

Fashion Retrieval using ResNet-50 and FAISS Indexing

Ayra, Arif Hakimi, Seung Hun Lee

2023320097, 2023320329, 2021320125

Abstract

001 Fashion image retrieval plays a central role in e-commerce
002 applications, allowing users to find visually similar cloth-
003 ing items based on example inputs. This project focuses
004 on developing a deep learning model that performs fashion
005 image retrieval by learning compact and discriminative em-
006 beddings. Our current model utilizes a ResNet50 backbone
007 trained with triplet loss to optimize for visual similarity in
008 an embedding space. The output embeddings are indexed
009 using FAISS to support efficient similarity search. Although
010 classification functionality is planned as an extension, the
011 current version is solely retrieval-focused. This work serves
012 as the foundation for a future dual-purpose fashion recog-
013 nition system capable of both retrieval and category predic-
014 tion.

015 1. Introduction

016 Image-based retrieval is a cornerstone of modern fashion-
017 AI applications: given a query photo, the system must sur-
018 face visually similar garments from a candidate gallery. Un-
019 like pure classification—which assigns a single category la-
020 bel—retrieval demands an embedding space whose geom-
021 etry captures subtle cues such as silhouette, texture, and
022 color.

023 Our baseline tackles this requirement by training a deep
024 convolutional network (ResNet-50) with triplet loss, which
025 explicitly pulls visually-similar items together while push-
026 ing dissimilar ones apart. The learned embeddings are
027 indexed by FAISS to enable sub-second top- K nearest-
028 neighbor search.

029 Because project time allowed, we also added a
030 lightweight classification head on the same backbone. Al-
031 though its current top-1 accuracy is only $\sim 57.5\%$, the
032 predicted category lets us filter the FAISS search to a
033 semantically-coherent subset, yielding cleaner retrieval re-
034 sults in many cases. Future iterations will focus on boosting
035 the classifier and jointly fine-tuning both objectives.

2. Problem Identification

036 Convolutional neural networks like ResNet50 are com-
037 monly used for image classification, producing softmax
038 probabilities over fixed labels. However, such models are
039 not optimized for measuring visual similarity, which is es-
040 sential in fashion image retrieval. Classification outputs do
041 not reflect how similar two items look, especially when they
042 belong to the same category but differ in style or appear-
043 ance.
044

3. Technical Soundness

3.1. Baseline Model and its Limitations

045 As a baseline, we consider a standard image classification
046 model based on ResNet50 trained with cross-entropy loss.
047 This model is designed to assign a single clothing category
048 (e.g., "t-shirt", "dress", "pants") to each input image. It
049 represents a simpler approach compared to our retrieval-
050 focused model, as it focuses solely on category prediction
051 without learning an embedding space or supporting image
052 similarity search.
053
054

055 This baseline fulfills the role of a minimal benchmark:

- 056 • It is **simpler** than our proposed model, with a single clas-
057 sification objective.
- 058 • It performs **only one task well**, categorical prediction but
059 does not support retrieval or visual similarity matching.
- 060 • It uses **standard components**: a ResNet50 backbone and
061 a fully connected classification head trained with cross-
062 entropy loss, which is a widely accepted setup in image
063 classification.
- 064 Despite its effectiveness in predicting class labels, this
065 model has several limitations:
066 • It cannot retrieve visually similar clothing items, limiting
067 its utility in recommendation or search systems.
- 068 • The output is a discrete label, which discards nuanced vi-
069 sual differences between garments within the same cate-
070 gory.
- 071 • It does not support efficient similarity indexing or rank-
072 ing, as there is no learned embedding space.

073 These limitations highlight the need for a more versatile
074 model that can simultaneously classify and retrieve fashion

075 items. Our current project addresses this by focusing on
076 metric learning and embedding-based retrieval, with plans
077 to later integrate classification for dual-task capability.

078 4. Experimental Methodology

079 This section describes the full technical pipeline used to
080 develop and evaluate deep neural network-based image re-
081 trieval models for fashion products. Our experiments are
082 grounded in four trained models that vary in terms of their
083 loss functions, preprocessing methods, architectural com-
084 ponents, and auxiliary objectives. Each model is designed
085 to learn a 512-dimensional embedding suitable for visual
086 similarity search. We evaluate their performance on the
087 DeepFashion2 dataset using cosine similarity and a top-K
088 retrieval protocol.

089 4.1. Dataset and Preprocessing

090 We use the DeepFashion2 dataset, a widely used bench-
091 mark for image-based fashion retrieval and recognition. It
092 contains over 800,000 images with annotations including
093 bounding boxes, clothing categories, landmark points, seg-
094 mentation masks, and most critically, pair IDs that connect
095 images of the same garment (e.g., shop and consumer pho-
096 tos). These pair IDs are used to define positive matches for
097 training and evaluation.

098 To study the impact of background clutter and preprocess-
099 ing, we construct two versions of the dataset:

- 100 • **Uncropped:** In the TripletMarginLoss model, the full im-
101 age is used without any spatial filtering, which includes
102 background, human pose, and lighting noise.
- 103 • **Cropped:** In the TripletMarginLossWithCropping,
104 CrossEntropyLoss, and ContrastiveLossandReLU mod-
105 els, each image is cropped to the ground-truth bounding
106 box before being passed to the model, encouraging the
107 network to focus only on garment-specific visual features.

108 In all settings, the images are resized to 224×224 , nor-
109 malized using ImageNet channel statistics, and augmented
110 with horizontal flipping, color jittering, and random affine
111 transformations. These augmentations help improve gener-
112 alization by simulating real-world variations in lighting and
113 pose.

114 4.2. Model Architectures

115 Each retrieval model uses a ResNet-50 backbone pretrained
116 on ImageNet. The original classification head is discarded
117 and replaced by a new feature embedding head. This head
118 consists of a global average pooling layer followed by a
119 fully connected layer that maps the feature maps to a 512-
120 dimensional vector. L2 normalization is applied to the out-
121 put embeddings so that all vectors lie on the unit hyper-
122 sphere, making cosine similarity a valid distance metric.

TripletMarginLoss. This model is trained with triplet
margin loss and uses the full, uncropped images. It includes
no classification head or other auxiliary tasks. The embed-
dings are directly optimized to satisfy the relative distance
constraint among anchor, positive, and negative samples.
Hard negative mining is applied within each batch to select
the most informative negative examples for training.

TripletMarginLossWithCropping. This model applies
the same architecture and loss function as TripletMargin-
Loss but uses images cropped to the garment bounding box.
This modification reduces the variance introduced by irrel-
evant background information and helps the model focus
more on the shape, texture, and structure of garments. It
retains a pure metric learning objective with no auxiliary
task.

CrossEntropyLoss. This model extends TripletMargin-
LossWithCropping by incorporating a secondary prediction
branch. In addition to the 512-dimensional embedding out-
put, a parallel classification head is attached to the penulti-
mate ResNet feature layer. This head is trained to predict
one of the 13 DeepFashion2 clothing categories using soft-
max and cross-entropy loss. The total loss is a weighted
combination of triplet margin loss and classification loss:

$$\mathcal{L}_{total} = \mathcal{L}_{triplet} + \lambda \cdot \mathcal{L}_{classification}$$

where λ is empirically set to 0.5. This dual-task learning
scheme allows the model to learn category-discriminative
features while preserving fine-grained similarity relation-
ships between garments.

ContrastiveLossandReLU. This model differs from the
others in two key aspects: it uses contrastive loss instead
of triplet loss, and introduces a ReLU activation before the
final embedding layer. The contrastive loss treats training
samples as positive or negative pairs and directly penalizes
their pairwise distance:

$$\mathcal{L}_{contrastive} = y \cdot D^2 + (1 - y) \cdot \max(0, m - D)^2$$

where y is 1 for positive pairs, 0 otherwise, and m is the
margin (set to 1.0). The inclusion of ReLU restricts the
embeddings to the positive orthant, reducing their represen-
tational freedom, which we later show can affect perfor-
mance. This model receives cropped images and focuses
purely on pairwise similarity learning without category su-
pervision.

515 4.3. Training Setup

516 All models are trained for 10 epochs using the Adam op-
517 timizer with a learning rate of 1×10^{-4} . We use a batch
518 size of 1084 to maximize utilization of GPU memory and

stabilize gradient estimates across large batches. No learning rate scheduling or early stopping is applied; each model is trained for a fixed number of epochs to maintain consistency in evaluation. All training runs are performed using a single GPU. Model checkpoints are generated during training, although intermediate checkpointing frequency is not a focus of this study.

4.4. Retrieval Pipeline

After training, we extract 512-dimensional embeddings from all query and gallery images. These embeddings are L2-normalized and indexed using FAISS with inner product similarity, which is equivalent to cosine similarity due to normalization. For each query image, we retrieve the top-K most similar gallery embeddings.

Following the DeepFashion2 evaluation setup, a query is considered successfully retrieved if at least one of the top-K retrieved gallery items shares the same pair ID and has an intersection-over-union (IoU) greater than 0.5 with the ground-truth bounding box of the query image. Where applicable, we report retrieval accuracy at Top-1 and Top-5 levels. In models where quantitative metrics are not reported, qualitative inspection of the top-5 retrieved results is used to assess embedding performance.

5. Results and Analysis

This section presents a quantitative and technical analysis of the four trained models under evaluation. Retrieval performance is assessed using top-K accuracy metrics on the DeepFashion2 dataset, where a retrieval is considered correct if any of the top-K retrieved gallery images share the same pair ID or have a low distance value. Each model differs in either its loss function, preprocessing pipeline, or auxiliary objectives. Retrieval performance is evaluated using Top-K accuracy where explicitly available. For models without quantitative accuracy output, we analyze qualitative retrieval results and failure cases based on visual inspection of the top-5 nearest neighbors.

5.1. TripletMarginLoss

This baseline model applies triplet margin loss without cropping. As such, full images are passed into the embedding model, which allows non-garment elements (e.g., background, lighting, pose, shadows) to influence training. Retrieval results show that the model frequently retrieves images with similar background settings rather than semantically similar garments. For example, in one case, a dress was recovered along with a shirt with similar lighting in the scene and wall color, but dissimilar structure. This indicates that the learned embeddings are heavily influenced by global image features and do not reliably capture garment-specific identity. As shown in the diagram below is one of the results of the TripletMarginLoss model.



Figure 1. TripletMarginLoss's result

5.2. TripletMarginLossWithCropping

This model uses the same triplet loss but includes a cropping step based on ground-truth bounding boxes during both training and inference. Qualitatively, the retrievals improve significantly. The top-5 matches for a query garment generally show strong coherence in silhouette, texture, and category. For instance, dresses are retrieved alongside other dresses of similar length and neckline, rather than being influenced by background context. These results confirm the importance of spatial preprocessing in fashion metric learning. Although no top-K accuracy metrics are printed, the visual consistency of the retrieved items supports the hypothesis that the cropping of the bounding box improves the focus of the embedding. The figure below shows the output result from TripletMarginLossWithCropping model.



Figure 2. TripletMarginLossWithCropping's result

5.3. CrossEntropyLoss

This variant extends by introducing an auxiliary classification head trained with cross-entropy loss to predict garment category. During inference, the model outputs both category predictions and retrieval results. This dual-headed structure allows the system to assign class labels to new, previously unseen fashion items (e.g., online test images). Retrievals from this model display strong category alignment: garments retrieved share not only visual features but also correct class identity. The added semantic regularization appears to constrain the embedding space more effectively. For example, a test image of a blouse is consistently matched with other blouses in the gallery. Despite the lack of numerical accuracy output, the observed retrievals are the most consistent and semantically aligned of all four models. Below is the output of combining the prediction implementation with retrieval function model.



Figure 3. CrossEntropyLoss's result

5.4. ContrastiveLossandReLU

This model uses contrastive loss instead of triplet loss and introduces a ReLU activation before the final embedding layer. Unlike triplet loss, which enforces relative distance constraints, contrastive loss focuses on absolute distances between pairs. The use of ReLU restricts the output embeddings to the nonnegative orthant, potentially limiting the model's representational power. The following is one of the results obtained from testing the model.



Figure 4. ContrastiveLossandReLU's result

Qualitative results show that retrievals are visually reasonable but less consistent than those from TripletMarginLossWithCropping or CrossEntropyLoss. Some top-5 results contain garments that differ in category or structure from the query, despite texture similarity. This suggests that the contrastive loss may be less effective at capturing fine-grained rank-sensitive relations compared to triplet-based methods.

5.5. Qualitative Insights

Across all models, we analyzed representative top-5 retrieval results for both validation and external test images. The uncropped TripletMarginLoss baseline shows the most inconsistent retrievals, with background influence evident in failure cases. Models trained with bounding box cropping consistently focus on relevant garment features, such as sleeve type, color, and cut. The CrossEntropyLoss model additionally retrieves items that belong to the same semantic category as the query, supporting the value of auxiliary classification supervision. In contrast, the ContrastiveLossandReLU model retrieves texture-similar but semantically inconsistent garments, reflecting the limitations of contrastive loss and ReLU constraint for fine-grained fashion retrieval.

6. Conclusion

This work investigates deep metric learning techniques for fashion image retrieval using the DeepFashion2 dataset. We designed and evaluated four model variants based on different combinations of loss functions, spatial preprocessing, and auxiliary supervision. Our findings emphasize the critical importance of training setup and architectural decisions in shaping retrieval behavior.

First, we demonstrate that preprocessing, specifically, the cropping of garments using ground truth bounding boxes, significantly improves the focus of the retrieval by eliminating distracting background features. Second, we show that triplet margin loss consistently outperforms contrastive loss in producing rank-sensitive embeddings suitable for top-K nearest neighbor search. Third, incorporating an auxiliary classification objective further enhances semantic alignment within the embedding space, enabling the model to generalize better to unseen examples and improving category consistency in retrievals.

Finally, we observe that architectural modifications such as restricting the embedding space via ReLU activation can negatively impact expressiveness and reduce retrieval accuracy. Our qualitative analysis across all models provides insight into how these design decisions interact and affect learned representations.

Future work will explore multi-task training strategies that jointly optimize for classification, retrieval, and attribute prediction, as well as domain adaptation methods to improve robustness to out-of-distribution fashion content.

DeepFashion2 [1] provides richer annotations such as landmarks, masks, and pair IDs, enabling fine-grained visual retrieval. Triplet loss [2], contrastive loss [3], and ProxyNCA [4] are core metric learning methods for embedding-based retrieval.

References

- [1] Y. Ge, R. Zhang, X. Wang, X. Tang, and P. Luo, "Deep-fashion2: A versatile benchmark for detection, pose estimation, segmentation and re-identification of clothing images," *CVPR*, pp. 5337–5345, 2019.
- [2] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *CVPR*, pp. 815–823, 2015.
- [3] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *CVPR*, pp. 1735–1742, 2006.
- [4] Y. Movshovitz-Attias, A. Toshev, T. Leung, S. Ioffe, and S. Singh, "No fuss distance metric learning using proxies," in *ICCV*, pp. 360–368, 2017.