# DeepMark++: Real-time Clothing Detection at the Edge

Alexey Sidnev[1,2], Alexander Krapivin[1], Alexey Trushkov[1], Ekaterina Krasikova[1], Maxim Kazakov[1,3], Mikhail Viryasov[1]

[1]Huawei Research Center, Nizhny Novgorod, Russia
[2]Lobachevsky State University of Nizhny Novgorod, Russia
[3]National Research University Higher School of Economics, Nizhny Novgorod, Russia

{sidnev.alexey, krapivin.alexander, trushkov.alexey, krasikova.ekaterina, kazakov.maxim, viryasov.mikhail}@huawei.com

## Abstract

*Clothing recognition is the most fundamental AI application challenge within the fashion domain. While existing solutions offer decent recognition accuracy, they are generally slow and require significant computational resources. In this paper we propose a single-stage approach to overcome this obstacle and deliver rapid clothing detection and keypoint estimation. Our solution is based on a multi-target network CenterNet [26], and we introduce several powerful post-processing techniques to enhance performance. Our most accurate model achieves results comparable to state-of-the-art solutions on the DeepFashion2 dataset [4], and our light and fast model runs at 17 FPS on the Huawei P40 Pro smartphone. In addition, we achieved second place in the DeepFashion2 Landmark Estimation Challenge 2020[1] with 0.582 mAP on the test dataset.*

## 1. Introduction

Image recognition is often the first step towards processing an image in many applications, and involves locating, classifying and, sometimes, identifying the objects presented in the image. An object's location can be estimated with different levels of precision (bounding boxes, keypoints, or segmentation masks), and each introduces a new challenge.

Object detection aims to locate an object by estimating its boundaries. A detected object is typically presented by a tight bounding box and a predicted class label. A bounding box can be used directly in applications such as pedestrian detection [24], or to crop an object from the original image for further processing.
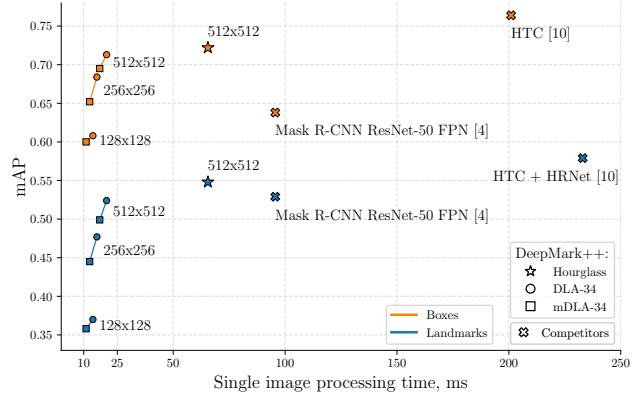


Figure 1. Speed/accuracy trade-off for object detection and landmark estimation on the DeepFashion2 validation dataset [4] using an RTX 2080ti. Mask R-CNN [4] time is estimated with Keypoint R-CNN model from Detectron2[2]. HTC [10] detection time is estimated with COCO detection model[3]. HTC + HRNet [10] is used without test-time augmentations in the single (aggregated) configuration.

Keypoint detection aims to estimate the precise location of an object's keypoints, which are sometimes referred to as landmarks. Given a set of correctly defined keypoints, one can transform an object in various ways – for example, position alignment for future processing. However, details of keypoint locations are not enough to estimate an object's boundaries due to a variety of object forms and its orientation in space. Instead, object detection and keypoint estimation are combined to provide enough information to locate and process objects.

The performance of models that operate on keypoints depends heavily on the number of unique keypoints defined in the task. The Deepfashion2 dataset [4], one of the largest and most diverse datasets in the fashion domain, provides annotation for 13 classes of clothing, each characterized by a certain set of keypoints (294 keypoints in total).

State-of-the-art methods used to solve keypoint estima-

---

[1]https://sites.google.com/view/cvcreative2020/deepfashion2

[2]https://github.com/facebookresearch/detectron2

[3]https://github.com/open-mmlab/mmdetection

tion problems on the DeepFashion2 dataset rely on heavy architectures such as Mask R-CNN [4] or HTC + HR-Net [10], which require large amounts of GPU resources and extended training time, and can not be considered to run on mobile devices. The recently proposed DeepMark [19] uses a one-stage anchor-free detection model called CenterNet [26] to simultaneously detect clothing items and estimate clothing landmark locations. While this approach is fast in the inference phase, it is still both time- and memory-consuming at the training stage. The massive number of unique keypoints in the DeepFashion2 dataset presents a significant bottleneck.

In this paper we propose a keypoint grouping strategy to drastically reduce the number of keypoints and to accelerate the training process. In addition, we propose and thoroughly study several post-processing techniques which can improve the model's accuracy without compromising its performance. Through this approach, we will demonstrate that our model is capable of achieving results comparable to the state-of-the-art.

Finally, we propose a light and fast version of our model capable of running on a mobile device. We will demonstrate optimal real-time performance of 50 FPS and near real-time performance of 17 FPS with an acceptable accuracy of 0.445 mAP for the keypoint estimation problem.

## 2. Related work

Object keypoints can typically be utilized in numerous applications. For example, they can be used to identify a human pose [16], locate facial landmarks [25], or to estimate 3D or 6D object locations [6]. However, the application we are most concerned with in this paper is clothing landmark estimation.

There are two open datasets capable of providing clothing images annotated with bounding boxes and keypoints. DeepFashion [13] defines between 4 and 8 keypoints for three classes of clothing items - top, bottom and full-body. Apart from a modest variety of clothing items and sparse keypoint definition, DeepFashion is limited to one object per image, making it less suitable for real-life applications. DeepFashion2 [4] improves on some of these drawbacks by providing annotations for 13 classes, each characterized by a unique set of keypoints, from 8 to 39, for 294 keypoints in total. Examples of keypoint definition are presented in Figure 3. The number of objects on an image ranges from 1 to 5.

Several papers are aiming to solve the landmark estimation challenge with regards to the DeepFashion2 dataset, with the original solution proposed in [4] built upon the Mask R-CNN architecture [5]. This solution uses the ResNet-FPN [11] backbone and RoIAlign to extract features from different levels of the pyramid map, which are then processed by two different branches, one each for ob-

ject detection and landmark estimation, in order to obtain final outputs.

In [2], authors introduce an attention module to aggregate features from different levels and to produce an output heatmap for landmark estimation. While this approach has delivered an improvement of 0.02 mAP for landmark estimation, it does not provide bounding boxes and instead offers only keypoint locations.

A recently proposed solution by Tzu-Heng Lin [10] offers current state-of-the-art capabilities with 0.614 mAP for landmark estimation and 0.764 mAP for object detection. This top-down approach first uses Hybrid Task Cascade [1] to perform object detection and obtain bounding boxes, and then runs HRNet-w48 [20] on a cropped image to estimate keypoint locations. This method also uses the aggregation strategy to reduce the number of keypoints and finetune the landmark estimation model for each class separately, thereby enhancing performance.

As existing solutions rely on slow and resource-consuming architecture, they are unsuitable for low-power devices. At the same time, keypoint-based object detection methods have gained significant popularity in recent papers as they are simpler, faster, and more accurate than the corresponding anchor-based detectors.

Previous approaches, such as [12, 17], required anchor boxes to be manually designed in order to train detectors. A series of anchor-free object detectors were then developed, aiming to predict the bounding box keypoints, rather than trying to fit an object to an anchor. Law and Deng proposed CornerNet [9], a novel anchor-free framework which detected objects like a pair of corners, and predicted class heatmaps, pair embeddings, and corner offsets on each position of the feature map. Class heatmaps calculated the probability of a corner, corner offsets were used to regress the corner location, and the pair embeddings grouped a pair of corners that belong to the same object. Without relying on manually designed anchors to match objects, CornerNet achieved a significant performance improvement on the MS COCO datasets. Subsequently, several other variants of keypoint detection-based one-stage detectors have been developed [21, 27].

In [19], authors naturally treat landmark estimation as a pose estimation problem, and use keypoint-based multi-target architecture CenterNet [26] with slight modifications to handle both object detection and landmark estimation. We further develop that idea in this paper, and propose several improvements to accelerate the training process and enhance performance.

Despite state-of-the-art performance in object detection and keypoint estimation tasks, CenterNet's approach can produce overconfident incorrect predictions, leading to a decrease in performance metrics. A separate line of research exists in the field of deep learning which aims to
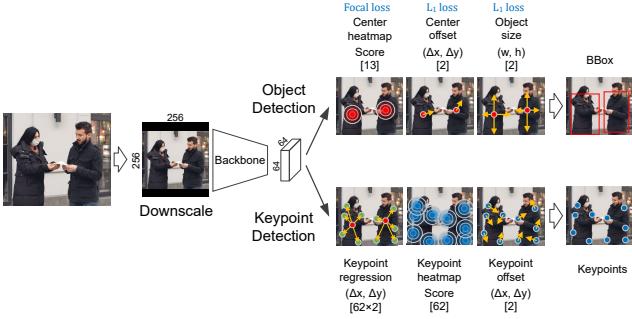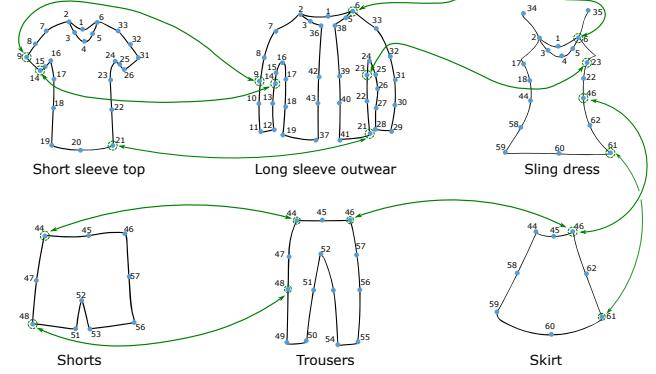
Figure 2. Scheme of the proposed approach.



Figure 3. The semantic grouping approach. A number indicates the group a keypoint belongs to, and the green arrows are examples of merged keypoints.
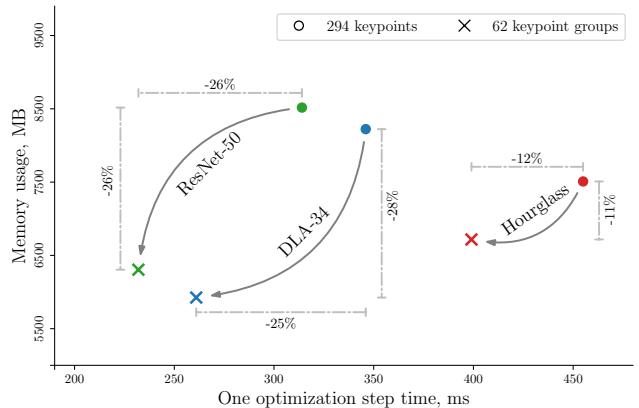


Figure 4. GPU memory consumption and training iteration time using an RTX 2080ti and PyTorch 1.4. The input resolution is 256×256, while the batch size is 32 for both DLA-34 and ResNet-50, and 8 for Hourglass. The time in ms was measured for 1 optimization step: batch loading to GPU, forward pass and backward pass. GPU memory was measured using the nvidia-smi tool.

solve the problem of estimating the uncertainty of model predictions. In this regard, a number of papers show how the use of various techniques, aimed at obtaining a more accurate estimate of the prediction uncertainty, can improve the quality of detection [14, 8, 22]. The most popular method for estimating model uncertainty, MC Dropout, requires several forward passes through a network to be performed, significantly increasing the inference time. In [22], a sample-free method for uncertainty estimation in anchor-based detectors was presented which uses a set of observations that relate to the same object, according to Jaccard overlap (before applying NMS). Motivated by this idea, we created a way to improve CenterNet's performance by considering a set of predictions related to a single object instead of point estimation.

## 3. Proposed approach

Our approach is based on the CenterNet [26] architecture, and simultaneously solves the following two tasks: object detection and keypoint location estimation. The architecture is presented in Figure 2. A backbone downsamples the input image by a factor of 4 to produce a feature map, which is then processed to obtain object bounding boxes and corresponding keypoints.

As the DeepFashion2 dataset contains 13 classes, 13 channels are used in the center heatmap to predict the probabilities for each pixel being the object center for each class. The center of an object is defined as the center of a bounding box, and a ground truth heatmap is generated by applying a Gaussian function at each object's center. Two additional channels in the output feature map – $\triangle x$ and $\triangle y$ – are used to refine the center coordinates, and both width and height are predicted directly.

Another branch handles the fashion landmark estimation task, which involves estimating 2D keypoint locations for each item of clothing in one image. The coarse locations of the keypoints are regressed as relative displacements from the box center (keypoint regression in Figure 2). Consequently, if a certain pixel has already been classified as an object center, one can take the values in the same spatial location from this tensor and interpret them as vectors to keypoints. The keypoint position obtained through regression is not entirely accurate, so an additional heatmap with probabilities is used for each keypoint type to refine a keypoint location. Here, a local maximum with high confidence in the heatmap is used as a refined keypoint position. Similar to the detection branch, two additional channels – $\triangle x$ and $\triangle y$ – are used to obtain more precise landmark coordinates. During model inference, each coarse keypoint location is replaced with the closest refined keypoint position. As such, keypoints that belong to the same object are grouped together.

### 3.1. Semantic keypoint grouping

One of the first steps involved in solving keypoint detection tasks is defining the model output. The number of key-

| Backbone | $mAP_{pt}$ | | $mAP_{box}$ | |
|---|---|---|---|---|
| | 294 kps | 62 kps | 294 kps | 62 kps |
| ResNet-50 | 0.377 | 0.365 | 0.629 | 0.627 |
| mDLA-34 | 0.434 | 0.422 | 0.649 | 0.651 |
| DLA-34 | 0.466 | 0.456 | 0.679 | 0.679 |

Table 1. The detection accuracy for different backbones with (62 kps) and without (294 kps) semantic keypoint grouping.

points for every category varies from 8 for a skirt, to 39 for a long sleeve outerwear in the DeepFashion2 dataset, and the total number of unique keypoints is 294. The straightforward approach is to concatenate keypoints from every category and deal with each one separately. Directly predicting 294 keypoints leads to a huge number of output channels: $901 = 13 + 2 + 2 + 294 \cdot 2 + 294 + 2$ (13 – center heatmap, 2 – center offset, 2 – object size, $294 \cdot 2$ – keypoint regression, 294 – keypoint heatmap, 2 – keypoint offset).

It is evident that certain clothing landmarks are actually a subset of others. For example, shorts do not require unique keypoints as they can be represented by a subset of trouser keypoints. The semantic grouping rule is defined as follows (Figure 3): identical semantic meaning keypoints (such as collar center and top sleeve edge) with different categories can be merged into one group. This approach enables the formation of 62 groups and reduces the number of output channels from 901 to $205 = 13 + 2 + 2 + 62 \cdot 2 + 62 + 2$ (see Figure 2 for details).

The semantic grouping approach reduces training time and memory consumption by up to 26% and 28%, respectively, without notable decrease in accuracy (see Figure 4 and Table 1). In addition, the latter reduction enables the use of larger batches during model training.

The semantic grouping approach also reduces inference time and memory consumption by up to 26% and 37%, respectively (see section 4.3 for details). An in-depth analysis of the grouping approach is presented in [18].

## 3.2. Post-processing techniques

We have developed 4 post-processing techniques that increase the model's accuracy without compromising performance.

### 3.2.1 Center rescoring

It is evident that keypoint scores for the wrong category may have low confidence (see Figure 5).

We use this insight to recalculate the detection confidence score using keypoint scores from the keypoint heatmap. Let $Score_{bbox}$ be the original detection confidence score from the center heatmap, and $Score_{kps}$ be the average score of the refined keypoints for the predicted category from the keypoint heatmap. The final detection con-



Figure 5. Detection results for two different classes. The keypoint score for the right class has a higher value. Short sleeved outwear on the left has $Score_{bbox} = 0.31$, $Score_{kps} = 0.28$. Long sleeved outwear on the right has $Score_{bbox} = 0.29$, $Score_{kps} = 0.36$.
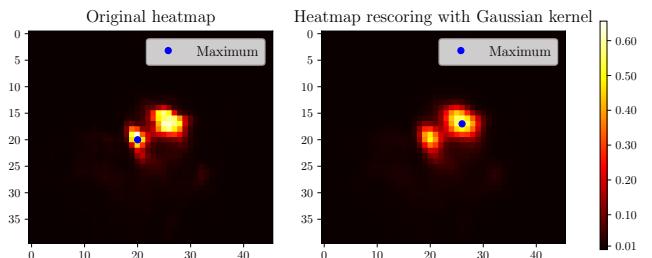


Figure 6. A heatmap with center scores and a heatmap after rescoring with Gaussian kernel ($\sigma = 0.45$).

fidence scores are calculated as a linear combination of the bounding box score and the average keypoint score:

$$Score = \alpha \cdot Score_{bbox} + (1 - \alpha) \cdot Score_{kps}, \quad (1)$$

where $\alpha \in [0, 1]$.

### 3.2.2 Heatmap rescoring with Gaussian kernel

Applying smoothing operators on heatmaps is quite a common technique. For instance, it is used in CenterNet to generate ground truth labels for a heatmap during the training stage. We propose a general approach that can be applied to any keypoint-based architecture during the inference stage as well.

Let $H$ be a heatmap with center or keypoint scores. Taking the training procedure into account, you can expect the 8-connected neighbors of each item to be related to the same object, and this fact can be used to improve the estimation of each heatmap value (see Figure 6). Consequently, we have applied the following formula:

$$\hat{H} = H \circledast G(\sigma), \quad (2)$$

where $\circledast$ is the convolution operation, and $G(\sigma)$ is the $3 \times 3$ Gaussian kernel with the standard deviation $\sigma$. Experimental results show that in our model, the proposed technique
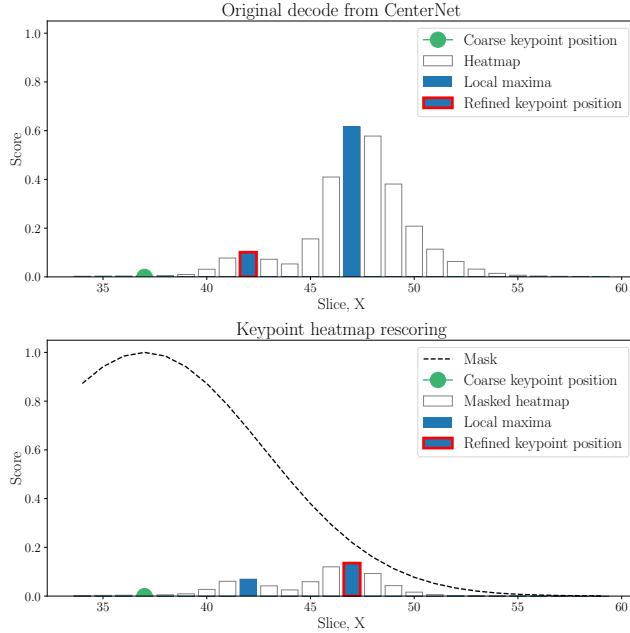
4

Figure 7. In CenterNet, the refined location is defined as the closest point to the coarse keypoint position, which in turn is defined as a local maximum on the heatmap. This can lead to errors when the closest prediction is not actually the best one (as evident in the diagram above). In the proposed technique, the refined location is determined as a global maximum of the keypoint heatmap rescored by the Gaussian mask, improving keypoint localization.

improves the localization of peaks and their values that correspond to object centers or keypoints and their scores. A similar operation has been considered in [23] as a part of the proposed method.

### 3.2.3 Keypoint location refinement

The third technique further corrects keypoint locations by calculating a linear combination of the refined and coarse keypoint positions (see Figure 8). Our experiments have shown a small improvement when using this technique, with no impact on performance.

Let $(x, y)_{refined}$ be the refined keypoint location from the heatmap and $(x, y)_{coarse}$ be coarse positions predicted as offsets from object centers. The final keypoint locations are calculated through the following expression:

$$(x, y) = \gamma \cdot (x, y)_{refined} + (1 - \gamma) \cdot (x, y)_{coarse}, \quad (3)$$

where $\gamma \in [0, 1]$.

### 3.2.4 Keypoint heatmap rescoring

In the fourth technique we propose the addition of a penalty to the keypoint score in proportion to the distance from a coarse keypoint position with the Gaussian function. Let
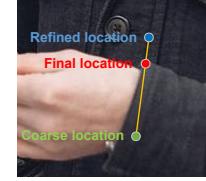


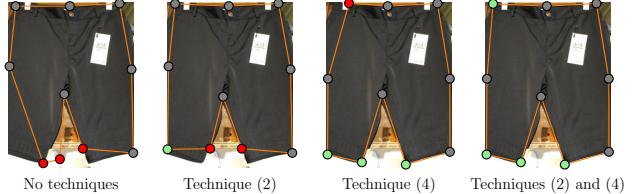Figure 8. A linear combination of the refined and coarse keypoint positions.



Figure 9. The impact of post-processing techniques on keypoint detection accuracy. Red dots highlight incorrect detections, and green dots highlight fixes by post-processing techniques.

$mask$ be a heatmap with zero values by default. We set 1 into the $mask$ in the coarse keypoint position and fill neighbor values with 2D Gaussian function with standard deviation $sigma = \min(width, height)/9$, where $width$ and $height$ are the object size (see the second image in Figure 7).

The keypoint heatmap is rescored through the following expression:

$$\hat{H}_{kps} = H_{kps} \cdot mask. \quad (4)$$

## 4. Results

All experiments were performed on the publicly available DeepFashion2 Challenge dataset [4], which contains 191,961 images in the training set and 32,153 images in the validation set.

### 4.1. Post-processing techniques

We considered 5 fast post-processing techniques: bounding box non-maximum suppression and 4 techniques from section 3.2. The individual and combined effectiveness of each technique is shown in Table 2 and Figure 9.

Certain techniques can increase $mAP_{pt}$ and reduce $mAP_{box}$ simultaneously. Note that bounding box detection and keypoint estimation results for the same object may have different IoU and OKS with the ground truth, for example, when bounding box was detected correctly, but keypoints were not. In this case, technique (1) involves lowering a score for false positive keypoints, which is advisable. The corresponding true positive bounding box also suffers from this lowered score.

We also evaluated proposed techniques on the COCO keypoint detection task. Here, we took the DLA-34 model for human keypoint detection on the COCO dataset with
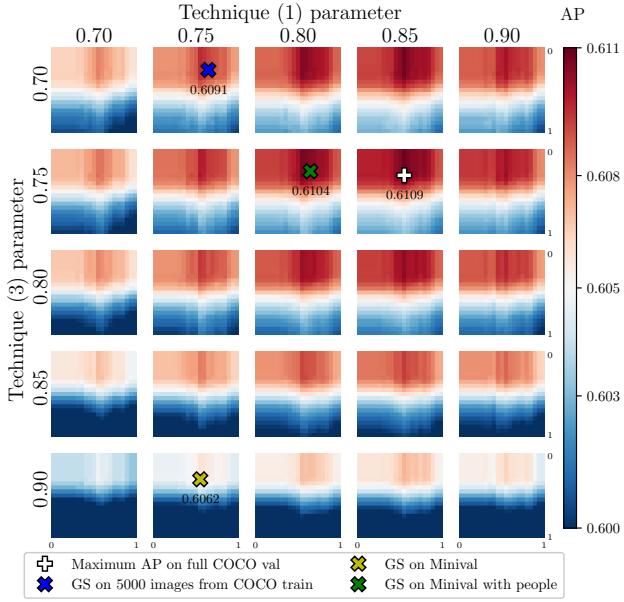
Figure 10. AP on the COCO validation dataset with different post-processing technique parameters. Every heatmap shows AP for technique (2) with step 0.05 in a range from 0 to 1 (center heatmap $\sigma$ goes on the X-axis and keypoint heatmap $\sigma$ goes on the Y-axis). Optimal values found with grid search (GS) on different datasets are shown.

| Post-processing | $mAP_{pt}$ | $mAP_{box}$ | Total time, ms |
|---|---|---|---|
| Baseline | 0.422 | 0.651 | 54.50 |
| NMS | 0.422 | 0.652 | 54.57 |
| Technique (1) | 0.428 | 0.650 | 54.61 |
| Technique (2) | 0.430 | 0.651 | 54.44 |
| Technique (3) | 0.431 | 0.651 | 54.81 |
| Technique (4) | 0.430 | 0.651 | 55.95 |
| Combined | 0.445 | 0.652 | 56.76 |

Table 2. Different post-processing techniques applied independently to mDLA-34 $256\times256$ on the DeepFashion2 validation dataset. The technique numbers correspond to the numbers in section 3.2. Total processing time is measured on a Huawei P40 Pro. Technique (2) works faster than the base version, as it significantly reduces the number of local maxima and increases the speed of further decoding. The last line contains metrics for combined techniques.

| Approach | $mAP_{pt}$ | $mAP_{box}$ |
|---|---|---|
| Mask R-CNN [4] | 0.529 | 0.638 |
| DeepMark [19] | 0.532 | 0.723 |
| DAFE [2] | 0.549 | — |
| Aggregation and Finetuning [10] | 0.614 | 0.764 |
| Hourglass $512\times512$ | 0.583 | 0.735 |
| Hourglass $768\times768$ | 0.591 | 0.737 |

Table 3. Accuracy comparison between the proposed and alternative approaches with the DeepFashion2 validation dataset. All models are trained with the released DeepFashion2 Challenge dataset.

0.589 AP from [26], and applied 4 techniques from section 3.2. Three datasets were selected for hyperparameter optimization: 5000 random images from the COCO train dataset; a subset of 100 images from the COCO validation dataset, of which only 61 images contain people (Minival); a subset of 100 images with people from the COCO validation dataset (Minival with people).

Each post-processing technique parameter was determined through grid searching with step 0.05 within the range 0 to 1, resulting in 21 experiments for a single parameter. Grid search on a representative subset enables accuracy very close to the global maximum (see Figure 10), and the accuracy was increased by 0.021 AP (from 0.589 AP to 0.6104 AP) without model retraining. On Minival with people we determined the following optimal parameters: $\alpha = 0.75$, center heatmap $\sigma = 0.25$, keypoint heatmap $\sigma = 0.65$, $\gamma = 0.8$.

## 4.2. DeepFashion2 Challenge

We used the CenterNet MS COCO model for object detection as the initial checkpoint and performed experiments with the Hourglass backbone and Adam the optimizer to achieve the best results (Table 3) for object detection and keypoint estimation tasks on the DeepFashion2 validation dataset. The Hourglass $512\times512$ model was trained for 100 epochs with a batch size of 46 images. The learning rate schedule is: 1e-3 - 65 epochs, 4e-4 - 20 epochs, 4e-5 -

15 epochs. The Hourglass $768\times768$ model was fine-tuned from Hourglass $512\times512$ for 25 epochs, with a batch size of 22 images: 2e-5 - 20 epochs, 1e-5 - 5 epochs.

To optimize post-processing technique parameters with grid search (subsection 4.1) we defined a small validation subset (1285 images) from the DeepFashion2 validation dataset. The following parameters were used for the Hourglass $768\times768$: $\alpha = 0.8$, center heatmap $\sigma = 0.45$, keypoint heatmap $\sigma = 0.90$, $\gamma = 0.75$.

With the Hourglass $768\times768$ model we achieved the second place in the DeepFashion2 Challenge 2020 with 0.582 $mAP$ on the test dataset.

During all experiments, our target was to increase the keypoint estimation accuracy instead of the object detection accuracy. Consequently, object detection increased by only 0.042 $mAP_{box}$, but all the techniques added more than 0.07 to $mAP_{pt}$ (see Table 4).

### 4.2.1 Multi-inference strategies

We have considered 2 extra inference strategies: fusing model outputs from original and flipped images with equal
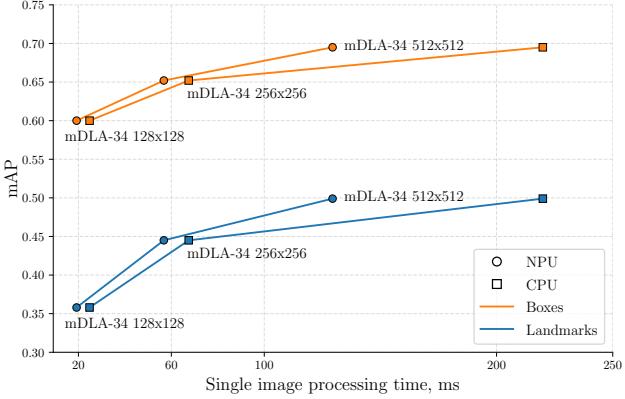
Figure 11. Speed/accuracy trade-off for object detection and landmark estimation on the DeepFashion2 validation dataset. Total processing time (pre-procesing + inference + post-processing) in ms was measured on a Huawei P40 Pro with MNN on CPU and with HiAI DDK on NPU.

| Metric | $mAP_{pt}$ | | $mAP_{box}$ | |
|---|---|---|---|---|
| Resolution | 512×512 | 768×768 | 512×512 | 768×768 |
| Baseline | 0.529 | 0.520 | 0.720 | 0.695 |
| + Post-processing | 0.545 | 0.540 | 0.713 | 0.698 |
| + NMS | 0.548 | 0.549 | 0.718 | 0.712 |
| + Flip | 0.561 | 0.563 | 0.731 | 0.727 |
| + Multiscale | 0.568 | 0.578 | **0.735** | **0.737** |
| + PoseFix | **0.583** | **0.591** | **0.735** | **0.737** |

Table 4. Clothing detection and landmark estimation. Hourglass with 512×512 and 768×768 resolution were used in the experiments. The next technique is added to each of the previous ones, and the post-processing refers to the application of all techniques from section 3.2. We applied the following multipliers for the multiscale technique: 0.85, 0.95, 1.1. Our target was to increase the keypoint estimation accuracy (also consider the note from 4.1), so $mAP_{box}$ accuracy is dropped after applying some of our techinques.

weights; and fusing model results with the original image downscaled/upscaled through certain multipliers. While these proposed techniques increase accuracy, they require several model inferences, and this significantly impacts processing time.

#### 4.2.2 Keypoint Refinement Network

At the final stage, detection results are refined with the PoseFix [15] model-agnostic pose refinement method, which learns the typical error distributions of any other pose estimation method and corrects mistakes at the testing stage.

We trained a set of 13 PoseFix models using the number of classes in the DeepFashion2 dataset, and the inference results of our method on the training set are used to train each of the 13 models. Subsequently, we applied the trained
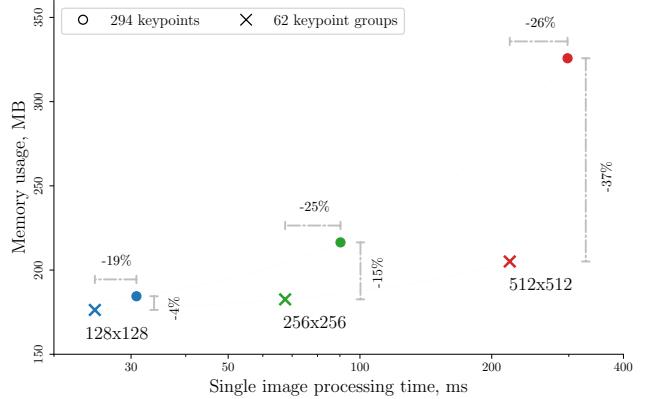


Figure 12. Single image processing time and memory consumption with mDLA-34 on a Huawei P40 Pro with MNN. Total processing time (pre-procesing + inference + post-processing) and memory usage during inference are shown for models with direct prediction of 294 keypoints and a prediction of 62 keypoint groups.

PoseFix models to the result.

### 4.3. Experiments on a mobile phone

The DLA-34 model contains DCN [3] layers by default. This increases model accuracy by 2-3 mAP, but most inference frameworks do not currently support DCN. In order to run the DLA-34 model on a mobile phone, we replaced all DCN layers with conventional convolution layers $3 \times 3$. This model is referred to as mDLA-34, and it contains 19.5 millions parameters. The number of floating point operations for the model is presented in Table 6.

The Huawei P40 Pro smartphone (complete with Kirin 990) is used for our experiments. To run the model on CPU we use MNN [7] v1.0.2 (number of threads = 4), and on NPU we use HiAI DDK v100.310.011. Both NPU and CPU can run mDLA-34 at 14 FPS with relatively high accuracy (see Figure 11).

We also estimated the influence on the keypoint grouping approach to memory consumption and execution time (see Figure 12). Here, the grouping approach only impacts the last part of a neural network and a post-processing step. The higher the input resolution, the greater the performance and memory gains. The keypoint grouping approach works 25% faster and requires 15% less memory with the mDLA-34 $256 \times 256$ model on a Huawei P40 Pro, compared to the same model without keypoint grouping.

mDLA-34 $256 \times 256$ performance on the different validation subsets and for each class is presented in Table 5 and Table 7. Figure 13 demonstrates examples of the model performance on several use-cases.

7

|  | (a) | (b) | (c) | (d) |

Figure 13. Examples of fashion landmark detection with mDLA-34 $256\times256$ with all post-processing techniques. A simple pose without any occlusion is not a major challenge for even the $256\times256$ model. Keypoints are detected quite accurately for images featuring a single person (a) or multiple people (c). However, uncommon hand positions (b) can lead to errors in the elbow area. As the subject's legs are not fully visible in image (b), this leads to both erroneous determination of the trouser type and errors in the keypoint positions. In addition, clothing detection on distant or blurred subjects (d) works poorly in some cases.

|  | Scale | | | Occlusion | | | Zoom-in | | | Viewpoint | | | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | Small | Moderate | Large | Slight | Medium | Heavy | No | Medium | Large | No wear | Frontal | Side or back | |
| $mAP_{pt}$ | 0.408 | 0.560 | 0.467 | 0.509 | 0.502 | 0.209 | 0.536 | 0.440 | 0.315 | 0.431 | 0.550 | 0.409 | 0.507 |
|  | 0.346 | 0.497 | 0.434 | 0.490 | 0.421 | 0.132 | 0.473 | 0.393 | 0.295 | 0.418 | 0.482 | 0.345 | 0.445 |
| $mAP_{box}$ | 0.574 | 0.683 | 0.667 | 0.691 | 0.636 | 0.285 | 0.663 | 0.641 | 0.592 | 0.669 | 0.669 | 0.614 | 0.652 |

Table 5. The results for landmark estimation with mDLA-34 $256\times256$ with all post-processing techniques on different validation subsets, including scale, occlusion, zoom-in, and viewpoint. Results of evaluation on visible landmarks only and evaluation on both visible and occluded landmarks are separately shown in each row for $mAP_{pt}$.

| Resolution | Model inference | Post-processing techniques | Decode | Total |
|---|---|---|---|---|
| $128 \times 128$ | 2.423 | 0.0008 | 0.0004 | 2.424 |
| $256 \times 256$ | 9.691 | 0.0028 | 0.0015 | 9.696 |
| $512 \times 512$ | 38.765 | 0.0113 | 0.0058 | 38.783 |

Table 6. The number of floating point operations in GFLOPs for mDLA-34 inference, post-processing techniques, and decode.

| Class | $AP_{pt}$ | $AP_{box}$ |
|---|---|---|
| Short sleeve top | 0.606 | 0.804 |
| Long sleeve top | 0.507 | 0.724 |
| Short sleeve outwear | 0.175 | 0.347 |
| Long sleeve outwear | 0.458 | 0.724 |
| Vest | 0.426 | 0.679 |
| Sling | 0.297 | 0.422 |
| Shorts | 0.545 | 0.721 |
| Trousers | 0.443 | 0.739 |
| Skirt | 0.481 | 0.740 |
| Short sleeve dress | 0.586 | 0.721 |
| Long sleeve dress | 0.375 | 0.542 |
| Vest dress | 0.481 | 0.710 |
| Sling dress | 0.400 | 0.605 |

Table 7. Accuracy by class for mDLA-34 $256\times256$ with all post-processing techniques. Accuracy for some classes is significantly lower than for others due to the fact that these classes contain smaller number of samples in train and validation datasets. For short sleeve outwear class it's 543 and 142 samples in train and validation subsets respectively, for sling class – 1985 and 322. For comparison, vest dress contains 17949 and 3352 samples.

## 5. Conclusion

This new approach is proposed as an adaptation of CenterNet [26] for clothing landmark estimation tasks. The accuracy is comparable to state-of-the-art was achieved on the DeepFashion2 dataset by applying several post-processing techniques as well as the semantic keypoint grouping approach: clothing detection hit 0.737 $mAP$ and clothing landmark estimation reached 0.591 $mAP$. Performance runs at 56 ms per image (more than 17 FPS) on a Huawei P40 Pro for mDLA-34 $256\times256$, and yields considerably high accuracy (0.445 $mAP_{pt}$ and 0.652 $mAP_{box}$) for clothing detection tasks.

# References

[1] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4974–4983, 2019.

[2] Ming Chen, Yingjie Qin, Lizhe Qi, and Yunquan Sun. Improving fashion landmark detection by dual attention feature enhancement. In *The IEEE International Conference on Computer Vision (ICCV) Workshops*, Oct 2019.

[3] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. *CoRR*, abs/1703.06211, 2017.

[4] Yuying Ge, Ruimao Zhang, Xiaogang Wang, Xiaoou Tang, and Ping Luo. DeepFashion2: A versatile benchmark for detection, pose estimation, segmentation and re-identification of clothing images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5337–5345, 2019.

[5] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. *CoRR*, abs/1703.06870, 2017.

[6] Yisheng He, Wei Sun, Haibin Huang, Jianran Liu, Haoqiang Fan, and Jian Sun. Pvn3d: A deep point-wise 3d keypoints voting network for 6dof pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[7] Xiaotang Jiang, Huan Wang, Yiliu Chen, Ziqi Wu, Lichuan Wang, Bin Zou, Yafeng Yang, Zongyang Cui, Yu Cai, Tianhang Yu, Chengfei Lv, and Zhihua Wu. Mnn: A universal and efficient inference engine. In *MLSys*, 2020.

[8] Florian Kraus and Klaus Dietmayer. Uncertainty estimation in one-stage object detection. *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, Oct 2019.

[9] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 734–750, 2018.

[10] Tzu-Heng Lin. Aggregation and finetuning for clothes landmark detection. *ArXiv*, abs/2005.00419, 2020.

[11] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017.

[12] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.

[13] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1096–1104, 2016.

[14] Dimity Miller, Lachlan Nicholson, Feras Dayoub, and Niko Sünderhauf. Dropout sampling for robust object detection in open-set conditions, 2017.

[15] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. Posefix: Model-agnostic general human pose refinement network, 2018.

[16] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation, 2016.

[17] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.

[18] Alexey Sidnev, Ekaterina Krasikova, and Maxim Kazakov. Efficient grouping for keypoint detection. *CoRR*, abs/2010.12390, 2020.

[19] Alexey Sidnev, Alexey Trushkov, Maxim Kazakov, Ivan Korolev, and Vladislav Sorokin. Deepmark: One-shot clothing detection. In *The IEEE International Conference on Computer Vision (ICCV) Workshops*, Oct 2019.

[20] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5693–5703, 2019.

[21] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[22] Michael Truong-Le, Frederik Diehl, Thomas Brunner, and Alois Knoll. Uncertainty estimation for deep neural object detectors in safety-critical applications. *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pages 3873–3878, 2018.

[23] Feng Zhang, Xiatian Zhu, Hanbin Dai, Mao Ye, and Ce Zhu. Distribution-aware coordinate representation for human pose estimation. *arXiv preprint arXiv:1910.06278*, 2019.

[24] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z. Li. Occlusion-aware r-cnn: Detecting pedestrians in a crowd. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.

[25] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Facial landmark detection by deep multi-task learning. In *European conference on computer vision*, pages 94–108. Springer, 2014.

[26] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019.

[27] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.