# CS331 Project: Retrieval Augmented Generation
## Assignment 6

**Team Members:**
Khushi Mandal - 2201108
Arya Sahu - 2201033
Anushka Srivastava - 2201030
Ahlad Pataparla - 2201017

March 20, 2025

# Q1. Core Functional Modules (Business Logic Layer)

This section details the core functional modules of the application, focusing on the business logic layer (BLL). It describes the purpose of each module, provides code snippets, and shows how they interact. The interaction diagram is provided at the end of this section.

## 1. Product Retrieval

**Module:** `get_item(item_id)`

    **Purpose:** Retrieves a single product from the preloaded Pandas DataFrame (which is loaded from the database during application startup) based on its ID. This function also constructs the image URL for the frontend.

    **Code:**

```python
def get_item(item_id: str) -> Dict[str, Any]:
    """Retrieve an item from the preloaded dataframe."""
    try:
        item = ml_model.df[ml_model.df["id"] == item_id].iloc[0].
            to_dict()
        item["id"] = int(item["id"])
        item["image_url"] = f"/static/images/{item['id']}.jpg"
        return item
    except IndexError:
        raise HTTPException(status_code=404, detail="Item not found")
```

    **Interaction:**

- Called by: `product_page`, `get_random_products`

- Interacts with: The preloaded `ml_model.df` (Pandas DataFrame). No direct database interaction, as the data is loaded at startup.

## 2. Random Product Retrieval

**Module:** `get_random_products(limit=10)`

**Purpose:** Retrieves a specified number of random products from the preloaded DataFrame. This is used for displaying a selection of products on the homepage or other sections where a diverse set of items is needed.

**Code:**

```python
@app.get("/api/products", response_model=ProductsResponse)
async def get_random_products(limit: int = 10):
    """Return a random selection of products."""
    try:
        all_product_ids = ml_model.df["id"].tolist()
        selected_ids = sample(all_product_ids, min(limit, len(
            all_product_ids)))
        products = [Item(**get_item(product_id)) for product_id in
            selected_ids]
        return ProductsResponse(products=products)
    except Exception as e:
        print(f"Error fetching random products: {e}")
        raise HTTPException(status_code=500, detail="An error occurred
            while fetching random products")
```

**Interaction:**

- Called by: Frontend (e.g., homepage).

- Interacts with: `get_item` (to retrieve details of each selected product), `ml_model.df`.

## 3. Attribute Prediction

**Module:** `predict_attributes(image)`

**Purpose:** Takes an image as input and uses the CLIP model to predict various attributes of the item in the image. These attributes include gender, article type, season, usage, master category, subcategory, and base color. The dominant color is extracted from the image, and the closest color name is found using a predefined color map.

**Code:**

```python
def get_dominant_color(image: Image.Image) -> np.ndarray:
    """Get the dominant color from an image."""
    image = image.resize((100, 100))
    img_array = np.array(image).reshape(-1, 3)
    kmeans = KMeans(n_clusters=3, random_state=0).fit(img_array)
    counts = np.bincount(kmeans.labels_)
    return kmeans.cluster_centers_[np.argmax(counts)]

def find_closest_color(target_color: np.ndarray, color_names: List[str
    ]) -> str:
    """Find the closest color name to the target RGB."""
    color_map = {
        "Navy Blue": (0, 0, 128),
        "Blue": (0, 0, 255),
        "Black": (0, 0, 0),
        "Silver": (192, 192, 192),
        "Grey": (128, 128, 128),
                # ... (rest of the color map)
        "Fluorescent Green": (127, 255, 0),
    }

    target_rgb = target_color.astype(int)
```

```
22      min_dist, closest = float("inf"), "Black"
23      for name, rgb in color_map.items():
24          dist = np.linalg.norm(np.array(rgb) - target_rgb)
25          if dist < min_dist:
26              min_dist, closest = dist, name
27      return closest
28
29  def predict_attributes(image: Image.Image) -> dict:
30      """Predict attributes from an image using CLIP."""
31      attributes = {}
32      for label_type, labels in [
33          ("gender", ["Men", "Women", "Boys", "Girls", "Unisex"]),
34          ("articleType", ml_model.df["articleType"].unique().tolist()),
35          ("season", ["Summer", "Winter", "Spring", "Fall"]),
36          ("usage", ["Casual", "Ethnic", "Formal", "Sports", "Smart
              Casual", "Travel", "Party", "Home"]),
37          ("masterCategory", ["Apparel", "Accessories", "Footwear", "
              Personal Care", "Free Items", "Sporting Goods", "Home"]),
38          ("subCategory", ["Topwear", "Bottomwear", "Watches", "Socks", "
              Shoes", "Belts", "Flip Flops", "Bags", "Innerwear", "Sandal"
              , "Shoe Accessories", "Fragrance", "Jewellery", "Lips", "
              Saree", "Eyewear", "Nails", "Scarves", "Dress", "Loungewear
              and Nightwear", "Wallets", "Apparel Set", "Headwear", "
              Mufflers", "Skin Care", "Makeup", "Free Gifts", "Ties", "
              Accessories", "Skin", "Beauty Accessories", "Water Bottle",
              "Eyes", "Bath and Body", "Gloves", "Sports Accessories", "
              Cufflinks", "Sports Equipment", "Stoles", "Hair", "Perfumes"
              , "Home Furnishing", "Umbrellas", "Wristbands", "Vouchers"])
39      ]:
40          inputs = ml_model.clip_processor(text=labels, images=image,
              return_tensors="pt", padding=True)
41          outputs = ml_model.clip_model(**inputs)
42          attributes[label_type] = labels[outputs.logits_per_image.
              softmax(dim=1).argmax().item()]
43
44      dominant_color = get_dominant_color(image)
45      attributes["baseColour"] = find_closest_color(dominant_color,
          ml_model.df["baseColour"].unique().tolist())
46      return attributes
```

**Interaction:**

- Called by: `recommend_from_image`

- Interacts with: `ml_model.clip_model`, `ml_model.clip_processor`, `get_dominant_color`, `find_closest_color`.

## 4. Outfit Recommendation

**Modules:**

- `product_page(item_id)`: Provides recommendations for a specific product page.

- `recommend_from_image(file)`: Provides recommendations based on an uploaded image.

- `get_ml_recommendations(...)`: The core recommendation engine (called by both of the above).

- `get_compatible_types(article_type)`: Retrieves compatible types from 'constants.py'.

- `get_accessory_types(usage, season)`: Retrieves accessory types from 'constants.py'.

- `check_negative_constraints(target_item, candidate_item)`: Filters out incompatible combinations.

- `maximal_marginal_relevance(...)`: Ensures diversity in recommendations.

- `color_compatibility(color1, color2)`: Calculates color compatibility.

**Purpose:** These modules work together to generate outfit recommendations. `product_page` and `recommend_from_image` are the API endpoints, while `get_ml_recommendations` performs the core recommendation logic using cosine similarity, MMR, and various filtering rules. The helper functions retrieve compatibility data and enforce constraints.

**Code:**

```python
def maximal_marginal_relevance(target_features, category_features,
    top_n: int, lambda_param: float = 0.1) -> np.ndarray:
    """Select highly diverse items using MMR."""
    selected_indices = []
    remaining_indices = list(range(category_features.shape[0]))

    similarities = cosine_similarity(target_features, category_features
        ).flatten()
    first_idx = np.argmax(similarities)
    selected_indices.append(first_idx)
    remaining_indices.remove(first_idx)

    # Dynamic lambda: more diversity for larger pools, min 0.05
    dynamic_lambda = max(0.05, lambda_param - (len(remaining_indices) /
        2000))

    for _ in range(min(top_n - 1, len(remaining_indices))):
        mmr_scores = []
        for idx in remaining_indices:
            relevance = similarities[idx]
            diversity = min(cosine_similarity(category_features[idx],
                category_features[selected_indices]).flatten())
            mmr_score = dynamic_lambda * relevance - (1 -
                dynamic_lambda) * diversity
            mmr_scores.append(mmr_score)

        next_idx = remaining_indices[np.argmax(mmr_scores)]
        selected_indices.append(next_idx)
        remaining_indices.remove(next_idx)

    return np.array(selected_indices)

def get_ml_recommendations(
    target_features,
    target_article_type: str,
    product_gender: str,
    target_color: str,
    target_id: Optional[str] = None,
    top_n: int = 3
) -> tuple[List[Dict], float, float, float]:
```

```python
36      """Generate recommendations with maximum diversity."""
37      df = ml_model.df
38      # Broaden pool with ARTICLE_TYPE_GROUPS and COMPATIBLE_TYPES
39      target_group = next((group for group, types in ARTICLE_TYPE_GROUPS.
            items() if target_article_type in types), None)
40      compatible_types = COMPATIBLE_TYPES.get(target_article_type, [])
41      if target_group:
42          candidate_types = list(set(ARTICLE_TYPE_GROUPS[target_group] +
                compatible_types))
43      else:
44          candidate_types = [target_article_type] + compatible_types
45
46      mask = (df["articleType"].isin(candidate_types)) & (df["gender"].
            isin([product_gender, "Unisex"]))
47      if target_id:
48          mask &= (df["id"] != target_id)
49
50      category_df = df[mask]
51      logger.info(f"Initial filter: {len(category_df)} items for {
            target_article_type}, gender: {product_gender}")
52
53      if category_df.empty:
54          logger.warning(f"No items found for {target_article_type}")
55          return [], 0.0, 0.0, 0.0
56
57      color_scores = category_df["baseColour"].apply(lambda x:
            color_compatibility(target_color, x))
58      min_threshold = 0.2 if len(category_df[color_scores >= 0.2]) >=
            top_n * 2 else 0.0
59      category_df = category_df[color_scores >= min_threshold]
60      logger.info(f"After color filter (threshold {min_threshold}): {len(
            category_df)} items")
61
62      if category_df.empty:
63          return [], 0.0, 0.0, 0.0
64
65      category_indices = category_df.index.tolist()
66      category_features = ml_model.combined_features[category_indices]
67      color_matches = (category_df["baseColour"] == target_color).astype(
            float)
68
69      similarities = cosine_similarity(target_features, category_features
            ).flatten()
70      similarities += similarities.max() * 0.05 * color_matches
71
72      # Use MMR with maximum diversity emphasis
73      top_indices = maximal_marginal_relevance(target_features,
            category_features, top_n, lambda_param=0.1)
74
75      # Ensure at least two distinct colors
76      results = []
77      color_set = set()
78      selected_positions = []
79      for idx in top_indices:
80          item = category_df.iloc[idx].to_dict()
81          if len(color_set) < 2 or item["baseColour"] in color_set:
82              results.append(item)
83              color_set.add(item["baseColour"])
```

```python
 84                 selected_positions.append(idx)
 85             if len(results) == top_n:
 86                 break
 87
 88     # Fill with randomized remaining indices for variety
 89     if len(results) < top_n:
 90         remaining_indices = [i for i in top_indices if i not in
                selected_positions]
 91         shuffle(remaining_indices)  # Randomize for diversity
 92         for idx in remaining_indices:
 93             item = category_df.iloc[idx].to_dict()
 94             results.append(item)
 95             if len(results) == top_n:
 96                 break
 97
 98     for item in results:
 99         item["image_url"] = f"/static/images/{item['id']}.jpg"
100         item["id"] = int(item["id"])
101
102     logger.info(f"Recommended items: {[item['id'] for item in results]}
            ")
103     logger.info(f"Item details: {[f'{item['articleType']} - {item['
            baseColour']}' for item in results]}")
104
105     novelty_score = inverse_popularity_score(results)
106     diversity_score = intra_list_diversity(results)
107     serendipity_score = serendipity_measure(results, ml_model.df.iloc[
            ml_model.id_to_index[target_id]].to_dict()) if target_id else
            0.0
108
109     logger.info(f"Recommended {len(results)} items with novelty: {
            novelty_score}, diversity: {diversity_score}, serendipity: {
            serendipity_score}")
110     return results, novelty_score, diversity_score, serendipity_score
111
112 def check_negative_constraints(target_item: dict, candidate_item: dict)
        -> bool:
113     """Check for incompatible combinations based on updated
            ARTICLE_TYPE_GROUPS."""
114     def get_group(article_type: str) -> Optional[str]:
115         for group_name, types in ARTICLE_TYPE_GROUPS.items():
116             if article_type in types:
117                 return group_name
118         return None
119
120     target_group = get_group(target_item["articleType"])
121     candidate_group = get_group(candidate_item["articleType"])
122
123     if target_group is None or candidate_group is None:
124         return True  # Assume compatible if not in any group
125
126     # General rule: items from the same group are incompatible, except
            for accessories
127     if target_group == candidate_group and target_group not in ["
            Accessories", "Jewellery", "Bags", "Makeup", "Skincare", "Bath
            and Body", "Haircare", "Fragrance", "Tech Accessories", "Home
            Decor"]:
128         return False
```

```python
129
130     # Specific rules
131     if target_group == "Trousers":
132         if candidate_group in ["Casual␣Shoes", "Formal␣Shoes"]:
133             if candidate_item["articleType"] in ["Sandals", "Flip␣Flops
                    "] and target_item["usage"] != "Casual":
134                 return False
135     if target_group in ["Shirts", "Tshirts", "Tops", "Dresses", "Suits"
            ]:
136         if candidate_group in ["Shirts", "Tshirts", "Tops", "Dresses",
                "Suits"]:
137             return False
138     if candidate_group == "Casual␣Shoes" and target_item["usage"] == "
            Formal":
139         if candidate_item["articleType"] in ["Flip␣Flops", "Sports␣
                Sandals"]:
140             return False
141     # Add more specific rules as needed
142
143     return True
144
145 def get_compatible_types(article_type: str) -> List[str]:
146     """Get␣compatible␣article␣types␣from␣COMPATIBLE_TYPES."""
147     return COMPATIBLE_TYPES.get(article_type, [])
148
149 def get_accessory_types(usage: str, season: str) -> List[str]:
150     """Get␣accessory␣types␣based␣on␣usage␣and␣season."""
151     accessory_list = ACCESSORY_COMBINATIONS.get(usage, []) +
            SEASONAL_ACCESSORIES.get(season, [])
152     return list(set(accessory_list))
153
154
155 def color_compatibility(color1: str, color2: str) -> float:
156     """Calculate␣color␣compatibility␣score␣using␣COLOR_COMPATIBILITY␣
            dictionary."""
157     if color1 == color2:
158         return 1.0
159     if color2 in COLOR_COMPATIBILITY.get(color1, []):
160         return 0.8
161     return 0.0
162 @app.get("/api/product/{item_id}", response_model=ProductPageResponse)
163 async def product_page(item_id: str):
164     """Get␣product␣details␣and␣outfit␣recommendations."""
165     product = get_item(item_id)
166     target_id = str(product["id"])
167     target_gender, target_usage, target_season = product["gender"],
            product["usage"], product["season"]
168     target_color, target_article_type = product["baseColour"], product[
            "articleType"]
169
170     target_features = ml_model.combined_features[ml_model.id_to_index[
            target_id]]
171     compatible_types_list = get_compatible_types(target_article_type)
172     accessory_types = get_accessory_types(target_usage, target_season)
173
174     recommendations_dict = {}
175     for compatible_type in compatible_types_list:
176         recs, _, _, _ = get_ml_recommendations(target_features,
```

```
                    compatible_type, target_gender, target_color, target_id)
177         filtered_recs = [item for item in recs if
                check_negative_constraints(product, item)]
178         if filtered_recs:
179             recommendations_dict[compatible_type] = [Item(**item) for
                    item in filtered_recs[:3]]
180
181     for accessory_type in accessory_types:
182         recs, _, _, _ = get_ml_recommendations(target_features,
                accessory_type, target_gender, target_color, target_id)
183         filtered_recs = [item for item in recs if
                check_negative_constraints(product, item)]
184         if filtered_recs:
185             recommendations_dict[accessory_type] = [Item(**item) for
                    item in filtered_recs[:3]]
186
187     return ProductPageResponse(product=Item(**product), recommendations
            =OutfitRecommendation(recommendations=recommendations_dict))
188
189 @app.post("/api/recommend-from-image", response_model=
        OutfitRecommendation)
190 async def recommend_from_image(file: UploadFile = File(...)):
191     """Recommend␣outfits␣based␣on␣an␣uploaded␣image."""
192     try:
193         contents = await file.read()
194         image = Image.open(BytesIO(contents)).convert("RGB")
195         attributes = predict_attributes(image)
196
197         synthetic_name = f"{attributes.get('gender',␣'Unisex')}'s␣{
                attributes.get('baseColour',␣'')}␣{attributes.get('
                articleType',␣'Fashion␣Item')}"
198         categorical_data = [attributes.get(col, "Unknown") for col in [
                "gender", "masterCategory", "subCategory", "articleType", "
                baseColour", "season", "usage"]]
199         onehot = ml_model.onehot_encoder.transform([categorical_data])
200         tfidf = ml_model.tfidf_vectorizer.transform([synthetic_name])
201         target_features = hstack([onehot, tfidf])
202
203         target_article_type = attributes.get("articleType", "Shirts")
204         target_gender, target_usage = attributes.get("gender", "Unisex"
                ), attributes.get("usage", "Casual")
205         target_season, target_color = attributes.get("season", "Summer"
                ), attributes.get("baseColour", "Black")
206
207         compatible_types_list = get_compatible_types(
                target_article_type)
208         accessory_types = get_accessory_types(target_usage,
                target_season)
209
210         recommendations_dict = {}
211         for compatible_type in compatible_types_list:
212             recs,_,_,_ = get_ml_recommendations(target_features,
                    compatible_type, target_gender, target_color)
213             filtered_recs = [item for item in recs[0] if
                    check_negative_constraints(attributes, item)]
214
215             if filtered_recs:
216                 recommendations_dict[compatible_type] = [Item(**item)
```

```
                           for item in filtered_recs]
217
218        for accessory_type in accessory_types:
219            recs,_,_,_ = get_ml_recommendations(target_features,
                   accessory_type, target_gender, target_color)
220            filtered_recs = [item for item in recs[0] if
                   check_negative_constraints(attributes, item)]
221
222            if filtered_recs:
223                recommendations_dict[accessory_type] = [Item(**item)
                       for item in filtered_recs]
224
225        return OutfitRecommendation(recommendations=
               recommendations_dict)
226    except Exception as e:
227        raise HTTPException(status_code=500, detail=f"Error processing
               image: {str(e)}")
```

**Interaction:** A complex interaction between multiple modules. Key interactions include:

- product_page and recommend_from_image call get_ml_recommendations to get recommendations.

- get_ml_recommendations uses ml_model.combined_features (precomputed feature vectors), cosine_similarity, and maximal_marginal_relevance for ranking.

- get_compatible_types and get_accessory_types (from constants.py) provide rule-based filtering.

- check_negative_constraints enforces additional compatibility rules.

## 5. Search

**Module:** search(query)

**Purpose:** Performs a semantic search using ChromaDB and OpenCLIP embeddings. This allows users to search for products using natural language queries. The results are returned as a list of image URLs and their corresponding distances (similarity scores).

**Code:**

```
1 @app.post("/api/search", response_model=SearchResult)
2 async def search(query: str = Form(...)):
3     """Search for images based on a text query."""
4     try:
5         fashion_collection = ml_model.chroma_client.get_collection("
              fashion", embedding_function=OpenCLIPEmbeddingFunction(),
              data_loader=ImageLoader())
6         results = fashion_collection.query(query_texts=[query],
              n_results=5, include=["uris", "distances"])
7
8         results["uris"] = [[uri.replace("/kaggle/input/fashion-product-
              images-dataset/fashion-dataset/", "") for uri in results["
              uris"][0]]]
9         image_data = [
10            {"id": results["ids"][0][i], "distance": results["distances
                 "][0][i], "image_url": f"/static/images/{os.path.
                 basename(results['uris'][0][i])}"}
```

```
11              for i in range(len(results["ids"][0]))
12              if os.path.exists(os.path.join(STATIC_DIR, "images", os.
                    path.basename(results["uris"][0][i])))
13          ]
14          return SearchResult(images=image_data)
15      except Exception as e:
16          print(f"Error during search: {e}")
17          raise HTTPException(status_code=500, detail="An error occurred 
                during search")
```

**Interaction:**

- Called by: Frontend (search bar).

- Interacts with: `ml_model.chroma_client` (ChromaDB).

## 6. Evaluation

**Module:** `evaluate_recommendations()`

**Purpose:** This module provides an endpoint to evaluate the performance of the recommendation system. It compares the ML-based recommendations with popularity-based and random baselines using novelty, diversity, and serendipity metrics.

**Code:**

```
1  def popularity_based_recommender(top_n=5):
2      popularity = ml_model.df['articleType'].value_counts().index[:top_n
           ]
3      return ml_model.df[ml_model.df['articleType'].isin(popularity)].
           sample(top_n).to_dict('records')
4
5  def random_recommender(top_n=5):
6      return ml_model.df.sample(top_n).to_dict('records')
7
8  def inverse_popularity_score(recommendations):
9      total_purchases = len(ml_model.df)
10     scores = [1 - (ml_model.df['articleType'].value_counts()[item['
           articleType']] / total_purchases) for item in recommendations]
11     return np.mean(scores)
12
13 def intra_list_diversity(recommendations):
14     features = [ml_model.combined_features[ml_model.id_to_index[str(
           item['id'])]] for item in recommendations]
15     stacked_features = vstack(features).toarray()
16     similarities = cosine_similarity(stacked_features)
17     return 1 - np.mean(similarities)
18
19 def serendipity_measure(recommendations, user_history):
20     user_features = ml_model.combined_features[ml_model.id_to_index[str
           (user_history['id'])]]
21     rec_features = [ml_model.combined_features[ml_model.id_to_index[str
           (item['id'])]] for item in recommendations]
22     distances = [cosine_similarity(user_features, rec.reshape(1, -1))
           [0][0] for rec in rec_features]
23     return np.mean(distances)
24
25 @app.get("/api/evaluate")
26 async def evaluate_recommendations():
```
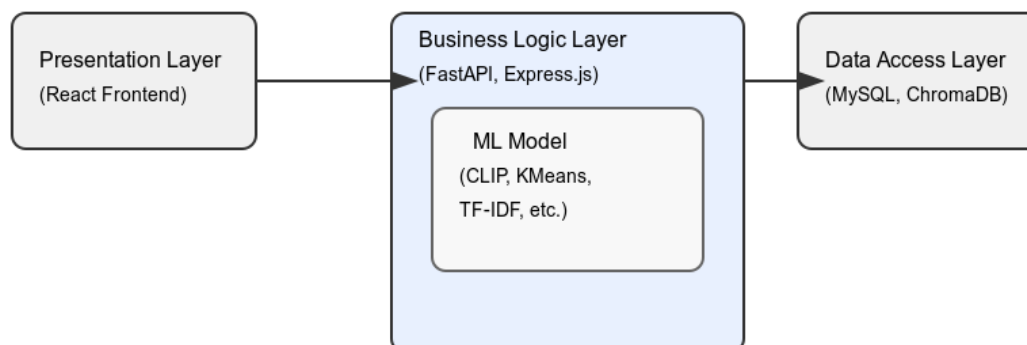
```
27      # Select a random product as the target
28      target_product = ml_model.df.sample(1).iloc[0]
29      target_id = str(target_product['id'])
30      target_features = ml_model.combined_features[ml_model.id_to_index[
            target_id]]
31
32      # Get recommendations from different methods
33      ml_recs, ml_novelty, ml_diversity, ml_serendipity =
            get_ml_recommendations(
34          target_features, target_product['articleType'], target_product[
                'gender'],
35          target_product['baseColour'], target_id
36      )
37
38      popularity_recs = popularity_based_recommender()
39      random_recs = random_recommender()
40
41      # Calculate metrics for baseline methods
42      pop_novelty = inverse_popularity_score(popularity_recs)
43      pop_diversity = intra_list_diversity(popularity_recs)
44      pop_serendipity = serendipity_measure(popularity_recs,
            target_product)
45
46      rand_novelty = inverse_popularity_score(random_recs)
47      rand_diversity = intra_list_diversity(random_recs)
48      rand_serendipity = serendipity_measure(random_recs, target_product)
49
50      return {
51          "ML Model": {
52              "Novelty": ml_novelty,
53              "Diversity": ml_diversity,
54              "Serendipity": ml_serendipity
55          },
56          "Popularity Baseline": {
57              "Novelty": pop_novelty,
58              "Diversity": pop_diversity,
59              "Serendipity": pop_serendipity
60          },
61          "Random Baseline": {
62              "Novelty": rand_novelty,
63              "Diversity": rand_diversity,
64              "Serendipity": rand_serendipity
65          }
66      }
```

**Interaction:**

- Called by: An external script or tool for evaluation.

- Interacts with: get_ml_recommendations, helper functions for calculating metrics.

## Interaction Diagram (Simplified)



**Flow:** A simplified representation of the core flow. Note that the diagram does not show every single function call, but rather the high-level interaction between layers and major components.

# Q2. Business Rules, Validation, and Data Transformation

## A. Business Rules

The application implements several business rules to ensure the quality and relevance of recommendations and to maintain data integrity. These rules are primarily enforced within the `get_ml_recommendations` function and related helper functions.

- **Compatibility Rules:** Ensures that recommended outfit components are compatible with each other. This is primarily handled by the `COMPATIBLE_TYPES` dictionary in `constants.py`, which defines allowed combinations of article types (e.g., "Shirts" are compatible with "Trousers", "Jeans", etc.). The `get_compatible_types` function retrieves these rules.

- **Gender, Usage, and Seasonal Filtering:** Recommendations are filtered based on the gender, usage, and season of the target item. This ensures that, for example, a formal shirt for men will not be recommended with casual shorts. This is implemented in `get_ml_recommendations` using a filtering mask and in the `get_accessory_types` function.

- **Accessory Rules:** The `get_accessory_types` function uses the `ACCESSORY_COMBINATIONS` and `SEASONAL_ACCESSORIES` dictionaries to filter the accessories based on the target item's usage and season.

- **Color Compatibility:** The `color_compatibility` function calculates a score based on the `COLOR_COMPATIBILITY` dictionary, preventing clashing colors in recommendations. Items with low color compatibility scores are filtered out.

- **Negative Constraints:** The `check_negative_constraints` function implements specific rules to prevent incompatible item pairings that might not be captured by the broader compatibility rules (e.g., preventing sandals from being recommended with

12

formal trousers, even if "Casual Shoes" are generally compatible with "Trousers"). This uses the `ARTICLE_TYPE_GROUPS` dictionary.

- **Diversity Enforcement (MMR):** The `maximal_marginal_relevance` function is used to ensure diversity among the top recommendations. This prevents the system from recommending very similar items (e.g., multiple shirts of the same color and style).

## B. Validation Logic

Validation is performed at multiple levels to ensure data integrity and prevent errors.

- **Input Checks:** Basic checks for required fields and data types are performed using Pydantic models in the API endpoints (e.g., `Item`, `OutfitRecommendation`, `ProductPageResponse`, etc.).

- **Dataframe-Level Validation:** The `get_item` function includes error handling (raising an `HTTPException`) if an item with the requested ID is not found in the preloaded DataFrame.

- **Empty Result Checks:** The `get_ml_recommendations` function includes checks to ensure that the filtered DataFrame is not empty *after* applying filters (e.g., compatibility, color, negative constraints). This prevents errors later in the recommendation process and returns an empty list if no suitable recommendations are found.

- **Logging:** The `get_ml_recommendations` function utilizes the `logger` object (from Python's `logging` module) extensively. This provides valuable insights into the filtering process, including:

  - The number of items remaining after each filtering step (initial filter, color filter).
  - The recommended item IDs and details (article type and color).
  - The calculated novelty, diversity, and serendipity scores.
  - Warnings if no items are found for a particular article type.

  This detailed logging is crucial for debugging, understanding the recommendation process, and identifying potential issues (e.g., overly restrictive filters).

## C. Data Transformation

Data transformation is a crucial part of the application, converting raw data into formats suitable for the ML model and the frontend.

- **SQL Result (Database) → Pandas DataFrame → Python Dictionary:** The `load_data` function retrieves data from the MySQL database using SQLAlchemy and converts it into a Pandas DataFrame. The `get_item` function then converts a single row of this DataFrame into a Python dictionary, making it easier to work with. This transformation also includes adding the `image_url` field.

- **One-Hot Encoding and TF-IDF Vectorization:** The `preprocess_data` function transforms categorical features (gender, masterCategory, etc.) into numerical vectors using one-hot encoding. It also transforms the product display name into a TF-IDF

vector. These numerical representations are essential for the cosine similarity calculations used in the recommendation engine. The combined features are stored in `ml_model.combined_features`.

- **CLIP Output → Predicted Attributes:** The `predict_attributes` function uses the CLIP model to predict attributes from an image. The raw output of the CLIP model (logits) is converted into predicted labels (e.g., "Men", "Summer", "Casual") using softmax and argmax. The dominant color is extracted and mapped to the closest color name using a predefined mapping.

- **Image Path → URL:** The `get_item` function and other recommendation functions construct the image URL (`/static/images/{item_id}.jpg`) from the item ID, which is used by the frontend to display the product images.

- **Implicit Transformations within Similarity Calculations:** The `cosine_similarity` function itself performs a transformation. It takes the pre-processed, numerical feature vectors (one-hot encoded and TF-IDF) and calculates the cosine similarity between them. This similarity score is then used to rank the recommendations. The `maximal_marginal_relevance` function builds upon this, adding a diversity component to the ranking.

- **ChromaDB Query Results:** In the `search` function, the results from ChromaDB (which include IDs, distances, and URIs) are transformed into a list of dictionaries, each containing the ID, distance, and a constructed `image_url`. This makes the data suitable for the `SearchResult` Pydantic model and, consequently, for sending as a JSON response to the frontend.