

基于强化学习的多人姿态检测算法优化

黄 铎 应 娜 蔡哲栋

(杭州电子科技大学通信工程学院 浙江 杭州 310018)

摘要 在多目标人体姿态检测算法过程中,人体的定位精度依然不够精确。针对该问题,采用速度与精度兼顾的SSD算法作为目标检测器获得人体的初步包围框,定义该包围框为智能体,引入强化学习。采用马尔科夫决策过程以及Q网络组成的目标精细模型对智能体训练其九种动作,分别为左上角与右下角两个点的四方向进行迭代调整以及终止策略,使得包围框达到更贴近人体的效果。结合先进的Stacked hourglass算法作为姿态检测器,对调整后的包围框进行姿态预测。该算法的引入使得多目标人体检测算法在MPII数据集上的精度提升了1.6 mAP,达到了73.7 mAP。

关键词 深度学习 强化学习 姿态检测 模式识别

中图分类号 TP3 文献标识码 A DOI:10.3969/j.issn.1000-386x.2019.04.029

OPTIMIZATION OF ESTIMATION ALGORITHM FOR THE MULTI-PERSON POSE BASED ON REINFORCEMENT LEARNING

Huang Duo Ying Na Cai Zhedong

(School of Communication Engineering, Hangzhou Dianzi University, Hangzhou 310018, Zhejiang, China)

Abstract The accuracy of locating people is imprecise in multi-objective human pose estimation. Aiming at this problem, we applied SSD with both speed and accuracy as an object detector to obtain the preliminary bounding box of human body. The algorithm defined the bounding box as an agent and introduced the reinforcement learning. Markov decision-making process and the objective fine model consisting of Q net were applied to train the agent for nine actions. The training actions included iteration adjustments and termination policy directing to two points from the top left corner and points from the bottom right corner respectively, which made the bounding box closer to the human body. The advanced Stacked hourglass algorithm was implemented as pose detector to predict the adjusted bounding box. The introduction of the referred algorithm improves the accuracy of multi-objective human pose estimation on MPII dataset by 1.6 mAP and reaches 73.7 mAP.

Keywords Deep learning Reinforcement learning Pose estimation Pattern recognition

0 引言

人体姿态检测问题是基本的计算机视觉问题,目前已经有了很多算法对这一问题进行了研究,并且取得了不错的成效。

在单人姿态检测问题中,由于人体占据了图片的大部分像素,所以对于检测器而言不会有许多干扰信

息。这在姿态检测问题中是比较简单的一个部分。一些传统的方法采用了图像结构模型。例如,在人体姿态检测问题中非常有效的树模型^[8]和随机森林模型^[9]。在文献[10]中广泛研究了一些基于图像的模型,例如随机场模型^[11]和图像依赖模型^[12]。

最近,深度学习成为了一种在模式识别方面非常有效的方法,利用深度学习的人体姿态检测同样取得了很好的成绩。一些有代表性的方法例如深度姿

态^[13],基于DNN的模型^[14]和基于CNN的模型^[1]。

以上所说的这些方法在人体信息足够完整的情况下可以有很好的表现,但是在多人姿态检测问题中常常不会有如此精确的人体信息。

而与单人姿态检测问题相比,多人姿态检测则更具有挑战性(意为在单张图片中有多个个体)。针对多人姿态检测,目前的算法大致分为两种。其一是使用两步检测框架,其二是采用肢体框架。两步检测框架中的第一步,是目标检测算法通过包围框的方式定位出人体,第二步是分别将每个包围框做姿态检测。肢体框架则是独立地检测人体的各个部分,再把各个部分拼接起来,以形成多个个体姿态。两种方法各有利弊,其中两步检测框架的检测精度高度依赖目标检测算法的精度,包围框对人体的契合程度会影响姿态检测算法的精度;而肢体框架的图片中若多个个体的距离相近,甚至互相遮挡,肢体框架的拼接工作就会变得困难。

1 相关工作

(1) 肢体检测框架 肢体检测框架有许多代表性的方法。Chen等^[4]提出了一种通过图像模型来解析被遮挡的人体的方法,该方法可以将人体的各个部分灵活组合。Gkioxari等^[3]使用k-poselets来全局检测人和预测人体姿势,通过加权平均来预测最终的姿态。Pishchulin等^[5]提出了DeepCut模型,该模型先检测身体的每个部分,然后通过整体线性程序来标记并且组装这些部分,最后达到预测人体姿态的目的。Insafutdinov等^[15]提出了一个更强的基于ResNet的肢体检测器和一个更好的增量优化策略^[6]。

虽然基于肢体检测方法有良好的表现,但是由于该方法只专注于小区域的检测,缺少全局信息,所以该方法依旧存在精度上的缺陷。

(2) 两步框架 在两步框架的研究中,几乎都由人体检测器和姿态检测器两部分组成。Pishchulin等^[2]采用了传统图像结构模型进行姿态预测。Insafutdinov等^[6]则采用了相似的两步框架,由人体检测器Faster R-CNN^[16]和姿态预测器DeeperCut构成。他们的方法在MPII数据集上最高只能达到51.0 mAP的精度。由于计算机技术不断发展,硬件水平快速上升的情况下,物体检测器和单人姿态检测器都达到了一个相当好的水平,而两步框架得益于两者的先进性,在多人姿态检测问题中也取得了非常优异的表现。

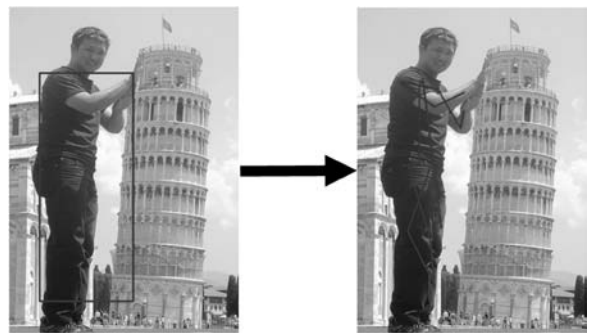
但如上文所说,两步框架中依旧存在着人体监测框架不精确的问题,本文致力于解决这个问题,在两步

框架中加入了强化学习的顺序决策过程,用于调整包围框以提升算法精度。

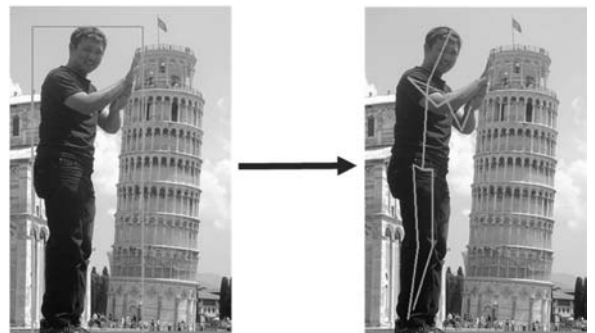
强化学习的概念从博弈论衍生而来,通过不断的试错学习,从而找到一种最佳的解决问题的方法。强化学习在诸多领域都是非常强有力的工具,如在雅达利2600游戏机和AlphaGo等项目都有出色应用。

在目标检测领域,强化学习也发挥出很强大的作用。Juan等^[17]设计了一些包围框的动作,通过对包围框的调整不断接近目标,最终定位目标。Míriam等^[18]将图片分为五个区域,挑选出目标存在的区域后重复这一步骤,不断聚焦目标,实现定位功能。

本文采用的算法基于两步检测框架。如同上文所述,包围框对人体的贴合程度会影响人体姿态检测的精度,图1通过SSD300^[7]和hourglass算法^[1]的结合说明了一个问题。其中,图1(a)表示在目标检测步骤中被判定为正确的情况(比如交并比大于0.6),在人体姿态检测时依然会出现很严重的错误。这是因为通过目标检测算法得到的包围框虽然检测到了人体,但是这些信息却不够精确,导致人体姿态检测算法的精度受到了负面影响。



(a) 包围框信息不精确导致的姿态不准问题



(b) 精确的包围框在姿态检测上可以有很好的效果

图1 不同包围框精度情况下的检测效果

为了解决人体检测的包围框信息不精确问题,本文算法在两步检测框架之间引入了强化学习,形成基于强化学习的两步检测框架,调整包围框使其更贴近人体,以提升精度。算法定义包围框为智能体,建立了一个顺序决策过程。智能体通过卷积网络获取图片信息进行包围框的调整,下一次迭代中智能体会根据调

整后的包围框信息决定该次迭代中的动作。算法不断重复这一过程,最终得到一个最适合人体姿态检测的包围框,以此来提升精度。其中,算法采用了马尔科夫决策过程,设计了8种变形动作以及一种终止动作,如图2所示,8种动作分别为左上角点和右下角点的上下左右平移。这8种变形动作涵盖了包围框所有的动作,以满足贴近人体的需求。算法在MPII数据集上做了训练和验证,结果达到了74 mAP。

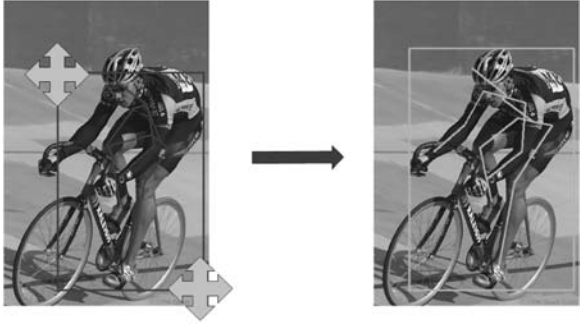


图2 通过调整包围框实现精度的提升

2 算法设计

本文的算法如图3所示,首先通过改进的目标检测器SSD回归出初步的人体包围框,再通过深度强化学习网络判断包围框是否准确,并对定位的包围框设计了8种变形操作以及一种终止动作,即采用马尔科夫策略对其左上角与右下角坐标进行迭代调整,调整动作集合为两个点的上下左右平移,从而提升精度。最后,由人体姿态检测器对调整后的包围框进行姿态检测。

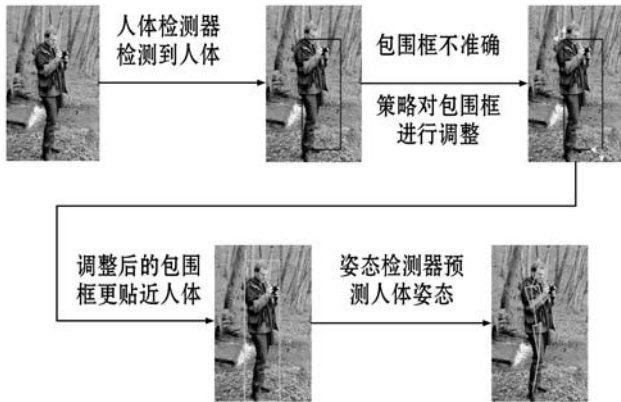


图3 本文算法的流程图

2.1 人体目标检测器

本文采用了兼顾精度与速度的目标检测算法SSD^[7] (Single Shot MultiBox Detector)。该算法采用VGG-16^[21]作为前置网络,删除了末尾的全连接层,并在卷积层后添加置信层和回归层。SSD在网络中选择6种不同尺度的特征图,每张特征图的每个像素点都

预设了一定数量的预设框。特征图经过回归层后会得到与预设框相同数量的包围框,并且与预设框一一对应。预设框与真实包围框(label)进行匹配,IOU大于0.5的定义为正样本,其余的为负样本。最后将与正样本对应的包围框与真实值进行loss计算。常用的SSD算法输入尺度为300,但是小目标检测的精度不高。在本文中SSD算法的输入尺度为512,置信层与回归层都为7层,相比于300尺度的SSD算法,更大的输入尺度在小目标检测的精度上有很大的改善。

本文通过SSD对图像进行特征提取,并回归得到目标的包围框坐标。

$$bbox = Conv_{SSD}(img) \quad (1)$$

式中:img是输入图片,Conv_{SSD}(·)是SSD算法,bbox是回归得到的包围框坐标。

2.2 基于强化学习的目标精细模型

针对目标检测器信息不够精确的问题,算法定义目标检测器所回归出来的包围框为智能体,该智能体会与环境交互获取信息,建立马尔科夫决策过程。在每一次迭代中,智能体需要获取信息来决定一次变形动作,在下一次的迭代中,智能体会根据上一次变性之后的信息来决定再下一次迭代的变形动作,直到确定目标最优或者达到限制的迭代次数为止。

2.2.1 马尔科夫决策过程参数

状态 智能体迭代至当前情况下的信息,状态是智能体决定动作的依据。

动作 动作来自于动作空间,智能体根据当前的状态和之前迭代的历史奖励信息,选择能达到最大化期望奖励的动作。

奖励 每次动作执行后,算法计算出该状态下执行该动作的奖励。奖励代表着智能体是否正朝着最终目标进行动作。

2.2.2 Q学习

智能体奖励的来源,动作A和状态st,受到函数Q(st,A)控制,该函数可以通过Q学习函数进行估计。智能体会通过函数Q(st,A)选择可以获得奖励的动作。Q学习函数使用Bellman方程(式(2))不断迭代更新动作选择策略,其中st和A是当前的相对应的动作和状态,R是当前的奖励,max_{a'}Q(st',A')表示未来的奖励,γ表示折扣因子。本文采用强化学习对Q网络进行训练,以使其能近似于Q函数^[19]。

$$Q(st,A) = R + \gamma \max_{a'} Q(st',A') \quad (2)$$

2.2.3 算法模型

当前智能体的信息由通过不同尺度的卷积层对bbox抽取的特征所组成,如式(3)所示。算法的状态st是由当前智能体的信息ft和历史动作hv所组成,如式

(4)所示,其中函数 $Cat(\cdot)$ 用于将 ft 和 hw 进行拼接。历史动作是一个向量,向量中包含了前4次迭代中智能体发生形变而所选择的动作,并入当前智能体的信息将有助于在训练的过程中稳定调整轨迹。算法将前4次迭代中的动作组成一个向量编入状态中,每次迭代会有9种不同的动作提供选择,所以一个历史动作向量是36维的。这种类型的向量也被文献[17]所采用。

$$ft = Conv(bbox) \quad (3)$$

$$st = Cat(ft, hw) \quad (4)$$

智能体通过由卷积神经网络提取到的特征进行决策,选择当前状态下所应该选择的动作,如式(5)所示,式中 $QNet$ 表示深度 Q 网络。算法设计了两种类型的动作:一是调整动作 $a(\cdot)$,该类型的动作会调整包围框的形状(式(6));二是终止动作 t ,该类型的动作一旦被选择,调整过程即终止。其中的调整动作数量有8种,分别是包围框左上角坐标的四个方向平移,和包围框右下角坐标的四个方向平移。这样设计的理由是这8种动作涵盖包围框的所有动作可能,相比于一般的包围框缩放和平移的规则动作,这样设计可以使得包围框做出不规则的动作,更有利于使包围框贴近人体。在迭代过程中智能体会不断根据当前的状态选择动作,每次调整包围框后会获得新的状态,再选择新的动作,直到选择为终止动作为止。

$$a = QNet(st) \quad (5)$$

$$bbox' = a(bbox) \quad (6)$$

算法选择由文献[17]中所提出的奖励公式作为该模型的奖励公式。调整动作的和终止动作奖励公式分别由式(9)和式(10)给出。因为智能体所选择的动作会产生新的包围框,人体姿态检测器 $Conv_{HAD}(\cdot)$ 会根据新的包围框产生新的精度 acc_1 (式(7));算法定义不加入强化学习的两步框架的精度 acc_0 作为真实值(式(8))。对于当前的状态 st ,智能体选择的调整动作得到新状态 st' ,产生的新精度 acc_1 ,如果大于真实值 acc_0 ,则会获得一个奖励(1),反之则会获得一个惩罚(-1)。对于终止动作而言,终止时最终的新精度 acc_1 若大于 acc_0 ,会获得一个比较大的奖励,反之会获得一个大惩罚。而对于真实值大于 τ 的目标,算法选择直接让智能体选择终止动作,获得奖励,这样做可以将训练时间缩短,只调整精度较差的目标。考虑原始精度,选择 $\tau = 0.5$ 。经过调参,在式(10)中 $\eta = 3$ 。

$$acc_1 = Conv_{HAD}(bbox') \quad (7)$$

$$acc_0 = Conv_{HAD}(bbox) \quad (8)$$

$$R_a(st, st') = sign(acc_1 - acc_0) \quad (9)$$

$$R_t(st, st') = \begin{cases} +\eta & acc_1 \geq \tau \\ -\eta & otherwise \end{cases} \quad (10)$$

在本文中,采用了比较先进的网络 $densenet^{[20]}$ 作为特征提取的方法。图4所示是深度强化学习网络的结构。

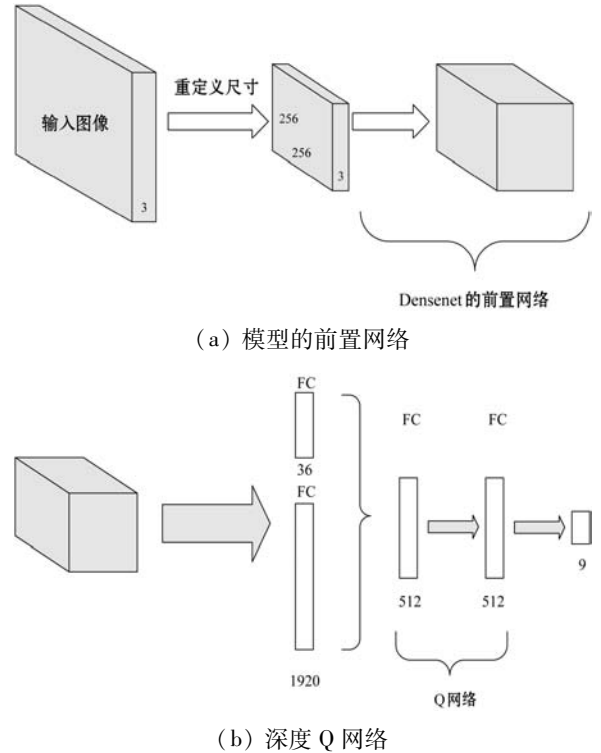


图4 深度强化学习网络的结构

在深度强化学习网络运行前,算法先把每个包围框裁剪好,将尺度调整为 256×256 ,图像不发生形变,多余的部分用0填充。算法将 $densenet$ 的特征提取部分取出作为前置网络,将该网络的输出重置成一个1920维的向量,与历史动作的36维向量组成了一个新的1956维的向量。深度 Q 网络是由两个全连接层组成,每层后都跟有一个 $ReLU$ 激活函数和一个 $dropout$ 层。最后,输出层对应于智能体可能选择的动作,在本文中数量为9个。

2.2.4 训练

• **探索-利用** 算法采用 ϵ -greedy 策略来训练深度 Q 网络。训练的时候会有 ϵ 的概率随机选择动作,即探索过程;有 $1 - \epsilon$ 的概率会通过深度 Q 网络选择一个动作,即利用过程。初始化 $\epsilon = 1$,之后随着训练次数增加 ϵ 会逐渐减小,直到 $\epsilon = 0.1$ 为止。事实上,随机选择动作会导致智能体学习终止动作更加困难,为了使得智能体更好地学习终止动作,当原始精度大于0.8时,算法强制智能体选择终止动作。这样可以使得训练的速度更快一些。由于算法一直都在做探索的过程,所以该方法不会被卡在局部最小值。

• **训练参数** 深度 Q 网络的初始化权重采用的是正态分布初始化, $densenet$ 则采用预训练权重。为了训练,本文采用了 Adam 优化器,学习率设置为 $1e - 6$ 以

避免梯度爆炸。本文在每个模型上训练 50 个 epoch。

• **经验回放** 上文提到, bellman 方程(式(2))从 (st, a, r, st') 的迁移中进行学习, 这也被称为经验。深度 Q 网络中的连续经验是非常相关且重要的, 处理不当可能导致学习速率低下和学习的不稳定, 这是 Q 学习中的一个传统性问题。算法采用了一个回放存储器使网络收敛, 以解决这个问题。回放存储器在每次迭代后会收集经验保存下来, 每次训练时会从保存在回放存储器中的经验中随机挑选一些进行训练。在本文中, 最大可保存的经验数为 1 000, 每次训练挑选的数量为 100。

• **折扣因子** 想要在长期的训练中表现出良好的效果, 只考虑当前的奖励是不够的, 所以算法将未来的奖励也加入训练。本文在式(3)中设置了折扣因子 $\gamma = 0.90$, 该数值可以很好地平衡当前的奖励和未来的奖励。

2.3 姿态检测器

算法采用了 stacked hourglass^[1] 模型作为人体姿态检测器。该模型以 Residual^[15] 模块为基础, 通过级联多个 Residual 模块和降采样层来获取不同尺度的信息, 随着阶数的增加, 级联的 Residual 模块的数量和特征图尺度的跨度也会增加, 最后将不同尺度的信息进行特征融合以预测人体姿态。

本文以 256×256 的图片作为姿态检测器的输入, 模型的阶数为四, 在五个不同的尺度上采集图像特征, 并且最后跳级融合。

3 实验

实验采用的数据集有两部分: 人体检测器所使用的数据集为 VOC0712; 单人姿态检测器所采用的数据集来自 MPII。目标精细模型所采用的数据集为原始的两步框架的结果。

3.1 数据集

MPII 数据集包含了 3 844 张训练图片和 1 758 张验证图片, 其中 1 758 张验证图片官网不提供标签数据。在训练图片中包含了 28 000 个左右的训练样本可供单人姿态检测器所训练。本次实验采用了 25 000 个目标作为训练样本, 3 000 个目标作为验证样本。

目标精细模型所采用的数据集测试由两步框架的精度结果构成的。预先将人体检测器所检测到的目标一一做精度计算, 保存其结果, 做成一个由图片名、包围框、原始精度所构成的数据集。训练时读取包围框以供调整, 读取原始精度作为真实值。

3.2 实验结果

在 MPII 数据集上验证了本文算法。完整的精度结果参见表 1。结果的演示图如图 5 所示。结果显示, 本文算法可以以很高的精度完成多人姿态检测任务。

表 1 在 MPII 数据集上的验证结果 mAP

算法	头	肩	肘	腕
Iqbal&Gall ^[22]	58.4	53.9	44.5	35.0
DeeperCut ^[6]	78.4	72.5	60.2	51.0
Levinkov et al. ^[23]	89.8	85.2	71.8	59.6
本文算法	67.0	77.5	76.6	76.9
算法	臀	膝	踝	整体
Iqbal&Gall ^[22]	42.2	36.7	31.1	43.1
DeeperCut ^[6]	57.2	52.0	45.4	59.5
Levinkov et al. ^[23]	71.1	63.0	53.5	70.6
本文算法	72.8	71.8	73.2	73.7



图 5 效果演示图一

本文算法虽然在精度上有改进, 但是却存在计算量巨大的缺点。一张图片的预测需要经过目标检测器、目标精细模型和姿态检测器三个模型。其中, 在目标精细模型中, 每迭代调整一次包围框都需要经历一次网络前传, 这导致了计算量的陡增。同时, 由于一张图片可能存在多个目标, 每个目标需要多次迭代, 所以时间开销也更大。在 GTX1080TI 显卡上进行预测, 目标检测器与姿态检测器构成的两步检测算法单个目标耗时 120 ms; 加入目标精细模型后的基于强化学习的两步检测算法单个目标耗时 220 ms。

同时, 对于人体信息丢失严重的情况, 检测效果不够理想, 如图 6 所示。



图 6 效果演示图二

4 结 语

本文提出了一种新的多人姿态检测算法,在其准确率方面优于以往的两步框架算法。算法加入了目标精细模型,该模型可通过对包围框的调整使得信息更加精确,使得两步框架的精度有了一定的提升,证明该方法是可靠的。但由于它需要不断对包围框进行迭代调整,对单个目标的处理速度偏高,在 GTX1080TI 显卡上达到了约 160 ms。在未来的工作中,可以期待强化学习在更多的领域参与到其中去。

参 考 文 献

- [1] Newell A, Yang K, Deng J. Stacked Hourglass Networks for Human Pose Estimation[C]//European Conference on Computer Vision, 2016:483 - 499.
- [2] Jain A. Articulated people detection and pose estimation: Reshaping the future [C]//Computer Vision and Pattern Recognition. IEEE, 2012:3178 - 3185.
- [3] Gkioxari G, Hariharan B, Girshick R, et al. Using k-Poselets for Detecting People and Localizing Their Keypoints [C]//IEEE Conference on Computer Vision and Pattern Recognition. IEEE Computer Society, 2014:3582 - 3589.
- [4] Chen X, Yuille A. Parsing occluded people by flexible compositions [C]//Computer Vision and Pattern Recognition. IEEE, 2015:3945 - 3954.
- [5] Pishchulin L, Insafutdinov E, Tang S, et al. DeepCut: Joint Subset Partition and Labeling for Multi Person Pose Estimation [J]. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, 2016:4929 - 4937.
- [6] Insafutdinov E, Pishchulin L, Andres B, et al. DeeperCut: A Deeper, Stronger, and Faster Multi-person Pose Estimation Model [C]//European Conference on Computer Vision. Springer International Publishing, 2016:34 - 50.
- [7] Liu W, Anguelov D, Erhan D, et al. SSD: Single Shot MultiBox Detector [C]//European Conference on Computer Vision. Springer International Publishing, 2016:21 - 37.
- [8] Sapp B, Toshev A, Taskar B. Cascaded models for articulated pose estimation [C]//European Conference on Computer Vision. Springer-Verlag, 2010:406 - 420.
- [9] Dantone M, Gall J, Leistner C, et al. Human Pose Estimation Using Body Parts Dependent Joint Regressors [C]//IEEE Conference on Computer Vision and Pattern Recognition. IEEE, 2013:3041 - 3048.
- [10] Savarese S, Lee H, Telaprolu M, et al. An efficient branch-and-bound algorithm for optimal human pose estimation [C]//Computer Vision and Pattern Recognition. IEEE, 2012:1616 - 1623.
- [11] Kiefel M, Gehler P V. Human Pose Estimation with Fields of Parts [M]//Computer Vision—ECCV 2014. Springer International Publishing, 2014:331 - 346.
- [12] Hara K, Chellappa R. Computationally Efficient Regression on a Dependency Graph for Human Pose Estimation [C]//Computer Vision and Pattern Recognition. IEEE, 2013:3390 - 3397.
- [13] Toshev A, Szegedy C. DeepPose: Human Pose Estimation via Deep Neural Networks [C]//Computer Vision and Pattern Recognition. IEEE, 2014:1653 - 1660.
- [14] Ouyang W, Chu X, Wang X. Multi-source Deep Learning for Human Pose Estimation [C]//Computer Vision and Pattern Recognition. IEEE, 2014:2337 - 2344.
- [15] He K, Zhang X, Ren S, et al. Deep Residual Learning for Image Recognition [C]//2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, 2015:770 - 778.
- [16] Ren S, He K, Girshick R, et al. Faster R-CNN: towards real-time object detection with region proposal networks [C]//International Conference on Neural Information Processing Systems. MIT Press, 2015:91 - 99.
- [17] Caicedo J C, Lazebnik S. Active Object Localization with Deep Reinforcement Learning [C]//Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV). IEEE Computer Society, 2015: 2488 - 2496.
- [18] Bellver M, Giro-I-Nieto X, Marques F, et al. Hierarchical Object Detection with Deep Reinforcement Learning [J]. Advances in Parallel Computing, 2016, 31.
- [19] Mnih V, Kavukcuoglu K, Silver D, et al. Human-level control through deep reinforcement learning. [J]. Nature, 2015, 518(7540):529.
- [20] Huang G, Liu Z, Laurens V D M, et al. Densely Connected Convolutional Networks [C]//2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017). IEEE, 2016:2261 - 2269.
- [21] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition [EB]. arXiv preprint arXiv:1409.1556, 2014.
- [22] Iqbal U, Gall J. Multi-Person Pose Estimation with Local Joint-to-Person Associations [C]//European Conference on Computer Vision, ECCV 2016 Workshops, 2016:627 - 642.
- [23] Levinkov E, Uhrig J, Tang S, et al. Joint graph decomposition & node labeling: Problem, algorithms, applications [C]//IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, 2017.