



GenAI: Best Practices

Release 1.0

Wenqiang Feng and Di Zhen

December 19, 2024

CONTENTS

1 Preface	3
1.1 About	3
1.1.1 About this book	3
1.1.2 About the authors	3
1.2 Feedback and suggestions	4
2 Preliminary	5
2.1 Math Preliminary	5
2.1.1 Vector	5
2.1.2 Norm	7
2.1.3 Distances	7
2.2 NLP Preliminary	9
2.2.1 Vocabulary	9
2.2.2 Tagging	10
2.2.3 Lemmatization	11
2.2.4 Tokenization	13
2.2.5 BERT Tokenization	14
2.3 Platform and Packages	16
2.3.1 Google Colab	16
2.3.2 HuggingFace	18
2.3.3 Ollama	18
2.3.4 langchain	20
3 Word and Sentence Embedding	23
3.1 Traditional word embeddings	23
3.1.1 One Hot Encoder	24
3.1.2 CountVectorizer	25
3.1.3 TF-IDF	26
3.2 Static word embeddings	30
3.2.1 Word2Vec	31
3.2.2 GloVe	33
3.2.3 Fast Text	34
3.3 Contextual word embeddings	35
3.3.1 BERT	35
3.3.2 gte-large-en-v1.5	36

3.3.3	bge-base-en-v1.5	37
4	Prompt Engineering	41
4.1	Prompt	41
4.2	Prompt Engineering	41
4.2.1	What's Prompt Engineering	41
4.2.2	Key Elements of a Prompt	41
4.3	Advanced Prompt Engineering	42
4.3.1	Role Assignment	42
4.3.2	Contextual Setup	43
4.3.3	Explicit Instructions	44
4.3.4	Chain of Thought (CoT) Prompting	45
4.3.5	Few-Shot Prompting	46
4.3.6	Iterative Prompting	47
4.3.7	Instructional Chaining	49
4.3.8	Use Constraints	50
4.3.9	Creative Prompting	50
4.3.10	Feedback Incorporation	50
4.3.11	Scenario-Based Prompts	50
4.3.12	Multimodal Prompting	51
5	Retrieval-Augmented Generation	53
5.1	Overview	54
5.2	Indexing	54
5.2.1	Naive Chunking	54
5.2.2	Late Chunking	58
5.2.3	Types of Indexing	61
5.2.4	Vector Database	62
5.3	Retrieval	64
5.3.1	Common retrieval methods	66
5.3.2	Reciprocal Rank Fusion	66
5.4	Generation	68
5.5	Advanced Topic	68
5.5.1	Corrective RAG	68
5.5.2	Self-RAG	68
5.5.3	Adaptive RAG	68
5.5.4	Agentic RAG	68
6	Fine Tuning	69
6.1	Cutting-Edge Strategies for LLM Fine-Tuning	70
6.1.1	LoRA (Low-Rank Adaptation)	70
6.1.2	QLoRA (Quantized Low-Rank Adaptation)	70
6.1.3	PEFT (Parameter-Efficient Fine-Tuning)	71
6.1.4	SFT (Supervised Fine-Tuning)	71
6.1.5	Summary Table	72
6.2	Key Early Fine-Tuning Methods	72
6.2.1	Full Fine-Tuning	72
6.2.2	Feature-Based Approach	73

6.2.3	Layer-Specific Fine-Tuning	73
6.2.4	Task-Adaptive Pre-training	73
6.3	Embedding Model Fine-Tuning	73
6.3.1	Prepare Dataset	74
6.3.2	Import and Evaluate Pretrained Baseline Model	75
6.3.3	Loss Function with Matryoshka Representation	77
6.3.4	Fine-tune Embedding Model	79
6.3.5	Evaluate Fine-tuned Model	80
6.3.6	Results Comparison	80
6.4	LLM Fine-Tuning	81
6.4.1	Load Dataset and Pretrained Model	81
6.4.2	Fine-tuning Configuration	82
6.4.3	Fine-tune model	83
7	Pre-training	85
8	LLM Evaluation Metrics	87
8.1	Statistical Scorers (Traditional Metrics)	88
8.2	Model-Based Scorers (Learned Metrics)	88
8.3	Human-Centric Evaluations (Augmenting Metrics)	89
8.4	GEval with DeepEval	89
8.4.1	G-Eval Algorithm	89
8.4.2	G-Eval with DeepEval	89
9	Main Reference	101
Bibliography		103



Welcome to our **GenAI: Best Practices!!!** For each chapter, we provide detailed Colab notebooks [!\[\]\(d84e7ea36f695d92cb39ec32c307ac93_img.jpg\) Open in Colab](#) that you can open and run directly in Google Colab. The PDF version can be downloaded from [HERE](#).

CHAPTER ONE

PREFACE

1.1 About

1.1.1 About this book

This is the book for our Generative AI: Best practices [GenAI]. The PDF version can be downloaded from [HERE](#). You may download and distribute it. Please beware, however, that the note contains typos as well as inaccurate or incorrect description.

In this book, we aim to demonstrate best practices for Generative AI through detailed demo code and practical examples. For each chapter, we provide detailed Colab notebooks  [Open in Colab](#) that you can open and run directly in Google Colab.

1.1.2 About the authors

- **Wenqiang Feng**

- Sr. Mgr Data Enginner and PhD in Mathematics
- University of Tennessee at Knoxville
- Webpage: <http://web.utk.edu/~wfeng1/>
- Email: von198@gmail.com

- **Biography**

Wenqiang Feng is the Senior Manager of Data Engineering and former Director of AI Engineering/Data Science at American Express (AMEX). Before his tenure at AMEX, Dr. Feng served as a Senior Data Scientist in the Machine Learning Lab at H&R Block and as a Data Scientist at Applied Analytics Group, DST (now SS&C). Throughout his career, Dr. Feng has focused on equipping clients with cutting-edge skills and technologies, including Big Data analytics, advanced modeling techniques, and data enhancement strategies.

Dr. Feng brings extensive expertise in data mining, analytic systems, machine learning algorithms, business intelligence, and the application of Big Data tools to solve complex, cross-functional industry challenges. Prior to his role at DST, Dr. Feng was an IMA Data Science Fellow at the Institute for Mathematics and its Applications (IMA) at the University of Minnesota. In this capacity, he collaborated with startups to develop predictive analytics solutions that informed strategic marketing decisions.

Dr. Feng holds a Ph.D. in Computational Mathematics and a Master's degree in Statistics from the University of Tennessee, Knoxville. He also earned a Master's degree in Computational Mathematics from Missouri University of Science and Technology (MST) and a Master's degree in Applied Mathematics from the University of Science and Technology of China (USTC).

- **Declaration**

The work of Wenqiang Feng was supported by the IMA, while working at IMA. However, any opinion, finding, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the IMA, UTK and DST.

Warning

ChatGPT has been extensively used in the creation of this book. If you notice that your work has not been cited or has been cited incorrectly, please notify us.

1.2 Feedback and suggestions

Your comments and suggestions are highly appreciated. I am more than happy to receive corrections, suggestions or feedback through email (Wenqiang Feng: von198@gmail.com and Di Zhen: dizhen318@gmail.com) for improvements.

CHAPTER TWO

PRELIMINARY

Chinese proverb

A journey of a thousand miles begins with a single step. – Lao Tzu

In this chapter, we will introduce some math and NLP preliminaries which are highly used in Generative AI.

Colab Notebook for This Chapter

- Math Preliminary:  [Open in Colab](#)
- Ollama in Colab:  [Open in Colab](#)
- BERT Tokenization:  [Open in Colab](#)

2.1 Math Preliminary

2.1.1 Vector

A vector is a mathematical representation of data characterized by both magnitude and direction. In this context, each data point is represented as a feature vector, with each component corresponding to a specific feature or attribute of the data.

```
import numpy as np
import gensim.downloader as api
# Download pre-trained GloVe model
glove_vectors = api.load("glove-twitter-25")

# Get word vectors (embeddings)
word1 = "king"
word2 = "queen"

# embedding
king = glove_vectors[word1]
```

(continues on next page)

(continued from previous page)

```
queen = glove_vectors[word2]

print('king:\n', king)
print('queen:\n', queen)
```

```
king:
[-0.74501 -0.11992  0.37329  0.36847 -0.4472  -0.2288   0.70118
 0.82872  0.39486 -0.58347  0.41488  0.37074 -3.6906  -0.20101
 0.11472 -0.34661  0.36208  0.095679 -0.01765  0.68498 -0.049013
 0.54049 -0.21005 -0.65397  0.64556 ]

queen:
[-1.1266  -0.52064  0.45565  0.21079 -0.05081 -0.65158  1.1395
 0.69897 -0.20612 -0.71803 -0.02811  0.10977 -3.3089  -0.49299
 -0.51375  0.10363 -0.11764 -0.084972  0.02558  0.6859  -0.29196
 0.4594  -0.39955 -0.40371  0.31828 ]
```

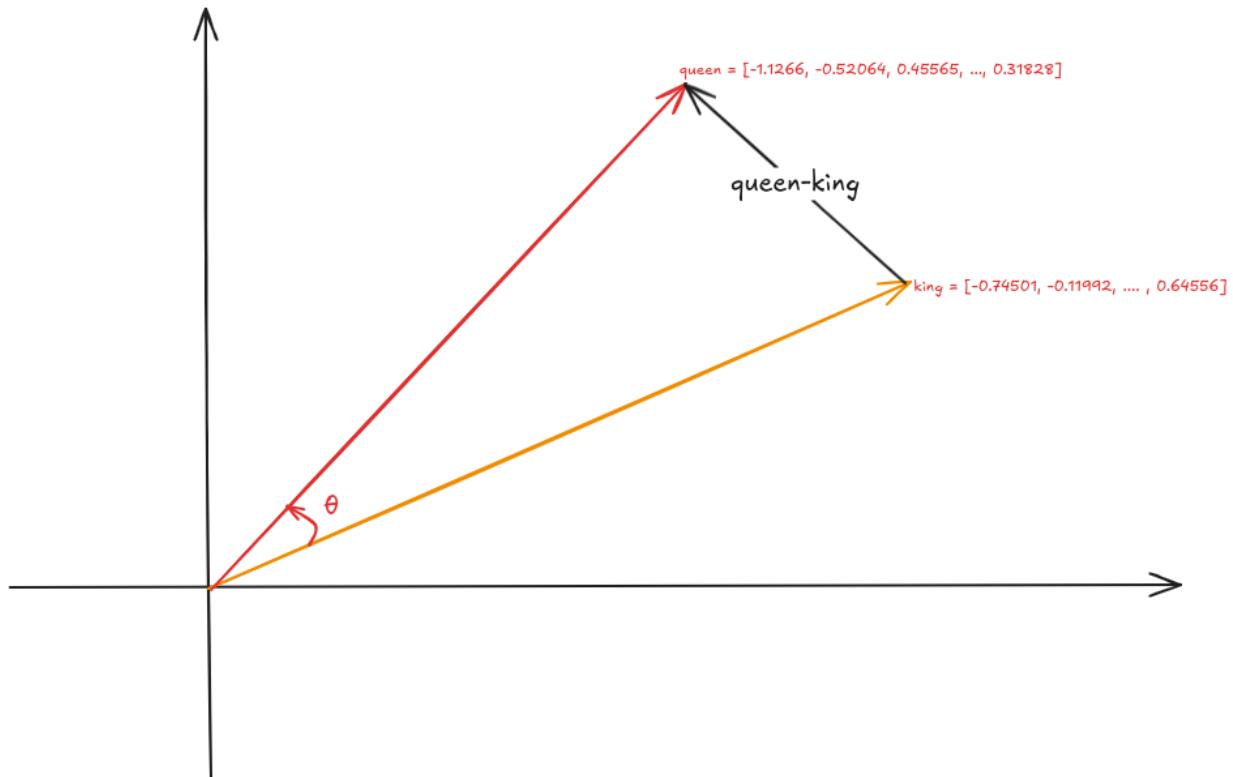


Fig. 1: Vector

2.1.2 Norm

A norm is a function that maps a vector to a single positive value, representing its magnitude. Norms are essential for calculating distances between vectors, which play a crucial role in measuring prediction errors, performing feature selection, and applying regularization techniques in models.

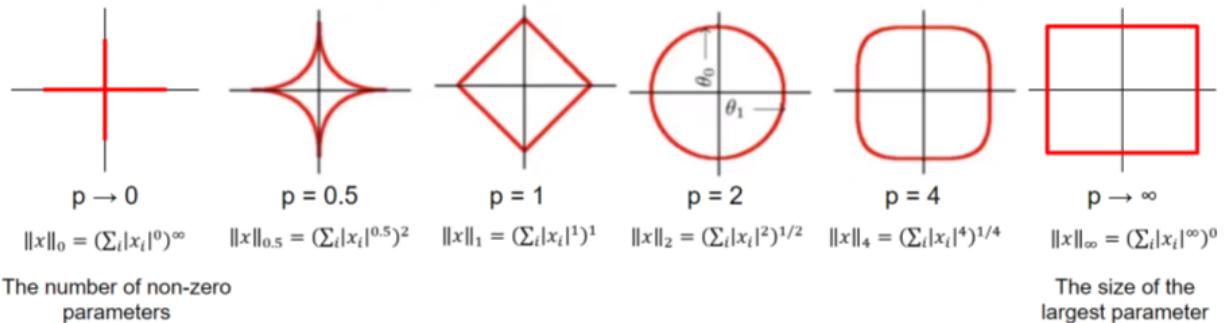


Fig. 2: Geometrical Interpretation of Norm (source_1)

- Formula:

The ℓ^p norm for $\vec{v} = (v_1, v_2, \dots, v_n)$ is

$$\|\vec{v}\|_p = \sqrt[p]{|v_1|^p + |v_2|^p + \dots + |v_n|^p}$$

- ℓ^1 norm: Sum of absolute values of vector components, often used for feature selection due to its tendency to produce sparse solutions.

```
# 11 norm
np.linalg.norm(king, ord=1) # max(sum(abs(x), axis=0))

### 13.188952
```

- ℓ^2 norm: Square root of the sum of squared vector components, the most common norm used in many machine learning algorithms.

```
# 12 norm
np.linalg.norm(king, ord=2)

### 4.3206835
```

- ℓ^∞ norm (Maximum norm): The largest absolute value of a vector component.

2.1.3 Distances

- Manhattan Distance (ℓ^1 Distance)

Also known as taxicab or city block distance, Manhattan distance measures the absolute differences between the components of two vectors. It represents the distance a point would travel along grid lines in a Cartesian plane, similar to navigating through city streets.

For two vector $\vec{u} = (u_1, u_2, \dots, u_n)$ and $\vec{v} = (v_1, v_2, \dots, v_n)$, the Manhattan Distance distance $d(\vec{u}, \vec{v})$ is

$$d(\vec{u}, \vec{v}) = \|\vec{u} - \vec{v}\|_1 = |u_1 - v_1| + |u_2 - v_2| + \dots + |u_n - v_n|$$

- Euclidean Distance (ℓ^2 Distance)

Euclidean distance is the most common way to measure the distance between two points (vectors) in space. It is essentially the straight-line distance between them, calculated using the Pythagorean theorem.

For two vector $\vec{u} = (u_1, u_2, \dots, u_n)$ and $\vec{v} = (v_1, v_2, \dots, v_n)$, the Euclidean Distance distance $d(\vec{u}, \vec{v})$ is

$$d(\vec{u}, \vec{v}) = \|\vec{u} - \vec{v}\|_2 = \sqrt{(u_1 - v_1)^2 + (u_2 - v_2)^2 + \dots + (u_n - v_n)^2}$$

- Minkowski Distance (ℓ^p Distance)

Minkowski distance is a generalization of both Euclidean and Manhattan distances. It incorporates a parameter, p , which allows for adjusting the sensitivity of the distance metric.

- Cos Similarity

Cosine similarity measures the angle between two vectors rather than their straight-line distance. It evaluates the similarity of two vectors by focusing on their orientation rather than their magnitude. This makes it particularly useful for high-dimensional data, such as text, where the direction of the vectors is often more significant than their magnitude.

The Cos similarity for two vector $\vec{u} = (u_1, u_2, \dots, u_n)$ and $\vec{v} = (v_1, v_2, \dots, v_n)$ is

$$\cos(\theta) = \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|}$$

- 1 means the vectors point in exactly the same direction (perfect similarity).
- 0 means they are orthogonal (no similarity).
- -1 means they point in opposite directions (complete dissimilarity).

```
# Compute cosine similarity between the two word vectors
np.dot(king,queen)/(np.linalg.norm(king)*np.linalg.norm(queen))

### 0.92024213
```

```
# Compute cosine similarity between the two word vectors
similarity = glove_vectors.similarity(word1, word2)
print(f"Word vectors for '{word1}': {king}")
print(f"Word vectors for '{word2}': {queen}")
print(f"Cosine similarity between '{word1}' and '{word2}':
    ↪{similarity}")
```

```

Word vectors for 'king': [-0.74501 -0.11992  0.37329  0.36847 -
                         ↵0.4472  -0.2288   0.70118
                         0.82872  0.39486 -0.58347  0.41488  0.37074 -3.6906  -0.20101
                         0.11472  -0.34661  0.36208  0.095679 -0.01765  0.68498  -0.049013
                         0.54049  -0.21005 -0.65397  0.64556 ]
Word vectors for 'queen': [-1.1266   -0.52064   0.45565  0.21079 -
                           ↵0.05081  -0.65158   1.1395
                           0.69897  -0.20612 -0.71803 -0.02811  0.10977 -3.3089  -0.49299
                           -0.51375  0.10363 -0.11764 -0.084972  0.02558  0.6859  -0.29196
                           0.4594   -0.39955 -0.40371  0.31828 ]
Cosine similarity between 'king' and 'queen': 0.920242190361023

```

2.2 NLP Preliminary

2.2.1 Vocabulary

In Natural Language Processing (NLP), **vocabulary** refers to the complete set of unique words or tokens that a model recognizes or works with during training and inference. Vocabulary plays a critical role in text processing and understanding, as it defines the scope of linguistic units a model can handle.

- Types of Vocabulary in NLP
 1. **Word-level Vocabulary:** - Each word in the text is treated as a unique token. - For example, the sentence “I love NLP” would generate the vocabulary: {I, love, NLP}.
 2. **Subword-level Vocabulary:** - Text is broken down into smaller units like prefixes, suffixes, or character sequences. - For example, the word “loving” might be split into {lov, ing} using techniques like Byte Pair Encoding (BPE) or SentencePiece. - Subword vocabularies handle rare or unseen words more effectively.
 3. **Character-level Vocabulary:** - Each character is treated as a token. - For example, the word “love” would generate the vocabulary: {l, o, v, e}.
- Importance of Vocabulary
 1. **Text Representation:** - Vocabulary is the basis for converting text into numerical representations like one-hot vectors, embeddings, or input IDs for machine learning models.
 2. **Model Efficiency:** - A larger vocabulary increases the model’s memory and computational requirements. - A smaller vocabulary may lack the capacity to represent all words effectively, leading to a loss of meaning.
 3. **Handling Out-of-Vocabulary (OOV) Words:** - Words not present in the vocabulary are either replaced with a special token like <UNK> or processed using subword/character-based techniques.
- Building a Vocabulary

Common practices include:

 1. Tokenizing the text into words, subwords, or characters.

2. Counting the frequency of tokens.
 3. Keeping only the most frequent tokens up to a predefined size (e.g., top 50,000 tokens).
 4. Adding special tokens like <PAD>, <UNK>, <BOS> (beginning of sentence), and <EOS> (end of sentence).
- Challenges
 - **Balancing Vocabulary Size:** A larger vocabulary increases the richness of representation but requires more computational resources.
 - **Domain-specific Vocabularies:** In specialized fields like medicine or law, standard vocabularies may not be sufficient, requiring domain-specific tokenization strategies.

2.2.2 Tagging

Tagging in NLP refers to the process of assigning labels or annotations to words, phrases, or other linguistic units in a text. These labels provide additional information about the syntactic, semantic, or structural role of the elements in the text.

- Types of Tagging
 1. **Part-of-Speech (POS) Tagging:**
 - Assigns grammatical tags (e.g., noun, verb, adjective) to each word in a sentence.
 - Example: For the sentence “The dog barks,” the tags might be: - The/DET (Determiner) - dog/NOUN (Noun) - barks/VERB (Verb).
 2. **Named Entity Recognition (NER) Tagging:**
 - Identifies and classifies named entities in a text, such as names of people, organizations, locations, dates, or monetary values.
 - Example: In the sentence “John works at Google in California,” the tags might be: - John/PERSON - Google/ORGANIZATION - California/LOCATION.
 3. **Chunking (Syntactic Tagging):**
 - Groups words into syntactic chunks like noun phrases (NP) or verb phrases (VP).
 - Example: For the sentence “The quick brown fox jumps,” a chunking result might be: - [NP The quick brown fox] [VP jumps].
 4. **Sentiment Tagging:**
 - Assigns sentiment labels (e.g., positive, negative, neutral) to words, phrases, or entire documents.
 - Example: The word “happy” might be tagged as positive, while “sad” might be tagged as negative.
 5. **Dependency Parsing Tags:**
 - Identifies the grammatical relationships between words in a sentence, such as subject, object, or modifier.

- **Example:** In “She enjoys cooking,” the tags might show:
 - * She/nsubj (nominal subject)
 - * enjoys/ROOT (root of the sentence)
 - * cooking/dobj (direct object).
- Importance of Tagging
 - **Understanding Language Structure:** Tags help NLP models understand the grammatical and syntactic structure of text.
 - **Improving Downstream Tasks:** Tagging is foundational for tasks like machine translation, sentiment analysis, question answering, and summarization.
 - **Feature Engineering:** Tags serve as features for training machine learning models in text classification or sequence labeling tasks.
- Tagging Techniques
 1. **Rule-based Tagging:** Relies on predefined linguistic rules to assign tags. Example: Using dictionaries or regular expressions to match specific patterns.
 2. **Statistical Tagging:** Uses probabilistic models like Hidden Markov Models (HMMs) to predict tags based on word sequences.
 3. **Neural Network-based Tagging:** Employs deep learning models like LSTMs, GRUs, or Transformers to tag text with high accuracy.
- Challenges
 - **Ambiguity:** Words with multiple meanings can lead to incorrect tagging. Example: The word “bank” could mean a financial institution or a riverbank.
 - **Domain-Specific Language:** General tagging models may fail to perform well on specialized text like medical or legal documents.
 - **Data Sparsity:** Rare words or phrases may lack sufficient training data for accurate tagging.

2.2.3 Lemmatization

Lemmatization in NLP is the process of reducing a word to its base or dictionary form, known as the **lemma**. Unlike stemming, which simply removes word suffixes, lemmatization considers the context and grammatical role of the word to produce a linguistically accurate root form.

- How Lemmatization Works
 1. **Contextual Analysis:**
 - Lemmatization relies on a vocabulary (lexicon) and morphological analysis to identify a word’s base form.
 - For example: - running → run - better → good
 2. **Part-of-Speech (POS) Tagging:**
 - The process uses POS tags to determine the correct lemma for a word.

- Example: - barking (verb) → bark - barking (adjective, as in “barking dog”) → barking.
- Importance of Lemmatization
 - 1. **Improves Text Normalization:**
 - Lemmatization helps normalize text by grouping different forms of a word into a single representation.
 - Example: - run, running, and ran → run.
 - 2. **Enhances NLP Applications:**
 - Lemmatized text improves the performance of tasks like information retrieval, text classification, and sentiment analysis.
 - 3. **Reduces Vocabulary Size:**
 - By mapping inflected forms to their base form, lemmatization reduces redundancy in text, resulting in a smaller vocabulary.
- Lemmatization vs. Stemming
 - **Lemmatization:**
 - * Produces linguistically accurate root forms.
 - * Considers the word’s context and POS.
 - * Example: - studies → study.
 - **Stemming:**
 - * Applies heuristic rules to strip word suffixes without considering context.
 - * May produce non-dictionary forms.
 - * Example: - studies → studi.
- Techniques for Lemmatization
 - 1. **Rule-Based Lemmatization:**
 - Relies on predefined linguistic rules and dictionaries.
 - Example: WordNet-based lemmatizers.
 - 2. **Statistical Lemmatization:**
 - Uses probabilistic models to predict lemmas based on the context.
 - 3. **Deep Learning-Based Lemmatization:**
 - Employs neural networks and sequence-to-sequence models for highly accurate lemmatization in complex contexts.
- Challenges
 - **Ambiguity:** Words with multiple meanings may result in incorrect lemmatization without proper context.
 - * Example: - left (verb) → leave - left (noun/adjective) → left.

- **Language-Specific Complexity:** Lemmatization rules vary widely across languages, requiring language-specific tools and resources.
- **Resource Dependency:** Lemmatizers require extensive lexicons and morphological rules, which can be resource-intensive to develop.

2.2.4 Tokenization

Tokenization in NLP refers to the process of splitting a text into smaller units, called **tokens**, which can be words, subwords, sentences, or characters. These tokens serve as the basic building blocks for further analysis in NLP tasks.

- Types of Tokenization

1. **Word Tokenization:**

- Splits the text into individual words or terms.
- **Example:**

- * Sentence: "I love NLP."
- * Tokens: ["I", "love", "NLP"].

2. **Sentence Tokenization:**

- Divides a text into sentences.

- **Example:**

- * Text: "I love NLP. It's amazing."
- * Tokens: ["I love NLP.", "It's amazing."].

3. **Subword Tokenization:**

- Breaks words into smaller units, often using methods like Byte Pair Encoding (BPE) or SentencePiece.

- **Example:**

- * Word: unhappiness.
- * Tokens: ["un", "happiness"] (or subword units like ["un", "happi", "ness"]).

4. **Character Tokenization:**

- Treats each character in a word as a separate token.

- **Example:**

- * Word: hello.
- * Tokens: ["h", "e", "l", "l", "o"].

- Importance of Tokenization

1. **Text Preprocessing:**

- Tokenization is the first step in many NLP tasks like text classification, translation, and summarization, as it converts text into manageable pieces.

2. Text Representation:

- Tokens are converted into numerical representations (e.g., word embeddings) for model input in tasks like sentiment analysis, named entity recognition (NER), or language modeling.

3. Improving Accuracy:

- Proper tokenization ensures that a model processes text at the correct granularity (e.g., words or subwords), improving accuracy for tasks like machine translation or text generation.

- Challenges of Tokenization

1. Ambiguity:

- Certain words or phrases can be tokenized differently based on context.
- Example: “New York” can be treated as one token (location) or two separate tokens (["New", "York"]).

2. Handling Punctuation:

- Deciding how to treat punctuation marks can be challenging. For example, should commas, periods, or quotes be treated as separate tokens or grouped with adjacent words?

3. Multi-word Expressions (MWEs):

- Some expressions consist of multiple words that should be treated as a single token, such as “New York” or “machine learning.”

- Techniques for Tokenization

1. **Rule-Based Tokenization:** Uses predefined rules to split text based on spaces, punctuation, and other delimiters.

2. **Statistical and Machine Learning-Based Tokenization:** Uses trained models to predict token boundaries based on patterns learned from large corpora.

3. **Deep Learning-Based Tokenization:** Modern tokenization models, such as those used in transformers (e.g., BERT, GPT), may rely on subword tokenization and neural networks to handle complex tokenization tasks.

2.2.5 BERT Tokenization

- Vocabulary: The BERT Tokenizer’s vocabulary contains 30,522 unique tokens.

```
from transformers import BertTokenizer, BertModel
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
# model = BertModel.from_pretrained("bert-base-uncased")

# vocabulary size
print(tokenizer.vocab_size)
```

(continues on next page)

(continued from previous page)

```
# vocabulary
print(tokenizer.vocab)
```

```
# vocabulary size
30522

# vocabulary
OrderedDict([('PAD', 0), ('unused0', 1)
              ....,
              ('writing', 3015), ('bay', 3016),
              ....,
              ('##?', 30520), ('##~', 30521)])
```

- Tokens and IDs

- Tokens to IDs

```
text = "Gen AI is awesome"
encoded_input = tokenizer(text, return_tensors='pt')

# tokens to ids
print(encoded_input)

# output
{'input_ids': tensor([[ 101,  8991,  9932,  2003, 12476,  102]]), \
'token_type_ids': tensor([[0, 0, 0, 0, 0, 0]]), \
'attention_mask': tensor([[1, 1, 1, 1, 1, 1]])}
```

You might notice that there are only four words, yet we have six token IDs. This is due to the inclusion of two additional special tokens [CLS] and [SEP].

```
print({x : tokenizer.encode(x, add_special_tokens=False) for x in [
    '[CLS]' + text.split() + '[SEP']])

### output
{'[CLS]': [101], 'Gen': [8991], 'AI': [9932], 'is': [2003], 'awesome': [12476], '[SEP]': [102]}
```

- Special Tokens

```
# Special tokens
print({x : tokenizer.encode(x, add_special_tokens=False) for x in [
    '[CLS]', '[SEP]', '[MASK]', '[EOS']])

# tokens to ids
{'[CLS]': [101], '[SEP]': [102], '[MASK]': [103], '[EOS]': [1031, 1041, 2891, 1033]}
```

- IDs to tokens

```
# ids to tokens
token_id = encoded_input['input_ids'].tolist()[0]
print({tokenizer.convert_ids_to_tokens(id, skip_special_
    ↪tokens=False):id \
        for id in token_id})

### output
{'[CLS]': 101, 'gen': 8991, 'ai': 9932, 'is': 2003, 'awesome': 12476,
    ↪'[SEP]': 102}
```

- Out-of-vocabulary tokens

```
text = "Gen AI is awesome "
encoded_input = tokenizer(text, return_tensors='pt')

print({x : tokenizer.encode(x, add_special_tokens=False) for x in [
    ↪'[CLS]' + text.split() + '[SEP']]}
print(tokenizer.convert_ids_to_tokens(100, skip_special_tokens=False))

### output
{'[CLS]': [101], 'Gen': [8991], 'AI': [9932], 'is': [2003], 'awesome': [
    ↪[12476], '': [100], '[SEP]': [102]}
[UNK]
```

- Subword Tokenization

```
# Subword Tokenization
text = "GenAI is awesome "
print({x : tokenizer.encode(x, add_special_tokens=False) for x in [
    ↪'[CLS]' + text.split() + '[SEP']]}
print(tokenizer.convert_ids_to_tokens(100, skip_special_tokens=False))

# output
{'[CLS]': [101], 'GenAI': [8991, 4886], 'is': [2003], 'awesome': [
    ↪[12476], '': [100], '[SEP]': [102]}
[UNK]
```

2.3 Platform and Packages

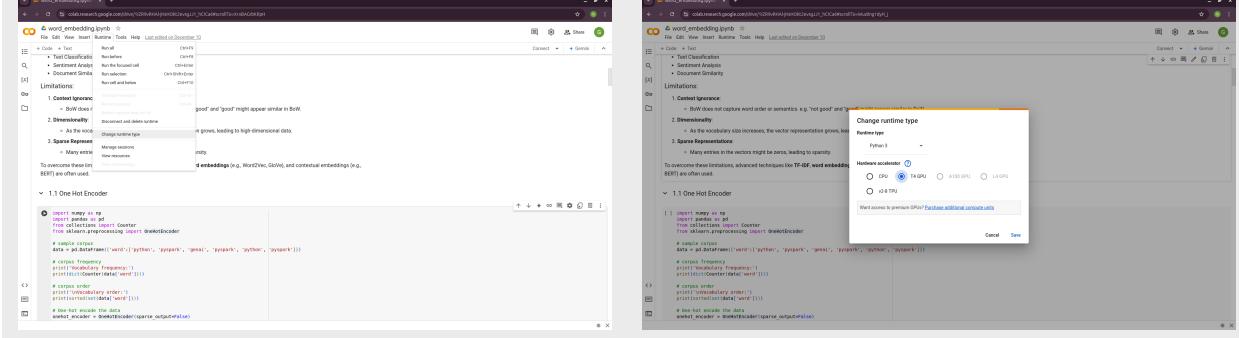
2.3.1 Google Colab

Google Colab (short for Colaboratory) is a free, cloud-based platform that provides users with the ability to write and execute Python code in an interactive notebook environment. It is based on Jupyter notebooks and is powered by Google Cloud services, allowing for seamless integration with Google Drive and other Google services. We will primarily use Google Colab with free T4 GPU runtime throughout this book.

- Key Features

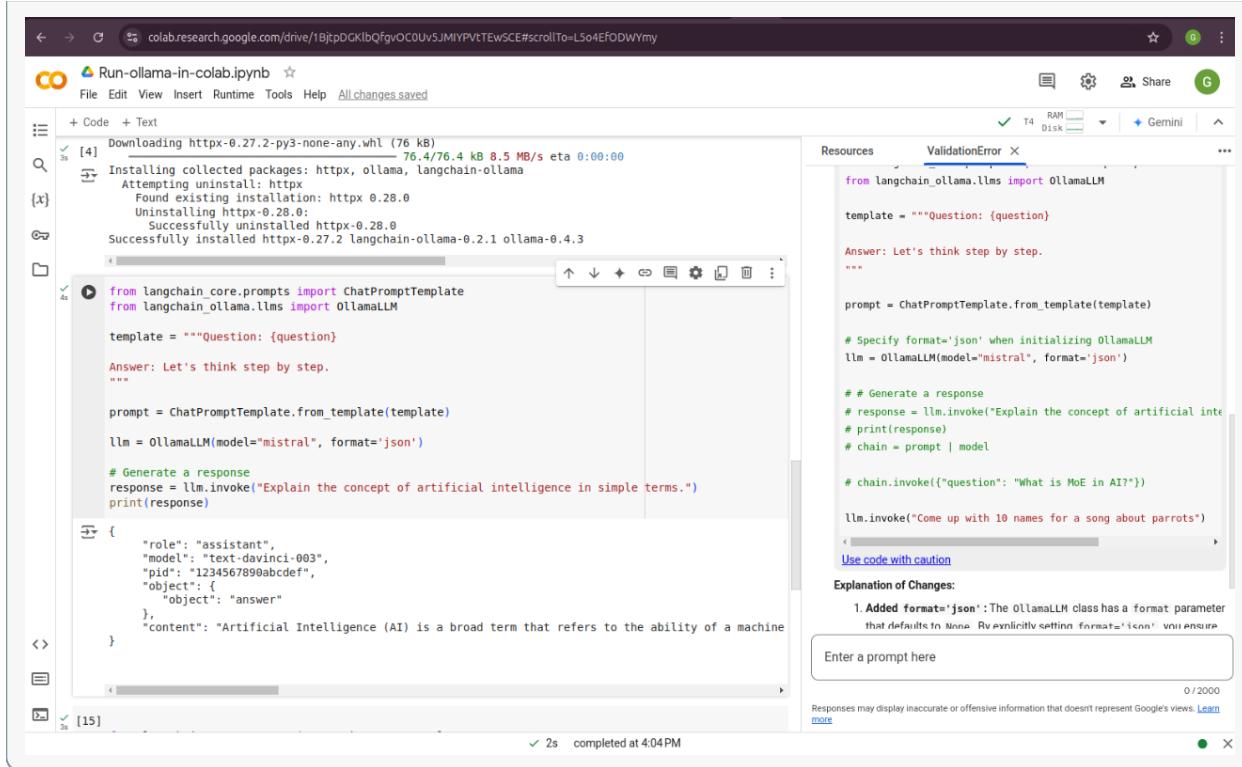
- 1. Free Access to GPUs and TPUs** Colab offers free access to Graphics Processing Units (GPUs) and Tensor Processing Units (TPUs), making it an ideal environment for machine learning, deep learning, and other computationally intensive tasks.
- 2. Integration with Google Drive** You can store and access notebooks directly from your Google Drive, making it easy to collaborate with others and keep your projects organized.
- 3. No Setup Required** Since Colab is entirely cloud-based, you don't need to worry about setting up an environment or managing dependencies. Everything is ready to go out of the box.
- 4. Support for Python Libraries** Colab comes pre-installed with many popular Python libraries, including TensorFlow, PyTorch, Keras, and OpenCV, among others. You can also install any additional libraries using *pip*.
- 5. Collaborative Features** Multiple users can work on the same notebook simultaneously, making it ideal for collaboration. Changes are synchronized in real-time.
- 6. Rich Media Support** Colab supports the inclusion of rich media, such as images, videos, and LaTeX equations, directly within the notebook. This makes it a great tool for data analysis, visualization, and educational purposes.
- 7. Easy Sharing** Notebooks can be easily shared with others via a shareable link, just like Google Docs. Permissions can be set for viewing or editing the document.

- GPU Activation Runtime --> change runtime type --> T4/A100 GPU



Tips

You can use the Gemini API for code troubleshooting in a Colab notebook for free.



The screenshot shows a Google Colab notebook titled "Run-ollama-in-colab.ipynb". The left pane displays Python code for interacting with an Ollama Large Language Model (LLM) using the langchain library. The right pane shows the "Resources" tab, which includes a "Validation Error" section and a "Gemini" section. The Gemini section lists hardware resources: T4 RAM (16GB), Disk (1TB), and Gemini (1 instance). A message indicates that the OllamaLLM class has a 'format' parameter that defaults to 'None'. A note also cautions against using 'code with caution'.

```

from langchain_core.prompts import ChatPromptTemplate
from langchain_ollama.llms import OllamaLLM

template = """Question: {question}

Answer: Let's think step by step.

"""

prompt = ChatPromptTemplate.from_template(template)

llm = OllamaLLM(model="mistral", format='json')

# Generate a response
response = llm.invoke("Explain the concept of artificial intelligence in simple terms.")
print(response)

{
    "role": "assistant",
    "model": "text-davinci-003",
    "pid": "1234567890abcdef",
    "object": {
        "object": "answer"
    },
    "content": "Artificial Intelligence (AI) is a broad term that refers to the ability of a machine
}

```

2.3.2 HuggingFace

Hugging Face is a company and open-source community focused on providing tools and resources for NLP and machine learning. It is best known for its popular **Transformers** library, which allows easy access to pre-trained models for a wide variety of NLP tasks. Moreover, Hugging Face's libraries provide simple Python APIs that make it easy to load models, preprocess data, and run inference. This simplicity allows both beginners and advanced users to leverage cutting-edge NLP models. We will mainly use the embedding models and Large Language Models (LLMs) from **Hugging Face Model Hub** central repository.

2.3.3 Ollama

Ollama is a package designed to run LLMs locally on your personal device or server, rather than relying on external cloud services. It provides a simple interface to download and use AI models tailored for various tasks, ensuring privacy and control over data while still leveraging the power of LLMs.

- Key features of Ollama:
 - Local Execution: Models run entirely on your hardware, making it ideal for users who prioritize data privacy.
 - Pre-trained Models: Offers a curated set of LLMs optimized for local usage.
 - Cross-Platform: Compatible with macOS, Linux, and other operating systems, depending on hardware specifications.
 - Ease of Use: Designed to make setting up and using local AI models simple for non-technical users.

- Efficiency: Focused on lightweight models optimized for local performance without needing extensive computational resources.

To simplify the management of access tokens for various LLMs, we will use Ollama in Google Colab.

- Ollama installation in Google Colab

1. colab-xterm

```
!pip install colab-xterm
%load_ext colabxterm
```

2. download ollama

```
/content# curl https://ollama.ai/install.sh | sh
```

```
%xterm # curl https://ollama.ai/install.sh | sh
Launching Xterm...
/content# curl https://ollama.ai/install.sh | sh
% Total    % Received % Xferd  Average Speed   Time     Time      Current
          Dload  Upload Total Spent   Spent    Left  Speed
100 13269    0 13269    0      0  47685      0 --:--:--:--:--:-- 47730
>>> Installing ollama to /usr/local
>>> Downloading Linux amd64 bundle
#####
>>> Creating ollama user...
>>> Adding ollama user to video group...
>>> Adding current user to ollama group...
>>> Creating ollama systemd service...
WARNING: systemd is not running
WARNING: Unable to detect NVIDIA/AMD GPU. Install lspci or lshw to automatically detect and install GPU dependencies.
>>> The ollama API is now available at 127.0.0.1:11434.
>>> Install complete. Run "ollama" from the command line.
/content#
```

3. launch Ollama serve

```
/content# ollama serve
```

```
%xterm
Launching Xterm...
/content# ollama serve
Couldn't find '/root/.ollama/id_ed25519'. Generating new private key.
Your new public key is:
ssh-ed25519 AAAAC3NzaC1lZDI1NTE5AAAIFSGi7ztCm6LqgqliyANCLCgZYCC5eGwEGsZBv1G/OfmJ
2024/12/14 03:40:54 routes.go:1195: INFO server config env="map[CUDA_VISIBLE_DEVICES: GPU_DEVICE_ORDINAL: HI
P_VISIBLE_DEVICES: HSA_OVERRIDE_GFX_VERSION: HTTPS_PROXY: HTTP_PROXY: NO_PROXY: OLLAMA_DEBUG:false OLLAMA_FL
ASH_ATTENTION:false OLLAMA_GPU_OVERHEAD:0 OLLAMA_HOST:http://127.0.0.1:11434 OLLAMA_INTEL_GPU:false OLLAMA_K
```

4. download models

```
/content# ollama pull mistral #llama3.2 #bge-m3
```

5. check

```
!ollama list
```

(continues on next page)

```
Launching Xterm...
/content# ollama pull mistral #llama3.2 #bge-m3
pulling manifest
pulling ff82381e2bea... 100% 4.1 GB
pulling 43070e2d4e53... 100% 11 KB
pulling 491dfa501e59... 100% 801 B
pulling ed11eda7790d... 100% 30 B
pulling 42347cd80dc8... 100% 485 B
verifying sha256 digest
writing manifest
success
/content# ollama pull llama3.2
pulling manifest
pulling dde5aa3fc5ff... 100% 2.0 GB
pulling 966de95ca8a6... 100% 1.4 KB
pulling fcc5a6bec9da... 100% 7.7 KB
pulling a70ff7e570d9... 100% 6.0 KB
pulling 56bb8bd477a5... 100% 96 B
pulling 34bb5ab01051... 100% 561 B
verifying sha256 digest
writing manifest
success
```

(continued from previous page)

####			
NAME	ID	SIZE	MODIFIED
llama3.2:latest	a80c4f17acd5	2.0 GB	14 seconds ago
mistral:latest	f974a74358d6	4.1 GB	About a minute ago

2.3.4 langchain

LangChain is a powerful framework for building AI applications that combine the capabilities of large language models with external tools, memory, and custom workflows. It enables developers to create intelligent, context-aware, and dynamic applications with ease.

It has widely applied in:

1. **Conversational AI** Create chatbots or virtual assistants that maintain context, integrate with APIs, and provide intelligent responses.
2. **Knowledge Management** Combine LLMs with external knowledge bases or databases to answer complex questions or summarize documents.
3. **Automation** Automate workflows by chaining LLMs with tools for decision-making, data extraction, or content generation.
4. **Creative Applications** Use LangChain for generating stories, crafting marketing copy, or producing artistic content.

We will primarily use LangChain in this book. For instance, to work with downloaded Ollama LLMs, the `langchain_ollama` package is required.

```
# chain of thought prompting
from langchain_ollama.llms import OllamaLLM
from langchain_core.prompts import ChatPromptTemplate
from langchain.output_parsers import CommaSeparatedListOutputParser

template = """Question: {question}
```

(continues on next page)

(continued from previous page)

Answer: Let's think step by step.

11

```
prompt = ChatPromptTemplate.from_template(template)
model = OllamaLLM(temperature=0.0, model='mistral', format='json')
output_parser = CommaSeparatedListOutputParser()

chain = prompt | model | output_parser

response = chain.invoke({"question": "What is Mixture of Experts(MoE) in AI?"})
print(response)
```


WORD AND SENTENCE EMBEDDING

Chinese proverb

Though the forms may vary, the essence remains unchanged. – old Chinese proverb

Word embedding is a method in natural language processing (NLP) to represent words as dense vectors of real numbers, capturing semantic relationships between them. Instead of treating words as discrete symbols (like one-hot encoding), word embeddings map words into a continuous vector space where similar words are located closer together.

Colab Notebook for This Chapter

- Word Embedding:  [Open in Colab](#)

3.1 Traditional word embeddings

Bag of Words (BoW) is a simple and widely used text representation technique in natural language processing (NLP). It represents a text (e.g., a document or a sentence) as a collection of words, ignoring grammar, order, and context but keeping their frequency.

Key Features of Bag of Words:

1. **Vocabulary Creation:** - A list of all unique words in the dataset (the “vocabulary”) is created. - Each word becomes a feature.
2. **Representation:** - Each document is represented as a vector or a frequency count of words from the vocabulary. - If a word from the vocabulary is present in the document, its count is included in the vector. - Words not present in the document are assigned a count of zero.
3. **Simplicity:** - The method is computationally efficient and straightforward. - However, it ignores the sequence and semantic meaning of the words.

Applications:

- Text Classification
- Sentiment Analysis

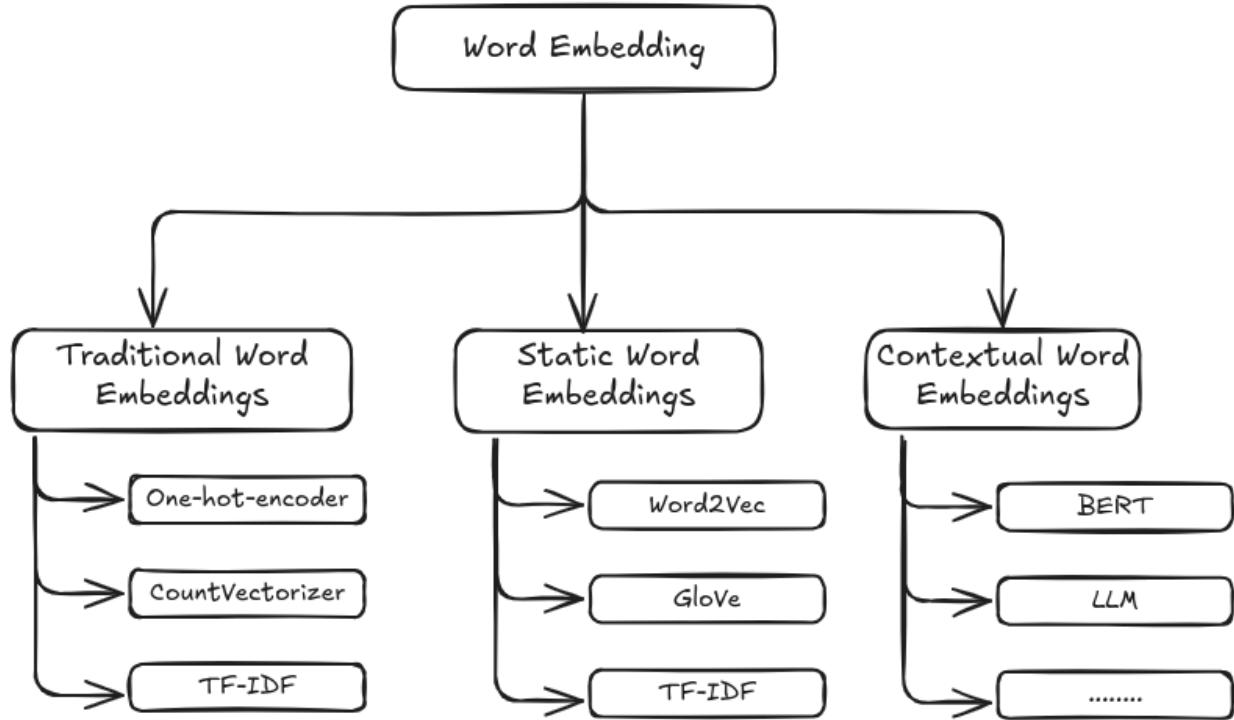


Fig. 1: Embedding Diagram

- Document Similarity

Limitations:

1. **Context Ignorance:** - BoW does not capture word order or semantics. - For example, “not good” and “good” might appear similar in BoW.
2. **Dimensionality:** - As the vocabulary size increases, the vector representation grows, leading to high-dimensional data.
3. **Sparse Representations:** - Many entries in the vectors might be zeros, leading to sparsity.

3.1.1 One Hot Encoder

```

import numpy as np
import pandas as pd
from collections import Counter
from sklearn.preprocessing import OneHotEncoder

# sample corpus
data = pd.DataFrame({'word':['python', 'pyspark', 'genai', 'pyspark', 'python',
                           'pyspark']})

# corpus frequency
print('Vocabulary frequency:')
  
```

(continues on next page)

(continued from previous page)

```

print(dict(Counter(data['word'])))

# corpus order
print('\nVocabulary order:')
print(sorted(set(data['word'])))

# One-hot encode the data
onehot_encoder = OneHotEncoder(sparse_output=False)
onehot_encoded = onehot_encoder.fit_transform(data[['word']])

# the encoded order base on the order of the copus
print('\nEncoded representation:')
print(onehot_encoded)

```

Vocabulary frequency:
`{'python': 2, 'pyspark': 3, 'genai': 1}`

Vocabulary order:
`['genai', 'pyspark', 'python']`

Encoded representation:
`[[0. 0. 1.]`
`[0. 1. 0.]`
`[1. 0. 0.]`
`[0. 1. 0.]`
`[0. 0. 1.]`
`[0. 1. 0.]]`

3.1.2 CountVectorizer

```

from sklearn.feature_extraction.text import CountVectorizer

# sample corpus
corpus = [
    'Gen AI is awesome',
    'Gen AI is fun',
    'Gen AI is hot'
]

# Initialize the CountVectorizer
vectorizer = CountVectorizer()

# Fit and transform
X = vectorizer.fit_transform(corpus)

```

(continues on next page)

(continued from previous page)

```
print('Vocabulary:')
print(vectorizer.get_feature_names_out())

print('\nEmbedded representation:')
print(X.toarray())
```

Vocabulary:

```
['ai' 'awesome' 'fun' 'gen' 'hot' 'is']
```

Embedded representation:

```
[[1 1 0 1 0 1]
[1 0 1 1 0 1]
[1 0 0 1 1 1]]
```

To overcome these limitations, advanced techniques like **TF-IDF**, **word embeddings** (e.g., Word2Vec, GloVe), and contextual embeddings (e.g., BERT) are often used.

3.1.3 TF-IDF

TF-IDF (Term Frequency-Inverse Document Frequency) is a statistical measure used in text analysis to evaluate the importance of a word in a document relative to a collection (or corpus) of documents. It builds upon the **Bag of Words (BoW)** model by not only considering the frequency of a word in a document but also taking into account how common or rare the word is across the corpus. The pyspark implementation can be found at [\[PySpark\]](#).

- Components of TF-IDF
- **t**: the term in corpus.
- **d**: the document.
- **D**: the corpus.
- **|D|**: the length of the corpus or total number of documents.
 - **Document Frequency (DF)**:
 - $DF(t, D)$: the number of documents that contains term t .
 - **Term Frequency (TF)**:
 - * Measures how frequently a term appears in a document. The higher the frequency, the more important the term is assumed to be to that document.
 - * Formula:

$$TF(t, d) = \frac{\text{Number of occurrences of term } t \text{ in document } d}{\text{Total number of terms in document } d}$$
 - **Inverse Document Frequency (IDF)**:

- * Measures how important a term is by reducing the weight of common terms (like “the” or “and”) that appear in many documents.

- * Formula:

$$IDF(t, D) = \log \left(\frac{|D| + 1}{DF(t, D) + 1} \right) + 1$$

- * Adding 1 to the denominator avoids division by zero when a term is present in all documents.
- * Note that the IDF formula above differs from the standard textbook notation that defines the IDF

Note

The IDF formula above differs from the standard textbook notation that defines the IDF as

$$IDF(t) = \log[|D|/(DF(t, D) + 1)].$$

- TF-IDF Score:

- * The final score is the product of TF and IDF.

- * Formula:

$$TF-IDF(t, d, D) = TF(t, d) \cdot IDF(t, D)$$

```
import pandas as pd
import numpy as np
from collections import Counter
from sklearn.feature_extraction.text import TfidfVectorizer

# sample corpus
corpus = [
    'Gen AI is awesome',
    'Gen AI is fun',
    'Gen AI is hot'
]

# Initialize the TfidfVectorizer
vectorizer = TfidfVectorizer() # norm default norm='l2'

# Fit and transform
X = vectorizer.fit_transform(corpus)

print('Vocabulary:')
print(vectorizer.get_feature_names_out())
```

(continues on next page)

(continued from previous page)

```
# [item for row in matrix for item in row]
corpus_flatted = [item for sub_list in [s.split(' ') for s in
                                          corpus]
                           for item in sub_list]

print('\nVocabulary frequency:')
print(dict(Counter(corpus_flatted)))

print('\nEmbedded representation:')
print(X.toarray())
```

Vocabulary:

```
['ai' 'awesome' 'fun' 'gen' 'hot' 'is']
```

Vocabulary frequency:

```
{'Gen': 3, 'AI': 3, 'is': 3, 'awesome': 1, 'fun': 1, 'hot': 1}
```

Embedded representation:

```
[[0.41285857 0.69903033 0.          0.41285857 0.          0.41285857]
 [0.41285857 0.          0.69903033 0.41285857 0.          0.41285857]
 [0.41285857 0.          0.          0.41285857 0.69903033 0.41285857]]
```

The above results can be validated by the following steps (IDF in document 1):

```
# Step 1: Vocabulary `['ai' 'awesome' 'fun' 'gen' 'hot' 'is']`

tf_idf = pd.DataFrame({'term':vectorizer.get_feature_names_out()})\
    .set_index('term')

# Step 2: |D|
tf_idf['|D|'] = [len(corpus)]*len(vectorizer.get_feature_names_
                  .out())

# Step 3: Compute TF for doc 1: Gen AI is awesome
# - TF for "ai" in Document 1 = 1 (appears once doc 1)
# - TF for "awesome" in Document 1 = 1 (appears once in doc 1)
# - TF for "fun" in Document 1 = 0 (does not appear in doc 1)
# - TF for "gen" in Document 1 = 1 (appear oncein doc 1 )
# - TF for "hot" in Document 1 = 0 (does not appear doc 1 )
# - TF for "is" in Document 1 = 1 (appear once in doc 1 )

tf_idf['TF'] = [1, 1, 0, 1, 0, 1]

# Step 4: Compute DF for doc 1
# - DF For "ai": Appears in all 3 documents.
```

(continues on next page)

(continued from previous page)

```

# - DF For "awesome": Appears in 1 document.
# - DF For "fun": Appears in 1 document.
# - DF For "Gen": Appears in all 3 documents.
# - DF For "hot": Appears in 1 document.
# - DF For "is": Appears in all 3 documents.

tf_idf['DF'] = [3, 1, 1, 3, 1, 3]

# Step 5: Compute IDF
tf_idf['IDF'] = np.log((tf_idf['|D|']+1)/(tf_idf['DF']+1))+1

# Step 6: Compute TF-IDF
tf_idf['TF-IDF'] = tf_idf['TF']*tf_idf['IDF']

# Step 7: l2 normalization
tf_idf['TF-IDF(12)'] = tf_idf['TF-IDF']/np.linalg.norm(tf_idf['TF-IDF'])

print(tf_idf)

```

	D	TF	DF	IDF	TF-IDF	TF-IDF(12)
term						
ai	3	1	3	1.000000	1.000000	0.412859
awesome	3	1	1	1.693147	1.693147	0.699030
fun	3	0	1	1.693147	0.000000	0.000000
gen	3	1	3	1.000000	1.000000	0.412859
hot	3	0	1	1.693147	0.000000	0.000000
is	3	1	3	1.000000	1.000000	0.412859

Fun Fact

TfidfVectorizer is equivalent to CountVectorizer followed by TfidfTransformer.

```

import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.pipeline import Pipeline

# sample corpus
corpus = [
    'Gen AI is awesome',
    'Gen AI is fun',
    'Gen AI is hot'
]

```

```
# pipeline
pipe = Pipeline([('count', CountVectorizer(lowercase=True)),
                 ('tfid', TfidfTransformer())]).fit(corpus)
print(pipe)

# TF
print(pipe['count'].transform(corpus).toarray())

# IDF
print(pipe['tfid'].idf_)

Pipeline(steps=[('count', CountVectorizer()), ('tfid',
                                                TfidfTransformer())])
[[1 1 0 1 0 1]
 [1 0 1 1 0 1]
 [1 0 0 1 1 1]]
[1.          1.69314718 1.69314718 1.          1.69314718 1.
 ]
```

- Applications of TF-IDF
 1. **Information Retrieval:** Ranking documents based on relevance to a query.
 2. **Text Classification:** Feature extraction for machine learning models.
 3. **Document Similarity:** Comparing documents by their weighted term vectors.
- Advantages
 - Highlights important terms while reducing the weight of common terms.
 - Simple to implement and effective for many tasks.
- Limitations
 - Does not capture semantic relationships or word order.
 - Less effective for very large corpora or when working with very short documents.
 - Sparse representation due to high-dimensional feature vectors.

For more advanced representations, embeddings like **Word2Vec** or **BERT** are often used.

3.2 Static word embeddings

Static word embeddings are word representations that assign a fixed vector to each word, regardless of its context in a sentence or paragraph. These embeddings are pre-trained on large corpora and remain unchanged during usage, making them “static.” These embeddings are usually pre-trained on large text corpora using algorithms like Word2Vec, GloVe, or FastText.

3.2.1 Word2Vec

- The Context Window
- CBOW and Skip-Gram Model

```

import gensim
from gensim.models import Word2Vec
from nltk.tokenize import sent_tokenize, word_tokenize

# sample corpus
corpus = [
'Gen AI is awesome',
'Gen AI is fun',
'Gen AI is hot'
]

def tokenize_gensim(corpus):

    tokens = []
    # iterate through each sentence in the corpus
    for s in corpus:

        # tokenize the sentence into words
        temp = gensim.utils.tokenize(s, lowercase=True, deacc=False, \
                                      errors='strict', to_lower=False, \
                                      lower=False)

        tokens.append(list(temp))

    return tokens

tokens = tokenize_gensim(corpus)

# Create Word2Vec model
# sg ({0, 1}, optional) - Training algorithm: 1 for skip-gram; otherwise CBOW.
# CBOW
model1 = gensim.models.Word2Vec(tokens, sg=0, min_count=1,
                                 vector_size=10, window=5)

# Vocabulary
print(model1.wv.key_to_index)

print(model1.wv.get_normed_vectors())

# Print results

```

(continues on next page)

(continued from previous page)

```

print("Cosine similarity between 'gen' " +
      "and 'ai' - Word2Vec(CBOW) : ",
      model1.wv.similarity('gen', 'ai'))


# Create Word2Vec model
# sg ({0, 1}, optional) - Training algorithm: 1 for skip-gram; otherwise CBOW.
# skip-gram
model2 = gensim.models.Word2Vec(tokens, sg=1, min_count=1,
                                 vector_size=10, window=5)

# Vocabulary
print(model2.wv.key_to_index)

print(model2.wv.get_normed_vectors())

# Print results
print("Cosine similarity between 'gen' " +
      "and 'ai' - Word2Vec(skip-gram) : ",
      model2.wv.similarity('gen', 'ai'))

```

```

{'is': 0, 'ai': 1, 'gen': 2, 'hot': 3, 'fun': 4, 'awesome': 5}
[[-0.02660277  0.0117296   0.25318226  0.44695902 -0.4615286  -0.35307196
  0.3204311   0.4451589   -0.24882038 -0.18670462]
 [ 0.41619968 -0.08647515 -0.2558276   0.3695945  -0.274073   -0.10240843
  0.1622154   0.05593351 -0.46721786 -0.5328355 ]
 [ 0.43418837  0.30108306  0.40128633  0.0453006   0.37712952 -0.20221795
 -0.05619935  0.34255028 -0.44665098 -0.2337343 ]
 [-0.41098067 -0.05088534  0.5218584   -0.40045303 -0.12768732 -0.10601949
  0.44194022 -0.32449666  0.00247097 -0.2600907 ]
 [-0.44081825  0.22984274 -0.40207896 -0.20159177 -0.00161115 -0.0135952
 -0.3516631   0.44133204  0.2286844   0.423816 ]
 [-0.42753762  0.23561442 -0.21681462  0.04321203  0.44539306 -0.23385239
  0.23675178 -0.35568893 -0.18596812  0.49255413]]
Cosine similarity between 'gen' and 'ai' - Word2Vec(CBOW) :  0.32937223
{'is': 0, 'ai': 1, 'gen': 2, 'hot': 3, 'fun': 4, 'awesome': 5}
[[-0.02660277  0.0117296   0.25318226  0.44695902 -0.4615286  -0.35307196
  0.3204311   0.4451589   -0.24882038 -0.18670462]
 [ 0.41619968 -0.08647515 -0.2558276   0.3695945  -0.274073   -0.10240843
  0.1622154   0.05593351 -0.46721786 -0.5328355 ]
 [ 0.43418837  0.30108306  0.40128633  0.0453006   0.37712952 -0.20221795
 -0.05619935  0.34255028 -0.44665098 -0.2337343 ]
 [-0.41098067 -0.05088534  0.5218584   -0.40045303 -0.12768732 -0.10601949
  0.44194022 -0.32449666  0.00247097 -0.2600907 ]
 [-0.44081825  0.22984274 -0.40207896 -0.20159177 -0.00161115 -0.0135952
 -0.3516631   0.44133204  0.2286844   0.423816 ]
 [-0.42753762  0.23561442 -0.21681462  0.04321203  0.44539306 -0.23385239
  0.23675178 -0.35568893 -0.18596812  0.49255413]]

```

(continues on next page)

(continued from previous page)

```
-0.3516631  0.44133204  0.2286844  0.423816  ]
[-0.42753762  0.23561442 -0.21681462  0.04321203  0.44539306 -0.23385239
 0.23675178 -0.35568893 -0.18596812  0.49255413]]
Cosine similarity between 'gen' and 'ai' - Word2Vec(skip-gram) : 0.32937223
```

3.2.2 GloVe

```
import gensim.downloader as api
# Download pre-trained GloVe model
glove_vectors = api.load("glove-wiki-gigaword-50")
# Get word vectors (embeddings)
word1 = "king"
word2 = "queen"
vector1 = glove_vectors[word1]
vector2 = glove_vectors[word2]
# Compute cosine similarity between the two word vectors
similarity = glove_vectors.similarity(word1, word2)
print(f"Word vectors for '{word1}': {vector1}")
print(f"Word vectors for '{word2}': {vector2}")
print(f"Cosine similarity between '{word1}' and '{word2}': {similarity}")
```

```
[=====] 100.0% 66.0/66.0MB ↵
↳ downloaded
Word vectors for 'king': [ 0.50451   0.68607  -0.59517  -0.022801  0.60046  -0.
 ↵13498  -0.08813
0.47377  -0.61798  -0.31012  -0.076666  1.493    -0.034189 -0.98173
0.68229  0.81722  -0.51874  -0.31503  -0.55809  0.66421   0.1961
-0.13495  -0.11476  -0.30344  0.41177  -2.223    -1.0756  -1.0783
-0.34354  0.33505   1.9927   -0.04234  -0.64319  0.71125   0.49159
0.16754  0.34344  -0.25663  -0.8523   0.1661   0.40102   1.1685
-1.0137  -0.21585  -0.15155   0.78321  -0.91241  -1.6106  -0.64426
-0.51042 ]
Word vectors for 'queen': [ 0.37854    1.8233    -1.2648    -0.1043    0.35829 ↵
 ↵ 0.60029
-0.17538   0.83767  -0.056798 -0.75795   0.22681   0.98587
0.60587  -0.31419   0.28877   0.56013  -0.77456   0.071421
-0.5741   0.21342   0.57674   0.3868  -0.12574   0.28012
0.28135  -1.8053   -1.0421   -0.19255 -0.55375  -0.054526
1.5574   0.39296  -0.2475   0.34251   0.45365   0.16237
0.52464  -0.070272 -0.83744  -1.0326   0.45946   0.25302
-0.17837  -0.73398  -0.20025   0.2347  -0.56095  -2.2839
0.0092753 -0.60284 ]
Cosine similarity between 'king' and 'queen': 0.7839043140411377
```

3.2.3 Fast Text

Fast Text incorporates subword information (useful for handling rare or unseen words)

```
from gensim.models import FastText

import gensim
from gensim.models import Word2Vec

# sample corpus
corpus = [
'Gen AI is awesome',
'Gen AI is fun',
'Gen AI is hot'
]

def tokenize_gensim(corpus):

    tokens = []
    # iterate through each sentence in the corpus
    for s in corpus:

        # tokenize the sentence into words
        temp = gensim.utils.tokenize(s, lowercase=True, deacc=False, \
                                      errors='strict', to_lower=False, \
                                      lower=False)

        tokens.append(list(temp))

    return tokens

tokens = tokenize_gensim(corpus)

# create FastText model
model = FastText(tokens, vector_size=10, window=5, min_count=1, workers=4)
# Train the model
model.train(tokens, total_examples=len(tokens), epochs=10)

# Vocabulary
print(model.wv.key_to_index)

print(model.wv.get_normed_vectors())

# Print results
print("Cosine similarity between 'gen' " +
      "and 'ai' - Word2Vec : ",
      model.wv.similarity('gen', 'ai'))
```

```

WARNING:gensim.models.word2vec:Effective 'alpha' higher than previous training
  ↵cycles
{'is': 0, 'ai': 1, 'gen': 2, 'hot': 3, 'fun': 4, 'awesome': 5}
[[-0.01875759  0.086543  -0.25080433  0.2824868  -0.23755953 -0.11316587
  0.473383   0.39204055  -0.30422893 -0.5566626 ]
 [ 0.5088161  -0.3323528  -0.128698  -0.11877266 -0.38699347  0.20977001
  0.05947014 -0.05622245  -0.36257952 -0.5177341 ]
 [ 0.18038039  0.51484865  0.40694886  0.05965518 -0.05985437 -0.10832689
  0.37992737  0.5992712   0.01503773  0.1192203 ]
 [-0.5694013   0.23560704  0.0265804   -0.41392225 -0.00285366 -0.3076269
  0.2076883  -0.425648   0.29903153  0.19965051]
 [-0.23892775  0.10744874 -0.03730153  -0.23521401  0.32083488  0.21598674
 -0.29570717 -0.03044808  0.75250715  0.26538488]
 [-0.31881964 -0.06544963 -0.44274488  0.15485793  0.39120612 -0.05415314
  0.15772066 -0.05987714 -0.6986104   0.03967094]]
Cosine similarity between 'gen' and 'ai' - Word2Vec : -0.21662527

```

3.3 Contextual word embeddings

Contextual word embeddings are word representations where the embedding of a word changes depending on its context in a sentence or document. These embeddings capture the meaning of a word as influenced by its surrounding words, addressing the limitations of static embeddings by incorporating contextual nuances.

3.3.1 BERT

```

from transformers import BertTokenizer, BertModel
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
model = BertModel.from_pretrained("bert-base-uncased")

text = "Gen AI is awesome"
encoded_input = tokenizer(text, return_tensors='pt')
embeddings = model(**encoded_input).last_hidden_state

print(encoded_input)
print({x : tokenizer.encode(x, add_special_tokens=False) for x in ['[CLS]']+
  ↵text.split()+'[SEP]', '[EOS']})

print(embeddings.shape)
print(embeddings)

```

```

t{'input_ids': tensor([[ 101,  8991,  9932,  2003, 12476,   102]]), 'token_type_
  ↵ids': tensor([[0, 0, 0, 0, 0, 0]]), 'attention_mask': tensor([[1, 1, 1, 1, 1,
  ↵1]])}
{'[CLS]': [101], 'Gen': [8991], 'AI': [9932], 'is': [2003], 'awesome': [12476],

```

(continues on next page)

(continued from previous page)

```

→ '[SEP]': [102], '[EOS]': [1031, 1041, 2891, 1033]}

torch.Size([1, 6, 768])
tensor([[[ -0.1129, -0.1477, -0.0056, ..., -0.1335, 0.2605, 0.2113],
         [-0.6841, -1.1196, 0.3349, ..., -0.5958, 0.1657, 0.6988],
         [-0.5385, -0.2649, 0.2639, ..., -0.1544, 0.2532, -0.1363],
         [-0.1794, -0.6086, 0.1292, ..., -0.1620, 0.1721, 0.4356],
         [-0.0187, -0.7320, -0.3420, ..., 0.4028, 0.1425, -0.2014],
         [ 0.5493, -0.1029, -0.1571, ..., 0.3503, -0.7601, -0.1398]]],
grad_fn=<NativeLayerNormBackward0>)

```

3.3.2 gte-large-en-v1.5

The `gte-large-en-v1.5` is a state-of-the-art text embedding model developed by Alibaba's Institute for Intelligent Computing. It's designed for natural language processing tasks and excels in generating dense vector representations (embeddings) of text for applications such as text retrieval, classification, clustering, and reranking.

It can handle up to 8192 tokens, making it suitable for long-context tasks. More details can be found at: <https://huggingface.co/Alibaba-NLP/gte-large-en-v1.5>.

```

# Requires transformers>=4.36.0

import torch.nn.functional as F
from transformers import AutoModel, AutoTokenizer

input_texts = [
    'Gen AI is awesome',
    'Gen AI is fun',
    'Gen AI is hot'
]

model_path = 'Alibaba-NLP/gte-large-en-v1.5'
tokenizer = AutoTokenizer.from_pretrained(model_path)
model = AutoModel.from_pretrained(model_path, trust_remote_code=True)

# Tokenize the input texts
batch_dict = tokenizer(input_texts, max_length=8192, padding=True, \
                       truncation=True, return_tensors='pt')

print(batch_dict)

outputs = model(**batch_dict)
embeddings = outputs.last_hidden_state[:, 0]

```

(continues on next page)

(continued from previous page)

```
# (Optional) normalize embeddings
embeddings = F.normalize(embeddings, p=2, dim=1)
scores = (embeddings[:1] @ embeddings[1:].T) * 100
print(embeddings)
print(scores.tolist())
```

```
{'input_ids': tensor([[ 101,  8991,  9932,  2003, 12476,   102],
       [ 101,  8991,  9932,  2003,  4569,   102],
       [ 101,  8991,  9932,  2003,  2980,   102]]), 'token_type_ids': tensor([[0,
→ 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0]]), 'attention_mask': tensor([[1, 1, 1, 1, 1, 1],
       [1, 1, 1, 1, 1, 1],
       [1, 1, 1, 1, 1, 1]])}

tensor([[ 0.0079,  0.0008, -0.0001, ...,  0.0418, -0.0138, -0.0236],
       [ 0.0079,  0.0218, -0.0171, ...,  0.0412, -0.0230, -0.0237],
       [ 0.0073, -0.0106, -0.0194, ...,  0.0711, -0.0204, -0.0036]],
       grad_fn=<DivBackward0>)
[[92.85284423828125, 92.81655883789062]]
```

3.3.3 bge-base-en-v1.5

The bge-base-en-v1.5 model is a general-purpose text embedding model developed by the Beijing Academy of Artificial Intelligence (BAAI). It transforms input text into 768-dimensional vector embeddings, making it useful for tasks like semantic search, text similarity, and clustering. This model is fine-tuned using contrastive learning, which helps improve its ability to distinguish between similar and dissimilar sentences effectively. More details can be found at: <https://huggingface.co/BAAI/bge-base-en-v1.5>.

```
from transformers import AutoTokenizer, AutoModel
import torch

# Sentences we want sentence embeddings for
sentences = [
    'Gen AI is awesome',
    'Gen AI is fun',
    'Gen AI is hot'
]
# Load model from HuggingFace Hub
tokenizer = AutoTokenizer.from_pretrained('BAAI/bge-large-zh-v1.5')
model = AutoModel.from_pretrained('BAAI/bge-large-zh-v1.5')
model.eval()

# Tokenize sentences
```

(continues on next page)

(continued from previous page)

```

encoded_input = tokenizer(sentences, padding=True, truncation=True, return_
↪tokens='pt')
print(encoded_input)

# Compute token embeddings
with torch.no_grad():
    model_output = model(**encoded_input)
    # Perform pooling. In this case, cls pooling.
    sentence_embeddings = model_output[0][:, 0]
# normalize embeddings
sentence_embeddings = torch.nn.functional.normalize(sentence_embeddings, p=2,_
↪dim=1)
print("Sentence embeddings:", sentence_embeddings)

```

```

{'input_ids': tensor([[ 101, 10234, 8171, 8578, 8310, 143, 11722, 9974,_
↪8505, 102],
      [ 101, 10234, 8171, 8578, 8310, 9575, 102, 0, 0, 0],
      [ 101, 10234, 8171, 8578, 8310, 9286, 102, 0, 0, 0]]),
↪'token_type_ids': tensor([[0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
      [0, 0, 0, 0, 0, 0, 0, 0, 0, 0], 'attention_mask': tensor([[1, 1, 1, 1, 1,
↪1, 1, 1, 1, 1],
      [1, 1, 1, 1, 1, 1, 0, 0, 0],
      [1, 1, 1, 1, 1, 1, 0, 0, 0]])}

Sentence embeddings: tensor([[ 0.0700,  0.0119,  0.0049, ...,  0.0428, -0.0475,_
↪ 0.0242],
      [ 0.0800, -0.0065, -0.0519, ...,  0.0057, -0.0770,  0.0119],
      [ 0.0740, -0.0185, -0.0369, ...,  0.0083, -0.0026,  0.0016]])}

```

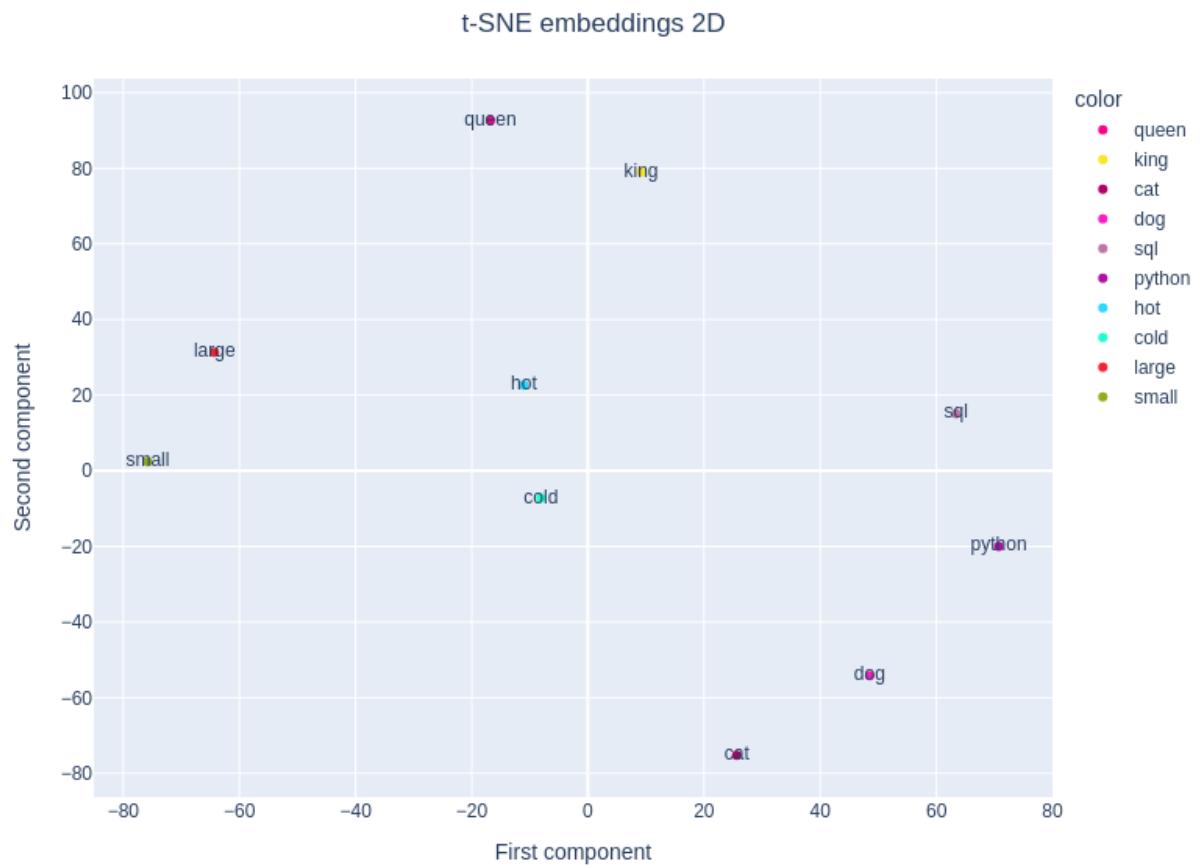


Fig. 2: t-SNE embeddings 2D

CHAPTER FOUR

PROMPT ENGINEERING

Proverb

Our master knows how to guide people skillfully and methodically. He broadens my mind with culture and restrains me with ritual. – Zi Han

Colab Notebook for This Chapter

- Prompt Engineering with Local LLM:  [Open in Colab](#)
- Prompt Engineering with OpenAI API:  [Open in Colab](#)

4.1 Prompt

A prompt is the input or query given to an LLM to elicit a specific response. It acts as the user's way of “programming” the model without code, simply by phrasing questions or tasks appropriately.

4.2 Prompt Engineering

4.2.1 What's Prompt Engineering

Prompt engineering is the practice of designing and refining input prompts to guide LLMs to produce desired outputs effectively and consistently. It involves crafting queries, commands, or instructions that align with the model's capabilities and the task's requirements.

4.2.2 Key Elements of a Prompt

- Clarity: A clear and unambiguous prompt ensures the model understands the task.
- Specificity: Including details like tone, format, length, or audience helps tailor the response.
- Context: Providing background information ensures the model generates relevant outputs.

4.3 Advanced Prompt Engineering

4.3.1 Role Assignment

- Assign a specific role or persona to the AI to shape its style and expertise.

- **Example:**

"You are a professional data scientist. Explain how to build a machine learning model to a beginner."

```
# You are a professional data scientist. Explain how to build a machine
# learning model to a beginner.
template = """Role: you are a {role}
task: {task}
Answer:
"""

prompt = ChatPromptTemplate.from_template(template)
model = OllamaLLM(temperature=0.0, model=MODEL, format='json')
output_parser = CommaSeparatedListOutputParser()

chain = prompt | model | output_parser

response = chain.invoke({'role': 'data scientist', \
                        "task": "Explain how to build a machine \
                        learning model to a beginner"})
print(response)
```

```
{
"role": "assistant",
"content": "To build a machine learning model, let's follow these steps as a
→beginner: \n\n1. **Define the Problem**: Understand what problem you are
→trying to solve. This could be anything from predicting house prices,
→recognizing images, or even recommending products. \n\n2. **Collect and
→Prepare Data**: Gather relevant data for your problem. This might involve web
→scraping, APIs, or using existing datasets. Once you have the data, clean it
→by handling missing values, outliers, and errors. \n\n3. **Explore and
→Visualize Data**: Understand the structure of your data, its distribution, and
→relationships between variables. This can help in identifying patterns and
→making informed decisions about the next steps. \n\n4. **Feature
→Engineering**: Create new features that might be useful for the model to make
→accurate predictions. This could involve creating interactions between
→existing features or using techniques like one-hot encoding. \n\n5. **Split
→Data**: Split your data into training, validation, and testing sets. The
→training set is used to train the model, the validation set is used to tune
→hyperparameters, and the testing set is used to evaluate the final performance
→of the model. \n\n6. **Choose a Model**: Select a machine learning algorithm,
```

(continues on next page)

(continued from previous page)

```

→that suits your problem. Some common algorithms include linear regression for
→regression problems, logistic regression for binary classification problems,
→decision trees, random forests, support vector machines (SVM), and neural
→networks for more complex tasks. \n\n7. **Train the Model**: Use your training
→data to train the chosen model. This involves feeding the data into the model
→and adjusting its parameters based on the error it makes. \n\n8. **Tune
→Hyperparameters**: Adjust the hyperparameters of the model to improve its
→performance. This could involve changing learning rates, number of layers in a
→neural network, or the complexity of a decision tree. \n\n9. **Evaluate the
→Model**: Use your testing data to evaluate the performance of the model.
→Common metrics include accuracy for classification problems, mean squared
→error for regression problems, and precision, recall, and F1 score for
→imbalanced datasets. \n\n10. **Deploy the Model**: Once you are satisfied with
→the performance of your model, deploy it to a production environment where it
→can make predictions on new data."
}
```

4.3.2 Contextual Setup

- Provide sufficient background or context for the AI to understand the task.
- **Example:**

"I am planing to write a book about GenAI best practice, help me draft the contents for the book."

```

# Contextual Setup

# I am planing to write a book about GenAI best practice, help me draft the
# contents for the book.
template = """Role: you are a {role}
task: {task}
Answer:
"""

prompt = ChatPromptTemplate.from_template(template)
model = OllamaLLM(temperature=0.0, model=MODEL, format='json')
output_parser = CommaSeparatedListOutputParser()

chain = prompt | model | output_parser

response = chain.invoke({'role': 'book writer', \
                        "task": "I am planing to write a book about \
                        GenAI best practice, help me draft the \
                        contents for the book."})
print(response)

```

```
{
  "1. Introduction": "Introduction to General Artificial Intelligence (GenAI) and its significance in today's world.",
  "2. Chapter 1 - Understanding AI": "Exploring the basics of Artificial Intelligence, its history, and evolution.",
  "3. Chapter 2 - Types of AI": "Detailed discussion on various types of AI such as Narrow AI, General AI, and Superintelligent AI.",
  "4. Chapter 3 - GenAI Architecture": "Exploring the architecture of General AI systems, including neural networks, deep learning, and reinforcement learning.",
  "5. Chapter 4 - Ethics in AI Development": "Discussing the ethical considerations involved in developing GenAI, such as privacy, bias, and accountability.",
  "6. Chapter 5 - Data Collection and Management": "Understanding the importance of data in AI development, best practices for data collection, and responsible data management.",
  "7. Chapter 6 - Model Training and Optimization": "Exploring techniques for training AI models effectively, including hyperparameter tuning, regularization, and optimization strategies.",
  "8. Chapter 7 - Testing and Validation": "Discussing the importance of testing and validation in ensuring the reliability and accuracy of GenAI systems.",
  "9. Chapter 8 - Deployment and Maintenance": "Exploring best practices for deploying AI models into production environments, as well as ongoing maintenance and updates.",
  "10. Case Studies": "Real-world examples of successful GenAI implementations across various industries, highlighting key takeaways and lessons learned.",
  "11. Future Trends in GenAI": "Exploring emerging trends in the field of General AI, such as quantum computing, explainable AI, and human-AI collaboration.",
  "12. Conclusion": "Summarizing the key points discussed in the book and looking forward to the future of General AI."
}
```

4.3.3 Explicit Instructions

- Clearly specify the format, tone, style, or structure you want in the response.
- **Example:**

"Explain the concept of word embeddings in 100 words, using simple language suitable for a high school student."

```
# Explicit Instructions
from langchain_ollama.llms import OllamaLLM
from langchain_core.prompts import ChatPromptTemplate
from langchain.output_parsers import CommaSeparatedListOutputParser

# Explain the concept of word embeddings in 100 words, using simple
# language suitable for a high school student
```

(continues on next page)

(continued from previous page)

```

template = """you are a {role}
task: {task}
instruction: {instruction}
Answer: Let's think step by step.
"""

prompt = ChatPromptTemplate.from_template(template)
model = OllamaLLM(temperature=0.0, model=MODEL, format='json')
output_parser = CommaSeparatedListOutputParser()

chain = prompt | model

response = chain.invoke({'role': 'AI engineer', \
                        'task': "Explain the concept of word embeddings in \
                                100 words", \
                        'instruction': "using simple \
                                language suitable for a high school student"})

print(response)

```

```
{
"assistant": {
    "message": "Word Embeddings are like giving words a special address in a big library. Each word gets its own unique location, and words that are used in similar ways get placed close together. This helps the computer understand the meaning of words better when it's reading text. For example, 'king' might be near 'queen', because they are both types of royalty. And 'apple' might be near 'fruit', because they are related concepts."
}
}
```

4.3.4 Chain of Thought (CoT) Prompting

- Encourage step-by-step reasoning for complex problems.
- **Example:**

“Solve this math problem step by step: A train travels 60 miles in 1.5 hours. What is its average speed?”

```

# CoT
from langchain_ollama.llms import OllamaLLM
from langchain_core.prompts import ChatPromptTemplate
from langchain.output_parsers import CommaSeparatedListOutputParser

```

(continues on next page)

(continued from previous page)

```
# Solve this math problem step by step: A train travels 60 miles in 1.5 hours.
# What is its average speed?

template = """you are a {role}
task: {task}
question: {question}
Answer: Let's think step by step.
"""

prompt = ChatPromptTemplate.from_template(template)
model = OllamaLLM(temperature=0.0, model=MODEL, format='json')
output_parser = CommaSeparatedListOutputParser()

chain = prompt | model

response = chain.invoke({'role': 'math student', \
                        'task': "Solve this math problem step by step: \
                                  A train travels 60 miles in 1.5 hours.", \
                        'question': "What is its average speed per minute?"})

print(response)
```

```
{
  "Solution": {
    "Step 1": "First, let's find the average speed of the train per hour.",
    "Step 2": "The train travels 60 miles in 1.5 hours. So, its speed per hour is\u2192 60 miles / 1.5 hours = 40 miles/hour.",
    "Step 3": "Now, let's find the average speed of the train per minute. Since\u2192 there are 60 minutes in an hour, the speed per minute would be the speed per\u2192 hour multiplied by the number of minutes in an hour divided by 60.",
    "Step 4": "So, the average speed of the train per minute is (40 miles/hour * \u2192 (1 hour / 60)) = (40/60) miles/minute = 2/3 miles/minute."
  }
}
```

4.3.5 Few-Shot Prompting

- Provide examples to guide the AI on how to respond.
- **Example:**
 - “Here are examples of loan application decision: ‘example’: {‘input’: {‘fico’:800, ‘income’:100000,’loan_amount’: 10000} ‘decision’: “accept” Now Help me to make a decision to accpet or reject the loan application and give the reason. ‘input’: “{‘fico’:820, ‘income’:100000, ‘loan_amount’: 1,000}” “*

```

# Few-Shot Prompting
from langchain_ollama.llms import OllamaLLM
from langchain_core.prompts import ChatPromptTemplate
from langchain.output_parsers import CommaSeparatedListOutputParser

# Here are examples of loan application decision:
# 'example': {'input': {'fico':800, 'income':100000,'loan_amount': 10000}}
# 'decision': "accept"
# Now Help me to make a decision to accpet or reject the loan application and
# give the reason.
# 'input': {"fico':820, 'income':100000, 'loan_amount': 1,000}"

template = """you are a {role}
task: {task}
examples: {example}
input: {input}
decision:
"""

prompt = ChatPromptTemplate.from_template(template)
model = OllamaLLM(temperature=0.0, model=MODEL, format='json')
output_parser = CommaSeparatedListOutputParser()

chain = prompt | model

response = chain.invoke({'role': 'banker', \
                        'task': "Help me to make a decision to accpet or \\\nreject the loan application ",\
                        'example': {'input': {'fico':800, 'income':100000,\n                           'loan_amount': 10000},\n                           'decision': "accept"}, \
                        'input': {'fico':820, 'income':100000, \
                               'loan_amount': 1000}}
                       )

print(response)

```

```
{"decision": "accept"}
```

4.3.6 Iterative Prompting

- Build on the AI's response by asking follow-up questions or refining the output.

- **Example:**

- *Initial Prompt:* “Help me to make a decision to accpet or reject the loan application.”

- *Follow-Up*: “give me the reason”

```
# Few-Shot Prompting
from langchain_ollama.llms import OllamaLLM
from langchain_core.prompts import ChatPromptTemplate
from langchain.output_parsers import CommaSeparatedListOutputParser

# Here are examples of loan application decision:
# 'example': {'input': {'fico':800, 'income':100000,'loan_amount': 10000}
# 'decision': "accept"
# Now Help me to make a decision to accpet or reject the loan application and
# give the reason.
# 'input': {"fico':820, 'income':100000, 'loan_amount': 1,000}"'

template = """you are a {role}
task: {task}
examples: {example}
input: {input}
decision:
reason:
"""

prompt = ChatPromptTemplate.from_template(template)
model = OllamaLLM(temperature=0.0, model=MODEL, format='json')
output_parser = CommaSeparatedListOutputParser()

chain = prompt | model

response = chain.invoke({'role': 'banker', \
                        'task': "Help me to make a decision to accpet or \
                                reject the loan application and \
                                give the reason.",\
                        'example': {'input': {'fico':800, 'income':100000,\\
                                             'loan_amount': 10000},\\
                                    'decision': "accept"}, \
                        'input': {'fico':820, 'income':100000, \
                                  'loan_amount': 1000}\
                      })

print(response)
```

```
{"decision": "accept", "reason": "The applicant has a high credit score (FICO ↗820), a stable income of $100,000, and is requesting a relatively small loan ↗amount ($1000). These factors indicate a low risk for the bank."}
```

4.3.7 Instructional Chaining

- Break down a task into a sequence of smaller prompts.

- **Example:**

- step 1: check the fico score
- step 2: check the income,
- step 3: check the loan amount,
- step 4: make a decision,
- step 5: give the reason.

```
# Instructional Chaining
from langchain_ollama.llms import OllamaLLM
from langchain_core.prompts import ChatPromptTemplate
from langchain.output_parsers import CommaSeparatedListOutputParser

# Now Help me to make a decision to accpet or reject the loan application and
# give the reason.
# "input": {"fico":320, "income":10000, "loan_amount": 100000}

template = """you are a {role}
task: {task}
instruction: {instruction}
input: {input}
decision:
reason:
"""

prompt = ChatPromptTemplate.from_template(template)
model = OllamaLLM(temperature=0.0, model=MODEL, format='json')
output_parser = CommaSeparatedListOutputParser()

chain = prompt | model

response = chain.invoke({'role': 'banker', \
                        'task': "Help me to make a decision to accpet or \
                                reject the loan application and \
                                give the reason.", \
                        'instruction': {'step 1': "check the fico score", \
                                      'step 2': "check the income", \
                                      'step 3': "check the loan amount", \
                                      'step 4': "make a decision", \
                                      'step 5': "give the reason"}, \
                        'input': {'fico':320, 'income':10000, \
                                'loan_amount': 100000}})
```

(continues on next page)

(continued from previous page)

```
        'loan_amount': 100000}  
    })  
  
print(response)
```

```
{  
    "decision": "reject",  
    "reason": "Based on the provided information, the applicant's FICO score is ↵320 which falls below our minimum acceptable credit score. Additionally, the ↵proposed loan amount of $100,000 exceeds the income level of $10,000 per year, ↵making it difficult for the borrower to repay the loan."  
}
```

4.3.8 Use Constraints

- Impose constraints to keep responses concise and on-topic.
- **Example:**

“List 5 key trends in AI in bullet points, each under 15 words.”

4.3.9 Creative Prompting

- Encourage unique or unconventional ideas by framing the task creatively.
- **Example:**

“Pretend you are a time traveler from the year 2124. How would you describe AI advancements to someone today?”

4.3.10 Feedback Incorporation

- If the response isn’t perfect, guide the AI to refine or retry.
- **Example:**

“This is too general. Could you provide more specific examples for the education industry?”

4.3.11 Scenario-Based Prompts

- Frame the query within a scenario for a contextual response.
- **Example:**

“Imagine you’re a teacher explaining ChatGPT to students. How would you introduce its uses and limitations?”

4.3.12 Multimodal Prompting

- Use prompts designed for mixed text/image inputs (or outputs if using models like DALL-E).

- **Example:**

“Generate an image prompt for a futuristic cityscape, vibrant, with flying cars and greenery.”

RETRIEVAL-AUGMENTED GENERATION

Colab Notebook for This Chapter

- Naive Chunking: [Open in Colab](#)
- Late Chunking: [Open in Colab](#)
- Reciprocal Rank Fusion: [Open in Colab](#)
- RAG in colab: [Open in Colab](#)

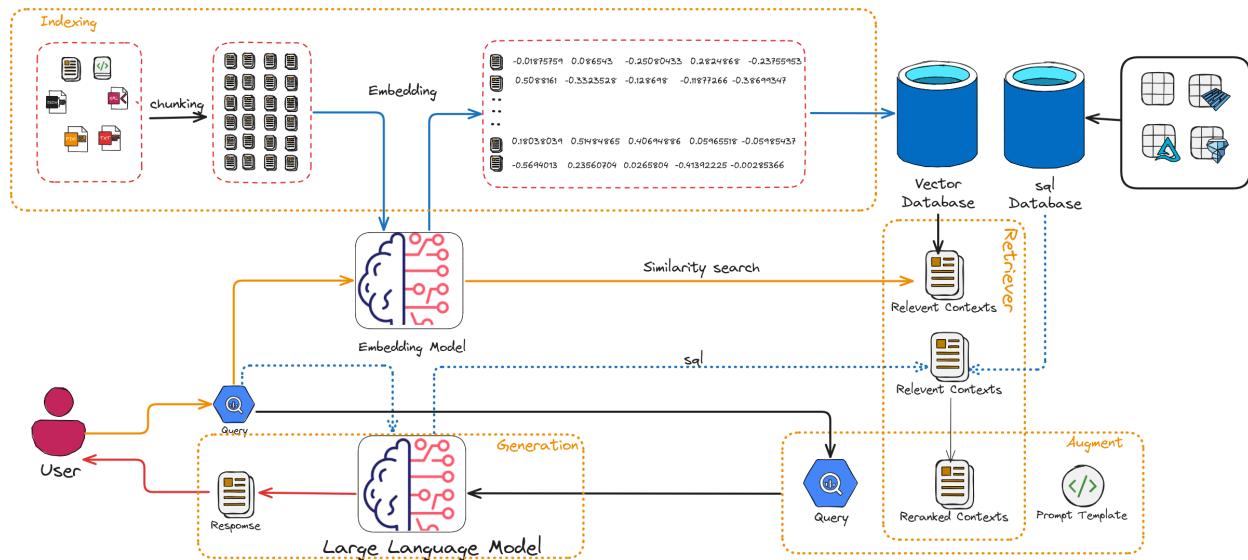


Fig. 1: Retrieval-Augmented Generation Diagram

Note

The naive chunking strategy was used in the diagram above. More advanced strategies, such as Late Chunking [lateChunking] (or Chunked Pooling), are discussed later in this chapter.

5.1 Overview

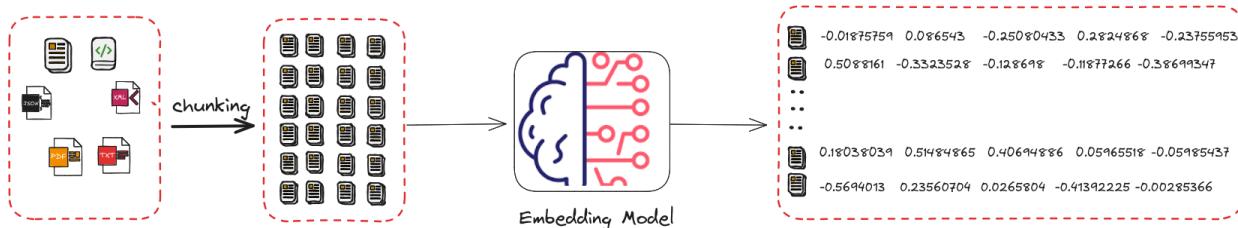
Retrieval-Augmented Generation (RAG) is a framework that enhances large language models (LLMs) by combining their generative capabilities with external knowledge retrieval. The goal of RAG is to improve accuracy, relevance, and factuality by providing the LLM with specific, up-to-date, or domain-specific context from a knowledge base or database during the generation process.

As you can see in [Retrieval-Augmented Generation Diagram](#), the RAG has three main components

- Indexer: The indexer processes raw text or other forms of unstructured data and creates an efficient structure (called an index) that allows for fast and accurate retrieval by the retriever when a query is made.
- Retriever: Responsible for finding relevant information from an external knowledge source, such as a document database, a vector database, or the web.
- Generator: An LLM (like GPT-4, T5, or similar) that uses the retrieved context to generate a response. The model is “augmented” with the retrieved information, which reduces hallucination and enhances factual accuracy.

5.2 Indexing

The indexing processes raw text or other forms of unstructured data and creates an efficient structure (called an index) that allows for fast and accurate retrieval by the retriever when a query is made.



5.2.1 Naive Chunking

Chunking in Retrieval-Augmented Generation (RAG) involves splitting documents or knowledge bases into smaller, manageable pieces (chunks) that can be efficiently retrieved and used by a language model (LLM).

Below are the common chunking strategies used in RAG workflows:

1. Fixed-Length Chunking

- Chunks are created with a predefined, fixed length (e.g., 200 words or 512 tokens).
- Simple and easy to implement but might split content mid-sentence or lose semantic coherence.

Example:

```
def fixed_length_chunking(text, chunk_size=200):
    words = text.split()
    return [
        " ".join(words[i:i + chunk_size])
```

(continues on next page)

(continued from previous page)

```

    for i in range(0, len(words), chunk_size)
]

# Example Usage
document = "This is a sample document with multiple sentences to demonstrate fixed-length chunking."
chunks = fixed_length_chunking(document, chunk_size=10)
for idx, chunk in enumerate(chunks):
    print(f"Chunk {idx + 1}: {chunk}")

```

Output:

- Chunk 1: This is a sample document with multiple sentences to demonstrate
- Chunk 2: fixed-length chunking.

2. Sliding Window Chunking

- Creates overlapping chunks to preserve context across splits.
- Ensures important information in overlapping regions is retained.

Example:

```

def sliding_window_chunking(text, chunk_size=100, overlap_size=20):
    words = text.split()
    chunks = []
    for i in range(0, len(words), chunk_size - overlap_size):
        chunk = " ".join(words[i:i + chunk_size])
        chunks.append(chunk)
    return chunks

# Example Usage
document = "This is a sample document with multiple sentences to demonstrate sliding window chunking."
chunks = sliding_window_chunking(document, chunk_size=10, overlap_size=3)
for idx, chunk in enumerate(chunks):
    print(f"Chunk {idx + 1}: {chunk}")

```

Output:

- Chunk 1: This is a sample document with multiple sentences to demonstrate
- Chunk 2: with multiple sentences to demonstrate sliding window chunking.
- Chunk 3: sliding window chunking.

3. Semantic Chunking

- Splits text based on natural language boundaries such as paragraphs, sentences, or specific delimiters (e.g., headings).
- Retains semantic coherence, ideal for better retrieval and generation accuracy.

Example:

```

import nltk
nltk.download('punkt_tab')

def semantic_chunking(text, sentence_len=50):
    sentences = nltk.sent_tokenize(text)

    chunks = []
    chunk = ""
    for sentence in sentences:
        if len(chunk.split()) + len(sentence.split()) <= sentence_len:
            chunk += " " + sentence
        else:
            chunks.append(chunk.strip())
            chunk = sentence
    if chunk:
        chunks.append(chunk.strip())
    return chunks

# Example Usage
document = ("This is a sample document. It is split based on semantic boundaries. "
            "Each chunk will have coherent meaning for better retrieval.")
chunks = semantic_chunking(document, 10)
for idx, chunk in enumerate(chunks):
    print(f"Chunk {idx + 1}: {chunk}")

```

Output:

- Chunk 1: This is a sample document.
- Chunk 2: It is split based on semantic boundaries.
- Chunk 3: Each chunk will have coherent meaning for better retrieval.

4. Dynamic Chunking

- Adapts chunk sizes based on content properties such as token count, content density, or specific criteria.
- Useful when handling diverse document types with varying information density.

Example:

```

from transformers import AutoTokenizer

def dynamic_chunking(text, max_tokens=200, tokenizer_name="bert-base-uncased"):
    tokenizer = AutoTokenizer.from_pretrained(tokenizer_name)
    tokens = tokenizer.encode(text, add_special_tokens=False)

```

(continues on next page)

(continued from previous page)

```

chunks = []
for i in range(0, len(tokens), max_tokens):
    chunk = tokens[i:i + max_tokens]
    chunks.append(tokenizer.decode(chunk))
return chunks

# Example Usage
document = ("This is a sample document to demonstrate dynamic chunking. "
             "The tokenizer adapts the chunks based on token limits.")
chunks = dynamic_chunking(document, max_tokens=10)
for idx, chunk in enumerate(chunks):
    print(f"Chunk {idx + 1}: {chunk}")

```

Output:

- Chunk 1: this is a sample document to demonstrate dynamic chunking
- Chunk 2: . the tokenizer adapts the chunks based on
- Chunk 3: token limits.

Comparison of Strategies

Strategy	Pros	Cons
Fixed-Length Chunking	Simple, fast	May split text mid-sentence or lose coherence.
Sliding Window Chunking	Preserves context	Overlapping increases redundancy.
Semantic Chunking	Coherent chunks	Requires NLP preprocessing.
Dynamic Chunking	Adapts to content	Computationally intensive.

Each strategy has its strengths and weaknesses. Select based on the task requirements, context, and available computational resources.

The optimal chunk length depends on the type of content being processed and the intended use case. Below are recommendations for chunk lengths based on different context types, along with their rationale:

Context Type	Chunk Length (Tokens)	Rationale
FAQs or Short Texts	100-200	Short enough to handle specific queries.
Articles or Blog Posts	300-500	Covers logical sections while fitting multiple chunks in the LLM context.
Research Papers or Reports	500-700	Captures detailed sections like methodology or results.
Legal or Technical Texts	200-300	Maintains precision due to dense information.

Thevaluating Chunking Strategies for Retrieval can be found at: <https://research.trychroma.com/>

evaluating-chunking

5.2.2 Late Chunking

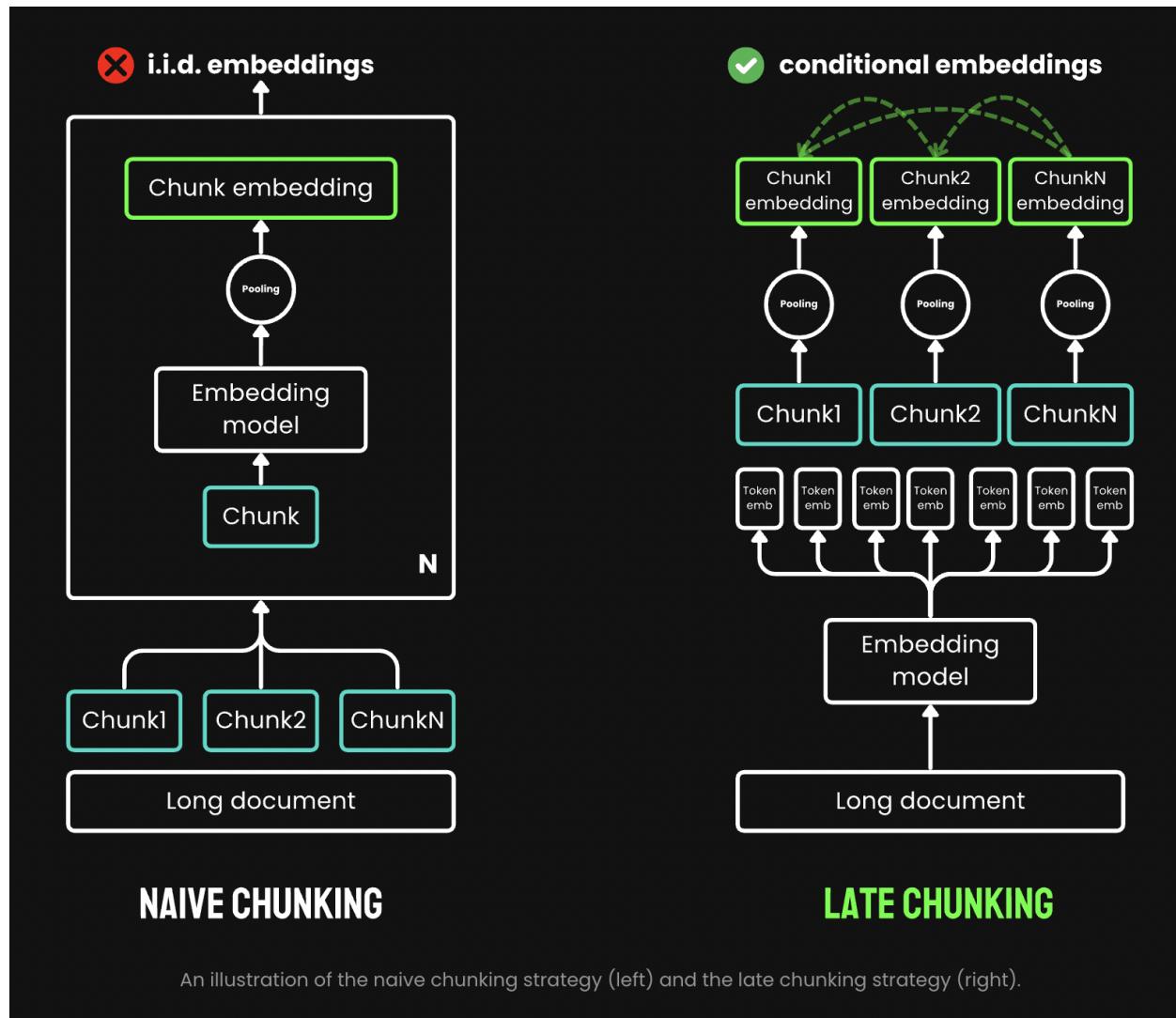


Fig. 2: An illustration of the naive chunking strategy (left) and the late chunking strategy (right). (Source Jina AI)

Late Chunking refers to a strategy in Retrieval-Augmented Generation (RAG) where chunking of data is **deferred until query time**. Unlike pre-chunking, where documents are split into chunks during preprocessing, late chunking dynamically extracts relevant content when a query is made.

- Key Concepts of Late Chunking
 - **Dynamic Chunk Creation:**
 - * Full documents or large sections are stored in the vector database.
 - * Relevant chunks are dynamically extracted at query time based on the query and similarity match.

- **Query-Time Optimization:**

- * The system identifies relevant content using similarity search or semantic analysis.
- * Only the most relevant content is chunked and passed to the language model.

- **Reduced Preprocessing Time:**

- * Eliminates extensive preprocessing and fixed chunking during data ingestion.
- * Higher computational cost occurs during query-time retrieval.

The following implementations are from [Jina AI](#), and the copyright belongs to the original author.

```
def chunk_by_sentences(input_text: str, tokenizer: callable):
    """
    Split the input text into sentences using the tokenizer
    :param input_text: The text snippet to split into sentences
    :param tokenizer: The tokenizer to use
    :return: A tuple containing the list of text chunks and their corresponding
    token spans
    """

    inputs = tokenizer(input_text, return_tensors='pt', return_offsets_
    ↪mapping=True)
    punctuation_mark_id = tokenizer.convert_tokens_to_ids('.')
    sep_id = tokenizer.convert_tokens_to_ids('[SEP]')
    token_offsets = inputs['offset_mapping'][0]
    token_ids = inputs['input_ids'][0]
    chunk_positions = [
        (i, int(start + 1))
        for i, (token_id, (start, end)) in enumerate(zip(token_ids, token_
    ↪offsets))
        if token_id == punctuation_mark_id
        and (
            token_offsets[i + 1][0] - token_offsets[i][1] > 0
            or token_ids[i + 1] == sep_id
        )
    ]
    chunks = [
        input_text[x[1] : y[1]]
        for x, y in zip([(1, 0)] + chunk_positions[:-1], chunk_positions)
    ]
    span_annotations = [
        (x[0], y[0]) for (x, y) in zip([(1, 0)] + chunk_positions[:-1], chunk_
    ↪positions)
    ]
    return chunks, span_annotations

def late_chunking(
    model_output: 'BatchEncoding', span_annotation: list, max_length=None
```

(continues on next page)

(continued from previous page)

```

): token_embeddings = model_output[0]
outputs = []
for embeddings, annotations in zip(token_embeddings, span_annotation):
    if (
        max_length is not None
    ): # remove annotations which go beyond the max-length of the model
        annotations = [
            (start, min(end, max_length - 1))
            for (start, end) in annotations
            if start < (max_length - 1)
        ]
        pooled_embeddings = [
            embeddings[start:end].sum(dim=0) / (end - start)
            for start, end in annotations
            if (end - start) >= 1
        ]
        pooled_embeddings = [
            embedding.detach().cpu().numpy() for embedding in pooled_embeddings
        ]
        outputs.append(pooled_embeddings)

return outputs

```

```

input_text = "Berlin is the capital and largest city of Germany, both by area\u202d\u202d and by population. Its more than 3.85 million inhabitants make it the European\u202d\u202d Union's most populous city, as measured by population within city limits. The\u202d\u202d city is also one of the states of Germany, and is the third smallest state in\u202d\u202d the country in terms of area."
# determine chunks
chunks, span_annotations = chunk_by_sentences(input_text, tokenizer)
print('Chunks:\n- ' + '\n- '.join(chunks) + '\n')

# chunk before
embeddings_traditional_chunking = model.encode(chunks)

# chunk afterwards (context-sensitive chunked pooling)
inputs = tokenizer(input_text, return_tensors='pt')
model_output = model(**inputs)
embeddings = late_chunking(model_output, [span_annotations])[0]

import numpy as np

```

(continues on next page)

(continued from previous page)

```

cos_sim = lambda x, y: np.dot(x, y) / (np.linalg.norm(x) * np.linalg.norm(y))

berlin_embedding = model.encode('Berlin')

for chunk, new_embedding, trad_embeddings in zip(chunks, embeddings, embeddings_
→traditional_chunking):
    print(f'similarity_new("Berlin", "{chunk}"):', cos_sim(berlin_embedding, new_
→embedding))
    print(f'similarity_trad("Berlin", "{chunk}"):', cos_sim(berlin_embedding,_
→trad_embeddings))

```

Quer Chunk		similar- ity_new	similar- ity_trad
Berli Berlin is the capital and largest city of Germany, both by area and by popu- lation.		0.849546	0.8486219
Berli Its more than 3.85 million inhabitants make it the European Union's most populous city, as measured by population within city limits."		0.8248902	0.70843387
Berli The city is also one of the states of Germany, and is the third smallest state in the country in terms of area."		0.8498009	0.75345534

5.2.3 Types of Indexing

The embedding methods we introduced in Chapter *Word and Sentence Embedding* can be applied here to convert each chunk into embeddings and create indexing. These indexings(embeddings) will be used to retrieve relevant documents or information.

- Sparse Indexing:

Uses traditional keyword-based methods (e.g., TF-IDF, BM25). Index stores the frequency of terms and their associations with documents.

- Advantages: Easy to understand and deploy and works well for exact matches or keyword-heavy queries.
- Disadvantages: Struggles with semantic understanding or paraphrased queries.

- Dense Indexing:

Uses vector embeddings to capture semantic meaning. Documents are represented as vectors in a high-dimensional space, enabling similarity search.

- Advantages: Excellent for semantic search, handling synonyms, and paraphrasing.
- Disadvantages: Requires more computational resources for storage and retrieval.

- Hybrid Indexing:

Combines sparse and dense indexing for more robust search capabilities. For example, Elasticsearch can integrate BM25 with vector search.

5.2.4 Vector Database

Vector databases are essential for Retrieval-Augmented Generation (RAG) systems, enabling efficient similarity search on dense vector embeddings. Below is a comprehensive overview of popular vector databases for RAG workflows:

1. FAISS (Facebook AI Similarity Search)

- **Description:** - An open-source library developed by Facebook AI for efficient similarity search and clustering of dense vectors.
- **Features:** - High performance and scalability. - Supports various indexing methods like Flat, IVF, and HNSW. - GPU acceleration for faster searches.
- **Use Cases:** - Research and prototyping. - Scenarios requiring custom implementations.
- **Limitations:** - File-based storage; lacks a built-in distributed or managed cloud solution.
- **Official Website:** [FAISS GitHub](#)

2. Pinecone

- **Description:** - A fully managed vector database designed for production-scale workloads.
- **Features:** - Scalable and serverless architecture. - Automatic scaling and optimization of indexes. - Hybrid search (combining vector and keyword search). - Integrates with popular frameworks like LangChain and OpenAI.
- **Use Cases:** - Enterprise-grade applications. - Handling large datasets with minimal operational overhead.
- **Official Website:** [Pinecone](#)

3. Weaviate

- **Description:** - An open-source vector search engine with a strong focus on modularity and customization.
- **Features:** - Supports hybrid search and symbolic reasoning. - Schema-based data organization. - Plugin support for pre-built and custom vectorization modules. - Cloud-managed and self-hosted options.
- **Use Cases:** - Applications requiring hybrid search capabilities. - Knowledge graphs and semantically rich data.
- **Official Website:** [Weaviate](#)

4. Milvus

- **Description:** - An open-source, high-performance vector database designed for similarity search on large datasets.
- **Features:** - Distributed and scalable architecture. - Integration with FAISS, Annoy, and HNSW indexing techniques. - Built-in support for time travel queries (searching historical data).
- **Use Cases:** - Video, audio, and image search applications. - Large-scale datasets requiring real-time indexing and retrieval.

- **Official Website:** [Milvus](#)

5. Qdrant

- **Description:** - An open-source, lightweight vector database focused on ease of use and modern developer needs.
- **Features:** - Supports HNSW for efficient vector search. - Advanced filtering capabilities for combining metadata with vector queries. - REST and gRPC APIs for integration. - Docker-ready deployment.
- **Use Cases:** - Scenarios requiring metadata-rich search. - Lightweight deployments with simplicity in mind.
- **Official Website:** [Qdrant](#)

6. Redis (with Vector Similarity Search Module)

- **Description:** - A popular in-memory database with a module for vector similarity search.
- **Features:** - Combines vector search with traditional key-value storage. - Supports hybrid search and metadata filtering. - High throughput and low latency due to in-memory architecture.
- **Use Cases:** - Applications requiring real-time, low-latency search. - Integrating vector search with existing Redis-based systems.
- **Official Website:** [Redis Vector Search](#)

7. Zilliz

- **Description:** - A cloud-native vector database built on Milvus for scalable and managed vector storage.
- **Features:** - Fully managed service for vector data. - Seamless scaling and distributed indexing. - Integration with machine learning pipelines.
- **Use Cases:** - Large-scale enterprise deployments. - Cloud-native solutions with minimal infrastructure management.
- **Official Website:** [Zilliz](#)

8. Vespa

- **Description:** - A real-time serving engine supporting vector and hybrid search.
- **Features:** - Combines vector search with advanced ranking and filtering. - Scales to large datasets with support for distributed clusters. - Powerful query configuration options.
- **Use Cases:** - E-commerce and recommendation systems. - Applications with complex ranking requirements.
- **Official Website:** [Vespa](#)

9. Chroma

- **Description:** - An open-source, user-friendly vector database built for LLMs and embedding-based applications.

- **Features:** - Designed specifically for RAG workflows. - Simple Python API for seamless integration with AI models. - Efficient and customizable vector storage for embedding data.
- **Use Cases:** - Prototyping and experimentation for LLM-based applications. - Lightweight deployments for small to medium-scale RAG systems.
- **Official Website:** [Chroma](#)

Comparison of Vector Databases:

Database	Open Source	Managed Service	Key Features	Best For
FAISS	Yes	No	High performance, GPU acceleration	Research, prototyping
Pinecone	No	Yes	Serverless, automatic scaling	Enterprise-scale applications
Weaviate	Yes	Yes	Hybrid search, modularity	Knowledge graphs
Milvus	Yes	No	Distributed, high performance	Large-scale datasets
Qdrant	Yes	No	Lightweight, metadata filtering	Small to medium-scale apps
Redis	No	Yes	In-memory performance, hybrid search	Real-time apps
Zilliz	No	Yes	Fully managed Milvus	Enterprise cloud solutions
Vespa	Yes	No	Hybrid search, real-time ranking	E-commerce, recommendations
Chroma	Yes	No	LLM-focused, simple API	Prototyping, lightweight apps

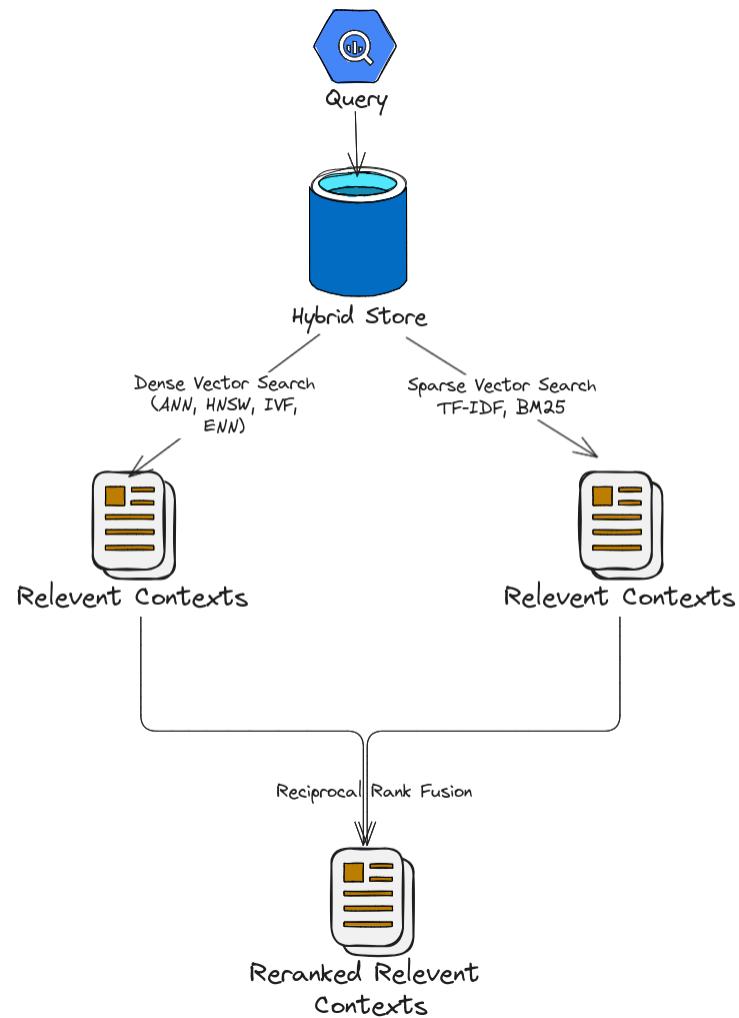
Choosing a Vector Database

- **For Research or Small Projects:** FAISS, Qdrant, Milvus, or Chroma.
- **For Enterprise or Cloud-Native Workflows:** Pinecone, Zilliz, or Weaviate.
- **For Real-Time Use Cases:** Redis or Vespa.

Each database has unique strengths and is suited for specific RAG use cases. The choice depends on scalability, integration needs, and budget.

5.3 Retrieval

The retriever selects “chunks” of text (e.g., paragraphs or sections) relevant to the user’s query.



5.3.1 Common retrieval methods

- Sparse Vector Search: Traditional keyword-based retrieval (e.g., TF-IDF, BM25).
- Dense Vector Search: Vector-based search using embeddings e.g.
 - **Approximate Nearest Neighbor (ANN) Search:**
 - * HNSW (Hierarchical Navigable Small World): Graph-based approach
 - * IVF (Inverted File Index): Clusters embeddings into groups and searches within relevant clusters.
 - Exact Nearest Neighbor Search: Computes similarities exhaustively for all vectors in the corpus
- Hybrid Search: the combination of Sparse and Dense vector search.

Summary of Common Algorithms:

Met- ric/Algorith m	Purpose	Common Use
TF-IDF	Keyword matching with term weighting.	Effective for small-scale or structured corpora.
BM25	Advanced keyword matching with term frequency saturation and document length normalization.	Widely used in sparse search; default in tools like Elasticsearch and Solr.
Cosine Similarity	Measures orientation (ignores magnitude).	Widely used; works well with normalized vectors.
Dot Product Similarity	Measures magnitude and direction.	Preferred in embeddings like OpenAI's models.
Euclidean Distance	Measures absolute distance between vectors.	Less common but used in some specific cases.
HNSW (ANN)	Fast and scalable nearest neighbor search.	Default for large-scale systems (e.g., FAISS).
IVF (ANN)	Efficient clustering-based search.	Often combined with product quantization.

5.3.2 Reciprocal Rank Fusion

Reciprocal Rank Fusion (RRF) is a ranking technique commonly used in information retrieval and ensemble learning. Although it is not specific to large language models (LLMs), it can be applied to scenarios where multiple ranking systems (or scoring mechanisms) produce different rankings, and you want to combine them into a single, unified ranking.

The reciprocal rank of an item in a ranked list is calculated as $\frac{1}{k+r}$, where

- r is the rank of the item (1 for the top rank, 2 for the second rank, etc.).
- k is a small constant (often set to 60 or another fixed value) to control how much weight is given to higher ranks.

Example:

Suppose two retrieval models give ranked lists for query responses:

- Model 1 ranks documents as: [A,B,C,D]
- Model 2 ranks documents as: [B,A,D,C]

RRF combines these rankings by assigning each document a combined score:

- Document A: $\frac{1}{60+1} + \frac{1}{60+2} = 0.03252247488101534$
- Document B: $\frac{1}{60+2} + \frac{1}{60+1} = 0.03252247488101534$
- Document C: $\frac{1}{60+3} + \frac{1}{60+4} = 0.03149801587301587$
- Document D: $\frac{1}{60+4} + \frac{1}{60+3} = 0.03149801587301587$

```
from collections import defaultdict

def reciprocal_rank_fusion(ranked_results: list[list], k=60):
    """
    Fuse rank from multiple retrieval systems using Reciprocal Rank Fusion.

    Args:
        ranked_results: Ranked results from different retrieval system.
        k (int): A constant used in the RRF formula (default is 60).

    Returns:
        Tuple of list of sorted documents by score and sorted documents
    """

    # Dictionary to store RRF mapping
    rrf_map = defaultdict(float)

    # Calculate RRF score for each result in each list
    for rank_list in ranked_results:
        for rank, item in enumerate(rank_list, 1):
            rrf_map[item] += 1 / (rank + k)

    # Sort items based on their RRF scores in descending order
    sorted_items = sorted(rrf_map.items(), key=lambda x: x[1], reverse=True)

    # Return tuple of list of sorted documents by score and sorted documents
    return sorted_items, [item for item, score in sorted_items]

# Example ranked lists from different sources
ranked_a = ['A', 'B', 'C', 'D']
ranked_b = ['B', 'A', 'D', 'C']

# Combine the lists using RRF
```

(continues on next page)

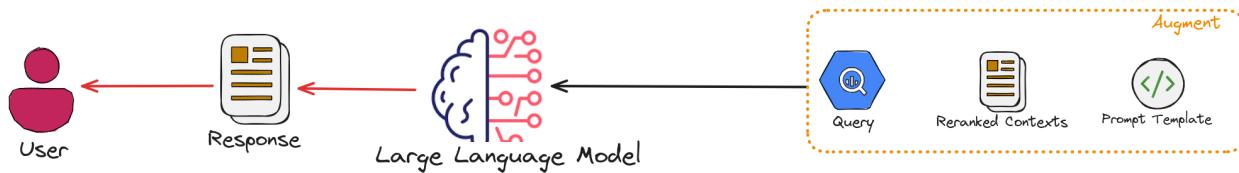
(continued from previous page)

```
combined_list = reciprocal_rank_fusion([ranked_a, ranked_b])
print(combined_list)
```

```
([('A', 0.03252247488101534), ('B', 0.03252247488101534), ('C', 0.
˓→03149801587301587), ('D', 0.03149801587301587)], ['A', 'B', 'C', 'D'])
```

5.4 Generation

Finally, the retrieved relevant information will be feed back into the LLMs to generate responses.



5.5 Advanced Topic

5.5.1 Corrective RAG

https://langchain-ai.github.io/langgraph/tutorials/rag/langgraph_crag_local/

5.5.2 Self-RAG

https://langchain-ai.github.io/langgraph/tutorials/rag/langgraph_self_rag_local/

5.5.3 Adaptive RAG

https://langchain-ai.github.io/langgraph/tutorials/rag/langgraph_adaptive_rag/ https://langchain-ai.github.io/langgraph/tutorials/rag/langgraph_adaptive_rag_local/

5.5.4 Agentic RAG

https://langchain-ai.github.io/langgraph/tutorials/rag/langgraph_agentic_rag/

CHAPTER SIX

FINE TUNING

Chinese proverb

Good tools are prerequisite to the successful execution of a job. – old Chinese proverb

Colab Notebook for This Chapter

- Embedding Model Fine-tuning: [!\[\]\(b7293d18f7684ef0ad2146a473d2c14a_img.jpg\) Open in Colab](#)
- LLM (Llama 2 7B) Model Fine-tuning: [!\[\]\(6f59378d55a5234cc18de57cdfb912e2_img.jpg\) Open in Colab](#)

Fine-tuning is a machine learning technique where a pre-trained model (like a large language model or neural network) is further trained on a smaller, specific dataset to adapt it to a particular task or domain. Instead of training a model from scratch, fine-tuning leverages the knowledge already embedded in the pre-trained model, saving time, computational resources, and data requirements.

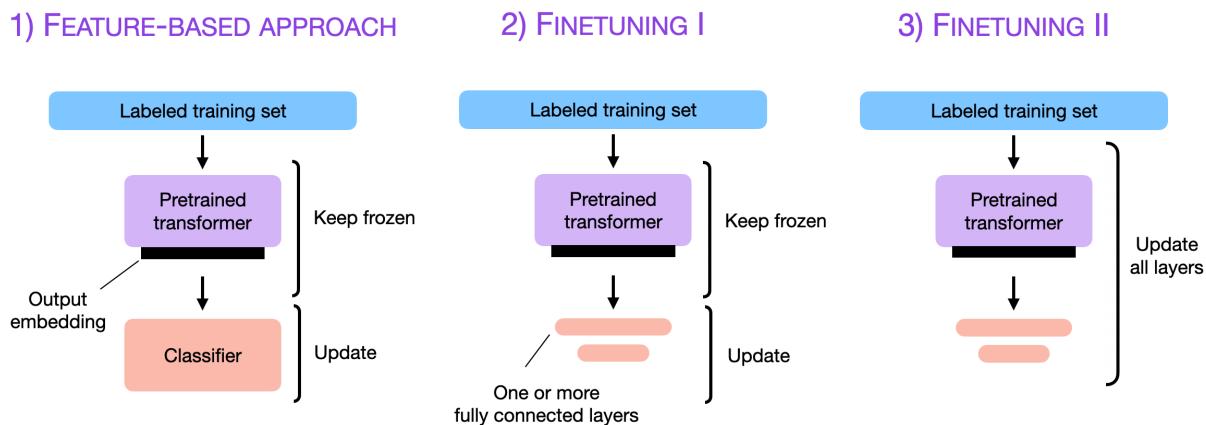


Fig. 1: The three conventional feature-based and finetuning approaches (Source Finetuning Sebastian).

6.1 Cutting-Edge Strategies for LLM Fine-Tuning

Over the past year, fine-tuning methods have made remarkable strides. Modern methods for fine-tuning LLMs focus on efficiency, scalability, and resource optimization. The following strategies are at the forefront:

6.1.1 LoRA (Low-Rank Adaptation)

LoRA reduces the number of trainable parameters by introducing **low-rank decomposition** into the fine-tuning process.

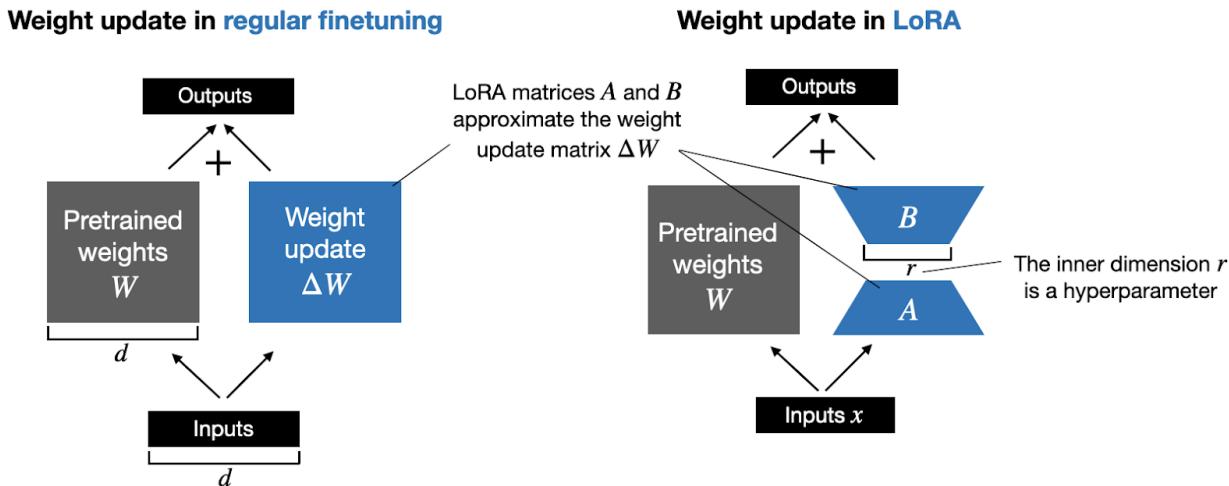


Fig. 2: Weight update matrix (Source [LORA Sebastian](#)).

How It Works:

- Instead of updating all model weights, LoRA injects **low-rank adapters** into the model's layers.
- The original pre-trained weights remain frozen; only the low-rank parameters are optimized.

Benefits:

- Reduces memory and computational requirements.
- Enables fine-tuning on resource-constrained hardware.

6.1.2 QLoRA (Quantized Low-Rank Adaptation)

QLoRA combines **low-rank adaptation** with **4-bit quantization** of the pre-trained model.

How It Works:

- The LLM is quantized to **4-bit precision** to reduce memory usage.
- LoRA adapters are applied to the quantized model for fine-tuning.
- Precision is maintained using methods like **NF4 (Normalized Float 4)** and double backpropagation.

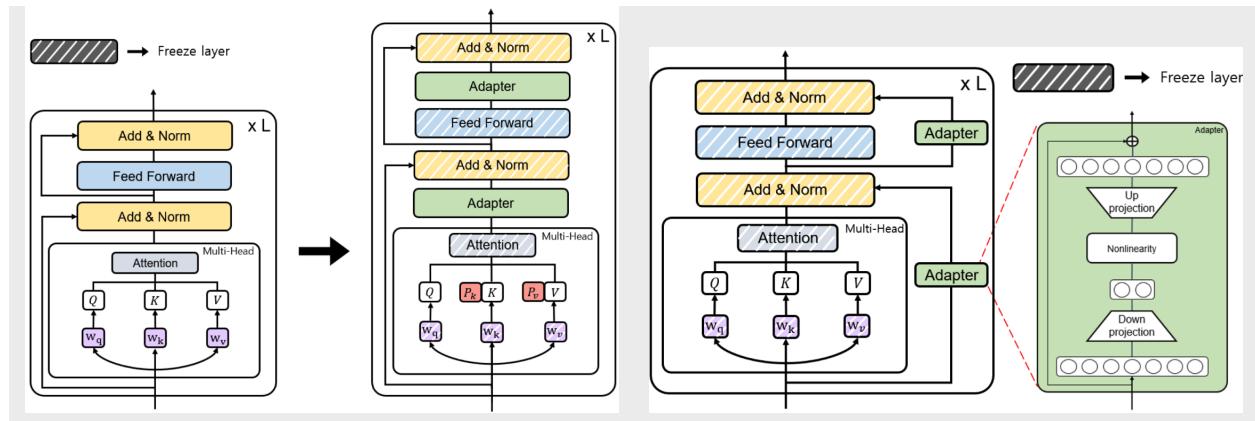
Benefits:

- Further reduces memory usage compared to LoRA.

- Enables fine-tuning of massive models on consumer-grade GPUs.

6.1.3 PEFT (Parameter-Efficient Fine-Tuning)

PEFT is a general framework for fine-tuning LLMs with minimal trainable parameters.



Source: [PEFT]

Techniques Under PEFT:

- **LoRA:** Low-rank adaptation of weights.
- **Adapters:** Small trainable layers inserted into the model.
- **Prefix Tuning:** Fine-tuning input prefixes instead of weights.
- **Prompt Tuning:** Optimizing soft prompts in the input space.

Benefits:

- Reduces the number of trainable parameters.
- Faster training and lower hardware requirements.

6.1.4 SFT (Supervised Fine-Tuning)

SFT adapts an LLM using a labeled dataset in a fully supervised manner.

How It Works:

- The model is initialized with pre-trained weights.
- It is fine-tuned on a task-specific dataset with a supervised loss function (e.g., cross-entropy).

Benefits:

- Achieves high performance on specific tasks.
- Essential for aligning models with labeled datasets.

6.1.5 Summary Table

Method	Description	Key Benefit
LoRA	Low-rank adapters for parameter-efficient tuning.	Reduces trainable parameters significantly.
QLoRA	LoRA with 4-bit quantization of the model.	Fine-tunes massive models on smaller hardware.
PEFT	General framework for efficient fine-tuning.	Includes LoRA, Adapters, Prefix Tuning, etc.
SFT	Supervised fine-tuning with labeled data.	High performance on task-specific datasets

These strategies represent the forefront of **LLM fine-tuning**, offering efficient and scalable solutions for real-world applications. To choose the most suitable strategy, consider the following factors:

- **Resource-Constrained Environments:** Use LoRA or QLoRA.
- **Large-Scale Models:** QLoRA for low-memory fine-tuning.
- **High Performance with Labeled Data:** SFT.
- **Minimal Setup:** Zero-shot or Few-shot learning.
- **General Efficiency:** Use PEFT frameworks.

6.2 Key Early Fine-Tuning Methods

Early fine-tuning methods laid the foundation for current approaches. These methods primarily focused on updating the entire model or selected components.

6.2.1 Full Fine-Tuning

All the parameters of a pre-trained model are updated using task-specific data *The three conventional feature-based and finetuning approaches (Souce Finetuning Sebastian)*. (right).

How It Works:

- The pre-trained model serves as the starting point.
- Fine-tuning is conducted on a smaller, labeled dataset using a supervised loss function.
- A low learning rate is used to prevent **catastrophic forgetting**.

Benefits:

- Effective at adapting models to specific tasks.

Challenges:

- Computationally expensive.
- Risk of overfitting on small datasets.

6.2.2 Feature-Based Approach

The pre-trained model is used as a **feature extractor**, while only a task-specific head is trained *The three conventional feature-based and finetuning approaches (Souce Finetuning Sebastian)*. (left).

How It Works:

- The model processes inputs and extracts features (embeddings).
- A separate classifier (e.g., linear or MLP) is trained on top of these features.
- The pre-trained model weights remain **frozen**.

Benefits:

- Computationally efficient since only the task-specific head is trained.

6.2.3 Layer-Specific Fine-Tuning

Only certain layers of the pre-trained model are fine-tuned while the rest remain frozen *The three conventional feature-based and finetuning approaches (Souce Finetuning Sebastian)*. (middle).

How It Works:

- Earlier layers (which capture general features) are frozen.
- Later layers (closer to the output) are fine-tuned on task-specific data.

Benefits:

- Balances computational efficiency and task adaptation.

6.2.4 Task-Adaptive Pre-training

Before fine-tuning on a specific task, the model undergoes additional **pre-training** on a domain-specific corpus.

How It Works:

- A general pre-trained model is further pre-trained (unsupervised) on domain-specific data.
- Fine-tuning is then performed on the downstream task.

Benefits:

- Provides a better starting point for domain-specific tasks.

6.3 Embedding Model Fine-Tuning

In the chapter *Retrieval-Augmented Generation*, we discussed how embedding models are crucial for the success of RAG applications. However, their general-purpose training often limits their effectiveness for company- or domain-specific use cases. Customizing embeddings with domain-specific data can significantly improve the retrieval performance of your RAG application.

In this chapter, we will demonstrate how to fine-tune embedding models using the `SentenceTransformersTrainer`, building on insights shared in the blog [fineTuneEmbedding] and

Sentence Transformer Training Overview. Our main contribution was introducing LoRA to enable functionality on NVIDIA T4 GPUs, while the rest of the pipeline and code remained almost unchanged.

Note

Please ensure that the package versions are set as follows:

```
pip install "torch==2.1.2" tensorboard

pip install --upgrade \
    sentence-transformers>=3 \
    datasets==2.19.1 \
    transformers==4.41.2 \
    peft==0.10.0
```

Otherwise, you may encounter the error.

6.3.1 Prepare Dataset

We are going to directly use the synthetic dataset `philschmid/finanical-rag-embedding-dataset`, which includes 7,000 positive text pairs of questions and corresponding context from the [2023_10 NVIDIA SEC Filing](#).

```
from datasets import load_dataset

# Load dataset from the hub
dataset = load_dataset("philschmid/finanical-rag-embedding-dataset", split="train")
dataset = dataset[0:7000]

# rename columns
dataset = dataset.rename_column("question", "anchor")
dataset = dataset.rename_column("context", "positive")

# Add an id column to the dataset
dataset = dataset.add_column("id", range(len(dataset)))

# split dataset into a 10% test set
dataset = dataset.train_test_split(test_size=0.1)

# save datasets to disk
dataset["train"].to_json("train_dataset.json", orient="records")
dataset["test"].to_json("test_dataset.json", orient="records")
```

Note

In practice, most dataset configurations will take one of four forms:

- **Positive Pair:** A pair of related sentences. This can be used both for symmetric tasks (semantic

textual similarity) or asymmetric tasks (semantic search), with examples including pairs of paraphrases, pairs of full texts and their summaries, pairs of duplicate questions, pairs of (query, response), or pairs of (source_language, target_language). Natural Language Inference datasets can also be formatted this way by pairing entailing sentences.

- **Triplets:** (anchor, positive, negative) text triplets. These datasets don't need labels.
- **Pair with Similarity Score:** A pair of sentences with a score indicating their similarity. Common examples are "Semantic Textual Similarity" datasets.
- **Texts with Classes:** A text with its corresponding class. This data format is easily converted by loss functions into three sentences (triplets) where the first is an "anchor", the second a "positive" of the same class as the anchor, and the third a "negative" of a different class.

Note that it is often simple to transform a dataset from one format to another, such that it works with your loss function of choice.

6.3.2 Import and Evaluate Pretrained Baseline Model

```
import torch
from sentence_transformers import SentenceTransformer
from sentence_transformers.evaluation import (
    InformationRetrievalEvaluator,
    SequentialEvaluator,
)
from sentence_transformers.util import cos_sim
from datasets import load_dataset, concatenate_datasets
from peft import LoraConfig, TaskType

model_id = "BAAI/bge-base-en-v1.5"
matryoshka_dimensions = [768, 512, 256, 128, 64] # Important: large to small

# Load a model
model = SentenceTransformer(
    model_id,
    trust_remote_code=True,
    device="cuda" if torch.cuda.is_available() else "cpu"
)

# load test dataset
test_dataset = load_dataset("json", data_files="test_dataset.json", split="train")
train_dataset = load_dataset("json", data_files="train_dataset.json", split="train")
corpus_dataset = concatenate_datasets([train_dataset, test_dataset])

# Convert the datasets to dictionaries
corpus = dict(
```

(continues on next page)

(continued from previous page)

```

    zip(corpus_dataset["id"], corpus_dataset["positive"])
) # Our corpus (cid => document)
queries = dict(
    zip(test_dataset["id"], test_dataset["anchor"])
) # Our queries (qid => question)

# Create a mapping of relevant document (1 in our case) for each query
relevant_docs = {} # Query ID to relevant documents (qid => set([relevant_cids]))
for q_id in queries:
    relevant_docs[q_id] = [q_id]

matryoshka_evaluators = []
# Iterate over the different dimensions
for dim in matryoshka_dimensions:
    ir_evaluator = InformationRetrievalEvaluator(
        queries=queries,
        corpus=corpus,
        relevant_docs=relevant_docs,
        name=f"dim_{dim}",
        truncate_dim=dim, # Truncate the embeddings to a certain dimension
        score_functions={"cosine": cos_sim},
    )
    matryoshka_evaluators.append(ir_evaluator)

# Create a sequential evaluator
evaluator = SequentialEvaluator(matryoshka_evaluators)

```

Note

If you encounter the error `Cannot import name 'EncoderDecoderCache' from 'transformers'`, ensure that the package versions are set to `peft==0.10.0` and `transformers==4.37.2`.

```

# Evaluate the model
results = evaluator(model)

# Print the main score
for dim in matryoshka_dimensions:
    key = f"dim_{dim}_cosine_ndcg@10"
    print
    print(f"{key}: {results[key]}")

```

`dim_768_cosine_ndcg@10: 0.754897248109794`

(continues on next page)

(continued from previous page)

```
dim_512_cosine_ndcg@10: 0.7549275773474213
dim_256_cosine_ndcg@10: 0.7454714780163237
dim_128_cosine_ndcg@10: 0.7116728650043451
dim_64_cosine_ndcg@10: 0.6477174937632066
```

6.3.3 Loss Function with Matryoshka Representation

```
from sentence_transformers import SentenceTransformerModelCardData,
    SentenceTransformer

# Hugging Face model ID: https://huggingface.co/BAAI/bge-base-en-v1.5
model_id = "BAAI/bge-base-en-v1.5"

# load model with SDPA for using Flash Attention 2
model = SentenceTransformer(
    model_id,
    model_kwarg={"attn_implementation": "sdpa"},
    model_card_data=SentenceTransformerModelCardData(
        language="en",
        license="apache-2.0",
        model_name="BGE base Financial Matryoshka",
    ),
)

# Apply PEFT with PromptTuningConfig
peft_config = LoraConfig(
    task_type=TaskType.FEATURE_EXTRACTION,
    inference_mode=False,
    r=8,
    lora_alpha=32,
    lora_dropout=0.1,
)
model.add_adapter(peft_config, "dense")

# train loss
from sentence_transformers.losses import MatryoshkaLoss,
    MultipleNegativesRankingLoss

matryoshka_dimensions = [768, 512, 256, 128, 64] # Important: large to small
inner_train_loss = MultipleNegativesRankingLoss(model)
train_loss = MatryoshkaLoss(model,
                           inner_train_loss,
                           matryoshka_dims=matryoshka_dimensions)
```

Note

Loss functions play a critical role in the performance of your fine-tuned model. Sadly, there is no “one size fits all” loss function. Ideally, this table should help narrow down your choice of loss function(s) by matching them to your data formats.

You can often convert one training data format into another, allowing more loss functions to be viable for your scenario. For example,

Inputs	La-bels	Appropriate Loss Functions
single sentences	class	BatchAllTripletLoss, BatchHardSoftMarginTripletLoss, BatchHardTripletLoss, BatchSemiHardTripletLoss
single sentences	none	ContrastiveTensionLoss, DenoisingAutoEncoderLoss
(anchor, anchor) pairs	none	ContrastiveTensionLossInBatchNegatives
(damaged_sent	none	DenoisingAutoEncoderLoss
original_sent pairs		
(sentence_A, sentence_B)	class	SoftmaxLoss
pairs		
(anchor, positive)	none	MultipleNegativesRankingLoss, CachedMultipleNegativesRankingLoss, MultipleNegativesSymmetricRankingLoss, CachedMultipleNegativesSymmetricRankingLoss, MegaBatchMarginLoss, GISTEmbedLoss, CachedGISTEmbedLoss
pairs		
(anchor, positive/negative)	1 if positive, 0 if negative	ContrastiveLoss, OnlineContrastiveLoss
pairs		
(sentence_A, sentence_B)	float similarity score	CoSENTLoss, AngleLoss, CosineSimilarityLoss
pairs		
(anchor, positive, negative)	none	MultipleNegativesRankingLoss, CachedMultipleNegativesRankingLoss, TripletLoss, CachedGISTEmbedLoss, GISTEmbedLoss
triplets		
(anchor, positive, negative_1, ..., negative_n)	none	MultipleNegativesRankingLoss, CachedMultipleNegativesRankingLoss, CachedGISTEmbedLoss

6.3.4 Fine-tune Embedding Model

```

from sentence_transformers import SentenceTransformerTrainingArguments
from sentence_transformers.training_args import BatchSamplers

# load train dataset again
train_dataset = load_dataset("json", data_files="train_dataset.json", split=
    "train")

# define training arguments
args = SentenceTransformerTrainingArguments(
    output_dir=output_dir, # output directory and hugging face model ID
    num_train_epochs=4, # number of epochs
    per_device_train_batch_size=32, # train batch size
    gradient_accumulation_steps=16, # for a global batch size of 512
    per_device_eval_batch_size=16, # evaluation batch size
    warmup_ratio=0.1, # warmup ratio
    learning_rate=2e-5, # learning rate, 2e-5 is a good
    value
    lr_scheduler_type="cosine", # use constant learning rate
    scheduler
    optim="adamw_torch_fused", # use fused adamw optimizer
    tf32=False, # use tf32 precision
    bf16=False, # use bf16 precision
    batch_sampler=BatchSamplers.NO_DUPLICATES, # MultipleNegativesRankingLoss
    benefits from no duplicate samples in a batch
    eval_strategy="epoch", # evaluate after each epoch
    save_strategy="epoch", # save after each epoch
    logging_steps=10, # log every 10 steps
    save_total_limit=3, # save only the last 3 models
    load_best_model_at_end=True, # load the best model when
    training ends
    metric_for_best_model="eval_dim_128_cosine_ndcg@10", # Optimizing for the
    best ndcg@10 score for the 128 dimension
    greater_is_better=True, # maximize the ndcg@10 score
)

from sentence_transformers import SentenceTransformerTrainer

trainer = SentenceTransformerTrainer(
    model=model, # bg-base-en-v1
    args=args, # training arguments
    train_dataset=train_dataset.select_columns(
        ["anchor", "positive"]
    ), # training dataset
    loss=train_loss,
    evaluator=evaluator,
)

```

(continues on next page)

(continued from previous page)

```
)
# start training
trainer.train()

# save the best model
#trainer.save_model()
trainer.model.save_pretrained("bge-base-finetuning")
```

6.3.5 Evaluate Fine-tuned Model

```
from sentence_transformers import SentenceTransformer

fine_tuned_model = SentenceTransformer(
    'bge-base-finetuning', device="cuda" if torch.cuda.is_available() else "cpu"
)
# Evaluate the model
results = evaluator(fine_tuned_model)

# # COMMENT IN for full results
# print(results)

# Print the main score
for dim in matryoshka_dimensions:
    key = f"dim_{dim}_cosine_ndcg@10"
    print(f"{key}: {results[key]}")
```

```
dim_768_cosine_ndcg@10: 0.7650276801072632
dim_512_cosine_ndcg@10: 0.7603951540556889
dim_256_cosine_ndcg@10: 0.754743133407988
dim_128_cosine_ndcg@10: 0.7205317098443929
dim_64_cosine_ndcg@10: 0.6609117856061502
```

6.3.6 Results Comparison

Although we did not observe the significant performance boost reported in the original blog, the fine-tuned model outperformed the baseline model across all dimensions using only 6.3k samples and partial parameter fine-tuning. More details can be found as follows:

Dimension	Baseline	Fine-tuned	Improvement
768	0.75490	0.76503	1.34%
512	0.75492	0.76040	0.73%
256	0.74547	0.75474	1.24%
128	0.71167	0.72053	1.24%
64	0.64772	0.66091	2.04%

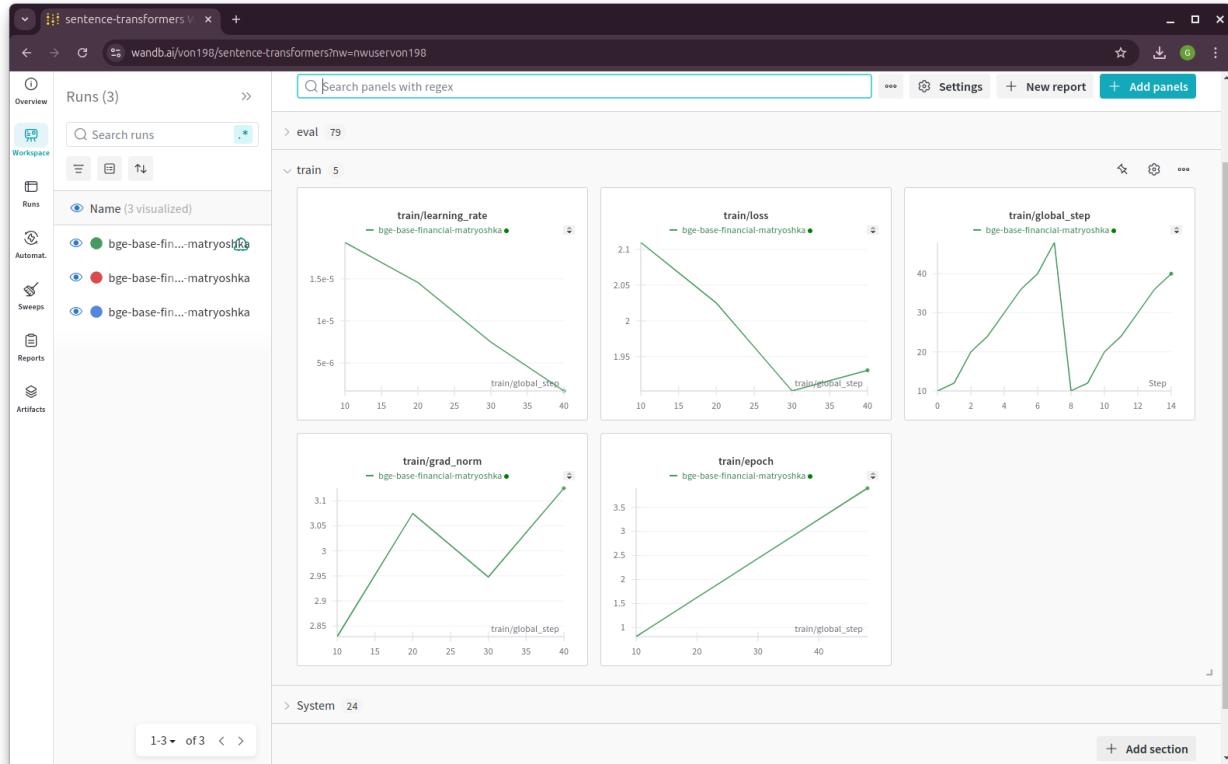


Fig. 3: Epoch, Training Loss/steps in Wandb

6.4 LLM Fine-Tuning

In this chapter, we will demonstrate how to fine-tune a Llama 2 model with 7 billion parameters using a T4 GPU with 16 GB of VRAM. Due to VRAM limitations, traditional fine-tuning is not feasible, making parameter-efficient fine-tuning (PEFT) techniques like LoRA or QLoRA essential. For this demonstration, we use QLoRA, which leverages 4-bit precision to significantly reduce VRAM consumption.

The following code is from notebook [fineTuneLLM], and the copyright belongs to the original author.

6.4.1 Load Dataset and Pretrained Model

```
# Step 1 : Load dataset (you can process it here)
dataset = load_dataset(dataset_name, split="train")

# Step 2 : Load tokenizer and model with QLoRA configuration
compute_dtype = getattr(torch, bnb_4bit_compute_dtype)

bnb_config = BitsAndBytesConfig(
    load_in_4bit=use_4bit,
    bnb_4bit_quant_type=bnb_4bit_quant_type,
    bnb_4bit_compute_dtype=compute_dtype,
    bnb_4bit_use_double_quant=use_nested_quant,
```

(continues on next page)

(continued from previous page)

```

)
# Step 3 :Check GPU compatibility with bfloat16
if compute_dtype == torch.float16 and use_4bit:
    major, _ = torch.cuda.get_device_capability()
    if major >= 8:
        print("=" * 80)
        print("Your GPU supports bfloat16: accelerate training with bf16=True")
        print("=" * 80)

# Step 4 :Load base model
model = AutoModelForCausalLM.from_pretrained(
    model_name,
    quantization_config=bnb_config,
    device_map=device_map
)
model.config.use_cache = False
model.config.pretraining_tp = 1

# Step 5 :Load LLaMA tokenizer
tokenizer = AutoTokenizer.from_pretrained(model_name, trust_remote_code=True)
tokenizer.add_special_tokens({'pad_token': '[PAD]'})
tokenizer.pad_token = tokenizer.eos_token
tokenizer.padding_side = "right"

```

6.4.2 Fine-tuning Configuration

```

# Step 6 :Load LoRA configuration
peft_config = LoraConfig(
    lora_alpha=lora_alpha,
    lora_dropout=lora_dropout,
    r=lora_r,
    bias="none",
    task_type="CAUSAL_LM",
)

# Step 7 :Set training parameters
training_arguments = TrainingArguments(
    output_dir=output_dir,
    num_train_epochs=num_train_epochs,
    per_device_train_batch_size=per_device_train_batch_size,
    gradient_accumulation_steps=gradient_accumulation_steps,
    optim=optim,
    save_steps=save_steps,
    logging_steps=logging_steps,

```

(continues on next page)

(continued from previous page)

```
learning_rate=learning_rate,
weight_decay=weight_decay,
fp16=fp16,
bf16=bf16,
max_grad_norm=max_grad_norm,
max_steps=max_steps,
warmup_ratio=warmup_ratio,
group_by_length=group_by_length,
lr_scheduler_type=lr_scheduler_type,
report_to="tensorboard"
)
```

6.4.3 Fine-tune model

```
# Step 8 :Set supervised fine-tuning parameters
trainer = SFTTrainer(
    model=model,
    train_dataset=dataset,
    peft_config=peft_config,
    dataset_text_field="text",
    max_seq_length=max_seq_length,
    tokenizer=tokenizer,
    args=training_arguments,
    packing=packing,
)

# Step 9 :Train model
trainer.train()

# Step 10 :Save trained model
trainer.model.save_pretrained(new_model)
```

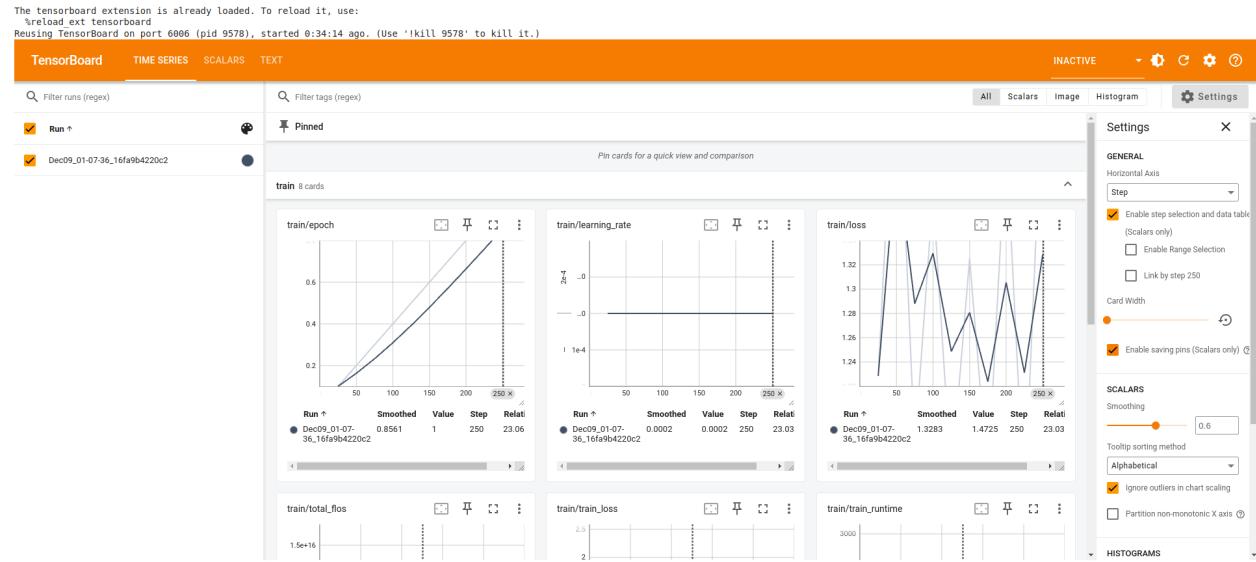


Fig. 4: Llama 2 Model Fine-Tuning TensorBoard

PRE-TRAINING

Proverb

Pre-training as we know it will end. – Ilya Sutskever at neurips 2024

Pre-training as we know it will end

Compute is growing:

- Better hardware
- Better algorithms
- Larger clusters

Data is not growing:

- We have but one internet
- **The fossil fuel of AI**

Fig. 1: Ilya Sutskever at neurips 2024

In industry, most companies focus primarily on prompt engineering, RAG, and fine-tuning, while advanced techniques like pre-training from scratch or deep model customization remain less common due to the significant resources and expertise required.

LLMs, like GPT (Generative Pre-trained Transformer), BERT (Bidirectional Encoder Representations from Transformers), and others, are large-scale models built using the transformer architecture. These models are

trained on vast amounts of text data to learn patterns in language, enabling them to generate human-like text, answer questions, summarize information, and perform other natural language processing tasks.

This chapter delves into transformer models, drawing on insights from [The Annotated Transformer](#) and [Tracing the Transformer in Diagrams](#), to explore their underlying architecture and practical applications.

CHAPTER EIGHT

LLM EVALUATION METRICS

Colab Notebook for This Chapter

GEval with DeepEval: [Open in Colab](#)

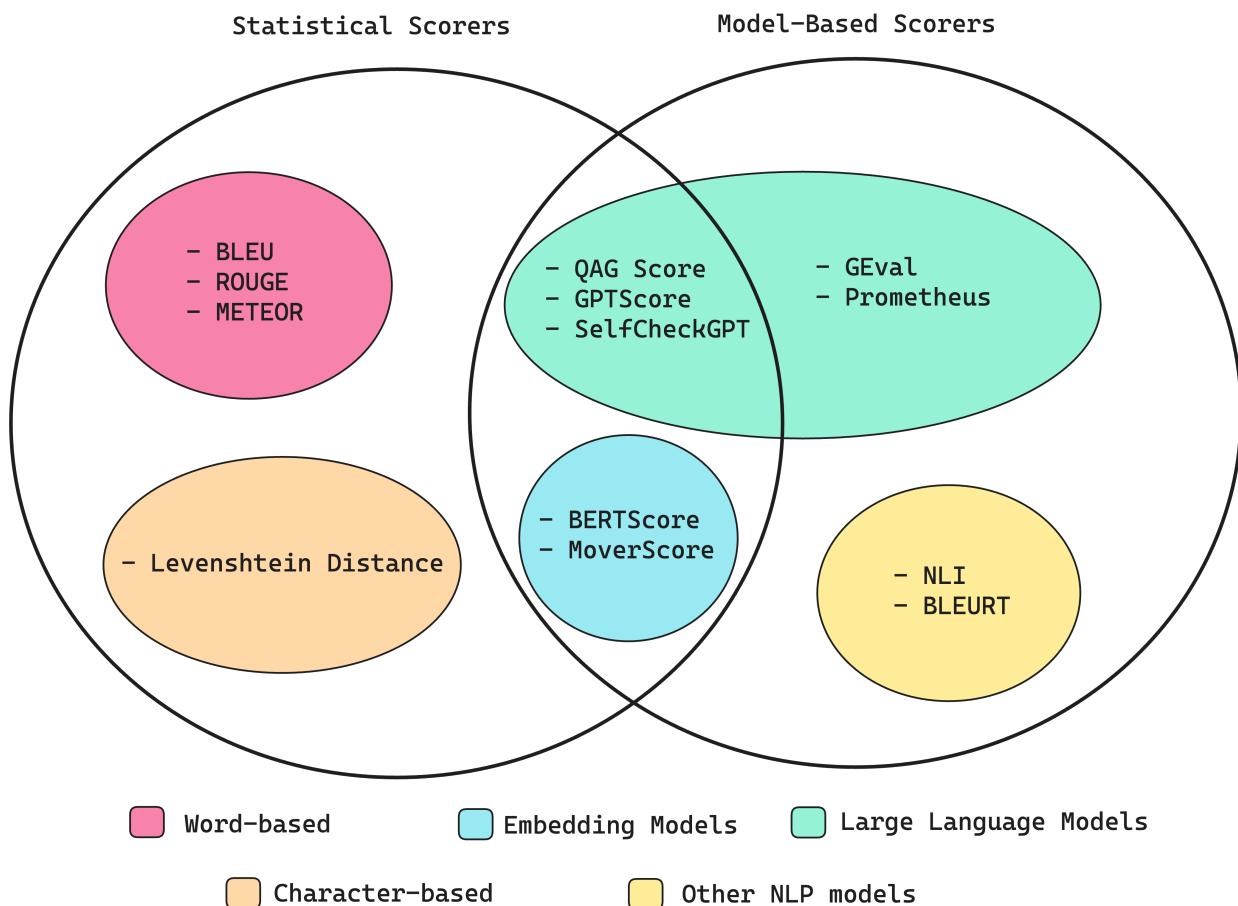


Fig. 1: Types of metric scorers (Source: LLM Evaluation Metrics: The Ultimate LLM Evaluation Guide)

8.1 Statistical Scorers (Traditional Metrics)

These metrics evaluate text outputs based on statistical comparisons to references or expected outputs.

Note

I completely agree with the author of [LLM Evaluation Metrics: The Ultimate LLM Evaluation Guide](#) that statistical scoring methods are, in my opinion, non-essential to focus on. These methods tend to perform poorly whenever reasoning is required, making them too inaccurate as scorers for most LLM evaluation criteria. Additionally, more advanced metrics, such as GEval [*GEval with DeepEval*](#), provide significantly better alternatives.

- **BLEU (Bilingual Evaluation Understudy):** Measures overlap of n-grams between generated and reference texts. Common for translation tasks.
- **ROUGE (Recall-Oriented Understudy for Gisting Evaluation):** Focuses on recall of n-grams (ROUGE-N), longest common subsequences (ROUGE-L), and skip bigrams (ROUGE-S). Popular in summarization tasks.
- **METEOR (Metric for Evaluation of Translation with Explicit ORdering):** Considers synonymy and paraphrasing via stemming and synonym matching.
- **TER (Translation Edit Rate):** Measures the number of edits required to turn the generated output into the reference text.
- **CIDEr (Consensus-based Image Description Evaluation):** Designed for image captioning, using TF-IDF weighting of n-grams.
- **BERTScore:** Leverages contextual embeddings (e.g., BERT) to compute similarity between generated and reference texts.
- **GLEU (Google BLEU):** A variation of BLEU designed for grammatical error correction tasks.

8.2 Model-Based Scorers (Learned Metrics)

These metrics employ models trained to assess the quality of generated text, often based on human annotations.

- **BLEURT:** Combines pre-trained models (e.g., BERT) with fine-tuning on human judgment data.
- **COMET (Cross-lingual Optimized Metric for Evaluation of Translation):** A neural network model trained on translation quality data.
- **PRISM:** Measures semantic similarity by paraphrasing both the hypothesis and reference into a shared space.
- **UniEval:** A unified framework for evaluation across multiple tasks, focusing on both factual accuracy and linguistic quality.
- **Perplexity:** Estimates the likelihood of generated text under the original model's probability distribution (lower is better).
- **GPTScore:** Uses a large pre-trained LLM (e.g., GPT-4) to rate the quality of outputs.

- **MAUVE:** Measures the divergence between the distribution of generated text and that of human-written text.
- **DRIFT:** Focuses on domain-specific evaluation, checking how well outputs align with domain-specific data distributions.

8.3 Human-Centric Evaluations (Augmenting Metrics)

While not automated, human evaluations are crucial for assessing subjective qualities such as:

- **Fluency**
- **Coherence**
- **Relevance**
- **Factuality**
- **Style Appropriateness**

Both statistical and model-based scorers are often used in tandem with human evaluation to ensure a holistic assessment of LLM outputs.

8.4 GEval with DeepEval

GEval is a recently developed evaluation framework developed from paper [GEval] to assess large language models (LLMs) using GPT-based evaluators. It leverages the capabilities of advanced LLMs (like GPT-4 or beyond) to rate and critique the outputs of other models, including themselves, across various tasks. This approach shifts the evaluation paradigm by relying on the intrinsic understanding and reasoning power of the models, rather than traditional metrics.

8.4.1 GEval Algorithm

8.4.2 GEval with DeepEval

In DeepEval, a metric serves as a standard for measuring the performance of an LLM's output based on specific criteria of interest. Essentially, while the metric functions as the “ruler”, a test case represents the subject being measured. DeepEval, provides a variety of default metrics to help you get started quickly, including:

- G-Eval
- Summarization
- Faithfulness
- Answer Relevancy
- Contextual Relevancy
- Contextual Precision
- Contextual Recall
- Ragas

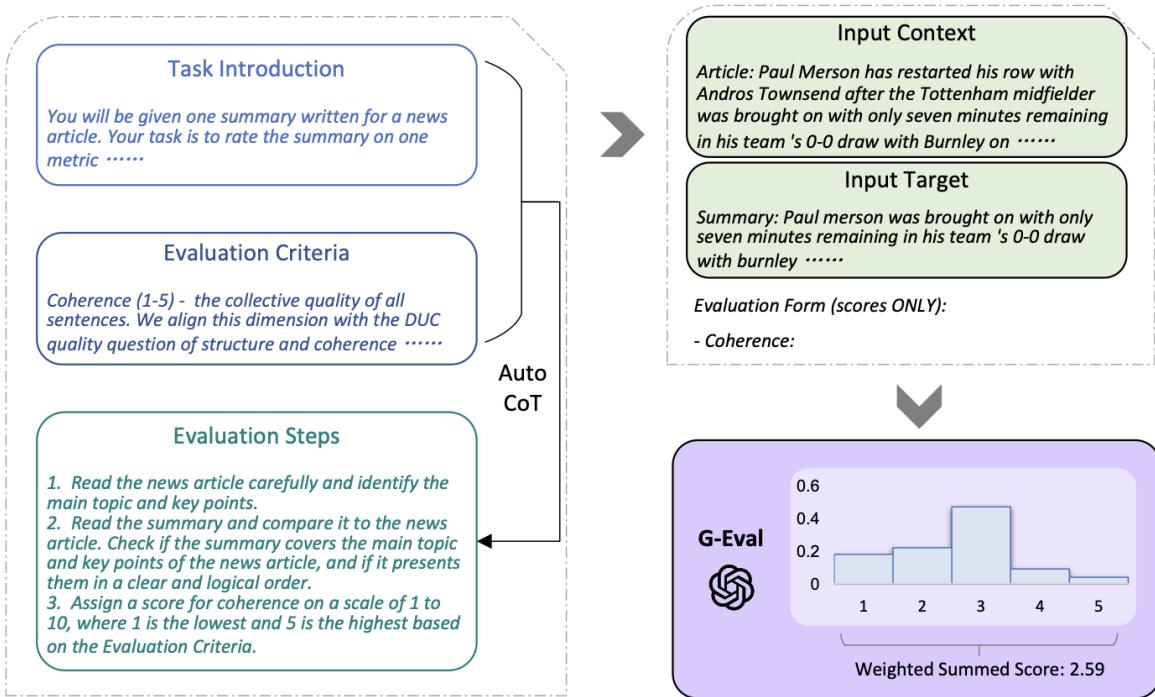


Figure 1: The overall framework of G-EVAL. We first input Task Introduction and Evaluation Criteria to the LLM, and ask it to generate a CoT of detailed Evaluation Steps. Then we use the prompt along with the generated CoT to evaluate the NLG outputs in a form-filling paradigm. Finally, we use the probability-weighted summation of the output scores as the final score.

Fig. 2: G-Eval Algorithm (Source: [GEval])

- Hallucination
- Toxicity
- Bias

DeepEval also provides conversational metrics, designed to evaluate entire conversations rather than individual, granular LLM interactions. These include:

- Conversation Completeness
- Conversation Relevancy
- Knowledge Retention
- **Set Up Local Model**

```
deepeval set-local-model --model-name='mistral' \
--base-url="http://localhost:11434/v1/" \
--api-key="ollama"
```

- **Default Metrics**

- AnswerRelevancyMetric

```
from deepeval import evaluate
from deepeval.metrics import AnswerRelevancyMetric
from deepeval.test_case import LLMTTestCase

answer_relevancy_metric = AnswerRelevancyMetric(threshold=0.7)
test_case = LLMTTestCase(
    input="What if these shoes don't fit?",
    # Replace this with the actual output from your LLM
    application="We offer a 30-day full refund at no extra costs.",
    retrieval_context=["All customers are eligible for a 30 day full refund at no extra costs."]
)
evaluate([test_case], [answer_relevancy_metric])
```

- * Metrics Summary

- Answer Relevancy (score: 1.0, threshold: 0.7, strict: False, evaluation model: local model, reason: The score is 1.00 because it directly and accurately answered the question about shoe fitting, making it highly relevant., error: None)

- * For test case:

- input: What if these shoes don't fit?
- actual output: We offer a 30-day full refund at no extra costs.
- expected output: None

- context: None
- retrieval context: ['All customers are eligible for a 30 day full refund at no extra costs.']}
- * Overall Metric Pass Rates

Answer Relevancy: 100.00% pass rate

```
EvaluationResult(test_results=[TestResult(name='test_case_0',  
    success=True, metrics_data=[MetricData(name='Answer Relevancy',  
        threshold=0.7, success=True, score=1.0, reason='The score is  
        1.00 because it directly and accurately answered the question  
        about shoe fitting, making it highly relevant.', strict_  
        mode=False, evaluation_model='local model', error=None,  
        evaluation_cost=0.0, verbose_logs='Statements:\n\n    "We  
        offer a 30-day full refund",\n    "The refund does not incur  
        any additional costs"\n    \n    Verdicts:\n    {\n        "verdict": "yes",\n        "reason": "The statements about the  
        refund policy are relevant to addressing the input, which asks  
        about what to do if the shoes don\'t fit."\n    },\n    {\n        "verdict": "yes",\n        "reason": "The statement that  
        the refund does not incur any additional costs is also  
        relevant as it provides further information about the refund  
        process.\n    }\n}]), conversational=False, multimodal=False,  
    input="What if these shoes don't fit?", actual_output='We  
offer a 30-day full refund at no extra costs.', expected_  
output=None, context=None, retrieval_context=['All customers  
are eligible for a 30 day full refund at no extra costs.']),  
    confident_link=None)
```

- FaithfulnessMetric

```
from deepeval import evaluate
from deepeval.metrics import FaithfulnessMetric
from deepeval.test_case import LLMTTestCase

# input
input = "What if these shoes don't fit?"

# Replace this with the actual output from your LLM application
actual_output = "We offer a 30-day full refund at no extra cost."

# Replace this with the actual retrieved context from your RAG  
pipeline
retrieval_context = ["All customers are eligible for a 30 day  
full refund at no extra cost."]
```

(continues on next page)

(continued from previous page)

```

metric = FaithfulnessMetric(
    threshold=0.7,
    #model="gpt-4",
    include_reason=True
)
test_case = LLMTestCase(
    input=input,
    actual_output=actual_output,
    retrieval_context=retrieval_context
)

metric.measure(test_case)
print(metric.score)
print(metric.reason)

# or evaluate test cases in bulk
evaluate([test_case], [metric])

```

* Metrics Summary

- Faithfulness (score: 1.0, threshold: 0.7, strict: False, evaluation model: local model, reason: The faithfulness score is 1.00 because there are no contradictions found between the actual output and the retrieval context., error: None)

* For test case:

- input: What if these shoes don't fit?
- actual output: We offer a 30-day full refund at no extra cost.
- expected output: None
- context: None
- retrieval context: ['All customers are eligible for a 30 day full refund at no extra cost.']}

* Overall Metric Pass Rates

Faithfulness: 100.00% pass rate

```

EvaluationResult(test_results=[TestResult(name='test_case_0',
success=True, metrics_data=[MetricData(name='Faithfulness',
threshold=0.7, success=True, score=1.0, reason='The
faithfulness score is 1.00 because there are no contradictions
found between the actual output and the retrieval context.',

strict_mode=False, evaluation_model='local model', error=None,
evaluation_cost=0.0, verbose_logs='Truths (limit=None):\n\n
"Customers are eligible for a 30 day full refund.",\n
"The
refund is at no extra cost."\n\n\nClaims:\n\n
"The

```

(continues on next page)

(continued from previous page)

```

→refund is offered for a period of 30 days.",\n      "The refund\n
→does not incur any additional costs."\n] \n\nVerdicts:\n[\n  ↳ {\n    "verdict": "yes",\n      "reason": null\n  },\n  ↳ {\n    "verdict": "yes",\n      "reason": null\n  }\n]\n)], conversational=False, multimodal=False, input="What\n→if these shoes don't fit?", actual_output='We offer a 30-day\n→full refund at no extra cost.', expected_output=None,\n→context=None, retrieval_context=['All customers are eligible\n→for a 30 day full refund at no extra cost.']), confident_\n→link=None)

```

- ContextualPrecisionMetric

```

from deepeval import evaluate
from deepeval.metrics import ContextualPrecisionMetric
from deepeval.test_case import LLMTTestCase

# input
input = "What if these shoes don't fit?"

# Replace this with the actual output from your LLM application
actual_output = "We offer a 30-day full refund at no extra cost."

# Replace this with the expected output from your RAG generator
expected_output = "You are eligible for a 30 day full refund at\n→no extra cost."

# Replace this with the actual retrieved context from your RAG
→pipeline
retrieval_context = ["All customers are eligible for a 30 day\n→full refund at no extra cost."]

metric = ContextualPrecisionMetric(
    threshold=0.7,
    #model="gpt-4",
    include_reason=True
)
test_case = LLMTTestCase(
    input=input,
    actual_output=actual_output,
    expected_output=expected_output,
    retrieval_context=retrieval_context
)

metric.measure(test_case)
print(metric.score)

```

(continues on next page)

(continued from previous page)

```
print(metric.reason)

# or evaluate test cases in bulk
evaluate([test_case], [metric])
```

* Metrics Summary

- Contextual Precision (score: 1.0, threshold: 0.7, strict: False, evaluation model: local model, reason: The contextual precision score is 1.00 because the node ranked first (with reason: ‘The text verifies that customers are indeed eligible for a 30 day full refund at no extra cost.’) is relevant and correctly placed as the highest-ranked response to the input ‘What if these shoes don’t fit?’. All other nodes, if present, should be ranked lower due to their irrelevance to the question., error: None)

* For test case:

- input: What if these shoes don’t fit?
- actual output: We offer a 30-day full refund at no extra cost.
- expected output: You are eligible for a 30 day full refund at no extra cost.
- context: None
- retrieval context: [‘All customers are eligible for a 30 day full refund at no extra cost.’]

* Overall Metric Pass Rates

- Contextual Precision: 100.00% pass rate

- ContextualRecallMetric

```
from deepeval import evaluate
from deepeval.metrics import ContextualRecallMetric
from deepeval.test_case import LLMTTestCase

metric = ContextualRecallMetric(
    threshold=0.7,
    model="gpt-4",
    include_reason=True
)
test_case = LLMTTestCase(
    input=input,
    actual_output=actual_output,
    expected_output=expected_output,
    retrieval_context=retrieval_context
)

metric.measure(test_case)
```

(continues on next page)

(continued from previous page)

```
print(metric.score)
print(metric.reason)

# or evaluate test cases in bulk
evaluate([test_case], [metric])
```

* Metrics Summary

- Contextual Recall (score: 1.0, threshold: 0.7, strict: False, evaluation model: local model, reason: The score is 1.00 because the expected output is exactly as stated in the retrieval context., error: None)

* For test case:

- input: What if these shoes don't fit?
- actual output: We offer a 30-day full refund at no extra cost.
- expected output: You are eligible for a 30 day full refund at no extra cost.
- context: None
- retrieval context: ['All customers are eligible for a 30 day full refund at no extra cost.]

* Overall Metric Pass Rates

- Contextual Recall: 100.00% pass rate

– HallucinationMetric

```
from deepeval import evaluate
from deepeval.metrics import HallucinationMetric
from deepeval.test_case import LLMTTestCase

# input

input = "What was the blond doing?"

# Replace this with the actual documents that you are passing as
→input to your LLM.
context=["A man with blond-hair, and a brown shirt drinking out
→of a public water fountain."]

# Replace this with the actual output from your LLM application
actual_output="A blond drinking water in public."

test_case = LLMTTestCase(
    input= input,
    actual_output=actual_output,
    context=context
```

(continues on next page)

(continued from previous page)

```

)
metric = HallucinationMetric(threshold=0.5)

metric.measure(test_case)
print(metric.score)
print(metric.reason)

# or evaluate test cases in bulk
evaluate([test_case], [metric])

```

* Metrics Summary

- Hallucination (score: 0.0, threshold: 0.5, strict: False, evaluation model: local model, reason: The score is 0.00 because the actual output correctly aligns with the provided context., error: None)

* For test case:

- input: What was the blond doing?
- actual output: A blond drinking water in public.
- expected output: None
- context: ['A man with blond-hair, and a brown shirt drinking out of a public water fountain.']}
- retrieval context: None

* Overall Metric Pass Rates

Hallucination: 100.00% pass rate

• Custom Metrics

```

from deepeval.metrics import GEval
from deepeval.test_case import LLMTTestCaseParams

correctness_metric = GEval(
    name="Correctness",
    criteria="Determine whether the actual output is factually correct based on the expected output.",
    # NOTE: you can only provide either criteria or evaluation_steps, and not both
    evaluation_steps=[
        "Check whether the facts in 'actual output' contradicts any facts in 'expected output'",
        "You should also heavily penalize omission of detail",
        "Vague language, or contradicting OPINIONS, are OK"
    ],
    evaluation_params=[LLMTTestCaseParams.INPUT, LLMTTestCaseParams.

```

(continues on next page)

(continued from previous page)

```

    ↵ACTUAL_OUTPUT, LLMTestCaseParams.EXPECTED_OUTPUT],
)

test_case = LLMTestCase(
    input="The dog chased the cat up the tree, who ran up the tree?
",
    actual_output="It depends, some might consider the cat, while
others might argue the dog.",
    expected_output="The cat."
)

correctness_metric.measure(test_case)
print(correctness_metric.score)
print(correctness_metric.reason)

```

Event loop **is** already running. Applying nest_asyncio patch to allow
→async execution...

0.1
The actual output does **not** match the expected output **and** omits
→specific details about which animal climbed the tree.

Evaluation Framework:

```

from deepeval import evaluate
from deepeval.metrics import GEval
from deepeval.test_case import LLMTestCase
from deepeval.test_case import LLMTestCaseParams

correctness_metric = GEval(
    name="Correctness",
    criteria="Determine whether the actual output is factually correct
based on the expected output.",
    # NOTE: you can only provide either criteria or evaluation_steps,
    # and not both
    evaluation_steps=[
        "Check whether the facts in 'actual output' contradicts any
facts in 'expected output'",
        "You should also heavily penalize omission of detail",
        "Vague language, or contradicting OPINIONS, are OK"
    ],
    evaluation_params=[LLMTestCaseParams.INPUT, LLMTestCaseParams.
ACTUAL_OUTPUT, LLMTestCaseParams.EXPECTED_OUTPUT],
)

test_case = LLMTestCase(
    input="The dog chased the cat up the tree, who ran up the tree?",
```

(continues on next page)

(continued from previous page)

```

    actual_output="It depends, some might consider the cat, while others might argue the dog.",
    expected_output="The cat."
)

evaluate([test_case], [correctness_metric])

```

- Metrics Summary
 - Correctness (GEval) (score: 0.2, threshold: 0.5, strict: False, evaluation model: local model, reason: Actual output omits the expected detail (the cat) and contradicts the expected output., error: None)
- For test case:
 - input: The dog chased the cat up the tree, who ran up the tree?
 - actual output: It depends, some might consider the cat, while others might argue the dog.
 - expected output: The cat.
 - context: None
 - retrieval context: None
- Overall Metric Pass Rates

Correctness (GEval): 0.00% pass rate

```

EvaluationResult(test_results=[TestResult(name='test_case_0',
success=False, metrics_data=[MetricData(name='Correctness (GEval)',
threshold=0.5, success=False, score=0.2, reason='Actual output omits the expected detail (the cat) and contradicts the expected output.', strict_mode=False, evaluation_model='local model', error=None,
evaluation_cost=0.0, verbose_logs='Criteria:\nDetermine whether the actual output is factually correct based on the expected output.\n\nEvaluation Steps:\n[\n    "Check whether the facts in \'actual output\' contradicts any facts in \'expected output\'",\n    "You should also heavily penalize omission of detail",\n    "Vague language, or contradicting OPINIONS, are OK"\n]\n'), conversational=False, multimodal=False, input='The dog chased the cat up the tree, who ran up the tree?', actual_output='It depends, some might consider the cat, while others might argue the dog.', expected_output='The cat.', context=None, retrieval_context=None)], confident_link=None)

```

**CHAPTER
NINE**

MAIN REFERENCE

BIBLIOGRAPHY

- [GenAI] Wenqiang Feng, Di Zhen. [GenAI: Best Practices](#), 2024.
- [PySpark] Wenqiang Feng. [Learning Apache Spark with Python](#), 2017.
- [lateChunking] Michael Gunther etc. [Late Chunking: Contextual Chunk Embeddings Using Long-Context Embedding Models](#), 2024.
- [PEFT] Yunho Mo etc. [Parameter-Efficient Fine-Tuning Method for Task-Oriented Dialogue Systems](#), 2023.
- [fineTuneEmbedding] Philipp Schmid. [Fine-tune Embedding models for Retrieval Augmented Generation \(RAG\)](#), 2024.
- [fineTuneLLM] Maxime Labonne. [Fine-Tune Your Own Llama 2 Model in a Colab Notebook](#), 2024.
- [GEval] Yang Liu. [G-EVAL: NLG Evaluation using GPT-4 with Better Human Alignment](#), 2023.