

Multi-Camera Trajectory Forecasting with Trajectory Tensors

Olly Styles, Tanaya Guha, and Victor Sanchez

Abstract—We introduce the problem of multi-camera trajectory forecasting (MCTF), which involves predicting the trajectory of a moving object across a network of cameras. While multi-camera setups are widespread for applications such as surveillance and traffic monitoring, existing trajectory forecasting methods typically focus on single-camera trajectory forecasting (SCTF), limiting their use for such applications. Furthermore, using a single-camera limits the field-of-view available, making long-term trajectory forecasting impossible. We address these shortcomings of SCTF by developing an MCTF framework that simultaneously uses object locations from several viewpoints and predicts the object's future location in all possible viewpoints. Our framework follows a *Which-When-Where* approach that predicts in *which* camera(s) the objects appear and *when* and *where* within the camera views they appear. To this end, we propose the concept of *trajectory tensors*: a new technique to encode trajectories across multiple camera views and the associated uncertainties. We develop several encoder-decoder MCTF models for trajectory tensors and present extensive experiments on our own database (comprising 600 hours of video data from 15 camera views) created particularly for the MCTF task. Results show that our trajectory tensor models outperform coordinate trajectory-based MCTF models and existing SCTF methods adapted for MCTF.

Index Terms—Trajectory forecasting, multi-camera tracking, person re-identification, multi-camera trajectory forecasting.

1 INTRODUCTION

PREDICTING the future trajectory of objects in videos is a challenging problem with multiple application domains such as intelligent surveillance [1], autonomous driving [2], person re-identification (RE-ID) [3], and traffic monitoring [4]. Existing works in this area focus on single-camera trajectory forecasting (SCTF), i.e., predicting a future trajectory of an object within the same camera view in which the object is observed [2], [5]–[11]. A critical drawback of the single-camera settings is that models cannot anticipate when new objects will enter the scene [12], as they are limited to the data from a single camera. Furthermore, SCTF methods are only suitable for short-term trajectory forecasting, typically 1 to 5 seconds [6]–[11], due to the limited field-of-view. The constraint of a single camera viewpoint must be removed to overcome these issues. To this end, we introduce *multi-camera trajectory forecasting* (MCTF) - a new framework within trajectory prediction. Given the information about an object's location in one or more camera views, we want to predict its future location across a camera network, in all possible camera views. Fig. 1 presents an overview of the MCTF framework.

Tracking an object-of-interest across a large camera network requires simultaneously running state-of-the-art algorithms for object detection, tracking, and RE-ID. These algorithms can be computationally demanding, particularly for large camera networks. Processing videos at a lower image resolution or frame-rate may reduce the computational demands, but this often results in missed detections [13]. Alternatively, we may choose to monitor only a subset of the cameras in a network, which may also result in missed detections. A successful MCTF model can mitigate this issue

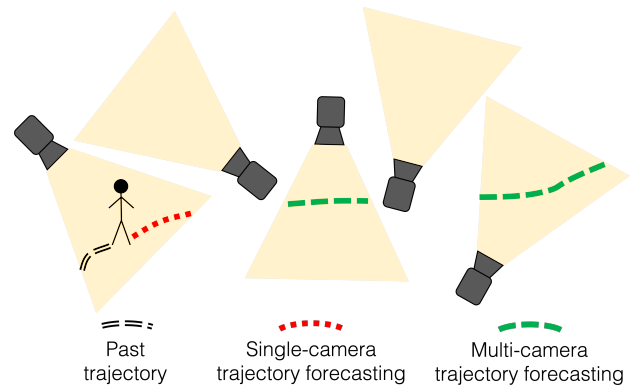


Fig. 1. **Multi-camera trajectory forecasting.** We introduce a novel formulation of the trajectory forecasting task which considers multiple camera views. Companion video: <https://youtu.be/H0gYJ-z8mVc>

by preempting an object's location in a distributed camera network, thereby allowing the system to monitor fewer cameras through an intelligent selection technique. This is distinct from previous works that use trajectory information for person RE-ID [3] or vehicle tracking [14] which are *reactive* to observations *after* an object has been observed in multiple camera views. Our proposed MCTF framework is *proactive* - it predicts the future location of an object *before* it enters the camera view.

We present a deep encoder-decoder approach to MCTF that introduces the idea of *trajectory tensors* - a new technique to encode trajectories across multiple camera views and the associated uncertainty in future locations. Trajectory tensors are an attractive alternative to the coordinate-based trajectory, which is the de facto representation in existing trajectory forecasting works [6]–[11], with only a few recent exceptions [15]. Coordinate trajectories represent

• The authors are with the Department of Computer Science, University of Warwick, United Kingdom. Contact email: o.c.styles@warwick.ac.uk

Manuscript received April 19, 2005; revised August 26, 2015.

the historical and predicted future locations of an object as a sequence of coordinates. This representation suffers from several drawbacks, such as the inability to represent multiple possible future locations and their associated uncertainty, both of which are intrinsic to the trajectory forecasting problem. Furthermore, coordinate trajectories can not be easily extended to multiple camera viewpoints. In contrast to the coordinate approach, proposed trajectory tensors divide viewpoints from several cameras into grid cells with values indicating an object's presence or absence. This representation enables us to intuitively model diverse future locations, their associated uncertainty, and object locations in an arbitrary number of viewpoints. Trajectory tensors also offer easy to interpret results which can be visualized easily.

Given an object tracklet in one or more camera view(s), our MCTF framework comprises the following three tasks: (i) In *which* cameras will the object appear next? (ii) *When* will the object appear in those cameras? (iii) *Where* will the object appear in the identified camera views? We use an object tracklet of 2 seconds to predict 12 seconds into the future. Owing to the wide body of complementary literature on pedestrian detection [16] and person RE-ID [17], we focus on pedestrians for our MCTF task. Nevertheless, the task, along with our proposed data representation and model, can be easily generalized to any moving object.

MCTF and the Warwick-NTU Multi-camera Forecasting (WNMF) dataset were first introduced in our earlier work [18]. This work extends our previous work as follows:

- We introduce three problem formulations within the MCTF framework: *which*, *when*, and *where* (Section 3). Each task predicts trajectories with more granularity than the previous.
- We propose trajectory tensors, a new data representation for multi-camera trajectories (Section 4.1). Trajectory tensors overcome many of the shortcomings of coordinate trajectories for MCTF.
- We present a deep encoder-decoder for MCTF (Section 4.2). Our novel approach, based on trajectory tensors, outperforms several relevant baselines (Section 5).

2 RELATED WORK

To the best of our knowledge, our work is the first to consider trajectory forecasting in a multi-camera setting. Therefore, we review literature in the most related fields.

Single-camera trajectory forecasting. Trajectory forecasting has seen considerable attention in recent years, generally focusing on forecasting pedestrian trajectories from a bird's-eye viewpoint [19], [20]. Many approaches use recurrent models such as Long Short-Term Memory (LSTM) networks, and focus on extracting features to encode social norms such as avoiding collisions with others and group movements [6]–[8]. Other works also consider environmental constraints [21]. A comparatively small number of works consider visual features for trajectory forecasting, such as those extracted from human pose [9], [22] or optical flow [10] from non-nadir viewpoints. Models are often optimized using loss functions such as mean squared error (MSE) [9], [10], [23]. As the distribution of future trajectories is multi-modal, using the MSE loss heavily penalizes reasonable but

incorrect forecasts, such as turning right rather than left at an intersection. Rather than using MSE, some works use an adversarial loss to train models capable of generating multiple possible future paths using generative adversarial networks (GANs) [7], [8] or variational auto-encoders (VAEs) [2], [24]. However, evaluating multi-output models can be challenging, as typically only a single ground truth trajectory exists despite multiple plausible futures. Liang et al. overcome this issue by generating a simulated dataset where each trajectory has multiple futures [15]. Note that all of the aforementioned works consider only a single camera viewpoint at a time.

Person RE-ID. Person RE-ID is the task of identifying an individual in a set of gallery images given a probe image, often in a multi-camera setting. Substantial progress has been made in this area over recent years. For example, on the commonly used Market [25] benchmark, rank 1 person RE-ID performance (i.e. correct RE-ID given 1 guess only) has improved from 47.3% in 2015 [25] to 96.8% in 2020 [26]. Most work on person RE-ID has focused on image-level matching, and recent state-of-the-art methods [27], [28] exploit this visual cue without other sources of information. However, image-level similarity matching is just a single component of a comprehensive RE-ID system. Persons must first be detected and tracked before matching, which can be computationally expensive to run simultaneously on a large network of cameras. Incorporating trajectory information for RE-ID in a multi-camera setting has seen comparatively little attention, in part due to a lack of publicly available datasets. One notable exception is the recent work of Wang et al. [3]. The authors demonstrate that by retrospectively utilizing trajectory information and visual features, their approach can attain state-of-the-art RE-ID results on the prevalent Duke-MCMT benchmark dataset [29], which is no longer available to the research community [30].

Multi-camera surveillance. A typical automated multi-camera surveillance system consists of detection, tracking, and RE-ID components in order to monitor objects-of-interest [31]–[33]. Scaling such systems to large camera networks can be challenging due to computational demands. Jain et al. [32] study the impact of scaling multi-camera tracking to large networks and show that filtering the search space to high traffic areas can considerably reduce the search space at only a small cost in recall. In addition to using the traffic level in each area, trajectory information has also been used for both RE-ID [3] and multi-camera tracking [14]. Using trajectories for these tasks can supplement existing appearance-based models by providing a second source of information for matching objects across camera views. However, previous works [3], [14], [31]–[33] reactively use trajectory information, i.e., the object must have already been observed in multiple camera views. In contrast, our MCTF framework predicts an object's future location *before* it is observed in another camera view.

3 THE MCTF FRAMEWORK

Consider a typical multi-camera surveillance setup where a set of k cameras $\mathcal{C} = \{c_i\}_{i=1}^k$ are mounted overhead monitoring objects of interest. Given an object's bounding boxes from the previous n timesteps up to the current

timestep t , $B_{(t-n:t)} = \{B_{(t-n)}, \dots, B_{(t-1)}, B_{(t)}\}$, we propose a hierarchy of MCTF problem formulations.

In which camera(s) will the object appear in the future?

Given $B_{(t-n:t)}$, our task is to identify a subset of \mathcal{C} in which the object may appear at any future timestep, up to a maximum of m timesteps. We cast this as a multi-class multi-label classification problem. Our goal is to estimate the probability of appearance of the object $P_a(c_i|B_{(t-n:t)})$ for each camera $c_i \in \mathcal{C}$ where a positive class is a camera in which the object re-appears. The output is a vector $[P_a(c_1|B_{(t-n:t)}), \dots, P_a(c_k|B_{(t-n:t)})]$ of length k .

When will the object appear? The task here is to predict when the object will reappear within the next m timesteps in a given camera. Similar to the *Which* problem, we also formulate this as a multi-class multi-label problem, where we compute the joint probability of appearance of the object $P_a(c_i, t_j|B_{(t-n:t)})$ for each camera $c_i \in \mathcal{C}$ at each timestep t_j with $j = 1, \dots, m$. The output is a matrix $\begin{bmatrix} P_a(c_1, t_1|B_{(t-n:t)}) & \dots & P_a(c_k, t_1|B_{(t-n:t)}) \\ \dots & \dots & \dots \\ P_a(c_1, t_m|B_{(t-n:t)}) & \dots & P_a(c_k, t_m|B_{(t-n:t)}) \end{bmatrix}$ of dimension $k \times m$.

Where will the object appear? This task aims at spatially localizing the object within a camera view in addition to *which* and *when*. To this end, we divide each camera view into a $w \times h$ grid and predict a probability of appearance score $P_a(c_i, t_j, g_{xy}|B_{(t-n:t)})$ for each grid cell g_{xy} , where $x = 1, \dots, w$ and $y = 1, \dots, h$. The output is a tensor \mathbf{Z} of dimension $k \times m \times w \times h$ containing the probability of appearance scores $P_a(c_i, t_j, g_{xy}|B_{(t-n:t)})$ for each camera $c_i \in \mathcal{C}$, timestep t_j with $j = 1, \dots, m$, and grid cell g_{xy} with $x = 1, \dots, w$ and $y = 1, \dots, h$.

4 PROPOSED APPROACH

In this section we introduce our data representation technique i.e., trajectory tensors, our models, and evaluation strategy for MCTF.

4.1 Trajectory Tensors

Existing works represent trajectories using a coordinate approach [6]–[8]. Coordinate trajectories are $(x, y)_t$ vectors in the image or world space, representing the location of an object at a particular timestep t . Sometimes, the width (w) and height (h) of the object may be also included in the representation, i.e., $(x, y, w, h)_t$ [10], [34].

There are three critical drawbacks to a coordinate trajectory representation: (i) Coordinates trajectories cannot represent a null trajectory. If no coordinates exist for an object, there is no defined representation for this absence. This can be problematic in real-world scenarios where object coordinates can become unavailable, such as when the object is occluded or the detection algorithm fails. It is particularly problematic in a multi-camera scenario where objects are not visible in all camera views simultaneously. (ii) Coordinate trajectories can only represent a trajectory as viewed from a single camera. Due to the null trajectory problem, coordinates cannot be easily generalized to multiple cameras, unless all objects are simultaneously visible in all cameras in \mathcal{C} , or separate models are created for each camera. Trajectories can instead be mapped to the world coordinate space to overcome this issue, however, this mapping requires

accurate measurements of the camera locations and intrinsic parameters which are not always available. (iii) Finally, coordinate trajectories do not inherently represent uncertainty, which is intrinsic to the trajectory forecasting task. The space of future trajectories is multi-modal, e.g., an object can travel either left or right at a junction. Existing works address this issue by using generative models to simulate multiple futures [7], [15], from which a probability distribution can be created. Our approach, in contrast, intuitively models uncertainty in one shot.

To overcome the shortcomings of coordinate trajectory representation, we introduce the idea of *trajectory tensors* – a novel technique for compact representation of multi-camera trajectories. As shown in Fig. 2, trajectory tensors are constructed in four steps:

- (i) We consider a set of cameras $\mathcal{C} = \{c_i\}_{i=1}^k$, where an object of interest may appear in any number of cameras in \mathcal{C} .
- (ii) Each camera c_i has an associated detection \mathbf{d}_i representing the object bounding box, if present.
- (iii) We convert each \mathbf{d}_i into a heatmap \mathbf{H}_i which is an alternative representation of the object's location. A heatmap is a matrix of size $w \times h$, where each entry is a binary value indicating the presence or absence of the object in this grid cell. A single object may span any number of grid cells, depending on its size.
- (iv) Finally, the heatmaps \mathbf{H}_i for each of the k cameras are stacked along the camera dimension, and computed for t timesteps. We also smooth each heatmap using a Gaussian kernel. The result is a trajectory tensor, \mathbf{Z} , of shape $k \times t \times w \times h$, where each entry is a value between 0 and 1. \mathbf{Z} represents the trajectory of an object in multiple camera views simultaneously and can be used to represent both past (inputs) and future (predicted) object locations.

Trajectory tensors share similarities with the grid-cell representation proposed recently by Liang et al. [15]; however, our proposed encoding includes multiple camera views. Besides, objects represented using trajectory tensors may span multiple grid cells, which allows us to account for variability in object scale. Object scale is not considered works involving birds-eye camera views [6]–[8], [15], but is a critical component in our framework.

4.2 MCTF models

Existing trajectory forecasting methods such as Social-LSTM [6], Social-GAN [7], and SoPhie [8] are designed for SCTF using datasets with birds-eye view cameras [19], [20] that do not include object scale. Although the methods proposed in [9] and [10] forecast object scale in addition to location, they are also designed for SCTF; hence direct comparison for MCTF is not possible. To compare with these existing works, we adapt the methods proposed in [9] and [10] to our MCTF framework. These models, along with new approaches based on trajectory tensors, are summarized in Fig. 3. Each neural-network-based model is comprised of a combination of fully-connected, recurrent, and convolutional layers. Inspired by fully-convolutional networks for semantic segmentation [35], we use transposed convolutional layers when predicting trajectory tensors as targets. Due to the large number of neural network models introduced in this section, we introduce each briefly and provide

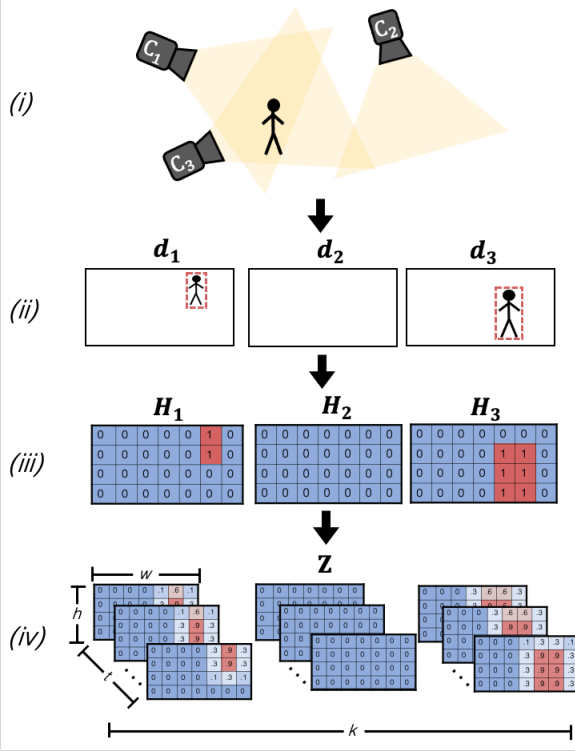


Fig. 2. **Trajectory tensors.** A traditional coordinate trajectory representation cannot represent object location information in multiple camera views, represent null trajectories, or represent uncertainty. Our proposed trajectory tensors overcome these shortcomings by representing object locations with heatmaps.

full details of the architectures needed to reproduce our results in the supplementary materials.

Coordinate trajectory approaches

Here we present our coordinate trajectory approaches which are shown in Fig. 3 (a) and (b). Due to the lack of data representations for coordinate trajectories in multiple camera views, our coordinate trajectory approaches each use a separate model for each camera, resulting in k models.

GRU. Recurrent networks have been a prevalent approach for trajectory forecasting in a single-camera. We use an adapted version of STED [10], a method for SCTF which uses gated recurrent units (GRUs) with an encoder-decoder architecture. The model proposed in the original work uses two encoders, one for bounding box coordinates, and another which uses a convolutional neural network (CNN) to extract motion information from optical flow. We use only the bounding box encoder for a fair comparison as other methods do not use visual features. For the *which* task, we use a fully-connected classification layer. For the *when* task, we use another GRU as a decoder with 128 hidden units followed by a fully connected classification layer. For the *where* task, the decoder GRU is followed by 2D transposed convolutional (tCNN) layers for spatial upsampling.

LSTM. Our LSTM model is the same as our GRU, with an alternative recurrent unit for both encoder and decoder.

1D-CNN. 1D-CNNs have received some attention as an alternative to the de facto recurrent models for tasks involving time series [36], including trajectory forecasting [9]. We use

the encoder architecture proposed by Yagi et al. [9], with modified decoders adapted for the MCTF formulation. For the *which* task, we use a fully-connected classification layer. For the *when* task, we use 1D transposed convolutional layers as a decoder with 4 layers. For the *where* task, the 1D transposed convolutional layers are followed by 2D transposed convolutional layers for spatial upsampling. Similarly to STED, we use only the trajectory feature extractor for a fair comparison with other methods.

Trajectory tensor approaches

Here we present our trajectory tensor approaches which are shown in Fig. 3 (c), (d), and (e). Our proposed representation enables efficient modeling of object trajectories across multiple camera views. Therefore, each approach consists of a single unified model for all cameras, which is less cumbersome and more scalable than the multi-model approach for coordinate trajectories.

3D-CNN. We use a 3D-CNN with four layers for spatio-temporal feature extraction. 3D convolutions can simultaneously extract spatial and temporal features and have seen some success for tasks such as action recognition [37].

2D-1D-CNN. We use a CNN consisting of three 2D convolutional and pooling layers for spatial feature extraction followed by three 1D convolutional and pooling layers for temporal feature extraction. Using separate 2D and 1D convolutional layers rather than 3D layers reduce the number of network parameters and is inspired by existing models for video classification [38].

CNN-GRU. We train a convolutional auto-encoder to reduce a trajectory tensor at a single timestep (size $k \times w \times h$) to a feature vector of size 512. This encoding is then used as the input to a GRU which predicts future feature vectors, which are decoded as shown in Fig. 3 (e). The auto-encoder is first pre-trained until convergence and then trained end-to-end with the GRU. This architecture is inspired by [39], which uses a similar strategy of forecasting future convolutional features for predicting future instance segmentation maps.

4.3 Evaluation strategy

Due to the high levels of uncertainty and the multi-modal nature of the MCTF, traditional SCTF metrics such as average and final displacement errors are not well-suited to MCTF. Some works generate multiple trajectories and select the one most similar to the ground truth [7]; however, this evaluation method is optimistic as performance comparable to sophisticated methods can be obtained using a simple constant velocity model that generates trajectories with high variance [40]. Owing to shortcomings of displacement error metrics, we instead compute the average precision (AP) for all problem formulations and plot precision-recall curves. We choose precision-recall curves over ROC curves. We find ROC to be an overly-optimistic metric due to the considerable class imbalance between the positive (object presence) and negative (no object presence) classes. We define AP_{which} , AP_{when} , and AP_{where} for the three problem definitions respectively.

The AP metrics provide a holistic interpretation of model performance, but are not easy to interpret. Therefore, we

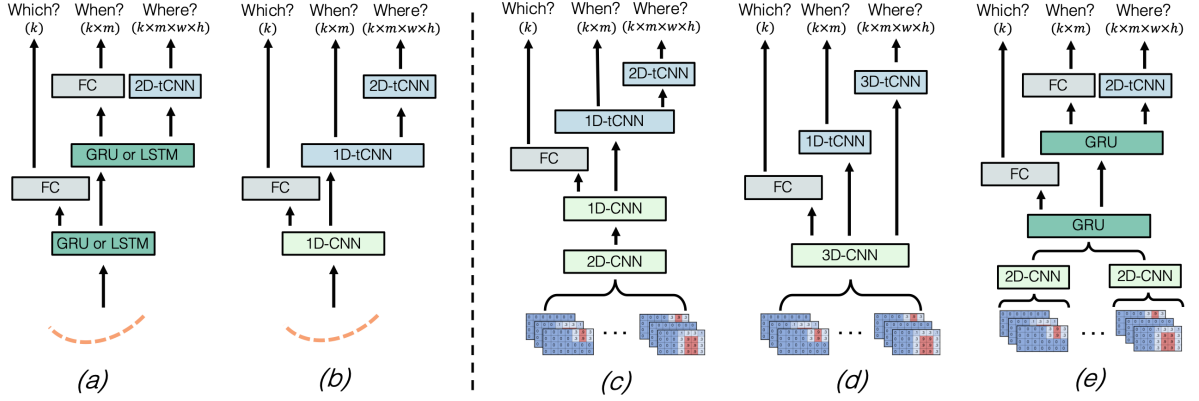


Fig. 3. **MCTF models.** We introduce 2 coordinate-trajectory based (left) and 3 trajectory-tensor based (right) approaches for MCTF. (a) A recurrent encoder-decoder adapted from [10]. (b) A 1D-CNN adapted from [9]. (c) A CNN approach with separated layers for spatial and temporal feature extraction. (d) A CNN approach with 3D convolutions for extracting spatial and temporal features simultaneously. (e) A hybrid CNN-GRU approach which uses a CNN to extract spatial features which are passed to an GRU for extracting temporal features. Note that for coordinate trajectory models, a separate model is created for each camera. In contrast, trajectory tensor models use a single unified model for all cameras. Each model is trained with either the *which*, *when*, or *where* MCTF formulation.

propose two new metrics for MCTF evaluation, the soft-intersection-over-union (*SIOU*) for the *when* and *where* problem formulations. The *SIOU* for evaluating the *when* formulation is as follows:

$$SIOU_{when} = \frac{1}{|\mathcal{C}^+|} \sum_c \frac{\sum_t^{\mathcal{T}^+} \hat{y}_t^c}{\sum_t^{\mathcal{T}} \hat{y}_t^c} \quad (1)$$

where \mathcal{C}^+ and \mathcal{T}^+ are sets of true positive cameras and timesteps, respectively. $\hat{y}_t^c \in (0, 1)$ is the predicted value for the presence the object in camera c at timestep t . Intuitively, $SIOU_{when}$ is a value between 0 and 1, representing the temporal overlap between the predicted and ground truth timesteps for all cameras where the individual appears.

We compute the *SIOU* for evaluating the *where* problem as follows:

$$SIOU_{where} = \frac{1}{|\mathcal{C}^+|} \sum_c \frac{1}{|\mathcal{T}^+|} \sum_t \frac{\sum_i^{\mathcal{I}^+} \hat{y}_t^c}{\sum_i^{\mathcal{I}} \hat{y}_t^c} \quad (2)$$

where \mathcal{I}^+ is the set of grid cells where the individual is present. $SIOU_{where}$ is similarly a value between 0 and 1 representing the spatial overlap between predicted and ground truth grid cell locations for all true positive cameras and timesteps.

5 PERFORMANCE EVALUATION

In this section, we introduce the evaluation dataset, baseline approaches, and assess the performance of each model for the *which*, *when*, and *where* tasks.

5.1 Warwick-NTU multi-camera forecasting dataset

We created the Warwick-NTU Multi-camera Forecasting (WNMF) dataset to train and evaluate each model. WNMF was collected specifically for MCTF using a set of 15 cameras in a building on the Nanyang Technological University campus and contains both overlapping and non-overlapping views recorded over 20 days. The dataset consists of cross-camera trajectories where an individual departs from one

camera view and then re-appears in another after no more than 12 seconds. An individual may also be visible in any number of other camera views during this tracking period. We add two extensions to the dataset compare to our previous release: (i) Multi-viewpoint departures. The original dataset consists of matches across pairs of cameras. We add new annotations for overlapping viewpoints where an individual is visible in multiple cameras simultaneously. (ii) Cleaned annotations. The original dataset contains 2.3K cross-camera matches; however, we found several erroneous matches that have been removed, leaving 2.0K matches. The erroneous matches were caused by high visual similarity between different persons during labelling which were removed by hand, referring to the camera network topology along with the appearance timestamps. Details on the semi-automated data annotation procedure can be found in our previous work [18]. We use a robust 5-fold cross-validation setup in all experiments that follow, where training and testing sets contain footage recorded on different days. The WNMF dataset is available to download at <https://github.com/olly-styles/Multi-Camera-Trajectory-Forecasting>.

5.2 Baseline approaches

In addition to the models introduced in Section 4.2 we also evaluate several baselines.

Shortest real-world distance. We use the physical distance between cameras in the real world and predict the camera closest to the current camera. This baseline applies to the *which* formulation only.

Training set mean. We consider all training set observations for a particular camera and take the mean of all ground truth labels in the training set.

Most similar trajectory. We find the most similar trajectory in terms of L_2 distance in the same camera from the training set to the observed trajectory and predict the same label.

Hand-crafted features. We extract some features from the bounding boxes and classify them with a fully-connected network. Our 10-dimensional hand-crafted feature vector contains velocity in x and y direction, acceleration in x and y

TABLE 1

Which results. Given observations from one camera, each model predicts which camera(s) the person will re-appear.

	Model	AP_{which}
Baselines	Shortest real-world distance	44.0
	Training set mean	67.4
	Most similar trajectory	64.4
	Hand-crafted features	77.0
Coordinate trajectories	LSTM	83.6
	GRU	84.7
	1D-CNN	83.3
Trajectory tensors	2D-1D-CNN	86.1
	3D-CNN	86.6
	CNN-GRU	86.0

direction, last observed bounding box height and width, and its four coordinates. We compute all features with respect to the 2D coordinate system as captured by the camera.

5.3 Which camera will they appear?

Experimental setup. We use our proposed network architectures (Fig. 3) and baselines (Section 5.2). For the coordinate trajectory approaches, we train a separate model for each camera. To adapt the SCTF models for MCTF, we leave encoders unchanged and change decoders to fully-connected output layers of size 15, the number of cameras in the WNMF dataset. We use a binary cross-entropy loss function with a sigmoid activation function at the output layer. The activation maps each output to a value between 0 and 1, representing the appearance of the individual in each camera. Coordinate trajectory and trajectory tensor models are trained with the Adam optimizer using learning rates of 1×10^{-3} and 1×10^{-4} , respectively. Coordinate trajectory encoders extract a feature vector of size 128, whereas trajectory tensors approaches extract a feature vector of size 512 as a single encoder is shared across all cameras and therefore requires larger representation capacity. All models are trained using a batch size of 64. We use a heatmap H of size either 16×9 , 32×18 , or 48×27 and a Gaussian smoothing kernel size between 0 and 4 as input, chosen using cross-validation. The impact of changing these parameters is investigated in Section 5.6.

Results. The AP_{which} for each model is shown in Table 1 and precision-recall plots are shown in Fig. 4 (left). Coordinate trajectory approaches outperform our 4 baselines, and trajectory tensor approaches perform the best.

5.4 When will they appear?

Experimental setup. We use the same encoders for evaluating *when*, with either 1D-transposed convolutional layers or recurrent unit for decoding as shown in Fig. 3. Each model is trained with the binary cross-entropy loss, and the trajectory tensor models are trained with the same hyperparameters as the *which* problem formulation.

Results. The AP_{when} and $SIOU_{when}$ for each model are shown in Table 2 and precision-recall plots are shown in Fig. 4 (center). We observe a similar trend to the results of the *which* task, with the 3D-CNN trajectory tensor model performing best. Some baseline methods perform well in terms of $SIOU_{when}$, but poorly in terms of AP_{when} . This

TABLE 2

When results. Given observations from one camera, each model predicts in which camera(s) the person will re-appear, and in which timesteps they will be present.

	Model	AP_{when}	$SIOU_{when}$
Baselines	Training set mean	60.6	78.9
	Most similar trajectory	45.7	62.3
	Hand-crafted features	69.1	80.4
Coordinate trajectories	LSTM	77.4	79.1
	GRU	77.4	78.6
	1D-CNN	77.4	80.3
Trajectory tensors	2D-1D-CNN	75.9	80.5
	3D-CNN	78.3	81.2
	CNN-GRU	76.7	80.8

TABLE 3

Where results. Given observations from one camera, each model predicts which camera(s) the person will re-appear, in which timesteps they will be present, and where in the camera view they will appear.

	Model	AP_{where}	$SIOU_{where}$
Baselines	Training set mean	27.9	45.3
	Most similar trajectory	11.4	65.1
	Hand-crafted features	24.8	60.9
Coordinate trajectories	LSTM	15.8	42.4
	GRU	15.9	42.6
	1D-CNN	29.8	54.1
Trajectory tensors	2D-1D-CNN	37.2	60.1
	3D-CNN	37.6	66.1
	CNN-GRU	29.8	45.8

result suggests that these approaches effectively predict the correct time window an individual will be visible in the target camera but often predict the wrong camera.

5.5 Where will they appear?

Experimental setup. We use the same target heatmap size of 16×9 regardless of the input heatmap size for a fair comparison. We use this small heatmap size to reduce computational complexity, and 16×9 is the smallest possible while maintaining the original aspect ratio of the video frames. The output trajectory tensor is, therefore, of size $15 \times 60 \times 16 \times 9$. We do not apply Gaussian smoothing to the ground truth trajectory tensor targets. We use 2D or 3D transposed convolutional layers to upsample extracted feature vectors to trajectory tensor outputs, representing an individual's future location. The trajectory tensor models are trained with the same hyperparameters as in Section 5.3, again using the binary cross-entropy loss.

Results. The AP_{where} and $SIOU_{where}$ for each model is shown in Table 3 and precision-recall plots are shown in Fig. 4 (right). Note that due to the considerably increased complexity of the *Where* problem formulation compared to *Which* and *When*, the MAP is much lower. Trajectory tensor-based approaches perform the best.

5.6 Ablation studies

Multi-view trajectory tensors. Our proposed trajectory tensor data representation enables us to model the same individual captured in multiple camera views. To study the impact of multiple camera views, we train our trajectory tensor-based models using only one of the available views,

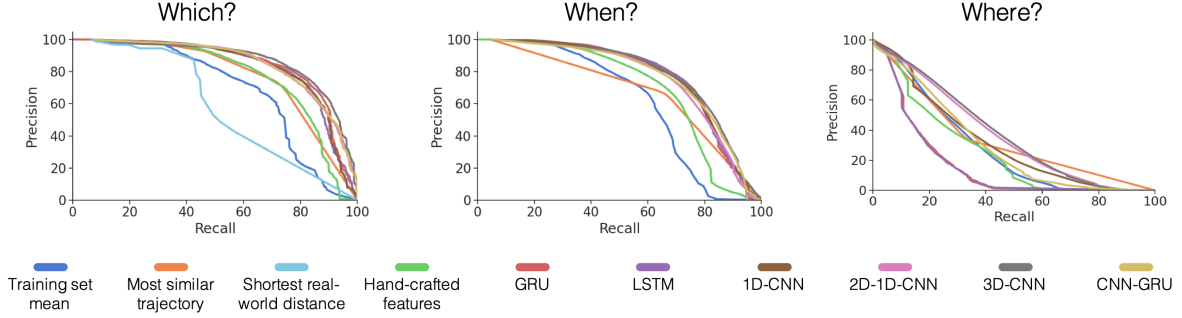


Fig. 4. **Precision-recall plots.** Precision and recall of each model for all three problem formulations. Best viewed in colour.

TABLE 4

Multi-view trajectory tensor results. Comparison of single-view and multi-view trajectories. Observing a trajectory in multiple camera views improves model performance in almost all cases.

Model	Multi-view trajectories	AP_{which}	AP_{when}	AP_{where}
2D-1D-CNN	No	81.1	74.9	37.0
	Yes	86.1	75.9	37.2
3D-CNN	No	83.1	76.1	37.2
	Yes	87.3	78.3	37.6
CNN-GRU	No	81.2	75.4	30.0
	Yes	86.0	76.7	29.8

i.e., we set each heatmap at each timestep to the zero matrix for all but one of the camera channels $c \in \mathbf{Z}$. The results in Table 4 show a comparison of our trajectory tensor models with a single-view trajectory (each heatmap is 0 for all but one of the cameras $c \in \mathbf{Z}$) and a multi-view trajectory (where more than one of the camera channels $c \in \mathbf{Z}$ may be non-zero). The results show that the models are able to make use of the location information available in multiple camera views.

Heatmap size and smoothing. We investigate the impact of heatmap size and standard deviation of the Gaussian filter for smoothing. Larger heatmap sizes afford models more representation power at the cost of a larger number of parameters and a tendency to overfit the training data. On the other hand, smoothing has a regularizing effect by reducing the impact of errors during the detection and tracking stages. We suggest, therefore, that both heatmap size and smoothing sigma should be tuned in tandem. Fig. 5 shows the impact of heatmap size and smoothing sigma.

5.7 Discussion and applications

Our results show that using trajectory tensors has several advantages over traditional coordinate trajectories for MCTF. Representing trajectories as tensors allows us to model relative object locations in multiple camera views simultaneously, which improves MCTF performance when used as inputs and provides an intuitive way to represent multi-modal futures when used as targets. Trajectory tensors are also considerably more efficient than using coordinate trajectories in an MCTF setting, as a single model can be used rather than creating separate models for each camera.

We find the 3D-CNN architecture consistently performs the best across MCTF tasks. 3D-CNNs have traditionally only seen major success in settings where vast amounts of data are available for training [37]. However, in our

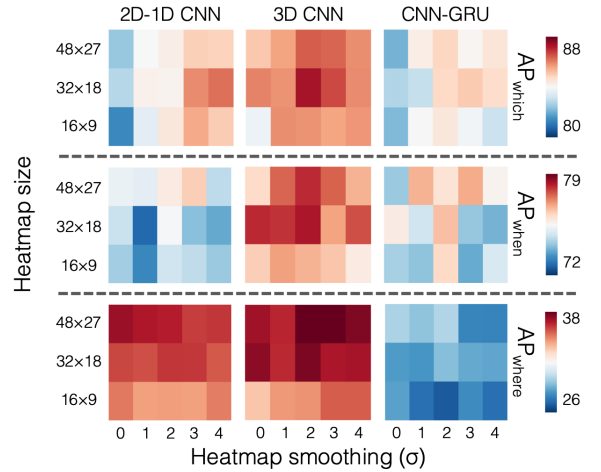


Fig. 5. **Impact of heatmap size and smoothing.** Here we show the impact of heatmap size and smoothing sigma using 5-fold cross validation. We observe that larger heatmap sizes are most beneficial for the *Where* problem formulation where more fine-grained predictions are required.

setting, we use resolutions of up to 48×27 rather than the 224×224 or higher commonly used in activity recognition, which considerably reduces the number of parameters and facilitates training on smaller datasets.

We anticipate three applications of the MCTF framework: (i) Long-term forecasting. This is possible by removing the constraint of a single camera viewpoint. (ii) Intelligent camera monitoring. When tracking a particular object-of-interest across a camera network, a subset of cameras may be monitored intelligently using the predictions from an MCTF model rather than continually monitoring all cameras. (iii) Enhanced tracking. Location predictions may be used in conjunction with a RE-ID model for more robust multi-camera tracking.

6 CONCLUSION

We have developed a complete framework for MCTF formulated in a hierarchy of three spatio-temporal localization tasks: In *which* camera, *when*, and *where* will the object(s) appear? Our work is the first to address the challenges of trajectory forecasting in a multi-camera environment. We introduced the idea of the *trajectory tensor* - a new trajectory representation that facilitates the encoding of multi-modal futures and associated uncertainty. Trajectory tensors are an attractive alternative to the traditional coordinate trajectory

representation used in previous works. Our proposed MCTF models based on trajectory tensors show promising results on our WNMf database, and we hope to evaluate in more settings when more datasets for MCTF become available. We envision our MCTF model will enhance multi-camera surveillance systems, complementing existing models for person RE-ID and cross-camera tracking.

ACKNOWLEDGMENTS

This work is funded by the UK EPSRC (grant no. E/L016400/1). Our thanks to NVIDIA for their generous hardware donation. We would also like to thank the ROSE lab at the Nanyang Technological University, Singapore for providing the data used in this research.

REFERENCES

- [1] Xiaogang Wang. Intelligent multi-camera video surveillance: A review. *Pattern recognition letters*, 34(1), 2013.
- [2] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B Choy, Philip HS Torr, and Manmohan Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents. In *CVPR*, 2017.
- [3] Guangcong Wang, Jianhuang Lai, Peigen Huang, and Xiaohua Xie. Spatial-temporal person re-identification. In *AAAI*, 2019.
- [4] Zheng Tang, Milind Naphade, Ming-Yu Liu, Xiaodong Yang, Stan Birchfield, Shuo Wang, Ratnesh Kumar, David Anastasiu, and Jenq-Neng Hwang. Cityflow: A city-scale benchmark for multi-target multi-camera vehicle tracking and re-identification. In *CVPR*, 2019.
- [5] Shuai Yi, Hongsheng Li, and Xiaogang Wang. Pedestrian behavior understanding and prediction with deep neural networks. In *ECCV*. Springer, 2016.
- [6] Alexandre Alahi, Kratharth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *CVPR*, 2016.
- [7] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *CVPR*, 2018.
- [8] Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, Hamid Reza Tofighi, and Silvio Savarese. Sophie: An attentive gan for predicting paths compliant to social and physical constraints. In *CVPR*, 2019.
- [9] Takuma Yagi, Karttikeya Mangalam, Ryo Yonetani, and Yoichi Sato. Future person localization in first-person videos. In *CVPR*, 2018.
- [10] Olly Styles, Tanaya Guha, and Victor Sanchez. Multiple object forecasting: Predicting future object locations in diverse environments. In *WACV*. IEEE, 2020.
- [11] Francesco Marchetti, Federico Becattini, Lorenzo Seidenari, and Alberto Del Bimbo. Mantra: Memory augmented networks for multiple trajectory prediction. In *CVPR*, 2020.
- [12] Hiroaki Minoura, Ryo Yonetani, Mai Nishimura, and Yoshitaka Ushiku. Crowd density forecasting by modeling patch-based dynamics. *arXiv preprint arXiv:1911.09814*, 2019.
- [13] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *CVPR*, 2017.
- [14] Hung-Min Hsu, Tsung-Wei Huang, Gaoang Wang, Jiarui Cai, Zhichao Lei, and Jenq-Neng Hwang. Multi-camera tracking of vehicles based on deep features re-id and trajectory-based camera link models. In *CVPRW*, 2019.
- [15] Junwei Liang, Lu Jiang, Kevin Murphy, Ting Yu, and Alexander Hauptmann. The garden of forking paths: Towards multi-future trajectory prediction. *arXiv preprint arXiv:1912.06445*, 2019.
- [16] Shanshan Zhang, Rodrigo Benenson, Mohamed Omran, Jan Hosang, and Bernt Schiele. Towards reaching human performance in pedestrian detection. *T-PAMI*, 40(4), 2017.
- [17] Liang Zheng, Yi Yang, and Alexander G Hauptmann. Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984*, 2016.
- [18] Olly Styles, Tanaya Guha, Victor Sanchez, and Alex Kot. Multi-camera trajectory forecasting: Pedestrian trajectory prediction in a network of cameras. In *CVPRW*, 2020.
- [19] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. You'll never walk alone: Modeling social behavior for multi-target tracking. In *ICCV*. IEEE, 2009.
- [20] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. Crowds by example. In *Computer graphics forum*, volume 26. Wiley Online Library, 2007.
- [21] Igor Gilitschenski, Guy Rosman, Arjun Gupta, Sertac Karaman, and Daniela Rus. Deep context map: Agent trajectory prediction using location-specific latent maps. *arXiv preprint arXiv:1912.06785*, 2019.
- [22] Junwei Liang, Lu Jiang, Juan Carlos Niebles, Alexander G Hauptmann, and Li Fei-Fei. Peeking into the future: Predicting future person activities and locations in videos. In *CVPR*, 2019.
- [23] Amir Sadeghian, Ferdinand Legros, Maxime Voisin, Ricky Vesel, Alexandre Alahi, and Silvio Savarese. Car-net: Clairvoyant attentive recurrent network. In *ECCV*, 2018.
- [24] Karttikeya Mangalam, Harshayu Girase, Shreyas Agarwal, Kuan-Hui Lee, Ehsan Adeli, Jitendra Malik, and Adrien Gaidon. It is not the journey but the destination: Endpoint conditioned trajectory prediction. *arXiv preprint arXiv:2004.02025*, 2020.
- [25] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *ICCV*, 2015.
- [26] Zhihui Zhu, Xinyang Jiang, Feng Zheng, Xiaowei Guo, Feiyue Huang, Xing Sun, and Weishi Zheng. Viewpoint-aware loss with angular regularization for person re-identification. In *AAAI*, 2020.
- [27] Ruijie Quan, Xuanyi Dong, Yu Wu, Linchao Zhu, and Yi Yang. Auto-reid: Searching for a part-aware convnet for person re-identification. *CVPR*, 2019.
- [28] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang. Bag of tricks and a strong baseline for deep person re-identification. In *CVPRW*, 2019.
- [29] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *ECCV*. Springer, 2016.
- [30] Jake Satsky. *A Duke study recorded thousands of students' faces. Now they're being used all over the world*, 2019 (accessed 20th February, 2020). <https://www.dukechronicle.com/article/2019/06/duke-university-facial-recognition-data-set-study-surveillance-video-students-china-uyghur>.
- [31] Niki Martinel, Gian Luca Foresti, and Christian Micheloni. Person reidentification in a distributed camera network framework. *IEEE transactions on cybernetics*, 47(11), 2016.
- [32] Samvit Jain, Ganesh Ananthanarayanan, Junchen Jiang, Yuanchao Shu, and Joseph Gonzalez. Scaling video analytics systems to large camera deployments. In *International Workshop on Mobile Computing Systems and Applications*, 2019.
- [33] Octavia Camps, Mengran Gou, Tom Hebble, Srikrishna Karanam, Oliver Lehmann, Yang Li, Richard J Radke, Ziyang Wu, and Fei Xiong. From the lab to the real world: Re-identification in an airport camera network. *IEEE transactions on circuits and systems for video technology*, 27(3), 2016.
- [34] Yu Yao, Mingze Xu, Chiho Choi, David J Crandall, Ella M Atkins, and Behzad Dariush. Egocentric vision-based future vehicle localization for intelligent driving assistance systems. In *ICRA*. IEEE, 2019.
- [35] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- [36] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhasane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. *Data Mining and Knowledge Discovery*, 33(4), 2019.
- [37] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *CVPR*, 2018.
- [38] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *ECCV*, 2018.
- [39] Pauline Luc, Camille Couprie, Yann Lecun, and Jakob Verbeek. Predicting future instance segmentation by forecasting convolutional features. In *ECCV*, 2018.
- [40] Christoph Schöller, Vincent Aravantinos, Florian Lay, and Alois Knoll. What the constant velocity model can teach us about pedestrian motion prediction. *IEEE Robotics and Automation Letters*, 5(2), 2020.

Supplementary: Multi-Camera Trajectory Forecasting With Trajectory Tensors

Olly Styles, Tanaya Guha, and Victor Sanchez

1 MODEL ARCHITECTURE DETAILS

In this section, we provide full details of each model architecture in our experiments. The code is available at: <https://github.com/olly-styles/Trajectory-Tensors>

TABLE 1
RNN architecture. Based on [1]

Block	Layer type	Kernel size	Stride	Output dimension
Encoder	Input	-	-	4x10
	FC + ReLU	-	-	128x10
	RNN + ReLU	-	-	128
Which decoder	FC + Sig.	-	-	15
When decoder	RNN + ReLU	-	-	128x60
	FC + Sig.	-	-	15x60
Where decoder	RNN + ReLU	-	-	128x60
	2D-TConv + ReLU	1x3	1	128x60x1x3
	2D-TConv + ReLU	4x3	2	128x60x4x7
	2D-TConv + Sig.	3x4	2	128x60x9x16

TABLE 2
1D-CNN architecture. Based on [2]

Block	Layer type	Kernel size	Stride	Output dimension
Encoder	Input	-	-	4x10
	1D-Conv + ReLU	3	1	32x8
	1D-Conv + ReLU	3	1	64x6
	1D-Conv + ReLU	3	1	128x4
Which decoder	1D-Conv + ReLU	3	1	128x2
	FC + Sig.	-	-	15
When decoder	1D-TConv + ReLU	4	2	128x6
	1D-TConv + ReLU	4	2	128x14
	1D-TConv + ReLU	4	2	128x30
	1D-TConv + Sig.	4	2	128x60
Where decoder	1D-TConv + ReLU	4	2	128x6
	1D-TConv + ReLU	4	2	128x14
	1D-TConv + ReLU	4	2	128x30
	1D-TConv + ReLU	4	2	128x60
	2D-TConv + ReLU	1x3	1	128x60x1x3
	2D-TConv + ReLU	4x3	2	128x60x4x7
	2D-TConv + Sig.	3x4	2	128x60x9x16

TABLE 3
2D-1D CNN architecture.

Block	Layer type	Kernel size	Stride	Output dimension
Encoder (9x16)	Input	-	-	15x10x9x16
	2D-Conv + ReLU	3x3	1	64x10x7x14
	2D Max pool	2x2	-	64x10x3x7
	2D-Conv + ReLU	2x3	1	128x10x2x5
	2D Max pool	2x5	1	128x10x1x1
	1D-Conv + ReLU	3	1	256x6
Encoder (18x32)	1D-Max pool + ReLU	2	1	256x3
	1D-Conv + ReLU	3	1	256x1
	Input	-	-	15x10x18x32
	2D-Conv + ReLU	5x5	1	64x10x14x28
Encoder (27x48)	2D Max pool	2x3	-	64x10x7x9
	2D-Conv + ReLU	3x5	1	256x10x5x5
	2D Max pool	5x5	1	256x10x1x1
	1D-Conv + ReLU	3	1	256x6
	1D-Max pool + ReLU	2	1	256x3
	1D-Conv + ReLU	3	1	256x1
Encoder (27x48)	Input	-	-	15x10x27x48
	2D-Conv + ReLU	5x5	1	32x10x23x44
	2D Max pool	2x2	-	32x10x11x22
	2D-Conv + ReLU	3x3	1	128x10x9x20
	2D Max pool	2x3	1	128x10x4x6
	2D-Conv + ReLU	3x3	1	256x10x2x3
Which decoder	2D Max pool	2x3	1	256x10x1x1
	1D-Conv + ReLU	3	1	256x6
	1D-Max pool + ReLU	2	1	256x3
	1D-Conv + ReLU	3	1	256x1
When decoder	FC + Sig.	-	-	15
	1D-TConv + ReLU	4	2	256x5
	1D-TConv + ReLU	4	2	128x13
	1D-TConv + ReLU	4	2	64x29
Where decoder	1D-TConv + ReLU	4	2	32x60
	1D-TConv + Sig.	1	1	15x60
	1D-TConv + ReLU	7	1	256x7
	1D-TConv + ReLU	7	2	256x19
Where decoder	1D-TConv + ReLU	6	3	256x60
	2D-TConv + ReLU	1x3	2x2	128x60x1x3
	2D-TConv + Sig.	4x3	2x2	64x60x4x7
	2D-TConv + Sig.	3x4	2x2	15x60x9x16

• The authors are with the Department of Computer Science, University of Warwick, United Kingdom. Contact email: o.c.styles@warwick.ac.uk

Manuscript received April 19, 2005; revised August 26, 2015.

TABLE 4
3D CNN architecture.

Block	Layer type	Kernel size	Stride	Output dimension
Encoder (9x16)	Input	-	-	15x10x9x16
	3D-Conv + ReLU	3x3x3	1	64x8x7x14
	3D Max pool	2x2x3	-	64x4x3x4
	3D-Conv + ReLU	3x3x3	1	256x2x1x2
	3D Max pool	1x1x2	-	256x1x1x1
Encoder (18x32)	Input	-	-	15x10x18x32
	3D-Conv + ReLU	3x3x3	1	64x8x16x30
	3D Max pool	1x2x2	-	64x8x8x15
	3D-Conv + ReLU	3x3x3	1	128x6x6x13
	3D Max pool	2x2x2	-	128x3x3x6
	3D-Conv + ReLU	3x3x3	1	256x1x1x4
Encoder (27x48)	3D Max pool	1x1x4	-	256x1x1x1
	Input	-	-	15x10x27x48
	3D-Conv + ReLU	1x3x3	1	64x10x25x46
	3D Max pool	1x1x2	-	64x10x25x23
	3D-Conv + ReLU	3x3x3	1	128x8x23x21
	3D Max pool	1x2x2	-	128x8x11x10
	3D-Conv + ReLU	3x3x3	1	256x6x9x8
	3D Max pool	2x2x2	-	256x3x4x4
Which decoder	3D-Conv + ReLU	3x3x3	1	256x1x2x2
	3D Max pool	1x2x2	-	256x1x1x1
When decoder	FC + Sig.	-	-	15
Where decoder	1D-TConv + ReLU	4	2	256x5
	1D-TConv + ReLU	4	2	128x13
	1D-TConv + ReLU	4	2	64x29
	1D-TConv + ReLU	4	2	32x60
	1D-TConv + Sig.	1	1	15x60
Where decoder	3D-TConv + ReLU	5x1x1	2	256x5x1x1
	3D-TConv + ReLU	5x1x3	2	256x13x1x3
	3D-TConv + ReLU	5x4x3	2	128x29x4x7
	3D-TConv + Sig.	4x3x4	2	15x60x9x16

TABLE 5
CNN-GRU architecture.

Block	Layer type	Kernel size	Stride	Output dimension
Encoder (9x16)	Input	-	-	15x10x9x16
	2D-Conv + ReLU	3x3	1	128x10x7x14
	2D Max pool	2x2	-	128x10x3x7
	2D-Conv + ReLU	2x3	1	256x10x2x5
	2D Max pool	2x5	1	256x10x1x1
	GRU + ReLU	-	-	256x1
Encoder (18x32)	Input	-	-	15x10x18x32
	2D-Conv + ReLU	5x5	1	64x10x14x28
	2D Max pool	2x3	-	64x10x7x9
	2D-Conv + ReLU	3x5	1	256x10x5x5
	2D Max pool	5x5	1	256x10x1x1
	GRU + ReLU	-	-	256x1
Encoder (27x48)	Input	-	-	15x10x27x48
	2D-Conv + ReLU	5x5	1	32x10x23x44
	2D Max pool	2x2	-	32x10x11x22
	2D-Conv + ReLU	3x3	1	128x10x9x20
	2D Max pool	2x3	1	128x10x4x6
	2D-Conv + ReLU	3x3	1	256x10x2x3
	2D Max pool	2x3	1	256x10x1x1
Which decoder	GRU + ReLU	-	-	256x1
	FC + Sig.	-	-	15
When decoder	GRU + ReLU	-	-	256x60
	FC + Sig.	-	-	15x60
Where decoder	GRU + ReLU	-	-	256x60
	2D-TConv + ReLU	1x3	2x2	128x60x1x3
	2D-TConv + Sig.	4x3	2x2	64x60x4x7
	2D-TConv + Sig.	3x4	2x2	15x60x9x16

2 EXAMPLE PREDICTIONS

In this section, we show example successful and unsuccessful predictions using our 3 problem formulations. To aid with visualization, please see the camera network topology shown in Figure 1. More visual examples can be found in our companion video, available online: <https://youtu.be/H0gYJ-z8mVc>

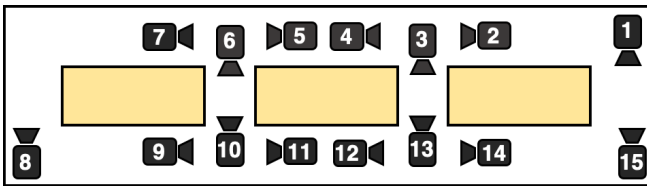


Fig. 1. **Camera network topology.** The WNMF dataset contains 15 cameras mounted overhead in the layout as shown.

2.1 Which

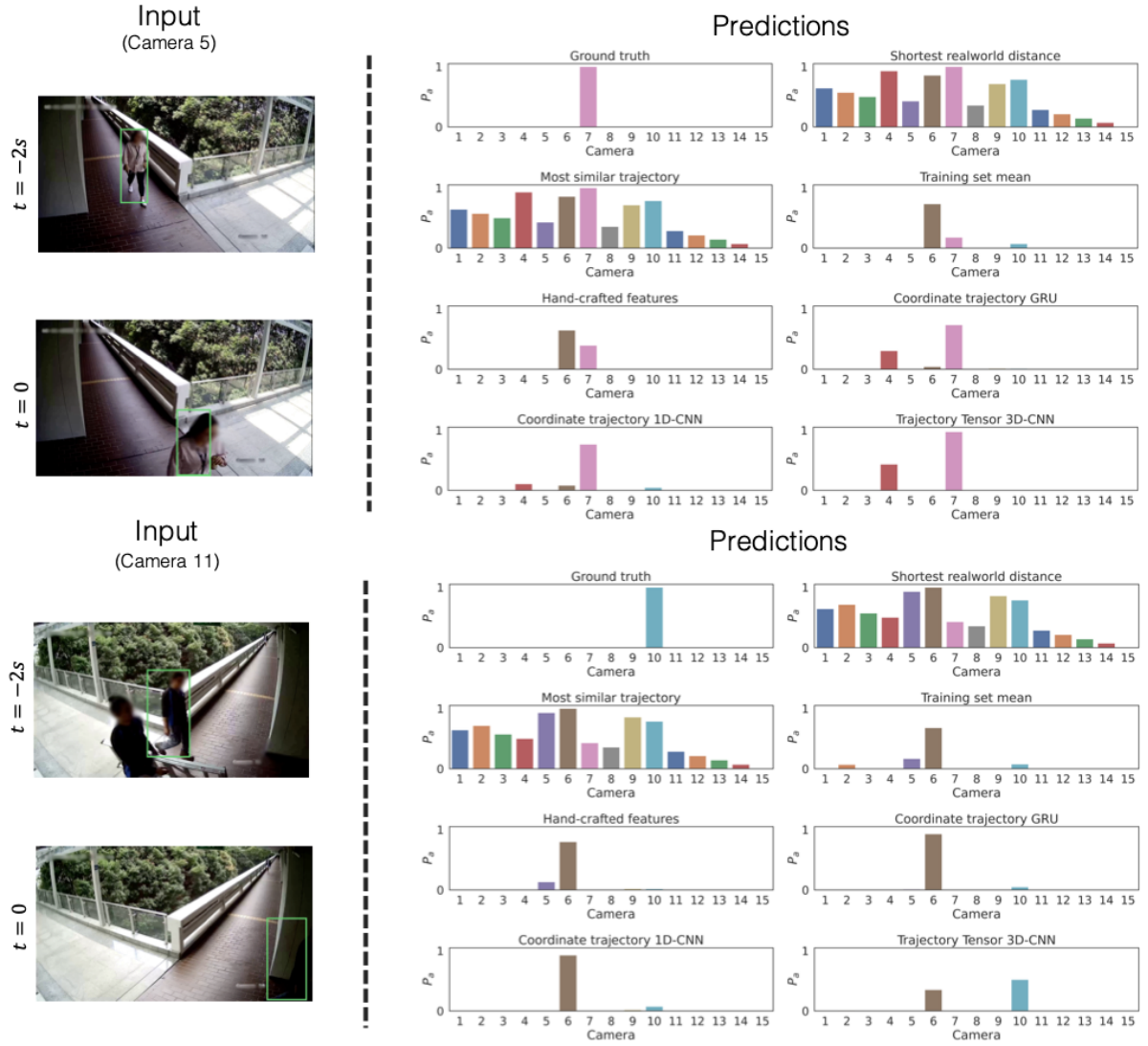


Fig. 2. **Example successful which predictions.** Our model successfully anticipates the next camera an individual will re-appear.

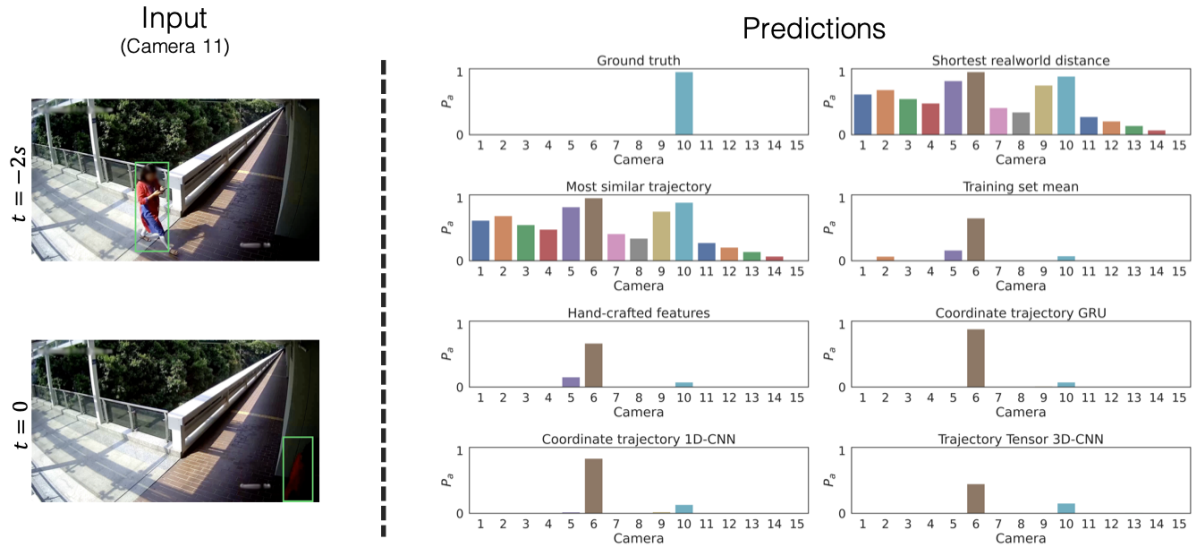


Fig. 3. **Example unsuccessful which predictions.** All models incorrectly assign camera 6 as the most likely next camera rather than camera 10. These two camera views are overlapping, so either may detect the individual.

2.2 When

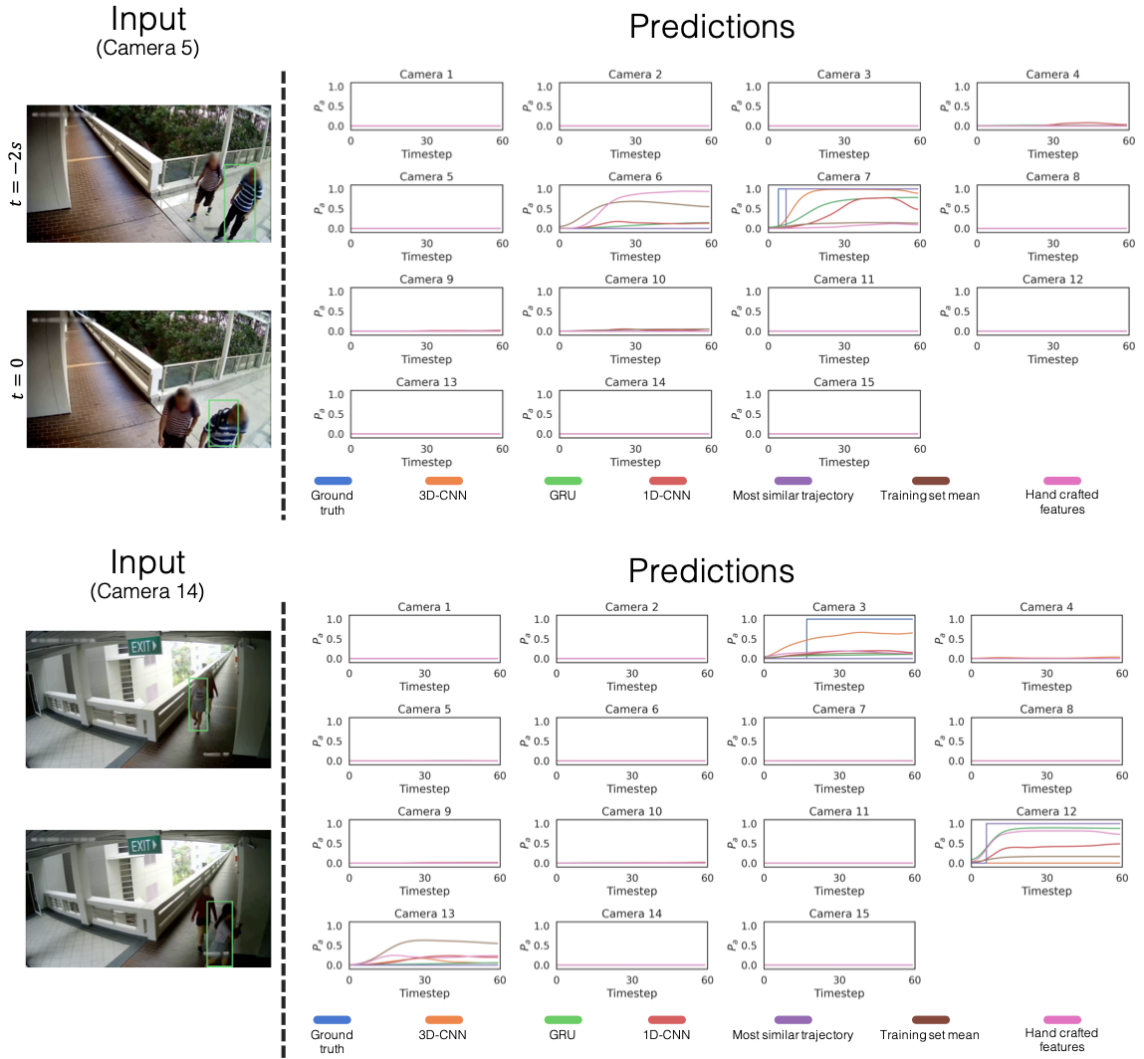


Fig. 4. **Example successful when predictions.** Our models correctly anticipates individuals turning at an intersections and approximates the time of re-appearance in another camera in the network.

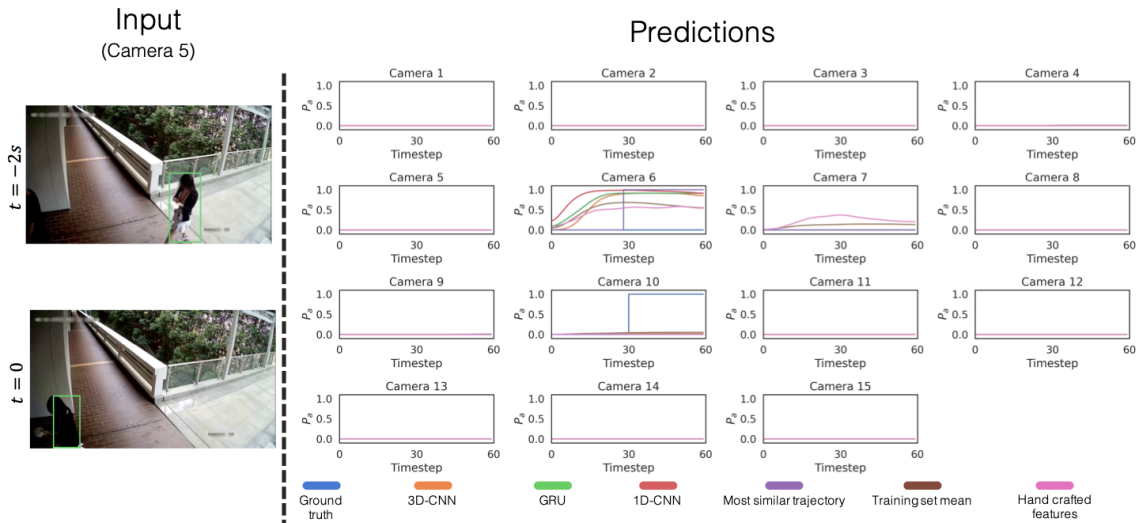


Fig. 5. **Example unsuccessful when predictions.** This is a similar scenario as the *which* prediction in Figure 3 where the individual is predicted to appear in the wrong camera out of two overlapping views.

2.3 Where

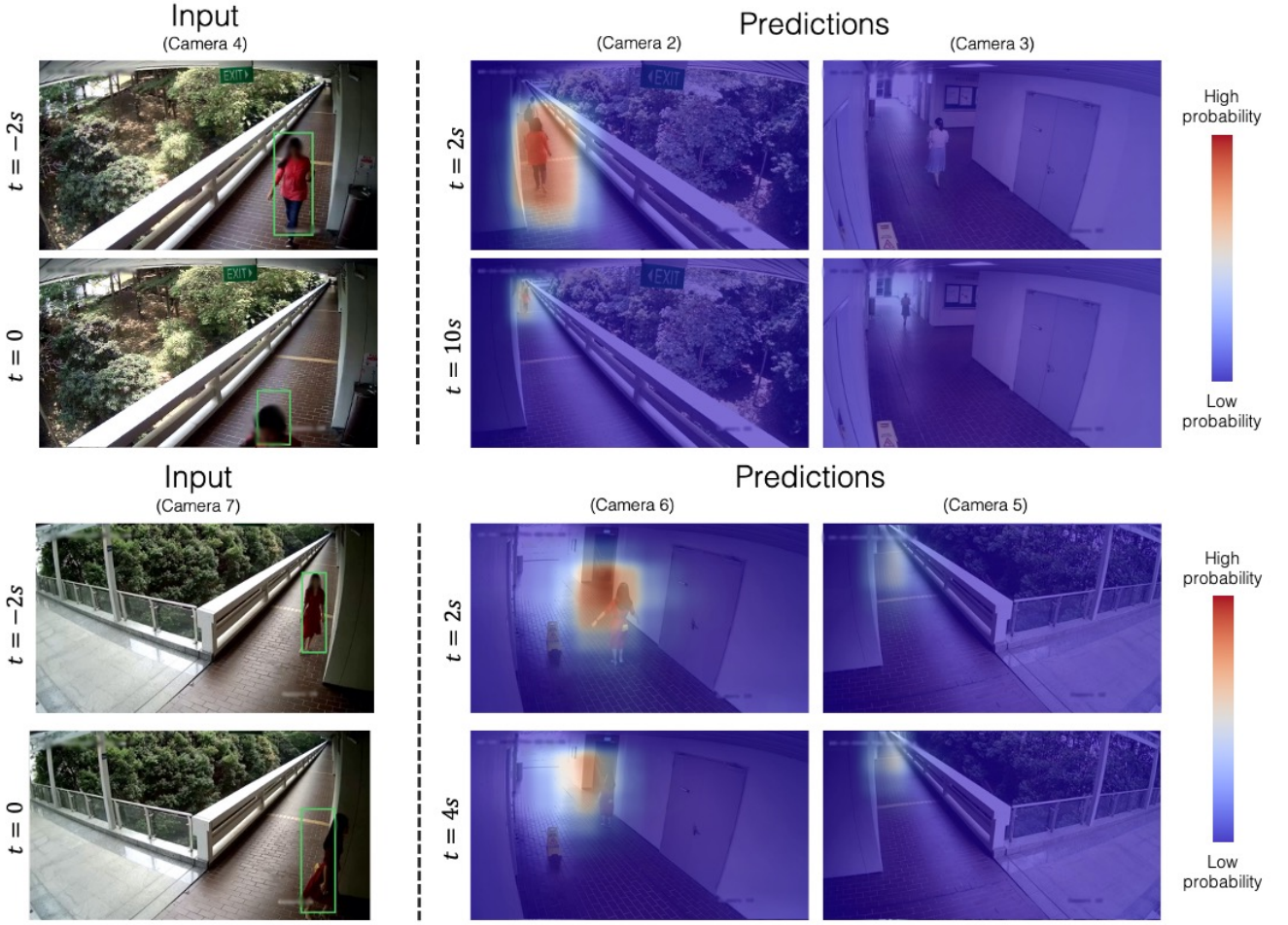


Fig. 6. **Example successful predictions.** Here we show example *where* predictions for our trajectory tensor model with 3D-CNN architecture. The upper example shows a long-term uni-modal example where the model assigns high probability to the individual continuing straight and low probability to other camera views (only camera 3 shown). The lower example shows a shorter-term multi-modal example, where the model predicts some chance that the individual will continue going straight (visible in camera 5) in addition to correctly predicting a left turn (visible in camera 6). The predictions are made in a 16×9 grid which are smoothed with a Gaussian kernel for cleaner visualization.

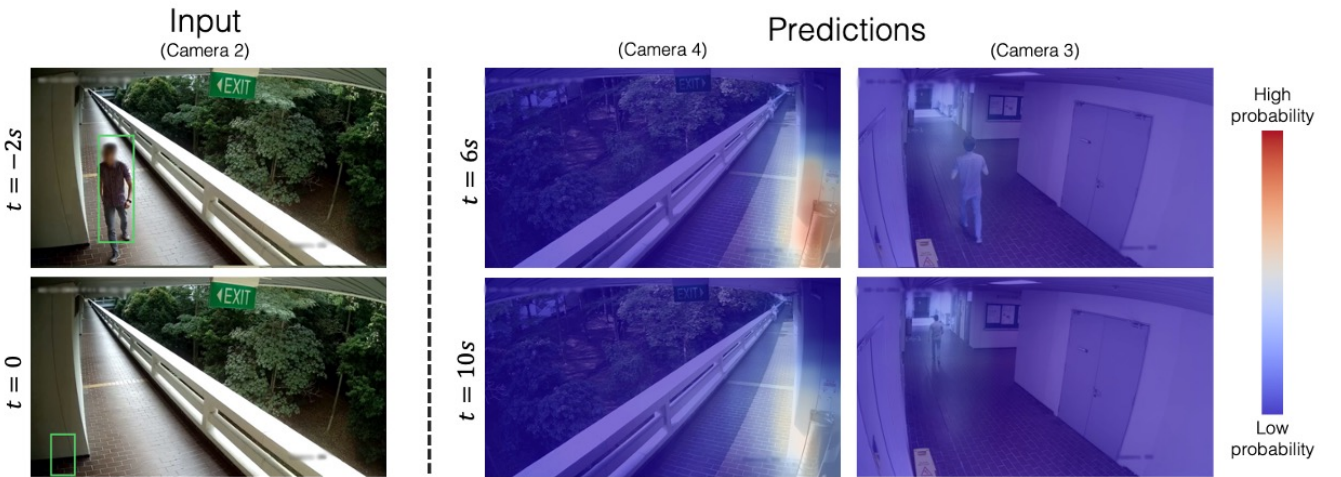


Fig. 7. **Example unsuccessful prediction.** Here the model predicts that the individual will stop at the water dispenser (visible in camera 4) rather than turning right (visible in camera 3). This result highlights the reliance on training data and importance therefore of gathering sufficient data to cover a variety of scenarios. The predictions are made in a 16×9 grid which are smoothed with a Gaussian kernel for cleaner visualization.

REFERENCES

- [1] Olly Styles, Tanaya Guha, and Victor Sanchez. Multiple object forecasting: Predicting future object locations in diverse environments. In *Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2020.
- [2] Takuma Yagi, Karttikeya Mangalam, Ryo Yonetani, and Yoichi Sato. Future person localization in first-person videos. In *Computer Vision and Pattern Recognition*, 2018.