# Take a NAP: Non-Autoregressive Prediction for Pedestrian Trajectories

**Hao Xue** , **Du. Q. Huynh** , **Mark Reynolds**
**The University of Western Australia**
hao.xue@research.uwa.edu.au, du.huynh@uwa.edu.au, mark.reynolds@uwa.edu.au

## Abstract

Pedestrian trajectory prediction is a challenging task as there are three properties of human movement behaviors which need to be addressed, namely, the social influence from other pedestrians, the scene constraints, and the multimodal (multi-route) nature of predictions. Although existing methods have explored these key properties, the prediction process of these methods is autoregressive. This means they can only predict future locations sequentially. In this paper, we present NAP, a non-autoregressive method for trajectory prediction. Our method comprises specifically designed feature encoders and a latent variable generator to handle the three properties above. It also has a time-agnostic context generator and a time-specific context generator for non-autoregressive prediction. Through extensive experiments that compare NAP against several recent methods, we show that NAP has state-of-the-art trajectory prediction performance.

## 1 Introduction

Pedestrian trajectory prediction is an important component in a range of applications such as social robots and self-driving vehicles, and plays a key role in understanding human movement behaviors. This task is not easy due to three key properties in pedestrian trajectory prediction: (i) **social interaction**: People are not always walking alone in a public space. Pedestrians often socially interact with others to avoid collisions, walk with friends or keep a certain distance from strangers; (ii) **environmental scene constraints**: Besides social interaction, pedestrians' routes also need to obey scene constraints such as obstacles and building layouts; and (iii) **multimodal nature of future prediction**: People can follow different routes as long as these routes are both socially and environmentally acceptable. For example, a person can choose to turn right or turn left to bypass an obstacle.

Recently, researchers have made progress in incorporating these properties into the trajectory prediction process. For example, the Social LSTM model [Alahi *et al.*, 2016] is one of the methods that can capture social influence information in a local neighborhood around each pedestrian. Based on the generative model GAN [Goodfellow *et al.*, 2014], the SGAN model proposed by Gupta *et al.* [Gupta *et al.*, 2018] can handle multimodality in the prediction process while also capturing the social influence from other pedestrians in the scene. To deal with the scene constraints, Convolutional Neural Networks (CNNs) are often used to extract scene information in the trajectory prediction network, such as SS-LSTM [Xue *et al.*, 2018], SoPhie [Sadeghian *et al.*, 2019], and Social-BiGAT [Kosaraju *et al.*, 2019] .

While other methods like [Su *et al.*, 2017; Vemula *et al.*, 2018; Hasan *et al.*, 2018; Zou *et al.*, 2018; Xue *et al.*, 2019; Li, 2019; Zhang *et al.*, 2019] miss one or two aforementioned key properties, SoPhie [Sadeghian *et al.*, 2019], Liang *et al.* [Liang *et al.*, 2019], and Social-BiGAT [Kosaraju *et al.*, 2019] are three typical papers that have taken all three properties into consideration. However, these methods predict the future locations recurrently. There are two main limitations in using autoregression to generate trajectory prediction: (i) the autoregressive prediction process works in a recursive manner and so the prediction errors accumulated from previous time steps are passed to the prediction for the next time step; (ii) the process cannot be parallelized, *i.e.*, predictions must be generated sequentially, even though one might be interested in generating only the final destination of the pedestrian rather than the entire trajectory.

To overcome the above limitations and inspired by the application of non-autoregressive decoder in other areas such as machine translation [Gu *et al.*, 2018; Guo *et al.*, 2019] and time series forecasting [Wen *et al.*, 2017], we propose a novel trajectory prediction method that can predict future trajectories non-autoregressively. We name our method *NAP* (short for *N*on-*A*utoregressive *P*rediction). Our research contributions are threefold: **(i)** To the best of our knowledge, we are the first to explore non-autoregressive trajectory prediction. The network architecture of NAP includes trainable context generators to ensure that context vectors are available for the non-autoregressive decoder to forecast good quality predictions. The state-of-the-art performance of NAP is demonstrated through the extensive experiments and ablation study conducted. **(ii)** Both the social and scene influences are handled by NAP through specially designed feature encoders. The social influence is captured by social graph features propagated through an LSTM; the scene influence is modeled by a CNN. The effectiveness of these encoders is confirmed from

the performance of NAP. **(iii)** Unlike previous work in the literature, NAP tackles multimodal predictions by training a latent variable generator to learn the latent variables of the sampling distribution for each pedestrian's trajectory. The generator is shown to give NAP top performance in multimodal predictions.

## 2 Background

### 2.1 Problem Definition

Pedestrian trajectory prediction is defined as forecasting the future trajectory of the person $i$ given his/her observed trajectory. We assume that trajectories have already been obtained in the format of time sequences of coordinates (*i.e.*, $\mathbf{u}_t^i = (x_t^i, y_t^i) \in \mathbb{R}^2, \forall i$). The length of the observed trajectory and predicted trajectory are represented by $T_\mathrm{o}$ and $T_\mathrm{p}$. Thus, considering the observed trajectory $X^i = \{\mathbf{u}_t^i \mid t = 1, \cdots, T_\mathrm{o}\}$, our target is to generate the prediction $\hat{Y}^i = \{\hat{\mathbf{u}}_t^i \mid t = T_\mathrm{o} + 1, \cdots, T_\mathrm{o} + T_\mathrm{p}\}$.

### 2.2 Autoregressive and Non-Autoregressive Predictions

Mathematically, to generate the predicted trajectory $\hat{Y}^i$ from a given observed trajectory $X^i$, the conditional probability $P_{\mathrm{AR}}(\hat{Y}^i|X^i; \boldsymbol{\theta})$ with parameter $\boldsymbol{\theta}$ for an autoregressive predictor is defined as:

$$P_{\mathrm{AR}}(\hat{Y}^i|X^i; \boldsymbol{\theta}) = \prod_{t=T_\mathrm{o}+1}^{T_\mathrm{o}+T_\mathrm{p}} P(\hat{\mathbf{u}}_t^i|\hat{\mathbf{u}}_{T_\mathrm{o}:t-1}^i, X^i; \boldsymbol{\theta}), \quad (1)$$

where generating the prediction of time step $t$ requires the prediction of all previous time steps in the prediction phase. This recursive prediction process can not be parallelized.

Different from the autoregressive prediction process, by treating $\hat{\mathbf{u}}_t^i$ and $\hat{\mathbf{u}}_{t'}^i$, for all $t \neq t' \geq T_\mathrm{o}$, as independent, the above conditional probability in non-autoregressive predictors becomes:

$$P_{\mathrm{NAR}}(\hat{Y}^i|X^i; \boldsymbol{\theta}) = \prod_{t=T_\mathrm{o}+1}^{T_\mathrm{o}+T_\mathrm{p}} P(\hat{\mathbf{u}}_t^i|X^i; \boldsymbol{\theta}). \quad (2)$$

Compared to the autoregression based prediction where the prediction at time step $t$ depends on the information at time step $t - 1$, non-autoregression based predictors do not need to generate predictions sequentially. However, the removal of this sequential dependency in non-autoregression based methods means that the time awareness ability is compromised in the prediction model, leading to poorer prediction performance. To compensate for the loss of time awareness ability, we design context generators (detailed in Section 3.2) that are trained on the training trajectories. This allows context vectors to be generated from the observed trajectories in the testing stage for the prediction process. As these context vectors can be computed before the start of the prediction phase, predictions at different time steps can be forecast in parallel.

## 3 Proposed Method

Our proposed NAP comprises four major parts (Fig. 1): (i) feature encoders which are used to encode the input information such as observed trajectories and scene images (Section 3.1); (ii) context generators to yield context vectors for prediction (Section 3.2); (iii) a latent variable generator for multimodality (Section 3.3); and (iv) a non-autoregressive decoder for predicting future trajectories (Section 3.4). Details of these parts are described in the following subsections.

### 3.1 Feature Encoders

In NAP, there are three feature encoders: a trajectory encoder, to learn the representation of the observed history movement of each pedestrian; a social encoder, to learn the representation of the influence from other pedestrians; and a scene encoder, to learn the representation of the scene features.

**Trajectory Encoder.** The coordinates of the $i^{\mathrm{th}}$ pedestrian in the observation phase ($t = 1, \cdots, T_\mathrm{o}$) are firstly embedded into a higher dimensional space through an embedding layer $\phi(\cdot)$. Then, across different time steps, the embedded features are used as inputs of an LSTM layer (denoted by $\mathrm{LSTM}_{\mathrm{ENC}}(\cdot)$) to get the encoded hidden state $\mathbf{h}_t^i$ which captures the observed path information. This trajectory encoding is given by:

$$\mathbf{e}_t^i = \phi(x_t^i, y_t^i; \mathbf{W}_{\mathrm{EMB}}), \quad (3)$$

$$\mathbf{h}_t^i = \mathrm{LSTM}_{\mathrm{ENC}}(\mathbf{h}_{t-1}^i, \mathbf{e}_t^i; \mathbf{W}_{\mathrm{ENC}}), \quad (4)$$

where $\mathbf{W}_{\mathrm{EMB}}$ and $\mathbf{W}_{\mathrm{ENC}}$ are trainable weights of the corresponding layers.

**Social Encoder.** At each time step $t$, NAP captures the social influence on the $i^{\mathrm{th}}$ pedestrian through a graph $\mathcal{G}_t^i = (V_t^i, E_t^i)$. The $i^{\mathrm{th}}$ pedestrian and all other pedestrians $\mathcal{N}_t^{(i)}$ at the same time step are considered as nodes in the set $V_t^i$. Edges linking the $i^{\mathrm{th}}$ pedestrian and pedestrians in $\mathcal{N}_t^{(i)}$ form the edge set $E_t^i$.

We then use a graph convolutional network (GCN) to process these graphs. In the $l^{\mathrm{th}}$ graph convolutional layer, the node feature of pedestrian $i$ is aggregated as follows:

$$\mathbf{a}_t^{i,l} = \mathrm{ReLU}\left(\mathbf{b}^l + \frac{1}{|\mathcal{N}_t^{(i)}|} \sum_{j \in \mathcal{N}_t^{(i)}} \mathbf{W}^l \mathbf{a}_t^{j,l-1}\right), \quad (5)$$

where $\mathbf{W}^l$ and $\mathbf{b}^l$ are the weight matrix and bias term. At the first layer, we initialize the node feature $\mathbf{a}_t^{i,0}$ as the location coordinates of the $i^{\mathrm{th}}$ pedestrian, *i.e.*, $\mathbf{a}_t^{i,0} = (x_i^t, y_i^t)$.

The social graph feature $\mathbf{g}_t^i$ (Eq. (6)) is designed to model the surrounding (or *social*) information of pedestrian $i$ at each time step $t$. To compute this feature, the node features, denoted by $\{\mathbf{a}_t^i \mid t = 1, \cdots, T_\mathrm{o}\}$, from the final GCN layer across all the time steps in the observation phase are passed through an LSTM layer with trainable weights $\mathbf{W}_{\mathrm{SG}}$, *i.e.*,

$$\mathbf{g}_t^i = \mathrm{LSTM}_{\mathrm{SG}}(\mathbf{g}_{t-1}^i, \mathbf{a}_t^i; \mathbf{W}_{\mathrm{SG}}). \quad (6)$$
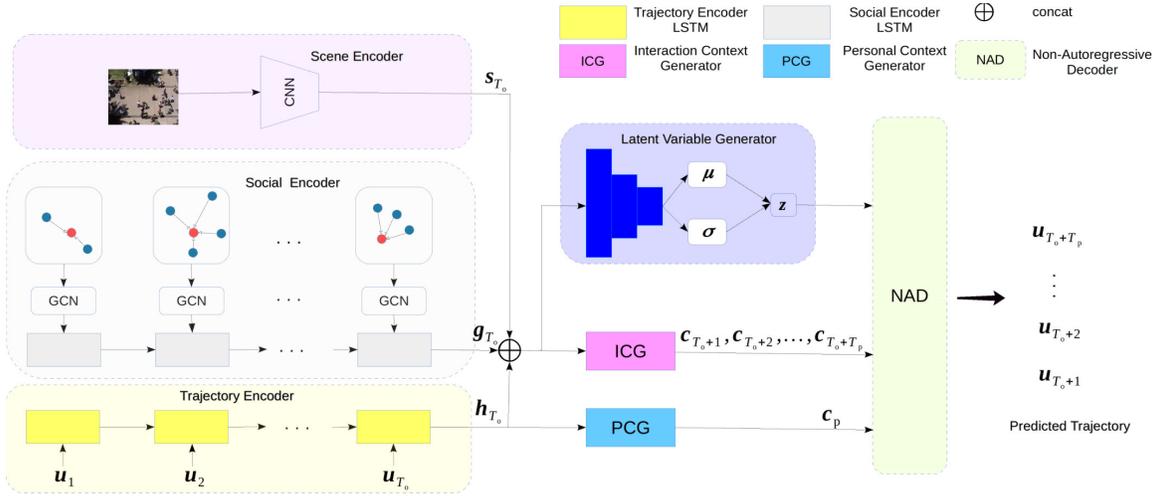
Figure 1: The network architecture of NAP. There are three encoders to extract features, two context generators to generate context vectors for the non-autoregressive decoder (NAD), and a latent variable generator to handle multimodal predictions. The embedding layer and the superscript $i$ (representing the pedestrian index) are not shown to simplify visualization.

**Scene Encoder.** Different from other methods (such as [Xue *et al.*, 2018; Sadeghian *et al.*, 2019]) that process each image frame in the observation phase, the scene encoder of NAP takes only image $I_{T_o}^i$ as input, since the scene encoder focuses on the static information like scene layouts and obstacles. Not only does this save the computation time, but it also supplies the most up-to-date and sufficient scene context before prediction kicks in at $t = T_o + 1$. We use a CNN to model the scene feature $\mathbf{s}_{T_o}^i$ as follows:

$$\mathbf{s}_{T_o}^i = \text{CNN}(I_{T_o}^i; \mathbf{W}_{\text{CNN}}). \tag{7}$$

We take the merit of DeepLabv3+ [Chen *et al.*, 2018], a state-of-the-art semantic segmentation architecture, and set the initial value of $\mathbf{W}_{\text{CNN}}$ to the weight matrix from DeepLabv3+ that is pre-trained on the Cityscapes dataset [Cordts *et al.*, 2016].

## 3.2 Context Generators

The role of the context generators is to aggregate the outputs of the feature encoders for the downstream decoder for trajectory forecasting. We use two context generators in NAP: (i) a *personal context generator* that is *time-agnostic*, as its input is the hidden state $\mathbf{h}_t^i$ computed from the $i^{\text{th}}$ pedestrian's own trajectory only; (ii) an *interaction context generator* that is *time-specific* as its input includes both social graph and scene interaction features also.

**Personal Context Generator (PCG).** We use a Multi-Layer Perceptron (MLP) to model this context generator. The output context vector $\mathbf{c}_p^i$ is computed as

$$\mathbf{c}_p^i = \text{MLP}_A(\mathbf{h}_{T_o}^i; \mathbf{W}_A), \tag{8}$$

where $\mathbf{W}_A$ is the corresponding weight matrix. The context $\mathbf{c}_p^i$ captures the "personal" cues such as the pedestrian's preferable walking speed and direction in the observation phase, oblivious of his/her surrounding. This context is *time-agnostic* because, without considering the social and scene

influences, such personal cues can remain the same for the entire trajectory, *e.g.*, the pedestrian can continue to walk in a straight line or at a fast pace with no penalty when bumping into obstacles or other pedestrians since both the social graph and scene features are not present in the equation.

**Interaction Context Generator (ICG).** This context generator incorporates both the social graph and scene features. These two types of influences allow the context generator to be time-specific, *e.g.*, while a pedestrian can walk at a fast pace in the initial part of his/her trajectory, he/she would need to slow down or detour at a later part of the trajectory in order to avoid other pedestrians or scene obstacles. Similar to PCG, we use an MLP to model ICG but its input, being the concatenation of $\mathbf{h}_{T_o}^i$, $\mathbf{g}_{T_o}^i$, and $\mathbf{s}_{T_o}^i$, contains richer information. The output of ICG comprises different context vectors for different time steps in the prediction phase, as given below:

$$(\mathbf{c}_{T_o+1}^i, \mathbf{c}_{T_o+2}^i, \cdots, \mathbf{c}_{T_o+T_p}^i) = \text{MLP}_B(\mathbf{h}_{T_o}^i \oplus \mathbf{g}_{T_o}^i \oplus \mathbf{s}_{T_o}^i; \mathbf{W}_B), \tag{9}$$

where $\mathbf{W}_B$ is the corresponding weight matrix.

## 3.3 Latent Variable Generator

For multimodal prediction, NAP is designed to generate multiple trajectories through the latent variables $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ (see Fig. 1). Although several existing trajectory prediction methods such as [Lee *et al.*, 2017; Gupta *et al.*, 2018; Li, 2019; Huang *et al.*, 2019] also use latent variables to handle multimodality, the latent variables in these methods are either directly sampled from the normal distribution or a multivariate normal distribution conditioned on the observed trajectories. To make our latent variables more aware of the social and scene cues, we design NAP to learn the parameters ($\boldsymbol{\mu}_i$ and $\boldsymbol{\sigma}_i$) of the sampling distribution from the observed trajectories, the social influence, and the scene influence features. To this end, the concatenated feature $\mathbf{h}_{T_o}^i \oplus \mathbf{g}_{T_o}^i \oplus \mathbf{s}_{T_o}^i$ is passed to two different MLPs (Eqs. (10)-(11)) to yield the mean vector
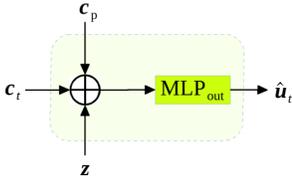
Figure 2: Details of the non-autoregressive decoder (NAD) of NAP.

$\boldsymbol{\mu}_i$ and variance $\boldsymbol{\sigma}_i$ and finally $\mathbf{z}_i$ for the downstream non-autoregressive decoder:

$$\boldsymbol{\mu}_i = \text{MLP}_\mu(\mathbf{h}_{T_o}^i \oplus \mathbf{g}_{T_o}^i \oplus \mathbf{s}_{T_o}^i; \mathbf{W}_\mu), \qquad (10)$$

$$\boldsymbol{\sigma}_i = \text{MLP}_\sigma(\mathbf{h}_{T_o}^i \oplus \mathbf{g}_{T_o}^i \oplus \mathbf{s}_{T_o}^i; \mathbf{W}_\sigma), \qquad (11)$$

$$\mathbf{z}_i \sim \mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i^2)), \qquad (12)$$

where $\mathbf{W}_\mu$ and $\mathbf{W}_\sigma$ are trainable weights of $\text{MLP}_\mu(\cdot)$ and $\text{MLP}_\sigma(\cdot)$. The reparameterization trick [Kingma and Welling, 2013] is applied to sample the latent variable $\mathbf{z}_i$.

### 3.4 Non-Autoregressive Decoder (NAD)

In the work of [Gu *et al.*, 2018], the authors introduce in their model a component that enhances the inputs passed to the decoder for their machine translation problem. The idea is to help the model learn the internal dependencies (which are absent in their non-autoregressive translator) within a sentence. In the work of [Guo *et al.*, 2019], the authors use a positional module to improve the decoder's ability to perform local re-ordering. The context generators in NAP play a similar role as these two approaches. In the testing stage, the trained PCG and ICG are able to generate context vectors for new (unseen) observed trajectories to help the decoder improve its time awareness for trajectory prediction. While the ICG, which generates time-specific contexts $\{\mathbf{c}_t \,|\, T_o + 1 \leqslant t \leqslant T_o + T_p\}$, is obviously more important than the PCG, the time-agnostic PCG is also needed so as to help the NAD tie with the specific trajectory being considered.

To make multimodal predictions at time step $t$, the NAD therefore takes the concatenation of the two contexts $\mathbf{c}_p^i$ and $\mathbf{c}_t^i$ and the latent variable $\mathbf{z}_i$ as inputs, modeled by the MLP below (see Fig. 2):

$$(\hat{x}_t^i, \hat{y}_t^i) = \text{MLP}_{\text{out}}(\mathbf{c}_t^i \oplus \mathbf{c}_p^i \oplus \mathbf{z}_i; \mathbf{W}_{\text{out}}), \qquad (13)$$

where $\text{MLP}_{\text{out}}(\cdot)$ is the MLP used for predicting the location coordinates. Its parameter $\mathbf{W}_{\text{out}}$ is shared across all the time steps in the prediction phase. Note that the input passed to the NAD in Fig. 2 is different for each time step $t$ in the prediction phase as $\mathbf{c}_t^i$ depends on $t$. We can consider that the contexts $\{\mathbf{c}_t^i\}$ function like the hidden states in the decoder of an LSTM except that they are not recursively defined.

### 3.5 Implementation Details

The embedding layer $\phi$ in Eq. (3) is modeled as a single-layer perceptron that outputs 32-dimensional embedded vectors for the input location coordinates. The dimensions of the hidden states of the LSTM layers for both the Trajectory and Social Encoders are 32. The GCN in the Social Encoder is a single graph convolutional layer (*i.e.*, $l = 1$ in Eq. (5)). For the ICG,

$\text{MLP}_B$ is a three-layer MLP with ReLU activations. All the other MLPs used in Eqs. (8), (10), (11), and (13) are single-layer MLPs. Except for in Section 4.4 where we explore the prediction performance for different prediction lengths, the observed length of input trajectories is 8 time steps ($T_o = 8$) and the prediction length is 12 time steps ($T_p = 12$) for all other experiments.

We implemented NAP and its variants (Section 4.3) using the PyTorch framework in Python. The Adam optimizer was used to train our models with the learning rate set to 0.001 and the mini-batch size to 128.

## 4 Experiments

### 4.1 Datasets and Metrics

We use the popular ETH [Pellegrini *et al.*, 2009] and UCY [Lerner *et al.*, 2007] datasets, which, altogether, include 5 scenes: ETH, HOTEL, UNIV, ZARA1, and ZARA2. Similar to [Sadeghian *et al.*, 2019; Zhang *et al.*, 2019], we normalize each pedestrian's coordinates and augment the training data by rotating trajectories. As raw scene images are used as inputs for extracting scene features, we also rotate the input images when input trajectories are rotated. Same as previous work in the literature [Alahi *et al.*, 2016; Gupta *et al.*, 2018; Sadeghian *et al.*, 2019; Huang *et al.*, 2019], the leave-one-out strategy is adopted for training and testing. All methods are evaluated based on two standard metrics: the Average Displacement Error (ADE) and the Final Displacement Error (FDE). Smaller errors indicate better prediction performance.

### 4.2 Comparison with Other Methods

We compare NAP against the following state-of-the-art trajectory prediction methods: Social-LSTM [Alahi *et al.*, 2016], SGAN [Gupta *et al.*, 2018], MX-LSTM [Hasan *et al.*, 2018], Nikhil & Morris [Nikhil and Morris, 2018], Liang *et al.* [Liang *et al.*, 2019], MATF [Zhao *et al.*, 2019], SR-LSTM [Zhang *et al.*, 2019], SoPhie [Sadeghian *et al.*, 2019], IDL [Li, 2019], STGAT [Huang *et al.*, 2019], and Social-BiGAT [Kosaraju *et al.*, 2019].

In Table 1, all the compared methods are put into two groups depending on whether they generate only one prediction (top half of the table) or multiple predictions (bottom half and indicated by a tick under the # column) for each input observed trajectory. The multimodal predictions being considered is 20. The reported ADEs and FDEs are computed from the best predictions out of the 20. Five methods report both single and multimodal prediction results and so they appear in both halves of the table: SGAN, MATF, STGAT, Liang *et al.*, and NAP.

Our proposed method is able to achieve results on par with the state-of-the-art methods for the single prediction setting. Specifically, NAP has the same smallest average ADE (0.45m) as SR-LSTM while outperforming all methods on the average FDE (0.89m). In addition to NAP, SR-LSTM, MX-LSTM, and STGAT 1V-1 also have the best performance on one or more scenes. In the lower half of Table 1, results of multimodal predictions are given and compared. On average, our NAP achieves the smallest ADE of 0.39m and the smallest FDE of 0.80m. For each scene, the best performers that

| Method | # | ETH & UCY scenes | | | | | |
|---|---|---|---|---|---|---|---|
| | | ETH | HOTEL | UNIV | ZARA1 | ZARA2 | Average |
| Social-LSTM [Alahi *et al.*, 2016][‡] | | 1.09 / 2.35 | 0.79 / 1.76 | 0.67 / 1.40 | 0.47 / 1.00 | 0.56 / 1.17 | 0.72 / 1.54 |
| SGAN 1V-1 [Gupta *et al.*, 2018][⋆] | | 1.13 / 2.21 | 1.01 / 2.18 | 0.60 / 1.28 | 0.42 / 0.91 | 0.52 / 1.11 | 0.74 / 1.54 |
| MX-LSTM [Hasan *et al.*, 2018][⋆] | | – | – | **0.49 / 1.12** | 0.59 / 1.31 | 0.35 / 0.79 | – |
| Nikhil & Morris [Nikhil and Morris, 2018][⋆] | | 1.04 / 2.07 | 0.59 / 1.17 | 0.57 / 1.21 | 0.43 / 0.90 | 0.34 / 0.75 | 0.59 / 1.22 |
| Liang *et al.* [Liang *et al.*, 2019][⋆] | | 0.88 / 1.98 | 0.36 / 0.74 | 0.62 / 1.32 | 0.42 / 0.90 | 0.34 / 0.75 | 0.52 / 1.14 |
| MATF [Zhao *et al.*, 2019][⋆] | | 1.33 / 2.49 | 0.51 / 0.95 | 0.56 / 1.19 | 0.44 / 0.93 | 0.34 / 0.73 | 0.64 / 1.26 |
| SR-LSTM [Zhang *et al.*, 2019][⋆] | | 0.63 / 1.25 | 0.37 / 0.74 | 0.51 / 1.10 | **0.41 / 0.90** | 0.32 / 0.70 | **0.45 / 0.94** |
| STGAT 1V-1 [Huang *et al.*, 2019][⋆] | | 0.88 / 1.66 | 0.56 / 1.15 | 0.52 / 1.13 | **0.41 / 0.91** | **0.31 / 0.68** | 0.54 / 1.11 |
| NAP (ours) | | **0.59 / 1.13** | **0.30 / 0.51** | 0.59 / 1.23 | **0.41 / 0.86** | 0.36 / 0.72 | **0.45 / 0.89** |
| SGAN 20V-20 [Gupta *et al.*, 2018][⋆] | ✔ | 0.81 / 1.52 | 0.72 / 1.61 | 0.60 / 1.26 | 0.34 / 0.69 | 0.42 / 0.84 | 0.58 / 1.18 |
| SoPhie [Sadeghian *et al.*, 2019][⋆] | ✔ | 0.70 / 1.43 | 0.76 / 1.67 | 0.54 / 1.24 | 0.30 / 0.63 | 0.38 / 0.78 | 0.54 / 1.15 |
| Liang *et al.* [Liang *et al.*, 2019][⋆] | ✔ | 0.73 / 1.65 | 0.30 / 0.59 | 0.60 / 1.27 | 0.38 / 0.81 | 0.31 / 0.68 | 0.46 / 1.00 |
| MATF GAN [Zhao *et al.*, 2019][⋆] | ✔ | 1.01 / 1.75 | 0.43 / 0.80 | **0.44 / 0.91** | 0.26 / 0.45 | 0.26 / 0.57 | 0.48 / 0.90 |
| IDL [Li, 2019][⋆] | ✔ | 0.59 / 1.30 | 0.46 / 0.83 | 0.51 / 1.27 | **0.22 / 0.49** | **0.23 / 0.55** | 0.40 / 0.89 |
| STGAT 20V-20 [Huang *et al.*, 2019][⋆] | ✔ | 0.65 / 1.12 | 0.35 / 0.66 | 0.52 / 1.10 | 0.34 / 0.69 | 0.29 / 0.60 | 0.43 / 0.83 |
| Social-BiGAT [Kosaraju *et al.*, 2019][⋆] | ✔ | 0.69 / 1.29 | 0.49 / 1.01 | 0.55 / 1.32 | 0.30 / 0.62 | 0.36 / 0.75 | 0.48 / 1.00 |
| NAP (ours) | ✔ | **0.53 / 1.08** | **0.26 / 0.46** | 0.58 / 1.22 | 0.30 / 0.65 | 0.28 / 0.60 | **0.39 / 0.80** |

Table 1: The ADEs / FDEs (in meters) of various methods. The settings are: $T_o = 8$ and $T_p = 12$. The results with a ⋆ are taken from the authors' papers. The result with ‡ is taken from [Gupta *et al.*, 2018].

achieve the smallest ADE/FDE in the lower half of the table include NAP, IDL, MATF, and GAN. Taken together, these results demonstrate the efficacy of our proposed method in both single and multimodal prediction settings.

### 4.3 Ablation Study

To explore the effectiveness of different contexts working together in our proposed method, we consider four variants of NAP listed below:

- **NAP-P**: This variant only uses the Personal Context Generator (time-agnostic context, the light blue PCG box in Fig. 1), *i.e.*, the interaction context $\mathbf{c}_t^i$ in Eq. (13) is removed.

- **NAP-ISS**: In contrast to NAP-P, NAP-ISS disables the personal context and forecasts predictions based on the time-specific interaction context (the pink box in Fig. 1). The personal context $\mathbf{c}_p^i$ in Eq. (13) is removed. The rest of NAP-ISS is the same as NAP.

- **NAP-ISg**: In order to further investigate the impact of removing the scene influence, we drop the scene feature $\mathbf{s}_{T_o}^i$ from the Interaction Context Generator in Eq. (9) to form this variant. That is, the interaction context $\mathbf{c}_t^i$ in NAP-ISS is computed using both the social graph and scene features, whereas the $\mathbf{c}_t^i$ in NAP-ISg is computed using the social graph feature only.

- **NAP-ISc**: Similar to NAP-ISg, this variant is designed to investigate the impact of removing the social influence. We drop the social graph feature $\mathbf{g}_{T_o}^i$ but keep the scene feature $\mathbf{s}_{T_o}^i$ in Eq. (9) so the context $\mathbf{c}_t^i$ is computed from the scene feature only.

In our ablation study, we compare only the single prediction performance (see Table 2) of these four variants, *i.e.*, the latent variable $\mathbf{z}_i$ for multimodality is removed from Eq. (13) in the experiments. In general, NAP-P, which uses only the personal context (time-agnostic), has a poorer performance

| | NAP-P | NAP-ISS | NAP-ISg | NAP-ISc |
|---|---|---|---|---|
| ETH | 0.87 / 1.63 | 0.66 / 1.22 | 0.69 / 1.31 | 0.74 / 1.52 |
| HOTEL | 0.43 / 0.77 | 0.34 / 0.61 | 0.37 / 0.70 | 0.38 / 0.73 |
| UNIV | 0.71 / 1.42 | 0.68 / 1.37 | 0.68 / 1.35 | 0.70 / 1.39 |
| ZARA1 | 0.46 / 0.95 | 0.45 / 0.94 | 0.47 / 0.96 | 0.45 / 0.95 |
| ZARA2 | 0.44 / 0.88 | 0.42 / 0.83 | 0.44 / 0.84 | 0.44 / 0.86 |
| Average | 0.58 / 1.13 | 0.51 / 0.99 | 0.53 / 1.03 | 0.54 / 1.09 |

Table 2: The ADEs / FDEs (in meters) of the four variants for single predictions in the ablation study.

than the other three variants. This is not unexpected as, without the time-specific context, the NAD is not able to forecast good predictions for different time steps in the prediction phase. Comparing the three interaction context based variants against each other, it is not surprising to see that NAP-ISS outperforms the other two variants due to the presence of both the social graph and scene features. As for NAP-ISg against NAP-ISc, we observe that NAP-ISg slightly outperforms NAP-ISc. This demonstrates that the social influence is more important than the scene influence. However, it should be noted that the five scenes in the ETH/UCY datasets do not have many obstacles scattered in the pedestrian pathways. The slightly better performance of NAP-ISg confirms that there are more social (pedestrian) interactions than scene interactions in these datasets.

Comparing the results from the four variants in Table 2 and from NAP in Table 1, we observe that NAP outperforms all the four variants. Our ablation study justifies the need for all the contexts to be present in NAP.

### 4.4 Different Prediction Lengths

In addition to the prediction length setting ($T_p = 12$ frames, corresponding to 4.8 seconds) used in Tables 1 & 2 and similar to previous work such as SGAN [Gupta *et al.*, 2018] and STGAT [Huang *et al.*, 2019], we conduct experiments for the prediction length $T_p = 8$ frames (or 3.2 seconds) to further evaluate the performance of NAP. Table 3 shows the average

| | $T_\mathrm{p} = 8$ | $T_\mathrm{p} = 12$ | Increment |
|---|---|---|---|
| Social-LSTM | 0.45 / 0.91 | 0.72 / 1.54 | 60.00% / 69.23% |
| SGAN 1V-1 | 0.49 / 1.00 | 0.74 / 1.54 | 51.02% / 54.00% |
| STGAT 1V-1 | 0.39 / 0.81 | 0.54 / 1.11 | 38.46% / 37.03% |
| NAP (ours) | 0.35 / 0.67 | 0.45 / 0.89 | **28.57% / 32.84%** |
| SGAN 20V-20* | 0.39 / 0.78 | 0.58 / 1.18 | 48.72% / 51.28% |
| STGAT 20V-20* | 0.31 / 0.62 | 0.43 / 0.83 | 38.71% / 33.87% |
| NAP (ours)* | 0.31 / 0.61 | 0.39 / 0.80 | **25.81% / 31.15%** |

Table 3: The ADEs / FDEs (in meters) of various methods for differ-ent prediction lengths. A method with * indicates that it generates 20 predictions for each input observed trajectory.



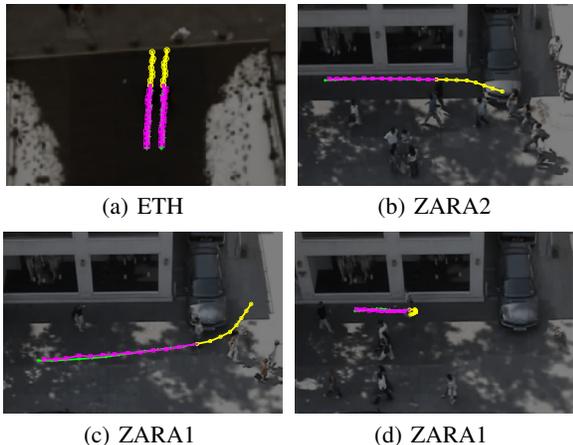(a) ETH      (b) ZARA2

(c) ZARA1      (d) ZARA1

Figure 3: Examples of predicted trajectories (shown in pink) gener-ated by NAP. The observed trajectories and ground truth trajectories are shown in yellow and green (figure best viewed in color).

ADE/FDE results for this prediction length setting. The fig-ures under the '$T_\mathrm{p} = 12$' column are copied from the *Average* column of Table 1. Each error increment (last column) due to the increase of $T_\mathrm{p}$ is calculated as: $(e_{\mathrm{p}12} - e_{\mathrm{p}8})/e_{\mathrm{p}8} \times 100\%$, where $e_{\mathrm{p}12}$ and $e_{\mathrm{p}8}$ are errors (ADEs or FDEs) for $T_\mathrm{p} = 12$ and $T_\mathrm{p} = 8$ of the same method.

As expected, all methods shown in Table 3 have better per-formance for the shorter prediction length. In the top half of the table, when generating a predicted trajectory for each in-put, the error increments of Social-LSTM and SGAN 1V-1 are over 50%. Compared to these two methods, STGAT 1V-1 has smaller error increments for both ADE and FDE. For the multimodal predictions (bottom half of the table), STGAT 20V-20 again outperforms SGAN 20V-20.

We observe from Table 3 that NAP consistently outper-forms all other methods for both prediction length settings and for both single and multimodal predictions. Furthermore, NAP also has the smallest error increments for both ADE and FDE when $T_\mathrm{p}$ increases. This demonstrates that NAP is more robust in generating long trajectories. The reason is due to the non-autoregressive nature of the decoder, which not only allows the location coordinates at different time steps to be independently forecast but also helps minimize the accumula-tion of prediction errors when the prediction length increases.



(a) HOTEL      (b) HOTEL

Figure 4: Examples of multiple predicted trajectories shown as heatmaps. The observed trajectories and ground truth trajectories are shown in yellow and green (figure best viewed in color).

## 4.5 Qualitative Results

Figure 3 illustrates some prediction examples generated by NAP. The observed and ground truth trajectories are shown in yellow and green; the best trajectory of the 20 predictions of each pedestrian is shown in pink. For better visualization, the video frames have been darkened and blurred. These exam-ples cover different movement behaviors of pedestrians. For example, Fig. 3(a) shows two simple straight path scenarios, Fig. 3(b) and (c) show a gentle turning scenario, and Fig. 3(d) shows a more difficult scenario in which an abrupt turning occurs near the end of the observation phase. Although the predicted trajectory (in pink) in Fig. 3(d) does not perfectly overlap with the ground truth trajectory, NAP is still able to correctly predict the trajectory from the late turning cue.

Figure 4 shows two more difficult scenarios where all the 20 predicted trajectories are displayed as a heatmap around each pedestrian. For the pedestrian in Fig. 4(a) and the right pedestrian in Fig. 4(b), each made an abrupt turn at almost the last frame of the observation phase. However, NAP is still able to give good predicted trajectories, as all the plausi-ble paths (including the ground truth trajectories (green)) are well covered by the heatmaps. The left pedestrian in Fig. 4(b) is a stopping scenario. After stopping, the pedestrian can re-main still or resume walking in any direction. The generated heatmap shows a good coverage of possible paths; however, it has a small dent in the bottom left hand region due to the presence of a bench there, showing that the pedestrian must bypass the obstacle. This example shows the importance of including scene influence in the method.

## 5 Conclusion

We have presented a novel method called NAP which can handle both social influence and scene influence in the pedes-trian trajectory prediction process. NAP captures these in-fluences using the trainable feature encoders in the network. In addition, NAP handles multimodal predictions via a la-tent variable generator which models the sampling distribu-tion that describes the multiple plausible paths of each pedes-trian. Unlike existing trajectory prediction methods, the de-coder of NAP is non-autoregressive. NAP is therefore able to forecast predictions for different time steps simultaneously or to forecast only for those time steps that are of interest. From our extensive experiments and ablation study, the context en-coders used in NAP have been demonstrated to be effective. Not only does NAP achieve state-of-the-art performance, it

also has smaller error increments when the prediction length increases.

# References

[Alahi *et al.*, 2016] Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social LSTM: Human trajectory prediction in crowded spaces. In *CVPR*, June 2016.

[Chen *et al.*, 2018] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, September 2018.

[Cordts *et al.*, 2016] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.

[Goodfellow *et al.*, 2014] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. In *NeurIPS*, pages 2672–2680, 2014.

[Gu *et al.*, 2018] Jiatao Gu, James Bradbury, Caiming Xiong, Victor O.K. Li, and Richard Socher. Non-autoregressive neural machine translation. In *ICLR*, 2018.

[Guo *et al.*, 2019] Junliang Guo, Xu Tan, Di He, Tao Qin, Linli Xu, and Tie-Yan Liu. Non-autoregressive neural machine translation with enhanced decoder input. In *AAAI*, volume 33, pages 3723–3730, 2019.

[Gupta *et al.*, 2018] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social GAN: Socially acceptable trajectories with generative adversarial networks. In *CVPR*, June 2018.

[Hasan *et al.*, 2018] Irtiza Hasan, Francesco Setti, Theodore Tsesmelis, Alessio Del Bue, Fabio Galasso, and Marco Cristani. MX-LSTM: Mixing tracklets and vislets to jointly forecast trajectories and head poses. In *CVPR*, June 2018.

[Huang *et al.*, 2019] Yingfan Huang, Huikun Bi, Zhaoxin Li, Tianlu Mao, and Zhaoqi Wang. STGAT: Modeling spatial-temporal interactions for human trajectory prediction. In *ICCV*, October 2019.

[Kingma and Welling, 2013] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[Kosaraju *et al.*, 2019] Vineet Kosaraju, Amir Sadeghian, Roberto Martín-Martín, Ian Reid, S. Hamid Rezatofighi, and Silvio Savarese. Social-bigat: Multimodal trajectory forecasting using bicycle-gan and graph attention networks. In *NeurIPS*, 2019.

[Lee *et al.*, 2017] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B. Choy, Philip H. S. Torr, and Manmohan Chandraker. DESIRE: distant future prediction in dynamic scenes with interacting agents. In *CVPR*, 2017.

[Lerner *et al.*, 2007] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. Crowds by example. In *Computer Graphics Forum*, volume 26, pages 655–664. Wiley Online Library, 2007.

[Li, 2019] Yuke Li. Which way are you going? imitative decision learning for path forecasting in dynamic scenes. In *CVPR*, June 2019.

[Liang *et al.*, 2019] 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*, June 2019.

[Nikhil and Morris, 2018] Nishant Nikhil and Brendan Tran Morris. Convolutional neural network for trajectory prediction. In *ECCV Workshop*, September 2018.

[Pellegrini *et al.*, 2009] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. You'll never walk alone: Modeling social behavior for multi-target tracking. In *ICCV*, pages 261–268, 2009.

[Sadeghian *et al.*, 2019] Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, Hamid Rezatofighi, and Silvio Savarese. SoPhie: An attentive gan for predicting paths compliant to social and physical constraints. In *CVPR*, June 2019.

[Su *et al.*, 2017] Hang Su, Jun Zhu, Yinpeng Dong, and Bo Zhang. Forecast the plausible paths in crowd scenes. In *IJCAI*, pages 2772–2778, 2017.

[Vemula *et al.*, 2018] A. Vemula, K. Muelling, and J. Oh. Social attention: Modeling attention in human crowds. In *ICRA*, pages 1–7, May 2018.

[Wen *et al.*, 2017] Ruofeng Wen, Kari Torkkola, Balakrishnan Narayanaswamy, and Dhruv Madeka. A multi-horizon quantile recurrent forecaster. In *NIPS Workshop*, 2017.

[Xue *et al.*, 2018] Hao Xue, Du Q Huynh, and Mark Reynolds. SS-LSTM: A hierarchical LSTM model for pedestrian trajectory prediction. In *WACV*, pages 1186–1194. IEEE, 2018.

[Xue *et al.*, 2019] Hao Xue, Du Q Huynh, and Mark Reynolds. Pedestrian trajectory prediction using a social pyramid. In *Pacific Rim International Conference on Artificial Intelligence*, pages 439–453. Springer, 2019.

[Zhang *et al.*, 2019] Pu Zhang, Wanli Ouyang, Pengfei Zhang, Jianru Xue, and Nanning Zheng. SR-LSTM: State Refinement for LSTM Towards Pedestrian Trajectory Prediction. In *CVPR*, June 2019.

[Zhao *et al.*, 2019] Tianyang Zhao, Yifei Xu, Mathew Monfort, Wongun Choi, Chris Baker, Yibiao Zhao, Yizhou Wang, and Ying Nian Wu. Multi-agent tensor fusion for contextual trajectory prediction. In *CVPR*, June 2019.

[Zou *et al.*, 2018] Haosheng Zou, Hang Su, Shihong Song, and Jun Zhu. Understanding human behaviors in crowds by imitating the decision-making process. In *AAAI*, 2018.