

# Thesis

Filip Kronström

January 20, 2021

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Medical background . . . . .	1
1.1.1	POEs . . . . .	1
1.2	Machine learning . . . . .	1
1.3	Human pose estimation . . . . .	1
1.3.1	Datasets . . . . .	2
1.3.2	Pose estimation models . . . . .	2
1.4	Time series classification . . . . .	5
<b>2</b>	<b>Methods</b>	<b>6</b>
2.1	Overview . . . . .	6
2.2	Data . . . . .	6
2.2.1	Single-leg squat, SLS . . . . .	6
2.2.2	Stair descending, SD . . . . .	6
2.2.3	Forward Lunge, FL . . . . .	7
2.3	Body part localization . . . . .	7
2.3.1	Preprocessing . . . . .	7
2.3.2	Pose estimation . . . . .	7
2.4	Classification . . . . .	7
2.4.1	Preprocessing and dataset blabla.. . . .	7
2.4.2	Models . . . . .	8
	<b>References</b>	<b>9</b>

# Chapter 1

## Introduction

skriv om risker med bias fr dataset osv...

### 1.1 Medical background

skriv om acl och varf;r detta arbete beh;vs related work, se ref i mendeley osv

#### 1.1.1 POEs

...

### 1.2 Machine learning

skriv om att ml blivit s[ stort sen alexnet och gpu osv

skriv om backprop o s[nt h'r?

The emergence of computing power discussed in Section 1.2 allowing deeper networks shown by AlexNet [12]

### 1.3 Human pose estimation

Human pose estimation is a well explored problem which, like many other computer vision tasks has developed rapidly in the recent years. The reasons behind this progress can mainly be explained by two factors. Firstly the emergence of computing power discussed in Section 1.2, allowing more powerful deep learning models. Secondly several datasets with images labeled with human body joints has been made available [3]. These datasets not only provide data, but also introduces competition in the research community making it possible to compare the results of different approaches.

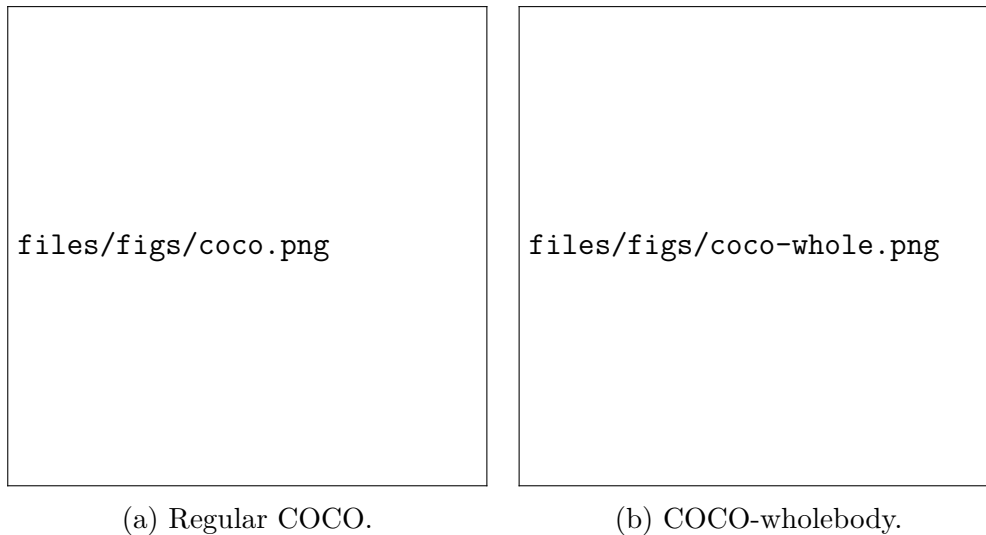


Figure 1.3.1: Keypoints for the two COCO datasets.

### 1.3.1 Datasets

Some of the widely used datasets today are Max Planck Institute for Informatics (**MPII**) [2], Microsoft Common Objects in Context (**COCO**) [15], AI Challenger Human Keypoint Detection (**AIC-HKD**) [24], and COCO-wholebody [11].

The COCO dataset consists of 328k images containing 91 different object types. The images come from Google, Bing, and Flickr image search and are mainly hand annotated through Amazon Mechanical Turk. The interesting part of the dataset for this work is the one with human poses. In total there are 250k instances of people labeled with joint locations [15]. The joints, 17 per person, in the dataset can be seen in Figure 1.3.1a. Along with the datasets containing body keypoints mentioned above there are also datasets with dense keypoints for specific bodyparts, e.g. OneHand10k [23]. COCO-wholebody is an attempt to combine these two types of datasets by extending COCO with dense keypoints at hands, feet, and faces. The resulting 133 joints can be seen in Figure 1.3.1b.

### 1.3.2 Pose estimation models

The human pose estimation (HPE) problem has been explored since long before the most recent deep learning era. Pictorial Structures were introduced by Fishler and Elschlager in the 1970s. This meant identifying individual parts or features in images modeled with pair-wise spring-like connections [6]. LeCun et al. showed that convolutional neural networks (CNN) could be trained with backpropagation [13] and how gradient based learning together with CNN outperformed other image analysis methods at the time [14]. Today CNNs form the basis in most computer vision methods, including the state-of-the-art (SOTA) HPE models [18, 20]. In 2014 Toshev and Szegedy successfully presented

DeepPose, the first HPE method based on deep neural networks (DNNs) [21]. Today's HPE methods are generally categorized as top-down or bottom-up approaches. This has to do with how they handle multiple persons. Bottom-up models starts by finding all keypoints for all persons in an image and then match them together to form persons. Top-down models on the other hand starts by finding bounding boxes for all individuals and then identifies keypoints for one person at a time. The sequential nature of the top-down methods and the fact that two models are needed means that bottom-up models scale better with the number of persons to analyze. However top-down models tends to be more accurate [4]. Since the application we present requires single person recognition only the top-down SOTA methods presented below are used.

### **High-Resolution Net (HRNet)**

Sun et al. [20] presented the HRNet architecture in 2019, initially for HPE, but also for other computer vision tasks such as semantic segmentation and object detection. Such problems had traditionally been solved using networks built on high-to-low resolution convolutions with increasing numbers of feature maps (e.g ResNet [8], VGGNet [19]). The classification task was solved in the low-resolution space and then transformed back to form the high-resolution representation needed for e.g. the HPE. Sun et al's proposed architecture preserves a high resolution representation throughout the network. It does so while also producing low-resolution/high dimensional representations suitable for classification.

The network architecture is shown in Figure 1.3.2 and consists of four stages (blue blocks in depth direction in Figure) with convolutional layers. After each stage a new low-resolution representation is created by performing strided convolutions. At the same time the existing representations also exchange information by either nearest neighbor upsampling or strided convolutions. The  $K$  estimated keypoints are represented as heatmaps,  $\{\mathbf{H}_1, \dots, \mathbf{H}_K\}$ , indicating the locations. These heatmaps are formed from the last high-resolution feature map (top right in Figure 1.3.2) [20]. Although a high resolution heatmap is desirable as it gives smaller quantization errors, the computational cost increases quadratically with the size [25]. Hence, to make the computations feasible the input image is downsampled through strided convolutions, resulting in a four times smaller heatmap [22].

### **Distribution-Aware coordinate Representation of Key-point (DARK)**

As discussed above a high resolution heatmap should result in higher accuracy, but is computationally expensive. Zhang et al. [25] propose a way to reduce the quantization error by i) analyzing the distributions of the predicted heatmaps, and ii) creating the training heatmaps in a slightly new fashion.

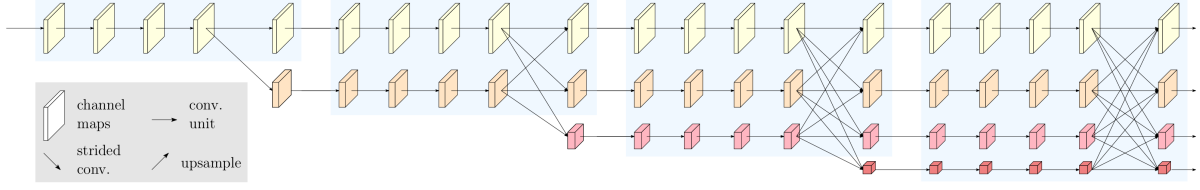


Figure 1.3.2: Network architecture for HRNet. The top row shows high resolution representations with fewer number of feature maps. Each step downwards reduces the resolution with a factor of two while the number of feature maps are doubled [22].

The actual keypoint location is found at the maximal activation of the heatmap. Since it is smaller than the actual image this turns into a sub-pixel localisation problem. Newel et al. [18] empirically found that a weighted average between the two highest activations, according to (1.1), yielded a good result.

$$\mathbf{p} = \mathbf{m} + \frac{1}{4} \frac{\mathbf{s} - \mathbf{m}}{\|\mathbf{s} - \mathbf{m}\|_2} \quad (1.1)$$

where:  $\mathbf{p}$  = predicted maximum  
 $\mathbf{m}$  = highest activation  
 $\mathbf{s}$  = second highest activation

This has been the de facto standard heatmap decoding, but Zhang et al. suggests using the fact that the heatmaps used for training usually are created as 2D Gaussian distributions, i.e. that the heatmaps can be expressed as (1.2).

$$\mathcal{G}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{2\pi |\Sigma|^{\frac{1}{2}}}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right) \quad (1.2)$$

where:  $\boldsymbol{\mu}$  = maximum of heatmap  
 $\mathbf{x}$  = pixel location  
 $\Sigma$  = diagonal covariance matrix

By Taylor expanding of the logarithm of (1.2) in the point  $\mathbf{m}$ , i.e. the point with the highest sampled activation, an expression for  $\boldsymbol{\mu}$  is obtained:

$$\boldsymbol{\mu} = \mathbf{m} - \left( \mathcal{D}''(\mathbf{m}) \right)^{-1} \mathcal{D}'(\mathbf{m}) \quad (1.3)$$

where:  $\mathcal{D}(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$ ,  
 i.e. the non constant term in the logarithm of (1.2)

The derivatives  $\mathcal{D}'(\mathbf{m})$  and  $\mathcal{D}''(\mathbf{m})$  are efficiently estimated from the heatmap. As this approach strongly assumes a Gaussian structure it is proposed to modulate the heatmap

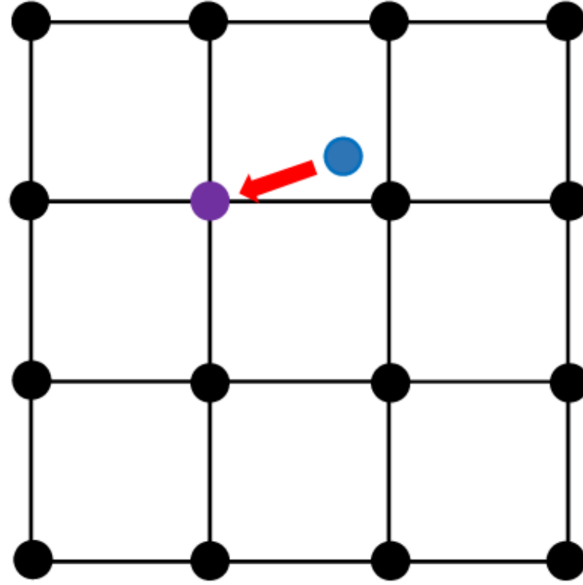


Figure 1.3.3: Quantization error due to off-grid keypoint location. Correct location (blue) represented by on-grid coordinate, here using floor quantization [25]. SKRIV ANTAGLIGEN NAGOT BATTRE HAR!!!!

before estimating the maximal activation. This is done by performing a convolution with a Gaussian kernel with the same covariance as the one used for the training data.

The second improvement suggested by Zhang et al. concerns the creation of the training heatmaps. Traditionally these have been created from the quantized keypoint locations, resulting in a slightly biased heatmap. In Figure 1.3.3 this would correspond to having the peak activation in the purple dot. By instead using the non-quantized location an unbiased heatmap is obtained. This would correspond to having the peak off-grid, in the blue dot in Figure 1.3.3.

## 1.4 Time series classification

[9]

# Chapter 2

## Methods

### 2.1 Overview

In this chapter the system for assessing POEs will be presented. This system is naturally divided into two parts where firstly the videos are analyzed. The first subsystem extracts body part coordinates of the subjects. This information is then passed to the second subsystem where it is used to calculate a score according to [17]. The data is presented in Section 2.2 and the two subsystems are described in Sections 2.3 and 2.4 respectively.

### 2.2 Data

The data available is in the form of videos each containing one subject, recorded from the front, performing three-five repetitions of specific motions. The motions are *Single-leg squat*, *Forward lunge*, *Stair descending*, *Forward lunge*, *Single leg hop for distance*, and *Side hop*. Each motion has a number of POE scores associated with them. The motion-POE combinations are shown in Table 2.2.1. In Sections 2.2.1-?? the motions and POEs evaluated in this project are described.

#### 2.2.1 Single-leg squat, SLS

The subject performed a squat standing on one leg to a knee angle of approximately  $60^\circ$ . The exercise was repeated five times and the entire movement was used to assess the POEs [1]. An illustration is shown in Figure ??.

#### 2.2.2 Stair descending, SD

The subject stepped down from a 30 cm step board. The exercise was repeated five times and POEs were evaluated for the loaded leg during loading phase [1]. An illustration is



Table 2.2.1: Motion-POE combinations available in the data.

POE \ Motion	Single leg squat	Stair descending	Forward lunge	Single leg hop for distance	Side hop
Trunk	x	x		x	x
Hip	x	x	x	x	x
Femoral valgus	x	x	x	x	x
Knee medial to foot position	x	x	x	x	x
Femur medial to shank	x	x	x	x	x
Foot	x				

shown in Figure ??.

### 2.2.3 Forward Lunge, FL

Helo I am now mispealing.

## 2.3 Body part localization

### 2.3.1 Preprocessing

... rotation, flip etc

### 2.3.2 Pose estimation

The pose estimation can be seen as a feature extraction and dimensionality reduction. The pose estimation is built around the open-source toolbox MMPose [16] from MMLab. Each frame is considered to be an independent image and is analyzed with the DARK-HRNet [20, 25] trained on the COCO-dataset [15] described in Section 1.3.

## 2.4 Classification

### 2.4.1 Preprocessing and dataset blabla..

Before assessing the POEs based on the body part positions a number of preprocessing steps are conducted. Firstly the data is resampled as the videos are recorded with a number of different frame rates (25, 30, and 60Hz). The resampling is performed using linear interpolation to a new sample frequency of 25Hz. This data is then low pass filtered through a fourth order Butterworth filter with a cutoff frequency of 2.5Hz. PLOT P[;VERF;RING F;R FILTER? ELLER P[ FILTRERAD DATA?

While the POE assessment, see Section ??, is performed on a per repetition basis the body part coordinates are extracted on a per video basis. Hence, each repetition should be extracted. This is done by finding peaks in the time series corresponding to certain body part positions. Which body part is used for this sequence splitting depends on which movement is analyzed. For SLS the  $y$ -coordinate of the right shoulder is used. The duration of each repetition varies significantly. With the duration of each repetition varying substantially between subjects padding in the time dimension is desirable. The reason for this is twofold, i) to simplify the handling of the data by storing it as a multidimensional array, and ii) to be able to train the eventual model in a more efficient manner using batches [7]. The padding is done by adding constant values of -1000 at the end of the sequences to some specified length. Details on how this is handled by the model are presented in Section ??

#### BESKRIV ALGORITM FÖR ATT DELA UPP I REPS

Finally the data is normalized to put the mean of the first five right hip-samples in the origin and the distance to the first five right shoulder-samples to one, according to (2.1).

$$\begin{aligned} (x, y)_i &= (x, y)_i - \overline{(x, y)}_{rh} \\ (x, y)_i &= \frac{(x, y)_i}{\|\overline{(x, y)}_{rs}\|_2}, \forall i \end{aligned} \tag{2.1}$$

where:  $\overline{(x, y)}_i$  = mean over first five samples for body part  $i$   
 $rh$  = right hip  
 $rs$  = right shoulder  
 $i$   $\in$  Available body parts

After these preprocessing steps a dataset  $\in \mathbb{R}^{N \times T \times F}$  is created consisting of  $N$  multivariate time series with  $F$  channels of length  $T$ . The input channels can be  $x$ - and  $y$ -coordinates of body parts as well as angles between body parts.

### 2.4.2 Models

The proposed models are neural networks with convolutional layers as feature extractors. The network is inspired by InceptionTime by Fawaz et al. [10], described in Section ??.

[5] lol

# Bibliography

- [1] Jenny Älmqvist Nae. “Is seeing just believing? Measurement properties of visual assessment of Postural Orientation Errors (POEs) in people with anterior cruciate ligament injury”. English. PhD thesis. Department of Health Sciences, June 2020. ISBN: 978-91-7619-940-4.
- [2] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. “2D Human Pose Estimation: New Benchmark and State of the Art Analysis”. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition*. 2014, pp. 3686–3693. DOI: 10.1109/CVPR.2014.471.
- [3] Yucheng Chen, Yingli Tian, and Mingyi He. “Monocular Human Pose Estimation: A Survey of Deep Learning-based Methods”. In: *Computer Vision and Image Understanding* 192 (June 2020). DOI: 10.1016/j.cviu.2019.102897. arXiv: 2006.01423.
- [4] Bowen Cheng, Bin Xiao, Jingdong Wang, Honghui Shi, Thomas S. Huang, and Lei Zhang. “HigherHRNet: Scale-Aware Representation Learning for Bottom-Up Human Pose Estimation”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (Aug. 2019), pp. 5385–5394. arXiv: 1908.10357.
- [5] Kevin Fauvel, Tao Lin, Véronique Masson, Élisabeth Fromont, and Alexandre Termier. “XCM: An Explainable Convolutional Neural Network for Multivariate Time Series Classification”. In: (Sept. 2020). arXiv: 2009.04796.
- [6] Martin A. Fischler and Robert A. Elschlager. “The Representation and Matching of Pictorial Structures Representation”. In: *IEEE Transactions on Computers* C-22.1 (1973), pp. 67–92. ISSN: 00189340. DOI: 10.1109/T-C.1973.223602.
- [7] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. Vol. 2016-December. IEEE Computer Society, Dec. 2016, pp. 770–778. ISBN: 9781467388504. DOI: 10.1109/CVPR.2016.90. arXiv: 1512.03385.
- [9] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre Alain Muller. “Deep learning for time series classification: a review”. In: *Data Mining and Knowledge Discovery* 33.4 (July 2019), pp. 917–963. ISSN: 1573756X. DOI: 10.1007/s10618-019-00619-1. arXiv: 1809.04356.

- [10] Hassan Ismail Fawaz, Benjamin Lucas, Germain Forestier, Charlotte Pelletier, Daniel F. Schmidt, Jonathan Weber, Geoffrey I. Webb, Lhassane Idoumghar, Pierre Alain Muller, and François Petitjean. “InceptionTime: Finding AlexNet for time series classification”. In: *Data Mining and Knowledge Discovery* 34.6 (Nov. 2020), pp. 1936–1962. ISSN: 1573756X. DOI: 10.1007/s10618-020-00710-y. arXiv: 1909.04939.
- [11] Sheng Jin, Lumin Xu, Jin Xu, Can Wang, Wentao Liu, Chen Qian, Wanli Ouyang, and Ping Luo. “Whole-Body Human Pose Estimation in the Wild”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 12354 LNCS (July 2020), pp. 196–214. arXiv: 2007.11858.
- [12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: ed. by Peter L. Bartlett, Fernando C. N. Pereira, Christopher J. C. Burges, Léon Bottou, and Kilian Q. Weinberger. 2012, pp. 1106–1114. URL: <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks>.
- [13] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. “Backpropagation Applied to Handwritten Zip Code Recognition”. In: *Neural Computation* 1.4 (Dec. 1989), pp. 541–551. ISSN: 0899-7667. DOI: 10.1162/neco.1989.1.4.541.
- [14] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2323. ISSN: 00189219. DOI: 10.1109/5.726791.
- [15] Tsung Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. “Microsoft COCO: Common objects in context”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. Vol. 8693 LNCS. PART 5. Springer Verlag, May 2014, pp. 740–755. DOI: 10.1007/978-3-319-10602-1\_48. arXiv: 1405.0312.
- [16] *MMPose - OpenMMLab Pose Estimation Toolbox and Benchmark*. URL: <https://github.com/open-mmlab/mmpose>.
- [17] Jenny Nae, Mark W. Creaby, Gustav Nilsson, Kay M. Crossley, and Eva Ageberg. “Measurement Properties of a Test Battery to Assess Postural Orientation During Functional Tasks in Patients Undergoing ACL Injury Rehabilitation”. In: *Journal of Orthopaedic & Sports Physical Therapy* 47.11 (Oct. 2017), pp. 1–42. ISSN: 0190-6011. DOI: 10.2519/jospt.2017.7270. URL: <http://www.jospt.org/doi/10.2519/jospt.2017.7270>.
- [18] Alejandro Newell, Kaiyu Yang, and Jia Deng. “Stacked hourglass networks for human pose estimation”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. Vol. 9912 LNCS. Springer Verlag, 2016, pp. 483–499. ISBN: 9783319464831. DOI: 10.1007/978-3-319-46484-8\_29. arXiv: 1603.06937.

- [19] Karen Simonyan and Andrew Zisserman. “Very deep convolutional networks for large-scale image recognition”. In: *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*. International Conference on Learning Representations, ICLR, Sept. 2015. arXiv: 1409.1556.
- [20] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. “Deep high-resolution representation learning for human pose estimation”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. Vol. 2019-June. IEEE Computer Society, June 2019, pp. 5686–5696. ISBN: 9781728132938. DOI: 10.1109/CVPR.2019.00584. arXiv: 1902.09212.
- [21] A. Toshev and C. Szegedy. “DeepPose: Human Pose Estimation via Deep Neural Networks”. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition*. 2014, pp. 1653–1660. DOI: 10.1109/CVPR.2014.214.
- [22] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang, W. Liu, and B. Xiao. “Deep High-Resolution Representation Learning for Visual Recognition”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020), pp. 1–1. DOI: 10.1109/TPAMI.2020.2983686.
- [23] Yangang Wang, Cong Peng, and Yebin Liu. “Mask-Pose Cascaded CNN for 2D Hand Pose Estimation from Single Color Image”. In: *IEEE Transactions on Circuits and Systems for Video Technology* 29.11 (Nov. 2019), pp. 3258–3268. ISSN: 15582205. DOI: 10.1109/TCSVT.2018.2879980.
- [24] Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipei Zhou, Guosen Lin, Yanwei Fu, Yizhou Wang, and Yonggang Wang. “AI Challenger : A Large-scale Dataset for Going Deeper in Image Understanding”. In: *arXiv* (Nov. 2017). arXiv: 1711.06475.
- [25] Feng Zhang, Xiatian Zhu, Hanbin Dai, Mao Ye, and Ce Zhu. “Distribution-Aware Coordinate Representation for Human Pose Estimation”. In: *Institute of Electrical and Electronics Engineers (IEEE)*, Aug. 2020, pp. 7091–7100. DOI: 10.1109/cvpr42600.2020.00712. arXiv: 1910.06278.