

# EDIT-BASED FLOW MATCHING FOR TEMPORAL POINT PROCESSES

Anonymous authors

Paper under double-blind review

## ABSTRACT

Temporal point processes (TPPs) are a fundamental tool for modeling event sequences in continuous time, but most existing approaches rely on autoregressive parameterizations that are limited by their sequential sampling. Recent non-autoregressive, diffusion-style models mitigate these issues by jointly interpolating between noise and data through event insertions and deletions in a discrete Markov chain. In this work, we generalize this perspective and introduce an Edit Flow process for TPPs that transports noise to data via insert, delete, and substitute edit operations. By learning the instantaneous edit rates within a continuous-time Markov chain framework, we attain a flexible and efficient model that effectively reduces the total number of necessary edit operations during generation. Empirical results demonstrate the generative flexibility of our unconditionally trained model in a wide range of unconditional and conditional generation tasks on benchmark TPPs.

## 1 INTRODUCTION

Temporal point processes (TPPs) capture the distribution over sequences of events in time, where both the continuous arrival-times and number of events are random. They are widely used in domains such as finance, healthcare, social networks, and transportation, where understanding and forecasting event dynamics and their complex interactions is crucial. Most (neural) TPPs capture the complex interactions between events *autoregressively*, parameterizing a conditional intensity/density of each event given its history (Daley & Vere-Jones, 2006; Shchur et al., 2021). While natural and flexible, this factorization comes with inherent limitations: sampling scales linearly with sequence length, errors can compound in multi-step generation, and conditional generation is restricted to forecasting tasks.

**Beyond autoregression.** Recent advances demonstrate that modeling event sequences *jointly* proposes a sound alternative to overcome these limitations. Inspired by diffusion, ADDTHIN (Lüdtke et al., 2023) and PSDIFF (Lüdtke et al., 2025) leverage the thinning and superposition properties of TPPs to construct a discrete Markov chain that learns to transform noise sequences  $t_0 \sim p_{\text{noise}}(t)$  into data sequences  $t_1 \sim q_{\text{target}}(t)$  through *insertions* and *deletions* of events. These methods highlight the promise of joint sequence modeling for TPPs by learning stochastic set interpolations and have shown state-of-the-art results, especially in forecasting.

In parallel, Havasi et al. (2025) introduced Edit Flow, a discrete flow-matching framework (Gat et al., 2024; Campbell et al., 2024; Shi et al., 2025) for variable-length sequences of tokens (e.g., language). Their approach models discrete flows in sequence space through *insertions*, *deletions*, and *substitutions*, formalized as a continuous-time Markov Chain (CTMC). To make the learning process tractable, they introduce an expanded auxiliary state space that aligns sequences, simultaneously

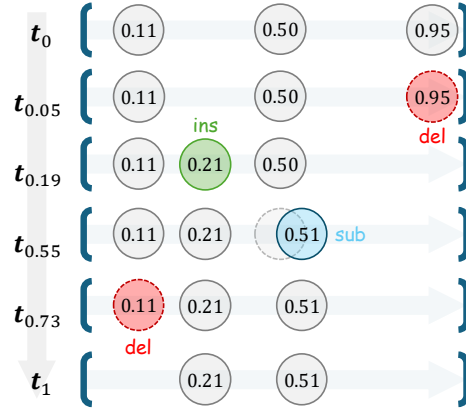


Figure 1: Edit process transporting  $t_0 \sim p_{\text{noise}}(t)$  to  $t_1 \sim q_{\text{target}}(t)$  by inserting, deleting and substituting events.

reducing the complexity of marginalizing over possible transitions and enabling efficient element-wise parameterization in sequence space.

In this paper, we unify these perspectives and propose EDITPP, an Edit Flow for TPPs that learns to transport noise sequences  $\mathbf{t}_0 \sim p_{\text{noise}}(\mathbf{t})$  to data sequences  $\mathbf{t}_1 \sim q_{\text{target}}(\mathbf{t})$  via atomic *edit operations* insertions, deletions, and substitutions (see figure 1). We define these operations specifically for TPPs, efficiently parameterize their instantaneous rates within a CTMC, propose an auxiliary alignment space for TPPs, and show that our unconditionally trained model can be flexibly applied to both unconditional and conditional tasks with adaptive complexity. Our main contributions are:

- We introduce EDITPP, the first generative framework that models TPPs via continuous-time edit operations, unifying stochastic set interpolation methods for TPPs with Edit Flows for discrete sequences.
- We propose a tractable parameterization of insertion, deletion, and substitution rates for TPPs within the CTMC framework, effectively reducing the number of edit operations for generation.
- We demonstrate empirically that EDITPP achieves state-of-the-art results in both unconditional and conditional tasks across diverse real-world and synthetic datasets.

## 2 BACKGROUND

### 2.1 TEMPORAL POINT PROCESSES

TPPs (Daley & Vere-Jones, 2006; 2007) are stochastic processes whose realizations are **almost surely** finite, ordered sets of random events in time. Let  $\mathbf{t} = \{t^{(i)}\}_{i=1}^n$ , with  $t^{(i)} \in [0, T]$ , denote a realization of  $n$  events on a bounded time interval, which can equivalently be represented by the *counting process*  $N(t) = \sum_{i=1}^n \mathbf{1}\{t^{(i)} \leq t\}$  counting the number of events up to time  $t$ . A TPP is uniquely characterized by its *conditional intensity function* (Rasmussen, 2018):

$$\lambda^*(t) = \lim_{\Delta t \downarrow 0} \frac{\mathbb{E}[N(t + \Delta t) - N(t) \mid \mathcal{H}_t]}{\Delta t}, \quad (1)$$

where  $\mathcal{H}_t = \{t^{(i)} : t^{(i)} < t\}$  denotes the history up to time  $t$ . Intuitively,  $\lambda^*(t)$  represents the instantaneous rate of events given the past. Two important properties of TPPs are superposition and thinning. Superposition, i.e., *inserting* one sequence into another,  $\mathbf{t} = \mathbf{t}_1 \cup \mathbf{t}_2$ , where  $\mathbf{t}_1$  and  $\mathbf{t}_2$  are realizations from TPPs with intensities  $\lambda_1$  and  $\lambda_2$ , results in a sample from a TPP with intensity  $\lambda = \lambda_1 + \lambda_2$ . Independent thinning, i.e., randomly *deleting* any event of a sequence from a TPP with intensity  $\lambda$  with probability  $p$ , results in an event sequence from a TPP with intensity  $(1 - p)\lambda$ .

The likelihood of observing an event sequence  $\mathbf{t}$  given the conditional intensity/density is:

$$p(\mathbf{t}) = \left( \prod_{i=1}^n p(t^{(i)} \mid \mathcal{H}_{t^{(i)}}) \right) (1 - F(T \mid \mathcal{H}_{\mathbf{t}})) = \left( \prod_{i=1}^n \lambda^*(t^{(i)}) \right) \exp \left( - \int_0^T \lambda^*(s) ds \right), \quad (2)$$

where  $F(T \mid \mathcal{H}_t)$  is the CDF of the conditional event density  $p(t \mid \mathcal{H}_t)$ . While this autoregressive formulation of TPPs provides a natural framework for modeling event dependencies, it also poses challenges. Parameterizing the conditional intensity or density is generally nontrivial, and the inherently sequential factorization can lead to inefficient sampling, error accumulation, and limits conditional tasks to forecasting (Lüdtke et al., 2023; 2025).

### 2.2 MODELING TPPS BY SET INTERPOLATION

Instead of explicitly modeling the intensity function, Lüdtke et al. (2023; 2025) leverage the thinning and superposition properties of TPPs to derive diffusion-like generative models that interpolate between data event sequences  $\mathbf{t}_1 \sim q_{\text{target}}(\mathbf{t})$  and noise  $\mathbf{t}_0 \sim p_{\text{noise}}(\mathbf{t})$  by inserting and deleting elements. ADDTHIN (Lüdtke et al., 2023) defines the noising Markov chain recursively over a fixed number of steps with size  $\Delta$  indexed by  $s \in [0, 1]$  as follows:

$$\lambda_s(t) = \underbrace{\alpha_s \lambda_{s-\Delta}(t)}_{\text{(i) Thin}} + \underbrace{(1 - \alpha_s) \lambda_0(t)}_{\text{(ii) Add}}, \quad (3)$$

where  $\lambda_1(t)$  is the unknown target intensity of the TPP and  $\alpha_s \in (0, 1)$ . Intuitively, this noising process increasingly deletes events from the data sequence, while inserting events from a noise TPP  $\lambda_0(t)$ . PSDIFF (Lüdtke et al., 2025) further separates the adding and thinning to yield a Markov chain for the forward process, that stochastically interpolates between  $t_0$  and  $t_1$  as follows:

$$p_s(t | t_1, t_0) = \prod_{t \in t} \begin{cases} \bar{\alpha}_s & \text{if } t \in t_1 \\ 1 - \bar{\alpha}_s & \text{if } t \in t_0 \end{cases} \quad (4)$$

or equivalently  $\lambda_s(t) = \bar{\alpha}_s \lambda_1(t) + (1 - \bar{\alpha}_s) \lambda_0(t)$ , with  $\bar{\alpha}_s$  being the product of  $\alpha_i$ 's. Eq. (4) defines an element-wise conditional path by independent insert and delete operations on TPPs, assuming  $t_0 \cap t_1 = \emptyset$ .

### 2.3 FLOW MATCHING WITH EDIT OPERATIONS

Havasi et al. (2025) introduce Edit Flows, a non-autoregressive generative framework for variable-length token sequences with a fixed, discrete vocabulary (e.g., language). They propose a discrete flow that transports a noisy sequence  $x_0 \sim p_{\text{noise}}(x)$  to a data sequence  $x_1 \sim q_{\text{data}}(x)$  via elementary *edit operations*: insertions, deletions, and substitutions. This is formalized via the discrete flow matching framework (Gat et al., 2024; Campbell et al., 2024; Shi et al., 2025) in an augmented space, yielding a CTMC  $\Pr(X_{s+h} = x | X_s = x_s) = \delta_{x_s}(x) + h u_s^\theta(x | x_s) + o(h)$  with transition rates  $u_s^\theta$  governed by the edit operations.

Directly defining a conditional rate  $u_s(x | x_1, x_0)$  to match  $u_s^\theta$  to, as in discrete flow matching, is very hard or even intractable, since all possible edits producing  $x$  must be considered. Thus, to train this CTMC, they rely on two major insights. First, a CTMC in a data space  $\mathcal{X}$  can be learned by introducing an augmented space  $\mathcal{X} \times \mathcal{Z}$  where the true dynamics are known. Second, designing the auxiliary space  $\mathcal{Z}$  to follow the element wise mixture probability path  $p_s(z | z_0, z_1) = \prod_n [(1 - \kappa_s) \delta_{z_0^{(i)}}(z^{(i)}) + \kappa_s \delta_{z_1^{(i)}}(z^{(i)})]$  with kappa schedule  $\kappa_s \in [0, 1]$  (Gat et al., 2024) enables training the CTMC directly in the data space  $\mathcal{X}$  of variable-length sequences.

Edit operations are encoded by **introducing a blank token  $\epsilon$  and mapping  $(x_0, x_1)$  into aligned sequences  $(z_0, z_1)$  in  $\mathcal{Z}$** , where pairs  $(z_0^{(i)}, z_1^{(i)})$  correspond to insertions  $(\epsilon, x)$ , deletions  $(x, \epsilon)$ , or substitutions  $(x, y)$ . Crucially, since the discrete flow matching dynamics in  $\mathcal{Z}$  are known, they can be transferred back to  $\mathcal{X}$  via  $p_s(x, z | z_0, z_1) = p_s(z | z_0, z_1) \delta_{\text{f\_rm-blanks}(z)}(x)$ , **by removing  $\epsilon$ 's with  $\text{f\_rm-blanks}$** . Then, the marginal rates  $u_s^\theta$  are learned in  $\mathcal{X}$  by marginalizing over  $z$  with the Bregman divergence

$$\mathcal{L} = \mathbb{E}_{\substack{(z_0, z_1) \sim \pi(z_0, z_1) \\ s, p_s(z_s, x_s | z_0, z_1)}} \left[ \sum_{x \neq x_s} u_s^\theta(x | x_s) - \sum_{z_s^{(i)} \neq z_1^{(i)}} \frac{\dot{\kappa}_s}{1 - \kappa_s} \log u_s^\theta(x(z_s, i, z_1^{(i)}) | x_s) \right], \quad (5)$$

where  $x(z_s, i, z_1^{(i)}) = \text{f\_rm-blanks}((z_s^{(1)}, \dots, z_s^{(i-1)}, z_1^{(i)}, z_s^{(i+1)}, \dots, z_s^{(n)}))$ .

## 3 METHOD

We introduce EDITPP, an Edit Flow process for TPPs that directly learns the joint distribution of event times. Our process leverages the three elementary edit operations *insert*, *substitute*, and *delete* to define a CTMC that continuously interpolates between two event sequences  $t_0 \sim p_{\text{noise}}(t)$  and  $t_1 \sim q_{\text{data}}(t)$ .

Let  $\mathcal{T} = [0, T]$  denote the support of the TPP. We define the state space as  $\mathcal{X}_{\mathcal{T}} = \bigcup_{0 \leq n < \infty} \{(0, t^{(1)}, \dots, t^{(n)}, T) : 0 < t^{(1)} < \dots < t^{(n)} < T\}$ , denoting the set of all possible *padded* TPP sequences with finitely many events. **Note that the padding values are introduced for notational simplicity when defining the edit operations on  $\mathcal{T}$ .**

### 3.1 EDIT OPERATIONS

Our model navigates the state space  $\mathcal{X}_{\mathcal{T}}$  through a set of atomic edit operations. While Edit Flow was originally defined for discrete state spaces, we can generalize the method to continuous state spaces provided that the set of edit operations remains discrete. We achieve this by defining a finite

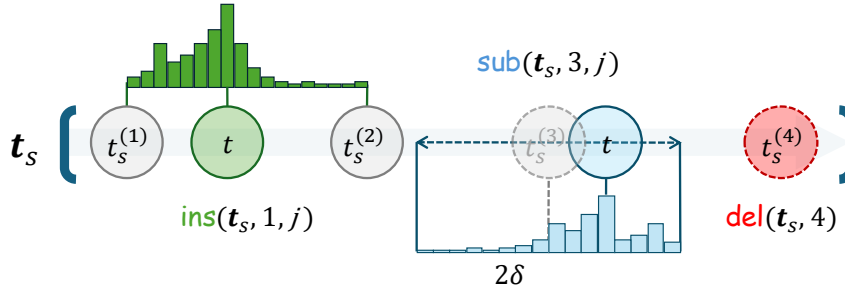


Figure 2: Our discrete edit operations transform continuous event sequences through insertions, substitutions and deletion.

set of edit operations on our continuous state space  $\mathcal{X}_{\mathcal{T}}$  that nonetheless allow us to transition from any sequence  $\mathbf{t}$  to any other  $\mathbf{t}'$  through repeated application.

Similar to Havasi et al. (2025), we design our operations to be mutually exclusive: if two sequences differ by exactly one edit, the responsible operation is uniquely determined. This simplifies the parameterization of the model and computation of the Bregman divergence in Eq. (5).

**Insertion:** To discretize the event insertion, we quantize the space between any two adjacent events  $t^{(i)}$  and  $t^{(i+1)}$  into  $b_{\text{ins}}$  evenly-spaced bins. Then, we define the insertion operation relative to the  $i$ th event as

$$\text{ins}(\mathbf{t}, i, j) = \left( t^{(0)}, \dots, t^{(i)}, t^{(i)} + \frac{j-1+\alpha}{b_{\text{ins}}}(t^{(i+1)} - t^{(i)}), t^{(i+1)}, \dots, t^{(n+1)} \right) \quad (6)$$

for  $i \in \{0, \dots, n\}$ ,  $j \in [b_{\text{ins}}]$ , where  $\alpha \sim \mathcal{U}(0, 1)$  is a dequantization factor inspired by uniform dequantization in likelihood-based generative models (Theis et al., 2016). The boundary elements  $t^{(0)} = 0$  and  $t^{(n+1)} = T$  ensure that insertions are possible across the entire support  $\mathcal{T}$ . Since the bins between different  $i$  are non-overlapping, insertions are mutually exclusive.

**Substitution:** We implement event substitutions by discretizing the continuous space around each event into  $b_{\text{sub}}$  bins. In this case, the bins are free to overlap, since a substitution is always uniquely determined by the substituted event. We choose a maximum movement distance  $\delta$  and define

$$\text{sub}(\mathbf{t}, i, j) = \text{sort} \left( \{t^{(0)}, \dots, t^{(i-1)}, t^{(i+1)}, \dots, t^{(n+1)}\} \cup \{\tilde{t}^{(i)}\} \right) \quad (7)$$

for  $i \in \{1, \dots, n\}$ ,  $j \in [b_{\text{sub}}]$ , where  $\tilde{t}^{(i)} = [t^{(i)} - \delta + \frac{j-1+\alpha}{b_{\text{sub}}} 2\delta]_0^T$  is the updated event restricted to the support  $\mathcal{T}$  and, again,  $\alpha \sim \mathcal{U}(0, 1)$  is a uniform dequantization factor within the  $j$ -th bin.

**Deletion:** Finally, we define removing event  $i \in \{1, \dots, n\}$  straightforwardly as

$$\text{del}(\mathbf{t}, i) = (t^{(0)}, \dots, t^{(i-1)}, t^{(i+1)}, \dots, t^{(n+1)}). \quad (8)$$

In combination, these operations facilitate any possible edit of an event sequence through insertions and deletions with substitutions as a shortcut for local delete-insert pairs. Note that we neither allow inserting after the last boundary event nor substituting or deleting the first or last boundary events, thus guaranteeing operations to stay in the state space  $\mathcal{X}_{\mathcal{T}}$ . We illustrate the edit operations in Fig. 2.

Our choice of ins, sub and del ensures three key properties: (i) the resulting event sequences remain valid TPPs, (ii) the number of valid operations, e.g.  $\text{ins}(\mathbf{t}, i, j)$ , is independent of the position  $i$ , which is necessary for efficient parameterization, and (iii) at most one unique operation can transition between any two states, which significantly reduces the complexity of the training loss in Section 3.3. While these properties are comparably simple to achieve for token sequences in language modeling (Havasi et al., 2025), where any token can replace any other, they require special care in the case of TPPs. del and sub are defined to ensure that the resulting event sequence remains in increasing order and that the padding events  $t^{(0)}$  and  $t^{(n+1)}$  remain in place. del transitions are unique because the removed event determines exactly which deletion occurred. Similarly, sub transitions are unique because the original position of the substituted event disambiguates the operation, even though two distinct sub operations may yield the same substituted event value. To achieve uniqueness for ins, the

insertion bins corresponding to  $\text{ins}(t, i, j)$  have to be mutually disjoint for any  $i, j$  since insertions lack a removed event to disambiguate them. We achieve this by sizing the bins relative to the distance between  $t^{(i)}$  and  $t^{(i+1)}$ .

**Parameterization** Generating a new event sequence in the Edit Flow framework then means to emit a continuous stream of edit operations by integrating a rate model  $u_s^\theta(\cdot | t)$  from  $s = 0$  to  $s = 1$ . The emitted operations transform a noise sequence  $t_0$  into a data sample  $t_1$  by transitioning through a series of intermediate states  $t$ . Given a current state  $t_s$ , we parameterize the transition rates as

$$u_s^\theta(\text{ins}(t_s, i, j) | t_s) = \lambda_{s,i}^{\text{ins}}(t_s) Q_{s,i}^{\text{ins}}(j | t_s), \quad (9)$$

$$u_s^\theta(\text{sub}(t_s, i, j) | t_s) = \lambda_{s,i}^{\text{sub}}(t_s) Q_{s,i}^{\text{sub}}(j | t_s), \quad (10)$$

$$u_s^\theta(\text{del}(t_s, i) | t_s) = \lambda_{s,i}^{\text{del}}(t_s), \quad (11)$$

where  $\lambda_{s,i}^{\text{del}}, \lambda_{s,i}^{\text{ins}}, \lambda_{s,i}^{\text{sub}}$  denote the total rate of each of the three basic operations at each event  $t^{(i)}$ . The distributions  $Q_{s,i}^{\text{ins}}$  and  $Q_{s,i}^{\text{sub}}$  are *categorical* distributions over the discretization bins  $j \in [b_{\text{ins}}]$  and  $j \in [b_{\text{sub}}]$ , respectively. They distribute the total insertion and substitution rates between the specific options.

### 3.2 AUXILIARY ALIGNMENT SPACE

Training our rate model  $u_s^\theta$  by directly matching a marginalized conditional rate  $u_s(t | t_1, t_0)$  generating a  $p_s(t | t_1, t_0)$ , as is common in discrete flow matching (Campbell et al., 2024; Gat et al., 2024), is challenging or even intractable for Edit Flows, since it would require accounting for all possible edits that could produce  $t$  (Havasi et al., 2025).

To address this, following Havasi et al. (2025), we introduce an auxiliary alignment space for TPPs, where every possible edit operation is uniquely defined in the element wise mixture path  $z_s \sim p_s(z_s | z_0, z_1)$ , making the learning problem tractable.

In language modeling, any token can appear in any position, so Havasi et al. (2025) achieve strong results even when training with a simple alignment that juxtaposes two sequences after shifting one of them by a constant number of places. In our case, for the alignments to correspond to possible edit operations, two events can only be matched, i.e.,  $z_0^{(i)} \neq \epsilon$  and  $z_1^{(i)} \neq \epsilon$ , if  $|z_0^{(i)} - z_1^{(i)}| < \delta$  since otherwise the resulting sub operation would be invalid. Furthermore,  $z_s$  have to correspond to sequences in  $\mathcal{X}_{\mathcal{T}}$ , so  $f_{\text{rm-blanks}}(z)$  has to be increasing, and in particular any mixing  $z_s$  between  $z_0$  and  $z_1$  needs to be valid, i.e.,  $z_s \sim p_s(z_s | z_0, z_1) \Rightarrow f_{\text{rm-blanks}}(z_s) \in \mathcal{X}_{\mathcal{T}}$ .

We find the minimum-cost alignment between the non-boundary events of  $t_0$  and  $t_1$  with the Needleman-Wunsch algorithm (Needleman & Wunsch, 1970), i.e.,

$$\text{align}(t_0, t_1) = \text{wrap-boundaries} \left( \text{Needleman-Wunsch} \left( t_0^{(1:n)}, t_1^{(1:m)}, c_{\text{ins}}, c_{\text{sub}}, c_{\text{del}} \right) \right) \quad (12)$$

and the cost functions

$$c_{\text{sub}}(i, j) = \begin{cases} |t_0^{(i)} - t_1^{(j)}| & \text{if } |t_0^{(i)} - t_1^{(j)}| < \delta \text{ and } t_0^{(i-1)} < t_1^{(j)} < t_0^{(i+1)} \\ \infty & \text{otherwise} \end{cases} \quad (13)$$

$$c_{\text{ins}}(i, j) = \begin{cases} \frac{\delta}{2} & \text{if } t_0^{(i)} < t_1^{(j)} \\ \infty & \text{otherwise} \end{cases} \quad c_{\text{del}}(i, j) = \begin{cases} \frac{\delta}{2} & \text{if } t_0^{(i)} > t_1^{(j)} \\ \infty & \text{otherwise} \end{cases}$$

where wrap-boundaries wraps the sequences with aligned boundary events 0 and  $T$ . The algorithm builds up the aligned sequences pair by pair. The operations corresponds to adding different pairs to the end of  $(z_0, z_1)$ , i.e., insertion  $(\epsilon, t_1^{(j)})$ , deletion  $(t_0^{(i)}, \epsilon)$  and substitution  $(t_0^{(i)}, t_1^{(j)})$  (see Fig. 3).

We carefully craft the cost functions in Eq. (13), to guarantee that the minimum-cost alignment corresponds to ins, sub and del operations as we define them in Section 3.1. With the  $|t_0^{(i)} - t_1^{(j)}| < \delta$

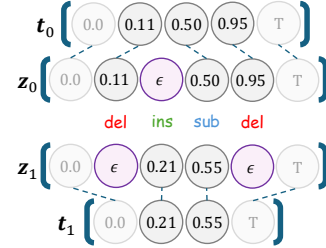


Figure 3: Illustration of the alignment space for  $t_0$  and  $t_1$ .



condition in  $c_{\text{sub}}$ , we ensure that the aligned sequences will never encode a sub operation for two events that are further than  $\delta$  apart. The costs for insertions and deletions and the additional condition on  $c_{\text{sub}}$  ensure that the aligned sequences are jointly sorted, i.e., for any  $i < j$  we have  $\max(z_0^{(i)}, z_1^{(i)}) < \min(z_0^{(j)}, z_1^{(j)})$  where  $\min$  and  $\max$  ignore  $\epsilon$  tokens. This means that any interpolated  $z_s$  is sorted by construction. The validity of encoded ins and del operations follows immediately.

### 3.3 TRAINING

We train our model  $u_s^\theta(\cdot | t_s)$  by optimizing the Bregman divergence in Eq. (5). This amounts to sampling from a coupling  $\pi(z_0, z_1)$  in the aligned auxiliary space and then matching the ground-truth conditional event rates. Note that the coupling  $\pi(z_0, z_1)$  is implicitly defined by its sampling procedure: sample  $t_0, t_1 \sim \pi(t_0, t_1)$  from a coupling of the noise and data distribution, e.g., the independent coupling  $\pi(t_0, t_1) = p(t_0)q(t_1)$ , and then align the sequences  $z_0, z_1 = \text{align}(t_0, t_1)$ . For our choice of operations, the divergence is

$$\mathcal{L} = \mathbb{E}_{(z_0, z_1) \sim \pi(z_0, z_1)} \left[ \sum_{\omega \in \Omega(t_s)} u_s^\theta(\omega | t_s) - \sum_{z_s^{(i)} \neq z_1^{(i)}} \frac{\dot{\kappa}_s}{1 - \kappa_s} \log u_s^\theta(\omega(z_s^{(i)}, z_1^{(i)}) | t_s) \right], \quad (14)$$

where  $\Omega(t_s)$  is the set of all edit operations applicable to  $t_s$  and  $\omega(z_s^{(i)}, z_1^{(i)})$  is the edit operation encoded in the  $i$ -th position of the aligned sequences  $z_s$  and  $z_1$ . To make it precise, we have

$$\Omega(t_s) = \bigcup \left\{ \begin{array}{l} \{\text{ins}(t_s, i, j) \mid i \in \{0\} \cup [n], j \in [b_{\text{ins}}]\} \\ \{\text{sub}(t_s, i, j) \mid i \in [n], j \in [b_{\text{sub}}]\} \\ \{\text{del}(t_s, i) \mid i \in [n]\} \end{array} \right\} \quad (15)$$

and

$$\omega(z_s^{(i)}, z_1^{(i)}) = \begin{cases} \text{ins}(t_s, i', j') & \text{if } z_s^{(i)} = \epsilon \text{ and } z_1^{(i)} \neq \epsilon, \\ \text{sub}(t_s, i', j') & \text{if } z_s^{(i)} \neq \epsilon \text{ and } z_1^{(i)} \neq \epsilon, \\ \text{del}(t_s, i') & \text{if } z_s^{(i)} \neq \epsilon \text{ and } z_1^{(i)} = \epsilon. \end{cases} \quad (16)$$

$i'$  is the index such that  $f_{\text{rm-blanks}}(z_s)$  maps  $z_s^{(i)}$  to  $x_s^{(i')}$  with the convention that  $\epsilon$  is mapped to the same  $i'$  as the last element of  $z_s$  before  $i$  that is not  $\epsilon$ .  $j'$  is the index of the insertion or substitution bin relative to  $x_s^{(i')}$  that  $z_1^{(i)}$  falls into.

### 3.4 SAMPLING

Sampling from our model is done by forward simulation of the CTMC from noise  $t_0 \sim p_{\text{noise}}(t)$  up to  $s = 1$ . We follow (Havasi et al., 2025; Gat et al., 2024) and leverage their Euler approximation, since exact simulation is intractable. Even though the rates are parameterized per element, sampling multiple edits within a time horizon can be done in parallel. At each step of length  $h$ , insertions at position  $i$  occur with probability  $h \lambda_{s,i}^{\text{ins}}(t)$  and deletions or substitutions occur with probability  $h(\lambda_{s,i}^{\text{del}}(t) + \lambda_{s,i}^{\text{sub}}(t))$ . Since they are mutually exclusive the probability of substitution vs deletion is  $\lambda_{s,i}^{\text{sub}}(t)/(\lambda_{s,i}^{\text{sub}}(t) + \lambda_{s,i}^{\text{del}}(t))$ . Lastly, the inserted or substituted events are drawn from the respective distributions  $Q$  to update  $t_s$ . For a short summary of the unconditional sampling step refer to the Euler update step depicted in algorithm Algorithm 1.

---

#### Algorithm 1: Conditional Sampling

---

**Input:**

condition  $t_1^c = C(t_1)$ , noise  $t_0 \sim p_{\text{noise}}$ ,  $h = 1/n_{\text{steps}}$

$(z_0^c, z_1^c) \leftarrow \text{align}(C(t_0), t_1^c)$

**while**  $s < 1$  **do**

**Euler update**

    Sample edits  $\omega_s \sim h u_s^\theta(\cdot | t_s)$

$t_{s+h} \leftarrow \text{apply } \omega_s \text{ to } t_s$

**Recondition**

$\tilde{z}_{s+h}^c \sim p_{s+h}(\cdot | z_0^c, z_1^c)$

$t_{s+h}^c \leftarrow f_{\text{rm-blanks}}(\tilde{z}_{s+h}^c)$

**Merge**

$t_{s+h} \leftarrow C'(t_{s+h}) \cup t_{s+h}^c$

$s \leftarrow s + h$

**end**

**Return:** forecast trajectory  $C'(t_{s=1})$

---

**Conditional sampling.** We can extend the unconditional model to conditional generation given a binary mask on time  $c : \mathcal{T} \rightarrow \{0, 1\}$  (e.g., for forecasting,  $c(t) = t \leq t_{\text{history}}$ ). For a sequence  $\mathbf{t}$ , we define the conditioned part  $C(\mathbf{t}) = \{t \in \mathbf{t} : c(t) = 1\}$  and its complement  $C'(\mathbf{t})$ . Then as depicted in algorithm Algorithm 1, for conditional sampling, we can simply enforce the conditional subsequence to follow a noisy interpolation between  $\mathbf{t}_0^c = C(\mathbf{t}_0)$  and  $\mathbf{t}_1^c = C(\mathbf{t}_1)$ , while the complement evolves freely in the sampling process.

### 3.5 MODEL ARCHITECTURE

For our rate model  $u_s^\theta(\cdot \mid \mathbf{x}_s)$ , we adapt the Llama architecture, a transformer widely applied for variable-length sequences in language modeling (Touvron et al., 2023). We employ FlexAttention in the Llama attention blocks, which supports variable-length sequences natively without padding (Dong et al., 2024). As a first step, we convert the scalar event sequence  $\mathbf{x}_s$  into a sequence of token embeddings by applying  $\text{MLP}(\text{SinEmb}(x_s^{(i)}/T))$  to each to each event, where MLP refers to a small multi-layer perceptron (MLP) and SinEmb is a sinusoidal embedding (Vaswani et al., 2017). We convert  $s$  and  $|\mathbf{x}_s|$  into two additional tokens in an equivalent way with separate MLPs and prepend them to the embedding sequence, which we then feed to the Llama. Lastly, we apply one more MLP to map the output embedding  $\mathbf{h}^{(i)}$  of each event to transition rates. In particular, we parameterize

$$\lambda_{s,i}^{\text{ins}} = \exp(\lambda_M \tanh(\mathbf{h}_{\text{ins}}^{(i)})), \quad \lambda_{s,i}^{\text{sub}} = \exp(\lambda_M \tanh(\mathbf{h}_{\text{sub}}^{(i)})), \quad \lambda_{s,i}^{\text{del}} = \exp(\lambda_M \tanh(\mathbf{h}_{\text{del}}^{(i)})), \quad (17)$$

$$\mathbf{Q}_{s,i}^{\text{ins}} = \text{softmax}(\mathbf{h}_{\text{ins}}^{(i)}), \quad \mathbf{Q}_{s,i}^{\text{sub}} = \text{softmax}(\mathbf{h}_{\text{sub}}^{(i)}). \quad (18)$$

We list the values of all relevant hyperparameters in Appendix A.

## 4 EXPERIMENTS

We evaluate our model on seven real-world and six synthetic benchmark datasets (Omi et al., 2019; Shchur et al., 2020b; Lüdke et al., 2023; 2025). In our experiments, we compare against IFTPP (Shchur et al., 2020a), an autoregressive baseline which consistently shows state-of-the-art performance (Bosser & Taieb, 2023; Lüdke et al., 2023; Kerrigan et al., 2025). We further compare to PSDIFF (Lüdke et al., 2025) and ADDTHIN (Lüdke et al., 2023), given their strong results in both conditional and unconditional settings and their methodological similarity to our approach. All models are trained with five seeds and we select the best checkpoint based on  $W_1$ -over- $d_{\text{IET}}$  against a validation set. EDITPP, ADDTHIN, and PSDIFF are trained unconditionally but can be conditioned at inference time.<sup>1</sup> We list the full results in Appendix E.3.

For forecasts, we compare predicted and target sequences by three metrics:  $d_{\text{Xiao}}$  introduced by Xiao et al. (2017), the mean relative error (MRE) of the event counts and  $d_{\text{IET}}$ , which compares inter-event times to quantify the relation between events such as burstiness. In unconditional generation, we compare our generated sequences to the test set in terms of maximum mean discrepancy (MMD) (Shchur et al., 2020b) and their Wasserstein-1 distance with respect to their counts ( $d_l$ ) and inter-event times ( $d_{\text{IET}}$ ). See Appendix C for details.

### 4.1 UNCONDITIONAL GENERATION

To evaluate how well samples from each TPP model follow the data distribution, we compute distance metrics between 4000 sampled sequences and a hold-out test set. We report the unconditional sampling results in Table 1. EDITPP achieves the best rank in unconditional sampling by strongly matching the test set distribution across all evaluation metrics, outperforming all baselines. The autoregressive baseline IFTPP shows very strong unconditional sampling capability, closely matching and on some dataset and metric combination outperforming the other non-autoregressive baselines ADDTHIN and PSDIFF.

<sup>1</sup>To stay comparable, we employ the conditioning algorithm from Lüdke et al. (2025) for ADDTHIN.

Table 1: Unconditional sampling performance. Bold is best, underlined second best. Ranking follows full results in Appendix E.3 and results are grouped if they fall within the std of the best member.

		H1	H2	NSP	NSR	SC	SR	PG	R/C	R/P	Tx	Tw	Y/A	Y/M
MMD	IFTTP	<u>1.6</u>	<b>1.2</b>	<u>3.2</u>	<u>3.9</u>	<b>6.7</b>	<b>1.2</b>	16.2	<b>7.5</b>	<u>2.0</u>	5.0	<u>2.6</u>	5.8	<b>2.9</b>
	ADDTHIN	2.4	<u>1.8</u>	<u>3.5</u>	15.7	24.6	<u>2.5</u>	4.6	<u>63.0</u>	10.2	<u>4.1</u>	4.4	11.8	<u>3.7</u>
	PSDIFF	3.3	<u>1.8</u>	<b>2.0</b>	5.9	<u>19.8</u>	<u>2.4</u>	<u>3.2</u>	<b>6.5</b>	<b>1.0</b>	<u>3.8</u>	3.4	<u>4.1</u>	<u>3.4</u>
	EDITPP	<b>1.1</b>	<b>1.2</b>	<b>1.7</b>	<b>3.5</b>	<b>7.7</b>	<b>1.0</b>	<b>1.4</b>	<b>8.2</b>	<u>2.4</u>	<b>3.1</b>	<b>1.3</b>	<b>3.7</b>	<u>4.0</u>
		$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-3}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$
$W_{1,d_i}$	IFTTP	<u>20.5</u>	<u>13.3</u>	11.5	14.1	<b>1.5</b>	<u>23.0</u>	294.6	3.9	<u>3.2</u>	<u>2.9</u>	<u>6.5</u>	<u>3.3</u>	2.5
	ADDTHIN	33.3	21.8	12.8	49.0	22.7	41.8	<u>24.5</u>	37.0	33.6	<b>2.3</b>	15.5	6.0	<b>1.6</b>
	PSDIFF	26.9	29.6	<u>5.5</u>	<u>13.3</u>	<u>10.6</u>	30.3	<u>16.1</u>	<b>1.3</b>	<b>2.5</b>	<b>2.8</b>	<u>6.3</u>	<b>1.5</b>	<b>1.5</b>
	EDITPP	<b>7.6</b>	<b>7.0</b>	<b>3.1</b>	<b>1.5</b>	<b>1.3</b>	<b>6.4</b>	<b>6.2</b>	<u>1.9</u>	5.7	<b>2.5</b>	<b>3.4</b>	<b>1.4</b>	<u>1.7</u>
		$\times 10^{-3}$	$\times 10^{-3}$	$\times 10^{-3}$	$\times 10^{-3}$	$\times 10^{-3}$	$\times 10^{-3}$	$\times 10^{-3}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-2}$	$\times 10^{-3}$	$\times 10^{-2}$	$\times 10^{-2}$
$W_{1,d_{\text{DET}}}$	IFTTP	<u>6.3</u>	<u>5.8</u>	<u>3.2</u>	<b>2.3</b>	<u>6.5</u>	<b>7.1</b>	30.3	1.8	7.1	17.4	<u>4.9</u>	<u>3.2</u>	2.8
	ADDTHIN	<u>6.6</u>	7.0	<u>3.2</u>	<u>3.9</u>	15.1	<u>9.4</u>	<u>8.0</u>	5.3	20.0	<b>8.8</b>	5.5	<u>3.2</u>	<u>2.4</u>
	PSDIFF	8.6	9.9	<b>3.0</b>	5.1	32.6	12.8	9.0	<u>1.6</u>	<b>4.3</b>	<u>11.1</u>	6.7	<b>2.4</b>	<b>2.3</b>
	EDITPP	<b>5.3</b>	<b>5.5</b>	<b>3.1</b>	<b>2.2</b>	<b>6.4</b>	<b>7.0</b>	<b>7.5</b>	<b>1.4</b>	<u>6.0</u>	<u>11.1</u>	<b>4.6</b>	<b>2.5</b>	<b>2.3</b>
		$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-2}$	$\times 10^{-1}$	$\times 10^{-2}$	$\times 10^{-1}$	$\times 10^{-3}$	$\times 10^{-2}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$

## 4.2 CONDITIONAL GENERATION (FORECASTING)

Predicting the future given some history window is a fundamental TPP task. For each test sequence, we uniformly sample 50 forecasting windows  $[T_0, T]$ ,  $T_0 \in [\Delta T, T - \Delta T]$ , with minimal history and forecast time  $\Delta T$ . While, this set-up is very similar to the one proposed by Lüdke et al. (2023), there are key differences: we do not fix the forecast window and do not enforce a minimal number of forecast or history events. In fact, even an empty history encodes the information of not having observed an event and a TPP should capture the probability of not observing any event in the future.

We report the forecasting results in Table 2. EDITPP shows very strong forecasting capabilities closely matching or surpassing the baselines across most dataset and metric combinations. Even though IFTPP is explicitly trained to auto-regressively predict the next event given its history, it shows overall worse forecasting capabilities compared to the unconditionally trained EDITPP, ADDTHIN and PSDIFF. This again, underlines previous findings (Lüdke et al., 2023), that autoregressive TPPs can suffer from error accumulation in forecasting. Similar to the unconditional setting, PSDIFF (transformer) outperforms ADDTHIN (convolution with circular padding), which showcases the improved posterior and modeling of long-range interactions.

Table 2: Forecasting accuracy up to  $T$ . Bold is best, underlined second best. Ranking follows full results in Appendix E.3 and results are grouped if they fall within the std of the best member.

		PG	R/C	R/P	Tx	Tw	Y/A	Y/M
$d_{\text{Xiao}}$	IFTTP	6.0	3.9	<u>6.3</u>	4.7	<b>2.6</b>	1.8	3.4
	ADDTHIN	<u>2.5</u>	8.8	<u>7.3</u>	<b>4.0</b>	2.8	<u>1.5</u>	<b>2.9</b>
	PSDIFF	<b>2.4</b>	<b>3.2</b>	<b>4.8</b>	<u>4.4</u>	<u>2.6</u>	<b>1.5</b>	<u>3.0</u>
	EDITPP	<u>2.5</u>	<u>3.4</u>	<b>4.9</b>	<u>4.5</u>	<u>2.7</u>	<b>1.5</b>	<u>3.0</u>
		$\times 10^1$	$\times 10^1$					
MRE	IFTTP	38.9	7.5	3.5	<u>3.2</u>	<b>2.1</b>	3.7	<u>3.9</u>
	ADDTHIN	3.7	14.8	4.6	<b>3.0</b>	3.0	<b>3.5</b>	<b>3.7</b>
	PSDIFF	<b>3.4</b>	<b>3.3</b>	<u>3.0</u>	11.4	2.4	<b>3.5</b>	9.2
	EDITPP	<u>3.5</u>	<u>3.6</u>	<b>2.8</b>	12.3	<u>2.3</u>	<u>3.5</u>	9.0
		$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$
$d_{\text{DET}}$	IFTTP	4.7	<u>6.8</u>	14.7	1.4	2.2	5.9	3.9
	ADDTHIN	<u>4.0</u>	<u>6.9</u>	<u>10.3</u>	<u>1.2</u>	<u>1.5</u>	<b>4.9</b>	<b>2.6</b>
	PSDIFF	<u>4.1</u>	<b>6.2</b>	<b>9.5</b>	<b>1.1</b>	<u>1.5</u>	<b>4.9</b>	<b>2.6</b>
	EDITPP	<b>4.0</b>	<u>6.8</u>	<u>10.1</u>	<b>1.1</b>	<b>1.4</b>	<u>5.0</u>	<u>2.7</u>
		$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-3}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$	$\times 10^{-1}$

## 4.3 EDIT EFFICIENCY

The sub operation allows our model to modify sequences in a more targeted way when compared to PSDIFF or ADDTHIN, which have to rely on just inserts and deletes.



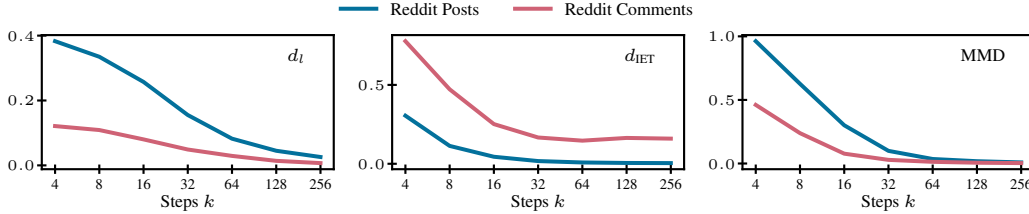


Figure 4: Changing the number of steps  $k$  allows trading off compute and sample quality in terms of  $d_l$ ,  $d_{IET}$  and  $MMD$  at inference time.

Note that one sub operation can replace an insert-delete pair. Table 3 shows this results in EDiTPP using fewer edit operations than PSDIFF on average even if one would count substitutions twice, as an insert and a delete.<sup>2</sup> This is further amplified by the fact, that unlike EDiTPP, PSDIFF and ADDTHIN only indirectly parameterizes the transition edit rates by predicting  $t_1$  by insertion and deletion at every sampling step.

In Table 4, we compare their actual sampling run-time for a batch size of 1024 on the two dataset with the longest sequences. Our implementation beats the reference implementations of ADDTHIN and PSDIFF by a large margin. Note, that for a fair comparison, we fixed the number of sampling steps to 100 in all previous evaluations. As a continuous-time model, EDiTPP can further trade off compute against sample quality at inference time without retraining, in contrast to discrete-time models like ADDTHIN and PSDIFF. Fig. 4 shows that sample quality improves as we increase the number of sampling steps and therefore reduce the discretization step size of the CTMC dynamics. At the same time, the figure also shows rapidly diminishing quality improvements, highlighting potential for substantial speedups with only minor quality loss.

Table 3: Average number of edit operations in unconditional sampling across datasets. Full result in Table 13.

	Ins	Del	Sub	Total
PSDIFF	173.48	61.04	0.00	234.52
EDiTPP	137.42	33.08	29.16	199.65

Table 4: Sample run-time (ms) on a H100 GPU.

	R/P	R/C
ADDTHIN	18,075.62	17,689.36
PSDIFF	7,776.35	3,913.78
EDiTPP	4,120.38	1,505.68

## 5 RELATED WORK

The statistical modeling of TPPs has a long history (Daley & Vere-Jones, 2007; Hawkes, 1971). Classical approaches such as the Hawkes process define parametric conditional intensities, but their limited flexibility has motivated the development of neurally parameterized TPPs:

**Autoregressive Neural TPP:** Most neural TPPs adopt an autoregressive formulation, modeling the distribution of each event conditional on its history. These models consist of two components: a *history encoder* and an *event decoder*. Encoders are typically implemented using recurrent neural networks (Du et al., 2016; Shchur et al., 2020a) or attention mechanisms (Zhang et al., 2020a; Zuo et al., 2020; Mei et al., 2022), with attention-based models providing longer-range context at the cost of higher complexity (Shchur et al., 2021). Further, some propose to encode the history of a TPP in a continuous latent stochastic processes (Chen et al., 2020; Enguehard et al., 2020; Jia & Benson, 2019; Hasan et al., 2023). For the *decoder*, a wide variety of parametrizations have been explored. Conditional intensities or related measures (e.g., hazard function or conditional density), can be modeled, parametrically (Mei & Eisner, 2017; Zuo et al., 2020; Zhang et al., 2020a), via neural networks (Omi et al., 2019), mixtures of kernels (Okawa et al., 2019; Soen et al., 2021; Zhang et al., 2020b) and mixture distributions (Shchur et al., 2020a). Generative approaches further enhance flexibility: normalizing flow-based (Shchur et al., 2020b), GAN-based (Xiao et al., 2017), VAE-based (Li et al., 2018), and diffusion-based decoders (Lin et al., 2022; Yuan et al., 2023) have all been

<sup>2</sup>Due to its recursive definition, ADDTHIN inserts and subsequently deletes some noise events during sampling, which results in additional edit operations compared to PSDIFF.

proposed. While expressive, autoregressive TPPs are inherently sequential, which makes sampling scale at least linearly with sequence length, can lead to error accumulation in multi-step forecasting and limit conditional generation to forecasting.

**Non-autoregressive Neural TPPs:** Similar to our method, these approaches model event sequences through a latent variable process that refines the entire sequence jointly. Diffusion-inspired (Lüdke et al., 2023; 2025) and flow-based generative models (Kerrigan et al., 2025) have recently emerged as promising alternatives to auto-regressive TPP models by directly modelling the joint distribution over event sequences.

## 6 CONCLUSION

We have presented EDITPP, an Edit Flow for TPPs that generalises diffusion-based set interpolation methods (Lüdke et al., 2023; 2025) with a continuous-time flow model introducing substitution as an additional edit operation. By parameterizing insertions, deletions, and substitutions within a CTMC, our approach enables efficient and flexible sequence modeling for TPPs. Empirical results demonstrate that EDITPP matches state-of-the-art performance in both unconditional and conditional generation tasks across synthetic and real-world datasets, while reducing the number of edit operations.

## REFERENCES

- Tanguy Bosser and Souhaib Ben Taieb. On the predictive accuracy of neural temporal point process models for continuous-time event data. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=3OSISBQPrM>. Survey Certification.
- Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design, 2024. URL <https://arxiv.org/abs/2402.04997>.
- Ricky TQ Chen, Brandon Amos, and Maximilian Nickel. Neural spatio-temporal point processes. *arXiv preprint arXiv:2011.04583*, 2020.
- Daryl J Daley and David Vere-Jones. *An introduction to the theory of point processes: volume II: general theory and structure*. Springer Science & Business Media, 2007.
- D.J. Daley and D. Vere-Jones. *An Introduction to the Theory of Point Processes: Volume I: Elementary Theory and Methods*. Probability and Its Applications. Springer New York, 2006.
- Juechu Dong, Boyuan Feng, Driss Guessous, Yanbo Liang, and Horace He. Flex Attention: A Programming Model for Generating Optimized Attention Kernels, December 2024.
- Nan Du, Hanjun Dai, Rakshit Trivedi, Utkarsh Upadhyay, Manuel Gomez-Rodriguez, and Le Song. Recurrent marked temporal point processes: Embedding event history to vector. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1555–1564, 2016.
- Joseph Enguehard, Dan Busbridge, Adam Bozson, Claire Woodcock, and Nils Hammerla. Neural temporal point processes for modelling electronic health records. In *Machine Learning for Health*, pp. 85–113. PMLR, 2020.
- Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky T. Q. Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching, 2024. URL <https://arxiv.org/abs/2407.15595>.
- Ali Hasan, Yu Chen, Yuting Ng, Mohamed Abdelghani, Anderson Schneider, and Vahid Tarokh. Inference and sampling of point processes from diffusion excursions. In *The 39th Conference on Uncertainty in Artificial Intelligence*, 2023.
- Marton Havasi, Brian Karrer, Itai Gat, and Ricky T. Q. Chen. Edit flows: Flow matching with edit operations, 2025. URL <https://arxiv.org/abs/2506.09018>.
- Alan G Hawkes. Spectra of some self-exciting and mutually exciting point processes. *Biometrika*, 58 (1):83–90, 1971.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Neural Information Processing Systems*, 2017.
- Junteng Jia and Austin R Benson. Neural jump stochastic differential equations. *Advances in Neural Information Processing Systems*, 32, 2019.
- Gavin Kerrigan, Kai Nelson, and Padhraic Smyth. Eventflow: Forecasting temporal point processes with flow matching, 2025. URL <https://arxiv.org/abs/2410.07430>.
- Diederik P. Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational Diffusion Models, April 2023.
- Shuang Li, Shuai Xiao, Shixiang Zhu, Nan Du, Yao Xie, and Le Song. Learning temporal point processes via reinforcement learning. *Advances in neural information processing systems*, 31, 2018.

- Marten Lienen, David Lüdke, Jan Hansen-Palmus, and Stephan Günnemann. From Zero to Turbulence: Generative Modeling for 3D Flow Simulation. In *International Conference on Learning Representations*, 2024.
- Marten Lienen, Marcel Kollovieh, and Stephan Günnemann. Generative Modeling with Bayesian Sample Inference, 2025.
- Haitao Lin, Lirong Wu, Guojiang Zhao, Liu Pai, and Stan Z Li. Exploring generative neural temporal point process. *Transactions on Machine Learning Research*, 2022.
- David Lüdke, Marin Biloš, Oleksandr Shchur, Marten Lienen, and Stephan Günnemann. Add and thin: Diffusion for temporal point processes. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=tn9Dldam9L>.
- David Lüdke, Enric Rabasseda Raventós, Marcel Kollovieh, and Stephan Günnemann. Unlocking point processes through point set diffusion. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=4anfpHj0wf>.
- Hongyuan Mei and Jason M Eisner. The neural hawkes process: A neurally self-modulating multivariate point process. In *Neural Information Processing Systems (NeurIPS)*, 2017.
- Hongyuan Mei, Chenghao Yang, and Jason Eisner. Transformer embeddings of irregularly spaced events and their participants. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=Rty5g9imm7H>.
- Saul B. Needleman and Christian D. Wunsch. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology*, 48(3): 443–453, March 1970. ISSN 0022-2836. doi: 10.1016/0022-2836(70)90057-4.
- Alex Nichol and Prafulla Dhariwal. Improved Denoising Diffusion Probabilistic Models. In *International Conference on Machine Learning*, 2021. doi: 10.48550/arXiv.2102.09672.
- Maya Okawa, Tomoharu Iwata, Takeshi Kurashima, Yusuke Tanaka, Hiroyuki Toda, and Naonori Ueda. Deep mixture point processes: Spatio-temporal event prediction with rich contextual information. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 373–383, 2019.
- Takahiro Omi, Kazuyuki Aihara, et al. Fully neural network based model for general temporal point processes. *Advances in neural information processing systems*, 32, 2019.
- Jakob Gulddahl Rasmussen. Lecture notes: Temporal point processes and the conditional intensity function, 2018. URL <https://arxiv.org/abs/1806.00221>.
- Oleksandr Shchur, Marin Biloš, and Stephan Günnemann. Intensity-free learning of temporal point processes. In *International Conference on Learning Representations (ICLR)*, 2020a.
- Oleksandr Shchur, Nicholas Gao, Marin Biloš, and Stephan Günnemann. Fast and flexible temporal point processes with triangular maps. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020b.
- Oleksandr Shchur, Ali Caner Türkmen, Tim Januschowski, and Stephan Günnemann. Neural temporal point processes: A review. *arXiv preprint arXiv:2104.03528*, 2021.
- Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis K. Titsias. Simplified and generalized masked diffusion for discrete data, 2025. URL <https://arxiv.org/abs/2406.04329>.
- Alexander Soen, Alexander Mathews, Daniel Grixti-Cheng, and Lexing Xie. Unipoint: Universally approximating point processes intensities. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 9685–9694, 2021.
- Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative models. In *International Conference on Learning Representations*. arXiv, 2016. doi: 10.48550/arXiv.1511.01844.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and Efficient Foundation Language Models, February 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *Neural Information Processing Systems*, 2017.
- Shuai Xiao, Mehrdad Farajtabar, Xiaojing Ye, Junchi Yan, Le Song, and Hongyuan Zha. Wasserstein learning of deep generative point process models. *Advances in neural information processing systems*, 30, 2017.
- Yuan Yuan, Jingtao Ding, Chenyang Shao, Depeng Jin, and Yong Li. Spatio-temporal Diffusion Point Processes. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 3173–3184, New York, NY, USA, 2023. Association for Computing Machinery.
- Qiang Zhang, Aldo Lipani, Omer Kirnap, and Emine Yilmaz. Self-attentive hawkes process. In *International conference on machine learning*, pp. 11183–11193. PMLR, 2020a.
- Wei Zhang, Thomas Panum, Somesh Jha, Prasad Chalasani, and David Page. Cause: Learning granger causality from event sequences using attribution methods. In *International Conference on Machine Learning*, pp. 11235–11245. PMLR, 2020b.
- Simiao Zuo, Haoming Jiang, Zichong Li, Tuo Zhao, and Hongyuan Zha. Transformer hawkes process. *arXiv preprint arXiv:2002.09291*, 2020.



## A MODEL PARAMETERS

Table 5: Hyperparameters of our  $u_s^\theta(\cdot | \mathbf{x}_s)$  model shared across all datasets.

Parameter	Value
Number of ins bins $b_{\text{ins}}$	64
Number of sub bins $b_{\text{sub}}$	64
Maximum sub distance $\delta$	$T/100$
Maximum log-rate $\lambda_M$	32
$\kappa(s)$	$1 - \cos\left(\frac{\pi}{2}s\right)^2$
<b>Llama architecture:</b>	
Hidden size $H$	64
Layers	2
Attention heads	4
Optimizer	Adam
Sample steps	100

All MLPs have input and output sizes of  $H$ , except for the final MLP whose output size is determined by the number of  $\lambda$  and  $Q$  parameters of the rate. The MLPs have a single hidden layer of size  $4H$ . The sinusoidal embeddings map a scalar  $s \in [0, 1]$  to a vector of length  $H$ . In contrast to Havasi et al. (2025), we choose a cosine  $\kappa$  schedule  $\kappa(s) = 1 - \cos\left(\frac{\pi}{2}s\right)^2$  as proposed by Nichol & Dhariwal (2021) for diffusion models as it improved results slightly compared  $\kappa(s) = s^3$ .

For evaluation, we use an exponential moving average (EMA) of the model weights. We also use low-discrepancy sampling of  $s$  in Eq. (14) during training to smooth the loss and thus training signal (Kingma et al., 2023; Lienen et al., 2025).

We train all models for 20 000 steps and select the best checkpoint by its  $W_1$ -over- $d_{\text{IET}}$ , which we evaluate on a validation set every 1000 steps.

## B DATA

### B.1 SYNTHETIC DATASETS

The six synthetic datasets were generated by Shchur et al. (2020b) following the simulation procedures detailed in Section 4.1 of Omi et al. (2019). Each dataset contains 1,000 sequences supported on the interval  $T = [0, 100]$ . They cover a diverse set of temporal dynamics, defined as follows:

**Hawkes Processes (H1, H2).** Hawkes processes capture self-exciting features of temporal point processes. The two Hawkes processes are parameterized as follows:

$$\lambda(t | \mathcal{H}_t) = \mu + \sum_{t_i < t} \sum_{j=1}^M \alpha_j \beta_j \exp\{-\beta_j(t - t_i)\},$$

with **H1** ( $M = 1, \mu = 0.2, \alpha_1 = 0.8, \beta_1 = 1.0$ ) and **H2** ( $M = 2, \mu = 0.2, \alpha_1 = 0.4, \beta_1 = 1.0, \alpha_2 = 0.4, \beta_2 = 20.0$ ).

**Non-stationary Poisson Process (NSP).** A periodic time-varying intensity:

$$\lambda(t | \mathcal{H}_t) = 0.99 \sin\left(\frac{2\pi t}{20000}\right) + 1.$$

**Stationary Renewal Process (SR).** Inter-event times  $\tau_i = t_{i+1} - t_i$  are i.i.d. from a log-normal distribution (mean 1.0, std. 6.0): this produces bursty patterns with short activity bursts followed by long silent periods.

Table 6: Summary statistics for all synthetic and real-world datasets.  $\tau$  is the average inter-event time.

Full Name	Abbrev.	# Seq.	Mean Len.	Support [0, T]	$\tau$
Hawkes 1	H1	1000	95.4	100	$1.01 \pm 2.38$
Hawkes 2	H2	1000	97.2	100	$0.98 \pm 2.56$
Nonstationary Poisson	NSP	1000	100.3	100	$0.99 \pm 2.22$
Nonstationary Renewal	NSR	1000	98.0	100	$0.98 \pm 1.83$
Self-Correcting	SC	1000	100.2	100	$0.99 \pm 0.71$
Stationary Renewal	SR	1000	109.2	100	$0.83 \pm 2.76$
PUBG	PG	3001	76.5	38 minutes	$0.41 \pm 0.56$
Reddit Comments	R/C	1356	295.7	24 hours	$0.07 \pm 0.28$
Reddit Submissions	R/P	1094	1129.0	24 hours	$0.02 \pm 0.03$
Taxi Pick-ups (Manhattan)	Tx	182	98.4	24 hours	$0.24 \pm 0.40$
Twitter Activity	Tw	2019	14.9	24 hours	$1.26 \pm 2.80$
Yelp Check-ins (Airport)	Y/A	319	30.5	24 hours	$0.77 \pm 1.10$
Yelp Check-ins (Mississauga)	Y/M	319	55.2	24 hours	$0.43 \pm 0.96$

**Non-stationary Renewal Process (NSR).** A stationary renewal process is first generated using a gamma distribution (mean 1.0, std. 0.5), then timestamps are time-warped by

$$t'_i = \int_0^{t_i} r(s) ds, \quad r(t) = 0.99 \sin\left(\frac{2\pi t}{20000}\right) + 1.$$

This induces temporally varying expected inter-event intervals while preserving local correlations.

**Self-correcting Process (SC).** The intensity grows with the time elapsed since the last event:

$$\lambda(t | \mathcal{H}_t) = \exp\left(t - \sum_{t_i < t} 1\right).$$

This discourages extended silent periods and promotes regular spacing.

## B.2 REAL-WORLD DATASETS

We use the seven real-world datasets proposed by (Shchur et al., 2020b):

**PG (PUBG)** represents death-event timestamps from matches of PUBG. **R/C (Reddit-Comments)** consists of comment timestamps within the first 24 hours of threads posted on `r/askscience`, covering 01.01.2018–31.12.2019. **R/P (Reddit-Submissions)** captures daily submission timestamps from `r/politics`, covering 01.01.2017–31.12.2019. **Tx (Taxi)** are taxi pick-up events in the southern part of Manhattan, New York. **Tw (Twitter)** covers tweet timestamps of user ID 25073877, collected over multiple years. **Y/A (Yelp-Airport)** consists of check-in events at McCarran International Airport (27 users, year 2018). Lastly, **Y/M (Yelp-Mississauga)** presents check-ins for businesses in the city of Mississauga (27 users, year 2018).

## C METRICS

A standard way in generative modeling to compare generated and real data is the Wasserstein distance (Heusel et al., 2017). It is the minimum average distance between elements of the two datasets under the optimal (partial) assignment between them,

$$W_p(\mathcal{X}, \mathcal{X}') = \left( \min_{\gamma \in \Gamma(\mathcal{X}, \mathcal{X}')} \mathbb{E}_{(\mathbf{x}, \mathbf{x}') \sim \gamma} [d(\mathbf{x}, \mathbf{x}')^p] \right)^{1/p} \quad (19)$$

where  $d$  is a distance that compares elements from the two sets. In the case of sequences of unequal length, one can choose  $d$  itself as a nested Wasserstein distance (Lienen et al., 2024). Xiao et al. (2017) were the first to design such a distance between TPPs. They exploit a special case of  $W_1$

for sorted sequences of equal length and assign the remaining events of the longer sequence to pseudo-events at  $T$  to define

$$d_{\text{Xiao}}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{|\mathbf{x}|} |t^{(i)} - t'^{(i)}| + \sum_{i=|\mathbf{x}|+1}^{|\mathbf{x}'|} |T - t'^{(i)}| \quad (20)$$

where  $\mathbf{x}'$  is assumed to be the longer sequence.  $d_{\text{Xiao}}$  captures a difference in both location and number of events between two sequences through its two terms.

(Shchur et al., 2020b) propose to compute the MMD between sets based on a Gaussian kernel and  $d_{\text{Xiao}}$ . In addition, we evaluate the event count distributions via a Wasserstein-1 distance with respect to a difference in event counts  $W_{1,d_l}$  where  $d_l(\mathbf{x}, \mathbf{x}') = ||\mathbf{x}| - |\mathbf{x}'||$ . Finally, we the distributions of inter-event times between our generated sequences and real sequences in  $W_{1,d_{\text{IET}}}$ , i.e., a Wasserstein-1 distance of  $d_{\text{IET}}$ .  $d_{\text{IET}}$  is itself the  $W_2$  distance between inter-event times of two sequences and quantifies how adjacent events relate to each other to capture more complex patterns.

## D ABLATIONS

We ablate the hyperparameters  $\delta$ ,  $b_{\text{ins}}$  and  $b_{\text{sub}}$  in Figs. 5 to 8.

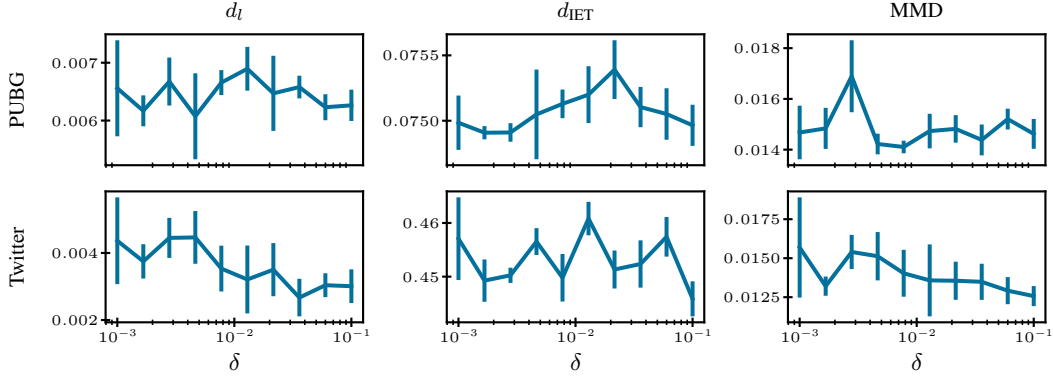


Figure 5: Mean and standard error of  $d_l$ ,  $d_{\text{IET}}$  and MMD on two datasets as we vary the  $\delta$  parameter for substitutions.

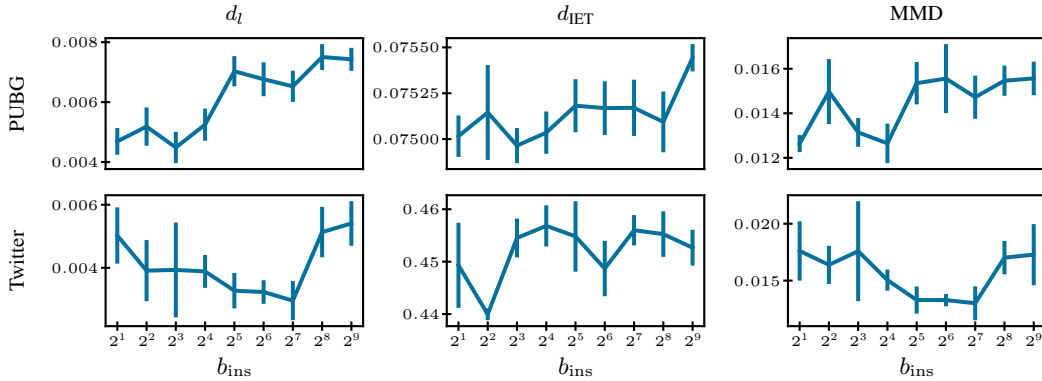


Figure 6: Mean and standard error of  $d_l$ ,  $d_{\text{IET}}$  and MMD on two datasets as we vary the number of insertion bins  $b_{\text{ins}}$ .

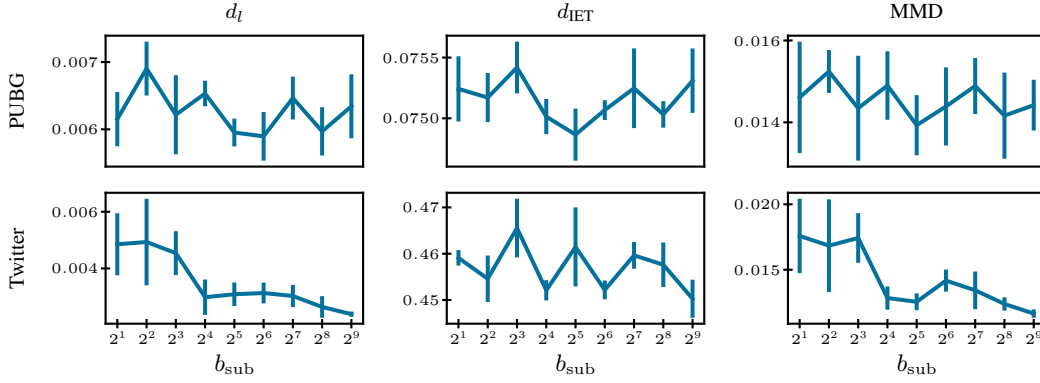


Figure 7: Mean and standard error of  $d_l$ ,  $d_{IET}$  and MMD on two datasets as we vary the number of substitution bins  $b_{sub}$ .

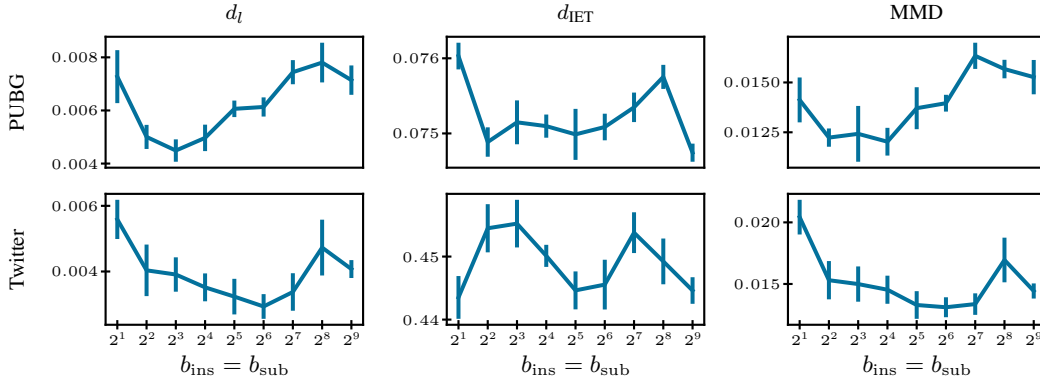


Figure 8: Mean and standard error of  $d_l$ ,  $d_{IET}$  and MMD on two datasets as we vary the number of insertion and substitution bins together.

## E DETAILED RESULTS

### E.1 INPAINTING VS. FORECASTING

To demonstrate the flexibility of EDITPP for conditional generation, we evaluate its performance when generating events on the interval  $[T/3, 2T/3]$ . In the *forecasting* setting, the model is conditioned only on events occurring before  $T/3$ , whereas in the *inpainting* setting, it is conditioned on both the past ( $t < T/3$ ) and the future ( $t > 2T/3$ ). As the results show, providing both past and future context substantially improves the quality of the generated middle segment compared to conditioning on the past alone.

	PG	R/C	R/P	Tx	Tw	Y/A	Y/M
$d_{xiao}$							
Inpainting	$2.22 \pm 0.04$	$22.66 \pm 0.77$	$13.13 \pm 0.35$	$3.02 \pm 0.20$	$1.58 \pm 0.03$	$0.59 \pm 0.02$	$1.11 \pm 0.04$
Forecasting	$2.27 \pm 0.05$	$25.53 \pm 1.16$	$18.09 \pm 0.81$	$3.27 \pm 0.33$	$1.65 \pm 0.04$	$0.63 \pm 0.03$	$1.13 \pm 0.05$
MRE							
Inpainting	$0.29 \pm 0.01$	$5.79 \pm 0.96$	$0.20 \pm 0.01$	$0.60 \pm 0.22$	$1.96 \pm 0.09$	$0.48 \pm 0.01$	$0.74 \pm 0.11$
Forecasting	$0.30 \pm 0.01$	$5.07 \pm 1.13$	$0.33 \pm 0.01$	$0.81 \pm 0.13$	$2.05 \pm 0.11$	$0.52 \pm 0.05$	$0.78 \pm 0.08$
$d_{IET}$							
Inpainting	$0.39 \pm 0.01$	$0.40 \pm 0.01$	$0.01 \pm 0.00$	$0.10 \pm 0.01$	$1.13 \pm 0.02$	$0.80 \pm 0.05$	$0.70 \pm 0.06$
Forecasting	$0.41 \pm 0.01$	$0.44 \pm 0.03$	$0.02 \pm 0.00$	$0.10 \pm 0.00$	$1.18 \pm 0.01$	$0.78 \pm 0.07$	$0.72 \pm 0.08$

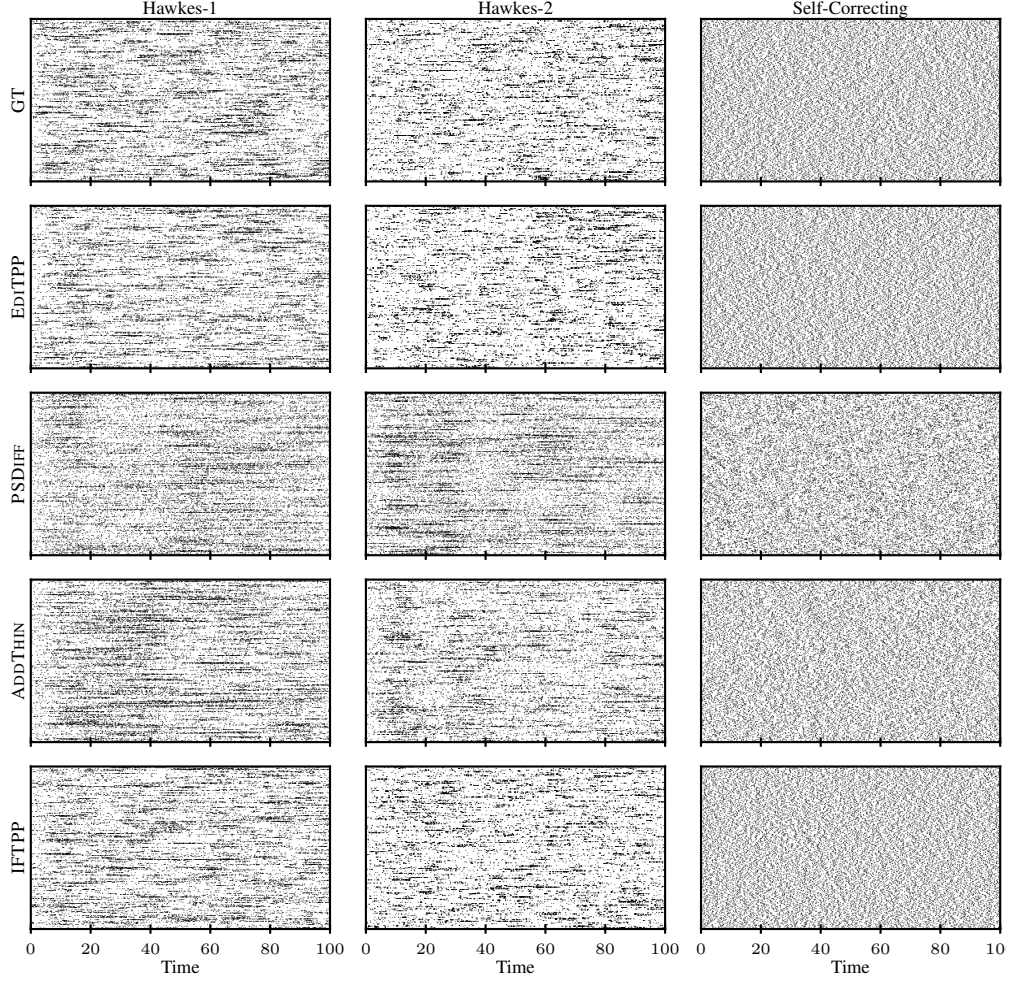
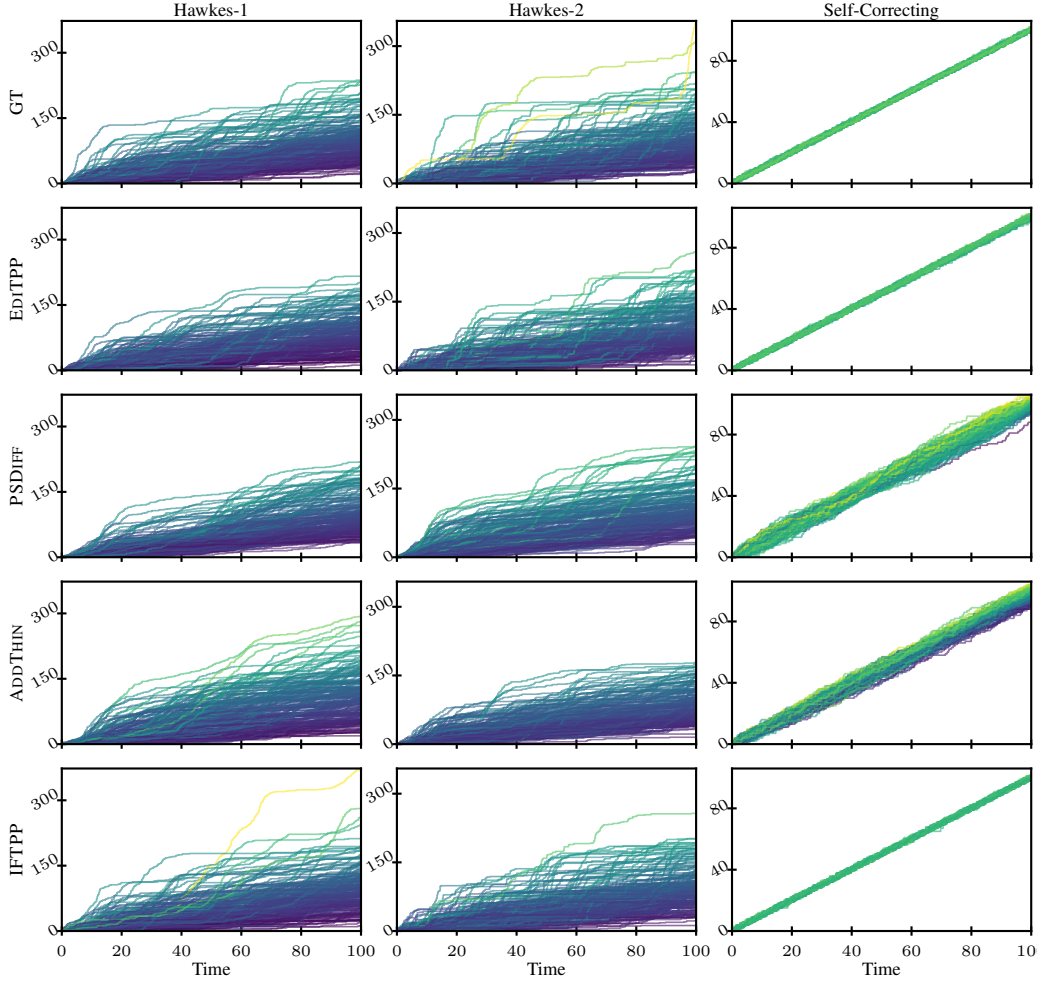


Figure 9: Event times for 200 samples from ground truth data (GT) and each model. Each event sequence is represented as a separate row.

## E.2 PARAMETRIC TPP SAMPLES

To illustrate how well each model captures parametric TPPs, we draw 200 samples for the Hawkes and Self-Correcting processes. In Fig. 10, we plot the cumulative count  $N(t)$  for each sample, while Fig. 9 shows each event sequence as a separate row, directly visualizing the events over time. These visualizations further highlight the strong unconditional sampling performance of EdiTPP demonstrated in Table 1.



Figure 10:  $N(t)$  for 200 samples from ground truth data (GT) and each model.

### E.3 FULL RESULTS

Table 7: Forecasting accuracy up to  $T$  measured by  $d_{\text{IET}}$ .

	EdITPP	PSDIFF	ADDTHIN	IFTTP
PUBG	<b><math>0.400 \pm 0.002</math></b>	$0.413 \pm 0.009$	$0.403 \pm 0.010$	$0.473 \pm 0.019$
Reddit Comments	$0.684 \pm 0.005$	<b><math>0.625 \pm 0.012</math></b>	$0.693 \pm 0.012$	$0.684 \pm 0.012$
Reddit Posts	$0.010 \pm 0.000$	<b><math>0.009 \pm 0.000</math></b>	$0.010 \pm 0.001$	$0.015 \pm 0.003$
Taxi	<b><math>0.113 \pm 0.003</math></b>	<b><math>0.113 \pm 0.001</math></b>	$0.116 \pm 0.001$	$0.145 \pm 0.009$
Twitter	<b><math>1.441 \pm 0.020</math></b>	$1.487 \pm 0.012$	$1.493 \pm 0.033$	$2.187 \pm 0.029$
Yelp Airport	$0.497 \pm 0.009$	<b><math>0.492 \pm 0.005</math></b>	<b><math>0.493 \pm 0.013</math></b>	$0.587 \pm 0.019$
Yelp Mississauga	$0.272 \pm 0.003$	<b><math>0.262 \pm 0.003</math></b>	<b><math>0.260 \pm 0.003</math></b>	$0.388 \pm 0.024$

Table 8: Forecasting accuracy up to  $T$  measured by mean relative error of event counts.

	EDITPP	PSDIFF	ADDTHIN	IFTPP
PUBG	$0.349 \pm 0.001$	<b><math>0.339 \pm 0.008</math></b>	$0.367 \pm 0.005$	$3.892 \pm 0.035$
Reddit Comments	$3.594 \pm 0.118$	<b><math>3.260 \pm 0.268</math></b>	$14.777 \pm 3.226$	$7.515 \pm 2.112$
Reddit Posts	<b><math>0.281 \pm 0.001</math></b>	$0.296 \pm 0.006$	$0.457 \pm 0.065$	$0.352 \pm 0.022$
Taxi	$1.234 \pm 0.036$	$1.140 \pm 0.043$	<b><math>0.301 \pm 0.014</math></b>	$0.321 \pm 0.018$
Twitter	$2.327 \pm 0.042$	$2.435 \pm 0.106$	$2.984 \pm 0.246$	<b><math>2.060 \pm 0.027</math></b>
Yelp Airport	$0.350 \pm 0.007$	<b><math>0.346 \pm 0.004</math></b>	<b><math>0.347 \pm 0.014</math></b>	$0.366 \pm 0.009$
Yelp Mississauga	$0.902 \pm 0.027$	$0.920 \pm 0.033$	<b><math>0.374 \pm 0.012</math></b>	$0.392 \pm 0.012$

Table 9: Forecasting accuracy up to  $T$  measured by  $d_{\text{Xiao}}$ .

	EDITPP	PSDIFF	ADDTHIN	IFTPP
PUBG	$2.478 \pm 0.007$	<b><math>2.400 \pm 0.007</math></b>	$2.466 \pm 0.024$	$5.954 \pm 0.195$
Reddit Comments	$34.135 \pm 0.382$	<b><math>32.467 \pm 0.534</math></b>	$87.666 \pm 20.184$	$39.010 \pm 7.508$
Reddit Posts	<b><math>48.776 \pm 0.355</math></b>	<b><math>47.829 \pm 1.050</math></b>	$72.754 \pm 12.134$	$63.256 \pm 9.695$
Taxi	$4.464 \pm 0.088$	$4.444 \pm 0.076$	<b><math>4.032 \pm 0.129</math></b>	$4.744 \pm 0.125$
Twitter	$2.669 \pm 0.022$	$2.635 \pm 0.078$	$2.802 \pm 0.132$	<b><math>2.557 \pm 0.055</math></b>
Yelp Airport	<b><math>1.524 \pm 0.013</math></b>	<b><math>1.512 \pm 0.016</math></b>	$1.548 \pm 0.026$	$1.795 \pm 0.015$
Yelp Mississauga	$3.027 \pm 0.046$	$3.005 \pm 0.046$	<b><math>2.895 \pm 0.039</math></b>	$3.430 \pm 0.047$

Table 10: Sample quality as measured by MMD.

	EDITPP	PSDIFF	ADDTHIN	IFTPP
Hawkes-1	<b><math>0.011 \pm 0.002</math></b>	$0.033 \pm 0.009$	$0.024 \pm 0.009$	$0.016 \pm 0.002$
Hawkes-2	<b><math>0.012 \pm 0.001</math></b>	$0.018 \pm 0.006$	$0.018 \pm 0.006$	<b><math>0.012 \pm 0.001</math></b>
Nonstationary Poisson	<b><math>0.017 \pm 0.003</math></b>	<b><math>0.020 \pm 0.005</math></b>	$0.035 \pm 0.011$	$0.032 \pm 0.008$
Nonstationary Renewal	<b><math>0.035 \pm 0.001</math></b>	$0.059 \pm 0.006$	$0.157 \pm 0.084$	$0.039 \pm 0.007$
PUBG	<b><math>0.014 \pm 0.001</math></b>	$0.032 \pm 0.012$	$0.046 \pm 0.025$	$0.162 \pm 0.010$
Reddit Comments	<b><math>0.008 \pm 0.001</math></b>	<b><math>0.006 \pm 0.002</math></b>	$0.063 \pm 0.012$	<b><math>0.007 \pm 0.003</math></b>
Reddit Posts	$0.024 \pm 0.001$	<b><math>0.010 \pm 0.002</math></b>	$0.102 \pm 0.004$	$0.020 \pm 0.007$
Self-Correcting	<b><math>0.077 \pm 0.004</math></b>	$0.198 \pm 0.002$	$0.246 \pm 0.018$	<b><math>0.067 \pm 0.011</math></b>
Stationary Renewal	<b><math>0.010 \pm 0.002</math></b>	$0.024 \pm 0.005$	$0.025 \pm 0.013$	<b><math>0.012 \pm 0.002</math></b>
Taxi	<b><math>0.031 \pm 0.002</math></b>	$0.038 \pm 0.005$	$0.041 \pm 0.004$	$0.050 \pm 0.003$
Twitter	<b><math>0.013 \pm 0.002</math></b>	$0.034 \pm 0.007$	$0.044 \pm 0.012$	$0.026 \pm 0.005$
Yelp Airport	<b><math>0.037 \pm 0.002</math></b>	$0.041 \pm 0.004$	$0.118 \pm 0.036$	$0.058 \pm 0.002$
Yelp Mississauga	$0.040 \pm 0.003$	$0.034 \pm 0.007$	$0.037 \pm 0.006$	<b><math>0.029 \pm 0.002</math></b>

Table 11: Sample quality as measured by  $W_1$ -over- $d_{\text{IET}}$ .

	EDITPP	PSDIFF	ADDTHIN	IFTPP
Hawkes-1	<b><math>0.526 \pm 0.020</math></b>	$0.865 \pm 0.035$	$0.655 \pm 0.081$	$0.628 \pm 0.030$
Hawkes-2	<b><math>0.546 \pm 0.005</math></b>	$0.991 \pm 0.038$	$0.703 \pm 0.049$	$0.582 \pm 0.009$
Nonstationary Poisson	<b><math>0.306 \pm 0.005</math></b>	<b><math>0.303 \pm 0.007</math></b>	$0.318 \pm 0.015$	$0.317 \pm 0.006$
Nonstationary Renewal	<b><math>0.224 \pm 0.006</math></b>	$0.511 \pm 0.016$	$0.393 \pm 0.064$	<b><math>0.229 \pm 0.027</math></b>
PUBG	<b><math>0.075 \pm 0.000</math></b>	$0.090 \pm 0.001$	$0.080 \pm 0.003$	$0.303 \pm 0.039$
Reddit Comments	<b><math>0.144 \pm 0.003</math></b>	$0.157 \pm 0.006$	$0.532 \pm 0.014$	$0.176 \pm 0.008$
Reddit Posts	$0.006 \pm 0.000$	<b><math>0.004 \pm 0.000</math></b>	$0.020 \pm 0.001$	$0.007 \pm 0.001$
Self-Correcting	<b><math>0.064 \pm 0.000</math></b>	$0.326 \pm 0.003$	$0.151 \pm 0.005$	$0.065 \pm 0.001$
Stationary Renewal	<b><math>0.697 \pm 0.018</math></b>	$1.281 \pm 0.049$	$0.941 \pm 0.145$	<b><math>0.714 \pm 0.028</math></b>
Taxi	$0.111 \pm 0.001$	$0.111 \pm 0.001$	<b><math>0.088 \pm 0.003</math></b>	$0.174 \pm 0.015$
Twitter	<b><math>0.460 \pm 0.004</math></b>	$0.672 \pm 0.007$	$0.545 \pm 0.024$	$0.492 \pm 0.023$
Yelp Airport	<b><math>0.246 \pm 0.002</math></b>	<b><math>0.244 \pm 0.004</math></b>	$0.316 \pm 0.046$	$0.318 \pm 0.017$
Yelp Mississauga	<b><math>0.226 \pm 0.003</math></b>	<b><math>0.225 \pm 0.003</math></b>	$0.236 \pm 0.004$	$0.276 \pm 0.017$

Table 12: Sample quality as measured by  $W_1$ -over- $d_l$ .

	EDiTPP	PSDIFF	ADDTHIN	IFTTP
Hawkes-1	<b><math>0.008 \pm 0.001</math></b>	$0.027 \pm 0.008$	$0.033 \pm 0.015$	$0.020 \pm 0.004$
Hawkes-2	<b><math>0.007 \pm 0.001</math></b>	$0.030 \pm 0.009$	$0.022 \pm 0.014$	$0.013 \pm 0.003$
Nonstationary Poisson	<b><math>0.003 \pm 0.001</math></b>	$0.006 \pm 0.001$	$0.013 \pm 0.005$	$0.012 \pm 0.003$
Nonstationary Renewal	<b><math>0.001 \pm 0.000</math></b>	$0.013 \pm 0.001$	$0.049 \pm 0.022$	$0.014 \pm 0.011$
PUBG	<b><math>0.006 \pm 0.000</math></b>	$0.016 \pm 0.008$	$0.024 \pm 0.014$	$0.295 \pm 0.007$
Reddit Comments	$0.019 \pm 0.002$	<b><math>0.013 \pm 0.003</math></b>	$0.370 \pm 0.081$	$0.039 \pm 0.023$
Reddit Posts	$0.057 \pm 0.003$	<b><math>0.025 \pm 0.003</math></b>	$0.336 \pm 0.045$	$0.032 \pm 0.011$
Self-Correcting	<b><math>0.001 \pm 0.000</math></b>	$0.011 \pm 0.001$	$0.023 \pm 0.002$	<b><math>0.001 \pm 0.001</math></b>
Stationary Renewal	<b><math>0.006 \pm 0.002</math></b>	$0.030 \pm 0.019$	$0.042 \pm 0.022$	$0.023 \pm 0.005$
Taxi	<b><math>0.025 \pm 0.002</math></b>	<b><math>0.028 \pm 0.004</math></b>	<b><math>0.023 \pm 0.006</math></b>	$0.029 \pm 0.003$
Twitter	<b><math>0.003 \pm 0.001</math></b>	$0.006 \pm 0.003$	$0.015 \pm 0.008$	$0.007 \pm 0.002$
Yelp Airport	<b><math>0.014 \pm 0.002</math></b>	<b><math>0.015 \pm 0.004</math></b>	$0.060 \pm 0.021$	$0.033 \pm 0.003$
Yelp Mississauga	$0.017 \pm 0.003$	<b><math>0.015 \pm 0.002</math></b>	<b><math>0.016 \pm 0.003</math></b>	$0.025 \pm 0.006$

Table 13: Average number of edit operations during unconditional sampling.

	EDiTPP			PSDIFF	
	Ins	Del	Sub	Ins	Del
H1	$55.05 \pm 28.6$	$65.93 \pm 10.5$	$37.74 \pm 10.6$	$92.58 \pm 34.3$	$102.50 \pm 10.6$
H2	$63.69 \pm 32.4$	$71.72 \pm 9.6$	$30.98 \pm 8.7$	$97.57 \pm 38.9$	$101.30 \pm 10.2$
NSP	$49.59 \pm 7.3$	$49.66 \pm 7.1$	$50.56 \pm 8.2$	$100.14 \pm 9.7$	$100.16 \pm 9.8$
NSR	$42.39 \pm 5.8$	$44.23 \pm 7.2$	$55.88 \pm 7.8$	$97.58 \pm 7.5$	$100.51 \pm 10.5$
PG	$56.69 \pm 7.2$	$19.71 \pm 4.4$	$19.96 \pm 4.8$	$76.32 \pm 8.7$	$40.90 \pm 6.4$
R/C	$247.83 \pm 251.3$	$8.69 \pm 6.6$	$15.38 \pm 7.4$	$274.49 \pm 254.3$	$24.45 \pm 5.7$
R/P	$972.89 \pm 300.2$	$0.11 \pm 0.3$	$23.74 \pm 5.0$	$1109.78 \pm 307.1$	$24.56 \pm 7.4$
SC	$34.93 \pm 5.7$	$34.37 \pm 6.9$	$66.18 \pm 8.6$	$98.95 \pm 7.8$	$99.61 \pm 10.2$
SR	$70.67 \pm 21.6$	$64.09 \pm 11.3$	$38.00 \pm 11.5$	$108.85 \pm 31.4$	$101.92 \pm 10.9$
Tw	$11.83 \pm 9.4$	$22.34 \pm 4.6$	$3.04 \pm 2.3$	$14.37 \pm 10.3$	$24.46 \pm 5.2$
Tx	$96.62 \pm 14.4$	$9.35 \pm 3.3$	$17.23 \pm 4.8$	$97.69 \pm 17.4$	$24.37 \pm 5.2$
Y/A	$29.43 \pm 6.4$	$22.64 \pm 5.0$	$9.00 \pm 3.4$	$30.42 \pm 6.3$	$24.30 \pm 4.8$
Y/M	$54.81 \pm 13.8$	$17.14 \pm 4.2$	$11.37 \pm 3.8$	$56.54 \pm 15.5$	$24.45 \pm 4.9$
Mean	137.42	33.08	29.16	173.48	61.04
Total	199.65			234.52	