

Seq2Seq Model

Technische Universität München

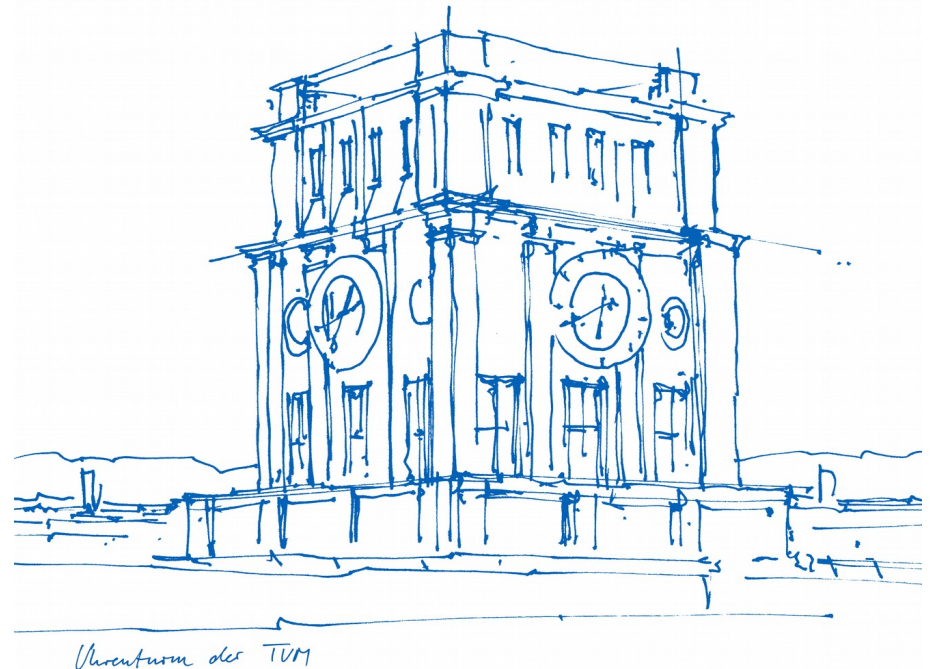
Department of Informatics

Seminar - Applied Deep Learning for NLP

Prof. Dr. Simon Hegelich

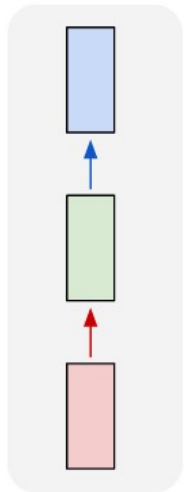
Francesco Cognolato

10th January 2019

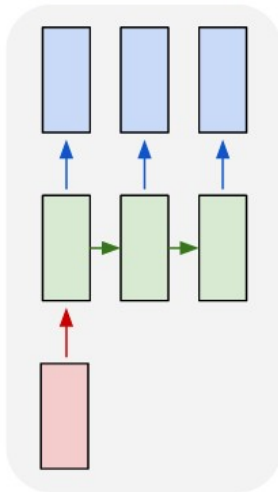


Different tasks require different models

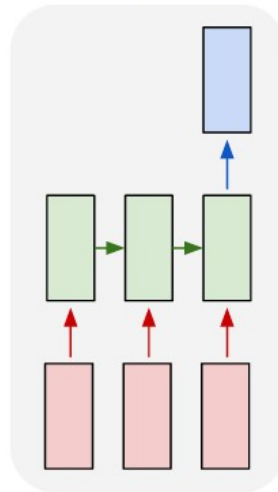
one to one



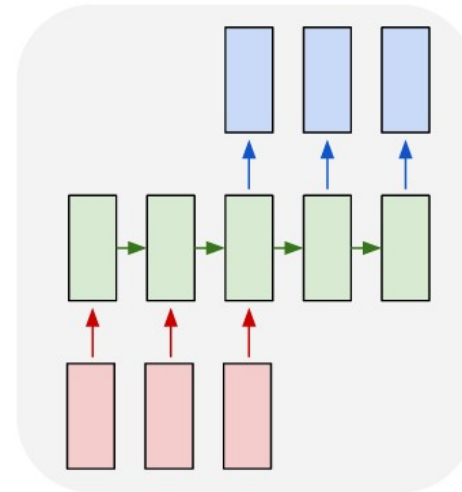
one to many



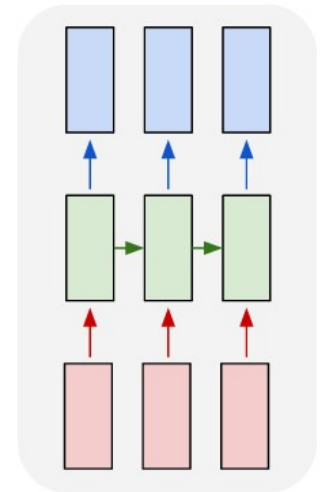
many to one



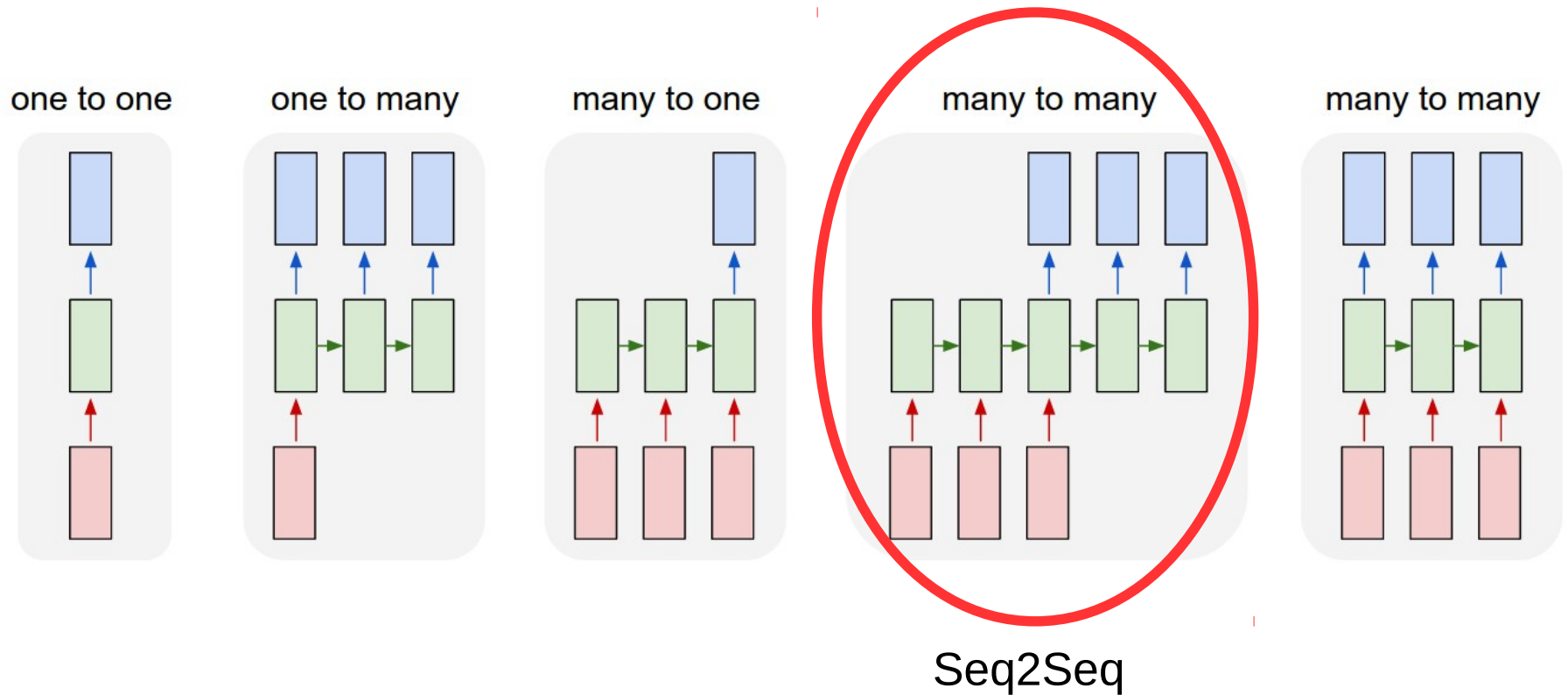
many to many



many to many



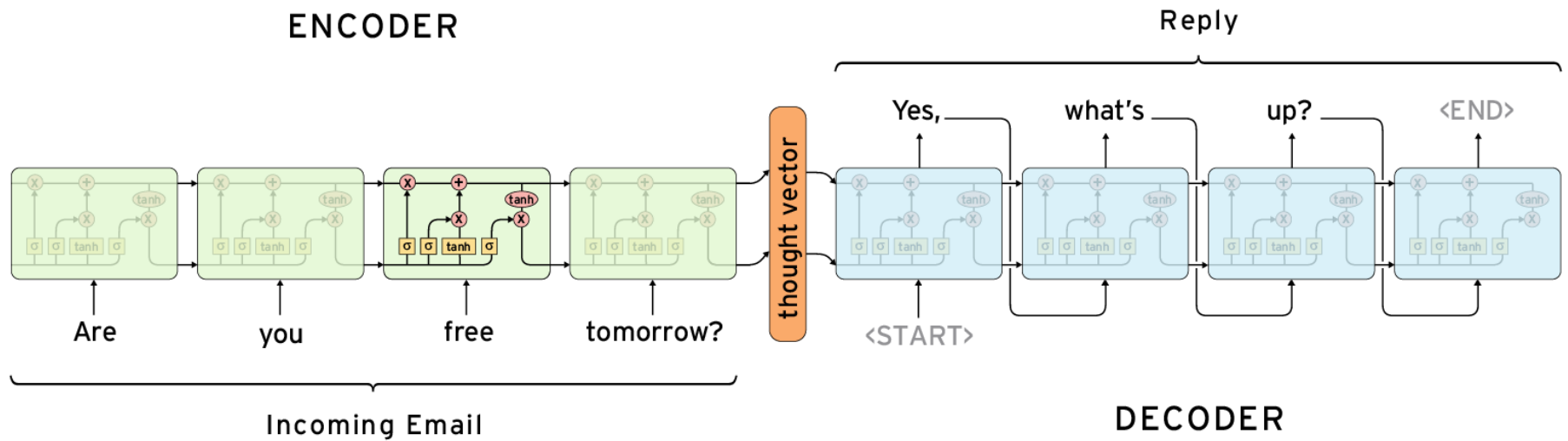
Different tasks require different models



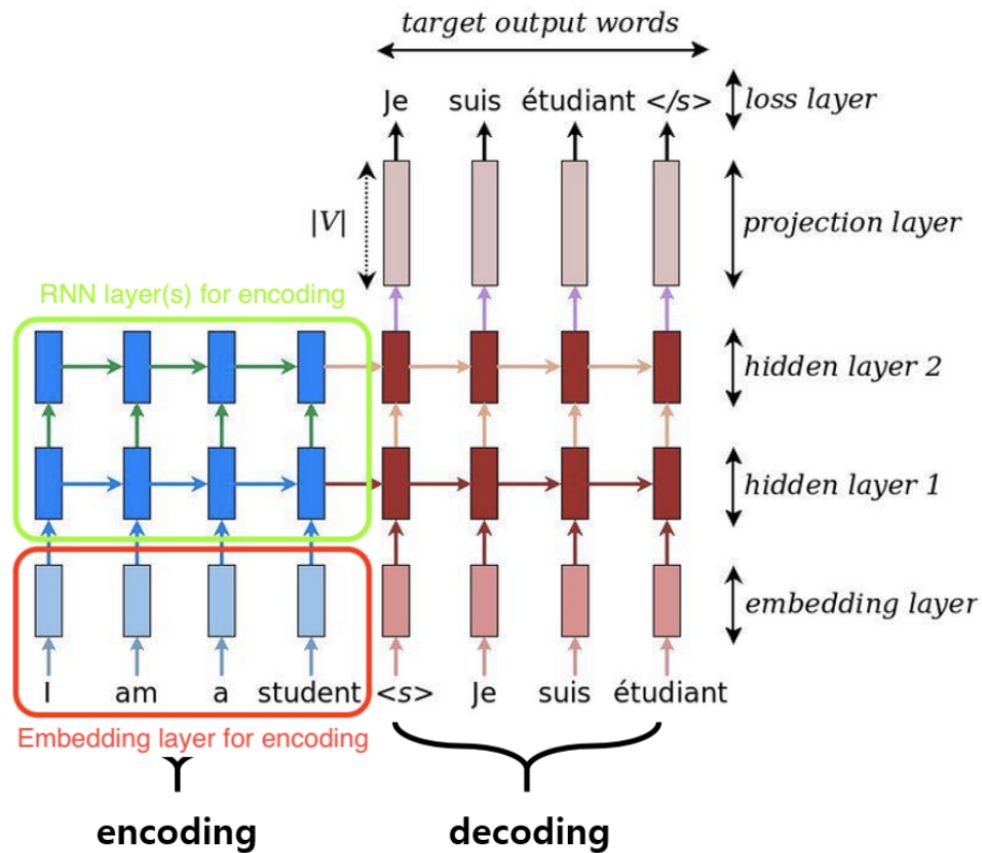
Some possible applications

- Neural Machine Translation : input is a sentence (sequence) in one language, output is one in another language
- Text summarization : input is a long text (sequence), output is summary of the text
- Chatbot : input is a sentence, output is a reply
- Speech recognition : input is sequence (digital waveform), output is text

Seq2Seq



Seq2Seq - Basic architecture

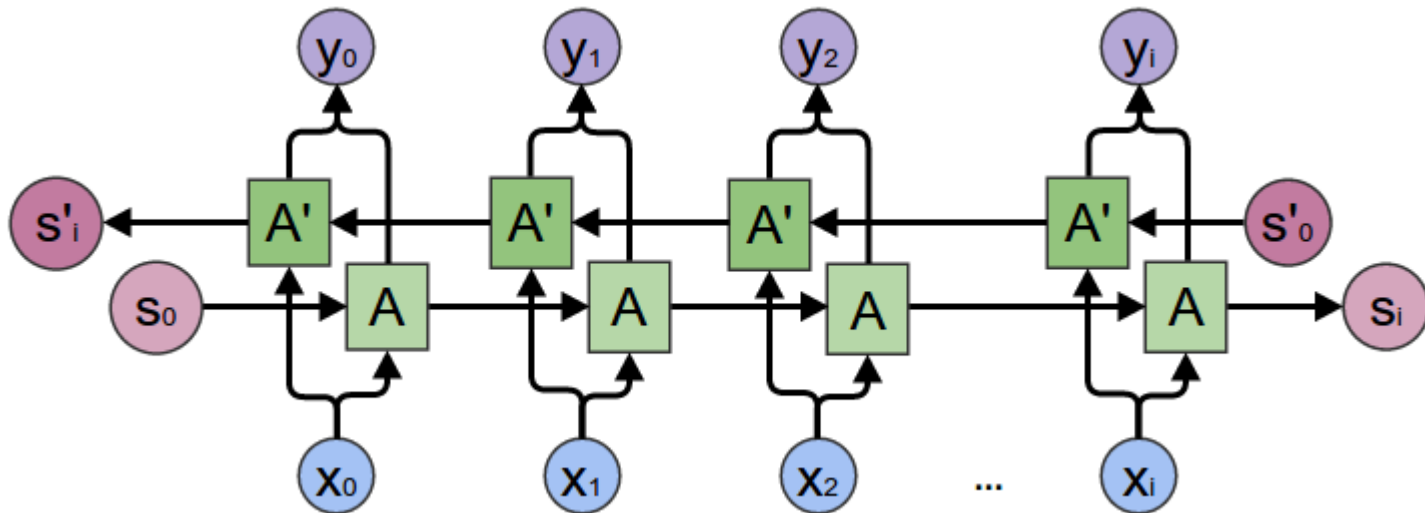
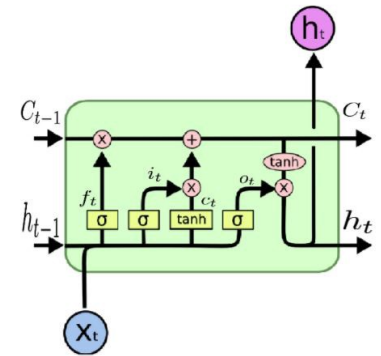


Embedding layer

- Nothing more than a matrix, $N \times F$ (N words, F embedding size)
- Both source and target language have their own embedding layer
- Unique mapping between a word and an integer
- Look up the word vector given the corresponding row index in the embedding matrix
- If enough data train embeddings from scratch, otherwise use GloVe pretrained embeddings provided from Stanford University

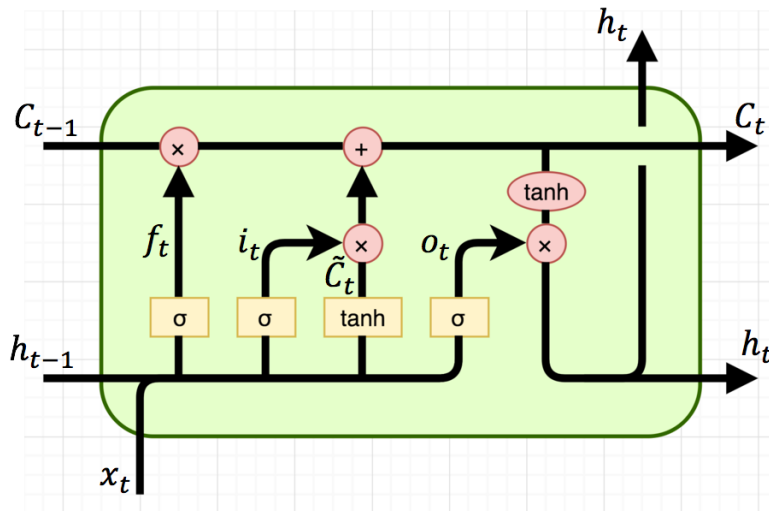
Encoder architecture

- (Deep) Bidirectional LSTM or GRU

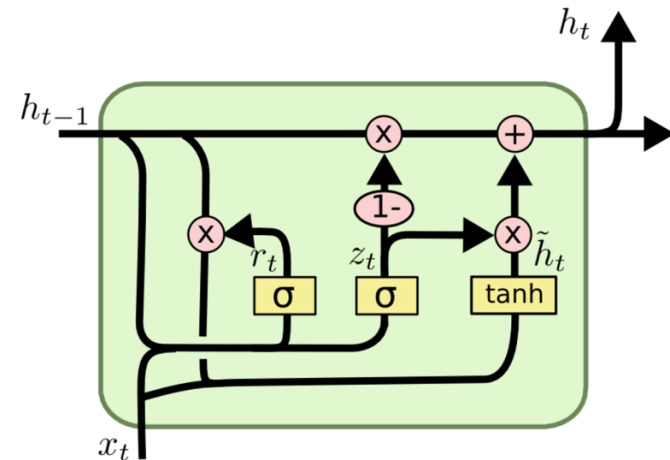


Decoder architecture

- (Deep) LSTM or GRU, initial state is set from last state of encoder



(a) Long Short-Term Memory



(b) Gated Recurrent Unit

Attention mechanism

- In the basic seq2seq architecture, we ignore encoder outputs, but just consider context vector (last state of encoder)

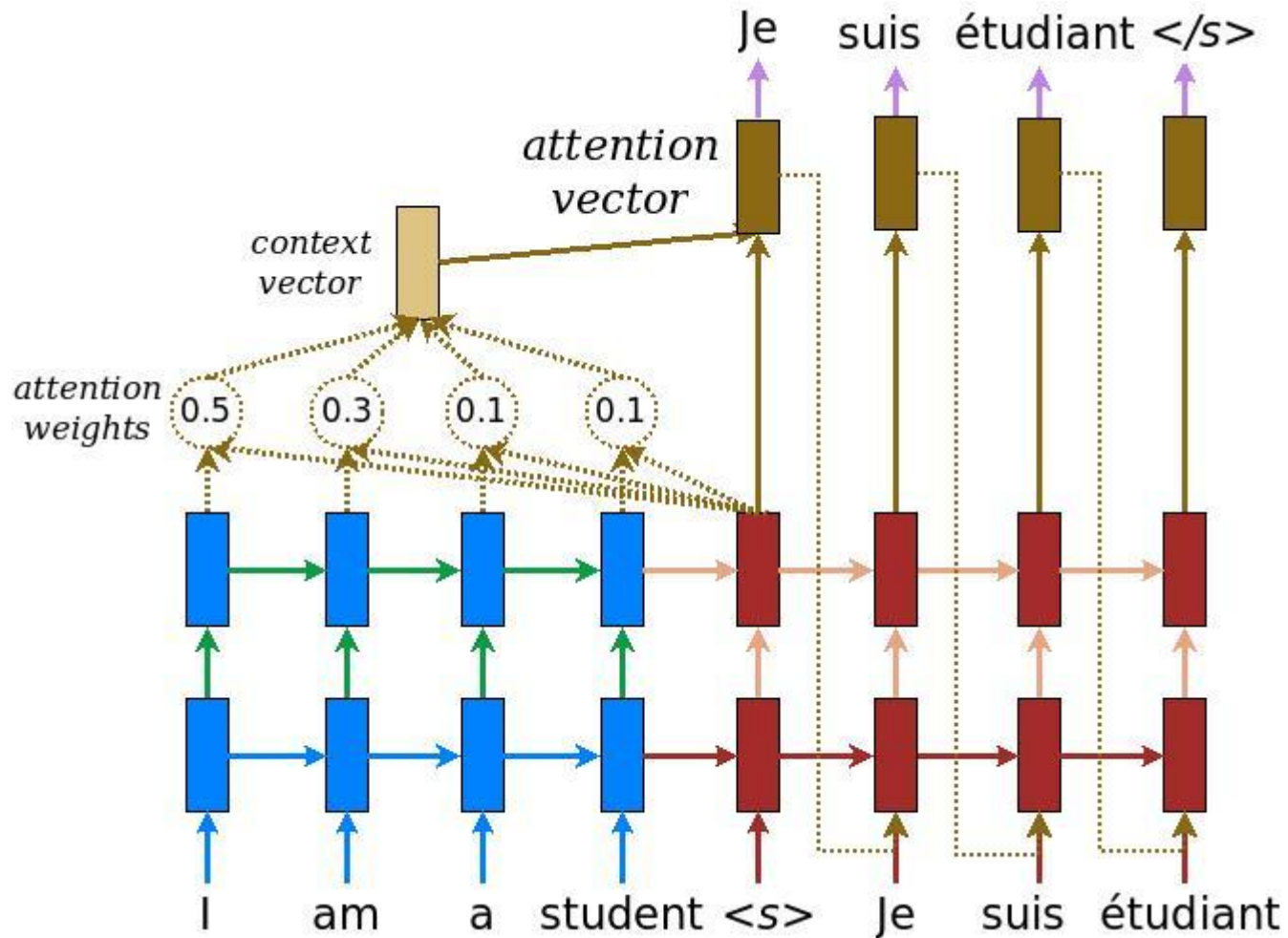
$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \rightarrow \text{Luong score} \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a [\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat}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

$$\alpha_{ts} = \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'=1}^S \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad [\text{Attention weights}]$$

$$\mathbf{c}_t = \sum_s \alpha_{ts} \bar{\mathbf{h}}_s \quad [\text{Context vector}]$$

$$\mathbf{a}_t = f(\mathbf{c}_t, \mathbf{h}_t) = \tanh(\mathbf{W}_c [\mathbf{c}_t; \mathbf{h}_t]) \quad [\text{Attention vector}]$$

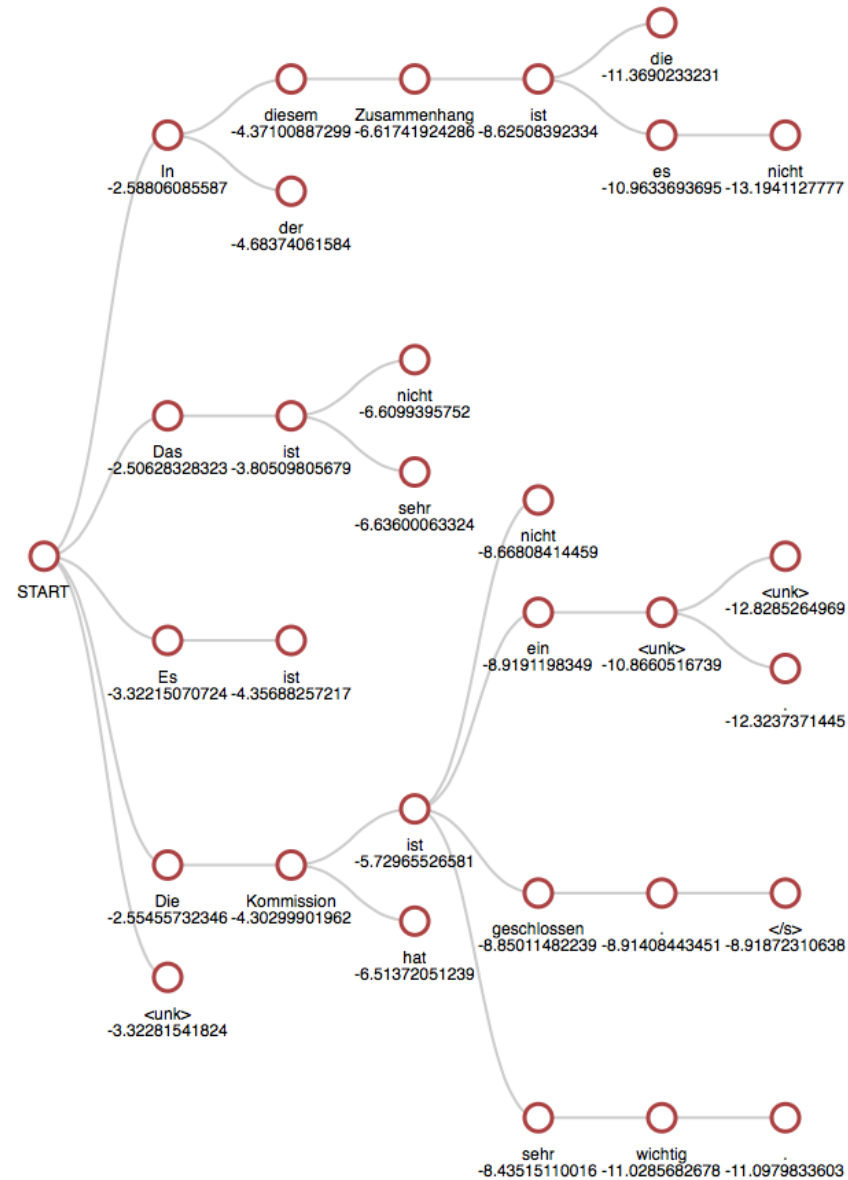
Attention mechanism



Remarks

- Add padding of zeros to all the sentences to have the same length (otherwise use some fancy things, like bucketing)
- After decoder, softmax over whole vocabulary and calculate cross entropy loss
- Testing time : source language sequence, decoder input will be only <START> token, generate prediction, append to <START>, generate next token and so on... Stop when <END> token generated (Greedy search)
- Use Beam search, pick best n (e.g. 5-10 beams) predictions and generate a tree of sentences. Final sentence is the one with highest final probabilities

Beam search



References

- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- <http://web.stanford.edu/class/cs224n/>

Questions?

Thank you