

基于LoRA微调的撒拉族建筑知识图谱构建

马咏梅^{1,2,3}, 耿生玲^{1,2,3*}, 赵维纳^{1,3}, 马龙龙⁴, 安波⁵

- 青海师范大学 计算机学院, 青海 西宁 810008;
- 青海省人民政府-北京师范大学 高原科学与可持续发展研究院, 青海 西宁 810008;
- 藏语智能信息处理及应用国家重点实验室, 青海 西宁 810008;
- 中国科学院 软件研究所, 北京 100190;
- 中国社会科学院 民族学与人类学研究所, 北京 100081

摘要:在弘扬中华优秀传统文化的时代背景下,青海省循化地区特有的撒拉族文化成了青海省旅游资源的重要组成部分。撒拉族建筑涵盖了撒拉族的人文历史、政治和宗教文化,但目前撒拉族建筑的构造手艺逐渐失传,对于撒拉族建筑的传承和保护迫在眉睫。知识图谱作为目前主流的知识数字化技术,将其与撒拉族建筑相结合能够更好地保护和传承撒拉族建筑。因此该文构建了独特的撒拉族建筑数据集以及本体结构,使用大语言模型LoRA微调方法以撒拉族建筑数据为对象进行实体和关系抽取,并与经典知识抽取模型进行对比实验,相比于ChatGPT,实体、关系抽取的F1值分别提高了13.37%、16.32%。该研究所构建的知识图谱可以用于撒拉族建筑知识的知识推荐、智能搜索以及基于检索增强生成(Retrieval Augmented Generation, RAG)技术的知识图谱增强问答系统。

关键词:大语言模型;LoRA微调;知识抽取

中图分类号:O436 文献标志码:A 文章编号:0253-2395(2025)05-0859-12

Construction of Knowledge Graph for Salar Architecture Based on LoRA Fine-tuning

MA Yongmei^{1,2,3}, GENG Shengling^{1,2,3*}, ZHAO Weina^{1,3}, MA Longlong⁴, AN Bo⁵

- School of Computer Science, Qinghai Normal University, Xining 810008, China;
- Academy of Plateau Science and Sustainability, People's Government of Qinghai Province & Beijing Normal University, Xining 810008, China;
- The State Key Laboratory of Tibetan Intelligent Information Processing and Application, Xining 810008, China;
- Institute of Software, Chinese Academy of Sciences, Beijing 100190, China;
- Institute of Ethnology and Anthropology, Chinese Academy of Social Sciences, Beijing 100081, China)

Abstract: In the context of promoting the excellent traditional Chinese culture, the unique Salar culture in the Methodist region of Qinghai Province has become an important part of the tourism resources in Qinghai Province. Salar architecture covers the Salar human history, politics and religious culture, but at present the construction technology of Salar architecture is gradually lost, and the inheritance and protection of Salar architecture is urgent. As a mainstream knowledge digitization technology, the combination of knowledge graph and Salar architecture can better protect and inherit the Salar architecture. Therefore, this paper constructs a unique

收稿日期:2024-11-12;接受日期:2025-03-08

基金项目:国家自然科学基金(61862055;62266036);国家社会科学基金(22BTQ010);科技援青合作专项-区域处(2022-QY-203);中国社会科学院专项(2024SJK017)

作者简介:马咏梅(2000-),女,青海西宁人,硕士研究生,研究方向为自然语言处理与知识图谱。E-mail:3421877930@qq.com

*通信作者:耿生玲(GENG Shengling),E-mail:geng_sl@126.com

引文格式:马咏梅,耿生玲,赵维纳,等.基于LoRA微调的撒拉族建筑知识图谱构建[J].山西大学学报(自然科学版),2025,48(5):859-870. DOI:10.13451/j.sxu.ns.2025012.

Salar architecture dataset and ontology structure, uses the large language model LoRA fine-tuning method to extract entities and relations from Salar architecture data, and conducts comparative experiments with the classical knowledge extraction model. Compared with ChatGPT, the $F1$ values of entity extraction and relation extraction are improved by 13.37% and 16.32%. The knowledge graph constructed by the institute can be used for knowledge recommendation, intelligent search and knowledge graph enhanced question answering system based on RAG technology.

Key words: large language model; LoRA fine-tuning; knowledge extraction

0 引言

青海省地处黄河中原文化、佛教文化、伊斯兰教文化的交融地带,循化撒拉族自治县是中国撒拉族的发祥地,撒拉族于八百年前自中亚远渡中国来到青海省循化撒拉族自治县扎根生活。几百年来,与当地的汉族、回族、藏族等民族进行了深刻的融合交流,形成了独特的多元民族文化^[1]。循化撒拉族自治县境内建筑文化旅游资源丰富,有古清真寺、古拱北、古篱笆楼、古传统民居遗存 58 处^[2],其中篱笆楼被列为国家非物质文化遗产,篱笆楼营造技艺是大美青海多彩循化建筑文化的重要组成部分^[3]。2013 年 3 月 5 日,国务院公布了第七批全国重点文物保护单位,其中白庄镇科哇清真寺、张尕清真寺、清水乡河东清真寺、孟达大庄清真寺等四座寺院作为撒拉族清真寺古建筑群跻身“国保”行列。循化撒拉族自治县撒拉族建筑的建造年代历史悠久,木雕图案醒目大方,样式精美别致,具有很高的历史意义和保存价值。但是因为朝代的更迭、战争等历史原因,撒拉族建筑的构造流程和手艺逐渐失传,对于撒拉族建筑的传承和保护也是政府近几年比较关注的问题。

目前,知识图谱在撒拉族建筑的研究中并没有得到广泛的应用,而现存的主流大模型对于撒拉族建筑此类特定领域的的数据,在实体、关系抽取方面存在数据需求量大、专业术语和复杂概念难以准确理解和抽取等问题。因此,本文以现有的循化地区撒拉族建筑资料入手,构建撒拉族建筑知识图谱,提出一种基于大模型的低秩适配(Low-Rank Adaptation, LoRA)微调知识抽取方法,进行实体和关系的抽取,实验结果表明,ChatGLM^[4]3-6B LoRA 微调的方法,比 ChatGPT^[5]等模型在精确率 P 、召回率 R 和 $F1$ 方面有了进一步提升。

1 相关研究

1.1 知识图谱

随着非物质文化遗产数字化保护的话题越来越受大众关注,知识图谱作为一门知识的可视化技术也引起了很多人的兴趣,知识图谱是 2012 年 5 月由 Google 公司提出的一种用于表示和组织知识的图形化数据结构,本质上是一种语义网络,网络中的节点代表实体或者概念,边代表实体或概念之间的各种语义关系^[6]。知识图谱旨在捕捉和表达现实世界中的知识,使计算机能够理解和推理这些知识的信息。

知识图谱被应用在众多领域。陈明等^[7]基于多重注意力机制的命名实体识别模型解析医疗问句中的实体,采用基于 BERT-BiLSTM 的关系抽取模型进行关系抽取, BERT-BiLSTM 是结合基于 Transformer 的双向编码器表示模型(Bidirectional Encoder Representations from Transformers, BERT)^[8]和双向长短期记忆网络(Bidirectional Long Short Term Memory Network, BiLSTM)两种技术的模型,生成查询语句从医疗知识图谱中获取问题的答案。黄伟春等^[9]构建军事术语本体,采用基于规则的方法和基于训练模型的方法抽取军事术语关系,构建高质量的军事术语知识图谱。韩子威等^[10]使用自顶向下构建的知识图谱对动车故障信息进行知识管理,提出动车组故障模型,提升动车组故障信息的智能化处理。杨硕等^[11]基于 CompIEEx (Complex Embedding) 图谱嵌入的方法,引入上下文注意机制和卷积神经网络(Convolutional Neural Network, CNN)网络获取多模态知识图谱的特征表示,将多模态知识图谱提供的背景知识与问答对的文本语义信息融合。

知识图谱被广泛使用于医疗、金融、法律等行业,对于多模态知识图谱、迁移学习和跨语言知识图谱的研究也处于最热门的话题,因此

知识图谱与撒拉族建筑的结合为撒拉族建筑文化提供了高效且新颖的保护模式。

1.2 建筑知识图谱

撒拉族人民因为社会历史、居住生活和宗教活动的需求,在明清时代的循化撒拉族聚集村落大兴土木,修建壮观的中国殿堂形制清真寺、楼阁式拱北建筑和秀丽的河湟传统民居家院、篱笆楼建筑等,融入了中华民族建筑文化艺术林。但由于社会历史等多种原因,撒拉族的精华建筑遭到了不同程度的肆意破坏。使用知识图谱方式对于建筑知识的传承与保护,已有很多人提出新观念,张琳等^[12]使用绿色建筑知识图谱进行分析,为绿色建筑的发展提供了理论参考。Chen等^[13]以Web of Science数据库的数据为样本,利用软件绘制建筑领域和大数据交叉领域的知识图谱,为建筑信息模型(Building Information Modeling, BIM)、建筑节能、智慧城市、防灾防损等研究热点进行定性分析。Cao等^[14]为了解决城市历史建筑(UHBs)知识难以被有效访问和管理的问题,提出一个从非结构化文本中提取和管理UHBs知识的总体框架,构建了城市居民住房的知识库以及北京居民住房知识图谱,并实现了其知识检索和可视化。Xie等^[15]将数字孪生(Digital Twin, DT)技术与知识图谱相结合,使用知识图谱可视化和分析相关科学领域文献,对于建筑行业数字孪生的当前焦点和未来方向进行定量分析。

1.3 ChatGLM

自从ChatGPT面世以来,以大模型为中心的知识抽取形式更为多样,“一问一答”交互形式的问答体系变得备受关注,其中对话生成语言模型(Chat Generative Language Model, ChatGLM)是使用中英双语预训练进行自然语言理解和生成的通用预训练框架,它有效地将来自中文问答反馈的自然语言响应与人类监督学习、反馈强化策略结合起来。以ChatGLM为基础,衍生而出ChatGLM-6B、ChatGLM2-6B以及ChatGLM3-6B模型,其中ChatGLM3-6B模型针对中文问答和对话进行优化,经过中英双语训练,辅以监督微调、反馈自助、人类反馈强化学习等技术的加持,对于信息抽取等自然语言处理任务有着优异的表现。尹娴等^[16]针对

水生动物疾病诊断智能对话系统存在的复杂的专业性知识和准确性低的问题,提出了一种基于ChatGLM模型的改进水生动物疾病诊断相关问题的优化方法,该方法通过在ChatGLM模型的中间层插入Adapter块针对相关的专业问题进行微调,提高了模型的专业性和准确性。Zeng等^[17]对一个具有1300亿个参数的双语(英语和汉语)预训练语言模型GLM-130B进行预训练,在相关基准测试中,GLM-130B始终显著优于最大的中文模型ERNIE TITAN 3.0 260b。Tao等^[18]提出了RoleCraft-GLM解决了对话AI中缺乏个性化交互的关键问题,并提供了一个具有详细和情感微妙的角色描绘的解决方案。

2 撒拉族建筑知识图谱构建

撒拉族建筑包括篱笆楼、松木大房以及宗教建筑拱北、清真寺,建筑样式多样、时间跨度大,实现数据的有效整合对于了解撒拉族社会结构和生活方式有着深远的意义,撒拉族知识图谱的构建流程如图1所示。

2.1 撒拉族建筑数据集构建

数据集构建是知识图谱构建的基础步骤和关键环节,高质量的语义数据能够提高知识抽取、知识融合的精确率和置信度。撒拉族建筑原生数据质量低,网上可查的电子资源有限且多为图书馆无法外借的纸质版数据,所以构建撒拉族建筑数据库存在一定的难度,因此使用Scrapy爬虫方法从《中国非物质文化遗产网》、百度百科词条等网站爬取数据,对于图书馆馆藏书籍例如《撒拉族》^[1]、《撒拉族古建筑》^[2]、《撒拉族篱笆楼》^[3]进行馆藏扫描,对于《浅谈撒拉族的古民居建筑》^[19]、《撒拉族松木大房子民居木构架营造技艺研究》^[20]、《青海撒拉族传统民居门窗研究》^[21]等学术论文进行OCR文字转换识别,对撒拉族建筑形成的历史、人文因素作为切入点进行建筑特征以及建筑建构等方面数据的筛选和处理,整理为TXT文本资料进行存储,分析撒拉族民居、宗教建筑之间的关系。撒拉族建筑数据集的构建框架如图2所示。

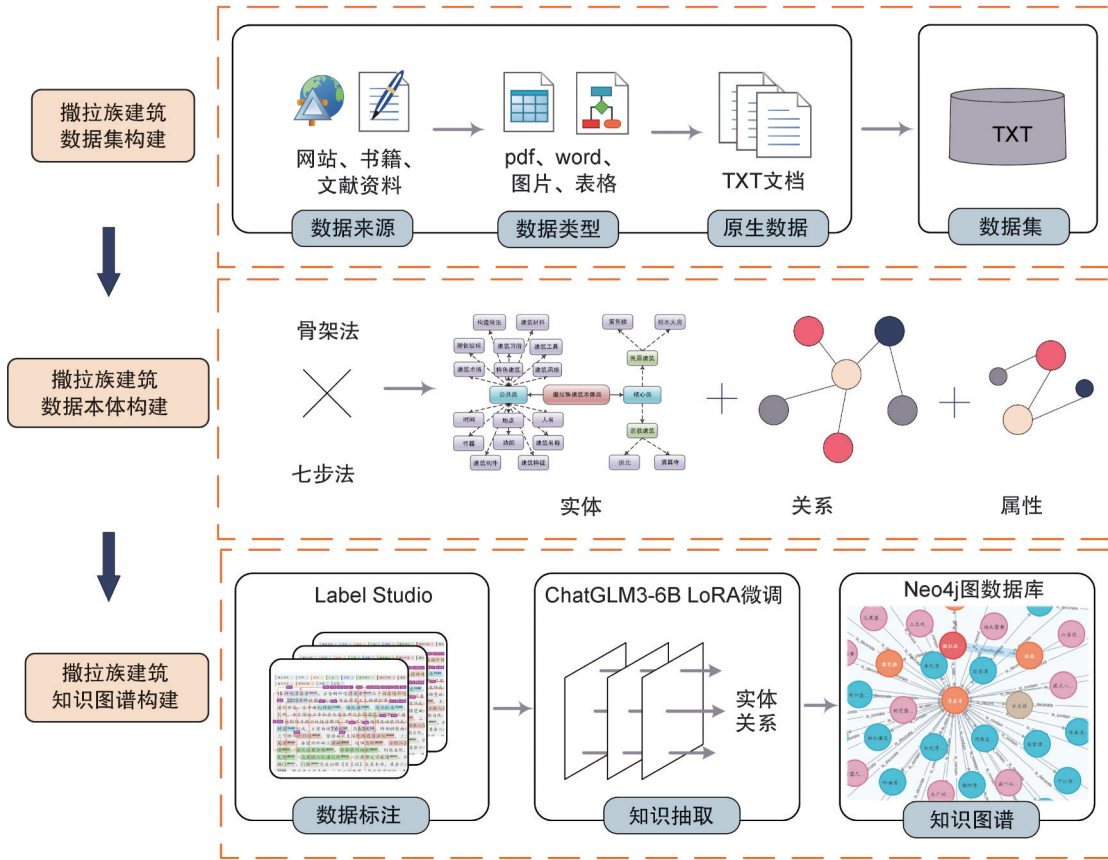


图1 撒拉族建筑知识图谱构建流程

Fig. 1 Construction process of Salar architectural knowledge graph

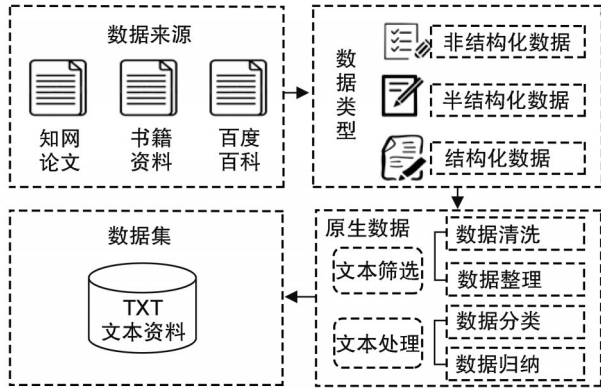


图2 撒拉族建筑数据集构建框架

Fig. 2 Framework for building salar architecture dataset

2.2 撒拉族建筑知识图谱本体构建

2.2.1 本体构建工具

构建知识图谱的本体工具多样,其中Protégé^[22]作为主流的本体构建工具是一款由斯坦福大学编写并维护的开源本体建模和编辑工具,使用网络本体语言(Web Ontology Language, OWL)对知识进行表示,其功能包括类建模、模型处理以及模型交换等,可以通过包

含、限制等方式,重用模块的结构库中的已有本体。

2.2.2 本体构建方法

常见的本体构建^[23]法有多伦多虚拟企业本体构建方法(TORonto Virtual Enterprise Methodology, TOVE)、METHONTOLOGY本体构建方法(METHONTOLOGY Methodology, METHONTOLOGY)、骨架法、七步法^[24]和基于叙词表的领域本体构建法,根据撒拉族建筑的特点和研究需求,基于骨架法和七步法的特点,结合提供开发本体基本框架的骨架法和构建本体需要明确本体应用范围的七步法,作为撒拉族建筑领域本体的构建方法,利用骨架法和七步法相结合,构建撒拉族建筑本体的步骤如下:

(1)明确本体应用的范围和目的:构建撒拉族建筑的本体领域意在撒拉族建筑数据进行结构化梳理。

(2)确定核心概念:分析确定撒拉族建筑领域核心类概念以及次级类概念,确定最底层的

本体。

(3)定义概念属性:为每个核心类、次级类概念定义撒拉族建筑属性,包括属性名称、属性类型、属性值等。

(4)定义概念关系:为每个核心类、次级类概念定义撒拉族建筑关系,包括关系名称、关系类型、关系值等。

(5)构建骨架:根据核心概念、属性和关系,构建撒拉族建筑本体的骨架,包括类、子类、属性关系等。

(6)本体评估:对于构建的撒拉族建筑本体,进行人工校验,确保满足要求。

(7)本体的建立:使用Protégé编辑工具建立本体,并保存本体。

2.2.3 本体构建过程

在确定其本体构建方法的基础之上,分析定义撒拉族的建筑核心概念,分为民居建筑和宗教建筑两大核心类,民居建筑分为篱笆楼建筑、松木大房建筑,宗教建筑又可分为拱北建筑和清真寺建筑。

撒拉族文化历史中,篱笆楼为明清时期撒拉族人民的生活场所,松木大房多为近现代撒拉族人居住的建筑,而拱北是撒拉族人民以圣贤者陵墓为中心建造的宗教性建筑,用于祭奠先贤。撒拉族信仰的宗教为伊斯兰教,所以清真寺不但承担着广大信教群众进行宗教活动的特殊职能,更具有开展政治、文化、教育等一切

社会活动的综合职能。

在两个核心类、四个次级类实体的基础上,撒拉族建筑本体再细分为建筑名称、时间、地点、人名、功能、建筑特征、建筑风格、建筑术语、建筑材料、建筑工具、特色建筑、构造做法、建筑构件、雕饰纹样、书籍、建筑习俗等16个实体公共类,如图3所示。

使用Protégé进行实体可视化展示,如图4所示。

表1选取部分核心类、公共类的实例,进行详细说明,篱笆楼、松木大房等建筑都具有各自的建筑时间、风格、材料等。

撒拉族建筑知识图谱类的层次结构只构成了本体的基本框架,必须借助属性来描述既定的知识,因此在本体构建的判定、分析过程中定义了面阔、壁高、进深、开间、间阔、坡度、占地面积和建筑面积等15个属性。表2对于部分属性的含义详细进行阐述。

表3在定义其类与属性的基础上进行了关系的定义,让实体与实体、实体与属性进行关联形成经典的三元组形式,多个三元组构成知识图谱。撒拉族建筑知识图谱经过数据的整理,总共抽取、定义了“记载”“建造”“位于”等10种关系。撒拉族建筑知识图谱中存在的三元组为(实体-关系-实体)、(实体-属性-实体)两类,例如(“红光清真寺”,is_locate,“查汗都斯乡赞卜乎村”)。

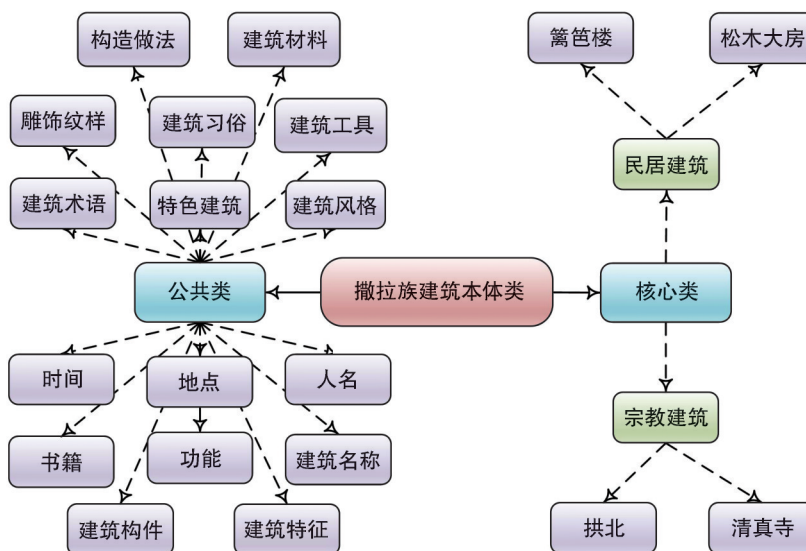


图3 撒拉族建筑领域本体结构层次

Fig. 3 The ontology structure hierarchy of Salar architecture field

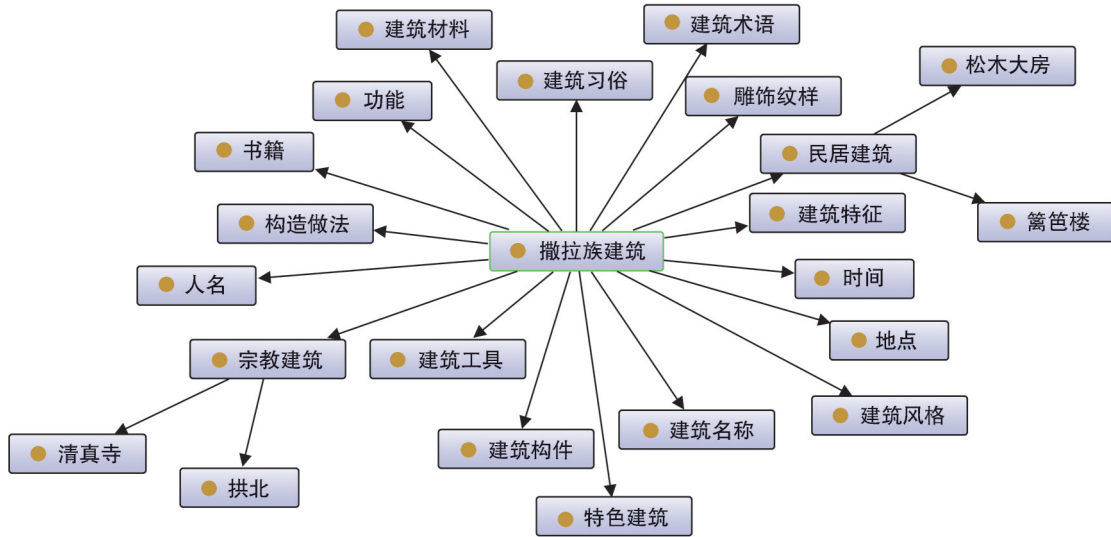


图4 核心类与公共类的可视化

Fig. 4 Visualization of core and public classes

表1 公共类实体实例表(部分)

Table 1 Public entity instance table (partial)

建筑名称	时间	特色建筑	建筑风格	雕饰纹样	建筑构建	建筑材料
篱笆楼	明清时期	篱笆墙	四合院	大丽花	旋转棋牌门	篱笆
松木大房	清代中期	次间	三合院	花斑纹	檐枋	松木
街子拱北	明代清初	拱北院	一出檐	石榴纹	角科	云杉
奄古录拱北	元明年代	静修室	悬山式	阿文清真言	拱北亭	松木
白庄拱北	大清道光年间	拱北厅座	阁楼式	两狮滚绣球	宝瓶花脊	花岗石
孟达清真寺	明朝中期	礼拜殿	拐角式	鹿鹤蝙蝠	昂斗	落叶松
张尕清真寺	清朝年间	唤礼楼	歇山式	牡丹秋菊	支摘窗	河光石

表2 撒拉族建筑属性表

Table 2 Salar architecture attribute table

Property	属性含义
Width of space	间阔:房间的宽度
Wall_height	壁高:建筑墙壁高度
A_surname	开间:相邻两个横向定位墙体间的距离
Breadth	面阔:度量建筑物平面宽度的单位
Depth	进深:建筑物的长度
Slope	坡度:指建筑物的屋面、地面或楼梯等各个水平承重面的倾斜程度

2.3 基于 ChatGLM3-6B 微调的撒拉族建筑文本数据知识抽取

2.3.1 数据标注处理

为了更好地进行实体和关系抽取任务,在构建的撒拉族建筑本体基础之上需要进行数据标注。图5中使用 Label Studio 标注平台将撒拉族建筑数据中的实体、属性和关系分别进行标注。在数据标注前需要将文本数据进行分类、整理,保证其格式符合平台需求,其后在标注

平台将预先定义的16个实体名、15个属性名以及10种关系,依次在文本数据中标注实体类对应的实体名、实体属性以及实体与实体之间的关系。

表3 撒拉族建筑知识图谱关系表

Table 3 Relationship table of salar architecture knowledge graph

relation	关系含义	关系实例
is_locate	位于	(“红光清真寺”,is_locate,“积石镇”)
is_style	风格	(“科哇清真寺”,is_style,“四合院”)
is_contain	包含	(“阿伊清真寺”,is_contain,“随梁”)
is_record	记载	(“《西海都市报》”,is_record,“篱笆楼”)
is_materials	材料	(“阿河滩清真寺”,is_materials,“落叶松”)
is_decorate	装饰	(“仙桃石榴”,is_decorate,“门前照壁”)

2.3.2 基于 LoRA 微调方法的知识抽取

ChatGLM3-6B的底层架构 ChatGLM 的基础原理如图6所示。BERT使用原理与完形填空类似,除去其中一词使用上下文联合预测,Chat-

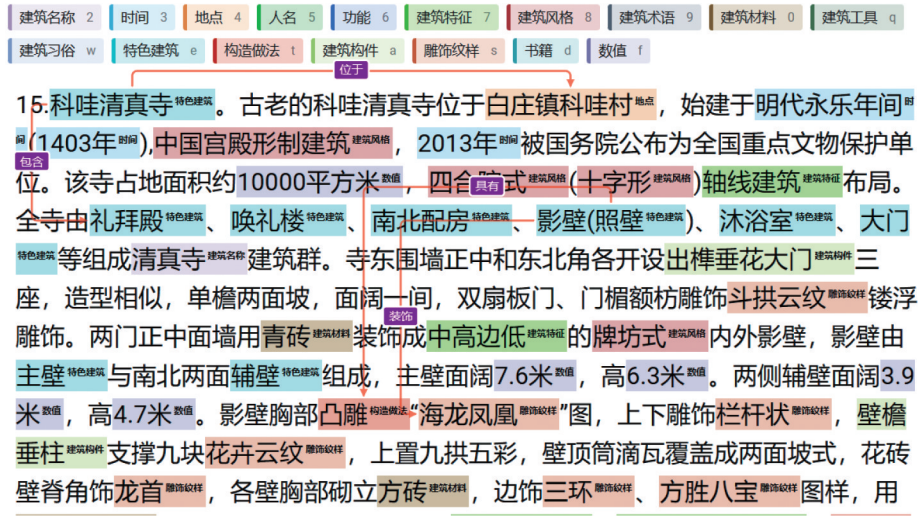


图5 Label Studio平台标注数据
Fig. 5 Annotates data of Label Studio platform

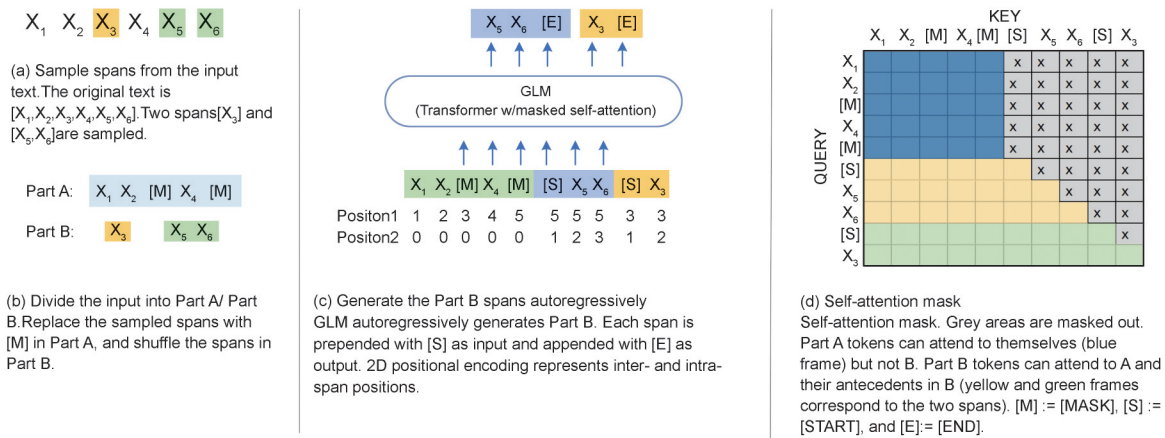


图6 ChatGLM基础架构
Fig. 6 ChatGLM infrastructure

GPT是根据上文预测下文，而 ChatGLM将 ChatGPT 优异的生成能力和BERT出色的双向注意力进行结合，在被隐藏的地方使用单项注意力，未被掩藏的地方使用双向注意力，是一种自回归填空。由于大模型训练所需的数据量庞大，如果对模型整体进行全量微调，无法满足需要的算力与数据量，ChatGLM3-6B目前主推的微调方式有 LoRA 和 P-Tuning V2。P-Tuning V2 微调方法引入了 Prefix-Tuning 的思想，在模型的每一层加入 Prefix，并采用多任务学习，但是所需 GPU 高达上百 GB，训练成本高，所以基于撒拉族建筑文本数据的独特性，采用 LoRA 方法进行微调实验。

ChatGLM3-6B 使用的激活函数可以近似实现为：

$$GELU(x) \approx$$

$$0.5x(1 + \tanh(\sqrt{\frac{2}{\pi}}(x + 0.044715x^3))). \quad (1)$$

LoRA 的基本原理是冻结预训练好的模型权重参数，然后在冻结的模型基础上加入额外的网络层。这些新增的网络层只包含少量的参数，因此训练这些参数的成本较低。同时，由于 LoRA 采用低秩分解来模拟参数更新量，可以在不引入额外推理延迟的情况下实现与全模型微调类似的效果，其本质是在 PLM (Pre-trained Language Model) 旁增加一个旁路，意在进行一个先降低维度再升高维度的操作来模拟内训练时固定 PLM，只训练旁路，推理时旁路先相乘合并再和 PLM 权重相加从而更新权重。LoRA 存储矩阵的更新方式如图 7(a) 与图 7

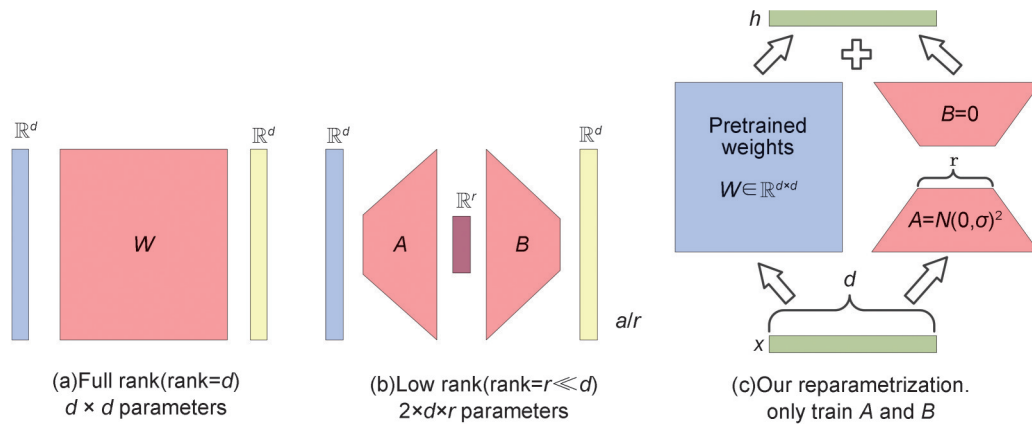


图7 ChatGLM3-6B 模型 LoRA 微调核心流程示意图

Fig. 7 Schematic diagram of ChatGLM3-6B model LoRA fine-tuning core workflow

(b)所示,完全微调需要 $d \times d$, 而 LoRA 微调只需要 $2 \times d \times r$ 个参数。ChatGLM3-6B LoRA 微调模型如图 7(c)所示,左边是预训练模型的权重,在训练期间被冻结,不接受梯度更新。右边部分对 A 使用随机的高斯初始化,B 在训练开始时为零, r 是秩,会对 ΔW_x 做缩放 d/r 。

2.3.3 评价指标

通过微调可以提高模型的精度、减少训练时间以及增加模型的泛化能力,将这种泛化能力应用到特定任务中,从而使模型能够更好地适应其他类型的数据。采用精确率(Precision, P)、召回率(Recall, R)和 $F1$ 值作为对于模型是否有效进行信息提取的评价指标,指标如下:

精确率:表示模型做出正确预测的数量占总体预测为正确数量的百分比。计算公式为:

$$P = \frac{TP}{TP + FP}, \quad (2)$$

其中 TP (True Positives) 表示模型正确预测的正例数量, FP (False Positives) 表示模型错误地将负例预测为正例的数量。

召回率:表示模型做出正确预测的数量占实际预测正确数量的百分比。计算公式为:

$$R = \frac{TP}{TP + FN}, \quad (3)$$

其中 TP (True Positives) 表示模型正确预测的正例数量, FN (False Negatives) 表示模型错误地将正例预测为负例的数量。

$F1$: 表示精确率和召回率的调和平均数,用于综合评价模型的性能,综合考虑了模型的精度和覆盖率。计算公式为:

$$F1 = \frac{2 \times P \times R}{P + R}, \quad (4)$$

其中 $F1$ 分数的取值范围在 0 到 1 之间,越接近 1 表示性能越好。

3 实验及可视化展示

3.1 对比实验

知识图谱分为通用领域和垂直领域知识图谱,垂直领域知识图谱知识质量要求更高,知识的应用形式也更加广泛且具有高度时效性和更深层次的推理需求等特点。撒拉族建筑知识图谱属于垂直领域的知识图谱,所以选择使用 ChatGPT、国内大语言模型文心一言和通用信息抽取模型(Universal Information Extraction, UIE)从通用领域和垂直领域的角度出发作对比试验,以及使用 ChatGLM3-6B LoRA 微调模型进行多次迭代训练得出的精确率、召回率和 $F1$,如表 4 所示。

3.2 错误分析

为了深入分析模型在撒拉族建筑领域的表

表4 对比实验结果

Table 4 Comparative experimental results

模型	任务	Precision	Recall	F1 ↑ better
ChatGPT	实体抽取	75.00%	69.23%	72.00%
	关系抽取	72.50%	65.91%	69.05%
文心一言	实体抽取	58.54%	63.16%	60.76%
	关系抽取	58.14%	54.35%	56.18%
UIE	实体抽取	71.79%	68.29%	70.00%
	关系抽取	85.19%	65.71%	74.19%
本文模型	实体抽取	83.33%	81.08%	82.19%
	关系抽取	87.50%	83.33%	85.37%

现,从撒拉族建筑数据中选择一段具有代表性的语料“全寺由礼拜殿……”为例,使用UIE模型只能抽取其中的部分实体和关系,无法实现完整的知识抽取。ChatGPT在抽取实体时,将整句输出,统一归类为建筑特征,未能抽取出“垂柱”为建筑构件、“两面坡式”为建筑风格。图8为未使用LoRA微调方式之前,ChatGLM3-6B模型对于撒拉族建筑特有的建筑实体如“礼拜殿”错误的抽取为“建筑名称”实体类,并且对于构造做法、建筑构件以及建筑材料等实体无法进行抽取。

图9为使用LoRA微调之后的结果,对于给定语料中的特色建筑、建筑特征和建筑构件等实体类型进行了正确的抽取。

从实验结果可以看出,对于撒拉族建筑知识的实体抽取、关系抽取方面,ChatGLM3-6B LoRA微调方法相比与UIE模型、ChatGPT、文心一言有更好的抽取性能。相比于ChatGPT在实体抽取方面,精确率、召回率、F1值提升了8.33%、11.85%、13.37%,在关系抽取的精确率、

召回率、F1值方面分别提升了15.00%、17.42%、16.32%;相比于UIE模型在实体抽取方面,F1值提升了12.19%,关系抽取方面F1值提升了11.17%。

3.3 消融实验

在本文的模型中,通过ChatGLM结构特征和ChatGLM3-6B生成式性能特征,使用LoRA模块进行实体、关系的微调训练。为了进一步研究所加入微调训练的有效性,设计消融实验:

本文模型为基线模型,以此为对照实验建立三个模型版本:A模型为ChatGLM-6B;B模型为ChatGLM3-6B;C模型为基准模型ChatGLM3-6B LoRA。

实验数据集包含了撒拉族建筑实体、关系的标注信息,设置数据集时考虑到数据集的代表性、平衡性和分布一致性,采用分层抽样方法将数据集划分为70%训练集、15%验证集和15%测试集。

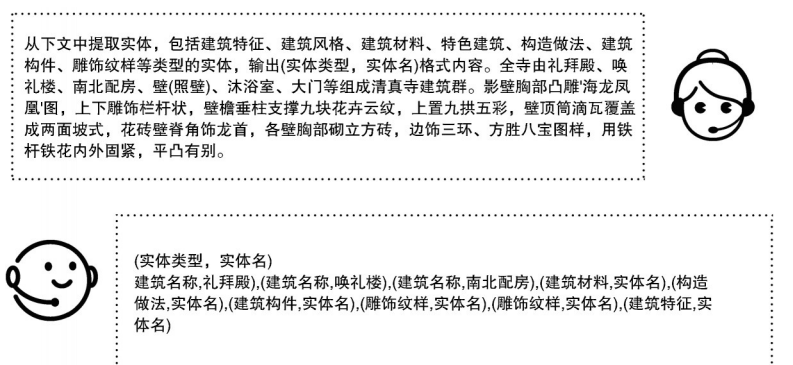


图8 LoRA微调前模型在对话场景下的展示图

Fig. 8 Display diagram of the model in conversation scenarios before LoRA fine-tuning

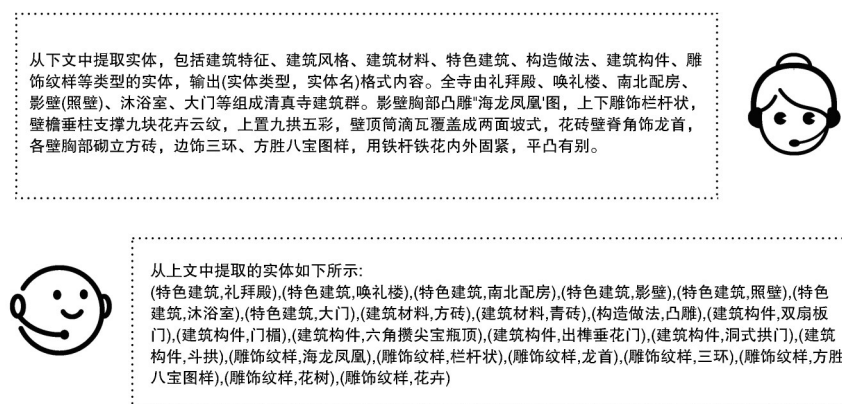


图9 LoRA微调后模型在对话场景下的展示图

Fig. 9 Display diagram of the model in conversation scenarios after LoRA fine-tuning

索和阅读理解两个阶段结合起来,以实现更准确和高效的问答的技术,基于RAG方式的问答系统可以支持用户与知识图谱进行交互式的探索和研究,为用户提供撒拉族建筑的相关知识。

5 结论

本文立足撒拉族建筑文化的传承与保护,通过多渠道获得的撒拉族建筑数据,建立撒拉族建筑数据集,基于本体构建方法、Protégé工具从原生数据质量低、无法被高效利用的数据集中构建出撒拉族建筑数据中独有的民居建筑和宗教建筑的本体结构,使用Label Studio平台将数据集中的实体、关系和属性按要求进行标注。由于撒拉族建筑涉及建筑艺术、民族文化等诸多领域内知识,所以在知识抽取方面存在一定难度,经过多个模型的实验结果验证使用ChatGLM3-6B LoRA微调方法对于撒拉族建筑知识领域的实体、关系联合抽取任务有着更优异的综合表现,最后利用Neo4j图数据库将抽取的实体、关系进行撒拉族建筑知识图谱的构建,实现了撒拉族建筑知识与知识图谱的结合,促进了非遗文化的传播和撒拉族建筑数据的保护。在撒拉族建筑知识图谱构建过程中发现撒拉族建筑知识融合、推理以及撒拉族建筑知识图谱的应用等方面的问题还有待完善与深入研究。因此后期的工作将着重关注以上问题,并且进行数据源的持续监控,进行本体演化、技术迭代等方式不断维护和更新撒拉族建筑知识图谱以及实现基于RAG方式的撒拉族建筑知识问答系统。

参考文献:

- [1] 禹规娥.撒拉族[M].乌鲁木齐:新疆美术摄影出版社;乌鲁木齐:新疆电子音像出版社,2010.
YU G E. Salar Tribe[M]. Urumqi: Xinjiang Fine Arts Photography Publishing House; Urumqi: Xinjiang Electronic Audio and Video Publishing House, 2010.
- [2] 马进明,马晓红.撒拉族古建筑[M].西宁:青海民族出版社,2014.
MA J M, MA X H. Salar Ancient buildings[M]. Xining: Qinghai Ethnic Publishing, 2014.
- [3] 马进明,马晓红.撒拉族篱笆楼[M].西宁:青海民族出版社,2013.
MA J M, MA X H. Salar Tribe Fence Building[M]. Xining: Qinghai Ethnic Publishing House, 2013.
- [4] DU Z X, QIAN Y J, LIU X, et al. GLM: General Language Model Pretraining with Autoregressive Blank Infilling[EB/OL]. (2021-03-18) [2024-09-07]. <https://arxiv.org/abs/2103.10360>.
- [5] ROUMELIOTIS K I, TSELIKAS N D. ChatGPT and Open-AI Models: a Preliminary Review[J]. *Future Internet*, 2023, **15**(6): 192. DOI:10.3390/fi15060192.
- [6] 王萌,王昊奋,李博涵,等.新一代知识图谱关键技术综述[J].计算机研究与发展,2022,**59**(9): 1947-1965. DOI: 10.7544/issn1000-1239.20210829.
WANG M, WANG H F, LI B H, et al. Survey on Key Technologies of New Generation Knowledge Graph[J]. *J Comput Res Dev*, 2022, **59**(9): 1947-1965. DOI: 10.7544/issn1000-1239.20210829.
- [7] 陈明,刘蓉,熊回香.基于医疗知识图谱的智能问答系统研究[J].情报科学,2023,**41**(12): 118-126. DOI: 10.13833/j.issn.1007-7634.2023.12.015.
CHEN M, LIU R, XIONG H X. Intelligent Question Answering System Based on Medical Knowledge Graph [J]. *Inf Sci*, 2023, **41**(12): 118-126. DOI: 10.13833/j.issn.1007-7634.2023.12.015.
- [8] DEVLIN J, CHANG M W, LEE K, et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding[EB/OL]. (2018-10-11) [2024-09-07]. <https://arxiv.org/abs/1810.04805>.
- [9] 黄伟春,肖刚,杨健,等.基于本体的军事术语知识图谱构建方法[J].指挥控制与仿真,2023,**45**(5): 10-17. DOI: 10.3969/j.issn.1673-3819.2023.05.002.
HUANG W C, XIAO G, YANG J, et al. Ontology-based Military Terminology Knowledge Graph Construction Method[J]. *Command Contr Simul*, 2023, **45**(5): 10-17. DOI: 10.3969/j.issn.1673-3819.2023.05.002.
- [10] 韩子威,朱建生.动车组故障知识图谱构建方法研究[J].铁道机车车辆,2023,**43**(4): 17-22. DOI: 10.3969/j.issn.1008-7842.2023.04.03.
HAN Z W, ZHU J S. Research on Construction Method of EMU Failure Domain Knowledge Map[J]. *Railw Locomot Car*, 2023, **43**(4): 17-22. DOI: 10.3969/j.issn.1008-7842.2023.04.03.
- [11] 杨硕,李书琴.多模态知识图谱增强葡萄种植问答对的答案选择模型[J].农业工程学报,2023,**39**(14): 207-214. DOI: 10.11975/j.issn.1002-6819.202304240.
YANG S, LI S Q. Enhancing Answer Selection Model of Grape Planting Using Multimodal Knowledge Graph [J]. *Trans Chin Soc Agric Eng*, 2023, **39**(14): 207-214. DOI: 10.11975/j.issn.1002-6819.202304240.
- [12] 张琳,陈立文,曹江红,等.基于CiteSpace软件的绿色

- 建筑知识图谱分析[J]. 山东建筑大学学报, 2018, **33**(3): 26-31. DOI: 10.12077/sdjz.2018.03.005.
- ZHANG L, CHEN L W, CAO J H, *et al.* Visualizing Analysis of Green Building Knowledge Map Based on Citespace Software[J]. *J Shandong Jianzhu Univ*, 2018, **33**(3): 26-31. DOI: 10.12077/sdjz.2018.03.005.
- [13] CHEN G X, HOU J, LIU C S, *et al.* Visualization Analysis of Cross Research between Big Data and Construction Industry Based on Knowledge Graph[J]. *Buildings*, 2022, **12**(11): 1812. DOI:10.3390/buildings12111812.
- [14] CAO X, GUO X, JIANG J. Knowledge Graph Enabled Representation and Exploration for Urban Historical Buildings: A Case Study in Beijing, China[J]. *Int Arch Photogramm Remote Sens Spatial Inf Sci*, 2022, XLVIII-3/W2-2022: 1-7. DOI: 10.5194/isprs-archives-xlvi-3-w2-2022-1-2022.
- [15] XIE H Y, XIN M Y, LU C W, *et al.* Knowledge Map and Forecast of Digital Twin in the Construction Industry: State-of-the-art Review Using Scientometric Analysis[J]. *J Clean Prod*, 2023, **383**: 135231. DOI:10.1016/j.jclepro.2022.135231.
- [16] 尹娴, 冯艳红, 叶仕根. 基于ChatGLM的水生动物疾病诊断智能对话系统的优化研究[J]. 现代电子技术, 2024, **47**(14): 177-181. DOI: 10.16652/j. issn. 1004-373x.2024.14.027.
- YIN X, FENG Y H, YE S G. Optimization of a ChatGLM-based Intelligent Dialogue System for Aquatic Animal Disease Diagnosis[J]. *Mod Electron Tech*, 2024, **47**(14): 177-181. DOI: 10.16652/j. issn.1004-373x.2024.14.027.
- [17] ZENG A, LIU X, DU Z, *et al.* GLM-130B: An Open Bilingual Pre-trained Model[EB/OL]. (2022-10-05) [2024-09-07]. <https://arxiv.org/abs/2210.02414>.
- [18] TAO M, LIANG X, SHI T, *et al.* RoleCraft-GLM: Advancing Personalized Role-Playing in Large Language Models[EB/OL]. (2024-01-18) [2024-10-15]. <https://arxiv.org/abs/2401.09432>.
- [19] 马鸣. 浅谈撒拉族的古民居建筑[J]. 中国土族, 2015(4): 68-70.
- MA M. On the Ancient Residential Buildings of Salar Nationality[J]. *China's TU Natio*, 2015(4): 68-70.
- [20] 邵超. 撒拉族松木大房民居木构架营造技艺研究[D]. 西安: 西安建筑科技大学, 2020. DOI: 10.27393/d.cnki.gxazu.2020.001448.
- SHAO C. Study on Wood Frame Construction Technology of Salar Pine House[D]. Xi'an: Xi'an University of Architecture and Technology, 2020. DOI: 10.27393/d.cnki.gxazu.2020.001448.
- [21] 由懿行. 青海撒拉族传统民居门窗研究[D]. 西安: 西安建筑科技大学, 2018.
- YOU Y X. Study on Windows and Doors of Salar Traditional Houses in Qinghai Province[D]. Xi'an: Xi'an University of Architecture and Technology, 2018.
- [22] SIVAKUMAR R, ARIVOLI P V. Ontology Visualization PROTÉGÉ Tools-A Review[J]. *Int J Adv Inf Technol (IJAIT)*, 2011, **1**: 1-7. DOI: 10.5121/ijait.2011.1401.
- [23] STUDER R, BENJAMINS V R, FENSEL D. Knowledge Engineering: Principles and Methods[J]. *Data Knowl Eng*, 1998, **25**(1/2): 161-197. DOI: 10.1016/S0169-023X(97)00056-6.
- [24] NOY N F, MCGUINNESS D L. Ontology Development 101: A Guide to Creating Your First Ontology: Tanford Knowledge Systems Laboratory Technical Report KSL-01-05 and Stanford Medical Informatics Technical Report SMI-2001-0880[R/OL]. (2001) [2024-11-12]. https://protege.stanford.edu/publications/ontology_development/ontology101.pdf.
- [25] MELZ E. Enhancing LLM Intelligence with ARM-RAG: Auxiliary Rationale Memory for Retrieval Augmented Generation[EB/OL]. (2023-11-07)[2024-10-15]. <https://arxiv.org/abs/2311.04177>.