

三 茴香豆搭建你的 RAG 智能助理笔记

第三章

一、RAG技术概述

RAG (Retrieval Augmented Generation) 技术，结合了信息检索和生成式问答的模型，能够通过检索用户输入相关的信息片段，并利用这些信息片段辅以外部知识库，以生成更准确、更丰富的回答。这样的技术可以解决大型语言模型 (LLMs) 在处理知识密集型任务时可能出现的问题，例如内容的虚构、信息的过时、以及缺乏清晰可追溯的推理过程。RAG的优势在于可以提升回答的准确度、降低推理的成本，并且能够实现外部知识的更新，无需再次对模型进行训练即可掌握新的知识领域。

二、环境配置与模型部署

在Intern Studio的服务器上配置和部署RAG技术。具体包括配置Conda环境、复制模型文件、以及git clone模型相关的代码。

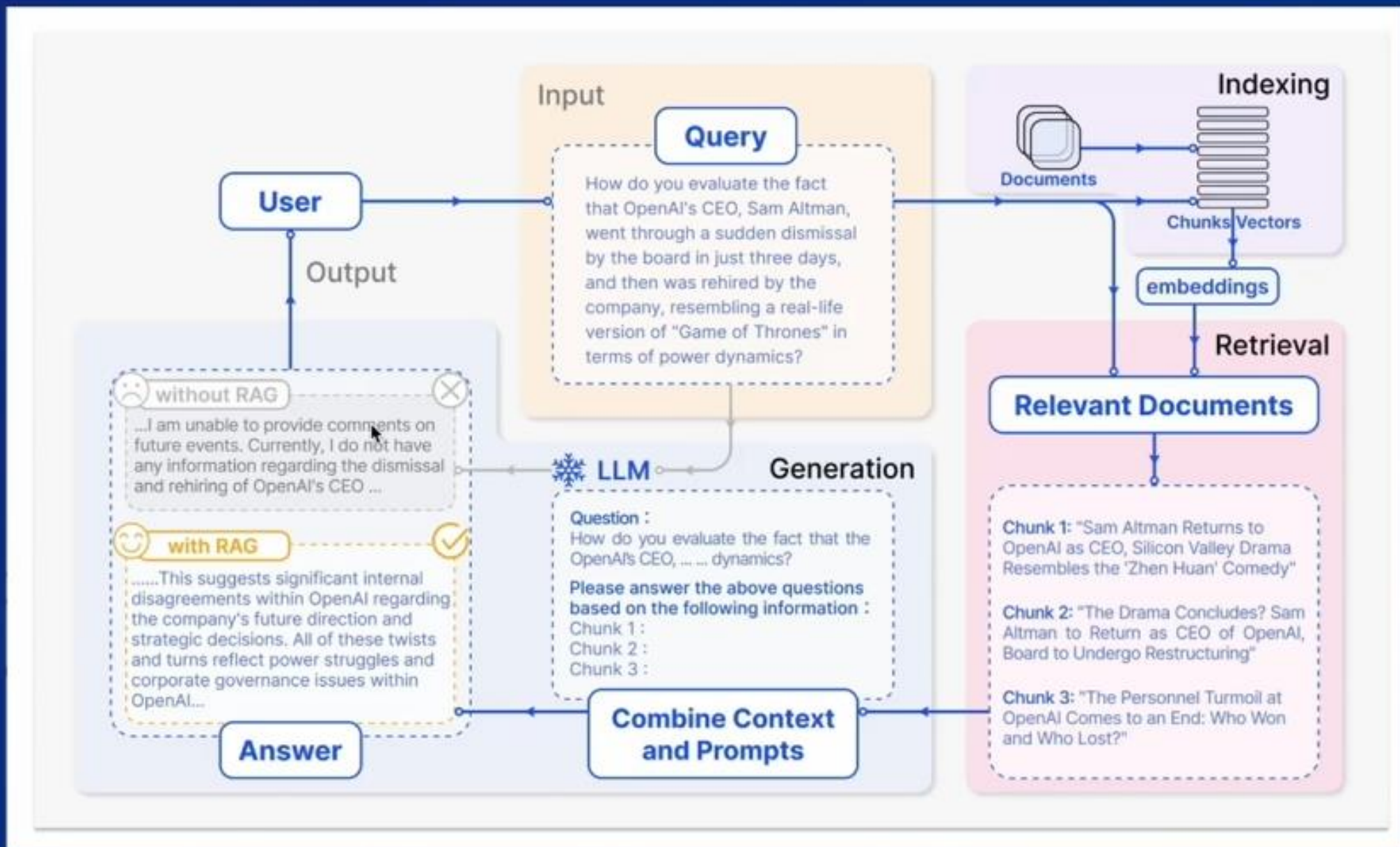
三、使用茴香豆搭建RAG助手

首先是修改配置文件，指定向量数据库、词嵌入模型、用于检索的重排序模型，以及所使用的语言模型（internlm2-chat-7b）。接下来创建知识库，使用InternLM的Huixiangdou文档作为新增知识的检索来源，以此建立技术问答助手。此外，还需要提取知识库特征来创建向量数据库，使用的是LangChain的相关模块和网易BCE双语模型。

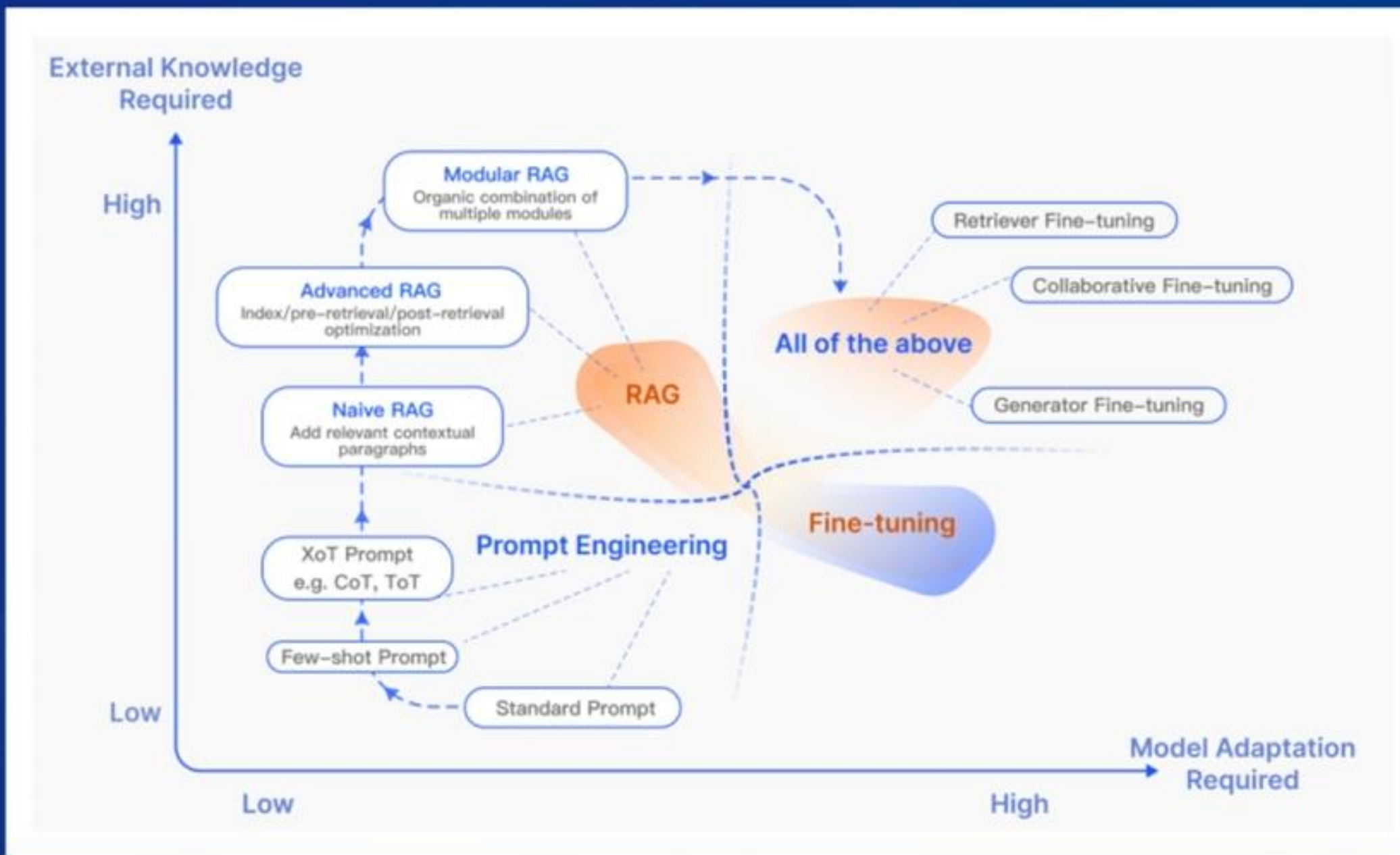
为了让RAG更加精准地判断提问的相关性，茴香豆还建立了接受问题和拒绝问题两个向量数据库，其中包括希望茴香豆助手回答或拒答的示例问题。例如，接受问题列表中的问题存储在huixiangdou/resource/good_questions.json中，而拒绝问题列表中的则存储在huixiangdou/resource/bad_questions.json中，主要涵盖了一些技术无关或闲聊性质的问题。完成语料来源设定后，执行命令创建RAG在检索过程中使用的向量数据库。

在实际检索过程中，茴香豆会将输入的问题与两个列表中的问题进行向量空间的相似性比较，以判断问题是否应当回答，从而避免了问答的泛滥。确定要回答的问题之后，基础模型会提取关键词，并在知识库中检索与之最相似的信息片段（chunk），最终结合这些信息生成回答。

RAG 流程示例



LLM 模型优化方法比较





复制完成后，在本地查看环境。

```
conda env list
```

结果如下所示。

```
# conda environments:
#
base                * /root/.conda
InternLM2_Huixiangdou /root/.conda/envs/InternLM2_Huix:
```

运行 **conda** 命令，激活 **InternLM2_Huixiangdou** **python** 虚拟环境：

```
conda activate InternLM2_Huixiangdou
```

环境激活后，命令行左边会显示当前（也就是 **InternLM2_Huixiangdou**）的环境名称，如下图所示：

```
=====
ALL DONE!
=====

(base) root@intern-studio-40059224:~# conda env list
# conda environments:
#
base                * /root/.conda
InternLM             /root/.conda/envs/InternLM
InternLM2_Huixiangdou /root/.conda/envs/InternLM2_Huixiangdou
internlm-demo        /root/.conda/envs/internlm-demo
opencompass          /root/.conda/envs/opencompass
xcomposer-demo       /root/.conda/envs/xcomposer-demo
xtuner0.1.9          /root/.conda/envs/xtuner0.1.9

(base) root@intern-studio-40059224:~# conda activate InternLM2_Huixiangdou
(InternLM2_Huixiangdou) root@intern-studio-40059224:~#
```

后续教程所有操作都需要在本地环境下进行，需向开发机或打开命令后重新激活环境。



File Edit View Run Kernel Tabs Settings Help

```
root@intern-studio-40059:~#
Using cached https://pypi.tuna.tsinghua.edu.cn/packages/80/03/6ea8b1b2a5ab40a7a60dc464d3daa7aa546e0a74d74a9f8ff551ea7905db/executing-2.0.1-py2.py3-none-any.whl (24 kB)
Collecting asttokens>=2.1.0 (from stack-data->ipython>=7.23.1->ipykernel)
Using cached https://pypi.tuna.tsinghua.edu.cn/packages/45/86/4736ac618d82a20d87d2f92ae19441ebc7ac9e7a581d7e58bbe79233b24a/asttokens-2.4.1-py2.py3-none-any.whl (27 kB)
Collecting pure-eval (from stack-data->ipython>=7.23.1->ipykernel)
Using cached https://pypi.tuna.tsinghua.edu.cn/packages/2b/27/77f9d5684e6bce929f5cfe18d6cfbe5133013c06cb2fbf5933670e60761d/pure_eval-0.2.2-py3-none-any.whl (11 kB)
Installing collected packages: wcwidth, pure-eval, ptyprocess, traitlets, tornado, six, pyzmq, pygments, psutil, prompt-toolkit, platformdirs, pexpect, parso, packaging, nest-asyncio, executing, exceptiongroup, decorator, debugpy, python-dateutil, matplotlib-inline, jupyter-core, jedi, comm, asttokens, stack-data, jupyter-client, ipython, ipykernel
Successfully installed asttokens-2.4.1 comm-0.2.2 debugpy-1.8.1 decorator-5.1.1 exceptiongroup-1.2.0 executing-2.0.1 ipykernel-6.29.4 ipython-8.23.0 jedi-0.19.1 jupyter-client-8.6.1 jupyter-core-5.7.2 matplotlib-inline-0.1.6 nest-asyncio-1.6.0 packaging-24.0 parso-0.8.4 pexpect-4.9.0 platformdirs-4.2.0 prompt-toolkit-3.0.43 psutil-5.9.8 ptyprocess-0.7.0 pure-eval-0.2.2 pygments-2.17.2 python-dateutil-2.9.0.post0 pyzmq-25.1.2 six-1.16.0 stack-data-0.6.3 tornado-6.4 traitlets-5.14.2 wcwidth-0.2.13
WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package manager. It is recommended to use a virtual environment instead: https://pip.pypa.io/warnings/venv
Installed kernelspec InternLM2_Huixiangdou in /root/.local/share/jupyter/kernels/internlm2_huixiangdou
conda环境：InternLM2_Huixiangdou安装成功！

=====
ALL DONE!
=====

(base) root@intern-studio-40059224:/opt/jupyterlab# conda env list
# conda environments:
#
base                * /root/.conda
InternLM             /root/.conda/envs/InternLM
InternLM2_Huixiangdou /root/.conda/envs/InternLM2_Huixiangdou
internlm-demo        /root/.conda/envs/internlm-demo
opencompass          /root/.conda/envs/opencompass
xcomposer-demo       /root/.conda/envs/xcomposer-demo
xtuner0.1.9          /root/.conda/envs/xtuner0.1.9

(base) root@intern-studio-40059224:/opt/jupyterlab#
```