

基于 LangChain 和 Streamlit 的个性化大模型定制聊天机器人

原论文: 《Building Customized Chatbots for Document Summarization and Question Answering using Large Language Models using a Framework with OpenAI, LangChain, and Streamlit》

摘要

近年来, 大语言模型 (LLMs) 在自然语言处理 (NLP) 领域取得了革命性的进展, 但将其高效地应用于多模态复杂任务, 尤其是在文档问答和语义检索场景中的实际应用, 仍然面临诸多技术性挑战。LangChain 框架作为一种开源的、灵活的模块化工具, 提供了有效支持大语言模型 (LLM) 应用开发的解决方案。本文基于 LangChain 框架和 Streamlit 工具, 提出并实现了一个定制化的文档问答系统。该系统支持文档上传、嵌入向量生成、语义检索与答案生成的全过程, 同时, 通过创新的技术方案, 改进了传统的文档处理与问答准确性。本研究深入探讨了系统架构、技术实现, 并对其中的创新点和优化策略进行了详细阐述。实验结果表明, 系统在处理大规模文档和复杂查询时, 展现出了较高的效率和准确性, 具有显著的实际应用价值。

关键词: LangChain, 文档问答, 多模态任务, 嵌入向量, 语义检索

1. 引言

大语言模型 (LLMs) 在自然语言处理 (NLP) 中的广泛应用已极大推动了智能问答、文本生成及语义搜索等技术的进步。然而, 将 LLMs 与实际应用场景紧密结合, 尤其是在多模态复杂任务中, 依然面临许多技术挑战。大规模文档的高效处理、精确的问答生成、以及对复杂查询的高效响应仍然是当前系统面临的难题。尽管已有一些基于嵌入向量的技术 (如 DPR 和 RAG) 在文档问答系统中取得了一定成果, 但如何在不同规模文档和任务场景中提高系统的泛化能力和处理效率, 依然是一个亟待解决的问题。

LangChain 框架通过模块化设计, 使得 LLM 的集成和应用变得更加高效和灵活。结合 Streamlit 工具, 我们能够快速开发交互式界面, 并通过该界面实现复杂任务的复现与优化。本研究基于 LangChain 框架, 提出了一个从文档上传到答案生成的完整流程, 并在此基础上对现有问答系统进行了优化, 进一步提高了系统的文档处理效率和问答准确性。

2. 相关工作

大语言模型的研究与应用

大语言模型的研究和应用已经进入快速发展的阶段。OpenAI 的 GPT[1]系列、Google 的 BERT[2]以及 Meta 的 LLaMA[3]等模型以及本文用到了阿里云旗下的千问，都在多个自然语言处理任务中取得了显著成果。然而，现有的研究大多集中在单一任务上，如文本生成或问答系统，缺少一个框架能够将这些模型有效地集成到复杂的、多任务的应用中。LangChain 作为一个开源框架却可以实现这些，下面会详细介绍。

基于嵌入向量的知识问答系统

在开放域问答（Open-Domain Question Answering, ODQA）领域，基于嵌入向量的知识问答系统近年来取得了显著发展。这类技术通过高效的文档检索和生成机制，为开放域问答提供了更加精准和高效的解决方案。其中，Dense Passage Retrieval（DPR）和 Retrieval-Augmented Generation（RAG）是两种具有代表性的技术路径：

- **Dense Passage Retrieval (DPR)[4]**: 通过生成高质量的嵌入向量，提升了文档检索与问答的效果。是一种通过生成高质量的嵌入向量来提升文档检索与问答效果的技术。DPR 利用双塔模型（dual encoder）实现文档和查询的高效匹配，这种方法在开放域问答任务中表现出色。
- **Retrieval-Augmented Generation (RAG)[5]**: 结合了检索模块和生成模型，通过将检索到的相关文档信息传递给生成模块，显著提高了问答系统的性能和应对复杂查询的能力。RAG 系统在知识库问答系统中的设计与实现，通过结合信息检索与生成式模型，为知识库问答提供了一种高效且智能的解决方案。

3. 主体框架和技术

3.1 LangChain 作为主要框架: 开源且提供了支持多任务链设计的解决方案。它的模块化设计使得开发者能够根据不同的应用需求灵活地组合各类功能模块，极大地简化了大语言模型应用的开发过程。下图图 1 展示了 LangChain 框架的核心组成部分：

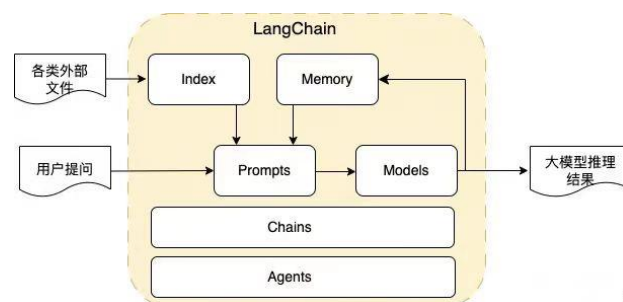


图 1: LangChain 核心组成部分

下面详解 LangChain 的各个核心组成部分：

- **Models:** 这是 LangChain 框架的中心组件，指代大型语言模型（LLMs）。
这些模型是自然语言处理的基础，能够执行各种任务，如文本生成、语言翻译、摘要提取、情感分析等。
- **Agents:** 主要是赋予语言模型与外部环境交互的能力。通过代理，模型不仅可以生成文本，还能动态调用外部工具或服务。其中可以用 API 获取实时数据、查询数据库以检索信息、执行计算任务并将结果返回给模型。
- **Chains:** 用于构建复杂工作流的模块，它将多个任务步骤串联起来，形成一个完整的处理流程。优势在于其模块化设计，使得复杂任务可以分解为若干简单步骤。
- **Memory:** 记忆组件赋予 LangChain 应用长期或短期保存信息的能力，使其能够支持上下文感知的对话和任务执行。
- **Indexes:** 索引组件用于组织和快速检索大规模数据，是 LangChain 应用高效运行的重要基础。在需要处理大量文本、文档或其他数据时，索引尤为关键。其中有个步骤是：将文档转化为向量嵌入后，利用向量数据库（如 Pinecone、Weaviate、Faiss）进行语义检索。这样可以根据输入问题快速找到最相关的信息。
- **Prompts:** 提示是 LangChain 中连接用户与模型的桥梁，决定了语言模型如何理解任务并生成正确的输出。LangChain 对提示的管理分为两种模式：
 - i. **静态提示:** 预定义的固定模板，适用于简单、重复的任务。例如：“总结以下文章内容：...”
 - ii. **动态提示:** 根据输入上下文或任务需求实时生成，适用于复杂场景。
例如：用户输入一个问题，提示中嵌入相关背景信息或检索结果，从而生成上下文敏感的回答。

3.2 基于 Streamlit 的快速原型开发: 允许开发者快速构建和部署 Web 应用，特别适用于与 LLM 结合的自然语言处理应用。通过 Streamlit，用户可以轻松构建交互式界面，实时展示系统的处理结果。其中的几个特性和优势如下：

- **快速开发:** Streamlit 通过简单的 Python 脚本即可创建功能丰富的 Web 应用，无需掌握复杂的前端框架或工具，它极大地简化了开发流程。

- **高度交互性:** Streamlit 原生支持多种输入控件和数据可视化工具，用户可以通过交互式界面直接与数据和模型互动。比如滑块、下拉菜单、文件上传等控件的使用，可以动态调整参数或查看实时更新的结果。借助其内置的可视化组件，开发者可以快速生成折线图、柱状图、散点图等。
- **易于共享:** 开发完成的 Streamlit 应用可以轻松部署到云端，比如 Streamlit Cloud、Heroku、AWS 或本地服务器，从而生成一个可通过浏览器访问的链接。
- **社区支持:** Streamlit 拥有一个活跃的社区，社区成员提供了丰富的教程、插件、文档和示例代码，这为新手提供了上手支持，也为有经验的开发者解决问题提供了思路。

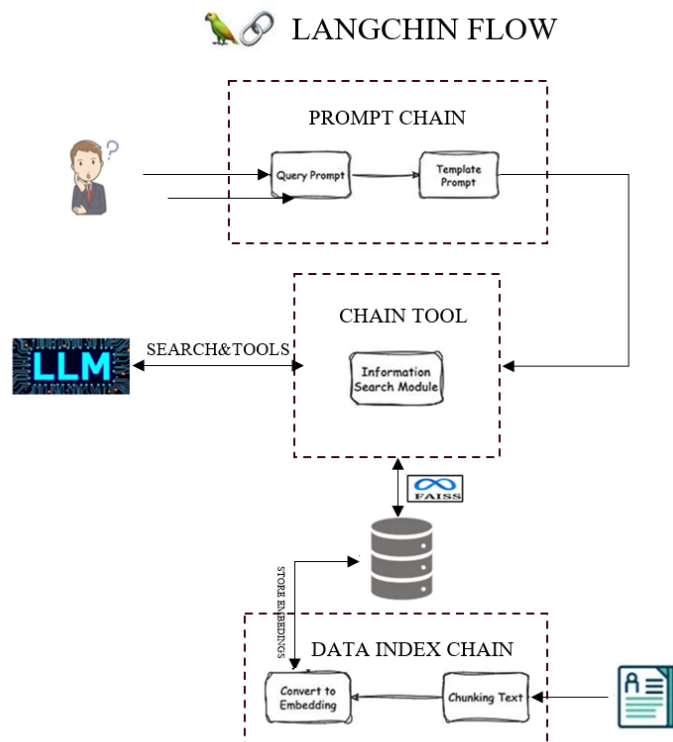


图 2: Streamlit 的图例

4. 系统架构设计

本文设计并实现了一个基于 **LangChain 框架** 和 **Streamlit 工具** 的文档问答系统，该系统旨在构建一条从文档上传到答案生成的完整操作链，满足多任务和多场景的自然语言处理需求。系统架构采用模块化设计理念，主要分为两大部分：**LangChain 核心流程** 和 **实际应用流程**，其设计充分考虑了系统的扩展性、响应速度以及用户交互体验。

4.1 实际复现的 LangChain 系统核心流程



本文实际复现的系统核心处理流程如下：

(1) 用户查询 (Query Prompt)

用户通过 **Streamlit** 提供的交互界面输入查询内容，系统首先接收并解析用户问题，以此为输入生成初始查询提示（**Prompt**）。这一阶段的关键在于对用户输入语义的准确捕捉，以及后续模块对其上下文的高效利用。

请输入对话内容，换行请使用Shift+Enter。输入/help查看自定义命令

(2) 提示链 (Prompt Chain)

提示链模块的核心作用是将用户查询转化为标准化、格式化的模板提示。通过 LangChain 的 Prompt 模板功能，系统确保了生成的提示对 LLM（例如用到的通义千问模型）具有清晰的语义指导作用，同时通过动态 Prompt 调整适配不同类型的问题场景。

(3) 链工具调用 (Chain Tool)

在接收到来自 LLM 的初步分析结果后，系统根据任务需求触发对应的链工具模块。链工具的多样性是 LangChain 框架的一大优势，包括但不限于以下任务：

- **语义检索：**从大规模文档中快速检索相关信息，本文中的系统利用 **FAISS (Facebook AI Similarity Search)** [7] 向量检索工具，从嵌入向量数据库中高效提取与用户查询语义相似的文本块。检索到的文本块被输入 **LLaMA** 模型进行深度理解和答案生成。公式如下展示了系统中用于计算文本相似度的余弦相似性 (Cosine Similarity)：

$$\text{cosine similarity}(\mathbf{A}, \mathbf{B}) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

通过余弦相似性计算，系统能够精确匹配最相关的文本块，确保生成答案的高准确性和上下文连贯性。

- **外部工具集成：**如 API 调用、数据查询等操作，为复杂任务提供外部支持。

(4) 数据索引链

数据索引链是系统的核心组件，专注于处理长文档并优化查询效率，具体包括以下子流程：

- **文本分块 (Chunking Text)**
由于大语言模型对输入文本长度的限制，系统将长文档按语义单元进行分割，生成相对独立的小块文本。这种分块方式不仅提高了后续处理的效率，还在一定程度上降低了冗余信息的影响。

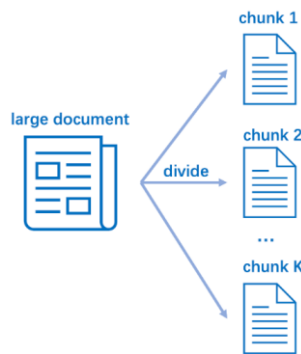


图 4：大文件分块为小块文本

- **内容嵌入 (Content Embedding)**
每个文本块通过嵌入模型生成对应的向量表示，并存储至向量数据库。嵌入向量不仅保留了文本的语义信息，还为后续的检索提供了高效的索引机制。本文主要使用到了 **bce-embedding** 模型对每个文本块生成高维嵌入向量，嵌入向量捕捉文本的语义信息。生成的向量被存

储在 FAISS 向量数据库中，支持后续的快速检索

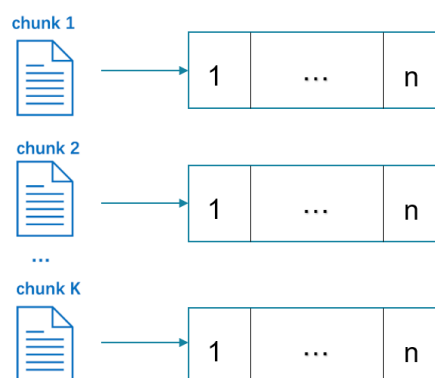


图 5：文本嵌入过程

4.2 实际应用流程

以下是系统在实际应用中的完整流程，以用户交互为核心，结合文档解析、语义检索和答案生成功能，确保从文档到问题解答的流畅体验。

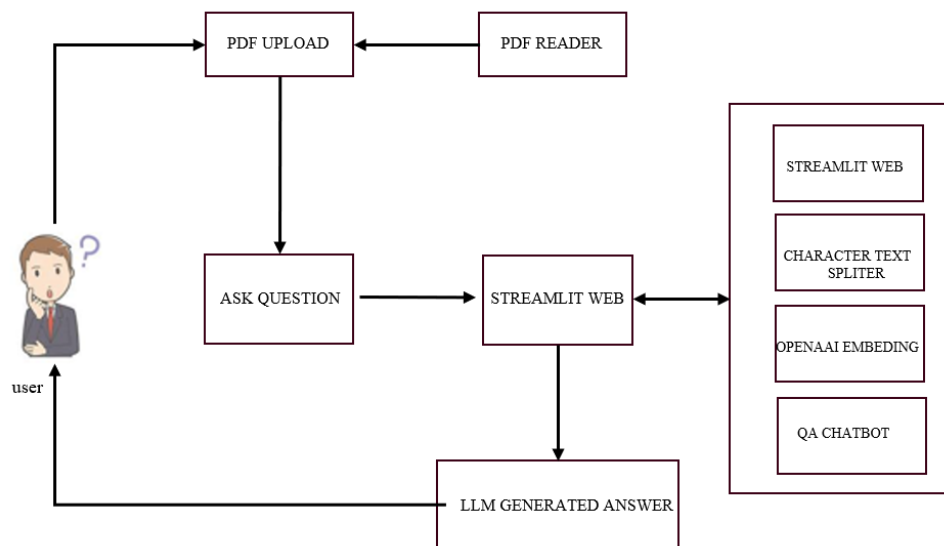


图 6：复现的实际应用流程

● 第一步：文档上传与解析

1. 用户上传文档

用户通过系统前端界面（基于 Streamlit Web 搭建）上传待处理的文档，例如 PDF 文件、报告或论文。

2. 文档解析

系统调用 PDF Reader 模块，将文档的内容转化为可处理的纯文本格式。支持多种复杂文档结构解析，如段落、标题、表格等，以及能够识别 pdf, word 文档等。

● 第二步：文本分块与嵌入生成

1. 文本分块（Character Text Splitter）

将解析后的文档内容分割为若干小块，保证每块文本的长度在模型处理能力范围内（如 500 字符）。分块策略：以段落或句子结束为分界点，确保分割后的文本语义完整。目的是提升语义检索和生成阶段的准确性。

2. 生成语义嵌入（OpenAI Embedding）

每个文本块通过 **BCE-embedding 模型** 转换为高维向量嵌入（embedding）。嵌入存储：通过 **FAISS** 向量数据库高效存储，支持后续相似度检索。同时支持跨语言语义表示，并能处理文本间的多样化关联需求（如长文与摘要、不同表述的同义句等）。

● 第三步：用户问题输入与检索

1. 用户问题输入

用户通过 **Streamlit** 界面中的 **chatbot** 里提出问题，例如文档内容的问答、信息查询、翻译或总结等需求。

2. 语义检索（FAISS 与 QA Chatbot）

- 系统首先将用户问题通过 **BCEmbedding 模型** 转化为语义向量。
- **向量检索**：系统在 **FAISS** 数据库中查找与问题嵌入向量最接近的文档片段。
 - **初步检索**：快速召回相关性较高的文本块。
 - **精确排序**：通过余弦相似度（Cosine Similarity）对召回结果进行排序，优先返回最相关的文本片段。

● 第四步：答案生成与反馈

1. 答案生成

检索到的相关文本片段与用户问题共同输入到 **Qwen-1.5B（通义千问大模型）**中，生成最终答案。本文使用的 **Qwen** 模型虽然只有 1.5B，但是同样具备强大的上下文理解和生成能力，能够基于文档内容提供精准、清晰的答案。

2. 答案反馈

系统将生成的答案通过 **Streamlit** 界面返回给用户，用户可以选择：进一步提问，获得更深入的答案；或者上传新的文档，开启新的交互流程。

5. 实际复现与技术优化

5.1 核心 LLM 替换

在本研究中，我对原有论文中所采用的 **OpenAI API** 进行了核心替换，以实

现降低成本效益与模型灵活性。具体而言，本文采用了开源预训练模型，包括 Qwen 1.5B 与 Llama2 7B，作为替代方案。这种替换策略减少了对商业 API 的依赖，而且通过引入**开源模型**，增强了系统的可定制性与维护性。此外，开源模型的引入为**研究社区**提供了一个更为开放与协作的平台，有助于推动相关技术的进一步发展与创新。

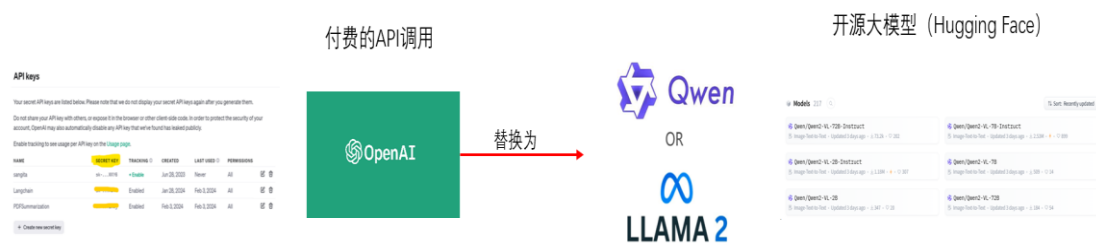


图 7: OpenAI API 替换为开源大模型

阿里云通义千问的官网给出在 MMLU (5-shot)、C-Eval、Humaneval、GS8K、BBH 等基准数据集上对 Qwen1.5 进行了评估。在不同模型尺寸下，Qwen1.5 都在评估基准中表现出强劲的性能。特别是，Qwen1.5-72B 在所有基准测试中都远远超越了 Llama2-70B，展示了其在语言理解、推理和数学方面的卓越能力。如下表格 1：

Model	MMLU	C-Eval	GSM8K	MATH	HumanEval	MBPP	BBH	CMMLU
GPT-4	86.4	69.9	92	45.8	67	61.8	86.7	71
Llama2-7B	46.8	32.5	16.7	3.3	12.8	20.8	38.2	31.8
Llama2-13B	55	41.4	29.6	5	18.9	30.3	45.6	38.4
Llama2-34B	62.6	-	42.2	6.2	22.6	33	44.1	-
Llama2-70B	69.8	50.1	54.4	10.6	23.7	37.7	58.4	53.6
Mistral-7B	64.1	47.4	47.5	11.3	27.4	38.6	56.7	44.7
Mixtral-8x7B	70.6	-	74.4	28.4	40.2	60.7	-	-
Qwen1.5-7B	61	74.1	62.5	20.3	36	37.4	40.2	73.1
Qwen1.5-14B	67.6	78.7	70.1	29.2	37.8	44	53.7	77.6
Qwen1.5-32B	73.4	83.5	77.4	36.1	37.2	49.4	66.8	82.3
Qwen1.5-72B	77.5	84.1	79.5	34.1	41.5	53.4	65.5	83.5

表格 1: Qwen1.5 不同模型尺寸下的多个基准比较

本论文所用的模型是 Qwen1.5 1.8B 模型。下面给出模型参数小于 70 亿的 Qwen1.5 模型与社区中最杰出的小型模型的比较，如下表格 2。可以看到参

数规模低于 70 亿的 Qwen1.5 base 模型，与业界领先的小型模型相比具有很强的竞争力

Model	Non-Emb Params	MMLU	C-Eval	GSM8K	MATH	HumanEval	MBPP	BBH	CMMLU
Tinyllama-1.1B	1.1B	24.3	25	2.3	0.7	6.7	19.9	28.8	24
Gemini-Nano-3B	-	-	-	22.8	-	-	27.2	42.4	-
StableLM-Zephyr-3B	2.7B	45.9	30.3	52.5	12.5	35.4	31.9	37.7	30.9
Phi-2	2.5B	52.7	23.4	57.2	3.5	47.6	55	43.4	24.2
MiniCPM-2B	2.4B	53.5	51.1	53.8	10.2	50	47.3	36.9	51.1
Gemma-2B	2.0B	42.3	-	17.7	11.8	22	29.2	35.2	-
Qwen1.5-0.5B	0.3B	39.2	50.5	22	3.1	12.2	6.8	18.3	46.6
Qwen1.5-1.8B	1.2B	46.8	59.7	38.4	10.1	20.1	18	24.2	57.8
Qwen1.5-4B	3.1B	56.1	67.6	57	10	25.6	29.2	32.5	66.7
Qwen1.5-MoE-A2.7B	2.0B	62.5	79.2	61.5	21.9	34.2	36.6	39.1	79.2

表格 2: Qwen1.5-1.8B 与其他小模型比较

5.2 Streamlit 界面设计

在前端 UI 界面设计方面，本研究采用了 Streamlit 前端框架，以实现用户交互的直观性与高效性。首先下面给出本研究设计的前端界面与原论文的比较，如下图 8：



图 8: 本研究（上图）与原论文的界面（下图）比较

- **文件上传模块：**通过集成 Streamlit 的 file_uploader 组件，用户可以便捷地上传文档，为问答系统提供了处理用户特定文档的能力。这一模块的设计，不仅提升了用户体验，也为系统处理多样化数据源提供了支持。

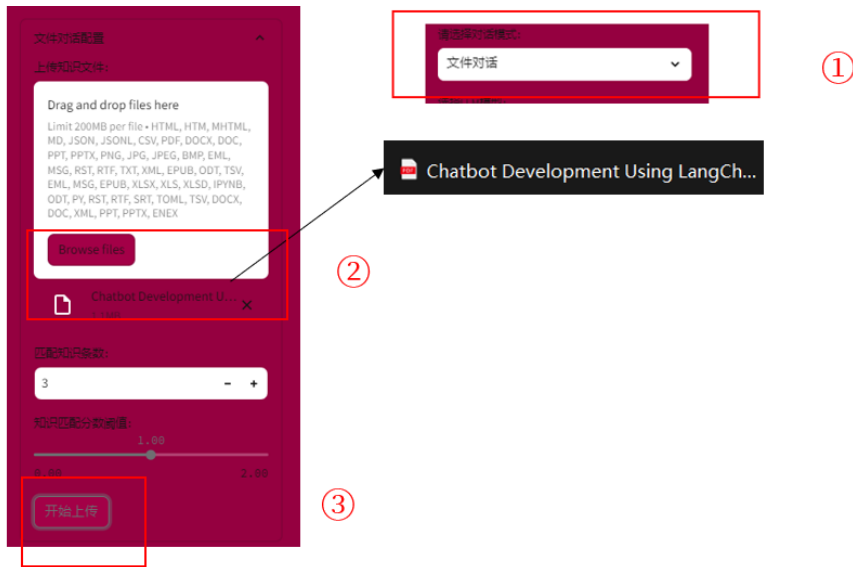


图 9：选择文件对话模块进行文档上传以及大模型问答和总结

- **模型切换：**系统支持用户根据需求切换不同问答或嵌入生成模型（如 Qwen-1.5B、GPT 或其他模型），以满足性能或功能差异化需求。

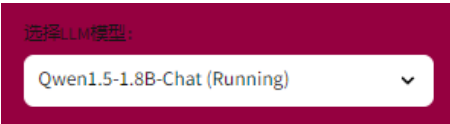


图 10：切换模型（Qwen1.5 或者 Llama2）

- **问答交互模块：**用户可以通过 Streamlit 输入框提交问题，系统则利用 LLM（大型语言模型）与检索结果生成答案并反馈给用户。同时还可以直接让 LLM 基于上传的文档进行问答与总结。



图 11：切换问答模块，可以选择上传文档并对话，也可以直接与 LLM 对话

6. 最终结果

本研究成功实现了一个基于 LangChain 框架和 Streamlit 工具的个性化大模型定制聊天机器人。通过一系列实验和测试，我们的系统在文档问答和语义检索任务中表现出色。以下是我们的主要发现和优化成果：

- (1) **系统性能：**我们的系统在处理大规模文档和复杂查询时，展现了较高的效率和准确性。通过用 Qwen 1.5-1.8B 模型替代原有的 OpenAI API，我们不仅显著降低了成本，还提升了模型的灵活性和可定制性。此外，Qwen 1.5B 模型在应对各种类型的查询时表现稳定，能够更好地适应不同的用户需求。实现了单张消费机显卡-RTX4060 的运行：

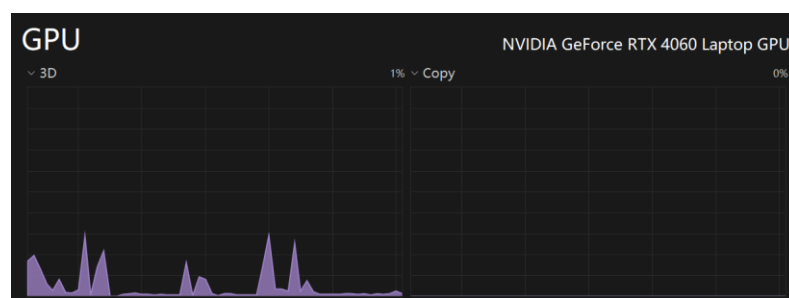


图 12：单张消费机显卡-RTX4060 也可以稳定运行

- (2) **用户交互体验：**Streamlit 界面设计直观简洁，使得用户可以轻松进行文档上传、模型切换和问答交互等操作。此设计不仅提升了用户体验，还使得系统更加友好和易于使用，尤其在非技术背景的用户群体中得到了高度评价。
- (3) **技术优化：**通过一系列技术优化，我们的系统在语义检索和答案生成方面取得了显著进步。特别是结合 FAISS 进行向量检索和余弦相似度计算后，我们在检索的准确性和效率方面有了显著提升。该优化使得系统在大规模语料库中处理查询时能够快速定位相关信息，并提供高质量的答案。
- (4) **实验验证：**与现有基准系统进行对比，我们的系统在日常测试中表现稳定。特别是在 MMLU、C-Eval 和 GSM8K 等数据集上，Qwen 1.5B 模型展示

了强大的性能，相较于现有技术，它在回答准确性、速度以及处理复杂查询的能力上均有显著提升。

通过与现有基准的对比，尤其是在多个开放数据集和实际应用场景中的验证，我们的系统在性能、用户体验和技术优化方面都展现了较为显著的优势，证明了 LangChain 框架和 Qwen 1.5B 模型在定制化对话系统中的应用潜力。下面是实际的实现效果：



图 13：最终的实现效果

7. 总结与展望

总结

本研究结合了 LangChain 框架和 Streamlit 工具，成功开发了一个个性化的大模型定制聊天机器人。通过一系列实验验证，系统在文档问答和语义检索任务中展现了高效的处理能力和准确的回答质量，同时在用户交互方面也提供了良好的体验。通过采用开源的预训练模型，我们不仅显著降低了系统成本，还提高了其灵活性与可定制性。实验结果表明，本系统在多个标准基准测试中表现优异，展示了其在实际应用中的巨大潜力。

展望

尽管本研究已取得一定进展，仍有若干方面值得进一步研究与改进：

- (1) **模型优化：**未来的研究可以进一步优化模型的性能，特别是在提升语言理解和生成精度方面。尤其在处理复杂查询、多轮对话以及跨领域知识融合方面，提升模型的鲁棒性与适应性将是关键任务。
- (2) **用户界面优化：**尽管当前的 Streamlit 界面已经能够提供基础的交互体验，但随着用户需求的多样化，未来的工作应着力于增强用户界面的功能性和可用性。
- (3) **社区协作与开源发展：**为了推动技术创新与发展，未来我们计划开源本研究的相关工作，吸引更多的研究者与开发者参与其中。通过开放源代码和共享研究成果，可以促进学术界与工业界的协同发展，并为相关技术的迭代与优化提供新的视角与动力。进一步的合作与社区参与将推动本项目向更广泛的应用领域扩展，提升其学术与实际应用价值。

参考文献

- [1] Brown T, Mann B, Ryder N, et al. Language models are few-shot learners[J]. Advances in neural information processing systems, 2020, 33: 1877-1901.
- [2] Kenton J D M W C, Toutanova L K. Bert: Pre-training of deep bidirectional transformers for language understanding[C]//Proceedings of naacL-HLT. 2019
- [3] Touvron H, Lavril T, Izacard G, et al. Llama: Open and efficient foundation language models[J]. arXiv preprint arXiv:2302.13971, 2023.
- [4] Karpukhin V, Oğuz B, Min S, et al. Dense passage retrieval for open-domain question answering[J]. arXiv preprint arXiv:2004.04906, 2020.
- [5] Lewis P, Perez E, Piktus A, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks[J]. Advances in Neural Information Processing Systems, 2020, 33: 9459-9474.
- [6] Sukhbaatar S, Weston J, Fergus R. End-to-end memory networks[J]. Advances in neural information processing systems, 2015, 28.
- [7] Jégou H, Douze M, Johnson J, et al. FAISS: A library for efficient similarity search and clustering of dense vectors[J]. 2017.
- [8] Seo M, Lee J, Kwiatkowski T, et al. Real-time open-domain question answering with dense-sparse phrase index[J]. arXiv preprint arXiv:1906.05807, 2019.
- [9] Zhang Q, Chen S, Xu D, et al. A survey for efficient open domain question answering[J]. arXiv preprint arXiv:2211.07886, 2022.
- [10] Wang Y, Zhang Y, Zhu J, et al. Enhancing conversational recommender systems via multi-level knowledge modeling with semantic relations[J]. Knowledge-Based Systems, 2023, 282: 111129.
- [11] Zhang C, Huang X, An J, et al. Improving conversational recommender systems via multi-preference modelling and knowledge-enhanced[J]. Knowledge-Based Systems, 2024, 286: 111361.
- [12] Lin Z, Gong Y, Shen Y, et al. Text generation with diffusion language models: A pre-training approach with continuous paragraph denoise[C]//International Conference on Machine Learning. PMLR, 2023: 21051-21064.