# Using the Web

**ISYS5002, School of Marketing and Management**

Curtin University

# ELECTRONIC WARNING NOTICE FOR COPYRIGHT STATUTORY LICENCES

I acknowledge the traditional custodians of the land on which I work and live, and recognise their continuing connection to land, water and community. I pay respect to elders past, present and emerging.

Curtin University

# Today

- Accessing/Download Data

    3 methods in Python

- Create Webpage

- Web Scraping

Curtin University

# urllib.request.urlretrieve

```python
from urllib import request


# Define the remote file to retrieve

remote_url = 'https://www.google.com/robots.txt'


# Define the local filename to save data

local_file = 'local_copy.txt'


# Download remote and save locally

request.urlretrieve(remote_url, local_file)
```

Curtin University

# requests.get + manual save

```
import requests
# Define the remote file to retrieve remote_url =
'https://www.google.com/robots.txt'


# Define the local filename to save data

local_file = 'local_copy.txt'


# Make http request for remote file data

data = requests.get(remote_url)


# Save file data to local copy

with open(local_file, 'wb')as file:

    file.write(data.content)
```

Curtin University

# wget.download

```python
import wget

# Define the remote file to retrieve

remote_url = 'https://www.google.com/robots.txt'


# Define the local filename to save data

local_file = 'local_copy.txt'


# Make http request for remote file data
wget.download(remote_url, local_file)
```

Curtin University

# Website Development

- Django

   Fast, secure

- Flask

   minimalistic

Curtin University

```python
from flask import Flask

app = Flask(__name__)

@app.route("/")
def home():
    return "Hello, World!"

if __name__ == "__main__":
    app.run(debug=True)
```

Curtin University

# What is Web Scraping



**HTML Websites** → **Web Scraping** → **Data**
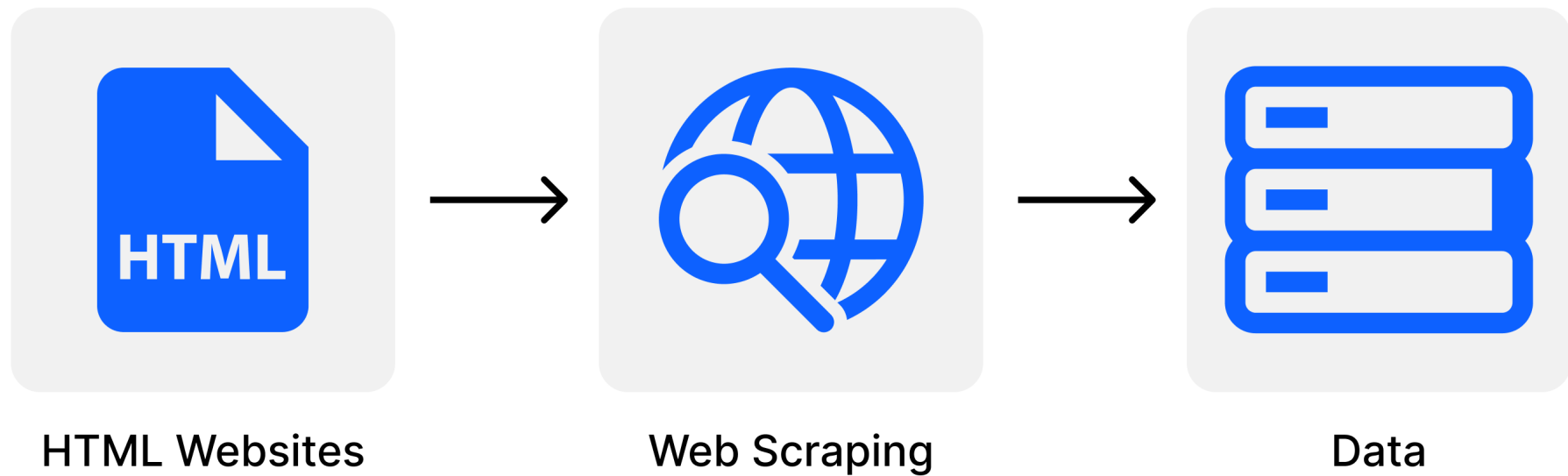
Image Source: https://i.imgur.com/6zM7JBq.png

Curtin University

# Why is Web Scraping Used?

- Price Comparison

- Email address gathering

- Social Media Scraping

- Research and Development

- Job Listings

Curtin University

# Is Web Scraping Legal?

- Copyright law protects original work

    Fair dealing

- Data not consider original work

- Read Website Terms of Service

- Check robots.txt

ISYSS5002 – Using the Web
Curtin University is a trademark of Curtin University of Technology
CRICOS Provider Code 00301J

Curtin University

# Is Python Good for Web Scraping

- Ease of Use

- Large Collection of Libraries

- Dynamically Typed

- Small code large Task

- Community

Curtin University

# How Scrape

- Find URL

- Inspect the Page

- Find the data want to extract

- Write the code

- Run code and extract the data

- Store data in the required format

Curtin University

# Libraries

- Selenium – automate browser activities

- BeautifulSoup – parsing HTML and XML documents

- Requests – make HTTP request simpler

ISYSS5002 – Using the Web
Curtin University is a trademark of Curtin University of Technology
CRICOS Provider Code 00301J

Curtin University

# Anatomy of URL

- Scheme

  **https**://developer.mozilla.org

- Domain

  https://**developer.mozilla.org**/en-US/docs/Learn/

- Path

  https://developer.mozilla.org**/en-US/docs/Learn/**

- Parameters

  https://developer.mozilla.org/en-US/**search?q=URL**

Curtin University

# HTML Tags

```
<!DOCTYPE html>
<html>
<head>
<title>Page Title</title>
</head>
<body>

<h1>This is a Heading</h1>
<p>This is a paragraph.</p>

</body>
</html>
```

Curtin University

# Can you

- Explain Download Data

- Understand Create Webpage

- Describe Web Scraping

Curtin University