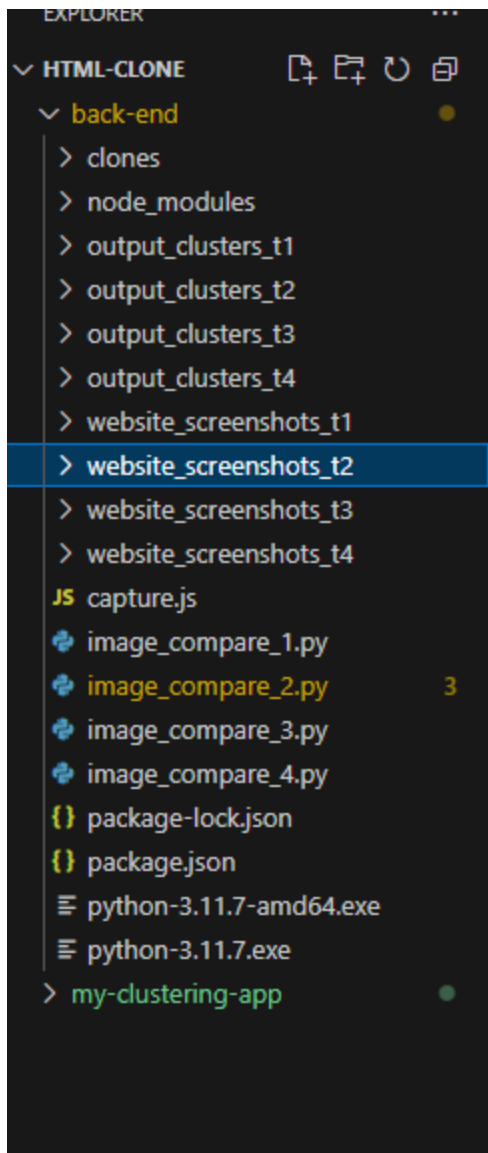**Why did I choose this task:** A while back, I participated in a hackathon (Hack-It-All) where I used sklearn and LLMs to create specific recommendations for users, and I loved every second of the learning and development process. As soon as I saw the task description, I thought that maybe sklearn and AI could help me, and I was right.

**My initial thought process:** At first, when I saw the Tier 1 websites, I decided to just compare the metadata (code) between websites, and after some tweaking, I got where I wanted. The first tier was pretty easy, but the second one was where I really got AI and LLMs involved in my code because it got difficult. Metadata comparison was no longer enough due to the different languages in the code and similar aspects of the content. I decided to use tools like sklearn, TensorFlow, and Selenium to get through it all. First, the code accesses all websites with the help of ChromeDriver, a powerful tool for automation. Then, Selenium is used to take screenshots of every website. After all this, every screenshot is compared using LLMs (Sentence-BER and others) and stored separately in clusters based on similarity. The obvious choice of language for the back-end was python, due to powerful tools used by pip3, and for front end I chose Node.js.
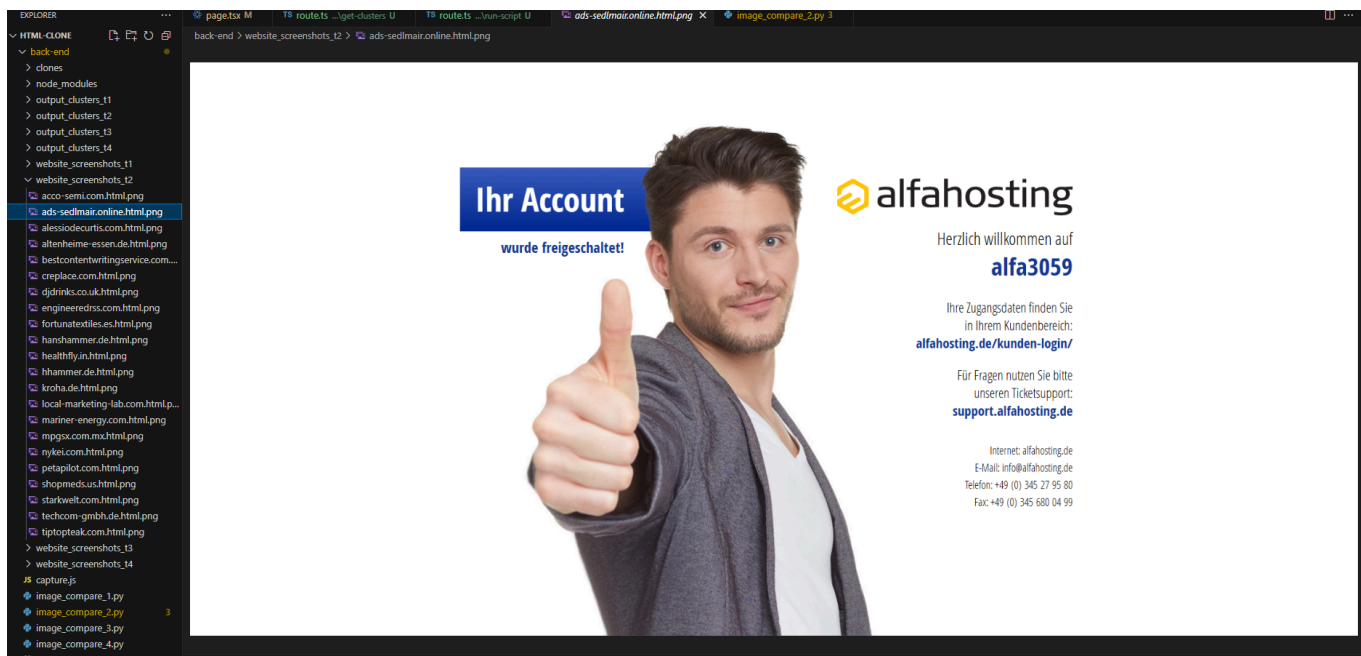
I applied this algorithm to all tiers, with the only difference being the output clusters and website screenshot names. The code is a bit of an overkill for the first tier, but it works like a charm for every tier. That being said, my code and algorithm only focus on the visual aspect of the page because that's what users actually see—the text isn't important. This approach is highly scalable and can work with very large datasets and HTML files; it just needs processing time.

For example, it took the algorithm 60 seconds to process all 100 websites in tier 1. That's 16 hours for 100,000 websites with the limited resources I have at home, but with better processing power, it could be reduced to just one hour for hundreds of thousands of websites. I added a front end design that generates this clusters with a press of a button, which makes the program slower because it needs to fetch information from back-end.

**How I've accomplished this:** As I said, at first, I wanted to compare every website with only the metadata that I have, but when tier 2 started to get more complex I needed something smarter. I will explain everything that I have done using photos of the code.
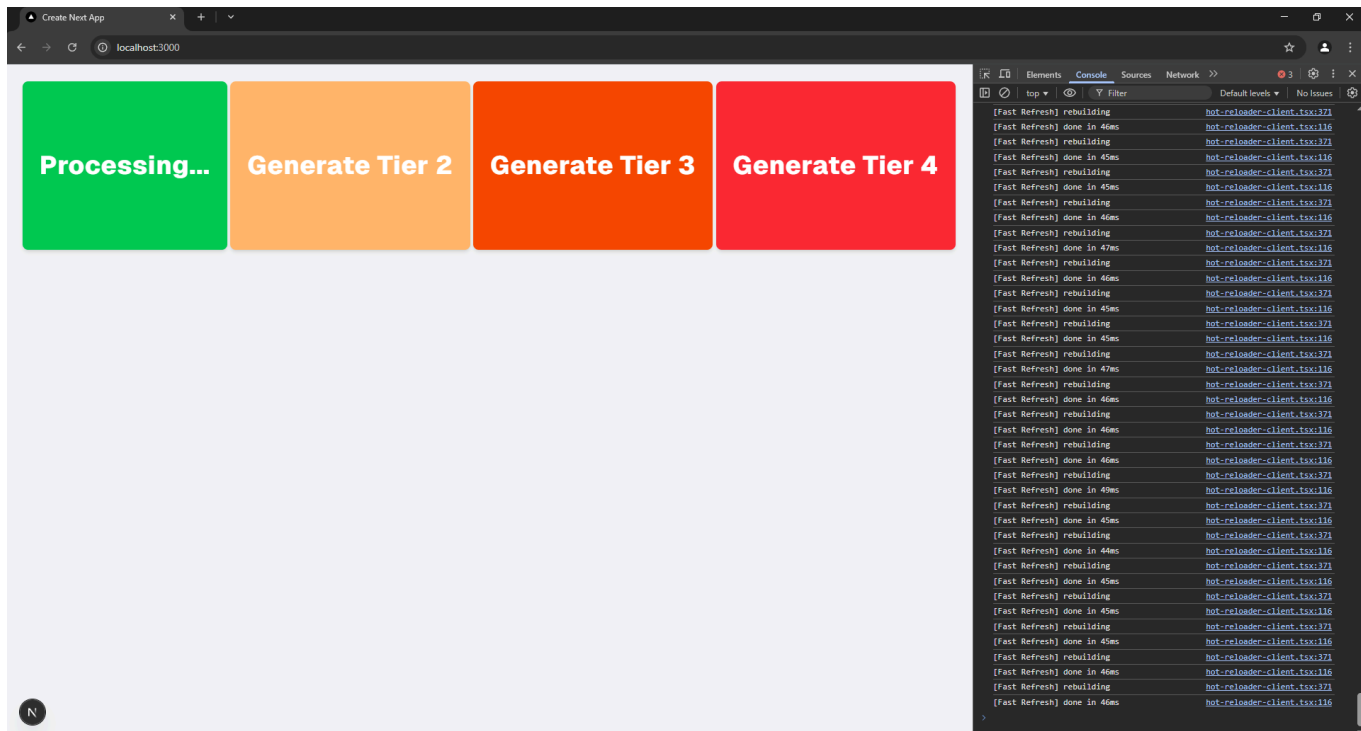
**Back-end:** This is my back-end project structure, with 4 scripts (image_compare_n.py), clones which store all the websites, website_screenshots_tn where web screenshots are stored, and output_clusters_tn where every similar website is stored in a cluster. First, the algorithm opens every website stored in a folder using chromedriver, which is an automated process whenever a function inside my code is called. When opened, it takes a screenshot that is stored in website_screenshots_tn. Then, using sklearn, TensorFlow, and selenium, each image is compared by visual similarities, just like a user sees them, and then places every website in clusters. The only drawback of this algorithm is that if the page is scrollable, the screenshot will only be of the top part of the page. Here is an example of a screenshot taken by the code:

Everything is sorted in the backend, whenever a image_compare_tn is ran it generated 2 things: website_screenshots_tn and output_clusters_tn. Every image_compare_tn is exactly the same, the only change is the websites is testing.

**Front-end:** I really wanted to add a visual representation of the algorithm that any user can use without the need of calling functions and running dev tools in code editor, so I made a Node app:



There are 4 buttons, each generating clusters for tier 1, 2, 3 and 4. When one button is pressed, a function is called that runs the image_compare_tn designated to the button number. The code runs in background and after generating it displays the results.

**For example this are the results for the first tier, and this is what is saved in the project:**



```
Cluster 1 (20 documents)
==========================================
- ..\clones\tier3\1stmortgages.london.html
- ..\clones\tier3\adeptohomes.com.html
- ..\clones\tier3\belledermaaesthetics.com.html
- ..\clones\tier3\columbiacouncilofneighborhoods.com.html
- ..\clones\tier3\deltadoula.com.html
- ..\clones\tier3\eastbourne.online.html
- ..\clones\tier3\fieldtech.info.html
- ..\clones\tier3\flyerfunnelsuk.com.html
- ..\clones\tier3\frankieswinebar.com.html
- ..\clones\tier3\furniturehutuk.com.html
- ..\clones\tier3\g3rcq.com.html
- ..\clones\tier3\hycareakarui.com.html
- ..\clones\tier3\lagustosavaldebebas.com.html
- ..\clones\tier3\lipolondon.com.html
- ..\clones\tier3\londonviptaxi.com.html
- ..\clones\tier3\okcis.info.html
- ..\clones\tier3\renautautomotive.com.html
- ..\clones\tier3\renewconsultants.com.html
- ..\clones\tier3\sophrologue-paris14.com.html
- ..\clones\tier3\thefoxinnbroadwell.com.html
```

And this, is for the 4th tier:

**And this is it:** In conclusion, this app access every html provided, takes a screenshot, compare it to others and put them in clusters. I really loved developing this app, I really hope it raises to the expectations and I want to take for the opportunity and for making such interesting task!

**References:**

1. LLM's like chatgpt and deepseek
2. Stackoverflow
3. ChromeForum
4. Selenium Documentation
5. Clustering Algorithms (DBSCAN)
6. My mind:)