

# **DOCUMENTATION FOR GLANCE ASSIGNMENT**

## **BY- M.MEYYAPPAN**

### **APPROACHES:**

Since we are developing a fashion retrieval system we need to consider the data reliability and precision along with the retrieval speed.

Here I have used the fashionpedia dataset and they have the metadata along with them , but they are useful when we need to perform a quick retrieval and not something with precision, which is the main drawback.

To understand the color , the vibe , the surroundings only metadata based retrieval is not useful so we use Vision Language Models to create descriptions for the images which is the synthetic metadata and then use a cross – encoder to finally perform a precise accurate retrieval of the images , using the natural language queries, but this approach comes with a tradeoff which is speed , but considering that we have ample VRAM in a Gpu we can use this method.

GITHUB LINK : <https://github.com/mm0177/GLANCE-ML/tree/main>

### **Version -1 ( Approach 1)**

So here with the first approach we have used a Decomposed Semantic Retrieval Flow , which separates it from the standard Clip Model Usage.

We have used Microsoft-Florence-2 Model to create a detailed caption for the images, it defines each of the elements like the background , the material texture etc.

Then we use the Vit-L/14 which encodes the image into 5 distinct vectors into a local Qdrant Database which includes ( Global , Clothing, Color, Scene, Style).

Then we have the Regex Based Query Parser , which analyses the user's text , where in it maps the query using Weights. This ensures that the search engine understands which part of the user's query is most important.

Then finally we have the Cross -encoder ( ms-marco-MiniLM-L-6-v2) ,where in after the Clip Model retrieves the top 50 images , the cross encoder takes the user query and the AI generated metadata from Florence-2 and compare them directly in a single neural pass

This is how the first approach is made.

Lets look at the output

Query 1: 'A person in a bright yellow raincoat'

#1 (Score: -5.721)



#2 (Score: -5.914)



#3 (Score: -5.965)



#4 (Score: -6.050)



#5 (Score: -6.292)



coat, pants, collar +4 more  
 tan  
 street

pocket, zipper, lapel +3 more  
 red  
 studio

zipper, lapel, shoe +5 more  
 red  
 urban

zipper, pocket, shoe +6 more  
 red  
 casual setting

pocket, coat, collar +3 more  
 tan  
 street

Query 2: 'Professional business attire inside a modern office'

#1 (Score: 4.704)



#2 (Score: 4.095)



#3 (Score: 4.071)



#4 (Score: 3.986)



#5 (Score: 3.534)



pocket, lapel, coat +4 more  
□ office

shoe, belt, skirt +3 more  
□ office

shoe, jacket, bag, wallet +6 more  
□ office

shoe, pants, bag, wallet +3 more  
□ office

zipper, pocket, lapel +7 more  
□ office

Query 3: 'Someone wearing a blue shirt sitting on a park bench'

#1 (Score: 1.790)



#2 (Score: 1.715)



#3 (Score: 1.531)



#4 (Score: 1.454)



#5 (Score: 1.374)



top, t-shirt, sweatshirt, neckline  
□ park

shirt, blouse, pocket, shoe +4 more  
□ park

shirt, blouse, dress, hat +2 more  
□ park

neckline, pocket, pants +2 more  
□ park

skirt, shirt, blouse  
□ park

Query 4: 'Casual weekend outfit for a city walk'

#1 (Score: 0.243)



#2 (Score: -0.376)



#3 (Score: -0.574)



#4 (Score: -0.647)



#5 (Score: -0.753)



shoe, pants, bag, wallet +3 more  
 city

pocket, shoe, jacket +5 more  
 city

pocket, rivet, zipper +7 more  
 city

shoe, bag, wallet, skirt +3 more  
 park

shorts, sleeve, cardigan +3 more  
 casual setting

Query 5: 'A red tie and a white shirt in a formal setting'

#1 (Score: 3.182)



#2 (Score: 2.893)



#3 (Score: 2.765)



#4 (Score: 2.716)



#5 (Score: 2.595)



fringe, skirt, top, t-shirt, sweatshirt +1 more  
 red  
 formal setting

shoe, bag, wallet, skirt +3 more  
 red  
 formal setting

shoe, ruffle, skirt +3 more  
 red  
 formal setting

shoe, belt, skirt +3 more  
 red  
 formal setting

shoe, tie, dress +5 more  
 red  
 formal setting

I have put only the images of 5 queries in the document , to check the output of all the 10 queries Kindly view the Notebook in Github.

Link: <https://github.com/mm0177/GLANCE-ML/blob/main/VERSION-1.ipynb>

## Version 2 ( Approach -2)

In the version 2 , the captions which we get from the Florence-2 are vectorized and saved in Qdrant DB ( compared to version 1 where the captions were used just as metadata).

For the CLIP model (ViT-L/14) , it performs a Zero Shot Color Detection , where in it uses prompt ensembling ( which means “ is it crimson” , “is it scarlet” ) , and verifies and corrects the dataset tags to ensure that even if an image is untagged it is correctly indexed.

For the query Parser , I have made it intent aware , where in it uses Synonym Mapping and Adaptive Weighting. It analyses the complexity of query and if the query is long , then it shifts the search power from simple tags to Florence-2 caption.

For the Cross-Encoder we are using an L-12 , to perform better.

I have also introduced a new function which is search by image which the CLIP Global and Clothing vectors to find "Visual Twins." It ignores the text pathway entirely to focus on the silhouette, cut, and "vibe" of an uploaded image.

Query 1: 'A person in a bright yellow raincoat'



Query 2: 'Professional business attire inside a modern office'



Query 3: 'Someone wearing a blue shirt sitting on a park bench'



Query 4: 'Casual weekend outfit for a city walk'



Query 5: 'A red tie and a white shirt in a formal setting'

#1 (Score: 1.975)



#2 (Score: 1.906)



#3 (Score: 1.748)



#4 (Score: 1.643)



#5 (Score: 1.476)



top, t-shirt, sweatshirt, fringe, neckline +1  
 red  
 formal setting

tights, stockings, ruffle, shirt, blouse +4  
 red  
 formal setting

ruffle, top, t-shirt, sweatshirt, hat +4  
 red  
 formal setting

ruffle, top, t-shirt, sweatshirt, skirt +2  
 red  
 formal setting

glasses, ruffle, top, t-shirt, sweatshirt +5  
 red  
 formal setting

To check all the query outputs kindly view the github notebook

Link : <https://github.com/mm0177/GLANCE-ML/blob/main/VERSION-2.ipynb>

### Version -3(Apporach-3)

In Version -3 we introduce Blip to work along side Florence-2 , it gives a narrative flow , which helps in getting better details of the clothing and the environment.

Now we have 7 distinct vectors per image to get a better understanding of the image.

For the color detection , in version 2 it was averaged color names (synonyms) into a single hardcoded template, but now For every color (e.g., "Crimson"), the system now generates 5+ different contextual sentences (e.g., "a photo of crimson fabric," "close-up of crimson garment," "crimson colored clothes").

Also the intent of the query is understood so like if there is more focus on a color then the search weight shifts to Florence and Clothing Vectors, but if its something more about the environment then it shifts to BLIP and Style Vectors.

Now the Reranker uses the Dataset Tags , Florence Technical Captions , Blip Narrative Captions and Clip Detected colors.

Link To Notebook: <https://github.com/mm0177/GLANCE-ML/blob/main/VERSION-3.ipynb>

Query 1: 'A person in a bright yellow raincoat'



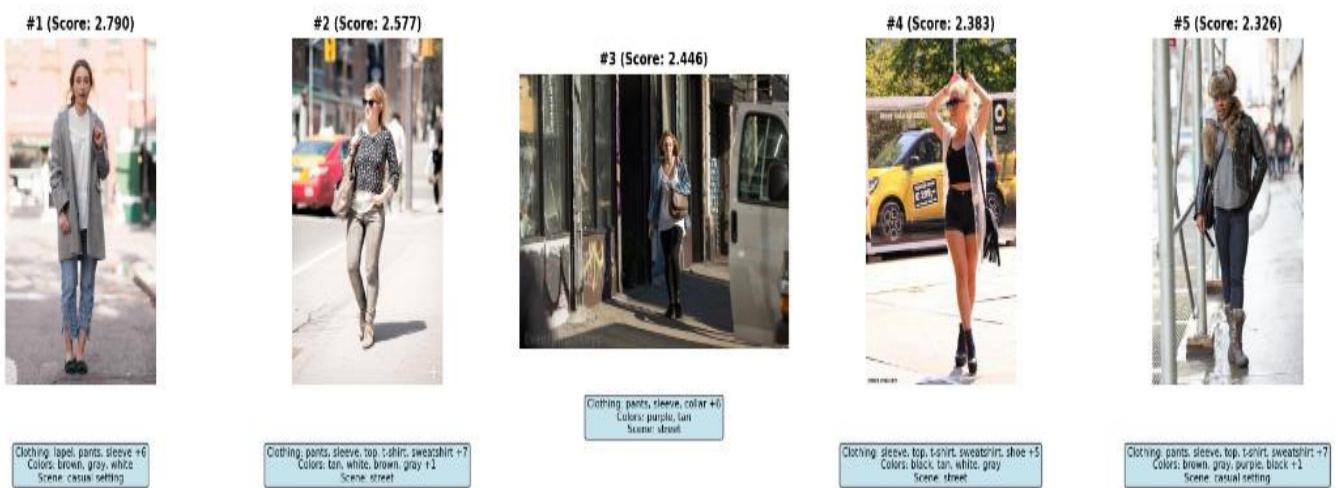
Query 2: 'Professional business attire inside a modern office'



Query 3: 'Someone wearing a blue shirt sitting on a park bench'



Query 4: 'Casual weekend outfit for a city walk'



Query 5: 'A red tie and a white shirt in a formal setting'



## FUTURE WORK AND IMPORVEMENTS

- 1) To incorporate location and weather awareness, you could enrich your system with **geo-temporal contextual embeddings** by training a specialized adapter that maps location-weather pairs (e.g., "Tokyo summer", "NYC winter") to appropriate fashion attributes. This could be achieved by collecting a dataset of fashion images tagged with location metadata and weather conditions, then fine-tuning a small auxiliary network that modulates your existing CLIP embeddings based on contextual factors
- 2) Currently the version 1 uses 10k images , the version 2 uses 15k images and the version 3 uses 20k images. To imporve the performace we can increase it to 35k images and also create a fine tuned model for the color expertise so that the color matching can be improved.
- 3) Also we can use user feedback so that using the query the model can understand the user preference over -time.