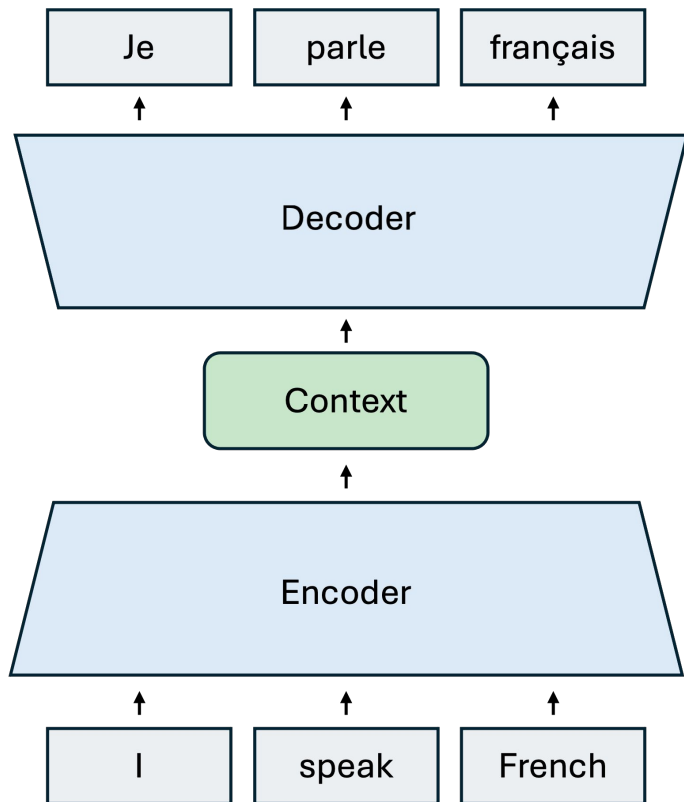


# Encoder-Decoder MT

# High Level



# Preparing the Data

# Tokenization

- We need to convert text  $\rightarrow$  vectors
- What is the right granularity?
  - Words?
  - Characters?

**Sentence:** The quick brown fox jumps over the lazy dog

**Word Tokenization:** ['The', 'quick', 'brown', 'fox', 'jumps', 'over', 'the', 'lazy', 'dog']

**Character Tokenization:** ['T', 'h', 'e', ' ', 'q', 'u', 'i', 'c', 'k', ' ', 'b', 'r', 'o', 'w', 'n', ' ', 'f', 'o', 'x', ' ', 'j', 'u', 'm', 'p', 's', ' ', 'o', 'v', 'e', 'r', ' ', 't', 'h', 'e', ' ', 'l', 'a', 'z', 'y', ' ', 'd', 'o', 'g']

# Subword Tokenization

- Using **subwords** solves a lot of the previous problems
- How do we get our subwords?

bert-base-cased

The handyman redecorated the farmhouse using sandpaper and paintbrushes|

TOKENS	CHARACTERS
18	71

[CLS] The handyman redecorated the farmhouse using sandpaper and paintbrushes [SEP]

# BPE

- Byte-Pair Encoding (BPE)
- Text = Sequence of bytes
- Bytes which co-occur together more frequently get merged to form tokens
- Merging stops when we reach the desired vocabulary size

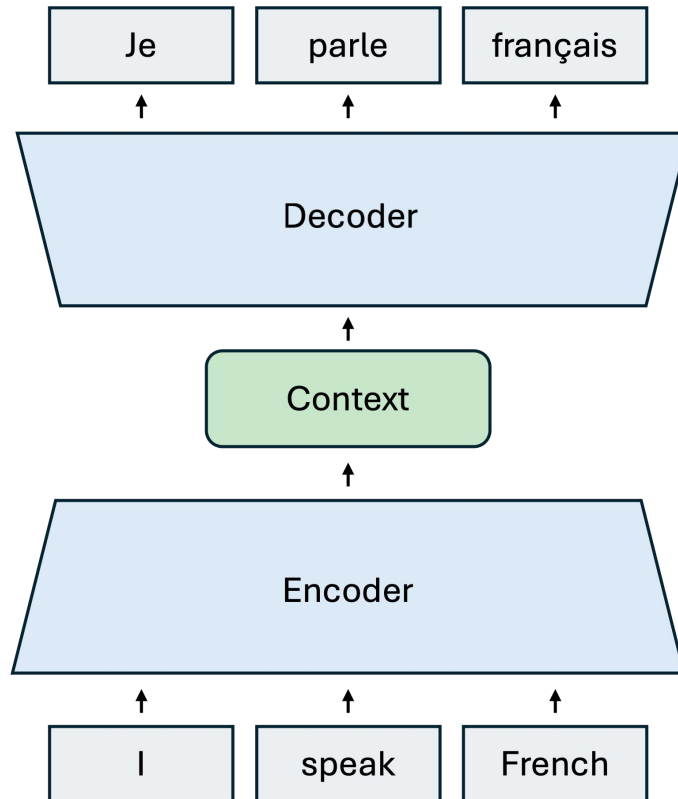
# Padding

- Now, we have a mapping Text  $\rightarrow$  Vectors
- We want to perform large amounts of computation in parallel  $\rightarrow$  We form **batches** of input text, or a large tensor with shape (BATCH\_SIZE, SEQ\_LEN, DIM)
- What if sequences are not the same length?
- **Pad** to the length of the max-length sequence in the batch

# Model Architecture

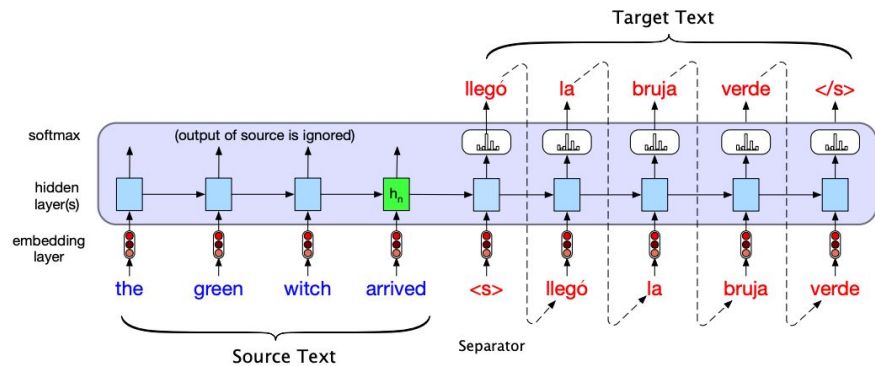


# High Level



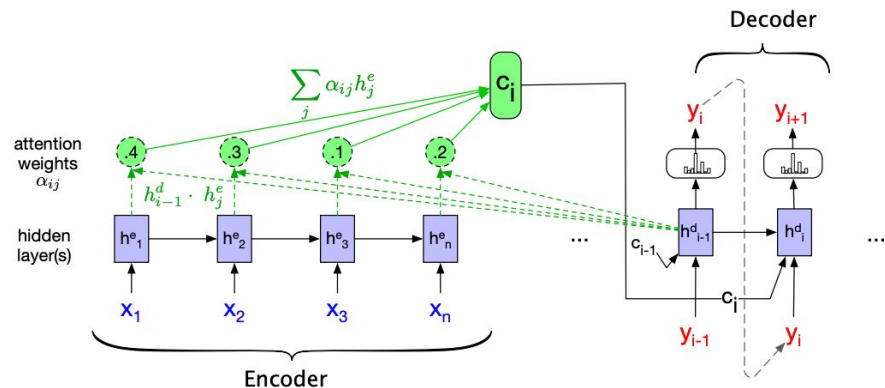
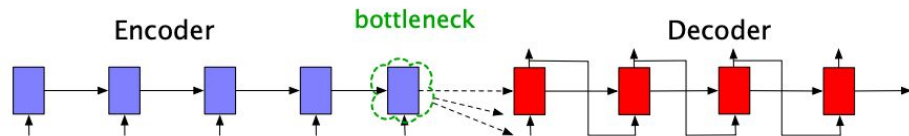
# Encoder-Decoder with RNNs

- The **encoder** RNNs produces a final hidden state after processing the source text
- The **decoder** RNN takes that hidden state, and begins decoding the target text

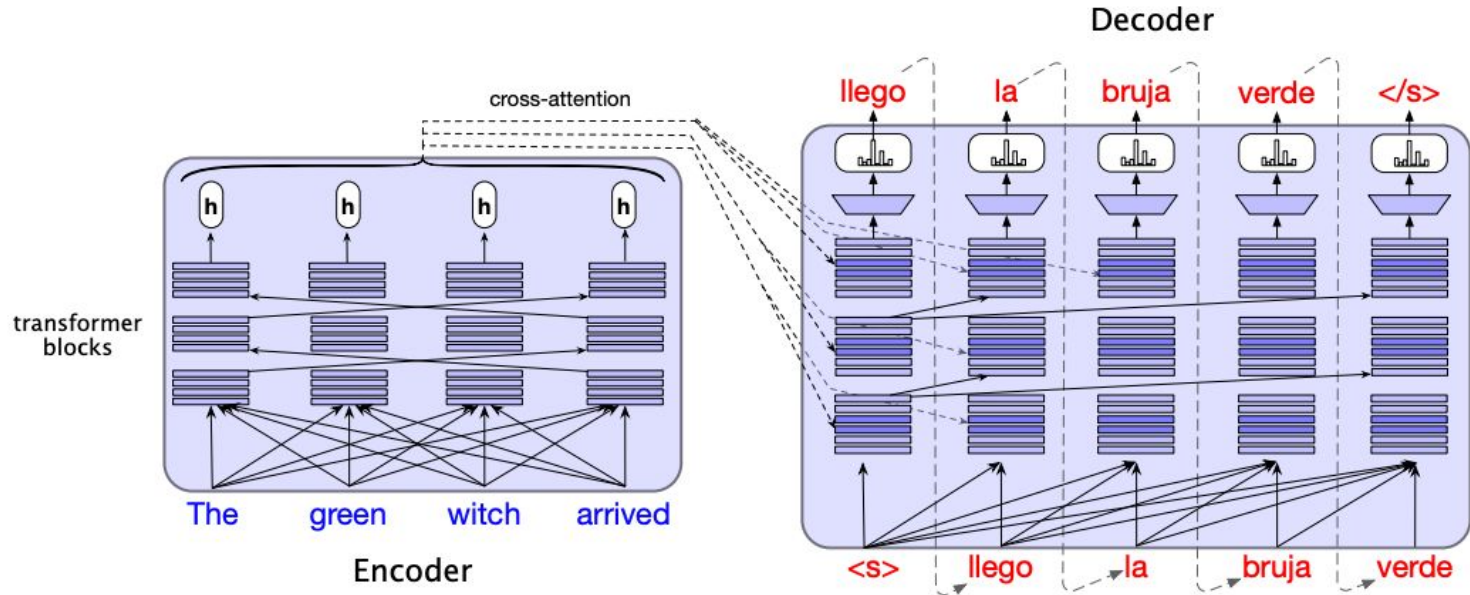


# Attention

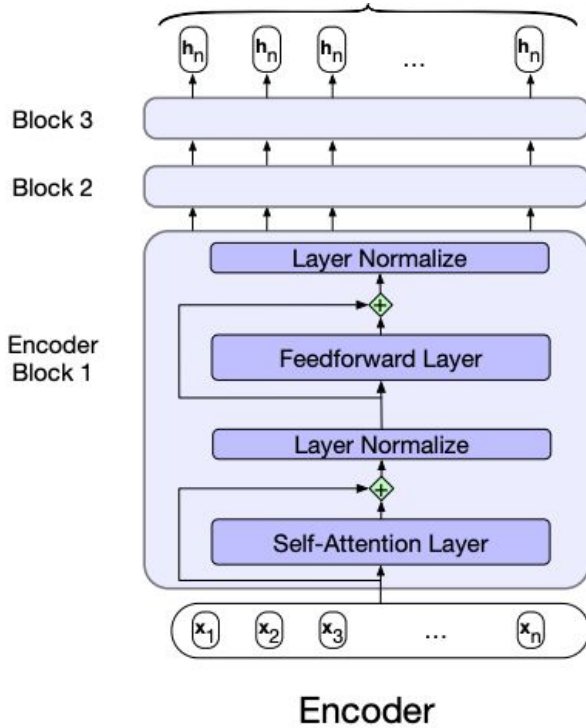
- If the decoder only uses the last hidden state, there is a bottleneck which can harm decoding
- Solution: **Attention**
- At each decoding step, we produce a context vector  $c_t$ , which is a function of all the encoder hidden states



# Transformer Encoder–Decoder

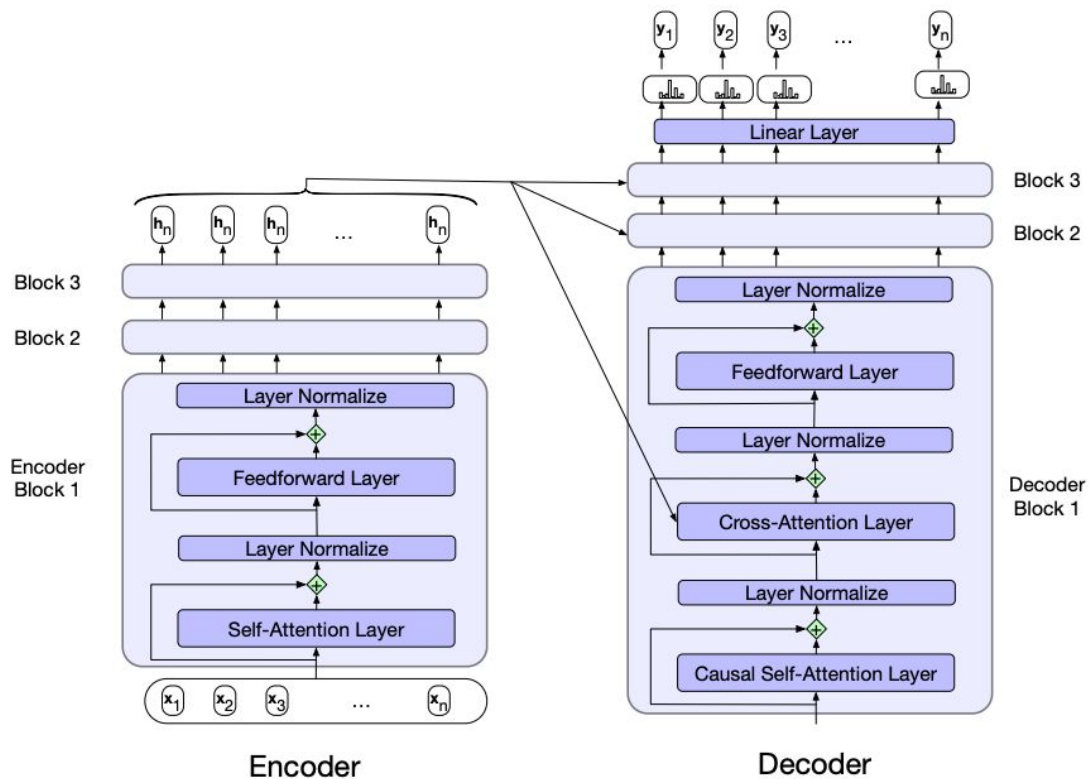


# Transformer Encoder



$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \text{3x3 grid} \end{matrix} \times \begin{matrix} \text{K}^T \\ \text{3x3 grid} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \text{3x3 grid} \end{matrix}$$
$$= \begin{matrix} \text{Z} \\ \text{3x3 grid} \end{matrix}$$

# Decoder



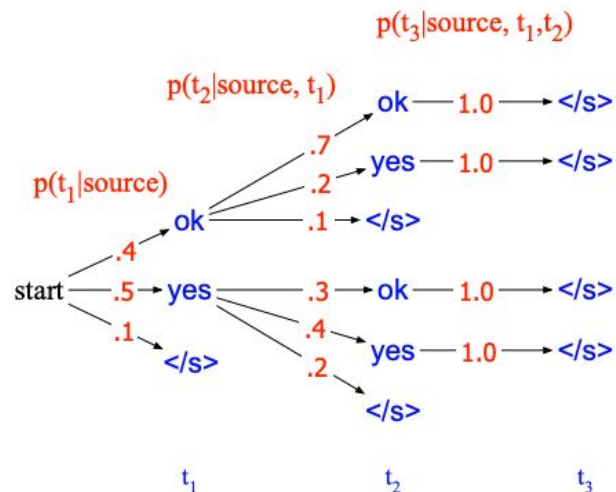
# Training

- Teacher Forcing
  - At each decoding step, we condition on the **true** values for previous context rather than the values we would have decoded from the decoder
- Loss function
  - At each step, we compute the cross entropy loss of the true tokens

$$-\sum_{c=1}^M y_{o,c} \log(p_{o,c})$$

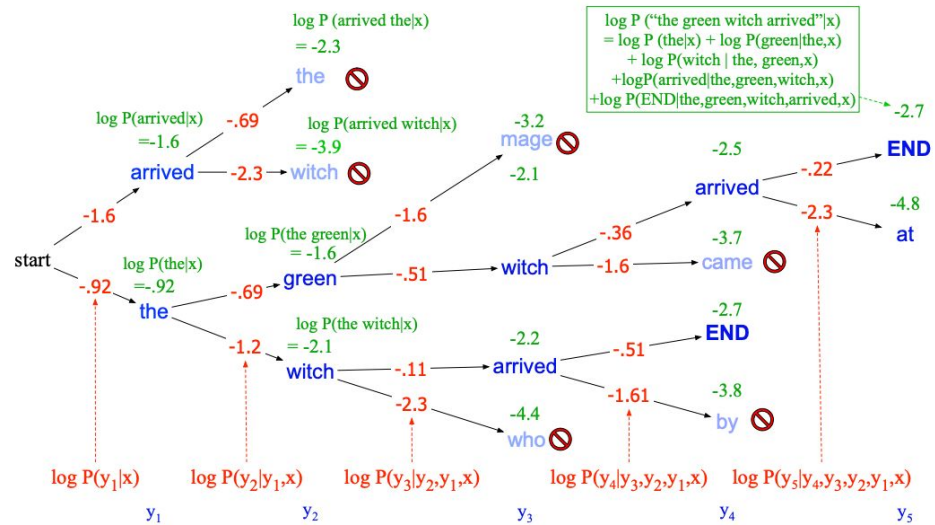
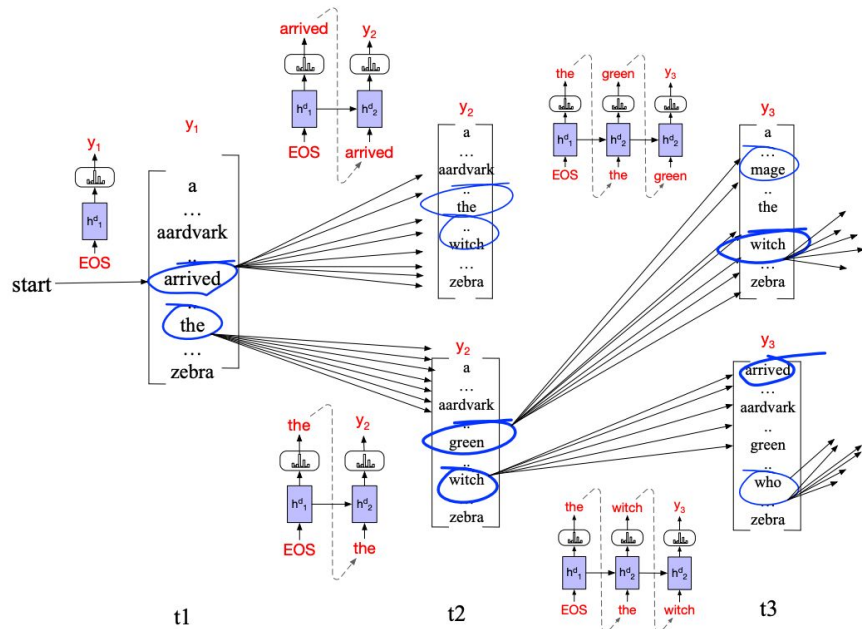
# Inference/Decoding – Simple

- Greedy Decoding: At each decoding step, take the most likely token
- Sampling: Sample from the output token distribution
- Problems:
  - The best/most likely sequence may not be the one where each token is most likely at each step





# Beam Search



# Notebook Link

[https://colab.research.google.com/drive/1z2eL6qtXmVRdsXCYMVuMp\\_wpUkEuB\\_xuA?usp=sharing](https://colab.research.google.com/drive/1z2eL6qtXmVRdsXCYMVuMp_wpUkEuB_xuA?usp=sharing)

# References

Speech and Language Processing. Chapter 10: Machine Translation and Encoder-Decoder Models. Daniel Jurafsky & James H. Martin.

[https://web.stanford.edu/~jurafsky/slp3/old\\_dec21/10.pdf](https://web.stanford.edu/~jurafsky/slp3/old_dec21/10.pdf)

[Loss Functions](#)

[The Illustrated Transformer](#)