






Lost and Found: Overcoming Detector Failures in Online Multi-Object Tracking

Lorenzo Vaquero^{1,2}, Yihong Xu³, Xavier Alameda-Pineda⁴,
V́ctor M. Brea², and Manuel Mucientes²

¹ Fondazione Bruno Kessler, Italy
lvaquero@fbk.eu

² CiTIUS, Univ. of Santiago de Compostela, Spain
{victor.brea, manuel.mucientes}@usc.es

³ Valeo.ai, France

yihong.xu@valeo.com

⁴ Inria Grenoble, Univ. Grenoble Alpes, France
xavier.alameda-pineda@inria.fr

Abstract. Multi-object tracking (MOT) endeavors to precisely estimate the positions and identities of multiple objects over time. The prevailing approach, tracking-by-detection (TbD), first detects objects and then links detections, resulting in a simple yet effective method. However, contemporary detectors may occasionally miss some objects in certain frames, causing trackers to cease tracking prematurely. To tackle this issue, we propose BUSCA, meaning ‘to search’, a versatile framework compatible with any online TbD system, enhancing its ability to persistently track those objects missed by the detector, primarily due to occlusions. Remarkably, this is accomplished without modifying past tracking results or accessing future frames, i.e., in a fully online manner. BUSCA generates proposals based on neighboring tracks, motion, and learned tokens. Utilizing a decision Transformer that integrates multimodal visual and spatiotemporal information, it addresses the object-proposal association as a multi-choice question-answering task. BUSCA is trained independently of the underlying tracker, solely on synthetic data, without requiring fine-tuning. Through BUSCA, we showcase consistent performance enhancements across five different trackers and establish a new state-of-the-art baseline across three different benchmarks. Code available at: <https://github.com/lorenzovaquero/BUSCA>.

Keywords: 2D Tracking · Multi-target tracking · Online

1 Introduction

Multi-object tracking (MOT) entails the process of locating and identifying multiple objects over time within a scene. It is a crucial task in computer vision with applications spanning various domains such as robotics [17], autonomous vehicles [20, 71], and video surveillance systems [44]. The prevalent MOT paradigm



Fig. 1: Due to occlusions, detectors fail to locate many relevant elements on a scene (e.g., the woman in red). Accordingly, online multi-object trackers may lose track of some objects. With BUSCA, we propose a fully online framework that can be integrated into any online TbD tracker to persistently track those objects missed by the detector. Box colors represent object identities.

is tracking-by-detection (TbD) [9], where object trajectories are obtained by (i) first detecting objects and (ii) then associating detections. Although alternative frameworks have been proposed in the literature [1, 34], TbD has surfaced capitalizing on significant progress in object detection. Notably, over the past few years, center- [70, 83] and Transformer-based architectures [55, 70] have emerged. More recently, the MOT performance has been further improved thanks to the adoption of YOLO-based detectors [19, 45] coupled with a straightforward intersection-over-union (IoU) matching. This simple yet effective approach has even contributed to the renewed popularity of SORT [7, 15, 45].

Meanwhile, significant efforts in the community have been also dedicated to improving identity consistency within a trajectory. This is achieved by devising better association schemes [7, 15, 77, 84] or through re-identification (Re-ID) [46, 51]. However, these methods remain highly dependent on the availability of detections, which makes them susceptible to trajectory fragmentations.

Current state-of-the-art detectors are not perfect and fail to detect all the objects in a video. To have an idea, 17% of the detections in MOT17 [35] validation set are still missed by the YOLOX detector [19], and the extremely occluded objects (visibility = 0, provided in the ground-truth annotations) contribute 11.0 points to the MOTA score based on the standard MOT evaluation [12, 35]. Meanwhile, modern online trackers pause or terminate the tracking process during these situations where an object fails to be detected, leading to suboptimal results. We argue that more care should be taken in this regard, avoiding premature terminations of objects that genuinely exist. In this work, we introduce **BUSCA** (**B**uilding **U**nmatched trajectories **S** Capitalizing on **A**ttention), which helps online TbD systems handle those objects, often highly occluded, overlooked

by the detector. BUSCA propagates unmatched tracks and, by design, can be applied to the outcome of any online TbD track assignment process.

Some works in the literature [13, 41, 50] focus on repairing fragmented tracks and improving trajectory continuity. However, these have so far been implemented through offline methods, as they alter decisions made on previous time steps (e.g., interpolating a trajectory after re-detection) and/or leverage future information. Thus, despite some of them claiming to be online, they should be considered as offline according to the widely accepted definition of ‘online’ in MOTChallenge [12, 35] where “*the solution has to be immediately available with each incoming frame and cannot be changed at any later time*”. The offline fashion makes them impractical for certain real-world applications and not comparable to online methods. Conversely, BUSCA is able to *persistently track undetected objects in a fully online setting*⁵.

As an example illustrated in Fig. 1, some objects are missed due to low visibility even by a highly performant detector [19], causing the tracker to lose them. With BUSCA, we can enhance any TbD online tracker to continuously track those undetected objects without resorting to offline methods. To this end, BUSCA is built on a multi-choice question-answering Transformer that finds undetected objects given (i) *candidate* generated with a motion model (independent of the detector), (ii) *contextual information* derived from neighboring objects, and (iii) *previous observations* from the object of interest. These inputs are composed of visual and spatiotemporal information. The visual component characterizes object appearances while the spatiotemporal element encapsulates the size, center location, and timing of the object in a condensed format using an innovative spatiotemporal encoder.

In summary, the main contributions and novelties of this work are as follows:

- BUSCA is a *general* framework to persistently track those objects missed by the detector, in a fully online manner, without (i) modifying past tracking predictions (ii) or accessing future frames.
- BUSCA entails (i) a novel *Decision Transformer* inspired by multi-choice question-answering tasks, (ii) a *Proposal Generator* that relies on neighboring tracks, motion, and learned tokens, and (iii) an innovative *Spatiotemporal Encoder* that captures the size, location, and time of the objects. The network is trained independently from the underlying tracker and using synthetic data [16], without any fine-tuning on real MOT sequences.
- BUSCA can be seamlessly integrated on top of any online TbD tracker, as demonstrated in our comprehensive experiments where we systematically enhance the performance of five distinct trackers on standard benchmarks [12, 35], defining a new state-of-the-art among online trackers.

2 Related Work

End-to-end MOT methods model detection, tracking, and their implicit matching within a unified architecture. The most common approaches tackle this

⁵ BUSCA strictly respects the ‘online’ definition, thus ‘fully online’.

through identity embeddings [65], regression [1, 61] or the recent use of attention mechanisms [6, 18, 34, 76, 84, 85]. Nonetheless, this holistic design can create challenges during the joint training process [18] and, prevent these methods from being applicable to other trackers and leveraging leading-edge detectors. Consequently, these models have not yet superseded TbD techniques.

Tracking by detection (TbD) is an effective paradigm that decouples the MOT task into object detection and data association. This decomposition enables TbD methods [21, 51, 55, 69, 70, 77, 80, 83] to benefit from classical [45, 47, 69], more advanced [26, 77] or self-constructed [55, 70, 83] detectors, coupled with diverse association processes such as hierarchical clustering [80], graph neural networks [21] or geometric cues [77].

In particular, center-based methods like CenterTrack [84] and TransCenter [70] alleviate the ambiguity in bounding boxes by predicting object center heatmaps in a CNN-based or Transformer-based architecture, respectively. Recently, ByteTrack [77] showcases remarkable results using a meticulously tuned YOLOX detector [19] paired with a simple IoU-based matching mechanism. This powerful detector has also revived SORT [4] with a stronger association mechanism in methods such as OC-SORT and StrongSORT [7, 15]. Nevertheless, these TbD trackers remain highly vulnerable to missed detections. This issue motivates us to introduce BUSCA, a framework designed to improve any online TbD tracker by persistently tracking those objects overlooked by the detector.

Improving trajectory consistency, i.e., maintaining consistent object identities over time, is one of the main challenges of online multi-object trackers. Most of these methods rely on frame-by-frame association of detections solved via Hungarian matching [28]. However, pure motion-based associations [4, 5, 77] often encounter difficulties in crowded environments or moving-camera scenarios. As a result, other works turn to appearance-based techniques [27, 43, 46, 52, 59, 67], hybrid cues [15, 29, 51, 58], or Transformer solvers [76, 84]. Notably, GHOST [51] redesigns the use of a ReID model and builds a simple yet strong baseline. In efforts to lessen the impact of occlusions, some methods aim to predict an object’s visibility in order to adjust its detections’ confidences [24] or re-weight the association matrix [79]. On the other hand, some strategies improve associations by hallucinating object trajectories [57] or by prompting re-detections in areas where occluders are present [31].

Nonetheless, unlike BUSCA, these more advanced association processes *re-main heavily dependent on the detector as they operate on available detections*. [29] is a rare exception but at the cost of MOT performance drop.

Ensuring trajectory continuity is a non-trivial task that attempts to repair the trajectory of an object from the instant it is lost until it is re-identified again. Thus, most current trackers perform an extra *offline* post-processing step based on linear [77] or Gaussian-smoothed [15] interpolation. Some more sophisticated methods involve implementing a probabilistic model to retroactively insert missed detections [50], learning an additional Refind Module [41] to bridge these gaps, or 2D-to-3D lifting and performing motion forecasting in a bird’s eye view [13]. Nevertheless, these strategies remain *offline* [12, 35] as they either al-

ter predictions on past time steps or take into account future frames, limiting their applicability in certain real-world scenarios. We introduce thus BUSCA, a framework that can be *built on top of any online TbD tracker* to enhance its continuity and consistency *in a fully online fashion*.

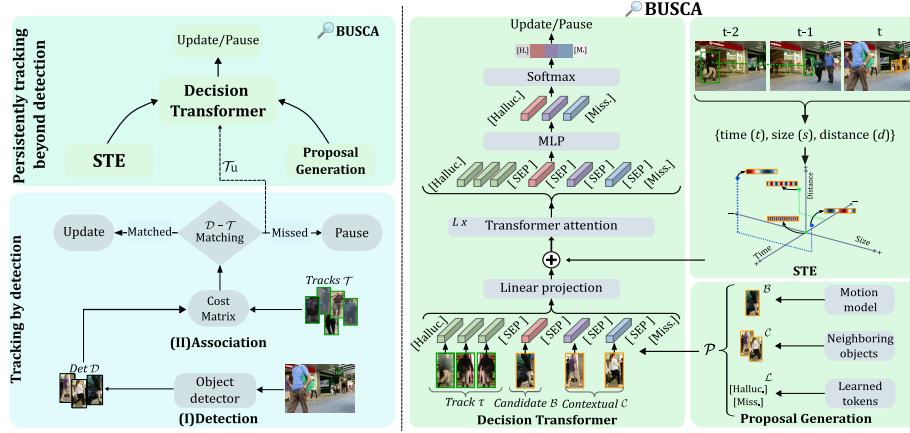


Fig. 2: The bottom-left panel depicts the tracking-by-detection (TbD) paradigm (Sec. 3), where a track is paused when the detector fails to locate the object. To address this issue, we integrate BUSCA into the online TbD tracker (Sec. 4) as shown in the top-left panel. This allows for the extension of trajectories of undetected objects by pairing them with proposals comprising *candidates* (\mathcal{B}), *contextual information* (\mathcal{C}) and *learned tokens* (\mathcal{L}) (Sec. 4.2) via an innovative decision Transformer (Sec. 4.1). Comprehensive details about the components of BUSCA are showcased in the right-hand panel. The *track* observations and proposals fed to the decision Transformer are made up of both appearance features (extracted with a convolutional backbone omitted here for clarity) and spatiotemporal cues for time, size, and distance encoded in a compact embedding through our novel spatiotemporal encoding (STE, Sec. 4.3).

3 TbD in a Nutshell

In the tracking by detection (TbD) paradigm, at a given frame a detector first produces a set $\mathcal{D} = \{\delta_1, \dots, \delta_M\}$ of M detections, with each detection $\delta_i = \{a_i, c_i, \omega_i\}$ is defined by its appearance a_i (i.e., features of the image contained in the coordinates), coordinates c_i (object size and center location) and confidence score ω_i . These detections are used to propagate the position of a set $\mathcal{T} = \{\tau_1, \dots, \tau_N\}$ of N active tracks, each represented by a time-ordered set $\tau_j = (o_{j,1}, \dots, o_{j,Z})$ of observations $o_k = \{a_k, c_k\}$ over the past Z frames.

\mathcal{D} is compared with \mathcal{T} , using coordinates and geometric cues [4, 77], appearance information [67], or both [51], yielding a cost matrix of size $N \times M$ whose optimal assignments are determined through Hungarian matching [28].

Thus, as shown in the bottom-left part of Fig. 2, correctly matched tracks are updated with the assigned detections, while those without a matching detection are paused. Having correct and sufficient detections for all tracks is critical, leading many trackers to resort to offline interpolation techniques to repair missing observations. In order to address this issue without resorting to offline interpolation, we present BUSCA, which tracks those undetected objects in a fully online fashion.

4 BUSCA: Finding Objects without Detections

Current detectors still fail to detect all the objects, especially in low-visibility situations i.e., heavy occlusions. Modern trackers heavily rely on the detection quality, thus naively stopping the tracking process whenever the detector fails. Therefore, BUSCA comes to help by saving those objects missed by the detector and finding where they are.

In particular, BUSCA is a fully online framework that can be coupled with any TbD tracker to persistently track those objects missed by the detector. As can be seen in the upper left part of Fig. 2, BUSCA receives unmatched tracks \mathcal{T}_u and compares them with a set of proposals generated through a proposal generation process (Sec. 4.2). This comparison is carried out through a novel decision Transformer (Sec. 4.1), which uses an innovative spatiotemporal encoding (STE, Sec. 4.3) to aggregate information of different nature. This way, BUSCA can update the coordinates of those unmatched tracks or determine whether they have really left the scene.

4.1 Decision Transformer: To Be or Not To Be

Deciding whether to pause an undetected track or propagate its identity can be formulated as a multiple-choice question-answering task [42]. That is, given a question (the track τ) and a set of possible options (the proposals $\mathcal{P} = \{p_1, \dots, p_J\}$, where $p_i = \{a_i, c_i\}$), the goal of the network is to find the correct answer (the decision of which proposal to match to the track) forming the assignment set $\mathcal{A} = \{\tau_j \mapsto p_i | \tau_j \in \mathcal{T}, p_i \in \mathcal{P}\}$. Inspired by this formulation, we propose to maintain undetected objects via a Transformer-based design that inputs different *proposals* and a *track*, outputting the best match, i.e., the proposal with the highest probability.

As shown on the right side of Fig. 2, our decision Transformer is implemented through an L -layer encoder model, which receives an input $\mathcal{I} = \{\tau, \mathcal{P}\}$, in which the past observations of the track are included. For each of the individual elements that make up the input (referred to as *tokens*), the appearance information a is processed by a convolutional backbone and projected to a lower dimensional space. This visual information of each token is then fused with its geometric cues c using our innovative spatiotemporal encoding (Sec. 4.3), to allow the Transformer to reason complex relationships between motion and visual features.

Within the decision Transformer, the input tokens are self-attended with each other, yielding refined tokens $\mathcal{J} = \{\bar{\tau}, \bar{\mathcal{P}}\}$ where the features most closely related to the track have been enhanced. Then, the elements of $\bar{\mathcal{P}}$ are fed to a shared-weight multi-layer perceptron (MLP) that generates one logit per token. After a Softmax operation, we output the probabilities that the track τ is assigned to each proposal p , allowing us to obtain \mathcal{A} by finding the maximum probability. Finally, we update τ when it is successfully matched with a candidate proposal (See Sec. 4.2) or pause it otherwise. It should be noted that the MLP is share-weight, so as not to be restricted to any fixed input size.

4.2 Proposal Generation: Missing Puzzle Pieces

As with textual question-answering problems, the composition of the proposals \mathcal{P} is one of the most critical aspects, and this is no different for our decision Transformer. $\mathcal{P} = \{\mathcal{B}, \mathcal{C}, \mathcal{L}\}$ is composed of candidates \mathcal{B} , contextual proposals \mathcal{C} , and learned proposals \mathcal{L} . As shown in the bottom-right of Fig. 2, \mathcal{B} and \mathcal{C} are extracted from the frame, while \mathcal{L} is learned. BUSCA will keep a track τ active and update it with the proposal information if it is associated with any element from \mathcal{B} and pause τ otherwise.

Generating the sets of proposals \mathcal{B} and \mathcal{C} is nontrivial given that none of the detections in \mathcal{D} can be associated with τ . Given its reasonable performance [4, 15, 77], we opt for a simple yet effective Kalman filter [25] to predict a new observation of τ at the current frame. To this end, it is possible to obtain $\mathcal{B} = \{\text{Kalman}(\tau)\}$ without adding extra complexity to BUSCA, all while effectively managing complex motion scenarios, as evidenced in the supplementary material. Regarding the contextual proposals \mathcal{C} , their goal is to provide BUSCA with more information about the scene. \mathcal{C} is composed of the Q closest observations within the neighborhood of τ , $V(\tau)$. Details for the computation of the maximum neighborhood distance for τ are given in the supplementary material.

The input proposals \mathcal{P} of BUSCA also comprise a set $\mathcal{L} = \{[\text{Halluc.}], [\text{Miss.}]\}$ of learned tokens that allow the Transformer to make complex decisions about the tracking process and pause τ if necessary. Specifically, $[\text{Halluc.}]$ is learned to capture whether any observation o is corrupted (i.e., belonging to a different object) whereas $[\text{Miss.}]$ handles if τ has left the scene or none of the elements of $\{\mathcal{B}, \mathcal{C}\}$ are suitable enough to be matched. Additionally, a separator token $[\text{SEP}]$ borrowed from textual Transformers [42] is also learned to delimit each of the elements of \mathcal{P} .

4.3 Spatiotemporal Encoding (STE): Merging Modalities

Along with appearance features, spatiotemporal information is also crucial for making correct assignments. This information is however more complex to be encoded due to its multi-dimensionality (i.e., time-stamp t at which observations are recorded, the size s of the bounding box, and their distance d in the 2D coordinate space). To this end, we propose the spatiotemporal encoding (STE)

depicted on the top-right part of Fig. 2, which models these relationships between observations and allows its fusion with visual features so BUSCA can effectively learn complex relationships. Our spatiotemporal encoding supersedes the conventional positional encoding often implemented in Transformer models [60]. This encoding is generated through a two-step process comprising the *interplay mapping* and subsequent the *embedding projection*.

Interplay mapping. The encodings employed in visual Transformers rely on absolute values, which limit the network’s overall adaptability and make them rely on interpolation techniques to handle diverse frame sizes [8,14,37]. Moreover, this method has consequential downsides for tracking tasks, as identical interactions might be represented differently depending on their specific occurrence (e.g. proximity between a track and an observation will be encoded differently depending on their absolute position within the frame or video).

To address this, our STE relies on a novel interplay mapping that models interactions relative to an anchor κ . In our specific use case, $\kappa = \{x_\kappa, y_\kappa, w_\kappa, h_\kappa, t_\kappa\}$ corresponds to the coordinates (i.e., object center, width, and height) and time-stamp of the last known observation of the track $o \in \tau$. To this end, we can compute a spatiotemporal embedding $\{E^t, E^s, E^d\}$ comprising time, size, and distance, respectively, for each token $\iota \in \mathcal{I}$ as:

$$E^t = \sigma^t (t_\iota - t_\kappa) \quad (1)$$

$$E^s = \sigma^s \left(\log \left(\frac{w_\iota}{w_\kappa} \right) + \log \left(\frac{h_\iota}{h_\kappa} \right) \right) \quad (2)$$

$$E^d = \sigma^d \log \sqrt{\left(\frac{x_\iota - x_\kappa}{w_\kappa} \right)^2 + \left(\frac{y_\iota - y_\kappa}{h_\kappa} \right)^2} \quad (3)$$

where $\sigma^t, \sigma^s, \sigma^d$ are scaling factors. This relative representation boosts the generalization capacity of BUSCA and improves convergence during training.

Embedding Projection. After computing the interplay mapping between input tokens and τ , it is essential to make this representation compatible with both the transformer and the visual features. However, adding multiple independent sinusoidal functions could lead to potentially ambiguous information, according to [64]. To this end, it is necessary to establish a joint spatiotemporal encoding by expanding the function used in [60] to a 3-dimensional space. Given the Transformer’s internal dimension of D^{Tr} channels, we equally distribute it among the three components of our spatiotemporal embedding $D = D^{\text{Tr}}/3$. Therefore, for a given dimension E^Δ where $\Delta \in \{t, s, d\}$ we can compute its projected embedding PE^Δ :

$$PE_{2i}^\Delta = \sin \left(\frac{E^\Delta}{10000^{2i/D}} \right) \quad PE_{2i+1}^\Delta = \cos \left(\frac{E^\Delta}{10000^{2i/D}} \right) \quad (4)$$

where $0 \leq i < D/2$. And subsequently concatenate the components of the different dimensions to create our compact spatiotemporal encoding $STE = (PE^t, PE^s, PE^d)$ for each one of the tokens $\iota \in \mathcal{I}$.

5 Experimental Results

In Sec. 5.1, we clarify the experimental settings along with the used datasets and metrics. In Sec. 5.2, we validate the necessity of BUSCA compared to the naive solutions and show that it can systematically extend tracks’ lifespan, improving trajectory continuity without losing consistency. Subsequently, we empirically demonstrate the effectiveness of each component of BUSCA and justify its design choices. Once validated, we show in Sec. 5.3 that BUSCA is a plug-and-play component that consistently improves various trackers, setting new state-of-the-art performance in all tested benchmarks compared to other online methods. Finally, some successful and failure cases are qualitatively shown in Sec. 5.4.

5.1 Experimental Settings

We conduct our experiments on the widely-used MOT16 [35], MOT17 [35] and the crowded MOT20 [12] datasets. In contrast to other methods, we train BUSCA using solely synthetic data from MOTSynth [16], which consists of 764 full-HD videos recorded at 20 fps. For each training sample, we construct a track of length $Z = 11$ and randomly select 5 objects near τ to form a proposal set (current observation of τ is the positive candidate while objects with an overlap smaller than 0.5 are negatives. Additionally, we set a 15% probability of not sampling any positives ([Miss.] will be considered the correct option) and a 1% chance of altering observations within τ ([Halluc.] will be the correct option). Our training process focuses only on bounding box annotations and does not require any fine-tuning towards particular datasets or tracking systems. The computational cost of BUSCA is relatively small, with only 8.7M parameters and a runtime of 45ms per frame on a single NVIDIA RTX GPU (when integrated with [77], the whole system runs at roughly 13fps).

For the ablation, we focus on MOT17 with the widely-adopted split [51, 77, 83] that evenly divides each video sequence into training and validation sets. Unless otherwise stated, we employ ByteTrack [77] as our baseline tracker due to its state-of-the-art performance, but we remove its offline interpolation and its per-sequence curated thresholds. For the comparison with the state-of-the-art, we submit our test set results to the MOTChallenge servers and compare our approach with current *online* methods as defined in the challenge [12, 35].

For evaluation, we report the standard metrics adopted by the MOTChallenge [11]. These include MOTA [2] reflecting the overall performance of a predicted trajectory; the recently introduced HOTA [33] that balances object coverage and identity preservation; IDF1 [49] focusing on association quality; IDentity SWitches (IDSW) to reflect identity consistency; and False Positives (FP) as well as False Negatives (FN) to assess detection performance. Additional experiments and implementation details can be found in the supplementary material.

5.2 Model Validation and Ablation

Naive approaches are not enough. Persistently tracking objects overlooked by the detector is not a trivial task and cannot be achieved with simpler naive ap-

Table 1: Comparison to different simpler solutions on MOT17 [35] val set. The difference with the baseline is depicted next to each metric. ByteTrack [77] is used as base tracker removing its offline interpolation and per-sequence thresholds, noted with \star .

	MOTA \uparrow	HOTA \uparrow	FN \downarrow	FP \downarrow
ByteTrack \star	76.5	67.4	9120	3410
+ LD	75.3 (-1.2)	65.6 (-1.8)	8854 (-266)	4196 (+786)
+ IoU	75.4 (-1.1)	67.0 (-0.4)	7588 (-1532)	5493 (+2083)
+ Mixed	76.6 (+0.1)	67.6 (+0.2)	8393 (-727)	4063 (+653)
+ BUSCA (ours)	77.1 (+0.6)	67.6 (+0.2)	8326 (-794)	3889 (+479)

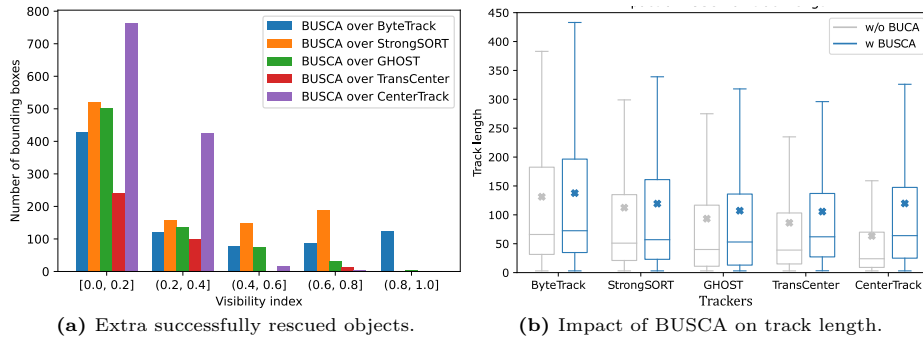


Fig. 3: (a) Analysis of the additional objects that BUSCA *successfully* locates when integrated with different trackers. The objects are grouped by their visibility [35]. (b) Analysis of the impact of BUSCA on the resulting track length in different trackers. Additional implementation details can be found in the supplementary material.

proaches. Specifically, ByteTrack [77] demonstrates that with a reliable detector, some low-score detections can be leveraged in a second-round association. One would then expect that **Lowering the Detection (LD)** threshold $\epsilon = 0.01$ would provide further benefits during the tracks-detections matching. Another direct approach similar to BUSCA consists of using a motion model (e.g., Kalman filter) to estimate the track future coordinates and perform an extra round of associations based on motion and geometry cues like **IoU**. Alternatively, we also propose an extra recovery round based on **Mixed** cues (i.e. both IoU and appearance), as shown important for more robust associations [67].

As shown in Tab. 1, lower-score detections are not reliable and **+LD** increases FP (+786) with a slight decrease in FN (-266), leading to a MOTA (-1.2) and HOTA (-1.8) drop. This demonstrates that the leftover detections in [77] are not reliable and insufficient for finding lost objects and it is therefore necessary to leverage a motion model providing better candidates. However, not every candidate is reliable, and relying solely on **+IoU** associations does not improve MOT performance (-1.1/-0.4 in MOTA/HOTA). Adding visual cues with our **+Mixed** approach brings improvements, but the limited increase in MOTA (+0.1) evidences that this simple method still struggles to make correct assignments. Differently, **BUSCA** considers visual and spatiotemporal information

Table 2: Ablation on MOT17 [35] val set of the different components that comprise BUSCA. HLC=[Halluc.] learned token, MSS=[Miss.] learned token, STE=spatiotemporal encoding, CTX=contextual proposals. The difference with the baseline is depicted next to each metric. ByteTrack [77] is used as base tracker removing its offline interpolation and per-sequence thresholds.

Line	HLC	MSS	STE	CTX	MOTA \uparrow	HOTA \uparrow	FN \downarrow	FP \downarrow
1					76.5	67.4	9120	3410
2	✓				75.0 (-1.5)	66.3 (-1.1)	8395 (-725)	4911 (+1501)
3		✓			76.4 (-0.1)	67.3 (-0.1)	8064 (-1056)	4513 (+1103)
4	✓	✓			76.5 (0.0)	67.1 (-0.3)	8656 (-464)	3853 (+443)
5	✓	✓	✓		76.7 (+0.2)	67.4 (0.0)	8528 (-592)	3851 (+441)
6	✓	✓	✓	✓	76.9 (+0.4)	67.6 (+0.2)	8387 (-733)	3884 (+474)
7	✓	✓	✓	✓	77.1 (+0.6)	67.6 (+0.2)	8326 (-794)	3889 (+479)

from the track, the candidate, and the context in a Transformer-based design, providing better decisions to prevent undetected tracks from being paused.

Longer trajectories with BUSCA. As illustrated in Fig. 3a, the efficacy of BUSCA is evident in its ability to *successfully* keep alive an extensive array of missing objects under different baselines. We observe that most of those saved objects have low visibility (i.e., under heavy occlusions), proving that BUSCA is particularly good at mitigating instances where the detector exhibits a proclivity for failure. Accordingly, BUSCA correctly extends the resulting track trajectories *in every tested tracker*, as demonstrated in Fig. 3b.

BUSCA component ablation. BUSCA relies on different components that ensure its proper operation and allow it to associate proposals and tracks accurately. In Table 2, we analyze the impact of the learned [Miss.] and [Halluc.] tokens, the spatiotemporal encoding, and the use of contextual information.

BUSCA may decide to pause a track either because it is a hallucinated track ([Halluc.] token), or because none of the candidates is suitable enough ([Miss.] token). Relying solely on the [Halluc.] token (Line 2) yields negative results, resulting in an additional +1501 false positives compared to the baseline. Conversely, if track termination is guided solely by the [Miss.] token (Line 3), the output remains marginally below the baseline with a decrease of -0.1 points in MOTA. The integration of these two learned tokens leads to improved performance (Line 4) because taking into account both conditions for whether to associate a track more accurately represents real-world situations.

By adding our spatiotemporal encoding *STE* (Line 5), the MOTA score is further increased by +0.2 points. Nonetheless, a high number of false negatives persist due to duplicated tracks occasionally kept alive. These tracks negatively impact the system when kept active, and so far BUSCA has had no way of identifying them. To address this issue, we integrate contextual proposals from nearby observations (Line 6), successfully reducing false negatives by -733 and resulting in a MOTA increase of +0.4 points. The best results are achieved when all components are integrated into BUSCA (Line 7).

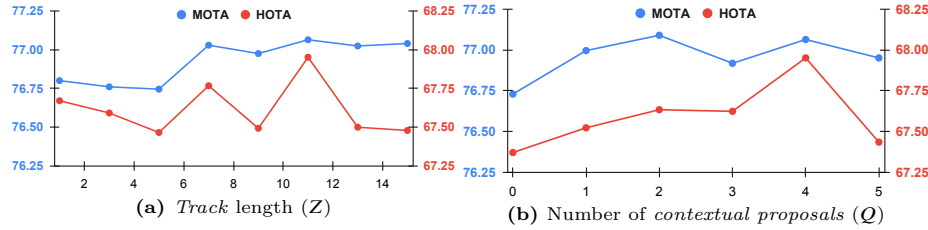


Fig. 4: Study of track length and number of contextual proposals used as input in our decision Transformer w.r.t. HOTA and MOTA performance.

Track length, contextual proposal size. Adhering to the definition of an online method, BUSCA considers the past observations of a track and its interaction with neighboring objects, learning deep relationships between motion and appearance. On Fig. 4a, we study the optimal amount of observations fed as input to BUSCA. The HOTA curve contains noisier observations, whereas MOTA displays an upward trend that starts to converge at $Z = 11$ where HOTA also achieves the best score. Regarding the maximum number of contextual proposals, from Fig. 4b, we observe that both curves have a positive slope which decays when $Q > 4$. We hypothesize this is due to the additional contextual proposals being too distant and uninformative on the track’s environment.

5.3 State-of-the-Art Comparisons

By design, BUSCA can be seamlessly incorporated into any existing online TbD tracker. To illustrate its performance, we extensively integrate BUSCA into five diverse state-of-the-art trackers and compare them against the current state-of-the-art in online MOT. Our base trackers include the center-based CenterTrack [82] (CNN network) and TransCenter [70] (Transformer network); as well as the YOLOX-based ByteTrack [77] (IoU matching), StrongSORT [15] (appearance-enhanced association), and GHOST [51] (attentive Re-ID scheme). Evaluations were conducted on the test sets of MOT16 [35], MOT17 [35], and MOT20 [12]. As shown in Tab. 3, *BUSCA consistently improves the performance of all trackers in every benchmark for nearly all metrics, without requiring training on any real MOT data nor necessitating to be fine-tuned for any tracker.*

Remarkably, BUSCA *drastically* enhances both CenterTrack and TransCenter without the necessity for a recent state-of-the-art detector. For instance, in CenterTrack, we achieve a boost of +12 HOTA and +21 IDF1 in MOT20. Similarly, TransCenter also gets significantly improved due to a marked reduction in IDSW, thereby bolstering HOTA (e.g., +5.1/+8.6 in MOT17/20) and IDF1 (e.g., +8.6/+15 in MOT17/20). When paired with high-performing trackers such as ByteTrack and StrongSORT that rely on a potent YOLOX detector [19], BUSCA sets a new state-of-the-art for online multi-object tracking. Furthermore, BUSCA can also join efforts with identity-preserving methods like the advanced Re-ID mechanism in GHOST [51] to further enhance its performance.

Table 3: State-of-the-art comparison on MOT16, MOT17, and MOT20 test sets. \star means that the offline interpolation and the per-sequence thresholds in ByteTrack [77] and OC-SORT [7] are removed for fair comparison. \dagger and \ddagger indicate reproduced results for GHOST [51] and StrongSORT [15] on MOT16 and for CenterTrack [83] on MOT20, respectively, due to their unavailability in the original works. Private detections are used. BUSCA consistently improves all baseline trackers in almost every metric, as shown in **bold**. Best results are highlighted in blue.

	MOT16				MOT17				MOT20			
	MOTA \uparrow	HOTA \uparrow	IDF1 \uparrow	IDSW \downarrow	MOTA \uparrow	HOTA \uparrow	IDF1 \uparrow	IDSW \downarrow	MOTA \uparrow	HOTA \uparrow	IDF1 \uparrow	IDSW \downarrow
TubeTK [38]	66.9	50.8	62.2	1236	63.0	48.0	58.6	5727	—	—	—	—
CTracker [40]	67.6	48.8	57.2	1897	66.6	49.0	57.4	5529	—	—	—	—
QDTrack [39]	69.8	54.5	67.1	1097	68.7	53.9	66.3	3378	—	—	—	—
TraDeS [68]	70.1	53.2	64.7	1144	69.1	52.7	63.9	3555	—	—	—	—
MTrack [74]	72.9	—	74.3	642	72.1	—	73.5	2028	63.5	—	69.2	6031
MeMOT [6]	72.6	57.4	69.7	845	72.5	56.9	69.0	2724	63.7	54.1	66.1	1938
MeMOTR [18]	—	—	—	—	72.8	58.8	71.5	1902	—	—	—	—
GSDT [63]	74.5	56.6	68.1	1229	73.2	55.2	66.5	3891	67.1	53.6	67.5	3230
Decode-MOT [29]	74.7	60.2	73.0	1094	73.2	59.6	72.0	3363	67.2	54.5	69.0	2805
MOTR [76]	—	—	—	—	73.4	57.8	68.6	2439	—	—	—	—
OUTrack [31]	74.2	59.2	71.1	1328	73.5	58.7	70.2	4122	68.6	56.2	69.4	2223
FairMOT [78]	75.7	61.6	75.3	621	73.7	59.3	72.3	3303	61.8	54.6	67.3	5243
TrackFormer [34]	—	—	—	—	74.1	57.3	68.0	2829	68.6	54.7	65.7	1532
TransTrack [55]	—	—	—	—	74.5	—	63.9	3663	64.5	—	59.2	3565
AOH [24]	—	—	—	—	75.1	59.6	72.6	3312	67.9	55.1	70.0	2698
GTR [84]	—	—	—	—	75.3	59.1	71.5	2859	—	—	—	—
CrowdTrack [53]	—	—	—	—	75.6	60.3	73.6	2544	70.7	55.0	68.2	3198
OC-SORT \star [7]	—	—	—	—	76.0	61.7	76.2	2199	73.1	60.5	74.4	1307
SGT [23]	76.8	61.2	73.5	1276	76.3	60.6	72.4	4578	72.8	56.9	70.5	2649
CorrTracker [62]	76.6	61.0	74.3	1709	76.5	60.7	73.6	3369	65.2	—	69.1	5183
ReMOT [73]	76.9	60.1	73.2	742	77.0	59.7	72.0	2853	—	—	—	—
Unicorn [72]	—	—	—	—	77.2	61.7	75.5	5379	—	—	—	—
MTracker [79]	—	—	—	—	77.3	—	75.9	3255	66.3	—	67.7	2715
MO3TR-YOLOX [85]	—	—	—	—	77.6	60.3	72.9	2847	72.3	57.3	69.0	2200
CountingMOT [48]	77.6	62.0	75.2	1087	78.0	61.7	74.8	3453	70.2	57.0	72.4	2795
CenterTrack \ddagger [83]	69.6	—	60.7	2124	67.8	52.2	64.7	3039	45.8	31.8	36.6	6296
+ BUSCA (ours)	70.4 (+0.8)	55.7 (-)	69.7 (+9.0)	927 (-1197)	68.9 (+1.1)	55.1 (+2.9)	68.8 (+4.1)	2817 (-222)	49.5 (+3.7)	44.2 (+12)	58.0 (+21)	1370 (-4926)
TransCenter [70]	75.7	56.9	65.9	1717	76.2	56.6	65.5	5427	72.9	50.2	57.7	2625
+ BUSCA (ours)	75.7 (+0.0)	61.9 (+5.0)	74.5 (+8.6)	1038 (-679)	76.2 (+0.0)	61.7 (+5.1)	74.1 (+8.6)	3282 (-2145)	73.9 (+1.0)	58.8 (+8.6)	72.4 (+15)	1756 (-869)
GHOST \dagger [51]	78.3	63.0	77.4	709	78.7	62.8	77.1	2325	73.7	61.2	75.2	1264
+ BUSCA (ours)	78.5 (+0.2)	63.2 (+0.2)	77.5 (+0.1)	707 (-2)	79.0 (+0.3)	62.9 (+0.1)	77.0 (-0.1)	2295 (-30)	74.2 (+0.5)	61.3 (+0.1)	75.1 (-0.1)	1251 (-13)
StrongSORT \dagger [15]	78.3	63.8	78.9	437	78.3	63.5	78.5	1446	72.2	61.5	75.9	1066
+ BUSCA (ours)	78.4 (+0.1)	64.2 (+0.4)	79.5 (+0.6)	434 (-3)	78.6 (+0.3)	63.9 (+0.4)	79.2 (+0.7)	1428 (-18)	72.7 (+0.5)	61.8 (+0.3)	76.3 (+0.4)	1006 (-60)
ByteTrack \star [77]	78.2	62.8	77.2	892	78.9	62.8	77.1	2363	74.2	60.4	74.5	925
+ BUSCA (ours)	78.5 (+0.3)	63.3 (+0.5)	77.9 (+0.7)	743 (-145)	79.3 (+0.4)	63.1 (+0.3)	77.7 (+0.6)	2358 (-5)	74.5 (+0.3)	60.5 (+0.1)	74.4 (-0.1)	920 (-5)

Lastly, recent tracking-by-attention methods [18, 34, 76, 84, 85] strive to create a fully end-to-end architecture that performs both object detection and track-detection matching within a single network. However, this streamlined process hinders their ability to easily incorporate new elements, such as a more powerful detector. This is illustrated by MOT3TR-YOLOX [85], a recent model that, despite adopting a more powerful YOLOX detection backbone, still underperforms TransCenter+BUSCA by -1.4 HOTA, -1.2 IDF1 in MOT17 and by -1.5 HOTA -3.4 IDF1 in MOT20. This underscores the superior performance of TbD methods and the opportunities that BUSCA brings, offering a plug-and-play module that systematically enhances state-of-the-art TbD trackers in a fully online manner without the need for retraining.

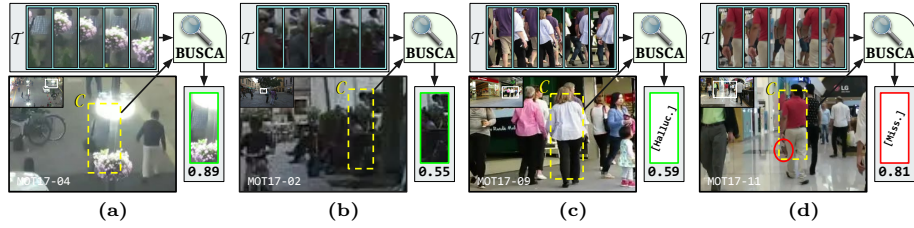


Fig. 5: Qualitative examples of BUSCA integrated into ByteTrack [77] for MOT17-val [11]. a, b, and c depict correct predictions while d shows a scenario where BUSCA incorrectly labels the pedestrian wearing a gray shirt as ‘missing’, even though the individual’s left foot (highlighted with a red circle) remains visible. The values indicate the assignment confidence.

5.4 Qualitative Results

Fig. 5 showcases a series of qualitative visualizations. In Fig. 5a, the YOLOX detector [19] fails to detect the person obscured by the street lamp and flowers due to substantial occlusion. However, with BUSCA, we can successfully preserve his identity. A similar scenario unfolds in Fig. 5b, where the pedestrian in the background is accurately identified by BUSCA despite his minimal size and the scarce visibility of only his head. Fig. 5c illustrates a clearly spurious track created by ByteTrack [77] that does not correlate to any specific person. BUSCA correctly identifies it as a hallucination and deactivates it, effectively preventing any further false positives. Lastly, in Fig. 5d, due to the noisy track and the almost total occlusion, the pedestrian wearing a gray shirt is incorrectly labeled as missing, even though his left foot can still be spotted behind the man in red. Additional videos are provided in the supplementary material.

6 Conclusion

In this work, we present BUSCA, an innovative and plug-and-play framework that can enhance any online tracking-by-detection system to persistently track undetected objects in a fully online fashion. This implies that BUSCA *does not* alter the outputs of previous time steps or access future frames. To achieve this, our novel Decision Transformer associates tracks with proposals having both visual and spatiotemporal information, maintaining the identity of tracks in a lightweight manner and without any need for fine-tuning.

We extensively validate our proposed method with five distinct trackers, bringing systematic performance improvements and setting new state-of-the-art results across different benchmarks. For future work, we aim to factor in extreme motions via nonlinear multi-candidate proposals, incorporate 3D multimodal cues, and explore the use of BUSCA to override previous tracking decisions and fix incorrect associations. We hope that BUSCA can inspire future research towards fully online trackers without overly relying on the detectors.

Acknowledgements

This work was partially supported by the EU ISFP PRECRISIS (ISFP-2022-TFI-AG-PROTECT-02-101100539) project, the EU WIDERA PATTERN (HORIZON-WIDERA-2023-ACCESS-04-01-101159751) project, MIAI@Grenoble Alpes (ANR-19-P3IA-0003), and the Spanish Ministerio de Ciencia e Innovación (grant numbers PID2020-112623GB-I00 and PID2021-128009OB-C32). We thank Eloi Zablocki from Valeo.ai for the meaningful discussion.

References

1. Bergmann, P., Meinhardt, T., Leal-Taixé, L.: Tracking without bells and whistles. In: IEEE Int. Conf. Comput. Vis. (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. Bertinetto, L., Valmadre, J., Henriques, J.F., Vedaldi, A., Torr, P.H.S.: Fully-convolutional siamese networks for object tracking. In: European Conf. Comput. Vis. (ECCV) Workshops. pp. 850–865 (2016)
4. Bewley, A., Ge, Z., Ott, L., Ramos, F.T., Upcroft, B.: Simple online and realtime tracking. In: IEEE Int. Conf. Image Process. (ICIP). pp. 3464–3468 (2016)
5. Bochinski, E., Eiselein, V., Sikora, T.: High-speed tracking-by-detection without using image information. In: IEEE Int. Conf. Adv. Video Signal Based Surveill. (AVSS). pp. 1–6 (2017)
6. Cai, J., Xu, M., Li, W., Xiong, Y., Xia, W., Tu, Z., Soatto, S.: Memot: Multi-object tracking with memory. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 8080–8090 (2022)
7. Cao, J., Weng, X., Khiroudkar, R., Pang, J., Kitani, K.: Observation-centric sort: Rethinking sort for robust multi-object tracking. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 9686–9696 (2023)
8. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: IEEE Int. Conf. Comput. Vis. (ICCV). pp. 9630–9640 (2021)
9. Dai, Y., Hu, Z., Zhang, S., Liu, L.: A survey of detection-based video multi-object tracking. Displays **75**, 102317 (2022)
10. Darcet, T., Oquab, M., Mairal, J., Bojanowski, P.: Vision Transformers Need Registers. In: Int. Conf. Learn. Repr. (ICLR) (2024)
11. Dendorfer, P., Osep, A., Milan, A., Schindler, K., Cremers, D., Reid, I., Roth, S., Leal-Taixé, L.: MOTChallenge: A benchmark for single-camera multiple target tracking. Int. J. Comput. Vis. **129**(4), 845–881 (2021)
12. Dendorfer, P., Rezatofighi, H., Milan, A., Shi, J., Cremers, D., Reid, I.D., Roth, S., Schindler, K., Leal-Taixé, L.: MOT20: A benchmark for multi object tracking in crowded scenes. CoRR **abs/2003.09003** (2020)
13. Dendorfer, P., Yugay, V., Osep, A., Leal-Taixé, L.: Quo vadis: Is trajectory forecasting the key towards long-term multi-object tracking? Adv. Neural Inf. Process. Syst. (NeurIPS) **35**, 15657–15671 (2022)
14. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: Int. Conf. Learn. Repr. (ICLR) (2021)

15. Du, Y., Zhao, Z., Song, Y., Zhao, Y., Su, F., Gong, T., Meng, H.: Strongsort: Make deepsort great again. *IEEE Trans. Multimedia* (2023)
16. Fabbri, M., Brasó, G., Maugeri, G., Ošep, A., Gasparini, R., Cetintas, O., Calderara, S., Leal-Taixé, L., Cucchiara, R.: Motsynth: How can synthetic data help pedestrian detection and tracking? In: *IEEE Int. Conf. Comput. Vis. (ICCV)*. pp. 10829–10839 (2021)
17. Gad, A., Basmaji, T., Yaghi, M., Alheeh, H., Alkhedher, M., Ghazal, M.: Multiple object tracking in robotic applications: Trends and challenges. *Applied Sciences* **12**(19) (2022)
18. Gao, R., Wang, L.: Memotr: Long-term memory-augmented transformer for multi-object tracking. In: *IEEE Int. Conf. Comput. Vis. (ICCV)*. pp. 9901–9910 (2023)
19. Ge, Z., Liu, S., Wang, F., Li, Z., Sun, J.: YOLOX: exceeding YOLO series in 2021. *CoRR abs/2107.08430* (2021)
20. Guo, S., Wang, S., Yang, Z., Wang, L., Zhang, H., Guo, P., Gao, Y., Guo, J.: A review of deep learning-based visual multi-object tracking algorithms for autonomous driving. *Applied Sciences* **12**(21) (2022)
21. He, J., Huang, Z., Wang, N., Zhang, Z.: Learnable graph matching: Incorporating graph partitioning with deep feature learning for multiple object tracking. In: *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*. pp. 5299–5309 (2021)
22. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*. pp. 770–778 (2016)
23. Hyun, J., Kang, M., Wee, D., Yeung, D.: Detection recovery in online multi-object tracking with sparse graph tracker. In: *IEEE Winter Conf. Appl. Comp. Vis. (WACV)*. pp. 4839–4848 (2023)
24. Jiang, M., Zhou, C., Kong, J.: AOH: online multiple object tracking with adaptive occlusion handling. *IEEE Signal Process. Lett.* **29**, 1644–1648 (2022)
25. Kalman, R.E.: A new approach to linear filtering and prediction theory. *J. Fluids. Eng.* **82**(1), 35–45 (1960)
26. Khan, A.H., Munir, M., van Elst, L., Dengel, A.: F2dnet: Fast focal detection network for pedestrian detection. In: *IEEE Int. Conf. Pattern Recognit. (ICPR)*. pp. 4658–4664 (2022)
27. Kim, C., Li, F., Alotaibi, M., Rehg, J.M.: Discriminative appearance modeling with multi-track pooling for real-time multi-object tracking. In: *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*. pp. 9553–9562 (2021)
28. Kuhn, H.W.: The hungarian method for the assignment problem. *Naval Research Logistics* **52**(1-2), 83–97 (1955)
29. Lee, S.H., Park, D.H., Bae, S.H.: Decode-mot: How can we hurdle frames to go beyond tracking-by-detection? *IEEE Trans. Image Process.* **32**, 4378–4392 (2023)
30. Liu, L., Jiang, H., He, P., Chen, W., Liu, X., Gao, J., Han, J.: On the variance of the adaptive learning rate and beyond. In: *Int. Conf. Learn. Repr. (ICLR)* (2020)
31. Liu, Q., Chen, D., Chu, Q., Yuan, L., Liu, B., Zhang, L., Yu, N.: Online multi-object tracking with unsupervised re-identification learning and occlusion estimation. *Neurocomputing* **483**, 333–347 (2022)
32. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: *Int. Conf. Learn. Repr. (ICLR)*. pp. 1–11 (2019)
33. Luiten, J., Osep, A., Dendorfer, P., Torr, P.H.S., Geiger, A., Leal-Taixé, L., Leibe, B.: HOTA: A higher order metric for evaluating multi-object tracking. *Int. J. Comput. Vis.* **129**(2), 548–578 (2021)
34. Meinhardt, T., Kirillov, A., Leal-Taixe, L., Feichtenhofer, C.: Trackformer: Multi-object tracking with transformers. In: *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*. pp. 8844–8854 (2022)

35. Milan, A., Leal-Taixé, L., Reid, I.D., Roth, S., Schindler, K.: MOT16: A benchmark for multi-object tracking. CoRR **abs/1603.00831** (2016)
36. Nasser, M.H., Moradi, H., Hosseini, R., Babaei, M.: Simple online and real-time tracking with occlusion handling. CoRR **abs/2103.04147** (2021)
37. Oquab, M., Darcet, T., Moutakanni, T., Vo, H., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., Assran, M., Ballas, N., Galuba, W., Howes, R., Huang, P., Li, S., Misra, I., Rabbat, M.G., Sharma, V., Synnaeve, G., Xu, H., Jégou, H., Mairal, J., Labatut, P., Joulin, A., Bojanowski, P.: Dinov2: Learning robust visual features without supervision. CoRR **abs/2304.07193** (2023)
38. Pang, B., Li, Y., Zhang, Y., Li, M., Lu, C.: Tubetk: Adopting tubes to track multi-object in a one-step training model. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 6307–6317 (2020)
39. Pang, J., Qiu, L., Li, X., Chen, H., Li, Q., Darrell, T., Yu, F.: Quasi-dense similarity learning for multiple object tracking. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 164–173 (2021)
40. Peng, J., Wang, C., Wan, F., Wu, Y., Wang, Y., Tai, Y., Wang, C., Li, J., Huang, F., Fu, Y.: Chained-tracker: Chaining paired attentive regression results for end-to-end joint multiple-object detection and tracking. In: European Conf. Comput. Vis. (ECCV). vol. 12349, pp. 145–161 (2020)
41. Qin, Z., Zhou, S., Wang, L., Duan, J., Hua, G., Tang, W.: Motiontrack: Learning robust short-term and long-term motions for multi-object tracking. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 17939–17948 (2023)
42. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al.: Improving language understanding by generative pre-training. OpenAI Research pp. 1–12 (2018)
43. Rafi, U., Doering, A., Leibe, B., Gall, J.: Self-supervised keypoint correspondences for multi-person pose estimation and tracking in videos. In: European Conf. Comput. Vis. (ECCV). pp. 36–52 (2020)
44. Rani, J.U., Raviraj, P.: Real-time human detection for intelligent video surveillance: An empirical research and in-depth review of its applications. SN Comput. Sci. **4**(3), 258 (2023)
45. Redmon, J., Divvala, S.K., Girshick, R.B., Farhadi, A.: You only look once: Unified, real-time object detection. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 779–788 (2016)
46. Ren, H., Han, S., Ding, H., Zhang, Z., Wang, H., Wang, F.: Focus on details: Online multi-object tracking with diverse fine-grained representation. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 11289–11298 (2023)
47. Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: towards real-time object detection with region proposal networks. IEEE Trans. Pattern Anal. Mach. Intell. **39**(6), 1137–1149 (2017)
48. Ren, W., Chen, B., Shi, Y., Jiang, W., Liu, H.: Countingmot: Joint counting, detection and re-identification for multiple object tracking. CoRR **abs/2212.05861** (2022)
49. Ristani, E., Solera, F., Zou, R., Cucchiara, R., Tomasi, C.: Performance measures and a data set for multi-target, multi-camera tracking. In: European Conf. Comput. Vis. (ECCV). pp. 17–35 (2016)
50. Saleh, F.S., Aliakbarian, S., Rezatofighi, H., Salzmann, M., Gould, S.: Probabilistic tracklet scoring and inpainting for multiple object tracking. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 14329–14339 (2021)

51. Seidenschwarz, J., Brasó, G., Elezi, I., Leal-Taixé, L.: Simple cues lead to a strong multi-object tracker. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 13813–13823 (2023)
52. Shuai, B., Berneshawi, A.G., Li, X., Modolo, D., Tighe, J.: Siammot: Siamese multi-object tracking. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 12372–12382 (2021)
53. Stadler, D., Beyerer, J.: On the performance of crowd-specific detectors in multi-pedestrian tracking. In: IEEE Int. Conf. Adv. Video Signal Based Surveill. (AVSS). pp. 1–12 (2021)
54. Sun, P., Cao, J., Jiang, Y., Yuan, Z., Bai, S., Kitani, K., Luo, P.: Dancetrack: Multi-object tracking in uniform appearance and diverse motion. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR) (2022)
55. Sun, P., Jiang, Y., Zhang, R., Xie, E., Cao, J., Hu, X., Kong, T., Yuan, Z., Wang, C., Luo, P.: Transtrack: Multiple-object tracking with transformer. *CoRR* **abs/2012.15460** (2020)
56. Sun, T., Segù, M., Postels, J., Wang, Y., Gool, L.V., Schiele, B., Tombari, F., Yu, F.: SHIFT: A synthetic driving dataset for continuous multi-task domain adaptation. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 21339–21350 (2022)
57. Tokmakov, P., Li, J., Burgard, W., Gaidon, A.: Learning to track with object permanence. In: IEEE Int. Conf. Comput. Vis. (ICCV). pp. 10840–10849 (2021)
58. Vaquero, L., Brea, V.M., Mucientes, M.: Real-time siamese multiple object tracker with enhanced proposals. *Pattern Recognit.* **135**, 109141 (2023)
59. Vaquero, L., Mucientes, M., Brea, V.M.: Tracking more than 100 arbitrary objects at 25 fps through deep learning. *Pattern Recognit.* **121**, 108205 (2022)
60. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Adv. Neural Inf. Process. Syst. (NeurIPS). pp. 5998–6008 (2017)
61. Wan, X., Cao, J., Zhou, S., Wang, J., Zheng, N.: Tracking beyond detection: Learning a global response map for end-to-end multi-object tracking. *IEEE Trans. Image Process.* **30**, 8222–8235 (2021)
62. Wang, Q., Zheng, Y., Pan, P., Xu, Y.: Multiple object tracking with correlation learning. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 3876–3886 (2021)
63. Wang, Y., Kitani, K., Weng, X.: Joint object detection and multi-object tracking with graph neural networks. In: IEEE Int. Conf. Rob. Autom. (ICRA). pp. 13708–13715 (2021)
64. Wang, Z., Liu, J.: Translating math formula images to latex sequences using deep neural networks with sequence-level training. *Int. J. Document Anal. Recognit.* **24**(1), 63–75 (2021)
65. Wang, Z., Zheng, L., Liu, Y., Li, Y., Wang, S.: Towards real-time multi-object tracking. In: European Conf. Comput. Vis. (ECCV). vol. 12356, pp. 107–122 (2020)
66. Williams, C.K.I., Rasmussen, C.E.: Gaussian processes for regression. In: Advances in Neural Information Processing Systems (NIPS). pp. 514–520 (1995)
67. Wojke, N., Bewley, A., Paulus, D.: Simple online and realtime tracking with a deep association metric. In: IEEE Int. Conf. Image Process. (ICIP). pp. 3645–3649 (2017)
68. Wu, J., Cao, J., Song, L., Wang, Y., Yang, M., Yuan, J.: Track to detect and segment: An online multi-object tracker. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 12352–12361 (2021)

69. Xiang, Y., Alahi, A., Savarese, S.: Learning to track: Online multi-object tracking by decision making. In: IEEE Int. Conf. Comput. Vis. (ICCV). pp. 4705–4713 (2015)
70. Xu, Y., Ban, Y., Delorme, G., Gan, C., Rus, D., Alameda-Pineda, X.: Transcenter: Transformers with dense representations for multiple-object tracking. IEEE Trans. Pattern Anal. Mach. Intell. **45**(6), 7820–7835 (2023)
71. Xu, Y., Chambon, L., Chen, M., Alahi, A., Cord, M., Perez, P., et al.: Towards motion forecasting with real-world perception inputs: Are end-to-end approaches competitive? In: IEEE Int. Conf. Rob. Autom. (ICRA) (2024)
72. Yan, B., Jiang, Y., Sun, P., Wang, D., Yuan, Z., Luo, P., Lu, H.: Towards grand unification of object tracking. In: European Conf. Comput. Vis. (ECCV). vol. 13681, pp. 733–751 (2022)
73. Yang, F., Chang, X., Sakti, S., Wu, Y., Nakamura, S.: Remot: A model-agnostic refinement for multiple object tracking. Image Vis. Comput. **106**, 104091 (2021)
74. Yu, E., Li, Z., Han, S.: Towards discriminative representation: Multi-view trajectory contrastive learning for online multi-object tracking. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 8824–8833 (2022)
75. Yu, F., Chen, H., Wang, X., Xian, W., Chen, Y., Liu, F., Madhavan, V., Darrell, T.: BDD100K: A diverse driving dataset for heterogeneous multitask learning. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 2633–2642 (2020)
76. Zeng, F., Dong, B., Zhang, Y., Wang, T., Zhang, X., Wei, Y.: MOTR: end-to-end multiple-object tracking with transformer. In: European Conf. Comput. Vis. (ECCV). vol. 13687, pp. 659–675 (2022)
77. Zhang, Y., Sun, P., Jiang, Y., Yu, D., Weng, F., Yuan, Z., Luo, P., Liu, W., Wang, X.: Bytetrack: Multi-object tracking by associating every detection box. In: European Conf. Comput. Vis. (ECCV). pp. 1–21 (2022)
78. Zhang, Y., Wang, C., Wang, X., Zeng, W., Liu, W.: Fairmot: On the fairness of detection and re-identification in multiple object tracking. Int. J. Comput. Vis. **129**(11), 3069–3087 (2021)
79. Zhang, Y., Wang, C., Wang, X., Zeng, W., Liu, W.: Robust multi-object tracking by marginal inference. In: European Conf. Comput. Vis. (ECCV). vol. 13682, pp. 22–40 (2022)
80. Zhao, K., Imaseki, T., Mouri, H., Suzuki, E., Matsukawa, T.: From certain to uncertain: Toward optimal solution for offline multiple object tracking. In: IEEE Int. Conf. Pattern Recognit. (ICPR). pp. 2506–2513 (2020)
81. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: IEEE Int. Conf. Comput. Vis. (ICCV). pp. 1116–1124 (2015)
82. Zhou, Q., Li, X., He, L., Yang, Y., Cheng, G., Tong, Y., Ma, L., Tao, D.: Transvod: End-to-end video object detection with spatial-temporal transformers. IEEE Trans. Pattern Anal. Mach. Intell. **45**(6), 7853–7869 (2023)
83. Zhou, X., Koltun, V., Krähenbühl, P.: Tracking objects as points. In: European Conf. Comput. Vis. (ECCV). vol. 12349, pp. 474–490 (2020)
84. Zhou, X., Yin, T., Koltun, V., Krähenbühl, P.: Global tracking transformers. In: IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR). pp. 8761–8770 (2022)
85. Zhu, T., Hiller, M., Ehsanpour, M., Ma, R., Drummond, T., Reid, I., Rezatofighi, H.: Looking beyond two frames: End-to-end multi-object tracking using spatial and temporal transformers. IEEE Trans. Pattern Anal. Mach. Intell. (2022)

Supplementary Material



Fig. 6: Additional qualitative results showing the benefit of using BUSCA on different online TbD trackers like ByteTrack [77]. We can see that BUSCA improves the trajectory consistency and continuity of the baseline trackers. Colors represent object identities. Results are shown for only one subject to ease the visualization.

A Introduction

In this supplementary material, we show additional qualitative results in Appendix B, which demonstrate the benefits of using BUSCA. Following this, we underline the fact that BUSCA is a generic framework (Appendix C) applicable to any online TbD method by definition, which is backed up by the experimental results in the main paper. We then provide implementation details about the network architecture, training, and inference parameters in Appendix D, and discuss in Appendix E the characteristics of the naive baselines introduced in the main paper. Subsequently, we further explain in Appendix F how to calculate the vicinity of a track and select its neighbors as contextual proposals

and further discuss the impact of BUSCA on object trajectories in Appendix G. Then, we detail how to encode spatiotemporal information for the learned tokens in Appendix H and demonstrate the effect of the [SEP] token in BUSCA’s performance in Appendix I.

BUSCA generalizes well in different trackers and scenarios without being trained on real MOT datasets. Nonetheless, we discuss in Appendix J the possible benefits of training and fine-tuning BUSCA on real in-domain data. Additionally, in Appendix K we show how our Kalman-based motion model handles complex motions like the dancing scenes found in DanceTrack [54], and in Appendix L we demonstrate the performance of BUSCA for categories different than humans.

Trackers like ByteTrack [77] and StrongSORT [15] employ *offline* interpolation methods to further improve their MOT performance by modifying past predictions with future frame information. Differently, BUSCA strictly respects the *online* definition. Although they are not directly comparable, to give an idea of the performance differences, we show in Appendix M the offline version of ByteTrack and StrongSORT versus their fully online versions with BUSCA, and the potential benefits of embedding BUSCA in an offline tracker. Lastly, to facilitate the analysis of BUSCA’s performance, we include its sequence-wise results in Appendix N.

B Additional Qualitative Results

The efficacy of online tracking by detection (TbD) is largely dependent on the underlying detectors. Issues such as object loss or identity switches often arise when these detectors miss some of the objects in the scene. In order to showcase how incorporating BUSCA into a TbD system can enhance track consistency and continuity, we provide several illustrative examples. More specifically, Fig. 6 elucidates the advantages of integrating BUSCA into ByteTrack [77], while Figs. 7, 8, 9, and 10 (found at the end of this supplementary material) visually demonstrate the application of BUSCA to StrongSORT [15], GHOST [51], TransCenter [70], and CenterTrack [83], respectively. The visualizations display the results for a single subject for ease of visualization. To view the results for every object in the sequences, the reader is referred to the videos included in this supplementary material.

C On the Generality of BUSCA

We claim that BUSCA’s design allows it to be applied to *any* online TbD tracker. This is because, by definition [9], online TbD trackers (i) detect objects in the current frame and (ii) link them to existing tracks. During this process, it is natural for some detections to remain unmatched, leading either to the creation of new tracks or their dismissal. Similarly, it is also common for some tracks to be unmatched and, therefore, paused. BUSCA introduces an additional step that (iii) propagates these unmatched tracks without requiring additional detections,

a feature applicable to the results of any online TbD track assignment procedure. Hence, we assert the adaptability of BUSCA *any* online TbD tracker.

D Network, Inference and Training Details

Network architecture. BUSCA’s decision Transformer is composed of $L = 4$ encoder blocks, each one of them comprising a 4-headed multi-attention layer and a 1024-size feed-forward layer. The internal dimension of the Transformer is $D^{Tr} = 512$ channels. For our spatiotemporal encoding, the scaling factors are set as $\sigma^t = 2$, $\sigma^s = 15$, and $\sigma^d = 15$.

The extraction of the appearance features a of each observation is performed with a ResNet-50 [22] with an extra fully-connected layer for downsampling and the domain adaptation mechanism described in [51] (i.e., samples are normalized using the mean and variance of the batch, instead of the learned ones). To this end, the coordinates c of each observation are cropped to 128×384 px and fed to the feature extractor, yielding a 512-channel embedding. This appearance model is pre-trained on the Re-ID dataset Market-1501 [81] and *is not trained* on any MOT data.

Inference parameters. Following the experimental results discussed in Sec. 5.2 of the main paper, candidate \mathcal{B} is generated using a simple-yet-effective Kalman filter [25] that forecasts a new observation for τ at the present frame, while the contextual proposal set \mathcal{C} comprises the $Q = 4$ closest observations within the neighborhood of τ . If \mathcal{B} is chosen as the correct candidate, we will keep τ active and update it with the new Kalman-based observation. Otherwise, or if τ has fewer than $Z = 11$ observations (indicating low reliability in the Decision Transformer’s prediction), we will not update the track and let the underlying base tracker handle it through its usual process (either deactivating the track or increasing its inactive counter).

Training parameters. BUSCA parameters are randomly initialized and trained via label-smoothed cross-entropy loss for 25 epochs using an AdamW optimizer [32] with a dropout regularization probability of 0.1 and a batch size of 256. We set the weight decay at 1×10^{-5} , with the initial learning rate established at 2×10^{-5} . Following the 20th epoch, we reduce the learning rate by a factor of 10.

For each training sample, we randomly choose an object identity from the dataset and construct a track τ by sampling $Z = 11$ observations, ensuring a maximum separation of 10 frames between consecutive observations. Subsequently, we randomly pick a frame within a range of 20 frames from the last observation and select 5 objects near τ to form a proposal set. From these proposals, we designate the ground truth annotation of τ as the positive candidate, while objects with an overlap smaller than 0.5 are marked as negatives (others are ignored). Additionally, we set a 15% probability of not sampling any positives, in which case the [Miss.] token will be considered the correct option within the proposal set, and a 1% chance of randomly eliminating some proposals. Furthermore, observations within τ are subject to alteration with a 1% probability,

either through removal or replacement with different object observations. If at least 5% of the observations in τ correspond to a different object, the track is deemed unreliable, and the [Halluc.] token is considered the correct option.

The entire training process took roughly 28 hours, utilizing a single NVIDIA Quadro RTX 8000. As highlighted in the main paper, BUSCA does not necessitate training aimed at any particular tracker or real MOT dataset. BUSCA is trained on a subset of 100 MOTSynth [16] videos, and does not need to be fine-tuned. Therefore, we consistently use *the same weights* across all experiments featured in our main paper.

E Implementation of Naive Approaches

In Sec. 5.2 of the main paper, we present the results of approaches proposed by us that follow the same philosophy as BUSCA (i.e., handling those tracks without matching detections in an online manner). Specifically, the IoU-based approach computes the intersection over union between the proposals generated by BUSCA and the object’s last known bounding box, assigning as correct proposal to the one with the highest overlap. On the other hand, the Mixed method uses overlap as a threshold to filter out the unrealistic proposals ($\text{IoU} < 0.7$) and, in a second step, it uses the cosine similarity between the ReID features (extracted with GHOST [51]) of proposals and the last observation of the track to determine the most suitable match.

F Neighborhood Computation

The intention behind our contextual proposals \mathcal{C} is to equip BUSCA with a broader understanding of the scene. To this end, as stated in the main paper, we pool the Q closest observations within the neighborhood of track τ , $V(\tau)$. We envision $V(\tau)$ possessing two characteristics: first, observations that are spatially adjacent to τ are deemed closer neighbors, and second, a clear demarcation is maintained between foreground and background objects. To meet these ends, we calculate the distance $\phi(\tau, o)$ between the last known coordinates of track τ and an observation o as follows:

$$\phi(\tau, o) = \text{Eucl}(\tau, o) * \text{Ratio}(\tau, o) \quad (5)$$

$$\text{Eucl}(\tau, o) = \sqrt{(x_\tau - x_o)^2 + (y_\tau - y_o)^2} \quad (6)$$

$$\text{Ratio}(\tau, o) = \max\left(\frac{\sqrt{w_\tau * h_\tau}}{\sqrt{w_o * h_o}}, \frac{\sqrt{w_o * h_o}}{\sqrt{w_\tau * h_\tau}}\right) \quad (7)$$

Thus, in defining this distance, we consider not only the Euclidean distance between centers $\text{Eucl}(\cdot, \cdot)$, but also penalize this value based on the difference in object sizes $\text{Ratio}(\cdot, \cdot)$. This approach is driven by the fact that object size serves as a strong indicator of depth in the objects within a scene [36].

Consequently, the neighborhood of τ can be defined as $V(\tau) = \{o \in \mathcal{T} \setminus \mathcal{T}_u \mid \phi(\tau, o) < \nu_\tau\}$, where ν_τ represents the maximum distance within which an observation is considered a neighbor. The maximum distance acts as a variable parameter, fluctuating about the area of the track as discussed in [3]. This can be computed using the following equation:

$$\nu_\tau = \sqrt{(w_\tau + \zeta(w_\tau + h_\tau)) * (h_\tau + \zeta(w_\tau + h_\tau))} \quad (8)$$

Here, $\zeta = 1$ is employed as a scaling factor that governs the growth of ν concerning τ .

G Impact of BUSCA on Object Trajectories

BUSCA can be incorporated into any online tracking-by-detection system, enhancing its capabilities to persistently track those objects missed by the detector. As can be seen in Fig. 3 of the main paper, BUSCA primarily focuses on those objects where the detector most frequently fails, specifically those with minimal visibility. This has the added advantage of extending the average lifespan of the tracks, thereby enhancing their trajectory consistency and continuity.

To conduct the experiment shown in Fig. 3a, we used the visibility attributes contained within the MOT17 [11] ground truth. Thus, for every object that BUSCA finds and that would otherwise have been paused by the tracker, we refer to its corresponding annotation and visibility attribute. This is performed for each combination of tracker+BUSCA studied, confirming that BUSCA is capable of identifying a substantial number of objects with extremely low visibility, due to occlusions.

About the effect of BUSCA on the track length, Fig. 3b illustrates the difference between using various standalone trackers and combining them with BUSCA. For this experiment, we evaluated the results produced by each combination, quantifying the frequency of each ID reported (i.e., the length of each track). As demonstrated, incorporating BUSCA engenders positive effects, amplifying both the median and average length of tracks for all five tested trackers.

H Learned Tokens

The appearance features a of the learned proposals $\mathcal{L} = \{\text{[Halluc.]}, \text{[Miss.]}\}$, along with the separator token [SEP] , are initialized using a random Gaussian distribution and are trained end-to-end alongside the rest of the architecture. There is no need for these features to pass through the appearance extractor, being directly fed to BUSCA’s decision Transformer.

Regarding the coordinates component c of the learned tokens, [Miss.] is given the same coordinates as the last known observation of τ and [Halluc.] , is computed by maximizing its distance w.r.t. τ in the spatiotemporal representation space (main paper, Sec. 4.3). Lastly, [SEP] tokens are given the same coordinates as the proposals they delimit.

Table 4: Ablation of the effect of the [SEP] token on MOT17 [11] validation data. BUSCA uses several [SEP] tokens to delimit every proposal $p \in \mathcal{P}$, for a total of $|\mathcal{P}|$ separator tokens. In $\times 1$ [SEP], we utilize a single token to delimit τ from \mathcal{P} . ByteTrack [77] is used as base tracker without its offline interpolation and per-sequence threshold, noted with \star .

	MOTA \uparrow	HOTA \uparrow	IDF1 \uparrow	IDSW \downarrow
ByteTrack \star	76.5	67.4	79.4	165
+ BUSCA w/ $\times 1$ [SEP]	76.8 (+0.3)	67.3 (-0.1)	78.8 (-0.8)	162 (-3)
+ BUSCA w/ $\times \mathcal{P} $ [SEP]	77.1 (+0.6)	67.6 (+0.2)	79.5 (+0.1)	166 (+1)

Table 5: Ablation of training BUSCA on in-domain data for MOT17 [11]. We test both training from scratch and fine-tuning BUSCA after training it on MOTSynth [16]. The difference with the baseline is depicted next to each metric. ByteTrack [77] is used as base tracker without its offline interpolation and per-sequence threshold, noted with \star .

	MOTA \uparrow	HOTA \uparrow	IDF1 \uparrow	IDSW \downarrow
ByteTrack \star	76.5	67.4	79.4	165
Val. + BUSCA (MOT17 train)	76.8 (+0.3)	67.4 (+0.0)	79.1 (-0.3)	167 (+2)
+ BUSCA (MOTSynth train)	77.1 (+0.6)	67.6 (+0.2)	79.5 (+0.1)	166 (+1)
+ BUSCA (MOT17 fine-tune)	77.2 (+0.7)	67.9 (+0.5)	79.8 (+0.4)	150 (-15)
Test ByteTrack \star	78.9	62.8	77.1	2363
+ BUSCA (MOT17 train)	79.3 (+0.4)	63.1 (+0.3)	77.7 (+0.6)	2358 (-5)
+ BUSCA (MOT17 fine-tune)	79.3 (+0.4)	63.1 (+0.3)	78.8 (+0.9)	2349 (-14)

I Effect of the [SEP] token

We incorporate [SEP] tokens to delimit different input segments, following standard practice in textual Transformers [42]. Nonetheless, track τ is the only element in input $\mathcal{I} = \{\tau, \mathcal{P}\}$ with variable length and, thus, separating each proposal $p \in \mathcal{P}$ using [SEP] is not necessarily mandatory. Still, Tab. 4 shows the advantages of retaining [SEP] for every proposal, as visual Transformers benefit from having additional registers to store, process, and retrieve global information [10].

J Impact of Training BUSCA on In-Domain Data

BUSCA aims to be as portable and generic as possible to facilitate its integration with any type of tracker by detection. This is why we train it on the MOTSynth [16] synthetic dataset, without making any adjustments for specific trackers or scenarios. Still, in-domain training from scratch is possible, as shown in Tab. 5 (trained on the first half of MOT17-train and validated on the second half). Despite being trained on less than 2 minutes of video, BUSCA still improves +0.3 MOTA over the baseline. Nevertheless, BUSCA still benefits from additional data, such as the 100 sequences from MOTSynth used in the primary experiments. Accordingly, to further boost the performance of BUSCA.

Table 6: BUSCA results on DanceTrack [51]. \star denotes reproduced results for GHOST [51] using the publicly available official code.

	HOTA \uparrow	IDF1 \uparrow	MOTA \uparrow	DetA \uparrow	AssA \uparrow
CenterTrack [83]	41.8	35.7	86.8	78.1	22.6
FairMOT [78]	39.7	40.8	82.2	66.7	23.8
QDTrack [39]	54.2	50.4	87.7	80.1	36.8
TransTrack [55]	45.5	45.2	88.4	75.9	27.5
TraDeS [68]	43.3	41.2	86.2	74.5	25.4
MOTR [76]	54.2	51.5	79.7	73.5	40.2
GTR [84]	48.0	50.3	84.7	72.5	31.9
ByteTrack [77]	47.7	53.9	89.6	71.0	32.1
GHOST [51]	56.7	57.7	91.3	81.1	39.8
GHOST \star [51]	54.8	55.5	91.3	81.1	37.1
+ BUSCA (ours)	55.5 (+0.7)	56.1 (+0.6)	91.5 (+0.2)	81.4 (+0.3)	38.0 (+0.9)

Table 7: BUSCA results on BDD100K [75] validation set. \star denotes reproduced results for GHOST [51] using the publicly available official code.

	mMOTA \uparrow	mHOTA \uparrow	mIDF1 \uparrow
D2TT2D-100K [75]	25.9	–	44.5
MOTR [76]	32.0	–	43.5
QDTrack [39]	36.3	41.7	51.5
TETer [30]	39.1	–	53.3
GHOST [51]	44.9	45.7	55.6
ByteTrack [77]	45.2	45.4	54.6
GHOST \star [51]	43.4	42.5	50.7
+ BUSCA (ours)	43.7 (+0.3)	43.1 (+0.6)	52.4 (+1.7)

K Behaviour under Complex Motions

BUSCA’s simple-yet-effective design is also able to model complex motions, like the ones found in DanceTrack [54]. Among the five trackers we employed, only GHOST [51] provides an official code for DanceTrack. We thus take it as an example and show here its test results with and without BUSCA. While the official code does not replicate the exact results in [51], BUSCA still yields improvements in DanceTrack, as shown in Tab. 6.

L Performance on Different Categories

The performance improvement that BUSCA yields is not limited to people only. In Tab. 7, we show how it can improve the mMOTA, mHOTA, and mIDF1 of GHOST [51] (other tested baseline trackers do not provide an official implementation for the dataset) in BDD100K [75], which comprises eight different categories. For this experiment, we freeze the appearance feature extractor and fine-tune BUSCA on SHIFT [56].

Table 8: State-of-the-art comparison on MOT17 and MOT20 test sets including offline versions of contemporary methods (colored gray). \star means that the offline interpolation and the per-sequence thresholds in ByteTrack [77] are removed.

	MOT17			MOT20		
	MOTA \uparrow	HOTA \uparrow	IDF1 \uparrow	MOTA \uparrow	HOTA \uparrow	IDF1 \uparrow
StrongSORT [15]	78.3	63.5	78.5	72.2	61.5	75.9
+ AFLink (StrongSORT+)	78.3 (+0.0)	63.7 (+0.2)	79.0 (+0.5)	72.2 (+0.0)	61.6 (+0.1)	76.3 (+0.4)
+ AFLink+GSI (StrongSORT++)	79.6 (+1.3)	64.4 (+0.9)	79.5 (+1.0)	73.8 (+1.6)	62.6 (+1.1)	77.0 (+1.1)
+ BUSCA (ours)	78.6 (+0.3)	63.9 (+0.4)	79.2 (+0.7)	72.7 (+0.5)	61.8 (+0.3)	76.3 (+0.4)
ByteTrack \star [77]	78.9	62.8	77.1	74.2	60.4	74.5
+ interp.+thresh. (ByteTrack)	80.3 (+1.4)	63.1 (+0.3)	77.3 (+0.2)	77.8 (+3.6)	61.3 (+0.9)	75.2 (+0.7)
+ BUSCA (ours)	79.3 (+0.4)	63.1 (+0.3)	77.7 (+0.6)	74.5 (+0.3)	60.5 (+0.1)	74.4 (-0.1)

M BUSCA vs. Offline Processing

BUSCA offers a comprehensive framework to persistently track objects missed by the detector in a *fully online* manner (i.e., without modifying past tracking predictions or accessing future frames). This feature renders it highly valuable for applications where the solution has to be immediately available with each incoming frame and cannot be changed at any later time. Still, certain algorithms in the literature recover objects post hoc through offline post-processing techniques. While these offline techniques are not directly comparable to BUSCA, their results serve as a possible insight into the theoretical upper bound for online methods like BUSCA. Specifically, Tab. 8 showcases the offline mechanisms employed in StrongSORT++ [15] (i.e., AFLink and GSI) and in ByteTrack [77] (i.e., offline linear interpolation and per-sequence thresholds).

AFLink, an appearance-free linking model that leverages spatiotemporal information to predict if two tracklets belong to the same object ID, provides a slight boost in HOTA and IDF1, albeit still inferior to BUSCA’s. With the addition of GSI, which employs Gaussian process regression [66] for bounding box interpolation, the achieved MOTA surpasses BUSCA by one point, highlighting the importance of handling extremely occluded objects. This effect is further emphasized with ByteTrack’s offline linear interpolation and tracking thresholds, which are adapted based on the evaluated test sequence. Nevertheless, BUSCA’s performance remains competitive, consistently enhancing the capabilities of TbD trackers in a fully online manner.

Lastly, despite being designed to enhance online multi-object trackers, BUSCA can also potentially improve batch-based and offline tracking algorithms. Exhaustive analysis in this regard falls out of the scope of this paper. However, for demonstration purposes, we show in Tab. 9 how integrating BUSCA with the interpolation-based offline version of [77] results in +0.7/+0.2 MOTA/HOTA in the MOT17 validation set.

Table 9: Compatibility of BUSCA with offline algorithms. \star means that the offline interpolation and the per-sequence thresholds in ByteTrack [77] are removed.

	MOTA \uparrow	HOTA \uparrow	IDF1 \uparrow	IDSW \downarrow
ByteTrack \star	76.5	67.4	79.4	165
+ BUSCA(ours)	77.1 (+0.6)	67.6 (+0.2)	79.5 (+0.1)	166 (+1)
ByteTrack	77.8	67.9	79.9	168
+ BUSCA(ours)	78.5 (+0.7)	68.1 (+0.2)	80.1 (+0.2)	166 (-2)

N Sequence-Wise Results

To facilitate the comparison of BUSCA with other approaches on a finer-grained level, we show sequence-wise results for the test sets of MOT16 [35], MOT17 [35], and MOT20 [12] in Tabs. 10, 11, and 12, respectively.

Table 10: Sequence-wise results on MOT16 test set. ★ means that the offline interpolation and the per-sequence thresholds in ByteTrack [77] are removed for fair comparison. Private detections are used.

	MOT16-01			MOT16-03			MOT16-06			MOT16-07			MOT16-08			MOT16-12			MOT16-14		
	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1
CenterTrack + BUSCA	60.5	45.1	53.1	86.6	63.8	80.6	56.9	45.9	59.5	55.2	41.5	50.7	39.2	42.0	46.7	48.9	48.8	60.4	41.7	37.5	50.4
TransCenter + BUSCA	55.3	49.0	58.5	90.7	71.9	85.8	58.8	45.0	56.0	67.7	49.9	60.0	50.6	44.5	51.6	54.9	53.7	67.1	48.1	42.5	60.5
GHOST + BUSCA	60.2	52.1	61.6	91.5	71.2	87.9	61.6	51.7	62.7	71.4	51.4	60.1	56.4	49.8	57.2	57.0	56.5	68.0	57.2	48.4	65.9
StrongSORT + BUSCA	61.7	52.0	63.2	91.4	72.0	89.0	62.3	51.6	63.5	71.3	55.2	69.7	56.5	48.8	56.3	62.9	59.8	73.1	54.3	48.6	68.1
ByteTrack★ + BUSCA	61.6	49.2	57.8	91.6	72.5	90.3	61.5	47.1	57.3	72.1	49.4	59.5	55.3	46.4	51.8	64.3	58.6	71.5	53.8	45.5	63.6

Table 11: Sequence-wise results on MOT17 test set. ★ means that the offline interpolation and the per-sequence thresholds in ByteTrack [77] are removed for fair comparison. Private detections are used.

	MOT17-01			MOT17-03			MOT17-06			MOT17-07			MOT17-08			MOT17-12			MOT17-14		
	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1
CenterTrack + BUSCA	60.3	45.1	53.0	86.8	64.0	80.7	57.9	45.9	59.3	55.4	41.0	50.1	32.0	37.6	40.4	48.0	48.2	59.4	41.9	37.5	50.4
TransCenter + BUSCA	57.3	49.5	59.0	90.9	72.2	85.9	60.7	45.1	55.8	68.1	49.4	59.3	58.0	43.1	50.6	54.4	53.1	66.2	48.1	42.5	60.5
GHOST + BUSCA	60.4	52.3	62.3	92.5	71.3	88.2	63.3	51.9	62.4	71.5	51.7	60.4	61.1	47.4	54.2	56.8	55.8	66.7	57.2	48.4	65.9
StrongSORT + BUSCA	61.8	52.1	63.2	92.3	72.3	89.5	62.9	51.5	63.0	71.0	55.0	69.3	58.9	46.4	54.2	62.0	59.1	72.4	54.3	48.6	68.1
ByteTrack★ + BUSCA	61.5	49.3	57.8	92.6	72.8	90.8	62.3	47.0	56.9	73.1	49.7	59.9	61.8	44.7	50.4	63.3	57.8	70.6	53.8	45.5	63.6

Table 12: Sequence-wise results on MOT20 test set. ★ means that the offline interpolation and the per-sequence thresholds in ByteTrack [77] are removed for fair comparison. Private detections are used.

	MOT20-04			MOT20-06			MOT20-07			MOT20-08		
	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1	MOTA	HOTA	IDF1
CenterTrack + BUSCA	57.1	48.3	65.6	42.8	37.7	48.1	73.9	56.1	69.3	23.8	34.5	44.2
TransCenter + BUSCA	84.9	66.1	82.2	61.1	49.0	58.6	77.7	59.3	71.1	55.6	44.9	57.2
GHOST + BUSCA	87.3	69.2	84.3	59.8	50.1	62.2	81.8	64.5	75.6	49.6	44.3	56.9
StrongSORT + BUSCA	87.1	69.4	84.7	56.8	51.1	65.2	81.9	67.6	79.8	45.3	43.5	56.1
ByteTrack★ + BUSCA	86.8	67.8	83.6	61.4	50.7	62.4	81.3	62.9	73.2	50.3	43.7	55.7

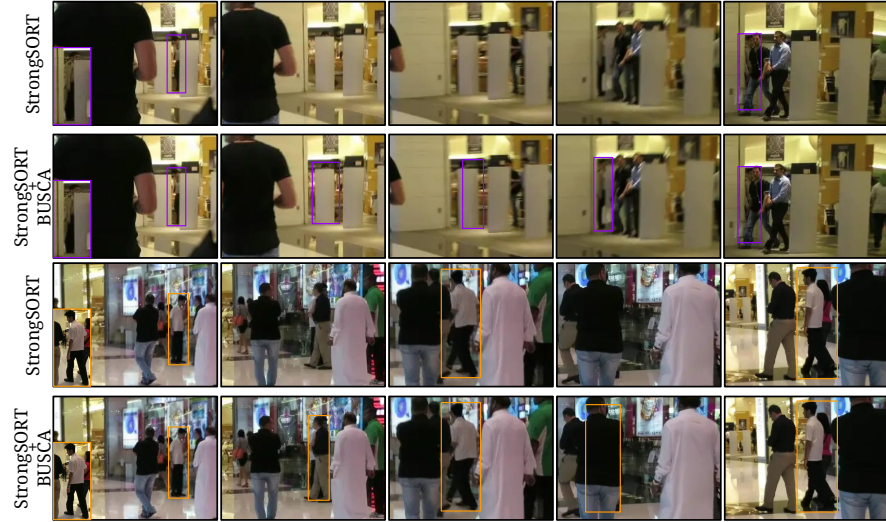


Fig. 7: Qualitative results showing the benefit of integrating BUSCA into StrongSORT [15]. Colors represent object identities. Results are shown for only one subject to ease the visualization.



Fig. 8: Qualitative results showing the benefit of integrating BUSCA into GHOST [51]. Colors represent object identities. Results are shown for only one subject to ease the visualization.



Fig. 9: Qualitative results showing the benefit of integrating BUSCA into TransCenter [70]. Colors represent object identities. Results are shown for only one subject to ease the visualization.



Fig. 10: Qualitative results showing the benefit of integrating BUSCA into CenterTrack [83]. Colors represent object identities. Results are shown for only one subject to ease the visualization.