

## REAL-TIME POSE ESTIMATION VIA STATE-OF-THE-ART COMPUTER VISION MODELS

Nghia Duong-Trung, Benjamin Paaßen, Milos Kravcik German Research Center for Artificial Intelligence (DFKI GmbH) consortium MILKI-PSY meeting, Aachen, 12.12.2022

SPONSORED BY THE



Federal Ministry of Education and Research

#### Content



- 1. Introduction to Pose Estimation
- 2. Introduction to YOLO Family Models
- 3. YOLO-Pose
- 4. YOLO-Pose: Demo
- 5. Remarks and Discussion

**Introduction to Pose Estimation** 



#### 1. Introduction to Pose Estimation

- 2. Introduction to YOLO Family Models
- 3. YOLO-Pose
- 4. YOLO-Pose: Demo
- 5. Remarks and Discussion

## What is pose estimation?



- Pose estimation and tracking is a computer vision task that includes detecting, associating, and tracking semantic key points.
  - ► human, animal, object.
- Offline (image, replayed video) or online (real-time video).
- SOTA methods are typically based on designing the CNN architecture tailored particularly for human pose inference.



Figure: Examples of pose predictions on sports, professional and casual photos. Image source: https://arxiv.org/abs/2103.02440.

- In tradition object detection, people are perceived as a bounding box
- Pose detection and tracking -> semantically understanding of human body
- ► High-performing real-time pose detection and tracking will drive some of the biggest trends in computer vision
  - sport, autonomous driving, animal research



Figure: Example of animal pose predictions. Image source: https://github.com/DeepLabCut.

## What is human pose estimation?



- Predict the locations of body parts and joints in images or videos: 2D and 3D
- ► Human pose modeling:
  - Kinematic: skeleton-based model: set of joint positions and limb orientations
  - Planar: contour-based model: the appearance and shape of a human body
  - ➤ Volumetric model: 3D estimation for recovering 3D human mesh

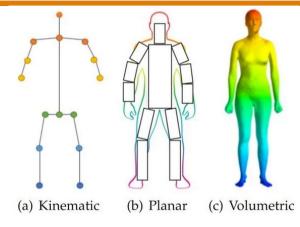


Figure: Human pose modeling. Image source: https://arxiv.org/abs/2012.13392.

### How does pose estimation work?



- Estimate spatial positions of a body in an image or video
  - ► multi-pose, single pose
- ► Human pose estimation on the popular MS COCO dataset<sup>a</sup> can detect 17 different keypoints. Each keypoint is annotated with a tuple (x, y, v), where x and y mark the coordinates, and v indicates if the keypoint is visible.



Figure: Human pose modeling. Image source: https://viso.ai/deep-learning/openpose/.

ahttps://cocodataset.org/#download

- 1. Deep Pose, CVPR 2014
- 2. PoseNet, ICCV 2015
- 3. DeepCut, CVPR 2016
- 4. OpenPose, CVPR 2017
- 5. Regional Multi-Person Pose Estimation (AlphaPose), CVPR 2017
- 6. Dense Pose, CVPR 2018
- 7. Deep High-Resolution Net (HRNet), CVPR 2019
- 8. BlazePose (MediaPipe), CVPR 2020

## Use cases and applications



- ► Human-computer interaction, action recognition, motion analysis, augmented reality, sports and fitness, and robotics.
  - ▶ Detect gestures
  - Support the analysis of football, basketball, and sports
  - ► Analyze dance techniques
  - Analyze posture for body works and finesses
  - ► Play games via virtual tutors
  - ► Human tracking for consoles
- ► Future trends: transfer learning and edge AI

## Introduction to YOLO Family Models



- 1. Introduction to Pose Estimation
- 2. Introduction to YOLO Family Models
- 3. YOLO-Pose
- 4. YOLO-Pose: Demo
- 5. Remarks and Discussion

## You only look once (YOLO)



- ▶ J Redmon et al. You only look once: Unified, real-time object detection¹. CVPR 2016. 29950 citations
  - ► YOLOv1 to YOLOv3
- ► YOLOv4, introduced by A. Bochkovskiy et al. in 2020
- ► Scaled-YOLOv4, introduced by CY Wang et al. in 2021
- ► YOLOR, YOLOX, NanoDet-Plus
- ► PP-YOLOE, an industrial object detector PaddlePaddle, 2019
- ► YOLOv5 model v6.1, Ultralytics, 2022
- ► YOLOv7, CVPR 07.2022, https://arxiv.org/abs/2207.02696
  - ► There are actually two YOLOv7, the other developed on Detectron2

https://www.cv-foundation.org/openaccess/content\_cvpr\_2016/papers/Redmon\_You\_ Only\_Look\_CVPR\_2016\_paper.pdf

## YOLOv7 vs other object detectors



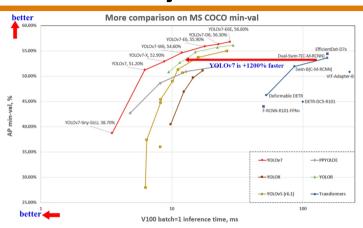


Figure: Image source: https://arxiv.org/abs/2207.02696, https://github.com/WongKinYiu/yolov7.

# YOLO-Pose



- 1. Introduction to Pose Estimation
- 2. Introduction to YOLO Family Models
- 3. YOLO-Pose
- 4. YOLO-Pose: Demo
- 5. Remarks and Discussion

### **YOLO-Pose paper**



- ▶ Maji, D., Nagori, S., Mathew, M., & Poddar, D. (2022). YOLO-Pose: Enhancing YOLO for Multi Person Pose Estimation Using Object Keypoint Similarity Loss. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 2637-2646) <sup>2</sup>.
  - ► https://github.com/ultralytics/yolov5/
  - ► https://github.com/TexasInstruments/edgeai-yolov5
  - ► https://github.com/TexasInstruments/edgeai-yolox
  - ► https://github.com/WongKinYiu/yolov7/blob/main/tools/keypoint.ipynb
  - ► PyTorch, ONNX, CoreML, TFLite
  - tested on image, video, webcam

<sup>&</sup>lt;sup>2</sup>https://openaccess.thecvf.com/content/CVPR2022W/ECV/papers/Maji\_Y0L0-Pose\_ Enhancing\_Y0L0\_for\_Multi\_Person\_Pose\_Estimation\_Using\_Object\_CVPRW\_2022\_paper.pdf

- Multi-person 2D pose estimation can be categorized into top-down and bottom-up approaches
- ► Top-down: or two-stage approaches first perform human detection using a heavy person detector followed by estimating 2D pose for each detected person
  - ► Computational complexity increases linearly with the number of persons
- ► Bottom-up: find out identity-free keypoints of all the persons in an image in a single shot followed by grouping them into individual person instances
  - ► Rely on heatmaps to detect all the keypoints; require a complex post-processing to group keypoints

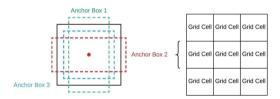
#### **YOLO-Pose idea**



- ➤ YOLO-Pose is a single shot approach like other bottom-up approaches but it does not use heatmaps.
- ► YOLO-Pose associates all keypoints of a person with anchors.
- Keypoints associated with an anchor are already grouped
- Independent of the number of persons in an image
- End-to-end training without independent post-processing
  - Extend the idea of IoU loss from box detection to keypoints
  - ► Object keypoint similarity (OKS) is used for evaluation and a loss for training

## What is anchor (box) in object detection Partification Partification

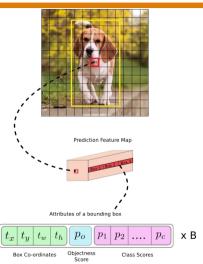
- Anchor box is a prior box that could have different pre-defined aspect ratios determined before training by running K-means on the entire dataset or a hyperparameter
- The output of convolution is a matrix as a "grid", then we assign anchor (boxes) to the grid cells, and they share the same centroid
- One we define those anchors, we can determine how much the ground truth box overlaps with the anchor boxes and pick the one with the best IoU and couple them together



## **Anchor example**



► Each cell is assigned 3 anchors containing some set of properties (x,y,w,h,object score, classes)



#### **Anchor in YOLO-Pose**



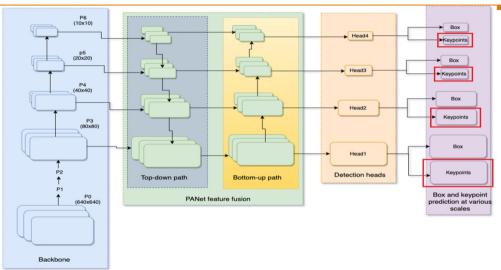
- ▶ It boils down to a single class person detection problem with each person having 17 associated keypoints, and each keypoint is again identified with a location and confidence: {x, y, conf}
- ► The overall prediction vector is defined as:

$$P_{v} = \{C_{x}, C_{y}, W, H, box_{conf}, class_{conf}, K_{x}^{1}, K_{y}^{1}, K_{conf}^{1}, \dots, K_{x}^{n}, K_{y}^{n}, K_{conf}^{n}\}$$
(1)

where a box head predicts 6 elements and a keypoint head predicts 51 elements.

### **YOLO-Pose architecture**





For a ground truth bounding box that is matched with  $k^{th}$  anchor at location (i, j) and scale s, loss is defined as:

$$\mathcal{L}_{box}(s, i, j, k) = (1 - CloU(Box_{gt}^{s, i, j, k}, Box_{pred}^{s, i, j, k}))$$
 (2)

where CloU is defined in another paper 3

https://ojs.aaai.org/index.php/AAAI/article/view/6999/6853

- Object keypoint similarity (OKS) is the most popular metric for evaluating keypoints
- Corresponding to each bounding box, YOLO-pose model stores the entire pose information
  - ▶ If a ground truth bounding box is matched with  $k^{th}$  anchor at location (i, j) and scale s, the model predicts the keypoints with respect to the center of the anchor
  - ► OKS is computed for each keypoint separately and then summed to give the final OKS loss or keypoint IoU loss

$$\mathcal{L}_{kpts}(s, i, j, k) = 1 - \frac{\sum_{n=1}^{N_{kpts}} \exp(\frac{2s^2 k_n^2}{d_n^2}) \delta(v_n > 0)}{\sum_{n=1}^{N_{kpts}} \delta(v_n > 0)}$$
(3)

where

 $d_n$  = Euclidean distance between predicted and ground truth location for  $n^{th}$  keypoint

 $k_n$  = keypoint specific weights

s =scale of an object

 $\delta(v_n)$  = visibility flag for each keypoint

## Human pose loss function formulation

Corresponding to each keypoint, we learn a confidence parameter that shows whether a keypoint is present for that person or not

$$\mathcal{L}_{kpts\_conf}(s, i, j, k) = \sum_{n=1}^{N_{kpts}} BCE(\delta(v_n > 0), p_{kpts}^n)$$
(4)

where

 $p_{kpts}^n =$ predicted confidence for  $n^{th}$  keypoint

Loss at location (i, j) is valid for  $k^{th}$  anchor at scale s if a ground truth bounding box is matched against that anchor. Finally, the total loss is summed over all scales, anchors and locations:

$$\mathcal{L}_{total} = \sum_{s.i.i.k} (\lambda_{cls} \mathcal{L}_{cls} + \lambda_{box} \mathcal{L}_{box} + \lambda_{kpts} \mathcal{L}_{kpts} + \lambda_{kpts\_conf} \mathcal{L}_{kpts\_conf})$$
 (5)

where  $\lambda_{\textit{cls}}, \lambda_{\textit{box}}, \lambda_{\textit{kpts}}, \lambda_{\textit{kpts\_conf}}$  are hyperparameters chosen to balance between losses at different scales

## **YOLO-pose pretrained model**



- ► The YOLO-pose model was trained on COCO dataset<sup>4</sup>
  - ▶ 200,000 images with 250,000 person instances with 17 keypoints
- ▶ train2017 (57k images), val2017(5k images), test-dev2017(20k images)
- ► https://github.com/WongKinYiu/yolov7/tree/pose
- https://github.com/WongKinYiu/yolov7/releases/download/v0.1/ yolov7-w6-person.pt

<sup>4</sup>https://cocodataset.org/#home

## YOLO-Pose: Demo



- 1. Introduction to Pose Estimation
- 2. Introduction to YOLO Family Models
- 3. YOLO-Pose
- 4. YOLO-Pose: Demo
- 5. Remarks and Discussion

#### Demo



- ► https://github.com/MILKI-PSY/yolov7-pose-bicep
- ▶ image, video
  - Examples of teacher and students: https://github.com/MILKI-PSY/yolov7-pose-bicep/tree/main/input
  - ► After applying YOLO-pose, see the output files at https://github.com/MILKI-PSY/yolov7-pose-bicep
- Live demo via webcam

## Remarks and Discussion



- 1 Introduction to Pose Estimation
- 2. Introduction to YOLO Family Models
- 3. YOLO-Pose
- 4. YOLO-Pose: Demo
- 5. Remarks and Discussion

#### Remarks and discussion



- ► SOTA object detection and pose estimation model, all-in-one model
- ▶ Privacy and sensible information are preserved at the partners server
- Flexible to integrate more exercises estimation
- Flexible to input data
- ► Further recommendation



Thank you for listening