# Project: Image Retrieval
# (TA Session)

**Nguyễn Đăng Nhã**

# Objectives

## Scraping URL of Images

**urllib³**

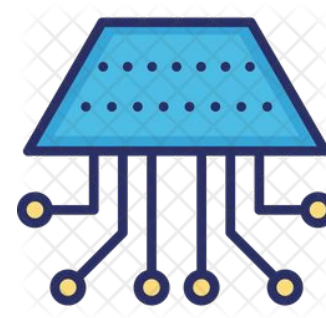Handle Request- Response

**BeautifulSoup**

Execute HTML content

Interact with
Web elements

## Getting Images from URLs

Multi-Threading for
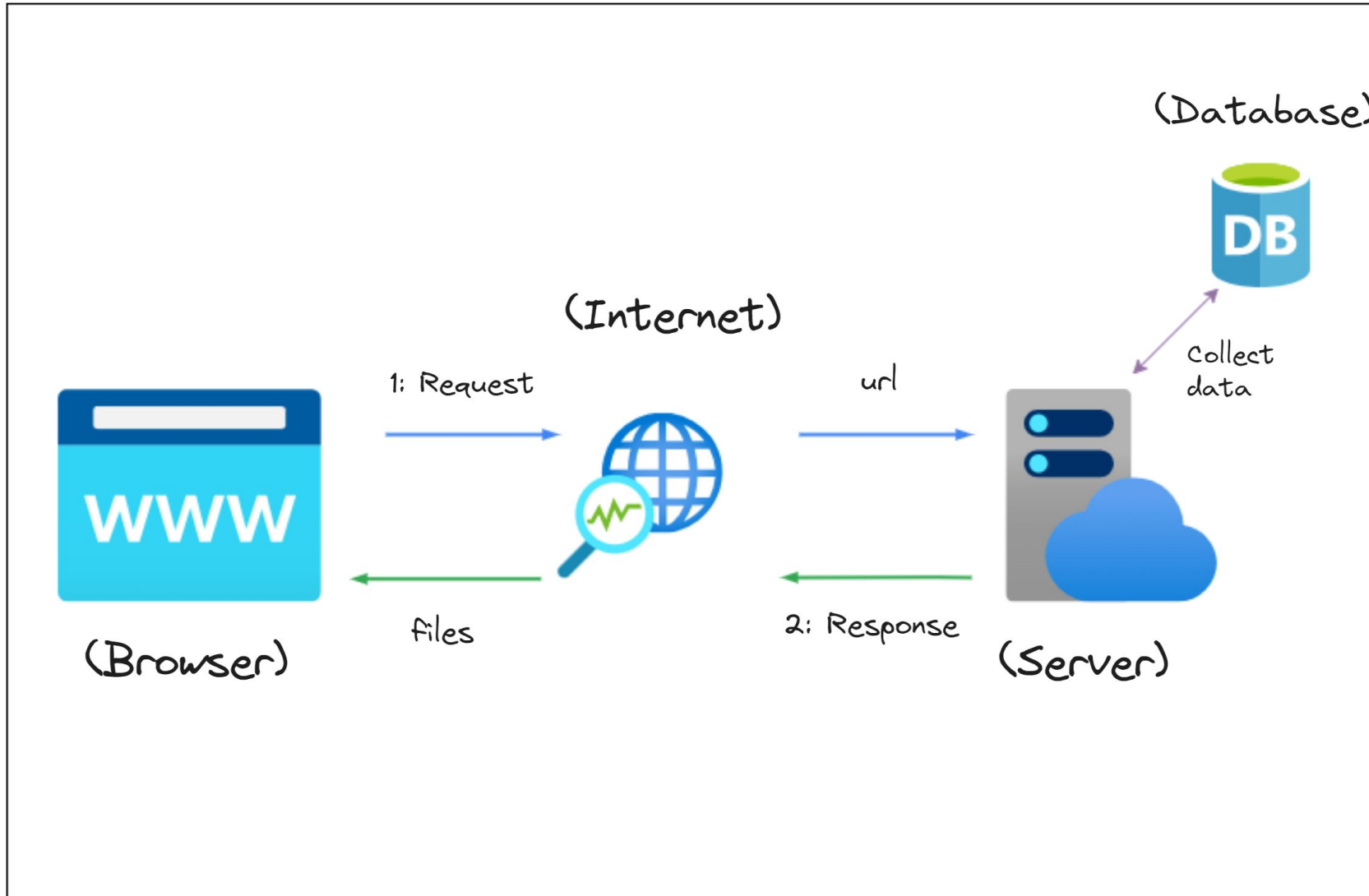Efficient Downloading

Polite Delay in
Multiple Requests

Process data to
create a Dataset

# Outline

- How can we scrape image urls for single class?
- How can we scrape image urls for multiple classes?
- Downloading images via urls with multi-threading
- Why we need polite delay?
- Clean and Organize our final dataset

# How can we scrape image urls for single class?

# Client – Server Protocol



(Database)

DB

(Internet)

1: Request

url

Collect
data

WWW

files

2: Response

(Browser)

(Server)

# Getting HTML from URL

The response of server

Request action

URL lead to website server

```
response = urllib.request.urlopen('https://aivietnam.edu.vn/')
response.status
```

200 → The status of response, 200 is mean successfully
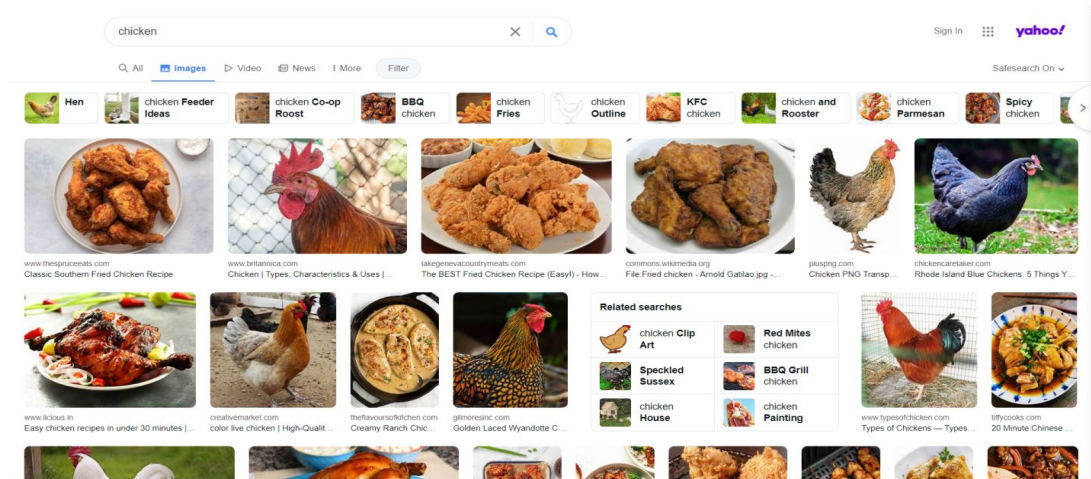
HTML source code of this website

Parse response into BeautifulSoup()

```
html_source = BeautifulSoup(response, 'html.parser')
html_source
```
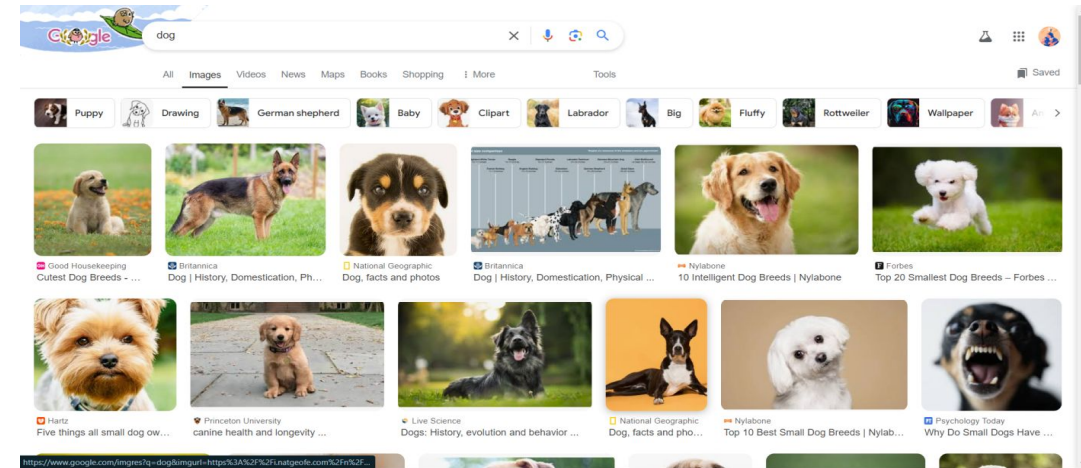
HTML file in plant text (string)

```
<head>
<title>AI Việt Nam - AI Việt Nam - Dịch vụ</title>
<meta content="initial-scale=1.0, width=device-width" name="v
<link href="https://fonts.googleapis.com/css?family=Open+Sans
```
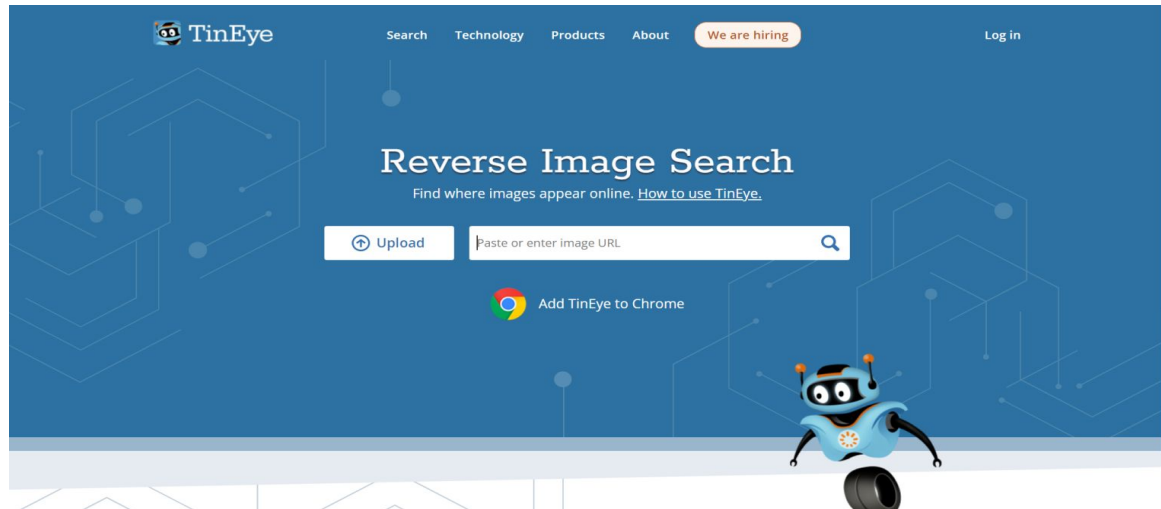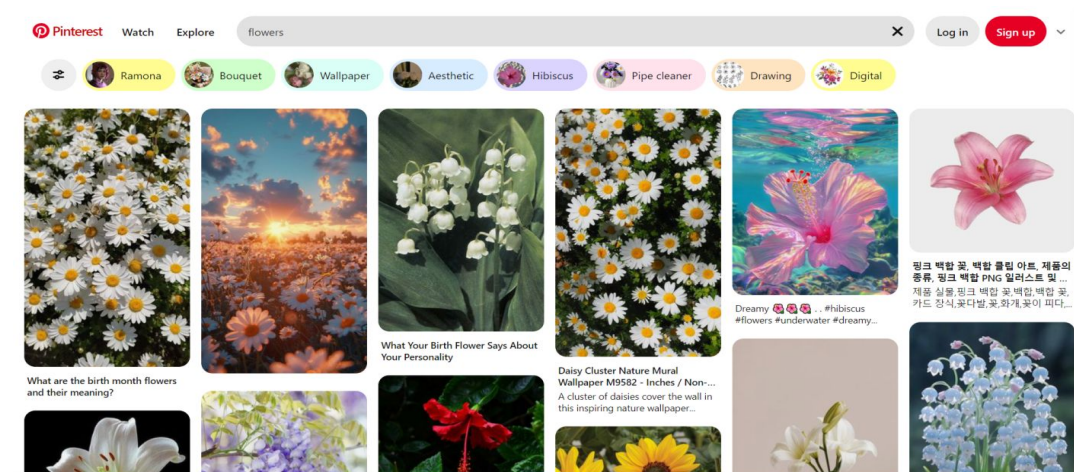
# Image search engine



Yahoo Image Search



Google Image Search



TinyEye



Pinterest

# Image search engine



Flickr.com

# Collect image's URL

URL route directly to the searching page of class

```
URL = "https://www.flickr.com/search/?text="
search_term = 'cat'
response = urllib.request.urlopen(URL+search_term)
html_source = BeautifulSoup(response, 'html.parser')
html_source.find_all("img")
```

Extract the \<img/> tag in html source

```
[<img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
 <img height="100%" loading="lazy" src="//live.staticflick
```

All img tags contain url of all images display in searching page

The image's url we can collect by **string manipulation**

Read detail manipulation in Colab (Brute force part)

# Quiz time

```
▶  # Manipulation code
   # ============================================================================================
   urls = []                                                                                  #=
   for img in img_tags:                                                                        #=
       if 'src' in img.attrs:                                                                  #=
           href = img.attrs['src']                                                             #=
           img_path = urljoin(URL, href)                                                       #=
           img_path = img_path.replace("_m.jpg", "_b.jpg").replace("_n.jpg", "_b.jpg").replace("_w.jpg", "_b.jpg") #=
           if img_path == "https://combo.staticflickr.com/ap/build/images/getty/IStock_corporate_logo.svg":        #=
               continue                                                                        #=
           urls.append(img_path)                                                               #=
   # ============================================================================================

   # Print number of urls already collected
   print(f"The total urls we collected = {len(urls)}")
```

```
⊒▾  The total urls we collected = 22
```

❖ Why we only collected  22 urls?

 Answer: …

❖ How can we collect more urls as we want?

 Answer: …

# Quiz time

```python
# Manipulation code
# ================================================================================
urls = []                                                                          #=
for img in img_tags:                                                               #=
    if 'src' in img.attrs:                                                         #=
        href = img.attrs['src']                                                    #=
        img_path = urljoin(URL, href)                                             #=
        img_path = img_path.replace("_m.jpg", "_b.jpg").replace("_n.jpg", "_b.jpg").replace("_w.jpg", "_b.jpg") #=
        if img_path == "https://combo.staticflickr.com/ap/build/images/getty/IStock_corporate_logo.svg":     #=
            continue                                                               #=
        urls.append(img_path)                                                      #=
# ================================================================================

# Print number of urls already collected
print(f"The total urls we collected = {len(urls)}")
```

```
The total urls we collected = 22
```
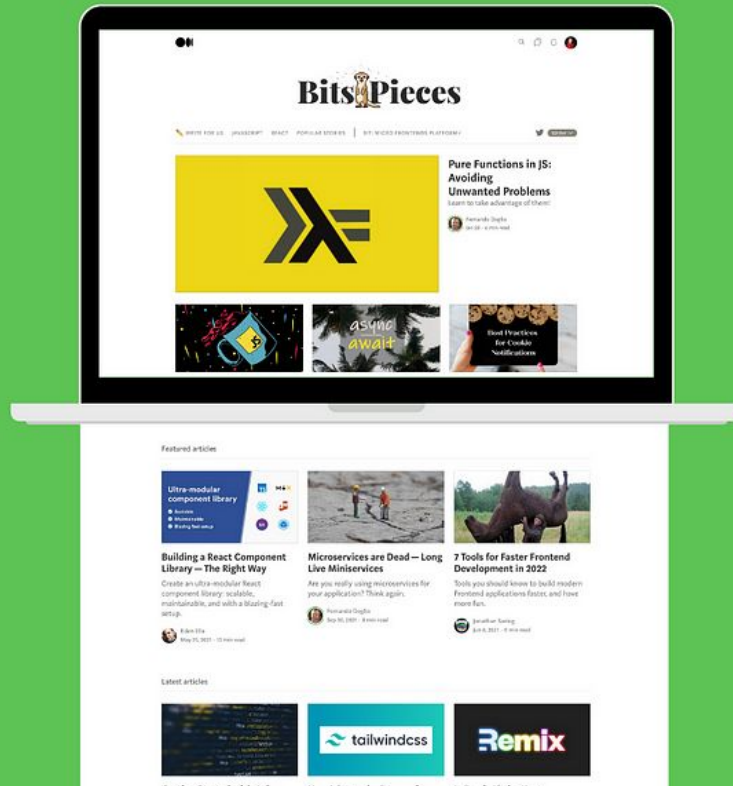
❖ Why we only collected  22 urls?

     Answer: Lazy loading and restrict loading more image.

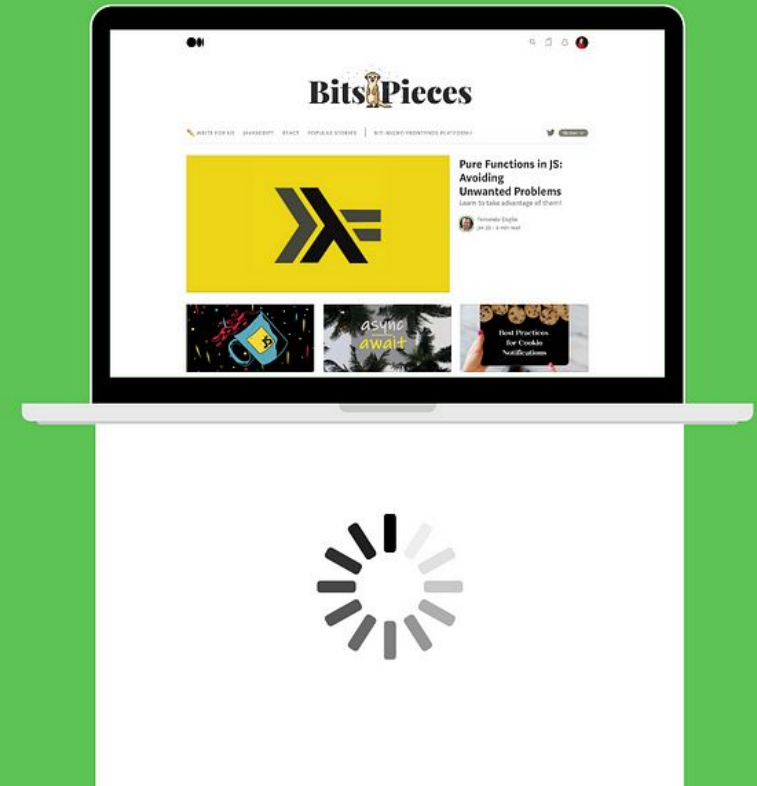❖ How can we collect more urls as we want?

     Answer: Use Selenium loading more content in searching page.
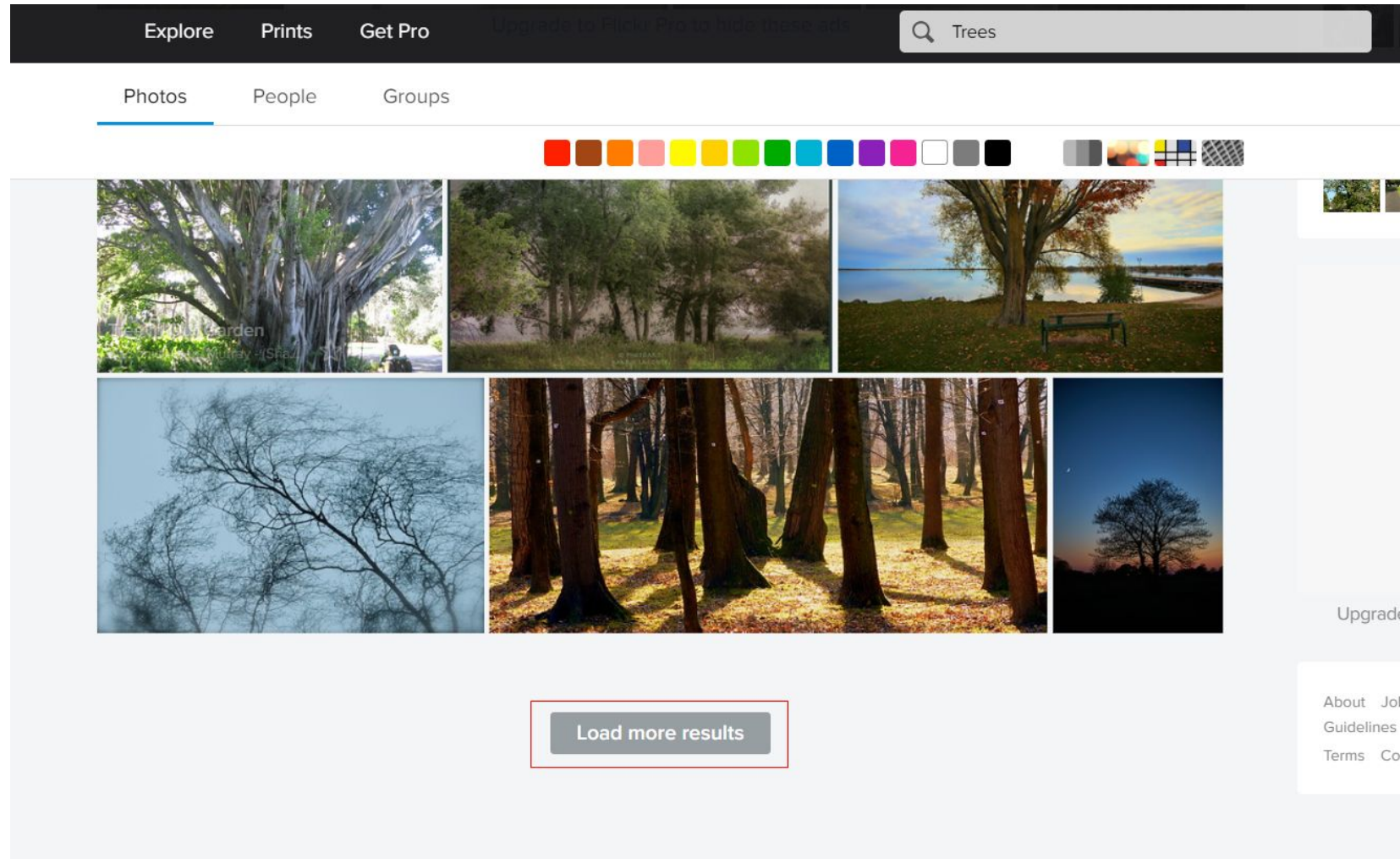
# Lazy loading



Without Lazy Loading

With Lazy Loading

Lazy loading is a design pattern commonly used in programming, especially in web development, to delay the loading of resources until they are actually needed.

# Load more result



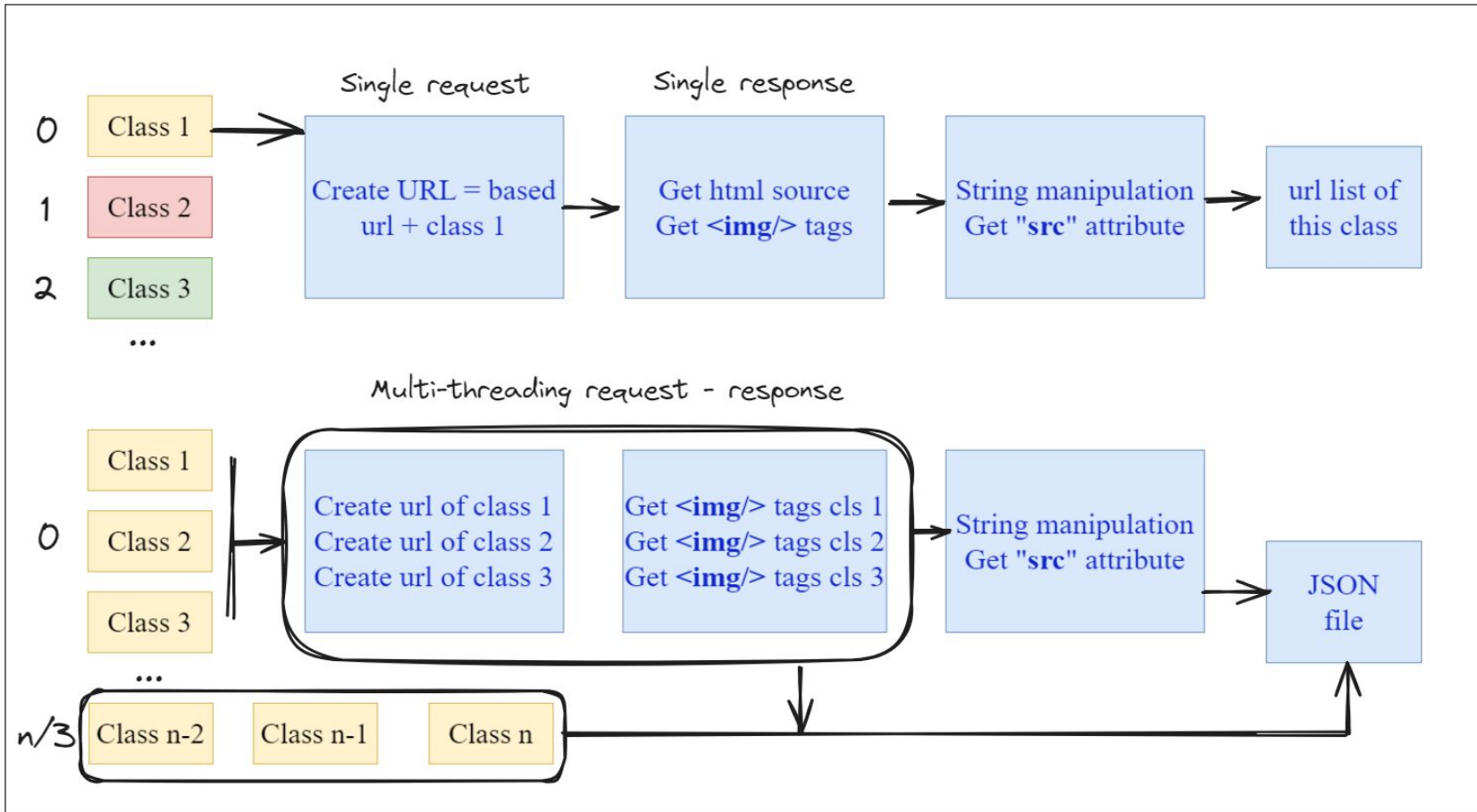We need to click this button in code view to load more content until we collect enough urls

# Handle with Selenium

```python
# Click load more button or scroll page for more image
try:
    load_more_button = WebDriverWait(driver, 10).until(
        EC.element_to_be_clickable((By.XPATH, '//button[@id="yui_3_16_0_1_1721642285931_28620"]'))
    )
    load_more_button.click()
    time.sleep(2) # Wait for generating content
except:
    driver.execute_script("window.scrollTo(0, document.body.scrollHeight);")
    time.sleep(2) # Wait for generating content

    # Check number of new generating image
    new_soup = BeautifulSoup(driver.page_source, "html.parser")
    new_img_tags = new_soup.find_all("img", loading_="lazy")
    if len(new_img_tags) == len(img_tags):
        more_content_available = False
    img_tags = new_img_tags
```

# How can we scrape image urls for multiple classes?

# Collect urls of all class



Single request

Single response

| 0 | Class 1 | → | Create URL = based url + class 1 | → | Get html source Get **<img/>** tags | → | String manipulation Get "**src**" attribute | → | url list of this class |
| 1 | Class 2 | | | | | | | |
| 2 | Class 3 | | | | | | | |

Multi-threading request - response

| 0 | Class 1, Class 2, Class 3 | → | Create url of class 1, Create url of class 2, Create url of class 3 | Get **<img/>** tags cls 1, Get **<img/>** tags cls 2, Get **<img/>** tags cls 3 | → | String manipulation Get "**src**" attribute | → | JSON file |

n/3 | Class n-2 | Class n-1 | Class n |

Combine process
in Class
UrlScraper

| **UrlScraper** |
| --- |
| - url_template: str |
| - max_images: int |
| - max_workers: int |
| + __init__(url_template, max_images, |
| + setup_environment() |
| + get_url_images(term: str) -> list |
| + scrape_images(categories: dict) -> dict |
| + save_to_file(data: dict, filename: str) -> None |

# Final JSON file

# Downloading images via urls with multi-threading

# Download 1 image

Dataset/animals

Dataset

animals

```python
category_dir = os.path.join(self.download_dir, category)
if not os.path.exists(category_dir):
    os.makedirs(category_dir)
```

scheme='https',

Dataset/animals/cat

```python
term_dir = os.path.join(category_dir, term)
if not os.path.exists(term_dir):
    os.makedirs(term_dir)
```

cat

netloc='live.staticflick r.com'

Dataset/animals/cat /7190755946_ea97 e85765_b.jpg

```python
filename = os.path.join(term_dir, os.path.basename(urlparse(url).path))

self.filename.add(filename)  # Record the filename directory

try:
    urllib.request.urlretrieve(url, filename)
    pbar.update(1)
    return f"Downloaded: {url}"
except Exception as e:
    pbar.update(1)
    return f"Failed to download {url}: {str(e)}"
```

path='/7073/71907559 46_ea97e85765_b.jpg'

params=''

query=''

fragment=''

url: "https://live.staticflickr.com/7313/9775005856_9b5e0ebe16_b.jpg",

**Download images with multi-threading**

# Quiz time

❖ What is difference between 100.000 request from many users and 1000 requests from one user? (in 5 seconds)

Multiple users

Single user

# Why we need polite delay?

# Server Overloading



Many users

simultaneously

(server)

One user

simultaneously

(server)

# Polite delay



❖ Polite delay in web scraping is the practice of adding a pause between consecutive requests to a website to avoid overloading the server and to respect the website's resources.

# Clean and Organize our final dataset

# Build final data folder

# Final dataset



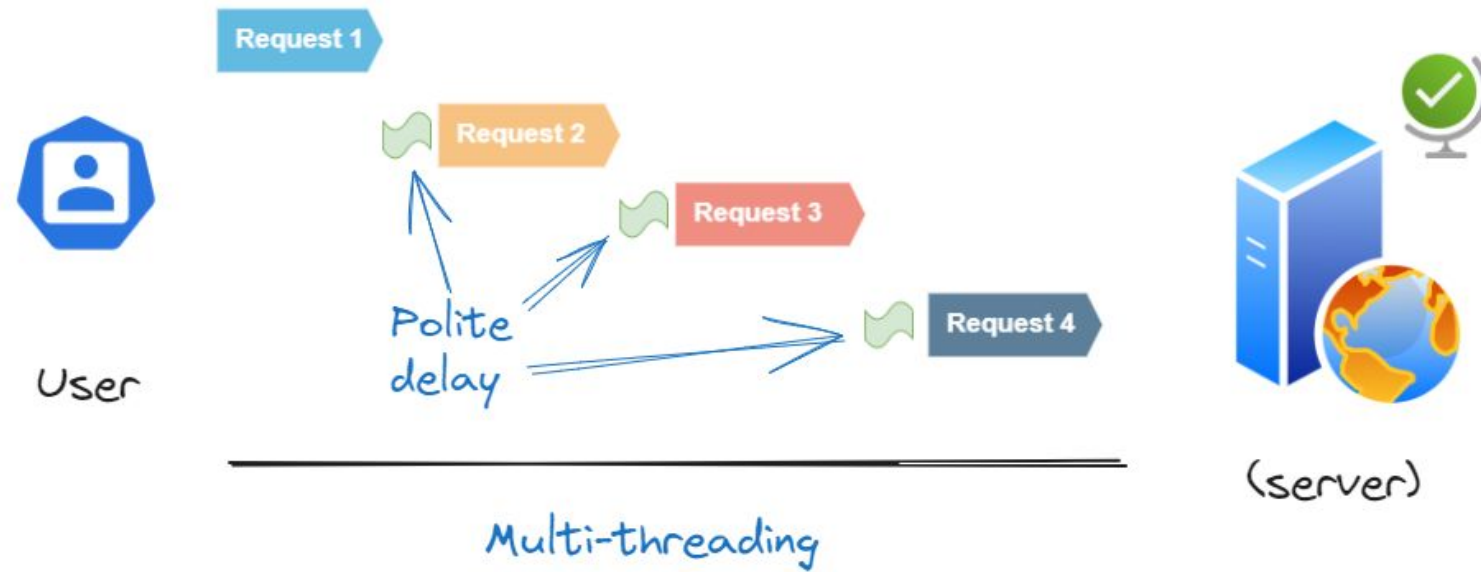| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Alligator | Alpaca | Anteater | Ants | Apple | Apricot | ball | Bamboo | Banana | Bats | Bean | bear |
| bed | Bee | bird | book | Bookcase | Bougainvillea | Bowl | Buffalo | Butterfly | cabinet | Camel | carrot |
| cart | Cat | Caterpillar | cave | chair | Cheetah | chests | Chicken | Civet | Cliff | clock | Clothes |
| cloud | Clove | Coast | Cock | Cockroaches | Coconut | Coffeeplant | Coral | Corn | Corneliantree | Cottonplant | cows |
| Crab | desert | desks | Dog | Dogwood | donkey | Dragonfly | Duck | Durian | Elephant | farmland | Fern |
| Ferns | Fig | Fish | flamingo | Flax | Flies | Flower | flowergarden | Forests | fox | Frangipani | fridge |