**BERT**

**B**idirectional **E**ncoder

**R**epresentations from
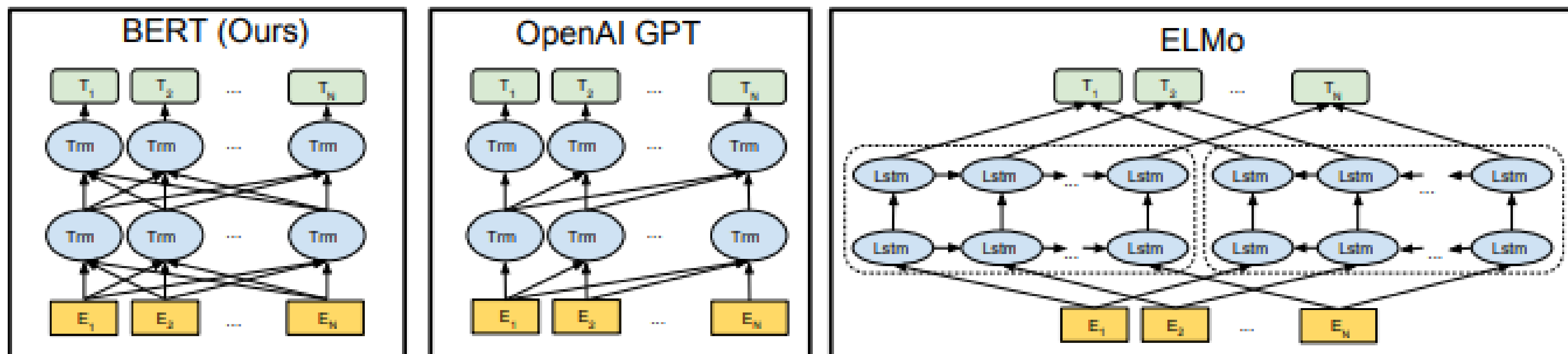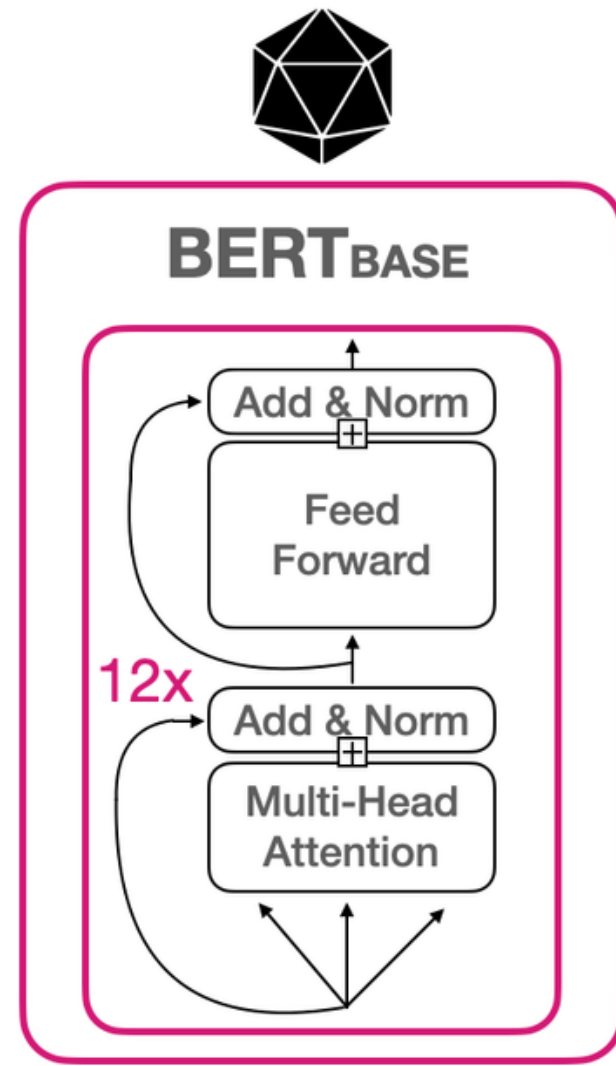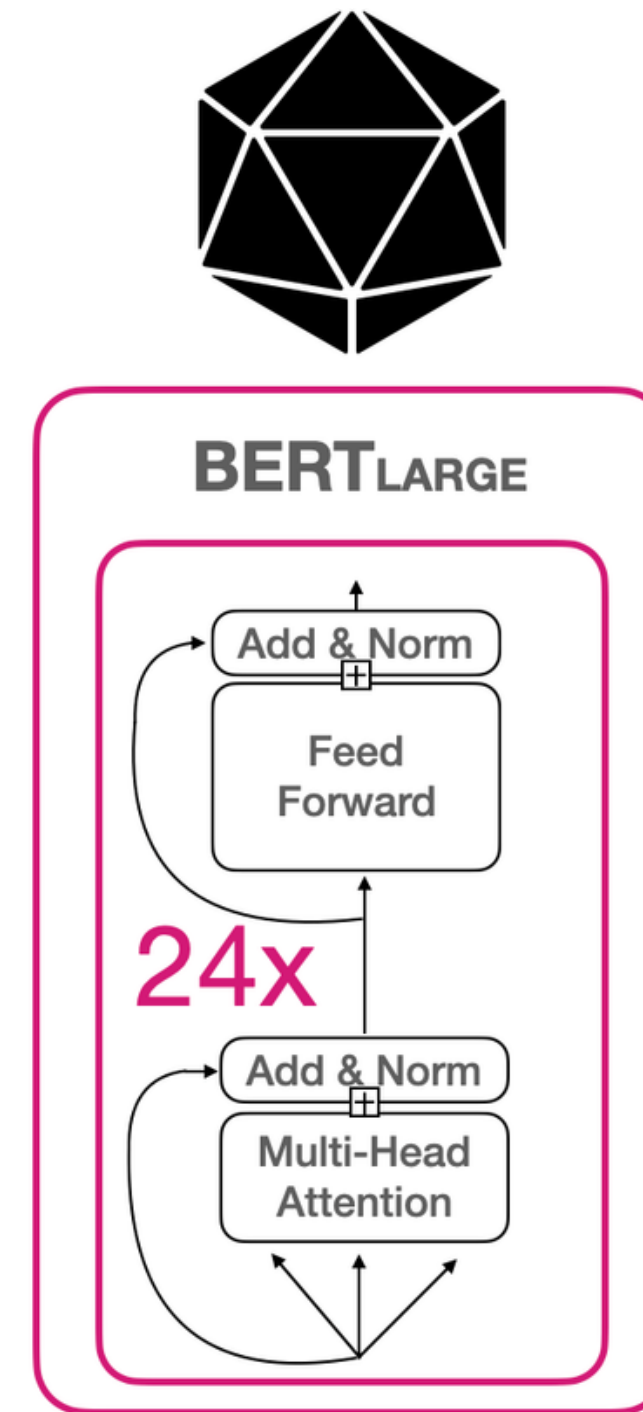
**T**ransformers

# Architecture



Figure 1: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTM to generate features for downstream tasks. Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

# Architecture



BERT Size & Architecture

**BERT**BASE

Add & Norm

Feed
Forward

12x

Add & Norm

Multi-Head
Attention

110M Parameters

**BERT**LARGE

Add & Norm

Feed
Forward

24x

Add & Norm

Multi-Head
Attention

340M Parameters

# Input Representation
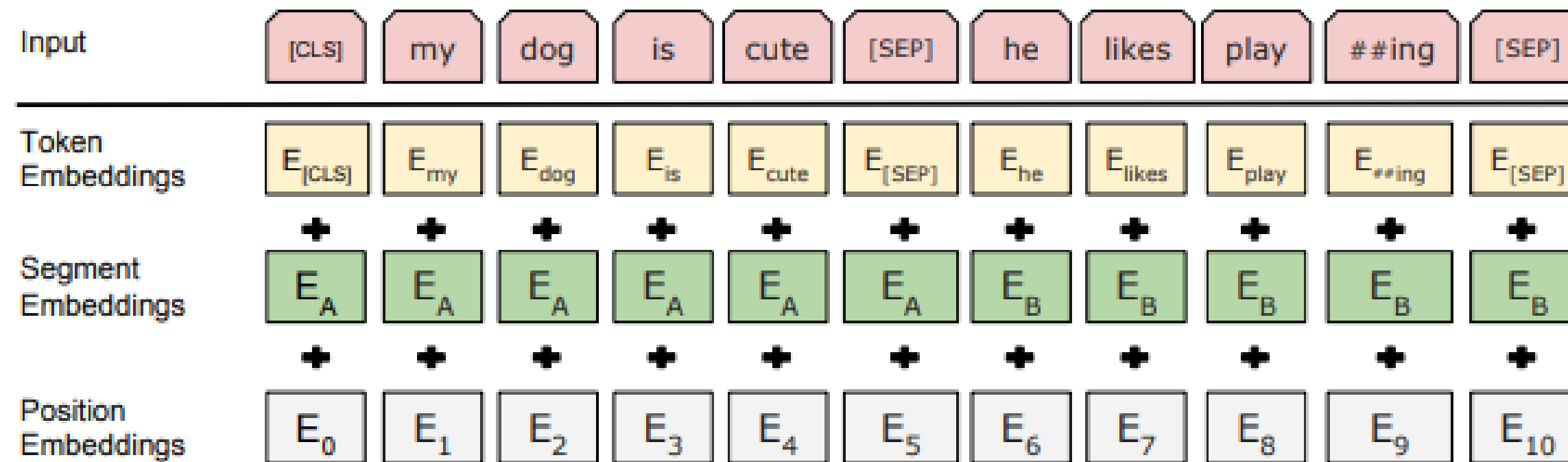
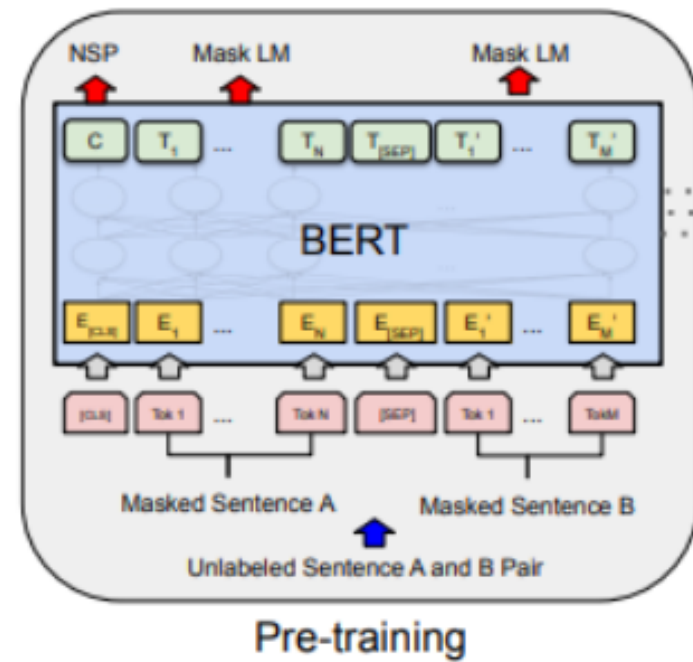

Figure 2: BERT input representation. The input embeddings is the sum of the token embeddings, the segmentati embeddings and the position embeddings.
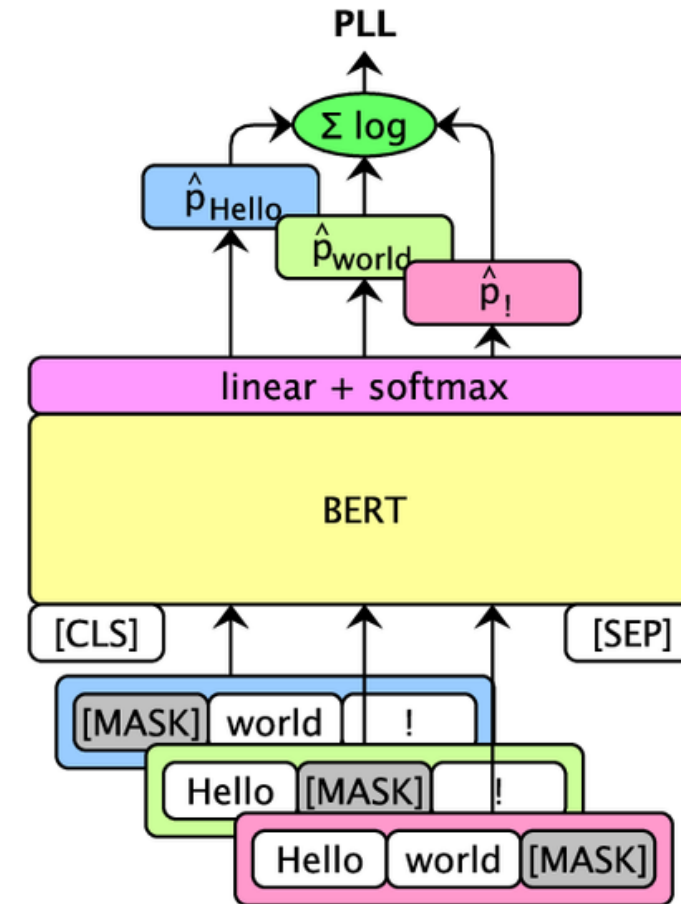
——→ Thêm Segment Embedding
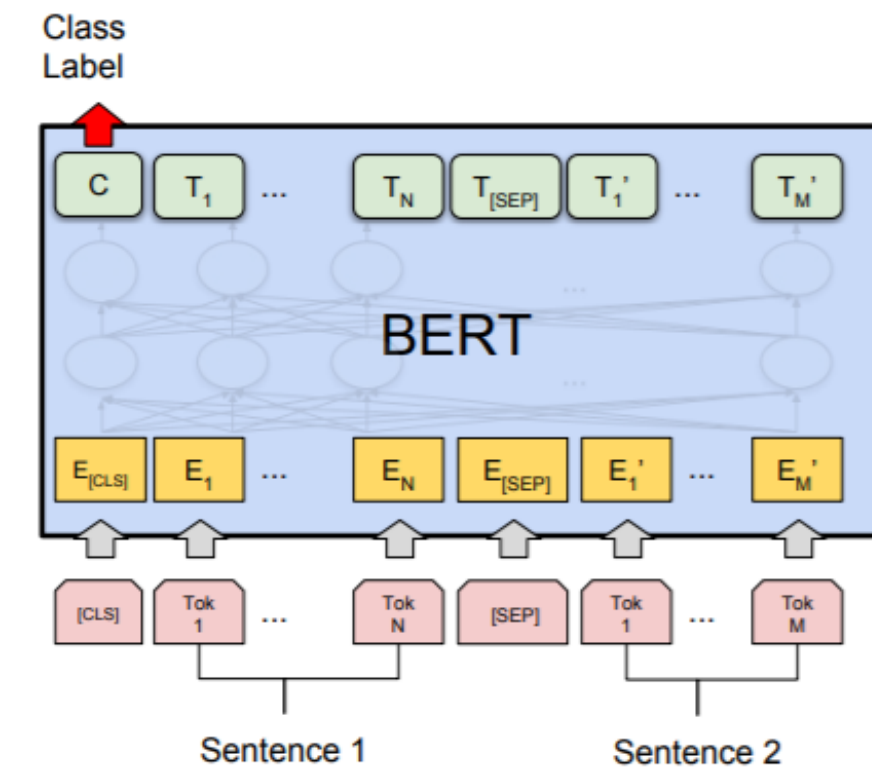
——→ Học thêm về mối tương quan của các câu
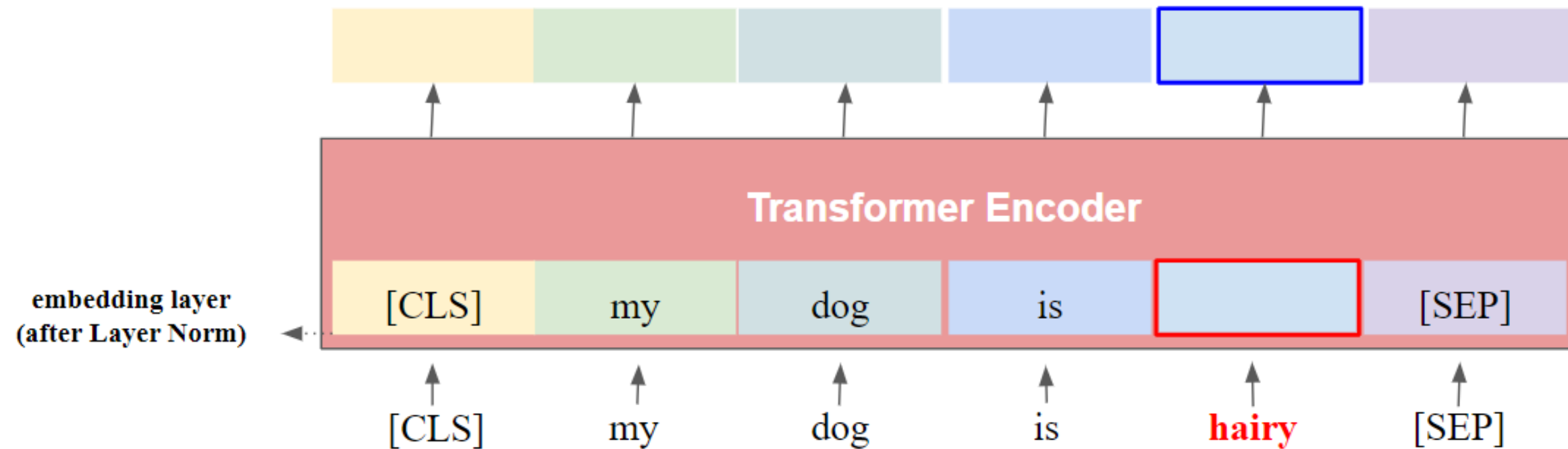
# Pre-training Tasks



Mask LM

Next Sentence Prediction

# MASK LM
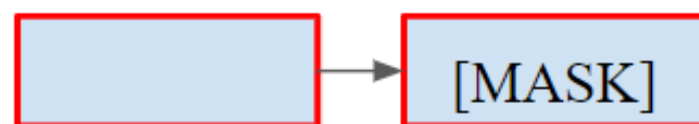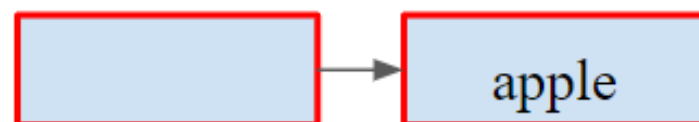
Target: Học cách dự đoán dựa trên bối cảnh ⟶ Che đi một phần dữ liệu
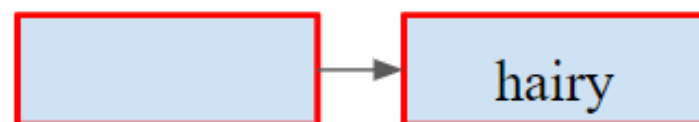


Mask **15%** of all WordPiece tokens in each sequence at **random**. ( e.g., **hairy** )
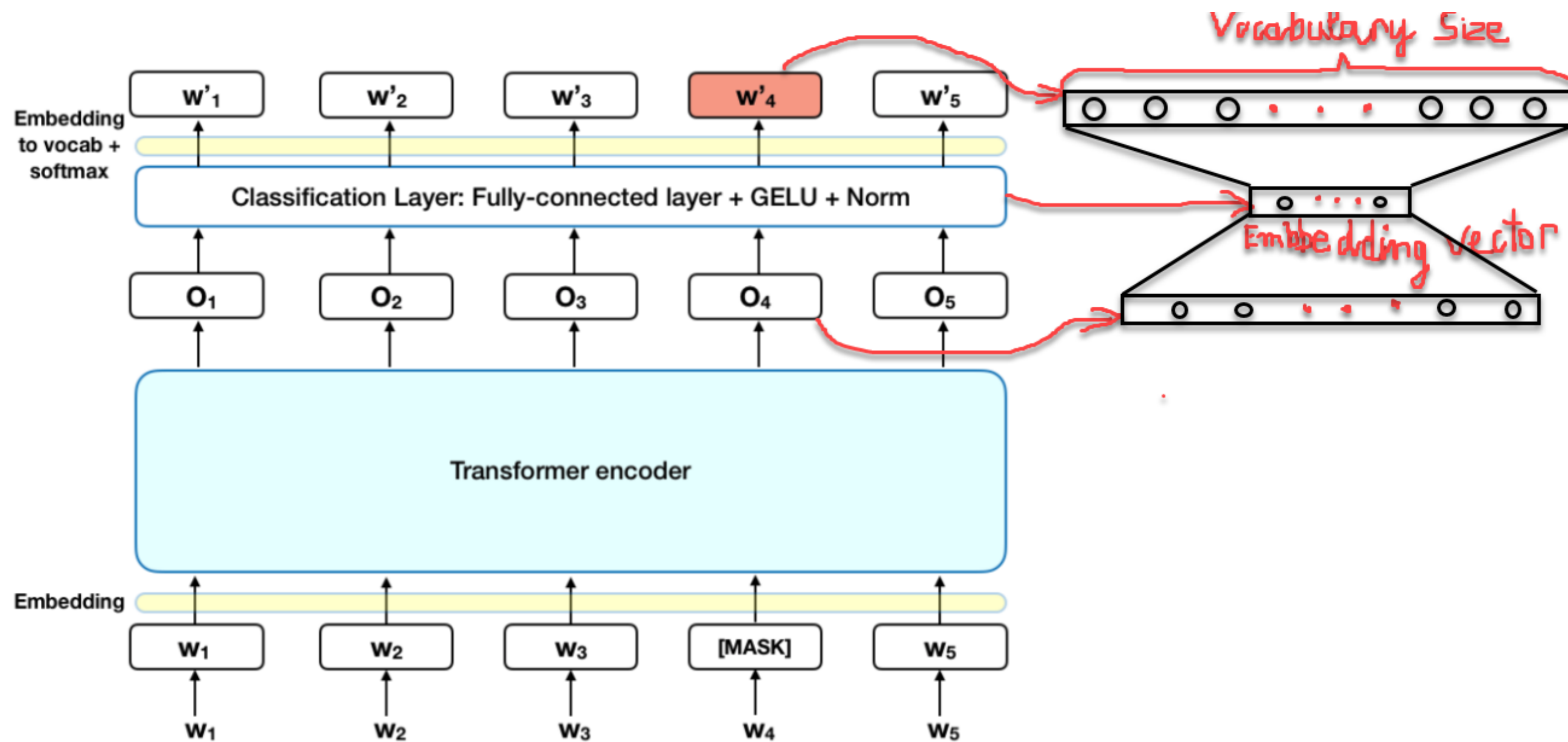
[ ] ⟶ [MASK]    80% of the time : Replace **[MASK]** token.

[ ] ⟶ apple    10% of the time : Replace the word with a **random** word

[ ] ⟶ hairy    10% of the time : Keep the word **unchanged**.

# MASK LM

# MASK LM

# Next Sentence Prediction



related (isNext)

not related (notNext)

# Reference

- https://www.miai.vn/2020/12/14/bert-series-chuong-1-bert-la-cai-chi-chi/

- https://viblo.asia/p/hieu-hon-ve-bert-buoc-nhay-lon-cua-google-eW65GANOZDO

Thanks for watching!