

用戶查詢意圖預測與 RAG 應用合作的方法

基於廣泛的研究分析，以下提供一系列用於預測用戶查詢意圖的方法，這些方法能夠與 RAG 應用有效協作：

1. 基於大語言模型的意圖識別方法

LLM 驅動的意圖分類

現代意圖識別系統利用 **大語言模型 (LLMs)** 的強大能力來理解用戶查詢背後的意圖[1][2][3]。這些系統採用：

- **自適應上下文學習 (Adaptive In-Context Learning)**：利用檢索到的相似查詢作為示例來改善分類效果[2][4]
- **思維鏈提示 (Chain-of-Thought Prompting)**：通過結構化推理提高複雜意圖的識別準確性[2][5]
- **混合路由策略**：結合傳統分類器與 LLM 的不確定性路由機制，在準確性與延遲之間取得平衡[2][5]

語義路由技術

語義路由 (Semantic Routing) 是一種基於語義相似性的決策層[6][7]：

- 將用戶查詢轉換為向量嵌入
- 與預定義的意圖表述進行相似性匹配
- 根據語義意義而非關鍵字規則來路由查詢
- 能夠在不依賴慢速 LLM 推理的情況下快速做出決策[6]

2. 檢索增強生成 (RAG) 協作方法

REIC：RAG 增強意圖分類

REIC 框架將 RAG 技術直接應用於意圖分類[8][9]：

- **索引構建**：創建包含（查詢，意圖）對的密集向量索引
- **候選檢索**：使用預訓練的句子變換器模型進行語義檢索
- **概率計算**：通過 LLM 對檢索到的意圖候選進行最終分類

DSRAG：雙流檢索增強生成

針對複雜意圖場景，**DSRAG 框架**結合兩種檢索策略[10][11]：

- **查詢到查詢 (Q2Q)**：快速匹配預構建的查詢模板庫
- **查詢到元數據 (Q2M)**：從元數據中檢索相關意圖並使用 LLM 進行選擇
- **雙流融合**：當 Q2Q 無法找到匹配時，自動轉向 Q2M 策略

混合 RAG 意圖分類

混合 RAG 方法[12][13] 將檢索與生成模型的優勢結合：

- **意圖管理系統**：將意圖、實體和話語嵌入並存儲在向量存儲中
- **認知檢索**：先在向量存儲中搜索，再通過精心設計的提示進行 LLM 驗證
- **動態路由**：基於信心度動態選擇預定義回應或 RAG 管道[13]

3. 查詢理解與擴展技術

RQ-RAG：查詢精煉方法

RQ-RAG 技術[14][15][16] 通過查詢精煉提升檢索效果：

- **查詢重寫**：將模糊查詢重新表述為更明確的形式
- **查詢分解**：將複雜查詢拆分為多個子查詢
- **歧義消解**：利用上下文信息澄清不明確的查詢意圖

查詢擴展策略

多種查詢擴展方法[17][18][19][20] 可改善檢索品質：

- **假設答案生成**：讓 LLM 生成假設性答案來豐富查詢上下文[21][22]
- **多查詢生成**：生成多個相關查詢以增加檢索覆蓋範圍[19][22]
- **上下文化嵌入**：使用 BERT、ELMo 等模型生成查詢感知的上下文嵌入[23][24]

多跳推理與查詢分解

針對需要多步推理的複雜查詢[25][26][27]：

- **逐步推理框架**：結合支持句識別和子問題生成[26][27]
- **推理鏈提取**：維護從查詢到答案的推理步驟序列[28]
- **上下文感知查詢表示**：整合結構性和關係性上下文信息[29]

4. 用戶偏好建模與適應性方法

個性化意圖識別

自適應意圖識別模型[30][31] 能夠學習個體用戶特徵：

- 行為模式識別：在線識別用戶的行為風格和偏好[30][31]
- 強化學習驅動：使用 Q-Learning 等方法動態適應用戶需求[32]
- 雙重偏好對齊：同時進行外部和內部偏好對齊[33][34]

RAGate：自適應檢索門控

RAGate 模型[35][36][37] 智能決定何時需要外部知識增強：

- 對話上下文建模：綜合考慮對話歷史和當前查詢
- 信心度評估：評估模型對回應的信心水平
- 動態增強決策：根據複雜度和信心度決定是否使用 RAG

5. 記憶增強與上下文建模

記憶增強神經網絡 (MANNs)

MANNs 架構[38][39][40][41][42] 為意圖識別提供長期記憶能力：

- 外部記憶模組：存儲和檢索長期上下文信息
- 注意力機制：選擇性關注記憶中的相關部分
- 動態讀寫操作：在處理序列過程中動態更新記憶內容

對話上下文建模

對話感知系統[43][44][45][46] 能夠維護豐富的對話狀態：

- 多粒度上下文：同時建模詞級和話語級的依賴關係[45]
- 時間敏感檢索：處理基於時間和事件順序的查詢[46]
- 異構圖建模：使用圖結構表示複雜的對話上下文[43]

6. 少樣本學習與對比學習

對比學習意圖檢測

CPFT 方法[47][48][49][50][51] 在少量標記數據下實現有效意圖分類：

- 自監督預訓練：在無標籤數據上學習語義辨別能力
- 監督對比學習：明確拉近相同意圖、推遠不同意圖的表示
- 語義相似性處理：特別適用於細粒度且語義相近的意圖識別

開放意圖檢測

自適應決策邊界方法[52][53] 能夠處理未見過的意圖：

- **距離感知表示**：學習有利於開放意圖檢測的特徵表示
- **球形決策邊界**：為每個已知意圖學習適應性的決策邊界
- **風險平衡**：在經驗風險和開放空間風險之間取得平衡

總結

這些方法形成了一個完整的生態系統，能夠與 RAG 應用深度整合。關鍵成功因素包括：

1. **動態適應性**：系統能根據查詢複雜度和用戶偏好動態調整策略
2. **多層次理解**：從語法到語義，從局部到全域的多維度查詢理解
3. **上下文感知**：充分利用對話歷史和長期記憶信息
4. **效率優化**：在準確性和回應時間之間取得最佳平衡
5. **持續學習**：通過用戶反饋和互動持續改善系統表現

這些方法的組合使用能夠構建出既準確又高效的智能查詢理解系統，為 RAG 應用提供強有力的意圖預測支持。