





# RAG 技术详解与实战应用

第12讲:实践:用缓存、异步与向量引擎加速你的RAG









- 1. 上节回顾
- 2. 使用向量数据库实现知识库持久化存储
- 3. 构建高效索引提升检索性能
- 4. 优化模型推理性能
- 5. Q&A



### 上节回顾







### 提升RAG启动效率

#### 持久化存储

- 利用向量数据库实现RAG系统 的数据持久化存储
- 了解了两种主流向量数据库Ch roma、Milvus的特点、适用 场景

#### 内存、硬盘与数据库

- 了解内存与硬盘的关系
- 内存缓存机制
- 标量数据库与向量数据库

### 构建索引加速检索

#### 索引的定义

- "以空间换时间",加速数据 检索的特殊数据结构
- 向量索引-高效的ANN搜索

#### Milvus

内置多种先进向量索引, 支持 复杂查询,助力企业级应用落 地

#### 系统工程优化

- 并行与并发
- 如何使用Python实现多任务并 行

### 模型推理性能优化

#### 模型推理性能指标

首字延迟、吞吐量

#### 性能优化策略

- 模型蒸馏
- 更好的推理框架
- 参数精度→模型量化











- 1. 上节回顾
- 2. 使用向量数据库实现知识库持久化存储
- 3. 构建高效索引提升检索性能
- 4. 优化模型推理性能
- 5. Q&A



### 向量数据库实现持久化存储







为了防止程序重启导致数据丢失,使用向量数据库,将存储在硬盘中,以实现数据的持久化存储。 LazyLLM 提供了多种存储组件同时支持基于内存和向量数据库两种存储方式,以供开发者灵活选 择。

组件名	存储介质	功能描述	适用场景
MapStore	内存	基于内存的存储,提供节点存储的基本能力,能以最快的速度实现数据读写	技术探索阶段,将少量文档处理数据保存在内存中,以技术选型与算法迭代的快速验证
ChromadbStore	硬盘	使用 Chroma向量数据库,实现数据的持久化存储。	轻量级实验,快速POC,缩短从构 思到可运行原型的研发周期
MilvusStore	硬盘	使用Milvus向量数据库,实现数据的本地存储/分布式远程存储。	企业级部署,适用于数据体量大,对并发、检索性能都有较高要求的场景

## LazyLLM支持的向量数据库







在基于LazyLLM实现RAG系统时,只需在实例化Document时传入存储配置,即可实现不同存储类型的切换。

doc = lazyllm.Document(dataset\_path=...,..., store conf=store conf)



map - 基于内存kv存储

chroma - 使用轻量级向量数据库ChromaDB

milvus - 基于高性能向量数据库Milvus

·kwargs:不同存储所需额外配置参数(字典的形式)

·chroma:

·dir(必填):存储数据目录

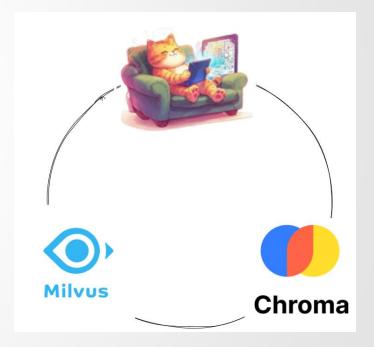
·milvus:

·uri(必填):存储数据地址,文件路径或远程连接的

url

·index kwargs (可选, dict/list[dict]): 内置索引配 置,主要包含索引类型index\_type和度量类型 metric\_type,支持指定不同向量模型使用不同索引配置。 目前支持 稠密和稀疏向量,感兴趣的开发者可以前往 官网查看详细参数







## LazyLLM支持的向量数据库







目前ChromaDB仅作为数据的持久化存储使用,而Milvus在支持数据持久化存储的同时,也支持

了多样内置向量索引的选择,以便开发者轻松实现高效的向量检索。

以下是使用不同存储时的store\_conf示例:

### MapStore

无需传入任何配置

### ChromadbStore

```
chroma_store_conf = {
   'type': 'chroma',
   'kwargs': {
      'dir': 'dbs/chroma1',
    }
}
```

请注意,目前 LazyLLM 支持 Milvus 2.4.x,请 ■ 注意您的 pymilvus >= 2.4.11 , milvus-lite == 2.4.10

#### MilvusStore

```
milvus_store_conf = {
   'type': 'milvus',
   'kwargs': {
      'uri': 'dbs/milvus1.db',
      'index_kwargs': {
         'index_type': 'HNSW',
      'metric_type': 'COSINE',
      }
   },
}
```

### 实践: 系统启动性能对比







分别使用内存、ChromaDB和Milvus作为数据存储,测试三者首次启动和二次启动耗时,以展示持久化存储对于系统二次启动效率的提升。

```
embedding_model = OnlineEmbeddingModule(
                                                                                                      # 线上embedding模型
      source="doubao",
      embed model name="doubao-embedding-large-text-240915"
1. def test store(store conf: dict=None):
                                                                                                      # 传入存储配置
      """接收存储配置,测试不同配置下系统启动性能"""
      st1 = time.time()
2.
      docs = lazyllm.Document(dataset path=DOC PATH, embed=embedding model, store conf=store conf)
                                                                                                      # 实例化Document, 传入存储配置
3.
      docs.create node group(name='sentence', parent="MediumChunk", transform=lambda x: x.split('.'))
                                                                                                      # 创建按句子划分节点组
4.
      if store conf and store conf.get('type') == "milvus":
                                                                                                      # milvus配置中已指定度量类型
5.
          retriever1 = lazyllm.Retriever(docs, group name="sentence", topk=3)
                                                                                                      # 无需similarity参数
6.
7.
      else:
          retriever1 = lazyllm.Retriever(docs, group name="sentence", similarity='cosine', topk=3)
                                                                                                      # similariy=cosine, 使用向量检索
      retriever1.start()
                                                                                                      # 启动检索器, 完成即系统启动完毕
9.
      et1 = time.time()
10.
11.
      st2 = time.time()
      res = retriever1("牛车水")
                                                                                                      # 测试单次检索耗时
12.
13.
      et2 = time.time()
      nodes = "\n=====\n".join([node.text for node in res])
14.
      msg = f"Init time: {et1 - st1}, retrieval time: {et2 - st2}s\n"
                                                                                                      # 输出系统启动耗时、检索耗
15.
      LOG.info(msg)
16.
      LOG.info(nodes)
17.
18.
      return msg
```

### 实践: 系统启动性能对比

```
# chroma db存储配置
1. chroma store conf = {
       'type': 'chroma',
       'kwargs': {
          'dir': 'dbs/chroma1',
  # milvus存储配置
2. milvus store conf = {
       'type': 'milvus',
       'kwargs': {
          'uri': 'dbs/milvus1.db',
          'index kwargs': {
              'index type': 'HNSW',
              'metric type': 'COSINE',
      },
  #测试集,依次测试使用内存、chroma、milvus时的系统启动性能
3. test conf = {
       "map": None,
      "chroma": chroma store conf,
      "milvus": milvus store conf
4. start times = ""
5. for store type, store conf in test conf.items():
      # 调用测试函数
      res = test store(store conf=store conf)
6.
      start times += f"Store type: {store type}: {res}"
8. print(start times)
```







Store type: map: Init time: 17.713993787765503, retrieval time: 0.3607618808746338s Store type: chroma: Init time: 17.690227270126343, retrieval time: 0.36835598945617676s Store type: milvus: Init time: 15.303590059280396, retrieval time: 0.02257847785949707s

#### 第一次启动

Store type: map: Init time: 17.303208112716675, retrieval time: 0.3642408847808838s Store type: chroma: Init time: 13.890714883804321, retrieval time: 0.3575289249420166s Store type: milvus: Init time: 1.8017945289611816, retrieval time: 0.020511150360107422s

#### 第二次启动

存储类型	第一次启动/检索耗时(s)	第二次启动/检索耗时(s)
Мар	17. 71/0. 36	17. 30/0. 36
ChromaDB	17. 69/0. 37	13.89 ( 121.5%) /0.36
Milvus	15. 30/0. 02	1.80(\$\.\dagge\$88.2%)/0.02

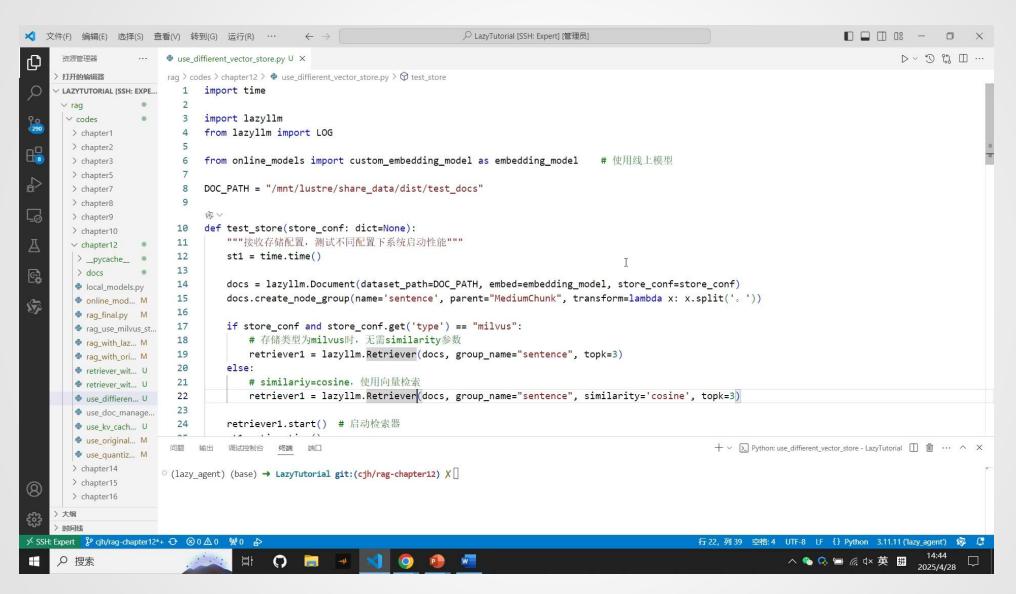


### 实践: 系统启动性能对比 - 实机演示



















- 1. 上节回顾
- 2. 使用向量数据库实现知识库持久化存储
- 3. 构建高效索引提升检索性能
- 4. 优化模型推理性能
- 5. Q&A



### 高效索引加速数据检索







建立索引能够使检索器高效定位和检索相关信息,本部分我们将会学习:

- ·如何使用LazyLLM的索引组件 IndexBase 编写并应用字典树索引
- ·了解高性能向量数据库Milvus的原生使用方式
- ·学习如何在LazyLLM中使用Milvus向量数据库的向量索引





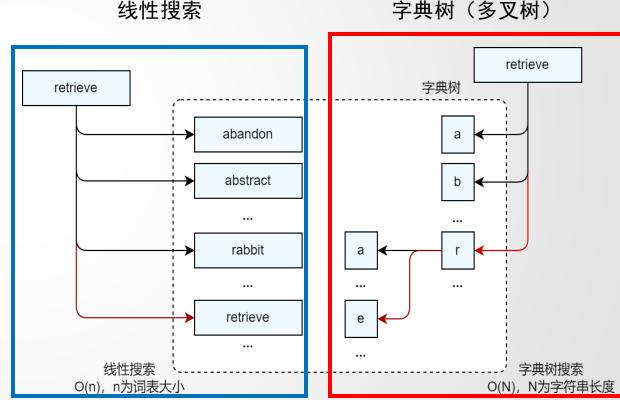


LazyLLM索引组件的基类,所有索引均使用IndexBase组件创建,随后索引组件可被注册在Document当中,并在后 续的检索过程中使用。

#### IndexBase需要实现的三个方法:

- · update:接收节点列表,为节点建立/更新索引。
- · remove:接收节点id和节点组名,在索引中删除 对应节点。
- ·query:根据接收的查询以及相关参数,在索引中 查询符合要求的节点。

注意: 使用默认存储 (MapStore) 的情况下, 节点 的更新与删除会自动关联Index的更新与删除。









```
class TrieNode:
       def init (self):
           self.children: Dict[str, TrieNode] = {}
           self.is end of word: bool = False
           self.uids: set[str] = set()
  class TrieTreeIndex(IndexBase):
       def init (self, store: 'StoreBase'):
           self.store = store
           self.root = TrieNode()
           self.uid to word: Dict[str, str] = {}
10.
       @override
11.
12.
       def update(self, nodes: List['DocNode']) -> None:
           if not nodes or nodes[0]._group != 'words':
13.
14.
               return
15
           for n in nodes:
16.
               uid = n. uid
17
               word = n.text
18
               self.uid to word[uid] = word
               node = self.root
19
20
               for char in word:
21.
                   node = node.children.setdefault(char, TrieNode())
22.
               node.is end of word = True
```

node.uids.add(uid)

23.

1. 定义字典树节点,包含叶节点字典集合、词尾判断、节点id集合

#### 2. 字典树索引类:

- ·继承IndexBase,接收store存储,用于对应节点 查找
- ·初始化方法中创建root节点作为检索入口

#### 3. 索引更新方法

- ·索引仅用于特定节点组检索,故设定如果节点组名不是tree则不建立索引。
- ·对于每个单词,从root节点进入,遍历每个字母:如果当前字母不在当前节点叶子节点集合,则创建叶子结点,key为当前字母。

随后转移至改叶子节点,检查下一字母 直到字母遍历完毕,叶子节点位置即为词尾,标记 并记录id







```
@override
      def remove(self, uids: List[str], group name: Optional[str] = None) -> None:
1.
          if group name != 'words':
              return
                                                                                   1. 实现索引删除方法,从root入口,递归删除对
          for uid in uids:
             word = self.uid to word.pop(uid, None)
                                                                                  应单词索引
             if not word:
6.
                 continue
             self. remove(self.root, word, 0, uid)
8.
9.
      def remove(self, node: TrieNode, word: str, index: int, uid: str) -> bool:
10.
         if index == len(word):
             if uid not in node.uids:
11.
12.
                 return False
                                                                                   2. 递归停止条件: 到达词尾, 删除索引
13.
             node.uids.remove(uid)
             node.is end of word = bool(node.uids)
14.
              return not node.children and not node.uids
15.
16.
          char = word index
          child = node.children.get(char)
17.
18.
         if not child:
19.
             return False
          should delete = self. remove(child, word, index + 1, uid)
                                                                                 3.如果子节点已经可以删了,就把这个 char 对应
20.
          if should delete:
21.
                                                                                  的 child 从 children 里删掉。再判断当前节点是否
             del node.children[char]
22.
             return not node.children and not node.uids
23.
                                                                                  也可以删(没有子节点且没有其他 uid)。
24.
          return False
```







```
1. @override
2. def query(self, query: str, group_name: str, **kwargs) -> List[str]:
3. node = self.root # 查询时以根节点为入口节点
4. for char in query: # 顺序遍历query的每个字符
5. node = node.children.get(char) # 若存在字符对应的叶子节点,则转移至叶子结点,直到完成遍历
6. if node is None:
7. return []
8. return self.store.get_nodes(group_name=group_name, uids=list(node.uids)) if node.is_end_of_word else []
```

#### 将自定义索引成功注册并使用:

```
1. docs = lazyllm.Document(dataset_path=DOC_PATH, embed=embedding_model)
2. docs.create_node_group(name'words', transform=(lambda d: d.split('\r\n'))) # 创建节点组
3. docs.register_index("trie_tree", TrieTreeIndex, docs1.get_store()) # 注册索引
4. retriever = lazyllm.Retriever(docs1, group_name="words", index="trie_tree", topk=1) # 创建检索器,指定索引类型
5. retriever.start() # 检索器初始化
```

注意:在注册自定义索引时,请关注使用的存储<mark>是否支持hook机制</mark>,以自动同步索引。例如 Milvus 这种专门做向量检索的数

据库,本身架构不允许外部代码主动注册额外索引,故只能使用预定义索引。







#### 为了方便性能对比,同样定义线性搜索索引。

```
1. class LinearSearchIndex(IndexBase):
       def init (self):
          self.nodes = []
      @override
      def update(self, nodes: List['DocNode']) -> None:
          if not nodes or nodes[0]. group != 'words':
8.
              return
          for n in nodes:
10.
               self.nodes.append(n)
                                                                                     # 顺序加入每个单词
11.
      @override
12.
       def remove(self, uids: List[str], group name: Optional[str] = None) -> None:
          if group name != 'words':
13.
14.
              return
15.
          for uid in uids:
16.
              for i, n in enumerate(self.nodes):
                  if n. uid == uid:
17.
                      del self.nodes[i]
18.
                      break
19.
20.
      @override
21.
       def query(self, query: str, **kwargs) -> List[str]:
22.
          res = []
23.
          for n in self.nodes:
                                                                                     # 线性搜索: 顺序遍历每个单词
24.
              if n.text == query:
25.
                  res.append(n)
                                                                                     # 假设每个单词出现一次, 精准匹配
26.
                  break
27.
          return res
```







我们准备一个庞大的词表,其中包含37w个英文字符串(大多数为英文单词)。

#### 对比使用线性搜索与字典树索引的检索耗时:

query: a, trie time: 7.152557373046875e-05, trie res: a

query: lazyllm, linear time: 0.0384678840637207, linear res: lazyllm query: lazyllm, trie time: 0.00010895729064941406, trie res: lazyllm

query: zwitterionic, linear time: 0.08076953887939453, linear res: zwitterionic query: zwitterionic, trie time: 0.00010061264038085938, trie res: zwitterionic

```
1. docs = lazyllm.Document(dataset path=DOC PATH, embed=embedding model)
 2. docs.create node group(name='words', transform=(lambda d: d.split('\r\n')))
                                                                                          # 创建节点组
 3. docs.register index("trie tree", TrieTreeIndex, docs.get store())
                                                                                          # 注册索引
 4. docs.register index("linear search", LinearSearchIndex)
 5. retriever1 = lazyllm.Retriever(docs, group name="words", index="linear search", topk=1) # 创建检索器,指定索引类型
 6. retriever2 = lazyllm.Retriever(docs, group name="words", index="trie tree", topk=1)
 7. retriever1.start()
                                                                                          # 检索器初始化
 8. retriever2.start()
 9. for query in ["a", "lazyllm", "zwitterionic"]:
                                                                                          # 分别检索词表中靠前、中间、靠后单词
       st = time.time()
 10.
       res = retriever1(query)
 11.
       et = time.time()
 12.
       LOG.info(f"query: {query}, linear time: {et - st}, linear res: {res[0].text}")
 13.
       st = time.time()
 14.
       res = retriever2(query)
 15.
       et = time.time()
 16.
        LOG.info(f"query: {query}, trie time: {et - st}, trie res: {res[0].text}")
 17.
query: a, linear time: 6.175041198730469e-05, linear res: a
```



线性搜索:搜索时间随单词位置增加

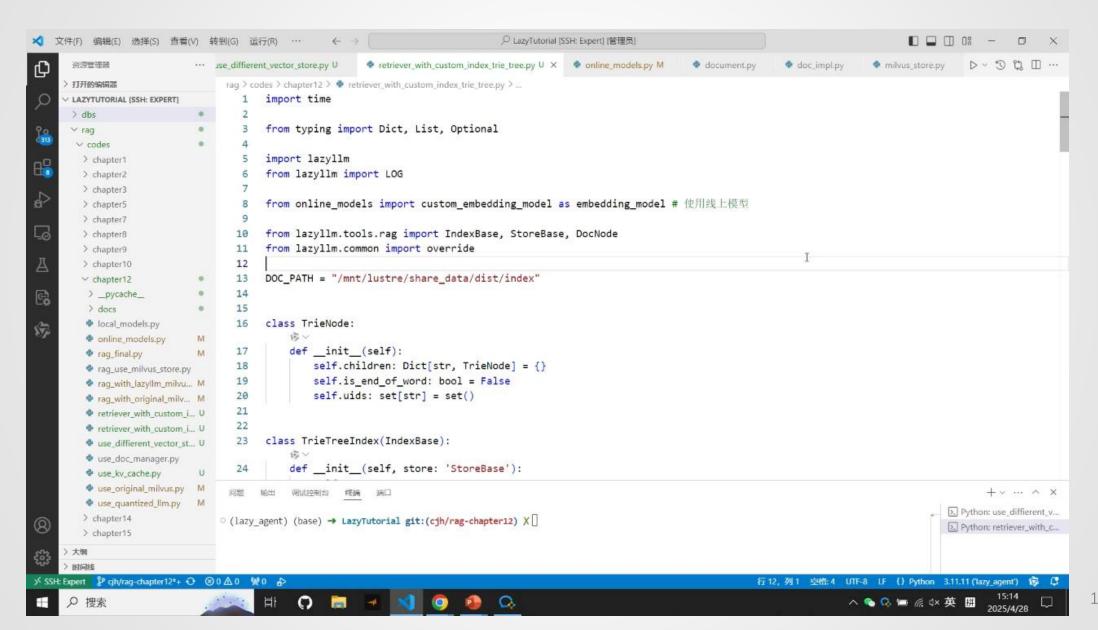
字典树:搜索效率高且稳定













### 使用Milvus内置向量索引







向量索引是向量数据库的核心能力之一,实际生产中往往使用向量数据库(如Milvus)实现高效向量检索。

我们首先了解一下本次的主角——Milvus的原生基础使用方法,包含Milvus的安装以及基本功能的使用(初始化、数据注入、数据检索)。

在本次实践中使用 Milvus Lite,它是pymilvus中包含的一个 python 库,可以嵌入到应用程序中。Milvus 还支持在Docker和Kubernetes上部署,适用于生产用例。

安装命令: pip install -U pymilvus

	索引方式	度量类型
浮点索引	FLAT IVF_FLAT,IVF_PQ,IVF_SQ8 HNSW, HNSW_PQ,SCANN 	欧氏距离 L2 内积 IP 余弦相似度 COSINE
二进制索引	BIN_FLAT、BIN_IVF_FLAT	Jaccard (JACCARD) 哈明 (HAMMING)
稀疏索引	SPARSE_INVERTED_INDEX、 SPARSE_WAND	内积 IP





相当于SQL数据库中建表





#### 设置向量数据库(实例化客户端以创建本地数据库):

- 1. from pymilvus import MilvusClient
- # 创建milvus客户端, 传入本地数据库的存储路径, 若路径不存在则创建
- 2. client = MilvusClient("dbs/origin\_milvus.db")

#### 创建Collections, 存储向量及相关元数据:

- 1. if client.has collection(collection name="demo collection"): # 如果已存在同名collection,则先删除
- client.drop collection(collection name="demo collection")
- 3. client.create collection(
- 4. collection name="demo collection",
- 5. dimension=1024,

6.)

# 设置存储向量维度

注: create\_collection支持更多的入参,如主键列名primary\_field\_name、向量列名vector\_field\_name、向量索引参数index\_params、相似度参数metric\_type等,此处仅介绍基本用法,其余均使用默认参数。











#### 数据入库:

```
    docs = [
    "Artificial intelligence was founded as an academic discipline in 1956.",
    "Alan Turing was the first person to conduct substantial research in Al.",
    "Born in Maida Vale, London, Turing was raised in southern England.",
    ]
    vecs = [embedding_model(doc) for doc in docs]
    data = [
    { "id" : i, "vector" : vecs[i], "text" : docs[i], "subject" : "history" }
    for i in range(len(vecs))
    ]
    11.res = client.insert(collection_name = "demo_collection" , data=data)
```

#### 向量检索:

```
1. query = "Who is Alan Turing?"
2. q_vec = embedding_model(query) # query向量化
3. res = client.search(
4. collection_name="demo_collection", # 指定collection
5. data=[q_vec],
6. limit=2, # 指定检索数量(top_k)
7. output_fields=["text", "subject"], # 指定检索结果中包含的字段
8.)
```

# 测试文段列表

- # 文段向量化
- # 将数据转换为入库所需格式
- # subject为自定义label,
- # insert方法数据入库









#### 标量过滤检索:

```
1. docs2 = [
    "Machine learning has been used for drug design.",
    "Computational synthesis with AI algorithms predicts molecular properties.",
    "DDR1 is involved in cancers and fibrosis.",
2. vecs2 = [embedding model(doc) for doc in docs2]
3. data2 = [
    {"id": 3 + i, "vector": vecs2[i], "text": docs2[i], "subject": "biology"}
    for i in range(len(vecs2))
4. client.insert(collection name= "demo collection", data=data2)
5. res = client.search(
    collection name="demo collection",
    data=[embedding model("tell me Al related information")],
    filter= "subject == 'biology' ",
    limit=2,
    output fields=["text", "subject"],
```

# 另一组文段

#数据插入同一个collection中

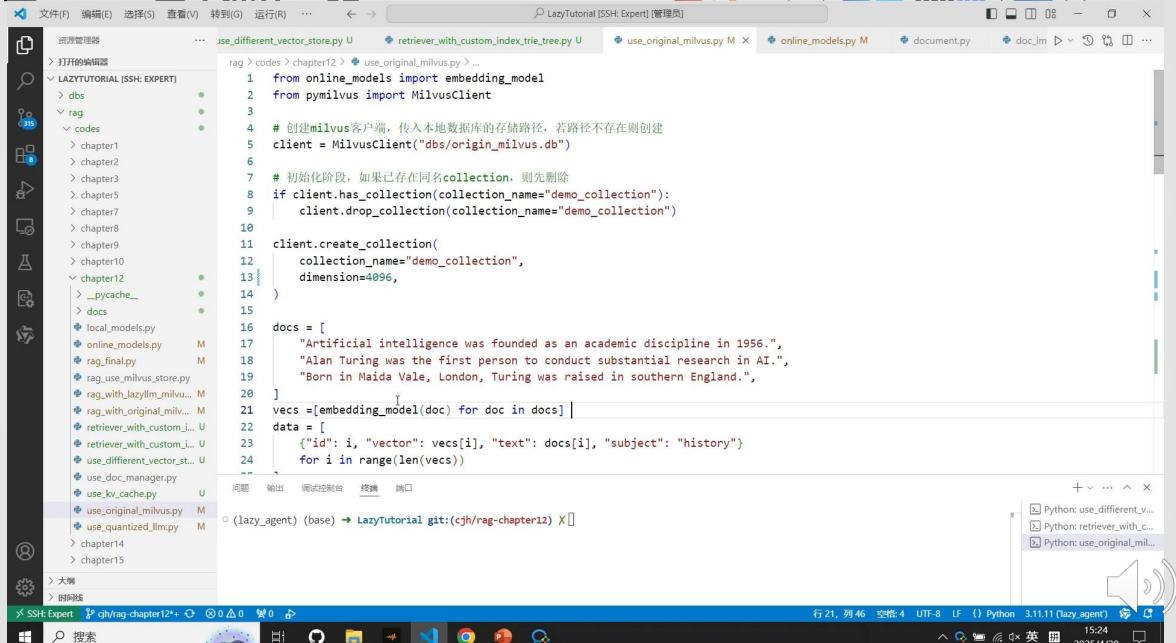
# 检索时使用filter实现标量过滤











### 使用原生Milvus搭建RAG









- 1. import lazyllm
- 2. from lazyllm.tools.rag import SimpleDirectoryReader, SentenceSplitter
- 3. from lazyllm.tools.rag.doc\_node import MetadataMode
- 4. from pymilvus import MilvusClient
- 5. from online models import embedding model, llm, rerank model

# Ilm、embedding、reranker均使用线上模型

```
6. DOC_PATH = "/xxx/xxx"
```

```
7. client = MilvusClient("dbs/rag_milvus.db")
```

# milvus client初始化

```
8. if client.has_collection(collection_name="demo_collection"):
```

```
9. client.drop collection(collection name="demo collection")
```

```
    client.create_collection(
        collection_name="demo_collection",
        dimension=1024,
```

11. docs = SimpleDirectoryReader(input dir=DOC PATH)()

#使用reader加载文件目录

#数据入库

12. block\_transform = SentenceSplitter(chunk\_size=256, chunk\_overlap=25) # 定义分块策略

13. nodes = **□** 

14. for doc in docs:

nodes.extend(block transform(doc))

```
16. vecs = [embedding_model(node.get_text(MetadataMode.EMBED)) for node in nodes] # 切片向量化
```

# 创建collection



### 使用原生Milvus搭建RAG

27. context\_str = "\n----\n".join(rerank\_contexts)
28. res = llm({ "query" : query, "context str" : context str})

29. print(res)









```
# 检索与生成
19. prompt = '你是一个友好的 AI 问答助手,你需要根据给定的上下文和问题提供答案。\
        根据以下资料回答问题: \
        {context str} \n
20. llm.prompt(lazyllm.ChatPrompter(instruction=prompt, extro keys=[ 'context str' ])) # 初始化prompt
21. query = "证券管理的基本规范?"
22. q vec = embedding model(query)
                                                                              # query向量化
23. res = client.search(
                                                                              # 向量检索
    collection name="demo collection",
    data=[q vec],
    limit=12,
    output fields=["text"],
24. contexts = [res[0][i].get( 'entity' ).get( "text" , "" ) for i in range(len(res[0]))]
                                                                              # 提取检索结果
25. rerank res = rerank model(text=query, documents=contexts, top n=3)
                                                                              # 重排序
26. rerank contexts = [contexts[res[0]] for res in rerank res]
```

# 召回信息拼接,大模型生成



### 使用原生Milvus搭建RAG

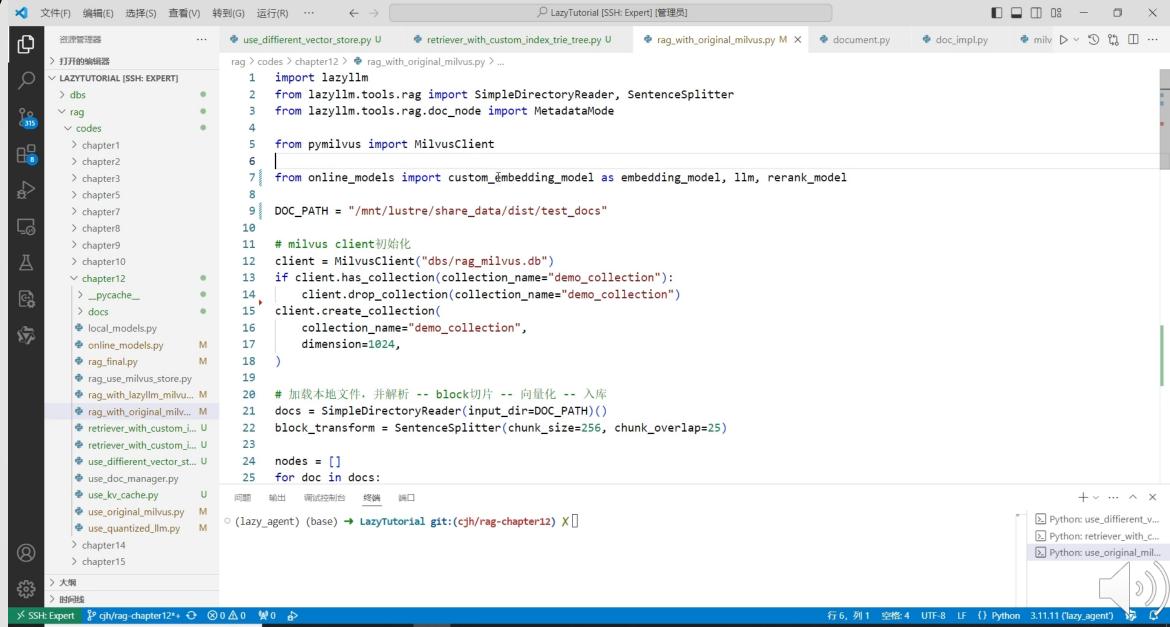
₽ 搜索







へ 🐼 🖆 🦟 🗘 英 拼



### Milvus in LazyLLM







从头学习如何使用Milvus向量数据库成本略大,使用LazyLLM一键接入Milvus。

配置方式: store\_conf中额外添加字段 indices, SmartEmbeddingIndex快速实现Milvus向量索引。

注:如果存储类型type=milvus的情况下,由于kwargs中已经包含了索引配置信息,因此无需额外配置indices,即可实现内置的向量检索。



### Milvus in LazyLLM

11. res = Ilm({ "query" : query, "context str" : context str})

12. print(res)



# 大模型生成







```
1. docs = lazyllm.Document(dataset path=DOC PATH, embed=embedding model, store conf=store conf)
2. retriever = lazyllm.Retriever(
    docs,
    group name="MediumChunk",
    topk=6,
    index='smart embedding index'
                                                                           #使用smart embedding index索引,调用milvus向量检索
retriever.start()
4. reranker = Reranker( 'ModuleReranker', model=rerank model, topk=3)
                                                                           # 定义reranker
5. prompt = '你是一个友好的 AI 问答助手,你需要根据给定的上下文和问题提供答案。\
      根据以下资料回答问题:\
      {context str} \n
6. llm.prompt(lazyllm.ChatPrompter(instruction=prompt, extro keys=['context str']))
7. query = "证券管理有哪些准则?"
8. nodes = retriever(query=query)
                                                                           # query检索
9. rerank nodes = reranker(nodes, query)
                                                                           # 节点重排
10. context str = "\n======\n" .join([node.get content() for node in rerank nodes]) # 召回节点拼接
```

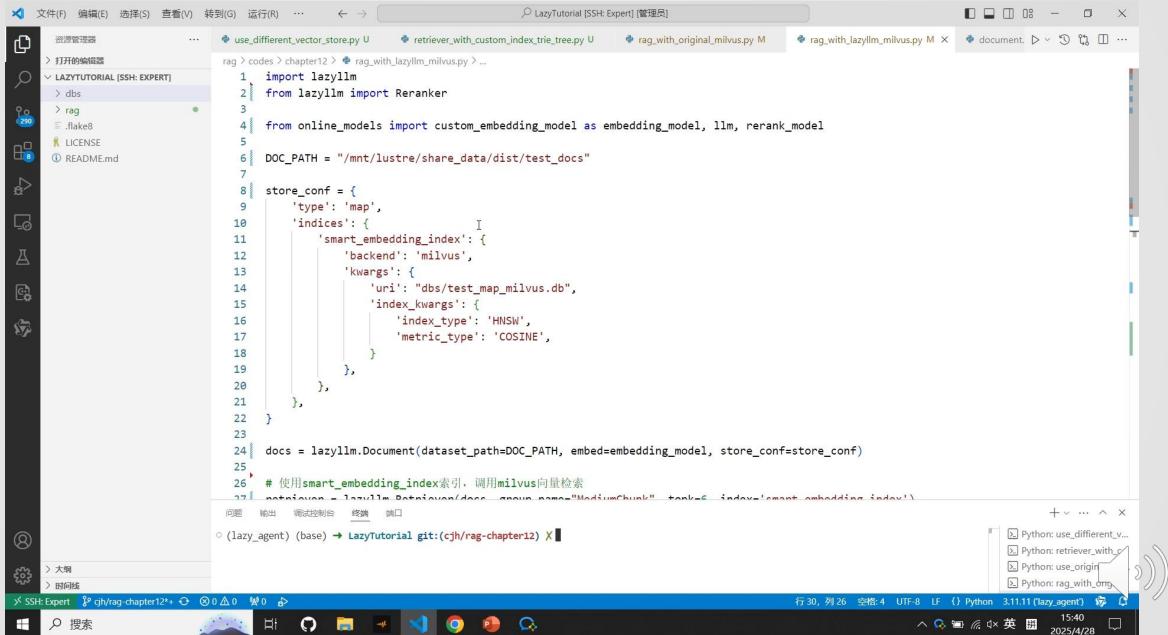


### Milvus in LazyLLM









### Milvus的分布式使用方式







Milvus不仅支持本地db文件链接,同时支持远程服务器端点的接入。

使用 Docker 完成服务部署(Milvus Standalone / Distributed)

直接使用云服务厂商(如Zilliz Cloud)提供的Milvus云服务



client = MilvusClient("dbs/origin\_milvus.db")



获取远端地址与身份验证



```
client = MilvusClient(
   uri="http://localhost:19530",
   token="username:password "
)
```



### Milvus的分布式使用方式







例如,我们期望在Linux系统上部署Milvus Standalone服务。

- 安装Docker, 并根据 Milvus 官方文档检查硬件与软件要求。
- · 在Docker中安装Milvus,安装脚本:

#### #下载安装脚本

curl -sfL https://raw.githubusercontent.com/milvus-io/milvus/master/scripts/standalone\_embed.sh -o standalone\_embed.sh

# 启动Docker容器 bash standalone embed.sh start

#### 安装结果:

- 名为 Milvus 的 docker 容器在19530端口启动。要更改 Milvus 的默认配置,需将设置添加到当前文件夹中的user.yaml文件,然后重新启动服务。
- Milvus 数据卷被映射到当前文件夹中的volumes/milvus。



### 工程优化: 使用Flow实现高效RAG管道







定义document 定义retriever1 定义retriever2 定义reranker

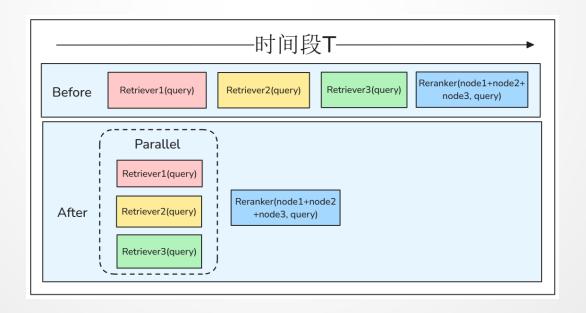
..

nodes1 = retriever1(query=query)
nodes2 = retriever2(query=query)
rerank\_nodes = reranker(nodes1 + nodes2, query)



Pipeline实现RAG系统整体数据管道

Parallel实现多路召回并行执行





### 使用Flow实现高效RAG管道







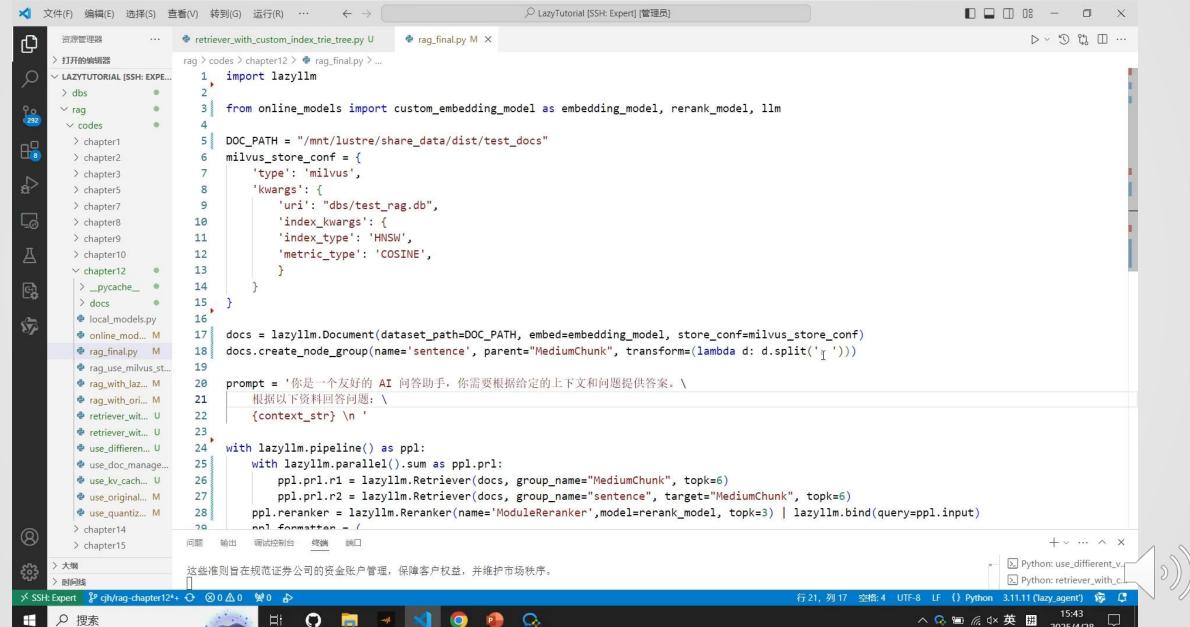
```
1. docs = lazyllm.Document(dataset path=DOC PATH, embed=embedding model, store conf=milvus store conf)
2. docs.create node group(name='sentence', parent="MediumChunk", transform=(lambda d: d.split('。')))
3. prompt = '你是一个友好的 AI 问答助手,你需要根据给定的上下文和问题提供答案。\
    根据以下资料回答问题: \
    {context str} \n '
5. with lazyllm.pipeline() as ppl:
    with lazyllm.parallel().sum as ppl.prl:
       ppl.prl.r1 = lazyllm.Retriever(docs, group name="MediumChunk", topk=6)
       ppl.prl.r2 = lazyllm.Retriever(docs, group name="sentence", target="MediumChunk", topk=6)
    ppl.reranker = lazyllm.Reranker(name='ModuleReranker',model=rerank model, topk=3) | lazyllm.bind(query=ppl.input)
     ppl.formatter = (
      lambda nodes, query: dict(
        context str="\n".join(node.get content() for node in nodes),
        query=query)
    ) | lazyllm.bind(query=ppl.input)
     ppl.llm = llm.prompt(lazyllm.ChatPrompter(instruction=prompt, extro keys=['context str']))
12. w = lazyllm.WebModule(ppl, port=23492, stream=True).start().wait()
```

### 使用Flow实现高效RAG管道



















- 1. 上节回顾
- 2. 使用向量数据库实现知识库持久化存储
- 3. 构建高效索引提升检索性能
- 4. 优化模型推理性能
- 5. Q&A

### 使用不同推理框架部署大模型







LazyLLM支持使用LightLLM、vLLM、LMDeploy三种推理框架实现模型部署,开发者可以根据自己的实际需求, 灵活选择推理框架,以实现最适合自己的模型推理体验。

```
灵活选择推理框架,以实现最适合自己的模型推理体验。
Ilm = lazyllm.TrainableModule('internlm2-chat-20b')
    Ilm = lazyllm.TrainableModule( 'internlm2-chat-20b' ).deploy_method(lazyllm.deploy.Lightllm)
    Ilm = lazyllm.TrainableModule( 'internlm2-chat-20b' ).deploy_method(lazyllm.deploy.Vllm)
    Ilm = lazyllm.TrainableModule( 'internlm2-chat-20b' ).deploy method(lazyllm.deploy.LMDeploy)
Ilm = TrainableModule('DeepSeek-R1').deploy method(
 (deploy.VIIm, {
   'tensor-parallel-size': 8,
                                                                          传入推理框架所需入参
    'pipeline-parallel-size': 2,
   'max-num-batched-tokens': 131072,
   'launcher': launchers.sco(nnode=2, ngpus=8, sync=False)
```



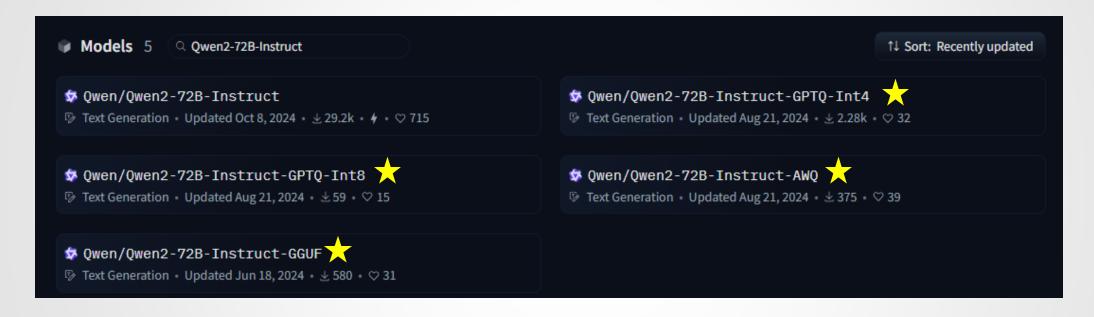
### 使用量化模型







#### 从开源社区寻找大模型的量化版本



运行BF16或FP16的 Qwen2-72B-Instruct 模型运行需要至少144GB显存 (例如2xA100-80G或5xV100-32G) 而运行 Int4 模型至少需要 48GB 显存 (例如1xA100-80G或2xV100-32G) , 是原来的33%



### 使用量化模型

1. import time

17. Ilm awq(query)

18. end time = time.time()

19. print("AWQ量化模型耗时: ", end time-start time)



# 量化模型生成





```
2. from lazyllm import TrainableModule, deploy
3. start time = time.time()
4. Ilm = TrainableModule( 'Qwen2-72B-Instruct' ).deploy method(
                                                                                  # 启动原始模型
       (deploy.VIIm, {
         'tensor-parallel-size' : 2,
       })).start()
5. end time = time.time()
6. print("原始模型加载耗时: ", end time-start time)
7. start time = time.time()
8. Ilm_awq = TrainableModule( 'Qwen2-72B-Instruct-AWQ' ).deploy method(
                                                                                  # 启动量化模型
       (deploy.VIIm, {
         'tensor-parallel-size' : 2,
       })).start()
9. end time = time.time()
10. print("AWQ量化模型加载耗时: ", end time-start time)
11. query = "生成一份1000字的人工智能发展相关报告"
                                                                                  # 测试query
12. start time = time.time()
13. Ilm(query)
                                                                                  #原始模型生成
14. end time = time.time()
15. print("原始模型耗时: ", end time-start time)
16. start time = time.time()
```

原始模型加载耗时: 129.6051540374756 AWQ量化模型加载耗时: 86.4980857372284

原始模型耗时: 13.104065895080566 AWQ量化模型耗时: 8.81701111793518











# Q&A

1. 面对性能测试中的不稳定性,有什么处理方法么









# 感谢聆听 **Thanks for Listening**