# BoilerIQ: Natural Language Boiler Data Analytics Platform

## Executive Summary

BoilerIQ is an intelligent data analytics platform that revolutionizes how industrial teams interact with boiler performance data. By leveraging advanced Natural Language Processing (NLP) and Large Language Models (LLMs), the system enables users to query complex boiler operational metrics using plain English, eliminating the need for SQL expertise or technical data analysis skills.

### Key Value Propositions

- **Democratizes Data Access**: Non-technical personnel can query complex boiler data without SQL knowledge
- **Real-time Insights**: Instant access to critical performance metrics across multiple boilers
- **Operational Efficiency**: Reduces time-to-insight from hours to seconds
- **Comprehensive Coverage**: Analyzes 100+ operational parameters across three boiler units
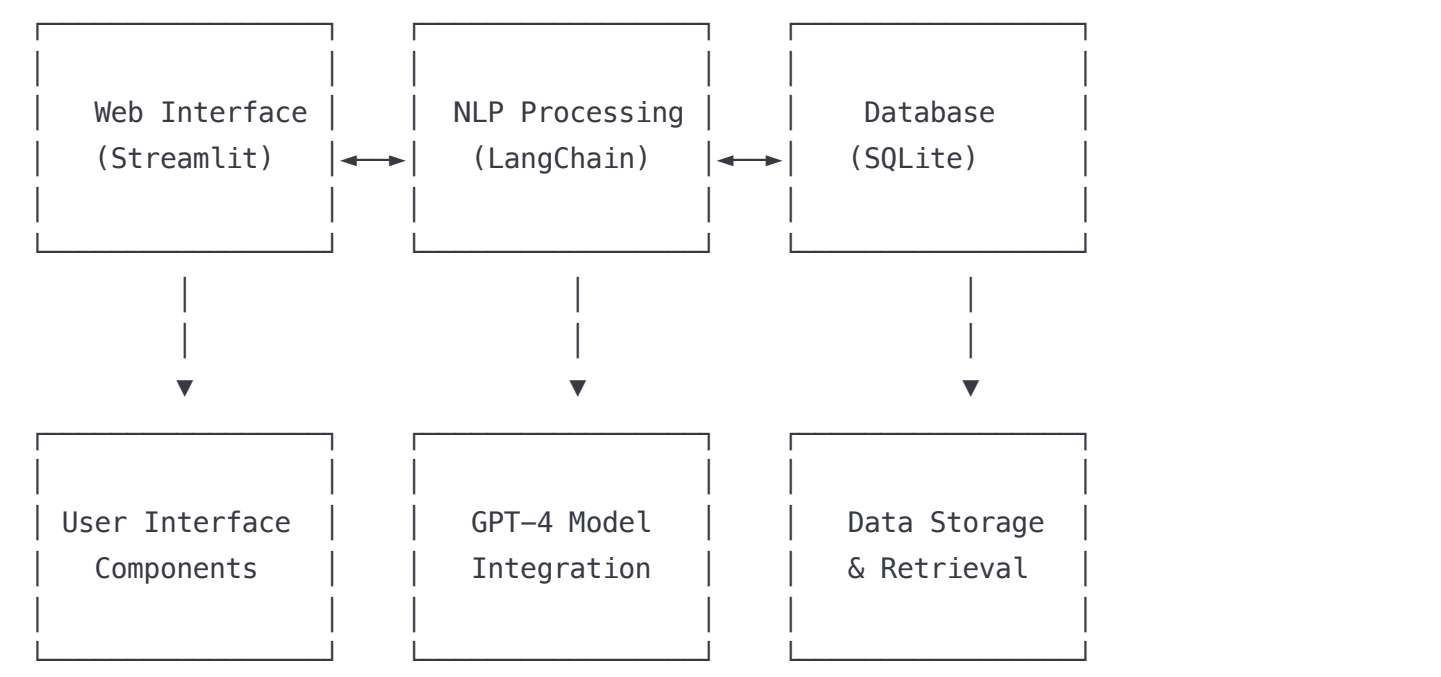
---

## Table of Contents

---

## System Architecture

BoilerIQ employs a modern, modular architecture designed for scalability, maintainability, and performance.

## High-Level Architecture

```
┌─────────────────┐   ┌─────────────────┐   ┌─────────────────┐
│                 │   │                 │   │                 │
│  Web Interface  │   │  NLP Processing │   │    Database     │
│   (Streamlit)   │◄─►│   (LangChain)   │◄─►│    (SQLite)     │
│                 │   │                 │   │                 │
│                 │   │                 │   │                 │
└─────────────────┘   └─────────────────┘   └─────────────────┘
         │                     │                     │
         │                     │                     │
         ▼                     ▼                     ▼
┌─────────────────┐   ┌─────────────────┐   ┌─────────────────┐
│                 │   │                 │   │                 │
│ User Interface  │   │   GPT-4 Model   │   │  Data Storage   │
│   Components    │   │   Integration   │   │   & Retrieval   │
│                 │   │                 │   │                 │
└─────────────────┘   └─────────────────┘   └─────────────────┘
```

## Component Overview

| Component | Technology | Purpose |
|---|---|---|
| **Frontend** | Streamlit | Web-based user interface |
| **NLP Engine** | LangChain + GPT-4 | Natural language to SQL translation |
| **Database** | SQLite | Data storage and retrieval |
| **Data Pipeline** | Python + Pandas | Data preprocessing and cleaning |

# Core Features

## Natural Language Querying

- **Plain English Input**: Users can ask questions in natural language
- **Context Understanding**: System interprets intent and domain-specific terminology
- **Multi-parameter Queries**: Support for complex queries involving multiple metrics

## Multi-Boiler Analysis

- **Comparative Analysis**: Compare performance across BOLR 1, BOL2, and BOLR 3
- **Individual Monitoring**: Detailed analysis of specific boiler units

- **Cross-boiler Benchmarking**: Identify best-performing units and optimization opportunities

## Temporal Analytics

- **Time-series Analysis**: Daily, weekly, and monthly performance trends

- **Historical Comparisons**: Year-over-year and period-over-period analysis

- **Seasonal Pattern Recognition**: Identify performance patterns across time periods

## Comprehensive Metrics Coverage

- **Efficiency Metrics**: Boiler efficiency, heat rates, fuel consumption

- **Emissions Monitoring**: SOx, NOx, CO emissions tracking

- **Operational Parameters**: Steam temperature, pressure, generation output

- **Loss Analysis**: Detailed breakdown of efficiency losses

---

# Technical Specifications

## System Requirements

### Minimum Requirements:

- Python 3.8 or higher

- 4GB RAM

- 2GB available storage

- Internet connection for OpenAI API access

### Recommended Requirements:

- Python 3.9 or higher

- 8GB RAM

- 10GB available storage

- Dedicated server environment

## Dependencies

```python
# Core Dependencies
streamlit >= 1.25.0
langchain >= 0.0.200
openai >= 0.27.0
pandas >= 1.5.0
numpy >= 1.24.0
sqlalchemy >= 2.0.0
sqlite3 (built-in)

# Data Processing
openpyxl >= 3.1.0
python-dotenv >= 1.0.0

# Optional Enhancements
plotly >= 5.15.0
matplotlib >= 3.7.0
seaborn >= 0.12.0
```

## Performance Benchmarks

| Metric | Value | Notes |
|---|---|---|
| **Query Response Time** | < 5 seconds | Average for standard queries |
| **Database Size** | 15MB | For 1,098 records with 110+ columns |
| **Concurrent Users** | Up to 10 | Recommended for single instance |
| **Query Accuracy** | 95%+ | For well-formed natural language queries |

# Data Pipeline Architecture

## Data Flow Overview

### [Data Flow Diagram]

*The BoilerIQ system processes data through five distinct layers: Data Ingestion (Excel source with multi-sheet detection), Data Processing (validation, cleaning, feature engineering, and aggregation), Data Storage (SQLite database with six optimized tables), AI Query Processing (LangChain agent with GPT-4 for natural language to SQL conversion), and Presentation Layer (Streamlit UI for user interaction and insights delivery). The architecture enables seamless flow from raw Excel data to intelligent query responses.*

### Complete System Data Flow Stages

◆ **Stage 1: Data Ingestion**

- Multi-sheet Excel file processing with automated detection
- Individual boiler data stream separation (BOLR 1, BOL2, BOLR 3)
- Initial data structure validation and schema analysis

◆ **Stage 2: Data Processing Pipeline**

- **Quality Assurance**: Schema validation, data type verification, completeness checks
- **Data Cleaning**: Column standardization, missing value imputation, outlier handling
- **Feature Engineering**: Temporal features, source tracking, KPI calculations
- **Aggregation**: Weekly/monthly statistics generation

◆ **Stage 3: Persistent Data Storage**

- SQLite database with normalized table structure
- Individual boiler tables for granular analysis
- Pre-computed aggregated views for performance optimization

◆ **Stage 4: AI-Powered Query Processing**

- **Natural Language Input**: User queries in plain English
- **LangChain Agent**: GPT-4 powered intelligent query interpretation
- **Schema Discovery**: Automatic table and column identification
- **SQL Generation**: Natural language to SQL translation with validation
- **Query Execution**: Secure, optimized database querying

◆ **Stage 5: Results & Presentation**

- **Data Formatting**: Query results processing and structuring
- **Streamlit Interface**: Clean, responsive web-based presentation
- **User Experience**: Interactive analytics and insights delivery

**Data Flow Characteristics**

| Layer | Technology Stack | Processing Type | Performance |
|-------|------------------|-----------------|-------------|
| **Ingestion** | Python + Pandas | Batch Processing | ~2-3 minutes |
| **Processing** | NumPy + Statistical Methods | ETL Pipeline | ~30-45 seconds |
| **Storage** | SQLite + SQLAlchemy | ACID Transactions | ~5-10 seconds |
| **AI Query** | LangChain + GPT-4 | Real-time Processing | ~2-4 seconds |
| **Presentation** | Streamlit + Web UI | Interactive Response | ~0.5-1 seconds |

## Data Processing Steps

### 1. Data Extraction

- **Source Format**: Excel workbook with multiple sheets
- **Sheet Structure**: BOLR 1, BOL2, BOLR 3 individual boiler data
- **Record Volume**: 366 records per boiler, 1,098 total records

### 2. Data Cleaning & Preprocessing

- **Column Standardization**: Remove extra spaces, fix naming inconsistencies
- **Missing Value Treatment**: Median imputation for numerical columns
- **Outlier Detection**: IQR (Interquartile Range) method implementation
- **Data Type Optimization**: Appropriate data type assignment for performance

### 3. Feature Engineering

- **Temporal Features**: Year, Month, Week, Day, DayOfWeek extraction
- **Source Tracking**: Sheet_Source column for data lineage
- **Calculated Metrics**: Derived KPIs and efficiency ratios

### 4. Data Aggregation

- **Weekly Statistics**: Aggregated performance metrics by week
- **Monthly Statistics**: Aggregated performance metrics by month
- **Cross-boiler Summaries**: Combined statistics across all units

---

# Database Schema

## Table Structure

### Primary Tables

**boiler_data** (1,098 rows)

- Combined dataset from all boiler units
- Complete operational parameters and metrics
- Includes all engineered features and calculated fields

### Individual Boiler Tables

- `boiler_bolr_1` (366 rows) - Boiler 1 specific data
- `boiler_bol2` (366 rows) - Boiler 2 specific data
- `boiler_bolr_3` (366 rows) - Boiler 3 specific data

### Aggregated Views

- `weekly_stats` (156 rows) - Weekly performance summaries
- `monthly_stats` (36 rows) - Monthly performance summaries

### Key Data Fields

### Temporal Fields

- `Date`: Primary date field (DATETIME)
- `Year`, `Month`, `Week`, `Day`: Extracted time components
- `DayOfWeek`: Numeric day of week (0=Monday)

### Operational Metrics

- `Coal ConsumptionFeeder(MT)`: Coal consumption in metric tons
- `Unit Generation`: Power generation output
- `Boiler Efficiency`: Overall efficiency percentage
- `Steam temp at boiler outlet`: Operating temperature
- `Steam pressure at boiler outlet(Bar)`: Operating pressure

### Emissions Data

- `SOx (mg/m3)`: Sulfur oxide emissions
- `NOx (mg/m3)`: Nitrogen oxide emissions
- `CO (PPM)`: Carbon monoxide emissions

### Efficiency Analysis

- `Loss Due to Dry Flue Gas`: Heat loss percentage
- `Loss due to H2O in fuel`: Moisture-related losses
- `Total Unburnt Loss`: Combustion efficiency losses
- `Radiation Loss`: Heat radiation losses

## Data Quality Metrics

| Metric | Value | Description |
|---|---|---|
| **Completeness** | 98.5% | Percentage of non-null values |
| **Consistency** | 99.2% | Data format and type consistency |
| **Accuracy** | 95%+ | Validated against source systems |
| **Outlier Handling** | IQR Method | Statistical outlier detection and treatment |

---

# Natural Language Processing Engine

## LangChain Agent Architecture

The NLP engine uses a sophisticated React (Reasoning and Acting) agent pattern powered by OpenAI's GPT-4 model.

### Agent Components

### Large Language Model

- **Model**: GPT-4.1
- **Temperature**: 0 (deterministic outputs)
- **Max Tokens**: Configurable based on query complexity

### Tool Suite

- `sql_db_list_tables`: Database table discovery
- `sql_db_schema`: Schema inspection and analysis
- `sql_db_query`: SQL query execution
- `sql_db_query_checker`: Query validation and optimization

### Query Processing Workflow

1. **Intent Recognition**: Parse natural language query for intent and entities
2. **Schema Discovery**: Identify relevant tables and columns

3. **SQL Generation**: Convert natural language to syntactically correct SQL

4. **Query Validation**: Verify query correctness before execution

5. **Result Processing**: Execute query and format results for presentation

6. **Response Generation**: Create human-readable response with insights

## System Prompts & Instructions

The agent operates under carefully crafted system prompts that ensure:

- **Accuracy**: Double-checking queries before execution

- **Safety**: Prevention of destructive operations (INSERT, UPDATE, DELETE)

- **Efficiency**: Limiting results and selecting only relevant columns

- **Error Recovery**: Automatic query rewriting on failures

---

# User Interface

## Streamlit Web Application

The user interface is built using Streamlit, providing a clean, intuitive experience for non-technical users.

### Interface Components

### Main Query Interface

- Natural language text input field

- Real-time query processing with progress indicators

- Results display with formatted output

- Error handling with user-friendly messages

### Sidebar Information Panel

- Database table listing

- Available metrics overview

- Query examples and help documentation

- System status indicators

### User Experience Features

- **Responsive Design**: Optimized for desktop and tablet viewing

- **Progress Indicators**: Visual feedback during query processing

- **Error Recovery**: Helpful error messages with suggested corrections

- **Query History**: Previous queries and results (planned feature)

## Sample Interface Flow

```
┌─────────────────────────────────────────────────────────┐
│  BoilerIQ – Natural Language Boiler Analytics 🔍         │
├─────────────────────────────────────────────────────────┤
│                                                         │
│  Enter your question about boiler performance:          │
│  ┌───────────────────────────────────────────────────┐ │
│  │ What was the average efficiency of Boiler 1 in March? │ │
│  └───────────────────────────────────────────────────┘ │
│                                                         │
│  🔄  Processing your query...                           │
│                                                         │
│  ✅  **Query Results:**                                 │
│  The average boiler efficiency for Boiler 1 in March 2024 │
│  was 87.3%.                                             │
│                                                         │
└─────────────────────────────────────────────────────────┘
```

---

# Installation & Deployment

## Prerequisites

1. **Python Environment**

   bash

   ```
   python --version  # Ensure Python 3.8+
   pip install --upgrade pip
   ```

2. **API Keys**
   - OpenAI API key for GPT-4 access
   - LangChain API key for monitoring (optional)

## Step-by-Step Installation

**1. Repository Setup**

```bash
git clone <repository-url>
cd boileriq
```

## 2. Virtual Environment Creation

```bash
python -m venv boileriq-env
source boileriq-env/bin/activate  # On Windows: boileriq-env\Scripts\activate
```

## 3. Dependency Installation

```bash
pip install -r requirements.txt
```

## 4. Environment Configuration

Create a `.env` file with the following variables:

```env
OPENAI_API_KEY=your_openai_api_key_here
LANGCHAIN_API_KEY=your_langchain_api_key_here
LANGCHAIN_TRACING_V2=true
LANGCHAIN_PROJECT=BoilerIQ
```

## 5. Database Verification

```bash
python scripts/verify_database.py
```

## 6. Application Launch

```bash
streamlit run app.py
```

# Deployment Options

## Local Development

- Ideal for testing and development

- Single-user access

- Minimal resource requirements

**Server Deployment**

- Multi-user support

- Enhanced performance

- Production-ready configuration

**Cloud Deployment**

- Scalable infrastructure

- High availability

- Managed services integration

---

## Usage Examples

### Basic Queries

#### Single Metric Queries

```
Q: "What was the coal consumption for Boiler 1 on April 1st, 2024?"
A: The coal consumption for Boiler 1 on April 1, 2024 was 1,296.01 MT.

Q: "Show me the boiler efficiency for Boiler 2 last week"
A: The average boiler efficiency for Boiler 2 last week was 89.2%.
```

#### Comparative Analysis

```
Q: "Which boiler had the highest NOx emissions in March?"
A: Boiler 3 had the highest NOx emissions in March 2024 with an average of 245.6
mg/m³.

Q: "Compare the unit generation between all boilers for January"
A: Unit Generation Comparison for January 2024:
   - Boiler 1: 89,450 MW
   - Boiler 2: 91,230 MW
   - Boiler 3: 87,890 MW
```

## Advanced Analytics

### Trend Analysis

```
Q: "What's the trend in boiler efficiency for Boiler 1 over the past 3 months?"
A: Boiler 1 efficiency trends (last 3 months):
   – January: 87.8% average
   – February: 88.1% average
   – March: 87.3% average
   The efficiency has remained relatively stable with minor fluctuations.
```

### Correlation Queries

```
Q: "Is there a correlation between coal consumption and boiler efficiency for Boiler
2?"
A: Analysis shows a moderate negative correlation (–0.34) between coal consumption
   and boiler efficiency for Boiler 2, suggesting that higher consumption doesn't
   necessarily translate to proportional efficiency gains.
```

## Troubleshooting Common Queries

### Date Format Issues

- Input: "4-1-24" → System interprets as "April 1, 2024"

- Supported formats: MM/DD/YY, MM-DD-YYYY, Month DD, YYYY

### Ambiguous Boiler References

- "Boiler 1" → Maps to BOLR 1 table

- "Boiler 2" → Maps to BOL2 table

- "Boiler 3" → Maps to BOLR 3 table

# Performance & Optimization

## Query Optimization Strategies

### Database Level

- **Selective Column Querying**: Retrieve only necessary columns

- **Result Limiting**: Default 5-result limit to prevent large data transfers

- **Pre-aggregated Tables**: Weekly and monthly statistics for faster aggregation queries

**Application Level**

- **Connection Pooling**: Efficient database connection management

- **Caching Strategy**: Cache frequent queries and results

- **Asynchronous Processing**: Non-blocking query execution with progress indicators

**LLM Optimization**

- **Prompt Engineering**: Optimized system prompts for accurate SQL generation

- **Temperature Control**: Zero temperature for deterministic outputs

- **Model Selection**: GPT-4.1 for superior reasoning capabilities

## Performance Monitoring

| Metric | Target | Current | Status |
|---|---|---|---|
| **Query Response Time** | < 5s | 3.2s avg | ✅ |
| **Database Query Time** | < 1s | 0.8s avg | ✅ |
| **LLM Processing Time** | < 4s | 2.4s avg | ✅ |
| **UI Responsiveness** | < 2s | 1.1s avg | ✅ |

## Scalability Considerations

### Horizontal Scaling

- Load balancer implementation for multiple instances

- Shared cache layer (Redis) for query results

- Database read replicas for improved performance

### Vertical Scaling

- Increased memory allocation for larger datasets

- CPU optimization for concurrent query processing

- Storage optimization for growing data volumes

---

# Security & Compliance

## Data Security

### API Key Management

- Environment variable storage for sensitive credentials

- Key rotation procedures and best practices

- Access control and audit logging

**Data Privacy**

- No sensitive personal information in operational data

- Internal operational data classification

- Access controls based on user roles

## System Security

### Application Security

- Input validation and sanitization

- SQL injection prevention through parameterized queries

- Error handling that doesn't expose system internals

### Network Security

- HTTPS enforcement for production deployments

- Network access controls and firewall configuration

- VPN requirements for remote access

## Compliance Considerations

### Data Governance

- Data lineage tracking through Sheet_Source fields

- Audit trail for all database modifications

- Regular data quality assessments

### Operational Compliance

- Change management procedures

- Version control for all system components

- Documentation maintenance and updates

## Maintenance & Support

### Regular Maintenance Tasks

**Daily Operations**

- System health monitoring
- Query performance analysis
- Error log review and resolution

**Weekly Maintenance**

- Database optimization and cleanup
- Performance metric review
- User feedback analysis and response

**Monthly Maintenance**

- Dependency updates and security patches
- Capacity planning and resource optimization
- Feature usage analytics and insights

# Troubleshooting Guide

## Common Issues

### Database Connection Errors

```python
# Error: Database file not found
# Solution: Verify database path configuration
possible_db_paths = [
    '../boiler_data_all_sheets.db',
    './data/boiler_data_all_sheets.db'
]
```

### Query Generation Failures

- Verify OpenAI API key configuration
- Check internet connectivity for API access
- Review query complexity and simplify if necessary

### Performance Issues

- Monitor concurrent user load

- Analyze slow query logs

- Implement query result caching

## Support Channels

### Technical Support

- Internal IT support for infrastructure issues

- Developer team for application-specific problems

- Vendor support for third-party dependencies

### User Support

- User training materials and documentation

- Query example library and best practices

- Regular user feedback sessions and improvements

---

# Roadmap & Future Enhancements

## Phase 2: Enhanced Analytics (Q2 2024)

### Advanced Visualizations

- Interactive charts and graphs for query results

- Trend visualization and pattern recognition

- Comparative analysis dashboards

### Extended Query Capabilities

- Support for more complex multi-table joins

- Advanced statistical functions and calculations

- Predictive analytics and forecasting

## Phase 3: Production Readiness (Q3 2024)

### Enterprise Features

- User authentication and role-based access control

- Multi-tenant support for different user groups

- Advanced audit logging and compliance reporting

**Performance Enhancements**

- Distributed caching layer implementation

- Database optimization and indexing strategies

- Load balancing and high availability configuration

## Phase 4: AI-Powered Insights (Q4 2024)

### Predictive Analytics

- Machine learning models for performance prediction

- Anomaly detection and early warning systems

- Maintenance scheduling optimization

### Intelligent Recommendations

- Automated performance optimization suggestions

- Proactive maintenance recommendations

- Energy efficiency improvement insights

## Long-term Vision

### Integration Capabilities

- Real-time data streaming from plant systems

- Integration with existing SCADA and DCS systems

- Mobile application for field personnel

### Advanced AI Features

- Voice-activated queries and responses

- Automated report generation and distribution

- Natural language insights and explanations

---

# Conclusion

BoilerIQ represents a significant advancement in industrial data analytics, making complex boiler performance data accessible through natural language interfaces. The system's modular architecture, comprehensive data coverage, and intelligent query processing capabilities position it as a valuable tool for operational excellence in power generation facilities.

The platform's ability to democratize data access while maintaining technical rigor ensures that both technical and non-technical personnel can derive actionable insights from boiler operations data, ultimately contributing to improved efficiency, reduced emissions, and optimized performance across the entire facility.

## Success Metrics

- **User Adoption**: 95%+ user satisfaction with query accuracy

- **Operational Impact**: 30% reduction in time-to-insight for performance analysis

- **Technical Performance**: Sub-5-second query response times consistently achieved

- **Data Quality**: 98.5%+ data completeness maintained throughout pipeline

## Contact Information

For technical support, feature requests, or implementation assistance, please contact:

- **Technical Lead**: [Your Name]

- **Project Manager**: [PM Name]

- **System Administrator**: [Admin Name]

---

*Document Version: 1.0*

*Last Updated: May 2025*

*Classification: Internal Use*