# Seminar II: Deep Learning-based Natural Language Processing

**"BERT: Pre-training of Deep Bidirectional Transformers
for Language Understanding"
Paper review**

# Contents

1. Introduction
   - Contextualized word representations - ELMo
   - OpenAI GPT-1
   - Key Contribution

2. Methodology
   - BERT: Masked Language model
   - BERT: Next sentence prediction
   - BERT: Input representation
   - BERT: Model and training details

3. Results and Discussion

4. Conclusions

# Introduction. Contextualized word representations - ELMo
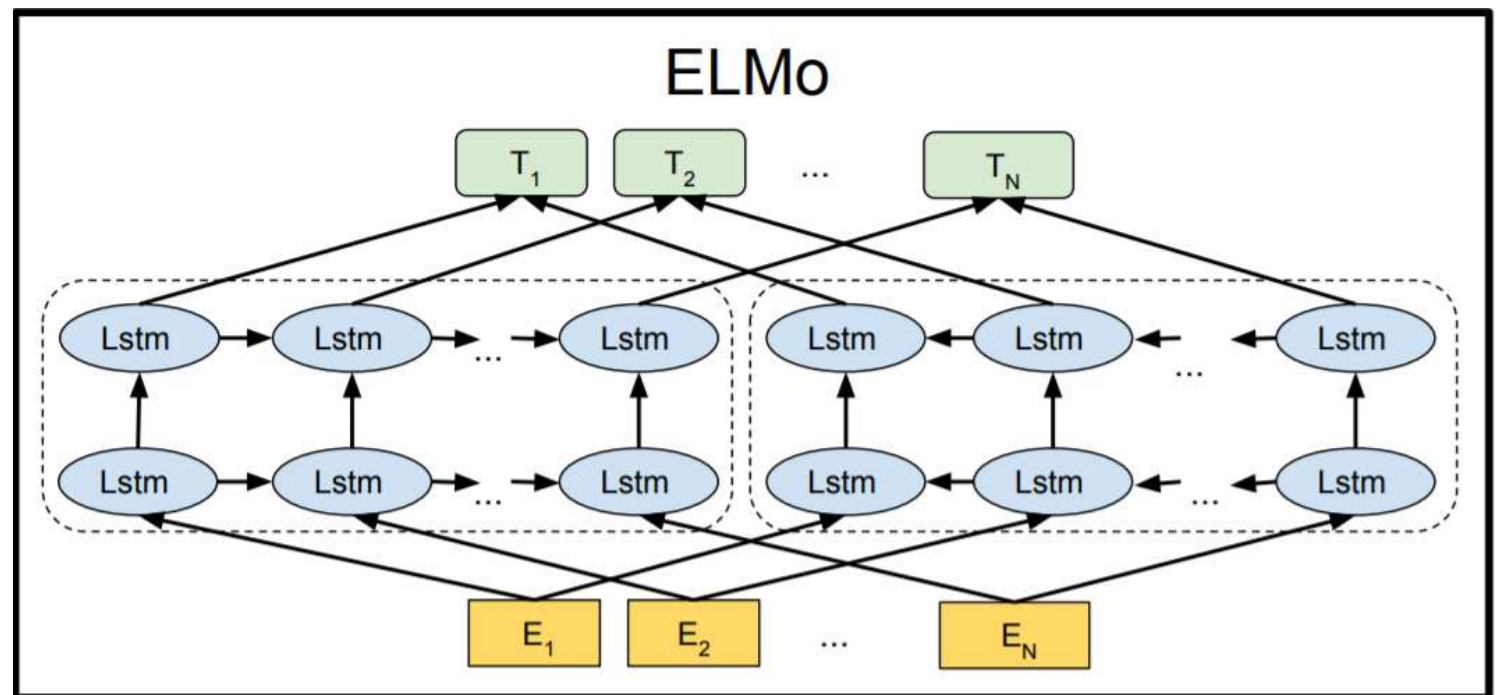
Words have several meanings based on the context.



*Contextualized word-embeddings can give words different embeddings based on the meaning they carry in the context of the sentence.*
*Source: https://jalammar.github.io/illustrated-bert/*

# Introduction. Contextualized word representations - ELMo
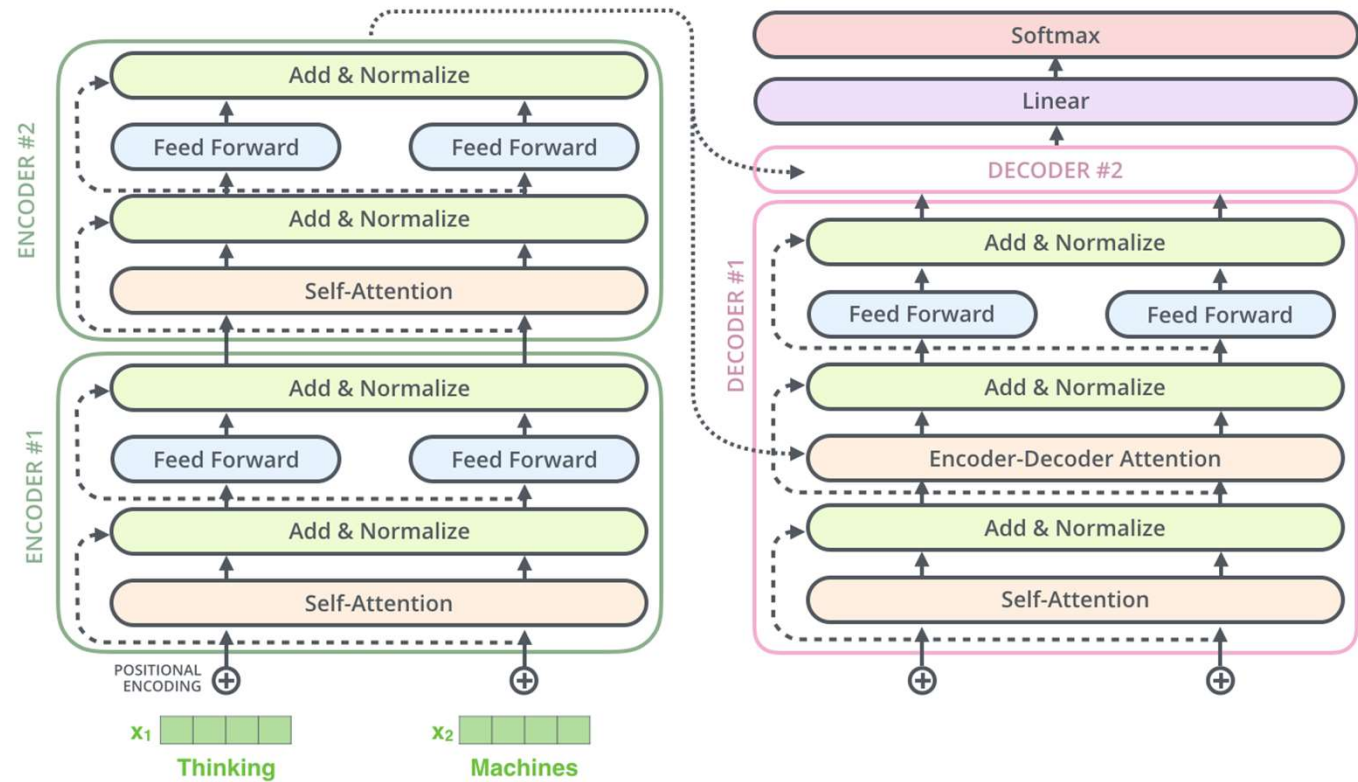
Bi-directional language modeling



*Bi-directional stacked LSTM.*
*Source: https://medium.com/saarthi-ai/elmo-for-contextual-word-embedding-for-text-classification-24c9693b0045*

# Introduction. OpenAI GPT-1

Transformer decoder without Enc-Dec attention



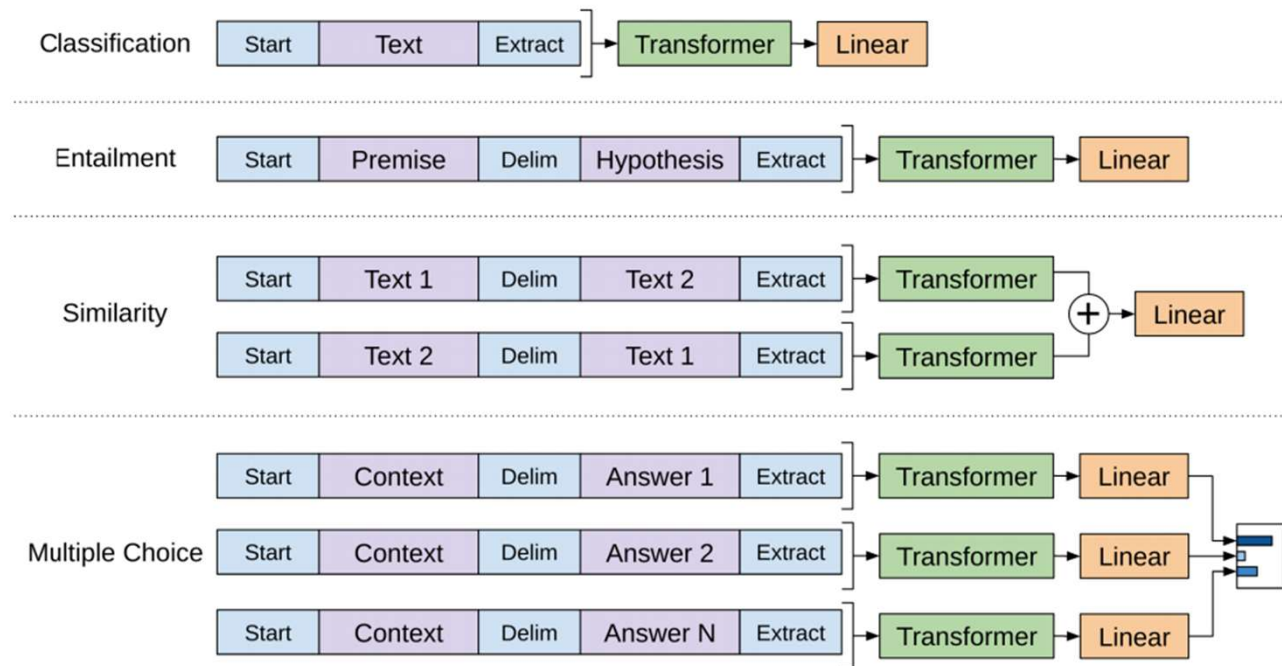Source: https://jalammar.github.io/illustrated-transformer/

# Introduction. OpenAI GPT-1

Pre-train on BooksCorpus

      - 12 layers, 768 hidden size, 12 attention heads (110M parameters)



*Source: Improving Language Understanding by Generative Pre-Training*

# Introduction. OpenAI GPT-1

| Method | MNLI-m | MNLI-mm | SNLI | SciTail | QNLI | RTE |
|---|---|---|---|---|---|---|
| ESIM + ELMo [44] (5x) | - | - | 89.3 | - | - | - |
| CAFE [58] (5x) | 80.2 | 79.0 | 89.3 | - | - | - |
| Stochastic Answer Network [35] (3x) | 80.6 | 80.1 | - | - | - | - |
| CAFE [58] | 78.7 | 77.9 | 88.5 | 83.3 | | |
| GenSen [64] | 71.4 | 71.3 | - | - | 82.3 | 59.2 |
| Multi-task BiLSTM + Attn [64] | 72.2 | 72.1 | - | - | 82.1 | **61.7** |
| Finetuned Transformer LM (ours) | **82.1** | **81.4** | **89.9** | **88.3** | **88.1** | 56.0 |

*Experimental results on natural language inference tasks*
*Source: Improving Language Understanding by Generative Pre-Training*

| Method | Story Cloze | RACE-m | RACE-h | RACE |
|---|---|---|---|---|
| val-LS-skip [55] | 76.5 | - | - | - |
| Hidden Coherence Model [7] | 77.6 | - | - | - |
| Dynamic Fusion Net [67] (9x) | - | 55.6 | 49.4 | 51.2 |
| BiAttention MRU [59] (9x) | - | 60.2 | 50.3 | 53.3 |
| Finetuned Transformer LM (ours) | **86.5** | **62.9** | **57.4** | **59.0** |

*Results on question answering and common-sense reasoning*
*Source: Improving Language Understanding by Generative Pre-Training*

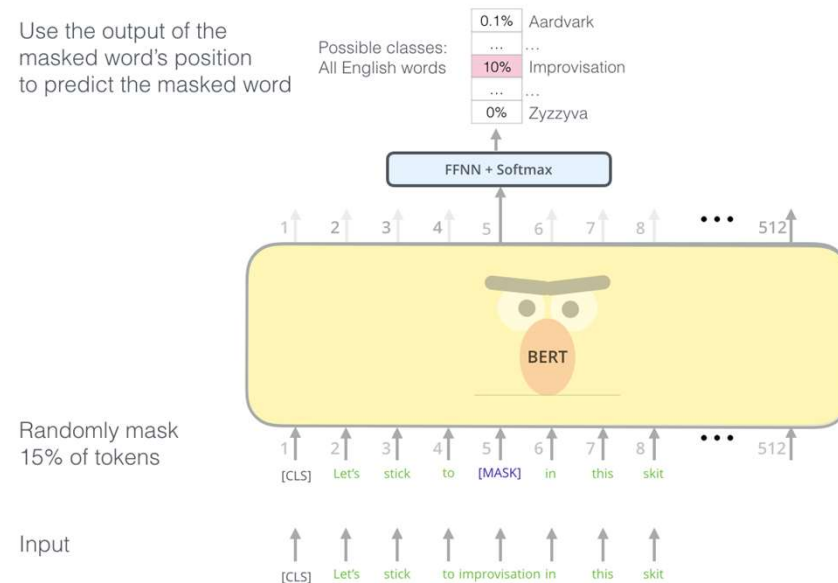# Introduction. Key Contribution

**OBJECTIVE**

To improve **the fine-tuning based LM approaches**
   **1. Bidirectional approach – "Masked LM" (MLM) <- ELMo**
   **2. Transformer – "GPT"**

# Methodology – BERT: Masked Language model

**Solution:** Mask out k% of the input words, and then predict the masked words
- We always use k = 15%



*BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word*
*Source: https://jalammar.github.io/illustrated-bert/*
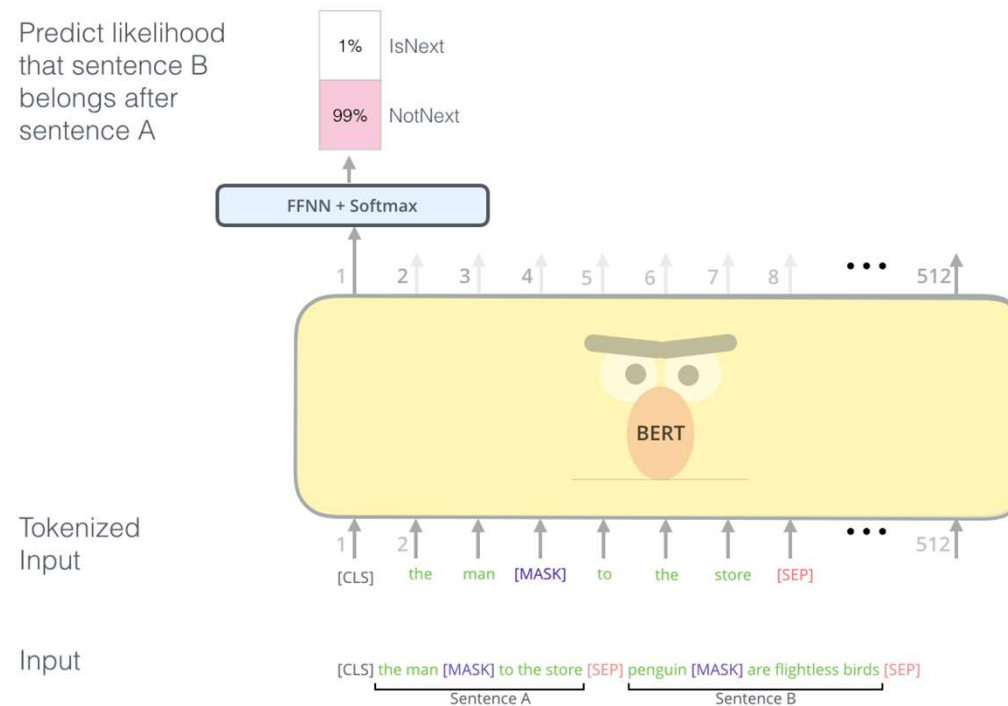
# Methodology – BERT: Masked Language model

**Problem:** Mask token never seen at fine-tuning

**Solution:** 15% of the words are *[mask]* but do not always replace with it. Instead:

- 80% of the time, replace with *[mask]*

went to the store -> went to the *[mask]*

- 10% of the time, replace random word

went to the store -> went to the *running*

- 10% of the time, keep same

went to the store -> went to the *store*

# Methodology – BERT: Next sentence prediction

Helpful for some downstream tasks



*The second task BERT is pre-trained on is a two-sentence classification task. The tokenization is oversimplified in this graphic as BERT actually uses WordPieces as tokens rather than words --- so some words are broken down into smaller chunks.*
*Source: https://jalammar.github.io/illustrated-bert/*

# Methodology – BERT: Input representation
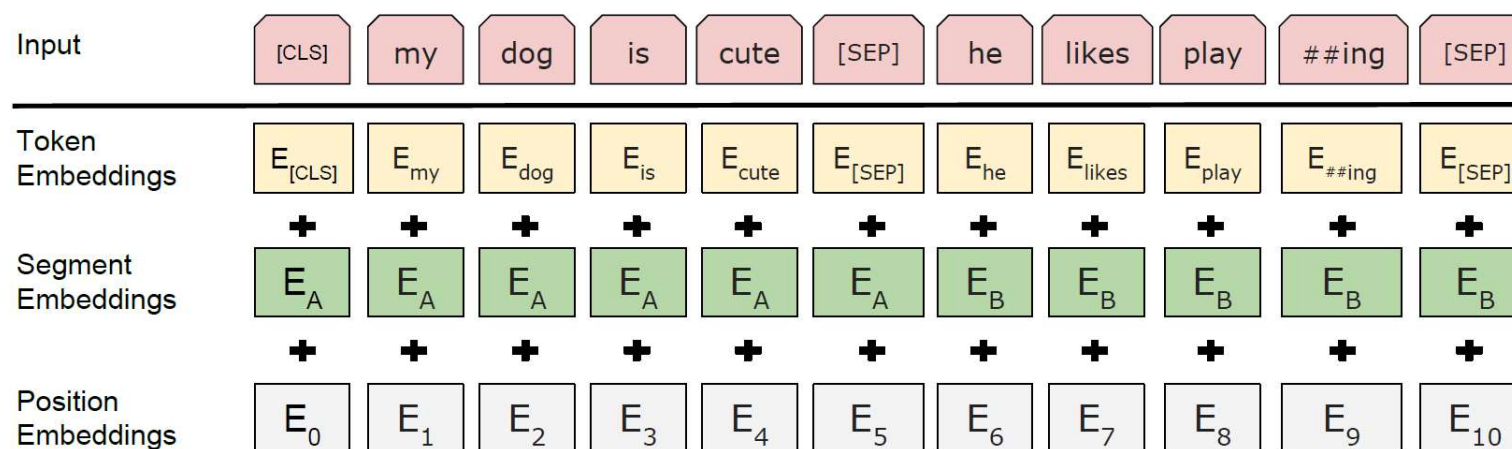
Use 30,000 WordPiece vocabulary on input

Each token is sum of three embeddings

Trainable Segment Embedding

CLS: Special tokens representing all input



*BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings and the position embeddings.*
*Source: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*

# Methodology – BERT: Model and training details

**Pre-Train:**

Data: Wikipedia (2.5B words) + BookCorpus (800M words)

BERT-base (Layer=12, Hidden=768, Head=12, Total Parameters= 110M)

BERT-large (Layer=24, Hidden=1024, Head=16, Total Parameters=340M)

**Fine-Tune:**

Different input/output

for different downstream task



(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG

(b) Single Sentence Classification Tasks:
SST-2, CoLA

(c) Question Answering Tasks:
SQuAD v1.1

(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

13

# Results and Discussions

| System | MNLI-(m/mm) 392k | QQP 363k | QNLI 108k | SST-2 67k | CoLA 8.5k | STS-B 5.7k | MRPC 3.5k | RTE 2.5k | Average - |
|---|---|---|---|---|---|---|---|---|---|
| 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 | 79.8 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 87.4 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.1 |
| BERT$_{BASE}$ | 84.6/83.4 | 71.2 | 90.5 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **92.7** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **82.1** |

GLUE Test results (Collection of NLP tasks)
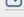*Source: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*

| Rank | Name | Model | URL | Score | CoLA | SST-2 | MRPC | STS-B | QQP | MNLI-m | MNLI-mm | QNLI | RTE | WNLI | AX |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | JDExplore d-team | Vega v1 | | 91.3 | 73.8 | 97.9 | 94.5/92.6 | 93.5/93.1 | 76.7/91.1 | 92.1 | 91.9 | 96.7 | 92.4 | 97.9 | 51.4 |
| 2 | Microsoft Alexander v-team | Turing NLR v5 | ⬈ | 91.2 | 72.6 | 97.6 | 93.8/91.7 | 93.7/93.3 | 76.4/91.1 | 92.6 | 92.4 | 97.9 | 94.1 | 95.9 | 57.0 |
| 3 | DIRL Team | DeBERTa + CLEVER | | 91.1 | 74.7 | 97.6 | 93.3/91.1 | 93.4/93.1 | 76.5/91.0 | 92.1 | 91.8 | 96.7 | 93.2 | 96.6 | 53.3 |
| 4 | ERNIE Team - Baidu | ERNIE | ⬈ | 91.1 | 75.5 | 97.8 | 93.9/91.8 | 93.0/92.6 | 75.2/90.9 | 92.3 | 91.7 | 97.3 | 92.6 | 95.9 | 51.7 |
| 5 | AliceMind & DIRL | StructBERT + CLEVER | ⬈ | 91.0 | 75.3 | 97.7 | 93.9/91.9 | 93.5/93.1 | 75.6/90.8 | 91.7 | 91.5 | 97.4 | 92.5 | 95.2 | 49.1 |
| 6 | DeBERTa Team - Microsoft | DeBERTa / TuringNLRv4 | ⬈ | 90.8 | 71.5 | 97.5 | 94.0/92.0 | 92.9/92.6 | 76.2/90.8 | 91.9 | 91.6 | 99.2 | 93.2 | 94.5 | 53.2 |
| 7 | HFL iFLYTEK | MacALBERT + DKM | | 90.7 | 74.8 | 97.0 | 94.5/92.6 | 92.8/92.6 | 74.7/90.6 | 91.3 | 91.1 | 97.8 | 92.0 | 94.5 | 52.6 |
| 8 | PING-AN Omni-Sinitic | ALBERT + DAAF + NAS | | 90.6 | 73.5 | 97.2 | 94.0/92.0 | 93.0/92.4 | 76.1/91.0 | 91.6 | 91.3 | 97.5 | 91.7 | 94.5 | 51.2 |
| 9 | T5 Team - Google | T5 | ⬈ | 90.3 | 71.6 | 97.5 | 92.8/90.4 | 93.1/92.8 | 75.1/90.6 | 92.2 | 91.9 | 96.9 | 92.8 | 94.5 | 53.1 |
| 10 | Microsoft D365 AI & MSR AI & GATECH | MT-DNN-SMART | ⬈ | 89.9 | 69.5 | 97.5 | 93.7/91.6 | 92.9/92.5 | 73.9/90.2 | 91.0 | 90.8 | 99.2 | 89.7 | 94.5 | 50.2 |

GLUE Benchmark Leaderboard
*Source: https://gluebenchmark.com/leaderboard*

# Results and Discussions

| System | Dev | | Test | |
|---|---|---|---|---|
| | EM | F1 | EM | F1 |
| Top Leaderboard Systems (Dec 10th, 2018) | | | | |
| Human | - | - | 82.3 | 91.2 |
| #1 Ensemble - nlnet | - | - | 86.0 | 91.7 |
| #2 Ensemble - QANet | - | - | 84.5 | 90.5 |
| Published | | | | |
| BiDAF+ELMo (Single) | - | 85.6 | - | 85.8 |
| R.M. Reader (Ensemble) | 81.2 | 87.9 | 82.3 | 88.5 |
| Ours | | | | |
| BERT$_{BASE}$ (Single) | 80.8 | 88.5 | - | - |
| BERT$_{LARGE}$ (Single) | 84.1 | 90.9 | - | - |
| BERT$_{LARGE}$ (Ensemble) | 85.8 | 91.8 | - | - |
| BERT$_{LARGE}$ (Sgl.+TriviaQA) | **84.2** | **91.1** | **85.1** | **91.8** |
| BERT$_{LARGE}$ (Ens.+TriviaQA) | **86.2** | **92.2** | **87.4** | **93.2** |

SQuAD 1.1 results
*Source: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*

**Paragraph:** The scientific revolution was a period when European ideas in classical Physics, Astronomy, Biology, Human Anatomy, Chemistry, and other classical sciences were rejected and led to doctrines supplanting those that had prevailed from ancient Greece to the middle ages which would lead to a transition to modern science. this period saw a fundamental transformation in scientific ideas across Physics, Astronomy, and Biology, in institutions supporting scientific investigation, and in the more widely held picture of the universe. individuals started to question all manners of things and it was this questioning that led to the scientific revolution, which in turn formed the foundations of contemporary sciences and the establishment of several modern scientific fields.

**Question:** What did the scientific revolution cause?

**Answer:** a transition to modern science

| System | Dev | | Test | |
|---|---|---|---|---|
| | EM | F1 | EM | F1 |
| Top Leaderboard Systems (Dec 10th, 2018) | | | | |
| Human | 86.3 | 89.0 | 86.9 | 89.5 |
| #1 Single - MIR-MRC (F-Net) | - | - | 74.8 | 78.0 |
| #2 Single - nlnet | - | - | 74.2 | 77.1 |
| Published | | | | |
| unet (Ensemble) | - | - | 71.4 | 74.9 |
| SLQA+ (Single) | - | | 71.4 | 74.4 |
| Ours | | | | |
| BERT$_{LARGE}$ (Single) | 78.7 | 81.9 | 80.0 | 83.1 |

SQuAD 2.0 results
*Source: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*

# Results and Discussions

| System | Dev | Test |
|--------|-----|------|
| ESIM+GloVe | 51.9 | 52.7 |
| ESIM+ELMo | 59.1 | 59.2 |
| OpenAI GPT | - | 78.0 |
| BERT$_{BASE}$ | 81.6 | - |
| BERT$_{LARGE}$ | **86.6** | **86.3** |
| Human (expert)[†] | - | 85.0 |
| Human (5 annotations)[†] | - | 88.0 |

SWAG Dev and Test accuracies

| | | Dev Set | | | | |
|-------|------|------|-------|------|-------|-------|
| Tasks | MNLI-m (Acc) | QNLI (Acc) | MRPC (Acc) | SST-2 (Acc) | SQuAD (F1) |
| BERT$_{BASE}$ | 84.4 | 88.4 | 86.7 | 92.7 | 88.5 |
| No NSP | 83.9 | 84.9 | 86.5 | 92.6 | 87.9 |
| LTR & No NSP | 82.1 | 84.3 | 77.5 | 92.1 | 77.8 |
| + BiLSTM | 82.1 | 84.1 | 75.7 | 91.6 | 84.9 |

Ablation over the pre-training tasks

| Hyperparams | | | | Dev Set Accuracy | | |
|-----|-----|-----|----------|--------|------|-------|
| #L | #H | #A | LM (ppl) | MNLI-m | MRPC | SST-2 |
| 3 | 768 | 12 | 5.84 | 77.9 | 79.8 | 88.4 |
| 6 | 768 | 3 | 5.24 | 80.6 | 82.2 | 90.7 |
| 6 | 768 | 12 | 4.68 | 81.9 | 84.8 | 91.3 |
| 12 | 768 | 12 | 3.99 | 84.4 | 86.7 | 92.9 |
| 12 | 1024 | 16 | 3.54 | 85.7 | 86.9 | 93.3 |
| 24 | 1024 | 16 | 3.23 | 86.6 | 87.8 | 93.7 |

Ablation over BERT model size

| System | Dev F1 | Test F1 |
|--------|--------|---------|
| ELMo (Peters et al., 2018a) | 95.7 | 92.2 |
| CVT (Clark et al., 2018) | - | 92.6 |
| CSE (Akbik et al., 2018) | - | **93.1** |
| Fine-tuning approach | | |
| BERT$_{LARGE}$ | 96.6 | 92.8 |
| BERT$_{BASE}$ | 96.4 | 92.4 |
| Feature-based approach (BERT$_{BASE}$) | | |
| Embeddings | 91.0 | - |
| Second-to-Last Hidden | 95.6 | - |
| Last Hidden | 94.9 | - |
| Weighted Sum Last Four Hidden | 95.9 | - |
| Concat Last Four Hidden | 96.1 | - |
| Weighted Sum All 12 Layers | 95.5 | - |

CoNLL-2003 Named Entity Recognition results

# Conclusions

- It works well.