# Using the Transformer

# BERT (Devlin et al., 2018)
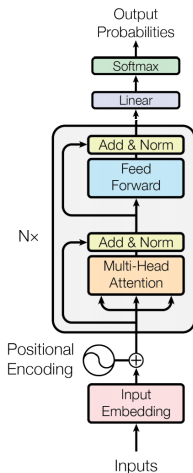


**Learning goals**

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

## KEY FACTS ON BERT ▸ DEVLIN ET AL. (2018)

*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 $4H$
    $\rightarrow$ 110M (340M) parameters in total for $BERT_{Base}$ ($BERT_{Large}$)

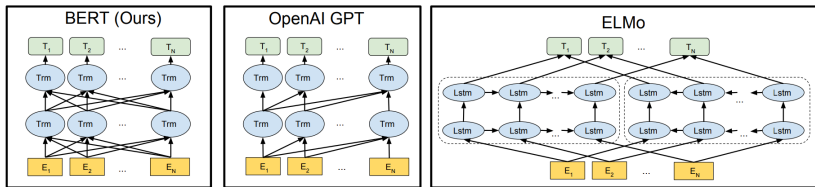# CORE OF BERT – THE TRANSFORMER ENCODER



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
  $\rightarrow$ 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
  $\rightarrow$ The model would be allowed to cheat!
  $\rightarrow$ 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:

$$Self\text{-}Attention + (Deep)\ Bidirectionality$$

# INPUT EMBEDDINGS



| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

Source: Devlin et al. (2018)

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