

HTML Content Scraper Development Guide

IODA Team

July 16, 2025

Contents

1	Core Architecture	3
1.1	Class Structure	3
1.2	Before You Start - Information to Collect	3
1.3	Critical Architectural Decisions	3
1.3.1	Processing Strategy	3
1.3.2	Error Handling Strategy	3
1.3.3	Directory Structure	4
2	HTML Parsing Patterns	4
2.1	Pattern 1: Table-Based Data	4
2.2	Pattern 2: Container-Based Data	4
2.3	Pattern 3: Mixed Content	4
3	Processor Requirements	5
3.1	Processor Structure	5
3.2	Command-Line Processing	5
3.3	Standard Output Format	5
3.3.1	DateTime Format Rules	5
4	Configuration Requirements	6
4.1	Ukraine-Specific Configuration	6
4.2	Regional Execution Differences	6
4.2.1	India, Nigeria, Pakistan	6
4.2.2	Ukraine	6
5	File Naming Standards	6
5.1	Actual Pattern in Use	6
5.2	Country Codes	6
6	Real-World Handling	7
6.1	Sites with No Data	7
6.2	Sites with Irregular HTML	7
7	Integration Requirements	7
7.1	Adding to daily_scraping.py	7
7.2	Adding to monthly_scraping.py	7
7.3	Adding to process_file.py	7

8	Development Process	8
8.1	Analysis Phase	8
8.2	Implementation Order	8
8.3	Testing Requirements	8
9	Success Criteria	8
10	Additional Resources	9
10.1	Repository Structure Overview	9

1 Core Architecture

1.1 Class Structure

```
1 class ProviderName:
2     def __init__(self):
3         # URLs, dates, paths
4
5     def check_folder(self, type):
6         # Directory creation (keep this name - it's used everywhere)
7
8     def scrape(self):
9         # Main entry point - does everything
```

1.2 Before You Start - Information to Collect

1. **Target URL** – The exact webpage containing outage data
2. **Data Location** – How data is structured (tables, divs, lists)
3. **Update Schedule** – How often the site updates (daily, weekly, etc.)
4. **Data Fields** – What information is available (dates, times, areas, etc.)
5. **Site Quirks** – Any special requirements (SSL issues, headers, etc.)

1.3 Critical Architectural Decisions

1.3.1 Processing Strategy

Always save raw HTML files AND process immediately. Both steps are required.

```
1 def scrape(self):
2     # Fetch data
3     response = requests.get(self.url)
4
5     # Save raw file - ALWAYS KEEP THIS
6     with open(file_path, "w", encoding="utf-8") as file:
7         file.write(response.text)
8
9     # Process immediately - ALWAYS DO THIS TOO
10    processor = Process_Provider(self.year, self.month, self.today, file_path)
11    processor.run()
```

1.3.2 Error Handling Strategy

```
1 # Use this pattern consistently
2 response = requests.get(self.url)
3 if response.status_code == 200:
4     # Process
5 else:
6     print(f"Failed to retrieve webpage. Status code: {response.status_code}")
7     return # Don't crash, just exit gracefully
```

1.3.3 Directory Structure

Non-negotiable pattern:

```
{country}/{provider}/raw/{year}/{month}/  
{country}/{provider}/processed/{year}/{month}/
```

Implementation:

```
1 def check_folder(self, type):  
2     self.folder_path = f"./{country}/{provider}/" + type + "/" + self.year + "/"  
3     " + self.month  
4     os.makedirs(self.folder_path, exist_ok=True)
```

Keep this exact signature - it's used across all scrapers.

2 HTML Parsing Patterns

2.1 Pattern 1: Table-Based Data

Used by sites with structured HTML tables:

```
1 def parse_table_data(self, soup):  
2     table = soup.find('table') # or find_all if multiple tables  
3     if not table:  
4         return []  
5  
6     rows = table.find_all('tr')  
7     data = []  
8     for row in rows[1:]: # Skip header  
9         cols = row.find_all('td')  
10        if len(cols) >= expected_columns:  
11            # Extract data from columns  
12            data.append(extracted_item)  
13    return data
```

2.2 Pattern 2: Container-Based Data

Used by sites with div containers or card layouts:

```
1 def parse_container_data(self, soup):  
2     containers = soup.find_all('div', class_='outage-item') # Adjust selector  
3     data = []  
4     for container in containers:  
5         # Extract from each container  
6         title = container.find('h3')  
7         if title: # Always check if elements exist  
8             # Process  
9     return data
```

2.3 Pattern 3: Mixed Content

Used by sites with irregular structure:

```
1 def parse_mixed_content(self, soup):  
2     # Find all potential data containers  
3     # Use flexible selectors  
4     # Validate each piece of data  
5     # Handle missing elements gracefully
```

3 Processor Requirements

3.1 Processor Structure

Every HTML scraper will have a processor class:

```

1 class Process_Provider:
2     def __init__(self, year, month, today, file_path):
3         # Standard signature - don't change this
4
5     def check_folder(self, type):
6         # Directory management for processed data
7
8     def get_data(self):
9         # Parse the HTML file
10
11    def save_json(self, data):
12        # Save to standard JSON format
13
14    def run(self):
15        # Orchestrate everything

```

3.2 Command-Line Processing

The process_file.py script supports batch processing:

```

python process_file.py \
    --country india \
    --provider npp \
    --start_date 2025-04-18 \
    --end_date 2025-04-22

```

For manual processor execution (mainly Ukraine):

```

1 # Modify file path and folder path in processor
2 file = "path/to/raw/file.html"
3 self.folder_path = "path/to/processed/output/"

```

3.3 Standard Output Format

Format:

```

1 {
2     "country": "Country Name",
3     "start": "YYYY-MM-DD_HH-MM-SS",
4     "end": "YYYY-MM-DD_HH-MM-SS",
5     "duration_(hours)": 2.5,
6     "event_category": "planned maintenance",
7     "area_affected": {
8         "Region/State": ["Area1", "Area2", "Area3"]
9     }
10 }

```

3.3.1 DateTime Format Rules

- Always use ISO date format for the date part
- Always use underscore separator: YYYY-MM-DD_HH-MM-SS
- Use 24-hour time format
- Duration in decimal hours (2.5 = 2 hours 30 minutes)

4 Configuration Requirements

4.1 Ukraine-Specific Configuration

For Ukraine providers, modify constants.py:

```
1 # Set root directory for data storage
2 root_dir = "/path/to/your/data/directory"
3 year = "2025"
4 month = "01"
5 date = "01"
```

Important: You must change the "root_dir" in constants.py to where you plan to save data.

4.2 Regional Execution Differences

4.2.1 India, Nigeria, Pakistan

- Processors are automatically called when running scraper files
- Scrapers can be directly executed
- Integrated processing approach

4.2.2 Ukraine

- Scrapers and processors run separately
- Uses constants.py for configuration
- Post-processor must be executed manually after scraping
- Oblast-based directory structure

5 File Naming Standards

5.1 Actual Pattern in Use

```
power_outages.{COUNTRY_CODE}.{provider}.raw.{YYYY-MM-DD}.html
power_outages.{COUNTRY_CODE}.{provider}.processed.{YYYY-MM-DD}.json
```

5.2 Country Codes

Country codes that actually exist:

- IND (India)
- NG (Nigeria)
- PK (Pakistan)
- UA (Ukraine)
- CM (Cameroon)

6 Real-World Handling

6.1 Sites with No Data

```

1 # Always handle this case
2 if not table or len(rows) == 0:
3     print(f"No outage data found for {self.today}")
4     # Still save an empty processed file
5     self.save_json([])
6     return

```

6.2 Sites with Irregular HTML

```

1 # Always validate before accessing
2 element = soup.find('div', class_='data')
3 if element:
4     text = element.get_text(strip=True)
5     if text: # Check if text exists and isn't empty
6         # Process

```

7 Integration Requirements

7.1 Adding to daily_scraping.py

```

1 # Import
2 from {country}.{provider}.{provider} import {ProviderClass}
3
4 # In scrape() function
5 try:
6     provider_scraper = {ProviderClass}()
7     provider_scraper.scrape()
8 except Exception as e:
9     print(f"Failed to scrape outage data from {Provider}.")

```

Note: It is recommended to run daily_scraping.py near noon time EST to avoid intermittent internet issues.

7.2 Adding to monthly_scraping.py

For providers that offer monthly outage data:

```

1 # Import
2 from {country}.{provider}.{provider} import {ProviderClass}
3
4 # In scrape() function
5 try:
6     provider_scraper = {ProviderClass}()
7     provider_scraper.scrape()
8 except Exception as e:
9     print(f"Failed to scrape monthly data from {Provider}.")

```

7.3 Adding to process_file.py

```

1 # Add to provider lists
2 {country}_providers = ["existing1", "existing2", "{new_provider}"]
3
4 # Add processing logic

```

```
5 elif provider == "{new_provider}":  
6     process = Process_{Provider}(year, month, date, file_path)  
7     process.run()
```

8 Development Process

8.1 Analysis Phase

- Manually visit the target site
- View source to understand HTML structure
- Check if data updates consistently
- Note any authentication/headers needed

8.2 Implementation Order

1. Create basic scraper class with `__init__` and `scrape`
2. Test URL fetching and HTML saving
3. Create processor class and test parsing
4. Integrate both pieces
5. Test with multiple dates

8.3 Testing Requirements

- Test with at least 3 different dates
- Test with a day that has no data
- Verify file naming matches convention
- Check output JSON format

9 Success Criteria

A successful HTML content scraper:

1. Follows naming conventions exactly
2. Produces standard JSON output
3. Handles missing data gracefully
4. Integrates with existing scripts
5. Works reliably across multiple dates

10 Additional Resources

- **Detailed workflow:** /docs/workflow.md
- **Notion documentation:** See repository README for link
- **Repository structure:** Follows country/provider pattern
- **Data analysis:** Available in Data Analysis directory with KeepItOn raw CSV data

10.1 Repository Structure Overview

```
data_analysis/          # Raw CSV data and analysis
{country}/              # Country-specific directories
  {provider}/           # Provider-specific code
    {provider}.py       # Main scraper
    process_{provider}.py # Processor
docs/
  workflow.md           # Detailed workflow
daily_scraping.py        # Daily automation
monthly_scraping.py      # Monthly automation
process_file.py          # Batch processing
```