

1、系统概述

1.1、系统简介

在我们日常学习和看论文的过程中，经常需要复制/翻译 PDF 中的内容用以询问 AI，在阅读比较长的文献时也会需要从文章中快速定位相关知识点。传统的解决方案是通过 OCR 获取纯文字，同时搜索也基于以上的纯文字，这样就难免出现低效率和不准确的情况，针对以上的情况，我们设计了学生文档助手这个程序。

我们的学生文档助手是一个前后端分离的项目，其主要分成前端和后端两个系统。其中，前端又可分为处理结果展示模块、嵌入结果导出、设置三个子模块；后端又可分为 PDF 转 Markdown 模块、数据库模块、设置相关模块和 RESTfulAPI 服务器模块。

1.2、系统目标

功能目标：pdf 预览和 markdown 预览，数据库存储和检索，资源下载

性能目标：转化时间<10s,

markdown 转换精度 100%,

数据库信息切割精度>97%

1.3、系统运行环境

硬件平台：

CPU：12th Gen Intel® Core(TM) i7-12700H (20) @ 4.70 GHz

GPU：NVIDIA GeForce RTX 3050 Laptop GPU (2048) @ 2.10GHz

操作系统：Windows® 11 / Fedora Linux 42 (KDE Plasma Desktop Edition) x86_64

数据库系统：SQLite3

编程平台：VScode / Trae (with Extensions)

网络协议：HTTP

1.4、开发环境

工程工具：VScode v1.1.00.0 / Trae v1.0.14321 (with Extensions)

开发语言：前端：Vue3；后端：Python3.12

前端打包工具：Vite v4

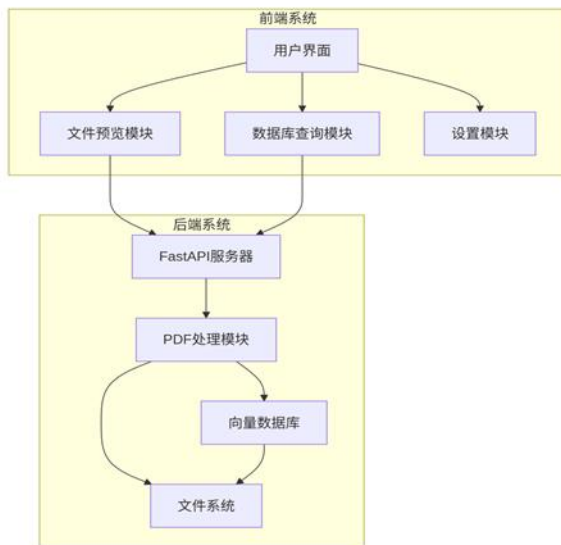
后端运行时：Uvicorn ASGI Server (Python 3.12 Conda Environment)

接口测试工具：Reqable v2.33.7

2、总体结构设计

2.1、软件结构

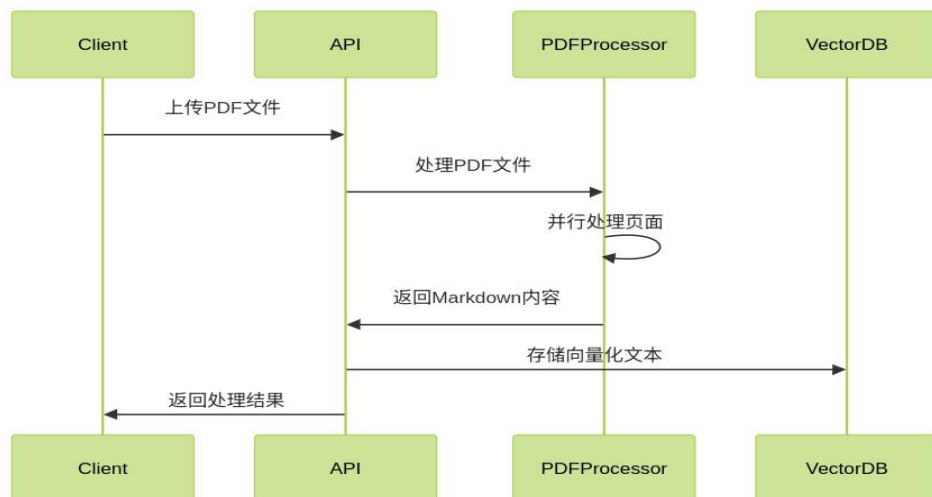
系统采用前后端分离的架构设计，主要分为以下模块：



3、模块设计

3.1、PDF 处理模块

3.1.1、设计图



3.1.2、功能描述

PDF 处理模块负责将 PDF 文档转换为 Markdown 格式，并支持图片提取和并行处理。

3.1.3、输入数据

PDF 文件（通过 multipart/form-data 上传）

处理参数：

- use_parallel: 是否使用并行处理
- num_workers: 并行处理的工作线程数

3.1.4、输出数据

Markdown 格式的文本内容

提取的图片文件

处理状态和错误信息

3.1.5、数据设计

给临时文件存储：

- temp/：临时 PDF 文件
- images/：提取的图片文件
- pdf/：已处理的 PDF 文件备份

3.1.6、算法和流程

1. 文件上传与验证
2. 并行页面处理：
 - 页面分批
 - OCR 文本识别
 - 图片提取
3. Markdown 转换
4. 向量化存储

3.1.7、函数说明

server/src/parallel_pdf_extract.py

`def process_pdf(pdf_path: str, use_parallel: bool = True, num_workers: int = 3):`

`"""`

处理 PDF 文件

参数:

`pdf_path`: PDF 文件路径

`use_parallel`: 是否使用并行处理

`num_workers`: 工作线程数

返回:

`InferenceResult` 对象

`"""`

3.1.8、全局数据结构与该模块的关系

访问 VectorDB 类进行向量存储

使用临时文件系统存储处理结果

4、数据库与数据结构设计

系统使用 SQLite3作为数据库管理系统，主要用于存储文档内容和向量数据。

4.1、数据库表设计

主要是存储文档内容和向量数据的主表

```
1 CREATE TABLE documents (  
2     id INTEGER PRIMARY KEY, -- 自增主键  
3     filename TEXT NOT NULL, -- 文档名称  
4     content TEXT NOT NULL, -- 文本内容片段  
5     vector BLOB -- 向量化后的内容数据  
6 )
```

字段说明：

- id: 自增主键，用于唯一标识每条记录
- filename: 文档名称，用于关联原始文件
- content: 文本内容片段，存储分割后的文本段落
- vector: BLOB 类型，存储向量化后的数据

4.2、文件系统结构

除数据库外，系统还维护以下目录结构：

```
├— temp/          # 临时文件目录  
├— images/        # 提取的图片文件  
├— pdf/           # 原始 PDF 文件  
└— source_docs/   # Markdown 文档存储
```

5、接口设计

5.1、用户接口

系统提供 Web 图形用户界面，主要包含以下功能区域：

1. PDF 文件上传和预览区
2. Markdown 预览和复制区
3. 向量检索交互区
4. 系统设置面板

5.2、外部接口

系统通过 RESTful API 提供以下主要接口

5.2.1、PDF 处理接口

```
POST /upload_pdf  
Content-Type: multipart/form-data  
  
请求参数:  
- file: PDF 文件  
  
返回数据:  
{
```

```
"status": "success|error",
"filename": "文件名",
"message": "处理结果信息"
}
```

5.2.2、向量检索接口

POST /search

Content-Type: application/json

请求参数:

```
{
  "text": "搜索文本",
  "top_k": 10  // 可选，默认返回前10条结果
}
```

返回数据:

```
{
  "status": "success|error",
  "results": ["匹配文本1", "匹配文本2", ...]
}
```

5.2.3、文件管理接口

从 server.py 中实际实现的 API 路由:

```
1  @app.get("/files")
2  async def get_files():
3      try:
4          files = vector_db.get_all_filenames()
5          return JSONResponse(content={
6              "status": "success",
7              "files": files
8          })
9      except Exception as e:
10         return JSONResponse(
11             status_code=500,
12             content={"status": "error", "message": str(e)}
13         )
14
15  @app.get("/file/content/{filename}")
16  async def get_file_content(filename: str):
17      try:
18          content = vector_db.get_markdown_text(filename)
19          if content is None:
20              return JSONResponse(
21                  status_code=404,
22                  content={"status": "error", "message": "File not found"}
23              )
```

```

24         return JsonResponse(content={
25             "status": "success",
26             "filename": filename,
27             "content": content
28         })
29     except Exception as e:
30         return JsonResponse(
31             status_code=500,
32             content={"status": "error", "message": str(e)}
33         )
34
35 @app.get("/file/pdf/{filename}")
36 async def get_pdf_file(filename: str):
37     try:
38         pdf_path = f"pdf/{filename}"
39         return FileResponse(
40             path=pdf_path,
41             media_type="application/pdf"
42         )
43     except Exception as e:
44         return JsonResponse(
45             status_code=500,
46             content={"status": "error", "message": str(e)}
47         )
48
49 @app.get("/images/{filename}")
50 async def get_image(filename: str):
51     try:
52         image_path = f"images/{filename}"
53         return FileResponse(
54             path=image_path,
55             media_type="image/jpeg"
56         )
57     except Exception as e:
58         return JsonResponse(
59             status_code=500,
60             content={"status": "error", "message": str(e)}
61         )

```

5.2.4、数据导出接口

GET /export

返回: ZIP 文件, 包含所有 PDF、Markdown 和图片文件

5.3、内部接口

5.3.1、模块间通信

以下展示了从 server.py 中提取的核心模块间通信代码:

1. PDF 处理模块 → 向量数据库模块

1 # 处理 PDF 并存储到向量数据库

2 result = process_pdf(temp_path, use_parallel=True, num_workers=10)

3 markdown_content = result.pipe_txt_mode(dr).get_markdown("output")

4 true_markdown = process_images(markdown_content)

5 vector_db.add_markdown(true_markdown, filename.replace(".pdf", ""))

2. 向量数据库模块 → 文件系统

1 def get_markdown_text(self, filename: str):

2 """获取原始 Markdown 内容"""

3 file_path = f'source_docs/{filename}.md'

4 if not os.path.exists(file_path):

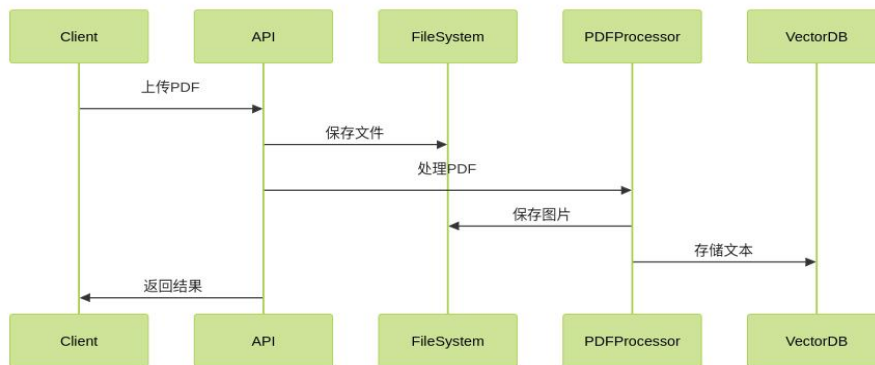
5 return None

6 with open(file_path, 'r', encoding='utf-8') as f:

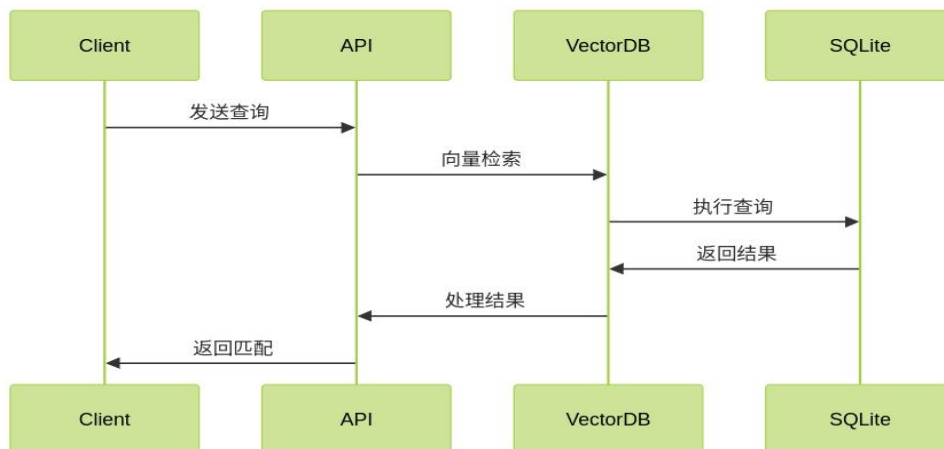
7 return f.read()

5.3.2、数据流

1. PDF 上传流程



2. 检索流程



6、其他设计

6.1、性能优化

6.1.1、并行处理机制

系统采用多线程并行处理机制来提高 PDF 处理速度。从 `parallel_pdf_extract.py` 中可以看到具体实现：

```
1 class ParallelBatchAnalyze:
2     def __init__(self, num_workers=3):
3         self.num_workers = num_workers
4         self.result_queue = queue.PriorityQueue()
5         self.lock = threading.Lock()
6         self.model_manager = ModelSingleton()
7
8     def process_batches(self, batches: List[List[Tuple]], ocr: bool, show_log: bool,
9                        layout_model=None, formula_enable=None, table_enable=None) -> List:
10        """使用线程池并行处理批次"""
11        with concurrent.futures.ThreadPoolExecutor(max_workers=self.num_workers) as executor:
12            futures = []
13            for batch_id, batch in enumerate(batches):
14                future = executor.submit(
15                    self.worker,
16                    batch_id,
17                    batch,
18                    ocr,
19                    show_log,
20                    layout_model,
21                    formula_enable,
22                    table_enable
23                )
24            futures.append(future)
```

关键优化点：

1. 使用 `ThreadPoolExecutor` 进行并行处理
2. 可配置的工作线程数（`num_workers`）
3. 使用优先级队列保证结果顺序
4. 线程锁保证并发安全

6.1.2、并行处理机制

系统在向量数据库模块中选择了 `all-MiniLM-L6-v2` 作为文本向量化模型，这个模型具有以下优势：

1. 选择384维向量表示，相比其他模型（如768维或更高）大幅降低存储和计算开销
2. 保持较好的语义表示能力，适合文本相似度检索
3. 批量处理文本分块，提高向量化效率

6.2、安全设计

系统采用异常处理机制来保证运行时的稳定性。对于文件操作、数据库访问等关键操作，都实现了异常捕获和错误响应处理，确保在出现异常情况时能够返回合适的错误信息，并进行必要的资源清理。