



News Aggregator with Python and Machine Learning 1st Audit

Client: Thomas Griffith

Team Members: Yuliang Zhang, Jiahua Liang, Chen Zhang, Luokun Gong, Vishnu Vardhan Jasti, Jun Yang

Background



News Channel Problem:

- Traditional Media Recession
- New Media Fragmentisation & Vulgar

Difficult to get **high-value** news **efficiently**



Solution:

- A News Aggregator Focusing on **High-Value** News
- Using Web Crawler and Machine Learning Tools

<http://gpu.jkl.io:5000/>



Continuous Project Problem:

- Vague User Scope
- Rely on Algorithm, difficult to estimate user satisfaction
- Insufficient news channels
- Existing Code Bugs
 - Bad BERT Results
 - Crawler Stops Occasionally
 - Only 1 news source could be stored in data base

Solution??

Client's Expectation & Requirements



More Canberra Focused:

- More Extraction Work on Canberra Based News Channels
- Improve Algorithm for Canberra residents



User Surveys:

- Conduct 2-3 User Surveys
- Improve our work according to user feedback



Improve Web UI

Stakeholders

Client: Thomas Griffith



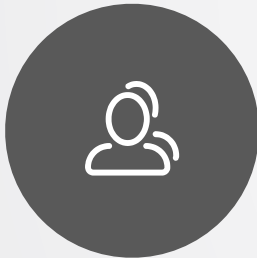
Provide: development resources

- server, web frame, data, etc

Expect:

- A Canberra Focused, operable news website
- 2-3 User Survey Reports

Team Members



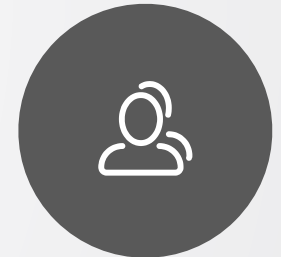
Provide:

- Time and effort

Expect:

- Improve code skills
- Improve teamwork skills
- Get good marks

Users



Provide:

- Time to use the website
- Time to do user survey

Expect:

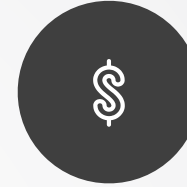
- A efficient website for getting high-value news

Resources & Costs



Resources:

- Kosmos Dataset (news dataset including more than 300,000 articles)
- GPU servers (used for training machine learning model)
- The code source and products from last semester's work
- The Bootstrap or other frames for web development
- Some Canberra focused news websites.



Costs:

Development related costs will be solved by our client.

Risks



Risks:

- Fail to achieve the expected accuracy and efficiency.
- Fail to achieve final integrate of different work groups.
- Timeouts for debugging and final deliverable.
- The risk of potential dissatisfaction from the users
- The risk of failing to respond to the feedback from the user survey



Risk Management:

- Review work every 2 weeks. If any progress is falling behind, try: switch methods/ add workload/ adjust team structure...
- Conduct first user survey early (week 5), so we can have sufficient time to response to survey results.
- Communicate with client and tutors timely.

Tools & Constrains

Development Tools



Github



Python3



Google Doc



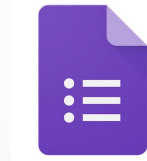
MongoDB



Flask



Bootstrap



Google Form



Gensim



Keras



Pytorch



Tensorflow

Communication Tools



Slack
With Client



WeChat
With Team members

- In-Person Client Meeting every 2 weeks.
- Online Client Meeting all other weeks.
- Team Meeting whenever necessary.

Process Management Tool



Trello

<https://trello.com/b/vxJ6hSCv/news-aggregator-sem2>

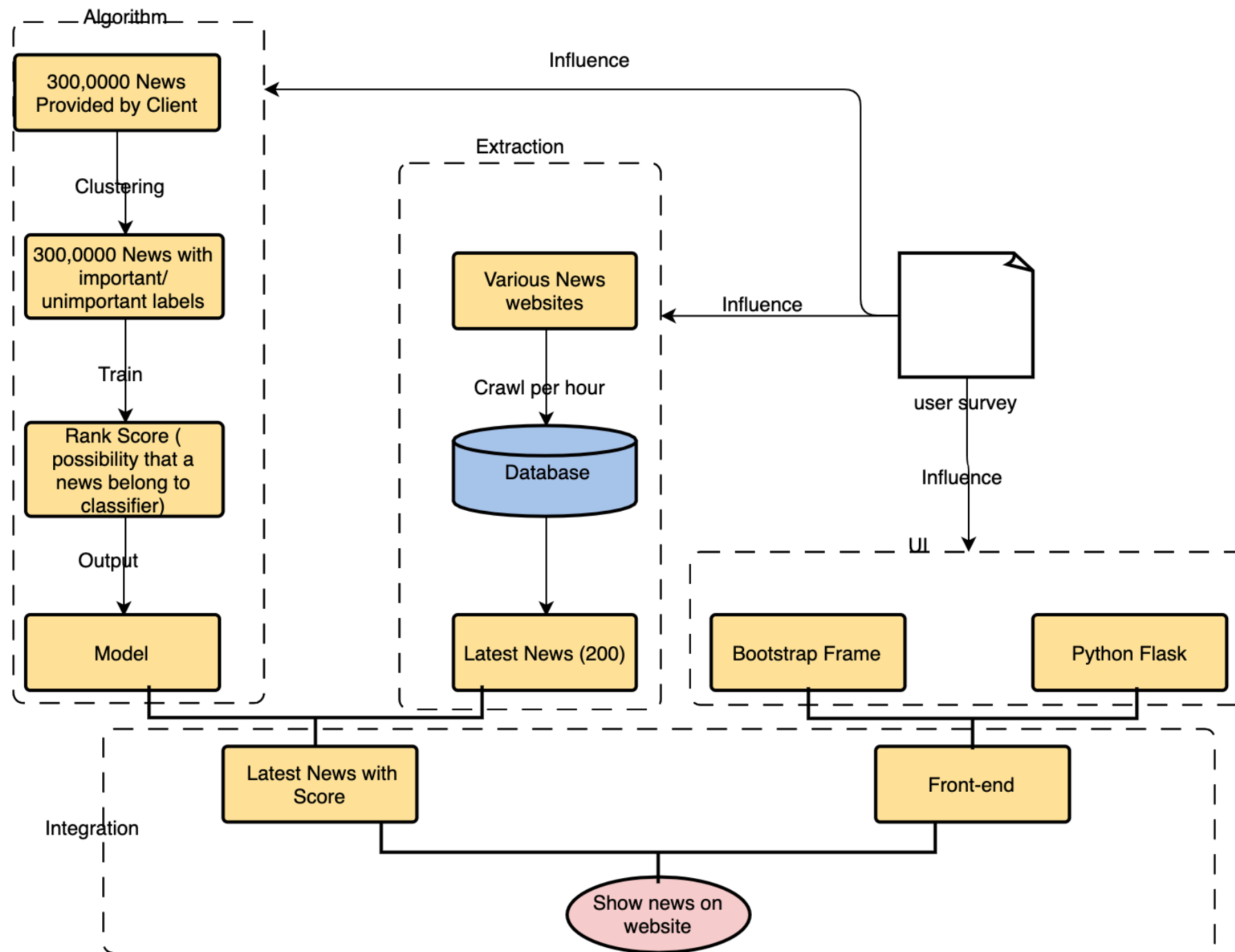
Teamwork

Name	uid	Technical Role	Progress Management Role
Jiahua Liang	u6162679	UI & User	Spokesperson 1
Luokun Gong	u5917339		Minutes Taker
Xiangyun Kong	u6556183		Progress Tracker
Yulinag Zhang	u6782445	Extraction & Algorithm	Spokesperson 2
Vishnu Vardhan Jasti	u6611697		Clerical Assistant
Jun Yang	u6767560		Clerical Assistant
Chen Zhang	u6745297		Progress Tracker

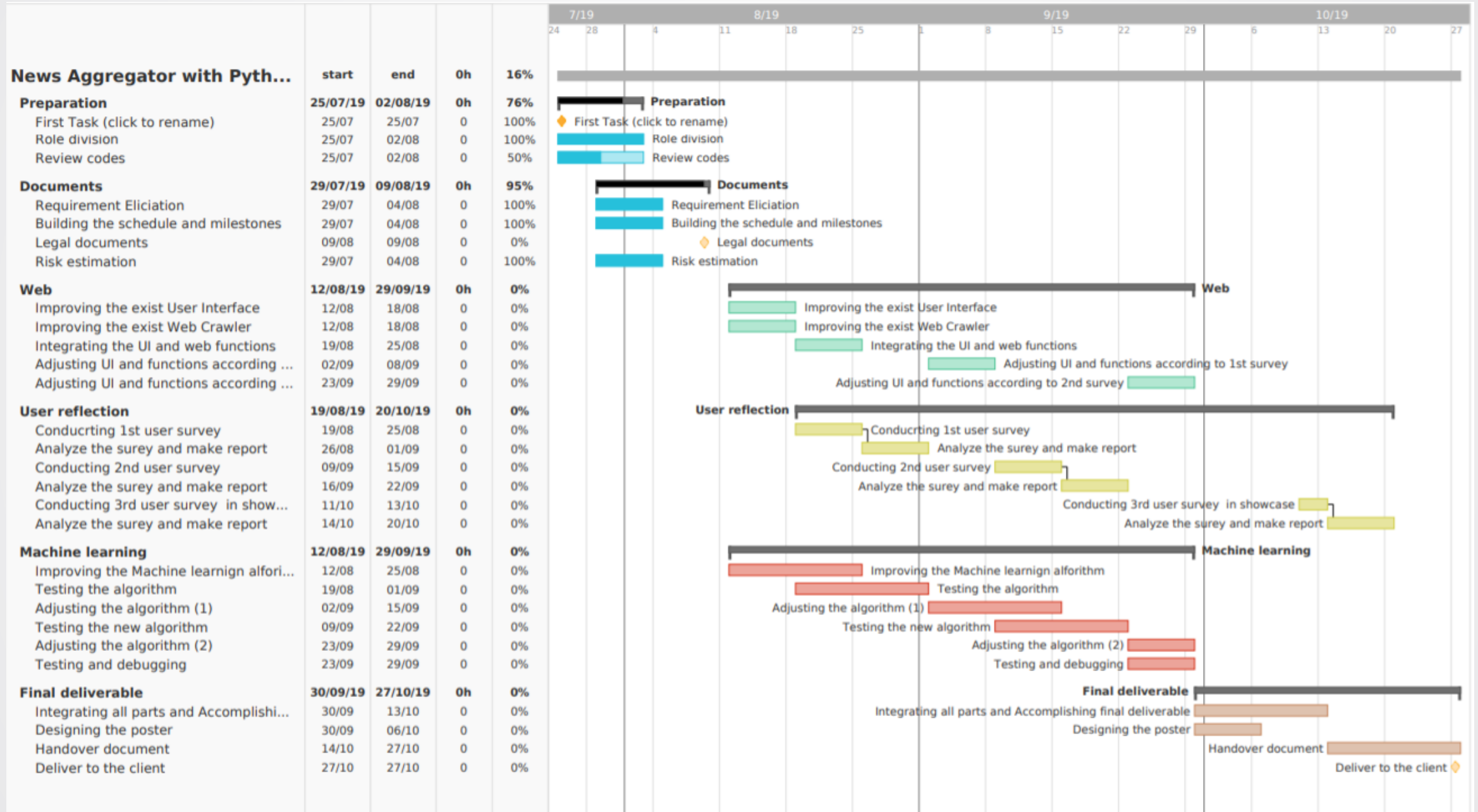
Schedule

Week 1-2	<ul style="list-style-type: none">• Form Team• Divide Roles• Set up development environment	Week 7-8	<ul style="list-style-type: none">• Second user survey• Reflect to user survey results
Week 3-4	<ul style="list-style-type: none">• Improve extraction work• Improve algorithm• Improve UI	Week 9-10	<ul style="list-style-type: none">• Poster• Final deliverable• Showcase (3rd user survey)
Week 5-6	<ul style="list-style-type: none">• First user survey• Reflect to user survey results	Week 11-12	<ul style="list-style-type: none">• Handover documents• Report to client

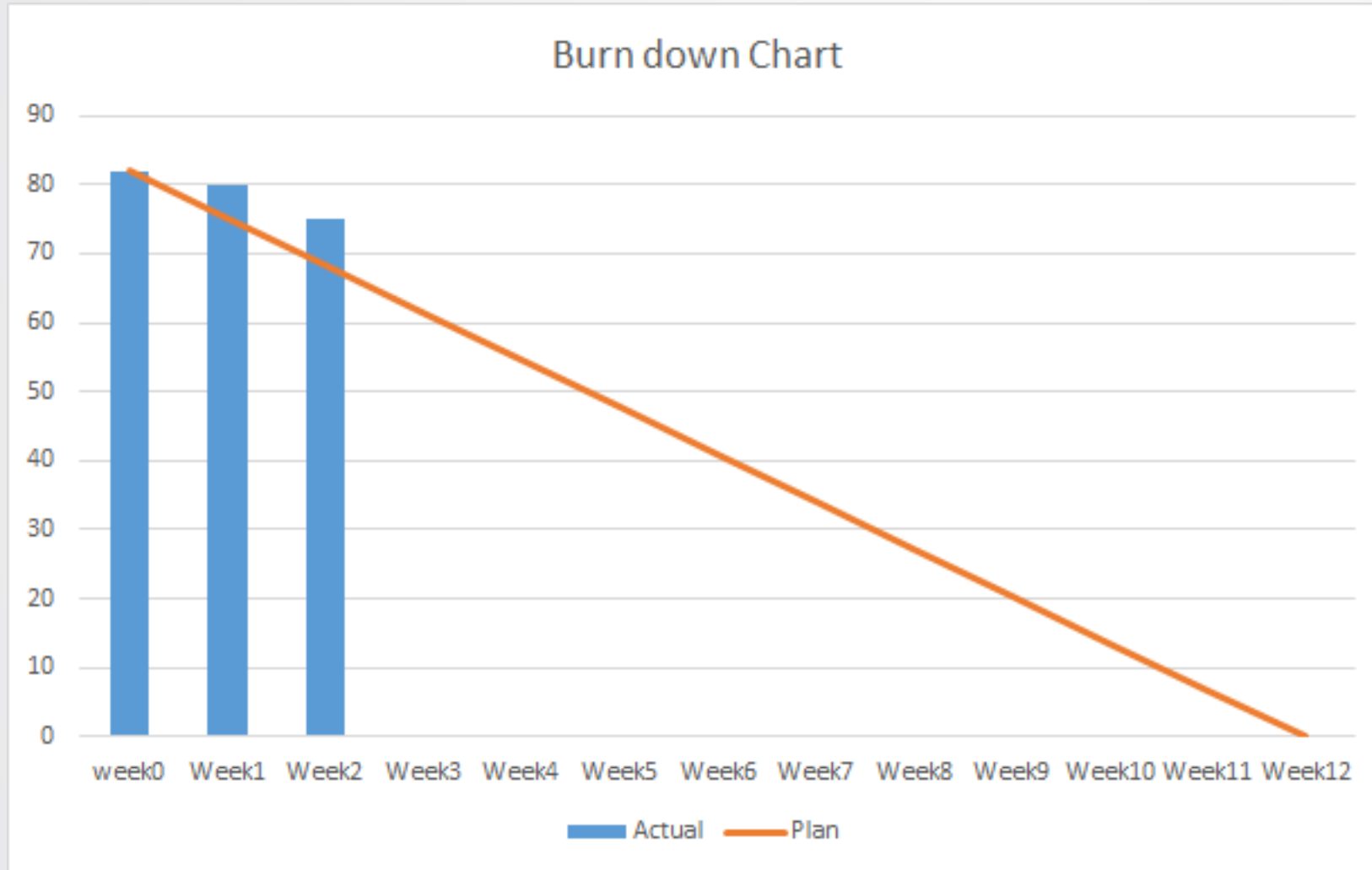
Workflow



Gantt chart



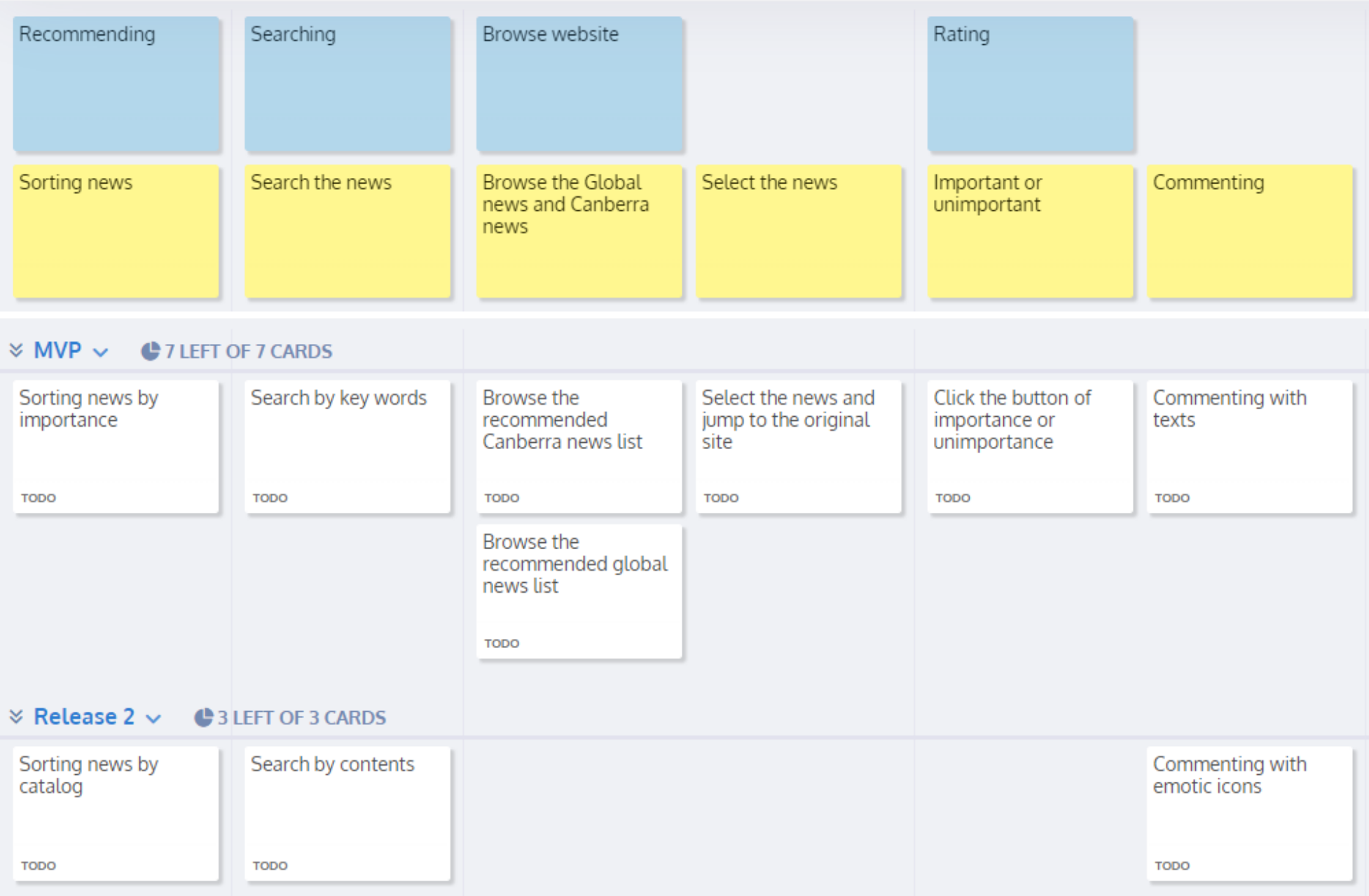
Burndown Chart



The Burndown Chart is calculated according to our Trello record:

<https://trello.com/b/vxJ6hSCv/news-aggregator-sem2>

User Story Map



User Story Point Matrix

Story points	User story
5	<ul style="list-style-type: none">• As a visitor, I would like to select the news so that I can read its details on its original site.• As a visitor, I would like to click the buttons of important and unimportant so that I can give feedback on the news
10	<ul style="list-style-type: none">• As a visitor, I would like to browse the Canberra news list, so that, I can know the events happening in Canberra. that more likely to be related to me.• As a visitor, I would like to browse the world news, so that I can know the important events happening in the world.• As a visitor, I would like to write texts to comment on the news, so that I can describe my feeling about the news more specifically.
15	<ul style="list-style-type: none">• As a visitor, I would like to search the specified news by keywords, so that I do not need to waste time• to scroll down the news list to find the news.
30	<ul style="list-style-type: none">• As a visitor, I would like to only read the most important news at the top of their lists, so that I can avoid wasting time on looking for the most important news.



Thank you