Byte Pair Encoding (BPE) is a **subword tokenization algorithm** used in natural language processing, particularly with large language models. It works by iteratively merging the most frequent pair of bytes (or characters) in a text into a single token, effectively breaking down words into smaller, meaningful units. This allows the model to handle rare or unknown words by

BPE starts from character-level representations and progressively merges frequently co-occurring patterns to build a vocabulary of meaningful subword tokens. This constructed vocabulary and the tokenization rules are then used as input for large language models, enabling efficient training and inference. representing them as combinations of these learned subword units

It works by iteratively merging the most frequent pair of adjacent characters (or subwords) in a text into a single, new token. This effectively breaks down words into smaller, more meaningful units.

Why do we need?

Managing Vocabulary Size:

Addressing Rare Word / Out-Of-Vocabulary (OOV) Traditional word-level tokenization struggles with words not present in its training vocabulary (e.g., new words, typos, compound nouns). These are often mapped to an 'unknown' token (`[UNK]`), leading to information loss.

\* BPE tokenizes words into subwords. For instance, "unfriendable" might be broken down into `['un', 'friend', 'able']`. Since `un`, `friend`, and `able` are common subwords, the model can still understand and process the word effectively, even if "unfriendable" itself wasn't in the original vocabulary.

There's an ever-growing number of words in any language. Including every single word in a vocabulary would result in an extremely large vocabulary, making model training and inference computationally expensive and inefficient. \* BPE helps manage vocabulary size by creating a reasonably sized set of subword tokens that can represent most words, balancing

better generalization and a more compact representation.

coverage and efficiency.

Preserving Meaning and Compression Words like "running," "runner," and "runs" share a common root: "run." BPE can identify and merge such common morphological units (e.g., "run") into a single subword token. This helps the model learn and exploit morphological similarities between words, leading to

1. Initial Vocabulary: The process starts with a vocabulary consisting of individual characters present in the training data.

{'l', 'o', 'w', 'e', 'r', 'n', 's', 't'}

extract all unique individual naracters present in the given text (training data) to form the initial vocabulary (token set). While BPE can start at the byte level ("Byte Pair Encoding"), it's often explained starting at the character level. We will proceed at the character level

def initialize\_vocabulary(text: str): Extracts all unique individual characters from the given text to form the initial vocabulary.

text (str): The text data to be used for training.

Returns: tuple: A tuple containing: - set: A set of unique characters present in the text (initial vocabulary). - list: A list of lists, where each inner list represents a word tokenized into characters.

# Split the text into words based on spaces words = text.split()

# Collect all unique characters to create the initial vocabulary initial\_vocab = set() for char in text: initial\_vocab.add(char)

# Tokenize each word into individual characters (needed for subsequent merging steps) tokenized\_words = [[char for char in word] for word in words]

print(f"Original Text: '{text}'") print(f"Initial Vocabulary (Unique Characters): {sorted(list(initial\_vocab))}") print(f"Initial Tokenized Words: {tokenized\_words}")

return initial\_vocab, tokenized\_words # Example Usage:

sample\_text = "low low low lower newest" initial\_vocab, tokenized\_words\_list = initialize\_vocabulary(sample\_text)

print("\n--- Additional Explanation ---") print("`initial\_vocab` represents all basic tokens available at the start of BPE.") print("`tokenized\_words\_list` stores each word broken down into individual characters,") print("which will be used in subsequent steps to find and merge the most frequent pairs.") 2. Frequency Counting:
The algorithm then counts the frequency of adjacent character pairs in the training data.

tokenized\_words\_list = [['l', 'o', 'w'], ['l', 'o', 'w'], ['l', 'o', 'w'], ['l', 'o', 'w', 'e', 'r'], ['n', 'e', 'w', 'e', 's', 't']]

from collections import defaultdict

def count\_pair\_frequencies(tokenized\_words\_list: list[list[str]]): Counts the frequency of adjacent token pairs within the tokenized words list.

tokenized\_words\_list (list[list[str]]): A list of lists, where each inner list represents a word tokenized into characters

dict: A dictionary where keys are (token1, token2) tuples and values are their frequencies.

pair\_frequencies = defaultdict(int) for word\_tokens in tokenized\_words\_list: # Iterate through adjacent pairs in each word # For a word like ['I', 'o', 'w'], pairs are ('I', 'o') and ('o', 'w') for i in range(len(word\_tokens) - 1):

pair = (word\_tokens[i], word\_tokens[i+1]) pair\_frequencies[pair] += 1 print("\n--- Step 2: Frequency Counting ---") print(f"Input tokenized words: {tokenized\_words\_list}")

print("Calculated Pair Frequencies:") # Sort for consistent output, though not strictly necessary for the algorithm sorted\_frequencies = sorted(pair\_frequencies.items(), key=lambda item: item[1], reverse=True) for pair, freq in sorted\_frequencies: print(f" {pair}: {freq}")

return pair\_frequencies

# Continue from the previous step's output: # initial\_vocab, tokenized\_words\_list = initialize\_vocabulary(sample\_text) # This line would be run first

# Example Usage: # Assuming tokenized\_words\_list is from the previous step: # sample\_text = "low low low lower newest"

# initial\_vocab, tokenized\_words\_list = initialize\_vocabulary(sample\_text) # Run this if not already run

pair\_counts = count\_pair\_frequencies(tokenized\_words\_list)

print("\n--- Additional Explanation ---") print("This function identifies the most common adjacent sequences of tokens.") print("The pair with the highest frequency will be chosen for merging in the next step.") 3. Merging:
The most frequent pair is merged into a new token, and the vocabulary is updated with this new token.

identify the most frequent adjacent token pair from the frequencies calculated in the previous step. This most frequent pair is then merged into a single new token. The vocabulary is subsequently updated to include this new token. This merging process is the core of BPE, as it systematically generates new subword tokens

from collections import defaultdict # Re-including previous functions for context and executability def initialize\_vocabulary(text: str): Extracts all unique individual characters from the given text to form the initial vocabulary. words = text.split()

initial\_vocab = set()

initial\_vocab.add(char)

print(f"Original Text: '{text}'")

return initial\_vocab, tokenized\_words

pair\_frequencies = defaultdict(int)

for word\_tokens in tokenized\_words\_list:

pair = (word\_tokens[i], word\_tokens[i+1])

for i in range(len(word\_tokens) - 1):

pair\_frequencies[pair] += 1

print("Calculated Pair Frequencies:")

for pair, freq in sorted\_frequencies:

# --- New Code for Step 3: Merging ---

merged\_tokenized\_words = []

new\_word\_tokens = []

while i < len(word\_tokens):

return merged\_tokenized\_words

tuple: A tuple containing:

if not pair\_frequencies:

# 2. Find the most frequent pair

print(f"\n--- Step 3: Merging ---")

set[str]: The updated vocabulary.

# 1. Count frequencies of all adjacent pairs

- tuple: The pair that was merged in this step.

- str: The new token created by the merge.

# with the highest value (frequency) in the dictionary.

new\_token\_str = "".join(most\_frequent\_pair)

print(f"New token created: '{new\_token\_str}'")

# 4. Update the vocabulary by adding the new token

updated\_vocab.add(new\_token\_str)

# --- Example Usage (connecting all steps) ---

# Store current state for subsequent iterations

# Step 3: Perform a single BPE merge step

print("\n--- Additional Explanation ---")

current\_tokenized\_words = tokenized\_words\_list

sample\_text = "low low low lower newest"

current\_vocab = initial\_vocab

# 3. Perform the merge operation across all tokenized words

Returns:

first, second = pair\_to\_merge

for word\_tokens in tokenized\_words\_list:

new\_word\_tokens.append(new\_token)

new\_word\_tokens.append(word\_tokens[i])

merged\_tokenized\_words.append(new\_word\_tokens)

i += 2 # Skip both merged tokens

print(f" {pair}: {freq}")

return pair\_frequencies

tokenized\_words = [[char for char in word] for word in words]

def count\_pair\_frequencies(tokenized\_words\_list: list[list[str]]):

print(f"Initial Tokenized Words: {tokenized\_words}")

print(f"Initial Vocabulary (Unique Characters): {sorted(list(initial\_vocab))}")

Counts the frequency of adjacent token pairs within the tokenized words list.

Merges the specified pair into a new token within the tokenized words list.

pair\_to\_merge (tuple[str, str]): The (token1, token2) pair to be merged.

list[list[str]]: The updated tokenized words list after merging.

# Check if the current and next token form the pair to merge

tokenized\_words\_list (list[list[str]]): The current list of words tokenized into subwords.

new\_token (str): The string representation of the new merged token (e.g., "lo" from ("l", "o")).

if i + 1 < len(word\_tokens) and word\_tokens[i] == first and word\_tokens[i+1] == second:

def perform\_bpe\_merge\_step(current\_tokenized\_words: list[list[str]], current\_vocab: set[str]):

Performs one step of BPE: finds the most frequent pair, merges it, and updates the vocabulary.

current\_tokenized\_words (list[list[str]]): The current list of words tokenized into subwords.

print("No pairs found for merging. BPE process might be complete or input is empty.")

print(f"Most frequent pair found: {most\_frequent\_pair} with frequency {pair\_frequencies[most\_frequent\_pair]}")

updated\_tokenized\_words = merge\_pair(current\_tokenized\_words, most\_frequent\_pair, new\_token\_str)

updated\_vocab = current\_vocab.copy() # Create a copy to avoid modifying the original set directly

# The `max` function with `key=pair\_frequencies.get` efficiently finds the key (pair)

current\_vocab (set[str]): The current set of available tokens (vocabulary).

list[list[str]]: The updated tokenized words list after merging.

pair\_frequencies = count\_pair\_frequencies(current\_tokenized\_words)

return current\_tokenized\_words, current\_vocab, None, None

most\_frequent\_pair = max(pair\_frequencies, key=pair\_frequencies.get)

# Construct the new token from the merged pair by joining its elements

print(f"Updated Tokenized Words after merge: {updated\_tokenized\_words}")

return updated\_tokenized\_words, updated\_vocab, most\_frequent\_pair, new\_token\_str

# This function internally calls count\_pair\_frequencies (from Step 2) and merge\_pair

print("In a complete BPE implementation, this `perform\_bpe\_merge\_step` function")

print("would be repeated for a fixed number of iterations or until a desired vocabulary size is reached.")

perform\_bpe\_merge\_step(current\_tokenized\_words, current\_vocab)

print("This step demonstrates one iteration of the BPE algorithm.")

updated\_words\_after\_merge, updated\_vocab\_after\_merge, merged\_pair, new\_token = \

print(f"Updated Vocabulary after merge: {sorted(list(updated\_vocab))}")

# Step 1: Initialize the vocabulary and tokenize words into characters

initial\_vocab, tokenized\_words\_list = initialize\_vocabulary(sample\_text)

for char in text:

4. Iteration: Steps 2 and 3 are repeated until a desired vocabulary size or number of iterations is reached.

When processing new text, the same merging rules learned during training are applied to break down words into the learned subword

Tokenization follows the training process closely, in the sense that new inputs are tokenized by applying the following steps: 1 Normalization 2 Pre-tokenization 3 Splitting the words into individual characters

4 Applying the merge rules learned in order on those print("\n--- Step 2: Frequency Counting (called internally by Step 3) ---") # Sort for consistent output, though not strictly necessary for the algorithm sorted\_frequencies = sorted(pair\_frequencies.items(), key=lambda item: item[1], reverse=True) def merge\_pair(tokenized\_words\_list: list[list[str]], pair\_to\_merge: tuple[str, str], new\_token: str):