# Lecture 20: Applications of Encoder Models

## Week 6, Lecture 3 - From Theory to Practice

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

Winter 2026

# Today's Agenda 📋

1. 🚀 **Real-World Applications**: Where BERT shines

2. 🧠 **Cognitive Neuroscience**: Brain-model parallels

3. 🤔 **Understanding vs. Pattern Matching**: The big debate

4. ⚠️ **Limitations**: What BERT can't do

5. 💡 **Practical Tips**: Deployment and optimization

6. 🔮 **Future Directions**: Where are we heading?

*Goal: Connect BERT to real applications and understand broader implications*

# BERT Applications 🚀

**BERT excels at understanding tasks:**

**Classification Tasks:**

- Sentiment Analysis

- Topic Classification

- Spam Detection

- Intent Recognition

**Token-Level Tasks:**

- Named Entity Recognition (NER)

- Part-of-Speech Tagging

- Word Sense Disambiguation

# Case Study: Google Search 🔍

## BERT revolutionized search in 2019

```
1# Why word order matters: BERT understands prepositions!
2query = "2019 brazil traveler to usa need a visa"
3
4# Before BERT (bag-of-words matching):
5keywords = ["brazil", "traveler", "usa", "visa"]
6# Matches both: "US traveler to Brazil" AND "Brazil traveler to US"
7
8# With BERT (contextual understanding):
9bert_understanding = {
10    "subject": "brazil traveler",       # WHO is traveling
11    "destination": "usa",               # WHERE they're going
12    "direction": "brazil → usa",        # The preposition "to" is key!
13    "intent": "visa requirements"
14}
15# BERT correctly ranks: "Brazil citizen visa requirements for USA"
```

# Question Answering with BERT 💬

**Extractive QA: Find answer span in passage**

Example

**Context:** "The Normans (Norman: Nourmands; French: Normands; Latin: Normanni) were the people who in the 10th and 11th centuries gave their name to Normandy, a region in France."

**Question:** "In what country is Normandy located?"

**Answer:** France

```
1 from transformers import pipeline
2
3 # Load QA pipeline with BERT
4 qa_pipeline = pipeline("question-answering", model="bert-large-uncased-whole-wo
5
```

# Named Entity Recognition 🏷️

**Token-level classification task**

Example

**Input:** "Apple Inc. is headquartered in Cupertino, California."

**Output:**

- Apple Inc. → {ORGANIZATION}

- Cupertino → {LOCATION}

- California → {LOCATION}

```
1 from transformers import pipeline
2
3 # Load NER pipeline
4 ner_pipeline = pipeline("ner", model="dslim/bert-base-NER")
5
```

# Sentiment Analysis 😊😐😢

**Sequence classification task**

Examples

- "This movie was absolutely amazing!" → {POSITIVE}

- "The product broke after one week." → {NEGATIVE}

- "The weather is cloudy today." → {NEUTRAL}

```
1from transformers import pipeline
2
3# Load sentiment analysis pipeline
4sentiment_pipeline = pipeline("sentiment-analysis", model="distilbert-base-unca
5
6# Analyze sentiments
7texts = [
8    "This movie was absolutely amazing!",
9    "The product broke after one week."
```

7

# Semantic Similarity 🔗

## Measuring sentence similarity with BERT embeddings

```
1from transformers import BertTokenizer, BertModel
2import torch
3import torch.nn.functional as F
4
5tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
6model = BertModel.from_pretrained('bert-base-uncased')
7
8def get_sentence_embedding(sentence):
9    inputs = tokenizer(sentence, return_tensors='pt', padding=True, truncation=
10    outputs = model(**inputs)
11    # Use [CLS] token embedding as sentence representation
12    return outputs.last_hidden_state[:, 0, :]
13
14# Compare sentences
15sent1 = "The cat is sleeping on the couch"
16sent2 = "A feline is resting on the sofa"
```

# Semantic Similarity 🔗

```
21 emb3 = get_sentence_embedding(sent3)
22
23 # Compute cosine similarities
24 sim_12 = F.cosine_similarity(emb1, emb2).item()
25 sim_13 = F.cosine_similarity(emb1, emb3).item()
26
27 print(f"Similarity (1–2): {sim_12:.3f}")  # High (paraphrases)
28 print(f"Similarity (1–3): {sim_13:.3f}")  # Low (different topics)
```

...continued

# Cognitive Neuroscience Perspective 🧠

**How do brains and models process language?**

**Predictive Processing in the Brain:**

- Brain constantly predicts upcoming input

- N400: Neural response to unexpected words

- P600: Syntactic anomaly detection

- Context shapes predictions

- Prediction errors drive learning

**Key Brain Regions:**

- **Left IFG**: Syntax processing

- **Left STG/MTG**: Semantic processing

# Prediction in Brains vs. Language Models 🧠 🤖

**Parallels between neural and artificial systems**

| Phenomenon | Human Brain | Transformer Models |
|---|---|---|
| **Surprise** | N400 amplitude (EEG) | Cross-entropy loss |
| **Hierarchy** | sounds → words → sentences | tokens → phrases → meaning |
| **Context** | Prior discourse, world knowledge | Self-attention over sequence |
| **Representation** | Population coding (neurons) | Distributed embeddings (vectors) |

```
1# Concrete example: Surprise/N400 parallel
2sentence_a = "I take my coffee with cream and sugar"  # Expected
3sentence_b = "I take my coffee with cream and socks"  # Surprising
4
5# Brain: N400 amplitude higher for "socks"
```

# Neural Encoding with Language Models 🔬

**Can we predict brain activity from language models?**

```python
1# Neural encoding experiment workflow
2import numpy as np
3from transformers import BertModel
4
5# 1. Participant reads sentences while in fMRI scanner
6sentences = ["The dog chased the cat", "She opened the door", ...]
7brain_activity = fmri_scanner.record(sentences)  # (n_sentences, n_voxels)
8
9# 2. Extract BERT representations for same sentences
10bert = BertModel.from_pretrained("bert-base-uncased")
11bert_embeddings = []
12for sent in sentences:
13    outputs = bert(tokenizer(sent, return_tensors="pt"))
14    # Use layer 8 (found to correlate best with semantic areas)
15    bert_embeddings.append(outputs.hidden_states[8].mean(dim=1))
16
```

# Neural Encoding with Language Models 🔬

```
21# 4. Predict brain activity for new sentences
22predictions = encoder.predict(bert_embeddings[80:])
23correlation = np.corrcoef(predictions.flat, brain_activity[80:].flat)[0,1]
24# Correlation ~ 0.3–0.5 in language areas (significant!)
```

...continued

**Key finding:** BERT layer 8 best predicts semantic areas; layers 2-4 predict phonological areas

*Reference: Caucheteux & King (2022) - "Brains and algorithms partially converge"*

# Discussion: What Does the Model "Understand"? 🤔

**Does BERT understand language?**

**Evidence FOR understanding:**

- Captures syntax and semantics

- Resolves ambiguity

- Handles long-range dependencies

- Generalizes to new examples

- Predicts brain activity

- Solves complex tasks

*"If it acts like it understands, maybe it does?"*

**Evidence AGAINST understanding:**

# Adversarial Examples and Brittleness ⚠️

**BERT can be fooled easily**

```
1 from transformers import pipeline
2 classifier = pipeline("sentiment-analysis")
3
4 # Works correctly
5 classifier("This movie was absolutely wonderful!")
6 # → [{'label': 'POSITIVE', 'score': 0.9998}]
7
8 # Adding irrelevant negative words flips prediction!
9 classifier("This movie was absolutely wonderful! [SEP] bad bad bad bad")
10 # → [{'label': 'NEGATIVE', 'score': 0.9234}]  # WRONG!
11
12 # Synonym substitution can break it
13 classifier("The food was good")   # → POSITIVE (0.99)
14 classifier("The food was fine")   # → POSITIVE (0.72)  # Less confident
15 classifier("The food was ok")     # → NEGATIVE (0.51)  # WRONG!
16
```

# Limitations of Current Models ⚠️

**Despite impressive performance, transformers have limitations:**

1. **Quadratic Complexity**
   - Self-attention scales as $O(n^2)$

- Limited context windows (512-4096 tokens)

- Cannot process very long documents efficiently

2. **No True Understanding**
   - Pattern matching vs. comprehension

- Lack of common sense

- No world model

3. **Data Efficiency**

- Requires massive training data

# Bias in Language Models ⚖️

## Models reflect and can amplify societal biases

```
1from transformers import pipeline
2unmasker = pipeline("fill-mask", model="bert-base-uncased")
3
4# Gender bias in occupations
5unmasker("The doctor said [MASK] would be late.")
6# → [('he', 0.62), ('she', 0.18), ('it', 0.08), ...]
7
8unmasker("The nurse said [MASK] would be late.")
9# → [('she', 0.71), ('he', 0.15), ('it', 0.06), ...]
10
11# Racial bias (different sentiment for names)
12classifier = pipeline("sentiment-analysis")
13classifier("Emily is a brilliant scientist.")  # POSITIVE: 0.98
14classifier("Jamal is a brilliant scientist.")  # POSITIVE: 0.94   # Lower!
15
16# Where does bias come from?
```

# **Practical Tips for Working with Transformers** 💡

1. **Start with Pre-trained Models**

   ○ Don't train from scratch (too expensive!)

- Use HuggingFace Model Hub

- Choose appropriate model size

2. **Fine-tuning Best Practices**

   ○ Use small learning rate (1e-5 to 5e-5)

- Add warmup steps

- Monitor for overfitting

- Freeze early layers if data is limited

3. **Computational Efficiency**

   ○ Use mixed precision training (FP16)

# Deployment Considerations 🚀

## Moving from research to production

```
1# Example: Optimizing BERT for production deployment
2from transformers import BertModel, BertTokenizer
3import torch
4import onnxruntime
5
6# Step 1: Load model
7model = BertModel.from_pretrained("bert-base-uncased")
8tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
9
10# Step 2: Quantize for speed (INT8 instead of FP32)
11quantized_model = torch.quantization.quantize_dynamic(
12    model, {torch.nn.Linear}, dtype=torch.qint8
13)
14# Result: 4x smaller, 2x faster on CPU
15
16# Step 3: Export to ONNX for production
```

# Deployment Considerations 🚀

```
21session = onnxruntime.InferenceSession("bert.onnx")
22# 1.5x faster than PyTorch, works on any platform
```

...continued

| Optimization | Size | Latency | Quality |
|---|---|---|---|
| Original (FP32) | 420MB | 50ms | 100% |
| Quantized (INT8) | 110MB | 25ms | 99.5% |
| ONNX + Quantized | 110MB | 20ms | 99.5% |
| DistilBERT + ONNX | 65MB | 12ms | 97% |

# Future Directions 🔮

**Where is the field heading?**

1. **Longer Context**
   - Efficient attention mechanisms (linear, sparse)

- Models with 100K+ token context

- Better long-document understanding

2. **Multimodal Models**
   - Vision + Language (CLIP, DALL-E)

- Audio + Language (Whisper)

- Grounded understanding

3. **Better Pre-training**
   - More efficient objectives

# Encoder vs Decoder Models Revisited 🔄

**Different models for different tasks**

| p{4.5cm}} | Encoder (BERT) | Decoder (GPT) |
|---|---|---|
| • Classification | | |
| • NER, QA | | |
| • Similarity | Generation tasks: | |
| • Text completion | | |
| • Dialogue | | |
| • Creative writing | | |

**Interesting observation:**

22

• Decoder-only models (GPT-3, LLaMA) can also do classification via prompting!

# Discussion Questions 💭

1. **Understanding vs. Pattern Matching:**
    - Where do you draw the line?

- Is there a test for "true" understanding?

- Does it matter for applications?

2. **Brain-Model Parallels:**
    - How useful are these comparisons?

- What can neuroscience learn from AI?

- What can AI learn from neuroscience?

3. **Bias and Fairness:**
    - Who is responsible for addressing bias?

- Can we ever have completely unbiased models?

# Assignment 4: Context-Aware Models 📝

**Hands-on experience with transformers!**

**Tasks:**

1. **Implement Attention Mechanism**
   - Build scaled dot-product attention from scratch

- Visualize attention weights

2. **Fine-tune BERT**
   - Load pre-trained BERT

- Fine-tune on sentiment analysis

- Compare to baseline models

3. **Analyze Contextual Embeddings**

# Summary: Weeks 5-6 🎯

**What we learned:**

1. **Evolution of Context**
   - Seq2Seq → Attention → Transformers

- From bottleneck to full parallelization

2. **Transformer Architecture**
   - Self-attention, multi-head attention

- Positional encoding, layer norm, residuals

- Encoder-only (BERT), Decoder-only (GPT), Both (T5)

3. **BERT & Variants**
   - Masked Language Modeling

# Resources & Further Reading 📚

**Key Papers:**

- Vaswani et al. (2017) – Attention Is All You Need

- Devlin et al. (2019) – BERT: Pre-training of Deep Bidirectional Transformers

- Liu et al. (2019) – RoBERTa

- Sanh et al. (2019) – DistilBERT

- Dao et al. (2022) – FlashAttention

**Cognitive Neuroscience:**

- Hagoort & Indefrey (2014) – The neurobiology of language beyond single words

- Kuperberg & Jaeger (2016) – What do we mean by prediction in language comprehension?

- Willems et al. (2016) – Prediction during natural language comprehension

# Looking Forward in the Course 🔮

**Where do we go from here?**

**Upcoming Topics:**

- **Week 7:** Decoder models and text generation (GPT family)

- **Week 8:** Scaling laws and large language models

- **Week 9:** Prompting, in-context learning, and instruction tuning

- **Week 10:** Alignment, RLHF, and ethical considerations

**The Journey Continues:**

- From understanding (BERT) to generation (GPT)

- From supervised learning to few-shot learning

- From narrow tasks to general-purpose models

# Questions? 🙋

**Discussion Time**

**Topics to discuss:**

- BERT applications

- Understanding vs. pattern matching

- Cognitive neuroscience connections

- Limitations and future work

- Assignment 4 questions

Thank you! 🙏

See you in Week 7 for GPT and text generation!