

Fashionpedia: Ontology, Segmentation, and an Attribute Localization Dataset

Menglin Jia^{*1}, Mengyun Shi^{*1,4}, Mikhail Sirotenko^{*3}, Yin Cui^{*3},
Claire Cardie¹, Bharath Hariharan¹, Hartwig Adam³, Serge Belongie^{1,2}

¹Cornell University

²Cornell Tech

³Google Research

⁴Hearst Magazines

1 Supplementary Material

In our work we presented the new task of instance segmentation with attribute localization. We introduced a new ontology and dataset, Fashionpedia, to further describe the various aspects of this task. We also proposed a novel evaluation metric and Attribute-Mask R-CNN model for this task. In the supplemental material, we provide the following items that shed further insight on these contributions:

- More comprehensive experimental results (§ 1.1)
- An extended discussion of Fashionpedia ontology and potential knowledge graph applications (§ 1.2)
- More details of dataset analysis (§ 1.3)
- Additional information of annotation process (§ 1.4)
- Discussion (§ 1.5)

1.1 Attribute-Mask R-CNN

Per-class evaluation. Fig. 1 presents detailed evaluation results per supercategory and per category. In Fig. 1, we follow the same metrics from COCO leaderboard (AP, AP50, AP75, AP1, APm, APs, AR@1, AR@10, AR@100, ARs@100, ARm@100, ARI@100), with τ_{IoU} and τ_{F1} if possible. Fig. 1 shows that metrics considering both constraint τ_{IoU} and τ_{F1} are always lower than using τ_{IoU} alone across all the supercategories and categories. This further demonstrates the challenging aspect of our proposed task.

In general, categories belong to “garment parts” have a lower AP and AR, comparing with “outerwear” and “accessories”.

A detailed breakdown of detection errors is presented in Fig. 2 for supercategories and three main categories. In terms of supercategories in Fashionpedia, “outerwear” errors are dominated by within supercategory class confusions (Fig. 2(a)). Within this supercategory class, ignoring localization errors would only raise AP slightly from 77.5 to 79.1 (+1.6). A similar trend can be observed in class “skirt”, which belongs to “outerwear” (Fig. 2(d)). Detection errors of “part” (Fig. 2(b) 2(e)) and “accessory” (Fig. 2(c) 2(f)) on the other hand, are

* equal contribution.

	AP		AP50		AP70		API		APm		APs		AR@1		AR@10		AR@100		ARs@100		ARm@100		ARl@100		
	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1	IoU	IoU + F1		
overall	33.35	43.15	42.28	60.26	37.60	47.65	35.38	50.01	27.03	40.22	9.39	17.30	41.17	47.35	44.41	53.97	44.43	56.06	12.20	22.46	33.12	49.67	47.50	68.28	
superdis-outwear	40.72	64.12	49.01	77.41	46.21	72.85	42.98	67.07	29.25	44.44	4.39	16.99	52.63	76.42	53.05	77.07	53.05	77.07	5.88	21.83	33.87	50.32	55.56	80.12	
superdis-part	13.41	19.28	20.81	35.51	14.37	18.45	14.52	28.30	16.40	23.86	9.83	12.49	15.01	22.27	22.38	30.05	22.45	30.27	13.11	18.30	24.84	37.53	24.90	53.20	
superdis-accessory	56.07		77.93		63.88		57.48		60.54		24.99		54.39		64.98		64.98		29.88		67.53		73.48		
superdis-clothes	shirt, blouse	35.04	53.52	46.65	72.49	36.71	58.86	42.84	62.26	17.44	36.43	0.00	5.05	45.54	64.31	46.91	67.35	46.91	67.35	0.00	5.00	25.46	48.46	55.72	75.08
	top*	46.09	69.24	56.19	85.22	52.49	78.72	50.92	75.55	33.86	49.67	5.54	30.91	54.45	76.27	55.50	77.97	55.50	77.97	11.54	42.31	42.47	59.90	60.50	84.01
	sweater	44.73	56.66	51.34	64.89	51.16	64.66	50.84	64.34	9.43	11.78			57.86	70.46	59.00	71.90	59.00	71.90			21.33	26.67	65.28	79.44
	cardigan	25.41	43.79	32.49	55.67	31.87	54.78	25.41	43.79					46.33	70.83	46.33	70.83	46.33	70.83					46.33	70.83
	jacket	39.63	68.66	50.10	66.83	46.14	79.91	40.26	69.55	23.42	46.63			49.03	76.94	49.53	77.70	49.53	77.70			23.33	46.67	50.42	76.76
	vest	38.13	53.08	47.25	65.80	42.38	59.38	42.13	58.62	0.00	0.00			45.05	60.00	45.05	60.00	45.05	60.00			0.00	0.00	49.55	66.00
	pants	56.29	82.86	64.76	95.29	62.48	91.96	57.21	84.26	39.67	49.07			63.68	87.17	64.22	87.83	64.22	87.83			46.87	57.33	65.09	89.36
	shorts	51.49	74.19	60.83	87.88	56.32	80.62	56.37	81.16	45.81	63.08	12.00	40.00	59.50	81.89	60.19	82.64	60.19	82.64	12.00	40.00	50.89	69.43	65.53	89.86
	skirt	41.02	70.30	45.18	77.78	43.93	75.98	41.39	70.92	50.65	81.46	0.00	0.00	53.65	83.64	53.72	83.70	53.72	83.70	0.00	0.00	57.36	89.09	54.17	84.43
	coat	41.77	69.10	53.01	87.93	48.26	79.64	41.66	69.97	48.00	80.00			52.35	79.71	52.35	79.71	52.35	79.71			48.00	80.00	52.39	79.71
	dress	43.26	81.73	49.47	93.40	47.72	90.06	43.65	82.44	5.50	10.77			52.47	87.26	52.53	87.32	52.53	87.32			8.80	16.00	52.96	88.03
	jumpsuit	29.01	46.03	34.54	54.56	33.89	53.67	28.59	45.60	48.00	60.00			48.48	70.95	48.48	70.95	48.48	70.95			48.00	60.00	48.50	71.50
	cape	37.44	64.39	45.36	78.55	45.36	78.55	37.44	64.39					55.80	84.00	55.80	84.00	55.80	84.00					55.80	84.00
superdis-hat	outfit	11.44	35.39	21.23	69.77	10.90	34.43	0.00	7.42	15.90	47.46	10.17	23.64	15.83	43.90	16.38	46.01	16.38	46.01	11.06	28.51	20.91	59.50	0.00	50.00
	lapel	3.96	46.86	5.88	79.86	4.92	52.20	3.19	56.33	10.12	47.11	3.37	13.75	12.44	52.52	14.53	55.11	14.53	55.11	3.33	16.67	12.98	50.48	18.91	68.89
	sleeve	42.21	70.81	54.81	93.71	48.54	81.20	45.89	79.74	45.13	70.95	19.78	32.71	29.17	41.22	53.35	75.12	53.35	75.12	24.72	40.62	51.88	74.43	63.25	85.27
	pocket	11.64	34.10	18.48	58.18	13.30	36.16	16.37	47.17	11.07	46.83	16.24	26.04	13.41	24.38	24.87	46.88	24.89	46.95	22.19	35.50	29.01	62.93	21.25	72.50
	neckline	0.08	14.33	0.22	47.23	0.03	3.60	0.00	3.48	0.01	13.74	1.23	19.61	1.39	22.98	1.84	25.18	1.84	25.18	2.16	22.07	0.98	33.71	0.00	33.33
	hood	27.44		63.30		16.63		32.82		34.15		9.74		39.06		44.06		44.06		12.50		42.38		67.14	
	epauleto	32.36		61.44		27.57				15.10		35.12		28.57		47.86		47.86		45.38		80.00			
	buckle	13.57		40.80		4.81				22.47		13.02		20.30		23.28		23.28		18.64		57.50			
	zipper	3.70		13.93		0.68				5.10		6.56		5.62		13.56		14.18		14.93		11.88			
	applique	11.54		17.27		12.34		29.16		12.25		10.54		19.84		30.82		31.31		22.69		36.33		46.00	
	bead	6.26		12.13		6.10		36.45		33.57		2.84		10.28		14.95		15.98		8.04		56.67		76.67	
	bow	12.48		18.23		16.83		0.00		0.00		29.01		20.00		20.00		20.00		40.00		0.00		0.00	
	flower	1.29		2.39		1.29		8.86		0.00		2.62		3.24		10.54		10.81		8.44		0.00		65.00	
	fringe	18.34		31.71		19.42		41.00		43.27		1.19		16.00		21.00		21.00		1.18		43.33		50.00	
	ribbon	0.99		2.97		0.99				3.74		0.00		2.22		6.67		6.67		0.00		12.00			
	rivet	7.41		20.32		2.96		0.00		35.35		7.29		8.04		18.81		20.56		20.29		35.00		0.00	
	ruffle	25.12		36.93		28.80		35.17		21.79		3.74		43.42		49.61		49.61		12.22		34.40		66.67	
	sequin	4.30		4.49		4.49		18.57		0.50		0.00		21.54		21.54		21.54		0.00		22.50		63.33	
	tassel	0.00		0.00		0.00		0.00		0.00		0.00		0.00		0.00		0.00		0.00		0.00			
superdis-accessory	glasses	68.06		94.75		84.48		88.00		75.94		44.20		71.54		71.77		71.77		47.06		80.11		90.00	
	hat	72.45		88.58		87.30		71.93		77.18		6.73		77.30		77.30		77.30		6.67		79.81		81.59	
	head acc*	36.76		72.11		29.97		34.04		49.75		30.97		44.22		46.61		46.61		35.25		63.97		43.33	
	tie	61.64		75.64		75.64				61.64				83.33		83.33		83.33						83.33	
	glove	52.42		84.45		69.67				72.85		2.88		36.45		58.39		58.39		8.75		75.65			
	watch	52.23		90.44		81.19				68.77		47.64		57.02		60.00		60.00		56.00		73.68			
	belt	38.94		71.17		41.30		32.59		45.02		30.69		48.54		50.98		50.98		35.84		54.30		64.44	
	leg warmer	40.59		47.13		47.13		54.40		25.64				35.71		55.71		55.71				32.00		68.89	
	lights, stockings	70.66		85.12		77.66		70.91		74.03		12.62		42.38		80.16		80.16		25.00		78.57		83.86	
	sock	42.38		59.07		49.97		23.33		58.77		30.21		28.51		50.80		50.80		38.98		65.95		78.00	
	shoe	61.10		91.89		69.17		58.49		67.22		43.03		36.63		68.07		68.07		45.38		74.33		85.21	
	bag, wallet	58.26		85.84		63.55		67.31		56.58		15.82		62.06		65.47		65.47		20.00		59.92		76.59	
	scarf	46.47		73.81		47.23		47.21		53.89		10.10		53.75		57.08		57.08		10.00		56.67		60.32	
	umbrella	83.07		100.00		100.00		83.07						84.00		84.00		84.00						84.00	

Fig. 1: Detailed results (for masks) using Mask R-CNN with SpineNet-143 backbone. We present the same metrics as COCO leaderboard for overall categories, three super categories for apparel objects, and 46 fine-grained apparel categories. We use both constraints (for example, AP_{IoU} and AP_{IoU+F1}) if possible. For categories without attributes, the value represents AP_{IoU} or AR_{IoU} . “top” is short for “top, t-shirt, sweatshirt”. “head acc” is short for “headband, head covering, hair accessory”

dominated by both background confusion and localization. “part” also has a lower AP in general, compared with other two super-categories. A possible reason is that objects belong to “part” usually have smaller sizes and lower counts.

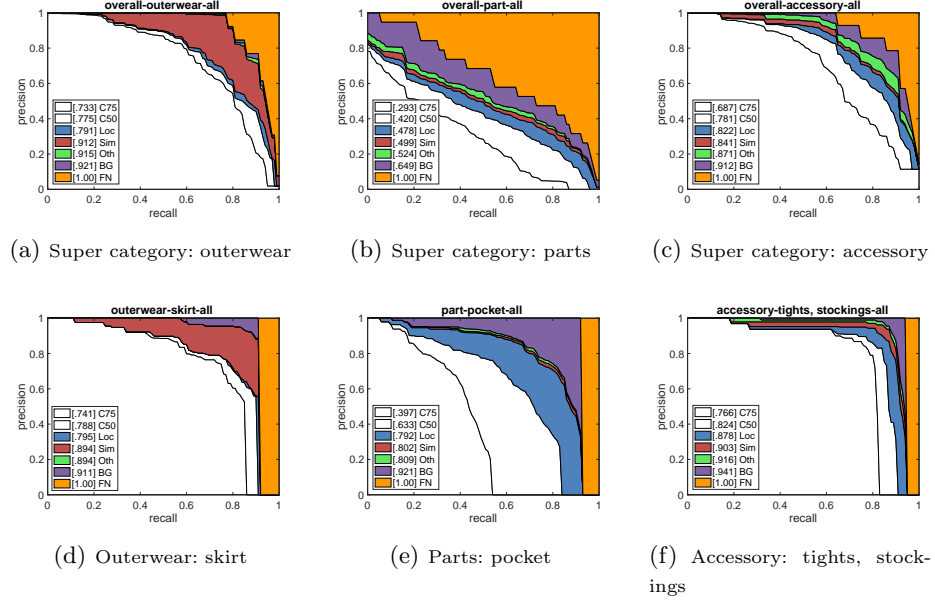


Fig. 2: Main apparel detectors analysis. Each plot shows 7 precision recall curves where each evaluation setting is more permissive than the previous. Specifically, **C75**: strict IoU ($\tau_{IoU} = 0.75$); **C50**: PASCAL IoU ($\tau_{IoU} = 0.5$); **Loc**: localization errors ignored ($\tau_{IoU} = 0.1$); **Sim**: supercategory False Positives (FPs) removed; **Oth**: category FPs removed; **BG**: background (and class confusion) FPs removed; **FN**: False Negatives are removed. The first row (*overall-[supercategory]-[size]*) contains results for three supercategories in Fashionpedia; the second row (*[supercategory]-[category]-[size]*) shows results for three fine-grained categories (one per supercategory). Legends present the area under each curve (corresponds to AP metric) in brackets as well

F1 score calculation. Since we measure the f1 score of predicted attributes and groundtruth attributes per mask, we consider the both options of multi-label multi-class classification with 294 classes for one instance, and binary classification for 294 instances. Multi-label multi-class classification is a straightforward task, as it is a common setting for most of the fine-grained classification tasks. In binary classification scenario, we consider the 1 and 0 of the multi-hot encoding of both results and ground-truth labels as the positive and negative classes respectively. There are also two averaging choices: “micro” and “macro”. “Micro”

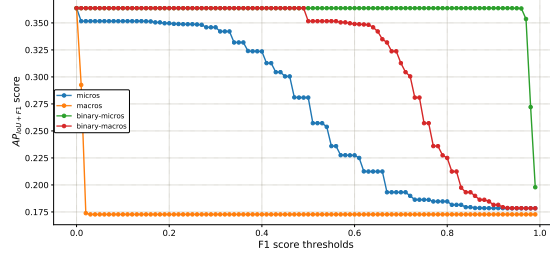


Fig. 3: $AP_{IoU+F1}^{F1=\tau_{F1}}$ score with different τ_{F1} . The value presented are average over $\tau_{IoU} \in [0.5 : 0.05 : 0.95]$. We use “binary-macro” as our main metric

averaging calculates the score globally by counting the total true positives, false negatives and false positives. “Macro” averaging calculates the metrics for each attribute class and reports the unweighted mean. In sum, there are four options of f1-score averaging methods: 1) “micro”, 2) “macro”, 3) “binary-micro”, 4) “binary-macro”.

As shown in Fig. 3, we present the $AP_{IoU+F1}^{F1=\tau_{F1}}$, with τ_{IoU} averaged in the range of $[0.5 : 0.05 : 0.95]$. τ_{F1} is increased from 0.0 to 1.0 with a increment of 0.01. Fig. 3 illustrates that as the value of τ_{F1} increases, $AP_{IoU+F1}^{F1=\tau_{F1}}$ decreases in different rates given different choices of f1 score calculation. There are 294 attributes in total, and an average of 3.7 attributes per mask in Fashionpedia training data. It’s not surprising to observe that “Binary-micro” produces high f1-scores in general (higher than 0.97), as the AP_{IoU+F1} score only decreases if the $\tau_{F1} \geq 0.97$. On the other hand, “macro” averaging in multi-label multi-class classification scenario gives us extremely low f1-scores (0.01 – 0.03). This further demonstrates the room for improvement for localized attribute classification task. We used “binary-macros” as our main metric.

Result visualization. Figure 4 shows that our simple baseline model can detect most of the apparel categories correctly. However, it also produces false positives sometimes. For example, it segments legs as tights and stockings (Figure 4(f) Tights). A possible reason is that both objects have the same shape and stockings are worn on the legs.

Predicting fine-grained attributes, on the other hand, is a more challenging problem for the baseline model. We summarize several issues: (1) predict more attributes than needed: (Figure 4(a) Neckline, (b) Collar, (c) Sleeve); (2) fail to distinguish among fine-grained attributes: for example, dropped-shoulder sleeve (ground truth) v.s. set-in sleeve (predicted) (Figure 4(e) Sleeve 1); (3) other false positives: Figure 4(e) Dress has a double-breasted opening, yet the model predicted it as the zip opening.

These results further show that there are rooms for improvement and future development of more advanced computer vision models on this instance segmentation with attribute localization task.

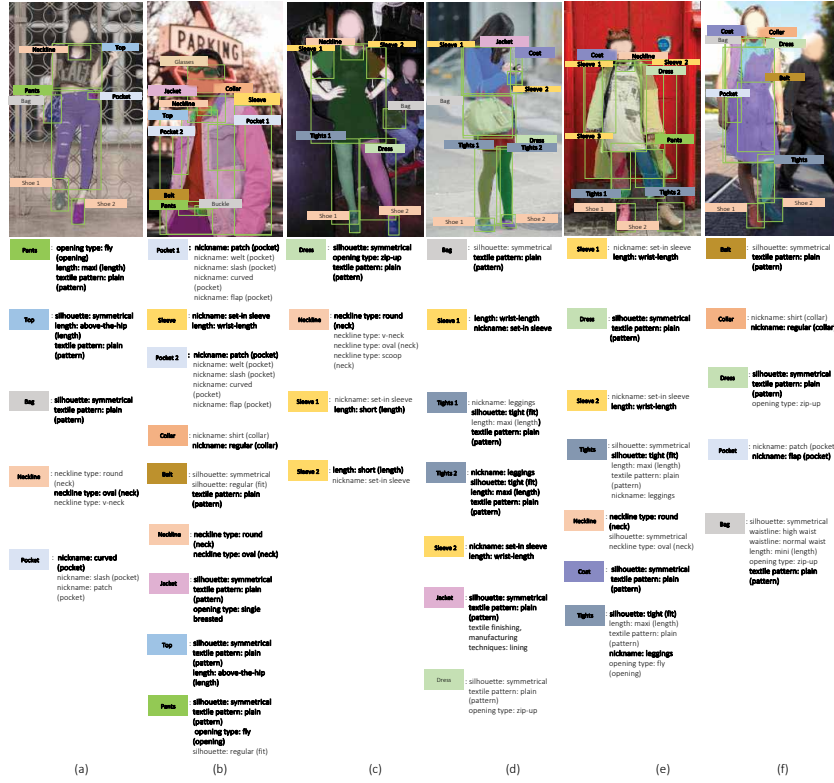


Fig. 4: Baseline results on the Fashionpedia validation set. Masks, bounding boxes, and apparel categories (category score > 0.6) are shown. Attributes from top 10 masks (that contain attributes) from each image are also shown. Correct predictions of objects and attributes are bolded

Result visualization on other datasets. Other fashion datasets such as ModaNet and DeepFashion2 also contain instance segmentation masks. Aside from the results presented in the main paper (see Table 3 of the main paper) on Fashionpedia, we present the qualitative analysis on the segmentation masks generated among Fashionpedia(Fig. 4), ModaNet(Fig. 5(a-f)) and DeepFashion2(Fig. 5(g-l)) datasets. Photos of the first row in Fig. 5 are from ModaNet. They show that the quality of the generated masks on ModaNet is fairly good and comparable to Fashionpedia in general (Fig. 5(a)). We also have a couple of observations of the failure cases: (1) fail to detect apparel objects: for example, the shoe from Fig. 5(c) is not detected. Parts of the pants (Fig. 5(c)) and coat (Fig. 5(d)) are not detected; (2) fail to detect some categories: Fig. 5(e) shows that the shoes on the shoe rack and right foot are not detected, possibly due to a lack of such instances in the ModaNet training dataset. Similar to



Fig. 5: Baseline results on ModaNet and DeepFashion2 validation set



Fig. 6: Generated masks on online-shopping images [2]. (a) and (b) show the same types of shoes in different settings. Our model correctly detects and categorizes the pair of shoes worn by a fashion model, yet mistakenly detects shoes as jacket and a bag in (b)

Fashionpedia, ModaNet mostly consist of street style images. See Fig. 6(b) for example predictions from model trained on Fashionpedia; 3) close-up images: ModaNet contains mostly full-body images. This might be the possible reason to the decreased quality of predicted masks on close-up shot like Fig. 5(f).

For DeepFashion 2 (Fig. 5(g,h,k)), the generated segmentation masks tends to not tightly follow the contours of garments in the images. The main reason possibly is that the average number of vertices per polygon is 14.7 for Deep-fashion2, which is lower than Fashionpedia and ModaNet (see Table 2 in the main text). Our qualitative analysis also shows that: 1) the model will generate the segmentation masks of pants (Fig. 5(i)) and tops (Fig. 5(j)) that are not visible in the images. Both of them are covered by a jacket. And we find that in DeepFashion 2, some part of the garments which is covered by other objects are indeed annotated with segmentation masks; 2) better performance on objects that are not on human body (Fig. 5(l)): DeepFashion 2 contains many commercial-customer image pairs (both images with and without human body) in the training dataset. In contrast, both Fashionpedia and ModaNet contain more images with human body than images without human body in the training datasets.

Generalization to the other image domains. For Fashionpedia, we also inference on images found in online shopping websites, which usually displays a single apparel category, with or without a fashion model. We found out that the learned model works reasonably well if the apparel item is worn by a model (Fig. 6).

1.2 Fashionpedia Ontology and Knowledge Graph

Fig. 7 presents our Fashionpedia ontology in detail. Fig. 8 and 9 displays the training data mask counts per category and per attributes. Utilizing the proposed ontology and the image dataset, a large-scale fashion knowledge graph can be constructed to represent the fashion world in the product level. Fig. 10 illustrates a subset of the Fashionpedia knowledge graph.

Apparel graphs. Integrating the main garments, garment parts, attributes, and relationships presented in one outfit ensemble, we can create an apparel graph representation for each outfit in an image. Each apparel graph is a structured representation of an outfit ensemble, containing certain types of garments. Nodes in the graph represent the main garments, garment parts, and attributes. Main garments and garment parts are linked to their respective attributes through different types of relationships. Figure 11 shows more image examples with apparel graphs.

Fashionpedia knowledge graph. While apparel graphs are localized representations of certain outfit ensembles in fashion images, we can also create a single Fashionpedia knowledge graph (Figure 10). The Fashionpedia knowledge graph is the union of all apparel graphs and includes entire main garments, garment parts, attributes, and relationships in the dataset. In this way, we are able to represent and understand fashion images in a more structured way.

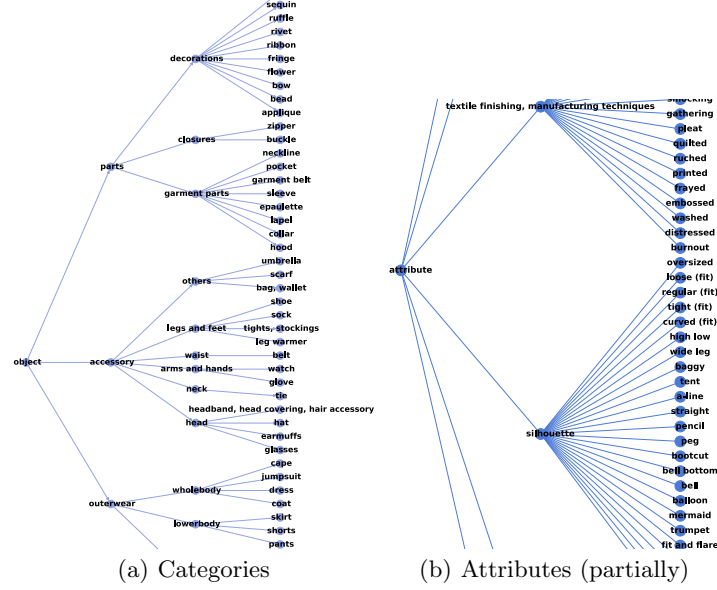


Fig. 7: Apparel categories (a) and fine-grained attributes (b) hierarchy in Fashionpedia

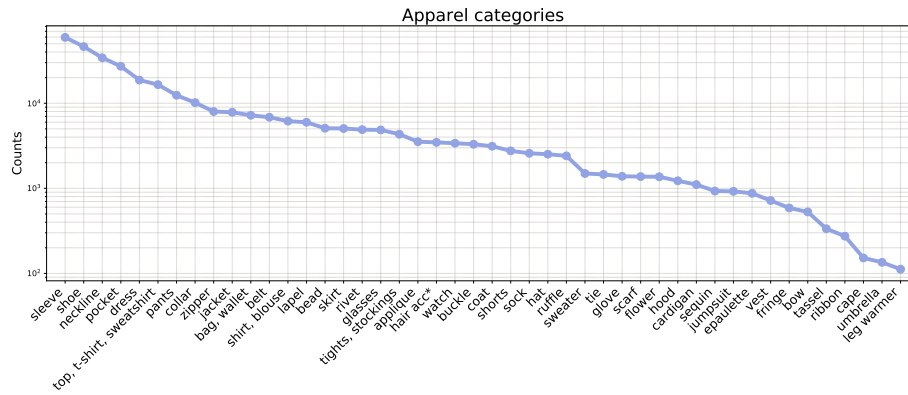


Fig. 8: Mask counts per apparel categories in training data. “head acc” is short for “headband, head covering, hair accessory”

We expect our Fashionpedia knowledge graph and the database to have applicability to extending the existing knowledge graph (such as WikiData [1]) with novel domain-specific knowledge, improving the underlying fashion product recommendation system, enhancing search engine’s results for fashion visual search, resolving ambiguous fashion-related words for text search, and more.

1.3 Dataset Analysis

Fig. 11 shows more annotation examples, represented in the exploded views of annotation diagrams. Table 1 displays the details about “not sure” and “not on the list” results during attribute annotation process. We present the result per super-categories of attributes. Label “not sure” means the expert annotator is uncertain about the choice given the segmentation mask. “Not on the list” means the annotator is certain that the given mask presents another attributes that is not presented in the Fashionpedia ontology. Other than “nicknames” (which is the specific name for a certain apparel category), less than 6% of the total masks account for the “not on the list” category.

Fig. 12 and 13 also compare Fashionpedia and other images datasets in terms of image size and vertices per polygons.

We compare image resolutions between Fashionpedia and four other segmentation datasets (COCO and LVIS share the same images). Fig. 12 shows that images in Fashionpedia have the most diverse image width and height. While ModaNet has the most consistent resolutions of images. Note that high resolution images will burden the data loading process of training. With that in mind, we will release our dataset in both the resized and the original versions.

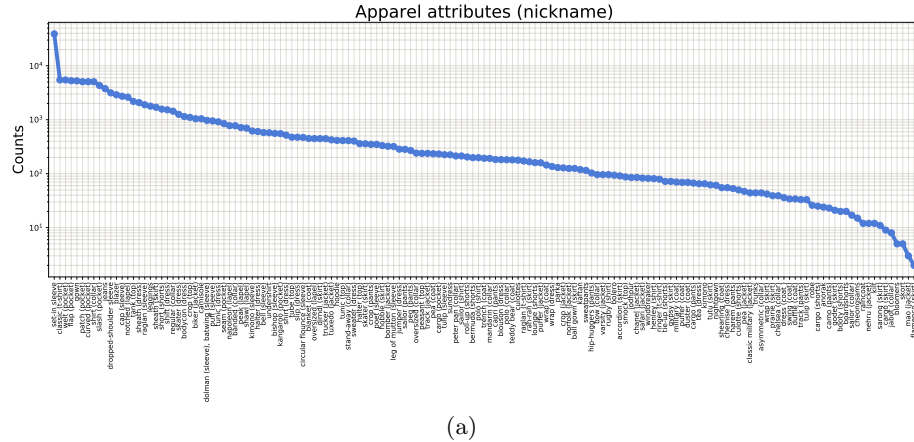


Fig. 9: Mask counts per attributes in training data, grouped by super categories. Best viewed digitally

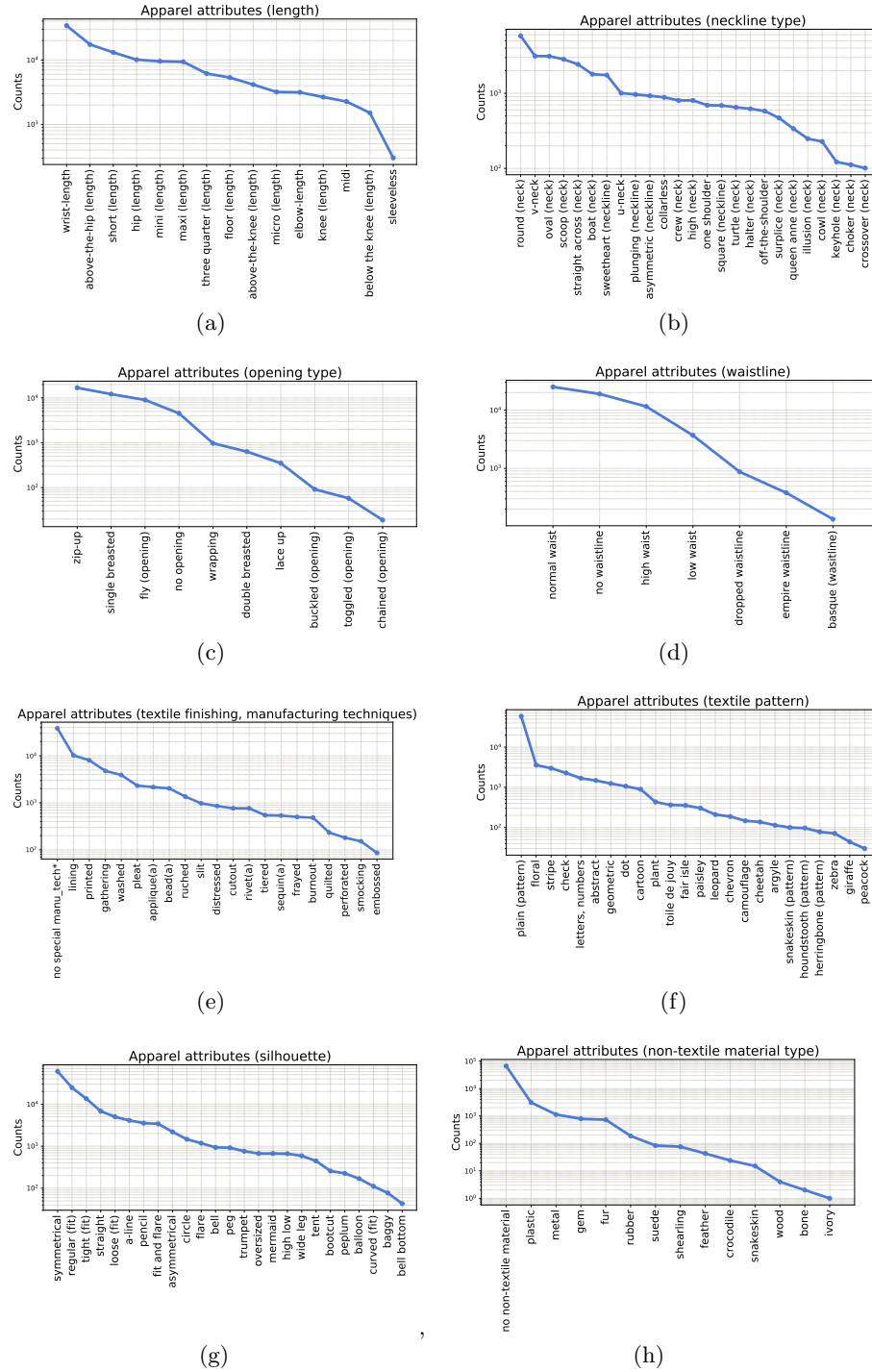


Fig. 9: Mask counts per attributes in training data, grouped by super categories (cont.). Best viewed digitally

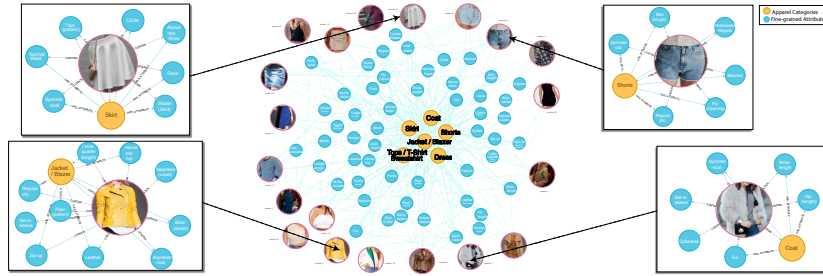


Fig. 10: Fashionpedia Knowledge Graph: we present a subset of the Fashionpedia knowledge graph by aggregating 20 annotated products. The knowledge graph can be used as a tool for generating structural information

We also report the distribution of number of vertices per polygons in Fig. 13. This measures the annotation effort in mask annotation. Fashionpedia has the second-widest range, next to LVIS.

1.4 Fashionpedia dataset creation details

Image collection. To avoid photo bias, all the images are randomly collected from Flickr and free license stock photo websites (Unsplash, Burst by Shopify, Free stocks, Kaboompics, and Pexels). The collected images are further verified manually by two fashion experts. Specifically, they check the scenes’ diversity

Table 1: Percentage of attributes in Fashionpedia broken down by super-class. “Tex finish, manu-tech.” is short for “Textile finishing, Manufacturing techniques”. Summaries of “not sure” and “not on the list” during attributes annotations are also presented. It was calculated by the counts divided by the total masks with attributes. “not sure” is mainly due to occlusion inside the images, which cause some super-classes (such as waistline, opening type, and length) are unidentifiable in the images. The percentage of “not on the list” is less than 15%. This demonstrates the comprehensiveness of our Fashionpedia ontology

Super-category	class count	not sure	not on the list
Length	15	12.79%	0.01%
Nickname	153	9.15 %	12.76%
Opening Type	10	32.69%	3.90%
Silhouettes	25	2.90%	0.27%
Tex finish, manu-tech	21	4.47%	1.34%
Textile Pattern	24	2.18%	5.30%
None-Textile Type	14	4.90%	4.07%
Neckline	25	9.57%	3.38%
Waistline	7	30.46%	0.17%

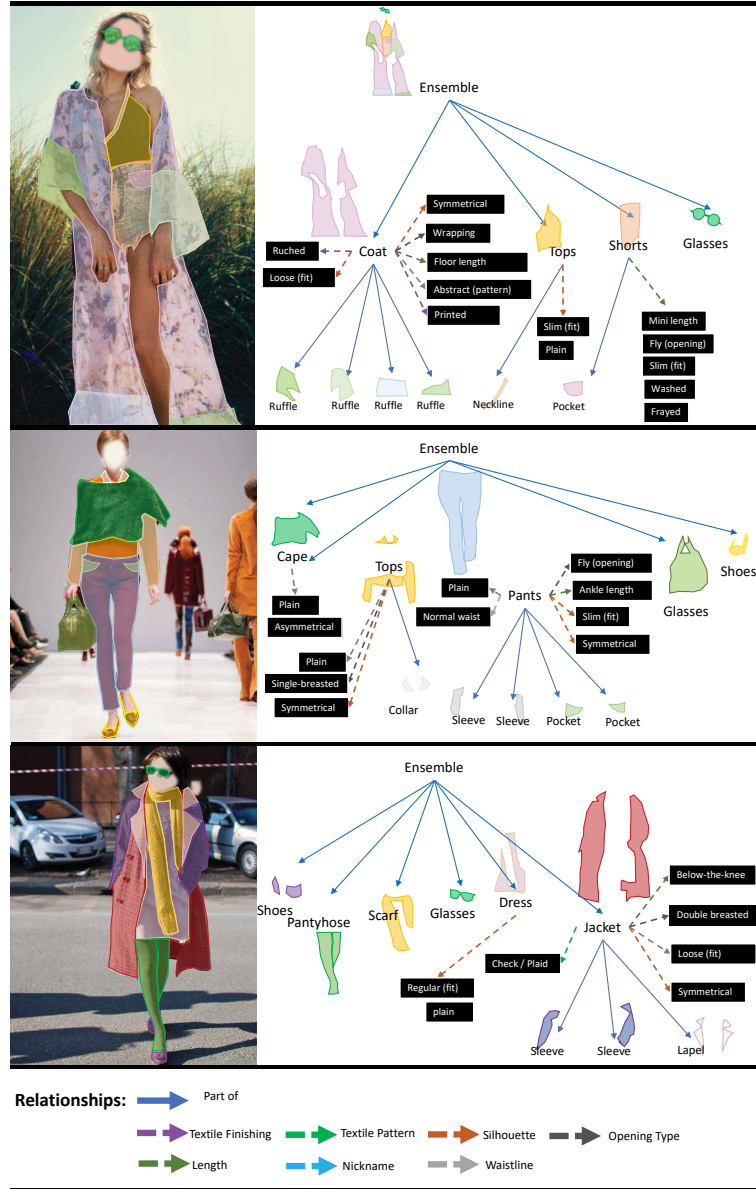


Fig. 11: Example images and annotations from our dataset: the images are annotated with both instance segmentation masks and fine-grained attributes (black boxes)

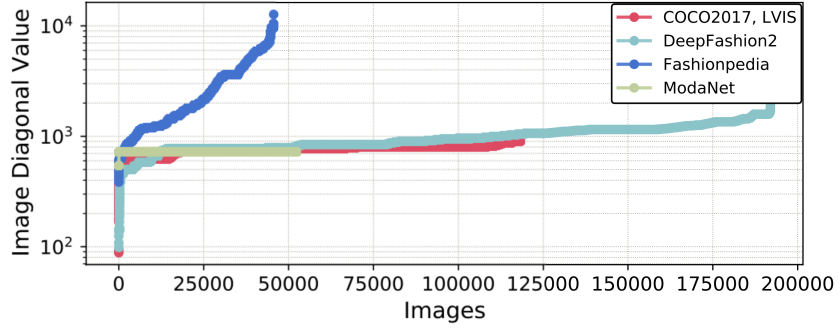


Fig. 12: Image size comparison among Fashionpedia, ModeNet, DeepFashion2, and COCO2017, LVIS. Only training images are shown. The Fashionpedia images has the most diverse resolutions. Note that COCO2017 and LVIS have higher resolution images for annotation. The distribution presented here are the publicly available photos

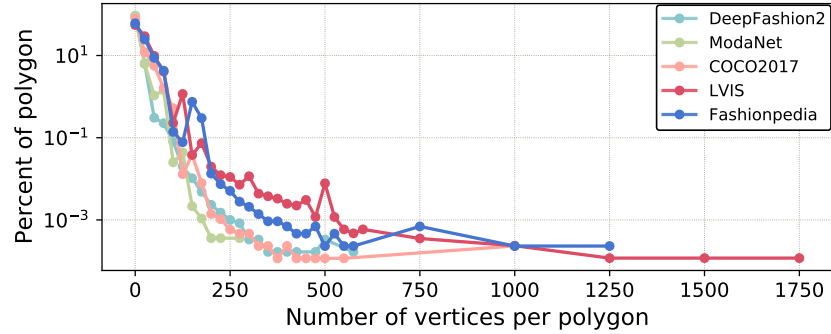


Fig. 13: The number of vertices per polygon. This represents the quality of masks and the efforts of annotators. Values in the x-axis were discretized for better visual effect. Y-axis is on log scale. Fashionpedia has the second widest range, next to LVIS



Fig. 15: Annotation tutorial example for shirt and top

the difference among different garment parts, such as ‘tassel’ and ‘fringe’. To help them understand the difference of these objects, we ask them to practice and identify more sample images before the annotation process. Fig. 16 shows our tutorials for these two categories. We specifically shows some correct and wrong examples of annotations.

Third, we ask for the quality of annotations. In particular, we ask the annotators to follow the contours of garments in the images as closely as possible. The polygon annotation process is monitored and verified for a few days before the workers started the actual annotation process.

Quality control of debatable apparel categories. During the annotation process, we allow annotators to ask questions about the uncertain categories. Two fashion experts monitored the annotation process by answering questions, checking the annotation quality, and providing weekly feedback to annotators.

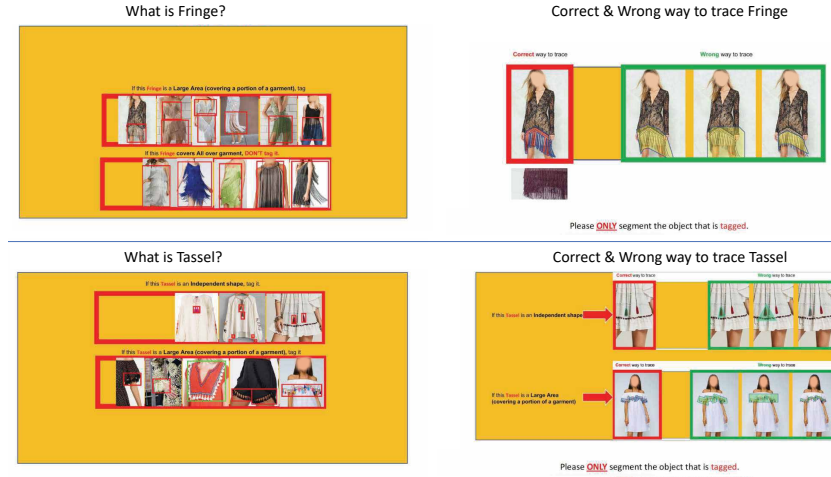


Fig. 16: Annotation tutorial for fringe and tassel

Instead of asking annotators to rate their confidence level of each segmentation mask, we asked them to send back all the uncertain masks to us during the annotation. The same two fashion experts made the final judgement and gave the feedback to the workers on these debatable or unsure fashion categories. Some examples of debatable or fuzzy fashion items that we have documented can be found in Figure 17.

1.5 Discussion

Does this dataset include the images or labels of previous datasets? We only include the previous datasets for comparison. Our dataset doesn't intentionally use any images or labels from previous datasets. All the images and labels from Fashionpedia are newly collected and annotated.

Who were the fashion experts annotating localized attributes in Fashionpedia dataset? The fashion experts are the 15 fashion graduate students that we recruited from one of the top fashion design institutes. For double-blind policy, we cannot mention the name of the university. But we will release the name of this university and the collaborators from this university in the final version of this paper.

Instance segmentation v.s. semantic segmentation. We didn't conduct semantic segmentation experiments on our dataset for the following two reasons: 1) Although semantic segmentation is a useful task, we believe instance segmentation is more meaningful for fashion images. For example, if we need to distinguish the different shoe style of a fashion image containing 3 pair of different shoes, instance segmentation (Figure 18(a)) can help us distinguish each shoe separately. However, semantic segmentation (Figure 18(b)) will mix all the

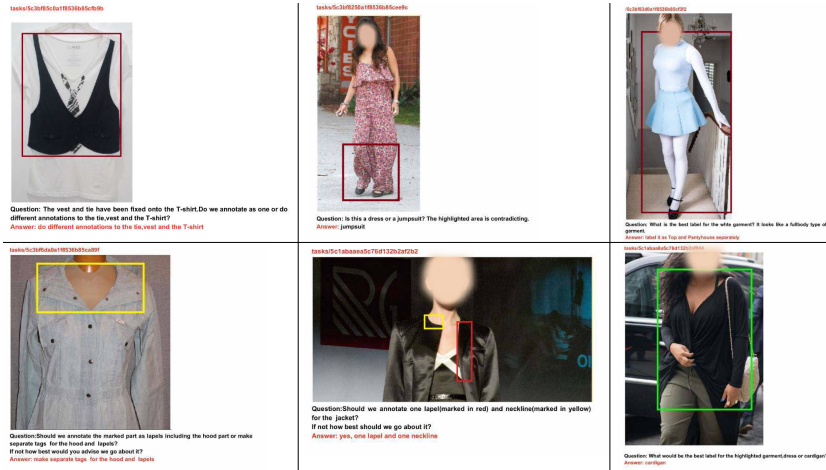


Fig. 17: Example of debatable fashion items in Fashionpedia dataset. The questions are asked by the crowdworkers. The answers are provided by two fashion experts



Fig. 18: Instance segmentation (left) and semantic segmentation (right)

shoe instances together. 2) Semantic segmentation is the sub-problem of instance segmentation. If we merge the same detected object class from our instance segmentation experimental result, it yields the results for semantic segmentation.

Which image has the most annotated masks? In Fashionpedia dataset, the maximum number of segmentation masks in an image is 74 (Fig. 19). We find that most of the masks are belonging to “rivets” (garment parts).

What’s the difference between Fashionpedia and other fine-grained datasets like CUB-200? We propose to localize fine-grained attributes within segmentation masks of images. This is a novel task with real-world application to the best of our knowledge. The differences between Fashionpedia and CUB are as follows: 1) CUB uses keypoints as annotation to indicate different locations on birds, while Fashionpedia has segmentation masks of garments, garment parts, and accessories; 2) Fashionpedia attributes are associated with garment or garment part instances in images, whereas CUB provides global attributes, not associated with any specific keypoints.

iMat-Fashion Kaggle challenges To advance state-of-the-art of visual analysis of clothing, we hosted two kaggle challenges (imaterialist-fashion) on Kaggle in 2019 ¹ and 2020 ² respectively.

¹ <https://www.kaggle.com/c/imaterialist-fashion-2019-FGVC6>

² <https://www.kaggle.com/c/imaterialist-fashion-2020-fgvc7>



Original image

Original image with masks

The detailed info of each mask with associated localized attributes:

<p>Segmentation 0: Category: pants Attributes: textile finishing, manufacturing techniques: no special manufacturing technique textile pattern: plain (pattern) waistline: normal waist non-textile material type: no non-textile material neckline: pos (point) length: maxi (length) silhouette: asymmetrical opening type: fly (opening) silhouette: peg</p> <p>Segmentation 1: Category: shoe Segmentation 2: Category: shoe Segmentation 3: Category: sleeve Attributes: length: wrist-length neckline: not-in sleeve Segmentation 4: Category: sleeve Attributes: length: wrist-length neckline: not-in sleeve Segmentation 5: Category: jacket Attributes: length: above-the-hip (length) textile pattern: plain (pattern) neckline: collar (point) silhouette: regular (fit) textile finishing, manufacturing techniques: lining waistline: no waistline opening type: zip-up Segmentation 6: Category: top, t-shirt, sweatshirt Attributes: length: above-the-hip (length) textile finishing, manufacturing techniques: no special manufacturing technique textile pattern: plain (pattern) non-textile material type: no non-textile material neckline: classic (t-shirt) silhouette: asymmetrical waistline: no waistline Segmentation 7: Category: zipper Segmentation 8: Category: pocket Attributes: neckline: collar (point) Segmentation 9: Category: zipper Segmentation 10: Category: zipper Segmentation 11: Category: zipper Segmentation 12: Category: zipper Segmentation 13: Category: collar Attributes: neckline: regular (collar)</p>	<p>Segmentation 14: Category: zipper Segmentation 15: Category: zipper Segmentation 16: Category: zipper Segmentation 17: Category: zipper Segmentation 18: Category: zipper Segmentation 19: Category: zipper Segmentation 20: Category: zipper Segmentation 21: Category: zipper Segmentation 22: Category: zipper Segmentation 23: Category: zipper Segmentation 24: Category: zipper Segmentation 25: Category: zipper Segmentation 26: Category: zipper Segmentation 27: Category: zipper Segmentation 28: Category: zipper Segmentation 29: Category: zipper Segmentation 30: Category: zipper Segmentation 31: Category: zipper Segmentation 32: Category: zipper Segmentation 33: Category: zipper Segmentation 34: Category: zipper Segmentation 35: Category: zipper Segmentation 36: Category: zipper Segmentation 37: Category: zipper Segmentation 38: Category: zipper Segmentation 39: Category: zipper Segmentation 40: Category: zipper Segmentation 41: Category: zipper Segmentation 42: Category: zipper Segmentation 43: Category: zipper Segmentation 44: Category: zipper Segmentation 45: Category: zipper</p>	<p>Segmentation 46: Category: zipper Segmentation 47: Category: zipper Segmentation 48: Category: zipper Segmentation 49: Category: zipper Segmentation 50: Category: zipper Segmentation 51: Category: zipper Segmentation 52: Category: zipper Segmentation 53: Category: zipper Segmentation 54: Category: zipper Segmentation 55: Category: zipper Segmentation 56: Category: zipper Segmentation 57: Category: zipper Segmentation 58: Category: zipper Segmentation 59: Category: zipper Segmentation 60: Category: zipper Segmentation 61: Category: zipper Segmentation 62: Category: zipper Segmentation 63: Category: zipper Segmentation 64: Category: zipper Segmentation 65: Category: zipper Segmentation 66: Category: zipper Segmentation 67: Category: zipper Segmentation 68: Category: zipper Segmentation 69: Category: zipper Segmentation 70: Category: zipper Segmentation 71: Category: zipper Segmentation 72: Category: zipper Segmentation 73: Category: zipper Segmentation 74: Category: zipper Segmentation 75: Category: zipper Segmentation 76: Category: zipper Segmentation 77: Category: zipper Segmentation 78: Category: zipper Segmentation 79: Category: zipper Segmentation 80: Category: zipper Segmentation 81: Category: zipper Segmentation 82: Category: zipper Segmentation 83: Category: zipper Segmentation 84: Category: zipper Segmentation 85: Category: zipper Segmentation 86: Category: zipper Segmentation 87: Category: zipper Segmentation 88: Category: zipper Segmentation 89: Category: zipper Segmentation 90: Category: zipper Segmentation 91: Category: zipper Segmentation 92: Category: zipper Segmentation 93: Category: zipper Segmentation 94: Category: zipper Segmentation 95: Category: zipper Segmentation 96: Category: zipper Segmentation 97: Category: zipper Segmentation 98: Category: zipper Segmentation 99: Category: zipper</p>
---	--	--

Fig. 19: The image with 74 masks in Fashionpedia dataset

References

1. Vrandečić, D., Krötzsch, M.: Wikidata: a free collaborative knowledge base (2014)
9
2. ZARA.com: Zara white leather flat ankle boots with top stitching size 5 bnwt,
retrieved May 9, 2019 from <https://www.zara.com> 6