CS2104 Problem Solving in Computer Science

Margaret Ellis, Naren Ramakrishnan, Sehrish Basir Nizamani





Web Scraping





What is Web Scraping

 An automatic way to retrieve unstructured data from a website and store them in a structured format.





tools and library for web scraping

















Scraping all up-front websites

- 1. Inspect the website HTML that you want to crawl
- 2. Access URL of the website using code and download all the HTML contents on the page
- 3. Format the downloaded content into a readable format
- Extract out useful information and save it into a structured format
- 5. For information displayed on multiple pages of the website, you may need to repeat steps 2–4 to have the complete information.





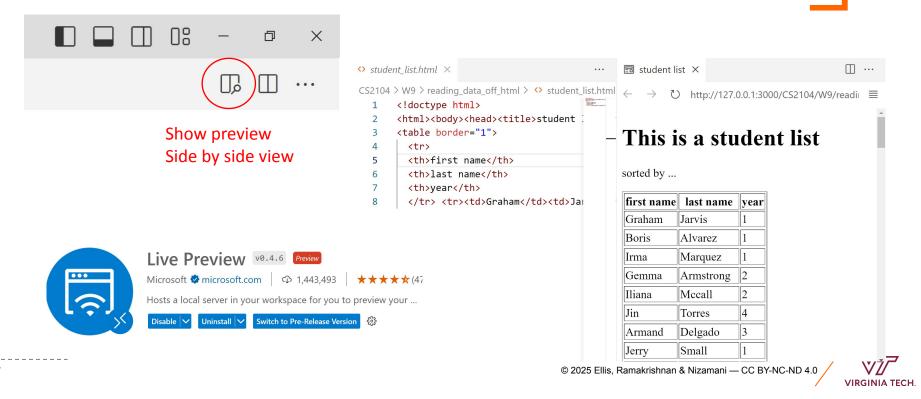
Scraping websites with API

- 1. Inspect the XHR network section of the URL that you want to crawl
- 2. Find out the request-response that gives you the data that you want
- 3. Depending on the type of request(post or get) and also the request header & payload, simulate the request in your code and retrieve the data from API. Usually, the data got from API is in a pretty neat format.
- 4. Extract out useful information that you need
- 5. For API with a limit on query size, you will need to use 'for loop' to repeatedly retrieve all the data





Classwork: reading data off html





Sample Code:

```
import re
# use regex to parse the html string and extract the target data
in filename = 'student list.html'
out filename = 'student list.txt'
# defines the name of the file that you're going to read and write into
content = open(in filename, "r")
output = open (out filename, "a")
# content and output are file stream object, you could use 'r'(read), 'w'(write) or 'a'(append)
str = content.read()
# read() method will turn a file stream into a regular string
# use matches = re.findall(pattern, str), you need to figure out what pattern to use
# pattern = r'(\w^*)<\/td>(\w^*)<\/td>(\d)'
# once you get a match, you access different components from the match object using indices.
# use output.write() to write content to the output file
# for output formatting, you can try rjust()
output.close()
content.close()
#close the stream
```



Introduction to Docker



What are containers?

Allows you to take development environment and configuration with you

- Containers are a method of operating system virtualization that allow you to run an application and its dependencies in resource-isolated processes.
- "A standardized unit of software"
- For example: https://www.docker.com/why-docker
- Amazon Web Services, Google Cloud Platform provide this







Why use containers?

Trouble:

- Path dependency
- Software version
- Environment variables

Benefits:

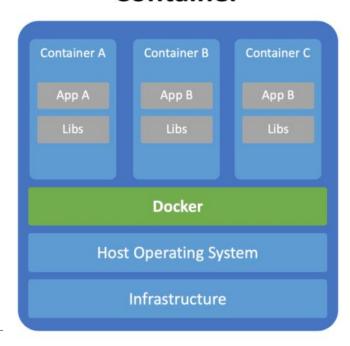
- Encapsulated-portable
- Lightweight-fast
- Demands less resource



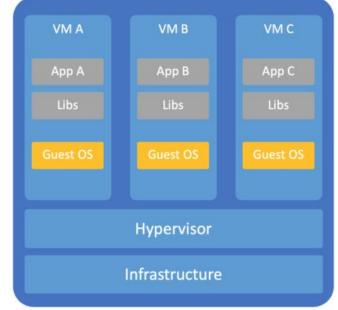


Container vs VM

Container



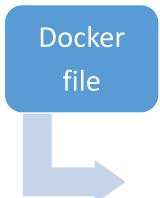
Virtual Machines







Containers *lightweight compared to virtual machines*



Configuration file (can start with templates from dockerhub such as for python ,ubuntu, node, nginx) -similar to source code

Docker image

Runnable image that is built from the Docker file and can be stored in a container registry *—similar to compiled code*

Docker container

Running instance of a docker image –similar to executable code





Classwork: Dockerize a website and deploy it on the cloud

Fork the project on Gitlab

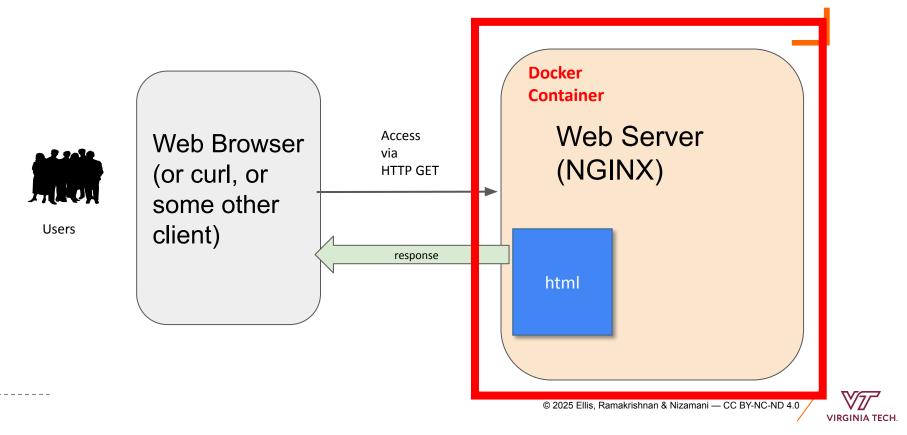


Name	Last commit	Last update
images images	Base code	1 week ago
□ site-content	Base code	1 week ago
→ Dockerfile	Base code	1 week ago
™ README.md	Update instructions to be more specific	1 day ago





Sample Static Web App Architecture

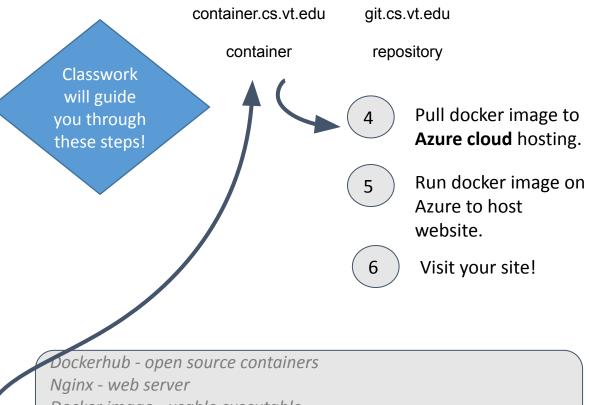


Edit the Dockerfile. Dockerfiles can use another image as a base and then build on top of it, on rlogin create a Dockerfile based on nginx:latest

> FROM nginx:latest COPY ...

Build the docker image. This takes Dockerfile source code and turns it into a Docker image, which is an executable

Push image to gitlab container **registry**. This makes it so your image can be used elsewhere without having to distribute source and rebuild.



Docker image - usable executable

Docker file - source code for Docker image

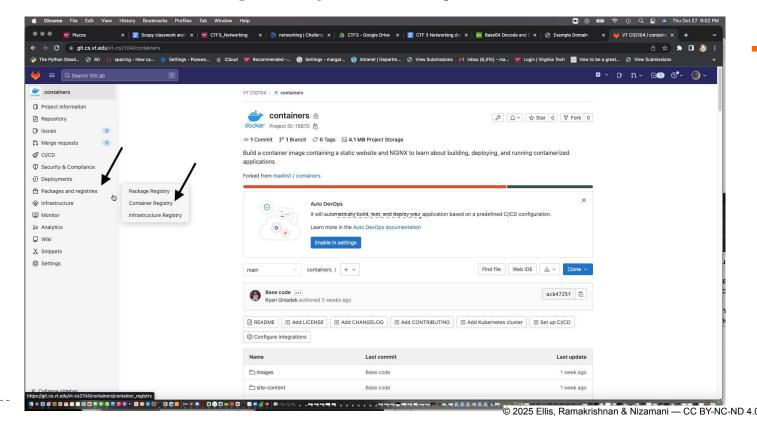
Docker container - running instance of the Docker image executable







Container registry example in GitLab

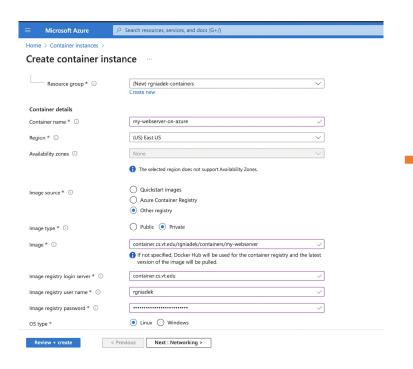








Pull container to Azure and host website!





Edit this in your index.html using vi on rlogin

