

**ext\_aeid5473-ChanelAsk\_v13**

**Analysis & Design Document**

**Date – 29 Oct 2022**

## **Table of Contents**

<b>Topic</b>	<b>Page No.</b>
1. Scope of work .....	3
2. Solution Approach.....	3-5
3. Script Development Flow.....	6
4. Technology Considerations.....	7
5. Base Collector Code.....	8-9
6. Template Parameters & Description.....	10
7. Risk & Dependencies.....	11

## 1. Scope of work

Scrap the below data from SITE: <https://www.chanel.com/> for all the operational countries.

1. Product ID
2. Product Name
3. Product Size
4. Product Price
5. Product Availability
6. Product Details (Optional)

## 2. Solution Approach

We are following the below steps to develop the script as per the requirement

- As the website is operational in many countries, we have created the collector code for 47 countries listed below.
- We checked which categories have the products that we need to scrap and then locate the URLs of that category.
- Checked the javascript data (the data we get from AJAX calls) with the help of view page source.
- We are fetching all category links then fetching the product links from where we are scraping the required product details.

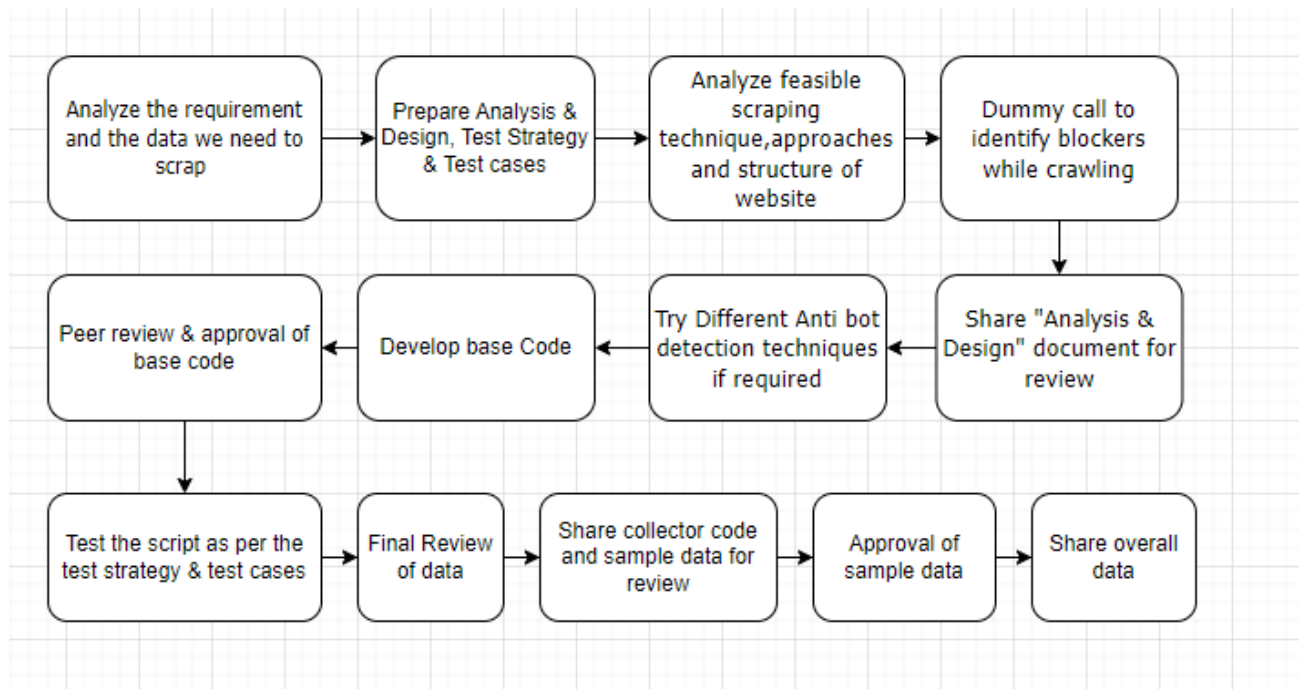
Country	Link	Spider
Albania , EN	<a href="https://www.chanel.com/al/">https://www.chanel.com/al/</a>	Done
Australia , EN	<a href="https://www.chanel.com/au/">https://www.chanel.com/au/</a>	Done
Austria , EN	<a href="https://www.chanel.com/at/">https://www.chanel.com/at/</a>	Done
Belgium , EN	<a href="https://www.chanel.com/be-fr/">https://www.chanel.com/be-fr/</a>	Done
Bosnia and Herzegovina ,EN	<a href="https://www.chanel.com/ba/">https://www.chanel.com/ba/</a>	Done
Brazil ,	<a href="https://www.chanel.com/br/">https://www.chanel.com/br/</a>	Done
Bulgaria, EN	<a href="https://www.chanel.com/bg/">https://www.chanel.com/bg/</a>	Done
Canada , EN	URL is not operational	
Croatia , EN	<a href="https://www.chanel.com/hr/">https://www.chanel.com/hr/</a>	Done
Czech Republic ,EN	<a href="https://www.chanel.com/cz/">https://www.chanel.com/cz/</a>	Done
Denmark , EN	<a href="https://www.chanel.com/dk/">https://www.chanel.com/dk/</a>	Done

## Analysis & Design Document

Estonia , EN	<a href="https://www.chanel.com/ee/">https://www.chanel.com/ee/</a>	Done
Finland , EN	<a href="https://www.chanel.com/fi/">https://www.chanel.com/fi/</a>	Done
France, EN	<a href="https://www.chanel.com/fr/">https://www.chanel.com/fr/</a>	Done
Germany, EN	<a href="https://www.chanel.com/de/">https://www.chanel.com/de/</a>	Done
Greece, EN	<a href="https://www.chanel.com/gr/">https://www.chanel.com/gr/</a>	Done
Hong Kong S.A.R ,EN	<a href="https://www.chanel.com/hk-en">https://www.chanel.com/hk-en</a>	Done
Hungary, EN	<a href="https://www.chanel.com/hu/">https://www.chanel.com/hu/</a>	Done
Italy,EN	<a href="https://www.chanel.com/it/">https://www.chanel.com/it/</a>	Done
Japan,	<a href="https://www.chanel.com/jp/">https://www.chanel.com/jp/</a>	Done
Kingdom Of Saudi Arabia , EN	<a href="https://www.chanel.com/sa/">https://www.chanel.com/sa/</a>	Done
Korea,	<a href="https://www.chanel.com/kr/">https://www.chanel.com/kr/</a>	Done
Kuwait, EN	<a href="https://www.chanel.com/kw/">https://www.chanel.com/kw/</a>	Done
Latin America,	<a href="https://www.chanel.com/lx/">https://www.chanel.com/lx/</a>	Done
Latvia , EN	<a href="https://www.chanel.com/lv/">https://www.chanel.com/lv/</a>	Done
Lithuania , EN	<a href="https://www.chanel.com/lt/">https://www.chanel.com/lt/</a>	Done
Luxembourg, EN	<a href="https://www.chanel.com/lu-fr/">https://www.chanel.com/lu-fr/</a>	Done
Malaysia, EN	<a href="https://www.chanel.com/my/">https://www.chanel.com/my/</a>	Done
Mexico	<a href="https://www.chanel.com/mx/">https://www.chanel.com/mx/</a>	Done
Netherlands - English (UK)	<a href="https://www.chanel.com/nl/">https://www.chanel.com/nl/</a>	Done
Norway,EN	<a href="https://www.chanel.com/no/">https://www.chanel.com/no/</a>	Done
Poland, EN	<a href="https://www.chanel.com/pl/">https://www.chanel.com/pl/</a>	Done
Portugal , EN	<a href="https://www.chanel.com/pt/">https://www.chanel.com/pt/</a>	Done
Qatar , EN	<a href="https://www.chanel.com/qa/">https://www.chanel.com/qa/</a>	Done
Romania , EN	<a href="https://www.chanel.com/ro/">https://www.chanel.com/ro/</a>	Done
Russia, EN	<a href="https://www.chanel.com/ru/">https://www.chanel.com/ru/</a>	Done
Serbia, EN	<a href="https://www.chanel.com/rs/">https://www.chanel.com/rs/</a>	Done
Singapore , EN	<a href="https://www.chanel.com/sg/">https://www.chanel.com/sg/</a>	Done
Slovakia , EN	<a href="https://www.chanel.com/sk/">https://www.chanel.com/sk/</a>	Done
Slovenia , EN	<a href="https://www.chanel.com/si/">https://www.chanel.com/si/</a>	Done
Spain , EN	<a href="https://www.chanel.com/es/">https://www.chanel.com/es/</a>	Done
Sweden , EN	<a href="https://www.chanel.com/se/">https://www.chanel.com/se/</a>	Done
Taiwan Region ,	<a href="https://www.chanel.com/tw/">https://www.chanel.com/tw/</a>	Done

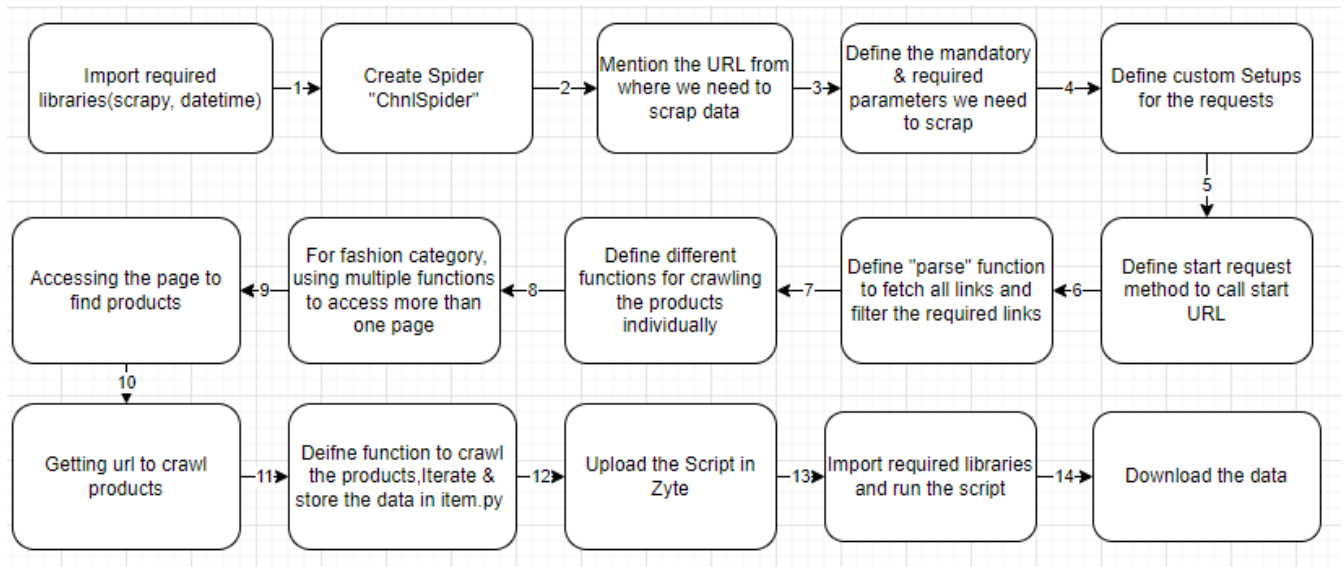
## Analysis & Design Document

Thailand	<a href="https://www.chanel.com/th/">https://www.chanel.com/th/</a>	Done
Turkey , EN	<a href="https://www.chanel.com/tr/">https://www.chanel.com/tr/</a>	Done
United Arab Emirate , EN	<a href="https://www.chanel.com/ae/">https://www.chanel.com/ae/</a>	Done
Middle East (EN)	<a href="https://www.chanel.com/ae/">https://www.chanel.com/ae/</a>	
UK , EN	<a href="https://www.chanel.com/gb/">https://www.chanel.com/gb/</a>	Done
US , EN	<a href="https://www.chanel.com/us/">https://www.chanel.com/us/</a>	Done
Vietnam	<a href="https://www.chanel.com/vn/">https://www.chanel.com/vn/</a>	Done



### 3. Script Development Flow

Below steps are followed to create spider



## 4. Technology Considerations

**Custom signup** - Not required

**Programming Language** - Python

**Framework** - Scrapy

**Tool** - Zyte

**Functions & Libraries used** - datetime, scrapy-user-agents

**Storage (Database)** - Zyte Cloud

### **Deployment Requirements**

- Install all the required libraries in Zyte Cloud

### **Logging considerations**

- No logging is required
- No CAPCTHA authentication required

### **Proxy Details**

- We are using user agent to avoid getting blocked, this is present in settings.py file.

## 5. Base Collector Code

**File name** - chnl.py

Here we are scraping the data as per the requirements

### Step 1 - Importing required libraries

```
import json
import scrapy
from ..items import ChanelItem
from datetime import datetime
```

### Step 2 - Here we created the spider "ChnlSpider" and added allowed domain, start url of all the product website URL are defined that we are crawling

```
class ChnlSpider(scrapy.Spider):
    name = 'chanel'
    site = 'https://www.chanel.com'
    start_urls = ['https://www.chanel.com/us/']
    frag_url = 'https://www.chanel.com/us/fragrance/women/c/7x1x1/page-'
    eyegls_url = 'https://www.chanel.com/us/eyewear/eyeglasses/c/2x1x2/page-'
    sungls_url = 'https://www.chanel.com/us/eyewear/sunglasses/c/2x1x1/page-'
    watch_url = 'https://www.chanel.com/us/watches/collection/c/4x2/page-'
    fine_jewellery = 'https://www.chanel.com/us/fine-jewelry/collection/c/3x2/page-'
    face_url = 'https://www.chanel.com/us/makeup/face/c/5x1x6/page-'
    eye_url = 'https://www.chanel.com/us/makeup/eyes/c/5x1x4/page-'
    lips_url = 'https://www.chanel.com/us/makeup/lips/c/5x1x1/page-'
    nails_url = 'https://www.chanel.com/us/makeup/nails/c/5x1x7/page-'
    brushes_url = 'https://www.chanel.com/us/makeup/brushes-and-accessories/c/5x1x3/page-'
    nxp = []
    execution_id = '621291'
    feed_code = 'aaid5473'
    record_create_by = 'aaid5473_chanel'
    record_create_date = datetime.now()
    source_country = 'USA'
    for i in range(10):
        nxp.append(1)
```

### Step 3 - Here we are defining the custom details

```
custom_settings = {
    'SCHEDULER_PRIORITY_QUEUE': 'scrapy.pqueues.DownloaderAwarePriorityQueue',
```



```
'REACTOR_THREADPOOL_MAXSIZE': '20',  
'LOG_LEVEL': 'INFO',  
'RETRY_ENABLED': 'False',  
'DOWNLOAD_TIMEOUT': '1000',  
'REDIRECT_ENABLED': 'False',  
'AJAXCRAWL_ENABLED': 'True',  
'CONCURRENT_REQUESTS_PER_DOMAIN': '2',  
'DNS_RESOLVER': 'scrapy.resolver.CachingThreadedResolver',  
'DUPEFILTER_CLASS': "scrapy.dupefilters.BaseDupeFilter",  
'AUTOTHROTTLING_ENABLED': 'False'  
}
```

**Step 4 - Here we are fetching all links from selector and filtering the required urls**

```
def parse(self, response):
```

**Step 5 - Here we are defining the functions for crawling the products category individually**

```
def sunglass(self , response):  
def eyeglass(self , response):  
def fragrance(self , response):  
def watches(self , response):  
def fine_jwlry(self , response):  
def mkp_face(self , response):  
def mkp_eye(self , response):  
def mkp_lips(self , response):  
def mkp_nails(self , response):  
def mkp_brushes(self , response):
```

**Step 6 - For fashion category using multiple functions to access more than one page**

```
def parse_page(self,response):
```

**Step 7 - Accessing the page to find products**

```
def parse_products(self,response):
```

**Step 9 - Getting url to crawl products finally**

```
def parse_categories(self, response):
```

**Step 10 - Defining the function to crawl the required products**

```
def scrape(self,response):
```

**Step 11 - yielding all items here**

```
yield item
```

## 6. Template Parameters & Description

The template contains the data that is scraped as per the ranking of newly listed products.

For the parameters where **mandatory** is mentioned, this is mandatory parameters as per the required template.

For the parameters where **Required** is mentioned, this is parameters needed as per the requirement document.

Below are the parameters that we are scraping and their description

1. **Context\_identifier (Mandatory)** - We are capturing the hierarchy of product in a website
2. **Execution\_id (Mandatory)** - Execution id will be taken automatically from zyte.
3. **Feed\_code (Mandatory)** - This is hardcoded as project name.
4. **Availability (Required)** - This we are getting from website
5. **Color (Required)** - This we are getting from website
6. **Description (Required)** - This we are getting from website
7. **Material (Required)** - This we are getting from the website.
8. **Name (Required)** - This we are getting from the website.
9. **Price (Required)** - This we are getting from the website.
10. **product\_id (Required)** - This we are getting from the website.
11. **Size (Required)** - This we are getting from the website.
12. **Record\_create\_by (Mandatory)** - This is hardcoded with spider name
13. **Record\_create\_dt (Mandatory)** - This is the timestamp for capturing the data.
14. **Site (Mandatory)**- This is hardcoded.
15. **Source (Mandatory)** - This is the link of the individual product.
16. **Source\_country (Mandatory)** -This is hardcoded as per the specific country.

## 7. Risks and Dependencies

Below are the identified risks and their possible solutions:

Risk	Mitigation
Risk of getting blacklisted/blocked/IP restrictions due to security/network policies on the web server.	we need to control the concurrency & use different proxy methods.
If the semantic code/markup of the website changes, the script will have a possibility of failure.	Identify the changes in the semantic code/markup of the website and modify the script accordingly.