# A few words about BERT
# (or what's up with NLP field)
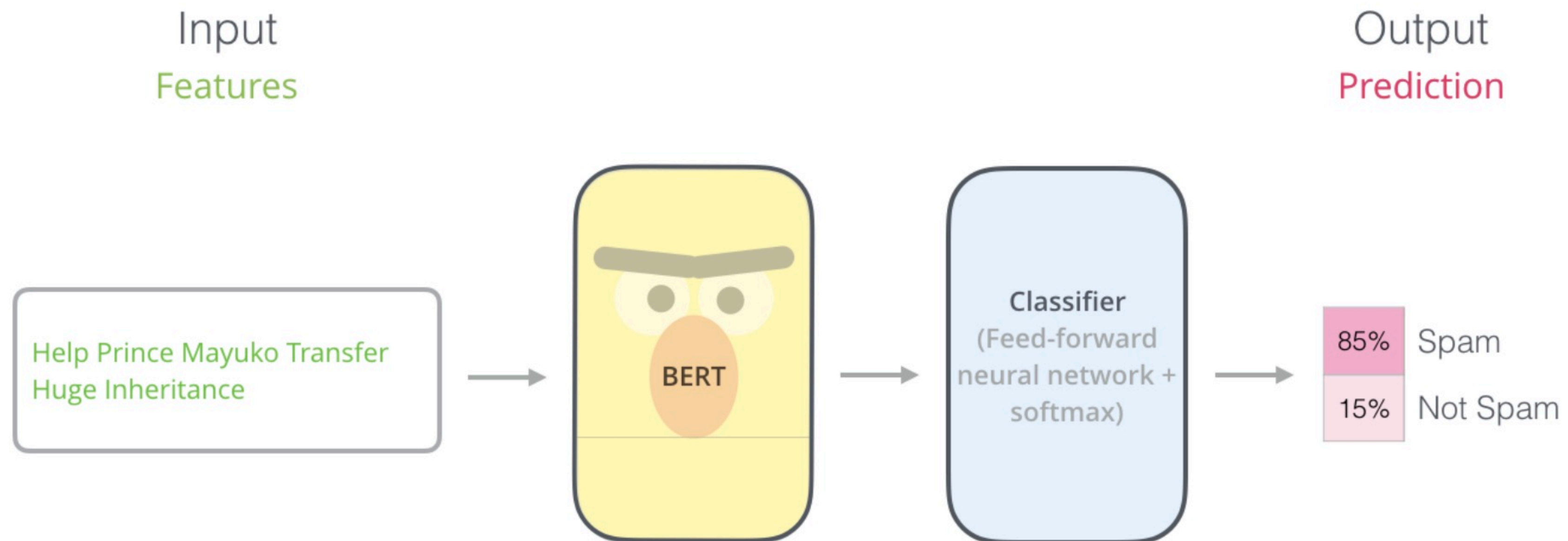
based on the original article https://arxiv.org/abs/1810.04805
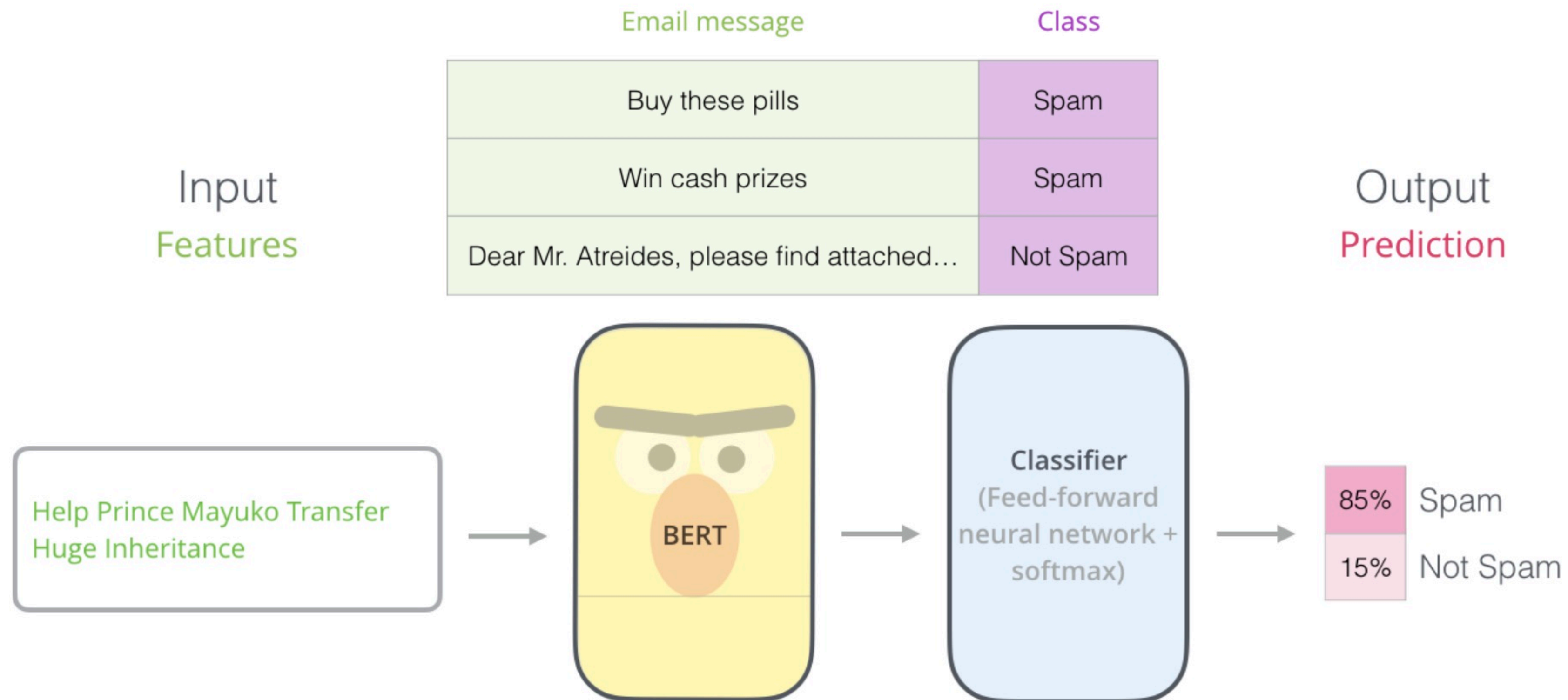
Sochi, 2019

# Outline

- BERT Model

- Benchmarks for BERT

- Details of BERT Model

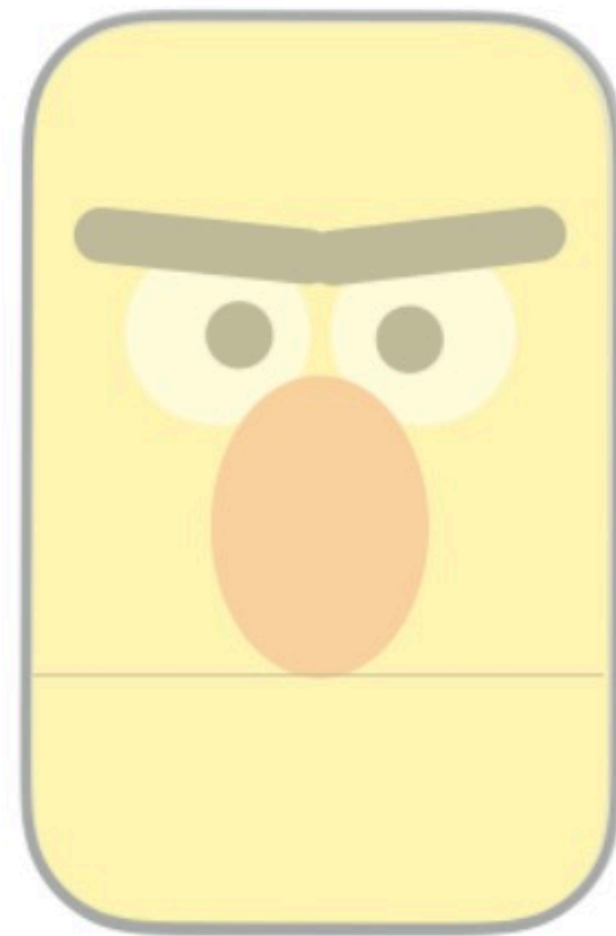- Training framework

- BERT vs XLNet

- RoBERTa by FAIR

https://arxiv.org/abs/1810.04805

# BERT Model



Source: http://jalammar.github.io/illustrated-bert/

# BERT Model



Source: http://jalammar.github.io/illustrated-bert/

# Benchmarks for BERT

| System | MNLI-(m/mm) | QQP | QNLI | SST-2 | CoLA | STS-B | MRPC | RTE | **Average** |
|---|---|---|---|---|---|---|---|---|---|
| | 392k | 363k | 108k | 67k | 8.5k | 5.7k | 3.5k | 2.5k | - |
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.9 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 88.1 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.2 |
| BERT$_{BASE}$ | 84.6/83.4 | 71.2 | 90.1 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **91.1** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **81.9** |

- **QQP** - Quora Question Pairs

- **MRPC** - Microsoft Research Paraphrase Corpus

- ...

https://arxiv.org/abs/1810.04805

4

# Details of BERT Model



BERT<sub>BASE</sub>

BERT<sub>LARGE</sub>

Source: http://jalammar.github.io/illustrated-bert/

# Details of BERT Model

With total parameters of **~340M**

With total parameters of **~110M**



Source: http://jalammar.github.io/illustrated-bert/

# Details of BERT Model

Use the output of the masked word's position to predict the masked word

Possible classes: All English words

| | |
|---|---|
| 0.1% | Aardvark |
| … | … |
| 10% | Improvisation |
| … | … |
| 0% | Zyzzyva |

FFNN + Softmax

1  2  3  4  5  6  7  8  •••  512

BERT

Randomly mask 15% of tokens

1  2  3  4  5  6  7  8  •••  512

[CLS]  Let's  stick  to  [MASK]  in  this  skit

Input

[CLS]  Let's  stick  to improvisation in  this  skit

Source: http://jalammar.github.io/illustrated-bert/
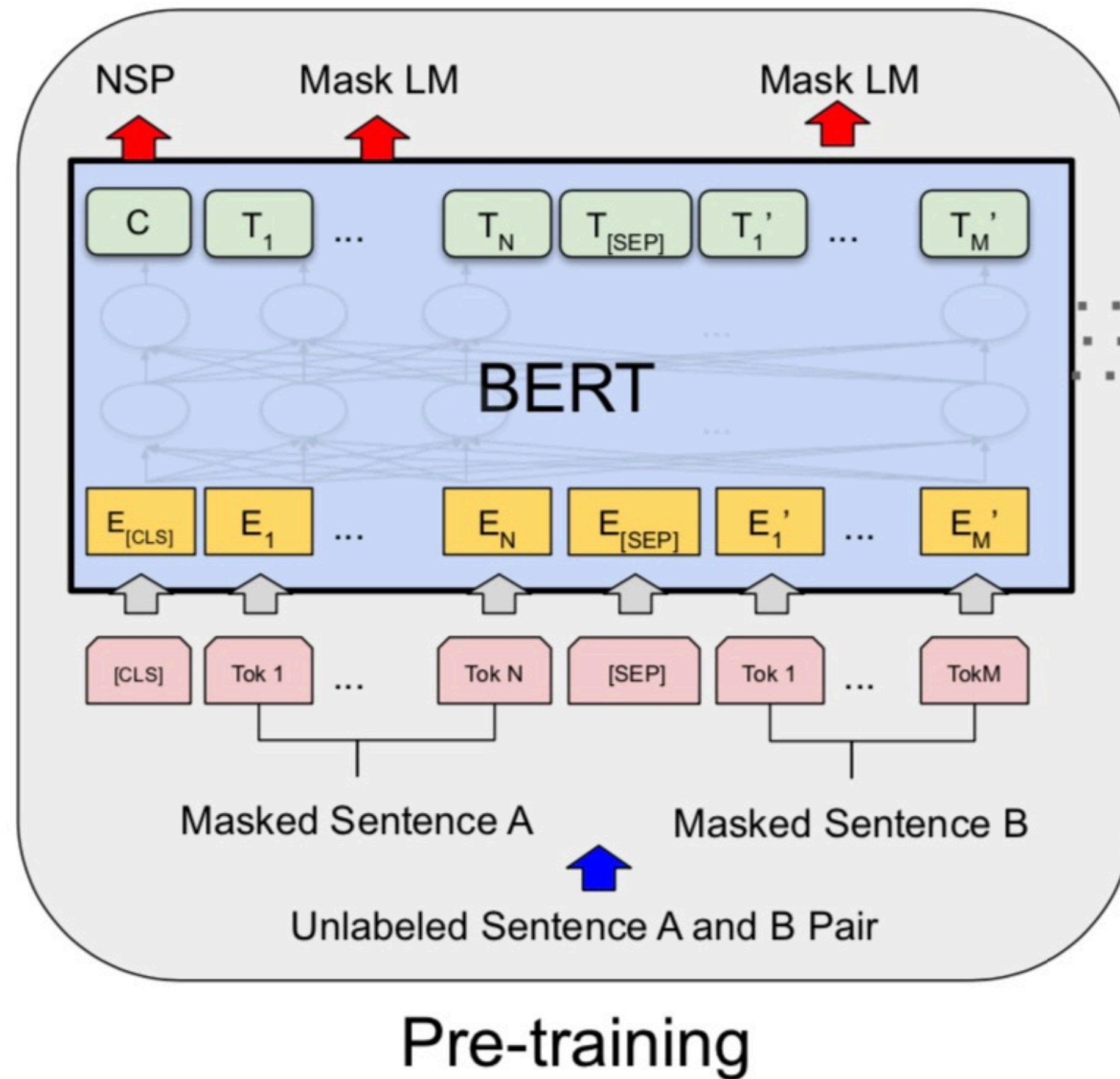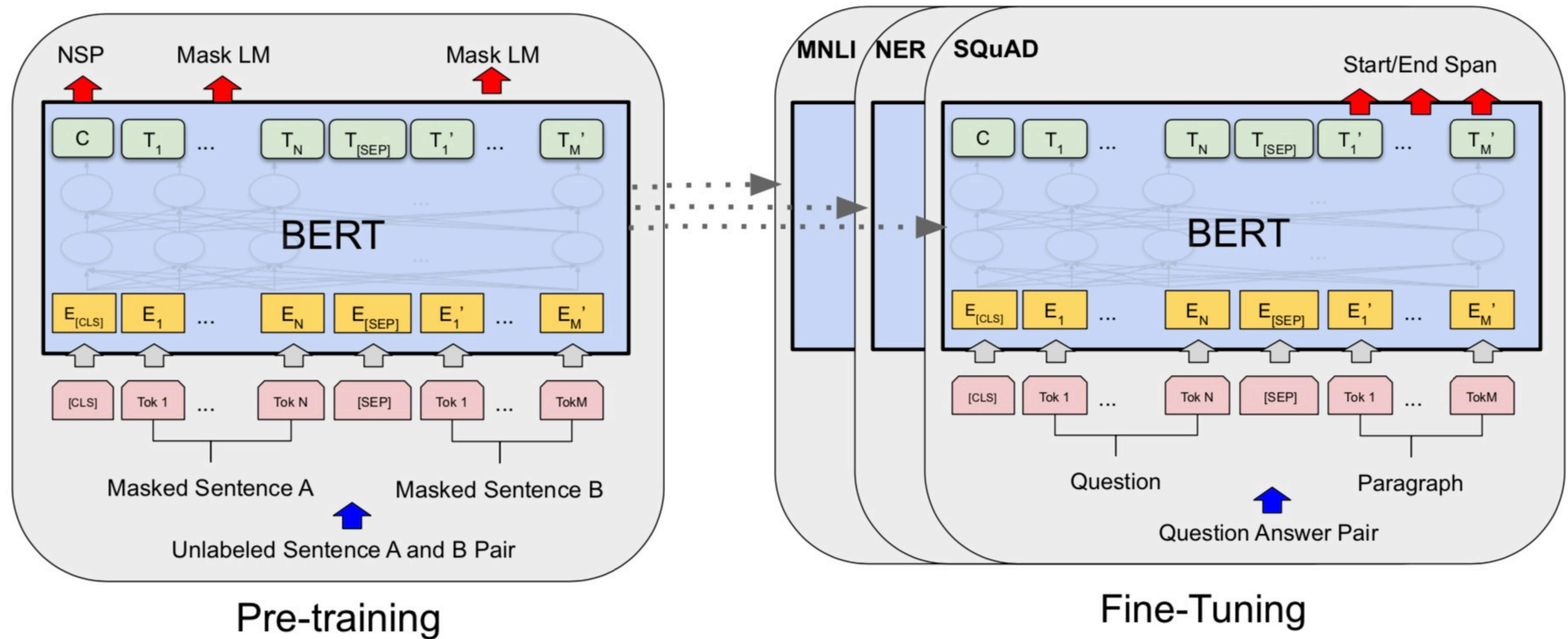
7

# Training framework

- **Pre-training.** Model is trained on unlabelled data over different pre-training tasks.

- **Fine-tuning.** Firstly, model is initialised with the pre-trained parameters. After that all the parameters are fine-tuned using labelled data from the downstream tasks.
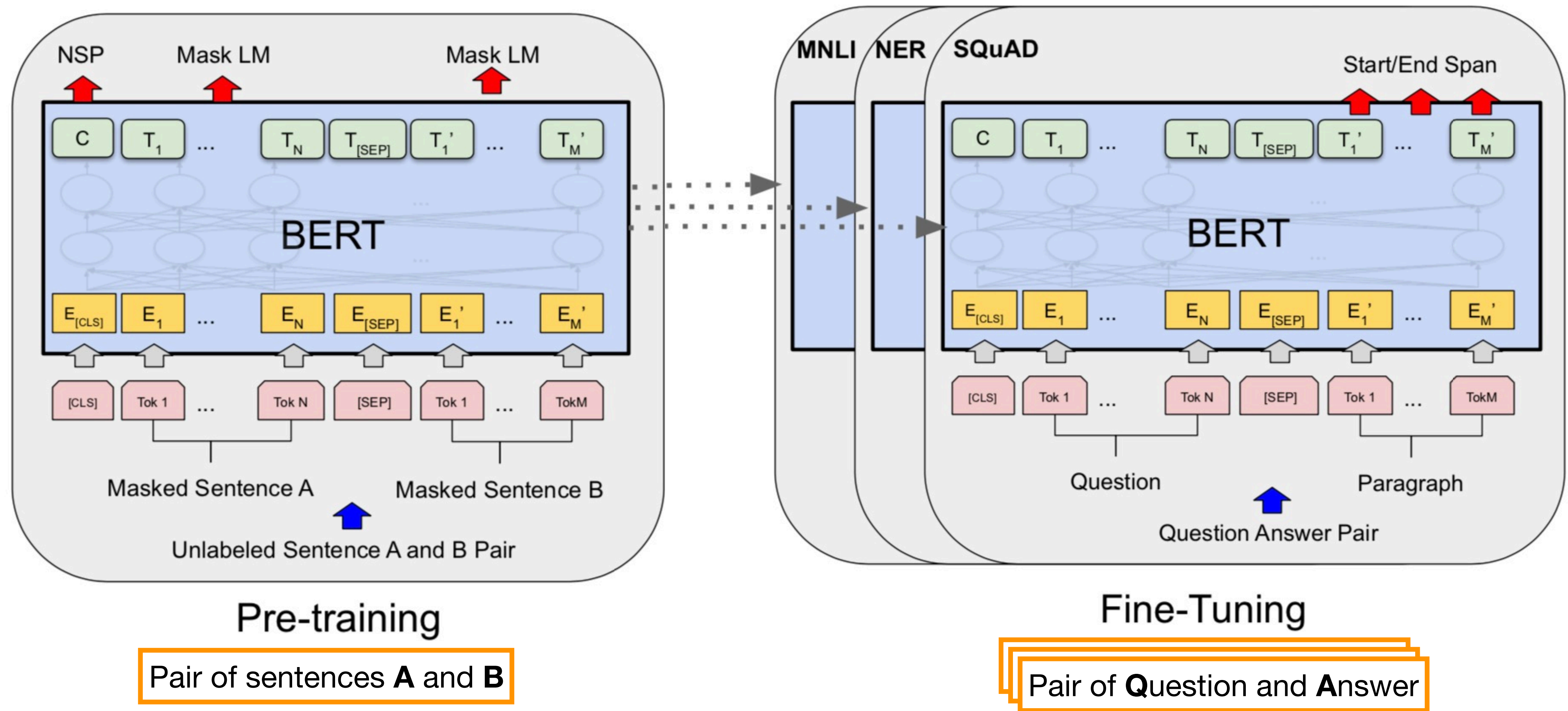
https://arxiv.org/abs/1810.04805

# Training framework



Pre-training

https://arxiv.org/abs/1810.04805

# Training framework

# Training framework



Pre-training

Fine-Tuning

Pair of sentences **A** and **B**

Pair of **Q**uestion and **A**nswer

9

https://arxiv.org/abs/1810.04805

# Training framework



Pre-training

Fine-Tuning

Pair of **Q**uestion and **A**nswer

https://arxiv.org/abs/1810.04805

9

# BERT vs XLNet

- BERT - **B**idirectional **E**ncoder **R**epresentations from **T**ransformers (looses connection between words)

- XLNet - AutoRegressive Language Modelling (word order does count)

# BERT vs XLNet



| Model | MNLI | QNLI | QQP | RTE | SST-2 | MRPC | CoLA | STS-B | WNLI |
|---|---|---|---|---|---|---|---|---|---|
| *Single-task single models on dev* | | | | | | | | | |
| BERT [2] | 86.6/- | 92.3 | 91.3 | 70.4 | 93.2 | 88.0 | 60.6 | 90.0 | - |
| XLNet | **89.8/-** | **93.9** | **91.8** | **83.8** | **95.6** | **89.2** | **63.6** | **91.8** | - |
| *Single-task single models on test* | | | | | | | | | |
| BERT [10] | 86.7/85.9 | 91.1 | 89.3 | 70.1 | 94.9 | 89.3 | 60.5 | 87.6 | 65.1 |
| *Multi-task ensembles on test (from leaderboard as of June 19, 2019)* | | | | | | | | | |
| Snorkel* [29] | 87.6/87.2 | 93.9 | 89.9 | 80.9 | 96.2 | 91.5 | 63.8 | 90.1 | 65.1 |
| ALICE* | 88.2/87.9 | 95.7 | **90.7** | 83.5 | 95.2 | 92.6 | **68.6** | 91.1 | 80.8 |
| MT-DNN* [18] | 87.9/87.4 | 96.0 | 89.9 | **86.3** | 96.5 | 92.7 | 68.4 | 91.1 | 89.0 |
| XLNet* | **90.2/89.7**$^{\dagger}$ | **98.6**$^{\dagger}$ | 90.3$^{\dagger}$ | **86.3** | **96.8**$^{\dagger}$ | **93.0** | 67.8 | **91.6** | **90.4** |

Source: https://arxiv.org/abs/1906.08237

11

# SQuAD 2.0
## The Stanford Question Answering Dataset

## What is SQuAD?

**S**tanford **Qu**estion **A**nswering **D**ataset (SQuAD) is a reading comprehension dataset, consisting of questions posed by crowdworkers on a set of Wikipedia articles, where the answer to every question is a segment of text, or *span*, from the corresponding reading passage, or the question might be unanswerable.

**New** **SQuAD2.0** combines the 100,000 questions in SQuAD1.1 with over 50,000 new, unanswerable questions written adversarially by crowdworkers to look similar to answerable ones. To do well on SQuAD2.0, systems must not only answer questions when possible, but also determine when no answer is supported by the paragraph and abstain from answering. SQuAD2.0 is a challenging natural language understanding task for existing models, and we release SQuAD2.0 to the community as the successor to SQuAD1.1. We are optimistic that this new dataset will encourage the development of reading comprehension systems that know what they don't know.

**Explore SQuAD2.0 and model predictions**

## Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph. How will your system compare to humans on this task?

| Rank | Model | EM | F1 |
|------|-------|-----|-----|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Mar 20, 2019 | BERT + DAE + AoA (ensemble)<br>*Joint Laboratory of HIT and iFLYTEK Research* | **87.147** | **89.474** |
| 2<br>Mar 15, 2019 | BERT + ConvLSTM + MTL + Verifier (ensemble)<br>*Layer 6 AI* | 86.730 | 89.286 |
| 3<br>Mar 05, 2019 | BERT + N-Gram Masking + Synthetic Self-Training (ensemble)<br>*Google AI Language*<br>https://github.com/google-research/bert | 86.673 | 89.147 |
| 4<br>May 21, 2019 | XLNet (single model)<br>*Google Brain & CMU* | 86.346 | 89.133 |

12

Source: https://rajpurkar.github.io/SQuAD-explorer/

# RoBERTa by FAIR

# Summary

- BERT Model

- Benchmarks for BERT

- Details of BERT Model

- Training framework

- BERT vs XLNet

- RoBERTa by FAIR

# References

- The Illustrated BERT, ELMo, and co.: http://jalammar.github.io/illustrated-bert/

- About BERT in Google AI Blog: https://ai.googleblog.com/2018/11/open-sourcing-bert-state-of-art-pre.html

- SQuAD 2.0: https://rajpurkar.github.io/SQuAD-explorer/

- GLUE Benchmark: https://gluebenchmark.com/leaderboard

https://arxiv.org/abs/1810.04805