

# 基于 RAG 的大语言模型优化方法研究

冯东 张琳 梁丙杰 王再庆 陈玥如

(北京国电通网络技术有限公司 北京 100000)

**摘要** 随着自然语言处理技术的快速发展,大语言模型在各类任务中展现出了强大的性能。检索-生成架构的提出,则在减少计算资源投入的同时,进一步提高了模型的生成质量。文中提出了一种基于 RAG 的大语言模型优化方法,通过系统复杂性与效率优化、知识库质量与维护、检索与生成结果一致性等,为大语言模型的优化提供了一种有效的解决方案,具有一定的应用价值。

**关键词:** 大语言模型;RAG 模型;优化方法

**中图分类号** TP274.2

## Research on Optimization Method of Large Language Models Based on RAG

FENG Dong, ZHANG Lin, LIANG Bingjie, WANG Zaiqing and CHEN Yueru

(Beijing Guodiantong Network Technology Co., Ltd., Beijing 100000, China)

**Abstract** With the rapid development of natural language processing technology, large language models have shown strong performance in various tasks. The retrieval-generation architecture is proposed to further improve the generation quality of the model while reducing the investment of computing resources. This paper proposes a RAG-based optimization method for large language models, including system complexity and efficiency optimization, knowledge base quality and maintenance, retrieval and generation results consistency, etc., which provides an effective solution for the optimization of large language models and has certain application value.

**Key words** Big language model, RAG model, Optimization method

## 0 引言

近年来,大语言模型(Large Language Models, LLMs)在自然语言处理(Natural Language Processing, NLP)领域得到了广泛的应用,展示了其卓越的文本生成、理解和翻译能力。然而,尽管大语言模型在许多任务中表现出色,但在处理特定领域知识时仍存在一定的局限性。为弥补这些不足,基于检索-生成架构(Retrieval-Augmented Generation, RAG)的研究应运而生<sup>[1]</sup>。RAG 结合了检索模型和生成模型的优点,通过引入外部知识库和检索机制,在生成文本时动态获取相关信息,提高了生成内容的准确性和相关性<sup>[2]</sup>。但是,引入 RAG 也带来了一些挑战。本文总结了基于 RAG 的大语言模型的现存问题及产生问题的原因,通过一系列方法对模型进行优化,希望能为大语言模型的进一步优化和实际应用提供有益的参考。

## 1 RAG 的研究现状

作为一种新兴的自然语言处理方法,RAG 得到了广泛的关注。目前,对 RAG 的研究主要集中在 4 个方面。(1)模型架构的优化。研究者致力于改进检索和生成模块的集成方

式,以提高模型的效率和效果。例如,如何更高效地选择和融合文档,避免冗余信息的干扰。(2)大规模知识库的构建和应用。有效的知识库是保证 RAG 性能的基础,有研究者在构建、更新和管理大规模知识库方面进行了大量探索。(3)RAG 的实际应用取得了重要进展,特别是在开放问答领域、个性化推荐和复杂任务的多轮对话系统中,RAG 模型具有显著的优势。通过引入动态检索机制,模型能实时获取最新信息,显著提升系统的动态响应能力和准确性。(4)RAG 的可解释性研究也逐渐成为热点,通过透明的检索-生成过程,能使模型的决策过程更加可控、可解释。

总之,作为一种创新的语言模型架构,RAG 通过引入检索机制有效地扩展了生成模型的知识边界,已成为当前自然语言处理领域的重要研究方向之一。

## 2 基于 RAG 的大语言模型优化方法

### 2.1 现存问题

引入 RAG 的大语言模型在提升性能的同时,也带来了系统复杂性增加、知识库依赖、结果一致性、安全和隐私、可解释性等方面的问题。

**作者简介:**冯东(1984—),硕士,工程师,研究方向为电力信息化建设。

(1)检索和生成的集成显著增加了系统的复杂性。RAG模型在生成过程中需要实时检索外部知识库,对检索速度和生成效率提出了更高的要求。如果检索过程耗时过长,可能会导致系统整体响应速度下降,影响用户体验。

(2)RAG模型依赖于外部知识库的质量和完整性。知识库中的信息必须是最新且准确的,但构建和维护一个高质量的知识库需要投入大量的人力和物力。此外,知识库的覆盖范围和更新频率直接影响着RAG模型的表现。如果知识库中的信息不够全面或未能及时更新,就可能降低生成结果的准确性和相关性。

(3)检索和生成结果的一致性。RAG模型通过检索外部文档来优化生成过程,但检索到的信息可能会与生成模型的内部知识产生冲突。如何有效融合检索结果和生成模型的内部表示,以确保生成文本的连贯性和一致性,仍是一个亟待解决的难题。

(4)RAG模型在隐私保护和安全性方面也面临挑战。该模型依赖于外部知识库进行检索,存在潜在的信息泄露风险。如果检索过程中涉及敏感数据或私人信息,还可能引发隐私和法律问题。因此,确保检索过程的安全性和生成结果的合规性至关重要。

(5)RAG模型的可解释性和可控性问题。虽然引入检索机制在一定程度上提高了模型的透明度,但检索过程和生成结果之间的关系仍不够清晰,导致用户难以完全理解模型的决策过程,尤其是在生成内容涉及复杂推理或跨领域知识时,模型的输出往往难以预测和控制。

## 2.2 模型优化

为有效解决引入RAG后的大语言模型面临的问题,可以从以下几个方面对模型进行优化。

### 2.2.1 系统复杂性与效率优化

#### (1)异步检索与生成

为提高系统的响应速度和处理效率,设计异步处理机制至关重要。该机制能让检索和生成模块并行工作,显著提高系统响应速度。具体而言,当用户输入查询后,系统会立即启动检索过程,同时生成模块会根据初步检索结果开始生成文本。检索结束后,生成模块会对文本进行细化和优化。为实现这种异步处理方法,需要采用任务队列机制。用户的查询请求首先被放入任务队列中,由调度系统分配检索和生成任务。调度系统可利用优先级调度算法,根据任务的重要性和紧急程度动态调整任务的处理顺序,以快速响应高优先级任务。这不仅增强了RAG模型在高并发环境下的性能,也为实现更快、更精确的文本生成提供了技术保障。

#### (2)缓存机制

通过识别并缓存用户常见的查询请求及检索结果,系统可以分析历史查询数据,建立常见查询的缓存列表,提高缓存命中率。此外,对于被高频访问的数据,如特定领域的常用知识和信息,可以进行预先缓存,以降低实时检索的负

载。为确保缓存数据的时效性和准确性,可以设定缓存刷新周期,定期更新缓存内容,避免生成过时的信息。当外部知识库或数据源发生变化时,其会触发缓存更新机制,及时同步最新的数据到缓存中,提高系统的检索和生成能力。

#### (3)分层检索架构

在分层检索架构的设计中,可以构建包含快速初筛层、语义匹配层和精细检索层的多层次检索模型,以提高检索效率和准确性。在初筛层,可通过构建倒排索引,将知识库中的每个文档与关键词关联起来,以便快速查找包含特定关键词的文档,并使用快速匹配算法对用户查询进行分词,筛选出包含查询关键词的文档集合。在语义匹配层,可使用预训练的双塔模型或BERT,将用户查询和初筛文档集合转换为向量表示,并计算查询向量与文档向量之间的相似度,再基于相似度得分筛选出语义相关性较高的文档。在精细检索层,可利用多层BERT或Transformer等复杂的语义理解模型,对语义匹配层筛选出的文档进行精细化检索,并将检索文档与查询上下文进行深度融合,评估文档在特定上下文中的相关性,最终筛选出最相关的文档集合。

### 2.2.2 知识库质量与维护

#### (1)数据多样化与融合

多来源数据集成和数据清洗是提升知识库质量的重要途径。通过整合不同领域的的数据,如维基百科、专业数据库和行业报告,可以确保内容的广泛覆盖和多样性,并利用自动化数据清洗工具去除重复和冗余信息,保持数据的一致性和准确性。

#### (2)数据质量验证与过滤

数据质量验证与过滤是确保数据可靠性的关键。采用机器学习算法和规则引擎,可以自动化地验证新加入的数据,结合人工审核和机器校验,对关键数据和高价值信息进行双重审核,提升数据质量。

#### (3)知识库结构优化

优化知识库结构有助于增强数据的关联性和可检索性。通过构建知识图谱,可以以实体和关系的形式来结构化地表示数据,并根据数据的使用频率和重要性实现分层存储和索引,提高数据的访问速度和检索效率。

#### (4)知识库动态维护

通过网络爬虫和API接口实时采集最新数据,并建立自动化数据同步机制,确保知识库内容的及时更新。同时,可基于变化检测的更新策略,利用变化检测算法监测外部数据源的变化,当检测到数据变化时自动触发更新流程。根据数据变化的类型和范围,采用全量更新或增量更新策略,以提高更新效率,降低系统负担。此外,用户反馈与数据修正机制可通过建立反馈渠道,收集用户对知识库内容的意见和建议,并根据反馈结果定期修正和补充知识库内容,确保数据的完整性和准确性。

### 2.2.3 检索与生成结果一致性

#### (1)协同优化机制

在训练过程中同时优化检索模型和生成模型,通过共

享部分参数或梯度信息,使两者在相互影响下共同提高性能。设计多任务学习框架,使得检索任务和生成任务可以在共享的表示空间中进行优化,提高模型的协同工作能力。另外,通过构建端到端的训练框架,可以优化从用户查询到最终生成结果的整体流程,确保检索结果与生成内容的一致性。同时,可利用生成结果的反馈信息对检索模型进行调整,使得生成内容更符合用户需求。

### (2)上下文融合技术

上下文融合技术可以有效融合检索结果与查询上下文,实现信息增强。例如,使用加权融合策略,根据检索结果与上下文的相关性,融合文档信息与上下文信息,以生成更为连贯的文本;在生成模型的编码层和解码层分别进行多层次融合,提高文本的一致性;在生成文本前对数据进行预处理,提取关键信息并嵌入生成模型输入,使其充分利用检索结果;设计上下文感知生成策略,根据检索结果和上下文动态调整生成内容,确保生成结果的相关性和一致性。

### (3)冲突检测与解决

冲突检测主要包括语义一致性检测和事实一致性检测。其中,语义一致性检测可利用语义匹配算法分析生成内容与检索结果之间的匹配度;事实一致性检测则通过事实校验工具确保生成内容中的事实信息与检索结果一致。此外,冲突解决策略可设定优先级规则,在发生冲突时决定优先采用检索结果或生成模型的内部知识,确保内容的连贯性;也可引入自动纠正机制,根据冲突检测结果实时调整生成内容,以保证数据的一致性。

#### 2.2.4 可解释性与可控性

##### (1)可解释性优化

通过保留详细的检索记录和开发可视化工具,用户能查看生成内容的具体来源,以理解模型的决策路径。引入解释模型和自解释生成模型,可提供生成结果的局部解释和注释。采用 Attention 机制,能以可视化工具展示模型在生成特定内容时关注的输入部分,并使用 LIME(Local

Interpretable Model-agnostic Explanations)生成特定输出的解释,帮助用户理解模型的具体行为。

##### (2)可控性优化

通过交互式生成模式,用户能在生成过程中提供反馈或调整参数,实时影响生成结果,并提供控制参数接口,如关键词强化和内容排除,以使用户控制生成内容。采用模板化生成和规则约束生成技术,在一定的结构和规则下生成内容,并对生成过程施加逻辑和语义约束,可以输出具备一定格式和规范的内容,防止生成不恰当或有害内容。引入内容审查和过滤机制,在生成内容发布前进行审核,可以过滤敏感或不合适的信息,并根据实际生成情况和用户反馈动态调整生成策略和约束规则,提升生成内容的可靠性。

## 3 结语

本文提出了一系列优化策略,包括系统复杂性与效率优化、知识库质量与维护、检索与生成结果一致性、安全与隐私保护、可解释性与可控性等,旨在解决 RAG 架构在实际应用中存在的问题,提升大语言模型的整体性能和应用价值,为大语言模型的优化提供一种有效的解决方案,为大语言模型的进一步发展和应用提供有益的参考和借鉴。未来,可以进一步探索如何结合更多类型的外部知识库,提高检索机制的智能性和适应性,增强大语言模型在特定领域的表现。此外,还需要在安全性与隐私保护方面进行深入研究,确保大语言模型在各种应用场景中的可靠性。

#### 参考文献

(上接第 213 页)

- [2] 胡黎明,王东伟.新型数字化居家式养老社区解决方案[J].智能建筑,2007(11):20-21.
- [3] 孙博.“数字鸿沟”背景下山东东明县D社区智慧养老存在的问题和对策研究[D].湖北:华中师范大学,2022.
- [4] 唐敏.太原市老年人智慧社区居家养老参与意愿及影响因素研究[D].太原:山西财经大学,2023.
- [5] 张云秋,张先庚,曹冰,等.基于 Andersen 模型的独居老人养老意愿及影响因素研究[J].包头医学院学报,2023,39(11):64-69.

- [1] 周扬,蔡霏涵,董振江.大模型知识管理系统[J].中兴通讯技术,2024,30(2):63-71.
- [2] 窦凤岐,胡珊,李佳隆,等.基于 LangChain 的 RAG 问答系统设计与实现——以 C 语言课程问答系统为例[J].信息与电脑(理论版),2024,36(6):101-103.
- [6] Kim Katherine K, Rudin Robert S, Wilson Machel D. Health information technology adoption in California community health centers.[J].The American journal of managed care, 2025(12): 677-683.
- [7] 周明珠.社会工作介入智慧社区居家养老康养共融模式研究[D].哈尔滨:黑龙江大学,2023.
- [8] Rosemary Isaacs Journal[J].Pathology. Volume 2024(S1):56.
- [9] Anna Carin Karlsson, Anna Karin Edberg, Malin Sundström. Annica Backman Journal[J].Nursing ethics, 2020(8):15-33.