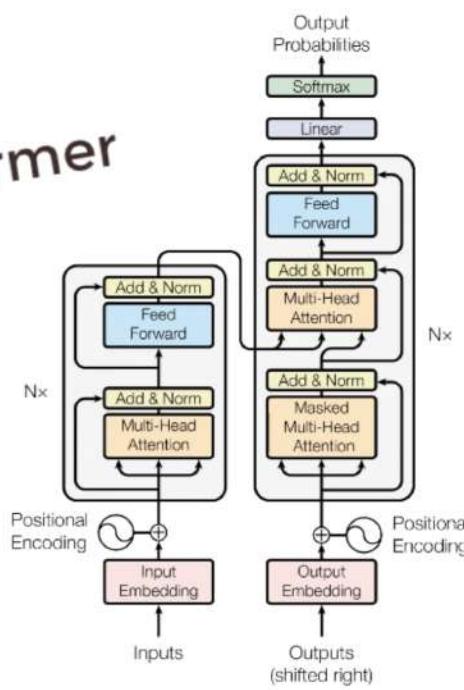


# BERT

## LSTM Vs Transformer

Transformer  
↔

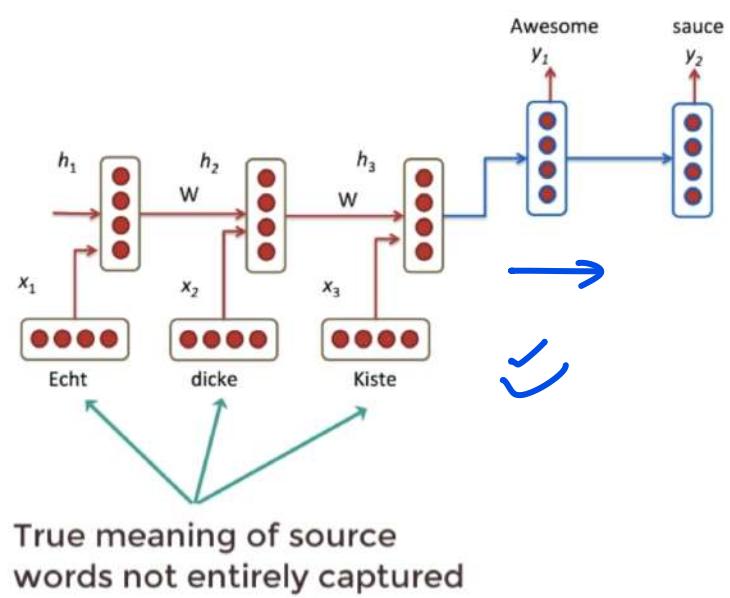


Source: Attention is All you Need (Vaswani et al., 2017)

# LSTM Vs Transformer

## LSTM Networks

1. Slow ✓
2. Not truly Bidirectional



Source [Blog]: Attention & Memory in Deep Learning (Britz, 2017)

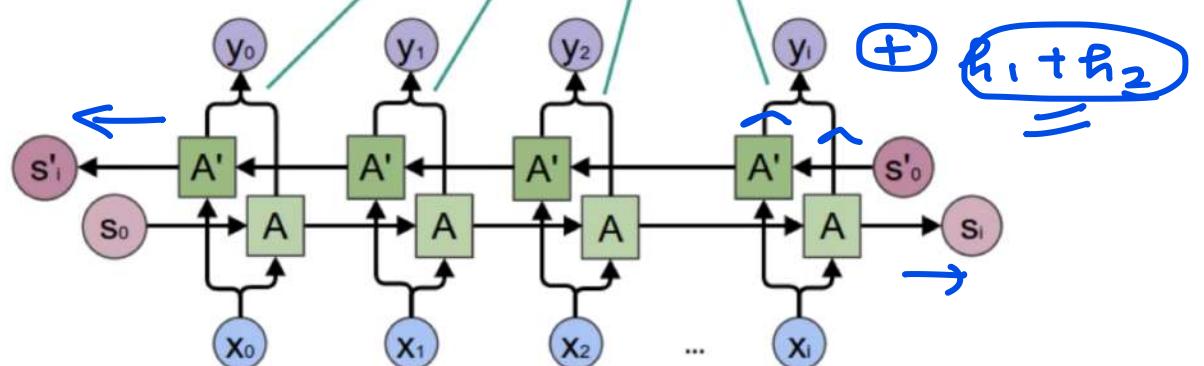
# LSTM Vs Transformer

## LSTM Networks

1. Slow 2  $1 \rightarrow 1 \text{ min}$   $\approx 2 \text{ min}$
2. Not truly Bidirectional

Simple concatenation

Bi LSTM  
= RNN



Source [Blog]: Understanding LSTMs (Colah, 2015)

# LSTM Vs Transformer

? ✓ →

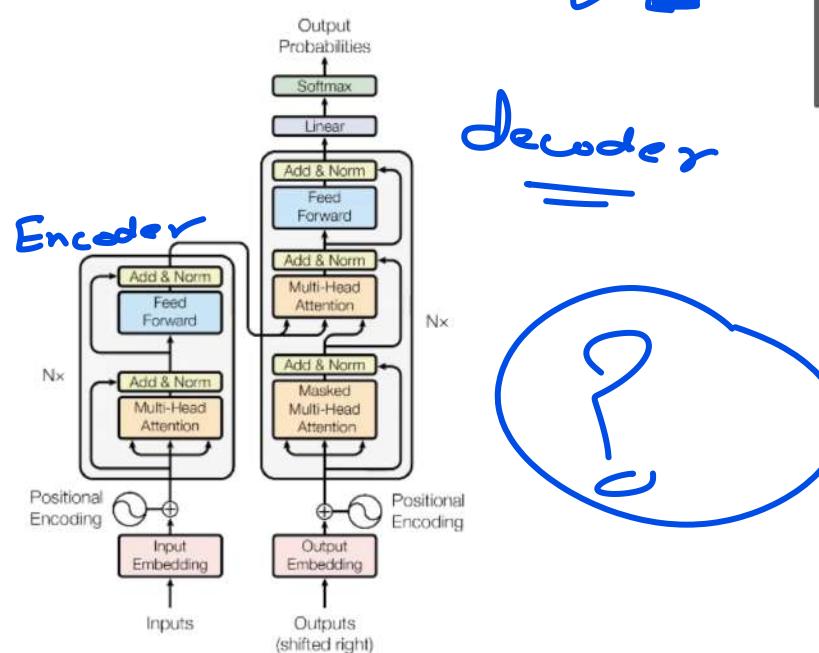
## Transformer

1. ~~Slow~~ Faster
2. ~~Not truly Bidirectional~~

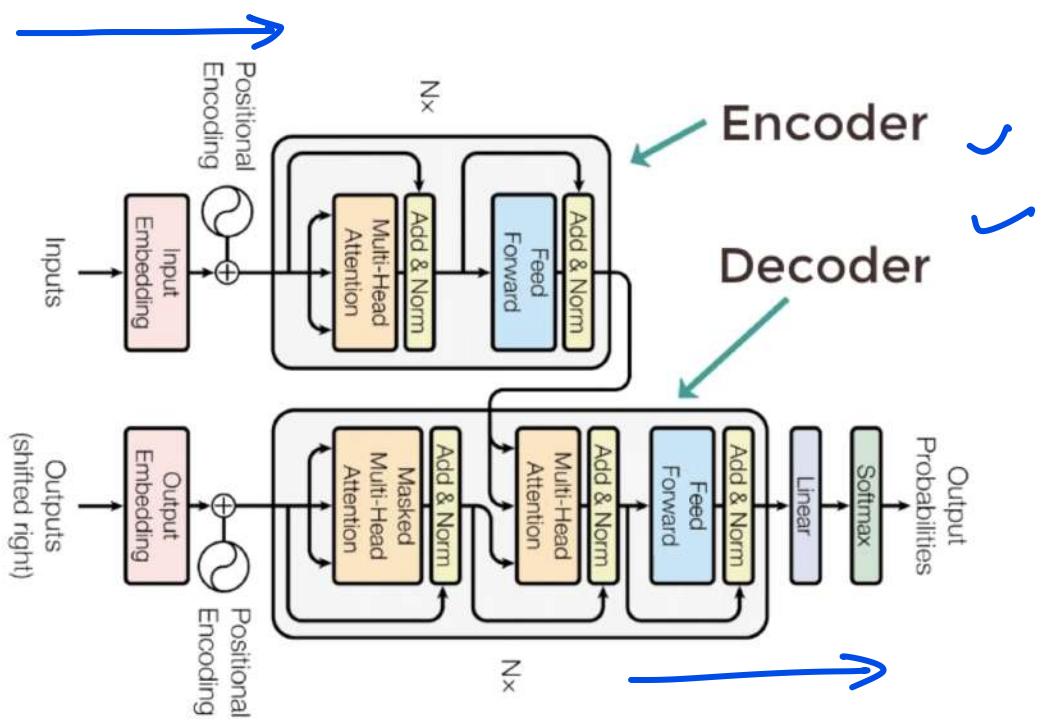
Deeply Bidirectional

$x_i$  Self attention

left context - - - - - right context



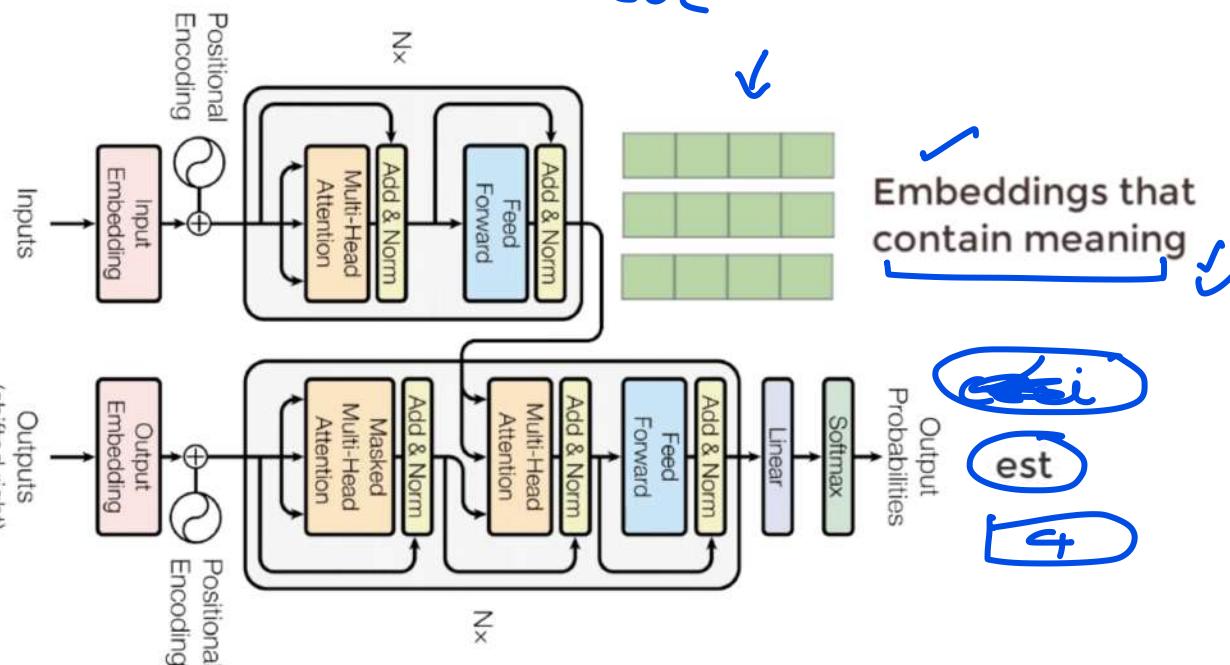
# Transformer Flow



# Transformer Flow

This  
is  
Amazing

<start>  
Ceci  
est

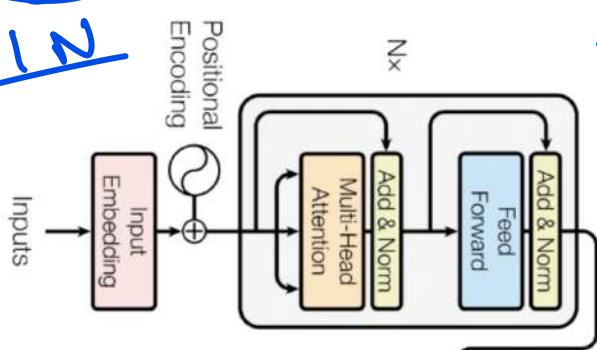


# Transformer Flow

Eng → french

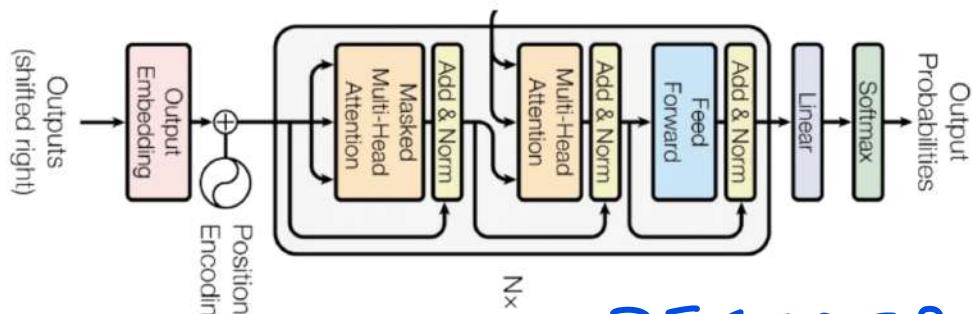
ENCODER training

R  
T  
Y



What is English? What is context?  
What is language!

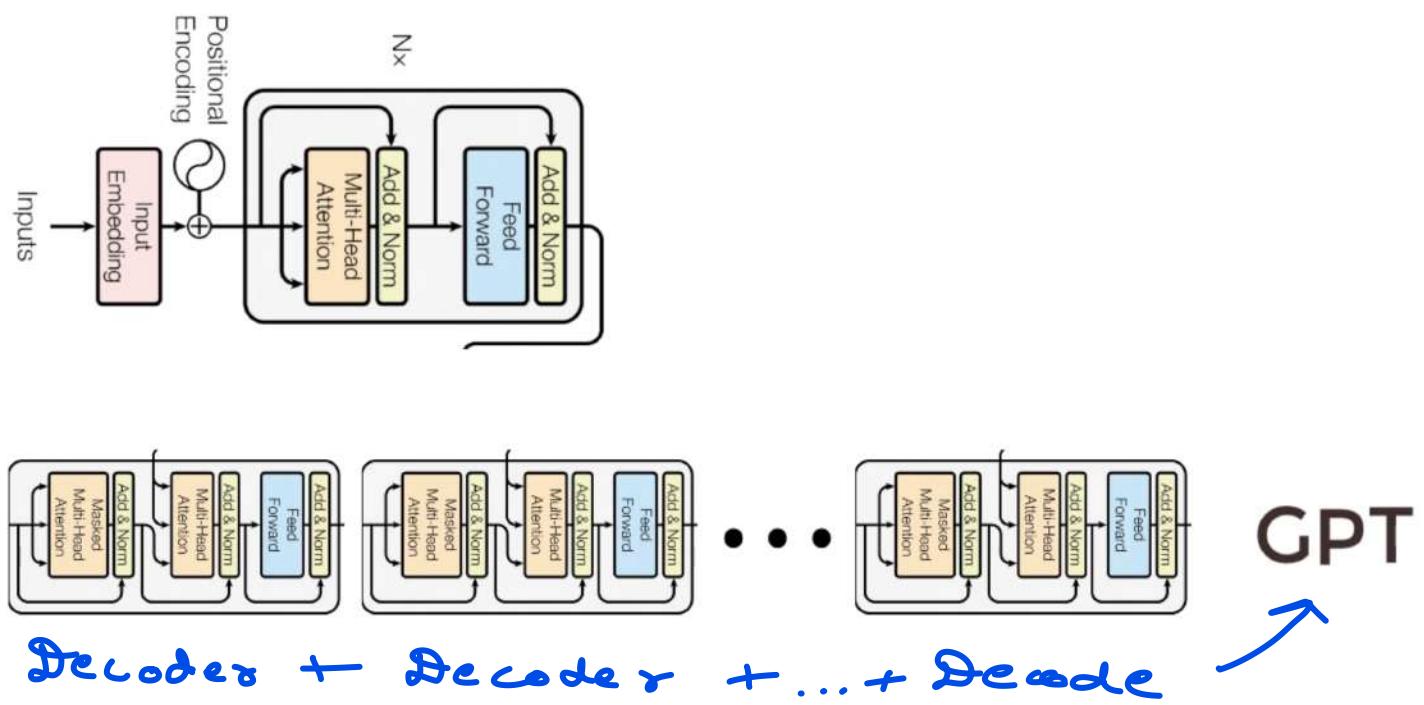
NLU



How to map English words to French words?  
What is language!

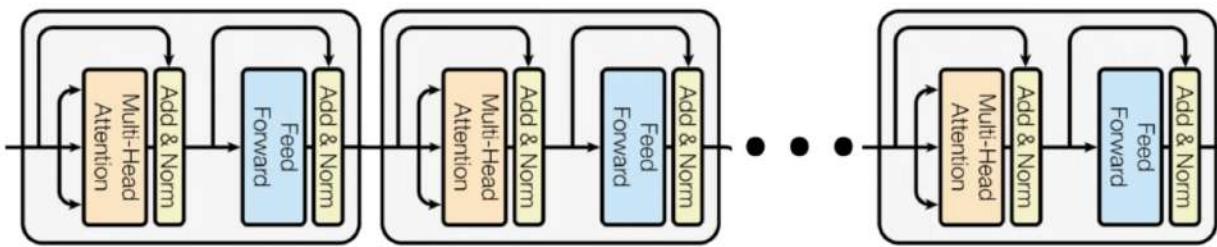
DECODER

# Transformer Flow



# Transformer Flow

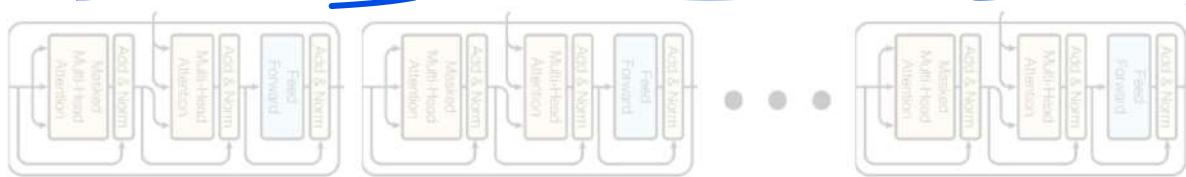
Encoder + Encoder + ... + Encoder



?

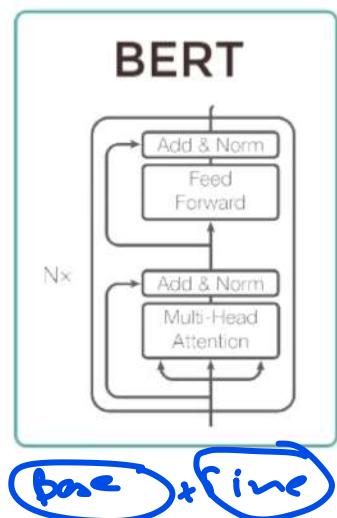
BERT

Bidirectional Encoder Representation from Transformers



GPT

# Bidirectional Encoder Representation from Transformers



## Problems to Solve

- Neural Machine Translation
- Question Answering
- Sentiment Analysis
- Text summarization

cross-farmish

## Encoder

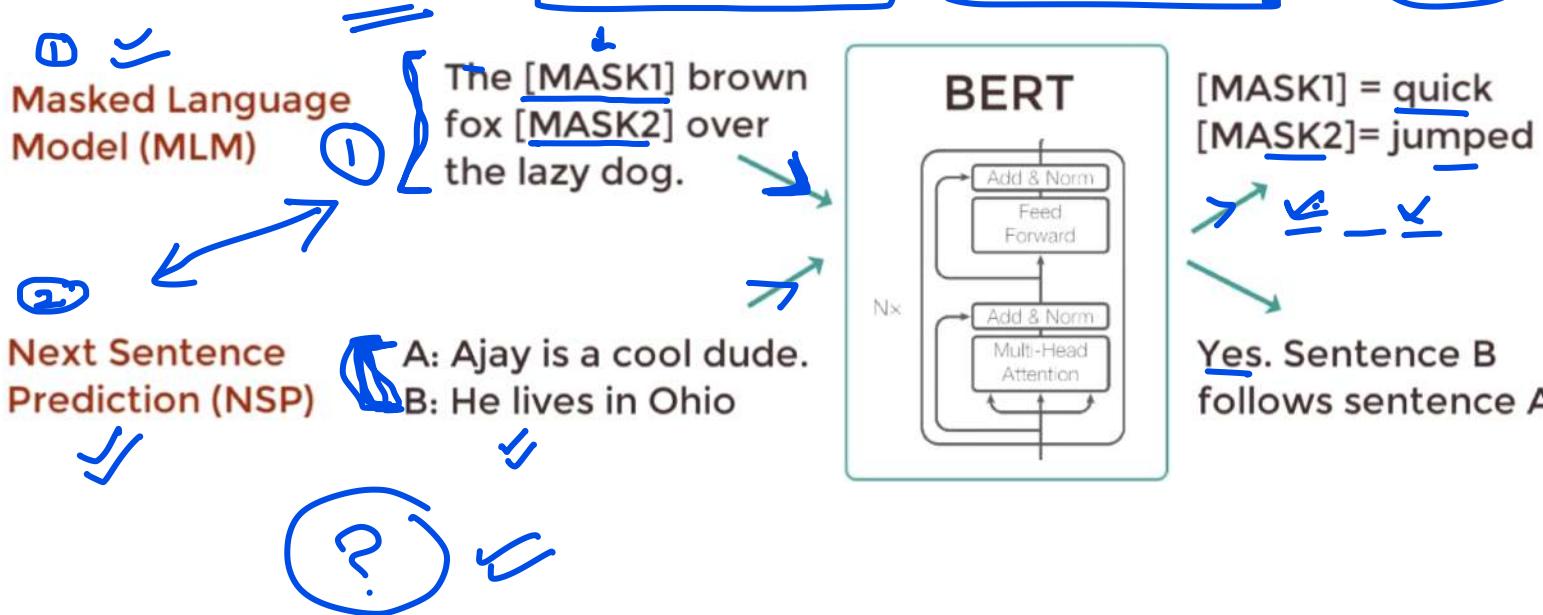
Needs Language understanding

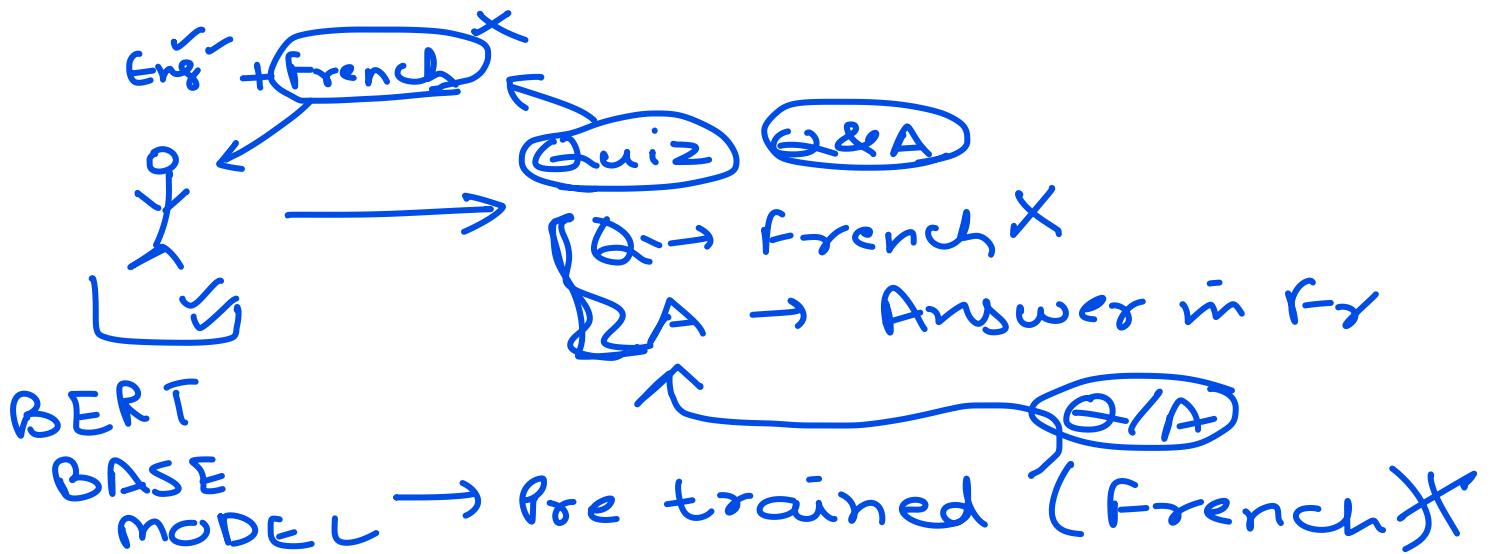
## How to solve Problems (BERT Training)

- Pretrain BERT to understand language
- Fine tune BERT to learn specific task

## Bidirectional Encoder Representation from Transformers

Pretraining (Pass 1) : "What is language? What is context?" = NW





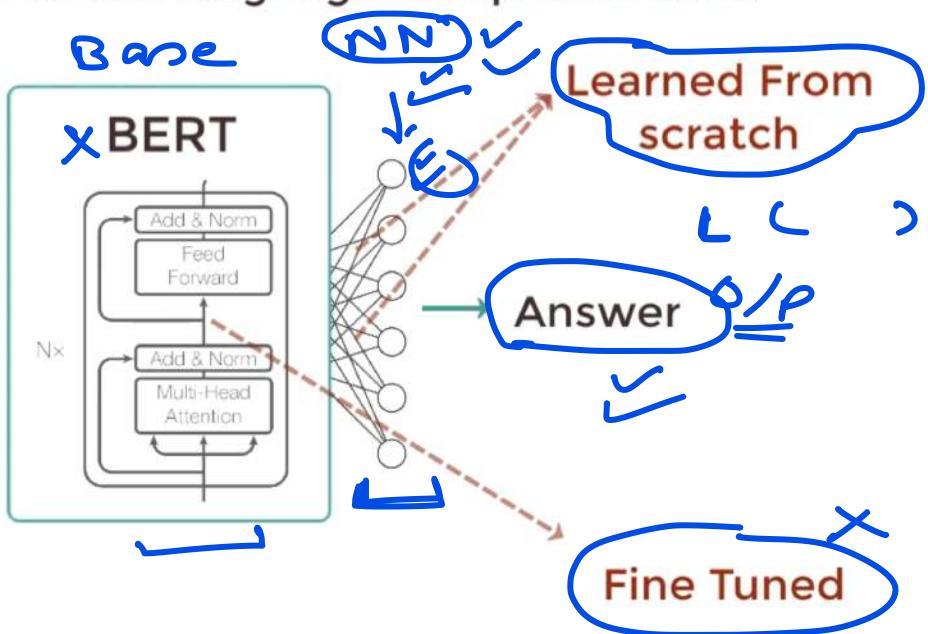
## Bidirectional Encoder Representation from Transformers

Fine Tuning (Pass 1): "How to use language for specific task?"

Fine tuned Q & A

Question  
Passage

**FAST!**

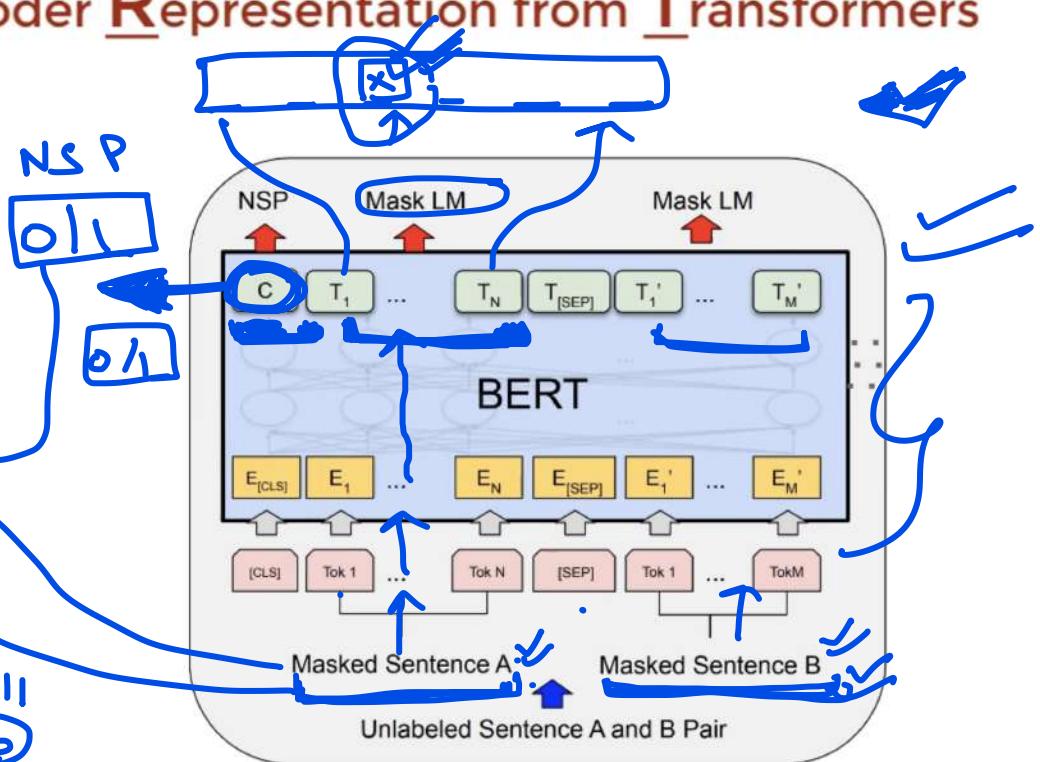


# Bidirectional Encoder Representation from Transformers

## Pretraining (Pass 2)

Problems to train on simultaneously:

- ✓ 1. Masked Language Modeling (Mask LM)
- ✓ 2. Next Sentence Prediction (NSP)



Source: BERT: Pre-training of deep bidirectional Transformers for language understanding (Devlin et al., 2019)

## Bidirectional Encoder Representation from Transformers

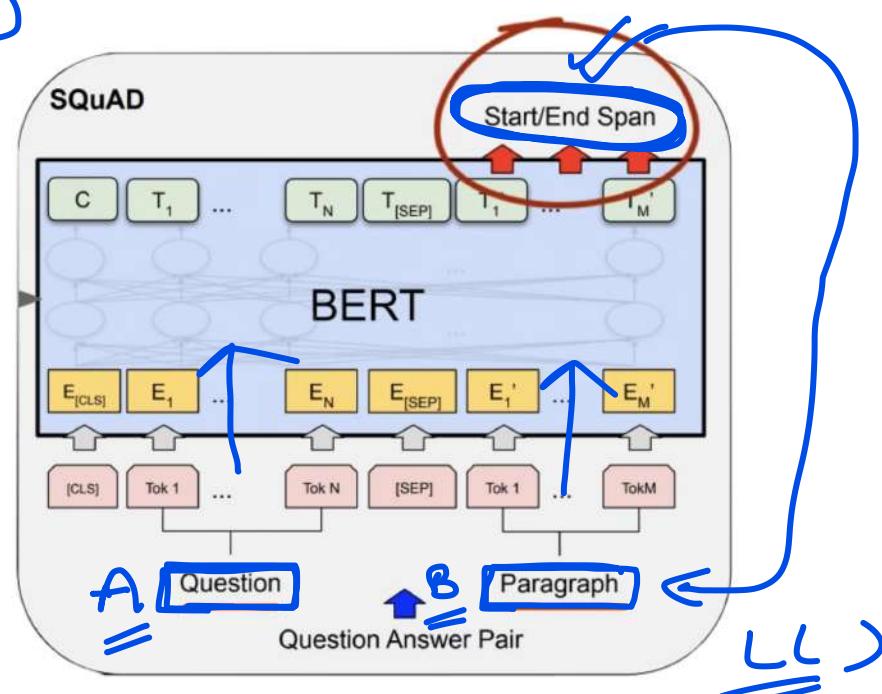
Fine Tuning (Pass 2)

Q/A

Change output to display  
text in which answer exists

$\equiv \exists Q \rightarrow A$

Change inputs to take in  
Question, Passage

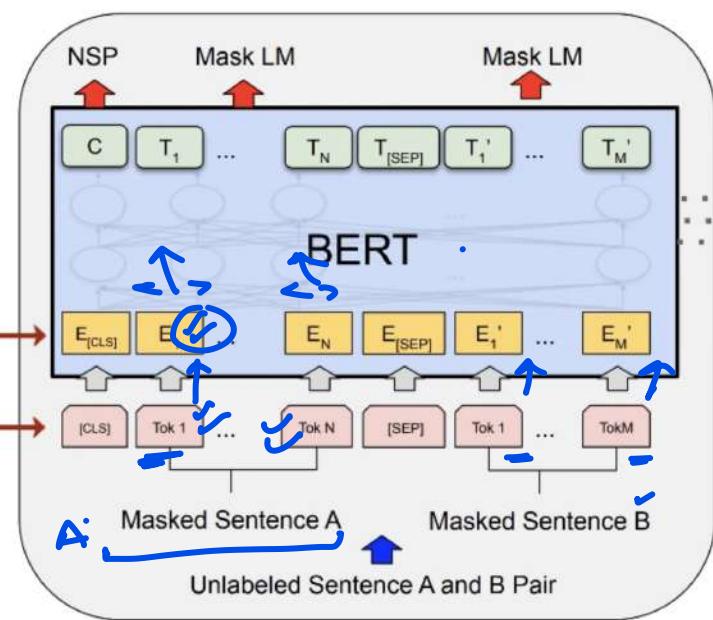


# Bidirectional Encoder Representation from Transformers

## Pretraining (Pass 3)

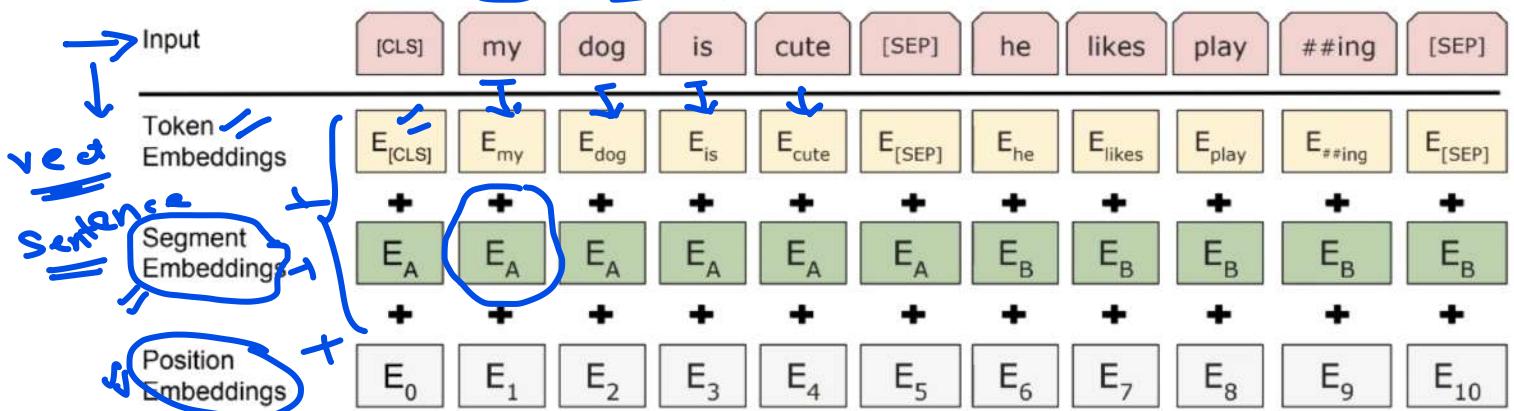
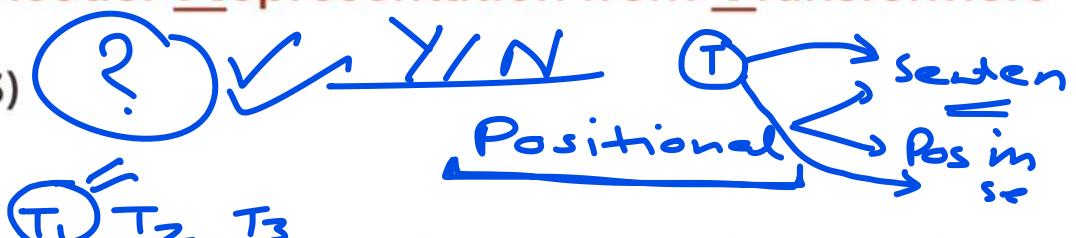
$$E_1 = T_1 + S_1 + P_1$$

initial embeddings  
word token inputs



## Bidirectional Encoder Representation from Transformers

Pretraining (Pass 3)



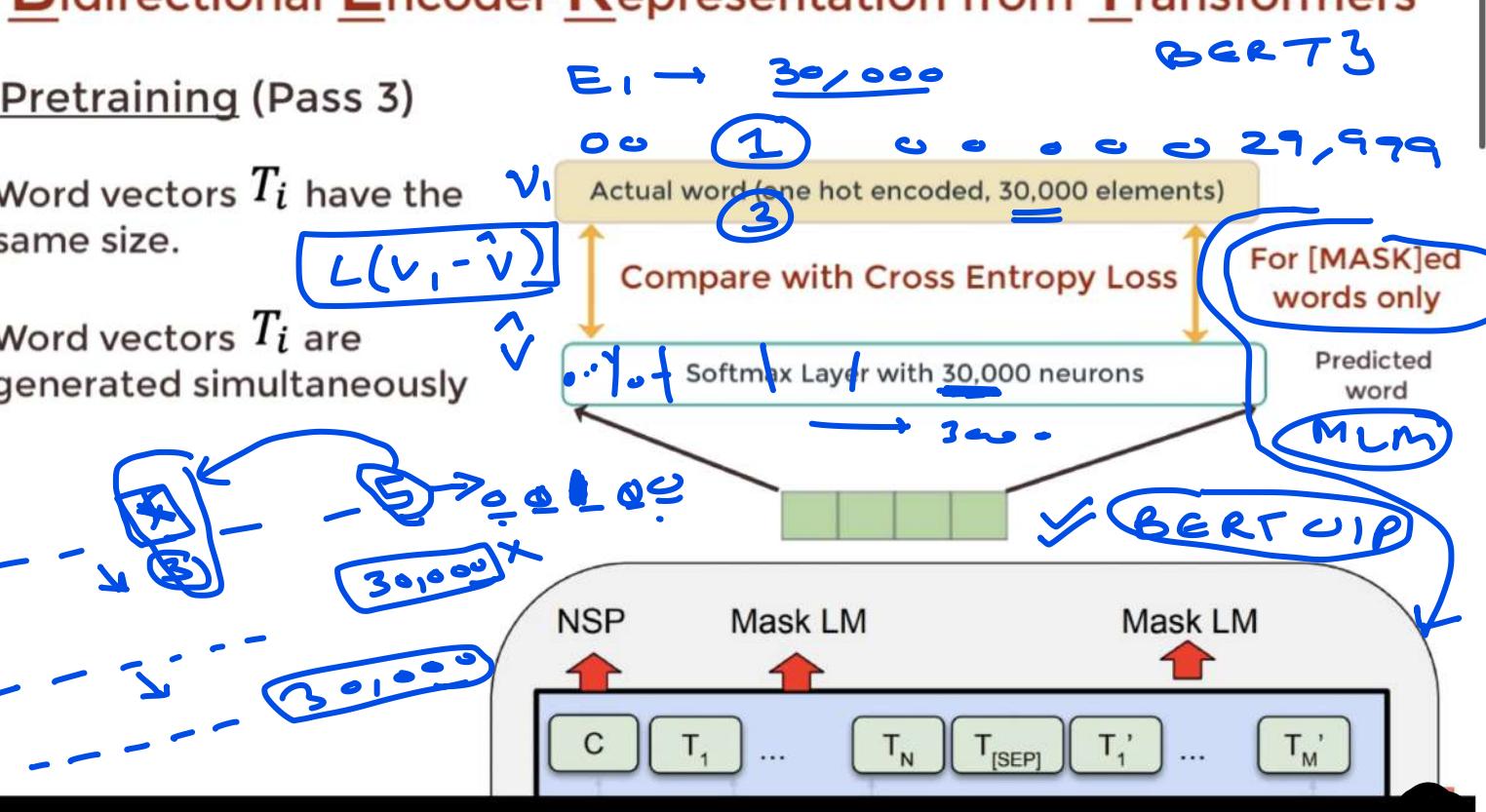
Segment + Position embeddings add “ordering” for inputs  $\Rightarrow$

# Bidirectional Encoder Representation from Transformers

## Pretraining (Pass 3)

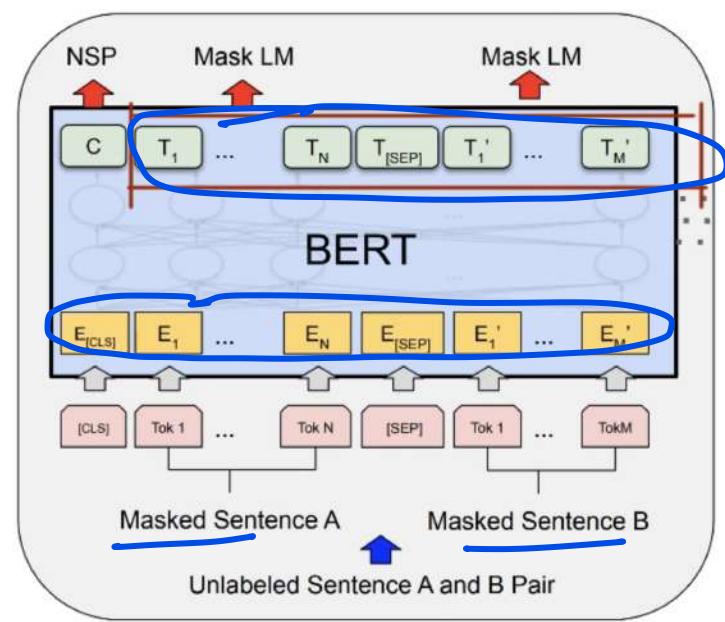
Word vectors  $T_i$  have the same size.

Word vectors  $T_i$  are generated simultaneously

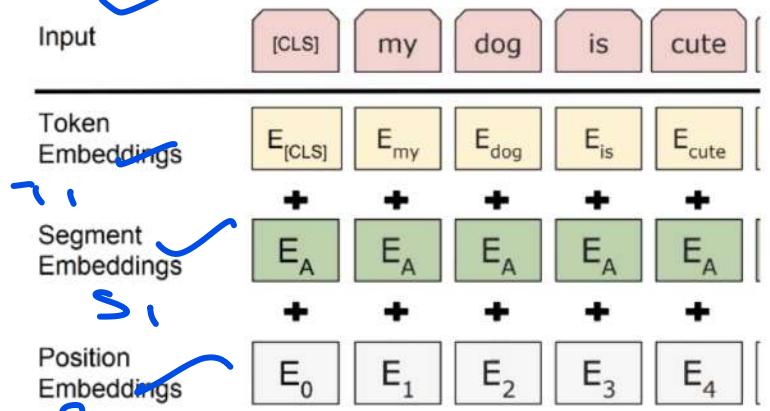


# Bidirectional Encoder Representation from Transformers

## Pretraining (Summary)



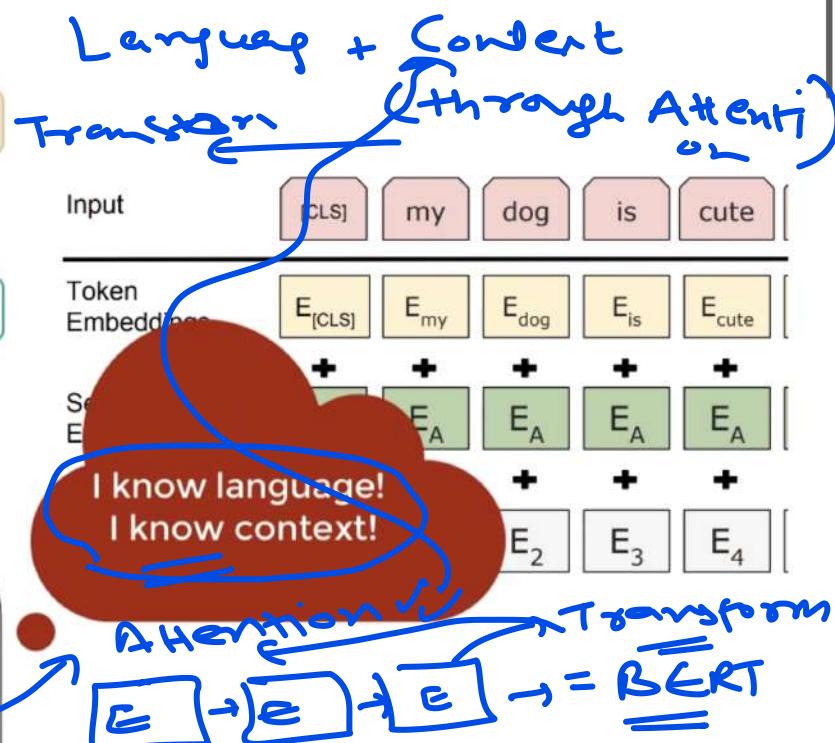
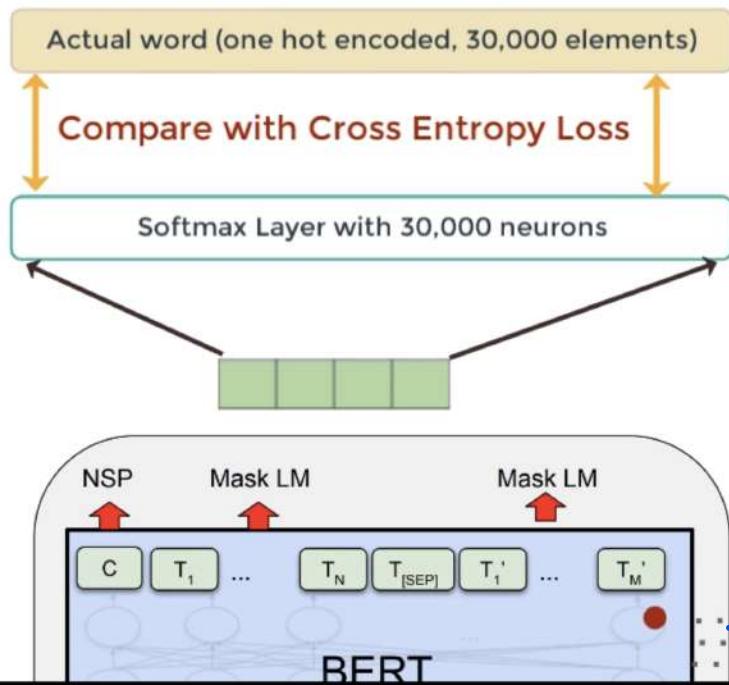
BERT → Pretrain → NLU



$$E_i = T_i + S_i + P_i$$

# Bidirectional Encoder Representation from Transformers

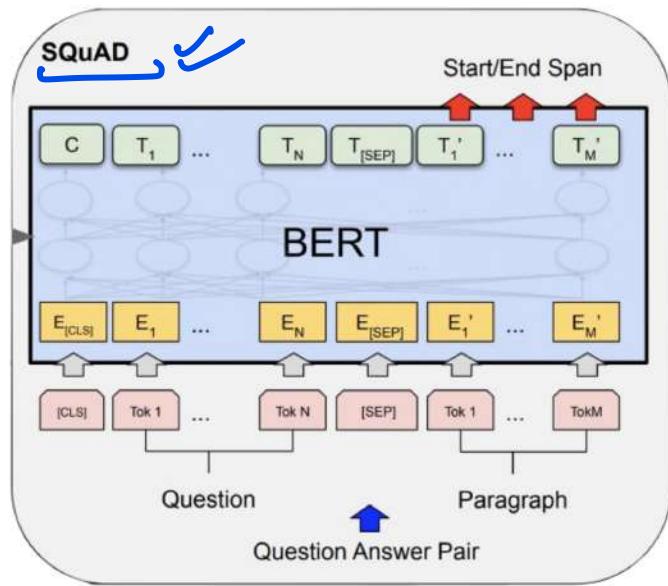
## Pretraining (Summary)



## Bidirectional Encoder Representation from Transformers

Fine Tuning (Summary)

Q/A = BERT 30min

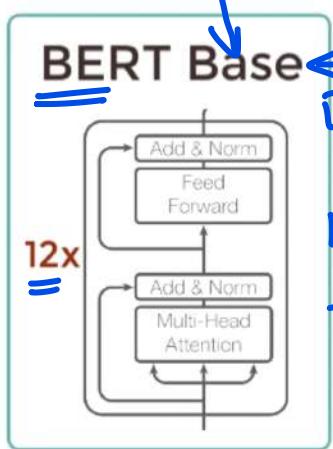
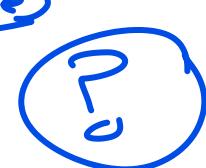


"Stanford Question & Answer Dataset"  
- 30 minute training on single cloud  
TPU with 91% F1 score.  
Tensor Processing Unit

## Bidirectional Encoder Representation from Transformers

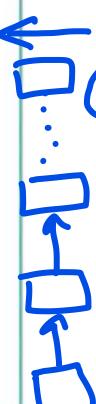
Performance

Score

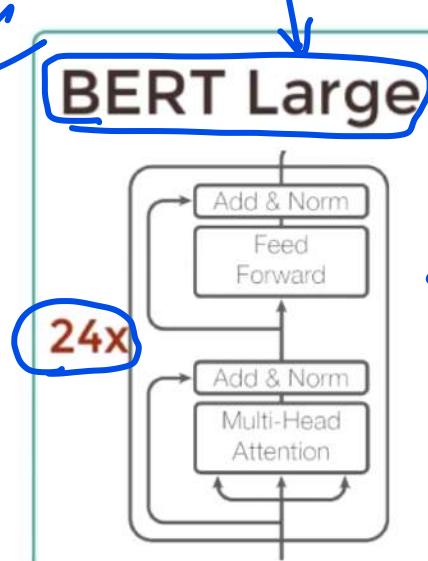


110M Parameters

12x



12



24x

340M Parameters

BERT Large



2<sup>n</sup> Encoder layers

Much more capable now

