

Chapter 1

Planning

1.1 Planning concept

So as to plan the project, a combination of the two well known planning frameworks scrum and RUP are used.

For a first rough planning, the assignment is split into working packages and assigned to milestones. Delivery objects are defined for each milestone.

This plan is then assigned to the given time table of about 12 weeks. The project effort is defined as 180 hours. This results in about 15 hours work load per week.

A more detailed planning is done for the incoming milestone / sprint. The predefined working packages are split into smaller packages. For the first draft, only the first milestone is split into smaller packages. The later milestones are going to be defined in more detail as soon as all needed information is available.

The effort needed for the documentation is not listed separately. All the tasks already contain additional time for updating the documentation.

The milestones dates are not finally defined, which means that the meeting dates can vary by up to some days.

1.2 Milestones overview

Name	Shortcut	Weeks	Estimated hours	Hours total	Closing date
Milestone one	m1	2.5	39	39	01.10.2014
Milestone two	m2	3	45	84	22.10.2014
Milestone three	m3	2	30	114	05.11.2014
Milestone four	m4	2	30	144	19.11.2014
Milestone five	m5	2.5	38	182	08.12.2014

1.3 Delivery objects

Milestone	Delivery objects
Milestone one	<ul style="list-style-type: none">• System specification• Sketch software architecture• Short presentation CI environment• Draft risk evaluation
Milestone two	<ul style="list-style-type: none">• Elaborated software architecture• Tested code of test framework (tbd: which components)• Interface definition for justext/boilerplate components• HTML test data
Milestone three	<ul style="list-style-type: none">• Working test environment with both justext and boilerplate components integrated
Milestone four	<ul style="list-style-type: none">• Evaluation environment for output data of test framework• First approach to new algorithm
Milestone five	<ul style="list-style-type: none">• Implementation of new algorithm• Final documentation• Final presentation

1.4 Milestone one - m1

- Closing date date: 1.10.2014
- Available time: ca. 39h

Story	Shortcut	Estimated time
Planning	s1	4h
Research HTML / Algorithms	s2	8h
System specification	s3	12h
Risk evaluation	s4	3h
Draft software architecture	s5	8h
Configuration CI environment	s6	4h
Total		39h

1.4.1 Stories m1

Title	Planning
Id	s0
Estimated time	4h
Description	As a project owner, you need to have a time schedule so that you can see when you will achieve which results. The PAWI project is split into several working packages which are then split into single stories. The working packages are assignment to milestones and for each milestone, delivery objects are defined. This can be a document, a piece of test or production code or some other kind of work.

Title	Research HTML / Algorithms
Id	s1
Estimated time	8h
Description	My knowledge of HTML and content extraction algorithms is still limited. In order to find out what challenges I will face and which aspects I will have to take into consideration for performing the first tasks, a short research on these topics is needed.

Title	System specification
Id	s2
Estimated time	12h
Description	The PAWI project is defined through a short project description. This description does not cover all necessary information to both plan and perform this project. The key features, interfaces and delivered objects have to be defined more closely. The system specification should cover all these requirements.

Title	Draft software architecture
Id	s3
Estimated time	8h
Description	A first rough software architecture should be made as soon as possible, so that any misunderstandings between tutors and student can be uncovered. Moreover, it is much easier to plan the further steps when the software is split into several parts.

Title	Risk evaluation
Id	s4
Estimated time	8h
Description	Potential risks should be uncovered with the knowledge that was gathered by defining the specification and the software architecture. What is more, further actions can be defined to minimize the above mentioned risks.

Title	Configuration CI environment
Id	s5
Estimated time	4h
Description	<p>To deliver high quality software a continuous integration environment is required. Following tools should be evaluated and configured for further use:</p> <ul style="list-style-type: none">• Version control (git)• Project build automation tool (gradle)• continuous integration service (Travis CI)

1.5 Milestone two - m2

- Closing date date: 22.10.2014
- Available time: ca. 45h

Story	Shortcut	Estimated time
Implementation test framework	s6	20h
Prototype Integration of justext/boilerpipe	s7	17h
Collection of test data	s8	8h
Total		45h

1.5.1 Stories m2

Title	Implementation Testframework
Id	s6
Estimated time	20h
Description	Implementation of a first part of the test framework. This story will be divided into smaller stories as soon as the software architecture and the system specification is reviewed.

Title	Prototype Integration of justext/boilerpipe
Id	s7
Estimated time	4h
Description	Implementation of a small prototype which uses the existing implementation of justext and boilerpipe. A final interface for both components needs to be defined for further use. This story will be divided into smaller stories as soon as the software architecture and the system specification is reviewed.

Title	Collection of test data
Id	s8
Estimated time	8h
Description	To evaluate the functionality of the text extraction algorithms, a certain amount of test data is needed. This test data contains HTML files of several web pages. The HTML code is categorized into content and boilerplate.

1.6 Milestone three - m3

- Closing date date: 5.11.2014
- Available time: ca. 30

Story	Shortcut	Estimated time
Implementation test framework	s9	20h
Final integration of justext / boilerplate	s10	10h
Total		30h

1.6.1 Stories m3

Title	Implementation test framework
Id	s9
Estimated time	20h
Description	Final implementation of the test framework. This story will be divided into smaller stories as soon as the software architecture and the system specification is reviewed.

Title	Prototype Integration of justext/boilerpipe
Id	s10
Estimated time	4h
Description	Complete integration of the justext and boilerplate algorithms into the test framework. This story will be divided into smaller stories as soon as the software architecture and the system specification is reviewed.

1.7 Milestone four - m4

- Closing date date: 19.11.2014
- Available time: ca. 30h

Story	Shortcut	Estimated time
Evaluation environment for results	s11	20h
Research on new algorithm	s12	10h
Total		30h

1.7.1 Stories m4

Title	Evaluation environment of results
Id	s11
Estimated time	20h
Description	The test framework will produce a lot of output data, which has to be reviewed using an evaluation environment. This should process this data and present the results in a descriptive way. This story will be divided into smaller stories as soon as the software architecture and the system specification is reviewed.

Title	Research on new algorithm
Id	s12
Estimated time	20h
Description	A first research on the new algorithm should be performed. After this research it should be possible to decide if this solution is possible and if an implementation with the remaining time resources is realistic. This story will be divided into smaller stories as soon as the software architecture and the system specification is reviewed.

1.8 Milestone five - m5

- Closing date date: 8.12.2014
- Available time: ca. 38h

Story	Shortcut	Estimated time
Implementation of new algorithm	s13	19h
Complete documentation	s14	15h
Prepare final presentation	s15	4h
Total		38h

1.8.1 Stories m5

Title	Implementation of new algorithm
Id	s13
Estimated time	19h
Description	Implementation of the new algorithm and analysis of the test results with the existing evaluation environment.

Title	Complete documentation
Id	s14
Estimated time	15h
Description	Complete and review all chapters of the documentation.

Title	Prepare final presentation
Id	s15
Estimated time	4h
Description	Prepare the final presentation and the final printed / digital version of the thesis.