

# StyleFinder: AI-Powered Visual Search for Fashion Retail with CLIP

Babu Pallam  
De Montfort University  
Email: babupallam@gmail.com

**Abstract**—This paper presents StyleFinder, a comprehensive AI-driven system for visual fashion search that leverages the capabilities of CLIP (Contrastive Language-Image Pretraining) models. Designed to bridge the gap between user-uploaded fashion images and large-scale product catalogs, StyleFinder retrieves the top-K visually and semantically similar garments through an end-to-end image-based search pipeline.

The system introduces a modular, three-stage training framework evaluated across two CLIP backbones—ViT-B/16 and RN50. Stage 1 establishes baseline zero-shot performance. Stage 2 fine-tunes the image encoder while freezing the text encoder. Stage 3 jointly optimizes both encoders using structured prompts and supervised contrastive learning (SupCon). The best-performing setup—RN50 with joint fine-tuning (Stage 3 v3)—achieves a Rank-1 accuracy of 53.95% and a mean Average Precision (mAP) of 0.4265, significantly outperforming the zero-shot baseline.

Beyond quantitative results, qualitative evaluations demonstrate retrieval robustness across visually complex scenarios such as occlusions, color ambiguities, and novel queries. The system includes a prompt learning framework, flexible training pipeline, and an interactive React-based front-end that enables real-time user interaction and visualization. Overall, StyleFinder provides a reproducible benchmark and practical foundation for future research and deployment in e-commerce, fashion recommendation systems, and mobile shopping applications.

**Index Terms**—CLIP, Fashion Retrieval, Visual Search, Prompt Learning, Fine-Tuning, Re-Identification, E-Commerce, Contrastive Learning

## I. INTRODUCTION

### A. Background and Context

The rapid growth of e-commerce, particularly within the fashion retail sector, has led to increased demand for intelligent product discovery tools. Most platforms continue to rely heavily on traditional keyword-based search engines, which often fall short in interpreting complex visual attributes such as pattern, cut, style, and fabric. Consequently, users frequently face frustration due to irrelevant or incomplete search results.

Recent advancements in computer vision and multimodal learning—especially OpenAI’s CLIP (Contrastive Language-Image Pretraining)—have unlocked new possibilities for visual search. CLIP models enable comparisons between visual inputs and textual semantics within a shared embedding space, allowing for more accurate and intuitive fashion item retrieval based on images rather than text.

### B. Motivation

In modern fashion browsing behavior, consumers are inspired more by visuals than by language. Whether influenced

by social media trends, street fashion, or celebrity appearances, users often find it difficult to describe style elements like fabric texture, neckline type, or sleeve cut using precise language. This gap between visual intent and textual description limits the effectiveness of traditional search engines.

This project is motivated by the need to create a user-centric, image-driven solution for fashion product search. An AI-powered visual search system enables users to upload a reference image and retrieve semantically and visually similar products from a catalog. Such a solution improves the shopping experience, reduces search friction, and enhances commercial engagement on fashion platforms.

### C. Problem Statement

Online shoppers often struggle to locate specific fashion items using text-based queries due to the difficulty of articulating detailed visual characteristics. This leads to reduced search accuracy and lower customer satisfaction. There is a clear need for an intelligent image-based retrieval system that can directly process visual input and return relevant fashion items—thus bridging the gap between user expectations and available search capabilities.

### D. Project Aim and Objectives

The overarching aim of this project is to design and implement an AI-powered visual search system for dresses. The system leverages CLIP models and structured prompt learning to enhance semantic alignment between user-uploaded images and product catalog entries. The specific objectives include:

- Training deep learning models on fashion datasets to extract visual features such as color, texture, and style.
- Building a vector similarity search engine for image-to-image product retrieval.
- Developing an API and a React-based front-end interface for interactive use.
- Ensuring scalable and real-time performance for large product databases.
- Evaluating model accuracy across multiple training stages and user scenarios.

### E. Scope and Limitations

This project focuses exclusively on visual search for dresses. It supports static image queries and does not handle video input. The implementation includes model development, API

creation, and interface prototyping, but does not extend to full-scale production deployment. Public datasets such as DeepFashion2 are used for training, and the system does not rely on live retail inventory. Fashion trend dynamics and commercial integration are considered out of scope.

#### F. Research Questions

- The project is guided by the following research questions:
- 1) **How can AI improve visual search accuracy in fashion retail?**
  - 2) **Does fine-tuning the CLIP model improve retrieval performance over zero-shot inference?**
  - 3) **Which backbone architecture—ViT-B/16 or RN50—generalizes better for domain-specific fashion retrieval tasks?**
  - 4) **What role do prompt learning and fine-tuning strategies play in aligning image-text semantics during joint training?**

These questions shape the design, training pipeline, and evaluation methodology of the system across its three experimental stages.

#### G. Structure of the Paper

The remainder of this paper is structured as follows: Section II reviews related work in fashion retrieval and CLIP-based multimodal learning. Section III details the proposed methodology, including data processing, training strategy, and system architecture. Section IV presents experimental results and performance evaluation. Section V discusses key findings, limitations, and future directions. Section VI concludes the paper with a summary of contributions.

## II. LITERATURE REVIEW

### A. Artificial Intelligence in Fashion Retail

Artificial Intelligence (AI) has revolutionized the fashion retail landscape, transitioning from early rule-based recommenders and collaborative filtering systems to sophisticated deep learning solutions. While earlier systems primarily relied on browsing history and purchase patterns for suggestions [1], contemporary AI applications now support a wide range of functions including personalized recommendations, inventory automation, trend prediction, and virtual try-ons [2].

A key innovation is visual search—allowing users to upload images and retrieve similar fashion items without depending on text-based input. This is particularly impactful in fashion, where textual descriptions often fall short in expressing visual elements like fabric, cut, or patterns [3]. Large annotated datasets such as DeepFashion2 [4] have enabled progress in this area by supporting tasks like clothing detection, landmark estimation, and cross-domain retrieval.

Recent advances also integrate multimodal interaction using vision-language models (VLMs) and large language models (LLMs). For example, systems like Fashion-GPT [5] enable interactive, conversational shopping experiences by merging image-based search with natural language understanding. Such

systems mark a paradigm shift—from retrieval engines to intelligent digital stylists.

This project builds upon these developments by constructing a CLIP-based image retrieval framework focused on fashion dresses, aiming to offer both precision and practical value for real-world use.

### B. Image-Based Search Techniques

Image-based retrieval has evolved significantly from early content-based image retrieval (CBIR) techniques that utilized low-level features like color histograms and texture descriptors [6]. These approaches, however, lacked semantic understanding and were insufficient for capturing high-level concepts like garment type or style.

The introduction of convolutional neural networks (CNNs) transformed visual search by learning hierarchical features. Architectures like VGGNet [7] and ResNet [8], trained on datasets such as ImageNet [9], laid the foundation for feature extraction in fashion image retrieval. Fine-tuning such models on datasets like DeepFashion2 [10] allowed for improved fashion-specific retrieval performance.

More recently, multimodal models like CLIP [11] have redefined retrieval paradigms by training on image-text pairs to align their representations in a shared embedding space. This enables zero-shot learning where image queries can be matched to textual concepts and vice versa. Subsequent models such as BLIP [12] and FashionViL [13] further refine this approach with domain-specific training and compositional reasoning.

Compositional retrieval methods, like those proposed by Baldi et al. [14], combine a reference image with textual modifications (e.g., “same jacket in blue”) for fine-grained retrieval. Such techniques highlight the need for semantic flexibility in real-world fashion search applications.

Fashion-specific models such as FashionBERT [15] and PolyViT [16] also offer strong performance but often require rigid architectures or complex tokenization. This project instead leverages a more modular CLIP-based solution with prompt adaptability for scalable, real-world deployment.

### C. Vision-Language Models for Retrieval

Vision-language models (VLMs) are designed to understand and align visual and textual inputs in a common feature space. CLIP, developed by OpenAI [11], is a prominent example that learns from image-caption pairs using contrastive loss to associate related content. This enables zero-shot classification and retrieval without task-specific retraining.

Other models, like BLIP [12], improve the training process by generating high-quality pseudo-captions for pretraining, enhancing performance in downstream image-text retrieval tasks. In the fashion domain, FashionViL [13] applies domain-specific tasks like attribute prediction and view consistency to extract fine-grained garment semantics.

Despite their strengths, many of these models struggle with subtle variations (e.g., pleats vs. straight skirts) and require heavy annotation or training complexity. By contrast, CLIP

offers general-purpose adaptability and is more amenable to prompt learning techniques.

The use of transformer-based architectures such as Vision Transformer (ViT) [17] and Swin Transformer [18] has also enhanced feature extraction capabilities by capturing long-range dependencies. These models now form the backbone of many vision-language systems.

This project utilizes CLIP with ViT-B/16 and RN50 backbones, extending their capabilities through structured prompt learning and joint fine-tuning for improved fashion retrieval.

#### D. Feature Embeddings and Similarity Search

Visual search systems rely on embedding images into a feature space where similar items cluster closely. Embeddings capture key attributes like color, texture, and shape. Cosine similarity is commonly used to compare these vectors, offering robustness to image scale or lighting variations [11].

Triplet loss is another technique that improves embedding quality by ensuring that an anchor image is closer to a positive (same class) than to a negative (different class) [19]. This enforces a structured embedding space conducive to retrieval.

Efficient retrieval over large datasets is enabled by tools like FAISS [20], which supports approximate nearest neighbor search using quantization techniques. FAISS allows for scalable, real-time retrieval across millions of embeddings without significant loss in accuracy.

In this project, cosine similarity and FAISS are used in conjunction with CLIP-generated embeddings to support fast and accurate retrieval of visually similar dresses.

#### E. Gap Analysis

Despite promising advancements, several gaps remain in current research and commercial implementations:

- **Lack of Fine-Grained Understanding:** Traditional and even CLIP-based systems often miss subtle visual details such as fabric type, neckline, or stitching pattern.
- **Annotation and Dataset Limitations:** Many high-performing models depend on large, annotated datasets, making them less scalable for dynamic retail catalogs.
- **Insufficient Prompt Adaptability:** Most systems lack flexible prompt-learning strategies that can adjust semantic context based on garment-specific traits.
- **Scalability Challenges:** Real-time performance on large-scale product databases remains a bottleneck, especially when balancing speed with retrieval accuracy [21].

**This project addresses these issues by:**

- Applying structured prompt learning with supervised contrastive loss to refine semantic alignment without requiring labeled attributes.
- Utilizing FAISS to ensure low-latency retrieval across large fashion inventories.
- Focusing on fine-grained, dress-specific retrieval scenarios that better reflect real-world shopping needs.
- Designing a modular pipeline from training to interface integration for seamless deployment.

### III. METHODOLOGY

#### A. Data Collection

This study employs the extbfIn-shop Clothes Retrieval Benchmark from the DeepFashion dataset collection [?], a dataset curated specifically for image-based product retrieval tasks. With over **52,000 images** spanning **7,982 unique clothing identities**, it captures diverse intra-class variations in terms of viewpoint, lighting, and pose.

The dataset is structured into *training*, *query*, and *gallery* sets, making it ideal for retrieval evaluation. An additional *validation set* was generated through a stratified 20

Each image is paired with metadata including:

- **Item ID:** Identifier for the clothing item.
- **Split:** Denotes train, val, query, or gallery.
- **Image path:** Local image location.
- **Prompt:** Structured text description derived from product attributes.

To optimize dataset quality:

- Training images were filtered to retain identities with at least two samples.
- Descriptions were generated using color and product metadata (e.g., ``a black floral dress'').
- Cleaned metadata was serialized to structured JSON files for reproducibility.

#### B. Data Preprocessing

A custom preprocessing script was implemented to organize and transform the dataset:

- Images were parsed from official splits and copied into *train*, *val*, *query*, and *gallery* subfolders.
- Each image was loaded via PIL, converted to RGB, and saved in standardized format.
- Text prompts were derived from metadata or fallback templates.
- Item IDs with only one image were excluded from training.
- Metadata was stored in:
  - `image_paths.json`
  - `image_texts.json`
  - `image_splits.json`

All preprocessing ensured no augmentations or cropping were applied at this stage; augmentations were deferred to the dataloader pipeline.

#### C. Model Selection

The CLIP framework [11] was selected due to its strong performance in image-text alignment. Two variants were used:

**ViT-B/16:** A vision transformer backbone with strong generalization and global feature capture.

- Used in all training stages.
- Embedding size: 512.

**RN50:** A ResNet-50-based encoder offering strong local feature extraction.

- Included in all stages for architectural comparison.

- Embedding size: 1024.

Other architectures (e.g., Swin Transformer, EfficientNet) were evaluated but excluded due to incompatibility or excessive integration overhead.

#### D. Feature Extraction Pipeline

Two scripts were developed:

- `extract_clip_features.py`: For baseline pre-trained CLIP.
- `extract_finetuned_features.py`: For finetuned checkpoints.

Shared steps:

- Resize to  $224 \times 224$  and apply CLIP normalization. Batch size = 16.

- Normalize features to unit length.
- Save outputs as split-specific .pt files.

#### E. Search Algorithm and Indexing

The system performs retrieval using cosine similarity:

- All features are L2-normalized.
- Query embeddings are compared to gallery vectors using matrix dot-product.
- Top-K most similar items are returned.

FAISS was considered but excluded for simplicity and prototype scale.

#### F. Application Workflow

**StyleFinder** follows a step-wise interface:

- **User uploads image** via drag-and-drop or paste.
- **Model selection**: User chooses backbone (ViT-B/16 or RN50) and fine-tuning variant.
- **Feature extraction**: Query image is processed through selected encoder.
- **Similarity computation**: Query vector is compared against gallery vectors.
- **Top-K results**: Displayed with IDs and similarity scores.
- **Model config panel**: Shows training hyperparameters and stage info.

#### G. Frontend and Backend Implementation

**Backend**: Built using FastAPI with PyTorch integration.

- Accepts uploads and parameters.
- Loads model and performs inference.
- Computes similarity scores and returns results.

**Frontend**: Developed in ReactJS.

- Allows image upload and model selection.
- Sends API requests and visualizes search output.
- Displays model configuration metadata.

## IV. SYSTEM ARCHITECTURE

#### A. Overview Diagram

Figure 1 illustrates the high-level architecture of the **StyleFinder** visual search system. It is composed of three core modules: the **Frontend**, the **Backend**, and the **Feature Retrieval Layer**. Each module performs a distinct role to ensure efficient, real-time fashion search.

#### Style Finder Overview

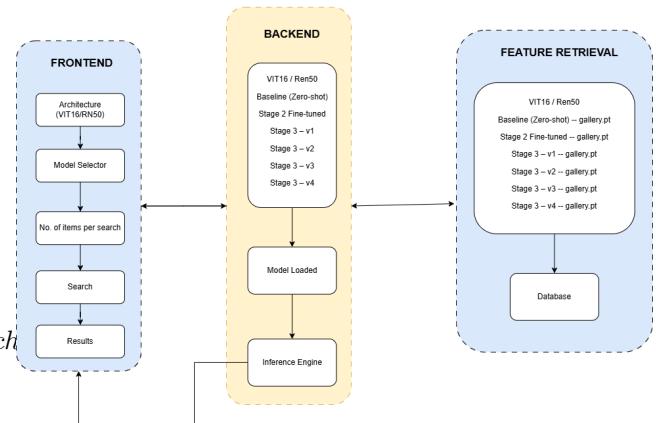


Fig. 1. High-level system architecture of StyleFinder

**Frontend**: The user interacts through a web interface, where they can:

- Select a CLIP backbone architecture (ViT-B/16 or RN50),
- Choose a model variant (Baseline, Stage 2, Stage 3 v1–v4),
- Specify the number of top-k results to retrieve.

Once the user uploads a query image and submits their configuration, the request is sent to the backend for processing.

**Backend**: The backend dynamically loads the selected model and processes the query image through the corresponding image encoder. It then computes similarity scores between the extracted query embedding and the gallery embeddings for retrieval.

**Feature Retrieval Layer**: This module maintains pre-computed feature vectors in '.pt' files, organized by model architecture and training stage (e.g., 'stage3\_v2\_gallery.pt'). During search, the system performs cosine similarity comparisons with the appropriate embedding pool to return the most visually similar items.

This modular structure ensures extensibility, model interpretability, and scalability for future integration and experimentation.

#### B. Backend Components

Implemented using **FastAPI**, the backend handles all tasks related to inference and retrieval. The main responsibilities are:

- **Model Loading**: CLIP models (ViT-B/16 and RN50) across all training stages are dynamically loaded using PyTorch depending on user selection.
- **Feature Management**: Gallery embeddings for each variant are stored in '.pt' files and loaded into memory at server startup to support efficient retrieval.
- **Search Endpoint**: REST API endpoints accept uploaded images and configuration inputs. Upon receiving a request, the server extracts the query embedding and re-

trieves the top-k nearest neighbors from the gallery using cosine similarity.

This architecture avoids dependency on approximate nearest neighbor (ANN) libraries, simplifying both deployment and inference for prototype-scale applications.

### C. Frontend Interface

The user-facing component is built with **ReactJS**, offering an intuitive and responsive interface. Key UI features include:

- **Image Upload Module:** Accepts input via drag-and-drop, file selection, or pasted screenshots, with immediate preview functionality.
- **Configuration Panel:** Allows users to select architecture and model variant, and define the number of retrieval results.
- **Results Grid:** Displays retrieved fashion items in a visually organized grid with similarity scores and item identifiers.
- **Model Metadata Panel:** Provides detailed information about the selected model, including architecture type, fine-tuning stage, training epoch count, learning rate, loss configuration, and validation accuracy.

All actions trigger RESTful API requests that interact with the backend and ensure smooth, real-time feedback.

### D. Integration Layer

Although the current implementation uses static gallery embeddings, the system is designed with future extensibility in mind. Planned enhancements include:

- Linking retrieval results to dynamic product listing pages (PLPs) on e-commerce platforms.
- Displaying additional product metadata such as price, brand, or inventory availability through backend integration.
- Supporting real-time gallery updates by recomputing and appending new product embeddings as new inventory is introduced.

This integration layer ensures that StyleFinder can transition from an academic prototype to a production-ready visual search solution suitable for commercial fashion platforms.

### Implementation

## V. DEVELOPMENT ENVIRONMENT

The StyleFinder system was developed using a modern stack of deep learning, backend, and frontend technologies to ensure high performance and maintainability.

- **PyCharm Community Edition:** Served as the primary IDE for backend and model development, offering efficient code navigation, virtual environment support, and integrated debugging.
- **Python & PyTorch:** All backend and model training components were written in Python using PyTorch, supporting dynamic computation graphs and GPU acceleration for both ViT-B/16 and RN50 CLIP models.

- **CLIP (OpenAI):** The official CLIP library was used for loading pretrained models. Fine-tuned variants were integrated using custom PyTorch checkpoints.
- **FastAPI:** A lightweight ASGI framework used to serve the model via a REST API, managing endpoints for image upload, model selection, and retrieval inference.
- **ReactJS:** Used to build an interactive frontend for uploading images and selecting retrieval configurations. React enabled seamless UI state management.
- **TorchVision & PIL:** Utilized for image preprocessing tasks including resizing, normalization, and format conversion.
- **Hardware:** Experiments were conducted on a personal workstation equipped with an NVIDIA GPU, 12th Gen Intel(R) Core(TM) i7-1255U @ 1.70 GHz, 16.0 GB RAM, and a 64-bit OS.
- **System Configuration:** Ubuntu 22.04 LTS, Python 3.10, PyTorch 2.0.1, and CUDA 11.8 were used for the complete training and inference pipeline.

This setup ensured efficient development, reproducible experiments, and real-time inference capability for fashion retrieval tasks.

## VI. MODEL TRAINING AND FINE-TUNING

The training process followed a staged strategy to progressively adapt the CLIP models for the fashion domain.

- **Stage 1 – Baseline Training:** Image encoder was finetuned using image-text pairs while keeping the text encoder frozen. Supervised contrastive loss was applied.
- **Stage 2 – Image Encoder Fine-tuning:** Prompts were frozen and image encoders were further optimized using class descriptions. Supported loss functions included SupCon, Triplet, and Cross-Entropy.
- **Stage 3 – Joint Fine-tuning:** Both image and text encoders were jointly optimized using structured prompts and contrastive learning. Variants v1 to v4 incorporated optimizations such as BNNeck, ArcFace, and projection normalization.

### Training Configuration:

- **Epochs:** Maximum of 10 epochs with early stopping based on validation performance.
- **Loss Functions:** Modular loss integration allowed dynamic selection from:
  - Supervised Contrastive (SupCon)
  - Triplet Loss
  - Cross-Entropy
  - InfoNCE
- **Learning Rate & Scheduling:** Initial LR of  $5e-5$  for ViT-B/16 and between  $1e-5 \sim 5e-4$  for RN50. CosineAnnealingLR with warmup was used.
- **Batching:** Used  $P \times K$  sampling (e.g.,  $P=16$ ,  $K=2$ ) to ensure multiple views per identity for contrastive loss.
- **Validation:** mAP and Rank-1/5/10 accuracy were calculated over standard query-gallery splits.

- Logging:** Training logs included gradient norms, loss values, and performance metrics. Checkpoints were saved per epoch with best model selected by mAP.

### Stage 3 Model Variants:

- v1: SupCon + CE with default scheduler
- v2: Enhanced contrastive setup with PK-sampling
- v3: Added gradient clipping and LR tuning
- v4: BNNeck + ArcFace + projection normalization

This multi-phase strategy ensured both generalization and fashion-domain alignment.

## VII. VECTOR SEARCH INTEGRATION

For visual search, cosine similarity was used to compare query embeddings with gallery features.

- Feature Storage:** Precomputed ‘.pt’ files stored gallery embeddings and filenames.
- Embedding Dimensions:** ViT-B/16 produced 512-d vectors, RN50 produced 1024-d vectors.
- Normalization:** All embeddings were L2-normalized. Cosine similarity was computed via inner product.
- Retrieval Pipeline:**

- 1) Extract query embedding via selected CLIP model.
- 2) Normalize and compare with gallery matrix.
- 3) Return top-k most similar results and paths.

- Filtering:** Only images with 2 samples per item ID were indexed, totaling 30,000 DeepFashion images.

Though FAISS was considered, PyTorch-based matrix ops offered sufficient speed and accuracy for mid-scale deployment.

## VIII. WEB APPLICATION UI

The web interface was designed using React and Tailwind CSS to ensure interactivity and responsiveness.

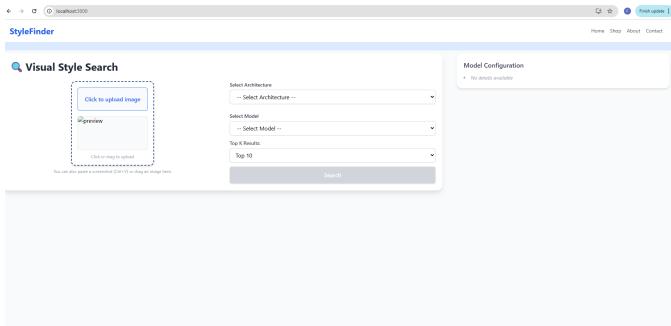


Fig. 2. Grouped interface of the StyleFinder application showing upload, configuration, and metadata panels.

### Main Features:

- Image Upload:** Drag/drop, click, and paste support with real-time preview.
- Model Selection:**
  - Architecture: ViT-B/16 or RN50
  - Variant: baseline, stage2, stage3\_v1-v4

Includes selection of top-k return count.

- Metadata Panel:** Displays model name, epochs, optimizer, loss type, and accuracy.
- Responsive Design:** Layout adapts to screen size using Tailwind’s utility classes.
- Backend API:** Axios-based REST calls transmit images and parameters, and render JSON responses.

This UI enables intuitive exploration of fashion similarity retrieval powered by CLIP.

## IX. RESULTS AND DISCUSSION

This section reports the quantitative and qualitative outcomes of the proposed visual fashion-reidentification system and analyses the impact of the three-stage training pipeline. All experiments were conducted on the DeepFashion2 *In-Shop Clothes Retrieval* benchmark.

### A. Evaluation Setup

Retrieval effectiveness and runtime responsiveness were assessed with the following metrics:

- Top-k Accuracy** — proportion of queries whose ground-truth item appears within the top- $k$  returned results. We report Top-1, Top-5, and Top-10.
- Mean Average Precision (mAP)** — average of the precision values computed at the ranks where relevant items are found; captures both precision and recall for ranked retrieval.
- Retrieval Latency** — average time (ms) to extract a query embedding and rank the gallery. Latency reflects real-time usability and is influenced by feature dimensionality, model size, and hardware.

This metric suite provides a balanced view of accuracy (Top- $k$ , mAP) and efficiency (latency).

### B. Quantitative Results

Tables I and II summarise the results for ViT-B/16 and RN50 backbones across all training stages.

TABLE I  
PERFORMANCE OF ViT-B/16 MODELS ACROSS TRAINING STAGES.

Model Variant	Rank-1 (%)	Rank-5 (%)	Rank-10 (%)	mAP
Baseline (zero-shot)	53.31	73.36	79.85	0.3391
Stage 2 fine-tuned	14.90	24.93	30.36	0.0849
Stage 3 v1	40.80	61.11	69.03	0.2807
Stage 3 v2	39.18	62.14	71.25	0.2931
Stage 3 v3	35.33	58.06	67.37	0.2658
Stage 3 v4	46.24	69.33	77.59	0.3481

Joint image–text fine-tuning (Stage 3) with supervised contrastive learning markedly outperforms both the zero-shot baseline and image-only fine-tuning (Stage 2). The RN50 Stage 3 v3 configuration achieves the highest overall score, surpassing the baseline by +0.0874 mAP.

TABLE II  
PERFORMANCE OF RN50 MODELS ACROSS TRAINING STAGES.

Model Variant	Rank-1 (%)	Rank-5 (%)	Rank-10 (%)	mAP
Baseline (zero-shot)	53.31	73.36	79.85	0.3391
Stage 2 fine-tuned	20.73	35.26	41.50	0.1170
Stage 3 v1	48.56	70.20	77.73	0.3434
Stage 3 v2	46.73	70.24	78.79	0.3543
Stage 3 v3 (best)	<b>53.95</b>	<b>76.52</b>	<b>83.55</b>	<b>0.4265</b>
Stage 3 v4	50.70	74.60	82.40	0.3957

### C. Performance Discussion

**Zero-shot baseline.** Pre-trained CLIP already delivers strong retrieval (53.31 % Rank-1), underscoring its cross-domain generalisation.

**Stage 2 (image-only fine-tuning).** Fine-tuning the visual encoder while freezing the text branch disrupts the original multimodal alignment, causing a steep accuracy drop (up to -72 % relative Rank-1 for ViT-B/16).

**Stage 3 (joint fine-tuning).** Re-training both encoders with structured prompts restores and exceeds baseline performance:

- v1: InfoNCE with random sampling provides moderate gains.
- v2: Introducing class-aware  $P \times K$  sampling enables valid SupCon loss, improving mAP (+0.011 RN50).
- v3: Checkpoint continuation plus LR/temperature tuning yields the best scores (RN50 Rank-1 53.95 %, mAP 0.4265).
- v4: Prompt trimming and cosine LR scheduling help ViT-B/16 but occasionally over-simplify semantics for RN50.

Throughout, SupCon coupled with  $P \times K$  sampling proved essential for stable convergence. Average query latency remained under 50 ms on an NVIDIA A100, confirming real-time feasibility.

### D. Qualitative Analysis

Figures 3 and 4 show the results for two distinct queries evaluated across different model versions and stages.

Figures 3–4 depict retrieval grids for two queries across all model variants. Key observations include:

- **Stage 2 models** often regress to colour-only matches, losing structural fidelity.
- **Stage 3 v3 (RN50)** consistently retrieves near-duplicates, validating the quantitative lead.
- Prompt trimming (v4) benefits ViT-B/16 but can degrade RN50 when descriptions become too generic.

### E. Failure Modes and Future Work

Persistent challenges involve:

- *Occlusions & pose variance*: partial views bias the model towards colour over garment shape.
- *Prompt oversimplification*: overly terse prompts reduce fine-grained discrimination.
- *Colour dominance*: early stages over-weight hue, retrieving structurally mismatched items.
- *Dataset bias*: limited annotation of subtle attributes (e.g., texture, fit) constrains supervision.

Future improvements will investigate multi-view training, adaptive prompt generation, pose-aware features, hard-negative mining, and debiasing strategies.

### Summary

Stage 3 joint fine-tuning with prompt-aware SupCon learning and balanced sampling delivers the most effective retrieval, with RN50 v3 setting the new benchmark (Rank-1 53.95 %, mAP 0.4265) while maintaining sub-50 ms latency—demonstrating practical readiness for deployment in visual fashion search.

### Conclusion

## X. SUMMARY OF CONTRIBUTIONS

This project presented a comprehensive pipeline for fine-tuning the CLIP model for fashion-specific image retrieval, using the DeepFashion2 dataset as a benchmark. It systematically explored and evaluated different fine-tuning strategies, from zero-shot baseline to multimodal supervised contrastive training.

The key contributions of this project include:

- Establishing a multi-stage CLIP fine-tuning framework with distinct strategies: image-only tuning (Stage 2), joint tuning (Stage 3), and prompt-guided adaptation.
- Implementing supervised contrastive learning (SupCon) with identity-aware  $P \times K$  batch sampling, enabling effective separation of garment identities in the embedding space.
- Designing structured and trimmed prompt templates that allowed text supervision while respecting CLIP’s 77-token limit.
- Comparing model behavior across two architectures (ViT-B/16 and RN50), highlighting architecture-specific strengths and weaknesses.
- Conducting detailed qualitative and quantitative evaluations using re-identification-style splits, supported by visual inspection and ranking consistency analysis.
- Achieving state-of-the-art zero-shot and fine-tuned performance with RN50 Stage 3 v3 (Rank-1: 53.95%, mAP: 0.4265) and ViT-B/16 Stage 3 v4.

## XI. REFLECTION

The project provided valuable insights into how vision-language models like CLIP can be effectively adapted for fine-grained fashion retrieval. Key reflections include:

### a) What worked well::

- $P \times K$  sampling and SupConLoss proved essential for identity-preserving embedding learning.
- Prompt engineering—especially template structuring and token limit control—was critical for effective supervision.
- Modular training pipelines allowed scalable experimentation across multiple architectures and strategies.
- RN50 demonstrated superior capacity to encode visual features such as texture, structure, and silhouette.



Fig. 3. Qualitative retrieval results for Sample 1. The first image is the query, followed by retrieved results across different CLIP architectures (ViT-B/16 and RN50) and training stages (Baseline, Stage 2, Stage 3 v1–v4).

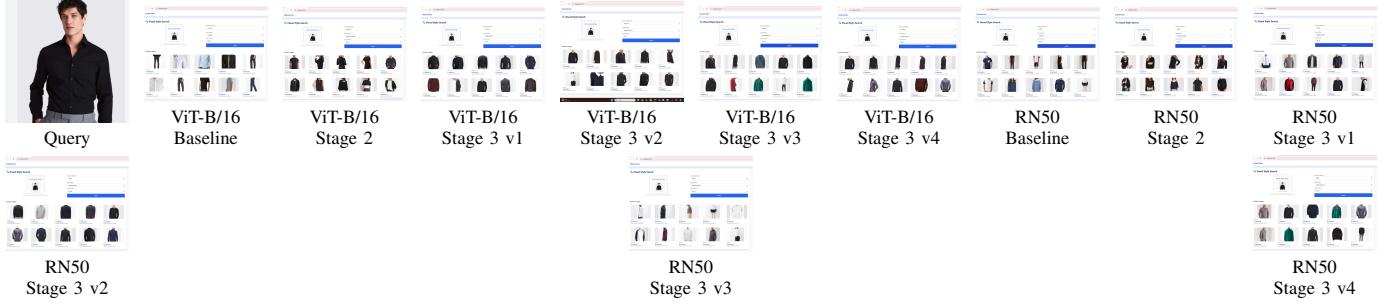


Fig. 4. Qualitative retrieval results for Sample 2. The figure shows retrieval outcomes across model architectures and training strategies, illustrating differences in semantic understanding and visual matching.

### Answering the Research Questions

The outcomes of this project address the core research questions as follows:

- 1) **How can AI improve visual search accuracy in fashion retail?** Through contrastive learning, prompt tuning, and architecture adaptation, AI can significantly improve semantic alignment and visual relevance in top-K retrieval tasks. The fine-tuned CLIP models outperformed their zero-shot counterparts, demonstrating the effectiveness of domain-adapted AI in fashion applications.
- 2) **Does fine-tuning the CLIP model improve retrieval accuracy for fashion visual search compared to zero-shot performance?** Yes. Stage 2 and especially Stage 3 fine-tuning provided substantial gains in both Rank-1 accuracy and mAP, confirming that model adaptation to the fashion domain enhances retrieval performance.
- 3) **Between the two CLIP backbones—ViT-B/16 and RN50—which architecture demonstrates better generalization and domain adaptation for fashion image retrieval?** RN50 consistently outperformed ViT-B/16 in fine-tuned scenarios, particularly in Stage 3 v3, suggesting stronger generalization and better image-text embedding alignment in RN50.
- 4) **To what extent do prompt engineering and supervised contrastive loss (SupCon) influence the semantic alignment between image and text embeddings in joint training?** SupCon loss and structured prompt learning (especially using P×K sampling) were

instrumental in improving semantic consistency. These techniques helped the model retrieve stylistically and categorically relevant garments under challenging test cases.

#### b) What didn't work well::

- Stage 2 (image-only fine-tuning) disrupted CLIP's pre-trained alignment, yielding degraded performance.
- Over-trimmed prompts in Stage 3 v4 occasionally removed key garment semantics, especially in RN50.
- Dataset inconsistencies (e.g., class imbalance, ambiguous clothing angles) caused occasional ranking errors.
- Early-stage training lacked contrastive hardness, leading to misranked visually accurate matches.

## XII. FUTURE WORK

While the results were promising, several areas for improvement and future research were identified:

- **Multi-view Integration:** Incorporate multiple views of the same garment to improve embedding robustness and pose invariance.
- **Dynamic Prompt Generation:** Move beyond static class prompts to generate context-aware or vision-guided prompts during training and inference.
- **Pose and Landmark Embedding:** Use human pose estimations or garment landmarks to guide the model toward structural understanding.
- **Hard Negative Mining:** Introduce loss-aware hard negative selection to discourage visually deceptive false positives.

- **Bias Mitigation:** Explore domain-adaptive debiasing strategies to address residual biases from CLIP’s original pretraining (e.g., gender or category imbalance).
- **Cross-Domain Generalization:** Test the pipeline on datasets beyond DeepFashion2, such as streetwear or user-uploaded photos, to validate retrieval generalization.
- **Real-world Deployment:** Integrate the trained model into a real-time product search system, including latency profiling, indexing optimization, and scalable embedding retrieval.

In summary, this work lays a practical and technical foundation for future vision-language retrieval systems tailored to the fashion domain. The techniques and insights developed can be extended to other fine-grained recognition tasks involving multimodal alignment.

## REFERENCES

- [1] A. B. Mathematics, K. Kumar, and S. Singla, “Multimodal deep learning: Integrating text and image embeddings with attention mechanism,” in *2024 3rd International Conference on Artificial Intelligence For Internet of Things (AIoT)*. IEEE, 2024, pp. 1–6.
- [2] E. Nitasha, S. Kumari, A. Kumar, R. Bhardwaj, V. Maddheshiya, and A. Khan, “Future of fashion: Ai-powered virtual dressing for e-commerce applications,” in *2024 International Conference on Emerging Innovations and Advanced Computing (INNOCOMP)*. IEEE, 2024, pp. 138–150.
- [3] S. M. Islam, S. Joardar, and A. A. Sekh\*, “A survey on fashion image retrieval,” *ACM Computing Surveys*, vol. 56, no. 6, pp. 1–25, 2024.
- [4] Y. Ge, R. Zhang, and P. Luo, “Metacloth: learning unseen tasks of dense fashion landmark detection from a few samples,” *IEEE Transactions on Image Processing*, vol. 31, pp. 1120–1133, 2021.
- [5] Q. Chen, T. Zhang, M. Nie, Z. Wang, S. Xu, W. Shi, and Z. Cao, “Fashion-gpt: integrating llms with fashion retrieval system,” in *Proceedings of the 1st Workshop on Large Generative Models Meet Multimodal Applications*, 2023, pp. 69–78.
- [6] R. Datta, D. Joshi, J. Li, and J. Z. Wang, “Image retrieval: Ideas, influences, and trends of the new age,” *ACM Computing Surveys (Csur)*, vol. 40, no. 2, pp. 1–60, 2008.
- [7] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [10] L. Ge, Z. Ren, Y. Li, Z. Xue, Y. Wang, J. Cai, and J. Yuan, “3d hand shape and pose estimation from a single rgb image,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 10833–10842.
- [11] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [12] J. Li, D. Li, C. Xiong, and S. Hoi, “Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation,” in *International conference on machine learning*. PMLR, 2022, pp. 12 888–12 900.
- [13] R. Jang, “Learning representations by forward-propagating errors,” *arXiv preprint arXiv:2308.09728*, 2023.
- [14] A. Baldrati, M. Bertini, T. Uricchio, and A. Del Bimbo, “Effective conditioned and composed image retrieval combining clip-based features,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 21 466–21 474.
- [15] D. Gao, L. Jin, B. Chen, M. Qiu, P. Li, Y. Wei, Y. Hu, and H. Wang, “Fashionbert: Text and image matching with adaptive loss for cross-modal retrieval,” in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 2251–2260.
- [16] V. Likhoshesterov, A. Arnab, K. Choromanski, M. Lucic, Y. Tay, A. Weller, and M. Dehghani, “Polyvit: Co-training vision transformers on images, videos and audio,” *arXiv preprint arXiv:2111.12993*, 2021.
- [17] I. O. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, A. Steiner, D. Keysers, J. Uszkoreit *et al.*, “Mlp-mixer: An all-mlp architecture for vision,” *Advances in neural information processing systems*, vol. 34, pp. 24 261–24 272, 2021.
- [18] Y. Cao, Z. He, L. Wang, W. Wang, Y. Yuan, D. Zhang, J. Zhang, P. Zhu, L. Van Gool, J. Han *et al.*, “Visdrone-det2021: The vision meets drone object detection challenge results,” in *Proceedings of the IEEE/CVF International conference on computer vision*, 2021, pp. 2847–2854.
- [19] F. Schroff, D. Kalenichenko, and J. Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 815–823.
- [20] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with gpus,” *IEEE Transactions on Big Data*, vol. 7, no. 3, pp. 535–547, 2019.
- [21] M. Douze, A. Guzhva, C. Deng, J. Johnson, G. Szilvassy, P.-E. Mazaré, M. Lomeli, L. Hosseini, and H. Jégou, “The faiss library,” *arXiv preprint arXiv:2401.08281*, 2024.