

# 基于Sentence-BERT的语义相似度计算技术方案

## 0. 摘要

本方案旨在构建一个基于语义相似度的智能问答（QA）系统。核心思路是：离线时，采用**Sentence-BERT (SBERT)**模型将知识库中的**FAQ（常见问题）文本转化为高质量的句子向量，并存入FAISS/Milvus**向量数据库构建索引；在线查询时，将用户问题同样编码为向量，在数据库中进行高效的相似度检索，找出最匹配的标准问题，并返回其对应答案。此方案通过结合先进的深度学习模型与向量检索技术，可实现快速、精准的语义匹配，有效提升问答系统的智能化水平。

## 1. 目标

构建一个基于语义相似度的智能问答系统。当用户提出问题时，系统能快速、准确地从现有的**FAQ**知识库中，找到与用户问题在语义上最相似的标准问法，并返回其对应的答案。

## 2. 核心技术选型

### • 文本表示模型: **Sentence-BERT (SBERT)**

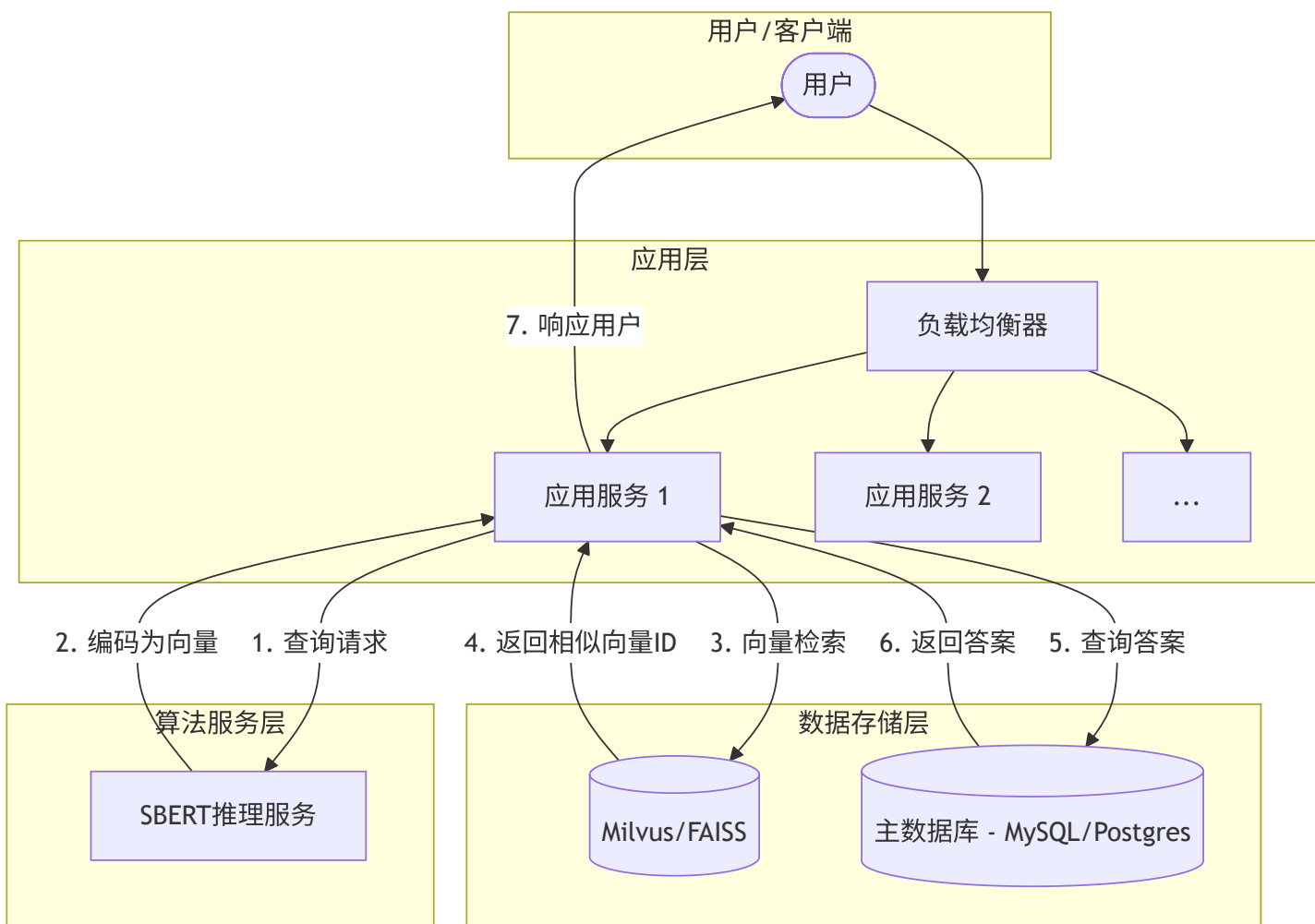
- **选型理由:** 传统的**BERT**模型输出的词向量或 [CLS] 向量未针对句子相似度任务进行优化，直接计算其相似度效果不佳。**SBERT**基于孪生网络架构在**BERT**上进行微调，专门用于生成语义信息丰富的句子级别向量。通过计算**SBERT**向量的**余弦相似度**，可以精准度量句子间的语义接近度，是本场景的最佳选择。
- **推荐模型:**
  - **高精度模型:** paraphrase-multilingual-mpnet-base-v2 - 效果最好，但推理速度稍慢。
  - **均衡模型:** paraphrase-multilingual-MiniLM-L12-v2 - 速度和效果的优秀平衡。
  - **中文优化模型:** shibing624/text2vec-base-chinese - 针对中文语料训练，表现优异。

### • 向量存储与检索引擎: **FAISS** 或 **Milvus**

- **选型理由:** 当知识库规模达到万级以上时，线性遍历计算相似度将产生不可接受的延迟。向量数据库通过构建高效索引（如**IVF-FLAT**, **HNSW**），能够将海量向量的相似度搜索时间缩短至毫秒级。
- **对比:**
  - **FAISS:** Facebook AI出品的相似度检索库，性能极致，适合嵌入到现有服务中。需要自行管理向量数据和ID映射。

- **Milvus**: 一个完整的分布式向量数据库服务，提供数据管理、索引构建、监控等全套功能，更适合大规模、企业级部署。
- **相似度度量: 余弦相似度 (Cosine Similarity)**
  - 通过测量向量间的夹角来判断方向上的一致性，忽略向量长度差异，是度量文本语义相似度的标准方法。

### 3. 部署架构



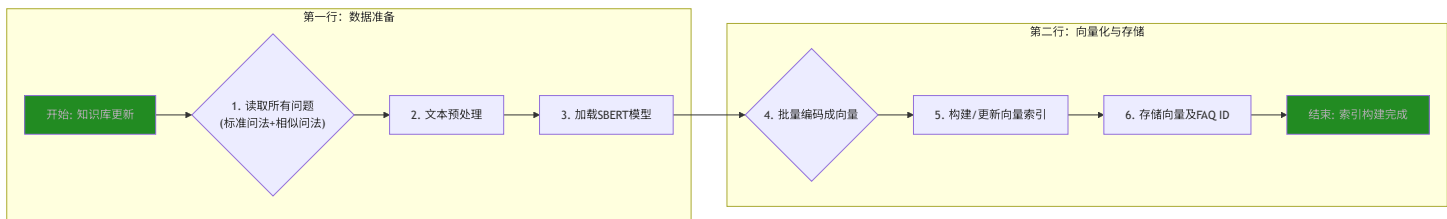
### 4. 技术方案流程

整个系统分为 离线索引构建 和 在线实时查询 两个阶段。

#### 4.1 离线索引构建 (Offline Indexing)

此阶段的目标是将知识库中的所有问题都转换成向量并存入向量数据库，构建起可供快速查询的索引。此过程应在后台异步执行，每当知识库内容更新时触发。

流程图:



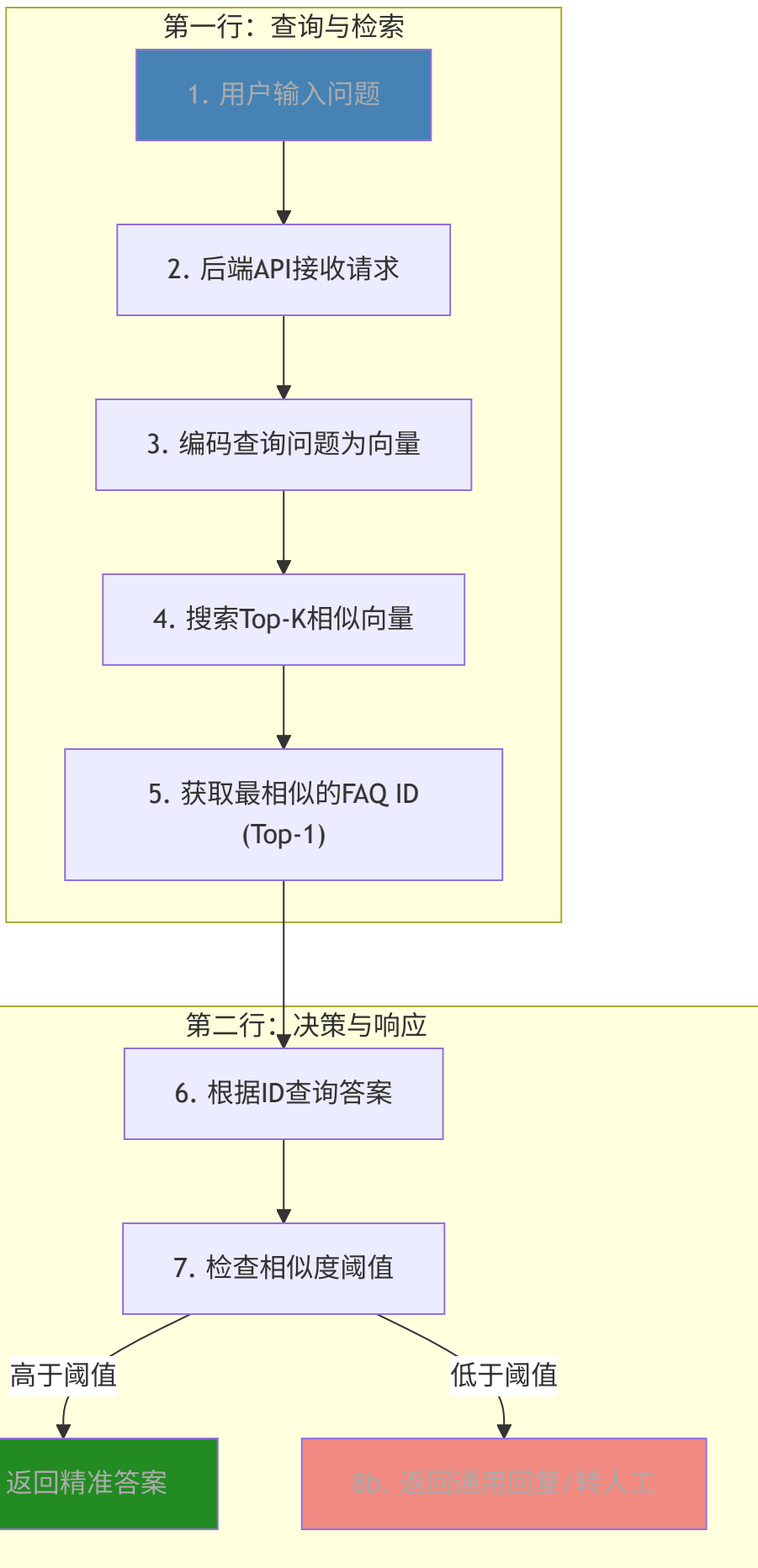
步骤详解:

- 1. **读取数据:** 从主数据库（如MySQL）中，提取所有FAQ的 标准问法 及其关联的 相似问法 ，形成一个完整的问题列表。
- 2. **文本预处理:** 对文本进行必要的清洗，如统一大小写、去除HTML标签等（此步骤可根据SBERT模型的预处理能力简化）。
- 3. **加载模型:** 在后台的算法服务中加载预训练好的SBERT模型。
- 4. **批量编码:** 将问题列表分批（Batching）送入SBERT模型，高效地生成每个问题对应的句子向量。
- 5. **构建/更新索引:** 将生成的向量批量添加（add/insert）到FAISS或Milvus中。向量数据库会自动为这些向量建立或更新索引。
- 6. **存储映射关系:** 确保每个向量都与其原始的 faq\_id 一一对应，以便检索后能找到原始条目。Milvus可直接存储，FAISS则需要外部自行维护此映射。

## 4.2 在线实时查询 (Online Querying)

此阶段处理用户的实时请求，要求在毫秒级内完成从提问到回答的全过程。

流程图:



步骤详解:

1. **接收请求:** 后端API接收到来自用户的查询字符串。
2. **查询编码:** 调用算法服务, 使用与离线阶段**完全相同**的SBERT模型, 将查询字符串编码成一个向量。
3. **向量检索:** 使用该查询向量, 在FAISS或Milvus中执行搜索 (**search**) 操作, 请求返回最相似的**Top-K**个结果 (例如**K=3**)。
4. **获取ID与相似度:** 检索引擎返回**K**个最相似向量的**ID**及其对应的相似度得分。
5. **决策与查询:**
  - 取相似度最高的**Top-1**结果。
  - **进行阈值判断** (见下一节)。如果相似度高于预设阈值 (如**0.85**), 则认为匹配成功。
  - 使用**Top-1**结果的 `faq_id`, 去主数据库中查询对应的标准答案。
6. **返回结果:** 将查询到的答案返回给用户。如果相似度低于阈值, 则返回一个预设的通用回复 (如 抱歉, 未能找到您问题的答案, 请尝试换个问法或联系人工客服 )。

## 5. 优化与考量

- **相似度阈值调优:** 阈值是平衡系统**召回率**与**精确率**的关键。
  - **阈值过高:** 可能导致很多相关的问题被拒绝 (低召回率), 但返回的答案非常精准 (高精确率)。
  - **阈值过低:** 会让系统回复更多问题 (高召回率), 但可能包含错误答案 (低精确率)。
  - **调优方法:** 准备一个标注好的测试集, 通过测试不同阈值下的表现, 选择一个在业务上最能接受的平衡点。
- **低相似度处理 (Fallback):** 当所有候选问题的相似度都低于阈值时, 应设计合理的兜底策略, 如引导用户、转接人工客服等, 避免无效回答。
- **Re-ranking (重排序):** 对于要求极高准确率场景, 可以采用两阶段检索策略:
  - i. **召回 (Recall):** 使用**SBERT+FAISS**快速从海量数据中召回**Top-K** (如**K=20**) 个候选。
  - ii. **排序 (Rank):** 使用一个更强大但更慢的模型 (如**Cross-Encoder**) 对这**K**个候选进行重新打分排序, 选择最终的**Top-1**。这能显著提升精度, 但会增加延迟。

## 6. 总结

本方案通过**SBERT**将非结构化文本问题转化为结构化的语义向量, 再利用高效的向量检索引擎, 构建了一个响应迅速、匹配精准的现代化智能问答系统。该架构具有良好的扩展性, 能够有效代替部分人工客服工作, 并为后续引入更复杂的**NLP**功能 (如对话管理、意图识别) 打下坚实基础。