
JOB AGGREGATOR LITE

PROJECT 6

Jack Kammerer

CS.4332.252 Introduction to Database Systems

Texas State University

March 2025

Contents

1	Introduction	1
2	Business Forms & Requirement Analysis	1
2.1	Business Forms and UI	1
2.2	Summary of Requirements	1
3	User View Schema	2
3.1	User Journey	2
3.2	Schema Diagram	2
4	Preliminary Database Formation	3
4.1	Database Schema	3
4.2	Normalization	3
4.3	Sample Data	3
4.4	Code Integration	4
5	Conclusion	5

1 Introduction

This document outlines the design and implementation details for **Job Aggregator Lite**, a web scraper application built for Project 6. The project includes both a functional web scraper and a preliminary database design to store scraped data.

2 Business Forms & Requirement Analysis

2.1 Business Forms and UI

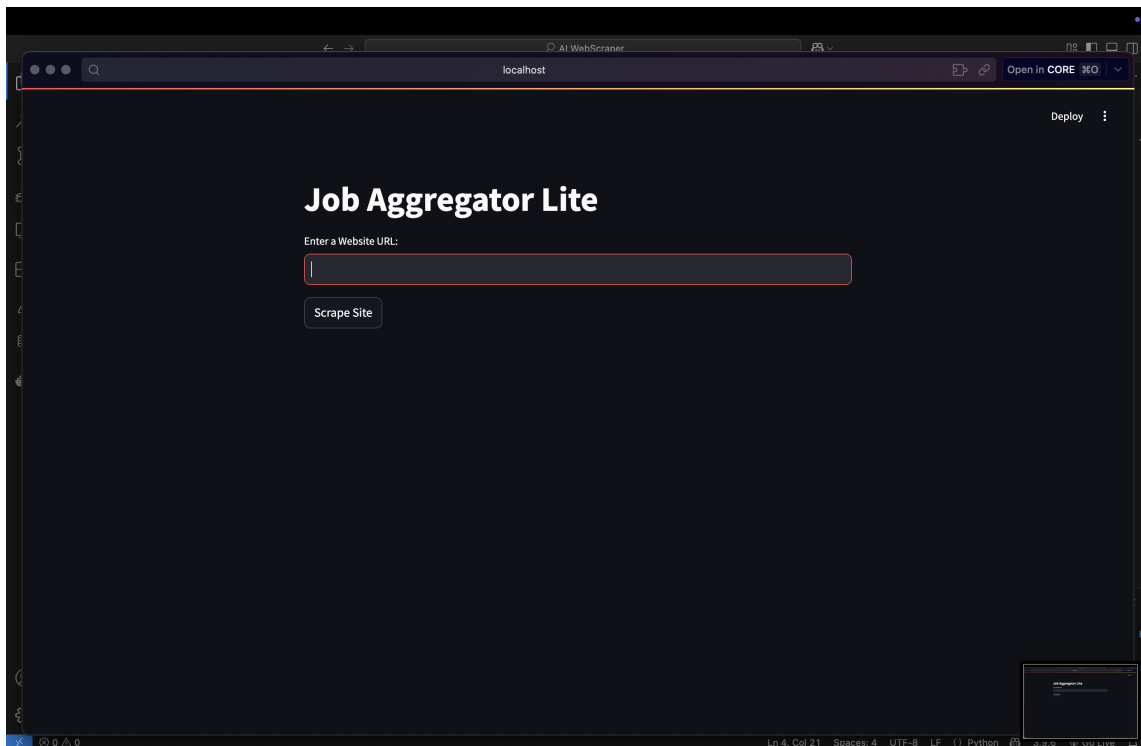


Figure 1: Webscraper UI

- **Input Form:** A simple UI where users enter a URL to be scraped.
- **Results Display:** Once the scraper runs, it shows the scraped content, including details like the URL, timestamp, and content type.
- **History/Report Page:** A page listing previous scrapes with their details.

2.2 Summary of Requirements

- The user enters a URL like Indeed or LinkedIn

- The system scrapes data (e.g., text, links) from the provided URL.
- The scraped data is stored in a relational database.
- A history page lets users view past scrapes.

3 User View Schema

3.1 User Journey

1. **Home Page:** User inputs the URL.
2. **Scraping Process:** The system scrapes the URL.
3. **Results Page:** Displays the scraped data including content type and timestamp.
4. **History Page:** Shows past scrapes for review.

3.2 Schema Diagram

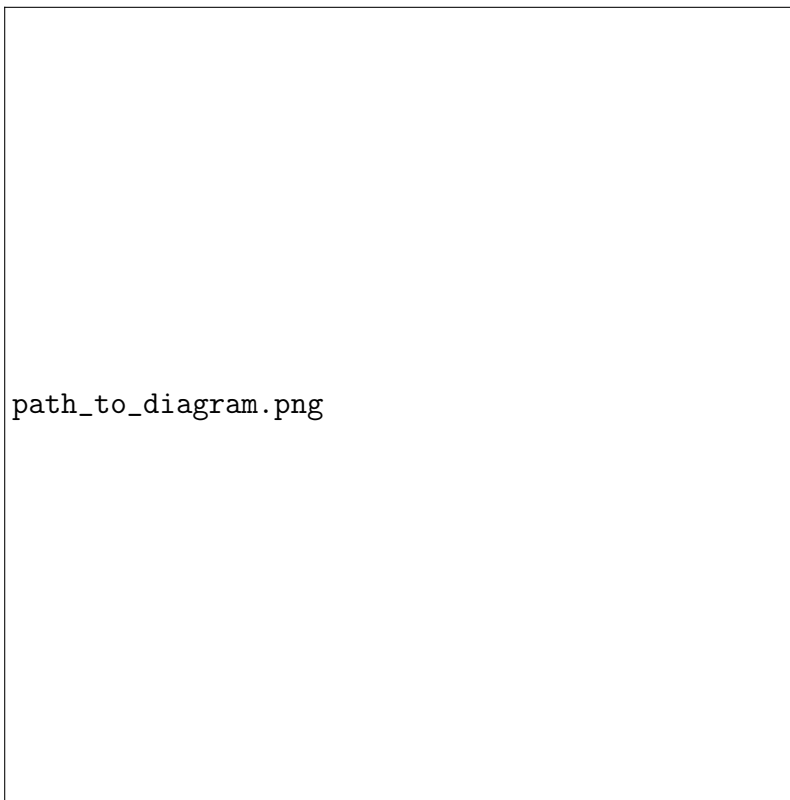


Figure 2: User View Schema Diagram

4 Preliminary Database Formation

4.1 Database Schema

The database consists of two main tables:

Websites Table

- **Columns:** website_id (PK), url, date_added
- **Functional Dependency:** website_id \rightarrow url, date_added

Scraped Data Table

- **Columns:** data_id (PK), website_id (FK), content_type, content, date_scraped
- **Functional Dependency:** data_id \rightarrow website_id, content_type, content, date_scraped

4.2 Normalization

Both tables are designed in Third Normal Form (3NF):

- All non-key attributes are fully functionally dependent on the primary key.
- There are no transitive dependencies.

4.3 Sample Data

Websites Table:

website_id	url	date_added
1	https://www.example.com	2025-03-27 10:00 AM
2	https://www.testsite.org	2025-03-28 02:15 PM

Scraped Data Table:

data_id	website_id	content_type	content
1001	1	text	"Welcome to Example.com, we offer..."
1002	1	link	"https://www.example.com/about"
1003	2	text	"Test Site: This is a sample paragraph..."

4.4 Code Integration

Below is an excerpt of the database integration code from `db.py`:

```

1 import sqlite3
2 from datetime import datetime
3
4 def init_db():
5     conn = sqlite3.connect("scraper.db")
6     c = conn.cursor()
7     c.execute('''CREATE TABLE IF NOT EXISTS websites (
8                 website_id INTEGER PRIMARY KEY AUTOINCREMENT,
9                 url TEXT UNIQUE,
10                date_added DATETIME)''')
11     c.execute('''CREATE TABLE IF NOT EXISTS scraped_data (
12                 data_id INTEGER PRIMARY KEY AUTOINCREMENT,
13                 website_id INTEGER,
14                 content_type TEXT,
15                 content TEXT,
16                 date_scraped DATETIME,
17                 FOREIGN KEY (website_id) REFERENCES websites (
18                     website_id))''')
19     conn.commit()
20     conn.close()
21
22 def insert_scrape(url, scraped_items):
23     conn = sqlite3.connect("scraper.db")
24     c = conn.cursor()
25     now = datetime.now().strftime("%Y-%m-%d %H:%M:%S")
26
27     # Insert into websites table if new URL
28     c.execute("INSERT OR IGNORE INTO websites (url, date_added) VALUES (?, ?)", (url, now))
29     conn.commit()
30
31     # Get website_id
32     c.execute("SELECT website_id FROM websites WHERE url = ?", (url,))
33     website_id = c.fetchone()[0]
34
35     # Insert scraped items (assumes scraped_items is a list of tuples: (
36     content_type, content))
37     for item in scraped_items:
38         c.execute("INSERT INTO scraped_data (website_id, content_type,
39             content, date_scraped) VALUES (?, ?, ?, ?)",

```

```
37         (website_id, item[0], item[1], now))
38     conn.commit()
39     conn.close()
```

Listing 1: Database Code in db.py

5 Conclusion

This document has detailed the design of Job Aggregator Lite, covering the UI, user flow, and preliminary database design. Future work involves fully integrating the UI with the database and performing thorough testing to ensure smooth functionality.