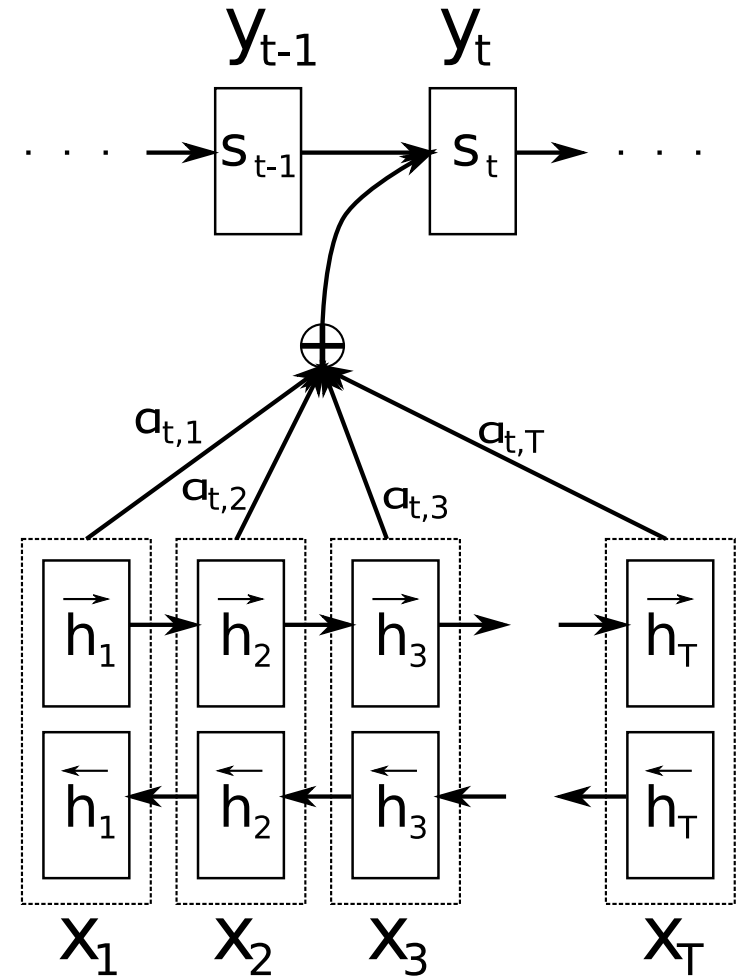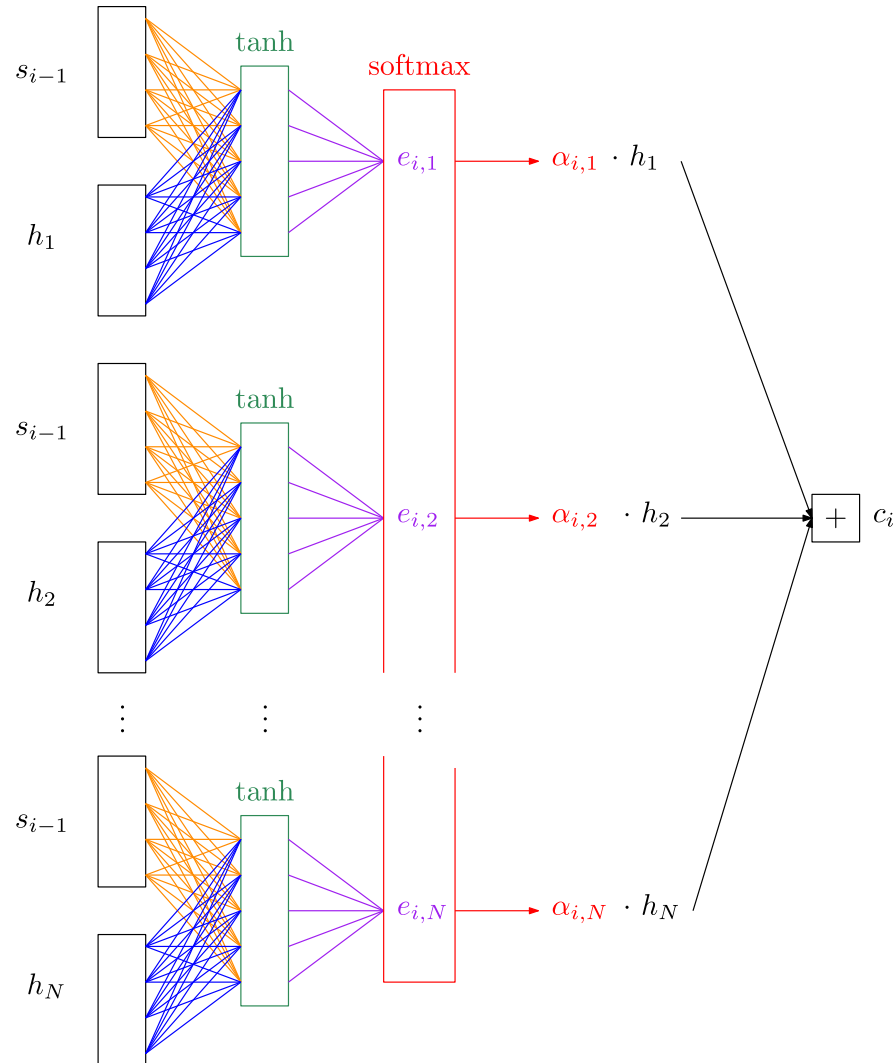$$e_{ij} = v^\top \tanh(V h_j + W s_{i-1} + b).$$



Figure 1 of "Neural Machine Translation by Jointly Learning to Align and Translate", https://arxiv.org/abs/1409.0473

(a)  (b)

(c)  (d)

*Figure 3 of "Neural Machine Translation by Jointly Learning to Align and Translate", https://arxiv.org/abs/1409.0473*

# Subword Units

Translate **subword units** instead of words. The subword units can be generated in several ways, the most commonly used are:

- **BPE**: Using the *byte pair encoding* algorithm. Start with individual characters plus a special end-of-word symbol $\cdot$. Then, merge the most occurring symbol pair $A, B$ by a new symbol $AB$, with the symbol pair never crossing word boundary (so that the end-of-word symbol cannot be inside a subword).

  Considering a dictionary with words *low, lowest, newer, wider*, a possible sequence of merges:

$$r \ \cdot \to r\cdot$$
$$l \ o \to lo$$
$$lo \ w \to low$$
$$e \ r\cdot \to er\cdot$$

- **Wordpieces**: Given a text divided into subwords, we can compute unigram probability of every subword, and then get the likelihood of the text under a unigram language model by multiplying the probabilities of the subwords in the text.

  When we have only a text and a subword dictionary, we divide the text in a greedy fashion, iteratively choosing the longest existing subword.

  When constructing the subwords, we again start with individual characters, and then repeatedly join such a pair of subwords that increases the unigram language model likelihood the most.

Both approaches give very similar results; the biggest difference is that during the inference:

- for BPE, the sequence of merges must be performed in the same order as during the construction of the BPE (because we use the output of BPE as training data),
- for Wordpieces, it is enough to find longest matches from the subword dictionary.

Usually quite little subword units are used (32k-64k), often generated on the union of the two vocabularies (the so-called *joint BPE* or *shared wordpieces*).
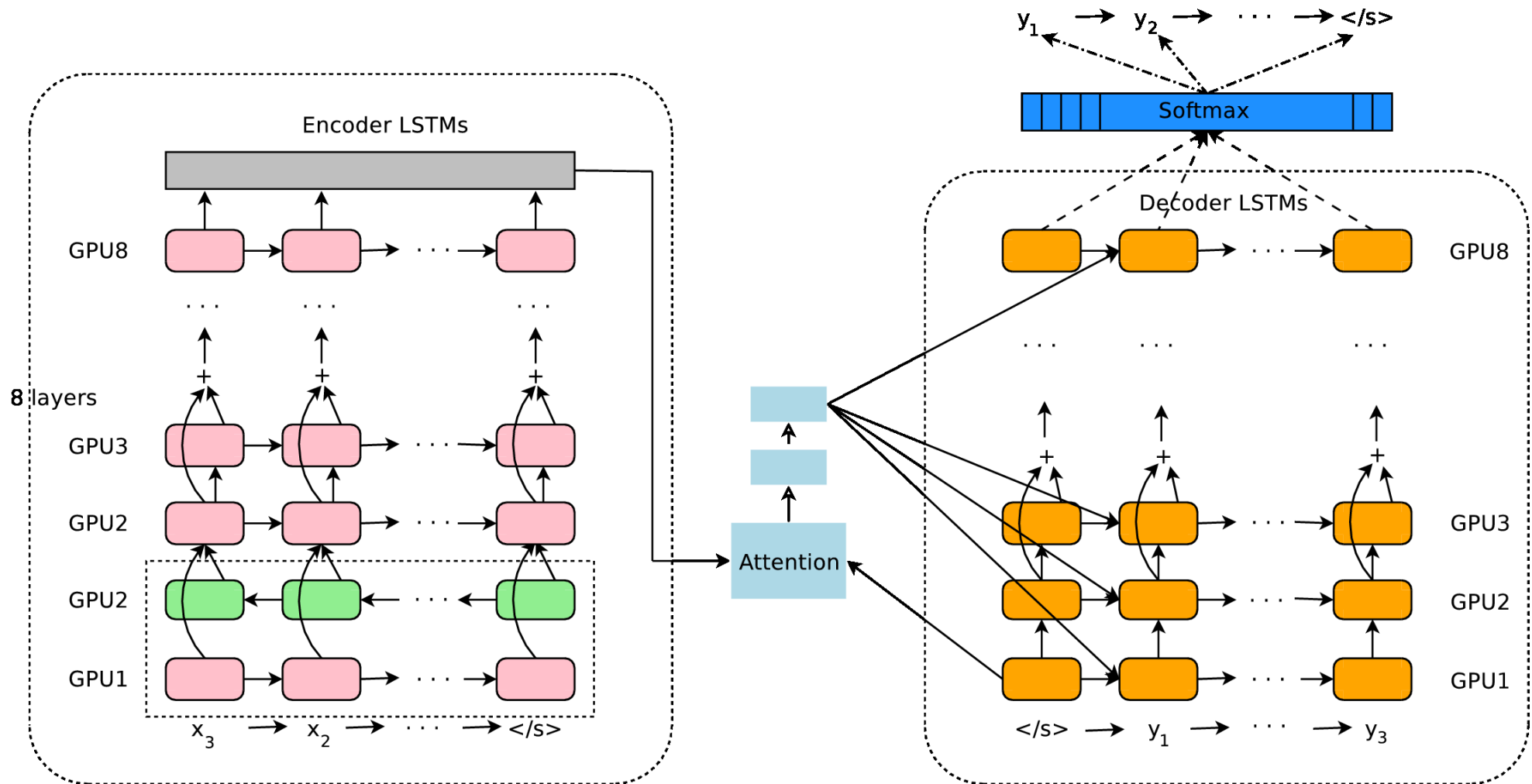
Figure 1 of "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation", https://arxiv.org/abs/1609.08144
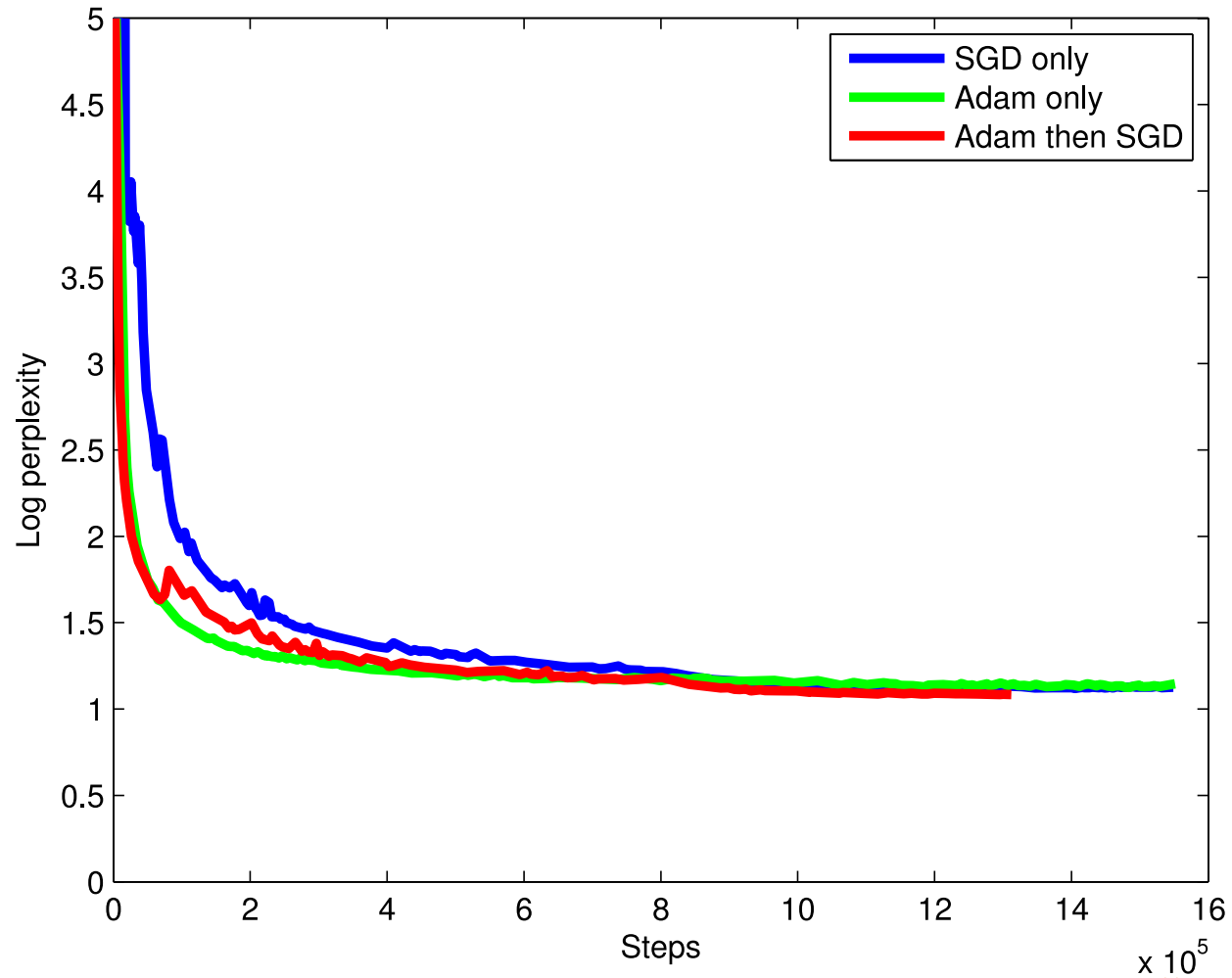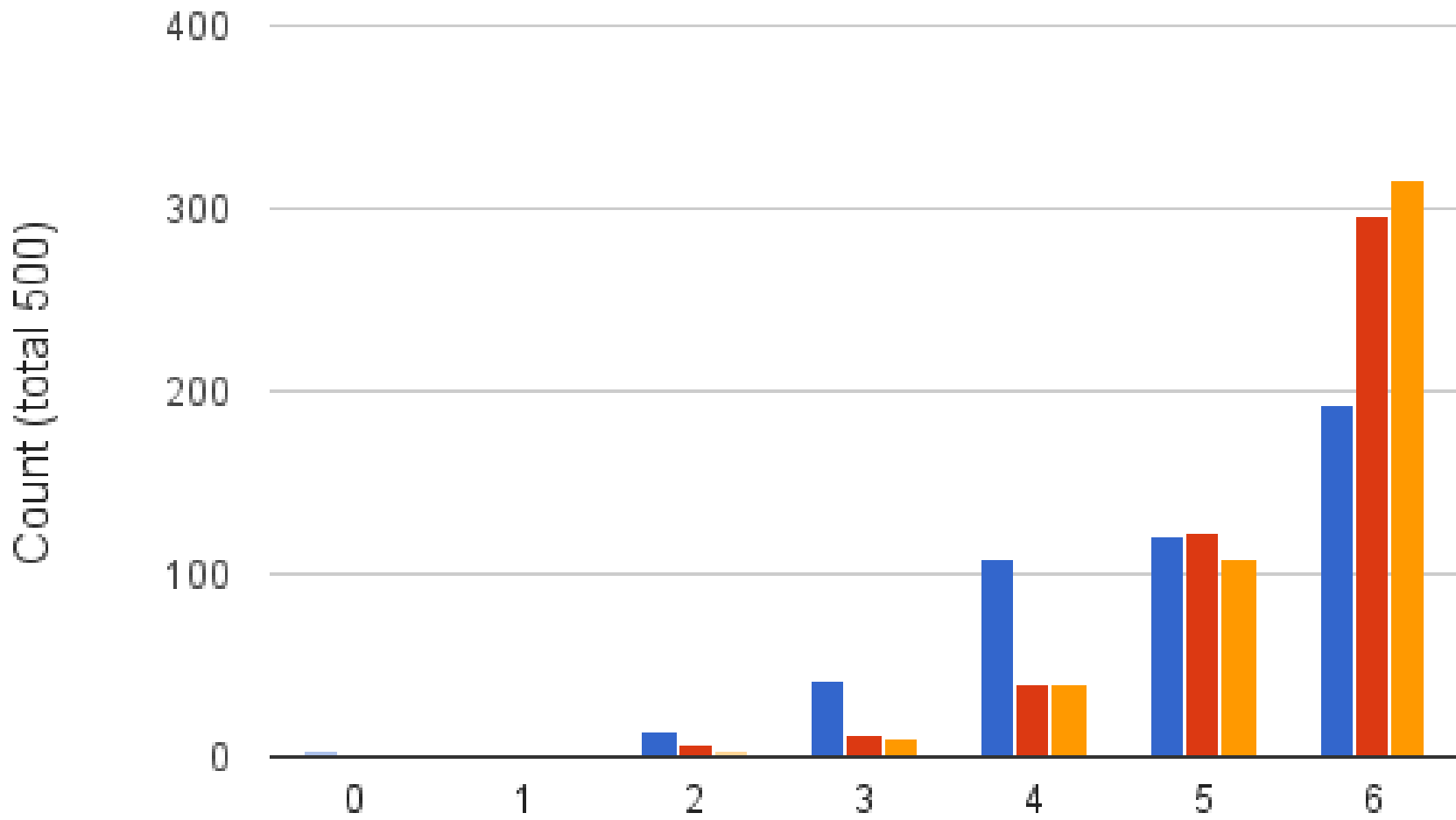
Figure 5 of "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation", https://arxiv.org/abs/1609.08144

Figure 6 of "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation", https://arxiv.org/abs/1609.08144
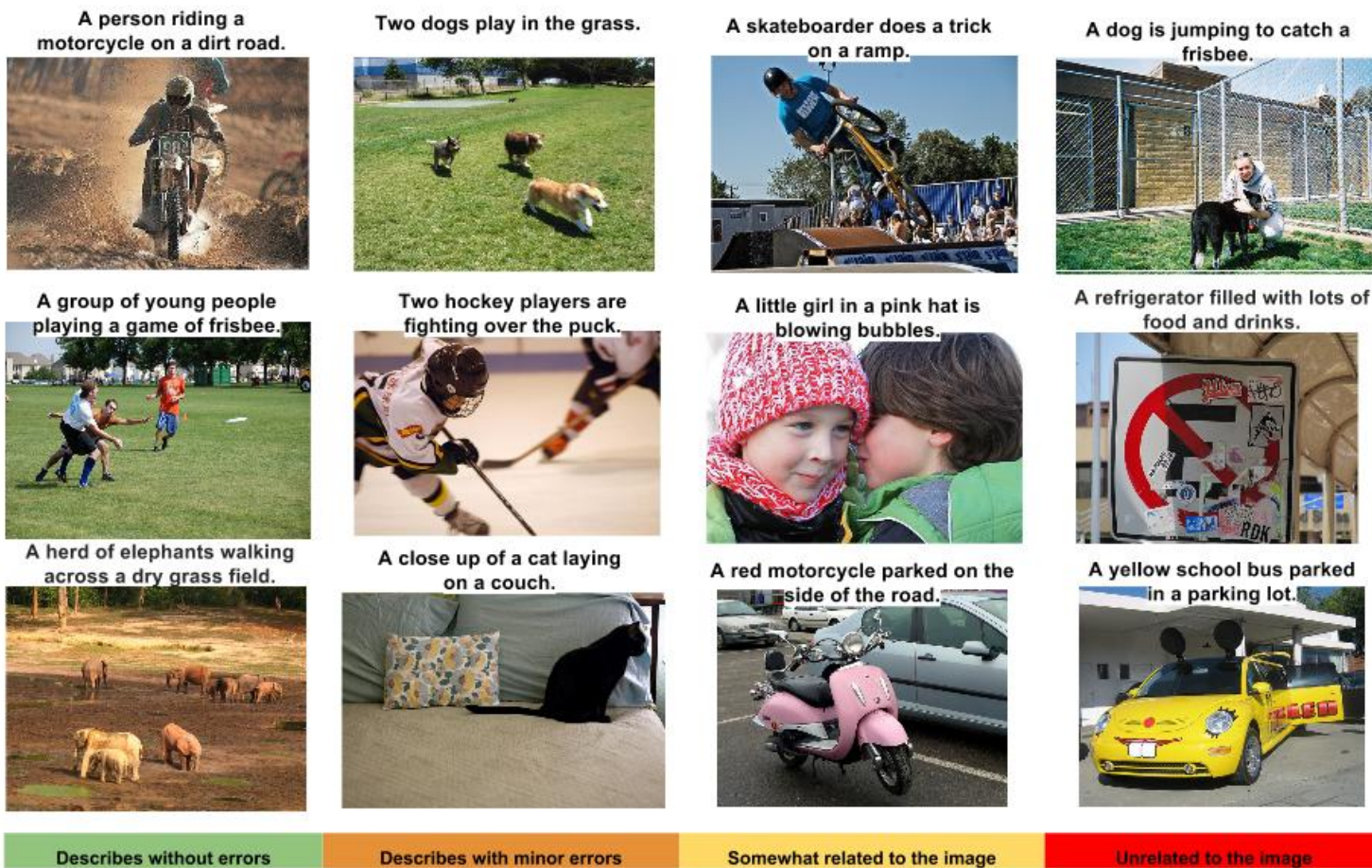
Fig. 5. A selection of evaluation results, grouped by human rating.

*Figure 5 of "Show and Tell: Lessons learned from the 2015 MSCOCO...", https://arxiv.org/abs/1609.06647*

What vegetable is the dog chewing on?
MCB: carrot
GT: carrot

What kind of dog is this?
MCB: husky
GT: husky

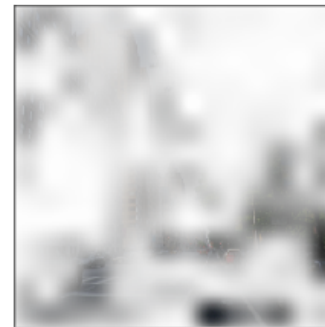What kind of flooring does the room have?
MCB: carpet
GT: carpet

What color is the traffic light?
MCB: green
GT: green

Is this an urban area?
MCB: yes
GT: yes

Where are the buildings?
MCB: in background
GT: on left

Figure 6 of "Multimodal Compact Bilinear Pooling for VQA and Visual Grounding", https://arxiv.org/abs/1606.01847

# Multilingual and Unsupervised Translation

Many attempts at multilingual translation.

- Individual encoders and decoders, shared attention.
- Shared encoders and decoders.

Surprisingly, even unsupervised translation is attempted lately. By unsupervised we understand settings where we have access to large monolingual corpora, but no parallel data.

In 2019, the best unsupervised systems were on par with the best 2014 supervised systems.

|  |  | WMT-14 | | | |
|---|---|---|---|---|---|
|  |  | fr-en | en-fr | de-en | en-de |
| Unsupervised | Proposed system | 33.5 | 36.2 | 27.0 | 22.5 |
|  | *detok. SacreBLEU*[*] | 33.2 | 33.6 | 26.4 | 21.2 |
| Supervised | WMT best[*] | 35.0 | 35.8 | 29.0 | 20.6[†] |
|  | Vaswani et al. (2017) | - | 41.0 | - | 28.4 |
|  | Edunov et al. (2018) | - | 45.6 | - | 35.0 |

Table 3: Results of the proposed method in comparison to different supervised systems (BLEU).

Table 3 of "An Effective Approach to Unsupervised Machine Translation", https://arxiv.org/abs/1902.01313

For some sequence processing tasks, *sequential* processing (as performed by recurrent neural networks) of its elements might be too restrictive.

Instead, we may want to be able to combine sequence elements independently on their distance.

Such processing is allowed in the **Transformer** architecture, originally proposed for neural machine translation in 2017 in *Attention is All You Need* paper.
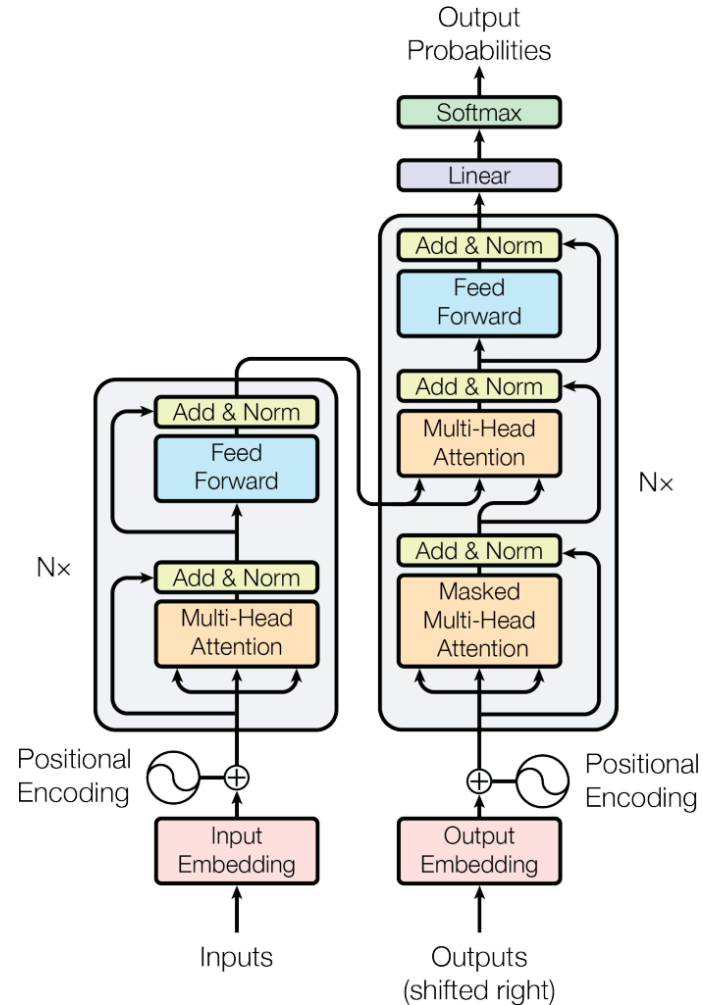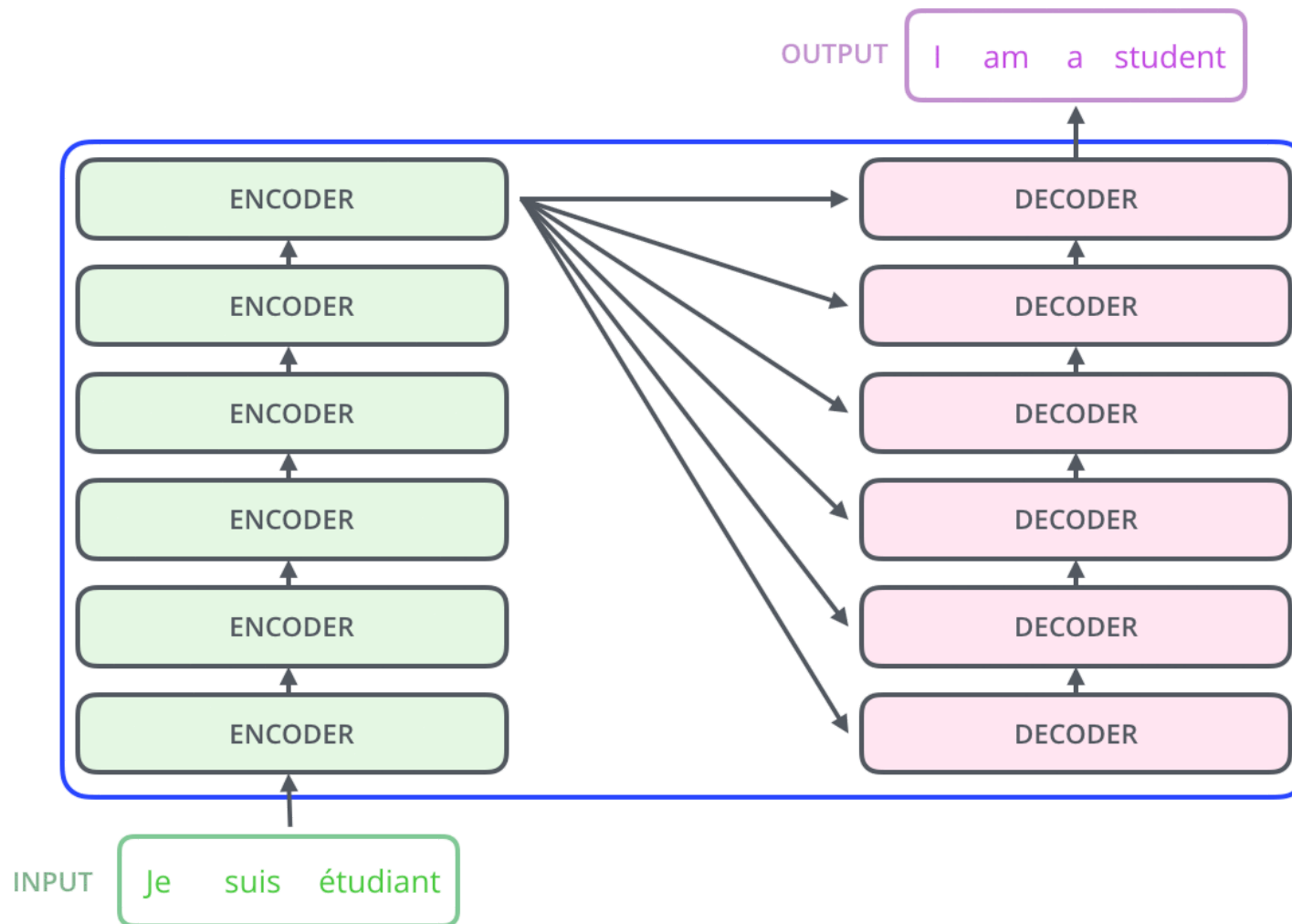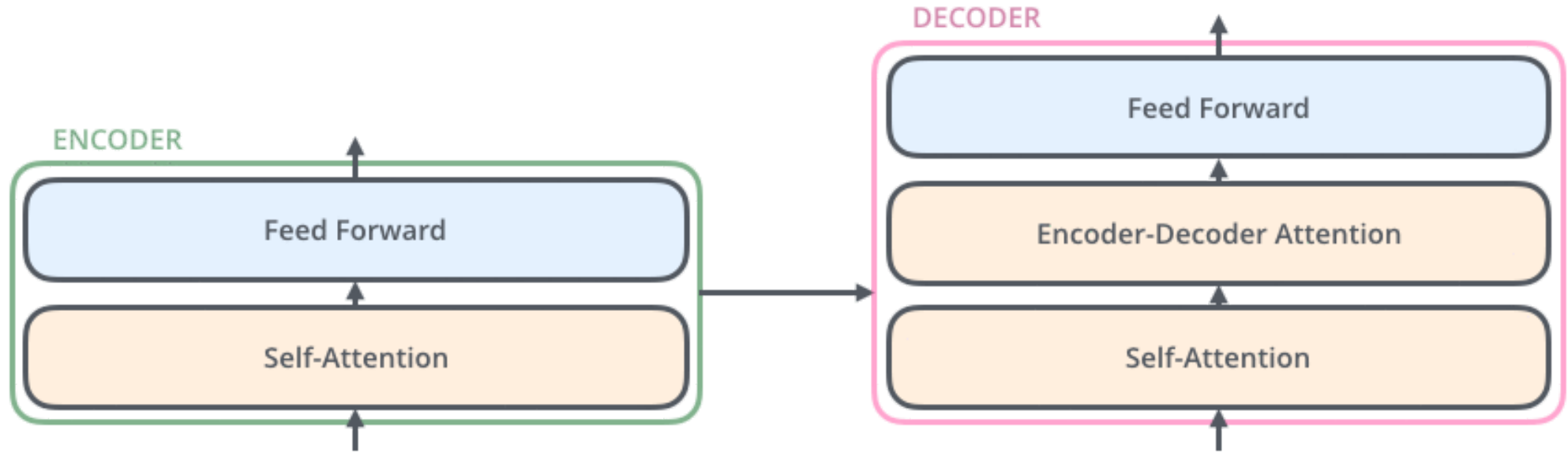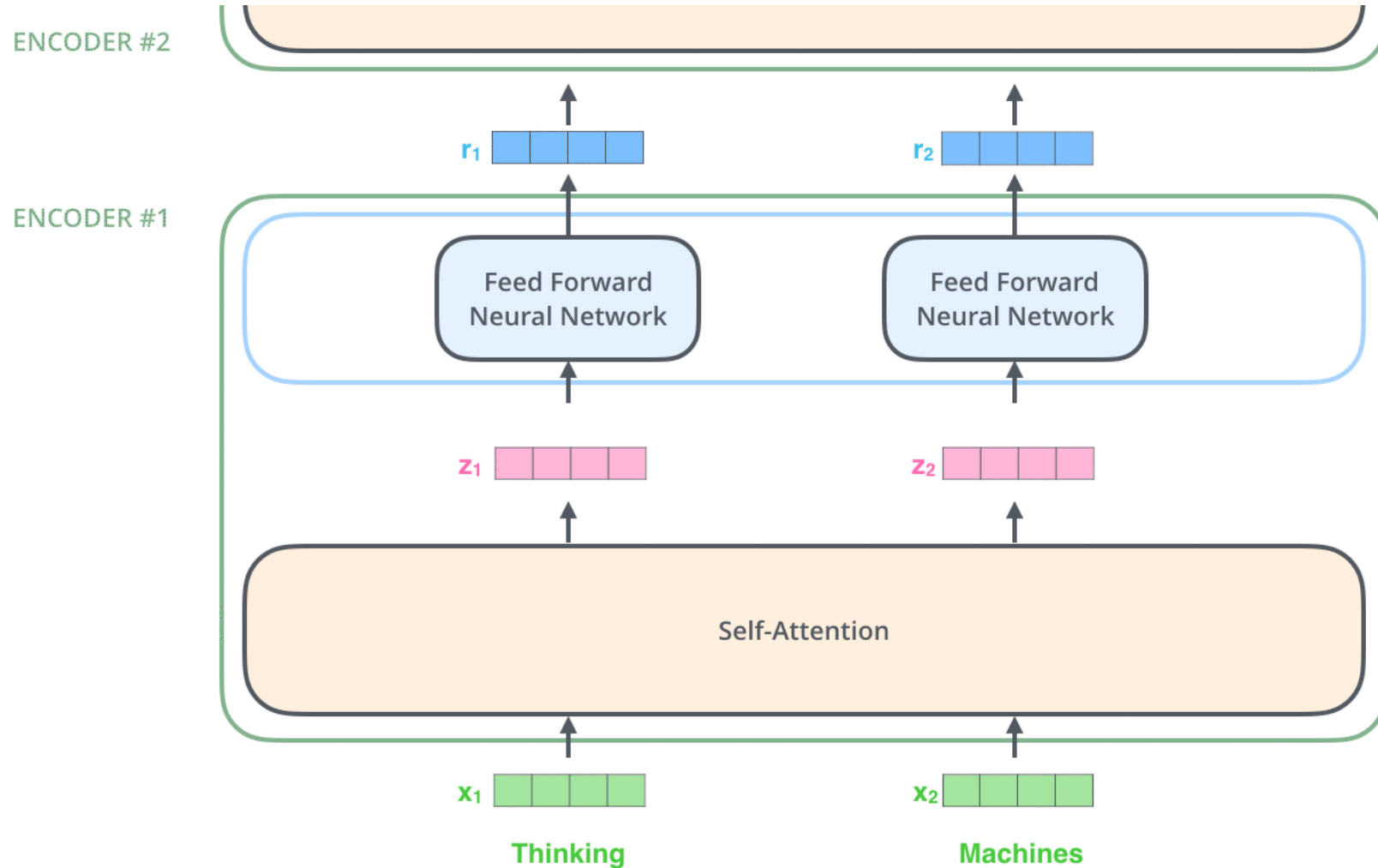
*Figure 1 of "Attention Is All You Need", https://arxiv.org/abs/1706.03762*

# Transformer

http://jalammar.github.io/images/t/The_transformer_encoder_decoder_stack.png

http://jalammar.github.io/images/t/encoder_with_tensors_2.png

# Transformer − Self-Attention

Assume that we have a sequence of $n$ words represented using a matrix $\boldsymbol{X} \in \mathbb{R}^{n \times d}$.

The attention module for queries $\boldsymbol{Q} \in \mathbb{R}^{n \times d_k}$, keys $\boldsymbol{K} \in \mathbb{R}^{n \times d_k}$ and values $\boldsymbol{V} \in \mathbb{R}^{n \times d_v}$ is defined as:

$$\mathrm{Attention}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \mathrm{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^\top}{\sqrt{d_k}}\right)\boldsymbol{V}.$$

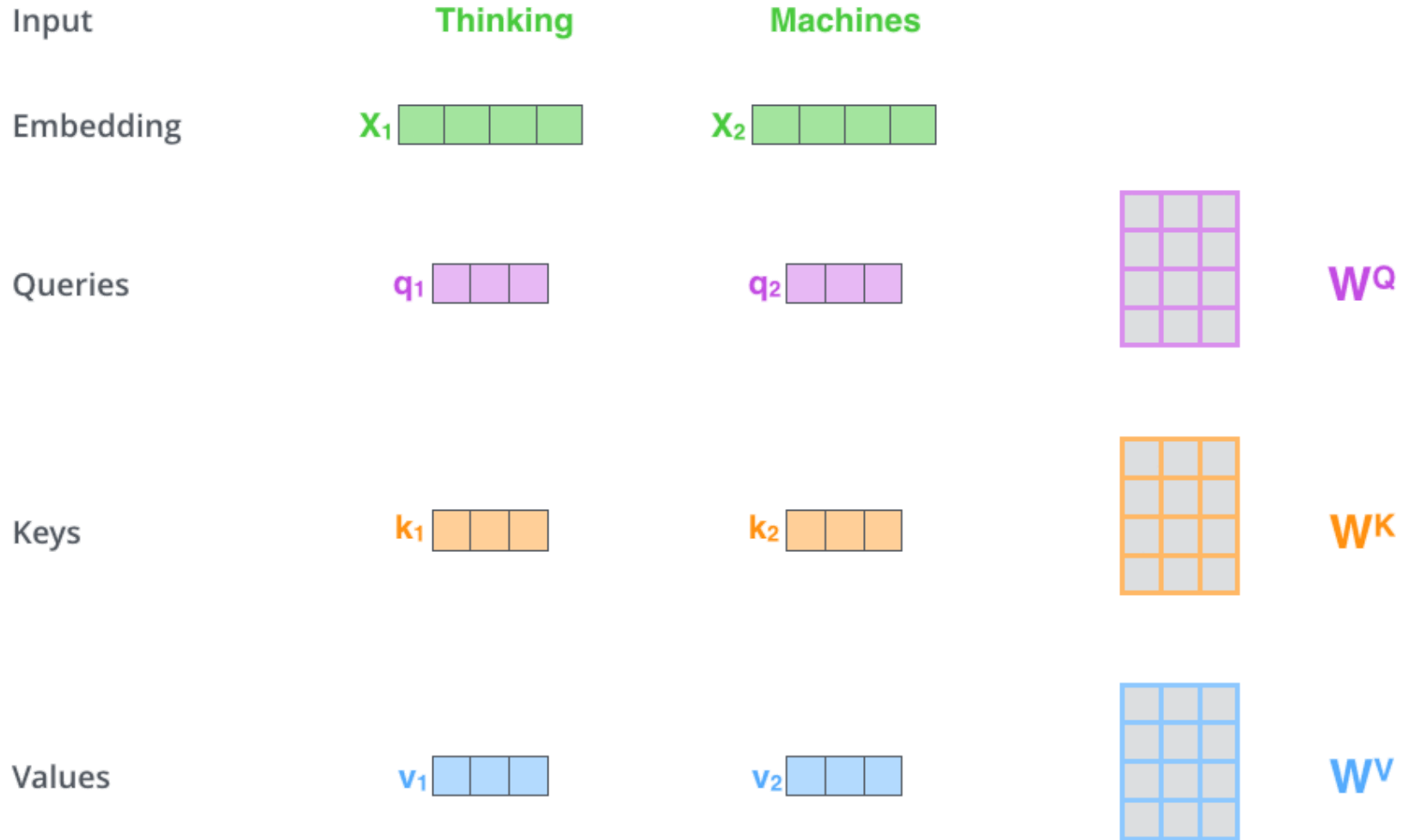The queries, keys and values are computed from the input word representations $\boldsymbol{X}$ using a linear transformation as

$$\boldsymbol{Q} = \boldsymbol{X}\boldsymbol{W}^Q$$
$$\boldsymbol{K} = \boldsymbol{X}\boldsymbol{W}^K$$
$$\boldsymbol{V} = \boldsymbol{X}\boldsymbol{W}^V$$

for trainable weight matrices $\boldsymbol{W}^Q, \boldsymbol{W}^K \in \mathbb{R}^{d \times d_k}$ and $\boldsymbol{W}^V \in \mathbb{R}^{d \times d_v}$.

https://miro.medium.com/max/2000/1*jBsfVNOOcJ-I3tsLVgni_w.png

http://jalammar.github.io/images/t/self-attention-matrix-calculation-2.png

http://jalammar.github.io/images/t/self-attention-matrix-calculation.png