# Spatio-temporal Diffusion Point Processes

苑苑

清华大学电子系城市科学与计算研究中心
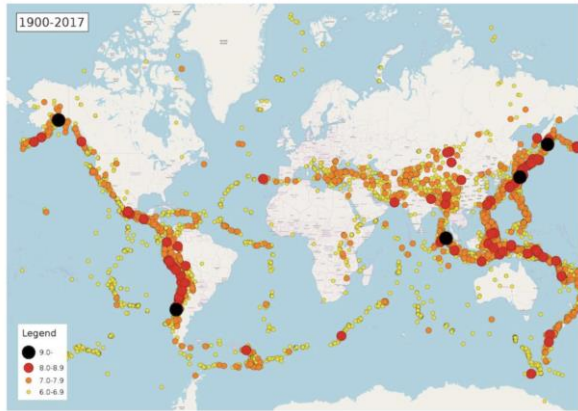
2023/06/16
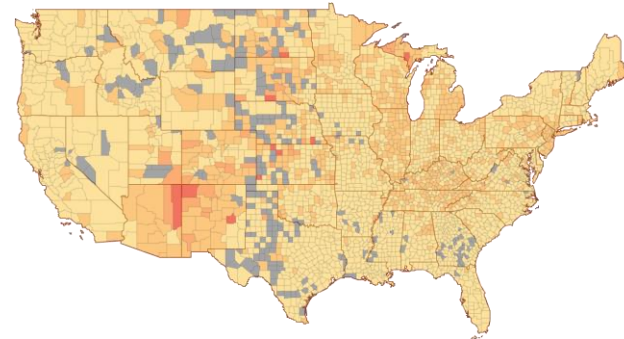
清华大学电子工程系

Department of Electronic Engineering, Tsinghua University

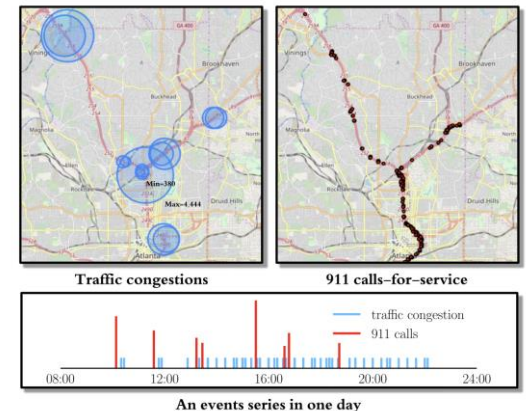# Spatio-temporal Events

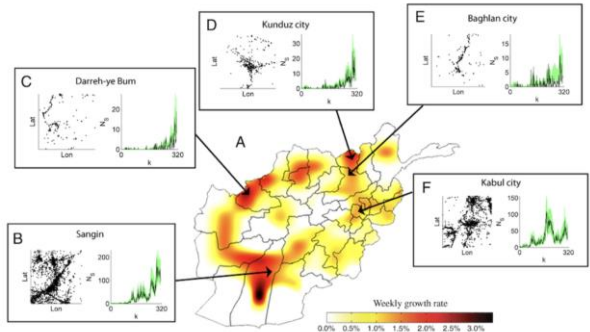## Events naturally come with time and location.


Earthquake


Covid-19


Traffic congestion


Conflict


Brain activity
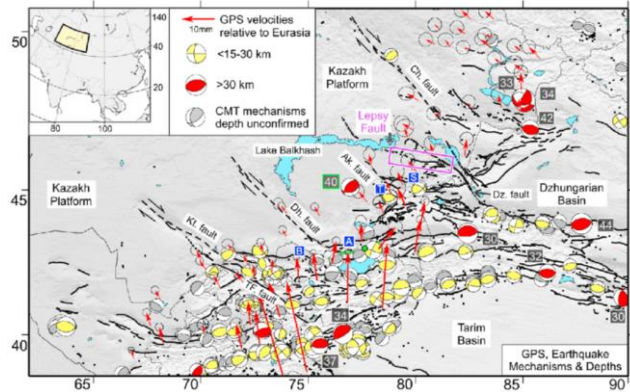

Checkin

时空事件**随机性强，时空关联复杂度高**，但其建模**价值高、意义大**

**Understanding**

Study the mechanisms that give rise to the dynamics of the recurrence of events



**Predicting**

Predict the dynamics of events in the future based on event history



**Controlling**

Design intervention and control measures to steer the dynamics of events to desirable outcomes



时空事件序列建模能有效**预测未来事件**，揭示其内在**时空关联**。

# Event modeling: point processes

**Discrete events localized in continuous time and space are usually modeled by point processes.**

Point process in 1D
(Temporal Point Process)

A point process in time.

Point process in 3D
(Spatio-temporal Point Process)

$t$

$(t_1, x_1)$

$(t_2, x_2)$

$x_{lat}$

$x_{lng}$

A point process in space and time.

- **Conditional Intensity function**

$$\lambda_k^*(t) = \lim_{\Delta t \downarrow 0} \frac{\Pr(\text{event of type } k \text{ in } [t, t+\Delta t) | \mathcal{H}_t)}{\Delta t},$$

- **Density**

$$f^*(t) = \lambda^*(t) e^{-\int_{t_{i-1}}^{t} \lambda^*(\tau) d\tau}$$

- **Likelihood**

$$L = f^*(t_1) \cdot f^*(t_2) \cdot \cdots \cdot f^*(t_n) \cdot S^*(T)$$

$$= \left( \prod_{i=1}^{n} \lambda^*(t_i) \cdot \exp\left(-\int_{t_{i-1}}^{t_i} \lambda^*(s) ds\right) \right) \cdot \exp\left(-\int_{t_n}^{T} \lambda^*(s) ds\right)$$

$$= \left( \prod_{i=1}^{n} \lambda^*(t_i) \right) \cdot \exp\left(-\int_{0}^{T} \lambda^*(s) ds\right)$$

- **A STPP is uniquely defined by its conditional intensity** $\lambda(t, x | H_t)$

  ➢ the rate of event occurrence given the history $H_t$

- **Likelihood**

  ➢ $$\log p\left(\mathcal{H}\right) = \sum_{i=1}^{n} \log \lambda^*(t_i, \boldsymbol{x}_i) - \int_0^T \int_{\mathbb{R}^d} \lambda^*(\tau, \boldsymbol{x}) \, d\boldsymbol{x} d\tau.$$

**Keep in mind when specifying the distribution $f^*$**

  ✓ *Closed-form likelihood* : the likelihood should be computed in closed form for efficient model training.

  ✓ *Closed-form sampling* : Can we draw samples from $f^*$ analytically?

  ✓ *Flexibility* : Does the given parametrization of $f^*$ allow us to approximate any distribution?

**Complex spatial and temporal correlations**

Spatial

Temporal

**The spatial and temporal interaction patterns have high variations in urban environment.**

**Neural Spatio-temporal Point Process,** *NIPS 2020*

$$\lambda^*(t, \boldsymbol{x}) = \lambda^*(t)\, p^*(\boldsymbol{x} \mid t)$$

$$\log p\left(\mathcal{H}\right) = \sum_{i=1}^{n} \log \lambda^*(t_i, \boldsymbol{x}_i) - \int_0^T \int_{\mathbb{R}^d} \lambda^*(\tau, \boldsymbol{x})\, d\boldsymbol{x}d\tau$$

**Simplify**

$$\log p(\mathcal{H}) = \underbrace{\sum_{i=1}^{n} \log \lambda^*(t_i) - \int_0^T \lambda^*(\tau)\, d\tau}_{\text{temporal log-likelihood}} + \underbrace{\sum_{i=1}^{n} \log p^*(\boldsymbol{x}_{t_i}^{(i)} \mid t_i)}_{\text{spatial log-likelihood}}$$

# Neural Spatio-temporal Point Process, *NIPS 2020*

$$\lambda^*(t, \boldsymbol{x}) = \lambda^*(t)\, p^*(\boldsymbol{x} \mid t)$$

$$\boldsymbol{h}_{t_0} = \boldsymbol{h}_0$$
$$\frac{d\boldsymbol{h}_t}{dt} = f_h(t, \boldsymbol{h}_t)$$
$$\lim_{\varepsilon \to 0} \boldsymbol{h}_{t_i + \varepsilon} = g_h\left(t_i, \boldsymbol{h}_{t_i}, \boldsymbol{x}_{t_i}^{(i)}\right)$$

$$\boldsymbol{x}_0 \sim p(\boldsymbol{x}_0)$$
$$\frac{d\boldsymbol{x}_t}{dt} = f_x(t, \boldsymbol{x}_t, \boldsymbol{h}_t)$$
$$\lim_{\varepsilon \to 0} \boldsymbol{x}_{t_i + \varepsilon} = g_x(t_i, \boldsymbol{x}_{t_i}, \boldsymbol{h}_{t_i})$$

ODE Solve$(z_{t_0}, f, \boldsymbol{\theta}_f, t_0, ..., t_M)$

$z_{t_0}$  $z_{t_1}$  $z_{t_N}$  $z_{t_{N+1}}$  $z_{t_M}$

$\hat{x}(t)$

Neural ODE

$t_0$  $t_1$  $t_N$  $t_{N+1}$  $t_M$

Continuous normalizing flow

**Limitation: require assumptions on the spatio-temporal dependence**

- **Normalizing flows**

  ➤ Get flexible distribution by transforming a simple density

  ➤ $p(\tau) = p_0(f_\theta^{-1}(\tau)) \left| \frac{d}{d\tau} f_\theta^{-1}(\tau) \right|$

  ➤ Invertible

- **Mixture distribution**

  ➤ Convex combination of simple densities

  ➤ $p(\tau) = \sum_k \pi_k p_0(\tau|\theta_k)$

*They both require a **certain structure** in the functional approximators.*

**Flexibility** ⟷ **Computation cost**

**Flexibility** ⟷ **Computation cost**

**We expect the model:**

✓ **Update based on new events**

✓ **Easy to train and sample**

✓ **Use powerful neural networks freely**

# From complex to simple

Joint distribution:
**difficult to learn**

Task

… 

Step 1    Step 2    Step K-1    Step K

K subtasks: learn a **simple** distribution

**Complex task**

**Simple subtasks**

16

**Diffusion Model**

**Algorithm 1** Training
1: **repeat**
2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3: $t \sim \text{Uniform}(\{1, \ldots, T\})$
4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5: Take gradient descent step on
$$\nabla_\theta \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) \right\|^2$$
6: **until** converged

- **DDPM** [1]

Use variational lower bound



$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is unknown

Diffusion process: $q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$

Denoising process: $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$

Training loss: $L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon}\left[ \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) \right\|^2 \right]$

[1] Ho, Jonathan, Ajay Jain, and Pieter Abbeel. "Denoising diffusion probabilistic models." NIPS 2020.

# Why diffusion model is feasible?

## Complex function

## Complex distribution



**Similar idea**

**Complex curves can be sufficiently fitted with piecewise linear function.**

**Minor changes can be sufficiently modeled by Gaussian distributions.**

# Spatio-temporal Diffusion Point Processes



**Spatio-temporal Self-att Encoder**
**Learn history representation**

**Condition**

**Spatio-temporal Diffusion Model**
**Learn ST joint distribution**

- **Diffusion**

  - Separately diffuse on the spatial and temporal domains

- **Denoise**

  - Conditionally independent at the current step

  - Dependent on the predicted values from the last step

$$q_{st}(\boldsymbol{x}_i^k|\boldsymbol{x}_i^{k-1}) := (q(\tau_i^k|\tau_i^{k-1}), q(s_i^k|s_i^{k-1})) \ ,$$
$$q(x^k|x^{k-1}) := \mathcal{N}(x^k; \sqrt{1-\beta_k}x^k, \beta_k\boldsymbol{I}) \ ,$$

$$p_\theta(x_i^{k-1}|x_i^k, h_{i-1}) = p_\theta(\tau_i^{k-1}|\tau_i^k, s_i^k, h_{i-1})p_\theta(s_i^{k-1}|\tau_i^k, s_i^k, h_{i-1})$$

Current step: independent

$\tau^{k-1}$     $s^{k-1}$

Last step $\tau^k, s^k$

$\tau^{k-1}$     $s^{k-1}$

$\tau^k$     $s^k$

Interstep: interaction

# Spatio-temporal Diffusion Point Processes

- **Parametrization of the denoising network**



Spatio-temporal
Co-attention

# Spatio-temporal Diffusion Point Processes

- **Comparison on important properties**

| Model | No Asmp.[1] | No Restr.[2] | Flexible[3] | Closed-form sampling[4] |
|---|---|---|---|---|
| Hawkes [17] | ✗ | ✗ | ✗ | ✗ |
| Self-correcting [21] | ✗ | ✗ | ✗ | ✗ |
| KDE [2] | - | - | ✗ | ✓ |
| CNF [3] | - | - | ✗ | ✓ |
| ST Hawkes [45] | ✗ | ✗ | ✗ | ✗ |
| RMTPP [9] | ✗ | ✗ | ✓ | ✗ |
| NHP [32] | ✗ | ✗ | ✓ | ✗ |
| THP [69] | ✗ | ✗ | ✓ | ✗ |
| SNAP [64] | ✗ | ✗ | ✓ | ✗ |
| LogNormMix [47] | ✗ | ✗ | ✗ | ✓ |
| NJSDE [22] | ✗ | ✗ | ✓ | ✗ |
| Neural STPP [3] | ✓ | ✗ | ✓ | ✗ |
| DeepSTPP [66] | ✗ | ✗ | ✓ | ✗ |
| DSTPP (ours) | ✓ | ✓ | ✓ | ✓ |

# Datasets and Metrics

- **Datasets**

  - ***Continuous-space***

    - Real-world: Earthquake, Citybikes, Covid19, Human mobility

    - Synthetic: Pinewheel, HawkesGMM

  - ***Discrete-space***

    - Taxi-NYC

    - Crime-Atlanta

- **Evaluation metrics**

  - ***Log-likelihood***

  - ***Prediction error***

    - Temporal: MAE, RMSE

    - Spatial: Euclidean distance, Accuracy

# Overall performance

- **Likelihood**

Table 2: Performance evaluation for negative log-likelihood per event on test data. ↓ means lower is better. Bold denotes the best results and underline denotes the second-best results.

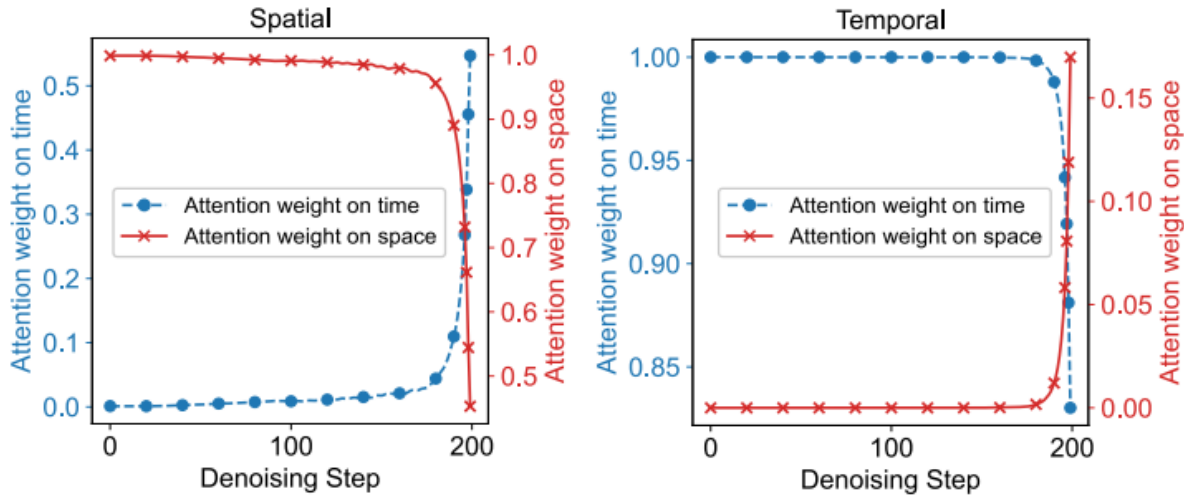| Model | Earthquake Spatial ↓ | Earthquake Temporal ↓ | COVID-19 Spatial ↓ | COVID-19 Temporal ↓ | Citibike Spatial ↓ | Citibike Temporal ↓ | HawkesGMM Spatial ↓ | HawkesGMM Temporal ↓ |
|---|---|---|---|---|---|---|---|---|
| Conditional KDE | $2.21_{\pm 0.105}$ | $-^{(1)}$ | $2.31_{\pm 0.084}$ | - | $2.74_{\pm 0.001}$ | - | $\underline{0.236_{\pm 0.001}}$ | - |
| CNF | $1.35_{\pm 0.000}$ | - | $2.05_{\pm 0.014}$ | - | $2.15_{\pm 0.000}$ | - | $0.427_{\pm 0.002}$ | - |
| TVCNF | $1.34_{\pm 0.008}$ | - | $2.04_{\pm 0.004}$ | - | $2.19_{\pm 0.025}$ | - | $0.431_{\pm 0.008}$ | - |
| Possion | - | $-0.146_{\pm 0.000}$ | - | $-0.876_{\pm 0.021}$ | - | $-0.626_{\pm 0.000}$ | - | $1.34_{\pm 0.000}$ |
| Hawkes | - | $-0.514_{\pm 0.000}$ | - | $-2.06_{\pm 0.000}$ | - | $-1.06_{\pm 0.001}$ | - | $0.880_{\pm 0.000}$ |
| Self-correcting | - | $13.8_{\pm 0.533}$ | - | $7.13_{\pm 0.062}$ | - | $7.11_{\pm 0.010}$ | - | $4.59_{\pm 0.135}$ |
| RMTPP | - | $0.0930_{\pm 0.051}$ | - | $-1.30_{\pm 0.022}$ | - | $1.24_{\pm 0.001}$ | - | $1.52_{\pm 0.002}$ |
| NHP | - | $-0.676_{\pm 0.001}$ | - | $-2.30_{\pm 0.001}$ | - | $-1.14_{\pm 0.001}$ | - | $0.580_{\pm 0.000}$ |
| THP | - | $\underline{-0.976_{\pm 0.011}}$ | - | $-2.12_{\pm 0.002}$ | - | $\underline{-1.49_{\pm 0.003}}$ | - | $-0.402_{\pm 0.001}$ |
| SAHP | - | $-0.229_{\pm 0.007}$ | - | $-1.37_{\pm 0.118}$ | - | $-1.02_{\pm 0.067}$ | - | $\underline{-1.25_{\pm 0.136}}$ |
| LogNormMix | - | $-0.341_{\pm 0.071}$ | - | $-2.01_{\pm 0.025}$ | - | $-1.06_{\pm 0.005}$ | - | $0.630_{\pm 0.004}$ |
| NJSDE | $1.65_{\pm 0.012}$ | $0.0950_{\pm 0.203}$ | $2.21_{\pm 0.005}$ | $-1.82_{\pm 0.002}$ | $2.63_{\pm 0.001}$ | $-0.804_{\pm 0.059}$ | $0.395_{\pm 0.001}$ | $1.77_{\pm 0.030}$ |
| NSTPP | $\underline{0.885_{\pm 0.037}}$ | $-0.623_{\pm 0.004}$ | $1.90_{\pm 0.017}$ | $\underline{-2.25_{\pm 0.002}}$ | $2.38_{\pm 0.053}$ | $-1.09_{\pm 0.004}$ | $0.285_{\pm 0.011}$ | $0.824_{\pm 0.005}$ |
| DeepSTPP | $4.92_{\pm 0.007}$ | $-0.174_{\pm 0.001}$ | $\underline{0.361_{\pm 0.01}}$ | $-1.09_{\pm 0.01}$ | $\underline{-4.94_{\pm 0.016}}$ | $-1.13_{\pm 0.002}$ | $0.519_{\pm 0.001}$ | $0.322_{\pm 0.002}$ |
| DSTPP (ours) | $\mathbf{0.308_{\pm 0.006}}$ | $\mathbf{-1.96_{\pm 0.020}}$ | $\mathbf{-1.02_{\pm 0.029}}$ | $\mathbf{-3.08_{\pm 0.003}}$ | $\mathbf{-5.41_{\pm 0.011}}$ | $\mathbf{-3.36_{\pm 0.010}}$ | $\mathbf{-2.95_{\pm 0.047}}$ | $\mathbf{-3.42_{\pm 0.002}}$ |

[1] Spatial baselines and temporal baselines can be combined freely for modeling spatio-temporal domains.
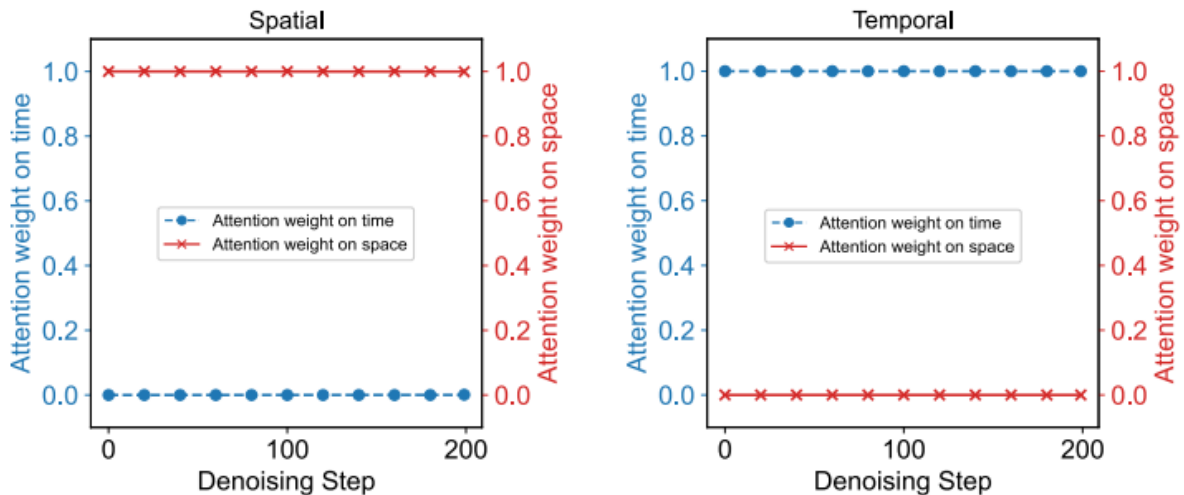
- **Event prediction**

Table 3: Performance evaluation for predicting both time and space of the next event. We use Euclidean distance to measure the prediction error of the spatial domain and use RMSE between real intervals and predicted intervals for time prediction.

| Model | Earthquake | | COVID-19 | | Citibike | | HawkesGMM | |
|---|---|---|---|---|---|---|---|---|
| | Spatial ↓ | Temporal ↓ | Spatial ↓ | Temporal ↓ | Spatial ↓ | Temporal ↓ | Spatial ↓ | Temporal ↓ |
| Conditional KDE | $11.3_{\pm 0.658}$ | – | $0.688_{\pm 0.047}$ | – | $0.718_{\pm 0.001}$ | – | $1.54_{\pm 0.006}$ | – |
| CNF | $8.48_{\pm 0.054}$ | – | $0.559_{\pm 0.000}$ | – | $0.722_{\pm 0.000}$ | – | $71663_{\pm 60516}$ | – |
| TVCNF | $8.11_{\pm 0.001}$ | – | $0.560_{\pm 0.000}$ | – | $0.705_{\pm 0.000}$ | – | $2.03_{\pm 0.000}$ | – |
| Possion | – | $0.631_{\pm 0.017}$ | – | $0.463_{\pm 0.021}$ | – | $0.438_{\pm 0.001}$ | – | $2.81_{\pm 0.070}$ |
| Hawkes | – | $0.544_{\pm 0.010}$ | – | $0.672_{\pm 0.088}$ | – | $0.534_{\pm 0.011}$ | – | $2.63_{\pm 0.002}$ |
| Self-correcting | – | $11.2_{\pm 0.486}$ | – | $2.83_{\pm 0.141}$ | – | $10.7_{\pm 0.169}$ | – | $9.72_{\pm 0.159}$ |
| RMTPP | – | $0.424_{\pm 0.009}$ | – | $1.32_{\pm 0.024}$ | – | $2.07_{\pm 0.015}$ | – | $3.38_{\pm 0.012}$ |
| NHP | – | $1.86_{\pm 0.023}$ | – | $2.13_{\pm 0.100}$ | – | $2.36_{\pm 0.056}$ | – | $2.82_{\pm 0.028}$ |
| THP | – | $2.44_{\pm 0.021}$ | – | $0.611_{\pm 0.008}$ | – | $1.46_{\pm 0.009}$ | – | $5.35_{\pm 0.002}$ |
| SAHP | – | $0.409_{\pm 0.002}$ | – | $0.184_{\pm 0.024}$ | – | $\underline{0.203_{\pm 0.010}}$ | – | $2.75_{\pm 0.049}$ |
| LogNormMix | – | $0.593_{\pm 0.005}$ | – | $0.168_{\pm 0.011}$ | – | $0.350_{\pm 0.013}$ | – | $2.79_{\pm 0.021}$ |
| WGAN | – | $0.481_{\pm 0.007}$ | – | $\underline{0.124_{\pm 0.002}}$ | – | $0.238_{\pm 0.003}$ | – | $2.83_{\pm 0.048}$ |
| NJSDE | $9.98_{\pm 0.024}$ | $0.465_{\pm 0.009}$ | $0.641_{\pm 0.009}$ | $0.137_{\pm 0.001}$ | $0.707_{\pm 0.001}$ | $0.264_{\pm 0.005}$ | $1.62_{\pm 0.003}$ | $2.25_{\pm 0.007}$ |
| NSTPP | $8.11_{\pm 0.000}$ | $0.547_{\pm 0.010}$ | $0.560_{\pm 0.000}$ | $0.145_{\pm 0.002}$ | $0.705_{\pm 0.000}$ | $0.355_{\pm 0.013}$ | $2.02_{\pm 0.000}$ | $3.30_{\pm 0.201}$ |
| DeepSTPP | $\underline{6.51_{\pm 0.000}}$ | $\underline{0.341_{\pm 0.000}}$ | $\underline{0.486_{\pm 0.000}}$ | $0.197_{\pm 0.000}$ | $\underline{0.0312_{\pm 0.000}}$ | $\underline{0.234_{\pm 0.000}}$ | $\underline{1.38_{\pm 0.000}}$ | $\underline{1.46_{\pm 0.000}}$ |
| DSTPP (ours) | $\mathbf{2.84_{\pm 0.193}}$ | $\mathbf{0.149_{\pm 0.001}}$ | $\mathbf{0.170_{\pm 0.001}}$ | $\mathbf{0.0243_{\pm 0.000}}$ | $\mathbf{0.00495_{\pm 0.000}}$ | $\mathbf{0.0301_{\pm 0.002}}$ | $\mathbf{0.136_{\pm 0.013}}$ | $\mathbf{0.0891_{\pm 0.009}}$ |

- **Comparison on important properties**



(a) Earthquakes

(b) Synthetic-Independent

*Spatio-temporal* *dependent*

**Different spatio-temporal relationships**
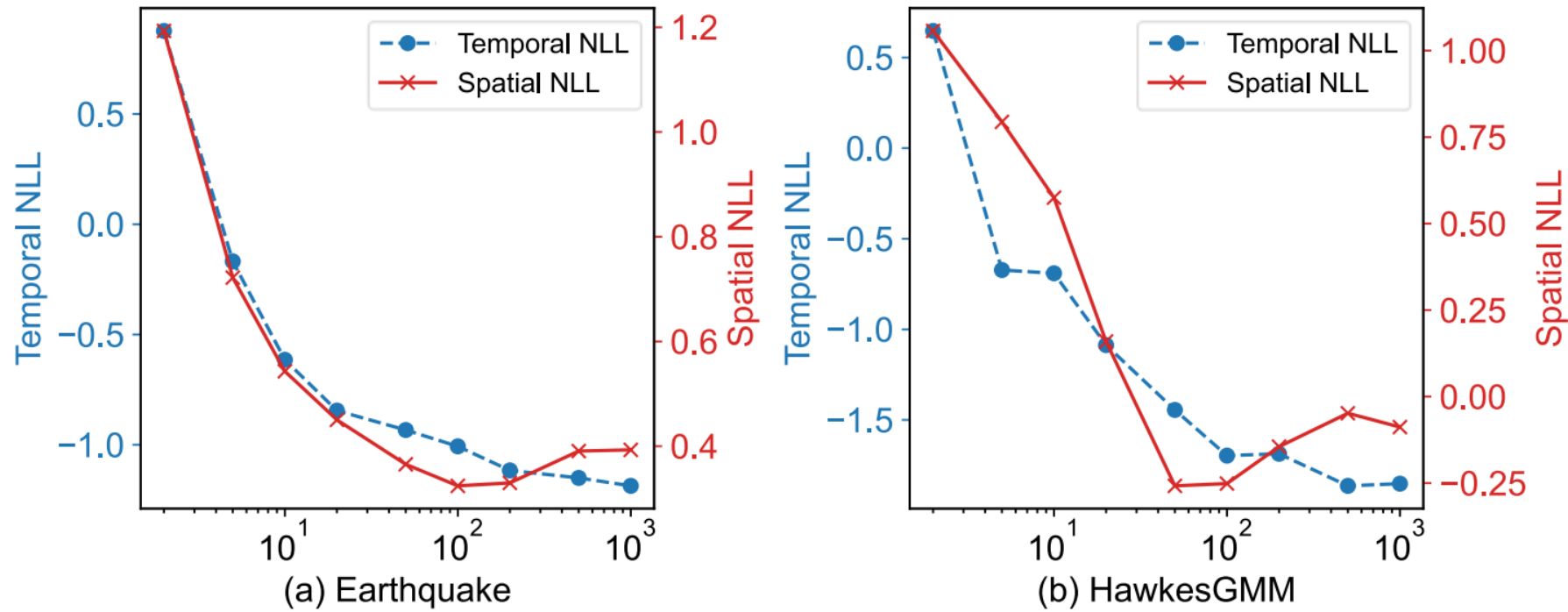
*Spatio-temporal* *independent*

26

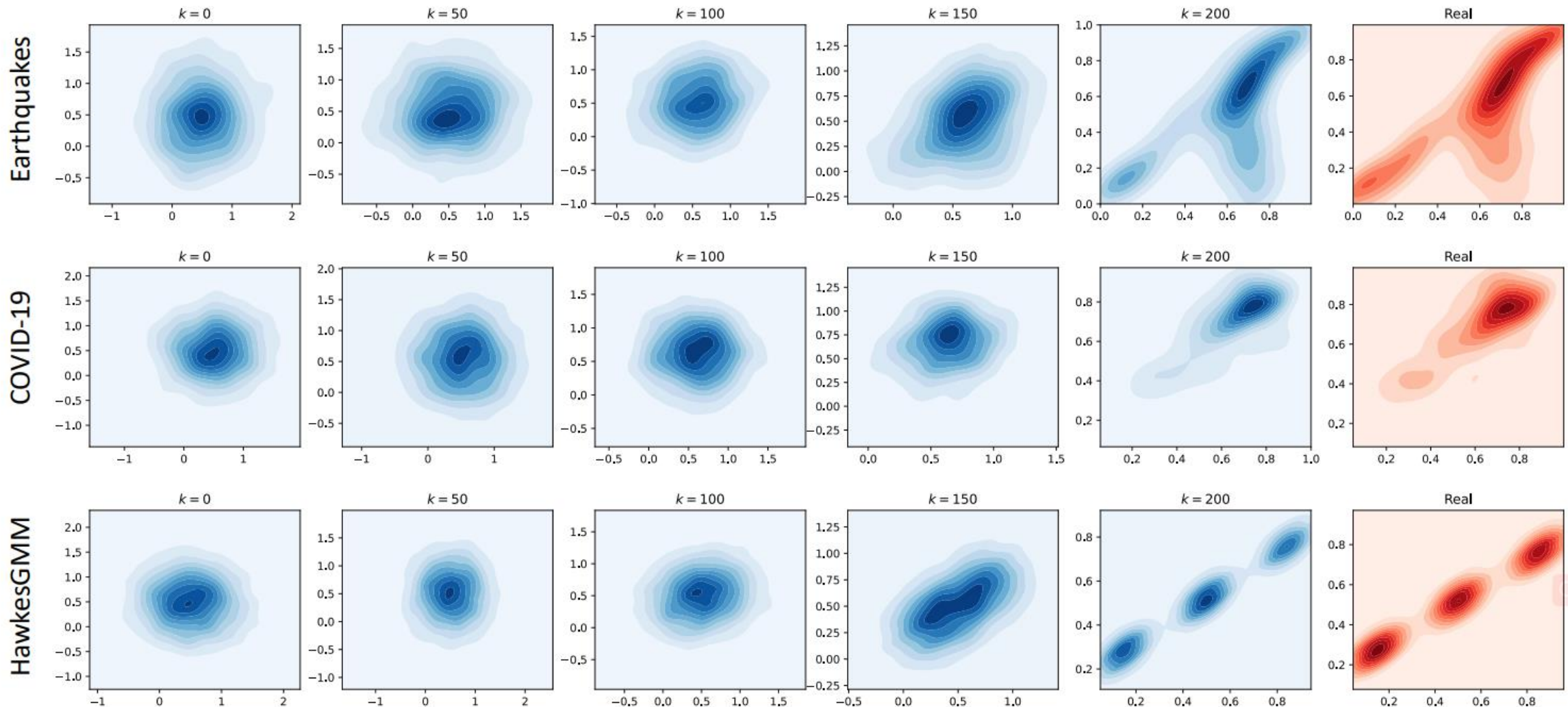# Impact of diffusion steps

- **The number of total steps K is a crucial hyperparameter**

With the increase of diffusion steps, the denoising network approximates more minimal changes between steps.



(a) Earthquake      (b) HawkesGMM

**Accuracy & Efficiency**

# Visualization of denoising processes

- **DSTPP is able to learn the generative process of distributions**

# Thanks for listening !

https://github.com/tsinghua-fib-lab/Spatio-temporal-Diffusion-Point-Processes

Lab Info: http://fi.ee.tsinghua.edu.cn

Contact: liyong07@tsinghua.edu.cn