# LITERATURE REVIEW

## AUTOMATED TAGGING AND DESCRIPTION FOR ECOMMERCE PRODUCTS

### TEAM MEMBERS

### VASUDHA RANI PATHEDA-(21MIS7121)

### GALETI YASWANTH SAI-(21BCE9674)

### S CHIRANJEEVI-(21BCE8950)

## Important Points and Limitations

1. **Automatic tagging and retrieval of E-Commerce products based on visual features.[1]**
   This paper present an approach of automatically tag-ecommerce product images based on the visual characters. Utilizing deep convolution neural network (CNNs), and the proposed method also applies inverse distance-weighted K-nearest neighbour classifiers to allocate tags and constructs a product retrieval system based on these tags. The authors tested the system using the Amazon product dataset and achieved promising results across categories like apparel, electronics, and sports equipment.
   **Limitations:**
   a. No Colour Feature Integration: Lacks extraction of dominant colours, missing key product attributes for tagging.
   b. Basic Feature Extraction: Relies on CNNs, missing detailed object type, shape, and texture features.
   c. No LLM for Description Generation: Doesn't use an LLM to generate product descriptions from extracted features.
   d. Simplistic Classifier: Uses KNN, which underutilizes combined visual and contextual data for classification.
   e. Limited Textual Attribute Focus: Misses generating product names and descriptions based on visual and colour features.

2. **Application of Improved k-means Algorithm in E-commerce Data Processing .[2]**
   The above-mentioned system, "Application of Improved K-means Algorithm in E-commerce Data Processing," enhances the quality of recommendations using the improved K-means algorithm and integrates it with genetic algorithms and SVD++. The test conducted on the Taobao dataset attained 85% precision, 87% recall, and an AUC of 0.83 values on this data set. Moreover, the optimized time for computation improves it in its performance of processing large databases, but these improvements have their price. Its heavy memory usage severely hinders its scalability in commercial use, which is a prime factor when efficiency is considered.

   Limitations:

   a. This work's enhanced K-means algorithm is a memory consumer - certainly not that efficient for large-scale real-time applications of e-commerce. My project lies on color extraction through K-Means with coupling on top of object detection in Grounding DINO. Feature extraction processes are optimized by extracting what it really cares for to be classified, i.e., colors and objects.

   b. Limited Evaluation of Visual Features: The authors are strengthening their recommendations rather than the user data but not digging further into the visual

feature of products. My approach would include unsupervised color extraction with K-Means and grounding DINO object feature extraction to ensure that the saliency of dominant colors along with parts of the product complementing its shape would boost its recognition and recommendations abilities.

3. **Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection.[3]**
The project described in the document proposes **Grounding DINO**, an open-set object detection model. Grounding DINO integrates the Transformer-based DINO detector with grounded pre-training to detect arbitrary objects specified by human inputs, such as category names or referring expressions. It excels in generalizing across unseen object categories through a tight fusion of language and vision modalities. The model is pre-trained on large datasets and evaluated on benchmarks like COCO, LVIS, and ODinW. Grounding DINO achieves impressive performance, including a new record on the ODinW zero-shot benchmark.
**Limitations:**
a. Fine-Grained Segmentation: Grounding DINO struggles with detecting small or intricate product features. By enhancing Grounding DINO's ability to focus on key areas and refining the extraction process, this project mitigates the limitations in fine-grained segmentation, ensuring that critical product details are captured effectively.
b. Limited Training Data: Grounding DINO's performance drops when encountering rare or novel object categories without extensive fine-tuning. This project strengthens Grounding DINO's feature extraction through targeted fine-tuning, allowing the model to recognize product-specific features even when data is limited or categories are rare.
c. False Positives and Hallucinations: Grounding DINO may produce false positives in dense scenes or complex backgrounds.The approach reduces false positives by enhancing bounding box accuracy and focusing Grounding DINO on key product features, ensuring better detection precision in challenging scenes.

4. **Multi-Feature Extraction from Product Images Using Deep Learning and Image Processing Techniques.[4]**
This paper discusses a method for extracting multiple features such as colour, texture, and shape from product images using a combination of deep learning models and traditional image processing algorithms. Convolutional Neural Networks (CNNs) are used for high-level feature extraction, while Gabor filters are applied to capture texture details. The combination of these features improves the accuracy of image-based product categorization and retrieval. The method was tested on a large dataset of product images and showed improvements in retrieval speed and accuracy across diverse categories.
**Limitations:**
a. **Over-reliance on CNN for Feature Extraction**: The approach depends heavily on CNNs, potentially missing finer details like intricate textures and edge information.
b. **Limited Integration of Domain-specific Features**: The model does not account for domain-specific visual features, which may affect the accuracy for niche product categories.
c. **No Integration of Colour Histograms**: The method lacks the use of colour histograms, which could enhance colour-based product categorization.
d. **Limited Use of External Data Sources**: Does not integrate additional metadata (e.g., tags or descriptions) that could further refine product categorization.
e. **Single-channel Processing**: Focuses primarily on image features without incorporating multimodal inputs like textual data for enhanced accuracy.

5. **Texture-Based Feature Extraction for Product Image Categorization.[5]**
This research focuses on texture-based feature extraction to improve the categorization of eCommerce product images. The approach employs Grey Level Co-occurrence Matrix (GLCM) and Local Binary Patterns (LBP) for extracting detailed texture features from product images. By combining these texture features with shape information obtained from boundary detection techniques, the proposed method achieves a higher classification accuracy for complex categories such as fabrics and accessories. The method was validated on a dataset of fashion and apparel images.
**Limitations:**
a. **Low Performance on Non-textured Products**: The approach is highly dependent on texture information, which may not be as useful for smooth, non-textured products like electronics or home goods.
b. **Computational Complexity**: The combination of GLCM and LBP increases the computational cost, limiting its applicability for large-scale datasets.
c. **Absence of Colour Feature Integration**: The method does not account for colour information, which is crucial for differentiating products in certain categories like fashion or home decor.
d. **Limited Generalization**: The method may not generalize well to categories beyond those with significant texture variation, such as accessories or furniture.
e. **No Deep Learning Involvement**: The reliance on traditional image processing methods may fall short in terms of adaptability and scalability compared to deep learning approaches.

6. **Automatic Product Description Generation Using Transformer-based Models.[6]**
This research explores the generation of automated product descriptions using transformer models like GPT-3 and BART. The approach trains the model on large datasets of product details to automatically create concise and informative product descriptions. The paper outlines a method for fine-tuning pre-trained models to generate descriptions based on product attributes such as title, price, and key features. Evaluations showed improved readability and relevance in comparison to traditional template-based generation methods.
**Limitations:**
a. **Repetitive Descriptions**: Generated descriptions can sometimes be repetitive, especially for similar product categories, reducing their uniqueness.
b. **Lack of Domain-Specific Knowledge**: Models may fail to incorporate highly technical or domain-specific knowledge, leading to generic descriptions for specialized products.
c. **Data Bias**: The model's training data can introduce bias, resulting in product descriptions that favour certain products or brands.
d. **Limited Customization Options**: The system offers limited control for retailers to adjust the tone, style, or length of the generated descriptions.
e. **Dependence on High-Quality Input Data**: The accuracy of the generated descriptions depends on the quality of input product attributes, meaning poorly labelled data can lead to irrelevant or incorrect descriptions.

7. **A Multimodal In-Context Tuning Approach for E-Commerce Product Description Generation [7]**
This paper proposes a new approach termed Multimodal In-Context Tuning (ModICT), which tunes the automatic generation of product descriptions from images using marketing keywords. Converse to traditional methods, ModICT overcomes the generic and often inaccurate descriptions by means of in-context learning-either by referring to similar samples of a product for actual generation. Encoders of visual and language were frozen, while emphasis was placed on optimizing modules responsible for creation of in-context references

and dynamic prompts. The strategy allows one to improve both the diversity and accuracy of product descriptions, as evidenced by experiments on three categories of E-commerce products:. For example, ModICT allows improving up to 3.3% accuracy (Rouge-L) and up to 9.4% diversity (D-5), thus showing a prospective ability to refine the generation of product descriptions for better use in practical applications.

**Limitations:**

**a. Generic Descriptions:** The previous work had explained that the descriptions for products were broadly generic, because again the same category products share the same type of description. However, your approach is trying to ground out more specific features like colors and object-oriented information making use of Grounding DINO as well as clustering by using K-Means where description is going to be much more specific and not so similar.

**b. Overreliance on Common Words:** The old approach had the disadvantage that the concentration was so much on common words that the models overlooked the unique product features. This project will be targeted on more focus on feature extraction by bounding boxes and color clustering, thus much more attention to product-specific attributes that lead to correct feature-based product predictions.

**c. Fixed Language Model Framework:** ModICT still leaves something in the bucket; it freezes the visual encoder and language model, and then little is left to fine-tune. Your approach allows for dynamic learning within the supervised classifier using the features extracted and, therefore, leaves room for more flexibility and improvement with time.

8. **Extracting Related Images from E-commerce Utilizing Supervised Learning**
This paper presents a Siamese deep convolutional neural network (CNN) in order to learn embeddings that can determine the visual similarity between images. It well distinguishes between similar and dissimilar objects by training with positive and negative image pairs. A new angular loss metric is proposed to efficiently measure the loss across multi-dimensional space in order to make the comparisons of embeddings more accurate. Finally, it integrates both low and high-level embeddings to produce a final image representation. In addition, the fractional distance matrix is applied to compute the distances of the embeddings so that the model may be able to make more precise computations. The architecture is examined on four datasets other deep CNN models didn't perform well with tasks related to image retrieval as well as fine-grained images comparison. This proved that the proposed network was better at its task of visual similarity capture.

**Limitations:**

**1. Resource-Intensive Architecture:**
The use of VGG19 for feature extraction and the use of multiple CNNs for triplet image processing makes this model highly memory-intensive and computationally expensive. However, by integrating Grounding DINO for feature extraction instead of VGG19 in my project, it concentrates on extracting object-specific features and bounding boxes, thus reducing the memory overhead but maintains their corresponding accuracy.

**2. Less Attention on Colour Features:**
The paper fundamentally relies on structural features through image similarity and puts no emphasis on extracting or using color information. My approach integrates unsupervised K-Means clustering to extract dominant colors from images in ensuring that color and object features contribute to classification and recommendation, with improved ability in the system to handle e-commerce images.

**3. Scalability Issues:**
For a model that combines deep CNNs with weak and poor real-time scalability especially when processing large datasets or real time applications such as recommending, my project has optimised the both feature extraction and classification using a combination of

unsupervised learning and supervised learning to make recommendations more efficient and fast in balancing accuracy efficiency.

9. **Recognize Anything: A Strong Image Tagging Model**

We present the Recognize Anything Model, as we refer to a state-of-the-art foundation model for image tagging with substantial zero-shot capabilities towards recognizing plain categories without requiring hand annotations. It does not utilize standard, hand-labeled datasets. Instead, it utilizes a tremendous amount of image-text pairs in its training process. Develop A unified captioning and tagging model Train the model under a four-step process: automated text semantic parsing for tag generation, followed by training the unified captioning and tagging model, and finally through a data engine that cleans and refines annotations in order to generate high-quality data. Through such processes, RAM has succeeded in terms of both accuracy and scope and has outperformed models such as CLIP and BLIP while being almost at par with some fully supervised models.

**Limitations:**

a. Dependance on Textual Data for Recognition: RAM depends hugely on large scale image-text pairs. Their training may not even allow them to reproduce actual visual characteristics of objects, especially those with complex textures and colors. Your approach goes directly to look at a product image at the pixel level through analysis based on unsupervised color extraction with Grounding DINO's feature extraction.

b. Low sensitivity to feature specificity: While being excellent in general, the recognition ability of RAM may not be sensitive to details regarding product features, such as textures and color patterns or parts of the object. In that case, with Grounding DINO focus on extracting dominant colors and detailed features of objects, this requirement for the sensitivity of feature specificity will certainly be met.

c. No Explicit Color Features Extraction: RAM does not explicitly extract color features, which may be a key aspect for products that nearly look similar in their visual aspect. Your project will fill that gap by including the K-Means clustering to extract color features so that the model can better differentiate between such products based on their color feature.

**Comparative Analysis:**

| Sno | Paper Title | Methodology | Datasets Used | Performance Metrics | Advantages | Disadvantages |
|---|---|---|---|---|---|---|
| 1 | Automatic tagging and retrieval of E-Commerce products based on visual features. | Uses VGG-19 pre-trained on ImageNet to extract high-level features from images. Employs Weighted K-Nearest Neighbors (KNN) for multi-label tag assignment based on the inverse distance of neighbouring images. Feature vectors are extracted from the final fully connected layers of VGG-19 and tags are assigned based on weighted averages of tag presence among K nearest neighbours. Tags are stored for fast retrieval. | Amazon e-commerce dataset with images and metadata for apparel, clothing, electronics, and sports equipment categories. | Precision, Recall, F1-score for multiple K values. Metrics were evaluated on different categories of products. | Effective multi-label tagging approach using transfer learning (VGG-19) and weighted KNN. Scales well to large datasets. Reduces training time using transfer learning. Efficient for fast product retrieval. | K-Nearest Neighbours-based tag assignment can be computationally expensive when scaling to very large datasets. No discussion on handling highly noisy labels. Performance varies with K values. |
| 2. | Application of Improved K-means Algorithm in E-commerce Data Processing | The paper introduces enhanced K-means by employing genetic algorithms and the coefficient of variation method for the clustering of e-commerce data on product content with a view to further improving the recommendation. Hidden features of SVD++ are extracted from the data. | Taobao e-commerce dataset (5842 users, 6447 items, 827,384 user rating records). | Performance metrics: Precision (85%), Recall (87%), AUC (0.83), Average computation time (54.2s). | High accuracy and computational efficiency in recommendations. Successfully identifies complex user-product interactions. | Higher memory consumption compared to other models (ISVD++_I_k-means: 16.8 MB). Model consumes more computational resources. Limited evaluation in real-world commercial applications. |
| 3. | Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection | Dual-encoder-single-decoder architecture: Uses Swin Transformer for multi-scale image feature extraction, BERT for text feature extraction, and a cross-modality feature fusion with a feature enhancer. Applies a language-guided query selection method to identify relevant queries for object detection, and a | COCO, LVIS, ODinW, O365, GoldG, RefCOCO/+/g | AP (Average Precision) for zero-shot, fine-tuning, and full-shot performance across datasets. | Effective in zero-shot and few-shot learning scenarios, scalable to larger datasets. Achieves state-of-the-art results on COCO and LVIS benchmarks. | High computational cost, especially during training. Performance on rare categories is lower without additional training data. |

| # | | | | | | |
|---|---|---|---|---|---|---|
| | | cross-modality decoder for refinement. | | | | |
| 4. | Multi-Feature Extraction from Product Images Using Deep Learning and Image | Combines CNNs for high-level feature extraction (colour, shape, texture) with Gabor filters for texture detection. Emphasizes deep learning-based image classification with traditional processing techniques. | Large-scale dataset of product images from various eCommerce categories. | Accuracy, retrieval speed, feature extraction precision | Improves accuracy and retrieval speed for diverse product categories. Leverages deep learning for comprehensive feature extraction. | Over-reliance on CNN, missing finer textures and edge details. Lacks integration of multimodal inputs like text for enhanced classification. No colour histogram integration. |
| 5. | Texture-Based Feature Extraction for Product Image Categorizatio n | Uses Grey Level Co-occurrence Matrix (GLCM) and Local Binary Patterns (LBP) for texture-based feature extraction. Combines texture features with shape information using boundary detection techniques. | Fashion and apparel images dataset, focusing on high-texture products. | Classification accuracy, texture feature precision | Effective for texture-rich product categories like fabrics and accessories. High accuracy in distinguishing texture-based categories. | Poor performance on non-textured products like electronics. High computational cost due to the combination of GLCM and LBP. No colour feature extraction. |
| 6 | Automatic Product Description Generation Using Transformer-based Models (Zhang & Xu, 2021) | Uses transformer models (GPT-3, BART) to generate product descriptions based on product attributes like title, price, and features. Fine-tunes pre-trained models on product datasets. | Large datasets of product descriptions and attributes from eCommerce platforms. | Readability, relevance of descriptions, uniqueness | Produces relevant and readable descriptions that improve on template-based methods. Suitable for large-scale eCommerce platforms. | Prone to generating repetitive descriptions for similar products. High dependency on well-labeled input data. Lacks domain-specific knowledge, leading to generic descriptions for technical products. |

| 7. | A Multimodal In-Context Tuning Approach for E-Commerce Product Description Generation | The paper introduces multimodal in-context tuning using pre-trained models like CLIP for the encoding of images and large language models like GPT for the processing of text. ModICT freezes most of the LM and the visual encoder and allows the tuning of only a small part of the parameters for task-specific learning, hence allowing relatively very few resources for efficient finetuning. The approach takes cross-modal in-context examples to guide the generation of product descriptions by combining image features and text inputs, thus making it adaptive to diverse multimodal tasks. | MD2T (Chinese E-commerce product summarization corpus) | BLEU-4: 34.2, ROUGE-L: 48.9, BERTScore: 85.6, Diversity (D-n): 91.4 | Reduces computational cost with fewer learnable parameters, adaptable to various large language models (LLMs), improves diversity in product descriptions by considering both text and images. | Requires a pre-trained CLIP model for image encoding, and the method is highly dependent on the quality of in-context references. |
|---|---|---|---|---|---|---|
| 8 | Extracting Related Images from E-commerce Utilizing Supervised Learning | The architecture of Siamese network is used in the study for learning embeddings from triplets, which are made up of a positive, a negative image along with its anchor. Deep CNN (VGG19) is taken as the feature extractor and is further learned through contrastive loss that takes it toward good-quality embeddings. The entire model uses some in-house training loops on TensorFlow for optimizing and validation purposes. | Fashion-MNIST, CIFAR-10, Exact Street2Shop dataset, Triplet dataset | Accuracy: 94.19% on validation set, Precision and loss metrics based on triplet image pairs. | High accuracy on the validation set, effective use of VGG19 for feature extraction, efficient learning using the Siamese network to learn visual similarity embeddings. | Memory-heavy due to VGG19 architecture, requires multiple CNNs to handle complex triplet data. Real-time deployment may face performance issues due to high resource consumption. |
| 9. | Recognize Anything: A Strong Image Tagging Model | The RAM architecture combines three main components: an image encoder, an image-tag recognition decoder, and a text generation encoder-decoder. The | COCO, OpenImages V6, ADE20k, Conceptual Captions, SBU Captions, Visual Genome, | COCO: mAP = 64.1%, Top-1 Accuracy = 70.4%; OpenImages: | Open-vocabulary recognition, can generalize to unseen categories, high tagging | Performance may degrade with large-sized categories, high inference time for large images. |

| | | image encoder extracts visual features from the input image using a pre-trained backbone, often based on a Vision Transformer (ViT) or Convolutional Neural Networks (CNNs). The image-tag recognition decoder assigns descriptive tags to detected objects using an open-vocabulary mechanism. This allows the model to recognize not just predefined categories but also new, unseen objects. A separate text generation encoder-decoder further refines the textual output, improving the quality of tag-based descriptions. The training process employs an asymmetric loss function to deal with the imbalance between seen and unseen categories, enhancing model generalization. The methodology also focuses on multi-modal training that links images and textual descriptions, helping the model align visual features with language representations. | Conceptual 12M | mAP = 61.3%, Top-1 Accuracy = 68.7%; Visual Genome: mAP = 65.2%, Top-1 Accuracy = 72.0% | accuracy | |

**References:**

1. Sharma, V., & Karnick, H. (2016, June). Automatic tagging and retrieval of E-Commerce products based on visual features. In *Proceedings of the NAACL Student Research Workshop* (pp. 22-28).
2. Chen, W., & Wang, Q. (2024). Application of Improved k-means Algorithm in E-commerce Data Processing. Informatica, 48(11).

3. Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., ... & Zhang, L. (2023). Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*.
4. Gupta, A., & Singh, P. (2020). Multi-Feature Extraction from Product Images Using Deep Learning and Image Processing Techniques. *Journal of Computer Vision*, 45(3), 123-139.
5. Kim, H., & Park, S. (2019). Texture-Based Feature Extraction for Product Image Categorization. *International Journal of Pattern Recognition and Artificial Intelligence*, 33(2), 101-115.
6. Zhang, Y., & Xu, J. (2021). Automatic Product Description Generation Using Transformer-based Models. *Advances in Information Retrieval*, 44(1), 112-125.
7. Li, Y., Hu, B., Luo, W., Ma, L., Ding, Y., & Zhang, M. (2024). A Multimodal In-Context Tuning Approach for E-Commerce Product Description Generation. *arXiv preprint arXiv:2402.13587*.
8. Rajest, S. S., Sharma, D. K., Regin, R., & Singh, B. (2021). Extracting Related Images from E-commerce Utilizing Supervised Learning. *Innovations in Information and Communication Technology Series*, *1*, 34-46.
9. Zhang, Y., Huang, X., Ma, J., Li, Z., Luo, Z., Xie, Y., ... & Zhang, L. (2024). Recognize anything: A strong image tagging model. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 1724-1732).