

Simple Unsupervised Multi-Object Tracking

Shyamgopal Karthik¹ Ameya Prabhu² Vineet Gandhi¹

¹ Center for Visual Information Technology
Kohli Center on Intelligent Systems, IIIT Hyderabad, India

² University of Oxford
{shyamgopal.karthik@research,vgandhi@}.iiit.ac.in
ameya@robots.ox.ac.uk

Abstract. Multi-object tracking has seen a lot of progress recently, albeit with substantial annotation costs for developing better and larger labeled datasets. In this work, we remove the need for annotated datasets by proposing an unsupervised re-identification network, thus sidestepping the labeling costs entirely, required for training. Given unlabeled videos, our proposed method (SimpleReID) first generates tracking labels using SORT [3] and trains a ReID network to predict the generated labels using crossentropy loss. We demonstrate that SimpleReID performs substantially better than simpler alternatives, and we recover the full performance of its supervised counterpart consistently across diverse tracking frameworks. The observations are unusual because unsupervised ReID is not expected to excel in crowded scenarios with occlusions, and drastic viewpoint changes. By incorporating our unsupervised SimpleReID with CenterTrack trained on augmented still images, we establish a new state-of-the-art performance on popular datasets like MOT16/17 without using tracking supervision, beating current best (CenterTrack) by 0.2-0.3 MOTA and 4.4-4.8 IDF1 scores. We further provide evidence for limited scope for improvement in IDF1 scores beyond our unsupervised ReID in the studied settings. Our investigation suggests reconsideration towards more sophisticated, supervised, end-to-end trackers [56,5] by showing promise in simpler unsupervised alternatives.

Keywords: Multi-Object Tracking, Re-Identification, Unsupervised Learning

1 Introduction

Understanding human interactions and behaviour over videos has been a fundamental problem in computer vision with applications in action recognition, sports video analytics, and assistive tech and requires tracking multiple people over time. Multi-object trackers broadly consist of two key components: (i) A spatio-temporal association model which associates boxes in nearby frames to create clusters of tracklets, and (ii) A re-identification model which associates tracklets over larger windows to deal with complexities in tracking such as occlusions and target interactions. Re-identification is a major challenge in tracking,

with sophisticated supervised approaches requiring expensive annotations to assign trajectories across frames to every single person in a video. Availability of labeled datasets [36,37] has alleviated the problem. For instance IDF1 (MOTA) scores have improved from 51.3(48.8) [49] to 59.9 (55.9) [5] on the MOT16 [37] benchmark in the past 3 years.

There has been a growing need to annotate larger tracking datasets with the aim of improving re-identification (ReID) models. However, annotating tracking datasets require hefty labeling costs and scale poorly with dataset size. To illustrate the effort and cost required, annotating 6 minutes worth of video of the MOT15 benchmark [27] using the standard annotation procedures would take at least 22 hours of annotation time [36]. Annotating just twenty-six hours of video data (VIRAT dataset [39]) with state-of-the-art protocols in place [39,50] costs tens of thousands of dollars. We propose to learn our model in an unsupervised manner in the free-labels paradigm (Section 6.3.2 in [21]) in a two-step manner. We first generate tracking labels given unlabeled videos and the corresponding set of detections. Then, we learn a ReID network to predict the generated label given an input image. To the best of our knowledge, ours is the first work to propose unsupervised ReID models for multi-object tracking and completely do away with the tremendous annotation costs for tracking datasets. Throughout the paper, we consider supervision only in the context of sidestepping trajectory-level annotations. Using off-the-shelf detectors [41,40,7] trained on COCO is not viewed as supervision in our context. The proposed ReID network complements the unsupervised spatio-temporal association models [53,1] proposed in the prior art, leading to a more complete unsupervised tracking framework.

We go one step further and aim to test the limits of our unsupervised tracking paradigm. We empirically test for two desiderata w.r.t IDF1 scores: (i) Our unsupervised ReID should perform significantly better than naive ReID methods when incorporated into any tracker; (ii) Our unsupervised ReID should achieve performance equivalent to the original supervised counterpart. We demonstrate that we are able to achieve these desiderata consistently across datasets, detectors, and diverse trackers. The resultant unsupervised tracker, when combined with CenterTrack [69] trained on single images, achieves state-of-the-art performance on the MOT16/17 test challenge server. We beat the latest supervised trackers by large margins, outperforming CenterTrack by 0.3 MOTA, and 4.8 IDF1 scores. We then demonstrate that there is limited scope for further improvement beyond our proposed unsupervised ReID by demonstrating that the Oracle counterpart of our ReID model makes only minor gains.

We would also like to highlight that while our work is conceptually simple, the contributions made are significant. We expect our investigation to be of significant interest to the MOT community by demonstrating that simple unsupervised ReID is sufficient even in crowded scenarios with occlusions and person interactions. Our investigation contrasts the current shift towards using more supervised, end-to-end trackers for MOT Challenge datasets. We hope our work spurs research in the unsupervised MOT paradigm, exploring extensions to other

tracking scenarios (3D/vehicles/pose tracking) and do away with the labeling effort wherever unnecessary.

2 Related Works

Monocular 2D multi-object tracking on videos is an extensively studied problem. [13] offers a comprehensive review of works on MOT Challenge datasets. A popular paradigm is to model the detections as a graph. Various approaches have been proposed here including using network flows [62], graph cuts [49], MCMC [60] and minimum cliques [61] if the entire video is provided beforehand (batch processing). In scenarios where we get frame-by-frame input, Hungarian matching [53,3], greedy matching [69] and Recurrent Neural Networks [15,43] are popularly used models for sequential prediction (online processing). The association metrics/cost functions used by these consist of (i) Spatio-temporal relations (ii) Re-identification.

Spatiotemporal relations: There has been much investigation into appearance-free methods for the spatio-temporal association. Basic methods proposed include using Intersection-Over-Union (IoU) between detections [4] or incorporating a velocity model using a Kalman filter [3]. The velocity model can also be learned using Recurrent Neural Networks [15,43]. The complexity of assigning pairwise costs can be further increased by incorporating additional cues from head/joint detectors [6,19], segmentation [38], activity recognition [10], or keypoint trajectories [9]. Recent approaches leverage appearance-reliant pre-trained bounding box regressors from object detection [1] or single object tracking [56,11] pipelines to regress the bounding box in the next frame. Since most of the above models are unsupervised (requiring no tracking annotations), they complement our work and can be incorporated with our proposed approach for creating efficacious unsupervised trackers.

ReID across multiple cameras: Supervised training of CNNs [68] on large labeled datasets [65,30] has given excellent results for ReID across multiple cameras. In addition to this, there have been approaches to exploit the pose information using off-the-shelf body pose detectors [47,48]. Attention mechanisms have also been explored to capture the important regions in the foreground [45,46]. Generative models have been employed to augment the training data for improved performance [66,31]. We recommend this excellent survey [58] for a complete review. In contrast, we work on tracking with a single camera, with reasonable frame-rates (no drastic appearance variations). Additionally, the objective is to distinguish the target pedestrian among a small set of different looking pedestrians in a given frame, with the aid of additional detection information. Hence, we believe our simple, noisy unsupervised re-identification model might suffice. Sophisticated unsupervised ReID networks [32,29] designed for multiple cameras ReID may not be required for MOT.

ReID for monocular 2D tracking: Re-identification has been a major challenge in tracking, with matching using similarity between CNN features being the dominant approach [42]. Past works have proposed different methods

to train the CNN ranging from using siamese networks [26] with triplet loss, further augmented by hard negative mining [1] or other metric learning losses like cosine loss [53]. Incorporating a combination of loss functions [34] or pose information [49] as well as fine-tuning the ReID model on the test sequence [34]. All the above ReID networks are supervised and fairly complex to train. We are the first work to demonstrate that simple unsupervised ReID networks are sufficient for this context. It is important to note that in most MOT pipelines, this is the only component that uses tracking annotations.

Evaluation metric for MOT: Multi-Object Tracking Accuracy (MOTA) is not a good metric to illustrate ReID performance because it focuses on object coverage and therefore is dominated by false negatives. An excellent detector can achieve high MOTA scores despite being a poor tracker with a large number of ID switches [69]. Identity-F1 (IDF1) has been shown to measure long consistent tracks without switches and widely shown [35,13] to be a better metric for tracking performance. We accordingly focus and emphasize on IDF1 scores.

End-to-end supervised MOT: Recent works circumvent the above paradigm either partially or completely by learning the MOT solver using end-to-end supervision. Early works [51,44] performed end-to-end learning in the min-cost flow data association framework. Recently, approaches like [56] and [5] perform end-to-end optimization by introducing differentiable forms of Hungarian matching and clustering formulation, respectively. Parallel works [69,64,52] attempt to perform simultaneous object detection, data association, and sometimes re-identification in a single network. Most notable among these, CenterTrack [69] is capable of training the detector using only augmentations of still images. These methods involving joint detection and tracking deliver high performance at real-time inference speeds but require high annotation costs. Our work differs in principle by removing and replacing supervised components yet outperforming these trackers, without incurring the associated labelling cost.

3 Approach

Our goal is to leverage the abundance of unlabeled videos to learn ReID models (without manual cost). Our unsupervised learning method can be categorized as learning by generating labels (Ref. Section 6.3.2 of [21]). In a nutshell, given unlabeled videos and corresponding bounding boxes, we first generate tracking labels. We then learn a ReID network by predicting the generated label given a detection.

3.1 Framework: Learning by generating tracking labels

Here, we describe the two parts of our proposed framework in detail: (i) Generating the labels, and (ii) Learning the network. **Generating labels:** Given a set of videos, each video is passed independently through an object detector. An **unsupervised spatiotemporal association model** from the list given in Table 1 (left) is then run through the detections to obtain short contiguous tracks or

Model	Ref	Training Strategy	Ref
Kalman filter+Hungarian matching	[3]	Crossentropy	[49]
IoU based tracking	[4]	Triplet+hard negative mining	[1]
Network Flow	[62]	Contrastive	[25]
Linear Programming	[28]	SymTriplet	[63]
Conditional Random Fields(CRFs)	[38]	Cosine Loss	[53]
Markov Decision Proceses(MDPs)	[54]	Joint Detections	[49]
Recurrent Neural Networks(RNNs)	[43]	Verification+Classification Loss	[34]
Bounding Box Regression	[1]		

Table 1. Approaches use for Spatiotemporal data association (Left). Loss functions and methods used to train CNNs for Appearance modeling (Right). We choose the simplest approach for both these components.

tracklets (set of associated detections of the same person over time). Examples of spatiotemporal models can range from tracking using a constant velocity assumption with Kalman filtering [3] (bounding box information only) to incorporating appearance features by using pre-trained bounding box regression from object detection pipelines to regress the bounding box [1] in the next frame. Now to cluster/associate detections, we can use online methods like greedy/Hungarian matching or expensive offline methods like graph-cuts. Ultimately, the output of this step is a set of **noisy track labels** for each video, resulting in a pool of labeled video tracklets.

Training ReID models: Now, given noisy track labels per video, the task is to learn a ReID model using any of the methods given in Table 1 (right). In absence of trajectory level supervision, the challenge here is to explore ways to harness the given regularities in data (in form of tracklets). There are two simple assumptions which can help the cause: (i) The videos are independent of each other (i.e., no common tracks between any two videos), and (ii) the tracklets within a video are independent of each other (i.e., each tracklet belongs to a different person). If both the assumptions are followed then each tracklet can be considered as an **independent class**. The simplest option which follows is to train at network to predict a label given an image, optimized with cross-entropy loss (with number of classes equalling to the number of tracklets). However, assumption (ii) may break in cases like missed detections and occlusions and may result into multiple tracklets for the same person in a video.

An alternate option (by relaxing assumption (ii)) could be to **form positive pairs from the same tracklet and negative pairs from across other videos** or simultaneous tracks from the same video. Such pairing can enable learning **Siamese networks** to compare two images and predict whether they are the same person or not. They can be trained with pairwise losses such as **contrastive loss** [25] or **triplet loss with hard-negative mining** [1], or more complex ones like **symtriplet** [63] or the **group loss** [14], resulting in a trained ReID network.

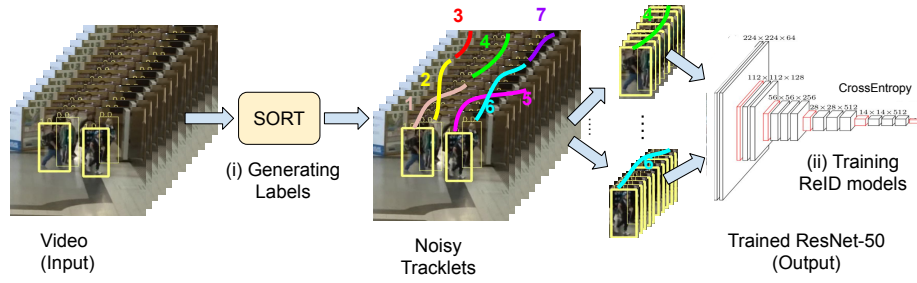


Fig. 1. Overview of our approach: Given a video with detections, we use SORT [3] to simulate noisy tracking labels. Then, we train the ReID network (ResNet50) to predict the track label for each input image.

3.2 Our method

We use simple methods to both simulate labels and learn the ReID network, as illustrated in Figure 1. In step (i), we only utilize the bounding boxes and use Kalman filtering combined with Hungarian matching to simulate labels. Since we use no appearance information, our tracking labels are noisy. In step (ii), we proceed by making **both the aforementioned assumptions** that no two videos or tracklets share common labels. We assign a unique label to each tracklet and train a network with cross-entropy loss to predict this label given any image from that tracklet. At inference time, we integrate our ReID model into existing frameworks by simply replacing their models with ours, with no other changes. In CenterTrack, we extract tracks using its unsupervised model and refine it with our ReID network using a DeepSORT framework. Although we are aware that some enhancements can be performed to our proposed process (e.g., using a siamese framework), we show in subsequent sections that simpler choices alone are sufficient to match the performance of supervised networks.

4 Experiments

In a nutshell, in this section we incorporate our developed unsupervised ReID model (**SimpleReID**) into various trackers and show compelling evidence for three results: (i) our unsupervised tracker obtains state-of-the-art tracking performance on MOT16/17, outperforming recent works (ii) naive unsupervised trackers can replace their supervised counterparts consistently (iii) there is **limited scope for improvement beyond our unsupervised ReID** complemented with better detectors in settings we tested.

4.1 Experimental Setup

Datasets: We evaluate our performance on the standard multi-object tracking benchmark – MOT Challenge – which consists of several challenging pedestrian

tracking sequences with frequent occlusions, crowded scenes with sequences varying in their angle of view, size of objects, camera motion, and frame rate. It contains two challenging tracking benchmarks, namely MOT16 and MOT17 [37]. They both share the same training and testing sequences, but MOT16 provides only DPM [16] detections, whereas MOT17 provides two additional sets of public detections (namely Faster R-CNN [41] and SDP [57]) and has more accurate ground truth. The primary metrics used for measuring performance are MOTA [2] and IDF1, which are a combination of simpler metrics like False Positives, False Negatives, and ID Switches.

Implementation details: We obtain our SimpleReID model by training a ResNet50 [17] backbone popularly used by trackers for a fair comparison. We train the model with tracklets generated by SORT [3] on the PathTrack [36] dataset to test generalization to unseen MOT16/17 data. We perform analysis studies on the entire training dataset and report results on MOT Challenge hidden test set ³. Our model was implemented using PyTorch and Torchreid [67] and trained on a GTX1080Ti GPU. For any tracker used [53,1], we utilize the implementations provided by the authors, leaving all the hyperparameters unchanged and simply replacing their supervised ReID model with SimpleReID. We use the CenterTrack model trained with single images w.r.t augmentations and incorporate the SimpleReID model using the DeepSORT framework. Our code and pretrained models will be released upon acceptance of the paper.

4.2 MOT Challenge Benchmark Evaluation

We submit our best performing unsupervised tracker to the MOT Challenge Benchmark. The submitted tracker consists of our proposed SimpleReID model incorporated with CenterTrack [69] for bounding box regression using public detections. We compare the performance on the MOT Challenge test set with state-of-the-art supervised trackers and provide results in Table 2. Surprisingly, we observe that our developed unsupervised tracker outperforms all supervised trackers on MOT16/17 setting a new state-of-the-art in terms of MOTA and IDF1 scores among all trackers on public detections.

We beat the previous best tracker (CenterTrack) by 0.2/0.3 MOTA and 4.4/4.8 IDF1 scores on MOT16/MOT17, respectively. The significant increase in IDF1 score can be entirely attributed to the efficacy of our SimpleReID model, because while CenterTrack is a good detector, it cannot maintain long tracks which is compensated by using our appearance features for Re-identification. We reduce ID switches made by CenterTrack by nearly 3x, achieving the lowest ID switches compared to other online trackers.

4.3 Analysis

Past literature [49,34] indicates that unsupervised ReID is unlikely to excel in crowded scenarios due to the complexities of tracking in such scenes. In this

³ The MOT Challenge web page: <https://motchallenge.net>.

Detector Method		Published	Unsup	MOTA \uparrow	IDF1 \uparrow	IDS \downarrow	FP \downarrow	FN \downarrow
MOT16								
Batch	GCRA [33]	ICME18	\times	48.2	48.6	821	5104	88586
	HCC [34]	ACCV18	\times	49.3	50.7	391	5333	86795
	LMP [49]	CVPR17	\times	48.8	51.3	481	6654	86245
	MPN [5]	CVPR20	\times	55.9	59.9	431	7086	72902
Online	AMIR [43]	ICCV17	\times	47.2	46.3	774	2681	92856
	KCF [11]	WACV19	\times	48.8	47.2	906	5875	86567
	RAR16 [15]	WACV18	\times	45.9	48.8	648	6871	91173
	MOTDT [8]	ICME18	\times	47.6	50.9	792	9253	85431
	STRN [55]	ICCV19	\times	48.5	53.9	747	9038	84178
	DeepMOT [56]	CVPR20	\times	54.8	53.4	645	2955	78765
	CenterTrack [69]	Arxiv20*	\checkmark	62.2	54.1	1677	5433	61767
	DMAN [70]	ECCV18	\times	46.1	54.8	532	7909	89874
	Tracktor++v2 [1]	ICCV19	\times	56.2	54.9	617	2394	76844
	MIFT	Arxiv20*	\times	60.1	56.9	739	6964	65044
	Ours	-	\checkmark	62.4	58.5	588	5909	61981
MOT17								
Batch	MHT [23]	CVPR15	\times	50.7	47.2	6543	46638	224955
	FWT [18]	CVPRW18	\times	51.3	47.6	2648	24101	247921
	MHT-bLSTM [24]	ECCV18	\times	47.5	51.9	2069	25981	268042
	jCC [22]	TPAMI18	\times	51.2	54.5	1802	25937	247822
	MPN [5]	CVPR20	\times	55.7	59.1	1433	25013	223531
Online	FAMNet [12]	ICCV19	\times	52.0	48.7	3072	14138	253616
	DeepMOT [56]	CVPR20	\times	56.7	52.1	2351	8895	233206
	MOTDT [8]	ICME18	\times	50.9	52.7	2474	24069	250768
	CenterTrack [69]	Arxiv20*	\checkmark	61.4	53.3	5326	15520	196886
	Tracktor++v2 [1]	ICCV19	\times	56.3	55.1	1987	8666	235449
	DMAN [70]	ECCV18	\times	48.2	55.7	2194	26218	263608
	MIFT [20]	Arxiv20*	\times	60.1	56.4	2556	14966	206619
	STRN [55]	ICCV19	\times	50.9	56.5	2397	25295	249365
	Ours	-	\checkmark	61.7	58.1	1864	16872	197632

Table 2. Results on the MOT Challenge test set benchmark using public detections. Unsup indicates approach does not need supervision (no tracking labels required). * are recent parallel works. Up/down arrows indicate higher/lower is better.

subsection, we provide two sets of evidence to demonstrate that SimpleReID indeed performs well across diverse scenarios: (i) We show that the test performance of SimpleReID (on unseen videos) is equivalent to that of a supervised ReID model, on its training set itself (ii) We show that SimpleReID achieves the above desiderata even with simple trackers which are highly reliant on the ReID component.

Limits of unsupervised ReID: We test the limits of SimpleReID by comparing the performance of our model with supervised models. We perform experiments across various weaker scenarios such as having no ReID, or using pretrained-ImageNet as-is, and show that these perform significantly worse than SimpleReID - proving that SimpleReID is important to match supervised performance. We first train another recent supervised tracker, Tracktor++v2[1], which

ReID	MOTA \uparrow	IDF1 \uparrow	ReID	MOTA \uparrow	IDF1 \uparrow
MOT16					
	DPM			POI	
None	57.6	62.0	None	68.3	67.6
ImageNet	57.6	62.0	ImageNet	68.3	67.7
Ours	57.6	62.6	Ours	68.5	69.5
Supervised	57.6	62.5	Supervised	68.5	69.4
MOT17					
	FRCNN			POI	
None	61.6	64.6	None	68.5	67.6
ImageNet	61.6	64.7	ImageNet	68.5	67.6
Ours	61.7	65.2	Ours	68.6	69.4
Supervised	61.7	65.2	Supervised	68.6	69.3

Table 3. Ablation study comparing the performance of different ReID models within the Tracktor [1] framework. We observe that our unsupervised SimpleReID achieves the same performance (IDF1 scores) as supervised ReID. DPM, FRCNN and POI correspond to different detectors.

uses bounding box regression along with a supervised ReID model to predict the position of an object in the next frame. We train the supervised ReID model on the training data for MOT16/ MOT17 and then **benchmark the performance on the same training set**. In contrast, this data is new to our SimpleReID models, i.e., have not seen these videos previously. Our experiment results are tabulated in Table 3. We observe that using ImageNet-pretrained ReID somewhat improves IDF1 scores compared to using no ReID network at all, but fails to achieve the upper bound by a considerable margin. Our SimpleReID approach successfully recovers the remaining performance gap. This is achieved consistently across different variations.

ReID-reliant unsupervised tracking: Due to the **low dependence** of Tracktor on its ReID model, one may argue that it might not be the best framework for evaluation of ReID models in tracking. Hence, we also perform the same experiments on a popular tracker DeepSORT [53] that is highly reliant on the ReID model used, since the only visual features it receives is from the ReID network. We replace the supervised ReID model used in DeepSORT with different ReID methods and tabulate results in Table 4. First, we observe that replacing supervised ReID with random features causes a severe drop in performance over supervised counterpart, with MOTA score decreasing by 9.4% and IDF1 decreasing by 31.3%, demonstrating the degree of reliance on ReID in the DeepSORT framework. When substituted with features from an ImageNet-pretrained ResNet, we get a similar result: a significant improvement over SORT, yet much lower than supervised ReID performance. We further benchmark with a supervised ReID model trained on **Market1501 dataset** [65] and observe lower performance compared to the ImageNet-pretrained model, indicating that features learned for

ReID	MOTA \uparrow	IDF1 \uparrow	ReID	MOTA \uparrow	IDF1 \uparrow
MOT16-POI			MOT17-POI		
No ReID	58.1	57.1	No ReID	57.9	56.9
Random	51	34.6	Random	50.7	34.3
ImageNet	60.3	62	ImageNet	59.9	61.6
Market1501	60.3	61.5	Market1501	59.9	61.1
Ours	60.5	65.9	Ours	60.1	65.5
Supervised	60.4	65.9	Supervised	60	65.5

Table 4. Ablation study comparing the performance of different ReID models within the DeepSORT [53] framework. We observe that our unsupervised SimpleReID achieves the same performance (IDF1 scores) as supervised ReID.

Detector	SimpleReID		Oracle ReID+Kill+MM		
	MOTA \uparrow	IDF1 \uparrow	MOTA \uparrow	IDF1 \uparrow	IDF1 Gain
YOLOv3 [40]	56.5	62.5	61.5	66.1	3.6
DPM [16]	58.5	62.9	62.4	66.3	3.4
Faster-RCNN [41]	61.7	65.2	65.5	68.5	3.3
HTC [7]	67.7	68.1	75.6	70.5	2.4
SDP [57]	67.7	68.1	73.0	70.6	2.5
POI [59]	68.6	69.4	73.5	71.4	2.0

Table 5. Ablation study comparing the difference between performance of SimpleReID across detectors on MOT17. We observe that the difference decreases from 3.6 to 2.0 with improved detectors.

cross-camera person-ReID datasets without trajectory annotations do not transfer to multi-object tracking. Lastly, we observe that our unsupervised SimpleReID covers the remaining performance gap, as seen above.

Scope for improvement in ReID: We further explore the best performance achievable by a ReID network using the Tracktor framework and explore the scope for further improvement of our SimpleReID. To obtain the possible best performance, we test Tracktor with an **Oracle ReID** [1] and observe that there is a 3.3 IDF1 score gap between SimpleReID and the Oracle. We repeat the same experiment with the latest off-the-shelf detectors and tabulate the results in Table 5. We observe that with modern detectors, the gap between SimpleReID and the corresponding oracles is small enough to limit the scope for further improvement.

Overall, we conclude that unsupervised SimpleReID counterintuitively matches the limiting performance of supervised counterparts in difficult MOT scenarios, by leveraging only unlabeled videos. Since our model works in extreme cases such as DeepSORT, where tracking is entirely reliant on the ReID model for encoding appearance information, we expect that the efficacy of SimpleReID will generalize to other trackers as well. We demonstrated the potential of unsupervised trackers by outperforming all supervised MOT16/17 trackers, setting a new state-of-the-

art in MOTA and IDF1 scores and performing close to the optimal ReID. If it is indeed generalizable, we believe that this work has significant implications for research in supervised ReID for tracking.

5 Conclusion

We propose the first step in the direction of developing unsupervised re-identification for MOT and demonstrate that our simple approach performs at par with supervised counterparts across diverse setups. When combined with recent unsupervised association models [56,1], we obtain accurate unsupervised trackers. The tracker we submit ranks first in the MOT Challenge, beating all the latest supervised approaches. Our investigation suggests reconsideration on whether the shift towards more complex, supervised, end-to-end MOT models is necessary. We hope our work is useful to sidestep high annotation costs otherwise thought to be a requirement necessary to feed the data-hungry supervised trackers.

References

1. Bergmann, P., Meinhardt, T., Leal-Taixe, L.: Tracking without bells and whistles. In: ICCV. pp. 941–951 (2019)
2. Bernardin, K., Stiefelhagen, R.: Evaluating multiple object tracking performance: the clear mot metrics. EURASIP Journal on Image and Video Processing **2008**, 1–10 (2008)
3. Bewley, A., Ge, Z., Ott, L., Ramos, F., Upcroft, B.: Simple online and realtime tracking. In: ICIP. pp. 3464–3468 (2016)
4. Bochinski, E., Eiselein, V., Sikora, T.: High-speed tracking-by-detection without using image information. In: AVSS. pp. 1–6 (2017)
5. Brasó, G., Leal-Taixé, L.: Learning a neural solver for multiple object tracking. In: CVPR (2020)
6. Chari, V., Lacoste-Julien, S., Laptev, I., Sivic, J.: On pairwise costs for network flow multi-object tracking. In: CVPR. pp. 5537–5545 (2015)
7. Chen, K., Pang, J., Wang, J., Xiong, Y., Li, X., Sun, S., Feng, W., Liu, Z., Shi, J., Ouyang, W., et al.: Hybrid task cascade for instance segmentation. In: CVPR. pp. 4974–4983 (2019)
8. Chen, L., Ai, H., Zhuang, Z., Shang, C.: Real-time multiple people tracking with deeply learned candidate selection and person re-identification. In: ICME. pp. 1–6 (2018)
9. Choi, W.: Near-online multi-target tracking with aggregated local flow descriptor. In: ICCV. pp. 3029–3037 (2015)
10. Choi, W., Savarese, S.: A unified framework for multi-target tracking and collective activity recognition. In: ECCV. pp. 215–230 (2012)
11. Chu, P., Fan, H., Tan, C.C., Ling, H.: Online multi-object tracking with instance-aware tracker and dynamic model refreshment. In: WACV. pp. 161–170 (2019)
12. Chu, P., Ling, H.: Famnet: Joint learning of feature, affinity and multi-dimensional assignment for online multiple object tracking. In: ICCV. pp. 6172–6181 (2019)
13. Ciaparrone, G., Sánchez, F.L., Tabik, S., Troiano, L., Tagliaferri, R., Herrera, F.: Deep learning in video multi-object tracking: A survey. Neurocomputing (2020)

14. Elezi, I., Vascon, S., Torcinovich, A., Pelillo, M., Leal-Taixe, L.: The group loss for deep metric learning. arXiv preprint arXiv:1912.00385 (2019)
15. Fang, K., Xiang, Y., Li, X., Savarese, S.: Recurrent autoregressive networks for online multi-object tracking. In: WACV. pp. 466–475 (2018)
16. Felzenszwalb, P., McAllester, D., Ramanan, D.: A discriminatively trained, multi-scale, deformable part model. In: CVPR. pp. 1–8 (2008)
17. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016)
18. Henschel, R., Leal-Taixé, L., Cremers, D., Rosenhahn, B.: Improvements to frank-wolfe optimization for multi-detector multi-object tracking. arXiv preprint arXiv:1705.08314 (2017)
19. Henschel, R., Leal-Taixé, L., Cremers, D., Rosenhahn, B.: Fusion of head and full-body detectors for multi-object tracking. In: CVPR-W. pp. 1428–1437 (2018)
20. Huang, P., Han, S., Zhao, J., Liu, D., Wang, H., Yu, E., Kot, A.C.: Refinements in motion and appearance for online multi-object tracking. arXiv preprint arXiv:2003.07177 (2020)
21. Jing, L., Tian, Y.: Self-supervised visual feature learning with deep neural networks: A survey. arXiv preprint arXiv:1902.06162 (2019)
22. Keuper, M., Tang, S., Andres, B., Brox, T., Schiele, B.: Motion segmentation & multiple object tracking by correlation co-clustering. TPAMI **42**(1), 140–153 (2018)
23. Kim, C., Li, F., Ciptadi, A., Rehg, J.M.: Multiple hypothesis tracking revisited. In: ICCV. pp. 4696–4704 (2015)
24. Kim, C., Li, F., Rehg, J.M.: Multi-object tracking with neural gating using bilinear lstm. In: ECCV. pp. 200–215 (2018)
25. Kim, M., Alletto, S., Rigazio, L.: Similarity mapping with enhanced siamese network for multi-object tracking. arXiv preprint arXiv:1609.09156 (2016)
26. Leal-Taixé, L., Canton-Ferrer, C., Schindler, K.: Learning by tracking: Siamese cnn for robust target association. In: CVPR-W. pp. 33–40 (2016)
27. Leal-Taixé, L., Milan, A., Reid, I., Roth, S., Schindler, K.: Motchallenge 2015: Towards a benchmark for multi-target tracking. arXiv preprint arXiv:1504.01942 (2015)
28. Leal-Taixé, L., Pons-Moll, G., Rosenhahn, B.: Everybody needs somebody: Modeling social and grouping behavior on a linear programming multiple people tracker. In: ICCV-W. pp. 120–127 (2011)
29. Li, M., Zhu, X., Gong, S.: Unsupervised person re-identification by deep learning tracklet association. In: ECCV. pp. 737–753 (2018)
30. Li, W., Zhao, R., Xiao, T., Wang, X.: Deepreid: Deep filter pairing neural network for person re-identification. In: CVPR (2014)
31. Li, X., Wu, A., Zheng, W.S.: Adversarial open-world person re-identification. In: ECCV. pp. 280–296 (2018)
32. Lin, Y., Dong, X., Zheng, L., Yan, Y., Yang, Y.: A bottom-up clustering approach to unsupervised person re-identification. In: AAAI. pp. 8738–8745 (2019)
33. Ma, C., Yang, C., Yang, F., Zhuang, Y., Zhang, Z., Jia, H., Xie, X.: Trajectory factory: Tracklet cleaving and re-connection by deep siamese bi-gru for multiple object tracking. In: ICME. pp. 1–6 (2018)
34. Ma, L., Tang, S., Black, M.J., Van Gool, L.: Customized multi-person tracker. In: ACCV. pp. 612–628 (2018)
35. Maksai, A., Fua, P.: Eliminating exposure bias and metric mismatch in multiple object tracking. In: CVPR. pp. 4639–4648 (2019)
36. Manen, S., Gygli, M., Dai, D., Van Gool, L.: Pathtrack: Fast trajectory annotation with path supervision. In: ICCV. pp. 290–299 (2017)

37. Milan, A., Leal-Taixé, L., Reid, I., Roth, S., Schindler, K.: Mot16: A benchmark for multi-object tracking. arXiv preprint arXiv:1603.00831 (2016)
38. Milan, A., Leal-Taixé, L., Schindler, K., Reid, I.: Joint tracking and segmentation of multiple targets. In: CVPR. pp. 5397–5406 (2015)
39. Oh, S., Hoogs, A., Perera, A., Cuntoor, N., Chen, C.C., Lee, J.T., Mukherjee, S., Aggarwal, J., Lee, H., Davis, L., et al.: A large-scale benchmark dataset for event recognition in surveillance video. In: CVPR. pp. 3153–3160 (2011)
40. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: CVPR. pp. 779–788 (2016)
41. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: NeurIPS. pp. 91–99 (2015)
42. Ristani, E., Tomasi, C.: Features for multi-target multi-camera tracking and re-identification. In: CVPR. pp. 6036–6046 (2018)
43. Sadeghian, A., Alahi, A., Savarese, S.: Tracking the untrackable: Learning to track multiple cues with long-term dependencies. In: ICCV. pp. 300–311 (2017)
44. Schuster, S., Vernaza, P., Choi, W., Chandraker, M.: Deep network flow for multi-object tracking. In: CVPR. pp. 6951–6960 (2017)
45. Si, J., Zhang, H., Li, C.G., Kuen, J., Kong, X., Kot, A.C., Wang, G.: Dual attention matching network for context-aware feature sequence based person re-identification. In: CVPR. pp. 5363–5372 (2018)
46. Song, C., Huang, Y., Ouyang, W., Wang, L.: Mask-guided contrastive attention model for person re-identification. In: CVPR. pp. 1179–1188 (2018)
47. Su, C., Li, J., Zhang, S., Xing, J., Gao, W., Tian, Q.: Pose-driven deep convolutional model for person re-identification. In: ICCV. pp. 3960–3969 (2017)
48. Suh, Y., Wang, J., Tang, S., Mei, T., Mu Lee, K.: Part-aligned bilinear representations for person re-identification. In: ECCV. pp. 402–419 (2018)
49. Tang, S., Andriluka, M., Andres, B., Schiele, B.: Multiple people tracking by lifted multicut and person re-identification. In: CVPR. pp. 3539–3548 (2017)
50. Vondrick, C., Ramanan, D.: Video annotation and tracking with active learning. In: NeurIPS. pp. 28–36 (2011)
51. Wang, S., Fowlkes, C.C.: Learning optimal parameters for multi-target tracking with contextual interactions. IJCV **122**(3), 484–501 (2017)
52. Wang, Z., Zheng, L., Liu, Y., Wang, S.: Towards real-time multi-object tracking. arXiv preprint arXiv:1909.12605 (2019)
53. Wojke, N., Bewley, A., Paulus, D.: Simple online and realtime tracking with a deep association metric. In: ICIP. pp. 3645–3649 (2017)
54. Xiang, Y., Alahi, A., Savarese, S.: Learning to track: Online multi-object tracking by decision making. In: ICCV. pp. 4705–4713 (2015)
55. Xu, J., Cao, Y., Zhang, Z., Hu, H.: Spatial-temporal relation networks for multi-object tracking. In: ICCV. pp. 3988–3998 (2019)
56. Xu, Y., Osep, A., Ban, Y., Horaud, R., Leal-Taixé, L., Alameda-Pineda, X.: How to train your deep multi-object tracker. In: CVPR (2020)
57. Yang, F., Choi, W., Lin, Y.: Exploit all the layers: Fast and accurate cnn object detector with scale dependent pooling and cascaded rejection classifiers. In: CVPR. pp. 2129–2137 (2016)
58. Ye, M., Shen, J., Lin, G., Xiang, T., Shao, L., Hoi, S.C.: Deep learning for person re-identification: A survey and outlook. arXiv preprint arXiv:2001.04193 (2020)
59. Yu, F., Li, W., Li, Q., Liu, Y., Shi, X., Yan, J.: Poi: Multiple object tracking with high performance detection and appearance feature. In: ECCV. pp. 36–42 (2016)
60. Yu, Q., Medioni, G., Cohen, I.: Multiple target tracking using spatio-temporal markov chain monte carlo data association. In: CVPR. pp. 1–8 (2007)

61. Zamir, A.R., Dehghan, A., Shah, M.: Gmcp-tracker: Global multi-object tracking using generalized minimum clique graphs. In: ECCV. pp. 343–356 (2012)
62. Zhang, L., Li, Y., Nevatia, R.: Global data association for multi-object tracking using network flows. In: CVPR. pp. 1–8 (2008)
63. Zhang, S., Gong, Y., Huang, J.B., Lim, J., Wang, J., Ahuja, N., Yang, M.H.: Tracking persons-of-interest via adaptive discriminative features. In: ECCV. pp. 415–433 (2016)
64. Zhang, Y., Wang, C., Wang, X., Zeng, W., Liu, W.: A simple baseline for multi-object tracking. arXiv preprint arXiv:2004.01888 (2020)
65. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: ICCV (2015)
66. Zheng, Z., Yang, X., Yu, Z., Zheng, L., Yang, Y., Kautz, J.: Joint discriminative and generative learning for person re-identification. In: CVPR. pp. 2138–2147 (2019)
67. Zhou, K., Xiang, T.: Torchreid: A library for deep learning person re-identification in pytorch. arXiv preprint arXiv:1910.10093 (2019)
68. Zhou, K., Yang, Y., Cavallaro, A., Xiang, T.: Omni-scale feature learning for person re-identification. In: ICCV. pp. 3702–3712 (2019)
69. Zhou, X., Koltun, V., Krähenbühl, P.: Tracking objects as points. arXiv:2004.01177 (2020)
70. Zhu, J., Yang, H., Liu, N., Kim, M., Zhang, W., Yang, M.H.: Online multi-object tracking with dual matching attention networks. In: ECCV. pp. 366–382 (2018)