

src_explain

src_explain.md - LLM-RAG 多輪對話框架程式碼詳解

整體架構概覽：智能客服中心

想像這個系統就像一個**高級智能客服中心**，而不是傳統的機械式電話客服。這個客服中心有以下特點：

1. **記憶力超強** - 記得你們之前聊過什麼
2. **會讀心術** - 能理解你真正想問什麼
3. **有百科全書** - 可以查詢各種知識
4. **會追問** - 如果不清楚你的需求，會聰明地問你更多問題
5. **會學習** - 從每次對話中變得更聰明

核心角色介紹

1. 狀態定義類 (Enum Classes)

```
class IntentClarityLevel(Enum):  
    HIGH = "high"    # >80分 - 像是顧客說「我要買一台 iPhone 15 Pro 256GB 黑色」  
    MEDIUM = "medium" # 50-80分 - 像是顧客說「我想買手機」  
    LOW = "low"       # <50分 - 像是顧客說「我想要那個」
```

比喻說明：

- 這就像判斷顧客表達的清晰程度
- **HIGH (高)**：顧客講得很清楚，客服知道該怎麼做
- **MEDIUM (中)**：大概知道方向，但需要問幾個問題
- **LOW (低)**：完全不知道顧客要什麼，需要重新開始

```

class ConversationState(Enum):
    INIT = "initialization"      # 剛接電話，說「您好」
    INTENT_ANALYSIS = "intent_analysis" # 聽顧客說話，理解需求
    CLARIFICATION = "clarification" # 詢問更多細節
    RETRIEVAL = "retrieval"      # 去倉庫找資料
    RESPONSE_GENERATION = "response_generation" # 組織回答
    FEEDBACK_COLLECTION = "feedback_collection" # 詢問滿意度
    COMPLETED = "completed"    # 通話結束

```

比喻說明：

- 這就像客服處理電話的不同階段
- 每個階段都有特定的任務和目標

2. 資料結構類 (Data Classes)

Message 類 - 對話記錄本

```

@dataclass
class Message:
    role: str      # 誰在說話："user" (顧客) 或 "assistant" (客服)
    content: str   # 說了什麼
    timestamp: datetime # 什麼時候說的
    metadata: Dict # 額外筆記 (如情緒、重要性等)

```

比喻：就像客服的通話記錄，每句話都要記下來，包括誰說的、說了什麼、什麼時候說的。

IntentAnalysisResult 類 - 理解報告單

```

@dataclass
class IntentAnalysisResult:
    clarity_score: float # 清晰度分數 (0-100分)
    main_intent: str # 主要目的是什麼
    key_entities: List[str] # 提到了哪些關鍵東西
    missing_info: List[str] # 還缺什麼資訊
    ambiguities: List[str] # 哪些地方不清楚
    suggested_strategy: str # 建議怎麼處理
    confidence_level: IntentClarityLevel # 信心等級

```

比喻：這就像客服聽完顧客第一句話後，在心裡做的分析筆記。

3. 核心功能模組

IntentAnalyzer 類 - 心理分析師

```

class IntentAnalyzer:
    def analyze(self, conversation: List[Message]) -> IntentAnalysisResult:
        # 分析對話，理解顧客想要什麼

```

運作原理：

1. 聽 - 讀取對話內容
2. 理解 - 找出關鍵詞、判斷意圖
3. 評分 - 給清晰度打分
4. 診斷 - 找出缺失的資訊和模糊的地方

打分機制詳解：

```

def _calculate_clarity_score(self, content: str) -> float:

    score = 50.0 # 基礎分，像考試的基本分

    # 關鍵詞覆蓋率 (30分)
    # 如果說了「想要」「需要」這些明確的詞，加分
    keywords = ["想要", "需要", "請問", "如何", ...]

    # 語句長度 (20分)
    # 太短說不清楚，太長可能廢話多

    # 具體性檢查 (20分)
    # 有沒有舉例子，有沒有說具體

```

KnowledgeRetriever 類 - 知識庫管理員

```

class KnowledgeRetriever:

    def retrieve(self, query: str, entities: List[str], top_k: int = 5):

        # 去知識倉庫找相關資料

```

比喻：

- 就像圖書館管理員，根據你的需求去找相關的書
- 會根據關鍵詞的匹配程度，把最相關的資料先拿出來
- `top_k=5` 表示最多拿5本最相關的書

ResponseGenerator 類 - 回答組裝師

```

class ResponseGenerator:

    def generate_response(self, query, retrieval_results, context):

        # 把找到的資料組織成完整的回答

    def generate_clarification_questions(self, missing_info, ambiguities):

        # 生成追問的問題

```

運作方式：

1. 整理資料 - 把找到的資料排序整理

2. 組織語言 - 用友善的方式表達
3. 加入建議 - 提供延伸閱讀或相關問題
4. 標註來源 - 告訴用戶資訊從哪來的

4. 主框架類 - 總指揮中心

```
class MultiRoundDialogueFramework:  
    def __init__(self):  
        self.intent_analyzer = IntentAnalyzer()    # 心理分析師  
        self.knowledge_retriever = KnowledgeRetriever() # 知識管理員  
        self.response_generator = ResponseGenerator() # 回答組裝師  
        self.sessions = {} # 所有進行中的對話
```



對話流程詳解

第一步：創建會話

```
def create_session(self, session_id: str) -> ConversationContext:  
    # 像是客服接起電話，開始新的服務
```

1. 創建一個新的對話記錄本
2. 寫下開場白（歡迎詞）
3. 準備好接收顧客的第一句話

第二步：處理用戶輸入

```
def process_user_input(self, session_id: str, user_input: str) -> str:  
    # 處理顧客說的每一句話
```

這個函數就像是客服的大腦，決定該如何回應：

1. 找到對話記錄 - 看看這是新顧客還是回頭客
2. 記錄顧客說的話 - 寫進記錄本

3. 判斷該怎麼做 - 根據當前狀態決定下一步

第三步：智能分流處理

情況A：高置信度路徑（顧客說得很清楚）

```
def _handle_high_confidence_path(self, context):  
    # 直接去找答案
```

流程：

1. 去知識庫查資料
2. 整理成完整回答
3. 加上延伸建議
4. 友善地回覆顧客

實際例子：

- 顧客：「我想知道2024年AI在醫療領域的應用」
- 系統：理解清晰，直接查找相關資料並回答

情況B：需要澄清（顧客說得不夠清楚）

```
def _handle_clarification_path(self, context):  
    # 聰明地問問題
```

流程：

1. 分析缺少什麼資訊
2. 設計3-5個引導問題
3. 友善地詢問顧客
4. 等待顧客補充

實際例子：

- 顧客：「推薦幾本好書」
- 系統：「請問您對哪類書籍感興趣？文學、科技還是商業？」

情況C：完全不理解（需要重新開始）

```
def _handle_low_confidence_path(self, context):  
    # 重新引導對話
```

流程：

1. 道歉並說明情況
2. 提供清晰的引導
3. 給出範例
4. 鼓勵顧客重新表達

關鍵設計理念

1. 狀態機設計

整個系統就像一個**自動販賣機**：

- 每個狀態是一個固定的步驟
- 根據輸入決定下一個狀態
- 確保對話有序進行

2. 評分機制

就像**考試評分**：

- 多個維度綜合評分
- 每個維度有不同權重
- 根據總分決定處理策略

3. 上下文管理

像是**連續劇**：

- 記住之前發生的事
- 理解當前的情節
- 預測接下來的發展

4. 知識檢索

像是搜尋引擎：

- 根據關鍵詞匹配
- 按相關度排序
- 返回最佳結果

擴展性設計

1. 介面設計 (Interface)

```
class LLMInterface(ABC):  
    @abstractmethod  
    def generate(self, prompt: str, **kwargs) -> str:  
        pass
```

比喻：這就像電器的插頭標準，不管是什麼品牌的LLM，只要符合這個介面，都可以接入使用。

2. 模組化設計

每個功能都是獨立的模組，就像樂高積木：

- 可以單獨替換某個模組
- 可以增加新的模組
- 模組之間通過標準介面連接

3. 配置化設計

關鍵參數都可以配置，就像汽車的設定：

- 可以調整評分標準
- 可以改變問題數量
- 可以切換不同的知識庫

實用技巧

1. 漏斗式提問

從寬到窄，就像**醫生問診**：

- 先問大方向：「哪裡不舒服？」
- 再問具體：「是頭痛還是頭暈？」
- 最後確認：「是持續性還是間歇性？」

2. 信心分級

根據理解程度採取不同策略，就像**導航系統**：

- 地址清楚：直接導航
- 地址模糊：先確認區域
- 完全不知：重新輸入

3. 友善互動

始終保持友善專業，就像**五星級酒店服務**：

- 主動打招呼
- 耐心聆聽
- 提供超預期的幫助
- 詢問滿意度

總結

這個系統就像是一個**智能管家**：

1. **聰明** - 能理解你的需求
2. **博學** - 知道很多知識
3. **體貼** - 會主動詢問和建議
4. **可靠** - 記住所有對話內容
5. **進化** - 不斷學習和改進

通過這種設計，我們創造了一個既智能又人性化的對話系統，能夠像真人客服一樣理解需求、提供幫助，甚至做得更好！