

Transformers

Motivation and Intuition

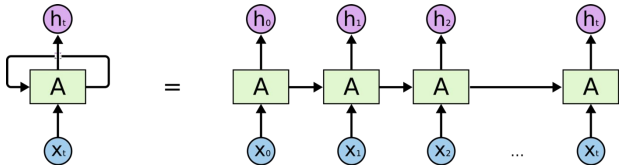
Overview

3 slide decks:

- **Motivation and Intuition**
- Architecture
- Popular Transformer-based models: BERT, GPT, RoBERTa, etc.

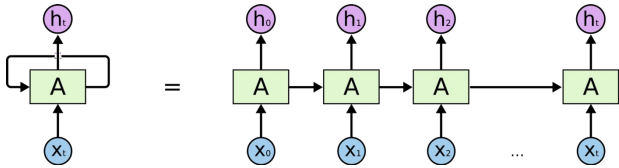
Motivations

- Which disadvantages of recurrent models can you think of?



Motivations

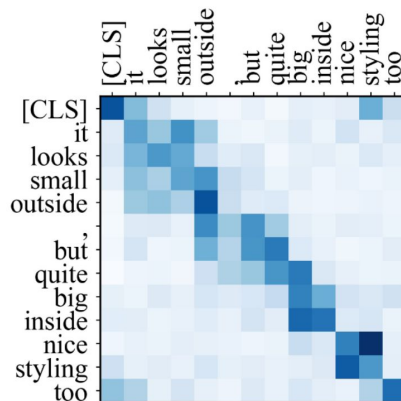
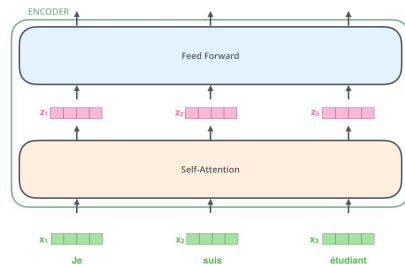
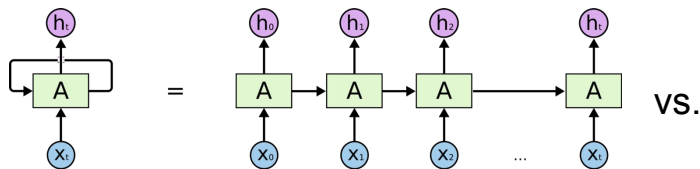
- Which disadvantages of recurrent models can you think of?



- Handling long distance dependencies: *Wer hat gestern Mittag kurz nach zwölf von meiner Palme die Kokosnuss geklaut?*
- Parallel processing: Faster implementations possible

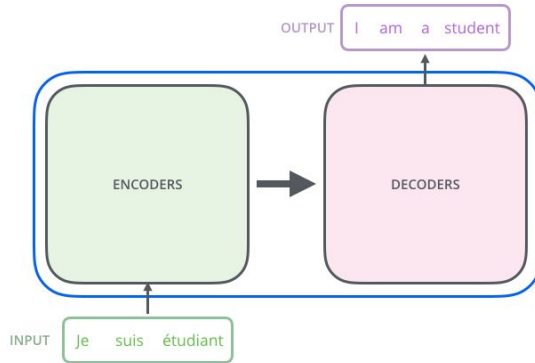
Motivation

- whole sentence is often available during processing
- So why go through the sentence word by word?
- Transformer is more efficient and uses more information!



Transformer Architecture

- Designed for Translation: Attention is all you need, Vaswani et al. 2017
- Encoder-Decoder



- Most LMs use only Encoder or Decoder

Most important hyperparameters

- representation size h (768, 1024, 12888)
- number of layers L (12, 24, 96)
- number of heads H (12, 16, 96)
- number of parameters P (110M, 340M, 175B)

