

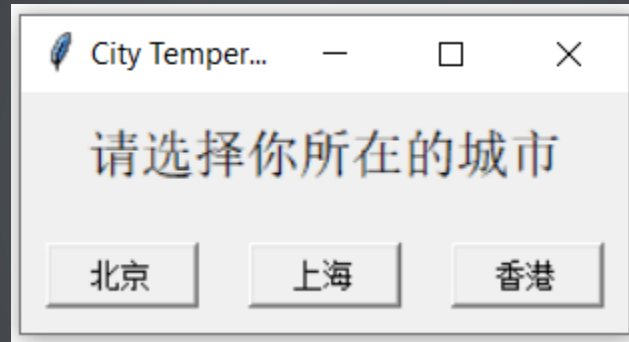
API



Demonstration and Vote!

This will be included as part of your in-class participation.

I developed an APP. Let's try it.



In your opinion, how does it get the temperature data?

This is done through an API.

API: application programming interface.

APIs are mechanisms that enable two software components to communicate with each other using a set of definitions and protocols.

You can view an API as a function (usually developed by a company), which helps you solve a task. You send your input to the function and it returns the output to you. You do not need to know how the function is computed.

This is done through an API.

You can view an API as a function (usually developed by a company), which helps you solve a task. You send your input to the function and it returns the output to you. You do not need to know how the function is computed. It is a black box.

In my example, there is a temperature API. I send the location to the API and it returns the temperature to me.

<https://www.youtube.com/embed/s7wmiS2mSXY?enablejsapi=1>

Let me show you the basics.

```
1 import requests
2 url = "http://apis.juhe.cn/simpleWeather/query"
3 params = {
4     "city": "北京",
5     "key": "your API key",
6 }
7 response = requests.get(url, params = params)
8 data = response.json()
9 temperature = data['result']['realtime']['temperature']
10 print(temperature)
```

You need to apply for your API key (password)

I will not share my key with you because a free key can only be used 50 times a day...

what is the request module in python

Completed 2 steps >

The `requests` module is a popular third-party Python library for sending HTTP requests to web servers easily and intuitively. It is **not** part of Python's standard library and must be installed separately with `pip install requests`. `realpython`

What It Does

`requests` abstracts the complexity of making HTTP connections behind simple, human-friendly function calls. It handles query strings, headers, authentication, cookies, and response parsing automatically. `pypi +1`

How to get the key?

1. Sign up an account [here](#) and activate it.
2. Sign up for an API key.

The screenshot shows a dashboard interface with a sidebar on the left and a main content area on the right. The sidebar contains navigation items: 我的聚合, 数据中心, 我的API, 我的数据块, 我的收藏, 安全中心, 账号管理, 工单管理, 费用管理, and 应用中心. The main content area is titled '我的API' and features a search bar with the text '已申请的接口名称/ID检索'. Below the search bar, there is a list of APIs. The first API is '天气预报' (Weather Forecast) with ID: 73 and a '免费' (Free) tag. It includes an AppKey: 38b0441a463012bee90bda232e0abedb and usage limits: 限每天50次请求, 开通「黑钻会员」最高每天100000次数. There is a '开通会员' (Upgrade to Member) button. At the bottom of the API entry, there are links for 统计 (Statistics), 测试 (Test), 文档 (Documentation), and 置顶 (Pin).

Exercise

Apply for your own API key and develop your own weather APP!

```

1 import tkinter as tk
2 import requests
3 def get_temperature(city):
4     url = "http://apis.juhe.cn/simpleWeather/query"
5     params = { "city": city, "key": "your key" }
6     response = requests.get(url, params=params)
7     data = response.json()
8     temperature = data['result']['realtime']['temperature']
9     spaces_to_add = 3 - len(temperature)
10    output_string = " " * spaces_to_add + temperature
11    return output_string
12 def update_display(city):
13    temperature = get_temperature(city)
14    display_label.config(text="{}当前温度为: {}摄氏度".format(city, temperature))
15 root = tk.Tk()
16 root.title("City Temperature App")
17 display_label = tk.Label(root, text="    请选择你所在的城市    ", font=("Arial", 16))
18 display_label.grid(row=0, column=0, columnspan=3, pady=10)
19 cities = ["北京", "上海", "香港"]
20 for index, city in enumerate(cities):
21    city_button = tk.Button(root, text=city, command=lambda c=city: update_display(c))
22    city_button.grid(row=1, column=index, padx=10, pady=10, sticky="ew")
23 for i in range(3):
24    root.grid_columnconfigure(i, weight=1)
25
26 root.grid_rowconfigure(1, weight=1)
27 root.mainloop()

```

Let's try the following [API](#) for currency exchange

If you don't understand the API, ask AI for help!

API can not only return data to you. It can also return images! Try the following [API](#) that returns a random dog image to you.

Try the following [API](#) for recognizing IP address.

AviationStack is a platform which provides an API to track the status of flights. Its website is here:

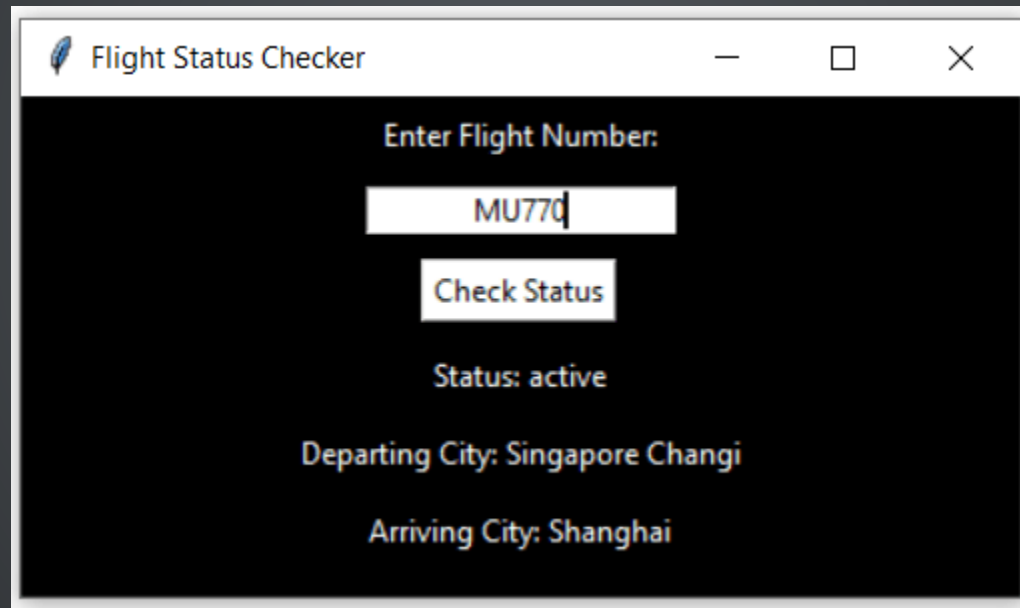
<https://aviationstack.com/>

You can find its API documentation [here](#).

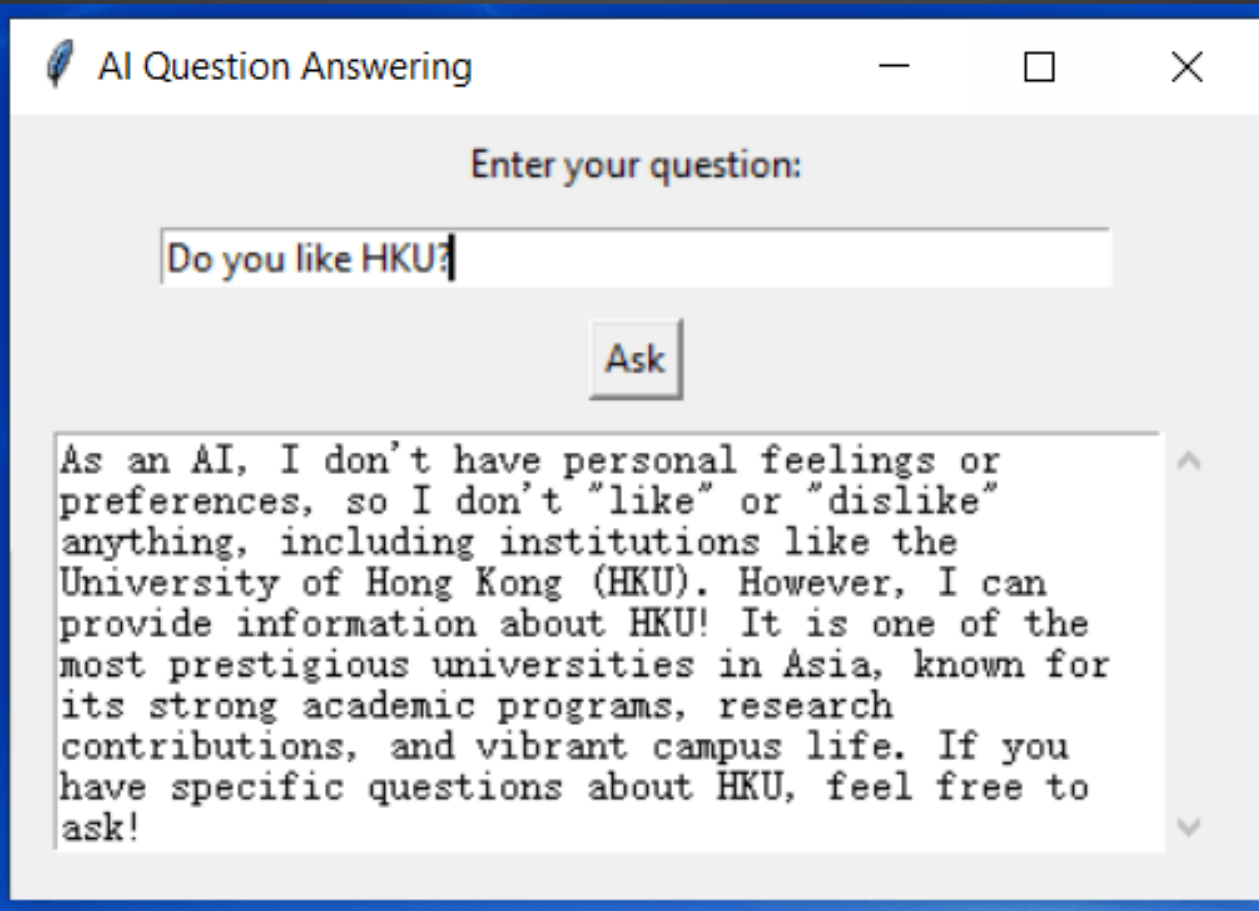
Exercise

- Acquire a free API Key from the platform
- Develop an APP, the APP takes flight number as input from the user, and displays its flight status, departing airport and arriving airport on the screen.
- **Try to do everything yourself (with help from AI)!**

This is my APP.



This is “my” AI (API)



Exercise

- Acquire an API Key from [DeepSeek](https://platform.deepseek.com/) (https://platform.deepseek.com/)
- This is not free, but very cheap. You need to top up \$2 or RMB 10 to use the API.
- Then, develop an APP which takes a question from user, sends the question to the API, and displays the answer from API to the user.

Top up

USD

CNY

Amount

\$2

\$5

\$10

\$20

\$50

\$100

\$500

Custom

Pricing

Total



\$2.12

Total excluding tax \$2.00 + VAT(6%) \$0.12 ⓘ

Off-Peak Discounts: DeepSeek-V3 with 50% off and DeepSeek-R1 with 75% off at off-peak hours (16:30-00:30 UTC daily). Optimize your workflow while enjoying these exclusive savings. [View Details](#)

API keys

Your API keys are listed below. The API key is only visible and can be copied once at creation. Save it securely. Do not share your API key with others, or expose it in the browser or other client-side code. In order to protect the security of your account, DeepSeek may also automatically disable any API key that we've found has leaked publicly. Usage of API keys created before April 25, 2024, was not tracked.

Name	Key	Created	Last used	
ximarketing	sk-16926*****2e3e	2025-02-26	2025-02-26	 

If you need help on how to use the API, check the document
here: <https://api-docs.deepseek.com/>

What are the main topics in modern economics?

Scraping made Simple

We will use Google Chrome as our browser.

Install one if you don't have it.

Scraping

Consider a simple **URL**:

<https://www.aeaweb.org/issues/785>

American Economic Review is an academic journal in economics (we don't care it anyway). Let's try to write a scraper to download the article titles from the page:

Decisions under Risk Are Decisions under Complexity

by Ryan Oprea

(pp. 3789–3811)

Curbing Leakage in Public Programs: Evidence from India's Direct Benefit Transfer Policy

by Prabhat Barnwal

(pp. 3812–46)

Aiming for the Goal: Contribution Dynamics of Crowdfunding

by Joyee Deb, Aniko Öry, and Kevin R. Williams

(pp. 3847–76)

Scraping

I am scraping article titles from <https://www.aeaweb.org/issues/785>. A title looks like Decisions under Risk Are Decisions under Complexity. Give me the python code

已完成 2 个步骤 >

To scrape the article titles from the provided American Economic Association (AEA) journal issue, you can use Python with the `requests` and `BeautifulSoup` libraries. The target webpage places article titles inside hyperlink tags that point to specific article IDs, which can be easily isolated .

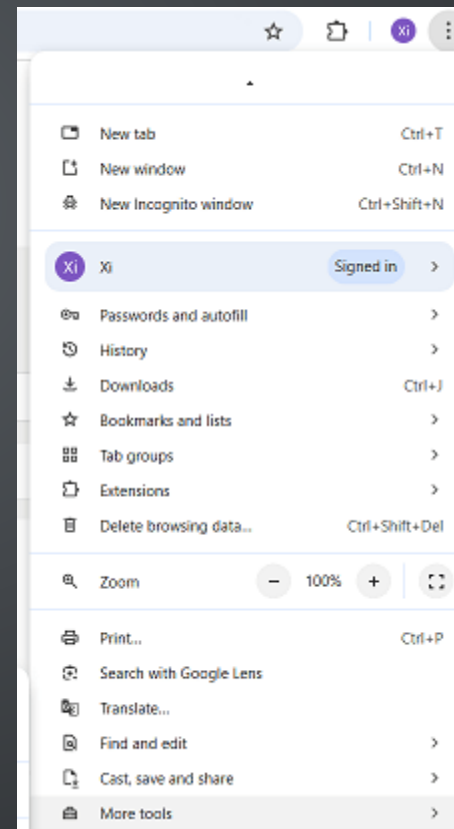
It sometimes works, but not always.

```
1 import requests
2 from bs4 import BeautifulSoup
3 # URL of the page to scrape
4 url = "https://www.aeaweb.org/issues/785"
5 # Send a GET request to fetch the page content
6 response = requests.get(url)
7 # Check if the request was successful (status code 200)
8 if response.status_code == 200:
9     # Parse the page content with BeautifulSoup
10    soup = BeautifulSoup(response.content, 'html.parser')
11
12    # Find all article titles
13    # Note: The exact HTML structure of the page needs to be inspected to
14    # find the correct tags and classes
15    titles = soup.find_all('h2', class_='title') # Example: Change these
16    # tags and classes based on actual HTML structure
17
18    # Print the titles
19    for title in titles:
20        print(title.get_text().strip())
21 else:
22    print(f"Failed to retrieve the page. Status code:
23    {response.status_code}")
```

Check HTML with Chrome

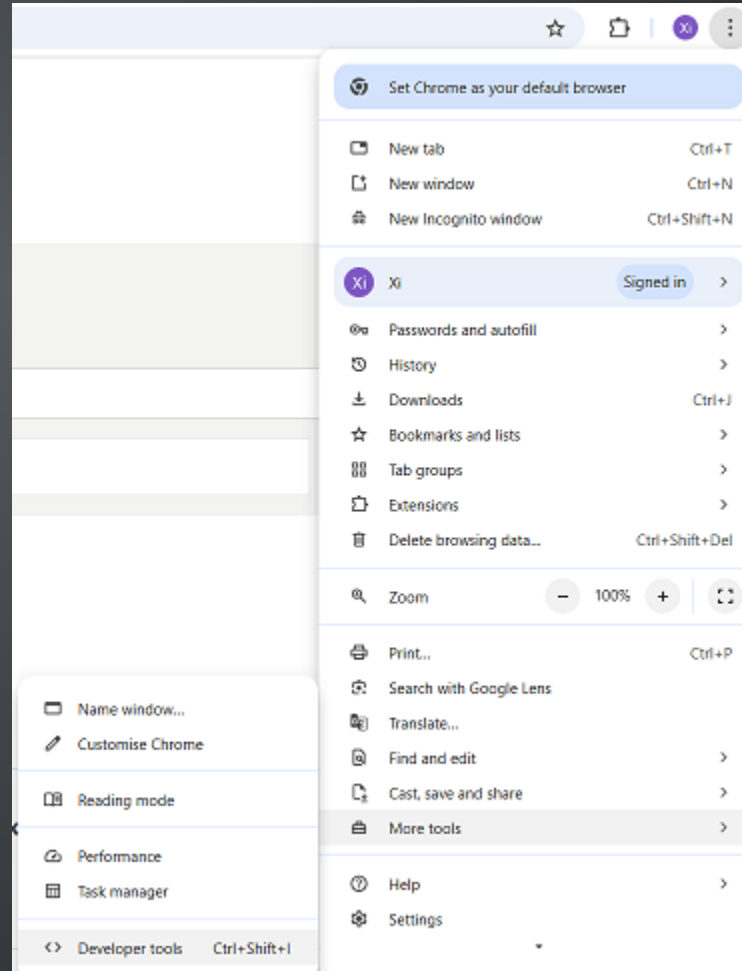
Open the webpage in your Chrome browser.

Click the upper right Chrome setting button of your browser and you will be directed here.




Check HTML with Chrome


Choose “More tools” ...
Choose “Developer tools” ...

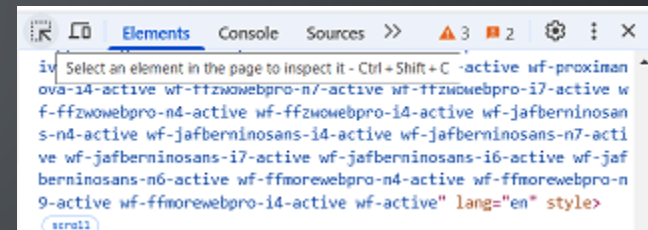


Check HTML with Chrome

Click the  button and you will get to “select an element in the page to inspect it”.

Alternatively, use “Ctrl + Shift + C.”

If needed, also click the  button to switch between desktop and mobile version.



Check HTML with Chrome

```
▼ <article class="journal-article art_18081" id="10.1257/aer.20221227">
  ▼ <h3 class="title">
    ... <a href="/articles?id=10.1257/aer.20221227">Decisions under Risk Are
      Decisions under Complexity</a> == $0
    </h3>
    ▶ <div class="article-item-authors">⋮</div>
      <div class="page-range">(pp. 3789–3811)</div>
    </article>
  ▶ <article class="journal-article art_18082" id="10.1257/aer.20161864">⋮
    </article>
  ▶ <article class="journal-article art_18083" id="10.1257/aer.20181851">⋮
    </article>
```

The title information is here.

Upload your screenshot to AI!

I am scraping article titles from <https://www.aeaweb.org/issues/785>. A title looks like Decisions under Risk Are Decisions under Complexity. Give me the python code



image.jpg

The code I got

```
1 import requests
2 from bs4 import BeautifulSoup
3 def scrape_aea_titles(url):
4     headers = {
5         "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/120.0.0.0 Safari/537.36"
6     }
7
8     response = requests.get(url, headers=headers)
9     if response.status_code != 200:
10        print(f"Failed to retrieve page. Status code:
{response.status_code}")
11        return
12    soup = BeautifulSoup(response.text, 'html.parser')
13    title_elements = soup.select('article.journal-article h3.title a')
14    print(f"Found {len(title_elements)} articles:\n")
15    for element in title_elements:
16        clean_title = ' '.join(element.text.split())
17        print(clean_title)
18 target_url = "https://www.aeaweb.org/issues/785"
19 scrape_aea_titles(target_url)
```

Scraping

Consider a simple URL:

<https://www.aeaweb.org/issues/785>

Try to download the links of the articles now, with the help of AI. Here is an example of the link:

<https://www.aeaweb.org/articles?id=10.1257/aer.20221227>

Scraping

I am scraping article links from <https://www.aeaweb.org/issues/785>. A link looks like <https://www.aeaweb.org/articles?id=10.1257/aer.20221227>. Give me the python code



image.jpg

My Code

```
1 import requests
2 from bs4 import BeautifulSoup
3 from urllib.parse import urljoin
4 def scrape_aea_article_links(url):
5     base_url = "https://www.aeaweb.org"
6     headers = {
7         "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/120.0.0.0 Safari/537.36"
8     }
9     response = requests.get(url, headers=headers)
10    response.raise_for_status()
11    soup = BeautifulSoup(response.text, 'html.parser')
12    article_links = []
13    for link_tag in soup.select("article.journal-article h3.title
a[href]"):
14        full_url = urljoin(base_url, link_tag['href'])
15        article_links.append(full_url)
16    return article_links
17 if __name__ == "__main__":
18    issue_url = "https://www.aeaweb.org/issues/785"
19    links = scrape_aea_article_links(issue_url)
20    print(f"Found {len(links)} links:\n")
21    for link in links:
22        print(link)
```

Exercise

- Go to <https://www.aeaweb.org/issues/785> to download the links to each article.
- Go to each of the above links to download the abstract of the article.
- Print the all abstracts on the screen.

Exercise

First, I asked AI to scrape the abstract for a given URL.

I am scraping abstract from URL like <https://www.aeaweb.org/articles?id=10.1257/aer.20221227>. Given the URL, output abstract like "We provide evidence that classic lottery anomalies..."



image.jpg

Exercise

I got the following code.

```
import requests
from bs4 import BeautifulSoup

def get_aea_abstract(url):
    response = requests.get(url)
    soup = BeautifulSoup(response.text, 'html.parser')
    abstract_section = soup.find('section', class_='article-information
abstract')
    if abstract_section:
        # The abstract text is directly inside the section, after the <h4> tag
        return abstract_section.get_text(strip=True).replace('Abstract', '',
1).strip()
    return None

url = "https://www.aeaweb.org/articles?id=10.1257/aer.20221227"
print(get_aea_abstract(url))
```

Exercise

Next, ask AI to merge the two codes into one.

I have two codes. Code 1 allows me to get URLs and Code 2 allows me to get abstract from each URL. Merge them into a single python code so that it prints each abstract of each URL one by one.

Code 1:

显示更多 ▾

Exercise

```
1 import requests
2 from bs4 import BeautifulSoup
3 from urllib.parse import urljoin
4
5 def get_aea_abstract(url):
6     headers = {
7         "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36"
8     }
9     response = requests.get(url, headers=headers)
10    soup = BeautifulSoup(response.text, "html.parser")
11    sec = soup.find("section", class_="article-information abstract")
12    if sec:
13        return sec.get_text(strip=True).replace("Abstract", "", 1).strip()
14    return "No abstract found"
15
16
17 issue_url = "https://www.aeaweb.org/issues/785"
18 base_url = "https://www.aeaweb.org"
19
20 response = requests.get(issue_url)
21 soup = BeautifulSoup(response.text, "html.parser")
22 links = [
23     urljoin(base_url, a["href"])
24     for a in soup.select("article.journal-article h3.title a[href]")
25 ]
26
27 for url in links:
28     print(url)
29     print(get_aea_abstract(url))
30     print("-" * 80)
```

We can do more!

- On page <https://www.aeaweb.org/journals/aer/issues>, we can scrape the URL of each issue.
- On the issue page, we can scrape the URL of each article.
- On the article page, we can scrape the abstract of each article (**it took me two hours**).
- Then, we can save all abstracts in one single file.
- After that, we can perform topic modeling!

We can do more!

I have already uploaded the txt file to the course website, which can be found [here](#):

http://ximarketing.github.io/data/AER_abstracts.txt

Please try the top words and LDA algorithm yourself and demonstrate the main topics in modern economics.

Exercise: Download the Photos of All HKU Business School
Professors from

[https://www.hkubs.hku.hk/people/faculty?
pg=1&staff_type=faculty&subject_area=all&track=professor
iate®ion=hku](https://www.hkubs.hku.hk/people/faculty?pg=1&staff_type=faculty&subject_area=all&track=professoriate®ion=hku)

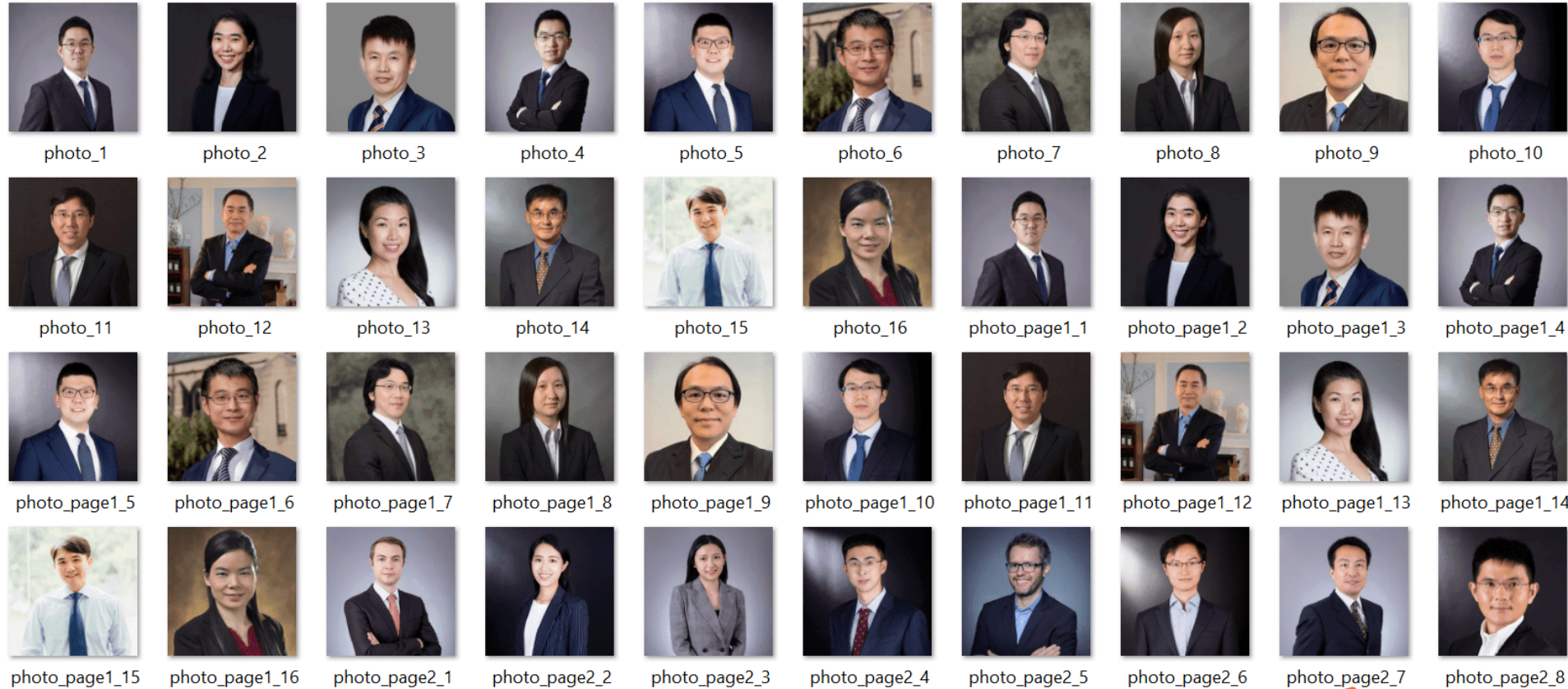
Downloading from a single page:

```
1 import requests
2 from bs4 import BeautifulSoup
3 import os
4 url = "https://www.hkubs.hku.hk/people/faculty?
pg=1&staff_type=faculty&subject_area=all&track=professoriate&region=hku"
5 response = requests.get(url)
6 if response.status_code == 200:
7     soup = BeautifulSoup(response.text, 'html.parser')
8
9     photo_divs = soup.find_all('div', class_='people-item')
10    photo_urls = []
11
12    for div in photo_divs:
13        img_tag = div.find('img')
14        if img_tag and 'src' in img_tag.attrs:
15            photo_url = img_tag['src']
16            photo_urls.append(photo_url)
17
18    os.makedirs('profile_photos', exist_ok=True)
19
20    for idx, photo_url in enumerate(photo_urls):
21        img_response = requests.get(photo_url)
22
23        if img_response.status_code == 200:
24            # Save the photo to the directory
25            with open(f'profile_photos/photo_{idx + 1}.jpg', 'wb') as f:
26                f.write(img_response.content)
27                print(f'Downloaded: photo_{idx + 1}.jpg')
28        else:
29            print(f'Failed to download photo from {photo_url}')
30 else:
31    print(f'Failed to retrieve the page. Status code: {response.status_code}')
```

I further asked the following question...

I have 10 URLs each containing 16 photos, and I want to download all these photos. The URLs are https://www.hkubs.hku.hk/people/faculty?pg=1&staff_type=faculty&subject_area=all&track=professoriate®ion=hku to https://www.hkubs.hku.hk/people/faculty?pg=10&staff_type=faculty&subject_area=all&track=professoriate®ion=hku

I am done!



The code I used

```
1 import requests
2 from bs4 import BeautifulSoup
3 import os
4 base_url = "https://www.hkubs.hku.hk/people/faculty?pg=
   {}&staff_type=faculty&subject_area=all&track=professoriate&region=hku"
5 os.makedirs('profile_photos', exist_ok=True)
6 for page in range(1, 11):
7     url = base_url.format(page)
8     print(f"Fetching page: {url}")
9     response = requests.get(url)
10    if response.status_code == 200:
11        soup = BeautifulSoup(response.text, 'html.parser')
12        photo_divs = soup.find_all('div', class_='people-item')
13        for idx, div in enumerate(photo_divs):
14            img_tag = div.find('img')
15            if img_tag and 'src' in img_tag.attrs:
16                photo_url = img_tag['src']
17                img_response = requests.get(photo_url)
18                if img_response.status_code == 200:
19                    photo_filename = f'profile_photos/photo_page{page}_{idx + 1}.jpg'
20                    with open(photo_filename, 'wb') as f:
21                        f.write(img_response.content)
22                    print(f'Downloaded: {photo_filename}')
23                else:
24                    print(f'Failed to download photo from {photo_url}')
25    else:
26        print(f'Failed to retrieve the page. Status code: {response.status_code}')
```

It becomes easier when I use more advanced LLMs: [link](#)

I would like to write Python codes to download photos from https://www.hkubs.hku.hk/people/faculty?pg=1&staff_type=faculty&subject_area=all&track=professoriate®ion=hku and save them in my local folder. Here is a sample link to photo: https://www.hkubs.hku.hk/wp-content/uploads/fly-images/196951/Andy_Back-scaled-940x1000-ct.jpg Moreover, when saving the photos, name them using the faculty name, e.g., Andy_Back.jpg

Try to scrape article titles from the following page:
<https://pubsonline.informs.org/toc/mnsc/71/1>

Try to scrape article titles from the following [page](https://pubsonline.informs.org/toc/mnsc/71/1):

<https://pubsonline.informs.org/toc/mnsc/71/1>

In many cases, the simple scraper does not work. Possible reasons:

- They adopt anti-scraping technologies to prevent you from scraping their websites.
- The website offers different content to human users and scrapers.
- The website is dynamic (e.g., it requires password)

What else can we do with scraping?

What else can we do with scraping?

- Scraping a website with selenium.
- Web scraping with concurrency (parallel scraping).
- Changing IP addresses with VPN (using proxies / VPNs).

<https://www.youtube.com/embed/cc26zFE8X1k?enablejsapi=1>

Scraping with Selenium

```
1 import undetected_chromedriver as uc
2 from selenium.webdriver.common.by import By
3 from selenium.webdriver.support.ui import WebDriverWait
4 from selenium.webdriver.support import expected_conditions as EC
5 url = "https://pubsonline.informs.org/toc/mnsc/71/1"
6 options = uc.ChromeOptions()
7 # options.add_argument("--headless=new") # keep headed; Cloudflare blocks
  headless harder
8
9 driver = uc.Chrome(options=options, version_main=147) # <-- pin to your
  Chrome major version
10
11 try:
12     driver.get(url)
13
14     WebDriverWait(driver, 30).until(
15         EC.presence_of_all_elements_located((By.CSS_SELECTOR, "h5.issue-
  item__title a")))
16     )
17
18     titles = driver.find_elements(By.CSS_SELECTOR, "h5.issue-item__title
  a")
19     for t in titles:
20         print(t.text.strip())
21 finally:
22     driver.quit()
```