

SmallPose: A Large-Scale Benchmark for Small Human Pose Estimation

Anonymous CVPR submission

Paper ID 9341

Abstract

001 *Various sizes of human body in the wild pose significant challenges for human pose estimation. A robust human*
002 *pose estimation model must possess not only cross-scenario generalizability but also the capability for stable predictions across different instance sizes. From the perspective of data-driven model performance, while prior dataset works*
003 *have mainly focused on scene complexity, we argue that even in simple scenarios, human body size is also a key*
004 *factor affecting human pose estimation, with small human instances being particularly critical. To fill this gap, we*
005 *present **SmallPose**, the first large-scale dataset that specifically targets the human body size with a focus on the challenge of small-size instances. Additionally, to thoroughly*
006 *assess model performance on small instances, we introduce the concept of **keypoint density** and establish three novel*
007 *fine-grained evaluation metrics based on it. By employing a systematic analysis with these metrics on SmallPose,*
008 *we dissect the efficacy of current pose estimation methodologies across various design paradigms and network ar-*
009 *chitectures. Experimental results demonstrate that models trained on SmallPose not only improve in performance on*
010 *small instances but also exhibit increased robustness. The SmallPose dataset and the models of previous works opti-*
011 *mized on it will be released soon.*

tion and training paradigm selection, the characteristics of the dataset significantly influence the model's performance and robustness.

However, widely used datasets such as MPII[1] and COCO[16] primarily focus on relatively simple everyday scenarios, while datasets like CrowdPose[13] and HiEve[17] concentrate on crowded or complex event scenes. Although these datasets include small human instances, the number of such instances is relatively low. Moreover, they lack appropriate metrics to evaluate model performance specifically on small-sized human instances, leading to a gap in assessing and improving models for these cases.

In this paper, we have conducted a thorough analysis of existing datasets, leading to the formulation of **SmallPose**, the first large-scale dataset dedicated to the diverse sizes of human instances, with an emphasis on small instances. Figure 1(b) depicts several samples from our dataset that showcase the range of human sizes across diverse real-world settings. Notably, SmallPose also encompasses a broad range of scenarios from simple to complex, but its primary aim is to highlight the ubiquity of small instances across diverse settings. We have also formulated a new taxonomy for small instances, defined by area thresholds, and proposed the performance metric specifically for small human instances. Additionally, we observed that small human instances, such as those where only the human head is visible, pose unique challenges for evaluation due to the incomplete representation of the human body. The coarse-grained metrics are inadequate in these cases, as keypoint density can vary significantly, affecting the estimation difficulty. To address this, we introduced the concept of **keypoint density**. This concept is delineated into three gradations: sparse, moderate, and dense, which categorize the density levels of small human poses. Importantly, the instances with lower keypoint density do not necessarily mean these are easier to estimate. It may contain partial bodies, which may increase the challenge in terms of recognition and localization. We've also established several fine-grained metrics based on it to evaluate model performance on small instances. Employing the SmallPose dataset and novel metrics, we conducted

1. Introduction

Recent advancements in human pose estimation have been substantially driven by deep neural networks[21, 24, 25, 28, 30, 31] trained on large-scale datasets[1, 2, 13, 16, 17, 29]. However, the effectiveness of these methods is often constrained by the complexity and diversity inherent in real-world scenarios. In practical applications such as autonomous driving[22, 33] and video surveillance[5, 6, 23], the challenge of accurately estimating human poses is further amplified due to the variability in human body sizes. This is because the human instances of varying sizes makes it challenging for networks to maintain consistent performance. In addition to strategies such as network customiza-

038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078

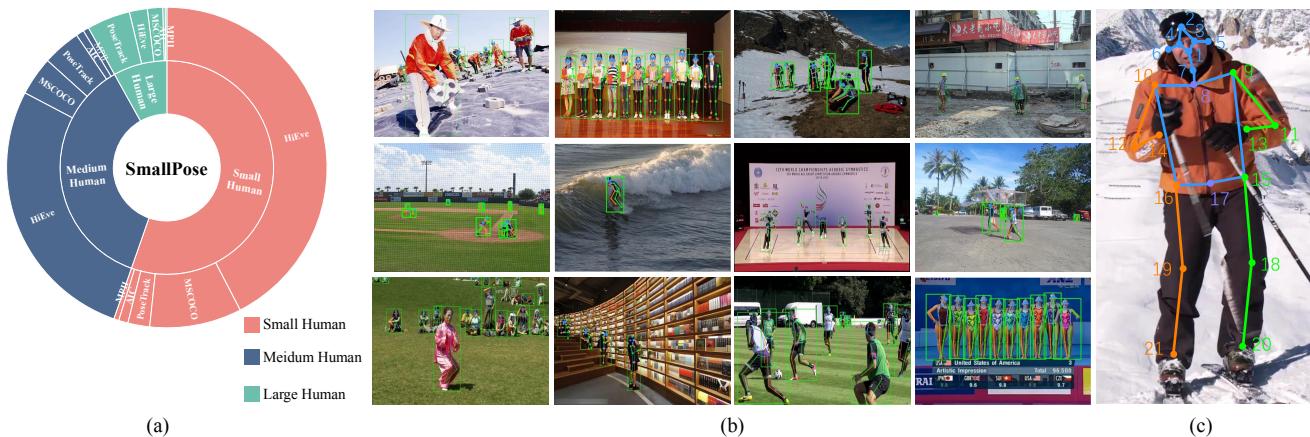


Figure 1. (a) The detailed breakdown of the SmallPose dataset composition. (b) The visualization of samples in the SmallPose dataset. Our dataset covers a wide range of human sizes from small to large, focusing on the ubiquity of small human bodies. (c) The proposed new human body skeleton. It is more comprehensive than the skeleton of the existing dataset, including 21 joints.

079 a comparative analysis of diverse models. The experimental results suggest that models trained on SmallPose exhibit 080 enhanced performance on small instances, with improved 081 robustness. Our analysis further reveals the subtle interplay 082 between paradigm choice, network architecture, and 083 the specific demands of human pose estimation, particularly 084 in scenarios involving human instances of different sizes. 085 We summarize the contributions as follows:

- 087 1. We introduce **SmallPose**, the first dataset specifically 088 designed to address the challenge of varying human instance 089 sizes, aimed at advancing the robustness of human 090 pose estimation.
- 091 2. We propose the novel concept of **keypoint density** and 092 establish three new fine-grained metrics based on it to 093 assess model performance in the real world.
- 094 3. We systematically analyzed existing human pose estimation 095 methods using SmallPose and new evaluation metrics, 096 revealing the impact of different network designs 097 and paradigms.
- 098 4. The SmallPose dataset and the models of previous works 099 optimized on it will be released soon.

100 2. Related Work

101 2.1. 2D Human Pose Estimation Dataset

102 In the nascent stages of human pose estimation, datasets 103 such as LSP[12], FashionPose[7], PASCAL Person 104 Layout[8], and J-HMDB[11] were pivotal, focusing on 105 single-person pose estimation within controlled environments. 106 These foundational datasets laid the groundwork for 107 understanding and predicting human pose in simplified 108 settings. With the development of this field, the demand for 109 estimating poses in multi-person scenarios has become 110 apparent, leading to the development of datasets like MPII[1],

COCO[16], and AI Challenger[29]. These datasets further enhance the real-world applicability of pose estimation models. Afterward, the advent of crowded scenes introduced a new set of challenges, prompting the creation of datasets such as PoseTrack[2], HiEve[17], and CrowdPose[13], designed to accurately identify and track human bodies under challenging conditions. While these datasets have propelled the field forward, a truly robust human pose estimation model demands not only cross-scenario generalizability but also stable predictions across a range of human sizes. To this end, we have introduced SmallPose, the first large-scale dataset to focus on body size. SmallPose encompasses human bodies of various sizes from simple to complex scenes, highlighting the ubiquity of small human instances across diverse settings. Empirical evidence demonstrates that models trained on SmallPose exhibit increased robustness.

128 2.2. 2D Human Pose Estimation Methods

129 The field of human pose estimation has progressed through 130 three primary categories: top-down, bottom-up, and single-stage 131 paradigms. The top-down paradigm[9, 24, 30, 31] has been dominant in public benchmarks, utilizing pre-trained 132 detectors to delineate individuals with bounding boxes, simplifying 133 the task by scaling the human body to a uniform size. In contrast, 134 the bottom-up paradigm[3, 4, 9, 19] shifts the focus to keypoint detection, independent of specific instances. By identifying and aggregating keypoints across 135 the entire image, it reconstructs human poses by emphasizing 136 the distribution of keypoints rather than individual bodies. The 137 single-stage paradigm[18, 20, 27, 35] embodies a simplified and 138 efficient methodology, conducting multi-person pose estimation 139 in a single forward pass. This approach directly regresses 140 keypoints from a predefined root,

111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143

	Image	Instance			
		Small	Medium	Large	Total
Train	68,311	351,722	225,950	49,453	627,125
Test	10,296	75,841	53,869	12,609	142,319
Total	78,607	427,563	279,819	62,062	769,444

Table 1. The number of images and the numbers for small, medium and large instances in each split.

144 aiming to enhance efficiency and speed without compromising accuracy. Leveraging the SmallPose dataset, we evaluated multi-person pose estimation methods grounded in these three dominant design paradigms. Our analysis highlights the importance of aligning the chosen paradigm with the network architecture to effectively tackle the unique challenges presented by diverse body sizes in human pose
145
146
147
148
149
150
151

152 3. Dataset

153 3.1. Dataset Organization and Description

154 **Data Collection.** To construct our dataset, we choose
155 the COCO[16], MPII[1], AIC[29], HiEve[17], and
156 Posetrack[2] datasets as the foundational sources. Upon
157 our verification, we confirmed that these five datasets permit
158 distribution, modification, and adjustment of the data.
159 During the phase of data collection, we devised a novel cri-
160 terion based on the bounding box area, as opposed to the
161 segmentation area used in COCO[16], to identify and fil-
162 ter samples containing small, medium, and large human
163 instances from these datasets. This strategic approach enabled
164 us to amass a substantial collection of images containing
165 small human instances while ensuring comprehensive cov-
166 erage of the entire spectrum of body sizes, culminating in
167 the creation of the SmallPose dataset. Figure 1(a) provides
168 a detailed breakdown of the dataset composition.

169 **Data Annotation.** In the process of annotating human
170 poses for the SmallPose dataset, we encountered the chal-
171 lenge of inconsistent label formats across existing datasets.
172 Specifically, the number of keypoints annotated varied, with
173 MPII[1] having 16 keypoints, AIC[29] and HiEve[17] using
174 14 keypoints, and PoseTrack[2] and COCO[16] employing
175 17 keypoints. To address this challenge, we con-
176 ducted a thorough analysis and proposed a comprehensive
177 human keypoint annotation scheme that integrated the la-
178 bel formats from the existing datasets. This unified scheme
179 encompasses a total of 21 keypoints, with the following
180 numbering and corresponding body parts: 1-nose, 2-head
181 top, 3-left eye, 4-right eye, 5-left ear, 6-right ear, 7-head
182 bottom, 8-chest, 9-left shoulder, 10-right shoulder, 11-left
183 elbow, 12-right elbow, 13-left wrist, 14-right wrist, 15-left
184 hip, 16-right hip, 17-pelvis, 18-left knee, 19-right knee, 20-
185 left ankle, and 21-right ankle. Additionally, we designed

a new set of body skeleton connection rules to more accu-
186 rately capture the human body structure, as illustrated in
187 Figure 1(c).
188

189 3.2. Dataset Statistics

Dataset Size. Our SmallPose dataset comprises over 70k
190 images and more than 700k human instances, with a sig-
191 nificant portion of over 420k small instances, substantially
192 exceeding the scale of other existing datasets. Specifically,
193 as detailed in Table 1, the training set consists of 68,311
194 images, with 351,722 small, 255,950 medium, and 49,453
195 large human instances, a ratio of approximately 56:36:8.
196 The testing set includes 10,296 images, with 75,841 small,
197 53,869 medium, and 12,609 large human instances, main-
198 taining a similar ratio of 53:39:8. Across the entire dataset,
199 small human instances account for 55.6%, medium hu-
200 man instances for 36.4%, and large human instances for
201 8%. Due to the abundance of small human samples in the
202 dataset, it is able to demonstrate the prevalence of small in-
203 stances in various scenarios.

Human Area Distribution. We conducted a comparative
204 analysis for the distributions of human body size across
205 six datasets: MPII[1], AIC[29], PoseTrack[2], COCO[16],
206 HiEve[17], and the proposed SmallPose dataset. The Fig-
207 ure 2 reveals that the MPII[1], AIC[29], and PoseTrack[2]
208 datasets exhibit a more concentrated distribution of sam-
209 ples in the large human size range [256 × 192, +∞). This
210 suggests that models trained on these datasets may struggle
211 to handle small human instances, potentially compromis-
212 ing their overall robustness. Additionally, while the COCO[16]
213 and HiEve[17] datasets have a wider distribution of samples
214 in the small human size range (0, 128 × 96), their coverage
215 and sample scale are not as comprehensive and extensive
216 as those of the SmallPose dataset. the SmallPose, on the
217 218

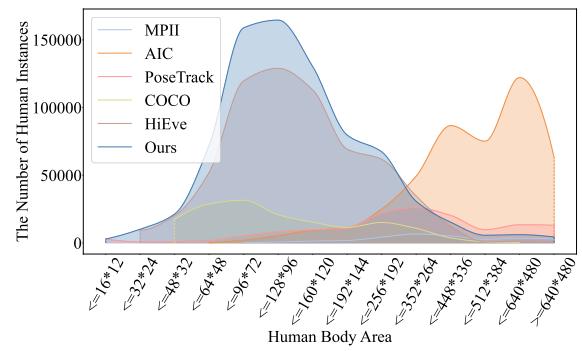


Figure 2. The distribution of human body area in six datasets:
MPII[1], AIC[29], PoseTrack[2], COCO[16], HiEve[17], and ours. The proposed SmallPose dataset features a significantly larger sample size of small and medium-sized human instances and a more comprehensive size coverage, showcasing its potential advantages in training robust pose estimation models.

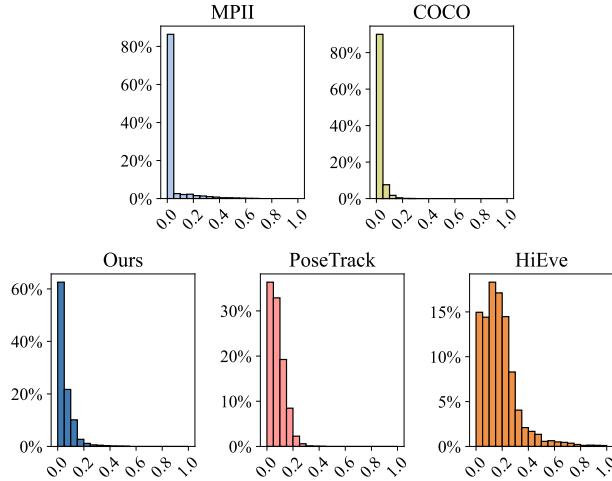


Figure 3. The CrowdIndex distribution of five datasets: MPII[1], COCO[16], PoseTrack[2], HiEve[17], and ours. The proportion of crowded scenes in SmallPose is not significant, indicating that there is no direct correlation between human body size and scene complexity.

other hand, excels with its extensive collection of small and medium human instances, highlighting its benefits for training robust pose estimation models. **CrowdIndex Distribution.** We utilized the CrowdIndex (CI) metric defined in the CrowdPose[13] dataset to quantify the degree of crowding across different datasets. For a given image, the CI is calculated using the formula:

$$CI = \frac{1}{n} \sum_{i=1}^n \frac{N_i^b}{N_i^a}, \quad (1)$$

where n is the total number of human instances, N_i^a represents the number of keypoints belonging to the i -th instance, and N_i^b denotes the number of keypoints within the i -th instance's bounding box but not belonging to that instance. We evaluated the CI distributions for five datasets: MPII[1], COCO[16], PoseTrack[2], HiEve[17], and the proposed SmallPose. As shown in Figure 3, the results indicate that the MPII[1], and COCO[16] datasets are dominated by simple scenarios, while PoseTrack[2], and HiEve[17] contain a higher proportion of complex scenes. In contrast, the SmallPose dataset focuses on the challenge of small human instances, yet the proportion of crowded scenes is not significantly high, suggesting that human body size and scene complexity are not directly correlated.

4. Evaluation Metrics

4.1. Size-based Metrics

Existing datasets and evaluation metrics lacked a focus on small instances. To address this, we introduced size-based evaluation metrics. Specifically, we defined small,

medium, and large human instances based on bounding box areas: small ($0, 128 \times 96$), medium [$128 \times 96, 256 \times 192$], and large [$256 \times 192, +\infty$]). We then proposed the AP^S , AP^M , and AP^L metrics to separately evaluate model performance on each size category. This evaluation approach allowed us to more comprehensively assess the models' capabilities in handling instances of varying sizes - a crucial aspect for developing robust and practical human pose estimation systems. For the sake of simplicity of the formula, we will write the defining intervals of small, medium, and large instances as R^S, R^M , and R^L . The detailed formulas are as follows:

$$\begin{aligned} AP^S &= \frac{\sum_p \delta(OKS_p > T)}{\sum_p \delta(S_p \in R^S)}, \\ AP^M &= \frac{\sum_p \delta(OKS_p > T)}{\sum_p \delta(S_p \in R^M)}, \\ AP^L &= \frac{\sum_p \delta(OKS_p > T)}{\sum_p \delta(S_p \in R^L)}, \end{aligned} \quad (2)$$

where S_p represents the ground-truth bounding box area of individual p , T is the threshold, and only an OKS score greater than T is considered accurate.

4.2. Definition of Keypoint Density for Small Human Instance

In the SmallPose dataset, we observed that small human instances, such as when only the head of a person is visible, present unique challenges for evaluation due to the incomplete human body. The coarse-grained metrics are not suitable for these cases, as keypoint density can vary significantly, affecting estimation difficulty. To address this issue, we introduce the keypoint density to provide a more fine-grained evaluation of small human instances under different situations. The specific formula is as follows:

$$KD_p = 2 \left[1 - \sigma \left(\frac{S_p}{S_{small}} \cdot \frac{K}{n_p} \right) \right] \in (0, 1), \quad (3)$$

where p represents the p -th human instance, S_p denotes the bounding box area of the p -th instance, n_p is the number of keypoint annotations for the p -th instance, S_{small} represents the maximum of the defined area range of small human body, which is 128×96 , and K is the number of keypoints in the defined skeleton, such as 17 for COCO[16] and 14 for HiEve[17]. The sigmoid function $\sigma(\cdot)$ is used to restrict the keypoint density value to the range of (0, 1).

4.3. Keypoint-density-based Metrics for Small Human Instance

Based on the definition of keypoint density, we analyzed the relationship between keypoint density and model performance, as shown in Figure 4. To provide further analysis, we defined three categories based on keypoint density: sparse (0-0.35), moderate (0.35-0.95), and dense (0.95-1.0).

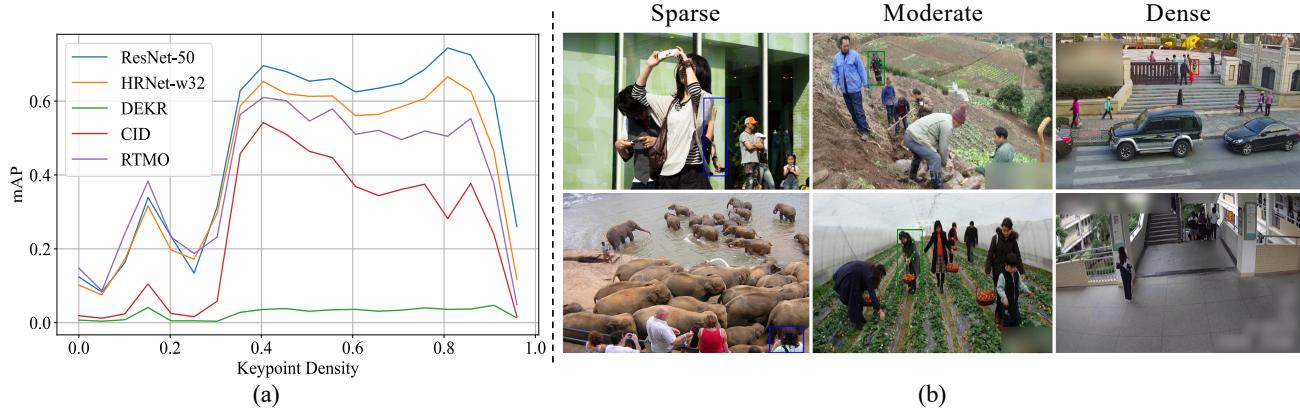


Figure 4. (a) The comparisons of mAP of widely used models, including ResNet[10], HRNet[24], DEKR[9], CID[27] and RTMO[18]. The x-axis represents the keypoint density for small human instances. The results consistently show that the accuracy is lower for both low and high keypoint density, while it is higher for a moderate keypoint density. (b) The visualization of the sparse, moderate, and dense small human instances. The sparse category generally consists of incomplete small instances, while the dense category predominantly contains small instances with dense keypoint annotations. These three types of small instances present varying levels of difficulty.

According to this classification, we introduced the AP_{Sp}^S , AP_{Mo}^S , and AP_{De}^S fine-grained metrics to evaluate model performance on each type of small instances. This fine-grained analysis provided deeper insights into the models' strengths and limitations in handling small human instances with varying keypoint densities. The detailed formulas are as follows:

$$\begin{aligned} AP_{Sp}^S &= \frac{\sum_p \delta(OKS_p > T)}{\sum_{area_p \in R^S} \delta(KD_p \in (0.0, 0.35))}, \\ AP_{Mo}^S &= \frac{\sum_p \delta(OKS_p > T)}{\sum_{area_p \in R^S} \delta(KD_p \in [0.35, 0.95])}, \\ AP_{De}^S &= \frac{\sum_p \delta(OKS_p > T)}{\sum_{area_p \in R^S} \delta(KD_p \in [0.95, 1.0])}, \end{aligned} \quad (4)$$

5. Experimental Setup

5.1. Datasets

SmallPose. Our proposed SmallPose dataset comprises a total of 78,607 images and 769,444 human instances, with a remarkable 427,563 small instances. The SmallPose dataset was designed to improve model performance and robustness in handling instances of diverse sizes, particularly small ones, by data samples of various human body sizes.

COCO. The COCO dataset[16], which we used for comparative evaluation, contains a training set of 118,287 images and 156,165 human instances, of which only 98,582 are small instances. Moreover, the COCO samples are predominantly from simple scenes.

5.2. Evaluation Metrics

In our evaluation, we adhere to the metrics defined in Section 4. We employ AP , AP^S , AP^M , AP^L , AP_{Sp}^S , AP_{Mo}^S , AP_{De}^S , and the average recall (AR) to evaluate the

performance. Notably, in setting A, due to the inconsistency between the number of keypoint types in the COCO dataset[16] and our testing set, we only report the accuracy for the keypoint types present in COCO.

5.3. Experimental Settings

To thoroughly analyze various human pose estimation paradigms and network architectures, we have established two distinct experimental settings:

Setting A: The models were trained on the COCO training set and then tested on our SmallPose testing set, with the keypoint categories aligned.

Setting B: The models were trained on our SmallPose training set and then tested on our SmallPose testing set.

This experimental design allowed us to systematically evaluate the models' performance and generalization capabilities across different training conditions, particularly in relation to handling instances of varying sizes.

5.4. Human Pose Estimation Methods

To comprehensively evaluate the performance of existing models, we employed some representative works from three paradigms (top-down, bottom-up, and single-stage), two types of backbones (single-branch and multi-branch), three prediction architectures (heatmap-based, regression-based, and classification-based), and two different network architectures (CNN-based and Transformer-based) to benchmark the dataset.

SimpleBaseline[30] simply utilizes ResNet[10] as its backbone to regress heatmaps for keypoint localization.

HRNet[24] begins with a high-resolution subnetwork and progressively adds high-to-low resolution subnetworks in later stages, connecting them in parallel. It performs re-

Table 2. The comparative results of various design paradigms and backbones under setting A and setting B.

Method	BackBone	Setting	AP	AP^S	AP^M	AP^L	AP_{Sp}^S	AP_{Mo}^S	AP_{De}^S	AR
Top-down methods										
SimpleBaseline [30]	ResNet-50	A	43.2	32.4	50.5	77.3	8.8	33.9	0.1	45.8
		B	68.2	65.9	69.1	79.7	32.8	69.4	26.0	69.1
HRNet [24]	HRNet-W32	A	45.0	34.3	52.3	79.6	10.3	36.0	0.1	47.7
		B	66.4	62.2	69.3	82.1	32.0	65.3	11.7	67.7
SimCC [15]	ResNet-50	A	43.2	32.9	50.3	76.9	18.7	33.8	0.0	46.2
		B	68.8	68.0	67.3	78.5	38.6	70.5	63.6	69.4
CAL [26]	HRNet-W32	A	45.9	34.4	54.1	81.2	11.0	35.8	0.1	48.4
		B	69.6	66.3	70.8	85.2	35.8	68.9	33.3	70.5
Bottom-up methods										
DEKR [9]	HRNet-W32	A	32.5	15.6	46.5	78.3	3.6	15.9	0.0	36.8
		B	23.2	20.1	32.4	22.8	4.0	20.0	1.3	46.6
Single-stage methods										
CID [27]	HRNet-W32	A	36.6	20.5	47.4	78.4	4.1	21.2	0.0	41.1
		B	60.9	49.4	71.6	83.3	11.8	52.1	1.7	63.0
RTMO [18]	CSPDarknet	A	33.2	17.2	43.7	79.4	7.7	17.3	0.0	37.4
		B	63.1	55.1	70.6	79.9	29.7	58.1	4.9	64.3

peated multi-scale fusion, allowing each resolution to exchange information with others, resulting in rich high-resolution representations for keypoint localization.

DEKR[9] proposes separating regression and adaptive convolution to enhance the quality of regression.

CID[27] utilizes attention mechanisms to separate the instance-level features, and generate instance-level heatmaps for keypoint localization.

RTMO[18] seamlessly integrates coordinate classification by representing keypoints using dual 1-D heatmaps, achieving accuracy comparable to top-down methods while maintaining high speed.

HRFormer[32] leverages the multi-resolution parallel design of HRNet[24] and local-window self-attention to improve memory and computation efficiency.

ViTPose[31] employs plain, non-hierarchical vision transformers as backbones to extract features and a lightweight decoder for pose estimation, demonstrating the strong capabilities of plain vision transformers for pose estimation.

RLE[14] introduces Residual Log-likelihood Estimation (RLE), a regression paradigm that captures the output distribution by learning its change, rather than the unreference underlying distribution, to facilitate training.

ViTPose-RLE integrates the ViTPose-S backbone with the RLE paradigm, forming a competitive regression-based method, ViTPose-S-RLE.

SimCC[15] reformulates HPE as two classification tasks for horizontal and vertical coordinates, dividing each pixel into bins to achieve sub-pixel localization and low quantization error. This approach eliminates the need for additional

refinement and upsampling layers in certain settings, resulting in a simpler and more efficient HPE pipeline.

CAL[26] introduces Confidence-Aware Learning (CAL) to address two key limitations of offset learning: inconsistent training and testing, decoupled heatmap and offset learning. This method significantly outperforms state-of-the-art approaches for low-resolution human pose estimation.

5.5. Implementation Details

The experiments utilized the PyTorch-based MMPoser framework. For setting A, we loaded model weights pre-trained on COCO[16] provided by MMPoser to test on the SmallPose testing dataset. For setting B, Specifically, the SimpleBaseline[30] model employed the Adam optimizer with a learning rate of 5e-4, a batch size of 128, and was trained for 210 epochs. HRNet[24] used the Adam optimizer with a learning rate of 5e-6, a batch size of 64, and was trained for 210 epochs. DEKR[9] adopted the Adam optimizer with a learning rate of 1e-3, a batch size of 20, and was trained for 140 epochs. CID[27] also used the Adam optimizer with a learning rate of 1e-3, a batch size of 16, and was trained for 140 epochs. RTMO[18] employed the AdamW optimizer with an initial learning rate of 4e-3, a batch size of 32, and was trained for 600 epochs. HRFormer-S[32], ViTPose-S[31] and ViTPose-S-RLE all adopted the AdamW optimizer with a learning rate of 5e-4, a batch size of 64, and was trained for 210 epochs. RLE[14] employed the Adam optimizer with an initial learning rate of 1e-3, a batch size of 128, and was trained for 210 epochs. SimCC[15] also employed the Adam optimizer with an ini-

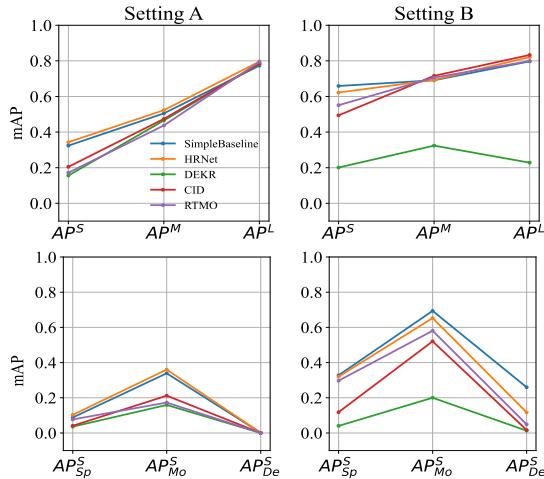


Figure 5. The performance variations across diverse design paradigms and backbones under setting A and setting B.

404 trial learning rate of 1e-3, a batch size of 192, and was
 405 trained for 210 epochs. All experiments were conducted
 406 on ten 16GB Nvidia V100 GPUs.

407 6. Results and Analysis

408 6.1. The Impact on Different Design Paradigms

409 Our experiments evaluated several popular models across
 410 different design paradigms (top-down, single-stage,
 411 bottom-up) trained on setting A and B. We present a visual
 412 analysis of the experimental results from Table 2, as de-
 413 picted in Figure 5. Our comparative analysis demonstrated
 414 that model performance consistently increased with larger
 415 human instances. However, notable performance gaps
 416 exist across various instance sizes, indicating the models
 417 prioritized larger instances while struggling with smaller
 418 ones, limiting robustness. In contrast, top-down and single-
 419 stage models trained on SmallPose achieved consistent
 420 gains across all sizes, particularly for small instances.
 421 The performance gap between different instance sizes
 422 was effectively narrowed, enhancing overall robustness.
 423 Notably, in the top-down approach, the SimpleBaseline
 424 method showed nearly equivalent accuracy on small and
 425 medium instances, while the SimCC performed even better
 426 on small instances than on medium instances. Conversely,
 427 bottom-up models exhibited an inverse trend, seemingly
 428 sacrificing the performance for larger instance to compen-
 429 sate for smaller instances. Analysis of keypoint density
 430 levels showed the performance of all models trained on
 431 SmallPose is improved across sparse, moderate, and dense
 432 ranges. In summary, SmallPose enhanced models’ ability
 433 to handle small instances and reduced performance gaps
 434 between different sizes, yielding more robust systems.

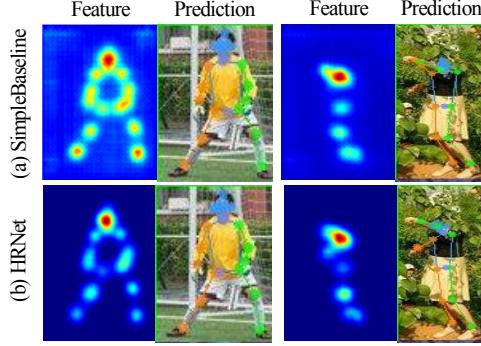


Figure 6. The feature visualization of SimpleBaseline and HRNet for small instances under setting B. It shows that SimpleBaseline offers a more confident and accurate keypoint prediction for small instances, contrasting with HRNet’s lower confidence and greater vulnerability to occlusions.

Top-down models exhibited superior performance for small instances and better robustness compared to single-stage. Bottom-up models face conflicts in handling instances of varying sizes, prioritizing learning one category based on the data distribution.

440 6.2. The Impact on Different Backbones

Our analysis distinguished between single-branch and multi-branch network architectures based on their handling of feature resolution. Single-branch networks like ResNet and CSPDarknet, which decrease in feature resolution with increasing depth, showed enhanced learning capabilities for small instances. This was evidenced by ResNet’s 3.7AP improvement over HRNet on the AP^S metric and the CSPDarknet-based RTMO’s 5.7AP lead over the HRNet-w32-based CID on the same metric. Features visualization for ResNet and HRNet in Figure 6 indicated that ResNet offers a more confident and accurate keypoint prediction for small instances, contrasting with HRNet’s lower confidence and greater vulnerability to occlusions. These insights emphasize the key role of backbone architectures in advancing model performance for specific tasks, notably the superiority of single-branch networks in detecting smaller human instances.

458 6.3. The Impact on Different Prediction Paradigms

We conducted a comparison of the performance of heatmap-based, regression-based and classification-based methods for small human pose estimation. As depicted in the Table 3, We highlight the distinct advantages of regression-based methods for the task of small human pose estimation. These methods excel by effectively utilizing global information, which is reflected in their superior performance across sparse and dense keypoint scenarios. Notably, the RLE shows an improvement over the SimpleBaseline,

435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457

Table 3. The comparison results of various prediction paradigms and network architectures under setting A and setting B.

Method	Backbone	Setting	AP	AP^S	AP^M	AP^L	AP_{Sp}^S	AP_{Mo}^S	AP_{De}^S	AR
Heatmap-based										
SimpleBaseline[30]	ResNet-50	A	43.2	32.4	50.5	77.3	8.8	33.9	0.1	45.8
		B	68.2	65.9	69.1	79.7	32.8	69.4	26.0	69.1
HRNet[24]	HRNet-W32	A	45.0	34.3	52.3	79.6	10.3	36.0	0.1	47.7
		B	66.4	62.2	69.3	82.1	32.0	65.3	11.7	67.7
HRFormer[32]	HRFormer-S	A	44.8	33.9	51.9	76.9	8.9	35.6	0.1	47.3
		B	68.7	66.1	69.1	81.6	33.7	69.1	35.4	69.4
ViTPose[31]	ViT-S	A	45.4	34.6	52.6	79.3	10.6	36.1	0.1	47.8
		B	68.9	66.9	69.1	80.2	33.0	70.3	27.1	69.6
Regression-based										
RLE[14]	ResNet-50	A	44.1	33.4	51.2	77.5	10.9	35.1	0.1	46.4
		B	67.8	65.5	68.3	79.0	36.8	68.7	31.4	68.4
ViTPose-S-RLE	ViT-S	A	29.2	19.8	34.9	63.8	8.9	19.9	0.0	38.9
		B	68.4	66.9	68.2	79.2	37.1	69.9	45.2	68.9
Classification-based										
SimCC[15]	ResNet-50	A	43.2	32.9	50.3	76.9	18.7	33.8	0.0	46.2
		B	68.8	68.0	67.3	78.5	38.6	70.5	63.6	69.4

and a similar trend is observed with ViTPose-S-RLE outperforming ViTPose-S. The integration of regression with transformer-based architectures further amplifies this advantage, as it allows for a more comprehensive capture of features vital for accurate estimation. This synergy underscores the importance of global information in enhancing pose estimation models, especially in the nuanced context of small human instances where other methods may falter. In addition, we found that the classification-based methods, SimCC, which is proposed for quantization error, achieves the best performance in small human instances, especially on dense small human instances, where the problem of quantization error caused becomes particularly serious due to the dense number of joints. This is sufficient to demonstrate that quantization error is an important factor affecting the performance of small human pose estimation. These findings highlight the potential advantages of classification-based methods in addressing the challenges posed by small instances in human pose estimation tasks.

6.4. The Impact on Different Network Architectures

We conducted a comparative analysis of CNN-based and Transformer-based human pose estimation models using the SmallPose dataset, as illustrated in the Table 3. The Transformer-based architecture, trained on the SmallPose dataset, demonstrated superior robustness across various keypoint density levels, particularly excelling at both sparse and dense levels. This indicates that the global attention mechanism utilized by the Transformer-based models pro-

vides a more precise detection for small-sized human instances compared to the local receptive fields of CNN-based models.

7. Conclusion

In this paper, we introduce SmallPose, the first dataset dedicated to the challenge of human pose estimation across varying instance sizes. The SmallPose dataset contains a variety of human sizes from simple to complex scenes, focusing on the ubiquity of small instances across diverse settings. With the introduction of the keypoint density concept and corresponding metrics, we have analyzed the human size robustness of different paradigms and different backbones and demonstrated that models trained on SmallPose show improved performance and robustness.

Limitations and Future Plans. While we utilized five benchmark datasets to construct the SmallPose dataset, it is possible that they may not fully represent the distribution of real-world scenarios. Additionally, the concept of keypoint density that we introduced is primarily focused on small human instances, but the notion of keypoint density is applicable to medium and large human bodies as well, albeit with less pronounced challenges. In the future, we can expand the keypoint density concept to encompass the full range of human sizes. This would enable a more comprehensive analysis and potentially provide additional insights for improving human pose estimation models.

522

References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE Conference on computer Vision and Pattern Recognition*, pages 3686–3693, 2014. 1, 2, 3, 4
- [2] Mykhaylo Andriluka, Umar Iqbal, Eldar Insafutdinov, Leonid Pishchulin, Anton Milan, Juergen Gall, and Bernt Schiele. Posetrack: A benchmark for human pose estimation and tracking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5167–5176, 2018. 1, 2, 3, 4
- [3] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7291–7299, 2017. 2
- [4] Bowen Cheng, Bin Xiao, Jingdong Wang, Honghui Shi, Thomas S Huang, and Lei Zhang. Higherhrnet: Scale-aware representation learning for bottom-up human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5386–5395, 2020. 2
- [5] Mickael Cormier, Aris Clepe, Andreas Specker, and Jürgen Beyerer. Where are we with human pose estimation in real-world surveillance? In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 591–601. IEEE, 2022. 1
- [6] Rita Cucchiara and Matteo Fabbri. Fine-grained human analysis under occlusions and perspective constraints in multi-media surveillance. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 18(1s):1–23, 2022. 1
- [7] Matthias Dantone, Juergen Gall, Christian Leistner, and Luc Van Gool. Human pose estimation using body parts dependent joint regressors. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3041–3048, 2013. 2
- [8] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338, 2010. 2
- [9] Zigang Geng, Ke Sun, Bin Xiao, Zhaoxiang Zhang, and Jingdong Wang. Bottom-up human pose estimation via disentangled keypoint regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14676–14686, 2021. 2, 5, 6
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [11] Hueihan Jhuang, Juergen Gall, Silvia Zuffi, Cordelia Schmid, and Michael J Black. Towards understanding action recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 3192–3199, 2013. 2
- [12] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *bmvc*, page 5. Aberystwyth, UK, 2010. 2

- [13] Jiefeng Li, Can Wang, Hao Zhu, Yihuan Mao, Hao-Shu Fang, and Cewu Lu. Crowdpose: Efficient crowded scenes pose estimation and a new benchmark. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10863–10872, 2019. 1, 2, 4
- [14] Jiefeng Li, Siyuan Bian, Ailing Zeng, Can Wang, Bo Pang, Wentao Liu, and Cewu Lu. Human pose regression with residual log-likelihood estimation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11025–11034, 2021. 6, 8
- [15] Yanjie Li, Sen Yang, Peidong Liu, Shoukui Zhang, Yunxiao Wang, Zhicheng Wang, Wankou Yang, and Shu-Tao Xia. Simcc: A simple coordinate classification perspective for human pose estimation. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022. 6, 8
- [16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755, 2014. 1, 2, 3, 4, 5, 6
- [17] Weiyao Lin, Huabin Liu, Shizhan Liu, Yuxi Li, Rui Qian, Tao Wang, Ning Xu, Hongkai Xiong, Guo-Jun Qi, and Nicu Sebe. Human in events: A large-scale benchmark for human-centric video analysis in complex events. *arXiv preprint arXiv:2005.04490*, 2020. 1, 2, 3, 4
- [18] Peng Lu, Tao Jiang, Yining Li, Xiangtai Li, Kai Chen, and Wenming Yang. Rtm: Towards high-performance one-stage real-time multi-person pose estimation. *arXiv preprint arXiv:2312.07526*, 2023. 2, 5, 6
- [19] Zhengxiong Luo, Zhicheng Wang, Yan Huang, Liang Wang, Tieniu Tan, and Erjin Zhou. Rethinking the heatmap regression for bottom-up human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13264–13273, 2021. 2
- [20] Debapriya Maji, Soyeb Nagori, Manu Mathew, and Deepak Poddar. Yolo-pose: Enhancing yolo for multi person pose estimation using object keypoint similarity loss. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2637–2646, 2022. 2
- [21] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VIII 14*, pages 483–499. Springer, 2016. 1
- [22] Arindam Sengupta, Feng Jin, Renyuan Zhang, and Siyang Cao. mm-pose: Real-time human skeletal posture estimation using mmwave radars and cnns. *IEEE Sensors Journal*, 20(17):10032–10044, 2020. 1
- [23] Pourya Shamsolmoali, Masoumeh Zareapoor, Huiyu Zhou, and Jie Yang. Amil: Adversarial multi-instance learning for human pose estimation. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 16(1s):1–23, 2020. 1
- [24] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF conference on*

- 636 *computer vision and pattern recognition*, pages 5693–5703,
637 2019. 1, 2, 5, 6, 8
- 638 [25] Alexander Toshev and Christian Szegedy. Deeppose: Human
639 pose estimation via deep neural networks. In *Proceedings of*
640 *the IEEE conference on computer vision and pattern recog-*
641 *nition*, pages 1653–1660, 2014. 1
- 642 [26] Chen Wang, Feng Zhang, Xiatian Zhu, and Shuzhi Sam Ge.
643 Low-resolution human pose estimation. *Pattern Recognition*,
644 126:108579, 2022. 6
- 645 [27] Dongkai Wang and Shiliang Zhang. Contextual instance de-
646 coupling for robust multi-person pose estimation. In *Pro-*
647 *ceedings of the IEEE/CVF Conference on Computer Vision*
648 *and Pattern Recognition*, pages 11060–11068, 2022. 2, 5, 6
- 649 [28] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser
650 Sheikh. Convolutional pose machines. In *Proceedings of the*
651 *IEEE conference on Computer Vision and Pattern Recog-*
652 *nition*, pages 4724–4732, 2016. 1
- 653 [29] Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming
654 Yan, Rui Liang, Wenjia Wang, Shipei Zhou, Guosen Lin,
655 Yanwei Fu, et al. Ai challenger: A large-scale dataset
656 for going deeper in image understanding. *arXiv preprint*
657 *arXiv:1711.06475*, 2017. 1, 2, 3
- 658 [30] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines
659 for human pose estimation and tracking. In *Proceedings of*
660 *the European conference on computer vision (ECCV)*, pages
661 466–481, 2018. 1, 2, 5, 6, 8, 3
- 662 [31] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vit-
663 pose: Simple vision transformer baselines for human pose
664 estimation. *Advances in Neural Information Processing Sys-*
665 *tems*, 35:38571–38584, 2022. 1, 2, 6, 8
- 666 [32] Yuhui Yuan, Rao Fu, Lang Huang, Weihong Lin, Chao
667 Zhang, Xilin Chen, and Jingdong Wang. Hrformer: High-
668 resolution vision transformer for dense predict. *Advances*
669 *in neural information processing systems*, 34:7281–7293,
670 2021. 6, 8
- 671 [33] Andrei Zanfir, Mihai Zanfir, Alex Gorban, Jingwei Ji,
672 Yin Zhou, Dragomir Anguelov, and Cristian Sminchisescu.
673 Hum3dil: Semi-supervised multi-modal 3d humanpose es-
674 timation for autonomous driving. In *Conference on Robot*
675 *Learning*, pages 1114–1124. PMLR, PMLR, 2023. 1
- 676 [34] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva,
677 and Antonio Torralba. Places: A 10 million image database
678 for scene recognition. *IEEE Transactions on Pattern Analy-*
679 *sis and Machine Intelligence*, 2017. 1
- 680 [35] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Ob-
681 jects as points. *arXiv preprint arXiv:1904.07850*, 2019. 2