

# 检索增强生成技术支持下的校园问答系统研究

贾春燕<sup>1</sup>, 方伟杰<sup>1</sup>, 谢宇威<sup>1</sup>, 凌在盈<sup>2</sup>

(1. 浙江大学信息技术中心, 浙江 杭州 310027; 2. 杭州师范大学信息科学与技术学院, 浙江 杭州 310027)

**摘要:** 针对高等学校师生用户从海量校园信息中获取有效信息的困难, 以校务领域知识为数据源, 基于检索增强生成技术, 设计了一个校园智能问答系统。融合大语言模型和垂直领域专业知识, 以学校百事通项目为依托, 将包括办事指南、常见问题、规范性文件等校务信息作为外挂数据语料库, 应用检索增强生成专用的 Infinity 数据库, 构建校务知识库, 采用提示词工程, 增强大语言模型生成答案。通过检索增强生成技术进行教育领域特定的校园问答, 旨在以互动方式为用户提供各种校务服务信息, 有助于解决校园常见问题, 简化师生咨询流程, 减轻学校管理工作负担。

**关键词:** 信息获取; 大语言模型; 检索增强生成; 校园问答

**中图分类号:** TP391

**文献标志码:** A

**DOI:** 10.11959/j.issn.1000-436x.2024259

## Research on campus question answering system supported by retrieval-augmented generation technology

JIA Chunyan<sup>1</sup>, FANG Weijie<sup>1</sup>, XIE Yuwei<sup>1</sup>, LING Zaiying<sup>2</sup>

1. Information Technology Center, Zhejiang University, Hangzhou 310027, China

2. School of Information Science and Engineering, Hangzhou Normal University, Hangzhou 310027, China

**Abstract:** To solve the problem of obtaining effective information from the vast amount of campus information for teachers and students, an intelligent campus question answering (QA) system based on RAG was designed. An approach that integrates large language models and domain knowledge for QA system construction was proposed, relying on the campus's *Everything You Need to Know* project, and using campus information such as procedural guides, frequently asked questions, and normative documents as an external data corpus. A campus knowledge database was constructed, with the RAG Infinity database. To improve the retrieval efficiency of domain knowledge and the accuracy of answers, the prompt approach was proposed. Using RAG for campus QA, the system provides users various service information in an interactive manner, which helps to solve common campus issues, simplify the consultation process for teachers and students, and alleviate the burden on campus management, and enrich campus knowledge resources.

**Keywords:** information retrieval, large language model, retrieval-augmented generation, campus QA

### 0 引言

当前, 伴随高等学校数字化转型和智慧校园工作的深入推进, 校园网络数据和信息呈爆炸式增长<sup>[1]</sup>, 师生用户一般通过学校自建的搜索系统获取有效信息, 但是由于搜索系统基于关键词匹配, 不

涉及语义, 用户无法使用自然语言准确表达信息需求, 导致搜索系统无法准确理解用户查询的真实意图, 并且搜索返回结果太多, 很难快速准确定位到用户所需信息。随着校园应用和数据的持续增加, 用户通过搜索获取有效信息的难度也在不断加大, 因此, 各个高校应该如何构建一个更加高效统一的

收稿日期: 2024-10-22

通信作者: 方伟杰, fangwj@zju.edu.cn

信息获取平台,成为目前亟待解决的问题<sup>[2-3]</sup>。

近年来,随着人工智能技术的不断成熟,基于大语言模型(LLM, large language model)的智能问答得到广泛应用,智能问答以智能化的人机交互方式帮助用户从海量信息中快速获取有效信息,可以弥补传统搜索的不足。融合外部知识的检索增强生成(RAG, retrieval-augmented generation)技术与提示词工程是当前 LLM 领域的主要研究热点<sup>[4-8]</sup>。

为了构建教育领域校务方面的智能问答系统,本文提出融合大语言模型和校务领域专业知识,探索一种基于检索增强生成的校园问答技术实现路径。研究以学校百事通项目为依托,将包括办事指南、常见问题、规范性文件等校务信息作为外挂数据语料库,应用 RAG 专用的 Infinity 数据库,构建校务知识库,采用提示词工程,增强 LLM 生成答案。通过 RAG 进行教育领域的校务问答,旨在建立一个统一高效的信息获取平台,以互动方式为用户提供各种校务服务信息,有助于简化师生咨询流程,减轻学校管理负担,丰富学校的校务知识资源<sup>[8-11]</sup>。

## 1 相关技术

### 1.1 大语言模型

大语言模型,如 GPT-4<sup>[12]</sup>、LLAMA2<sup>[13]</sup> 和 ChatGLM-2<sup>[14]</sup>,是基于拥有十亿至数千亿参数的 Transformer 架构,在计算语言学中占据前沿地位。这些模型依赖 Transformer 内部的自注意力机制。LLM 擅长理解和生成人类语言,从而改变了自然语言处理的格局。它们利用各种 Transformer 架构和预训练目标,包括仅解码器模型(如 GPT2、GPT3)、仅编码器模型(如 BERT<sup>[15]</sup>)以及编码器-解码器架构(如 BART)。这些架构能够有效处理序列数据,捕捉文本内部的复杂依赖关系,同时支持有效的并行化。LLM 结合了提示或上下文学习,通过整合上下文信息来增强文本生成能力,从而促进连贯且与上下文相关的回答,推动互动问答的参与度<sup>[16]</sup>。大语言模型展现了强大的自然语言理解能力和解决复杂任务的能力,许多传统的自然语言处理任务在 LLM 的帮助下正在变得更加简化,如问答系统。

大语言模型作为一种预训练的语言生成模型,

在通用领域智能问答中展现出了强大的能力,但是,将其直接应用于教育垂直领域尚存在一些限制和挑战。

1) 安全性问题。目前通用领域的 LLM 问答系统,将训练数据统一输入 LLM 进行训练,存在敏感信息泄露风险。

2) 知识更新慢。LLM 的结构化参数通过训练方法更新,不会实时更新知识库,存在信息过时问题,在处理需要最新信息的任务时会受限。

3) 计算资源需求高。训练和运行大语言模型需要大量的计算资源,会导致高昂的经济成本和环境影响。

4) 数据依赖度高。LLM 的性能在很大程度上依赖于训练数据的质量和多样性。如果训练数据存在偏差或不均衡,模型的表现也会受到影响。

5) 多模态能力的挑战。尽管 LLM 在文本处理方面表现出色,但它们在处理图像、音频等非文本模态数据时仍面临挑战。

6) 生成幻觉。LLM 无法完成超出其训练数据范围的内容生成,可能会生成听起来合理但实际上是错误的信息,这种现象被称为“幻觉”。如果模型训练得当,幻觉可以最小化,但无法完全消除<sup>[17-18]</sup>。

### 1.2 检索增强生成

RAG 由 Meta AI 于 2020 年提出<sup>[19]</sup>,结合检索领域专业知识库和 LLM,在 LLM 生成答案之前,先从外部知识库中检索相关信息,提高生成内容的可信度与准确度。RAG 使大模型能够通过外部知识库获得额外的知识扩充,增强模型对背景知识的综合理解,可以有效缓解 LLM 的生成幻觉,提高知识更新的速度,增强内容生成的可追溯性。RAG 技术的核心包括 3 个主要阶段:索引、检索和生成。索引阶段构建高效的数据结构以快速存储和访问海量信息;检索阶段基于向量相似性原理提取相关信息;生成阶段利用 LLM 结合检索到的信息生成高质量的输出<sup>[20-22]</sup>。

RAG 于 2024 年 4 月 1 日开源了新一代的引擎 RAGFlow<sup>[23]</sup>,采用深度文档理解和光学字符识别(OCR)技术,具备深度文档理解、兼容异构数据源、可控可解释文本切片、自动化 RAG workflow 等能力,是一个完整的端到端 RAG 解决方案。RAG-Flow 通过一系列精心设计的组件,实现了对复杂

查询的快速响应和精准处理,包括文档解析器、问题分析器、知识检索、重排序和 LLM 这 5 个核心组件。

文档解析器负责将各种格式的文档进行解析,从中提取出文本、图像和表格等关键内容,是 RAGFlow 系统的“大脑”。

问题分析器是 RAGFlow 系统的“神经系统”,对用户输入的问题进行深入分析,识别并提取出查询中的关键信息,通过分析,系统能够更准确地理解用户的需求,为检索工作提供精确的指导。

知识检索是 RAGFlow 系统的“搜索引擎”,使用问题分析器提供的关键信息,从海量文档中快速检索与之相关的信息。

重排序对检索到的信息进行排序和过滤,确保最终呈现给用户的信息是最相关、最有价值的。通过这种方式,系统能够去除冗余和不相关的数据,提高信息的准确性和可用性,这个组件是 RAGFlow 系统的“过滤器”。

大语言模型是 RAGFlow 系统的“语言生成器”,负责将排序后的信息整合并生成最终的答案或输出,其强大的生成能力不仅能够确保答案的准确性,还能够使答案表达得更加自然和流畅。

这些组件共同构成了 RAGFlow 系统的强大架构,使得它能够高效地处理用户的查询,快速地从文档中检索信息,并生成准确、有用的答案。

### 1.3 知识库基础

知识库的建设是校园问答系统建设的重点。浙江大学于 2019 年底开始,由学校统筹,机关党委协同各机关单位及学院各部门打造校园百事通项目。项目除了通过爬虫采集浙大校内的网站资讯信息外,还依托校务信息库,以“三张清单”为重要工作基础,结合各单位办事指南、常见问题、部门通讯录、规范性文件、报告讲座、教室信息查询等内容的统一采集与管理,整合校园网络信息资源,建立统一的校园信息检索库,目前总检索次数已近 500 万。此项目可以为构建校务知识库提供数据基础,项目前期建立的数据建设维护组织体系也可为校务问答知识库提供组织和制度保障。

根据项目前期数据梳理,可以根据信息来源和表现形式将校务知识进行分类,针对不同的知识类型,构建不同的查询方法,为后续知识库构建和知识检索做准备。

问答类知识。学校各部门办事指南、常见问题、规范性文件等“一问一答”(FAQ)型知识,来源于各部门发布维护,可以采用问题-答案对策略,基于关键词、语义和句法结构的方法计算问句相似度,查询获取答案。

文本类知识。校园网站上的网页包含的内容属于文本检索类知识,来源广泛,以自然语言文本形式表示,通过信息检索和答案抽取获取答案。

服务类知识。教室空闲状态等服务类知识来源于管理部门的业务管理系统中,是格式固定的结构化数据,存储于关系型数据库,通过 SQL 查询语句转化成查询指令获取答案。

## 2 系统架构及问答实现

### 2.1 系统架构

基于 RAG 的校园问答系统架构共分四层,分别为基础设施层、数据资源层、平台层和展现层,如图 1 所示。

系统最底层为基础设施层,由学校数字基础建设提供,包括高速稳定的校园网络、安全智能的云计算平台和标准领先的算力中心等设备。

数据资源层主要负责结构化、非结构化源数据的处理和存储,包括数据采集、数据预处理和存储 3 个模块。数据采集模块利用业务同步以及网络爬虫从校务信息库、协同办公、教务等相关系统或网站抽取可用的问答语料,形成问答数据,再由数据预处理模块对获取的数据做清洗、文本解析和格式转换、切分、索引向量化等处理,为后续算法实现做准备。

平台层由问答模块和管理模块组成,问答模块由问题分类、问题聚类、知识检索、重排序、提示词、答案生成 6 个子模块组成,利用从数据资源层获取的数据协作完成问答功能。系统管理模块包含用户管理、数据管理、推送消息、模型更新子模块。模型更新子模块负责问答系统中各个子模块的更新。

最顶层是展现层,提供用户 Web、APP 以及业务应用程序接口(API)等多终端的用户界面,实现系统与用户的交互。

### 2.2 基于 RAGFlow 的问答实现

采用 RAGFlow 提供的组件,可以快速便捷地构建校园问答模块,具体实现流程如图 2 所示,核

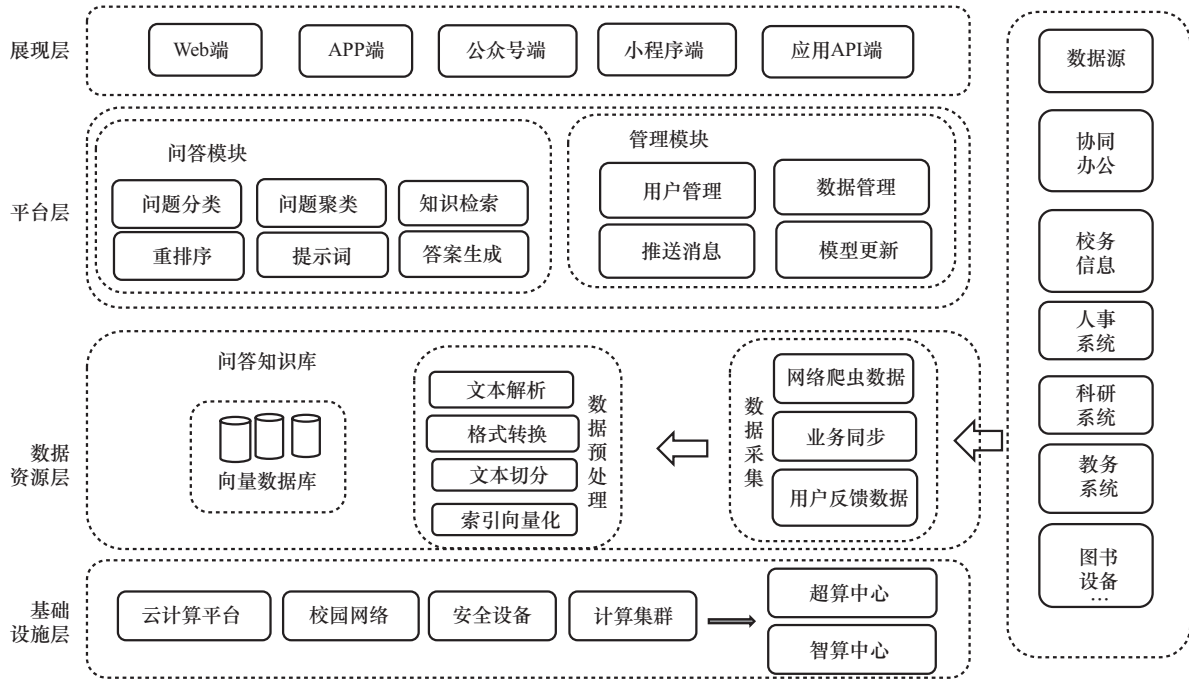


图1 基于RAG的校园问答系统架构

心步骤包括数据准备、知识库构建、知识检索和答案生成。

### 1) 数据准备

采集校务信息相关的结构化和非结构化数据，以构建校务领域知识库。利用RAGFlow可以兼容各类异构数据源的能力，对文档系统中的各类校务相关文件（包含word、ppt、excel、txt、PDF、图片、影印件、复印件、结构化数据、网页等格式）进行预处理，通过解析技术对文档进行格式转换、

解析，然后进行内容抽取。对于无序文本数据，模型可以自动提取其中的关键信息并转化为结构化表示，对于结构化数据可以灵活切入，挖掘内在语义联系。最终将2种不同来源的数据统一进行索引和检索，为用户提供一站式的数据处理和问答体验。

### 2) 知识库构建

应用RAG专用的Infinity数据库，构建校务知识库，实现检索多路召回。主要包括文档加载、切分、向量化及索引入库4个模块。利用RAGFlow

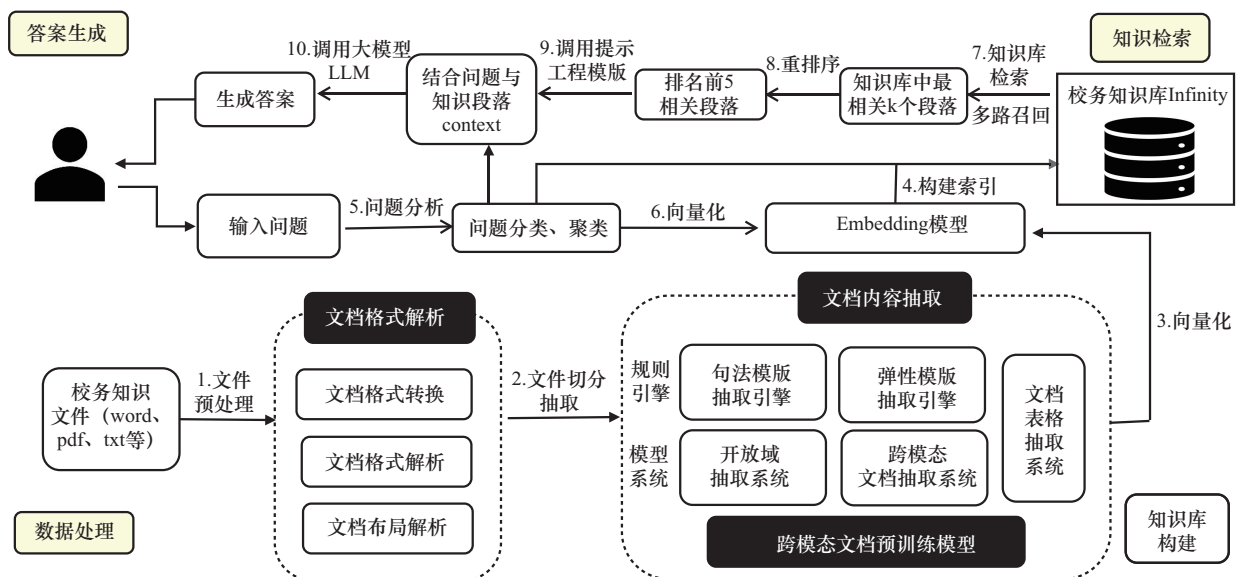


图2 问答实现流程

可控可解释的文本切片能力, 根据不同需求选择合适的文本模板, 确保结果的可控性和可解释性。校务知识的存储包含结构化和非结构化知识库 2 种方式。表格类型数据存储在数据库中, 通过查询语句完成检索。文本型数据存储在非结构化数据库中, 采用稀疏检索和稠密向量检索相结合的方式建立非结构化数据库的索引, 采用 BM25 算法和关键词检索相结合的方式稀疏检索构建向量库。

### 3) 知识检索

使用问题分析后提供的关键信息, 从预先构建的知识库中快速检索出与用户问题相关的信息。针对用户提出的校务信息相关问题进行深入分析, 识别并提取问题中的关键信息, 可以进行问题改写、拓展优化, 在知识库中进行文本块向量检索, 多路召回最相关知识库中向量检索的文本段落。BM25 是一种词袋检索功能, 给定用户问题  $Q$ , 包含词  $q_1, \dots, q_n$ , 一个文本片段  $D$  的 BM25 分数为  $\text{score}(D, Q) =$

$$\sum_{i=1}^n \text{IDF}(q_i) \frac{f(q_i, D)(k_1 + 1)}{f(q_i, D) + k_1(1 - b + b \frac{|D|}{\text{avgdl}})} \quad (1)$$

其中,  $f(q_i, D)$  为  $q_i$  在文本片段  $D$  中出现的频数,  $|D|$  是文本片段  $D$  的长度,  $\text{avgdl}$  为文本片段集合的平均长度。  $k_1$  和  $b$  为超参数。  $\text{IDF}(q_i)$  是每个词  $q_i$  的逆文档频率。计算方法为

$$\text{IDF}(q_i) = \ln \left( \frac{N - n(q_i) + 0.5}{n(q_i) + 0.5} + 1 \right) \quad (2)$$

其中,  $N$  为文本片段集合的总数量,  $n(q_i)$  为包含  $q_i$  的文档数。

此外, 系统借助大模型生成每一个文本片段  $D$  中的关键词  $\text{Key}_D$  来更加简要地表征文本片段的内容。使用 Embedding 模型对每个文本片段  $D$  进行编码, 得到每个文本片段对应的向量  $\mathbf{v}_D$ 。

检索过程中, 用户问题  $Q$  通过 BM<sub>25</sub> 检索得到每一个文本片段对应的相关分数  $\text{score}_{\text{bm25}}(D, Q)$ , 此外将用户问题通过句子向量模型进行编码得到向量  $\mathbf{v}_Q$ , 然后  $\mathbf{v}_Q$  与所有文本片段对应的向量  $\mathbf{v}_D$  通过内积计算相关分数  $\text{score}_{\text{dense}}(D, Q)$ 。根据 2 种检索方法, 算混合分数为

$$\text{score}_{\text{hybrid}} = \beta \text{score}_{\text{bm25}}(D, Q) + (1 - \beta) \text{score}_{\text{dense}}(D, Q) \quad (3)$$

根据混合分数选取分数最高的  $k$  个文本片段。

上述检索过程中, 算法更加侧重于共现词汇,

对于语义的相关性的检索能力较差, 因此本文使用 ROM (representation overlap model) 语义相关性模型对检索到的文本片段进行重排序, 将问题分别与 50 个文本片段输入 ROM 中, 计算相关性分数, 然后选取分数最高的 5 个文本片段。

对于结构化数据库的索引, 对表格进行预处理, 对于普通表格, 直接构造相应的数据库表格; 对于存在合并列的表格, 对表格进行拆分, 重组多个表格后再录入数据库中。

对于结构化数据库的检索, 首先使用大模型将用户问题转化为 SQ 语句, 然后通过查询语句查询相应的数据库得到相应的答案。

### 4) 答案生成

采用提示词工程模板, 合并知识检索后的相关段落与用户问题形成 context, 调用大语言模型生成符合上下文检索的答案。RAGFlow 的文本切片过程可视化, 让用户随时查看 LLM 是基于哪些原文生成答案, 同时生成原文的引用链接, 并允许用户的鼠标悬停上去即可调出原文内容。支持手动调整, 答案提供关键引用的快照并支持追根溯源, 降低幻觉风险。

本文将用户问题分为两大类: 非结构化和结构化问答。针对校务领域场景, 根据用户意图识别结果, 设计直接生成、非结构化查询、结构化查询三类意图以及相应问答流程。直接生成: 用户提问不通过校务知识检索, 直接经过提示词工程构建指令后回答。非结构化查询: 用户问题首先通过混合检索, 从非结构化知识库中检索和问题最相关的  $k$  个文本片段; 然后问题与最相关的文本片段通过提示词工程生成指令输入 LLM, LLM 生成答案给用户。结构化查询: 用户问题首先通过混合检索方式, 从结构化知识库中检索和问题最相关的  $k$  个表格; 然后问题最相关表格信息通过提示词工程生成 SQL 语句指令输入 LLM, 再根据 SQL 语句去关系型数据库中查询, 再次警告提示词工程构建回复指令输入 LLM 中生成相应答案回复给用户。

最后将检索和生成模块无缝集成, 并通过持续的训练和优化提高系统的准确性和效率。

## 3 结束语

为了解决师生用户在海量校园信息中准确获取有效信息的困难, 针对大语言模型在教育领域落地

面临的生成幻觉、知识更新慢、计算资源高等问题, 本文通过对相关文献和开源项目调研, 探索了一种基于检索增强生成的校园问答技术实现路径。融合垂直领域专业知识和大语言模型, 采用 RAG-Flow 开源框架, 依托浙江大学的校园百事通项目, 将办事指南、常见问题、规范性文件等校务信息作为外部数据语料库, 建立统一高效的校园问答系统, 以互动方式为用户提供各种校务服务信息, 简化师生咨询流程, 减轻学校管理负担。

检索增强生成是大模型落地具体应用场景的最佳路径, 本文基于检索增强生成技术实现校园问答平台的研究, 该技术路线同样适用于其他高校或单位, 特别是对于拥有自有知识管理体系, 不方便对外共享, 想应用大模型但自身算力有限的单位。在应用检索增强生成技术构建问答系统时, 需要重点关注两方面工作: 首先应重视知识库的构建, 并且有机制定期更新知识库, 确保单位有一个高质量的知识库或文档集合; 其次要确保检索系统和生成模型无缝集成, 实现端到端的系统整合。

大语言模型落地教育领域尚处于起步阶段, 采用检索增强技术, 将大语言模型与校务领域专业知识融合是一种崭新的尝试, 未来如何正式应用于智慧校园并形成完善的服务体系尚需进一步研究实践。

## 参考文献:

- [1] Tony H, Stewart T, Kristin T. 第四范式: 数据密集型科学发现[M]. 潘教峰, 张晓林, 等译. 北京: 科学出版社, 2012.
- [2] 陈帅朴, 刘芳霖, 钱宇星, 等. 检入新境: 大语言模型引领的信息检索主题与知识关联演化分析[J/OL]. 图书情报知识, (2024-06-27)[2024-10-20].  
CHEN S P, LIU F L, QIAN Y X, et al. Topic and Knowledge Association Evolution in the Field of Large Language Model-enabled Information Retrieval[J/OL]. Documentation, Information & Knowledge, (2024-06-27)[2024-10-20].
- [3] 赵鑫, 窦志成, 文继荣. 大语言模型时代下的信息检索研究发展趋势[J]. 中国科学基金, 2023, 37(5): 786-792.  
ZHAO X, DOU Z C, WEN J R. The development of information retrieval in the era of large language model[J]. Bulletin of National Natural Science Foundation of China, 2023, 37(5): 786-792.
- [4] 曹培杰, 谢阳斌, 武卉紫, 等. 教育大模型的发展现状、创新架构及应用展望[J]. 现代教育技术, 2024, 34(2): 5-12.  
CAO P J, XIE Y B, WU H Z, et al. The development status, innovation architecture and application prospects of educational big models[J]. Modern Educational Technology, 2024, 34(2): 5-12.
- [5] 苗逢春. 生成式人工智能技术原理及其教育适用性考证[J]. 现代教育技术, 2023, 33(11): 5-18.
- [6] MIAO F C. Examination of the technique principle of generative AI and its educational applicability[J]. Modern Educational Technology, 2023, 33(11): 5-18.
- [7] 余胜泉, 熊莎莎. 基于大模型增强的通用人工智能教师架构[J]. 开放教育研究, 2024, 30(1): 33-43.  
YU S Q, XIONG S S. General artificial intelligence teacher architecture based on enhanced pre-trained large models[J]. Open Education Research, 2024, 30(1): 33-43.
- [8] 齐思洋, 胡慧云, 李洪冰, 等. 融合大语言模型的领域问答系统构建方法[J]. 北京邮电大学学报, 2024: doi.org/10.13190/j.jbupt.2023-279.  
QI S Y, HU H Y, LI H B, et al. Domain question answering system construction approach integrated with large language model[J]. Journal of Beijing University of Posts and Telecommunications, 2024: doi.org/10.13190/j.jbupt.2023-279.
- [9] 文森, 钱力, 胡懋地, 等. 基于大语言模型的问答技术研究进展综述[J]. 数据分析与知识发现, 2024, 8(6): 16-29.  
WEN S, QIAN L, HU M D, et al. Review of research progress on question-answering techniques based on large language models[J]. Data Analysis and Knowledge Discovery, 2024, 8(6): 16-29.
- [10] 张金营, 王天堃, 么长英, 等. 基于大语言模型的电力知识库智能问答系统构建与评价[J/OL]. 计算机科学, 2024, https://link.cnki.net/urlid/50.1075.TP.20240528.0931.002(网络首发地址)(网络首发日期: 2024-05-28)  
ZHANG J Y, WANG T K, YAO C Y, et al. Construction and Evaluation of Intelligent Question Answering System for Electric Power Knowledge Base based on Large Language Model[J/OL]. Computer Science, 2024, https://link.cnki.net/urlid/50.1075.TP.20240528.0931.002
- [11] 竹倩叶, 鄂海红. 基于大语言模型的垂直领域问答系统研究[J]. 新一代信息技术, 2023, 6(17): 8-16.  
ZHU Q Y, E H H. Research on Vertical Domain Dialogue Systems Based on Large Language Model[J]. New Generation Of Information Technology, 2023, 6(17): 8-16.
- [12] 卢宇, 余京蕾, 陈鹏鹤, 等. 生成式人工智能的教育应用与展望: 以 ChatGPT 系统为例[J]. 中国远程教育, 2023(4): 24-31, 51.  
LU Y, YU J L, CHEN P H, et al. Educational application and prospect of generative artificial intelligence: taking ChatGPT system as an example[J]. Chinese Journal of Distance Education, 2023(4): 24-31, 51.
- [13] OpenAI. GPT-4 technical report[J]. arXiv Preprint, arXiv: 2303.08774, 2023.
- [14] TOUVRON H, MARTIN L, STONE K R, et al. Llama 2: open foundation and fine-tuned chat models[J]. arXiv Preprint, arXiv: 2307.09288, 2023.
- [15] ZENG A H, LIU X, DU Z X, et al. GLM-130B: an open bilingual pre-trained model[J]. arXiv Preprint, arXiv: 2210.02414v2, 2022.
- [16] DEVLIN J, CHANG M, LEE K, et al. BERT: pre-training of deep bidirectional transformers for language understanding[C]//Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Stroudsburg: ACL Press, 2019: 4171-4186.
- [17] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need[J]. Advances in Neural Information Processing Systems, 2017 (30): 5998-6008.
- [18] CHANG Y P, WANG X, WANG J D, et al. A survey on evaluation of large language models[J]. ACM Transactions on Intelligent Systems and Technology, 2024, 15(3): 1-45.
- [19] NEUPANE S, HOSSAIN E, KEITH J, et al. From questions to insightful answers: building an informed chatbot for university resources[J].

arXiv Preprint, arXiv: 2405.08120, 2024.

- [19] LEWIS P, PEREZ E, PIKTUS A, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks[C]//Proceedings of the 34th International Conference on Neural Information Processing Systems. Massachusetts: MIT Press, 2020: 9459-9474.
- [20] OpenAI. Our approach to alignment research[EB/OL]. (2023-07-05) [2024-10-22].
- [21] JI Z, LEE N, FRIESKE R, et al. Survey of hallucination in natural language generation[J].ACM Computing Surveys, 2023,12(55):1-38.
- [22] PETRONI F, ROCKTÄSCHEL T, LEWIS P, et al. Language models as knowledge bases? [J]. arXiv Preprint, arXiv: 1909.01066, 2019.
- [23] RAGFlow. 端到端的检索增强生成引擎[EB/OL]. (2024-04-01) [2024-10-22].  
RAGFlow. End-to-end retrieval-augmented generation engine[EB/OL]. (2024-04-01)[2024-10-22].

[作者简介]



贾春燕 (1982-), 女, 山西昔阳人, 浙江大学工程师, 主要研究方向为教育信息化、信息检索、教育大模型研究及应用等。



方伟杰 (1975-), 男, 浙江杭州人, 浙江大学高级工程师, 主要研究方向为教育信息化、教育数字化转型等。



谢宇威 (1988-), 男, 浙江杭州人, 浙江大学工程师, 主要研究方向为人工智能、网络通信、用户服务数据挖掘分析、网络服务分析、运维服务一体化等。



凌在盈 (1982-), 男, 山东临沂人, 杭州师范大学工程师, 主要研究方向为生态遥感、水色遥感、遥感工程等。