# DistilBERT and ALBERT: Efficient Variants of BERT

## 1. Introduction

The introduction of BERT (Bidirectional Encoder Representations from Transformers) revolutionized natural language processing (NLP) by providing powerful contextual embeddings that significantly improved performance across tasks like sentiment analysis, text classification, and question answering. However, BERT's size and computational demands limit its deployment in real-time and resource-constrained environments.

To address these challenges, researchers introduced more efficient variants:

- **DistilBERT** – a compact, faster version of BERT developed using knowledge distillation.

- **ALBERT** – a parameter-efficient redesign of BERT that achieves drastic reductions in model size through factorization and weight sharing.

This paper summarizes their architectures, techniques, and applications, highlighting how they complement each other in addressing efficiency and scalability.

## 2. DistilBERT

### 2.1 Overview

DistilBERT, developed by HuggingFace, compresses BERT-base by:

- Reducing model size by 40%.

- Achieving 60% faster inference.

- Retaining about 97% of BERT's performance.

This balance makes it a strong candidate for deployment on edge devices and interactive systems.
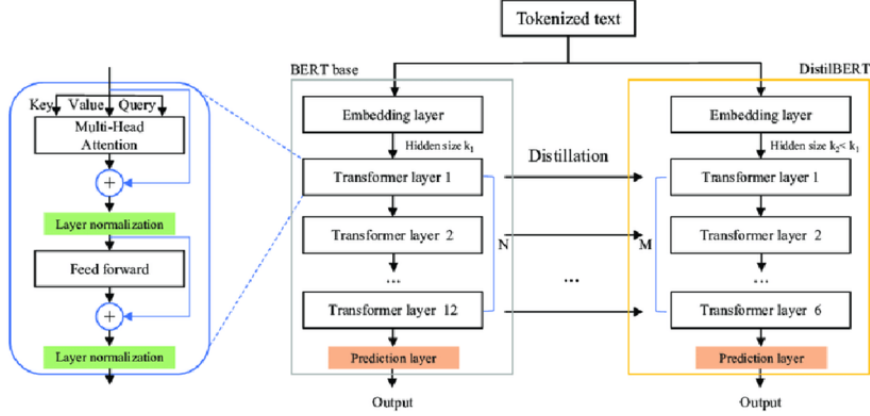
Figure 1: DistilBERT knowledge distillation process: a smaller student model learns from a larger teacher model.

## 2.2 Knowledge Distillation

The key idea behind DistilBERT is knowledge distillation:

- The teacher is a pretrained BERT model.

- The student is a smaller transformer with fewer layers.

- A triple loss function is used: distillation loss, masked language modeling loss, and cosine embedding loss.

This allows the student to capture both the accuracy and the uncertainty in teacher predictions, making it more efficient than training from scratch.

## 2.3 Key Features and Applications

- Uses 6 layers (vs. 12 in BERT-base).

- Omits token-type embeddings and the pooler layer.

- Compatible with text classification, sentiment analysis, question answering, and information retrieval.

- Optimized training: large batch sizes, dynamic masking, and no next sentence prediction objective.

# 3. ALBERT

## 3.1 Overview

ALBERT (A Lite BERT), developed by Google Research and TTI-Chicago, reduces the parameter count drastically while preserving depth. For instance, BERT-base has 110M

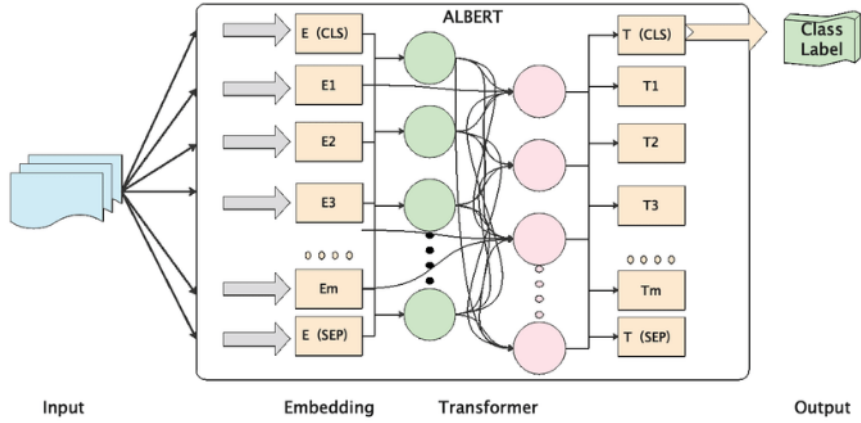parameters, whereas ALBERT-base has only 12M, achieved without sacrificing accuracy on major NLP benchmarks.



Figure 2: ALBERT architecture with parameter sharing across layers and factorized embeddings.

## 3.2 Techniques Used

- **Factorized Embedding Parameterization:** Splits embeddings into smaller dimensions, reducing memory usage.

- **Cross-Layer Parameter Sharing:** Encoder layers reuse the same parameters, drastically lowering model size.

- **Sentence-Order Prediction (SOP):** Improves inter-sentence understanding, replacing BERT's Next Sentence Prediction.

## 3.3 Additional Insights

The backbone of ALBERT remains transformer encoder layers with GELU activation, but it introduces three important innovations:

- Embedding factorization reduces parameters by nearly 80%.

- Parameter sharing across layers cuts redundancy by about 70%.

- SOP encourages coherence-focused learning rather than topic prediction.

This enables ALBERT models (Base, Large, XL, XXL) to achieve state-of-the-art results on benchmarks like GLUE, SQuAD, and RACE with far fewer parameters than BERT.

## 3.4 Applications

- Scalable pretraining on massive corpora.

- Reading comprehension benchmarks such as RACE and SQuAD.

- General NLP tasks requiring robust contextual understanding.

# 4. Comparison

- **DistilBERT:** Prioritizes speed and inference efficiency. Ideal for deployment on devices with limited compute resources.

- **ALBERT:** Prioritizes parameter reduction and scalability. Ideal for large-scale research and memory-constrained environments.

# 5. Conclusion

DistilBERT and ALBERT demonstrate two complementary strategies for improving BERT efficiency:

- DistilBERT compresses the model using knowledge distillation, balancing performance and speed.

- ALBERT redesigns the architecture to eliminate parameter redundancy and improve scalability.

Together, they enable broader adoption of transformer-based models in both industrial and research contexts.