# Pedestrian Trajectory Prediction with Structured Memory Hierarchies

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

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

**Abstract.** This paper presents a novel framework for human trajectory prediction based on multimodal data (video and radar). Motivated by recent neuroscience discoveries, we propose incorporating a structured memory component in the human trajectory prediction pipeline to capture historical information to improve performance. We introduce structured LSTM cells for modelling the memory content hierarchically, preserving the spatiotemporal structure of the information and enabling us to capture both short-term and long-term context. We demonstrate how this architecture can be extended to integrate salient information from multiple modalities to automatically store and retrieve important information for decision making without any supervision. We evaluate the effectiveness of the proposed models on a novel multimodal dataset that we introduce, consisting of 40,000 pedestrian trajectories, acquired jointly from a radar system and a CCTV camera system installed in a public place. The performance is also evaluated on the publicly available New York Grand Central pedestrian database. In both settings, the proposed models demonstrate their capability to better anticipate future pedestrian motion compared to existing state of the art.

**Keywords:** Human Trajectory Prediction · Structured Memory Networks · Multimodal Information Fusion · long-term Planing.

## 1 Introduction

Understanding and predicting crowd behaviour is an important topic due to its myriad applications (surveillance, event detection, traffic flow, etc). However this remains a challenging problem due to the complex nature of human behaviour and the lack of attention that researchers pay to human navigational patterns when developing machine learning models.

Recent neuroscience studies have revealed that humans utilise map and grid like structures for navigation [12, 26]. The human brain builds a unified representation of the spatial environment, which is stored in the hippocampus [13] and guides the decision making process. Further studies [6] provide strong evidence towards a hierarchical spatial representation of these maps. Additionally in [11, 20] authors have observed multiple representations of structured maps instead of one single map in the long-term memory. This idea was explored in [28]

using structured memory for Deep Reinforcement Learning. To generate an output at a particular time step, the system passes the memory content through a series of convolution layers to summarise the content. We argue this is inefficient and could lead to a loss of information when modelling large spatial areas.

Motivated by recent neuroscience [12,26] and deep reinforcement leaning [28] studies, we utilise a structured memory to predict human navigational behaviour. In particular such a memory structure allows a machine learning algorithm to exploit historical knowledge about the spatial structure of the environment, and reason and plan ahead, instead of generating reflexive behaviour based on the current context. Novel contributions of this paper are summarised as follows:

- We introduce a novel neural memory architecture which effectively captures the spatiotemporal structure of the environment.
- We propose structured LSTM (St-LSTM) cells, which model the structured memory hierarchically, preserving the memories' spatiotemporal structure.
- We incorporate the neural memory network into a human trajectory prediction pipeline where it learns to automatically store and retrieve important information for decision making without any supervision.
- We introduce a novel multimodal dataset for human trajectory prediction, containing more than 40,000 trajectories from Radar and CCTV streams.
- We demonstrate how the semantic information from multiple input streams can be captured through multiple memory components and propose an effective fusion scheme that preserves the spatiotemporal structure.
- We provide extensive evaluations of the proposed method using multiple public benchmarks, where the proposed method is capable of imitating human navigation behaviour and outperforms state-of-the-art methods.

## 2    Related Work

The related literature can be broadly categorised into human behaviour prediction approaches, introduced in Sec 2.1; neural memory architectures, presented in Sec. 2.2; and multimodal information fusion which we review in Sec. 2.3.

### 2.1    Human Behaviour Prediction

Before the dawn of deep learning, Social Force models [33,34] had been extensively applied for modelling human navigational behaviour. They rely on the attractive and repulsive forces between pedestrians to predict motion. However as shown in [1,14,16] these methods ill represent the structure of human decision making by modelling the behaviour with just a handful of parameters.

One of the most popular deep learning methods for predicting human behaviour is the Social LSTM model of [1], which removed the need for hand-crafted features by using LSTMs to encode and decode trajectory information. This method is further augmented in [16] where the authors incorporate the

entire trajectory of the pedestrian of interest as well as the neighbouring pedestrians and extract salient information from these embeddings through a combination of soft and hardwired attention. Similar to [16] the works in [3,32,36] also highlight the importance of fully capturing context information. However these methods all consider short-term temporal context in the given neighbourhood, completely discarding scene structure and the longterm scene context.

## 2.2 Neural Memory Architectures

External memory modules are used to store historic information, and learn to automatically store and retrieve important facts to aid future predictions. Many approaches across numerous domains [14,15,17,18,22,27] have utilised memory modules to aid prediction, highlighting the importance of stored knowledge for decision making. However existing memory structures are one dimensional modules which completely ignore the environmental spatial structure. This causes a significant hindrance when modelling human navigation, since they are unable to capture the map-like structures humans use when navigating [12,26].

The work of Parisotto et al. [28] proposes an interesting extension to memory architectures where they structure the memory as a 3D block, preserving spatial relationships. However, when generating memory output they rely on a static convolution kernel to summarise the content, failing to generate dynamic responses and propagate salient information from spatial locations to the trajectory prediction module, where multiple humans can interact in the environment.

Motivated by the hierarchical sub-map structure humans use to navigate [11,20], we model our spatiotemporal memory with gated St-LSTM cells, which are arranged hierarchically in a grid structure.

## 2.3 Multimodal information fusion

Multimodal information fusion addresses the task of integrating inputs from various modalities and has shown superior performance compared to unimodal approaches [4,10] in variety of applications [2,21,35]. The simplest approach is to concatenate features to obtain a single vector representation [23,29]. However it ignores the relative correlation between the modalities [2].

More complex fusion strategies include concatenating higher level representations from individual modalities separately and then combining them together, enabling the model to learn the salient aspects of individual streams. In this direction, attempts were made using Deep Boltzmann Machines [31] and neural network architectures [9,25].

In [15] the authors explore the importance of capturing both short and longterm temporal context when performing feature fusion, utilising separate neural memory units for individual feature streams and aggregating the temporal representation during fusion. Yet this fails to preserve the spatial structure, restricting its applicability when modelling human navigation.

## 3  Architecture

In this section we introduce the encoding scheme utilised to embed the trajectory information of the pedestrian of interest and their neighbours; the structure and the operations of the proposed hierarchical memory; how to utilise memory output to enhance the future trajectory prediction; and an architecture for effectively coupling multimodal information streams through structured memories.

### 3.1  Embedding local neighbourhood context

In order to embed the current short-term context of the pedestrian of interest and the local neighbourhood, we utilise the trajectory prediction framework proposed in [16]. Let the observed trajectory of pedestrian $k$ from frame 1 to frame $T_{obs}$ be given by,

$$X^k = [(x_1, y_1), (x_2, y_2), \ldots, (x_{T_{obs}}, y_{T_{obs}})], \tag{1}$$

where the trajectory is composed of points in a 2D Cartesian grid. Similar to [16] we utilise the soft attention operation [8] to embed the trajectory information from the pedestrian of interest $(k)$ and generate a vector embedding $C_t^{s,k}$. To embed neighbouring trajectories the authors in [16,19] have shown that distance based hardwired attention is efficient and effective. We denote the hardwired context vector as $C_t^{h,k}$.

Now we define the combined context vector, $C_t^{*,k}$, representing the short-term context of the local neighbourhood of the $k^{th}$ pedestrian as,

$$C_t^{*,k} = \tanh([C_t^{s,k}, C_t^{h,k}]), \tag{2}$$

where $[.,.]$ denotes a concatenation operation. Please see [16,19] for details.

### 3.2  Structured Memory Network (SMN)

Let the structured memory, $M$, be a $l \times W \times H$ block where $l$ is the embedding dimension of $p_t^k$. $W$ is the vertical extent of the map and $H$ is the horizontal extent. We define a function $\psi(x, y)$ which maps spatial coordinates $(x, y)$ with $x \in \mathbb{R}$ and $y \in \mathbb{R}$ to a map grid $(x', y')$ where $x' \in 0, \ldots, W$ and $y' \in 0, \ldots, H$. The works of [16,19] have shown that the context embeddings $C_t^{*,k}$ capture the short-term context of the pedestrian of interest and the local neighbourhood. Hence we store these embeddings in our structured memory as it represents the temporal context in that grid cell. The operations of the proposed structured memory network (SMN) can be summarised as follows,

$$h_t = \text{read}(M_t), \tag{3}$$

$$\beta_{t+1}^{(x',y')} = \text{write}(C_t^{*,k}, M_t^{(x',y')}), \tag{4}$$

$$M_{t+1} = \text{update}(M_t, w_{t+1}^{(x',y')}). \tag{5}$$

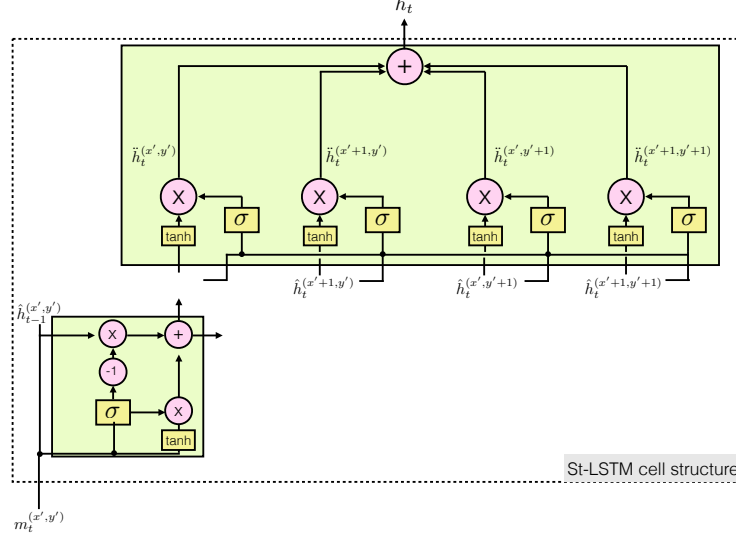The following subsections explain these three operations.

**Fig. 1.** The operations of the proposed St-LSTM cell. It considers the current representation of the respective memory cell and the 3 adjacent neighbours as well as the previous time step outputs and utilises gated operations to render the output in the present time step.

**Hierarchical Read Operation** The read operation outputs output a vector, $h_t$, capturing the most salient information from the entire memory for decision making in the present state. We define a hierarchical read operation which passes the current memory representation, $M_t$, through a series of gated, structured LSTM (St-LSTM) cells arranged in a grid like structure. Fig. 1 depicts the operations of the proposed St-LSTM cells.

Let the content of $(x', y')$ memory cell at time $t$ be represented by $m_t^{(x',y')}$ and the three adjacent cells be represented by $m_t^{(x'+1,y')}, m_t^{(x',y'+1)}$ and $m_t^{(x'+!,y'+1)}$. As shown in Fig. 2, we first pass the current state of the memory cell through an input gate to decide how much information to pass through the gate and how much information to gather from the previous hidden state of the that particular cell, $\hat{h}_{t-1}^{(x',y')}$. This operation is given by,

$$
\begin{aligned}
z_t^{(x',y')} &= \sigma(w_z^{(x',y')}[m_t^{(x',y')}, \hat{h}_{t-1}^{(x',y')}]), \\
\hat{o}_t^{(x',y')} &= \tanh([m_t^{(x',y')}, \hat{h}_{t-1}^{(x',y')}]).
\end{aligned}
\tag{6}
$$

Then we generate the new hidden state of the cell using,

$$
\hat{h}_t^{(x',y')} = z_t^{(x',y')}\hat{o}_t^{(x',y')} + (1 - z_t^{(x',y')})\hat{h}_{t-1}^{(x',y')},
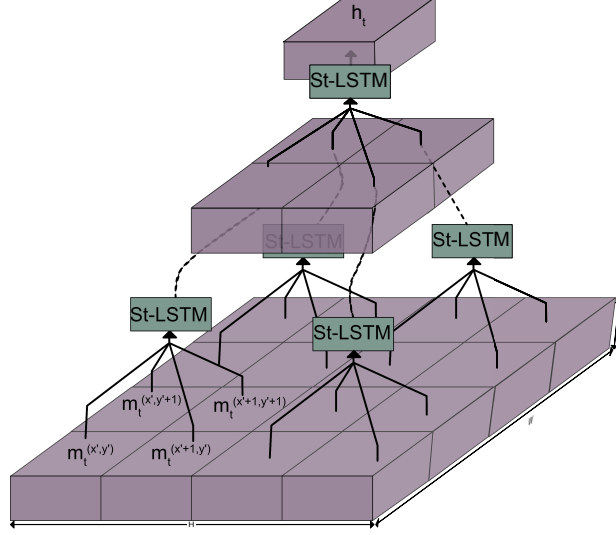\tag{7}
$$

**Fig. 2.** Utilisation of proposed St-LSTM cell to generate a hierarchical embedding of the structured memory. In each layer we summarise the content of 4 adjacent neighbours via propagating the most salient information to the layer above. The process is repeated until we generate a single vector representation of the entire memory block.

and pass the hidden state of that particular cell as well as the hidden states of the adjacent cell through a composition gate function which determines the amount of information to be gathered from each of the cells as,

$$q_t^{(x',y')} = \sigma(w_q^{(x',y')}[\hat{h}_t^{(x',y')}, \hat{h}_t^{(x'+1,y')}, \hat{h}_t^{(x',y'+1)}, \hat{h}_t^{(x'+1,y'+1)}]). \tag{8}$$

Now we can generate the augmented state of the cell $(x', y')$ as,

$$\ddot{h}_t^{(x',y')} = \tanh(\hat{h}_t^{(x',y')})q_t^{(x',y')}. \tag{9}$$

We perform the above operations to the rest of the group of 3 cells: $(x' + 1, y'), (x', y'+1)$ and $(x'+1, y'+1)$; and generate the representations $\ddot{h}_t^{(x'+1,y')}, \ddot{h}_t^{(x',y'+1)}$ and $\ddot{h}_t^{(x'+1,y'+1)}$ respectively. Then the feature embedding representation, $h_t$, of the merged 4 cells in the next layer of the memory is given by,

$$h_t = \ddot{h}_t^{(x',y')} + \ddot{h}_t^{(x'+1,y')} + \ddot{h}_t^{(x',y'+1)} + \ddot{h}_t^{(x'+1,y'+1)}. \tag{10}$$

We repeat this process for all cells $(x', y')$ where $x' \in 0, \dots, W$ and $y' \in 0, \dots, H$. Note that as we are merging four adjacent cells we have $W/2$ and $H/2$ St-LSTM cells in the immediate next layer of the memory block. We continue merging cells until we are left with one cell summarising the entire memory block. We denote the hidden state of this cell as $h_t$.

**Write Operation** Given the current position of the pedestrian of interest at $(x, y)$, we first evaluate the associated location in the map grid by passing $(x, y)$ through function, $\psi$, such that,

$$(x', y') = \psi(x, y). \tag{11}$$

Then we retrieve the current memory state of $(x', y')$ as,

$$m_t^{(x',y')} = M_t^{(x',y')}. \tag{12}$$

Then by utilising the above stated vector and the short-term context of the pedestrian of interest, $C_t^*$, we define a write function which generates a write vector for memory update,

$$\beta_{t+1}^{(x',y')} = \text{LSTM}_w(c_t^*, m_t^{(x',y')}). \tag{13}$$

**Update Operation** We update the memory map for the next time step by,

$$M_{t+1}^{(a,b)} = \begin{cases} \beta_{t+1}^{(x',y')} & \text{for } (a,b) = (x', y') \\ M_t^{(a,b)} & \text{for } (a,b) \neq (x', y') \end{cases} \tag{14}$$

The new memory map is equal to to the memory map at the previous time step except for the current location of the pedestrian where we completely update the content with the generated write vector.

### 3.3   Trajectory prediction with structured memory hierarchies

We utilise the combined context vector $C_t^{*,k}$ representing the short-term context of the local neighbourhood of the $k^{th}$ pedestrian, and the generated memory output $h_t$ to generate an augmented vector for the context representation,

$$\bar{c}_t^{(k)} = \tanh([C_t^{*,k}, h_t]), \tag{15}$$

which is used to predict the future trajectory of the pedestrian $k$,

$$Y_t = \text{LSTM}(p_{t-1}^k, \bar{c}_t^{(k)}, Y_{t-1}). \tag{16}$$

### 3.4   Coupling multimodal information to improve prediction

Using multimodal input streams allows us to capture different semantics that are present in the same scene, and compliment individual streams. For instance, in a surveillance setting, radar and video fusion is widely utilised [5,30] as radar offers better coverage in the absence of visual light, however has a lower frame rate (<5 fps) compared to video ($\sim$25fps) which records more fine grained motion.

As pointed out in [15], simply concatenating data from both streams leads to information loss as they contain information at different granularities. Hence it is vital to jointly back propagate among the information streams to learn

the important aspects of each. Therefore, we capture streams through separate memory modules, and perform gated coupling of memory hierarchies.

We denote the two synchronised input modalities as $I$ and $R$, where the trajectory of pedestrian $k$ observed in stream $I$ is denoted as $X_I^k$ and the same pedestrian trajectory observed in stream $R$ is given as $X_R^k$. We pass each stream separately through the local neighbourhood embedding mechanism proposed in Sec. 3.1 and generate vector embeddings $C_{t,I}^{*,k}$ and $C_{t,R}^{*,k}$ respectively. We embed these through individual memory blocks denoted as $M_{t,I}$ and $M_{t,R}$. In the absences of such trajectories (i.e due to poor coverage, occlusions, ... ), we evaluate only the neighbourhood embeddings $C_{t,R}^{h,k}$ and use them as $C_{t,R}^{*,k}$.

After the hierarchical gated operations, let the memory output generated using Eq. 10 from memory $M_{t,I}$ at time instance $t$ be denoted by $h_{t,I}$ and memory $M_{t,R}$ be denoted as $h_{t,R}$. For simplicity, in Fig. 3 we consider 2 input streams, however the proposed coupling mechanism is flexible and is able to handle any number of modalities.

Motivated by [2, 24] we perform gated modality fusion such that,

$$
\begin{aligned}
\bar{h}_{t,I} &= \tanh(W_I h_{t,I}), \\
\bar{h}_{t,R} &= \tanh(W_R h_{t,R}), \\
\nu &= \sigma(W_\nu[\bar{h}_{t,I}, \bar{h}_{t,R}]),
\end{aligned}
\tag{17}
$$

where $W_I$ and $W_R$ are the weights for the respective memories and $W_\nu$ is the weight of the fusion gate. This can be seen as performing attention from one modality over the other where each modality determines the amount of information to flow from the other. We then obtain the combined feature vector,

$$
h_t = \nu \bar{h}_{t,I} + (1-\nu)\bar{h}_{t,R},
\tag{18}
$$

and augment Eq. 15 to utilise information from both streams,

$$
\bar{C}_t^{(k)} = \tanh([C_{t,I}^{*,k}, C_{t,R}^{*,k}, h_t]),
\tag{19}
$$

and predict the future trajectory using Eq. 16. In contrast to [5, 30] where simple concatenation of multimodal data is used, the proposed multi-memory architecture allows the model to store salient information of individual streams separately and propagate it effectively to the decision making process. We denote this model as $SMN(I+R)$ as it couples $I$ and $R$ streams to the $SMN$ model.

## 4    Evaluation and Discussion

### 4.1    Datasets

We present the experimental results for the single modal framework on the publicly available New York Grand Central (GC) [34] dataset. The Grand Central dataset consist of 12,600 trajectories. For training, testing and validation we use the same splits defined in [16]. Due to the unavailability of public multimodal
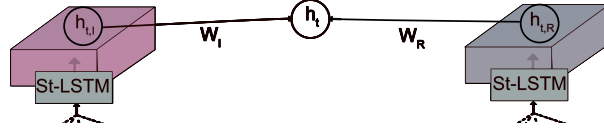
**Fig. 3.** Coupling multimodal information through multiple memory modules. The information from each modality is stored separately. Note that the figure shows only the top most layer in each memory.

pedestrian trajectory data, we introduce a new large scale dataset [1]. Pedestrian trajectories from a CCTV surveillance feed (I) and Radar (R) streams, for 32 hours, were collected and synchronised. Please refer to the supplementary material for statistics, calibration and synchronisation details of the dataset.

### 4.2  Evaluation Metrics

Following [1,16] we evaluate the performance with the following 3 error metrics: Average displacement error (ADE), Final displacement error (FDE) and Average non-linear displacement error (n-ADE). Please refer to [1,16] for details.

### 4.3  Evaluation of trajectory prediction with single modal data

The evaluation of single modal trajectories is conducted on the GC dataset [34]. We compare our model against 6 state of the art baselines. The first baseline is the Social Force (SF) model of [33]. It requires the destination of the pedestrian as input, and a linear SVM is trained with ground truth destination areas for this task. The next baseline is the Social LSTM (So-LSTM) model of [1]. It requires the neighbourhood size as a hyper-parameter and is set to 32px. The soft + hardwired attention model from [16] (SHA) does not posses any memory and computes the trajectory prediction by modelling the local neighbourhood of the pedestrian of interest. We also consider the Tree Memory Network (TMN) [14] which models the memory as a tree structure. This model uses the hyper parameter $\delta$, which defines the length of the memory as it structures a flat memory vector as a tree. We also evaluate the Neural Map (NM) model introduced in [28]. The pedestrian of interest's trajectory is embedded using a soft attention mechanism as defined in Sec. 3.1 and is stored in the memory. To provide a fair comparison, we also augment the NM module with the neighbourhood embeddings, $C_t^{s,k}$ and $C_t^{h,k}$, combine these with the memory output vector generated from the NM as in Eq. 15. We define this model as NMA.

To provide a direct comparison among baselines we set the hidden state dimensions of So-LSTM, SHA, TMN, NM, NMA and the proposed SMN model to be 30 units. As the models NM, NMA and SMN have a map width (W) and map height (H) as hyper-parameters, we evaluate different memory sizes. Similarly, for TMN we evaluate different memory lengths $\delta$. Please refer to the

---

[1] available at https://github.com/qutsaivt/SAIVTMultiSpectralTrajectoryDataset

supplementary material for those evaluations. Best results are shown in Tab. 1. To evaluate the relative performance of each model, we observe the trajectory for 20 frames and predict the future trajectory for the next 20 frames.

| Method | Metric | | |
|---|---|---|---|
| | ADE | FDE | n-ADE |
| SF [33] | 3.364 | 5.808 | 3.983 |
| So-LSTM [1] | 1.990 | 4.519 | 1.781 |
| SHA [16] | 1.096 | 3.011 | 0.985 |
| TMN ($\delta$=64) [14] | 2.982 | 4.989 | 2.780 |
| NM (W=H=64) [28] | 2.505 | 4.151 | 2.432 |
| NMA (W=H=64) | 1.466 | 3.811 | 1.445 |
| SMN (W=H=128) | **0.891** | **2.899** | **0.814** |

**Table 1.** Quantitative results with the GC dataset [34] for Social Force (SF) [33], Social LSTM (So-LSTM) [1], Soft + Hardwired Attention (SHA) [16], Tree Memory Network (TMN) [14], Neural Map (NM) [28], Neural Map Augmented (NMA) and the proposed Structured Memory Network (SMN) models. In all the methods forecast trajectories are of length 20 frames. The measured error metrics are as in Sec. 4.2.

From the results tabulated in Tab. 1 we observe poor performance in the SF model due to its lack of capacity to model long-term history. Models So-LSTM and SHA utilise short-term history from the pedestrian of interest and the local neighbourhood and generate improved predictions accordingly.

The lack of spatial structure and context modelling in the TMN module leads to it's poor performance despite it's long-term history modelling capacity. Comparing the NM and NMA models, the performance increase from NM and NMA is due to the addition of local context, highlighting the importance of capturing both long and short-term context. The NMA model attains improved performance due to the improved modelling of the local neighbourhood, and the structured memory; however when compared to the SHA model it fails to propagate salient spatiotemporal information from the structured memory to aid the decision making. This is due to the static kernel used when generating the memory output. In contrast, we map the memory output hierarchically using the proposed St-LSTM cells and propagate salient information to the upper layer, enabling efficient information transfer to the prediction model. The proposed gated architecture considers the evolution of memory over time, where multiple humans can interact with the environment, changing the state of multiple spatial locations. Hence we are able to generate dynamic responses instead of passing the information through a static convolution kernel as in NM and NMA; enabling superior performance even with large memory sizes.

We present a qualitative evaluation of the proposed SMN model with the SHA and NMA baselines in Fig. 4. We selected these baselines as they provide the highest comparative results. The trajectories are shown in the first column where the observed part of the trajectory is denoted in green, the ground truth

observations in blue, neighbouring trajectories are in purple and the predicted trajectories are shown in red (SMN), yellow (SHA) and orange (NMA).
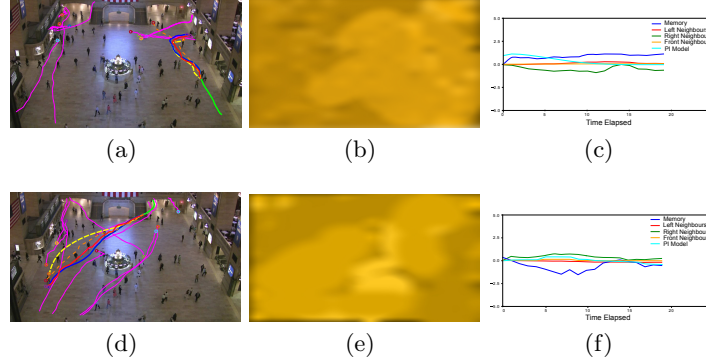


**Fig. 4.** Qualitative results for the GC dataset [34]: Given (in green), Ground Truth (in blue), Neighbouring (in purple) and Predicted trajectories from SMN model (in red), from SHA model (in yellow), from NMA model (in orange) along with the respective structured memory activations and relative activation contribution of each component in the prediction module. Please note that in the structured memory activations the intensity of the colour represents the degree of the activation and has been manually aligned with the figure in the first column for the clarity of visualisation.

When observing the qualitative results it can be clearly seen that the proposed SMN model generates better trajectory predictions compared to the state-of-the-arts. For instance in Fig. 4 (a) and (d) we observe significant deviation of the predictions of SHA and NMA models from the ground truth. However the proposed SMN model has been able to anticipate the pedestrian motion more accurately with the improved context modelling.

From the memory activation visualisations, it is evident more attention is given to cells surrounding the trajectory of the pedestrian of interest and the neighbours. Varying levels of attention are given to the cells occupied by the neighbours. However by passing this information through the proposed gated St-LSTM cells the proposed model is able to learn salient information among the passed activations from the layer below. This can be verified by observing the relative activation plots presented in the 3rd column of Fig. 4. While in general more attention is given to the encoded trajectory information from the pedestrian of interest (PI model), in cases such as Fig. 4 (c) more attention is given to the historic neighbourhood embeddings present in memory, where as in Fig. 4 (f) the model gives more attention to the neighbours. This verifies our hypothesis that both current context information encoded within the motion of pedestrian of interest and the neighbouring trajectories as well as the information

from the long-term history that preserves the structural integrity, is vital for prediction. Refer to the supplementary material for more qualitative evaluations.

### 4.4  Evaluation of trajectory prediction with multimodal data

The evaluation of multi-modal trajectories is conducted on the proposed multi-modal dataset. We compared our proposed model, SMN(I+R), with 4 state of the art baselines. In the first baseline, SHA(I+R), we concatenate the embeddings $C_{I,t}^{*,k}$ and $C_{R,t}^{*,k}$ for the $I$ and $R$ modalities directly to generate the augmented vector representation, $\bar{C}_t^{(k)}$, and use it in Eq. 16 to generate the prediction. The work of Fernando et al. [15] introduces a multi-modal extension to the TMN module. We use this model, TMN(I+R), as our next baseline. We extend the NM and NMA architectures (see Sec. 4.3) to handle multi-modal data. Similar to TMN(I+R) model we use multiple memories to store each input streams and pass the memory outputs through Eq. 17 to generate predictions. The augmented models are denoted NM(I+R) and NMA(I+R) in the evaluations. For models NM(I+R), NMA(I+R) and SMN(I+R) we set the map width (W) and map height (H) to be 128 and for TMN(I+R) we set the memory length $\delta$=64, as this provided the best accuracies in Sec 4.3.

Following the previous experiment, we observe the trajectory for 20 frames and predict the trajectory for the next 20 frames. After filtering out short and fragmented trajectories we are left with 40,800 trajectories. We randomly selected 28,560 trajectories for training, 10,200 for testing and 2,040 for validation.

| Method | Metric | | |
|---|---|---|---|
| | ADE | FDE | n-ADE |
| SHA(I+R) [16] | 1.245 | 1.654 | 1.454 |
| TMN(I+R) [15] | 2.901 | 3.169 | 3.001 |
| NM(I+R) [28] | 2.015 | 2.741 | 2.344 |
| NMA(I+R) | 1.325 | 1.814 | 1.558 |
| SMN(I+R) | **0.979** | **0.998** | **1.036** |

**Table 2.** Quantitative results with the proposed multimodal dataset for, Soft + Hard-wired Attention (SHA(I+R)) [16], Tree Memory Network (TMN(I+R)) [14], Neural Map (NM(I+R)) [28], Neural Map Augmented (NMA(I+R)) and the proposed Structured Memory Network (SMN(I+R)) models. In all the methods forecast trajectories are of length 20 frames. Error metrics are defined in Sec. 4.2.

Similar to the evaluations in Sec. 4.3, we observe poor performance from TMN(I+R) and NM(I+R) due to their inability to capture local neighbourhood information. However we observe a significant reduction in the performance gap between SHA(I+R) and NMA(I+R), compared to the that in Tab.1, which is a result of the naive fusion method used in the former model. SHA(I+R) simply concatenates the two modes together, and as such the model lacks the capacity

to capture salient information from individual modes. In contrast, by capturing long-term temporal dependencies of the two modalities, the memory based coupling mechanism yields better predictions. We further augment this process in SMN(I+R) by utilising the St-LSTM cells to hierarchically capture salient information from each mode. This enables the model to jointly back propagate through the two modalities and learn the strengths and weaknesses of each, effectively complimenting the prediction module with the additional information stream. Please refer to supplementary material for qualitative evaluations of the proposed SMN(I+R) model with the SHA(I+R) and NMA(I+R) baselines.

### 4.5   Ablation Experiments

To further demonstrate the effectiveness of our proposed fusion approach, we conduct a series of ablation experiments, identifying the crucial components of the proposed architecture. In the same settings as the experiment in Sec. 4.4, we compare the SMN(I+R) (proposed) method to a series of counterparts constructed by removing components of the model as follows:

- **SA(I)**: Uses only the soft attention context vector, $C_t^{s,k}$, and data from the image stream (I) for trajectory prediction.
- **SHA(I)**: Uses both soft ($C_t^{s,k}$) and hardwired ($C_t^{h,k}$) attention vectors and data from image stream (I) for trajectory prediction
- **SMN(I)**: Uses the proposed SMN model and data from Image (I) stream.
- **SA(R)**: Similar to SA-I but uses data from the Radar (R) stream.
- **SHA(R)**: Similar to SHA-I but uses data from the Radar (R) stream.
- **SMN(R)**: Similar to SMN-I but uses data from the Radar (R) stream.
- **SA(I+R)**: SA model that directly concatenates $C_{t,I}^{s,k}$ and $C_{t,R}^{s,k}$ and generates a vector embedding for Eq. 16.
- **SHA(I+R)**: SHA model that directly concatenates $C_{t,I}^{*,k}$ and $C_{t,R}^{*,k}$ and generates a vector embedding for Eq. 16.
- **SMN(I+R)**: Uses the model proposed in Sec. 3.4.

Note that for all $SMN$ models we used $W = H = 128$.

The results of our ablation study are presented in Tab. 3. Models SA(I) and SA(R) perform poorly due to their inability to oversee the neighbourhood context. We observe improved performance in SHA(I) and SHA(R) with the introduction of information from neighbouring pedestrians. The combined information from both modalities contributes to the performance gain we observe in SHA(I+R) over the unimodal counterparts, verifying the observations in [4,10].

Comparing the unimodal SMN(I) and SMN(R) models with the multimodal SHA(I+R) model, the former outperforms the latter by a significant margin, emphasising the importance of capturing long-term spatial context, and propagating the information effectively to the prediction model. The introduction of a secondary modality in SMN(I+R) further improves the prediction accuracy.

We would like to further compare the results obtained from the individual models in the $I$ and $R$ streams. We observe a performance boost in modality $I$,

| Method | Metric | | |
|---|---|---|---|
| | ADE | FDE | n-ADE |
| SA(I) | 2.012 | 3.011 | 2.190 |
| SHA(I) | 1.235 | 2.731 | 1.442 |
| SMN(I) | 1.029 | 1.104 | 1.092 |
| SA(R) | 2.259 | 3.312 | 2.261 |
| SHA(R) | 1.613 | 3.070 | 1.892 |
| SMN(R) | 1.198 | 1.330 | 1.288 |
| SA(I+R) | 1.334 | 1.813 | 1.579 |
| SHA(I+R) | 1.245 | 1.654 | 1.454 |
| SMN(I+R) | **0.979** | **0.998** | **1.036** |

**Table 3.** Ablation experiment evaluations

due to the finer granularity present in the CCTV stream due to the higher frame rate, compared to the radar stream. Hence extracted trajectories are smoother compared to the trajectories from modality $R$, making it easier to model.

### 4.6   Implementation Details

We use Keras [7] for our implementation. The SMN and SMN(I+R) modules do not require any special hardware (i.e. GPUs) to run. The SMN (W=H=128) model has 152K trainable parameters, and SMN(I+R) (W=H=128) has 358K. We ran the test set in Sec. 4.3 on a single core of an Intel Xeon E5-2680 2.50GHz CPU and the SMN algorithm was able to generate 1000 predicted trajectories with 40, 2 dimensional data points (i.e. using 20 observations to predict the next 20 data points) in 2.791 seconds. In a similar experiment with the test set in Sec. 4.4 we were able to generate 1000 predicted trajectories in 11.722 seconds.

## 5   Conclusions

In this paper we propose a method to anticipate complex human motion by analysing structural and temporal accordance. We extend the standard pedestrian trajectory prediction framework by introducing a novel model, Structured Memory Network (SMN), which is able to oversee the long-term history, preserving the structural integrity and improving prediction of pedestrian motion. As an extension to the proposed SMN model, we contribute a novel data driven method to capture salient information from multiple modalities and demonstrate how to incorporate this to enhance prediction. Additionally, we introduce a novel multi-modal pedestrian trajectory dataset, collected from synchronised CCTV and Radar streams, and consisting of 40,000 pedestrian trajectories. Our evaluations on both single and multi-modal datasets demonstrate the capacity of the proposed SMN method to learn complex real world human navigation behaviour.

**Acknowledgement**

## References

1. Alahi, A., Goel, K., Ramanathan, V., Robicquet, A., Fei-Fei, L., Savarese, S.: Social lstm: Human trajectory prediction in crowded spaces. In: CVPR. pp. 961–971 (2016)
2. Arevalo, J., Solorio, T., Montes-y Gómez, M., González, F.A.: Gated multimodal units for information fusion. ICLR (2017)
3. Bartoli, F., Lisanti, G., Ballan, L., Del Bimbo, A.: Context-aware trajectory prediction. arXiv preprint arXiv:1705.02503 (2017)
4. Bhatt, C.A., Kankanhalli, M.S.: Multimedia data mining: state of the art and challenges. Multimedia Tools and Applications $51$(1), 35–76 (2011)
5. Boström, M., Claesson, T.: Reducing false triggers in surveillance systems using sensor fusion. Master's Theses in Mathematical Sciences (2017)
6. Brun, V.H., Solstad, T., Kjelstrup, K.B., Fyhn, M., Witter, M.P., Moser, E.I., Moser, M.B.: Progressive increase in grid scale from dorsal to ventral medial entorhinal cortex. Hippocampus $18$(12), 1200–1212 (2008)
7. Chollet, F.: Keras. URL http://keras. io, 2017 (2017)
8. Chorowski, J.K., Bahdanau, D., Serdyuk, D., Cho, K., Bengio, Y.: Attention-based models for speech recognition. In: NIPS. pp. 577–585 (2015)
9. Coates, A., Ng, A.Y.: The importance of encoding versus training with sparse coding and vector quantization. In: ICML. pp. 921–928 (2011)
10. Deng, L., Yu, D., et al.: Deep learning: methods and applications. Foundations and Trends® in Signal Processing $7$(3–4), 197–387 (2014)
11. Derdikman, D., Moser, E.I.: A manifold of spatial maps in the brain. Trends in cognitive sciences $14$(12), 561–569 (2010)
12. Epstein, R.A., Patai, E.Z., Julian, J.B., Spiers, H.J.: The cognitive map in humans: spatial navigation and beyond. Nature neuroscience $20$(11),  1504 (2017)
13. Fanselow, M.S., Dong, H.W.: Are the dorsal and ventral hippocampus functionally distinct structures? Neuron $65$(1), 7–19 (2010)
14. Fernando, T., Denman, S., McFadyen, A., Sridharan, S., Fookes, C.: Tree memory networks for modelling long-term temporal dependencies. Neurocomputing $304$, 64–81 (2018)
15. Fernando, T., Denman, S., Sridharan, S., Fookes, C.: Going deeper: Autonomous steering with neural memory networks. In: ICCV. pp. 214–221 (2017)
16. 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)
17. Fernando, T., Denman, S., Sridharan, S., Fookes, C.: Learning temporal strategic relationships using generative adversarial imitation learning. IFAAMAS (2018)
18. Fernando, T., Denman, S., Sridharan, S., Fookes, C.: Task specific visual saliency prediction with memory augmented conditional generative adversarial networks. In: WACV. pp. 1539–1548. IEEE (2018)

19. Fernando, T., Denman, S., Sridharan, S., Fookes, C.: Tracking by prediction: A deep generative model for multi-person localisation and tracking. WACV (2018)
20. Gobet, F., Lane, P.C., Croker, S., Cheng, P.C., Jones, G., Oliver, I., Pine, J.M.: Chunking mechanisms in human learning. Trends in cognitive sciences **5**(6), 236–243 (2001)
21. Huang, Y., Wu, Q., Wang, L.: Learning semantic concepts and order for image and sentence matching. arXiv preprint arXiv:1712.02036 (2017)
22. Kaiser, L., Sutskever, I.: Neural gpus learn algorithms. ICLR (2016)
23. Kiela, D., Bottou, L.: Learning image embeddings using convolutional neural networks for improved multi-modal semantics. In: EMNLP. pp. 36–45 (2014)
24. Kiela, D., Grave, E., Joulin, A., Mikolov, T.: Efficient large-scale multi-modal classification. arXiv preprint arXiv:1802.02892 (2018)
25. Kiros, R., Salakhutdinov, R., Zemel, R.: Multimodal neural language models. In: ICML. pp. 595–603 (2014)
26. Madl, T., Franklin, S., Chen, K., Trappl, R., Montaldi, D.: Exploring the structure of spatial representations. PloS one **11**(6), e0157343 (2016)
27. Malinowski, M., Fritz, M.: A multi-world approach to question answering about real-world scenes based on uncertain input. In: NIPS. pp. 1682–1690 (2014)
28. Parisotto, E., Salakhutdinov, R.: Neural map: Structured memory for deep reinforcement learning. In: ICLR (2018)
29. Pei, D., Liu, H., Liu, Y., Sun, F.: Unsupervised multimodal feature learning for semantic image segmentation. In: IJCNN. pp. 1–6. IEEE (2013)
30. Roy, A., Gale, N., Hong, L.: Automated traffic surveillance using fusion of doppler radar and video information. Mathematical and Computer Modelling **54**(1-2), 531–543 (2011)
31. Srivastava, N., Salakhutdinov, R.R.: Multimodal learning with deep boltzmann machines. In: NIPS. pp. 2222–2230 (2012)
32. Varshneya, D., Srinivasaraghavan, G.: Human trajectory prediction using spatially aware deep attention models. arXiv preprint arXiv:1705.09436 (2017)
33. Yamaguchi, K., Berg, A.C., Ortiz, L.E., Berg, T.L.: Who are you with and where are you going? In: CVPR. pp. 1345–1352. IEEE (2011)
34. Yi, S., Li, H., Wang, X.: Understanding pedestrian behaviors from stationary crowd groups. In: CVPR. pp. 3488–3496 (2015)
35. Yuan, A., Li, X., Lu, X.: Ffgs: Feature fusion with gating structure for image caption generation. In: CCF Chinese Conference on Computer Vision. pp. 638–649. Springer (2017)
36. Zou, H., Su, H., Song, S., Zhu, J.: Understanding human behaviors in crowds by imitating the decision-making process. arXiv preprint arXiv:1801.08391 (2018)