

An Effective Pipeline for a Real-world Clothes Retrieval System

Yang-Ho Ji^{1*}, HeeJae Jun^{2*}, Insik Kim^{3*}, Jongtack Kim^{4*}, Youngjoon Kim^{5*},
Byungsoo Ko^{6*}, Hyong-Keun Kook^{7*}, Jingeun Lee^{8*}, Sangwon Lee^{9*}, Sanghyuk Park^{10*}
NAVER/LINE Vision[†], Search Solution Vision[‡]

{arnilone¹, kobiso62⁶, hyongkuen63⁷, jglee0206⁸, shine0624¹⁰}@gmail.com
{heejae.jun², insik.kim³, jongtack.kim⁴, kim.youngjoon⁵, leee.sangwon⁹}@navercorp.com

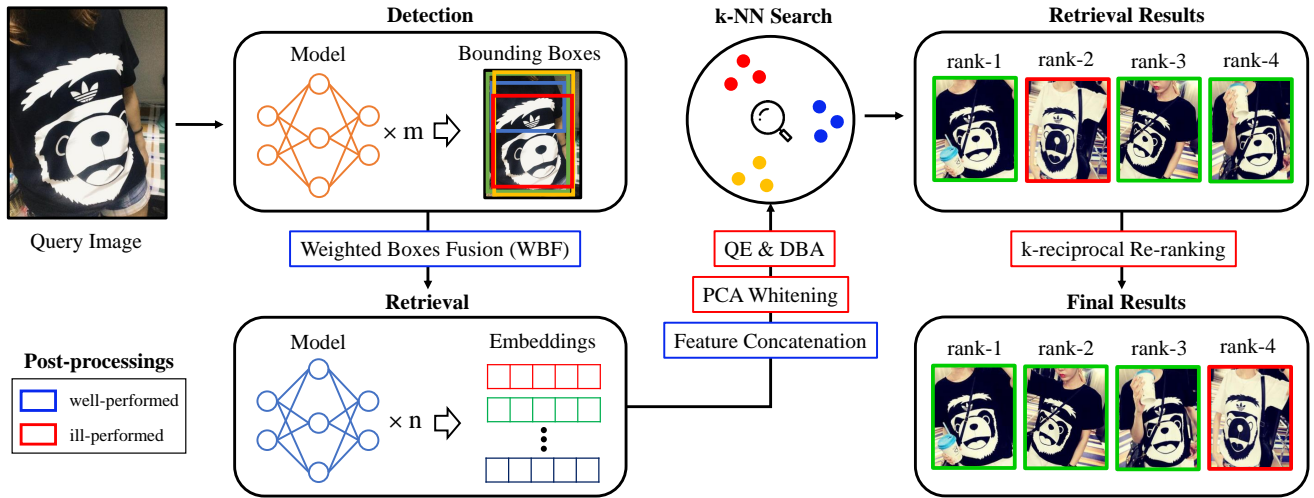


Figure 1: The proposed pipeline for real-world clothes retrieval that consists of three components: detection, retrieval and post-processing. Blue boxes indicate the post-processing methods that performed well while red is for those that did not.

Abstract

In this paper, we propose an effective pipeline for clothes retrieval system which has sturdiness on large-scale real-world fashion data. Our proposed method consists of three components: detection, retrieval, and post-processing. We firstly conduct a detection task for precise retrieval on target clothes, then retrieve the corresponding items with the metric learning-based model. To improve the retrieval robustness against noise and misleading bounding boxes, we apply post-processing methods such as weighted boxes fusion and feature concatenation. With the proposed methodology, we achieved 2nd place in the DeepFashion2 Clothes Retrieval 2020 challenge.

*All authors from Visual Search & Recommendation team at NAVER/LINE corporation contributed equally to this work and are listed in alphabetical order by last name.

1. Introduction

Recently, fashion domain has attracted much attention in computer vision research. Similarly, demand for online shopping for ‘fashionable’ items has also emerged to be one of the world’s largest industries. Thus, it has become crucial for advanced visual clothing retrieval systems to provide the finest shopping experience for online customers.

In this paper, we present an effective pipeline for a real-world clothes retrieval system using large-scale fashion data. The clothes retrieval system includes two vision tasks: detection and retrieval. Since the advent of deep convolutional neural networks (DCNN), vision tasks have improved to show astonishing results thus we employ modern variants of these DCNN models in our detection and retrieval network. In addition, we take advantage of popular post-processing methods to further improve the performance of the retrieval system: weighted boxes fusion (WBF) methods to effectively combine the predicted boxes of multiple detectors, and feature concatenation to elevate the similarities

Dataset	# of pairs	# of images
DeepFashion2 [4]	873,234	491,895
DARN [9]	91,390	182,780
Street2Shop [13]	39,479	416,840
Zalando [5]	16,253	32,642
MVC [12]	36,656	156,449
MPV [3]	13,524	50,290

Table 1: The utilized fashion-related datasets.

Model	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
ATSS	72.8	83.6	80.6	52.7	53.2	73.0
Cascade Mask R-CNN	72.7	83.9	80.4	42.6	52.0	68.4
CenterNet	70.4	81.4	77.9	30.1	44.3	70.8
RetinaNet-R101-FPN	67.4	81.1	76.2	32.7	46.1	67.6
RetinaNet-X101-FPN	67.4	81.2	76.1	32.7	51.9	67.6

Table 2: Detection results on DeepFashion2 *validation set*.

of the feature embeddings from different DCNN models.

2. Methodology

In this section we present our method of clothes retrieval which consists of three components: detection, retrieval, and post-processing. Figure 1 depicts the whole pipeline.

2.1. Detection

To retrieve fashion items correctly on an image, it is necessary to detect each fashion item via a detection network. We leverage five state-of-the-art object detection models to get well-localized and classified clothing boxes: adaptive training sample selection (ATSS) [24], cascaded mask R-CNN [1], CenterNet [26], RetinaNet-R101-FPN, and RetinaNet-X101-FPN [11]. We report detection performances in detail on Section 3.2.

2.2. Retrieval

To perform retrieval tasks based on the detected clothing items, we evaluate prevalent DCNN models (e.g. ResNet152 [6] and SE-ResNeXt101 [8]), optimize loss functions and fine-tune hyperparameters to decide on the well performing options. In addition, we leveraged a unique framework called combination of multiple global descriptors (CGD) [10] which has been proven to increase performance of a model by combining multiple global descriptors, such as maximum activation of convolutions (MAC) [19] and generalized mean pooling (GeM) [15]. We report the performances of the utilized models in Section 3.3.

2.3. Post-processing

To boost the performance of the clothes retrieval task, we exploit several post-processing methods. We create ensembles using detection and retrieval results as well as experiment with other widely-used techniques, such as feature manipulation and re-ranking. We report the effectiveness of these post-processing methods on Section 3.4.

Model	Description	Acc@10
m0	used DeepFashion2 train set only to train	0.868
m1	+Street2Shop	0.873
m2	+DARN+MVC	0.873
m3	+Street2Shop+DARN	0.878
m4	+Street2Shop+DARN+MVC	0.879
m5	+Street2Shop+DARN+MVC+MPV	0.877
m6	+Street2Shop+DARN+MVC+MPV+Zalando	0.878
m7	used DeepFashion2 train and valid set to train	-
m8	used hparams by automl for msloss on m0	0.871
m9	changed backbone to SENext101 on m8	0.858
m10	increased resol. to 224x384 on m0	0.870
m11	used GeM+MAC descriptors on m4	0.870
m12	changed lr from 0.2 to 0.17 on m11	0.877
m13	changed lr from 0.2 to 0.17 on m4	0.878
m14	used hparams by automl for msloss on m13	0.881
m15	changed lr from 0.2 to 0.17 on m6	0.882
m16	used hparams by automl for msloss on m15	0.880
m17	used GeM+SPoC descriptors on m4	0.870

Table 3: Top-10 accuracy of retrieval models with ground truth bounding boxes on DeepFashion2 *validation set*.

Weighted Boxes Fusion In order to enhance the retrieval performance even further, we adopt the WBF [18] technique of creating well-cropped query and database images from the results of multiple detectors. WBF collates all box information from our detection network and provides more accurate predictions which we found to correct inconsistencies in results.

Feature Concatenation As each model captures different features of a given image, we ultimately concatenate l_2 -normalized feature embeddings to increase our accuracy when retrieving images. This method has a trade-off with computation time and memory being used; however, each addition of features helped to achieve better performance in the task of image retrieval.

3. Experiments

In this section, we describe the datasets that were used and report experimental result details. Included are the performances of each component of our retrieval pipeline and discussions about the effectiveness of post-processing methods that were assessed.

3.1. Dataset

Table 1 shows various fashion-related datasets that were used in detection and retrieval tasks. Specifically, original images were used to train detection models while cropped and deduplicated images were for the retrieval models.

3.2. Detection Results

We train our detection models on DeepFashion2 *train set* only, then validate their performances on DeepFashion2 *validation set*. On Table 2, we summarize the evaluated detection performances on validation set.

Param	Value	Param	Value
Backbone	ResNet152-D [7]	Descriptor	GeM [14]
Input size	224×224	Output dim.	1,024
Total epochs	200	Optimizer	SGD
Init. LR	0.2	LR decay	cosine
Losses	Softmax + Center [22] + Triplet [16] + MS [21]		

Table 4: Default settings for training retrieval models.

Models	WBF_1	WBF_2	WBF_3	WBF_4	WBF_5
ATSS	✓	✓	✓	✓	✓
+ Cascade Mask R-CNN		✓	✓	✓	✓
+ CenterNet			✓	✓	✓
+ RetinaNet-R101-FPN				✓	✓
+ RetinaNet-X101-FPN					✓
AP_{50} (%)	83.6	84.0	84.9	85.4	85.8

Table 5: Weighted boxes fusion (WBF) results on DeepFashion2 *validation set*. WBF_n refers WBF result with n models. There are performance gains as more detection models are added.

3.3. Retrieval Results

Table 3 shows the performances of retrieval models using ground-truth bounding boxes on the DeepFashion2 validation set. In addition to the DeepFashion2 *train set*, we accumulated publicly available datasets as listed in Table 1 to train our retrieval models. The default configurations used to train retrieval models is described in the Table 4.

3.4. Post-processing Results

We summarize the aforementioned post-processing methods with a description of how it was applied and whether it worked to improve image retrieval result.

3.4.1 What worked well

WBF Using ATSS bounding boxes showed the best score on the Acc@10 so it is chosen to be the base when fusing other results from detection models. As shown in Table 5, the detection performance increments as more boxes are merged together.

Feature Concatenation We summarize performances of ensemble retrieval models using WBF_5 bounding boxes on the DeepFashion2 validation/test set in Table 6. By involving more models, we gain 0.019/0.014 increases on the Acc@10 results as compared to that of the single baseline retrieval model on DeepFashion2 validation/test set.

3.4.2 What did not work well

PCA Whitening In our experiments, dimensionality reduction by principal component analysis (PCA) [23] always decreased performance by a small margin. Using the findings of our prior experiment, we applied whitening [17] without

Ensemble Model	Validation (Acc@10)	Test (Acc@10)
m4 (baseline)	0.869435	0.840060
m0-13 (w/o m6-9)	0.875656	0.849722
m0-13 (w/o m6-7)	0.876141	0.853057
m0-16 (w/o m7)	0.876868	0.851860
m1-16 (w/o m7)	0.877192	0.851945
m0-17 (w/o m7)	0.876222	0.851945
m1-16	0.888341	0.854168
m1-17	0.888180	0.853356

Table 6: Top-10 accuracy of retrieval ensemble models with WBF_5 bounding boxes on DeepFashion2 *validation/test set*.

Model	Detection	Search Method	Acc@1	Acc@10
m4	WBF_5 (590k)	NN	0.668740	0.869435
m4	WBF_5 (126k)	NN	0.650885	0.853761
m4	WBF_5 (126k)	NN + re-ranking	0.661954	0.858124

Table 7: The experimental results using k-reciprocal re-ranking on DeepFashion2 *validation set*.

Year	Rank	Team	Acc@10
2020	1	Alibaba	0.872082
	2	NAVER/LINE Vision (Ours)	0.854168
	3	DeepBlueAI	0.848012
2019	1	Hydra@ViSenze	0.840658
	2	MM AI kakao	0.823258
	3	DeepBlueAI	0.816460

Table 8: The final results of DeepFashion2 clothes retrieval challenge in 2020 and 2019.

dimensionality reduction, but that still did not add any improvements to the overall performance. Thus, we decided to use the full dimensions of the concatenated features to maximize the score.

QE and DBA Query expansion (QE) [2] and database-side augmentation (DBA) [20] replaces each feature point with a weighted sum of its top k nearest neighbors and the point itself. The purpose of these techniques is to obtain rich and distinctive image representations by exploiting the nearest neighbors. However, QE decreased in overall performance while DBA gave a negligible increase. We speculate that this is because there were too many similar boxes proposed for the high recall query, hence the weighted sum of those boxes were not an informative image representation.

Re-ranking In Table 7, we present the experiment of k-reciprocal encoding [25], a well-known re-ranking technique in the Re-ID task, based on the DeepFashion2 validation set. This method requires much less bounding boxes to be present thus we used WBF to reduce the bounding boxes from 590k (WBF_5) to 126k (WBF_5). Interestingly, the k-reciprocal encoding helped increase the acc@10 by 0.02 but overall, the performance was still less than simply having reduced the boxes to 440k (WBF_5) using the WBF.

4. Conclusion

In this paper, we presented an effective pipeline for a real-world clothes retrieval system, including detection and retrieval task. And also, we investigated various post-processing methods such as weighted boxes fusion, feature concatenation, and other techniques. Table 8 shows the final result of the DeepFashion2 clothes retrieval challenge in 2020 and 2019. With the proposed pipeline, we achieved 0.854168 on Acc@10 in the test phase and ranked on 2nd place in the DeepFashion2 clothes retrieval challenge 2020. Furthermore, we believe the considered pipeline and methods described in this paper can be generalized to fulfill visual search for images not only about fashion but also furniture, beauty products and toys.

References

- [1] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [2] Ondrej Chum, James Philbin, Josef Sivic, Michael Isard, and Andrew Zisserman. Total recall: Automatic query expansion with a generative feature model for object retrieval. In *2007 IEEE 11th International Conference on Computer Vision*, pages 1–8. IEEE, 2007. 3
- [3] Haoye Dong, Xiaodan Liang, Xiaohui Shen, Bochao Wang, Hanjiang Lai, Jia Zhu, Zhiting Hu, and Jian Yin. Towards multi-pose guided virtual try-on network. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 2
- [4] Yuying Ge, Ruimao Zhang, Lingyun Wu, Xiaogang Wang, Xiaoou Tang, and Ping Luo. A versatile benchmark for detection, pose estimation, segmentation and re-identification of clothing images. *CVPR*, 2019. 2
- [5] Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S Davis. Viton: An image-based virtual try-on network. In *CVPR*, 2018. 2
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 2
- [7] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 558–567, 2019. 3
- [8] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Enhua Wu. Squeeze-and-excitation networks, 2017. 2
- [9] Junshi Huang, Rogerio S. Feris, Qiang Chen, and Shuicheng Yan. Cross-domain image retrieval with a dual attribute-aware ranking network. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015. 2
- [10] HeeJae Jun, ByungSoo Ko, Youngjoon Kim, Insik Kim, and Jongtaek Kim. Combination of multiple global descriptors for image retrieval, 2019. 2
- [11] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 2
- [12] Kuan-Hsien Liu, Ting-Yen Chen, and Chu-Song Chen. Mvc: A dataset for view-invariant clothing retrieval and attribute prediction. In *ICMR '16*, 2016. 2
- [13] Svetlana Lazebnik Alexander C. Berg Tamara L. Berg M. Hadi Kiapour, Xufeng Han. Where to buy it: matching street clothing photos in online shops. In *International Conference on Computer Vision*, 2015. 2
- [14] Filip Radenović, Giorgos Tolias, and Ondřej Chum. Fine-tuning cnn image retrieval with no human annotation. *IEEE transactions on pattern analysis and machine intelligence*, 41(7):1655–1668, 2018. 3
- [15] Filip Radenovi, Giorgos Tolias, and Ondej Chum. Fine-tuning cnn image retrieval with no human annotation, 2017. 2
- [16] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015. 3
- [17] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 806–813, 2014. 3
- [18] Roman Solovyev and Weimin Wang. Weighted boxes fusion: ensembling boxes for object detection models, 2019. 2
- [19] Giorgos Tolias, Ronan Sifre, and Herv Jgou. Particular object retrieval with integral max-pooling of cnn activations, 2015. 2
- [20] Panu Turcot and David G Lowe. Better matching with fewer features: The selection of useful features in large database recognition problems. In *2009 IEEE 12th International Conference on Computer Vision Workshops, ICCV Workshops*, pages 2109–2116. IEEE, 2009. 3
- [21] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott. Multi-similarity loss with general pair weighting for deep metric learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5022–5030, 2019. 3
- [22] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *European conference on computer vision*, pages 499–515. Springer, 2016. 3
- [23] Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52, 1987. 3
- [24] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z. Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. In *CVPR*, 2020. 2
- [25] Zhun Zhong, Liang Zheng, Donglin Cao, and Shaozi Li. Re-ranking person re-identification with k-reciprocal encoding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1318–1327, 2017. 3
- [26] Xingyi Zhou, Dequan Wang, and Philipp Krhenbhl. Objects as points, 2019. 2