

# 关于Tokenizer你所需要知道的一切

初学者请着重看1,2,3,6章节-----Mingming Zhang

## 1. 引言：连接人类认知与机器智能的桥梁

在当今人工智能的宏大叙事中，大语言模型（Large Language Models, LLMs）如GPT-4、Claude 3和Llama 3等，以其惊人的生成能力和逻辑推理能力占据了舞台的中心。然而，在这些参数量高达数千亿的神经网络开始处理任何信息之前，必须先经过一道至关重要的、却往往被忽视的工序——分词（Tokenization）。分词器（Tokenizer）是连接人类自然语言与机器二进制世界的第一道关卡，它将连续的文本流转化为离散的数字序列（Token IDs），使得计算机能够进行计算、统计和理解。

尽管在表面上，分词似乎只是一个简单的字符串预处理步骤，但其背后的数学原理、算法选择以及工程实现细节，对最终模型的性能、推理成本、多语言支持能力甚至安全性都有着深远的影响 [1]。分词策略的微小差异，可能导致模型在处理Python代码时缩进混乱，在进行数学运算时位数错误，或者对某些非拉丁语系语言表现出极高的“Token通胀率” [3]。

本报告将对Tokenizer进行穷尽式的剖析。我们将从最基础的文本表示概念讲起，深入探讨Byte-Pair Encoding (BPE)、WordPiece和Unigram三种主流算法的数学原理与异同；我们将以一份标准的BPE代码实现为蓝本，逐行解析其Python代码逻辑，揭示训练（Training）与编码（Encoding）的微观运作机制；我们还将深入对比GPT-2与GPT-4分词器的演进，特别是其正则表达式（Regex）预分词模式的精妙设计与字节级处理的工程智慧。

## 2. 文本表示的演变：从字符到子词的必然选择

在深入代码和具体算法之前，我们必须理解为什么自然语言处理（NLP）领域最终汇聚于“子词（Subword）”分词这一路线上。这是计算效率、词汇覆盖率与语义表达能力之间漫长博弈的结果。

### 2.1 词级别（Word-level）分词的局限

早期的NLP系统，如基于统计机器翻译（SMT）的模型，往往采用词级别分词。这种方法最符合人类的直觉：以空格或标点符号为界，将句子切分为单词。

- **机制：**将 "I love AI." 切分为 ["I", "love", "AI", "."]。
- **优点：**保留了词的语义完整性，每个Token都对应一个具有明确定义的人类语言概念。
- **致命缺陷——词表爆炸（Vocabulary Explosion）：**人类语言具有极其丰富的形态变化。以英语为例，动词 "run" 有 "runs", "running", "ran" 等多种变形；在形态丰富的语言（如土耳其语、芬兰语）中，一个词根通过黏着语素可以生成成千上万种变体。如果将每种变形都视为独立Token，模型需要维护一个数百万量级的巨大词表，这在计算资源上是不可接受的。
- **未登录词（OOV）问题：**无论词表多大，总会有未见过的词（如新造词 "uninstagrammable" 或人名）。在词级别模型中，这些词只能被映射为通用的 <UNK> (Unknown) 标记，导致该位置的信息完全丢失，这对翻译或理解任务是灾难性的 [2]。

### 2.2 字符级别（Character-level）分词的尝试

为了彻底解决OOV问题，研究者一度转向字符级别分词，将文本拆解为单个字符。

- **机制：**将 "love" 切分为 ["l", "o", "v", "e"]。

- 优点：词表极小（仅需包含字母表、数字和符号，通常在100-1000量级），理论上消灭了OOV问题。
- 缺陷——序列过长与语义稀疏：
  1. 计算成本：Transformer模型的注意力机制（Self-Attention）的时间复杂度与序列长度的平方成正比（ $O(N^2)$ ）。字符级分词会导致序列长度增加5-10倍，使得训练和推理成本指数级上升。
  2. 语义缺失：单个字符（如 "t"）通常不承载独立的语义信息。模型需要花费大量的层数和参数去组合字符以识别出“词”的概念，这浪费了模型的表达能力 [2]。

## 2.3 子词（Subword）分词的崛起

子词分词（Subword Tokenization）是上述两种方法的辩证统一，也是目前大语言模型的标准配置。其核心哲学是：常用词保持完整，罕见词拆分为有意义的子部件（Subword units）。

例如，单词“tokenization”在子词分词器中可能被拆分为“token”和“ization”。

- “token”是高频词根，作为一个整体被保留，模型可以直接获取其语义嵌入。
- “ization”是常见的高频后缀，作为一个整体被保留。
- 这种机制使得模型既能高效处理常见词，又能通过组合词根和词缀来泛化理解未见过的复合词（如“modernization”，“optimization”），从而在有限的词表大小下实现了无限的词汇表达能力 [4]。

目前主流的子词算法三巨头包括：

1. **BPE (Byte-Pair Encoding)**：基于频率的合并策略，广泛用于GPT系列、Llama、RoBERTa [7]。
2. **WordPiece**：基于概率（似然度）的合并策略，起源于BERT [6]。
3. **Unigram**：基于概率的剪枝策略（从大词表删减），用于SentencePiece (ALBERT, T5) [10]。

## 3. 深度解析：Byte-Pair Encoding (BPE) 算法与代码实现

Byte-Pair Encoding（字节对编码）最初是作为一种数据压缩算法由Philip Gage在1994年提出的 [12]。Sennrich等人于2015年将其引入NLP领域，用于解决神经机器翻译中的稀有词问题。如今，它已成为GPT系列模型的基石。

为了彻底理解BPE，我们将不仅仅停留在理论层面，而是通过逐行解析代码的方式，通过一个参考实现（基于Karpathy的minbpe逻辑）来剖析其内部运作。

### 3.1 算法核心逻辑与代码结构

BPE的训练过程本质上是一个迭代的数据压缩过程。

1. 初始化：将所有文本拆解为基础单元（通常是字节）。
2. 统计：统计所有相邻单元对（Pair）在数据中出现的频率。
3. 合并：找到频率最高的那个对（例如 ('e', 's')），将其合并为一个新的符号（'es'），并分配一个新的ID。
4. 迭代：重复步骤2和3，直到达到预设的词表大小（Vocabulary Size） [13]。

我们将代码分为三个核心功能模块进行讲解：统计（Statistics）、合并（Merge）和训练主循环（Training Loop）。

### 3.2 代码详解：统计频率（get\_stats）

这是BPE算法中最基础的原子操作：扫描当前的Token序列，统计每一对相邻Token出现的次数。

```

def get_stats(ids):
    """
    输入:
    ids (list of integers): 当前的Token ID列表。
    输出:
    counts (dict): 映射 (id1, id2) -> frequency
    """
    counts = {}
    # zip(ids, ids[1:]) 是一个Python技巧
    # zip后生成: [(1, 2), (2, 3), (3, 4)]
    # 这正是我们需要统计的所有相邻对
    for pair in zip(ids, ids[1:]):
        counts[pair] = counts.get(pair, 0) + 1
    return counts

```

深入解析与知识点：

- **输入数据的本质：**这里的ids最初是什么？在现代LLM（如GPT-4）中，我们使用**字节级（Byte-level）BPE**。这意味着ids的初始状态是文本UTF-8编码后的字节序列，范围在0-255之间。例如，英文字母a是97，中文你是三个字节。这种设计极其重要，因为它保证了Tokenizer可以处理任何Unicode字符串，哪怕是Emoji或从未见过的外文，因为一切皆为字节[7]。
- **时间复杂度：**此函数的时间复杂度为 $O(N)$ ，其中 $N$ 是序列长度。在训练过程中，每次合并都需要重新扫描整个语料库，这使得朴素实现的BPE训练非常慢( $O(NcdotV)$ ，其中 $V$ 是合并次数)。工业级实现（如Rust编写的Tiktoken）会使用链表或优先队列来优化更新过程，避免全量扫描[15]。

### 3.3 代码详解：执行合并(merge)

一旦我们找到了出现频率最高的对（比如(101, 115)对应('e', 's')），我们就需要将序列中所有的(101, 115)替换为新的Token ID（比如256）。

```

def merge(ids, pair, idx):
    """
    输入:
    ids (list): 当前的Token列表
    pair (tuple): 要合并的一对Token，例如 (101, 115)
    idx (int): 分配给这个新Token的ID，例如 256

    输出:
    newids (list): 合并后的新列表
    """
    newids =
    i = 0
    while i < len(ids):
        # 检查是否刚好碰到了我们要合并的pair，且没有越界
        # 注意: i < len(ids) - 1 是为了防止检查 ids[i+1] 时越界
        if i < len(ids) - 1 and ids[i] == pair and ids[i+1] == pair:
            newids.append(idx) # 替换为新的ID
            i += 2 # 跳过接下来两个元素，因为它们已经被合并了
        else:

```

```

        newids.append(ids[i]) # 保持原样
        i += 1
    return newids

```

深入解析与知识点：

- **贪婪算法 (Greedy Strategy)**：BPE是贪婪的。只要选定了频率最高的Pair，它就会在全局范围内将所有的实例都替换掉。这不同于某些动态规划算法。
- **序列变短**：每次合并操作都会使 ids 列表的长度变短。这正是“压缩”的体现。对于大模型，更短的序列意味着能放入更多的上下文信息。
- **新ID的分配**：基础词表是0-255。第一次合并产生的Token ID是256，第二次是257，依此类推。GPT-4的cl100k\_base 词表大小约为100,277，这意味着这个合并过程在训练阶段重复了约10万次 [17]。

## 3.4 代码详解：训练主循环 (train)

将上述两个函数结合，就构成了BPE的训练器。

```

def train(text, vocab_size, verbose=False):
    """
    输入:
    text (str): 训练语料文本
    vocab_size (int): 目标词表大小 (例如 50257)
    """

    assert vocab_size >= 256
    num_merges = vocab_size - 256 # 需要进行的合并次数

    # 1. 预处理: 将文本转换为UTF-8字节流
    text_bytes = text.encode("utf-8")
    ids = list(text_bytes) # 初始列表, 元素范围 0-255

    merges = {} # 记录合并规则: (p0, p1) -> idx

    print(f"Original length: {len(ids)}")

    for i in range(num_merges):
        # 2. 统计当前频率
        stats = get_stats(ids)
        if not stats:
            break # 如果没有可以合并的对, 提前退出

        # 3. 找到频率最高的对
        # key=stats.get 表示按字典的值(频率)排序
        pair = max(stats, key=stats.get)

        # 4. 分配新ID
        idx = 256 + i

        # 5. 执行合并
        ids = merge(ids, pair, idx)

```

```

# 6. 记录规则
merges[pair] = idx

if verbose:
    print(f"Merge {i+1}/{num_merges}: {pair} -> {idx} (count: {stats[pair]})")

return merges

```

深入解析与知识点：

- **核心产物**：训练结束后，最重要的产物不是 ids，而是 merges 字典。这个字典定义了Tokenizer的“知识”。当我们下载一个预训练模型时，tokenizer.json 或 merges.txt 存储的就是这个字典。
- **确定性与Tie-breaking**：如果两个Pair频率相同怎么办？Python的 max 函数在值相同时会返回第一个遇到的键。为了保证Tokenizer的可复现性（Deterministic），在工业级实现中通常会规定：当频率相同时，按照 Pair中字符的字典序（Lexicographical order）进行选择。
- **词表大小的权衡**：vocab\_size 是一个关键超参数。
  - **太小**：导致序列过长，模型推理慢，无法捕捉长距离依赖。
  - **太大**：导致Embedding矩阵（vocab\_size x hidden\_dim）参数量激增，增加训练负担；且由于稀有词频次太低，其Embedding可能训练不充分（Undertrained） [1]。

## 3.5 代码详解：推理阶段的编码 (encode)

有了训练好的 merges 规则，如何将新的文本转化为Token IDs？这是一个极易出错的环节。初学者常犯的错误是：在推理时依然每次找文本中频率最高的对。这是错误的。推理时，必须严格按照训练时确定的优先级顺序进行合并。

```

def encode(text, merges):
    # 1. 转为字节流
    ids = list(text.encode("utf-8"))

    while len(ids) >= 2:
        # 获取当前文本中所有相邻对
        stats = get_stats(ids)

        # 寻找“在merges规则表中存在，且ID最小（即最早被训练出来）”的对
        # 因为ID越小，说明它在训练集中频率越高，优先级越高
        pair_to_merge = None
        min_rank = float("inf") # rank即ID

        for pair in stats:
            if pair in merges:
                rank = merges[pair]
                if rank < min_rank:
                    min_rank = rank
                    pair_to_merge = pair

        # 如果当前序列中没有任何对在我们的规则表中，停止

```

```

if pair_to_merge is None:
    break

# 执行合并
ids = merge(ids, pair_to_merge, min_rank)

return ids

```

深入解析与知识点：

- **优先级逻辑**: 代码中的 min\_rank 逻辑至关重要。假设我们有规则  $(a, b) \rightarrow X$  (Rank 1) 和  $(b, c) \rightarrow Y$  (Rank 2)。对于输入 abc，我们必须先合并 Rank 1 的规则，得到 Xc。如果我们先合并 bc 得到 aY，就破坏了BPE的构造一致性，导致生成的Token序列与训练时分布不一致。
- **递归与迭代**: 上述代码使用了类似迭代的方式。在Python的 minbpe 实现中，也可以通过递归来实现，或者使用更高效的Regex拆分后独立处理每个块（详见后文GPT-2/4部分）。

## 3.6 代码详解：解码 (decode)

解码是编码的逆过程，相对简单，但有一个关键细节：处理无效字节。

```

def decode(ids, vocab):
    """
    ids: token ID列表
    vocab: 映射 {idx: bytes} (由基础字节和merges反推得到)
    """
    # 将所有ID映射回字节串，并拼接
    tokens = b''.join(vocab[idx] for idx in ids)

    # errors="replace" 是关键
    text = tokens.decode("utf-8", errors="replace")
    return text

```

深入解析与知识点：

- **errors="replace"**: 在GPT-3或GPT-4的输出中，我们偶尔会看到 (Replacement Character)。这是因为大模型有时会生成不完整的Token序列。例如，一个中文字符由3个字节组成，如果模型只输出了前2个字节就停止了（比如因为 max\_tokens 限制），那么这2个字节无法构成合法的UTF-8字符。使用 replace 策略可以防止程序崩溃。[7]。

## 4. 算法三巨头深度对比：BPE, WordPiece 与 Unigram

虽然BPE占据了GPT系列的主导地位，但在BERT、T5等模型中，WordPiece和Unigram算法同样重要。理解它们的区别是深入NLP底层逻辑的关键。

### 4.1 WordPiece：从频率到概率的跨越

WordPiece算法由Google开发，是BERT模型的核心。它的整体流程与BPE非常相似（也是自底向上的合并），但在选择合并哪一对的标准上完全不同。

- **BPE标准**: 选择频次 (**Frequency**) 最高的对。
  - 目标: 最大化数据压缩比。
- **WordPiece标准**: 选择合并后能使训练数据似然度 (**Likelihood**) 增加最多的对 6。
  - 这等价于选择点互信息 (**Pointwise Mutual Information, PMI**) 最高的对。
  - WordPiece 得分公式:

$$textScore(A, B) = \frac{P(AB)}{P(A)P(B)}$$

其中  $P(AB)$  是对  $AB$  出现的概率,  $P(A)$  和  $P(B)$  是各自的概率。

#### 深度洞察: 为什么PMI比频率更优?

WordPiece的评分机制考虑了子词的独立概率。

- 假设 A="the" 和 B="book" 都是极高频词, 那么它们连在一起 thebook 出现的次数可能很高。但在BPE中它们可能会被合并。
- 在WordPiece中, 由于分母  $P("the") \times P("book")$  非常大, Score 会被拉低。这防止了两个本该独立的常见词被意外合并。
- 相反, 假设 A="Z" 和 B="qa" 都很罕见, 但只要出现就总是粘在一起 (如 "Zqa")。此时  $P(AB) \approx P(A) \approx P(B)$ , Score 会非常高 (接近  $1/P(A)$ )。
- 结论: WordPiece倾向于合并那些**内在关联性强** (比随机组合更紧密) 的词对, 而不仅仅是高频词对。这使得WordPiece在处理词缀 (如 "un-", "-ing") 时通常比BPE在语言学上更合理 [19]。

## 4.2 Unigram: 基于概率图模型的自顶向下剪枝

Unigram Tokenization (主要在SentencePiece库中实现) 的思路与前两者截然相反。BPE和WordPiece是**自底向上 (Bottom-up)** 的构造, Unigram是**自顶向下 (Top-down)** 的。

数学原理与EM算法:

Unigram模型假设每个Subword的出现是独立的。一个句子  $X$  被切分为序列  $\mathbf{x} = (x_1, \dots, x_m)$  的概率为:

$$P(\mathbf{x}) = \prod_{i=1}^m P(x_i)$$

其中  $P(x_i)$  是子词  $x_i$  的发生概率。

训练流程 (EM算法) :

1. **初始化**: 构建一个极其巨大的词表 (例如包含语料中所有出现过的子串, 可能有几百万个)。
  2. **E-step (期望步)** : 固定当前词表, 使用**Viterbi算法**计算语料库中每个句子的最优切分路径。
    - 对于单词 "tokenization", 可能有多种切分方式 (["token", "ization"] vs ["t", "o", "ken", "..."])。Viterbi算法能找到使得联合概率  $P(\mathbf{x})$  最大的那条路径 21。
  3. **M-step (最大化步)** : 重新计算每个子词的出现概率  $P(x_i)$ 。
  4. 计算Loss并剪枝: 计算如果从词表中移除某个子词  $x$ , 总似然度  $L$  会下降多少 (Loss)。
- $DeltaL = L_{new} - L_{old}$
5. **剪枝策略**: 移除那些对总似然度贡献最小 (Loss最小) 的Token (通常每轮移除20%)。
  6. **循环**: 重复上述过程直到词表缩小到预定大小。

独特优势——子词正则化 (Subword Regularization) :

Unigram不仅仅是一个分词器，它本身就是一个微型的语言模型。在训练LLM时，我们可以利用Unigram的概率特性进行数据增强。

- 对于同一个词 "New York"，我们不总是输出最优切分 ["New", " York"]。
- 我们可以根据概率采样 (Sampling)，有时将其切分为["N", "ew", " Yo", "rk"]。
- 这种技术迫使模型学习不同切分下的语义，显著提升了模型在处理拼写错误和噪声文本时的鲁棒性 [23]。

## 4.3 算法特性对比总结表

特性	BPE (GPT-2/3/4, Llama)	WordPiece (BERT)	Unigram (T5, ALBERT)
构建方向	自底向上 (Bottom-up)	自底向上 (Bottom-up)	自顶向下 (Top-down)
核心指标	频率 (Frequency)	似然度增益 / PMI	似然度损失 (Loss)
训练复杂度	较低	高 (每步需重算所有Pair得分)	较高 (需EM迭代)
OOV处理	字节回退 (Byte fallback)	字符回退 (需UNK Token)	字符回退
分词确定性	确定性 (Deterministic)	确定性	概率性 (可采样多种结果)
适用场景	生成式模型 (Generative)	理解式模型 (NLU)	需正则化的任务

## 5. GPT分词器的演进：从GPT-2到GPT-4的工程飞跃

OpenAI的GPT系列一直沿用BPE算法，但在细节上进行了极其关键的优化。了解这些演进，能让我们明白为什么GPT-4在写代码和处理非英语语言时比GPT-2强得多。

### 5.1 预分词 (Pre-tokenization) 与Regex的秘密

BPE算法本身是“盲目”的。如果不加干预，它可能会跨越标点和单词的边界进行合并。例如，它可能会把句子结尾的 "dog." 中的 "g." 合并为一个Token。这通常是不希望看到的，因为标点符号通常有独立的语法意义。

因此，在运行BPE合并之前，需要先用正则表达式 (Regex) 将文本切分为一个个基础的“单词块”。BPE只能在这些块内部进行合并，不能跨块。

#### 5.1.1 GPT-2 Regex 模式分析

GPT-2使用的Regex模式如下 (Python re 语法) [25]:

```
r"""\s*(?:[sdmt]|ll|ve|re)|?p{L}+|?p{N}+|?[^sp{L}p{N}]+|s+(?!s)|s+""
```

逐段解析：

- '(?:[sdmt]|ll|ve|re)': 处理缩写。它将 's, 't, 're 等缩写从单词中剥离出来。例如 "don't" 会被拆分为 "don" 和 "'t"。这保证了模型能理解否定、所有格等概念。
- ?p{L}+: 处理单词。匹配前导空格 (可选) 加上一串字母。注意，GPT-2将空格视为单词的一部分 (通常在开头)，如 token。

3. ?p{N}+: 处理数字。匹配连续的数字。
4. ?[^sp{L}p{N}]+: 处理标点。匹配非空格、非字母、非数字的字符序列。
5. s+(?!S): 处理尾部空格。

GPT-2的缺陷：

这个正则在处理多空格时表现不佳，且对代码中的缩进（通常是大量空格）处理效率极低。更重要的是，它对大小写的处理不够完美。例如，它能识别's为缩写，但如果用户输入大写的 HOW'S，由于正则中没有忽略大小写的标志，'S可能不会被单独切分，从而导致分词不一致 27。

### 5.1.2 GPT-4 (cl100k\_base) Regex 模式分析

GPT-4（以及Tiktoken库）使用了更复杂的Regex [18]：

```
r"""(?i:'s|'t|'re|'ve|'m|'ll|'d)|[^rnp{L}p{N}]?p{L}+|p{N}{2,}|[^rnp{L}p{N}]?[^sp{L}p{N}]++[rn]*|s*[rn]+|s+(?!S)|s+"""
```

重大改进与洞察：

1. **大小写不敏感 (?i:...)**: 解决了GPT-2中 Don't 和 DON'T 切分不一致的问题。这看似微小，实则极大提升了模型对大写输入的理解能力。
2. **数字合并策略 p{N}{2,}**: GPT-4要求数字匹配至少两位。这倾向于将数字进行更激进的切分，或者限制数字合并的长度。这有助于提升数学计算能力，防止非常长的数字串被合成一个极其稀有的Token，导致模型无法理解其数值含义。
3. **空格与代码优化**: GPT-4的正则允许合并更多的连续空格。这对于**Python**代码（依赖缩进）极其重要。GPT-2往往把4个空格切成4个Token，浪费了宝贵的上下文窗口（Context Window），而GPT-4可以将其压缩为一个Token，大幅提升了处理长代码的能力 3。

## 5.2 字节级BPE的精妙实现 (bytes\_to\_unicode)

在GPT-2的原始代码中，有一个神秘的函数 bytes\_to\_unicode [29]。很多开发者在初次阅读时会感到困惑。

问题背景：

BPE运行在字节（0-255）上。但是，有些字节是控制字符（Control Characters），例如换行符、回车符、或者不可见的二进制数据。如果直接打印这些Token，会导致终端乱码或显示为空白，这对于调试非常不友好。

OpenAI的解决方案：

他们建立了一个可逆的映射表，将256个字节映射到256个可打印的Unicode字符上。

- ASCII的可打印字符（如 a, B, %）保持不变。
- 不可见字符（如字节0-32）被映射到Unicode中较高位的字符（如 U+0120, U+010A）。
- 这就是为什么我们在查看GPT Tokenizer输出时，经常会看到 world。这里的空格实际上就是空格（Space, 字节32）。
- **代码意义**：这个函数保证了分词过程中的所有中间状态都是人类可读的字符串，同时没有丢失任何底层字节信息。

## 5.3 词表大小的演进逻辑

模型	词表大小	含义	影响分析
GPT-2	50,257	较小	英语为主，其他语言效率低。
GPT-4	100,277	翻倍	显著提升多语言压缩率，代码效率提升。
Llama 3	128,000	更大	进一步优化多语言支持。

为什么词表越来越大？

更大的词表意味着同一个句子被切分成更少的Token。

- 优点：
  - 推理更快：生成的步数变少了。
  - 上下文更“大”：同样的Token限制下能塞进更多的实际文本。
  - 多语言公平性：大词表允许容纳更多非英语语言的常用词（如中文单字），减少非英语文本的“Token膨胀”率 1。
- 代价：Embedding层参数量剧增（128k \* 4096 维度的矩阵非常巨大），训练收敛更难，需要更多的数据来填满这些稀疏的Token嵌入。

## 6. 特殊Token (Special Tokens) 的工程处理机制

在代码实现中，特殊Token（如 <| endoftext |>, <PAD>, <MASK>, <| im\_start |>）的处理往往是初学者最容易混淆的地方，也是导致模型产生幻觉或安全漏洞的常见原因。

### 6.1 为什么不能用普通BPE处理特殊Token？

假设我们的特殊Token是 <| endoftext |>，用于标记文本结束。如果我们把它当作普通文本传给BPE算法：

- BPE会先将其按字符拆分。
- 然后按照合并规则，它可能会被切分为 ['<', '|', 'endo', 'ft', 'ext', '|', '>']。
- 后果：模型将无法把这个序列识别为一个统一的“停止信号”，而是一堆无意义的碎片。模型可能无法正确停止生成。

### 6.2 解决方案：Tiktoken vs HuggingFace

Tiktoken (GPT-4) 的处理方式：

Tiktoken 要求用户显式传递 allowed\_special 参数。这是为了安全。

- Prompt注入防御：如果不允许特殊Token，当用户输入 "Hello <| endoftext |>" 时，Tiktoken会强制将其作为普通文本进行分词（即切碎），防止用户伪造系统指令。
- 实现机制：Tiktoken 内部维护了一个独立的字典来存储特殊Token。在分词开始前，它会先用正则把这些特殊字符串“抠”出来，不参与BPE合并，直接赋予特定的ID [13]。

HuggingFace Tokenizers 的处理方式：

HF引入了 AddedToken 对象，并区分 special\_tokens（具有特殊语义，如EOS）和 additional\_special\_tokens。HF的分词流程中包含一个 Normalization 步骤，但在处理特殊Token时，它同样会优先保护这些Token不被切分 [31]。

代码实战：如何在自制Tokenizer中处理？

如果你在使用 minBPE 或自己写 tokenizer，必须在 encode 函数的最开始加入对特殊 Token 的正则匹配。

```
# 伪代码逻辑：处理特殊Token
def encode(text, special_tokens):
    # special_tokens 是一个字典 {'<|endoftext|>': 100257}

    # 1. 创建一个正则模式，匹配所有特殊Token
    # 例如 pattern = "<|endoftext|>"
    special_pattern = create_pattern(special_tokens.keys())

    # 2. 将文本切分为“普通部分”和“特殊部分”
    # text: "Hello <|endoftext|> world"
    # splits: ["Hello ", "<|endoftext|>", " world"]
    splits = re.split(special_pattern, text)

    final_ids =
        for part in splits:
            if part in special_tokens:
                # 如果是特殊Token，直接追加ID
                final_ids.append(special_tokens[part])
            else:
                # 如果是普通文本，应用BPE算法
                final_ids.extend(bpe_encode(part))

    return final_ids
```

这是保证模型控制流正确的唯一方法 [7]。

## 7. 分词对模型性能的深远影响

Tokenizer 不仅仅是数据的搬运工，它实际上重塑了模型眼中的世界。许多LLM的怪异行为都可以追溯到分词阶段。

### 7.1 算术与数字的“盲区”

LLM通常在算术任务上表现挣扎，部分原因在于分词。

- **不一致性**: 1000 可能是一个Token。1001 可能被切分为["100", "1"]。1002 可能是["10", "02"]。
- **位值丢失**: 由于数字被切分得支离破碎且不规律，模型很难学习到统一的“位值（Place Value）”规则（即个位、十位、百位的关系）。
- **GPT-4的改进**: 通过正则限制数字的合并，尽量保持数字切分的一致性，但这依然是基于文本模型的硬伤 [3]。

### 7.2 编程语言的缩进噩梦

在Python代码中，缩进即逻辑。

- 在GPT-2 Tokenizer中，4个空格通常被切分为4个 G Token。这意味着一段深度缩进的代码会消耗大量的 Token配额，且模型必须精确数出有多少个 G 才能确定代码块层级。

- GPT-4通过将连续空格合并，显著缓解了这个问题，使得模型在生成代码时逻辑更严密，上下文利用率更高。

## 7.3 “Glitch Tokens” (故障Token)

研究人员发现，某些Token（如SolidGoldMagikarp, Dragonbound）会导致模型生成乱码或崩溃。

- **原因：**这些词通常来自Reddit等网络论坛的用户名，在爬取数据时被BPE统计为高频词并加入了词表。
- **灾难：**然而，在后续的训练数据清洗中，这些词可能被当作“噪声”过滤掉了。
- **结果：**导致这些Token存在于词表中，但在Embedding层的训练中从未被更新过（处于初始随机状态）。当模型在推理时偶尔遇到这些词，就会激活一个未经训练的随机向量，导致输出崩坏。这提醒我们在构建Tokenizer时，数据的预处理和后处理必须严格对齐 28。

# 8. 性能与生态：Tiktoken, SentencePiece与HuggingFace实战指南

在工程实践中，我们很少直接运行纯Python版的BPE，因为效率太低。我们通常使用高度优化的库。

## 8.1 核心库对比

库	核心语言	支持算法	特点	适用模型
<b>Tiktoken</b> (OpenAI)	Rust	BPE	极速（比HF快3-6倍）。专为OpenAI模型优化。API简单，仅支持推理，不支持训练新模型（通常）。	GPT-3.5, GPT-4
<b>SentencePiece</b> (Google)	C++	BPE, Unigram	无损处理（Lossless）。将空格视为特殊字符 _，不需要预分词正则。完全可逆。对多语言支持极好。	Llama, ALBERT, T5
<b>Tokenizers</b> (HuggingFace)	Rust	All	大一统。功能最全，支持训练和推理。集成了上述两者的优点，但API相对复杂。	BERT, RoBERTa, Mistral

## 8.2 SentencePiece 的“空格”哲学

传统的Tokenizer（如BERT）先按空格分词，这就丢失了“这里原来有几个空格”的信息。SentencePiece将空格视作为一个普通字符（用下划线 \_ 或 (U+2581) 表示）。

- 输入：Hello World (两个空格)
- SP分词：["\_Hello", "\_", "World"] (保留了所有信息)

这就是为什么Llama等模型可以直接处理原始文本，而不需要复杂的预处理规则。这种“Raw stream in, Token stream out”的设计理念是目前开源大模型的主流 [10]。

# 9. 结论与未来展望

Tokenization 是自然语言处理中一个古老而又充满活力的领域。从早期的空格切分，到统计学的BPE，再到如今融合了语言学规则与概率图模型的复杂系统，Tokenizer的每一次进化都推动了模型性能的边界。

未来的方向在哪里？

1. **Token-free Models (Byte-level Transformers)**: 研究者正在探索直接在字节层面 (MegaByte, MambaByte) 进行训练。虽然序列长度增加了4倍，但随着Flash Attention等线性注意力机制的发展，这正在成为可能。这将彻底消灭Tokenizer带来的所有偏见和问题。
2. **多模态融合**: 随着GPT-4o等模型的出现，Tokenizer不仅要处理文本，还要处理图像Patch和音频帧。未来的Tokenizer将是“万物皆Token”的统一体。

通过理解Tokenizer的每一个字节、每一行代码、每一个正则符号，我们不仅是在学习一个预处理工具，更是在窥探大语言模型认知世界的“第一眼”。这一眼，决定了它能看多远。

---

注：本报告中涉及的代码逻辑主要基于karpathy/minbpe及OpenAI tiktoken的公开实现原理。引用来源包括但不限于arXiv论文、官方文档及技术博客。

## 引用的著作

1. Tokenization Disparities as Infrastructure Bias: How Subword Systems Create Inequities in LLM Access and Efficiency - arXiv, 访问时间为一月 21, 2026, <https://arxiv.org/html/2510.12389v1>
2. The Role of Tokenization in LLMs: Does It Matter? - DZone, 访问时间为一月 21, 2026, <https://dzone.com/articles/the-role-of-tokenization-in-langs-does-it-matter>
3. LLM Tokenization - Hundred Blocks, 访问时间为一月 21, 2026, [https://hundredblocks.github.io/transcription\\_demo/](https://hundredblocks.github.io/transcription_demo/)
4. Byte Pair Encoding vs. Unigram Tokenization: A Deep Dive into Subword Models - Medium, 访问时间为一月 21, 2026, <https://medium.com/@xiamaxi/byte-pair-encoding-vs-unigram-tokenization-a-deep-dive-into-subword-models-4963246e9a34>
5. A Comparative Analysis of Byte-Level and Token-Level Transformer Models in Natural Language Processing - Greg Robison, 访问时间为一月 21, 2026, <https://gregrobison.medium.com/a-comparative-analysis-of-byte-level-and-token-level-transformer-models-in-natural-language-9fb4331b6acc>
6. Subword Tokenization: BPE, WordPiece, and Unigram - Kaggle, 访问时间为一月 21, 2026, <https://www.kaggle.com/code/danishmahdi/subword-tokenization-bpe-wordpiece-and-unigram>
7. karpathy/minbpe: Minimal, clean code for the Byte Pair Encoding (BPE) algorithm commonly used in LLM tokenization. - GitHub, 访问时间为一月 21, 2026, <https://github.com/karpathy/minbpe>
8. Byte-Pair Encoding tokenization - Hugging Face LLM Course, 访问时间为一月 21, 2026, <https://huggingface.co/learn/llm-course/chapter6/5>
9. 访问时间为一月 21, 2026, [https://arxiv.org/html/2411.17669v1#:~:text=WordPiece%20is%20similar%20to%20BPE,to%20obtain%20a%20smaller%20vocabulary\\_](https://arxiv.org/html/2411.17669v1#:~:text=WordPiece%20is%20similar%20to%20BPE,to%20obtain%20a%20smaller%20vocabulary_)
10. Linguistic Laws Meet Protein Sequences: A Comparative Analysis of Subword Tokenization Methods - arXiv, 访问时间为一月 21, 2026, <https://arxiv.org/html/2411.17669v1>
11. Unigram tokenizer: how does it work? - Data Science Stack Exchange, 访问时间为一月 21, 2026, <https://datascience.stackexchange.com/questions/88824/unigram-tokenizer-how-does-it-work>
12. Byte-pair encoding - Wikipedia, 访问时间为一月 21, 2026, [https://en.wikipedia.org/wiki/Byte-pair\\_encoding](https://en.wikipedia.org/wiki/Byte-pair_encoding)

13. Implementing A Byte Pair Encoding (BPE) Tokenizer From Scratch - Sebastian Raschka, 访问时间为一月 21, 2026, <https://sebastianraschka.com/blog/2025/bpe-from-scratch.html>
14. Let's Build the GPT Tokenizer: A Complete Guide to Tokenization in LLMs - Fast.ai, 访问时间为一月 21, 2026, <https://www.fast.ai/posts/2025-10-16-karpathy-tokenizers>
15. [D] SentencePiece, WordPiece, BPE... Which tokenizer is the best one? : r/MachineLearning, 访问时间为一月 21, 2026,  
[https://www.reddit.com/r/MachineLearning/comments/rprm93/d\\_sentencepiece\\_wordpiece\\_bpe\\_which\\_tokenizer\\_is/](https://www.reddit.com/r/MachineLearning/comments/rprm93/d_sentencepiece_wordpiece_bpe_which_tokenizer_is/)
16. Rs-bpe tokenizer [PyPI | Python] - Outperforms tiktoken & tokenizers - Hugging Face Forums, 访问时间为一月 21, 2026, <https://discuss.huggingface.co/t/rs-bpe-tokenizer-pypi-python-outperforms-tiktoken-tokenizers/146410>
17. tiktoken is a fast BPE tokeniser for use with OpenAI's models. - GitHub, 访问时间为一月 21, 2026, <https://github.com/openai/tiktoken>
18. 01. Let's build the GPT Tokenizer v2 - Kaggle, 访问时间为一月 21, 2026, <https://www.kaggle.com/code/shravankumar147/01-let-s-build-the-gpt-tokenizer-v2>
19. Probabilistic Interpretation of WordPiece | Saibo Geng, 访问时间为一月 21, 2026, [https://saibo-creator.github.io/post/2024\\_04\\_03\\_probabilistic\\_interpretation\\_wordpiece/](https://saibo-creator.github.io/post/2024_04_03_probabilistic_interpretation_wordpiece/)
20. Pointwise mutual information - Wikipedia, 访问时间为一月 21, 2026, [https://en.wikipedia.org/wiki/Pointwise\\_mutual\\_information](https://en.wikipedia.org/wiki/Pointwise_mutual_information)
21. Which Pieces Does Unigram Tokenization Really Need? - arXiv, 访问时间为一月 21, 2026, <https://arxiv.org/html/2512.12641>
22. Unigram Tokenization Algorithm - Emergent Mind, 访问时间为一月 21, 2026, <https://www.emergentmind.com/topics/unigram-tokenization-algorithm>
23. Distributional Properties of Subword Regularization - arXiv, 访问时间为一月 21, 2026, <https://arxiv.org/html/2408.11443v1>
24. SentencePiece Tokenizer Demystified - Towards Data Science, 访问时间为一月 21, 2026, <https://towardsdatascience.com/sentencepiece-tokenizer-demystified-d0a3aac19b15/>
25. Words and Tokens - Stanford University, 访问时间为一月 21, 2026, <https://web.stanford.edu/~jurafsky/slp3/2.pdf>
26. Study Notes: Stanford CS336 Language Modeling from Scratch [4], 访问时间为一月 21, 2026, <https://bearbearyu1223.github.io/cs336/2025/08/10/cs336-gpt2-regex-for-pretokenization-explain.html>
27. Let's Build the GPT Tokenizer: A Complete Guide to Tokenization in LLMs - Fast.ai, 访问时间为一月 21, 2026, <https://www.fast.ai/posts/2025-10-16-karpathy-tokenizers.html>
28. Tokenization All You Need. What is the meaning of the tokenizer's... | by Mujeeb Ur Rehman | Medium, 访问时间为一月 21, 2026, <https://medium.com/@mujeeb.merwat/tokenization-all-you-need-88e0fe798cfe>
29. Why u0120 (Ğ) is in so many pairs? · Issue #80 · openai/gpt-2 - GitHub, 访问时间为一月 21, 2026, <https://github.com/openai/gpt-2/issues/80>
30. GPT2 From Scratch In C++ - Part 1 - Tokenization - DaoPlays, 访问时间为一月 21, 2026, [https://www.daoplays.org/blog/gpt2\\_p1](https://www.daoplays.org/blog/gpt2_p1)

31. Module2 Graded Lab1 padding and eos tokens - Fine-tuning & RL for LLMs: Intro to Post-training - DeepLearning.AI Community, 访问时间为 一月 21, 2026,  
<https://community.deeplearning.ai/t/module2-graded-lab1-padding-and-eos-tokens/884684>
32. Tokenizer - Hugging Face, 访问时间为 一月 21, 2026,  
[https://huggingface.co/docs/transformers/main\\_classes/tokenizer](https://huggingface.co/docs/transformers/main_classes/tokenizer)
33. Summary of the tokenizers - Hugging Face, 访问时间为 一月 21, 2026,  
[https://huggingface.co/docs/transformers/en/tokenizer\\_summary](https://huggingface.co/docs/transformers/en/tokenizer_summary)