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

Supplementary Material

Anonymous Author(s)

Affiliation

Address

email

1 Experimental Settings

2 Our PyTorch source code will be released. We briefly present the details of training and testing as
3 follows.

4 **Training.** We train two LOGO-CAP networks with the ImageNet pretrained HRNet-W32 and
5 HRNet-W48 [9] as the feature backbone respectively on the COCO-train-2017 dataset [6]. Common
6 training specifications are used for simplicity in experiments. The Adam optimizer [4] is used with
7 default coefficients $\beta_1 = 0.9$ and $\beta_2 = 0.999$. For both the backbones, the total number of epochs
8 is set to 140 and the batch size is set to 12 images per GPU card. The same learning rate schedule is
9 used for both models. The learning rate is initially set to 0.001 and then decayed to 10^{-4} and 10^{-5}
10 at the 90-th and 120-th epoch respectively. We use 4 and 8 V100 GPUs to accelerate the training
11 for the two LOGO-CAP models with HRNet-W32 and HRNet-W48 respectively. The resolution
12 of training images is 512×512 and 640×640 for the two models respectively. Following the
13 widely adopted experimental settings in [10], the data augmentations in training include (1) random
14 rotation with the rotation degree from -30° to 30° , (2) random scaling with the factor in the range
15 of $[0.75, 1.5]$, (3) random translation in the range $[-40\text{pix}, 40\text{pix}]$ along both x and y directions, and
16 (4) random horizontal flipping with the probability of 0.5.

17 Similarly, for different hyperparameters such as the trade-off parameter λ in the total loss (Section
18 3.2.2 in the submission), we did not run computationally expensive hyperparameter optimization for
19 simplicity.

20 **Testing.** We focus on the single-scale testing protocol in the COCO keypoint benchmark for the
21 sake of efficient human pose estimation. In the testing phase, the short side of input images is resized
22 to a specific length (*e.g.*, 384, 512, or 640 pixels) and keep unchanged the aspect ratio between the
23 height and the width. As commonly adopted in many bottom-up pose estimation approaches (*e.g.*,
24 AE [7], HrHRNet [1], DEKR [2]), the flip testing is used as our default setting for the fair com-
25 parison. In the implementation, we feed the stacked tensor with an input image and a horizontally-
26 flipped one together to get the global heatmaps and the offset fields. The flipped outputs are then
27 averaged (according to the flip index) to get the final global heatmaps and the offset fields. For the
28 computation of local heatmaps and the local-global adaptation, only the non-flipped outputs are used
29 for the final predictions.

30 2 The Computation of the Empirical Upper Bound

31 We elaborate on the details of computing the empirical upper bound of performance for a vanilla
32 center-offset pose estimation method (Table.1 in the submission).

33 **Network Architecture.** The vanilla center-offset regression baseline uses the ImageNet pretrained
 34 HRNet-W32 [9] as the backbone, and the same modules as in our LOGO-CAP+HRNet-W32 for the
 35 center heatmap regression and the offset vector regression. See Fig.3 and Section 3 in the submission
 36 for detailed specifications. We present the details of computing keypoint expansion maps (KEMs)
 37 that are used in calculating the empirical upper bound as follows.

38 **Computation of Keypoint Expansion Maps.** Denoted by $\mathcal{P}_{N \times 17 \times 2}$ the initial pose parameters
 39 (i.e., the 2-D locations for the 17 keypoints of the N pose instances) estimated by the vanilla center-
 40 offset method, we expand each of the estimated keypoints with a local 11×11 mesh grid, that is to
 41 lift a keypoint to a 2-D mesh to counter the estimation uncertainty. As shown in the Alg. 1, we use
 42 the COCO benchmark provided keypoint sigmas to scale the unit length of the meshgrid for differ-
 43 ent types (e.g., nose, eyes, hips) of keypoints. After getting the expanded keypoint meshes \mathcal{M} of
 44 the initial poses, we compute their keypoint similarities $\mathcal{S}_{N \times 11 \times 11 \times K \times 17}$ between the groundtruth
 45 keypoints $\mathcal{G}_{K \times 17 \times 3}$ and the keypoint expansion maps. After applying the sum reduction on the simi-
 46 larity tensor $\mathcal{S}_{N \times 11 \times 11 \times K \times 17}$ along the 2-nd, 3-rd and the last axes, we have known the optimal cor-
 47 respondence (including the low-quality matches) for each center anchor, denoted by $\mathcal{S}_{N \times 11 \times 11 \times 17}$.
 48 Then, the pose with the maximal similarity in the 11×11 local window for each center anchor are
 49 used as the best one to compute the empirical upper bound on the fully-annotated COCO-val-2017
 50 dataset.

Algorithm 1: Computation of Keypoint Expansion Maps in a PyTorch-like style

```

1 coco_sigmas = torch.tensor([0.026, 0.025, 0.025, 0.035, 0.035, 0.079, 0.079, 0.072, 0.072,
  0.062, 0.062, 0.107, 0.107, 0.087, 0.087, 0.089, 0.089 ]) #Keypoint sigmas for the
  COCO dataset.
2 def KptsExpansionCoco (P, ksize=11)
  #Initial poses P: Nx17x2
3   radius = ksize // 2
4   dy, dx = torch.meshgrid(torch.arange(-radius, radius), torch.arange(-radius, radius))
  #dx, dy: ksize x ksize
5   dy, dx = dy.reshape(1, 1, ksize, ksize), dx.reshape(1, 1, ksize, ksize)
6   scale = coco_sigmas.reshape(1, 17, 1, 1) / coco_sigmas.min() #using different
  expansion rate for the different keypoint categories.
7   dy = dy * scale #dy: 1x17x11x11
8   dx = dx * scale #dx: 1x17x11x11
9   dxy = torch.stack((dx, dy), dim=-1) #dxy: 1x17x11x11x2
10  M = P.reshape(N, 17, 1, 1, 2) + dxy #M: Nx17x11x11x2
11  return M

```

51 3 The Speed-Accuracy Comparisons

52 We elaborate on the speed-accuracy comparisons in Fig. 2 in the submission, which is reproduced
 53 in Fig. 1.

54 We compare our LOGO-CAP with state-of-the-art bottom-up pose estimation approaches in terms
 55 of the speed-accuracy trade-off using different input resolutions of testing images.

- 56 • Both the proposed LOGO-CAP method and the DEKR [2] method have two versions of
 57 models trained with the HRNet-W32 backbone at the resolution of 512×512 , and the
 58 HRNet-W48 backbone at the resolution of 640×640 , respectively. The two versions of
 59 models are denoted by $W32$ and $W48$ respectively in Fig. 1. In comparing the speed-
 60 accuracy trade-offs, we evaluate each of the two versions of models with two different
 61 resolutions: For $W32$, we test the models, LOGO-CAP and DEKR, under the resolutions
 62 of 384 and 512 for short side of the testing images, which are denoted by $W32 - 384$ and
 63 $W32 - 512$ in Fig. 1 respectively. Similarly, it is done for the $W48$ models. The longer
 64 side of the testing images are computed according to their original aspect ratios.

- 65 • For the CenterNet method [11], we report the result obtained using the large Hourglass
66 backbone [8] due to its better AP score. The testing resolution is 512 in the short side of
67 the images.
- 68 • For the PifPaf method [5], we report their results using three different backbones of ResNet-
69 50, ResNet-101 and ResNet-152 [3]. We follow the image resolution setting presented in
70 the original paper [5] that resizes the long side of the testing images to be 641 and keep the
71 original aspect ratio.
- 72 • For the associative embedding approach [7], the result using the Higher-HRNet-W32 [1] as
73 the backbone is reported for the concern of inference speed. The testing resolution is also
74 512 in the short side of the testing image.

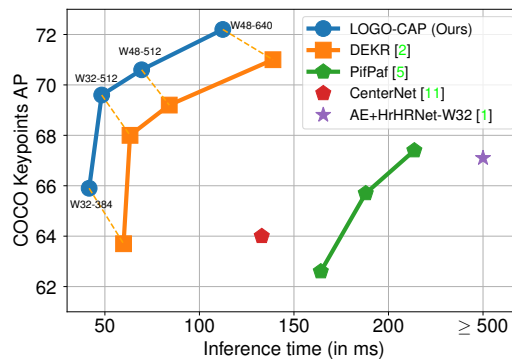


Figure 1: Speed-accuracy comparisons between our LOGO-CAP and prior arts on the COCO val-2017 dataset. $Wx-Y$ (e.g. W32-384) means that a model uses the backbone HRNet- Wx (HRNet-W32) and is tested with the image resolution Y in the short side.

75 As shown in Fig. 1, our approach obtains the best speed-accuracy trade-off with different backbones
76 and image resolutions. When using HRNet-W32 as the backbone, our LOGO-CAP uses the low-
77 resolution images (with 384 pixels in the short side) obtains the AP of 65.9 while approaching
78 real-time performance with the FPS of 24.0, which surpasses DEKR [2] by 2.2 AP and
79 7.3 FPS. Benefitting from our design rationales of simplicity and fully end-to-end learning, our
80 approach also obtains the best performance in both aspects of speed and accuracy when increasing
81 the image resolution to 512 and 640.

82 4 More Qualitative Results

83 **Results on the COCO-val-2017 and the OCHuman Datasets.** Fig. 2 shows examples of pose
84 estimation in the two datasets by the proposed LOGO-CAP with the HRNet-W32 backbone. Our
85 proposed LOCO-CAP is able to handle large structural and appearance variations in human pose
86 estimation.

87 **Fast pose estimation for video frames.** To justify the potential of our proposed approach in prac-
88 tical applications, we run our LOGO-CAP (W32 model) on two videos that have the resolution of
89 1280×720 from YouTube. We follow our testing protocol to resize the short side of the video frames
90 to 512 pixels and keep their original aspect ratios for inference. Without using any pose tracking
91 techniques, our LOGO-CAP achieves fast and accurate human pose estimation. Please click the
92 following anonymous google drive links for the demo videos:

- 93 - <https://bit.ly/3z2t8fA> (video credit: <https://youtu.be/2DiQUX11YaY>)
- 94 - <https://bit.ly/3cghWT4> (video credit: <https://youtu.be/kTzvU1sGSyA>)

95 In these two demo videos, the instantaneous FPS for each video frame is marked in the left corner
96 of the video.



Figure 2: Qualitative results of our LOGO-CAP (HRNet-W32). All images were picked thematically without considering our algorithms performance. The first three rows display our approach on the COCO-val-2017 dataset and the last three ones show our results on the OCHuman test dataset.

97 References

- 98 [1] Bowen Cheng, Bin Xiao, Jingdong Wang, Honghui Shi, Thomas S. Huang, and Lei Zhang. Higherhrnet:
 99 Scale-aware representation learning for bottom-up human pose estimation. In *IEEE/CVF Conference on*
 100 *Computer Vision and Pattern Recognition (CVPR)*, pages 5385–5394. IEEE, 2020. 1, 3
- 101 [2] Zigang Geng, Ke Sun, Bin Xiao, Zhaoxiang Zhang, and Jingdong Wang. Bottom-up human pose es-
 102 timation via disentangled keypoint regression. In *IEEE Conference on Computer Vision and Pattern*
 103 *Recognition (CVPR)*, 2021. 1, 2, 3
- 104 [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
 105 In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 3
- 106 [4] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and
 107 Yann LeCun, editors, *International Conference on Learning Representations (ICLR)*, 2015. 1

- 108 [5] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. Pifpaf: Composite fields for human pose estimation.
109 In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11977–11986, 2019. 3
- 110 [6] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
111 and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In David J. Fleet, Tomás Pajdla,
112 Bernt Schiele, and Tinne Tuytelaars, editors, *European Conference on Computer Vision (ECCV)*, volume
113 8693, pages 740–755, 2014. 1
- 114 [7] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint
115 detection and grouping. In *Advances in Neural Information Processing Systems 30 (NeurIPS)*, pages
116 2277–2287, 2017. 1, 3
- 117 [8] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In
118 *European Conference on Computer Vision (ECCV)*, pages 483–499, 2016. 3
- 119 [9] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for
120 human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
121 5693–5703, 2019. 1, 2
- 122 [10] Feng Zhang, Xiatian Zhu, Hanbin Dai, Mao Ye, and Ce Zhu. Distribution-aware coordinate representa-
123 tion for human pose estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*
124 *(CVPR)*, pages 7091–7100, 2020. 1
- 125 [11] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *CoRR*, abs/1904.07850, 2019. 3