

Fine-Grained Fashion Similarity Learning by Attribute-Specific Embedding Network

Zhe Ma^{1*}, Jianfeng Dong^{2,3*†}, Yao Zhang^{1*},
Zhongzi Long¹, Yuan He⁴, Hui Xue⁴, Shouling Ji^{1,3†}

¹Zhejiang University, ²Zhejiang Gongshang University,

³Alibaba-Zhejiang University Joint Institute of Frontier Technologies, ⁴Alibaba Group

{maryeon, y.zhang, akasha, sji}@zju.edu.cn, dongjf24@gmail.com, {heyuan.hy, hui.xueh}@alibaba-inc.com

Abstract

This paper strives to learn fine-grained fashion similarity. In this similarity paradigm, one should pay more attention to the similarity in terms of a specific design/attribute among fashion items, which has potential values in many fashion related applications such as fashion copyright protection. To this end, we propose an *Attribute-Specific Embedding Network (ASEN)* to jointly learn multiple attribute-specific embeddings in an end-to-end manner, thus measure the fine-grained similarity in the corresponding space. With two attention modules, *i.e.*, *Attribute-aware Spatial Attention* and *Attribute-aware Channel Attention*, ASEN is able to locate the related regions and capture the essential patterns under the guidance of the specified attribute, thus make the learned attribute-specific embeddings better reflect the fine-grained similarity. Extensive experiments on four fashion-related datasets show the effectiveness of ASEN for fine-grained fashion similarity learning and its potential for fashion reranking. Code and data are available at <https://github.com/Maryeon/asen>.

Introduction

Learning the similarity between fashion items is essential for a number of fashion-related tasks including in-shop clothes retrieval (Liu et al. 2016; Ak et al. 2018b), cross-domain fashion retrieval (Huang et al. 2015; Ji et al. 2017), fashion compatibility prediction (He, Packer, and McAuley 2016; Vasileva et al. 2018) and so on. The majority of methods are proposed to learn a general embedding space so the similarity can be computed in the space (Zhao et al. 2017; Ji et al. 2017; Han et al. 2017b). As the above tasks aim to search for identical or similar/compatible fashion items w.r.t. the query item, methods for these tasks tend to focus on the overall similarity. In this paper, we aim for the fine-grained fashion similarity. Consider the two fashion images in Fig. 1, although they appear to be irrelevant overall, they actually present similar characteristics over some attributes, *e.g.*, both of them have the similar lapel design. We consider such similarity in terms of a specific attribute as the fine-grained similarity.

*Zhe Ma, Jianfeng Dong and Yao Zhang are the co-first authors.

†Corresponding authors: Jianfeng Dong and Shouling Ji.

Copyright © 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

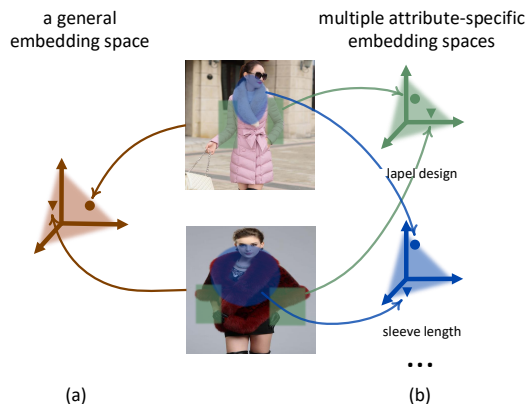


Figure 1: As fashion items typically have various attributes, we propose to learn multiple attribute-specific embeddings, thus the fine-grained similarity can be better reflected in the corresponding attribute-specific space.

There are cases where one would like to search for fashion items with certain similar designs instead of identical or overall similar items, so the fine-grained similarity matters in such cases. In the fashion copyright protection scenario (Martin 2019), the fine-grained similarity is also important to find items with plagiarized designs. Hence, learning the fine-grained similarity is necessary. However, to the best of our knowledge, such a similarity paradigm has been ignored by the community to some extent, only one work focuses on it. In (Veit, Belongie, and Karaletsos 2017), they first learn an overall embedding space, and then employ a fixed mask to select relevant embedding dimensions w.r.t. the specified attribute. The fine-grained similarity is measured in terms of the masked embedding feature. In this work, we go further in this direction. As shown in Fig. 1, we propose to learn multiple attribute-specific embedding spaces thus measure the fine-grained similarity in the corresponding space. For example, from the perspective of *neckline design*, the similarity between two clothes can be measured in the embedding space of *neckline design*. To this end, we propose an *Attribute-Specific Embedding Network (ASEN)* to jointly learn multiple attribute-specific embeddings in an end-to-end manner. Specifically, we introduce the novel *attribute-aware spatial attention (ASA)* and *attribute-aware channel*

attention (ACA) modules in the network, allowing the network being able to locate the related regions and capture the essential patterns w.r.t. the specified attribute. It is worth pointing out that fine-grained similarity learning is orthogonal to overall similarity learning, allowing us to utilize ASEN to facilitate traditional fashion retrieval, such as in-shop clothes retrieval. In sum, this paper makes the following contributions:

- Conceptually, we propose to learn multiple attribute-specific embedding spaces for fine-grained fashion similarity prediction. As such, a certain fine-grained similarity between fashion items can be measured in the corresponding space.
- We propose a novel ASEN model to effectively realize the above proposal. Combined with ACA and ASA, the network extracts essential features under the guidance of the specified attribute, which benefits the fine-grained similarity computation.
- Experiments on FashionAI, DARN, DeepFashion and Zappos50k datasets demonstrate the effectiveness of proposed ASEN for fine-grained fashion similarity learning and its potential for fashion reranking.

Related Work

Fashion Similarity Learning To compute the similarity between fashion items, the majority of existing works (Liu et al. 2016; Gajic and Baldrich 2018; Shankar et al. 2017; Ji et al. 2017; Huang et al. 2015) learn a general embedding space thus the similarity can be measured in the learned space by standard distance metric, *e.g.*, cosine distance. For instance, in the context of in-shop clothes retrieval, (Liu et al. 2016) employs a Convolutional Neural Network (CNN) to embed clothes into a single compact feature space. Similarly, for the purpose of fashion compatibility prediction, (Veit et al. 2015) also utilize a CNN to map fashion items in an embedding space, thus predict whether two input fashion items are compatible in the space. Different from the above methods that focus on the overall similarity (identical or overall similar/compatible), we study the fine-grained similarity in the paper. (Veit, Belongie, and Karaletsos 2017) have made a first attempt in this direction. In their approach, an overall embedding space is first learned, and the fine-grained similarity is measured in this space with the fixed mask w.r.t. a specified attribute. By contrast, we jointly learn multiple attribute-specific embedding spaces, and measure the fine-grained similarity in the corresponding attribute-specific space. It is worth noting that (Vasileva et al. 2018; He, Packer, and McAuley 2016) also learn multiple embedding spaces, but they still focus on the overall similarity.

Attention Mechanism Recently attention mechanism has become a popular technique and showed superior effectiveness in various research areas, such as computer vision (Woo et al. 2018; Wang et al. 2017a; Qiao, Dong, and Xu 2018) and natural language processing (Vaswani et al. 2017; Bahdanau, Cho, and Bengio 2014). To some extent, attention can be regarded as a tool to bias the allocation of the input information. As fashion images always present with

complex backgrounds, pose variations, etc., attention mechanism is also common in the fashion domain (Ji et al. 2017; Wang et al. 2017b; Han et al. 2017a; Ak et al. 2018a; Ak et al. 2018b). For instance, (Ak et al. 2018b) use the prior knowledge of clothes structure to locate the specific parts of clothes. However, their approach can be only used for upper-body clothes thus limits its generalization. (Wang et al. 2017b) propose to learn a channel attention implemented by a fully convolutional network. The above attentions are in a self-attention manner without explicit guidance for attention mechanism. In this paper, we propose two attribute-aware attention modules, which utilize a specific attribute as the extra input in addition to a given image. The proposed attention modules capture the attribute-related patterns under the guidance of the specified attribute. Note that (Ji et al. 2017) also utilize attributes to facilitate attention modeling, but they use all attributes of fashion items and aim for learning a better discriminative fashion feature. By contrast, we employ each attribute individually to obtain more fine-grained attribute-aware feature for fine-grained similarity computation.

Proposed Method

Network Structure

Given an image I and a specific attribute a , we propose to learn an attribute-specific feature vector $f(I, a) \in \mathbb{R}^c$ which reflects the characteristics of the corresponding attribute in the image. Therefore, for two fashion images I and I' , the fine-grained fashion similarity w.r.t. the attribute a can be expressed by the cosine similarity between $f(I, a)$ and $f(I', a)$. Moreover, the fine-grained similarity for multiple attributes can be computed by summing up the similarity scores on the individual attributes. Note that the attribute-specific feature vector resides in the corresponding attribute-specific embedding space. If there are n attributes, n attribute-specific embedding spaces can be learned jointly. Fig. 2 illustrates the structure of our proposed network. The network is composed of a feature extraction branch combined with an attribute-aware spatial attention and an attribute-aware channel attention. For the ease of reference, we name the two attention modules as ASA and ACA, respectively. In what follows, we first detail the input representation, followed by the description of two attribute-aware attention modules.

Input Representation To represent the image, we employ a CNN model pre-trained on ImageNet (Deng et al. 2009) as a backbone network, *e.g.*, ResNet (He et al. 2016). To keep the spatial information of the image, we remove the last fully connected layers in the CNN. So the image is represented by $I \in \mathbb{R}^{c \times h \times w}$, where $h \times w$ is the size of the feature map, c indicates the number of channels. For the attribute, we represent it with a one-hot vector $a \in \{0, 1\}^n$, where $n \in \mathbb{N}$ indicates the number of different attributes.

Attribute-aware Spatial Attention (ASA) Considering the attribute-specific feature is typically related to the specific regions of the image, we only need to focus on the certain related regions. For instance, in order to extract the

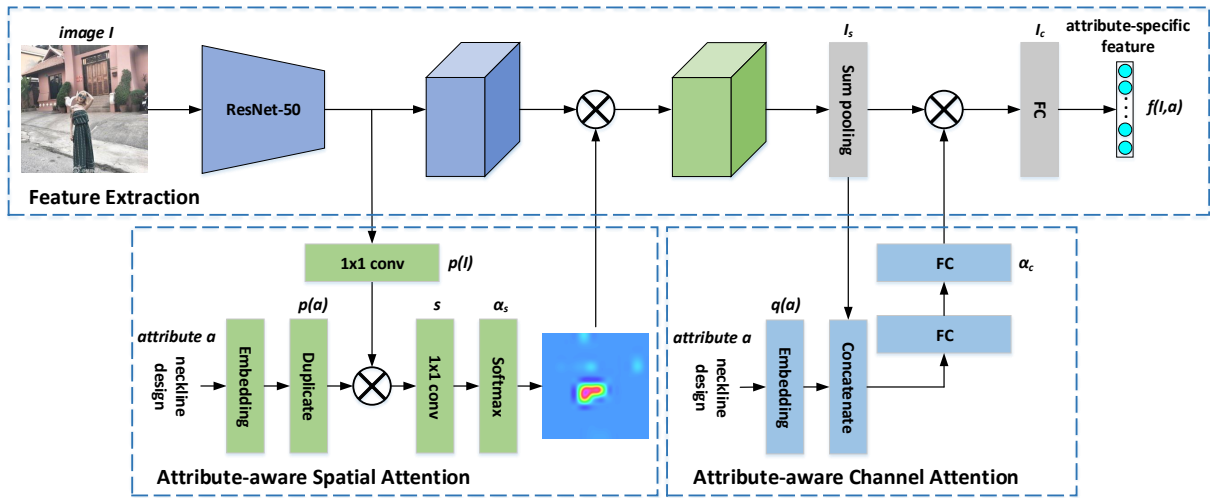


Figure 2: The structure of our proposed Attribute-Specific Embedding Network (ASEN). Mathematical notations by the side of function blocks (e.g., α_c on the right of FC layer) denotes their output.

attribute-specific feature of the *neckline design* attribute, the region around *neck* is much more important than the others. Besides, as fashion images always show up in large variations, e.g., various poses and scales, using a fixed region with respect to a specific attribute for all images is not optimal. Hence, we propose an attribute-aware spatial attention which adaptively attends to certain regions of the input image under the guidance of a specific attribute. Given an image I and a specific attribute a , we obtain the spatially attended vector w.r.t. the given attribute a by $I_s = Att_s(I, a)$, where the attended vector is computed as the weighted average of input image feature vectors according to the given attribute. Specifically, we first transform the image and the attribute to make their dimensionality same. For the image, we employ a convolutional layer followed by a nonlinear \tanh activation function. Formally, the mapped image $p(I) \in \mathbb{R}^{c' \times h \times w}$ is given by

$$p(I) = \tanh(Conv_{c'}(I)), \quad (1)$$

where $Conv_{c'}$ indicates a convolutional layer that contains c' 1×1 convolution kernels. For the attribute, we first project it into a c' -dimensional vector through an attribute embedding, implemented by a Fully Connected (FC) layer, then perform spatial duplication. Hence, the mapped attribute $p(a) \in \mathbb{R}^{c' \times h \times w}$ is

$$p(a) = \tanh(W_a a) \cdot 1, \quad (2)$$

Where $W_a \in \mathbb{R}^{c' \times n}$ denotes the transformation matrix and $1 \in \mathbb{R}^{1 \times h \times w}$ indicates spatially duplicate matrix. After the feature mapping, the attention weights $\alpha_s \in \mathbb{R}^{h \times w}$ is computed as

$$\begin{aligned} s &= \tanh(Conv_1(p(a) \odot p(I))), \\ \alpha_s &= softmax(s), \end{aligned} \quad (3)$$

where \odot indicates the element-wise multiplication, $Conv_1$ is a convolutional layer only containing one 1×1 convolution kernel. Here, we employ a $softmax$ layer to normalize

the attention weights. With adaptive attention weights, the spatially attended feature vector of the image I w.r.t. a specific attribute a is calculated as:

$$I_s = \sum_j^{h \times w} \alpha_{sj} I_j. \quad (4)$$

where $\alpha_{sj} \in \mathbb{R}$ and $I_j \in \mathbb{R}^c$ are the attention weight and the feature vector at location j of α_s and I , respectively.

Attribute-aware Channel Attention (ACA) Although the attribute-aware spatial attention adaptively focuses on the specific regions in the image, the same regions may still be related to multiple attributes. For example, attributes *collar design* and *collar color* are all associated with the region around *collar*. Hence, we further employ attribute-aware channel attention over the spatially attended feature vector I_s . The attribute-aware channel attention is designed as an element-wise gating function which selects the relevant dimensions of the spatially attended feature with respect to the given attribute. Concretely, we first employ an attribute embedding layer to embed attribute a into an embedding vector with the same dimensionality of I_s , that is:

$$q(a) = \delta(W_c a) \quad (5)$$

where $W_c \in \mathbb{R}^{c \times n}$ denotes the embedding parameters and δ refers to $ReLU$ function. Note we use separated attribute embedding layers in ASA and ACA, considering the different purposes of the two attentions. Then the attribute and the spatially attended feature are fused by simple concatenation, and further fed into the subsequent two FC layers to obtain the attribute-aware channel attention weights. As suggested in (Hu, Shen, and Sun 2018), we implement the two FC layers by a dimensionality-reduction layer with reduction rate r and a dimensionality-increasing layer, which have fewer parameters than one FC layer. Formally, the attention weights $\alpha_c \in \mathbb{R}^c$ is calculated by:

$$\alpha_c = \sigma(W_2 \delta(W_1 [q(a), I_s])), \quad (6)$$

where $[,]$ denotes concatenation operation, σ indicates *sigmoid* function, $W_1 \in \mathbb{R}^{\frac{c}{r} \times 2c}$ and $W_2 \in \mathbb{R}^{c \times \frac{c}{r}}$ are transformation matrices. Here we omit the bias terms for description simplicity. The final output of the ACA is obtained by scaling I_s with the attention weight α_c :

$$I_c = I_s \odot \alpha_c. \quad (7)$$

Finally, we further employ a FC layer over I_c to generate the attribute-specific feature of the given image I with the specified attribute a :

$$f(I, a) = WI_c + b, \quad (8)$$

where $W \in \mathbb{R}^{c \times c}$ is the transformation matrix, $b \in \mathbb{R}^c$ indicates the bias term.

Model Learning

We would like to achieve multiple attribute-specific embedding spaces where the distance in a particular space is small for images with the same specific attribute value, but large for those with the different ones. Consider the *neckline design* attribute for instance, we expect the fashion images with *Round Neck* near those with the same *Round Neck* in the *neckline design* embedding space, but far away from those with *V Neck*. To this end, we choose to use the triplet ranking loss which is consistently found to be effective in multiple embedding learning tasks (Vasileva et al. 2018; Dong et al. 2019). Concretely, we first construct a set of triplets $\mathcal{T} = \{(I, I^+, I^-|a)\}$, where I^+ and I^- indicate images relevant and irrelevant with respect to image I in terms of attribute a . Given a triplet of $\{(I, I^+, I^-|a)\}$, triplet ranking loss is defined as

$$\mathcal{L}(I, I^+, I^-|a) = \max\{0, m - s(I, I^+|a) + s(I, I^-|a)\}, \quad (9)$$

where m represents the margin, empirically set to be 0.2, $s(I, I'|a)$ denotes the fine-grained similarity w.r.t. the attribute a which can be expressed by the cosine similarity between $f(I, a)$ and $f(I', a)$. Finally, we train the model to minimize the triplet ranking loss on the triplet set \mathcal{T} , and the overall objective function of the model is as:

$$\operatorname{argmin}_{\theta} \sum_{(I, I^+, I^-|a) \in \mathcal{T}} \mathcal{L}(I, I^+, I^-|a), \quad (10)$$

where θ denotes all trainable parameters of our proposed network.

Evaluation

Experimental Setup

To verify the viability of the proposed attribute-specific embedding network for fine-grained fashion similarity computation, we evaluate it on the following two tasks. (1) Attribute-specific fashion retrieval: Given a fashion image and a specified attribute, its goal is to search for fashion images of the same attribute value with the given image. (2) Triplet relation prediction: Given a triplet of $\{I, I', I''\}$ and a specified attribute, the task is asked to predict whether the relevance between I and I' is larger than that between I and I'' in terms of the given attribute.

Datasets As there are no existing datasets for attribute-specific fashion retrieval, we reconstruct three fashion datasets with attribute annotations to fit the task, *i.e.*, FashionAI (Zou et al. 2019), DARN (Huang et al. 2015) and DeepFashion(Liu et al. 2016). For triplet relation prediction, we utilize Zappos50k (Yu and Grauman 2014). *FashionAI* is a large scale fashion dataset with hierarchical attribute annotations for fashion understanding. We choose to use the FashionAI dataset, because of its high-quality attribute annotations. As the full FashionAI has not been publicly released, we utilize its early version released for the FashionAI Global Challenge 2018¹. The released FashionAI dataset consists of 180,335 apparel images, where each image is annotated with a fine-grained attribute. There are 8 attributes, and each attribute is associated with a list of attribute values. Take the attribute *neckline design* for instance, there are 11 corresponding attribute values, such as *round neckline* and *v neckline*. We randomly split images into three sets by 8:1:1, which is 144k / 18k / 18k images for training / validation / test. Besides, for every epoch, we construct 100k triplets from the training set for model training. Concretely, for a triplet with respect to a specific attribute, we randomly sample two images of the same corresponding attribute values as the relevant pair and an image with different attribute value as the irrelevant one. For validation or test set, 3600 images are randomly picked out as the query images, with remaining images annotated with the same attribute as the candidate images for retrieval. Additionally, we reconstruct DARN and DeepFashion in the same way as FashionAI. Details are included in the supplementary material.

Zappos50k is a large shoe dataset consisting of 50,025 images collected from the online shoe and clothing retailer Zappos.com. For the ease of cross-paper comparison, we utilize the identical split provided by (Veit, Belongie, and Karalestos 2017). Specifically, we use 70% / 10% / 20% images for training / validation / test. Each image is associated with four attributes: the type of the shoes, the suggested gender of the shoes, the height of the shoes' heels and the closing mechanism of the shoes. For each attribute, 200k training, 20k validation and 40k testing triplets are sampled for model training and evaluation.

Metrics For the task of attribute-specific fashion retrieval, we report the Mean Average Precision (MAP), a popular performance metric in many retrieval-related tasks (Awad et al. 2018; Dong, Li, and Xu 2018). For the triplet relation prediction task, we utilize the prediction accuracy as the metric.

Due to the limited space of the paper, we present results on DeepFashion, implementation details and efficiency evaluation of our proposed model in the supplementary material.

Attribute-Specific Fashion Retrieval

Table 1 summarizes the performance of different models on FashionAI, and performance of each attribute type are also reported. As a sanity check, we also give the performance of a random baseline which sorts candidate images randomly. All the learning methods are noticeably better than the random result. Among the five learning based

¹<https://tianchi.aliyun.com/markets/tianchi/FashionAI>

Table 1: Performance of attribute-specific fashion retrieval on FashionAI. Our proposed ASEN model consistently outperforms the other counterparts for all attribute types.

Method	MAP for each attribute								MAP
	skirt length	sleeve length	coat length	pant length	collar design	lapel design	neckline design	neck design	
Random baseline	17.20	12.50	13.35	17.45	22.36	21.63	11.09	21.19	15.79
Triplet network	48.38	28.14	29.82	54.56	62.58	38.31	26.64	40.02	38.52
CSN	61.97	45.06	47.30	62.85	69.83	54.14	46.56	54.47	53.52
ASEN w/o ASA	62.65	49.98	49.02	63.48	69.10	61.65	50.88	57.10	56.35
ASEN w/o ACA	58.12	43.30	42.30	60.03	65.98	49.95	46.86	52.06	50.87
ASEN	64.44	54.63	51.27	63.53	70.79	65.36	59.50	58.67	61.02

Table 2: Performance of attribute-specific fashion retrieval on DARN. AESN with both ACA and ASA again performs best.

Method	MAP for each attribute									MAP
	clothes category	clothes button	clothes color	clothes length	clothes pattern	clothes shape	collar shape	sleeve length	sleeve shape	
Random baseline	8.49	24.45	12.54	29.90	43.26	39.76	15.22	63.03	55.54	32.26
Triplet network	23.59	38.07	16.83	39.77	49.56	47.00	23.43	68.49	56.48	40.14
CSN	34.10	44.32	47.38	53.68	54.09	56.32	31.82	78.05	58.76	50.86
ASEN w/o ASA	33.94	45.37	48.56	54.36	53.83	57.33	32.78	77.77	59.32	51.39
ASEN w/o ACA	30.39	42.37	49.14	50.18	53.63	48.84	26.03	75.28	57.99	48.02
ASEN	36.69	46.96	51.35	56.47	54.49	60.02	34.18	80.11	60.04	53.31

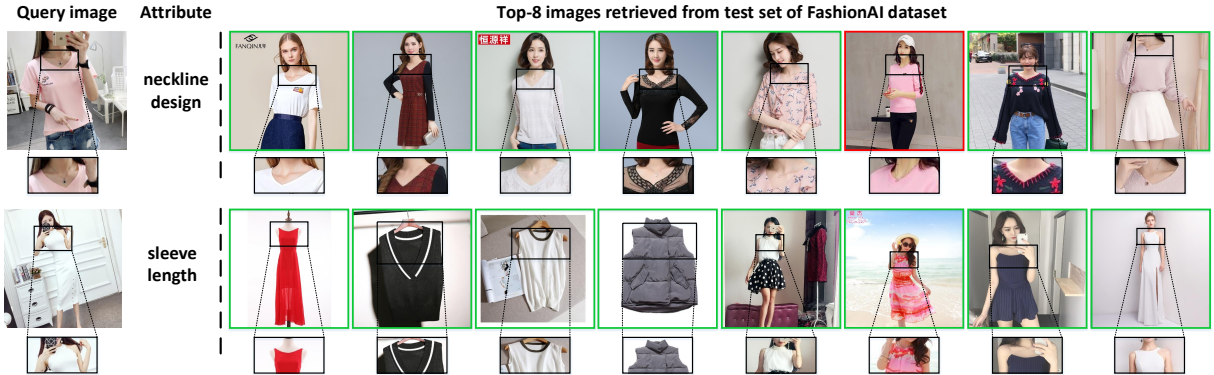


Figure 3: Attribute-specific fashion retrieval examples on FashionAI. Green bounding box indicates the image has the same attribute value with the given image in terms of the given attribute, while the red one indicates the different attribute values. The results demonstrate that our ASEN is good at capturing the fine-grained similarity among fashion items.

models, the triplet network which learns a general embedding space performs the worst in terms of the overall performance, scoring the overall MAP of 38.52%. The result shows that a general embedding space is suboptimal for fine-grained similarity computation. Besides, our proposed ASEN outperforms CSN (Veit, Belongie, and Karaletsos 2017) with a clear margin. We attribute the better performance to the fact that ASEN adaptively extracts feature w.r.t. the given attribute by two attention modules, while CSN uses a fixed mask to select relevant embedding dimensions. Moreover, we investigate ASEN with a single attention, resulting in two reduced models, *i.e.*, ASEN w/o ASA and ASEN w/o ACA. These two variants obtain the overall MAP of 56.35 and 50.87, respectively. The lower scores justify the necessity of both ASA and ACA attentions. The result also suggests that attribute-aware channel attention is more beneficial. Table 2 shows the results on

the DARN dataset. Similarly, our proposed ASEN outperforms the other counterparts. The result again confirms the effectiveness of the proposed model for fine-grained fashion similarity computation. Additionally, we also try the verification loss (Zheng, Zheng, and Yang 2017) in ASEN, but find its performance (MAP=50.63) worse than the triplet loss counterpart (MAP=61.02) on FashionAI. Some qualitative results of ASEN are shown in Fig. 3. Note that the retrieved images appear to be irrelevant to the query image, as ASEN focuses on the fine-grained similarity instead of the overall similarity. It can be observed that the majority of retrieved images share the same specified attribute with the query image. Consider the second example for instance, although the retrieved images are in various fashion category, such as *dress* and *vest*, all of them are *sleeveless*. These results allow us to conclude that our model is able to figure out fine-grained patterns in images.

Table 3: Performance of triplet relation prediction on Zappos50k. Our proposed ASEN is the best.

Method	Prediction Accuracy(%)
Random baseline	50.00
Triplet network	76.28
CSN	89.27
ASEN w/o ASA	90.18
ASEN w/o ACA	89.01
ASEN	90.79

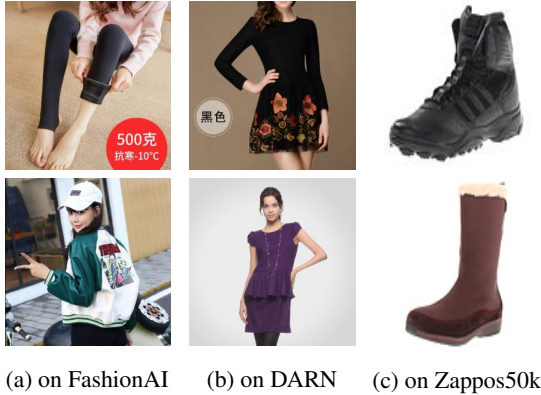


Figure 4: Images from FashionAI, DARN and Zappos50k, showing the images of Zappos50k are less challenging.

Triplet Relation Prediction

Table 3 shows the results on the Zappos50k dataset. Unsurprisingly, the random baseline achieves the worst performance as it predicts by random guess. Among the four embedding learning models, our proposed model variants again outperform triplet network which only learns a general embedding space with a large margin. The result verifies the effectiveness of learning attribute-specific embeddings for triplet relation prediction. Although ASEN is still better than its counterparts ASEN without ASA or ACA, its performance improvement is much less than that on FashionAI and DARN. We attribute it to that images in Zappos50k are more iconic and thus easier to understand (see Fig. 4), so only ASA or ACA is enough to capture the fine-grained similarity for such “easy” images.

What has ASEN Learned?

t-SNE Visualization In order to investigate what has the proposed ASEN learned, we first visualize the obtained attribute-specific embedding spaces. Specifically, we take all test images from FashionAI, and use t-SNE (Maaten and Hinton 2008) to visualize their distribution in 2-dimensional spaces. Fig. 5 presents eight attribute-specific embedding spaces w.r.t. *coat*, *pant*, *sleeve*, *skirt length* and *lapel*, *neck*, *neckline*, *collar design* respectively. It is clear that dots with different colors are well separated and dots with the same color are more clustered in the particular embedding

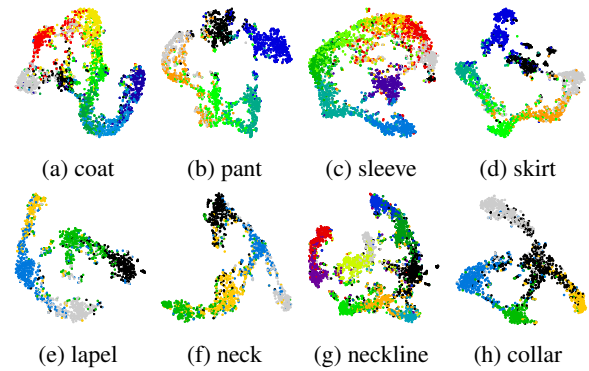


Figure 5: t-SNE visualization of attribute-specific embedding spaces obtained by our proposed ASEN on FashionAI dataset. Dots with the same color indicate images annotated with the same attribute value. Best viewed in zoom in.

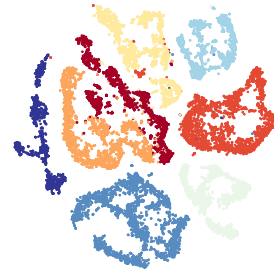


Figure 6: t-SNE visualization of a whole embedding space comprised of eight attribute-specific embedding spaces learned by ASEN. Dots with the same color indicate images in the same attribute-specific embedding space.

space. In other words, images with the same attribute value are close while images with different attribute value are far away. The result shows the good discriminatory ability of the learned attribute-specific embeddings by ASEN. One may ask what is the relationship between the attribute-specific embedding spaces? To answer this question, we visualize eight attribute-specific embeddings into a whole 2-dimensional space. As shown in Fig. 6, different attribute-specific embeddings learned by ASEN are well separated. The result is consistent with the fact that different attributes reflect different characteristics of fashion items.

Attention Visualization To gain further insight of our proposed network, we visualize the learned attribute-aware spatial attention. As shown in Fig. 7, the learned attention map gives relative high responses on the relevant regions while low responses on irrelevant regions with the specified attribute, showing the attention is able to figure out which regions are more important for a specific attribute. An interesting phenomenon can be observed that attention maps for length-related attributes are more complicated than that for design-related attributes; multiple regions show high response for the former. We attribute it to that the model requires to locate the start and end of a fashion item thus spec-

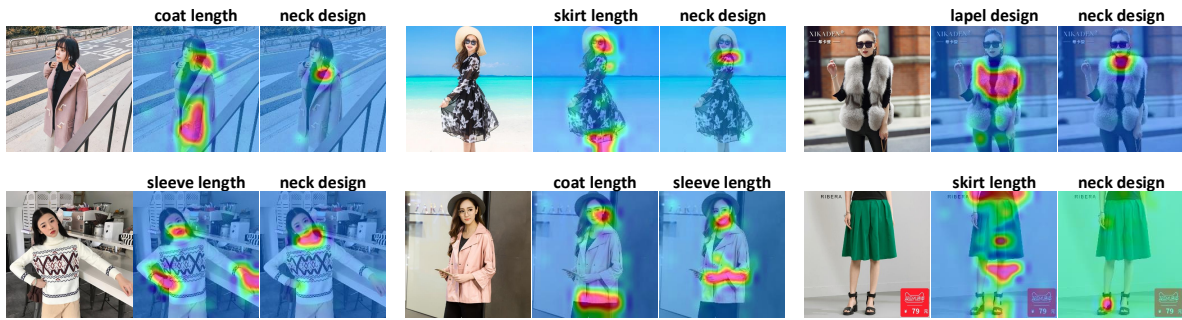


Figure 7: Visualization of the attribute-aware spatial attention with the guidance of a specified attribute (above the attention image) on FashionAI.

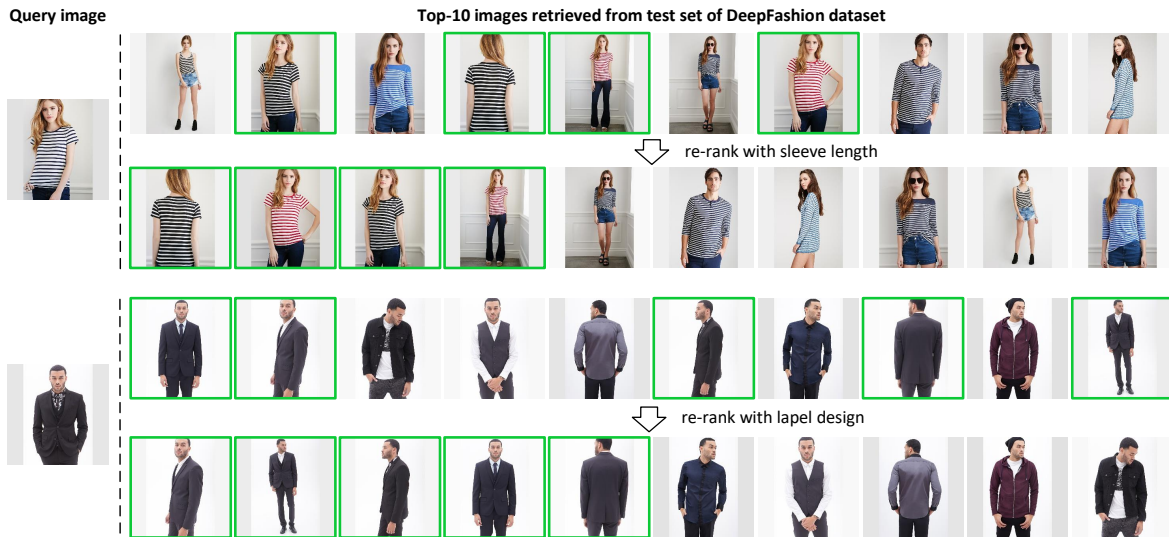


Figure 8: Reranking examples for in-shop clothes retrieval on DeepFashion dataset. The ground-truth images are marked with green bounding box. After reranking by our proposed ASEN, the retrieval results become better.

ulate its length. Besides, consider the last example with respect to *neck design* attribute, when the specified attribute *neck design* can not be reflected in the image the attention response is almost uniform, which further demonstrates the effectiveness of attribute-aware spatial attention.

The Potential for Fashion Reranking

In this experiment, we explore the potential of ASEN for fashion reranking. Specifically, we consider in-shop clothes retrieval task, in which given a query of in-shop clothes, the task is asked to retrieve the same items. Triplet network is used as the baseline to obtain the initial retrieval result. The initial top 10 images are reranked in descending order by the fine-grained fashion similarity obtained by ASEN. We train the triplet network on the official training set of DeepFashion, and directly use ASEN previously trained on FashionAI for the attribute-specific fashion retrieval task. Fig. 8 presents two reranking examples. For the first example, by reranking in terms of the fine-grained similarity of *sleeve length*, images have the same short sleeve with the query image are ranked higher, while the others with mid or long

sleeve are ranked later. Obviously, after reranking, the retrieval results become better. The result shows the potential of our proposed ASEN for fashion reranking.

Summary and Conclusions

This paper targets at the fine-grained similarity in the fashion scenario. We contribute an *Attribute-Specific Embedding Network (ASEN)* with two attention modules, *i.e.*, ASA and ACA. ASEN jointly learns multiple attribute-specific embeddings, thus measuring the fine-grained similarity in the corresponding space. ASEN is conceptually simple, practically effective and end-to-end. Extensive experiments on various datasets support the following conclusions. For fine-grained similarity computation, learning multiple attribute-specific embedding spaces is better than learning a single general embedding space. ASEN with only ACA is more beneficial when compared with its counterpart with only ASA. For state-of-the-art performance, we recommend ASEN with both attention modules.

Acknowledgments

This work was partly supported by the National Key Research and Development Program of China under No. 2018YFB0804102, NSFC under No. 61772466, 61902347 and U1836202, the Zhejiang Provincial Natural Science Foundation for Distinguished Young Scholars under No. LR19F020003, the Provincial Key Research and Development Program of Zhejiang, China under No. 2017C01055, the Zhejiang Provincial Natural Science Foundation under No. LQ19F020002, and the Alibaba-ZJU Joint Research Institute of Frontier Technologies.

References

- [Ak et al. 2018a] Ak, K. E.; Kassim, A. A.; Hwee Lim, J.; and Yew Tham, J. 2018a. Learning attribute representations with localization for flexible fashion search. In *CVPR*, 7708–7717.
- [Ak et al. 2018b] Ak, K. E.; Lim, J. H.; Tham, J. Y.; and Kassim, A. A. 2018b. Efficient multi-attribute similarity learning towards attribute-based fashion search. In *WACV*, 1671–1679.
- [Awad et al. 2018] Awad, G.; Butt, A.; Curtis, K.; Lee, Y.; Fiscus, J.; Godil, A.; Joy, D.; Delgado, A.; Smeaton, A.; Graham, Y.; et al. 2018. Trecvid 2018: Benchmarking video activity detection, video captioning and matching, video storytelling linking and video search. In *TRECVID Workshop*.
- [Bahdanau, Cho, and Bengio 2014] Bahdanau, D.; Cho, K.; and Bengio, Y. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- [Deng et al. 2009] Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*, 248–255.
- [Dong et al. 2019] Dong, J.; Li, X.; Xu, C.; Ji, S.; He, Y.; Yang, G.; and Wang, X. 2019. Dual encoding for zero-example video retrieval. In *CVPR*, 9346–9355.
- [Dong, Li, and Xu 2018] Dong, J.; Li, X.; and Xu, D. 2018. Cross-media similarity evaluation for web image retrieval in the wild. *IEEE Transactions on Multimedia* 20(9):2371–2384.
- [Gajic and Baldrich 2018] Gajic, B., and Baldrich, R. 2018. Cross-domain fashion image retrieval. In *CVPR Workshop*, 1869–1871.
- [Han et al. 2017a] Han, X.; Wu, Z.; Huang, P. X.; Zhang, X.; Zhu, M.; Li, Y.; Zhao, Y.; and Davis, L. S. 2017a. Automatic spatially-aware fashion concept discovery. In *ICCV*, 1463–1471.
- [Han et al. 2017b] Han, X.; Wu, Z.; Jiang, Y.-G.; and Davis, L. S. 2017b. Learning fashion compatibility with bidirectional lstms. In *ACM Multimedia*, 1078–1086.
- [He et al. 2016] He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *CVPR*, 770–778.
- [He, Packer, and McAuley 2016] He, R.; Packer, C.; and McAuley, J. 2016. Learning compatibility across categories for heterogeneous item recommendation. In *ICDM*, 937–942.
- [Hu, Shen, and Sun 2018] Hu, J.; Shen, L.; and Sun, G. 2018. Squeeze-and-excitation networks. In *CVPR*, 7132–7141.
- [Huang et al. 2015] Huang, J.; Feris, R. S.; Chen, Q.; and Yan, S. 2015. Cross-domain image retrieval with a dual attribute-aware ranking network. In *ICCV*, 1062–1070.
- [Ji et al. 2017] Ji, X.; Wang, W.; Zhang, M.; and Yang, Y. 2017. Cross-domain image retrieval with attention modeling. In *ACM Multimedia*, 1654–1662.
- [Kingma and Ba 2014] Kingma, D. P., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [Liu et al. 2016] Liu, Z.; Luo, P.; Qiu, S.; Wang, X.; and Tang, X. 2016. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *CVPR*, 1096–1104.
- [Maaten and Hinton 2008] Maaten, L. v. d., and Hinton, G. 2008. Visualizing data using t-sne. *Journal of Machine Learning Research* 9(Nov):2579–2605.
- [Martin 2019] Martin, T. 2019. Fashion law needs custom tailored protection for designs. *University of Baltimore Law Review* 48:453.
- [Qiao, Dong, and Xu 2018] Qiao, T.; Dong, J.; and Xu, D. 2018. Exploring human-like attention supervision in visual question answering. In *AAAI*, 7300–7307.
- [Shankar et al. 2017] Shankar, D.; Narumanchi, S.; Ananya, H.; Kompalli, P.; and Chaudhury, K. 2017. Deep learning based large scale visual recommendation and search for e-commerce. *arXiv preprint arXiv:1703.02344*.
- [Vasileva et al. 2018] Vasileva, M. I.; Plummer, B. A.; Dusad, K.; Rajpal, S.; Kumar, R.; and Forsyth, D. 2018. Learning type-aware embeddings for fashion compatibility. In *ECCV*, 390–405.
- [Vaswani et al. 2017] Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *NIPS*, 5998–6008.
- [Veit, Belongie, and Karaletsos 2017] Veit, A.; Belongie, S. J.; and Karaletsos, T. 2017. Conditional similarity networks. In *CVPR*, 830–838.
- [Veit et al. 2015] Veit, A.; Kovacs, B.; Bell, S.; McAuley, J.; Bala, K.; and Belongie, S. 2015. Learning visual clothing style with heterogeneous dyadic co-occurrences. In *ICCV*, 4642–4650.
- [Wang et al. 2017a] Wang, F.; Jiang, M.; Qian, C.; Yang, S.; Li, C.; Zhang, H.; Wang, X.; and Tang, X. 2017a. Residual attention network for image classification. In *CVPR*, 3156–3164.
- [Wang et al. 2017b] Wang, Z.; Gu, Y.; Zhang, Y.; Zhou, J.; and Gu, X. 2017b. Clothing retrieval with visual attention model. In *VCIP*, 1–4.
- [Woo et al. 2018] Woo, S.; Park, J.; Lee, J.-Y.; and

- So Kweon, I. 2018. Cbam: Convolutional block attention module. In *ECCV*, 3–19.
- [Yu and Grauman 2014] Yu, A., and Grauman, K. 2014. Fine-grained visual comparisons with local learning. In *CVPR*, 192–199.
- [Zhao et al. 2017] Zhao, B.; Feng, J.; Wu, X.; and Yan, S. 2017. Memory-augmented attribute manipulation networks for interactive fashion search. In *CVPR*, 1520–1528.
- [Zheng, Zheng, and Yang 2017] Zheng, Z.; Zheng, L.; and Yang, Y. 2017. A discriminatively learned cnn embedding for person reidentification. *ACM Transactions on Multimedia Computing, Communications, and Applications* 14(1):13.
- [Zou et al. 2019] Zou, X.; Kong, X.; Wong, W.; Wang, C.; Liu, Y.; and Cao, Y. 2019. Fashionai: A hierarchical dataset for fashion understanding. In *CVPR Workshop*.

Experimental Details

Dataset

DARN (Huang et al. 2015) is constructed for attribute prediction and street-to-shop retrieval task. The dataset contains 253,983 images, each of which is annotated with 9 attributes. As some images' URLs have been broken, only 214,619 images are obtained for our experiments. Images are randomly divided by 8:1:1 for training, validation and test, resulting in 171k, 43k, 43k images respectively. Triplets are similarly sampled as on FashionAI dataset where two images with the same attribute values for an attribute form a positive pair and an image with different attribute value is randomly sampled as a negative one. We randomly choose 100k triplets for training. For validation set and test set, images are split as query and candidate images by 1:4.

DeepFashion (Liu et al. 2016) is a large dataset which consists of four benchmarks for various tasks in the field of clothing including category and attribute prediction, in-shop clothes retrieval, fashion landmark detection and consumer-to-shop clothes retrieval. In our experiments, we use the category and attribute prediction benchmark for attribute-specific retrieval task and in-shop clothes retrieval benchmark for fashion reranking task.

The category and attribute prediction benchmark contains 289,222 images, 6 attributes and 1050 attribute values, and each image is annotated with several attributes. Similar to *DARN*, we randomly split the images into training, validation and test set by 8:1:1 and construct 100k triplets for training. For validation set and test set, images are also split as query and candidate images by 1:4.

The in-shop clothes retrieval benchmark has a total of 52,712 images, consisting of 7,982 different items. Each item corresponds to multiple images by changing the color or angle of shot. The images are officially divided into training, query, and gallery set, with 25k, 14k and 12k images. We keep the training set unchanged in the experiment, and respectively divide the gallery and the query images to validation and test by 1:1.

Table 5 summarizes the attribute annotations and corresponding attribute values on the FashionAI, Zappos50k, *DARN* and category and attribute prediction benchmark of *DeepFashion* dataset.

Compared Models

The conceptual structures of the five models are demonstrated in Fig. 9. All of them are based on the same CNN backbone.

- *Triplet network*: This model learns a general embedding space to measure the fine-grained fashion similarity. It simply ignores attributes and performs mean pooling on feature map generated by CNN. The standard triplet ranking loss is used to train the model.

- *Conditional similarity network (CSN)* (Veit, Belongie, and Karaletsos 2017): This model first learns an overall embedding space and then employs a fixed mask to select relevant embedding dimensions w.r.t. the specified attribute.

- *ASEN w/o ASA*: The model employs mean pooling instead of attribute-aware spatial attention module to aggregate

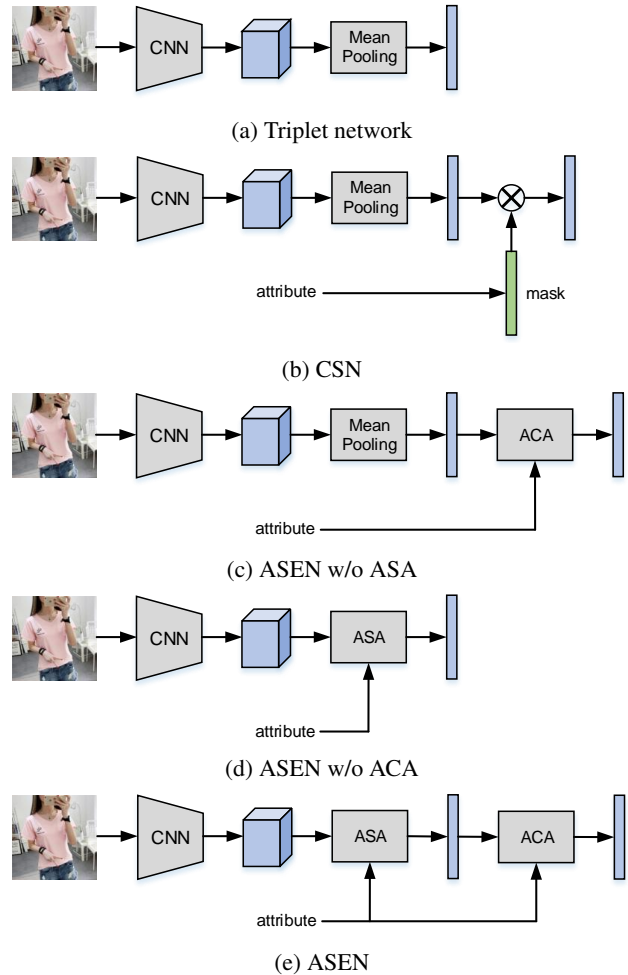


Figure 9: The conceptual structure of compared models and our proposed ASEN.

gate features.

- *ASEN w/o ACA*: The model utilizes the vector I_s as the attribute-specific feature vector without employing the attribute-aware channel attention.

- *ASEN*: It is our proposed model for fine-grained fashion similarity learning.

Implementation Details

On the FashionAI, *DARN* and *DeepFashion* dataset, we adopt the same setting. We use ResNet-50 network pre-trained on ImageNet (Deng et al. 2009) as the backbone network. Input images are first resized to 224 by 224. The dimension of the attribute-specific embedding is set to be 1024.

On Zappos50K, we follow the same experimental setting as (Veit, Belongie, and Karaletsos 2017) for a fair comparison. Concretely, we utilize ResNet-18 network pre-trained on ImageNet as the backbone network, resize input images to 112 by 112 and set the dimension of the attribute-specific embedding to 64.

Table 4: Performance of attribute-specific fashion retrieval on DeepFashion. Our proposed ASEN model consistently outperforms the other counterparts for all attribute types.

Method	MAP for each attribute type					MAP
	texture-related	fabric-related	shape-related	part-related	style-related	
Random baseline	6.69	2.69	3.23	2.55	1.97	3.38
Triplet network	13.26	6.28	9.49	4.43	3.33	7.36
CSN	14.09	6.39	11.07	5.13	3.49	8.01
ASEN w/o ASA	14.73	6.90	11.73	5.25	3.43	8.41
ASEN w/o ACA	14.71	6.78	11.71	5.03	3.42	8.32
ASEN	15.13	7.11	12.39	5.51	3.56	8.74

Table 5: Summary of attribute and its corresponding attribute values on FashionAI, Zappos50k, DARN and DeepFashion.

Dataset	Attribute	Values	Total
FashionAI	skirt length	short length, knee length, midi length, ankle length, floor length, invisible	6
	sleeve length	sleeveless, cup sleeves, short sleeves...	9
	coat length	high waist length, regular length, long length...	8
	pant length	short pant, mid length, 3/4 length, cropped pant, full length, invisible	6
	collar design	shirt collar, peter pan, puritan collar, rib collar, invisible	5
	lapel design	notched, collarless, shawl collar, plus size shawl, invisible	5
	neckline design	strapless neck, deep V neckline, straight neck, V neckline, invisible...	11
	neck design	invisible, turtle neck, ruffle semi-high collar, low turtle neck, draped collar	5
Zappos50k	category	shoes, boots, sandals, slippers	4
	gender	women, men, girls, boys	4
	heel height	1-4in, 5in&over, flat, under 1in	7
	closure	buckle, pull on, slip on, hook and loop...	19
DARN	clothes category	formal skirt, cotton clothes, lace shirt, small suit, shirt, knitwear, fur_clothes...	20
	clothes button	pullover, zipper, single breasted type1, single breasted type2...	13
	clothes color	yellow, apricot, flower colors, red, green, white, rose_red...	55
	clothes length	mid-long, long, short, normal, ultra-long, ultra-short	7
	clothes pattern	solid color, lattice, flower, animal, abstract, floral...	28
	clothes shaoe	slim, straight, cloak, loose, high_waist, shape1, A-shape...	11
	collat shaoe	polo, shape1, round, shape2, V_shape, boat_neck, ruffle_collar...	26
	sleeve length	long sleeve, three-quarter sleeve, sleeve1, sleeveless...	8
	sleeve shape	puff sleeve, regular, lantern sleeve, pile sleeve...	17
DeepFashion	texture-related	abstract, animal, bandana, baroque, bird, breton, butterfly...	156
	fabric-related	acid, applique, bead, bejeweled, cable, canvas, chenille, chino...	218
	shape-related	a-line, ankle, asymmetric, baja, bermuda, bodycon, box...	180
	part-related	arrow collar, back bow, batwing, bell, button, cinched, collared...	216
	style-related	americana, art, athletic, barbie, beach, bella, blah, boho...	230

On all datasets, the model is trained by ADAM optimizer (Kingma and Ba 2014) with an initial learning rate of $1E-4$. The learning rate is decayed by multiplying it with 0.985 after each epoch. We train every model for 200 epochs, and the snapshot with the highest performance on validation set is reserved for evaluation on the test set.

More Experiment Results

Attribute-Specific Retrieval on DeepFashion Table 4 summarizes the results on DeepFashion. Our proposed network outperforms the other counterparts. The result again verifies the effectiveness of our proposed ASEN. But the MAP scores on DeepFashion of all models are relatively worse compared to that on FashionAI and DARN. We attribute it to the relative low annotation quality of DeepFash-

ion. For instance, only 77.8% of images annotated with *A-line* of Shape type are correctly labelled (Zou et al. 2019).

The Qualitative Results We demonstrate more attribute-specific fashion retrieval examples in Fig. 10, attention visualization examples in Fig. 11 and fashion reranking examples in Fig. 12 and Fig. 13.

Efficiency Evaluation Once our model is trained, given an image and a specific attribute, it takes approximately 0.3 milliseconds to extract the corresponding attribute-aware feature. The speed is fast enough for attribute-based fashion applications. The performance is tested on a server with an Intel Xeon E5-2680 v4 CPU, 256 G RAM and a GeForce GTX 1080 Ti GPU.



Figure 10: Attribute-specific fashion retrieval examples on FashionAI. Green bounding box indicates the image has the same attribute value with the given image in terms of the given attribute, while the red one indicates the different attribute values.



Figure 11: Visualization of attribute-specific spatial attention with the guidance of a specified attribute (above the attention image) on FashionAI.

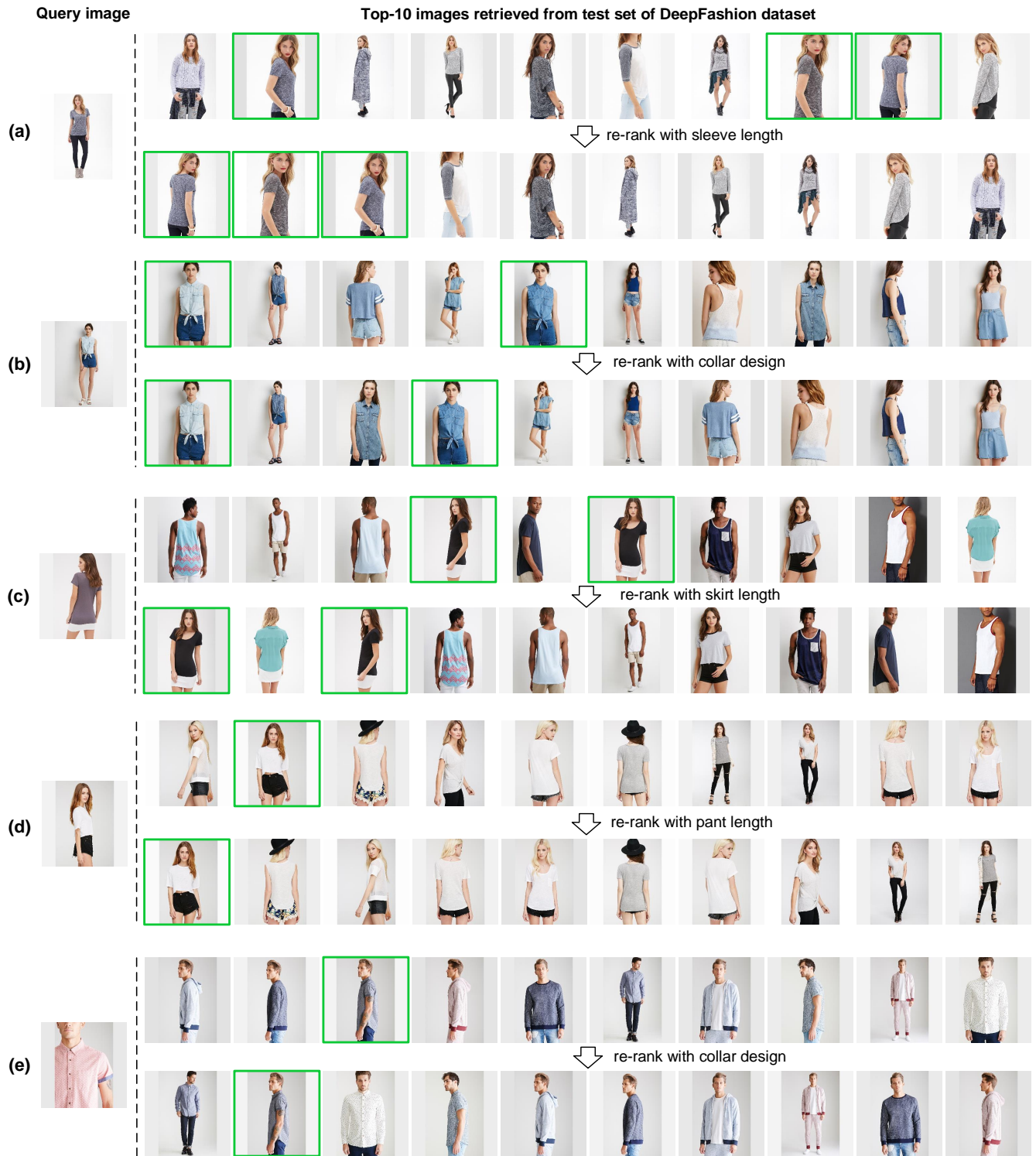


Figure 12: Reranking examples for in-shop clothes retrieval on DeepFashion. After the fashion reranking by our proposed ASEN, the results look better.

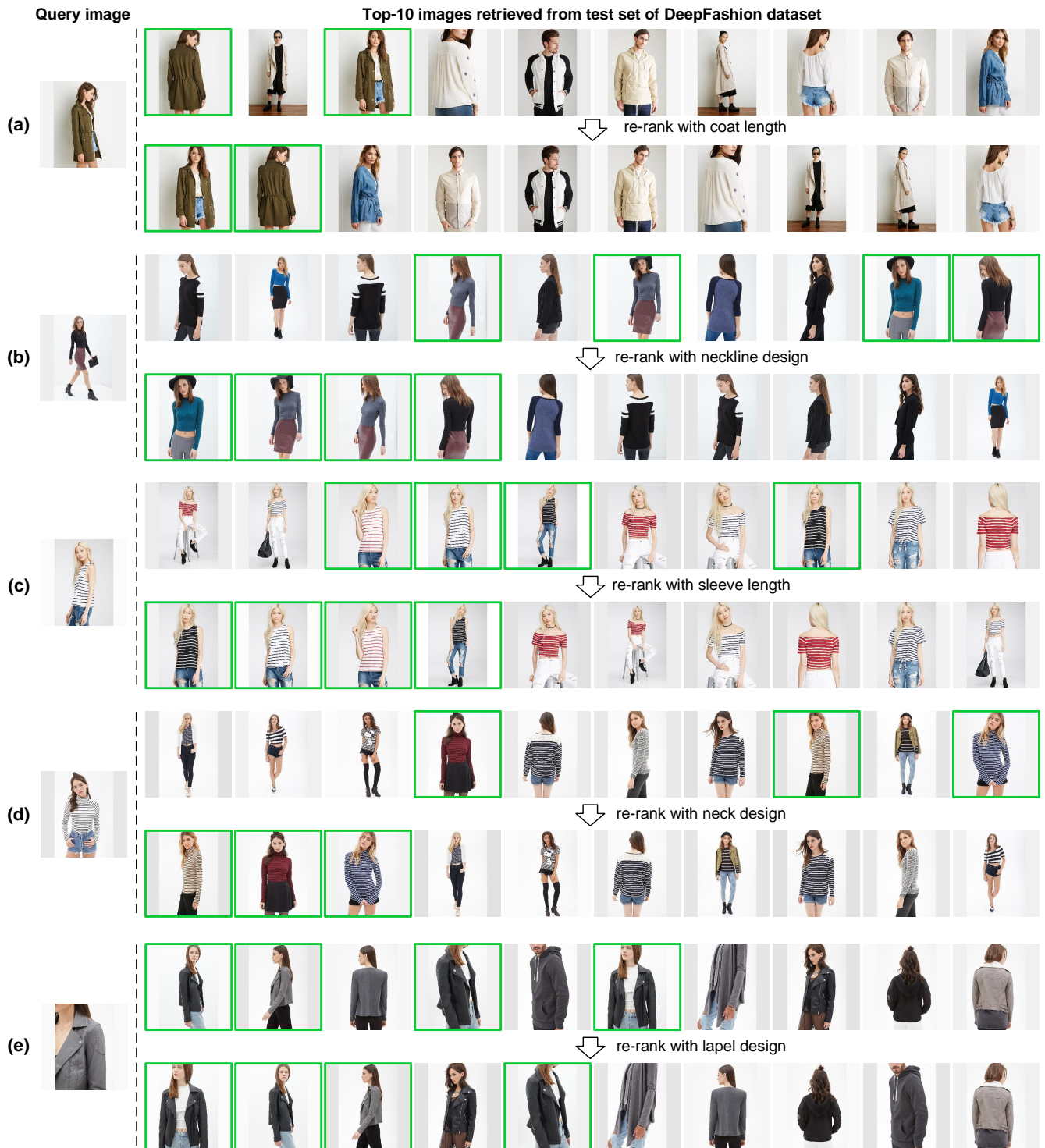


Figure 13: Reranking examples for in-shop clothes retrieval on DeepFashion. After the fashion reranking by our proposed ASEN, the results look better.