

# **Lecture 18: BERT Deep Dive**

**Week 6, Lecture 1 - Bidirectional Encoder Representations**

**PSYC 51.07: Models of Language and Communication**

**Winter 2026**

# Today's Agenda



1. **BERT Introduction:** What makes it special?
2. **Masked Language Modeling:** The key training objective
3. **BERT Architecture:** Model sizes and specifications
4. **Pre-training & Fine-tuning:** The two-stage paradigm
5. **Contextual Embeddings:** Seeing polysemy in action
6. **Using BERT:** Practical code examples

*Goal: Deep understanding of BERT and how it revolutionized NLP*

# BERT: Bidirectional Encoder Representations

**BERT = Encoder-only Transformer**

**Key Innovation: Masked Language Modeling (MLM)**

**Traditional Language Models:**

- Left-to-right (GPT)
- Or right-to-left
- Can't see full picture

**Example:**

- "The cat sat on the \_\_"
- Only sees left context

# Why BERT Was Revolutionary

## Before BERT (2018):

- Feature-based approaches (use Word2Vec/GloVe as features)
- Task-specific architectures
- Limited transfer learning
- Unidirectional or shallow bidirectional models

## BERT's Contributions:

### 1. Deep Bidirectionality

- True bidirectional context at every layer
- Not just concatenating left-to-right and right-to-left

# Masked Language Modeling (MLM)



## BERT's Pre-training Objective

### Training Procedure:

1. Take a sentence
2. Randomly mask 15% of tokens
3. Of the masked tokens:
  - 80%: Replace with [MASK]
  - 10%: Replace with random word
  - 10%: Keep unchanged
4. Predict the original tokens

# MLM Example: Step by Step



**Sentence:** "The quick brown fox jumps over the lazy dog"

## Step 1: Select Tokens (15%)

9 tokens total, mask ~1-2

```
1tokens = ["The", "quick", "brown", "fox",
2          "jumps", "over", "the", "lazy", "dog"]
3# Randomly select: "quick" (idx 1), "over" (idx 5)
```

## Step 2: Apply 80/10/10 Strategy

```
1"quick" → 80% → [MASK]
2"over"   → 10% → "under" (random)
```

# Next Sentence Prediction (NSP)

**BERT's second pre-training objective (debated usefulness)**

**Task:** Given two sentences A and B, predict if B follows A

Positive Example (IsNext)

**Sentence A:** "The man went to the store."

**Sentence B:** "He bought a gallon of milk."

**Label:** IsNext ✓

Negative Example (NotNext)

**Sentence A:** "The man went to the store."

**Sentence B:** "Penguins are flightless birds."

**Label:** NotNext ✗

# BERT Architecture Variants



BERT-Large	24	1024	16	340M
------------	----	------	----	------

## Architecture Details (BERT-Base):

- 12 transformer encoder layers
- 768-dimensional hidden states
- 12 attention heads per layer (64 dims each)
- 3072-dimensional feed-forward intermediate size (4x expansion)
- Maximum sequence length: 512 tokens
- Vocabulary size: 30,000 WordPiece tokens

## Special Tokens:

# BERT Input Representation

abc

Three types of embeddings are summed:

```
1# Example: Sentence pair for NSP
2sentence_a = "My dog is cute"
3sentence_b = "He likes playing"
4
5# Tokenization
6tokens = ["[CLS]", "my", "dog", "is", "cute", "[SEP]", "he", "likes", "playing"]
7
8# Three embedding types (each is a 768-dim vector):
9token_emb = [E_CLS, E_my, E_dog, E_is, E_cute, E_SEP, E_he, E_likes, E_playing]
10segment_emb = [E_A, E_A, E_A, E_A, E_A, E_A, E_B, E_B, E_B,
11position_emb= [E_0, E_1, E_2, E_3, E_4, E_5, E_6, E_7, E_8,
12
13# Final input = token + segment + position (element-wise sum)
14input_embedding = token_emb + segment_emb + position_emb
```

# WordPiece Tokenization: Worked Example

abc

## How BERT handles unknown words

```
1from transformers import BertTokenizer
2tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
3
4# Common words stay intact
5tokenizer.tokenize("The cat sat on the mat")
6# → ['the', 'cat', 'sat', 'on', 'the', 'mat']
7
8# Rare/unknown words get split into subwords
9tokenizer.tokenize("unbelievably")
10# → ['un', '##believable', '##ly'] # "##" means continuation
11
12tokenizer.tokenize("ChatGPT is transformative")
13# → ['chat', '##g', '##pt', 'is', 'transform', '##ative']
```

# BERT Pre-training



**Massive scale pre-training on unlabeled text**

**Pre-training Data:**

- **BooksCorpus:** 800M words (novels, fiction)
- **English Wikipedia:** 2,500M words
- Total: 3.3 billion words
- Diverse, high-quality text

**Training Details:**

- Batch size: 256 sequences (128,000 tokens)
- Training steps: 1M steps

# Fine-tuning BERT

Two-stage process: Pre-train then Fine-tune

Stage 1: Pre-training (done once)

```
1# Expensive: weeks on TPUs
2# Data: 3.3B words (books + Wikipedia)
3# Task: MLM + NSP
4# Result: General language understanding
5
6model = pretrain_bert(
7    data=["BooksCorpus", "Wikipedia"],
8    steps=1_000_000,
9    hardware="16 TPUs"
10)
```

Stage 2: Fine-tuning (per task)

# Fine-tuning for Different Tasks

Minimal architecture changes needed!

## 1. Single Sentence Classification

- Input: [CLS] sentence [SEP]
- Output: [CLS] representation → classifier
- Example: Sentiment analysis

## 2. Sentence Pair Classification

- Input: [CLS] sentence A [SEP] sentence B [SEP]
- Output: [CLS] representation → classifier
- Example: Natural Language Inference

## 3. Question Answering

- Input: [CLS] question [SEP] passage [SEP]

# Fine-tuning BERT: Code Example



## Using HuggingFace Transformers

```
1from transformers import BertForSequenceClassification, Trainer, TrainingArguments
2
3# Load pre-trained BERT with classification head
4model = BertForSequenceClassification.from_pretrained(
5    'bert-base-uncased',
6    num_labels=2 # Binary classification
7)
8
9# Define training arguments
10training_args = TrainingArguments(
11    output_dir='./results',
12    num_train_epochs=3,
13    per_device_train_batch_size=16,
14    learning_rate=2e-5,
15    warmup_steps=500,
16)
```

## Fine-tuning BERT: Code Example



```
21     args=training_args,  
22     train_dataset=train_dataset,  
23     eval_dataset=eval_dataset,  
24)  
25  
26trainer.train()
```

...continued

Reference: *HuggingFace Course - Chapters 1.5, 7.3*

# Contextual Embeddings in Action



Remember "bank"? Let's see BERT handle it!

```
1from transformers import BertTokenizer, BertModel
2import torch
3
4tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
5model = BertModel.from_pretrained('bert-base-uncased')
6
7# Two different contexts for "bank"
8sent1 = "I deposited money at the bank"
9sent2 = "We sat by the river bank"
10
11# Get embeddings
12def get_embedding(sentence, target_word):
13    inputs = tokenizer(sentence, return_tensors='pt')
14    outputs = model(**inputs)
15    # Find position of target word
16    tokens = tokenizer.tokenize(sentence)
```

# Contextual Embeddings in Action

```
21emb2 = get_embedding(sent2, "bank") # River bank
22
23# Compare similarity
24similarity = torch.cosine_similarity(emb1, emb2, dim=0)
25print(f"Similarity: {similarity:.3f}") # Low! (~0.3-0.5)
26# Different contexts → Different embeddings!
```

...continued

# Visualizing BERT's Contextual Embeddings



Same word, different meanings, different vectors

```
1# Find nearest neighbors for "bank" in each context
2from sklearn.neighbors import NearestNeighbors
3
4# Financial "bank" context
5neighbors_financial = find_nearest_words(emb1, vocabulary)
6# → ["banks", "financial", "account", "deposit", "loan", "credit"]
7
8# River "bank" context
9neighbors_river = find_nearest_words(emb2, vocabulary)
10# → ["shore", "riverside", "banks", "stream", "water", "edge"]
```

Concrete Measurements

Word Pair	Word2Vec Similarity	BERT Similarity
bank (fin) vs bank (river)	1.00 (same vector!)	0.42

# BERT's Impressive Results

State-of-the-art on 11 NLP tasks when released (2018)

SQuAD 2.0 (QA)	F1	66.3	83.1
MNLI (NLI)	Accuracy	80.6	86.7
SST-2 (Sentiment)	Accuracy	93.2	94.9
CoNLL-2003 (NER)	F1	92.6	92.8

## Key Observations:

- Largest gains on tasks requiring understanding (QA, NLI)
- Improvements even on well-studied benchmarks
- BERT-Large generally better than BERT-Base
- Fine-tuning is simple but very effective

# What Does BERT Learn?

Probing BERT's internal representations

Research has shown BERT captures:

## 1. Syntactic Information

- Part-of-speech tags
- Constituent structure
- Dependency relations
- Lower layers encode more syntax

## 2. Semantic Information

- Word sense disambiguation

# BERT Layer Analysis

Different layers capture different linguistic properties

```
1# Probing experiment: Train linear classifiers on each layer's representations
2from transformers import BertModel
3import numpy as np
4
5model = BertModel.from_pretrained('bert-base-uncased', output_hidden_states=True)
6
7# Get hidden states for all 12 layers
8outputs = model(**inputs)
9hidden_states = outputs.hidden_states # (13 layers: embedding + 12 transformer
10
11# Results from probing studies (Tenney et al., 2019):
12layer_specialization = {
13    "Layers 0-2": ["POS tagging", "Word boundaries"],      # Surface
14    "Layers 3-6": ["Parse trees", "Dependencies"],        # Syntax
15    "Layers 7-9": ["Semantic roles", "Coreference"],      # Semantics
16    "Layers 10-12": ["Task-specific representations"]}    # Task
```

# Discussion Questions

## 1. MLM vs. Autoregressive:

- Why is MLM better for understanding tasks?
- Can BERT generate text like GPT?
- What are the trade-offs?

## 2. The 80/10/10 Masking Strategy:

- Why not just use 100% [MASK]?
- What problem does the random replacement solve?
- Could we improve this strategy?

## 3. Pre-training Data:

- Why use books and Wikipedia?
- Would social media text work as well?

# Looking Ahead



## Today we learned:

- BERT architecture and innovations
- Masked Language Modeling
- Pre-training and fine-tuning
- Contextual embeddings
- What BERT learns

## Next lecture (Lecture 16 - BERT Variants):

- : Optimized BERT training
- : Parameter-efficient BERT
- : Smaller, faster BERT
- : ELECTRA, DeBERTa, and more

# Summary

## Key Takeaways:

### 1. BERT = Encoder-only Transformer

- Bidirectional self-attention

- Trained with Masked Language Modeling

### 2. Pre-train + Fine-tune Paradigm

- Expensive pre-training on unlabeled data (once)

- Cheap fine-tuning on task-specific data (per task)

### 3. Contextual Embeddings

- Different representations based on context

# References

## Essential Papers:

- **Devlin et al. (2019)** - "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding"

- The original BERT paper
- Introduced MLM and NSP

\item **Tenney et al. (2019)** - "BERT RedisCOVERS the Classical NLP Pipeline"

- Analysis of what BERT learns
- Layer-wise linguistic properties

\item **Clark et al. (2019)** - "What Does BERT Look At? An Analysis of BERT's

# Questions?

## Discussion Time

### Topics for discussion:

- Masked Language Modeling
- Pre-training vs fine-tuning
- Contextual embeddings
- BERT architecture details
- Implementation questions

Thank you! 

Next: BERT Variants and Improvements!