

High-Resolution Network with Attention Module for Human Pose Estimation

Tien-Dat Tran, Xuan-Thuy Vo, Duy-Linh Nguyen and Kang-Hyun Jo
School of Electrical Engineering, University of Ulsan

Ulsan (44610), South Korea

Email: (tdat,xthuy)@islab.ulsan.ac.kr, ndlinh301@mail.ulsan.ac.kr, acejo@ulsan.ac.kr

Abstract—Convolution neural networks (CNNs) have achieved the highest performance today not only for human posture prediction but also for other machine vision tasks (e.g., object identification, semantic segmentation, images classification). Furthermore, the Attention Module demonstrates their superiority over other conventional networks (AM). As a result, this work focuses on a useful feed-forward AM for CNNs. First, following a stage in the backbone network, feed the feature map into the attention module, which is separated into two dimensions: channel and spatial. The AM then multiplies these two feature maps and passes them on to the next level in the backbone. The network can collect more information in long-distance dependencies (channels) and geographical data, resulting in higher precision efficiency. Our experimental results would also show a difference between the employment of the attention module and current methodologies. As a result of the switch to a High-resolution network (HRNet), the predicted joint heatmap keeps accuracy while reducing the number of parameters compared to the baseline-CNN backbone. In terms of AP, the suggested design outperforms the baseline-HRNet by 2.0 points. Furthermore, the proposed network was trained using the COCO 2017 benchmarks, which are currently available as an open dataset.

Index Terms—deep learning, attention module, high-resolution network, human pose estimation.

I. INTRODUCTION

In today's contemporary world, 2D human pose estimate plays an important but challenging function in computer vision, serving numerous objectives such as human pose estimation [1], [2], activity recognition [3], [4], human re-identification [5], [6], or 3D human pose estimation [7], [8]. Human pose's main goal is to identify bodily sections for human body joints. Spatial and channel data are vital in improving the precision of key point regression. As a result, this research will concentrate on how to teach the network more about attention information.

Deep convolutional neural networks have recently attained state-of-the-art performance, according to recent breakthroughs. Most existing approaches route the input through a network, which is generally made up of high-to-low resolution subnetworks connected in series, before increasing the resolution. Hourglass [9], for example, restores high resolution using a symmetric low-to-high process. SimpleBaseline [10] generates high-resolution representations using a few transposed convolution layers. Furthermore, dilated convolutions are employed to enlarge the latter layers of a high-to-low resolution network (e.g., VGGNet or ResNet) [11], [12].

Deep convolution of neural networks has now stored significant advances in human posture [13], [14]. These networks, however, still have a lot of challenges to sort out. First and foremost, how can accuracy be improved in various types of networks? (e.g., real-time network, accuracy network). Second, while updating or modifying a network, it is frequently necessary to examine its speed. Last but not least, the present network must improve accuracy while remaining as quick as feasible. This research describes an unique network and the attention module's dependability in terms of speed and accuracy. The suggested experiment compares using and not using the attention module. The experiment also differs from the Simple Baseline [10] experiment, which did not employ the attention mechanism and instead used the transpose convolution [15] for upsampling. Our experiment would focus on how efficient and cost-effective each network situation is.

In particular, our technique was based on a simple fine-tune attention module [16], which demonstrated a considerable improvement in mean Average Precision (mAP). Inspired by VGG16 [11], the suggested network attempts to enhance the spatial attention module (SAM) by employing two 3×3 convolution layers rather than a 7×7 convolution layer. The network maintains the mAP while lowering the implementation cost by using 3×3 kernel. Furthermore, the number of parameters was reduced, resulting in an increase in network speed. To make clear about modify AM, our network increase 4.7 point in AP for accuracy and only increase around 16.5 percent of parameters compared with the Attention mechanism baseline [16] when used High-Resolution Network [17] as a backbone network. This research offers a novel network attention module that can readily react to a variety of difficulties in numerous applications, such as object identification, images classification, and human position estimation. The proposed method computes joint human pose estimations based on feature map recovery using an up-sampling network.

II. RELATED WORK

2D-Human Pose Estimation The most important aspect of human pose estimate is joint detection and its interaction with spatial space, as seen in Fig.1. DeepPose [18], Simple baseline makes use of joint prediction using an end-to-end network with a larger parameter. Later, Newell with the Stacked hourglass network [9] reduces the amount of settings while maintaining great accuracy. To represent local joints, all of the approaches

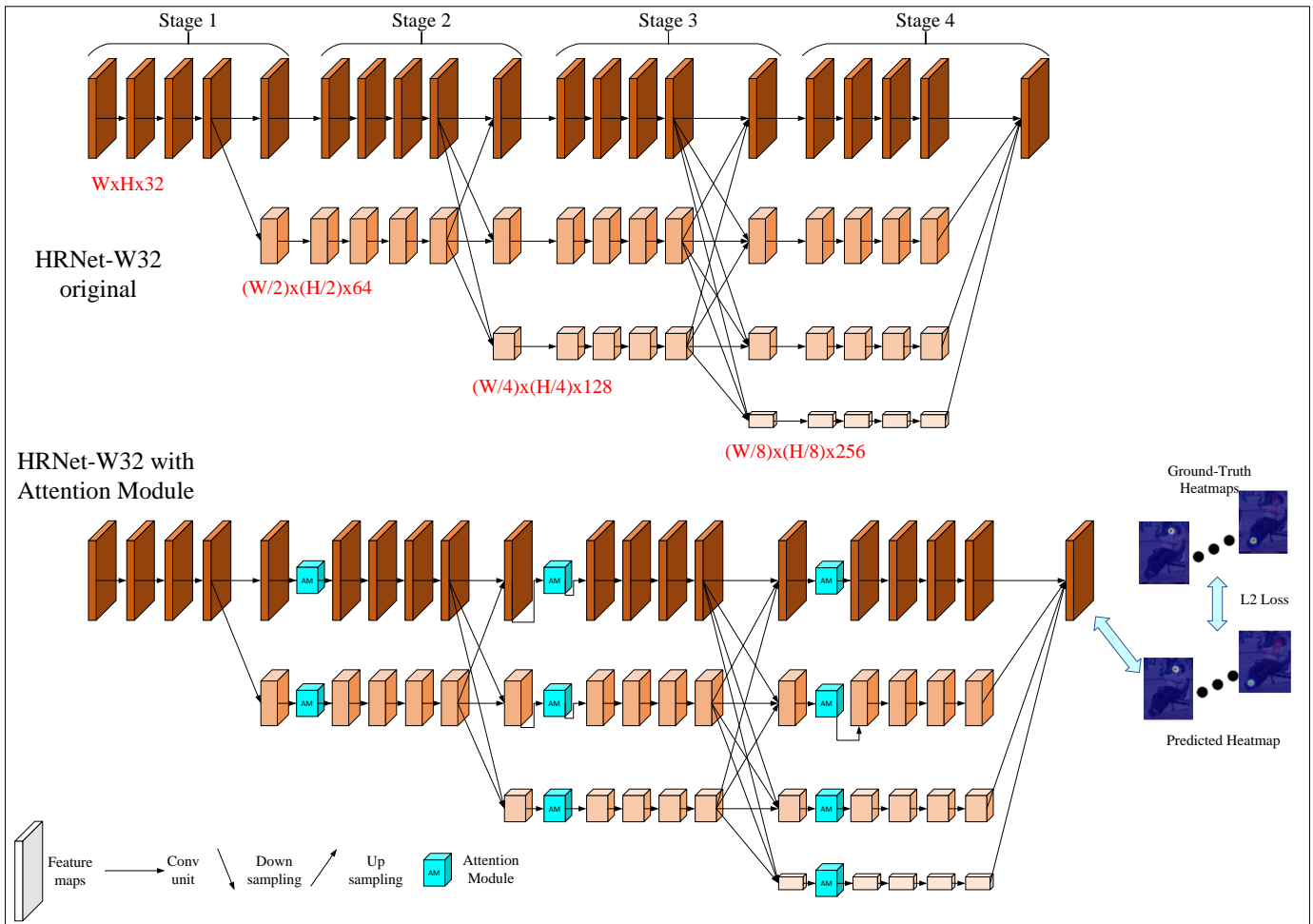


Fig. 1. Illustrating the design of the proposed 2D-human-pose estimation network. The suggested approach split the system 4 stages, each stage was connected by an attention module.

employed Gaussian distributions. After that, a convolution neural network was utilized to estimate human posture estimation. To minimize employment costs, they must reduce the number of parameters, and using appropriate attention approaches will reduce the network's parameter. As a result, the suggested strategy focuses on the employed attention module while increasing accuracy and decreasing the number of parameters.

On the other side, for increasing network performance, a 3×3 kernel size outperforms a 7×7 kernel size. However, in certain more sophisticated and expensive architectures, the 7×7 kernel size provides more precision. In comparison, our attention module gives a sufficient perspective for network design, with a limited number of parameters and high speed or a larger number of parameters and lower speed. The article then demonstrates how the attention module will function in each procedure and outcome.

High resolution network: Most convolutional neural networks for keypoint heatmap estimation are composed of a stem subnetwork, similar to a classification network, that decreases the resolution, a main body that produces representations with

the same resolution as its input, and a regressor that estimates the heatmaps where the keypoint positions are estimated and then transformed in full resolution. The main body primarily employs a high-to-low and low-to-high structure, which may be supplemented by multi-scale fusion and intermediate (deep) supervision.

In parallel, High Resolution network connects high-to-low subnetworks. It keeps high-resolution representations throughout the process, allowing for spatially exact heatmap estimate. It produces consistent high-resolution representations by repeatedly merging the representations created by the high-to-low subnetworks. Our technique differs from most previous efforts in that it requires a distinct low-to-high upsampling procedure as well as aggregate low-level and high-level representations. Without the need of intermediate heatmaps supervision, the technique is superior in keypoint identification accuracy and efficient in computing complexity and parameters.

Attention mechanism: Human visualization is vital in computer vision, and a variety of focus processing methods are being made to improve the efficiency of CNNs. Wang et al.

[19] also proposed a non-local network for gathering long-distance interdependence. SKNet [20] integrated the SENet Channel Focus Module with the Inception Multi-Branch Convolution, which was inspired by SENet [21] and Inception [22]. Furthermore, the Module for Spatial Focus is derived from Google’s STN [23], which gathers the background data of the feature maps. Furthermore, the attention module provides several benefits for saliency detection, multi-label categorization, and individual identification.

The suggested approach in this research was inspired by the CBAM network [24] to create the effective between both channel and spatial module by employing element-wise multiplication. Following that, the feature map adds to the previous feature map to merge the old and new information from the AT module.

III. METHODOLOGY

A. Network architecture

Backbone network The backbone network includes HRNet-W32 and HRNet-W48 [17], as shown in Figure 1 for a full architecture. Each HRNet has four phases, which include residual blocks and connections. The input RGB image shrinks the size to 256×192 (HRNet-W32, HRNet-W48), the feature maps traverse each column block, and the resolution of $W \times H$ drops twice for each stage. Finally, after travelling down the spine, the function map’s size is decreased to $\frac{W}{16} \times \frac{H}{16}$ with 256 channels at the last bottom layer of network. However, the backbone network will only use the first subnetwork which keeps the size is $W \times H$ until the end of regression. Furthermore, the size of the channels would be doubled at each stage. It progresses from 32 after the first block to 256 in the last layer. The backbone network’s job is to collect information and feature maps from the input picture and transmit them to the Training System, which uses cross entropy loss to predict human joints.

After extracting the information using the backbone network, the upsampling network recovers the information by taking the feature map from the final layer of the backbone network and upsampling it. Following that, the feature map will practice with Ground-truth Heat Maps, as shown in Fig.1. The default heat map size is same with the original images 256×192 for images worth 256×192 and 384×288 for images worth 384×288 . In order to match the size of the feature maps throughout the training phase, the heat maps must grasp the image’s scale. For regression, the network will utilize these heat maps and the ground truth heatmap to generate the predicted main point. This article employs the up-sampling module for the up-sampling network, which consists of one bilinear [25] layer and one convolution layer for the down-sampling (in Figure 1). The residual block contains both batch normalization and ReLU [26].

Attention Module The Attention Mechanism is made up of two primary components, as shown in Fig.2. First, the feature map was sent to the channel attention module following block one in the backbone network (CAM). The feature map in CAM uses global average pooling to reduce the feature map from

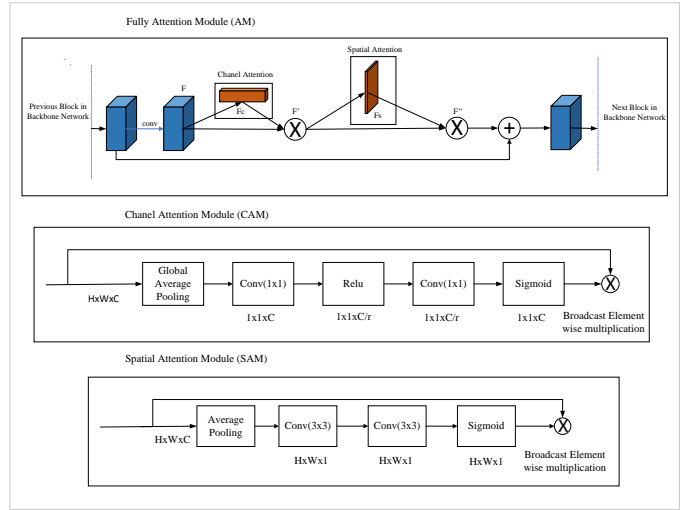


Fig. 2. Channel Attention Module (CAM) and Spatial Attention Module (SAM) Architecture (SAM). At comparison, this picture depicts the description of the attention module, which includes the channel and the spatial module in the center and bottom of the list, respectively, and the whole attention module at the top.

$H \times W \times C$ to $1 \times 1 \times C$. It first passes through the convolution layer, which converts the feature map to $1 \times 1 \times \frac{C}{r}$, where r is the reduction ratio and r is set to 16. The weight was then triggered by the CAM using the ReLU. The last stage in CAM is to employ a 1×1 convolution layer to restore the channel to $1 \times 1 \times C$ and to normalize the feature map using the sigmoid. The information for CAM were then combined using element-wise multiplication.

The feature map will be supplied into the Spatial Attention Module after passing through the CAM (SAM). The feature map in SAM takes the average channel pooling from $H \times W \times C$ to $H \times W \times 1$. Following pooling, two 3×3 convolution layers were utilized to extract the spatial information attribute diagram, and the final step in SAM is identical to the CAM shown in Figure 2. Finally, the intended solution employed element-wise extensions to the original feature map and the feature map after AT to be merged, as well as a new feature map for the next backbone network block.

B. Loss Function

Heat maps are used in this work to illustrate body joint locations for the loss function. As the ground-truth position in Fig. 1 by $a = \{a_k\} k = 1^K$, where $x_k = (x_k, y_k)$ is the spatial coordinate of the k th body joint in the image. The ground-truth heat map value H_k is then constructed using the Gaussian distribution with the mean a_k and variance \sum as shown below.

$$H_k(p) \sim N(a_k, \sum) \quad (1)$$

where $p \in \mathbb{R}^2$ represents the coordinate, and \sum is experimentally defined as an identity matrix \mathbf{I} . The last layer of the neural network predicts K heat maps, i.e., $\hat{S} = \{\hat{S}^k\} k = 1^K$

With out Attention Module



With Attention Module



Fig. 3. Predicted Heat-map before and after used Attention Module

for K body joints. A loss function is defined by the mean square error, which is calculated as follows:

$$L = \frac{1}{NK} \sum_{n=1}^N \sum_{k=1}^K \left\| s_k - \hat{s}_k \right\|^2 \quad (2)$$

N denotes the number of samples in the training session. Using information from the backbone network’s last layer, the network generated prediction heat maps using ground-truth heat maps.

IV. EXPERIMENTS

A. Experiment Setup

Dataset. The proposed technique uses the Microsoft COCO 2017 dataset [27] throughout the experiments. This dataset comprises around 200K pictures and 250K human samples, each with 17 keypoint labels. The study’s data collection includes three folders: train set, validation set, and test-dev set, each having training, validation, and testing photos. Furthermore, the validation and training annotations are open to the public and are accompanied by the original. **Evaluation metrics.** This paper utilized Object Keypoint Similarity (OKS) for COCO [27] with $OKS = \frac{\sum_i \exp(-d_i^2/2s^2k_i^2)\delta(v_i>0)}{\sum_i \delta(v_i>0)}$. In this case, d_i is the Euclidean distance between the predicted keypoint and the groundtruth, v_i is the target’s visibility flag, s is the object scale, and k_i is a keypoint for each joint. The standard average accuracy and recall score are then computed. In table I, AP and AR are the averages from OKS=0.5 to OKS=0.95, with AP^M representing medium objects and AP^L representing large objects.

Implementation details The suggested technique employed data increase in model training, such as flip, rotation at 40 degrees by design, and scale, which put the factor at 0.3. For training images, set the batch size to 4 and utilize the shuffle function. The total number of epochs in our experiment is 210, with the baseline learning-rate set at 0.001 and multiplied by 0.1 (learning rate factor) at the 170-th and 200-th epoch. The momentum is 0.9, and the Adam optimizer [28] was employed.

All experiments are carried out using the Pytorch framework and tested on two datasets. The picture input resolution was reduced to 256x192. The model was trained using CUDA 10.2 and CuDNN 7.3 on a single NVIDIA GTX 1080Ti GPU.

B. Experiment Result

TABLE I
THE RESULT FOR APPLY THE ATTENTION MODULE FOR EACH STAGE OF HRNET

Backbone	Stage	#Param	mAP
HRNet-W32	-	28.5M	74.4
HRNet-W32	1	30.2M	75.5
HRNet-W32	1+2	32.9M	76.0
HRNet-W32	1+2+3	36.4M	76.4

The suggested technique compares each circumstance while adding the attention module for each step from stage 1 to stage 3, as shown in Table 1. The Average Precision (AP) demonstrates that using AM in the first stage gains 1.1 in mAP, which boosts accuracy more than using AM in the second and third stages. Furthermore, the AP is enhanced by 1.5 percent, 2.2 percent, and 2.7 percent, respectively, while the number of parameters grows by 5.96 percent, 15.4 percent, and 27.7 percent for adding AM with stages 1, 2, and 3. In our proposed network, we used only 2 blocks of AM in stage 1, 3 blocks for stage 2 and 4 blocks for stage 3.

TABLE II
THE RESULT FOR APPLY THE ATTENTION MODULE FOR EACH SUB-NETWORK OF HRNET

Backbone	Sub-network	#Param	mAP
HRNet-W32	-	28.5M	74.4
HRNet-W32	1	31.1M	75.4
HRNet-W32	1+2	33.8M	75.9
HRNet-W32	1+2+3	35.5M	76.3
HRNet-W32	1+2+3+4	36.4M	76.4

As shown in Table 2, the proposed approach compares each case while adding the attention module for each step from sub-network 1 to sub-network 4. The Average Precision (AP) shows that utilizing AM in the first sub-network results in a 1.0 increase in mAP, which improves accuracy more than using AM in the second, third, and fourth sub-networks. Furthermore, the AP increases by 1.3 percent, 2.0 percent, 2.6 percent, and 2.7 percent, respectively, while the number of parameters increases by 9.1 percent, 18.6 percent, 24.5 percent, and 27.7 percent when AM with sub-stages 1, 2, 3, and 4 is included. In our suggested network, we employed three blocks of AM in the first sub-network, three blocks in

the second sub-network, two blocks in the third sub-network, and one block in the final sub-network.

COCO datasets result Our result was estimate on COCO validation dataset. The AP in the suggested approach is greater than the Basic High-Resolution benchmark in all situations of 1.7 AP, 1.3 AP in HRNet-32, HRNet-W48, respectively. Furthermore, the average recall (AR) is 1.4 points higher in the case of HRNet-W32 and 1.2 points higher in the case of HRNet-W48. The visualize result can see in Fig.3 which show that used attention module make the predicted heat map get more accurate. Figure 4 show the qualitative result for the COCO 2017 dataset.

However, human pose estimation, like many other designs today, has a number of issues that must be addressed. The first issue was that the images had hidden joints that were hard to train and anticipate. Second, low-resolution human photos must be correctly removed for human body joints. Following that are images of crowd scenarios, in which it is frequently difficult to determine all of the locations of the joints for all participants. Finally, there is a scarcity of information on images with incomplete parts for evaluating human postures.

V. CONCLUSION

This research shows the effect of the attention module on CNNs, with a focus on High-Resolution networks. Furthermore, our work demonstrates that by not increasing the amount of parameters, the attention module utilized has a bigger effect. On the other hand, the Attention Module highlighted the critical feature map rather than the other component. As a result, the network will improve efficiency, notably for various activities in the field of computer vision. Future research will focus on defining specific applications or settings to be included in our study, such as the surveillance system and the 3D human pose estimation. Another challenge is related to the limitations in assessing human exposure, which restricts the network's accuracy.

ACKNOWLEDGEMENT

REFERENCES

- [1] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," 2016.
- [2] C.-J. Chou, J.-T. Chien, and H.-T. Chen, "Self adversarial training for human pose estimation," 2017.
- [3] Z. Hussain, M. Sheng, and W. E. Zhang, "Different approaches for human activity recognition: A survey," 2019.
- [4] E. Kim, S. Helal, and D. Cook, "Human activity recognition and pattern discovery," *IEEE Pervasive Computing*, vol. 9, no. 1, pp. 48–53, Jan 2010.
- [5] X. Yang, M. Wang, and D. Tao, "Person re-identification with metric learning using privileged information," *CoRR*, vol. abs/1904.05005, 2019. [Online]. Available: <http://arxiv.org/abs/1904.05005>
- [6] W. Li, R. Zhao, and X. Wang, "Human reidentification with transferred metric learning," in *Asian Conference on Computer Vision (ACCV)*, 11 2012, pp. 31–44.
- [7] C. Chen and D. Ramanan, "3d human pose estimation = 2d pose estimation + matching," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017, pp. 5759–5767.
- [8] S. Li, L. Ke, K. Pratama, Y. Tai, C. Tang, and K. Cheng, "Cascaded deep monocular 3d human pose estimation with evolutionary training data," *CoRR*, vol. abs/2006.07778, 2020. [Online]. Available: <https://arxiv.org/abs/2006.07778>

- [9] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," *CoRR*, vol. abs/1603.06937, 2016. [Online]. Available: <http://arxiv.org/abs/1603.06937>
- [10] B. Xiao, H. Wu, and Y. Wei, "Simple baselines for human pose estimation and tracking," *CoRR*, vol. abs/1804.06208, 2018. [Online]. Available: <http://arxiv.org/abs/1804.06208>
- [11] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2015.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2015.
- [13] G. Moon, J. Y. Chang, and K. M. Lee, "Posefix: Model-agnostic general human pose refinement network," 2018.
- [14] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, "Deepcruc: A deeper, stronger, and faster multi-person pose estimation model," 2016.
- [15] V. Dumoulin and F. Visin, "A guide to convolution arithmetic for deep learning," 2016.
- [16] T.-D. Tran, X.-T. Vo, M.-A. Russo, and K.-H. Jo, "Simple fine-tuning attention modules for human pose estimation," in *International Conference on Computational Collective Intelligence*. Springer, 2020, pp. 175–185.
- [17] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," 2019.
- [18] A. Toshev and C. Szegedy, "DeepPose: Human pose estimation via deep neural networks," *CoRR*, vol. abs/1312.4659, 2013. [Online]. Available: <http://arxiv.org/abs/1312.4659>
- [19] X. Wang, R. B. Girshick, A. Gupta, and K. He, "Non-local neural networks," *CoRR*, vol. abs/1711.07971, 2017. [Online]. Available: <http://arxiv.org/abs/1711.07971>
- [20] X. Li, W. Wang, X. Hu, and J. Yang, "Selective kernel networks," 2019.
- [21] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, "Squeeze-and-excitation networks," 2017.
- [22] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," 2016.
- [23] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial transformer networks," 2015.
- [24] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," 2018.
- [25] M. Mastyło, "Bilinear interpolation theorems and applications," *Journal of Functional Analysis*, vol. 265, p. 185–207, 07 2013.
- [26] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," 2015.
- [27] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," *CoRR*, vol. abs/1405.0312, 2014. [Online]. Available: <http://arxiv.org/abs/1405.0312>
- [28] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations*, 12 2014.
- [29] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," 2017.

TABLE III
COMPARISON ON COCO VALIDATION DATASET. AM IS MEAN ATTENTION MODULE

Method	Backbone	Input size	#Params	AP	AP^{50}	AP^{75}	AP^M	AP^L	AR
8-Stage Hourglass [9]	8-Stage Hourglass	256×192	25.1M	66.9	-	-	-	-	-
Mask-RCNN [29]	ResNet-50-FPN	256×192	-	63.1	87.3	68.7	57.8	71.4	-
SimpleBaseline [10]	ResNet-50	256×192	34.0M	70.4	88.6	78.3	67.1	77.2	76.3
SimpleBaseline [10]	ResNet-101	256×192	53.0M	71.4	89.3	79.3	68.1	78.1	77.1
SimpleBaseline [10]	ResNet-152	256×192	68.6M	73.7	91.9	81.1	70.3	80.0	79.0
Fine-tuning AM [16]	ResNet-50	256×192	31.2M	71.4	91.6	78.6	68.2	75.7	76.3
Fine-tuning AM [16]	ResNet-101	256×192	50.2M	72.3	92.0	79.4	68.3	77.1	77.1
HRNetBaseline [17]	HRNet-W32	256×192	28.5M	74.4	90.5	81.9	70.8	81.0	79.8
HRNetBaseline [17]	HRNet-W48	256×192	63.6M	75.1	90.6	82.2	71.5	81.8	80.4
HRNet + our AM	HRNet-W32	256×192	36.4M	76.1	91.0	82.7	71.5	82.9	81.2
HRNet + our AM	HRNet-W48	256×192	71.8M	76.4	91.1	83.1	72.2	83.3	81.4



Fig. 4. Qualitative result for human pose estimation in COCO2017 test-dev set