



# Fine-Tuned Sentence Transformer Model for Question Answering Task

**Ercong Nie**

e.nie@campus.lmu.de

**Master Student**

Center for Information and Language Processing, Faculty of Languages and Literatures, LMU of Munich

May 26, 2022





# Contents

Introduction

Transformer-Based Language Models

An Application of Sentence Transformers in Question Answering Task

Summary

References











## Different Types of Language Models

- **N-Gram Language Models**
- **Neural Network Based Language Models**
  - Predicting the next word in a sequence by its previous words based on neural network structures, like **RNN**, **LSTM**...
  - Adopting a more powerful method to represent words, i.e. **the word embedding**, a semantically meaningful vector.
- **Transformer-Based Language Models**

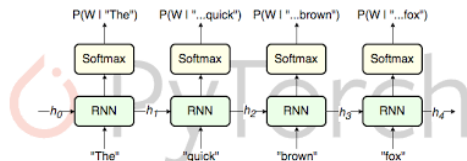


Figure 1: An Example of RNN language model structure<sup>a</sup>

<sup>a</sup>[https://medium.com/@florijan.stamenkovic\\_99541/](https://medium.com/@florijan.stamenkovic_99541/)



## Different Types of Language Models

- N-Gram Language Models
- Neural Network Based Language Models
- Transformer-Based Language Models
  - **Transformer**: A network structure based on the **self-attention** mechanism.
  - An improvement compared with NN-based models: A sequence processing method that eliminates recurrent connections.

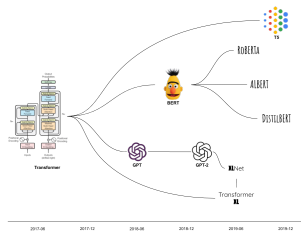


Figure 2: Collection of some transformer-based models<sup>a</sup>

<sup>a</sup><https://www.factored.ai/2021/09/21/>



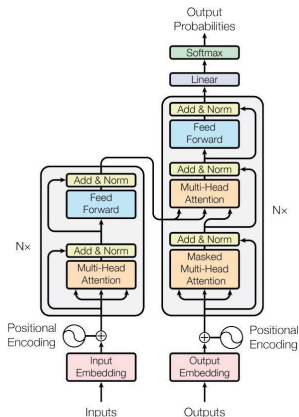


## Transformer-Based Language Models



## Structure of Transformer Model<sup>3</sup>

- Composed of two parts, **Encoder** (left) and **Decoder** (right).
- Both parts contain modules that can be stacked on top of each other multiple times.
- Core module parts: Multi-head attention and feed forward layers.
- Inputs and Outputs first embedded into an n-dim space.
- Positional embedding offers positional information.





## Structure of BERT

- BERT: Bidirectional Encoder Representations from Transformers
- Has two pre-training tasks:
  1. **Masked Language Modeling (MLM)**
  2. **Next Sentence Prediction (NSP)**
- An autoencoder model, self-supervised training objectives (no need for labeled data in pre-training)

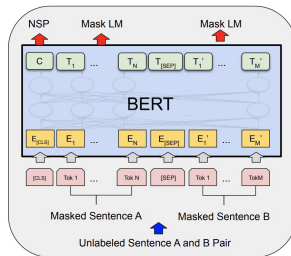


Figure 3: The structure of BERT<sup>ab</sup>

<sup>a</sup><https://paperswithcode.com/lib/allennlp>

<sup>b</sup>Devlin et al. (2016)



## Sizes of Transformer-Based Models

Year	Model	# of Parameters	Dataset Size
2019	BERT [39]	3.4E+08	16GB
2019	DistilBERT [113]	6.60E+07	16GB
2019	ALBERT [70]	2.23E+08	16GB
2019	XLNet (Large) [150]	3.40E+08	126GB
2020	ERNIE-GEN (Large) [145]	3.40E+08	16GB
2019	RoBERTa (Large) [74]	3.55E+08	161GB
2019	MegatronLM [122]	8.30E+09	174GB
2020	T5-11B [107]	1.10E+10	745GB
2020	T-NLG [112]	1.70E+10	174GB
2020	GPT-3 [25]	1.75E+11	570GB
2020	GShard [73]	6.00E+11	–
2021	Switch-C [43]	1.57E+12	745GB

Figure 4: Overview of sizes of recent large language models<sup>4</sup>

<sup>4</sup>Bender et al. (2021)



## Pre-Training Fine-Tuning Paradigm

- First **pre-training** model on large **general raw** language data.
- Then **fine-tuning** the model's parameters by using relatively small amount of **labeled** data so that the model can be applied in a **specific** field or task.
- Basically a **transfer learning** method.
- Most transform-based models adopt this paradigm, so they are also called **Pretrained Language Models (PLMs)**.

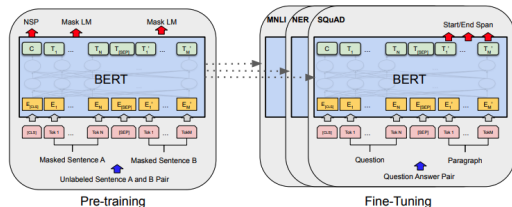


Figure 5: Pre-training Fine-tuning Paradigm  
Devlin et al. (2016)





## Structure of SBERT

- Add a **pooling** operation to the output of BERT.
- Derive a **fixed** size sentence embedding.
- Create siamese and triplet networks to update the weights in the fine-tuning.

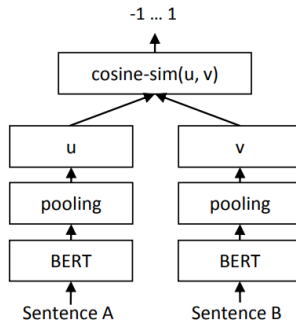


Figure 6: SBERT Architecture. Reimers and Gurevych (2019)



## An Application of Sentence Transformers in Question Answering Task









## Dataset Introduction I

- **Dataset:** **InsuranceQA**, created by IBM Watson researchers.
- **Contents:** Questions and answers in the **insurance** domain, original data from the website Insurance Library.
- Composed of training and evaluation parts. Each question has a **pool** of **100** possible answers, some of which are true answers (ground truth).



## Dataset Introduction II

- **Data Structure:** <Domain><TAB><QUESTION><TAB><Groundtruth><TAB><Pool>  
<QUESTION> is represented by the tokens in the form of **index**. <Groundtruth> and <POOL> are represented by the **labels** of answer.

```
renters-insurance  idx_1285 idx_1010 idx_999 idx_136 idx_65807 22542 4380 2235 26739 24916 17855 3406 21201 70 19553 22008
5768 18118 2105 20821 2316 25072 6805 9846 24262 6317 6250 13690 1467 4770 10917 18784 7229 8030 26792 15729 13179 2464 3884
23403 24493 6292 24533 26 13983 15294 26821 22449 20057 10641 16504 8153 14453 6276 16349 27141 14698 13650 12175 836 17050
15911 5230 10139 5955 18903 2710 471 20269 11888 17073 16128 8026 26441 19693 11405 7745 14596 13353 22175 5530 17982 10727
18225 12703 8782 26875 16985 1324 16967 25308 7420 17279 11137 26769 22418 18674 21014 14694 20737 22135 3346 7342 5099 12388
20032 12029
```

Figure 7: Example of a data item

- 12,889 questions in training set, 4,000 questions in evaluation set, 27,413 answers totally.



## Data Preparation for Fine-Tuning

- Convert the original training set into the form of sentence pairs with a label.
- For each question in the original training set, create two sentence pairs as the inputs for model fine-tuning.
  1. Select the first answer from the <Groundtruth> to compose the **positive** sentence pair (labeling 1).
  2. Randomly select one answer from the <POOL> to compose the **negative** pair (labeling 2).











## Evaluation and Result II

<b>Pre-trained Model without Fine-Tuning</b>	
<b>Name</b>	paraphrase-TinyBERT-L6-v2
<b>Accuracy.</b>	0.2664
<b>Runtime</b>	5069.734 s
<b>Pre-trained Model with Fine-Tuning</b>	
<b>Name</b>	paraphrase-TinyBERT-L6-v2
<b>Accuracy.</b>	0.4867
<b>Runtime</b>	5651.248 s

Table 2: Result of fine-tuned Sentence Transformer applied to the insurance QA task.

**Result** By fine-tuning with the domain data, the accuracy of the QA system has been significantly improved, rising from **0.2664** to **0.4867** with an increasing of 82.7%.



# Summary





## References I

- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big?. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2016). Bert: Bidirectional encoder representations from transformers.
- Feng, M., Xiang, B., Glass, M. R., Wang, L., and Zhou, B. (2015). Applying deep learning to answer selection: A study and an open task. In *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, pages 813–820. IEEE.





