



Synthetic image translation for football players pose estimation

Michał Sypetkowski, Grzegorz Sarwas and Tomasz Trzcíński

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

February 20, 2019

Synthetic image translation for football players pose estimation

Michał Sypetkowski

(Institute of Computer Science, Warsaw University of Technology, Poland;
Sport Algorithmics and Gaming Sp. z o. o., Poland
m.sypetkowski@gmail.com)

Grzegorz Sarwas

(Institute of Control and Industrial Electronics,
Warsaw University of Technology, Poland;
Sport Algorithmics and Gaming Sp. z o. o., Poland
grzegorz.sarwas@ee.pw.edu.pl)

Tomasz Trzciński

(Institute of Computer Science, Warsaw University of Technology, Poland
t.trzcinski@ii.pw.edu.pl)

Abstract: In this paper, we present an approach for football players pose estimation on very low-resolution images. The camera recording the football match is far away from the pitch in order to register at least half of it. As a result, even using very high resolution cameras, the image area presenting every single player is very small. Additionally, variable weather conditions or shadows and reflections, make this aim very hard. Such images are very hard to annotate by human. In our research we assume lack of manually annotated training data from our target distribution. Instead of manual annotation of large dataset, we create simple python script for rendering synthetic images with perfect annotations. Then we train vanilla CycleGAN for transformation of raw synthetic images into more realistic. We use transformed images to train CPN model. Without bells and whistles, we achieve similar precision on our images as the same CPN model trained with COCO keypoints dataset.

Key Words: pose estimation, deep convolutional neural network, image translation, synthetic dataset

Category: I.3.3, I.2.10, I.4.8

1 Introduction

Visual analysis of football matches and training sessions is a demanding task, consisting of multiple aspects such as proper video acquisition, tracking in a multi-view system with occlusions, 3D calibration and human behavior analysis. The latter can be split in various conceptual and algorithmic problems, one of each is player's pose estimation. Human pose helps football analysts to validate players' mobility during match and ability to properly perform various game interceptions. In particular analysts check how often player uses non-dominant

leg during ball repossession. Accurate pose estimation is also a key step for higher level tasks as analysis of visibility of action for each player, or having an open position while ball pass receiving.

Visual tracking systems installed in football academies uses wide view cameras, spanning on whole pitch or near half. Depending on installation site cameras could be positioned near ground, producing substantial occlusions, or on a high pylons giving non-standard human view from above. Moreover, wide view cameras imposes very low-quality human visuals even for top tier recording hardware. We did not find any literature nor the databases with annotated human pose for the high view and low resolution scenario, what imposed the presented research problem. All successful pose estimation approaches concern high or medium resolution images. The literature presents two generalized approaches in that case. The first one is called bottom-up and the second is top-down. We tested multiple known state-of-the-art algorithms for pose estimation with our custom test images. The images have been acquired from real system with four high-view and high-class wide-view cameras. In next subsections we present analysis of related work in different pose estimation approaches.

1.1 2D multi-person bottom-up approaches

Bottom-up approach predicts all keypoints, which are considered as skeleton model parts in a single scene. Those are further assembled into full skeleton by assigning the parts to appropriate place in the model. In [4] the multiple-stage fully convolutional networks for estimating Part Confidence Map (heat map) and PAF (Part Affinity field - 2D vector field) have been considered. This solution uses multi-stage convolutional network that generates heat map and 2D vector field for each body part (e.g. right elbow, left wrist, neck). The affinity graph is build using 2D vector field part. Based on it, the 2D skeleton with a particular heuristic graph relaxation technique proposed in the article can be constructed. The approach presented in [21] achieved the best result in COCO 2016 Keypoint Detection Task, being valid proposition for solving our problem. Along with work of Simon [29], this approach has publicly available implementation called OpenPose [11]. Highest score on MPII multi-person pose dataset [1] got an approach presented in work of Newell et al.[25]. Authors trained a network to simultaneously output detections and group assignments. Output of their neural network consist of detection heatmaps with respective associative embeddings. Grouping body parts is performed by an algorithm based on thresholding the parts embeddings distances. This approach differs from other bottom-up approaches by the lack of separation between detection and grouping. An entire prediction is done at once by a single-stage, generic network based on a stacked hourglass architecture [26].

1.2 2D multi-person top-down approaches and single person pose estimation

Top-down approaches localize and crop all persons from an image at first, then solve the single person pose estimation problem (which becomes the main difficulty). Modern single person pose estimation techniques incorporate priors about a structure of human bodies. Best results in COCO 2017 Keypoint Detection Task [21] were achieved by Cascaded Pyramid Network [6]. This algorithm focuses on the "hard" keypoints (i.e. occluded, invisible and with non-trivial background). It is achieved by explicitly selecting the hard keypoints and backpropagating the gradients only from the selected keypoints.

Approach called Mask R-CNN [15], extends Faster R-CNN [28] by adding a branch for predicting an object mask in parallel with bounding box recognition. Using this simple modification the Mask R-CNN can be applied to keypoints detection. This approach achieves high results in all COCO 2017 challenges (i.e. object detection, object segmentation, keypoint detection).

Simon et al.[29] presents precise hand 2D keypoint detector. It introduces a semi-supervised training algorithm called Multiview Bootstrapping. Initially, the algorithm needs a set of annotated examples. The model is trained using only these examples at the beginning. Then, the model detects keypoints on unannotated examples with multiple camera views. Each multi-view example is then robustly 3D triangulated, and reprojected creating additional training set.

Stacked hourglass [26] achieves state-of-the-art result on MPII [1]. It presents a CNN architecture for bottom-up and top-down inference with residual blocks. Approach introduced by Ke et al.[19] aims to improve stacked hourglass [26] achieving the best score on MPII single person pose dataset.

1.3 Other approaches

Modern pose estimation approaches are already robust to blurring and low-resolution in general. Significantly improving their performance with simple methods, like heuristic data augmentation or upscaling the images with generic upscaling algorithms may be extremely hard with limited training data. A straightforward solution for improving the results on images from a specific distribution may be manual annotation of some examples (e.g. a few thousands) for training or fine-tuning existing state-of-the-art models. Manual annotations on low resolution images not only require immense amount of work, but also may be hard to be done precisely in our case.

In our approach we consider single person pose estimation on large dataset of blobs detected with external tracking system. We create synthetic dataset and improve it using modern achievements of Generative Adversarial Networks. In the following sections, we discuss selected existing approaches that use synthetic dataset, and use of GANs for pose estimation.

2 Synthetic datasets

In this section we present successful approaches that focus on generating large (practically infinite, but every distribution have it's effective variety limit) annotated synthetic or partially synthetic (e.g. [9]) datasets with minimal effort. They show that limited realism may provide enough training signal for current state-of-the-art object/keypoints detector models.

In [9] to create the dataset authors propose simply 'cut' real object instances and 'paste' them on random backgrounds (without any perspective or lighting adjustment). This process implemented in naive way would give the trained model possibility of exploiting subpixel discrepancies at the boundaries. To address this problem, the approach blends 'cut' objects into the background with heuristic methods. Additionally, it blends in the same way the distractor objects along with the correct ones. Synthetic data is then feed into the model along with the real data. In the end, such training set gives significantly higher performance, than the non-augmented dataset. (e.g. 51 AP instead of 42 AP on GMU Kitchen Scenes [10] dataset).

In [14] Gupta et al. introduce a method of 'pasting' synthetically rendered text into the real images with respect to the local region cues, i.e. surface geometry predicted with other models and local colors. Models trained with such dataset achieve high accuracy in the task of text detection in the wild.

In [23] authors create fully synthetic dataset using ray-trace rendered scenes — interiors of buildings. It shows that large-scale high-quality synthetic RGB datasets with task-specific labels can be more effective for pre-training than the large-scale real-world images dataset like ImageNet [7].

Many modern successful approaches concerning synthetic dataset use in some way real images to create it:

- [9] uses both real examples and synthetic. Synthetic images are made by simple editions of the real images, therefore they are not dependent on graphics renderings.
- [24] uses real examples (unpaired with synthetic) for improving the distribution of 3D rendered synthetic images.
- [14] 'pastes' rendered text into the real images (background).
- [23] does not use real images explicitly, but the 3D models of furniture may be using textures that are created at some degree using real photos.

In paper [9] authors suggest that state-of-the art detection methods like Faster-RCNN [28] care more about local region-based features for detection than the global scene layout. This fact somehow justifies their result.

In our previous article [30] we have shown robustness to low resolution and small distortions, of CPN [6] trained on large datasets. One may suspect, that in our case artifacts of rendered synthetic dataset will not cause the optimizer to find significant exploits, or stuck in a bad local minimum.

3 Generative Adversarial Networks for image generation and pose estimation

In recent years, we observe a rapid progress in results achieved by Generative Adversarial Networks for image generation. In this section, we review selected approaches and discuss application of GANs in pose estimation task.

Initial research concerning image generation using GANs was done by [12]. In recent years, there was many improvements for loss functions, model architecture and overall training process (DCGAN [27], LSGAN [22], SRGAN [20], StackGAN [31], Wasserstein GAN [2], Improved Wasserstein GAN [13]). Modern state-of-the-art approaches (StackGANv2 [32], Progressive growing of GANs [17]) can infer photo-realistic high-resolution images using multi-stage generator architecture. Recent paper [18] introduced an architecture that allows unsupervised separation and control of high, mid, and low-level attributes of high-resolution, photo-realistic generated images.

Adversarial PoseNet [5] presents an interesting approach that trains a GAN, with multi-task pose generator and two discriminator networks. It achieves state-of-the-art results on MPII [1] single person pose estimation dataset. The model consists of the generator network, the pose discriminator network and the confidence discriminator. Half of generated heatmaps represent keypoint locations and the other half occlusion predictions. The generator architecture is based on stacked hourglass architecture [26].

In pose estimation from low-resolution images, an idea worth consideration is generative upscaling. Modern generic upscaling deep learning methods are focused on minimizing the mean squared reconstruction error (MSE) [8]. SRGAN [20] is capable of inferring photo-realistic natural images for 4x upscaling factors. The approach uses GAN, trained using a perceptual loss function consisting both of an adversarial loss and a content loss. Such generic upscaling algorithms like these will not improve results on our dataset as we have shown in our preliminary article version [30], because of too low resolution and characteristic distortions caused during the scene recording (e.g. compression). Super-FAN [3] addresses the problem of generative upscaling of very low resolution images. It focuses on improving the quality of low resolution facial images and locating the facial landmarks on such images. The idea is to connect third network (Face Alignment Network) to GAN. This third network detects facial landmarks on the upscaled image. Generator loss includes additional component – landmark detection loss. Therefore it learns to generate face that fits geometrically.

3D Hand pose estimation approach [24] focuses on enhancing a synthetic dataset to make their distribution more like the distribution of the real images. It uses CycleGAN with an additional geometric consistency loss. The paper shows, that training with generated images significantly outperform standard augmentation techniques. Similar approach may be applied to pose estimation.

In our approach we propose to use CycleGAN for enhancing synthetic dataset for pose estimation. Our images are very low-resolution and the human body details are not visible on our images, therefore GANs are easier to learn their distribution.

4 Proposed approach

For our tests we gathered data using 4 cameras placed at the field corners. Because of resolution limits, in practice we can assume that for a given player only 2 cameras are close enough to produce usable visuals. The cameras are production class CCTV devices with 4K resolution and high compression bandwidth. Even though the crop factor around single player magnifies compression and optics artifacts, which renders high frequency data unusable. Low quality and viewing angle creates uncommon characteristics of the images. Comparing this scenario with the standard pose estimation datasets like COCO [21] and MPII [1] we can list main problems:

- Human based annotations are much more difficult and time consuming for our images. Some images have practically indistinguishable joint locations, even with much human time and effort spent
- Border areas of the pitch generates almost top down views, where the human parts are mostly occluded by upper body
- Images are blurred with non-deterministic distribution, which makes generic upscaling algorithms useless
- All players wear single-color clothes, which makes it harder to distinguish limbs (especially hands) from the body

In our preliminary article [30] we selected most efficient network architecture trained with external data (see Table 1). In all experiments, we used CPN [6] model (smaller version – with input resolution of 256x192 and based on ResNet50 [16]).

Our new approach consists of 4 steps:

¹ <https://github.com/CMU-Perceptual-Computing-Lab/openpose>

² <https://github.com/umich-vl/pose-hg-demo>

³ <https://github.com/wbenbihi/hourglassstensorflow>

⁴ <https://github.com/bearpaw/pytorch-pose>

⁵ <https://github.com/chenyilun95/tf-cpn>

Approach	Implementation / experiment	Training set	Language Library	corr. pose	corr. legs	N / A
PAF [4]	OpenPose ¹	COCO [21]	C++, Caffe	58	106	29
Stacked hourglass [26]	original implementation ² , 8-stack model	MPII [1]	Lua, Torch	142	203	-
	alternative implementation ³ , hg_refined_200, 4-stack model,	MPII [1]	Python, Tensorflow	29	90	
	alternative implementation - not official ⁴ , 8-stack model,	MPII [1]	Python, Pytorch	135	186	
CPN [6]	original implementation ⁵ , COCO.res50.256x192, snapshot_350.ckpt	COCO [21]	Python, Tensorflow	171	224	-
	SRGAN for upscaling			90	167	
	blurred images, 50 more epochs, lr 1.6e-5 (from COCO.res50.256x192)			155	206	
	COCO.res101.384x288, snapshot_350.ckpt			158	223	

Table 1: Selected human pose estimation implementation results (original and our experiments). The table contains results of experiments from our previous paper. Measured implementations vary in skeleton structure used as a reference, therefore the measurements are done without annotated testset. We’ve taken into account few the easiest football aspects for automation. We measured precision on 300 test images with human based decision, whether the answer is one of 4 classes: correct, only correct legs pose estimation, wrong pose, N/A. The human-based bias has been lowered by cross-checkup with industry football analyst but still may produce significant variance, opposed to keypoint-based difference metrics.

1. rendering synthetic dataset (see section 4.1),
2. training CycleGAN [33], using generated synthetic dataset as first distribution examples, and real players blobs for the second (see section 4.2),
3. training CPN [6], with synthetic dataset, cycled-synthetic dataset, and mixed with COCO [21] dataset,
4. measuring pose estimation accuracy on our benchmark (see section 4.3).

4.1 Rendering synthetic dataset

We use blender ⁶ for scene modelling and rendering, and ManuelbastioniLAB ⁷ for creating human 3D models. We use blender ray-trace rendering engine – cycles. We design armature pose distribution empirically – by randomizing bones Inverse Kinematics (IK) targets transform (with respect to the rest pose - A-pose) with normal distributions. One character armature has 8 IK targets in total: 2 for hands, legs, elbows, 1 for body center and a head look-at position. Each IK target, has hard-coded means and standard deviations for each axis, e.g.:

- hands IK targets have higher standard deviation on backward-forward axis than on left-right axis, because hands are moving usually switching between front and back position during running
- similarly feet IK targets have higher standard deviation on backward-forward axis than on left-right axis, because running is usually for forward movement
- mean of body center IK is lower than in the rest pose, because it is usually lower during dynamic actions like running or kicking

We randomize each IK target independently. Additionally, we constraint randomized pose with heuristics:

- in general, feet are standing on the ground (jumping positions are rare)
- if left foot is moved forward and standing on the ground, then probably right knee is bent backwards (as it is during the run)
- foot in the air (not standing on the ground) is usually rotated similarly to the corresponding calf
- heel may be lifted when leg is standing

We rendered 1000 football fields with constant camera position, various lighting angle, and randomly (with uniform distribution) placed and rotated 100 players. Each image has resolution 4000x3000 (like the original cameras). Then, we cut 100k training blobs from the large images. Example synthetic blobs are shown in figure 1.

4.2 CycleGAN-ing synthetic dataset

In experiments we train vanilla CycleGAN [33] architecture with 256x256 input/output. We use 100k synthetic blobs for first distribution, and 186K real

⁶ <https://www.blender.org/>

⁷ <http://www.manuelbastioni.com/>

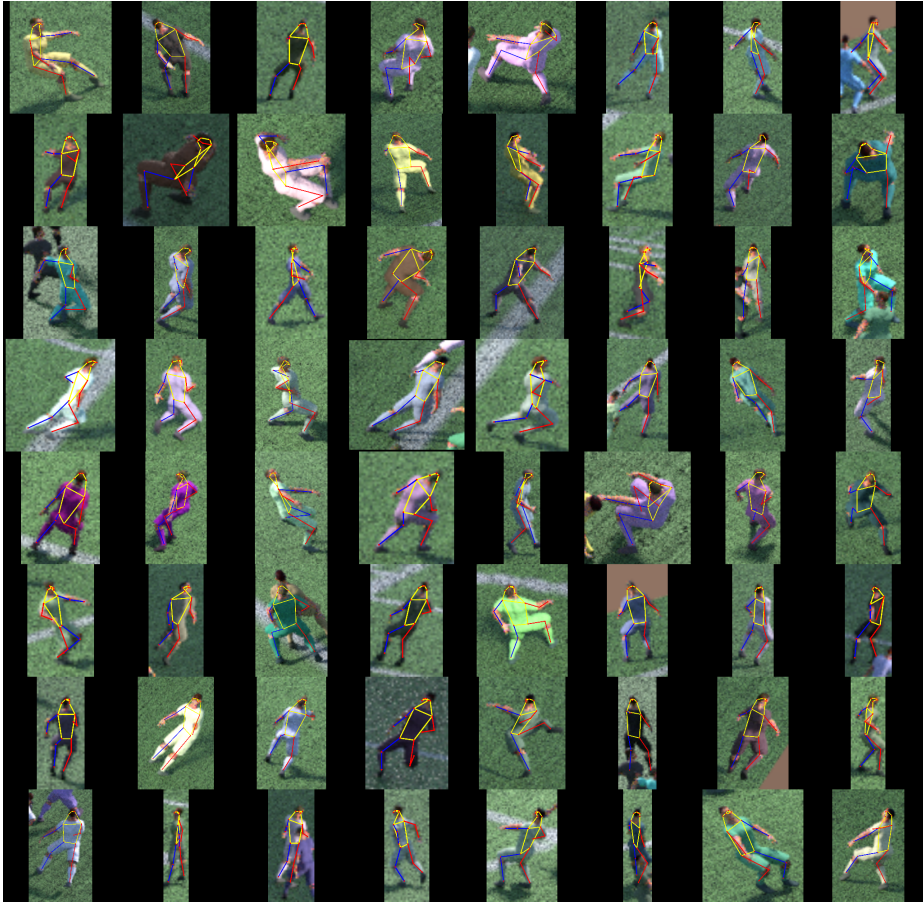


Figure 1: Example synthetic blobs with drawn ground truth skeletons.

blobs from 2.1k captured sequences for the second distribution. Example CycleGAN-ed training images are shown in figure 2. The model is trained with batch size of 1.

4.3 Benchmark

We annotated 400 real blobs with full skeletons. First, all testing blobs are fit into 256x192 rectangles. For each blob we measure OKS (Object Keypoint Similarity) given by:

$$OKS = \frac{\sum_i^n \exp(-d_i^2)}{n}, \quad (1)$$

$$d_i = \sqrt{\left(\frac{x_i - \bar{x}_i}{32}\right)^2 + \left(\frac{y_i - \bar{y}_i}{32}\right)^2} \quad (2)$$

where x_i, y_i are ground truth keypoint coordinates in pixel space, \bar{x}_i, \bar{y}_i predicted coordinates in pixel space, and n is number of keypoints (17 in our case). In COCO human keypoint annotations, head has 5 defined keypoints (eyes, ears and nose). For our benchmark we merge it into one keypoint (averaging coordinates of these 5 keypoints, both in prediction and ground truth), because such details are not visible at all on our images. We measure mean OKS over all test blobs. We consider our testset sufficiently large for measurement (see figure 3).

Additionally we create raw synthetic benchmark on raw synthetic images to better illustrate capabilities of our trained models. This testset consists of 10k blobs (it is not included CycleGAN training set).

4.4 Results

In experiments, we use CycleGAN-ed images made using checkpoint after 5k, and 18k iterations of training CycleGAN. As source raw synthetic images for transformation we use the same 100k samples that were used during the training. Example detections for selected experiments are shown in figure 5. Mean OKS value of these experiments over training epoch is shown in figure 4. COCO Minival and raw synthetic benchmark scores for selected checkpoints are shown in table 2.

Usually, best score is achieved in early epochs of training. It is possible, because our training and evaluation images are from different distributions.

Clearly, long training with only raw synthetic data causes the model to exploit synthetic artifacts and assumptions based on imperfect artificial heuristic pose distribution. In this case the model learns artificial distribution very easily – achieves almost perfect results on this distribution after only 5 epochs (see table 2).

Training with CycleGAN-ed data achieves high results (close to training on COCO) in early training epochs, therefore our augmentation method of raw synthetic dataset makes its distribution more similar to the real images distribution. Training with CycleGAN-ed images from early checkpoint (5k iterations) shows somewhat averaged results between training with raw synthetic and later checkpoint (18k iterations) CycleGAN-ed dataset. In general, models trained with artificially created (for our domain) data doesn't work at all on COCO benchmark – their score is close to 0. Moreover, mixed dataset training decreases the score.

Despite our experiments are not exhaustive (e.g. in this paper we try only one option in terms of selecting various model parameters), mixed dataset training achieves high scoring checkpoints on our benchmark faster than training on

training set	training epoch	mean OKS (our benchmark)	mean OKS (our synthetic benchmark)	COCO (Minival) AP @0.5:0.95
coco	163 (best)	0.725	0.824	0.691
coco	400	0.700	0.836	0.700
coco + CycleGAN-ed	46 (best)	0.725	0.926	0.595
CycleGAN-ed	23 (best)	0.691	0.923	0.009
raw synthetic	5 (best)	0.572	0.966	0.006
raw synthetic	100	0.303	0.977	0.004

Table 2: Summary of selected checkpoints scores. Epoch marked with ”(best)”, is the one after which the model achieves best score (among the other checkpoints from this experiment) on our real images benchmark.

COCO only. We suspect that detailed experiments on various stages of our experiments may achieve even higher results.

5 Conclusions

In this paper, we focus on football players pose estimation on very low-resolution images, received from the actual High Quality CCTV system located on lighting spots in the corners of the football pitch. To omit the need for the manual annotation of many thousands of training examples we create simple python script for rendering synthetic images. In order to give more realism to our raw synthetic images we used vanilla CycleGAN. Conducted experiments proved, that training neural networks for pose estimation without manually annotated data can (in some cases) achieve as good results as training with large, manually annotated, generic datasets (like COCO keypoints). With more exhaustive experiments, it may be possible to achieve even better results by changing synthetic dataset generation method, various hyperparameters, and architectures of both image translation and pose estimation models.

Precise annotation of a large training set requires many hours of human labor, while script for rendering synthetic dataset for a specific task using heuristics, can be created by one person and much faster.

ACKNOWLEDGEMENTS

This work was co-financed by the European Union within the European Regional Development Fund.

References

1. ANDRILUKA, M., PISHCHULIN, L., GEHLER, P., AND SCHIELE, B. 2D human pose estimation: New benchmark and state of the art analysis. In Proc. IEEE Conf. Computer Vision and Pattern Recognition (June 2014), pp. 3686–3693.
2. ARJOVSKY, M., CHINTALA, S., AND BOTTOU, L. Wasserstein gan.
3. BULAT, A., AND TZIMIROPOULOS, G. Super-fan: Integrated facial landmark localization and super-resolution of real-world low resolution faces in arbitrary poses with gans.
4. CAO, Z., SIMON, T., WEI, S. E., AND SHEIKH, Y. Realtime multi-person 2D pose estimation using part affinity fields. In Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR) (July 2017), pp. 1302–1310.
5. CHEN, Y., SHEN, C., WEI, X. S., LIU, L., AND YANG, J. Adversarial posenet: A structure-aware convolutional network for human pose estimation. In Proc. IEEE Int. Conf. Computer Vision (ICCV) (Oct. 2017), pp. 1221–1230.
6. CHEN, Y., WANG, Z., PENG, Y., ZHANG, Z., YU, G., AND SUN, J. Cascaded pyramid network for multi-person pose estimation.
7. DENG, J., DONG, W., SOCHER, R., LI, L., LI, K., AND FEI-FEI, L. Imagenet: A large-scale hierarchical image database. In Proc. IEEE Conf. Computer Vision and Pattern Recognition (June 2009), pp. 248–255.
8. DONG, C., LOY, C. C., HE, K., AND TANG, X. Image super-resolution using deep convolutional networks. IEEE Transactions on Pattern Analysis and Machine Intelligence 38, 2 (Feb. 2016), 295–307.
9. DWIBEDI, D., MISRA, I., AND HEBERT, M. Cut, paste and learn: Surprisingly easy synthesis for instance detection. In Proc. IEEE Int. Conf. Computer Vision (ICCV) (Oct. 2017), pp. 1310–1319.
10. GEORGAKIS, G., REZA, M. A., MOUSAVIAN, A., LE, P., AND KOŠECKÁ, J. Multiview rgb-d dataset for object instance detection. In Proc. Fourth Int. Conf. 3D Vision (3DV) (Oct. 2016), pp. 426–434.
11. GINES HIDALGO, ZHE CAO, T. S. S.-E. W. H. J. Y. S. Openpose. <https://github.com/CMU-Perceptual-Computing-Lab/openpose>, June 2017.
12. GOODFELLOW, I. J., POUGET-ABADIE, J., MIRZA, M., XU, B., WARDE-FARLEY, D., OZAIR, S., COURVILLE, A., AND BENGIO, Y. Generative adversarial networks.
13. GULRAJANI, I., AHMED, F., ARJOVSKY, M., DUMOULIN, V., AND COURVILLE, A. Improved training of wasserstein gans.
14. GUPTA, A., VEDALDI, A., AND ZISSERMAN, A. Synthetic data for text localisation in natural images. In Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR) (June 2016), pp. 2315–2324.
15. HE, K., GKIOXARI, G., DOLLÁR, P., AND GIRSHICK, R. Mask r-CNN. In Proc. IEEE Int. Conf. Computer Vision (ICCV) (Oct. 2017), pp. 2980–2988.
16. HE, K., ZHANG, X., REN, S., AND SUN, J. Deep residual learning for image recognition.
17. KARRAS, T., AILA, T., LAINE, S., AND LEHTINEN, J. Progressive growing of gans for improved quality, stability, and variation.
18. KARRAS, T., LAINE, S., AND AILA, T. A style-based generator architecture for generative adversarial networks.
19. KE, L., CHANG, M.-C., QI, H., AND LYU, S. Multi-scale structure-aware network for human pose estimation.
20. LEDIG, C., THEIS, L., HUSZÁR, F., CABALLERO, J., CUNNINGHAM, A., ACOSTA, A., AITKEN, A., TEJANI, A., TOTZ, J., WANG, Z., AND SHI, W. Photo-realistic single image super-resolution using a generative adversarial network. In Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR) (July 2017), pp. 105–114.

21. LIN, T.-Y., MAIRE, M., BELONGIE, S., BOURDEV, L., GIRSHICK, R., HAYS, J., PERONA, P., RAMANAN, D., ZITNICK, C. L., AND DOLLÁR, P. Microsoft coco: Common objects in context.
22. MAO, X., LI, Q., XIE, H., LAU, R. Y. K., WANG, Z., AND SMOLLEY, S. P. Least squares generative adversarial networks. In Proc. IEEE Int. Conf. Computer Vision (ICCV) (Oct. 2017), pp. 2813–2821.
23. MCCORMAC, J., HANDA, A., LEUTENEGGER, S., AND DAVISON, A. J. Scenenet rgb-d: 5m photorealistic images of synthetic indoor trajectories with ground truth.
24. MUELLER, F., BERNARD, F., SOTNYCHENKO, O., MEHTA, D., SRIDHAR, S., CASAS, D., AND THEOBALT, C. Gnerated hands for real-time 3d hand tracking from monocular rgb.
25. NEWELL, A., HUANG, Z., AND DENG, J. Associative embedding: End-to-end learning for joint detection and grouping.
26. NEWELL, A., YANG, K., AND DENG, J. Stacked hourglass networks for human pose estimation.
27. RADFORD, A., METZ, L., AND CHINTALA, S. Unsupervised representation learning with deep convolutional generative adversarial networks.
28. REN, S., HE, K., GIRSHICK, R., AND SUN, J. Faster r-CNN: Towards real-time object detection with region proposal networks. IEEE Transactions on Pattern Analysis and Machine Intelligence 39, 6 (June 2017), 1137–1149.
29. SIMON, T., JOO, H., MATTHEWS, I., AND SHEIKH, Y. Hand keypoint detection in single images using multiview bootstrapping. In Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR) (July 2017), pp. 4645–4653.
30. SYPETKOWSKI, M., KURZEJAMSKI, G., AND SARWAS, G. Football players pose estimation. In Image Processing and Communications Challenges 10 (Cham, 2019), M. Choraś and R. S. Choraś, Eds., Springer International Publishing, pp. 63–70.
31. ZHANG, H., XU, T., LI, H., ZHANG, S., WANG, X., HUANG, X., AND METAXAS, D. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In Proc. IEEE Int. Conf. Computer Vision (ICCV) (Oct. 2017), pp. 5908–5916.
32. ZHANG, H., XU, T., LI, H., ZHANG, S., WANG, X., HUANG, X., AND METAXAS, D. N. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. IEEE Transactions on Pattern Analysis and Machine Intelligence (2018), 1.
33. ZHU, J.-Y., PARK, T., ISOLA, P., AND EFROS, A. A. Unpaired image-to-image translation using cycle-consistent adversarial networks.



Figure 2: Example images from CycleGAN with drawn ground truth skeletons. The skeletons are drawn with thin lines, so that visual artifacts are visible. First column shows original synthetic images, others correspond to training iterations – from left: 5k, 10k, 15k, 50k, 100k, 180k. Transformed images are not perfectly consistent geometrically, but characteristic distortions occurring in the real conditions are performed in appropriate parts of the image – it enables the network comprehensible inference on the real images.

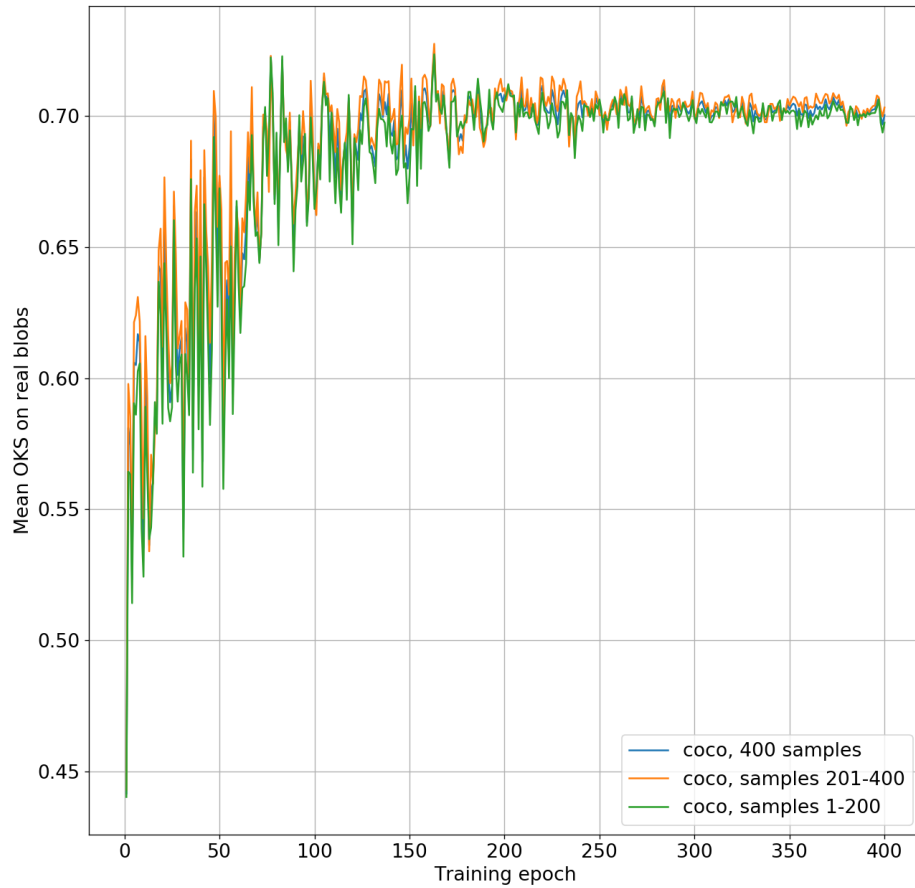


Figure 3: The same model trained on COCO gives similar plot shape when evaluated on 2 separated parts of our small testset. Basing on this observation, we assume that our benchmark allows meaningful comparison between different models.

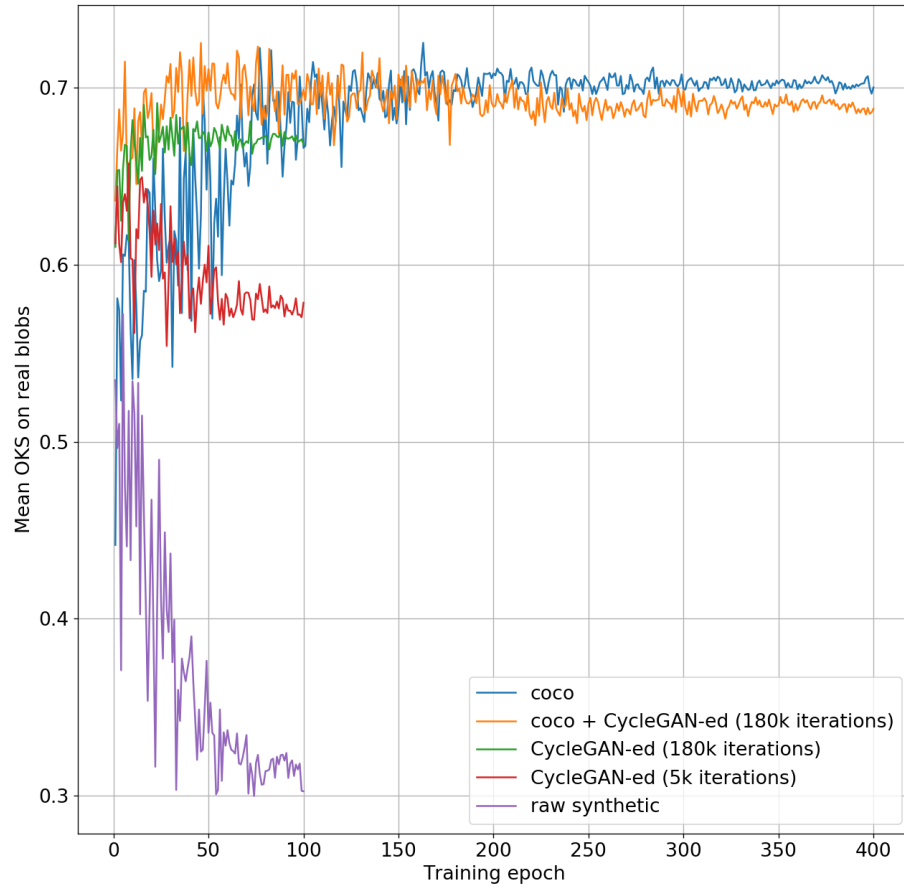


Figure 4: Mean OKS over training epoch. Initial learning rate is 0.0005, and it is decreased by 50% after each 60 epochs in experiments where COCO dataset is used, and after 15 epochs in the other experiments. We use batch size of 32. Each epoch is 60k randomly sampled images without returning (with original CPN data augmentation).

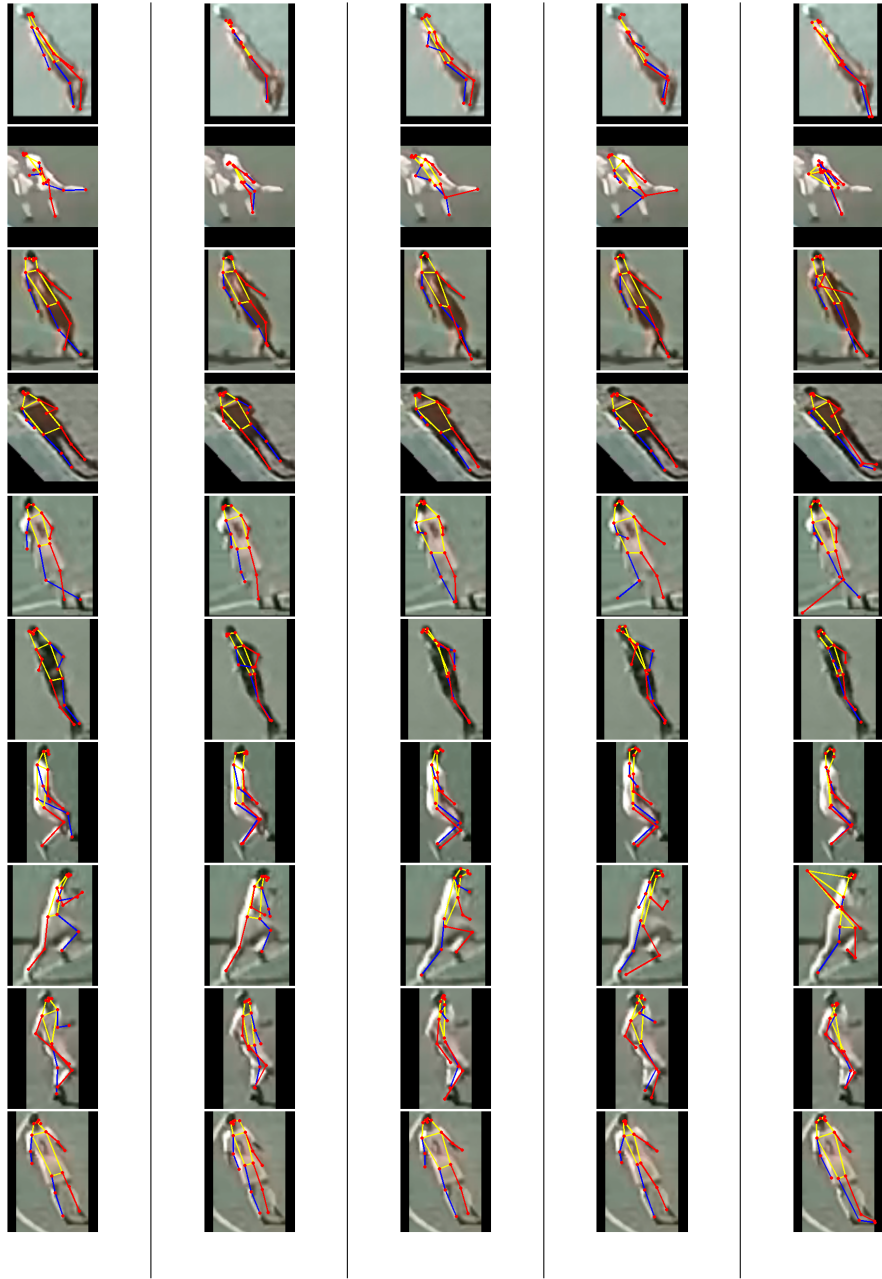


Figure 5: Example pose estimations for best checkpoints. Columns from left: ground truth, coco, coco + CycleGAN-ed, CycleGan-ed only , raw synthetic.