

# Word Embedding with BERT

Chien-An Chiu <sup>1</sup>

<sup>1</sup>Department of Electrophysics, Department of Electrical Engineering, National Yang Ming Chiao Tung University, Hsinchu 300, Taiwan

March 19, 2025

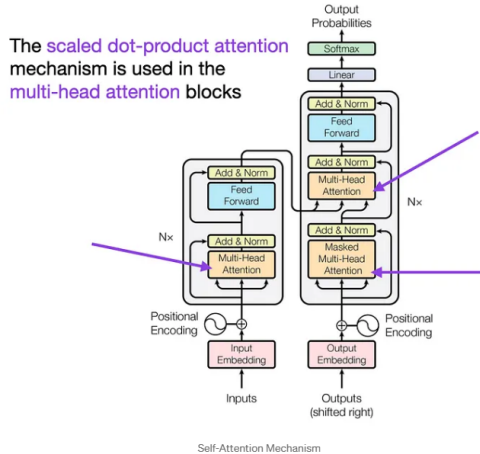
## Abstract

This challenge explores the mechanism of attention in BERT and its impact on sentiment classification using the IMDB dataset. We investigate attention scores, fine-tuning behavior under partial and contrastive data, the structure of contextual embeddings, and how they can be clustered to reflect sentiment categories. Embeddings are further analyzed using dimensionality reduction techniques such as PCA and t-SNE.

## 1 Introduction

Transformer-based models like BERT rely on the attention mechanism to understand contextual relationships within text. In this challenge, we explore how BERT performs sentiment classification and how the learned word embeddings reflect semantic structures.

### 1.1 Multi-Head Attention Architecture



**Figure 1:** Multi-Head Attention in Transformer Architecture [1]

Figure 1 illustrates the attention mechanism within the Transformer architecture, specifically in BERT. Each input token is first projected into

three separate matrices: the Query  $Q$ , Key  $K$ , and Value  $V$ .

In each multi-head attention layer:

- Attention scores are computed via scaled dot-product attention.
- Multiple heads (typically 12 or more in BERT) run attention in parallel.
- Their outputs are concatenated and linearly projected.

This design allows different heads to specialize in capturing different types of linguistic dependencies such as syntactic structures (e.g., subjectverb) and semantic patterns (e.g., adjectivenoun pairing). Each head independently computes:

$$\text{Attention}(Q_i, K_i, V_i) = \text{softmax} \left( \frac{Q_i K_i^T}{\sqrt{d_k}} \right) V_i$$

The final output is:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

This module is essential to BERT's ability to model contextual meaning in both directions simultaneously, as all positions can attend to one another within the same sequence.

## 2 BERT Architecture and Pretraining

BERT (Bidirectional Encoder Representations from Transformers) is composed entirely

of Transformer **encoder blocks** there is no decoder in its architecture. In the base version (BERT-base), the model consists of:

- 12 encoder layers
- 12 attention heads per layer
- Hidden size of 768
- Total of 110 million parameters

Each encoder layer contains two main sub-components:

1. **Multi-Head Self-Attention Layer:** Enables each token to attend to all other tokens in the sequence in parallel. Multiple attention heads allow the model to capture different linguistic dependencies simultaneously.
2. **Position-wise Feed-Forward Network:** Applies two linear transformations with a non-linearity (typically GELU) in between, independently to each token position.

Both subcomponents are wrapped with **residual connections** followed by **layer normalization**:

$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

This stabilizes training and helps gradients flow across deep layers.

BERTs input representation is formed by summing:

- **Token Embeddings:** WordPiece tokenized input.
- **Segment Embeddings:** Indicates sentence A or B (used in NSP).
- **Positional Embeddings:** Since the Transformer has no recurrence, position is encoded explicitly.

## 2.1 Pretraining Objectives

BERT is pretrained on unlabeled text from Wikipedia and BookCorpus using two unsupervised objectives:

1. **Masked Language Modeling (MLM)** Randomly masks 15% of input tokens and trains the model to predict the masked words. Unlike standard left-to-right models, this enables BERT to use bidirectional context:

$$\text{Prediction: } p(w_i | w_{<i}, w_{>i})$$

## 2. Next Sentence Prediction (NSP)

Given two segments A and B, the model predicts whether B is the actual next sentence that follows A in the original document. This trains the model to understand inter-sentence coherence, which is critical for tasks like QA and NLI.

## 2.2 Token Output

The final output of BERT is a contextualized vector for each token. The special [CLS] token's final hidden state is used as the aggregate sequence representation and typically passed to a classifier head for downstream tasks.

## 3 Fine-Tuning on IMDB Dataset

We fine-tuned BERT using a combination of a small custom dataset and a subset of the IMDB movie review dataset. Instead of training on the entire corpus, we selected 2500 positive and 2500 negative samples, and combined them with 8 handcrafted samples (4 positive, 4 negative) containing sarcasm and contrastive patterns.

```
# Load IMDB
imdb = load_dataset("imdb")
df_train_full = pd.DataFrame(imdb["train"])

# --- Select 2500 positive + 2500 negative samples
df_pos = df_train_full[df_train_full["label"] == 1].sample(2500, random_state=42)
df_neg = df_train_full[df_train_full["label"] == 0].sample(2500, random_state=42)
df_balanced = pd.concat([df_pos, df_neg], ignore_index=True)

# Add 'sentiment' column for consistency
df_balanced["sentiment"] = df_balanced["label"]

# Combine with custom data
df_total = pd.concat([df_custom, df_balanced], ignore_index=True)

# Split into train and validation
train, val = train_test_split(df_total, test_size=0.2, stratify=df_total["sentiment"], random_state=42)
```

**Figure 2:** Partial dataset selection and preprocessing for fine-tuning

## 3.1 Effect of Contrastive Data on Sarcasm Understanding

We evaluated the model's ability to handle sarcasm by comparing training with and without contrastive data. When sarcastic or ambiguous examples were excluded, the model exhibited poor generalization on these cases.

```
Epoch 1/6
8/8 [====- 43s 3s/step - loss: 0.6652 - accuracy: 0.6250 - val_loss: 0.6775 - val_accuracy: 0.6923
Epoch 2/6
8/8 [====- 20s 2s/step - loss: 0.5442 - accuracy: 0.8214 - val_loss: 0.6114 - val_accuracy: 0.6923
Epoch 3/6
8/8 [====- 17s 2s/step - loss: 0.3584 - accuracy: 0.9643 - val_loss: 0.4530 - val_accuracy: 0.9231
Epoch 4/6
8/8 [====- 16s 2s/step - loss: 0.2068 - accuracy: 0.9643 - val_loss: 0.3395 - val_accuracy: 0.9462
Epoch 5/6
8/8 [====- 16s 2s/step - loss: 0.1048 - accuracy: 0.9821 - val_loss: 0.2471 - val_accuracy: 0.9231
Epoch 6/6
8/8 [====- 18s 2s/step - loss: 0.0515 - accuracy: 1.0000 - val_loss: 0.9795 - val_accuracy: 0.5385
```

	precision	recall	f1-score	support
Negative	1.00	0.14	0.25	7
Positive	0.50	1.00	0.67	6
accuracy			0.54	13
macro avg	0.75	0.57	0.46	13
weighted avg	0.77	0.54	0.44	13

**Figure 3:** Low validation performance without contrastive examples (e.g., sarcasm)

The F1-score for the negative class dropped significantly despite high precision. The model overfitted on literal sentiment but failed to recognize cases such as: *"Blew my mind how bad it was. What a waste."*

### 3.2 Improved Performance with Full Dataset

When we included contrastive examples in training, performance across both positive and negative classes became more balanced and robust. As shown in the figure below, both accuracy and F1-scores improved across all metrics:

```
Epoch 1/6
16/16 [====- 49s 2s/step - loss: 0.6629 - accuracy: 0.6667 - val_loss: 0.4203 - val_accuracy: 0.7333
Epoch 2/6
16/16 [====- 22s 2s/step - loss: 0.4561 - accuracy: 0.9242 - val_loss: 0.4472 - val_accuracy: 0.8667
Epoch 3/6
16/16 [====- 24s 2s/step - loss: 0.2027 - accuracy: 0.9848 - val_loss: 0.2838 - val_accuracy: 0.9333
Epoch 4/6
16/16 [====- 23s 2s/step - loss: 0.0745 - accuracy: 1.0000 - val_loss: 0.3854 - val_accuracy: 0.9333
Epoch 5/6
16/16 [====- 24s 2s/step - loss: 0.0362 - accuracy: 1.0000 - val_loss: 0.2183 - val_accuracy: 0.9333
Epoch 6/6
16/16 [====- 22s 2s/step - loss: 0.0155 - accuracy: 1.0000 - val_loss: 0.2375 - val_accuracy: 0.9333
```

	precision	recall	f1-score	support
Negative	1.00	0.88	0.93	8
Positive	0.88	1.00	0.93	7
accuracy			0.93	15
macro avg	0.94	0.94	0.93	15
weighted avg	0.94	0.93	0.93	15

**Figure 4:** Improved validation performance after including contrastive (sarcastic) samples

This result suggests that BERT benefits significantly from seeing linguistically challenging data such as sarcasm. Even partial fine-tuning of only the classification head was sufficient to learn subtle sentiment features, demonstrating BERT’s strong generalization capacity when guided by appropriate examples.

## 4 Embedding Analysis and Visualization

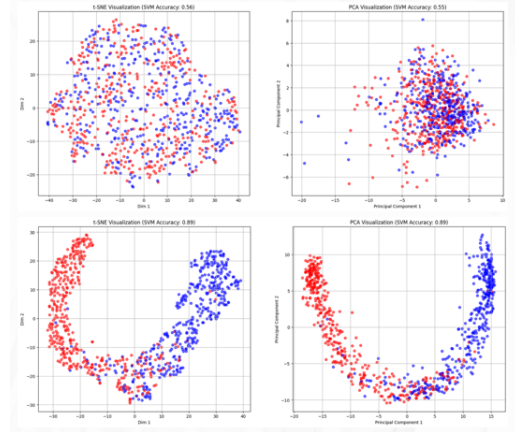
In BERT, each input sentence is transformed into a sequence of contextualized token embeddings through multiple self-attention and feed-forward layers. For sentence-level tasks, we extract the embedding corresponding to the special [CLS] token from the final layer, which ag-

gregates semantic information from the entire sentence.

### 4.1 Dimensionality Reduction using PCA and t-SNE

Since each embedding vector lies in a high-dimensional space ( $\mathbb{R}^{768}$ ), we employ dimensionality reduction techniques to visualize the geometric structure of the learned representation:

- **Principal Component Analysis (PCA):** A linear projection technique that maximizes variance along principal axes, revealing global structure.
- **t-distributed Stochastic Neighbor Embedding (t-SNE):** A nonlinear method that preserves local neighborhood structure, suitable for clustering analysis.



**Figure 5:** PCA and t-SNE visualizations of sentence embeddings before (top) and after (bottom) fine-tuning

As shown in Figure 4, embeddings before fine-tuning exhibit overlapping clusters with poor separation between sentiment labels. After fine-tuning, the red (positive) and blue (negative) embeddings are mapped to more distinct regions in the feature space, suggesting that BERT has learned to encode sentiment polarity into geometric structure.

### 4.2 Visualizing Extreme Sentiments

To better understand the output behavior, we also examined the most confidently classified examples. After applying softmax on prediction logits, we collected the top 5 most positive and most negative reviews based on prediction confidence:

Most Positive Reviews:  
(0.997) This movie is about human relationships. Charming, funny, and well written.  
(0.997) Garam Masala is one of the funniest film I've seen in ages. Akshay Kumar is  
(0.997) I saw the movie recently during the Boston Film festival. The movie was very  
(0.997) This is possibly the best short crime drama I've ever seen. The acting is su  
(0.997) Great entertainment from start to the end. Wonderful performances by Belushi  
Most Negative Reviews:  
(0.001) It's been a long time since I last saw a movie this bad.. The acting is very  
(0.001) This is the worst ripoff of Home Alone movies that I have EVER seen! Watch p  
(0.001) This is not Michael Madsen's fault, he was hardly in it. This movie was just  
(0.001) How the hell did they get this made?! Presenting itself as a caper comedy, t  
(0.001) I must say that during my childhood I'm quite proud of a lot of the movies I

**Figure 6:** Examples of reviews with maximum and minimum sentiment scores

Positive samples typically contain emotionally rich, descriptive language (e.g., Charming, funny, and well written), whereas negative samples are more blunt and emphatic (e.g., This is the worst ripoff I have EVER seen!). These outputs indicate that BERT not only learns general sentiment direction, but also captures strength and extremeness in tone.

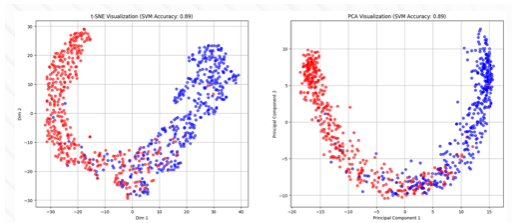
The combination of embedding visualization and extreme case inspection confirms that the fine-tuned model organizes sentence semantics meaningfully and effectively within the embedding space.

## 5 Clustering Embeddings

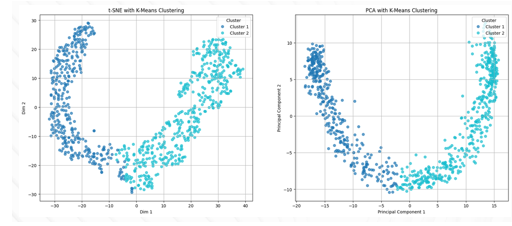
To analyze the semantic structure of BERT sentence embeddings, we apply K-Means clustering to the [CLS] token embeddings after fine-tuning. The goal is to determine whether sentence embeddings naturally organize according to sentiment.

### 5.1 Clustering and Visualization

We applied K-Means with  $k = 2$  and visualized the clusters using both t-SNE and PCA. As shown in Figure 7 and Figure 8, the embeddings are projected into two distinct regions.



**Figure 7:** t-SNE and PCA of sentence embeddings colored by ground truth sentiment



**Figure 8:** t-SNE and PCA of sentence embeddings colored by K-Means clustering results

We observe that the embeddings separate clearly into two distinct regions corresponding to positive and negative sentiment, even when clustering is done without label supervision.

### 5.2 Cluster Evaluation with Confusion Matrix

To evaluate the correspondence between the discovered clusters and the true sentiment labels, we compute a confusion matrix, shown in Figure 9.



**Figure 9:** Confusion matrix comparing clusters with true sentiment labels

Although not perfect, the alignment demonstrates BERTs capability to encode sentiment information geometrically in embedding space.

### 5.3 Cluster Interpretation through Samples and Keywords

To interpret the content of each cluster, we examine representative reviews (Figure 10) and the most frequent keywords (Figure 11).

Cluster 1 Sample Reviews ---  
- The year 1934 was when Shirley Temple played three major movies and really began to make a name for herself. Unfortun  
- This is a great film for McDermott's and Beaulieu family's splendid time is guaranteed for all. The audience (East come  
- I saw this for Gary Bussey and Fred Williamson thinking they were buddy cops. They are but Bussey is in the opening ac  
- The film Forrest was a first and a last for Greta Garbo. It was her first American made film at MGM, the only studio  
- A year or so ago, I was watching the TV news when a story was broadcast about a mobile movie being filmed in my area.  
- Shakespeare said that we are actors put into a great stage. But when this stage is Israel, the work that we interpret  
- This is a very good movie. Do you want to know the real reason why so many here are knocking this movie? I will tel  
- GONE IN 60 SECONDS / (2000) \*\*\* (out of four)Grr />Grr /> "Gone in 60 Seconds" is an energetic, slick, stylish action  
- Modern, original, romantic story.Grr />Grr />Very good acting of both Nicole Kidman and Ben Chaplin.Grr />Grr />Miss  
- I agree with the other comments. I saw this movie years ago. Christopher Plummer is hilarious as a dandy. The rival

Cluster 2 Sample Reviews ---  
- I am 2 home entertainment stores and I've seen a lot of bad movies in my time but this one was so bad it compelled  
- Okay, I saw this movie as a child and really loved it. My parents never purchased the movie for me, but I think I'll  
- Cassidy(Maria Brady)puts a gun in her mouth blowing the back of her head out on boyfriend Neal(Iason Dittler). Cassidy  
- America needs the best man possible to win "The game" so who do they hire? A runaway (Oh brother!) played by Rusty T  
- A terrible movie containing a berry of D-list Canadian actors who seem so self-conscious about the fact they are on-sc  
- You know all those letters to "Father Christmas" and "Joan" that are sent every year? Well, it turns out that they  
- This movie will confuse you to death. Furthermore, if you're a Beanie Richards' fan, don't even think of renting this  
- Okay the promos promised a comedy and people(few) went to watch it Being the first release of 2006 is not a bad thing  
- This is not my fanny and gory as the DVD box claims. I really love twisted and weird movies, but this one is really  
- I am a college student studying a-levels and need help and comments from someone who has any stress at all about the

**Figure 10:** Sample reviews from Cluster 1 and Cluster 2

```

Cluster 1
Top words: good, great, life, people, seen, best, love, did, think
Cluster 2
Top words: bad, don, good, make, plot, acting, worst, watching, know, thing

```

**Figure 11:** Top keywords extracted from each cluster

Cluster 1 is dominated by strongly positive words such as *good*, *great*, *love*, while Cluster 2 contains typical negative sentiment indicators such as *bad*, *worst*, *boring*. These results confirm that BERTs embeddings reflect sentiment polarity in a spatially meaningful way, and that unsupervised clustering is able to recover it to a large extent.

## 6 Results and Evaluation

Our fine-tuned BERT model demonstrates strong performance in sentiment classification, with balanced precision, recall, and F1-scores across both classes. Incorporating even a small number of contrastive examples especially sarcastic or ambiguous ones significantly improved generalization to challenging cases.

Without such data, models struggled with sentiment reversals, but their inclusion boosted validation accuracy and robustness. Dimensionality reduction techniques (PCA and t-SNE) revealed clear separation between sentiment classes, further confirmed by over 89% accuracy using an SVM classifier in 2D space.

K-Means clustering on [CLS] embeddings also produced distinct sentiment groups, with high alignment to true labels and polarized keywords (e.g., *love* vs. *boring*), highlighting

BERTs ability to encode sentiment geometrically and semantically.

## 7 Conclusion

This challenge demonstrates BERTs strong ability to capture both local and global semantic relationships through its attention-based architecture. Even with limited data, fine-tuning on contrastive and challenging examples enables robust performance, particularly in handling sentiment and sarcasm.

Visualizations using PCA, t-SNE, and K-Means clustering show that BERTs embeddings naturally separate sentiment classes and encode interpretable structures. In summary, fine-tuned BERT models not only excel at sentiment classification but also reveal meaningful internal representations that support deeper analysis and interpretability.

## References

- [1] R. Shaikh, “A comprehensive guide to understanding BERT: from beginners to advanced,” Medium, Apr. 2023. [Online]. Available: <https://medium.com/@shaikhrayyan123/a-comprehensive-guide-to-understanding-bert-from-beg>
- [2] A. Vaswani *et al.*, “Attention is all you need,” in *NeurIPS*, 2017.
- [3] J. Devlin *et al.*, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *NAACL*, 2019.