

# Using the Transformer

## BERT – Architecture



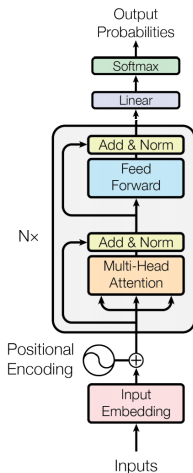
### Learning goals

- Understand the use of the transformer encoder in this model
- Understand the architectural components

## *Bidirectional Encoder Representations from Transformers:*

- Bidirectionally contextual model
  - Introduces new self-supervised objective(s)
  - Completely replaces recurrent architectures by Self-Attention  
+ simultaneously able to include bidirectionality
  - Transformer *encoder* as backbone of the architecture
    - 12 (24) Transformer encoder blocks
    - Embedding size of  $E = 768$  (1024)
    - Hidden layer size  $H = E$
    - $A = H/64 = 12$  (16) attention heads
    - Feed-forward size is set to  $4H$
- 110M (340M) parameters in total for  $BERT_{Base}$  ( $BERT_{Large}$ )

# CORE OF BERT – THE TRANSFORMER ENCODER



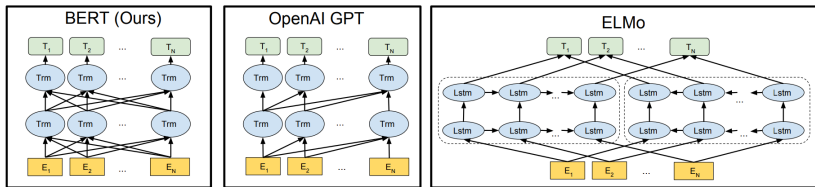
Source: Vaswani et al. (2017)

# A REMARK ON "CAUSALITY"

## Causality is an issue!

- Goal: Learn contextual representations for words/tokens
- *Self-Supervision*: Input and target sequence are the same  
→ We modify the input to create a meaningful task
- Unconstrained Self-Attention makes using the LM objective infeasible
- Bidirectionality at a lower layer would allow a word to see itself at later hidden layers  
→ The model would be allowed to cheat!  
→ This would not lead to meaningful internal representations

# ELMO VS. GPT VS. BERT



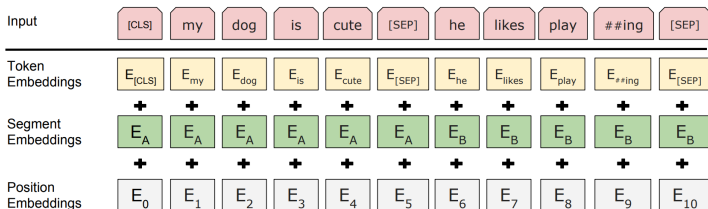
Source: Devlin et al. (2018)

## Major architectural differences:

- ELMo uses two separate unidirectional models to achieve bidirectionality → Only "*shallow*" bidirectionality
- GPT is not bidirectional, thus no issues concerning causality
- BERT combines the best of both worlds:

*Self-Attention + (Deep) Bidirectionality*

# INPUT EMBEDDINGS



Source: Devlin et al. (2018)

- Two concatenated sentences as input
- WordPiece tokenization (Wu et al., 2016) for the inputs  
→ Vocabulary of 30.000 tokens
- Learned segment + position embeddings
- Special [CLS] and [SEP] tokens