

# Deep High-Resolution Representation Learning for Visual Recognition

Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao

**Abstract**—High-resolution representations are essential for position-sensitive vision problems, such as human pose estimation, semantic segmentation, and object detection. Existing state-of-the-art frameworks first encode the input image as a low-resolution representation through a subnetwork that is formed by connecting high-to-low resolution convolutions *in series* (e.g., ResNet, VGGNet), and then recover the high-resolution representation from the encoded low-resolution representation. Instead, our proposed network, named as High-Resolution Network (HRNet), maintains high-resolution representations through the whole process. There are two key characteristics: (i) Connect the high-to-low resolution convolution streams *in parallel*; (ii) Repeatedly exchange the information across resolutions. The benefit is that the resulting representation is semantically richer and spatially more precise. We show the superiority of the proposed HRNet in a wide range of applications, including human pose estimation, semantic segmentation, and object detection, suggesting that the HRNet is a stronger backbone for computer vision problems. All the codes are available at <https://github.com/HRNet>.

**Index Terms**—HRNet, high-resolution representations, low-resolution representations, human pose estimation, semantic segmentation, object detection.

## 1 INTRODUCTION

DEEP convolutional neural networks (DCNNs) have achieved state-of-the-art results in many computer vision tasks, such as image classification, object detection, semantic segmentation, human pose estimation, and so on. The strength is that DCNNs are able to learn richer representations than conventional hand-crafted representations.

Most recently-developed classification networks, including AlexNet [77], VGGNet [126], GoogleNet [133], ResNet [54], etc., follow the design rule of LeNet-5 [81]. The rule is depicted in Figure 1 (a): gradually reduce the spatial size of the feature maps, connect the convolutions from high resolution to low resolution in series, and lead to a *low-resolution representation*, which is further processed for classification.

*High-resolution representations* are needed for *position-sensitive tasks, e.g., semantic segmentation, human pose estimation, and object detection*. The previous state-of-the-art methods adopt the high-resolution recovery process to raise the representation resolution from the low-resolution representation outputted by a classification or classification-like network as depicted in Figure 1 (b), e.g., Hourglass [105], SegNet [3], DeconvNet [107], U-Net [119], SimpleBaseline [152], and encoder-decoder [112]. In addition, *dilated convolutions* are used to remove some down-sample layers and thus yield medium-resolution representations [19], [181].

We present a novel architecture, namely High-Resolution Net (HRNet), which is able to *maintain high-resolution representations* through the whole process. We start from a high-resolution convolution stream, gradually add high-to-low resolution convolution streams one by one, and connect the

multi-resolution streams in parallel. The resulting network consists of several (4 in this paper) stages as depicted in Figure 2, and the  $n$ th stage contains  $n$  streams corresponding to  $n$  resolutions. We conduct repeated multi-resolution fusions by exchanging the information across the parallel streams over and over.

The high-resolution representations learned from HRNet are not only *semantically strong* but also *spatially precise*. This comes from two aspects. (i) Our approach connects high-to-low resolution *convolution streams in parallel rather than in series*. Thus, our approach is able to *maintain the high resolution instead of recovering high resolution* from low resolution, and accordingly the learned representation is potentially spatially more precise. (ii) Most existing fusion schemes aggregate high-resolution low-level and high-level representations obtained by upsampling low-resolution representations. Instead, we repeat multi-resolution fusions to boost the high-resolution representations with the help of the low-resolution representations, and vice versa. As a result, all the high-to-low resolution representations are semantically strong.

We present two versions of HRNet. The first one, named as HRNetV1, only outputs the high-resolution representation computed from the high-resolution convolution stream. We apply it to human pose estimation by following the heatmap estimation framework. We empirically demonstrate the superior pose estimation performance on the COCO keypoint detection dataset [94].

The other one, named as HRNetV2, combines the representations from all the high-to-low resolution parallel streams. We apply it to semantic segmentation through estimating segmentation maps from the combined high-resolution representation. The proposed approach achieves state-of-the-art results on PASCAL-Context, Cityscapes, and

• J. Wang is with Microsoft Research, Beijing, P.R. China.  
E-mail: [jingdw@microsoft.com](mailto:jingdw@microsoft.com)

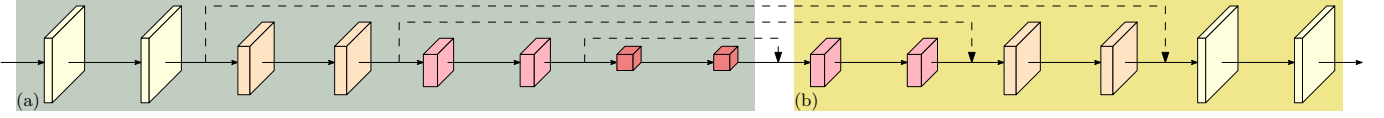


Fig. 1. The structure of recovering high resolution from low resolution. (a) A low-resolution representation learning subnetwork (such as VGGNet [126], ResNet [54]), which is formed by connecting high-to-low convolutions in series. (b) A high-resolution representation recovering subnetwork, which is formed by connecting low-to-high convolutions in series. Representative examples include SegNet [3], DeconvNet [107], U-Net [119] and Hourglass [105], encoder-decoder [112], and SimpleBaseline [152].

LIP with similar model sizes and lower computation complexity. We observe similar performance for HRNetV1 and HRNetV2 over COCO pose estimation, and the superiority of HRNetV2 to HRNet1 in semantic segmentation.

In addition, we construct a multi-level representation, named as HRNetV2p, from the high-resolution representation output from HRNetV2, and apply it to state-of-the-art detection frameworks, including Faster R-CNN, Cascade R-CNN [12], FCOS [136], and CenterNet [36], and state-of-the-art joint detection and instance segmentation frameworks, including Mask R-CNN [53], Cascade Mask R-CNN, and Hybrid Task Cascade [16]. **The results show that our method gets detection performance improvement and in particular dramatic improvement for small objects.**

## 2 RELATED WORK

We review closely-related representation learning techniques developed mainly for human pose estimation [57], semantic segmentation and object detection, from three aspects: low-resolution representation learning, high-resolution representation recovering, and high-resolution representation maintaining. Besides, we mention about some works related to multi-scale fusion.

**Learning low-resolution representations.** The fully-convolutional network approaches [99], [124] compute low-resolution representations by removing the fully-connected layers in a classification network, and estimate their **coarse segmentation maps**. The estimated segmentation maps are improved by combining the fine segmentation score maps estimated from intermediate low-level medium-resolution representations [99], or iterating the processes [76]. Similar techniques have also been applied to edge detection, e.g., holistic edge detection [157].

The fully convolutional network is extended, by replacing a few (typically two) strided convolutions and the associated convolutions with dilated convolutions, to the dilation version, leading to medium-resolution representations [18], [19], [86], [168], [181]. The representations are further augmented to multi-scale contextual representations [19], [21], [181] through feature pyramids for segmenting objects at multiple scales.

**Recovering high-resolution representations.** An upsample process can be used to gradually recover the high-resolution representations from the low-resolution representations. The upsample subnetwork could be a symmetric version of the downsample process (e.g., VGGNet), with skipping connection over some mirrored layers to transform the pooling indices, e.g., SegNet [3] and DeconvNet [107], or copying the feature maps, e.g., U-Net [119] and Hourglass [8], [9], [27], [31], [68], [105], [134], [163], [165], encoder-decoder [112], and so on. An extension of U-Net, full-resolution residual

network [114], introduces an extra full-resolution stream that carries information at the full image resolution, to replace the skip connections, and each unit in the downsample and upsample subnetworks receives information from and sends information to the full-resolution stream.

The asymmetric upsample process is also widely studied. RefineNet [90] improves the combination of upsampled representations and the representations of the same resolution copied from the downsample process. Other works include: light upsample process [7], [24], [92], [152], possibly with dilated convolutions used in the backbone [63], [89], [113]; light downsample and heavy upsample processes [141], recombinator networks [55]; improving skip connections with more or complicated convolutional units [64], [111], [180], as well as sending information from low-resolution skip connections to high-resolution skip connections [189] or exchanging information between them [49]; studying the details of the upsample process [147]; combining multi-scale pyramid representations [22], [154]; stacking multiple DeconvNets/U-Nets/Hourglass [44], [149] with dense connections [135].

**Maintaining high-resolution representations.** Our work is closely related to several works that can also generate high-resolution representations, e.g., convolutional neural fabrics [123], interlinked CNNs [188], GridNet [42], and multi-scale DenseNet [58].

The two early works, convolutional neural fabrics [123] and interlinked CNNs [188], lack careful design on when to start low-resolution parallel streams, and how and where to exchange information across parallel streams, and do not use batch normalization and residual connections, thus not showing satisfactory performance. GridNet [42] is like a combination of multiple U-Nets and includes two symmetric information exchange stages: the first stage passes information only from high resolution to low resolution, and the second stage passes information only from low resolution to high resolution. This limits its segmentation quality. Multi-scale DenseNet [58] is not able to learn strong high-resolution representations as there is no information received from low-resolution representations.

**Multi-scale fusion.** Multi-scale fusion<sup>1</sup> is widely studied [11], [19], [24], [42], [58], [66], [122], [123], [157], [161], [181], [188]. The straightforward way is to feed multi-resolution images separately into multiple networks and aggregate the output response maps [137]. Hourglass [105], U-Net [119], and SegNet [3] combine low-level features in the high-to-low downsample process into the same-resolution high-level features in the low-to-high upsample process progressively through skip connections. PSPNet [181] and

1. In this paper, Multi-scale fusion and multi-resolution fusion are interchangeable, but in other contexts, they may not be interchangeable.

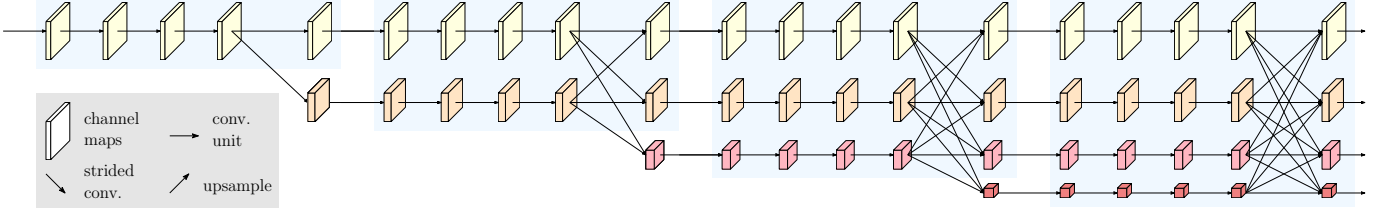


Fig. 2. An example of a high-resolution network. Only the main body is illustrated, and the stem (two stride-2  $3 \times 3$  convolutions) is not included. There are four stages. The 1st stage consists of high-resolution convolutions. The 2nd (3rd, 4th) stage repeats two-resolution (three-resolution, four-resolution) blocks. The detail is given in Section 3.

DeepLabV2/3 [19] fuse the pyramid features obtained by pyramid pooling module and atrous spatial pyramid pooling. Our multi-scale (resolution) fusion module resembles the two pooling modules. The differences include: (1) Our fusion outputs four-resolution representations other than only one, and (2) our fusion modules are repeated several times which is inspired by deep fusion [129], [143], [155], [178], [184].

**Our approach.** Our network connects high-to-low convolution streams in parallel. It maintains high-resolution representations through the whole process, and generates reliable high-resolution representations with strong position sensitivity through repeatedly fusing the representations from multi-resolution streams.

This paper represents a very substantial extension of our previous conference paper [130] with an additional material added from our unpublished technical report [131] as well as more object detection results under recently-developed start-of-the-art object detection and instance segmentation frameworks. The main technical novelties compared with [130] lie in threefold. (1) We extend the network (named as HRNetV1) proposed in [130], to two versions: HRNetV2 and HRNetV2p, which explore all the four-resolution representations. (2) We build the connection between multi-resolution fusion and regular convolution, which provides an evidence for the necessity of exploring all the four-resolution representations in HRNetV2 and HRNetV2p. (3) We show the superiority of HRNetV2 and HRNetV2p over HRNetV1 and present the applications of HRNetV2 and HRNetV2p in a broad range of vision problems, including semantic segmentation and object detection.

### 3 HIGH-RESOLUTION NETWORKS

We input the image into a stem, which consists of two stride-2  $3 \times 3$  convolutions decreasing the resolution to  $\frac{1}{4}$ , and subsequently the main body that outputs the representation with the same resolution ( $\frac{1}{4}$ ). The main body, illustrated in Figure 2 and detailed below, consists of several components: parallel multi-resolution convolutions, repeated multi-resolution fusions, and representation head that is shown in Figure 4.

#### 3.1 Parallel Multi-Resolution Convolutions

We start from a high-resolution convolution stream as the first stage, gradually add high-to-low resolution streams one by one, forming new stages, and connect the multi-resolution streams in parallel. As a result, the resolutions for the parallel streams of a later stage consists of the resolutions from the previous stage, and an extra lower one.

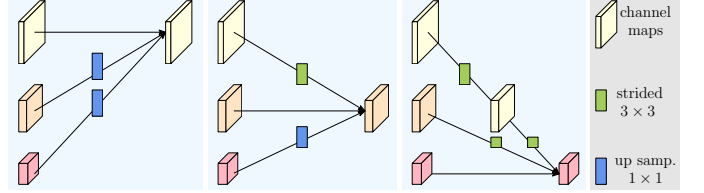


Fig. 3. Illustrating how the fusion module aggregates the information for high, medium and low resolutions from left to right, respectively. Right legend: strided  $3 \times 3 =$  stride-2  $3 \times 3$  convolution, up samp.  $1 \times 1 =$  bilinear upsampling followed by a  $1 \times 1$  convolution.

An example network structure illustrated in Figure 2, containing 4 parallel streams, is logically as follows,

$$\begin{array}{ccccccc} \mathcal{N}_{11} & \rightarrow & \mathcal{N}_{21} & \rightarrow & \mathcal{N}_{31} & \rightarrow & \mathcal{N}_{41} \\ & \searrow & \mathcal{N}_{22} & \rightarrow & \mathcal{N}_{32} & \rightarrow & \mathcal{N}_{42} \\ & & & \searrow & \mathcal{N}_{33} & \rightarrow & \mathcal{N}_{43} \\ & & & & & \searrow & \mathcal{N}_{44}, \end{array} \quad (1)$$

where  $\mathcal{N}_{sr}$  is a sub-stream in the  $s$ th stage and  $r$  is the resolution index. The resolution index of the first stream is  $r = 1$ . The resolution of index  $r$  is  $\frac{1}{2^{r-1}}$  of the resolution of the first stream.

#### 3.2 Repeated Multi-Resolution Fusions

The goal of the fusion module is to exchange the information across multi-resolution representations. It is repeated several times (e.g., every 4 residual units).

Let us look at an example of fusing 3-resolution representations, which is illustrated in Figure 3. Fusing 2 representations and 4 representations can be easily derived. The input consists of three representations:  $\{\mathbf{R}_r^i, r = 1, 2, 3\}$ , with  $r$  is the resolution index, and the associated output representations are  $\{\mathbf{R}_r^o, r = 1, 2, 3\}$ . Each output representation is the sum of the transformed representations of the three inputs:  $\mathbf{R}_r^o = f_{1r}(\mathbf{R}_1^i) + f_{2r}(\mathbf{R}_2^i) + f_{3r}(\mathbf{R}_3^i)$ . The fusion across stages (from stage 3 to stage 4) has an extra output:  $\mathbf{R}_4^o = f_{14}(\mathbf{R}_1^i) + f_{24}(\mathbf{R}_2^i) + f_{34}(\mathbf{R}_3^i)$ .

The choice of the transform function  $f_{xr}(\cdot)$  is dependent on the input resolution index  $x$  and the output resolution index  $r$ . If  $x = r$ ,  $f_{xr}(\mathbf{R}) = \mathbf{R}$ . If  $x < r$ ,  $f_{xr}(\mathbf{R})$  downsamples the input representation  $\mathbf{R}$  through  $(r - x)$  stride-2  $3 \times 3$  convolutions. For instance, one stride-2  $3 \times 3$  convolution for  $2\times$  downsampling, and two consecutive stride-2  $3 \times 3$  convolutions for  $4\times$  downsampling. If  $x > r$ ,  $f_{xr}(\mathbf{R})$  upsamples the input representation  $\mathbf{R}$  through the bilinear upsampling followed by a  $1 \times 1$  convolution for aligning the number of channels. The functions are depicted in Figure 3.

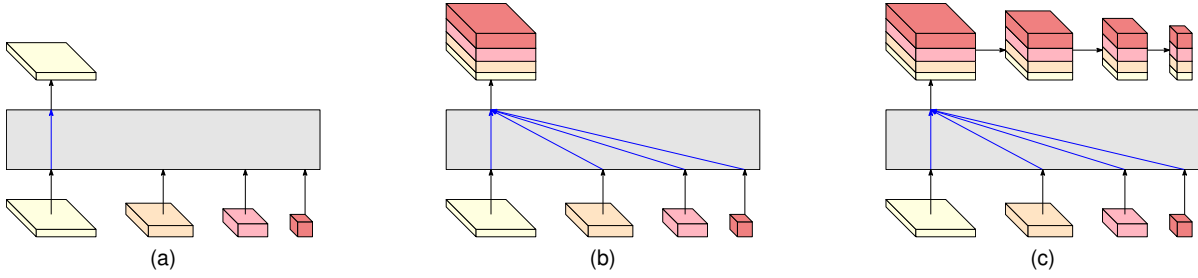


Fig. 4. (a) HRNetV1: only output the representation from the high-resolution convolution stream. (b) HRNetV2: Concatenate the (upsampled) representations that are from all the resolutions (the subsequent  $1 \times 1$  convolution is not shown for clarity). (c) HRNetV2p: form a feature pyramid from the representation by HRNetV2. The four-resolution representations at the bottom in each sub-figure are outputted from the network in Figure 2, and the gray box indicates how the output representation is obtained from the input four-resolution representations.

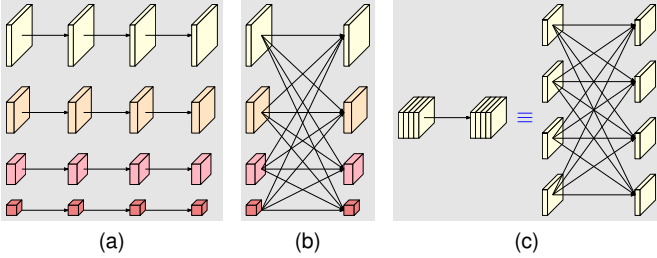


Fig. 5. (a) Multi-resolution parallel convolution, (b) multi-resolution fusion. (c) A normal convolution (left) is equivalent to fully-connected multi-branch convolutions (right).

### 3.3 Representation Head

We have three kinds of representation heads that are illustrated in Figure 4, and call them as HRNetV1, HRNetV2, and HRNetV1p, respectively.

**HRNetV1.** The output is the representation only from the high-resolution stream. Other three representations are ignored. This is illustrated in Figure 4 (a).

**HRNetV2.** We rescale the low-resolution representations through bilinear upsampling without changing the number of channels to the high resolution, and concatenate the four representations, followed by a  $1 \times 1$  convolution to mix the four representations. This is illustrated in Figure 4 (b).

**HRNetV2p.** We construct multi-level representations by downsampling the high-resolution representation output from HRNetV2 to multiple levels. This is depicted in Figure 4 (c).

In this paper, we will show the results of applying HRNetV1 to human pose estimation, HRNetV2 to semantic segmentation, and HRNetV2p to object detection.

### 3.4 Instantiation

The main body contains four stages with four parallel convolution streams. The resolutions are  $1/4$ ,  $1/8$ ,  $1/16$ , and  $1/32$ . The first stage contains 4 residual units where each unit is formed by a bottleneck with the width 64, and is followed by one  $3 \times 3$  convolution changing the width of feature maps to  $C$ . The 2nd, 3rd, 4th stages contain 1, 4, 3 modularized blocks, respectively. Each branch in multi-resolution parallel convolution of the modularized block contains 4 residual units. Each unit contains two  $3 \times 3$  convolutions for each resolution, where each convolution is followed by batch normalization and the nonlinear activation ReLU. The widths (numbers of channels) of the

convolutions of the four resolutions are  $C$ ,  $2C$ ,  $4C$ , and  $8C$ , respectively. An example is depicted in Figure 2.

### 3.5 Analysis

We analyze the modularized block that is divided into two components: multi-resolution parallel convolutions (Figure 5 (a)), and multi-resolution fusion (Figure 5 (b)). The multi-resolution parallel convolution resembles the group convolution. It divides the input channels into several subsets of channels and performs a regular convolution over each subset over different spatial resolutions separately, while in the group convolution, the resolutions are the same. This connection implies that the multi-resolution parallel convolution enjoys some benefit of the group convolution.

The multi-resolution fusion unit resembles the multi-branch full-connection form of the regular convolution, illustrated in Figure 5 (c). A regular convolution can be divided as multiple small convolutions as explained in [178]. The input channels are divided into several subsets, and the output channels are also divided into several subsets. The input and output subsets are connected in a fully-connected fashion, and each connection is a regular convolution. Each subset of output channels is a summation of the outputs of the convolutions over each subset of input channels. The differences lie in that our multi-resolution fusion needs to handle the resolution change. The connection between multi-resolution fusion and regular convolution provides an evidence for exploring all the four-resolution representations done in HRNetV2 and HRNetV2p.

## 4 HUMAN POSE ESTIMATION

Human pose estimation, a.k.a. keypoint detection, aims to detect the locations of  $K$  keypoints or parts (e.g., elbow, wrist, etc) from an image  $I$  of size  $W \times H \times 3$ . We follow the state-of-the-art framework and transform this problem to estimating  $K$  heatmaps of size  $\frac{W}{4} \times \frac{H}{4}$ ,  $\{\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_K\}$ , where each heatmap  $\mathbf{H}_k$  indicates the location confidence of the  $k$ th keypoint.

We regress the heatmaps over the high-resolution representations output by HRNetV1. We empirically observe that the performance is almost the same for HRNetV1 and HRNetV2, and thus we choose HRNetV1 as its computation complexity is a little lower. The loss function, defined as the mean squared error, is applied for comparing the predicted heatmaps and the groundtruth heatmaps. The groundtruth heatmaps are generated by applying 2D Gaussian with





Fig. 6. Qualitative COCO human pose estimation results over representative images with various human size, different poses, or clutter background.

TABLE 1

Comparisons on COCO val. Under the input size  $256 \times 192$ , our approach with a small model HRNetV1-W32, trained from scratch, performs better than previous state-of-the-art methods. Under the input size  $384 \times 288$ , our approach with a small model HRNetV1-W32 achieves a higher AP score than SimpleBaseline with a large model. In particular, the improvement of our approach for  $AP^{75}$ , a strict evaluation scheme, is more significant than  $AP^{50}$ , a loose evaluation scheme. Pretrain = pretrain the backbone on ImageNet. OHKM = online hard keypoints mining [24]. #Params and FLOPs are calculated for the pose estimation network, and those for human detection and keypoint grouping are not included.

Method	Backbone	Pretrain	Input size	#Params	GFLOPs	AP	$AP^{50}$	$AP^{75}$	$AP^M$	$AP^L$	AR
8-stage Hourglass [105]	8-stage Hourglass	N	$256 \times 192$	25.1M	14.3	66.9	—	—	—	—	—
CPN [24]	ResNet-50	Y	$256 \times 192$	27.0M	6.20	68.6	—	—	—	—	—
CPN + OHKM [24]	ResNet-50	Y	$256 \times 192$	27.0M	6.20	69.4	—	—	—	—	—
SimpleBaseline [152]	ResNet-50	Y	$256 \times 192$	34.0M	8.90	70.4	88.6	78.3	67.1	77.2	76.3
SimpleBaseline [152]	ResNet-101	Y	$256 \times 192$	53.0M	12.4	71.4	89.3	79.3	68.1	78.1	77.1
SimpleBaseline [152]	ResNet-152	Y	$256 \times 192$	68.6M	15.7	72.0	89.3	79.8	68.7	78.9	77.8
HRNetV1	HRNetV1-W32	N	$256 \times 192$	28.5M	7.10	73.4	89.5	80.7	70.2	80.1	78.9
HRNetV1	HRNetV1-W32	Y	$256 \times 192$	28.5M	7.10	74.4	90.5	81.9	70.8	81.0	79.8
HRNetV1	HRNetV1-W48	Y	$256 \times 192$	63.6M	14.6	75.1	90.6	82.2	71.5	81.8	80.4
SimpleBaseline [152]	ResNet-152	Y	$384 \times 288$	68.6M	35.6	74.3	89.6	81.1	70.5	79.7	79.7
HRNetV1	HRNetV1-W32	Y	$384 \times 288$	28.5M	16.0	75.8	90.6	82.7	71.9	82.8	81.0
HRNetV1	HRNetV1-W48	Y	$384 \times 288$	63.6M	32.9	<b>76.3</b>	<b>90.8</b>	<b>82.9</b>	<b>72.3</b>	<b>83.4</b>	<b>81.2</b>

TABLE 2

Comparisons on COCO test-dev. The observations are similar to the results on COCO val.

Method	Backbone	Input size	#Params	GFLOPs	AP	$AP^{50}$	$AP^{75}$	$AP^M$	$AP^L$	AR
Bottom-up: keypoint detection and grouping										
OpenPose [15]	—	—	—	—	61.8	84.9	67.5	57.1	68.2	66.5
Associative Embedding [104]	—	—	—	—	65.5	86.8	72.3	60.6	72.6	70.2
PersonLab [108]	—	—	—	—	68.7	89.0	75.4	64.1	75.5	75.4
MultiPoseNet [72]	—	—	—	—	69.6	86.3	76.6	65.0	76.3	73.5
Top-down: human detection and single-person keypoint detection										
Mask-RCNN [53]	ResNet-50-FPN	—	—	—	63.1	87.3	68.7	57.8	71.4	—
G-RMI [109]	ResNet-101	$353 \times 257$	42.6M	57.0	64.9	85.5	71.3	62.3	70.0	69.7
Integral Pose Regression [132]	ResNet-101	$256 \times 256$	45.0M	11.0	67.8	88.2	74.8	63.9	74.0	—
G-RMI + extra data [109]	ResNet-101	$353 \times 257$	42.6M	57.0	68.5	87.1	75.5	65.8	73.3	73.3
CPN [24]	ResNet-Inception	$384 \times 288$	—	—	72.1	91.4	80.0	68.7	77.2	78.5
RMPE [38]	PyraNet [165]	$320 \times 256$	28.1M	26.7	72.3	89.2	79.1	68.0	78.6	—
CFN [60]	—	—	—	—	72.6	86.1	69.7	78.3	64.1	—
CPN (ensemble) [24]	ResNet-Inception	$384 \times 288$	—	—	73.0	91.7	80.9	69.5	78.1	79.0
SimpleBaseline [152]	ResNet-152	$384 \times 288$	68.6M	35.6	73.7	91.9	81.1	70.3	80.0	79.0
HRNetV1	HRNetV1-W32	$384 \times 288$	28.5M	16.0	74.9	92.5	82.8	71.3	80.9	80.1
HRNetV1	HRNetV1-W48	$384 \times 288$	63.6M	32.9	75.5	92.5	83.3	71.9	81.5	80.5
HRNetV1 + extra data	HRNetV1-W48	$384 \times 288$	63.6M	32.9	<b>77.0</b>	<b>92.7</b>	<b>84.5</b>	<b>73.4</b>	<b>83.1</b>	<b>82.0</b>

standard deviation of 2 pixel centered on the groundtruth location of each keypoint. Some example results are given in Figure 6.

**Dataset.** The COCO dataset [94] contains over 200,000 images and 250,000 person instances labeled with 17 keypoints. We train our model on the COCO train2017 set, including 57K images and 150K person instances. We evaluate our approach on the val2017 and test-dev2017 sets, containing 5000 images and 20K images, respectively.

**Evaluation metric.** The standard evaluation metric is based on Object Keypoint Similarity (OKS):  $OKS = \frac{\sum_i \exp(-d_i^2 / 2s^2 k_i^2) \delta(v_i > 0)}{\sum_i \delta(v_i > 0)}$ . Here  $d_i$  is the Euclidean distance between the detected keypoint and the corresponding

ground truth,  $v_i$  is the visibility flag of the ground truth,  $s$  is the object scale, and  $k_i$  is a per-keypoint constant that controls falloff. We report standard average precision and recall scores<sup>2</sup>:  $AP^{50}$  (AP at OKS = 0.50),  $AP^{75}$ , AP (the mean of AP scores at 10 OKS positions, 0.50, 0.55, ..., 0.90, 0.95);  $AP^M$  for medium objects,  $AP^L$  for large objects, and AR (the mean of AR scores at 10 OKS positions, 0.50, 0.55, ..., 0.90, 0.95).

**Training.** We extend the human detection box in height or width to a fixed aspect ratio: height : width = 4 : 3, and then crop the box from the image, which is resized to a fixed size,  $256 \times 192$  or  $384 \times 288$ . The data augmenta-

2. <http://cocodataset.org/#keypoints-eval>

tion scheme includes random rotation ( $[-45^\circ, 45^\circ]$ ), random scale ( $[0.65, 1.35]$ ), and flipping. Following [146], half body data augmentation is also involved.

We use the Adam optimizer [71]. The learning schedule follows the setting [152]. The base learning rate is set as  $1e-3$ , and is dropped to  $1e-4$  and  $1e-5$  at the 170th and 200th epochs, respectively. The training process is terminated within 210 epochs. The models are trained on 4 V100 GPUs and it takes around 60 (80) hours for HRNet-W32 (HRNet-W48).

**Testing.** The two-stage top-down paradigm similar as [24], [109], [152] is used: detect the person instance using a person detector, and then predict detection keypoints.

We use the same person detectors provided by SimpleBaseline<sup>3</sup> for both the val and test-dev sets. Following [24], [105], [152], we compute the heatmap by averaging the heatmaps of the original and flipped images. Each keypoint location is predicted by adjusting the highest heatmap location with a quarter offset in the direction from the highest response to the second highest response.

**Results on the val set.** We report the results of our method and other state-of-the-art methods in Table 1. The network - HRNetV1-W32, trained from scratch with the input size  $256 \times 192$ , achieves an AP score 73.4, outperforming other methods with the same input size. (i) Compared to Hourglass [105], our network improves AP by 6.5 points, and the GFLOP of our network is much lower and less than half, while the numbers of parameters are similar and ours is slightly larger. (ii) Compared to CPN [24] w/o and w/ OHKM, our network, with slightly larger model size and slightly higher complexity, achieves 4.8 and 4.0 points gain, respectively. (iii) Compared to the previous best-performed method SimpleBaseline [152], our HRNetV1-W32 obtains significant improvements: 3.0 points gain for the backbone ResNet-50 with a similar model size and GFLOPs, and 1.4 points gain for the backbone ResNet-152 whose model size (#Params) and GFLOPs are twice as many as ours.

Our network can benefit from (i) training from the model pretrained on the ImageNet: The gain is 1.0 points for HRNetV1-W32; (ii) increasing the capacity by increasing the width: HRNetV1-W48 gets 0.7 and 0.5 points gain for the input sizes  $256 \times 192$  and  $384 \times 288$ , respectively.

Considering the input size  $384 \times 288$ , our HRNetV1-W32 and HRNetV1-W48, get the 75.8 and 76.3 AP, which have 1.4 and 1.2 improvements compared to the input size  $256 \times 192$ . In comparison to SimpleBaseline [152] that uses ResNet-152 as the backbone, our HRNetV1-W32 and HRNetV1-W48 attain 1.5 and 2.0 points gain in terms of AP at 45% and 92.4% computational cost, respectively.

**Results on the test-dev set.** Table 2 reports the pose estimation performances of our approach and the existing state-of-the-art approaches. Our approach is significantly better than bottom-up approaches. On the other hand, our small network, HRNetV1-W32, achieves an AP of 74.9. It outperforms all the other top-down approaches, and is more efficient in terms of model size (#Params) and computation complexity (GFLOPs). Our big model, HRNetV1-W48, achieves the highest AP score 75.5. Compared to

TABLE 3

Semantic segmentation results on Cityscapes val (single scale and no flipping). The GFLOPs is calculated on the input size  $1024 \times 2048$ . The small model HRNetV2-W40 with the smallest GFLOPs performs better than two representative contextual methods (DeepLab and PSPNet). Our approach combined with the recently-developed object contextual (OCR) representation scheme [170] gets further improvement. D-ResNet-101 = Dilated-ResNet-101.

	backbone	#param.	GFLOPs	mIoU
UNet++ [189]	ResNet-101	59.5M	748.5	75.5
Dilated-ResNet [54]	D-ResNet-101	52.1M	1661.6	75.7
DeepLabv3 [20]	D-ResNet-101	58.0M	1778.7	78.5
DeepLabv3+ [22]	D-Xception-71	43.5M	1444.6	79.6
PSPNet [181]	D-ResNet-101	65.9M	2017.6	79.7
HRNetV2	HRNetV2-W40	45.2M	493.2	80.2
HRNetV2	HRNetV2-W48	65.9M	696.2	81.1
HRNetV2 + OCR [170]	HRNetV2-W48	70.3M	1206.3	<b>81.6</b>

TABLE 4

Semantic segmentation results on Cityscapes test. We use HRNetV2-W48, whose parameter complexity and computation complexity are comparable to dilated-ResNet-101 based networks, for comparison. Our results are superior in terms of the four evaluation metrics. The result from the combination with OCR [170] is further improved. D-ResNet-101 = Dilated-ResNet-101.

	backbone	mIoU	iIoU cla.	IoU cat.	iIoU cat.
<i>Model learned on the train set</i>					
PSPNet [181]	D-ResNet-101	78.4	56.7	90.6	78.6
PSANet [182]	D-ResNet-101	78.6	-	-	-
PAN [82]	D-ResNet-101	78.6	-	-	-
AAF [69]	D-ResNet-101	79.1	-	-	-
HRNetV2	HRNetV2-W48	<b>80.4</b>	<b>59.2</b>	<b>91.5</b>	<b>80.8</b>
<i>Model learned on the train+val set</i>					
GridNet [42]	-	69.5	44.1	87.9	71.1
LRR-4x [46]	-	69.7	48.0	88.2	74.7
DeepLab [19]	D-ResNet-101	70.4	42.6	86.4	67.7
LC [84]	-	71.1	-	-	-
Piecewise [91]	VGG-16	71.6	51.7	87.3	74.1
FRRN [114]	-	71.8	45.5	88.9	75.1
RefineNet [90]	ResNet-101	73.6	47.2	87.9	70.6
PEARL [65]	D-ResNet-101	75.4	51.6	89.2	75.1
DSSPN [88]	D-ResNet-101	76.6	56.2	89.6	77.8
LKM [111]	ResNet-152	76.9	-	-	-
DUC-HDC [144]	-	77.6	53.6	90.1	75.2
SAC [176]	D-ResNet-101	78.1	-	-	-
DepthSeg [73]	D-ResNet-101	78.2	-	-	-
ResNet38 [151]	WResNet-38	78.4	59.1	90.9	78.1
BiSeNet [166]	ResNet-101	78.9	-	-	-
DFN [167]	ResNet-101	79.3	-	-	-
PSANet [182]	D-ResNet-101	80.1	-	-	-
PADNet [159]	D-ResNet-101	80.3	58.8	90.8	78.5
CFNet [173]	D-ResNet-101	79.6	-	-	-
Auto-DeepLab [95]	-	80.4	-	-	-
DenseASPP [181]	WDenseNet-161	80.6	59.1	90.9	78.1
SVCNet [33]	ResNet-101	81.0	-	-	-
ANN [195]	D-ResNet-101	81.3	-	-	-
CCNet [61]	D-ResNet-101	81.4	-	-	-
DANet [43]	D-ResNet-101	81.5	-	-	-
HRNetV2	HRNetV2-W48	81.6	<b>61.8</b>	<b>92.1</b>	<b>82.2</b>
HRNetV2 + OCR [170]	HRNetV2-W48	<b>82.5</b>	61.7	<b>92.1</b>	81.6

SimpleBaseline [152] with the same input size, our small and big networks receive 1.2 and 1.8 improvements, respectively. With the additional data from AI Challenger [148] for training, our single big network can obtain an AP of 77.0.

3. <https://github.com/Microsoft/human-pose-estimation.pytorch>



Fig. 7. Qualitative segmentation examples from Cityscapes (left two), PASCAL-Context (middle two), and LIP (right two).

TABLE 5

Semantic segmentation results on PASCAL-Context. The methods are evaluated on 59 classes and 60 classes. Our approach performs the best for 60 classes, and performs worse for 59 classes than APCN [51] that developed a strong contextual method. Our approach, combined with OCR [170], achieves significant gain, and performs the best. D-ResNet-101 = Dilated-ResNet-101.

	backbone	mIoU (59)	mIoU (60)
FCN-8s [125]	VGG-16	-	35.1
BoxSup [29]	-	-	40.5
HO_CRF [2]	-	-	41.3
Piecewise [91]	VGG-16	-	43.3
DeepLab-v2 [19]	D-ResNet-101	-	45.7
RefineNet [90]	ResNet-152	-	47.3
UNet++ [189]	ResNet-101	47.7	-
PSPNet [181]	D-ResNet-101	47.8	-
Ding et al. [32]	ResNet-101	51.6	-
EncNet [172]	D-ResNet-101	52.6	-
DANet [43]	D-ResNet-101	52.6	-
ANN [195]	D-ResNet-101	52.8	-
SVCNet [33]	ResNet-101	53.2	-
CFNet [173]	D-ResNet-101	54.0	-
APCN [51]	D-ResNet-101	55.6	-
HRNetV2	HRNetV2-W48	54.0	48.3
HRNetV2 + OCR [170]	HRNetV2-W48	<b>56.2</b>	<b>50.1</b>

TABLE 6

Semantic segmentation results on LIP. Our method doesn't exploit any extra information, e.g., pose or edge. The overall performance of our approach is the best, and the OCR scheme [170] further improves the segmentation quality. D-ResNet-101 = Dilated-ResNet-101.

	backbone	extra.	pixel acc.	avg. acc.	mIoU
Attention+SSL [47]	VGG16	Pose	84.36	54.94	44.73
DeepLabV3+ [22]	D-ResNet-101	-	84.09	55.62	44.80
MMAN [100]	D-ResNet-101	-	-	-	46.81
SS-NAN [183]	ResNet-101	Pose	87.59	56.03	47.92
MuLA [106]	Hourglass	Pose	<b>88.50</b>	60.50	49.30
JPPNet [87]	D-ResNet-101	Pose	86.39	62.32	51.37
CE2P [98]	D-ResNet-101	Edge	87.37	63.20	53.10
HRNetV2	HRNetV2-W48	N	88.21	67.43	55.90
HRNetV2 + OCR [170]	HRNetV2-W48	N	88.24	<b>67.84</b>	<b>56.48</b>

## 5 SEMANTIC SEGMENTATION

Semantic segmentation is a problem of assigning a class label to each pixel. Some example results by our approach are given in Figure 7. We feed the input image to the HRNetV2 (Figure 4 (b)) and then pass the resulting 15C-dimensional representation at each position to a linear classifier with the softmax loss to predict the segmentation maps. The segmentation maps are upsampled (4 times) to the input size by bilinear upsampling for both training and testing. We report the results over two scene parsing datasets, PASCAL-Context [103] and Cityscapes [28], and a human parsing dataset, LIP [47]. The mean of class-wise intersection over union (mIoU) is adopted as the evaluation metric.

**Cityscapes.** The Cityscapes dataset [28] contains 5,000 high

quality pixel-level finely annotated scene images. The finely-annotated images are divided into 2,975/500/1,525 images for training, validation and testing. There are 30 classes, and 19 classes among them are used for evaluation. In addition to the mean of class-wise intersection over union (mIoU), we report other three scores on the test set: IoU category (cat.), iIoU class (cla.) and iIoU category (cat.).

We follow the same training protocol [181], [182]. The data are augmented by random cropping (from  $1024 \times 2048$  to  $512 \times 1024$ ), random scaling in the range of  $[0.5, 2]$ , and random horizontal flipping. We use the SGD optimizer with the base learning rate of 0.01, the momentum of 0.9 and the weight decay of 0.0005. The poly learning rate policy with the power of 0.9 is used for dropping the learning rate. All the models are trained for 120K iterations with the batch size of 12 on 4 GPUs and syncBN.

Table 3 provides the comparison with several representative methods on the Cityscapes val set in terms of parameter and computation complexity and mIoU class. (i) HRNetV2-W40 (40 indicates the width of the high-resolution convolution), with similar model size to DeepLabv3+ and much lower computation complexity, gets better performance: 4.7 points gain over UNet++, 1.7 points gain over DeepLabv3 and about 0.5 points gain over PSPNet, DeepLabv3+. (ii) HRNetV2-W48, with similar model size to PSPNet and much lower computation complexity, achieves much significant improvement: 5.6 points gain over UNet++, 2.6 points gain over DeepLabv3 and about 1.4 points gain over PSPNet, DeepLabv3+. In the following comparisons, we adopt HRNetV2-W48 that is pretrained on ImageNet and has similar model size as most Dilated-ResNet-101 based methods.

Table 4 provides the comparison of our method with state-of-the-art methods on the Cityscapes test set. All the results are with six scales and flipping. Two cases w/o using coarse data are evaluated: One is about the model learned on the train set, and the other is about the model learned on the train+val set. In both cases, HRNetV2-W48 achieves the superior performance.

**PASCAL-Context.** The PASCAL-Context dataset [103] includes 4,998 scene images for training and 5,105 images for testing with 59 semantic labels and 1 background label.

The data augmentation and learning rate policy are the same as Cityscapes. Following the widely-used training strategy [32], [172], we resize the images to  $480 \times 480$  and set the initial learning rate to 0.004 and weight decay to 0.0001. The batch size is 16 and the number of iterations is 60K.

We follow the standard testing procedure [32], [172]. The image is resized to  $480 \times 480$  and then fed into our network. The resulting  $480 \times 480$  label maps are then resized to the original image size. We evaluate the performance of our approach and other approaches using six scales and flipping.



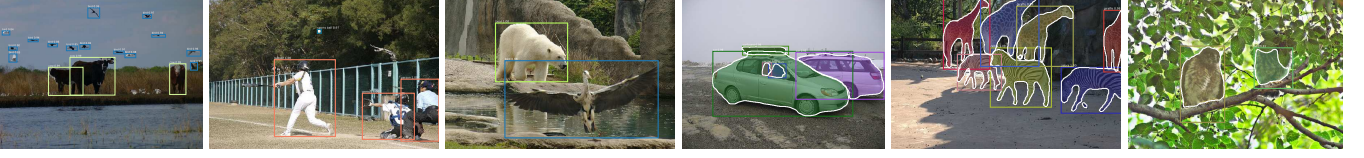


Fig. 8. Qualitative examples for COCO object detection (left three) and instance segmentation (right three).

TABLE 7

GFLOPs and #parameters for COCO object detection. The numbers are obtained with the input size  $800 \times 1200$  and if applicable 512 proposals fed into R-CNN except the numbers for CenterNet are obtained with the input size  $511 \times 511$ . R- $x$  = ResNet- $x$ -FPN, X-101 = ResNeXt-101-64 $\times$ 4d, H- $x$  = HRNetV2p-W $x$ , and HG-52 = Hourglass-52.

	Faster R-CNN [53]						Cascade R-CNN [13]						FCOS [136]				CenterNet [36]			
	R-50	H-18	R-101	H-32	X-101	H-48	R-50	H-18	R-101	H-32	X-101	H-48	R-50	H-18	R-101	H-32	HG-52	H-48	HG-104	H-64
#param. (M)	39.8	26.2	57.8	45.0	94.9	79.4	69.4	55.1	88.4	74.9	127.3	111.0	32.0	17.5	51.0	37.3	104.8	73.6	210.1	127.7
GFLOPs	172.3	159.1	239.4	245.3	381.8	399.1	226.2	207.8	298.7	300.8	448.3	466.5	190.0	180.3	261.2	273.3	227.0	217.1	388.4	318.5
	Cascade Mask R-CNN [13]						Hybrid Task Cascade [16]						Mask R-CNN [53]							
	R-50	H-18	R-101	H-32	X-101	H-48	R-50	H-18	R-101	H-32	X-101	H-48	R-50	H-18	R-101	H-32				
#param. (M)	77.3	63.1	96.3	82.9	135.2	118.9	80.3	66.1	99.3	85.9	138.2	121.9	44.4	30.1	63.4	49.9				
GFLOPs	431.7	413.1	504.1	506.2	653.7	671.9	476.9	458.3	549.2	551.4	698.9	717.0	266.5	247.9	338.8	341.0				

TABLE 8

Object detection results on COCO *val* in the Faster R-CNN and Cascade R-CNN frameworks. LS = learning schedule.  $1\times = 12e$ ,  $2\times = 24e$ . Our approach performs better than ResNet and ResNeXt. Our approach gets more significant improvement for  $2\times$  than  $1\times$  and for small objects ( $AP_S$ ) than medium ( $AP_M$ ) and large objects ( $AP_L$ ).

backbone	LS	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
Faster R-CNN [92]							
ResNet-50-FPN	$1\times$	36.7	58.3	39.9	20.9	39.8	47.9
HRNetV2p-W18	$1\times$	36.2	57.3	39.3	20.7	39.0	46.8
ResNet-50-FPN	$2\times$	37.6	58.7	41.3	21.4	40.8	49.7
HRNetV2p-W18	$2\times$	38.0	58.9	41.5	22.6	40.8	49.6
ResNet-101-FPN	$1\times$	39.2	61.1	43.0	22.3	42.9	50.9
HRNetV2p-W32	$1\times$	39.6	61.0	43.3	23.7	42.5	50.5
ResNet-101-FPN	$2\times$	39.8	61.4	43.4	22.9	43.6	52.4
HRNetV2p-W32	$2\times$	40.9	61.8	44.8	24.4	43.7	53.3
X-101-64 $\times$ 4d-FPN	$1\times$	41.3	63.4	45.2	24.5	45.8	53.3
HRNetV2p-W48	$1\times$	41.3	62.8	45.1	25.1	44.5	52.9
X-101-64 $\times$ 4d-FPN	$2\times$	40.8	62.1	44.6	23.2	44.5	53.7
HRNetV2p-W48	$2\times$	41.8	62.8	45.9	25.0	44.7	54.6
Cascade R-CNN [13]							
ResNet-50-FPN	20e	41.1	59.1	44.8	22.5	44.4	54.9
HRNetV2p-W18	20e	41.3	59.2	44.9	23.7	44.2	54.1
ResNet-101-FPN	20e	42.5	60.7	46.3	23.7	46.1	56.9
HRNetV2p-W32	20e	43.7	61.7	47.7	25.6	46.5	57.4
X-101-64 $\times$ 4d-FPN	20e	44.7	63.1	49.0	25.8	48.3	58.8
HRNetV2p-W48	20e	44.6	62.7	48.7	26.3	48.1	58.5

Table 5 provides the comparison of our method with state-of-the-art methods. There are two kinds of evaluation schemes: mIoU over 59 classes and 60 classes (59 classes + background). In both cases, HRNetV2-W48 achieves state-of-the-art results except that the result from [51] is higher than ours without using the OCR scheme [170].

**LIP.** The LIP dataset [47] contains 50,462 elaborately annotated human images, which are divided into 30,462 training images, and 10,000 validation images. The methods are evaluated on 20 categories (19 human part labels and 1 background label). Following the standard training and testing settings [98], the images are resized to  $473 \times 473$  and the performance is evaluated on the average of the segmentation maps of the original and flipped images.

The data augmentation and learning rate policy are the same as Cityscapes. The training strategy follows the recent setting [98]. We set the initial learning rate to 0.007 and the momentum to 0.9 and the weight decay to 0.0005. The batch

TABLE 9

Object detection results on COCO *val* in the FCOS and CenterNet frameworks. The results are obtained using the implementations provided by the authors. Our approach performs superiorly to ResNet and Hourglass for similar parameter and computation complexity. Our HRNetV2p-W64 performs slightly worse than Hourglass-104, and the reason is that Hourglass-104 is much more heavier than HRNetV2p-W64. See Table 7 for #parameters and GFLOPs.

backbone	LS	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
FCOS [136]							
ResNet-50-FPN	$2\times$	37.1	55.9	39.8	21.3	41.0	47.8
HRNetV2p-W18	$2\times$	37.7	55.3	40.2	22.0	40.8	48.8
ResNet-101-FPN	$2\times$	41.4	60.3	44.8	25.0	45.6	53.1
HRNetV2p-W32	$2\times$	41.9	60.3	45.0	25.1	45.6	53.2
CenterNet [36]							
Hourglass-52	-	41.3	59.2	43.9	23.6	43.8	55.8
HRNetV2p-W48	-	43.4	61.8	45.6	23.8	47.1	59.3
Hourglass-104	-	44.8	62.4	48.2	25.9	48.9	58.8
HRNetV2p-W64	-	44.0	62.5	47.3	23.9	48.2	60.2

size is 40 and the number of iterations is 110K.

Table 6 provides the comparison of our method with state-of-the-art methods. The overall performance of HRNetV2-W48 performs the best with fewer parameters and lighter computation cost. We also would like to mention that our networks do not use extra information such as pose or edge.

## 6 COCO OBJECT DETECTION

We perform the evaluation on the MS COCO 2017 detection dataset, which contains about 118k images for training, 5k for validation (*val*) and  $\sim 20$ k testing without provided annotations (*test-dev*). The standard COCO-style evaluation is adopted. Some example results by our approach are given in Figure 8.

We apply our multi-level representations (HRNetV2p)<sup>4</sup>, shown in Figure 4 (c), for object detection. The data is augmented by standard horizontal flipping. The input images are resized such that the shorter edge is 800 pixels [92]. Inference is performed on a single image scale.

We compare our HRNet with the standard models: ResNet [54] and ResNeXt [156]. We evaluate the de-

4. Same as FPN [93], we also use 5 levels.



TABLE 10

Object detection results on COCO val in the Mask R-CNN and its extended frameworks. The overall performance of our approach is superior to ResNet except that HRNetV2p-W18 sometimes performs worse than ResNet-50. Similar to detection (bbox), the improvement for small objects ( $AP_S$ ) in terms of mask is also more significant than medium ( $AP_M$ ) and large objects ( $AP_L$ ). The results are obtained from MMDetection [17].

backbone	LS	mask				bbox			
		AP	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>	AP	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
Mask R-CNN [53]									
ResNet-50-FPN	1×	34.2	15.7	36.8	50.2	37.8	22.1	40.9	49.3
HRNetV2p-W18	1×	33.8	15.6	35.6	49.8	37.1	21.9	39.5	47.9
ResNet-50-FPN	2×	35.0	16.0	37.5	52.0	38.6	21.7	41.6	50.9
HRNetV2p-W18	2×	35.3	16.9	37.5	51.8	39.2	23.7	41.7	51.0
ResNet-101-FPN	1×	36.1	16.2	39.0	53.0	40.0	22.6	43.4	52.3
HRNetV2p-W32	1×	36.7	17.3	39.0	53.0	40.9	24.5	43.9	52.2
ResNet-101-FPN	2×	36.7	17.0	39.5	54.8	41.0	23.4	44.4	53.9
HRNetV2p-W32	2×	37.6	17.8	40.0	55.0	42.3	25.0	45.4	54.9
Cascade Mask R-CNN [13]									
ResNet-50-FPN	20e	36.6	19.0	37.4	50.7	42.3	23.7	45.7	56.4
HRNetV2p-W18	20e	36.4	17.0	38.6	52.9	41.9	23.8	44.9	55.0
ResNet-101-FPN	20e	37.6	19.7	40.8	52.4	43.3	24.4	46.9	58.0
HRNetV2p-W32	20e	38.5	18.9	41.1	56.1	44.5	26.1	47.9	58.5
X-101-64×4d-FPN	20e	39.4	20.8	42.7	54.1	45.7	26.2	49.6	60.0
HRNetV2p-W48	20e	39.5	19.7	41.8	56.9	46.0	27.5	48.9	60.1
Hybrid Task Cascade [16]									
ResNet-50-FPN	20e	38.1	20.3	41.1	52.8	43.2	24.9	46.4	57.8
HRNetV2p-W18	20e	37.9	18.8	39.9	55.2	43.1	26.6	46.0	56.9
ResNet-101-FPN	20e	39.4	21.4	42.4	54.4	44.9	26.4	48.3	59.9
HRNetV2p-W32	20e	39.6	19.1	42.0	57.9	45.3	27.0	48.4	59.5
X-101-64×4d-FPN	20e	40.8	22.7	44.2	56.3	46.9	28.0	50.7	62.1
HRNetV2p-W48	20e	40.7	19.7	43.4	59.3	46.8	28.0	50.2	61.7
X-101-64×4d-FPN	28e	40.7	20.0	44.1	59.9	46.8	27.5	51.0	61.7
HRNetV2p-W48	28e	41.0	20.8	43.9	59.9	47.0	28.8	50.3	62.2

tection performance on COCO val. under two anchor-based frameworks: Faster R-CNN [118] and Cascade R-CNN [12], and two recently-developed anchor-free frameworks: FCOS [136] and CenterNet [36]. We train the Faster R-CNN and Cascade R-CNN models for both our HRNetV2p and the ResNet on the public MMDetection platform [17] with the provided training setup, except that we use the learning rate schedule suggested in [52] for 2x, and FCOS [136] and CenterNet [36] from the implementations provided by the authors. Table 7 summarizes #parameters and GFLOPs. Table 8 and Table 9 report detection scores.

We also evaluate the performance of joint detection and instance segmentation, under three frameworks: Mask R-CNN [53], Cascade Mask R-CNN [13], and Hybrid Task Cascade [16]. The results are obtained on the public MMDetection platform [17] and are in Table 10.

There are several observations. On the one hand, as shown in Tables 8 and 9, the overall object detection performance of HRNetV2 is better than ResNet under similar model size and computation complexity. In some cases, for 1x, HRNetV2p-W18 performs worse than ResNet-50-FPN, which might come from insufficient optimization iterations. On the other hand, as shown in Table 10, the overall object detection and instance segmentation performance is better than ResNet and ResNeXt. In particular, under the Hybrid Task Cascade framework, the HRNet performs slightly worse than ResNeXt-101-64×4d-FPN for 20e, but better for 28e. This implies that our HRNet benefits more from longer

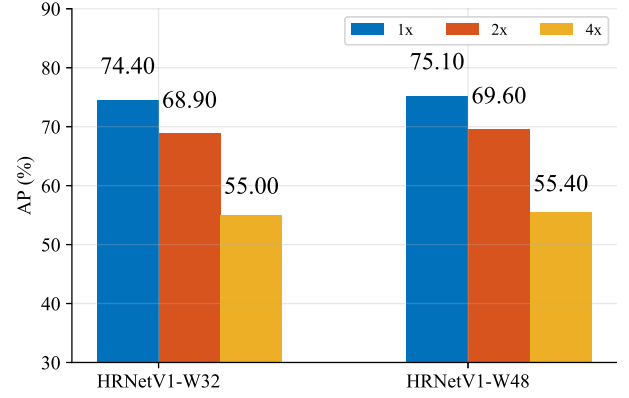


Fig. 9. Ablation study about the resolutions of the representations for human pose estimation. 1x, 2x, 4x correspond to the representations of the high, medium, low resolutions, respectively. The results imply that higher resolution improves the performance.

training.

Table 11 reports the comparison of our network to state-of-the-art single-model object detectors on COCO test-dev without using multi-scale training and multi-scale testing that are done in [85], [97], [110], [115], [127], [128]. In the Faster R-CNN framework, our networks perform better than ResNets with similar parameter and computation complexity: HRNetV2p-W32 vs. ResNet-101-FPN, HRNetV2p-W40 vs. ResNet-152-FPN, HRNetV2p-W48 vs. X-101-64 × 4d-FPN. In the Cascade R-CNN and CenterNet framework, our HRNetV2 also performs better. In the Cascade Mask R-CNN and Hybrid Task Cascade frameworks, the HRNet gets the overall better performance.

## 7 ABLATION STUDY

We perform the ablation study for the components in HRNet over two tasks: human pose estimation on COCO validation and semantic segmentation on Cityscapes validation. We mainly use HRNetV1-W32 for human pose estimation, and HRNetV2-W48 for semantic segmentation. All results of pose estimation are obtained over the input size  $256 \times 192$ . We also present the results for comparing HRNetV1 and HRNetV2.

**Representations of different resolutions.** We study how the representation resolution affects the pose estimation performance by checking the quality of the heatmap estimated from the feature maps of each resolution from high to low.

We train two HRNetV1 networks initialized by the model pretrained for the ImageNet classification. Our network outputs four response maps from high-to-low resolutions. The quality of heatmap prediction over the lowest-resolution response map is too low and the AP score is below 10 points. The AP scores over the other three maps are reported in Figure 9. The comparison implies that the resolution does impact the keypoint prediction quality.

**Repeated multi-resolution fusion.** We empirically analyze the effect of the repeated multi-resolution fusion. We study three variants of our network. (a) W/o intermediate fusion units (1 fusion): There is no fusion between multi-resolution streams except the final fusion unit. (b) W/ across-stage fusion units (3 fusions): There is no fusion between parallel streams within each stage. (c) W/ both across-stage

TABLE 11

Comparison with the state-of-the-art single-model object detectors on COCO  $\text{test-dev}$  with BN parameters fixed and without multi-scale training and testing. \* means that the result is from the original paper [12]. GFLOPs and #parameters of the models are given in Table 7. The observations are similar to those on COCO  $\text{val}$ , and show that the HRNet performs better than ResNet and ResNeXt under state-of-the-art object detection and instance segmentation frameworks.

	backbone	size	LS	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
MLKP [142]	VGG16	-	-	28.6	52.4	31.6	10.8	33.4	45.1
STDN [187]	DenseNet-169	513	-	31.8	51.0	33.6	14.4	36.1	43.4
DES [179]	VGG16	512	-	32.8	53.2	34.6	13.9	36.0	47.6
CoupleNet [194]	ResNet-101	-	-	33.1	53.5	35.4	11.6	36.3	50.1
DeNet [139]	ResNet-101	512	-	33.8	53.4	36.1	12.3	36.1	50.8
RFBNet [96]	VGG16	512	-	34.4	55.7	36.4	17.6	37.0	47.6
DFPR [74]	ResNet-101	512	1×	34.6	54.3	37.3	-	-	-
PPFNet [70]	VGG16	512	-	35.2	57.6	37.9	18.7	38.6	45.9
RefineDet [177]	ResNet-101	512	-	36.4	57.5	39.5	16.6	39.9	51.4
Relation Net [56]	ResNet-101	600	-	39.0	58.6	42.9	-	-	-
C-FRCNN [25]	ResNet-101	800	1×	39.0	59.7	42.8	19.4	42.4	53.0
RetinaNet [93]	ResNet-101-FPN	800	1.5×	39.1	59.1	42.3	21.8	42.7	50.2
Deep Regionlets [160]	ResNet-101	800	1.5×	39.3	59.8	-	21.7	43.7	50.9
FitnessNMS [140]	ResNet-101	768	-	39.5	58.0	42.6	18.9	43.5	54.1
DetNet [86]	DetNet59-FPN	800	2×	40.3	62.1	43.8	23.6	42.6	50.0
CornerNet [79]	Hourglass-104	511	-	40.5	56.5	43.1	19.4	42.7	53.9
M2Det [185]	VGG16	800	~ 10×	41.0	59.7	45.0	22.1	46.5	53.8
Faster R-CNN [92]	ResNet-101-FPN	800	1×	39.3	61.3	42.7	22.1	42.1	49.7
Faster R-CNN	HRNetV2p-W32	800	1×	39.5	61.2	43.0	23.3	41.7	49.1
Faster R-CNN [92]	ResNet-101-FPN	800	2×	40.3	61.8	43.9	22.6	43.1	51.0
Faster R-CNN	HRNetV2p-W32	800	2×	41.1	62.3	44.9	24.0	43.1	51.4
Faster R-CNN [92]	ResNet-152-FPN	800	2×	40.6	62.1	44.3	22.6	43.4	52.0
Faster R-CNN	HRNetV2p-W40	800	2×	42.1	63.2	46.1	24.6	44.5	52.6
Faster R-CNN [17]	X-101-64×4d-FPN	800	2×	41.1	62.8	44.8	23.5	44.1	52.3
Faster R-CNN	HRNetV2p-W48	800	2×	42.4	63.6	46.4	24.9	44.6	53.0
Cascade R-CNN [12]*	ResNet-101-FPN	800	~ 1.6×	42.8	62.1	46.3	23.7	45.5	55.2
Cascade R-CNN	ResNet-101-FPN	800	~ 1.6×	43.1	61.7	46.7	24.1	45.9	55.0
Cascade R-CNN	HRNetV2p-W32	800	~ 1.6×	43.7	62.0	47.4	25.5	46.0	55.3
Cascade R-CNN	X-101-64 × 4d-FPN	800	~ 1.6×	44.9	63.7	48.9	25.9	47.7	57.1
Cascade R-CNN	HRNetV2p-W48	800	~ 1.6×	44.8	63.1	48.6	26.0	47.3	56.3
FCOS [136]	ResNet-50-FPN	800	2×	37.3	56.4	39.7	20.4	39.6	47.5
FCOS	HRNetV2p-W18	800	2×	37.8	56.1	40.4	21.6	39.8	47.4
FCOS [136]	ResNet-101-FPN	800	2×	39.2	58.8	41.6	21.8	41.7	50.0
FCOS	HRNetV2p-W32	800	2×	40.5	59.3	43.3	23.4	42.6	51.0
CenterNet [36]	Hourglass-52	511	-	41.6	59.4	44.2	22.5	43.1	54.1
CenterNet	HRNetV2-W48	511	-	43.5	62.1	46.5	22.2	46.5	57.8
Cascade Mask R-CNN [13]	ResNet-101-FPN	800	~ 1.6×	44.0	62.3	47.9	24.3	46.9	56.7
Cascade Mask R-CNN	HRNetV2p-W32	800	~ 1.6×	44.7	62.5	48.6	25.8	47.1	56.3
Cascade Mask R-CNN [13]	X-101-64 × 4d-FPN	800	~ 1.6×	45.9	64.5	50.0	26.6	49.0	58.6
Cascade Mask R-CNN	HRNetV2p-W48	800	~ 1.6×	46.1	64.0	50.3	27.1	48.6	58.3
Hybrid Task Cascade [16]	ResNet-101-FPN	800	~ 1.6×	45.1	64.3	49.0	25.2	48.0	58.2
Hybrid Task Cascade	HRNetV2p-W32	800	~ 1.6×	45.6	64.1	49.4	26.7	47.7	58.0
Hybrid Task Cascade [16]	X-101-64 × 4d-FPN	800	~ 1.6×	47.2	66.5	51.4	27.7	50.1	60.3
Hybrid Task Cascade	HRNetV2p-W48	800	~ 1.6×	47.0	65.8	51.0	27.9	49.4	59.7
Hybrid Task Cascade [16]	X-101-64 × 4d-FPN	800	~ 2.3×	47.2	66.6	51.3	27.5	50.1	60.6
Hybrid Task Cascade	HRNetV2p-W48	800	~ 2.3×	47.3	65.9	51.2	28.0	49.7	59.8

TABLE 12

Ablation study for multi-resolution fusion units on COCO  $\text{val}$  human pose estimation (AP) and Cityscapes  $\text{val}$  semantic segmentation (mIoU). Final = final fusion immediately before representation head, Across = intermediate fusions across stages, Within = intermediate fusions within stages. We can see that the three fusions are beneficial for both human pose estimation and semantic segmentation.

Method	Final	Across	Within	Pose (AP)	Segmentation (mIoU)
(a)	✓			70.8	74.8
(b)	✓	✓		71.9	75.4
(c)	✓	✓	✓	73.4	76.4

and within-stage fusion units (totally 8 fusions): This is our proposed method. All the networks are trained from scratch. The results on COCO human pose estimation and Cityscapes semantic segmentation (validation) given in Table 12 show that the multi-resolution fusion unit is helpful and more fusions lead to better performance.

We also study other possible choices for the fusion design: (i) use bilinear downsample to replace strided convolutions, and (ii) use the multiplication operation to replace the sum operation. In the former case, the COCO pose estimation AP score and the Cityscapes segmentation mIoU score are reduced to 72.6 and 74.2. The reason is that downsam-

pling reduces the volume size ( $\text{width} \times \text{height} \times \text{\#channels}$ ) of the representation maps, and strided convolutions learn better volume size reduction than bilinear downsampling. In the later case, the results are much worse: 54.7 and 66.0, respectively. The possible reason might be that multiplication increases the training difficulty as pointed in [145].

**Resolution maintenance.** We study the performance of a variant of the HRNet: all the four high-to-low resolution streams are added at the beginning and the depths of the four streams are the same; the fusion schemes are the same to ours. Both the HRNets and the variants (with similar #Params and GFLOPs) are trained from scratch.

The human pose estimation performance (AP) on COCO *val* for the variant is 72.5, which is lower than 73.4 for HRNetV1-W32. The segmentation performance (mIoU) on Cityscapes *val* for the variant is 75.7, which is lower than 76.4 for HRNetV2-W48. We believe that the reason is that the low-level features extracted from the early stages over the low-resolution streams are less helpful. In addition, another simple variant, only the high-resolution stream of similar #parameters and GFLOPs without low-resolution parallel streams shows much lower performance on COCO and Cityscapes.

**V1 vs. V2.** We compare HRNetV2 and HRNetV2p, to HRNetV1 on pose estimation, semantic segmentation and COCO object detection. For human pose estimation, the performance is similar. For example, HRNetV2-W32 (w/o ImageNet pretraining) achieves the AP score 73.6, which is slightly higher than 73.4 HRNetV1-W32.

The segmentation and object detection results, given in Figure 10 (a) and Figure 10 (b), imply that HRNetV2 outperforms HRNetV1 significantly, except that the gain is minor in the large model case (1 $\times$ ) in segmentation for Cityscapes. We also test a variant (denoted by HRNetV1h), which is built by appending a  $1 \times 1$  convolution to align the dimension of the output high-resolution representation with the dimension of HRNetV2. The results in Figure 10 (a) and Figure 10 (b) show that the variant achieves slight improvement to HRNetV1, implying that aggregating the representations from low-resolution parallel convolutions in our HRNetV2 is essential for improving the capability.

## 8 CONCLUSIONS

In this paper, we present a high-resolution network for visual recognition problems. There are three fundamental differences from existing low-resolution classification networks and high-resolution representation learning networks: (i) Connect high and low resolution convolutions in parallel other than in series; (ii) Maintain high resolution through the whole process instead of recovering high resolution from low resolution; and (iii) Fuse multi-resolution representations repeatedly, rendering rich high-resolution representations with strong position sensitivity.

The superior results on a wide range of visual recognition problems suggest that our proposed HRNet is a stronger backbone for computer vision problems. Our research also encourages more research efforts for designing network architectures directly for specific vision problems

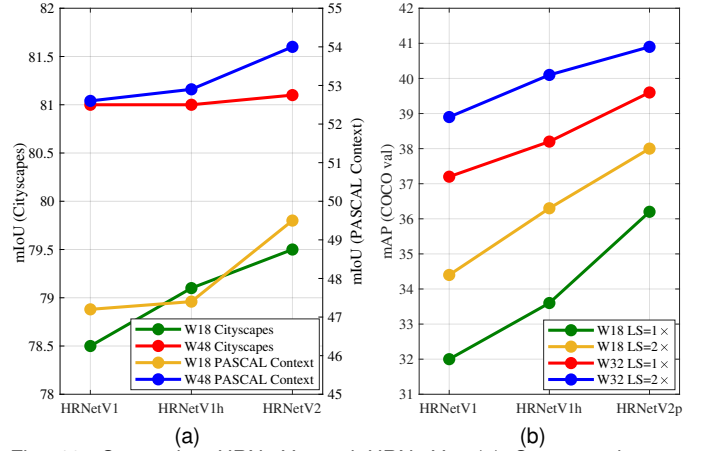


Fig. 10. Comparing HRNetV1 and HRNetV2. (a) Segmentation on Cityscapes *val* and PASCAL-Context for comparing HRNetV1 and its variant HRNetV1h, and HRNetV2 (single scale and no flipping). (b) Object detection on COCO *val* for comparing HRNetV1 and its variant HRNetV1h, and HRNetV2p (LS = learning schedule). We can see that HRNetV2 is superior to HRNetV1 for both semantic segmentation and object detection.

other than extending, remediating or repairing representations learned from low-resolution networks (e.g., ResNet or VGGNet).

**Discussions.** There is a possible misunderstanding: the memory cost of the HRNet is larger as the resolution is higher. In fact, the memory cost of the HRNet for all the three applications, human pose estimation, semantic segmentation and object detection, is comparable to state-of-the-arts except that the training memory cost in object detection is a little larger.

In addition, we summarize the runtime cost comparison on the PyTorch 1.0 platform. The training and inference time cost of the HRNet is comparable to previous state-of-the-arts except that (1) the inference time of the HRNet for segmentation is much smaller and (2) the training time of the HRNet for pose estimation is a little larger, but the cost on the MXNet 1.5.1 platform, which supports static graph inference, is similar as SimpleBaseline. We would like to highlight that for semantic segmentation the inference cost is significantly smaller than PSPNet and DeepLabv3. Table 13 summarizes memory and time cost comparisons<sup>5</sup>.

**Future and followup works.** We will study the combination of the HRNet with other techniques for semantic segmentation and instance segmentation. Currently, we have results (mIoU), which are depicted in Tables 3 4 5 6, by combining the HRNet with the object-contextual representation (OCR) scheme [170]<sup>6</sup>, a variant of object context [59], [171]. We will conduct the study by further increasing the resolution of the representation, e.g., to  $\frac{1}{2}$  or even a full resolution.

The applications of the HRNet are not limited to the above that we have done, and are suitable to other position-sensitive vision applications, such as facial landmark de-

5. The detailed comparisons are given in the supplementary file.

6. We empirically observed that the HRNet combined with ASPP [20] or PPM [181] did not get a performance improvement on Cityscape, but got a slight improvement on PASCAL-Context and LIP.



TABLE 13

Memory and time cost comparisons for pose estimation, semantic segmentation and object detection (under the Faster R-CNN framework) on PyTorch 1.0 in terms of training/inference memory and training/inference time. We also report inference time (in ()) for pose estimation on MXNet 1.5.1, which supports static graph inference that multi-branch convolutions used in the HRNet benefits from. The numbers for training are obtained on a machine with 4 V100 GPU cards. During training, the input sizes are  $256 \times 192$ ,  $512 \times 1024$ , and  $800 \times 1333$ , and the batch sizes are 128, 8 and 8 for pose estimation, segmentation and detection respectively. The numbers for inference are obtained on a single V100 GPU card. The input sizes are  $256 \times 192$ ,  $1024 \times 2048$ , and  $800 \times 1333$ , respectively. The score means AP for pose estimation on COCO val (Table 1) and detection on COCO val (Table 8), and mIoU for cityscapes segmentation (Table 3). Several observations are highlighted. Memory: The HRNet consumes similar memory for both training and inference except that it consumes smaller memory for training in human pose estimation. Time: The training and inference time cost of the HRNet is comparable to previous state-of-the-arts except that the inference time of the HRNet for segmentation is much smaller. SB-ResNet-152 = SimpleBaseline with the backbone of ResNet-152. PSPNet and DeepLabV3 use dilated ResNet-101 as the backbone (Table 3).

	Pose estimation		Segmentation			Detection			
	SB-ResNet-152	HRNetV1-W48	PSPNet	DeepLabV3	HRNetV2-W48	ResNet-101	ResNeXt-101	HRNetV2p-W32	HRNetV2p-W48
training memory	14.8G	7.3G	14.4G	13.3G	13.9G	5.4G	9.5G	8.5G	11.3G
inference memory/image	0.29G	0.27G	1.60G	1.15G	1.79G	0.62G	0.77G	0.51G	0.79G
training second/iteration	1.085	1.231	0.837	0.850	0.692	0.550	1.183	0.690	0.965
inference second/image	0.030 (0.012)	0.058 (0.017)	0.397	0.411	0.150	0.087	0.144	0.101	0.116
score	72.0	75.1	79.7	78.5	81.1	39.8	40.8	40.9	41.8

tection<sup>7</sup>, super-resolution, optical flow estimation, depth estimation, and so on. There are already followup works, e.g., image stylization [83], inpainting [50], image enhancement [62], image dehazing [1], temporal pose estimation [6], and drone object detection [190].

It is reported in [26] that a slightly-modified HRNet combined with ASPP achieved the best performance for Mapillary panoptic segmentation in the single model case. In the COCO + Mapillary Joint Recognition Challenge Workshop at ICCV 2019, the COCO DensePose challenge winner and almost all the COCO keypoint detection challenge participants adopted the HRNet. The OpenImage instance segmentation challenge winner (ICCV 2019) also used the HRNet.

## REFERENCES

- [1] C. O. Ancuti, C. Ancuti, R. Timofte, L. Van Gool, L. Zhang, and M.-H. Yang. Ntire 2019 image dehazing challenge report. In *CVPR Workshops*, June 2019. 12
- [2] A. Arnab, S. Jayasumana, S. Zheng, and P. H. S. Torr. Higher order conditional random fields in deep neural networks. In *ECCV*, pages 524–540, 2016. 7
- [3] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(12):2481–2495, 2017. 1, 2
- [4] T. Baltrusaitis, P. Robinson, and L. Morency. Constrained local neural fields for robust facial landmark detection in the wild. In *ICCVW*, pages 354–361, 2013. 21
- [5] P. N. Belhumeur, D. W. Jacobs, D. J. Kriegman, and N. Kumar. Localizing parts of faces using a consensus of exemplars. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(12):2930–2940, 2013. 19
- [6] G. Bertasius, C. Feichtenhofer, D. Tran, J. Shi, and L. Torresani. Learning temporal pose estimation from sparsely-labeled videos. *CoRR*, abs/1906.04016, 2019. 12
- [7] A. Bulat and G. Tzimiropoulos. Human pose estimation via convolutional part heatmap regression. In *ECCV*, pages 717–732, 2016. 2
- [8] A. Bulat and G. Tzimiropoulos. Binarized convolutional landmark localizers for human pose estimation and face alignment with limited resources. In *ICCV*, pages 3726–3734, 2017. 2
- [9] A. Bulat and G. Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230, 000 3d facial landmarks). In *ICCV*, pages 1021–1030, 2017. 2
- [10] X. P. Burgos-Artizzu, P. Perona, and P. Dollár. Robust face landmark estimation under occlusion. In *ICCV*, pages 1513–1520, 2013. 19, 20
- [11] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos. A unified multi-scale deep convolutional neural network for fast object detection. In *ECCV*, pages 354–370, 2016. 2
- [12] Z. Cai and N. Vasconcelos. Cascade R-CNN: delving into high quality object detection. In *CVPR*, pages 6154–6162, 2018. 2, 9, 10
- [13] Z. Cai and N. Vasconcelos. Cascade R-CNN: high quality object detection and instance segmentation. *CoRR*, abs/1906.09756, 2019. 8, 9, 10
- [14] X. Cao, Y. Wei, F. Wen, and J. Sun. Face alignment by explicit shape regression. In *CVPR*, pages 2887–2894, 2012. 20, 21
- [15] Z. Cao, T. Simon, S. Wei, and Y. Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, pages 1302–1310, 2017. 5
- [16] K. Chen, J. Pang, J. Wang, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Shi, W. Ouyang, C. C. Loy, and D. Lin. Hybrid task cascade for instance segmentation. *CoRR*, abs/1901.07518, 2019. 2, 8, 9, 10
- [17] K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu, Z. Zhang, D. Cheng, C. Zhu, T. Cheng, Q. Zhao, B. Li, X. Lu, R. Zhu, Y. Wu, J. Dai, J. Wang, J. Shi, W. Ouyang, C. C. Loy, and D. Lin. Mmdetection: Open mmlab detection toolbox and benchmark. *CoRR*, abs/1906.07155, 2019. 9, 10, 22
- [18] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *CoRR*, abs/1412.7062, 2014. 2
- [19] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(4):834–848, 2018. 1, 2, 6, 7
- [20] L. Chen, G. Papandreou, F. Schroff, and H. Adam. Rethinking atrous convolution for semantic image segmentation. *CoRR*, abs/1706.05587, 2017. 6, 11
- [21] L. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille. Attention to scale: Scale-aware semantic image segmentation. In *CVPR*, pages 3640–3649, 2016. 2
- [22] L. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, pages 833–851, 2018. 2, 6, 7
- [23] Y. Chen, C. Shen, X. Wei, L. Liu, and J. Yang. Adversarial posenet: A structure-aware convolutional network for human pose estimation. In *ICCV*, pages 1221–1230, 2017. 19, 21
- [24] Y. Chen, Z. Wang, Y. Peng, Z. Zhang, G. Yu, and J. Sun. Cascaded pyramid network for multi-person pose estimation. *CoRR*, abs/1711.07319, 2017. 2, 5, 6
- [25] Z. Chen, S. Huang, and D. Tao. Context refinement for object detection. In *ECCV*, pages 74–89, 2018. 10
- [26] B. Cheng, M. D. Collins, Y. Zhu, T. Liu, T. S. Huang, H. Adam, and L. Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. *CoRR*, abs/1911.10194, 2019. 12
- [27] X. Chu, W. Yang, W. Ouyang, C. Ma, A. L. Yuille, and X. Wang. Multi-context attention for human pose estimation. In *CVPR*, pages 5669–5678, 2017. 2
- [28] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, pages 3213–3223, 2016. 7

7. We provide the facial landmark detection results in the supplementary file.

- [29] J. Dai, K. He, and J. Sun. Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *ICCV*, pages 1635–1643, 2015. 7
- [30] J. Deng, Q. Liu, J. Yang, and D. Tao.  $M^3$  CSR: multi-view, multi-scale and multi-component cascade shape regression. *Image Vision Comput.*, 47:19–26, 2016. 21
- [31] J. Deng, G. Trigeorgis, Y. Zhou, and S. Zafeiriou. Joint multi-view face alignment in the wild. *CoRR*, abs/1708.06023, 2017. 2, 21
- [32] H. Ding, X. Jiang, B. Shuai, A. Q. Liu, and G. Wang. Context contrasted feature and gated multi-scale aggregation for scene segmentation. In *CVPR*, pages 2393–2402, 2018. 7
- [33] H. Ding, X. Jiang, B. Shuai, A. Q. Liu, and G. Wang. Semantic correlation promoted shape-variant context for segmentation. In *CVPR*, June 2019. 6, 7
- [34] X. Dong, Y. Yan, W. Ouyang, and Y. Yang. Style aggregated network for facial landmark detection. In *CVPR*, pages 379–388, 2018. 20, 21
- [35] X. Dong, S. Yu, X. Weng, S. Wei, Y. Yang, and Y. Sheikh. Supervision-by-registration: An unsupervised approach to improve the precision of facial landmark detectors. In *CVPR*, pages 360–368, 2018. 20, 21
- [36] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian. Centernet: Keypoint triplets for object detection. *CoRR*, abs/1904.08189, 2019. 2, 8, 9, 10
- [37] H. Fan and E. Zhou. Approaching human level facial landmark localization by deep learning. *Image Vision Comput.*, 47:27–35, 2016. 21
- [38] H. Fang, S. Xie, Y. Tai, and C. Lu. RMPE: regional multi-person pose estimation. In *ICCV*, pages 2353–2362, 2017. 5
- [39] Z. Feng, G. Hu, J. Kittler, W. J. Christmas, and X. Wu. Cascaded collaborative regression for robust facial landmark detection trained using a mixture of synthetic and real images with dynamic weighting. *IEEE Trans. Image Processing*, 24(11):3425–3440, 2015. 20
- [40] Z. Feng, J. Kittler, M. Awais, P. Huber, and X. Wu. Wing loss for robust facial landmark localisation with convolutional neural networks. In *CVPR*, pages 2235–2245, 2018. 19, 20
- [41] Z. Feng, J. Kittler, W. J. Christmas, P. Huber, and X. Wu. Dynamic attention-controlled cascaded shape regression exploiting training data augmentation and fuzzy-set sample weighting. In *CVPR*, pages 3681–3690, 2017. 20
- [42] D. Fourure, R. Emonet, É. Fromont, D. Muselet, A. Trémeau, and C. Wolf. Residual conv-deconv grid network for semantic segmentation. In *BMVC*, 2017. 2, 6
- [43] J. Fu, J. Liu, H. Tian, Z. Fang, and H. Lu. Dual attention network for scene segmentation. *CoRR*, abs/1809.02983, 2018. 6, 7
- [44] J. Fu, J. Liu, Y. Wang, and H. Lu. Stacked deconvolutional network for semantic segmentation. *CoRR*, abs/1708.04943, 2017. 2
- [45] G. Ghiasi and C. C. Fowlkes. Occlusion coherence: Localizing occluded faces with a hierarchical deformable part model. In *CVPR*, pages 1899–1906, 2014. 20
- [46] G. Ghiasi and C. C. Fowlkes. Laplacian pyramid reconstruction and refinement for semantic segmentation. In *ECCV*, pages 519–534, 2016. 6
- [47] K. Gong, X. Liang, X. Shen, and L. Lin. Look into person: Self-supervised structure-sensitive learning and A new benchmark for human parsing. *CoRR*, abs/1703.05446, 2017. 7, 8
- [48] R. A. Güler, G. Trigeorgis, E. Antonakos, P. Snape, S. Zafeiriou, and I. Kokkinos. Densereg: Fully convolutional dense shape regression in-the-wild. In *CVPR*, pages 2614–2623, 2017. 21
- [49] J. Guo, J. Deng, N. Xue, and S. Zafeiriou. Stacked dense u-nets with dual transformers for robust face alignment. In *BMVC*, page 44, 2018. 2
- [50] Z. Guo, Z. Chen, T. Yu, J. Chen, and S. Liu. Progressive image inpainting with full-resolution residual network. *CoRR*, abs/1907.10478, 2019. 12
- [51] J. He, Z. Deng, L. Zhou, Y. Wang, and Y. Qiao. Adaptive pyramid context network for semantic segmentation. In *CVPR*, June 2019. 7, 8
- [52] K. He, R. B. Girshick, and P. Dollár. Rethinking imagenet pre-training. *CoRR*, abs/1811.08883, 2018. 9
- [53] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick. Mask R-CNN. In *ICCV*, pages 2980–2988, 2017. 2, 5, 8, 9
- [54] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 1, 2, 6, 8, 16
- [55] S. Honari, J. Yosinski, P. Vincent, and C. J. Pal. Recombinator networks: Learning coarse-to-fine feature aggregation. In *CVPR*, pages 5743–5752, 2016. 2, 20, 21
- [56] H. Hu, J. Gu, Z. Zhang, J. Dai, and Y. Wei. Relation networks for object detection. In *CVPR*, pages 3588–3597, 2018. 10
- [57] P. Hu and D. Ramanan. Bottom-up and top-down reasoning with hierarchical rectified gaussians. In *CVPR*, pages 5600–5609, 2016. 2
- [58] G. Huang, D. Chen, T. Li, F. Wu, L. van der Maaten, and K. Q. Weinberger. Multi-scale dense convolutional networks for efficient prediction. *CoRR*, abs/1703.09844, 2017. 2
- [59] L. Huang, Y. Yuan, J. Guo, C. Zhang, X. Chen, and J. Wang. Interlaced sparse self-attention for semantic segmentation. *CoRR*, abs/1907.12273, 2019. 11
- [60] S. Huang, M. Gong, and D. Tao. A coarse-fine network for keypoint localization. In *ICCV*, pages 3047–3056. IEEE Computer Society, 2017. 5
- [61] Z. Huang, X. Wang, L. Huang, C. Huang, Y. Wei, and W. Liu. Ccnet: Criss-cross attention for semantic segmentation. *CoRR*, abs/1811.11721, 2018. 6
- [62] A. Ignatov and R. Timofte. Ntire 2019 challenge on image enhancement: Methods and results. In *CVPRW*, June 2019. 12
- [63] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele. Deepcrut: A deeper, stronger, and faster multi-person pose estimation model. In *ECCV*, pages 34–50, 2016. 2
- [64] M. A. Islam, M. Roohan, N. D. B. Bruce, and Y. Wang. Gated feedback refinement network for dense image labeling. In *CVPR*, pages 4877–4885, 2017. 2
- [65] X. Jin, X. Li, H. Xiao, X. Shen, Z. Lin, J. Yang, Y. Chen, J. Dong, L. Liu, Z. Jie, J. Feng, and S. Yan. Video scene parsing with predictive feature learning. In *ICCV*, pages 5581–5589, 2017. 6
- [66] A. Kanazawa, A. Sharma, and D. W. Jacobs. Locally scale-invariant convolutional neural networks. *CoRR*, abs/1412.5104, 2014. 2
- [67] V. Kazemi and J. Sullivan. One millisecond face alignment with an ensemble of regression trees. In *CVPR*, pages 1867–1874, 2014. 20, 21
- [68] L. Ke, M. Chang, H. Qi, and S. Lyu. Multi-scale structure-aware network for human pose estimation. *CoRR*, abs/1803.09894, 2018. 2
- [69] T. Ke, J. Hwang, Z. Liu, and S. X. Yu. Adaptive affinity fields for semantic segmentation. In *ECCV*, pages 605–621, 2018. 6
- [70] S. Kim, H. Kook, J. Sun, M. Kang, and S. Ko. Parallel feature pyramid network for object detection. In *ECCV*, pages 239–256, 2018. 10
- [71] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. 6
- [72] M. Kocabas, S. Karagoz, and E. Akbas. Multiposenet: Fast multi-person pose estimation using pose residual network. In *ECCV*, pages 437–453. Springer, 2018. 5
- [73] S. Kong and C. C. Fowlkes. Recurrent scene parsing with perspective understanding in the loop. In *CVPR*, pages 956–965, 2018. 6
- [74] T. Kong, F. Sun, W. Huang, and H. Liu. Deep feature pyramid reconfiguration for object detection. In *ECCV*, pages 172–188, 2018. 10
- [75] M. Köstinger, P. Wohlhart, P. M. Roth, and H. Bischof. Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In *ICCV*, pages 2144–2151, 2011. 19
- [76] M. Kowalski, J. Naruniec, and T. Trzcinski. Deep alignment network: A convolutional neural network for robust face alignment. *CoRR*, abs/1706.01789, 2017. 2, 21
- [77] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1106–1114, 2012. 1
- [78] A. Kumar and R. Chellappa. Disentangling 3d pose in a dendritic CNN for unconstrained 2d face alignment. In *CVPR*, pages 430–439. IEEE Computer Society, 2018. 21
- [79] H. Law and J. Deng. Cornernet: Detecting objects as paired keypoints. In *ECCV*, pages 765–781, 2018. 10
- [80] V. Le, J. Brandt, Z. Lin, L. D. Bourdev, and T. S. Huang. Interactive facial feature localization. In *ECCV (3)*, volume 7574 of *Lecture Notes in Computer Science*, pages 679–692. Springer, 2012. 19
- [81] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324, 1998. 1
- [82] H. Li, P. Xiong, J. An, and L. Wang. Pyramid attention network for semantic segmentation. In *BMVC*, page 285, 2018. 6
- [83] M. Li, C. Ye, and W. Li. High-resolution network for photorealistic style transfer. *CoRR*, abs/1904.11617, 2019. 12
- [84] X. Li, Z. Liu, P. Luo, C. C. Loy, and X. Tang. Not all pixels are equal: Difficulty-aware semantic segmentation via deep layer cascade. In *CVPR*, pages 6459–6468, 2017. 6
- [85] Z. Li, Y. Chen, G. Yu, and Y. Deng. R-FCN++: towards accurate

- region-based fully convolutional networks for object detection. In *AAAI*, pages 7073–7080, 2018. [9](#)
- [86] Z. Li, C. Peng, G. Yu, X. Zhang, Y. Deng, and J. Sun. Detnet: Design backbone for object detection. In *ECCV*, pages 339–354, 2018. [2](#), [10](#)
- [87] X. Liang, K. Gong, X. Shen, and L. Lin. Look into person: Joint body parsing & pose estimation network and A new benchmark. *CoRR*, abs/1804.01984, 2018. [7](#)
- [88] X. Liang, H. Zhou, and E. Xing. Dynamic-structured semantic propagation network. In *CVPR*, pages 752–761, 2018. [6](#)
- [89] I. Lifshitz, E. Fetaya, and S. Ullman. Human pose estimation using deep consensus voting. In *ECCV*, pages 246–260, 2016. [2](#)
- [90] G. Lin, A. Milan, C. Shen, and I. D. Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *CVPR*, pages 5168–5177, 2017. [2](#), [6](#), [7](#)
- [91] G. Lin, C. Shen, A. van den Hengel, and I. D. Reid. Efficient piecewise training of deep structured models for semantic segmentation. In *CVPR*, pages 3194–3203, 2016. [6](#), [7](#)
- [92] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie. Feature pyramid networks for object detection. In *CVPR*, pages 936–944, 2017. [2](#), [8](#), [10](#)
- [93] T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *ICCV*, pages 2999–3007, 2017. [8](#), [10](#)
- [94] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. In *ECCV*, pages 740–755, 2014. [1](#), [5](#)
- [95] C. Liu, L.-C. Chen, F. Schroff, H. Adam, W. Hua, A. Yuille, and L. Fei-Fei. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. *arXiv preprint arXiv:1901.02985*, 2019. [6](#)
- [96] S. Liu, D. Huang, and Y. Wang. Receptive field block net for accurate and fast object detection. In *ECCV*, pages 404–419, 2018. [10](#)
- [97] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia. Path aggregation network for instance segmentation. In *CVPR*, pages 8759–8768, 2018. [9](#)
- [98] T. Liu, T. Ruan, Z. Huang, Y. Wei, S. Wei, Y. Zhao, and T. Huang. Devil in the details: Towards accurate single and multiple human parsing. *CoRR*, abs/1809.05996, 2018. [7](#), [8](#)
- [99] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015. [2](#)
- [100] Y. Luo, Z. Zheng, L. Zheng, T. Guan, J. Yu, and Y. Yang. Macro-micro adversarial network for human parsing. In *ECCV*, pages 424–440, 2018. [7](#)
- [101] J. Lv, X. Shao, J. Xing, C. Cheng, and X. Zhou. A deep regression architecture with two-stage re-initialization for high performance facial landmark detection. In *CVPR*, pages 3691–3700, 2017. [20](#)
- [102] X. Miao, X. Zhen, X. Liu, C. Deng, V. Athitsos, and H. Huang. Direct shape regression networks for end-to-end face alignment. In *CVPR*, pages 5040–5049, 2018. [20](#), [21](#)
- [103] R. Mottaghi, X. Chen, X. Liu, N. Cho, S. Lee, S. Fidler, R. Urtasun, and A. L. Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, pages 891–898, 2014. [7](#)
- [104] A. Newell, Z. Huang, and J. Deng. Associative embedding: End-to-end learning for joint detection and grouping. In *NIPS*, pages 2274–2284, 2017. [5](#)
- [105] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, pages 483–499, 2016. [1](#), [2](#), [5](#), [6](#), [19](#)
- [106] X. Nie, J. Feng, and S. Yan. Mutual learning to adapt for joint human parsing and pose estimation. In *ECCV*, pages 519–534, 2018. [7](#)
- [107] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *ICCV*, pages 1520–1528, 2015. [1](#), [2](#)
- [108] G. Papandreou, T. Zhu, L.-C. Chen, S. Gidaris, J. Tompson, and K. Murphy. Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model. In *ECCV*, 2018. [5](#)
- [109] G. Papandreou, T. Zhu, N. Kanazawa, A. Toshev, J. Tompson, C. Bregler, and K. Murphy. Towards accurate multi-person pose estimation in the wild. In *CVPR*, pages 3711–3719, 2017. [5](#), [6](#)
- [110] C. Peng, T. Xiao, Z. Li, Y. Jiang, X. Zhang, K. Jia, G. Yu, and J. Sun. Megdet: A large mini-batch object detector. In *CVPR*, pages 6181–6189, 2018. [9](#)
- [111] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun. Large kernel matters - improve semantic segmentation by global convolutional network. In *CVPR*, pages 1743–1751, 2017. [2](#), [6](#)
- [112] X. Peng, R. S. Feris, X. Wang, and D. N. Metaxas. A recurrent encoder-decoder network for sequential face alignment. In *ECCV* (1), volume 9905, pages 38–56, 2016. [1](#), [2](#)
- [113] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele. Deepcut: Joint subset partition and labeling for multi person pose estimation. In *CVPR*, pages 4929–4937, 2016. [2](#)
- [114] T. Pohlen, A. Hermans, M. Mathias, and B. Leibe. Full-resolution residual networks for semantic segmentation in street scenes. In *CVPR*, pages 3309–3318, 2017. [2](#), [6](#)
- [115] L. Qi, S. Liu, J. Shi, and J. Jia. Sequential context encoding for duplicate removal. In *NeurIPS*, pages 2053–2062, 2018. [9](#)
- [116] S. Ren, X. Cao, Y. Wei, and J. Sun. Face alignment at 3000 FPS via regressing local binary features. In *CVPR*, pages 1685–1692, 2014. [20](#), [21](#)
- [117] S. Ren, X. Cao, Y. Wei, and J. Sun. Face alignment via regressing local binary features. *IEEE Trans. Image Processing*, 25(3):1233–1245, 2016. [19](#)
- [118] S. Ren, K. He, R. B. Girshick, and J. Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(6):1137–1149, 2017. [9](#)
- [119] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241, 2015. [1](#), [2](#)
- [120] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and F. Li. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. [16](#)
- [121] C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic. 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *ICCV Workshops*, pages 397–403, 2013. [19](#)
- [122] M. Samy, K. Amer, K. Eissa, M. Shaker, and M. ElHelw. Nu-net: Deep residual wide field of view convolutional neural network for semantic segmentation. In *CVPRW*, June 2018. [2](#)
- [123] S. Saxena and J. Verbeek. Convolutional neural fabrics. In *NIPS*, pages 4053–4061, 2016. [2](#)
- [124] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *CoRR*, abs/1312.6229, 2013. [2](#)
- [125] E. Shelhamer, J. Long, and T. Darrell. Fully convolutional networks for semantic segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(4):640–651, 2017. [7](#)
- [126] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. [1](#), [2](#)
- [127] B. Singh and L. S. Davis. An analysis of scale invariance in object detection SNIP. In *CVPR*, pages 3578–3587, 2018. [9](#)
- [128] B. Singh, M. Najibi, and L. S. Davis. SNIPER: efficient multi-scale training. In *NeurIPS*, pages 9333–9343, 2018. [9](#)
- [129] K. Sun, M. Li, D. Liu, and J. Wang. IGCv3: interleaved low-rank group convolutions for efficient deep neural networks. In *BMVC*, page 101. BMVA Press, 2018. [3](#)
- [130] K. Sun, B. Xiao, D. Liu, and J. Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019. [3](#), [17](#)
- [131] K. Sun, Y. Zhao, B. Jiang, T. Cheng, B. Xiao, D. Liu, Y. Mu, X. Wang, W. Liu, and J. Wang. High-resolution representations for labeling pixels and regions. *CoRR*, abs/1904.04514, 2019. [3](#)
- [132] X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei. Integral human pose regression. In *ECCV*, pages 536–553, 2018. [5](#)
- [133] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, pages 1–9, 2015. [1](#)
- [134] W. Tang, P. Yu, and Y. Wu. Deeply learned compositional models for human pose estimation. In *ECCV*, September 2018. [2](#)
- [135] Z. Tang, X. Peng, S. Geng, L. Wu, S. Zhang, and D. N. Metaxas. Quantized densely connected u-nets for efficient landmark localization. In *ECCV*, pages 348–364, 2018. [2](#)
- [136] Z. Tian, C. Shen, H. Chen, and T. He. FCOS: fully convolutional one-stage object detection. *CoRR*, abs/1904.01355, 2019. [2](#), [8](#), [9](#), [10](#)
- [137] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler. Efficient object localization using convolutional networks. In *CVPR*, pages 648–656, 2015. [2](#)
- [138] G. Trigeorgis, P. Snape, M. A. Nicolaou, E. Antonakos, and S. Zafeiriou. Mnemonic descent method: A recurrent process applied for end-to-end face alignment. In *CVPR*, pages 4177–4187, 2016. [21](#)
- [139] L. Tychsen-Smith and L. Petersson. Denet: Scalable real-time object detection with directed sparse sampling. In *ICCV*, pages 428–436, 2017. [10](#)
- [140] L. Tychsen-Smith and L. Petersson. Improving object localization with fitness NMS and bounded iou loss. In *CVPR*, pages 6877–



- 6885, 2018. **10**
- [141] R. Valle, J. M. Buenaposada, A. Valdés, and L. Baumela. A deeply-initialized coarse-to-fine ensemble of regression trees for face alignment. In *ECCV*, pages 609–624, 2018. **2, 19, 20, 21**
- [142] H. Wang, Q. Wang, M. Gao, P. Li, and W. Zuo. Multi-scale location-aware kernel representation for object detection. In *CVPR*, pages 1248–1257, 2018. **10**
- [143] J. Wang, Z. Wei, T. Zhang, and W. Zeng. Deeply-fused nets. *CoRR*, abs/1605.07716, 2016. **3**
- [144] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. W. Cottrell. Understanding convolution for semantic segmentation. In *WACV*, 2018. **6**
- [145] Y. Wang, L. Xie, C. Liu, S. Qiao, Y. Zhang, W. Zhang, Q. Tian, and A. L. Yuille. SORT: second-order response transform for visual recognition. In *ICCV*, pages 1368–1377, 2017. **11**
- [146] Z. Wang, W. Li, B. Yin, Q. Peng, T. Xiao, Y. Du, Z. Li, X. Zhang, G. Yu, and J. Sun. Mscoco keypoints challenge 2018. In *Joint Recognition Challenge Workshop at ECCV 2018*, 2018. **6**
- [147] Z. Wojna, J. R. R. Uijlings, S. Guadarrama, N. Silberman, L. Chen, A. Fathi, and V. Ferrari. The devil is in the decoder. In *BMVC*, 2017. **2**
- [148] J. Wu, H. Zheng, B. Zhao, Y. Li, B. Yan, R. Liang, W. Wang, S. Zhou, G. Lin, Y. Fu, et al. Ai challenger: A large-scale dataset for going deeper in image understanding. *arXiv preprint arXiv:1711.06475*, 2017. **6**
- [149] W. Wu, C. Qian, S. Yang, Q. Wang, Y. Cai, and Q. Zhou. Look at boundary: A boundary-aware face alignment algorithm. In *CVPR*, pages 2129–2138, 2018. **2, 19, 20, 21**
- [150] W. Wu and S. Yang. Leveraging intra and inter-dataset variations for robust face alignment. In *CVPR*, pages 2096–2105, 2017. **20**
- [151] Z. Wu, C. Shen, and A. van den Hengel. Wider or deeper: Revisiting the resnet model for visual recognition. *CoRR*, abs/1611.10080, 2016. **6**
- [152] B. Xiao, H. Wu, and Y. Wei. Simple baselines for human pose estimation and tracking. In *ECCV*, pages 472–487, 2018. **1, 2, 5, 6**
- [153] S. Xiao, J. Feng, J. Xing, H. Lai, S. Yan, and A. A. Kassim. Robust facial landmark detection via recurrent attentive-refinement networks. In *ECCV*, pages 57–72, 2016. **20**
- [154] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun. Unified perceptual parsing for scene understanding. In *ECCV*, pages 432–448, 2018. **2**
- [155] G. Xie, J. Wang, T. Zhang, J. Lai, R. Hong, and G. Qi. Interleaved structured sparse convolutional neural networks. In *CVPR*, pages 8847–8856. IEEE Computer Society, 2018. **3**
- [156] S. Xie, R. B. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *CVPR*, pages 5987–5995, 2017. **8**
- [157] S. Xie and Z. Tu. Holistically-nested edge detection. In *ICCV*, pages 1395–1403, 2015. **2**
- [158] X. Xiong and F. D. la Torre. Supervised descent method and its applications to face alignment. In *CVPR*, pages 532–539. IEEE Computer Society, 2013. **20, 21**
- [159] D. Xu, W. Ouyang, X. Wang, and N. Sebe. Pad-net: Multi-tasks guided prediction-and-distillation network for simultaneous depth estimation and scene parsing. In *CVPR*, pages 675–684, 2018. **6**
- [160] H. Xu, X. Lv, X. Wang, Z. Ren, N. Bodla, and R. Chellappa. Deep regionlets for object detection. In *ECCV*, pages 827–844, 2018. **10**
- [161] Y. Xu, T. Xiao, J. Zhang, K. Yang, and Z. Zhang. Scale-invariant convolutional neural networks. *CoRR*, abs/1411.6369, 2014. **2**
- [162] J. Yan, Z. Lei, D. Yi, and S. Z. Li. Learn to combine multiple hypotheses for accurate face alignment. In *ICCVW*, pages 392–396, 2013. **21**
- [163] J. Yang, Q. Liu, and K. Zhang. Stacked hourglass network for robust facial landmark localisation. In *CVPR*, pages 2025–2033, 2017. **2**
- [164] S. Yang, P. Luo, C. C. Loy, and X. Tang. WIDER FACE: A face detection benchmark. In *CVPR*, pages 5525–5533, 2016. **19**
- [165] W. Yang, S. Li, W. Ouyang, H. Li, and X. Wang. Learning feature pyramids for human pose estimation. In *ICCV*, pages 1290–1299, 2017. **2, 5**
- [166] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *ECCV*, pages 334–349, 2018. **6**
- [167] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang. Learning a discriminative feature network for semantic segmentation. In *CVPR*, pages 1857–1866, 2018. **6**
- [168] F. Yu, V. Koltun, and T. A. Funkhouser. Dilated residual networks. *CoRR*, abs/1705.09914, 2017. **2**
- [169] X. Yu, J. Huang, S. Zhang, W. Yan, and D. N. Metaxas. Pose-free facial landmark fitting via optimized part mixtures and cascaded deformable shape model. In *ICCV*, pages 1944–1951. IEEE Computer Society, 2013. **20**
- [170] Y. Yuan, X. Chen, and J. Wang. Object-contextual representations for semantic segmentation. *CoRR*, abs/1909.11065, 2019. **6, 7, 8, 11**
- [171] Y. Yuan and J. Wang. Ocnet: Object context network for scene parsing. *CoRR*, abs/1809.00916, 2018. **11**
- [172] H. Zhang, K. J. Dana, J. Shi, Z. Zhang, X. Wang, A. Tyagi, and A. Agrawal. Context encoding for semantic segmentation. In *CVPR*, pages 7151–7160, 2018. **7**
- [173] H. Zhang, H. Zhang, C. Wang, and J. Xie. Co-occurrent features in semantic segmentation. In *CVPR*, June 2019. **6, 7**
- [174] J. Zhang, M. Kan, S. Shan, and X. Chen. Occlusion-free face alignment: Deep regression networks coupled with de-corrupt autoencoders. In *CVPR*, pages 3428–3437, 2016. **20**
- [175] J. Zhang, S. Shan, M. Kan, and X. Chen. Coarse-to-fine auto-encoder networks (CFAN) for real-time face alignment. In *ECCV (2)*, volume 8690, pages 1–16. Springer, 2014. **21**
- [176] R. Zhang, S. Tang, Y. Zhang, J. Li, and S. Yan. Scale-adaptive convolutions for scene parsing. In *ICCV*, pages 2050–2058, 2017. **6**
- [177] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li. Single-shot refinement neural network for object detection. In *CVPR*, pages 4203–4212, 2018. **10**
- [178] T. Zhang, G. Qi, B. Xiao, and J. Wang. Interleaved group convolutions. In *ICCV*, pages 4383–4392, 2017. **3, 4**
- [179] Z. Zhang, S. Qiao, C. Xie, W. Shen, B. Wang, and A. L. Yuille. Single-shot object detection with enriched semantics. In *CVPR*, pages 5813–5821, 2018. **10**
- [180] Z. Zhang, X. Zhang, C. Peng, X. Xue, and J. Sun. Exfuse: Enhancing feature fusion for semantic segmentation. In *ECCV*, pages 273–288, 2018. **2**
- [181] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *CVPR*, pages 6230–6239, 2017. **1, 2, 6, 7, 11**
- [182] H. Zhao, Y. Zhang, S. Liu, J. Shi, C. C. Loy, D. Lin, and J. Jia. Psnnet: Point-wise spatial attention network for scene parsing. In *ECCV*, pages 270–286, 2018. **6, 7**
- [183] J. Zhao, J. Li, X. Nie, F. Zhao, Y. Chen, Z. Wang, J. Feng, and S. Yan. Self-supervised neural aggregation networks for human parsing. In *CVPRW*, pages 1595–1603, 2017. **7**
- [184] L. Zhao, M. Li, D. Meng, X. Li, Z. Zhang, Y. Zhuang, Z. Tu, and J. Wang. Deep convolutional neural networks with merge-and-run mappings. In *IJCAI*, pages 3170–3176, 2018. **3**
- [185] Q. Zhao, T. Sheng, Y. Wang, Z. Tang, Y. Chen, L. Cai, and H. Ling. M2det: A single-shot object detector based on multi-level feature pyramid network. *CoRR*, abs/1811.04533, 2018. **10**
- [186] E. Zhou, H. Fan, Z. Cao, Y. Jiang, and Q. Yin. Extensive facial landmark localization with coarse-to-fine convolutional network cascade. In *ICCVW*, pages 386–391, 2013. **21**
- [187] P. Zhou, B. Ni, C. Geng, J. Hu, and Y. Xu. Scale-transferrable object detection. In *CVPR*, pages 528–537, 2018. **10**
- [188] Y. Zhou, X. Hu, and B. Zhang. Interlinked convolutional neural networks for face parsing. In *ISNN*, pages 222–231, 2015. **2**
- [189] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang. Unet++: A nested u-net architecture for medical image segmentation. In *MICCAI*, pages 3–11, 2018. **2, 6, 7**
- [190] P. Zhu, D. Du, L. Wen, X. Bian, H. Ling, Q. Hu, T. Peng, J. Zheng, X. Wang, Y. Zhang, L. Bo, H. Shi, R. Zhu, B. Dong, D. Reddy Pailla, F. Ni, G. Gao, G. Liu, H. Xiong, J. Ge, J. Zhou, J. Hu, L. Sun, L. Chen, M. Lauer, Q. Liu, S. Saketh Chennamsetty, T. Sun, T. Wu, V. Alex Kollerathu, W. Tian, W. Qin, X. Chen, X. Zhao, Y. Lian, Y. Wu, Y. Li, Y. Li, Y. Wang, Y. Song, Y. Yao, Y. Zhang, Z. Pi, Z. Chen, Z. Xu, Z. Xiao, Z. Luo, and Z. Liu. Visdrone-vid2019: The vision meets drone object detection in video challenge results. In *ICCV Workshop*, Oct 2019. **12**
- [191] S. Zhu, C. Li, C. C. Loy, and X. Tang. Face alignment by coarse-to-fine shape searching. In *CVPR*, pages 4998–5006, 2015. **19, 20, 21**
- [192] S. Zhu, C. Li, C. C. Loy, and X. Tang. Unconstrained face alignment via cascaded compositional learning. In *CVPR*, pages 3409–3417, 2016. **20**
- [193] X. Zhu and D. Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *CVPR*, pages 2879–2886, 2012. **19**
- [194] Y. Zhu, C. Zhao, J. Wang, X. Zhao, Y. Wu, and H. Lu. Couplenet: Coupling global structure with local parts for object detection. In *ICCV*, pages 4146–4154, 2017. **10**
- [195] Z. Zhu, M. Xu, S. Bai, T. Huang, and X. Bai. Asymmetric non-local neural networks for semantic segmentation. *CoRR*, abs/1908.07678, 2019. **6, 7**

## APPENDIX A

### NETWORK INSTANTIATION

Our current design (except the standard stem and the head,) contains four stages, as shown in Table 14. Each stage consists of modularized blocks, repeated 1, 1, 4, and 3 times, respectively for the four stages. The modularized block consists of 1 (2, 3 and 4) branches for the 1st (2nd, 3rd and 4th) stages. Each branch corresponds to different resolution, and is composed of four residual units and one multi-resolution fusion unit (See Figure 3 in the main paper).

TABLE 14

The architecture of the HRNet (main body). There are four stages. Each stage consists of modularized blocks, repeated 1, 1, 4, and 3 times, respectively for the four stages. The modularized block consists of 1 (2, 3 and 4) branches for the 1st (2nd, 3rd and 4th) stages. Each branch corresponds to a different resolution, and is composed of four residual units and one multi-resolution fusion unit. For clarity, the fusion unit (after each modularized block) is not depicted in the table, and could be understood from Figure 3 in the main paper. In the table, each cell consists of three components: the first one ( $\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$ ) is the residual unit, the second number is the repetition times of the residual units, and the last number is the repetition times of the modularized blocks.  $C$  in each residual unit is the number of channels.

Resolution	Stage 1	Stage 2	Stage 3	Stage 4
4×	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 4 \times 1$	$\begin{bmatrix} 3 \times 3, C \\ 3 \times 3, C \end{bmatrix} \times 4 \times 1$	$\begin{bmatrix} 3 \times 3, C \\ 3 \times 3, C \end{bmatrix} \times 4 \times 4$	$\begin{bmatrix} 3 \times 3, C \\ 3 \times 3, C \end{bmatrix} \times 4 \times 3$
8×		$\begin{bmatrix} 3 \times 3, 2C \\ 3 \times 3, 2C \end{bmatrix} \times 4 \times 1$	$\begin{bmatrix} 3 \times 3, 2C \\ 3 \times 3, 2C \end{bmatrix} \times 4 \times 4$	$\begin{bmatrix} 3 \times 3, 2C \\ 3 \times 3, 2C \end{bmatrix} \times 4 \times 3$
16×			$\begin{bmatrix} 3 \times 3, 4C \\ 3 \times 3, 4C \end{bmatrix} \times 4 \times 4$	$\begin{bmatrix} 3 \times 3, 4C \\ 3 \times 3, 4C \end{bmatrix} \times 4 \times 3$
32×				$\begin{bmatrix} 3 \times 3, 8C \\ 3 \times 3, 8C \end{bmatrix} \times 4 \times 3$

## APPENDIX B

### NETWORK PRETRAINING

We pretrain our network, which is augmented by a classification head shown in Figure 11, on ImageNet [120]. The classification head is described as below. First, the four-resolution feature maps are fed into a bottleneck and the output channels are increased from  $C$ ,  $2C$ ,  $4C$ , and  $8C$  to 128, 256, 512, and 1024, respectively. Then, we downsample the high-resolution representation by a 2-strided  $3 \times 3$  convolution outputting 256 channels and add it to the representation of the second-high-resolution. This process is repeated two times to get 1024 feature channels over the small resolution. Last, we transform the 1024 channels to 2048 channels through a  $1 \times 1$  convolution, followed by a global average pooling operation. The output 2048-dimensional representation is fed into the classifier.

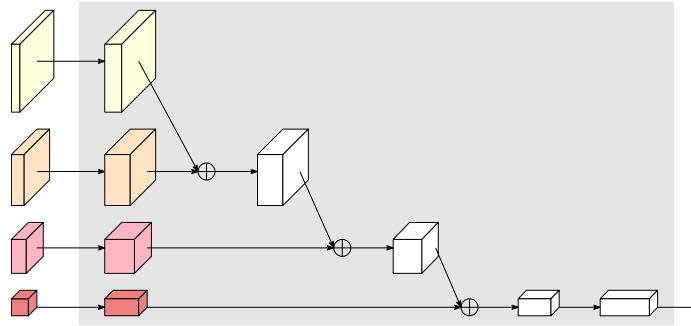


Fig. 11. Representation for ImageNet classification. The input of the box is the representations of four resolutions.

We adopt the same data augmentation scheme for training images as in [54], and train our models for 100 epochs with a batch size of 256. The initial learning rate is set to 0.1 and is reduced by 10 times at epoch 30, 60 and 90. We use SGD with a weight decay of 0.0001 and a Nesterov momentum of 0.9. We adopt standard single-crop testing, so that  $224 \times 224$  pixels are cropped from each image. The top-1 and top-5 error are reported on the validation set.

Table 15 shows our ImageNet classification results. As a comparison, we also report the results of ResNets. We consider two types of residual units: One is formed by a bottleneck, and the other is formed by two  $3 \times 3$  convolutions. We follow the PyTorch implementation of ResNets and replace the  $7 \times 7$  convolution in the input stem with two 2-strided  $3 \times 3$  convolutions decreasing the resolution to  $1/4$  as in our networks. When the residual units are formed by two  $3 \times 3$

TABLE 15  
ImageNet Classification results of HRNet and ResNets. The proposed method is named HRNet-W $x$ -C.  $x$  means the width.

	#Params.	GFLOPs	top-1 err.	top-5 err.
<i>Residual branch formed by two <math>3 \times 3</math> convolutions</i>				
ResNet-38	28.3M	3.80	24.6%	7.4%
HRNet-W18-C	21.3M	3.99	<b>23.1%</b>	<b>6.5%</b>
ResNet-72	48.4M	7.46	23.3%	6.7%
HRNet-W30-C	37.7M	7.55	<b>21.9%</b>	<b>5.9%</b>
ResNet-106	64.9M	11.1	22.7%	6.4%
HRNet-W40-C	57.6M	11.8	<b>21.1%</b>	<b>5.6%</b>
<i>Residual branch formed by a bottleneck</i>				
ResNet-50	25.6M	3.82	23.3%	6.6%
HRNet-W44-C	21.9M	3.90	<b>23.0%</b>	<b>6.5%</b>
ResNet-101	44.6M	7.30	21.6%	<b>5.8%</b>
HRNet-W76-C	40.8M	7.30	<b>21.5%</b>	<b>5.8%</b>
ResNet-152	60.2M	10.7	21.2%	5.7%
HRNet-W96-C	57.5M	10.2	<b>21.0%</b>	<b>5.6%</b>

convolutions, an extra bottleneck is used to increase the dimension of output feature maps from 512 to 2048. One can see that under similar #parameters and GFLOPs, our results are comparable to and slightly better than ResNets.

In addition, we look at the results of two alternative schemes: (i) the feature maps on each resolution go through a global pooling separately and then are concatenated together to output a  $15C$ -dimensional representation vector, named HRNet-W $x$ -Ci; (ii) the feature maps on each resolution are fed into several 2-strided residual units (bottleneck, each dimension is increased to the double) to increase the dimension to 512, and concatenate and average-pool them together to reach a 2048-dimensional representation vector, named HRNet-W $x$ -Cii, which is used in [130]. Table 16 shows such an ablation study. One can see that the proposed manner is superior to the two alternatives.

TABLE 16  
Ablation study on ImageNet classification by comparing our approach (abbreviated as HRNet-W $x$ -C) with two alternatives: HRNet-W $x$ -Ci and HRNet-W $x$ -Cii (residual branch formed by two  $3 \times 3$  convolutions).

	#Params.	GFLOPs	top-1 err.	top-5 err.
HRNet-W27-Ci	21.4M	5.55	26.0%	7.7%
HRNet-W25-Cii	21.7M	5.04	24.1%	7.1%
HRNet-W18-C	21.3M	3.99	<b>23.1%</b>	<b>6.5%</b>
HRNet-W36-Ci	37.5M	9.00	24.3%	7.3%
HRNet-W34-Cii	36.7M	8.29	22.8%	6.3%
HRNet-W30-C	37.7M	7.55	<b>21.9%</b>	<b>5.9%</b>
HRNet-W45-Ci	58.2M	13.4	23.6%	7.0%
HRNet-W43-Cii	56.3 M	12.5	22.2%	6.1%
HRNet-W40-C	57.6M	11.8	<b>21.1%</b>	<b>5.6%</b>



## APPENDIX C

### TRAINING/INFERENCE COST

Tables 17, 18 and 19 provide GPU memory comparisons between HRNets and other standard networks for both training and inference in the PyTorch platform. Compared to state-of-the-arts for human pose estimation, the training and inference memory costs of the HRNet are similar or lower for similar parameter complexity (Table 17). Compared to state-of-the-arts for semantic segmentation, the training and inference memory costs are similar (Table 18) for similar parameter complexity. Compared to state-of-the-arts for object detection for similar parameter complexity, the training and inference memory costs are similar or slightly higher (Table 19).

In addition, we provide the runtime cost comparison. (1) For semantic segmentation, the time cost of the HRNet for training is slightly smaller and for inference significantly smaller than PSPNet and DeepLabv3 (Table 18). (2) For object detection, the time cost of the HRNet for training is larger than ResNet based networks and smaller than ResNext based networks, and for inference the HRNet is smaller for similar GFLOPs (Table 19). (3) For human pose estimation, the time cost of the HRNet for training is similar and for inference larger; and the time cost of the HRNet for training and inference in the MXNet platform is similar as SimpleBaseline (Table 17).

TABLE 17

Human pose estimation complexities on PyTorch 1.0 in terms of #parameters, GFLOPs, training/inference memory. We also report training/inference time on Pytorch 1.0 and MXNet 1.5.1, shown as  $time(Pytorch)/time(MXNet)$ . The reason that the runtime cost is smaller than PyTorch is that MXNet supports dynamic graph which the HRNet benefits from. We compare the HRNet and the previous state-of-the-art, Simplebaseline. Training: 4 V100 GPU cards, input size  $256 \times 192$ , and batch size 128. Inference: a single V100 GPU card, input size  $256 \times 192$ , and batch size 1, 4, 8 and 16. Two observations are highlighted here: (1) The HRNet consumes smaller memory for both training and inference for similar #parameters, and for similar AP scores; (2) The HRNet takes slightly higher training runtime cost and a little higher inference runtime cost on Pytorch and similar on MXNet for similar GFLOPs, and the inference efficiency of HRNet for is improved for larger batch size. sec.= seconds, iter. = iteration, mem. = memory, bs = batchsize.

backbone	GFLOPs	#params	train sec./iter.	train mem.	infer sec./batch				infer mem./batch				AP
					bs = 1	bs = 4	bs = 8	bs = 16	bs = 1	bs = 4	bs = 8	bs = 16	
SB-Res-50	8.90	34.0M	0.946/0.211	11.9G	0.012/0.005	0.013/0.010	0.015/0.017	0.024/0.027	0.16G	0.17G	0.21G	0.30G	70.4
SB-Res-101	12.4	53.0M	1.008/0.320	13.2G	0.020/0.009	0.021/0.014	0.024/0.023	0.035/0.038	0.23G	0.24G	0.28G	0.37G	71.4
SB-Res-152	15.7	68.6M	1.085/0.415	14.8G	0.030/0.012	0.033/0.019	0.035/0.031	0.048/0.051	0.29G	0.31G	0.35G	0.43G	72.0
HRNetV1-W32	7.10	28.5M	1.153/0.389	5.7G	0.057/0.015	0.059/0.017	0.061/0.020	0.062/0.031	0.13G	0.15G	0.19G	0.28G	74.4
HRNetV1-W48	14.6	63.6M	1.231/0.507	7.3G	0.058/0.017	0.060/0.021	0.062/0.033	0.066/0.051	0.27G	0.30G	0.32G	0.44G	75.1

TABLE 18

Semantic segmentation complexities on PyTorch 1.0. Training: 4 V100 GPU cards, input size  $512 \times 1024$ , and batch size 8. Inference: a single V100 GPU card, input size  $1024 \times 2048$ , and batch size 1, 4, and 8. Several observations are highlighted: (1) The training memory costs are similar, and the inference memory costs are similar but larger for our approach with larger batch size; (2) The training and inference time costs of our approach are much smaller. The results are obtained on Cityscapes val.

backbone	GFLOPs	#params	train sec./iter.	train mem.	infer sec./batch			infer mem./batch			mIoU
					bs = 1	bs = 4	bs = 8	bs = 1	bs = 4	bs = 8	
Dilated-ResNet	1661.6	52.1M	0.6611	12.4G	0.3351	1.2882	2.7039	1.13G	3.92G	7.64G	75.7
PSPNet	2017.6	65.9M	0.8368	14.4G	0.3972	1.5296	3.2003	1.60G	5.04G	9.81G	79.7
DeepLabv3	1778.7	58.0M	0.8502	13.3G	0.4113	1.5307	3.2000	1.15G	3.95G	7.67G	78.5
HRNetV2-W48	696.2	65.9M	0.6920	13.9G	0.1502	0.05421	1.1032	1.79G	6.37G	12.5G	81.1

TABLE 19

COCO object detection complexities on PyTorch 1.0. Training: 4 V100 GPU cards, input size  $800 \times 1333$ , and batch size 8. Inference: a single V100 GPU card, input size  $800 \times 1333$ , and batch size 1, 4, and 8. The performance are reported on the COCO 2017val for each model with the learning schedule of  $2 \times$ . Several observations are as follows: (1) The training memory for the HRNet is a little larger for similar #params, but the inference memory are similar with higher AP scores; (2) The training runtime costs of the HRNet are a little larger than ResNet based networks and smaller than ResNext based networks, and the inference runtime costs is smaller for similar GFLOPs.

backbone	GFLOPs	#params	sec./iter.	train mem.	infer sec./batch			infer mem./batch			AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
					bs = 1	bs = 4	bs = 8	bs = 1	bs = 4	bs = 8						
ResNet-50	172.34	39.77M	0.4167	3.4G	0.0739	0.2495	0.4921	0.55G	1.68G	2.93G	37.6	58.7	41.3	21.4	40.8	49.7
ResNet-101	239.37	57.83M	0.5502	5.4G	0.0870	0.3018	0.6089	0.62G	1.76G	3.25G	39.8	61.4	43.4	22.9	43.6	52.4
X-101-64 $\times$ 4d	381.83	94.85M	1.1828	9.5G	0.1438	0.4832	0.9764	0.77G	1.92G	3.41G	40.8	62.1	44.6	23.2	44.5	53.7
HRNetV2p-W18	159.06	26.18M	0.6166	6.2G	0.0968	0.2320	0.4471	0.33G	1.01G	1.91G	38.0	58.9	41.5	22.6	40.8	49.6
HRNetV2p-W32	245.33	45.04M	0.6901	8.5G	0.1014	0.2738	0.5546	0.51G	1.51G	2.83G	40.9	61.8	44.8	24.4	43.7	53.3
HRNetV2p-W48	399.12	79.42M	0.9648	11.3G	0.1162	0.3762	0.7296	0.79G	2.16G	3.99G	41.8	62.8	45.9	25.0	44.7	54.6

## APPENDIX D

### FACIAL LANDMARK DETECTION

Facial landmark detection a.k.a. face alignment is a problem of detecting the keypoints from a face image. We perform the evaluation over four standard datasets: WFLW [149], AFLW [75], COFW [10], and 300W [121]. We mainly use the normalized mean error (NME) for evaluation. We use the inter-ocular distance as normalization for WFLW, COFW, and 300W, and the face bounding box as normalization for AFLW. We also report area-under-the-curve scores (AUC) and failure rates.

We follow the standard scheme [149] for training. All the faces are cropped by the provided boxes according to the center location and resized to  $256 \times 256$ . We augment the data by  $\pm 30$  degrees in-plane rotation,  $0.75 - 1.25$  scaling, and randomly flipping. The base learning rate is 0.0001 and is dropped to 0.00001 and 0.000001 at the 30th and 50th epochs. The models are trained for 60 epochs with the batch size of 16 on one GPU. Different from semantic segmentation, the heatmaps are not upsampled from  $1/4$  to the input size, and the loss function is optimized over the  $1/4$  maps.

At testing, each keypoint location is predicted by transforming the highest heatvalue location from  $1/4$  to the original image space and adjusting it with a quarter offset in the direction from the highest response to the second highest response [23].

We adopt HRNetV2-W18 for face landmark detection whose parameter and computation cost are similar to or smaller than models with widely-used backbones: ResNet-50 and Hourglass [105]. HRNetV2-W18: #parameters = 9.3M, GFLOPs = 4.3G; ResNet-50: #parameters = 25.0M, GFLOPs = 3.8G; Hourglass: #parameters = 25.1M, GFLOPs = 19.1G. The numbers are obtained on the input size  $256 \times 256$ . It should be noted that the facial landmark detection methods adopting ResNet-50 and Hourglass as backbones introduce extra parameter and computation overhead.

**WFLW.** The WFLW dataset [149] is a recently-built dataset based on the WIDER Face [164]. There are 7,500 training and 2,500 testing images with 98 manual annotated landmarks. We report the results on the test set and several subsets: large pose (326 images), expression (314 images), illumination (698 images), make-up (206 images), occlusion (736 images) and blur (773 images).

Table 20 provides the comparison of our method with state-of-the-art methods. Our approach is significantly better than other methods on the test set and all the subsets, including LAB that exploits extra boundary information [149] and PDB that uses stronger data augmentation [40].

**AFLW.** The AFLW [75] dataset is a widely used benchmark dataset, where each image has 19 facial landmarks. Following [149], [191], we train our models on 20,000 training images, and report the results on the AFLW-Full set (4,386 testing images) and the AFLW-Frontal set (1314 testing images selected from 4386 testing images).

Table 21 provides the comparison of our method with state-of-the-art methods. Our approach achieves the best performance among methods without extra information and stronger data augmentation and even outperforms DCFE with extra 3D information. Our approach performs slightly worse than LAB that uses extra boundary information [149] and PDB [40] that uses stronger data augmentation.

**COFW.** The COFW dataset [10] consists of 1,345 training and 507 testing faces with occlusions, where each image has 29 facial landmarks.

Table 22 provides the comparison of our method with state-of-the-art methods. HRNetV2 outperforms other methods by a large margin. In particular, it achieves the better performance than LAB with extra boundary information and PDB with stronger data augmentation.

**300W.** The dataset [121] is a combination of HELEN [80], LFPW [5], AFW [193], XM2VTS and IBUG datasets, where each face has 68 landmarks. Following [117], we use the 3,148 training images, which contains the training subsets of HELEN and LFPW and the full set of AFW. We evaluate the performance using two protocols, full set and test set. The full set contains 689 images and is further divided into a common subset (554 images) from HELEN and LFPW, and a challenging subset (135 images) from IBUG. The official test set, used for competition, contains 600 images (300 indoor and 300 outdoor images).

Table 23 provides the results on the full set, and its two subsets: common and challenging. Table 24 provides the results on the test set. In comparison to Chen et al. [23] that uses Hourglass with large parameter and computation complexity as the backbone, our scores are better except the  $AUC_{0.08}$  scores. Our HRNetV2 gets the overall best performance among methods without extra information and stronger data augmentation, and is even better than LAB with extra boundary information and DCFE [141] that explores extra 3D information.

TABLE 20

Facial landmark detection results (NME) on WFLW  $_{test}$  and 6 subsets: pose, expression (expr.), illumination (illu.), make-up (mu.), occlusion (occu.) and blur. LAB [149] is trained with extra boundary information (B). PDB [40] adopts stronger data augmentation (DA). Lower is better.

	backbone	test	pose	expr.	illu.	mu	occu.	blur
ESR [14]	-	11.13	25.88	11.47	10.49	11.05	13.75	12.20
SDM [158]	-	10.29	24.10	11.45	9.32	9.38	13.03	11.28
CFSS [191]	-	9.07	21.36	10.09	8.30	8.74	11.76	9.96
DVLN [150]	VGG-16	6.08	11.54	6.78	5.73	5.98	7.33	6.88
Our approach	HRNetV2-W18	<b>4.60</b>	<b>7.94</b>	<b>4.85</b>	<b>4.55</b>	<b>4.29</b>	<b>5.44</b>	<b>5.42</b>
<i>Model trained with extra info.</i>								
LAB (w/ B) [149]	Hourglass	5.27	10.24	5.51	5.23	5.15	6.79	6.32
PDB (w/ DA) [40]	ResNet-50	5.11	8.75	5.36	4.93	5.41	6.37	5.81

TABLE 21

Facial landmark detection results (NME) on AFLW. DCFE [141] uses extra 3D information (3D). Lower is better.

	backbone	full	frontal
RCN [55]	-	5.60	5.36
CDM [169]	-	5.43	3.77
ERT [67]	-	4.35	2.75
LBF [116]	-	4.25	2.74
SDM [158]	-	4.05	2.94
CFSS [191]	-	3.92	2.68
RCPR [10]	-	3.73	2.87
CCL [192]	-	2.72	2.17
DAC-CSR [41]	-	2.27	1.81
TSR [101]	VGG-S	2.17	-
CPM + SBR [35]	CPM	2.14	-
SAN [34]	ResNet-152	1.91	1.85
DSRN [102]	-	1.86	-
LAB (w/o B) [149]	Hourglass	1.85	1.62
Our approach	HRNetV2-W18	<b>1.57</b>	<b>1.46</b>
<i>Model trained with extra info.</i>			
DCFE (w/ 3D) [141]	-	2.17	-
PDB (w/ DA) [40]	ResNet-50	1.47	-
LAB (w/ B) [149]	Hourglass	1.25	1.14

TABLE 22

Facial landmark detection results on COFW  $_{test}$ . The failure rate is calculated at the threshold 0.1. Lower is better for NME and  $FR_{0.1}$ .

	backbone	NME	$FR_{0.1}$
Human	-	5.60	-
ESR [14]	-	11.20	36.00
RCPR [10]	-	8.50	20.00
HPM [45]	-	7.50	13.00
CCR [39]	-	7.03	10.90
DRDA [174]	-	6.46	6.00
RAR [153]	-	6.03	4.14
DAC-CSR [41]	-	6.03	4.73
LAB (w/o B) [149]	Hourglass	5.58	2.76
Our approach	HRNetV2-W18	<b>3.45</b>	<b>0.19</b>
<i>Model trained with extra info.</i>			
PDB (w/ DA) [40]	ResNet-50	5.07	3.16
LAB (w/ B) [149]	Hourglass	3.92	0.39



TABLE 23  
Facial landmark detection results (NME) on 300W: common, challenging and full. Lower is better.

	backbone	common	challenging	full
RCN [55]	-	4.67	8.44	5.41
DSRN [102]	-	4.12	9.68	5.21
PCD-CNN [78]	-	3.67	7.62	4.44
CPM + SBR [35]	CPM	3.28	7.58	4.10
SAN [34]	ResNet-152	3.34	6.60	3.98
DAN [76]	-	3.19	5.24	3.59
Our approach	HRNetV2-W18	<b>2.87</b>	<b>5.15</b>	<b>3.32</b>
<i>Model trained with extra info.</i>				
LAB (w/ B) [149]	Hourglass	2.98	5.19	3.49
DCF (w/ 3D) [141]	-	2.76	5.22	3.24

TABLE 24  
Facial landmark detection results on 300W test. DCF [141] uses extra 3D information (3D). LAB [149] is trained with extra boundary information (B). Lower is better for NME,  $FR_{0.08}$  and  $FR_{0.1}$ , and higher is better for  $AUC_{0.08}$  and  $AUC_{0.1}$ .

	backbone	NME	$AUC_{0.08}$	$AUC_{0.1}$	$FR_{0.08}$	$FR_{0.1}$
Balt. et al. [4]	-	-	19.55	-	38.83	-
ESR [14]	-	8.47	26.09	-	30.50	-
ERT [67]	-	8.41	27.01	-	28.83	-
LBF [116]	-	8.57	25.27	-	33.67	-
Face++ [186]	-	-	32.81	-	13.00	-
SDM [158]	-	5.83	36.27	-	13.00	-
CFAN [175]	-	5.78	34.78	-	14.00	-
Yan et al. [162]	-	-	34.97	-	12.67	-
CFSS [191]	-	5.74	36.58	-	12.33	-
MDM [138]	-	4.78	45.32	-	6.80	-
DAN [76]	-	4.30	47.00	-	2.67	-
Chen et al. [23]	Hourglass	3.96	<b>53.64</b>	-	2.50	-
Deng et al. [30]	-	-	-	47.52	-	5.50
Fan et al. [37]	-	-	-	48.02	-	14.83
DReg + MDM [48]	ResNet101	-	-	52.19	-	3.67
JMFA [31]	Hourglass	-	-	54.85	-	1.00
Our approach	HRNetV2-W18	<b>3.85</b>	52.09	<b>61.55</b>	<b>1.00</b>	<b>0.33</b>
<i>Model trained with extra info.</i>						
LAB (w/ B) [149]	Hourglass	-	-	58.85	-	0.83
DCF (w/ 3D) [141]	-	3.88	52.42	-	1.83	-

## APPENDIX E

## MORE OBJECT DETECTION AND INSTANCE RESULTS ON COCO val2017

TABLE 25: More object detection and instance segmentation results on COCO val.  $AP^b$  and  $AP^m$  denote box mAP and mask mAP respectively. Most results are taken from [17] except that the results using HRNet are obtained by running the code at <https://github.com/open-mmlab/mmdetection>.

Backbone	LS	$AP^b$	$AP^b_{50}$	$AP^b_{75}$	$AP^b_S$	$AP^b_M$	$AP^b_L$	$AP^m$	$AP^m_{50}$	$AP^m_{75}$	$AP^m_S$	$AP^m_M$	$AP^m_L$
FCOS													
R-50 (c)	1x	36.7	55.8	39.2	21.0	40.7	48.4	-	-	-	-	-	-
R-101 (c)	1x	39.1	58.5	41.8	22.0	43.5	51.1	-	-	-	-	-	-
R-50 (c)	2x	36.9	55.8	39.1	20.4	40.1	49.2	-	-	-	-	-	-
R-101 (c)	2x	39.1	58.6	41.7	22.1	42.4	52.5	-	-	-	-	-	-
HRNetV2-W18	1x	35.2	52.9	37.3	20.4	37.8	46.1	-	-	-	-	-	-
HRNetV2-W32	1x	38.2	56.2	40.9	22.2	41.8	50.0	-	-	-	-	-	-
HRNetV2-W18	2x	37.7	55.9	40.1	22.0	40.8	48.5	-	-	-	-	-	-
HRNetV2-W32	2x	40.3	58.7	43.3	23.6	43.4	52.9	-	-	-	-	-	-
FCOS (mstrain)													
R-50 (c)	2x	38.7	58.0	41.4	23.4	42.8	49.0	-	-	-	-	-	-
R-101 (c)	2x	40.8	60.1	43.8	24.5	44.5	52.8	-	-	-	-	-	-
HRNetV2-W18	2x	38.1	56.3	40.6	22.9	41.1	48.6	-	-	-	-	-	-
HRNetV2-W32	2x	41.4	60.3	44.2	25.2	44.8	52.3	-	-	-	-	-	-
HRNetV2-W48	2x	42.9	61.9	45.9	26.4	46.7	54.6	-	-	-	-	-	-
X-101-64x4d	2x	42.8	62.6	45.7	26.5	46.9	54.5	-	-	-	-	-	-
Mask R-CNN													
R-50 (c)	1x	37.4	58.9	40.4	21.7	41.0	49.1	34.3	55.8	36.4	18.0	37.6	47.3
R-101 (c)	1x	39.9	61.5	43.6	23.9	44.0	51.8	36.1	57.9	38.7	19.8	39.8	49.5
R-50	1x	37.3	59.0	40.2	21.9	40.9	48.1	34.2	55.9	36.2	18.2	37.5	46.3
R-101	1x	39.4	60.9	43.3	23.0	43.7	51.4	35.9	57.7	38.4	19.2	39.7	49.7
HRNetV2-W18	1x	37.3	58.2	40.7	22.1	40.2	47.6	34.2	55.0	36.2	18.4	36.7	46.0
HRNetV2-W32	1x	40.7	61.9	44.6	25.1	44.4	51.8	36.8	58.7	39.5	20.9	40.0	49.3
X-101-32x4d	1x	41.1	62.8	45.0	24.0	45.4	52.6	37.1	59.4	39.8	19.7	41.1	50.1
X-101-64x4d	1x	42.1	63.8	46.3	24.4	46.6	55.3	38.0	60.6	40.9	20.2	42.1	52.4
R-50	2x	38.5	59.9	41.8	22.6	42.0	50.5	35.1	56.8	37.0	18.9	38.0	48.3
R-101	2x	40.3	61.5	44.1	22.2	44.8	52.9	36.5	58.1	39.1	18.4	40.2	50.4
HRNetV2-W18	2x	39.2	60.1	42.9	24.2	42.1	50.8	35.7	57.3	38.1	17.6	37.8	52.3
HRNetV2-W32	2x	42.3	62.7	46.1	26.1	45.5	54.7	37.6	59.7	40.3	21.4	40.5	51.2
X-101-32x4d	2x	41.4	62.5	45.4	24.0	45.4	54.5	37.1	59.4	39.5	19.9	40.6	51.3
X-101-64x4d	2x	42.0	63.1	46.1	23.9	45.8	55.6	37.7	59.9	40.4	19.6	41.3	52.5
Cascade Mask R-CNN													
R-50	1x	41.2	59.1	45.1	23.3	44.5	54.5	35.7	56.3	38.6	18.5	38.6	49.2
R-101	1x	42.6	60.7	46.7	23.8	46.4	56.9	37.0	58.0	39.9	19.1	40.5	51.4
X-101-32x4d	1x	44.4	62.6	48.6	25.4	48.1	58.7	38.2	59.6	41.2	20.3	41.9	52.4
X-101-64x4d	1x	45.4	63.7	49.7	25.8	49.2	60.6	39.1	61.0	42.1	20.5	42.6	54.1
R-50	20e	42.3	60.5	46.0	23.7	45.7	56.4	36.6	57.6	39.5	19.0	39.4	50.7
R-101	20e	43.3	61.3	47.0	24.4	46.9	58.0	37.6	58.5	40.6	19.7	40.8	52.4
HRNetV2-W18	20e	41.9	59.6	45.7	23.8	44.9	55.0	36.4	56.8	39.3	17.0	38.6	52.9
HRNetV2-W32	20e	44.5	62.3	48.6	26.1	47.9	58.5	38.5	59.6	41.9	18.9	41.1	56.1
HRNetV2-W48	20e	46.0	63.7	50.3	27.5	48.9	60.1	39.5	61.1	42.8	19.7	41.8	56.9
X-101-32x4d	20e	44.7	63.0	48.9	25.9	48.7	58.9	38.6	60.2	41.7	20.9	42.1	52.7
X-101-64x4d	20e	45.7	64.1	50.0	26.2	49.6	60.0	39.4	61.3	42.9	20.8	42.7	54.1
Hybrid Task Cascade													
R-50	1x	42.1	60.8	45.9	23.9	45.5	56.2	37.3	58.2	40.2	19.5	40.6	51.7
R-50	20e	43.2	62.1	46.8	24.9	46.4	57.8	38.1	59.4	41.0	20.3	41.1	52.8
R-101	20e	44.9	63.8	48.7	26.4	48.3	59.9	39.4	60.9	42.4	21.4	42.4	54.4
HRNetV2-W18	20e	43.1	61.5	46.8	26.6	46.0	56.9	37.9	59.0	40.6	18.8	39.9	55.2
HRNetV2-W32	20e	45.3	63.6	49.1	27.0	48.4	59.5	39.6	61.2	43.0	19.1	42.0	57.9
HRNetV2-W48	20e	46.8	65.3	51.1	28.0	50.2	61.7	40.7	62.6	44.2	19.7	43.4	59.3

Continued on next page

TABLE 25 – *Continued from previous page*

Backbone	Lr Schd	$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP_S^b$	$AP_M^b$	$AP_L^b$	$AP^m$	$AP_{50}^m$	$AP_{75}^m$	$AP_S^m$	$AP_M^m$	$AP_L^m$
HRNetV2-W48	28e	47.0	65.5	51.0	28.8	50.3	62.2	41.0	63.0	44.7	20.8	43.9	59.9
X-101-32x4d	20e	46.1	65.1	50.2	27.5	49.8	61.2	40.3	62.2	43.5	22.3	43.7	55.5
X-101-64x4d	20e	46.9	66.0	51.2	28.0	50.7	62.1	40.8	63.3	44.1	22.7	44.2	56.3
X-101-64x4d	28e	46.8	65.6	50.9	27.5	51.0	61.7	40.7	63.1	43.9	20.0	44.1	59.9