# Temporal Pyramid Network for Pedestrian Trajectory Prediction with Multi-Supervision

**Rongqin Liang,**[1, 2] **Yuanman Li,**[1, 2,*] **Xia Li,**[1, 2] **yi tang,**[1, 2] **Jiantao Zhou,**[3, 4] **Wenbin Zou**[1, 2]

[1] College of Electronics and Information Engineering, Shenzhen University
[2] Guangdong Key Laboratory of Intelligent Information Processing
[3] Department of Computer and Information Science, University of Macau
[4] State Key Laboratory of Internet of Things for Smart City
1810262064@email.szu.edu.cn, {lixia, yuanmanli, yitang, wzou}@szu.edu.cn, jtzhou@um.edu.mo

## Abstract

Predicting human motion behavior in a crowd is important for many applications, ranging from the natural navigation of autonomous vehicles to intelligent security systems of video surveillance. All the previous works model and predict the trajectory with a single resolution, which is rather inefficient and difficult to simultaneously exploit the long-range information (e.g., the destination of the trajectory), and the short-range information (e.g., the walking direction and speed at a certain time) of the motion behavior. In this paper, we propose a temporal pyramid network for pedestrian trajectory prediction through a squeeze modulation and a dilation modulation. Our hierarchical framework builds a feature pyramid with increasingly richer temporal information from top to bottom, which can better capture the motion behavior at various tempos. Furthermore, we propose a coarse-to-fine fusion strategy with multi-supervision. By progressively merging the top coarse features of global context to the bottom fine features of rich local context, our method can fully exploit both the long-range and short-range information of the trajectory. Experimental results on several benchmarks demonstrate the superiority of our method.

## 1  Introduction

Modeling the behaviors of pedestrians is an essential step for many applications, including self-driving platforms for safe decision making (Liang et al. 2019), socially-aware robots for natural navigation (Monfort, Liu, and Ziebart 2015) and surveillance systems to identify suspicious activities (Bastani, Marcenaro, and Regazzoni 2016). Trajectory prediction as one of the most important future behavior modeling tasks, aims to predict possible future trajectories according to historical paths in the last few seconds (Alahi et al. 2016; Gupta et al. 2018; Mohamed et al. 2020; Xu, Yang, and Du 2020). Despite its importance, predicting the trajectory is very challenging due to the inherent properties of pedestrians. First, human motions are highly *multimodal*, which means that there could be several socially-acceptable and distinct future behaviors under the same trajectory history. Second, human motions are highly affected by the people around them, Jointly modeling the complex social behaviors is rather challenging in reality.

Traditional pedestrian trajectory prediction algorithms heavily rely on the handcrafted rules to describe human motions (Helbing and Molnar 1995; Pellegrini et al. 2009), which are difficult to generalize in complex new scenes. Recently, the data-driven based algorithms have received significant attention in the community. Among them, RNN and its variant LSTM have been widely adopted. Social-LSTM (Alahi et al. 2016) as one of the earliest works on pedestrian trajectory prediction, encoded the motion information using a recurrent network. CIDNN (Xu, Piao, and Gao 2018) considered different importance of persons to a target pedestrian in a crowd interaction module. The recent works PIF (Liang et al. 2019) and SR-LSTM (Zhang et al. 2019) enhanced the performance of Social-LSTM by taking the scene context as side information. Though the RNN architecture endowed above methods to learn and predict trajectories in a data-driven manner, they failed to capture the multimodal nature of human.

In order to produce multiple socially-acceptable trajectories, some researchers suggested constructing the recurrent models with generative settings, which led to learning the distribution of the future trajectory rather than directly generating a deterministic path (Gupta et al. 2018; van der Heiden et al. 2019; Li, Ma, and Tomizuka 2019; Zhao et al. 2019). Social-GAN (Gupta et al. 2018) is the pioneering trajectory prediction work incorporating the LSTM model with the generative adversarial networks (GANs) (Goodfellow et al. 2014), permitting to produce multiple plausible trajectories. SoPhie (Sadeghian et al. 2019) improved social-GAN through a scene feature extraction component. Some researchers proposed to use graph to model the social interactions (Ma et al. 2019; Vemula, Muelling, and Oh 2018; Kosaraju et al. 2019; Huang et al. 2019a; Ivanovic and Pavone 2019; Shi et al. 2020; Mohamed et al. 2020). For example, the most recent work Social-STGCNN (Mohamed et al. 2020) modeled trajectories using the spatio-temporal graph convolution neural network, and achieved promising performance.

Though trajectory prediction has been studied from many aspects, all the existing methods encoded and decoded the trajectory with a single resolution (*i.e.*, a fixed length of time steps). This makes them fail to fully exploit the temporal relations of the motion behavior. We argue that simultaneously modeling the global context (*e.g.*, where the pedestrian plans
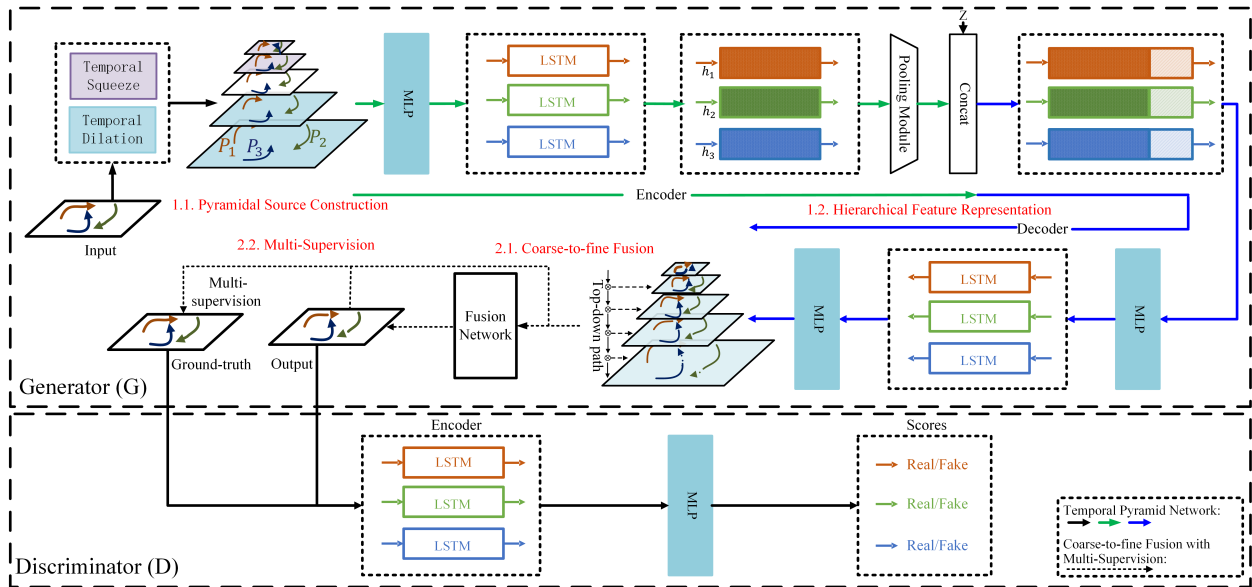
---

*The corresponding author.

Figure 1: The framework of our proposed TPNMS. The network consists of a generator and a discriminator. The input of the generator is the historical trajectories of pedestrians, and the output is the corresponding predicted future trajectories. The pyramidal source is first constructed through the temporal squeeze modulation and the temporal dilation modulation. Then, an encoder-decoder network is adopted for hierarchical feature learning. Features are finally fed into a fusion network (presented in Fig. 3) to generate the future trajectories with multi-supervision.

to go) and the local context (*e.g.*, the direction and speed at a certain time) with a single resolution is inefficient or rather difficult, if possible.

To alleviate the above limitation, in this work, we propose a novel Temporal Pyramid Network with Multi-Supervision (TPNMS) for pedestrian trajectory prediction. As shown in Fig. 1, our framework consists of a generator $G$ and a discriminator $D$, which are trained in opposition to each other. First, we devise a pyramid feature extractor composed of a squeeze module and a dilation module for multi-scale feature generation from a fixed length input trajectory. The pyramidal features are then fed into an RNN based Encoder-Decoder to generate hierarchical representations of the motion. To ensure effective representations of all pyramid levels, we further propose a coarse-to-fine fusion strategy with multi-supervision through progressively combining higher pyramid levels with lower ones. Finally, similar to Social-GAN, our network is trained in an adversarial manner to produce multiple socially-acceptable motion trajectories, conforming to the multimodal behavior of pedestrians.

It should be noted that most of the previous pyramid representation methods were designed in spatial domain and only for detection or recognition tasks. To the best of our knowledge, this is the *first* attempt that models trajectories in a scene as temporal pyramids. As will be shown later, our method outperforms previous approaches by a big margin on several datasets. The main contributions of our work are:

- A novel temporal pyramid network is proposed to capture the motion behaviors of pedestrians at various tempos. With our hierarchical design, both short-range and long-range motion behaviors can be effectively exploited.

- By progressively combining the global context with the local one, we further propose a coarse-to-fine trajectory modeling in a multi-supervised fashion.

- Our hierarchical design can be regarded as auxiliary modules, and easily extended to other sequence prediction frameworks, thus bringing performance improvements.

## 2 Proposed Temporal Pyramid Network with Multi-Supervision (TPNMS)

### Problem Formulation

Given a set of $N$ pedestrians in a scene with observed positions over a fixed duration, the trajectory prediction algorithm aims to jointly reason and forecast the future trajectories of all pedestrians. Let $(x_i^t, y_i^t)$ be the position of the $i$-th pedestrian at the time step $t$, where $i \in \{1, ..., N\}$. Denote $X_i^{(t_1:t_2)} = [(x_i^{t_1}, y_i^{t_1}), ..., (x_i^{t_2}, y_i^{t_2})]$ as the observed historical trajectory of the $i$-th pedestrian from the time step $t_1$ to $t_2$. Similarly, we define $Y_i^{(t_1:t_2)}$ as the future trajectory of the $i$-th pedestrian from the time $t_1$ to $t_2$. The trajectory prediction algorithm takes as input the previous trajectories with $t_o$ time steps of all pedestrians in a scene, denoted by

$$\mathcal{X} = \{X_1^{(1:t_o)}, ..., X_N^{(1:t_o)}\}, \qquad (1)$$

and aims to predict their trajectories in the next $t_p$ time steps simultaneously. We use $\mathcal{Y}$ to represent the true future trajectories, i.e.,

$$\mathcal{Y} = \{Y_1^{(t_o+1:t_o+t_p)}, ..., Y_N^{(t_o+1:t_o+t_p)}\}. \qquad (2)$$

For the sake of brevity, we hereafter will drop the superscript when there is no ambiguity, i.e., $X_i \triangleq X_i^{(1:t_o)}$ and

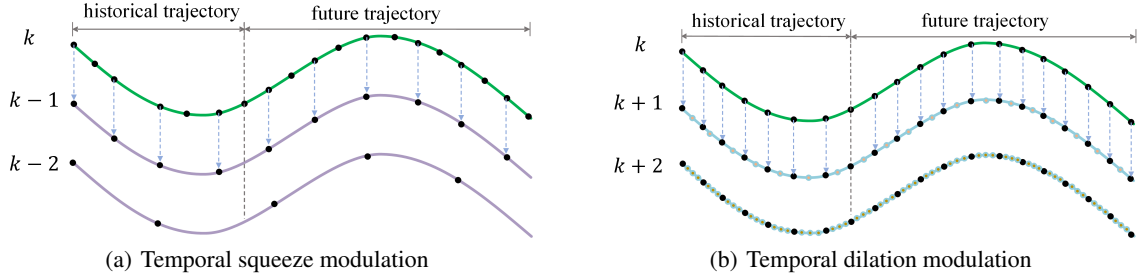(a) Temporal squeeze modulation          (b) Temporal dilation modulation

Figure 2: Illustration of the proposed (a) temporal squeeze modulation, and (b) temporal dilation modulation.

$Y_i \triangleq Y_i^{(t_o+1:t_o+t_p)}$. We further use $X$ and $Y$ to represent a generic historical trajectory and the corresponding future trajectory, respectively.

## TPN for Trajectory Prediction

Feature pyramids play a significantly important role in the field of computer vision for recognizing objects at vastly different scales (He et al. 2015). For example, the popular hand-engineered feature extractors such as SIFT (Lowe 2004) were designed to compute features in a multi-scale space. Lin *et. al.* (Lin et al. 2017) accommodated the idea of pyramid representation to deep convolutional neural networks, achieving quite promising performance in the detection task. Most of the previous approaches were designed in spatial domain. More recently, some works proposed to extract hierarchical features in temporal domain, and demonstrated its effectiveness in action recognition (Yang et al. 2020) and scene classification (Huang et al. 2019b).

Motivated by the great success of the pyramid representation, we propose a temporal pyramid network (TPN) tailored for pedestrian trajectory prediction. Compared with the existing algorithms which model the trajectory with a single resolution, our temporal pyramid architecture is effective in exploiting the motion behaviors at various tempos, and the hierarchical generation process could greatly facilitate the joint modeling of both global and local contexts. Besides, benefiting from the LSTM network, all levels of pyramids share the same parameters. This allows our method to operate on a single-branch backbone regardless of how many levels are adopted, then avoid to increase model complexity.

For better illustration, we decompose our TPN into the following two components: 1) pyramidal source construction, and 2) hierarchical feature representation.

**Pyramidal source construction**    For each trajectory of the input $\mathcal{X}$, we propose to generate a set of $L$ hierarchical features with multi-resolution, and then construct a feature pyramid, having increasingly richer temporal information from top to bottom. With the aid of the pyramid framework, our method can fully exploit the short-range behavior and the long-range behavior in a hierarchical way. As depicted in Fig. 1, this process can be summarized as two procedures, i.e., 1) the temporal squeeze modulation, and 2) the temporal dilation modulation.

*Temporal squeeze modulation:* Assume that there are to-

tally $L$ scales of the temporal pyramid network for each trajectory. Denote the feature of the $k$-th scale as $X_i^k$, which is identical to $X_i$. The goal of the temporal squeeze modulation is to reduce the impact of the local context, and generate a set of features with increasingly stronger global context from $\mathcal{X}$. To this end, we propose to gradually produce the top $k-1$ scales through uniformly sampling from the scale below with an interval factor 2. In this work, we refer to the above process as the temporal squeeze modulation.

Fig. 2(a) illustrates the procedure of the temporal squeeze modulation. For the $\ell$-scale ($\ell < k$), the feature can be represented as

$$X_i^\ell = [(\tilde{x}_i^1, \tilde{y}_i^1), ..., (\tilde{x}_i^{m_\ell}, \tilde{y}_i^{m_\ell})], \qquad (3)$$

where $m_\ell = \lceil t_o/2^{k-\ell} \rceil$, and

$$\tilde{x}_i^j = x_i^{1+2^{k-\ell}(j-1)}, \;\; \tilde{y}_i^j = y_i^{1+2^{k-\ell}(j-1)}. \qquad (4)$$

We can see that the detailed short-range information of the motion is gradually weakened by the temporal squeeze modulation, which encourages the higher scales to capture more long-range motion behaviors of pedestrians.

*Temporal dilation modulation:* Note that the observed trajectories are usually of short duration, then the number of scales generated by the temporal squeeze modulation cannot fully capture the motion behaviors. To handle this issue, we further introduce a complementary procedure called temporal dilation modulation, which is similar to the dilated convolution operator widely used in various vision tasks. The temporal dilation modulation could generate more dense trajectories for hierarchical feature representation, then exploit richer short-range information of the motion.

We propose to conduct the temporal dilation modulation through trajectory interpolation. It should be noted that pedestrians usually walk/run at varying speeds, accelerations and in different directions over time. In order to generate smooth dense trajectories, in this work, we adopt the cubic spline algorithm for the trajectory interpolation. For simplicity, we rewrite the observed trajectory of the $i$-th pedestrian as a series of time-position pairs

$$TX_i = \left[(1, (x_i^1, y_i^1)), ..., (t_o, (x_i^{t_o}, y_i^{t_o}))\right]. \qquad (5)$$

We adopt the cubic spline algorithm to seek for a piecewise-

cubic function $f(t)\colon \mathbb{R} \to \mathbb{R}^2$

$$f(t) = \begin{cases} f_1(t), & 1 \le t < 2 \\ \vdots & \\ f_{t_o-1}(t), & t_o - 1 \le t \le t_o \end{cases}, \quad (6)$$

where

$$f_k(t) = a_k + b_k(t-k) + c_k(t-k)^2 + d_k(t-k)^3 \quad (7)$$

represents the curve between the time steps $k$ and $k + 1$, and $a_k, b_k, c_k, d_k \in \mathbb{R}^2$ are parameters of the cubic spline. According to the cubic spline algorithm, given the trajectory $X_i$, there exists a unique set of parameters $\{a_k, b_k, c_k, d_k\}_{k=1,\cdots,t_o-1}$ such that the resulting trajectory curve passes through all the positions in $X_i$ with continuous velocity and acceleration at each position. Upon having $f(t)$, the feature $X_i^\ell$ at the $\ell$-scale ($\ell > k$) can be calculated by interpolating positions of the in-between and unobserved time steps as shown in Fig. 2(b). Mathematically, we have

$$X_i^\ell = [f(1), f(1+\tfrac{1}{c}), f(1+\tfrac{2}{c}), ..., f(2), ..., f(t_o - \tfrac{1}{c}), f(t_o)], \quad (8)$$

where $c = 2^{\ell-k}$. The interpolated dense trajectories offer more local information for the lower scales, permitting to capture more short-range motion behaviors of pedestrians. With above two modulations, we finally construct the pyramidal source as shown in Fig. 1.

**Hierarchical feature representation** For simplicity, we use a similar network architecture proposed in (Gupta et al. 2018) as the backbone to extract hierarchical features from the constructed pyramid. As shown in Fig. 1, the backbone network consists of two components, *i.e.*, the encoder and the decoder. At the encoder side, we embed the position of each pedestrian as

$$e_i^t = MLP(x_i^t, y_i^t; \Theta_{me}), \quad (9)$$

where $t \le t_o$ and $\Theta_{me}$ represents the parameters of MLP. The embedded feature $e_i^t$ is then fed into an LSTM block, which produces the hidden state at the time step $t$

$$h_i^t = LSTM(h_i^{t-1}, e_i^t; \Theta_{le}). \quad (10)$$

Note that the parameters of LSTM ($\Theta_{le}$) are shared among all the pedestrians and the scales of the pyramid.

In order to generate multiple socially-acceptable trajectories, our model is designed under the framework of GANs. According to GANs, at the decoder side, we concatenate a noise vector $z$ sampled from the standard normal distribution to the hidden state $h_i^{t_o}$. Further, we use the pooling module proposed by (Gupta et al. 2018) to encode the influence caused by pedestrians around. We write

$$h_i^{t_o} := [(h_i^{t_o}, P_i); z]. \quad (11)$$

Then for each $z$, we recurrently decode the hierarchical features of the trajectory as follows

$$\begin{aligned} \hat{e}_i^{t-1} &= MLP(\hat{x}_i^{t-1}, \hat{y}_i^{t-1}; \Theta_{md}) \\ h_i^t &= LSTM(h_i^{t-1}, \hat{e}_i^{t-1}; \Theta_{ld}), \quad (12) \\ (\hat{x}_i^t, \hat{y}_i^t) &= MLP(h_i^t; \Theta'_{md}) \end{aligned}$$

where $t \ge t_o + 1$, $(\hat{x}_i^{t_o}, \hat{y}_i^{t_o}) = (x_i^{t_o}, y_i^{t_o})$, and $\Theta_{md}, \Theta_{ld}, \Theta'_{md}$ are parameters to be learned.
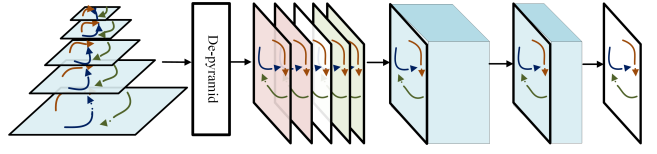


Figure 3: The framework of the fusion network, where the convolutional layers have the $1 \times 1$ kernel size, and the number of channels are $8, 4$ and $1$, respectively.

**Coarse-to-fine Fusion with Multi-supervision**

To merge and exploit the information of hierarchical features generated by the above TPN, we further propose a coarse-to-fine fusion strategy with multi-supervision.

**Coarse-to-fine fusion** As shown in Fig. 1, the coarse-to-fine fusion organizes features in a top-down pathway, where the top coarse features with long-range context are progressively merged to the bottom fine features with rich local-range context. Denote the extracted feature of the $\ell$-scale as $\hat{X}_i^\ell$, which is updated by merging the information of the above scale. We write

$$\hat{X}_i^\ell := \frac{1}{2}(\hat{X}_i^\ell \oplus \hat{X}_{i,\uparrow}^{\ell-1}), \quad (13)$$

where $\hat{X}_{i,\uparrow}^{\ell-1}$ means upsampling $\hat{X}_i^{\ell-1}$ by a factor of 2, and $\oplus$ serves as the element-wise addition. This process is iterated until the finest resolution feature is merged. With the coarse-to-fine fusion strategy, the long-range and local-range motion information is collaborated, and then can complete to each other.

**Multi-supervision** To ensure effective representation of each pyramid level, all the scales are supervised during training, where the corresponding loss function is formulated as

$$\mathcal{L}_s = \frac{1}{NL} \sum_{i=1}^{N} \sum_{\ell=1}^{L} \lambda_\ell \|\hat{X}_i^\ell - Y_i^\ell\|_2^2. \quad (14)$$

Here $Y_i^\ell$ is the ground-truth pyramidal source of the future trajectory, which can be constructed from $Y_i$ in the same way detailed in Section 2. The hyper-parameter $\lambda_\ell$ is inversely proportional to the feature length of $\hat{X}_i^\ell$, which we empirically set

$$\lambda_\ell = \frac{t_p}{m'_\ell}. \quad (15)$$

Here $m'_\ell$ represents the length of $\hat{X}_i^\ell$.

The final predicted trajectory is produced by a fusion layer as shown in Fig. 3, where a de-pyramid layer is first adopted to down-sample or up-sample the hierarchical features $\{\hat{X}_i^\ell\}_{\ell=1}^{L}$ to a fixed length $t_p$. The results are then concatenated as a tensor of size $L \times 2 \times t_p$, which is further processed through three convolutional layers to fuse information across the whole pyramid. The fusion layer finally generates the predicted trajectory $\hat{Y}_i$. We supervise the final output using the loss

$$\mathcal{L}_f = \frac{1}{N} \sum_{i=1}^{N} \|\hat{Y}_i - Y_i\|_2^2. \quad (16)$$

| Datasets | ETH | | Hotel | | Univ | | Zara1 | | Zara2 | | AVG | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Methods | ADE | FDE | ADE | FDE | ADE | FDE | ADE | FDE | ADE | FDE | ADE | FDE |
| Linear | 1.33 | 2.94 | 0.39 | 0.72 | 0.82 | 1.59 | 0.62 | 1.21 | 0.77 | 1.48 | 0.79 | 1.59 |
| S-LSTM | 1.09 | 2.35 | 0.79 | 1.76 | 0.67 | 1.40 | 0.47 | 1.00 | 0.56 | 1.17 | 0.72 | 1.54 |
| S-GAN | 0.81 | 1.52 | 0.72 | 1.61 | 0.60 | 1.26 | 0.34 | 0.69 | 0.42 | 0.84 | 0.58 | 1.18 |
| S-GAN-P | 0.87 | 1.62 | 0.67 | 1.37 | 0.76 | 1.52 | 0.35 | 0.68 | 0.42 | 0.84 | 0.61 | 1.21 |
| PIF | 0.73 | 1.65 | 0.30 | 0.59 | 0.60 | 1.27 | 0.38 | 0.81 | 0.31 | 0.68 | 0.46 | 1.00 |
| SoPhie | 0.70 | 1.43 | 0.76 | 1.67 | 0.54 | 1.24 | **0.30** | 0.63 | 0.38 | 0.78 | 0.54 | 1.15 |
| STGAT | 0.65 | 1.12 | 0.35 | 0.66 | 0.52 | 1.10 | 0.34 | 0.69 | 0.29 | 0.60 | 0.43 | 0.83 |
| SR-LSTM | 0.63 | 1.25 | 0.37 | 0.74 | 0.51 | 1.10 | 0.41 | 0.90 | 0.32 | 0.70 | 0.45 | 0.94 |
| Social-BiGAT | 0.69 | 1.29 | 0.49 | 1.01 | 0.55 | 1.32 | **0.30** | 0.62 | 0.36 | 0.75 | 0.48 | 1.00 |
| Social-STGCNN | 0.64 | 1.11 | 0.49 | 0.85 | **0.44** | **0.79** | 0.34 | **0.53** | 0.30 | **0.48** | 0.44 | 0.75 |
| **TPNMS** | **0.52** | **0.89** | **0.22** | **0.39** | 0.55 | 1.13 | 0.35 | 0.70 | **0.27** | 0.56 | **0.38** | **0.73** |

Table 1: The performance of different methods in terms of ADE / FDE metrics.

**Adversarial training**   The architecture of the discriminator is shown in Fig. 1, which consists of an LSTM component and an MLP component. The LSTM component takes as input the ground-truth trajectory $[X, Y]$ or the generated trajectory $[X, \hat{Y}]$. The last hidden state of the LSTM is fed into the MLP, which outputs the classification score. Let $\hat{Y} = G(z, X)$. The adversarial loss is defined as

$$
\begin{aligned}
\mathcal{L}_{avd} = \; & \mathbb{E}_{X,Y \sim P_{data}(X,Y)}[\log D(X, Y)] \\
& + \mathbb{E}_{X \sim P_{data}(X), z \sim P_z(z)}[\log(1 - D(X, G(z, X)))].
\end{aligned}
\tag{17}
$$

Finally, training the network is cast into a two-player min-max game with the following objective function

$$
\min_{G} \max_{D} \mathcal{L}_{avd} + \mathcal{L}_s + \mathcal{L}_f.
\tag{18}
$$

Above problem can be solved by alternatively updating the generator $G$ and the discriminator $D$.

## 3    Experimental Results

We implement our model TPNMS using the PyTorch framework with an NVIDIA TITAN Xp GPU. All the source code and models will be publicly available upon the acceptance.

**Implementation Details**

The number of pyramid scales is empirically set as 5, and the dimensions of the hidden state for the encoder and decoder are 32. Each input coordinate $(x, y)$ is embedded as a 16-dimensional vector. The length of the noise vector $z$ is 8. We adopt Adam algorithm (Kingma and Ba 2014) to optimize the loss function (18) and train our network with the following hyper-parameter settings: batch size is 64; learning rates for the Generator and Discriminator are set to be 1e-4 and 2e-4, respectively; betas are 0.9 and 0.999; weight decay is 1e-4 and the number of epochs is 400.

**Datasets and Metrics**

*Datasets*: We evaluate our method on two public datasets, *i.e.*, ETH(Pellegrini et al. 2009) and UCY(Lerner, Chrysanthou, and Lischinski 2007). These datasets consist of 5 unique scenes: ETH, HOTEL, UNIV, ZARA1 and ZARA2 with 4 different scenes. There are totally 1536 pedestrians

with thousands of trajectories containing challenging behaviors such as walking together, crossing each other, forming groups and dispersing.

*Metrics*: For the sake of fairness, we use the widely adopted leave-one-out approach evaluation methodology. The number of observed time steps is 8 (3.2 seconds) of each person and the upcoming trajectory of 12 time steps (4.8 seconds) is used to predict. Following prior works, we use two error metrics to evaluate the performance of different pedestrian trajectory prediction models.

1. *Average Displacement Error* (ADE): The average Euclidean distance between the ground-truth trajectory and the predicted one,

$$
\text{ADE} = \frac{\sum_{i=1}^{N} \sum_{t=t_o+1}^{t_o+t_p} \left\| Y_i^{(t)} - \hat{Y}_i^{(t)} \right\|_2}{N \times t_p}.
\tag{19}
$$

2. *Final Displacement Error* (FDE): The Euclidean distance between the ground-truth destination and the predicted one,

$$
\text{FDE} = \frac{\sum_{i=1}^{N} \left\| Y_i^{(t_o+t_p)} - \hat{Y}_i^{(t_o+t_p)} \right\|_2}{N}.
\tag{20}
$$

**Baselines**

We compare our method *TPNMS* with following approaches: *Linear* (Alahi et al. 2016): a linear regressor that predicts the next coordinates based on previous points. *S-LSTM* (Alahi et al. 2016): a method based on LSTM and social pooling. *S-GAN* and *S-GAN-P* (Gupta et al. 2018): a model that employs GAN to generate multimodal pedestrian trajectories and the latter with a global pooling module. *PIF* (Liang et al. 2019): a multi-task method using both visual features and interaction information. *SoPhie* (Sadeghian et al. 2019): an improved GAN based model considering the physical constraints. *SR-LSTM* (Zhang et al. 2019): a state refinement method for LSTM based pedestrian trajectory prediction. *Social-BiGAT* (Kosaraju et al. 2019) and *STGAT* (Huang et al. 2019a): methods based on GAN and graph attention. *Social-STGCNN* (Mohamed et al. 2020): an approach that models the social behavior of pedestrians

using a graph. Similar to previous works, we generate 20 samples based on the predicted distribution.

## Quantitative Analysis

Table 1 summarizes the results of different algorithms, where we report the average results for each method in the last two columns. We can see that all the algorithms perform much better than the linear model. Based on the results, we further draw the following conclusions:

- Overall, our method TPNMS outperforms all the previous approaches in terms of the average ADE and FDE.

- Compared with the baseline approach S-GAN (Gupta et al. 2018), TPNMS achieves significant performance gains. For example, S-GAN has an average error of 0.58 on ADE, and 1.18 on FDE, while TPNMS has much lower ADE (0.38) and FDE (0.73), corresponding to $34\%$ and $38\%$ relative improvements, respectively. This demonstrates that our proposed temporal pyramid network with multi-supervision indeed helps for pedestrian trajectory prediction.

- For the previous state-of-the-art method Social-STGCNN (Mohamed et al. 2020), TPNMS still achieves noticeable performance gains. For instance, TPNMS decreases the error of about $14\%$ on ADE and about $3\%$ on FDE compared with Social-STGCNN.

- Even without using any side information, TPNMS outperforms those methods utilizing the scene context, such as PIF (Liang et al. 2019), Sophie (Sadeghian et al. 2019) and Social-BiGAT (Kosaraju et al. 2019). This implies that the performance of TPNMS could potentially be improved by considering the scene context.

## Qualitative Analysis

In this subsection, we provide some examples to show how our TPNMS successfully captures complex motion behaviors of pedestrians. We qualitatively compare the prediction results between Social-GAN and TPNMS.

**Results in different interaction scenarios** We visualize examples from 4 scenarios in Figure 4.

**Walking in parallel** When people are walking side by side, they usually have tight connection to each other, and their relative positions tend to be preserved and motion behaviors tend to change consistently. In the Fig. 4(a), two target pedestrians A and B are walking in parallel. It can be noticed that S-GAN incorrectly predicts that these two pedestrians will walk across each other, and have a high possibility of collision. Compared with S-GAN, the predictions by our TPNMS show that these two pedestrians will keep walking in parallel, which is close to the ground-truth trajectories marked by green lines. This demonstrates the superiority of modeling motion behavior at various tempos.

**Meeting from opposite directions** People avoiding each other when moving in opposite direction is common in reality. Fig. 4(b) presents a scenario where two groups are meeting from opposite directions. We can see that the local behaviors of persons A, B and C are adjusted slightly to avoid collision. Compared with S-GAN, the trajectory of the person A predicted by TPNMS is more accurate after meeting. Further, TPNMS successfully predicts that persons B and C will keep walking in parallel, while the forecasts of S-GAN deviate from their true behaviors.

**Following people** When a person is following someone, he or she might want to draw attention to the person ahead, and maintain a safe distance between them. Fig. 4(c) shows a situation where the person A is walking behind the person B. We can see that S-GAN tends to decrease the speed of person A even when the distance between others is sufficiently large, while our proposed models TPNMS can more accurately forecast the speed at each time step, and still preserve a safe distance to avoid collision.

**Walking with complex social interactions** The complex interactions drive people using various ways to avoid collisions. As shown in Fig. 4(d), many trajectories generated by S-GAN have large deviations from the ground-truth ones, e.g., the persons A, F, H and I. Besides, we can see that S-GAN fails to adjust the behaviors of persons A, B and C, and then causes a collision at the end of the predicted trajectories. As for TPNMS, the predicted trajectories match better with the ground-truth one. For instance, persons I and J are maintained walking in parallel. Furthermore, we observe that the speed of person A is clearly slowed down by TPNMS, to avoid collision with persons B and C. For the trajectories of persons F and H, TPNMS achieves much better prediction accuracy than S-GAN. This demonstrates that our method can effectively capture the motion behaviors of pedestrians in scenarios of complex social interactions.

**Results of diverse predictions** Our model is capable of producing multiple plausible and diverse trajectories conforming to the multimodal behavior of pedestrians. In Fig. 5, we show some examples of diverse predictions by sampling the noise vector $z$ from the standard normal distribution. We can see from Fig. 5(b) and Fig. 5(c) that our model generates two socially acceptable and distinct trajectories with different $z$, including changing the direction and speed. For instance, the top image of Fig. 5(b) shows that the person is walking toward the car, where the direction is different from the true trajectory but the predicted path is still acceptable. Similar phenomenon can also be observed from the bottom image of Fig. 5(b). Besides, images presented in Fig. 5(c) show that $z$ can also affect the speed of pedestrians. In Fig. 5(d), we draw the density of the predicted trajectory by 20 randomly generated samples. The purple area constructs a plausible area that each pedestrian may pass. The position of darker color indicates a higher probability that the person will pass through. Furthermore, in Fig. 5(d), we also plot the best predicted trajectory from 20 samples for each scenario, and we can see that it closely matches the true trajectory shown in Fig. 5(a).

## Ablation Experiments

In Table 2, we systematically evaluate our method through a series of ablation experiments, where we consider the following variants of our method:
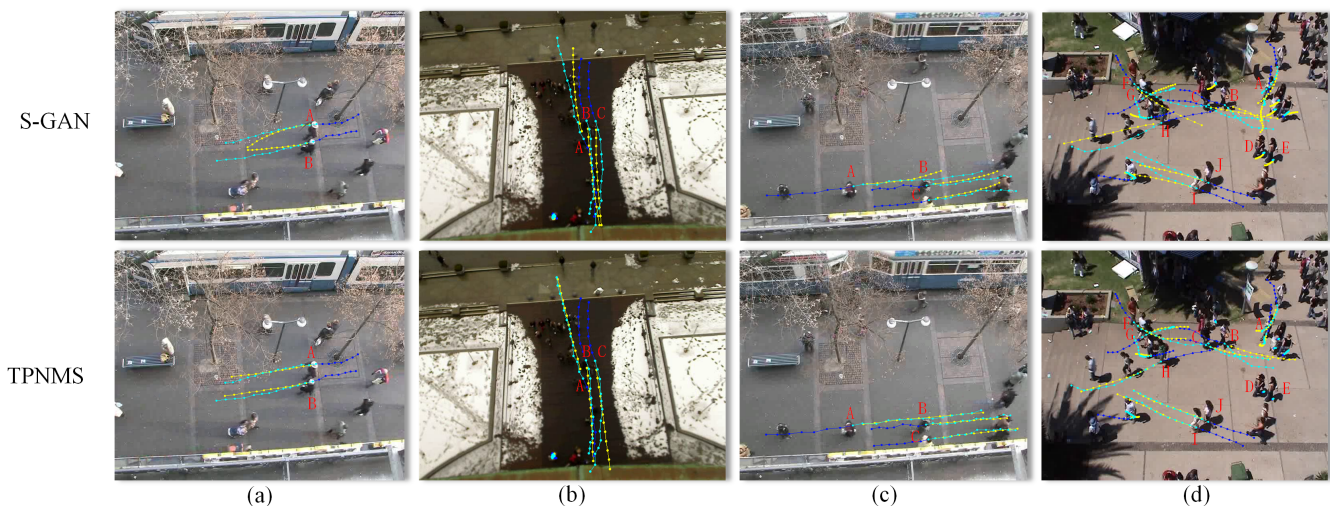*S-GAN-P*: the method without the temporal pyramid mod-

Figure 4: Examples of predicted trajectories by different methods. (a) walk in parallel; (b) meet from opposite directions; (c) follow people and (d) walk with complex social interactions. Blue line represents the historical trajectory; green line denotes the true future trajectory; yellow line shows the predicted future trajectory, and dots are the locations at different time steps.
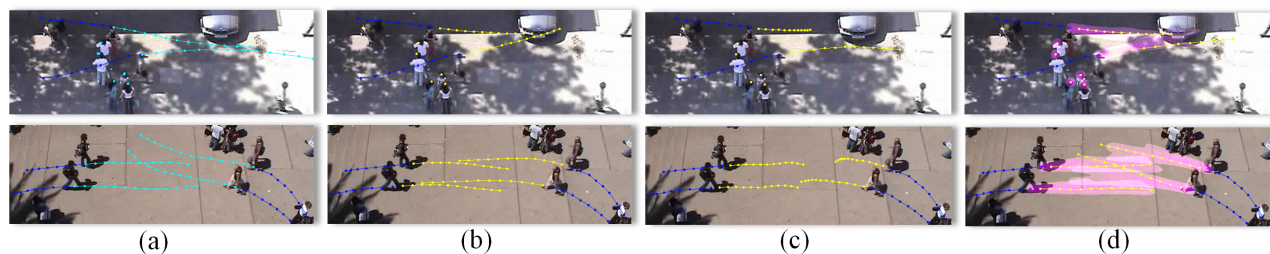


Figure 5: Examples of diverse predictions. The historical trajectory and predicted trajectory are marked in blue and yellow, respectively. (a) The ground-truth future trajectory (green). (b, c) two examples of diverse predictions. (d) The density of the prediction, where the purple area is the visualization result of the predicted 20 pedestrian trajectories after mean filtering.

| Models | Modules | | |
| | TP Layer | Multi-Supervision | AVG (ADE/FDE) |
| --- | --- | --- | --- |
| S-GAN-P | × | × | 0.61/1.21 |
| TPN | ✓ | × | 0.41/0.79 |
| TPNMS | ✓ | ✓ | **0.38/0.73** |

Table 2: The average ADE/FDE performance of variants.

that TPNSM can further improve the prediction accuracy in terms of ADE/FDE metrics, which demonstrates the importance of multi-supervision to ensure effective hierarchical representations.

## 4 Conclusion

In this paper, we have proposed a novel pyramid architecture for pedestrian trajectory prediction, which outperforms state-of-the-art methods on several benchmark datasets. First, we have devised a temporal pyramid network through squeeze and dilation modulations, which encodes and decodes the trajectory at multiple resolutions. This enables our method to capture both short-range and long-range motion behaviors of pedestrians. By resorting to a coarse-to-fine fusion strategy and the multi-supervision, our method can progressively merge high-scale global context with low-scale local context, finally resulting in an accurate trajectory prediction. Finally, with a GAN based framework, our method can generate multiple socially-acceptable trajectories conditioned on the same trajectory history, obeying the multimodal property of pedestrians. Both quantitative and qual-

ule and the multi-supervision module. With this setting, our method degrades to S-GAN-P;
*TPN*: the method only considers the temporal pyramid module without the multi-supervision;
*TPNMS*: the method considers both the temporal pyramid module and the multi-supervision.

Comparing TPN with S-GAN-P, we can see that TPN significantly reduces the ADE from 0.61 to 0.41, and the FDE from 1.21 to 0.79, which indicates that our temporal pyramid architecture can more effectively model the global context and local context of trajectories. Further, we observe

itative experimental results demonstrate the promising performance of our method under various situations.

# References

Alahi, A.; Goel, K.; Ramanathan, V.; Robicquet, A.; Fei-Fei, L.; and Savarese, S. 2016. Social LSTM: Human trajectory prediction in crowded spaces. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 961–971.

Bastani, V.; Marcenaro, L.; and Regazzoni, C. S. 2016. Online Nonparametric Bayesian Activity Mining and Analysis From Surveillance Video. *IEEE Trans. Image Process.* 25(5): 2089–2102.

Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. In *Advances in neural information processing systems*, 2672–2680.

Gupta, A.; Johnson, J.; Fei-Fei, L.; Savarese, S.; and Alahi, A. 2018. Social GAN: Socially acceptable trajectories with generative adversarial networks. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 2255–2264.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2015. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 37(9): 1904–1916.

Helbing, D.; and Molnar, P. 1995. Social force model for pedestrian dynamics. *Phy. rev. E* 51(5): 4282.

Huang, Y.; Bi, H.; Li, Z.; Mao, T.; and Wang, Z. 2019a. STGAT: Modeling Spatial-Temporal Interactions for Human Trajectory Prediction. In *Proc. IEEE Int. Conf. Comput. Vis.*, 6272–6281.

Huang, Y.; Cao, X.; Zhen, X.; and Han, J. 2019b. Attentive temporal pyramid network for dynamic scene classification. In *Proc. AAAI Conf. Art. Intel.*, volume 33, 8497–8504.

Ivanovic, B.; and Pavone, M. 2019. The Trajectron: Probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs. In *Proc. IEEE Int. Conf. Comput. Vis.*, 2375–2384.

Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* .

Kosaraju, V.; Sadeghian, A.; Martín-Martín, R.; Reid, I.; Rezatofighi, S. H.; and Savarese, S. 2019. Social-BiGAT: Multimodal Trajectory Forecasting using Bicycle-GAN and Graph Attention Networks. *arXiv preprint arXiv:1907.03395* .

Lerner, A.; Chrysanthou, Y.; and Lischinski, D. 2007. Crowds by example. *Computer graphics forum* 26(3): 655–664.

Li, J.; Ma, H.; and Tomizuka, M. 2019. Conditional Generative Neural System for Probabilistic Trajectory Prediction. In *Proc. IEEE Int. Conf. Intel. Robots and Sys.*

Liang, J.; Jiang, L.; Niebles, J. C.; Hauptmann, A. G.; and Fei-Fei, L. 2019. Peeking into the future: Predicting future person activities and locations in videos. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 5725–5734.

Lin, T.-Y.; Dollár, P.; Girshick, R.; He, K.; Hariharan, B.; and Belongie, S. 2017. Feature pyramid networks for object detection. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 2117–2125.

Lowe, D. G. 2004. Distinctive image features from scale-invariant keypoints. *Int. J. of Comp. Vision* 60(2): 91–110.

Ma, Y.; Zhu, X.; Zhang, S.; Yang, R.; Wang, W.; and Manocha, D. 2019. Trafficpredict: Trajectory prediction for heterogeneous traffic-agents. In *Proc. AAAI Conf. Art. Intel.*, volume 33, 6120–6127.

Mohamed, A.; Qian, K.; Elhoseiny, M.; and Claudel, C. 2020. Social-STGCNN: A Social Spatio-Temporal Graph Convolutional Neural Network for Human Trajectory Prediction. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 14424–14432.

Monfort, M.; Liu, A.; and Ziebart, B. D. 2015. Intent Prediction and Trajectory Forecasting via Predictive Inverse Linear-Quadratic Regulation. In *Proc. AAAI Conf. Art. Intel.*, 3672–3678.

Pellegrini, S.; Ess, A.; Schindler, K.; and Van Gool, L. 2009. You'll never walk alone: Modeling social behavior for multi-target tracking. In *Proc. IEEE Int. Conf. Comput. Vis.*, 261–268.

Sadeghian, A.; Kosaraju, V.; Sadeghian, A.; Hirose, N.; Rezatofighi, H.; and Savarese, S. 2019. Sophie: An attentive gan for predicting paths compliant to social and physical constraints. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 1349–1358.

Shi, X.; Shao, X.; Fan, Z.; Jiang, R.; Zhang, H.; Guo, Z.; Wu, G.; Yuan, W.; and Shibasaki, R. 2020. Multimodal Interaction-Aware Trajectory Prediction in Crowded Space. In *Proc. AAAI Conf. Art. Intel.*, 11982–11989.

van der Heiden, T.; Nagaraja, N. S.; Weiss, C.; and Gavves, E. 2019. SafeCritic: Collision-Aware Trajectory Prediction. *arXiv preprint arXiv:1910.06673* .

Vemula, A.; Muelling, K.; and Oh, J. 2018. Social attention: Modeling attention in human crowds. In *Proc. IEEE int. Conf. Robot. and Auto.*, 1–7.

Xu, Y.; Piao, Z.; and Gao, S. 2018. Encoding crowd interaction with deep neural network for pedestrian trajectory prediction. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 5275–5284.

Xu, Y.; Yang, J.; and Du, S. 2020. CF-LSTM: Cascaded Feature-Based Long Short-Term Networks for Predicting Pedestrian Trajectory. In *Proc. AAAI Conf. Art. Intel.*, 12541–12548.

Yang, C.; Xu, Y.; Shi, J.; Dai, B.; and Zhou, B. 2020. Temporal pyramid network for action recognition. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 591–600.

Zhang, P.; Ouyang, W.; Zhang, P.; Xue, J.; and Zheng, N. 2019. Sr-lstm: State refinement for lstm towards pedestrian trajectory prediction. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 12085–12094.

Zhao, T.; Xu, Y.; Monfort, M.; Choi, W.; Baker, C.; Zhao, Y.; Wang, Y.; and Wu, Y. N. 2019. Multi-agent tensor fusion for contextual trajectory prediction. In *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 12126–12134.