

告别复杂与低效：Chonkie让RAG中的文本分块变得简单又快速！

原创 南七无名氏 PyTorch研习社 2025年02月16日 09:13 安徽

在自然语言处理（NLP）领域，Chunking（文本分块）指的是**将一段长文本拆分成更小的部分（chunks），这些部分可以是词语、句子、段落或语义单元**。这种拆分通常是为了简化后续处理，提高计算效率，并使得模型能够更好地理解文本的结构和语义。

在处理大规模文本时，直接将整个文本传入模型会带来很多问题。例如，模型的输入大小往往是有限制的（GPT-2 的最大输入长度是 1024 个 token）。如果输入的文本超出了这个限制，模型就无法正确处理。因此，Chunking 成为了必不可少的步骤，它有以下几个重要作用：

1. 提高模型的效率：

大部分深度学习模型都有固定的输入长度限制，过长的文本会导致内存不足或性能下降。通过将文本分块，我们可以保证每一块的长度都在模型的输入限制范围内，从而提高模型的效率。

2. 增强上下文理解：

文本中的语义有时会跨越多个句子或段落，而某些任务（如 RAG 应用）需要根据不同的 chunk 进行信息检索。通过合理的 Chunking，我们可以保留更清晰的上下文，使得模型能够更好地理解语义关系，避免信息丢失。

3. 提升RAG的效果：

在 RAG 应用中，信息检索是关键步骤。Chunking 使得我们可以更精细地控制检索范围，提升检索和生成结果的准确性和相关性。通过将文本按一定规则分块，模型能够更加精准地从分块中提取关键信息，提升生成结果的质量。

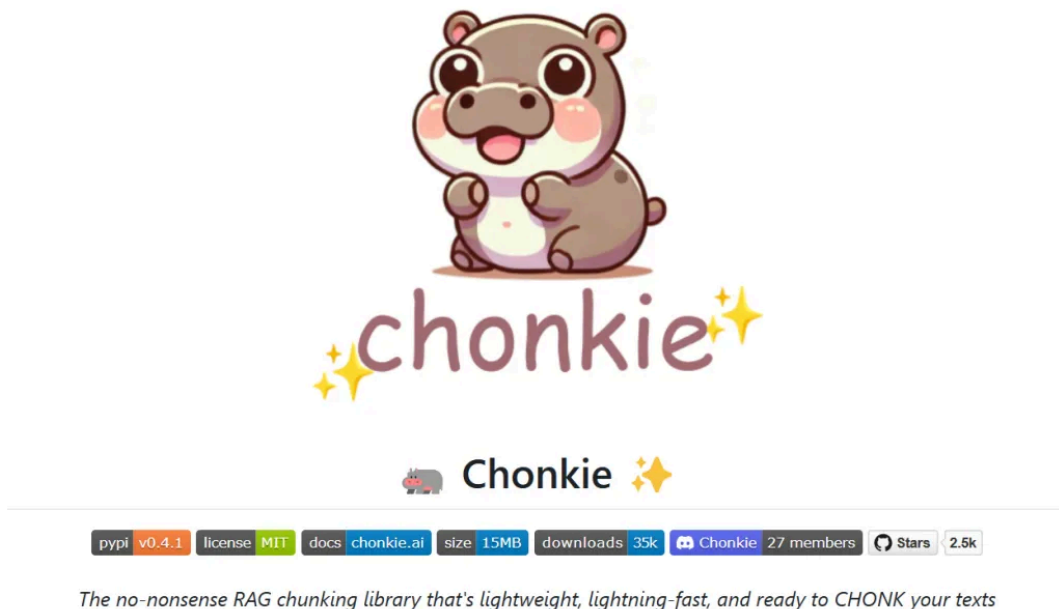
4. 减少计算负担：

对于大量文本数据，直接将其作为整体传给模型进行处理不仅低效，还可能导致性能瓶颈。通过分块处理，模型可以并行处理多个较小的文本片段，大大提高计算效率。

在RAG模型中，信息检索和生成是两个核心步骤。首先，系统会根据查询从一个大的文档集合中检索相关内容。然后，模型将基于检索到的信息生成回答。而为了让检索的内容更加精确，Chunking 在这里发挥着重要作用。通过将文本分割成合适的块，RAG 模型能够更好地理解文本的结构，进而提高信息检索的效果，使得生成的答案更加准确和相关。

但在构建 RAG 应用时，Chunking 操作常常让人头疼。尤其是在面对各种复杂的库和工具时，开发者往往不得不在繁琐的安装步骤和低效的实现之间做出选择。你可能已经厌倦了编写自己手动实现的 Chunker，但又无法找到一个**既简单又高效**的解决方案。那么，今天介绍的 **Chonkie** 库，可能正是你所需要的。





Chonkie：一款简单、快速、轻量的 Chunking 工具

Chonkie 库的设计目标是简化 Chunking 操作，让你可以轻松地将集成到RAG应用中。它的特点如下：

- **功能强大**：支持多种 Chunking 方式，可以满足你不同的需求。
- **易于使用**：安装、导入、调用，简单直白。
- **高速运行**：速度堪比闪电，让你高效处理文本数据。
- **轻量级**：仅提供必需的功能，没有冗余代码。
- **广泛支持**：兼容各种常见的 tokenizer，扩展性强。

✳ Chonkie支持的分块方法

- **TokenChunker**：按固定大小的 token 进行分块。
- **WordChunker**：按单词进行分块。
- **SentenceChunker**：按句子进行分块。
- **RecursiveChunker**：基于可定制的规则进行层次化分块，适合语义分块。
- **SemanticChunker**：根据语义相似性进行分块，适用于需要理解上下文的 RAG 应用。
- **SDPMChunker**：基于语义双重合并方法（Semantic Double-Pass Merge）进行分块。
- **LateChunker**（实验性）：先进行文本嵌入，再进行分块，从而得到更高质量的 chunk 嵌入。

每种 chunker 都有其独特的优势，可以根据需求灵活选择。

🚀 性能表现：Chonkie 真的很快！

Chonkie 不仅轻量，而且在性能上也表现不凡。下面是它在几个常见任务中的速度对比：

- **大小**：默认安装包仅为 15MB（相比之下，其他工具包可能需要 80MB 到 171MB）。
- **速度**：
 - Token Chunking：比最慢的替代方案快 33 倍。
 - Sentence Chunking：几乎是竞争对手的两倍速。
 - Semantic Chunking：比其他工具快最多 2.5 倍。

如此出色的性能，足以应对大规模的 RAG 应用场景。

安装与使用

🚀 如何安装 Chonkie

安装 Chonkie 非常简单，你只需要执行以下命令：

```
1 pip install chonkie
```

如果你需要更多的 chunker 支持，也可以使用：

```
1 pip install chonkie[all]
```

不过，推荐仅安装必要的 chunker，以避免包体积过大。

⚡ 如何使用 Chonkie

Chonkie 提供了多种 chunking 方法，下面是一个简单的示例，展示如何用它来处理文本：

```
1 from chonkie import TokenChunker
2 from tokenizers import Tokenizer
3
4 # 初始化tokenizer
5 tokenizer = Tokenizer.from_pretrained("gpt2")
6
7 # 创建chunker
8 chunker = TokenChunker(tokenizer)
9
10 # 对文本进行分块
11 chunks = chunker("Woah! Chonkie, the chunking library is so cool! I love
12
13 # 输出每个chunk的信息
14 for chunk in chunks:
15     print(f"Chunk: {chunk.text}")
16     print(f"Tokens: {chunk.token_count}")
```



在上面的代码中，我们首先导入了 **TokenChunker**，然后使用预训练的 GPT-2 模型作为 tokenizer 进行分块。通过调用 chunker 方法，我们可以轻松将文本分割成多个 token chunk。

更多信息请前往官方仓库：

<https://github.com/chonkie-ai>



PyTorch研习社

打破知识壁垒，做一名知识的传播者

670篇原创内容

公众号

