

# The iMaterialist Fashion Attribute Dataset

Sheng Guo<sup>1</sup> Weilin Huang<sup>1</sup> Xiao Zhang<sup>2</sup> Prasanna Srikhanta<sup>3</sup> Yin Cui<sup>4</sup> Yuan Li<sup>5</sup>

Hartwig Adam<sup>2</sup> Matthew R. Scott<sup>1</sup> Serge Belongie<sup>4</sup>

<sup>1</sup>Malong Technologies <sup>2</sup>Google AI <sup>3</sup>Wish <sup>4</sup>Cornell University <sup>5</sup>Horizon Robotics

## Abstract

Many large-scale image databases such as ImageNet were constructed only for single-label and coarse object-level classification, while multiple labels and fine-grained categories are often needed in real-world applications, yet very few such datasets exist publicly. In this work, we contribute to the community a new dataset called iMaterialist Fashion Attribute (iFashion) to address this problem in the fashion domain. The dataset was constructed from over one million fashion images with a label space that includes 8 groups of 228 fine-grained attributes in total. The result is the first known million-scale multi-label and fine-grained image dataset. We conduct experiments and provide baseline with various CNN models. Importantly, we demonstrate models pre-trained on iFashion can achieve better transfer learning performance on fashion-related tasks than ImageNet or other fashion datasets. Data is available at: [https://github.com/visipedia/imat\\_fashion\\_comp](https://github.com/visipedia/imat_fashion_comp).

## 1. Introduction

Recent deep learning models trained on large-scale datasets have significantly advanced various computer vision tasks, and the performance on existing image classification benchmarks such as ImageNet [2] has reached the saturation point [5, 15, 6]. New datasets need to be created to tackle more challenging problems, such as multi-label classification and fine-grained recognition. On the other hand, domain-specific datasets have raised a lot of interest, especially in fashion domain [18, 12, 4]. In light of this, we introduce an iMaterialist Fashion Attribute Dataset (iFashion), which includes over one million annotated fashion images where the labels are curated by fashion experts. The label space includes 8 groups and a total of 228 fashion attributes, as described in Table 1.

iFashion presents a few unique challenges. Firstly, it is a multi-label prediction problem and the models are evaluated by precision and recall. Most existing datasets created for multi-label image recognition are limited in scale, such as PASCAL VOC [3], COCO [11] and NUS-WIDE [1], which

have about 6K, 80K and 160K training images from 20, 80 and 81 categories, respectively. Both learning difficulty and annotation effort would be increased considerably when the number of categories increases.

Secondly, many fashion attributes in iFashion are fine-grained labels and have very similar visual patterns. For example, as shown in Fig. 1, in the group of *Neckline*, identifying fine-grained visual difference on the defined fashion pattern (*Neckline*) between classes of *U-necks* and *Shoulder* is particularly challenging because the images often have large visual diversity within each class. This is much more significant than the subtle distinctions between different classes on the defined fashion pattern, resulting in significantly larger intra-class diversity than inter-class variance. This gives rise to new challenges compared to existing benchmarks for fine-grained recognition where images often have similar visual appearance with low intra-class diversity, such as CUB-200-2011 database [16] and Stanford Cars database [9]. More details on related studies with full comparisons between existing datasets are presented in the supplementary material.

The goal of iFashion is to encourage research on a more complex task toward real-world applications, by jointly considering multi-label and fine-grained image recognition with a hierarchical label structure. Our major contributions are: (i) the first known million-scale image dataset with multiple fine-grained attribute labels curated by experts; (ii) extensive experiments were conducted by using recent CNN models for multi-label and fine-grained recognition tasks, providing meaningful baseline results; (iii) we demonstrate empirically that iFashion is valuable for transfer learning on other fashion related datasets and applications.

## 2. iFashion Dataset

We describe the details of iFashion database. All images in iFashion are provided by Wish. We collected 1M+ fashion images by randomly sampling across individual attribute classes. All the images were pre-tagged by humans using an organically grown taxonomy. Then these tags were mapped to our taxonomy. Please refer to the supplementary material for post-processing steps we applied to improve

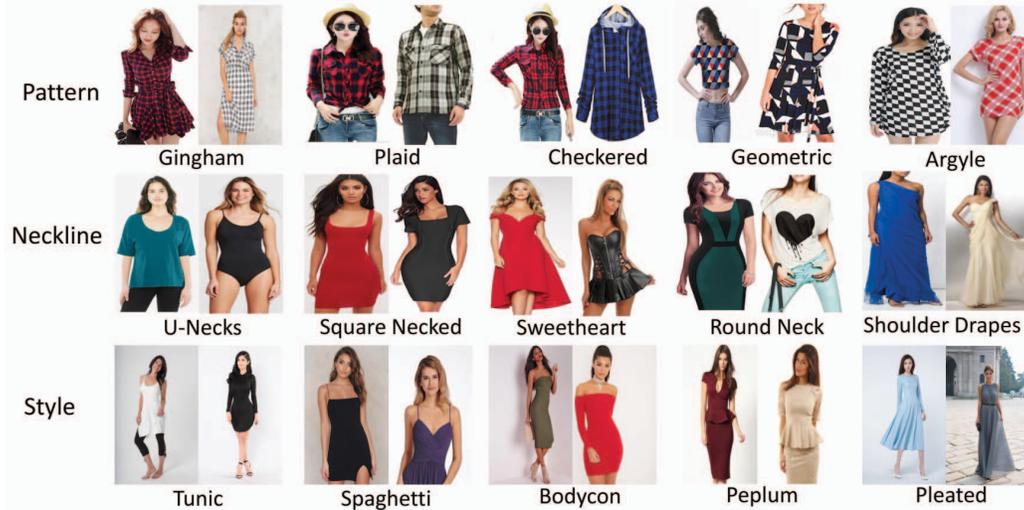


Figure 1. Examples from iFashion dataset for the attribute groups of *Pattern*, *Neckline*, *Style*.

Attribute	# Class	Type	# Label	# Image	Example
Category	105	S	913,857	913,857 - 90.2%	Athletic Pants, Bikinis, Cargo Pants, Heels, Petticoats ...
Color	21	M	894,904	467,137 - 46.1%	Black, Bronze, Gold, Gray, Green ...
Gender	3	M	1,012,947	935,265 - 92.3%	Male, Female, Neutral.
Material	34	M	701,197	591,175 - 58.4%	Nylon, Organze, Patent, Plush, Rayon ...
Neckline	11	S	721,908	721,908 - 71.3%	Racerback, Shoulder Drapes, Square Necked, Turtlenecks, U-Necks ...
Pattern	28	M	325,361	311,676 - 30.8%	Argyle, Camouflage, Checkered, Floral, Galaxy ...
Sleeve	5	S	733,501	733,501 - 72.4%	Long Sleeved, Puff Sleeves, Short Sleeves, Sleeveless, Strapless.
Style	21	S	610,442	610,443 - 60.3%	Asymmetric, Summer, Tunic, Vintage Retro, Wrap ...

Table 1. The number of classes, type (single-label or multi-label), number of labels and images for each attribute group in iFashion.

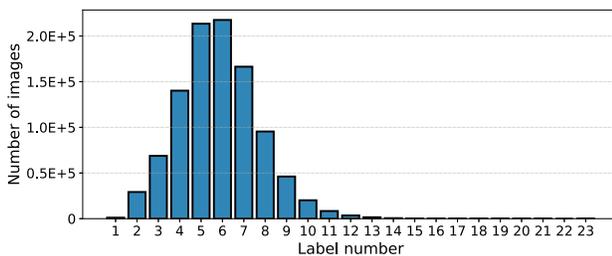


Figure 2. Histogram of number of labels per image, with an average of 5.8 and 8 per image in the training and validation sets.

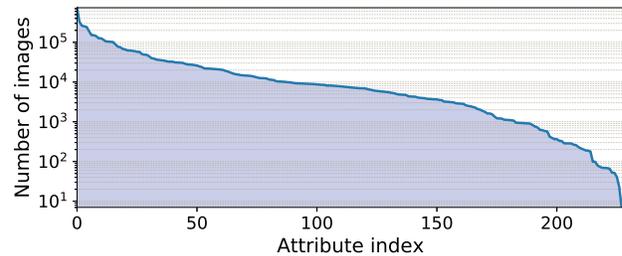


Figure 3. Number of images per attribute label, demonstrating the long tail nature of the dataset.

dataset quality. This results in iFashion database having 228 fine-grained attribute-level classes form 8 high-level groups defined professionally from the fashion industry. It contains 1,012,947 images for training, 9,897 and 39,706 manually-cleaned images for validation and testing.

**Dataset statistics.** The numbers of images and labels provided for each group are listed in Table 1. As can be found, the “Gender” group has a label in 92.3% images of the training set, while the “Pattern” group just has labels in 30.8% images. Histogram for the number of labels per image is shown in Fig. 2, where the number of labels per image is ranged from 1 to 23, with an average

of 5.8. Furthermore, the number of images per attribute-level class is shown in Fig. 3, where 31 classes have <500 training images, while 88 classes are with >10K images, indicating significant data imbalance. In addition, recognition difficulty is changed significantly over different high-level groups or attribute-level classes. Fig. 4 shows top-8 recalls for attribute-level classes, and the average top-1 recalls for 8 groups correspond to Table 1 are: 58.5%, 48.3%, 97.3%, 52.2%, 66.0%, 43.1%, 86.2%, and 28.8%.

Our database considers large scale (million level), multiple labels (with group structure), and fine-grained recognition jointly for fashion recognition, setting it apart from



Method	Top1	Top3	Top5
WTBI [4]	–	43.7	66.3
DARN [7]	–	59.5	79.6
Yang <i>et al.</i> [19]	–	75.3	84.9
FashionNet [12]	–	82.6	90.2
Inception-BN (DeepFashion)	64.6	85.4	91.6
Inception-BN (Clothes-1M)	65.6	85.9	91.9
Inception-BN (ImageNet)	67.6	87.3	92.9
Inception-BN (iFashion)	69.2	<b>88.2</b>	<b>93.3</b>

Table 4. Transfer learning on DeepFashion.

CNNs by using the Clothes-1M and Clothes-50K sequentially, by following previous approaches implemented on the Clothes-1M and Clothes-50K. Results on the validation set of Clothes-50K are reported in Table 3.

As shown in Table 3, by using Clothes-50K as training data, the pre-trained model from iFashion obtains the best performance with 78.9% average accuracy. It outperforms the other three pre-trained models by large margins, particularly ImageNet with 74.9%. This suggests that with a similar data scale, our database has stronger generalization capability to fashion-related tasks than the object-centralized ImageNet. Compared with fashion-related DeepFashion or Clothes-1M, iFashion is larger in scale and has higher label quality, resulting in better generalization performance. More detailed comparisons and discussions are presented in the supplementary material.

In the second group of experiments, iFashion pre-trained models consistently outperform ImageNet and DeepFashion, but the impact of pre-trained models is decreased when the amount of training data is increased from 50K to 1M+. Furthermore, our result of 80.5% is better than those of recent approaches specifically designed to handle noisy data in Clothes-1M, while our model, empowered by iFashion, just employs a simple and straightforward fine-tuning method, with an off-the-shelf CNN.

**Transfer learning on DeepFashion.** DeepFashion [12] has 46 classes with 209,222 training images and 40,000 validation images, which were manually cleaned and annotated. We investigate the transfer capability of three *million-level* databases: ImageNet, Clothes-1M and iFashion. We train Inception-BN models individually on each of the three databases, and then fine-tune them on DeepFashion. Results are compared in Table 4.

The results are consistent with those on Clothes-50K: (i) all pre-trained models improved the performance over that of training from scratch; (ii) iFashion obtains the best performance on all terms, demonstrating its stronger capability for transfer learning; (iii) with iFashion pre-training, we can achieve state-of-the-art results on DeepFashion, by simply using an off-the-shelf Inception-BN. Interestingly, ImageNet pre-training has better performance than Clothes-1M, which may due to high-quality data with a number of overlapped fashion categories between ImageNet and DeepFashion, as analyzed in the supplementary material.

## 4. Conclusion

We present the iFashion dataset, which is the first known million-scale expertly curated image dataset with multi-label and fine-grained attributes. The aforementioned characteristics of the iFashion enable it to be relevant for real-world applications, particularly in fashion domain. The introduction of iFashion allows us to compare different approaches for multi-label learning, which we provide several baselines. Our experiments show that there is still large room to improve in this space. We also demonstrated the value of iFashion for transfer learning, where it outperforms the other well-known datasets on fashion recognition.

## References

- [1] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y.-T. Zheng. Nus-wide: A real-world web image database from national university of singapore. In *CIVR*, 2009.
- [2] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [3] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *IJCV*, 2015.
- [4] M. Hadi Kiapour, X. Han, S. Lazebnik, A. C. Berg, and T. L. Berg. Where to buy it: Matching street clothing photos in online shops. In *ICCV*, 2015.
- [5] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition, 2016. *CVPR*.
- [6] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *CVPR*, 2018.
- [7] J. Huang, R. S. Feris, Q. Chen, and S. Yan. Cross-domain image retrieval with a dual attribute-aware ranking network. In *ICCV*, 2015.
- [8] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [9] J. Krause, M. Stark, J. Deng, and L. Fei-Fei. 3d object representations for fine-grained categorization. In *3DRR*, Sydney, Australia, 2013.
- [10] K.-H. Lee, X. He, L. Zhang, and L. Yang. Cleannet: Transfer learning for scalable image classifier training with label noise. *CoRR*, abs/1711.07131, 2017.
- [11] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [12] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *CVPR*, 2016.
- [13] G. Patrini, A. Rozza, A. K. Menon, R. Nock, and L. Qu. Making deep neural networks robust to label noise: A loss correction approach. In *CVPR*, 2017.
- [14] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- [15] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Re-thinking the inception architecture for computer vision. In *CVPR*, 2016.
- [16] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [17] X.-Z. Wu and Z.-H. Zhou. A unified view of multi-label performance measures. *CoRR*, abs/1609.00288, 2016.
- [18] T. Xiao, T. Xia, Y. Yang, C. Huang, and X. Wang. Learning from massive noisy labeled data for image classification. In *CVPR*, 2015.
- [19] Y. Yang and D. Ramanan. Articulated pose estimation with flexible mixtures-of-parts. In *CVPR*, 2011.