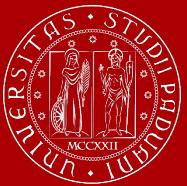


UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Modular Fashion Item Detection and Semantic Matching for Smart Outfit Retrieval

Mattia Roccatello, Ulaşcan Akbulut, Ayşenur Oya Özen



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

The project's purpose



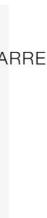
SEQUINNE... ⚡
39,95 EUR
DONNA



ASYMMET... ⚡
59,95 EUR
UOMO



FEW ITEMS LEFT
SATIN MIDI... ⚡
39,95 EUR



MATTIA • SUPPORTO

CARRELLO [0]

COCOA SU... ⚡
19,95 EUR

HIGH-WAIS... ⚡
25,95 EUR



STRETCH ... ⚡



LINEN BLE... ⚡



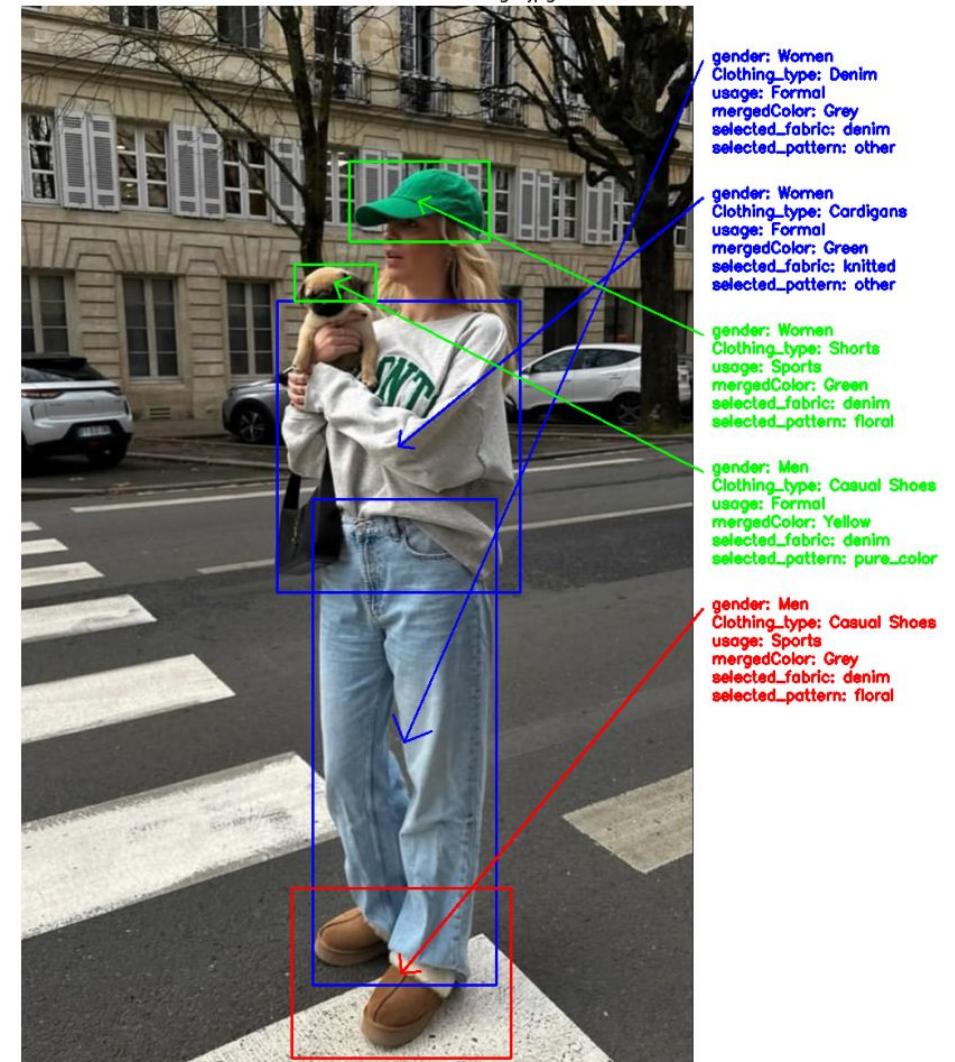
HALTER SA... ⚡



POLYAMID... ⚡



ZW COLLE... ⚡





Why Fashion Items? Why This Work?

Fashion understanding requires

fine-grained detection + semantic understanding + retrieval
→ ALL TOGETHER

Traditional end-to-end models fail to generalize across heterogeneous fashion domains

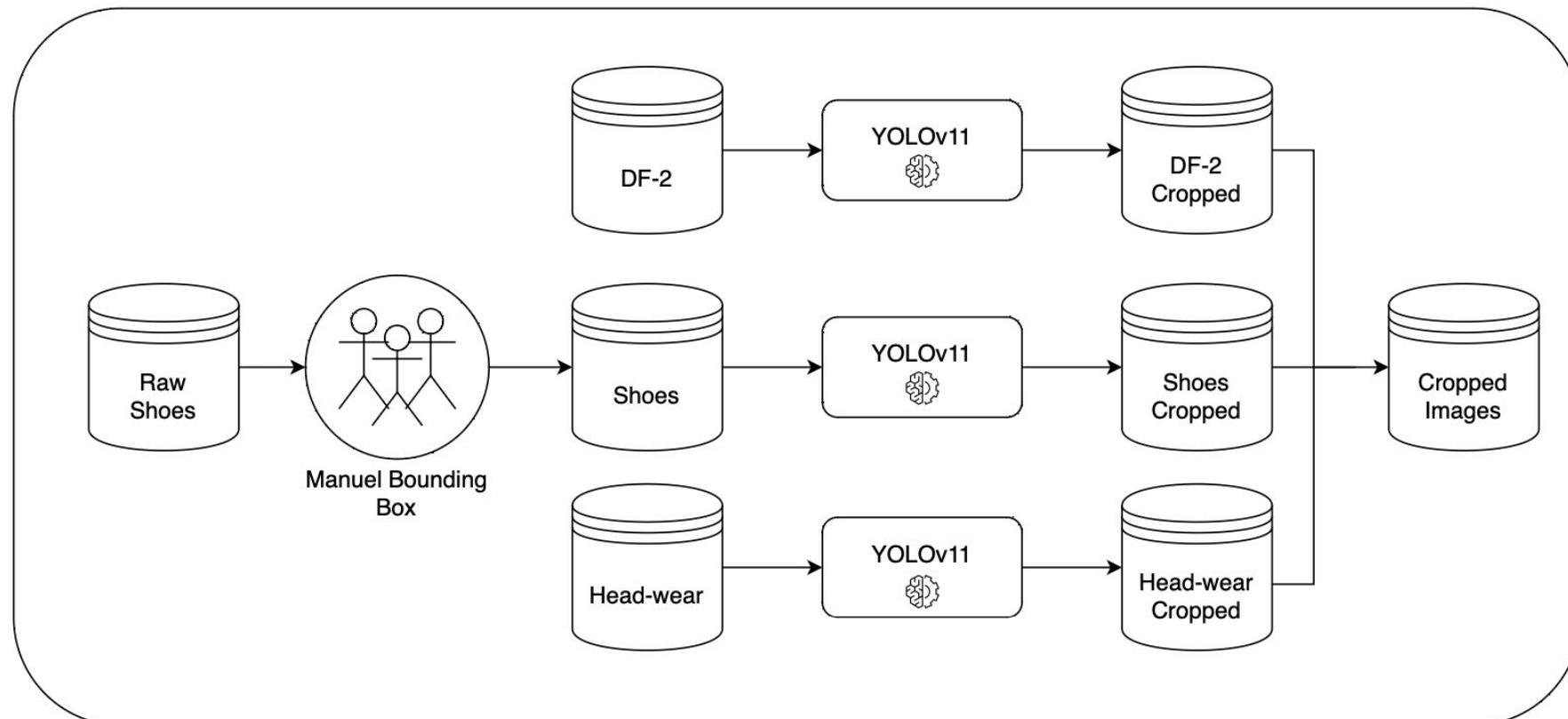
We need a **modular**, **scalable**, and **interpretable** solution that can:

- Detect clothing items
- Classify attributes
- Perform semantic retrieval



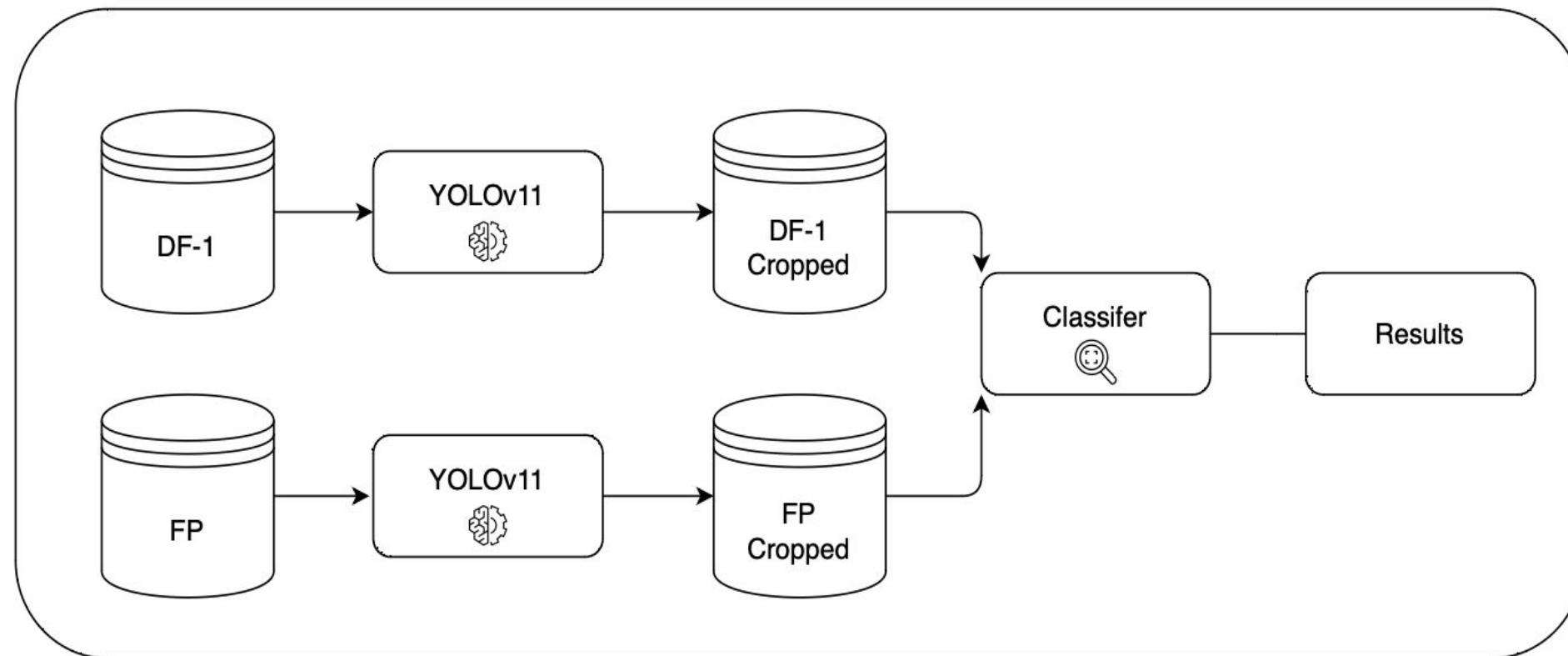


Modular Multi-Stage Pipeline





Modular Multi-Stage Pipeline



Detection → Classification → Matching



Our Contributions



Modular pipeline

A two-stage system with separate models for detection and classification to overcome incomplete datasets

Dataset-specific detectors

Trained separate YOLOv11x models on DeepFashion2, Shoes, and Headwear datasets to ensure full clothing category coverage

Multi-head classifier

Predicts attributes like gender, color, fabric, pattern, and usage using interleaved learning from DeepFashion1 and Fashion Product datasets

Custom matching logic

Bridges instance-level crops with image-level labels using filename structure and predicted types

Contrastive pretraining

Self-supervised SimCLR-based encoder improves representation robustness for subtle features like fabric and pattern

Semantic retrieval

Uses embedding similarity to match query outfit items with similar pieces in a virtual wardrobe



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Tackling Dataset - Detection Data

Scale: 2, Viewpoint: 2, Occlusion: 1, Zoom_in: 1, Set: train, Category: trousers

DeepFashion2

Headwear

Shoes

DeepFashion1

Fashion Products



800K images, 13 categories → Merged to 5 superclasses



Tackling Dataset - Detection Data

DeepFashion2
Headwear
Shoes
DeepFashion1
Fashion Products



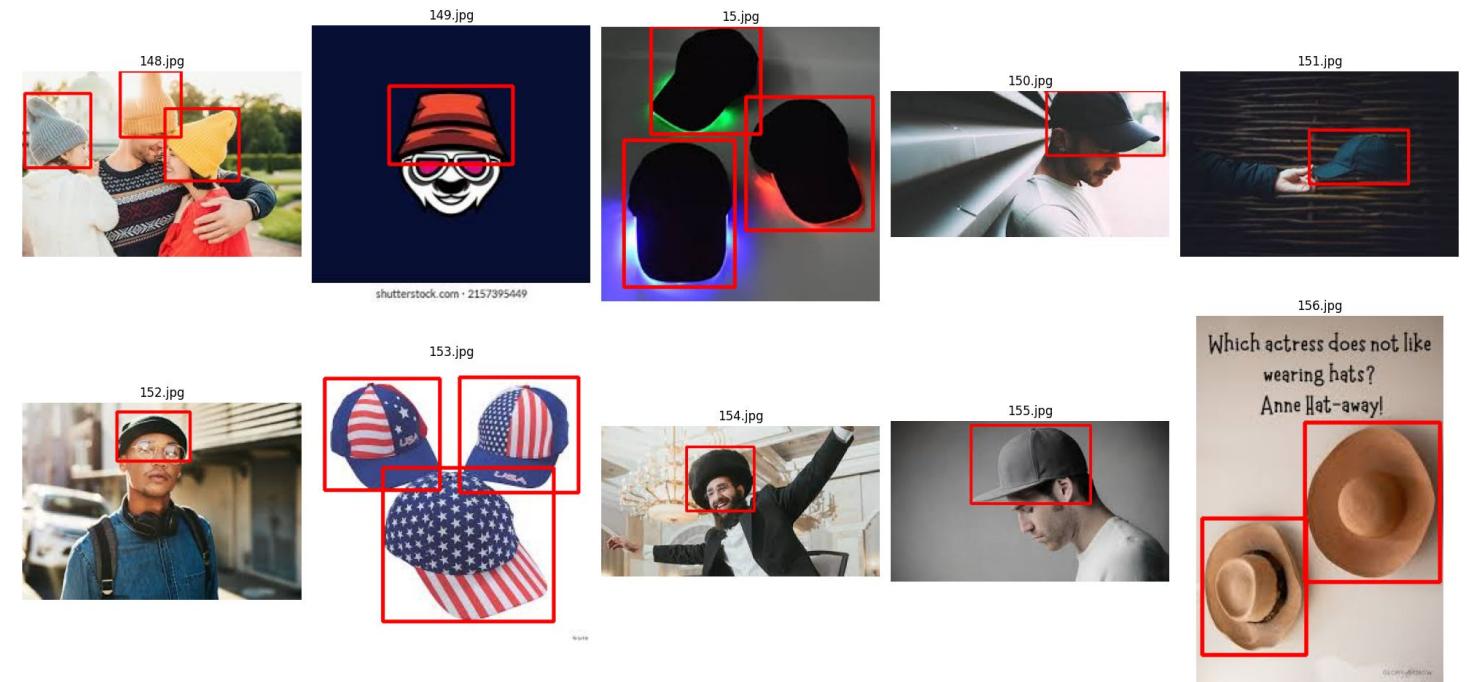
800K images, 13 categories → Merged to 5 superclasses

long sleeve outwear, short sleeve outwear	outwear
short sleeve dress, long sleeve dress, vest dress, sling, sling dress	dress
trousers, shorts, skirt	bottom
long sleeve top, short sleeve top	top
vest	vest



Tackling Dataset - Detection Data

DeepFashion2
Headwear
Shoes
DeepFashion1
Fashion Products



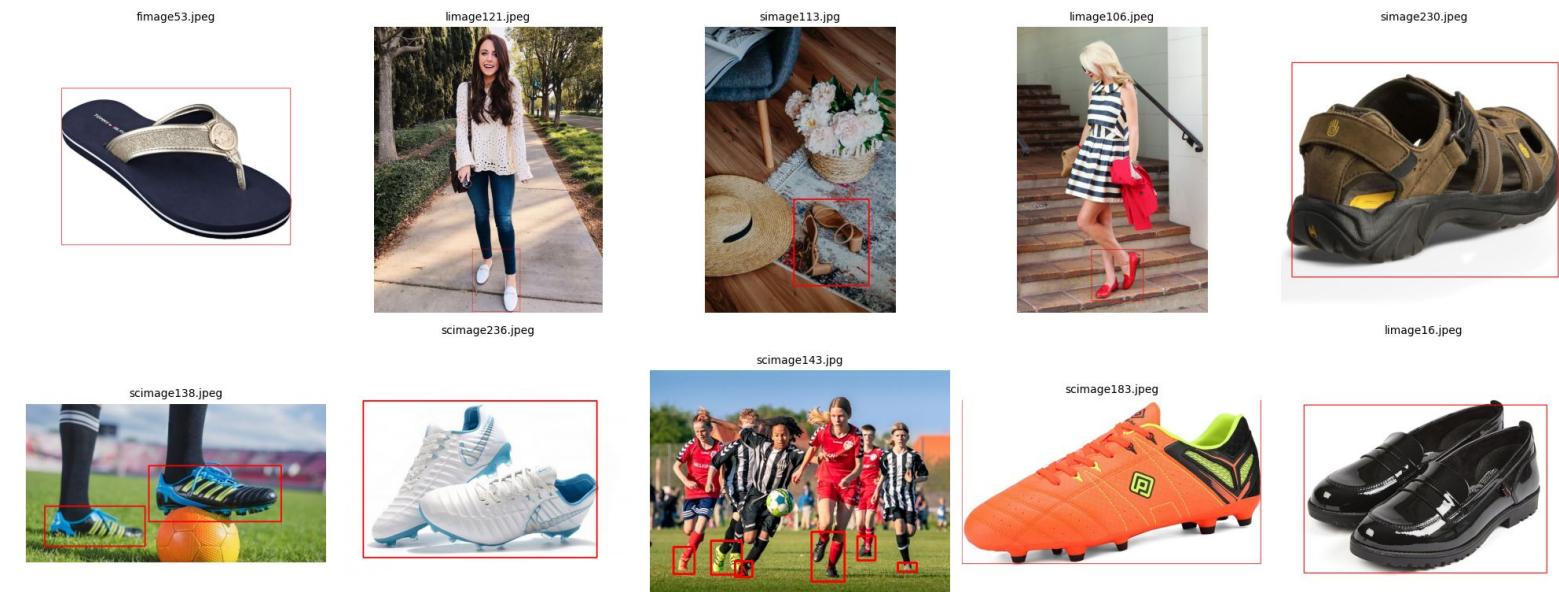
Pascal VOC to YOLO format, 650 images



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Tackling Dataset - Detection Data

DeepFashion2
Headwear
Shoes
DeepFashion1
Fashion Products



Manually labeled, ~12K images



Tackling Dataset - Classification Data

DeepFashion2

Headwear

Shoes

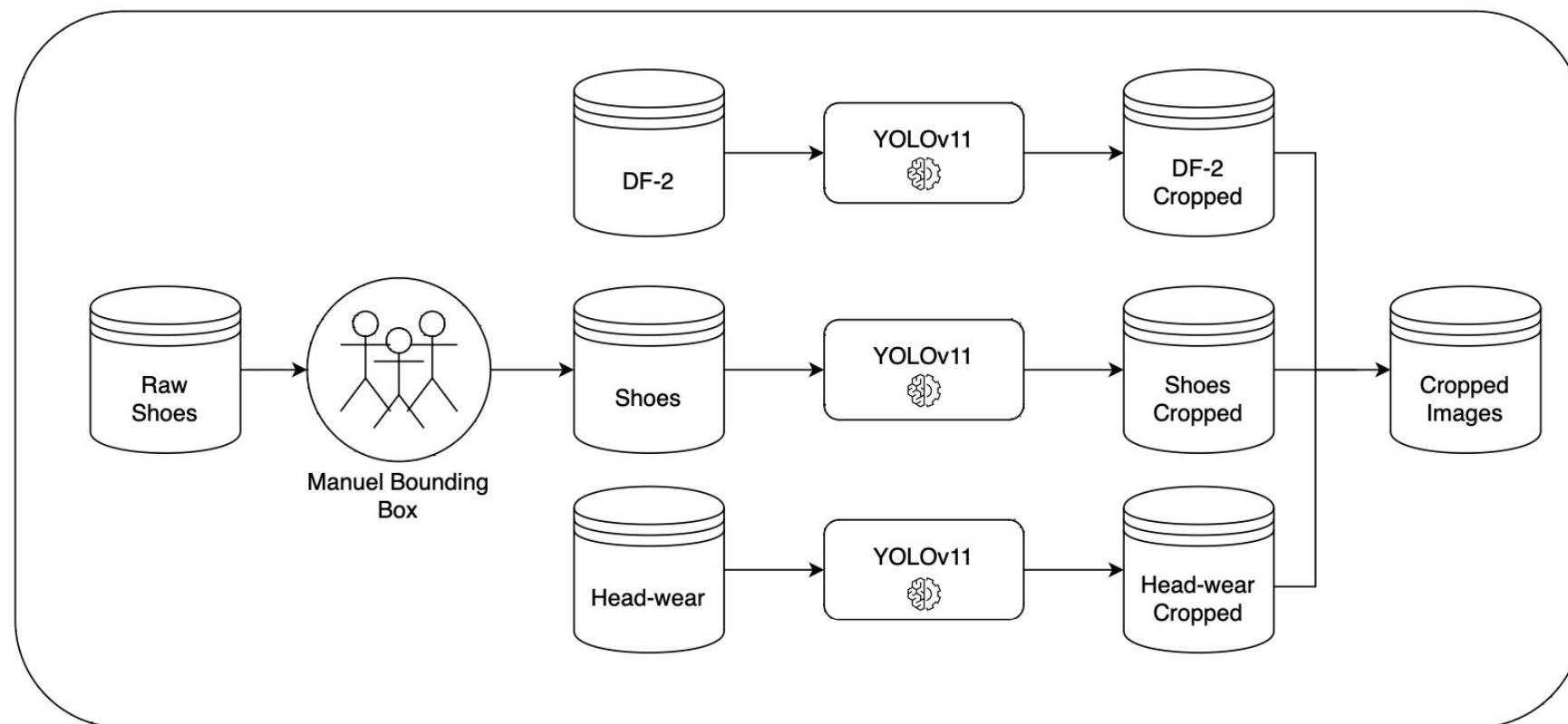
DeepFashion1- Rich attributes but weak crop-to-label mapping

Fashion Products - Strong structured metadata

Attributes: gender, color, usage, article type



Modular Multi-Stage Pipeline - Detection Module

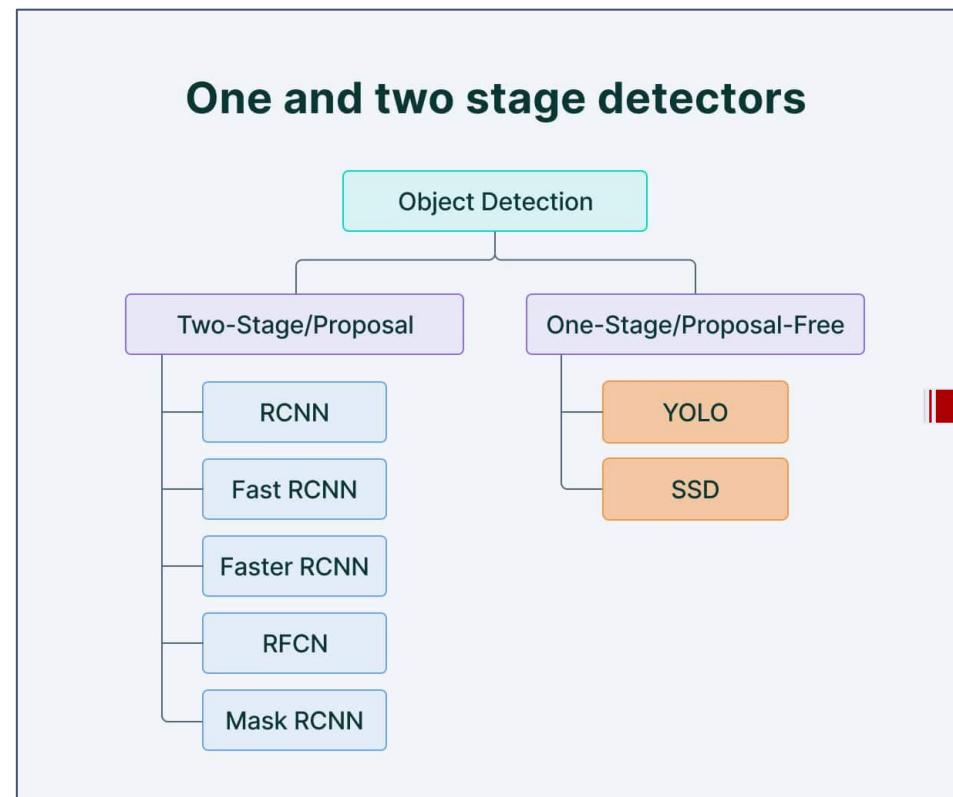




YOLOv11x Detection Models - Why YOLO?

Detects **what** objects are in an image and **where** they are, in a single forward pass

You Only Look Once !



It divides the input image into an $S \times S$ grid

Each grid cell directly predicts bounding boxes + class probabilities



YOLOv11x Detection Models - YOLOv11x Version

Version	Key Improvements
YOLOv1–v3	Basic object detection, faster than R-CNN
YOLOv4–v5	Added better backbones, anchor boxes, better mAP
YOLOv7 / v8	Ultralytics versions, modular, optimized
YOLOv11x	Our choice: high-capacity version optimized by Ultralytics, handles multiple datasets and small objects well

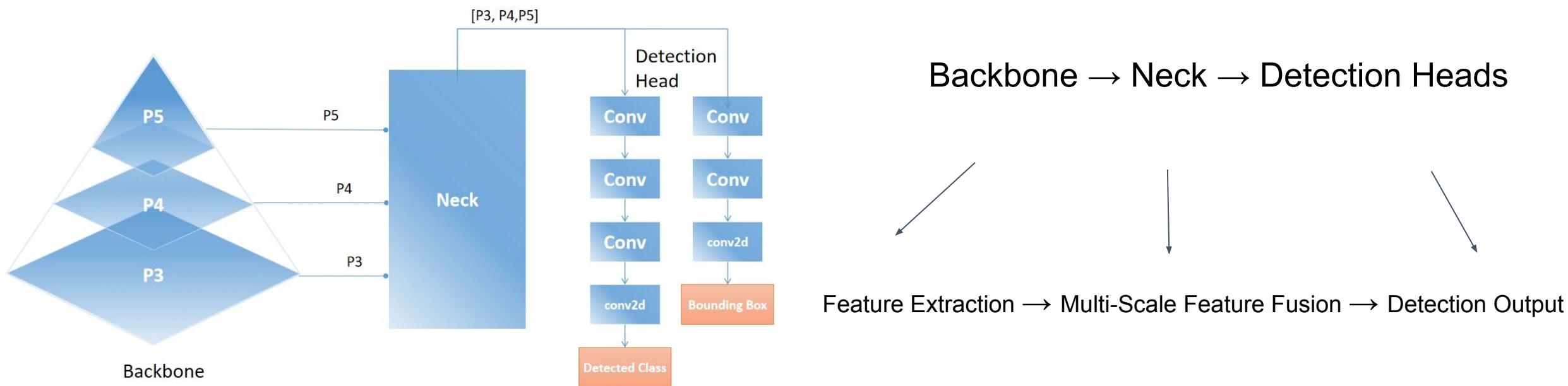
The "x" indicates it's the largest version

- More layers → better feature extraction
- Suited for many small or overlapping objects (like clothing)

Pretrained on COCO, fine-tuned on our datasets



YOLOv11x Detection Models - Architecture





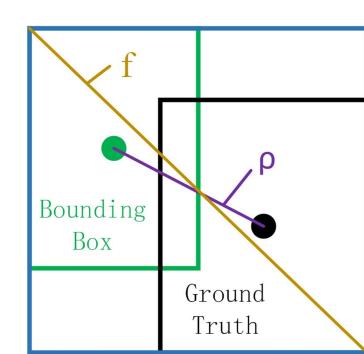
YOLOv11x Detection Models - Our Training Setup

Component	Configuration
Model	YOLOv11x (pretrained on COCO)
Image Size	256 × 256
Epochs	57 (DeepFashion2), 20 (Cap - dataset), 150 (Shoes - dataset)
Batch Size	32 (auto-adjusted for memory)
Optimizer	SGD (momentum = 0.937)
LR Scheduler	Cosine Annealing
Losses	CIoU (bbox), BCE (objectness), CE (classes)
Validation Metrics	Accuracy, precision, recall, mAP@0.5, mAP@0.5-0.9
Data Augmentation	Mosaic, HSV, Flip, Affine
Logging	Weights&Biases
Early Stopping	Enabled

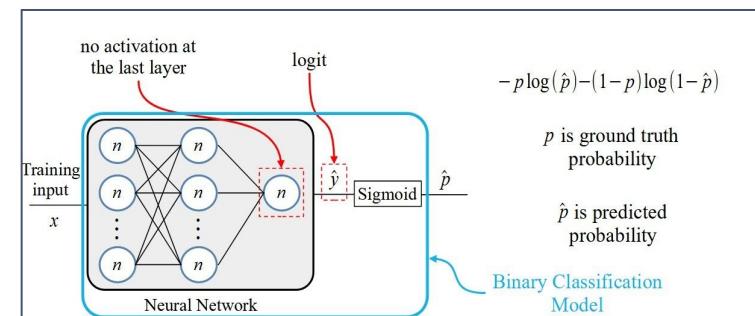
Cross-Entropy Loss

$$H(p, q) = - \sum_{x \in \text{classes}} p(x) \log q(x)$$

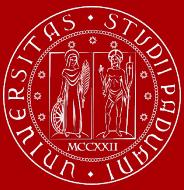
True probability distribution (one-shot)
Your model's predicted probability distribution



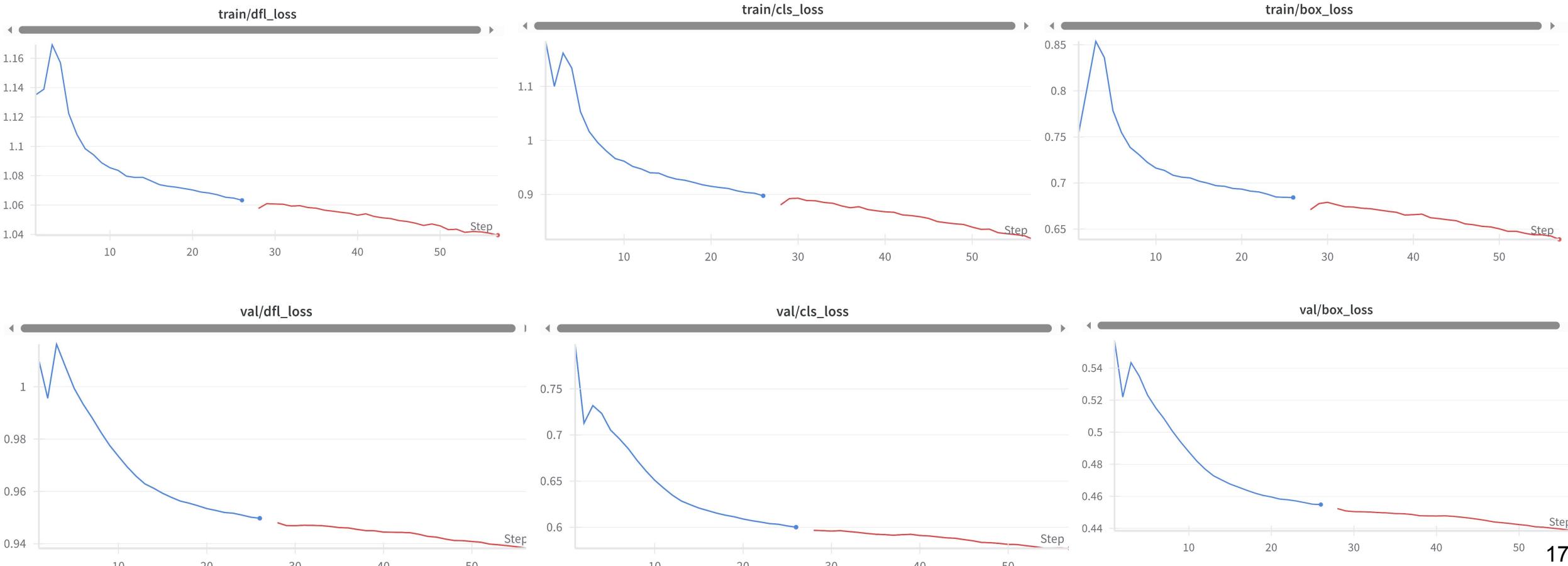
Complete Intersection over Union



Binary Cross-Entropy Loss

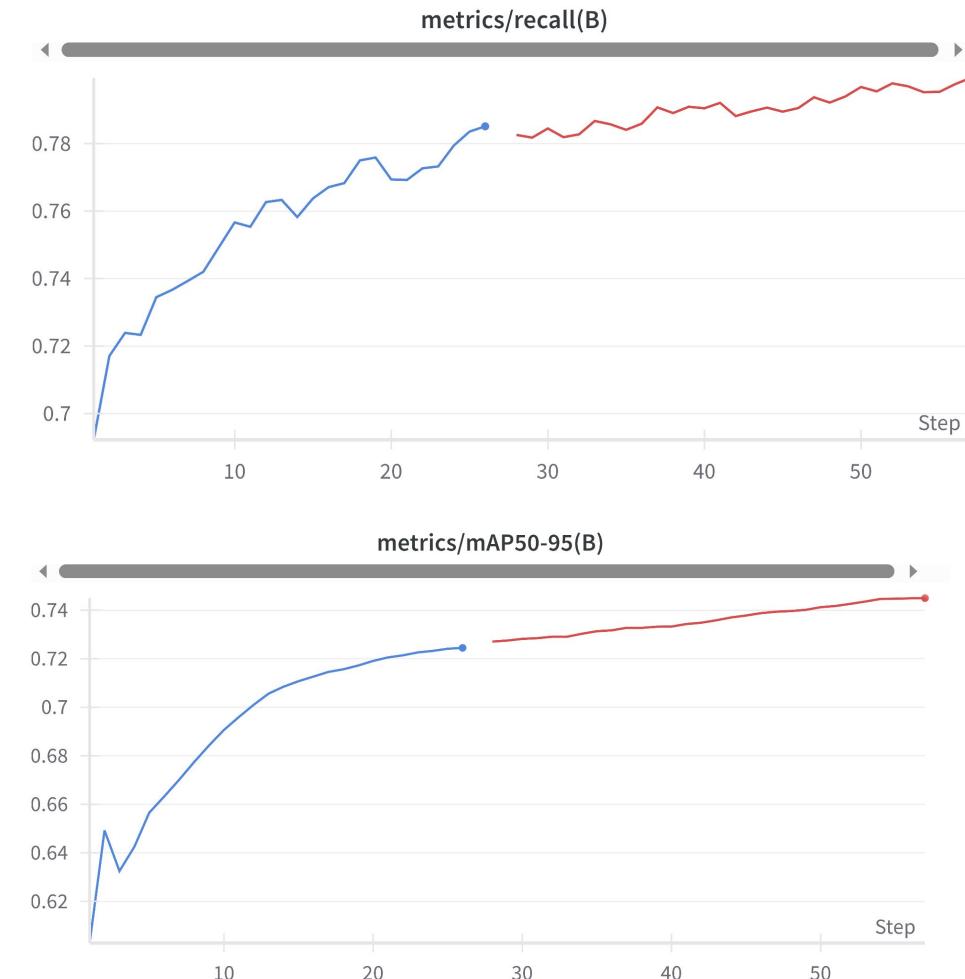
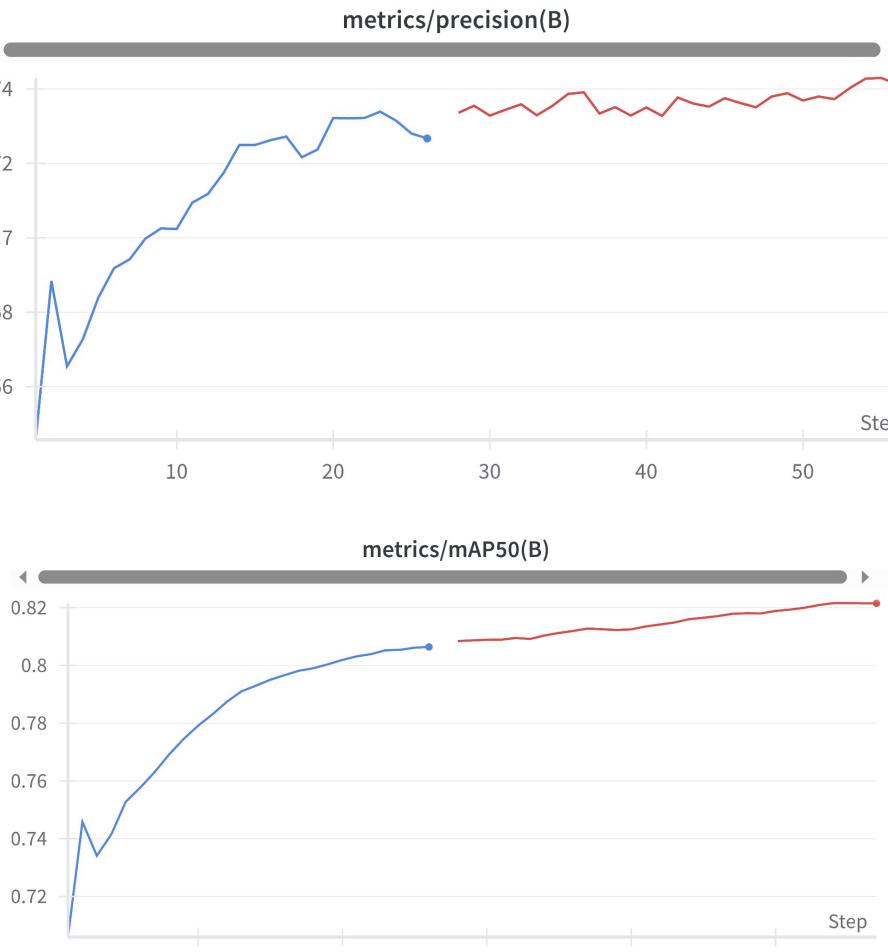


YOLOv11x Detection Models - DF2 Dataset Loss Curves



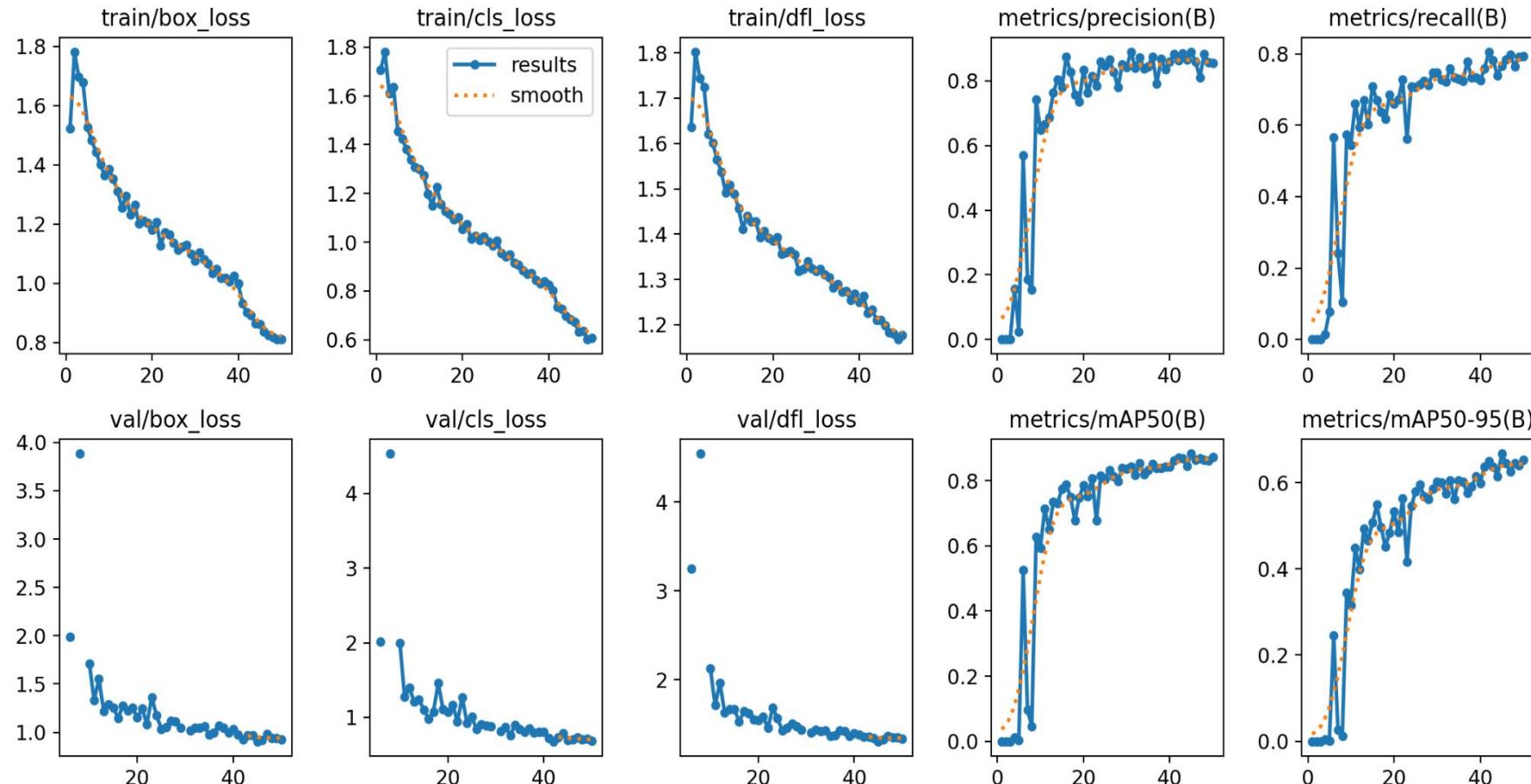


YOLOv11x Detection Models - DF2 Dataset Evaluation Metrics



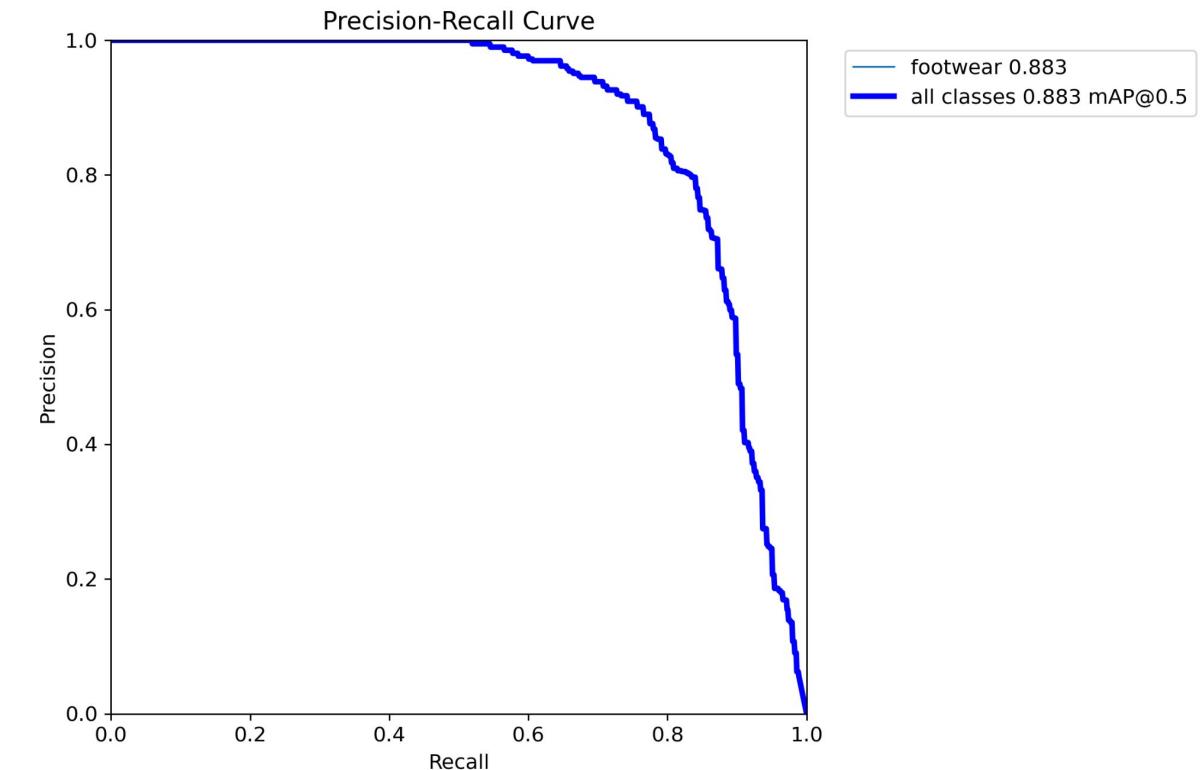
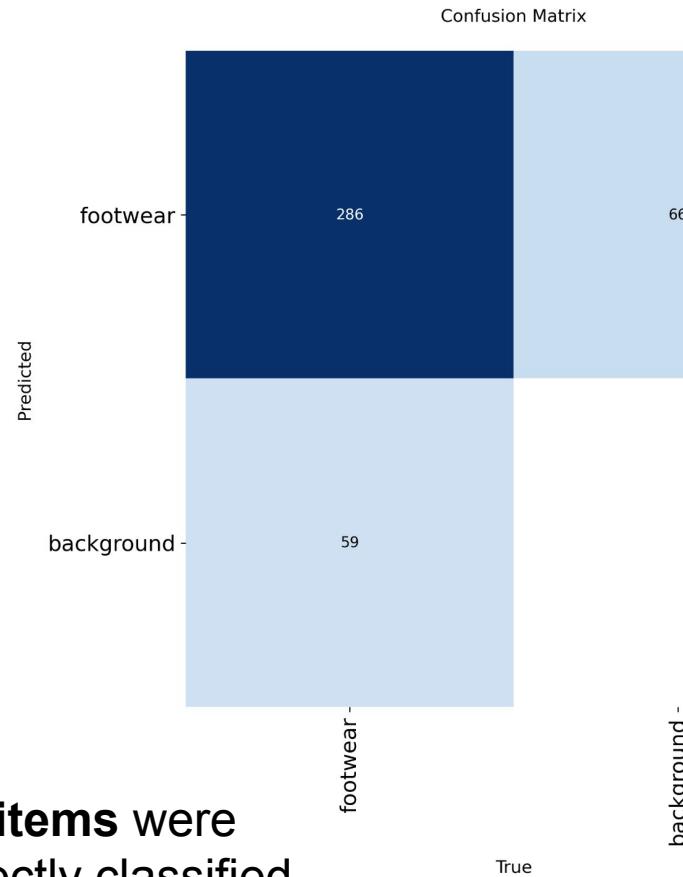


YOLOv11x Detection Models - Shoes Dataset Loss Curves



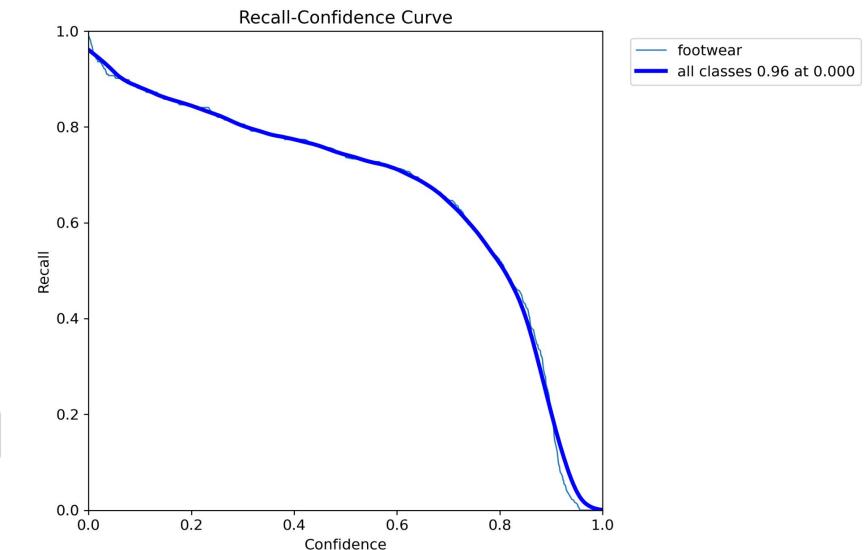
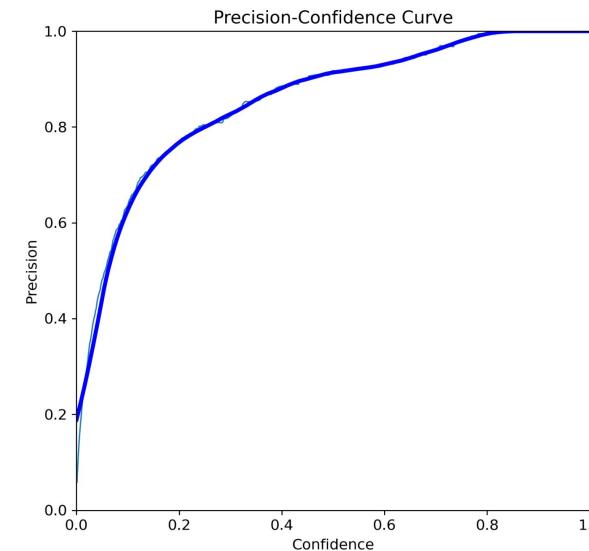
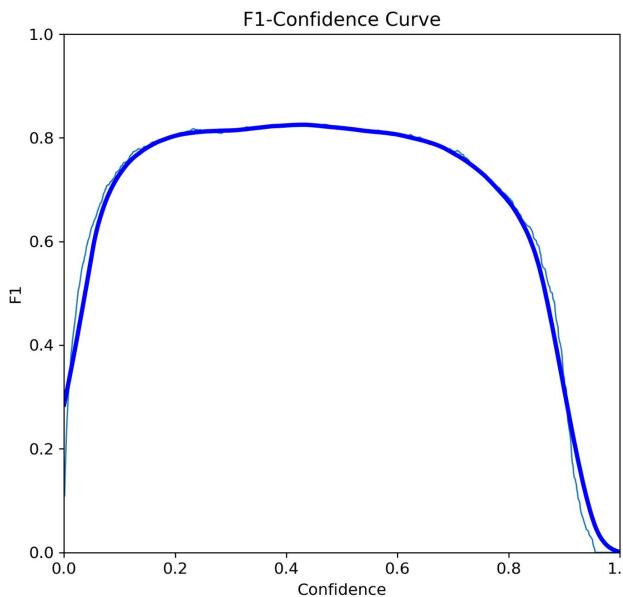


YOLOv11x Detection Models - Shoes Dataset Confusion Matrix & PR Curve





YOLOv11x Detection Models - Shoes Dataset Confidence-Based Metric Curves

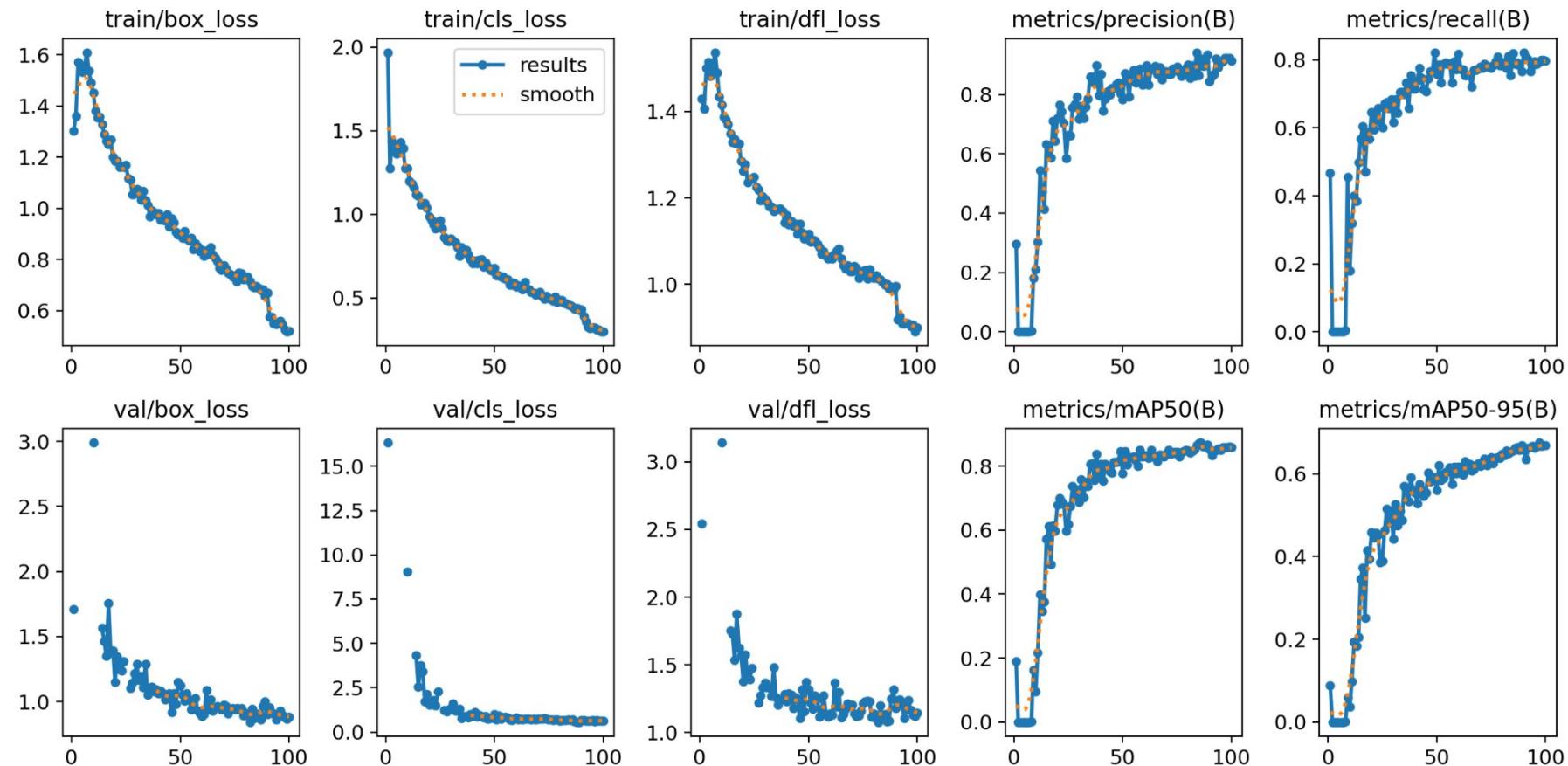


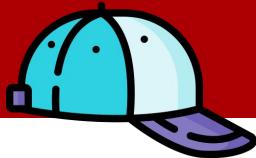


YOLOv11x Detection Models - Shoes Dataset Sample Detection The grid displays 20 shoe images arranged in five rows and four columns. Each image shows a different type of footwear, such as loafers, sneakers, flip-flops, and boots. A blue bounding box is drawn around each shoe, and the text "footwear" followed by a confidence score is displayed in the top-left corner of each image. The confidence scores range from 0.3 to 0.9. | Image | Confidence Score | |----------------|------------------| | limage268.jpeg | footwear 0.8 | | snimage265.jpg | footwear 0.6 | | snimage297.jpg | footwear 0.9 | | fimage270.jpg | footwear 0.9 | | simage255.jpg | footwear 0.4 | | simage264.jpg | footwear 0.3 | | limage288.jpeg | footwear 0.9 | | limage261.jpg | footwear 0.9 | | fimage276.jpeg | footwear 0.9 | | snimage286.jpg | footwear 0.9 | | snimage269.jpg | footwear 0.7 | | sfimage290.jpg | footwear 0.9 | | snimage274.jpg | footwear 0.9 | | simage288.jpg | footwear 0.9 | | fimage291.jpg | footwear 0.9 | | bimage291.jpg | footwear 0.9 | The model successfully detects various shoe types with **high confidence scores!** 22

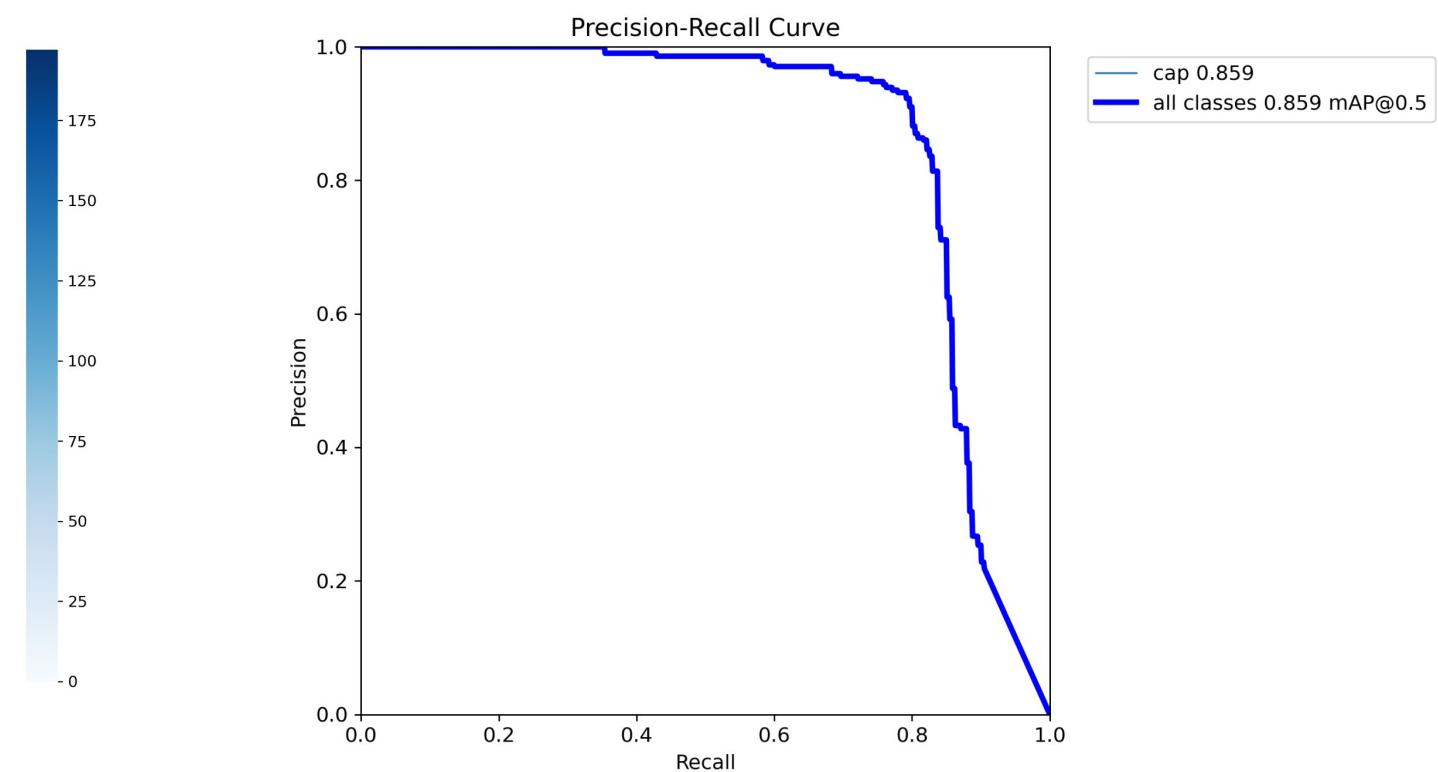
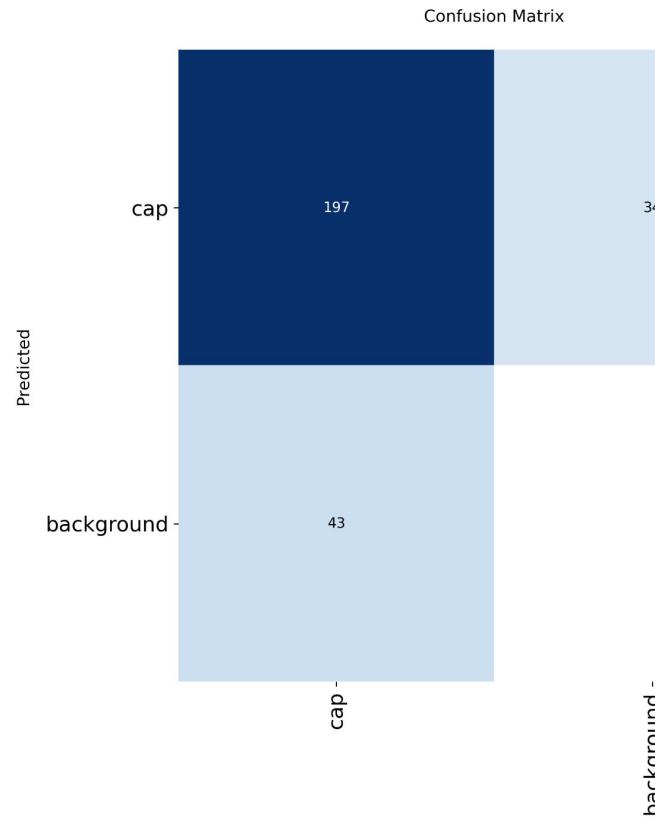


YOLOv11x Detection Models - Headwear Dataset Training Curves





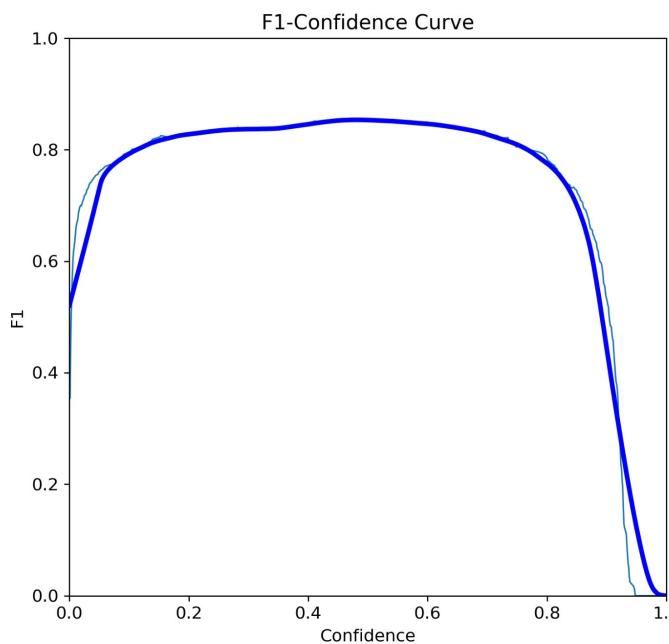
YOLOv11x Detection Models - Headwear Dataset Confusion Matrix & PR Curve



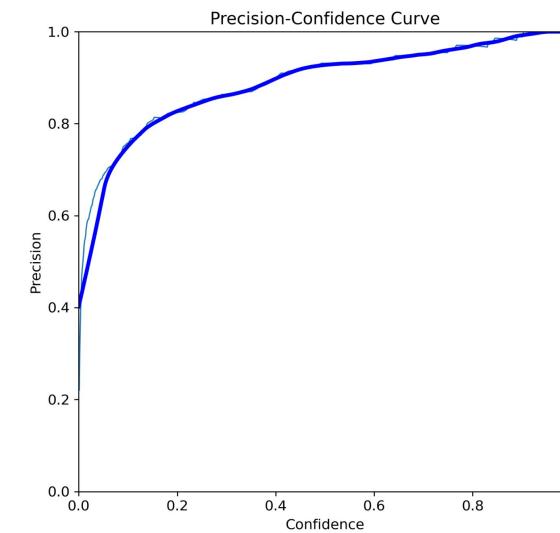
197 items were correctly classified as caps!



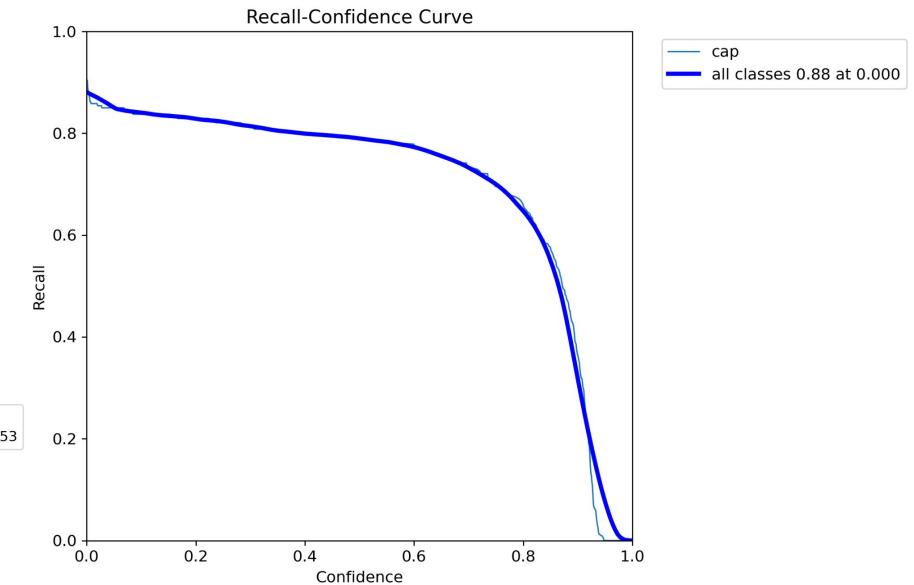
YOLOv11x Detection Models - Headwear Dataset Confidence-Based Metric Curves



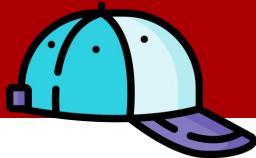
cap
all classes 0.85 at 0.476



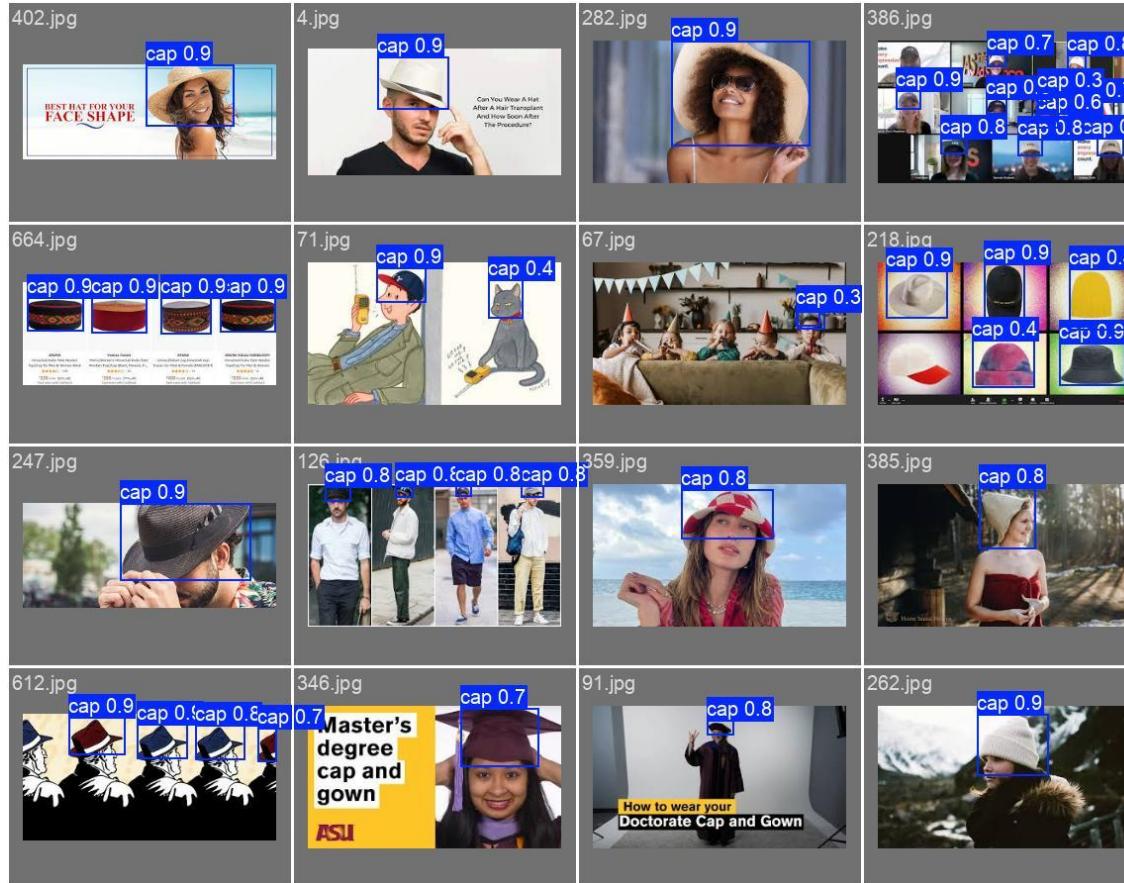
cap
all classes 1.00 at 0.953



cap
all classes 0.88 at 0.000



YOLOv11x Detection Models - Headwear Dataset Sample Detection



The model successfully
detects various headwear
types with **high confidence
scores!**



What We Have At The End with YOLOv11x?

Thanks to these reliable detections:

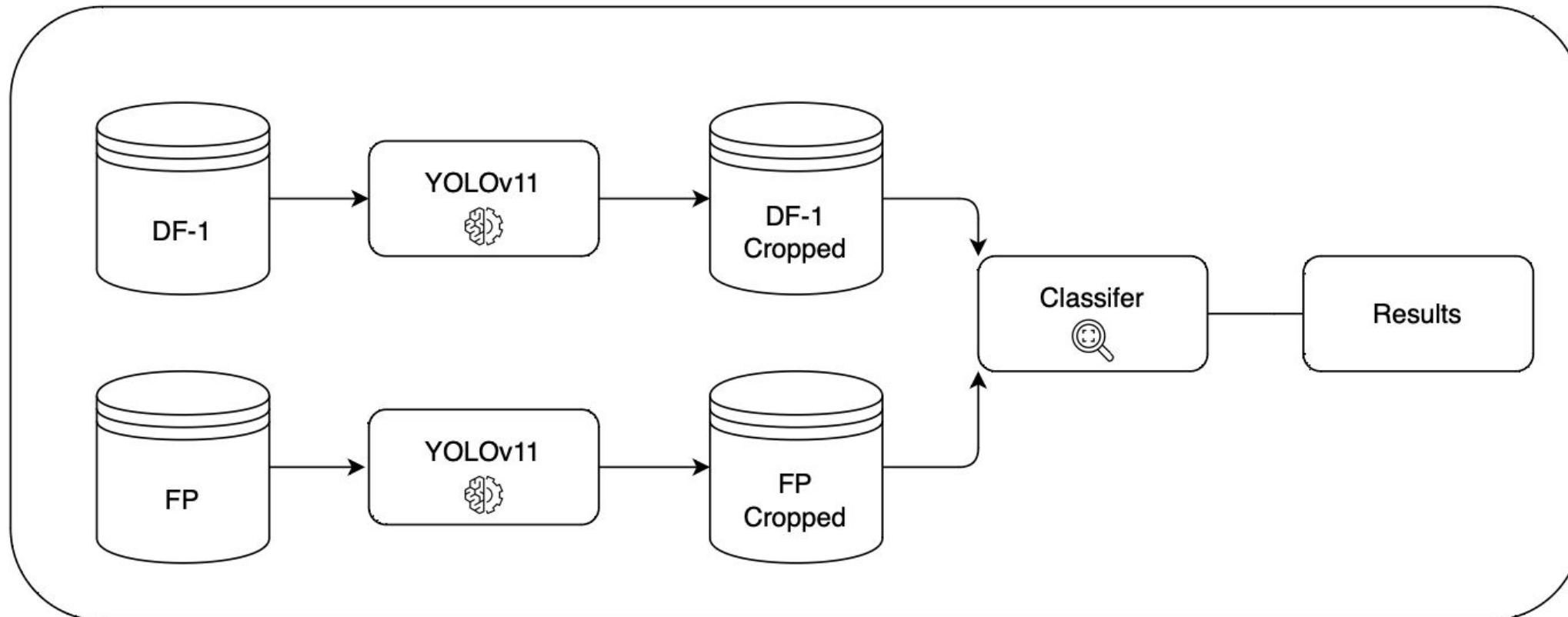
- We can **crop fashion items** from full-outfit images with high precision.
- These cropped items are then used as inputs for our **semantic attribute classifier**, which attempts to identify features like gender, usage, fabric, and color.
- Furthermore, it sets the stage for **matching detected items** with visually or semantically similar products.

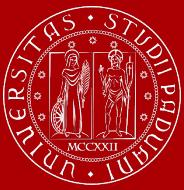
with YOLOv11x, we're not just detecting fashion items!!!

we're **preparing clean, object-level inputs** for the next stages of our pipeline! —> —>



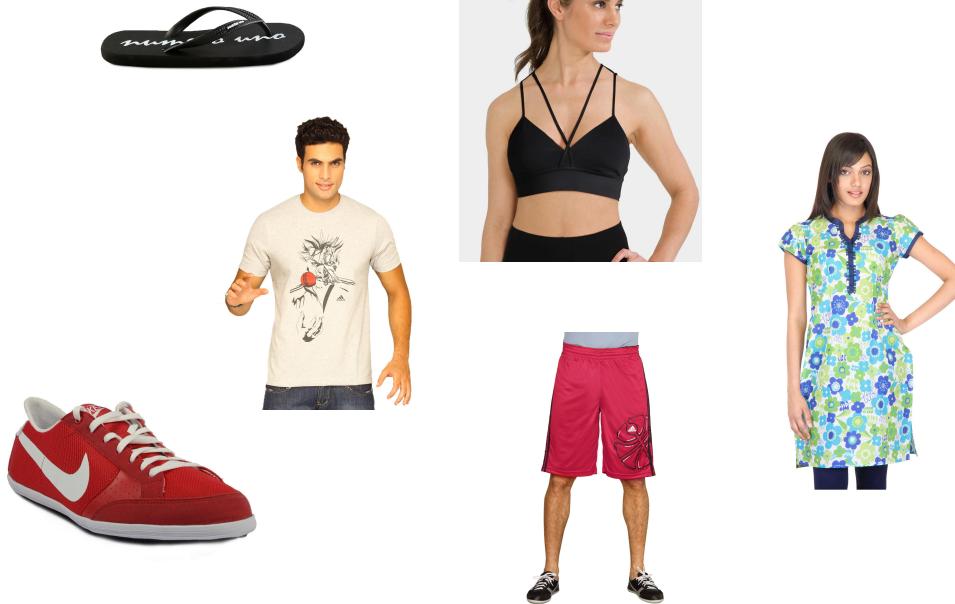
Dataset Preparation for Classification





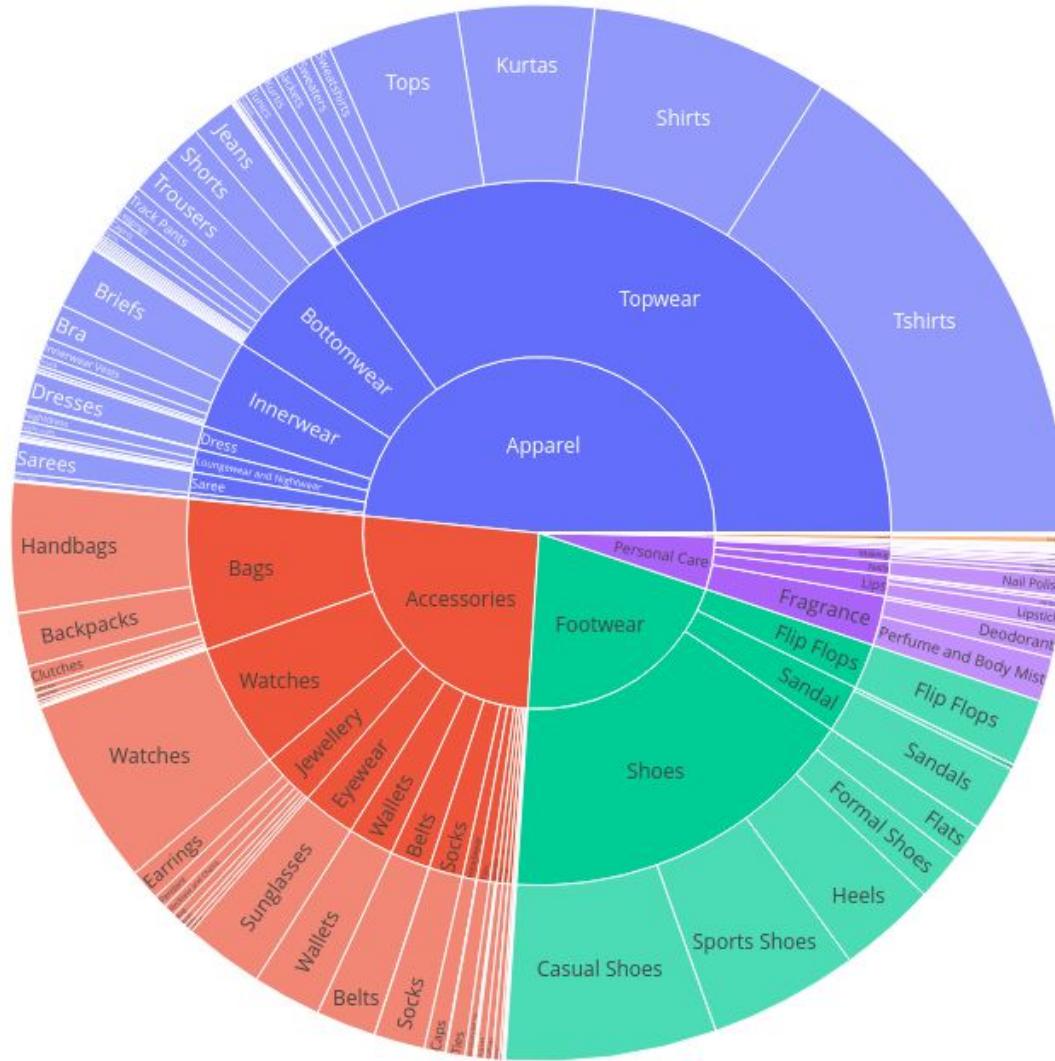
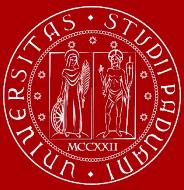
Dataset Preparation for Classification

Fashion Products Images (44.4k images)



DeepFashion (300k items)





Fashion Product Images: data cleaning



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Dropped classes: Year and Season



Fall



Summer



Dropped labels: Unisex, Kids, Girls dropped from Gender

ID: 2710 | subCategory: Dress



ID: 2697 | subCategory: Dress

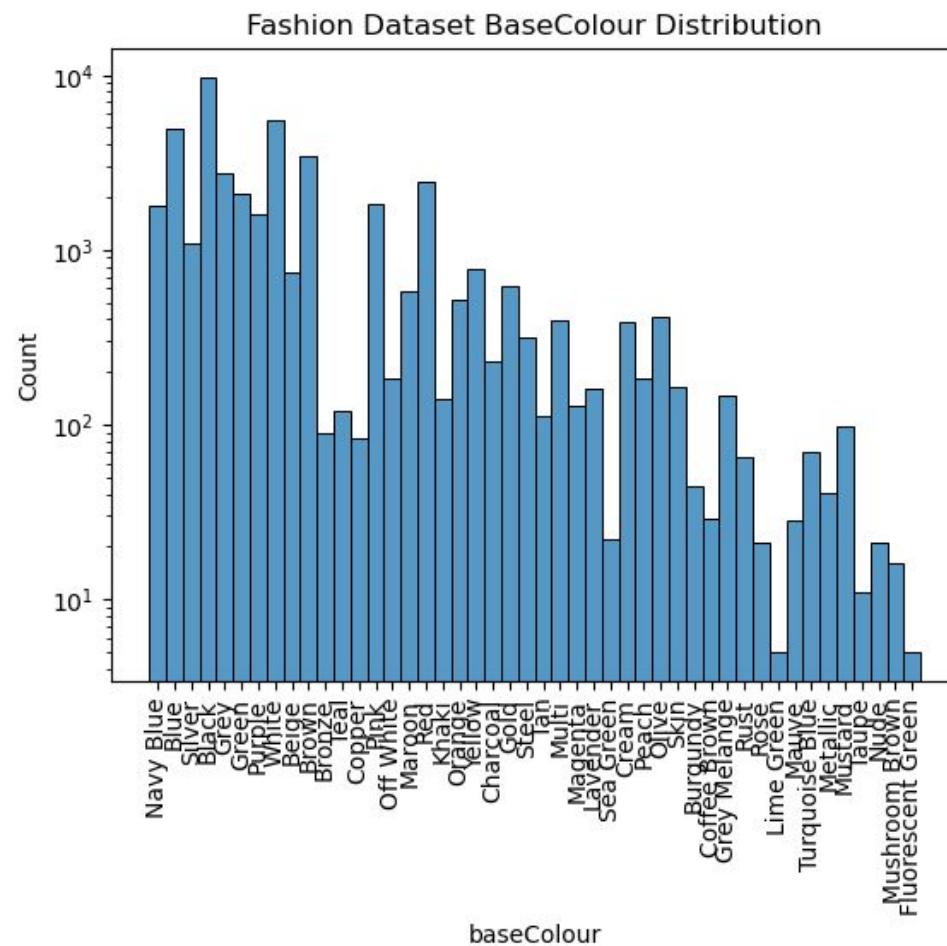


ID: 2699 | subCategory: Dress





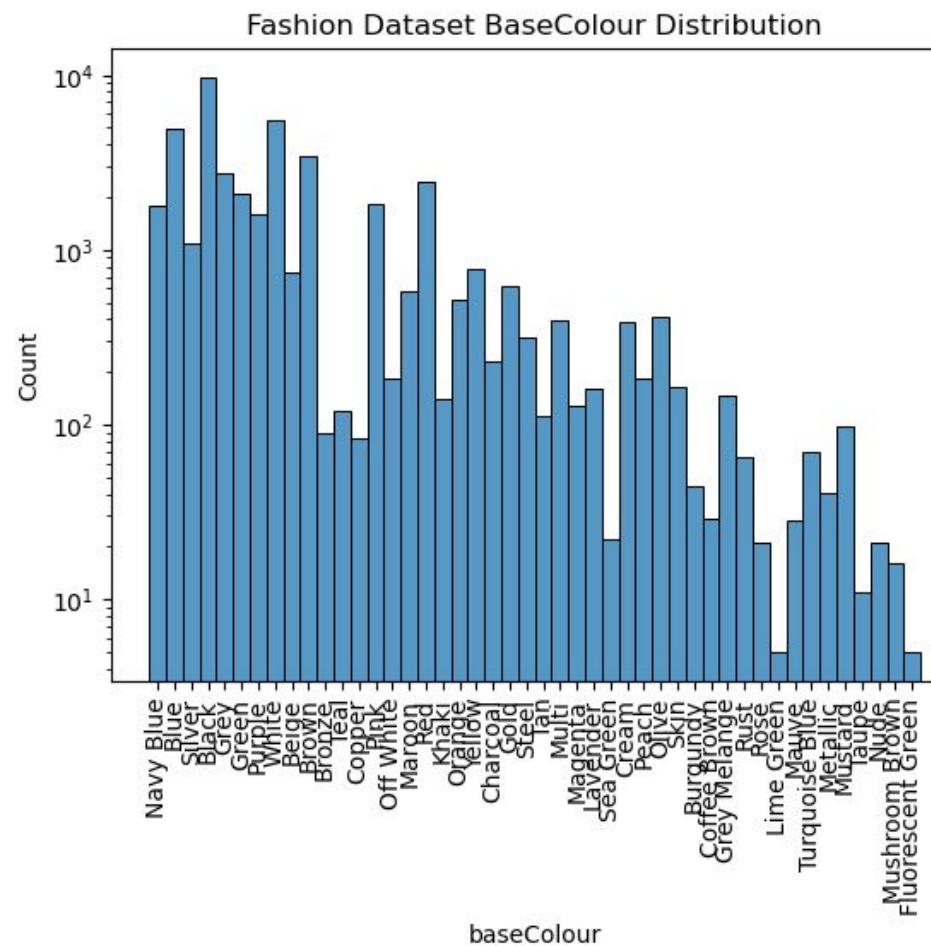
Merged: Color



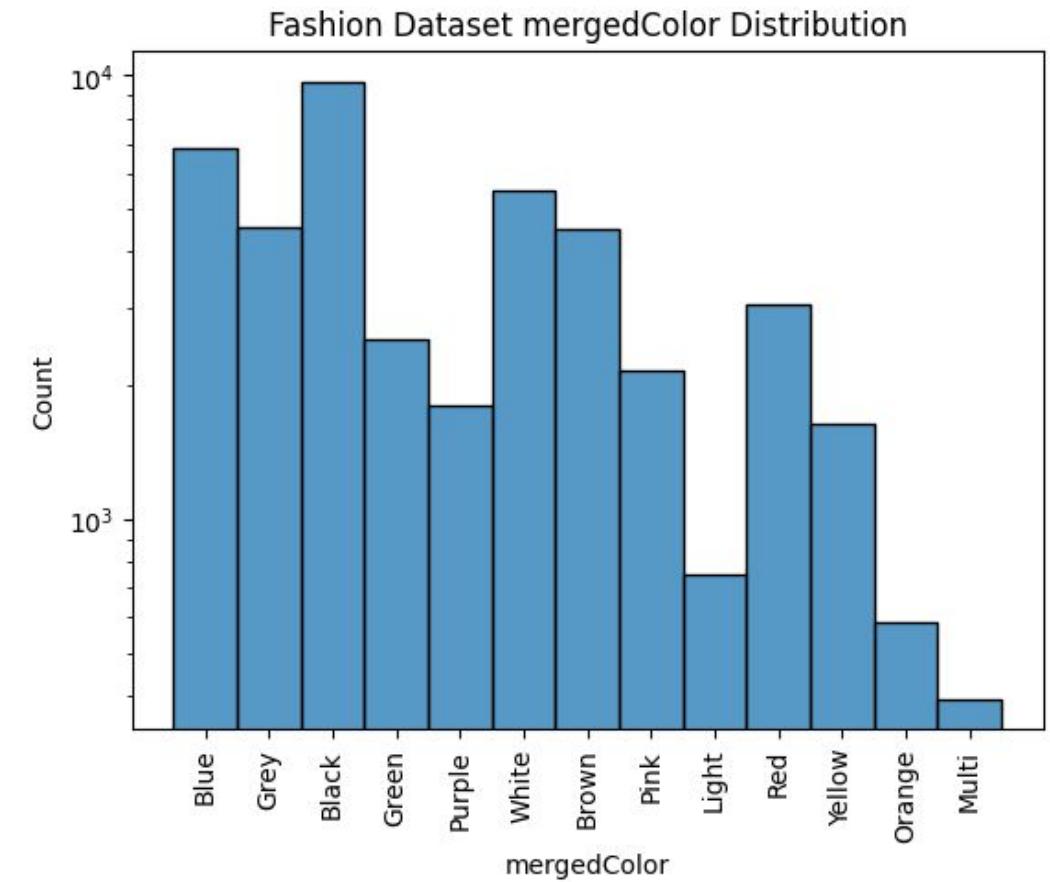
Blue	Navy Blue	Turquoise Blue	Teal	Blue
Green	Lime Green	Fluorescent Green	Sea Green	Olive
Grey	Charcoal	Grey Melange	Steel	Metallic
Light	Off White	Cream	Nude	Skin
Brown	Beige	Tan	Taupe	Mushroom Brown
Red	Maroon	Burgundy	Red	
Pink	Peach	Rose	Magenta	Pink
Purple	Lavender	Mauve	Purple	
Yellow	Mustard	Gold	Khaki	Yellow
Orange	Rust	Orange		
Black	Black			
White	White			
Multi	Multi			

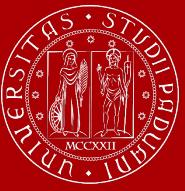


Merged: Color

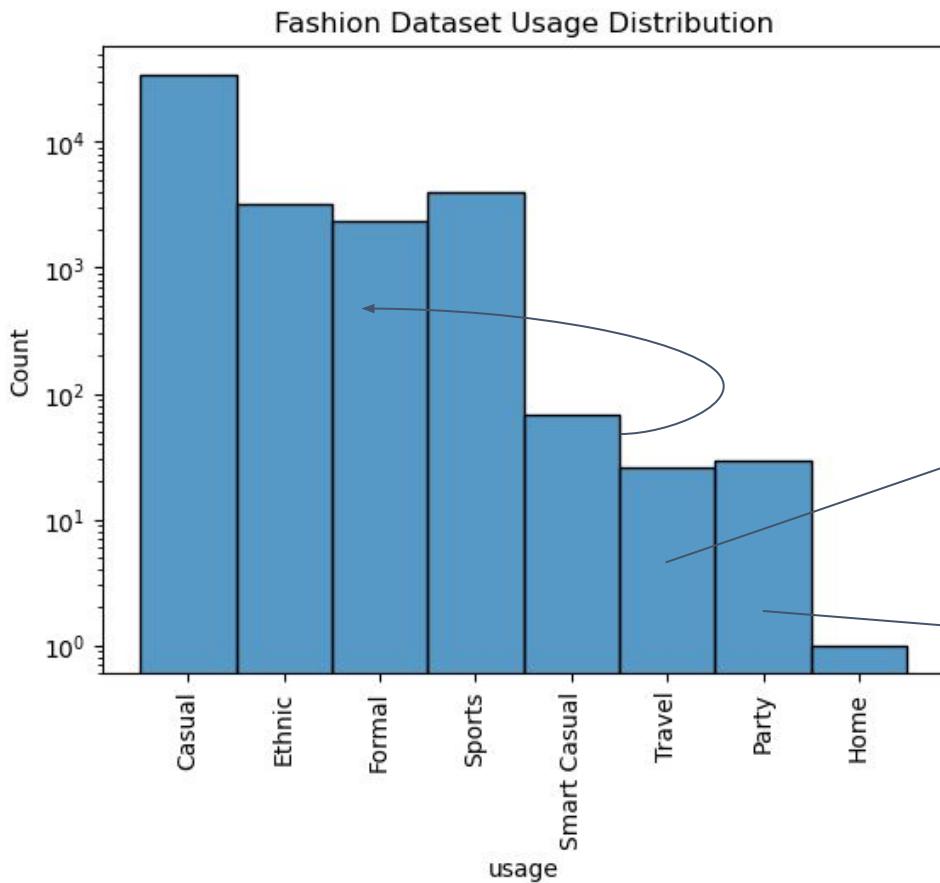


after merging
equivalent colors



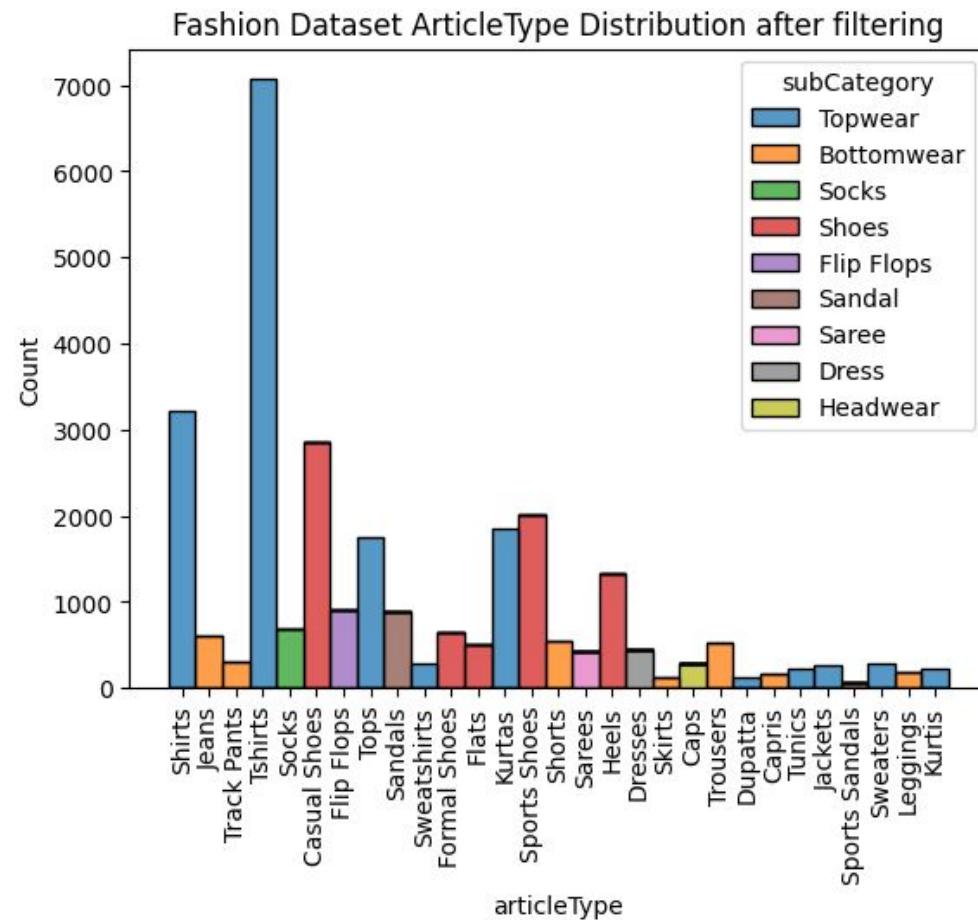


Cleaned: Usage





Matching articleType to the Image



Merged into:

- Topwear
- Bottomwear
- Footwear
- Dress
- Headwear



Matching articleType to the Image

Merged subCategory	Detector applied
'Headwear'	Trained on cap-dataset
'Footwear'	Trained on shoes-dataset
'Topwear', 'Bottomwear', 'Dress'	Trained on df2-dataset



Matching articleType to the Image

ID: 27224
Dress



ID: 42223
Topwear



ID: 7791
Topwear



ID: 42418
Topwear



ID: 3440
Topwear



ID: 15038
Topwear



ID: 33967
Topwear



ID: 11096
Topwear



ID: 31063
Topwear



ID: 11031

ID: 59596

ID: 25843

ID: 31244

ID: 3446
Topwear



ID: 31167

Cleaned out
mismatched labels



Matching articleType to the Image

The curse of multiple detections

Image ID: 26114



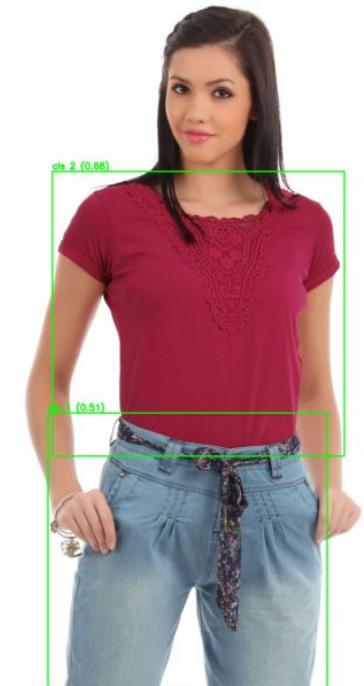
Image ID: 29182



Image ID: 9952



Image ID: 57174





Matching articleType to the Image: making predicted label and subCategory match

ID: 6863 | subCategory: Bottom



ID: 6863 | subCategory: Bottom

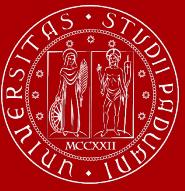


ID: 6982 | subCategory: Bottom

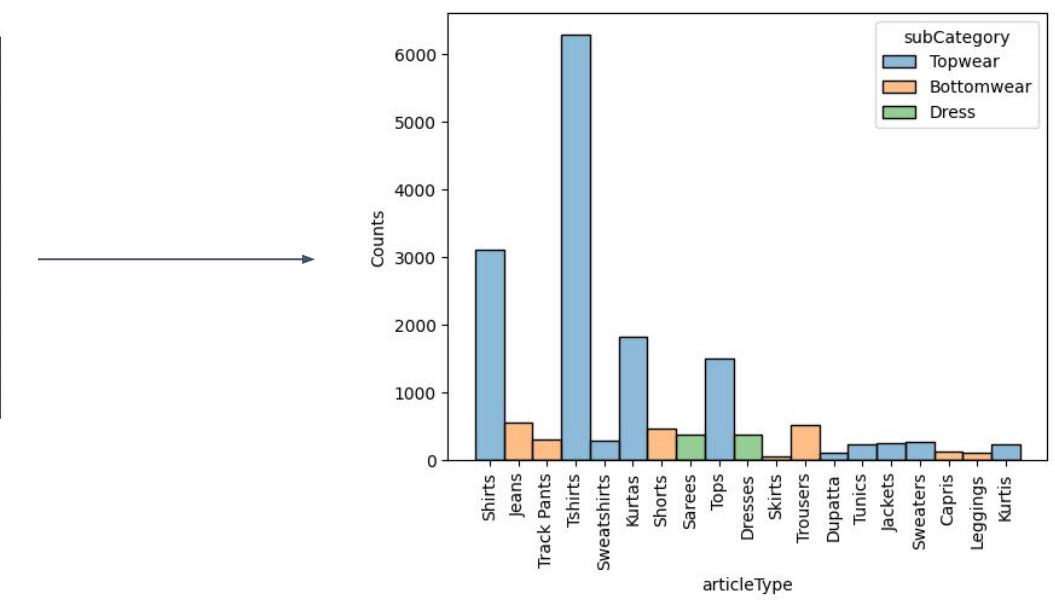
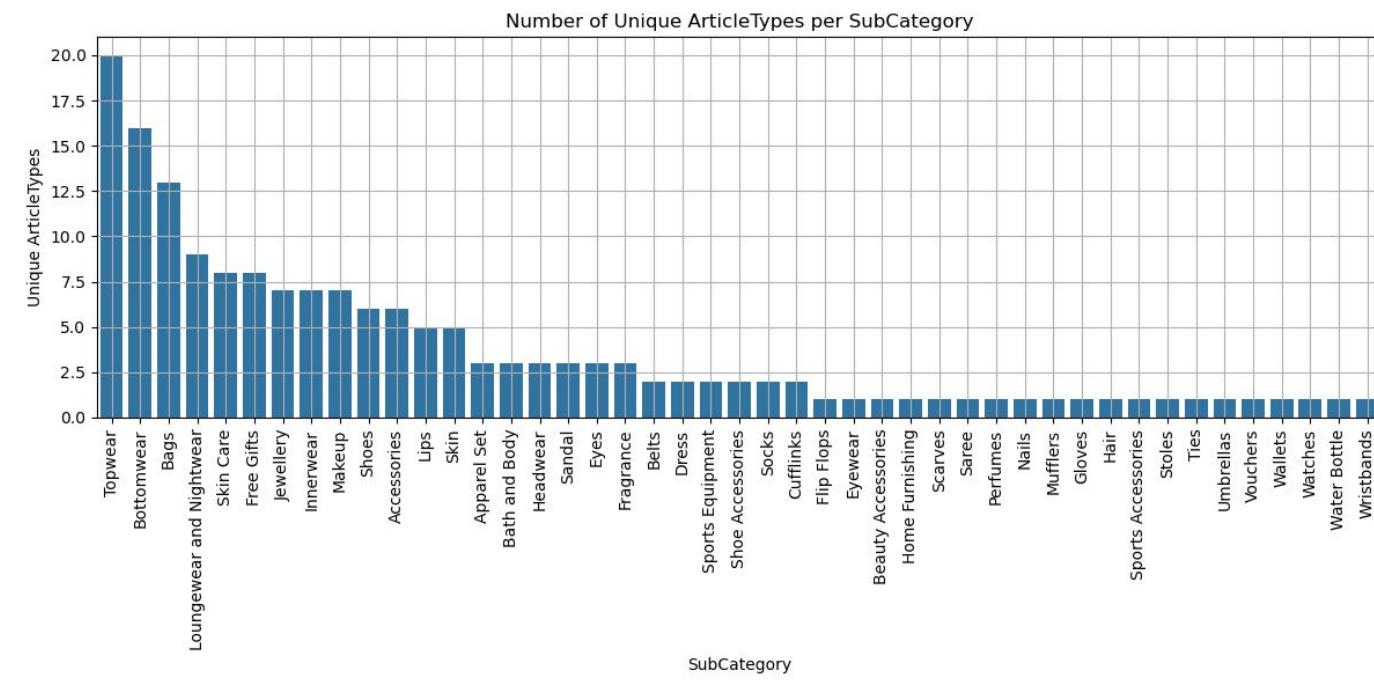


ID: 6982 | subCategory: Bottom



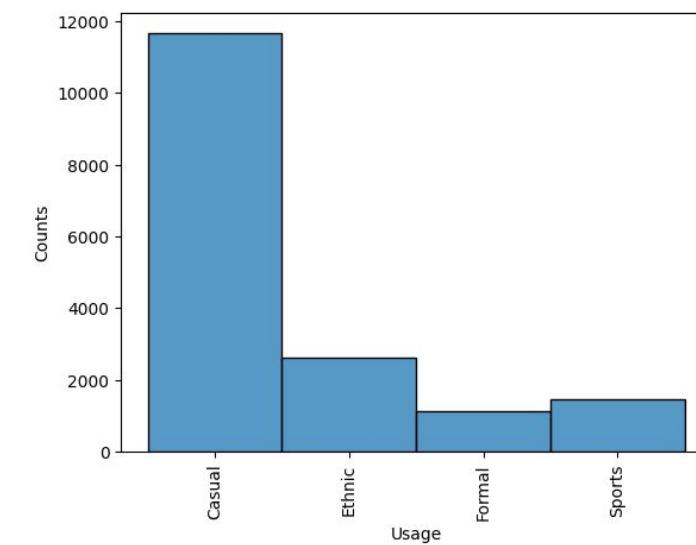
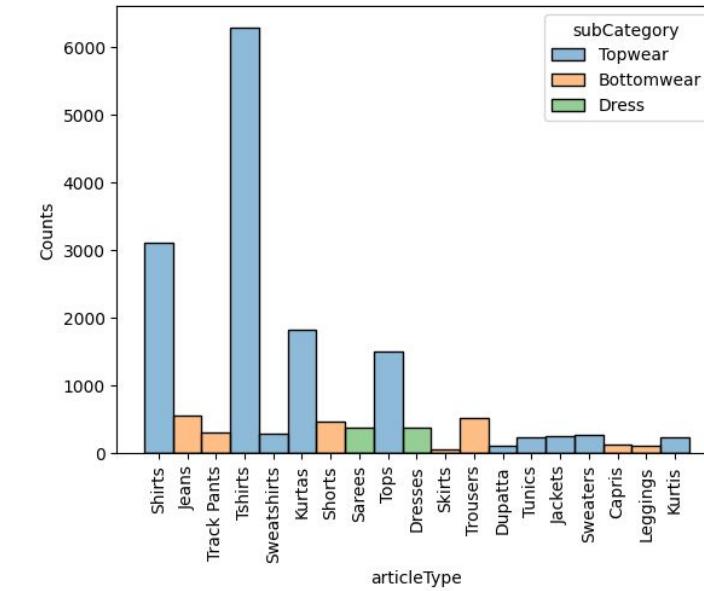
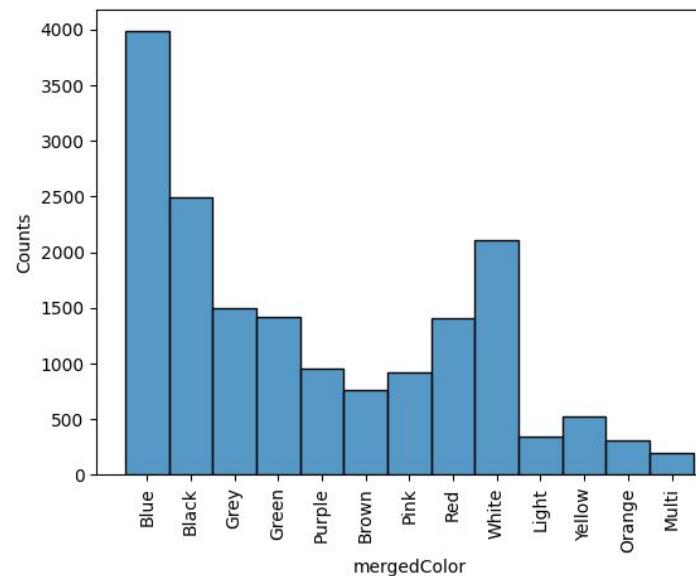
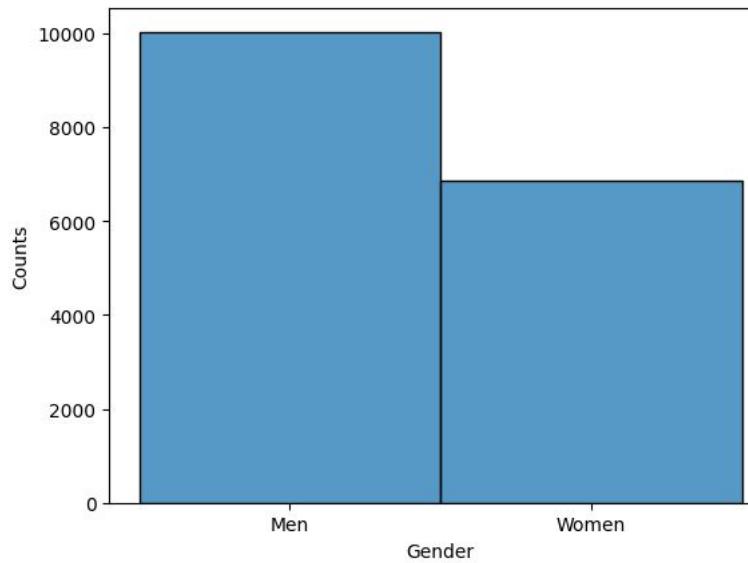


articleType: merged





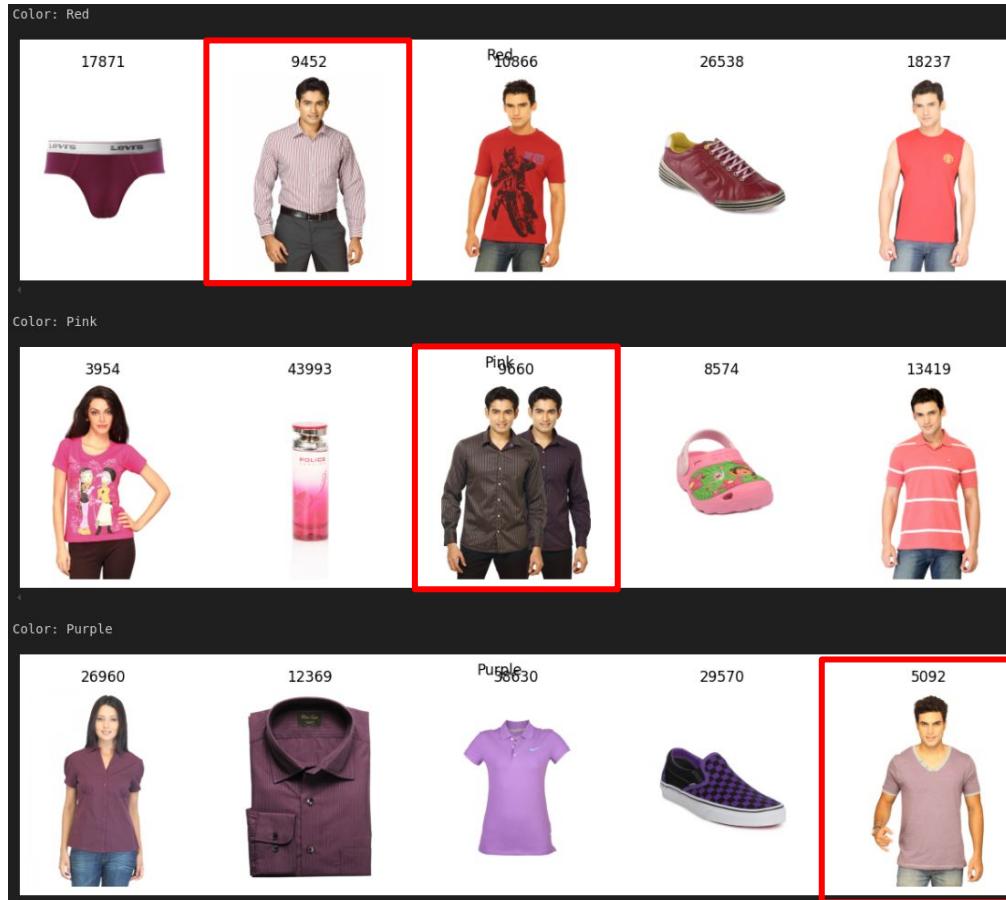
Statistics on the cleaned dataset (~18000) images





Weaknesses of FP Dataset (an example)

Red



Pink

Purple

ID: 42418
Topwear



Topwear



ID: 2179
Topwear



Labels inconsistency and mismatches



DeepFashion dataset structure

Images for id: 00000062 (n1 changes)

gender: WOMEN
type: Sweaters
id: 00000062
n1: 01
n2: 1
pos: front



gender: WOMEN
type: Sweaters
id: 00000062
n1: 01
n2: 2
pos: side



gender: WOMEN
type: Sweaters
id: 00000062
n1: 02
n2: 3
pos: back



gender: WOMEN
type: Sweaters
id: 00000062
n1: 01
n2: 7
pos: additional



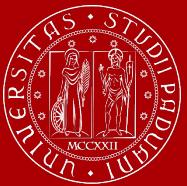
gender: WOMEN
type: Sweaters
id: 00000062
n1: 02
n2: 1
pos: front



gender: WOMEN
type: Sweaters
id: 00000062
n1: 02
n2: 7
pos: additional



naming convention:
{gender}-{clothing_item}-id_{id}-{n1}_{n2}_{pos}.jpg



DeepFashion dataset structure

Images with id=00000065, n1=02, n2=1, different types

type: Blouses_Shirts
pos: front



type: Pants
pos: front



- id: the model
- id + n1: outfit (of the model)
- id+n1+n2: model, outfit, position
- id+n1+n2+clothing_item: single item in the outfit



YOLO detector applied: same strategy to associate type to the correct bounding box

ID: 65, n1: 2, Type: Blouses_Shirts, Labels: 2



ID: 181, n1: 3, Type: Blouses_Shirts, Labels: 0



ID: 181, n1: 3, Type: Jackets_Coats, Labels: 0



ID: 181, n1: 3, Type: Jackets_Coats, Labels: 3





YOLO detector applied: same strategy to associate type to the correct bounding box

ID: 291, n1: 5, Type: Cardigans, Labels: 1



ID: 533, n1: 4, Type: Sweaters, Labels: 1



Mismatches: cleaned out



Image-specific features:dropped

ID: 533, n1: 4, Type: Sweaters, Labels: 1



sleeve lenght	long
lower_clothing_length	long
socks	NA
hat	NA
glasses	NA
neckwear	yes
wrist_weraing	no
ring	no
waist_accessories	yes
neckline	round
outer_clothing_a_cardigan	yes
upper_clothing_covering_navel	yes



The challenge of fabric and patterns

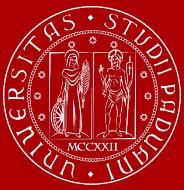
ID: 533, n1: 4, Type: Sweaters, Labels: 1



upper_fabric	lower_fabric	outer_fabric
cotton	denim	cotton

upper_pattern	lower_pattern	outer_pattern
pure_color	other	other

Image-related features: how to associate the label to the correct detected item?
a.k.a.: what is upper, lower and outer?



Priority types retrieval

Bottom-types: easy

shorts
skirts
pants
leggings
denim

Outer-types: guess

jackets_coats
cardigans
jackets_vests

Upper-types: guess

Dresses
Tees_Tanks
Sweaters
Suiting
Blouses_Shirts

:

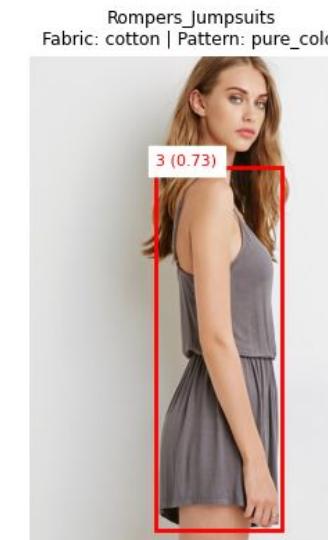
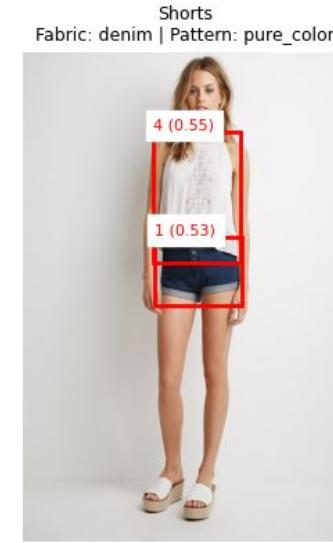


The policy applied

- If, for an image, the match is clear (one outer_type and one upper_type), retrieving the fabric/pattern label is clear.
- Otherwise:
 - Check for missing labels
 - Same label (no confusion)
 - Otherwise:
 - Image (and detected items): DROPPED

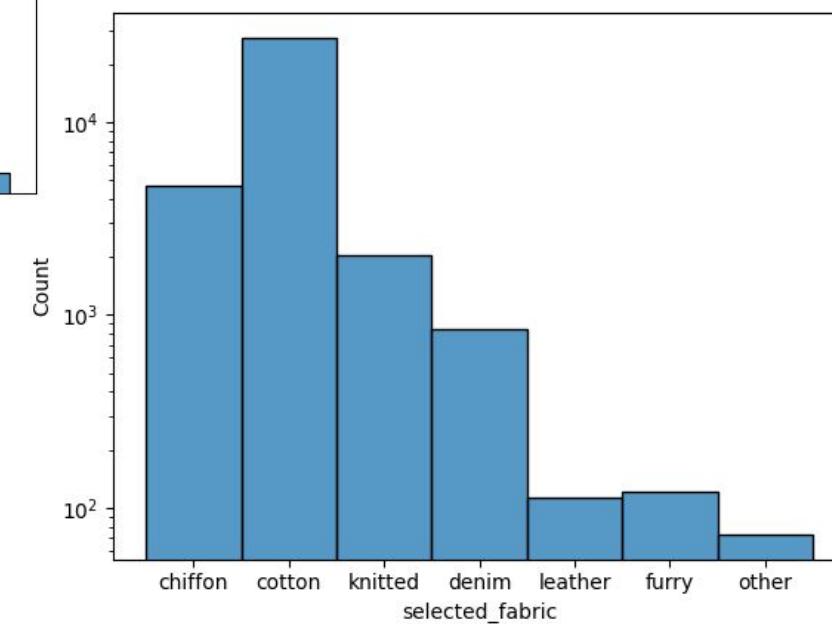
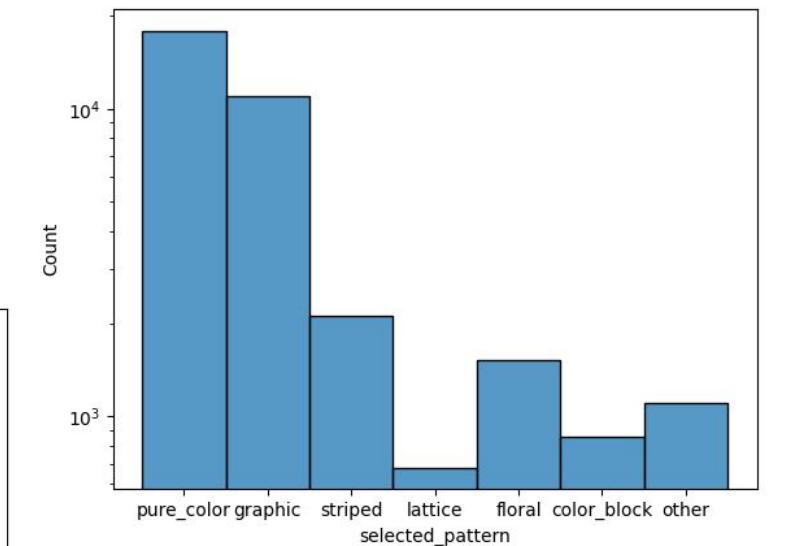
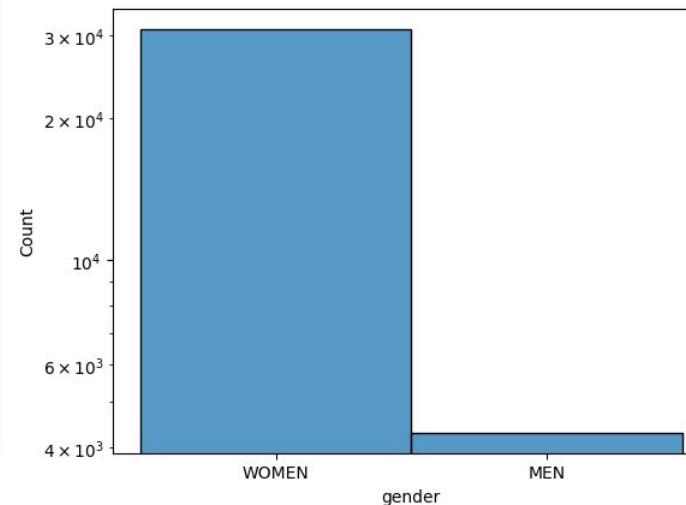
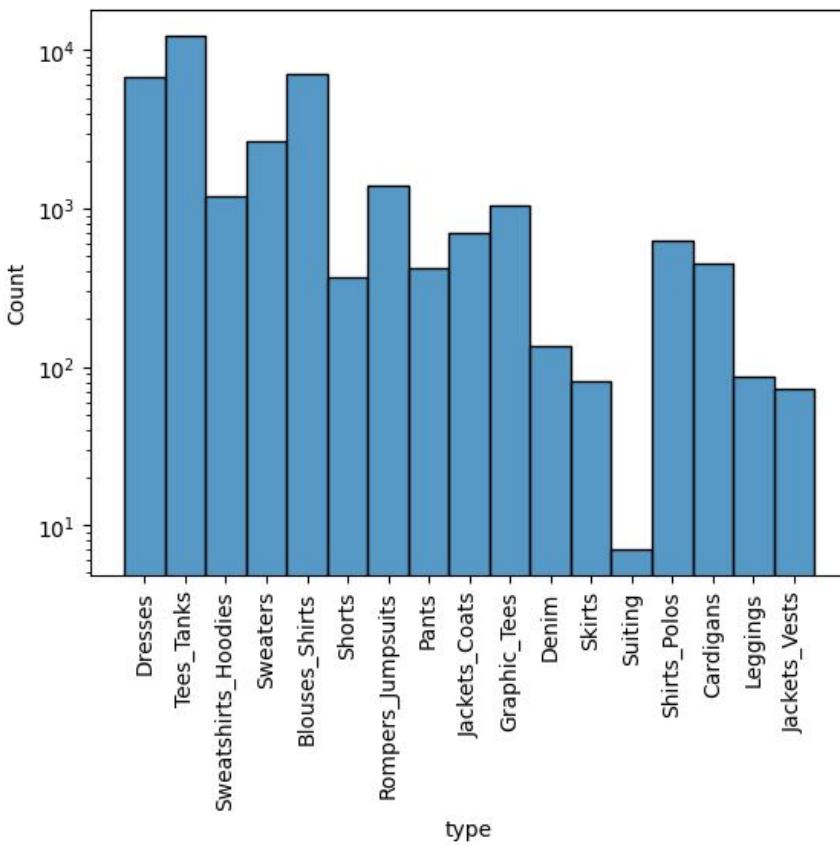


Some examples





Final statistics on DeepFashion dataset





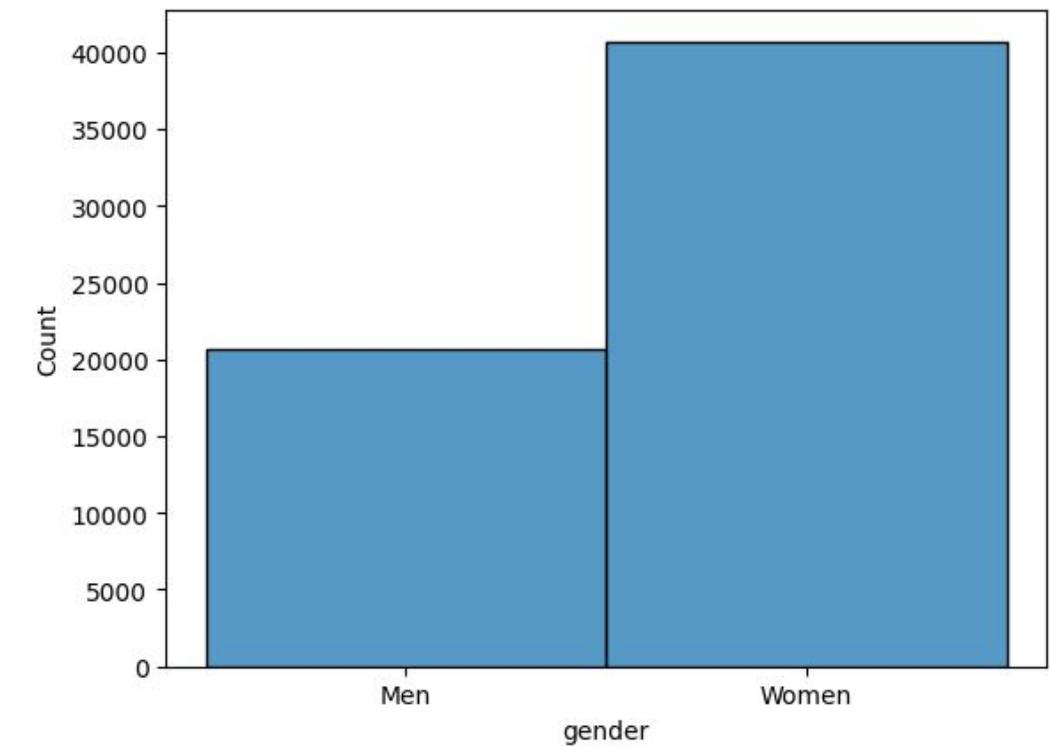
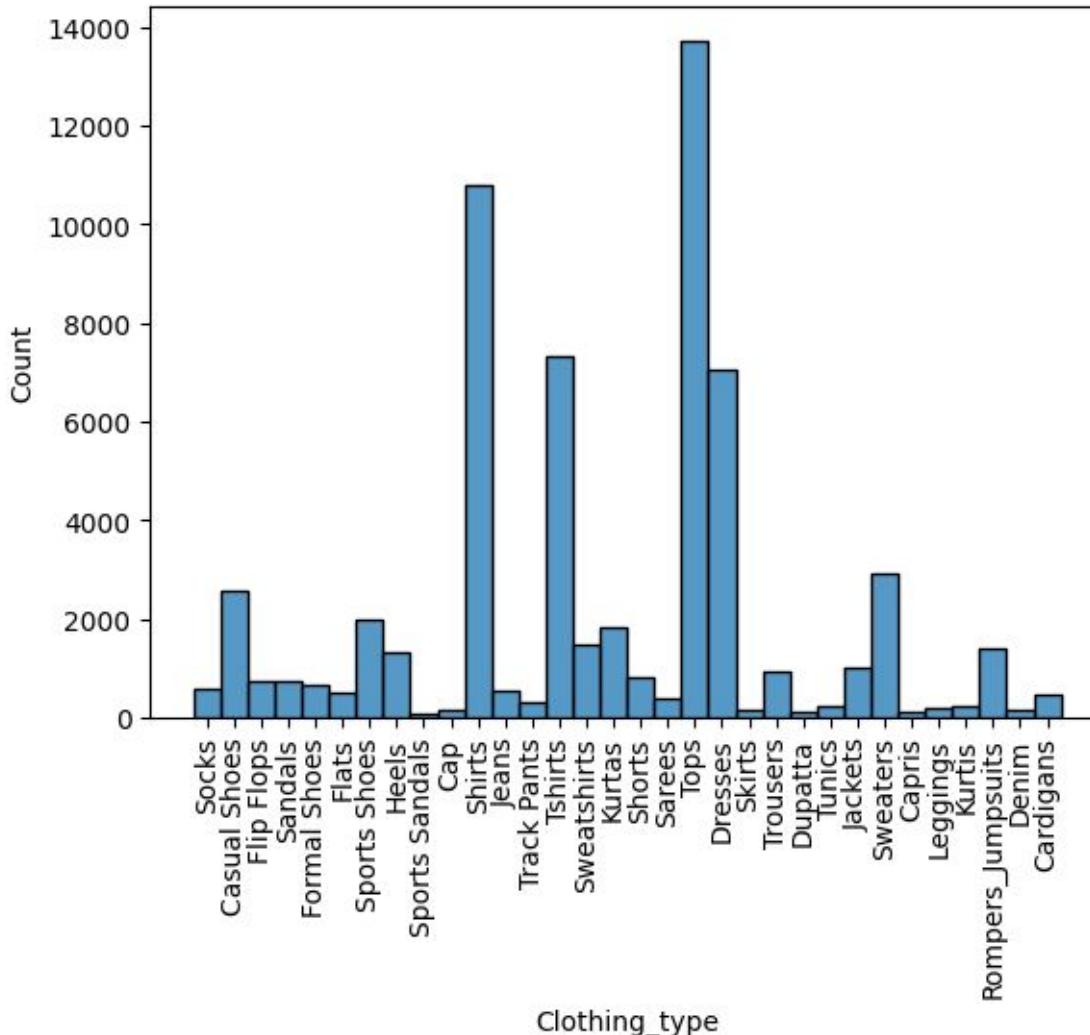
Finally

Feature	Source Dataset
Gender	FP & DF
Clothing Type	FP & DF
Usage	FP
Color	FP
Fabric	DF
Pattern	DF

Each cropped image were assigned with its corresponding labels.



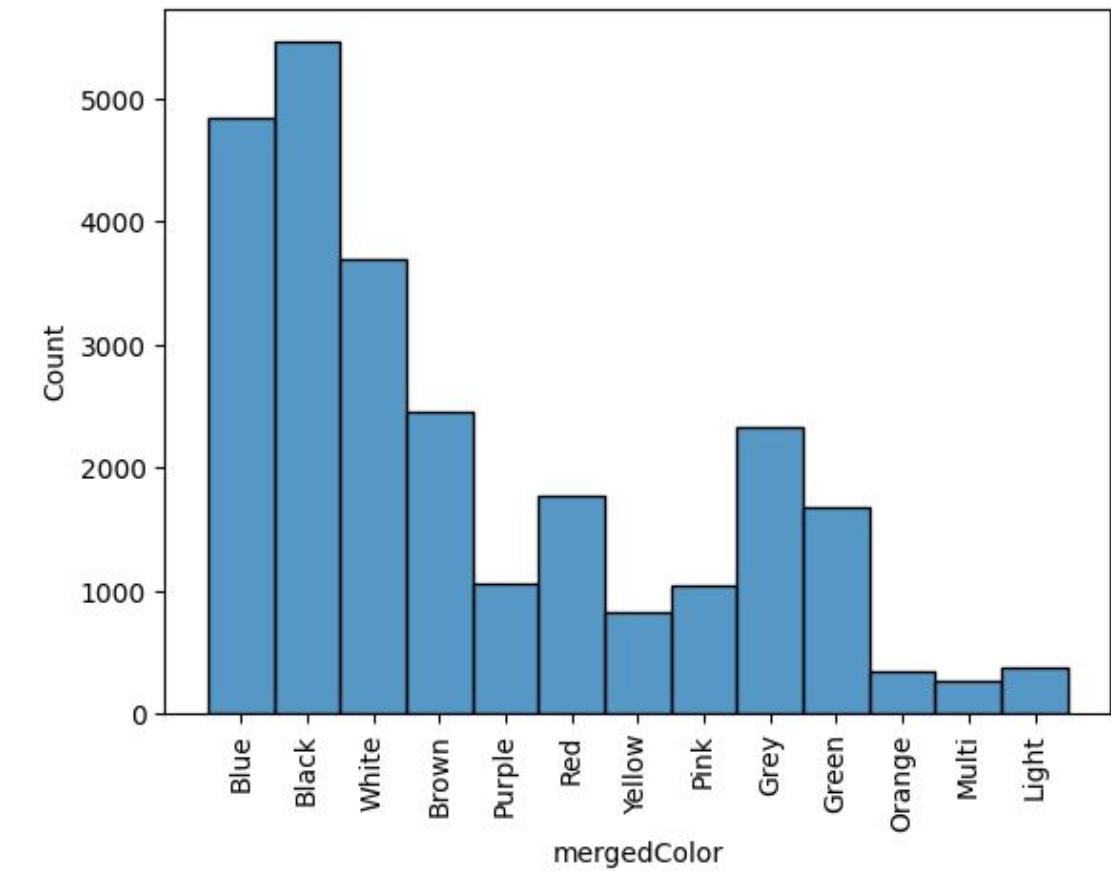
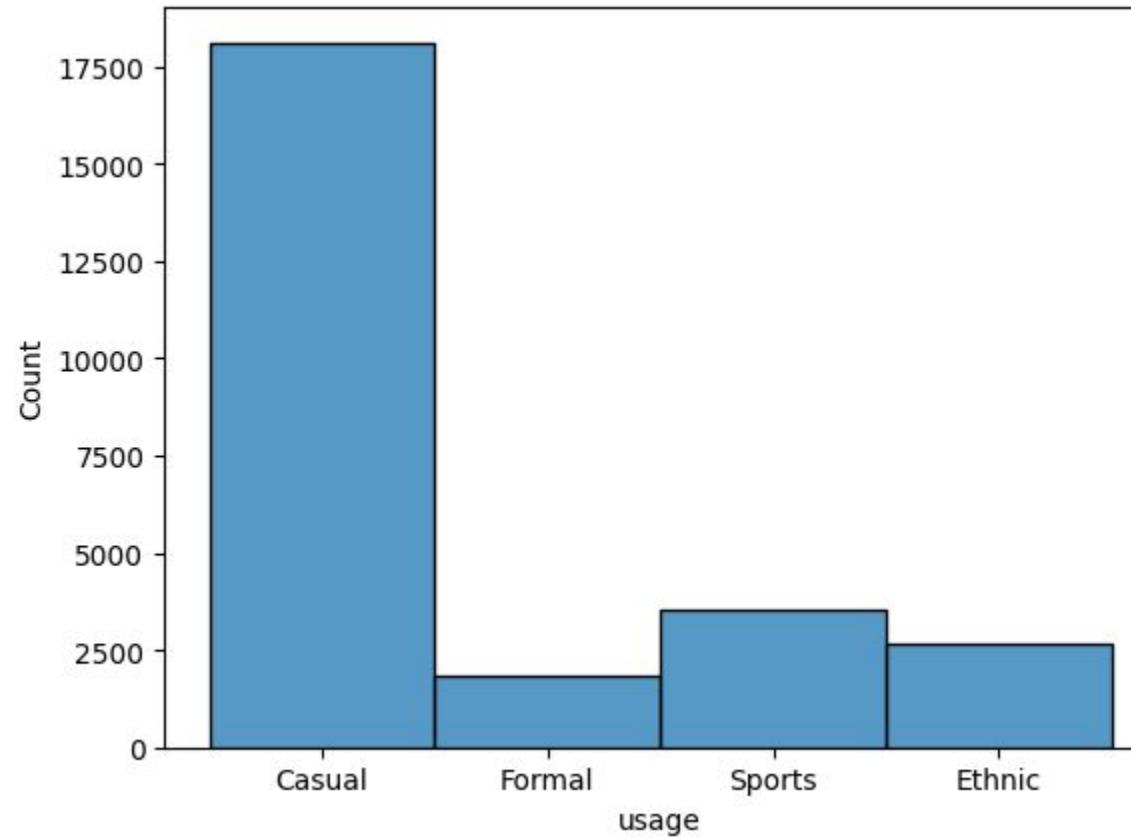
Final statistics the merged Dataset

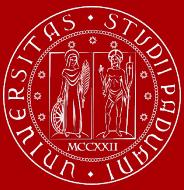




Final statistics the merged Dataset

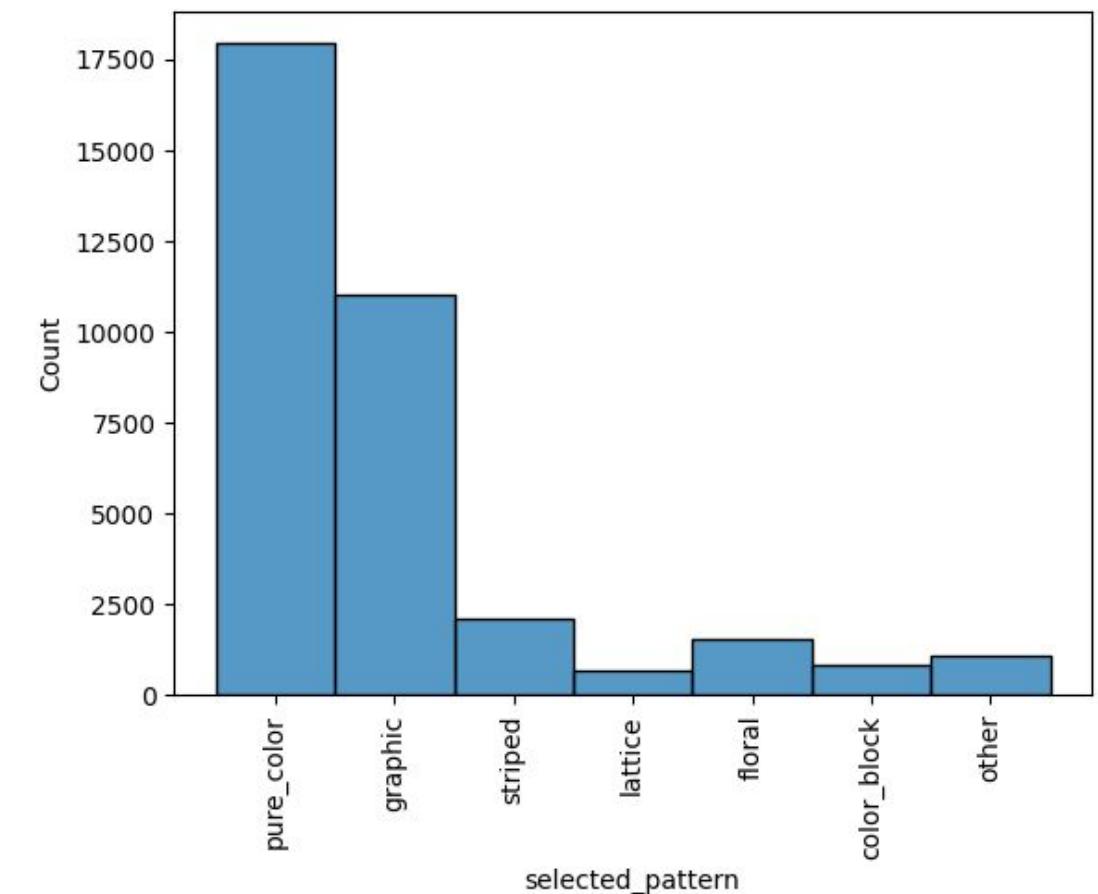
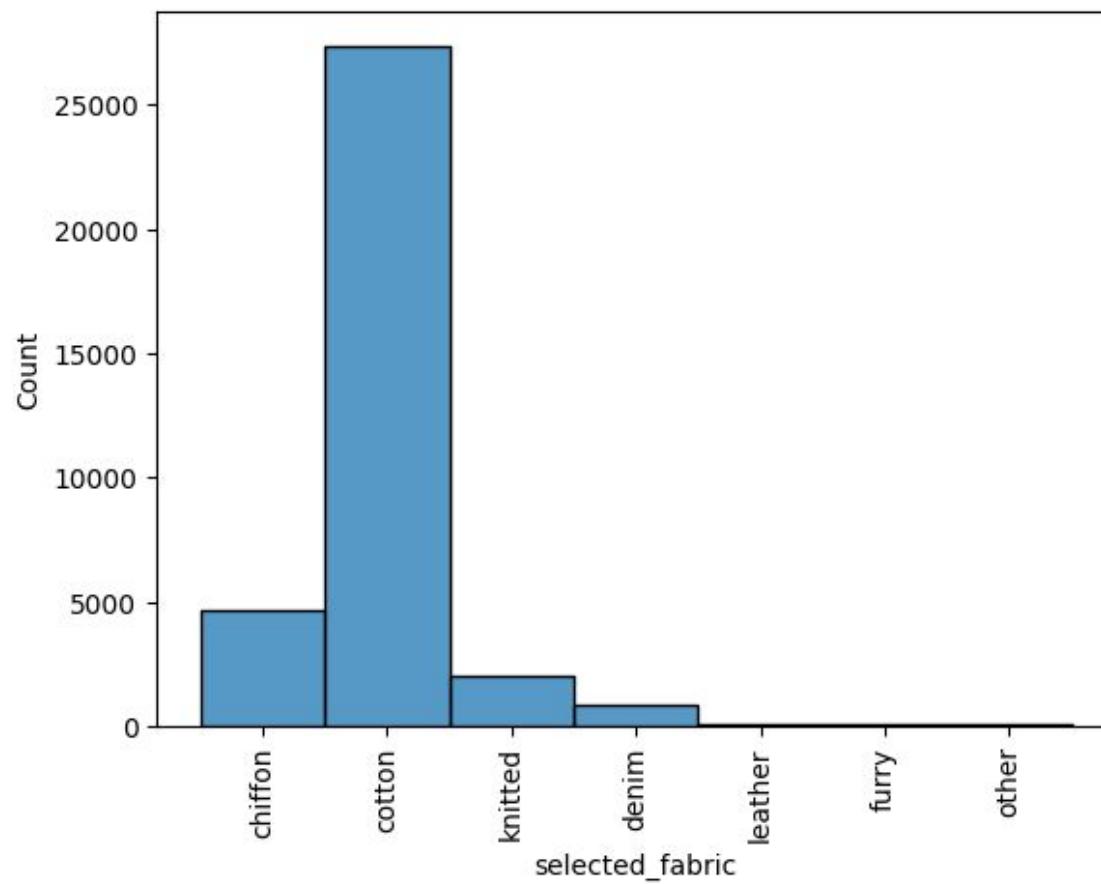
Only from Fashion Product





Final statistics the merged Dataset

Only from DeepFashion

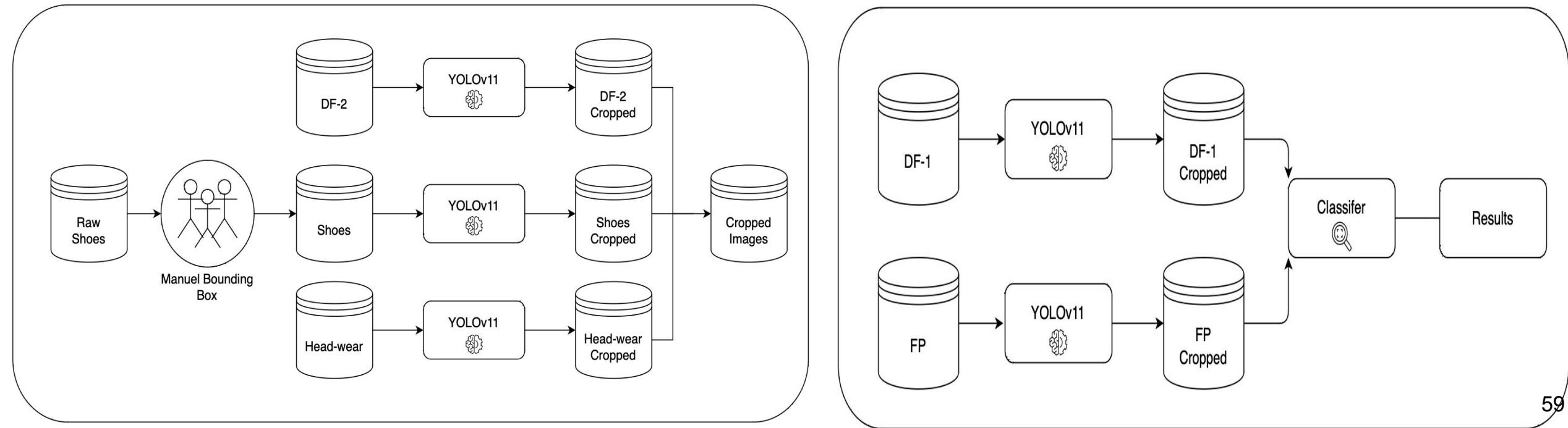




Important remark: Footwear and Headwear classes belong only to FP Dataset. Therefore, we don't have any information on the fabric and the pattern associated with these kind of items. It's all up to the generalizing power of the Classifier.

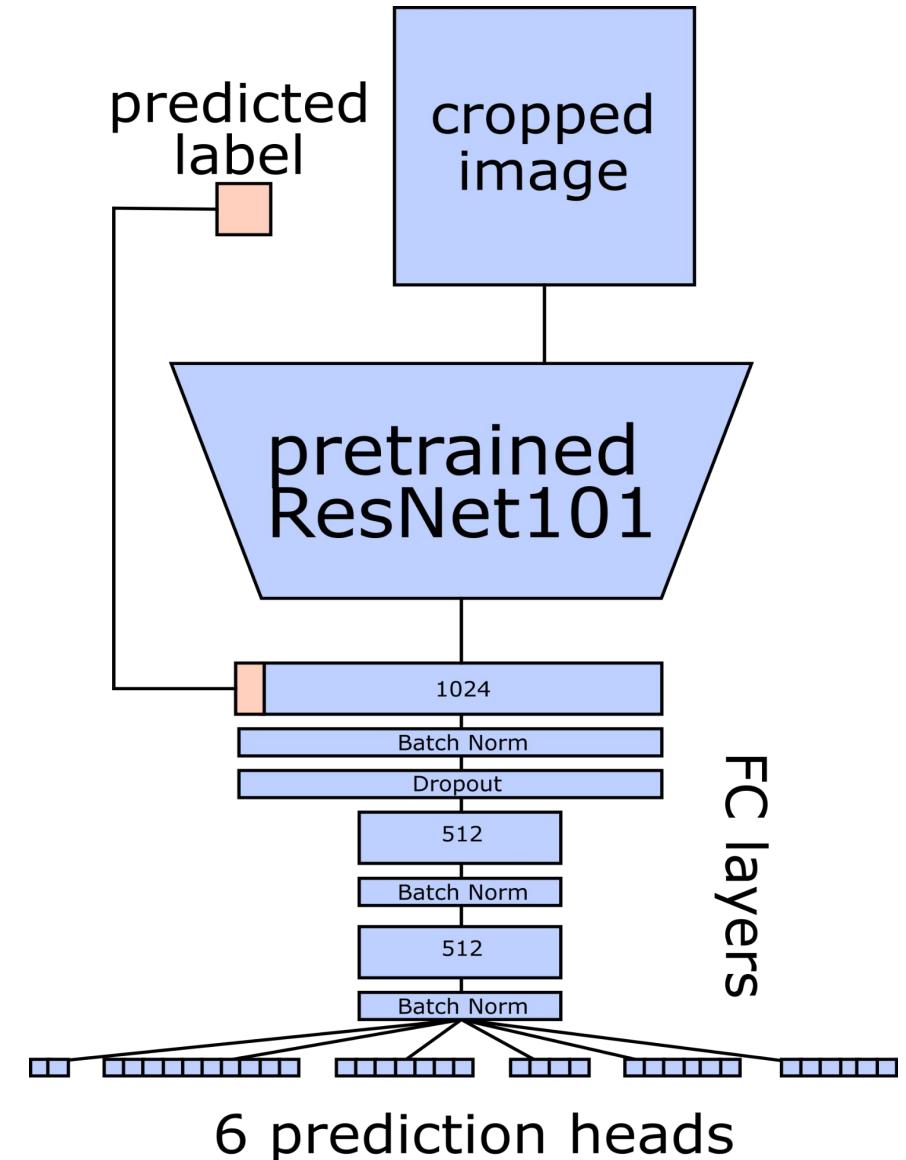


Multi-Head Classifier Model-A





Component	Description
Backbone	ResNet101 pretrained with self-supervised learning
Input	Cropped images + predicted label from detectors
FC Layers	[1024 + 1 → 512 → 512] with BatchNorm, LeakyReLU, Dropout, where +1 is the predicted label value from the applied YOLO detectors
Task Heads	6 linear heads for: gender (2), type (32), usage (4), color (13), fabric (7), pattern (7)
Loss	Focal Loss + Masked Cross Entropy with class weights
Augmentation	Blur, brightness/contrast shift, JPEG compression
Training Strategy	Interleaved learning based on dataset origin (FP or DF)
Optimizer	Adam: LR = $1e-4$ (heads/FC), $5e-5$ (backbone)
Epochs	20
Logging	W&B + per-task confusion matrices





Data Augmentation

Technique	Probability	Parameter Range	Description
Gaussian Blur	30%	Radius: 0.5 – 1.5	Applies a mild blur to simulate focus variation or motion blur.
Brightness Adjustment	10%	Factor: 0.8 – 1.2	Randomly increases or decreases image brightness.
Contrast Adjustment	10%	Factor: 0.8 – 1.2	Randomly adjusts the contrast of the image.
JPEG Compression	30%	Quality: 40 – 70	Simulates lossy compression artifacts by saving and reloading the image.
Resize	100%	Output Size: (256, 256)	Resizes the final output to the specified dimensions for model compatibility.



Multi-Head Classifier (Model-A)

The Final model

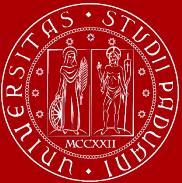
- 1) ResNet101 backbone, pre-trained with self-supervised learning
- 2) A unique feature of our approach is that we concatenate the predicted clothing type label (obtained from the detector) as an additional input alongside image features. This adds semantic info and enforces dependence on detection model.
- 3) The model includes a fully connected multi-layer block followed by independent heads, one for each semantic label (e.g., gender, fabric, usage)



Multi-Head Classifier (Model-A)

The Final model

- 4) For each item, only the relevant heads are updated based on the source dataset (DF1 or FP).
- 5) Interleaved learning applied on Data Loaders to prevent forgetting of information learned through the weights.
- 6) To handle strong class imbalance, we combine masked cross-entropy and Focal Loss, and apply class-specific weighting.



What is focal loss?

Focal Loss is a loss function designed to address class imbalance during training especially in classification problems where one class significantly outnumbers others.



Challenges in Multi-Head Classifier (Model-A)

- Validation accuracy for the Clothing Type class
 - (We enriched the input by combining the cropped item image with its predicted category label from the YOLO detector.)
- This injects prior knowledge derived from outfit detection phase
- Classes (Color, Fabric, and Pattern)
- Pre-trained backbones such as ResNet18 and ConvNeXt, struggled to extract sufficiently rich representations for these appearance-centric attributes.
- To overcome this limitation, we trained a ResNet101 encoder from scratch in a self-supervised manner.



Encoder

- Uses ResNet101
- SimCLR (Simple Contrastive Learning of Representations)

Component	Value
Backbone	ResNet101 (ImageNet, headless)
Proj. Head	MLP: $2048 \rightarrow 2048 \rightarrow 512$
Loss	NT-Xent (contrastive, temp. scaled)
Batch Size	32
Image Size	256×256
Epochs	30
Augmentations	Flip, rotation, masking
Color Aug.	Disabled (to retain color cues)
Optimizer	SGD ($lr=2e-2$, $mom=0.9$, $wd=5e-4$)
Scheduler	Cosine annealing
Validation	KNN on frozen embeddings
Framework	PyTorch Lightning + W&B



SimCLR Pretraining

- 1) SimCLR (Simple Contrastive Learning of Representations)
- 2) It's a self-supervised learning method
- 3) It can learn meaningful image features without needing labels.



SimCLR Pretraining How is applied?

- 1) For each image, two different augmentations are created (random horizontal flip and rotation).
- 2) These are treated as a positive pair different views of the same image.
- 3) All other images in the batch are negative pairs.
- 4) The model learns to: Pull positive pairs closer in feature space.
Push negative pairs apart.



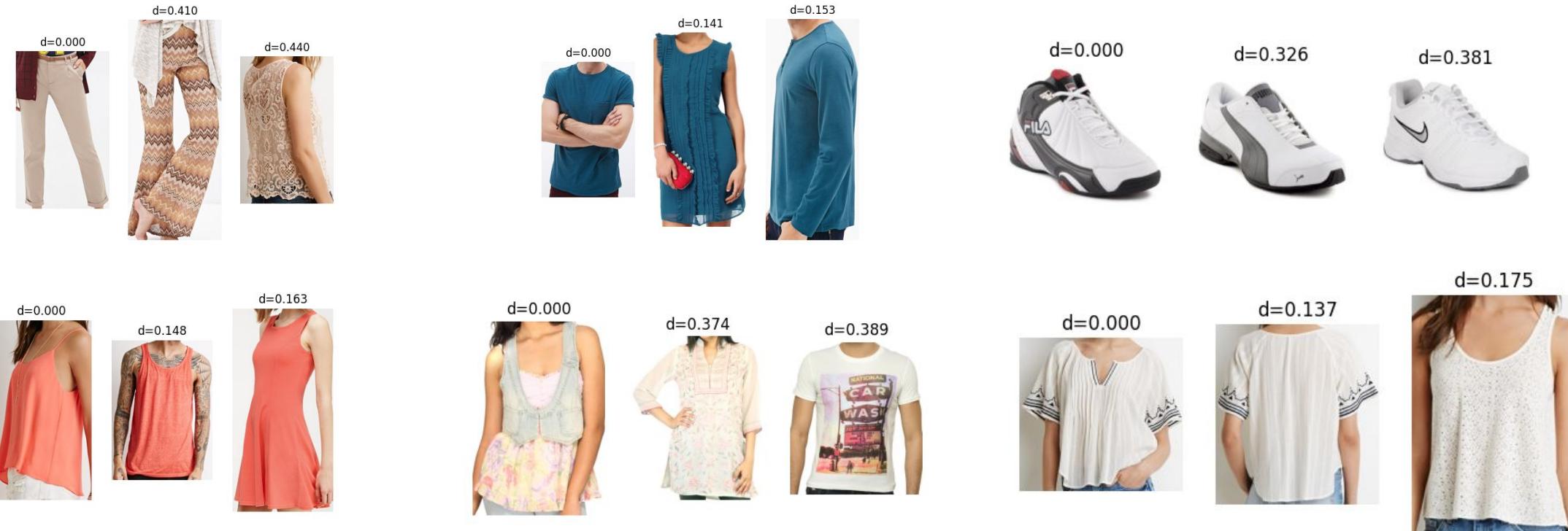
SimCLR Pretraining Method

- 5) This is done using a contrastive loss called NT-Xent (Normalized Temperature-scaled Cross Entropy).
- 6) To ensure the final representations preserved color fidelity, augmentations involving color jitter and grayscale conversion were explicitly disabled, unlike standard Sim-CLR pipelines
- 7) Once trained, the projection head is discarded, and the ResNet backbone is used to generate embeddings.



SimCLR Pretraining Method

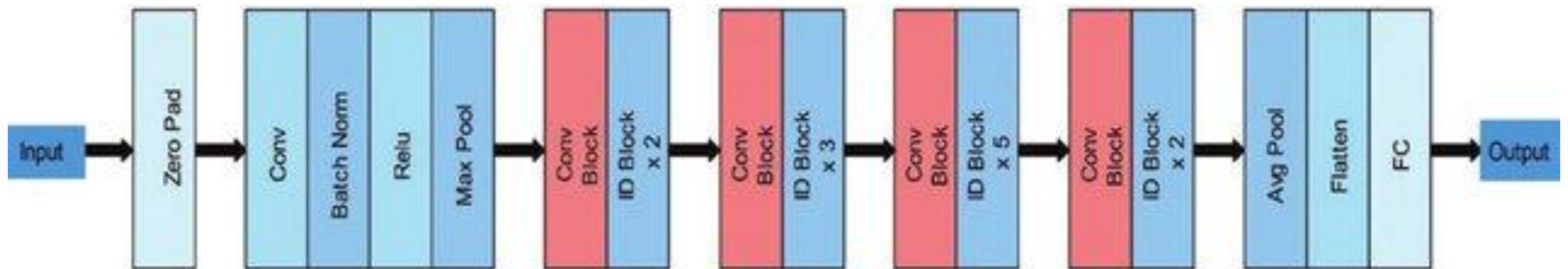
8) These were evaluated using K-Nearest Neighbors in the latent space, confirming that visually similar garments clustered closely, even without labels. Weights are saved for classification.





ResNet 101

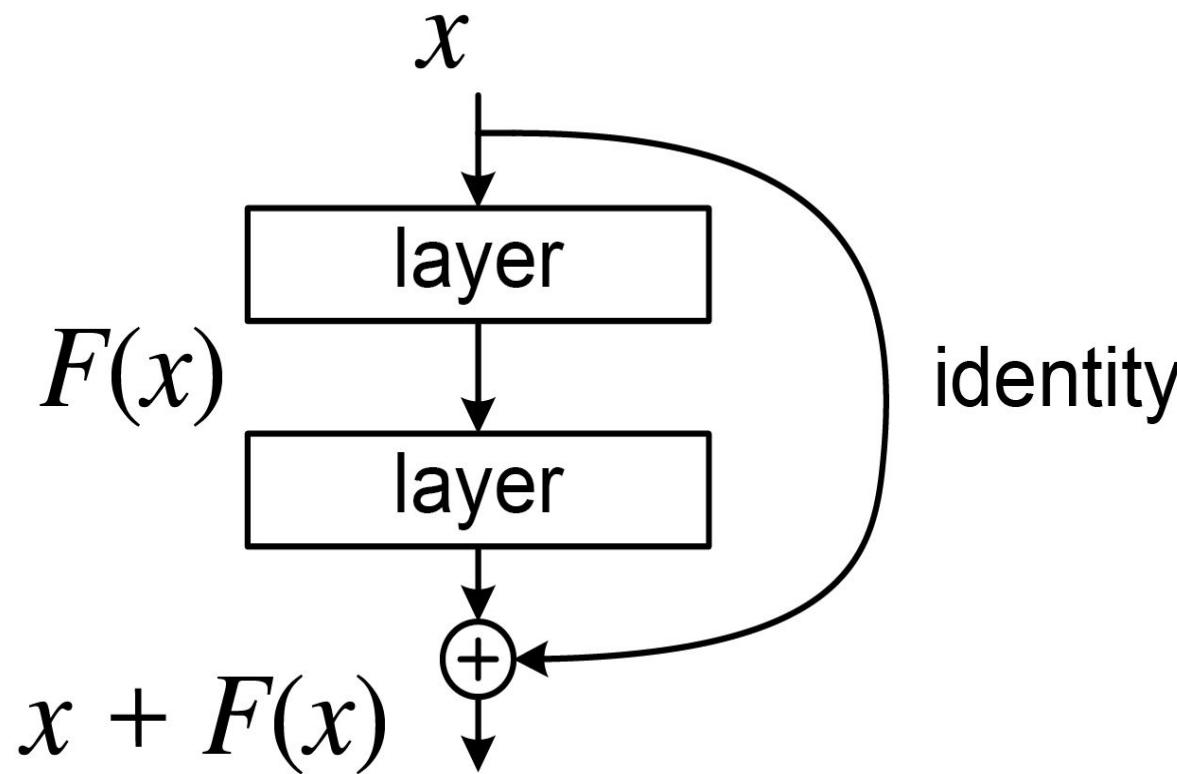
- Deep CNN
- It has residual connections or (skip connections)



Sonavane, R., Ghonge, P., & Patil, S. U. (2024). Exploring resnet101, inceptionv3, and xception for modi script character classification. *International Journal of Intelligent Systems and Applications in Engineering*, 12(17s), 117-124.



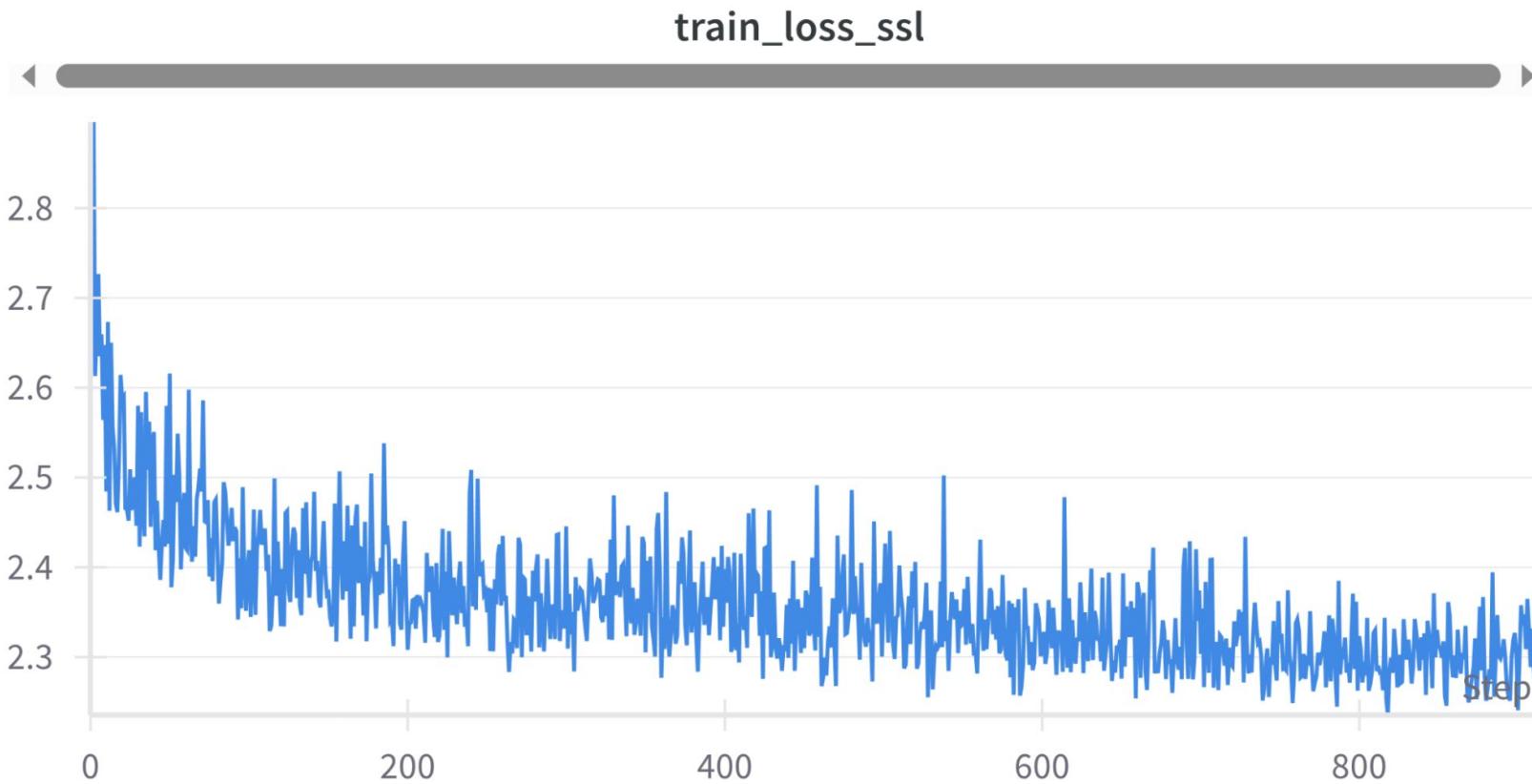
Residual Connections



- Effective for learning color, Fabric, and texture.
- They allow information (especially gradients) to flow more easily during training.
- This avoids the vanishing gradient problem
- Preserve high-frequency details. (color shifts, texture details.)



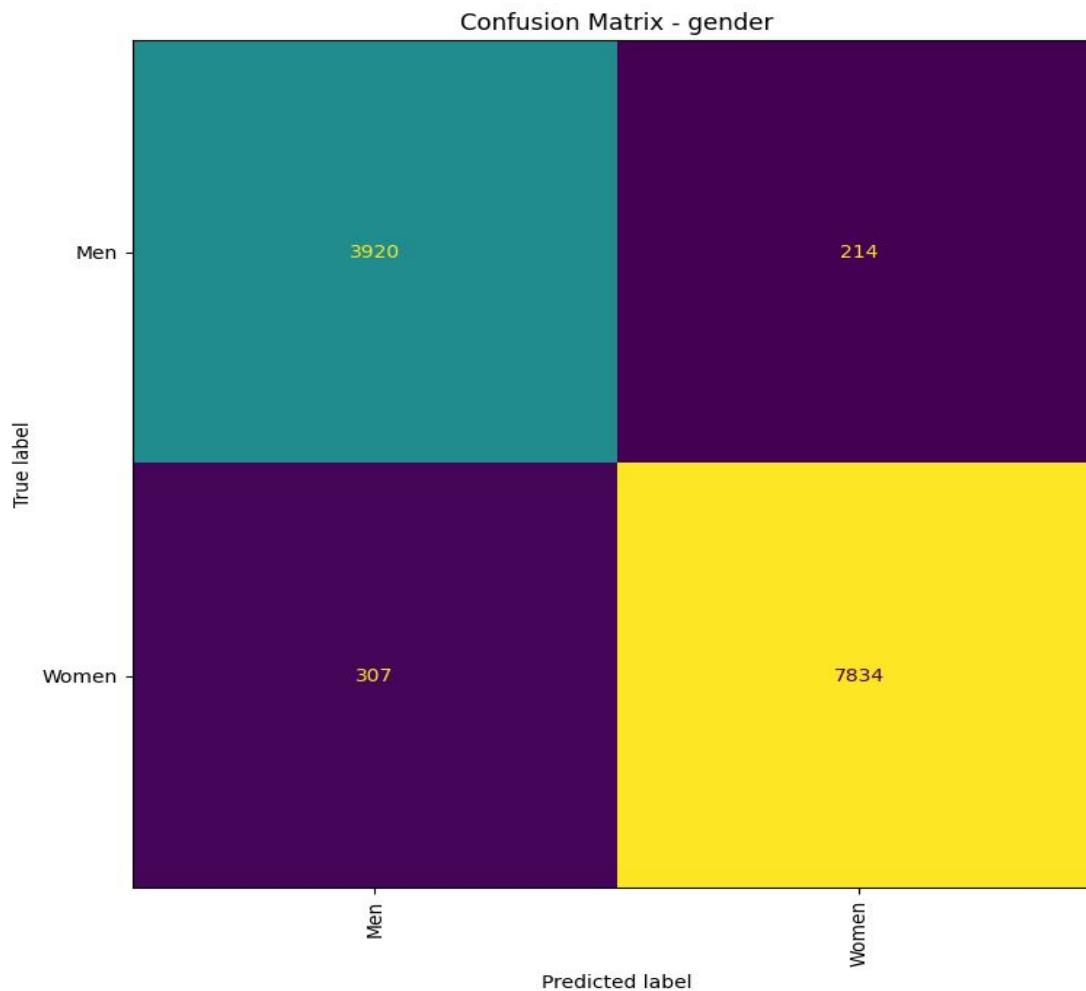
ResNet Training Loss



Loss function of the ResNet 101 architecture training in a self-supervised way using the SimCLR implementation



Confusion Matrix for Classification Model



Positive Class = Woman Negative Class = Man

True Positives (TP) = 7834

False Negatives (FN) = 307

True Negatives (TN) = 3920

False Positives (FP) = 214



Classification Metrics for Gender

Accuracy: 95.76%

Precision: 97.34%

Recall: 96.27%

F1 Score: 96.80%



True label

Confusion Matrix - mergedColor														
True label	Blue	Black	White	Brown	Purple	Red	Yellow	Pink	Grey	Green	Orange	Multi	Light	
Predicted label	Blue	Black	White	Brown	Purple	Red	Yellow	Pink	Grey	Green	Orange	Multi	Light	
Blue	21	416	14	110	26	21	15	3	336	65	5	63	15	
Black	11	2	418	6	23	21	12	15	75	12	0	17	152	
White	427	35	18	3	71	34	13	12	173	79	5	57	11	
Brown	0	7	3	236	3	5	118	4	23	13	13	19	53	
Purple	12	8	36	12	8	7	10	8	283	17	3	6	59	
Red	0	0	3	6	26	256	3	38	3	1	24	17	7	
Yellow	0	0	0	0	0	4	0	8	0	0	40	2	1	
Pink	1	0	0	4	2	2	31	1	7	233	2	12	14	
Grey	0	0	1	0	17	9	3	154	0	0	13	9	8	
Green	0	0	1	0	0	1	6	1	3	0	2	4	48	
Orange	1	1	0	6	0	1	119	0	1	3	10	2	8	
Multi	2	1	2	4	157	3	0	8	3	0	2	3	1	
Light	0	0	0	1	3	1	5	5	3	0	2	33	4	

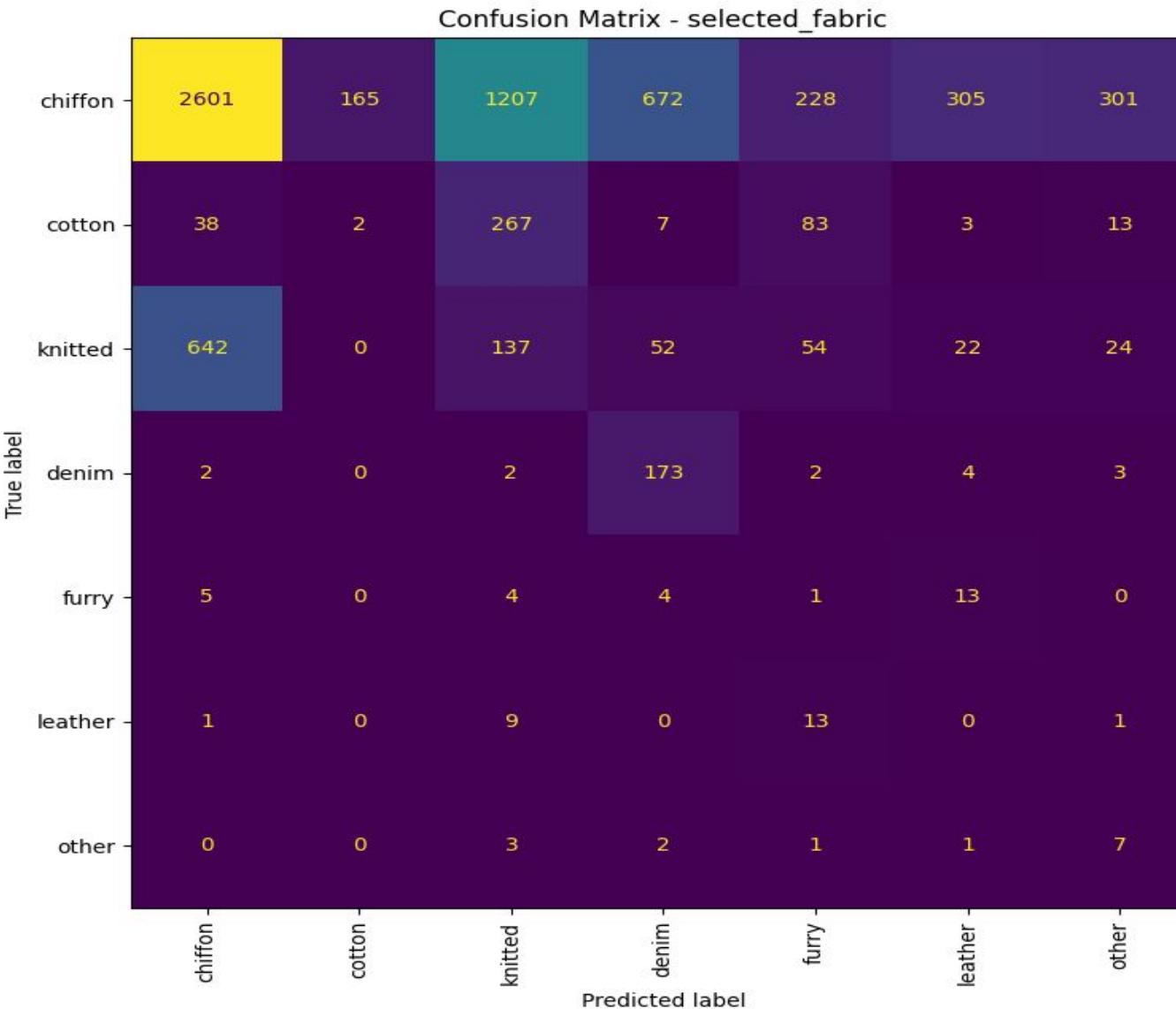
True (Actual Class)	Commonly-Misclassified
Blue	Black
White	Black
Black	White
Purple	Grey
Yellow	Orange
Pink	Green
Grey	Pink
Green	Light
Orange	Yellow
Multi	Purple
Light	Multi



Mislabelled samples



True (Actual Class)	Commonly-Misclassified
Blue	Black
White	Black
Black	White
Purple	Grey
Yellow	Orange
Pink	Green
Grey	Pink
Green	Light
Orange	Yellow
Multi	Purple
Light	Multi

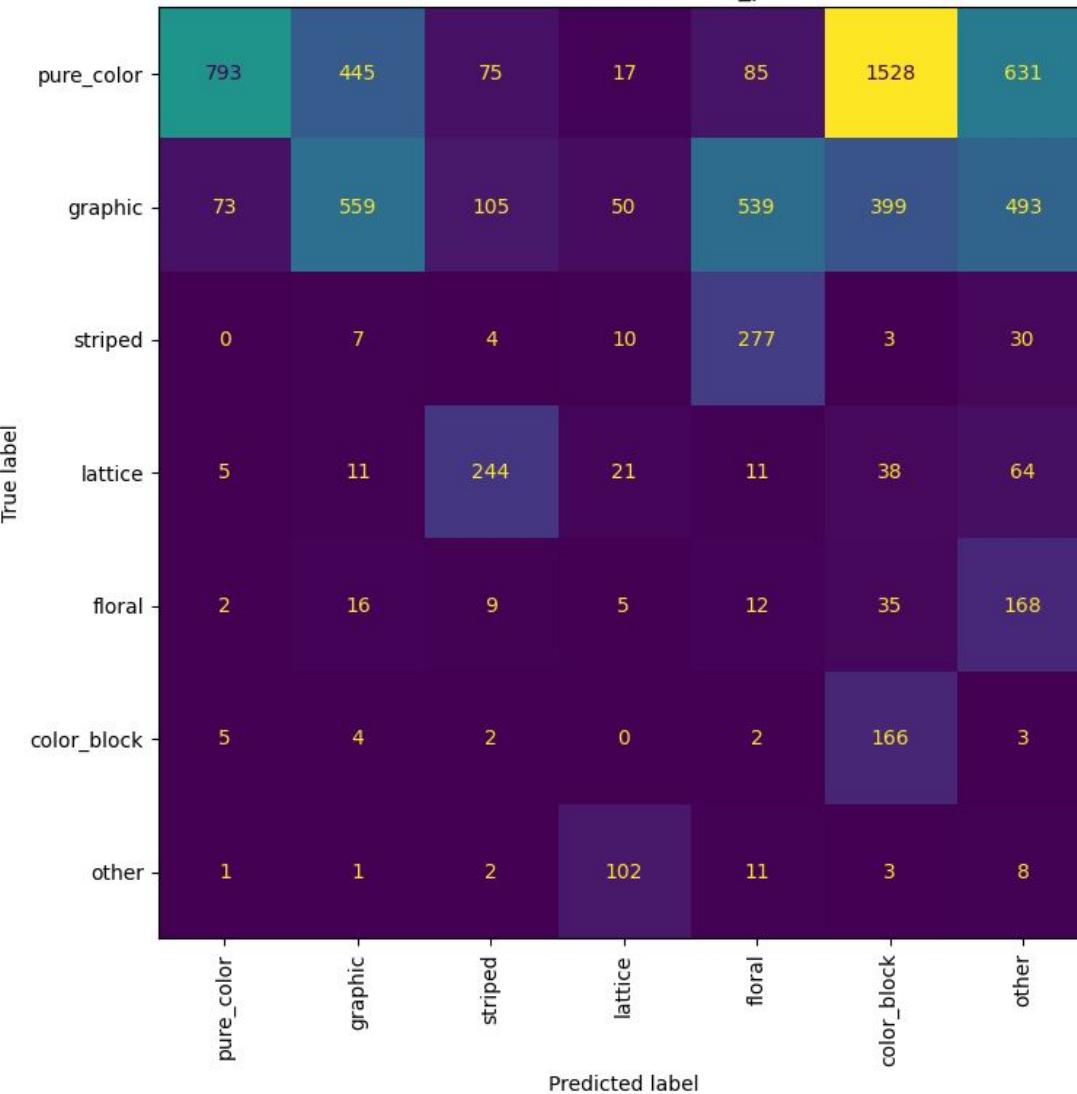


Problematic Classes

True (Actual Class)	Commonly-Misclassified
Cotton	Knitted
Knitted	Chiffon
Furry	Leather
Leather	Furry



Confusion Matrix - selected_pattern

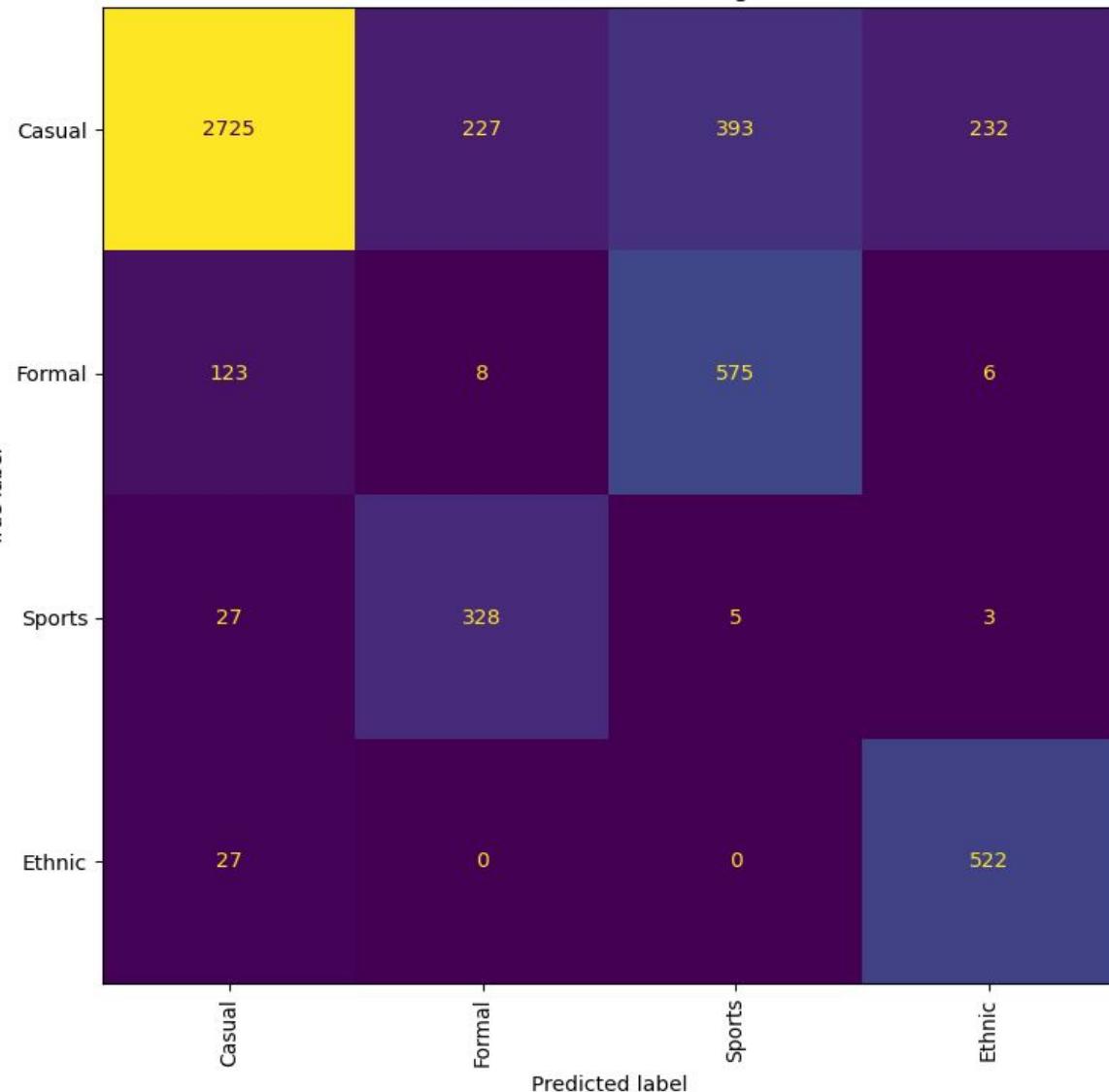


Problematic Classes

True (Actual Class)	Commonly-Misclassified
pure-color	color-block
striped	floral
lattice	striped
floral	other
other	lattice

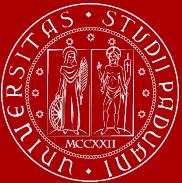


Confusion Matrix - usage



Problematic Classes

True (Actual Class)	Commonly-Misclassified
Formal	Sports
Sports	Formal



Confusion Matrix - Clothing_type

Problematic Classes

- 1) Jeans confused with other bottom wear
 - 2) Socks and Flip Flops



Confusion Matrix - Clothing_type

We have some structure in the matrix: Footwear and Clothes are unlikely to be confused.

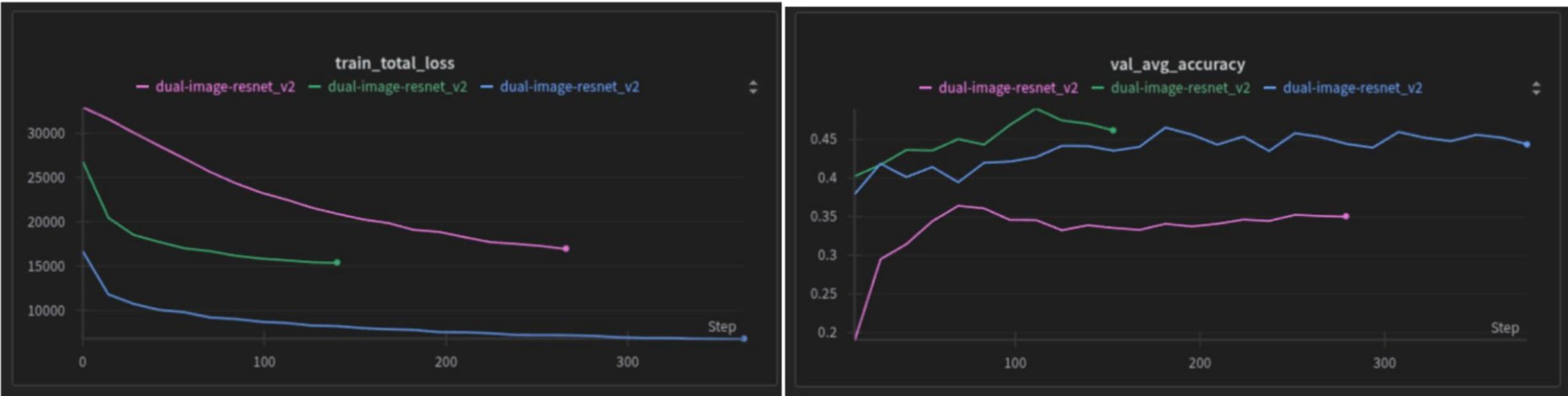


The final model performance on the validation set

	Gender	Usage	Fabric	Pattern	Color	Clothing type
Recall	0.95	0.42	0.26	0.21	0.12	0.07
Precision	0.95	0.24	0.16	0.24	0.12	0.07
Accuracy	0.96	0.60	0.41	0.30	0.13	0.04



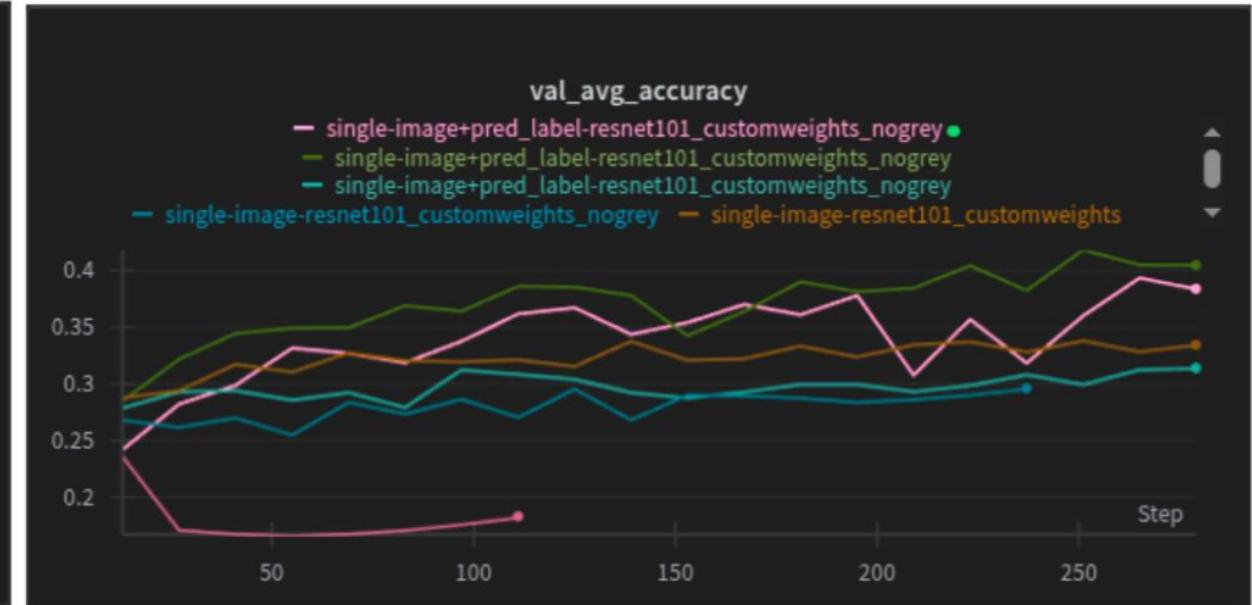
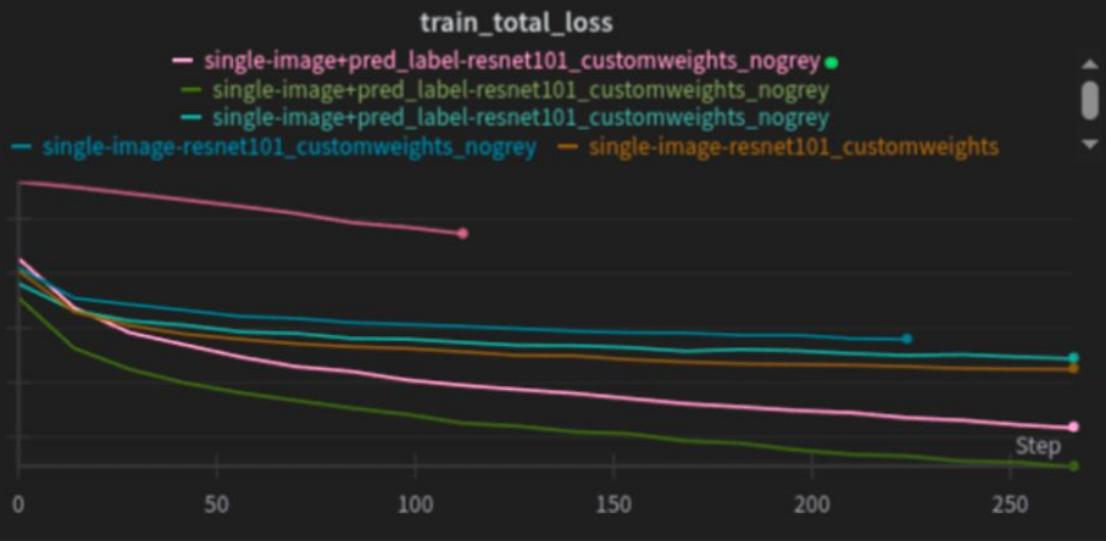
Experiments



Architecture	Val. avg. acc.
Starting structure	0.35
ResNet50, 3 FC layers, Dropout	0.44*
BatchNorm, masked loss	0.46*



Experiments



SimCLR ResNet101 for custom pre-training, 2 FC layers, no Dropout, Input: cropped im- ages	0.29
3 FC layers, Dropout, weight decay	0.31
Input: cropped images + YOLO predicted la- bel	0.33
Backbone unfreezed: slow training	0.38
Reduced regularization (Dropout and L2)	0.41



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

The Matcher



= ZARA

CERCA

LOGIN SUPPORTO CARRELLO [0]

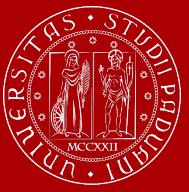
VISUALIZZA TUTTO T-SHIRTS Camicie POLO PANTALONI DENIM GIACCHE | TRENCH FELPE SCARPE & BORSE ACCESSORI | PROFUMI VISTA 2 3 FILTRI

CAMICIA DI LINO E COTONE ■+6 26,95 EUR ESPADRILLAS MOCASSINO... ■+2 49,95 EUR CAMICIA 100% LINO ■+4 45,95 EUR MAGLIETTA BOXY FIT SLAV... ■+3 17,95 EUR

BLAZER DOPPIOPETTO DA COM... 99,95 EUR BLAZER DA COMPLETO SLI... ■+1 79,95 EUR MOCASSINI IN PELLE CAS... ■+3 59,95 EUR BERMUDA 100% LINO ■+4 39,95 EUR

ESPADRILLAS IN TESSUTO 29,95 EUR PANTALONI FLUIDI RELAXE... ■+2 39,95 EUR BERMUDA 100% LINO ■+4 39,95 EUR BERMUDA DENIM BAGGY FIT ■+2 39,95 EUR

86



The Matching strategy



‘query’



The Matching strategy

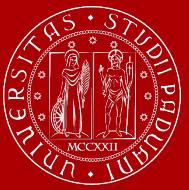


detectors

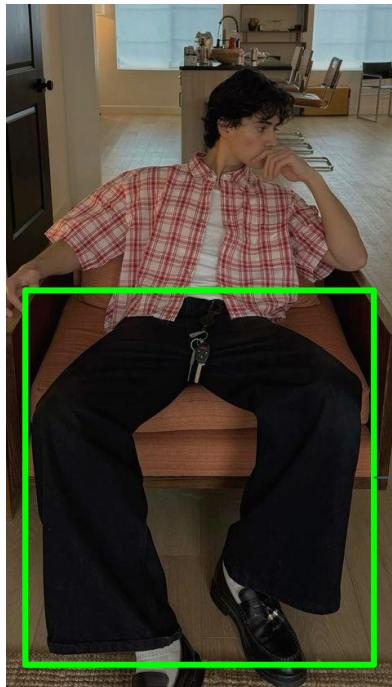


‘query’

predicted label	1 ('bottomwear')
confidence	0.34



The Matching strategy

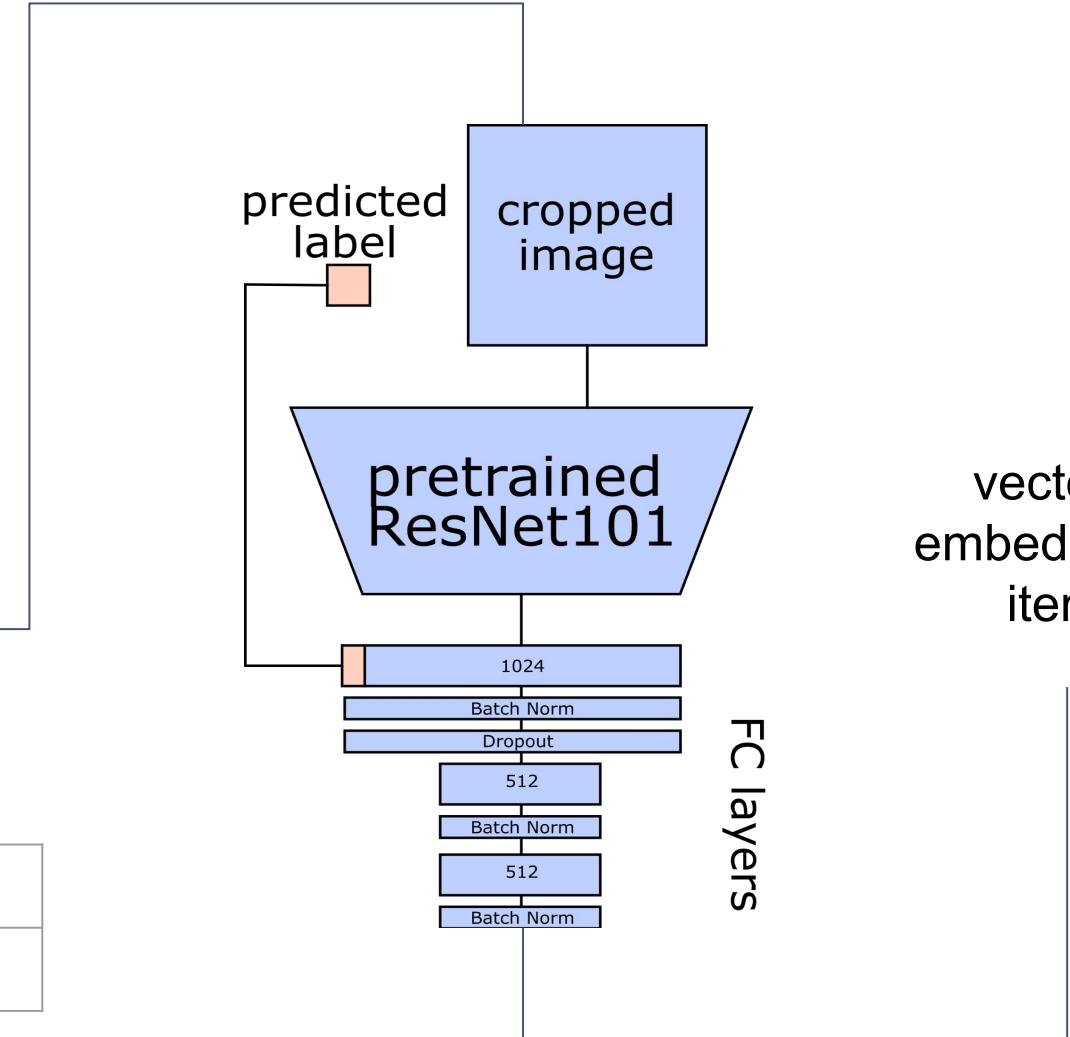


detectors



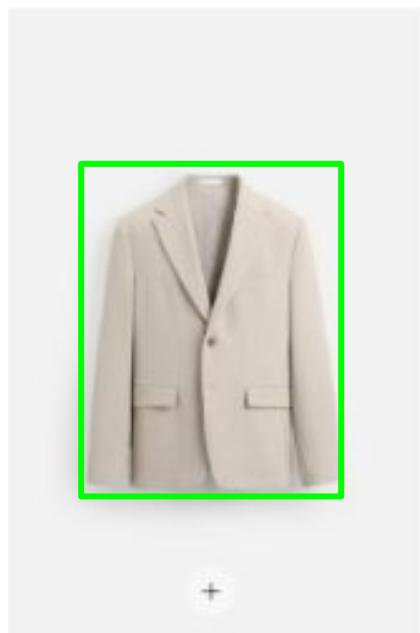
‘query’

predicted label	1 ('bottomwear')
confidence	0.34





The Matching strategy

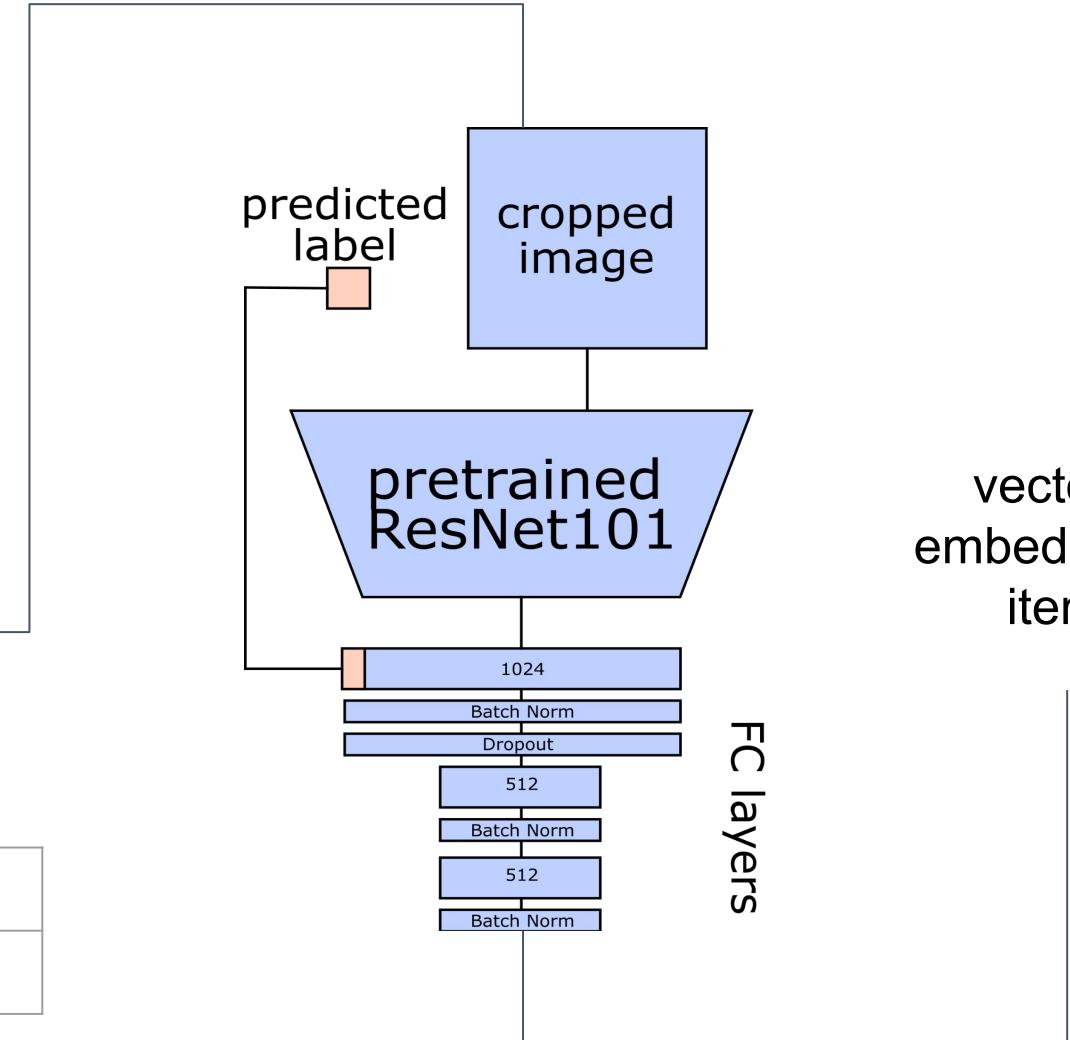


detectors



'wardrobe'

predicted label	0 (outwear)
confidence	0.71





The Matching strategy



Each item is mapped into a
512-dimensional,
batch-normalized **vector.**
The matching assumption:
similar items lay close in the
embedding space

How to
evaluate the
matcher
performance?



Again on DeepFashion Dataset

'pos' information: images on outfit are wore-on, while images in 'wardrobe' are flat





Again on DeepFashion Dataset

Again we used the naming convention to collect clothes with a ‘flat’ and a not-‘flat’ version

Set of Queries



Wardrobe





Matcher Evaluation Strategy

In this scenario we have the **ground truth**

Set of Queries



Wardrobe





Matcher Evaluation Strategy

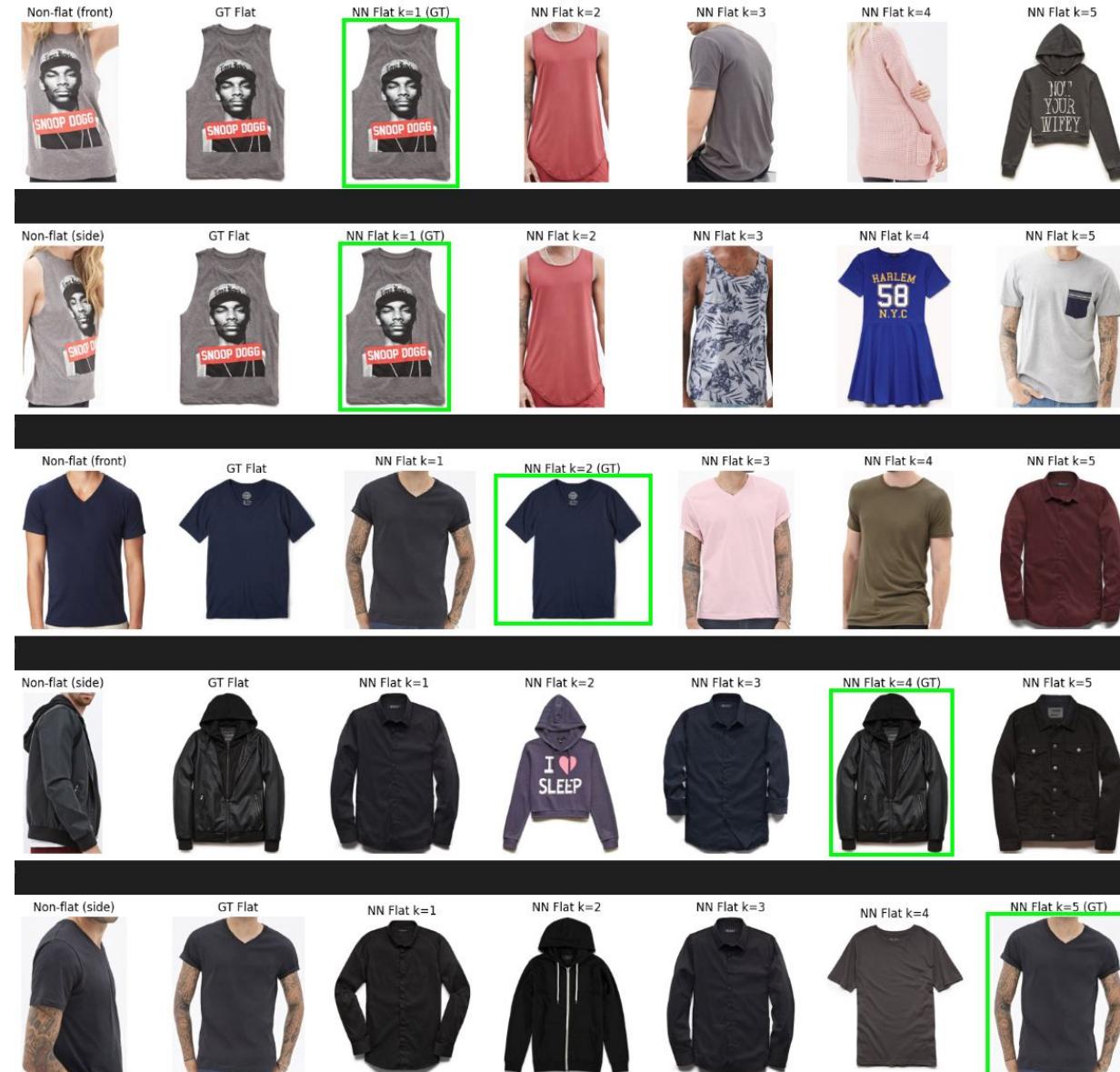
Is the ground-truth within
the k -closest neighbors
of the query encoding, in
the wardrobe space?





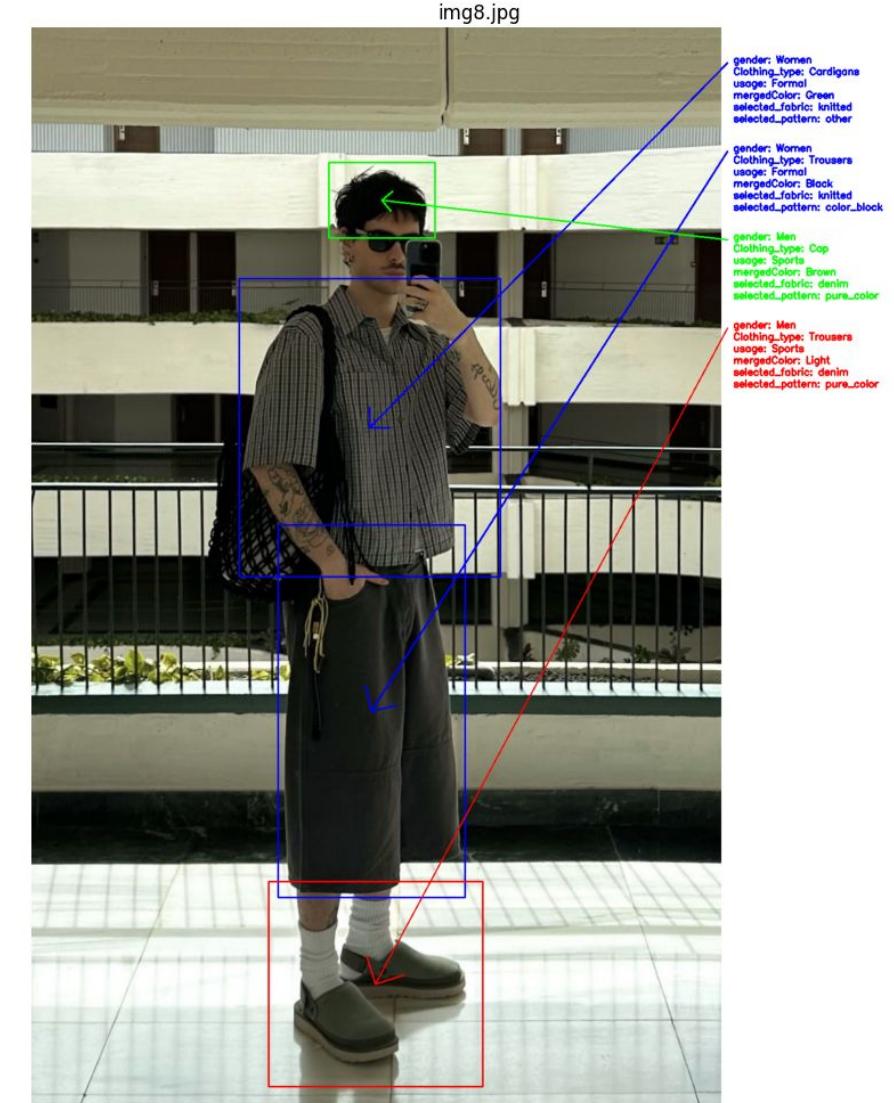
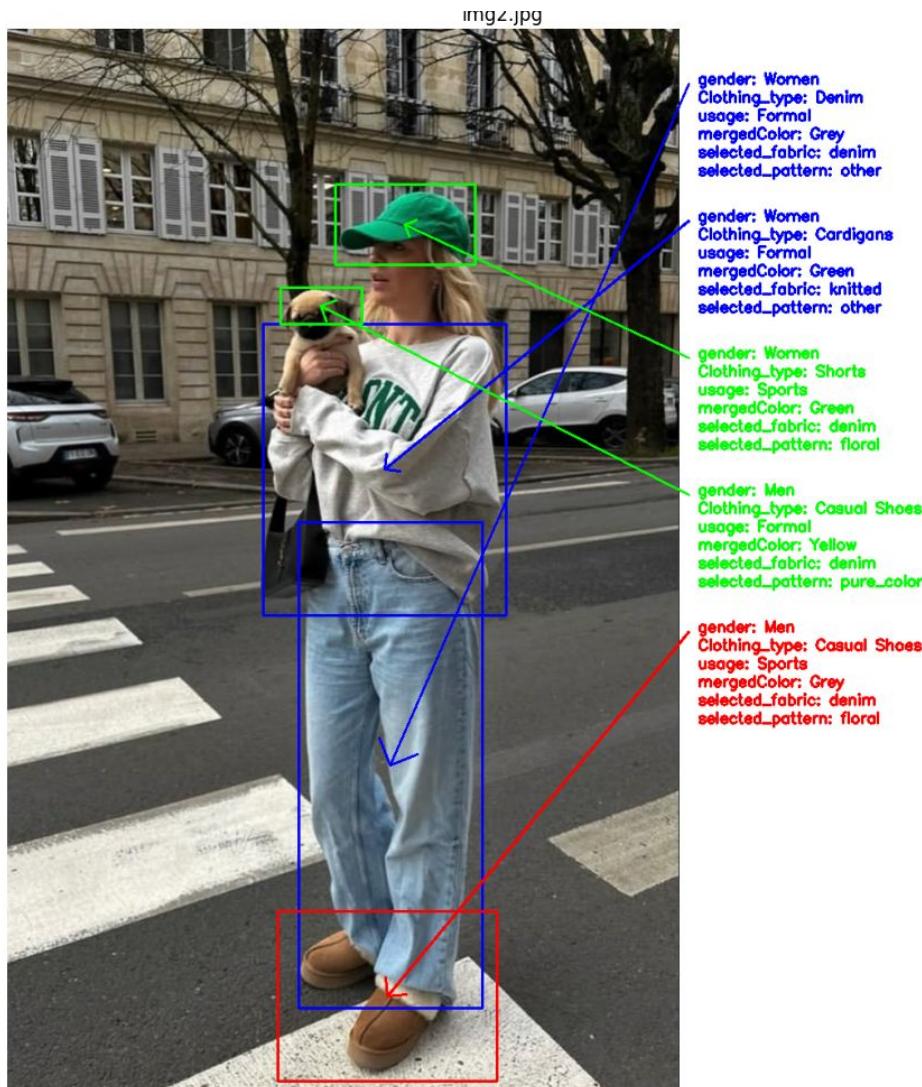
Matcher Evaluation: Results

Top- k	Accuracy
1	37.36%
2	48.28%
3	56.90%
4	63.79%
5	67.24%





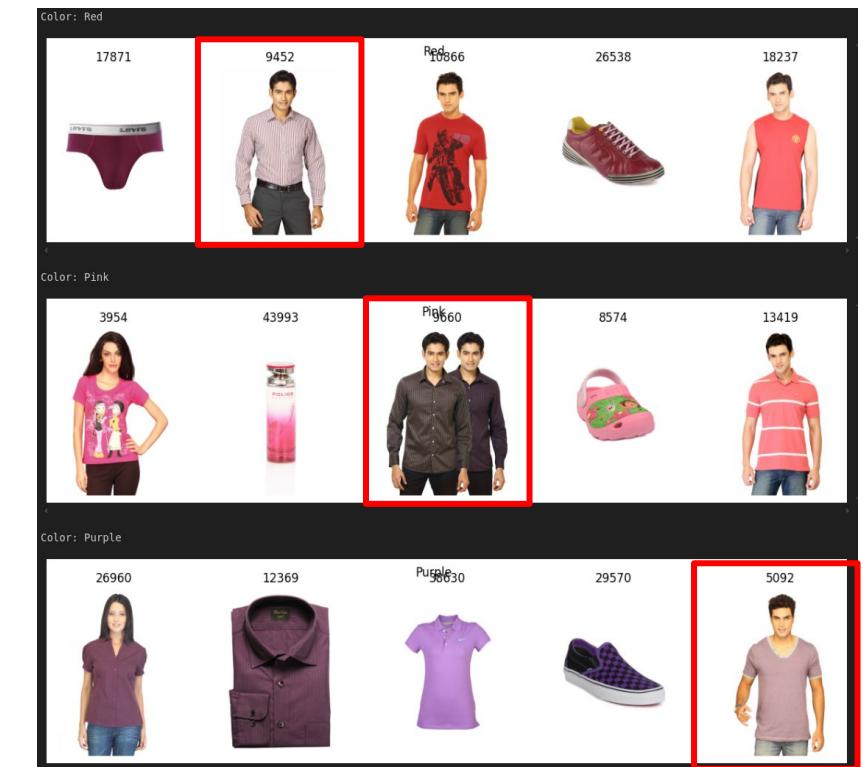
Conclusion





Further Improvements

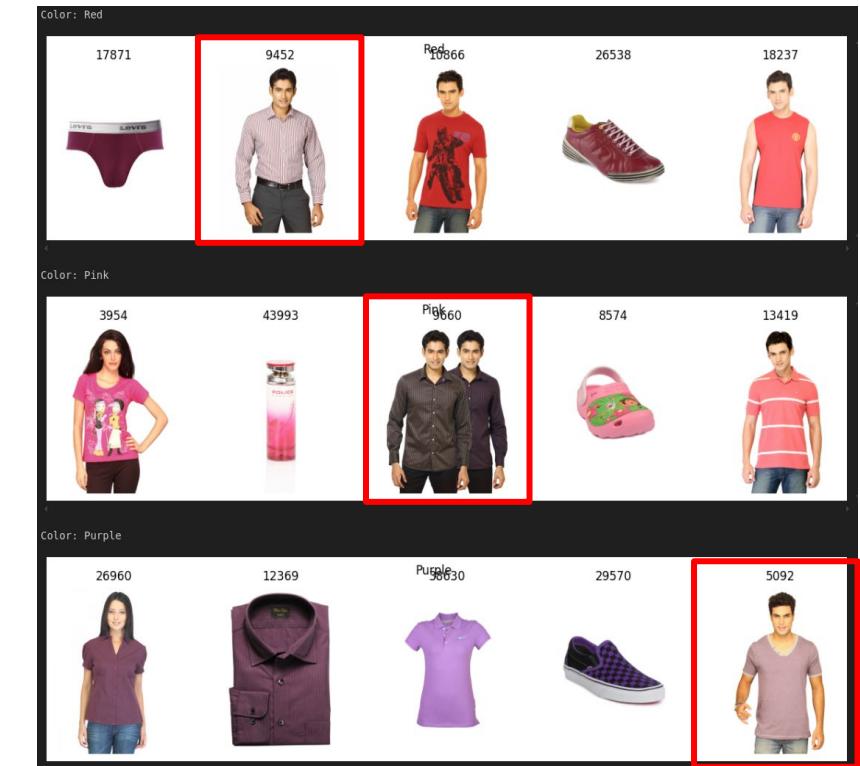
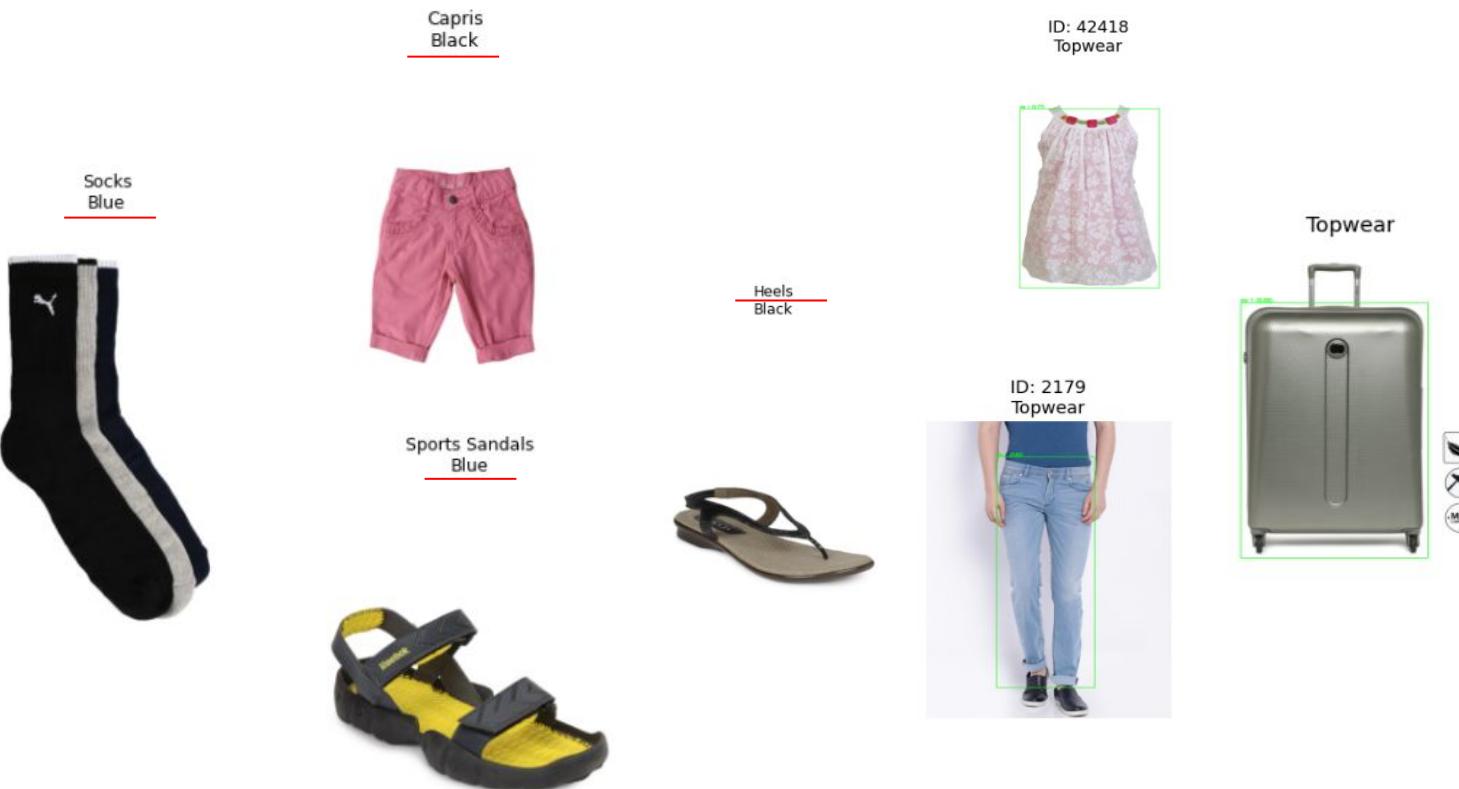
rubbish in - rubbish out: Fashion Product Images Dataset





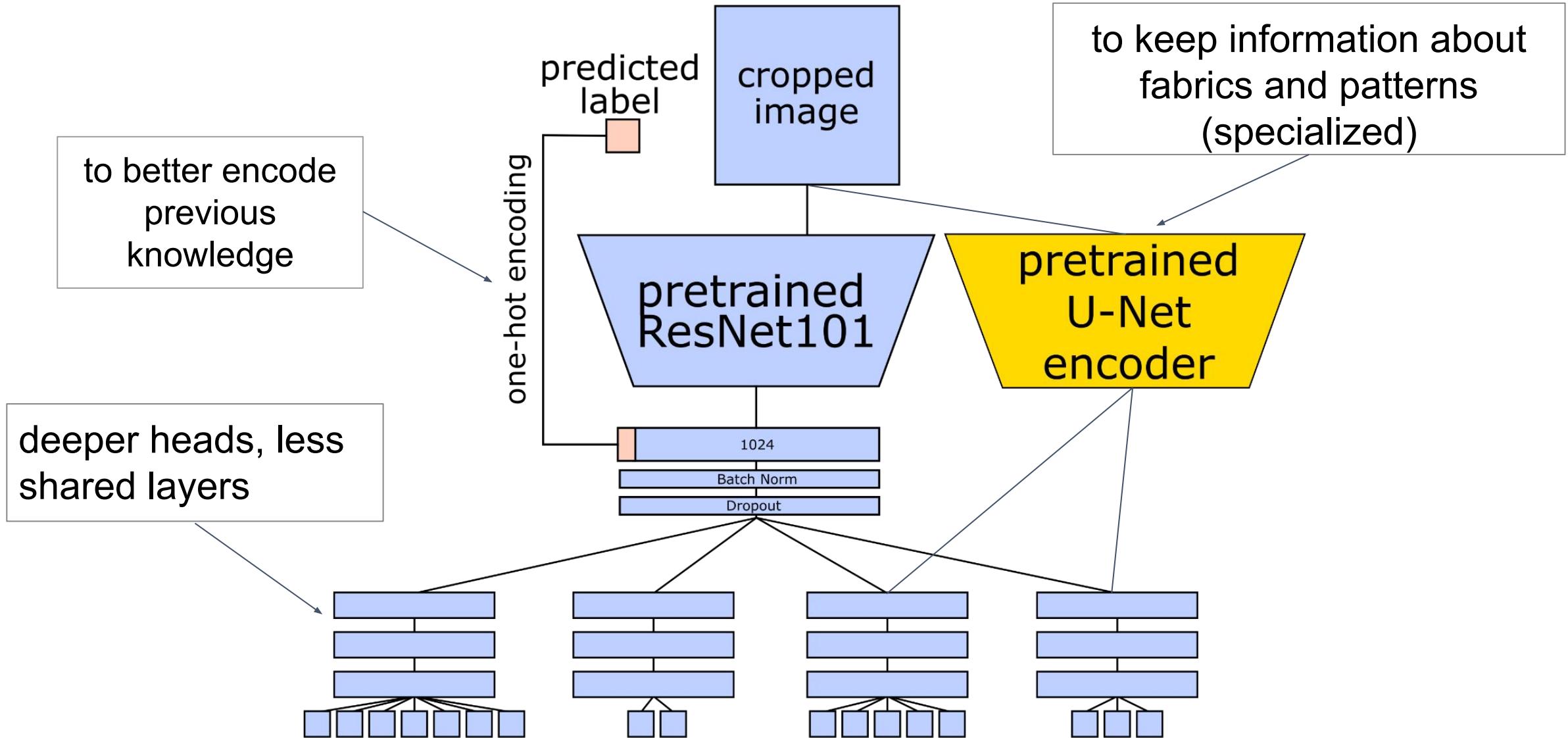
Further Improvements

Different dataset (complex to find) or weak-classifier(s) to re-label the dataset through unsupervised learning (by clustering features)





Further Improvements

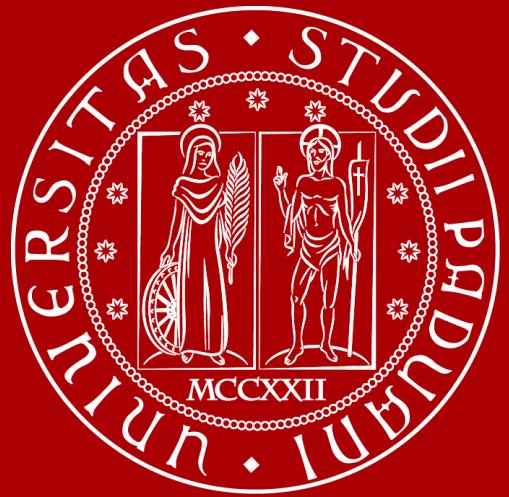




UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Thank you for your attention!

Mattia Roccatello
Ulaşcan Akbulut
Ayşenur Oya Özen



UNIVERSITÀ
DEGLI STUDI
DI PADOVA