
Self-Supervised Multi-Object Tracking with Cross-Input Consistency

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we propose a self-supervised learning procedure for training a robust
2 multi-object tracking (MOT) model given only unlabeled video. While several
3 self-supervisory learning signals have been proposed in prior work on single-object
4 tracking, such as color propagation and cycle-consistency, these signals are not
5 effective for training RNN models, which are needed to achieve accurate MOT:
6 they yield degenerate models that, for instance, always match new detections to
7 tracks with the closest initial detections. We propose a novel self-supervisory
8 signal that we call cross-input consistency: we construct two distinct inputs for
9 the same sequence of video, by hiding different information about the sequence in
10 each input. We then compute tracks in that sequence by applying an RNN model
11 independently on each input, and train the model to produce consistent tracks
12 across the two inputs. We evaluate our unsupervised method on the MOT17 and
13 KITTI benchmarks — remarkably, we find that, despite training only on a corpus of
14 unlabeled video, our unsupervised approach *outperforms* four supervised methods
15 published in the last 1–2 years, including Tractor++ [1], FAMNet [5], GSM [18],
16 and mmMOT [29].

17 1 Introduction

18 Multi-object trackers identify all instances of a particular object type in video, and track each instance
19 through the segment of video in which it is visible in the camera frame. Annotating training data
20 for multi-object tracking is tedious and costly; for example, annotation of pedestrian tracks in just
21 six minutes of video in the training set of the MOT15 Challenge [14] requires an estimated 22
22 hours [20] of human labeling time using LabelMe [28]. While unsupervised, heuristic detect-to-track
23 methods [2, 4] have been proposed that group detections into tracks by estimating motion using a
24 combination of spatial and visual cues, these methods suffer low-accuracy in scenarios with frequent
25 occlusion where heuristics are insufficient.

26 Recent work has proposed applying self-supervised learning for training *single*-object tracking
27 models on unlabeled video [24, 25]. These approaches train a model to propagate instance labels
28 from a reference frame through the rest of a video sequence. In contrast to work on self-supervised
29 representation learning from video, these fully unsupervised approaches do not require fine-tuning to
30 apply the model for single-object tracking.

31 However, a significant limitation in prior work is that the model independently compares pairs of
32 frames at a time. In multi-object tracking, a key challenge is robustly re-localizing tracks across
33 potentially long occlusions, especially when an object instance is occluded by other instances of the
34 same object type. Pairwise frame comparisons are thus insufficient for high-accuracy multi-object
35 tracking; instead, learning recurrent features that encode the history of a track is crucial for enabling
36 robust re-localization. However, extending prior work to learn RNN parameters is challenging. For

37 example, Wang et al. [25] propose training using forward-backward consistency: from a patch in an
38 initial frame, after tracking forwards through video and then backwards to return to the initial frame,
39 the final patch should align with the original patch. Training an RNN in this way would be ineffective
40 as the RNN could simply memorize the features of the original patch.

41 To address this challenge, we propose a novel self-supervised learning method, *cross-input con-*
42 *sistency*. We first compute object detections in each frame of unlabeled video (like unsupervised,
43 heuristic detect-to-track methods, we assume that a robust detector is available). Then, we derive a
44 learning signal from the unlabeled video by sampling a short sequence of video, constructing two
45 input variations of that sequence, and training the tracker to produce consistent tracking outputs when
46 applied independently on each of the two inputs. We propose two alternative input-hiding schemes for
47 computing the input variations: occlusion-based hiding and visual-spatial hiding. Occlusion-based
48 hiding eliminates information about object detections in random intermediate subsequences of frames
49 to simulate occlusion incidents; thus, it constructs two inputs by eliminating different subsequences.
50 Visual-spatial hiding applies the tracker once when only observing spatial inputs (bounding box
51 coordinates), and once when only observing visual inputs (pixel values inside detection boxes). After
52 sampling a sequence of video and computing the two input variations under the chosen input-hiding
53 scheme, we apply the tracker model independently on each input, and back-propagate a learning
54 signal that measures the consistency between tracks computed across the two inputs. To attain high
55 consistency, the model must accurately group detections that correspond to the same object: if the
56 model were to instead arbitrarily group detections into tracks, then variations in the inputs would
57 cause the tracker to produce inconsistent outputs.

58 To implement cross-input consistency, we adapt a now standard RNN model and tracker architecture
59 from prior work [12]: the tracker processes each frame in sequence by matching detections in the
60 current frame with tracks computed up to the previous frame. In prior work, this model is trained
61 under a supervised procedure: they sample a video sequence $\langle I_0, \dots, I_n \rangle$ and a track t in that
62 sequence, and apply the tracker on t over the sequence. On each frame I_j , the RNN outputs a
63 probability distribution indicating the likelihood that the prefix of a track t up to I_j matches with
64 each detection in I_j . Prior work back-propagates the label (i.e., the correct detection of t in I_j) under
65 cross entropy loss.

66 In contrast, under our method, on each training iteration, we propose to sample a sequence
67 $\langle I_0, \dots, I_n \rangle$ from a corpus of unlabeled video, and apply the RNN model to compute a transi-
68 tion matrix that specifies the probability that each detection in I_0 (rows) matches with each detection
69 in I_n (columns). We select the sequence length n so that most (but not all) objects in I_0 are still
70 visible in I_n . Then, when applying the tracker on two input variations extracted from the sequence,
71 we obtain two transition matrices (one for each input). We compute the dot-product similarity to
72 measure the consistency between these matrices, and back-propagate the negative similarity as a loss
73 function.

74 We evaluate our approach on the MOT17 and KITTI benchmarks against 8 state-of-the-art methods,
75 including both unsupervised and supervised methods. We train our tracker model using cross-input
76 consistency over a corpus of unlabeled video, which can be cheaply obtained. Like other unsupervised
77 methods, we use an object detector trained on image-level bounding box annotations in COCO [17],
78 but do not use any expensive video-level annotations. We find that our approach improves both IDF1
79 and MOTA accuracy over the unsupervised baselines by 14% to 18%. Moreover, remarkably, our
80 fully unsupervised approach *outperforms* five of the seven supervised methods we compared, even
81 though these methods train on expensive video-level bounding box and track annotations.

82 2 Related Work

83 Self-supervised learning over video has been studied extensively in many contexts. Most work focuses
84 on learning representations of video that can be applied through fine-tuning for tasks such as activity
85 recognition, image classification, and object detection [6, 7, 9, 15, 23, 26]. More closely related to
86 our work, several recent approaches have proposed leveraging widely available unlabeled video to
87 directly train *single-object* tracking models, without needing fine-tuning [13, 16]. Vondrick et al. [24]
88 train a model to colorize gray-scale video by propagating colors from a colored reference frame. The
89 model is then applied to track objects at inference time by propagating instance IDs instead of colors.
90 Wang et al. [25] train a model to capture correspondence by applying a cycle-consistent loss: from

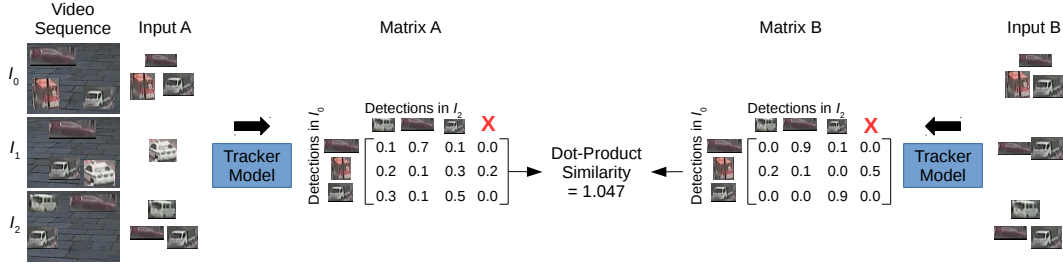


Figure 1: Overview of cross-input consistency. An input-hiding scheme produces two inputs, A and B, from one video sequence; these inputs contain identical information about objects detected in the first and last frames of the sequence, but vary in intermediate frames. We apply the tracker model on each input to derive two transition matrices that match detections between the first and last frames to represent tracker outputs. We then back-propagate a similarity score between the matrices that encourages the model to produce consistent outputs across both inputs.

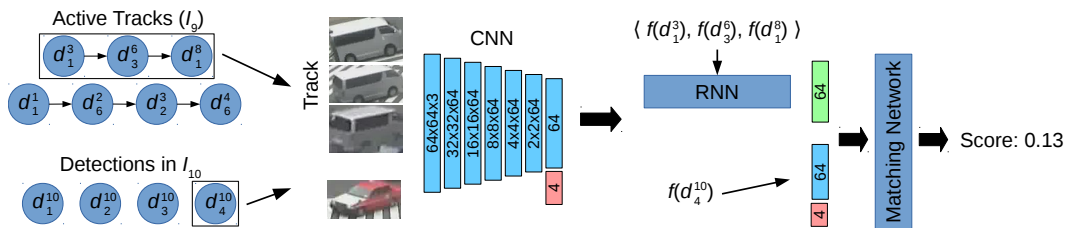


Figure 2: Tracker model architecture. The model scores the likelihood that each detection in the current frame matches with each track.

91 a patch in an initial frame, after tracking forwards through video and then backwards to return to
 92 the initial frame, the final patch should align with the original patch. As we discussed in Section 1,
 93 self-supervisory signals used in prior work such as color propagation and cycle-consistency are not
 94 effective for training RNN models, which are needed to achieve accurate MOT.

95 Our work is also related to unsupervised, heuristic detect-to-track multi-object tracking methods
 96 such as SORT [2] and V-IOU [4]. These methods group detections across different frames using a
 97 combination of heuristic spatial cues (e.g., Kalman filter over bounding box coordinates) and visual
 98 cues (e.g., optical flow) to track objects. Like our approach, these methods assume that a robust
 99 detector is available; however, because they rely on heuristics to group detections into tracks, they
 100 suffer low-accuracy in scenarios with frequent occlusion.

101 Multi-object tracking has been studied extensively in supervised settings, where methods are trained
 102 on video-level bounding box and track annotations [1, 5, 8, 12, 18, 30]. However, such annotations
 103 are expensive to hand-label, and so these methods are costly to extend to new types of video.

104 Other work explores using unsupervised and self-supervised learning to further improve the perfor-
 105 mance of fully supervised methods. SimpleReID [11] proposes improving the performance of one
 106 supervised method, CenterTrack [30], by training a re-identification model through unsupervised
 107 learning. However, while the model can in principle be trained only on image-level annotations
 108 through hallucinated motion techniques, their SimpleReID+CenterTrack tracking method depends
 109 on expensive video-level annotations to attain high-accuracy. In contrast, our method achieves
 110 competitive results without any video-level supervision.

111 3 Cross-Input Consistency

112 In our novel cross-input consistency method, we derive a learning signal for training an RNN tracker
 113 model through a three-step process. We assume that a corpus of unlabeled video is available, along
 114 with an object detector for the object category of interest. During pre-processing, we apply the
 115 detector on each frame of unlabeled video to compute object detections. Then, during training, we

116 first repeatedly randomly sample a video sequence $\langle I_0, \dots, I_n \rangle$, where each I_k is a video frame. Let
 117 D_k by the detections automatically computed in I_k by the detector, and let $d_i^k = (im, x, y, w, h)$ be
 118 a detection in D_k , where (x, y, w, h) are the 4D spatial coordinates (center point and lengths) of the
 119 detection bounding box, and im is the window of I_k corresponding to that box. We apply an input-
 120 hiding scheme to select two input variations $A(D), B(D)$ for the video segment, where each variation
 121 is a modified sequence of detections in the frames, $A(D) = \langle D_0^A, \dots, D_n^A \rangle, B(D) = \langle D_0^B, \dots, D_n^B \rangle$.
 122 For example, some detections may be removed entirely, while others may be partially hidden. Second,
 123 we apply the tracker model independently on each input variation to derive two probabilistic tracking
 124 outputs (one per input), represented as transition matrices. Third, we compare the transition matrices
 125 with dot-product similarity to update the RNN parameters.

126 Figure 1 summarizes our approach.

127 Below, we first introduce the model architecture that we adapt from prior work in Section 3.1. We
 128 then detail our novel training procedure, including the computation of the transition matrices and
 129 dot-product loss, in Section 3.2. Finally, we propose two input-hiding schemes for selecting the input
 130 variations required by our approach in Section 4.

131 3.1 Background: Tracker Model

132 We adopt a tracker model that is similar to prior work [12]. We summarize the architecture in Figure
 133 2. Given a video sequence $\langle I_0, \dots, I_n \rangle$, and sets of detections $D_k = \{d_1^k, \dots, d_{m_k}^k\}$ detected in each
 134 frame I_k , to initialize the tracking process, we create a length-1 track $t_i = \langle d_i^0 \rangle$ for each detection d_i^0
 135 in the first video frame I_0 . When processing subsequent frames, we will match the new detections
 136 with existing tracks, extending existing tracks if there is a match and initializing new tracks otherwise.
 137 Specifically, on each subsequent frame I_k , the model outputs a probability $p_{i,j}$ that each track t_i
 138 corresponds to each detection $d_j^k \in D_k$. At inference time, we formulate the problem of matching
 139 tracks with detections in I_k as a bipartite matching problem, where the cost of matching t_i with
 140 d_j^k is $1 - p_{i,j}$. We solve this problem and compute a minimum-cost matching using the Hungarian
 141 algorithm; for each pair (t_i, d_j^k) in the matching, we append d_j^k to t_i . For each detection in I_k that no
 142 track matches to, we create a new track for that detection.

143 The model consists of a CNN, RNN, and matcher network. Together, these components score the
 144 likelihood that the i th track, $t_i = \langle d_1, \dots, d_m \rangle$, matches with the j th detection in I_k , d_j^k . We first
 145 apply the CNN to derive detection-level features. Given a detection $d = (im, x, y, w, h)$, the CNN
 146 inputs im resized to 64×64 , and consists of 6 strided convolutional layers, with ReLU activation in
 147 the first 5 layers and linear activation in the last layer. It outputs a 64-vector, which we concatenate
 148 with the 4D spatial coordinates to derive a 68-vector detection representation $f(d)$. Then, we compute
 149 track-level features $f(t_i)$ by applying the RNN (an LSTM with 64 hidden states) over the sequence
 150 of detection-level features of detections in the track, $\langle f(d_1), \dots, f(d_m) \rangle$. We use the output of the
 151 RNN on the last timestep as the track-level features $f(t_i)$. Finally, we apply a matching network to
 152 score the likelihood that t_i matches d_j^k . The matching network inputs the concatenation of $f(t_i)$ and
 153 $f(d_j^k)$, applies four fully-connected layers, and outputs a match score.

154 3.2 Training Procedure

155 We develop a novel self-supervised learning method for training the model parameters on unlabeled
 156 video. During training, we repeatedly sample sequences of video $\langle I_0, \dots, I_n \rangle$. We apply one of
 157 two input-hiding schemes, which we will detail in the following section, to extract two distinct
 158 input variations $A(D)$ and $B(D)$ from a sampled video sequence, where each input is a sequence
 159 of detections. We then apply the tracker independently on $A(D)$ and $B(D)$ to derive two tracking
 160 outputs for the same video sequence. In cross-input consistency, we train the model by enforcing
 161 similarity between these two outputs.

162 To represent tracker outputs, we compute an $|D_0| \times |D_n| + 1$ transition matrix $M^{(0,n)}$, where
 163 $M_{i,j}^{(0,n)}, j < |D_n|$ is the probability that the track t_i matches d_j^n . We use the last column to represent
 164 tracks that are no longer visible in I_n , i.e., $M_{i,|D_n|}^{(0,n)}$ is the probability that the track t_i has exited the
 165 camera frame. When applying the model over video sequences during training, we update tracks with
 166 new detections based on the scores output by the model on intermediate frames, but do not create

167 additional tracks on frames after I_0 ; thus, each track t_i corresponds directly to a detection d_i^0 in I_0
 168 (i.e., $t_i = \langle d_i^0, \dots \rangle$). Thus, we can also think of $M^{(0,n)}$ as the probability that a detection in the first
 169 frame d_i^0 matches a detection in the last frame d_j^n .

170 Applying the tracker on both input variations yields two transition matrices $A^{(0,n)}$ and $B^{(0,n)}$ that
 171 match objects detected in I_0 with those in I_n . We train the model (CNN, RNN, and matching network)
 172 to maximize the dot-product similarity between these matrices. In addition to pushing the model to
 173 produce consistent outputs across both inputs, we also design our training method so that the model
 174 cannot attain a high similarity score by, for example, saying that all objects visible in I_0 are no longer
 175 visible in I_n .

176 In our method, it is important that the training sequence length n be chosen so that, in most sequences,
 177 most (but not all) objects in I_0 are still visible in I_n , but that objects nevertheless move non-trivially
 178 during the sequence (so that the tracking task is not too easy). In general, we find that setting n to
 179 one-half of the average time that objects linger in the camera frame works well.

180 Below, we detail our method to compute transition matrices, and discuss dot-product similarity loss.

181 **Transition Matrix.** We propose computing a transition matrix $M^{(0,k)}$ on each frame I_k to represent
 182 the tracker outputs, where $M_{i,j}^{(0,k)}$ is the probability that the track t_i matches the detection d_j^k . On
 183 intermediate frames, we apply the Hungarian method on $M^{(0,k)}$ to match detections in D_k with tracks,
 184 updating each track with the matched detection (if any). On the last frame I_n , we use the $M^{(0,n)}$
 185 matrix produced under different inputs (denoted $A^{(0,n)}$ and $B^{(0,n)}$) to compute and back-propagate
 186 a consistency score. Because we do not create new tracks after I_0 during training, $M_{i,j}^{(0,n)}$ is the
 187 likelihood that d_i^0 and d_j^n match (since t_i begins with d_i^0).

188 We first construct a score matrix $S^{(0,k)}$, by computing $S_{i,j}^{(0,k)}$ as the score (any real number) output
 189 by the tracker model given the track t_i and detection d_j^k . We then transform the score matrix into a
 190 probability matrix to derive $M^{(0,k)}$. We could simply compute $M^{(0,k)}$ by taking softmax along rows
 191 in $S^{(0,k)}$. However, computing the transition matrix in this way would allow the tracker to cheat and
 192 maximize similarity between $A^{(0,n)}$ and $B^{(0,n)}$ by simply matching all detections in I_0 to a single
 193 detection $d_j^n \in D_n$. Indeed, we find that in practice this yields degenerate models.

194 Thus, instead, we compute $M^{(0,k),\text{row}}$ and $M^{(0,k),\text{col}}$ by applying softmax along rows and columns,
 195 respectively, and compute $M^{(0,k)} = \min(M^{(0,k),\text{row}}, M^{(0,k),\text{col}})$:

$$M_{i,j}^{\text{row}} = \frac{\exp(S_{i,j})}{\sum_k \exp(S_{i,k})} \quad M_{i,j}^{\text{col}} = \frac{\exp(S_{i,j})}{\sum_k \exp(S_{k,j})} \quad (1)$$

$$M_{i,j} = \min(M_{i,j}^{\text{row}}, M_{i,j}^{\text{col}}) \quad (2)$$

196 This produces a transition matrix $M^{(0,k)}$ that is almost doubly stochastic: rows and columns sum
 197 to at most 1, but not necessarily exactly 1. The operation ensures that the model must match each
 198 detection in I_0 to unique detections in I_n to maximize the consistency score between $A^{(0,n)}$ and
 199 $B^{(0,n)}$: if two detections in I_0 are matched to the same detection in I_n , then the columnar softmax
 200 would reduce those probabilities in the corresponding matrix to at most 0.5, thereby reducing any
 201 dot-products involving those rows.

Dot-Product Similarity. We train the RNN tracker by computing two transition matrices $A^{(0,n)}$
 and $B^{(0,n)}$ over different input variations, and then back-propagating a loss that measures the
 inconsistency between the matrices. In particular, we use the dot-product to measure the similarity of
 corresponding rows in the matrices. We define the loss as:

$$L = - \sum_i \log \sum_j A_{i,j}^{(0,n)} B_{i,j}^{(0,n)}$$

202 Here, L is computed by taking the logarithm of the dot product of corresponding rows in $A^{(0,n)}$
 203 and $B^{(0,n)}$, averaged across rows. Note that this is equivalent to the cross-entropy loss between the
 204 diagonal matrix and the matrix product of $A^{(0,n)}$ and the transpose of $B^{(0,n)}$.

205 This loss function has several desirable properties. First, the dot-product is maximized when the
 206 matrices computed based on different inputs are most similar. This pushes the matching model
 207 to learn reasonable visual and spatial tracking constraints, because arbitrarily matching tracks to
 208 detections will lead to dissimilarity. Second, the dot-product pushes each matrix to be almost doubly
 209 stochastic rather than leaving some rows and columns summing to much less than 1. The model can
 210 only produce doubly stochastic $A^{(0,n)}$ and $B^{(0,n)}$ matrices by finding unique detections in I_n for
 211 each detection in I_0 .

212 4 Input-Hiding Schemes

213 In this section, we detail two alternative input-hiding schemes for selecting the two input variations,
 214 denoted $A(D) = \langle D_0^A, \dots, D_n^A \rangle$ and $B(D) = \langle D_0^B, \dots, D_n^B \rangle$. Recall that D is the original set of
 215 all objects detected in a video sequence $\langle I_0, \dots, I_n \rangle$.

216 4.1 Occlusion-based Hiding

217 At a high level, occlusion-based hiding produces the variations $A(D)$ and $B(D)$ by simulating
 218 random occlusion incidents where all detections in occluded frames are eliminated from the input, i.e.,
 219 if I_k is occluded for $A(D)$, then D_k^A is empty. We only occlude intermediate frames I_k , $0 < k < n$,
 220 so that the transition matrices still compare detections in I_0 with those in I_n . When processing
 221 an occluded I_k , the tracker is forced to match all tracks to the absent column in that frame, and
 222 re-localize the tracks after the occlusion.

223 We first introduce two schemes that do not work in isolation, and then show that we can combine
 224 these schemes to produce input variations that result in effective training.

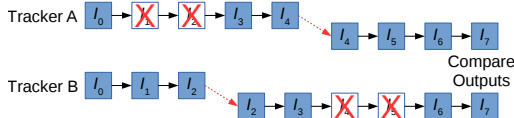
225 **Only-Occlusion.** For each training sequence $\langle I_0, \dots, I_n \rangle$, Only-Occlusion randomly selects four
 226 indexes $0 < k_1 \leq k_2 < k_3 \leq k_4 < n$ to construct two disjoint frame subsequences $\langle I_{k_1}, \dots, I_{k_2} \rangle$
 227 and $\langle I_{k_3}, \dots, I_{k_4} \rangle$. In $A(D)$, we occlude each frame I_k such that $k_1 \leq k \leq k_2$, and in $B(D)$, we
 228 occlude I_k if $k_3 \leq k \leq k_4$.

229 When training under Only-Occlusion, the tracker is forced to leverage track features computed
 230 through the RNN to re-localize tracks after each simulated occlusion ends. Learning to merely
 231 compare detection features across consecutive frames would yield low accuracy since features in
 232 occluded frames are not observed. Furthermore, because one tracker observes the detections and the
 233 other tracker does not, the model must make similar tracking decisions when re-localizing across
 234 occluded frames as it does when observing detections in each frame.

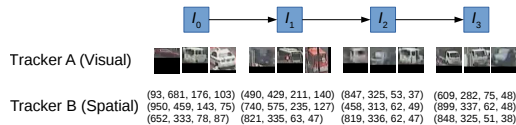
235 However, in practice, Only-Occlusion produces a model that simply memorizes the detections in I_0 ,
 236 and computes $A^{(0,n)}$ and $B^{(0,n)}$ by comparing the detections in I_n against memorized detections.
 237 This strategy yields high similarity because it is unaffected by simulated occlusions in intermediate
 238 frames. Thus, we must prevent the propagation of features directly from I_0 to I_n to make this scheme
 239 effective.

240 **RNN Hand-off.** RNN Hand-off prevents simple memorization by cutting off the propagation of
 241 RNN features through the application of two separate RNN executions. We select two indexes
 242 $0 < k_5, k_6 < n$. Instead of computing $A^{(0,n)}$ directly, we first apply the tracker on the frame
 243 sequence $\langle I_0, \dots, I_{k_5} \rangle$ to derive a transition matrix $A^{(0,k_5)}$ that matches detections in I_0 with
 244 detections in I_{k_5} . We then independently apply the tracker on $\langle I_{k_5}, \dots, I_n \rangle$ to derive another matrix
 245 $A^{(k_5,n)}$ that matches detections in I_{k_5} with detections in I_n . We combine these matrices through
 246 the matrix product to compute $A^{(0,n)}$: we compute $A^{(0,n)} = A^{(0,k_5)} A^{(k_5,n)}$. Similarly, we compute
 247 $B^{(0,n)} = B^{(0,k_6)} B^{(k_6,n)}$.

248 This scheme forces the tracker to find the same unique detection in I_{k_5} (and I_{k_6}) for two detections
 249 of the same object in I_0 and I_n in order to maximize similarity between the matrix products. On
 250 the other hand, though, a tracker that learns to match tracks to detections by comparing only the
 251 detection features in consecutive frames will exhibit high similarity between $A^{(0,n)}$ and $B^{(0,n)}$ under
 252 this scheme.



(a) Occlusion-based hiding produces input variations with different subsequences of occluded frames where all detections are hidden from the tracker. It also independently applies the tracker before and after a hand-off frame (I_4 and I_2), and merges the outputs through the matrix product.



(b) Visual-spatial hiding. One input includes only visual information about the detections, while the other includes only spatial bounding box coordinates.

Method	IDF1	MOTA	MT	ML	FP	FN	IDSW	Frag
Occlusion (ours)	52.4	56.7	28.7	16.9	1K	5K	136	172
Visual-Spatial (ours)	57.3	60.2	27.0	20.1	734	5K	83	118
Spatial-Only (ours)	56.5	57.8	27.1	18.3	1K	5K	144	169

Table 1: Ablation study on the MOT17 training set.

253 **Combined Scheme.** Only-Occlusion and RNN Hand-off have opposite advantages and drawbacks.
 254 Thus, we combine these in our occlusion-based hiding scheme. We first select the two sequences
 255 for simulated occlusion, $\langle I_{k_1}, \dots, I_{k_2} \rangle$ and $\langle I_{k_3}, \dots, I_{k_4} \rangle$. Then, we randomly pick k_5 and k_6 such
 256 that $k_3 \leq k_5 \leq k_4$ and $k_1 \leq k_6 \leq k_2$, i.e., the hand-off for one tracker occurs when the other tracker
 257 observes a simulated occlusion. We summarize the scheme in Figure 3a.

258 Under this scheme, neither memorizing features in I_0 nor comparing detections solely in a pairwise
 259 frame-by-frame manner is an effective tracking strategy. Instead, the tracker must learn to leverage
 260 RNN features for re-localizing across simulated occlusion, while still ensuring the tracking decisions
 261 reflect intermediate outputs.

262 4.2 Visual-Spatial Hiding

263 Under visual-spatial hiding, we apply one tracker instance that observes only visual features and one
 264 tracker instance that observes only spatial features: in $A(D)$, we set $x = 0, y = 0, w = 0, h = 0$ for
 265 all detections, and in $B(D)$, we set $im = 0$. By encouraging similarity in the tracking outputs, each
 266 tracker learns to leverage its input features and track objects as robustly as possible.

267 To prevent the visual tracker from examining the background to extract an approximation of the
 268 spatial features, we do not use a recurrent unit for the visual tracker; instead, its matching network
 269 inputs visual features for detections in I_0 and for detections in I_n , and scores each pair of detections
 270 in $A^{(0,n)}$ without observing intermediate frames. We compute $B^{(0,n)}$ through the spatial tracker,
 271 which inputs 4D spatial coordinates for each detection, and processes frames with the matching
 272 network and recurrent unit. Figure 3b illustrates the training procedure.

273 5 Evaluation

274 We compare our method and nine baselines on the MOT17 [21] and KITTI [10] benchmarks.

275 **Baselines.** We compare with two unsupervised methods (SORT [2] and V-IOU [4]) and seven
 276 fully supervised methods (Tracktor++ [1], MHT-BLSTM [12], FAMNet [5], LSST [8], GSM [18],
 277 mmMOT [29], and CenterTrack [30]). Like our approach, SORT and V-IOU require an object
 278 detector, but do not train on any video-level bounding box and track annotations in the MOT17
 279 and KITTI training sets. The fully supervised methods train on video-level annotations; Tracktor++
 280 incorporates a core component that uses only the detector regression network, but requires video-level
 281 annotations for training a re-identification network. Results for 8 baselines are available on MOT17,
 282 and results for 4 baselines are available on KITTI.

283 **Dataset.** MOT17 [21] consists of 14 video sequences of pedestrians in a wide range of contexts,
 284 including a moving camera inside a shopping mall and a fixed, elevated view of an outdoor plaza.

Method	IDF1	MOTA	MT	ML	FP	FN	IDSW	Frag
Visual-Spatial (ours)	58.3	56.8	538	880	11,567	230,645	1,320	2,061
SORT [2]	39.8	43.1	295	997	28,398	287,582	4,852	7,127
IOU [3]	39.4	45.5	369	953	19,993	281,643	5,988	7,404
Tracktor++ [1]	52.3	53.5	459	861	12,201	248,047	2,072	4,611
MHT-BLSTM [12]	51.9	47.5	429	981	25,981	268,402	2,069	3,124
FAMNet [5]	48.7	52.0	450	787	14,138	253,616	3,072	5,318
LSST [8]	62.3	54.7	480	944	26,091	228,434	1,243	3,726
GSM [18]	57.8	56.4	523	813	14,379	230,174	1,485	2,763
CenterTrack [30]	59.6	61.5	621	752	14,076	200,672	2,583	4,965

Table 2: Performance on the MOT17 test set. We show unsupervised methods, including our approach, at the top, and methods that require video-level annotations at the bottom.

Method	HOTA	DetA	AssA	DetRe	DetPr	AssRe	AssPr	LocA	MOTA
Visual-Spatial (ours)	62.5	61.1	65.3	67.7	73.8	69.1	83.1	80.3	71.6
SORT [2]	42.5	44.0	41.3	47.3	73.9	42.8	83.0	80.8	53.1
FAMNet [5]	52.6	61.0	45.5	64.4	78.7	48.7	77.4	81.5	75.9
mmMOT [29]	62.1	72.3	54.0	76.2	84.9	59.0	82.4	86.6	83.2
CenterTrack [30]	73.0	75.6	71.2	80.1	84.6	73.8	89.0	86.5	88.8

Table 3: Performance on the KITTI test set (tracking cars). We show unsupervised methods, including our approach, at the top, and methods that require video-level annotations at the bottom.

285 The dataset is split into 7 training sequences and 7 test sequences; each split includes approximately
286 11 minutes of video. KITTI [10] consists of 48 video sequences captured from vehicle-mounted
287 cameras, split into 20 for training and 28 for testing, and the objective is to track cars.

288 **Training.** The supervised baselines train on video-level bounding box and track annotations provided
289 by MOT17 and KITTI. In contrast, our method trains only on a corpus of unlabeled video. Because
290 video-level annotations are expensive to label, our method requires substantially less annotation time,
291 and thus greatly reduces the effort needed to apply multi-object tracking on new datasets.

292 For MOT17, we collect unlabeled video from two sources: we use five hours of video from seven
293 YouTube walking tours, and all train and test sequences from the PathTrack dataset [20] (we do not
294 use the PathTrack ground truth annotations). For KITTI, we use both the 46 minutes of video in the
295 KITTI dataset together with 7 hours of video from Berkeley DeepDrive [27]. We train our tracker
296 model on an NVIDIA Tesla V100 GPU; training time varies between 4 and 24 hours depending on
297 the input-hiding scheme. During training, we randomly select sequence lengths n between 4 and 16,
298 and apply stochastic gradient descent one sequence at a time. We apply the Adam optimizer with
299 learning rate 0.0001, decaying to 0.00001 after plateau.

300 In contrast to MOT17, KITTI does not provide object detections for use by tracking methods. We
301 extract detections from video using a YOLOv5 model trained on COCO. On MOT17, we pre-process
302 detections with classification and regression following the method in Tracktor++ [1].

303 **Metrics.** We use Multi-Object Tracking Accuracy (MOTA) [21] and ID F1 Score (IDF1) [22] on
304 MOT17, and Higher Order Tracking Accuracy (HOTA) [19] for KITTI. Broadly, these metrics
305 measure the accuracy of inferred tracks against ground truth tracks, and penalize both when an
306 inferred track contains a detection that doesn’t match to some ground truth detection (or vice versa),
307 and when a ground truth track is split into two or more inferred tracks (or vice versa).

308 **Ablation Study.** We first compare occlusion-based hiding and visual-spatial hiding on the MOT17
309 training set in Table 1. Visual-spatial hiding yields higher performance on both MOTA and IDF1
310 — occlusion-based hiding performs poorly; objects are often visible in the video for only a short
311 duration, making it challenging to learn to re-localize objects over simulated occlusion since the
312 simulated occlusion must then also be short. Under Spatial-Only, we show results for visual-spatial

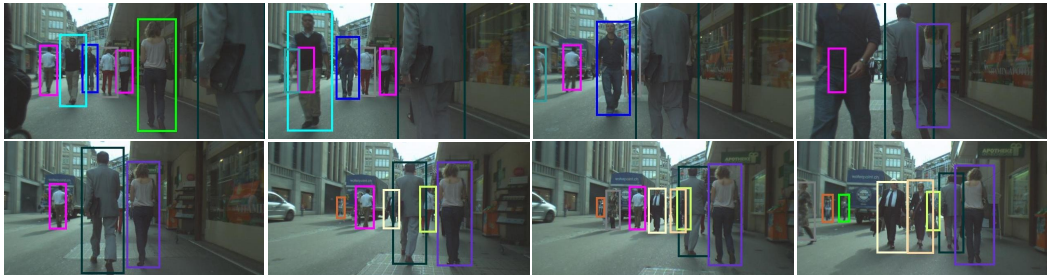


Figure 4: Output of Visual-Spatial on a portion of an MOT17 sequence. Our method tracks objects through several instances of occlusion.

313 hiding when inputting only the spatial bounding box coordinates of detections during inference (no
 314 image features).

315 **Quantitative Results.** Table 2 shows results on the MOT17 test set¹, and Table 3 shows results on
 316 the KITTI test set². Metrics are automatically computed by the challenge websites. Per the challenge
 317 policies, we only submit the best method, and thus show Visual-Spatial performance.

318 On MOT17, our approach substantially outperforms both of the unsupervised baselines. Moreover,
 319 despite training only on unlabeled video, our method *outperforms* Tracktor++ [1], MHT-BLSTM [12],
 320 FAMNet [5], and GSM [18], even though these baselines (all of which except MHT-BLSTM were
 321 published in the last 1–2 years) are supervised methods that train on expensive video-level annotations
 322 in the MOT17 training set. Our approach is also competitive with LSST [8]: our method improves
 323 in MOTA but yields lower IDF1. Nevertheless, CenterTrack [30] yields slightly higher accuracy on
 324 both metrics.

325 Similarly, on KITTI, our approach outperforms SORT [2], FAMNet [5], and mmMOT [29], but yields
 326 lower performance than CenterTrack [30].

327 **Qualitative Results.** We show qualitative results in Figure 4 and in the supplementary video.

328 6 Conclusion

329 In this paper, we have shown that a robust, fully unsupervised multi-object tracker can be trained
 330 through a novel self-supervisory learning signal, cross-input consistency, that enforces consistency in
 331 the tracking outputs across different input variations of one video sequence. Despite training only on
 332 unlabeled video, our approach outperforms four supervised trackers published in the last 1–2 years
 333 (Tracktor++ [1], FAMNet [5], GSM [18], and mmMOT [29]), which train on expensive video-level
 334 bounding box and track annotations.

335 **Social Impact.** By enabling a robust multi-object tracker to be trained given only unlabeled video,
 336 our work promises to greatly reduce the effort for users to apply multi-object tracking on new datasets
 337 without sacrificing accuracy. Thus, we believe that our novel self-supervised MOT method can open
 338 up new video analytics tasks that were previously too costly. This impact may be positive or negative
 339 depending on the nature of these tasks — however, in general, we believe that tasks with greater
 340 potential for negative impact such as surveillance and pedestrian tracking would not benefit from the
 341 reduction in annotation cost associated with our method.

¹These results are taken from <https://motchallenge.net/results/MOT17/>, where our method is denoted UNS20regress. Baselines are denoted SORT17, IOU17, Tracktor++, MHT_bLSTM, FAMNet, LSST17, GSM_Tracktor, and CTTrackPub.

²These results are taken from http://www.cvlibs.net/datasets/kitti/eval_tracking.php.

References

- 342
- 343 [1] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles.
344 In *IEEE International Conference on Computer Vision (ICCV)*, 2019.
- 345 [2] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Upcroft. Simple Online and
346 Realtime Tracking. In *IEEE International Conference on Image Processing (ICIP)*, 2016.
- 347 [3] Erik Bochinski, Volker Eiselein, and Thomas Sikora. High-Speed Tracking-by-Detection
348 Without Using Image Information. In *IEEE International Conference on Advanced Video and
349 Signal Based Surveillance (AVSS)*, 2017.
- 350 [4] Erik Bochinski, Tobias Senst, and Thomas Sikora. Extending IOU Based Multi-Object Tracking
351 by Visual Information. In *IEEE International Conference on Advanced Video and Signal Based
352 Surveillance (AVSS)*, 2018.
- 353 [5] Peng Chu and Haibin Ling. FAMNet: Joint Learning of Feature, Affinity and Multi-dimensional
354 Assignment for Online Multiple Object Tracking. In *IEEE International Conference on Com-
355 puter Vision (ICCV)*, 2019.
- 356 [6] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised Visual Representation
357 Learning by Context Prediction. In *IEEE International Conference on Computer Vision (ICCV)*,
358 2015.
- 359 [7] Carl Doersch and Andrew Zisserman. Multi-task Self-Supervised Visual Learning. In *IEEE
360 International Conference on Computer Vision (ICCV)*, 2017.
- 361 [8] Weitao Feng, Zhihao Hu, Wei Wu, Junjie Yan, and Wanli Ouyang. Multi-object tracking with
362 multiple cues and switcher-aware classification. *arXiv preprint arXiv:1901.06129*, 2019.
- 363 [9] Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould. Self-supervised video
364 representation learning with odd-one-out networks. In *Proceedings of the IEEE Conference on
365 Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- 366 [10] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving?
367 The KITTI Vision Benchmark Suite. In *IEEE Conference on Computer Vision and Pattern
368 Recognition (CVPR)*, 2012.
- 369 [11] Shyamgopal Karthik, Ameya Prabhu, and Vineet Gandhi. Simple unsupervised multi-object
370 tracking. *CoRR*, abs/2006.02609, 2020.
- 371 [12] Chanhok Kim, Fuxin Li, and James M. Rehg. Multi-Object Tracking with Neural Gating using
372 Bilinear LSTM. In *European Conference on Computer Vision (ECCV)*, 2018.
- 373 [13] Zihang Lai, Erika Lu, and Weidi Xie. Mast: A memory-augmented self-supervised tracker. In
374 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
375 June 2020.
- 376 [14] Laura Leal-Taixé, Anton Milan, Ian Reid, Stefan Roth, and Konrad Schindler. MOTChallenge
377 2015: Towards a Benchmark for Multi-Target Tracking. *arXiv preprint arXiv:1504.01942*,
378 2015.
- 379 [15] Hsin-Ying Lee, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Unsupervised represen-
380 tation learning by sorting sequences. In *Proceedings of the IEEE International Conference on
381 Computer Vision (ICCV)*, Oct 2017.
- 382 [16] Xueting Li, Sifei Liu, Shalini De Mello, Xiaolong Wang, Jan Kautz, and Ming-Hsuan Yang.
383 Joint-task self-supervised learning for temporal correspondence. In *Advances in Neural Infor-
384 mation Processing Systems*, volume 32, 2019.
- 385 [17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,
386 Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common Objects in Context. In
387 *European Conference on Computer Vision (ECCV)*, 2014.

- 388 [18] Qiankun Liu, Qi Chu, Bin Liu, and Nenghai Yu. GSM: Graph Similarity Model for Multi-Object
389 Tracking. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2020.
- 390 [19] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip Torr, Andreas Geiger, Laura Leal-
391 Taixé, and Bastian Leibe. HOTA: A Tigher Order Metric for Evaluating Multi-Object Tracking.
392 *International Journal of Computer Vision*, 129(2):548–578, 2021.
- 393 [20] Santiago Manen, Michael Gygli, Dengxin Dai, and Luc Van Gool. PathTrack: Fast Trajectory
394 Annotation with Path Supervision. In *IEEE International Conference on Computer Vision*
395 *(ICCV)*, 2017.
- 396 [21] Anton Milan, Laura Leal-Taixé, Ian Reid, Stefan Roth, and Konrad Schindler. MOT16: A
397 Benchmark for Multi-Object Tracking. *arXiv preprint arXiv:1603.00831*, 2016.
- 398 [22] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance
399 Measures and a Data Set for Multi-Target, Multi-Camera Tracking. In *European Conference on*
400 *Computer Vision (ECCV)*, 2016.
- 401 [23] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised Learning of Video
402 Representations using LSTMs. In *International Conference on Machine Learning (ICML)*,
403 2015.
- 404 [24] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy.
405 Tracking Emerges by Colorizing Videos. In *European Conference on Computer Vision (ECCV)*,
406 2018.
- 407 [25] Xiaolong Wang, Allan Jabri, and Alexei A Efros. Learning Correspondence from the Cycle-
408 Consistency of Time. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
409 2019.
- 410 [26] Donglai Wei, Joseph J. Lim, Andrew Zisserman, and William T. Freeman. Learning and using
411 the arrow of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
412 *Recognition (CVPR)*, June 2018.
- 413 [27] Huazhe Xu, Yang Gao, Fisher Yu, and Trevor Darrell. End-to-end learning of driving models
414 from large-scale video datasets. In *Proceedings of the IEEE Conference on Computer Vision*
415 *and Pattern Recognition (CVPR)*, July 2017.
- 416 [28] Jenny Yuen, Bryan Russell, Ce Liu, and Antonio Torralba. LabelMe video: Building a Video
417 Database with Human Annotations. In *IEEE International Conference on Computer Vision*
418 *(ICCV)*, 2009.
- 419 [29] Wenwei Zhang, Hui Zhou, Shuyang Sun, Zhe Wang, Jianping Shi, and Chen Change Loy.
420 Robust multi-modality multi-object tracking. In *Proceedings of the IEEE/CVF International*
421 *Conference on Computer Vision (ICCV)*, October 2019.
- 422 [30] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. Tracking Objects as Points. In *European*
423 *Conference on Computer Vision (ECCV)*, 2020.

424 Checklist

- 425 1. For all authors...
- 426 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
427 contributions and scope? [Yes] Our contribution is a self-supervised learning method
428 to train a multi-object tracker on unlabeled video, which we discuss in both the abstract
429 and the introduction.
- 430 (b) Did you describe the limitations of your work? [Yes] First, we have shown that while
431 our method outperforms four fully supervised MOT methods published in the last one
432 or two years, one fully supervised technique (CenterTrack) outperforms it on both
433 MOT17 and KITTI. Second, we discuss limitations in selection of the sequence length
434 n in Section 3.2.

- 435 (c) Did you discuss any potential negative societal impacts of your work? [Yes] In the
436 Conclusion.
- 437 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
438 them? [Yes]
- 439 2. If you are including theoretical results...
- 440 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
441 (b) Did you include complete proofs of all theoretical results? [N/A]
- 442 3. If you ran experiments...
- 443 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
444 mental results (either in the supplemental material or as a URL)? [Yes] We include the
445 code in the supplementary file, and the training and test data is available from MOT17,
446 KITTI, PathTrack, and BDD.
- 447 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
448 were chosen)? [Yes] See Section 5.
- 449 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
450 ments multiple times)? [No] We are only able to submit one final result to the MOT17
451 and KITTI challenge websites.
- 452 (d) Did you include the total amount of compute and the type of resources used (e.g., type
453 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5.
- 454 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 455 (a) If your work uses existing assets, did you cite the creators? [N/A]
456 (b) Did you mention the license of the assets? [N/A]
457 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
458
- 459 (d) Did you discuss whether and how consent was obtained from people whose data
460 you're using/curating? [N/A] We use data from four existing datasets: MOT17, KITTI,
461 PathTrack, and BDD. The data is provided by the respective organizations or authors
462 for research use.
- 463 (e) Did you discuss whether the data you are using/curating contains personally identifiable
464 information or offensive content? [N/A]
- 465 5. If you used crowdsourcing or conducted research with human subjects...
- 466 (a) Did you include the full text of instructions given to participants and screenshots, if
467 applicable? [N/A]
468 (b) Did you describe any potential participant risks, with links to Institutional Review
469 Board (IRB) approvals, if applicable? [N/A]
470 (c) Did you include the estimated hourly wage paid to participants and the total amount
471 spent on participant compensation? [N/A]