

WorkZone: Job Application Management System

Technical Documentation

Information Retrieval Project

May 22, 2025

Contents

1	Introduction	2
1.1	Project Overview	2
1.2	Motivation	2
2	System Architecture	2
2.1	Component Overview	2
2.2	Data Flow	3
3	Implementation Details	3
3.1	Web Scrapers	3
3.1.1	Design Choices	3
3.1.2	Concurrent Execution	4
3.1.3	Platform-Specific Implementations	4
3.2	Data Management	5
3.2.1	Job Model	5
3.2.2	Data Storage	5
3.2.3	Deduplication	6
3.3	Automation Engine	6
3.3.1	Platform-Specific Appliers	6
3.3.2	Automation Techniques	7
3.4	User Interface	8
3.4.1	Framework Selection	8
3.4.2	Main Application Structure	8
3.4.3	Advanced Search Implementation	9
3.4.4	Resume Management	10
4	Technical Challenges and Solutions	11
4.1	Challenge: Concurrent Database Access	11
4.2	Challenge: Duplicate Job Listings	11
4.3	Challenge: Widget Key Conflicts in Streamlit	11
4.4	Challenge: Date Parsing in Timestamp Filtering	12
5	Conclusion and Future Work	12
5.1	Achievements	12
5.2	Future Enhancements	12
5.3	Lessons Learned	13

1 Introduction

1.1 Project Overview

WorkZone is a comprehensive job application management system designed to streamline the job search and application process. It integrates web scraping, data management, and application automation into a cohesive platform. This document details the system's architecture, implementation choices, and core functionality.

1.2 Motivation

The job application process often involves:

- Searching across multiple job platforms
- Manually tracking applications
- Repeatedly entering the same information
- Managing resume versions

We developed WorkZone to address these pain points by automating and centralizing the entire process, providing a single interface for managing all job search activities.

2 System Architecture

2.1 Component Overview

The system is built around five core components:

System Components

1. **Web Scrapers:** Collection of specialized scrapers for different job platforms
2. **Data Storage:** SQLite database with job and resume information
3. **Automation Engine:** Platform-specific application submission automation
4. **Web Interface:** Streamlit-based user interface
5. **Utilities:** Support functions for logging and data management

2.2 Data Flow

1. Web scrapers collect job listings from multiple sources
2. Jobs are stored in a centralized SQLite database
3. Users search and filter jobs via the Streamlit interface
4. Users upload and manage resumes via the interface
5. The automation engine applies to selected jobs with user credentials

3 Implementation Details

3.1 Web Scrapers

We implemented specialized scrapers for seven different job platforms, each tailored to the specific HTML structure and data patterns of its target site.

3.1.1 Design Choices

We chose to implement individual scraper classes for each platform to maximize flexibility and maintainability:

```
1 # From crawler.py
2 scraper_classes = [
3     LinkedInScraper,
4     FreelancerScraper,
5     WuzzufScraper,
6     RemoteOKScraper,
```

```

7     WeWorkRemotelyScraper ,
8     UpworkScraper ,
9     PeoplePerHourScraper
10 ]

```

Listing 1: Scraper Class Structure

Each scraper implements a common interface with a `scrape()` method, allowing for consistent integration with the main crawler.

3.1.2 Concurrent Execution

To optimize data collection, we implemented concurrent scraping using Python's `ThreadPoolExecutor`:

```

1 # From crawler.py
2 with concurrent.futures.ThreadPoolExecutor(max_workers=len(
3     scraper_classes)) as executor:
4     futures = {executor.submit(run_scraper, scraper_class,
5         query, db_name): scraper_class.__name__
6         for scraper_class in scraper_classes}
7     for future in concurrent.futures.as_completed(futures):
8         scraper_name = futures[future]
9         try:
10             future.result()
11             logging.info(f"{scraper_name} completed
12             successfully")
13         except Exception as e:
14             logging.error(f"{scraper_name} generated an
15             exception: {e}")

```

Listing 2: Concurrent Scraping Implementation

This approach offers several advantages:

- Parallel execution reduces total scraping time
- Each scraper runs independently with its own database connection
- Error handling isolates failures to individual scrapers

3.1.3 Platform-Specific Implementations

Each scraper is tailored to its target platform. For example, the LinkedIn scraper handles authentication, pagination, and job extraction specific to LinkedIn's structure:

```

1 # From linkedin.py
2 class LinkedInScraper:

```

```

3     def __init__(self, storage, query="software engineer"):
4         self.storage = storage
5         self.query = query
6         self.driver = webdriver.Chrome()
7
8     def scrape(self, max_pages=5):
9         # LinkedIn-specific implementation
10        # ...

```

Listing 3: LinkedIn Scraper (Excerpt)

3.2 Data Management

3.2.1 Job Model

We created a standardized Job class to normalize data across platforms:

```

1 # From models.py
2 class Job:
3     def __init__(self, title=None, description=None, link=
4         None, company=None,
5             source=None, timestamp=None, location=None):
6         self.id = str(uuid.uuid4())
7         self.title = title or 'non'
8         self.description = description or 'non'
9         self.link = link or 'non'
10        self.company = company or 'non'
11        self.source = source or 'non'
12        self.timestamp = timestamp or datetime.now().
isoformat()
        self.location = location or 'non'

```

Listing 4: Job Data Model

Key design decisions:

- UUID-based primary key for database integrity
- Default values to ensure data consistency
- ISO format timestamp for standardization
- Common fields across all job platforms

3.2.2 Data Storage

We chose SQLite for data persistence due to its simplicity, portability, and performance:

```

1 # From models.py
2 class DataStorage:
3     def __init__(self, output_format='sqlite', db_name='jobs.
4         db'):
5         self.jobs = []
6         self.output_format = output_format
7         self.db_name = db_name
8         self.conn = sqlite3.connect(self.db_name)
9         self.create_table()

```

Listing 5: DataStorage Implementation

The `DataStorage` class provides an abstraction layer over the database operations:

- Table creation and schema management
- Job insertion with duplicate handling
- Connection management
- Support for multiple output formats (CSV, JSON)

3.2.3 Deduplication

To prevent duplicate job listings, we implemented a deduplication utility:

```

1 # From utils.py
2 def deduplicate_jobs(conn):
3     with conn:
4         conn.execute("""
5             DELETE FROM jobs
6             WHERE rowid NOT IN (
7                 SELECT MIN(rowid)
8                 FROM jobs
9                 GROUP BY link
10            )
11        """)

```

Listing 6: Deduplication Implementation

This SQL query efficiently removes duplicates while keeping the earliest record (by rowid) for each unique job link.

3.3 Automation Engine

3.3.1 Platform-Specific Appliers

We developed automation modules for three major platforms:

- LinkedIn
- Wuzzuf
- Freelancer

Each automation module handles the specific workflow of its target platform:

```

1 # From job_applier.py
2 def apply_to_job(job, resume_path, labels, credentials):
3     source = job['source'].lower()
4     if source == 'linkedin':
5         return apply_to_linkedin_job(job['link'], resume_path
6                                     , labels,
7                                     credentials['
8                                     linkedin_email'],
9                                     credentials['
10                                    linkedin_password'])
11     elif source == 'wuzzuf':
12         return apply_to_wuzzuf_job(job['link'], resume_path,
13                                   labels,
14                                   credentials['wuzzuf_email
15                                   '],
16                                   credentials['
17                                   wuzzuf_password'])
18     # ...

```

Listing 7: Job Application Automation

3.3.2 Automation Techniques

The automation modules use Selenium WebDriver to interact with job platforms:

- Automated form filling
- Resume upload
- Authentication handling
- Application status tracking

```

1 # From linkedin_applier.py
2 def apply_to_linkedin_job(job_url, resume_path, labels, email
3                           , password):
4     driver = webdriver.Chrome()
5     try:

```



```

5         # Login to LinkedIn
6         driver.get("https://www.linkedin.com/login")
7         # ... authentication logic ...
8
9         # Navigate to job
10        driver.get(job_url)
11        # ... application logic ...
12
13        return True, "Application submitted successfully."
14    except Exception as e:
15        return False, f"Failed to apply: {str(e)}"
16    finally:
17        driver.quit()

```

Listing 8: LinkedIn Application Automation (Excerpt)

3.4 User Interface

3.4.1 Framework Selection

We chose Streamlit for the user interface due to its:

- Rapid development capabilities
- Interactive widgets
- Python integration
- Minimal frontend code requirements

3.4.2 Main Application Structure

The main application provides navigation and authentication:

```

1  # From main.py
2  import streamlit as st
3  from streamlit_option_menu import option_menu
4  from auth import login, logout, get_current_user
5
6  st.set_page_config(page_title="Job Application Manager")
7
8  # User authentication
9  user = get_current_user()
10 if not user:
11     login()
12     st.stop()
13
14 # Navigation menu

```

```

15 selected = option_menu(
16     "Main Menu",
17     ["Job Search", "Saved Jobs", "Applications", "Resume
18     Manager",
19     "Analytics", "Settings"],
20     icons=["search", "bookmark", "check2-circle", "file-
21     earmark-person",
22     "bar-chart", "gear"],
23     menu_icon="cast",
24     default_index=0,
25     orientation="horizontal"
26 )
27
28 # Page routing
29 if selected == "Job Search":
30     from my_pages.job_search import job_search_page
31     job_search_page()
32 # ... other pages ...

```

Listing 9: Main Streamlit Application

3.4.3 Advanced Search Implementation

We implemented an advanced search interface with multiple filtering options:

```

1 # From job_search.py
2 # Advanced Filters
3 with st.expander("Advanced Search Options", expanded=True):
4     col1, col2, col3 = st.columns(3)
5     with col1:
6         title_filter = st.text_input("Job Title")
7         # Multi-source selection
8         all_sources = sorted(df["source"].unique().tolist())
9         source_filter = st.multiselect("Source(s)", ["All"] +
10                                     all_sources,
11                                     default=["All"])
12
13     # ... other filters ...
14
15     # Date range filter
16     parsed_timestamps = pd.to_datetime(df["timestamp"],
17                                     errors="coerce",
18                                     format='mixed')
19
20     min_date = parsed_timestamps.min()
21     max_date = parsed_timestamps.max()
22     date_range = st.slider(
23         "Date Range (Timestamp)",
24         min_value=min_date,
25         max_value=max_date,
26         value=(min_date, max_date),

```

```

23         format="YYYY-MM-DD"
24     )

```

Listing 10: Advanced Search Implementation

3.4.4 Resume Management

The resume manager page enables:

- PDF/DOCX resume upload
- Resume preview
- Metadata management
- Information categorization

```

1  # From resume_manager.py
2  def resume_manager_page():
3      st.header("Resume Manager")
4      user = get_current_user()
5
6      # ... database setup ...
7
8      # Upload resume section
9      uploaded_file = st.file_uploader("Upload or Update Resume
10                                     ",
11                                     type=["pdf", "docx"],
12                                     key=f"resume_upload_{user}
13                                     }")
14
15     # ... resume processing ...
16
17     # Display form for labels
18     st.subheader("Resume Information")
19
20     # Group labels by category
21     categories = {
22         "Personal Information": ["name", "age", "national_id"
23         , ...],
24         "Professional Summary": ["overview", "
25         career_objective"],
26         # ... other categories ...
27     }
28
29     # Create tabs for different categories
30     category_tabs = st.tabs(list(categories.keys()))

```

Listing 11: Resume Manager Implementation (Excerpt)

4 Technical Challenges and Solutions

4.1 Challenge: Concurrent Database Access

Problem: Multiple scrapers attempting to write to the same database simultaneously caused connection conflicts.

Solution: We implemented a separate database connection for each scraper thread:

```
1 # From crawler.py
2 def run_scraper(scraper_class, query, db_name):
3     storage = DataStorage(output_format='sqlite', db_name=
4     db_name)
5     scraper = scraper_class(storage, query=query)
6     scraper.scrape(max_pages=5)
7     del storage # Ensure connection is closed
8     return scraper_class.__name__
```

Listing 12: Independent Database Connections

4.2 Challenge: Duplicate Job Listings

Problem: The same job might appear on multiple platforms or in multiple search results.

Solution: We implemented SQL-based deduplication using the job URL as the unique identifier:

```
1 DELETE FROM jobs
2 WHERE rowid NOT IN (
3     SELECT MIN(rowid)
4     FROM jobs
5     GROUP BY link
6 )
```

Listing 13: Deduplication SQL

4.3 Challenge: Widget Key Conflicts in Streamlit

Problem: Multiple identical widgets caused ‘DuplicateWidgetID’ errors in the UI.

Solution: We implemented dynamic key generation for all widgets:

```
1 # Unique file uploader key
2 uploaded_file = st.file_uploader(
3     "Upload or Update Resume",
4     type=["pdf", "docx"],
5     key=f"resume_upload_{user}"
```

```

6 )
7
8 # Unique button keys with UUID
9 unique_id = str(uuid.uuid4())
10 st.button("Apply Now", key=f"apply_{row['id']}_{idx}_{unique_id}")

```

Listing 14: Dynamic Widget Keys

4.4 Challenge: Date Parsing in Timestamp Filtering

Problem: Inconsistent timestamp formats caused parsing errors.

Solution: We implemented robust timestamp parsing:

```

1 parsed_timestamps = pd.to_datetime(df["timestamp"],
2                                   errors="coerce",
3                                   format='mixed')

```

Listing 15: Robust Timestamp Parsing

5 Conclusion and Future Work

5.1 Achievements

The WorkZone system successfully:

- Integrates job data from seven major platforms
- Provides advanced search and filtering capabilities
- Manages resume uploads and metadata
- Automates application submission for three major platforms
- Centralizes job application tracking

5.2 Future Enhancements

Potential future developments include:

- Machine learning for job matching and recommendations
- Additional platform integrations
- Enhanced resume parsing and optimization
- API development for mobile applications
- Interview scheduling and preparation tools

5.3 Lessons Learned

This project reinforced several key principles:

- The importance of modular design for system extensibility
- Effective error handling for robust scraping
- Database optimization techniques for performance
- UI/UX considerations for complex data interaction
- The power of automation in streamlining repetitive tasks