

Application of Multi-Object Tracking with Siamese Track-RCNN to the Human in Events Dataset

Bing Shuai
Amazon
bshuai@amazon.com

Andrew Berneshawi
Amazon
bernea@amazon.com

Manchen Wang
Amazon
manchenw@amazon.com

Chunhui Liu
Amazon
chunhliu@amazon.com

Davide Modolo
Amazon
dmodolo@amazon.com

Xinyu Li
Amazon
xxnl@amazon.com

Joseph Tighe
Amazon
tighej@amazon.com

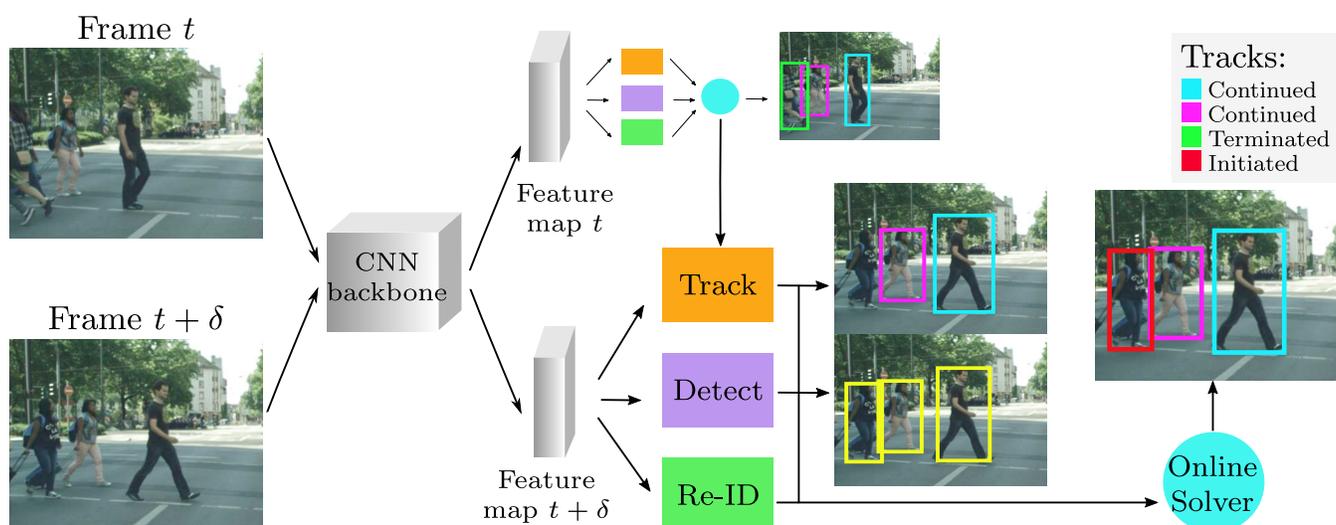


Figure 1: Siamese Track-RCNN unifies detection, tracking and re-id in a single network architecture. Importantly, these features share the same backbone, which results in low computation and efficient runtime.

ABSTRACT

Multi-object tracking systems often consist of a combination of a detector, a short term linker, a re-identification feature extractor and a solver that takes the output from these separate components and makes a final prediction. Differently, this work aims to unify all these in a single tracking system. Towards this, we propose Siamese Track-RCNN, a two stage detect-and-track framework which consists of three functional branches: (1) the detection branch localizes object instances; (2) the Siamese-based track branch estimates the

object motion and (3) the object re-identification branch re-activates the previously terminated tracks when they re-emerge. We used this design and apply it to the Human in Events [8] dataset.

CCS CONCEPTS

• Computing methodologies → Tracking.

KEYWORDS

Multi-object tracking, Siamese Track-RCNN

ACM Reference Format:

Bing Shuai, Andrew Berneshawi, Manchen Wang, Chunhui Liu, Davide Modolo, Xinyu Li, and Joseph Tighe. 2020. Application of Multi-Object Tracking with Siamese Track-RCNN to the Human in Events Dataset. In *Proceedings of the 28th ACM International Conference on Multimedia (MM '20)*, October 12–16, 2020, Seattle, WA, USA. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3394171.3416297>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '20, October 12–16, 2020, Seattle, WA, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7988-5/20/10...\$15.00

<https://doi.org/10.1145/3394171.3416297>

Challenging Success Cases

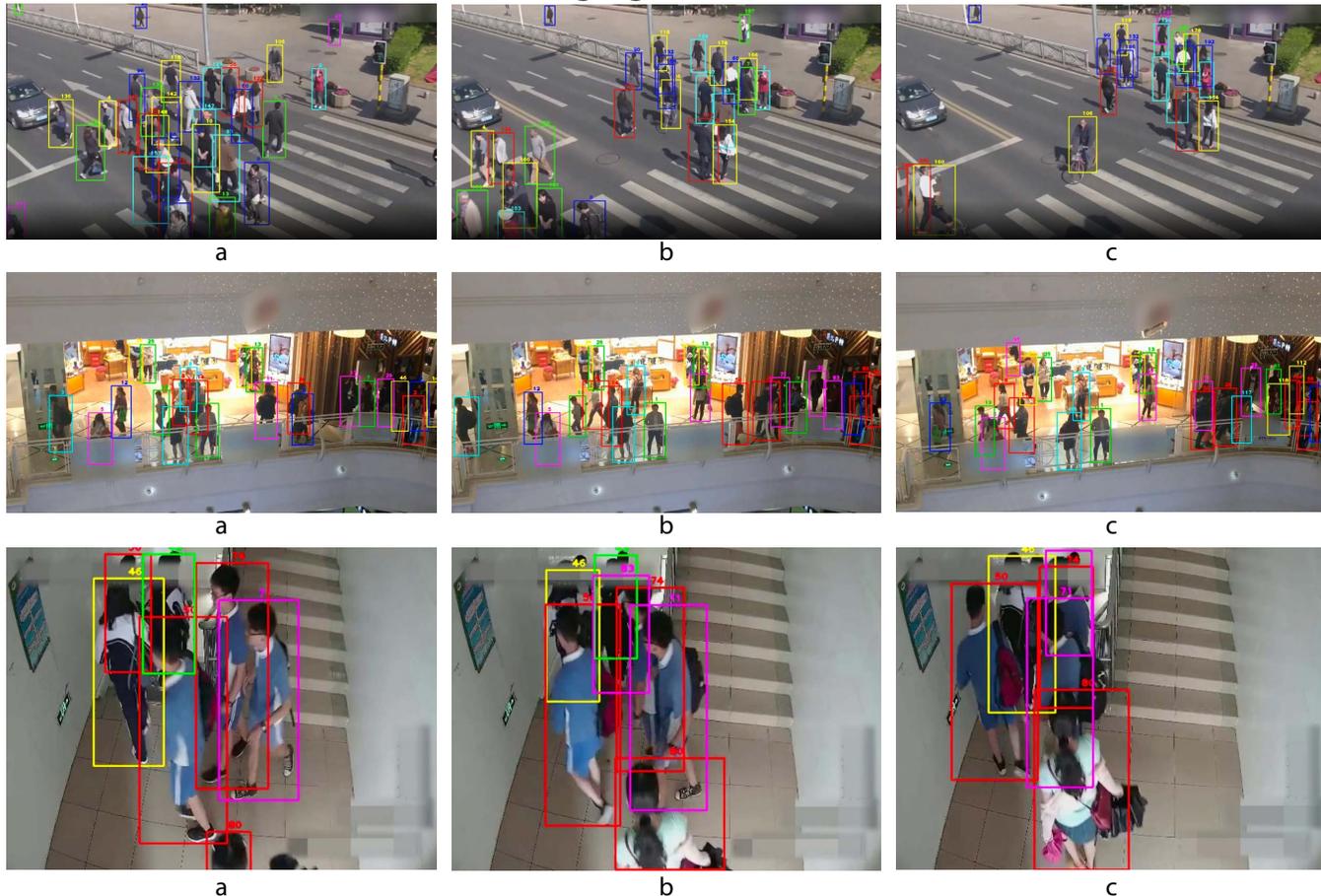


Figure 2: Three examples of challenging scenes where our tracker in the public detection setting is able to maintain tracks. We show three frames from each sequence non-uniformly sampled to highlight the strength of our tracker. In the first row we see a bicyclist (yellow track 106) enter the crowd and exit with the same track maintained. In the second row, we see generally impressive tracking with people tracked successfully through various challenging lighting conditions (see tracks 12, 26, and 29 in particular). In the third row, three people wearing the same clothes travel through occlusions but maintain all three tracks correctly (tracks 50, 71 and 74).

1 INTRODUCTION

This report outlines the method of Shuai et al. [11] that was applied to the Human in Events [8] challenge. This multi-object tracking method achieved a 3rd place ranking on the public detection leaderboard and a 4th place ranking on the private detection leaderboard.

Our tracking system is detailed in Shuai et al. [11] but we briefly outline how the system works here, for complete details please refer to [11]. Figure 1 shows an overview of the Track-RCNN system. Our MOT approach combines all the typical tracking sub-tasks (detection, short-term tracking and re-identification) in a single, unified network architecture. Our system performs both detection and association in a single forward pass of the network, which consists of a shared backbone and three task specific heads: detection, tracking and re-identification. At a high level, the detection branch localizes people entering the field of view, the track branch follows

them in the proceeding frames and the re-id branch is responsible for associations across longer periods of time. The detection branch is based on the popular Faster-RCNN architecture [4, 9] with a Region Proposal Network (RPN), followed by classification and regression of the generated proposals. The track branch matches the visual content from the last frames tracks to regions in roughly the same area in the following frame. The matching is performed internal to the network and is trained end-to-end. The matching mechanism is inspired by the GOTURN [5] single object tracking network. The tracking branch generates both an new estimated location for the object and a score for if the object is still in visible. During training binary cross entropy loss is used for the visibility output and a smooth ℓ_1 loss is used for the motion estimation. Finally, the reID branch learns to generate an embedding that is close

Table 1: Detailed result summary on for both public and private detection tasks.

Detection	MOTA (↑)	wMOTA (↑)	IDF1 (↑)	MT (↑)	ML (↓)	FP (↓)	FN (↓)	IDsw (↓)	IDsw DT (↓)	Frag (↓)	MOTP (↑)
Private	47.81	42.47	46.30	30.53	23.81	6399	27050	2913	93	2368	76.88
Public	50.55	45.38	47.21	30.64	26.96	4060	28071	2322	90	1942	77.65

in features space for two crops of the same person and more distant for different people. This branch is trained using a triplet loss.

Our online solver manages the initialization, continuation, termination and reinstatement of tracks. The detection branch is used to initialize a new track, unless the reID embedding is close to a previous track in which case it is reinstated. The track branch controls the continuation of tracks as well as the termination when it predicts a track is no longer visible. For more details of how the network is structured see [11].

In section 2 we present the details of how we applied [11] on the Human in Events [8] dataset and in section 3 we present our results from the challenge.

2 IMPLEMENTATION DETAILS

Network. For our challenge submission employ a DLA-34 [13] backbone with FPN [6]. The detection branch has the same structure of the popular Faster-RCNN [9], while the track and re-id branches consist of two fully connected layers of 1024 and 512 features, respectively. The re-id branch outputs an embedding of 128 features. We follow [5] and set the search extension ratio $r = 2$, practically doubling the target size, and we set the distance margin α empirically to 0.2. Finally, we use a ROI Align layer [4] and pool feature maps of 7×7 from frame regions. In the Track branch, we extract these features from both the target bounding box at time t and the enlarged search area at time $t + \delta$, even though the latter is twice as large. In other experiments we observed that enlarging the feature maps of the search area does not bring any improvement in performance.

Training. We pre-trained Siamese Track-RCNN backbone, detection, track branch (not the reID head) on a diverse set of person detection datasets: COCO [7], CrowdHuman [10], CUHK-SYSU [12], ETH Pedestrian [2], PRW [14] and TownCentre [3]. During the pre-training, we create image pairs with one sampled image and its augmented counterpart, which has been shifted and scaled (limit to 5% of image size), to simulate video data which is used to train the track branch. The model is trained over 50K iterations with a batch size of 16 image pairs and initial learning rate of 0.02. We decrease the learning rate by a factor of 10 after 30K and 40K iterations.

We use stochastic gradient descent to optimize our network on the Human in Events [8] training set for 15K iterations with an initial learning rate of 0.2. The learning rate is decayed by a factor of 10 after 7.5K and 12.5K iterations respectively. We use a fixed weight decay of 10^{-4} . Finally, we augment our training data of pairs by sampling them randomly within a 1 second temporal window, which is equivalent to setting $\delta = 30$ frames for 30fps videos (range: $[t, t + 30]$).

Inference. At inference we instead set $\delta = 1$ frame, as we aim to keep computation low and run Siamese Track-RCNN as a sliding

window. We feed the rich outputs of the three branches into an online solver that finalizes the ID of the localized people. Specifically, given a set of person localizations, our solver first merges those shared by both the detection and track branches (i.e., those that have intersection-over-union (IoU) > 0.3) and then terminates tracks that have a visibility score lower than 0.3 ($\hat{v} < 0.3$). Furthermore, it reinstates a previously terminated track when its embedding features are very similar to those of a newly localized person. In practice, the solver postpones this decision to after the new localized person has been tracked for a few frames. Then, it computes the average ℓ_2 embedding distance between the 5 most similar bounding boxes from the new and old tracks. If this value is less than 0.5, the track gets reinstated. Otherwise, a new track is initiated. To enable this feature comparison, our solver maintains and updates a small buffer that caches the embeddings of the terminated and ongoing tracks. We set the size of this buffer to 30 seconds, which offers a good trade-off between low memory consumption and enough temporal information.

3 RESULTS

Table 1 shows the results of our method as reported by the testing server on both public ¹ and private ² detection tracking tasks. For the public detection task we disable the detection head in our network and instead inject the provided public detections. We don't perform any type of bounding box re-scoring or refinement for the provided detections. For the private detections task, we mimic the public detection setup by training a heavy weight detector based on the DLA-169 [13] backbone with FPN [6] and deformable convolutions [1]. This detector provides detections to our tracker in the same way the public detections are used. We find that our heavy weight detector is still unable to outperform the provided detections and thus we achieve lower performance in the private settings. Our tracker is still strong enough to generate competitive results in this setting and achieve a 4th place ranking.

In figure 2 we show some challenging scenes where our tracker is able to maintain correct tracks. See the caption for the specific tracks that are successfully maintained. The first and third sequences in figure 2 show that the tracker head is able to successfully track each person foreword through changes between partial and no occlusion found in these crowded scenes. The second sequence highlights the trackers ability to track through large changes in lighting which often prove challenging for trackers as a persons appearance can change drastically in these cases.

In figure 3, we show hard cases that our tracker can not handle well. The first and third sequences show both poor quality detection

¹public detection results: <http://humanevents.org/tracker.html?tracker=1id=169>

²private detection results: <http://humanevents.org/tracker.html?tracker=2id=139>

Failure Cases

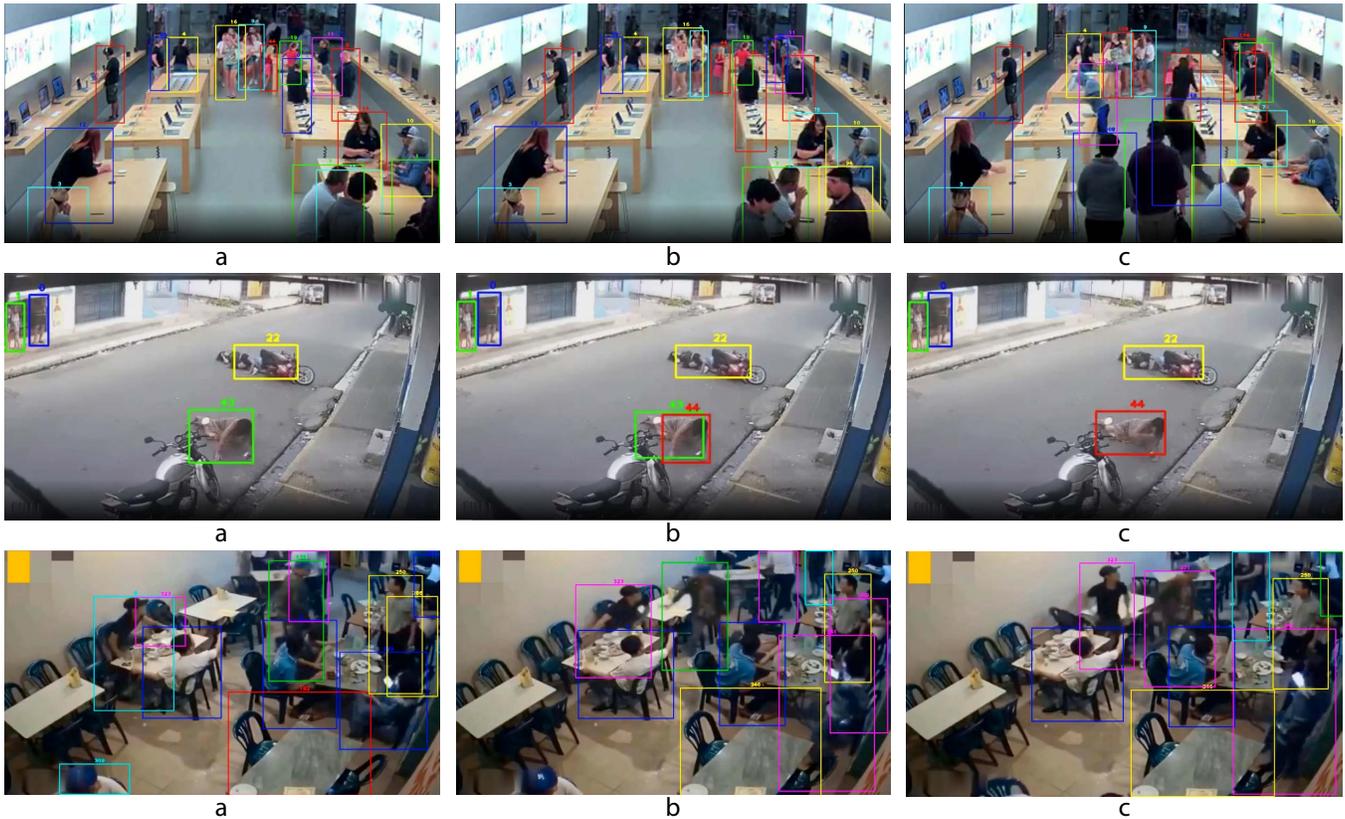


Figure 3: Three examples of places where our tracker in the public detection setting fails to detect or track. We show three frames from each sequence non-uniformly sampled to highlight the strength of our tracker. In the first row, the group at the top of the image is only detected as two people and the track jump to different people (track 9) as well as get split into two ids (track 16 and 158). The trend can be found through out the frame. In the second row, the track of the lower person is split due to a double detection in frame b. The second detection is continued to the next frame as track 44 and the original track 44 is terminated. In the third row, the green track in the center is split across two ids (175 and 371). There are also a number of false positives on chairs with one even fragmenting (tracks 182 and 346).

and tracking. We believe these errors a likely due to the low resolution and high compression observed in these videos. We believe training our model to better handle these compression signatures will help elevate these issues. The second sequence, shows a failure in our final solver. In frame b we get both a detection from our tracker and a partial detection from the detector. These two aren't considered duplicates by the solver and so both are attempted to be tracked forward. The red track (44) wins and thus the green track (43) gets terminated.

REFERENCES

- [1] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. 2017. Deformable convolutional networks. In *ICCV*.
- [2] Piotr Dollár, Christian Wojek, Bernt Schiele, and Pietro Perona. 2011. Pedestrian detection: An evaluation of the state of the art. *IEEE transactions on pattern analysis and machine intelligence* 34, 4 (2011), 743–761.
- [3] Jules. Harvey, Adam. LaPlace. 2019. *MegaPixels: Origins, Ethics, and Privacy Implications of Publicly Available Face Recognition Image Datasets*. <https://megapixels.cc/>
- [4] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *ICCV*.
- [5] David Held, Sebastian Thrun, and Silvio Savarese. 2016. Learning to track at 100 fps with deep regression networks. In *ECCV*.
- [6] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.
- [7] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *ECCV*.
- [8] Weiyao Lin, Huabin Liu, Shizhan Liu, Yuxi Li, Guo-Jun Qi, Rui Qian, Tao Wang, Nicu Sebe, Ning Xu, Hongkai Xiong, and Mubarak Shah. 2020. Human in Events: A Large-Scale Benchmark for Human-centric Video Analysis in Complex Events. *arXiv preprint arXiv:2005.04490* (2020).
- [9] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*.
- [10] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. 2018. CrowdHuman: A Benchmark for Detecting Human in a Crowd. *arXiv preprint arXiv:1805.00123* (2018).
- [11] Bing Shuai, Andrew G Berneshawi, Davide Modolo, and Joseph Tighe. 2020. Multi-Object Tracking with Siamese Track-RCNN. *arXiv preprint arXiv:2004.07786* (2020).

- [12] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. 2016. End-to-end deep learning for person search. *arXiv preprint arXiv:1604.01850* 2, 2 (2016).
- [13] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. 2018. Deep layer aggregation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2403–2412.
- [14] Liang Zheng, Hengheng Zhang, Shaoyan Sun, Manmohan Chandraker, Yi Yang, and Qi Tian. 2017. Person re-identification in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1367–1376.