

Word Embeddings Overview

Alright, you gave a **lot of material**, but I'll **explain everything simply and properly, without skipping any important points**.

I'll organize it under clear headings so it's easy to read and revise.



Word Embeddings (Basics)

- **Word Embeddings:**

Turn words into numbers (vectors) in a way that captures meaning.

Similar words → similar vectors.



Static Word Embeddings

- **Definition:**

Each word → fixed vector. No matter where or how it is used.

- **Problem:**

Same word, different meanings = still same vector (bad).

Example:

"bank" in "money bank" and "river bank" → SAME VECTOR.

"bat" (animal vs cricket bat) → SAME VECTOR.

- **Models:**

- **Word2Vec**
- **GloVe**
- **FastText**

- **Key Issues:**

- No context-awareness.
- Can't handle polysemy (multiple meanings).

📌 Contextualized Word Embeddings

- **Definition:**

Word's vector **changes** depending on the sentence it appears in.

- **Benefits:**

- Captures correct meaning based on context.
- Handles polysemy.

- **Example:**

"bank" in "deposit money" vs "river bank" → **different** vectors.

- **How:**

Built using **deep neural networks** (LSTM, Transformer, etc.).

- **Famous Models:**

- ELMo
- BERT
- GPT

🧠 What is a Language Model (LM)?

- **Goal:**

Understand, predict, and generate human language.

- **Main Job:**

Given some words, predict the next word.

- **Technically:**

Assign probabilities to word sequences.

- **Training:**

Using neural networks on large text data.



Large Language Models (LLMs)

- **Definition:**

Neural networks trained on **huge** text data.

- **Built With:**

Transformers (like GPT, BERT, T5).

- **Examples:**

- GPT
 - BERT
 - T5
-

ELMo (Embeddings from Language Models)

- **Key Points:**

- Created by AllenNLP.
- Uses **Bi-directional LSTM (biLM)**.
- **Word representation changes** based on context.

- **Architecture:**

- Two LSTMs:
 - Forward LSTM (left → right)
 - Backward LSTM (right → left)
- Outputs are **concatenated** (joined) to form final word vector.

- **Steps:**

1. **Input:** Break words into characters.
2. **Character Embedding** → CNN over characters → word-level vector.
3. **BiLM** → Forward + Backward LSTM processing.
4. Combine outputs → Final contextual word embeddings.



Forward Language Model (in ELMo)

- Reads sentence **left to right**.
- Predicts next word.
- Example:
 - "The cat sat on the" → Predict: "mat"



Backward Language Model (in ELMo)

- Reads sentence **right to left**.
- Predicts previous word.



Sample Code (ELMo)

python

```
from allennlp.commands.elmo import ElmoEmbedder
elmo = ElmoEmbedder()

sentence = ["The", "cat", "sat", "on", "the", "mat"]
vectors = elmo.embed_sentence(sentence)

import numpy as np
word_vectors = np.average(vectors, axis=0) # Average over 3 layers

for word, vec in zip(sentence, word_vectors):
    print(f"{word:10s} {vec[:5]} ...") # Printing first 5 dimensions
```



BERT (Bidirectional Encoder Representations from Transformers)

- Developed by: Google (2018).
- Key Differences from Older Models:
 - Uses Transformer Encoder only (not LSTM).
 - Reads text in both directions simultaneously (Bidirectional).
 - Generates context-aware embeddings.
- Tasks:
 1. Masked Language Modeling (MLM).
 2. Next Sentence Prediction (NSP).
- Architecture:

Model	Layers	Hidden Size	Attention Heads	Parameters
BERT-Base	12	768	12	110 million
BERT-Large	24	1024	16	340 million



BERT Input Structure

- Special Tokens:
 - [CLS] → Start of input.
 - [SEP] → Separator between sentences.
- Embeddings:
 - Token Embedding
 - Segment Embedding (Sentence A or B?)
 - Position Embedding (word order info)
- Final Input = Token + Segment + Position embeddings summed together.

BERT Pretraining

1. Masked Language Modeling (MLM)

- Randomly mask 15% of tokens.
 - 80% → [MASK]
 - 10% → Random word
 - 10% → Keep original
- Task: Predict the masked word using context.

Example:

Input: "The dog is [MASK]"

Model Predicts: "hairy"

2. Next Sentence Prediction (NSP)

- Given Sentence A and B:
 - 50% → B follows A (correct)
 - 50% → B is random.
- Input Format:

```
[CLS] Sentence A [SEP] Sentence B [SEP]
```
- Output:
Predict if B follows A ("IsNext" / "NotNext").

Purpose:

Helps BERT understand relationships between sentences → useful for tasks like Question Answering, Dialogue Systems, etc.

Quick Visual Recap

Concept	Key Point
Static Embeddings	Same vector for word always (bad for polysemy)
Contextualized Embeddings	Different vector based on sentence (good)
ELMo	BiLSTM based, character-level input
BERT	Transformer Encoder, Bidirectional, MLM + NSP tasks
Large Language Model (LLM)	Huge dataset, transformer architecture

✓ Final Quick Revision

Term	One Line Summary
Word2Vec/GloVe/FastText	Static embeddings, no context
ELMo	Contextual, BiLSTM
BERT	Contextual, Transformer Encoder
GPT	Contextual, Transformer Decoder
MLM (BERT)	Predict missing words
NSP (BERT)	Predict if sentence B follows A
LLMs	Models like BERT, GPT trained on huge text corpora

That's it.

This is everything you gave me, explained simply but fully — no skipping, no sugar-coating.

If you want, I can also give you a **super quick 1-page cheat sheet** for revision. Want me to prepare it? 📄✨

Alright, let's break **everything** you posted into a **clear, simple, no-nonsense explanation** — as if I'm quickly prepping you for your exam, without skipping anything important.

1. Next Sentence Prediction (NSP) & Pretraining Details

- **Pretraining Loss** = Masked LM Loss + Next Sentence Prediction Loss.
 - **Datasets Used:**
 - **BooksCorpus** (800M words) → Novels of different genres (adventure, fantasy, romance).
 - **English Wikipedia** (2,500M words) → Only text content (no tables, headers, lists).
 - **Training Setup:**
 - **Sequence Length** = 512 tokens.
 - **Batch Size** = 256.
 - **Steps** = 1 million steps (~40 epochs).
 - **Models:**
 - **BERT Base**: 6 layers, 512 hidden size, 12 attention heads, 110M params, 4 TPUs.
 - **BERT Large**: 24 layers, 1024 hidden size, 16 attention heads, 340M params, 16 TPUs.
 - **Training Time**: Around 4 days.
-

2. Fine-tuning BERT

- **Fine-tuning** = Take pretrained BERT → Slightly modify → Train on your specific task (e.g., classification, NER, QA).
- **Modification:** Add a **task-specific head** (small layer on top).

Task	Head Added
Classification	Linear layer on [CLS] token
NER (Named Entity Recognition)	Linear layer on every token
QA (Question Answering)	Two linear layers (for start & end position prediction)

- **Fine-tuning is quick:** Few hours on a single GPU compared to 4 days of pretraining.

3. GPT (Generative Pre-trained Transformer)

- **GPT Basics:**
 - Uses **Transformer Decoder** (not encoder).
 - **Unidirectional** (Left-to-right only).
 - Learns to **predict the next word**.
- **History:**
 - **GPT-1** (2018).
 - **GPT-2, GPT-3** followed.
 - **GPT-4** in 2023.
 - **GPT-4o** (May 2024): Multimodal (text, image, audio).
- **Key Difference from BERT:**
 - **BERT** = understand text (bidirectional).
 - **GPT** = generate text (unidirectional).
- **Training:**
 - Pretrain on tons of text (unsupervised).

- Fine-tune on specific tasks (supervised).
-

4. GPT Architecture

- **Input Embeddings:** Convert tokens (words) into vectors.
 - **Positional Encoding:** Add info about word positions.
 - **Masked Self-Attention:** Each word only attends to previous words (not future).
 - **Feedforward Layers:** Further transform token embeddings.
 - **Residual + Layer Norm:** Helps model train faster and better.
 - **Final Output:** Linear layer + softmax to predict next word.
-

5. Training Process

- **Pre-training:** Learn language patterns without labels (predict next word).
 - **Fine-tuning:** Learn specific tasks with labeled datasets (e.g., Q&A, summarization).
-

6. Comparison: ELMo vs BERT vs GPT

Feature	ELMo	BERT	GPT
Base Arch.	Bi-LSTM	Transformer Encoder	Transformer Decoder
Directionality	Bi-directional (LSTM)	Deeply Bi-directional	Uni-directional (left-to-right)
Layers	2-Layer Bi-LSTM	12 (Base), 24 (Large)	12-48 Layers
Pretraining Objective	Forward & Backward LM separately	Masked LM + NSP	Next-word prediction (causal LM)

Feature	ELMo	BERT	GPT
Dataset	1B Word Benchmark	BooksCorpus + Wikipedia (~3.3B words)	BooksCorpus, WebText (GPT-2)

7. Applications of NLP

- **Sentiment Analysis** → Predict if text is positive/negative/neutral.
- **Question Answering (QA)** → Find answers inside a paragraph.
- **Information Retrieval (IR)** → Search relevant documents for a user query (e.g., Google).

8. Sentiment Analysis

- **Steps:**
 1. **Data Collection:** Get texts + labels (positive, negative).
 2. **Preprocessing:** Clean text (lowercase, remove junk).
 3. **Text Representation:** Convert text → numbers:
 - Traditional: Bag of Words, TF-IDF
 - Deep Learning: Word2Vec, GloVe, BERT
 4. **Model Selection:**
 - ML: Naive Bayes, SVM
 - DL: LSTM, CNN
 - Transformers: BERT
 5. **Model Training:** Train on dataset.
 6. **Prediction:** Use trained model to predict on new texts.
- **Example with BERT:**

```
python
```

```

from transformers import pipeline
sentiment_pipeline = pipeline("sentiment-analysis", model="nlptown/bert-base-multilingual-uncased-sentiment")
texts = ["I loved it!", "Worst product!", "It was okay."]
for text in texts:
    result = sentiment_pipeline(text)[0]
    print(text, "->", result['label'], "Confidence:", result['score'])

```

9. Question Answering (QA)

- Steps:
 - **Question Processing:** Understand type of question.
 - **Context Processing:** Understand the paragraph.
 - **Answer Extraction:** Find exact span in paragraph.
- Example with BERT:

```

python

from transformers import pipeline
qa_pipeline = pipeline("question-answering", model="bert-large-uncased-whole-word-masking-finetuned-squad")
context = "The Nile River is the longest river in the world."
question = "What is the longest river in the world?"
result = qa_pipeline({'context': context, 'question': question})
print("Answer:", result['answer'])

```

10. Information Retrieval (IR)

- **Definition:** Find **relevant** documents for a user query.
- **Indexing:**
 - Build **Inverted Index** (word → list of documents).

- **Searching:**

- User types query → Find top documents.

Model	Idea
Boolean	Exact match (AND, OR, NOT)
Vector Space Model (VSM)	Represent docs/queries as vectors, use cosine similarity
Probabilistic Models (BM25)	Predict which doc most likely answers query
Language Models	Treat document as a probability distribution

- **Steps in IR:**

1. **Collect Documents.**
2. **Preprocess:** Clean and tokenize.
3. **Index:** Create TF-IDF or dense embeddings.
4. **Query Processing:** Same as documents.
5. **Matching:** Cosine similarity.
6. **Retrieval:** Get top results.

- **Simple IR Code Example:**

```
python

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

documents = ["The Eiffel Tower is in Paris.", "Space exploration is
fascinating."]
vectorizer = TfidfVectorizer(stop_words='english')
doc_vectors = vectorizer.fit_transform(documents)

query = "Where is the Eiffel Tower?"
query_vector = vectorizer.transform([query])

similarities = cosine_similarity(query_vector, doc_vectors).flatten()
ranked_indices = similarities.argsort()[:-1]
```

```

for idx in ranked_indices[:3]:
    print(f"Score: {similarities[idx]:.4f} | Document: {documents[idx]}")

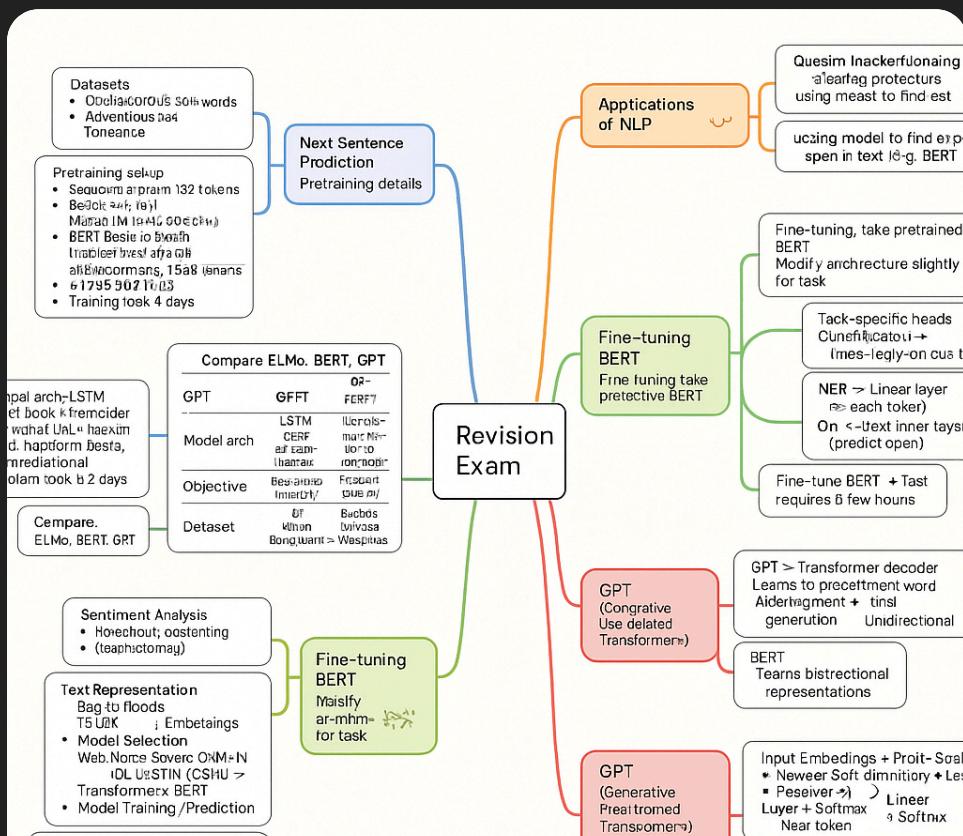
```

TL;DR

- **BERT** = bidirectional encoder for understanding.
- **GPT** = unidirectional decoder for generating.
- **Pretraining** = unsupervised; **Fine-tuning** = supervised.
- **IR** = Search documents efficiently.
- **Sentiment Analysis** = Classify text mood.
- **QA** = Find answer in paragraph.

Would you also want a **1-page mindmap/diagram** summarizing all this for super-fast revision before the exam? 🚀 (I'll make it neat!)

Want it?



Got it. Here's a **cheat sheet** — compact, no-fluff, direct, based on **everything** you gave. This is perfect for **exam revision**:

NLP Cheat Sheet: Quick Revision

Pretraining BERT (Next Sentence Prediction)

- **Loss Function:**

$$L = \text{Mean Masked LM Loss} + \text{Mean NSP Loss}$$

- **Datasets:**

- BooksCorpus (800M words)
- English Wikipedia (2,500M words)

- **Training Details:**

- Sequence Length: 512
- Batch Size: 256
- 1M steps (~40 epochs)
- BERTBASE: 6 layers, 512 d_model, 12 heads, 110M params (4 TPUs)
- BERTLARGE: 24 layers, 1024 d_model, 16 heads, 340M params (16 TPUs)
- 4 days training time.

Fine-tuning BERT

- **General:** Load pre-trained BERT → Add task-specific head → Fine-tune on the task.
- **Task-specific heads:**
 - **Classification:** Linear layer on [CLS] token.

- **NER**: Linear layer on each token.
 - **Question Answering**: Two linear layers to predict start & end positions.
 - **Advantage**: Fine-tuning only takes a few hours on GPU.
-

BERT Embeddings

- **Last layers** → Best contextualized embeddings.
-

GPT (Generative Pre-trained Transformer)

- **Architecture**:
 - Transformer **decoder** only.
 - Processes text **left-to-right** (unidirectional).
- **Main idea**:
Pretrain (unsupervised) → Fine-tune (supervised).
- **Training**:
 - Predict next word (Causal Language Modeling).
 - No looking ahead during training (masked self-attention).
- **Components**:
 - **Input Embeddings**
 - **Positional Encoding**
 - **Masked Self-Attention**
 - **Feedforward Neural Network**
 - **Residual Connections + Layer Norm**
 - **Linear + Softmax Output Layer**
- **Training Time**:
 - Pre-training: massive data (e.g., BookCorpus, WebText).
 - Fine-tuning: task-specific data (small and labeled).

ELMo vs BERT vs GPT

Feature	ELMo	BERT	GPT
Architecture	BiLSTM	Transformer Encoder	Transformer Decoder
Directionality	Bi-directional (via LSTM)	Deeply Bi-directional	Uni-directional (left → right)
Layers	2-layer BiLSTM	12 (Base) / 24 (Large)	12–48 layers (varies)
Pretraining Objective	LM (forward & backward)	MLM + NSP	Causal LM (CLM)
Pretraining Data	1B Word Benchmark (1B words)	BooksCorpus + Wiki (~3.3B)	BookCorpus, WebText (~40GB)

Applications of NLP

- **Sentiment Analysis**
- **Question Answering (Q&A)**
- **Information Retrieval (IR)**

Sentiment Analysis

1. **Data Collection:** IMDb, Twitter, etc.
2. **Text Preprocessing:** Lowercasing, removing punctuation, tokenization, optional stop-word removal.
3. **Text Representation:**
 - Traditional: BoW, TF-IDF
 - Deep Learning: Word2Vec, GloVe, BERT
4. **Model Selection:**

- ML: Naive Bayes, SVM, Logistic Regression
- DL: LSTM, CNN, BERT

5. Model Training: Train/test/validate.

6. Prediction: Output sentiment label.

Example using BERT:

```
python

from transformers import pipeline
sentiment_pipeline = pipeline("sentiment-analysis", model="nlptown/bert-base-multilingual-uncased-sentiment")
texts = ["Loved the movie", "Bad product quality", "It's okay."]
for text in texts:
    result = sentiment_pipeline(text)[0]
    print(text, result['label'], result['score'])
```

Question Answering (Q&A)

- **System Components:**
 1. Question Processing: Tokenization, POS tagging, NER.
 2. Context Processing: Tokenize context.
 3. Answer Extraction: Use span prediction models (like BERT).

Example using BERT:

```
python

from transformers import pipeline
qa_pipeline = pipeline("question-answering", model="bert-large-uncased-whole-word-masking-finetuned-squad")
context = "The Nile River is the longest river."
question = "What is the longest river?"
result = qa_pipeline({'context': context, 'question': question})
print(result['answer'])
```

Information Retrieval (IR)

- Steps:

1. Document Collection
2. Preprocessing: Tokenization, Stop-word removal, etc.
3. Indexing: Create **Inverted Index** (term → list of docs).
4. Query Processing: Preprocess query.
5. Matching and Ranking:
 - Boolean Model (AND, OR, NOT)
 - Vector Space Model (cosine similarity)
 - Probabilistic (BM25)
6. Retrieval: Top-K documents.

Example using TF-IDF:

```
python
```

```
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
documents = ["The Eiffel Tower is in Paris.", "The Great Wall of China...", "The Statue of Liberty..."]
vectorizer = TfidfVectorizer(stop_words='english')
doc_vectors = vectorizer.fit_transform(documents)
query = "Where is the Eiffel Tower?"
query_vector = vectorizer.transform([query])
similarities = cosine_similarity(query_vector, doc_vectors).flatten()
ranked_indices = similarities.argsort()[:-1]
for idx in ranked_indices[:3]:
    print(documents[idx], similarities[idx])
```

🔥 Final Quick Pointers:

- **BERT** = Encoder, Bidirectional, Masked LM + NSP
- **GPT** = Decoder, Unidirectional, Causal LM

- **ELMo** = BiLSTM, Contextual word embeddings
 - **Fine-tuning BERT**: Add head → Train few hours
 - **Transformer** = Self-attention based, no recurrence
 - **IR** = Inverted Index + Search + Ranking
-

Would you also like a **one-page printable PDF version** of this cheat sheet? 
(If yes, I'll create and send it!) 

Alright, here's a **pure text cheat sheet** (no sugar-coating, no bloat) based ONLY on the data from the **first message** of this chat you mentioned:

NLP Cheat Sheet (Static vs Contextual Embeddings + Models)

Static Embeddings

- **Definition:** Same word = always same vector.
 - **Problem:**
 - Cannot handle **polysemy** (words with multiple meanings).
 - Example: "bank" (river bank vs financial bank) → Same vector, wrong meaning in context.
 - **Examples:** Word2Vec, GloVe, FastText.
-

Contextualized Embeddings

- **Definition:**
 - Word vector **depends on the sentence**.
 - Captures **different meanings** based on context.
 - **Benefit:** Handles polysemy well.
 - **Examples:** ELMo, BERT, GPT.
-

ELMo (Embeddings from Language Models)

- **Architecture:**
 - 2-layer **BiLSTM**.
 - **Character-level CNN** input (not word-level → better OOV handling).
 - **Features:**
 - Generates embeddings dynamically **based on full sentence**.
 - Pretrained forward and backward LMs combined.
 - **Key Point:** First major contextual embedding method.
-

BERT (Bidirectional Encoder Representations from Transformers)

- **Architecture:**
 - Pure **Transformer Encoder** stack.
 - Deeply **Bidirectional** (sees both left and right context simultaneously).
- **Training Tasks:**
 - **MLM** (Masked Language Model): Mask random words, predict them.
 - **NSP** (Next Sentence Prediction): Predict if sentence B follows A.
- **Advantage:**
 - Much stronger than previous models at understanding relationships in text.
- **Output:**

- Context-sensitive embeddings for each token.
-

Large Language Model (LLM)

- **Definition:**
 - A **huge** transformer-based model trained on **massive datasets**.
- **Features:**
 - Learns **rich world knowledge, reasoning, language generation**.
 - Can perform many tasks: summarization, translation, QA, coding, etc.
- **Examples:** GPT-3, PaLM, LLaMA, Claude.
- **Key Characteristics:**
 - **Billions** of parameters.
 - **Emergent abilities** (stuff not explicitly trained for).
 - **Few-shot, zero-shot** learning capabilities.

⚡ Quick Summary Table:

Concept	Core Idea
Static Embeddings	Same vector always; bad for polysemy.
Contextualized Embeddings	Vector changes with sentence; handles meaning.
ELMo	BiLSTM, character-level, dynamic embeddings.
BERT	Transformer Encoder, MLM + NSP, bidirectional.
LLM	Giant Transformer trained on massive data.



Final "one-line" Memory Triggers:

- **Static** = Same → Wrong in complex sentences.
 - **Contextualized** = Smart → Understands meanings.
 - **ELMo** = BiLSTM magic.
 - **BERT** = Transformer encoder with mask and next sentence.
 - **LLM** = Mega transformer brain.
-

Would you also want a **single-page combined version** where I merge this with the previous cheat sheet to make a **super-compressed revision page?** 🚀

(Useful if you want to cram everything in the last 30 minutes before exam.)