

Progress Report #1

<https://cosc4p02.tpgc.me/>

Members:

Maulik Mann	6994214	mm20as@brocku.ca
Adrian Binu	6970677	ab19xe@brocku.ca
Thanushan Pirapakaran	6890206	tp19uy@brocku.ca
Vishant Khamar	7062920	vk20if@brocku.ca
Devaan Zastrow	6932842	dz19xr@brocku.ca
Michael Boulos	6973523	mb19ep@brocku.ca
Jason Hunter	6744247	jh19pp@brocku.ca
Brett Terpstra	6920201	bt19ex@brocku.ca

PROJECT NAME	Web Summarizer and URL Shortener		OVERALL PROGRESS
	PROJECT MANAGER	START DATE	
		01-29	
SCRUM MASTER	Vishant	04-28	≈35%

TASK NAME	FEATURE TYPE	RESPONSIBLE	START	FINISH	DURATION in days	STATUS	COMMENTS
SPRINT 1			01-29	02-14	17		
Deploy Website		Brett and Michael	01-29	02-14	17	Complete	
Create URL Shortener		Maulik	01-29	02-14	17	Complete	URL shortener created, adding to website pending
Use LLM to summarize		Brett, Michael, Deevan, Vishant	01-29	02-14	17	In Progress	LLM is currently being trained and adjusted, accuracy of the summarization provided by the LLM is a challenge
Create UI		Jason, Adrian, Thanushan	01-29	02-14	17	Complete	
SPRINT 2			02-14	02-27	14		
Log in with Google		Maulik	02-14	02-27	14	Complete	
Pro User Dashboard		Jason, Adrian, Thanushan	02-14	02-27	14	In Progress	
Incorporate Web Summarizer into site		Brett, Michael, Deevan, Vishant	02-14	02-27	14	In Progress	
Incorporate URL shortener into site		Maulik	02-14	02-27	14	In Progress	Difficulty with implementing the url redirection on site
SPRINT 3			02-27	03-12	15		
Pro Features: Summarization/Url Shortener history, etc.		Jason, Adrian, Thanushan, Maulik	02-27	03-12	15	Not Started	
Pro feature: Summarization Adjustment		Brett, Michael, Deevan, Vishant	02-27	03-12	15	Not Started	
SPRINT 4			03-12	03-26	15		

Implementation of SE Process:

Our SE Process is Scrum. We hold weekly scrum meetings to review and demonstrate what we have done in the last week. Everyone shows their progress and discusses any challenges they are facing. We are using the release planning document as our backlog and to set deliverables.

Detail of system:

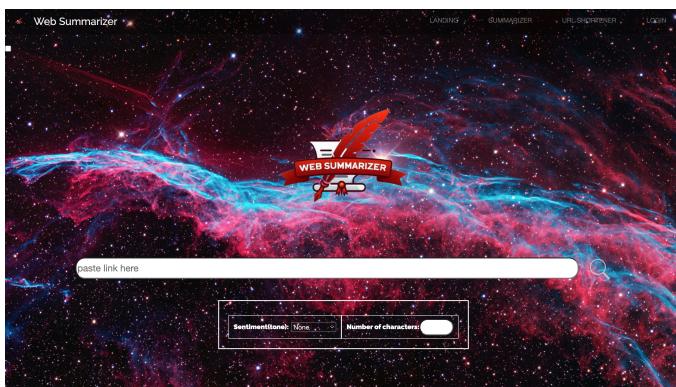
It is a website with pages for the different functionalities and features. Pro features are hidden until the user signs in with Google.

Web Team Log:

Challenges:

- UI: Splitting and distributing the work, while trying to make a seamless website that is consistent throughout
- URL Shortener: Integrating the URL shortener into the website and having the shortened urls redirect the user to the actual site

Screenshots:



The image consists of four vertically stacked screenshots of a web application with a dark, space-themed background featuring nebulae.

- Summary:** A modal window titled "Summary:" is open, showing a "Sample Result" area which is currently empty. At the top of the modal are buttons for "untitled" and "Save".
- URL Shortener:** The main dashboard shows a search bar with "paste link here" placeholder text and a magnifying glass icon. The "WEB SUMMARIZER" logo is centered above the search bar.
- Shortened URL:** A modal window titled "Shortened URL:" is open, displaying a "Sample URL" field containing "http://tinyurl.com/4t9mzg". Below it is a "Number of Uses:" field set to "0000".
- Login:** A login screen comparing "Basic Version" and "Pro Version". The "Basic Version" features:
 - Summarize videos
 - Summarize web pages
 - Shorten URLs
 - Text to Speech
 - Dark modeThe "Pro Version" features:
 - Control the character count of summarizations
 - Control the sentiment tone of summarizations
 - Save summarizations and shortened URLs
 - Access saved files in the dashboard
 - See the number of times a shortened link has been usedAt the bottom, there are "Sign in as moekki" and "Google" login buttons, and a "Login for Pro Features" link.

LLM Team Log:

The majority of the LLM (large language model) team's development has occurred on our private server, which is not visible on the GitHub (DNS + DDNS config, reverse proxy setup, hypervisor + VM config, LLM weights, etc). Some efforts are being made to cause our work to be more visible, such as supplying the docker run script. However, as much of the initial LLM setup involves configuring the underlying host this is not always possible. More elements will be added to the GitHub over the course of the next sprint.

Challenges:

- The large storage requirements of the LLM weights

```
du -hs /mnt/archive/public/llms
269M    COSC-4P02
732G      llms

du -hs /mnt/cache/vms/michael_server.qcow2
65G        /mnt/cache/vms/michael_server.qcow2
```

- High memory and CPU usage of the LLM

CONTAINER ID	NAME	CPU %	MEM USAGE / LIMIT	MEM %	NET I/O	BLOCK I/O	PIDS
24d386e223f1	llms	1492.00%	13.59GiB / 62.72GiB	21.67%	1.79MB / 84.3kB	35.1GB / 4.42GB	269

- Slow generation / response time of the LLM due to lack of specialized hardware

```
time curl -X POST http://localhost:6980/v1/completions -H "Content-Type: application/json" -d '{"model": "llama-2-13b-chat/gptj-model-05.1.qpuf", "prompt": "in a little tea pot", "temperature": 0.7}
{"created":1707872898,"object":"text_completion","id":"fb5ca0f6-a204-492e-a734-83662442df73","model":"llama-2-13b-chat/gptj-model-05.1.qpuf","choices":[{"index":0,"finish_reason":"stop","text":"\u2022\vn\ni have a little handle \u2022\vn\ni ca
n pour a cup \u2022\vn\ni'm a little tea pot \u2022\vn\ni have a little spout \u2022\vn\ni can make you a pot \u2022\vn\ni'm a little tea pot \u2022\vn"}],"usage":{"prompt_tokens":0,"completion_tokens":0,"total_tokens":0}}curl http://localhost:6980/v1/completions -H "Content-Type: application/json"
0.01s user 0.01s system 0% cpu 1:15.94 total
```

(1 minute 16 seconds on a small prompt.)

- Resource constraints mean we must use the smallest Llama which suffers from accuracy and generation issues.
- Ensuring a balance between the accuracy and speed of the LLM, requiring both prompt engineering and understanding of the underlying Llama model

Brett / System Administrator's Note:

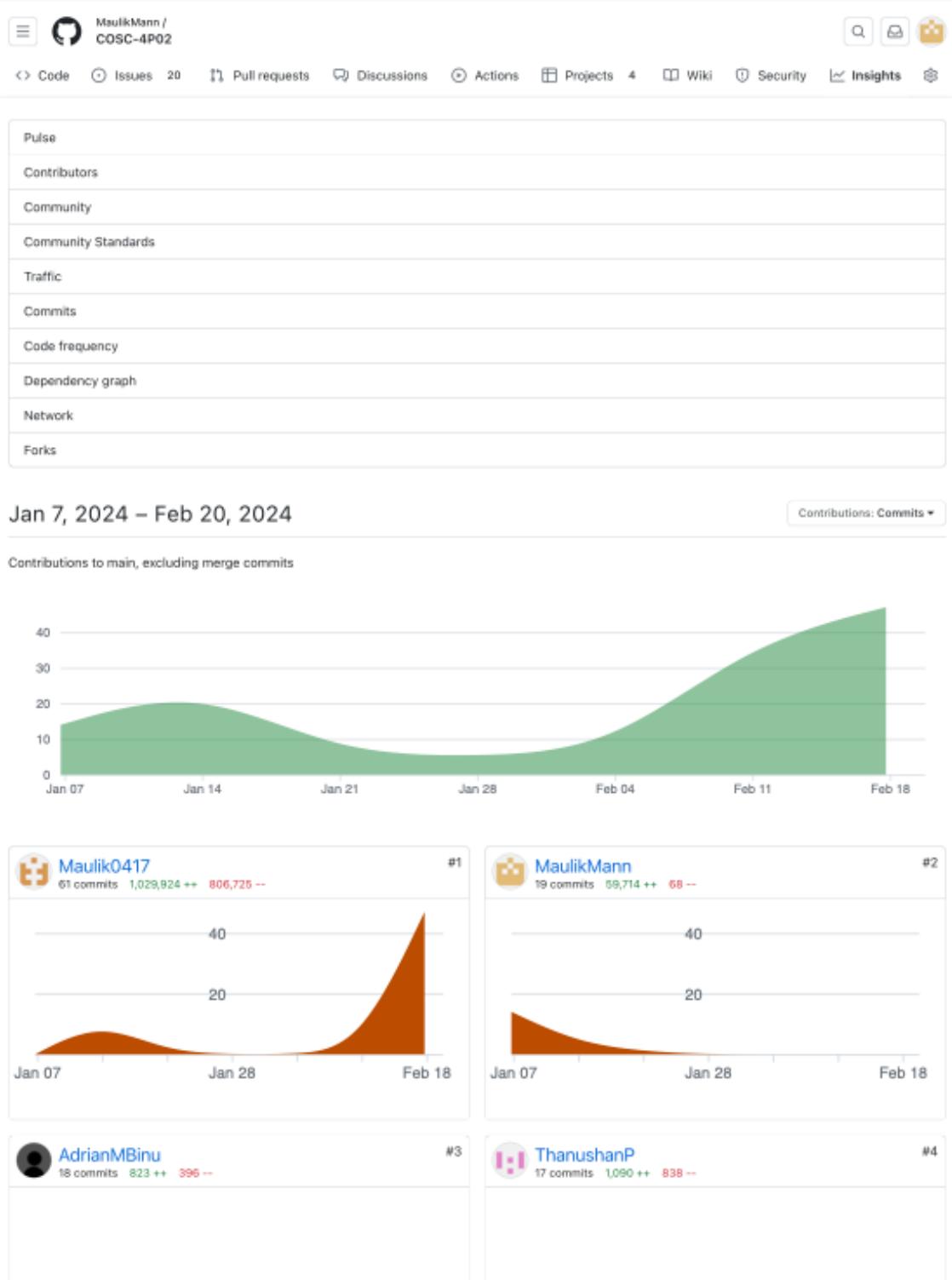
Access to the underlying host which manages the shared VM and executes the LLM will not be granted to anyone under any circumstances. Access to the cosc4p02 VM can be granted to anyone involved in the project, including stakeholders (the TA), but I am unwilling to compromise the security of my hardware.

GitHub contribution log:

(Next page)

2/20/24, 1:13 PM

Contributors to MaulikMann/COSC-4P02

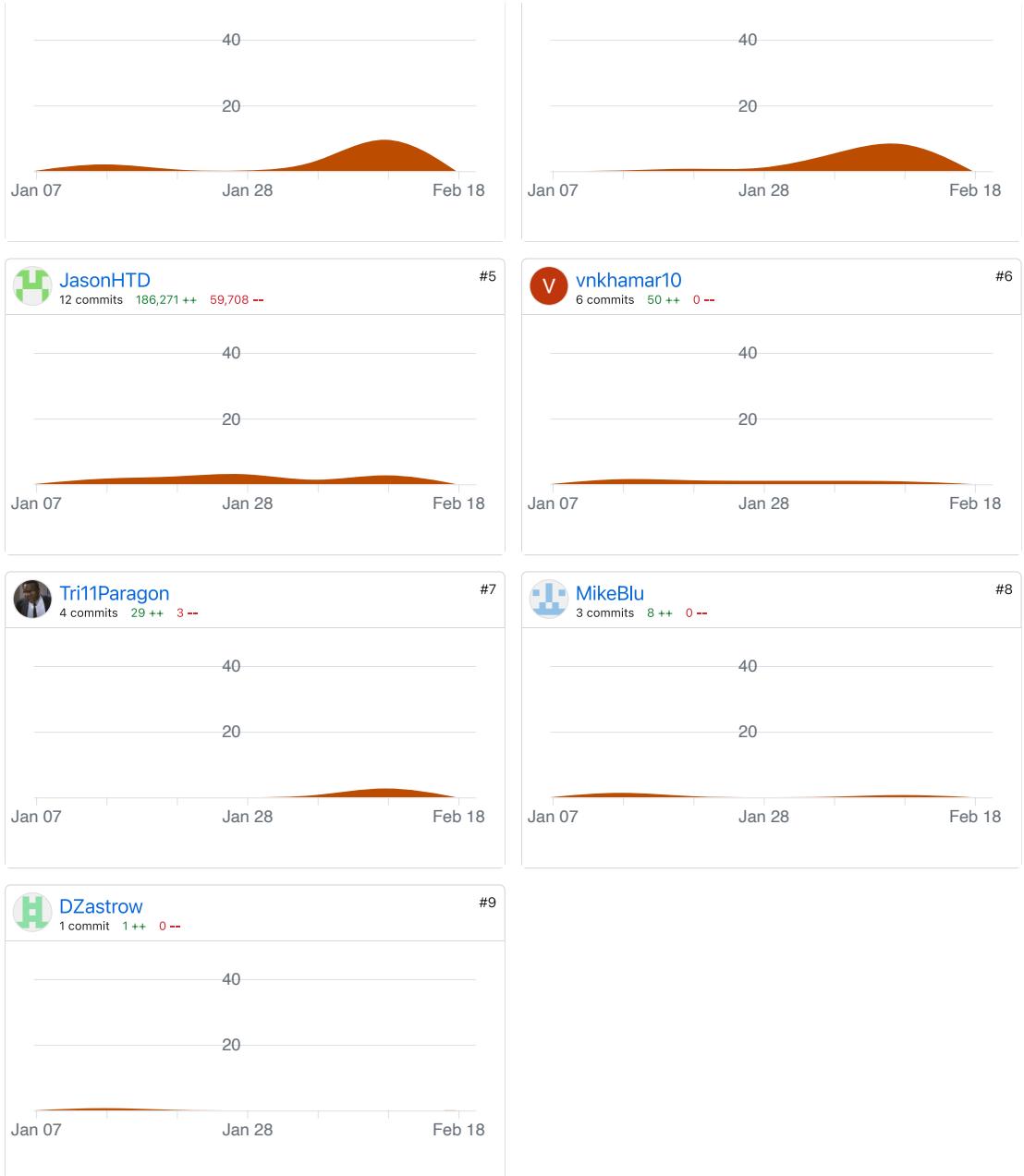


<https://github.com/MaulikMann/COSC-4P02/graphs/contributors>

1/2

2/20/24, 1:13 PM

Contributors to MaulikMann/COSC-4P02



Code frequency over the history of MaulikMann/COSC-4P02

