# BT5153 Group Project Report
# E-Commerce Product Matching

**Group 14**

Ji Lin Cheng (A0231886W)  |  Li Qiaojun (A0236012X)  |  Liu Binghui (A0232064R)

Ye Li (A0043499U)  |  Zhao Heng (A0236342L)

**Github Link:** https://github.com/Lesterhuihuihui/BT5153-Group-14-Shopee-Matching-Guarantee

## Abstract

In e-commerce scenarios, recommendation systems based on hundreds of thousands of product information are common and necessary. And product matching is an important feature of a recommendation system. In this report, we study the product matching algorithms with the product data alone, which is common in cold start. The data is product information from a popular e-commerce platform, Shopee, in 2021. In our study, we explore the textual and graphic data of products, and build models based on the matching performance of different products. Our matching methods are evaluated to be applicable to real e-commerce platforms.

**Keywords**: e-commerce, product matching, image transformation, text vectorization

## 1.  Introduction

E-commerce, the activity of buying and selling products over the internet, has grown rapidly in the 21st century, supported by the global digital revolution. After decades of increases and evolving, a study published by the IMF (International Monetary Fund) indicated that there's still a high growth potential of e-commerce in Asia[1]. In this growing e-commerce sector, there have been many online platforms that have appeared in the market, such as Amazon, Ebay, Taobao, Shopee and Lazada. These platforms have since employed different methods to improve their e-service quality to acquire and retain customers. One of the popular features on the platforms is the recommender system that suggests products to the customers.

To construct a good recommender system, it is important to identify and group similar products from the large pool of product data on the platforms and feed them into the recommendation filtering algorithms. However, on the online platforms, sellers are free to add the product name and upload any image for product display where different images of similar goods may represent the same product or completely different products. To avoid the misrepresentation issues where different items are recommended to the customers as similar ones, machine learning models could be used to analyze the text and image of the products and correctly match and group them. Thus, we have identified a Kaggle dataset of consumer products from the ecommerce platform Shopee for our study and development of the machine learning models to compare product similarity and predict the product groupings.

## 2.  Project Objective

We aim to discover the problem of product similarity on online platforms. In e-commerce, the information of products, such as product names and illustrations, often varies among different retailers to attract customers. Thus in recommendation system design, recognizing the similarity between products by their basic data can enhance the overall efficiency by providing more suitable options for users. We aim to build predictive models based on Shopee product data to recognize identical or similar products with different names and illustrations.

In the following sections of the report, we first explore our dataset, followed by the methodologies we have adopted, and finally our model results and the business applications that can be applied.

## 3.  Dataset

### 3.1  Data Overview

We obtained our dataset from a Kaggle competition, that is Shopee Price Match Guarantee held in April 2021. The dataset contains two parts, one is a csv file with 34,250 records, the other is the set of 32,412 images corresponding

---

[1] Kinda, Tidiane. "E-commerce as a Potential New Engine for Growth in Asia." *IMF Working Papers*, 1 July 2019.

to the csv file records. The data fields available in the csv file are summarized in Table 1.

| Column Name | Data Type | Description |
|---|---|---|
| posting_id | String | Product posting's unique id |
| image | String | Product image's name |
| image_phash | String | Perceptual hash, acts as a fingerprint |
| title | String | Product's title |
| label_group | Integer | Product's grouping id |

*Table 1.* Summary of csv file data.

For each posting, it has an unique posting_id, a title for describing the product, an image name that links to the actual image file in the image dataset, and a label_group id that indicates the product's group belonging. For identical or similar product postings, they will be grouped together with the same label_group id, and this will then be used in the models for matching the similar products. Some examples of the posting are shown in Figure 1.



| | posting_id | image | image_phash | title | label_group |
|---|---|---|---|---|---|
| 0 | train_129225211 | 0000a68812bc7e98c42888dfb1c07da0.jpg | 94974f937d4c2433 | Paper Bag Victoria Secret | 249114794 |
| 1 | train_3386243561 | 00039780dfc94d01db8676fe789ecd05.jpg | af3f9460c2838f0f | Double Tape 3M VHB 12 mm x 4,5 m ORIGINAL / DO... | 2937985045 |
| 2 | train_2288590299 | 000a190fdd715a2a36faed16e2c65df7.jpg | b94cb00ed3e50f78 | Maling TTS Canned Pork Luncheon Meat 397 gr | 2395904891 |
| 3 | train_2406599165 | 00117e4fc239b1b641ff08340b429633.jpg | 8514fc58eafea283 | Daster Batik Lengan pendek - Motif Acak / Camp... | 4093212188 |
| 4 | train_3369186413 | 00136d1cf4edede0203f32f05f660588.jpg | a6f319f924ad708c | Nescafe \xc3\x89clair Latte 220ml | 3648931069 |

*Figure 1.* First 5 postings

For the image dataset, all images are saved as jpg files and the file size ranges from 2 KB to 1.3 MB. Some images only show the product itself, while some show the product on human models, and some also show with text printed on the image. Examples of the images are shown in Figure 2.
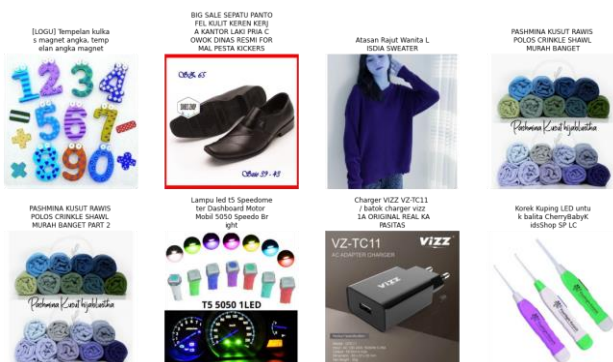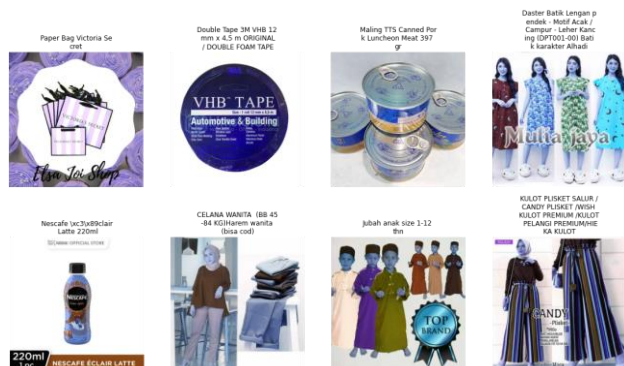




*Figure 2.* Random images display

## 3.2 Data Exploration

Within the csv dataset, we have noticed that all posting_id are unique with a total of 34,250 postings. However, the product images posted are not all unique. Of the 34,250 postings, only 32,412 unique images are found, indicating that there are 1,838 postings with duplicated images. When grouping the postings with duplicated images, we have noted that 1,246 unique images have been used for more than one posting. In Figure 3, the top 10 most used images are displayed with the number of usage.
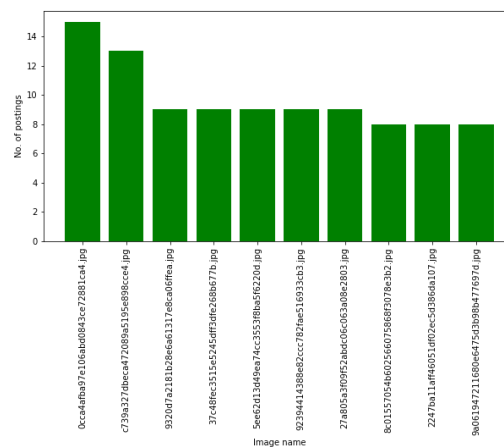


*Figure 3.* Top 10 most used image

The most used image (Figure 4) has been appearing in 15 postings with different product titles. All product titles contain the words 'bubble wrap', which is the product shown in the image, although some misspelled it as 'buble wrap' or 'bubble warp'. And each of the product titles has put in additional words to describe the product. Of the 15 postings, we have also observed that they're not all put in the same label_group, 12 of them are grouped under 4198148727 and 3 are under 2403374241 (as shown in Table 2).

0cca4afba97e106abd0843ce72881ca4.jpg



*Figure 4.* Most used image

| posting_id | title | lebel_group |
|---|---|---|
| train_2514153495 | Tambahan Bubble wrap / Plastik Bubble Pelindun... | 4198148727 |
| train_1236710293 | Bubble Warp Pengaman Pengiriman | 4198148727 |
| train_419018435 | Tambahan Bubble Wrap | 4198148727 |
| train_1381575164 | Bubble Wrap untuk ekstra packaging | 4198148727 |
| train_2420615645 | BUBBLE WRAP - EXTRA PACKING UNTUK BARANG ANDA | 4198148727 |
| train_3068759534 | BUBBLE PACK UNTUK PACKING TAMBAHAN 1BUBBLE UNT... | 4198148727 |
| train_2951822530 | Tambahan Extra Bubble Wrap Pengaman Packingan | 4198148727 |
| train_493140267 | PACKING TAMBAHAN BUBBLE WRAP | 4198148727 |
| train_1437764574 | Extra Bubble Wrap | 4198148727 |
| train_584097694 | bubble wrap - BUBLE WRAP | 4198148727 |
| train_3993385953 | Extra Bubble Wrap Pengaman Packingan | 2403374241 |
| train_1049463374 | BUBBLE WARP | 2403374241 |
| train_443869273 | Bubble Wrap | 4198148727 |
| train_4226152332 | EXTRA BUBBLE WRAP UNTUK PACKING | 2403374241 |
| train_2085280992 | Buble Wrap | 4198148727 |

*Table 2.* Postings using the most used image.

Within the dataset, the total of 34,250 postings are grouped into 11,014 distinct labels, with the largest label group comprising 51 postings and smallest label group

comprising only 2 postings. Figure 5 shows that most of the label groups contain less than 10 postings and only a small portion of the label groups contains a double digit number of postings.
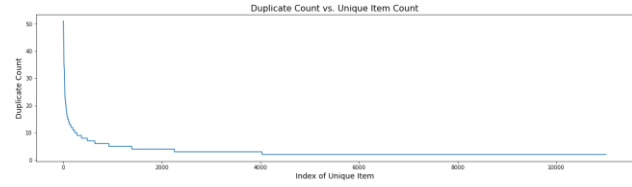


*Figure 5.* Label_group on postings count

A closer look at the most postings' label groups, we can see that the top 50 groups all have more than 20 postings within the group (Figure 6). And the top 7 groups all contain 51 postings. They're lip tint, serum, hair band, moisturizer, face tonic, faucet filter and soap respectively (Figure 7), which could also indicate that most sales on the platform are in the beauty and household categories.
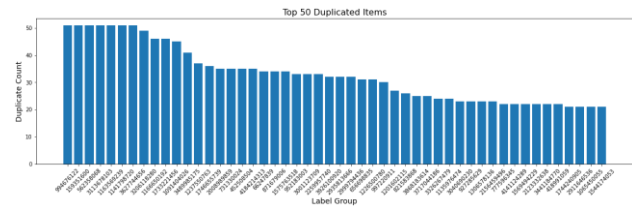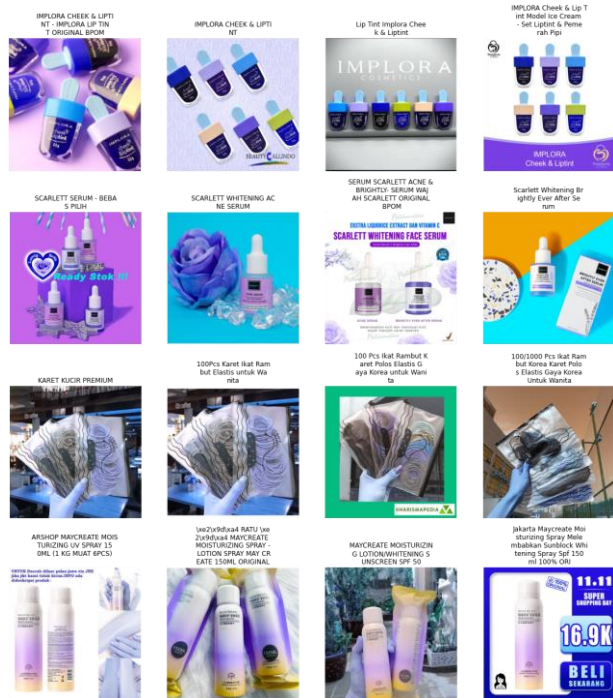


*Figure 6.* Top 50 label groups

*Figure 7.* Images of the top 7 label groups

To perform text analysis on the product titles, we have cleaned up the title text by removing the English stopwords and the punctuations to only keep the text that is meaningful for product description. Then a word cloud is used to get us an initial view of the frequently used words for product description and identification. As shown in Figure 8 word cloud, there are similar high frequency words in the product titles. Words like 'original', 'premium' are used by the sellers to describe the products as authentic, and Malay words like 'murah' (cheap) and 'termurah' (cheapest) are used by the sellers to position the product as budget bargains.



*Figure 8.* Word cloud of title words

We also counted the words used in the product titles and observed that the title length ranged from 1 to 61 words. And most of the postings have the title length between 3 to 15 words (distribution is shown Figure 9) with the mode at length of 6 words.
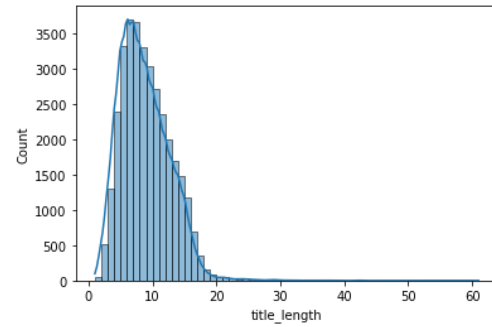


*Figure 9.* Product title length distribution

## 4. Data Preprocessing

As explored in the previous parts. Our dataset contains numerous image data and text data. To make our model more robust, we need to conduct critical data processing for both parts. For image data, all our images have corresponding hash value, helping me perform matching score computation. However, for loads of text data, we must apply some cleaning techniques to it.

As shown in the figure, our data contains many numbers, punctuations, and consecutive characters, which will disturb our data training. Therefore, we dropped all the alphanumeric tokens, fixed consecutive characters, and applied regex to filter non-alphabetic content from tokens.
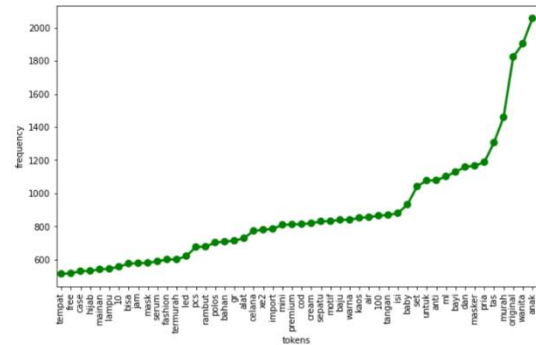


*Figure 10.* Text Frequency

## 5. 1st Baseline Clustering Solution

To compare with our model, we build several simpler baseline models based on clustering.

### 5.1 Model Introduction

DBSCAN is the abbreviation of Density-Based Spatial Clustering of Applications with Noise. The model finds core data points with high density and expands clusters from them. Therefore, the model is suitable for clustering similar products, which form clusters of similar density.

Compared with a more popular clustering model, K-nearest Neighbor, DBSCAN has better performance when there are outliers and when the samples are sparse in space. Also, DBSCAN does not require setting of cluster numbers, which requires much tuning in K-nearest Neighbor.

### 5.2  Data Preparation

To build the model, we only use the image data and apply embedding through EfficientNet to vectorize the images into high dimensional data.

EfficientNet is a convolutional neural network model which transforms graphic or textual data into lower dimensions without losing most of the information. The model is also a transfer learning model, which enables us to apply pretrained model weights to the embedding model. The model we use is the EfficientNet-B0, which is a popular image embedding pretrained model.

In image embedding, to make the model run in the main memory, we load data in chunks and process images in batches of 32 to generate images as matrices. Then we embed the image matrix with EfficientNet to generate data for modeling.

### 5.3  DBSCAN Clustering

We build a DBSCAN model to cluster similar products and evaluate the matching performance with F1 score. The metrics are suitable for machining. We make clustering on each of the samples and for every sample in a cluster, we add the others to its matching groups.

In real circumstances of recommendations based on product matching, there should not be a large number of similar products to be displayed to users. In the model, we should control the numbers of similar samples. In the baseline model, we simply remove redundant samples to meet the number requirements instead of detailed comparison of product similarity.

For the DBSCAN model, we control the general sizes of clusters by setting the maximum distance between two samples for one to be considered as the neighborhood of the other. In our baseline model, we tune the parameter from 1 to 10 and compare the corresponding F1 scores and set the ultimate distance as 6.

Also, we combine the clustering result of the DBSCAN model with the result from grouping the perceptual hashing value of products to generate an extra result for the baseline model. The best F1 score of perceptual hashing, DBSCAN and their combination reaches 0.6518.
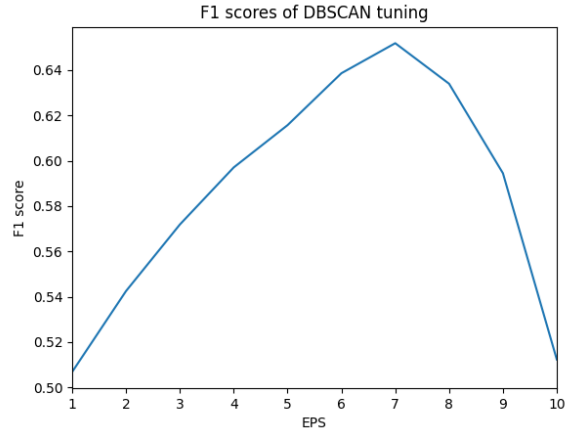


*Figure 11.* Tuning result of DBSCAN model

## 6.  2nd Deep Learning Baseline Solution

Given clustering methods show the low performance in price matching for Shopee, we transferred our focus to Computer Vision, NLP and phash matching phase. Therefore, we combined three methods including Perceptual hashing, TF IDF and ResNet18. For perceptual hash, we match them by map functions. For text analysis, we use TFIDF to transform text data into vectors and text embedding and then analyze text chunk by chunk. When it comes to image data, we load pretrain Pytorch models and set ResNet18 as our main framework. Our modeling performance is demonstrated in the following table, which is better than first machine learning models, namely clustering method.

| Type | F1 Score |
|---|---|
| Combined Score | 0.7342 |
| Image hash score | 0.5530 |
| TFIDF | 0.6153 |
| ResNet18 for images | 0.6527 |

*Table 3.* Deep Learning &NLP baseline

## 7.  Modeling Performance Enhancement

### 7.1  Data Argumentation

There exist loads of duplicated images and concentrated titles in our dataset. As shown in the figure, some images are overused while other images are only used once. This is similar in product titles, where the same tiles or expression words are utilized mostly.
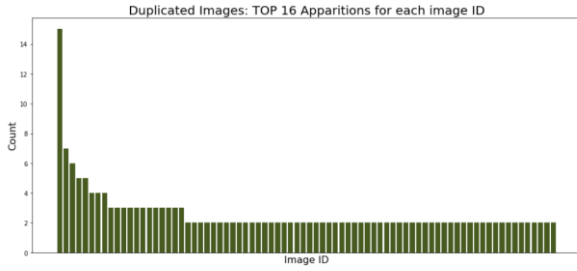
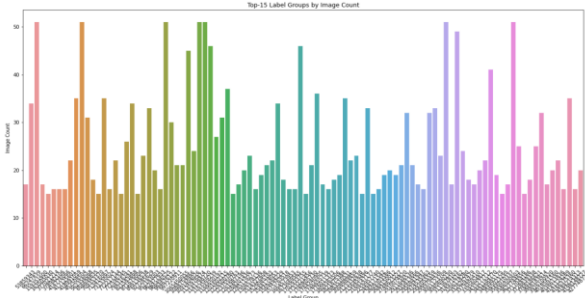*Figure 12.* Counts of Top 16 duplicated images



*Figure 13.* Counts of Top 13 overused titles

For image data: Under the circumstance, we need to perform data augmentation for existing imbalanced dataset. We use OpenCV to argument our image data including Horizontal Flip, Shift Scale Rotate, Optical Distortion and self-defined function to transform original images and create more data. The samples of original images and augmented images are shown below:
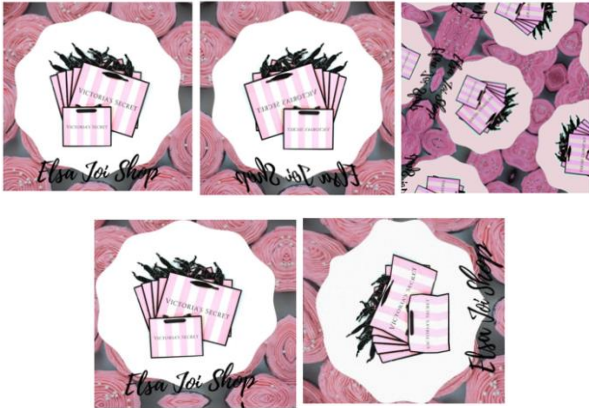


*Figure 14.* Original Image and Augmented Images

For text data: we use 6 methods to transform existing data, which are Swap character, Delete character, Word Augmenter, BERT Augmenter, TFIDF Augmenter and Sentence Augmenter respectively.

## 7.2 Model Adjustment

For image hash, each image owns their own hash value in our dataset. Therefore, how to better use hash to match

products will be significant for our modeling. We applied four hashing methods including average hashing, Perceptual hashing (pHash), difference hashing (dHash), wavelet hashing (wHash) as our main pipeline, and calculate corresponding F1 score. After consulting Kaggle Gold winners, we realize that the pHash will provide us with higher F1 score.

For Image color, we visualize the main color steam in our dataset in the following figure. As shown, most sample groups have similar color, demanding us to pay more attention to the differences of colors.
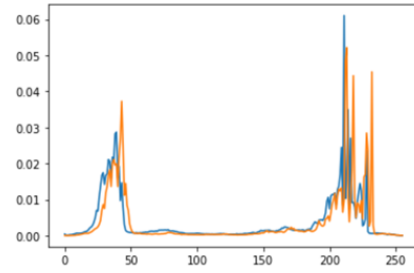


*Figure 15.* Color hist of our image dataset

For local features: Feature matching refers to the act of recognizing features of the same object across images with slightly different viewpoints. Some key points will be crucial in identifying specific images. There are a few different ways to find key points in an image. In this part, we use OpenCV KeyPoint to seize most local instances and features. Its affiliated working principle is shown below:
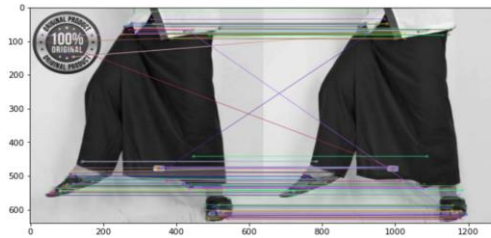


*Figure 16.* KeyPoints to identify products

## 8. 3rd Final ShopeeNet Solution

Based on previous models and data argumentation, we proposed our final model for final matching. We decide to further adopt NLP and deep learning methods, which are NFNet, Swin Transformer, ArcFace and TF-IDF. The aim of our target is to improve our modeling performance compared with our baseline models. The pipelines we worked are illustrated in the figures (1 represents preprocess data; 2 represents CNN construction; 3 represents CNN configuration; 4 represents model running and performance).
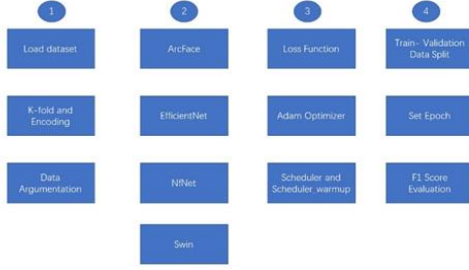
*Figure 17.* Image model workflow pipeline



*Figure 18.* Price matching workflow pipeline

## 8.1 For image matching

### 8.1.1 NFNet: Normalizer-Free Networks

In convention, ResNets without normalization are often unstable for large learning rates or strong data augmentation. NFNet proposes an adaptive gradient clipping technique which overcomes these instabilities for ResNets without normalization.
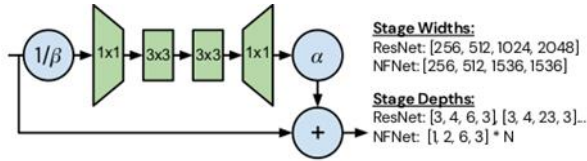


*Figure 19.* Working mechanism of NFNet

### 8.1.2 ArcFace-Additive Angular Margin Loss

ArcFace, or Additive Angular Margin Loss, is a loss function used in face recognition tasks. The softmax is traditionally used in these tasks. However, the softmax loss function does not explicitly optimize the feature embedding to enforce higher similarity for intraclass samples and diversity for inter-class samples, which results in a performance gap for deep face recognition under large intra-class appearance variations.

### 8.1.3 Swin Transformer

The Swin Transformer is a type of Vision Transformer. It builds hierarchical feature maps by merging image patches (shown in gray) in deeper layers and has linear computation complexity to input image size due to computation of self-attention only within each local window (shown in red). It can thus serve as a general-purpose backbone for both image classification and dense recognition tasks. In contrast, previous vision Transformers produce feature maps of a single low resolution and have quadratic computation complexity to input image size due to computation of self-attention globally.
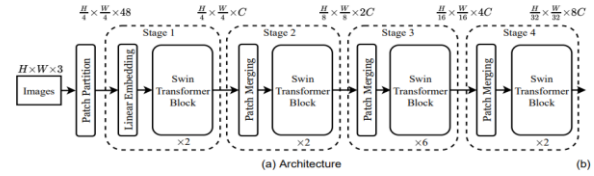


*Figure 20.* Working mechanism of NFNet

### 8.1.4 Image Embedding Models

Based on previous methods, we design two image matching models helping us extract useful features. Their working pipelines are as follows, which will supplement each other by extracting image features together. At the end of layers, we design a normalizer and cross entropy loss function to control modeling performance.



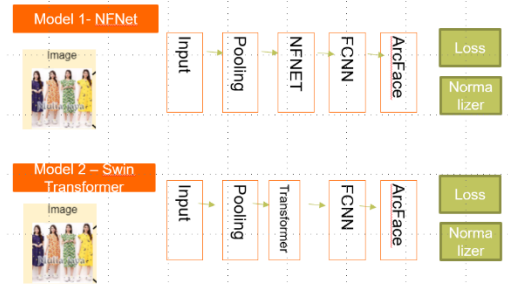*Figure 21.* Final image model working pipeline

## 8.2 For title matching – TF-IDF

TF-IDF (term frequency-inverse document frequency) is a statistical measure that evaluates how relevant a word is to a document in a collection of documents. This is done by multiplying two metrics: how many times a word appears in a document, and the inverse document frequency of the word across a set of documents.

## 8.3 For Phash matching

Phash (Perceptual hashes) must be robust enough to consider transformations or "attacks" on a given input and yet be flexible enough to distinguish between dissimilar files. Such attacks can include rotation, skew, contrast adjustment and different compression/formats. All these challenges make perceptual hashing an interesting field of study and at the forefront of computer science research.

## 8.4 Validation Data

We also split our data into training data, validation data and test data into modeling. A validation dataset is a sample of data held back from training your model that is used to give an estimate of model skill while tuning model's hyperparameters.

According to previous matching methods, we just combine their matching results using the intersection method, and then acquire the final matching method.

## 8.5 Training Performance Evaluation

In this part we use the same performance evaluation method F1 score as previous parts we did. We set epoch equals to 20 and compute the best threshold and corresponding scores. When the epoch is greater than 8, then the best score will decrease. The best threshold is 0.55 which appears from the $5^{th}$ epoch.

| Epoch | Threshold | Best F1 Score |
|-------|-----------|---------------|
| 0 | 0.25 | 0.7294 |
| 1 | 0.35 | 0.7883 |
| 2 | 0.45 | 0.8154 |
| 3 | 0.5 | 0.8247 |
| 4 | 0.5 | 0.8275 |
| 5 | 0.5 | 0.8319 |
| 6 | 0.55 | 0.8335 |
| 7 | 0.55 | 0.8336 |
| 8 | 0.55 | 0.8325 |

*Table 4.* Best score and threshold of each epoch in training data
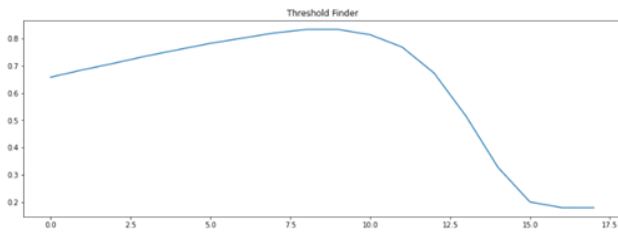


*Figure 22.* Accuracy and Threshold of Epoch 7 in training data

## 8.6 Test Performance Enhancement Process

In retrospect to our modeling performance enhancement process of this part, we primarily perform the following techniques in the table 5.

| Action | Final F1 Score |
|--------|----------------|
| Combine image, text, and hash | 0.721 |
| Add Normalization | 0.732 |
| With Best Threshold | 0.743 |
| With Validation | 0.757 |
| Final intersection | 0.776 |

*Table 5.* Our modeling performance enhancement process

## 9. Business Applications

Matching a seller listed item to an appropriate product has become a fundamental and one of the most significant steps for e-commerce platforms for product based experience. It has a huge impact for providing better user experiences by making search effective, supporting product reviews and product price estimation as well as many other advantages.

With the machine learning model we built, we can now easily identify and group the similar products from the large dataset of the e-commerce platform based on just the product image and product title description. Then the categorized product listing can be used by the platform to gather the product price insights, understand if the product is overpriced or underpriced among the peers and provide market analysis to the suppliers on the platform. This allows the platform to encourage the suppliers to stay active and competitive for attracting more customers.

The grouped product listing can also be used as the source information for the product search engine on the platform to improve the search results. From the traditional exact match of terms search, our processed dataset would be able to add the similar products into the search results. This not only allows the customers to get a more complete listing for their search, but also aids the customers to compare and find the best deals.

With the product groupings assigned, it allows the platform to rank the products within the group based on attributes like sales, reviews or even the promotion tactics. Then on the specific product pages, the higher ranked product from the same grouping can be selected to be displayed for customer's comparison and consideration, which enriches the customer's user experience, and in the event of promotional season, attracts the customer's attention to shop on the promoted items in the same product groupings.

Lastly, the product grouping can be further analyzed for understanding the relationships between groups. Then it could be used to construct recommender systems for recommending similar products to the customers based on user profiles and browsing histories, and recommending products in complementary groups based on purchases.

## 10. Limitations and Further Studies

Our product matching is only based on the product image and product title, which is the limited information we get from the dataset. For further studies, other attributes of the products can be considered, like the product brand, weight, material, etc.

For further studies, our solution can be extended to identify similarity among entities with richer information such as video or audio, and support different types of product embedding such as triplet model or graph-based models to fit more production use cases.

The inputs to our models are structured texts that it is not very easy to scale up to millions of classes. And computing pairwise similarities among billions of products are resource demanding which makes it challenging in large-scale data processing. To provide a flexible and consolidated solution at scale, distributed computing technology may be considered in the future studies or for actual business applications.

 Lastly, the definition of similarity varies across different applications in different eyes, e.g. substitute products could be viewed as similar, products with common design patterns could be viewed as similar and products with near colors could be viewed as similar. Hence, the grouping of the products could be further improved by incorporating feedbacks from the downstrain models or even directly from the users or domain experts.

## References

E. Shieh, S. Simhon, G. Aluri, G. Papachristoudis, D. Yakut and D. Raghu, "Attribute Similarity and Relevance-Based Product Schema Matching for Targeted Catalog Enrichment," 2021 IEEE International Conference on Big Knowledge (ICBK), 2021, pp. 261-270.

Gingold, D., (2021, August 11). Face Recognition and ArcFace: Additive Angular Margin Loss for Deep Face Recognition. https://medium.com/analytics-vidhya/face-recognition-and-arcface-additive-angular-margin-loss-for-deep-face-recognition-44abc56916c

Huang, B., Juaneda, C., Sénécal, S., & Léger, P. (2021). "Now you see me": The attention-grabbing effect of product similarity and proximity in online shopping.

Journal of Interactive Marketing, 54(1), 1-10. doi:10.1016/j.intmar.2020.08.004

Jbene, M., Tigani, S., Saadane, R., & Chehri, A. (2021). Deep neural network and boosting based hybrid quality ranking for e-commerce product search. Big Data and Cognitive Computing, 5(3), 35. doi:10.3390/bdcc5030035

Khan, S., (2021, March 12). NFNet-High-Performance Large-Scale Image Recognition Without Normalization (paper explained). https://samreenkhan498.medium.com/nfnet-high-performance-large-scale-image-recognition-without-normalization-paper-explained-53d43b5905c4

Payne, M., (2021, November 10). 5 Best Use Cases For Product Matching in Ecommerce & How You Can Implement Each One. https://www.width.ai/post/product-matching-in-ecommerce

Shieh, E., Simhon, S., Aluri, G., Papachristoudis, G., Yakut, D., & Raghu, D. (2021). Attribute similarity and relevance-based product schema matching for targeted catalog enrichment. Paper presented at the 261-270. doi:10.1109/ICKG52313.2021.00043

Shah, Kashif, Selcuk Kopru, and Jean David Ruvini. "Neural network based extreme classification and similarity models for product matching." Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers). 2018.

Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, B. Guo. (2021) Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. doi:10.48550/arXiv.2103.14030

Zuo, Zhen, et al. "A flexible large-scale similar product identification system in e-commerce." KDD Workshop on Industrial Recommendation Systems. 2020.