# GD-GAN: Generative Adversarial Networks for Trajectory Prediction and Group Detection in Crowds

Tharindu Fernando[1], Simon Denman[1], Sridha Sridharan[1], and Clinton Fookes[1]

Image and Video Research Laboratory, SAIVT, Queensland University of Technology (QUT), Australia.
{t.warnakulasuriya, s.denman, s.sridharan, c.fookes}@qut.edu.au

**Abstract.** This paper presents a novel deep learning framework for human trajectory prediction and detecting social group membership in crowds. We introduce a generative adversarial pipeline which preserves the spatio-temporal structure of the pedestrian's neighbourhood, enabling us to extract relevant attributes describing their social identity. We formulate the group detection task as an unsupervised learning problem, obviating the need for supervised learning of group memberships via hand labeled databases, allowing us to directly employ the proposed framework in different surveillance settings. We evaluate the proposed trajectory prediction and group detection frameworks on multiple public benchmarks, and for both tasks the proposed method demonstrates its capability to better anticipate human sociological behaviour compared to the existing state-of-the-art methods. [1]

**Keywords:** Group detection · Generative Adversarial Networks · Trajectory Prediction.

## 1 Introduction

Understanding and predicting crowd behaviour plays a pivotal role in video based surveillance; and as such is becoming essential for discovering public safety risks, and predicting crimes or patterns of interest. Recently, focus has been given to understanding human behaviour at a group level, leveraging observed social interactions. Researchers have shown this to be important as interactions occur at a group level, rather than at an individual or whole of crowd level.

As such we believe group detection has become a mandatory part of an intelligent surveillance system; however this group detection task presents several new challenges [31,32]. Other than identifying and tracking pedestrians from video, modelling the semantics of human social interaction and cultural gestures over a short sequence of clips is extremely challenging. Several attempts [27,31,32,34]

---

have been made to incorporate handcrafted physics based features such as relative distance between pedestrians, trajectory shape and motion based features to model their social affinity. Hall et. al [16] proposed a proxemic theory for such physical interactions based on different distance boundaries; however recent works [31, 32] have shown these quantisations fail in cluttered environments.

Furthermore, proximity doesn't always describe the group membership. For instance two pedestrians sharing a common goal may start their trajectories in two distinct source positions, however, meet in the middle. Hence we believe being reliant on a handful of handcrafted features to be sub-optimal [1, 10, 19].
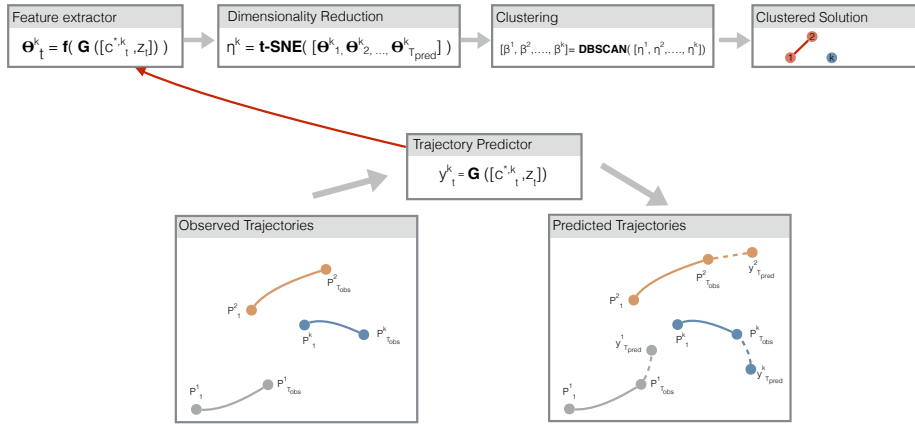


**Fig. 1.** Proposed group detection framework: After observing short segments of trajectories for each pedestrian in the scene, we apply the proposed trajectory prediction algorithm to forecast their future trajectories. The context representation generated at this step is extracted and compressed using t-SNE dimensionality reduction. Finally, the DBSCAN clustering algorithm is applied to detect the pedestrian groups.

To this end we propose a deep learning algorithm which automatically learns these group attributes. We take inspiration from the trajectory modelling approaches of [8] and [11], where the approaches capture contextual information from the local neighbourhood. We further augment this approach with a Generative Adversarial Network (GAN) [10, 15, 28] learning pipeline where we learn a custom, task specific loss function which is specifically tailored for future trajectory prediction, learning to imitate complex human behaviours.

Fig. 1 illustrates the proposed approach. First, we observe short segments of trajectories from 1 to $T_{obs}$ for each pedestrian, $p^k$, in the scene. Then, we apply the proposed trajectory prediction algorithm to forecast their future trajectories from $T_{obs+1} - T_{pred}$. This step generates hidden context representations for each pedestrian describing the current environmental context in the local neighbourhood of the pedestrian. We then apply t-SNE dimensionality reduction to

extract the most discriminative features, and we detect the pedestrian groups by clustering these reduced features.

The simplistic nature of the proposed framework offers direct transferability among different environments when compared to the supervised learning approaches of [27,31,32,34], which require re-training of the group detection process whenever the surveillance scene changes. This ability is a result of the proposed deep feature learning framework which learns the required group attributes automatically and attains commendable results among the state-of-the-art.

Novel contributions of this paper can be summarised as follows:

- We propose a novel GAN pipeline which jointly learns informative latent features for pedestrian trajectory forecasting and group detection.
- We remove the supervised learning requirement for group detection, allowing direct transferability among different surveillance scenes.
- We demonstrate how the original GAN objective could be augmented with sparsity regularisation to learn powerful features which are informative to both trajectory forecasting and group detection tasks.
- We provide extensive evaluations of the proposed method on multiple public benchmarks where the proposed method is able to generate notable performance, especially among unsupervised learning based methods.
- We present visual evidence on how the proposed trajectory modelling scheme has been able to embed social interaction attributes into its encoding scheme.

## 2  Related Work

Related literature is categorised into human behaviour prediction approaches (see Sec. 2.1); and group detection architectures (see Sec. 2.2).

### 2.1  Human Behaviour Prediction

Social Force models [17, 34], which rely on the attractive and repulsive forces between pedestrians to model their future behaviour, have been extensively applied for modelling human navigational behaviour. However with the dawn of deep learning, these methods have been replaced as they have been shown to ill represent the structure of human decision making [7, 8, 15].

One of the most popular deep learning methods is the social LSTM [1] model which represents the pedestrians in the local neighbourhood using LSTMs and then generates their future trajectory by systematically pooling the relevant information. This removes the need for handcrafted features and learns the required feature vectors automatically through the encoded trajectory representation. This architecture is further augmented in [8] where the authors propose a more efficient method to embed the local neighbourhood information via a soft and hardwired attention framework. They demonstrate the importance of fully capturing the context information, which includes the short-term history of the pedestrian of interest as well as their neighbours.

Generative Adversarial Networks (GANs) [10, 15, 28] propose a task specific loss function learning process where the training objective is a minmax game between the generative and discriminative models. These methods have shown promising results, overcoming the intractable computation of a loss function, in tasks such as autonomous driving [9, 23], saliency prediction [10, 25], image to image translation [19] and human navigation modelling [15, 28].

Even though the proposed GAN based trajectory modelling approach exhibits several similarities to recent works in [15, 28], the proposed work differs in multiple aspects. Firstly, instead of using CNN features to extract the local structure of the neighbourhood as in [28], pooling out only the current state of the neighbourhood as in [15], or discarding the available historical behaviour which is shown to be ineffective [7, 8, 28]; we propose an efficient method to embed the local neighbourhood context based on the soft and hardwired attention framework proposed in [8]. Secondly, as we have an additional objective of localising the groups in the given crowd, we propose an augmentation to the original GAN objective which regularises the sparsity of the generator embeddings, generating more discriminative features and aiding the clustering processes.

## 2.2   Group Detection

Some earlier works in group detection [5, 29] employ the concept of F-formations [20], which can be seen as specific orientation patterns that individuals engage in when in a group. However such methods are only suited to stationary groups.

In a separate line of work researchers have analysed pedestrian trajectories to detect groups. Pellegrinin et. al [27] applied Conditional Random Fields to jointly predict the future trajectory of the pedestrian of interest as well as their group membership. [34] utilises distance, speed and overlap time to train a linear SVM to classify whether two pedestrians are in the same group or not. In contrast to these supervised methods, Ge et. al [13] proposed using agglomerative clustering of speed and proximity features to extract pedestrian groups.

Most recently Solera et. al [31] proposed proximity and occupancy based social features to detect groups using a trained structural SVM. In [32] the authors extend this preliminary work with the introduction of sociologically inspired features such as path convergence and trajectory shape. However these supervised learning mechanisms rely on hand labeled datasets to learn group segmentation, limiting the methods applicability. Furthermore, the above methods all utilise a predefined set of handcrafted features to describe the sociological identity of each pedestrian, which may be suboptimal. Motivated by the impressive results obtained in [8] with the augmented context embedding, we make the first effort to learn group attributes automatically and jointly through trajectory prediction.

## 3   Architecture

### 3.1   Neighbourhood Modelling

We use the trajectory modelling framework of [8] (shown in Fig. 2) for modelling the local neighbourhood of the pedestrian of interest.
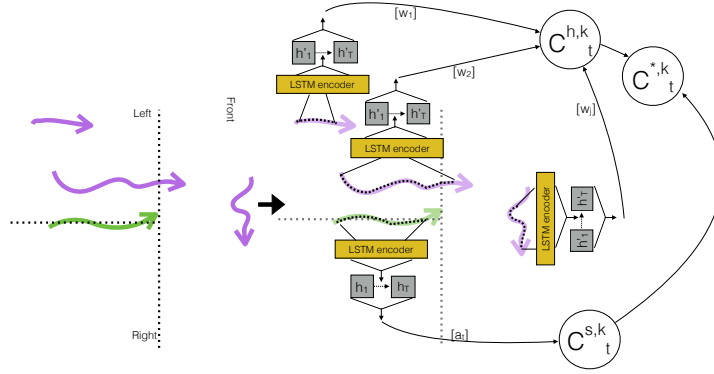
**Fig. 2.** Proposed neighbourhood modelling scheme [8]: A sample surveillance scene is shown on the left. The trajectory of the pedestrian of interest, $k$ , is shown in green, and has two neighbours (in purple) to the left, one in front and none on right. The neighbourhood encoding scheme shown on the right: Trajectory information is encoded with LSTM encoders. A soft attention context vector $C_t^{s,k}$ is used to embed trajectory information from the pedestrian of interest, and a hardwired attention context vector $C_t^{h,k}$ is used for neighbouring trajectories. In order to generate $C_t^{s,k}$ we use a soft attention function denoted $a_t$ in the above figure, and the hardwired weights are denoted by $w$. The merged context vector $C_t^{*,k}$ is then generated by merging $C_t^{s,k}$ and $C_t^{h,k}$.

Let the trajectory of the pedestrian $k$, from frame 1 to $T_{obs}$ be given by,

$$p^k = [p_1, \ldots, p_{T_{obs}}], \tag{1}$$

where the trajectory is composed of points in a Cartesian grid. Then we pass each trajectory through an LSTM [18] encoder to generate its hidden embeddings,

$$h_t^k = LSTM(p_t^k, h_{t-1}^k), \tag{2}$$

generating a sequence of embeddings,

$$h^k = [h_1^k, \ldots, h_{T_{obs}}^k]. \tag{3}$$

Following [8], the trajectory of the pedestrian of interest is embedded with soft attention such that,

$$C_t^{s,k} = \sum_{j=1}^{T_{obs}} \alpha_{tj} h_j^k, \tag{4}$$

which is the weighted sum of hidden states. The weight $\alpha_{tj}$ is computed by,

$$\alpha_{tj} = \frac{exp(e_{tj})}{\sum_{l=1}^{T} exp(e_{tl})}, \tag{5}$$

$$e_{tj} = a(h_{t-1}^k, h_j^k). \tag{6}$$

The function $a$ is a feed forward neural network jointly trained with the other components.

To embed the effect of the neighbouring trajectories we use the hardwired attention context vector $C_t^{h,k}$ from [8]. The hardwired weight $w$ is computed by,

$$w_j^n = \frac{1}{\text{dist}(n,j)},$$

(7)

where $\text{dist}(n,j)$ is the distance between the $n^{th}$ neighbour and the pedestrian of interest at the $j^{th}$ time instant. Then we compute $C_t^{h,k}$ as the aggregation for all the neighbours such that,

$$C_t^{h,k} = \sum_{n=1}^{N} \sum_{j=1}^{T_{obs}} w_j^n h_j^n,$$

(8)

where there are $N$ neighbouring trajectories in the local neighbourhood, and $h_j^n$ is the encoded hidden state of the $n^{th}$ neighbour at the $j^{th}$ time instant. Finally we merge the soft attention and hardwired attention context vectors to represent the current neighbourhood context such that,

$$C_t^{*,k} = \tanh([C_t^{s,k}, C^{h,k}]).$$

(9)

### 3.2   Trajectory Prediction

Unlike [8], we use a GAN to predict the future trajectory. There exists a minmax game between the generator (G) and the discriminator (D) guiding the model G to be closer to the ground truth distribution. The process is guided by learning a custom loss function which generates an additional advantage when modelling complex behaviours such as human navigation, where multiple factors such as human preferences and sociological factors influence behaviour.

Trajectory prediction can be formulated as observing the trajectory from time 1 to $T_{obs}$, denoted as $[p_1, \ldots, p_{T_{obs}}]$, and forecasting the future trajectory for time $T_{obs+1}$ to $T_{pred}$, denoted as $[y_{T_{obs+1}}, \ldots, y_{T_{pred}}]$. The GAN learns a mapping from a noise vector $z$ to an output vector y, $G : z \to y$ [10]. Adding the notion of time, the output of the model $y_t$ can be written as $G : z_t \to y_t$.

We augment the generic GAN mapping to be conditional on the current neighbourhood context $C_t^*$, $G : (C_t^*, z_t) \to y_t$, such that the synthesised trajectories follow the social navigational rules that are dictated by the environment.

This objective can be written as,

$$V = \mathbb{E}_{y_t, C_t^* \sim p_{data}}([logD(C_t^*, y_t)]) + \mathbb{E}_{C_t^* \sim p_{data}, z_t \sim noise}([1 - logD(C_t^*, G(C_t^*, z_t))]).$$

(10)

Our final aim is to utilise the hidden state embeddings from the trajectory generator to discover the pedestrian groups via clustering those embeddings. Hence having a sparse feature vector for clustering is beneficial as they are more

discriminative compared to their dense counterparts [12]. Hence we augment the objective in Eq. 10 with a sparsity regulariser such that,

$$L_1 = ||f(G(C_t^*, z_t))||_1, \tag{11}$$

and

$$V^* = V + \lambda L_1, \tag{12}$$

where $f$ is a feature extraction function which extracts the hidden embeddings from the trajectory generator $G$, and $\lambda$ is a weight vector which controls the tradeoff between the GAN objective and the sparsity constraint.



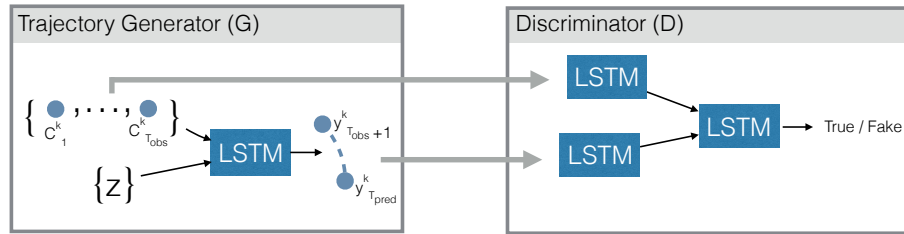**Fig. 3.** Proposed trajectory prediction framework: The generator model $G$ samples from the noise distribution $z$ and synthesises a trajectory $y_t$, which is conditioned upon the local neighbourhood context $C_t^*$. The discriminator $D$ considers both $y_t$ and $C_t^*$ when classifying the authenticity of the trajectory.

The architecture of the proposed trajectory prediction framework is presented in Fig. 3. We utilise LSTMs as the Generator $(G)$ and the Discriminator $(D)$ models. $G$ samples from the noise distribution, $z$, and synthesises a trajectory for the pedestrian motion which is conditioned upon the local neighbourhood context, $C_t^*$, of that particular pedestrian. Utilising these predicted trajectories, $y_t$, and the context embeddings, $C_t^*$, $D$ tries to discriminate between the synthesised and ground truth human trajectories.

### 3.3   Group Detection

Fig. 1 illustrates the proposed group detection framework. We pass each trajectory in the given scene through Eq. 2 to Eq. 9 and generate the neighbourhood embeddings, $C_t^{*,k}$. Then using the feature extraction function $f$ we extract the hidden layer activations for each pedestrian $k$ such that,

$$\theta_t^k = f(G(C_t^{*,k}, z_t)). \tag{13}$$

Then we pass the extracted feature vectors through a t-SNE [24] dimensionality reduction step. The authors in [12] have shown that it is inefficient to cluster dense deep features. However they have shown the t-SNE algorithm to generate

discriminative features capturing the salient aspects in each feature dimension. Hence we apply t-SNE for the $k^{th}$ pedestrian in the scene such that,

$$\eta^k = \text{t-SNE}([\theta_1^k, \ldots, \theta_{T_{obs}}^k]). \tag{14}$$

As the final step we apply DBSCAN [6] to discover similar activation patterns, hence segmenting the pedestrian groups. DBSCAN enables us to cluster the data on the fly without specifying the number of clusters. The process can be written as,

$$[\beta^1, \ldots, \beta^N] = \text{DBSCAN}([\eta^1, \ldots, \eta^N]), \tag{15}$$

where there are $N$ pedestrians in the given scene and $\beta^n \in [\beta^1, \ldots, \beta^N]$ are the generated cluster identities.

## 4    Evaluation and Discussion

### 4.1    Implementation Details

When encoding the neighbourhood information, similar to [8], we consider the closest 10 neighbours from each of the left, right, and front directions of the pedestrian of interest. If there are more than 10 neighbours in any direction, we take the closest 9 trajectories and the mean trajectory of the remaining neighbours. If a trajectory has less than 10 neighbours, we created dummy trajectories with hardwired weights (i.e Eq. 7) of 0, such that we always have 10 neighbours.

For all LSTMs, including LSTMs for neighbourhood modelling (i.e Sec. 3.1), the trajectory generator and the discriminator (i.e Sec 3.2), we use a hidden state embedding size of 300 units. We trained the trajectory prediction framework iteratively, alternating between a generator epoch and a discriminator epoch with the Adam [21] optimiser, using a mini-batch size of 32 and a learning rate of 0.001 for 500 epochs. The hyper parameter $\lambda = 0.2$, and the hyper parameters of DBSCAN, epsilon= 0.50, minPts= 1, are chosen experimentally.

### 4.2    Evaluation of the Trajectory Prediction

**Datasets** We evaluate the proposed trajectory predictor framework on the publicly available walking pedestrian dataset (BIWI) [26], Crowds By Examples (CBE) [22] dataset and Vittorio Emanuele II Gallery (VEIIG) dataset [3]. The BIWI dataset records two scenes, one outside a university (ETH) and one at a bus stop (Hotel). CBE records a single video stream with a medium density crowd outside a university (Student 003). The VEIIG dataset provides one video sequence from an overhead camera in the Vittorio Emanuele II Gallery (gall). The training, testing and validation splits for BIWI, CBE and VEIIG are taken from [26], [31] and [32] respectively.

These datasets include a variety of pedestrian social navigation scenarios including collisions, collision avoidance and group movements, hence presenting challenging settings for evaluation. Compared to BIWI which has low crowd densities, CBE and VEIIG contain higher crowd densities and as a result more challenging crowd behaviour arrangements, continuously varying from medium to high densities.

**Evaluation Metrics** Similar to [15, 28] we evaluated the trajectory prediction performance with the following 2 error metrics: Average Displacement Error (ADE) and Final Displacement Error (FDE). Please refer to [15, 28] for details.

**Baselines and Evaluation** We compared our trajectory prediction model to 5 state-of-the-art baselines. As the first baseline we use the Social Force (SF) model introduced in [34], where the destination direction is taken as an input to the model and we train a linear SVM model similar to [8] to generate this input. We use the Social-LSTM (So-LSTM) model of [1] as the next baseline and the neighbourhood size hyper-parameter is set to 32 px. We also compare to the Soft + Hardwired Attention (SHA) model of [8] and similar to the proposed model we set the embedding dimension to be 300 units and consider a 30 total neighbouring trajectories. We also considered the Social GAN (So-GAN) [15] and attentive GAN (SoPhie) [28] models. To provide fair comparisons we set the hidden state dimensions for the encoder and decoder models of So-GAN and SoPhie to be 300 units. For all models we observe the first 15 frames (i.e 1- $T_{obs}$) and predicted the future trajectory for the next 15 frames (i.e $T_{obs+1}$ - $T_{pred}$).

**Table 1.** Quantitative results for the BIWI [26], CBE [22] and VEIIG [3] datasets. In all methods the forecast trajectories are of length 15 frames. Error metrics are as in Sec. 4.2. '-' refers to unavailability of that specific evaluation. The best values are denoted in bold.

| Metric | Dataset | SF [34] | So-LSTM [1] | SHA [8] | So-GAN [15] | SoPhie [28] | Proposed |
|---|---|---|---|---|---|---|---|
| ADE | ETH (BIWI) | 1.42 | 1.05 | 0.90 | 0.92 | 0.81 | **0.63** |
| | Hotel (BIWI) | 1.03 | 0.98 | 0.71 | 0.65 | 0.76 | **0.55** |
| | Student 003 (CBE) | 1.83 | 1.22 | 0.96 | - | - | **0.72** |
| | gall (VEIIG) | 1.72 | 1.14 | 0.91 | - | - | **0.68** |
| FDE | ETH (BIWI) | 2.20 | 1.84 | 1.43 | 1.52 | 1.45 | **1.22** |
| | Hotel (BIWI) | 2.45 | 1.95 | 1.65 | 1.62 | 1.77 | **1.43** |
| | Student 003 (CBE) | 2.63 | 1.97 | 1.80 | - | - | **1.65** |
| | gall (VEIIG) | 2.55 | 1.83 | 1.65 | - | - | **1.45** |

When observing the results tabulated in Tab. 1 we observe poor performance for the SF model due to it's lack of capacity to model history. Models So-LSTM and SHA utilise short term history from the pedestrian of interest and the local neighbourhood and generate improved predictions. However we observe a significant increase in performance from methods that optimise generic loss functions such as So-LSTM and SHA to GAN based methods such as So-GAN and So-Phie. This emphasises the need for task specific loss function learning in order to imitate complex human social navigation strategies. In the proposed method we further augment this performance by conditioning the trajectory generator on the proposed neighbourhood encoding mechanism.

We present a qualitative evaluation of the proposed trajectory generation framework with the SHA and So-GAN baselines in Fig. 4 (selected based on the availability of their implementations). The observed portion of the trajectory is

denoted in green, the ground truth observations in blue and predicted trajectories are shown in red (proposed), yellow (SHA) and brown (So-GAN). Observing the qualitative results it can be clearly seen that the proposed model generates better predictions compared to the state-of-the-art considering the varying nature of the neighbourhood clutter. For instance in Fig. 4 (c) and (d) we observe significant deviations between the predictions for SHA and So-GAN and the ground truth. However the proposed model better anticipates the pedestrian motion with the improved context modelling and learning process. It should be noted that the proposed method has a better ability to anticipate stationary groups compared to the baselines, which is visible in Fig. 4 (c).
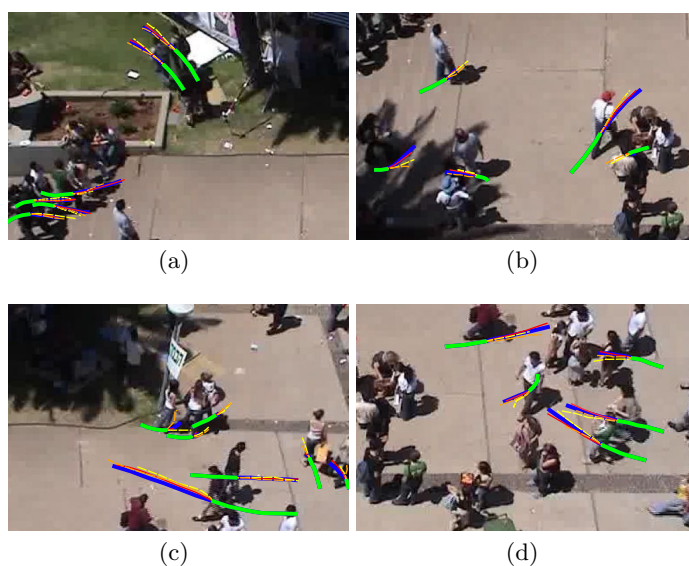


**Fig. 4.** Qualitative results for the proposed trajectory prediction framework for sequences from the CBE dataset. Given (in green), Ground Truth (in blue) and Predicted trajectories from proposed (in red), SHA model (in yellow) crom So-GAN (in brown). For visual clarity, we show only the trajectories for some of the pedestrians in the scene.

### 4.3   Evaluation of the Group Detection

**Datasets** Similar to Sec. 4.2 we use the BIWI, CBE and VEIIG datasets in our evaluation. Dataset characteristics are reported in Tab. 2.

**Evaluation Metrics** One popular measure of clustering accuracy is the pairwise loss $\Delta_{pw}$ [35], which is defined as the ratio between the number of pairs

**Table 2.** Dataset characteristics for different sequences in BIWI [26], CBE [22] and VEIIG [3] datasets

| Dataset | ETH (BIWI) | Hotel (BIWI) | Student-003 (CBE) | gall (VEIIG) |
|---------|-----------|--------------|-------------------|--------------|
| Frames | 1448 | 1168 | 541 | 7500 |
| Pedestrian | 360 | 390 | 434 | 630 |
| Groups | 243 | 326 | 288 | 207 |

on which $\beta$ and $\hat{\beta}$ disagree on their cluster membership and the number of all possible pairs of elements in the set.

However as described in [31, 32] $\Delta_{pw}$ accounts only for positive intra-group relations and neglects singletons. Hence we also measure the Group-MITRE loss, $\Delta_{GM}$, introduced in [31], which has overcome this deficiency. $\Delta_{GM}$ adds a fake counterpart for singletons and each singleton is connected with it's counterpart. Therefore $\delta_{GM}$ also takes singletons into consideration.

**Baselines and Evaluation** We compare the proposed Group Detection GAN (GD-GAN) framework against 5 recent state-of-the-art baselines, namely [13, 30, 32, 34, 35], selected based on their reported performance in public benchmarks.

In Tab. 3 we report the Precision ($P$) and Recall ($R$) values for $\Delta_{pw}$ and $\Delta_{GM}$ for the proposed method along with the state-of-the-art baselines. The proposed GD-GAN method has been able to achieve superior results, especially among unsupervised grouping methods. It should be noted that methods [30, 32, 34, 35] utilise handcrafted features and use supervised learning to separate the groups. As noted in Sec. 1 these methods cannot adapt to scene variations and require hand labeled datasets for training. Furthermore we would like to point out that the supervised grouping mechanism in [32] directly optimises $\Delta_{GM}$. However, without such tedious annotation requirements and learning strategies, the proposed method has been able to generate commendable and consistent results in all considered datasets, especially in cluttered environments [2].

In Fig. 5 we show groups detected by the proposed GD-GAN method for sequences from the CBE and VEIIG datasets. Regardless of the scene context, occlusions and the varying crowd densities, the proposed GD-GAN method generates acceptable results. We believe this is due to the augmented features that we derive through the automated deep feature learning process. These features account for both historical and future behaviour of the individual pedestrians, hence possessing an ability to detect groups even in the presence of occlusions such as in Fig. 5 (c).

We selected the first 30 pedestrian trajectories from the VEIIG test set and in Fig. 6 we visualise the embedding space positions before (in blue) and after (in red) training of the proposed trajectory generator (G). Similar to [2] we extracted the activations using the feature extractor function $f$ and applied PCA

---

[2] see the supplementary material for the results for using supervised learning to separate the groups on proposed context features.

**Table 3.** Comparative results on the BIWI [26], CBE [22] and VEIIG [3] datasets using the $\Delta_{GM}$ [31] and $\Delta_{PW}$ [35] metrics. '-' refers to unavailability of that specific evaluation. The best results are shown in bold and the second best results are underlined.

| | | Shao et. al [30] | | zanotto et. al [35] | | Yamaguchi et. al [34] | | Ge et. al [13] | | Solera et al. [32] | | GD-GAN | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $P$ | $R$ | $P$ | $R$ | $P$ | $R$ | $P$ | $R$ | $P$ | $R$ | $P$ | $R$ |
| BIWI | $\Delta_{GM}$ | 67.3 | 64.1 | - | - | 84.0 | 51.2 | 89.2 | 90.9 | <u>97.3</u> | **97.7** | **97.5** | **97.7** |
| Hotel | $\Delta_{PW}$ | 51.5 | 90.4 | 81.0 | 91.0 | 83.7 | **93.9** | 88.9 | 89.3 | <u>89.1</u> | 91.9 | **90.2** | <u>93.1</u> |
| BIWI | $\Delta_{GM}$ | 69.3 | 68.2 | - | - | 60.6 | 76.4 | 87.0 | 84.2 | <u>91.8</u> | **94.2** | **92.5** | **94.2** |
| ETH | $\Delta_{PW}$ | 44.5 | 87.0 | 79.0 | 82.0 | 72.9 | 78.0 | 80.7 | 80.7 | <u>91.1</u> | <u>83.4</u> | **91.3** | **83.5** |
| CEB | $\Delta_{GM}$ | 40.4 | 48.6 | - | - | 56.7 | 76.0 | 77.2 | 73.6 | **81.7** | **82.5** | <u>81.0</u> | <u>81.8</u> |
| Student-003 | $\Delta_{PW}$ | 10.6 | **76.0** | 70.0 | 74.0 | 63.9 | 72.6 | 72.2 | 65.1 | **82.3** | <u>74.1</u> | <u>82.1</u> | 63.4 |
| VEIIG | $\Delta_{GM}$ | - | - | - | - | - | - | - | - | **84.1** | **84.1** | <u>83.1</u> | <u>79.5</u> |
| gall | $\Delta_{PW}$ | - | - | - | - | - | - | - | - | **79.7** | **77.5** | <u>77.6</u> | <u>73.1</u> |

[33] to plot them in 2D. The respective ground truth group IDs are indicated in brackets. This helps us to gain an insight into the encoding process that $G$ utilises, which allows us to discover groups of pedestrians. Considering the examples given, it can be seen that trajectories from the same cluster become more tightly grouped. This is due to the model incorporating source positions, heading direction, trajectory similarity, when embedding trajectories, allowing us to extract pedestrian groups in an unsupervised manner.

### 4.4 Ablation Experiment

To further demonstrate the proposed group detection approach, we conducted a series of ablation experiments identifying the crucial components of the proposed methodology [3]. In the same setting as the previous experiment we compare the proposed GD-GAN model against a series of counter parts as follows:

- GD-GAN / GAN: removes $D$ and the model $G$ is learnt through supervised learning as in [8].
- GD-GAN / cGAN: optimises the generic GAN objective defined in [14].
- GD-GAN / $L_1$: removes sparsity regularisation and optimises Eq. 10.
- GD-GAN + hf: utilises features from $G$ as well as the handcrafted features defined in [32] for clustering.

**Table 4.** Ablation experiment evaluations

| | | GD-GAN / GAN | | GD-GAN / cGAN | | GD-GAN / $L_1$ | | GD-GAN + hf | | GD-GAN | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | P | R | P | R | P | R | P | R | P | R |
| CEB | $\Delta_{GM}$ | 73.6 | 75.1 | 76.7 | 76.2 | 77.3 | 78.0 | **79.0** | **79.2** | 78.7 | **79.2** |
| Student-003 | $\Delta_{PW}$ | 74.1 | 52.8 | 75.5 | 60.2 | 78.1 | 65.1 | **80.4** | 68.0 | **80.4** | **68.4** |

The results of our ablation experiment are presented in Tab. 4. Model GD-GAN / GAN performs poorly due to the deficiencies in the supervised learning process. It optimises a generic mean square error loss, which is not ideal to guide

---

[3] see the supplementary material for an ablation study for the trajectory prediction

(a) GVEII - Frame 2127          (b) GVEII- Frame 2320

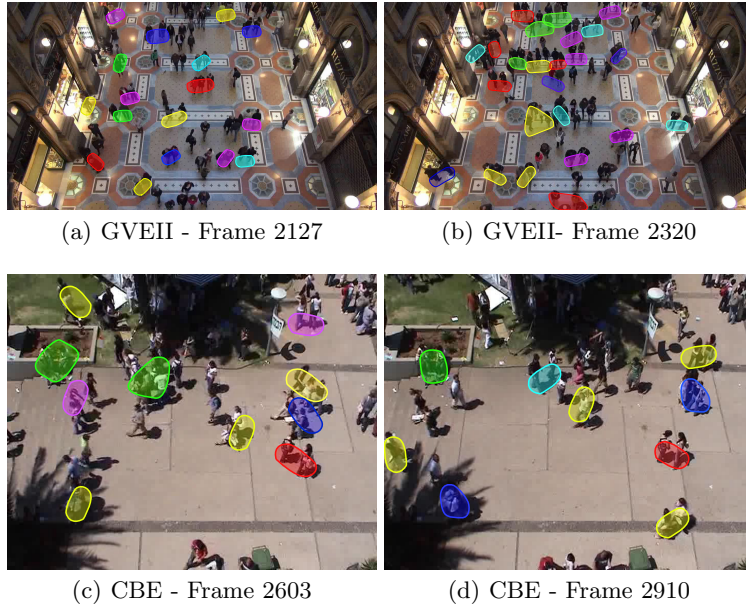(c) CBE - Frame 2603          (d) CBE - Frame 2910

**Fig. 5.** Qualitative results from the proposed GD-GAN methods for sequences from the CBE and GVEII datasets. Connected coloured blobs indicate groups of pedestrians.

the model through the learning process when modelling a complex behaviour such as human navigation. Therefore the resultant feature vectors do not capture the full context which contributes to the poor group detection accuracies. We observe an improvement in performance with GD-GAN / cGAN due to the GAN learning process which is further augmented and improved through GD-GAN / $L_1$ where the model learns a conditional behaviour depending on the neighbourhood context. $L_1$ regularisation further assists the group detection process via making the learnt feature distribution more discriminative.

In order to demonstrate the credibility of the learnt group attributes from the proposed GD-GAN model, we augment the feature vector extracted in Eq. 13 together with the features proposed in [32] and apply subsequent process (i.e Eq. 14 and 15) to discover the groups. We utilise the public implementation [4] released by the authors for the feature extraction.

We do not observe a substantial improvement with the group detection performance being very similar, indicating that the proposed GD-GAN model is sufficient for modelling the social navigation structure of the crowd.
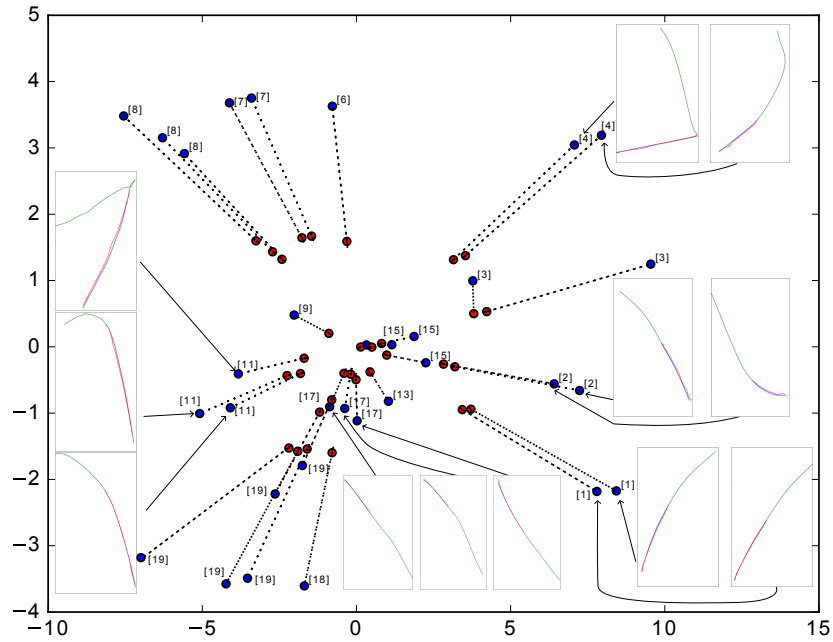
---

[4] https://github.com/francescosolera/group-detection

**Fig. 6.** Projections of the trajectory generator (G) hidden states before (in blue) and after (in red) training. Ground truth group IDs are in brackets. Each insert indicates the trajectory associated with the embedding. The given portion of the trajectory is in green, and the ground truth and prediction are in blue and red respectively

### 4.5   Time efficiency

We use the Keras [4] deep learning library for our implementation. The GD-GAN module does not require any special hardware such as GPUs to run and has 41.8K trainable parameters. We ran the test set in Sec. 4.3 on a single core of an Intel Xeon E5-2680 2.50GHz CPU and the GD-GAN algorithm was able to generate 100 predicted trajectories with 30, 2 dimensional data points in each trajectory (i.e. using 15 observations to predict the next 15 data points) and complete the group detection process in 0.712 seconds.

## 5   Conclusions

In this paper we have proposed an unsupervised learning approach for pedestrian group segmentation. We avoid the the need to handcraft sociological features by automatically learning group attributes through the proposed trajectory prediction framework. This allows us to discover a latent representation accounting for both historical and future behaviour of each pedestrian, yielding a more efficient platform for detecting their social identities. Furthermore, the unsupervised

learning setting grants the approach the ability to employ the proposed framework in different surveillance settings without tedious learning of group memberships from a hand labeled dataset. Our quantitative and qualitative evaluations on multiple public benchmarks clearly emphasise the capacity of the proposed GD-GAN method to learn complex real world human navigation behaviour.

## References

1. Alahi, A., Goel, K., Ramanathan, V., Robicquet, A., Fei-Fei, L., Savarese, S.: Social lstm: Human trajectory prediction in crowded spaces. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 961–971 (2016)
2. Aubakirova, M., Bansal, M.: Interpreting neural networks to improve politeness comprehension. arXiv preprint arXiv:1610.02683 (2016)
3. Bandini, S., Gorrini, A., Vizzari, G.: Towards an integrated approach to crowd analysis and crowd synthesis: A case study and first results. Pattern Recognition Letters **44**, 16–29 (2014)
4. Chollet, F., et al.: Keras (2015) (2017)
5. Cristani, M., Bazzani, L., Paggetti, G., Fossati, A., Tosato, D., Del Bue, A., Menegaz, G., Murino, V.: Social interaction discovery by statistical analysis of f-formations. In: BMVC. vol. 2, p. 4 (2011)
6. Ester, M., Kriegel, H.P., Sander, J., Xu, X., et al.: A density-based algorithm for discovering clusters in large spatial databases with noise. In: Kdd. vol. 96, pp. 226–231 (1996)
7. Fernando, T., Denman, S., McFadyen, A., Sridharan, S., Fookes, C.: Tree memory networks for modelling long-term temporal dependencies. Neurocomputing **304**, 64–81 (2018)
8. Fernando, T., Denman, S., Sridharan, S., Fookes, C.: Soft+ hardwired attention: An lstm framework for human trajectory prediction and abnormal event detection. arXiv preprint arXiv:1702.05552 (2017)
9. Fernando, T., Denman, S., Sridharan, S., Fookes, C.: Learning temporal strategic relationships using generative adversarial imitation learning. arXiv preprint arXiv:1805.04969 (2018)
10. Fernando, T., Denman, S., Sridharan, S., Fookes, C.: Task specific visual saliency prediction with memory augmented conditional generative adversarial networks. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 1539–1548. IEEE (2018)
11. Fernando, T., Denman, S., Sridharan, S., Fookes, C.: Tracking by prediction: A deep generative model for mutli-person localisation and tracking. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 1122–1132. IEEE (2018)
12. Figueroa, J.A., Rivera, A.R.: Learning to cluster with auxiliary tasks: A semi-supervised approach. In: Graphics, Patterns and Images (SIBGRAPI), 2017 30th SIBGRAPI Conference on. pp. 141–148. IEEE (2017)
13. Ge, W., Collins, R.T., Ruback, R.B.: Vision-based analysis of small groups in pedestrian crowds. IEEE Transactions on Pattern Analysis and Machine Intelligence **34**(5), 1003–1016 (May 2012). https://doi.org/10.1109/TPAMI.2011.176
14. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in neural information processing systems. pp. 2672–2680 (2014)

15. Gupta, A., Johnson, J., Fei-Fei, L., Savarese, S., Alahi, A.: Social gan: Socially acceptable trajectories with generative adversarial networks. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). No. CONF (2018)
16. Hall, E.T.: The hidden dimension (1966)
17. Helbing, D., Molnar, P.: Social force model for pedestrian dynamics. Physical review E **51**(5),  4282 (1995)
18. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation **9**(8), 1735–1780 (1997)
19. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. arXiv preprint (2017)
20. Kendon, A.: Conducting interaction: Patterns of behavior in focused encounters, vol. 7. CUP Archive (1990)
21. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
22. Lerner, A., Chrysanthou, Y., Lischinski, D.: Crowds by example. In: Computer Graphics Forum. vol. 26, pp. 655–664. Wiley Online Library (2007)
23. Li, Y., Song, J., Ermon, S.: Infogail: Interpretable imitation learning from visual demonstrations. In: Advances in Neural Information Processing Systems. pp. 3815–3825 (2017)
24. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. Journal of machine learning research **9**(Nov), 2579–2605 (2008)
25. Pan, J., Canton, C., McGuinness, K., O'Connor, N.E., Torres, J., Sayrol, E., Giro-i Nieto, X.: Salgan: Visual saliency prediction with generative adversarial networks. arXiv preprint arXiv:1701.01081 (2017)
26. Pellegrini, S., Ess, A., Schindler, K., Van Gool, L.: You'll never walk alone: Modeling social behavior for multi-target tracking. In: Computer Vision, 2009 IEEE 12th International Conference on. pp. 261–268. IEEE (2009)
27. Pellegrini, S., Ess, A., Van Gool, L.: Improving data association by joint modeling of pedestrian trajectories and groupings. In: European conference on computer vision. pp. 452–465. Springer (2010)
28. Sadeghian, A., Kosaraju, V., Sadeghian, A., Hirose, N., Savarese, S.: Sophie: An attentive gan for predicting paths compliant to social and physical constraints. arXiv preprint arXiv:1806.01482 (2018)
29. Setti, F., Lanz, O., Ferrario, R., Murino, V., Cristani, M.: Multi-scale f-formation discovery for group detection. In: Image Processing (ICIP), 2013 20th IEEE International Conference on. pp. 3547–3551. IEEE (2013)
30. Shao, J., Loy, C.C., Wang, X.: Scene-independent group profiling in crowd. In: 2014 IEEE Conference on Computer Vision and Pattern Recognition. pp. 2227–2234 (June 2014). https://doi.org/10.1109/CVPR.2014.285
31. Solera, F., Calderara, S., Cucchiara, R.: Structured learning for detection of social groups in crowd. In: Advanced Video and Signal Based Surveillance (AVSS), 2013 10th IEEE International Conference on. pp. 7–12. IEEE (2013)
32. Solera, F., Calderara, S., Cucchiara, R.: Socially constrained structural learning for groups detection in crowd. IEEE transactions on pattern analysis and machine intelligence **38**(5), 995–1008 (2016)
33. Wold, S., Esbensen, K., Geladi, P.: Principal component analysis. Chemometrics and intelligent laboratory systems **2**(1-3), 37–52 (1987)
34. Yamaguchi, K., Berg, A.C., Ortiz, L.E., Berg, T.L.: Who are you with and where are you going? In: Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on. pp. 1345–1352. IEEE (2011)

35. Zanotto, M., Bazzani, L., Cristani, M., Murino, V.: Online bayesian nonparametrics for group detection. In: Proc. of BMVC (2012)