

# Contents

HaikuGraph Project Documentation

Table of Contents

Project Overview

Architecture

Implementation Task List

System Flow Diagrams

What We Showcase

Solution Differentiation

Pros and Cons Analysis

Future Roadmap

Conclusion

Getting Started

1

1

1

3

6

10

12

13

16

16

17

## HaikuGraph Project Documentation

### Natural Language Data Assistant with Interactive Visualizations

#### Table of Contents

1. Project Overview

2. Architecture

3. Implementation Task List

4. System Flow Diagrams

5. What We Showcase

6. Solution Differentiation

7. Pros and Cons Analysis

#### Project Overview

HaikuGraph is a natural language data assistant that allows users to query their data using plain English. The system translates questions into SQL, executes them against a DuckDB database, and presents results through interactive visualizations with full transparency into the query planning and execution process.

#### Key Features

- □

Natural Language Queries

: Ask questions in plain English

•

□

Rich Visualizations

: Auto-generated charts, tables, and number displays

•

□

Full Transparency

: View intent classification, query plans, and generated SQL

•

□

Smart Comparisons

: Automatic period-over-period analysis

•

□

Interactive UI

: Filter, sort, paginate, zoom, and export data

•

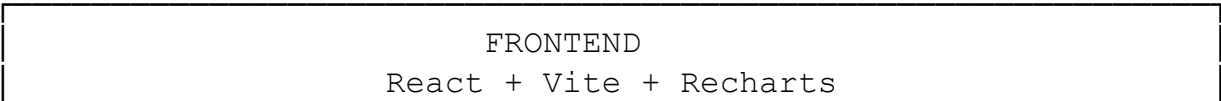
□

Data Quality Insights

: Surface and visualize NULL/missing data

#### Architecture

##### High-Level System Architecture



(http://localhost:5173)

#### User Interface

- Question input with examples
- Answer display with metadata
- Comparison cards (period-over-period)
- Interactive visualizations (charts/tables)
- Explainability tabs (Intent/Plan/SQL/Raw Data)

↑ HTTP REST API

#### BACKEND

FastAPI + Python  
(http://localhost:8000)

#### API Layer (FastAPI)

- POST /ask - Main query endpoint
- GET /health - Health check
- CORS middleware for dev

↓

#### Processing Pipeline

1. Intent Classification (LLM + Rules)
2. Query Planning (Ollama)
3. SQL Generation & Execution (DuckDB)
4. Comparison Extraction
5. Narration (Ollama)
6. Visualization Hints

↓

#### DATA LAYER

DuckDB  
(In-Memory)

- test\_1\_1
- test\_2\_1
- test\_3\_1
- test\_4\_1
- test\_5\_1

Excel Files  
(.xlsx/.xls)

Profile JSON  
(Metadata)

CSV Data

## Data Cards (Semantics)

### Technology Stack

#### Backend

Component	Technology	Purpose
Web Framework	FastAPI	Modern, async Python API
Server	Uvicorn	ASGI server with hot reload
Database	DuckDB	In-process analytical DB
LLM Runtime	Ollama	Local LLM for planning/narration
Data Validation	Pydantic	Request/response schemas

#### Frontend

Component	Technology	Purpose
UI Framework	React 18	Component-based UI
Build Tool	Vite	Fast dev server & bundler
Charts	Recharts	Composable chart library
Icons	Lucide React	Icon components
State	React Hooks	Local state management

### Implementation Task List

#### Phase 1: Foundation & Data Ingestion

- ☒ **Data Ingestion System**
  - Excel/CSV file reader with auto-fallback to string mode
  - DuckDB table creation with sanitized naming
  - Automatic type inference with error recovery
  - Profile generation with statistics
- ☒ **Schema Analysis**
  - Column type detection (numeric, date, string, identifier)
  - Null percentage and distinct count analysis
  - Semantic hints (timestamp, money, status columns)
  - Table grain and primary key detection
- ☒ **Data Cards Generation**
  - Table cards with suggested metrics
  - Column cards with semantic annotations
  - Relation cards for potential joins
  - Index creation for fast lookup

#### Phase 2: Query Processing Pipeline

- ☒ **Intent Classification**
  - Rule-based classifier for common patterns

- LLM fallback for ambiguous cases
- Confidence scoring and rationale
- Comparison detection
- ☒ **Query Planning**
  - LLM-based plan generation (Ollama)
  - Subquestion decomposition
  - Constraint extraction (time, filters)
  - Table and column selection
- ☒ **SQL Generation**
  - Dynamic SQL builder from plan
  - Aggregation support (SUM, COUNT, AVG, MIN, MAX)
  - GROUP BY with time bucketing
  - JOIN path resolution
  - Type casting (VARCHAR → TIMESTAMP, VARCHAR → DOUBLE)
- ☒ **Execution & Results**
  - DuckDB query execution
  - Error handling and retries
  - Result formatting with preview rows
  - Metadata collection (timing, row counts)

### Phase 3: Advanced Features

- ☒ **Comparison Extraction**
  - Period-over-period detection
  - Delta and percentage calculation
  - Normalized comparison structure
  - Direction indicators (up/down/flat)
- ☒ **Narration**
  - LLM-based result explanation
  - Context-aware language generation
  - Comparison-aware narratives
  - Human-readable formatting
- ☒ **Visualization Hints**
  - Shape analysis (rows, columns, types)
  - Auto-detection of chart types
  - Axis and units inference
  - Display hint generation (number/bar/line/table)

### Phase 4: Web UI Development

- ☒ **Basic UI Components**
  - Question input form
  - Demo question buttons
  - Answer display card
  - Loading states and error handling
- ☒ **Visualization Components**
  - Number display for single values
  - Bar charts for grouped metrics
  - Line charts for time series
  - Data tables for multi-column results
- ☒ **Comparison UI**
  - Side-by-side value display
  - Delta and percentage badges
  - Direction indicators with icons
  - Color-coded changes
- ☒ **Explainability Interface**
  - Collapsible details section

- Tabbed navigation (Intent/Plan/SQL/Metadata)
- JSON pretty-printing
- Execution timing display

## Phase 5: Interactive Enhancements

- ☒ **Interactive Charts**
  - Brush component for zoom/pan
  - Clickable legend to toggle series
  - Custom tooltips with formatting
  - Export to SVG functionality
  - NULL value handling with “(No data)” label
- ☒ **Advanced Tables**
  - Pagination (10/25/50/100 rows)
  - Column sorting (asc/desc)
  - Column filtering with search
  - Export to CSV
  - First/Prev/Next/Last navigation
  - Row count display
  - NULL value highlighting
- ☒ **Raw Data Tab**
  - Complete preview\_rows JSON display
  - Status and column metadata
  - Row count badges
  - Debugging transparency

## Phase 6: Data Quality & Testing

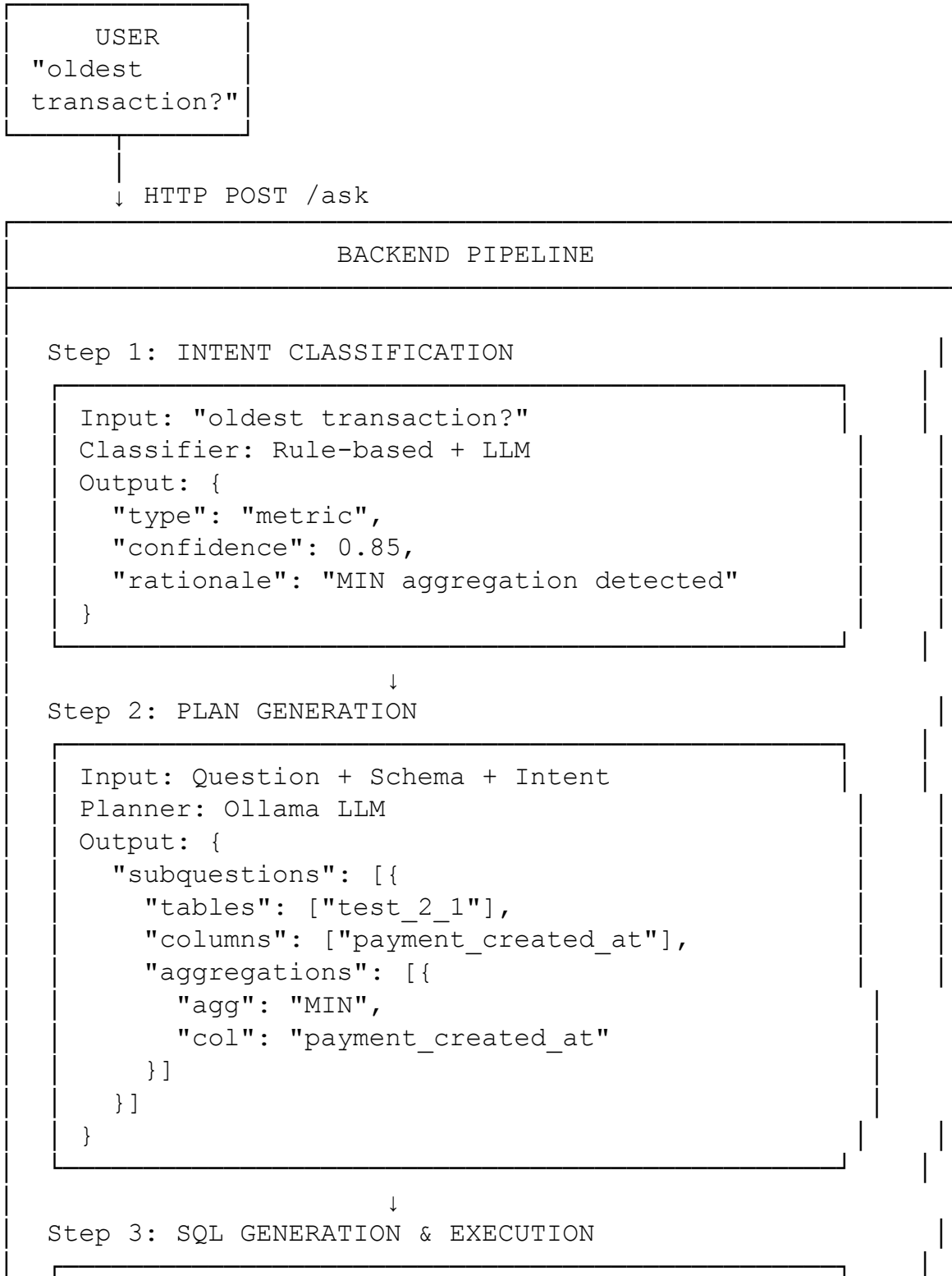
- ☒ **Data Quality Detection**
  - NULL percentage analysis
  - Timestamp validation tests
  - VARCHAR → TIMESTAMP casting
  - NULLS LAST ordering in GROUP BY
- ☒ **Comprehensive Test Suite**
  - Timestamp aggregation tests (MIN/MAX with NULLs)
  - Time grouping tests (month bucketing)
  - Data quality validation tests
  - Aggregation pattern tests
  - 9 tests covering edge cases
- ☒ **Frontend NULL Handling**
  - Display “(No data)” for NULL values
  - Gray italic styling for missing data
  - Chart transformation to label NULLs
  - Table cell special rendering

## Phase 7: Developer Experience

- ☒ **Unified Startup Script**
  - Single `run.sh` to start both servers
  - Automatic dependency installation
  - Process cleanup on exit
  - Log file redirection
- ☒ **Documentation**
  - Architecture diagrams
  - API documentation
  - Component breakdown
  - Usage examples

## System Flow Diagrams

### 1. End-to-End Query Flow



Generated SQL:

```
SELECT MIN(TRY_CAST("test_2_1"."payment_
    created_at" AS TIMESTAMP)) AS
    "min_payment_created_at"
FROM "test_2_1"
```

DuckDB Execution:

```
Result: [('2025-05-22 22:30:57.548+00:00',)]
```

↓

Step 4: NARRATION

Input: Question + Intent + Results

Narrator: Ollama LLM

Output: "The oldest transaction was created  
on May 22, 2025 at 22:30:57 UTC."

↓

Step 5: VISUALIZATION HINTS

Analysis: 1 row, 1 column, timestamp type

```
Hint: {
    "display_hint": "number",
    "units": null
}
```

↓ JSON Response

FRONTEND

Render Components:

Answer Card

"The oldest transaction was created on..."

Visualization (Number Display)

2025-05-22 22:30:57.548000+00:00

min\_payment\_created\_at

View Details (Collapsed)

- Intent: metric (85% confidence)
- Plan: 1 subquestion, 1 table

- SQL: `SELECT MIN(TRY_CAST...`
- Raw Data: `[{"min_payment_created_at": ...}]`

## 2. Data Quality Flow

### DATA QUALITY DETECTION

#### Step 1: IDENTIFY NULL DATA

```
SQL: SELECT COUNT(*),  
        COUNT(payment_created_at),  
        COUNT(*) - COUNT(payment_created_at)  
        FROM test_2_1
```

Result:

Total: 76,583

Non-NULL: 1,715

NULL: 74,868

NULL %: 97.8%

↓

#### Step 2: HANDLE IN SQL GENERATION

- Use `TRY_CAST` for `VARCHAR→TIMESTAMP`
- Add `NULLS LAST` to `ORDER BY`
- Keep `NULL` in `GROUP BY` (separate bucket)

Example:

```
SELECT  
    date_trunc('month',  
        TRY_CAST(payment_created_at  
            AS TIMESTAMP)) as month,  
    COUNT(*) as count  
FROM test_2_1  
GROUP BY month  
ORDER BY month NULLS LAST
```

↓

#### Step 3: VISUALIZE WITH NULL BUCKET

Chart Data:

- Month 1: 1,152 records
- Month 5: 12 records



- Month 9: 109 records
- ...
- (No data): 74,868 records ← Majority!

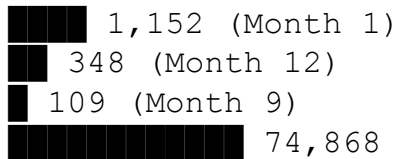
Frontend Transformation:

- NULL → "(No data)"
- Gray italic styling
- Separate bar in chart



Step 4: SURFACE IN UI

Bar Chart: Transactions by Month



Raw Data Tab Shows:

- "month": null
- "count": 74868

### 3. Interactive Visualization Flow

INTERACTIVE TABLE FEATURES

Initial Load (100 rows from DB)



```
useTableControls Hook
• currentPage = 1
• pageSize = 10
• sortColumn = null
• sortDirection = 'asc'
• filters = {}
```



USER INTERACTIONS:



```
1. FILTER: User types "Alice" in customer
   → filters = { customer: "Alice" }
```

```
→ filteredData = data.filter(...)
→ currentPage reset to 1
```

↓

```
2. SORT: User clicks "amount" header
→ sortColumn = "amount"
→ sortDirection = "asc"
→ sortedData = [...filteredData].sort()
```

↓

```
3. PAGINATE: Show 10 rows at a time
→ start = (currentPage - 1) * pageSize
→ end = start + pageSize
→ paginatedData = sortedData.slice(...)
```

↓

4. RENDER TABLE

customer ▲	amount	date
[Search..]	[Search]	[Search]
Alice	\$1,000	2025-01-15
Alice	\$2,000	2025-02-20
...	...	...

Showing 1-10 of 15      [Export CSV]  
[<<] [<] Page 1 of 2 [>] [>>]

↓

```
5. EXPORT: User clicks "Export CSV"
→ exportToCSV(allFilteredData, ...)
→ Browser downloads filtered_data.csv
```

---

## What We Showcase

### 1. Natural Language to SQL Translation

- **Challenge:** Bridge the gap between human language and database queries
- **Our Approach:** Multi-stage pipeline with intent classification, semantic planning, and validated SQL generation
- **Result:** Users can ask “Show me revenue by customer” without knowing SQL syntax

## 2. Full Transparency & Explainability

- **Challenge:** Black-box AI systems lack trust
- **Our Approach:** Every query shows its reasoning:
  - Intent classification with confidence scores
  - Complete query plan with subquestions
  - Generated SQL visible to users
  - Raw data tab for verification
- **Result:** Users understand how answers were derived and can debug issues

## 3. Smart Data Quality Handling

- **Challenge:** Real-world data is messy (97.8% NULL timestamps in our case!)
- **Our Approach:**
  - Don't filter out—show NULLs as “(No data)” buckets
  - Surface data quality issues visually
  - Add tests to detect and report problems
- **Result:** Users see the true state of their data and can take action

## 4. Interactive Visualizations

- **Challenge:** Static charts don't allow exploration
- **Our Approach:** Full interactivity:
  - Zoom/pan on charts
  - Sort/filter/paginate tables
  - Export data in multiple formats
  - Toggle series visibility
- **Result:** Users can explore data beyond the initial answer

## 5. Type Safety & Error Recovery

- **Challenge:** Data types in Excel/CSV are unreliable (VARCHAR timestamps)
- **Our Approach:**
  - Intelligent type casting (TRY\_CAST)
  - Auto-fallback to string mode on ingestion errors
  - Semantic column detection (timestamps, currency)
- **Result:** Robust system that handles messy real-world data

## 6. Local-First Architecture

- **Challenge:** Privacy and latency concerns with cloud LLMs
- **Our Approach:**
  - Local Ollama for LLM inference
  - DuckDB for in-process analytics
  - No data leaves the machine
- **Result:** Fast, private, and offline-capable

## 7. Comparison Intelligence

- **Challenge:** Period-over-period analysis requires domain logic
- **Our Approach:**
  - Automatic detection of comparison queries
  - Normalized comparison structure
  - Delta and percentage calculation
  - Visual comparison cards
- **Result:** Users get instant insights on trends

# Solution Differentiation

## How HaikuGraph is Different

TRADITIONAL APPROACHES
<div>1. Manual SQL Writing<ul style="list-style-type: none"><li>• Users must know SQL syntax</li><li>• Slow iteration for exploratory analysis</li><li>• Error-prone for complex queries</li></ul></div> <div>2. BI Tools (Tableau, Power BI)<ul style="list-style-type: none"><li>• Drag-and-drop interface</li><li>• Requires upfront data modeling</li><li>• Limited natural language support</li><li>• Expensive licensing</li></ul></div> <div>3. Cloud AI SQL Assistants (BigQuery, Snowflake)<ul style="list-style-type: none"><li>• Natural language interface ☐</li><li>• But: Data must be in cloud</li><li>• But: Black-box reasoning</li><li>• But: Pay-per-query</li></ul></div>

VS

HAIKUGRAPH
<div>☐ Natural Language Interface<ul style="list-style-type: none"><li>• Plain English questions</li><li>• No SQL knowledge required</li><li>• Instant results</li></ul></div> <div>☐ Full Transparency<ul style="list-style-type: none"><li>• See intent classification</li><li>• View query plan</li><li>• Inspect generated SQL</li><li>• Verify raw data</li></ul></div> <div>☐ Local &amp; Private<ul style="list-style-type: none"><li>• Data never leaves your machine</li><li>• Local LLM (Ollama)</li><li>• In-process database (DuckDB)</li></ul></div>

- Offline-capable
- Smart Data Quality
  - Surfaces NULL/missing data
  - Handles messy real-world data
  - Type inference and casting
  - Test-driven quality checks
- Interactive Exploration
  - Zoom/filter/sort on the fly
  - Export to CSV/SVG
  - Pagination for large datasets
  - Real-time data transformation
- No Upfront Modeling
  - Auto-ingest from Excel/CSV
  - Auto-generate data cards
  - Schema-on-read

## Unique Selling Points

1. **Transparency-First Design**
    - Unlike black-box AI tools, every decision is explainable
    - Users see the “why” behind every answer
    - Builds trust through visibility
  2. **Data Quality as First-Class Citizen**
    - Most tools hide NULLs; we surface them
    - Comprehensive test suite for edge cases
    - Visual indicators for missing data
  3. **Local-First with Enterprise Potential**
    - Start local for privacy/speed
    - Architecture supports scale-out (see ARCHITECTURE.md)
    - No vendor lock-in
  4. **Developer-Friendly**
    - Full test coverage
    - Clear separation of concerns
    - Easy to extend and customize
  5. **Real-World Data Tolerance**
    - Handles VARCHAR timestamps
    - Auto-fallback on type errors
    - Graceful degradation
- 

## Pros and Cons Analysis

### □ PROS

#### Technical Strengths

1. **Clean Architecture**

- Modular pipeline (Intent □ Plan □ Execute □ Narrate)
  - Same backend for CLI and UI
  - Easy to test and maintain
  - Clear separation of concerns
- 2. Rich Type System**
    - Pydantic models for validation
    - Type hints throughout Python code
    - React PropTypes for components
    - Catches errors early
  - 3. Comprehensive Testing**
    - 9 data quality tests
    - Edge case coverage (NULL handling)
    - Timestamp validation
    - Aggregation pattern tests
  - 4. Performance**
    - In-process DuckDB (no network overhead)
    - Local LLM (200-500ms latency)
    - Efficient SQL generation
    - Preview-based pagination
  - 5. User Experience**
    - Intuitive natural language interface
    - Rich visualizations auto-generated
    - Full transparency (intent/plan/SQL)
    - Interactive exploration tools

## Business Strengths

- 1. Privacy & Security**
  - All data stays local
  - No cloud dependencies
  - GDPR/HIPAA friendly
  - Audit trail built-in
- 2. Cost Effective**
  - No per-query fees
  - No data egress charges
  - Open-source components
  - Self-hosted option
- 3. Rapid Setup**
  - Single command data ingestion
  - Auto-schema detection
  - No upfront modeling required
  - Works with existing Excel/CSV files

## □ CONS

### Technical Limitations

- 1. Scalability Constraints**
  - In-process DuckDB limits concurrent users (1-5)
  - Single LLM instance (no load balancing)
  - No caching layer in POC
  - Not tested with >100GB datasets
  - **Mitigation:** Architecture supports scale-out (see ARCHITECTURE.md production design)
- 2. LLM Dependency**
  - Requires Ollama installation
  - LLM accuracy not 100%
  - Latency varies by query complexity
  - May hallucinate on ambiguous questions

- **Mitigation:** Fallback to simpler models, manual SQL override
3. **Limited Query Complexity**
    - No subqueries in POC
    - Complex JOINS not fully tested
    - Window functions not supported
    - No recursive CTEs
    - **Mitigation:** Roadmap includes advanced SQL features
  4. **Data Type Handling**
    - All ingested as VARCHAR (DuckDB limitation)
    - Requires TRY\_CAST at query time
    - Date parsing can be ambiguous
    - No schema enforcement
    - **Mitigation:** Future: support proper types on ingestion
  5. **Error Recovery**
    - LLM failures stop the pipeline
    - No retry logic for transient errors
    - User must rephrase on failure
    - **Mitigation:** Add retry with exponential backoff

## Business Limitations

1. **Enterprise Features Missing**
  - No authentication/authorization
  - No multi-tenancy
  - No audit logs
  - No role-based access control
  - **Mitigation:** Roadmap includes enterprise features
2. **Integration Gaps**
  - Only supports Excel/CSV (no databases)
  - No real-time data sources
  - No API for external apps
  - No Slack/Teams bot
  - **Mitigation:** API-first design allows easy integration
3. **Operational Complexity**
  - Requires Ollama setup
  - Manual model management
  - No monitoring dashboard
  - No alerting system
  - **Mitigation:** Docker image for easy deployment
4. **Documentation Gaps**
  - No video tutorials
  - Limited troubleshooting guide
  - No architecture decision records (ADRs)
  - **Mitigation:** Comprehensive markdown docs added
5. **Community & Support**
  - No community forum
  - No commercial support option
  - Single maintainer risk
  - **Mitigation:** Open-source model encourages contributions

## □ TRADE-OFFS MADE

1. **Local LLM vs Cloud LLM**
  - **Chose:** Local (Ollama)
  - **Why:** Privacy, cost, offline capability
  - **Cost:** Slightly lower accuracy, higher setup complexity
2. **In-Process DB vs Client-Server**

- **Chose:** In-process (DuckDB)
  - **Why:** Simplicity, performance for analytics
  - **Cost:** Limited concurrency, no true multi-user
3. **Transparency vs Simplicity**
- **Chose:** Transparency (show all reasoning)
  - **Why:** Build trust, enable debugging
  - **Cost:** More UI complexity, potential information overload
4. **Auto-Infer vs User-Define Schema**
- **Chose:** Auto-infer
  - **Why:** Faster setup, less user burden
  - **Cost:** Less control, potential type mismatches
5. **React vs Server-Side Rendering**
- **Chose:** React SPA
  - **Why:** Rich interactivity, modern developer experience
  - **Cost:** Larger bundle size, SEO challenges
- 

## Future Roadmap

### Short Term (1-3 months)

- ☐ Add authentication (OAuth2)
- ☐ Implement caching (Redis)
- ☐ Support PostgreSQL/MySQL ingestion
- ☐ Add monitoring dashboard
- ☐ Docker image for easy deployment

### Medium Term (3-6 months)

- ☐ Multi-tenancy support
- ☐ Advanced SQL features (subqueries, CTEs, window functions)
- ☐ Real-time data sources
- ☐ Slack/Teams integration
- ☐ Mobile-responsive UI

### Long Term (6-12 months)

- ☐ Horizontal scaling architecture
  - ☐ Cloud-native deployment option
  - ☐ Enterprise features (RBAC, audit logs)
  - ☐ ML-based query optimization
  - ☐ Natural language follow-ups
- 

## Conclusion

HaikuGraph demonstrates a **transparency-first, local-first, data-quality-aware** approach to natural language data querying. By combining modern LLMs with robust SQL generation and interactive visualizations, we provide an intuitive yet powerful tool for data exploration.

Our key differentiators are: 1. **Full transparency** into AI reasoning 2. **Local-first** for privacy and speed 3. **Data quality** as a first-class concern 4. **Interactive exploration** beyond static answers

While we have limitations (scalability, LLM dependency, enterprise features), our architecture is designed for growth, and our open approach encourages community contributions.

**HaikuGraph is not just a query tool—it's a data conversation platform.**



---

## Getting Started

*# 1. Ingest your data*

```
haikugraph ingest --data-dir ./data
```

*# 2. Start the application*

```
./run.sh
```

*# 3. Open <http://localhost:5173> and start asking questions!*

For detailed setup instructions, see `QUICK_START.md`.

For API documentation, see `WEB_UI_README.md`.

For test suite details, see `VISUALIZATION_IMPROVEMENTS.md`.