

- **RoBERTa (Facebook):** Improved training procedures, larger datasets, and dynamic masking. Generally outperformed BERT.
  - **ALBERT (Google):** Parameter reduction techniques to make models smaller and faster.
  - **ELECTRA (Google):** Used a more efficient pre-training objective (replaced token detection instead of masked token prediction).
  - **DeBERTa (Microsoft):** Introduced disentangled attention mechanisms for better contextual understanding. Often considered one of the best BERT-style models.
  - **Longformer and other Long-Context Models:** Addressed BERT's limitation with handling long sequences of text.
- 7.5 (Text-to-Text Transfer Transformer):** Reframed all NLP tasks as text-to-text problems.
- **GPT-3 and later LLMs (OpenAI):** While different in architecture (decoder-only), these models, and especially their successors, represent a major shift in scale and capabilities. They're often used for zero-shot or few-shot learning, reducing the need for extensive fine-tuning.
  - **The Rise of Large Language Models (LLMs):** The current focus is heavily on extremely large language models (LLMs) with billions or even trillions of parameters. These models, like GPT-4, PaLM, and others, demonstrate emergent abilities and can perform a wide range of tasks with minimal task-specific training. They often use variations of the Transformer architecture but are trained on massive datasets.

- **ELMo: Shallow Concatenation** - ELMo has two separate LSTMs, one left-to-right and one right-to-left. It then combines the outputs of these two LSTMs, but

- two-stage framework: pre-training and fine-tuning. During pre-training
- Pre-training**
- the model learns from unlabeled data using tasks like masked language modeling and next sentence prediction.
- Fine-Tuning**

BERTLARGE (24 layers) be the same size as One

bidirectional self-attention,

**CLS**  
Definition : The [CLS] token starts as a special token with a random (or pre-trained) embedding. As it passes through the BERT network, its vector representation is updated based on the entire input sequence, becoming a powerful representation of the sequence as a whole.

- Sentences in a pair are separated by a semicolon
- To distinguish between the two

- The final input representation for a token is the sum of its WordPiece embedding, **segment embedding**, and **position embedding**.

**Output:**

- The final hidden state vector corresponding to the [CLS] token (denoted as C) is used as the aggregate representation of the *entire sequence* for classification tasks.
- The final hidden state vectors for the individual input tokens (denoted as Ti) can be used for token-level tasks.

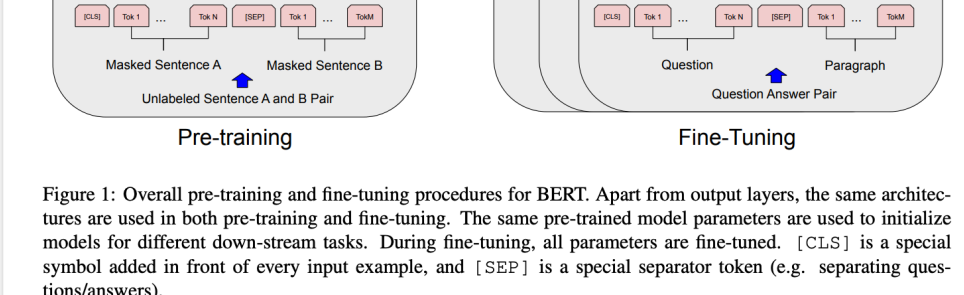
- **Method:** Mask a cer
- **Advantages:**

- The final hidden vector
- Instead of always repl

- ## MLM Limitations
- **Mismatch between Pre-training and Fine-tuning:**
  - The [MASK] token is used during pre-training, but it does not appear during fine-tuning.
  - This means the model's learning during pre-training may not be fully applicable during fine-tuning, potentially leading to performance degradation.
  - **Mitigation:**
  - Instead of always replacing masked words with the [MASK] token, they are sometimes replaced with random tokens or the original token to prevent the model from over-relying on the specific token.

- **Goal:** To train a model that understands

- For each pre-training example, select two sentences, A and B.
- 50% of the time, B is the actual sentence that follows A (labeled "IsNext").
- 50% of the time, B is a random sentence from the corpus (labeled "NotNext").
- The model then predicts whether B is the next sentence after A.
- 
- 
- **Benefit:** Pre-training on this task is shown to be very beneficial for both QA and NLI tasks.



- **End-to-End Fine-tuning:**

## GLUE

- **Dominant Performance:** BERT significantly outperformed previous state-of-the-art models on the GLUE benchmark capabilities in various natural language understanding tasks. Notably, it showed impressive performance gains even in data-scarce
- **Importance of Model Size:** A larger BERT model (BERT-LARGE) leads to improved performance. This indicates that natural language understanding capabilities.
- **Structural Similarity:** BERTBASE shares a nearly identical model structure with OpenAI GPT, except for the attention suggests that attention masking is a key factor contributing to BERT's performance.

- | System                                 |      |      |           |
|--|------|------|-----------|
|  | Dev  | Test |           |
|  | EM   | F1   | F1        |
| Top Leadboard Systems (Dec 10th, 2018) |      |      |           |
| Human                                  | -    | -    | 82.3 81.2 |
| #1 Ensemble - dist                     | -    | -    | 80.0 91.7 |
| #2 Ensemble - GANet                    | -    | -    | 84.5 90.9 |
| Published                              |      |      |           |
| DPM-HL-Ms (Single)                     | -    | 85.6 | 85.8      |
| M Reader (Ensemble)                    | 81.2 | 87.9 | 88.5      |
| Ours                                   |      |      |           |
| ERTNet (Single)                        | 80.8 | 88.5 | -         |
| ERTNet (Ensemble)                      | 84.1 | 90.9 | -         |
| ERTNet (Single)                        | 85.8 | 91.8 | -         |
| ERTNet (Ensemble)                      | 88.8 | 94.8 | -         |

**SQuAD v1.1+SQuAD v2.0**

Top Leaderboard Systems (Dec 10th, 2018)			
Huzaifa	-	-	82.3 91.2
#1 Ensemble - nhat	-	-	86.0 91.7
#2 Ensemble - qNNet	-	-	85.4 90.6
Published			
BiDAFv2+FLMs (Single)	-	85.6	85.8
R.M. Reader (Ensemble)	81.2	87.9	82.3 85.5
Ours			
BERT <sub>BASE</sub> (Single)	80.8	88.5	-
BERT <sub>LARGE</sub> (Single)	84.1	90.9	-
BERT <sub>LARGE</sub> (Ensemble)	85.8	91.8	-
BERT <sub>LARGE</sub> (Sp+TriTQAQ)	84.2	91.1	85.1 91.8
BERT <sub>LARGE</sub> (En+TriTQAQ)	86.2	92.2	87.4 93.2

**Excellent Performance of BERT:** BERT's performance is significantly better than traditional models, demonstrating the importance of Fine-tuning: By fine-tuning

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 <sub>LARGE</sub> (Single)	78.7	81.9	80.0	83.1

FRASE, BERTI, ARGE) demonstrated outstanding

- **Handling Unanswerable Questions:** SQuAD v2.0 includes questions with no answer, requiring models to determine whether a question is answerable. BERT effectively handled these questions by outputting the [CLS] token.
- **Surpassing Previous Best Models:** BERT also achieved higher F1 scores than previous state-of-the-art models on SQuAD v2.0, demonstrating its superior performance.

**SWAG**

## System

## size

args	Dev Set Accuracy			
#A	LM (jpl)	MNLI-m	MRPC	SST-2
12	5.84	77.9	79.8	88.4
3	5.24	80.6	82.2	90.7
12	4.68	81.9	84.8	91.3
12	3.99	84.4	86.7	92.9
16	3.54	85.7	86.9	93.3
16	3.23	86.6	87.8	93.7

blation over BERT model size. #L = the layers; #H = hidden size; #A = number of attention heads. “LM (ppl)” is the masked LM perplexity on training data.

## Feature-Based Approach

System	Dev F1	Test F1
ELMo (Peters et al., 2018a)	95.7	92.2
CVT (Clark et al., 2018)	-	92.6
CSE (Araki et al., 2018)	-	93.1
<b>★ Fine-tuning approach</b>		
BERT <sub>task1</sub>	96.6	92.8
BERT <sub>task2</sub>	96.4	92.4
<b>Feature-based approach (BERT<sub>base</sub>)</b>		
Embedding	83.0	-

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	-

- BERTLARO

- a Dev F1 score of 96.6 and a Test F1 score of 92.8. This is a superior result compared to existing models such as ELMo, CVT, and CSE.

BERT is effe  
tuning provi

fine-tuning provides the highest performance, but the feature-based approach, which using a combination of features from multiple layers, can achieve performance comparable to fine-tuning. This demonstrates that BERT can be flexibly applied to a variety of NLP tasks.

- BERTLARGE achieved very high performance when using fine-tuning, with a Dev F1 score of 96.6 and a Test F1 score of 92.8. This is a superior result compared to existing models such as ELMo, CVT, and CSE.