# BERT

**Bidirectional Encoder Representations from Transformer**

Team 7

March 13, 2020
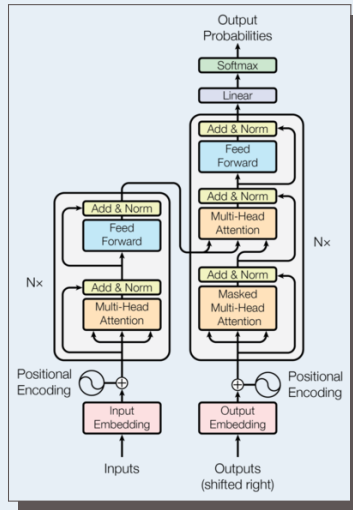
# Outline

# Transformer



Figure: **Transformer**: Attention is all you need
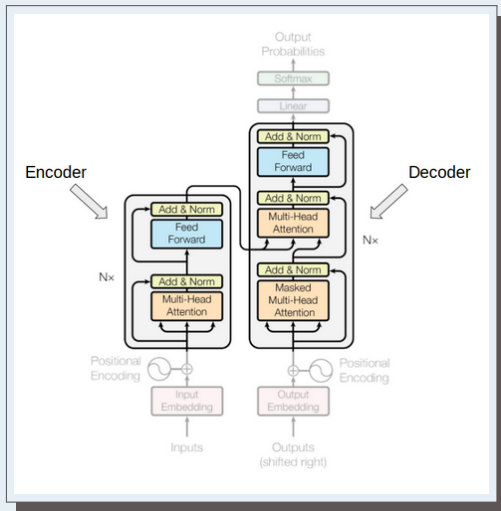
Figure: **Encoder-Decoder**

# Encoder
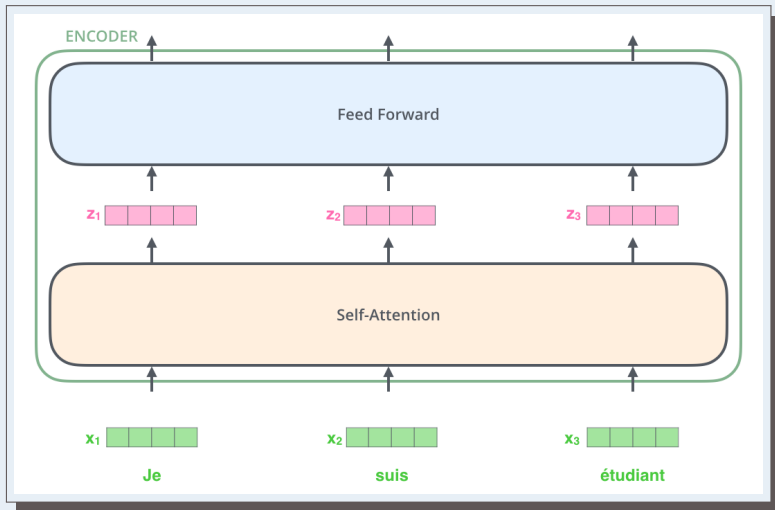


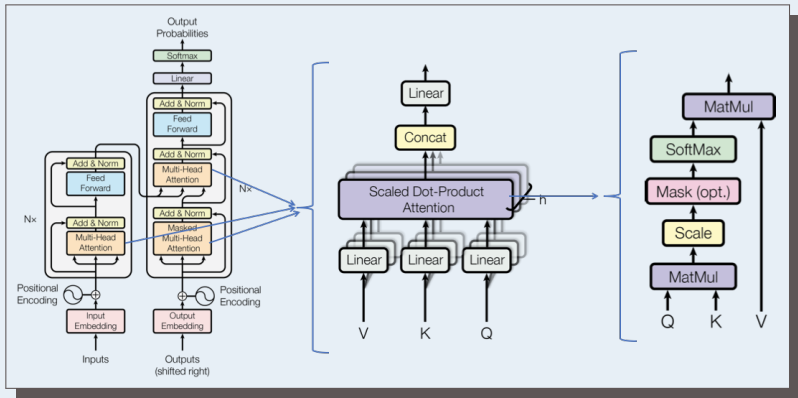Figure: **1 Encoder-layer**

# Self-Attention



Figure: **Self-Attention**
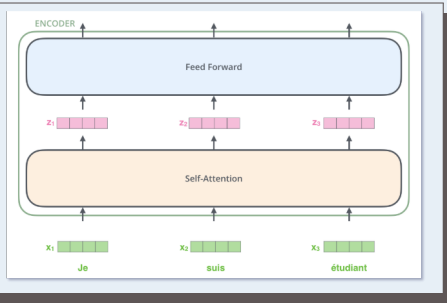
Figure: **BERT**

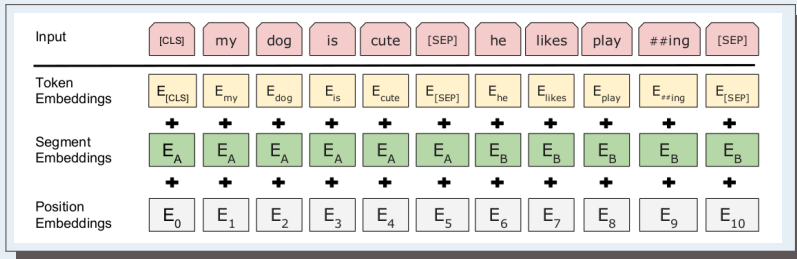Figure: **BERT-Architecture**

# Input



Figure: **BERT-Input**

# 2-stage-training

## Pre-training

Self-supervised tasks.

## Finetune

3 finetune tasks.

# Pre-Training

## Pretraining tasks

2 pretraining tasks

## NSP

Next Sentence Prediction

## MLM

Masked Language Model - self supervised learning

Figure: **Next Sentence Prediction**



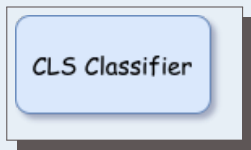Figure: **[CLS] token classifier**

What is self-supervised learning ?



Figure: **Self-supervised learning definition**

# MLM

[CLS] 做人如果沒夢想 [SEP] 那跟 [MASK] [MASK] 有什麼分別啊？[SEP]

鹹魚

Figure: **Masked LM**

MLM Classifier

Figure: **MLM classifier**

# Fine-Tune

Downstream tasks:

## NER
Name Entity Recognition

## NLI
Natural Language Inference

## QA
Question Answering

Figure: **Natural Language Inference**



Figure: **[CLS] token classifier**

[CLS] 彼得住在哪裡 ？ [SEP] 彼得是一個居住在歐洲的聰明男孩，
最近新冠肺炎在歐洲疫情十分險峻，聰明的彼得想到了一個方法，去減緩疫情的發展。[SEP]

Ans: 歐洲
start index: 17 end index: 18

Figure: **Question Answering Example**

start classifier     end classifier

Figure: **Random initialize two classifiers**
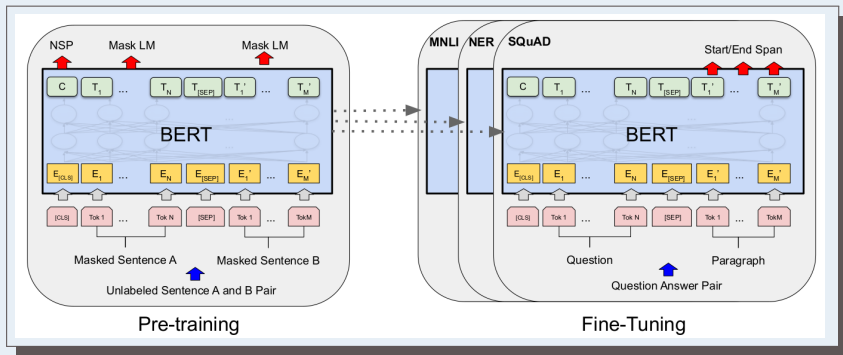
# Conclusion 2-stage training



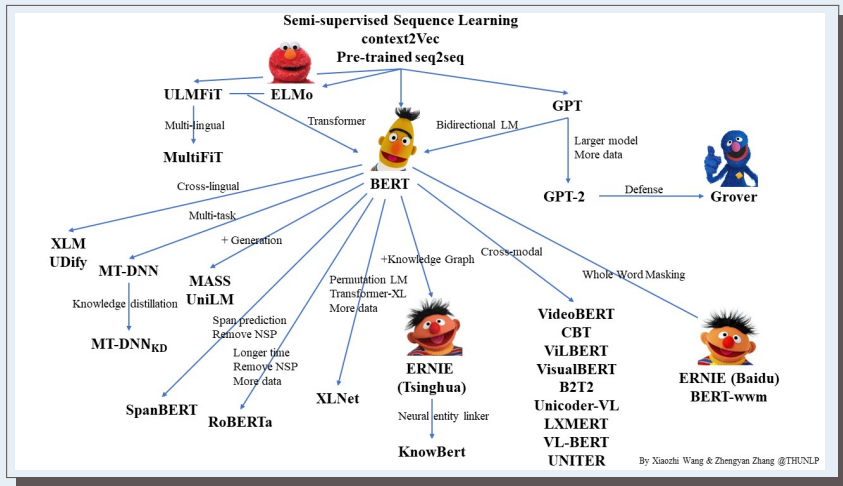Figure: **2-stage training**

Figure: **BERT Family**

# BERT-Family (2)

## SpanBERT - EMNLP2019

Improving Pre-training by Representing and Predicting Spans

## Roberta - ICLR2020 Reject

A Robustly Optimized BERT Pretraining Approach

## ALBERT - ICLR2020 Spotlight

A Lite BERT for Self-supervised Learning of Language Representations

# SpanBERT, Roberta, ALBERT

## Difference with BERT - SpanBERT

1. change NSP task to SBO (span boundary objective) task.
2. mask random consecutive spans.

## Difference with BERT - Roberta

1. More big batch size when training bert.
2. More optimal hyperparameters.
3. More data.
4. Remove NSP task.

## Difference with BERT - ALBERT

1. change NSP task to SOP (Sentence Order Prediction) task.
2. One layer recall 12 times to build 12 layers. (share parameters)
3. factorize embedding table. initize small and project to big dimension.