
Learning Local-Global Contextual Adaptation for Fully End-to-End Bottom-Up Human Pose Estimation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper presents a method of learning *Local-Global* Contextual Adaptation for
2 fully end-to-end and fast bottom-up human *Pose* estimation, dubbed as *LOGO-*
3 *CAP*. It is built on the conceptually simple center-offset formulation that lacks
4 inaccuracy for pose estimation. When revisiting the bottom-up human pose es-
5 timation with the thought of “thinking, fast and slow” by D. Kahneman, we in-
6 troduce a “slow keypointer” to remedy the lack of sufficient accuracy of the “fast
7 keypointer”. In learning the “slow keypointer”, the proposed *LOGO-CAP* lifts the
8 initial “fast” keypoints by offset predictions to keypoint expansion maps (KEMs)
9 to counter their uncertainty in two modules. Firstly, the local KEMs (e.g. 11×11)
10 are extracted from a low-dimensional feature map. A proposed convolutional mes-
11 sage passing module learns to “re-focus” the local KEMs to the keypoint attraction
12 maps (KAMs) by accounting for the structured output prediction nature of human
13 pose estimation, which is directly supervised by the object keypoint similarity
14 (OKS) loss in training. Secondly, the global KEMs are extracted, with a suffi-
15 ciently large region-of-interest (e.g., 97×97), from the keypoint heatmaps that
16 are computed by a direct map-to-map regression. Then, a local-global contextual
17 adaptation module is proposed to convolve the global KEMs using the learned
18 KAMs as the kernels. This convolution can be understood as the learnable offsets
19 guided deformable and dynamic convolution in a pose-sensitive way. The pro-
20 posed method is end-to-end trainable with near real-time inference speed, obtain-
21 ing state-of-the-art performance on the COCO keypoint benchmark for bottom-up
22 human pose estimation. With the COCO trained model, our *LOGO-CAP* also
23 outperforms prior arts by a large margin on the challenging OCHuman dataset.

24 1 Introduction

25 1.1 Motivation and Objective

26 Human pose is highly articulated with large structural and appearance variations. 2D human pose
27 estimation in images is a classic structured output prediction problem, and remains a challenging one
28 in computer vision and machine learning. Human pose estimation has numerous applications such
29 as people-centered image understanding, autonomous driving and Augmented Reality (AR). With
30 the recent resurgence of deep neural networks (DNNs), the performance of human pose estimation
31 has witnessed remarkable improvement [12, 3, 15, 22, 11]. This paper focuses on the deep learning
32 based problem formulation.

33 There are two deep learning based paradigms for human pose estimation in the literature. The top-
34 down paradigm consists of human detection and single human pose estimation in each detected
35 human bounding box [12]. The bottom-up paradigm also includes two components: human pose
36 keypoint detection and keypoint grouping [3]. The top-down paradigm often obtains better accuracy
37 performance, but suffers from its inferior efficiency since the computational cost of the single human

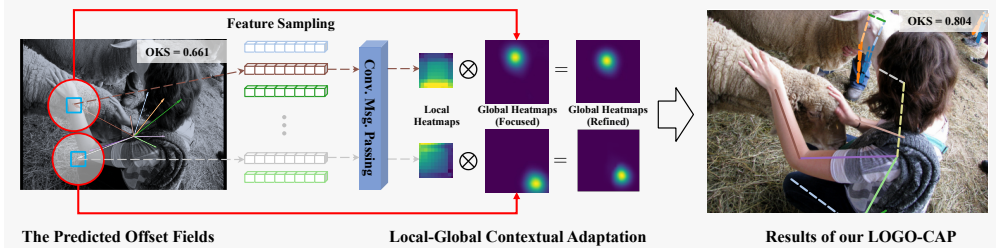


Figure 1: Illustration of the proposed LOGO-CAP for bottom-up human pose estimation. It is built on the center-offset representation. See text for detail.

pose estimation component is linearly scaled with respect to the number of detected human bounding boxes in an image. It is also largely affected by the performance of the human detection component (e.g., not handling occlusion very well). Thanks to its efficiency, especially in real-time applications, the bottom-up paradigm becomes more and more attractive. For both paradigms, state-of-the-art methods often are not fully end-to-end trained and utilize different post-hoc processing modules to improve the performance. This paper is interested in developing a fully end-to-end bottom-up paradigm and aims at bridging its performance gap with the top-down paradigm.

For the bottom-up paradigm, the recently proposed center-offset approach [6, 28, 26, 11] is a conceptually simple formulation (see the left of Fig. 1 for an illustrative example and Fig. 3 for the detailed workflow). It alleviates the need of sophisticated keypoint grouping. When introducing human keypoints centers (i.e., anchors) by treating objects as points [35], it encodes a human pose as a star structure using the offset vectors of human keypoints relative to the anchors. The main challenge of the center-offset regression paradigm lies in the difficulty of accurately learning offset vectors with large structural variations, especially the long-range ones, which also leads to inferior performance. This paper builds on the center-offset approach and addresses its drawback.

1.2 Method Overview

To address the drawback of the center-offset formulation, we build the intuitive idea of “**Keypointing, fast and slow**”, by analogy to the modes of thought suggested by Daniel Kahneman in “*Thinking, fast and slow*” [14]: (i) *Fast Keypointer*: We treat the vanilla center-offset based estimation [35] as the *Fast Initializer* of pose estimation. (ii) *Slow Keypointer*: The lack of localization accuracy in the Fast Initializer entails a *Slow Solver* that learns to refine the “fast” keypoints. By slow, it is only relatively speaking. The Slow Keypointer is actually fast with near real-time speed.

To realize the Slow Keypointer, as illustrated in Fig. 1 and Fig. 3, this paper presents a method of learning **LOcal-GIObal Contextual Adaptation** for fully end-to-end and fast bottom-up human Pose estimation, dubbed as **LOGO-CAP**. To quantitatively motivate the proposed method, we first present a surprisingly strong observation for a vanilla center-offset regression method (Table 1) in the fully-annotated subset of the COCO val-2017 dataset. Specifically, the vanilla regression method utilizes the HRNet-W32 [27] as the feature backbone to directly predict keypoints center heatmap and the offset vectors. This vanilla center-offset model obtains 60.1 average precision (AP), which is not great, but reasonably good. It clearly shows that the pose keypoints center and the offset vectors can be learned reasonably well. Instead of directly utilizing the learned offset vectors for human pose estimation, we treat them as human pose keypoint initialization and do a local window search to compute the empirical upper-bound of performance. More detailed, based on the predicted human poses, by introducing a local window (e.g., 11×11) centered at each detected key point and by computing the single keypoint similarity with the ground-truth keypoint, an empirical upper-bound of 88.9 AP is obtained, which is significantly higher than the state of the art and shows the potential of improving the vanilla center-offset regression paradigm.

Table 1: The performance of a vanilla center-offset regression approach, its empirical upper bound, and the performance of our proposed LOGO-CAP using HRNet-W32 [27] as the feature backbone. See text for detail.

	Baseline	Emp. Bound	LOGO-CAP
AP	60.1	88.9	70.0
AP ⁵⁰	85.2	93.1	88.2
AP ⁷⁵	66.7	90.6	76.4
AP ^M	53.7	87.7	64.4
AP ^L	71.5	90.2	78.4

Motivated by the above observation, a straightforward way is just to learn a local heatmap (e.g., 11×11) for each human pose keypoint based on the learned center and offset vectors, and then to

83 compute the refined keypoints by taking $\arg \max$ within the local heatmap. Although appealing,
 84 this does not work as observed during our development of the LOGO-CAP. The underlying reason
 85 is easy to understand: if this can work, the original offset vector regression should work at the
 86 first place since no additional information is introduced through learning the local heatmap. *We*
 87 *hypothesize* that on the one hand, on top of the local heatmap, the structural relationship between
 88 different keypoints of a human pose needs to be taken into account, and on the other hand, the
 89 intrinsic uncertainty of the local information in a local heatmap needs to be resolved. The former
 90 is the key challenge of structured output prediction problems. Many message passing algorithms
 91 have been developed in the literature. The latter can not be addressed by simply increasing the local
 92 window size. It entails learning stronger local-global information interaction and adaptation.,

93 Along with the two hypotheses, the proposed LOGO-CAP lifts the initial keypoints via the center-offset pre-
 94 diction to keypoint expansion maps (KEMs) to counter their lack of localization accuracy in two modules (Sec-
 95 tion 3.2). The KEMs extend the star-structured representation of the center-offset formulation to the pictorial
 96 structure representation [10, 8]. The first module computes local KEMs and learns to account for the struc-
 97 tured output prediction nature of the human pose estimation problem, leading to the keypoint attraction maps
 98 (KAMs). The second computes global KEMs and learns to refine the global KEMs by leveraging the KAMs.
 100

105 Our LOGO-CAP is a fully end-to-end bottom-up human pose estimation method with near real-time infer-
 106 ence speed. It obtains 70.0 AP in the fully-annotated subset of the COCO val-2017 dataset, which is an absolute
 107 increase of 9.9 AP compared to the vanilla center-offset
 108 method, making a significant step forward. Fig. 1 shows a
 109 pose estimation example. Fig. 2 shows the advantage of the proposed LOGO-CAP in terms of over-
 110 all speed-accuracy comparisons between our LOGO-CAP and prior arts. Meanwhile, we should
 111 notice that there is also a significant gap compared to the empirical upper bound (Table 1), which
 112 encourages more work to be investigated.
 114

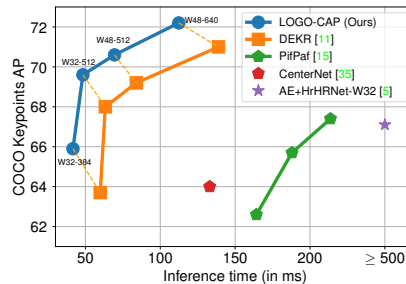


Figure 2: Speed-accuracy comparisons between our LOGO-CAP and prior arts on the COCO val-2017 dataset. Wx - Y (e.g. W32-384) means that a model uses the backbone HRNet- Wx (HRNet-W32) and is tested with the image resolution Y in the short side.

115 2 Related Works and Our Contributions

116 There is a vast body of literature for human pose estimation. Many elegant representation schema
 117 have been developed for modeling articulated human pose in the traditional approaches such as the
 118 well-known pictorial structure model [10, 8] and its many variants [24, 1, 23, 33, 25]. Most of them
 119 focused on single person pose estimation. They perform inference over a combination of local ob-
 120 servations on body parts (i.e., the data term) and the spatial dependencies between them (i.e., the
 121 spring or clique term). The spatial dependencies are captured either using directed and acyclic struc-
 122 tures that facilitate the global optimization by dynamic programming [2, 9], or using structures with
 123 loop introduced (for high-order part relationship modeling) which resort to approximate inference
 124 by loopy belief propagation [19]. The bottleneck of the traditional methods lies in the data term
 125 which is often based on hand-crafted features. With the resurgence of DNNs and the end-to-end
 126 learning, the data term has been largely improved. We briefly review the recent deep learning based
 127 approaches for bottom-up human pose estimation.

128 **Limb-based Grouping Approaches** have been extensively developed due to the naturalness of
 129 modeling limbs based on keypoints. Given a predefined limb configuration (e.g., the COCO person
 130 skeleton template consisting of 19 limbs based on 17 keypoints), the grouping can be addressed by
 131 Part affinity field (PAF) [4, 3], Associative Embedding (AE) [20], mid-range offset fields in Person-
 132 Lab [22] and the fields of Part Intensity and Association [15]. Typically, sophisticated designs are
 133 entailed to achieve good performance. For example, a bipartite graph matching is used in Open-
 134 Pose [3]. In addition to be computationally expensive, another drawback of these methods is not
 135 fully end-to-end trainable. More recently, the differentiability issue was studied by the Hierarchical
 136 Graph Clustering (HGG) method [13], which utilizes graph convolution networks to repeatedly de-
 137 lineate pose parameters of multiple persons from a keypoint graph. HGG improves the performance
 138 compared to its baseline, the Associative Embedding method [20] at the expense of significantly

139 increased computational cost. In contrast to those approaches, our proposed LOGO-CAP is fully
 140 end-to-end trainable and achieves near real-time inference speed.

141 **Direct Regression based Approaches** have attracted much attention due to their conceptually simple
 142 formulation [6, 28, 26, 11, 30]. These center-offset based formulations are inspired by the recent
 143 remarkable success of direct bounding box regression in object detection such as the FCOS
 144 method [29] and CenterNets [35, 6]. As aforementioned, one main challenge is the difficulty of accurately
 145 regress the offset vectors, especially for the long-range keypoints with respect to the center.
 146 Sophisticated post-processing schemas are often entailed to improve the performance. For example,
 147 a method of matching the directly regressed poses to the nearest keypoints that are extracted from
 148 the global keypoint heatmaps is used in [35]. Although being simple, the performance of this line of
 149 work is usually inferior to the limb-based approaches. The mixture regression network [30] alleviated
 150 the issue of regression quality to some extent, but still remained an indispensable performance
 151 gap comparing with the grouping-based approaches. Most recently, Geng *et al.* presented the first
 152 competitive direct method, DEKR [11] with a novel pose-specific neural architecture for disentangled
 153 keypoint regression. To improve the performance, the DEKR method utilizes a lightweight
 154 rescore network to recalibrate the pose scores that are computed based on the keypoint heatmaps.
 155 Despite good performance, the DEKR method entails the additional rescore stage in both training
 156 and testing, and thus is not fully end-to-end. The proposed LOGO-CAP retains the simplicity of the
 157 vanilla center-offset formulation and enjoys fully end-to-end training and fast inference speed.

158 **Our Contributions.** The proposed LOGO-CAP makes three main contributions to the field of
 159 bottom-up human pose estimation: (i) It addresses the drawback of the vanilla center-offset formulation
 160 while retaining its efficiency. It proposes the key idea of lifting a keypoint to a keypoint expansion
 161 map to counter the lack of localization accuracy. To our knowledge, it is the first fully end-to-end
 162 trainable method that achieves state-of-the-art performance. (ii) It presents a novel local-global
 163 contextual adaptation formulation that accounts for the nature of structured output prediction in
 164 human pose estimation and harnesses local-global structural information integration. (iii) It obtains
 165 state-of-the-art performance in the COCO val-2017 and test-2017 datasets. It also shows state-of-the-art
 166 transferability performance in the OCHuman dataset.

167 3 Approach

168 3.1 Problem Formulation

169 We follow the COCO protocol of defining the human pose. It consists of 17 human pose keypoints:
 170 8 pairs of symmetric keypoints (hips, ankles, knees, shoulders, elbows, wrists, ears and eyes) and
 171 the nose keypoint. Let $P = \{1, \dots, 17\}$ be the set of keypoint indexes using a predefined order.
 172 Let Λ be an image lattice of the spatial size $H \times W$ (e.g., 512×512), and I be an image defined
 173 on Λ . Let P_I^n be the set of keypoint indexes for a human pose instance n in an image I and we
 174 have $P_I^n \subseteq P$. For example, in COCO, we typically have $1 \leq n \leq 30$, and different human pose
 175 instances have different number of visible keypoints due to occlusion and/or truncation. Denote by
 176 $L_I^n = \{(x_i, y_i); i \in P_I^n\}$ the keypoint locations of a human pose instance n in an image I , where
 177 $(x_i, y_i) \in \Lambda$. In the center-offset formulation, we introduce the keypoints center (i.e., the anchor),
 178 (x_c, y_c) based on a given L_I^n and we have,

$$x_c = 1/|L_I^n| \cdot \sum_{i \in P_I^n} x_i, \quad y_c = 1/|L_I^n| \cdot \sum_{i \in P_I^n} y_i. \quad (1)$$

179 With the anchor, a keypoint (x_i, y_i) is equivalently defined by its offset/displacement, denoted by
 180 $(\Delta x_i, \Delta y_i)$ with $\Delta x_i = x_i - x_c$ and $\Delta y_i = y_i - y_c$. So, L_I^n can also be equivalently expressed as
 181 $L_I^n = \{(x_c, y_c), (\Delta x_i, \Delta y_i); i \in P_I^n\}$.

182 The objective of human pose estimation is to recover $L_I^n = \{(x_i, y_i); i \in P_I^n\}$ for all human pose
 183 instances in an image. Denote by $\hat{L}_I^n = \{(\hat{x}_i, \hat{y}_i); i \in P_I^n\}$ the estimated human pose. Following
 184 the COCO protocol, the object keypoint similarity (OKS) is used to evaluate the accuracy,

$$\ell_{OKS}(\hat{L}_I^n, L_I^n) = 1/|P_I^n| \cdot \sum_{i \in P_I^n} \exp(-d_i^2/2s^2\kappa_i^2), \quad (2)$$

185 where d_i is the Euclidian distance between the ground-truth keypoint (x_i, y_i) and the predicted one
 186 (\hat{x}_i, \hat{y}_i) . s is the square root of the human segment area, and κ per-keypoint constant that controls
 187 fall-off in evaluation. We have $\ell_{OKS}(\hat{L}_I^n, L_I^n) \in [0, 1]$. The OKS metric is to evaluate the distance
 188 between predicted keypoints and ground-truth keypoints normalized by the scale of the person with
 189 the importance of keypoints equalized. In benchmarking different methods, the average precision

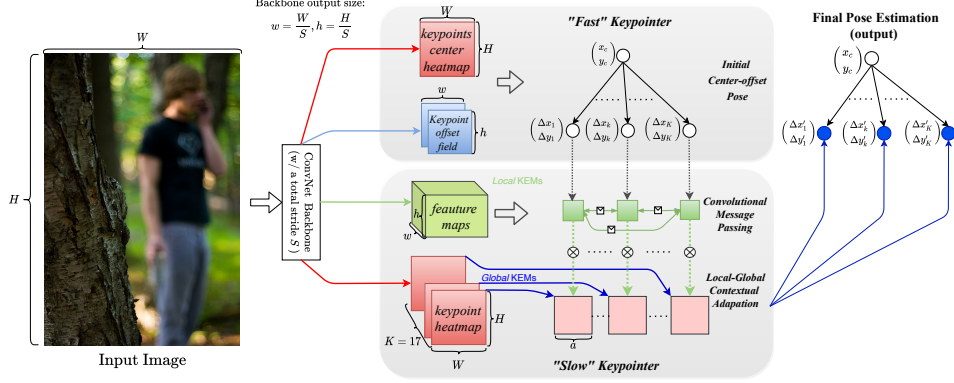


Figure 3: Illustration of the network and algorithmic flow of the proposed LOGO-CAP for bottom-up human pose estimation. See text for detail.

190 (AP) at OKS= 0.50 : 0.05 : 0.95 is used as the primary metric, together with AP^{50} at OKS= 0.50,
 191 AP^{75} at OKS= 0.75, and AP across medium and large scales, AP^M and AP^L respectively.

192 3.2 The Proposed LOGO-CAP

193 We first present the network and the inference of LOGO-CAP, and then give details of the train-
 194 ing. We keep different modules of the proposed LOGO-CAP simple, which in turn highlights the
 195 effectiveness of the proposed representation and algorithmic flow.

196 3.2.1 The Network and the Inference

197 As illustrated in Fig. 3, the proposed LOGO-CAP consists of four components as follows.

198 **i) A convolution neural network feature backbone.** Given an input image I , the output of the fea-
 199 ture backbone is a C -dim feature map, denoted by $F \in R^{C \times h \times w}$, where C is the feature dimension
 200 of the last convolutional layer in the feature backbone, and the spatial size $h \times w$ depends on the
 201 total stride in the feature backbone. We use off-the-shelf HRNets [27] in our experiments.

202 **ii) A parallel keypoint-offset regression module.** Given the feature map F , the output of keypoint
 203 regression is an 18-dim feature map (i.e., heatmaps) for the 17 keypoints and the keypoints center
 204 respectively. Denote by $\mathcal{H} \in R^{18 \times h \times w}$ the heatmaps, and by $\mathcal{H}^\uparrow \in R^{18 \times H \times W}$ the up-sampled
 205 heatmaps (using bi-linear interpolation in our experiments). The output of offset regression is a
 206 34-dim feature map (i.e., the offset vector fields) for the 17 keypoints. Denote by $\mathcal{O} \in R^{34 \times h \times w}$
 207 the offset fields. We adopt a minimally-simple design in realizing the regression modules using a
 208 channel-wise multi-layer perceptron (MLP). In implementation, we first apply dimension reduction
 209 to the feature map F using a 1×1 convolution followed by a Batch Normalization (BN) and a Rec-
 210 tified Linear Unit (ReLU). Then, the output is computed by a 1×1 convolution. More specifically,
 211 we have the two parallel branches as follows,

$$F_{C \times h \times w} \xrightarrow[C \times 1 \times 1 \times C_1]{Conv+BN+ReLU} F_{C_1 \times h \times w}^{\mathcal{H}} \xrightarrow[C_1 \times 1 \times 1 \times 18]{Conv} \mathcal{H}_{18 \times h \times w} \xrightarrow[\text{bi-linear}]{UpSampling} \mathcal{H}_{18 \times H \times W}^\uparrow, \quad (3)$$

$$F_{C \times h \times w} \xrightarrow[C \times 1 \times 1 \times C_2]{Conv+BN+ReLU} F_{C_2 \times h \times w}^{\mathcal{O}} \xrightarrow[C_2 \times 1 \times 1 \times 34]{Conv} \mathcal{O}_{34 \times h \times w}, \quad (4)$$

212 where C_1 and C_2 are predefined (e.g., $C_1 = 32$ and $C_2 = 256$ are typically used).

213 *Initial pose estimation via the center-offset approach.* Based on the computed keypoints center
 214 heatmap $\mathcal{H}_{(18)}^\uparrow$ and offset fields \mathcal{O} , a predefined maximum number of pose candidates is computed
 215 as done in the vanilla center-offset approach. A non-maximum suppression (NMS) with a 3×3
 216 window is applied in $\mathcal{H}_{(18)}^\uparrow$ and then the top- N keypoints centers are selected (e.g., $N = 30$ in our
 217 experiments). The N pose instances are computed by retrieving their offset vectors in \mathcal{O} based on
 218 the selected N keypoints centers. The N pose instances are further pruned by thresholding their
 219 confidence scores in $\mathcal{H}_{(18)}^\uparrow$ with a predefined threshold (e.g., 0.01 used in our experiments). Without
 220 confusion in the context, we still use N to denote the number of poses instances by this initial pose
 221 estimation step. We obtain the set of estimated keypoints centers, denoted by $\mathcal{C}_{N \times 3}$ each row of
 222 which represents the position coordinates and the confidence score.

223 **Lifting a keypoint to a keypoint expansion map (KEM) by imposing a mesh.** For each of the N
 224 pose instances, each of the 17 keypoints are placed in a local geometric mesh (e.g., 11×11) with the
 225 estimated location as the mesh center, capturing the uncertainty of the center-offset pose estimation
 226 as aforementioned in the introduction. This mesh can thus be interpreted as keypoint expansion
 227 map (KEM), accounting for competency-aware representations. The entire mesh is denoted by
 228 $\mathcal{M}_{N \times 17 \times 11 \times 11 \times 2}$, which is used in computing the empirical upper bound in Table 1. We have,

$$\{\mathcal{H}_{(18)}^\uparrow, \mathcal{O}_{34 \times h \times w}\} \xrightarrow[\text{center-offset}]{\text{initial pose estimation}} \{\mathcal{C}_{N \times 3}, \mathcal{M}_{N \times 17 \times 11 \times 11 \times 2}\} \quad (5)$$

229 **iii) A convolution message passing module.** We first encode the geometric mesh $\mathcal{M}_{N \times 17 \times 11 \times 11 \times 2}$
 230 in a latent space with the dimensionality C_3 (e.g., 64 in our experiments), computed based on the
 231 feature backbone output. Then, a keypoint is represented by a $C_3 \times 11 \times 11$ local feature map. A
 232 pose instance is represented by concatenating all the 17 keypoints. We have,

$$F_{C \times h \times w} \xrightarrow[\frac{C \times 1 \times 1 \times C_3}{\text{Conv+BN+ReLU}}]{\text{Conv+BN+ReLU}} F_{C_3 \times h \times w}^{\mathcal{M}} \xrightarrow[\text{bi-linear}]{\mathcal{M}_{N \times 17 \times 11 \times 11 \times 2}} \mathcal{K}_{N \times (17 \times C_3) \times 11 \times 11}, \quad (6)$$

233 where the bi-linear interpolation is used due to the sub-pixel based locations in the mesh and for
 234 better feature alignment.

235 To facilitate the structural information flow between different latent codes of the keypoints of a pose
 236 instance, we propose a simple convolutional message passing (CMP) module with three layers of
 237 Conv+BN+ReLU operations,

$$\mathcal{K}_{N \times (17 \times C_3) \times 11 \times 11} \Rightarrow \left[\frac{\text{Conv+BN+ReLU}}{C_{in} \times 3 \times 3 \times C_{out}} \right]_{\times 3} \Rightarrow \cdot \xrightarrow[\frac{C_6 \times 1 \times 1 \times 17}{\text{Conv}}]{\text{Conv}} \mathcal{K}_{N \times 17 \times 11 \times 11}, \quad (7)$$

238 where $C_{in} \in \{(17 \times C_3), C_4, C_5\}$ and $C_{out} \in \{C_4, C_5, C_6\}$ (e.g., $C_4 = 512, C_5 = 256, C_6 = 128$
 239 in our experiments). The resulting $\mathcal{K}_{N \times 17 \times 11 \times 11}$ can be interpreted as keypoint attraction maps
 240 (KAMs) which are “re-focused” based on the KEMs by the CMP. To account for the specificity of
 241 different pose instances in the CMP, we adopt the Attention Normalization [17] to replace the BN in
 242 the second Conv+BN+ReLU layer, which further improves the performance in our experiments.

243 Through the CMP, we obtain the dynamic (a.k.a., data-driven) kernels for the 17 keypoints in a pose
 244 instance-sensitive way, which are used to refine the global heatmaps \mathcal{H}^\uparrow for the 17 keypoints.

245 **iv) A local-global contextual adaptation module.** We first compute another geometric mesh with
 246 enlarged mesh window $a \times a$ (e.g., $a = 97$) for each keypoint of the N pose instances, and the entire
 247 mesh is denoted by $\mathcal{M}_{N \times 17 \times a \times a \times 2}$, as done in Eqn. 5. The mesh can be interpreted as the global
 248 KEM. It is then instantiated with appearance features extracted from the global heatmaps $\mathcal{H}_{(1:17)}^\uparrow$,
 249 similar to Eqn. 6, and we have,

$$\mathcal{H}_{(1:17)}^\uparrow \xrightarrow[\text{bi-linear}]{\mathcal{M}_{N \times 17 \times a \times a \times 2}} \mathbb{H}_{N \times 17 \times a \times a} \xrightarrow[\text{reweighing}]{\mathcal{G}_{a \times a}(0, \sigma)} \tilde{\mathbb{H}}_{N \times 17 \times a \times a}. \quad (8)$$

250 where to encode the Gaussian prior of keypoint heatmaps, the resulting pose-guided heatmaps \mathbb{H} is
 251 reweighed by a Gaussian kernel $\mathcal{G}_{a \times a}(0, \sigma = \frac{a-1}{2 \times 3})$ (e.g., $\sigma = 16$ when $a = 97$) in an element-wise
 252 way. By doing so, it means that the enlarged mesh follows the 3σ principle.

253 Then, we apply the learned keypoint 11×11 kernels $K_{n,i}$ ’s (Eqn. 7) to convolve the reweighed
 254 $a \times a$ heatmap $\tilde{\mathbb{H}}_{n,i}$ (Eqn. 8) in a pose instance-sensitive and keypoint-specific way, leading to
 255 **Local-Global Contextual Adaptation**,

$$\tilde{\mathbb{H}}_{N \times 17 \times a \times a} \xrightarrow[\text{LOGO-CA}]{K_{N \times 17 \times 11 \times 11}} \tilde{\tilde{\mathbb{H}}}_{N \times 17 \times a \times a}, \quad (9)$$

256 which represents the refined heatmaps for the 17 human pose keypoints.

257 **The Pose Estimation Output.** With the local-global contextually adapted heatmaps $\tilde{\tilde{\mathbb{H}}}_{N \times 17 \times a \times a}$,
 258 we maintain the top-2 locations for each keypoint within the $a \times a$ heatmap, and then utilize a convex
 259 average of the top-2 locations as the final predicted offset vectors (i.e. $(\Delta x'_i, \Delta y'_i)$ ’s in Fig. 3), and
 260 of their confidence scores as the prediction score, with a predefined weight λ for the top-1 location
 261 (0.75 in our experiments). Together with the predicted keypoints centers $\mathcal{C}_{N \times 3}$ (Eqn. 5), the final
 262 prediction score for each keypoint is the product between the convex average confidence score and
 263 the center confidence score. We keep the keypoints whose final scores are greater than 0. We have,

$$\{\mathcal{C}_{N \times 3}, \tilde{\tilde{\mathbb{H}}}_{N \times 17 \times a \times a}\} \xrightarrow[\text{Score thresholding}]{\text{Output}} \{\hat{L}_I^n; n = 1, \dots, N'\}, \quad (10)$$

264 where N' is the number of the final predicted pose instances in an image I .

265 **3.2.2 Loss Functions in Training**

266 In the fully end-to-end training, we need to define loss functions for the global heatmap \mathcal{H} (Eqn. 3),
 267 the refined local heatmap $\tilde{\mathbb{H}}$ (Eqn. 9), the offset field \mathcal{O} (Eqn. 4), and the keypoint kernels (Eqn. 7).

268 **The Heatmap Loss.** The widely adopted mean squared error (MSE) loss is used. Denoted by
 269 $\mathcal{H}_{18 \times h \times w}^{GT}$ the ground truth heatmaps in which each keypoint (including the center) is modeled by a
 270 2-D Gaussian with dataset-provided mean and variance. Let $\mathbf{p} = (i, \mathbf{x})$ be the index of the domain
 271 D of dimensions $18 \times h \times w$. For the predicted heatmaps $\mathcal{H}_{18 \times h \times w}$, the MSE loss is defined by,

$$\mathcal{L}_{\mathcal{H}} = 1/|D| \cdot \sum_{\mathbf{p} \in D} \|w(\mathbf{x})(\mathcal{H}(\mathbf{p}) - \hat{\mathcal{H}}(\mathbf{p}))\|_2^2, \quad (11)$$

272 where $w(\mathbf{x})$ represents the weight for the foreground and the background pixels. The foreground
 273 mask is provided by the dataset annotation. In our experiment, we set $w(\mathbf{x}) = 1$ for a foreground
 274 pixel and $w(\mathbf{x}) = 0.1$ for a background pixel.

275 In defining the loss function $\mathcal{L}_{\tilde{\mathbb{H}}}$ for the refined local heatmap $\tilde{\mathbb{H}}$ (Eqn. 9), the ground-truth heatmap
 276 $\tilde{\mathbb{H}}^{GT}$ is generated on-the-fly based on the mesh $\mathcal{M}_{N \times 17 \times a \times a}^L$ (Eqn. 8) and the ground-truth key-
 277 points using a Gaussian model with mean being the displacement between the current predicted
 278 keypoints and the ground-truth ones, and variance σ (i.e., the standard deviation of the reweighing
 279 Gaussian prior model in Eqn. 8).

280 **The Offset Field Loss.** The widely adopted SmoothL1 loss $[\]$ is used. Let $\mathcal{O}_{34 \times h \times w}^{GT}$ be the ground-
 281 truth offset field, and \mathcal{C}^{GT} be the non-empty set of ground-truth keypoints centers (Eqn. 1). For the
 282 predicted offset field $\mathcal{O}_{34 \times h \times w}$ (Eqn. 4), we have,

$$\mathcal{L}_{\mathcal{O}} = 1/|\mathcal{C}^{GT}| \cdot \sum_{\mathbf{p} \in \mathcal{C}^{GT}} \mathcal{A}(\mathbf{p}) \cdot \text{SmoothL1}(\mathcal{O}(\cdot, \mathbf{p}), \mathcal{O}^{GT}(\cdot, \mathbf{p}); \beta), \quad (12)$$

283 where $\mathcal{A}(\mathbf{p})$ is the area of the person centered at the pixel \mathbf{p} , and β the cutting-off threshold (e.g., $\frac{1}{9}$ in
 284 our experiments), and $\text{SmoothL1}(a, b; \beta) = 0.5 \times |a - b|^2 / \beta$ if $|a - b| \leq \beta$, otherwise $|a - b| - 0.5 \times \beta$.

285 **The OKS Loss for the Keypoint Kernels.** Consider a single predicted pose instance, learning the
 286 keypoint kernels, $K_{17 \times 11 \times 11}$ (Eqn. 7) is the key to facilitate the local-global contextual adaptation.
 287 To that end, the figure of merits of the KEF, $\mathcal{M}_{17 \times 11 \times 11 \times 2}$ (Eqn. 5) needs to directly reflect the task
 288 loss, i.e., the OKS loss (Eqn. 2). With respect to the N^{GT} ground-truth pose instances in an image,
 289 we can compute the similarity score per keypoint candidate in the KEF, and obtain the score tensor
 290 $S_{17 \times 11 \times 11 \times N^{GT}}$. The score tensor is further clamped with a threshold 0.5, i.e., $S_{17 \times 11 \times 11 \times N^{GT}} =$
 291 $\max(S_{17 \times 11 \times 11 \times N^{GT}}, 0.5)$. A mean reduction is applied to the first three dimensions of the clamped
 292 score tensor to compute the matching score for each of the N^{GT} pose instance. Then, the best
 293 ground-truth pose instance indexed by n^* is selected in terms of the matching score, and its matching
 294 score is denoted by s_{n^*} . Based on the selected ground-truth pose instance, we compute the per-
 295 keypoint similarity score for the predicted pose instance at hand, denoted by s_k ($k \in [1, 17]$). Then,
 296 the loss function fo the keypoint kernels are defined by,

$$\mathcal{L}_K = s_{n^*} \cdot \sum_{k, i, j} s_k \cdot |K_{k, i, j} - S_{k, i, j, n^*}|^2. \quad (13)$$

297 **The Total Loss** is then defined by $\mathcal{L} = \mathcal{L}_{\mathcal{H}} + \mathcal{L}_{\tilde{\mathbb{H}}} + \lambda \cdot (\mathcal{L}_{\mathcal{O}} + \mathcal{L}_K)$, where the trade-off parameter
 298 λ is used to balance the different loss items ($\lambda = 0.01$ in our experiments).

299 **4 Experiments**

300 In this section, we present detailed experimental results and analyses of the proposed LOGO-CAP.
 301 **Our PyTorch source code will be released for reproducibility.**

302 **Datasets.** We use two datasets in our experiments: **The COCO dataset [18]** is the most popu-
 303 lar testbed for human pose estimation. It consists of 65k, 5k and 20k images with human pose
 304 well-annotated in the training, validation and testing datasets respectively. In all experiments, the
 305 proposed LOGO-CAP is trained using the 65k training images. **The OCHuman dataset [34]** is one
 306 popular *testing-only* dataset for evaluating human pose estimation under the occlusion scenarios. It
 307 consists of a total number of 4713 images with 8110 detailed annotated human pose instances using
 308 the COCO keypoint configuration. All the annotated 8110 human pose instances have occlusions
 309 with the $\text{maxIOU} \geq 0.5$. Furthermore, 32% instances are more challenging with the $\text{maxIOU} \geq 0.75$.

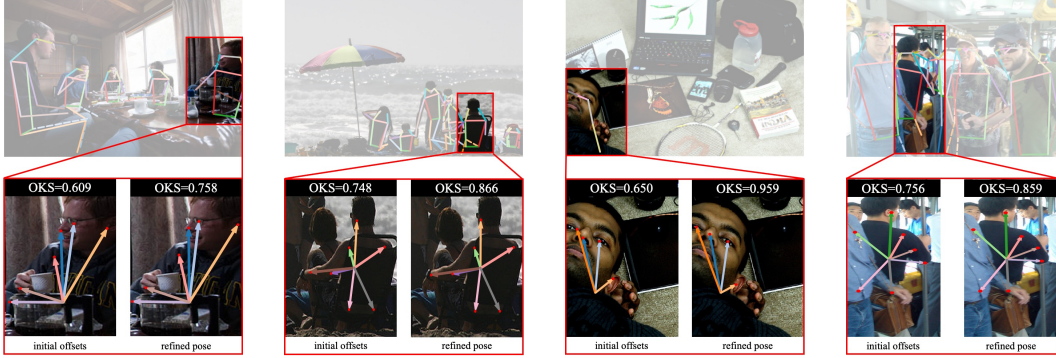


Figure 4: Examples of human pose estimation in the COCO val-2017 dataset by the proposed LOGO-CAP with the HRNet-W32 backbone. *Top*: The COCO skeleton template based visualization. *Bottom*: The close-up visualization and OKS comparisons between the initial center-offset estimation and the refined keypoints.

Table 2: Evaluation results on the COCO-val-2017 and COCO-testdev-2017 dataset. For HGG [13] and SimplePose [16], the multi-scale inference[†] is applied on the testdev-2017 dataset. For DEKR [11] that uses an rescoring network to get the final predictions, we report both the performance with and without rescoring (which is the fair baseline for our LOGO-CAP). The numbers of SPM [21] and HGG [13] are extracted from their papers.

	Method	Backbone	COCO-val-2017					COCO-testdev-2017					
			AP [%]	AP ⁵⁰ [%]	AP ⁷⁵ [%]	AP ^M [%]	AP ^L [%]	AP [%]	AP ⁵⁰ [%]	AP ⁷⁵ [%]	AP ^M [%]	AP ^L [%]	
Grouping	OpenPose [35]	VGG-19	61.0	84.9	67.5	56.3	69.3	61.8	84.9	67.5	57.1	68.2	
	PifPaf [15]	ResNet-152	67.4	86.9	73.8	63.1	74.1	66.7	87.8	73.6	62.4	72.9	
	PersonLab [22]	ResNet-152	66.5	86.2	71.9	62.3	73.2	66.5	88.0	72.6	62.4	72.3	
	AE [20, 5]	HrHRNet-W32	67.1	86.2	73.0	61.5	76.1	66.4	87.5	72.8	61.2	74.2	
		HrHRNet-W48	69.9	87.2	76.1	65.4	76.4	68.4	88.2	75.1	64.4	74.2	
		HGG [13]	Hourglass	60.4	83.0	66.2	—	—	67.6 [†]	85.1 [†]	73.7 [†]	62.7 [†]	74.6 [†]
		SimplePose [16]	IMHN	66.1	85.9	71.6	59.8	76.2	68.5 [†]	86.7 [†]	74.9 [†]	66.4 [†]	71.9 [†]
Direct	SPM [21]	Hourglass	—	—	—	—	—	66.9	88.5	72.9	62.6	0.731	
	CenterNet [35]	Hourglass	64.0	85.6	70.2	59.4	72.1	63.0	86.8	69.6	58.9	70.4	
	DEKR [11]	HRNet-W32	68.0	86.7	74.5	62.1	77.7	67.3	87.9	74.1	61.5	76.1	
	(w. Rescoring)	HRNet-W48	71.0	88.3	77.4	66.7	78.5	70.0	89.4	77.3	65.7	76.9	
	DEKR [11]	HRNet-W32	67.2	86.3	73.8	61.7	77.1	66.6	87.6	73.5	61.2	75.6	
	(w.o. Rescoring)	HRNet-W48	70.3	87.9	76.8	66.3	78.0	69.3	89.1	76.7	65.3	76.4	
	LOGO-CAP (Ours)	HRNet-W32	69.6	87.5	75.9	64.1	78.0	68.2	88.7	74.9	62.8	76.0	
	HRNet-W48	72.2	88.9	78.9	68.1	78.9	70.8	89.7	77.8	66.7	77.0		

310 4.1 Results on the COCO dataset

311 Fig. 4 shows some qualitative examples of human pose estimation by the proposed LOGO-CAP.
 312 More examples will be provided in the supplementary material.

313 The proposed LOGO-CAP is compared with prior arts including OpenPose [3], PifPaf [15], Person-
 314 Lab [22], AE [20] and DEKR [11]. As reported in Table 2, the proposed LOGO-CAP outperforms
 315 all of them on both both validation and test-dev datasets.

316 In comparisons to the best-performing grouping approach, AE [20] with a larger backbone
 317 HrHRNet-W48 [5], our LOGO-CAP obtains competitive performance with a smaller HRNet-32
 318 backbone, and improves the AP score with HRNet-W48 backbone on the validation and testdev
 319 datasets by 2.3 and 2.5 points, respectively. For the fully differentiable grouping approach
 320 HGG [13], our LOGO-CAP achieves better performance by a significantly large margin, more than
 321 9.2 points on the validation set under the single-scale testing. Although the performance of HGG
 322 is improved by the multi-scale testing on the test-dev set, the performance of our LOGO-CAP is still
 323 significantly better without using the multi-scale testing scheme.

324 In comparisons to the direct regression based approaches, our LOGO-CAP obtains the *best results*
 325 without incurring either the matching scheme used in CenterNet [35] or the additional rescoring
 326 network used in DEKR [11]. When we disable the rescoring network for DEKR [11] for fair com-
 327 parisons, our LOGO-CAP significantly improves the AP on the validation and testdev datasets by
 328 2.4 points and 1.6 points respectively when HRNet-W32 is used as backbone. The larger back-
 329 bone is beneficial for both DEKR and our method, which further improves the AP score of our
 330 LOGO-CAP to 72.2 and 70.8 on the validation and test-dev dataset respectively, outperforming
 331 DEKR by 1.9 and 1.5 respectively.

Table 3: Results on the OCHuman validation and testing datasets [34].

	Methods	Backbone	Val. AP [%]	Test AP [%]
Top-down	RMPE [7]	Hourglass	38.8	30.7
	SBL [32]	ResNet-50	37.8	30.4
	SBL [32]	ResNet-152	41.0	33.3
Bottom-up	AE [20]	Hourglass	32.1	29.5
	HGG [20]	Hourglass	35.6	34.8
	DEKR [11]	HRNet-W32	37.9	36.5
		HRNet-W48	38.8	38.2
	LOGO-CAP (Ours)	HRNet-W32	39.0	38.1
	HRNet-W48	41.2	40.4	

Table 4: The single image inference speed comparison for bottom-up human pose estimation approaches.

Method	AP [%]	Backbone	Time ↓ [ms]	FPS ↑
PifPaf [15]	67.4	ResNet-152	213	4.68
AE [20, 5]	67.1	HrHRNet-W32	560	1.78
CenterNet [35]	64.0	Hourglass	147	6.80
DEKR [11]	68.0	HRNet-W32	63	15.8
DEKR [11]	71.0	HRNet-W48	139	7.21
LOGO-CAP	69.6	HRNet-W32	48	20.7
LOGO-CAP	72.2	HRNet-W48	112	8.95

332 4.2 Results on the OCHuman dataset

333 Table 3 shows that our LOGO-CAP achieves the best AP performance on both the validation and
 334 testing datasets by significant margins of 2.4 and 2.2 points in comparing with the bottom-up
 335 approaches. For the top-down approaches, although they obtain strong AP scores on the validation
 336 split, there exists a large performance gap between the validation and testing sets. In comparisons
 337 to DEKR [11] (with the rescoreing network), our LOGO-CAP improves the performance from 37.9
 338 to 39.0 and from 36.5 to 38.1 on the validation and testing splits with the same backbone HRNet-
 339 W32, respectively. The similar improvement is observed when the HRNet-W48 backbone is used,
 340 outperforming both bottom-up and top-down approaches.

341 4.3 Inference Speed

342 In comparing the inference speed, we test all the models on a single TITAN RTX GPU for its
 343 popularity in practice. The average inference speed, FPS (frames per second), over the 5000 images
 344 in COCO-val-2017 is used for the comparison. For DEKR [11], we re-implement their inference
 345 code with better speed obtained for fair comparisons at the algorithm level. For methods that have
 346 post-processing schema on CPU, only one thread is used. As shown in Table 4, our LOGO-CAP
 347 runs significantly faster than PifPaf [15] and AE [20]. The CenterNet [35] runs slower than DEKR
 348 and our LOGO-CAP as it requires a post-processing scheme to match the predicted offsets to the
 349 keypoints obtained from heatmaps. Comparing with DEKR, the speed improvement of our LOGO-
 350 CAP is from the lightweight design of head modules since the same backbones are used. For the
 351 comparisons in Table 2, we run the models with different resolutions of testing images.

352 4.4 Potentials and Limitations of the Proposed LOGO-CAP

353 Consider the generic applicability of the center-offset formulation to many computer vision tasks as
 354 demonstrated in [35], we hypothesize that the proposed LOGO-CAP has a great potential to remedy
 355 the lack of sufficient accuracy using the vanilla center-offset method in those tasks. We also notice
 356 that the minimally-simple design in learning the “Slow Keypointer” can be relaxed for different
 357 accuracy-speed trade-offs in practice. For example, for the convolutional message passing module,
 358 an alternative method could be the Transformer model [31], which potentially will further improve
 359 the performance at the expense of inference speed. We leave these for future work.

360 5 Conclusion

361 This paper focuses on deep learning based formulation for bottom-up human pose estimation. It
 362 presents a method of learning Local-Global Contextual Adaptation for Pose estimation, dubbed as
 363 LOGO-CAP. The proposed LOGO-CAP is built on the conceptually simple center-offset paradigm
 364 and addresses its drawback of lacking the capability of accurately localizing human pose keypoints.
 365 The key idea of our LOGO-CAP is to lift the center-offset predicted keypoints to keypoint expansion
 366 maps (KEMs), which counters the inaccuracy and uncertainty of the initial keypoints. Two types of
 367 KEMs are introduced in two parallel modules on top of the feature backbone. Local KEMs are used
 368 to learn keypoint attraction maps (KAMs) via a convolutional message passing module that accounts
 369 for the structured output prediction nature of human pose estimation. Global KEMs are used to
 370 learn local-global contextual adaptation which convolves global KEMs using the KAMs as kernels.
 371 The refined global KEMs are used in computing the final human pose estimation. The proposed
 372 LOGO-CAP obtains state-of-the-art performance in COCO val-2017 and test-dev 2017 datasets for
 373 bottom-up human pose estimation. It also achieves state-of-the-art transferability performance in
 374 the OCHuman dataset with the COCO trained models.

References

- 375
- 376 [1] Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele. Pictorial structures revisited: People detection and
377 articulated pose estimation. In *2009 IEEE conference on computer vision and pattern recognition*, pages
378 1014–1021. IEEE, 2009. 3
- 379 [2] Richard Bellman. Dynamic programming. *Science*, 153(3731):34–37, 1966. 3
- 380 [3] Zhe Cao, Gines Hidalgo Martinez, Tomas Simon, Shih-En Wei, and Yaser A. Sheikh. Openpose: Realtime
381 multi-person 2d pose estimation using part affinity fields. *IEEE Trans. on Pattern Analysis and Machine
382 Intelligence (PAMI)*, 2019. 1, 3, 8
- 383 [4] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using
384 part affinity fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
385 1302–1310, 2017. 3
- 386 [5] Bowen Cheng, Bin Xiao, Jingdong Wang, Honghui Shi, Thomas S. Huang, and Lei Zhang. Higherhrnet:
387 Scale-aware representation learning for bottom-up human pose estimation. In *IEEE/CVF Conference on
388 Computer Vision and Pattern Recognition (CVPR)*, pages 5385–5394. IEEE, 2020. 3, 8, 9
- 389 [6] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint
390 triplets for object detection. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages
391 6568–6577, 2019. 2, 4
- 392 [7] Haoshu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. RMPE: regional multi-person pose estimation. In
393 *IEEE International Conference on Computer Vision (ICCV)*, pages 2353–2362, 2017. 9
- 394 [8] Pedro F Felzenszwalb and Daniel P Huttenlocher. Pictorial structures for object recognition. *International
395 journal of computer vision*, 61(1):55–79, 2005. 3
- 396 [9] Pedro F Felzenszwalb and Ramin Zabih. Dynamic programming and graph algorithms in computer vision.
397 *IEEE transactions on pattern analysis and machine intelligence*, 33(4):721–740, 2010. 3
- 398 [10] Martin A Fischler and Robert A Elschlager. The representation and matching of pictorial structures. *IEEE
399 Transactions on computers*, 100(1):67–92, 1973. 3
- 400 [11] Zigang Geng, Ke Sun, Bin Xiao, Zhaoxiang Zhang, and Jingdong Wang. Bottom-up human pose es-
401 timation via disentangled keypoint regression. In *IEEE Conference on Computer Vision and Pattern
402 Recognition (CVPR)*, 2021. 1, 2, 3, 4, 8, 9
- 403 [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. In *IEEE International
404 Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017. 1
- 405 [13] Sheng Jin, Wentao Liu, Enze Xie, Wenhai Wang, Chen Qian, Wanli Ouyang, and Ping Luo. Differentiable
406 hierarchical graph grouping for multi-person pose estimation. In *European Conference on Computer
407 Vision (ECCV)*, volume 12352, pages 718–734, 2020. 3, 8
- 408 [14] Daniel Kahneman. *Thinking, fast and slow*. Macmillan, 2011. 2
- 409 [15] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. Pifpaf: Composite fields for human pose estimation.
410 In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11977–11986, 2019. 1,
411 3, 8, 9
- 412 [16] Jia Li, Wen Su, and Zengfu Wang. Simple pose: Rethinking and improving a bottom-up approach for
413 multi-person pose estimation. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 11354–11361,
414 2020. 8
- 415 [17] Xilai Li, Wei Sun, and Tianfu Wu. Attentive normalization. In Andrea Vedaldi, Horst Bischof, Thomas
416 Brox, and Jan-Michael Frahm, editors, *European Conference on Computer Vision (ECCV)*, volume 12362,
417 pages 70–87, 2020. 6
- 418 [18] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
419 and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In David J. Fleet, Tomás Pajdla,
420 Bernt Schiele, and Tinne Tuytelaars, editors, *European Conference on Computer Vision (ECCV)*, volume
421 8693, pages 740–755, 2014. 7, 12
- 422 [19] Kevin Murphy, Yair Weiss, and Michael I Jordan. Loopy belief propagation for approximate inference:
423 An empirical study. *arXiv preprint arXiv:1301.6725*, 2013. 3

- 424 [20] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint
425 detection and grouping. In *Advances in Neural Information Processing Systems 30 (NeurIPS)*, pages
426 2277–2287, 2017. 3, 8, 9
- 427 [21] Xuecheng Nie, Jiashi Feng, Jianfeng Zhang, and Shuicheng Yan. Single-stage multi-person pose ma-
428 chines. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6950–6959, 2019.
429 8
- 430 [22] George Papandreou, Tyler Zhu, Liang-Chieh Chen, Spyros Gidaris, Jonathan Tompson, and Kevin Mur-
431 phy. Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geomet-
432 ric embedding model. In *European Conference on Computer Vision (ECCV)*, pages 282–299, 2018. 1, 3,
433 8
- 434 [23] Leonid Pishchulin, Mykhaylo Andriluka, Peter Gehler, and Bernt Schiele. Poselet conditioned pictorial
435 structures. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
436 588–595, 2013. 3
- 437 [24] Deva Ramanan, David A Forsyth, and Andrew Zisserman. Strike a pose: Tracking people by finding
438 stylized poses. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*
439 (*CVPR'05*), volume 1, pages 271–278. IEEE, 2005. 3
- 440 [25] Brandon Rothrock, Seyoung Park, and Song-Chun Zhu. Integrating grammar and segmentation for human
441 pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
442 pages 3214–3221, 2013. 3
- 443 [26] Ke Sun, Zigang Geng, Depu Meng, Bin Xiao, Dong Liu, Zhaoxiang Zhang, and Jingdong Wang.
444 Bottom-up human pose estimation by ranking heatmap-guided adaptive keypoint estimates. *CoRR*,
445 abs/2006.15480, 2020. 2, 4
- 446 [27] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for
447 human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
448 5693–5703, 2019. 2, 5
- 449 [28] Zhi Tian, Hao Chen, and Chunhua Shen. Directpose: Direct end-to-end multi-person pose estimation.
450 *CoRR*, abs/1911.07451, 2019. 2, 4
- 451 [29] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: fully convolutional one-stage object detection.
452 In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9626–9635, 2019. 4
- 453 [30] Ali Varamesh and Tinne Tuytelaars. Mixture dense regression for object detection and human pose esti-
454 mation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13083–
455 13092, 2020. 4
- 456 [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz
457 Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 9
- 458 [32] Bin Xiao, Haiping Wu, and Yichen Wei. Simple Baselines for Human Pose Estimation and Tracking.
459 *Computer Vision and Pattern Recognition*, 2018. 9
- 460 [33] Yi Yang and Deva Ramanan. Articulated human detection with flexible mixtures of parts. *IEEE transac-*
461 *tions on pattern analysis and machine intelligence*, 35(12):2878–2890, 2012. 3
- 462 [34] Song-Hai Zhang, Ruilong Li, Xin Dong, Paul L. Rosin, Zixi Cai, Xi Han, Dingcheng Yang, Haozhi
463 Huang, and Shi-Min Hu. Pose2seg: Detection free human instance segmentation. In *IEEE Conference*
464 *on Computer Vision and Pattern Recognition (CVPR)*, pages 889–898, 2019. 7, 9, 12
- 465 [35] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *CoRR*, abs/1904.07850, 2019.
466 2, 3, 4, 8, 9

467 **Checklist**

- 468 1. For all authors...
- 469 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
470 contributions and scope? [Yes]
- 471 (b) Did you describe the limitations of your work? [Yes] See Section 4.4.
- 472 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 473 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
474 them? [Yes]
- 475 2. If you are including theoretical results...
- 476 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 477 (b) Did you include complete proofs of all theoretical results? [N/A]
- 478 3. If you ran experiments...
- 479 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
480 mental results (either in the supplemental material or as a URL)? [Yes] See Section 4.
- 481 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
482 were chosen)? [Yes] See Section 4.
- 483 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
484 ments multiple times)? [No]
- 485 (d) Did you include the total amount of compute and the type of resources used (e.g., type
486 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4 and the supplemen-
487 tary material.
- 488 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 489 (a) If your work uses existing assets, did you cite the creators? [Yes] We cited the COCO
490 dataset [18] and the OCHuman dataset [34].
- 491 (b) Did you mention the license of the assets? [Yes] We mention the licenses in our source
492 code.
- 493 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
494
- 495 (d) Did you discuss whether and how consent was obtained from people whose data
496 you're using/curating? [Yes] We briefly discussed it in Section 4.
- 497 (e) Did you discuss whether the data you are using/curating contains personally identifi-
498 able information or offensive content? [No]
- 499 5. If you used crowdsourcing or conducted research with human subjects...
- 500 (a) Did you include the full text of instructions given to participants and screenshots, if
501 applicable? [N/A]
- 502 (b) Did you describe any potential participant risks, with links to Institutional Review
503 Board (IRB) approvals, if applicable? [N/A]
- 504 (c) Did you include the estimated hourly wage paid to participants and the total amount
505 spent on participant compensation? [N/A]