

## URL-Archiver

### Final Presentation

January 6, 2024

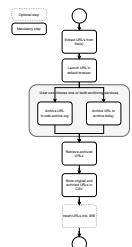
N. Dora, A. Vejseli, K. Wampfler | School of Engineering and Computer Science

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

## Table of Content

Topics Covered

- ▶ Problem Statement Recap
- ▶ Solution Overview
- ▶ Live-Demo
- ▶ Conclusion
- ▶ Questions?
  - Project Goal
  - Requirements
- ▶ Solving The Problem
  - Process Model
  - Architecture
  - Data Model
  - Technologies
- ▶ Project Management with SCRUM
  - Scrum Team
  - Our Scrum Adoptions
  - Scrum implementation**
- ▶ Questions?



2

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

## Template

Template with columns



3

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

## Problem Statement Recap

## Problem Statement Recap "URL-Archiver"

The ever changing internet

Kilian:

- Rückblick Project Setup: Erfassung der Ausgangssituation, Themen-Analyse, Stakeholder(-Management), Organisation, Installationen, etc.
- Eine kurze Répetition der Problemstellung und des Lösungswegs?

About what was this project again?



## Problem Statement Recap "URL-Archiver"

The ever changing Internet

### ■ problem statement

- ever-changing Internet (websites are taken down or changed)
- Forensic (e.g. legal report about a statement on the Internet)



7

## Project Goal

### ■ Develop a software that...

- ...searches hyperlinks within documents
- ...is capable of archiving websites
- ...can store current and archived links in a file
- ...is easy and intuitive to use



8

## Project Setup Review

Initial Situation

### ■ What we knew:

- the end result has to be a platform independent java application
- we can only use FLOSS-licensed libraries
- the software must support archive.ph and WayBackMachine
- the code must be publicly accessible



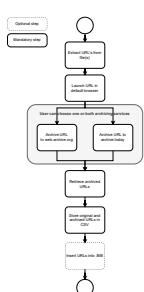
9

## Project Setup Review

### Initial Situation

#### ■ What we did first:

- project setup
- we researched both archiving services and what's possible
- we designed the rough process for user interaction
- we decided on a software design pattern



## Project Setup Review

### Installation / Project Setup

- creation of a new **github repository** for public accessibility
- setting up a **jira project** (with epics and first stories)
- installing the **latest java version**



## Solution Overview

## Command line application

Why a command line and not a GUI application?

- **Focus on Logic:** Emphasis on functional depth over visual design.
- **Technical Audience:** Suited for users proficient in command-line environments.
- **Efficiency in Development:** Streamlines the development process by focusing less on UI design.
- **Flexibility and Scalability:** Provides a stable base for future feature expansion.

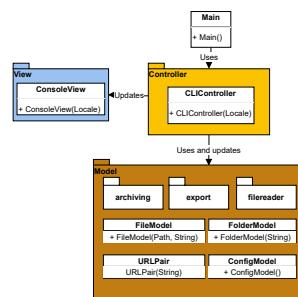


1. We have decided in favour of a command line program and against a GUI. For various reasons.
2. On the one hand, because we could focus on the logic and didn't have to build a responsive GUI.
3. On the other hand, our stakeholders are in the technical world and should therefore have few problems with a CLI application.
4. With a CLI application we create a stable basis to build on.
5. With the MVC pattern we can add a GUI in the future with little effort.

## Architecture

Employing the MVC Pattern

- Modular design for **easy extension**
- Separate data, view, and control flow
- Facilitates the potential addition of a GUI

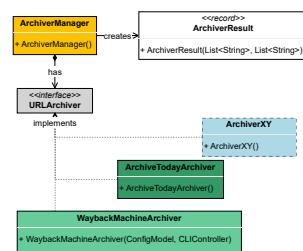


1. Speaking of the MVC pattern. As already mentioned, this pattern allows us to easily extend the application in the future.
2. As the view is separated from the model by the controller.
3. This means that these parts can be easily exchanged without having to make major changes to the others.

## Architecture

Factory Pattern

- **File Reader:** More input file types? No problem.
- **Archiving Services Integration:** An additional service can be added very easily.
- **Selenium Web Drivers:** Support another browser? Nothing easier than that!
- **Export functionality:** Easily add support for additional file types.



1. We have also used the factory pattern where appropriate.
2. Specifically with the archiving services, the Selenium web drivers, when exporting to a CVS and Bibtex file and when reading the files.
3. As you can see on the right, using the archiving services as an example, this allows us a simple extension.

## Enhancing Functionality and Flexibility

Configuration, Internationalization, and Testing

1. We have also made various decisions to improve the functionality and flexibility of the application.
2. We save configurations in a configuration file.
3. We have laid the foundation for multilingualism...
4. ... and are testing all relevant classes with unit tests; there are currently 107 unit tests in total.

- **Configuration File Usage:** Preserves settings between program sessions.
  - **Stored Settings:** Includes Wayback Machine API keys and default browser.
- **i18n Implementation:** Lays groundwork for multilingual support.
- **Unit Testing:** Comprehensive tests written for all relevant classes.



## Ease of Use: Installation Simplified

Custom Scripts for Cross-Platform Compatibility

1. To make it as easy as possible for users to get started with the application, we have created compilation and run scripts.
2. Of course bash scripts for Linux and macOS and PowerShell scripts for Windows.

- **Simplified Compilation and Installation:** Tailored scripts for Windows, Linux, and macOS.
- **User-Friendly Scripts:** Enable straightforward application compilation and execution.
- **Flexible Installation Options:** Users can choose to compile then run, or do both in one step.

```
#!/bin/bash

# Check for Java installation
if [ ! -e /usr/bin/java ]; then
    echo "Java is not installed or not in the PATH. Please install Java 11."
    exit 1
fi

# Check for Maven installation
if [ ! -e /usr/bin/mvn ]; then
    echo "Maven is not installed or not in the PATH. Please install Maven."
    exit 1
fi

# Build the project with Maven
mvn clean package
# Check for build success
if [ $? -ne 0 ]; then
    echo "Build failed."
    exit 1
fi

# Create 'bin' directory if it doesn't exist
mkdir -p bin

# Move the JAR file to the 'bin' directory
mv target/My-Application-1.0-Main.jar bin/
# If move fails
else
    echo "Build failed."
fi
```

## Live-Demo

## Conclusion

## Conclusion: Project Goal

- ✓ Develop a software that...
  - ✓ ...searches hyperlinks within documents
  - ✓ ...is capable of archiving websites
  - ✓ ...can store current and archived links in a file
  - ✓ ...is easy and intuitive to use



## Future Work

What can be done to improve the application further

- Add more **archiving services** (such as Memento Time Travel)
- Support more **input file types** (e.g. .docx)
- Integrate a **graphical interface**
- Implement **multilingual support**



## Lessons learned

Team

- **Basic code structure**
- **Code Reviews** more often
- **Weeklys** are important



## Lessons learned (personal)

Kilian Wampfler

- good refresher for java development
- the team worked well together
- my junit tests could be improved (consistency & quality)
- my git commit messages according to convention



## Lessons learned (personal)

Nicolin Dora



## Lessons learned (personal)

Abidin Vejseli



Questions?

## Problem Statement Recap "URL-Archiver"

The ever changing internet

- **ever-changing internet** (websites are taken down or changed)
- not a problem for everyday internet users but...
- ...what about **documents** (e. g. theses, documentations etc.)
  - invalid links to content-relevant information
  - invalid quote links
  - old documents may contain numerous broken hyperlinks
- another case: **Forensic**
  - e.g. **legal report** about a statement on the Internet



## Problem Statement

Simple Solution

- Archive the **actual state** of the linked websites
  - Wayback machine
  - archive.today
  - Memento Time Travel
  - and more...
- But who has the time and **motivation**?
  - ...to search each hyperlink
  - ...to manually archive each website



## Project Goal

- Develop a software that...
  - ...searches **hyperlinks** within documents
  - ...is capable of **archiving websites**
  - ...can store current and archived links in a file
  - ...is **easy** and **intuitive** to use



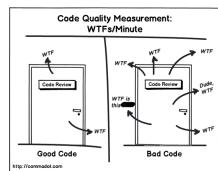
## Functional Requirements

- The software must be...
  - ...a Java application that is compatible with multiple platforms
  - ...Free and Open Source Software licensed (FLOSS)
  - ...capable of archiving websites on archive.ph or/and WayBackMachine
  - ...capable of generating a CSV-file with key-value (URL, archived URL)



## Non-Functional requirement

- Coding Norm
- publicly accessible
- Project-Method: SCRUM
- public documents in English
- project report with LaTeX

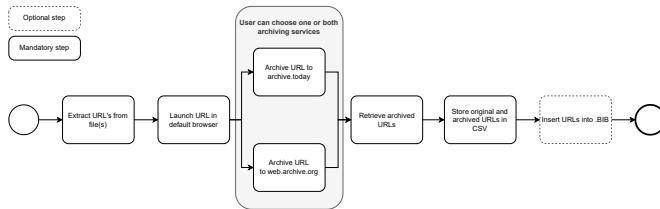


## Solving The Problem

## Process Model

Workflow for URL Extraction and Archiving

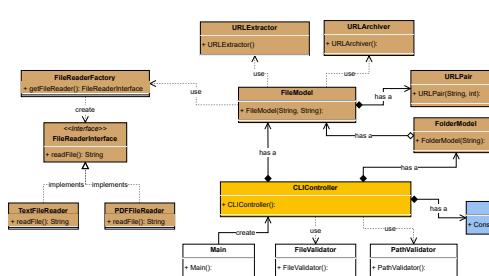
- User can skip URLs, launch them, access help or quit the program at any time



- Thank you for coming from my side too.
- First lets talk about how an user can archive urls.
- 
- For this you can see here a rough overview of how a user navigates through the application to archive URLs.
- The user can open the current URL in their default browser, access the help section, or view the archived URLs at any time.
- The user does not need to archive a url, he can skip the archiving process.
- Now to the more technical part.

## High Level Class Diagram

Overview

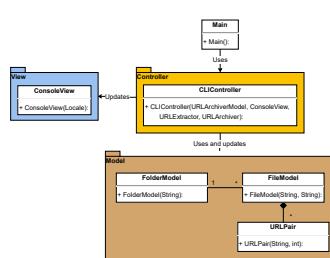


- To demonstrate how we meet the requirements for our application, which Kilian presented earlier, I will now provide a high-level overview of our structure and its individual components.
- It is important to note that this is our current structure, subject to expansion and adaptation as the project advances.
- The Main class serves as the application's entry point.
- CLIController interprets user commands, validates paths and types, and selects the data model.
- FileModel manages individual file data, and stores URL associations.
- FolderModel deals with file groupings within folders, integrating multiple FileModel instances.
- URLPair links original to the archived URL.
- I18n enables the application to support multiple languages, enhancing global usability.
- UserChoice defines possible user actions with corresponding I18n keys, aiding in international user interaction.
- ConsoleView acts as the communication hub, utilizing UserChoice for navigation and I18n for language support.
- Now I would like to discuss some aspects of the architecture in greater detail.

## Architecture

Employing the MVC Pattern

- Modular design for easy extension
- Separate data, view, and control flow
- Facilitates the potential addition of a GUI

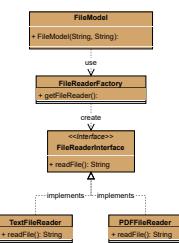


- As you may have observed, we utilize the MVC pattern, among other techniques, to ensure the scalability of the application.
- This involves separating the business logic from the presentation using a controller.
- This enables us to easily switch out the Command Line Interface with a Graphical User Interface.
- But that's not all!

## Architecture

Factory Pattern

- Scalable and robust design
- Ensures maintainability and clear code structure
- Commitment to best programming practices

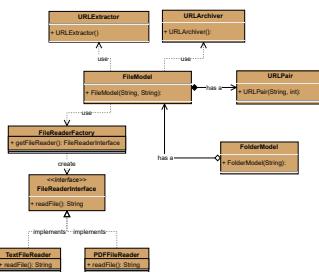


1. To enhance our application's flexibility further, we've implemented a factory class that reads and processes different file formats.
2. With a single call to the `getFileReader()` method from the `FileReaderFactory`, specifying the file type, we retrieve an appropriate `FileReader`.
3. Currently, our application can handle all UNICODE text files and PDFs.

## Key Components of the Data Model

Overview

- **FileModel:** Represents and handles individual files
- **FolderModel:** Manages a collection of `FileModel` instances
- **URLPair:** Links extracted URLs with their archived counterparts
- **FileReaderFactory:** Organizes and maintains `URLPairs`
- **URLExtractor:** Used for URL extraction
- **URLArchiver:** Used for URL archiving



Now, let us discuss our data model. It consists of the following classes:

1. The `FileModel` class manages the details of individual files. Each `FileModel` potentially containing multiple `URLPair` objects.
2. The `FolderModel` is applied when a user provides a folder path, storing the folder's structure and containing a list of `FileModel` objects.
3. The `FileReaderFactory` is our system for providing the right reader for a file, based on its type. This ensures that each file is read correctly and efficiently.
4. With the `URLExtractor`, our model extracts URLs from text, and the `URLArchiver` takes care of the URL archiving process.
5. This setup not only facilitates a clear division of tasks among the different components but also streamlines the process of archiving and retrieving URLs.
6. And now to the technologies we use for development.

## Technologies I/II

- **Java:** The primary programming language for our application
- **LaTeX:** Used for documentation and presentation
- **JUnit 5:** Utilized for unit testing
- **PDFTextStripper (PDFBox library):** Used for extracting text from PDF documents
- **Selenium:** Web automation tool used for tasks like inputting URLs, handling captchas, and retrieving archived URLs due to the absence of an official API from archive.today
- **Maven:** Essential for compilation, dependency management, and building the project



1. The application is developed in Java, while the report, presentations, and user manual are authored in LaTeX.
2. For testing, we utilize JUnit 5, and for text extraction from PDFs, we rely on PDFTextStripper.
3. The archiving process on Archive.today is executed using Selenium to navigate the lack of an API and to address captcha challenges.
4. Additionally, we use Maven for building the application and managing dependencies.

## Technologies II/II

- **Archive.today:** One of the services used by the URL-Archiver to archive URLs
- **Wayback Machine:** One of the services used by the URL-Archiver to archive URLs
- **JetBrains IntelliJ IDEA:** Used to develop the program and create the report and presentation in Latex
- **GitHub:** A platform for code hosting and collaboration, allowing us to manage our program's development and version control
- **Atlassian Jira:** A tool for tracking progress and managing our project goals and tasks



1. Furthermore, we rely on Archive.today and additionally on the Wayback Machine to archive URLs.
2. As our development environment, we use JetBrains' IntelliJ IDEA.
3. We use Github to store our code and for version control.
4. To manage and plan our project, we use Jira.

## Project Management with SCRUM

## Our Scrum Team



**Nicolin Dora**  
Product Owner & Developer



**Abidin Vejseli**  
Scrummaster & Developer



**Kilian Wampfler**  
Developer

## Other Roles



**Dr. Simon Kramer**  
Stakeholder & Customer



**Frank Helbling**  
PM-Advisor

## Our Scrum Adoptions

Sprint

Sprint	Sprint planning	Daily Scrum	Sprint Review	Sprint Retro
Two Weeks Sprint Goal (SMART)				

## Our Scrum Adoptions

Sprint Planning

Sprint	Sprint planning	Daily Scrum	Sprint Review	Sprint Retro
	First Monday Estimate User Stories Define Goal Choose User Stories			

## Our Scrum Adoptions

### Daily Scrum

Sprint	Sprint planning	Daily Scrum	Sprint Review	Sprint Retro
		<ul style="list-style-type: none"><li>• Daily</li><li>• Twice a week</li><li>• Current state</li><li>• MS Teams</li></ul>		

## Our Scrum Adoptions

### Sprint Review

Sprint	Sprint planning	Daily Scrum	Sprint Review	Sprint Retro
			<ul style="list-style-type: none"><li>• Last Day</li><li>• Completed?</li><li>• Problems?</li><li>• Testing</li><li>• Product increment</li></ul>	

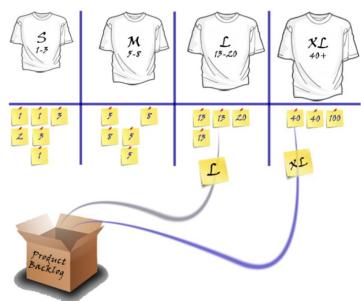
## Our Scrum Adoptions

### Sprint Retro

Sprint	Sprint planning	Daily Scrum	Sprint Review	Sprint Retro
				<ul style="list-style-type: none"><li>• Last Day</li><li>• Feedback</li><li>• Do's &amp; Don'ts</li></ul>

## Estimation Method

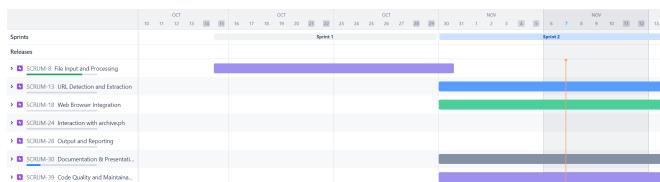
T-Shirt Sizes



Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

49

## Velocity



Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

50

## Backlog

SCRUM-0 File Input and Output or Context	SCRUM-1 URL Detection and Extraction	SCRUM-2 Web-Browser Integration	SCRUM-3 Interaction with antispam	SCRUM-4 Output and Reporting	SCRUM-5 Documentation & Presentat...	SCRUM-6 Code-Quality and Maintenance	SCRUM-7 User Decisions for Reporting	SCRUM-8 Automatic URL Submission	SCRUM-9 User Interaction for Capitalis...	SCRUM-10 Automatic Retrieval of Archived URLs	SCRUM-11 Generate CSV file	SCRUM-12 Integrate Archived URLs into Supported Sys...	SCRUM-13 User Manual	SCRUM-14 Intermediate presentation	SCRUM-15 Final presentation	SCRUM-16 User diagram	SCRUM-17 Installation manual & script	SCRUM-18 Requirements document	SCRUM-19 Implementation documentation
10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	
10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	
10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	
10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	10 days	

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

51

## Epics

The screenshot shows the Jira interface for the 'File Input and Processing' epic. On the left, there's a sidebar with a tree view of epics: 'Issues without epic' (selected), 'File Input and Processing', 'URL Detection and Extraction', 'Web Browser Integration', 'Interaction with archive.ph', 'Output and Reporting', 'Documentation & Presentation', and 'Code Quality and Maintainability'. The main area is titled 'File Input and Processing' and contains a 'Description' section with the text: 'As an user I want to input various file types via the command line so that they can be prepared for further processing.' Below this is a 'Child issues' table:

Issue Key	Name	Status	Progress
SCRM-10	Automatic File Type Detection	DONE	13%
SCRM-12	Processing Feedback	DONE	1%
SCRM-9	Prompt for File Path Input	NOT DONE	3%
SCRM-11	Processing of Directories	NOT DONE	3%
SCRM-48	Add user option for processing new path	TO DO	0%

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

52

## Sprint 2

### Sprint Goal

#### Sprint 2

The screenshot shows the Jira 'Sprint 2' board. At the top, there's a header with search, filter, and epic dropdown buttons. The board has four columns: 'TO DO', 'IN PROGRESS', 'IN REVIEW', and 'DONE'. Tasks are categorized by epic:

- Scan Files for URLs (URL DETECTION AND EXTRACTION)**:
  - SCRM-14 (In Progress)
  - SCRM-15 (In Progress)
  - SCRM-19 (To Do)
- Intermediate presentation (DOCUMENTATION & PRESENTATION)**:
  - SCRM-46 (In Progress)
- Creation of base structure (CODE QUALITY AND MAINTAINABILITY)**:
  - SCRM-45 (In Progress)
- Intermediate documentation (DOCUMENTATION & PRESENTATION)**:
  - SCRM-47 (To Do)

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

53

## Sprint 2

### Categories

#### Sprint 2

This screenshot is identical to the one above, showing the Jira 'Sprint 2' board with tasks in 'TO DO', 'IN PROGRESS', 'IN REVIEW', and 'DONE' states across four categories. The tasks and their status are the same as in the previous screenshot.

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

54

## User Story

Description, Acceptance Criteria, DOR and DOD

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

Scan Files for URLs

Checklist for Jira On-the-Fly

To Do Actions

Definition of Ready 100/100%

Add new checklist item

Details

Assignee Aladin Vepiel

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.

Priority High

Business Value 8

Acceptance Criteria:

- System can detect URLs in a variety of file formats including BB, 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.

Story point estimate 8

Sprint Sprint 2

Definition of Done 0/100%

Add new checklist item

- Coding standards and best practices are implemented
- Unit tests for the feature are written and pass
- Any changes to the code or functionality are documented

Berner Fachhochschule | Haute école spécialisée bernoise | Bern University of Applied Sciences

55

## Our learnings

Coding	Scrum
Branch per User Story	Create Tasks
Commit Messages	User Story good size
Code Reviews	Weeklys are valuable
Basic code structure	



## Questions?