



URL-Archiver

Project report

Course of study

Project 1

Author

Nicolin Dora, Abidin Vejseli & Kilian Wampfler

Advisor

Dr. Simon Kramer

Co-advisor

Frank Helbling

Version 1.0 of January 14, 2024

Abstract

The URL-Archiver is a platform-independent Java application which specializes in extracting and archiving URLs from Unicode text or PDF files. It supports archiving on Archive Today and the Wayback Machine, with options to export archived URLs into CSV or the original Bibtex files. The application was developed with a focus on simplicity, modularity and self-explanatory code.

This report details the development process of the URL-Archiver, explaining the implementation strategies, SCRUM methodology and challenges overcome, and concludes with a outlook on potential future enhancements.

Contents

| | |
|---|------|
| Abstract | ii |
| List of Tables | v |
| List of Figures | vi |
| Listings | viii |
| 1. Introduction | 1 |
| 1.1. Initial Situation | 1 |
| 1.2. Poduct Goal | 2 |
| 1.3. Priorities | 3 |
| 2. Specification | 4 |
| 2.1. System Delimitation | 4 |
| 2.1.1. System Environment (statics) | 4 |
| 2.1.2. Process Environment (dynamics) | 6 |
| 2.2. Requirements | 8 |
| 2.2.1. Epics and User Stories | 8 |
| 2.2.2. Functional requirements | 13 |
| 2.2.3. Boundary | 13 |
| 2.2.4. Pre-conditions | 14 |
| 2.3. Usability | 15 |
| 2.3.1. Personas | 15 |
| 2.3.2. Persona 1 - Alex Frei (Investigative Journalist) | 15 |
| 2.3.3. Persona 2 - Dr. Emma Katrin Winter (Academic Researcher) | 16 |
| 2.3.4. Persona 3 - Michael Schwarz (Content Marketing Specialist) | 17 |
| 2.3.5. Storyboard | 17 |
| 2.3.6. UX-Prototyping | 23 |
| 3. Implementation | 27 |
| 3.1. Architecture | 27 |
| 3.1.1. Frontend | 27 |
| 3.1.2. Backend | 29 |
| 3.2. Processes | 35 |
| 3.2.1. Allocation of roles | 35 |
| 3.2.2. Scrum roles | 35 |

| | |
|--|----|
| 3.2.3. Additional roles | 35 |
| 3.2.4. Sprint Goals | 36 |
| 3.2.5. Requirements | 37 |
| 3.2.6. Scrum Adaptionen | 60 |
| 3.2.7. Review Project Setup | 65 |
| 3.2.8. Review | 67 |
| 3.2.9. Retrospective I: Scrum Method | 69 |
| 3.2.10. Retrospektive II: Tools/Instruments | 72 |
| 3.2.11. Lessons learned | 74 |
| 4. Deployment/Integration | 76 |
| 4.1. Licensing and Compliance | 76 |
| 4.1.1. Licensing Overview | 76 |
| 4.1.2. Compliance with Open Source Licenses | 76 |
| 4.1.3. Purpose of Compliance | 77 |
| 4.2. Installation (Sysadmin) Manual & Script | 77 |
| 4.2.1. Requirements | 77 |
| 4.2.2. Clone the repository | 77 |
| 4.2.3. Build and run scripts | 77 |
| 4.3. User Manual | 79 |
| 4.3.1. Getting Started | 79 |
| 4.3.2. Operating Instructions | 79 |
| 5. Conclusion | 81 |
| 5.1. Discussion | 81 |
| 5.2. Bottom Line | 82 |
| 5.2.1. Conclusion | 82 |
| 5.2.2. Productivity Increase Calculation | 82 |
| 5.3. Future Work | 84 |
| Glossary | 85 |
| Bibliography | 88 |
| A. Appendix | 89 |
| A.1. Original Project Description | 90 |
| A.2. List of Used Libraries and Their Licenses | 91 |
| A.3. Class Diagram | 92 |

List of Tables

| | |
|--|----|
| 3.1. Scrum Roles | 35 |
| 3.2. Additional Scrum Roles | 36 |
| A.1. List of Used Libraries in the Project | 91 |

List of Figures

| | | |
|-------|--|----|
| 2.1. | High-level MVC-Pattern from URL-Archiver | 5 |
| 2.2. | Extension with new archiving service XY (using the Factory Pattern) | 5 |
| 2.3. | High-level operational process | 7 |
| 2.4. | The console greets Bob with a straightforward interface, presenting the archiving options. Bob opts to begin the archiving process. | 18 |
| 2.5. | Bob is prompted to enter the file or directory path. He types in the path to the PDF document that contains the URLs he wishes to archive. | 18 |
| 2.6. | Upon processing the input, the URL-Archiver quickly finds the first URL within the document. The console inquires if Alex wants to open it in a web browser or proceed with archiving. Alex decides not to open the URL in the browser and instead chooses to archive it directly. | 19 |
| 2.7. | The Archiver initiates the archiving process. A loading screen appears, indicating that the Archive Today platform is working on capturing the webpage. | 20 |
| 2.8. | The first URL is archived, and the Archiver outputs the archived link. | 20 |
| 2.9. | The application moves to the next URL. Without hesitation, Bob decides to archive this one as well. | 21 |
| 2.10. | The last URL is processed, and the Archiver asks Bob if he wants to save the list of archived URLs as a CSV file. Bob types "Y" to confirm. | 22 |
| 2.11. | Bob feels accomplished. The Archiver saved time, and the important URLs are now securely archived. | 22 |
| 2.12. | Starting point of the URL-Archiver where the user is prompted to begin archiving or exit the application. | 23 |
| 2.13. | Upon choosing to open a URL in the browser, the user is presented with the live web page. | 24 |
| 2.14. | The user opts to archive the second URL, initiating the archiving process through the interface. | 24 |
| 2.15. | A security step requiring captcha verification to proceed with the archiving process. | 25 |
| 2.16. | Archive Today informs that the website has been archived before. The application proceeds to re-archive the site. | 25 |
| 2.17. | The archiving process is actively underway on the Archive Today platform. | 26 |
| 2.18. | After successful archiving, the system provides the user with a confirmation message and the archived URL. | 26 |
| 3.1. | Screenshot of the Console View | 28 |
| 3.2. | Highlevel Diagram of implemented MVC Pattern | 30 |

| | | |
|-------|--|----|
| 3.3. | Diagram of the FileReader Factory | 32 |
| 3.4. | Diagram of the Exporter Factory | 32 |
| 3.5. | Product Backlog | 37 |
| 3.6. | Sprint 1 Backlog | 38 |
| 3.7. | User Story Detail for "Processing Feedback" | 39 |
| 3.8. | User Story Detail for "Prompt for file path input" | 39 |
| 3.9. | User Story Detail for "Automatic File Type Detection" | 40 |
| 3.10. | User Story Detail for "Processing of Directories" | 41 |
| 3.11. | Sprint 1 Burn Up Chart | 42 |
| 3.12. | Sprint 2 Backlog | 43 |
| 3.13. | User Story Detail for "Intermediate Documentation" | 43 |
| 3.14. | User Story Detail for "Creation of Base Structure" | 44 |
| 3.15. | User Story Detail for "Scan Files for URLs" | 44 |
| 3.16. | User Story Detail for "Sequential URL Preview" | 45 |
| 3.17. | User Story Detail for "Implement temporary store for extracted URLs" | 45 |
| 3.18. | User Story Detail for "Intermediate presentation" | 46 |
| 3.19. | Sprint 2 Burn Down Chart | 47 |
| 3.20. | Sprint 3 Backlog | 48 |
| 3.21. | User Story Detail for "Intermediate Documentation" | 48 |
| 3.22. | User Story Detail for "Creation of Base Structure" | 49 |
| 3.23. | User Story Detail for "Scan Files for URLs" | 49 |
| 3.24. | Sprint 3 Burn Down Chart | 50 |
| 3.25. | Sprint 4 Backlog | 51 |
| 3.26. | User Story Detail for "Generate CSV File" | 51 |
| 3.27. | User Story Detail for "Fix bug with selenium on linux" | 52 |
| 3.28. | User Story Detail for "Fix bug with selenium on MacOS" | 52 |
| 3.29. | User Story Detail for "Fix bug With selenium and Edge browser" | 53 |
| 3.30. | User Story Detail for "Document licences" | 53 |
| 3.31. | User Story Detail for "Periodic complete code review" | 53 |
| 3.32. | User Story Detail for "Document SCRUM sprint 3" | 54 |
| 3.33. | User Story Detail for "Creation of a config file" | 54 |
| 3.34. | User Story Detail for "Progress indicator archiving" | 55 |
| 3.35. | Bug Detail for "No URL found in file is not handled" | 55 |
| 3.36. | Sprint 4 Burn Down Chart | 56 |
| 3.37. | Sprint 5 Backlog | 57 |
| 3.38. | User Story Detail for "Integrate Archived URLs into Supported Files" | 57 |
| 3.39. | User Story Detail for "Create / Fix Tests" | 58 |
| 3.40. | User Story Detail for "Show archived urls path" | 58 |
| 3.41. | Sprint 5 Burn Down Chart | 59 |
| 3.42. | Screenshot from the user story template. | 62 |
| 3.43. | Illustration of T-Shirt Sizes | 63 |
| 3.44. | Velocity Report | 68 |

Listings

1. Introduction

1.1. Initial Situation

The Internet is constantly evolving, which means that there is no guarantee that a website as it exists today will still exist in a few years' time, let alone contain the same information. While this might not be a concern that the average Internet user has to grapple with, it poses a challenge to the academic demographic, where it becomes crucial to reference sources and potentially integrate links to additional data. If links become inactive, verifying the sources becomes challenging, if not impracticable.

Archiving the existing status of a website is achievable, but it currently necessitates a manual and hence time-intensive operation, which not many people take the time to do. The objective of this project is to devise an automated solution to this predicament that is independent of platforms.

The stakeholders for this solution include:

- ▶ Legal professionals and researchers who need to preserve web content as evidence or for case study references.
- ▶ Journalists and media agencies that require archiving web pages for future reporting or fact-checking.
- ▶ Librarians and archivists tasked with the digital preservation of online materials for historical records.
- ▶ Content creators and marketers who wish to maintain records of web content for portfolio or audit purposes.
- ▶ Educators and students who need to collect and cite online resources for academic projects and research.
- ▶ Organizations and businesses that need to archive their web presence for compliance and record-keeping.

1.2. Product Goal

The product goal is a platform independent Java application called “URL-Archiver“. The application must be Free/Libre and Open Source Software (FLOSS) licensed and fulfil the following functionalities:

1. The software should be CLI¹-based and offer a clear command line.
2. The software should allow the user to input a path, which can be a folder or any Unicode Text File.
3. The software examines the contents of a file or folder to extract any web URLs using a standard regular expression or similar method.
4. If desired, URLs can be automatically opened in a web browser.
5. The extracted URLs are archived on Archive Today and/or Wayback Machine as per the user’s preference.
6. The software outputs the resulting archive URLs to the user.
7. The software generates a CSV file containing the original URL and the archived Version of the URL.
8. Optionally, the archived Versions are written back into the provided BibTeX file.

The product goal is achieved if the software covers all the functionality listed above. Furthermore, the code should be minimalistic, modular, and self-explaining. In addition to the code, it is essential that the following documents are provided:

- ▶ User manual
- ▶ Installation instructions (including installation script)
- ▶ Software documentation

¹Command Line Interface

1.3. Priorities

The following priorities are listed in order of importance:

1. **Functionality:** The primary priority is the accurate extraction and archiving of URLs. The software should reliably identify URLs in varied file types and ensure their successful archiving on Archive Today² or Wayback Machine³.
2. **Usability:** Given the diverse potential user base, the program should be platform-independent and possess a user-friendly interface. While the underlying mechanisms may be complex, the user experience should be seamless and intuitive.
3. **Code Quality:** Emphasis should be placed on writing clean, minimal, and modular code. This not only aids in potential future enhancements but also in debugging and troubleshooting.
4. **Documentation:** As with any software project, proper documentation is paramount. The project report should be concise, adhering to the principle of being “maximally informative, minimally long,” ensuring clarity of information without overwhelming the reader.
5. **Integration with Existing File Types:** The ability to seamlessly insert archived URLs into BibTeX files is a priority, given the potential academic applications of the software.

²<https://archive.today>

³<https://web.archive.org/save/>

2. Specification

2.1. System Delimitation

2.1.1. System Environment (statics)

System Overview

The primary purpose of the URL-Archiver is to extract URLs from Unicode text files and PDFs, and archive them on supported platforms: Archive Today and the Wayback Machine. The system provides the archived URL versions to the user via a CSV file. Additionally, when a BibTeX file is provided by the user, the original BibTeX file is updated with a note field containing these archived URLs for each entry.

Hardware Specifications

The URL-Archiver does not impose any special hardware requirements. However, an internet connection is essential for the archiving process to function.

Software Components

The URL-Archiver is platform-independent , operating on major systems such as Windows (tested on Windows 10, version 22H2 and Windows 11, version 23H2), macOS (tested on macOS Sonoma), and Linux (tested on Ubuntu 20.04.3 LTS). The system has varying browser dependencies based on the operating system: Chrome is required for macOS, Edge for Windows, and Firefox for Ubuntu/Linux (Latest stable versions of the browsers are recommended). Users can change the default Browser in the configuration of the application. Other dependencies are installed with the URL-Archiver and do not require separate installation.

System Architecture

The URL-Archiver uses the Model-View-Controller (MVC) pattern, as illustrated in figure 2.1, to enable future enhancements, such as adding a GUI interface. The Factory Pattern is applied where appropriate to simplify the extension of functionalities. For instance, adding extra archiving services can be easily accomplished by introducing a new archiving service, as shown in figure 2.2.

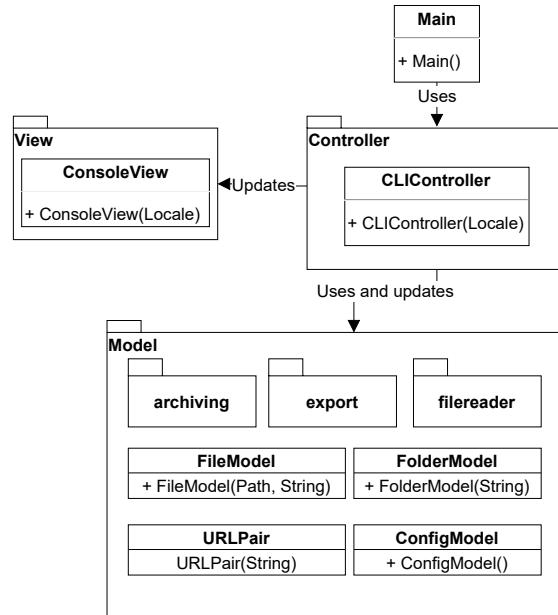


Figure 2.1.: High-level MVC-Pattern from URL-Archiver

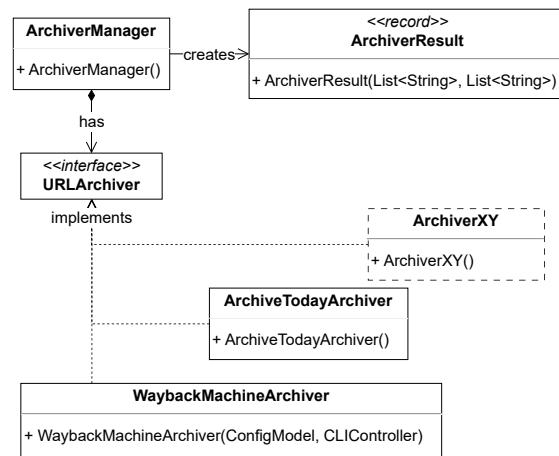


Figure 2.2.: Extension with new archiving service XY (using the Factory Pattern)

Data Management

Upon completion of its execution, the URL-Archiver generates a CSV file where each line contains an extracted URL and its archived versions, separated by semicolons. For example, a line like `https://xy.com;https://web.archive.org/xy;https://archive.ph/xy` shows the original URL and its archives. This simple format makes it easy to track and manage archived URLs.

Optionally, if the user provided a BibTeX file, URLs are integrated into the note field of each entry. If there's no existing note field, a new one is created with the format `note = Archived Versions: url1, url2`. If a note field already exists, the archived URLs are appended to it in the format `note = <current note>, Archived Versions: url1, url2`. This approach ensures that the archived URLs are neatly added to the BibTeX entries, maintaining the integrity of the original data.

User Interface

Currently, the system uses a command-line interface. The MVC-Pattern lays the groundwork for potential future implementation of a GUI interface.

Integration with Other Systems

The system integrates with the Wayback Machine via API, with certain limitations detailed in their API documentation¹. For archiving on Archive.today, which lacks an API, Selenium is used to automate the process as much as possible. However, users must manually complete captchas.

Maintenance and Support

Currently, there are no specified maintenance requirements or a support framework for the URL-Archiver.

2.1.2. Process Environment (dynamics)

Operational Processes

The URL-Archiver is initiated by the user, who provides a path to Unicode text or PDF files or a directory that contains such files. The application extracts URLs from these files and presents them sequentially to the user. The user then has options to open or archive that URL. He has also the following other options:

- ▶ **s:** Show a list of previously archived URLs.

¹<https://archive.org/details/spn-2-public-api-page-docs/mode/2up>

- ▶ **u:** Update and view pending archive jobs.
- ▶ **n:** Navigate to the next URL.
- ▶ **q:** Quit the application.
- ▶ **c:** Change application settings.
- ▶ **h:** Access the Help Menu for assistance.

Upon completion, the user is prompted to save URL pairs to a CSV file and, if a Bibtex file is provided, to write the archived URLs back into it.

Event Handling

In the URL Archiver, user actions are efficiently facilitated through the main menu. When archiving a URL, users can select either the Wayback Machine or Archive.today. In addition, the 'c' option in the menu allows users to configure settings, including setting up API keys for the Wayback Machine and selecting a default browser. The application cleverly handles unsupported input and incorrect paths by prompting the user for the correct information or action. This ensures smooth operation and user guidance throughout the process.

Life Cycle

The URL-Archiver's life cycle begins with launch and path input, proceeding to URL extraction and user interactions via the menu options, and ends with prompts for data saving upon completion. See the high-level process in figure 2.3.

Error Management

Errors within the URL-Archiver are caught and handled, typically prompting the user with a customized error message asking to retry the action. No system stack traces are shown

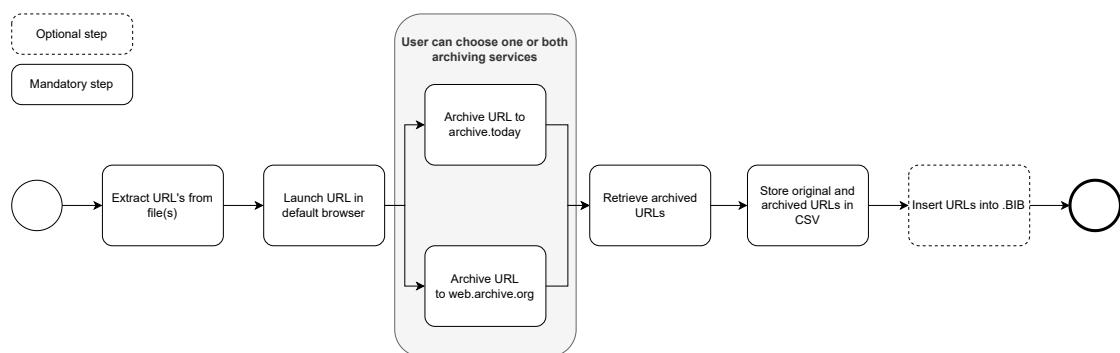


Figure 2.3.: High-level operational process

to the user.

Backup and Recovery

Currently, there are no backup features; progress is not saved if the application ends unexpectedly, which is slated for future improvement.

Update and Upgrade Policies

Software updates require manual download and recompilation from the Git repository. The system does not provide automatic updates or an in-built feature for update checks.

2.2. Requirements

2.2.1. Epics and User Stories

In this section, we outline the main features (Epics) of the project and break them down into detailed user tasks (User Stories). This helps provide a clear understanding of the desired functions and behaviors of our software.

Epic 1: File Input and Processing

Goal: Allow the user to input various file types via the command line and prepare these files for further processing.

1. Prompt for File Path Input

► **Description:** As a user, when I start the tool, I want to be prompted to input the path to my file, so the tool knows which file to process.

► **Acceptance Criteria:**

- Upon starting the tool, it prompts the user to enter a file path.
- On inputting an invalid path or if there are permissions issues, the tool provides a relevant error message.

2. Automatic File Type Detection

► **Description:** As a user, I want the tool to automatically detect the file type (based on file extension) and treat it accordingly so that I don't need to specify the file type separately.

► **Acceptance Criteria:**

- The tool automatically identifies if the file is a .BIB, .TEX, .HTML, or .PDF.

- For unrecognized file types, the tool provides an appropriate error message.

3. Processing of Directories

► **Description:** As a user, I want to input a whole directory, so the tool processes all supported files contained within.

► **Acceptance Criteria:**

- The tool can accept directory paths after the prompt.
- It processes all supported file types within the directory.
- The tool gives a message if files within the directory are skipped due to their type.

4. Processing Feedback

► **Description:** As a user, I want to receive feedback when the tool starts processing the file and when it finishes, to know the status.

► **Acceptance Criteria:**

- A message is displayed when the processing of a file starts.
- Upon completion, a confirmation message is shown, which also includes any potential errors or warnings.

Epic 2: URL Detection and Extraction

Goal: Accurately detect and extract URLs from input files for further processing.

1. Scan Files for URLs

► **Description:** As a user, I want the system to scan my input files and identify any embedded URLs so that they can be extracted for archiving.

► **Acceptance Criteria:**

- System can detect URLs in a variety of file formats including .BIB, .TEX, .HTML, and .PDF.
- Detected URLs are listed without any duplication.

2. Use Regular Expressions for Extraction

► **Description:** As a user, I want the system to use regular expressions or other reliable techniques to extract URLs so that all valid URLs are captured without error.

► **Acceptance Criteria:**

- System uses a robust regular expression pattern that matches most URL formats.
- Extracted URLs are validated to ensure they are in the correct format.

3. Store URL Line Number or Context

► **Description:** As a user, when a URL is detected and extracted, I want the system to also store its line number or contextual information from the original file, enabling precise placement of its archived counterpart later on.

► **Acceptance Criteria:**

- Upon URL detection, the system captures and stores the line number or relevant context of the URL from the source file.
- This information is utilized later if archived URLs need to be placed back into the original files.

4. Compile a List of URLs

► **Description:** After extraction, I want all URLs to be compiled into a single list, eliminating any duplicates, so that I have a clean list for archiving.

► **Acceptance Criteria:**

- The list contains all the unique URLs found in the input files.
- Invalid or broken URLs are flagged or removed from the list.

Epic 3: Web Browser Integration

Goal: Seamlessly open detected URLs, one at a time, in a web browser for user verification, and immediately initiate the archiving process upon user decision.

1. Sequential URL Preview

► **Description:** As a user, I want to preview each detected URL in my default browser sequentially to verify its content.

► **Acceptance Criteria:**

- System opens one URL at a time in the default browser.
- Immediately after the URL is displayed, the system presents the user with the option to archive.

2. Immediate Archiving Upon Decision

► **Description:** After reviewing a URL in the browser, I want to decide if it should be archived. If I decide to archive, the system should immediately initiate the archiving process.

► **Acceptance Criteria:**

- System provides a prompt to accept or decline the archiving of the displayed URL.
- If the user chooses to archive, the system directly begins the archiving process, and the user may need to manually solve captchas.

3. Track Archiving Progress

► **Description:** As a user, I want a clear indicator of how many URLs have been displayed, archived, and how many are left to process.

► **Acceptance Criteria:**

- The system displays a counter indicating the number of URLs already shown to the user.
- Another counter indicates how many URLs have been chosen for archiving.
- Yet another counter shows how many URLs remain to be processed/displayed.

4. Store User Decisions for Reporting

► **Description:** As a user, after making a decision about archiving each URL, I want the system to store my choices so that they can be referred to or reported on later.

► **Acceptance Criteria:**

- The system maintains a record of each URL and the user's decision (archived or not archived).
- The stored decisions are available for any subsequent reporting needs.

Epic 4: Interaction with archive.ph

Goal: Automate the process of archiving URLs via archive.ph while ensuring user interaction is seamless and all necessary data is captured for later use.

1. Automated URL Submission

► **Description:** As a user, I want the system to automatically fill in the URL into the archive.ph input field and submit it for archiving.

► **Acceptance Criteria:**

- Upon initiation, system opens the archive.ph website in a browser.
- System auto-fills the given URL into the appropriate input field.
- System automatically triggers the submission process for archiving.

2. User Interaction for Captchas

► **Description:** If required, I want to manually solve captchas to ensure the URL gets archived.

► **Acceptance Criteria:**

- If archive.ph presents a captcha, the system allows the user to solve it manually.
- The archiving process proceeds once the captcha is successfully solved.

3. Automatic Retrieval of Archived URL

► **Description:** Once a URL is archived, I want the system to automatically retrieve and display the archived URL to me.

► **Acceptance Criteria:**

- System captures the new archived URL from archive.ph after the process completes.
- The archived URL is displayed to the user immediately.
- The archived URL is stored for later processing and reporting.

Epic 5: Output and Reporting

Goal: Provide the user with an organized CSV file detailing URLs and their archived counterparts. Also, allow for integration of archived URLs back into supported input files.

1. Generate CSV File

► **Description:** As a user, I want the system to produce a CSV file containing all original URLs and their corresponding archived URLs.

► **Acceptance Criteria:**

- A CSV file is generated upon completion of the archiving process.
- Each row in the CSV contains the original URL and its archived counterpart.

2. Integrate Archived URLs into Supported Files

► **Description:** If desired, I want the system to insert the archived URL back into the original file, following its corresponding original URL.

► **Acceptance Criteria:**

- The system recognizes supported file types for this integration process.
- Upon user approval, the archived URL is inserted in the appropriate location (e.g., following its original URL) within the file.

2.2.2. Functional requirements

The URL-Archiver project, as described in the original project description, which can be found in the Appendix (A.1), is designed to be a platform-independent Java program that performs several key functions:

- ▶ **Input Processing:** It accepts directories or files (including any Unicode-text- and PDF-Files) as input and scans them for URLs.
- ▶ **URL Extraction:** The program extracts all URLs using regular expressions.
- ▶ **Optional Browser Interaction:** It has an optional feature to open URLs in a web browser.
- ▶ **URL Archiving:** The program posts URLs to an archiving service (Archive Today & The Wayback Machine) and retrieves the archived URLs.
- ▶ **Output Generation:** It outputs a CSV file containing the original and archived URL pairs.
- ▶ **Optional File Integration:** The program can optionally insert archived URLs into a .BIB file.

Additionally, the project emphasizes minimal, modular, and self-explanatory code. The report for the project is expected to be concise and include the original project description.

2.2.3. Boundary

The following limits apply:

- ▶ **Input:** The application processes only Unicode-text-files (e.g.: .BIB, .TEX; .HTML; etc.) and PDF-files
- ▶ **Output:** The application exports the extracted and archived URLs in either CSV format or directly in a BIB file.
- ▶ **Archiving:** At present, the application should only support Archive Today and The Wayback Machine. Please refer to the documentation of the Wayback Machine² or Archive Today³ for any restrictions on their APIs or services.

²API Documentation Wayback Machine

³FAQ of Archive Today

2.2.4. Pre-conditions

For the application to run smoothly, the following pre-conditions must be met:

- ▶ **Java Environment:** The application necessitates the installation of Java Runtime Environment (JRE) version 21 on the system where it is intended to be executed.
- ▶ **File Formats:** The software is designed to operate with specific file formats, primarily Unicode text files and PDF files. It is necessary to provide these file types for successful operation.
- ▶ **Archiving Service Access:** The application requires internet connectivity and access to web archiving services, such as Archive Today, to enable its URL archiving functionality.
- ▶ **Login + S3 Credentials (Wayback Machine only):** To archive with the Wayback Machine, a login and S3 credentials are required.

2.3. Usability

2.3.1. Personas

This section introduces personas for the URL-Archiver, detailing diverse user profiles and their interactions with the application, to provide a understanding of its usability and functionality.

2.3.2. Persona 1 - Alex Frei (Investigative Journalist)

Persona Overview

Name: Alex Frei

Age: 32

Occupation: Investigative Journalist

Industry: Media and News Reporting

Background

Alex is an experienced journalist specializing in political reporting. Known for in-depth investigative pieces, Alex often covers high-stakes political events, requiring rapid access to historical web data for fact-checking.

Goals

- ▶ Quickly archive web pages for timely reporting on political events.
- ▶ Maintain an archive of web pages for in-depth analysis and future reference.
- ▶ Retrieve archived URLs efficiently for referencing in articles and reports.

Challenges

- ▶ Finding a tool that archives web pages instantly and reliably.
- ▶ Ensuring easy access and shareability of archived content.
- ▶ Managing a collection of archived URLs efficiently and effectively.

Interactions with URL-Archiver

- ▶ Utilizes URL-Archiver's CLI for quick and efficient archiving.
- ▶ Values the open-source nature for trustworthiness in journalistic integrity.

- ▶ Relies on CSV output for organized referencing of archived content.

2.3.3. Persona 2 - Dr. Emma Katrin Winter (Academic Researcher)

Persona Overview

Name: Dr. Emma Katrin Winter
Age: 40
Occupation: Academic Researcher
Industry: Higher Education

Background

Dr. Winter is a tenured professor focused on social sciences, often citing digital journals and data repositories. She requires reliable digital archiving for scholarly citations and utilizes BibTeX files for managing bibliographic data.

Goals

- ▶ Ensure accurate citation of web content in academic papers.
- ▶ Maintain and manage a digital archive of resources for teaching and research.
- ▶ Share archived content with academic peers for collaboration.
- ▶ Integrate archived URLs into BibTeX files for streamlined academic referencing.

Challenges

- ▶ Needs a reliable method for archiving web pages.
- ▶ Seeks efficient integration of archived URLs into BibTeX files.
- ▶ Manages extensive digital archives for large-scale research projects.

Interactions with URL-Archiver

- ▶ Uses URL-Archiver for capturing and archiving academic papers.
- ▶ Values the integration of archived URLs into BibTeX files for academic referencing.
- ▶ Relies on CSV file generation for organizing digital resources.

2.3.4. Persona 3 - Michael Schwarz (Content Marketing Specialist)

Persona Overview

Name: Michael Schwarz

Age: 35

Occupation: Content Marketing Specialist

Industry: Digital Marketing

Background

Michael is a skilled content marketer, responsible for digital campaigns and web presence analytics in a dynamic marketing agency. His role involves archiving web content for brand reputation management and competitive analysis.

Goals

- ▶ Regularly archive web content for marketing audits and compliance.
- ▶ Manage a digital archive for campaign retrospectives.
- ▶ Ensure easy access to archived content for reviews and strategic planning.

Challenges

- ▶ Needs a tool capable of archiving web content efficiently.
- ▶ Requires a streamlined process for accessing archived material.
- ▶ Seeks an archive system that is user-friendly and reliable.

Interactions with URL-Archiver

- ▶ Regularly uses URL-Archiver for archiving online marketing content.
- ▶ Appreciates the CLI interface for ease and speed of use.
- ▶ Relies on the tool's capability to create organized audit trails.

2.3.5. Storyboard

In this section, we present the storyboard for the URL-Archiver, illustrating the user's interaction with the application. Through the sketches 2.4 to 2.11, we show the user's journey from the initial greeting in the console to the final archiving of URLs. This story concentrate on the user's experience, the steps they take, and the decisions they make, rather than the precise UI design details.

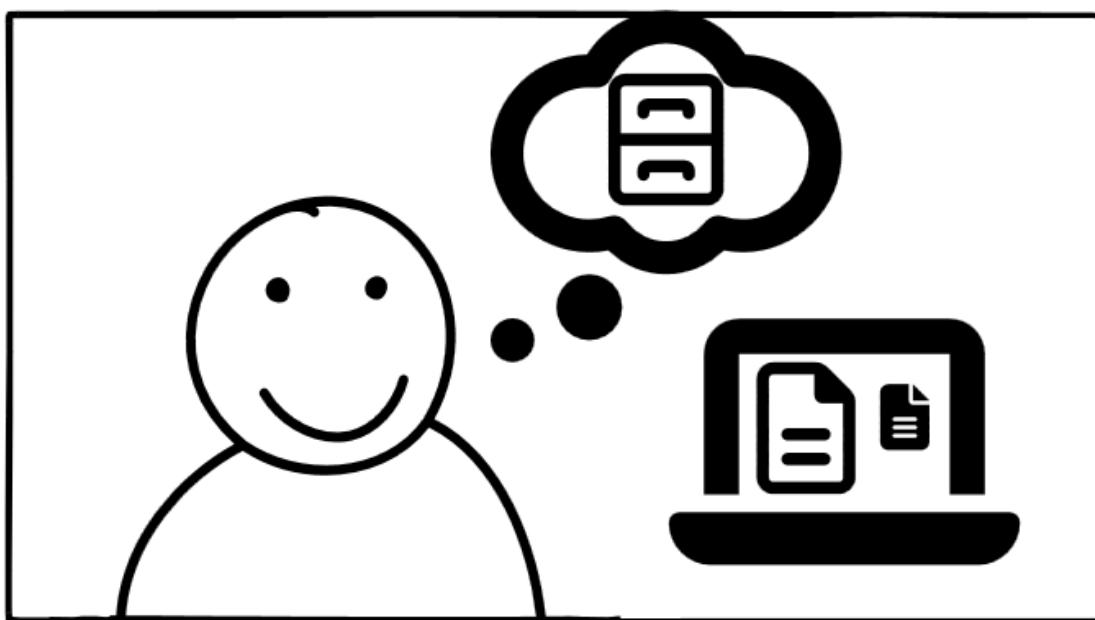


Figure 2.4.: The console greets Bob with a straightforward interface, presenting the archiving options. Bob opts to begin the archiving process.

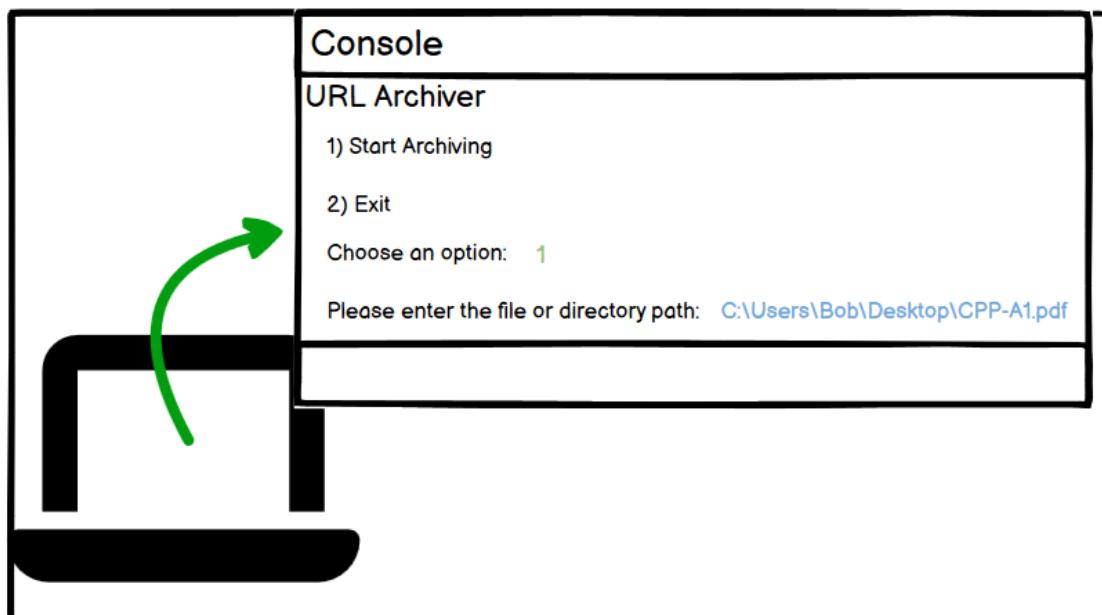


Figure 2.5.: Bob is prompted to enter the file or directory path. He types in the path to the PDF document that contains the URLs he wishes to archive.

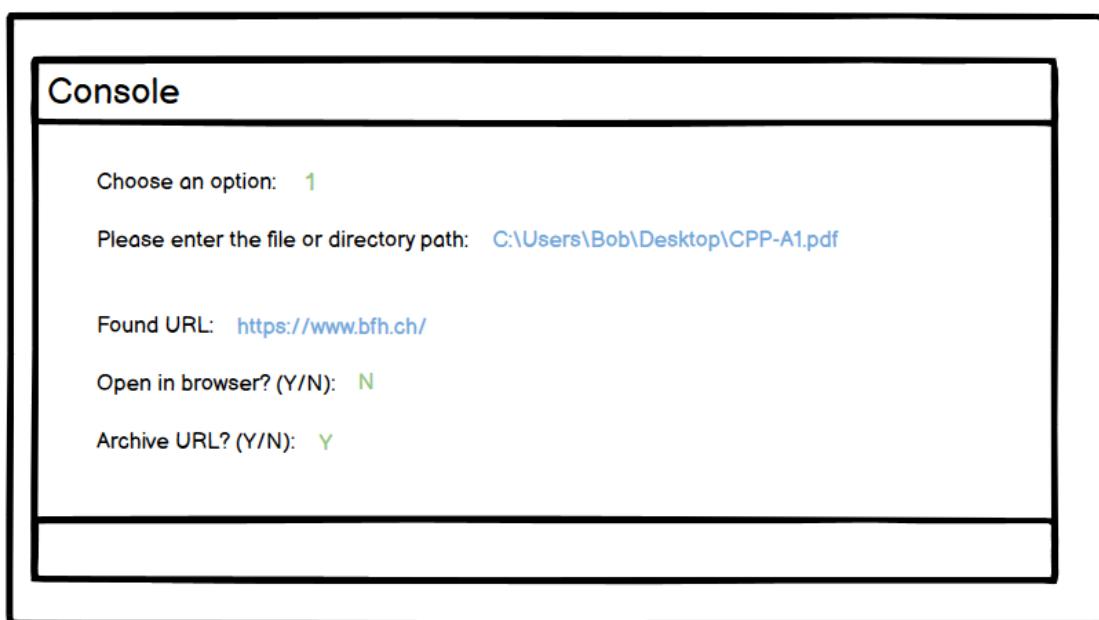


Figure 2.6.: Upon processing the input, the URL-Archiver quickly finds the first URL within the document. The console inquires if Alex wants to open it in a web browser or proceed with archiving. Alex decides not to open the URL in the browser and instead chooses to archive it directly.

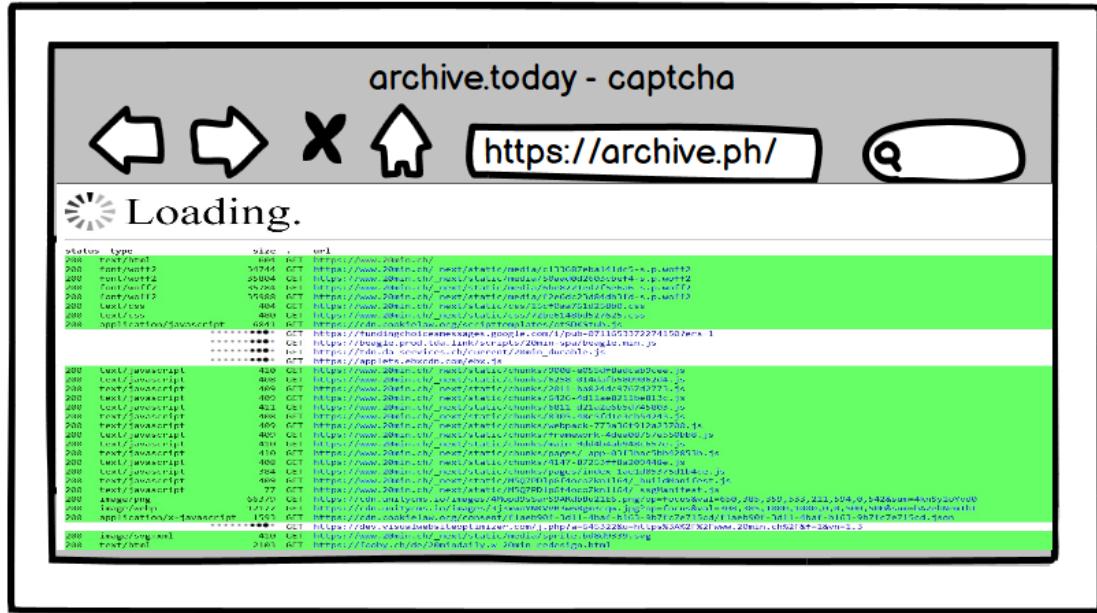


Figure 2.7.: The Archiver initiates the archiving process. A loading screen appears, indicating that the Archive Today platform is working on capturing the webpage.



Figure 2.8.: The first URL is archived, and the Archiver outputs the archived link.



Figure 2.9.: The application moves to the next URL. Without hesitation, Bob decides to archive this one as well.



Figure 2.10.: The last URL is processed, and the Archiver asks Bob if he wants to save the list of archived URLs as a CSV file. Bob types "Y" to confirm.

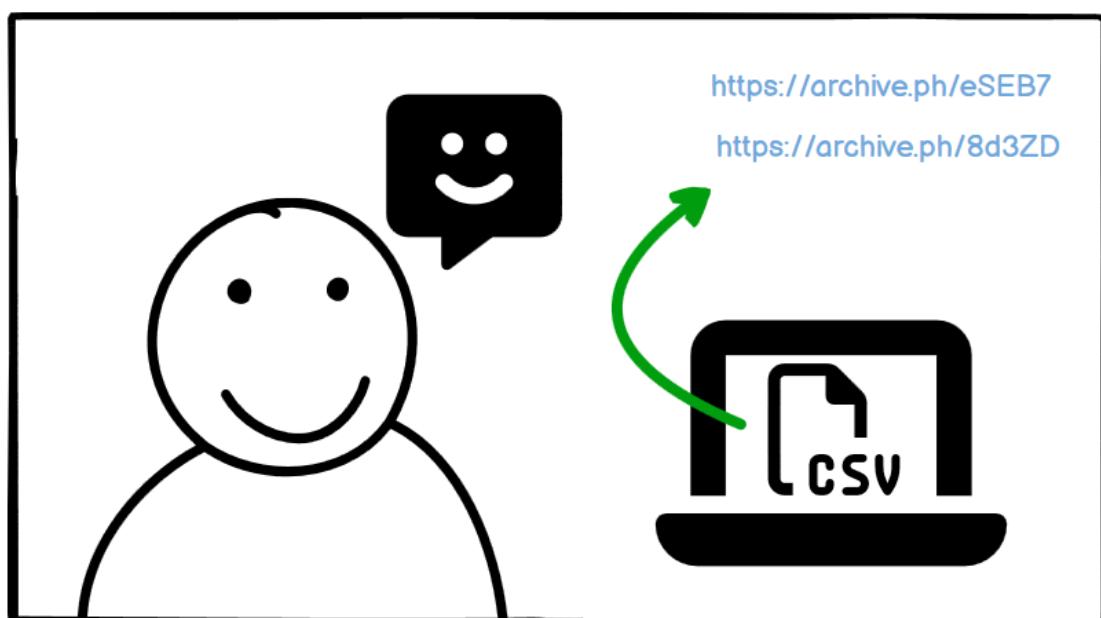


Figure 2.11.: Bob feels accomplished. The Archiver saved time, and the important URLs are now securely archived.

2.3.6. UX-Prototyping

This section presents the initial version of our UX prototype for the URL-Archiver, created using Balsamiq⁴. The series of UI sketches, from Figure 2.12 through Figure 2.18, represent an early conceptualization of the application's user experience. This prototype does not reflect the current state of the application, but serves as an important step in our design process. Each figure illustrates the key stages of user interaction, from starting the application to completing the archiving process.



Figure 2.12.: Starting point of the URL-Archiver where the user is prompted to begin archiving or exit the application.

⁴Balsamiq is a rapid wireframing tool that helps you work faster and smarter. It reproduces the experience of sketching on a whiteboard, but using a computer.



Figure 2.13.: Upon choosing to open a URL in the browser, the user is presented with the live web page.



Figure 2.14.: The user opts to archive the second URL, initiating the archiving process through the interface.

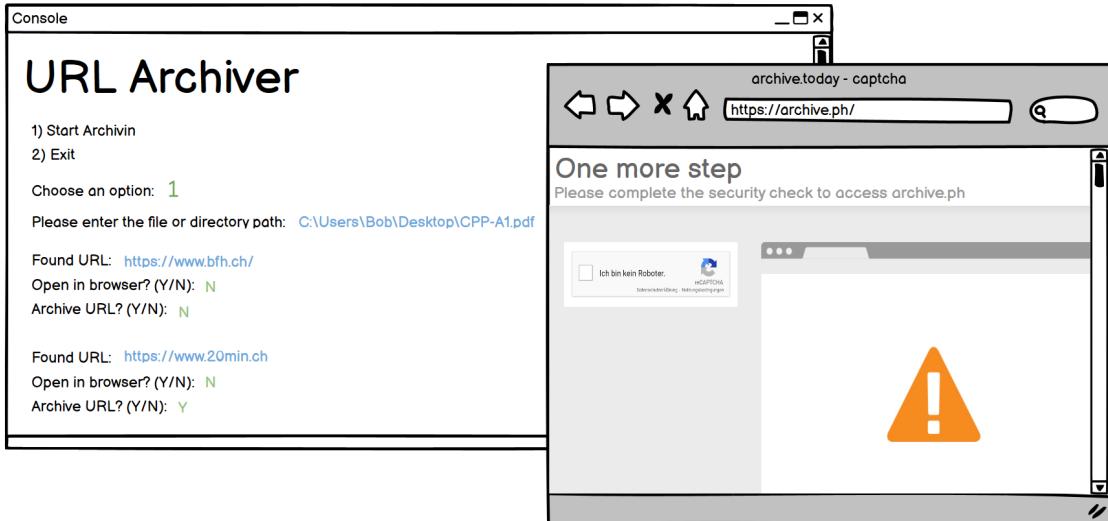


Figure 2.15.: A security step requiring captcha verification to proceed with the archiving process.

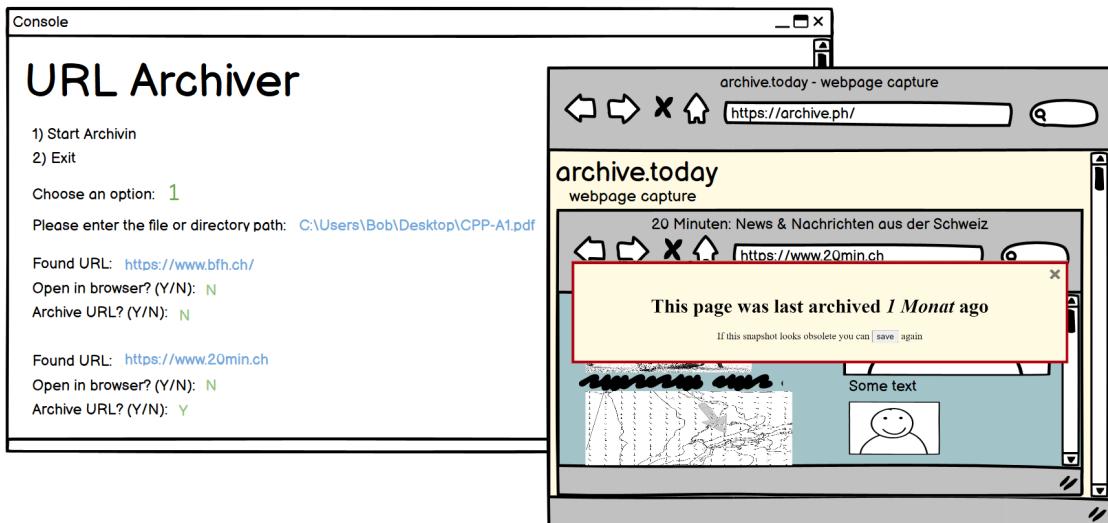


Figure 2.16.: Archive Today informs that the website has been archived before. The application proceeds to re-archive the site.

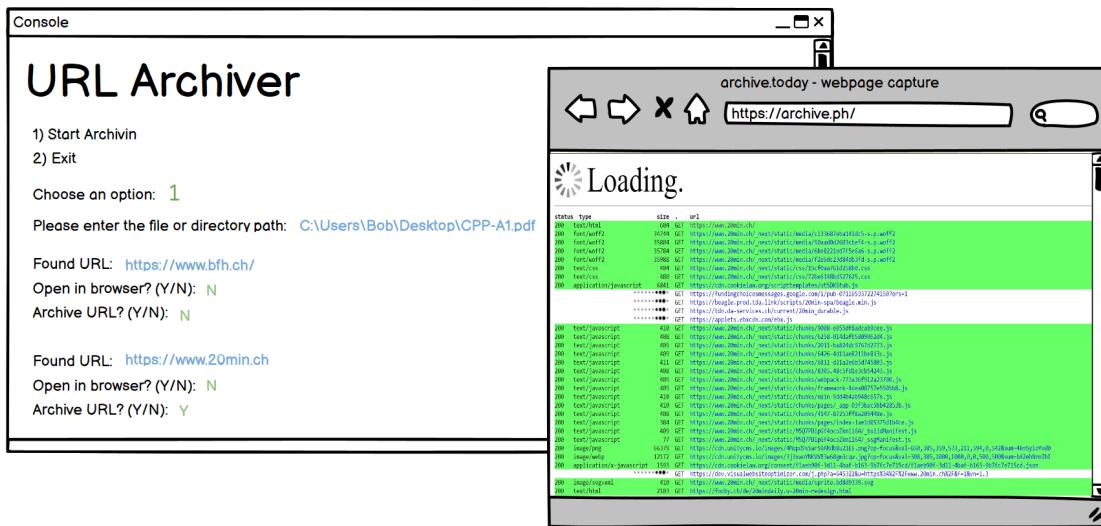


Figure 2.17.: The archiving process is actively underway on the Archive Today platform.



Figure 2.18.: After successful archiving, the system provides the user with a confirmation message and the archived URL.

3. Implementation

3.1. Architecture

This chapter describes the architecture of the URL-Archiver.

3.1.1. Frontend

The frontend architecture of the URL-Archiver is centered on a console-based interface, as evidenced by the `ConsoleView` and `CLIController` classes featured in Figure 3.2.

The `ConsoleView` class is essential for presenting information and managing user input in a console setting. It is specifically designed to display data and messages clearly and in a user-friendly way. In parallel, the `CLIController`, a key component of the Controller segment in the MVC (Model-View-Controller) framework, serves as an intermediary between the `ConsoleView` and the application's backend. This class efficiently processes user inputs from the `ConsoleView`, liaises with the model to retrieve or alter data, and then updates the console interface with these changes. This architecture is crafted to enable streamlined and effective user interactions within a command-line environment. The Figure 3.1 illustrates the welcome message that the URL archiver displays to users.

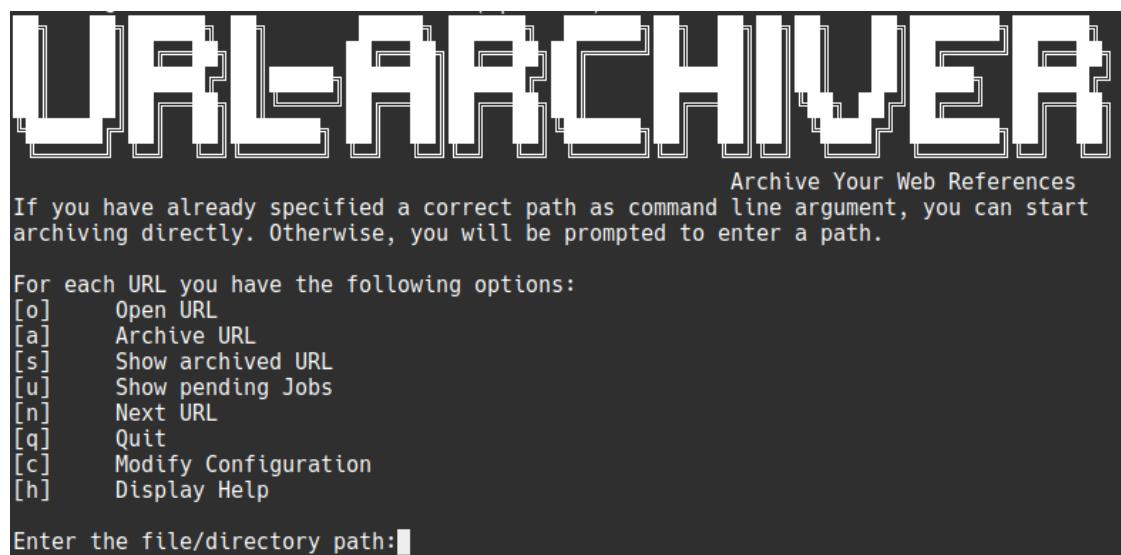


Figure 3.1.: Screenshot of the Console View

3.1.2. Backend

Software Architectural Design

The URL-Archiver is developed following the principles of the MVC (Model-View-Controller) software architectural pattern. The MVC framework, as detailed in Figure 3.2, is a structured approach to building user interfaces in software applications. It organizes an application into three interrelated components:

- ▶ **Model:** This layer embodies the URL-Archiver's data structures, business logic, and operational rules. Components such as `FileModel`, `FolderModel`, `URLPair`, and `ConfigModel`, as shown in Figure 3.2, represent the model.
- ▶ **View:** Tasked with rendering data from the model to the user and forwarding user commands to the controller. The `ConsoleView` in Figure 3.2 serves as a demonstration of this component.
- ▶ **Controller:** It is responsible for handling user inputs, engaging with the model to process data, and selecting the appropriate view for presentation. The `CLIController`, depicted in Figure 3.2, is indicative of this segment.

This modular structure allows for an efficient separation of code, which aids in easier maintenance and streamlined development. The `Main` class operates as the launching point for the URL-Archiver, effectively coordinating the MVC pattern. Interaction between the Controller and the Model results in updates to the Model, with the View being refreshed in tandem to maintain a clear demarcation of responsibilities.

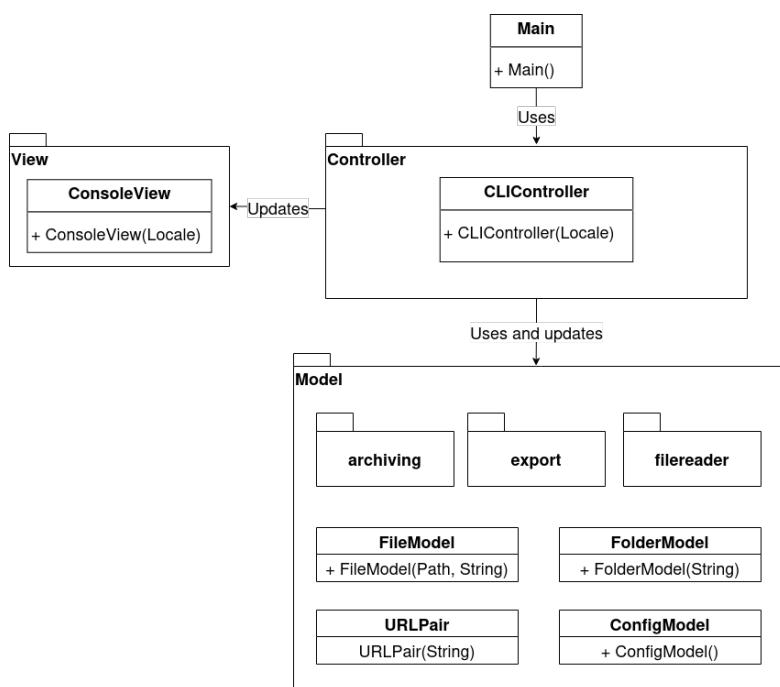


Figure 3.2.: Highlevel Diagram of implemented MVC Pattern

Factory Pattern

The URL-Archiver makes extensive use of the Factory Pattern in its architecture, a key strategy in software engineering recognised for its role in the creational design paradigm. This pattern is central to the encapsulation of the object instantiation process. Rather than creating objects directly using the “new” operator, the Factory Pattern specifies a factory for this purpose. This abstraction improves the flexibility and scalability of the code.

The factory pattern is particularly useful when dealing with a variety of object types to be created, or when specific logic is required in their creation. Its implementation in the URL Archive provides a clear separation of concerns and improves code reusability, in line with the core principles of object-oriented design and programming.

In the URL-Archiver, the Factory Pattern is implemented in several key areas:

- ▶ **File Readers:** Streamlines the addition of more input file types, thanks to its capability to produce objects conforming to the `FileReaderInterface`. This factory is pivotal in handling diverse file types.
- ▶ **Archiving Services:** Adding an additional service is straightforward, enhancing the application’s ability to integrate different web archiving services.
- ▶ **Selenium Web Drivers:** Supports the integration of additional browsers with ease, showcasing the pattern’s adaptability in browser interactions.
- ▶ **Export Functionality:** The pattern’s application here simplifies the inclusion of support for various file types in the export functionality.

These factories contribute significantly to the modularity and extensibility of the URL Archive by encapsulating the creation logic of various components, thus promoting loose coupling and scalability. The alignment of these approaches with the principles of the Factory Design Pattern underscores the maintainability of the application and its adaptability to new file types and data export formats.

The application of the factory pattern in the URL Archive is illustrated in figures 3.3 and 3.4, which show the `FileReaderFactory` and `ExporterFactory` respectively. Similar design principles are applied to other factories within the application, ensuring a consistent and modular approach across all components.

Configuration File Management

The URL-Archiver uses a JSON configuration file to effectively manage various settings and parameters. Key classes such as `ConfigModel` and `ConfigFileHelper` are integral to this process. The `ConfigModel` is used to define the structure of the configuration data, ensuring that all essential settings are systematically represented and easily accessible within the application. Conversely, the `ConfigFileHelper` is responsible for reading and writing the JSON configuration file, decoupling these tasks from the rest of the application code. This methodology not only streamlines the management of configuration settings, but also

3. Implementation

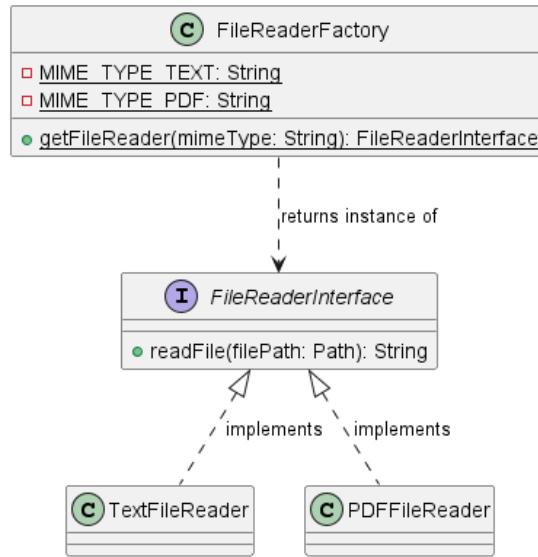


Figure 3.3.: Diagram of the FileReader Factory

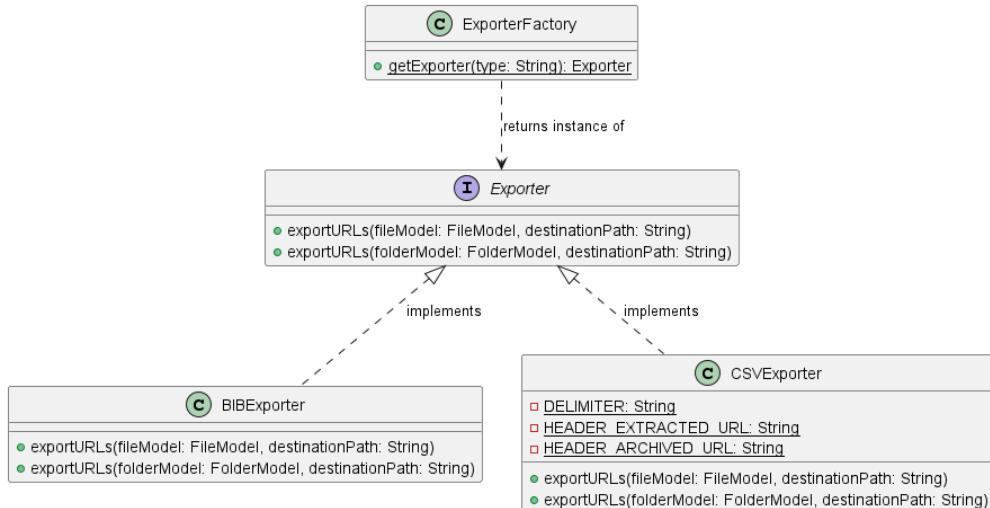


Figure 3.4.: Diagram of the Exporter Factory

increases the adaptability of the application, as configuration changes do not require code-base changes.

An example of the configuration file structure is as follows:

```
{
  "accessKey": "[Access Key]",
  "secretKey": "[Secretkey]",
```

```
    "browser": "FIREFOX"  
}
```

Users have the flexibility to modify the configuration file at runtime, or edit the file directly. Currently, the application uses this file to store credentials for the Wayback Machine and to select the correct browser for the Archive Today service. This feature underlines the URL-Archiver's commitment to user customisation and operational efficiency.

Adherence to SOLID Principles

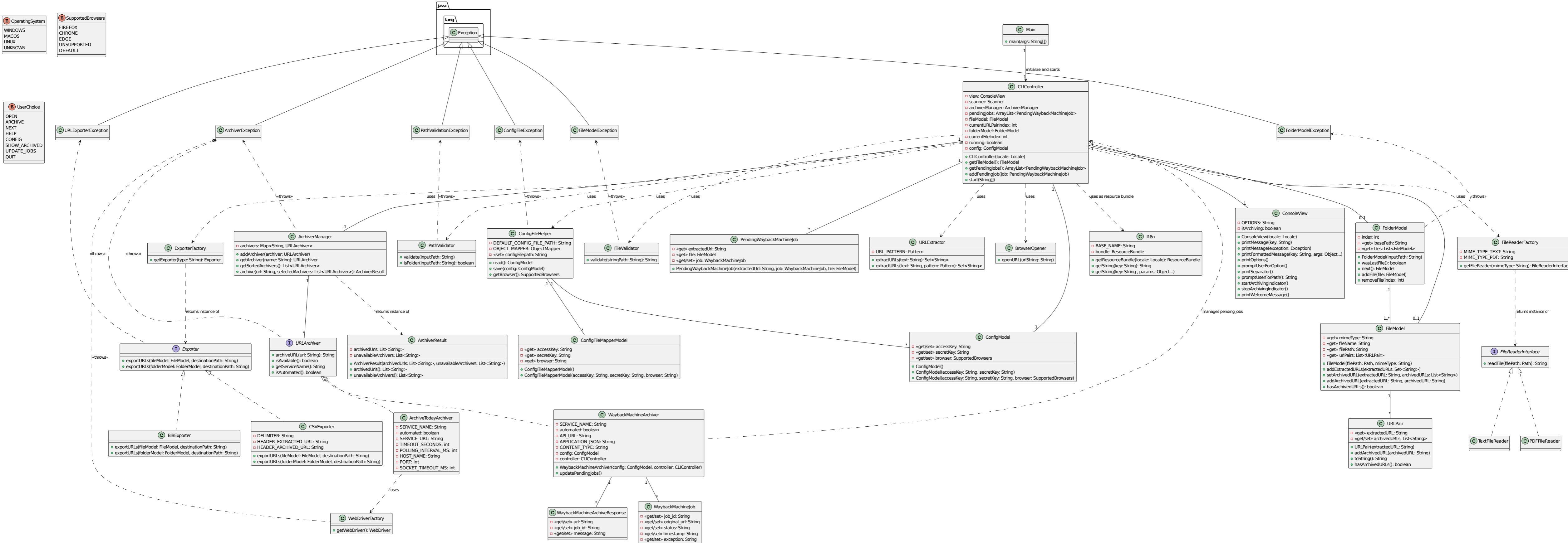
The URL-Archiver has been carefully designed with a strong emphasis on SOLID principles, ensuring a robust, maintainable and scalable architecture.

1. **Single Responsibility Principle:** Each class, such as `FileReaderFactory`, `ConfigModel` or `CLIController`, is dedicated to a single responsibility. This principle ensures that each module or class focuses on a single aspect of the application's functionality.
2. **Open/Closed Principle:** The application uses interfaces such as `FileReaderInterface` and factory classes to facilitate extensibility while keeping changes to existing code to a minimum. This approach allows new functionality to be added seamlessly.
3. **Liskov Substitution Principle:** The architectural design allows subclasses (such as specific file readers) to replace their parent class or interface without compromising the integrity of the application, thereby maintaining a robust design.
4. **Interface separation principle:** The design strategy of the URL Archive favours lean interfaces, ensuring that classes implement only the methods essential to their functionality. This is evident in the focused methods of interfaces such as `URLArchiver` and `Exporter`.
5. **Dependency Inversion Principle:** High-level modules such as `CLIController` rely on abstractions rather than concrete implementations. This is emphasised by the use of interfaces and factory classes for object creation, which reduces dependency on specific implementations.

The URL-Archiver's strict adherence to these SOLID principles is a testament to its well-structured, easily maintainable and scalable design. This approach underlines the application's commitment to high quality software engineering practices.

Class Diagram

Below the complete class diagram of this application is displayed.



3.2. Processes

3.2.1. Allocation of roles

In this chapter, the Scrum roles (Product Owner, Scrum Master, Developer) and additional roles such as Customer, Stakeholder, etc. are defined.

3.2.2. Scrum roles

We have decided to structure our Scrum team in the following manner:

| Role | Person |
|---------------|---|
| Product Owner | Nicolin Dora |
| Scrum Master | Abidin Vejseli |
| Developer | Nicolin Dora, Abidin Vejseli, Kilian Wampfler |

Table 3.1.: Scrum Roles

Nolin took on the role of Product Owner as he had concrete ideas and visions for the product at the start of the project. Additionally, he took on this role because he wanted to deal with the subjects surrounding the product backlog.

Abidin took on the role of the Scrum Master as he has the most experience with the agile way of working. He has already had the opportunity to perform this role professionally on several smaller projects in the past.

Kilian took on the role of a Developer, as he is an active programmer in his job and has already gained some experience with Scrum. Therefore, self-organization is not a foreign concept to him.

Besides Kilian, all the other members of the group were also assigned the role of Developer, as otherwise the project would not have been feasible in the given time. This is due to the fact that we all work alongside the university.

3.2.3. Additional roles

In addition to the Scrum roles, we have assigned the following roles to our specialist lecturer and PM-coach.

| Role | Person |
|-------------|------------------|
| Stakeholder | Dr. Simon Kramer |
| Customer | Dr. Simon Kramer |
| PM-Advisor | Frank Helbling |

Table 3.2.: Additional Scrum Roles

3.2.4. Sprint Goals

We have defined the goals of our past and current sprints in the best possible way according to the SMART¹ criteria. The goals of our sprints are listed below:

Sprint 1 Implement input handler for files (any unicode file e.g. .bib, .txt, .html and .pdf) and basic user guidance (Menu, Error messages).

Sprint 2 Implement a function to scan a provided text in order to identify and extract any URLs contained within it and upgrade our current console-based interface to enable users to easily open any extracted URL using their default web browser.

Sprint 3 Develop and implement a fully automated URL submission system that integrates with the Wayback Machine and Archive Today to ensure at least a 98% success rate in URL archiving.

Sprint 4 Enhance the system's stability and usability by resolving identified Selenium bugs across Linux, macOS, and Edge browsers, documenting the sprint process and licenses, conducting a thorough code review, and establishing a new configuration management file, aiming for zero critical bugs at sprint closure and readying the system for seamless URL archiving integration in subsequent sprints.

Sprint 5 Complete application refactoring for asynchronous archiving and .BIB file URL integration, ensuring no critical bugs and preparing for seamless future enhancements.

Sprint 6 Deliver a finalized application design, improved code quality, and complete documentation, with all components ready for review.

¹Atlassian blogpost: "How to write smart goals"

3.2.5. Requirements

In this chapter, we present our product and sprint backlogs, structured according to Scrum.

Product Backlog

Our product backlog consists of user stories and epics created by our Product Owner. The user stories are prioritised and represent a set of initial requirements that must be met to achieve our product goal. The product backlog is maintained in Jira.

| AV | ND | K | Type | Version | Epic | Status | AV |
|---|----|---|------------------------|---------|------|--------|----|
| Backlog (23 issues) | | | | | | | |
| SCRUM-4 Template | | | TO DO | | | | = |
| SCRUM-16 Store URL Line Number or Context | | | TO DO | | | | ~ |
| SCRUM-17 Compile a List of URLs | | | TO DO | | | | = |
| SCRUM-20 Immediate Archiving Upon Decision | | | WEB BROWSER INTEGR... | | | | ~ |
| SCRUM-21 Track Archiving Progress | | | WEB BROWSER INTEGR... | | | | ~ |
| SCRUM-22 Store User Decisions for Reporting | | | WEB BROWSER INTEGR... | | | | ~ |
| SCRUM-23 Automated URL Submission | | | INTERACTION WITH AR... | | | | ~ |
| SCRUM-25 User Interaction for Captchas | | | INTERACTION WITH AR... | | | | ~ |
| SCRUM-26 Automatic Retrieval of Archived URL | | | OUTPUT AND REPORTI... | | | | ~ |
| SCRUM-27 Generate CSV File | | | OUTPUT AND REPORTI... | | | | = |
| SCRUM-29 Integrate Archived URLs into Supported Files | | | DOCUMENTATION & P... | | | | ~ |
| SCRUM-31 User Manual | | | DOCUMENTATION & P... | | | | ~ |
| SCRUM-32 Intermediate presentation | | | DOCUMENTATION & P... | | | | = |
| SCRUM-33 Final presentation | | | DOCUMENTATION & P... | | | | ~ |
| SCRUM-34 UML diagram | | | DOCUMENTATION & P... | | | | ~ |
| SCRUM-35 Installation manual & script | | | DOCUMENTATION & P... | | | | ~ |
| SCRUM-36 Requirements document | | | DOCUMENTATION & P... | | | | = |

Figure 3.5.: Product Backlog

The prioritisation of user stories in the backlog is based on the business value field, which has a value between one and ten. The business value is a vague estimate of how much value the individual user story has to the business, or in our case, to our stakeholders. We endeavour to estimate the business value based on the expected importance of the function to the stakeholder. The following priorities are possible:

- ▶ Highest
- ▶ High
- ▶ Medium
- ▶ Low
- ▶ Lowest

Sprint Backlogs

Below we describe our recent and current sprints. During the sprint planning we fill the respective sprint backlog with user stories that serve the sprint goal. A user story must satisfy our Definition of Ready before it can be included in the sprint. In addition, the stories must be estimated and the total number of story points must not exceed our defined velocity.

Sprint 1 Below is a screenshot of our board from the first sprint with the corresponding sprint goal.

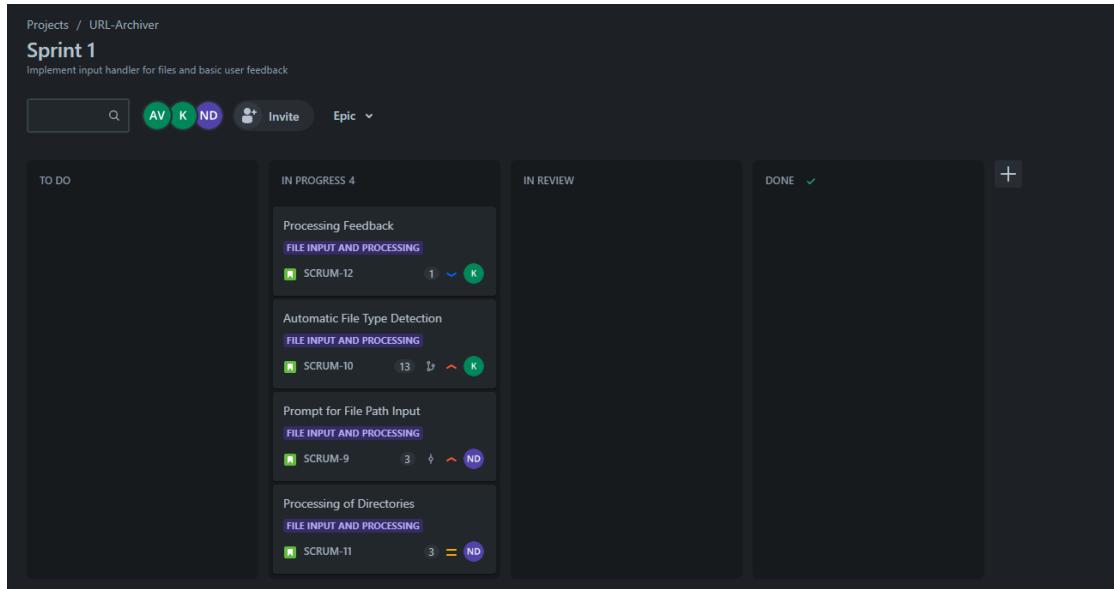


Figure 3.6.: Sprint 1 Backlog

The user stories from the first sprint are shown below. The stories have been estimated and prioritised.

The screenshot shows a user story detail page. At the top, there are two project names: SCRUM-8 and SCRUM-12. On the right side, there are several icons: a lock, a person, a thumbs up, a share icon, and three dots. Below the project names, the title "Processing Feedback" is displayed. Underneath the title are four small icons: a pen, a person, an envelope, and three dots. The main content area starts with a "Description" section, followed by a "Description:" heading. The text reads: "As a user, I want to get feedback when the tool starts processing the file and when it finishes reading it, to know the status." Below this is an "Acceptance Criteria:" section with two bullet points: "Once the file has been read, a confirmation message is displayed." and "If the file cannot be read an error or warning is thrown." To the right of the main content is a sidebar titled "Details". It contains the following information: Assignee (Kilian), Priority (Low), Business Value (3), and Story point estimate (1). There are also "Assign to me" and "Actions" buttons.

Figure 3.7.: User Story Detail for "Processing Feedback"

The screenshot shows a user story detail page. At the top, there are two project names: SCRUM-8 and SCRUM-9. On the right side, there are several icons: a lock, a person, a thumbs up, a share icon, and three dots. Below the project names, the title "Prompt for File Path Input" is displayed. Underneath the title are four small icons: a pen, a person, an envelope, and three dots. The main content area starts with a "Description" section, followed by a "Description:" heading. The text reads: "As a user, when I start the tool, I want to be prompted to input the path to my file, so the tool knows which file to process." Below this is an "Acceptance Criteria:" section with two bullet points: "Upon starting the tool, it prompts the user to enter a file path." and "On inputting an invalid path or if there are permissions issues, the tool provides a relevant error message." To the right of the main content is a sidebar titled "Details". It contains the following information: Assignee (Nicolin Dora), Priority (High), Business Value (8), and Story point estimate (3). There are also "Assign to me" and "Actions" buttons.

Figure 3.8.: User Story Detail for "Prompt for file path input"

The screenshot shows a user story detail page. At the top left, there are two labels: 'SCRUM-8 / SCRUM-10'. On the right side, there are several icons: a lock, a refresh, a person icon, a share icon, and a more options icon. Below these are two buttons: a green 'Done' button with a checkmark and a grey 'Done' button with a checkmark. To the right of these buttons is a 'Actions' dropdown menu. The main content area has a header 'Automatic File Type Detection'. Below the header are four small icons: a clipboard, a document, a gear, and three dots. The next section is titled 'Description' with a bold 'Description:' heading. The text reads: 'As a user, I want the tool to automatically detect the file type (based on file extension) and treat it accordingly so that I don't need to specify the file type separately.' Below this is a bold 'Acceptance Criteria:' heading followed by a bulleted list:

- The tool automatically identifies any unicode text file for example .BIB, .TEX, .HTML, or .PDF.
- For unrecognized file types, the tool provides an appropriate error message.
- The tool transforms the input file into a usable format.

On the right side of the page, there is a sidebar with a 'Details' tab. Under 'Assignee', it shows 'Kilian' with a green profile icon. There is a blue link 'Assign to me'. Under 'Priority', it shows 'High' with a red arrow icon. Under 'Business Value', it shows the number '8'. Under 'Story point estimate', it shows the number '13'.

Figure 3.9.: User Story Detail for "Automatic File Type Detection"

The screenshot shows a user story detail page. At the top left, there are two project names: SCRUM-8 and SCRUM-11. On the right side, there are several icons: a lock, a person icon with a '1', a thumbs up, a share icon, and a three-dot menu. Below these are two buttons: a green 'Done' button with a dropdown arrow and a checked green checkmark, followed by a 'Done' status indicator. To the right of these buttons is a 'Actions' dropdown with a gear icon. The main content area has a title 'Processing of Directories'. Below the title are four small icons: a clipboard, a document, a link, and a three-dot menu. A 'Description' section follows, containing a bolded 'Description:' heading and a text block: 'As a user, I want to input a whole directory, so the tool processes all supported files contained within.' Underneath this is a bolded 'Acceptance Criteria:' heading, followed by a bulleted list of five items. To the right of the main content area is a vertical sidebar with a grey background. On the far right edge of this sidebar, there is a horizontal scroll bar.

SCRUM-8 / SCRUM-11

Processing of Directories

Done ✓ Done

Actions

Description

Description:

As a user, I want to input a whole directory, so the tool processes all supported files contained within.

Acceptance Criteria:

- The tool can accept directory paths after the prompt.
- It processes all supported file types within the directory.
- The tool gives a message if files within the directory are skipped due to their type.
- For unrecognized directory paths, the tool provides an appropriate error message.

Details

Assignee: Nicolin Dora

Priority: Medium

Business Value: 6

Story point estimate: 3

Figure 3.10.: User Story Detail for "Processing of Directories"

In the first sprint, the burn down chart looks like this:



Figure 3.11.: Sprint 1 Burn Up Chart

The reason for this is that we did not create tasks for our user stories as they were small enough. Furthermore, a code review for corresponding user stories could only be conducted towards the end of the sprint, resulting in the finalisation of user stories at that point.

3. Implementation

Sprint 2 Below is a screenshot of our board from the second sprint with the corresponding sprint goal.

Sprint 2

Implement a function to scan a provided text in order to identify and extract any URLs contained within it and upgrade our current console-based interface to enable users to easily open any extracted URL using their default web browser.

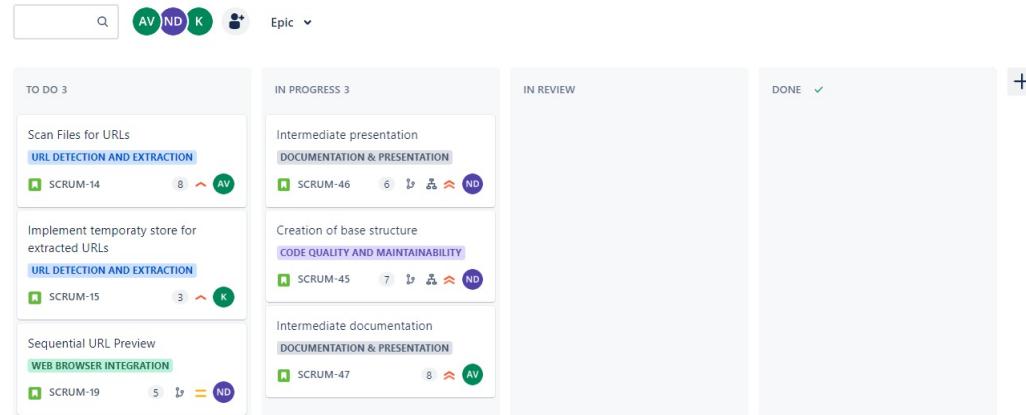


Figure 3.12.: Sprint 2 Backlog

The user stories from the second sprint are shown below. The stories have been estimated and prioritised.

Intermediate documentation

Description
User Story

As a project group member, I need to create a PDF document summarizing our project's presentation content, especially regarding Scrum, to provide detailed insights to stakeholders who might not attend the live presentation.

Acceptance Criteria

1. The PDF document effectively summarizes the main points of the presentation.
2. The documentation effectively communicates the project's problem statement, including goals, requirements, and constraints.
3. Information specific to Scrum, such as roles, Sprint objectives, and Scrum adaptations, is highlighted and detailed in the document.
4. The document is professionally formatted and easy to understand.
5. All relevant stakeholders, such as Mr. Helbling and the project supervisor, are emailed the document in advance.
6. The documentation effectively communicates the project's problem statement, including goals, requirements, and constraints.
7. The documentation is well-structured with parts that are coherent and coordinated.
8. The initial situation related to the topic is clearly analyzed, stakeholders identified, and any hypotheses formulated.
9. The goal of the deliverable is defined, allowing an assessment of goal achievement, with boundaries set.

Figure 3.13.: User Story Detail for "Intermediate Documentation"

3. Implementation

The screenshot shows a user story detail page. At the top, there are navigation links for SCRUM-39 and SCRUM-45. Below that, the title "Creation of base structure" is displayed, along with standard issue management buttons: Attach, Add a child issue, Link issue, and three dots. The main content area starts with a "Description" section containing the text: "As a developer, I want to review and refine the spontaneously created base structure of our project to ensure it aligns with the MVC pattern, so that we can have a solid foundation for our application and avoid future technical debt." Below this is an "Acceptance Criteria" section with the following list:

1. The current base structure is thoroughly reviewed to identify any deviations or inefficiencies in relation to the MVC pattern.
2. Any non-aligned components or code snippets are refactored to fit the MVC design pattern.
3. All model, view, and controller components are clearly separated and interact seamlessly.
4. Existing functionalities of the application remain intact post-refinement.
5. All code changes are documented, highlighting the reasons for modifications.
6. Peer review is conducted, ensuring that at least two other team members agree on the changes made.
7. After refactoring, the application should successfully pass all existing unit tests.

To the right of the main content is a sidebar titled "Details" which contains the following information:

| | |
|----------------------|--------------|
| Assignee | Nicolin Dora |
| Priority | Highest |
| Business Value | 9 |
| Story point estimate | 7 |

Figure 3.14.: User Story Detail for "Creation of Base Structure"

The screenshot shows a user story detail page. At the top, there are navigation links for SCRUM-13 and SCRUM-14. Below that, the title "Scan Files for URLs" is displayed, along with standard issue management buttons: Attach, Add a child issue, Link issue, and three dots. The main content area starts with a "Description" section containing the text: "As a user, I want the system to scan my input files and identify any embedded URLs so that they can be extracted for archiving." Below this is an "Acceptance Criteria" section with the following list:

- System can detect URLs in a variety of file formats including .BIB, .TEX, .HTML, and .PDF.
- System uses a robust regular expression or other reliable techniques to extract URLs pattern that matches most URL formats.
- Extracted URLs are validated to ensure they are in the correct format.

To the right of the main content is a sidebar titled "Details" which contains the following information:

| | |
|----------------------|----------------|
| Assignee | Abidin Vejseli |
| Priority | High |
| Business Value | 8 |
| Story point estimate | 8 |

Figure 3.15.: User Story Detail for "Scan Files for URLs"

3. Implementation

The screenshot shows a user story detail page. At the top, there are navigation links for SCRUM-18 and SCRUM-19, and a header with a lock icon, a comment count (1), a like icon, a share icon, and a more options icon. Below the header, the title is "Sequential URL Preview". There are buttons for "Attach", "Add a child issue", "Link issue", and three dots. A "Done" button is highlighted in green. To the right, there's a "Details" section with the following fields:

| | |
|----------------------|--------------|
| Assignee | Nicolin Dora |
| Priority | Medium |
| Business Value | 6 |
| Story point estimate | 5 |

Description:
As a user, I want to preview each detected URL in my default browser sequentially to verify its content.

Acceptance Criteria:

- System opens one URL at a time in the default browser.
- Immediately after the URL is displayed, the system presents the user with the option to archive.

Figure 3.16.: User Story Detail for "Sequential URL Preview"

The screenshot shows a user story detail page. At the top, there are navigation links for SCRUM-13 and SCRUM-15, and a header with a lock icon, a comment count (1), a like icon, a share icon, and a more options icon. Below the header, the title is "Implement temporary store for extracted URLs". There are buttons for "Attach", "Add a child issue", "Link issue", and three dots. A "Done" button is highlighted in green. To the right, there's a "Details" section with the following fields:

| | |
|----------------------|--------|
| Assignee | Kilian |
| Priority | High |
| Business Value | 7 |
| Story point estimate | 3 |

Description:
As a user, I want the system to store the captured URLs so that they can be subsequently archived.

Acceptance Criteria:

- The system efficiently stores collected URLs in a data structure
- The data structure is implemented as a set

Figure 3.17.: User Story Detail for "Implement temporary store for extracted URLs"

3. Implementation

The screenshot shows a user story detail view for an "Intermediate presentation". The top navigation bar includes links for SCRUM-30 and SCRUM-46, and icons for lock, refresh, and more. The main interface has tabs for "Done" (green), "In Progress" (orange), and "Upcoming" (blue). Below the tabs are sections for "Description", "User Story", and "Acceptance Criteria". The "User Story" section contains a detailed description of the requirement. The "Acceptance Criteria" section lists 10 numbered points. To the right, a "Details" panel is open, showing fields for "Assignee" (Nicolin Dora), "Priority" (Highest), "Business Value" (8), and "Story point estimate" (6). There are also "Assign to me" and "Actions" buttons.

Intermediate presentation

SCRUM-30 / SCRUM-46

Done In Progress Upcoming

Attach Add a child issue Link issue ...

Description

User Story

As a project group member, I need to deliver a 20-minute presentation about our project so that the audience understands our project's problem statement, our proposed solutions, and our use of Scrum.

Acceptance Criteria

1. The presentation lasts for a total of 20 minutes, leaving room for questions and discussions.
2. All group members actively participate in the presentation.
3. Cameras are turned on throughout the presentation.
4. The presentation effectively communicates the project's problem statement, including goals, requirements, and constraints.
5. Clearly depicts the proposed solution(s) showcasing architecture, data model, process model, technologies, etc.
6. Explains how the project is being managed using Scrum methodology.
7. The content of the presentation is relevant, factually correct, supported by arguments, and reflects deep engagement with the topic.
8. The presentation is well-structured with parts that are coherent and coordinated.
9. Any media used, such as slides or videos, is appealing, meticulously prepared, and professionally utilized.
10. The language used is convincing, the speech is clear, and both contribute to audience comprehension.

Done ✓ Done Actions

Details

| | |
|----------------------|--------------|
| Assignee | Nicolin Dora |
| Priority | Highest |
| Business Value | 8 |
| Story point estimate | 6 |

Assign to me

Figure 3.18.: User Story Detail for "Intermediate presentation"

In the second sprint, the burn down chart looks like this:

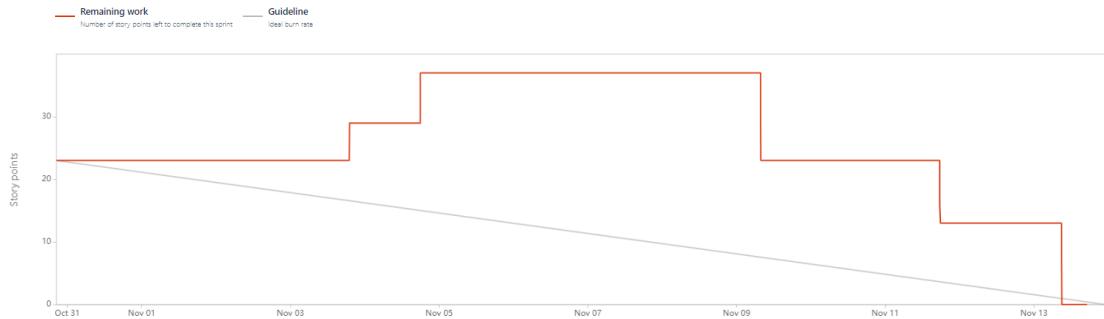


Figure 3.19.: Sprint 2 Burn Down Chart

Compared to the initial sprint, our workload significantly increased, and we introduced new user stories after the sprint began. Despite these additions, we maintained a strong pace and seamlessly managed the extra tasks that arose from the intermediate presentation and documentation requirements.

3. Implementation

Sprint 3 Below is a screenshot of our board from the third sprint with the corresponding sprint goal.

Projects / URL-Archiver
Sprint 3
Develop and implement a fully automated URL submission system that integrates with the Wayback Machine and Archive Today to ensure at least a 98% success rate in URL archiving.

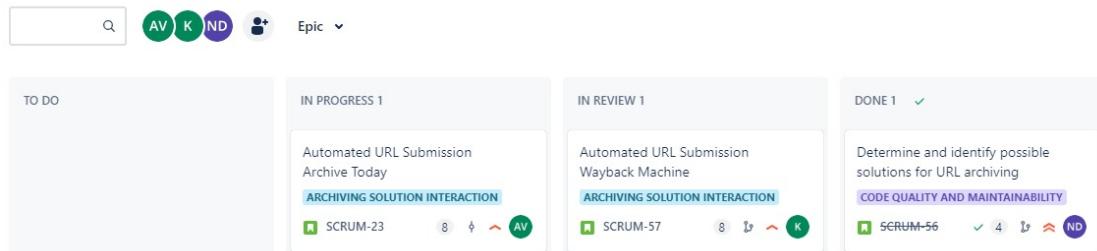


Figure 3.20.: Sprint 3 Backlog

The user stories from the third sprint are shown below. The stories have been estimated and prioritised.

The image shows a Jira user story detail page for 'Automated URL Submission Archive Today'. The story has a priority of 'High' and a business value of 8. The assignee is 'Abidin Vejseli' and the story point estimate is 8. The acceptance criteria listed are:

- Upon initiation, system opens the [archive.ph](#) website in a browser.
- System auto-fills the given URL into the appropriate input field.
- System automatically triggers the submission process for archiving.

Figure 3.21.: User Story Detail for "Intermediate Documentation"

3. Implementation

The screenshot shows a Jira user story detail page. At the top, there are navigation links for SCRUM-39 and SCRUM-56. Below the title "Determine and identify possible solutions for URL archiving", there are buttons for Attach, Add a child issue, Link issue, and more. The main content area includes sections for Description, Description, Acceptance criteria, and a detailed description of the requirements. On the right side, there is a sidebar titled "Details" showing project metadata: Assignee (Nicolin Dora), Priority (Highest), Business Value (7), Story point estimate (4), and a Done button.

Description

Description:

As a software developer, I want to lay the groundwork for a URL archiving system that can be expanded with various archiving services to meet future requirements.

Acceptance criteria

1. The existence of an `URLArchiver` interface that specifies the necessary methods for archiving services.
2. Placeholder classes for `waybackMachineArchiver` and `ArchiveTodayArchiver` indicating the intent to support these services.
3. An `ArchiverManager` class capable of managing a collection of archivers and orchestrating the archiving process.
4. A `CLIController` class designed to interact with the user and control the archiving flow based on user input.
5. A result-handling mechanism through the `ArchiverResult` class, to group the outcomes of archiving operations.
6. A `messages.properties` file to handle system messages, ensuring they are ready for future localization and user-friendly communication.
7. Preliminary JUnit tests for the `ArchiverManager` and `ArchiverResult` to validate the system's core functionality.
8. JavaDoc documentation for key interfaces and classes to provide an understanding of their purpose and guide future development.

Figure 3.22.: User Story Detail for "Creation of Base Structure"

The screenshot shows a Jira user story detail page. At the top, there are navigation links for SCRUM-24 and SCRUM-57. Below the title "Automated URL Submission Wayback Machine", there are buttons for Attach, Add a child issue, Link issue, and more. The main content area includes sections for Description, Description, Acceptance Criteria, and a detailed description of the requirements. On the right side, there is a sidebar titled "Details" showing project metadata: Assignee (Kilian), Priority (High), Business Value (8), Story point estimate (8), and a Done button.

Description

Description:

As a user, I want the system to automatically archive the URL via the WaybackMachine API.

Acceptance Criteria:

- Upon initiation, system connects to the WaybackMachine API.
- The URL is archived on the WaybackMachine
- The URL to the archived URL is returned to the system.

Figure 3.23.: User Story Detail for "Scan Files for URLs"

In the third sprint, the burn down chart looks like this:

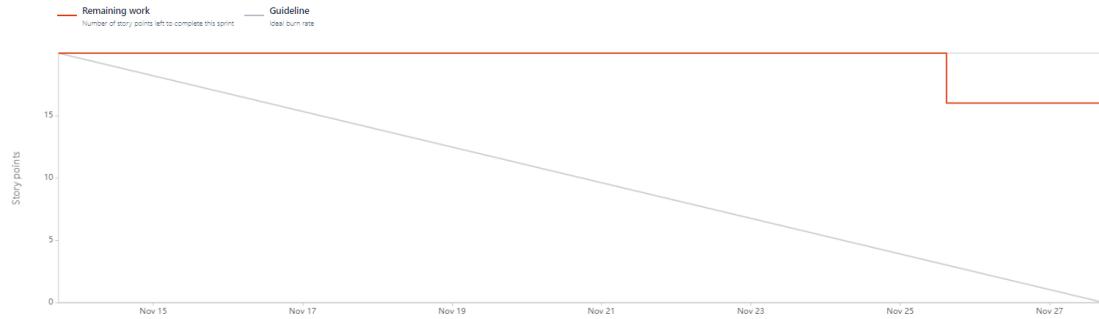


Figure 3.24.: Sprint 3 Burn Down Chart

During the third sprint, our progress on user stories was limited to the completion of only three due to the 'special week 3'. This meant that these stories had to be completed in the second half of the sprint due to the heavy workload of this special week. The impact of this adaptation to our sprint schedule due to 'special week 3' is reflected in the trends seen in our burndown chart.

3. Implementation

Sprint 4 Below is a screenshot of our board from the forth sprint.

| Key | Summary | Issue type | Epic | Status | Assignee | Story points |
|----------|--|------------|----------------|--------|----------|--------------|
| SCRUM-27 | Generate CSV File | Story | OUTPUT AN... | DONE | K | 4 |
| SCRUM-58 | No URL found in file is not handled | Bug | FILE INPUT ... | DONE | ND | 5 |
| SCRUM-60 | Fix bug with selenium on linux | Story | ARCHIVING ... | DONE | ND | 3 |
| SCRUM-61 | Fix bug with selenium on MacOS | Story | ARCHIVING ... | DONE | K | 3 |
| SCRUM-62 | Fix bug with selenium and Edge browser | Story | ARCHIVING ... | DONE | AV | 3 |
| SCRUM-63 | Document licences | Story | DOCUMENT... | DONE | ND | 4 |
| SCRUM-64 | Periodic complete code review | Story | CODE QUALI... | DONE | ND | 4 |
| SCRUM-65 | Document SCRUM sprint 3 | Story | DOCUMENT... | DONE | AV | 4 |
| SCRUM-66 | Creation of a config file | Story | CODE QUALI... | DONE | K | 4 |
| SCRUM-75 | Progress indicator archiving | Story | ARCHIVING ... | DONE | ND | 4 |

Figure 3.25.: Sprint 4 Backlog

The user stories from the forth sprint are shown below. The stories have been estimated and prioritised.

The screenshot shows the details of a Jira issue titled "Generate CSV File". The issue key is SCRUM-28. The summary is "Generate CSV File". The issue type is "Story". The status is "Done". The assignee is "Kilian". The priority is "High". The business value is 8. The story point estimate is 4. The description states: "As a user, I want the system to produce a CSV file containing all original URLs and their corresponding archived URLs." The acceptance criteria are: "A CSV file is generated upon completion of the archiving process." and "Each row in the CSV contains the original URL and its archived counterpart." There are buttons for "Attach", "Add a child issue", "Link issue", and "Actions".

Figure 3.26.: User Story Detail for "Generate CSV File"

3. Implementation

The screenshot shows a user story detail page. At the top, there are navigation links for 'SCRUM-24' and 'SCRUM-60'. Below the title 'Fix bug with selenium on linux', there are buttons for 'Attach', 'Add a child issue', 'Link issue', and more. The 'Description' section contains the following text:
As a developer, I want Selenium to function without errors on Linux, enabling Linux and Firefox users to utilise the application.
The 'Acceptance Criteria' section lists:

- Archiving a URL on Linux should no longer throw errors.
- Critical errors should still be thrown.

On the right, there is a 'Details' panel with the following information:

| Details | |
|----------------------|--------------|
| Assignee | Nicolin Dora |
| Priority | Medium |
| Business Value | 4 |
| Story point estimate | 3 |

Figure 3.27.: User Story Detail for "Fix bug with selenium on linux"

The screenshot shows a user story detail page. At the top, there are navigation links for 'SCRUM-24' and 'SCRUM-61'. Below the title 'Fix bug with selenium on Mac OS', there are buttons for 'Attach', 'Add a child issue', 'Link issue', and more. The 'Description' section contains the following text:
As a developer, I want Selenium to function without errors on Mac OS, enabling Mac OS users to utilise the application.
The 'Acceptance Criteria' section lists:

- Archiving a URL on Mac OS should no longer throw errors.
- Critical errors should still be thrown.
- Archiving through Chrome or Firefox is achievable.
- Archiving through Safari is not supported.

On the right, there is a 'Details' panel with the following information:

| Details | |
|----------------------|--------|
| Assignee | Kilian |
| Priority | Medium |
| Business Value | 4 |
| Story point estimate | 3 |

Figure 3.28.: User Story Detail for "Fix bug with selenium on Mac OS"

3. Implementation

The screenshot shows a user story detail page. At the top, there are navigation links for 'SCRUM-24 / SCRUM-62'. On the right, there are icons for lock, refresh, like, share, and more. Below the title 'Fix bug with selenium and Edge browser', there are buttons for 'Attach', 'Add a child issue', 'Link issue', and three dots. The main content area has a 'Description' section with a 'Description:' heading and a text block: 'As a developer, I want Selenium to function without errors on Windows, enabling Windows users to utilise the application.' Below this is an 'Acceptance Criteria:' section with a bulleted list: 'Archiving a URL on Windows should no longer throw errors.', 'Critical errors should still be thrown.', and 'Archiving through Edge, Chrome or Firefox is achievable.' To the right, a sidebar titled 'Details' shows fields: 'Assignee' (Abidin Vejseli), 'Priority' (Medium), 'Business Value' (4), and 'Story point estimate' (3). A green 'Done' button is at the top of the sidebar.

Figure 3.29.: User Story Detail for "Fix bug With selenium and Edge browser"

The screenshot shows a user story detail page. At the top, there are navigation links for 'SCRUM-30 / SCRUM-63'. On the right, there are icons for lock, refresh, like, share, and more. Below the title 'Document licences', there are buttons for 'Attach', 'Add a child issue', 'Link issue', and three dots. The main content area has a 'Description' section with a 'Description:' heading and a text block: 'As a developer, I want to document the licenses of the libraries used so that it is clear to future readers what we are using and under which license.' Below this is an 'Acceptance Criteria:' section with a bulleted list: 'All libraries used are documented.' and 'The license is documented for each library.' To the right, a sidebar titled 'Details' shows fields: 'Assignee' (Nicolin Dora), 'Priority' (Medium), 'Business Value' (4), and 'Story point estimate' (4). A green 'Done' button is at the top of the sidebar.

Figure 3.30.: User Story Detail for "Document licences"

The screenshot shows a user story detail page. At the top, there are navigation links for 'SCRUM-39 / SCRUM-64'. On the right, there are icons for lock, refresh, like, share, and more. Below the title 'Periodic complete code review', there are buttons for 'Attach', 'Add a child issue', 'Link issue', and three dots. The main content area has a 'Description' section with a 'Description:' heading and a text block: 'As a software development team, we need to conduct periodic complete code reviews, so that we can ensure our code remains minimal, self-explanatory, and modular.' Below this is an 'Acceptance Criteria' section with a bulleted list: 'The code is organised into well-defined and logically distinct modules or components.', 'The code must have descriptive naming conventions and provide an appropriate level of commentary.', and 'Avoid unnecessary complexity and focus on streamlined solutions.' To the right, a sidebar titled 'Details' shows fields: 'Assignee' (Nicolin Dora), 'Priority' (Medium), 'Business Value' (4), and 'Story point estimate' (4). A green 'Done' button is at the top of the sidebar.

Figure 3.31.: User Story Detail for "Periodic complete code review"

3. Implementation

The screenshot shows a user story detail page. At the top, there are navigation links for 'SCRUM-30 / SCRUM-65'. On the right, there are icons for lock, eye, like, share, and more. Below the title 'Document SCRUM sprint 3', there are buttons for 'Attach', 'Add a child issue', 'Link issue', and three dots. The main content area has a 'Description' section with a bold 'Description:' heading. The text reads: 'As a developer, I want to document Sprint 3 so that future readers of our documentation can see what we achieved and how efficient we were.' Below this is an 'Acceptance Criteria:' section with a bulleted list: 'All user stories are included', 'The burnup or burndown chart is present', and 'The sprint goal is defined'. To the right, there is a sidebar titled 'Details' with sections for 'Assignee' (Abidin Vejseli), 'Priority' (Low), 'Business Value' (4), and 'Story point estimate' (4). There are also 'Done' and 'Actions' buttons at the top of the sidebar.

Figure 3.32.: User Story Detail for "Document SCRUM sprint 3"

The screenshot shows a user story detail page. At the top, there are navigation links for 'SCRUM-39 / SCRUM-66'. On the right, there are icons for lock, eye, like, share, and more. Below the title 'Creation of a config file', there are buttons for 'Attach', 'Add a child issue', 'Link issue', and three dots. The main content area has a 'Description' section with a bold 'Description:' heading. The text reads: 'As a user, I want the application to manage a configuration file that stores my user-specific data such as API keys. During start-up, the application must verify the existence of the config file. If absent, the application must ask for the required information, and create the file using those details. If the configuration file already exists, the application must automatically read and use the information it contains, streamlining my experience by eliminating the need to repeatedly enter the same information.' Below this is an 'Acceptance Criteria:' section with a bulleted list: 'Check if config file exists' (with sub-points 'no config file exists: prompt user for the needed information and create the config file' and 'config file already exists: read the config file'), 'the config file contains the api keys for the wayback machine', 'the config file can easily be extended in the future', and 'write a manual about how to get the api keys for the wayback machine'. To the right, there is a sidebar titled 'Details' with sections for 'Assignee' (Kilian), 'Assign to me', 'Priority' (Medium), 'Business Value' (6), and 'Story point estimate' (4). There are also 'Done' and 'Actions' buttons at the top of the sidebar.

Figure 3.33.: User Story Detail for "Creation of a config file"

3. Implementation

The screenshot shows a user story detail page. At the top, there are navigation links: SCRUM-24 / SCRUM-75. Below that is a header for "Progress indicator archiving". A toolbar contains "Attach", "Add a child issue", "Link issue", and a "Done" button. The main content area has sections for "Description" and "Acceptance Criteria". The "Description" section contains a text block about wanting to see a progress indicator during the archiving process. The "Acceptance Criteria" section lists several bullet points regarding the behavior of the progress indicator. To the right, there is a sidebar titled "Details" which includes fields for "Assignee" (Nicolin Dora), "Priority" (Medium), "Business Value" (6), and "Story point estimate" (4). There are also "Actions" and "Done" buttons.

Figure 3.34.: User Story Detail for "Progress indicator archiving"

The screenshot shows a bug detail page. At the top, there are navigation links: SCRUM-8 / SCRUM-58. Below that is a header for "No URL found in file is not handled". A toolbar contains "Attach", "Add a child issue", "Link issue", and a "Done" button. The main content area has a "Description" section with a text block explaining that the `processFileModel` method fails if a file contains no URLs, leading to an error. To the right, there is a sidebar titled "Details" which includes fields for "Assignee" (Nicolin Dora) and "Actions" buttons.

Figure 3.35.: Bug Detail for "No URL found in file is not handled"

In the forth sprint, the burn down chart looks like this:

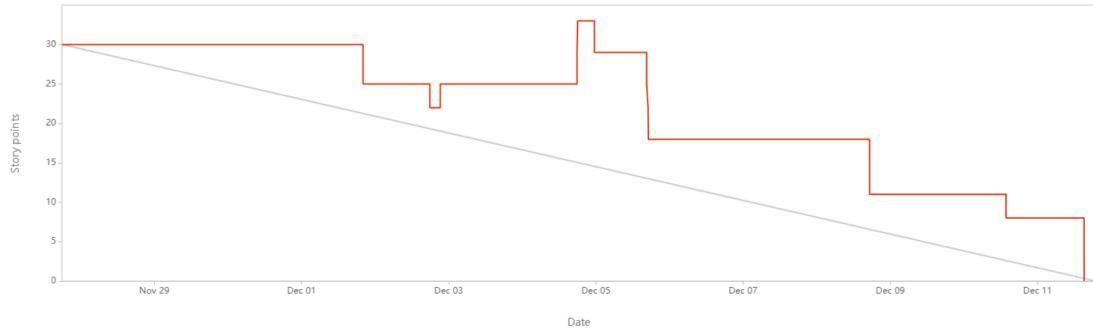


Figure 3.36.: Sprint 4 Burn Down Chart

Compared to the previous sprint, we were more efficient and successfully accommodated additional user stories that were initiated after the sprint began. The completion of all allocated user stories within the sprint timeframe demonstrates our solid teamwork and sprint management capabilities.

3. Implementation

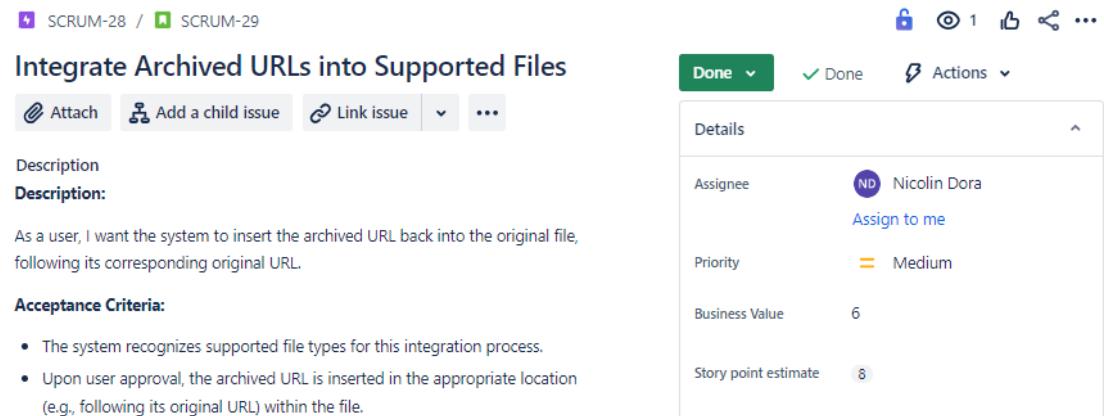
Sprint 5 Below is a screenshot of our board from the fifth sprint.

| Incomplete issues | | | | | | | View in issue navigator |
|-------------------|---|--------------|------------------------|-------------|------------|--------------|---|
| Key : | Summary : | Issue type : | Epic: | Status : | Assignee : | Story points | |
| SCRUM-79 | Refactor quit path | Story | ARCHIVING SOLUTION ... | IN PROGRESS | AV | 5 | |
| SCRUM-68 | If the browser gets manually closed by the user, there is a big error. Can... | Bug | | TO DO | AV | 3 | |

| Completed issues | | | | | | | View in issue navigator |
|------------------|---|--------------|------------------------|----------|------------|--------------|---|
| Key : | Summary : | Issue type : | Epic: | Status : | Assignee : | Story points | |
| SCRUM-67 | The application crashes if no valid file is found in the specified folder | Bug | | DONE | AV | 3 | |
| SCRUM-82 | Show archived urls path | Story | OUTPUT AND REPORTI... | DONE | ND | 6 | |
| SCRUM-72 | Create / Fix Tests | Story | CODE QUALITY AND M... | DONE | AV | 13 | |
| SCRUM-81 | Status update jobs | Story | OUTPUT AND REPORTI... | DONE | K | 6 | |
| SCRUM-80 | Refactor Archive path | Story | ARCHIVING SOLUTION ... | DONE | K | 13 | |
| SCRUM-29 | Integrate Archived URLs into Supported Files | Story | OUTPUT AND REPORTI... | DONE | ND | 8 | |

Figure 3.37.: Sprint 5 Backlog

The user stories from the fifth sprint are shown below. The stories have been estimated and prioritised.



The screenshot shows a Jira user story detail page. At the top, there are two story keys: SCRUM-28 and SCRUM-29. Below the title "Integrate Archived URLs into Supported Files", there are buttons for "Attach", "Add a child issue", "Link issue", and a three-dot menu. The main content area includes sections for "Description", "Acceptance Criteria:", and "Details". The "Description" section contains the story text: "As a user, I want the system to insert the archived URL back into the original file, following its corresponding original URL.". The "Acceptance Criteria:" section lists two bullet points: "The system recognizes supported file types for this integration process." and "Upon user approval, the archived URL is inserted in the appropriate location (e.g., following its original URL) within the file.". The "Details" sidebar on the right shows the following information: Assignee (Nicolin Dora), Priority (Medium), Business Value (6), and Story point estimate (8). There are also buttons for "Done" (with a green checkmark) and "Actions".

Figure 3.38.: User Story Detail for "Integrate Archived URLs into Supported Files"

3. Implementation

SCRUM-39 / SCRUM-72

Create / Fix Tests

Attach Add a child issue Link issue ...

Description

Description:

As a developer, I would like to have a test for every class where it makes sense, so that the probability of bugs is minimised.

Acceptance Criteria:

- There is a test for each class where it makes sense.
- The tests are clearly structured.
- The tests cover a large part of the functionalities.
- Where tests are not passed, bugs are created in Jira.

Done ✓ Done Actions ...

Details

Assignee AV Abidin Vejseli

Priority Medium

Business Value 4

Story point estimate 13

Figure 3.39.: User Story Detail for "Create / Fix Tests"

SCRUM-28 / SCRUM-82

Show archived urls path

Attach Add a child issue Link issue ...

Description

Description:

As a user, I would like to be able to view the archive URLs of already archived URLs at any time so that I can open them.

Acceptance Criteria:

- The user must have an option with which he can display the archive URLs.
- The user must be able to navigate through the archive URLs and open them.
- The user must be able to return to the main menu.

Done ✓ Done Actions ...

Details

Assignee ND Nicolin Dora Assign to me

Priority High

Business Value 7

Story point estimate 6

Figure 3.40.: User Story Detail for "Show archived urls path"

In the fifth sprint, the burn down chart looks like this:

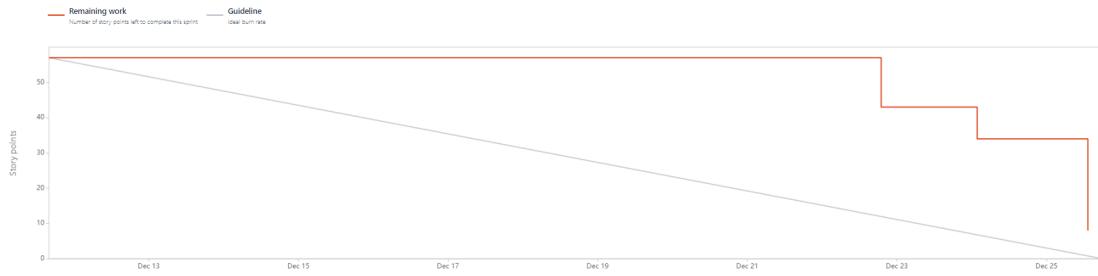


Figure 3.41.: Sprint 5 Burn Down Chart

In sprint 5, our time was limited due to commitments in other modules and the upcoming Christmas period, which resulted in team member absences. We encountered an unexpected challenge with the extensive time required for test refactoring, as we prioritise quality, which often demands more time. As a result, we were unable to complete all planned user stories and had to carry them over to the next sprint. These challenges are reflected in the burn down chart.

3.2.6. Scrum Adaptionen

As part of Project 1, we have adjusted Scrum in order to use it in the best possible way. The adjustments are explained in this chapter.

Definition of Ready (DOR)

Our DOR includes conditions that ensure that all team members understand the user stories and know when a user story can be included in a sprint. The DOR was set in line with the INVEST² criteria. A user story in the product backlog must meet the DOR before it can be included in a sprint.

Definition of Ready

- ▶ Ensure a clear definition
- ▶ Define the functionality or requirement to be implemented
- ▶ Clearly defined and testable acceptance criteria
- ▶ Ensure there are no or minimal dependencies
- ▶ Understood by the whole team
- ▶ The user story has been estimated
- ▶ The scope of the user story is small enough that it can be implemented in a single sprint.

²XP123 article: Invest in good stories and smart tasks

Definition of Done (DOD)

Our DOD contains all the characteristics and standards that a user story must meet to be considered complete. Once it satisfies the necessary quality requirements (acceptance criteria), the story can be considered complete and can be closed. The goal of our DOD is to create transparency so that everyone has a common understanding of when a story can be closed. A story that does not comply with the DOD may not be finalised.

Definition of Done

- ▶ Coding standards and best practices are implemented
- ▶ Unit tests for the feature are written and passed
- ▶ Any changes to the code or functionality are documented
- ▶ The code and functionality are reviewed by peers
- ▶ The feature works across multiple platforms
- ▶ Code is integrated with master branch
- ▶ Documentation has been updated
- ▶ Acceptance criteria are met

User Story Template

For the creation of a user story, we have defined a template so that the user stories contain all the necessary information. Below is a screenshot of our template. It includes all the relevant fields for us: Assignee, Priority, Business Value, Story Points estimate and assigned Sprint. Furthermore, we describe the user story in the "Description" field in the format "AS A <user role> I WANT TO <the goal> [SO THAT <reason>]" as well as the Acceptance Criteria. To ensure that we always have the DOR and DOD to hand, we also work with the Jira On-the-Fly add-on, which enables us to record both for each user story and tick off the individual points accordingly when they have been completed. This allows us to immediately recognise whether a user story can be included in a sprint and whether a story has been fully completed.

The screenshot shows the Jira 'Template' view for a user story titled 'User Story'. The left sidebar contains buttons for 'Attach', 'Add a child issue', 'Link issue', and a three-dot menu. The main area is divided into sections:

- Description:** Contains the title 'User Story'.
- Acceptance Criteria:** Contains the heading 'Checklist for Jira On-the-Fly' with a '+' button and a three-dot menu.
- Definition of Ready:** Contains a progress bar at 0/100% and a list of items:
 - Ensure a clear definition
 - Define the functionality or requirement to be implemented
 - Clearly defined and testable acceptance criteria
 - Ensure there are no or minimal dependencies
 - Understood by the whole team
 - The user story has been estimated
 - The scope of the user story is small enough that it can be implemented in a single sprint.
- Definition of Done:** Contains a progress bar at 0/100% and a list of items:
 - Coding standards and best practices are implemented
 - Unit tests for the feature are written and passed

The right side of the screen displays detailed project information under 'Details' and 'Automation' sections, including:

- Details:**
 - Assignee: Unassigned (with 'Assign to me' link)
 - Priority: Medium
 - Business Value: None
 - Story point estimate: None
 - Sprint: None
 - Fix versions: None
 - Original estimate: 0m
 - Development: Create branch, Create commit
 - Reporter: Abidin Vejseli
- Automation:** Rule executions

At the bottom, it shows the creation date ('Created October 11, 2023 at 5:38 PM') and the last update ('Updated 2 minutes ago').

Figure 3.42.: Screenshot from the user story template.

Estimation method

We have chosen the "T-shirt sizes" method because it is a simple way to estimate effort (story points). This method is based on the fact that everyone knows T-shirt sizes and that large sizes mean more work than small sizes. As a result, this method enables us to make efficient estimations, despite the lack of shared experience in the team.

Below is the scale we use:

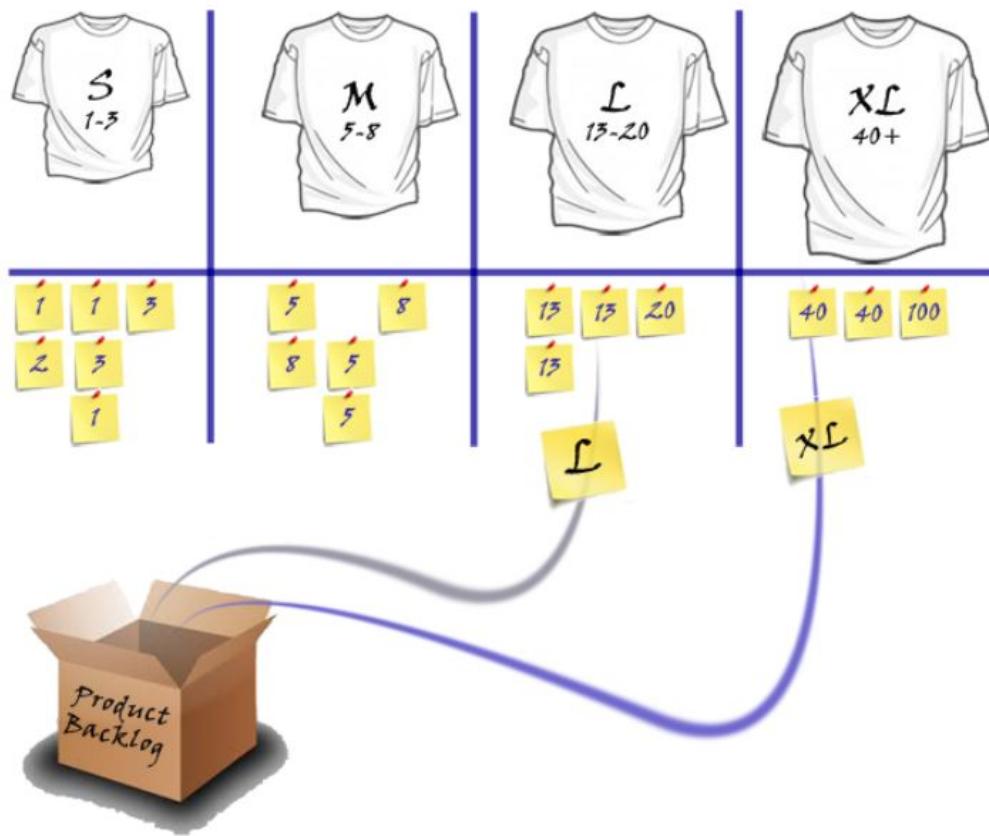


Figure 3.43.: Illustration of T-Shirt Sizes

Velocity

To select the user stories and tasks to be worked on, a suitable criterion, velocity, is applied to estimate what can be completed in the upcoming sprint.

We have decided that we have a velocity of 30 story points per sprint. Therefore, our workload per sprint should not exceed this threshold. We consciously take this into account during sprint planning.

Sprint

As a team, we have decided that our sprints will take place at two-week intervals. For each sprint we define a SMART sprint goal, which specifies the relevant user stories.

We decided in favour of the two-week rhythm because regular feedback is important to us and thus creates a greater learning effect. Additionally, we ascertained that one-week sprints would result in excessive overheads due to the administrative work involved in Scrum. Likewise, we consider sprints longer than two weeks to be impractical, as the interaction would suffer.

Sprint Planning

As part of sprint planning, we make decisions about which user stories can be implemented based on the sprint goal, the story points and the velocity. Before a user story can be included in the sprint, it must be estimated by the Scrum team. This task is always carried out at the start of our sprint planning. A sprint goal is then defined based on the estimated and prioritised user stories. The user stories we select for the sprint are based on the business value, priority, story points and velocity of the team. The sprint planning takes place on the first Monday of each sprint. We have decided not to have a second sprint planning as we have already defined in our DOR that the user stories should be as small as possible. In addition, each developer has the opportunity to divide their user stories into tasks within the sprint. This allows a better overview of the progress in the sprint.

Daily Scrum

As a team, we have decided not to have daily Scrum meetings, as this is not possible because all team members do work part times. Instead, we have two weekly meetings (weeklys), on Wednesday and Friday at 17:00, which last a maximum of 15 minutes. In addition, we have chosen to hold the meetings through Microsoft Teams as it is easier to organise. The goal of these meetings is to share the current progress, address issues, update the team and briefly discuss the next steps.

Sprint Review

In the sprint review, we check the intermediate result of the processed user stories. We check whether all the stories that should be completed meet the DOD. Furthermore, we discuss in the team what went well, what problems we encountered and how we solved them. Based on these results, the product increment is created. In our case, the product owner, who represents the customer, tests the product increment against the requirements. The review is done from the customer's perspective by testing the product increment. The outcomes are utilized to update the product backlog. The sprint review meeting takes place on the last day of the sprint.

Sprint Retrospective

In the sprint retrospective, we gather information as a team about what went well and what didn't go as planned in the previous sprint. We then derive specific improvements and plan their implementation. Our goal is to improve the efficiency, quality, communication and speed within our team. To achieve this, we give ourselves constructive criticism and are open to feedback. The sprint retrospective takes place on the last day of the sprint.

3.2.7. Review Project Setup

This chapter reflects on the initial setup and organization of our project.

Initial Situation

Our project began with an in-depth analysis of the provided functional and non-functional requirements. We used the initial project description as a guiding framework and scheduled a meeting with our coach, Dr. Simon Kramer, to clarify any ambiguities. This analysis led to several key insights:

- ▶ The application must be platform-independent and developed in Java.
- ▶ The code should be minimalistic, modular, and self-explanatory.
- ▶ A public repository is required for code availability.
- ▶ We depend on external resources like Archive Today and Apache PDFBox.
- ▶ A console application is preferred, eliminating the need for a graphical interface.

This systematic approach allowed us to create a modular and extendable codebase. When the integration of the Wayback Machine was later introduced as a requirement by our stakeholder, our well-structured codebase enabled us to agilely accommodate this new feature. This addition was assessed to ensure it added value to the project and didn't become a blocker, demonstrating the effectiveness of our initial analysis and flexible project approach.

Subject analysis

In our subject analysis, we thoroughly researched necessary external resources, focusing on Archive Today and the Wayback Machine. Our investigation included their automation interfaces, REST API support, and potential challenges such as Captcha and login requirements. We also ensured these resources were compatible with our FLOSS licensing requirement. This study allowed us to design an initial user interaction model and provided a clear roadmap for our application's development, facilitating a smooth transition into the coding phase.

Stakeholder(-Management)

Our stakeholder management focused on effective communication with our primary stakeholder, Dr. Simon Kramer. We efficiently addressed any uncertainties by scheduling regular meetings with Dr. Kramer, utilizing our designated Project 1 timeslot. This approach allowed for straightforward and effective resolution of questions and alignment of project deliverables.

Organisation

Our team's organization was based on agile principles, with scheduled Monday and Friday meetings serving as our Scrum events. These structured sessions were complemented by self-organized meetings among team members as needed, facilitating code reviews, quick resolution of questions, or addressing any issues. This agile approach ensured effective communication and collaboration, which was vital for the project's smooth progression.

Installations

For our project installations, we established a GitHub repository and a Jira project on Atlassian.net, ensuring efficient collaboration. Our first week involved defining objectives and creating epics in Jira. We also updated our IDEs and installed the latest Java version for consistency. Reflecting back, an early setup of a basic code structure would have been beneficial, potentially minimizing the need for later refactoring.

3.2.8. Review

This chapter provides a review of our project, discussing the management of our sprints and backlogs, and our progress towards achieving our product goal.

Product goal

Achieving the product goal for the URL Archiver application was a journey of iterative development and unwavering commitment to our goals. We successfully developed a platform-independent, CLI-based Java application that met all the specified functionalities, including efficient URL extraction, archiving options and CSV file generation. The application adheres to FLOSS licensing, ensuring open-source availability. Accompanying this, we created detailed user manuals, installation instructions and thorough software documentation. Our code is characterised by minimalism and modularity, with self-explanatory sections that underline our commitment to quality and maintainability.

Sprint goals

Throughout the sprints, we've rigorously set our goals in alignment with the SMART criteria, ensuring clarity and measurable targets. From the foundational work of implementing a robust input handler in Sprint 1 to enhancing system stability in Sprint 4, each goal was carefully crafted and successfully achieved. Sprint 5's focus on asynchronous archiving refined our application's functionality, while Sprint 6's emphasis on design and documentation prepared us for the final review. This structured approach to goal-setting has been instrumental in our consistent delivery of sprint objectives.

Delimitation

The delimitation in our project was implicitly conducted through disciplined backlog management and sprint planning. We confined our efforts to the most critical features, implicitly setting boundaries that guided our development focus. This implicit delimitation allowed us to concentrate resources on the highest priority tasks, ensuring the project's scope remained tightly aligned with our key objectives. While a formal delimitation review was not performed, our adherence to Agile principles effectively served the same purpose, keeping the project streamlined and on target.

Delivery objects

We successfully developed a platform-independent Java application that follows the principles of FLOSS. The application includes a command-line interface for ease of use, can process various file inputs, has URL extraction and archiving functionalities, and generates a CSV file documenting the archived URLs. In addition, we provided user manuals,

installation scripts, and software documentation to improve user experience and understanding. Our attention to detail and commitment to our goals throughout the project's lifecycle resulted in successfully meeting our delivery objectives.

Product Backlog

Our product backlog management was an exercise in adaptability and strategic foresight. In collaboration with our stakeholder, we regularly assessed and updated our backlog, ensuring it remained relevant and reflective of current project needs. This iterative process, led by the Product Owner, involved introducing new user stories and de-prioritizing or removing those that had lost their relevance. Our focus was on high-priority tasks that directly contributed to achieving the product goal, while medium and low-priority items were cataloged for potential future development. This approach maintained a balance between addressing immediate project requirements and envisioning future enhancements. The careful maintenance of the backlog provided a clear roadmap for our efforts, crucial for the effectiveness of our sprint reviews and overall project progress.

Sprint Backlogs

Throughout the sprints, our team consistently displayed impressive skills in managing and successfully completing the sprint backlogs. We consistently improved our workload management with each sprint, adapting our strategies based on previous experiences. There was only one instance of user stories being carried over to the next sprint, reflecting our overall planning efficiency and responsiveness to unforeseen challenges. The velocity report illustrates this progress, with an upward trend in completed story points, underlining our growing proficiency in sprint management.



Figure 3.44.: Velocity Report

3.2.9. Retrospective I: Scrum Method

This chapter reflects on the distribution and implementation of Scrum roles as previously described. It also considers our insights from Scrum Events and Scrum Artifacts.

Scrum Roles

Nicolin fulfilled the role of Product Owner exceptionally well. He consistently maintained and prioritized the Product Backlog. Moreover, he ensured we were in close communication with our customer, Dr. Kramer, guaranteeing that stakeholder requirements were discussed and incorporated. This was instrumental in the business success of the product being developed.

As Scrum Master, Abidin carried out his responsibilities with diligence, as evidenced by the team's efficient Scrum operation. Acting as a liaison between the Product Owner and the Developers, any questions regarding the agile methodology were promptly clarified. Additionally, obstacles beyond the team's capacity were resolved with the involvement of our coaches.

Kilian held the role of developer and played a key role in shaping the product. As there were only three of us in the project team, Abidin and Nicolin also contributed to the technical development.

Our goal to satisfy our client, Dr. Simon Kramer, was achieved through the early and continuous delivery of our software. This ensured the software developed matched Dr. Kramer's vision and requirements. Moreover, we were able to implement new requirements during the project, as clearly demonstrated by integrating the Wayback Machine. Throughout the project, we maintained self-organization, a consistent pace, reflection, and expeditious work.

Scrum Events

Sprint Our two-week sprint cycles have been effective in balancing workload and facilitating regular feedback. Setting SMART goals for each sprint provided us with a clear and focused direction and contributed significantly to our team's progress. Although we occasionally faced challenges in estimating workload, these instances offered valuable lessons in capacity planning. Exploring more precise estimation techniques could improve our ability to align tasks with our team's capacity, building on our already solid sprint planning process.

Sprint Planning Sprint planning was a strong point for our team, particularly with the upfront estimation of user stories and prioritization based on business value and team velocity. We ensured that the most impactful tasks were addressed each sprint. Encountering larger-than-expected user stories at times provided learning opportunities for better

task evaluation. A more comprehensive review process for estimating user stories could strengthen our sprint planning further.

Daily Scrum Adjusting to bi-weekly meetings complemented our part-time work schedules and kept our team communication robust. Using Microsoft Teams for meetings added flexibility and convenience. While the format sometimes delayed the addressing of immediate issues, it generally promoted focused and efficient discussions. A mid-week check-in could enhance our responsiveness without burdening the team's schedule.

Sprint Review Our sprint reviews were invaluable, offering objective evaluations that kept product increments in line with customer needs and expectations. Responding to customer feedback, though sometimes challenging, was a dynamic learning opportunity that led to significant product enhancements. More frequent reviews during the sprint could make the feedback integration process even smoother.

Sprint Retrospective The sprint retrospectives were characterized by constructive, open discussions that effectively pinpointed our successes and growth areas. These sessions were instrumental in fostering a culture of continuous improvement, with actionable steps consistently identified to enhance team processes. While implementing these improvements took effort, the retrospectives were central to our team's development and unity. A structured follow-up mechanism for retrospective suggestions could maximize these sessions' impact.

Although we attempted to provide each other with honest and constructive feedback, we realised that it can be challenging to do so. This is often due to fear of the recipient's reaction. However, we believe that with time, giving feedback will become easier for a team that has worked together for an extended period.

Scrum Artifacts

Product Backlog The Product Backlog was a dynamic, well-organized tool that effectively captured our project's requirements and priorities. Regular updates and refinements ensured it remained in sync with our evolving project goals and customer needs. Clear categorization and prioritization of items significantly aided our planning and decision-making. We are considering incorporating more frequent stakeholder feedback sessions to ensure that the backlog continues to reflect the most current and relevant project needs.

Sprint Backlog Our Sprint Backlog was used effectively to map out user stories, providing a clear outline of the sprint goals. Generally, we managed without creating detailed tasks for each story, favoring a broader view for simplicity and clarity. Introducing specific tasks

for user stories in future sprints could improve granularity and task management. This approach could enhance tracking and aid in early obstacle identification.

Product Increment Our approach to product increments has highlighted our team's dedication to delivering high-quality results at the end of each sprint. The increments met the set criteria and client expectations, showcasing our ability to effectively transform backlog items into tangible, valuable product features. Regular reviews and refinements ensured alignment with project goals and customer needs.

3.2.10. Retrospektive II: Tools/Instruments

Insights tool use

During the project, we used various tools that were essential to our success. In this section, we will discuss these tools and our experiences with them.

Java Java, our chosen primary programming language, facilitated a smooth development process due to the team's pre-existing knowledge. This familiarity allowed us to focus on delivering a high-quality application efficiently.

Maven Maven's role in our project was critical, handling compilation, dependency management, and the building process. Our prior experience with Maven enabled its seamless integration into our workflow, contributing to a smooth development cycle.

JUnit 5 JUnit 5 played a pivotal role in our testing framework, guaranteeing the stability and reliability of our application, especially following code refactoring activities, thus ensuring the integrity of our application's logic.

JetBrains IntelliJ IDEA The selection of IntelliJ IDEA as our integrated development environment provided an array of tools that streamlined our development processes, enhancing productivity, and facilitating a high level of code quality.

GitHub GitHub was the backbone of our version control system, offering a robust platform for collaborative code management. It enabled efficient teamwork and was essential in maintaining a consistent codebase.

Atlassian Jira Atlassian Jira was a cornerstone of our project management, enabling us to track progress, manage priorities, and streamline our workflow effectively, proving itself as an invaluable asset to our project organization.

LaTeX LaTeX was used for creating our project documentation and presentations. Despite initial challenges due to varying levels of experience within the team, it ultimately enabled us to produce professional and consistent documentation.

Communication

Microsoft Teams With Microsoft Teams, we were able to streamline our communication, enabling an effective platform for meetings and Scrum events that fostered rapid problem resolution and team agility. The platform's flexibility and immediacy were crucial to our effective problem-solving process.

Email and BigBlueButton For external communications, particularly with coaches, we used email for asynchronous exchanges and BigBlueButton for real-time video conferencing, ensuring clear and direct dialogue with our stakeholders.

Controlling

Atlassian Jira Jira served as the central control tool for our project, providing a comprehensive overview of project status and task distribution. It facilitated clear categorisation and prioritisation of tasks, thereby enhancing our workflow and project management.

Burndown Charts We utilized Burndown Charts to monitor our sprint efficiency and to inform future sprint planning. This allowed for a balanced workload management and realistic goal-setting based on our team's velocity and past performance.

3.2.11. Lessons learned

Team

Our team gained valuable experience from Project 1, learning to balance part-time study with software development. Implementing the Scrum methodology was challenging due to our limited time, but it highlighted the importance of adaptability and efficient planning. We tailored our approach, focusing on effective communication and workload distribution.

In summary, Project 1 was a valuable experience that taught us the importance of balancing teamwork and individual responsibilities. It has prepared us for future challenges in software development.

Nicolin Dora

Developing and managing my first software project from scratch was a great experience. However, working with SCRUM was challenging as we could only work on weekends or evenings due to part-time study. This made it difficult to use SCRUM effectively, and planning and report preparation took almost as much time as product development. In principle, SCRUM is a useful method as it allows for easy tracking of project progress. However, the students could have been given more freedom in this project to simplify the method.

I enjoyed the teamwork and appreciated the goal-oriented approach. It was interesting to balance the different demands on quality. Dr. Kramer and Mr. Helbling were always available to answer questions, which made collaboration easier. For future projects, I plan to use a slimmed-down version of SCRUM, depending on the complexity of the project. I am determined to improve my planning and communication within the team right from the start.

Abidin Vejseli

Project 1 has been an enriching yet demanding journey, requiring meticulous time management and organization from my side. The Scrum methodology, which I believe our team adapted effectively, and the complexities of part-time study contributed to this challenge. Personally, I felt the initial module introduction was rather lengthy and could be condensed to the essential material. It was unfortunate not to have both coaches at the presentation, which was an aspect that I missed.

Working within the team was a positive experience, and I'm proud of how we distributed the workload, ensuring valuable learning for everyone. The interaction with our coaches was a highlight, their readiness to assist was something I greatly appreciated. For the next project, I aim to start and ideally complete my user stories during the first half of the sprint.

Killian Wampfler

For me, Project 1 was a great refresher for my programming skills and software development techniques. As I already use the Scrum method in my company, this part wasn't new to me. However, it was interesting to use the Scrum method in such a small team (3 people). Also, the team worked really well together and I think we could all learn something from this project.

Personally, I learned two things from this project. First, my unit tests have a lot of room for improvement (in quality and quantity). And more importantly, my git commit messages weren't the norm. Because of this, our git repository doesn't look consistent, and as this project needs to be publicly accessible, this is not desired.

Lastly I really appreciated how much freedom Dr. Kramer gave us and that he always took the time to answer our questions.

4. Deployment/Integration

4.1. Licensing and Compliance

4.1.1. Licensing Overview

The project, along with its original source code, is licensed under the permissive and open MIT License. This license aligns with the project's goal of accessibility and ease of use, allowing for free use, modification, distribution, and private use of the software.

4.1.2. Compliance with Open Source Licenses

Although the project itself is licensed under the MIT License, it uses several open-source libraries that are subject to their respective licenses. Please refer to Appendix A.2 for more information. It is worth noting that this project employs libraries licensed under the Eclipse Public License v2.0 (EPL-2.0), including JUnit Jupiter API.

Using EPL-2.0 licensed libraries in a MIT-licensed project is compliant with open source licensing standards, as long as certain conditions are met.

- 1. Attribution and Notices:** The project includes a 'Licenses and Attributions' section in the README file. This section acknowledges the use of open-source libraries, specifying their licenses and providing due credits. Additionally, a 'NOTICES.txt' file is included in the project's resources, detailing the open-source components used, their licenses, and where to find the full license texts. This approach is a crucial step to respect and recognize the work of open-source contributors and to maintain transparency about the software's composition.
- 2. Separation of Licenses:** It is important to clarify that the MIT License applies to the original code developed for this project. In contrast, the libraries used retain their original licenses (EPL-2.0 in the case of JUnit).
- 3. No Modification of EPL-2.0 Libraries:** This project uses the EPL-2.0 licensed libraries in their unmodified form. Any modification to such libraries would require adherence to the specific terms and conditions of the EPL-2.0.

4.1.3. Purpose of Compliance

Ensuring compliance with the licensing terms of used libraries is not only a legal requirement but also a commitment to the open-source community's ethical standards. It ensures that the project respects the rights and efforts of other developers and contributes to the sustainable and responsible use of open-source software.

4.2. Installation (Sysadmin) Manual & Script

The URL-Archiver enables the extraction of URLs from any Unicode text or PDF file and allows for interactive archiving on one of the supported archiving services.



The application was designed to be platform-independent . However, it has only been tested on the following systems, so it cannot be guaranteed to work without restrictions on other platforms.

- ▶ Windows 11 (Version 23H2)
- ▶ Windows 10 (Version 22H2)
- ▶ macOS (Ventura)
- ▶ Ubuntu (20.04.3 LTS)

4.2.1. Requirements

To build and start the application, ensure that the following dependencies are installed on your system:

- ▶ Git: Latest stable version recommended.
- ▶ Maven: Version 3.8 or higher.
- ▶ Java: Version 21.

4.2.2. Clone the repository

To clone the repository, run the following command in a terminal:

```
git clone https://github.com/devobern/URL-Archiver.git
```

4.2.3. Build and run scripts

The build and run scripts are provided for Windows (`build.ps1`, `run.ps1`, `build_and_run.ps1`), Linux, and MacOS (`build.sh`, `run.sh`, `build_and_run.sh`). The scripts are located in the root directory of the project.

The scripts need to be executable. To make them executable, run the following command in a terminal:

- Linux / MacOS: `chmod +x build.sh run.sh build_and_run.sh`
- Windows:
 - Open PowerShell as an Administrator.
 - Check the current execution policy by running: `Get-ExecutionPolicy`.
 - If the policy is Restricted, change it to RemoteSigned to allow local scripts to run. Execute: `Set-ExecutionPolicy RemoteSigned`.
 - Confirm the change when prompted.
 - This change allows you to run PowerShell scripts that are written on your local machine. **Be sure to only run scripts from trusted sources.**



Windows

Build the application

To build the application, open a command prompt and run the following script:

```
./build.ps1
```

Run the application

To run the application, open a command prompt and run the following script:

```
./run.ps1
```

Build and run the application

To build and run the application, open a command prompt and run the following script:

```
./build_and_run.ps1
```

Linux and macOS

Build the application

To build the application, open a command prompt and run the following script:

```
./build.sh
```

Run the application

To run the application, open a command prompt and run the following script:

```
./run.sh
```

Build and run the application

To build and run the application, open a command prompt and run the following script:

```
./build_and_run.sh
```

4.3. User Manual



To follow the instructions in this section, the application must be built. See 4.2.

The URL-Archiver is a user-friendly application designed for extracting and archiving URLs from text and PDF files. Its intuitive interface requires minimal user input and ensures efficient management of URLs.

4.3.1. Getting Started

Windows

Open Command Prompt, navigate to the application's directory, and execute:

```
./run.ps1
```

Linux / MacOS

Open Terminal, navigate to the application's directory, and run:

```
./run.sh
```

4.3.2. Operating Instructions

Upon launch, provide a path to a text or PDF file, or a directory containing such files. The application will process and display URLs sequentially.

Navigation

Use the following keys to navigate through the application:

- ▶ **o:** Open the current URL in the default web browser.
- ▶ **a:** Access the Archive Menu to archive the URL.
- ▶ **s:** Show a list of previously archived URLs.
- ▶ **u:** Update and view pending archive jobs.
- ▶ **n:** Navigate to the next URL.

- ▶ **q:** Quit the application.
- ▶ **c:** Change application settings.
- ▶ **h:** Access the Help Menu for assistance.

Archiving URLs

Choose between archiving to Wayback Machine, Archive.today, both services, or canceling.

When opting to use Archive.today for archiving, an automated browser session will initiate, requiring you to complete a captcha. Once resolved, the URL is archived, and the corresponding archived version is then collected and stored within the application.

Configuration

Customize Access/Secret Keys and the default browser. Current settings are shown with default values in brackets.

Exiting

To exit, press **q**. If a BibTex file was provided, you'll be prompted to save the archived URLs in the BibTex file. Otherwise, or after saving the URLs in the BibTex file, you'll be prompted to save the archived URLs in a CSV file.

For BibTex entries:

- ▶ Without an existing note field, URLs are added as: `note = {Archived Versions: \url{url1}, \url{url2}}`
- ▶ With a note field, they're appended as: `note = {<current note>, Archived Versions: \url{url1}, \url{url2}}`

5. Conclusion

5.1. Discussion

In the course of our academic endeavor, the URL-Archiver project, we have successfully designed and implemented a Java-based application capable of extracting and archiving URLs from Unicode text and PDF documents. This development demonstrates the practical application of our academic learning in real-world scenarios. The application is intended to assist professionals in research and journalism by providing a streamlined approach to archive URLs and organize them.

Our project's outcome highlights the application's proficiency in fulfilling all requirements and reaching the goal to deliver a FLOSS-licensed, platform-independent Java-program called URL-Archiver. As young professionals in the IT field, our efforts not only gave us valuable hands-on experience, but also contributed to the larger dialogue regarding digital data management.

However, like any academic project, our application is not without its limitations. Currently, the tool has a limited compatibility with certain file formats, which may limit its applicability. Moreover, the necessity for manual captcha solving, while ensuring security, does pose an inconvenience and limits the tool's automation capabilities.

Looking ahead, we have a number of enhancements in mind for the URL-Archiver. Enhancing the application's compatibility with a broader range of file formats would increase its utility. Furthermore, integrating support for additional archiving services, particularly those with accessible APIs, stands as a significant upgrade. This would not only automate interactions with a wider range of archiving solutions but also improve both the utility and user experience of the application.

In conclusion, the URL-Archiver effectively serves its intended purpose. However, there is still room for optimisation. The findings and challenges identified during this project lay the groundwork for future improvements. We believe that with continued research and development, the URL-Archiver can become an even more versatile and valuable tool.

5.2. Bottom Line

5.2.1. Conclusion

The journey of developing the URL-Archiver has been a mix of unexpected challenges and rewarding experiences. Initially assumed to be straightforward, the project revealed its complexity as we delved deeper. This experience highlighted the importance of effective team communication, a critical factor in overcoming the unforeseen challenges of the project. Reflecting on our approach, establishing a more robust project structure from the beginning would have been beneficial. Despite these challenges, the project was rewarding and a great learning experience.

One of the key takeaways was the value of practical experience with Git. This tool not only enhanced our collaborative efforts but also contributed to our personal skill development. Remarkably, the URL-Archiver stands as our first project built entirely from scratch, marking a significant milestone in our journey as developers.

We firmly believe that the URL-Archiver has the potential to be a valuable asset to its users. By significantly reducing the time and effort required for extracting and archiving URLs, the application promises an increase in productivity for its users. This improvement is not just theoretical; it's a practical solution addressing a real need in the digital world.

In conclusion, while the project journey had its ups and downs, the collective learning and the potential impact of the URL-Archiver make it a fulfilling experience. We are optimistic about its utility and look forward to seeing its adoption and evolution in the professional world.

5.2.2. Productivity Increase Calculation

This section calculates the estimated productivity increase for users of our URL-Archiver software. The calculation factors include task duration without and with the software, the resulting time economy, cost per time unit, and overall cost economy.

Assumptions

The following assumptions were made for our calculations:

- ▶ The file contains 60 URLs.
- ▶ Without URL-Archiver, manually finding and copying each URL takes approximately 12 seconds.
- ▶ The URL-Archiver takes an average of 2 seconds to find all URLs in a file.
- ▶ Archiving one URL using Wayback Machine manually takes 60 seconds.
- ▶ Archiving one URL using Archive Today manually takes 90 seconds.

- ▶ With URL-Archiver, initiating the archiving process for one URL in Wayback Machine takes 10 seconds. The actual archiving in the background takes in average 30 seconds.
- ▶ Archiving one URL using Archive Today remains at 90 seconds even with the URL-Archiver.
- ▶ The average hourly wage for a professional in Switzerland is CHF 40.
- ▶ Parallel archiving with Wayback Machine allows for multiple URLs to be archived concurrently, reducing the total time significantly.

Wayback Machine Analysis

Without URL-Archiver:

$$\begin{aligned}\text{Task duration without SW} &= \text{URL extraction} + \text{URL archiving} \\ &= (60 \times 12 \text{ seconds}) + (60 \times 60 \text{ seconds}) \\ &= 720 \text{ seconds} + 3600 \text{ seconds} \\ &= 4320 \text{ seconds (72 minutes)}\end{aligned}$$

With URL-Archiver:

$$\begin{aligned}\text{Task duration with SW} &= \text{URL extraction} + \text{URL archiving} \\ &= 2 \text{ seconds} + (60 \times 40 \text{ seconds}) \\ &= 2 \text{ seconds} + 2400 \text{ seconds} \\ &= 2402 \text{ seconds (40 minutes)}\end{aligned}$$

Productivity Increase:

$$\begin{aligned}\text{Task duration economy} &= 72 - 40 = 32 \text{ minutes} \\ \text{Cost per time unit} &= \frac{40 \text{ CHF}}{60 \text{ minutes}} \\ \text{Cost economy} &= 32 \times \left(\frac{40 \text{ CHF}}{60} \right) = 21.33 \text{ CHF per document}\end{aligned}$$

Archive Today Analysis

Without URL-Archiver:

$$\begin{aligned}\text{Task duration without SW} &= \text{URL extraction} + \text{URL archiving} \\ &= (60 \times 12 \text{ seconds}) + (60 \times 90 \text{ seconds}) \\ &= 720 \text{ seconds} + 5400 \text{ seconds} \\ &= 6120 \text{ seconds (102 minutes)}\end{aligned}$$

With URL-Archiver:

$$\begin{aligned}\text{Task duration with SW} &= \text{URL extraction} + \text{URL archiving} \\ &= 2 \text{ seconds} + (60 \times 90 \text{ seconds}) \\ &= 2 \text{ seconds} + 5400 \text{ seconds} \\ &= 5402 \text{ seconds (90 minutes)}\end{aligned}$$

Productivity Increase:

$$\begin{aligned}\text{Task duration economy} &= 102 - 90 = 12 \text{ minutes} \\ \text{Cost per time unit} &= \frac{40 \text{ CHF}}{60 \text{ minutes}} \\ \text{Cost economy} &= 12 \times \left(\frac{40 \text{ CHF}}{60} \right) = 8.00 \text{ CHF per document}\end{aligned}$$

5.3. Future Work

The future vision for the URL-Archiver project includes several enhancements to extend its capabilities and improve the user experience. These potential additions, formulated within our project's limited timeframe, include:

- ▶ Adding more archiving services for a broader reach.
- ▶ Expanding support for various input file types.
- ▶ Implementing a user-friendly graphical interface.
- ▶ Enabling multilingual support for global accessibility.
- ▶ Automatically archiving all URLs in a file for efficiency.
- ▶ Providing more detailed setting options for user customization.
- ▶ Publishing the application in package repositories to simplify installation.
- ▶ Improving the code layout, like breaking up the controller for better clarity.

Integrating services such as Memento Time Travel would expand our archiving capabilities and provide a wider historical perspective. Adapting the application to handle various file formats like .docx could increase its relevance in various other settings. The introduction of a graphical user interface is another key area, which would dramatically improve user interaction, making the tool more accessible and engaging, particularly for those less versed in CLI environments. Adding multilingual support is also crucial, as it would break down language barriers, enhancing the tool's usability. The proposed enhancements, such as automatic URL archiving, detailed settings, application publishing in repositories, and code layout improvements like refactoring the controller, collectively aim to not only improve the tool's functionality and user experience but also to ensure better maintainability and ease of future extensions.

Glossary

Archive Today **A web archiving service that stores snapshots of web pages** for preservation and retrieval. 2

BibTeX **Program for the creation of bibliographical references** and directories in TeXor L^AT_EX documents 2

CLI **A Command Line Interface (CLI)** is a type of user interface navigated entirely using text commands. It allows users to interact with software or operating systems by typing commands into a console or terminal. 2

CSV **CSV (Comma-Separated Values)** is a file format used to store tabular data, such as a spreadsheet or database. Data fields in CSV files are separated using commas. 2

Factory Pattern **A design pattern in software development** used to create objects, allowing interfaces to define object creation but letting subclasses alter the type of objects that will be created. 4

FLOSS **Free/Libre and Open Source Software (FLOSS)** refers to software that is both free in the sense of freedom (libre) and open source. It allows users the freedom to run, study, modify, and distribute the software for any purpose. 2

GUI **GUI, or Graphical User Interface**, refers to a type of user interface that allows users to interact with electronic devices through graphical icons and visual indicators, as opposed to text-based interfaces, typed command labels, or text navigation. It simplifies the user experience by employing graphical representations of the commands and actions available. 4

Model-View-Controller (MVC) pattern **Model-View-Controller**, a **software design pattern** for implementing user interfaces, data, and controlling logic. 4

Unicode Text File **A Unicode Text File is a file that uses the Unicode standard** for encoding its content. Unicode enables the representation of a vast array of characters from various scripts and symbol sets, making it suitable for internationalization and localization. 2

URL A **URL (Uniform Resource Locator)** is an internet address that directs to a specific resource, like a webpage. 2

URL-Archiver **A software tool designed to extract, archive, and manage URLs found within digital documents.** It supports various file formats and integrates with services like the Wayback Machine for archiving web pages. 2

Wayback Machine **A digital archive of the World Wide Web**, allowing users to see older versions of web pages. 2

Index

Factory Pattern, 4, 5
internet connection, 4
MVC-Pattern, 4–6
platform-independent, 4, 77
priorities, 3
product goal, 2

Bibliography

A. Appendix

A.1. Original Project Description

URL-Archiver

Description

The goal of this project is to deliver a FLOSS-licensed, platform-independent Java-program (called "URL-Archiver") that

- (1) takes as input (the path of) a directory or any Unicode-text- (e.g.: .BIB, .TEX; .HTML; etc.) or .PDF-file (<https://www.baeldung.com/java-curl>);
- (2) scans it for any URLs (<https://stackoverflow.com/questions/4026614/extract-text-from-pdf-files> , <https://librepdf.github.io/OpenPDF> , <https://pdfbox.apache.org> ; see also https://en.wikipedia.org/wiki/List_of_PDF_software);
- (3) extracts all URLs (regular expression ;-) from the text;
- (4) optionally spring-loads all URLs in a Web-browser;
- (5) posts all URLs to <https://archive.ph> ;
- (6) gets the resulting archived URLs;
- (7) outputs a CSV-file of the resulting key-value (URL, archived URL) pairs; and
- (8) optionally inserts the archived URLs into a .BIB-file.

The program code should be minimal, modular, and self-explaining.

The project report should be concise (maximally informative, minimally long).

It must contain this project description as a quotation.

Technologies

Java, LaTeX

Advisor

Dr. Simon Kramer

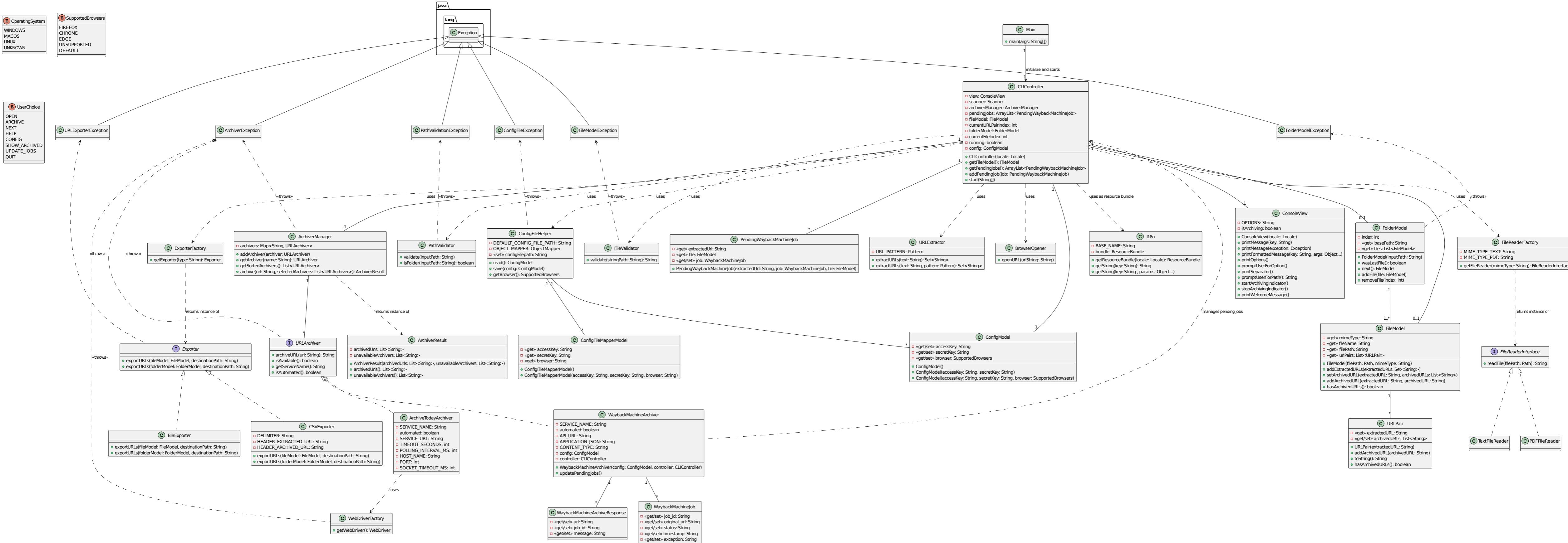
A.2. List of Used Libraries and Their Licenses

Below is the list of libraries used in the project, along with a short description and their versions.

| Library | Version | Short Description | Used License |
|------------------------|---------|--|-----------------------------|
| JUnit Jupiter API | 5.9.2 | Unit testing framework for Java applications. | Eclipse Public License v2.0 |
| JUnit Jupiter Engine | 5.9.2 | The test engine for running JUnit tests. | Eclipse Public License v2.0 |
| Selenium Java | 4.15.0 | Automation framework for web applications testing. | Apache-2.0 |
| Selenium Logger | 2.3.0 | A wrapper for enhanced Selenium log management. | MIT |
| Mockito Core | 5.4.0 | Mocking framework for unit tests in Java. | MIT |
| Mockito JUnit Jupiter | 5.4.0 | Integration of Mockito with JUnit Jupiter. | MIT |
| System Lambda | 1.2.1 | Utilities for testing Java code that uses system properties and environment variables. | MIT |
| Apache PDFBox | 3.0.0 | Library for creating and manipulating PDF documents. | Apache-2.0 |
| Jackson Core | 2.16.0 | Core part of Jackson that defines common low-level features. | Apache-2.0 |
| Jackson Dataformat XML | 2.15.2 | Support for reading and writing XML encoded data via Jackson abstractions. | Apache-2.0 |

Table A.1.: List of Used Libraries in the Project

A.3. Class Diagram



Declaration of Authorship

I hereby declare that I have written this thesis independently and have not used any sources or aids other than those acknowledged.

All statements taken from other writings, either literally or in essence, have been marked as such.

I hereby agree that the present work may be reviewed in electronic form using appropriate software.

January 14, 2024

N. Dora

N. Dora

C. Example

A. Vejseli

C. Example

K. Wampfler