

Extract Trends from social media data

Team Name: Yorozyua

Institute Name: PSG College of Technology (PSGCT),
Coimbatore

Team members details

Team Name	Yorozuya		
Institute Name	PSG College of Technology (PSGCT), Coimbatore		
Team Members >	1 (Leader)	2	3
Name	Sairam Vaidya M.	Guna M.	Lohith Sowmiyan P. S.
Batch	4 th Year	4 th Year	4 th Year

Deliverables/Expectations for Level 2 (Idea + Code Submission)

Deliverable 1:

Identification of trends from social media

1. Identify trends on social media based on category. Can restrict to Fashion as a category for the project. Ex: Polka dots dresses are trending on twitter.
2. Ranking/scoring logic for trends extracted.
3. Outcome format:
 - a. Option1: List of trending keyword(s) along with list of sample images and respective links from which the trend is derived with most trending first:
Example: Trends:[{Polka dot dresses, <list of links/images>,trending score}, {Bellbottom Jeans, <list of links/images>,trending score}..]
 - b. Option 2: structured data according to flipkart category, sub category, vertical and product attributes
Example: {category: Fashion, Sub-category: Women Western, vertical: Women dresses, trending attribute type: Pattern, trending attribute value: Polka Print, list of sample images and links from which the trend is derived}.
Outcome with Option 2 format will be given bonus points.

Deliverable 2:

Mapping trends with Flipkart products:

1. Create mapping of extracted trending keyword(s) with Flipkart category, sub category, vertical and product attribute(s), search page links.
Example:{category: Fashion, Sub-category: Women Western, vertical: Women dresses, trending attribute type: Pattern, trending attribute value: Polka Print}
Note: Use category, Subcategory combination from the Flipkart Website
2. From a trending keyword, creating a corresponding searchable term on Flipkart which will lead to matching products.
Example: Tropical Tops keywords will not give right results directly on Flipkart but we can construct search query for it using some intelligence.
3. Points will be given based on similarity between sample images for trends and product results on Flipkart.

Glossary

1. **VGG-16** – Visual Geometry Group - 16
2. **API** – Application Programming Interface
3. **YOLOv5** – You Only Look Once Version 5
4. **NLU** – Natural Language Understanding
5. **BERT** – Bidirectional Encoder Representations from Transformers
6. **D2V** – Doc2Vec
7. **W2V** – Word2Vec
8. **SQ** – Search Query
9. **PR** – Possible Results
10. **FR** – Fuzzy Ratio

Use-cases

Highest Priority Use Case(s)

1. **Automated Bots:** These bots are used to perform asynchronous scraping of images and metadata from various platforms such as Instagram, Pinterest, etc.
2. **Trend Identification:** This is achieved using deep neural networks (based on VGG16 Architecture) to identify the trends present in the data collected.
3. **Smart Mapper:** This is built using a combination of Fuzzy techniques, Rule based learning and the Wordnets library to map a given phrase to Flipkart relevant format.

Medium Priority Use Case(s)

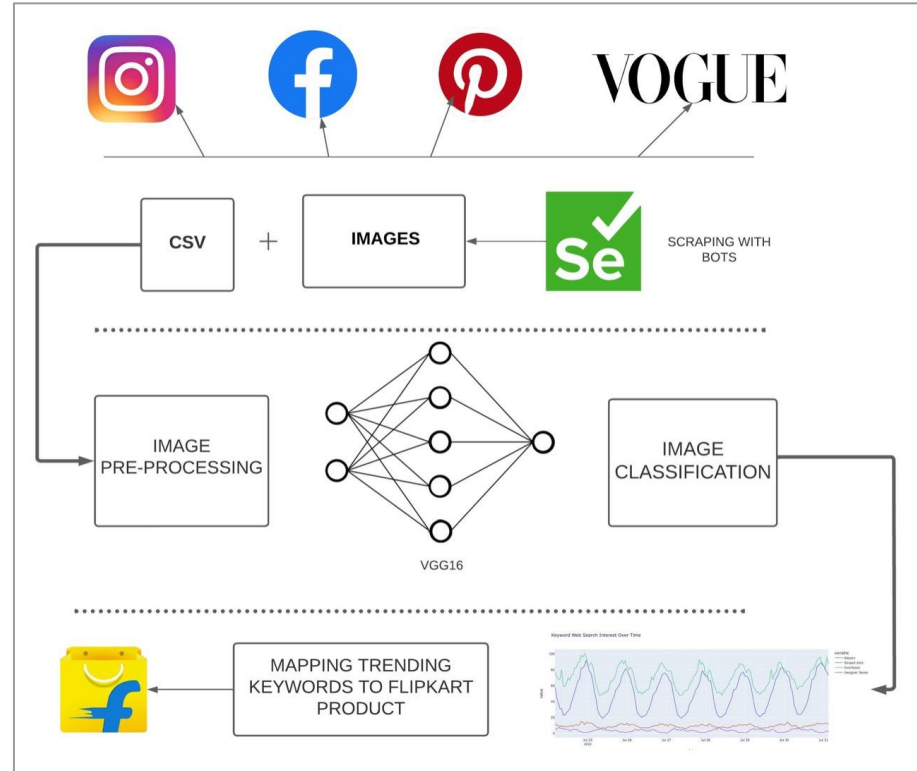
1. **Trends Corroboration:** The classified images need to be tested, and corroborated with actual trends observed; for this, a combination of Google Trends API and PyTrends API is used.

Proposed Approach

A **three-tier architecture** is used, with the layers,

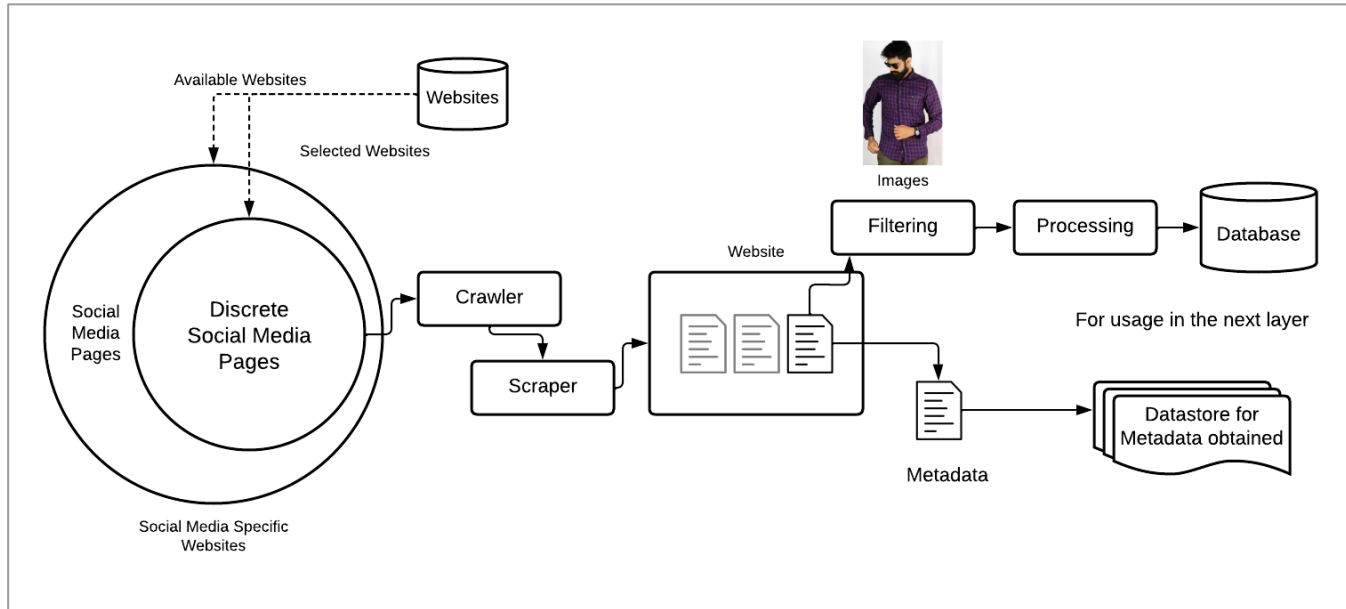
1. Crawling for Dataset Creation
2. Trends Identification and Classification
3. Flipkart Relevant Smart Mapper

A top view of the approach is shown on the image,



Proposed Approach

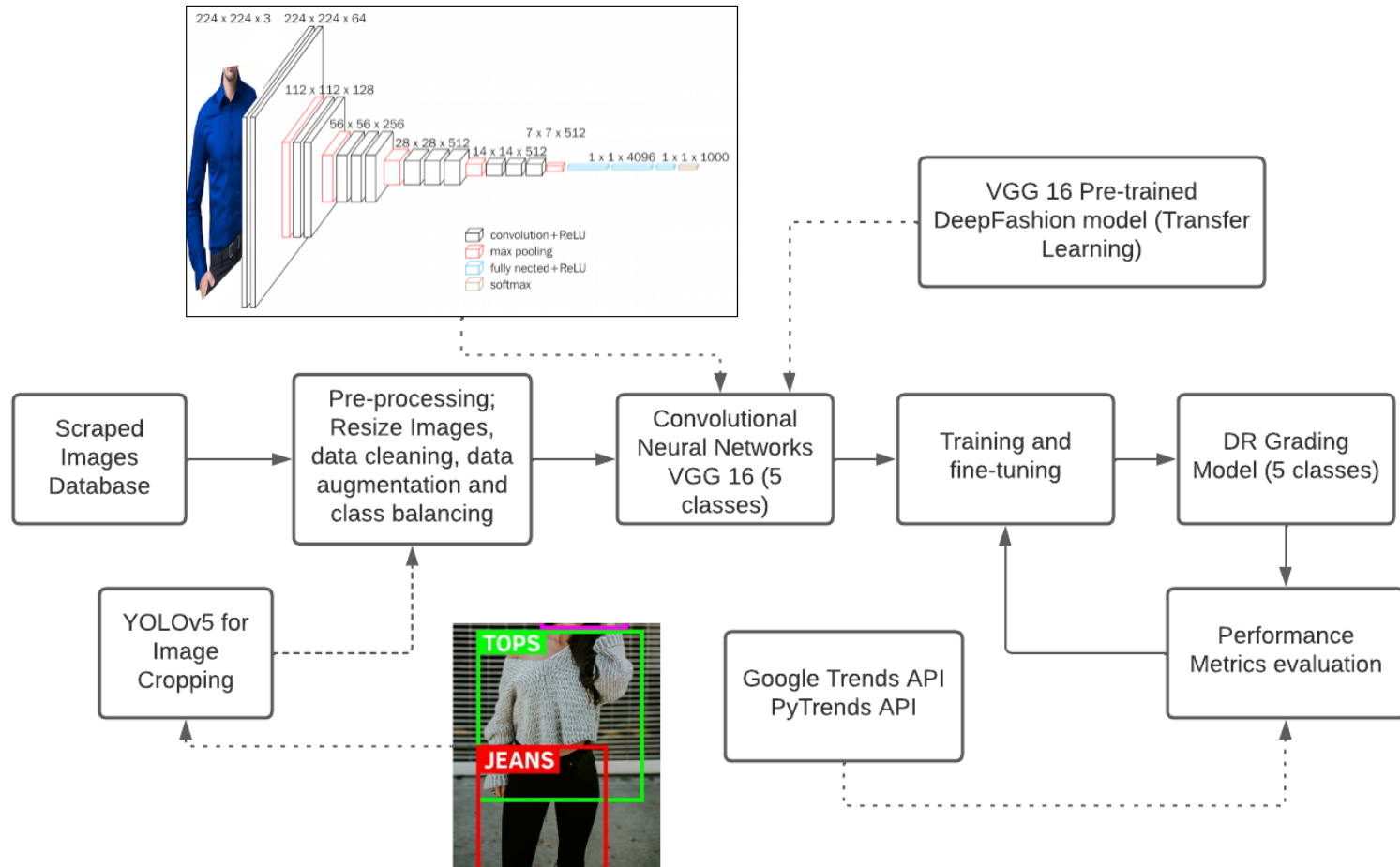
1. **Crawling and Scrapping** social media websites for images and metadata. The workflow used for achieved this is displayed below,



Proposed Approach

2. Trends Identification and Classification using a deep neural network based on the VGG16 architecture, and corroborating the trends with PyTrends API and Google Trends API.

- We make use of the VGG-16 architecture because it is one of the best models when it comes to image classifications, and consists of roughly 136 million parameters.
- As a result, we are able to classify and identify a large number of trends in the images – which can be easily scaled in the future.
- For images with unnecessary features such as human faces, arms, pets, etc., we use the YOLOv5 architecture to identify the “fashion” (clothes, footwear, etc.) and crop them out.
- We obtain the top trending fashion keywords from Google Trends API and PyTrends API and use it to corroborate/evaluate the result obtained.
- The described second tier is shown in the next slide,



Proposed Approach

3. **Smart Mapper** that is applicable to the categories, and verticals available on the Flipkart website was made after considering various techniques,

1. **NLU with BERT Embeddings**: This is a classic approach to identify the semantic meaning present in phrases/sentences, but performed poorly for the given use cases given a possible lack of semantic meaning.

E.g. “floral top red”, “men shirt floral”, “PINK FLORAL TOPS!!!” are possible search terms that are commonly observed but lack an inherent semantic meaning, which leads to the poor performance of the considered approach.

2. **D2V and W2V**: When it comes to sentence embedding, D2V is a ubiquitous model that is considered – however to similar reasons as mentioned above, it performs underwhelmingly for the resources it demands. An argument can be made for W2V to be used since most search queries are a combination of few words, but it also performs underwhelmingly in multiple cases, especially since most of the product titles on e-commerce websites are very long in comparison.

Proposed Approach

E.g. Search Query: "floral top red"

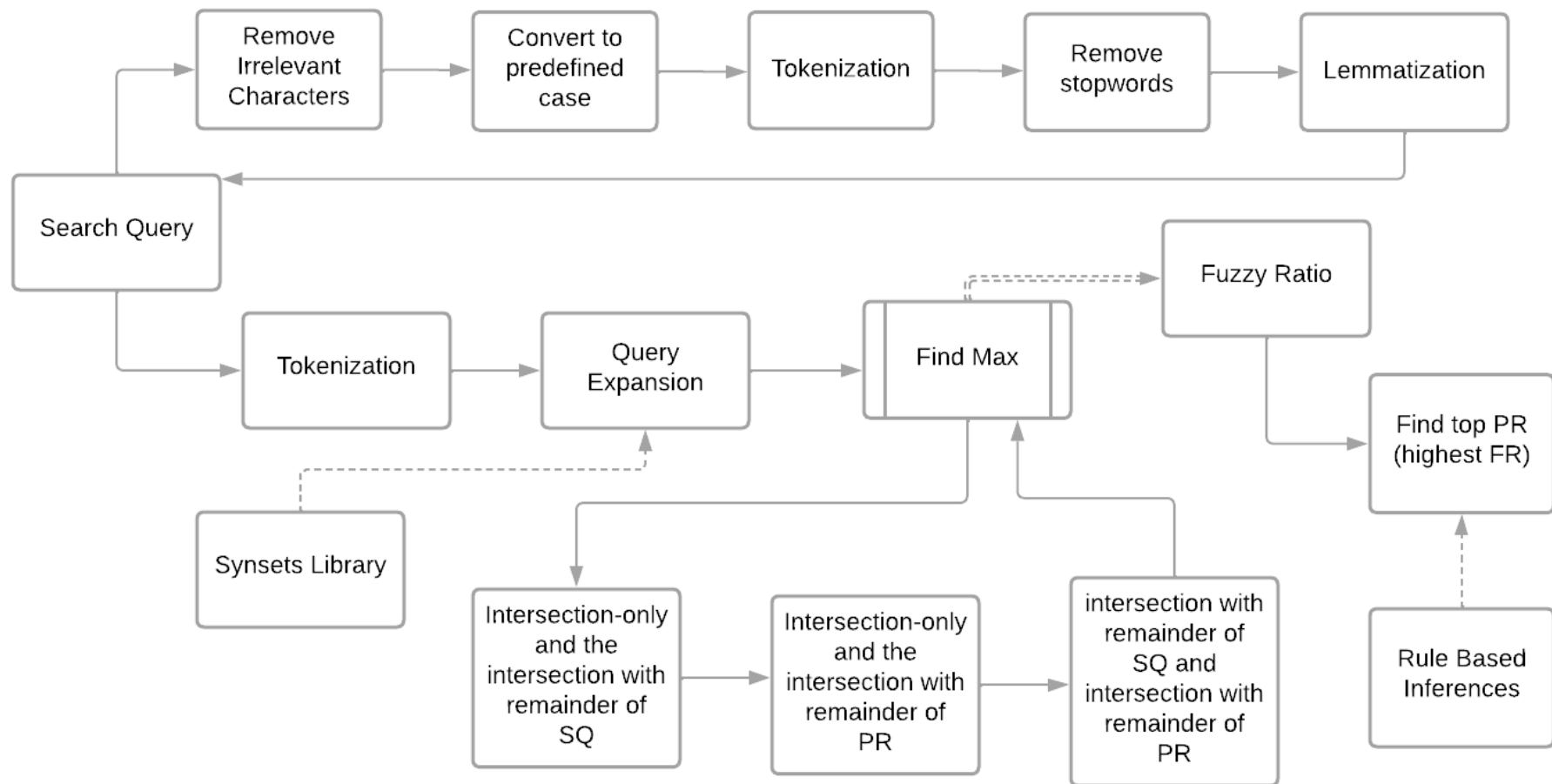
Possible Product Title: Women's 100% Floral Super Soft Pure Cotton
Top | Limited Edition

3. Distance Based Measures: Distance based measures such as Sørensen-Dice, Levenshtein and Needleman-Wunsch perform poorly when a large number of unnecessary words are present during the comparison (or) if any synonym, specific feature is present in the dataset.

E.g. Search Query: "magenta large shirt"

Possible Product Title: Big Red Shirt

4. Modified Fuzzy Synsets: Hence, considering the shortcomings of the previous approaches we use Occam's razor and pick the simplest one in principle. We expand the search queries using Synsets, the token-set fuzzy ratio, along with rule based inferences, i.e., if the search query contains "pyjama", the result should also contain "pyjama". This approach is illustrated in the next slide.



Limitations

1. **Asynchronous Operations:** The bots used for scraping might be unable to work properly at times due to a websites policies since they work asynchronously. As a result, manual supervision might be needed.
2. **Identification of Intricate Trends:** Given the plethora of designs, and ever-growing/changing pace of the fashion industry, it might be difficult to cover every single patter/trend on the uptake. However, to combat this issue we need to collect a very high number of images in high resolution to ensure that the intricate patterns are captured.
3. **Semantic Persistence in Queries:** Even though the semantic approaches have their downsides, in the extreme situation where a user inputs "I don't want shirt; I want jeans" it is impossible to approach such situations without semantic understanding. Hence, it cannot be completely eliminated, and needs to be integrated (even if at a low level) into the smart mapper.

Future Scope

1. **Increase the Generality in Scraping:** Currently specific scrapers need to be written, for specific websites. For extending the volume of images scraped in the future, a generalized scraper needs to be created that is able to extract relevant images and data from any website.
2. **Combination of Search Queries:** Currently, the model supports minimal combination of queries. In the future, we would like to further work on our model to be able to classify images with multiple pieces of fashion such as, "Red Beanie, Black Bomber Jacket, Blue Jeans, Brown Leather Shoes, White Watch"
3. **Extending Approach to other Categories:** After fashion, electronics is the most sought after category and on arriving upon a set method, the idea used and approaches can be extended to other categories.