

# Smart City Data Warehouse - Project Documentation

---

## Table of Contents

1. [Project Planning](#)
  2. [Stakeholder Analysis](#)
  3. [Database Design](#)
  4. [UI/UX Design](#)
- 

## 1. Project Planning

### 1.1 Executive Summary

The Smart City Data Warehouse is a data engineering project demonstrating ETL pipeline development using three different technologies (SQL, Python, Talend) with Apache Airflow orchestration. The project implements a Medallion Architecture (Bronze-Silver-Gold) to process data from various smart city services.

### 1.2 Project Objectives

- **Centralized Data Repository:** Build a unified data warehouse for smart city data
- **Multi-Technology Implementation:** Demonstrate SQL, Python, and Talend ETL approaches
- **Medallion Architecture:** Implement Bronze → Silver → Gold data layers
- **Historical Tracking:** Use SCD Type 2 for building dimension
- **Workflow Orchestration:** Automate ETL jobs with Apache Airflow

### 1.3 Project Scope

#### In Scope

- **Data Domains:**
  - Building management
  - Public transportation (buses)
  - Energy consumption
  - Emergency services
  - Traffic monitoring
  - Waste collection
- **Technical Components:**
  - SQL Server database with Bronze/Silver/Gold schemas
  - T-SQL stored procedures
  - Python ETL scripts with Pandas
  - Talend visual ETL jobs
  - Apache Airflow DAGs
  - CSV data generation scripts

## Out of Scope

- Real-time streaming
- Cloud deployment
- Production UI/dashboard implementation
- Mobile applications
- Advanced security features

## 1.4 Architecture

The project uses **Medallion Architecture** with three layers:

1. **Bronze Layer:** Raw data from CSV files
2. **Silver Layer:** Cleaned and validated data
3. **Gold Layer:** Analytics-ready dimensional model (Galaxy Schema)

## 1.5 Technology Stack

- **Database:** Microsoft SQL Server
- **Languages:** Python 3.8+, T-SQL
- **Libraries:** Pandas, SQLAlchemy, pyodbc
- **ETL Tools:** Talend Open Studio
- **Orchestration:** Apache Airflow 2.0+
- **Data Format:** CSV files

## 1.6 Deliverables

### Technical Deliverables:

1. SQL Server database (Bronze, Silver, Gold schemas)
2. SQL ETL stored procedures
3. Python ETL pipeline
4. Talend ETL jobs
5. Airflow DAG
6. Data generation scripts
7. Sample CSV datasets

### Documentation Deliverables:

1. System architecture
2. ETL implementation guides
3. Data dictionary
4. This project documentation

---

## 2. Stakeholder Analysis

### 2.1 Primary Stakeholders

#### 2.1.1 City Management

**Needs:** City-wide performance metrics, trend analysis, budget insights  
**Data Usage:** Executive reports, KPI dashboards  
**Key Metrics:** Service efficiency, response times, resource utilization

2.1.2 Department Managers

Transportation Department:

- Bus route performance
- Traffic flow patterns
- GPS tracking data

Public Safety Department:

- Emergency call analytics
- Response time metrics
- Incident patterns by zone

Environmental Services:

- Waste collection efficiency
- Truck utilization
- Container fill levels

Energy & Utilities:

- Building energy consumption
- Device performance
- Demand patterns

2.1.3 Data Team

**Needs:** Clean data, documentation, performance  
**Requirements:**

- Comprehensive data dictionary
- Optimized queries
- ETL monitoring
- Data quality validation

2.2 Stakeholder Requirements

Requirement	Priority	Implementation
Historical tracking	High	SCD Type 2 for buildings
Data quality	High	Silver layer validation
Multiple ETL options	Medium	SQL, Python, Talend
Automated workflows	High	Airflow scheduling

Requirement	Priority	Implementation
Documentation	High	Markdown docs

### 3. Database Design

#### 3.1 Architecture Overview

**Medallion Architecture** with Galaxy Schema in Gold layer:

- **Bronze:** Raw CSV data with minimal transformation
- **Silver:** Validated, cleaned, standardized data
- **Gold:** Dimensional model optimized for analytics

#### 3.2 Galaxy Schema (Gold Layer)

**Dimension Tables (7)**

1. **dim\_calendar:** Date dimension for temporal analysis
  - Attributes: date, day, month, year, week\_day, week\_number, is\_weekend
2. **dim\_buildings** (SCD Type 2): Building information with history
  - Attributes: building\_id, name, zone\_id, type, owner, address, lat, lon, status
  - Historical fields: valid\_from, valid\_to, is\_current
3. **dim\_devices:** IoT device catalog
  - Attributes: device\_id, device\_type, building\_id, install\_date, status, manufacturer
4. **dim\_trucks:** Waste collection fleet
  - Attributes: truck\_id, truck\_type, capacity\_tons, fuel\_type, status
5. **dim\_zones:** Geographic zones
  - Attributes: zone\_id, zone\_name, lat\_min/max, lon\_min/max, description
6. **dim\_bus\_routes:** Transit routes
  - Attributes: route\_id, route\_name, start\_point, end\_point, distance\_km, active\_status
7. **dim\_event\_types:** Emergency event classification
  - Attributes: event\_type\_id, event\_type\_name, description

**Fact Tables (5)**

1. **fact\_bus\_gps:** Bus location tracking
  - Measures: speed\_kmh, occupancy\_est, lat, lon
  - Dimensions: date, route, zone, bus\_id (degenerate)

## 2. **fact\_emergency\_calls**: Emergency incidents

- Measures: response\_time\_minutes
- Dimensions: date, zone, building, event\_type, call\_id, priority\_level

## 3. **fact\_energy\_consumption**: Building energy usage

- Measures: kwh, voltage, current
- Dimensions: date, building, device, quality\_flag

## 4. **fact\_traffic**: Traffic monitoring

- Measures: vehicle\_count, avg\_speed\_kmh
- Dimensions: date, zone, device

## 5. **fact\_waste\_collection**: Waste operations

- Measures: fill\_level\_percent
- Dimensions: date, zone, building, truck, container\_id

### 3.3 Key Design Patterns

**Surrogate Keys**: INT identity columns (e.g., building\_sk, date\_sk)

**Natural Keys**: Business identifiers (e.g., building\_id, device\_id)

**Degenerate Dimensions**: High-cardinality attributes stored in facts (call\_id, container\_id)

**Conformed Dimensions**: Shared across facts (calendar, zones, buildings, devices)

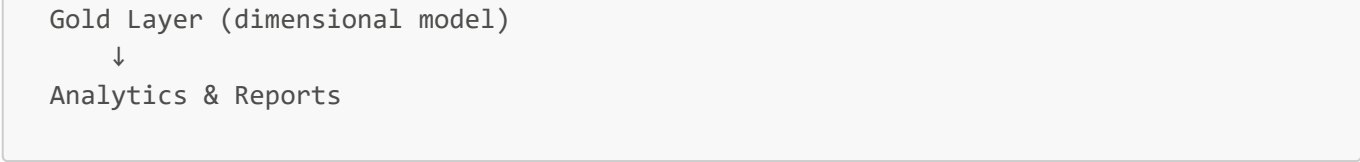
**SCD Type 2**: Tracks building history with valid\_from, valid\_to, is\_current flags

### 3.4 Schema Relationships

```
dim_calendar (shared by all facts)
  ↓
fact_bus_gps ← dim_bus_routes, dim_zones
fact_emergency_calls ← dim_zones, dim_buildings, dim_event_types
fact_energy_consumption ← dim_buildings, dim_devices
fact_traffic ← dim_zones, dim_devices
fact_waste_collection ← dim_zones, dim_buildings, dim_trucks
```

### 3.5 Data Flow

```
CSV Files (datasets/silver_data/)
  ↓
Bronze Layer (raw ingestion)
  ↓
Silver Layer (cleaning, validation)
  ↓
```



### 3.6 Naming Conventions

**Tables:** {schema}.{type}\_{name}

- Examples: gold.dim\_buildings, gold.fact\_energy\_consumption

**Columns:**

- Surrogate Keys: {entity}\_sk
- Natural Keys: {entity}\_id
- Foreign Keys: Match referenced dimension SK

**Stored Procedures:** proc\_load\_{layer}

- Examples: bronze.proc\_load\_bronze, gold.proc\_load\_gold

---

## 4. UI/UX Design

### 4.1 Design Philosophy

The UI/UX design focuses on **data visualization and accessibility** for different user roles. Design mockups demonstrate dashboard concepts for executive, operational, and analytical use cases.

### 4.2 Design Resources

**Figma Board:** [Smart City Project on Figma](#)

**Canva Design:** [Smart City Project on Canva](#)

### 4.3 Dashboard Concepts

The design prototypes include concepts for:

**Executive Dashboard**

- City-wide KPI cards
- Department performance comparison
- Trend charts
- Geographic heat maps

**Transportation Dashboard**

- Interactive map with bus locations
- Route performance metrics
- Time-series analysis
- Route comparison tables

## Emergency Services Dashboard

- Alert feed
- Response time analytics
- Event type distribution
- Zone-based heat maps

## Energy Dashboard

- Building consumption rankings
- Consumption trends
- Device performance monitoring
- Building detail views

## 4.4 Design System

### Color Palette:

- Primary: #2C3E50 (Navy Blue)
- Secondary: #3498DB (Sky Blue)
- Success: #27AE60 (Green)
- Warning: #F39C12 (Orange)
- Danger: #E74C3C (Red)

### Typography:

- Headings: Inter (Sans-serif)
- Data/Numbers: Roboto Mono (Monospace)

### Components:

- KPI cards
- Interactive charts (Chart.js/D3.js)
- Data tables with sorting
- Maps with zone overlays
- Filters and date pickers

## 4.5 Key Features

- **Role-based views:** Different layouts for executives vs. analysts
- **Interactive visualizations:** Drill-down capabilities
- **Responsive design:** Desktop and tablet support
- **Export functionality:** PDF and CSV downloads
- **Real-time updates:** Configurable refresh rates

## 4.6 Accessibility

- Color contrast compliance
- Keyboard navigation
- Screen reader support

- Clear visual hierarchy
- 

## 5. Implementation Status

### 5.1 Completed

- ☒ Database schema (Bronze, Silver, Gold)
- ☒ SQL ETL stored procedures
- ☒ Python ETL pipeline
- ☒ Talend ETL jobs
- ☒ Airflow orchestration
- ☒ Data generation scripts
- ☒ Sample datasets
- ☒ Technical documentation
- ☒ UI/UX design mockups

### 5.2 Future Enhancements

- Dashboard implementation
  - Real-time data streaming
  - Advanced analytics
  - Cloud deployment
  - Enhanced security
  - Performance optimization at scale
- 

## 6. Getting Started

### For Database Developers

1. Review [Database Design](#)
2. Check [Meta Data Definition](#)
3. Follow [SQL ETL Documentation](#)

### For Python Developers

1. Review [Architecture](#)
2. Follow [Python ETL Documentation](#)
3. Check requirements.txt for dependencies

### For ETL Developers

1. Review [Talend Documentation](#)
2. Check data flow diagrams in [docs folder](#)
3. Review [Airflow Orchestration](#)

### For Business Users

1. Review [Stakeholder Analysis](#)



2. Check UI/UX design mockups on Figma/Canva
  3. Understand available data domains
- 

## 7. Related Documentation

For detailed technical implementation guides, please refer to the documentation files in the [docs](#) directory:

- [System Architecture](#) - Technical architecture overview
  - [Meta Data Definition](#) - Complete data catalog with table/column specifications
  - [SQL ETL Documentation](#) - T-SQL stored procedures implementation
  - [Python ETL Documentation](#) - Python/Pandas pipeline implementation
  - [Talend Documentation](#) - Visual ETL jobs guide
  - [Airflow Orchestration](#) - Workflow automation and scheduling
- 

**Document Version:** 1.0

**Last Updated:** November 1, 2025

**Authors:** Smart City Data Engineering Team

**Status:** Complete