

构建先进知识管理系统的RAG系统最新进展与技术深度解析

第一部分: RAG系统概述与核心价值 (RAG System Overview and Core Value)

1.1 RAG系统定义与基本原理 (Definition and Basic Principles of RAG Systems)

检索增强生成 (Retrieval-Augmented Generation, RAG) 系统代表了大型语言模型 (Large Language Models, LLM) 应用领域的一项重要进展。其核心定义在于, RAG是一种旨在优化LLM输出的方法, 它使LLM在生成响应之前, 能够参考其自身训练数据源之外的、通常是更具权威性 or 实时性的知识库¹。本质上, RAG框架巧妙地融合了传统信息检索系统(例如搜索引擎和数据库技术)的精准定位能力与现代生成式LLM强大的自然语言理解和内容创造能力³。

RAG系统的基本工作原理可以概括为几个关键步骤: 首先, 系统接收用户的输入查询; 接着, 利用该查询通过检索模块从外部数据源(即知识库)中定位并提取出最相关的信息片段; 最后, 将用户的原始查询与这些检索到的信息片段整合后, 一同作为增强提示(augmented prompt)输入给LLM。LLM在充分理解这一增强上下文的基础上, 结合其自身的参数化知识(训练数据中学习到的知识)与新获取的外部知识, 生成一个信息更丰富、事实更准确、上下文更相关的响应²。

这一机制的重要性在于, 它直接解决了LLM固有的一些核心局限。LLM的知识库通常是静态的, 受限于其训练截止日期, 难以获取最新的信息或特定领域的细微知识。RAG通过动态地从外部知识库检索信息, 为LLM注入了实时性和领域特异性, 这对于构建需要提供即时、准确信息的知识管理系统而言至关重要。可以认为, RAG的本质是一种动态的知识注入机制。与通过成本高昂且耗时的微调(fine-tuning)来更新LLM的静态知识不同, RAG在模型推理(inference)阶段动态地为LLM提供与当前查询紧密相关的上下文知识¹。这意味着知识管理系统可以依赖RAG来提供最新的信息, 而无需频繁地对核心LLM进行重新训练, 从而显著降低了维护成本和响应延迟。

此外, RAG系统为提升LLM的“可解释性”和“可控性”提供了一条有效途径。传统LLM的决策过程往往如同一个“黑箱”, 其生成内容的依据难以追溯。RAG通过明确指出其响应所依据的外部知识来源(例如, 引用的文档片段或数据条目)², 极大地提高了生成内容的可追溯性和可信度。用户或系统管理员可以方便地验证信息的来源, 这对于知识管理系统中的事实核查、错误修正以及建立用户信任具有不可估量的价值。当LLM的回答出现偏差时, RAG提供的溯源能力有助于快速定位问题是源于检索质量不高还是LLM的理解偏差, 从而实现更精准的系统调优。

1.2 RAG在知识管理系统中的核心作用与优势 (Core Role and Advantages of RAG in Knowledge Management Systems)

在知识管理领域, RAG系统扮演着至关重要的角色, 它通过将强大的LLM与组织内部的专

有信息资源(如文档库、数据库、内部网站等)有效连接,能够将原本分散、孤立的机构知识转化为一个易于访问、可供智能分析和利用的动态情报层,从而在各个层面驱动更明智的决策制定⁶。

RAG在知识管理系统中的核心优势体现在以下几个方面:

- 获取最新与权威信息: LLM的内部知识库受其训练数据的时效性限制,可能无法反映最新的发展或组织内部的即时更新。RAG通过连接到实时更新的外部或内部知识库,确保系统能够提供当前最准确的信息³。这对于需要处理动态信息的知识管理应用(如政策更新、市场动态分析)尤为关键。
- 增强事实基础,显著减少幻觉: LLM在生成内容时有时会产生“幻觉”,即编造不准确或无事实根据的信息。RAG通过为LLM提供相关的、真实的上下文信息作为生成依据,能够有效地将模型的输出“锚定”在事实上,从而大幅减少幻觉现象的发生³。对于知识管理系统而言,提供准确可靠的信息是其核心使命,RAG在这方面提供了强大的技术保障。
- 提升成本效益: 传统的LLM知识更新方式(如完全重新训练或大规模微调)成本高昂且周期漫长。RAG允许组织充分利用其现有的数据资产和知识库,通过检索增强的方式提升LLM在特定领域的表现,而无需对LLM本身进行频繁或大规模的重新训练,这显著降低了AI知识管理系统的开发、部署和长期维护成本²。
- 促进知识发现与高效共享: RAG系统能够理解自然语言查询的深层语义,并从多样化的知识源中整合信息,从而提高信息检索的准确性和相关性。它不仅能找到用户明确提出的问题的答案,还可能揭示信息间的潜在联系,促进新知识的发现。同时,通过LLM将复杂信息综合成易于理解的答案,也极大地促进了组织内部的知识共享效率⁹。
- 实现个性化与上下文相关的知识服务: RAG能够根据用户的具体查询、历史交互记录、角色权限甚至当前任务的上下文,从知识库中检索并生成高度定制化的信息和答案,提供千人千面的知识服务体验⁹。
- 保障知识系统的可扩展性: 随着组织知识量的不断增长,RAG系统可以通过扩展其外部知识库的覆盖范围和深度来提升服务能力,而核心LLM的升级压力相对较小,保证了知识管理系统良好的可扩展性⁹。

从更深层次来看,RAG技术是解决企业长期存在的“知识孤岛”问题的关键赋能技术。企业内部的宝贵知识往往散落在不同的部门、不同的业务系统以及海量的各类文档之中,形成一个个信息壁垒⁶。RAG系统通过其强大的信息检索和整合能力,能够将这些分散的、独立的结构化和非结构化知识源(如数据库、文档库、CRM系统、ERP系统、协同平台等)连接起来,构建成一个统一的、可被自然语言查询的智能知识体系。这不仅打破了部门或信息系统之间的壁垒,还极大地提高了组织整体知识的可访问性和利用效率。

更进一步,RAG的引入标志着知识管理从传统的信息检索模式向知识服务模式的深刻变革。传统的知识管理系统通常基于关键词搜索,返回给用户的是一个相关文档或信息的列表,用户需要自行阅读、筛选、理解和综合这些信息才能得到最终答案。而RAG系统则通

过LLM对检索到的相关信息进行深度的理解、分析、综合、提炼，并最终以自然语言的形式直接生成条理清晰、逻辑严谨的答案或解决方案⁹。这种从“授人以鱼”到“授人以渔”再到直接“端上烹饪好的鱼”的转变，将“信息获取”提升到了“知识服务”的全新层面，直接为用户提供可操作的见解和决策支持，从而极大地提升了用户体验和工作效率。

然而，RAG系统的高效运作也对其所依赖的外部知识库的质量和时效性提出了新的、更高的要求²。虽然RAG降低了对LLM自身知识更新频率的要求，但其性能表现与外部知识库的准确性、完整性、新鲜度以及组织结构的合理性密切相关。因此，一个成功的RAG驱动的知识管理系统，必须辅以一套有效的机制来持续地、动态地更新和维护其核心知识库，包括文档内容的增删改、元数据的标注与管理、嵌入表示的定期刷新等²，以确保RAG系统能源源不断地获取高质量的“燃料”。

1.3 RAG系统基本架构与演进范式 (Basic Architecture and Evolving Paradigms of RAG Systems)

理解RAG系统的基本架构是掌握其更高级技术和应用的前提。一个典型的RAG应用流程，无论其复杂程度如何，通常都包含以下三个核心阶段¹¹：

1. **索引 (Indexing)**: 此阶段的目标是预处理原始数据，并将其转化为RAG系统可以高效检索的格式。具体步骤包括：
 - **数据加载与清洗**: 从各种数据源(如文档、网页、数据库)加载原始数据，并进行必要的清洗(如去除无关字符、格式转换)。
 - **文本分块 (Chunking)**: 将长文本分割成较小的、语义上相对独立的文本块(chunks)。这是因为LLM通常有上下文长度限制，而且更小的块有助于提高检索的精确度。
 - **向量化编码 (Encoding)**: 使用预训练的嵌入模型(embedding model)将每个文本块转换成高维向量表示(vector embeddings)。这些向量能够捕捉文本块的语义信息。
 - **存入向量数据库 (Vector Database Storage)**: 将生成的文本块向量及其对应的原始文本内容(有时还包括元数据)存储在专门的向量数据库中。向量数据库支持高效的基于向量相似性的搜索。
2. **检索 (Retrieval)**: 当用户提出查询时，此阶段负责从索引好的知识库中找出与查询最相关的文本块。
 - **查询向量化**: 使用与索引阶段相同的嵌入模型，将用户的自然语言查询也转换成一个查询向量。
 - **相似性搜索**: 在向量数据库中，计算查询向量与所有已存储文本块向量之间的相似度(常用的有余弦相似度、点积等)。
 - **Top-K选择**: 根据相似度得分，检索出得分最高的K个文本块作为与查询最相关的上下文信息。
3. **生成 (Generation)**: 此阶段利用LLM结合用户的原始查询和检索到的上下文信息来生成最终的答案。

- **提示构建 (Prompt Construction)**: 将用户的原始查询和检索到的Top-K文本块组合成一个结构化的提示(prompt)。提示的设计对生成答案的质量有显著影响。
- **LLM推理**: 将构建好的提示输入给LLM。LLM在理解查询意图和参考所提供上下文的基础上, 生成一个连贯、相关且尽可能准确的答案。LLM在此过程中可以利用其内部的参数化知识, 但更侧重于依据提供的上下文进行回答。

RAG技术并非一成不变, 其研究和应用范式在不断演进, 以克服早期方法的局限性并追求更高的性能。我们可以将RAG的演进大致归纳为以下几个阶段¹¹:

- **Naive RAG (朴素RAG)**: 这是RAG最初的、最基础的形态, 严格遵循上述“索引-检索-生成”的线性流程。它直接将检索到的文本块与用户问题拼接后输入LLM。虽然简单有效, 但Naive RAG常常面临一些挑战, 例如: 检索到的文本块可能包含噪声或不完全相关的内容, 导致检索精度和召回率不高; LLM在整合多个检索片段时可能出现逻辑不连贯或事实性错误(幻觉); 对于复杂问题, 单次检索可能无法提供足够的上下文。
- **Advanced RAG (高级RAG)**: 为了克服Naive RAG的局限性, 研究者们提出了一系列增强技术, 主要集中在优化检索过程和生成过程。这些技术可以分为:
 - **检索前优化 (Pre-retrieval Process Optimization)**: 在检索发生之前对索引或查询进行优化。例如, 改进分块策略(如语义分块、层级分块), 优化索引结构(如引入元数据、构建知识图谱索引), 以及对用户查询进行重写、转换或扩展(如使用HyDE、多查询分解)。
 - **检索后优化 (Post-retrieval Process Optimization)**: 在检索到初步的上下文信息后, 对其进行处理以提升质量, 再送入LLM。关键方法包括对检索到的文本块进行重排序(re-ranking), 将最相关的块放在更突出的位置; 以及进行上下文压缩(context compression), 去除冗余信息, 提取核心内容, 以适应LLM的上下文窗口限制并减少噪声干扰。
- **Modular RAG (模块化RAG)**: 这是当前RAG发展的一个重要趋势, 它将RAG系统视为一个由多个可插拔、可配置的模块组成的灵活框架。Modular RAG不仅继承了Advanced RAG中的诸多优化技术, 更强调系统的适应性和多功能性。它引入了许多新的模块和模式, 例如:
 - **新增功能模块**: 如专门的搜索模块(可直接对接企业搜索引擎或Web搜索引擎)、RAG-Fusion模块(用于融合多个查询策略的结果)、记忆模块(利用LLM的短期或长期记忆能力)、路由模块(根据查询类型将其导向不同的RAG流水线或知识源)、预测模块(让LLM预测可能需要的上下文类型)以及任务适配器(针对特定下游任务优化RAG流程)。
 - **新的执行模式**: Modular RAG不再局限于固定的“检索-阅读”序列, 而是支持更复杂的交互模式, 如“重写-检索-阅读”(先优化查询再检索)、“生成-阅读”(先让LLM基于初步理解生成假设性答案或子问题, 再用其指导检索)、迭代式检索(多次检索, 逐步逼近答案, 如ITER-RETGEN), 以及混合检索策略(结合稀疏检索和稠密检索)。
 - **与其他技术的协同**: Modular RAG更容易与微调(fine-tuning)、强化学习(

reinforcement learning)等其他AI技术结合, 实现端到端的系统优化。

RAG技术的这种演进路径, 清晰地展现了其从一个相对简单的信息增强流程, 逐步发展成为一个日益复杂、智能且具备自适应能力的系统级解决方案的趋势。Naive RAG奠定了基础, Advanced RAG针对其核心痛点进行了专项优化和深度打磨, 而Modular RAG则标志着RAG系统设计理念的进一步成熟。模块化的思想强调了灵活性、可扩展性以及与其他AI技术的深度协同, 这意味着未来的知识管理系统开发者可以根据自身业务场景的具体需求和可用资源, 像搭建乐高积木一样, 灵活地组合和配置不同的RAG模块、策略与技术, 构建出高度定制化的智能知识服务平台。

更值得注意的是, Modular RAG的出现也预示着RAG系统本身正在经历一场“智能化”的深刻变革。例如, 其中引入的路由模块能够基于对用户查询的理解, 智能地选择最合适的处理流程或知识源; 记忆模块则尝试赋予RAG系统利用LLM记忆(无论是短期上下文记忆还是通过特定机制实现的长期记忆)的能力。这些特性表明, RAG系统正从过去相对固定的、被动响应式的流水线作业, 向能够根据具体情境动态调整自身行为、主动进行信息探索与整合的智能系统演进。这对于构建能够高效处理日益多样化、复杂化的知识查询任务, 并提供深度洞察的下一代知识管理系统, 无疑具有里程碑式的意义。

第二部分: RAG系统的核心组件与关键技术 (Core Components and Key Technologies of RAG Systems)

RAG系统的效能高度依赖于其核心组件的精密设计与关键技术的有效运用。本部分将深入剖析RAG系统的两大支柱——智能检索与增强生成, 并探讨其中涉及的关键技术环节。

2.1 智能检索 (Intelligent Retrieval)

智能检索是RAG系统的基石, 其目标是从海量知识库中精准、高效地找出与用户查询最相关的上下文信息。这一过程涉及数据索引、分块、嵌入模型以及查询处理等多个环节。

2.1.1 数据索引与分块策略 (Data Indexing and Chunking Strategies)

数据索引是RAG系统处理知识的第一步。原始文档首先需要被分割成更小的单元, 即文本块(**chunks**)。随后, 这些文本块通过嵌入模型(embedding models)被转换成数值向量(vector embeddings), 这些向量能够捕捉文本的语义信息。最终, 这些向量连同其对应的原始文本块内容及可能的元数据, 被存储在一个专门的**向量数据库(vector database)**中, 以便进行高效的相似性检索¹¹。

****分块(Chunking)****在此过程中扮演着至关重要的角色。它是指将大型文档或数据集分割成更小、更易于管理且语义上相对完整的片段的过程¹³。选择合适的分块策略直接影响RAG系统的整体性能, 包括检索速度、准确性以及最终生成答案的质量¹³。不恰当的分块可能导致重要信息的割裂、上下文的丢失, 或是检索到大量不相关的噪声信息, 这些都会严重

影响LLM的理解和生成效果。

目前业界探索并应用了多种分块策略, 各有其优缺点¹³:

- **固定长度分块 (Fixed-size chunking)**: 这是最简单直接的方法, 将文本按照预先设定的字符数或token数(chunk_size)进行切割。通常会设置一定的重叠部分(chunk_overlap), 即相邻块之间共享一部分内容, 以期在一定程度上保持上下文的连续性¹⁴。
 - 优点: 实现简单, 计算开销小, 处理速度快, 易于控制块的大小。
 - 缺点: 容易在句子中间或语义单元的边界处强行切分, 可能破坏文本的语义完整性, 导致上下文信息丢失。对于结构和密度不均的文本, 效果可能不佳。
- **基于章节/递归分块 (Section-based/Recursive chunking)**: 这类方法试图利用文档的自然结构(如段落、标题、列表项)或通过递归地使用一系列预定义的分隔符(如换行符、句号等)来进行分块。例如, LangChain中的RecursiveCharacterTextSplitter会按优先级尝试不同的分隔符, 以尽可能保持语义单元(如段落、句子)的完整性¹⁴。
 - 优点: 相比固定长度分块, 更能尊重原文的结构和语义, 有助于保持上下文的连贯性。
 - 缺点: 分块大小可能不均匀, 依赖于文档的固有结构。对于缺乏清晰结构标记的文本, 效果可能退化。
- **句子分块 (Sentence chunking)**: 将文本精确地分割成单个句子。每个句子通常被认为是一个相对完整的思想单元。NLTK、spaCy等NLP库以及LangChain中的SpacyTextSplitter都提供了句子分割功能¹⁵。
 - 优点: 能较好地保持每个块的语义完整性, 因为句子是表达意义的基本单位。
 - 缺点: 句子长度变化很大, 导致块大小不一。过短的句子可能导致块过于零碎, 信息量不足; 而过长的复杂句本身可能包含多个语义焦点。
- **语义分块 (Semantic chunking)**: 这是一种更高级的分块策略, 它不仅仅依赖于文本的长度或结构标记, 而是试图根据文本内容的实际语义关系进行划分。例如, 可以通过计算句子或小段文本之间嵌入向量的相似度, 将语义相近的内容聚合为一块, 而在语义发生较大转变的地方进行切分。目标是确保每个块在内部语义上高度相关, 而块与块之间则相对独立¹³。
 - 优点: 能够最大程度地保留每个块的上下文相关性和语义内聚性, 从而提高检索的准确性和相关性, 尤其适合处理语义复杂或主题多变的文档。
 - 缺点: 计算成本较高, 因为它通常需要对文本进行嵌入并计算相似度。实现也相对复杂。
- **文档级分块 (Document-based chunking)**: 在某些情况下, 特别是当文档本身较短或者整个文档的上下文对于回答特定类型问题至关重要时, 可以将整个文档视为一个单独的块, 或者进行非常粗粒度的划分(例如, 按主要章节)。这适用于需要对整个文档进行理解和综合的任务¹⁶。
 - 优点: 完整保留了文档的上下文和结构, 避免了信息割裂。
 - 缺点: 对于长文档, 可能会超出LLM的上下文窗口限制, 或者因为信息密度过低而

影响检索效率和生成效果。

- **智能体分块 (Agentic chunking)**: 这是一种前沿的、实验性的分块方法, 它尝试利用 LLM 自身的智能来决定如何分割文档。LLM 被赋予一定的“自主性”, 通过分析文本的语义内容、逻辑结构(如段落类型、章节标题、步骤说明等), 模拟人类在阅读和理解长文档时的认知过程来进行分块¹³。
 - 优点: 潜力巨大, 理论上能够实现最符合语义和任务需求的分块效果, 特别适合处理结构复杂、内容多样的文档。
 - 缺点: 目前尚处于探索阶段, 实现复杂, 计算成本可能很高, 且其效果的稳定性和可控性有待进一步验证。

在实践中, 分块大小(**chunk size**)和块间重叠(**chunk overlap**)是两个需要仔细调整的关键参数¹⁴。较小的块有助于实现更细粒度的检索和更精确的关键词匹配, 但单个块可能因为缺乏足够的上下文而难以被 LLM 有效利用, 或者导致答案过于碎片化。较大的块能保留更丰富的上下文信息, 有助于 LLM 生成更连贯、更全面的答案, 但同时也可能稀释掉特定查询的精确匹配信号, 增加检索噪声, 并消耗更多的 LLM 处理资源(token 用量)¹⁵。块间重叠则可以在一定程度上缓解因分块导致的信息割裂问题, 但过大的重叠会增加索引的冗余度和存储空间。

此外, 对于包含特殊结构化内容的文档, 如代码片段、表格、Markdown 文本等, 需要采用**自适应的文档分块(adaptive document chunking)**技术, 以确保这些结构在分块过程中不被破坏, 从而保持其原有的意义和可读性, 提高后续的检索效率¹⁵。

一个重要的认知是, 并不存在一种普遍适用于所有场景的“最优”分块策略¹³。开发者需要根据其知识库的具体特性(如文档的平均长度、结构化程度、内容类型、主题多样性等)以及用户查询的典型特点(如查询长度、复杂度、对答案细节的要求等), 通过实验和评估来选择或组合最合适的分块方法。例如, 对于充斥着大量技术术语和代码片段的软件开发文档库, 可能需要结合基于结构的分块(如按函数或类定义分块)和语义分块。

值得注意的是, 分块技术本身也在经历从“规则驱动”向“智能驱动”的演进。早期的分块方法更多地依赖于固定的长度阈值或预定义的结构标记。而语义分块, 尤其是前瞻性的 Agentic 分块的出现¹³, 清晰地表明分块过程本身也开始融入更高级的人工智能能力, 系统试图更智能地理解文档内容和结构, 并据此进行更合理的切分。这对于知识管理系统未来处理日益增长的高度异构和极端复杂的知识文档, 无疑具有深远的战略意义。

最终, 分块策略的选择是在检索准确性、系统效率(索引大小、检索速度、LLM 处理成本)和实现复杂度之间进行的权衡。知识管理系统的开发者需要在这些相互制约的因素之间找到一个适合自身业务需求的最佳平衡点。

2.1.2 嵌入模型选择与微调 (Embedding Model Selection and Fine-tuning)

嵌入模型(Embedding Models)在 RAG 系统中扮演着核心的桥梁角色。它们负责将离

散的文本信息(无论是知识库中的文本块还是用户的查询)转换为稠密的数值向量(vector embeddings)。这些向量被设计用来捕捉文本的深层语义含义,使得在向量空间中,语义上相似的文本片段其对应的向量在空间位置上更为接近,而语义不相关的文本片段则相距较远¹⁷。正是基于这种特性,RAG系统才能够通过计算查询向量与文档块向量之间的相似度来实现高效的语义检索。

然而,并非所有嵌入模型都能胜任所有任务。通用嵌入模型,即那些在海量通用语料库(如互联网文本、维基百科)上预训练的模型,虽然具备广泛的语言理解能力,但在处理特定领域或高度专业的知识时,往往会显得力不从心¹⁷。这些模型可能无法准确理解特定行业术语的细微差别、组织内部的专有缩写或特定上下文中的隐含意义。例如,在金融领域的文档中,“alpha”可能指代超额收益,而在物理学中则可能指代某种粒子,通用模型可能难以根据上下文准确区分并给出合适的向量表示。

因此,为了提升RAG系统在特定知识领域(如法律、医疗、金融、工程等)的表现,**对嵌入模型进行微调(fine-tuning)**变得至关重要。微调是指在一个预训练好的通用嵌入模型的基础上,使用特定领域的语料库进行进一步的训练。通过这个过程,模型能够学习到该领域的特定语言模式、专业词汇、概念关系和上下文依赖,从而使其生成的向量表示更能反映领域内的语义相似性¹⁷。微调后的嵌入模型能够更准确地将用户的领域相关查询匹配到知识库中最相关的文档块,显著提高检索的精准度和召回率,进而提升整个RAG系统的答案质量。

那么,何时应该考虑对嵌入模型进行微调呢?以下是一些关键的信号¹⁸:

1. 处理专业化内容:当知识库主要由特定领域的文档构成,如法律文书、医学研究报告、金融分析、工程规范等,这些文档中充满了通用模型可能不熟悉的专业术语和概念。
2. 系统召回率低下:如果发现RAG系统经常遗漏知识库中实际存在的相关文档,即使用户查询与文档内容在语义上是匹配的,这可能表明当前的嵌入模型未能准确捕捉到这种领域内的语义相似性。
3. 生成不相关或含噪声的响应:如果LLM生成的答案经常偏离主题,或者包含了许多从不相关文档中提取的噪声信息,这可能间接反映了上游检索环节返回的上下文质量不高,而嵌入模型的匹配能力是其中的一个关键因素。

在这些情况下,通过微调嵌入模型来弥合通用模型与组织特定知识之间的语义鸿沟,往往能带来显著的性能提升。正如AWS的相关研究指出,针对特定任务的嵌入模型微调能够极大地提高专业化检索的准确率¹⁸。

选择一个合适的基础嵌入模型进行微调是成功的第一步。目前业界有一些表现优异且广受欢迎的基础模型可供选择¹⁸:

- **OpenAI text-embedding-3-large**: 以其卓越的性能和较长的上下文窗口(支持高达8,192个token)著称,适合处理复杂和较长的文本。

- **NVIDIA NV-Embed-v2**: 这是一个开源模型, 特别擅长处理长文档, 其上下文窗口可达32,768个token, 为需要深度理解长篇技术文档或报告的场景提供了有力支持。
- **Cohere embed-english-v3.0**: 该模型在成本效益和多语言支持方面取得了较好的平衡, 适合需要处理多种语言或对预算有一定要求的应用。

在选择基础模型时, 除了考虑模型本身的性能指标(如在标准benchmark上的得分), 还应关注其上下文窗口大小(对于处理长文档非常重要)、多语言支持能力(如果知识库包含多种语言)、以及许可协议(是否允许商业用途及修改)等因素。通常, 更大的上下文窗口更有利于处理企业内部常见的长篇文档内容¹⁸。

高质量的微调数据集是确保微调效果的基石。构建这样的数据集需要精心策划¹⁷:

- **收集“查询-文档”对 (Query-Document Pairs)**: 从组织内部的知识来源(如内部Wiki、知识库文章、技术支持工单系统、FAQ等)中收集真实的或模拟的用户查询及其对应的最相关文档(或文档块)。这些数据对的质量直接决定了模型能学到什么样的领域特定关联。
- **应用难负例采样 (Hard Negative Sampling)**: 在训练过程中, 除了正例(相关的查询-文档对), 还需要提供负例。简单随机采样得到的负例可能与查询的语义差异过大, 模型很容易区分。难负例则是那些与查询在表面上(如关键词重合)或语义上(如主题相近)有一定相似性, 但实际上并不相关的文档。通过让模型学习区分查询与这些难负例, 可以显著提升其判别能力和检索精度。
- **保持数据集平衡与多样性**: 微调数据集应覆盖领域内的各种主题、概念和查询类型, 并保持一定的复杂性梯度。避免数据集过于偏向某些热门主题或简单查询, 以确保微调后的模型具有良好的泛化能力。

微调过程通常包括以下步骤¹⁸:

1. **选择基础模型**: 根据前述考量选择一个合适的预训练嵌入模型。
2. **收集与预处理数据**: 按照上述原则构建“查询-文档”对数据集, 并进行必要的文本清洗(如去除HTML标签、特殊字符)、标准化(如大小写转换、词形还原)等预处理操作。
3. **模型训练**: 使用选定的深度学习框架(如Hugging Face Transformers, PyTorch, TensorFlow)和硬件资源(通常需要GPU)进行模型训练。在训练过程中, 模型会根据特定的损失函数(见下文)来调整其参数, 以最小化查询向量与相关文档向量之间的距离, 同时最大化与不相关文档向量之间的距离。Amazon SageMaker等云平台也提供了便捷的模型训练和管理服务。
4. **模型评估**: 在独立的验证集上评估微调后模型的性能。常用的检索评估指标包括 Recall@K(如Recall@10, 表示在前K个检索结果中找到相关文档的比例)、Mean Reciprocal Rank (MRR)、Normalized Discounted Cumulative Gain (NDCG)等。
5. **部署与监控**: 将性能达标的微调模型部署到RAG系统中, 并持续监控其在实际应用中的表现, 根据需要进行迭代优化。

在模型训练阶段, 选择合适的损失函数 (**Loss Functions**) 对于引导模型学习到期望的向量表示至关重要。常见的损失函数包括¹⁷:

- **Triplet Loss**: 适用于包含(锚点anchor, 正例positive, 负例negative)三元组的数据。目标是使锚点与正例的距离小于锚点与负例的距离, 并保持一定的边际(margin)。
- **Contrastive Loss**: 适用于正负样本对。目标是拉近正样本对的距离, 推远负样本对的距离。
- **Cosine Similarity Loss**: 适用于包含句子对和它们之间人工标注的相似度得分的数据。目标是使模型预测的余弦相似度与标注得分一致。
- **MultipleNegativesRankingLoss**: 特别适用于只有正样本对(如同义句对、查询-答案对)的数据集。它在每个训练批次中, 将批内的其他样本视为当前正样本对的负例, 从而构造出隐式的负采样。

一些最佳实践有助于提高微调效果并避免常见陷阱¹⁸:

- 避免过拟合 (**Overfitting**): 在独立的验证集上监控模型性能, 并在验证集性能不再提升时停止训练(early stopping), 或者使用正则化技术。过拟合会导致模型在训练数据上表现很好, 但在未见过的数据上表现糟糕。
- 使用多样化的语料库 (**Diverse Corpora**): 确保微调数据集覆盖领域内的广泛主题和语言风格, 以防止模型产生“隧道视野”, 即只在非常狭窄的子领域表现良好。
- 在生产环境中评估 (**Evaluate in Production**): 除了离线评估指标, 更重要的是观察模型在真实用户查询和实际知识库上的表现。有时离线指标的提升并不完全等同于线上效果的改善。

一个常见的错误是使用过时或过于狭窄的数据进行微调, 这会限制模型的泛化能力。因此, 持续地用最新的、有代表性的数据对模型进行验证和迭代至关重要。

总结而言, 嵌入模型的“领域适应性”是构建高性能RAG知识管理系统的关键驱动力之一。通过精心选择基础模型、细致构建高质量的领域特定微调数据集, 并采用恰当的训练策略和工具, 可以显著提升嵌入模型对专业知识的理解和表征能力。这不仅能提高检索的准确性和相关性, 还能减少下游LLM生成答案时的不确定性和幻觉, 最终为用户提供更值得信赖的知识服务。同时, 也应认识到嵌入模型的微调并非一劳永逸的解决方案, 而是一个需要根据领域知识的演变和用户需求的变化而持续迭代和优化的动态过程。知识管理系统需要具备相应的机制来支持这种持续学习和进化的能力。

2.1.3 查询理解与优化 (Query Understanding and Optimization)

用户向知识管理系统提出的查询往往具有多样性和复杂性, 可能包含模糊的表述、隐含的意图、多个知识点, 或者与知识库中信息的措辞方式不完全一致。如果直接使用这些原始查询进行检索, 很可能无法找到最相关或最全面的信息, 从而影响RAG系统的整体表现。因此, **查询理解与优化 (Query Understanding and Optimization)** 成为提升RAG系统

性能,特别是检索模块效能的关键环节。其核心目标是弥合用户自然语言表达与知识库内部信息组织和表述之间的潜在鸿沟。

目前,业界已经发展出多种行之有效的查询优化技术:

- **查询重写/转换 (Query Rewriting/Transformation)**:这类技术旨在将用户的原始查询改写成一个或多个更适合检索系统处理的形式。
 - **假设性文档嵌入 (Hypothetical Document Embeddings, HyDE)**:这是一种创新的查询转换方法。它首先利用LLM根据用户查询生成一个“假设性的”文档片段,这个片段被认为是可能包含答案的理想文本。然后,系统对这个假设性文档进行嵌入,并使用其向量表示去知识库中检索与之语义最相似的真实文档块¹²。其背后的逻辑是,一个包含答案的理想文档的嵌入,可能比原始的、有时较为简短或模糊的查询的嵌入,更能准确地定位到相关的真实知识。
 - **回退提示 (Step-back Prompting)**:这种技术通过引导LLM先从用户提出的具体问题“回退一步”,生成一个关于该问题更宏观、更泛化的问题或概念,然后再针对这个泛化问题和原始具体问题进行搜索或思考。这有助于系统首先把握问题的核心本质和背景,再深入到细节,从而可能找到更全面的上下文信息¹²。
- **查询扩展 (Query Expansion)**:当用户查询中的关键词过于狭窄或使用了不常见的同义词时,查询扩展技术可以通过增加相关的查询词来扩大搜索范围,以期召回更多潜在相关的文档。
 - **基于同义词/相关词的扩展**:利用词典、知识图谱或LLM生成与查询中关键词同义或语义相关的词语,将它们加入到原始查询中。
 - **多查询生成 (Multi-query Generation)**:让LLM根据原始查询生成多个不同角度、不同表述的查询变体。然后对这些变体分别执行检索,最后汇总和重排序所有检索结果¹²。这种方法可以从多个方面捕捉用户的潜在意图。
- **查询分解 (Query Decomposition)**:对于那些包含多个子问题或涉及多个实体及其关系的复杂查询,查询分解技术试图将其拆分成一系列更简单、更原子化的子查询。然后对每个子查询分别进行检索,最后将各子查询的结果综合起来,用于生成最终答案¹²。例如,在LevelRAG框架中,其高层搜索器(high-level searcher)的一个核心职责就是将用户的复杂查询分解为若干个原子查询(atomic queries),以便后续的低层搜索器(如稀疏搜索器、稠密搜索器)能够更精确地处理²¹。
- **查询路由 (Query Routing)**:在拥有多个异构知识源或多种RAG处理流程(pipeline)的复杂知识管理系统中,查询路由技术能够根据用户查询的类型、意图或预测的答案来源,智能地将其导向最合适的知识库或处理路径。例如,Adaptive RAG框架会首先通过一个分类器判断查询的复杂度,然后根据复杂度级别(如无需检索、需要单步检索、需要多步迭代检索)选择不同的RAG策略来响应查询²²。这种机制避免了对所有查询都采用“一刀切”的处理方式,从而提高了系统的整体效率和效果。

查询优化是提升RAG系统鲁棒性和改善用户体验的关键环节。许多RAG系统的失败,并非源于知识库内容的匮乏或是LLM生成能力的不足,而往往是因为未能准确、全面地理解用

用户的真实查询意图。通过应用上述查询重写、扩展、分解和路由等技术¹²，RAG系统能够显著增强其对模糊查询、复杂查询或隐晦查询的理解与处理能力。这使得系统能够更精准地定位到知识库中最相关的上下文信息，为后续的增强生成环节打下坚实的基础，最终提升输出答案的质量和用户满意度。

一个值得注意的趋势是，大型语言模型(LLM)自身越来越多地被直接应用于查询优化的各个环节中。例如，HyDE技术¹⁹完全依赖LLM来生成假设性文档；查询重写、查询扩展中的同义词生成、以及复杂查询的分解，都可以通过精心设计的提示(prompts)和强大的LLM(如GPT-4)来高效完成²¹。这表明LLM在RAG系统中的角色已经从最初单纯的“文本生成器”，扩展到了积极参与“检索过程优化”的智能体。

查询优化的复杂性和必要性，通常与知识库的规模、异构性以及用户查询的多样性成正比。对于一个规模较小、内容相对同质化的知识库，可能简单的原始查询直接检索就能取得不错的效果。然而，对于一个需要管理海量文档、包含多种数据类型(文本、表格、图像等)、覆盖多个专业主题的大型企业级知识管理系统而言，引入高级的查询优化技术(如基于主题模型的查询路由²⁴、多跳逻辑规划的查询分解²¹)则变得至关重要。这些技术能够确保在复杂的信息海洋中，系统依然能够高效、精准地为用户导航到他们真正需要的知识。

2.1.4 检索结果重排序与融合 (Retrieval Result Re-ranking and Fusion, including Hybrid Search)

在RAG系统中，初步的检索阶段(通常基于向量相似性或关键词匹配)可能会返回大量的候选文档或文本块。然而，这些初步结果的排序可能并非最优，其中可能混杂着部分不完全相关或重要性较低的内容。如果直接将这些未加筛选和优化的结果全部或按原始顺序提供给LLM，可能会因为信息噪声、上下文窗口限制或“迷失在中间(lost in the middle)”效应(即LLM更关注上下文开头和结尾部分的信息)而影响最终生成答案的质量。因此，**检索结果重排序(Re-ranking)与融合(Fusion)**技术成为确保LLM获得最优质、最相关上下文的关键后续步骤。

重排序 (Re-ranking):

重排序的目标是在初步检索(也称为一级检索或召回阶段)之后，对召回的候选文档或文本块进行更精细、更准确的相关性评估，并据此重新调整它们的顺序，将最有可能包含答案或对生成答案最有帮助的内容置于更靠前的位置¹²。重排序通常采用比初步检索更复杂、计算成本也更高的模型，因为此时处理的文档数量已经大大减少。常用的重排序方法包括：

- **基于交叉编码器的重排序 (Cross-encoder based Re-ranking)**: 交叉编码器模型(如某些BERT变体)会同时接收用户查询和每个候选文档作为输入，然后输出一个单一的相关性得分。与双编码器(用于初步检索，分别编码查询和文档再计算相似度)相比，交叉编码器能够更充分地建模查询与文档之间的细粒度交互信息，因此通常能提供更准确的相关性判断。

- **基于LLM的重排序 (LLM-based Re-ranking)**: 利用大型语言模型自身的理解和推理能力来对检索到的候选文档进行打分或排序。例如, 可以设计特定的提示, 让LLM评估每个文档与查询的相关性, 或者直接要求LLM从给定的文档列表中选出最相关的几个。

结果融合 (Fusion) 与混合搜索 (Hybrid Search):

当RAG系统采用多种检索策略或从多个知识源获取信息时, 就需要有效的融合机制来整合这些不同来源的结果。

- **混合搜索 (Hybrid Search)**: 这是一种非常重要的融合策略, 它结合了至少两种不同类型的检索方法, 最常见的是稀疏检索(**Sparse Retrieval**) (如基于关键词的BM25算法)和稠密检索(**Dense Retrieval**) (如基于深度学习嵌入模型的向量语义搜索)。稀疏检索擅长精确匹配关键词和专有名词, 而稠密检索则能更好地理解查询的整体语义和上下文, 捕捉同义词、近义词等关系。通过结合两者的优势, 混合搜索通常能够比任何单一方法提供更鲁棒、更全面的检索结果, 有效应对更多样化的用户查询³。LevelRAG框架就是一个很好的例子, 它在其低层搜索器中就集成了稀疏搜索器(利用Lucene语法)、稠密搜索器以及Web搜索器²¹。
- **RAG-Fusion**: 这是一种更具体的融合策略, 它首先通过生成用户原始查询的多个变体(例如, 使用LLM进行查询改写或分解), 然后对每个查询变体分别执行检索。最后, 将所有检索到的结果集合并, 并通过一种智能的重排序算法(如Reciprocal Rank Fusion, RRF)来融合这些结果, 得到最终的排序列表¹²。这种方法通过多角度的查询来提高召回相关信息的概率。

多模态RAG中的相关性评分 (Relevancy Score - RS):

对于处理包含图像、文本等多种模态数据的多模态RAG(MRAG)系统, 传统的基于CLIP的相似度计算在判断跨模态相关性时可能存在不足。为此, 研究者提出了更专门的**相关性评分(Relevancy Score, RS)**机制²⁶。RS模型通常基于视觉语言模型(VLM)进行微调, 旨在更精确地评估一个检索条目(无论是文本还是图像)与用户(可能是多模态的)查询之间的真实相关性, 特别是在区分不相关或仅表面相似的条目方面, RS表现优于传统的CLIP相似度。通过RS进行重排序, 可以显著提升MRAG系统中检索上下文的质量。

重排序和融合技术是RAG系统中计算成本和检索质量之间的一个重要平衡点。初步检索阶段追求的是速度和尽可能广的召回, 可能会牺牲一定的精度(precision)。而重排序阶段则可以在一个相对较小的候选集上, 使用计算更密集但判断更精准的模型(如交叉编码器或LLM)²⁶, 从而在可接受的额外计算开销下, 显著提升最终输入给LLM用于生成的上下文信息的质量。

混合搜索的兴起, 清晰地揭示了单一检索范式(无论是纯粹的关键词匹配还是纯粹的语义相似度搜索)在应对复杂多样的用户需求时的局限性。关键词搜索对于包含特定术语、实

体名称或代码片段的查询非常有效,但难以处理同义替换或概念性的查询。相反,语义搜索能够理解更广泛的查询意图和概念关联,但有时可能在精确匹配上表现不足。混合搜索³通过取长补短,结合了两者的优点,能够更稳健地处理来自用户的各种类型的输入,这对于需要具备普适性的企业级知识管理系统来说,是一个非常重要的能力。

展望未来,随着大型语言模型能力的持续增强,它们可能会在检索结果的融合与推理决策阶段扮演越来越核心的角色。目前像RAG-Fusion¹²和基于LLM的重排序等方法,已经初步展示了LLM在理解和整合多源异构信息方面的巨大潜力。未来,LLM可能不仅仅是简单地检索结果进行加权、排序或选择,而是能够进行更深层次的逻辑推理、信息校验、矛盾消解和知识综合,从而真正实现从“信息检索”到“知识构建”的飞跃,为用户提供高度提炼和富有洞察的答案。

2.2 增强生成 (Augmented Generation)

在RAG系统中,检索到的相关上下文信息为大型语言模型(LLM)提供了生成高质量答案的基础。然而,仅仅将检索到的文本块与原始查询简单拼接,并不足以保证LLM能够充分理解和利用这些信息。**增强生成 (Augmented Generation)**阶段的核心任务,就是通过精心的设计和技术手段,确保LLM能够以最有效的方式整合检索到的知识,并生成满足用户需求的、准确且忠实的回答。这一阶段的关键技术包括提示工程和减少幻觉、提升忠实度的策略。

2.2.1 RAG中的提示工程 (Prompt Engineering in RAG)

提示工程 (Prompt Engineering)在RAG系统中扮演着连接“检索”与“生成”两个核心环节的关键桥梁角色。其根本目标是设计出能够清晰、准确地引导LLM如何利用检索到的上下文信息来回应用户原始查询的指令(即提示)³⁰。一个精心设计的提示,能够显著提升LLM生成答案的准确性、相关性、连贯性和风格的恰当性。

RAG场景下的提示工程,相比于通用的LLM提示工程,具有其独特性和更高的复杂性。它不仅要向LLM传达用户的查询意图(例如,是要求回答问题、进行总结、还是比较分析),更重要的是,它必须有效地将检索到的外部知识(上下文)融入到LLM的“思考”过程中。这通常涉及到以下几个关键要素³⁰:

1. 明确数据需求与上下文的引入 (Specifying Data Needs and Context

Introduction):提示需要清晰地告知LLM,它应该依赖哪些信息来生成答案。这通常通过在提示中明确划定一个区域来放置检索到的上下文文本块来实现。例如,提示可以包含类似“请根据以下提供的上下文信息来回答问题:[此处插入检索到的文本块]”的指令。

2. 引导LLM的解释与信息聚焦 (Guiding LLM Interpretation and Information

Focus):仅提供上下文是不够的,还需要引导LLM如何去理解和使用这些上下文。提示可以包含一些约束条件,例如指示LLM“只依据提供的上下文回答,如果上下文中没有相关信息,请明确指出”、“请重点关注上下文中的XX方面”、“请忽略上下文中的

不相关信息”等。这有助于LLM在生成答案时，聚焦于最相关的细节，避免受到噪声信息的干扰或产生不必要的联想³⁰。

3. **定制语言风格、语气与输出格式 (Tailoring Language Style, Tone, and Output Format)**: 根据知识管理系统的应用场景和用户群体的不同，可能需要LLM以特定的风格(如正式、非正式、专业、通俗)、语气(如客观、友好、权威)或格式(如点列式、段落式、表格)来生成答案。这些要求都可以在提示中进行明确规定³⁰。

一个普遍的共识是，即使RAG系统检索到了高质量的上下文信息，一个设计拙劣或含糊不清的提示也可能导致LLM生成质量低下、不相关甚至错误的回答³¹。因此，RAG系统中的提示工程需要精确地指导LLM如何定位和利用检索到的数据，并确保LLM理解其在当前任务中的具体角色和行为准则。

RAG中的提示工程比通用LLM提示工程更为复杂，因为它需要在“指令的清晰性”和“上下文利用的充分性”之间取得精妙的平衡。一方面，指令需要足够清晰，让LLM明白任务目标和行为边界；另一方面，指令又不应过于死板，以至于限制了LLM对上下文信息进行合理推理和综合的能力。设计这样的提示，需要开发者对LLM的行为模式有深入的理解，并进行大量的实验和迭代优化。

此外，RAG系统中的提示工程并非一个孤立的环节，它与RAG的其他组件(如分块策略、重排序算法等)是相互影响、需要协同优化的。例如，如果分块策略倾向于生成较短的文本块，那么提示可能需要引导LLM更主动地去综合来自多个不同文本块的信息，以形成一个完整的答案。如果重排序算法将最相关的文本块置于上下文的开头，提示也可以利用这一点来引导LLM优先关注这些信息。

随着RAG系统向更高级的模块化(Modular RAG)和智能体化(Agentic RAG)方向发展¹¹，提示工程的角色也在发生深刻的演变。在这些更复杂的系统中，LLM可能需要根据提示来执行一系列的多步骤操作，或者调用外部的工具和服务。此时的提示工程，就不再仅仅是简单的问答指令设计，而更像是为LLM智能体编写一套复杂的“任务执行脚本”或“行为规则手册”，其中包含了条件判断、循环、工具调用描述等更高级的指令元素。这对提示工程技术提出了更高的要求，也为其开辟了更广阔的应用空间。

2.2.2 减少幻觉与提升忠实度的技术 (Techniques for Reducing Hallucination and Improving Faithfulness)

****幻觉(Hallucination)****是大型语言模型(LLM)在实际应用中面临的一个核心挑战，它指的是LLM生成了听起来合理但实际上不准确、与事实不符、或者在提供的上下文中找不到依据的内容⁷。对于知识管理系统而言，信息的准确性和可靠性是其存在的基石，因此，有效减少幻觉至关重要。

忠实度(Faithfulness)，在RAG的语境下，特指LLM生成的答案是否真实、准确地反映了所提供的检索上下文中的信息，没有进行歪曲、夸大、遗漏关键细节或添加上下文中未包含

的外部信息³⁵。提升忠实度是确保RAG系统输出值得信赖的关键。

RAG系统通过引入外部知识源,其核心目的之一就是将LLM的回答“锚定”在可验证的事实基础上,从而从根本上减少LLM仅凭其内部参数化知识进行“自由发挥”而导致幻觉的可能性³。然而,仅仅引入外部知识并不足以完全消除幻觉或保证绝对的忠实度。检索到的上下文可能本身就包含噪声、过时信息甚至错误;LLM在理解和综合上下文时也可能出现偏差。因此,研究者们开发了一系列专门的技术来进一步减少幻觉并提升RAG系统的忠实度:

1. **Corrective Retrieval Augmented Generation (CRAG)**: CRAG框架通过引入一个轻量级的**检索评估器(retrieval evaluator)**来主动评估初始检索到的文档集的质量。根据评估结果(置信度), CRAG会触发不同的知识检索和修正动作³²:
 - 如果检索质量高(**Correct**), 则可能对检索到的文档进行知识提炼(如“分解-再组合”算法), 以提取关键信息并滤除无关部分。
 - 如果检索质量低(**Incorrect**), 则会丢弃这些不可靠的文档, 并通过查询重写(**query rewriting**)后, 利用大规模Web搜索作为补充的知识来源。
 - 如果检索质量介于两者之间(**Ambiguous**), 则可能会结合使用提炼后的初始检索知识和Web搜索获取的知识。通过这种动态的、基于评估的纠错机制, CRAG旨在确保LLM在生成答案时所依赖的上下文尽可能准确和相关, 从而提高最终答案的鲁棒性和事实性。
2. **Self-RAG (Self-Reflective Retrieval-Augmented Generation)**: Self-RAG赋予LLM一种“自我反思”的能力。它通过训练LLM生成一系列特殊的“反思token(**reflection tokens**)”来实现对检索和生成过程的动态控制和批判性评估³⁸。这些反思token可以指示:
 - 是否需要进行检索(**Retrieve token**)。
 - 检索到的段落是否与当前查询相关(**Critique token: IsRelevant**)。
 - LLM生成的某个片段是否得到了所引用段落的支持(**Critique token: IsSupported**)。
 - 生成的回复对于用户查询的整体效用如何(**Critique token: IsUseful**)。在生成过程中, LLM会基于这些自我产生的反思token来调整其行为, 例如, 如果判断检索到的段落不相关, 它可能会选择不使用该段落, 或者如果判断生成的某个说法缺乏证据支持, 它可能会尝试修改或重新生成。这种内在的批判和调整机制有助于显著提升生成内容的忠实度和事实准确性, 并减少幻觉。
3. **知识蒸馏 (Knowledge Distillation - KD)** 用于减少幻觉: 知识蒸馏是一种模型压缩和性能提升技术。在RAG的背景下, 它可以被用来训练一个更“谨慎”、更不容易产生幻觉的学生LLM。具体做法是, 使用一个更强大、更可靠的教师模型(teacher model)为训练数据生成“平滑的软标签(smoothed soft labels)”, 这些软标签代表了教师模型对下一个token的概率分布, 而不仅仅是一个确定的token(硬标签)。学生模型通过学习拟合这些软标签, 可以减少其预测的过分自信(overconfidence), 并更好地理解答案的不确定性和多样性, 从而在生成时更倾向于依赖事实依据, 减少凭空捏造³⁴。
4. **Post-Hoc RAG** (主要用于多模态场景): 在多模态RAG(如处理视频问答)中, 一种减

少幻觉的策略是在LLM初步生成内容(如视频字幕)之后,动态地构建一个与当前输入视频内容高度相关的、临时的知识库(ad-hoc knowledge base)。然后,利用一个RAG模块(可以看作是一个事实核查模块)来检索这个临时知识库,以审查和验证LLM生成的输出是否与视频中的视觉和元数据信息保持事实一致性。如果不一致,则对输出进行修正⁷。

评估RAG系统幻觉和忠实度的常用评估指标包括上下文相关性(Context Relevance)、答案忠实度(Answer Faithfulness)、答案相关性(Answer Relevance)等³⁵。一些专门的基准测试如RAGTruth⁴²也被开发出来用于评估RAG系统在幻觉检测方面的能力。

从这些技术的发展可以看出,减少幻觉和提升忠实度已经成为RAG研究的核心议题之一。这不再仅仅依赖于提供外部知识,而是演变成一个涉及RAG系统多个环节的系统工程。从对检索质量的智能评估与主动纠正(如CRAG),到生成过程中引入模型的自我反思与批判机制(如Self-RAG),再到从模型训练层面直接优化其生成行为(如知识蒸馏),这些技术共同致力于构建更可靠、更值得信赖的RAG系统。

一个重要的趋势是,“可验证性(verifiability)”正成为衡量高级RAG系统质量的一个关键特性。例如,Self-RAG不仅尝试生成准确的答案,还强调生成引文(citations),并自我评估其生成的内容是否确实得到了所引用证据的支持³⁸。CRAG也强调对检索来源的评估和选择³²。这种对答案来源和证据链的关注,对于知识管理系统建立用户信任、支持决策制定以及满足合规性要求具有至关重要的意义。

此外,高级RAG系统的一个重要标志是其识别和处理“未知”或“不确定”情况的能力。当检索到的信息不足以回答问题,或者信息存在矛盾时,系统不应强行生成答案,而应能够识别这种不确定性,并采取适当的补救措施。例如,CRAG在判断初始检索结果质量不高时,会主动触发Web搜索以寻求更广泛的信息³²。Self-RAG则可以学习在某些情况下决定不进行检索,或者判断现有信息不足以支持一个高质量的回答³⁸。这种“自知己之所不知”并据此调整策略的能力,是RAG系统从简单的信息传递工具向真正智能的知识服务伙伴迈进的关键一步。

2.3 知识增强与整合 (Knowledge Augmentation and Integration)

在RAG系统中,检索到的外部知识需要与LLM固有的内部参数化知识进行有效结合,才能最终生成高质量的响应。仅仅将检索到的信息简单地堆砌给LLM是不够的,尤其是在处理复杂查询或需要综合多方面信息的场景下。因此,**知识增强与整合(Knowledge Augmentation and Integration)**成为RAG流程中一个至关重要的环节,其核心目标是确保外部知识能够以最优的方式被LLM理解和利用,从而最大化RAG系统的性能⁴³。

知识整合可以发生在RAG流程的不同层面⁸:

- **输入层整合 (Input Layer Integration)**: 这是最常见和直接的方式。检索到的相关文本块被直接拼接到用户的原始查询之后,或者按照特定的模板格式组织起来,形成一

个增强的提示(augmented prompt),然后将这个完整的提示输入给LLM。提示工程(Prompt Engineering)在此扮演了关键角色,负责设计有效的提示结构,以引导LLM关注和利用这些外部知识。

- **中间层整合 (Intermediate Layer Integration):**一些更复杂的RAG架构尝试在LLM的内部处理层(例如Transformer的某些中间层)引入检索到的知识。这种方式理论上可以让外部知识更深入地影响模型的内部表征学习和信息处理流,而不仅仅是作为初始输入。然而,这种方法通常需要对LLM的结构进行修改或采用特定的模型架构,实现起来更为复杂。
- **输出层整合 (Output Layer Integration):**在这种策略下,可能会先让LLM基于不同来源的知识(或对同一知识的不同解读)生成多个候选的初步响应。然后,再结合检索到的知识(或其他评估标准),对这些候选响应进行重排序、筛选、融合或修正,以得到最终的输出。例如,可以训练一个模型来判断哪个候选响应与检索到的证据最一致。

为了实现更智能和高效的知识整合,研究者们提出了一些创新的方法:

- **Self-Selection RAG:**这个框架赋予LLM一种“自我选择”的能力,即在面对一个查询时,LLM可以自行判断是仅依赖其内部参数化知识生成答案更好,还是结合外部检索到的知识生成答案更优,或者如何组合这两种知识来源⁴³。为了实现这一点,该框架通常需要在—个特殊的“检索-生成偏好(Retrieval-Generation Preference, RGP)”数据集上进行训练。这个数据集包含了(查询,仅内部知识生成的答案,结合外部知识生成的答案,偏好标签)这样的样本。通过使用诸如直接偏好优化(Direct Preference Optimization, DPO)等对齐技术进行训练,LLM能够学习到在不同情况下如何更有效地整合或选择知识来源,以期达到更高的准确性和可靠性。
- **利用知识图谱 (Knowledge Graphs, KGs) 进行增强整合:**知识图谱以其结构化的方式存储了实体及其之间的丰富关系,这为RAG系统提供了超越非结构化文本的知识维度。在知识整合阶段,KG可以发挥重要作用⁴⁷:
 - 提供结构化上下文:KG可以为检索到的文本块中的实体提供更丰富的背景信息和关系链接,帮助LLM更深入地理解上下文。
 - 指导信息综合:在需要综合多个检索片段以回答复杂问题时,KG中的路径和关系可以为LLM提供“推理路径”的指导,帮助其更有条理地组织和生成答案。
 - 事实校验与补充:KG可以作为一种事实来源,用于校验LLM基于文本生成的初步结论,或者补充文本中可能缺失的关键关系信息。**KG²RAG**框架就是一个很好的例子,它不仅在检索阶段利用KG进行块扩展,还在后续的上下文组织阶段利用KG来过滤和排列信息,确保LLM接收到的是结构良好、关系明确的知识输入⁴⁸。

从这些发展可以看出,RAG中的知识整合策略正在从早期相对简单的“信息拼接”向更高级的“智能决策与融合”演进。早期的RAG系统主要是在输入提示的层面将检索到的文本与原始查询进行简单的串联。而像Self-Selection RAG这样的新方法⁴³,则通过让LLM主动学习何时以及如何最有效地整合外部知识,甚至在来自不同来源(内部参数化知识 vs. 外部检

索知识)的潜在答案之间进行明智的选择,这无疑代表了知识整合过程的智能化趋势。

与此同时,结构化知识,特别是知识图谱,在RAG系统中的作用也日益受到重视⁴⁷。非结构化文本虽然信息量大,但往往缺乏对实体间复杂关系的显式表达。知识图谱正好弥补了这一不足,它能够清晰地揭示实体之间的层级、属性、依赖等多种关系。将KG融入RAG的知识整合环节,不仅可以帮助LLM更深刻地理解单个信息点的上下文,还能支持其进行更复杂的多跳推理(multi-hop reasoning),从而生成更具深度、更富有逻辑性的答案。这对于构建需要处理专业领域知识、回答复杂关联问题的知识管理系统来说,具有非常重要的应用价值。

最后,知识整合的有效性也与更广泛的LLM对齐(alignment)技术和偏好学习(preference learning)研究密切相关。例如,DPO等偏好学习方法被成功应用于训练LLM更好地整合和选择知识,如在Self-Selection RAG中所展示的⁴³。这意味着RAG系统的发展,正越来越多地受益于那些旨在使LLM的行为更符合人类期望、更能满足特定任务需求的前沿AI研究成果。未来的知识整合模块可能会更加依赖这些先进的对齐技术,以确保外部知识能够以最符合任务目标和用户偏好的方式被LLM所用。

第三部分:前沿RAG技术研究进展 (Advances in Cutting-Edge RAG Technologies)

随着RAG技术的不断成熟,研究者们不再满足于基础的“检索-生成”模式,而是积极探索更智能、更鲁棒、更高效的RAG新范式。本部分将聚焦于几类代表性的前沿RAG技术,包括具备自我进化与纠错能力的RAG、动态适应性RAG、知识图谱增强RAG、多模态RAG以及针对大规模知识库的优化策略。

3.1 自我进化与纠错的RAG (Self-Evolving and Corrective RAG)

传统的RAG系统在面对检索结果质量参差不齐或与用户查询不完全匹配时,其性能可能会受到较大影响。为了提升RAG系统在复杂和不确定环境下的鲁棒性和准确性,研究者们开始赋予RAG系统一定的“自我意识”和“纠错能力”,使其能够主动评估信息质量、调整检索策略并修正潜在错误。

- Self-RAG (Self-Reflective Retrieval-Augmented Generation) 19:

Self-RAG的核心思想是训练一个单一的LLM,使其能够通过生成一系列特殊的“反思token(reflection tokens)”来对自身的检索和生成过程进行实时的、细粒度的控制和批判。这些反思token嵌入在模型的生成序列中,指导着模型的行为。

其关键组件虽然在推理时统一由核心LLM(Generator)实现,但在训练阶段涉及到:

1. 一个外部的**Retriever**,负责根据指令获取文档。
2. 一个**Critic**模型(通常是另一个LLM),在训练数据准备阶段用于自动标注反思token,教会Generator何时以及如何进行反思。
3. 核心的**Generator** LLM,它学习预测正常的文本序列以及这些特殊的反思token。

Self-RAG的工作流程大致如下：

4. **检索决策 (Retrieve)**: 在生成每个片段之前, Generator会首先预测一个“检索”类型的反思token(例如, „ ``), 以判断当前是否有必要从外部知识库获取信息, 或者是否可以继续使用之前检索到的信息。如果决定检索, 则调用外部Retriever。
 5. **生成 (Generate)**: 如果不需要检索, Generator就像标准LLM一样直接预测下一个输出段。如果进行了检索, Generator会首先针对每个检索到的文档块生成一个“相关性批判”的反思token(例如, `` [Irrelevant]), 评估该文档块对于当前生成任务的价值。然后, 它会基于被认为是相关的文档块来生成后续的文本内容。
 6. **批判 (Critique)**: 在生成了一个文本片段后(尤其是在利用了检索信息的情况下), Generator会进一步生成“支持度批判”的反思token(例如, „ ``), 以评估其刚刚生成的这个片段是否得到了所引用证据的充分支持。最后, 还可能生成一个“整体效用”的反思token(例如, [Useful], [NotUseful]), 来评价整个回复对于用户原始查询的帮助程度。Self-RAG的主要优势在于其高度的自适应性和可控性。它能够根据具体查询的需要动态调整检索行为(可以进行多次检索, 也可以完全跳过检索), 并通过内部的批判机制来提升生成内容的事实性和与证据的一致性。实验结果表明, Self-RAG在开放域问答、复杂推理和事实验证等多种任务上, 其性能显著优于当时的SOTA LLM(如ChatGPT)以及传统的RAG模型, 尤其在提升事实准确性和引文精度方面表现突出。
- CRAG (Corrective Retrieval Augmented Generation) 19:
CRAG的设计目标是提升RAG系统在初始检索结果质量不佳(例如, 检索到不相关、不准确或过时的文档)时的鲁棒性。它引入了一种“纠错”机制。
其关键组件包括:
 1. 一个轻量级的检索评估器(**retrieval evaluator**)(例如, 一个经过微调的T5-large模型), 用于快速评估一批初始检索到的文档与用户查询的整体相关性质量。
 2. 一个查询重写器(**query rewriter**)(例如, 可以利用ChatGPT等强大LLM实现), 用于在需要进行Web搜索时优化原始查询。
 3. **Web搜索集成模块**, 用于从互联网上获取补充信息。
 4. 一种知识提炼算法(**knowledge refinement algorithm**), 通常采用“分解-再组合(decompose-then-recompose)”的策略, 将文档(无论是初始检索的还是Web搜索的)分解成更小的知识片段(如句子), 过滤掉不相关的片段, 然后将相关的片段重新组合成更精炼的上下文。CRAG的工作流程如下:
 5. **评估初始检索质量**: 当RAG系统进行初步检索后, 检索评估器会对返回的文档集给出一个整体的置信度评分, 并将其归类为三种状态之一: 正确(Correct)、不正确(Incorrect)或模糊(Ambiguous)。
 6. **执行纠错动作**:
 - **正确 (Correct)**: 如果评估器认为初始检索的文档中至少有一个质量较高、与查询高度相关, 系统会选择对这些高质量文档应用知识提炼算法, 提取核心信息, 然后用于生成。
 - **不正确 (Incorrect)**: 如果评估器认为所有初始检索的文档质量都很低或不相

关，系统会果断丢弃这些文档。然后，它会调用查询重写器对原始用户查询进行优化，使其更适合通用搜索引擎。接着，利用优化后的查询进行大规模的Web搜索，并将搜索结果作为新的知识来源，再经过知识提炼后用于生成。

- **模糊 (Ambiguous)**: 如果评估器对初始检索结果的判断介于上述两者之间(例如，有一些相关性但不够好)，系统可能会采取一种混合策略，即结合使用对初始文档进行提炼后得到的知识和通过Web搜索获取的补充知识，共同作为LLM生成的上下文。
- 7. **最终生成**: 利用经过上述纠错和提炼过程得到的优化知识，LLM最终生成对用户查询的响应。CRAG的主要优势在于其“即插即用”的特性，它可以方便地集成到各种现有的RAG框架中(包括标准的RAG甚至Self-RAG)。实验表明，CRAG能够显著提升这些RAG方法在多种生成任务(覆盖短文本问答和长文本生成)上的性能，尤其是在初始检索结果包含错误或噪声时，其鲁棒性优势更为明显。
- **Self-Routing RAG (SR-RAG) 49**:
SR-RAG的核心创新在于将选择性检索(selective retrieval)与知识口述化(knowledge verbalization)紧密结合起来。它赋予LLM一种动态的“自路由”能力，使其能够根据当前情境，自主决定是应该从外部知识源进行检索，还是直接“口述”其自身参数化模型中已有的相关知识来回答问题。
其关键机制包括：
 1. **知识源选择 (Knowledge Source Selection)**: SR-RAG将选择性检索问题重新表述为一个知识源选择问题。LLM在接收到查询后，会首先判断是应该依赖其内部已有的知识(参数化知识)，还是需要求助于外部检索。
 2. **知识口述化 (Knowledge Verbalization)**: 如果LLM决定依赖内部知识，它会尝试将其参数化知识中与查询相关的内容明确地“口述”或表达出来，形成一段内部生成的知识文本。这种能力扩展了LLM在没有外部检索或外部检索结果不佳时独立回答问题的能力。同时，清晰的知识口述也有助于后续更准确地判断是否仍需外部检索。
 3. **多任务联合优化 (Multi-task Objective)**: SR-RAG通过设计一个多任务的学习目标，来联合优化LLM在知识源选择、知识口述化以及最终响应生成这三个方面的能力。
 4. **动态知识源推理 (Dynamic Knowledge Source Inference)**: 为了应对在不同领域或任务(domain shifts)以及LLM自身能力因微调而发生变化时，知识源选择决策准确性可能下降的问题，SR-RAG还引入了一种基于最近邻搜索的动态知识源推理机制，利用与当前情境在隐空间中相似的过往决策案例来辅助判断。SR-RAG的主要优势在于，它通过更智能地利用LLM的内部知识，并将其与外部检索相结合，从而在提升响应准确率和推理效率(减少不必要的检索)方面取得了显著效果。实验表明，与基线的选择性检索方法相比，SR-RAG能够在显著减少检索次数的同时，进一步提升整体性能。

这些自我进化与纠错的RAG技术，如Self-RAG、CRAG和SR-RAG，共同揭示了一个重要的

发展趋势: RAG系统正在从一个相对被动的、按固定流程整合信息的工具, 进化成为一个更主动的、具备一定自我意识和纠错能力的智能体。它们不再盲目地信任和使用所有检索到的信息, 而是学会了评估信息质量、判断信息需求、选择信息来源, 甚至在必要时修正或补充信息。这种进化使得RAG系统在面对真实世界中复杂、多变、甚至充满噪声的信息环境时, 能够表现出更强的鲁棒性、更高的智能水平和更优的性能。

“元认知(meta-cognition)”能力, 即模型对其自身知识状态、信息处理过程以及输出结果质量的认知和反思能力, 正成为提升高级RAG系统性能的关键。Self-RAG通过引入反思 token, 让模型能够“思考”其检索行为的必要性以及生成内容的好坏与证据支持度³⁸。SR-RAG则让模型“思考”当前问题是更适合依赖其内部已存储的知识来回答, 还是必须求助于外部的检索过程⁵⁰。这种元认知能力的引入, 使得RAG系统能够更有效地规划其信息获取和利用策略, 避免不必要的计算资源浪费(如冗余检索), 并显著提高最终答案的准确性和可靠性。

对于构建先进的知识管理系统而言, 这些高级RAG技术带来了巨大的潜力, 同时也提出了更高的智能化要求。一个能够自我纠错、自我反思的知识管理系统, 不仅能够为用户提供更精准、更可信的答案, 更重要的是, 当其内部知识库存在暂时性的缺陷、不完整或者用户提出的查询本身较为模糊、刁钻时, 这样的系统能够展现出更强的适应性和问题解决能力。例如, CRAG在判断内部知识不足以回答问题时, 能够主动扩展其信息来源至Web搜索³²; SR-RAG则可能通过口述内部知识先给出一个初步的、基于已有理解的回答, 并可能指出其不确定性, 这都远胜于简单地返回“未找到答案”或生成一个不可靠的臆测。这标志着知识管理系统正在向更深层次的智能化迈进。

3.2 动态适应性RAG (Dynamic Adaptive RAG)

传统的RAG流程往往采用固定的检索和生成策略, 无论用户查询的简单与复杂, 都执行大致相同的操作。然而, 用户的查询千差万别, 有些可能非常直接, LLM仅凭自身知识就能回答; 有些可能需要简单的外部信息佐证; 而另一些则可能涉及复杂的多步骤推理和多源信息整合。针对所有查询都采用“一刀切”的RAG策略, 不仅可能对简单查询造成不必要的计算资源浪费(过度检索和处理), 也可能对复杂查询因处理深度不足而无法给出满意答案。动态适应性RAG(Adaptive RAG)正是为了解决这一问题而提出的, 其核心思想是让RAG系统能够根据输入查询的特性(主要是复杂度)动态地选择和调整最合适的处理策略。

- Adaptive-RAG 19:

Adaptive-RAG框架的核心在于引入了一个查询复杂度分类器(query complexity classifier)。这个分类器通常是一个相对小型的语言模型(LM), 它被专门训练用来预测输入查询的复杂性级别。例如, 可以将查询分为三个级别:

1. 级别A(无需检索): 非常直接或常识性的问题, LLM可以利用其内部参数化知识直接回答。
2. 级别B(单步检索): 中等复杂度的问题, 需要通过一次性的外部知识检索来获取必要的上下文信息, 然后LLM结合此上下文生成答案。

3. 级别**C**(多步/迭代检索):非常复杂的问题,可能需要进行多次迭代的检索、信息整合和推理步骤,LLM才能给出完整和准确的答案。在工作流程上,当用户查询输入后,首先由这个小型LM分类器对其进行分析并预测其复杂度级别。然后,Adaptive-RAG框架会根据分类结果,动态地选择并执行最适合该复杂度级别的(检索增强) **LLM**策略。例如,如果判断为级别A,则可能直接调用LLM生成答案,完全跳过检索环节;如果为级别B,则执行标准的单步RAG流程;如果为级别C,则可能启动一个更复杂的多步迭代式RAG流程(例如,先分解问题,再对子问题分别检索和推理,最后汇总)。Adaptive-RAG的主要优势在于它能够在系统的整体效率和回答的准确性之间取得更好的平衡。它避免了对简单查询进行不必要的、耗时的复杂处理流程,从而节省了计算资源并加快了响应速度;同时,它又能确保那些真正复杂的查询能够得到足够深入和细致的处理,从而保证答案的质量。实验表明,Adaptive-RAG相比于简单的RAG方法(如总是进行单步检索或从不检索)以及一些其他的自适应检索方法,能够在多种问答数据集上取得整体更优的性能表现,尤其是在混合了不同复杂度查询的场景下。
- AT-RAG (Adaptive Topic-RAG) 19:
AT-RAG是另一种形式的自适应RAG,它特别关注于通过主题过滤(topic filtering)和迭代推理(iterative reasoning)来提升RAG系统处理复杂多跳查询(multi-hop queries)的效率和准确性。多跳查询通常需要从多个文档或知识片段中获取信息,并将这些信息串联起来才能得到最终答案。
其关键机制包括:
 1. 动态主题建模(**Dynamic Topic Modeling**):AT-RAG利用诸如BERTopic这样的主题建模技术,对用户查询进行实时分析,并为其动态地分配一个或多个最相关的主题。
 2. 基于主题的检索优化:在检索阶段,系统会优先在与查询主题相匹配的文档子集或知识区域内进行搜索,或者利用主题信息来指导检索权重的分配。这有助于更精准地定位到与查询核心意图相关的初始上下文,提高检索的准确性和效率。
 3. 迭代推理(**Iterative Reasoning**):对于判断为需要多跳推理的复杂查询(可能基于主题分析或其他复杂度评估),AT-RAG会采用一种多步骤的推理过程。这可能涉及到逐步构建答案、在每一步根据已有的中间结果调整后续的检索方向或推理焦点,直到最终形成完整的答案。AT-RAG的主要优势在于其处理需要深度领域知识和复杂逻辑链条的查询(例如在医疗问答等专业领域)时所展现出的潜力。通过结合主题分析来缩小检索范围和指导迭代推理,它旨在以更高的精度和效率解决传统RAG难以应对的多跳问题。

动态适应性RAG的理念清晰地指出了一个重要方向:即RAG系统不应再是僵化的、一成不变的处理流程,而应具备根据具体输入(特别是用户查询的特性)灵活调整自身行为的能力。“一刀切”的RAG策略在面对真实世界中形形色色的用户查询时,其效率和效果都难以达到最优。无论是Adaptive-RAG通过显式的查询复杂度分类来选择不同的RAG路径²²,还是AT-RAG通过主题分析来指导更精细的检索和迭代推理²⁴,其核心思想都是对输入进行

更深入的理解, 并基于这种理解来匹配最合适的系统资源和处理逻辑。这种“量体裁衣”式的处理方式, 既避免了用“牛刀杀鸡”(即对简单问题进行过度复杂的检索和推理), 也防止了“杯水车薪”(即对复杂问题因处理深度不足而无法给出满意答案)的情况发生。

在这些自适应RAG架构中, 前端的查询分类/分析模块扮演了至关重要的“调度员”角色。Adaptive-RAG依赖一个小型LM分类器来判断查询的复杂度级别²², 而AT-RAG则使用主题模型来洞察查询的核心主题和潜在的关联范围²⁴。这些分析模块的准确性和效率, 直接决定了后续RAG流程能否被正确地选择和引导, 从而对整个系统的最终性能产生深远影响。

更进一步, 动态适应性RAG的理念完全可以被借鉴并扩展到知识管理系统的整体用户交互设计之中。未来的知识管理系统可以不仅仅根据单次查询的文本特性来调整RAG策略, 还可以综合考虑更丰富的上下文信息, 例如: 用户的历史查询行为、用户的专业背景或角色权限、用户当前正在执行的任务流程、甚至是用户明确表达的偏好(如希望得到简略答案还是详细解释)。基于这些多维度信息的综合判断, 知识管理系统可以动态地配置其RAG引擎, 以提供真正个性化、场景化、且最高效的知识服务。例如, 对于系统中的新手用户或进行探索性查询的用户, 系统可以自动采用更广泛的检索范围和提供包含多种视角解释的答案; 而对于领域专家或目标非常明确的查询, 则可以采用更精确、更深入的检索策略, 并直接给出高度凝练的专业答案。这种深层次的自适应能力, 将是下一代智能知识管理系统的重要特征。

3.3 知识图谱增强RAG (Knowledge Graph-Enhanced RAG - KG²RAG)

传统的RAG系统主要依赖于从非结构化的文本文档中检索信息。然而, 现实世界中的知识很多是以结构化的形式存在的, 其中, **知识图谱(Knowledge Graphs, KGs)**作为一种能够清晰表达实体(entities)及其之间复杂关系(relations)的强大工具, 为RAG系统提供了超越纯文本语义相似性检索的潜力。知识图谱增强RAG(KG-RAG)正是旨在利用KG的结构化优势来改进RAG系统的检索和生成能力。

然而, 在实际应用中, 构建和维护一个完整、准确且覆盖广泛的知识图谱本身就是一个巨大的挑战。许多真实世界的KG都存在**不完整性(incompleteness)**的问题, 即可能缺失回答某些问题所需的关键实体或关系三元组⁴⁷。这就要求KG-RAG方法不仅要能利用KG, 还要能一定程度上应对其不完整性。

KG²RAG (Knowledge Graph-Guided Retrieval Augmented Generation) 框架⁴⁸是KG-RAG领域一个值得关注的代表性工作。它提出了一套系统性的方法, 通过知识图谱来指导文本块(chunks)的扩展检索和上下文组织, 以期提高检索结果的多样性、连贯性, 并最终改善LLM生成答案的质量。其核心流程包括:

1. 文档离线处理与KG-块关联(Document Offline Processing and KG-Chunk Association):
 - 与标准RAG类似, 首先将原始文本文档分割成较小的文本块。

- 关键步骤是将这些文本块与一个知识图谱进行关联。这可以通过两种方式实现：
 - 如果已有一个现成的、相关的KG(例如, 领域知识图谱如WebQSP、CWQ等), 则可以使用实体识别与链接(Entity Recognition and Linking, ERL)技术, 将文本块中提及的实体链接到KG中对应的节点, 并将文本块与包含这些实体的KG子图或路径关联起来。
 - 如果缺乏现成的KG, KG²RAG提出了一种更为灵活的方法, 即直接从文本块自身中提取实体和它们之间的关系(例如, 可以利用LLM通过精心设计的提示来完成这项任务, 如论文图3所示), 从而为每个文本块构建一个局部的子图。然后, 这些从不同文本块中提取出的子图可以被合并起来, 形成一个覆盖整个文档集的、自底向上构建的知识图谱。
- 无论采用哪种方式, 最终的目标是建立起KG中的每个三元组(头实体, 关系, 尾实体)与其来源文本块之间的明确链接。这个KG-块关联的过程是独立于用户查询的, 在离线阶段完成。
- 2. **KG增强的块检索(KG-Enhanced Chunk Retrieval)**: 这是KG²RAG的核心检索机制, 它分为两个子阶段:
 - 基于语义的初步检索(**Semantic-based Retrieval**): 首先, 像标准RAG一样, 对用户查询进行嵌入, 并与知识库中所有文本块的嵌入进行比较, 通过计算语义相似度(如余弦相似度)来检索出Top-K个与查询最相似的文本块。这些块被称为“种子块(seed chunks)”, 构成了初始的检索集 D_q。
 - 图引导的扩展检索(**Graph-guided Expansion**): 单纯依赖语义相似性检索可能会导致召回的文本块在内容上较为孤立, 缺乏上下文关联。为了克服这一局限, KG²RAG利用之前构建的KG-块关联信息来进行扩展。具体做法是:
 1. 根据种子块 D_q 中包含的实体和关系, 从总的关联KG中抽取出与这些种子块直接相关的子图 G_{q0}。
 2. 然后, 在这个子图 G_{q0} 上应用图遍历算法(例如广度优先搜索BFS), 从种子实体出发, 探索其在KG中的m跳邻域内的其他实体和关系。
 3. 将这些在m跳邻域内新发现的实体和关系所对应的原始文本块也加入到检索结果中。通过这种图引导的扩展, KG²RAG能够将那些与种子块在知识图谱层面存在直接或间接关联(例如, 共享实体、处于同一关系路径上)的文本块也检索出来, 即使这些扩展块与原始查询的直接语义相似度不高。这大大增加了检索结果的多样性和上下文的丰富性, 有助于回答需要多方面信息或进行多跳推理的复杂问题。经过扩展后得到的文本块集合被称为 D_{qm}。
- 3. **基于KG的上下文组织(KG-based Context Organization)**: 在检索并扩展了相关的文本块之后, KG²RAG还引入了一个后处理阶段, 旨在对这些检索到的上下文进行优化和组织, 然后再将其提供给LLM进行最终的答案生成。这个阶段主要包含两个步骤:
 - 过滤(**Filtering**): 首先, 利用扩展检索到的文本块集 D_{qm} 及其在KG中的关联关系, 构建一个局部的、与当前查询相关的无向加权图 U_{qm}。在这个图中, 节点代表实体, 边代表实体间的关系, 边的权重可以根据查询与该关系所在源文本块的语义相似度来设定。然后, 这个图会被分解为若干个连通分量(connected

components)。对于每个连通分量, 系统会计算其最大生成树(Maximum Spanning Tree, MST)。通过这种方式, 可以保留那些最能连接相关信息的关键路径, 同时去除冗余的边和关系, 从而达到过滤噪声、提炼核心上下文的目的。

- 排列(**Arranging**): 接下来, 对于每个筛选出来的MST, KG²RAG会为其生成两种表示: 一种是文本表示, 通过深度优先搜索(DFS)遍历MST中的边(从权重最高的边开始), 并将这些边所对应的原始文本块按照遍历顺序拼接起来, 形成一个或多个逻辑连贯的段落; 另一种是三元组表示, 即MST中包含的所有知识三元组的集合。然后, 利用一个交叉编码器(cross-encoder)或其他重排序模型, 计算每个MST的三元组表示与用户原始查询之间的相关性得分。最后, 将所有MST按照这个相关性得分从高到低进行排序, 并依次将其对应的文本表示(即组织好的段落)作为最终的上下文输入给LLM, 直到达到预设的上下文长度或块数量的限制。

KG²RAG的主要优势在于, 它通过显式地利用知识图谱中蕴含的实体间事实级关系, 有效地克服了传统RAG仅依赖文本语义相似性进行检索的局限性。这使得KG²RAG能够检索到更具多样性、上下文更连贯、逻辑关系更清晰的知识片段, 从而为LLM提供更高质量的输入, 最终显著改善生成答案的质量, 特别是在处理需要多跳推理和复杂信息综合的问答任务时。在HotpotQA等标准数据集上的实验结果也充分证明了KG²RAG相比于现有的多种RAG基线方法的优越性, 无论是在答案的F1值、精确率和召回率, 还是在检索结果本身的质量方面, 都展现出了显著的提升。

从更深层次来看, 知识图谱为RAG系统引入了超越纯文本语义相似性的、宝贵的结构化上下文(**structured context**)。传统的基于向量的检索主要依赖于文本内容在语义空间中的邻近程度, 而KG则清晰地揭示了文本中提及的实体之间存在的各种显式关系, 例如“公司A是公司B的子公司”、“人物C发明了技术D”、“事件E发生在地点F”等等⁴⁸。这种结构化的信息对于RAG系统进行更深层次的理解和推理至关重要, 它能够帮助系统将看似孤立的信息片段有效地串联起来, 形成一个更完整的知识网络。

当然, KG-RAG的成功在很大程度上依赖于所用知识图谱的质量和维持。正如研究所指出的, 现实世界中的KG往往是不完整的⁴⁷。KG²RAG框架通过允许从文本块中自动构建KG的方式, 在一定程度上缓解了对预先存在的高质量、完整KG的强依赖性⁴⁸。然而, 如何保证自动构建的KG的准确性和覆盖范围, 以及如何更有效地处理KG固有的不完整性问题(例如, 通过知识图谱补全技术或在推理时对缺失信息进行鲁棒处理), 仍然是KG-RAG领域需要持续研究和探索的重要方向。

对于知识管理系统而言, KG-RAG技术尤其适用于那些需要处理包含复杂实体关系、需要进行多跳推理才能得到答案的场景。例如, 在企业知识管理中, 回答“公司X的最新产品Y所使用的核心技术Z是由哪个供应商提供的, 该供应商与公司X的历史合作关系如何?”这类问题, 就需要整合来自多个文档(产品手册、供应商合同、历史项目报告等)的信息, 并理解其中涉及的公司、产品、技术、供应商等实体间的复杂关系。KG-RAG通过其图引导的检索扩展和基于KG的上下文组织能力⁴⁸, 能够比传统RAG更有效地找到和组织这些盘根错节的

相关信息,从而生成更全面、更准确、更有深度的答案,极大地提升知识管理系统解决复杂问题的能力。

3.4 多模态RAG (Multimodal RAG - MRAG)

随着信息时代的飞速发展,知识的载体早已不再局限于纯文本。图像、视频、音频、图表等**多模态(multimodal)**数据在信息传播和知识表达中扮演着越来越重要的角色。传统的RAG系统主要针对文本数据进行设计,难以有效利用这些丰富的非文本信息。为了应对这一挑战,**多模态RAG(Multimodal RAG, MRAG)**应运而生,它旨在将RAG的理念从单一的文本模态扩展到能够整合和利用多种数据类型(如文本、图像、视频、音频等)的检索和生成过程⁸。

MRAG的核心思想是通过检索与用户查询相关的多模态信息,并利用多模态大型语言模型(**Multimodal Large Language Models, MLLMs**) (如OpenAI的GPT-4V、Google的Gemini、LLaVA等)来结合这些检索到的多模态上下文,生成更为全面和准确的答案。这不仅能够减少LLM产生幻觉的可能性,还能显著增强问答系统在需要联合理解视觉和文本信息等复杂场景下的表现。

MRAG技术的发展也经历了一个演进过程:

- **MRAG 1.0 的局限性**⁶²:早期的MRAG尝试通常是将非文本的多模态数据(如图像)预先转换成文本描述(例如,使用图像字幕生成模型为图像生成caption),然后再将这些文本描述纳入传统的文本RAG流程中。这种做法虽然在一定程度上实现了多模态信息的利用,但也带来了显著的局限性:
 - 信息损失:将图像等丰富的多模态信息强行压缩成简短的文本描述,不可避免地会导致大量细粒度细节和复杂语义信息的丢失。例如,一张包含复杂场景或微妙情感的图片,其对应的文本字幕往往只能概括其主要内容,难以完全传达原始图像的全部信息。
 - 系统复杂性与计算开销增加:需要为不同模态的数据分别设计和维护转换模型(如图像字幕模型、语音识别模型等),这增加了系统的整体复杂性和计算开销。
 - 检索瓶颈:基于转换后文本的检索,仍然受到传统文本分块策略和嵌入模型能力的限制。而且,转换过程本身引入的信息损失,也会进一步影响后续检索的准确性和召回率。
 - 生成挑战:如何有效地将来自不同模态(即使已转换为文本)的信息片段组织成连贯、有意义的提示,以供LLM使用,本身就是一个难题。输入质量的下降(GIGO原则)也使得LLM更容易生成不可靠的答案。
- **MRAG 2.0 (真正的多模态时代)**⁶²:为了克服MRAG 1.0的不足,新一代的MRAG系统(可称为MRAG 2.0)致力于在整个RAG流程中尽可能地保留和利用多模态数据的原始形式。其关键特征包括:
 - 多模态文档解析与索引:系统能够直接处理和索引包含混合内容(如文本、图像、表格并存的PDF文档)的多模态文档,而不是先将其完全文本化。这可能涉及到对不

同模态内容分别进行特征提取和嵌入，或者采用统一的多模态嵌入模型。

- 多模态检索 (**Multimodal Retrieval**): 检索模块需要具备根据用户(可能是文本、图像甚至混合模态的)查询, 从多模态知识库中检索出最相关的多模态信息片段(例如, 一段包含关键图表的文字, 或者一张与查询图片内容相似的图片及其附带说明)的能力。这要求有强大的跨模态相似性计算和检索算法。
- 多模态生成 (**Multimodal Generation**): 利用先进的MLLM作为生成器。这些MLLM本身就具备同时理解和处理多种输入模态(如文本和图像)并生成相应输出(通常是文本, 但未来也可能是多模态的)的能力。检索到的多模态上下文可以直接(或经过适当格式化后)作为MLLM的输入提示。

MRAG的关键组件通常包括⁵⁸:

- 多模态知识库 (**Multimodal Knowledge Base**): 存储包含文本、图像、视频、音频等多种数据类型的原始信息及其对应的嵌入表示。
- 多模态嵌入模型 (**Multimodal Embedding Models**): 能够将不同模态的数据(或其组合)映射到同一个或多个相互关联的语义向量空间, 以便进行跨模态的相似性比较和检索。
- 多模态检索器 (**Multimodal Retriever**): 负责根据用户查询(可以是单模态或多模态的)从知识库中检索最相关的多模态信息片段。
- 多模态大型语言模型 (**MLLM**): 作为生成器, 它接收用户查询和检索到的多模态上下文作为输入, 并生成最终的答案。

构建和评估MRAG系统需要专门的数据集, 这些数据集应包含多模态的问答对、文档或场景。评估MRAG系统也比评估纯文本RAG更为复杂, 需要综合考虑系统对不同模态信息的理解、整合能力以及生成答案的准确性、相关性和忠实度, 有时还需要评估生成内容与特定模态(如图像)的一致性⁵⁸。

尽管MRAG展现出巨大的潜力, 但其发展仍面临诸多挑战⁵⁸:

- 跨模态幻觉 (**Cross-modal Hallucination**): MLLM在处理和融合多种模态信息时, 仍有可能产生不一致或错误的关联, 例如错误地描述图像内容或将不相关的文本与图像联系起来。
- 复杂的评估需求 (**Complex Evaluation Requirements**): 如何全面、准确地评估MRAG系统的性能, 特别是在涉及主观判断(如图像与文本的关联是否自然)或细粒度跨模态理解的方面, 仍然是一个难题。
- 高质量多模态数据集的稀缺 (**Scarcity of High-Quality Multimodal Datasets**): 训练和评测强大的MRAG系统需要大规模、高质量、多样化的标注多模态数据集, 而这类资源的获取成本较高。
- 有效的多模态提示工程 (**Effective Multimodal Prompt Engineering**): 如何设计能够有效引导MLLM利用检索到的多模态上下文的提示, 是一个需要专门研究的问题。
- 可扩展性与成本效益 (**Scalability and Cost-effectiveness**): 处理和存储多模态数

据(尤其是视频等)通常比纯文本需要更多的计算资源和存储空间,这对MRAG系统的可扩展性和部署成本提出了挑战。

- **伦理考量(Ethical Considerations)**: MRAG系统在处理 and 生成包含图像、音频等内容时,更容易涉及到偏见放大、隐私泄露、不当内容生成等伦理风险,需要格外关注。

MRAG是RAG技术发展的一个必然且重要的方向,它顺应了现实世界知识日益多媒体化、信息获取和交互日益追求丰富感官体验的趋势。传统的纯文本RAG系统在面对包含图表、照片、流程图、视频演示等非文本知识时显得捉襟见肘,而MRAG的出现⁵⁸则使得知识管理系统能够真正开始整合和利用这些过去难以触及的多模态数据金矿,从而提供远比以往更全面、更直观、更具吸引力的知识服务。

在MRAG成功的背后,**多模态嵌入(multimodal embeddings)和跨模态对齐(cross-modal alignment)**是其核心的技术挑战。如何有效地将来自不同感官通道的信息(例如,一段描述性文字和一张与之相关的产品图片)表示在同一个或多个能够相互理解的语义向量空间中,并基于这种表示实现它们之间的准确、高效的相似性匹配和检索,是决定MRAG系统性能好坏的关键所在⁶¹。这不仅需要更强大的多模态嵌入模型(如CLIP及其后续发展模型),还需要更先进的跨模态检索算法和索引结构。

MRAG的出现为知识管理带来了诸多前所未有的机遇。想象一下,未来的知识管理系统用户可以用一张设备故障的照片来提问,系统则能通过MRAG检索相关的维修手册(其中可能包含详细的文字说明、部件分解图和操作视频),并最终给出一个图文并茂、甚至包含短视频演示的解决方案。或者,当用户查询一个复杂的业务流程时,系统不再仅仅返回冗长的文字描述,而是能够生成一个包含关键节点图示、数据流向箭头和核心指标解读的综合性视图。这些都极大地拓展了知识管理的应用场景和用户价值。然而,要充分释放MRAG的潜力,还需要在克服上述诸多挑战方面取得持续的突破。

3.5 面向大规模知识库的RAG优化 (Optimizing RAG for Large-Scale Knowledge Bases)

随着企业和组织积累的数据量呈指数级增长,知识管理系统需要处理的知识库规模也日益庞大。当RAG系统面对包含数百万甚至数千万文档的**大规模知识库(Large-Scale Knowledge Bases)**时,如何保持检索的效率和准确性,以及如何确保LLM能够有效地从可能海量的检索结果中筛选和利用真正有

Works cited

1. aws.amazon.com, accessed May 7, 2025, <https://aws.amazon.com/what-is/retrieval-augmented-generation/#:~:text=Augmented%20Generation%20requirements%3F-,What%20is%20Retrieval%2DAugmented%20Generation%3F,sources%20before%20generating%20a%20response.>
2. What is RAG? - Retrieval-Augmented Generation AI Explained - AWS, accessed May 7, 2025, <https://aws.amazon.com/what-is/retrieval-augmented-generation/>

3. What is Retrieval-Augmented Generation (RAG)? | Google Cloud, accessed May 7, 2025, <https://cloud.google.com/use-cases/retrieval-augmented-generation>
4. Retrieval augmented generation (RAG) | 🦜 LangChain, accessed May 7, 2025, <https://python.langchain.com/docs/concepts/rag/>
5. 5 key features and benefits of retrieval augmented generation (RAG) | The Microsoft Cloud Blog, accessed May 7, 2025, <https://www.microsoft.com/en-us/microsoft-cloud/blog/2025/02/13/5-key-features-and-benefits-of-retrieval-augmented-generation-rag/>
6. Retrieval-Augmented Generation (RAG): Transform Enterprise Knowledge - Stratechi.com, accessed May 7, 2025, <https://www.stratechi.com/retrieval-augmented-generation-ai-rag-knowledge-management/>
7. Addressing ResNetVLLM's Multi-Modal Hallucinations University of Windsor. - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2504.14429v1>
8. A Survey on Knowledge-Oriented Retrieval-Augmented Generation - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2503.10677v2>
9. AI and knowledge management: Why RAG is essential - Outshift | Cisco, accessed May 7, 2025, <https://outshift.cisco.com/blog/using-ai-knowledge-management-why-rag-is-essential>
10. How RAG Enhances Knowledge Management in IT Support Teams - Algomox, accessed May 7, 2025, https://www.algomox.com/resources/blog/how_rag_enhances_knowledge_management_in_it_support_teams.html
11. Retrieval-Augmented Generation for Large Language Models: A Survey - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2312.10997v5>
12. arxiv.org, accessed May 7, 2025, <https://arxiv.org/abs/2312.10997>
13. Chunking Strategies For Production-Grade RAG Applications, accessed May 7, 2025, <https://www.helicone.ai/blog/rag-chunking-strategies>
14. Chunking strategies for RAG tutorial using Granite | IBM, accessed May 7, 2025, <https://www.ibm.com/think/tutorials/chunking-strategies-for-rag-with-langchain-watsonx-ai>
15. Enhancing RAG performance with smart chunking strategies - IBM ..., accessed May 7, 2025, <https://developer.ibm.com/articles/awb-enhancing-rag-performance-chunking-strategies/>
16. 7 Chunking Strategies in RAG You Need To Know - F22 Labs, accessed May 7, 2025, <https://www.f22labs.com/blogs/7-chunking-strategies-in-rag-you-need-to-know/>
17. Get better RAG by fine-tuning embedding models - Redis, accessed May 7, 2025, <https://redis.io/blog/get-better-rag-by-fine-tuning-embedding-models/>
18. How to Fine-Tune Embedding Models for RAG Systems - Dataworkz ..., accessed May 7, 2025, <https://www.dataworkz.com/blog/how-to-fine-tune-embedding-models-for-rag-systems/>

19. FareedKhan-dev/all-rag-techniques: Implementation of all ... - GitHub, accessed May 7, 2025, <https://github.com/FareedKhan-dev/all-rag-techniques>
20. LLM-QE: Improving Query Expansion by Aligning Large Language Models with Ranking Preferences - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2502.17057v1>
21. arxiv.org, accessed May 7, 2025, <https://arxiv.org/abs/2502.18139>
22. arxiv.org, accessed May 7, 2025, <https://arxiv.org/abs/2403.14403>
23. Advanced RAG Technique : Langchain ReAct and Cohere, accessed May 7, 2025, <https://www.analyticsvidhya.com/blog/2024/05/advanced-rag-technique-langchain-react-and-cohere/>
24. [2410.12886] AT-RAG: An Adaptive RAG Model Enhancing Query Efficiency with Topic Filtering and Iterative Reasoning - arXiv, accessed May 7, 2025, <https://arxiv.org/abs/2410.12886>
25. MrRezaeiUofT/AT-RAG: AT-RAG (Adaptive Retrieval ... - GitHub, accessed May 7, 2025, <https://github.com/MrRezaeiUofT/AT-RAG>
26. Re-ranking the Context for Multimodal Retrieval Augmented Generation - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2501.04695v1>
27. arxiv.org, accessed May 7, 2025, <https://arxiv.org/abs/2501.04695>
28. Advanced RAG Implementation using Hybrid Search - Athina AI Hub, accessed May 7, 2025, <https://hub.athina.ai/athina-originals/advanced-rag-implementation-using-hybrid-search/>
29. 15 Best Open-Source RAG Frameworks in 2025 - Firecrawl, accessed May 7, 2025, <https://www.firecrawl.dev/blog/best-open-source-rag-frameworks>
30. RAG Prompt Engineering Makes LLMs Super Smart - K2view, accessed May 7, 2025, <https://www.k2view.com/blog/rag-prompt-engineering/>
31. RAG vs fine-tuning vs. prompt engineering - IBM, accessed May 7, 2025, <https://www.ibm.com/think/topics/rag-vs-fine-tuning-vs-prompt-engineering>
32. Corrective Retrieval Augmented Generation - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2401.15884v2>
33. Corrective Retrieval Augmented Generation - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2401.15884v3>
34. Mitigating LLM Hallucination with Smoothed Knowledge Distillation - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2502.11306v1>
35. Evaluating the evaluators: know your RAG metrics - Tweag, accessed May 7, 2025, <https://tweag.io/blog/2025-02-27-rag-evaluation/>
36. Retrieval Augmented Generation Evaluation in the Era of Large Language Models: A Comprehensive Survey - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2504.14891v1>
37. arxiv.org, accessed May 7, 2025, <https://arxiv.org/abs/2401.15884>
38. Self-RAG: Learning to Retrieve, Generate and Critique through Self-Reflection, accessed May 7, 2025, <https://selfrag.github.io/>
39. Self-Rag: Self-reflective Retrieval augmented Generation - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2310.11511>
40. [2502.11306] Smoothing Out Hallucinations: Mitigating LLM Hallucination with

- Smoothed Knowledge Distillation - arXiv, accessed May 7, 2025, <https://arxiv.org/abs/2502.11306>
41. accessed January 1, 1970, <https://arxiv.org/pdf/2502.11306.pdf>
 42. LettuceDetect: A Hallucination Detection Framework for RAG Applications - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2502.17125v1>
 43. [2502.06148] Optimizing Knowledge Integration in Retrieval-Augmented Generation with Self-Selection - arXiv, accessed May 7, 2025, <https://arxiv.org/abs/2502.06148>
 44. Optimizing Knowledge Integration in Retrieval-Augmented Generation with Self-Selection, accessed May 7, 2025, <https://arxiv.org/html/2502.06148v1>
 45. accessed January 1, 1970, <https://arxiv.org/pdf/2503.10677.pdf>
 46. [2503.10677] A Survey on Knowledge-Oriented Retrieval-Augmented Generation - arXiv, accessed May 7, 2025, <https://arxiv.org/abs/2503.10677>
 47. [2504.05163] Evaluating Knowledge Graph Based Retrieval Augmented Generation Methods under Knowledge Incompleteness - arXiv, accessed May 7, 2025, <https://arxiv.org/abs/2504.05163>
 48. arxiv.org, accessed May 7, 2025, <https://arxiv.org/abs/2502.06864>
 49. Self-Routing RAG: Binding Selective Retrieval with Knowledge ..., accessed May 7, 2025, <https://www.aimodels.fyi/papers/arxiv/self-routing-rag-binding-selective-retrieval-knowledge>
 50. [2504.01018] Self-Routing RAG: Binding Selective Retrieval with Knowledge Verbalization, accessed May 7, 2025, <https://arxiv.org/abs/2504.01018>
 51. GreatHayat/langgraph-corrective-rag: A repository to learn ... - GitHub, accessed May 7, 2025, <https://github.com/GreatHayat/langgraph-corrective-rag>
 52. Kirouane-Ayoub/Corrective-RAG - GitHub, accessed May 7, 2025, <https://github.com/Kirouane-Ayoub/Corrective-RAG>
 53. Self-Routing RAG: Binding Selective Retrieval with Knowledge Verbalization, accessed May 7, 2025, https://www.researchgate.net/publication/390404551_Self-Routing_RAG_Binding_Selective_Retrieval_with_Knowledge_Verbalization
 54. Self-Routing RAG: Binding Selective Retrieval with Knowledge Verbalization - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2504.01018v1>
 55. accessed January 1, 1970, <https://arxiv.org/pdf/2504.01018>
 56. accessed January 1, 1970, <https://arxiv.org/pdf/2504.01018.pdf>
 57. sksarvesh007/adaptive-rag - GitHub, accessed May 7, 2025, <https://github.com/sksarvesh007/adaptive-rag>
 58. [2504.08748] A Survey of Multimodal Retrieval-Augmented Generation - arXiv, accessed May 7, 2025, <https://arxiv.org/abs/2504.08748>
 59. AlzheimerRAG: Multimodal Retrieval Augmented Generation for PubMed articles - arXiv, accessed May 7, 2025, <https://arxiv.org/abs/2412.16701>
 60. Computer Science Apr 2025 - arXiv, accessed May 7, 2025, <https://www.arxiv.org/list/cs/2025-04?skip=4025&show=25>
 61. Multimodal LLM Guide: Addressing Key Development Challenges ..., accessed May 7, 2025, <https://www.galileo.ai/blog/multimodal-llm-guide-evaluation>

62. A Survey on Multimodal Retrieval-Augmented Generation - arXiv, accessed May 7, 2025, <https://arxiv.org/html/2504.08748v1>
63. accessed January 1, 1970, <https://arxiv.org/pdf/2504.08748>
64. accessed January 1, 1970, <https://arxiv.org/pdf/2504.08748.pdf>