

# A Comparative Study of DistilBERT and ALBERT: Efficient Transformer Models for Natural Language Processing

Prepared by: Hassan Gebrill

## Abstract

Transformer-based architectures such as BERT have revolutionized natural language processing (NLP). However, the high computational cost and memory requirements of BERT limit its deployment in real-world, resource-constrained environments. To address these challenges, smaller yet powerful models such as **DistilBERT** and **ALBERT** were introduced. This research provides a detailed comparative analysis of these two models, exploring their architectures, training strategies, efficiency improvements, and performance on various NLP tasks. The paper concludes with an evaluation of their practical applications and future directions for model compression in deep learning.

## 1. Introduction

The introduction of the Bidirectional Encoder Representations from Transformers (BERT) model by Google in 2018 marked a breakthrough in NLP. BERT achieved state-of-the-art results in numerous tasks including question answering, sentence classification, and named entity recognition. Despite its success, BERT's large size (110M parameters for BERT-Base and 340M for BERT-Large) made it computationally expensive for deployment in production systems or on devices with limited resources.

To mitigate these limitations, researchers developed lighter and faster variants of BERT, among which **DistilBERT** and **ALBERT** are prominent. Both models aim to maintain the representational power of BERT while reducing the computational footprint, but they achieve this goal through distinct strategies.

## 2. DistilBERT

### 2.1. Overview

DistilBERT, introduced by Hugging Face in 2019, is a *distilled* version of BERT. It uses the concept of **knowledge distillation**, a model compression technique where a smaller *student model* learns to mimic the behavior of a larger *teacher model*.

### 2.2. Architecture

DistilBERT reduces the number of layers in BERT by half while retaining the hidden size and attention heads. Specifically, DistilBERT has 6 Transformer layers compared to the 12 in BERT-Base. The model maintains the same hidden dimension (768) and number of attention heads (12), allowing it to preserve the core representational structure of BERT.

### 2.3. Training Method

The key idea behind DistilBERT is **knowledge distillation**. During training, the student model learns not only from the original training data but also from the *soft targets* (probability distributions) generated by the teacher model (BERT). The loss function combines three components:

- The distillation loss (Kullback-Leibler divergence between student and teacher logits)
- The traditional cross-entropy loss on true labels
- A cosine embedding loss between hidden states

This multi-objective training enables DistilBERT to retain around 97% of BERT's performance while being 40% smaller and 60% faster.

### 2.4. Performance and Efficiency

DistilBERT achieves near state-of-the-art performance across major benchmarks such as GLUE and SQuAD. It offers a good trade-off between accuracy and computational efficiency, making it suitable for applications like real-time NLP systems and mobile inference.

### 3. ALBERT

#### 3.1. Overview

ALBERT (A Lite BERT) was introduced by Google Research and Toyota Technological Institute at Chicago in 2019. Unlike DistilBERT, which focuses on reducing model depth, ALBERT focuses on reducing model *parameter redundancy*.

#### 3.2. Architecture

ALBERT employs two major parameter reduction techniques:

1. **Factorized Embedding Parameterization:** Instead of directly mapping one-hot word vectors to the hidden size (e.g., 768), ALBERT first projects them into a smaller embedding space (e.g., 128) and then to the hidden space. This significantly reduces the size of the embedding matrix.
2. **Cross-layer Parameter Sharing:** ALBERT shares parameters (weights) across layers, meaning that all Transformer layers use the same set of parameters. This reduces the total number of parameters dramatically without significantly degrading performance.

#### 3.3. Training Enhancements

ALBERT introduces a new pre-training objective called **Sentence Order Prediction (SOP)**, replacing BERT’s Next Sentence Prediction (NSP). SOP is designed to improve the model’s ability to understand inter-sentence coherence, leading to better performance in tasks such as question answering and natural language inference.

#### 3.4. Performance and Efficiency

ALBERT achieves competitive or even superior performance compared to BERT-Large, despite having significantly fewer parameters. For instance, ALBERT-Base has only 12M parameters compared to BERT-Base’s 110M. The model also scales better for distributed training and uses memory more efficiently.

### 4. Comparison Between DistilBERT and ALBERT

Table 1 provides a side-by-side summary of the two models.

Table 1: Comparison between DistilBERT and ALBERT

Aspect	DistilBERT	ALBERT
Compression Method	Knowledge distillation	Parameter sharing & factorization
Parameter Count	~66M	~12M (Base)
Training Objective	Distillation + MLM	MLM + Sentence Order Prediction
Architecture	6-layer Transformer	Shared 12-layer Transformer
Speed Improvement	~60% faster than BERT	Efficient memory use, faster training
Performance Retention	~97% of BERT	Comparable or higher than BERT
Use Case	Inference on limited devices	Large-scale distributed tasks

While DistilBERT focuses on reducing computation for inference, ALBERT emphasizes parameter efficiency and scalability during training. In practical scenarios, DistilBERT is preferred for lightweight applications such as chatbots and real-time systems, whereas ALBERT is used for tasks demanding high accuracy and large context understanding.

## 5. Applications

Both DistilBERT and ALBERT have found wide applications in NLP:

- **Text Classification:** Spam detection, sentiment analysis, and toxicity detection.
- **Question Answering:** Efficient deployment of SQuAD-like systems.
- **Named Entity Recognition (NER):** Faster and memory-efficient entity tagging.
- **Conversational AI:** Used in chatbots, virtual assistants, and dialogue generation.

## 6. Conclusion

DistilBERT and ALBERT represent two distinct yet complementary approaches to improving the efficiency of transformer models. DistilBERT leverages knowledge distillation to create smaller, faster models suitable for deployment, while ALBERT focuses on reducing parameter redundancy for scalable training. Both achieve impressive results with minimal compromise in performance, demonstrating the power of model compression techniques in modern deep learning. As research continues, combining these methods could yield even more compact and efficient transformer models.

## References

1. Sanh, V. et al. (2019). *DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter*. Hugging Face.
2. Lan, Z. et al. (2020). *ALBERT: A Lite BERT for Self-supervised Learning of Language Representations*. ICLR.
3. Devlin, J. et al. (2018). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. arXiv:1810.04805.