

# SkiLL: Skipping Color and Label Landscape: Self Supervised Design Representations for Products in E-commerce

Vinay K Verma, Dween Rabiuss Sanny, Shreyas Sunil Kulkarni, Prateek Sircar  
Abhishek Singh and Deepak Gupta  
International Machine Learning, Amazon

vinayugc@gmail.com, drsanny@amazon.com, sskshreyas@gmail.com, sircarp@amazon.com  
p15abhisheks@iima.ac.in and deepakgupta.cbs@gmail.com

## Abstract

*Understanding the design of a product without human supervision is a crucial task for e-commerce services. Such a capability can help in multiple downstream e-commerce tasks like product recommendations, design trend analysis, image-based search, and visual information retrieval, etc. For this task, getting fine-grain label data is costly and not scalable for the e-commerce product. In this paper, we leverage knowledge distillation based self-supervised learning (SSL) approach to learn design representations. These representations do not require human annotation for training and focus on only design related attributes of a product and ignore attributes like color, orientation, etc. We propose a global and task specific local augmentation space which captures the desired image information and provides robust visual embedding. We evaluated our model for the three highly diverse datasets, and also propose and measure a quantitative metric to evaluate the model’s color invariant feature learning ability. In all scenarios, our proposed approach outperforms the recent SSL model by upto 8.6% in terms of accuracy.*

## 1. Introduction

The automatic identification of visual properties, such as texture, pattern, style, and shape, is highly critical for e-commerce stores. This capability has many practical applications, including simulating an offline shopping experience for customers, enabling sellers and brands to understand and respond to current trends, and aiding in the recommendation of relevant products to customers. The visual feature-based identification should not be limited to clothing or home decor items, but should widely apply to a range of products, including fashion, electronics, furniture, beauty, and bags. One primary challenge is capturing the pattern, design, and shape while remaining agnostic to other properties, such as color and image orientation (e.g. image 2).

An usual e-commerce store may have millions of images



Figure 1. Same pattern label (striped) but visually different pattern



Figure 2. Similar design but different orientation and color



Figure 3. An example of 2 different products with similar pattern

for each product type, but labelling them is extremely challenging, costly and time consuming. Sometimes, sellers or brands provide the product attribute (e.g. ‘floral’ or ‘solid’ patterns) but these pattern labels are often too generic and coarse-grained, resulting in products with the same label (e.g., ‘striped’) that may not necessarily “look similar” in terms of pattern (as depicted in Figure 1). Self-supervised learning (SSL) has recently gained significant attention due to its ability to learn visual embeddings without labeled data [4, 5, 8, 13, 14, 21]. The SSL approaches, such as SimCLR [5], MoCo [17], and MoCoV2 [6], use image augmentation techniques (such as Gaussian Blur, color jittering, and random cropping) to generate positive pairs of the same image to learn visual similarity embeddings. However, these approaches also require negative samples for contrastive loss which is error prone. Recent works, BYOL [14], DINO [4], and SimSiam [8] have introduced self-knowledge distillation-based architectures that eliminate the use of negative pairs.

In the proposed model, we leverage the SSL approach to learn the desired feature. However, we observed that SSL has several shortcomings when used directly. Specifically, it fails to learn fine-grained design representations and is not agnostic to visual properties like color. To address these

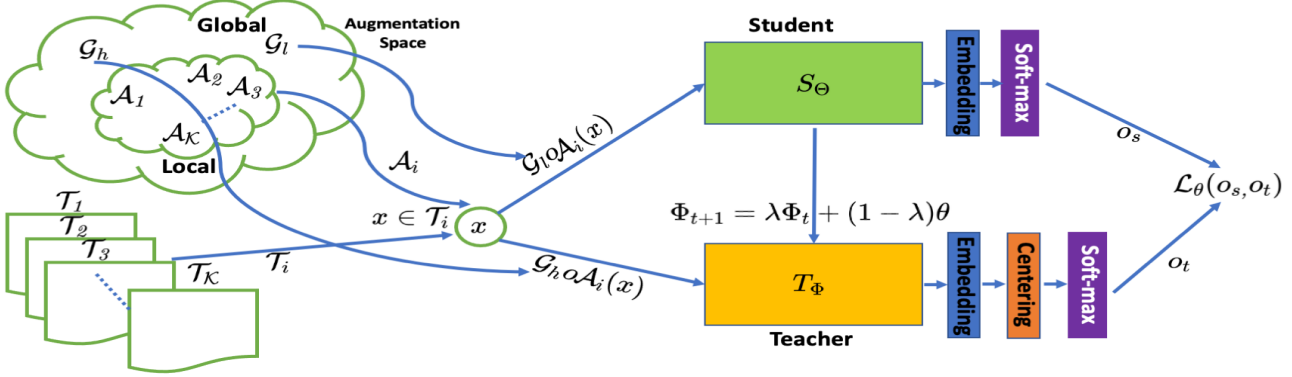


Figure 4. The block diagram of the proposed model, augmentation space contains the local and global augmentation. The local augmentation are selected based on the product type  $\mathcal{T}_i$  and  $\mathcal{G}_l$  and  $\mathcal{G}_h$  augmentation are applied to achieve two pair.

shortcomings, we used additional augmentations to learn the task specific and color invariant feature. The e-commerce products are diverse and same augmentation does not work well across multiple product therefore we propose task specific augmentation. The proposed model is robust and shows the promising result for a diverse set of task as compared to vanilla SSL model. Also we propose a new metric to measure the model’s color invariant learning ability. Our model has been extensively tested on diverse datasets, e.g. women’s dresses, furniture, and shirts, demonstrating its effectiveness. On tested classification tasks, our models show increase by upto 8.6% in accuracy.

## 2. Related Work

The goal of this research is to develop a model that learns an embedding, such that images with similar visual properties are closer in the embedding space, also the embedding should be agnostic to the color [1, 11]. To address this problem, several proposals [3, 15, 19, 20, 22] have been made. Recently, self-supervised approach allows learning without the need for supervisory signals, has gained significant attention due to its ability to eliminate the costly annotations. The initial SSL work RotNet [13], CFN [21] shows good performance, however the recent work in the SSL focus on the contrastive [16] or knowledge distillation based learning and shows the promising performance. The approach [5, 7] leverages the contrastive learning [10, 16] between the positive and negative pair. These pairs are made using the augmentation and random sampling respectively. Here sampling the negative data is error prone therefore model achieves degraded results.

Another set of recent research BOYL [14], SimSiam [8], DINO [4] discard the requirement of negative samples and trains the model only using the original sample and its augmentation as a positive samples. BOYL and SimSiam model use CNN architecture as backbone and tries to learn the similar embedding of the augmented version of the same samples. Other work DINO [4] incorporates the transformer architecture as backbone and minimize the cross entropy loss between the two embedding of the teacher and student model.

In these approach the model collapse is the key bottleneck, to overcome the above challenge knowledge distillation [2, 18] methods are used. During the distillation only one network [12] is trained and other learns using the exponential moving average.

## 3. Proposed Method

The e-commerce products encompass a wide range of diverse categories, including apparel (e.g. dresses, shirts, pants, jackets), curtains, furniture, toys, wall art, and shoes. Gathering labeled data for these categories can be a challenging and expensive process. The proposed approach is based on the self-supervised learning approaches [4, 5, 8, 14]. Similar to DINO [4] our approach leverages the vision transformer [9] as the base architecture for image encoding. The fixed augmentation technique does not work well for the diverse category since some key attribute in a category may be not required for others, also in the recommendation of design or pattern color are not important therefore model should learn the color invariant feature. To achieve this, we propose a data augmentation technique that captures the desired information and produces the robust output.

### 3.1. Notations

We consider a set of tasks  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_K$ , where each task corresponds to a specific product type. Let  $(S_\theta)$  be the student model and  $(T_\phi)$  be the teacher model, where  $\theta$  and  $\phi$  are the parameters of the student and teacher, respectively.

### 3.2. Augmentation Space

Data augmentation is crucial in SSL models, and can greatly impact the performance of the model. In our approach, we define two types of augmentations: global augmentations, which capture generic information that is common across tasks, and local augmentations, which capture task-specific information. Let  $\mathcal{G} = \mathcal{G}_h, \mathcal{G}_l$  represent global augmentations and let  $\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \dots, \mathcal{A}_K$  represent local augmentations. Here,  $\mathcal{G}_{\{h,l\}} = g_1, g_2, \dots, g_g$  and  $\mathcal{A}_i = a_1, a_2, \dots, a_{l_i}$ , where  $a_i \in \mathcal{A}$  and  $\mathcal{A}$  represents the local augmentation space. Each local  $\mathcal{A}_i$  contains a set of augmentations from the local augmentation space. During the training of the  $i^{th}$  task  $\mathcal{T}_i$ , we select the global augmentation

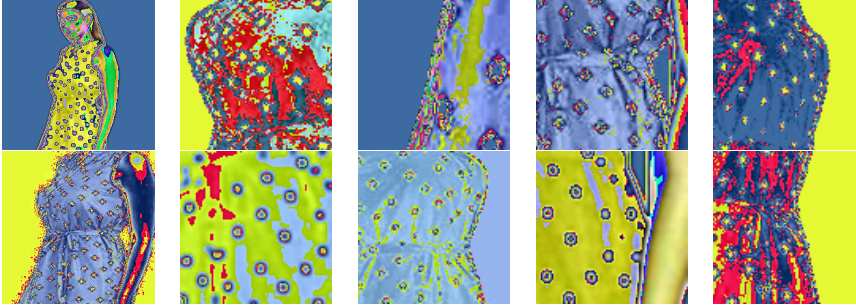


Figure 5. Global augmentation space,  $\mathcal{G}_h$  (first column),  $\mathcal{G}_l$  (remaining). Here  $\mathcal{G}_h$  are of high resolution ( $224 \times 224$ ) and  $\mathcal{G}_l$  are of low resolution ( $96 \times 96$ ).

$\mathcal{G}_h, \mathcal{G}_l$  and the corresponding local augmentation  $\mathcal{A}_i$ . More details around the augmentations is in the supplementary sheet. Also for a product, in e-commerce, there are typically multiple images (5-8) of a single product. Let  $\mathbf{x}'$  and  $\mathbf{x}''$  be two random images from the same product. The positive pair for the  $\mathbf{x}'\mathbf{x}'' \in \mathcal{T}_i$  is defined as follows:

$$\mathbf{x}_1, \mathbf{x}_2 = \mathcal{G}_h \circ \mathcal{A}_i(\mathbf{x}'), \mathcal{G}_h \circ \mathcal{A}_i(\mathbf{x}'') \quad (1)$$

$$\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_c = \mathcal{G}_l \circ \mathcal{A}_i(\mathbf{x}') \quad (2)$$

In our experiments,  $\mathcal{G}_h$  and  $\mathcal{G}_l$  are high and low-resolution images.  $\mathcal{G}_h$  generates two positive pairs and  $\mathcal{G}_l$  generates  $c$  positive pairs.  $\circ$  is a composite operation.  $\mathcal{G}_h \circ \mathcal{A}_i(\mathbf{x})$  and  $\mathcal{G}_l \circ \mathcal{A}_i(\mathbf{x})$  give the output after the composite operation. From the set of local augmentations, we can choose any task-specific augmentation and find images with positive pairs  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . See Figure 4, 5, 6 for a visual illustration.

### 3.3. Teacher-Student framework for SSL

We have a set of tasks  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_K$  where each task is a different product type. During the training of the  $i^{th}$  task  $\mathcal{T}_i$ , we select  $\mathcal{G}_h$  and  $\mathcal{G}_l$  as the global and  $\mathcal{A}_i$  as the local augmentation. Say, after applying the augmentation ( $\mathcal{G}_h \circ \mathcal{A}_i(\mathbf{x})$ ) we get  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , which are the high resolution augmented positive pair and ( $\mathcal{G}_l \circ \mathcal{A}_i(\mathbf{x})$ ) gives  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_c$  which are low resolution positive pair. The augmentation obtained by  $\mathcal{G}_h$  and  $\mathcal{G}_l$  passed to the network  $T_\phi$  and  $S_\theta$  respectively, which is given as:

$$\mathbf{o}_{s_i} = SM(S_\theta(\mathbf{x}_i)); \quad \& \quad \mathbf{o}_{t_i} = SM(C_e(T_\phi(\mathbf{y}_i))) \quad (3)$$

Where  $C_e$  is the centering function and  $SM$  is the Softmax operation. The Softmax operation for the input  $S_\theta(\mathbf{x}_i)$  is:

$$\mathbf{o}_{s_i} = \frac{\exp(S_\theta(\mathbf{x}_i)/\tau)}{\sum_{i=1}^N \exp(S_\theta(\mathbf{x}_i)/\tau)} \quad (4)$$

Here  $\tau$  is the temperature parameter which is used to flatten the output space. The loss between the output  $\mathbf{o}_{s_i}$  and  $\mathbf{o}_{t_i}$  is given as:  $\mathcal{L}(\mathbf{o}_{s_i}, \mathbf{o}_{t_i}) = -\mathbf{o}_{t_i} \log \mathbf{o}_{s_i}$  We optimize jointly the high and low resolution augmentation as follows:

$$\mathcal{L}(\theta, \phi) = \sum_{\{\mathbf{x}_1, \mathbf{x}_2\}} \sum_{\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_c\}} \mathcal{L}(\mathbf{o}_{s_i}, \mathbf{o}_{t_i}) \quad (5)$$



Figure 6. Local augmentation space,  $\mathcal{A}_i$  for the women-dress.

We obtained  $\mathbf{o}_{s_i}$  and  $\mathbf{o}_{t_i}$  from the Eq-3. The Eq-5 is optimized w.r.t. the student parameter  $\theta$  i.e.:  $\min_\theta \mathcal{L}(\theta, \phi)$  Note that the teacher model is not learned, and it uses exponential moving averaging to update its parameter and it is given as:  $\phi_t \leftarrow \lambda \theta_t + (1 - \lambda) \phi_t$ ; where  $\lambda$  is hyperparameter. The cross-entropy loss minimize the oscillate [4], and the model converges faster.

## 4. Data and Statistics

Product Category	Product Count	Training Images	Test Images
Women's Dress	94,393	534,057	176,993
Shirt	84,775	520,474	157,829
Furniture	56,161	334,589	80,024

Table 1. Datasets and their statistics for the three diverse datasets *Women-dress*, *Shirt* and *Furniture*.

To evaluate the proposed framework at scale, we considered three product categories - shirts, women's dress and furniture of a popular e-commerce service. The statistics are provided in the Table 1, For more details, refer to supplementary sheet.

## 5. Experiment and Results

In order to evaluate the performance of our model, we conducted rigorous experiments on a highly diverse dataset. The details of implementation, baselines and Evaluation metric are provided in the supplementary.

### 5.1. Evaluation Metric

To evaluate the models, we selected the attributes with good fill rate and collected the data for these attributes for each product type (task), to train a classification model. We measure the accuracy, Recall@90 (R@90) and Recall@95 (R@95) of the product over the selected attribute individually and report the result. Similarly, we define the evaluation metric to measure the color invariance. To our best knowledge, this is the first metric to measure the color invariant property for the model. We provide the details in the supplementary.

### 5.2. Results

We have evaluated the learned self-supervised embedding for the two scenario, that are given as follows:

		Pattern						Sleeve					
		Linear			Non-linear			Linear			Non-linear		
		Acc	R@90	R@95	Acc	R@90	R@95	Acc	R@90	R@95	Acc	R@90	R@95
Shirt	BYOL	56.31	04.10	02.24	59.91	6.494	01.70	81.30	71.14	60.15	83.37	73.51	65.46
	SimCLR	78.02	54.88	10.20	78.51	60.15	18.40	81.00	66.47	40.44	83.29	71.89	59.06
	DINO	80.29	<b>61.21</b>	25.87	83.19	71.62	50.71	91.39	78.81	77.81	93.93	88.59	81.32
	<b>SkiLL</b>	<b>81.01</b>	<b>59.90</b>	<b>31.07</b>	<b>85.66</b>	<b>76.76</b>	<b>63.52</b>	<b>92.81</b>	<b>79.36</b>	<b>78.75</b>	<b>95.14</b>	<b>91.99</b>	<b>86.93</b>
Women Dress	BYOL	66.77	14.96	01.24	66.95	15.55	02.43	34.46	00.03	00.03	44.99	00.87	00.34
	SimCLR	81.29	64.08	44.01	82.47	67.10	50.93	50.23	01.71	00.16	53.92	08.20	02.65
	DINO	83.11	67.55	48.28	90.87	87.92	79.96	62.26	17.65	04.08	78.19	58.34	37.09
	<b>SkiLL</b>	<b>83.99</b>	<b>68.04</b>	<b>52.30</b>	<b>92.32</b>	<b>91.34</b>	<b>84.46</b>	<b>70.97</b>	<b>38.32</b>	<b>18.67</b>	<b>85.57</b>	<b>75.71</b>	<b>61.21</b>

Table 2. Results for *Shirt* and *Women-dress* product type, evaluated for 2 attributes: Pattern & Sleeve which has 20 and 6 classes respectively.

	Other Hard-Line (OHL) Subcategories					
	Linear			Non-linear		
	Acc	R@90	R@95	Acc	R@90	R@95
BYOL	56.45	27.46	20.39	62.14	33.44	27.37
SimCLR	72.30	49.88	38.53	75.19	54.61	43.08
DINO	83.10	70.52	58.25	85.04	75.34	64.23
<b>SkiLL</b>	<b>83.77</b>	<b>71.79</b>	<b>61.35</b>	<b>86.02</b>	<b>77.45</b>	<b>65.32</b>

Table 3. Hard-line Subcategories classification result, there are 19 subcategories on the hard-line dataset.

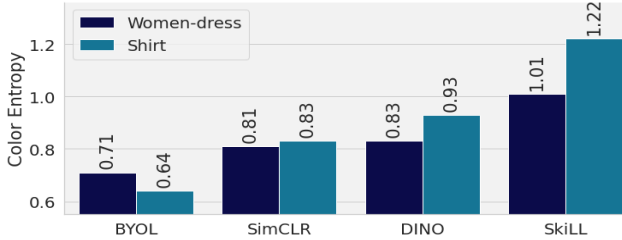


Figure 7. The result for the color invariant measurement, similar in style but agnostic to color shows high color entropy.

### 5.2.1 Class Discrimination

The results for class discrimination on the Shirt, Women-Dress, and Hard-Line datasets were evaluated using the metric discussed in Section 12.2. We have shown the results for these categories in Table 2 and 5. The CNN network-based models BYOL and SimCLR showed reasonable results for the easy categories of hard-line and shirt sleeves. However, for the pattern and women-dress sleeve categories, the models performance degraded significantly. These models struggled to capture the fine-grain details and discriminate between small variations in the complex categories. However, the proposed model easily capture these information and outperforms the CNN baselines.

### 5.2.2 Color Invariant

We evaluated the effectiveness of our learned embeddings in terms of color invariance by using the metrics outlined in Section-12.2. The results depicting the average color entropy across all clusters appear in Figure-12. We can see that SkiLL, demonstrates the highest level of color agnosticism compared to the other methods.

## 6. Ablations

**Local Augmentation and Preprocessing:** The local augmentation plays a key role in the specific product type. We

have conducted the ablation over the proposed task specific augmentation for the challenging dataset women-dress and pattern attribute type. We observe that if from *SkiLL* we remove the local augmentation ( $SkiLL \setminus \mathcal{L}$ ) the model performance degraded. The *SkiLL* without pre-processing ( $SkiLL \setminus \mathcal{H}$ ) significantly reduce the model performance. The results for the women dress are shown in the Table-6 for the linear and non-linear evaluation.

**Color Agnostic:** To learn the color agnostic embedding, we incorporate the *RandomGray*, *ColorJitter* and *GaussianBlur* as a global transformation with the probability  $p_1 = 0.5$ ,  $p_2 = 0.8$  and  $p_3 = 0.5$  respectively. We observe that for small value of  $p_1 = p_2 = p_3 = 0.1$  the color agnostic property of the model significantly degraded and color entropy reduce from 1.01 to 0.90. Therefore, these global augmentations are necessary to achieve the color agnostic embedding.

	Linear			Non-linear		
	Acc	R@90	R@95	Acc	R@90	R@95
$SkiLL \setminus \mathcal{H}$	83.67	<b>68.73</b>	50.85	90.92	87.81	78.59
$SkiLL \setminus \mathcal{L}$	83.74	68.15	52.27	91.91	89.25	79.02
<b>SkiLL</b>	<b>83.99</b>	68.04	<b>52.30</b>	<b>92.32</b>	<b>91.34</b>	<b>84.46</b>

Table 4. Ablation on the women-dress for the pattern classes, without local augmentation ( $SkiLL \setminus \mathcal{L}$ ) model performance degraded.

## 7. Conclusions

In this paper, we addressed the challenge of learning visual similarity in terms of style and pattern-type, agnostic to other properties, such as color, without the use of labeled data. We proposed a SSL approach using knowledge distillation with the transformer architecture as the backbone for the teacher and student networks. To solve the desired problem we carefully designed an augmentation space containing both local and global augmentations. The local augmentations were task-specific, allowing us to select the appropriate ones for a particular task in combination with the global augmentations. We incorporated the *RandomGray*, *ColorJitter*, and *GaussianBlur* global augmentations to achieve a color agnostic feature space. Our proposed augmentation outperformed the recent baseline by a significant margin. We proposed a novel metric to measure the model’s color agnostic property. The ablation shows the importance of the proposed component.

## References

- [1] Paul Baltescu, Haoyu Chen, Nikil Pancha, Andrew Zhai, Jure Leskovec, and Charles Rosenberg. Itemsage: Learning product embeddings for shopping recommendations at pinterest. *arXiv preprint arXiv:2205.11728*, 2022. 2
- [2] Cristian Bucilua, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541, 2006. 2
- [3] Xuefei Cao, Bor-Chun Chen, and Ser-Nam Lim. Unsupervised deep metric learning via auxiliary rotation loss. 2019. 2
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 3, 8
- [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning, ICML'20*. JMLR.org, 2020. 1, 2, 8
- [6] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 1
- [7] Xinlei Chen, Haoqi Fan, Ross B. Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *CoRR*, abs/2003.04297, 2020. 2
- [8] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. *arXiv preprint arXiv:2011.10566*, 2020. 1, 2
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. 2, 8
- [10] Alexey Dosovitskiy, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with convolutional neural networks. *Advances in neural information processing systems*, 27, 2014. 2
- [11] Ujjal Kr Dutta, Sandeep Repakula, Maulik Parmar, and Abhinav Ravi. A tale of color variants: Representation and self-supervised learning in fashion e-commerce. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 12482–12488, 2022. 2
- [12] Zhiyuan Fang, Jianfeng Wang, Lijuan Wang, Lei Zhang, Yezhou Yang, and Zicheng Liu. Seed: Self-supervised distillation for visual representation. *arXiv preprint arXiv:2101.04731*, 2021. 2
- [13] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations*, 2018. 1, 2
- [14] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent a new approach to self-supervised learning. In *Proceedings of the 34th NeurIPS*, 2020. 1, 2, 8
- [15] Geonmo Gu and Byungsoo Ko. Symmetrical synthesis for deep metric learning. 01 2020. 2
- [16] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006. 2
- [17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 1
- [18] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7), 2015. 2
- [19] Elad Levi, Tete Xiao, Xiaolong Wang, and Trevor Darrell. Reducing class collapse in metric learning with easy positive sampling. 06 2020. 2
- [20] Yang Li, Shichao Kan, and Zhihai He. Unsupervised deep metric learning with transformed attention consistency and contrastive clustering loss. 08 2020. 2
- [21] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European conference on computer vision*, pages 69–84. Springer, 2016. 1, 2
- [22] Yifan Sun, Changmao Cheng, Yuhan Zhang, Chi Zhang, Liang Zheng, Zhongdao Wang, and Yichen Wei. Circle loss: A unified perspective of pair similarity optimization. pages 6397–6406, 06 2020. 2

## Supplementary

### 8. Augmentation Space

The use of data augmentation is crucial in self-supervised learning models, as it can greatly impact the performance of the model. In e-commerce, where there is a wide variety of product categories, it is not feasible to apply a single set of augmentations to all tasks. For example, augmentations designed for learning design representations in the fashion category may not be effective when applied to the furniture category. Therefore, it is necessary to use different augmentations to capture the relevant information corresponding to each task. The model also needs to be agnostic to certain visual attributes, such as being color invariant. This means that if two items, such as clothes or furniture, are significantly different in color but share similar style or pattern, their embedding distance must be low.

#### 8.0.1 Intuition of Augmentation Space

To understand the intuition behind our proposed augmentation space, consider the example of a category of women’s dresses. Our goal is to learn features that are invariant to color, pose, etc. and focus on style, pattern, and texture. To achieve this, we divide the local and global features and define the augmentation accordingly. To make the model agnostic to color, we apply global augmentations such as *RandomGray*, *ColorJitter*, and *GaussianBlur*. These augmentations also include random size crop and random scale. When we augment an image using *RandomGray*, the teacher model receives a color image while the student network receives the grayscale equivalent. In order to minimize the embedding distance, the model must learn features that are invariant to color. Similarly, *ColorJitter* and *GaussianBlur* help the model learn color-invariant embeddings. In addition to these global augmentations, we also define lower-resolution augmentations to capture pattern. A close embedding distance is only possible if the model focuses on patterns.

One of the properties of e-commerce that we can exploit is the availability of multiple images for a single product, which we can use as positive pairs. These images typically show the same product from different angles or perspectives. By making them or their augmentations positive pairs, the model learns to ignore some noisy variations of the same image (e.g. different angles at which the image is shot). For example, if there are 5 images of a person wearing a dress in different poses, making these images positive pairs allows the model to focus on the dress and makes it agnostic to the pose of the person or the angle at which the image is shot. Besides these global and lower-resolution augmentations, we also devise local augmentations that capture task-specific properties such as neck style, sleeve length, and knee-length for women’s dresses. To determine image crops that contain

information about these attributes, we use head detection to identify the head of the person wearing the dress<sup>1</sup>. We then use the image crops below the head to capture information about sleeve length and knee-length, and we crop a rectangle below the detected head to capture neck style. By using these crops as positive pairs among local and with the global augmentations, we can train a model that focuses on the required attributes. If we move to another product category, such as shirts, we can discard any augmentations that do not exist, like knee-length. Similarly, if our use case is limited to detecting designs and not capturing any other attributes, we can drop the corresponding augmentations. In image 8 we show examples where product pairs were identified as “similar” by our model and “dissimilar” by the closest baseline.

### 9. Datasets and Evaluation Details

The collection of the dataset are discussed in the main paper. Here we are discussing the details of the evaluation metric and classes used for the evaluation.

#### 9.1. Women-dress

The women-dress are evaluated over the two attribute *Sleeve* and *Pattern* with 9 and 18 classes respectively. The class details are giving below:

##### 9.1.1 Sleeve

['3/4 sleeve', 'short sleeve', 'sleeveless', 'long sleeve', 'half sleeve', 'cap sleeve', 'puff sleeve', 'bell sleeve', 'flare sleeve']

##### 9.1.2 Pattern

['solid', 'floral', 'striped', 'polka', 'checkered', 'geometric', 'graphic', 'animal', 'tribal', 'paisley', 'printed', 'animal', 'plain', 'chevron', 'screen print', 'embroidered', 'quilted', 'print']

#### 9.2. Shirt

Shirt category also evaluated for the two attribute *Sleeve* and *Pattern* with 6 and 17 classes respectively. The class details are giving below:

##### 9.2.1 Sleeve

['long sleeve', 'half sleeve', 'short sleeve', 'cap sleeve', '3/4 sleeve', 'sleeveless']

<sup>1</sup>We use <https://github.com/deepinsight/insightface>



Figure 8. Examples of samples identified as “similar” (embedding distance  $\leq 0.15$ ) by our model. Nearest baseline (DINO), trained on the same data, identified them as dissimilar (distance  $\geq 0.8$ ). Women-dress (left), Shirt (middle) and Furniture (right).

### 9.2.2 Pattern

[‘solid’, ‘striped’, ‘checkered’, ‘printed’, ‘geometric’, ‘floral’, ‘tribal’, ‘quilted’, ‘animal’, ‘graphic’, ‘paisley’, ‘plain’, ‘polka’, ‘tie and die’, ‘print’, ‘chevron’, ‘embroidered’]

### 9.3. Furniture

The furniture class are evaluated for the sub-category. We have 19 subcategory that are discussed below.

#### 9.3.1 Sub-category

[‘Sofas- Large’, ‘Chairs- Large’, ‘Tables- Large’, ‘Chairs- Large’, ‘Sofa Sets- Large’, ‘Desks- Large’, ‘Dining Sets- Large’, ‘Patio Furniture Sets- Large’, ‘Arm Chairs- Large’, ‘Bedside Cabinets - Large’, ‘Bean Bags- Large’, ‘Dining Tables- Large’, ‘Ottomans & Footstools- Large’, ‘Dining Chairs- Large’, ‘Closet Organization’, ‘Folding Chairs- Large’, ‘Beds- Large’, ‘Desks & Tables- Large’, ‘Recliners & Loungers- Large’]

## 10. Evaluation Matrices

We have evaluated our model for the two scenario, the first focus over the discriminative feature learned by the model, however second evaluate the model’s ability to learn the color agnostic feature. The discriminative power are evaluated using the standard matrices like, accuracy, Recall@P. However to evaluate the color invariant feature we don’t have any standard metric. Therefore in this work we are proposing the evaluation metric for the same. The details of the color agnostic evaluation metric are discussed below.

### 10.1. Color Agnostic Metric

#### 10.1.1 Color ID

For learning color representations, we use an open-source detectron model to detect “person” in the image, then we crop the image to maximize the presence of the “person” in the image. After that, we map the color from RGB to LUV space. Finally, we get frequency distribution of color in

each image. Reason for moving from RGB to LUV space is that for 2 colors, correlation of Euclidean distance between the corresponding coordinates in LUV space and perceptual similarity that humans feel between them is higher than RGB space. We compare the histogram of the color coordinates with each color from a pre-defined list of colors. We assign the hex-code of the color in the list with which the apparel is closest. The hex-code is the color representation of the apparel.

#### 10.1.2 Measuring the color invariant

we collected the embeddings of the test set from our proposed method and cluster them using agglomerative clustering. After we get the clusters, we calculate the histogram of colours based on their frequency. Then we calculate the entropy of the clusters using equation 10. We then take the average of the entropy across all clusters.

### 10.2. Recall@P

We calculate the micro average precision and recall. We are reporting the recall at a predefined value for precision. We calculate the class probabilities and using precision-recall curve we calculate the the threshold for predefined precision value. After we get the threshold we calculate the true-positives and false-positives and false-negatives. We repeat this for every class and aggregate the values. At the end, we calculate the recall by:

$$Recall = \frac{True-positives}{(True-positives + False-negatives)} \quad (6)$$

## 11. Experiment and Results

In order to evaluate the performance of our model, we conducted rigorous experiments on a highly diverse dataset. In this section, we will provide information on the implementation details, the baseline model, the evaluation metric, and the results of the experiments.

	Other Hard-Line (OHL) Subcategories					
	Linear			Non-linear		
	Acc	R@90	R@95	Acc	R@90	R@95
BYOL	56.45	27.46	20.39	62.14	33.44	27.37
SimCLR	72.30	49.88	38.53	75.19	54.61	43.08
DINO	83.10	70.52	58.25	85.04	75.34	64.23
<b>SkiLL</b>	<b>83.77</b>	<b>71.79</b>	<b>61.35</b>	<b>86.02</b>	<b>77.45</b>	<b>65.32</b>

Table 5. Hard-line Subcategories classification result, there are 19 subcategories on the hard-line dataset.

### 11.1. Implementation Details

## 12. Implementation Details

We used patch sizes of eight in our models, which are more costly to use compared to patch sizes of 16, but provided better results. To aggregate information for the entire sequence, we added a learnable token to our model [9] without using label information. Our transformer module consisted of self-attention, residual connections, and fully connected layers. We trained our model with a batch size of 32 distributed over four GPUs. We employed a similar learning rate schedule as in [4], which utilized linear scaling for the first 10 epochs and then used cosine scheduling for the decay of the learning rate. The temperature parameter also played a crucial role in our model, and we set it to a lower value in the range of  $\tau \in [0.04, 0.07]$ . As the epochs progressed, we increased the  $l_2$  penalty from an initial value of 0.04 to  $10\times$  its original value. We used Pytorch libraries such as torchvision and opencv to train our model, which was done on an Nvidia V100 GPU.

### 12.1. Baseline

In our experiments, we compared our model (SkiLL) to several strong baselines in the self-supervised learning field. One of these baselines, the Transformer-based DINO model [4], is similar to our approach in that it uses a teacher-student framework with a transformer architecture. However, our approach differs in the use of local and global task-specific augmentation. Another baseline, BYOL [14], also uses a teacher-student framework with a convolutional network, but it measures the cosine similarity between positive pairs. Finally, SimCLR [5] employs a contrastive loss and negative samples for self-supervision. We evaluated all of these approaches, against SkiLL, using linear and non-linear classifier for the learned features and the proposed color invariant metric.

### 12.2. Evaluation Metric

We collected the data from a popular e-commerce service where corresponding to each product we have a set of attribute value for e.g., the women-dress has attribute sleeve length, neck style, pattern, etc. We collected the data from a

popular e-commerce service where corresponding to each product we have a set of attribute value for e.g., the women-dress has attribute sleeve length, neck style, pattern, etc. To evaluate the models, we decided on available attributes with good fill rate collected the data for these attributes for each product type (task), to train a classification model. For example, in the women-dress the *pattern* attribute has 18 labels [*solid*, *floral*, *striped*, *polka*,...]. The self-supervised learned embeddings are classified into one of the 18 classes. We measure the accuracy, Recall@90 (R@90) and Recall@95 (R@95) of the product over the selected attribute individually and report the result. The accuracy for the print-type attribute over the product  $\mathcal{T}_i$  is given as:

$$A_{print-type}(\mathcal{T}_i) = \frac{1}{N_i} \sum_{j=1}^{N_i} Acc(\hat{\mathbf{y}}, \mathbf{y}) \quad (7)$$

Here  $\hat{\mathbf{y}}$  and  $\mathbf{y}$  are the predicted and ground truth value  $Acc(\hat{\mathbf{y}}, \mathbf{y}) = 1$  if both belong to same class else zero.

Similarly, we define the evaluation metric to measure the color invariance, to our best knowledge, this is the first metric to measure the color invariant property for the model. We first use the embeddings to perform an agglomerative clustering, such that “visually similar” product as per the model get grouped into one cluster. Then we convert the image into LUV space, in this space we define each image by a unique color ID. See the supplement for the process to assign color ID to each product. Once we have color ID, for each cluster we can find the histogram of based on the frequency of the color id in a particular cluster. Let  $h_i$  is the histogram of a particular cluster  $j$ , then its entropy of the task is defined as:

$$C_i = \frac{1}{L} \sum_{j=1}^L -h_j \log h_j \quad (8)$$

More the color diversity in a cluster shows the higher average entropy. In this way, we can compare the two model’s color invariant learning ability, however having the same number of cluster  $L$ . A model which is not agnostic to color, would tend to have a peaky histogram and hence lower entropy.

## 12.3. Results

We have evaluated the learned self-supervised embedding for the two scenario, that are given as follows:

### 12.3.1 Class Discrimination

The results for class discrimination on the Shirt, Women-Dress, and Hard-Line datasets were evaluated using the metric discussed in Section 12.2. The pattern and sleeve attributes of the Shirt and Women-Dress datasets were evaluated. The Women-Dress dataset had 18 classes for the pattern attribute and 9 classes for the sleeve attribute, while the Shirt dataset had 17 and 6 classes, respectively, for the



same attributes. The Hard-Line dataset, which included products such as sofas and furniture, had 19 subcategories that were used as targets for classification. More details on the dataset are provided in the supplementary sheet. The results for these categories are shown in Table 3 and 2. The convolutional network-based models BYOL and SimCLR showed reasonable results for the easy categories of hard-line and shirt sleeves. However, for the pattern and women-dress sleeve categories, the models showed significantly degraded performance. These models struggled to capture the fine-grain details and discriminate between small variations in the complex categories. Table 2 demonstrates that the sleeve category for the shirt showed better results because it is simpler, while the women-dress sleeve category with more complex classes resulted in significantly degraded model performance. The DINO model shows promising results for most product types and outperformed the BYOL and SimCLR models by a significant margin. However, the proposed model, SkiLL, consistently outperforms the state-of-the-art DINO model and showed promising results even when all other models failed, such as the women-dress sleeve category. Figures 9, 10 and 11 are some examples from women’s dress sleeves classification task, where other baseline fails to properly predict the sleeve type. Some more examples can be seen in the supplementary sheet. While other models showed



Figure 9. DINO: long sleeve, SkiLL: puff sleeve  
 Figure 10. DINO: long sleeve, SkiLL: bell sleeve  
 Figure 11. DINO: three quarter sleeve, SkiLL: short sleeve

good accuracy, they struggled with high recall scores, particularly at the Recall@90 (R@90) and Recall@95 (R@95) levels. The proposed SkiLL model, however, shows high recall scores. Finally, the model was evaluated using both linear and non-linear classifiers, and the linear classifier also showed promising results, indicating that the learned embedding can be used for downstream tasks, such as clustering.

### 12.3.2 Color Invariant

We evaluated the effectiveness of our learned embeddings in terms of color invariance by using the metrics outlined in Section-12.2. The results, which depict the average color entropy across all clusters, are presented in Figure-12. It

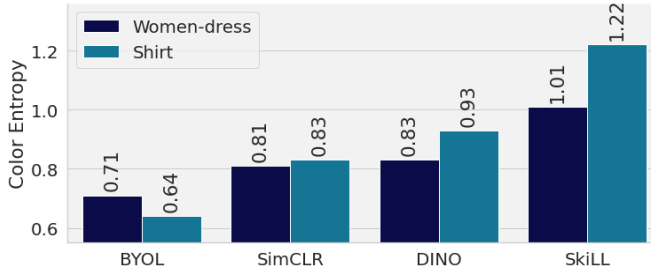


Figure 12. The result for the color invariant measurement, similar in style but agnostic to color shows high color entropy.

can be seen that our proposed model, SkiLL, demonstrates the highest level of color agnosticism compared to the other methods. Further, the embedding for shirts exhibits greater color agnosticism than for women’s dresses. This may be attributed to the higher level of color variation present in the training data for shirts, leading to a more robust model that can effectively ignore color information.

## 13. Ablations

**Local Augmentation and Preprocessing:** The local augmentation plays a key role in the specific product type. We have conducted the ablation over the proposed task specific augmentation for the challenging dataset women-dress and pattern attribute type. We observe that if from *SkiLL* we remove the local augmentation ( $SkiLL \setminus \mathcal{L}$ ) the model performance degraded. Also, the image pre-processing like object detection head cropping overcome the model bias towards the background or undesired feature. The *SkiLL* without pre-processing ( $SkiLL \setminus \mathcal{H}$ ) significantly reduce the model performance. The results for the women dress are shown in the Table-6 for the linear and non-linear evaluation.

**Color Agnostic:** As we discussed earlier to learn the color agnostic embedding, we incorporate the *RandomGray*, *ColorJitter* and *GaussianBlur* as a global transformation with the probability  $p_1 = 0.5$ ,  $p_2 = 0.8$  and  $p_3 = 0.5$  respectively. We observe that for small value of  $p_1 = p_2 = p_3 = 0.1$  the color agnostic property of the model significantly degraded and color entropy reduce from 1.01 to 0.90. Therefore, these global augmentations are necessary to achieve the color agnostic embedding.

	Linear			Non-linear		
	Acc	R@90	R@95	Acc	R@90	R@95
<i>SkiLL</i> \ $\mathcal{H}$	83.67	<b>68.73</b>	50.85	90.92	87.81	78.59
<i>SkiLL</i> \ $\mathcal{L}$	83.74	68.15	52.27	91.91	89.25	79.02
<b>SkiLL</b>	<b>83.99</b>	68.04	<b>52.30</b>	<b>92.32</b>	<b>91.34</b>	<b>84.46</b>

Table 6. Ablation on the women-dress for the pattern classes, without local augmentation ( $SkiLL \setminus \mathcal{L}$ ) model performance degraded.

#### 14. Examples where SkiLL classification is better than Baseline



Figure 13. DINO: three-fourth, SkiLL: bell sleeve  
 Figure 14. DINO: sleeveless, SkiLL: long sleeve  
 Figure 15. DINO: short sleeve, SkiLL: sleeveless

formed via agglomerative clustering of the SkiLL embeddings of the images. Note that the embeddings focus primarily on style of the apparels and is agnostic to the color of the product or pose of the human. These us analyse design trends and recommend a different color for the same design to customers.

#### 15. Qualitative Clusters with SkiLL embeddings



Figure 16. Cluster example 1



Figure 17. Cluster example 2

In figures 16 and 17 we show examples of clusters