

A DiffTraj Details & Hyperparameters

In this section, we cover the specific details of DiffTraj, including the DiffTraj framework, the Traj-UNet structure, and the implementation details.

A.1 Architecture

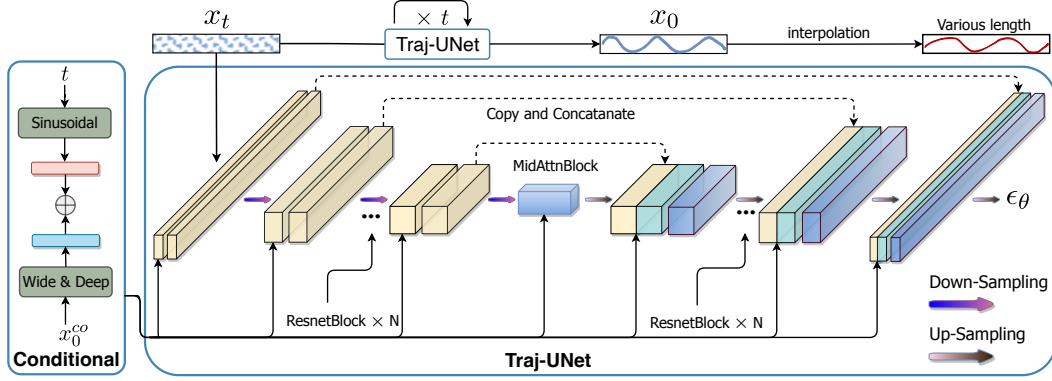


Figure 4: The network architecture used by DiffTraj in modeling $\epsilon_\theta(x_t^s, t | x_0^{co})$ is divided into two modules, down-sampling and up-sampling, each containing multiple Resnet blocks.

As illustrated in Fig. 4, DiffTraj is divided into two modules, i.e., down-sampling and up-sampling, and conditional module. Each down-sampling and up-sampling module consists of multiple stacked Resnet blocks. Between the two of them, a transitional module based on the attention mechanism is integrated. To better learn the noise of each time step and guide the generation, DiffTraj integrates a conditional module to embed the time step and external traffic information, later fed to each block. Since the CNN structure can only accept data of fixed shape, we first sample each trajectory as a tensor of $[2, \text{length}]$. Specifically, if the trajectory is below the set length, it is added using linear interpolation, and if it is greater than that, the redundant portion is removed using linear interpolation. Thus, the model will generate fixed-length trajectories, which will then be tailored to the desired length by the conditional information.

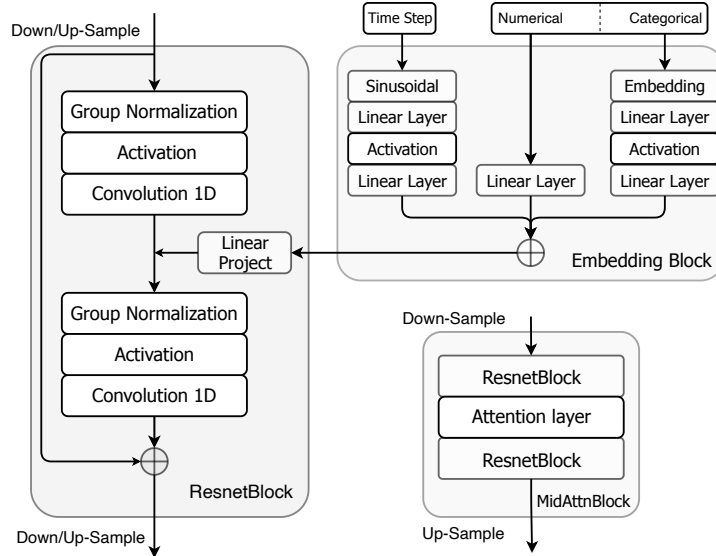


Figure 5: The main components of Traj-UNet including Resnet block, Embedding block, and middle attention block.

The main blocks of Traj-UNet are presented in Fig. 5, i.e., Resnet block, Embedding block, and middle attention block. Among them, each sampling block (down-sampling and up-sampling) consists of multiple Resnet blocks, each containing a series of group normalization, nonlinear activation, and 1D-CNN layers. Then, Traj-UNet applies up-sampling or down-sampling to the output, where down-sampling uses max pooling and up-sampling uses interpolation. After this, Traj-UNet integrates a middle attention block, which consists of two Resnet blocks and an attention layer. Note that there are no additional down/up-sampling operations in the Resnet block. Finally, we integrate a conditional embedding block to learn the diffusion time step and conditional information. For the diffusion step, we employ Sinusoidal embedding to represent each t as a 128-dimensional vector, and then apply two shared-parameter fully connected layers. For conditional information, such as distance, speed, departure time, travel time, trajectory length, and starting and ending locations, we use the Wide & Deep module for embedding. After getting the diffusion step embedding and conditional embedding, we sum them up and add them to each Resnet block.

A.2 Implementation Details

For the proposed DiffTraj framework, we summarize the adopted hyperparameter settings in Table 3.

Table 3: Hyperparameters setting for DiffTraj.

Hyperparameter	Setting value
Diffusion Steps	500
Skip steps	5
Guidance scale	3
β (linear)	0.0001 \sim 0.05
Batch size	1024
Sampling blocks	4
Resnet blocks	2
Input Length	200

500

501 The training and sampling phase of the proposed framework is summarized in Algorithm 1 and
502 Algorithm 2, respectively. The detailed implementation code is attached in Supplementary.

Algorithm 1 Diffusion Training Phase

```

1: for  $i = 1, 2, \dots$ , do
2:   Sample  $x_0 \sim q(x)$ ,
3:    $t \sim \text{Uniform}\{1, \dots, T\}$ 
4:    $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ 
5:    $\mathcal{L} = \|\epsilon - \epsilon_\theta(\sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon, t)\|^2$ 
6:    $\theta = \theta - \eta \nabla_\theta \mathcal{L}$ 
7: end for
```

Algorithm 2 Diffusion Sampling Phase

```

1: Sample  $x_T^s \sim \mathcal{N}(0, \mathbf{I})$ 
2: for  $t = T, T - S, \dots, 1$  do
3:   Compute  $\mu_\theta(x_t^s, t \mid x_0^{\text{co}})$  according to Eq. (5)
4:   Compute  $p_\theta(x_{t-1}^s \mid x_t^s, x_0^{\text{co}})$  according to Eq. (4)
5: end for
6: return  $x_0$ 
```

B Details of the Experimental Setup

B.1 Dataset

We evaluate the performance of DiffTraj and all baselines methods on two datasets with different cities, **Chengdu** and **Xi'an**¹. Both datasets are collected from cab trajectory data starting from November 1, 2016, to November 30, 2016. Table 4 summarizes the statistical information of these two datasets, and Fig. 6 shows the trajectory distribution and heat map of these two datasets, where the deeper color indicates the more concentrated trajectory in the region. For all datasets, we remove all trajectories with lengths less than 120 and sample them to a set fixed length.

Table 4: Statistics of Two Real-world Trajectory Datasets.

Dataset	Trajectory Number	Average Time	Average Distance
Chengdu	3 493 918	11.42 min	7.42 km
Xian	2 180 348	12.58 min	5.73 km

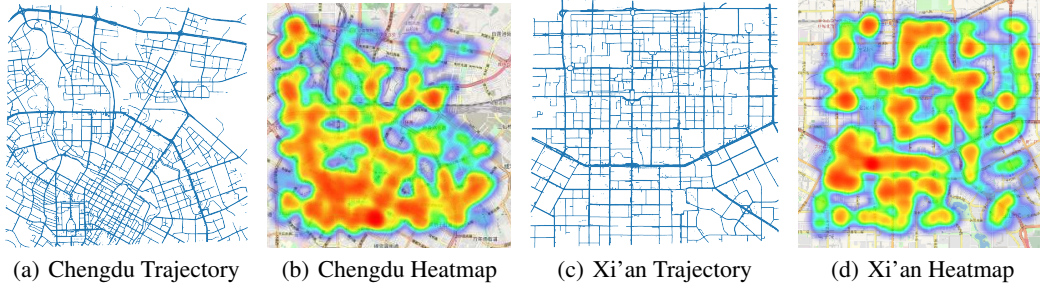


Figure 6: Origin trajectory distribution of two cities.

B.2 Evaluation Metrics

As trajectory generation aims to generate trajectories that can replace real-world activities and further benefit downstream tasks, we need to evaluate the “similarity” between the generated trajectories and real ones. In this work, we follow the common practice in previous studies [8] and measure the quality of the generated ones by Jensen-Shannon divergence (JSD). JSD compares the distribution of the real and generated trajectories, and a lower JSD indicates a better match with the original statistical features. Suppose that the original data has a probability distribution P and the generated data has a probability distribution G , the JSD is calculated as follows:

$$\text{JSD}(P, G) = \frac{1}{2} \mathbb{E}_P \left[\log \frac{P}{P+G} \right] + \frac{1}{2} \mathbb{E}_G \left[\log \frac{G}{G+P} \right]. \quad (8)$$

For the evaluation, we divided each city into 16×16 size grids and recorded the corresponding values for each grid. Based on this, we adopt the following metrics to evaluate the quality of the generated trajectories from four perspectives:

- **Density error:** This is a global level metric that is used to evaluate the geo-distribution between the entire generated trajectory $\mathcal{D}(\mathcal{T}_{\text{gen}})$ and the real trajectory $\mathcal{D}(\mathcal{T})$.

$$\text{Density Error} = \text{JSD}(\mathcal{D}(\mathcal{T}), \mathcal{D}(\mathcal{T}_{\text{gen}})), \quad (9)$$

where $\mathcal{D}(\cdot)$ denotes the grid density distribution in a given trajectory set, and $\text{JSD}(\cdot)$ represents the Jensen-Shannon divergence between two distributions.

- **Trip error:** This is a trajectory level metric that measures the correlation between the starting and ending points of a travel trajectory. Specifically, we calculate the probability distribution of the start/end points in the original and generated trajectories and use JSD to measure the difference between them.

¹These datasets can be downloaded at <https://outreach.didichuxing.com/>

- 530 • **Length error:** This is a trajectory level metric to evaluate the distribution of travel distances. It can
531 be obtained by calculating the Euclidean distance between consecutive points.
- 532 • **Pattern score:** This is a semantic level metric defined as the top- n grids that occur most frequently
533 in the trajectory. We define P and P_{gen} to denote the original and generated pattern sets, respectively,
534 and compute the following metrics:

$$\text{Pattern score} = 2 \times \frac{\text{Precision}(P, P_{\text{gen}}) \times \text{Recall}(P, P_{\text{gen}})}{\text{Precision}(P, P_{\text{gen}}) + \text{Recall}(P, P_{\text{gen}})} \quad (10)$$

535 B.3 Baselines

536 In this section, we introduce the implementation of baseline methods.

- 537 • **Random Perturbation (RP):** We generate $-0.01 \sim 0.01$ random noise to perturb the original
538 trajectory. This degree of noise ensures that the maximum distance between the perturbed points
539 and the original does not exceed 2 km
- 540 • **Gaussian Perturbation (GP):** We generate Gaussian noise perturbed original trajectories with
541 mean 0 and variance 0.01.
- 542 • **Variational AutoEncoder (VAE) [15, 30]:** In this work, trajectories are first embedded as a
543 hidden distribution through two consecutive convolutional layers and a linear layer. Then, we
544 generate the trajectories by a decoder consisting of a linear layer and two deconvolutional layers.
545 The size of convolution kernels in convolutional and deconvolutional layers is set to 4 to ensure
546 that input and output trajectories have the same size.
- 547 • **TrajGAN [9]:** The trajectory is first combined with random noise and then passes through a
548 generator consisting of two linear layers and two convolution layers. Subsequently, a convolutional
549 layer and a linear layer are adopted as the discriminator. The generator and the discriminator are
550 trained in an alternating manner.
- 551 • **Diffwave [13]:** Diffwave is a Wavenet structure model designed for sequence synthesis, which
552 employs extensive dilated convolution. Here, we use 16 residual connected blocks, each consisting
553 of a bi-directional dilated convolution. Then they are summed using sigmoid and tanh activation,
554 respectively, and fed into the 1D CNN.
- 555 • **Diff-scatter:** We randomly sample GPS scatter points from trajectories and generate scatter points
556 using a 4-layer MLP (neurons: $\{128, 256, 256, 128\}$) and the diffusion model.
- 557 • **Diff-wo/UNet:** This model uses only two Resnet blocks combined with a single attention layer
558 between them. Compared with DiffTraj, this model does not have a UNet-type structure, which
559 can be used to evaluate the necessity of the UNet structure.
- 560 • **Diff-LSTM:** This model has the same UNet-type structure and number of Resnet blocks as
561 DiffTraj, and the difference is that Diff-LSTM replaces the CNN with LSTM in Resnet block.
- 562 • **DiffTraj-wo/Con:** DiffTraj-wo/Con represents that the Wide & Deep conditional information
563 embedding module is removed. The rest is the same as DiffTraj.

C Additional Experiments

C.1 Data Utility Setup

In this paper, we use inflow/outflow traffic forecasting to test the utility of the generated data. Inflow/outflow traffic forecasting is a critical task in urban traffic management that involves predicting the volume of vehicles entering (inflow) or leaving (outflow) a specific region within a certain period of time. In this experiment, we divided a city into 16×16 grids, where each grid represents a specific region. The traffic volume entering (inflow) or leaving (outflow) each region within a certain period is predicted. The primary goal of this experiment is to train various prediction models using both original and generated trajectory data, comparing their prediction performance. This evaluation provides an important perspective on the real-world applicability of the data generated by DiffTraj, assessing not just the fidelity of the generated trajectories, but also their utility in downstream tasks. In the experimental setup, we train the prediction models using the generated data and the original data separately, and then test their prediction performance on real data. Advanced neural network models, such as AGCRN, Graph WaveNet, DCRNN, and MTGNN, have been employed to handle this task due to their ability to capture complex spatial-temporal dependencies in multivariate time series data. All the above models and code in this section are followed the publicly available code² provided in the literature [12].

- **AGCRN (Adaptive Graph Convolutional Recurrent Network):** AGCRN is a sophisticated model for spatial-temporal forecasting, which leverages both spatial and temporal features of data. It uses graph convolution to capture spatial dependencies and RNNs to model temporal dynamics, making it capable of handling complex spatial-temporal sequences.
- **GWNet (Graph WaveNet):** GWNet is designed for high-dimensional, structured sequence data. It incorporates a Graph Convolution Network (GCN) to model spatial correlations and a WaveNet-like architecture to model temporal dependencies. The combination allows for capturing both the spatial and temporal complexities present in high-dimensional data.
- **DCRNN (Diffusion Convolutional Recurrent Neural Network):** DCRNN is a deep learning model designed for traffic forecasting, which handles the spatial and temporal dependencies in traffic flow data. It uses a diffusion convolution operation to model spatial dependencies and a sequence-to-sequence architecture with scheduled sampling and residual connections to model temporal dependencies.
- **MTGNN (Multivariate Time-series Graph Neural Network):** MTGNN is a model that captures complex spatial-temporal relationships in multivariate time series data. The model leverages a graph neural network to model spatial dependencies and an auto-regressive process to capture temporal dependencies. It also uses a gating mechanism to adaptively select the relevant spatial-temporal components, thus improving the forecasting performance.

For accuracy comparison, we use the mean square error (MSE), root mean square error (RMSE), and mean absolute error (MAE) as metrics to assess the prediction accuracy of all methods. These three metrics are defined as follows:

$$\text{MSE}(\mathbf{X}, \hat{\mathbf{X}}) = \frac{1}{N} \sum_i \left(X^{(i)} - \hat{X}^{(i)} \right)^2, \quad (11)$$

$$\text{RMSE}(\mathbf{X}, \hat{\mathbf{X}}) = \sqrt{\frac{1}{N} \sum_i \left(X^{(i)} - \hat{X}^{(i)} \right)^2}, \quad (12)$$

$$\text{MAE}(\mathbf{X}, \hat{\mathbf{X}}) = \frac{1}{N} \sum_i \left| X^{(i)} - \hat{X}^{(i)} \right|, \quad (13)$$

where $X^{(i)}$ and $\hat{X}^{(i)}$ are the ground truth and predicted inflow/outflow value at time i , respectively.

²<https://github.com/deepkashiwa20/DL-Traff-Graph>

C.2 Conditional Generation

As we described in Sec. 4.2 and Sec. 4.3, the DiffTraj framework takes into account various external factors that influence real-world trajectories, such as road network structure and departure time. These factors are used to guide the generation process, ensuring the synthetic trajectories mimic real-world patterns and behaviors. The model employs a Wide & Deep network to effectively embed this conditional information, enhancing the capabilities of the Traj-UNet model.

To evaluate the conditional generation capability of DiffTraj, we investigate the case of generated trajectories where the start and end regions of the trajectories were pre-defined. The model was tasked to generate 20 random trajectories adhering to these conditions. The results depicted in Fig. 7 and Fig. 8 (The red and blue boxes indicate the starting and ending regions, respectively), effectively demonstrate proficiency in generating trajectories that align with the specified start and end regions of DiffTraj. This is observed consistently across both cities under study, reinforcing the model’s ability to accurately incorporate and adhere to conditional information. This robust capability underscores the versatility of DiffTraj in generating meaningful trajectories under specific conditions, and its applicability in real-world scenarios where such conditions often exist.

C.3 Ensuring Generation Diversity

In addition, DiffTraj is designed to generate high-quality trajectories and ensure a level of diversity that prevents overly deterministic behavior patterns, thereby upholding intended privacy protections. By integrating a classifier-free diffusion guidance method, DiffTraj can strike a calculated balance between sample quality and diversity. To validate the capacity for generating diverse trajectories of DiffTraj, we devised an experiment that manipulates the guiding parameter, ω . This experiment aimed to examine the quality-diversity balance in trajectories generated by DiffTraj, and how this equilibrium responds to variations in ω . In this experimental setup, we studied the trajectories yielded under different ω settings (specifically $\omega \in 0.1, 1, 10$) while keeping the conditional information the same.

The results, as illustrated in Fig. 7 and Fig. 8, reveal an unambiguous link between an increase in ω and a rise in trajectory diversity. This finding affirms that DiffTraj adeptly manages the balance between diversity and quality. As ω increases, the model demonstrates a tendency to spawn more varied trajectories. This is because a higher ω prompts the model to place more emphasis on unconditional noise prediction and reduce the sway of the conditional information. Thus, the model grows more proficient at creating diverse trajectories, albeit potentially compromising some quality. This greater diversity is a consequence of the model having fewer constraints from specific conditions, providing more latitude to explore a wider range of possible trajectories. This experiment underscores the flexibility and control inherent in DiffTraj in balancing trajectory quality and diversity, vital characteristics for generating realistic and diverse trajectories. Therefore, ω serves as a control knob for modulating the trade-off between trajectory quality and diversity, providing a powerful tool for users to align the generated trajectories with specific application requirements.

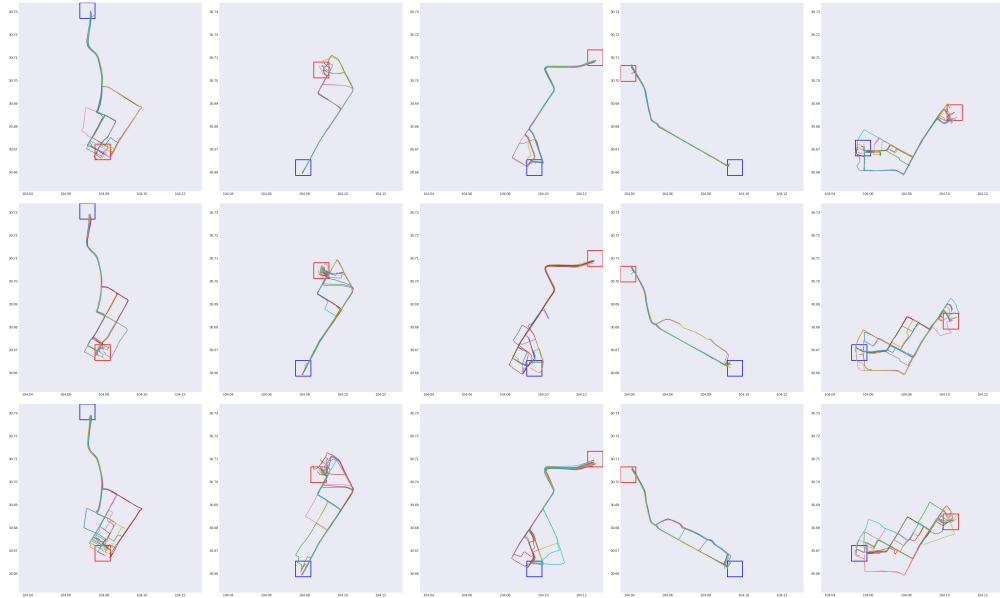


Figure 7: Conditional trajectory generation on Chengdu. The guidance scales ω of the first, second and third rows are 0.1, 1, 10, respectively. The rectangular box indicates the area of the assigned start and end points.

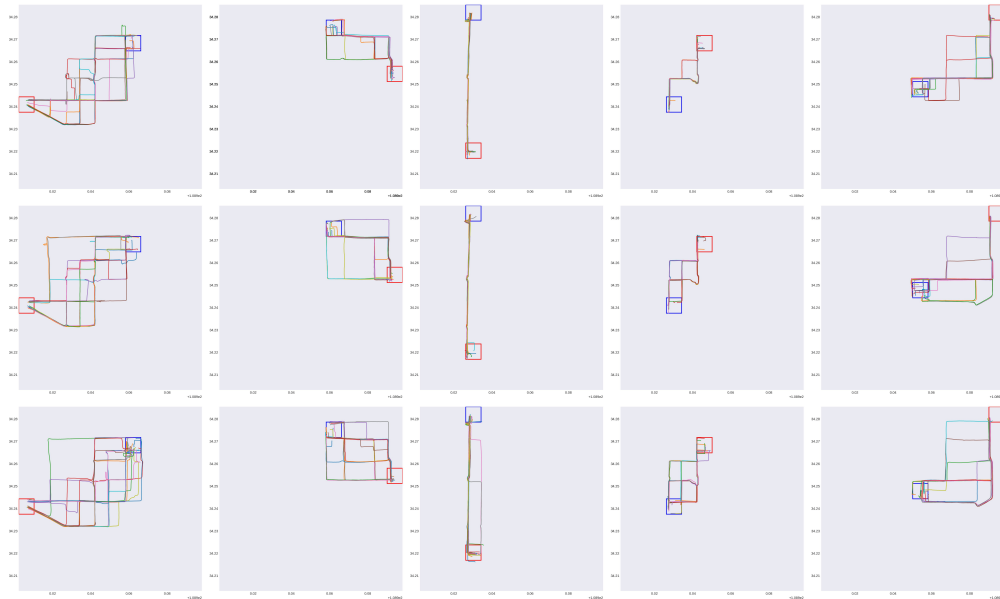


Figure 8: Conditional trajectory generation on Xi'an. The guidance scales ω of the first, second and third rows are 0.1, 1, 10, respectively. The rectangular box indicates the area of the assigned start and end regions.

D Visualization Results

We append a series of experimental results in this section due to space constraints. As shown in Fig. 9, we visualize the heat map of the trajectory distribution with multiple resolutions. Specifically, we divide the whole city into 32×32 , 16×16 , and 8×8 grids, and then count the distribution of trajectories in each grid. The comparison clearly indicates that the distributions are highly consistent from all resolutions. The visualized results also verify the effectiveness of metrics in Table 1, revealing that the proposed model can generate high-quality trajectories with remarkable accuracy and retain the original distribution.

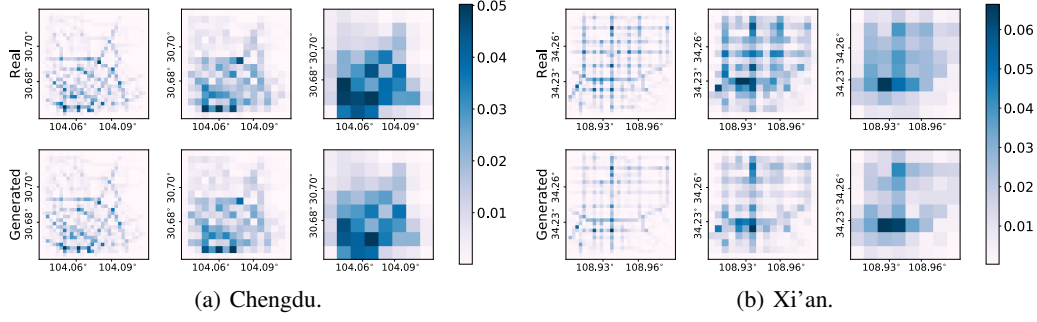


Figure 9: Comparison of the real and generated trajectory distributions with different resolutions. The city is divided into different size grids (32×32 , 16×16 , and 8×8 grids).

In addition, we also show the geographic results of the trajectories generated by different generation methods for two cities, Chengdu and Xi'an. The visualization results are concluded in Fig. 10 and Fig. 11.

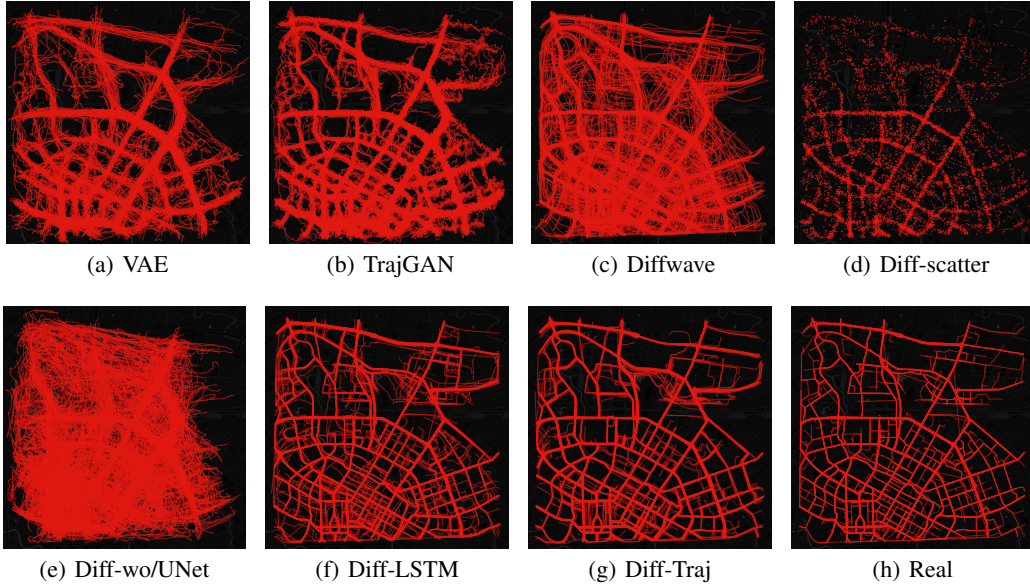


Figure 10: Geographic visualization of generated trajectory in Chengdu.

The rest visualize the forward trajectory addition noise process and reverse trajectory denoising process of Diff-Traj.

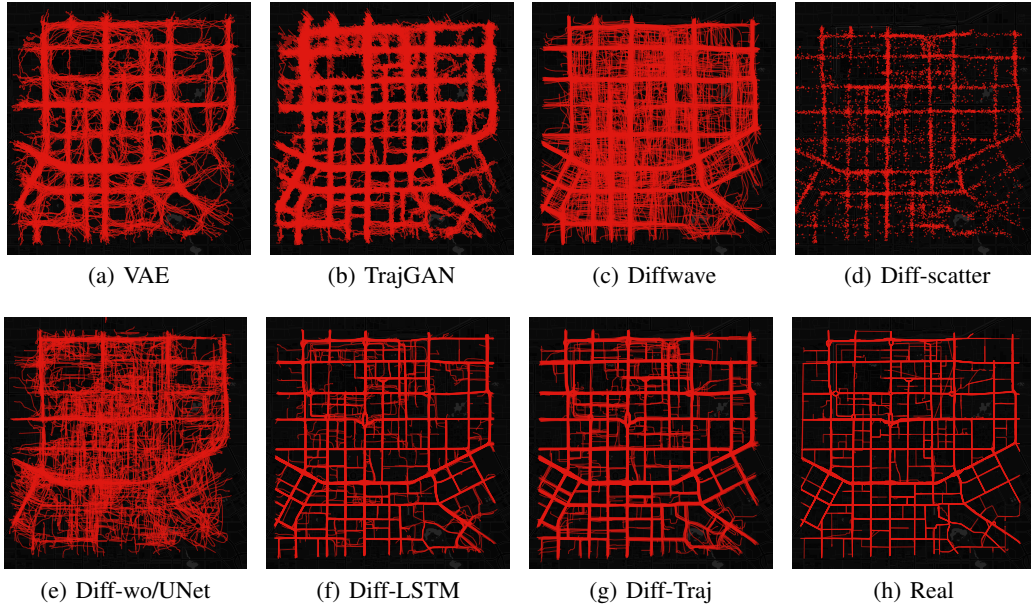


Figure 11: Geographic visualization of generated trajectory in Xi'an.

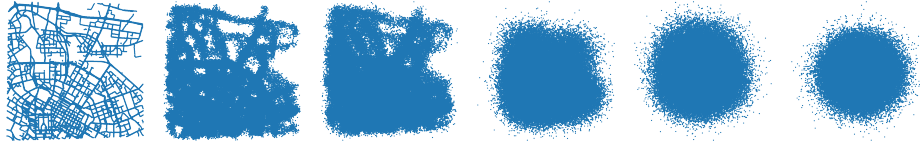


Figure 12: Forward trajectory noising process (Chengdu).

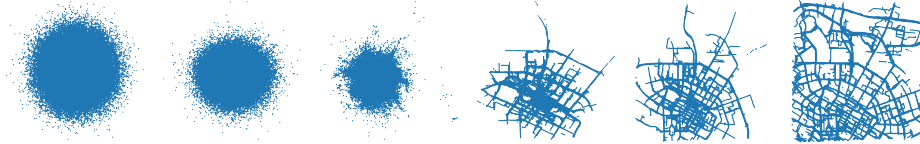


Figure 13: Reverse trajectory denoising process (Chengdu).



Figure 14: Forward trajectory noising process (Xi'an).

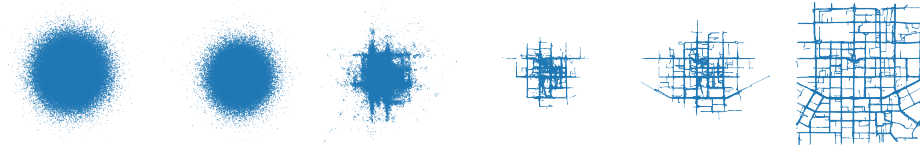


Figure 15: Reverse trajectory denoising process (Xi'an).