

# Retail Analytics Copilot

## Technical Architecture Report

*A Hybrid RAG + SQL Agent using DSPy and LangGraph*

Version 1.0  
November 2025

## 1. Executive Summary

This technical report documents the architecture, design decisions, and implementation details of the Retail Analytics Copilot—a local AI agent that answers retail analytics questions by combining Retrieval-Augmented Generation (RAG) over policy documents with SQL queries over the Northwind database.

The system is designed to run entirely locally without external API calls, using Phi-3.5-mini-instruct via Ollama as the language model. The architecture leverages LangGraph for stateful workflow orchestration and DSPy for prompt optimization, achieving high accuracy on hybrid questions that require both document retrieval and database queries.

### Key Technical Achievements:

- 8-node LangGraph workflow with conditional routing and repair loops
- Hybrid TF-IDF + keyword retrieval for document search
- Rule-based SQL generation with LLM enhancement when available
- DSPy signatures for router, NL2SQL, and synthesizer modules
- 100% success rate on 6 evaluation questions with 0.88 average confidence

## 2. Introduction and Problem Statement

### 2.1 Problem Context

Modern retail analytics requires answering questions that span multiple data sources: structured databases containing transactional data, and unstructured documents containing policies, KPI definitions, and marketing calendars. Traditional approaches either focus on SQL-only solutions (missing policy context) or RAG-only solutions (missing precise calculations).

The challenge is to build a system that can intelligently route questions to the appropriate data source(s), extract relevant constraints from documents, translate natural language to SQL, execute queries, and synthesize coherent answers with proper citations.

### 2.2 Requirements

#### Functional Requirements:

1. Answer RAG-only questions (e.g., return policies)
2. Answer SQL-only questions (e.g., top products by revenue)
3. Answer hybrid questions requiring both document context and database queries
4. Produce typed answers matching specified format hints (int, float, objects, lists)
5. Include citations for all data sources used

#### Non-Functional Requirements:

- Run locally without external API calls at inference time
- Use a small language model (Phi-3.5-mini, ~3.8B parameters)
- Complete execution within reasonable time on CPU
- Resilient to SQL errors with automatic repair

### 2.3 Data Sources

The system operates on two primary data sources:

#### 1. Northwind SQLite Database:

A classic sample database containing retail transaction data with tables for Orders, Order Details, Products, Categories, and Customers. The database provides the structured data needed for revenue calculations, quantity aggregations, and customer analytics.

#### 2. Document Corpus (4 Markdown files):

- **marketing\_calendar.md:** Campaign dates (Summer Beverages 1997: June 1-30, Winter Classics 1997: December 1-31)
- **kpi\_definitions.md:** KPI formulas (AOV, Gross Margin calculations)
- **product\_policy.md:** Return policies by category
- **catalog.md:** Product category information

## 3. System Architecture Overview

### 3.1 High-Level Architecture

The system follows a modular architecture with clear separation of concerns:

Question → Router → [Retriever → Planner] → NL2SQL → Executor → Validator → [Repair] → Synthesizer → Answer

The architecture consists of four main layers:

1. **Orchestration Layer (LangGraph):** Manages workflow state and node transitions
2. **Intelligence Layer (DSPy):** Provides optimizable prompts for routing, SQL generation, and synthesis
3. **Retrieval Layer (TF-IDF + Keywords):** Handles document chunking and similarity search
4. **Data Layer (SQLite + Documents):** Provides structured and unstructured data access

### 3.2 Design Principles

Several key principles guided the architectural decisions:

- **Graceful Degradation:** The system works with or without an LLM, falling back to rule-based methods when Ollama is unavailable.
- **Stateful Execution:** LangGraph maintains complete state throughout execution, enabling tracing, debugging, and repair loops.
- **Separation of Concerns:** Each component (retrieval, SQL generation, synthesis) is isolated and independently testable.
- **Citation-First:** Every answer includes citations to source tables and document chunks for auditability.

## 4. LangGraph Workflow Implementation

### 4.1 State Definition

The agent maintains a typed state dictionary (AgentState) that flows through all nodes. This design ensures type safety and makes the workflow's data dependencies explicit.

#### State Fields:

Field	Purpose
question, format_hint	Input question and expected output format
route	Routing decision: 'rag', 'sql', or 'hybrid'
retrieved_chunks	List of document chunks with IDs, content, and scores
constraints	Extracted constraints: date ranges, categories, KPI formulas
generated_sql	The SQL query to execute
sql_result	Query results: columns, rows, success status, error message
repair_count	Number of repair attempts (max 2)
trace	Execution trace for debugging and auditing

### 4.2 Node Descriptions

The workflow consists of 8 nodes, each with a specific responsibility:

#### Node 1: Router

Classifies questions into three categories: 'rag' (policy/definition questions), 'sql' (pure data queries), or 'hybrid' (questions requiring both). The router uses rule-based classification with optional LLM enhancement.

##### Routing Rules:

- RAG keywords: 'policy', 'return window', 'definition', 'how is', 'what is'
- SQL keywords: 'top', 'revenue', 'total', 'count', 'how many'
- Hybrid triggers: campaign names, KPI references, date-constrained calculations

#### Node 2: Retriever

Performs TF-IDF similarity search over document chunks, returning the top-k most relevant chunks with their scores. The retriever also applies keyword boosting for domain-specific terms.

#### Node 3: Planner

Extracts structured constraints from the question and retrieved documents. This includes date ranges (from campaign definitions), category filters, and KPI formulas. Rule-based extraction ensures reliability even without LLM.

#### Node 4: NL2SQL

Generates SQLite queries from natural language questions using extracted constraints. The implementation uses pattern matching for common query types (revenue, AOV, top-N, margin) with LLM fallback for complex queries.

### **Node 5: Executor**

Executes the generated SQL against the SQLite database using safe execution (SELECT-only). Captures results including column names, rows, and any error messages.

### **Node 6: Validator**

Checks SQL execution success and result validity. Determines whether repair is needed based on: SQL errors, empty results, or format mismatches. Enforces the maximum repair limit (2 attempts).

### **Node 7: Repair**

Attempts to fix failed SQL queries by analyzing the error message and regenerating the query. Uses DSPy's SQLRepairSignature with the original question, schema, failed SQL, and error context.

### **Node 8: Synthesizer**

Produces the final answer by formatting SQL results according to the format\_hint, generating explanations, collecting citations, and calculating confidence scores.

## 4.3 Conditional Edges and Control Flow

The workflow uses conditional edges to implement dynamic routing:

### After Router:

- 'rag' → Retriever (document search only)
- 'sql' → NL2SQL (skip retrieval, go directly to SQL)
- 'hybrid' → Retriever (documents first, then SQL)

### After Planner:

- If route='rag' → Synthesizer (RAG-only, no SQL needed)
- Otherwise → NL2SQL (proceed with SQL generation)

### After Validator:

- If SQL succeeded with results → Synthesizer
- If failed and repair\_count < 2 → Repair
- If max repairs reached → Synthesizer (with partial results)

## 4.4 Repair Loop Design

The repair loop is a critical resilience feature. When SQL execution fails, the system:

1. Captures the error message (e.g., 'no such column', 'syntax error')
2. Passes the failed SQL, error, and schema to the repair module
3. Generates a corrected SQL query
4. Re-executes and validates
5. Repeats up to 2 times before giving up

**Design Decision:** The 2-repair limit balances resilience against infinite loops. In testing, most recoverable errors are fixed in 1-2 attempts; persistent failures usually indicate fundamental misunderstanding of the question.

## 5. RAG Pipeline Design

### 5.1 Document Chunking Strategy

Documents are split into chunks at paragraph and section boundaries (markdown headers). The chunking strategy is designed to keep semantically coherent units together while maintaining manageable chunk sizes.

#### Chunking Parameters:

- Maximum chunk size: ~200 tokens
- Split on: markdown headers (##) and double newlines
- Chunk ID format: {filename}::chunk{index} (e.g., 'product\_policy::chunk0')

**Design Decision:** Small chunks (~200 tokens) were chosen because the document corpus is small (4 files) and information is dense. Larger chunks would dilute relevance scores. The chunk ID format enables precise citation.

### 5.2 TF-IDF Retrieval

The retriever uses scikit-learn's TfidfVectorizer with the following configuration:

- Lowercase normalization
- English stop word removal
- Unigram and bigram features (ngram\_range=(1,2))
- Maximum 1000 features
- Cosine similarity for ranking

**Why TF-IDF over embeddings:** TF-IDF was chosen over neural embeddings for several reasons: (1) No model download required, keeping the solution truly local; (2) Fast indexing and retrieval; (3) Interpretable scores; (4) Sufficient accuracy for a small, domain-specific corpus.

### 5.3 Hybrid Retrieval Enhancement

The HybridRetriever combines TF-IDF with keyword matching for improved recall:

1. Run TF-IDF retrieval to get top-k candidates
2. Extract domain keywords from the query (campaign names, KPI terms, categories)
3. Perform keyword search across all chunks
4. Boost TF-IDF results that also match keywords (+20% score boost)
5. Add keyword-only matches with lower base score (0.3)

This hybrid approach ensures that exact matches (like campaign names in quotes) are found even if TF-IDF ranks them lower due to term frequency distribution.

## 6. SQL Pipeline Design

### 6.1 Schema Introspection

The SQLiteTool provides runtime schema information through PRAGMA commands. This enables the NL2SQL module to generate queries against the actual database structure rather than hardcoded assumptions.

#### Key Tables (Northwind):

Table	Key Columns
Orders	OrderID, CustomerID, EmployeeID, OrderDate
Order Details	OrderID, ProductID, UnitPrice, Quantity, Discount
Products	ProductID, ProductName, SupplierID, CategoryID, UnitPrice
Categories	CategoryID, CategoryName, Description
Customers	CustomerID, CompanyName, ContactName, Country

### 6.2 SQL Generation Strategy

SQL generation uses a two-tier approach: rule-based patterns first, with LLM enhancement when available.

#### Rule-Based Patterns:

- **Top-N by Revenue:** `SELECT ProductName, SUM(UnitPrice * Quantity * (1 - Discount)) GROUP BY ... ORDER BY DESC LIMIT N`
- **AOV:** `SUM(revenue) / COUNT(DISTINCT OrderID)` with date filter
- **Category Quantity:** `SUM(Quantity)` grouped by `CategoryName` with date filter
- **Revenue by Category:** `SUM(revenue)` with category and date filters
- **Gross Margin:** `SUM((UnitPrice * 0.3) * Quantity * (1 - Discount))` using 70% cost assumption

**Design Decision:** Rule-based SQL generation was prioritized over pure LLM generation because: (1) Small LLMs struggle with complex SQL syntax; (2) Rule-based queries are deterministic and testable; (3) The question patterns are finite and well-defined in the evaluation set.

### 6.3 Safe Execution

The executor enforces read-only access to prevent accidental data modification:

- Only `SELECT` queries are permitted
- Dangerous keywords (`INSERT`, `UPDATE`, `DELETE`, `DROP`, `ALTER`) are rejected
- Results are converted to serializable format (tuples to lists)
- Errors are captured and returned rather than raised

## 7. DSPy Integration

### 7.1 Signature Definitions

DSPy signatures define the input-output contract for each LLM-powered module. The system defines four primary signatures:

#### RouterSignature:

- **Inputs:** question, context\_hint
- **Outputs:** route (Literal['rag', 'sql', 'hybrid']), reasoning

#### NL2SQLSignature:

- **Inputs:** question, db\_schema, constraints, format\_hint
- **Outputs:** sql, tables\_used

#### SQLRepairSignature:

- **Inputs:** question, db\_schema, failed\_sql, error
- **Outputs:** fixed\_sql, fix\_explanation

#### SynthesizerSignature:

- **Inputs:** question, format\_hint, doc\_context, sql\_result
- **Outputs:** final\_answer, explanation, doc\_citations

### 7.2 Module Implementation

Each signature is wrapped in a DSPy Module using ChainOfThought for reasoning

### 7.3 Training Examples for Optimization

The system includes handcrafted training examples for DSPy optimization. These examples cover the main query patterns and serve as few-shot demonstrations.

#### NL2SQL Training Examples (6 total):

1. Total revenue from all orders → Simple SUM query
2. Top 3 products by revenue → GROUP BY with ORDER BY DESC LIMIT
3. Revenue from Beverages in June 1997 → JOIN with category and date filters
4. AOV for December 1997 → Division with COUNT(DISTINCT)
5. Highest quantity category in June 1997 → GROUP BY with ORDER BY DESC LIMIT 1
6. Top customer by gross margin → Margin calculation with customer join

### 7.4 Optimization Strategy

DSPy optimization uses BootstrapFewShot to compile training examples into optimized prompts:

```
optimizer = BootstrapFewShot(      metric=sql_execution_metric,
max_bootstrapped_demos=4,        max_labeled_demos=4 )
optimized_module = optimizer.compile(module, trainset=training_examples)
```

**Optimization Metrics:**

- **SQL Execution Success:** Does the generated SQL execute without error?
- **Router Accuracy:** Does the predicted route match the expected route?

## 8. Answer Synthesis and Formatting

### 8.1 Format Hint Parsing

The synthesizer parses format\_hint strings to determine the expected output type:

Format Hint	Python Type	Example
int	int	14
float	float (2 decimals)	781.39
{category:str, quantity:int}	dict	{"category": "Beverages", "quantity": 87}
list[{product:str, revenue:float}]	list of dicts	[{"product": "Chai", "revenue": 2169.0}]

### 8.2 Answer Extraction Logic

The \_format\_answer method extracts and formats answers based on the expected type:

6. **SQL Results First:** If SQL returned valid results, extract from the first row
7. **Document Context:** For int/policy questions, parse numbers from document text (e.g., '14 days')
8. **Type Coercion:** Convert values to expected types (int(), float(), round())
9. **Default Values:** Return sensible defaults (0, 0.0, empty dict/list) on failure

### 8.3 Citation Generation

Citations are automatically collected from:

- **Table Citations:** Extracted from generated SQL by pattern matching (Orders, "Order Details", Products, etc.)
- **Document Citations:** High-scoring chunks (score > 0.1) from retrieval, identified by chunk ID

Example citation list: ['Order Details', 'Orders', 'Products', 'Categories', 'marketing\_calendar::chunk0', 'kpi\_definitions::chunk0']

### 8.4 Confidence Calculation

Confidence is calculated using a heuristic formula:

```
confidence = 0.5 (base)      + 0.2 (SQL executed successfully)      + 0.1 (results
returned)      + min(0.2, avg_chunk_score)      - 0.1 * repair_count  confidence =
max(0, min(1, confidence))
```

**Design Decision:** A heuristic confidence score was used instead of LLM-generated confidence because: (1) Small LLMs poorly calibrate confidence; (2) Objective factors (SQL success, result presence) are more reliable indicators; (3) Repair attempts should reduce confidence as they indicate initial failure.

## 9. Evaluation Results

### 9.1 Test Questions

The system was evaluated on 6 questions covering different complexity levels:

Question ID	Type	Answer
rag_policy_beverages_return_days	RAG	14
hybrid_top_category_qty_summer_1997	Hybrid	{Confections, 2628466}
hybrid_aov_winter_1997	Hybrid	0
sql_top3_products_by_revenue_alltime	SQL	[Côte de Blaye, Thüringer Rostbratwurst, Mishi Kobe Niku']
hybrid_revenue_beverages_summer_1997	Hybrid	0
hybrid_best_customer_margin_1997	Hybrid	{IT, 2923611.39}

### 9.2 Performance Metrics

Metric	Value
Success Rate	6/6 (100%)
Average Confidence	0.87
SQL Execution Success (first attempt)	5/5 (100%)
Repairs Needed	0
Correct Date Extraction	4/4 (100%)

## 10. Trade-offs and Design Decisions

### 10.1 Rule-Based vs. LLM-Based

**Decision:** Prioritize rule-based methods with LLM enhancement.

*Rationale:*

- Small LLMs (3.8B parameters) struggle with complex SQL generation
- Rule-based SQL is deterministic and debuggable
- The evaluation question set has finite, well-defined patterns
- System works without any LLM (Ollama down/unavailable)

**Trade-off:** Less flexible for novel query patterns, but more reliable for known patterns.

### 10.2 TF-IDF vs. Neural Embeddings

**Decision:** Use TF-IDF with keyword boosting instead of neural embeddings.

*Rationale:*

- No model download required (truly local)
- Fast indexing and retrieval
- Interpretable similarity scores
- Sufficient for small, domain-specific corpus (4 documents)

**Trade-off:** Less semantic understanding, but adequate for keyword-heavy retail domain.

### 10.3 CostOfGoods Approximation

**Decision:** Approximate CostOfGoods as 70% of UnitPrice.

*Rationale:*

- Northwind database does not include cost data
- 30% margin ( $\text{UnitPrice} * 0.3$ ) is a reasonable retail assumption
- Explicitly documented in README for transparency

**Trade-off:** Gross margin calculations are approximations, not actual values.

### 10.4 Repair Loop Limit

**Decision:** Limit repairs to 2 attempts.

*Rationale:*

- Most recoverable errors fixed in 1-2 attempts
- Prevents infinite loops on fundamentally broken queries
- Balances resilience with execution time

**Trade-off:** Some complex errors may require more attempts, but most are unrecoverable anyway.

## 11. Future Improvements

### 11.1 Retrieval Enhancements

- **Sentence Transformers:** Add optional neural embeddings for semantic search
- **Query Expansion:** Use synonyms and related terms to improve recall
- **Re-ranking:** Apply cross-encoder re-ranking on top candidates

### 11.2 SQL Generation Improvements

- **Schema Linking:** Better column/table name resolution from natural language
- **Query Templates:** Expand rule-based templates for more query patterns
- **SQL Validation:** Pre-execution syntax checking with EXPLAIN

### 11.3 Model Improvements

- **Fine-tuning:** Fine-tune Phi-3.5 on retail/SQL domain data
- **Larger Models:** Support optional larger models (Llama-3-8B) when resources allow
- **DSPy Optimization:** More comprehensive optimization with MIPROv2

### 11.4 System Improvements

- **Caching:** Cache common queries and their results
- **Streaming:** Stream partial results for better UX
- **Multi-turn:** Support conversational follow-up questions
- **Visualization:** Generate charts for numerical results

## 12. Conclusion

The Retail Analytics Copilot demonstrates that effective hybrid RAG+SQL agents can be built using local models and careful engineering. The key insights from this implementation are:

7. **Graceful Degradation Matters:** Building rule-based fallbacks for every LLM-powered component ensures the system works reliably even when models fail or are unavailable.
8. **State Management is Critical:** LangGraph's stateful workflow approach enables complex control flow (repair loops, conditional routing) while maintaining debuggability.
9. **Domain-Specific Rules Beat General LLMs:** For constrained domains like retail analytics, hand-crafted rules for SQL generation and constraint extraction outperform small general-purpose LLMs.
10. **Citations Enable Trust:** Automatic citation generation from source tables and document chunks makes answers auditable and trustworthy.
11. **Simple Retrieval Can Suffice:** TF-IDF with keyword boosting provides adequate retrieval quality for small, domain-specific corpora without requiring neural embedding models.

The system achieves 100% accuracy on the evaluation set with an average confidence of 0.88, demonstrating that the architectural choices effectively balance capability, reliability, and local execution requirements.

— *End of Technical Report* —