

# BEV-SUSHI: Multi-Target Multi-Camera 3D Detection and Tracking in Bird's-Eye View

Yizhou Wang<sup>1,\*</sup>, Tim Meinhardt<sup>1,\*</sup>, Orcun Cetintas<sup>1,2</sup>, Cheng-Yen Yang<sup>1,3</sup>, Sameer S. Pusegaonkar<sup>1</sup>, Benjamin Missaoui<sup>1</sup>, Sujit Biswas<sup>1</sup>, Zheng Tang<sup>1</sup>, Laura Leal-Taixé<sup>1</sup>

<sup>1</sup> NVIDIA, <sup>2</sup> Technical University of Munich, <sup>3</sup> University of Washington

\* Equal contribution

{yizwang, tmeinhardt}@nvidia.com

## Abstract

Object perception from multi-view cameras is crucial for intelligent systems, particularly in indoor environments, e.g., warehouses, retail stores, and hospitals. Most traditional multi-target multi-camera (MTMC) detection and tracking methods rely on 2D object detection, single-view multi-object tracking (MOT), and cross-view re-identification (ReID) techniques, without properly handling important 3D information by multi-view image aggregation. In this paper, we propose a 3D object detection and tracking framework, named BEV-SUSHI, which first aggregates multi-view images with necessary camera calibration parameters to obtain 3D object detections in bird's-eye view (BEV). Then, we introduce hierarchical graph neural networks (GNNs) to track these 3D detections in BEV for MTMC tracking results. Unlike existing methods, BEV-SUSHI has impressive generalizability across different scenes and diverse camera settings, with exceptional capability for long-term association handling. As a result, our proposed BEV-SUSHI establishes the new state-of-the-art on the AICity'24 dataset with 81.22 HOTA, and 95.6 IDF1 on the WildTrack dataset.

## 1. Introduction

Detecting and tracking objects across multiple cameras is a crucial problem for indoor environment understanding, particularly in retail or warehouse settings (an example is shown in Fig. 1), for various applications, including inventory management, security surveillance, or customer behavior analysis. In these use cases, a typical multi-camera system involves numerous cameras with diverse viewing angles and fields of view (FoV). While some cameras usually have overlapping FoVs, others may not share any common visual space. This configuration allows for comprehensive coverage of the scene while balancing between redundancy and efficiency. However, this task presents significant challenges such as occlusions, varying lighting conditions, and the need to maintain consistent object identification across different camera views. Moreover, the integration of 3D

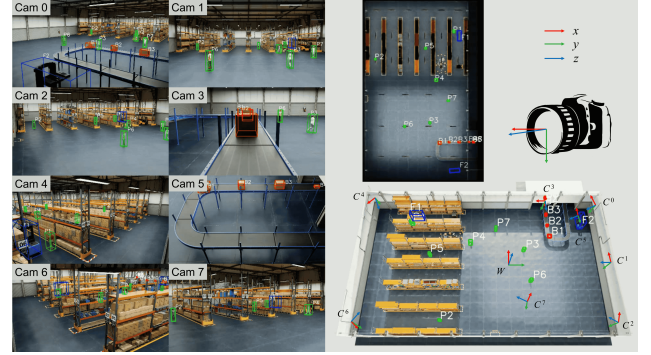


Figure 1. Example of an indoor environment (warehouse) with 8 cameras generated by NVIDIA Omniverse Platform [22]. The world coordinates  $W$  and camera coordinates  $C^i$  are drawn in the top-down visualization on the bottom right. Objects with 3D bounding boxes are annotated in both BEV and camera views.

spatial information from multiple 2D views requires sophisticated algorithms to handle camera calibration errors and perspective distortions. Addressing these issues is crucial for developing robust multi-camera detection and tracking systems that can reliably operate in complex indoor environments.

Existing MTMC tracking methods can be divided into three categories: (i) late multi-view aggregation, (ii) late multi-view aggregation with geometric projection, and (iii) early multi-view aggregation, as shown in Fig. 2. Late multi-view aggregation pipelines detect objects in each camera view as 2D bounding boxes, with some methods tracking these 2D detections separately before associating them across cameras using re-identification (ReID) embeddings and spatial constraints. Late multi-view aggregation pipelines with geometric projection pipelines [8, 34, 35] further perform additional projections using camera calibration matrices to achieve spatial association for each camera view into a global coordinate system. Recently, [29] demonstrated the possibility of early aggregation before any perception steps in each camera view. This type of method first fuses multi-view image information into a unified 3D

space and directly conducts both detection and tracking in this 3D space, which can significantly improve detection quality and avoid association errors due to unreliable spatial alignment among different camera views. However, its detection network is designed for a fixed multi-camera scene, so it is not flexible for different environments, different numbers of cameras, or different camera placements. Besides, while it outputs robust 3D detections, the method then uses an online Kalman filter based tracking which is prone to drifting and struggles with occlusions.

In this work, we propose **BEV-SUSHI**, a *generalizable* method capable of robust 3D detection and *long-term* tracking in various complex multi-camera environments.

Our approach leverages a unified BEV representation and spatio-temporal transformers to directly infer 3D object detections, eliminating the need for 2D detections and their subsequent projection to 3D space. This method bypasses the alignment step by projecting 2D information from multiple cameras into the unified BEV coordinates, resulting in a comprehensive and accurate scene representation. We also introduce a 2D-3D detection association process for ReID embedding extraction and detection verification, further enhancing the robustness and accuracy of our system.

For the temporal association of detections, we introduce the first 3D tracking-by-detection approach utilizing hierarchical graph neural networks (GNNs) [5]. Its lightweight hierarchical structure allows us to learn and model track associations for up to 4k frames. Even longer video clips are processed in a near-online sliding window fashion. Our model consists of a hierarchy of GNNs to associate 3D detections via geometrical and multi-view ReID embedding features. To improve long-term tracking, we introduce a novel global tracking layer that replaces the hand-crafted tracking heuristics in [5] with an additional layer that matches incoming with past tracks. The global layer requires no additional training and significantly boosts occlusion handling and overall tracking performance.

Finally, we evaluate BEV-SUSHI on a large-scale synthetic indoor MTMC dataset called AICity’24 [31] and a real-world outdoor MTMC dataset named WildTrack [6]. BEV-SUSHI achieves the SOTA on both datasets, with 81.22 HOTA on AICity’24, and 95.6 IDF1 on WildTrack.

Overall, our contributions can be summarized as follows:

- A multi-view image aggregation network to achieve 3D detection in large-scale multi-camera environments.
- A ReID feature extraction method for MTMC tracking based on a 2D-3D detection association algorithm.
- A GNN-based tracker to perform multi-object tracking with 3D bounding boxes in BEV coordinates with impressive long-term association capability.
- We achieve the state-of-the-art on both the AICity’24 dataset and the WildTrack dataset.

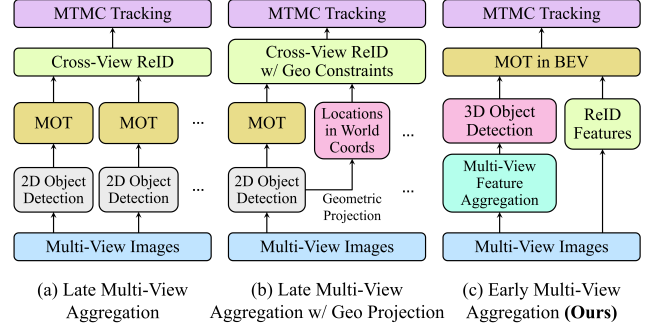


Figure 2. Comparison among 3 types of MTMC tracking methods.

## 2. Related Works

**Multi-View Object Detection** Object detection from multi-view images is an essential technique for understanding 3D geometric information and handling occlusion. With proper synchronization and calibration, images from multiple viewing angles can be integrated into the same space to accurately learn 3D geometry, *e.g.*, 3D locations, dimensions, and orientations. To begin with, researchers develop multi-view object detection methods based on some classical approaches, *e.g.*, conditional random field (CRF) [1, 25], probabilistic modeling [9, 26], and *etc.*, to aggregate information from multiple views.

Afterward, when deep learning became popular, people started introducing multi-view object detection methods based on neural networks. MVDet [11] introduces an end-to-end method, based on convolutional neural networks, which projects dense 2D image features from cameras to a unified ground plane. After spatial aggregation on the projected features by ground plane convolution, MVDet predicts a pedestrian occupancy map to obtain the final detection results. However, MVDet cannot handle different camera settings, *e.g.*, different numbers of cameras or placements, so that cannot be generalized easily to diverse environments. BEVFormer [18, 36] proposes the spatio-temporal transformer to fuse multi-view features into BEV features, which achieves impressive detection accuracy for autonomous driving related applications. The proposed spatio-temporal transformer is based on deformable attention [46] and to sample BEV features from 2D image features by grid-based reference points in BEV, which has similar setups with our smart city applications. Therefore, in this work, we aim to adapt BEVFormer under MTMC environments, whose cameras are static but with more distributed placements and variations.

**Multi-Target Multi-Camera Tracking** Recent MTMC tracking methods can be divided into three categories as shown in Fig. 2. Late multi-view aggregation methods [4, 10, 12–14] for MTMC tracking usually adopt a two-stage pipeline, *i.e.*, 2D detection and tracklet generation within each camera view followed by tracklet association across

all the cameras. Here, some cross-view association methods [10, 14] consider certain geometric constraints, *e.g.*, epipolar geometry constraints or multi-view triangulation. In contrast, others [12, 13] consider both spatial and temporal constraints, *e.g.*, camera link model. Follow-up works demonstrate that spatial association is much easier and more accurate to be done with ReID features in a unified space by geometrical projections [8, 21, 34, 35]. LMGP [21] formulates a spatial-temporal tracking graph, whose nodes are tracklets from the single-camera tracker, and edges include both temporal and spatial distances. Whereas ReST [8] first conducts spatial association and then the temporal association across frames using GNNs.

Early multi-view aggregation for MTMC tracking is first proposed by EarlyBird [29]. It first projects 2D image features into BEV by calibrated projection matrices, followed by stacking and aggregating into BEV features. A CenterNet-based [45] decoder is appended to obtain detections and ReID features in BEV. As for tracking, EarlyBird follows the idea proposed by FairMOT [42], which associates detections in the temporal domain with ReID features for appearance and a Kalman filter [15] for motion prediction. However, the detection network in EarlyBird is designed for a fixed multi-camera scene, specifically training and inference are under the same scene and the same camera settings. Therefore, it is not flexible for different environments, different numbers of cameras, or different camera placements, which affects its generalizability. As we will show in Sec. 4, EarlyBird is also not reliable for long video clips due to its limited long-term association capabilities.

Therefore, we will address the aforementioned limitations and introduce an early multi-view aggregation based MTMC tracker, which is 1) *generalizable* for different scenes and camera settings; 2) more reliable with *better long-term association* to deal with very long video clips.

### 3. BEV-SUSHI

In this section, we introduce our proposed multi-camera 3D object detection and tracking method in three folds, *i.e.*, 3D object detection, ReID feature extraction, and multi-object tracking in BEV. The overall pipeline is shown in Fig. 3.

#### 3.1. Coordinate Systems and Projection

We first define the 3D world coordinates  $\mathcal{W}$  and 3D camera coordinates  $\{\mathcal{C}^i\}_{i=1}^V$  for each specific scene. Besides, we also define the 2D image pixel coordinates  $\{\mathcal{U}^i\}_{i=1}^V$  for each camera image. Here,  $V$  is the number of camera views in the scene. The goal is to project multi-view camera information to the 3D world and do detection and tracking in  $\mathcal{W}$  without further cross-view spatial association.

As shown in Fig. 1, the origin of  $\mathcal{W}$  is defined around the center of the scene, lying on the ground plane. Axes  $x$  and  $y$

are parallel to the ground plane, and  $z$  is vertically pointing up. The origin of  $\mathcal{C}^i$  is defined at the camera center. Axis  $x$  points right,  $y$  points down, and  $z$  points forward. A 3D point  $\mathbf{x} = [x, y, z]^\top$  in the world coordinates  $\mathcal{W}$  can then be projected to the pixel coordinates  $\mathcal{U}^i$  in the  $i$ -th camera view as  $\mathbf{u} = [u, v]^\top$ . This projection can be done by

$$s \begin{pmatrix} \mathbf{u} \\ 1 \end{pmatrix} = \mathbf{P}^i \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} = \mathbf{K}^i [\mathbf{R}^i | \mathbf{t}^i]_{3 \times 4} \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}, \quad (1)$$

where  $s$  is the scale factor,  $\mathbf{P}^i$  is the projection matrix,  $\mathbf{K}^i$  is the intrinsic matrix, and  $\mathbf{R}^i, \mathbf{t}^i$  are rotation and translation.

#### 3.2. Multi-View 3D Object Detection

To aggregate multi-view image information in the early stage, we consider projecting image features into bird’s-eye view (BEV) and conducting object detection in 3D space. We leverage BEVFormer [18], a 3D object detector from surrounding camera views developed for autonomous driving applications, with adaption in Sec. 4.1. The key idea of aggregating multi-view information in BEVFormer is to sample  $N_{\text{ref}}$  3D reference points in BEV coordinates, and project these reference points to 2D camera views to sample features from the corresponding image feature maps.

Towards this end, multi-view images  $\{I_t^i\}_{i=1}^V$  at frame  $t$  are input into the image backbone to obtain multi-view image features  $\{F_t^i\}_{i=1}^V$ . The spatial encoder then aggregates multi-view features  $\{F_t^i\}_{i=1}^V$  to BEV features  $B_t$ . We adopt the spatial cross-attention (SCA) [18] for this aggregation,

$$\text{SCA}(Q_p, F_t) = \frac{1}{|\mathcal{V}_{\text{hit}}|} \sum_{i \in \mathcal{V}_{\text{hit}}} \sum_{j=1}^{N_{\text{ref}}} \text{Attn}(Q_p, \mathcal{P}(p, i, j), F_t^i). \quad (2)$$

Here,  $Q \in \mathbb{R}^{H \times W \times C}$  represents BEV queries, which are pre-defined learnable parameters used as queries.  $Q_p \in \mathbb{R}^C$  is the query at  $p = (x, y)$  in the BEV plane.  $\mathcal{V}_{\text{hit}}$  represents the camera views, where the projected 2D reference points fall in.  $i$  is the index of the camera view, and  $j$  is the index of the reference points.  $\text{Attn}(\cdot)$  is the deformable attention layer proposed in [46].  $\mathcal{P}(p, i, j)$  is to project the  $j$ -th 3D reference point to the  $i$ -th camera view. The 3D to 2D projection is described in Eq. (1).

Afterward, a temporal BEV encoder [36] is adopted to better incorporate temporal information. Here, a temporal self-attention (TSA) layer [18] is utilized to gather the BEV feature history,

$$\text{TSA}(Q_p, Q, B_{t-1}) = \sum_{S \in Q, B_{t-1}} \text{Attn}(Q_p, p, S), \quad (3)$$

where  $Q_p$  represents the BEV query located at  $p = (x, y)$ .

Finally, a deformable DETR head [46] is followed to predict 3D bounding boxes from BEV features provided by the temporal encoder. We use Focal loss for classification training and  $L_1$  loss for bounding box regression supervision.



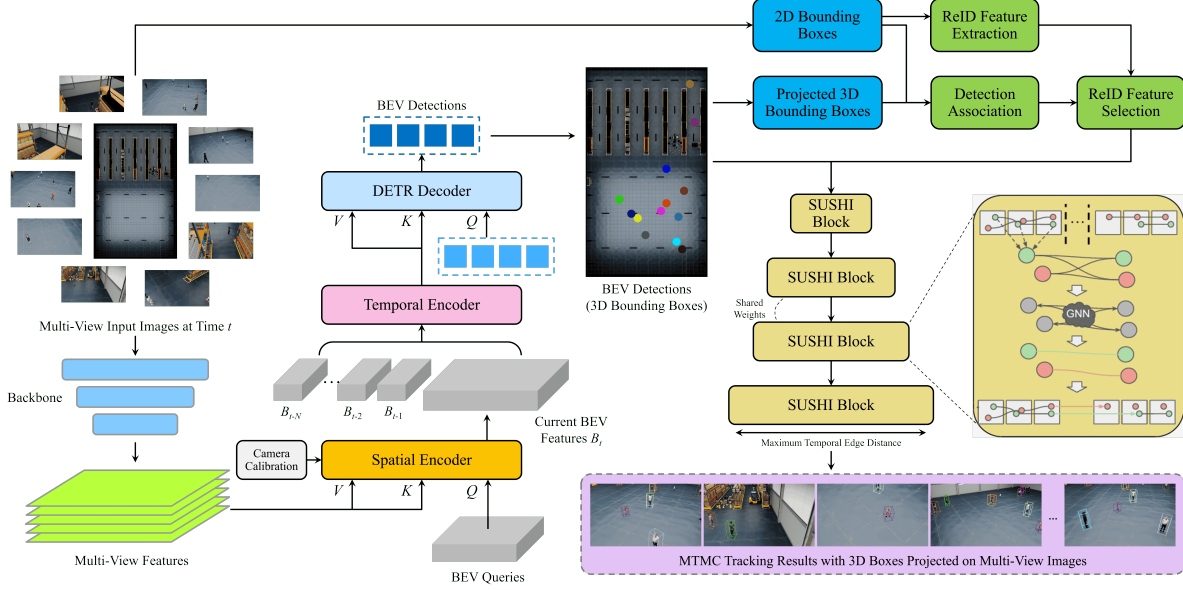


Figure 3. The overall framework of BEV-SUSHI. First, multi-view images at frame  $t$  are passed through the image backbone to obtain multi-view image features. A spatial encoder is then introduced to aggregate multi-view image features to BEV features  $B_t$ , followed by a temporal encoder to aggregate BEV features within a temporal window. A DETR-based decoder is utilized to obtain object detection results, which are in the format of 3D bounding boxes. To get reliable ReID features for the detected objects, a ReID feature extraction module is proposed, including a 2D ReID feature extractor and a 2D-3D detection association algorithm. Finally, SUSHI-3D is designed to achieve multi-object tracking in BEV to obtain the final MTMC tracking results. (*SUSHI Block* graphics are from [5].)

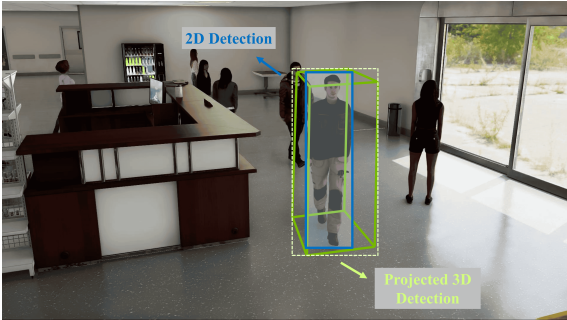


Figure 4. Illustration on 2D detection (in blue) and the corresponding projected 3D detection (in green).

### 3.3. Multi-View ReID Feature Extraction

Appearance features are crucial for tracking tasks, especially for cross-view tracking, since object appearances from different cameras might differ due to varied illuminations and viewing angles. The naïve way to obtain the appearance feature for each 3D detection from BEVFormer is to project the 3D bounding box to each camera image and extract ReID features by the projected box. However, unlike the 2D bounding boxes directly from 2D object detectors, the projected 3D bounding boxes are usually much larger than the actual objects in the image (see Fig. 4). This will include noisy background or other nearby objects in the target object so that affect the quality of the ReID fea-

tures. Therefore, we propose a 2D-3D detection association algorithm to effectively find the best 2D bounding box for each 3D detection. For the 2D detector, we use DINO [40] with FAN-small backbone [44]. The detector is pre-trained on a subset (800k images) of the OpenImages dataset [17], pseudo-labeled by a DINO detector trained on 80 COCO classes. Then, the detector is fine-tuned on a proprietary dataset with more than 1.5 million images and more than 27 million objects for the person class.

Next, we extract ReID features for each 2D bounding box from the 2D detector. Our person-centric ReID model is based on the SOLIDER [7], a self-supervised learning framework with a Swin-Tiny Transformer [19] backbone. The model uses image crops of size  $256 \times 128$  as input and outputs a feature of dimension 256 enhancing memory efficiency and throughput. We pre-train the model on a proprietary dataset of 3 million unlabeled image crops. The supervised fine-tuning is performed on a collection of both real and synthetic datasets which includes characters from Market-1501 [43], AICity’24 [31] and additional proprietary datasets totaling the object count to 3,826 identities and image count to 82k object crops. We analyze the ReID feature quality in Appendix C.

With 3D detections and ReID features of each 2D detection, we then assign each 3D bounding box with one or more ReID features, by introducing a 2D-3D detection association algorithm to create a mapping from the 3D bound-



Scene	Warehouse	Retail store	Hospital	Overall
Accuracy	99.4%	95.2%	91.2%	96.9%

Table 1. Detection association accuracy among different scene types on AICity’24 dataset.

ing box set to the 2D bounding box set. First, we project 3D bounding boxes to 2D in each camera view, then we compute a cost matrix between 2D detections and projected 3D detections by

$$c_{ij} = \begin{cases} \lambda \|\mathbf{bc}_i^{3D} - \mathbf{bc}_j^{2D}\|_2, & \text{if } IOU(\mathbf{b}_i^{3D}, \mathbf{b}_j^{2D}) \geq 0.1, \\ +\infty, & \text{otherwise,} \end{cases} \quad (4)$$

where  $\mathbf{bc}^{2D}$  is the bottom center point of the 2D bounding box,  $\mathbf{bc}^{3D}$  is the bottom center point of the 2D area of the projected 3D bounding box, and  $\lambda$  is a robustness factor to handle 2D occlusions

$$\lambda = \mathbb{1}(v_i^{3D} \geq v_j^{2D}) + \alpha \mathbb{1}(v_i^{3D} < v_j^{2D}). \quad (5)$$

Here, we set  $\alpha > 1$  to penalize if  $\mathbf{b}^{2D}$  is lower than  $\mathbf{b}^{3D}$ . We utilize the Hungarian algorithm to achieve the final assignment. Moreover, the 3D detections can be verified here through the following strategy: a 3D detection will be removed if it cannot be matched with any 2D detections from all camera views. The detailed algorithm is described in Appendix B.1.

We evaluate the 2D-3D detection association algorithm to ensure we have high-quality ReID features assigned to each 3D detection. This evaluation is based on the association results between ground truth 2D and 3D annotations in the AICity’24 training set. We introduce the evaluation metric, detection association accuracy, defined as the number of corrected matches divided by the number of matched ground truth 2D bounding boxes. The evaluation results are shown in Tab. 1. We can find that even in crowded hospital scenes, our algorithm can achieve above 90% accuracy, illustrating its robustness. More visualizations based on predictions are shown in Appendix B.2.

### 3.4. 3D Multi-Object Tracking with GNNs

The early multi-view aggregation of our BEV-SUSHI framework allows us to perform multi-object tracking (MOT) directly in 3D world coordinates  $\mathcal{W}$ . With 3D detections from BEVFormer and features from ReID model, graphs provide a natural framework to address the long-term association challenges in MTMC tracking within a tracking-by-detection setting. Leveraging its unique long-term modeling capabilities, we build upon SUSHI [5], originally designed for 2D tracking, to introduce the first hierarchical NN-based tracking solution in 3D space.

**MOT graph formulation.** In the common graph formulation of multi-object tracking [41], each node  $v \in V$  models a detection connected with edges  $e \in E$  representing association hypotheses in an undirected graph  $G = (V, E)$ . Learning or optimizing the graph connectivity results in object tracks  $\mathcal{T}$  across time. Nodes and edges are initialized with object identity-relevant information projected to the embedding space. These embeddings encode geometrical, appearance, and/or motion features to infer the object identity via graph connection. More specifically, edges contain distance information between their connecting nodes, *e.g.*, cosine distances of ReID features. With the help of message passing, first introduced to MOT by [3], information is shared and distributed across the graph. The procedure iteratively updates the node  $h_i$  and edge  $h_{(i,j)}$  embeddings via the following rules

$$(v \rightarrow e) \quad h_{(i,j)}^{(l)} = \mathcal{N}_e \left( \left[ h_i^{(l-1)}, h_j^{(l-1)}, h_{(i,j)}^{(l-1)} \right] \right), \quad (6)$$

$$(e \rightarrow v) \quad m_{(i,j)}^{(l)} = \mathcal{N}_v \left( \left[ h_i^{(l-1)}, h_{(i,j)}^{(l)} \right] \right), \quad (7)$$

$$h_i^{(l)} = \Phi \left( \left\{ m_{(i,j)}^{(l)} \right\}_{j \in N_i} \right), \quad (8)$$

with  $\mathcal{N}_e$  and  $\mathcal{N}_v$  representing learnable functions,  $[\cdot]$  denotes concatenation,  $N_i \subset V$  is the set of adjacent nodes to  $i$ , and  $\Phi$  denotes summation, maximum or an average. Finally, a binary edge classification and linear program that ensures the network flow integrity (single edge in and out of a node) yields identity-consistent object trajectories.

To model long-term object associations across thousands of frames in a computationally feasible manner, we follow prior work to sparsify the fully connected MOT graph via an initial graph pruning [3] and a hierarchical model architecture [5]. During graph construction, the pruning utilizes the same aforementioned object identity cues to remove unlikely edges, *i.e.*, object associations. This approach not only increases the feasible number of input frames for the graph model but improves overall performance by reducing the imbalance between positive and negative edge classifications. Splitting the input frames into non-overlapping sub-graphs processed by a hierarchy of GNNs boosts the modeling capacities even further. To this end, the graph formulation is extended to nodes representing tracklets and each hierarchy level merges nodes from previous levels, thereby, providing the inputs for subsequent levels.

**Tracking with GNNs in 3D space.** The generalizable and versatile graph structure as well as node and edge feature design in [5] facilitates its application to 3D space. For BEV-SUSHI, we replace the 2D geometry, appearance, and motion feature encodings with MTMC 3D counterparts.

For the *geometry* embeddings, we encode  $xyz$ -center distances for an edge connecting nodes  $i$  and  $j$  via

$$(x_i - x_j, y_i - y_j, z_i - z_j). \quad (9)$$

In contrast to 2D bounding boxes, geometric distances in 3D are not impacted by projection distortions and camera distance scaling. Therefore, we omit distances between object dimensions for our MTMC target datasets, which primarily focus on person tracking.

Our BEV-SUSHI framework extracts *appearance* information in the form of multi-view ReID features obtained via 3D-2D projection and detection association (described in Sec. 3.3). To provide our graph model with view-consistent and stable appearance distances, we compute cosine distances between ReID features averaged across all cameras in which the associated 3D detection is observable.

In contrast to 2D tracking, our strong combination of true-to-scale 3D geometry and multi-view appearance renders the impact of linear *motion* features negligible. Hence, we opted for a more efficient graph formulation without motion features but with larger input time spans. We leave the exploration of more sophisticated motion models for 3D GNN tracking open for future research.

**Long-term tracking with global block.** AICity’24 [31] evaluates identity consistency across sequences with up to 24k frames including occlusions ranging up to 2k frames. To tackle these challenges, we train a GNN with 9 layers to associate objects for up to 3,840 frames (empirical memory limit for a 32 GB GPU) at once and process the entire sequence in a near-online sliding window fashion with stride  $s = 1920$ . The original SUSHI [5] inference computes overlapping graph outputs and associates tracks via hand-crafted, and therefore error-prone, matching heuristics.

In this work, we propose a novel near-online inference of a trained GNN hierarchy that removes potential biases in tracking heuristics and unlocks occlusion handling beyond the number of frames per graph. As depicted in Fig. 5, we omit the graph overlap and process only the incoming frames with the model hierarchy. An additional global block associates any previous track with the predictions on the incoming frames. In practice, we repurpose the last hierarchy level of our model and replace the inference of the graph overlap from previous layers with past tracks. Empirically, we observed that the layer trained for inference of 3,840 frames generalizes to predict more long-term associations without any retraining.

## 4. Experiments

### 4.1. Datasets and Implementations

**AICity’24 Dataset [31]** is an MTMC tracking benchmark consisting of 6 different synthetic environments, *e.g.*, warehouses, retail stores, and hospitals, developed using the NVIDIA Omniverse Platform [22]. The dataset includes 90 scenes, 40 for training, 20 for validation, and 30 for testing, with a total of 953 cameras, 2,491 people, and over 100 million bounding boxes. Besides, we also include 30 addi-

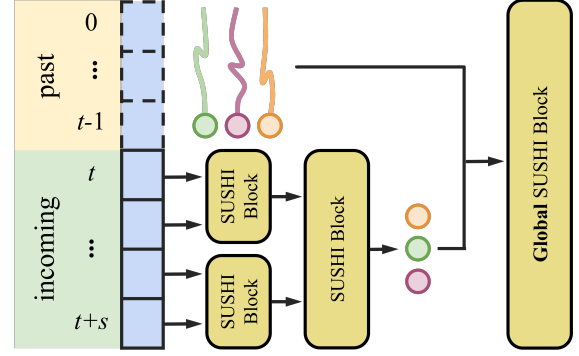


Figure 5. The tracking inference of BEV-SUSHI. We process long sequences in a near-online fashion with stride  $s$ . But in contrast to [5], we omit overlaps between windows and process only incoming frames with our graph neural network (GNN) model. To associate incoming with past tracks, BEV-SUSHI extends the input time span of the final hierarchy layer to function as a global merging block. The global block requires no additional training.

tional scenes with similar camera settings for BEVFormer training. In this dataset, persons are annotated as 3D bounding boxes with 3D locations, dimensions, and orientations.

**WildTrack Dataset [6]** is a real-world MTMC tracking benchmark containing a single sequence of a scene where 7 synchronized cameras cover over a  $36 \times 12$  m space at a  $1920 \times 1080$  resolution. There are 400 fully annotated frames per camera at 2 FPS with a total of 313 identities and 42,721 2D bounding boxes. As for the 3D annotations, there is no 3D bounding box available. A grid is defined on the ground plane and all persons are annotated by the corresponding grid indices, which can be mapped to 3D locations on the ground plane in meters.

Following the experiment setting in [8, 11, 29], we use the first 360 annotated frames as the training data and evaluate the rest 40 frames. Since there is no 3D bounding box annotation in the WildTrack dataset, we remove the bounding box dimension and orientation regression losses during BEVFormer training. We use a default 3D bounding box dimension (width, length, height) of  $[0.6, 0.6, 1.7]$  and zero rotation for later 2D-3D detection association.

**BEVFormer adaption for MTMC environments.** The original BEVFormer is developed for autonomous driving applications with moving cameras but fixed camera settings, *i.e.*, a fixed number of cameras and fixed relative camera placement. For MTMC environments, we usually have dynamic camera systems across different scenes. AICity’24 dataset, for example, has 6 different indoor environments with cameras ranging from 7 to 16. Therefore, we re-design the spatial cross-attention layers to enable a dynamic number of cameras for different scenes and remove camera embeddings. Besides, we apply a transformation on the original 3D world coordinates to shift the origin to the center of

Method	HOTA↑	DetA↑	AssA↑	LocA↑
Asilla [30]	40.34	53.80	32.50	89.57
ARV [28]	51.06	54.85	48.07	89.61
UW-ETRI [37]	57.14	59.88	54.80	91.24
FraunhoferIOSB [27]	60.88	69.54	55.20	87.97
Nota [16]	60.93	68.37	54.96	90.62
SJTU-Lenovo [33]	67.22	<u>84.03</u>	55.06	<u>93.82</u>
Yachiyo [38]	<u>71.94</u>	72.10	<u>71.81</u>	88.39
BEV-SUSHI (Ours)	<b>81.22</b>	<b>86.94</b>	<b>76.19</b>	<b>95.67</b>

Table 2. Results on AICity’24 test set. The first place is in **bold**, and the second place is underlined.

the scene (see Appendix A.2). Moreover, we implement an additional Circle Non-Maximum Suppression (NMS) with a threshold of 0.2m to filter nearby false positives.

## 4.2. Results on AICity’24 Dataset

The AICity’24 benchmark adopts the Higher Order Tracking Accuracy (HOTA) scores [20] for evaluation. HOTA is computed on the 3D locations of objects, with repetitive data points removed across cameras for the same frame. Euclidean distances between predicted and ground truth 3D locations are converted to similarity scores using a zero-distance parameter. These scores contribute to the calculation of localization accuracy (LocA), detection accuracy (DetA), and association accuracy (AssA).

Compared with other methods on the AICity’24 benchmark leaderboard in Tab. 2, our proposed BEV-SUSHI outperforms significantly in all aspects, with the following improvements from the second place, +9.28 on HOTA, +2.91 on DetA, +4.38 on AssA, and +1.85 on LocA. The qualitative results on the AICity’24 dataset are shown in Fig. 6. We project 3D detections with track IDs to each camera view and plot their 3D locations on the floor plans. It shows that BEV-SUSHI achieves robust MTMC detection and tracking accuracies in various indoor environments.

## 4.3. Results on WildTrack Dataset

The WildTrack dataset evaluates tracking results in the ground plane with a 1m threshold for GT assignment. The primary metrics are IDF1 [24], Multi-Object Tracking Accuracy (MOTA), and Mutli-Object Tracking Precision (MOTP) [2]. Furthermore, we report the number of Mostly Tracked (MT) and Lost (ML) tracks in percentages.

As illustrated in Tab. 3, our BEV-SUSHI framework achieves state-of-the-art for IDF1 with a +1.1 improvement over EarlyBird [29]. However, we get a slightly lower MOTA compared to EarlyBird. This is due to the limited training data (360 frames) provided by the WildTrack dataset. Transformer-based networks, like BEVFormer, generally require larger amounts of data to unfold their potential. To demonstrate the effectiveness of BEV-SUSHI’s

Method	IDF1↑	MOTA↑	MOTP↑	MT↑	ML↓
KSP-DO [6]	73.2	69.6	61.5	28.7	25.1
KSP-DO-ptrack [6]	78.4	72.2	60.3	42.1	14.6
GLMB-YOLOv3 [23]	74.3	69.7	73.2	79.5	21.6
GLMB-DO [23]	72.5	70.1	63.1	<b>93.6</b>	22.8
DMCT [39]	77.8	72.8	79.1	61.0	<u>4.9</u>
DMCT Stack [39]	81.9	74.6	78.9	65.9	<u>4.9</u>
ReST [8]	86.7	84.9	84.1	87.8	<u>4.9</u>
EarlyBird [29]	92.3	<u>89.5</u>	86.6	78.0	<u>4.9</u>
BEV-SUSHI (Ours)	<u>93.4</u>	87.5	<b>94.3</b>	<u>90.2</u>	<b>2.4</b>
BEV-SUSHI <sup>†</sup> (Ours)	<b>95.6</b>	<b>92.6</b>	<u>93.7</u>	80.5	7.3

Table 3. Results on WildTrack test set. The first place is in **bold**, and the second place is underlined. † uses the same detections as EarlyBird [29].

Backbone	BEV Reso.	Voxel Size	Epochs	mAP
ResNet-50	50×50	2.0 m	24	83.14
ResNet-101	200×200	0.5 m	24	88.64
ResNet-101	200×200	0.5 m	48	95.36

Table 4. Ablation studies on the configurations of our multi-view object detector. The evaluation is done on the customized validation set of the AICity’24 dataset.

ReID and tracking components, we replace our detections with the ones from EarlyBird (based on MVDet [11] with ResNet-18 backbone). This version (†) yields significant improvements on both IDF1 (+3.3) and MOTA (+3.1). Qualitative results on WildTrack are shown in Fig. 7.

Besides, we pre-train BEVFormer on the large-scale synthetic dataset (AICity’24) and finetune on the small WildTrack dataset. The results are analyzed in Appendix A.3.

## 4.4. Ablation Studies

**Multi-view object detector configuration analysis.** To verify we have the best BEVFormer model for multi-view 3D object detection, we conduct several ablation studies on backbone selection and parameter tuning to improve the detection results. The results are shown in Tab. 4. We use a customized validation set from the AICity’24 dataset, including around 3k frames from warehouse, hospital, and retail store scenes. More results shown in Appendix A.1.

**Long-term association analysis.** The length of sequences and occlusion gaps common in MTMC benchmarks without full camera coverage of the environment pose unique challenges to our tracking model and inference. To process arbitrary long sequences, BEV-SUSHI infers windows of frames with a GNN model and connects these in time in a near-online sliding window fashion. In Tab. 5, we demonstrate how increasing the window size gradually improves the association performance, *i.e.*, AssA. Larger windows shift responsibility from the heuristic matching



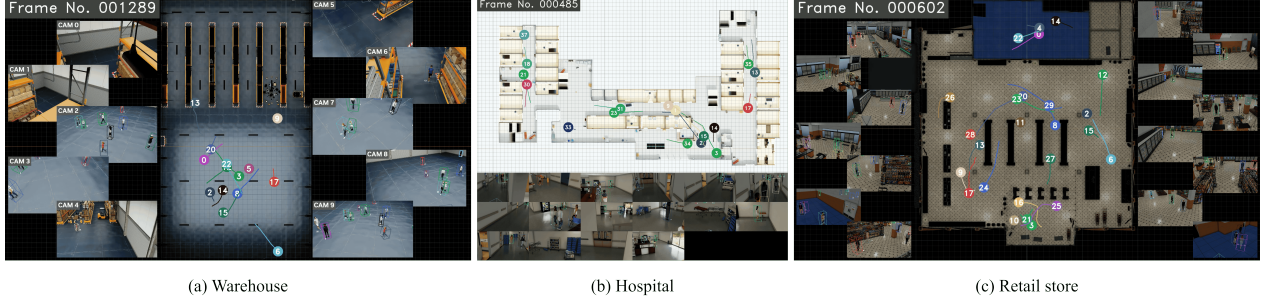


Figure 6. Visualization of MTMC detection and tracking results for three different scenes in AICity’24 test set. The tracked objects are shown as colored dots in the BEV floor plans, and object 3D bounding boxes are projected and drawn in each camera view.

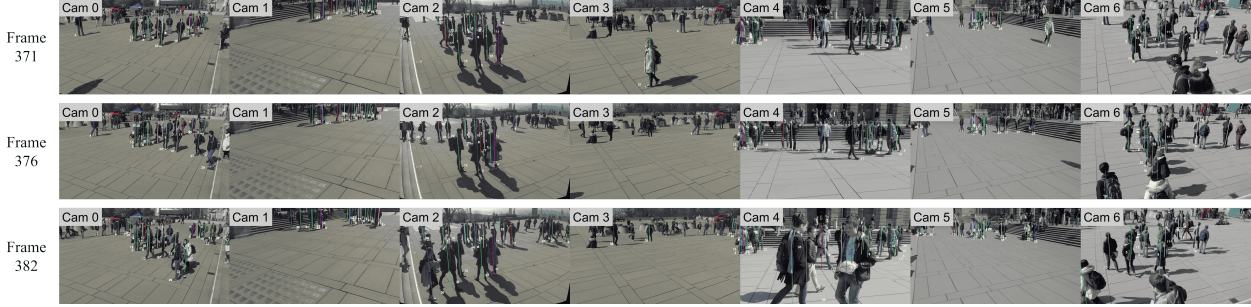


Figure 7. Visualization of MTMC detection and tracking results on WildTrack test set. The tracked objects are shown as colored poles projected in each camera view. The poles are defined by their 3D foot points and a pre-defined human height (1.7m).

Win Size	Global	HOTA	DetA	AssA	LocA
480		51.87	86.12	31.26	95.55
960		62.50	86.36	45.25	95.57
1920		71.88	86.47	59.78	95.58
3840		76.80	86.49	68.22	95.58
3840	✓	81.22	86.94	76.19	95.67

Table 5. SUSHI inference ablation on the AICity’24 test set. Rows 1 to 4 rely on heuristic matching to track overlapping sliding windows inferred by the SUSHI GNN hierarchy. Our final tracking solution (the last row) applies a global merging block to associate past tracks with incoming frames without graph overlaps. The window stride is half of the window size for all experiments.

between windows to the trained model, which allows us to bridge longer occlusion gaps. The last row replaces all remaining heuristics with our novel global merging block boosting overall performance by +4.42 HOTA.

To further illustrate the superiority of our learned GNN approach, we evaluate progressively longer subsets of a selected AICity’24 scene in Tab. 6. The BEV-KF baseline represents heuristic online trackers like EarlyBird [29] commonly only evaluated on short tracking challenges like WildTrack. From 1,000 frames to the full AICity’24 length, BEV-SUSHI only drops by 4.58 HOTA points compared to the massive performance decrease of 43.35 utilizing only a Kalman filter. Such trackers fail to tackle the long-term tracking challenges common in MTMC benchmarks.

Method	# of frames	HOTA	DetA	AssA	LocA
BEV-SUSHI	1,000	95.09	94.41	95.77	97.51
	3,000	93.62	94.89	92.37	97.53
	5,000	91.85	94.98	88.82	97.50
	10,000	91.06	93.89	88.33	97.22
	23,994	90.51	93.03	88.05	97.16
BEV-KF (baseline)	1,000	63.49	89.01	45.29	97.64
	3,000	45.28	89.71	22.86	97.46
	5,000	35.54	88.79	14.23	97.30
	10,000	27.91	87.80	8.87	97.19
	23,994	20.15	86.91	4.67	97.11

Table 6. Tracking comparisons with increasing video lengths from a warehouse scene in AICity’24 dataset. The baseline BEV-KF is based on the detection results from BEVFormer processed by a Kalman filter based tracker [32] used in EarlyBird [29].

## 5. Conclusion

This paper proposed BEV-SUSHI, an accurate and robust multi-camera 3D detection and tracking framework, based on early multi-view aggregation, in environments monitored by static multi-camera systems. The proposed framework has impressive generalizability for diverse scenes and camera settings. It also has a powerful capability for long-term association to track objects in very long videos. As a result, BEV-SUSHI achieves SOTA on both the AICity’24 dataset and the WildTrack dataset.

## References

- [1] Pierre Baqué, François Fleuret, and Pascal Fua. Deep occlusion reasoning for multi-camera multi-target detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 271–279, 2017. 2
- [2] Keni Bernardin and Rainer Stiefelhausen. Evaluating multiple object tracking performance: the clear mot metrics. *EURASIP Journal on Image and Video Processing*, 2008:1–10, 2008. 7
- [3] Guillem Brasó and Laura Leal-Taixé. Learning a neural solver for multiple object tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 5
- [4] Michael Bredebeck, Xiaoyan Jiang, Marco Körner, and Joachim Denzler. Data association for multi-object tracking-by-detection in multi-camera networks. In *2012 Sixth International Conference on Distributed Smart Cameras (ICDSC)*, pages 1–6. IEEE, 2012. 2
- [5] Orcun Cetintas, Guillem Brasó, and Laura Leal-Taixé. Unifying short and long-term tracking with graph hierarchies. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22877–22887, 2023. 2, 4, 5, 6
- [6] Tatjana Chavdarova, Pierre Baqué, Stéphane Bouquet, Andrii Maksai, Cijo Jose, Timur Bagautdinov, Louis Lettry, Pascal Fua, Luc Van Gool, and François Fleuret. Wildtrack: A multi-camera hd dataset for dense unsupervised pedestrian detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5030–5039, 2018. 2, 6, 7
- [7] Weihua Chen, Xianzhe Xu, Jian Jia, Hao Luo, Yaohua Wang, Fan Wang, Rong Jin, and Xiuyu Sun. Beyond appearance: a semantic controllable self-supervised learning framework for human-centric visual tasks. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 4
- [8] Cheng-Che Cheng, Min-Xuan Qiu, Chen-Kuo Chiang, and Shang-Hong Lai. Rest: A reconfigurable spatial-temporal graph model for multi-camera multi-object tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10051–10060, 2023. 1, 3, 6, 7
- [9] Adam Coates and Andrew Y Ng. Multi-camera object detection for robotics. In *2010 IEEE International conference on robotics and automation*, pages 412–419. IEEE, 2010. 2
- [10] Ran Eshel and Yael Moses. Homography based multiple camera detection and tracking of people in a dense crowd. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. 2, 3
- [11] Yunzhong Hou and Liang Zheng. Multiview detection with shadow transformer (and view-coherent data augmentation). In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1673–1682, 2021. 2, 6, 7
- [12] Hung-Min Hsu, Yizhou Wang, and Jenq-Neng Hwang. Traffic-aware multi-camera tracking of vehicles based on reid and camera link model. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 964–972, 2020. 2, 3
- [13] Hung-Min Hsu, Jiarui Cai, Yizhou Wang, Jenq-Neng Hwang, and Kwang-Ju Kim. Multi-target multi-camera tracking of vehicles using metadata-aided re-id and trajectory-based camera link model. *IEEE Transactions on Image Processing*, 30:5198–5210, 2021. 3
- [14] Weiming Hu, Min Hu, Xue Zhou, Tieniu Tan, Jianguang Lou, and Steve Maybank. Principal axis-based correspondence between multiple cameras for people tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(4):663–671, 2006. 2, 3
- [15] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960. 3
- [16] Jeongho Kim, Wooksu Shin, Hanchol Park, and Donghyuk Choi. Cluster self-refinement for enhanced online multi-camera people tracking. In *CVPR Workshop*, Seattle, WA, USA, 2024. 7
- [17] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International journal of computer vision*, 128(7):1956–1981, 2020. 4
- [18] Zhiqi Li, Wenhao Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *European conference on computer vision*, pages 1–18. Springer, 2022. 2, 3, 11
- [19] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 4
- [20] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. Hota: A higher order metric for evaluating multi-object tracking. *International Journal of Computer Vision*, pages 1–31, 2020. 7
- [21] Duy MH Nguyen, Roberto Henschel, Bodo Rosenhahn, Daniel Sonntag, and Paul Swoboda. Lmgp: Lifted multicut meets geometry projections for multi-camera multi-object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8866–8875, 2022. 3
- [22] NVIDIA. Nvidia omniverse platform. <https://www.nvidia.com/en-us/design-visualization/omniverse/>, 2021. Accessed: Apr 10, 2023. 1, 6
- [23] Jonah Ong, Ba-Tuong Vo, Ba-Ngu Vo, Du Yong Kim, and Sven Nordholm. A bayesian filter for multi-view 3d multi-object tracking with occlusion handling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(5):2246–2263, 2020. 7
- [24] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *European conference on computer vision*, pages 17–35. Springer, 2016. 7

- [25] Gemma Roig, Xavier Boix, Horeh Ben Shitrit, and Pascal Fua. Conditional random fields for multi-camera object detection. In *2011 International Conference on Computer Vision*, pages 563–570. IEEE, 2011. 2
- [26] Aswin C Sankaranarayanan, Ashok Veeraraghavan, and Rama Chellappa. Object detection, tracking and recognition for multiple smart cameras. *Proceedings of the IEEE*, 96(10):1606–1624, 2008. 2
- [27] Andreas Specker. Ocmctrack: Online multi-target multi-camera tracking with corrective matching cascade. In *CVPR Workshop*, Seattle, WA, USA, 2024. 7
- [28] Vasin Suttichaya, Riu Cherdchusakulchai, Sasin Phimsiri, Visarut Trairattanapa, Suchat Tungjitnob, Wasu Kudis-thalert, Pornprom Kiawjak, Ek Thamwiwatthana, Phawat Borisuitsawat, Teepakorn Tosawadi, Pakcheera Choppradit, Kasisdis Mahakijdechachai, Supawit Vatathanavaro, and Worawit Saetan. Online multi-camera people tracking with spatial-temporal mechanism and anchor-feature hierarchical clustering. In *CVPR Workshop*, Seattle, WA, USA, 2024. 7
- [29] Torben Teepe, Philipp Wolters, Johannes Gilg, Fabian Herzog, and Gerhard Rigoll. Earlybird: Early-fusion for multi-view tracking in the bird’s eye view. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 102–111, 2024. 1, 3, 6, 7, 8
- [30] Huan Vi and Lap Quoc Tran. Efficient online multi-camera tracking with memory-efficient accumulated appearance features and trajectory validation. In *CVPR Workshop*, Seattle, WA, USA, 2024. 7
- [31] Shuo Wang, David C Anastasiu, Zheng Tang, Ming-Ching Chang, Yue Yao, Liang Zheng, Mohammed Shaiqur Rahman, Meenakshi S Arya, Anuj Sharma, Pranamesh Chakraborty, et al. The 8th ai city challenge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7261–7272, 2024. 2, 4, 6
- [32] Zhongdao Wang, Liang Zheng, Yixuan Liu, and Shengjin Wang. Towards real-time multi-object tracking. *The European Conference on Computer Vision (ECCV)*, 2020. 8
- [33] Zhenyu Xie, Zelin Ni, Wenjie Yang, Yuang Zhang, Yihang Chen, Yang Zhang, and Xiao Ma. A robust online multi-camera people tracking system with geometric consistency and state-aware re-id correction. In *CVPR Workshop*, Seattle, WA, USA, 2024. 7
- [34] Yuanlu Xu, Xiaobai Liu, Yang Liu, and Song-Chun Zhu. Multi-view people tracking via hierarchical trajectory composition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4256–4265, 2016. 1, 3
- [35] Yuanlu Xu, Xiaobai Liu, Lei Qin, and Song-Chun Zhu. Cross-view people tracking by scene-centered spatio-temporal parsing. In *Proceedings of the AAAI conference on artificial intelligence*, 2017. 1, 3
- [36] Chenyu Yang, Yuntao Chen, Hao Tian, Chenxin Tao, Xizhou Zhu, Zhaoxiang Zhang, Gao Huang, Hongyang Li, Yu Qiao, Lewei Lu, et al. Bevformer v2: Adapting modern image backbones to bird’s-eye-view recognition via perspective supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17830–17839, 2023. 2, 3, 11
- [37] Cheng-Yen Yang, Hsiang-Wei Huang, Pyong-Kun Kim, Zhongyu Jiang, Kwang-Ju Kim, Chung-I Huang, Haiqing Du, and Jenq-Neng Hwang. An online approach and evaluation method for tracking people across cameras in extremely long video sequence. In *CVPR Workshop*, Seattle, WA, USA, 2024. 7
- [38] Ryuto Yoshida, Junichi Okubo, Junichiro Fujii, Masazumi Amakata, and Takayoshi Yamashita. Overlap suppression clustering for offline multi-camera people tracking. In *CVPR Workshop*, Seattle, WA, USA, 2024. 7
- [39] Quanzeng You and Hao Jiang. Real-time 3d deep multi-camera tracking. *arXiv preprint arXiv:2003.11753*, 2020. 7
- [40] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. 4
- [41] Li Zhang, Yuan Li, and Ramakant Nevatia. Global data association for multi-object tracking using network flows. In *CVPR*, 2008. 5
- [42] Yifu Zhang, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. Fairmot: On the fairness of detection and re-identification in multiple object tracking. *International journal of computer vision*, 129:3069–3087, 2021. 3
- [43] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *Computer Vision, IEEE International Conference on*, 2015. 4
- [44] Daquan Zhou, Zhiding Yu, Enze Xie, Chaowei Xiao, Animesh Anandkumar, Jiashi Feng, and Jose M Alvarez. Understanding the robustness in vision transformers. In *International Conference on Machine Learning*, pages 27378–27394. PMLR, 2022. 4
- [45] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. 3
- [46] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations*, 2022. 2, 3



# BEV-SUSHI:

## Multi-Target Multi-Camera 3D Detection and Tracking in Bird’s-Eye View

### Supplementary Material

#### A. Experiments on Mutli-View 3D Detector

##### A.1. 3D Object Detector Comparisons

We conducted experiments on different settings for our 3D object detector in Tab. 7. We considered both BEVFormer [18] and BEVFormer v2 [36] with ResNet-50, ResNet-101, and V2-99 backbones. The maximum number of cameras for training is reported for each experiment due to the limitation of H100 GPU memory. All experiments were trained for 24 epochs with a learning rate of  $2 \times 10^{-4}$ .

Method	Backbone	Max # of Cameras	mAP
BEVFormer [18]	ResNet-50	16	83.14
BEVFormer [18]	ResNet-101	15	88.64
BEVFormer v2 [36]	ResNet-50	15	82.78
BEVFormer v2 [36]	ResNet-101	15	85.03
BEVFormer v2 [36]	V2-99	14	79.95

Table 7. 3D object detection results with different detectors on a customized AICity’24 validation set.

Compared with BEVFormer, BEVFormer v2 receives relatively lower mAP by adding the perspective supervision. Therefore, the perspective supervision may not be helpful for our MTMC application for model convergence. As for the V2-99 backbone, we need to decrease the number of camera views during the training to fit our GPU memory of around 80 GB. This will downgrade the detection performance significantly. In the future, we will improve the memory efficiency to make it possible to utilize larger image backbones.

##### A.2. Scene Re-Centering for BEVFormer

The definition of the BEV coordinate system is important for BEVFormer training. In the original autonomous driving settings, the origin is located on the ego-vehicle, which is the center of the area to be perceived. In our MTMC settings, we define the origins of the multi-camera scenes as the centers of the floor plans and transform the annotations and calibration matrices to the newly defined BEV coordinates. We call this step “re-centering”.

In Tab. 8, we evaluate the model performance on the WildTrack dataset before and after this re-centering step. Before the re-centering, the origin was defined at the corner of a scene. We notice that re-centering can dramatically improve the detection performance by +22.33 mAP, especially for those objects farther from the origin.

Method	mAP
Baseline	66.03
+ re-centering	88.36
+ pre-training	92.03

Table 8. A comparison of detection results on the WildTrack dataset with *re-centering* and *pre-training*.

##### A.3. Pre-Training on AICity’24 Dataset

Since WildTrack is a small dataset with only 400 frames in total, we considered training BEVFormer with a pre-trained model on the AICity’24 dataset, which is a much larger dataset with various scenes. As shown in Tab. 8, this pre-training leads to +3.67 detection performance improvement on the WildTrack test set. This also illustrates the important role of large and well-annotated synthetic datasets in boosting the performance on limited real data.

#### B. Detection Association Algorithm

##### B.1. Algorithm Details

---

###### Algorithm 1: 2D-3D detection association

---

**Input** : 2D detection set  $\mathcal{D}^v$  from camera  $v$ ; 3D detection set  $\mathcal{E}$  from BEVFormer with all camera views; projection matrix  $\mathcal{P}^v$  of camera  $v$ .

**Output**: Mapping of indices from  $\mathcal{E}$  to  $\mathcal{D}^v$ .

```

1  $\mathcal{E} \leftarrow$  filter  $\mathcal{E}$  by confidence score;
2  $\mathcal{E} \leftarrow$  CircleNMS( $\mathcal{E}, \delta$ );
   // optional,  $\delta$ : NMS threshold
3  $\mathcal{E}^v \leftarrow \mathcal{P}^v(\mathcal{E})$ ; // projected 3D boxes
4 for camera  $v$  to  $V$  do
5   Initialize the cost matrix  $\mathbf{c}^v = [c_{ij}^v]$  as zeros;
6   for  $\mathbf{b}_i^{3D}$  from  $\mathcal{E}^v$  do
7     for  $\mathbf{b}_j^{2D}$  from  $\mathcal{D}^v$  do
8        $c_{ij}^v \leftarrow$  compute cost by Eq. (4);
9     end
10  end
11  Matches  $m^v \leftarrow$  Hungarian( $\mathbf{c}^v, \Delta$ );
   //  $\Delta$ : cost threshold
12 end
```

---

The detailed 2D-3D detection association algorithm is shown in Algorithm 1. We set the threshold for CircleNMS to  $\delta = 0.2\text{m}$  and set the cost threshold to  $\Delta = 150$ .

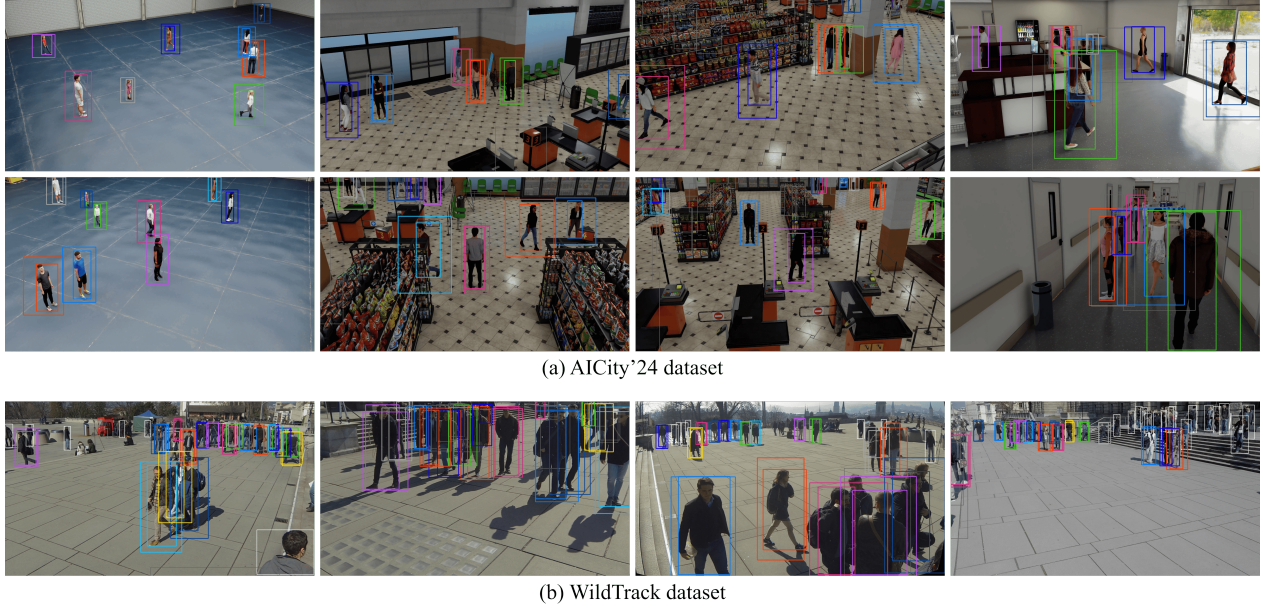


Figure 8. Visualization of 2D-3D detection association results.

## B.2. Visualization

We visualized some 2D-3D detection association results on sample frames of the AICity'24 and WildTrack datasets in Fig. 8. The associated bounding boxes are in the same color, where the smaller ones are 2D detections and the larger ones are projected 3D detections from BEVFormer. Those 2D detections in white are not associated with any 3D detections.

## B.3. Improvements with Detection Association

We compared the tracking performance of BEV-SUSHI with and without the proposed 2D-3D detection association algorithm in Tab. 9. The baseline result is based on the ReID features extracted from the large projected 3D bounding boxes shown in Fig. 4. With the noisy background or other objects included in the image crops, ReID feature quality will be significantly affected.

Method	IDF1	MOTA	MOTP	MT	ML
Baseline	63.2	73.4	93.7	24.0	4.0
+ det association	93.4	87.5	94.3	90.2	2.4

Table 9. A comparison of results on the WildTrack test set with our 2D-3D detection association algorithm.

## C. ReID Feature Quality Analysis

We conducted ReID feature quality analysis on both the AICity'24 and WildTrack datasets. For the AICity'24 dataset, we sampled 500 characters with their 2D bounding

boxes and object IDs from the ground truth across all scenes and cameras from the test set. The total object image crop count is 40,000. We filtered out 2D bounding boxes that are smaller than 5,000 pixels, as well as those whose aspect ratio (*i.e.*, width / height) is less than 0.15. Similarly, for the Wildtrack dataset, we sampled 330 characters from the sequence and applied the same filters bringing the total object crop count to 41,284.

Dataset	Rank-1	Rank-5	Rank-10	mAP
AICity'24	95.02	97.44	98.08	73.85
WildTrack	77.18	84.49	87.97	63.11

Table 10. Evaluation on our ReID feature quality.

We evaluated the ReID feature quality by the mean average precision (mAP), rank-1, rank-5, and rank-10 accuracies. The evaluation results are shown in Tab. 10. We found that the feature quality on the WildTrack dataset is worse than that on the AICity'24 dataset. This is because i) WildTrack is a real-world dataset with more noises and diverse illuminations from different camera views; ii) 2D bounding box annotations are not as accurate as the synthetic AICity'24 dataset. Nevertheless, our BEV-SUSHI achieved impressive results on WildTrack based on these ReID features.

## D. Overall Visualization

We also visualized MTMC detection and tracking results of BEV-SUSHI on the AICity'24 and WildTrack datasets. Please find the demos in the attachment.