

# OmniPose: A Multi-Scale Framework for Multi-Person Pose Estimation

Bruno Artacho

Rochester Institute of Technology

Rochester, NY

bmartacho@mail.rit.edu

Andreas Savakis

Rochester Institute of Technology

Rochester, NY

andreas.savakis@rit.edu

## Abstract

We propose OmniPose, a single-pass, end-to-end trainable framework, that achieves state-of-the-art results for multi-person pose estimation. Using a novel waterfall module, the OmniPose architecture leverages multi-scale feature representations that increase the effectiveness of backbone feature extractors, without the need for post-processing. OmniPose incorporates contextual information across scales and joint localization with Gaussian heatmap modulation at the multi-scale feature extractor to estimate human pose with state-of-the-art accuracy. The multi-scale representations, obtained by the improved waterfall module in OmniPose, leverage the efficiency of progressive filtering in the cascade architecture, while maintaining multi-scale fields-of-view comparable to spatial pyramid configurations. Our results on multiple datasets demonstrate that OmniPose, with an improved HRNet backbone and waterfall module, is a robust and efficient architecture for multi-person pose estimation that achieves state-of-the-art results.

## 1. Introduction

Human pose estimation is an important task in computer vision that has generated high interest for methods on 2D pose estimation [33], [23], [34], [4], [31] and 3D [29], [41], [1]; on a single frame [5] or a video sequence [13]; for a single [34] or multiple subjects [6]. The main challenges of pose estimation, especially in multi-person settings, are due to the large number of degrees of freedom in the human body mechanics and the occurrence of joint occlusions. To overcome the difficulty of detecting joints under occlusion, it is common for methods to rely on statistical and geometric models to estimate occluded joints [26], [24]. Anchor poses are also used as a resource to overcome occlusion [29], but this could limit the generalization power of the model and its ability to handle unforeseen poses.

Inspired by recent advances in the use of multi-scale approaches for semantic segmentation [9], [37], and expand-



Figure 1. Pose estimation examples with our OmniPose method.

ing upon state-of-the-art results on 2D pose estimation by HRNet [31] and UniPose [4], we propose OmniPose, an expanded single-stage network that is end-to-end trainable and generates state-of-the-art results without requiring multiple iterations, intermediate supervision, anchor poses or post-processing. A main aspect of our novel architecture is an expanded multi-scale feature representation that combines an improved HRNet feature extractor with advanced Waterfall Atrous Spatial Pooling (WASPv2) module. Our improved WASPv2 module combines the cascaded approach for atrous convolution with larger Field-of-View (FOV) and is integrated with the network decoder offering significant improvements in accuracy.

Our OmniPose framework predicts the location of multiple people’s joints based on contextual information due to the multi-scale feature representation used in our network. The contextual approach allows our network to include the information from the entire frame, and consequently does not require post analysis based on statistical or geometric methods. In addition, the waterfall atrous module, allows a better detection of shapes, resulting in a more accurate estimation of occluded joints. Examples of pose estimation

obtained with our OmniPose method are shown in Figure 1. The main contributions of this paper are the following.

- We propose the novel OmniPose framework, a single-pass, end-to-end trainable, multi-scale approach that produces state-of-the-art results for multi-person pose estimation.
- We propose an improved Waterfall module that increases the performance of the network by using a larger field view while maintaining the high resolution of feature maps through the branches of the module. In addition, the WASPv2 module acts simultaneously as feature extractor and decoder, reducing the computational cost and size of the network.
- The OmniPose framework achieves an increase in performance by incorporating Gaussian heatmap modulation that enhances deconvolution operations in the multi-scale encoder-decoder architecture for a more accurate representation of joint locations and reduction of the quantization error in the network.
- We propose the novel lightweight OmniPose-Lite architecture that achieves high accuracy results while dramatically decreasing the number of parameters and computational cost of the network by leveraging the size reduction of separable convolutions throughout the network.

## 2. Related Work

In recent years, deep learning methods relying on Convolutional Neural Networks (CNNs) have achieved superior results in human pose estimation [33], [34], [6], [31], [29] over early works [28], [36]. The popular Convolutional Pose Machine (CPM) [34] proposed an architecture that refined joint detection via a set of stages in the network. Building upon [34], Yan et al. integrated the concept of Part Affinity Fields (PAF) resulting in the OpenPose method [6].

Multi-scale representations have been successfully used in backbone structures for pose estimation. Stacked hourglass networks [23] use cascaded structures of the hourglass method for pose estimation. Expanding on the hourglass structure, the multi-context approach in [14] relies on an hourglass backbone to perform pose estimation. The original backbone is augmented by the Hourglass Residual Units (HRU) with the goal of increasing the receptive FOV. Post-processing with Conditional Random Fields (CRFs) is used to assemble the relations between detected joints. However, the drawback of CRFs is increased complexity that requires high computational power and reduces speed.

The High-Resolution Network (HRNet) [31] includes both high and low resolution representations. HRNet benefits from the larger FOV of multi resolution, a capability that we achieve in a simpler fashion with our WASPv2 module. An analogous approach to HRNet is used by the Multi-

Stage Pose Network (MSPN) [19], where the HRNet structure is combined with cross-stage feature aggregation and coarse-to-fine supervision.

UniPose [4] combined the bounding box generation and pose estimation in a single one-pass network. This was achieved by the use of WASP module that increases significantly the multi-scale representation and FOV of the network, allowing the method to extract a greater amount of contextual information.

More recently, the HRNet structure was combined with multi-resolution pyramids in [12] to further explore multi-scale features. The Distribution-Aware coordinate Representation of Keypoints (DARK) method [38] aims to reduce loss during the inference processing of the decoder stage when using an HRNet backbone.

Aiming to use contextual information for pose estimation, the Cascade Prediction Fusion (CPF) [39] uses graphical components in order to exploit the context for pose estimation. Similarly, the Cascade Feature Aggregation (CFA) [30] uses semantic information for pose with a cascade approach. In a related context, Generative Adversarial Networks (GANs) were used in [7] to learn dependencies and contextual information for pose.

### 2.1. Multi-Scale Feature Representations

A challenge with CNN-based pose estimation, as well as semantic segmentation methods, is a significant reduction of resolution caused by pooling. Fully Convolutional Networks (FCN) [21] addressed this problem by deploying upsampling strategies across deconvolution layers that increase the size of feature maps back to the dimensions of the input image. In DeepLab, dilated or atrous convolutions [8] were used to increase the size of the receptive fields in the network and avoid downsampling in a multi-scale framework. The Atrous Spatial Pyramid Pooling (ASPP) approach assembles atrous convolutions in four parallel branches with different rates, that are combined by fast bilinear interpolation with an additional factor of eight. This configuration recovers the feature maps in the original image resolution. The increase in resolution and FOV in the ASPP network can be beneficial for a contextual detection of body parts during pose estimation.

Improving upon [8], the waterfall architecture of the WASP module incorporates multi-scale features without immediately parallelizing the input stream [3], [4]. Instead, it creates a waterfall flow by first processing through a filter and then creating a new branch. The WASP module goes beyond the cascade approach by combining the streams from all its branches and average pooling of the original input to achieve a multi-scale representation.

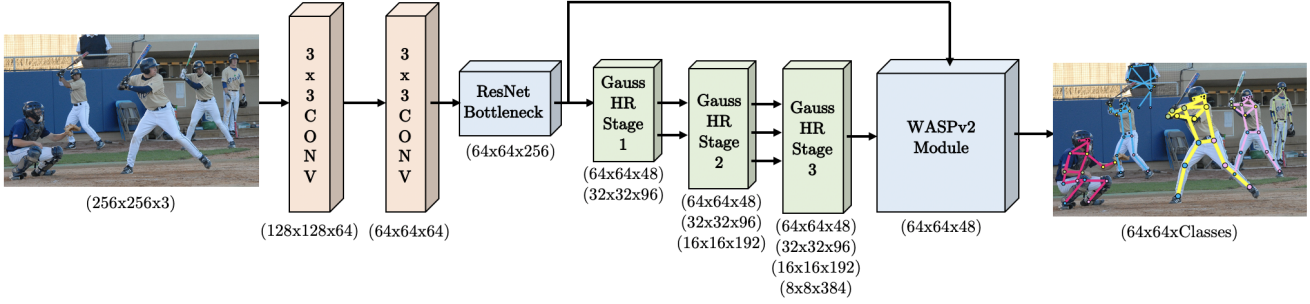


Figure 2. OmniPose framework for multi-person pose estimation. The input color image is fed through the improved HRNet backbone and WASPv2 module to generate one heatmap per joint or class.

### 3. OmniPose Architecture

The proposed OmniPose framework, illustrated in Figure 2, is a single-pass, single output branch network for pose estimation of multiple people instances. OmniPose incorporates improvements in feature representation from multi-scale approaches [31], [38] and an encoder-decoder structure combined with spatial pyramid pooling [10] and our proposed advanced waterfall module (WASPv2).

The processing pipeline of the OmniPose architecture is shown in Figure 2. The input image is initially fed into a deep CNN backbone, consisting of our modified version of HRNet [31]. The resultant feature maps are processed by our WASPv2 decoder module that generates  $K$  heatmaps, one for each joint, with the corresponding confidence maps. The integrated WASPv2 decoder in our network generates detections for both visible and occluded joints while maintaining the image high resolution through the network.

Our architecture includes several innovations to increase accuracy. The first is the application of atrous convolutions and waterfall architecture of the WASPv2 module, that increases the network’s capacity to compute multi-scale contextual information. This is accomplished by the probing of feature maps at multiple rates of dilation during convolutions, resulting in a larger FOV in the encoder. Our architecture integrates the decoding process within the WASPv2 module without requiring a separate decoder. Additionally, our network demonstrates good ability to detect shapes by the use of spatial pyramids combined with our modified HRNet feature extraction, as indicated by state-of-the-art (SOTA) results. Finally, the modularity of the OmniPose framework enables easy implementation and training.

OmniPose leverages the large number of feature maps at multiple scales in the proposed WASPv2 module. In addition, we improved the results of the backbone network by incorporating gaussian modulated deconvolutions in place of the upsampling operations during transition stages of the original HRNet architecture. The modified HRNet feature extractor is followed by the improved and integrated

multi-scale waterfall configuration of the WASPv2 decoder, which further improves the efficiency of the joint detection with the incorporation of Gaussian heatmap modulation of the decoder stage, and full integration with the decoder module.

Targeting the reduction of computational cost and number of parameters, we implement separable convolutions replacing the initial two layers of strided convolutions in our model and the atrous convolutions in the WASPv2 module. Figure 3 demonstrates the implementation of the strided convolution that consists of a spatial (or depth-wise) convolution through the individual channels of the feature maps, followed by a rectified linear unit (ReLU) activation function, and a point-wise convolution to incorporate all the layers of the feature maps.

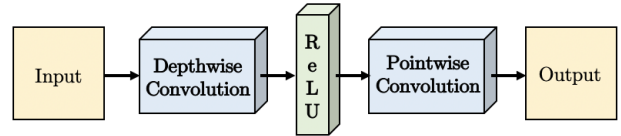


Figure 3. Implementation of our separable convolution. The cascade of depth-wise convolution, ReLU activation, and point-wise convolution replace the standard convolution in order to reduce the number of parameters and computations in the network.

#### 3.1. WASPv2 Module

The proposed advanced “Waterfall Atrous Spatial Pyramid” module, or WASPv2, shown in Figure 4, generates an efficient multi-scale representation that helps OmniPose achieve SOTA results. Our improved WASPv2 module expands the feature extraction through its multi-level architecture. It increases the FOV of the network with consistent high resolution processing of the feature maps in all its branches, which contributes to higher accuracy. In addition, WASPv2 generates the final heatmaps for joint localization without the requirement of an additional decoder module,

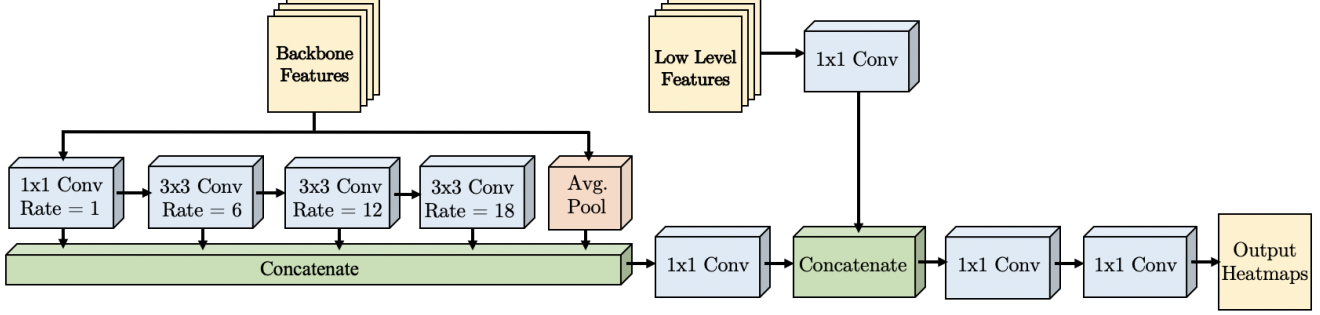


Figure 4. The proposed WASPv2 advanced waterfall module. The inputs are 48 features maps from the modified HRNet backbone and low-level features from the initial layers of the framework.

interpolation or pooling operations.

The WASPv2 architecture relies on atrous convolutions to maintain a large FOV, performing a cascade of atrous convolutions at increasing rates to gain efficiency. In contrast to ASPP [10], WASPv2 does not immediately parallelize the input stream. Instead, it creates a waterfall flow by first processing through a filter and then creating a new branch. In addition, WASPv2 goes beyond the cascade approach by combining the streams from all its branches and average pooling of the original input to achieve a multi-scale representation.

Expanding upon the original WASP module [4], WASPv2 incorporates the decoder in an integrated unit shown in Figure 4, and processes both of the waterfall branches with different dilation rates and low-level features in the same higher resolution, resulting in a more accurate and refined response. The WASPv2 module output  $f_{WASPv2}$  is described as follows:

$$f_{Waterfall} = K_1 \otimes \left( \sum_{i=1}^4 (K_{d_i} \otimes f_{i-1}) + AP(f_0) \right) \quad (1)$$

$$f_{WASPv2} = K_1 \otimes (K_1 \otimes (K_1 \otimes f_{LLF} + f_{Waterfall})) \quad (2)$$

where  $\otimes$  represents convolution,  $f_0$  is the input feature map,  $f_i$  is the feature map resulting from the  $i^{th}$  atrous convolution,  $AP$  is the average pooling operation,  $f_{LLF}$  are the low-level feature maps,  $K_1$  and  $K_{d_i}$  represent convolutions of kernel size  $1 \times 1$  and  $3 \times 3$  with dilations of  $d_i = [1, 6, 12, 18]$ , as shown in Figure 4. After concatenation, the feature maps are combined with low level features. The last  $1 \times 1$  convolution brings the number of feature maps down to the final number of joints for the pose estimation.

Differently than the previous version of WASP, our WASPv2 integrates in the same resolution the feature maps from the low-level features and the first part of the waterfall module, converting the score maps from the WASPv2 module to heatmaps corresponding to body joints. Due to the higher resolution afforded by the modified HRNet backbone, the WASPv2 module directly outputs the final

heatmaps without requiring an additional decoder module or need for bilinear interpolations to resize the output to the original input size.

Aiming to reduce the computational complexity and size of the network, and inspired by [10], our WASPv2 module implements separable atrous convolutions to its feature extraction waterfall branches. The inclusion of separable atrous convolutions in the WASPv2 module further reduces the number of parameters and computation cost of the framework.

### 3.2. Gaussian Heatmap Modulation

Conventional interpolation or upsampling methods during the decoding stage of the network result in an inevitable loss in resolution and consequently accuracy, limiting the potential of the network. Motivated by recent results with distribution aware modulation [38], we include Gaussian heatmap modulation to all interpolation stages of our network, resulting in a more accurate and robust network that diminishes the localization error due to interpolation.

Gaussian interpolation allows the network to achieve sub-pixel resolution for peak localization following the anticipated Gaussian pattern of the feature response. This method results in a smoother response and more accurate peak prediction for joints, by eliminating false positives in noisy responses during the joint detection.

Figure 5 demonstrates the modularization of a feature map response in our improved HRNet feature extractor. We utilize a transposed convolution operation of the feature map  $f_D$  with a Gaussian kernel  $K$ , shown in Equation (2), aiming to approximate the response shape to the expected label of the dataset during training. The feature maps  $f_G$  after the Gaussian convolution operation are:

$$f_G = K \otimes f_D \quad (3)$$

This behavior is learned and reproduced by the network during all parts of training, validation, and testing.

Following convolution with the Gaussian kernel, the modulation of the interpolation output is scaled to  $f_{G_s}$  by



mapping  $f_G$  to the range of the response of the original feature map  $f_D$  using:

$$f_{G_s} = \frac{f_G - \min(f_G)}{\max(f_G) - \min(f_G)} * \max(f_D) \quad (4)$$

Our Gaussian heatmap modulation approach allows for better localization of the coordinates during the transposed convolutions, by overcoming the quantization error naturally inherited from the increase in resolution.

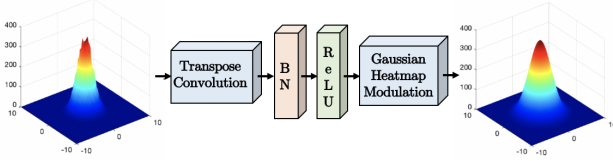


Figure 5. Illustration of the proposed transpose convolution with Gaussian modulation replacing upsampling stages of the network.

### 3.3. OmniPose-Lite

We introduce OmniPose-Lite, a lightweight version of OmniPose that is suitable for mobile and embedded platforms, as it achieves a drastic reduction in memory requirements and operations required for computation. The proposed OmniPose-Lite leverages the reduced computational complexity and size of separable convolutions, inspired by results obtained by MobileNet [17].

We implemented separable strided convolutions, as shown in Figure 3, for all convolutional layers of the original HRNet backbone, and implemented atrous separable convolutions in the WASPv2 decoder, resulting in a reduction of 74.3% of the network GFLOPs, from 22.6 GFLOPs to 5.8 GFLOPs required to process an image of size 256x256. In addition, OmniPose-Lite also reduces the number of parameters by 71.4%, from 67.9M to 19.4M.

The small size of the proposed OmniPose-Lite architecture, in combination with the reduced number of parameters allows the implementation of the OmniPose architecture for mobile applications without a large computational burden.

## 4. Datasets

We performed multi-person experiments on two datasets: Common Objects in Context (COCO) [20] and MPII [2]. The COCO dataset [20] is composed of over 200K images containing over 250K instances of the person class. The labelled poses contain 17 keypoints. The dataset is considered a challenging dataset due to the large number of images in a diverse set of scales and occlusion for poses in the wild.

The MPII dataset [2] contains approximately 25K images of annotated body joints of over 40K subjects. The

images are collected from YouTube videos in 410 everyday activities. The dataset contains frames with joints annotations, head and torso orientations, and body part occlusions.

In order to better train our network for joint detection, ideal Gaussian maps were generated at the joint locations in the ground truth. These maps are more effective for training than single points at the joint locations, and were used to train our network to generate Gaussian heatmaps corresponding to the location of each joint. Gaussians with different  $\sigma$  values were considered and a value of  $\sigma = 3$  was adopted, resulting in a well defined Gaussian curve for both the ground truth and predicted outputs. This value of  $\sigma$  also allows enough separation between joints in the image.

## 5. Experiments

OmniPose experiments were based on the metrics set by each dataset, and procedures applied by [31], [12], and [38].

### 5.1. Metrics

For the evaluation of OmniPose, various metrics were used depending on previously reported results and the available ground truth for each dataset. The first metric used is the Percentage of Correct Keypoints (PCK). This metric considers the prediction of a keypoint correct when a joint detection lies within a certain threshold distance of the ground truth. The commonly used threshold of PCKh@0.5 was adopted for the MPII dataset, which refers to a threshold of 50% of the head diameter.

In the case of the COCO dataset, the evaluation is done based on the Object Keypoint Similarity metric (OKS).

$$OKS = \frac{(\sum_i e^{-d_i^2/2s^2k_i^2})\delta(v_i > 0)}{\sum_i \delta(v_i > 0)} \quad (5)$$

where,  $d_i$  is the Euclidian distance between the estimated keypoint and its ground truth,  $v_i$  indicates if the keypoint is visible,  $s$  is the scale of the corresponding target, and  $k_i$  is the falloff control constant.

Since OKS is measured in an analogous form as the intersection over the union (IOU), and following the evaluation framework set by [20], we report OKS as the Average Precision (AP) for the IOUs for all instances between 0.5 and 0.95 (AP), at 0.5 ( $AP^{50}$ ) and 0.75 ( $AP^{75}$ ), as well as instances of medium ( $AP^M$ ) and large size ( $AP^L$ ). We also report the Average Recall between 0.5 and 0.95 (AR).

### 5.2. Parameter Selection

We process the input image in a set of different resolutions, reporting the trade-off of network size and accuracy performance. For that reason, the batch size varied depending on the size of the dataset images. We considered different rates of dilation on the WASP module and found that



Figure 6. Pose estimation examples using OmniPose with the MPII dataset.

Method	Params (M)	GFLOPs	Head	Shoulder	Elbow	Wrist	Hip	Knee	Ankle	PCKh @0.2
<b>OmniPose (WASPV2)</b>	68.1	22.6	<b>97.4%</b>	<b>97.1%</b>	<b>92.4%</b>	<b>88.7%</b>	<b>91.2%</b>	<b>89.9%</b>	<b>85.8%</b>	<b>92.3%</b>
<b>OmniPose (WASP)</b>	68.2	23.0	97.4%	96.6%	91.9%	87.2%	90.1%	88.0%	83.9%	91.2%
DarkPose [38]	63.6	19.5	97.2%	95.9%	91.2%	86.7%	89.7%	86.7%	84.0%	90.6%
HRNet [31]	63.6	19.5	97.1%	95.9%	90.3%	86.5%	89.1%	87.1%	83.3%	90.3%
<b>OmniPose-Lite</b>	<b>19.4</b>	<b>5.8</b>	96.6%	95.8%	89.1%	84.3%	89.0%	84.1%	79.6%	89.0%
CMU Pose [6]	-	-	92.4%	90.4%	80.9%	70.8%	79.5%	73.1%	66.5%	79.1%
SPM [25]	-	-	92.0%	88.5%	78.6%	69.4%	77.7%	73.8%	63.9%	77.7%
RMPE [15]	-	-	88.4%	86.5%	78.6%	70.4%	74.4%	73.0%	65.8%	76.7%

Table 1. OmniPose results and comparison with SOTA methods for the MPII dataset for validation.

larger rates result in better prediction. A set of dilation rates of  $r = \{1, 6, 12, 18\}$  was selected for the WASPV2 module.

We calculate the learning rate based on the step method, where the learning rate started at  $10^{-3}$  and was reduced in two steps by an order of magnitude at each steps at 170 and 200 epochs, following procedures set by [38]. All experiments were performed using PyTorch on Ubuntu 16.04. The workstation has an Intel i5-2650 2.20GHz CPU with 16GB of RAM and an NVIDIA Tesla V100 GPU.

## 6. Results

We present OmniPose results on two large datasets and provide comparisons with state-of-the-art methods.

### 6.1. Experimental results on the MPII dataset

During our experiments on the MPII dataset, we performed a series of ablation studies to analyze the gains due to different aspects of our method. Table 2 demonstrates the results for the inclusion of the Gaussian deconvolution modulation (GDM) in the HRNet backbone, and improvements gained by initially using the original WASP module

[3], [4], and then our proposed advanced WASPV2 in combination with the improved HRNet feature extractor.

Method	GDM	WASP	WASPV2	PCKh @0.2
DarkPose [38]				90.6%
OmniPose	✓			91.0%
OmniPose	✓	✓		91.2%
OmniPose	✓		✓	92.3%

Table 2. Results using different versions of OmniPose and comparison with SOTA for the MPII dataset for validation. GDM represents the use of Gaussian Deconvolution Modulation in the modified HRNet backbone, and WASP and WASPV2 indicates the use of the waterfall modules in the network.

Our OmniPose method progressively increases its performance with the addition of innovations, resulting in 1.9% improvement over DarkPose [38]. Most significantly, the integration of the enhanced multi-scale extraction with WASPV2 substantially increases the keypoints detection, particularly for occluded joints.

Following our experiments evaluating the individual





Figure 7. Pose estimation examples using OmniPose with the COCO dataset.

Method	Input Size	Params (M)	GFLOPs	AP	$AP^{50}$	$AP^{75}$	$AP^M$	$AP^L$	AR
<b>OmniPose (WASPv2)</b>	384x288	68.1	37.9	<b>79.5%</b>	<b>93.6%</b>	<b>85.9%</b>	<b>76.0%</b>	<b>84.6%</b>	<b>81.9%</b>
OmniPose (WASP)	384x288	68.2	38.6	79.2%	93.6%	85.7%	75.9%	84.2%	81.6%
DarkPose [38]	384x288	63.6	32.9	76.8%	90.6%	83.2%	72.8%	84.0%	81.7%
HRNet [31]	384x288	63.6	32.9	76.3%	90.8%	82.9%	72.3%	83.4%	81.2%
EvoPose2D [22]	384x288	7.3	5.6	75.1%	90.2%	81.9%	71.5%	81.7%	81.0%
Simple Baseline [35]	384x288	68.6	35.6	74.3%	89.6%	81.1%	70.5%	79.7%	79.7%

Table 3. OmniPose results and comparison with SOTA methods for the COCO dataset for validation.

Method	Input Size	Params (M)	GFLOPs	AP	$AP^{50}$	$AP^{75}$	$AP^M$	$AP^L$	AR
<b>OmniPose (WASPv2)</b>	384x288	68.1	37.9	<b>76.4%</b>	92.6%	<b>83.7%</b>	<b>72.6%</b>	<b>82.6%</b>	81.2%
DarkPose [38]	384x288	63.6	32.9	76.2%	92.5%	83.6%	72.5%	82.4%	81.1%
MSPN [19]	384x288	120	19.9	76.1%	<b>93.4%</b>	83.8%	72.3%	81.5%	<b>81.6%</b>
HRNet [31]	384x288	63.6	32.9	75.5%	92.5%	83.3%	71.9%	81.5%	80.5%
Simple Baseline [35]	384x288	68.6	35.6	73.7%	91.9%	70.3%	81.1%	80.0%	79.0%
RMPE [15]	320x256	28.1	26.7	72.3%	89.2%	79.1%	68.0%	78.6%	-
CPN [11]	384x288	-	-	72.1%	91.4%	80.0%	68.7%	77.2%	78.5%
IPR [32]	256x256	45.1	11.0	67.8%	88.2%	74.8%	63.9%	74.0%	-
G-RMI [27]	256x256	42.6	57.0	64.9%	85.5%	71.3%	62.3%	70.0%	69.7
Mask-RCNN [16]	-	-	-	63.1%	87.3%	68.7%	57.8%	71.4%	-
CMU Pose [6]	-	-	-	61.8%	84.9%	57.1%	67.5%	68.2%	66.5%

Table 4. OmniPose results and comparison with SOTA methods for the COCO dataset for test.

contributions of this work, we compared the results of OmniPose with other methods, as shown in Table 1. OmniPose achieved a overall PCKh@0.2 of 92.3%, showing significant gains in comparison to state-of-the-art. It is significant that OmniPose results in a improvement to previous SOTA methods in all individuals groups of joints for pose estimation, demonstrating the robustness and performance of our framework, particularly to harder to detect joints such as ankles (2.1% improvement from previous state-of-the-art)

and wrists (2.3% above previous state-of-the-art). Figure 6 demonstrates successful detections on the main person in MPII images. These examples illustrate that OmniPose deals effectively with occlusion, e.g. in the case of the skier.

OmniPose-Lite achieves accuracy of 89.0% while reducing computational cost by 74.3% for the MPII validation dataset (Table 1). This demonstrates its ability to significantly reduce size and computational cost, while maintaining good performance compared to heavier SOTA methods.

## 6.2. Experimental results on the COCO dataset

We next performed training and testing on the COCO dataset, which is more challenging due to the large number of diverse images with multiple people in close proximity, as well as images lacking a person instance.

We performed experiment to compare the proposed improvements of OmniPose with the original HRNet framework. OmniPose outperforms HRNet in terms of average precision for different input resolutions, as shown in Figure 8 for 3 different versions of OmniPose: small ( $128 \times 96$ ), medium ( $256 \times 192$ ), and large ( $384 \times 288$ ); as well as lower resolution versions of OmniPose-Lite. OmniPose demonstrates an increase in performance for all resolutions compared with the original HRNet architecture. The accuracy of the OmniPose framework steadily increases with the increase of the input resolution, but there is a trade-off with processing time due to the larger number of image pixels that are processed in the network.

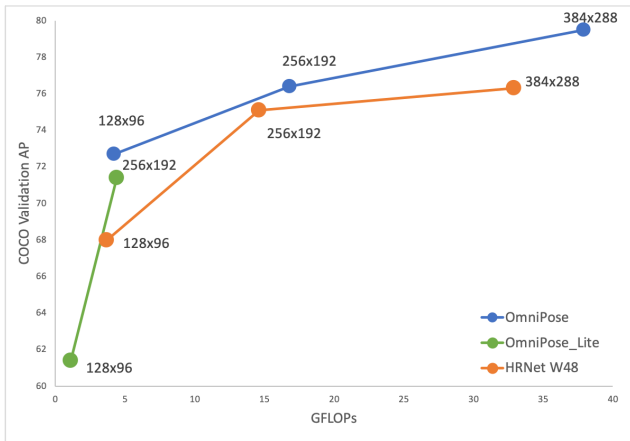


Figure 8. Average Precision comparison of OmniPose to the original HRNet method for different input resolutions.

OmniPose was compared with SOTA methods for the validation set of the COCO dataset. The results in Table 3 demonstrate that OmniPose shows significant improvement over the previous SOTA. The modification of the HRNet backbone, combined with the WASPv2 module results in an improved accuracy of 79.5%, a significant increase of 4.2% compared with the original HRNet, and 7.0% compared with the baseline model.

OmniPose improves accuracy for all detection metric sizes and IOU for COCO, as was the case for MPII. Most significantly, in harder detections the AP for person instances of medium size obtained by OmniPose shows an increase of 4.4% over the previous state-of-the-art. These results demonstrate the increased capability of OmniPose to estimate harder poses using a reduced number of pixels due to the multi-scale features from the WASPv2 module.

Comparing OmniPose-Lite to lightweight architectures, OmniPose-Lite shows a reduction of size of 12% while increasing the performance on the COCO validation set by 8.7% compared to the popular MobileNetV2 [17], as shown in Table 5. This establishes a significant improvement for lightweight pose estimation methods.

Method	Input Size	GFLOPs	AP
<b>OmniPose-Lite</b>	256x192	<b>4.4</b>	<b>71.4%</b>
MobileNetV2 [17]	256x192	5.0	65.7%

Table 5. Lightweight comparison for the COCO validation dataset.

Example results for the validation COCO dataset are shown in Figure 7. It is noticeable from these examples that our method identifies the location of symmetric body joints with high precision, providing high accuracy for challenging scenarios, that include multiple instances of people in near proximity and occluded joints, such as ankles and wrists, that are harder to detect. Challenging conditions include the detection of joints when limbs are not sufficiently separated or occlude each other, where OmniPose demonstrates a robust ability to detect.

We also compared OmniPose with SOTA methods using the COCO test-dev dataset, which contains a significantly larger number of images. OmniPose achieved a new state-of-the-art performance compared with other methods without the use of additional training data or postprocessing, achieving an average precision of 76.4%. Confirming our findings from previous datasets, OmniPose shows the most significant improvements in smaller targets.

## 6.3. Single person and video datasets

We further tested OmniPose on the Leeds Sports Pose (LSP) [18] dataset, for single person pose estimation. OmniPose achieved a significant improvement of 5% from the previous state-of-the-art achieved by UniPose [4], resulting in a PCK@0.2 of 99.5% and saturating the pose estimation performance for the LSP dataset. Similarly, running OmniPose on the PennAction dataset for pose estimation in short sports videos [40] shows saturation in performance by achieving state-of-the-art accuracy of 99.4% PCK.

## 7. Conclusion

We presented the OmniPose framework for multi-person pose estimation. The OmniPose pipeline utilizes the improved WASPv2 module that features a waterfall flow with a cascade of atrous convolutions and multi-scale representations. The OmniPose framework achieves state-of-the-art performance, with an improved HRNet feature extractor utilizing transposed convolutions with Gaussian heatmap modulation, replacing interpolations. OmniPose is an end-to-end trainable architecture that does not require anchor



poses or postprocessing. The results of the OmniPose framework demonstrated state-of-the-art performance on several datasets using various metrics.

## References

- [1] Riza Alp Guler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7297–7306, 2018. 1
- [2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE International Conference on Computer Vision (ICCV)*, 2014. 5
- [3] Bruno Artacho and Andreas Savakis. Waterfall atrous spatial pooling architecture for efficient semantic segmentation. *Sensors*, 19(24):5361, 2019. 2, 6
- [4] Bruno Artacho and Andreas Savakis. Unipose: Unified human pose estimation in single images and videos. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1, 2, 4, 6, 8
- [5] Vasileios Belagiannis and Andrew Zisserman. Recurrent human pose estimation. In *IEEE International Conference on Automatic Face & Gesture Recognition*, pages 468–475, 2017. 1
- [6] Z. Cao, T. Simon, S. E. Wei, and Y. Sheikh. Realtime multi-person 2D pose estimation using part affinity fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 2, 6, 7
- [7] Zhongzheng Cao, Rui Wang, Xiangyang Wang, Zhi Liu, and Xiaoqiang Zhu. Improving human pose estimation with self-attention generative adversarial networks. In *2019 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, pages 567–572. IEEE, 2019. 2
- [8] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution and fully connected cfrs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–845, 2018. 2
- [9] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam. Rethinking atrous convolution for semantic image segmentation. *CoRR*, abs/1706.05587, 2017. 1
- [10] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *European conference on computer vision (ECCV)*, pages 801–818, 2018. 3, 4
- [11] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun. Cascaded pyramid network for multi-person pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7103–7112, 2018. 7
- [12] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 5
- [13] Yu Cheng, Bo Yang, Bo Wang, Wending Yan, and Robby T. Tan. Occlusion-aware networks for 3d human pose estimation in video. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 1
- [14] Xiao Chu, Wei Yang, Wanli Ouyang, Cheng Ma, Alan L. Yuille, and Xiaogang Wang. Multi-context attention for human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1831–1840, 2017. 2
- [15] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. Rmpe: Regional multi-person pose estimation. In *IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 6, 7
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2961–2969, 2017. 7
- [17] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 5, 8
- [18] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *British Machine Vision Conference*, 2010. 8
- [19] Wenbo Li, Zhicheng Wang, Binyi Yin, Qixiang Peng, Yuming Du, Tianzi Xiao, Gang Yu, Hongtao Lu, Yichen Wei, and Jian Sun. Rethinking on multi-stage networks for human pose estimation. *arXiv preprint arXiv:1901.00148*, 2019. 2, 7
- [20] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. 5
- [21] J. Long, E. Shelhamer, and T. Darrel. Fully convolutional networks for semantic segmentation. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2
- [22] William McNally, Kanav Vats, Alexander Wong, and John McPhee. EvoPose2D: Pushing the boundaries of 2D human pose estimation using neuroevolution. *arXiv preprint arXiv:2011.08446*, 2020. 7
- [23] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *European Conference on Computer Vision (ECCV)*, pages 483–499. Springer, 2016. 1, 2
- [24] Aiden Nibali, Zhen He, Stuart Morgan, and Luke Prendergast. 3d human pose estimation with 2D marginal heatmaps. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1477–1485, 2019. 1
- [25] Xuecheng Nie, Jiashi Feng, Jianfeng Zhang, and Shuicheng Yan. Single-stage multi-person pose machines. In *IEEE International Conference on Computer Vision*, pages 6951–6960, 2019. 6
- [26] George Papandreou, Tyler Zhu, Liang-Chieh Chen, Spyros Gidaris, Jonathan Tompson, and Kevin Murphy. Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model. In *European Conference on Computer Vision (ECCV)*, 2018. 1

- [27] George Papandreou, Tyler Zhu, Nori Kanazawa, Alexander Toshev, Jonathan Tompson, Chris Bregler, and Kevin Murphy. Towards accurate multi-person pose estimation in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4903–4911, 2017. 7
- [28] Leonid Pishchulin, Mykhaylo Andriluka, Peter Gehler, and Bernt Schiele. Poselet conditioned pictorial structures. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 588–595, 2013. 2
- [29] Gregory Rogez, Philippe Weinzaepfel, and Cordelia Schmid. LCR-Net: Localization-classification-regression for human pose. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, HI, USA, 2017. 1, 2
- [30] Zhihui Su, Ming Ye, Guohui Zhang, Lei Dai, and Jianda Sheng. Cascade feature aggregation for human pose estimation. *arXiv preprint arXiv:1902.07837*, 2019. 2
- [31] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2, 3, 5, 6, 7
- [32] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. Integral human pose regression. In *European Conference on Computer Vision (ECCV)*, September 2018. 7
- [33] Alexander Toshev and Christian Szegedy. DeepPose: Human pose estimation via deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1653–1660, 2014. 1, 2
- [34] Shih-En Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional pose machines. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2
- [35] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *European Conference on Computer Vision (ECCV)*, 2018. 7
- [36] Yi Yang and Deva Ramanan. Articulated human detection with flexible mixtures of parts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2878–2890, 2012. 2
- [37] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. In *International Conference on Learning Representations*, 2016. 1
- [38] Feng Zhang, Xiatian Zhu, Hanbin Dai, Mao Ye, and Ce Zhu. Distribution-aware coordinate representation for human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2, 3, 4, 5, 6, 7
- [39] Hong Zhang, Hao Ouyang, Shu Liu, Xiaojuan Qi, Xiaoyong Shen, Ruigang Yang, and Jiaya Jia. Human pose estimation with spatial contextual information. *arXiv preprint arXiv:1901.01760*, 2019. 2
- [40] Weiyu Zhang, Menglong Zhu, and Konstantinos G Derpanis. From actemes to action: A strongly-supervised representation for detailed action understanding. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2248–2255, 2013. 8
- [41] Xiaowei Zhou, Menglong Zhu, Georgios Pavlakos, Spyridon Leonardos, Konstantinos G Derpanis, and Kostas Daniilidis. MonoCap: Monocular human motion capture using a cnn coupled with a geometric prior. *IEEE Transactions on*

*Pattern Analysis and Machine Intelligence*, 41(4):901–914, 2018. 1