# When Human Pose Estimation Meets Robustness: Adversarial Algorithms and Benchmarks

Jiahang Wang [1†]    Sheng Jin [2,3]    Wentao Liu[4]    Weizhong Liu [1]    Chen Qian [4]    Ping Luo[2]

[1] Huazhong University of Science and Technology    [2] The University of Hong Kong

[3] SenseTime Research    [4] SenseTime Research and Tetras.AI

jiahangwangchn@gmail.com   {jinsheng, liuwentao, qianchen}@sensetime.com

liuweizhong@mail.hust.edu.cn   pluo@cs.hku.hk

## Abstract

*Human pose estimation is a fundamental yet challenging task in computer vision, which aims at localizing human anatomical keypoints. However, unlike human vision that is robust to various data corruptions such as blur and pixelation, current pose estimators are easily confused by these corruptions. This work comprehensively studies and addresses this problem by building rigorous robust benchmarks, termed COCO-C, MPII-C, and OCHuman-C, to evaluate the weaknesses of current advanced pose estimators, and a new algorithm termed AdvMix is proposed to improve their robustness in different corruptions. Our work has several unique benefits. (1) AdvMix is model-agnostic and capable in a wide-spectrum of pose estimation models. (2) AdvMix consists of adversarial augmentation and knowledge distillation. Adversarial augmentation contains two neural network modules that are trained jointly and competitively in an adversarial manner, where a generator network mixes different corrupted images to confuse a pose estimator, improving the robustness of the pose estimator by learning from harder samples. To compensate for the noise patterns by adversarial augmentation, knowledge distillation is applied to transfer clean pose structure knowledge to the target pose estimator. (3) Extensive experiments show that AdvMix significantly increases the robustness of pose estimations across a wide range of corruptions, while maintaining accuracy on clean data in various challenging benchmark datasets.*

## 1. Introduction

Human pose estimation (HPE) is a fundamental task for action recognition and video surveillance [28, 38, 25]. Although convolutional neural networks (CNNs) achieved great progress [39, 40, 36, 5, 30, 7] on challenging datasets
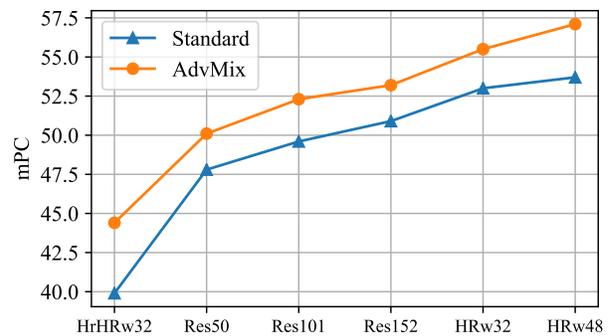


Figure 1. **Improvements** of model robustness (mPC) when AdvMix is applied to the state-of-the-art methods.

[26, 1, 47], which only contain clean and high-resolution images, deploying models in the real world requires not only good performance on clean data, but also robustness to commonly occurring image corruptions. For example, while tracking and estimating the keypoints of a moving person in outdoor environments, current pose estimators suffer severe performance drop due to the noise or blur caused by weather conditions or camera systems. Therefore, analyzing and enhancing the robustness of pose estimators are important and are the purposes of this work.

Unlike previous studies on common robustness for classification, detection and segmentation [17, 29, 23], human pose estimation uses a blend of classification and regression methods to model the structures of the human body, making it a challenging and collaborative field that is worthy of special investigations. The key challenges of robust human pose estimation are three folds. First, the lack of a benchmark for evaluating the robustness of state-of-the-art human pose estimation methods makes it difficult to construct rigorous comparisons between different models, not to mention to improve model robustness. Second, accuracy of clean data and corrupted data are trade-offs. Improving the robustness of the model while maintaining its performance

---

[†]The work was done during an internship at SenseTime Research.

on clean data is a non-trivial problem. Third, based on the proposed benchmark, we have examined the effectiveness of some data augmentation methods. However, we find that simply applying them sequentially does not achieve desirable performance. How can we effectively combine existing data augmentation techniques to improve the generalization of human pose estimators towards unforeseen corruptions?

Inspired by [17], we establish the robust pose benchmarks, consisting of three challenging datasets including COCO-C, MPII-C, and OCHuman-C. The benchmark datasets are constructed based on a full spectrum of **unforeseen** corruption types that are not encountered in model training (*i.e.* CNNs are trained on *clean* images, while evaluated on corrupted images). Extensive evaluations on these benchmarks show the weakness of both existing top-down and bottom-up pose estimators. (1) The state-of-the-art pose estimators suffer severe performance drop on corrupted images. (2) Models are generally more robust to brightness and weather changes, while less robust to motion and zoom blur. (3) The model robustness would increase by increasing model capacity.

Empirical evaluations on the proposed benchmarks help us screen a collection of useful data augmentation techniques to improve model robustness under severe corruptions. In order to make full use of these techniques and achieve optimal performance on **unforeseen** noisy data, we propose an augmentation generator, which *learns* to automatically combine augmented images. Specifically, we jointly train two neural networks in an adversarial manner, *i.e.* an augmentation generator and a human pose estimator. The generator produces weights to mix up randomly augmented images, while the pose estimator attempts to learn robust visual representation from harder training samples.

It is worth noting that the compositions produced by the augmentation generator may drift far from original images and such induced noise patterns may be harmful to performance on clean data. To reduce this negative impact, we propose to use a pre-trained teacher pose estimator to transfer structure knowledge learned from entire **clean** training data towards the target human pose estimator. Different from previous knowledge distillation methods that use a stronger network as the teacher model, our teacher pose estimator shares the same architecture as the target pose estimator. Extensive evaluations demonstrate that AdvMix significantly improves model robustness on diverse image corruptions while maintaining performance on clean data. The augmentation generator and the teacher pose estimator are only used for training and will be discarded at the inference stage, and thus introducing no computational overhead at inference time. Meanwhile, as shown in Fig. 1, our method is model-agnostic and is proved to be effective for various state-of-the-art pose estimation models.

Our main **contributions** can be summarized as follows.

- We propose three robust benchmarks COCO-C, MPII-C, and OCHuman-C, and demonstrate that both top-down and bottom-up pose estimators suffer severe performance drop on corrupted images, drawing the community's attention to this problem.

- With extensive experiments, we have many interesting conclusions that would help improve the accuracy and robustness of future works.

- We propose a novel adversarial data augmentation method together with knowledge distillation, termed AdvMix, which is model-agnostic and easy-to-implement. It significantly improves the robustness of pose estimation models while maintaining or slightly improving the performance on the clean data, without extra inference computational overhead.

## 2. Related Work

### 2.1. Human Pose Estimation

Human pose estimation (HPE) can be generally categorized into top-down and bottom-up methods. Top-down methods [39, 31, 13, 36, 22, 27, 11] divide the task into two stages: person detection and keypoint detection. SBL [40] presents a simple yet strong baseline network with several deconvolutional layers. HRNet [36] maintains high resolution representation combined with multi-scale feature fusions, and achieves the state-of-the-art performance on clean COCO dataset [26]. Bottom-up methods [5, 30, 24, 7, 21, 20] first detect all the keypoints and then group them into person instances. PifPaf [24] utilized a Part Intensity Field (PIF) to localize body parts and a Part Association Field (PAF) [5] to associate body parts to form full human poses. HigherHRNet [7] learns scale-aware representations using high-resolution feature pyramids, and groups keypoints with associative embeddings [30]. In this paper, we establish the benchmark and extensively evaluate the robustness of these state-of-the-art top-down and bottom-up methods.

### 2.2. Corruption Robustness

Recent studies have explored the corruption robustness of image classification [10, 17], object detection [29], and segmentation [23]. In comparison, the task of HPE is more comprehensive, which requires a blend of classification and regression approaches to model human body structure. Data denoising [4] *e.g.* sparse and redundant representations [12], non-local algorithm [3], and denoising auto-encoder [41] are effective in removing noise. However, such methods are noise-specific, thus are not applicable to improving robustness towards unforeseen noises. To improve general robustness, recent works have explored several useful techniques such as pre-training [18], stability

| Clean | Gaussian Noise | Defocus Blur | Snow | Contrast |

Figure 2. **Visualization of examples in our benchmark datasets.** Our proposed benchmarks contain 15 different types of corruptions with different severity levels for a single clean image. The image corruption types are grouped into four main categories: noise, blur, weather, and digital. We sample one corruption type from each category in this figure.

training [49], stylized image [15], NoisyStudent [42], and histogram equalization [17].

## 2.3. Data Augmentation

Data augmentation has been widely utilized as an effective method to improve model generalization. However, improving the general model robustness to unseen image corruptions is difficult. According to [37, 16], augmenting with one specific type of noise enhances the performance on the target noise, but it does not generalize to other unseen distortions while degrading the performance on clean data. Information dropping methods [50, 35, 9, 6] and multi-image mixing methods [46, 44, 19] have gained decent improvements on clean data for image classification. Learned augmentation methods [48, 8] have also been proposed to improve performance. For pose estimation, adversarial data augmentation methods [33, 2] are leveraged to optimize augmentation hyper-parameters, *e.g.* the rotation angle. However, they only focus on searching for common augmentation hyper-parameters and simply combine different augmentations sequentially to improve performance on clean data. Instead, we adversarially learn attention weights for mixing randomly augmented images to enhance model robustness. AugMix [19] is proposed to mix augmented images using beta or dirichlet coefficients to boost the robustness on image classification. Instead of simply using a fixed augmentation sampler(*e.g.* dirichlet distribution) to generate mixing weights, our AdvMix uses adversarial training to learn to generate appropriate mixing attention weights for each training sample.

## 3. Methods

### 3.1. Robust Pose Benchmark

#### 3.1.1 Benchmark Datasets

The robust pose estimation benchmark is composed of three benchmark datasets: COCO-C, MPII-C, OCHuman-C, which are constructed by applying 15 different types of image corruptions [17] with 5 severity levels to the official

*validation* set of COCO [26], MPII [1] and OCHuman [47]. Therefore, the total number of each benchmarking dataset is $15 \times 5$ times that of the corresponding validation dataset. The types of image corruption are sorted into four main categories (noise, blur, weather, and digital), which are diverse and enormous enough to cover real-world corruptions.

**COCO-C Dataset.** COCO-C dataset is constructed from COCO [26] val2017 set. The COCO dataset contains over 200,000 images and 250,000 person instances, where there are 5000 images for the val2017 set.

**MPII-C Dataset.** MPII-C dataset is constructed from MPII [1] test set. The MPII dataset consists of images taken from a wide range of real-world activities with full-body pose annotations. There are around 25K images with 40K subjects, where there are 12K subjects for testing and the remaining subjects for the training set.

**OCHuman-C Dataset.** OCHuman-C is constructed from OCHuman [47]. The OCHuman dataset focuses on heavily occluded human with comprehensive annotations including bounding-box, humans pose and instance mask, which contains 13360 elaborately annotated human instances within 5081 images.

#### 3.1.2 Evaluation Metrics

For COCO [26] and OCHuman [47], the standard average precision (AP) is used to evaluate the models. In this paper, we use $AP^*$ to denote the performance on **clean** data. To evaluate model robustness, we follow [29] to use mean performance under corruption (mPC):

$$mPC = \frac{1}{N_c} \sum_{c=1}^{N_c} \frac{1}{N_s} \sum_{s=1}^{N_s} AP_{c,s}. \qquad (1)$$

Here, $AP_{c,s}$ is the average performance measured on corruption type $c$ under severity level $s$. $N_c = 15$ and $N_s = 5$ are the numbers of corruption types and severity levels, respectively. To assess the robustness of models, we define the relative performance under corruption (rPC) as follows:
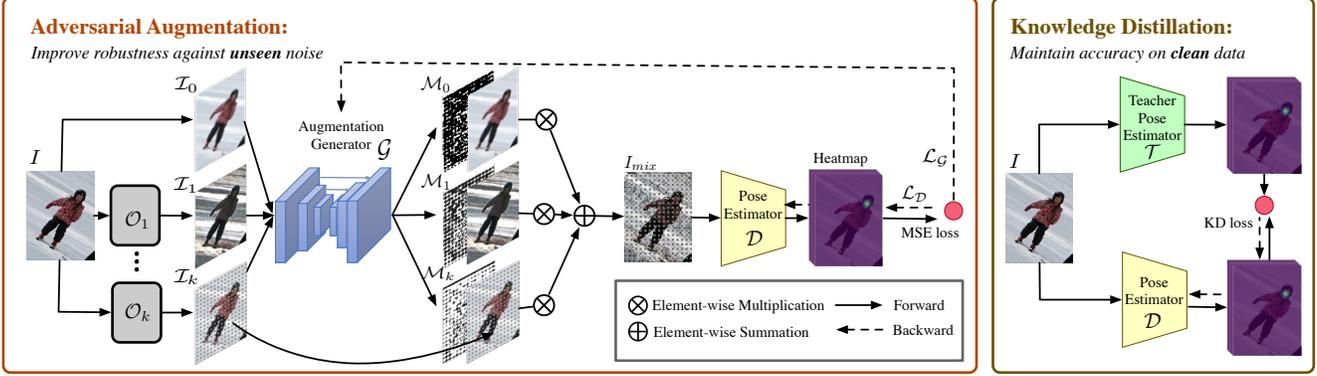
Figure 3. **Overview of AdvMix.** Our framework consists of two modules: the adversarial augmentation module and the knowledge distillation module. The adversarial augmentation module contains an augmentation generator and a pose estimator. For example, given two (*i.e.* $K = 2$) differently augmented images, the augmentation generator estimates the corresponding attention maps $\mathcal{M}$, and mixes them up to the final image $I_{mix}$, while the pose estimator generates keypoint heatmaps. They are trained in an adversarial manner following Eq.(10). The robustness of the pose estimator can be significantly improved if the generator cannot confuse it. To reduce the effect of induced noise patterns, we use a teacher pose estimator for transferring pose structure knowledge in adversarial training. The teacher pose estimator shares the same architecture as the target pose estimator and is pre-trained on the entire clean data.

$$rPC = \frac{mPC}{AP^*}. \tag{2}$$

For MPII [1] dataset, the official evaluation metric is PCKh. Similarly, we use PCKh$^*$ to denote the performance on clean data. Similar to mPC and rPC in Eq.(1,2), mean performance under corruption (mPC) and relative mean performance under corruption (rPC) are introduced as follows:

$$mPC = \frac{1}{N_c} \sum_{c=1}^{N_c} \frac{1}{N_s} \sum_{s=1}^{N_s} PCKh_{c,s} \tag{3}$$

$$rPC = \frac{mPC}{PCKh^*}. \tag{4}$$

### 3.2. Adversarial Augmentation Mix (AdvMix)

#### 3.2.1 Augmentation Generator

As shown in Figure 3, given an image $I$, we first randomly generate $K$ differently augmented images $\mathcal{O}_k(I)$ using parallel augmentation strategies $\mathcal{O}_k$. We use $K = 2$ by default in our implementation. Together with the original image $I$, we get a set of $K + 1$ proposal images. The augmentation generator $\mathcal{G}(\cdot, \theta)$ is applied to output the normalized attention maps $\mathcal{M}^{(K+1) \times H \times W}$, where $H$ and $W$ are the image resolution. The attention maps $\mathcal{M}$ are used as the weights to mix up the proposal images, following Eq.(5), where $\odot$ is the Hadamard product.

$$I_{mix} = \mathcal{M}_0 \odot I + \sum_{k=1}^{K} \mathcal{M}_k \odot \mathcal{O}_k(I). \tag{5}$$

Based on the benchmark (Table.1), we examine and find a collection of useful techniques to improve the model ro-

bustness for pose estimation. Grid-Mask [6] applies information removal to achieve state-of-the-art results in various computer vision tasks. AutoAugment [8] automatically searches for a mixture of augmentation policies. It improves both the clean performance and the robustness for image classification [43]. To avoid over-fitting to the test set of held-out corruptions, we manually exclude some operations, such as contrast, color, brightness, and sharpness sub-policies (as they appear in the benchmark). Given an image $I$, we choose Grid-Mask and AutoAugment to generate randomly augmented images $\mathcal{O}$, and then mix these images through mixing weight from augmentation generator. We adopt the U-Net [34] architecture to build the augmentation generator $\mathcal{G}(\cdot, \theta)$, which generates the attention maps for mixing up randomly augmented images. Please refer to the supplementary for more implementation details.

#### 3.2.2 Heatmap Regression

The human body pose is encoded with 2D Gaussian confidence heatmaps, where each channel corresponds to one body keypoint. The objective of human pose estimation network $\mathcal{D}(\cdot, \phi)$ is to minimize the MSE loss between the predicted heatmaps and the ground truth heatmaps $\mathcal{H}_{gt}$:

$$\mathcal{L}_{\mathcal{D}^*} = \|\mathcal{D}(I_{mix}, \phi) - \mathcal{H}_{gt}\|_2^2. \tag{6}$$

To enhance the supervision to the target pose estimator $\mathcal{D}$, we introduce a teacher pose estimator $\mathcal{T}$ for knowledge distillation and providing softened heatmap labels. Note that $\mathcal{T}$ and $\mathcal{D}$ share the same architecture and $\mathcal{T}$ is pre-trained on the clean training data. The parameters of the $\mathcal{T}$ will be fixed while training. We select MSE loss for knowledge distillation loss as:

$$\mathcal{L}_{\mathcal{D}_{kd}} = \|\mathcal{D}(I, \phi) - \mathcal{T}(I)\|_2^2. \qquad (7)$$

We formulate the overall loss function of pose estimator $\mathcal{D}$ while training as:

$$\mathcal{L}_{\mathcal{D}} = (1 - \alpha)\mathcal{L}_{\mathcal{D}^*} + \alpha\mathcal{L}_{\mathcal{D}_{kd}}. \qquad (8)$$

where $\alpha$ is the loss weight balancing between MSE loss and knowledge distillation loss. We observe that our model is not sensitive to $\alpha$. Thus we choose $\alpha = 0.1$ as the default setting in our experiments.

### 3.2.3 Adversarial Training

The augmentation generator $\mathcal{G}(\cdot, \theta)$ and human pose estimator $\mathcal{D}(\cdot, \phi)$ are trained in an adversarial manner. $\mathcal{G}(\cdot, \theta)$ tries to find the most confusing way to mix up randomly augmented images, while $\mathcal{D}(\cdot, \phi)$ learns more robust features from harder training samples. The optimization objective of augmentation generator is defined as:

$$\mathcal{L}_{\mathcal{G}} = -\mathcal{L}_{\mathcal{D}^*}. \qquad (9)$$

Overall, the whole learning process can be defined as two-player zero-sum game with value function $\mathcal{V}(\mathcal{D}, \mathcal{G})$.

$$\mathcal{V}(\mathcal{D}, \mathcal{G}) = \min_{\phi} \max_{\theta} \mathop{\mathbb{E}}_{\mathcal{I} \sim \Omega} \mathcal{L}(\mathcal{D}(\mathcal{G}(\mathcal{O}(\mathcal{I}), \theta), \phi), \mathcal{H}_{gt}). \qquad (10)$$

## 4. Evaluation on Robust Pose Benchmark

### 4.1. Experimental Setup

We extensively evaluate the performance of the state-of-the-art methods on the proposed Robust Human Pose Benchmark. To assess the robustness of human pose estimators, the models are only trained on clean data (*e.g.* COCO) and then evaluated on corrupted data (*e.g.* COCO-C). Note that OCHuman dataset is only designed for validation, we follow the common settings [47] to train the models on COCO and evaluate on OCHuman-C. For top-down methods (*i.e.* **SBL** [40], **HRNet** [36]) with input size $256 \times 192$ and $384 \times 288$, we use the officially trained checkpoints [1], and then directly evaluate their robustness performance on on COCO-C and MPII-C datasets. Since the pre-trained models with image size $128 \times 96$ are not publicly available, we retrain these models following the official training settings. For bottom-up methods (*i.e.* **PifPaf** [24], **HigherHR-Net** [7]), we also follow the official codes and settings. The pre-trained models of PifPaf [2] and HigherHRNet [3] are directly used for evaluating the robustness using on COCO-C and OCHuman-C. The results are reported in Table 1.

---

[1]https://github.com/leoxiaobin/deep-high-resolution-net.pytorch
[2]https://github.com/vita-epfl/openpifpaf/tree/v0.10.0
[3]https://github.com/HRNet/HigherHRNet-Human-Pose-Estimation

### 4.2. Benchmarking Conclusions

#### 4.2.1 Pose estimation methods performance

**Top-down and bottom-up models show similar performance degradation tendency on different corruptions.** As shown in Figure 4, the model robustness to different types of corruption varies a lot. However, performance degradation across different models are similar. For example, all models are more robust to weather or brightness changes, while less robust to motion or zoom blur.

#### 4.2.2 Backbone network performance

**Robustness increases with model capacity.** With the same pose estimation methods, the robustness of different backbone networks varies a lot and increases with model capacity. In Figure 5, the line graph shows the performance on clean data (AP*), while the bar graph shows the robustness score (rPC). The color denotes different input image resolutions. We observe that with the same input resolution, a model with higher capacity (HRNet) is generally more accurate and more robust than that with lower capacity (ResNet). This indicates that robustness can be improved by using a stronger backbone network.

## 5. Robustness Enhancement with AdvMix

### 5.1. Implementation Details

For adversarial training, we simply reuse the existing publicly available well-trained checkpoints and retrain the models with AdvMix on clean data. We simply follow the same training settings (*i.e.* the same hyper-parameters, learning rate, and total training epochs) as the official codes. By default, the initial learning rates for the augmentation generator and human pose estimator is 0.001. We decay the learning rate by the factor of 10 at the 170-th epoch, and 200-th epoch. The adversarial training ends at 210-th epoch. We adopt ADAM optimizer to train both the augmentation generator and the human pose estimator. All experiments are conducted using PyTorch [32] on NVIDIA TITAN X Pascal GPUs.

### 5.2. Quantitative Results

As shown in Table 2 and Table 3, we find that the proposed method significantly improves the model robustness while the performance on clean data is maintained or slightly improved. We also observe that AdvMix significantly boosts the robustness performance for bottom-up methods (39.9 to 45.4 mPC) and the robustness gain is larger than that of top-down methods shown in Table 2. Figure 4(a) shows the performance degradation and the robustness improvement of different models across different corruptions. Performance degradation among different models

Table 1. Pose robustness benchmark for both top-down and bottom-up models on COCO-C, MPII-C and OChuman-C. AP* and PCKh* represent the performance on clean data. mPC represents mean performance under all corruptions, while rPC measures the relative performance. The remaining columns are the breakdown APs for different corruptions. We see that performances of existing advanced pose estimators significantly drop when corruptions are presented.

**Results of top-down methods on COCO-C with the same detection bounding boxes as [36]**

| Method | Backbone | Input size | AP* | mPC | rPC | Gauss | Shot | Impulse | Defoucs | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SBL [40] | ResNet-50 | 128 × 96 | 59.0 | 40.7 | 69.0 | 37.8 | 39.8 | 36.9 | 39.5 | 38.1 | 35.3 | 13.3 | 39.1 | 43.0 | 48.1 | 54.5 | 38.5 | 49.2 | 50.9 | 47.1 |
| | | 256 × 192 | 70.4 | 47.8 | 67.9 | 45.8 | 48.1 | 45.6 | 43.4 | 42.1 | 38.8 | 16.3 | 49.1 | 52.5 | 58.9 | 65.5 | 47.2 | 56.7 | 55.2 | 51.7 |
| | | 384 × 288 | 72.2 | 47.7 | 66.1 | 45.2 | 47.8 | 45.9 | 42.9 | 42.1 | 38.1 | 16.3 | 49.9 | 53.2 | 60.0 | 66.8 | 48.2 | 57.3 | 53.2 | 49.2 |
| | ResNet-101 | 128 × 96 | 61.1 | 42.6 | 69.7 | 40.4 | 42.6 | 39.6 | 40.7 | 39.9 | 36.8 | 14.5 | 41.3 | 45.2 | 49.6 | 56.8 | 38.6 | 51.1 | 53.2 | 49.1 |
| | | 256 × 192 | 71.4 | 49.6 | 69.5 | 47.8 | 50.1 | 47.2 | 45.1 | 43.8 | 40.3 | 17.6 | 50.9 | 54.9 | 60.9 | 67.0 | 49.7 | 58.1 | 57.0 | 53.8 |
| | | 384 × 288 | 73.6 | 50.4 | 68.4 | 49.2 | 51.7 | 49.1 | 44.8 | 43.9 | 40.0 | 17.7 | 52.5 | 56.4 | 62.7 | 69.0 | 51.2 | 59.0 | 55.5 | 52.7 |
| HRNet [36] | HRNet-W32 | 128 × 96 | 66.9 | 47.2 | 70.6 | 42.7 | 45.7 | 43.0 | 44.0 | 43.1 | 40.9 | 15.9 | 47.1 | 51.4 | 57.4 | 63.0 | 47.5 | 55.7 | 57.2 | 53.4 |
| | | 256 × 192 | 74.4 | 53.0 | 71.3 | 51.3 | 54.2 | 52.6 | 46.9 | 46.3 | 43.5 | 19.2 | 55.9 | 59.1 | 65.2 | 70.3 | 54.1 | 60.5 | 59.4 | 56.9 |
| | | 384 × 288 | 75.7 | 53.7 | 70.9 | 51.9 | 54.7 | 53.7 | 47.8 | 47.1 | 43.8 | 19.8 | 57.9 | 60.3 | 66.5 | 71.6 | 55.4 | 61.1 | 58.1 | 55.8 |
| | HRNet-W48 | 128 × 96 | 68.6 | 49.3 | 71.8 | 45.3 | 44.4 | 42.4 | 16.5 | 49.6 | 54.1 | 59.8 | 65.0 | 49.4 | 57.3 | 59.2 | 55.2 |
| | | 256 × 192 | 75.1 | 53.7 | 71.6 | 52.5 | 55.2 | 53.4 | 46.8 | 46.7 | 43.5 | 19.1 | 57.0 | 60.1 | 66.4 | 71.4 | 55.2 | 61.1 | 60.0 | 57.6 |
| | | 384 × 288 | 76.3 | 54.2 | 71.1 | 52.8 | 55.8 | 54.2 | 47.6 | 47.3 | 43.4 | 19.5 | 58.3 | 60.9 | 67.5 | 72.3 | 56.3 | 61.6 | 59.2 | 57.1 |

**Results of top-down methods on MPII-C**

| Method | Backbone | Input size | PCKh* | mPC | rPC | Gauss | Shot | Impulse | Defoucs | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SBL [40] | Resnet-50 | 256 × 256 | 88.5 | 77.5 | 87.6 | 68.9 | 71.6 | 68.9 | 84.2 | 85.3 | 83.6 | 53.7 | 70.0 | 75.6 | 83.1 | 86.6 | 69.1 | 87.8 | 87.8 | 86.8 |
| | Resnet-101 | 256 × 256 | 89.1 | 78.6 | 88.3 | 70.4 | 73.1 | 70.6 | 84.8 | 86.0 | 84.1 | 53.9 | 72.2 | 76.9 | 84.1 | 87.0 | 71.9 | 88.3 | 88.4 | 87.3 |
| | Resnet-152 | 256 × 256 | 89.6 | 79.6 | 88.8 | 74.1 | 76.3 | 74.1 | 85.3 | 86.5 | 84.8 | 54.5 | 72.1 | 77.6 | 84.4 | 87.7 | 71.4 | 88.8 | 88.8 | 87.9 |
| HRNet [36] | HRNet-W32 | 256 × 256 | 90.3 | 80.1 | 88.7 | 71.5 | 74.1 | 72.4 | 86.0 | 87.0 | 85.3 | 56.0 | 73.8 | 78.6 | 86.1 | 88.5 | 74.5 | 89.5 | 89.5 | 88.7 |

**Results of bottom-up methods on COCO-C**

| Method | Backbone | Input size | AP* | mPC | rPC | Gauss | Shot | Impulse | Defoucs | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PifPaf [24] | ShuffleNet V2 | 641 × 641 | 60.7 | 32.9 | 54.2 | 25.7 | 27.3 | 24.6 | 29.1 | 28.5 | 26.4 | 8.8 | 33.5 | 38.5 | 47.6 | 54.3 | 34.1 | 45.7 | 33.4 | 35.7 |
| | Resnet-50 | 641 × 641 | 64.8 | 34.4 | 53.1 | 30.4 | 32.2 | 29.5 | 30.9 | 27.9 | 26.5 | 8.8 | 36.4 | 40.2 | 49.9 | 57.3 | 34.9 | 47.8 | 29.9 | 34.1 |
| | Resnet-101 | 641 × 641 | 68.3 | 40.6 | 59.5 | 38.0 | 39.4 | 36.6 | 36.9 | 34.6 | 31.6 | 11.2 | 42.6 | 46.6 | 56.0 | 62.0 | 42.4 | 52.6 | 37.2 | 41.3 |
| HrHRNet [7] | HrHRNet-W32 | 512 × 512 | 67.1 | 39.9 | 59.4 | 34.2 | 37.0 | 35.2 | 35.1 | 32.5 | 34.0 | 12.5 | 43.3 | 47.6 | 54.9 | 60.6 | 43.4 | 50.3 | 42.2 | 35.0 |
| | | 640 × 640 | 68.5 | 39.6 | 57.8 | 31.5 | 38.9 | 37.6 | 34.8 | 32.6 | 33.8 | 12.1 | 44.0 | 47.9 | 55.1 | 61.4 | 42.9 | 50.6 | 37.8 | 33.2 |
| | HrHRNet-W48 | 512 × 512 | 68.5 | 41.9 | 61.2 | 39.3 | 42.2 | 40.1 | 36.1 | 33.5 | 35.1 | 13.1 | 45.0 | 49.1 | 56.7 | 62.1 | 45.5 | 51.4 | 42.2 | 36.9 |
| | | 640 × 640 | 69.8 | 40.8 | 58.4 | 36.7 | 39.8 | 37.9 | 35.6 | 33.1 | 34.2 | 12.5 | 44.7 | 49.0 | 56.3 | 62.8 | 43.7 | 51.3 | 39.9 | 34.3 |

**Results of bottom-up methods on OChuman-C**

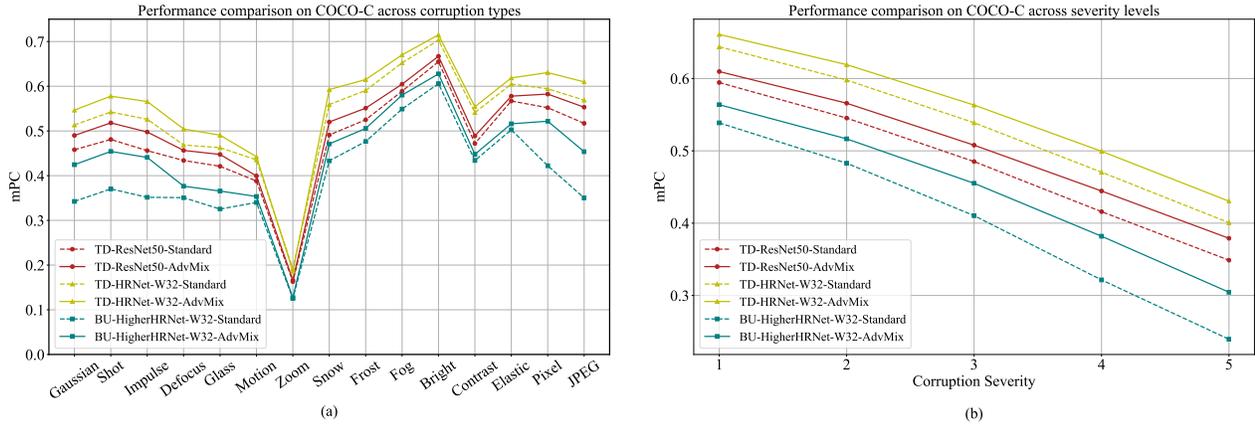| Method | Backbone | Input size | AP* | mPC | rPC | Gauss | Shot | Impulse | Defoucs | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PifPaf [24] | ShuffleNet V2 | 641 × 641 | 37.8 | 30.8 | 81.5 | 28.7 | 28.8 | 28.1 | 32.7 | 32.2 | 32.9 | 17.2 | 26.3 | 30.5 | 36.4 | 36.1 | 30.5 | 36.2 | 32.8 | 32.6 |
| | Resnet-50 | 641 × 641 | 37.6 | 30.8 | 82.0 | 30.1 | 30.0 | 29.4 | 33.3 | 31.6 | 32.4 | 16.5 | 27.0 | 30.5 | 36.5 | 36.5 | 29.9 | 36.1 | 31.8 | 30.7 |
| | Resnet-101 | 641 × 641 | 39.6 | 34.9 | 88.1 | 35.6 | 35.6 | 35.3 | 37.2 | 35.3 | 36.3 | 20.2 | 31.4 | 34.4 | 39.2 | 38.8 | 34.8 | 38.9 | 35.5 | 34.6 |
| HrHRNet [7] | HrHRNet-W32 | 512 × 512 | 40.0 | 35.1 | 87.6 | 32.9 | 33.1 | 33.2 | 37.3 | 36.5 | 36.9 | 21.7 | 31.6 | 36.0 | 39.9 | 39.4 | 36.9 | 39.2 | 37.1 | 34.4 |
| | | 640 × 640 | 39.3 | 33.6 | 85.4 | 30.9 | 31.1 | 31.8 | 36.2 | 35.0 | 36.1 | 19.6 | 31.1 | 35.1 | 39.3 | 38.1 | 35.8 | 38.3 | 33.8 | 31.1 |
| | HrHRNet-W48 | 512 × 512 | 41.7 | 36.7 | 87.9 | 35.5 | 35.6 | 35.7 | 38.8 | 37.5 | 38.5 | 22.0 | 33.5 | 37.6 | 41.5 | 40.5 | 39.0 | 40.6 | 38.0 | 35.6 |
| | | 640 × 640 | 40.9 | 35.4 | 86.4 | 33.0 | 33.3 | 33.5 | 37.7 | 36.6 | 37.6 | 21.1 | 32.4 | 37.4 | 41.5 | 39.9 | 37.2 | 39.9 | 36.6 | 32.6 |



Figure 4. (a) Performance improvements of AdvMix on COCO-C across different corruption types. (b) Performance improvements of AdvMix on COCO-C across different corruption severities. The results are obtained with input sizes as 256 × 192. Input sizes for top-down (TD) methods are 256 × 192, while for bottom-up (BU) methods (*e.g.* HigherHRNet) is 512 × 512.

shows a similar tendency. Models are generally more robust to brightness and weather changes, while less robust to motion and zoom blur. AdvMix performs better than baseline by different corruption types and the gain of noise and digital distortion are larger than other corruption types. Comparing the performance across different corruption severity

in Figure 4(b), we observe that AdvMix consistently improves over the baseline, and the gain gets larger for severer corruption.

We compare AdvMix with other state-of-the-art methods. FPD [45] proposes to improve the performance of lightweight models through knowledge distillation, and the
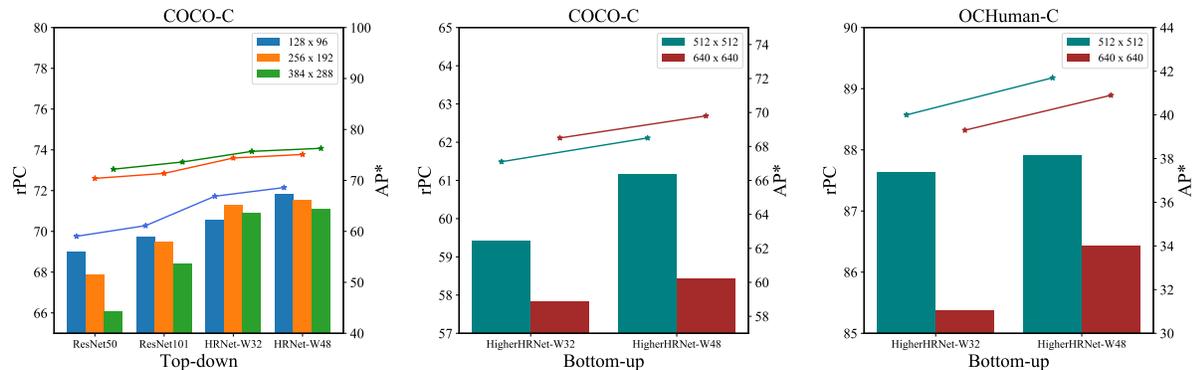
Figure 5. Performance of top-down and bottom-up models with different backbones and input sizes on COCO-C and OCHuman-C. The line graph shows performance on clean data (AP*), and the bar graph denotes robustness score (rPC).

Table 2. **Comparisons** between standard training and AdvMix on COCO-C. For top-down approaches, results are obtained with detected bounding boxes of [36]. We see that mPC and rPC are greatly improved, whilst clean performance AP* can be preserved.

| Method | Backbone | Input size | AP* | mPC | rPC |
|--------|----------|------------|-----|-----|-----|
| Standard | ResNet-50 | $256 \times 192$ | **70.4** | 47.8 | 67.9 |
| **AdvMix** | ResNet-50 | $256 \times 192$ | 70.1 | **50.1** | **71.5** |
| Standard | ResNet-101 | $256 \times 192$ | **71.4** | 49.6 | 69.5 |
| **AdvMix** | ResNet-101 | $256 \times 192$ | 71.3 | **52.3** | **73.3** |
| Standard | ResNet-152 | $256 \times 192$ | 72.0 | 50.9 | 70.7 |
| **AdvMix** | ResNet-152 | $256 \times 192$ | **72.3** | **53.2** | **73.6** |
| Standard | HRNet-W32 | $256 \times 192$ | 74.4 | 53.0 | 71.3 |
| **AdvMix** | HRNet-W32 | $256 \times 192$ | **74.7** | **55.5** | **74.3** |
| Standard | HRNet-W48 | $256 \times 192$ | 75.1 | 53.7 | 71.6 |
| **AdvMix** | HRNet-W48 | $256 \times 192$ | **75.4** | **57.1** | **75.7** |
| Standard | HrHRNet-W32 | $512 \times 512$ | 67.1 | 39.9 | 59.4 |
| **AdvMix** | HrHRNet-W32 | $512 \times 512$ | **68.3** | **45.4** | **66.5** |

Table 3. **Comparisons** between standard training and AdvMix of representative top-down methods on MPII-C. We draw similar conclusion that AdvMix is effective to increase mPC and rPC, whilst maintaining clean performance PCKh*.

| Method | Backbone | Input size | PCKh* | mPC | rPC |
|--------|----------|------------|-------|-----|-----|
| Standard | ResNet-50 | $256 \times 256$ | 88.5 | 77.5 | 87.6 |
| **AdvMix** | ResNet-50 | $256 \times 256$ | **88.9** | **82.0** | **92.3** |
| Standard | ResNet-101 | $256 \times 256$ | 89.1 | 78.6 | 88.3 |
| **AdvMix** | ResNet-101 | $256 \times 256$ | **89.4** | **82.8** | **92.5** |
| Standard | HRNet-W32 | $256 \times 256$ | 90.3 | 80.1 | 88.7 |
| **AdvMix** | HRNet-W32 | $256 \times 256$ | **90.5** | **83.9** | **92.7** |

teacher net is more sophisticated than the student net. FPD† is our implementation with the same setting as FPD except that the teacher net uses the same architecture as the student net. Stylized dataset [14] is also used as a technique to improve robustness. As shown in Table 5, we can see that though FPD can boost the performance on clean data, the improvement of robustness is limited. We also find that only using knowledge distillation (teacher use the same architecture in FPD†) without AdvMix will not improve the performance on clean data and model robustness. AdvMix can not only improve the mPC, but also slightly improve the performance on clean data, outperforming op-

erating AutoAugment or Grid-Mask separately. Finally, we also prove that AdvMix can be combined with the data stylizing technique and further enhance the robustness on corrupted images while almost maintaining the performance of clean data.

## 5.3. Ablation Studies

We conduct ablation studies on COCO-C dataset using HRNet-W32 backbone with input size $256 \times 192$ to validate the effectiveness of the augmentation composition method in AdvMix. SequentialMix simply composes augmentation operations in a chain and applies them sequentially. EqualMix mixes the augmented images from various augmentation chains using equal weights. DirichletMix follows [19] to sample the mix weights from Dirichlet distribution to mix the augmented images. In comparison, AdvMix uses the augmentation generator to *learn* to generate the per-pixel mix weights adversarially. Different from AdvMix, the augmentation generator of AdvMix-Image outputs per-image mix weights rather than per-pixel mix weights. We observe that AdvMix significantly outperforms EqualMix and DirichletMix on model robustness (mPC and rPC) by a large margin, demonstrating the effectiveness of the proposed adversarial weights learning. Meanwhile, AdvMix with learnt per-pixel composition is more robust than learnt per-image composition (AdvMix-Image) since per-pixel composition could generate more diverse and fine-grained training samples to boost robustness performance, *i.e.*, per-image mixing can hardly composite image in region-level, *e.g.*, person and background, while per-pixel mixing with learnt weights can pay more attention to important regions. As illustrated in Table 7, knowledge distillation prevents over-fitting to the induced noise patterns from creeping into the feature space and helps maintain or improve the performance on clean data.

Since we focus on how to improve robustness on **unseen** data, the corruption images in the benchmark should *not* be encountered while training. However, to investigate the

Table 4. Results of directly augmenting with the test-time corruption types (Transfer). We use † to denote the selected corruption types (*e.g.* Gaussian noise, Defocus blur, Snow, and Contrast) are used as data augmentation during training.

| Method | AP* | mPC | Gauss | Shot | Impulse | Defoucs | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Standard | 76.6 | 53.3 | 51.1 | 54.1 | 52.5 | 46.7 | 46.2 | 43.4 | 18.8 | 56.2 | 59.7 | 66.2 | 71.8 | 54.5 | 61.1 | 59.8 | 57.1 |
| Transfer | 75.4 | 56.2 | 58.7† | 60.3 | 59.7 | 55.0† | 48.3 | 44.3 | 19.0 | 65.6† | 61.7 | 65.8 | 70.9 | 58.9† | 60.3 | 61.8 | 52.2 |
| AdvMix | **77.1** | 55.9 | 54.7 | 57.9 | 56.7 | 50.4 | 49.1 | 44.2 | 18.2 | 59.8 | 62.4 | 68.2 | 73.2 | 56.0 | 62.6 | 63.4 | 61.5 |
| AdvMix+Stylized | 76.5 | **56.5** | 54.5 | 57.5 | 56.4 | 50.4 | 49.3 | 46.1 | 20.0 | 62.3 | 63.0 | 67.6 | 72.9 | 59.0 | 62.5 | 63.9 | 61.8 |

Table 5. **Comparisons** with other techniques on COCO-C with HRNet-W32 backbone. The results are obtained using ground truth bounding boxes. We can observe that AdvMix significantly outperforms the state-of-the-art methods on robustness, while preserving the performance of clean AP.

| Method | AP* | mPC | rPC |
|---|---|---|---|
| Standard | 76.6 | 53.3 | 69.6 |
| FPD [45] | 77.3 | 54.0 | 69.9 |
| FPD† [45] | 76.6 | 53.4 | 69.7 |
| Grid-Mask [6] | 76.4 | 54.8 | 71.7 |
| AutoAugment [8] | 76.2 | 54.2 | 71.1 |
| Stylized Only [17, 29] | 67.5 | 46.7 | 69.1 |
| Stylized [17, 29] | 76.1 | 54.8 | 72.0 |
| **AdvMix** | **77.1** | 55.9 | 72.5 |
| Stylized + **AdvMix** | 76.5 | **56.5** | **73.8** |

Table 6. **Ablation** study on COCO-C with HRNet-W32 backbone. The results are obtained using ground truth bounding boxes. We see that AdvMix performs better than other augmentation composing approaches on robustness by large margin.

| Method | AP* | mPC | rPC |
|---|---|---|---|
| SequentialMix | 76.2 | 54.6 | 71.6 |
| EqualMix | 76.6 | 54.2 | 70.7 |
| DirichletMix [19] | 76.5 | 54.4 | 71.1 |
| **AdvMix**-Image | **77.5** | 55.1 | 71.2 |
| **AdvMix** | 77.1 | **55.9** | **72.5** |

Table 7. Effect of knowledge distillation.

| Method | dataset | AP*/ PCKh* | mPC | rPC |
|---|---|---|---|---|
| **AdvMix** w/o KD | COCO-C | 76.8 | 55.9 | 72.7 |
| **AdvMix** | COCO-C | **77.1** | 55.9 | 72.5 |
| **AdvMix** w/o KD | MPII-C | 89.9 | 81.7 | 90.9 |
| **AdvMix** | MPII-C | **90.5** | **83.9** | **92.7** |

phenomenon if we augment the images with the test-time corruption type in the benchmark, we select one corruption type (*i.e.* Gaussian noise, Defocus blur, Snow and Contrast) from each corruption category (*i.e.* noise, blur, weather, and digital) as augmentation types. As illustrated in Table 4, we can observe that augmenting the training samples with test-time corruptions can boost the robustness, but the performance on clean data decreases a lot. Meanwhile, augmenting with some specific types of noises improves the performance on the target noises, but it does not always generalize to other unseen corruption types, even within the same corruption category (*e.g.* augmentation of Snow does not contribute to improving the robustness of Fog and Brightness). By contrast, AdvMix training with stylized data achieves the best mPC. It consistently improves the performance across all corruption types, while maintaining similar clean data performance as standard training.

## 5.4. Qualitative Comparison

In Figure 8, we visualize the results of images with different types of image corruptions, *i.e.* impulse noise, Motion blur, Brightness, and JPEG compression. We observe that 1) the standard pose models suffer a large performance drop on corrupted data, and 2) models trained with AdvMix perform consistently better than the baseline methods on various corruptions.



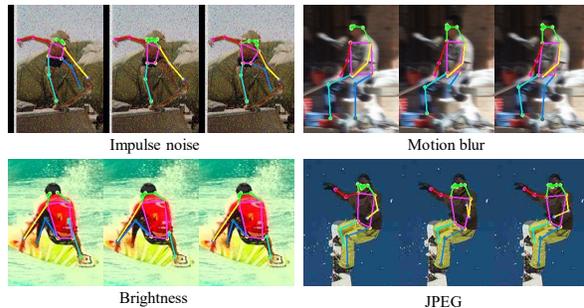Impulse noise      Motion blur

Brightness      JPEG

Figure 6. Qualitative comparisons between HRNet without and with AdvMix. For each image triplet, the images from left to right are ground truth, predicted results of Standard HRNet-W32, and predicted results of HRNet-W32 with AdvMix.

## 6. Conclusion

In this paper, we propose the Robust Pose Benchmark (COCO-C, MPII-C, OCHuman-C) and rigorously evaluate the performance of current state-of-the-art models on corrupted images. Based on the benchmark, we have some interesting conclusions and examine some useful techniques to improve model robustness. We envision this work will draw the community's attention to this challenging problem and promote the development of robust pose estimators. To improve the model robustness, we propose AdvMix, a novel model-agnostic data augmentation method, to learn to mix up randomly augmented images. Our method significantly improves the robustness of most existing pose estimation models across a wide range of common corruptions while maintaining performance on clean data without extra inference computational overhead.

# 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 (CVPR)*, 2014. 1, 3, 4

[2] Yanrui Bin, Xuan Cao, Xinya Chen, Yanhao Ge, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, Changxin Gao, and Nong Sang. Adversarial semantic data augmentation for human pose estimation. *arXiv preprint arXiv:2008.00697*, 2020. 3

[3] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 2, pages 60–65. IEEE, 2005. 2

[4] Antoni Buades, Bartomeu Coll, and Jean-Michel Morel. A review of image denoising algorithms, with a new one. *Multiscale Modeling & Simulation*, 4(2):490–530, 2005. 2

[5] 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 (CVPR)*, 2017. 1, 2

[6] Pengguang Chen. Gridmask data augmentation. *arXiv preprint arXiv:2001.04086*, 2020. 3, 4, 8

[7] 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. 1, 2, 5, 6

[8] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018. 3, 4, 8

[9] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. 3

[10] Samuel Dodge and Lina Karam. Understanding how image quality affects deep neural networks. In *2016 eighth international conference on quality of multimedia experience (QoMEX)*, pages 1–6. IEEE, 2016. 2

[11] Haodong Duan, Kwan-Yee Lin, Sheng Jin, Wentao Liu, Chen Qian, and Wanli Ouyang. Trb: a novel triplet representation for understanding 2d human body. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9479–9488, 2019. 2

[12] Michael Elad and Michal Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image processing*, 15(12):3736–3745, 2006. 2

[13] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. Rmpe: Regional multi-person pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2334–2343, 2017. 2

[14] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 7

[15] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-Trained Cnns Are Biased Towards Texture; Increasing Shape Bias Improves Accuracy and Robustness. . *ICLR*, 2019. 3

[16] Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A Wichmann. Generalisation in humans and deep neural networks. In *Advances in neural information processing systems*, pages 7538–7550, 2018. 3

[17] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. 1, 2, 3, 8

[18] Dan Hendrycks, Kimin Lee, and Mantas Mazeika. Using pre-training can improve model robustness and uncertainty. *arXiv preprint arXiv:1901.09960*, 2019. 2

[19] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781*, 2019. 3, 7, 8

[20] Sheng Jin, Wentao Liu, Wanli Ouyang, and Chen Qian. Multi-person articulated tracking with spatial and temporal embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5664–5673, 2019. 2

[21] Sheng Jin, Wentao Liu, Enze Xie, Wenhai Wang, Chen Qian, Wanli Ouyang, and Ping Luo. Differentiable hierarchical graph grouping for multi-person pose estimation. In *European Conference on Computer Vision*, pages 718–734. Springer, 2020. 2

[22] Sheng Jin, Lumin Xu, Jin Xu, Can Wang, Wentao Liu, Chen Qian, Wanli Ouyang, and Ping Luo. Whole-body human pose estimation in the wild. In *European Conference on Computer Vision*, pages 196–214. Springer, 2020. 2

[23] Christoph Kamann and Carsten Rother. Benchmarking the robustness of semantic segmentation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8828–8838, 2020. 1, 2

[24] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. Pifpaf: Composite fields for human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11977–11986, 2019. 2, 5, 6

[25] Dangwei Li, Xiaotang Chen, Zhang Zhang, and Kaiqi Huang. Pose guided deep model for pedestrian attribute recognition in surveillance scenarios. In *2018 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2018. 1

[26] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014. 1, 2, 3

[27] Wentao Liu, Jie Chen, Cheng Li, Chen Qian, Xiao Chu, and Xiaolin Hu. A cascaded inception of inception network with attention modulated feature fusion for human pose estima-

tion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018. 2

[28] Diogo C Luvizon, David Picard, and Hedi Tabia. 2d/3d pose estimation and action recognition using multitask deep learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5137–5146, 2018. 1

[29] Claudio Michaelis, Benjamin Mitzkus, Robert Geirhos, Evgenia Rusak, Oliver Bringmann, Alexander S Ecker, Matthias Bethge, and Wieland Brendel. Benchmarking robustness in object detection: Autonomous driving when winter is coming. *arXiv preprint arXiv:1907.07484*, 2019. 1, 2, 3, 8

[30] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint detection and grouping. In *Advances in Neural Information Processing Systems*, 2017. 1, 2

[31] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. 2

[32] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 5

[33] Xi Peng, Zhiqiang Tang, Fei Yang, Rogerio S Feris, and Dimitris Metaxas. Jointly optimize data augmentation and network training: Adversarial data augmentation in human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2226–2234, 2018. 3

[34] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 4, 11

[35] Krishna Kumar Singh, Hao Yu, Aron Sarmasi, Gautam Pradeep, and Yong Jae Lee. Hide-and-seek: A data augmentation technique for weakly-supervised localization and beyond. *arXiv preprint arXiv:1811.02545*, 2018. 3

[36] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. *arXiv preprint arXiv:1902.09212*, 2019. 1, 2, 5, 6, 7, 11

[37] Igor Vasiljevic, Ayan Chakrabarti, and Gregory Shakhnarovich. Examining the impact of blur on recognition by convolutional networks. *arXiv preprint arXiv:1611.05760*, 2016. 3

[38] Chunyu Wang, Yizhou Wang, and Alan L Yuille. An approach to pose-based action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 915–922, 2013. 1

[39] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2

[40] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 1, 2, 5, 6

[41] Junyuan Xie, Linli Xu, and Enhong Chen. Image denoising and inpainting with deep neural networks. In *Advances in neural information processing systems*, pages 341–349, 2012. 2

[42] Qizhe Xie, Minh Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3

[43] Dong Yin, Raphael Gontijo Lopes, Jon Shlens, Ekin Dogus Cubuk, and Justin Gilmer. A fourier perspective on model robustness in computer vision. In *Advances in Neural Information Processing Systems*, pages 13276–13286, 2019. 4

[44] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6023–6032, 2019. 3

[45] Feng Zhang, Xiatian Zhu, and Mao Ye. Fast human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3517–3526, 2019. 6, 8

[46] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 3

[47] Song-Hai Zhang, Ruilong Li, Xin Dong, Paul Rosin, Zixi Cai, Xi Han, Dingcheng Yang, Haozhi Huang, and Shi-Min Hu. Pose2seg: Detection free human instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 889–898, 2019. 1, 3, 5

[48] Xinyu Zhang, Qiang Wang, Jian Zhang, and Zhao Zhong. Adversarial autoaugment. *arXiv preprint arXiv:1912.11188*, 2019. 3

[49] Stephan Zheng, Yang Song, Thomas Leung, and Ian Goodfellow. Improving the robustness of deep neural networks via stability training. In *Proceedings of the ieee conference on computer vision and pattern recognition*, pages 4480–4488, 2016. 3

[50] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, pages 13001–13008, 2020. 3

# Appendix

## A. Architecture of Augmentation Generator

We adopt the U-Net [34] architecture to build the augmentation generator, which generates the attention maps for mixing up randomly augmented images. As shown in Figure 7, the augmentation generator is an encoder-decoder with skip connections in between layers, which consists of 6 convolution blocks and 6 transposed convolution blocks. To make sure the size of down-sampled features are equal to up-sampled features for concatenation, we only utilize 5 convolution blocks and 5 transposed convolution blocks for models of input size $384 \times 288$ and $128 \times 96$.
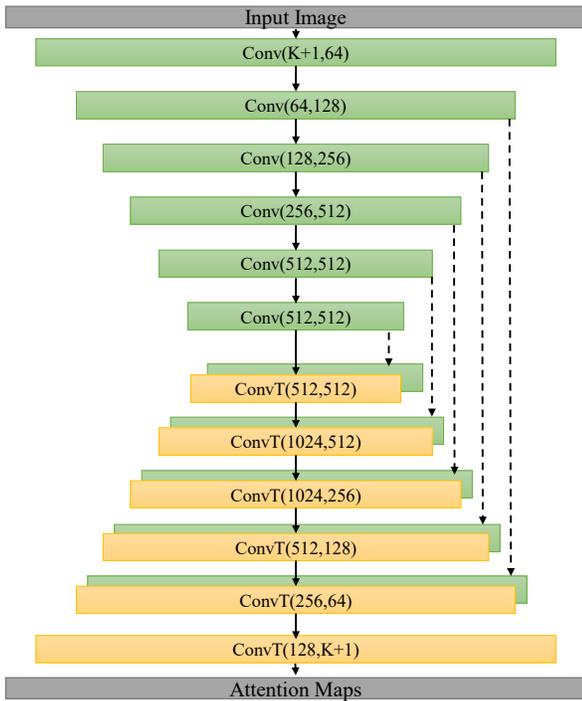


Figure 7. Architecture of Augmentation Generator with input size $256 \times 192$. Conv(c,k) means the convolution block with the output channel of c, the kernel size of k. ConvT(c,k) means the transposed convolution block with the output channel of c, the kernel size of k. The kernel size and stride for all blocks are 4 and 2. The activation layer for convolution layers is ReLU, while for transposed convolution layers is LeakyReLU. The black dotted arrow lines between Conv and ConvT denote feature concatenation.

## B. Robustness Enhancement for Different Input Sizes

Taking HRNet-W32 as the backbone network, we conduct more experiments with different input size $128 \times 96$, $256 \times 192$, and $384 \times 288$ on COCO-C to verify the effectiveness of AdvMix for different input resolutions. As shown in Table 8, AdvMix improves both mPC and rPC significantly for all the three different input sizes.

Table 8. **Comparisons** of same backbone with different input sizes between standard training and AdvMix on COCO-C. Results are obtained with the same detection bounding boxes as [36]. We observe that both mPC and rPC are greatly improved, while almost maintaining performance of clean data.

| Method | Backbone | Input size | AP* | mPC | rPC |
|---|---|---|---|---|---|
| Standard | HRNet-W32 | $128 \times 96$ | 66.9 | 47.2 | 70.6 |
| **AdvMix** | HRNet-W32 | $128 \times 96$ | 66.3 | 48.9 | 73.8 |
| Standard | HRNet-W32 | $256 \times 192$ | 74.4 | 53.0 | 71.3 |
| **AdvMix** | HRNet-W32 | $256 \times 192$ | 74.7 | 55.5 | 74.3 |
| Standard | HRNet-W32 | $384 \times 288$ | 75.7 | 53.7 | 70.9 |
| **AdvMix** | HRNet-W32 | $384 \times 288$ | 76.2 | 56.8 | 74.5 |

## C. Visualization Results

In Figure 8, we provide more human pose results of images with different types of image corruptions, *i.e.* gaussian noise, motion blur, frost and contrast. For each triplet, we visualize the ground-truth (the left column), the prediction of the Standard HRNet-W32 (the middle column), and the prediction of HRNet with AdvMix (the right column). We observe that 1) the standard pose models suffer large performance drop on corrupted data, and 2) models trained with AdvMix perform consistently better than the baseline methods on various corruptions.
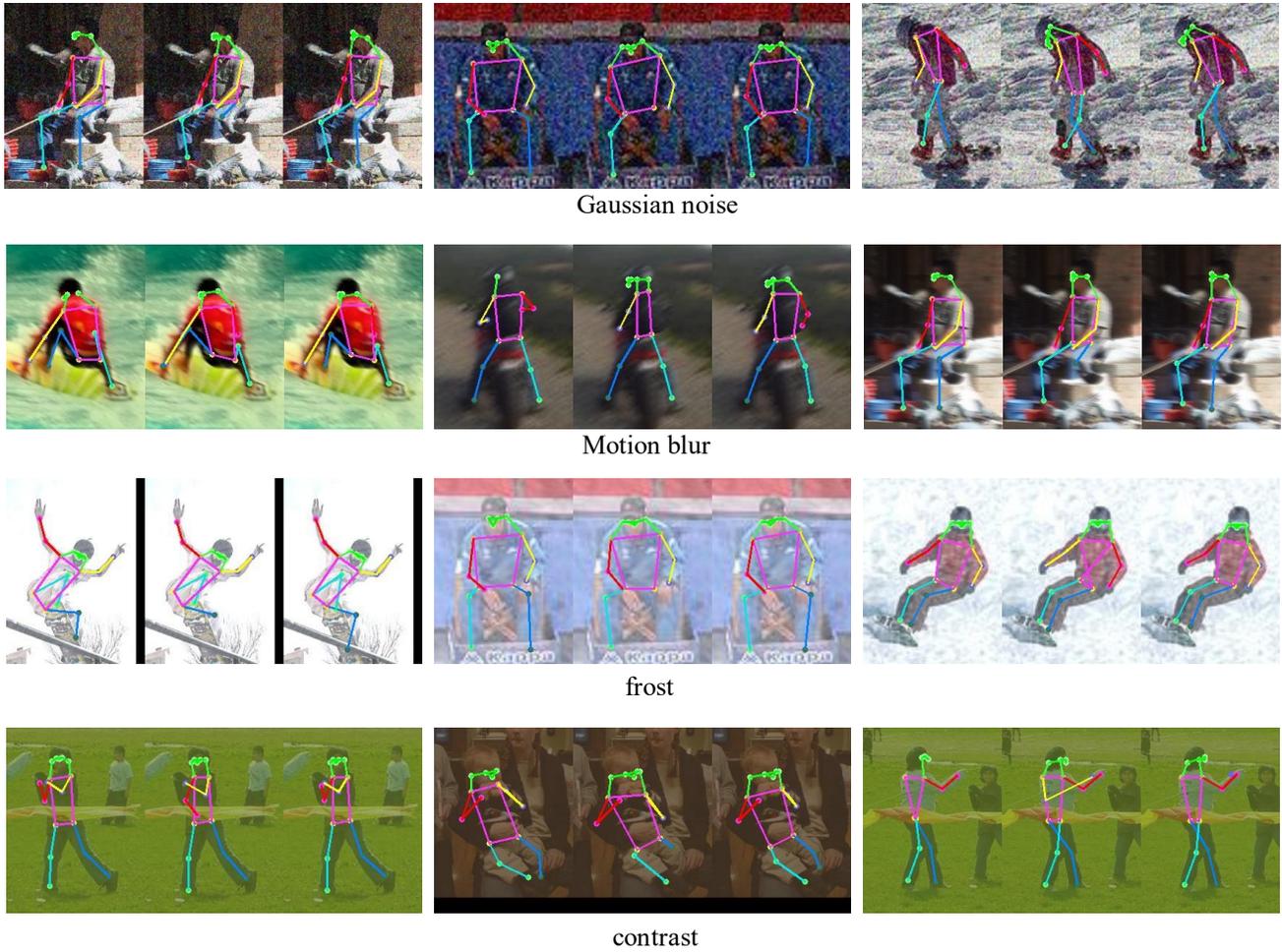
Gaussian noise

Motion blur

frost

contrast

Figure 8. Qualitative comparison between HRNet without and with AdvMix. For each image triplet, the images from left to right are ground truth, predicted results of Standard HRNet-W32, and predicted results of HRNet-W32 with AdvMix.