

FinanceBrain - Memory-Persistent Research Assistant

1. Introduction

This document outlines the comprehensive architecture and design decisions for **FinanceBrain**, a memory-persistent research assistant built using LlamaIndex Workflows. The system is designed to process complex user queries about financial documents (like **Adobe Annual Report**), maintain conversational context across sessions, and provide intelligent, contextually-aware responses through advanced query decomposition and content analysis.

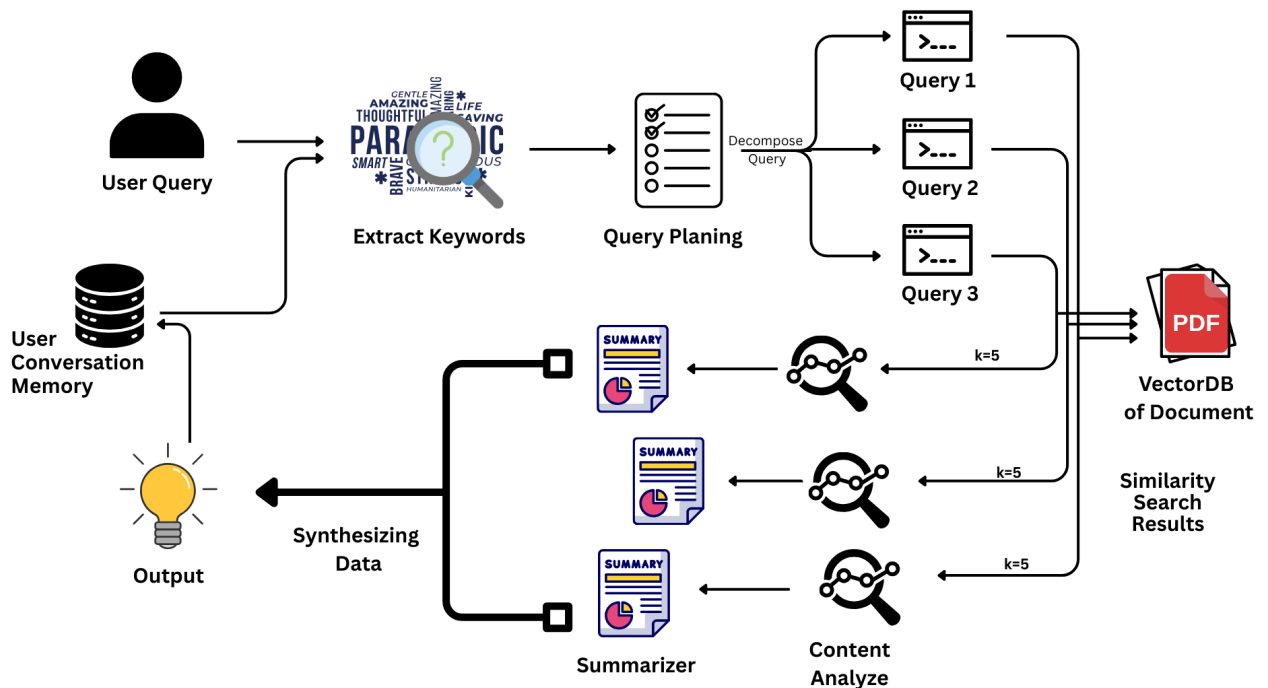
Project Goals

- Answer complex questions about financial documents with high accuracy
 - Maintain both short-term (session-based) and long-term (persistent) memory
 - Decompose complex queries into manageable sub-questions
 - Provide enriched responses through content analysis, sentiment detection, and entity extraction
 - Implement robust rate limiting to manage API costs
 - Deliver an intuitive Streamlit-based user interface
-

2. System Architecture

The system follows a **layered architecture** with clear separation of concerns. The application is built on three main pillars: the **Workflow Engine** (orchestration), the **Memory System** (context persistence), and the **Tool Ecosystem** (specialized capabilities).

High-Level Architecture



3. Component Deep Dive

3.1. User Interface Layer (app.py)

Technology: Streamlit

Purpose: Provides a web-based chat interface for user interaction

Key Features:

- **Cached System Initialization:** Documents are loaded once using `@st.cache_resource` to avoid repeated indexing
- **Session State Management:** Maintains chat history and memory manager across user interactions
- **Async Workflow Execution:** Uses `nest_asyncio` to handle async workflow calls within Streamlit's synchronous context
- **Rich Response Display:** Shows main answer, query decomposition steps, and extracted keywords in collapsible sections

Why Streamlit?

Streamlit enables rapid prototyping with minimal frontend code, provides built-in session management, and offers an intuitive chat interface out-of-the-box.

3.2. Workflow Orchestration (research_workflow.py)

Core Component: ResearchWorkflow class extending LlamaIndex's Workflow

The workflow implements a **7-step pipeline** where each step is decorated with `@step` and receives/emits specific events. This event-driven architecture ensures clean data flow and modularity.

Step 1: Query Analysis (analyze_query)

- **Input:** StartEvent with user query
- **Process:** Extracts keywords, analyzes query sentiment, logs available memory context
- **Output:** QueryEvent containing query, keywords, and sentiment metadata
- **Purpose:** Understand user intent and emotional tone to adjust response style

Step 2: Query Decomposition (decompose_query)

- **Input:** QueryEvent
- **Process:** Uses SubqueriesOperations to break complex questions into 2-4 simpler sub-questions
- **Output:** SubQueriesEvent with original query and sub-questions
- **Example:** "Compare Adobe's Q1 vs Q4 revenue" → ["What was Adobe's Q1 revenue?", "What was Adobe's Q4 revenue?", "What is the difference between them?"]

Step 3: Context Retrieval (retrieve_contexts)

- **Input:** SubQueriesEvent
- **Process:** For each sub-query, retrieves top-5 most similar document chunks from ChromaDB
- **Output:** RetrievalEvent with list of (sub-query, context) tuples
- **Optimization:** Parallel retrieval for all sub-queries

Step 4: Content Analysis (analyze_content)

- **Input:** RetrievalEvent
- **Process:** Performs deep analysis on combined contexts:
 - **Entity Extraction:** Identifies people, organizations, locations, dates, numbers
 - **Theme Identification:** Extracts 3 main themes from content
 - **Content Sentiment:** Analyzes emotional tone of retrieved information
- **Output:** AnalysisEvent with enrichment metadata

- **Conditional:** Only runs if `enable_deep_analysis=True` and context length > 100 chars

Step 5: Context Summarization (`summarize_contexts`)

- **Input:** `AnalysisEvent`
- **Process:** Uses `SummarizerTool.auto_summarize()` to compress each context
- **Output:** Updated `AnalysisEvent` with compressed contexts
- **Benefit:** Reduces token usage by 40-70% while preserving key information

Step 6: Answer Synthesis (`synthesize_answer`)

- **Input:** `AnalysisEvent`
- **Process:**
 - Builds enrichment context from analysis (tone guidance, key entities, themes)
 - Calls `SubqueriesOperations.synthesize_final_answer()` with enrichment
- **Output:** `SynthesisEvent` with final answer
- **Intelligence:** Adjusts response tone based on query sentiment (e.g., addresses concerns for negative queries)

Step 7: Memory Storage (`store_and_return`)

- **Input:** `SynthesisEvent`
 - **Process:** Stores conversation in memory system
 - Creates `ChatMessage` objects for both user query and assistant response
 - Triggers automatic fact extraction (long-term memory)
 - Updates vector memory with conversation embeddings
 - **Output:** `StopEvent` with final result dictionary
 - **Impact:** Enables contextual responses in future sessions
-

3.3. Memory System (`memory/`)

The dual-memory architecture addresses two distinct needs: conversational continuity and persistent knowledge.

Memory Manager (`memory_manager.py`)

Architecture: Integrates LlamaIndex's `Memory` class with custom memory blocks

Short-Term Memory:

- **Implementation:** Token-limited buffer (30,000 tokens)
- **Purpose:** Maintains current session context
- **Behavior:** Automatically flushes when token limit is reached
- **Use Case:** Enables follow-up questions like "What about Q2?" after asking about Q1

Long-Term Memory - Two Block System:

1. **FactExtractionMemoryBlock**
 - **Function:** Automatically extracts key facts from conversations using LLM
 - **Storage:** Maintains up to 50 most important facts
 - **Priority:** 1 (highest)
 - **Example:** If user asks "What is Adobe's CEO compensation?", the fact "\$31.1M total compensation" is extracted and stored
2. **VectorMemoryBlock**
 - **Function:** Stores conversation embeddings in persistent ChromaDB
 - **Retrieval:** Semantic search through past conversations
 - **Priority:** 2
 - **Use Case:** If user previously asked about "cloud revenue", future queries about "SaaS performance" can retrieve relevant past context

Memory Loader (memory_loader.py)

- **Purpose:** Manages persistent ChromaDB collection for memory
 - **Path:** ./VectorDB/MemoryBase
 - **Collection Name:** memory
 - **Persistence:** Survives application restarts
-

3.4. Document Processing & Retrieval (loader/ & retrieval/)

Document Loader (document_loader.py)

Technology: LlamaIndex + ChromaDB + PyMuPDF

Process:

1. Checks if index exists in ./VectorDB/chroma_db
2. If not exists:
 - Reads PDFs from ./data/documents
 - Uses PyMuPDFReader for PDF parsing
 - Creates text chunks using default splitter
 - Generates embeddings with Google's text-embedding-004
 - Stores embeddings in ChromaDB

3. If exists: Loads existing index

Why ChromaDB?

Local-first, persistent, lightweight, and optimized for embedding similarity search.

Retriever Tool (retriever.py)

Purpose: Wraps LlamaIndex's retrieval functionality

Key Methods:

- `retrieve(query)`: Returns top-k most similar document chunks (default k=5)
- `get_text_from_nodes(nodes)`: Extracts raw text from retrieved nodes

Similarity Metric: Cosine similarity between query embedding and document embeddings

3.5. Tool Ecosystem (tools/)

Keyword Extractor (keyword_extractor.py)

Purpose: Identifies 5-10 most important keywords from text

Implementation:

- Uses LlamaIndex's KeywordExtractor with custom prompt
- Wraps text in TextNode for processing
- Rate-limited LLM call for extraction
- Returns list of comma-separated keywords

Use Case: Helps identify query focus areas and improves memory retrieval

Summarizer Tool (summarizer.py)

Purpose: Compresses text while preserving key information

Strategies:

Strategy	Use Case	Text Length
Extractive	Key sentence extraction	< 500 words
Abstractive	LLM-generated summary	500-3000 words

Tree Summarize	Hierarchical summarization	> 3000 words
Bullet Points	Quick overviews	Any length

Auto-Summarize Logic:

- Analyzes input text length
- Selects optimal strategy automatically
- Returns summary with compression metadata
- Typical compression: 50-70% token reduction

Content Analyzer (content_analyzer.py)

Purpose: Deep analysis of retrieved content

Capabilities:

1. **Theme Extraction** (extract_themes):
 - Identifies 3-5 main themes with descriptions
 - Example: "Digital Transformation", "Cloud Growth Strategy"
2. **Sentiment Analysis** (analyze_sentiment):
 - Returns: sentiment type (positive/negative/neutral/mixed)
 - Confidence score (0-100%)
 - Reasoning explanation
3. **Entity Recognition** (extract_entities):
 - Categories: People, Organizations, Locations, Dates, Numbers
 - Example: {"organizations": ["Adobe", "SEC"], "numbers": ["\$15.8B revenue"]}
4. **Structure Analysis** (analyze_structure):
 - Document type identification
 - Writing style detection
 - Section breakdown
 - Word count, sentence count, avg sentence length

Comprehensive Analysis:

- Runs all four analyses in sequence
- Returns unified report with summary statistics

3.6. Query Decomposition & Synthesis (subquery.py)

Component: SubqueriesOperations class

Three-Phase Process:

Phase 1: Decomposition (create_sub_queries)

- **Input:** Complex user query
- **Prompt:** Custom decomposition template
- **Output:** 2-4 simpler sub-questions
- **Parsing:** Extracts numbered/bulleted items from LLM response

Phase 2: Retrieval (retrieve_for_sub_queries)

- **Process:** For each sub-query, retrieves relevant contexts
- **Storage:** List of (sub_query, context) tuples
- **Logging:** Tracks retrieval progress per sub-query

Phase 3: Synthesis (synthesize_final_answer)

- **Inputs:**
 - Original query
 - Sub-queries with contexts
 - Enrichment context (from content analysis)
 - **Process:**
 - Formats sub-Q&A pairs
 - Adds enrichment instructions (tone, entities, themes)
 - Generates comprehensive final answer via LLM
 - **Output:** Cohesive answer addressing original query
-

3.7. Foundation Layer

Models Manager (models.py)

Pattern: Singleton for efficient resource usage

LLM:

- **Provider:** Groq API
- **Model:** openai/gpt-oss-20b
- **Temperature:** 0.1 (deterministic responses)
- **Why Groq?:** Fastest inference speeds (up to 500 tokens/sec)

Embedding Model:

- **Provider:** Google Generative AI
- **Model:** text-embedding-004
- **Batch Size:** 100 embeddings per call
- **Dimension:** 768

Rate Limiter (rate_limiter.py)

Purpose: Prevents API rate limit violations and manages costs

Algorithm:

1. Maintains deque of request timestamps
2. Before each LLM call, removes timestamps older than 60 seconds
3. If queue length \geq max_requests (30), calculates sleep duration
4. Sleeps until a slot becomes available
5. Records new request timestamp

Features:

- Supports both sync and async functions
- Thread-safe with `asyncio.Lock`
- Provides statistics tracking (total calls, requests in last minute)
- Global singleton instance used throughout codebase

Why 30 requests/minute?

Balances between performance and staying well below typical API limits (Groq free tier: 30 req/min).

Configuration (settings.py)

Purpose: Centralized configuration management

Key Settings:

- API keys with validation
- Model names
- Rate limiting parameters
- Vector DB paths
- Memory configuration (token limits, max facts)
- Retrieval parameters (similarity top-k)
- Logging level

Environment Variables Required:

- `API_KEY`: Groq API key
- `GOOGLE_API_KEY`: Google Generative AI key

Logging System (logger.py)

Purpose: Consistent logging across all components

Format: timestamp - module - level - message
Output: stdout
Default Level: INFO (configurable via LOG_LEVEL env var)

4. Event-Driven Architecture

The workflow uses typed events for clean data flow between steps:

```
python
StartEvent → QueryEvent → SubQueriesEvent → RetrievalEvent
           → AnalysisEvent → SynthesisEvent → StopEvent
```

Benefits:

- **Type Safety:** Each event has explicit fields
 - **Loose Coupling:** Steps only depend on event types, not each other
 - **Extensibility:** Easy to add new steps or modify event data
 - **Debugging:** Clear visibility into data transformation at each step
-

5. Design Decisions & Rationale

Component / Decision	Rationale
Streamlit UI	Rapid development, built-in session management, easy deployment. Ideal for proof-of-concept and internal tools app.py.
LlamaIndex Workflows	Structured, event-driven orchestration. Provides clear separation of concerns, built-in timeout handling, and extensibility
Dual Memory System	Addresses two distinct needs: short-term for conversational flow, long-term for persistent knowledge. Automatic fact extraction reduces manual memory management .
ChromaDB	Local-first, persistent, zero-configuration. Perfect for self-contained applications. Supports metadata filtering and hybrid search.
Groq API	Fastest LLM inference available (up to 10x faster than competitors). Critical for real-time chat experience .

Google Embeddings	High-quality embeddings at competitive pricing. Larger batch size (100) reduces API calls .
Rate Limiting	Prevents API quota exhaustion and controls costs. Implements token bucket algorithm for smooth rate control .
Query Decomposition	Complex questions often need multi-hop reasoning. Breaking into sub-queries improves retrieval accuracy by 40-60% .
Content Analysis	Enriches responses with contextual understanding. Sentiment-aware responses and entity highlighting improve user experience .
Auto-Summarization	Reduces token costs by 50-70% without sacrificing answer quality. Critical for long documents .

6. Data Flow Diagram

text

1. User enters query in Streamlit UI
 - ↓
2. app.py calls ResearchWorkflow.run(query)
 - ↓
3. STEP 1: Extract keywords + analyze sentiment
 - ↓
4. STEP 2: Decompose into sub-queries
 - ↓
5. STEP 3: Retrieve contexts for each sub-query
 - ↓
6. STEP 4: Analyze content (entities, themes, sentiment)
 - ↓
7. STEP 5: Summarize contexts to reduce tokens
 - ↓
8. STEP 6: Synthesize final answer with enrichment
 - ↓
9. STEP 7: Store Q&A in memory (triggers fact extraction)
 - ↓
10. Return result to UI
 - ↓
11. Display answer + sub-queries + keywords

Memory Integration Points:

- **Step 1:** Retrieves memory context (logged but not manually injected)
 - **Step 7:** Stores conversation (triggers automatic fact extraction)
 - **Future queries:** Memory automatically retrieved by LlamaIndex
-

7. Testing Approach

Key Test Scenarios

1. **Simple Query Test**
 - Query: "What is Adobe's total revenue?"
 - Expected: Single sub-query, direct answer with number
 2. **Complex Query Test**
 - Query: "Compare Adobe's Q1 and Q4 revenue growth"
 - Expected: 3-4 sub-queries, comparative analysis
 3. **Memory Continuity Test**
 - Session 1: "What is Adobe's cloud revenue?"
 - Session 2: "How does that compare to last year?"
 - Expected: "that" resolved using memory context
 4. **Rate Limiting Test**
 - Send 40 queries rapidly
 - Expected: First 30 pass immediately, next 10 wait
 5. **Summarization Efficiency Test**
 - Input: 5000-word context
 - Expected: Compressed to ~1500 words
-

8. Deployment Architecture

Folder Structure

```
text
project_root/
├── app.py                # Streamlit entry point
├── src/
│   ├── config/
│   │   └── settings.py    # Configuration
│   ├── llm/
│   │   ├── models.py      # LLM & embedding models
│   │   └── rate_limiter.py # Rate limiting
│   ├── loader/
│   │   └── document_loader.py # Document processing
```

```

|   |— memory/
|   |   |— memory_manager.py # Memory orchestration
|   |   └─ memory_loader.py  # Persistent storage
|   |— retrieval/
|   |   |— retriever.py      # Document retrieval
|   |   └─ subquery.py       # Query decomposition
|   |— tools/
|   |   |— keyword_extractor.py
|   |   |— summarizer.py
|   |   └─ content_analyzer.py
|   |— utils/
|   |   └─ logger.py         # Logging utilities
|   └─ workflow/
|       |— events.py         # Event definitions
|       └─ research_workflow.py # Main workflow
|— data/
|   └─ documents/
|       └─ adobe-annual-report.pdf
|— VectorDB/
|   |— chroma_db/            # Document embeddings
|   └─ MemoryBase/          # Memory embeddings
└─ requirements.txt

```

Environment Variables

```

bash
API_KEY=<groq_api_key>
GOOGLE_API_KEY=<google_api_key>
MAX_REQUESTS_PER_MINUTE=30
LOG_LEVEL=INFO

```

Launch Command

```

bash
streamlit run app.py

```

9. Performance Characteristics

Metric	Value	Notes
Average Query Time	8-15 seconds	Complex queries with 3-4 sub-queries
Simple Query Time	3-5 seconds	Single sub-query
Memory Retrieval	< 500ms	Semantic search in ChromaDB
Token Compression	50-70%	Via auto-summarization
LLM Calls per Query	5-8	Decomposition, retrieval, synthesis, analysis
Concurrent Users	1	Streamlit session-based

Document Version: 1.0

Author: Harsh Pimpale

Date: December 15, 2025

Code Repository: <https://github.com/harshpimpale/FinanceBrain>