# Towards High Performance Human Keypoint Detection

**Jing Zhang** · **Zhe Chen** · **Dacheng Tao**

**Abstract** Human keypoint detection from a single image is very challenging due to occlusion, blur, illumination, and scale variance. In this paper, we address this problem from three aspects by devising an efficient network structure, proposing three effective training strategies, and exploiting four useful postprocessing techniques. First, we find that context information plays an important role in reasoning human body configuration and invisible keypoints. Inspired by this, we propose a cascaded context mixer (CCM), which efficiently integrates spatial and channel context information and progressively refines them. Then, to maximize CCM's representation capability, we develop a hard-negative person detection mining strategy and a joint-training strategy by exploiting abundant unlabeled data. It enables CCM to learn discriminative features from massive diverse poses. Third, we present several sub-pixel refinement techniques for postprocessing keypoint predictions to improve detection accuracy. Extensive experiments on the MS COCO keypoint detection benchmark demonstrate the superiority of the proposed method over representative state-of-the-art (SOTA) methods. Our single model achieves comparable performance with the winner of the 2018 COCO Keypoint Detection Challenge. The final ensemble model sets a new SOTA on this benchmark. The source code will be released at `https://github.com/chaimi2013/CCM`.

✉ Zhe Chen (zhe.chen1@sydney.edu.au)

✉ Dacheng Tao (dacheng.tao@sydney.edu.au)

Authors are with School of Computer Science, Faculty of Engineering, The University of Sydney, Darlington, NSW 2008, Australia
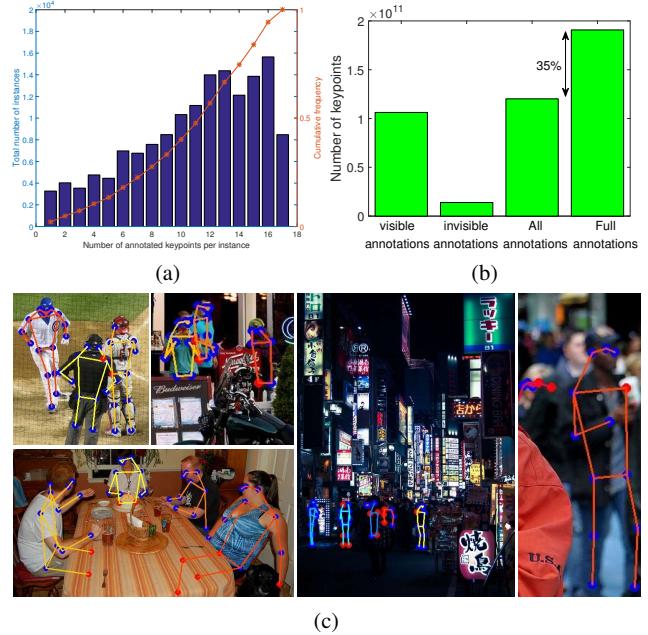
(a)          (b)

(c)

Fig. 1: (a)-(b) Statistics of the annotated keypoints in the MS COCO training dataset (Lin et al., 2014). (c) Some examples from the MS COCO dataset, where occluded, underexposed, and blurry person instances are very common. Blue and red dots denote the annotated visible and invisible keypoints, respectively.

## 1 Introduction

Human keypoint detection is also known as human pose estimation (HPE) refers to detecting human body keypoint location and recognizing their categories for each person instance from a given image. It is very useful in many downstream applications such as activity recognition (Ni et al., 2017; Baradel et al., 2018; Liu et al., 2018), human-robot

interaction (Mazhar et al., 2018; Zhang and Tao, 2020), and video surveillance (Hattori et al., 2018; Varadarajan et al., 2018). However, HPE is very challenging even for human annotators. For example, almost 50% of instances[1] in the COCO training dataset (Lin et al., 2014) have at least six unannotated keypoints (Figure 1(a)) and 35% keypoints are unannotated (Figure 1(b)) due to various factors including occlusion, truncation, under-exposed imaging, blurry appearance and low-resolution of person instances. Some examples from the MS COCO dataset are shown in Figure 1.

Prior methods have made significant progress in this area with the success of deep convolutional neural networks (DCNN) (Toshev and Szegedy, 2014), which can be categorized into two groups: top-down methods and bottom-up methods. Top-down methods are composed of two phases, i.e., person detection and keypoint detection, while bottom-up methods directly detect all keypoints from the image and associate them with corresponding person instances. Although bottom-up methods (Cao et al., 2017) are usually fast, top-down methods still dominate the leaderboard of public benchmark datasets like MS COCO due to their high accuracy. Using heatmaps to represent keypoint locations and fully CNNs with an encoder-decoder structure to learn features has gained prominence in recent studies (Newell et al., 2016; Huang et al., 2017), since spatial correspondence can be preserved between feature maps and heatmaps. Recently, to learn strong feature representations at multiple resolutions, pyramid-based networks have been proposed (Yang et al., 2017; Chen et al., 2018b). Although these methods improve detection accuracy by learning multi-scale features, there are still some issues to be addressed.

First, detecting invisible keypoints is more difficult than visible ones due to ambiguous appearance and inconsistent context bodies. How to effectively model multi-source context information to infer the hard keypoints is still under-explored. Second, external datasets such as the AI Challenger dataset[2], have been used to learn more discriminative feature representations (Xiao et al., 2018; Sun et al., 2019). However, they may have different annotation formats with the target set, for example, 17 keypoints for the MS COCO dataset and 14 keypoints for the AI Challenger dataset, or even have no annotations like the MS COCO unlabeled dataset. How to effectively leverage these external datasets to learn human pose configuration and discriminative feature representations for recognizing diverse poses remains challenging. Third, the ground truth keypoint location is annotated in pixel (or sub-pixel) in the high-resolution image plane, while the regression target is usually in the low-resolution heatmap, e.g., 1/4 size of the input image. This scale mismatch of representation will degrade the keypoint detection performance. To address this issue, sub-pixel

representation or post-processing techniques (Chen et al., 2018b; Zhang et al., 2020) have been proposed. Nonetheless, a systematic study of post-processing techniques is still absent and worth further discovering.

In this paper, we address these issues to improve human keypoint detection by devising an efficient network structure, proposing three effective training strategies, and exploiting four useful postprocessing techniques. First, inspired by the keypoint detection process carried out by humans, where the "context" contributes to the perception and inference process, we advance the research by studying the role of context information for human keypoint detection. Specifically, we adopt an encoder-decoder network structure and propose a novel cascaded context mixer (CCM) in the decoder. It can efficiently integrate both spatial and channel context information and progressively refine them. Then, we propose a joint training strategy and a knowledge-distilling approach to exploit abundant unlabeled data. Besides, we also propose a hard-negative person detection mining strategy to migrate the inconsistency of person instances between training and testing. These strategies endow the detection network with the capability of learning discriminative features. Third, we present and comprehensively evaluate four sub-pixel refinement techniques for postprocessing keypoint predictions. Extensive experiments on the MS COCO keypoint detection benchmark validate the effectiveness of the proposed CCM model, the training strategies, and the sub-pixel techniques. Our single model achieves comparable performance with the winner of the 2018 COCO Keypoint Detection Challenge. The final ensemble model sets a new SOTA on this benchmark[3].

The contributions of this work can be summarized as follows:

• We devise an effective cascaded context mixer in the decoder which can learn both spatial and channel context information to infer the human body and hard keypoints.

• We propose several efficient training strategies to guide our model to deal with false-positive person detections and learn discriminative features from diverse poses.

• We present some sub-pixel refinement techniques to enhance location accuracy and comprehensively evaluate their performance and complementarity.

## 2 Related work

HPE methods can be grouped into 2D pose estimation (Rogez et al., 2012; Toshev and Szegedy, 2014; Newell et al., 2016; Fang et al., 2017; Huang et al., 2017; Cao et al., 2017; Yang et al., 2017; Xiao et al., 2018; Sun et al., 2018; Chen et al., 2018b; Sun et al., 2019; Zhang et al., 2019a; Li et al.,

---

[1] The instances with at least one annotated keypoint are counted.
[2] https://challenger.ai/competition/keypoint/

[3] A video demo can be found in https://github.com/chaimi2013/CCM/video

Table 1: A summary of the human keypoint detection methods based on DCNN.

| Category | Method | Backbone | Decoder | Extra Data | Post-processing | Performance |
|---|---|---|---|---|---|---|
| Bottom-up | Pishchulin et al. (2016) | VGG | - | - | offset regression | 82.4PCK$_h$@MPII |
| | Cao et al. (2017) | VGG-19 | multi-stage CNN | - | - | 61.8AP@COCO |
| | Newell et al. (2017) | Hourglass | - | - | multiscale average | 65.5AP@COCO |
| Top-down | He et al. (2017) | ResNet-50-FPN | conv+deconv | - | offset regression | 63.1AP@COCO |
| | Fang et al. (2017) | STN+Hourglass | - | - | parametric NMS | 63.3AP@COCO |
| | Papandreou et al. (2018) | ResNet-101 | 1x1 conv | ✓ | offset regression | 68.5AP@COCO |
| | Chen et al. (2018b) | ResNet-Inception | GlobalNet | - | flip/GF | 72.1AP@COCO |
| | Xiao et al. (2018) | ResNet-152 | deconv | - | flip/sub-pixel shift | 76.5AP@COCO |
| | Sun et al. (2019) | HRNet-w48 | 1x1 conv | - | flip/sub-pixel shift | 77.0AP@COCO |
| | Li et al. (2019) | 4 x ResNet-50 | GlobalNet | ✓ | flip/GF/sub-pixel shift | 78.1AP@COCO |
| | **Our CCM** | ResNet-152 HRNet-w48 | CCM | ✓ | sub-pixel refinement | 78.9AP@COCO |

2019; Girdhar et al., 2018) and 3D pose estimation (Rogez et al., 2012; Pavlakos et al., 2018b,a; Rhodin et al., 2018; Hossain and Little, 2018; Yang et al., 2018) according to the dimension of the coordinates of the keypoint locations. In this paper, we focus on 2D pose estimation, specifically, the multi-person pose estimation (MPPE) problem, which is more challenging than the single-person pose estimation (SPPE) problem. MPPE approaches can be further divided into bottom-up and top-down approaches. A summary of these approaches is presented in Table 1, where we outline their features according to the network structure, whether using extra data or not, and the postprocessing techniques used to improve detection accuracy. The details are presented in the following part.

Bottom-up approaches first detect all human keypoints and then associate them with each detected person instance (Cao et al., 2017; Newell et al., 2017; Pishchulin et al., 2016). This approach is usually faster than top-down approaches, but the assembly step can become intractable when person instances are ambiguous due to occlusions, blur, etc., degrading accuracy compared to top-down approaches. Top-down approaches first detect all person instances in the image and then apply SPPE on each detected person. Benefiting from recent progress in DCNN-based object detection (Ren et al., 2015; Ouyang et al., 2016; He et al., 2017; Lin et al., 2017; Liu et al., 2020; Chen et al., 2020), person detectors have achieved promising detection accuracy. Further, different neural networks have been proposed to learn strong feature representations based on multi-scale feature fusion and multi-level supervisions, *e.g.*, the Pyramid Residual Module in (Yang et al., 2017), the Cascaded Pyramid Network in (Chen et al., 2018b), the simple baseline model in (Xiao et al., 2018), and the High-resolution Net in (Sun et al., 2019), which detect keypoint locations with high accuracy. Our proposed method follows the top-down scheme and has an encoder-decoder structure. However, in contrast to the above methods, we study the role of context information for human keypoint detection by devising a cascaded context mixer module in the decoder to sequentially capture both spatial and channel context information.

Deep neural models benefit from a large scale of training data and efficient training strategies. Trained on more examples with diverse poses, the keypoint detection model can learn to infer the occluded or blurry keypoints from similar poses (Xiao et al., 2018; Li et al., 2019). For example, all the top entries of the COCO keypoint detection leaderboard[4] leverage external data such as the AI Challenger human keypoint detection dataset. In this paper, we also validate the benefit of using external data. However, instead of using transfer learning, we propose a more effective joint training strategy to harvest external data with heterogeneous labels. Further, we make use of unlabeled data, e.g., the unlabeled MS COCO dataset, referring to the knowledge distilling idea, where the pseudo-labels are generated by a teacher model. Besides, few of the above approaches deal with the mismatch problem of person detections during the training phase and testing phase, i.e., training with ground truth person instances while testing with detected ones, which may be false positives. In this paper, we propose an effective hard-negative person detection mining strategy in the training phase, adapting the model to predict no keypoints for those false person instances.

Dominant methods adopt a heatmap to represent the keypoint location where a Gaussian density map is placed on the corresponding pixel. However, the heatmap is usually in low-resolution compared with the input, which has a side effect on the location accuracy. Increasing the heatmap resolution means to decode high-resolution features, which may incur extra computational cost and model complexity. Instead, prior methods adopt computationally efficient postprocessing techniques to refine the predictions, for example, shifting the detected location by 0.25 pixels according to the local gradient directions (Chen et al., 2018b; Xiao et al.,

---

[4] http://cocodataset.org/index.htm#keypoints-leaderboard

2018; Sun et al., 2019). In contrast to them, we present a sub-pixel refinement technique using the second-order approximation, it turns out to be more effective than the above one and boost the detection accuracy by a large margin. Besides, we comprehensively study the sub-pixel refinement techniques for postprocessing including the second-order approximation (SOA), the Soft Non-Maximum Suppression (Soft-NMS), the sub-pixel shift of flipped heatmaps (SSP), and the Gaussian filtering on heatmaps (GF). Experiment results validate that these techniques can significantly boost the keypoint detection performance. Moreover, they are complementary to each other, and so the combination of them will boost the performance further.

## 3 Cascaded Context Mixer for Human Keypoint Detection

In this part, we propose a novel human keypoint detection model based on a cascaded context mixer module to explicitly and simultaneously model spatial and channel context information. To further exploit the representation capacity of the model, we propose three efficient training strategies including a hard-negative person detection mining strategy to migrate the mismatch between training and testing, a joint-training strategy to use abundant unlabeled samples by knowledge distilling, and a joint-training strategy to exploit external data with heterogeneous labels. To improve the detection accuracy, we also present four postprocessing techniques to refine predictions at the sub-pixel level.

### 3.1 Motivation

Since occlusions are ubiquitous in and between human bodies, we take it as an example to show how humans carry out the detection process when dealing with the aforementioned hard cases. As shown in Figure 2(a), occlusions include self-occlusion (A, C, D, E, F), occlusion by others (C, G), and truncation (B, H). Note that all the faces are self-occluded in Figure 2(a), *i.e.*, half of each face is invisible. Humans can easily recognize a complete object and its boundaries, even if it is occluded. According to Gestalt psychology in visual perception, we can group fragmented contours under the law of closure (Wagemans et al., 2012), e.g., D and E. We have also seen numerous human body positions and acquired the common sense that a body consists of symmetric legs, hands, and face and that different body parts move and form different poses. Therefore, we can easily infer the occluded parts like A, B, C, G, and H. For the more difficult case F, we may infer the invisible arms by judging the boy's intention and rehearsing the same action psychologically.

The Recognition-by-Components (RBC) theory tells us that humans quickly recognize objects even under occlu-



(a)



(b)

Fig. 2: (a) Ubiquitous occlusions in an image from the MS COCO dataset. (b) Illustration of the keypoint detection process carried out by humans when faced with occlusions.

sions by characterizing the object's components (Biederman, 1987). One possible path of the keypoint detection process carried out by humans could be divided into four main stages as shown in Figure 2(b). First, we recognize and locate a human body (The top-down approaches follow this paradigm (Chen et al., 2018b; Xiao et al., 2018; Sun et al., 2019)). Then, we recognize each body part to further locate the keypoints belonging to it (Some graph-based models complete this step explicitly (Felzenszwalb et al., 2008; Holt et al., 2011; Wang and Li, 2013; Yang and Ramanan, 2013)). Here, we can easily identify some distinct and visible keypoints. Finally, we infer the remaining keypoints which may be invisible or ambiguous. How do we accomplish this? By reviewing the inference process and the occlusion cases in Figure 2(a), we think that the "*context*" plays an important role when we associate separate body parts into a whole or infer an invisible keypoint (Chen et al., 2020; Ma et al., 2020). Another key factor may be that we have *a priori* knowledge of human body configurations in all possible poses. The context tells us about the surrounding visible body parts, and the *a priori* knowledge helps us to determine what the category and location of the invisible part should be. This motivates us to design a context-aware model that can efficiently learn useful feature representation from diverse poses.

Fig. 3: Visualization of the feature maps learned by the ResNet-50 encoder on two test images as shown in (a)-(b) and (c)-(d), respectively.

## 3.2 Cascaded Context Mixer-based Decoder

In this paper, we tackle the multi-person pose estimation problem by following a topdown scheme. First, a human detector is used to detect the bounding box for each person instance. Then, the proposed CCM model detects keypoints for each person instance. After aggregating the detections using Object Keypoint Similarity (OKS)-based Non-Maximum Suppression (NMS), we obtain the final pose estimation. As shown in Figure 4, the CCM model has an encoder-decoder structure where we use ResNet (He et al., 2016) or HRNet (Sun et al., 2019) as the backbone encoder network. The decoder consists of three Context Mixer (CM) modules in a cascaded manner. The details of CM are presented as follows.

### 3.2.1 Context Mixer

One the one hand, as we know that DCNN is able to learn distributed semantics at each feature channel, which leads to powerful representation ability to recognize a vast number of categories. In Figure 3, we visualized some feature channels from the ResNet-50 encoder. As can be seen, the encoder learned to activate one or multiple explicit body parts. On the other hand, as discussed in Section 3.1, context information could be useful to infer invisible keypoints and help to learn human body configuration. Thereby, we exploit two types of context information in this paper, i.e., global context and local context.

First, since different feature channels correspond to certain body parts (Figure 3(a) and Figure 3(c)), the relationship between them can be modeled to further enhance the

feature representation. To this end, we use a channel attention module to encode the relationship between different channels, which inherits the idea of squeeze-and-excitation (SE) network (Hu et al., 2018). The attention vector is used to reweigh features in a channel-wise manner, which is a global operation from the perspective of spatial dimension. In this way, we can extract the global context to enhance the feature representation. Second, since some of the feature responses at different body parts have strong correlations (Figure 3(b) and Figure 3(d)), we can encode the spatial relationship by using convolution layers with large receptive fields such that they can cover different body parts. To this end, we leverage hybrid-dilated convolutions to capture multi-scale spatial context information within different receptive fields, which inherits the idea of atrous spatial pyramid pooling (ASPP) (Chen et al., 2018a). These two types of context information collaborate with each other to learn discriminative feature representations. Besides, we also use a residual branch to reuse the learned features from the previous stage. Thereby, there are three branches in CM as shown in the middle part of Figure 4. We present the details of each branch as follows.

In the residual branch of the $k^{th}$ CM, feature maps $f_{k-1}^{CM}$ from the previous stage are first up-sampled two times before being fed into a $1 \times 1$ convolutional layer to output feature maps $f_k^{RES}$ of size $H_k \times W_k \times C_k$, i.e.,

$$f_k^{RES} = \phi_k^{RES} \left( f_{k-1}^{CM} \uparrow_2 \right), \tag{1}$$

where $\uparrow_2$ denotes the up-sampling operation, $\phi_k^{RES} (\cdot)$ denotes the function learned by the convolutional layer.

In the SE branch of the $k^{th}$ CM, the feature maps $f_{k-1}^{CM}$ first go through a global pooling layer. Then, the obtained feature vector is fed into a bottle-neck layer with $1 \times 1$ convolutions. The feature dimension is reduced to $1 \times 1 \times C_k/4$. Then, it is fed into a subsequent $1 \times 1$ convolutional layer to increase the feature dimension to $1 \times 1 \times C_k$. A sigmoid function is used to squeeze the feature vector $f_k^{SE}$ into the range $[0, 1]$, i.e.,

$$\alpha_k^{SE} = \sigma \left( \phi_k^{SE} \left( GP \left( f_{k-1}^{CM} \right) \right) \right), \tag{2}$$

where $GP (\cdot)$ denotes the global pooling operation, $\phi_k^{SE} (\cdot)$ denotes the function learned by those two $1 \times 1$ convolutional layers, and $\sigma (\cdot)$ denotes the sigmoid activation function.

In the HDC branch of the $k^{th}$ CM, the feature maps $f_{k-1}^{CM}$ go through four $3 \times 3$ convolutional layers with different dilated rates, i.e., 1, 2, 3, and 4. Each convolutional layer has $C_k/4$ kernels. These feature maps are then concatenated and fed into a deconvolutional layer of stride 2. The output feature maps $f_k^{HDC}$ are of size $H_k \times W_k \times C_k$, i.e.,

$$f_k^{HDC} = \phi_k^{HDC} \left( \left[ \phi_k^{d1} \left( f_{k-1}^{CM} \right) ; \dots ; \phi_k^{d4} \left( f_{k-1}^{CM} \right) \right] \right), \tag{3}$$

where $\phi_k^{d1} (\cdot) \sim \phi_k^{d4} (\cdot)$ denote functions learned by the dilated convolutional layers, $\phi_k^{HDC} (\cdot)$ denotes the function
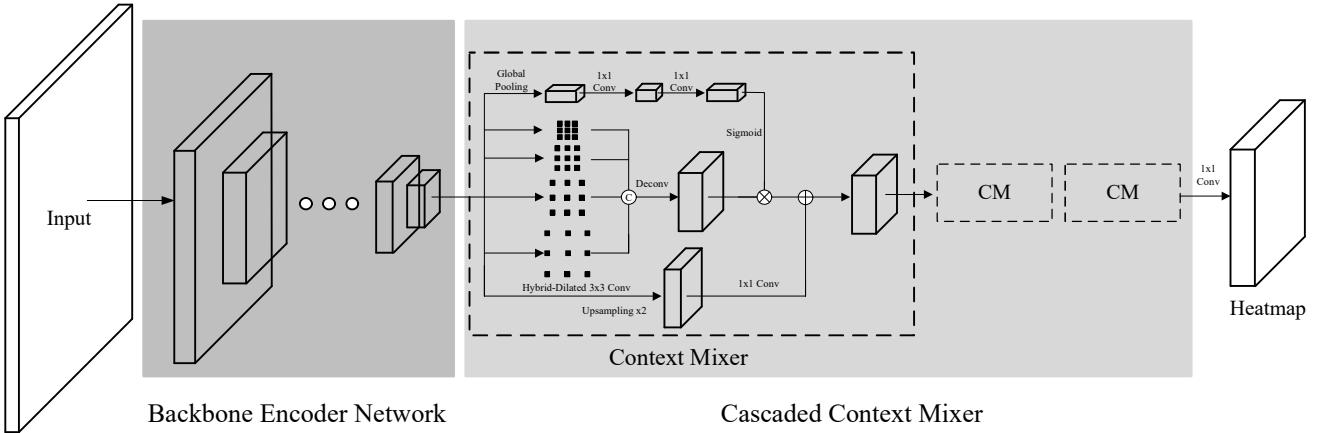
Fig. 4: The diagram of the proposed human keypoint detection model based on cascaded context mixer (CCM).

learned by the deconvolutional layer, and $[;]$ denotes the concatenation operation.

Then, the output of the $k^{th}$ CM is calculated as:

$$f_k^{CM} = f_k^{HDC} \odot \alpha_k^{SE} + f_k^{RES}$$
$$\triangleq \phi_k^{CM}\left(f_{k-1}^{CM}\right), \quad (4)$$

where $\odot$ denotes the channel-wise multiplication and $\phi_k^{CM}\left(\cdot\right)$ denotes the mapping function learned by the $k^{th}$ CM. Batch normalization (Ioffe and Szegedy, 2015) is used after each convolutional layer and deconvolutional layer. ReLU is used after the first convolutional layer in the SE branch, all the convolutional layers in the HDC branch, and the output of each CM module (Krizhevsky et al., 2012).

### 3.2.2 Cascaded Context Mixer

To capture the context information at multiple resolutions and learn multi-scale feature representation, we stack $K$ CMs sequentially to decode the features from the encoder step by step and increase their resolutions accordingly. Mathematically, it can be written as:

$$f^{ENC} = \phi^{ENC}\left(img\right), \quad (5)$$

$$f^{DEC} = \phi_K^{CM}\left(\ldots\left(\phi_1^{CM}\left(f^{ENC}\right)\right)\right)$$
$$\triangleq \phi^{DEC}\left(f^{ENC}\right) \quad (6)$$

where $img$ represents the input image, $\phi^{ENC}\left(\cdot\right)$ denotes the function learned by the encoder, $f^{ENC}$ is the encoded feature, $\phi^{DEC}\left(\cdot\right)$ denotes the function learned by the decoder, $f^{DEC}$ is the decoded feature. Note that we denote $f_0^{CM} \triangleq f^{ENC}$ for consistency.

The decoded feature $f^{DEC}$ is fed into a final $1 \times 1$ convolutional layer to predict the target heatmaps:

$$h = \phi^{PRE}\left(f^{DEC}\right), \quad (7)$$

where $\phi^{PRE}\left(\cdot\right)$ denotes the function learned by the prediction layer and $h$ represents the predicted heatmaps.

### 3.2.3 Auxiliary Decoder and Intermediate Supervision

Deep supervision (Lee et al., 2015) refers to the technique that adds auxiliary supervision on some intermediate layers within a deep neural network. It facilitates multi-scale and multi-level feature learning by allowing error information back-propagation from multiple paths and alleviating the problem of vanishing gradients in deep neural networks. Leveraging the deep supervision idea, we also add an auxiliary decoder $\phi_{aux}^{DEC}\left(\cdot\right)$ after the penultimate stage of the encoder. Its structure is identical to $\phi^{DEC}\left(\cdot\right)$ described above. These two decoders do not share weights.

### 3.3 Training objective

The ground truth heatmap is constructed by placing a Gaussian peak at each keypoint's location in the image plane. The number of heatmaps is identical to the number of keypoints predefined in the dataset, for example, 17 for the MS COCO dataset (Lin et al., 2014) and 14 for the AI Challenger dataset. We use an MSE loss to supervise the network during training. Mathematically, it is defined as:

$$L_{main} = \frac{1}{H_K W_K C} \sum_{i,j,c} \left\| h\left(i,j,c\right) - h^{GT}\left(i,j,c\right) \right\|^2, \quad (8)$$

where $C$ is the number of heatmaps, $i$, $j$, and $c$ are the spatial and channel index, $h$ and $h^{GT}$ are the predicted and ground truth heatmaps, respectively. Similar to Eq. (8), an extra MSE loss $L_{aux}$ is also added to the auxiliary decoder as an intermediate supervision. The final training objective is defined as the weighted sum of both losses:

$$L = L_{main} + \lambda L_{aux}, \quad (9)$$

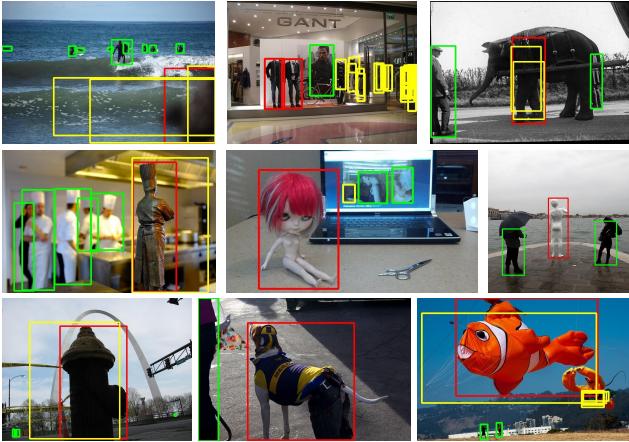where $\lambda$ is the weight for the auxiliary loss.

Fig. 5: Illustration of the hard-negative detections. Green: ground truth person instances. Red: hard-negative detections. Yellow: low-score false positive detections.

## 4 Learning from diverse poses with efficient strategies

CCM's capacity to model context information and learn discriminative feature representation can be exploited by learning from massive training samples with diverse poses. In this paper, we propose a hard-negative person detection mining strategy to migrate the mismatch problem of person detections during the training phase and testing phase, a joint-training strategy on unlabeled samples by knowledge distilling, and a joint-training strategy on external data with heterogeneous labels.

### 4.1 Hard-Negative Person Detection Mining (HNDM)

Top-down approaches detect person instances before detecting keypoints on them. On the one hand, although modern detection models have achieved a good detection performance, they may still produce some false positive detections due to occlusion, similar appearances, etc. On the other hand, the keypoint detection model is usually trained with ground truth bounding box annotations enclosing exact person instances. It has never seen any false positive detections during the training phase. Therefore, there is a mismatch between training and testing, which may lead to incorrect keypoint predictions for those false positive person detections. To address this issue, we propose a hard-negative person detection mining strategy.

First, we trained a Mask R-CNN on the MS COCO training set containing only the category of person. ResNeXt152 was used as the backbone network. It achieved a mean average precision (AP) of 60.4 on the COCO minival dataset. Then, we evaluated the training set using the detection model and screened out those detections with sufficiently high

scores, *e.g.*, $\geq 0.5$, but no intersections with ground truth person instances. These were treated as hard-negative detections in this paper. Some examples are shown in Figure 5. During training CCM, these hard-negative detections could be added to the training set. Their keypoint heatmaps are set to all-zero maps. In this way, CCM learns to predict no keypoints on those false "person instances".

### 4.2 Joint-Training on Unlabeled Samples by Knowledge Distilling

To increase pose diversity in the training samples, we leverage the MS COCO unlabeled dataset, which contains over 120k images. Since there are no person and keypoint annotations, we generate pseudo labels by referring to the knowledge distilling idea. First, we used the above-trained person detector to detect all possible person instances within the unlabeled images. Then, we screened out those detections with scores above a predefined threshold, which is determined by guaranteeing the number of average person instances per image to be identical to the one calculated from the ground truth annotations of the MS COCO training set. In our case, this threshold was 0.9924. Next, we trained several keypoint detection models using ResNet152 and HRNet-w48 as the backbone encoder network and used them to detect keypoints on the person detections obtained from the previous stage. We fused the predictions by different models, kept all keypoints with scores above 0.9 as the pseudo labels, and treated the rest as unlabeled. In this way, we distill the learned "knowledge" in the keypoint detection model to the pseudo labels of unlabeled training samples and use them to supervise the network.

### 4.3 Joint-Training on External Data with Heterogeneous Labels

Different datasets may not share the same annotation norms, even if they are used for the same purpose. For instance, 17 keypoints are used to define a human skeleton in the MS COCO dataset (Lin et al., 2014), but only 14 in AI Challenger (AIC). Five keypoints correspond to the eyes, ears, and nose in MS COCO, while only the "top of head" is annotated in AIC. AIC has a keypoint annotation for the neck, which is absent in MS COCO. The other 12 keypoints corresponding to limbs are the same in both datasets. To use AIC with MS COCO, a common practice is to train a network on AIC and then change the number of channels in the final prediction layer and fine-tune this network on MS COCO. In this paper, we propose a simple but effective joint-training strategy that mixes the training samples in both datasets to leverage the diverse poses simultaneously. To make the training tractable, we align the labels in AIC

with the ones in MS COCO by keeping the 12 common annotations and discarding the others.

To further adapt the trained model to the MS coco dataset, we can also add a finetune stage. In conclusion, the training strategies described above can be summarized as:

$$Train\,|\Phi \rightarrow Finetune\,|\Theta\ ,\tag{10}$$

where $\Phi \in \{A, AC, ACH, ACHU, CHU\}, \Theta \in \{C, CH\}$, $A$ denotes the AIC training dataset, $C$ denotes the COCO training set, $H$ denotes the hard-negative training samples, and $U$ denotes the reprocessed unlabeled training set.

## 5 Sub-pixel Refinement Techniques for Postprocessing

To migrate the scale mismatch of representation between the high-resolution ground-truth keypoint location and the corresponding position on the low-resolution heatmap, different sub-pixel refinement techniques have been introduced, e.g., the 0.25-pixel shift of the maximum response pixel (Chen et al., 2018b; Xiao et al., 2018; Sun et al., 2019; Li et al., 2019), the one-pixel shift of flipped heatmap (Xiao et al., 2018; Sun et al., 2019), and Gaussian filtering of predicted heatmap (Chen et al., 2018b; Li et al., 2019). With this regard, we devise four sub-pixel refinement techniques by (1) exploring the second-order approximation to locate the maximum response at the sub-pixel level accuracy, (2) introducing the Soft Non-Maximum Suppression (Soft-NMS) to refine the maximum response's location instead of the original NMS, (3) extending the one-pixel shift of flipped heatmap to a general sub-pixel form, and (4) applying the Gaussian filter on predicted heatmaps. They are complementary (or orthogonal) to each other and can be conducted sequentially at different stages of the inference phase.

### 5.1 Sub-pixel Refinement by the Second-Order Approximation

Sub-pixel refinement by the second-order approximation has ever been used in the depth super-resolution and disparity calculation literature (Yang et al., 2007). They use the cost volume at different disparities to find the optimal disparity that minimizes the inconsistency between the left and right views. Therefore, they face the issue to estimate the sub-pixel disparity given the costs at integer disparities. Likewise, we can adopt such a technique to estimate the sub-pixel keypoint location given the predicted heatmaps. In contrast to the one-dimensional cost volume, the heatmaps are in the two-dimensional space. Therefore, we present two kinds of second-order approximation by either using a parabola approximation for each dimension or using a paraboloid approximation simultaneously for the two dimensions.
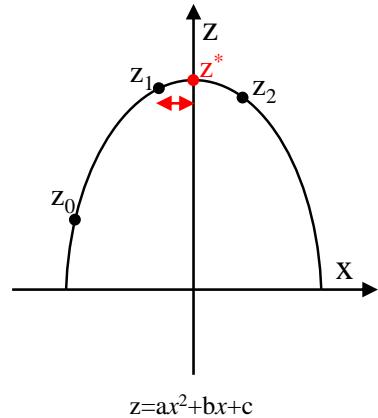


Fig. 6: The sub-pixel refinement by the second order approximation, i.e., the case of parabola.

#### 5.1.1 The Parabola Approximation

As we know that the heatmap is a Gaussian function w.r.t. the pixel dimension, nevertheless, it is reasonable to approximate it with a parabola function in a local neighborhood of the maximum pixel as shown in Figure 6:

$$z\,(x) = ax^2 + bx + c,\tag{11}$$

where $a$, $b$, and $c$ are the coefficients of the parabola. The maximum $z$ is subject to the first order condition:

$$\frac{\partial z}{\partial x} = 2ax + b = 0.\tag{12}$$

Therefore, the maximum $z$ is reached at $x^* = -\frac{b}{2a}$. Given $z_1$ is the maximum pixel at $x_0$, $z_0$ and the $z_2$ is the heatmap response at $x_0 - 1$ and $x_0 + 1$, we can calculate $x^*$ as:

$$x^* = x_0 + \frac{z_0 - z_2}{2\,(z_0 + z_2 - 2z_1)}.\tag{13}$$

The second term in the right-hand side (RHS) of Eq. (13) is the sub-pixel shift from the detected maximum response pixel $x_0$ to the underlying maximum one $x^*$. Similarly, we can calculate the optimal $y^*$ along the vertical dimension.

#### 5.1.2 The Paraboloid Approximation

Since the heatmap is represented as a two-dimension Gaussian function, we can approximate it with a paraboloid function in a local neighborhood of the maximum pixel as shown in Figure 7:

$$z\,(x, y) = ax^2 + by^2 + cxy + dx + ey + f,\tag{14}$$

where $a, b, c, d, e$, and $f$ are the coefficients of the paraboloid. The maximum $z$ is subject to the first order condition:

$$\begin{cases} \frac{\partial z}{\partial x} = 2ax + cy + d = 0 \\ \frac{\partial z}{\partial y} = 2by + cx + e = 0 \end{cases}.\tag{15}$$
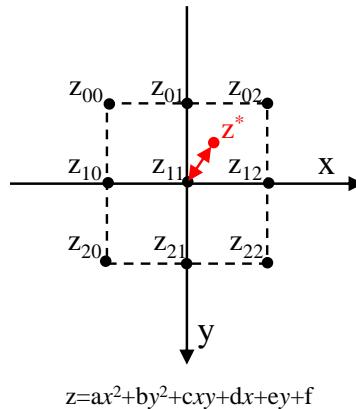
$$z = ax^2 + by^2 + cxy + dx + ey + f$$

Fig. 7: The sub-pixel refinement by the second order approximation, i.e., the case of paraboloid.

Therefore, the maximum $z$ is reached at:

$$\begin{cases} x^* = \frac{2bd - ce}{c^2 - 4ab} \\ y^* = \frac{2ae - cd}{c^2 - 4ab} \end{cases} . \tag{16}$$

As shown in Figure 7, $z_{11}$ is the maximum in the heatmap at $(x_0, y_0)$, $z_{00} \sim z_{22}$ are its 8-neighbor heatmap responses. Assuming that the coordinate original is at $(x_0, y_0)$, we can calculate the coefficients $a \sim e$ as:

$$\begin{aligned} a = \frac{1}{8} [ & 2(z_{12} + z_{10} - 2z_{11}) \\ & + (z_{02} + z_{00} - 2z_{01}) \\ & + (z_{22} + z_{20} - 2z_{21})] \end{aligned} \tag{17}$$

$$\begin{aligned} b = \frac{1}{8} [ & 2(z_{01} + z_{21} - 2z_{11}) \\ & + (z_{00} + z_{20} - 2z_{10}) \\ & + (z_{02} + z_{22} - 2z_{12})] \end{aligned} \tag{18}$$

$$c = \frac{1}{4} (z_{00} + z_{22} - z_{02} - z_{20}), \tag{19}$$

$$d = \frac{1}{8} [(z_{02} - z_{00}) + (z_{22} - z_{20}) + 2(z_{12} - z_{10})], \tag{20}$$

$$e = \frac{1}{8} [(z_{20} - z_{00}) + (z_{22} - z_{02}) + 2(z_{21} - z_{01})], \tag{21}$$

Therefore, the maximum $z$ is reached at:

$$\begin{cases} x^* = x_0 + \frac{2bd - ce}{c^2 - 4ab} \\ y^* = y_0 + \frac{2ae - cd}{c^2 - 4ab} \end{cases} . \tag{22}$$

The second terms in the RHS of Eq. (22) is the sub-pixel shift from the detected maximum response pixel $x_0$ (resp. $y_0$) to the underlying maximum one $x^*$ (resp. $y^*$). Given the heatmap, after locating the maximum pixel, we calculate the optimal location according to Eq. (13) or Eq. (22).
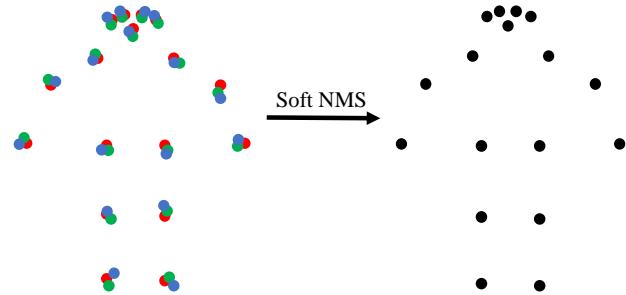


Fig. 8: The sub-pixel refinement by Soft NMS.

### 5.2 Sub-pixel Refinement by Soft-NMS

During the inference phase of top-down approaches, there may be several bounding boxes detected around a person instance. Consequently, we may estimate several poses on them. Similar to the Non-Maximum Suppression (NMS) postprocessing step used in object detection, NMS is also applied to the detected poses (Chen et al., 2018b; Xiao et al., 2018; Sun et al., 2019). Different from the Intersection over Union (IoU) used to compare the overlap between two bounding boxes, the Object Keypoint Similarity (OKS)-based IOU (OKS-IOU) is used to compare two poses. The original NMS is to filter out all the poses that have sufficient large OKS-IOUs with top-ranked detection. We argue that those poses can also be treated as reasonable estimates, which can be used to get a more stable result. Therefore, we present Soft-NMS for sub-pixel refinement as follows.

As shown in Figure 8, the poses marked in red, blue, and green dots are three estimates and the red pose has the highest score. We can calculate the fusion results by leveraging the OKS-IOU as the fusion weight, i.e.,

$$p_i^* = \frac{\sum\limits_{j \in \Lambda_i} IOU_{ij}^{OKS} p_j}{\sum\limits_{j \in \Lambda_i} IOU_{ij}^{OKS}}, \tag{23}$$

where $\Lambda_i$ is the index set of poses to be filtered given the top-ranked detection $p_i$, $IOU_{ij}^{OKS}$ is the OKS-IOU between $p_i$ and $p_j$. Note that we treat $i \in \Lambda_i$ and set $IOU_{ii}^{OKS} = 1$. We can re-write Eq. (23) as:

$$p_i^* = p_i + \frac{\sum\limits_{j \in \Lambda_i} IOU_{ij}^{OKS} (p_j - p_i)}{\sum\limits_{j \in \Lambda_i} IOU_{ij}^{OKS}}, \tag{24}$$

where the second term in the RHS is the sub-pixel shift from the top-ranked detection $p_i$ to the underlying best one $p_i^*$.

### 5.3 Gaussian Filtering on Heatmaps

The regressed heatmap may not be as smooth as the ground truth Gaussian density map. A Gaussian filter-based post-
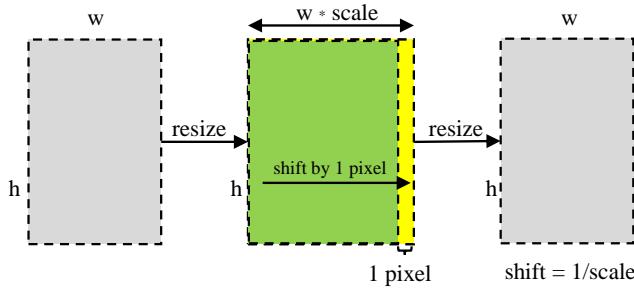
Fig. 9: The sub-pixel shift of flipped heatmaps.

processing technique is proposed to smooth it and minimize the variance of the prediction (Chen et al., 2018b; Li et al., 2019). We evaluate this technique and compare it with other sub-pixel refinement techniques.

### 5.4 Sub-pixel Shift of Flipped Heatmaps

During the inference phase, some methods predict the pose from the flipped image and average the flipped heatmap with the original one to get the final prediction (Chen et al., 2018b; Xiao et al., 2018). However, since the flipped heatmap is not aligned with the original one, one common practice is to shift the flipped heatmap by one pixel.

In this paper, we extend the one-pixel shift to the sub-pixel shift as shown in Figure 9. First, we resize the flipped heatmap by a scale ratio $r$ ($r \geq 1$) along the horizontal axis. Then, we shift it by one pixel to the right. Then, we resize it back to the original size. As can be seen, the effective shift becomes $1/r$ pixel. We name it as the sub-pixel shift (SSP) in this paper. Note the above process is equivalent to a linear interpolation of the original heatmap and its one-pixel shifted version, i.e.,

$$h_f^*(i,j) = \left(1 - \frac{1}{r}\right) h_f(i,j) + \frac{1}{r} h_f(i, j - j_0), \quad (25)$$

where $h_f$ is the flipped heatmap, $j_0$ is the maximum shifted pixels, i.e., $j_0 = 1$ in this paper. $h_f^*$ is the sub-pixel estimate. It becomes the one-pixel shift technique when $r = 1$.

## 6 Experiments

We conducted extensive experiments to demonstrate the effectiveness of the proposed model. First, comprehensive ablation studies on the components of CM were presented, followed by the comparative studies of the proposed training strategies and sub-pixel refinement techniques. Next, we compared the proposed model with representative state-of-the-art methods in terms of detection accuracy, model complexity, and computational cost. Then, we presented some visual examples of the detection results by our model and

explained the detection process by inspecting the learned features at each stage. Finally, we empirically studied the impact of visible and invisible annotations in our model and obtained useful some insights.

### 6.1 Experimental settings

**Datasets**: The COCO Keypoint Challenge addresses multi-person pose estimation in challenging uncontrolled conditions (Lin et al., 2014). The dataset is split into training, minival, test-dev, and test-challenge sets. The training set includes 118k images and 150k person instances, the minival dataset includes 5000 images, and the test-dev set includes 20k images. It also provides an unlabeled dataset containing 123k images. 110k person instances and corresponding keypoints were detected using the method described in Section 4.2. The external dataset from AIC contains a training set with 237k images and 440k person instances and a validation set with 3000 images. We also evaluated CCM and the baseline model on the recently proposed OCHuman benchmark (Zhang et al., 2019b) comprising heavily-occluded human instances to compare their performance on handling occluded cases. This dataset contains 8110 human instances with detailed keypoint annotations like COCO. It is divided into two subsets: OCHuman-Moderate and OCHuman-Hard. The first subset contains instances with MaxIoU in the range of 0.5 and 0.75, while the second contains instances with MaxIoU larger than 0.75. MaxIoU denotes the max IoU of a person with others in an image.

**Evaluation metrics**: We report the main results based on the object keypoint similarity (OKS)-based mean average precision (AP) over 10 OKS thresholds, where OKS defines the object keypoint similarity between different human poses. They are calculated as follows (Lin et al., 2014):

$$AP = mean \left\{ AP_{@(0.50:0.05:0.95)} \right\}, \quad (26)$$

$$AP_{@s} = \frac{\sum_p \delta\left(OKS_p > s\right)}{\sum_p 1}, \quad (27)$$

$$OKS_p = \frac{\sum_i \exp\left(-d_{pi}^2 / \left(2a_p^2 \sigma_i^2\right)\right) \delta\left(v_{pi} > 0\right)}{\sum_i \delta\left(v_{pi} > 0\right)}, \quad (28)$$

where $p$ is the person instance index, $i$ is the keypoint index, $\delta\left(\cdot\right)$ is the Kronecker function. $\delta\left(\cdot\right) = 1$ if the condition holds, otherwise 0. $s$ is a threshold, $d_{pi}$ is the Euclidean distance between the predicted $i^{th}$ keypoint of the person instance $p$ and its ground truth, $a_p$ is the area of the person instance $p$, $\sigma_i$ is the normalization factor predefined for each keypoint type, and $v_{pi}$ is the visible status.

**Implementation details**: The feature dimension $C_i$ of each CM was set to 256 for ResNet-50 and 128, 96, 64 for ResNet-152, 32 for HRNet-w32, and 48 for HRNet-w48. All backbone networks were pre-trained on the ImageNet

dataset (Deng et al., 2009). Gaussian initialization was used for convolutional and deconvolutional layers in the decoder. The weights and bias in BatchNorm layers were initialized as 1 and 0, respectively. CCM was implemented in Pytorch (Paszke et al., 2017) and trained on four NVIDIA Tesla V100 GPUs using the Adam optimizer. Hyper-parameters were set by following (Xiao et al., 2018; Sun et al., 2019). We used the detection results on the minival set and test-dev set released in (Xiao et al., 2018) for a fair comparison, which were obtained by a faster-RCNN detector (Ren et al., 2015) with detection AP 56.4 for the person category on COCO val2017. We obtained the final predictions by averaging the heatmaps of the original and flipped image as in (Chen et al., 2018b; Newell et al., 2016; Xiao et al., 2018).

## 6.2 Ablation Studies

### 6.2.1 Ablation Study of the Components of CM

Table 2: Ablation study on the components of CCM. AD: auxiliary decoder, D: dilated convolutions. AP/AR: mean average precision/recall on COCO minival set. Backbone network: ResNet-50 (R50) and HRNet-w32 (HR32).

| Method | SE | HDC | AD | D | $AP$ | $AR$ |
|---|---|---|---|---|---|---|
| Baseline (R50) (Xiao et al., 2018) | | | | | 70.4 | 76.3 |
| **CCM (R50), K=3** | ✓ | | | | 71.7 | 77.9 |
| | | ✓ | | | 71.8 | 78.0 |
| | | | ✓ | | 72.1 | 78.2 |
| | | | | ✓ | 72.7 | 78.7 |
| | ✓ | ✓ | ✓ | ✓ | 73.5 | 79.1 |
| HRNet-w32 (Sun et al., 2019) | | | | | 74.4 | 79.8 |
| **CCM (HR32), K=1** | ✓ | ✓ | ✓ | | 75.5 | 80.9 |
| **CCM (HR32), K=2** | ✓ | ✓ | ✓ | | 75.5 | **81.0** |
| **CCM (HR32), K=3** | ✓ | ✓ | ✓ | | **75.6** | **81.0** |

First, we conducted an ablation study on the components of CM by training different variants on the COCO training set and calculating the AP and AR on the minival set. The backbone network was ResNet-50 (R50) and HRNet-w32, and the input size was $256 \times 192$. The results are shown in Table 2. For the ResNet-50 backbone, we used three CM modules in the decoder to keep consistent with the baseline model which used three deconvolutional layers. In this way, the feature map size increased from $8 \times 6$ at the end of the encoder to $64 \times 48$ at the end of the decoder. As can be seen, each component achieved gains over the baseline model (Xiao et al., 2018), for example, adding the SE branch or HDC branch improved the detection accuracy by a margin

of 1.3 AP or 1.4 AP over the baseline model. The auxiliary decoder and intermediate supervision also benefited the detection model and achieved a gain of 1.7 AP. Using dilated convolutions in the encoder could produce feature maps with $2\times$ higher resolution (i.e., $128 \times 96$), thereby benefiting the localization accuracy. As can be seen, it achieved a gain of 2.3 AP over the baseline model. These components are complementary to each other that the combination of them improved the detection accuracy further, i.e., a gain of $0.8 \sim 1.8$ over the individual component.

For the HRNet-w32 backbone, since it could produce high-resolution features (i.e., $64 \times 48$), we only attached one CM module after HRNet-w32 and did not use any dilated convolutions, thereby increasing the feature map size to $128 \times 96$. Nonetheless, we also investigated whether extra CM modules were useful or not. To this end, we attached one or two extra CM modules (i.e., K=2 or 3) accordingly. To avoid huge computations for processing larger feature maps, we replaced the deconvolutional layers in the extra CM modules with convolutional layers to keep the feature map size. As can be seen from Table 2, CCM (K=1) outperformed the vanilla HRNet-w32 by a gain of 1.1 AP. Besides, adding extra CM modules only led to marginally better results while the parameters and computations increased from 25.56M and 7.92 GFLOPs (K=1) to 28.58M and 8.16 GFLOPs (K=2), and 28.6M and 8.4 GFLOPs (K=3), respectively. To make a trade-off between accuracy and computational efficiency, we chose K=1 as the default setting for the HRNet-w32 backbone.

Table 3: Comparison of CCM and the Baseline model on the OCHuman dataset (Zhang et al., 2019b). Backbone network: ResNet-50 (R50).

| Method | Dataset | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ |
|---|---|---|---|---|---|---|
| Baseline | $val$ | 59.9 | 79.3 | 66.1 | **64.7** | 59.9 |
| | $val_{[0.5-0.75]}$ | 62.7 | 81.5 | 69.4 | **64.7** | 62.7 |
| | $val_{[0.75-1]}$ | 42.9 | 63.1 | 46.4 | - | 42.9 |
| | $test$ | 52.1 | 70.6 | 57.0 | **72.6** | 52.1 |
| | $test_{[0.5-0.75]}$ | 61.1 | 80.4 | 66.7 | **89.2** | 61.0 |
| | $test_{[0.75-1]}$ | 43.4 | 60.8 | 47.4 | **40.0** | 43.4 |
| CCM(R50) | $val$ | **63.2**$_{(+3.3)}$ | 81.9 | 68.1 | 53.7 | 63.2 |
| | $val_{[0.5-0.75]}$ | **66.2**$_{(+3.5)}$ | 84.5 | 71.8 | 53.7 | 66.3 |
| | $val_{[0.75-1]}$ | **43.5**$_{(+0.6)}$ | 66.0 | 45.6 | - | 43.5 |
| | $test$ | **54.9**$_{(+2.8)}$ | 74.5 | 59.3 | 72.6 | 54.9 |
| | $test_{[0.5-0.75]}$ | **65.4**$_{(+4.3)}$ | 83.7 | 71.5 | 89.2 | 65.4 |
| | $test_{[0.75-1]}$ | **44.6**$_{(+1.2)}$ | 65.0 | 47.6 | 40.0 | 44.6 |

Besides, to compare the generalization ability of the proposed CCM and the baseline models when dealing with unseen heavily-occluded cases from other benchmarks, e.g., the OCHuman benchmark (Zhang et al., 2019b) and the PoseTrack benchmark (Andriluka et al., 2018), we evaluated

Table 4: Comparison of CCM and the Baseline model on the PoseTrack validation set (Andriluka et al., 2018). Backbone network: ResNet-50 (R50) and HRNet-w32 (HR32).

| Method | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ |
|---|---|---|---|---|---|---|
| Baseline(R50) | 69.6 | 86.4 | 75.7 | 40.6 | 73.8 | 72.2 |
| HRNet-w32 | 73.0 | 87.6 | 78.9 | 42.7 | 77.3 | 75.4 |
| **CCM(R50)** | **71.3$_{(+1.7)}$** | **87.1** | **76.5** | **42.0** | **75.5** | **74.1** |
| **CCM(HR32)** | **73.8$_{(+0.8)}$** | **87.7** | **79.0** | **44.8** | **78.2** | **76.2** |

them in a zero-shot manner, where all models were trained on the COCO training set without further fine-tuning. It is noteworthy that we only include the results of the baseline models and our model for the following several reasons. Firstly, the goal is of this paper is to investigate the key factors that have strong impacts on the performance of human keypoint detection. Since it has already been evidenced by prior excellent work (Xiao et al., 2018; Sun et al., 2019) that a better keypoint detector benefits pose tracking more effectively, we only focus on human keypoint detection in this paper. Here, we use the OCHuman dataset and Pose-Track dataset to validate the effectiveness of the proposed method for handling occluded cases. Secondly, the detection and tracking performance heavily depends on the person detector, but the detector as well as the implementation of joint propagation in (Xiao et al., 2018; Sun et al., 2019) are not available. As a result, it will be unfair to compare with their results in a different setting. We plan to investigate the role of context for pose tracking in our future work, where spatial-, temporal- and channel-wise context could be exploited together to enhance feature representation. Thirdly, the Simple baseline method (Xiao et al., 2018) and the High-resolution Network method (Sun et al., 2019) are already two strong baseline models. It is noteworthy that many recent methods use HRNet as the backbone and have achieved SOTA performance on public datasets or challenges. Thereby, it is representative to compare our model with them on these two datasets. Besides, we use the zero-shot transfer setting to evaluate how the model trained on a representative dataset (e.g., MS COCO) generalizes to unseen samples from different data distributions, especially those samples containing large occlusions.

The results on the OCHuman benchmark are listed in Table 3. As can be seen, CCM outperformed the baseline model for handling occlusions by large margins, e.g., 3.3 points of AP gain on the validation set and 2.8 points of AP gain on the test set. The gains mainly arise from the occlusion cases $val_{[0.5-0.75]}$ and $test_{[0.5-0.75]}$, which contain occluded instances with MaxIoU in the range of 0.5 and 0.75. The results on the PoseTrack validation set are listed in Table 4. Our model outperforms the the baseline model with a ResNet-50 backbone and the HRNet-w32 model by

1.7 AP and 0.8 AP, confirming the superiority of CCM when dealing with heavily-occluded human instances. For a person instance with occluded body parts, it is challenging to detect both visible and invisible keypoints due to the self-occlusion, incomplete body, and perplexity with adjacent overlapped bodies. The proposed CM enables the detection model to learn discriminative feature representation for diverse poses and help it to recognize occluded keypoints. After the network sees diverse poses, it "memorizes" different poses with/without occlusions in the form of feature mapping. Inferring an occluded keypoint thus becomes easier by associating it with similar poses. More discussions will be presented in Section 6.4 and Section 6.5.

**Remarks**: 1) CM has better representative capacity than the plain deconvolution layer in the simple baseline method (Xiao et al., 2018) because it leverages spatial and channel context information explicitly; 2) CM is also complementary to the high-resolution module in (Sun et al., 2019) and improves the performance of the stronger HRNet-w32; and 3) CM effectively handles occlusions by learning context features to infer the occluded keypoints.

### 6.2.2 Comparison of Training Strategies

Next, we present the results of using different training strategies described in Section 4 in Table 5, where A→C denotes the transfer learning strategy, AC→C denotes the joint-training strategy, i.e., training CCM on both AIC and COCO datasets then fine-tuning it on the COCO dataset. Other symbols have a similar meaning. As can be seen, leveraging the external AIC dataset increased the AP by 1.3 compared with the model trained on the COCO dataset in Table 2. The improvement became 1.5 AP when using the proposed joint-training strategy. The proposed HNDM method increased the AP further by an extra gain of 0.3, while the AR remained the same. This is reasonable since HNDM aims to suppress the keypoints of the false-positive person detections, meaning that it can increase the precision but has little influence on the recall. After exploiting the unlabeled dataset, CCM obtained a final AP of 75.6, a gain of 2.1 over the same model trained on the COCO dataset.

We also evaluated the impact of the detection score threshold in HNDM. Specifically, we set its value to 0.5, 0.7, and 0.9 in the training strategy "ACH→CH". The results are summarized in Table 6. As can be seen, the performance decreased with the increase of the threshold. Note that fewer hard negative detections were included in the training set when a larger threshold was used. For example, there were 11,762 detections at the threshold of 0.5 while only 4,359 and 952 detections were left at the threshold of 0.7 and 0.9, respectively. On the one hand, since there were many person detection proposals with low scores from the person detectors to increase the recall, thereby leveraging hard neg-

Table 5: Comparisons of CCM trained with the different strategies described in Section 4. A: the AIC training dataset; C: the COCO training set; H: the hard-negative training samples; U: the reprocessed unlabeled training set.

| Method | Training Strategy | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ | $AR^{@.5}$ | $AR^{@.75}$ | $AR^M$ | $AR^L$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **CCM (ResNet-50)** | A→C | 74.8 | **90.6** | 81.7 | 70.5 | 81.1 | 80.1 | **94.0** | 86.1 | 75.7 | 86.5 |
| | AC→C | 75.0 | 90.3 | 82.1 | 70.8 | 81.3 | 80.4 | **94.0** | 86.7 | 76.0 | 86.6 |
| | ACH→CH | 75.3 | **90.6** | 82.1 | 71.0 | 81.5 | 80.4 | **94.0** | 86.5 | 76.0 | 86.6 |
| | ACHU→CH | **75.6** | 90.5 | **82.5** | **71.4** | **81.8** | **80.7** | **94.0** | **86.8** | **76.3** | **87.0** |

Table 6: Study of the threshold setting in HNDM for the training strategy "ACH→CH" described in Section 4.1.

| Threshold | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ |
|---|---|---|---|---|---|---|
| 0.5 | **75.3** | 90.6 | 82.1 | 71.0 | **81.5** | 80.4 |
| 0.7 | 75.2 | **90.9** | **82.2** | **71.1** | 81.3 | 80.4 |
| 0.9 | 74.9 | 90.5 | 82.2 | 70.9 | 81.0 | 80.3 |

ative person detections with low scores (e.g., $0.5 \sim 0.7$) and predicting all-zero heatmaps could reduce false positive keypoints on those person detection proposals. On the other hand, due to the annotation policy, human-like objects such as model, sculpture, and doll (please see Figure 5) were not annotated as the person category, which however could be detected with high scores by the person detector. Since they shared similar appearance with real human bodies, only using these hard negative person detections with high scores (e.g., $\geq 0.9$) increased the difficulty for learning discriminative feature representation for keypoint detection.

**Remarks**: 1) The proposed joint-training strategy is more effective than transfer learning by reducing the domain gap during the pretraining phase; 2) HNDM can deal with the false-positive person detections, thereby improving the detection precision; and 3) exploiting extra unlabeled data using knowledge distilling enables the network to learn more discriminative features from abundant and diverse samples.

### 6.2.3 Comparison of Sub-pixel Refinement Techniques

We conducted the contrastive experiments by using different sub-pixel refinement techniques in both the simple baseline method (Xiao et al., 2018) and the High-Resolution Network. The results are summarized in Table 7 and Table 8.

First, the sub-pixel refinement techniques by second-order approximation described in Section 5.1 consistently improved the performance of both methods. Note that shifting towards the gradient directions by 0.25 pixels was effective and improved the performance of ResNet-50 and HRNet-w32 by 1.9 AP and 2.1 AP, respectively. Nevertheless, it only used the first-order derivative information and the shift was a fixed value, limiting its performance. In contrast, the proposed SOA refinement further increased the AP from 68.5 to 71.0, and 71.7 to 74.3, for ResNet-

50 and HRNet-w32, respectively. With the second-order approximation, SOA could adaptively calculate the shift vector for each heatmap. The paraboloid-based SOA and parabola-based SOA performed similarly. It is reasonable because the target heatmap is a 2D Gaussian density map where the density along the x-axis is independent of the density along the y-axis. The predicted heatmap had a similar pattern. Therefore, the paraboloid-based SOA had no obvious advantage over the parabola-based SOA. Nevertheless, the paraboloid-based SOA may be useful for those scenarios where the joint densities along different axes are not independent of each other, i.e., there is an elliptical response with a bias direction in the heatmap.

Second, Soft-NMS was effective, which improved the performance by 0.8 AP and 0.9 AP for ResNet-50 and HRNet-w32, respectively. The weighted fusion defined by Eq. (23) shifted the keypoint locations in sub-pixels by considering the reasonable estimations rather than filtering them out as done in standard NMS. Besides, the SOA technique was complementary to Soft-NMS. For example, the parabola-based SOA technique combined with Soft-NMS improved the performance further by 0.4 AP compared with using the parabola-based SOA technique individually and improved the performance further by 2.1 AP compared with using Soft-NMS individually for both baselines.

Third, the flip test together with the one-pixel shift of flipped heatmaps improved the performance consistently, i.e., by a margin of 1.2 AP and 1.0 AP ResNet-50 and HRNet-w32, respectively. The SOA technique was complementary to the flip test. We leave the analysis on the sub-pixel shift later.

Fourth, Gaussian filtering was beneficial for improving detection performance. A gain of 0.3 AP and 0.6 AP was achieved for ResNet-50 and HRNet-w32, respectively. It was complementary to the SOA technique. Using them together outperformed using each of them individually. To show the complementarity among all the techniques, we used them together in both methods. They boosted the vanilla baseline's performance by a large margin, e.g., 4.1 AP for ResNet-50 and 4.0 AP for HRNet-w32.

We evaluated the influence of the shifted pixel for the flipped heatmap described in Section 5.4. As can be seen from the bottom rows in Table 7 and Table 8, the performance dropped significantly without shifting the flipped

Table 7: Experiments on different sub-pixel refinement techniques using the simple baseline method (Xiao et al., 2018). SOA: sub-pixel refinement by the second-order approximation. Soft-NMS: Soft Non-Maximum Suppression. SSP: the sub-pixel shift of flipped heatmaps. GF: Gaussian filtering on predicted heatmaps. Backbone network: ResNet-50 (R50).

| SOA | Soft-NMS | SSP | GF | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ | $AR^{@.5}$ | $AR^{@.75}$ | $AR^M$ | $AR^L$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| - | - | - | - | 68.5 | 89.1 | 77.2 | 64.6 | 74.4 | 76.1 | 93.1 | 83.4 | 71.2 | 82.9 |
| 0.25 | - | - | - | 70.4 | 89.3 | 77.9 | 66.3 | 76.5 | 77.2 | 93.1 | 83.7 | 72.4 | 83.8 |
| parabola | - | - | - | 71.0 | 89.3 | 78.2 | 66.9 | 77.2 | 77.5 | 93.2 | 83.8 | 72.8 | 84.1 |
| paraboloid | - | - | - | 71.0 | 89.3 | 78.5 | 66.8 | 77.2 | 77.5 | 93.1 | 84.0 | 72.8 | 84.1 |
| - | ✓ | - | - | 69.3 | 89.1 | 77.3 | 65.2 | 75.4 | 76.5 | 93.1 | 83.4 | 71.6 | 83.3 |
| 0.25 | ✓ | - | - | 70.9 | 89.3 | 78.0 | 66.7 | 77.2 | 77.4 | 93.1 | 83.7 | 72.6 | 84.1 |
| parabola | ✓ | - | - | 71.4 | 89.4 | 78.4 | 67.1 | 77.6 | 77.7 | 93.2 | 83.9 | 72.9 | 84.4 |
| paraboloid | ✓ | - | - | 71.4 | 89.3 | 78.5 | 67.1 | 77.6 | 77.7 | 93.1 | 84.0 | 72.9 | 84.4 |
| - | - | 1 | - | 69.7 | 89.5 | 78.5 | 65.9 | 75.4 | 76.7 | 93.3 | 84.1 | 72.2 | 83.1 |
| 0.25 | - | 1 | - | 71.6 | 89.8 | 79.4 | 67.7 | 77.5 | 77.8 | 93.5 | 84.5 | 73.4 | 84.0 |
| parabola | - | 1 | - | 72.1 | 89.8 | 79.7 | 68.2 | 78.1 | 78.1 | 93.5 | 84.7 | 73.8 | 84.3 |
| paraboloid | - | 1 | - | 72.2 | 89.7 | 79.7 | 68.2 | 78.2 | 78.2 | 93.4 | 84.7 | 73.8 | 84.4 |
| - | - | - | ✓ | 68.8 | 88.9 | 77.5 | 64.9 | 74.7 | 76.2 | 93.0 | 83.6 | 71.3 | 82.9 |
| 0.25 | - | - | ✓ | 70.7 | 89.4 | 78.2 | 66.6 | 76.8 | 77.3 | 93.2 | 84.0 | 72.6 | 83.9 |
| parabola | - | - | ✓ | 71.3 | 89.4 | 78.6 | 67.0 | 77.4 | 77.6 | 93.1 | 84.0 | 72.9 | 84.1 |
| paraboloid | - | - | ✓ | 71.2 | 89.4 | 78.5 | 67.0 | 77.4 | 77.6 | 93.2 | 84.0 | 72.9 | 84.1 |
| - | ✓ | 1 | ✓ | 70.8 | 89.7 | 78.8 | 66.9 | 76.6 | 77.2 | 93.5 | 84.2 | 72.8 | 83.5 |
| 0.25 | ✓ | 1 | ✓ | 72.2 | 89.9 | 79.6 | 68.3 | 78.2 | 78.1 | 93.6 | 84.6 | 73.7 | 84.2 |
| parabola | ✓ | 1 | ✓ | 72.6 | 89.9 | **79.8** | 68.6 | 78.5 | 78.3 | 93.5 | 84.7 | 74.0 | 84.4 |
| paraboloid | ✓ | 1 | ✓ | 72.6 | 89.9 | 79.7 | 68.6 | 78.5 | 78.3 | 93.5 | 84.6 | 74.0 | 84.4 |
| parabola | ✓ | 0.8 | ✓ | **72.8** | **89.9** | 79.7 | **68.7** | 78.9 | **78.5** | **93.6** | **84.7** | **74.1** | 84.7 |
| parabola | ✓ | 0.6 | ✓ | 72.7 | 89.9 | 79.7 | 68.4 | **79.1** | 78.4 | 93.5 | 84.7 | 73.8 | **84.9** |
| parabola | ✓ | 0.4 | ✓ | 72.3 | 89.9 | 79.5 | 67.9 | 78.9 | 78.1 | 93.5 | 84.6 | 73.4 | 84.8 |
| parabola | ✓ | 0.2 | ✓ | 71.7 | 89.9 | 79.2 | 67.1 | 78.3 | 77.5 | 93.4 | 84.3 | 72.7 | 84.3 |
| parabola | ✓ | 0 | ✓ | 70.7 | 89.8 | 78.6 | 66.1 | 77.5 | 76.7 | 93.4 | 84.0 | 71.7 | 83.7 |

heatmap, i.e., from 72.6 AP to 70.7 AP for ResNet-50, and from 75.7 AP to 73.8 AP for HRNet-w32, since the flipped heatmap was not aligned with the original one. Generally, increasing the shifted pixel from zero to 0.8 consistently improved the detection accuracy. It saturated at a shift of 0.8 pixels and then dropped at a shift of one pixel. We used the sub-pixel refinement techniques in our submission to the 2019 COCO Keypoint Detection Challenge.

**Remarks**: 1) Each sub-pixel refinement technique has a positive but slightly different influence on the performance, i.e., SOA ≥ SSP ≥ Soft-NMS ≥ GF; 2) the proposed SOA and SSP are much better than their vanilla counterparts due to the closed-form and sub-pixel level approximation; and 3) these techniques are complementary since they are carried out in subsequent steps for unique and explicit purposes.

### 6.2.4 Inference Time Analysis

We also compared the inference time of different sub-pixel refinement techniques and the proposed CM module. Specifically, we chose HRNet-w32 as the backbone network. We divided the inference process into three parts, i.e., 1) the network forward process (denoting "Network Forward"); 2) the process of getting final keypoint coordinates from predicted heatmaps (denoting "Map → Coord"); and 3) the process of calculating the evaluation metrics of the keypoint detection results (denoting "Evaluation"), and recorded their inference time separately. For each setting, we ran the model three times and calculated the average inference time for single person instance. The results are summarized in Table 9.

As can be seen, the inference time of different post-processing techniques mainly differs in the process of "Map → Coord". For example, replacing the NMS in the default setting of the baseline model (Xiao et al., 2018) with the proposed soft-NMS only increased by 0.04 millisecond (ms). SOA and SSP increased the inference time slightly, i.e., from 3.85 ms to 4.13 ms and from 4.73 ms to 4.99 ms. However, Gaussian filtering required much more computations and increased the inference time from 4.13 ms to 4.73 ms. It is noteworthy that the process of "Map → Coord" cost more inference time than "Network Forward" since it was mainly carried out on CPU while network forward computation was carried on GPU. Adding a CM module on HRNet-w32 only increased the network forward time slightly, i.e., about 0.05 ms. However, since the heatmaps generated by CCM were two times larger than those by the baseline model,

Table 8: Experiments on different settings of sub-pixel refinement tricks using the High-Resolution Network (Sun et al., 2019) (HRNet-w32). SOA: sub-pixel refinement by the second-order approximation. Soft-NMS: sub-pixel refinement by Soft Non-Maximum Suppression (Soft-NMS). SSP: the sub-pixel shift of flipped heatmaps. GF: Gaussian filtering on heatmaps.

| SOA | Soft-NMS | SSP | GF | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ | $AR^{@.5}$ | $AR^{@.75}$ | $AR^M$ | $AR^L$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| - | - | - | - | 71.7 | 90.0 | 80.0 | 67.9 | 77.6 | 78.8 | 94.0 | 85.4 | 74.2 | 85.1 |
| 0.25 | - | - | - | 73.8 | 90.4 | 80.6 | 69.7 | 79.9 | 80.0 | 94.0 | 85.9 | 75.6 | 86.3 |
| parabola | - | - | - | 74.3 | 90.4 | 81.1 | 70.3 | 80.4 | 80.4 | 94.0 | 86.2 | 76.0 | 86.5 |
| paraboloid | - | - | - | 74.3 | 90.4 | 81.2 | 70.3 | 80.4 | 80.4 | 94.1 | 86.2 | 76.0 | 86.5 |
| - | ✓ | - | - | 72.6 | 90.0 | 80.0 | 68.5 | 78.7 | 79.2 | 94.0 | 85.5 | 74.6 | 85.6 |
| 0.25 | ✓ | - | - | 74.3 | 90.4 | 80.7 | 70.1 | 80.5 | 80.3 | 94.0 | 86.0 | 75.8 | 86.6 |
| parabola | ✓ | - | - | 74.7 | 90.4 | 81.2 | 70.5 | 80.9 | 80.5 | 94.0 | 86.2 | 76.1 | 86.7 |
| paraboloid | ✓ | - | - | 74.7 | 90.4 | 81.3 | 70.5 | 80.9 | 80.6 | 94.1 | 86.3 | 76.1 | 86.8 |
| - | - | 1 | - | 72.7 | 90.6 | 81.1 | 68.8 | 78.6 | 79.4 | 94.3 | 86.3 | 75.0 | 85.6 |
| 0.25 | - | 1 | - | 74.7 | 90.7 | 82.3 | 70.6 | 80.8 | 80.5 | 94.3 | 87.1 | 76.1 | 86.6 |
| parabola | - | 1 | - | 75.2 | 90.8 | 82.5 | 71.1 | 81.3 | 80.8 | 94.3 | 87.1 | 76.5 | 86.8 |
| paraboloid | - | 1 | - | 75.2 | 90.8 | 82.4 | 71.1 | 81.3 | 80.8 | 94.3 | 87.0 | 76.5 | 86.8 |
| - | - | - | ✓ | 72.3 | 90.3 | 80.3 | 68.4 | 78.3 | 79.1 | 94.1 | 85.6 | 74.5 | 85.5 |
| 0.25 | - | - | ✓ | 74.3 | 90.4 | 81.2 | 70.1 | 80.5 | 80.2 | 94.0 | 86.1 | 75.8 | 86.4 |
| parabola | - | - | ✓ | 74.8 | 90.4 | 81.4 | 70.7 | 81.0 | 80.5 | 94.0 | 86.1 | 76.2 | 86.7 |
| paraboloid | - | - | ✓ | 74.7 | 90.5 | 81.4 | 70.7 | 80.9 | 80.5 | 94.0 | 86.2 | 76.2 | 86.6 |
| - | ✓ | 1 | ✓ | 74.2 | 90.6 | 81.5 | 70.1 | 80.3 | 80.1 | 94.2 | 86.5 | 75.7 | 86.3 |
| 0.25 | ✓ | 1 | ✓ | 75.5 | 90.8 | 82.6 | 71.4 | 81.7 | 80.9 | 94.3 | 87.2 | 76.7 | 86.9 |
| parabola | ✓ | 1 | ✓ | 75.7 | 90.8 | 82.8 | 71.6 | 81.9 | 81.1 | 94.3 | 87.3 | 76.8 | 87.1 |
| paraboloid | ✓ | 1 | ✓ | 75.7 | 90.8 | 82.7 | 71.6 | 81.9 | 81.1 | 94.3 | 87.3 | 76.8 | 87.1 |
| parabola | ✓ | 0.8 | ✓ | **76.0** | 90.8 | 82.8 | **71.8** | 82.3 | **81.3** | **94.3** | **87.3** | **77.0** | 87.4 |
| parabola | ✓ | 0.6 | ✓ | 76.0 | **90.9** | **82.9** | 71.6 | **82.5** | 81.3 | 94.3 | 87.3 | 76.9 | **87.5** |
| parabola | ✓ | 0.4 | ✓ | 75.6 | 90.8 | 82.7 | 71.1 | 82.2 | 80.9 | 94.2 | 87.2 | 76.4 | 87.3 |
| parabola | ✓ | 0.2 | ✓ | 74.9 | 90.8 | 82.2 | 70.3 | 81.6 | 80.3 | 94.2 | 87.0 | 75.7 | 86.9 |
| parabola | ✓ | 0 | ✓ | 73.8 | 90.7 | 81.8 | 69.2 | 80.7 | 79.4 | 94.2 | 86.6 | 74.6 | 86.2 |

Table 9: Inference time (millisecond, i.e., ms) of different sub-pixel refinement techniques and the proposed CM module. Network Forward: the forward process of different networks. Map → Coord: the process of getting final keypoint coordinates from predicted heatmaps. Evaluation: the process of calculating the evaluation metrics of the keypoint detection results. Backbone network: HRNet-w32.

| Model | SOA | Soft-NMS | SSP | GF | Network Forward | Map → Coord | Eval-uation |
|---|---|---|---|---|---|---|---|
| | parabola | ✓ | 0.4 | ✓ | 2.44 | 4.99 | 0.43 |
| | parabola | ✓ | 1 | ✓ | **2.43** | 4.73 | 0.43 |
| Baseline | parabola | ✓ | 1 | - | 2.45 | 4.13 | 0.43 |
| | 0.25 | ✓ | 1 | - | 2.46 | 3.85 | 0.43 |
| | 0.25 | - | 1 | - | 2.45 | **3.81** | 0.41 |
| **CCM** | parabola | ✓ | 0.4 | ✓ | 2.50 | 9.30 | **0.39** |
| | 0.25 | - | 1 | - | 2.51 | 7.53 | **0.39** |

the process of "Map → Coord" became much slower, i.e., about two times. It is worth further study to find a fast GPU implementation for this process.

## 6.3 Comparison with state-of-the-art methods

The results of CCM and SOTA methods on the COCO minival are summarized in Table 10. CCM was trained on the COCO dataset without using the external AI Challenger dataset. We evaluated the proposed approach on three groups of backbone networks, i.e., the small ones including ShuffleNet-v2 and MobileNet-v2, the medium ones including ResNet-50 and HRNet-w32, and the large ones including ResNet-152 and HRNet-w48, respectively. As can be seen, our small model based on ShuffleNet-v2 and MobileNet-v2 significantly improved the detection accuracy compared with the baseline model (Xiao et al., 2018), i.e., a gain of 4.0 AP and 4.1 AP, respectively. The improvement is at the cost of 5% ∼ 10% more parameters and about 15% more GFLOPs, which are affordable. CCM also outperformed the Hourglass model (Newell et al., 2016) and was comparable with the CPN (Chen et al., 2018b) based on ResNet-50.

As for the medium backbone networks, simple baseline method (Xiao et al., 2018) achieved similar performance using ResNet-50 and ResNeXt-50, but they were inferior to the High-Resolution Network (Sun et al., 2019) based on HRNet-w32. The proposed CCM based on ResNet-50

Table 10: Comparisons of CAPE-Net and SOTA methods on the COCO minival set.

| Method | | Backbone | #Params | GFLOPs | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Input size: $256 \times 192$ | | | | | |
| Baseline (Xiao et al., 2018) | | ShuffleNet-v2 | 8.78M | 4.07 | 63.9 | 86.6 | 71.4 | 60.8 | 70.1 | 70.3 |
| Baseline (Xiao et al., 2018) | Small | MobileNet-v2 | 9.57M | 4.17 | 64.3 | 86.3 | 72.2 | 60.9 | 70.9 | 70.5 |
| **CCM** | | ShuffleNet-v2 | 9.72M | 4.75 | 67.9 | **88.4** | 75.3 | 63.8 | 74.0 | 74.2 |
| **CCM** | | MobileNet-v2 | 10.1M | 4.80 | **68.4** | 88.0 | **75.7** | **64.2** | **74.5** | **74.3** |
| | | | | | Input size: $256 \times 192$ | | | | | |
| Hourglass (Newell et al., 2016) | | 8xHourglass | 25.1M | 14.3 | 66.9 | - | - | - | - | - |
| CPN (Chen et al., 2018b) | | ResNet-50 | 27.0M | 6.20 | 69.4 | - | - | - | - | - |
| Baseline (Xiao et al., 2018) | | ResNet-50 | 34.0M | 8.20 | 70.4 | 88.6 | 78.3 | 67.1 | 77.2 | 76.3 |
| Baseline (Xiao et al., 2018) | Medium | ResNeXt-50 | 33.5M | 8.35 | 70.6 | 88.9 | 77.9 | 67.2 | 77.5 | 76.5 |
| HRNet (Sun et al., 2019) | | HRNet-w32 | 28.5M | 7.68 | 74.4 | 90.5 | 81.9 | 70.8 | 81.0 | 79.8 |
| **CCM** | | ResNet-50 | 40.7M | 45.4 | 74.3 | 90.3 | 81.3 | 70.0 | 80.6 | 79.6 |
| **CCM** | | ResNeXt-50 | 40.1M | 45.4 | 74.5 | 90.4 | 81.2 | 70.0 | 81.0 | 79.6 |
| **CCM** | | HRNet-w32 | 28.6M | 7.92 | **76.7** | **91.1** | **83.5** | **72.5** | **83.0** | **81.8** |
| | | | | | Input size: $384 \times 288$ | | | | | |
| Baseline (Xiao et al., 2018) | | ResNet-152 | 68.6M | 35.9 | 74.3 | 89.6 | 81.1 | 70.5 | 79.7 | 79.7 |
| HRNet (Sun et al., 2019) | Large | HRNet-w48 | 63.6M | 35.4 | 76.3 | 90.8 | 82.9 | 72.3 | 83.4 | 81.2 |
| **CCM** | | ResNet-152 | 63.5M | 40.1 | 76.7 | 91.2 | 83.4 | 72.4 | 83.2 | 81.7 |
| **CCM** | | HRNet-w48 | 63.7M | 36.6 | **77.5** | **91.2** | **83.6** | **73.0** | **84.0** | **82.3** |

achieved a 74.3 AP and outperformed other models with the same input size and backbone network, for example, a gain of 3.9 AP over the simple baseline method (Xiao et al., 2018). Replacing ResNet-50 to ResNeXt-50 leads to a slightly better result, i.e., from 74.3 AP to 74.5 AP. When using the HRNet-w32 as the backbone network, our CCM model increased the AP from 74.4 to 76.7 and achieved the best performance among all the models. Note that CCM based on ResNet-50 or ResNeXt-50 has more parameters and GFLOPs than the baseline model since it uses dilated convolutions to increase the feature map size and three CMs in the decoder. As for the one based on HRNet-w32, it has roughly the same amount of parameters and GFLOPs as its counterpart since only one CM was used.

Our large model based on ResNet-152 with input size $384 \times 288$ achieved a gain of 2.4 AP over the simple baseline model (Xiao et al., 2018) and a gain of 0.4 AP over the recent HRNet-w48 model (Sun et al., 2019). For example, CCM outperformed HRNet-w48 by a margin of 0.4 AP and 0.5 AP at the threshold 0.5 and 0.75. However, CCM was inferior to HRNet-w48 for large person instances, i.e., a drop of 0.2 AP. One possible explanation is that HRNet-w48 learned a high-resolution and discriminative feature representation by integrating the features from different scales. We attached a CM to HRNet-w48 and used the sub-pixel refinement techniques for postprocessing. It further improved the detection accuracy of HRNet-w48 from 76.3 AP to 77.5 AP. Besides, the result of large person instances was improved by 0.6 AP. Our CCM model achieved the best performance among all other models based on comparable backbone networks. Besides, both models have nearly the

same parameters and GFLOPs as the baseline models since we decreased the number of filters in the CM for ResNet-152 and only used one CM for HRNet-w48. These results demonstrate the effectiveness of the proposed CCM, training strategies, and sub-pixel refinement techniques.

The results of CCM and SOTA methods on the COCO test-dev set are summarized in Table 11. The CCM using ResNet-152 as the backbone network trained on the COCO dataset outperformed the baseline model (Xiao et al., 2018) by 2.1 AP. It also outperformed the recent HRNet-w48 model (Sun et al., 2019) by 0.3 AP. Using HRNet-w48 as the backbone network, CCM outperformed the vanilla HRNet-w48 model, MSPN (Li et al., 2019), and DARK (Zhang et al., 2020) by a margin of 1.1 AP, 0.5 AP, and 0.4 AP, respectively. It was even better than the ensemble simple baseline models trained with the external AI Challenger dataset. After joint-training with this external dataset, CCM based on ResNet-152 outperformed both the simple baseline method and HRNet-w48. Besides, replacing the backbone network from ResNet-152 to HRNet-w48, the performance was further improved by 0.7 AP. It was better than the recently proposed method DARK* (Zhang et al., 2020) by a margin of 0.6 AP using the same person detection results and comparable with the ensemble models MSPN+* (Li et al., 2019) (the champion of the 2018 COCO Keypoint Challenge), i.e., 78.0 v.s. 78.1. Generally, the external dataset brought about 1.5 AP improvement, which mainly arose from the AP at the larger threshold, demonstrating that training on more diverse poses helped the model to learn discriminative features and improve the location accuracy. Our final ensemble models brought another 0.9 AP and set a

Table 11: Comparisons of CCM and SOTA methods on the COCO test-dev set. Input size: $353 \times 257$ for G-RMI; $320 \times 256$ for RMPE; $384 \times 288$ for CPN, Baseline, HRNet, MSPN, DARK, RSN, and CCM. The symbol "*" denotes external data, "+" denotes ensemble models, "†" and "‡" denote the champion of the 2018 and 2019 COCO Keypoint Challenge, respectively.

| Method | Backbone | #Params | GFLOPs | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ |
|---|---|---|---|---|---|---|---|---|---|
| Mask-RCNN (He et al., 2017) | ResNet-50 | - | - | 63.1 | 87.3 | 68.7 | 57.8 | 71.4 | - |
| G-RMI (Papandreou et al., 2017) | ResNet-101 | 42.6M | 57.0 | 64.9 | 85.5 | 71.3 | 62.3 | 70.0 | 69.7 |
| G-RMI* (Papandreou et al., 2017) | ResNet-101 | 42.6M | 57.0 | 68.5 | 87.1 | 75.5 | 65.8 | 73.3 | 73.3 |
| RMPE (Fang et al., 2017) | Hourglass | 28.1M | 36.7 | 72.3 | 89.2 | 79.1 | 68.0 | 78.6 | - |
| CPN (Chen et al., 2018b) | ResNet-Inception | - | - | 72.1 | 91.4 | 80.0 | 68.7 | 77.2 | 78.5 |
| CPN+ (Chen et al., 2018b) | ResNet-Inception | - | - | 73.0 | 91.7 | 80.9 | 69.5 | 78.1 | 79.0 |
| Baseline (Xiao et al., 2018) | ResNet-152 | 68.6M | 35.9 | 73.7 | 91.9 | 81.1 | 70.3 | 80.0 | 79.0 |
| Baseline+* (Xiao et al., 2018) | ResNet-152 | - | - | 76.5 | 92.4 | 84.0 | 73.0 | 82.7 | 81.5 |
| HRNet (Sun et al., 2019) | HRNet-w48 | 63.6M | 35.4 | 75.5 | 92.5 | 83.3 | 71.9 | 81.5 | 80.5 |
| HRNet* (Sun et al., 2019) | HRNet-w48 | 63.6M | 35.4 | 77.0 | 92.7 | 84.5 | 73.4 | 83.1 | 82.0 |
| MSPN (Li et al., 2019) | 4xResNet-50 | 71.9M | 58.7 | 76.1 | 93.4 | 83.8 | 72.3 | 81.5 | 81.6 |
| MSPN* (Li et al., 2019) | 4xResNet-50 | 71.9M | 58.7 | 77.1 | 93.8 | 84.6 | 73.4 | 82.3 | 82.3 |
| MSPN+* (Li et al., 2019)† | 4xResNet-50 | - | - | 78.1 | 94.1 | 85.9 | 74.5 | 83.3 | 83.1 |
| DARK (Zhang et al., 2020) | HRNet-w48 | 63.6M | 32.9 | 76.2 | 92.5 | 83.6 | 72.5 | 82.4 | 81.1 |
| DARK* (Zhang et al., 2020) | HRNet-w48 | 63.6M | 32.9 | 77.4 | 92.6 | 84.6 | 73.6 | 83.7 | 82.3 |
| RSN (Cai et al., 2020) | 4xRSN-50 | 111.8M | 65.9 | 78.6 | 94.3 | 86.6 | 75.5 | 83.3 | 83.8 |
| RSN+ (Cai et al., 2020)‡ | 4xRSN-50 | - | - | **79.2** | **94.4** | **87.1** | **76.1** | **83.8** | **84.1** |
| **CCM** | ResNet-152 | 63.5M | 40.1 | 75.8 | 92.7 | 83.4 | 71.8 | 81.5 | 80.9 |
| **CCM** | HRNet-w48 | 63.7M | 36.6 | 76.6 | 92.8 | 84.1 | 72.6 | 82.4 | 81.7 |
| **CCM*** | ResNet-152 | 63.5M | 40.1 | 77.3 | 93.0 | 84.8 | 73.3 | 83.1 | 82.3 |
| **CCM*** | HRNet-w48 | 63.7M | 36.6 | 78.0 | 93.4 | 85.1 | 74.0 | 83.6 | 83.0 |
| **CCM+*** | HRNet-w48 | - | - | **78.9** | **93.8** | **86.0** | **75.0** | **84.5** | **83.6** |

new state-of-the-art on this benchmark, i.e., 78.9 AP. It was comparable with RSN+ (Cai et al., 2020) (the champion of the 2019 COCO Keypoint Challenge), which used a better person detector and had 59.8 AP for the person category on COCO val2017.

**Remarks**: 1) Our CCM model consistently outperforms the baseline models using small, medium, and large backbone networks, but the gain becomes smaller with the increasing of the model capacity, i.e., MobileNet ≈ ShuffleNet ≈ ResNet-50 ≈ ResNeXt-50 ≥ HRNet-w32 ≈ ResNet-152 ≥ HRNet-w48. It is reasonable because "bigger" backbone networks themselves have stronger representation capacity, thereby the impact of CM decreases accordingly; and 2) our CCM model benefits from the effective CM modules to model the context information, the efficient training strategies to learn discriminative features, and the sub-pixel refinement techniques to locate keypoints accurately, in a collaborative and complementary manner.

6.4 Subjective visual inspection and discussion

We presented some visual examples of the keypoint detection results by using the CCM model based on HRNet-w48 on the COCO minival set and PoseTrack validation set in Figure 10 and Figure 11, respectively. As can be seen from, our model could handle various poses. Besides, it also successfully inferred the occluded keypoints (self-occluded

or occluded by other objects). In the bottom two rows, we presented the detection results on multiple person instances within each image, which were also promising. Our model could handle small instances, blurry ones, low-light images as well as various occlusions. To see how CCM achieved the performance, we conducted an experiment to visually inspect the learned features by the network.

We overlaid the output feature maps from the CMs in CCM on the input images to inspect what has been learned by the network. Figure 12 shows the feature maps learned by the CMs at different levels. "HDC k" stands for the feature maps from the $k$th dilated convolutional layer in the hybrid-dilated convolutional branches. "HDC x SE" stands for the first term in Eq. (4). 'Deconv' stands for the outputted feature maps from the CMs, i.e., the left side of Eq. (4). "Level k" stands for the index of CM in CCM. The keypoints belonging to the left (right) body were connected by red (blue) lines. The left-right symmetric keypoints were connected by yellow lines. The predicted invisible keypoints of occluded body parts were indicated by red arrows.

As can be seen, HDC paid more attention to the context with the increase of the dilation rate. CCM learned to identify easy keypoints and inferred hard keypoints progressively through the cascaded CMs. Please check the arrows and ellipses. Moreover, 1) CCM probably learned the body configuration by inferring the invisible keypoints and potential poses conditioned on the visible keypoints as shown in the first row; 2) CCM also learned the body symmetry, e.g.,

Fig. 10: Some visual examples of the keypoint detection results on the COCO minival set (Lin et al., 2014).

there were two legs and arms in the human body, as shown in the left part of the second row; 3) CCM could also handle blurry images and infer the head and right arms as shown in the right part of the second row.

To illustrate the effectiveness of CCM for handling occlusions, we manually added masks on some body parts, *e.g.*, the head, hip, and ankle as supplements to the existing occlusions, as shown in Figure 13. CCM first recognized and located the human bodies using the encoder (Res5), then detected different body parts (Deconv1) to help identify some distinct keypoints (Deconv2). In the final stage (Deconv3), CCM predicted the difficult occluded keypoints using the context information of identified body parts and keypoints. To further investigate the effectiveness of CCM for handling severe occlusions, we manually masked out either the upper or lower body of each person as shown in the first column of Figure 13. As can be seen, although the predicted human poses were different from the ground truth annotations at the masked regions, CCM showed the ability to learn human body configuration and infer reasonable invisible keypoints. Moreover, the results also confirm that CCM detected those

distinct keypoints in the earlier stage while predicting the difficult occluded ones in the later stage.

**Remarks**: 1) Empirically, CCM's detection process follows a "Localization → Componentization → Identification → Prediction" routine, similar to the procedure that humans detect human keypoints in occluded settings (Section 3.1); and 2) CCM probably has learned the body configuration such as symmetric body parts and reasonable distances between adjacent keypoints, evidenced by visual examples.

6.5 Empirical Studies on Annotations

As we know that occluded keypoints are more difficult to be detected than visible ones, we are wondering whether the invisible keypoint annotations have the same impact on the model or not, compared with the same amount of visible keypoint annotations? To this end, we conducted an experiment to gain some insight into the keypoint annotations.

We constructed two training sets based on the COCO training set. The first one was COCO-I which was obtained by removing all the annotations of invisible keypoints in the

Fig. 11: Some visual examples of the keypoint detection results on the PoseTrack validation set (Andriluka et al., 2018).

Table 12: Comparison between CCM-I and CCM-V on the COCO minival set. Please refer to Section 6.5.

| Model | Backbone | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ | $AR^{@.5}$ | $AR^{@.75}$ | $AR^M$ | $AR^L$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| CCM-I | ResNet-50 | 70.4 | 89.1 | 76.7 | 66.0 | 77.0 | 77.1 | 93.4 | 82.9 | 72.2 | 83.8 |
| CCM-V | ResNet-50 | 73.7 | 90.0 | 80.9 | 69.5 | 79.9 | 79.2 | 93.5 | 85.5 | 74.7 | 85.5 |
| CCM | ResNet-50 | 73.8 | 90.2 | 80.9 | 69.6 | 80.1 | 79.3 | 93.7 | 85.6 | 74.8 | 85.7 |

Table 13: Comparison between CCM-I and CCM-V on the visible and invisible keypoints in the COCO minival set.

| Model | Backbone | Kpt. Type | $AP$ | $AP^{@.5}$ | $AP^{@.75}$ | $AP^M$ | $AP^L$ | $AR$ | $AR^{@.5}$ | $AR^{@.75}$ | $AR^M$ | $AR^L$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CCM-I | ResNet-50 | Visible | 79.3 | 94.6 | 86.7 | 77.4 | 82.2 | 82.4 | 95.5 | 88.7 | 80.2 | 85.5 |
| | | Invisible | 40.2 | 64.2 | 39.5 | 40.2 | 42.3 | 47.8 | 69.8 | 48.0 | 45.1 | 53.2 |
| CCM-V | ResNet-50 | Visible | 79.4 | 94.7 | 87.0 | 76.9 | 82.4 | 82.0 | 95.3 | 88.1 | 79.5 | 85.4 |
| | | $\pm\Delta$ | 0.1 | 0.1 | 0.3 | -0.5 | 0.2 | -0.4 | -0.2 | -0.6 | -0.7 | -0.1 |
| | | Invisible | 55.0 | 79.4 | 57.1 | 54.7 | 57.8 | 61.4 | 82.1 | 63.9 | 58.8 | 66.9 |
| | | $\pm\Delta$ | 14.8 | 15.2 | 17.6 | 14.5 | 15.5 | 13.6 | 12.3 | 15.9 | 13.7 | 13.7 |

Table 14: Comparison between CCM-I and CCM-V on different types of instances w.r.t. the numbers of annotated keypoints.

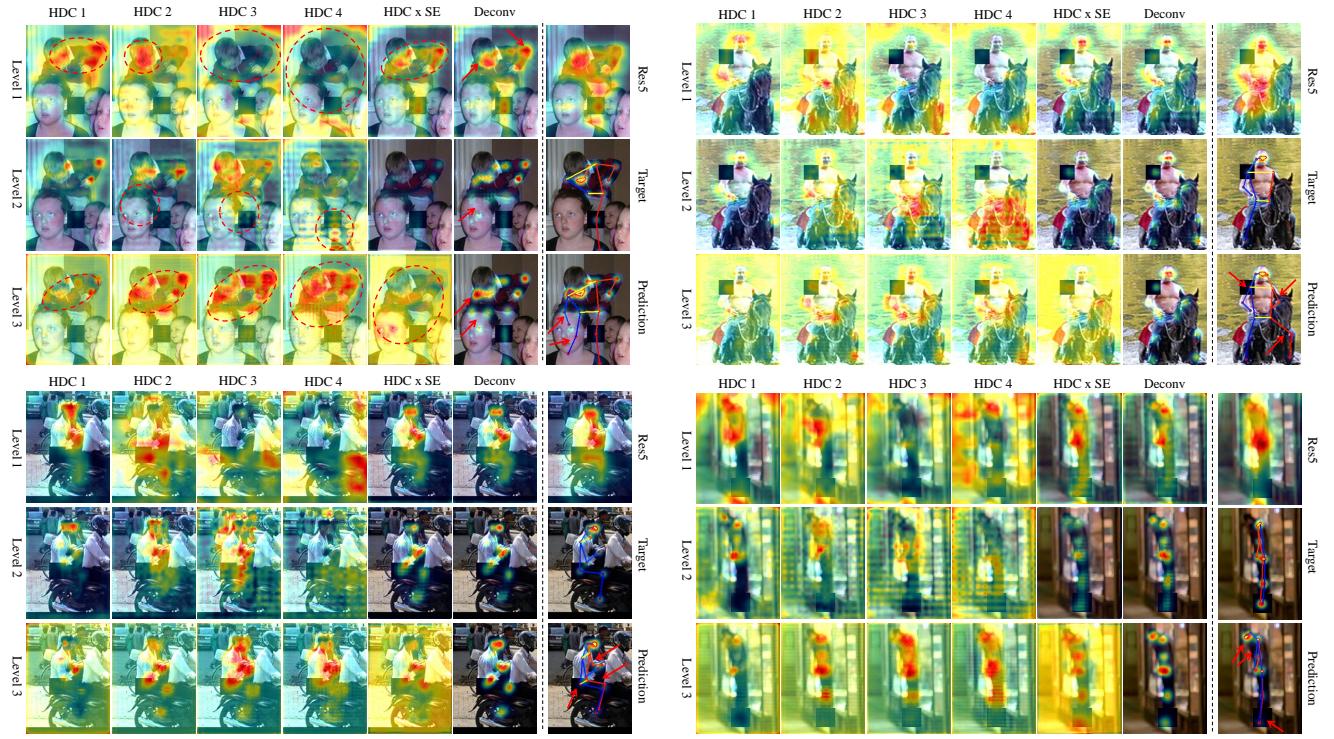| Model | Metric | Number of annotated keypoints per instance | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
| CCM-I | AP | 32.7 | 25.0 | 33.4 | 38.9 | 48.1 | 52.5 | 60.2 | 63.4 | 67.1 | 68.0 | 74.2 | 77.8 | 79.0 | 81.6 | 82.6 | 88.8 | 90.9 |
| CCM-V | AP | 33.9 | 23.9 | 37.6 | 43.6 | 48.2 | 58.1 | 64.2 | 70.0 | 72.6 | 72.8 | 77.9 | 81.5 | 81.7 | 83.4 | 84.1 | 89.3 | 91.2 |
| | $\pm\Delta$ | 1.2 | -1.1 | **4.2** | **4.7** | 0.1 | **5.6** | **4.0** | **6.6** | **5.5** | **4.8** | 3.7 | 3.7 | 2.7 | 1.8 | 1.5 | 0.5 | 0.3 |

Fig. 12: Visualization of the feature maps from the CMs in CCM based on ResNet-50. "HDC k" stands for the feature maps from the $k$th dilated convolutional layer in the hybrid-dilated convolutional branches. "HDC x SE" stands for the first term in Eq. (4). "Deconv" stands for the outputted feature maps from the CMs, i.e., the left side of Eq. (4).
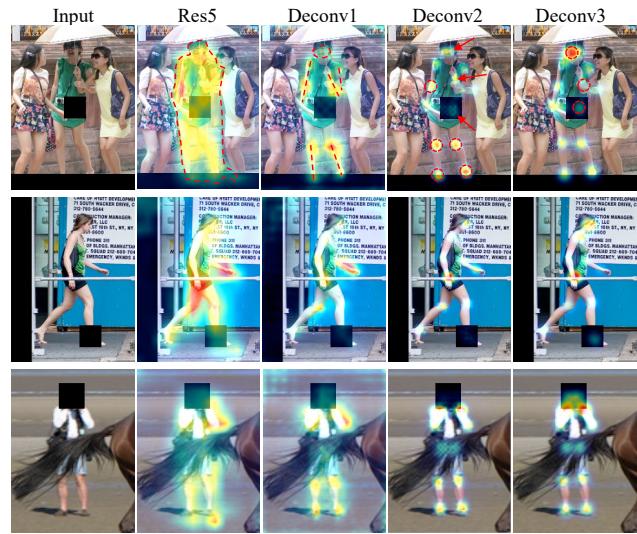


Fig. 13: Visualization of the feature maps learned by CCM based on ResNet-50 at different stages. "Deconv" stands for the outputted feature maps from the CMs, i.e., the left side of Eq. (4). One or several of keypoints of each person is manually occluded by a mask.
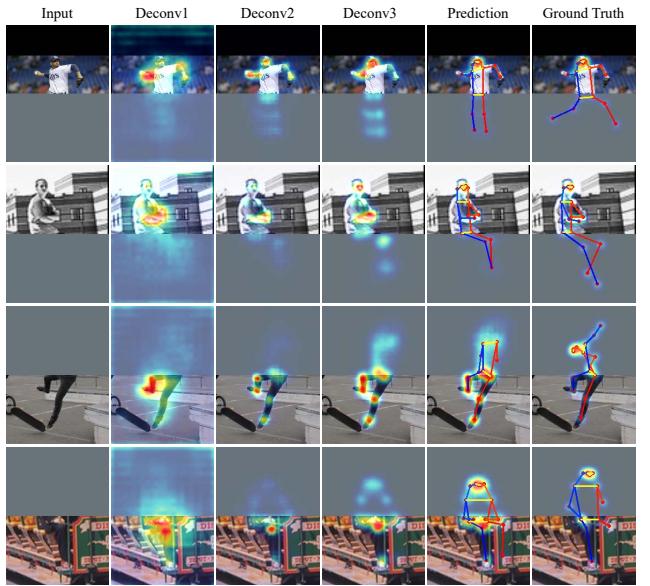


Fig. 14: Visualization of the feature maps learned by CCM based on ResNet-50 at different stages. Either the upper or lower body of each person is manually masked out.

original training set. The second one was COCO-V which was obtained by removing the same amount of visible keypoint annotations randomly. The resulting keypoints without annotations were treated as unlabeled. Then, we trained the proposed CCM using the ResNet-50 as the backbone encoder on COCO-I and COCO-V. They were denoted as "CCM-I" and "CCM-V", respectively. Their results on the COCO minival set are summarized in the first two rows of Table 12. As a reference, we also listed the model trained on the original COCO training set in the bottom row.

As can be seen, the scores of CCM-V dropped marginally compared with CCM. However, the scores of CCM-I dropped significantly by a large margin compared with CCM-V and CCM. These results confirm that *annotating invisible keypoints matters*, *i.e.*. The invisible keypoint annotations contributed more to the model than the same amount of visible ones. Since it is more difficult to detect invisible keypoints, the invisible keypoint annotations provide stronger supervisory signals to the model than the visible ones did. To infer an invisible keypoint with such supervision, it probably learned useful features from the context since the keypoint itself was invisible. In this way, the model could learn the knowledge of body configuration implicitly, e.g., as a form of discriminative feature for each category of keypoints.

To further analyze the impact of invisible keypoint annotations, we calculated the indexes on visible and invisible keypoints, respectively. We removed all the invisible keypoint annotations from the COCO minival set. The resulting minival set was used to evaluate the model's performance on the visible keypoints. Similarly, we also removed all the visible keypoint annotations from the COCO minival set to evaluate the model's performance on the invisible keypoints. We used the ground truth bounding boxes as the person detection results. The results are listed in Table 13. $\pm\Delta$ denoted the gain of CCM-V over CCM-I.

Unsurprisingly, CCM-V achieved better results on invisible keypoints compared with CCM-I, demonstrating that the gains of CCM-V over CCM-I in Table 12 mainly arose from the invisible ones. Note that although CCM-I was trained without using the annotations of invisible keypoints, it still obtained 40.2 AP on them, demonstrating that it had a bit of generalization on predicting the invisible keypoints. Since there were distinct appearance differences between visible keypoints and invisible ones of the same category and the model did not get any supervisory signal from the invisible keypoints, it implies that CCM probably learned useful features from contextual keypoints to infer the invisible ones. Using the invisible annotations, CCM achieved better performance by exploiting its representation capacity.

We also reported APs of the two models on instances with different numbers of annotated keypoints from the COCO minival set. The results are shown in Table 14. As can be seen, the gains mainly arose from the occluded

instances, for example, instances with less than 10 annotated keypoints. It implies that the invisible keypoint annotations helped the model to learn the body configuration for inferring the invisible keypoints and handling occlusions. Some visual results were presented in Figure 12 and Figure 13.

Although annotating invisible keypoints is more difficult than visible ones, the above results confirm that the invisible keypoint annotations are more valuable. Consequently, we could improve our model by 1) annotating more occluded keypoints to train a better model; 2) exploiting the active learning strategy to identify hard keypoints that should be annotated to continuously improve the model; 3) developing effective multi-view learning algorithms to utilize the complementary information between different views of data.

**Remarks**: 1) The invisible keypoint annotations have a larger impact on the model than the same amount of visible ones; and 2) even without the invisible annotations, the proposed CCM model is still able to generalize to the invisible keypoints.

# 7 Conclusion

In this paper, we address the human keypoint detection problem by devising a new cascaded context mixer-based neural network (CCM), which has a strong representative capacity to simultaneously model the spatial and channel context information. We also propose three efficient training strategies including a hard-negative person detection mining strategy to migrate the mismatch between training and testing, a joint-training strategy to use abundant unlabeled samples by knowledge distilling, and a joint-training strategy to exploit external data with heterogeneous labels. They collaboratively enable CCM to learn discriminative features from abundant and diverse poses. Besides, we present four post-processing techniques to refine predictions at the sub-pixel level accuracy. These complementary techniques are carried out sequentially during the inference phase for unique and explicit purposes, which further improve the detection accuracy. Our CCM model consistently outperforms public state-of-the-art models with various backbone networks by a large margin. We empirically show that CCM's detection process is similar to humans in occluded settings and probably learns human body configuration. Moreover, we also identify that invisible keypoint annotations have a larger impact on the model than the same amount of visible ones and present some promising research topics.

# References

Andriluka M, Iqbal U, Ensafutdinov E, Pischulin L, Milan A, Gall J, B S (2018) PoseTrack: A benchmark for human pose estimation and tracking. In: Proceedings of the IEEE

Conference on Computer Vision and Pattern Recognition (CVPR)

Baradel F, Wolf C, Mille J, Taylor GW (2018) Glimpse clouds: Human activity recognition from unstructured feature points. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 469–478

Biederman I (1987) Recognition-by-components: a theory of human image understanding. Psychological review 94(2):115

Cai Y, Wang Z, Luo Z, Yin B, Du A, Wang H, Zhou X, Zhou E, Zhang X, Sun J (2020) Learning delicate local representations for multi-person pose estimation. In: Proceedings of the European Conference on Computer Vision (ECCV)

Cao Z, Simon T, Wei SE, Sheikh Y (2017) Realtime multi-person 2d pose estimation using part affinity fields. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7291–7299

Chen LC, Papandreou G, Kokkinos I, Murphy K, Yuille AL (2018a) Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE transactions on pattern analysis and machine intelligence 40(4):834–848

Chen Y, Wang Z, Peng Y, Zhang Z, Yu G, Sun J (2018b) Cascaded pyramid network for multi-person pose estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7103–7112

Chen Z, Zhang J, Tao D (2020) Recursive context routing for object detection. International Journal of Computer Vision pp 1–19

Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L (2009) Imagenet: A large-scale hierarchical image database. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 248–255

Fang HS, Xie S, Tai YW, Lu C (2017) Rmpe: Regional multi-person pose estimation. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV), pp 2334–2343

Felzenszwalb P, McAllester D, Ramanan D (2008) A discriminatively trained, multiscale, deformable part model. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, pp 1–8

Girdhar R, Gkioxari G, Torresani L, Paluri M, Tran D (2018) Detect-and-track: Efficient pose estimation in videos. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 350–359

Hattori H, Lee N, Boddeti VN, Beainy F, Kitani KM, Kanade T (2018) Synthesizing a scene-specific pedestrian detector and pose estimator for static video surveillance. International Journal of Computer Vision 126(9):1027–1044

He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 770–778

He K, Gkioxari G, Dollár P, Girshick R (2017) Mask r-cnn. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV), pp 2961–2969

Holt B, Ong EJ, Cooper H, Bowden R (2011) Putting the pieces together: Connected poselets for human pose estimation. In: Proceedings of the IEEE International Conference on Computer Vision Workshops (ICCVW), IEEE, pp 1196–1201

Hossain MRI, Little JJ (2018) Exploiting temporal information for 3d human pose estimation. In: Proceedings of the European Conference on Computer Vision (ECCV), Springer, pp 69–86

Hu J, Shen L, Sun G (2018) Squeeze-and-excitation networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7132–7141

Huang S, Gong M, Tao D (2017) A coarse-fine network for keypoint localization. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV), pp 3028–3037

Ioffe S, Szegedy C (2015) Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: Proceedings of the International Conference on Machine Learning (ICML), pp 448–456

Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems, pp 1097–1105

Lee CY, Xie S, Gallagher P, Zhang Z, Tu Z (2015) Deeply-supervised nets. In: Artificial intelligence and statistics, pp 562–570

Li W, Wang Z, Yin B, Peng Q, Du Y, Xiao T, Yu G, Lu H, Wei Y, Sun J (2019) Rethinking on multi-stage networks for human pose estimation. arXiv preprint arXiv:190100148

Lin TY, Maire M, Belongie S, Hays J, Perona P, Ramanan D, Dollár P, Zitnick CL (2014) Microsoft coco: Common objects in context. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 740–755

Lin TY, Dollár P, Girshick R, He K, Hariharan B, Belongie S (2017) Feature pyramid networks for object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 2117–2125

Liu J, Shahroudy A, Xu D, Kot AC, Wang G (2018) Skeleton-based action recognition using spatio-temporal lstm network with trust gates. IEEE transactions on pattern analysis and machine intelligence 40(12):3007–3021

Liu L, Ouyang W, Wang X, Fieguth P, Chen J, Liu X, Pietikäinen M (2020) Deep learning for generic object

detection: A survey. International journal of computer vision 128(2):261–318

Ma B, Zhang J, Xia Y, Tao D (2020) Auto learning attention. Advances in Neural Information Processing Systems 33

Mazhar O, Ramdani S, Navarro B, Passama R, Cherubini A (2018) Towards real-time physical human-robot interaction using skeleton information and hand gestures. In: Proceedings of the 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), IEEE, pp 1–6

Newell A, Yang K, Deng J (2016) Stacked hourglass networks for human pose estimation. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 483–499

Newell A, Huang Z, Deng J (2017) Associative embedding: End-to-end learning for joint detection and grouping. In: Advances in Neural Information Processing Systems, pp 2277–2287

Ni B, Li T, Yang X (2017) Learning semantic-aligned action representation. IEEE transactions on neural networks and learning systems 29(8):3715–3725

Ouyang W, Zeng X, Wang X (2016) Learning mutual visibility relationship for pedestrian detection with a deep model. International Journal of Computer Vision 120(1):14–27

Papandreou G, Zhu T, Kanazawa N, Toshev A, Tompson J, Bregler C, Murphy K (2017) Towards accurate multi-person pose estimation in the wild. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 4903–4911

Papandreou G, Zhu T, Chen LC, Gidaris S, Tompson J, Murphy K (2018) Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 269–286

Paszke A, Gross S, Chintala S, Chanan G, Yang E, DeVito Z, Lin Z, Desmaison A, Antiga L, Lerer A (2017) Automatic differentiation in pytorch. In: Advances in neural information processing systems workshops

Pavlakos G, Zhou X, Daniilidis K (2018a) Ordinal depth supervision for 3d human pose estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7307–7316

Pavlakos G, Zhu L, Zhou X, Daniilidis K (2018b) Learning to estimate 3d human pose and shape from a single color image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 459–468

Pishchulin L, Insafutdinov E, Tang S, Andres B, Andriluka M, Gehler PV, Schiele B (2016) Deepcut: Joint subset partition and labeling for multi person pose estimation. In: Proceedings of the IEEE Conference on Computer Vision

and Pattern Recognition (CVPR), pp 4929–4937

Ren S, He K, Girshick R, Sun J (2015) Faster r-cnn: Towards real-time object detection with region proposal networks. In: Advances in neural information processing systems, pp 91–99

Rhodin H, Salzmann M, Fua P (2018) Unsupervised geometry-aware representation for 3d human pose estimation. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 750–767

Rogez G, Rihan J, Orrite-Uruñuela C, Torr PH (2012) Fast human pose detection using randomized hierarchical cascades of rejectors. International Journal of Computer Vision 99(1):25–52

Sun K, Xiao B, Liu D, Wang J (2019) Deep high-resolution representation learning for human pose estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 5693–5703

Sun X, Xiao B, Wei F, Liang S, Wei Y (2018) Integral human pose regression. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 529–545

Toshev A, Szegedy C (2014) Deeppose: Human pose estimation via deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 1653–1660

Varadarajan J, Subramanian R, Bulò SR, Ahuja N, Lanz O, Ricci E (2018) Joint estimation of human pose and conversational groups from social scenes. International Journal of Computer Vision 126(2-4):410–429

Wagemans J, Elder JH, Kubovy M, Palmer SE, Peterson MA, Singh M, von der Heydt R (2012) A century of gestalt psychology in visual perception: I. perceptual grouping and figure–ground organization. Psychological bulletin 138(6):1172

Wang F, Li Y (2013) Beyond physical connections: Tree models in human pose estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 596–603

Xiao B, Wu H, Wei Y (2018) Simple baselines for human pose estimation and tracking. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 466–481

Yang Q, Yang R, Davis J, Nistér D (2007) Spatial-depth super resolution for range images. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, pp 1–8

Yang W, Li S, Ouyang W, Li H, Wang X (2017) Learning feature pyramids for human pose estimation. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV), pp 1281–1290

Yang W, Ouyang W, Wang X, Ren J, Li H, Wang X (2018) 3d human pose estimation in the wild by adversarial learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp

5255–5264

Yang Y, Ramanan D (2013) Articulated human detection with flexible mixtures of parts. IEEE transactions on pattern analysis and machine intelligence 35(12):2878–2890

Zhang F, Zhu X, Dai H, Ye M, Zhu C (2020) Distribution-aware coordinate representation for human pose estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 7093–7102

Zhang H, Ouyang H, Liu S, Qi X, Shen X, Yang R, Jia J (2019a) Human pose estimation with spatial contextual information. arXiv preprint arXiv:190101760

Zhang J, Tao D (2020) Empowering things with intelligence: A survey of the progress, challenges, and opportunities in artificial intelligence of things. IEEE Internet of Things Journal

Zhang SH, Li R, et al (2019b) Pose2seg: Detection free human instance segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)