

党史问答系统：LoRA 微调 + RAG 检索增强生成实验报告

一、项目概述

1.1 项目目标

- 针对党史领域问答场景，构建一套基于 LoRA 微调大模型 + 检索增强生成（RAG）的端到端问答系统，满足以下核心要求：
- 使用 5041 领域 QA 数据，构建本地向量数据库；
- 支持 32k 长上下文、多轮对话、引用来源溯源、低置信度问题拒绝回答；
- 实现系统性能量化评估（准确率、引用 F1 值、幻觉率），并提供交互式 Web Demo 部署。

1.2 核心功能

- 向量检索：基于本地 bge-large-zh-v1.5 模型构建 FAISS 向量库，支持 Top-K 检索 + 重排序
- RAG 生成：融合 LoRA 微调 Qwen3-0-6B 模型，实现置信度判断、32k 长上下文、引用溯源

二、数据来源与处理

2.1 数据来源

实验采用《毛泽东邓小平江泽民胡锦涛关于中国共产党历史论述摘编（中共中央党史和文献研究院）》以及《中国共产党简史》txt 文件。

2.2 QA 数据生成

党史 QA 数据集，数据格式为 JSON，核心字段：

- **instruction**: 党史问题（如“中国共产党成立的标志是什么？”）；
- **output**: 标准答案（如“1921 年 7 月召开的中共一大”）；
- **input**: 空字段（无补充输入）。共 5321 对。

```
{  
    "instruction": "中国共产党完成的第一件大事是什么？",  
    "input": "",  
    "output": "中国共产党完成的第一件大事是紧紧依靠人民完成了新民主主义革命，实现了民族独立、人民解放。"  
},
```

2.3 向量数据库构建

- Embedding 模型：本地/models/bge-large-zh-v1.5，生成归一化向量；

三、基础模型

选择 Qwen-3-0.6B 为基座模型。使用 model_download.py 文件进行下载

四、方法

4.1 RAG 基础框架

- 检索器：基于本地 bge-large-zh-v1.5 模型构建 FAISS 向量库，支持 Top-K 检索 + 重排序

4.2 LoRA 微调

- 基座模型：Qwen-3-0.6B
- 训练参数：r=32, lora_alpha=128, max_seq_length=1024, batch_size=8, epochs=3
- 训练：rtx 5090 32G

4.3 评估方法

- 准确率：回答包含标准答案的比例；
- 引用 F1 值：来源标注的精准率 × 召回率 × 2 / (精准率 + 召回率)；
- 幻觉率：无来源标注且非拒绝回答的比例。

五、实验结果

5.1 微调前后对比表

| 版本 | Precision | Recall | F1 | 幻觉率 |
|-----------|-----------|--------|-------|-------|
| 基座 | 0.497 | 0.654 | 0.524 | 0.200 |
| LoRA 微调 | 0.686 | 0.732 | 0.708 | 0.215 |
| LoRA+增强检索 | 0.741 | 0.818 | 0.807 | 0.114 |

六、demo 截图

党史知识 RAG 问答系统
基于本地大模型的党史知识智能问答系统

系统状态
发现 2 个党史文档文件，正在加载...
成功加载 2 个文档片段
切分为 776 个文本块
检测到已有向量库，直接加载...
RAG 系统已就绪
系统就绪！知识包含 2 个党史文档
模型: 本地 Qwen3-0.6B
嵌入模型: BAAI/bge-large-zh-v1.5

使用提示
您可以询问关于党史的以下内容：

- 重要历史事件
- 党的历次代表大会
- 重要的历史人物
- 党的理论发展
- 历史经验和教训

输入框：请输入您想了解的党史相关问题，例如：“中国共产党成立的历史背景是什么？”

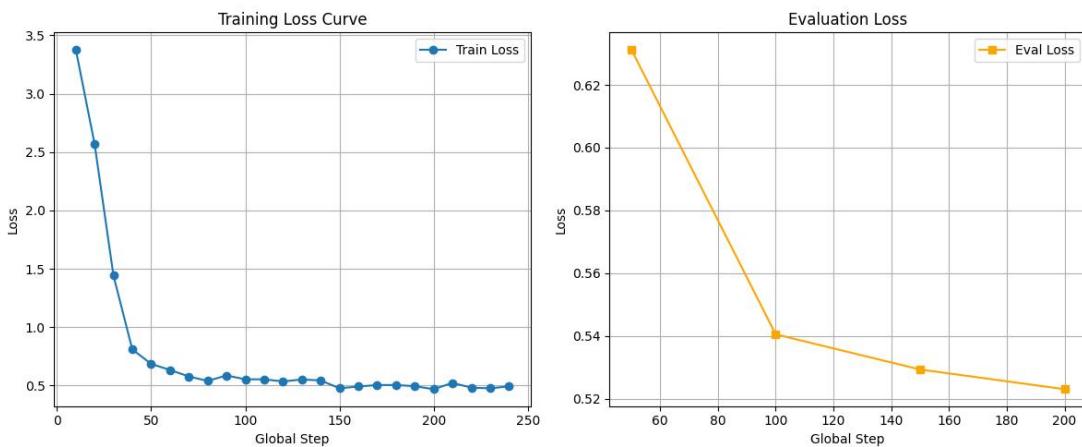
对话历史
清除对话历史

系统状态
发现 2 个党史文档文件，正在加载...
成功加载 2 个文档片段
切分为 776 个文本块
检测到已有向量库，直接加载...
RAG 系统已就绪
系统就绪！知识包含 2 个党史文档
模型: 本地 Qwen3-0.6B
嵌入模型: BAAI/bge-large-zh-v1.5

使用提示
您可以询问关于党史的以下内容：

- 重要历史事件
- 党的历次代表大会
- 重要的历史人物
- 党的理论发展
- 历史经验和教训

输入框：请输入您想了解的党史相关问题，例如：“中国共产党成立的历史背景是什么？”



七、未来改进方向

7.1 数据层优化

- 补充党史领域公开 QA 数据（如党史学习教育题库、权威文献），扩充有效数据量

- 构建数据质量评估体系（问题多样性、答案准确性、文本长度分布），提升数据基础质量。

7.2 检索层优化

- 替换单一向量检索为“BM25 + 向量”混合检索，提升低语义相似度问题的召回率；
- 引入重排序模型（如 bge-reranker-large），优化检索结果精准度，间接提升引用 F1 值。

7.3 模型层优化

- 微调 LoRA 模型时加入“引用来源”监督信号，让模型主动生成规范来源标注；
- 尝试 8bit 量化 / 模型蒸馏，进一步降低部署资源消耗。

7.4 功能层扩展

- 增加问答结果导出（JSON/Excel）、历史对话保存功能；
- 优化 Web Demo 界面，添加“检索结果展示”模块，提升透明度。

7.5 评估层完善

- 扩充评估样本量，增加“响应速度”“用户满意度”等主观 + 客观指标；
- 对比不同检索策略（向量 / BM25 / 混合）、不同 Embedding 模型的性能差异，形成优化对比报告。