**Introduction**

This document provides a detailed explanation of the automation architecture implemented for our capstone project. The system was designed to ensure the web scraping pipelines run automatically without requiring constant manual intervention. All automation is currently hosted on an AWS EC2 instance, which was created under my personal AWS Free Tier account.

To enable continuity beyond my account and ensure that future students can continue this work, I have created a public Amazon Machine Image (AMI) of the configured environment. This allows anyone with an AWS account to launch an identical EC2 instance without needing to repeat the initial setup. The following sections outline the steps required to use the AMI, connect to the instance, understand the configuration, and manage the automated workflows.

**AMI Access**

To avoid loss of work when my Free Tier expires, I created a public AMI from the configured EC2 instance. Future students can use this AMI to replicate the environment.

- **AMI ID:** ami-07a95230fff01681f

- **Region:** ap-southeast-2 (Asia Pacific – Sydney)

- **Direct Console Link:** Launch this AMI

**Instructions:**

1. Log in to your AWS account.

2. Open the link above (or search for the AMI ID in EC2 → AMIs).

3. Click **Launch Instance**.

4. Choose an instance type (e.g., t2.micro for testing).

5. Configure networking and storage, then launch the instance.

**Connecting to the Instance**

Once the instance is launched, connect via SSH using the .pem key pair you generated during setup.

1. Allow SSH access to your IP by adjusting security group settings in the AWS console. (Check your current IP at whatismyip.com).

2. Note the **Public IPv4 address** of your running instance.

3. Open PowerShell on your local machine and navigate to the directory where your .pem file is stored.

4. Connect using the following command:

5. ssh -i your-key.pem ubuntu@<public-ip-address>

6. When prompted with *"Are you sure you want to continue connecting?"*, type yes and press Enter.

You will now be connected to the EC2 instance.

## Project Files and Scripts

Inside the EC2 instance:

- Running ls will display the available project directories.

- The main scraping logic is stored inside the **scraper-runner** directory.

Within this directory, the script **run_all_scrapers.sh** controls the automation. This script handles:

- Logging

- Exporting environment paths

- Activating the Python virtual environment

- Running the scraping jobs

- Automatically shutting down the instance after completion to save costs

If modifications are needed (e.g., adding scrapers, changing paths), open this script with:

```
 ▣  ubuntu@ip-172-31-20-121: ~        ×    +   ∨                                               —   ▢   ✕

 * Documentation:  https://help.ubuntu.com
 * Management:     https://landscape.canonical.com
 * Support:        https://ubuntu.com/pro

 System information as of Wed Sep 24 08:00:34 UTC 2025

  System load:  0.0              Temperature:             -273.1 C
  Usage of /:   73.6% of 8.65GB  Processes:               117
  Memory usage: 14%              Users logged in:         0
  Swap usage:   0%               IPv4 address for ens5:   172.31.20.121

 * Ubuntu Pro delivers the most comprehensive open source security and
   compliance features.

   https://ubuntu.com/aws/pro

 Expanded Security Maintenance for Applications is not enabled.

 17 updates can be applied immediately.
 To see these additional updates run: apt list --upgradable

 Enable ESM Apps to receive additional future security updates.
 See https://ubuntu.com/esm or run: sudo pro status

 Last login: Wed Sep 24 07:57:34 2025 from 43.227.109.80
 ubuntu@ip-172-31-20-121:~$ ls
 DiscountMate_new  aws  awscliv2.zip  requirements.txt  scrape-env  scraper-runner  snap
 ubuntu@ip-172-31-20-121:~$
```

Our scraping code lives here which we pulled from remote GitHub repo:

```
 ubuntu@ip-172-31-20-121:~$ cd DiscountMate_new/Scrapping
 ubuntu@ip-172-31-20-121:~/DiscountMate_new/Scrapping$ ls
 'Aldi Scraping From a Catalogue and Item Classification (Attempt 1)'   Drakes_2025-09-17_16-08-53.csv
 'Aldi Scraping from price-reductions Page'                             Drakes_2025-09-17_16-08-53.json
  Australia_GroceriesScraper                                           'Foodland Scraper'
  Drake                                                                 README.md
  Drakes_2025-09-02_16-08-45.csv                                        __pycache__
  Drakes_2025-09-09_16-08-46.csv                                        chrome-user-data-IGA
  Drakes_2025-09-11_07-34-50.csv                                        costco_scrapeing_tool
  Drakes_2025-09-11_07-34-50.json                                       db-config.json
  Drakes_2025-09-11_08-10-44.csv                                        scrape-env
  Drakes_2025-09-11_08-10-44.json                                       scraper_IGA_catalogue.py
  Drakes_2025-09-11_08-35-07.csv                                        scraper_IGA_specials.py
  Drakes_2025-09-11_08-35-08.json                                       scraper_adelaidesfinest.py
  Drakes_2025-09-12_16-07-51.csv                                        scraper_drake.py
  Drakes_2025-09-12_16-07-51.json                                       scraper_foodland.py
  Drakes_2025-09-12_23-47-34.csv                                        test-read-data.py
  Drakes_2025-09-12_23-47-34.json                                       test-write-data.py
  Drakes_2025-09-14_00-37-32.csv                                        utils.py
  Drakes_2025-09-14_00-37-32.json
 ubuntu@ip-172-31-20-121:~/DiscountMate_new/Scrapping$
```

Now, if you go inside scraper-runner and type ls, you will see a shell script file called **run_all_scrapers.sh**

```
 ubuntu@ip-172-31-20-121:~/DiscountMate_new$ cd ..
 ubuntu@ip-172-31-20-121:~$ ls
 DiscountMate_new  aws  awscliv2.zip  requirements.txt  scrape-env  scraper-runner  snap
 ubuntu@ip-172-31-20-121:~$ cd scraper-runner
 ubuntu@ip-172-31-20-121:~/scraper-runner$ ls
 logs  run_all_scrapers.sh  venv
 ubuntu@ip-172-31-20-121:~/scraper-runner$
```

Type nano **run_all_scrapers.sh** and it will open the script that we use to run our all scrapers.

```bash
GNU nano 7.2                                          run_all_scrapers.sh
#!/bin/bash

LOG_FILE="/home/ubuntu/scraper-runner/logs/cron_log.txt"
echo "==== Cron started at $(date) ====" >> "$LOG_FILE"

# Set PATH in case cron doesn't have it
export PATH="/usr/bin:/bin:/usr/local/bin:/usr/sbin:/home/ubuntu/.local/bin"

# Activate virtual environment if needed
source /home/ubuntu/scraper-runner/venv/bin/activate

# Sample scraper commands (update to actual)
cd /home/ubuntu/DiscountMate_new/Scrapping || {
  echo "Failed to cd into /home/ubuntu/DiscountMate_new/Scrapping" >> "$LOG_FILE"
  exit 1
}

echo "Inside scraper directory, starting Drake scraping" >> "$LOG_FILE"

# Run the actual scraping (modify this as per your actual scripts)
python3 -u scraper_drake.py >> "$LOG_FILE" 2>&1 || echo " Drakes scraping failed" >> "$LOG_FILE"

echo "✅ Drakes Scraper Completed" >> "$LOG_FILE"

# Run IGA Scraper
echo "== Starting IGA Scraper ==" >> "$LOG_FILE"

                                          [ Read 47 lines ]
^G Help      ^O Write Out   ^W Where Is   ^K Cut      ^T Execute   ^C Location    M-U Undo   M-A Set Mark
^X Exit      ^R Read File   ^\ Replace    ^U Paste    ^J Justify   ^/ Go To Line  M-E Redo   M-6 Copy
```

If you need to add other scrapers or export path or anything, this is the place.

In this file, we configure logs, export path, set up the scraping virtual environment, give path and command to run scraping, and include command at last to auto stop the EC2 instance after the scraping is completed so that we can save costs.

If you go inside logs directory inside scraper_runner, you will see log files.

Our logs are in the cron_log.txt file. Type cat cron_log.txt and you will see the logs.

This will help you debug whether the scraper ran successfully and if they encountered any error.

**Cron Job Automation**

The EC2 instance is configured with a cron job that automatically runs run_all_scrapers.sh at the scheduled time and writes logs to cron_log.txt.

To view or edit the cron job:

**crontab -e**

You will see the cron job set up to automatically run the run_all_scrapers.sh file and send logs to cron_log.txt



This is the part on the EC2 instance. But, to make sure this cron runs successfully on this instance, we need to make sure that the EC2 instance is up and running at least 15 minutes before the above scheduled cron job.

**EventBridge Integration**

Since the cron job requires the EC2 instance to be running beforehand, Amazon EventBridge is used to automatically **start the EC2 instance 15 minutes before the scheduled cron job**.

**Setup Summary:**

1. Create an IAM role with permissions for EC2 start/describe actions.

2. Attach a policy with the following JSON:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "ec2:StartInstances",
        "ec2:DescribeInstances"
      ],
      "Resource": "arn:aws:ec2:ap-southeast-2:718975141381:instance/i-0a0483386e8bf25e6"
    }
  ]
}
```
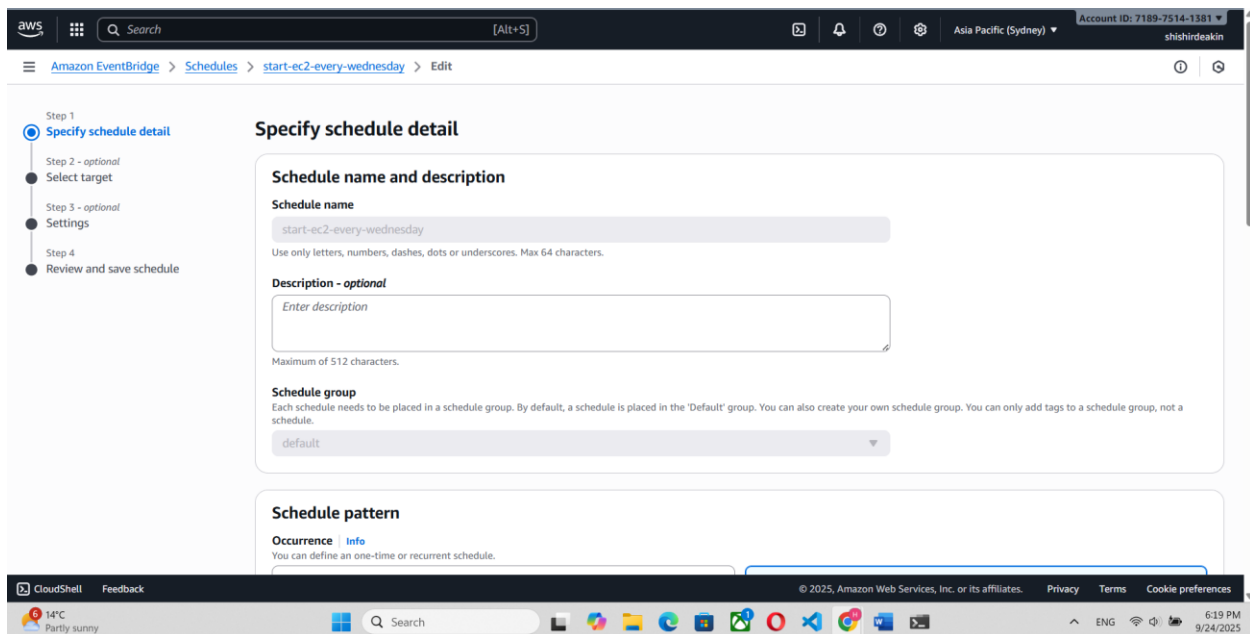
3. Assign this IAM role to the EventBridge schedule.

4. Create a schedule that triggers 15 minutes before the cron job.

5. Configure EventBridge to target the EC2 instance and execute **StartInstances**.

6. In addition, a separate IAM role with **AmazonEC2FullAccess** policy is attached to the EC2 instance, enabling it to stop itself once the scrapers finish execution.

The detailed steps with screenshot are provided below:

You can see I have this EventBridge rule currently on AWS that aligns with our cron job above.



So, you need to create a similar EventBridge rule on your AWS account.

But note that the EventBridge should have permission to auto start the EC2 instance. So, you should first create the required IAM roles and attach this so that it has permission to start instance.

So, go to your IAM console. And create role.

You need to create role like below

And attach Policy to it

Go to Policy and create policy with the following JSON:

```json
{

  "Version": "2012-10-17",

  "Statement": [

    {

      "Effect": "Allow",

      "Action": [

        "ec2:StartInstances",

        "ec2:DescribeInstances"

      ],

      "Resource": "arn:aws:ec2:ap-southeast-2:718975141381:instance/i-0a0483386e8bf25e6"

    }

  ]

}
```

Then you have to attach this policy to the EventBridge role you created like below:



Once that is done, come to Amazon EventBridge Schedule and create Schedule

Specify schedule name and details

Create cron pattern like below:



In Target, find Amazon EC2 and find **StartInstances**

Then provide your instance-id like below:



Next, in the Permissions, select Use existing role and select the permission you created above.

Finally, save the schedule.

Now, you have to create another IAM role that allows EC2 instance to stop it by itself once the scraping is complete.

In this role, attach AmazonEC2FullAccess policy.

**Conclusion**

With the AMI, cron job, and EventBridge integration in place, this automation framework ensures the scrapers run reliably on a weekly basis with minimal manual intervention. All results are stored in MongoDB collections for further processing. Future students can continue building on this foundation by extending the scraping logic or improving automation while reusing the existing infrastructure.