

# 基于LLM和向量数据库的自动文档分类系统设计

## 1. 项目整体概述

现代用户在 OneDrive 等云盘中往往存放了海量文档，但由于缺乏有效的整理，文件命名混乱、分类不清，导致查找和管理非常困难。针对这一痛点，本项目旨在开发一个基于 **LLM（大型语言模型）** 和 **向量数据库** 的自动文档分类系统，帮助用户对 OneDrive 文件夹中的文档进行智能归类和重命名，缓解文档混乱和分类交叉难管理的问题。

该系统的主要目标是：利用 AI 自动分析文档内容，根据语义将文档归入合理的类别或打上标签，并按规则重命名和组织文件夹结构，从而让用户的云盘文件 **结构清晰、命名规范**。这样用户能更快捷地检索资料，避免重复文件和遗漏。目标用户包括：需要管理大量文档的知识工作者、研究人员，以及个人资料较多的用户等。主要使用场景涵盖：

- **日常文件管理**：用户将杂乱无章的文件批量交由系统整理。系统自动识别文件内容（如工作文档、财务票据、个人照片等），移动到相应分类文件夹，并按内容生成规范名称。
- **持续文件监控**：系统可后台监听 OneDrive 同步文件夹，新文件加入或修改时自动触发分类和归档，不中断用户工作。
- **跨类别组织**：针对属于多个类别的文件（例如某项目报告同时涉及财务和技术），系统支持交叉标签，使文件能出现在多个虚拟目录视图中，方便多角度管理。

总之，本项目通过引入 LLM 智能理解文本以及向量相似度匹配，实现对用户文档的 **智能分类、规范命名和多维度组织**，解决用户文档管理难题，提高工作效率。

## 2. 功能列表

系统功能涵盖从文件读取解析，到分类决策、文件操作，再到用户反馈的一整套流程，具体包括：

- **文件夹监听**：持续监控指定的 OneDrive 本地同步文件夹。当有新增文件、文件修改或移入时，自动触发处理流程。用户也可手动运行整理命令。监听支持递归子文件夹，确保整个目录树都被监控。通过实时监控，保证**及时分类**新文档，减少人工介入。
- **内容抽取**：对新检测到的文件执行内容提取，支持多种文档类型。针对 PDF、Word、PPT 等文本型文件，使用解析工具提取全文文本；对图片扫描件或照片，可集成 OCR 获取文字；对音频视频等可选用转录（未来扩展）。内容抽取模块输出标准化的文本内容或摘要，用于后续分类判断。通过这一步骤，系统能够**理解文档的实际内容**，而不仅仅依据文件名。
- **自动分类**：基于提取的内容，利用向量嵌入和 LLM 推理对文档进行自动分类。系统将文档内容转换为向量，通过与已有**类别样本**或**已分类文档**进行相似度比较，初步判断所属类别，并调用 LLM 结合上下文进一步确认分类。支持单标签分类和多标签分类（交叉分类）：对于内容涉及多个领域的文档，可同时赋予多个标签。分类结果包括类别标签及置信度评分。整个过程无需人工干预，实现**智能语义分类**。

- **交叉标签支持**：系统允许一个文档归属于多个标签，使得分类体系不再局限于单棵文件夹树。例如，同一文件可被标记为「项目A」和「财务」两种类别。为了支持这种交叉分类，系统提供**多视图机制**（见后述），让用户从不同维度都能找到该文件。交叉标签支持避免了硬拷贝文件，实现**同一文件多处呈现**而无冗余。
- **自动命名**：根据文档内容和分类结果，为文件生成**规范化的文件名**。系统可提取文档中的标题、日期等关键信息，或通过 LLM 从内容摘要中生成简洁标题，并结合预设模板重命名文件。例如，原始名称“Document\_123.pdf”可自动更名为“项目A-季度报告-2023Q4.pdf”。命名规则可配置，确保文件命名风格统一，便于识别和排序。
- **交互修正**：提供用户交互界面或 CLI 提示，让用户参与修正系统决策。对于置信度不足或存在歧义的分类，系统可以进入**审查模式**：暂停自动处理，提示用户候选类别或标签，允许用户确认或选择正确分类。用户也可以查看日志，对错误分类进行反馈（例如指定正确类别）。交互修正保证系统**可控且可学习**：用户的反馈可用于调整模型或生成新规则，逐步优化分类准确率。
- **规则管理**：支持用户定义和管理自定义规则（规则引擎）以指导分类和命名。规则采用配置文件（例如 YAML）描述，内容包括匹配条件和执行动作。例如，用户可定义规则：“如果文件名包含‘发票’，则添加类别标签‘财务’”；又如“如果文档内容出现‘机密’，则标记为‘保密’并在名称后加‘CONF’标记”。规则管理功能允许用户将行业知识或个人偏好融入系统，实现**可定制的分类行为**。系统提供默认规则模板，用户可修改扩展，并支持加载/保存、多版本管理等。
- **审计与回滚**：系统将所有分类决策和文件操作记录到日志或数据库中，便于审计追踪。记录信息包括时间戳、文件原路径和新路径、更名前后名称、分配的类别标签、置信度及所用规则等。基于审计日志，系统支持**回滚**功能：当出现误分类或不当移动时，用户可一键将文件恢复到原始位置及名称。回滚机制确保用户对文件有最终控制权，增加操作的可靠性和信任度。

以上功能协同工作，使系统既能自动化处理常规情况，又能在特殊或不确定情形下让用户介入调整，保证最终分类结果准确、符合用户预期。

### 3. 技术栈选择

实现该自动文档分类系统，需要选择合适的技术栈以满足文本解析、向量计算、LLM 推理以及流程编排等需求。下面是本项目各部分的技术选型及其简要理由：

- **语言与环境（Python）**：后端采用 Python 开发。Python 拥有丰富的文本处理和机器学习生态，方便调用 **PDF 解析、向量计算、LLM 接口**等库。此外，Python 的跨平台特性有助于在不同操作系统的 OneDrive 环境下运行。本项目定位于本地运行的工具，Python 的脚本式部署和可扩展性非常契合需求。
- **文档内容解析工具**：使用成熟的开源库提取各种格式文档的正文内容。主要考虑 **pdfminer.six**（针对 PDF 文本精细提取），**Texttract**（支持多种文件类型的文本抽取）等工具。对于 Office 文档，可用诸如 `python-docx`、`pptx` 库解析文字，对图片扫描件可结合 Tesseract OCR。选择这些工具是因为它们**支持格式广泛且社区成熟**，能最大程度覆盖用户 OneDrive 内的常见文件类型，准确获取文本而不丢失信息。
- **嵌入模型（bge-m3 / E5）**：在生成文本向量表示时，选用高性能的预训练嵌入模型。例如 **BGE-M3**（北大智源的多语种多功能通用嵌入模型）或 **E5 系列模型**（开源通用文本 Embedding 模型）。这些模型能将文档内容映射为高维向量，向量相近表示语义相近<sup>① ②</sup>。选择它们的原因在于：①支持**中英文等多语言**

(OneDrive中文件可能包含英文报告和中文资料混杂)；②对语义检索和分类任务有良好表现，可提升**分类准确度**。模型可以通过 HuggingFace Transformers 在本地加载，结合 GPU 加速以处理大量文本嵌入计算。

- **向量数据库 (Chroma)**：使用 **Chroma** 作为向量数据库存储和检索文档向量。Chroma是轻量级的本地向量存储，提供简单的Python接口和持久化方案，适合在个人电脑上运行。将每份文档内容嵌入后存入 Chroma，方便后续**近似最近邻搜索**，快速找到内容相似的已分类文档或类别向量。相比远程向量服务，Chroma开源且**本地化部署**，保证数据隐私的同时响应迅速，无需网络依赖。
- **RAG 框架 (LlamaIndex)**：采用 **LlamaIndex** (又称 GPT Index) 构建检索增强生成 (RAG) 的索引框架。LlamaIndex 可以将文档内容与其分类标签一起建立索引，方便利用 LLM 查询。当需要让 LLM 决策分类或回答用户查询时，LlamaIndex充当中介，从向量库中检索相关文档/知识作为上下文提供给 LLM<sup>3</sup><sup>4</sup>。这一框架简化了“**内容-检索-生成**”的流程，使我们能够容易地实现诸如“根据已归档知识判断新文档类别”等推理任务，同时也为将来扩展问答功能打下基础。
- **编排引擎 (LangGraph)**：整个分类流程由多步骤组成，引入 **LangGraph** 编排引擎来组织这些步骤的执行。LangGraph 是 LangChain 生态中的底层框架，支持用有向图 (DAG) 定义复杂工作流<sup>5</sup>。借助 LangGraph，可以将文档处理流程拆解为节点：监听触发节点、内容提取节点、向量生成节点、分类决策节点、文件操作节点、反馈处理节点等，并定义节点之间的数据流转和条件分支。LangGraph 提供 **有状态、可控**的执行环境，支持出错重试、分支逻辑 (如低置信度时走人工确认分支) 和并行处理多个文件等，确保系统在复杂场景下依然可靠、清晰。
- **数据库与规则配置 (SQLite + YAML + Jinja2)**：系统会使用 **SQLite** 作为本地关系型数据库，用于保存元数据和日志 (例如文件索引、分类结果、操作记录、用户反馈等)，以便持久存储和高效查询。同时，选择 **YAML** 文件格式保存用户配置和规则，例如分类规则和命名模板。YAML易于阅读编辑，方便用户自定义类别、关键词和模板。**Jinja2** 模板引擎则用于在规则中编写动态命名或路径模式：用户可以在 **YAML** 中使用 Jinja2 语法插入变量 (如文档标题、日期、分类名等) 来定义文件重命名或路径规则。通过 **SQLite** 保证数据一致性，通过 **YAML+Jinja2** 提供**灵活的配置能力**，两者结合使系统既可靠又可定制。
- **LLM Provider 自动切换 (Ollama / OpenAI / Claude)**：为了在不同环境下取得最佳性能，系统设计支持多种 LLM 提供方式的自动切换。**Ollama** 是本地部署LLM的引擎，可运行诸如 Llama2 等模型，用于离线或保护隐私的场景；**OpenAI API** (如 GPT-4) 和 **Anthropic Claude API** 则可在需要更高质量或更专业语言理解时调用。系统通过抽象LLM调用接口，根据配置或运行时判断自动选择Provider：例如，当本地有足够计算资源时优先使用Ollama以节省成本，如遇复杂任务或本地模型信心不足时切换调用云端GPT-4或Claude，以提高**准确率**。还可设定容错策略：某一服务不可用时，自动改用替代模型，保证分类流程**稳定连续**。

以上技术栈的组合，奠定了系统实现的基础：Python 提供开发便利，解析库和嵌入模型实现内容理解，向量库与 LLM框架负责智能决策，而LangGraph串联起整个流程，SQLite/YAML/Jinja2保障配置和数据管理，最终通过灵活调用不同LLM提供商，系统将在本地高效运行并产出可靠结果。

## 4. 功能模块和交互关系

本系统采用模块化设计，将自动分类流程拆解为若干功能模块。各模块通过 **LangGraph** 工作流编排引擎串联，形成有序的流水线，每个模块对应一个或多个 **LangGraph** 节点。下面详细说明各模块的职责、它们之间的交互关

系，以及如何利用 LlamaIndex 构建“已归类知识库”辅助推理。同时阐述标签机制与多视图设计如何支持交叉分类。

## 内容提取模块（LangGraph节点：DocumentParser）

**职责：** 接收文件监听模块传来的文件路径，解析文件内容并输出纯文本或结构化文本。该模块根据文件类型选择合适的解析方法：- 对于 **PDF** 文档，使用 pdfminer.six 提取页面文字，保持段落顺序。- 对于 **Word/PowerPoint** 等 Office 文件，利用相应库或调用 Textract 获取文字。- 针对 **文本文件** (txt、markdown 等)，直接读取内容。- 若文件为 **图片或扫描件PDF**，可调用 OCR 引擎提取文字（如使用 Tesseract，必要时结合版面分析提高准确率）。- 对于 **音频/视频文件**（如出现此类，需要分类时），可选用语音识别模型将语音转文字（此为扩展功能）。

**输出：** 标准化的文本内容（可能附带元数据，如文件名、扩展名、基本属性）。文本过长时模块可截取摘要或分段输出：例如对100页的PDF，可产出前几页摘要供分类，也可将全文分块发送给后续模块逐块处理。

**LangGraph实现：** DocumentParser节点封装上述逻辑，对每个文件产生一个内容提取任务。该节点将结果（文本或解析失败错误）输出到后续节点，并记录解析情况。错误情况下（如文件加密或格式不支持），LangGraph 可以捕获异常，将文件标记为“无法解析”并跳过分类或通知用户。内容提取模块是整个流程理解文档的**基础**，它为后续嵌入和分类提供必要的原料。

## 嵌入生成模块（LangGraph节点：Embedder）

**职责：** 接收文档文本内容，调用嵌入模型将其编码为高维向量表示。通过向量，系统可利用数学空间中的相似度来衡量文档语义相近性。具体步骤：- 使用预选的 **bge-m3 或 E5** 模型，将文本输入模型得到 embedding 向量（如 512维或1024维浮点数组）。- 若文本较长超过模型输入限制，可能进行**摘要或分段嵌入**：例如提取内容中最具代表性的部分，或对每个段落嵌入再取平均/最大池化得到文档整体嵌入。- 为了兼顾分类精度和效率，可对文本预处理：去除常见无关词，确保embedding更多反映主题关键词。

**输出：** 文档的向量表示，以及可选的一个**简要文本摘要**（摘要可由内容提取模块直接提供，如文档标题/首段）。摘要可在后续生成命名或提示 LLM 时使用。

**LangGraph实现：** Embedder节点对从解析节点传来的文本执行同步的向量化计算。由于嵌入模型调用可能比较耗时，LangGraph可以并行处理多个文件的Embedder节点（在资源允许情况下），提高速度。Embedder节点完成后，将向量和摘要一并传递给分类决策模块。值得注意的是，该节点也可以调用不同模型：如首选使用本地模型（通过 Ollama 或 HuggingFace），如果失败则LangGraph分支切换到调用OpenAI等API节点获取嵌入。这样通过LangGraph的条件节点实现**Provider自动切换**，使嵌入生成稳定可靠。

## 分类判定模块（LangGraph节点：ClassifierAgent + RetrievalActions）

**职责：** 根据文档向量及文本摘要，确定文档所属的类别标签。此模块结合**向量检索**和**LLM推理**两方面，具体流程：  
1. **相似检索：** 将新文档的向量提交给 **Chroma 向量数据库**，检索出与之语义相似的若干已知项。已知项可以是**已有分类的文档或者类别典型向量**。例如，系统最初可使用用户提供的示例文件或描述生成每个类别的参考向量；随着运行，Chroma 中逐渐积累已分类文档向量。当新文档到来时，检索返回 top-N 相近向量项及其关联信息（比如对应的类别标签、文件名、内容摘要）。  
2. **LLM综合判断：** 将新文档的摘要和检索到的相似文档信息一起喂给 LLM，让其综合分析**最合适的类别**。可以设计一个提示模板，例如：

系统：你是一个文件分类助手。现在有一份新文档的内容摘要，以及数个相似的已分类文档供参考，请根据语义判断新文档属于哪些类别。

已有类别及示例：

- 类别A：（示例：文档X 摘要...）

- 类别B：（示例：文档Y 摘要...）

新文档摘要：...

请给出新文档最可能的类别，并说明是否需要多个标签。

LLM 根据提供的知识库上下文和自身常识，输出分类建议。如模型明确认为属于某一类别，会返回该类别；如果内容跨领域，可能返回多个标签及置信度评估。

1. **规则校验**：在LLM建议基础上，模块结合用户定义的规则进一步校验或调整。比如，如果用户规则规定某关键词必须归属特定类别，则无论LLM结果如何都应满足该约束（除非人工 override）。规则校验模块可在 LLM 判定前或后执行：简单条件（如文件扩展名/文件名规则）可先匹配以直接得出分类，提高效率；内容相关规则可在LLM结果后检查冲突并修正。

**输出**：最终确定的类别标签列表（通常1个主类别，可能附带其他次要标签）以及每个标签的置信度评分或理由。如果分类置信度不足，本模块应在输出中标记需要人工复核。对于不确定的情况，Classifier模块可输出一个特殊类别如“Uncategorized”或“NeedsReview”作为占位。

**LangGraph实现**：分类判定往往拆成多个节点：

- **RetrievalNode**：负责与Chroma交互获取相似项。
- **LLM Classification Node (ClassifierAgent)**：接收相似项和摘要，通过Prompt调用LLM得到分类结果。
- **RuleCheckNode**：（可选）应用规则调整LLM输出。

LangGraph 可以让 RetrievalNode 和 ClassifierAgent 串联，或者封装为一个复合智能体。在低置信度分支，LangGraph 会将流程引向“人工确认”节点（InteractiveReview），暂停自动执行等待用户反馈。**已归类知识库的构建**：这里利用 LlamaIndex 将所有已分类文档及其标签建立索引，使 RetrievalNode 更有效：LlamaIndex可以存储文档文本、embedding及元数据标签，支持更丰富的查询。例如，将新文档内容作为查询，直接让 LlamaIndex 找到语义相关的几个文档段落及其类别供 LLM 使用。这相当于动态地从“已归类知识库”中取出**案例参考**，提高 LLM 分类准确度和解释性（类似于让模型“以例为鉴”）。因此，ClassifierAgent节点每次决策前都可通过 LlamaIndex 检索，将返回的内容附加在Prompt里。分类判定模块是系统智能的核心，借助向量相似度和LLM理解，实现了**基于内容的自动决策**。

## 路径决策模块（LangGraph节点：PathPlanner）

**职责**：根据分类结果决定文件在 OneDrive 中的新存储路径（即应移动到的目标文件夹）。主要步骤：

- **目标文件夹确定**：每个类别通常对应一个文件夹。如用户预先定义了类别与文件夹的映射，则直接使用：例如类别“财务文件”映射到 OneDrive/分类/财务文件/ 目录下。如果映射未预定义，系统可按默认规则在主目录下创建同名文件夹（如 OneDrive/分类/项目A/ ）。对于多标签情况，需要确定一个**主存储路径**：可以以主类别为顶层文件夹，其余标签通过软链接方式处理（后述）。如果分类层次有多级（如“客户/项目/文档类型”三层分类），PathPlanner 应根据规则生成多级子目录路径。
- **文件夹创建与存在性检查**：检查目标路径是否存在，不存在则创建相应文件夹。创建时需要考虑并发和命名冲突（例如Windows下不区分大小写，同名文件夹判断等）。对于用户已有的部分分类文件夹，系统应识别并重用，不重复创建。必要时可以扫描OneDrive已有目录结构，对照类别列表匹配。
- **特殊文件处理**：某些文件可能被标记为多种特殊属性，如“归档”“重要”等，可能对应特殊路径（例如“重要”标签的文件也复制或链接到“Important”汇总文件夹）。PathPlanner在规划路径时需考虑这些**全局视图**文件夹，将链接创建纳

入计划。 - **移动策略**：PathPlanner不直接执行移动，但生成移动所需的信息：包括**目标完整路径**（含文件名）和所需的所有衍生路径（如次要标签的链接路径）。它会将这些路径信息交给下游模块。此模块也可在这里处理**命名冲突**：如果目标目录下已存在同名文件，PathPlanner可以选择在文件名后添加区分后缀（例如“\_1”）或提示Rename模块调整名字，防止覆盖。

**输出**：建议的文件目标路径结构，包括： - `primary_path`：主文件实际存放的新路径（目录 + 新文件名）。 - `link_paths`：若多标签，需在其它标签目录下创建的快捷方式/软链接路径列表。 - 可能还包括 `old_path` 和标记动作类型（移动/复制/跳过等）。

**LangGraph实现**：PathPlanner节点是一个纯逻辑节点，根据分类结果和配置规则（例如 `rules.yaml` 中定义的路径模板）计算路径。它通过读取配置中定义的模式（可能包含 Jinja2 模板，如 `"{category}/{year}/"`）来拼装路径字符串。LangGraph 在此可以检查输出，例如若某文件被分类为“未分类”或需要人工处理，则PathPlanner可以选择输出特殊路径（如 `OneDrive/分类/待整理/` 目录）或者不输出路径并将流程引向人工处理分支。正常情况下，PathPlanner将把决策的路径发送给后续节点执行实际移动。**路径决策模块**确保每个文件有恰当的“归宿”，构建出合理的文件结构架构。

## 命名生成模块（LangGraph节点：Renamer）

**职责**：按照命名规则为文件生成新的文件名。良好的命名可以直观体现文件内容、时间等关键信息。本模块利用文档内容或元数据生成有意义的名称： - **基于内容的命名**：如果文档内容包含标题（例如报告的标题、邮件主题行等），优先提取作为文件名的一部分。解析模块可能已经提取出标题行供使用。若无明显标题，可利用 LLM 对文档摘要生成一个短语作为标题。例如模型Prompt：“根据以下内容摘要，生成5-10字的中文标题：...”。LLM返回的标题可用于命名。 - **元数据补充**：文件的创建日期、作者（如果可获取）等信息可按照模板插入。例如在培训课件文件名中加上年份，或在会议记录中加上日期。若规则要求，Renamer可从文件的EXIF/属性中读取日期，或使用文件修改时间等。 - **模板应用**：用户可在配置中指定命名模板（使用Jinja2语法）。Renamer模块会将可用变量（如 `{{title}}`、`{{category}}`、`{{date}}`、`{{ext}}` 等）填充入模板，生成最终文件名字符串。举例：模板定义 `"{{category}}-{{title}}-{{date|strftime('%Y%m%d')}}.{{ext}}"`，在类别=财务，title=报销单，date=2023-08-15，ext=pdf时，生成文件名“财务-报销单-20230815.pdf”。模板提供了高度灵活性。 - **字符合法化**：对生成的文件名进行清理，确保不含文件系统不允许的字符（如 `\/:*?` 等），以及避免过长。中文环境下注意编码问题以及OneDrive对路径长度的限制。 - **避免冲突**：如果PathPlanner已指明目标目录下存在同名文件，Renamer需要调整名称避免覆盖。例如可在名称末尾添加序号或区分符。也可依据策略：如如果文件内容高度相似于现有文件，可能是重复文档，系统可以提醒用户而不是简单重命名。

**输出**：最终的新文件名，通常会直接整合到 PathPlanner提供的目标路径中（PathPlanner输出可能包含占位符文件名或原名，这一步将确定最终名称）。Renamer也可以输出旧名与新名的映射供日志记录。

**LangGraph实现**：Renamer节点获取 PathPlanner 给出的路径和所需变量，从数据库或前序节点中拿到文档解析信息（如提取的标题、日期）和分类结果，然后应用模板生成名称。LangGraph保证 Renamer在 Move 执行前完成。如果Naming依赖LLM（生成标题），Renamer节点内部可能再次调用 LLM API。这可以通过LangGraph在Renamer节点内部触发一个子流程（Prompt节点）完成，然后将结果拼入模板。整个命名生成模块确保**文件名称语义清晰**且符合用户习惯，与分类目录一起提升文件管理友好度。

## 移动执行模块（LangGraph节点：FileMover）

**职责：** 执行文件系统操作，将文件从原位置移动/重命名到新路径，并创建所需的链接或快捷方式。执行过程中注重事务性和错误处理：

- **主文件移动：** 调用操作系统文件操作接口（`os.rename` 或 `shutil.move`）将文件从 `old_path` 移动/改名到由 PathPlanner/Renamer 确定的 `primary_path`。如果目标目录不存在则先创建，已在 PathPlanner 中处理。移动操作成功后，在 OneDrive 同步客户端的作用下，云端文件也会相应更新。
- **多视图链接创建：** 对于分类判定中附加的次要标签，在对应标签文件夹下创建**软链接（symlink）**或**快捷方式（Windows .lnk / Mac Alias）**指向主文件位置。软链接允许文件出现在多个目录。例如文件A存储于“项目A”目录下，同时分类标签含“财务”，则在“财务”目录下创建一个指向文件A的新链接。这使用户无论浏览“项目A”还是“财务”文件夹，都能访问该文件。创建链接时考虑操作系统：Unix/macOS 下可用符号链接，Windows 下可通过 Python 调用 shell 接口创建快捷方式文件。链接的命名可与原文件名相同，或包含其他标签标识。
- **原位置处理：** 若文件原先已存在于某分类文件夹（例如这是一次重分类），需要判断是否删除残留的旧快捷方式或空文件夹。如果原位置在监控范围内，则移动后原位置变空可能需要清理（可延后统一清理空目录）。
- **错误回滚：** 如果移动过程中发生错误（例如权限不足、目标磁盘空间不足等），模块应及时捕获异常并采取措施：可以撤销部分已完成的操作（例如已经移动了一半的文件则移动回原位），然后报告失败，避免文件丢失或状态不一致。对于链接创建失败，则记录哪些链接未成功，以便稍后重试或提示用户手动处理。

**输出：** 操作结果状态（成功/失败）以及文件的新路径信息。如果全部成功，新路径会用于更新索引；如果失败，输出错误详情，LangGraph 可根据此将流程转到错误处理或回滚节点。

**LangGraph实现：** FileMover节点负责调用系统IO操作，是最终产生实际文件改动的节点。为确保安全，LangGraph可以在FileMover之前插入一个**用户确认节点**（如果运行在交互模式且设置为迁移前需用户批准），然后再执行移动。移动后，如果LangGraph检测到部分链路创建失败，可调用一个compensation节点记录待办。值得注意的是，通过LangGraph，我们可以批量处理多个文件的移动，但为了避免干扰，每个FileMover动作最好逐个执行并序列化关键操作。此外，LangGraph 还可在此节点上挂接审计记录（详见下文索引回写模块）。移动执行模块直接影响用户云盘文件，设计上要求**小心稳健**。

## 索引回写模块（LangGraph节点：IndexUpdater）

**职责：** 在文件成功分类移动后，更新系统的内部索引、日志和知识库，为未来的检索和推理提供支持，并记录操作审计：

- **向量库更新：** 将新文档的嵌入向量及其分类标签写入 **Chroma** 向量数据库。这样，新文档也成为“已归类知识库”的一部分。当后续有新文件到来时，这份文档就能被检索到作为相似参考。此外，如果有维护类别中心向量的机制，可更新对应类别的中心点（例如对类别向量取所有成员向量的平均）。
- **LlamaIndex 知识库更新：** 使用 LlamaIndex 将新文档内容（或其摘要）及元数据添加到索引。在索引中，将该文档标记上所属类别/标签。未来无论是在分类时提供给 LLM 参考，还是用户查询自己的文档内容，都可以通过该索引高效检索到这份文档的信息<sup>6</sup><sup>7</sup>。这一步保证知识库**与时俱进**，越用系统越聪明。
- **日志审计记录：** 将此次对文件的处理结果记录到 **SQLite 数据库** 或日志文件中。一条完整的记录包括：文件唯一ID或路径、旧路径、最终新路径、旧文件名、新文件名、分配的类别标签、分类置信度、使用的规则命中情况、处理时间戳、操作者（自动/人工）等。日志用于后续审计和可能的回滚。若文件需要人工审核未动，也应记录状态以便提醒用户后续处理。
- **状态标记：** 在数据库中更新文件状态，例如标记此文件已分类，以避免下次重复处理；若多轮分类（文件可能再次修改或规则调整导致需要重新分类），可以比较文件的上次分类时间和修改时间决定是否重新进入流程。也可以将“人工确认”状态写入，方便 CLI 查询有哪些文件待用户确认。

**输出：** 确认索引和日志更新完成的标志。如果出现数据库或索引写入错误，也需要报告（但通常这些操作出错可能性较低，属于本地操作，可重试或在日志中告警）。

**LangGraph实现：** IndexUpdater通常在FileMover成功后执行。可以将IndexUpdater拆成两个并行分支：一条用于更新Chroma/LlamaIndex（可能需要写文件或调用embedding存储API），另一条用于写审计DB，两者独立执行提高吞吐。在LangGraph的图中，这些可以作为FileMover后的两个子节点。如果任一失败，LangGraph可以汇总结果，在流程结尾引出一个通知或错误汇报节点。**审计与回滚机制：** 一旦日志写入SQLite，每条记录都有唯一ID，用户若发起回滚，可根据该ID查到原始位置等信息，再调用FileMover执行逆向移动（从分类目录移回原位，并删掉创建的链接）。LangGraph可以为回滚专门设计一个分支或者独立工作流，读取日志记录批量撤销。本模块使系统具备**自我反省**能力，每整理一步都有迹可循、可恢复。

## 标签机制与多视图设计

为支持文件的交叉分类，本系统引入**标签(Tag)**概念，将传统单一层级的文件归类扩展为多维度归类。具体机制和多视图呈现方式如下：

- **标签与类别：** 在本系统中，**类别**可视为一种特殊的标签。每个文件可以拥有一个或多个标签，标签可能来自不同分类维度。例如，设置两个维度的标签：“项目”和“文档类型”，一个文件可以同时标记为「项目A」（项目维度）和「合同」（类型维度）。用户可以在规则配置中定义标签集合（taxonomy）。若不区分维度，也可以将所有标签放在同一级，系统仍支持多选。但建议对不同性质的分类建立独立标签组，以利于管理（如 topic vs. document type vs. year）。
- **主标签与主存储：** 鉴于文件系统实体只能有一个物理存储位置，系统需要确定一个**主标签**来决定文件实际所在的文件夹。主标签通常对应文件最主要的归属类别。确定策略可以是：预先由用户指定某个标签维度为主（例如所有文件主目录按“项目”分类，其它标签作为辅助）；或由系统根据标签重要性/置信度选一个主标签。例如如果分类判定的多个标签有不同置信度，则选择最高者为主。同时考虑用户配置：某些标签被标记为“仅标签，不移动”（例如“重要”标签不改变文件主要位置，只作为标记）。一旦选定主标签，文件移动至对应的主标签文件夹下存放。
- **软链接与快捷方式 (多视图)：** 为实现其他标签的可见分类，系统在每个次要标签对应目录下创建指向主存储文件的引用：
  - 在 **Windows** 上，创建扩展名为 `.lnk` 的快捷方式文件，命名与原文件一致，并放置在次要标签的文件夹中。用户双击快捷方式将打开实际文件。
  - 在 **macOS/Linux** 上，使用符号链接 (`ln -s`) 指向目标文件。Finder中会将别名显示为类似原文件的图标。或者利用macOS Alias（.alias文件）机制创建效果更好的跨卷链接（OneDrive本地通常是同一卷，符号链接足矣）。
  - **注意：** OneDrive对符号链接的支持有限，符号链接文件本身会同步，但在云端可能不会被解析。因此，这种多视图主要用于本地使用体验。如果要求跨设备看到相同分类，可考虑使用 OneDrive 内的“**添加到快捷方式**”功能（但这更多是在云端完成，自动化困难）。本地多视图仍然带来管理便利，但用户需了解其他设备可能只能看到快捷方式文件而非真正的多副本。
  - 对于**虚拟集合**的需求，系统也可以在应用层实现：如提供一个UI或命令展示某标签下所有文件列表（通过查询数据库），而不依赖文件系统显示。比如生成一个HTML报告或交互CLI命令 `ods list --tag 财务` 来列出所有财务标签文件。这属于虚拟视图，用数据驱动而非文件夹结构。
- 系统支持可选启用/禁用实际软链接的创建。如果用户更关注检索查询，可以关闭物理链接，仅通过搜索功能获取多标签结果。



- **标签一致性维护**：当文件移动或重命名后，相关的快捷方式需同步更新指向。如果文件被重新分类（主标签改变），系统应：

- 删除旧主标签目录下的文件（或转为链接）并将实际文件移到新主标签目录。
- 更新所有受影响的链接：旧次要标签的链接删除，新次要标签链接创建。
- 在数据库中更新文件的标签集合记录。
- （可选）将上述变更记录到日志，以便审计和将来回滚。

LangGraph 可以通过触发“标签更新”子流程来处理这些维护工作。例如，当分类标签列表与数据库中上次记录不同时，引发UpdateTagViews操作，对差异部分进行调整，保持标签多视图与实际标签数据一致。

- **多视图示例**：假设有文件 **Proposal.pdf** 被标记标签：项目A，招标，2023。主标签规则设为项目优先，则文件实际路径可能为 `OneDrive/分类/项目A/Proposal.pdf`。然后系统在：
  - `OneDrive/分类/招标/` 下创建指向Proposal.pdf的链接；
  - `OneDrive/分类/2023/` 下也创建一份指向Proposal.pdf的链接（或者2023可以被视为日期分类，不一定实现为真实标签目录，而可能通过文件日期属性查询）。这样，从“项目A”文件夹能找到Proposal，从“招标”聚合文件夹也能找到它。用户若在“招标”目录打开文件编辑并保存，实际修改的是主存储文件；OneDrive会同步主文件，链接本身无需同步内容（链接文件本身可能被视为一个几KB的文件同步，但不重复保存主体）。

通过上述标签机制，系统实现了一对多的文件呈现：每个文件物理上唯一，但可以通过多个分类路径访问。多视图设计解决了传统文件夹分类中“此文件应该放在哪个文件夹”两难的问题，用户不必纠结归类路径，因为无论放哪，都能通过其它视图找到。同时，也避免了在多个文件夹保存文件副本造成的冗余和不同步风险。配合数据库索引的虚拟查询视图，用户可以轻松按标签筛选文件，真正做到**交叉分类、灵活浏览**。

## 5. 详细分阶段实现计划

为确保项目顺利实施，采用迭代开发方式，分阶段逐步交付功能。以下列出各阶段的计划、侧重点和产出内容，包括主要模块、接口工具、配置示例以及关键策略。

### 阶段1：MVP 实现

**目标**：构建最小可用产品，实现基本的自动分类整理功能。在此阶段，聚焦单一标签的分类（每文件一个类别），确保核心流程跑通，包括文件解析、简单分类、移动和重命名，初步具备日志记录和CLI操作。

**涵盖功能模块**：文件夹监听、内容抽取、嵌入生成、自动分类（单标签）、路径决策、自动命名（基础版）、文件移动、索引写入和日志。交互修正和复杂规则暂不实现或简化。

**关键实现要点**：

- **类别体系**：MVP预设一组基础类别，可根据用户主要文件类型定制。例如提供默认类别：`["工作", "个人", "财务", "其他"]`。这些可在初始化时写入配置文件供用户修改。系统假定每个文件**只能归入其中一个类别**（无交叉）。分类判定采取**简单策略**：利用 LLM 或embedding 最近邻来从这几个类别中挑选最佳匹配。
- **分类判定实现**：由于初始缺少已分类语料，MVP的分类主要依赖 LLM 对内容的直观判断。实现方式：
  - 构造一个提示模板请 LLM 直接分类，例如：

提示：请阅读以下文本内容，并判断它属于以下哪种类别：工作文档 / 个人文档 / 财务票据 / 其他。文本内容："{content\_excerpt}"。只需给出类别名称。

LLM 返回一个类别名称。系统将其映射到预设类别集合验证有效性。

- 或采用关键词匹配辅助：配置每类一些关键词，解析文本后计算出现频率作为参考。如“报销”“发票”命中则偏向财务类别。
- 嵌入模型的作用在MVP可选：如有足够参考文本，可计算新文档向量与每类示例向量的余弦相似度，选取最大者作为分类。
- 无论哪种方法，输出一个确定的类别。如果模型信心不足，可以输出“其他”类别表示未明。
- **命名与路径**：默认将文件移动至 `OneDrive/分类/{类别}` 文件夹下。命名规则在此阶段可以很简单，例如：
  - 若原文件名包含有意义的词则保留主要部分，只是添加类别前缀或后缀。
  - 或干脆保持原文件名不变（MVP阶段可不强制改名，防止过度变化干扰用户）。可以提供一项选项控制是否重命名，默认关闭自动命名，仅在文件名过于杂乱（如 `Document (1).pdf` 这类）时启用。
  - 如果启用，使用一个简单模板，例如 `"{{category}}-{{orig_name}}"` 或 `"{{category}}_{{date}}.{{ext}}"`。例如“财务-IMG001.pdf”或“工作\_20230801.docx”。这样在同一目录下按名称也可排序归类。
- **CLI 工具**：开发基本的命令行界面：
  - `ods init`：初始化配置和环境。具体行为：
    - 创建默认配置文件 `rules.yaml`（其中包含默认类别列表和简单规则模板）。例如：

```
categories: ["工作", "个人", "财务", "其他"]
naming:
  template: "{{category}}/{orig_name}"
```

（MVP阶段 `rules.yaml` 结构简单，详细规则引擎留待后续）。

- 初始化 SQLite 数据库文件（例如 `.ods/db.sqlite`），建立所需表（文件表、日志表等）。
- 如首次运行，可扫描OneDrive根目录已有子文件夹名映射到默认类别（比如发现有“财务”文件夹则将“财务”作为分类名）。
- 下载或加载所需模型（embedding 模型文件，确保LLM API可用性等）。
- 输出提示让用户检查/编辑规则配置。
- `ods apply`：执行一次分类整理：
  - 可选参数：`--dry-run` 模式下只模拟输出变更，不实际移动，用于审阅。
  - 扫描目标目录（或数据库中记录的自上次运行后的新增文件）列表，将每个文件送入处理流水线。
  - 显示处理进度，并在结束后输出汇总，如“共处理X个文件，其中自动分类Y个，其它Z个未确定”。
  - 在非 `dry-run` 模式，执行实际文件移动，并在控制台列出变动结果，例如：

```
[Moved] 报告.doc -> 分类/工作/报告.doc
[Renamed] Document(1).pdf -> 分类/财务/财务-银行账单.pdf
[Uncategorized] draft.txt -> 待整理/
```

其中标注未分类的文件留在“待整理”或报告用户。

- `ods apply` 结束后也会提示用户如有任何错误或需要人工确认的文件。

- **日志和审计**：MVP实现基本日志：
- 在 SQLite 的日志表中记录每次 `ods apply` 处理过的文件及动作。或简单起见，写入一个 `ods.log` 文本文件追加记录。格式例如：

```
2025-08-18 09:15: File "报告.doc" -> Category: 工作 (Auto)
2025-08-18 09:16: File "Document(1).pdf" -> Category: 财务 (Auto-renamed to
"财务-银行账单.pdf")
2025-08-18 09:16: File "draft.txt" -> Category: 未确定 (Skipped)
```

这样用户可以事后查看发生了什么变更。

- 回滚在MVP阶段不会实现完整自动操作，但通过日志用户至少知道原路径，可手动恢复。`ods undo` 命令可作为未来功能占位，目前可以提示“此版本暂不支持自动回滚，请手动移动需要恢复的文件”。
- **错误处理**：MVP注重主要路径，异常情况如解析失败、移动失败简单记录错误，不中断整个流程。例如无法解析的文件标记为“未处理”并跳过。
- **用户反馈**：MVP中交互 minimal：`ods apply` 默认非交互全自动，对于低信心分类统一作为“其他”处理。用户反馈通过查看报告后手动调整：用户可以编辑配置（例如在 `rules.yaml` 把某文件名加入规则指定类别）然后重跑 `ods apply` 修正，或直接手工移动文件到正确目录（之后系统可识别下次不再移动这个文件）。这些在MVP未自动检测，但在文档中指引用户这样做。

**阶段产出**：- 一个基本可运行的 `ods` 命令行工具，具备初始化和应用功能。- `rules.yaml` 默认配置示例:

```
categories:
  - 工作
  - 个人
  - 财务
  - 其他
naming:
  template: "{{category}}-{{original_name}}"
confidence_threshold:
  auto: 0.0    # MVP 阶段暂不使用置信度门槛，全部尝试自动分类
  review: 0.0
```

（以上阈值0表示不启用复核机制，一律自动分类）。- 系统说明文档，指导用户如何运行 init，检查/修改配置，然后 apply。- 该版本实现了**基本整理功能**：用户可以看到OneDrive下出现新的分类文件夹，文件被移动归类，显著减少了根目录的混乱。

## 阶段2：多标签支持与改进

**目标**：扩展系统为多标签分类框架，引入交叉分类功能。完善交互和部分规则支持，提高分类精准度。此阶段后，用户可为文件赋予多个标签，系统通过软链接提供多视图访问。同时引入用户复核机制以处理不确定情形。

**新增/改进功能模块**：分类判定模块升级（支持输出多个标签），路径决策和移动模块升级（处理多标签链接），标签一致性维护模块，新增交互反馈节点。部分简单规则支持（如基于关键词的自动标签）。

## 关键实现要点：

- **标签体系配置：** 修改规则配置结构，支持定义标签集合而非单一类别列表。可以让用户在 `rules.yaml` 中配置多个 tag 组。例如：

```
taxonomies:  
  主类别:  
    - 工作  
    - 个人  
    - 财务  
    - 其他  
  文档类型:  
    - 合同  
    - 发票  
    - 报告  
    - 照片
```

这表示存在两个维度标签。如果用户不需要明确维度，也可都放在一个列表里，但内部实现仍能支持区分。系统需更新解析该配置的逻辑，并在分类输出时相应标记标签所属类别组。

- **分类判定改进：** 让 LLM 可以输出多标签：

- Prompt 模板调整，如：“如果文档内容涉及多个类别，请给出所有相关标签。” LLM可返回一个列表。例如“财务, 报告”表示主类别财务，文档类型为报告。
- 利用向量相似度也可产生多建议：对于每个标签组分别决策。例如先判断主类别，再判断文档类型（这可通过为不同taxonomy各训练一个分类模型或提示分别进行）。
- 为控制标签数量，可设定最多标签数或最低次要标签置信度。例如仅当次要标签的置信度  $\geq$  主标签置信度的某比例时才包含，否则忽略。这样避免每份文档挂太多标签。
- 例如新文档embedding与“项目A”“项目B”都较接近，可全部考虑。如果规则配置某些标签互斥（比如“个人”与“工作”应该只能择一），则在LLM结果出来后强制互斥约束：选择分值更高的那个，另一个丢弃。
- 分类输出的数据结构从之前单一category字符串变为数组tags[]。同时记录哪个是主标签（可按顺序或额外标注）。

- **多视图实现：**

- 在 PathPlanner 中，根据主标签生成primary\_path，同步生成每个次要标签的link路径。OneDrive分类目录结构可能调整：例如所有主类别仍在 `分类/主类别名/`，而其他标签维度可以在 `分类/标签组名/标签名/`，也可能平铺。设计上，为避免混乱，可以将不同taxonomy放在不同顶层文件夹：
  - 如 `OneDrive/分类/主类别/工作/...`，`OneDrive/分类/文档类型/报告/...` 等。这需要rules.yaml提供每个taxonomy对应的文件夹路径前缀或模式。
  - 如果不想多层，可以也平铺，但要区分不同来源标签可能重名（比如有“项目A”既是主类别又是可能出现在其他组，区分较复杂）。因此建议以taxonomy区分路径。

- 移动执行FileMover需要支持创建多个链接。MVP里FileMover可能只移动文件，这里扩展FileMover或在其后增加Linker子节点：
  - Linker遍历link\_paths列表：逐个创建symlink/shortcut文件。如果失败，记录在日志但不影响主文件已移动结果。
  - 数据库也应存储每个文件的tag集合，方便重建链接或查询。
- 增加同步维护工具：为防止链接不同步，提供CLI命令比如ods sync-links，根据数据库重新生成所有标签链接视图，清理无效链接。每次大改规则或手动调整后可运行一次保证一致性。
- 处理删除情况：若用户删除了某文件或取消某标签，需要相应删掉过期的链接。这可以在下一次ods apply时检测（比如数据库有记录但文件不在了，则删链接）。

#### • 交互复核机制：

- 引入分类置信度指标或者根据模型返回的logprob/评分。由于LLM未必给概率，这里可以继续利用向量相似度或内置启发式。例如：
  - 若LLM输出“类别X”和“类别Y”两个标签，可能对顺序有暗示主次，但不定。可以结合embedding：计算文档向量与X类中心的距离和与其他类的距离，估算个置信度。
  - 或让LLM在输出中附加一个自信程度说明（需要解析）。
- 设置阈值：引入配置，比如：

```
confidence_threshold:
  auto: 0.8
  review: 0.5
```

含义：当分类决策置信度  $\geq 0.8$  时自动执行；介于0.5-0.8时标记需要复核； $< 0.5$ 时视为无法确定，归入“待分类”。

- LangGraph实现：Classifier节点输出结果后，增加一个判断节点：
  - 如果  $\geq$  auto 阈值：沿主路径执行移动等。
  - 如果  $<$  auto 但  $\geq$  review：走人工确认节点。此节点可以是一个等待用户输入的CLI对话：列出文件摘要和模型建议的标签，用户可按编号确认或选择另一个标签。实现上，可暂停流程（或将该文件处理加入待审核列表，不立即整理）。
  - 如果  $<$  review：直接将文件移动到未分类文件夹或者标记跳过，由用户稍后处理。
- CLI改进：提供ods review命令，列出待人工确认的文件及建议，允许用户交互处理。例如交互流程：

```
待审核文件1: "notes.txt", 模型建议标签: 工作? (置信0.6)
-> 选项: [Y]接受 [N]选择其他类别 [S]跳过
```

用户选N则列出所有类别让其选，或输入自定义标签。用户处理完所有待审项后，系统应用他们的决定（移动文件）。

- 交互也可在ods apply加--interactive时即时进行，逐个文件询问，但大量文件可能不便，故默认将模糊结果累积后统一交互更好。

- 规则引擎初步：阶段2可引入简单规则来辅助分类：

- 例如，引入 keywords->tag 规则：用户可在 rules.yaml 增加：``yaml rules:
  - if: "文件名包含 发票" add\_tag: "财务"
  - if: "content ~ /项目\s+\d+/" category: "工作" `` 解析：第一条表示如果文件名含“发票”，则添加“财务”标签。第二条示例使用正则内容找“项目 数字”模式，则直接分类为“工作”类。
- 规则处理顺序：可以在Classifier之前预处理，如满足规则则**锁定**某些标签或提高相应置信，甚至跳过LLM直接分类确定。这对明显的情况减少开销。
- 阶段2规则功能有限度开放，只支持简单条件，以防复杂性。高级复合规则、模板运算等留在阶段3。
- 在CLI上，可有 `ods test-rule "<filename>"` 命令，输出该文件会触发的规则，以便用户调试规则配置是否生效。

#### • 用户反馈学习：

- 阶段2在交互中收集的用户更改，除了立即作用于当前文件，也用于**学习**：
  - 若用户改正了一个错误分类，比如文件X模型标为A但用户选了B，则系统可以记录一个反馈条目。下次Classifier对类似文件时，可参考：例如将该文件向量加入类别B样本，或生成规则“如果包含X的某特征则偏向B”。
  - 暂无自动调整模型参数，但可以更新类别向量库或规则库，让此反馈影响后续行为。比如Chroma库中本来也存了X的embedding（错误标A），我们应修正其标签为B，避免误导下一次检索。或者在日志标记这是一例分类错误，防止用于指导LLM决策。
- 若多个文件发生类似错误，项目后续阶段也许考虑微调embedding模型或LLM的Few-shot提示加入这些案例。在阶段2，可以主要采用**规则**形式消化反馈：由开发者或高级用户根据积累的错误模式，补充规则避免重复犯错。

**阶段产出：** - 更新的配置文件结构和示例（支持多标签、基本规则）。例如修改后的 `rules.yaml` 示意：

```
taxonomies:
  主类别:
    - 工作
    - 个人
    - 财务
    - 其他
  文档类型:
    - 报告
    - 合同
    - 发票
    - 图片
  naming:
    template: "{{主类别}}/{{title|default(original_name)}}.{{ext}}"
    rules:
      - match_taxonomy: 文档类型
        value: 发票
        template: "{{主类别}}/财务/{{category}}-{{date|strftime('%Y%m%d')}}.{{ext}}"
        # 对于发票类文件，放入财务目录下，命名为 类别-日期
  rules:
    - if_filename: "发票"
```

```

    add_tag: "发票"
- if_content_regex: "CONFIDENTIAL"
    add_tag: "保密"
- if_tag_combo: ["保密", "财务"]
    action: "require_review" # 如果同时触发保密和财务标签，则标记需要人工审核
confidence_threshold:
    auto: 0.85
    review: 0.6

```

（以上仅为示例，真实格式会详细定义条件语法。这展示了多taxonomy和一些规则。） - CLI 工具增强： - `ods apply` 默认开启根据阈值的分类结果处理，未决的进入待审列表。 - 新命令 `ods review` 供用户处理待审文件。如果实现复杂，可暂时让 `ods apply` 在结尾自动调用审核流程。 - 或增加 `ods watch`，启动后台监控模式，不断应用分类（需注意与交互配合，一般监控模式下将低置信文件直接跳过或按策略处理，不可能弹交互）。 - 文件系统预期变化：除了前述 分类/主类别/... 目录，现在会出现其他标签目录，例如 分类/文档类型/合同/... 里面全是指向实际文件的快捷方式或链接。用户能够从多个入口找到文件验证多标签效果。 - 文档更新：提供使用指南，解释多标签的工作方式和注意事项（例如在其他电脑上看到快捷方式的情况）。

### 阶段3：规则引擎扩展与高级功能

**目标：** 引入完善的规则引擎机制，提高系统可定制性和自动化处理复杂情形的能力。同时增强可靠性和用户信心，包括全面的审计回滚支持。此阶段完成后，系统应当允许用户通过配置灵活定义分类和命名逻辑，并能安全地处理多样化场景。

**扩展功能模块：** 规则引擎模块（解析和执行复杂规则）、高级命名模板、冲突处理策略、完整的审计回滚实现、以及潜在的界面优化。

**关键实现要点：**

- **规则引擎设计：**
- 采用面向对象或清晰的数据驱动方式实现规则解析执行。可将每条规则建模为 Rule 对象，包含触发条件（Condition）和行动（Action）。系统在分类流程中某一节点调用规则引擎，顺序评估规则列表，对于匹配的规则执行相应操作。需要明确执行时机：有的规则适合分类前（如根据文件名/路径初步决定分类），有的适合分类后（如对LLM结果做修改或对已分类结果追加标签），甚至在移动后（如针对路径调整）。规则可按阶段分类，在配置中加入字段表示应用阶段。
- **条件类型：** 扩展支持更复杂的匹配：
  - 文件名/路径模式（通配符、正则表达式）。
  - 文档内容包含词汇/短语、正则（可设置大小写敏感、全文/标题等范围）。
  - 文档元数据（大小、创建日期、扩展名等）。
  - 分类结果（标签组合、置信度范围）。例如 `if: tag=="财务" and confidence<0.7 then require_review`。
  - 上下文条件：如“当前时间”（可用于定期分类的特殊处理）等。
- **行动类型：** 除了前阶段提到的 `add_tag`、`set_category`，还有：
  - `force_category`：强制将文件归为某类别（跳过LLM）。
  - `exclude`：跳过此文件不处理（例如规则指定某文件夹不自动整理）。
  - `move_to`：覆写目标路径（比如特定文件类型一律移动到某固定文件夹备份）。

- `rename_template`：为特定条件指定不同的命名模板。
- `require_review`：标记此文件必须人工确认（无论置信度）。
- `notify`：仅提示用户某事件（比如文件太大不处理）。
- **规则冲突与优先级**：当多条规则作用在同一文件时，需确定优先顺序：
  - 在配置中可引入 `priority` 字段或者依照顺序先后。高优先规则可以设置跳出后续规则检查（比如一个exclude规则匹配，则后续不再管这个文件）。
  - 或根据规则性质：例如explicit分类规则的优先级高于LLM输出调整规则。
  - 用户文档需说明规则评估顺序，以免混淆。
- **配置结构**：可能将规则单独列出，或按taxonomy/phase组织。可行的结构如：

```
rules:
  pre_classification:
    - condition: file_ext in ["mp3", "wav"]
      action: exclude
    - condition: filename regex "(?i)CV"
      action: set_tag: "简历"
  post_classification:
    - condition: tag=="保密"
      action: move_to: "OneDrive/分类/保密文档/"
    - condition: tags.contains("财务") and tags.contains("保密")
      action: require_review
```

这样将规则按阶段分组，提高可读性。

- **Jinja2 在规则中的应用**：可用于复杂条件或动态动作。
  - 例如 condition 用 Jinja2 表达式：`"{{ '重要' in content and size>1048576 }}"` 来判断内容和大小。
  - 或 action 中构造路径：`move_to: "{{OneDrive}}/Archive/{{year}}/"` 等。引擎需要安全地渲染这些模板，提供上下文变量。
  - 这要求谨慎处理，因用户写的模板可能出错，需要在 `ods init` 或专门的 `ods lint` 检查规则合法性。

#### • 高级命名模板：

- 在阶段2我们已有每个taxonomy的模板支持，这里更灵活：
  - 支持按照**标签组合**定义命名：如同时有“保密”标签则在文件名末尾加 “[Confidential]”字样。
  - 支持条件逻辑：比如大小写转换、截断过长标题。
  - 模板变量扩展：例如自动编号（若同一类别当天已有多份，则文件名加序号001等），或者在命名时访问文件元数据字段（作者、ID）。
  - 甚至调用自定义Python函数处理（这可通过Jinja2 filter机制或预置一些filter实现）。
- 这些高级功能主要供高级用户，默认模板仍简单为主防止出错。

#### • 冲突处理逻辑完善：



- **分类冲突**：指多标签系统中标签互斥或规则与模型冲突。
  - 例如规则要求文件A分类为X但LLM判断为Y。如果相信规则，直接用X；如果规则被设置为soft的，可以标记冲突状态，在日志提示“模型建议Y但规则偏好X，已采用X”。可让用户配置某规则是否严格。
  - **Tag互斥**：在配置中提供互斥组定义，如 `mutually_exclusive: [["个人", "工作"], ["公开", "机密"]]`。Classifier后检查，如果同时出现互斥标签，则按照优先顺序保留一个，其余移除并日志记录。也可设为需要人工决策的冲突，将文件放待审核，提示用户“文件A同时被标记为个人和工作，请选择其一”。
- **文件名/路径冲突**：更完善的处理：
  - 引入一个记录结构保存某目录下已存在文件名清单（或每次在移动时现查），发现冲突先尝试按规则处理。如配置可指定冲突策略：“rename\_new”（给新文件改名）、“rename\_old”或“skip”。默认采用rename\_new加序号。
  - 对于链接名称冲突（极少发生，除非不同文件在次要标签目录重名），也可加区分比如加文件ID或主标签在链接名字里区分。
- **版本冲突**：OneDrive自身有版本控制，若系统移动文件时与用户同步操作冲突，可能出现同步冲突副本。无法完全避免，但通过在安静状态下运行 apply 降低发生概率。如检测到OneDrive正在大量同步，可推迟执行。
- **性能冲突**：当大量文件一起分类时，系统资源占用与OneDrive同步带宽要权衡。可配置每分钟移动文件数上限，或者分批处理。例如每运行100个文件等待几秒，让OneDrive消化，防止过载。这些作为可选优化参数提供。
- **审计与回滚实现**：
  - 审计数据库结构在前阶段已有，这里完善 **ods undo** 等操作：
    - `ods log`：输出最近的操作日志，便于用户查看编号或细节。
    - `ods undo [N]`：回滚最近N步操作或指定某一次操作。实现逻辑：
      - 从日志中找到目标记录（或多条）。
      - 按相反顺序执行逆向操作：将文件从现位置移回原路径；若原路径已被占用（例如用户在那里放了别的文件），则需要提醒无法回滚或者备份现有文件后再回滚。
      - 删除回滚涉及的链接文件。
      - 更新数据库：移除那次记录或标记为已回滚，更新文件当前状态记录。
      - 支持批量回滚：如用户想撤销整个 `ods apply` 运行，可以按session id记录每次apply的操作集合，一键撤销该session。
      - 错误处理：若某些文件已不存在（用户可能手动删了），跳过并警告。回滚尽量幂等，失败则不要破坏其他文件。
  - **快照方案（可选）**：为安全起见，可在每次批量整理前执行一次文件结构快照（比如记录当前所有文件路径->分类的映射），存在数据库。回滚时若日志不全，可用快照做全面还原。但记录快照开销大，只对关键里程碑或用户主动要求时做即可。
- **用户界面改进**：
  - 尽管 CLI 已能用，或许提供一个简单图形/网页界面更直观：
    - 可选阶段3开发一个基于例如 Electron 或 Textual (Python TUI) 的界面，列出分类结果，让用户拖拽调整或勾选确认，所见即所得地配置规则。

- 这属于锦上添花，如果资源允许可做，否则CLI已满足基本功能。
- Documentation和帮助的改进，在CLI加 `--help` 详述命令和配置格式说明。
- **示例配置和模板**：提供更多示例让用户参考：
  - 如常见职场使用场景的分类模板、家庭相册场景的分类模板等等，在文档附录给出。
  - 规则.yaml 可以有注释说明每段作用，帮助用户二次开发。

**阶段产出：- 最终版 rules.yaml 模板**：包含前述各种配置点的示例和注释。比如：

```
taxonomies:
  项目: ["项目A", "项目B", "项目C"]
  文件类型: ["报告", "计划", "会议记录", "发票", "合同", "照片"]
  敏感级别: ["公开", "内部", "机密"]
naming:
  template: "{{项目}}/{{文件类型}}/{{title|default('未命名')}}({{date|
strftime('%Y%m%d')}}).{{ext}}"
  # 默认命名：项目/文件类型/标题(日期).扩展名
rules:
  - condition: "敏感级别 == '机密'"
    template: "{{项目}}/{{文件类型}}/[机密]{{title}}.{{ext}}"
  - condition: "文件类型 == '照片'"
    template: "图库/{{year}}/{{month}}/{{original_name}}"
rules:
  pre_classification:
    - condition: "ext in ['exe','tmp','log']"
      action: exclude # 排除无需分类的文件类型
    - condition: "path.startswith('OneDrive/相册')"
      action: set_tag: ["照片"] # 相册文件夹下的都标记为照片
  post_classification:
    - condition: "tags.contains('机密')"
      action: move_to: "OneDrive/分类/机密文件/" # 所有机密文件集中存放
    - condition: "tags.contains('财务') and year < 2020"
      action: set_tag: ["归档"] # 2020年前的财务文件加归档标签
    - condition: "tags.count > 1 and not tags.contains('机密')"
      action: require_review # 凡是跨两个以上标签且不涉及机密的，提示人工确认
confidence_threshold:
  auto: 0.9
  review: 0.75
```

上述配置展示了多维度标签定义、高级命名模式、不同阶段的规则和复杂条件，以及调整了置信度阈值。实际文件中会有更严格的语法和说明。

- **完善的 CLI**：除 `ods init/apply/review` 之外，增加 `ods undo`, `ods log`, `ods config` (打开配置文件供编辑), `ods lint` (检查配置正确性)等命令。 `ods help` 提供各命令用法。可能还有 `ods simulate <file>` 用于测试特定文件的分类决定（输出各模块推理信息）。
- **稳定可靠的系统**：

经过阶段3，系统在各方面达成熟：- **准确性**：借助规则+LLM+向量综合，以及用户反馈迭代，分类准确率显著提高，对特殊情况（如极端多义文件）能有妥善处理（比如标记人工决策）。- **性能**：可处理大型OneDrive目录，通过批处理和高效检索支撑上万文件规模。关键操作（如embedding计算）可在初始化阶段预处理缓存，加快增量运行。- **安全性**：日志完备，回滚方便，用户可以放心地让工具自动整理而无后顾之忧，因为一切改变都可追溯和恢复。- **可扩展性**：新增加文件类型（比如支持Email.msg解析）或适应其他存储（比如Google Drive）只需更换监听源，其余组件可重用。规则引擎让系统易于定制到不同领域，无需修改代码。

总结而言，经过三个阶段的迭代，本基于 LLM+向量数据库的自动文档分类系统将实现从**基础可用**到**多标签智能**再到**高度可定制可靠**的演进。它将大大减轻用户管理云端文档的负担，实现真正的智能归档。同时通过稳健的设计和用户掌控机制，确保过程透明可控，符合实际应用需求。

---

1 BAAI/bge-large-en-v1.5 - Hugging Face

<https://huggingface.co/BAAI/bge-large-en-v1.5>

2 E5 Text Embedding – Vertex AI - Google Cloud Console

<https://console.cloud.google.com/vertex-ai/publishers/intfloat/model-garden/e5>

3 4 6 7 Automate Document Classification in Azure - Azure Architecture Center | Microsoft Learn

<https://learn.microsoft.com/en-us/azure/architecture/ai-ml/architecture/automate-document-classification-durable-functions>

5 【LangGraph 入门指南】为智能代理打造灵活可控的工作流框架

<https://aillm.csdn.net/6836d03d965a29319f242d42.html>