

Neural Temporal Point Processes: A Review

Oleksandr Shchur^{1,2}, Ali Caner Türkmen², Tim Januschowski² and Stephan Günnemann¹

¹Technical University of Munich, Germany

²Amazon Research

{shchur,guennemann}@in.tum.de, {atturkm,tjnsch}@amazon.com,

Abstract

Temporal point processes (TPP) are probabilistic generative models for continuous-time event sequences. Neural TPPs combine the fundamental ideas from point process literature with deep learning approaches, thus enabling construction of flexible and efficient models. The topic of neural TPPs has attracted significant attention in recent years, leading to the development of numerous new architectures and applications for this class of models. In this review paper we aim to consolidate the existing body of knowledge on neural TPPs. Specifically, we focus on important design choices and general principles for defining neural TPP models. Next, we provide an overview of application areas commonly considered in the literature. We conclude this survey with the list of open challenges and important directions for future work in the field of neural TPPs.

1 Introduction

Many applications in science and industry are concerned with collections of events with *timesteps*. Earthquake occurrences in seismology, neural spike trains in neuroscience, trades and orders in a financial market, and user activity logs on the web, can all be represented as sequences of discrete (instantaneous) events observed in continuous time.

Temporal point processes (TPP) are probabilistic models for such event data [Daley and Vere-Jones, 2007]. More specifically, TPPs are generative models of variable-length point sequences observed on the real half-line—here interpreted as arrival times of events. TPPs are built on rich theoretical foundations, with early work dating back to the beginning of the 20th century, where they were used to model the arrival of insurance claims and telephone traffic [Brockmeyer *et al.*, 1948; Cramér, 1969]. The field underwent rapid development in the second half of the century, and TPPs were applied to a wide array of domains including seismology, neuroscience, and finance.

Nevertheless, TPPs entered the mainstream of machine learning research only very recently. One of the exciting ideas developed at the intersection of the fields of point processes

and machine learning were *neural* TPPs [Du *et al.*, 2016; Mei and Eisner, 2017]. Classical (non-neural) TPPs can only capture relatively simple patterns in event occurrences, such as self-excitation [Hawkes, 1971]. In contrast, neural TPPs are able to learn complex dependencies, and are often even computationally more efficient than their classical counterparts. As such, the literature on neural TPPs has witnessed rapid growth since their introduction.

Scope and structure of the paper. The goal of this survey is to provide an overview of neural TPPs, with focus on models (Sections 3–5) and their applications (Section 6). Due to limited space, we do not attempt to describe every existing approach in full detail, but rather focus on general principles and building blocks for constructing neural TPP models.

We also discuss the main challenges that the field currently faces and outline some future research directions (Section 7). For other reviews of TPPs for machine learning, we refer the reader to the tutorial by [Gomez-Rodriguez and Valera, 2018]; and two recent surveys by [Yan, 2019], who also covers non-neural approaches, and [Enguehard *et al.*, 2020] who experimentally compare neural TPP architectures in applications to healthcare data. Our work provides a more detailed overview of neural TPP architectures and their applications compared to the above papers.

2 Background and Notation

A TPP [Daley and Vere-Jones, 2007] is a probability distribution over variable-length sequences in some time interval $[0, T]$. A realization of a *marked* TPP can be represented as an event sequence $X = \{(t_1, m_1), \dots, (t_N, m_N)\}$, where N , the number of events, is itself a random variable. Here, $0 < t_1 < \dots < t_N \leq T$ are the arrival times of events and $m_i \in \mathcal{M}$ are the marks. Categorical marks (*i.e.*, $\mathcal{M} = \{1, \dots, K\}$) are most commonly considered in practice, but other choices, such as $\mathcal{M} = \mathbb{R}^D$, are also possible. Sometimes, it is convenient to instead work with the inter-event times $\tau_i = t_i - t_{i-1}$, where $t_0 = 0$ and $t_{N+1} = T$. For a given X , we denote the history of past events at time t as $\mathcal{H}_t = \{(t_j, m_j) : t_j < t\}$.

A distribution of a TPP with K categorical marks can be characterized by K conditional intensity functions $\lambda_k^*(t)$ (one

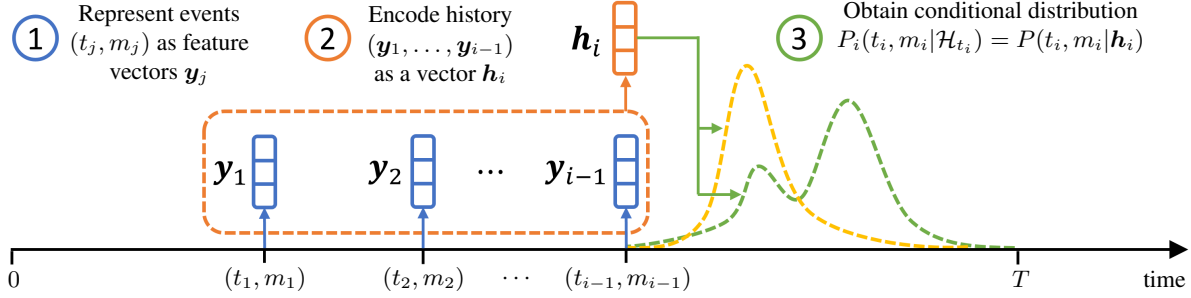


Figure 1: Schematic representation of an autoregressive neural TPP model.

for each mark k) that are defined as

$$\lambda_k^*(t) = \lim_{\Delta t \downarrow 0} \frac{\Pr(\text{event of type } k \text{ in } [t, t + \Delta t] | \mathcal{H}_t)}{\Delta t}, \quad (1)$$

where the $*$ symbol is used as a shorthand for conditioning on the history \mathcal{H}_t . Note that the above definition also applies to unmarked TPPs if we set $K = 1$. While the conditional intensity is often mentioned in the literature, it is not the only way to characterize a TPP—we will consider some alternatives in the next section.

3 Autoregressive Neural TPPs

Neural TPPs can be defined as autoregressive models, as done in the seminal work by [Du *et al.*, 2016]. Such autoregressive TPPs operate by sequentially predicting the time and mark of the next event (t_i, m_i) . Usually, we can decompose this procedure into 3 steps (see Figure 1):

1. Represent each event (t_j, m_j) as a feature vector \mathbf{y}_j .
2. Encode the history \mathcal{H}_{t_i} (represented by a sequence of feature vectors $(\mathbf{y}_1, \dots, \mathbf{y}_{i-1})$) into a fixed-dimensional history embedding \mathbf{h}_i .
3. Use the history embedding \mathbf{h}_i to parametrize the conditional distribution over the next event $P_i(t_i, m_i | \mathcal{H}_{t_i})$.

We will now discuss each of these steps in more detail.

3.1 Representing Events as Feature Vectors

First, we need to represent each event (t_j, m_j) as a feature vector \mathbf{y}_j that can then be fed into an encoder neural network (Section 3.2). We consider the features $\mathbf{y}_j^{\text{time}}$ based on the arrival time t_j (or inter-event time τ_j) and $\mathbf{y}_j^{\text{mark}}$ based on the mark m_j . The vector \mathbf{y}_j is obtained by combining $\mathbf{y}_j^{\text{time}}$ and $\mathbf{y}_j^{\text{mark}}$, *e.g.*, via concatenation.

Time features $\mathbf{y}_j^{\text{time}}$. Earlier works used the inter-event time τ_j or its logarithm $\log \tau_j$ as the time-related feature [Du *et al.*, 2016; Omi *et al.*, 2019]. Recently, [Zuo *et al.*, 2020] and [Zhang *et al.*, 2020a] proposed to instead obtain features from t_j using trigonometric functions, which is based on positional encodings used in transformer language models [Vaswani *et al.*, 2017].

Mark features $\mathbf{y}_j^{\text{mark}}$. Categorical marks are usually encoded with an embedding layer [Du *et al.*, 2016]. Real-valued marks

can be directly used as $\mathbf{y}_j^{\text{mark}}$. This ability to naturally handle different mark types is one of the attractive properties of neural TPPs (compared to classical TPP models).

3.2 Encoding the History into a Vector

The core idea of autoregressive neural TPP models is that event history \mathcal{H}_{t_i} (a variable-sized set) can be represented as a *fixed-dimensional* vector \mathbf{h}_i [Du *et al.*, 2016]. We review the two main families of approaches for encoding the past events $(\mathbf{y}_1, \dots, \mathbf{y}_{i-1})$ into a history embedding \mathbf{h}_i next.

Recurrent encoders start with an initial hidden state \mathbf{h}_1 . Then, after each event (t_i, m_i) they update the hidden state as $\mathbf{h}_{i+1} = \text{Update}(\mathbf{h}_i, \mathbf{y}_i)$. The hidden states \mathbf{h}_i are then used as the history embedding. The Update function is usually implemented based on the RNN, GRU or LSTM update equations [Du *et al.*, 2016; Xiao *et al.*, 2017b].

The main advantage of recurrent models is that they allow us to compute the history embedding \mathbf{h}_i for all N events in the sequence in $O(N)$ time. This compares favorably even to classical non-neural TPPs, such as the Hawkes process, where the likelihood computation in general scales as $O(N^2)$. One downside of recurrent models is their inherently sequential nature. Because of this, such models are usually trained via truncated backpropagation through time, which only provides an approximation to the true gradients [Sutskever, 2013].

Set aggregation encoders directly encode the feature vectors $(\mathbf{y}_1, \dots, \mathbf{y}_{i-1})$ into a history embedding \mathbf{h}_i . Unlike recurrent models, here the encoding is done independently for each i . The encoding operation can be defined, *e.g.*, using self-attention [Zuo *et al.*, 2020; Zhang *et al.*, 2020a]. It is postulated that such encoders are better at capturing long-range dependencies between events compared to recurrent encoders. However, more thorough evaluation is needed to validate this claim (see Section 7). On the one hand, set aggregation encoders can compute \mathbf{h}_i for each event in parallel, unlike recurrent models. On the other hand, usually the time of this computation scales as $O(N^2)$ with the sequence length N , since each \mathbf{h}_i depends on all the past events (and the model does not have a Markov property). This problem can be mitigated by restricting the encoder to only the last L events $(\mathbf{y}_{i-L}, \dots, \mathbf{y}_{i-1})$, thus reducing the time complexity to $O(NL)$.

3.3 Predicting the Time of the Next Event

For simplicity, we start by considering the unmarked case and postpone the discussion of marked TPPs until the next section. An autoregressive TPP models the distribution of the next arrival time t_i given the history \mathcal{H}_{t_i} . This is equivalent to considering the distribution of the next inter-event time τ_i given \mathcal{H}_{t_i} , which we denote as $P_i^*(\tau_i)$.¹ The distribution $P_i^*(\tau_i)$ can be represented by any of the following functions:

1. probability density function $f_i^*(\tau_i)$
2. cumulative distribution function $F_i^*(\tau_i) = \int_0^{\tau_i} f_i^*(u)du$
3. survival function $S_i^*(\tau_i) = 1 - F_i^*(\tau_i)$
4. hazard function $\phi_i^*(\tau_i) = f_i^*(\tau_i)/S_i^*(\tau_i)$
5. cumulative hazard function $\Phi_i^*(\tau_i) = \int_0^{\tau_i} \phi_i^*(u)du$.

In an autoregressive neural TPP, we pick a parametric form for one of the above functions and compute its parameters using the history embedding \mathbf{h}_i . For example, the conditional PDF f_i^* might be obtained as

$$f_i^*(\tau_i) = f(\tau_i|\boldsymbol{\theta}_i), \quad \text{where } \boldsymbol{\theta}_i = \sigma(\mathbf{W}\mathbf{h}_i + \mathbf{b}). \quad (2)$$

Here $f(\cdot|\boldsymbol{\theta})$ is some parametric density function over $[0, \infty)$ (e.g., exponential density) and \mathbf{W} , \mathbf{b} are learnable parameters. A nonlinear function $\sigma(\cdot)$ can be used to enforce necessary constraints on the parameters, such as non-negativity.

It is important to ensure that the chosen parametrization defines a valid probability distribution. For instance, the PDF f_i^* must be non-negative and satisfy $\int_0^\infty f_i^*(u)du = 1$. This corresponds to the cumulative hazard function Φ_i^* being strictly increasing, differentiable and satisfying $\Phi_i^*(0) = 0$ and $\lim_{\tau \rightarrow \infty} \Phi_i^*(\tau) = \infty$.

Some parametrizations of $P_i^*(\tau)$ proposed in the literature fail to satisfy the above conditions. For example, the hazard function $\phi_i^*(\tau) = \exp(w\tau + b)$ [Du *et al.*, 2016] fails to satisfy $\lim_{\tau \rightarrow \infty} \Phi_i^*(\tau) = \infty$ if the parameter w is negative. One more example is the cumulative hazard function defined by a single-hidden-layer neural network with positive weights $\mathbf{w}, \mathbf{v} \in \mathbb{R}_+^D$ as $\Phi_i^*(\tau) = \text{softplus}(\sum_{d=1}^D v_d \tanh(w_d\tau + b_d))$ [Omi *et al.*, 2019] that fails to satisfy both $\Phi_i^*(0) = 0$ and $\lim_{\tau \rightarrow \infty} \Phi_i^*(\tau) = \infty$. Since above parametric functions do not define a valid distribution over the inter-event times, sampling methods have a non-zero probability of failing or producing invalid event sequences. Therefore, it's crucial to pick a valid parametrization when designing a neural TPP model.

Specifying one of the functions (1) – (5) listed above uniquely identifies the conditional distribution $P_i^*(\tau_i)$, and thus the other four functions in the list. This, however, does not mean that choosing which function to parametrize is unimportant. In particular, some choices of the hazard function ϕ_i^* cannot be integrated analytically, which becomes a problem when computing the log-likelihood (as we will see in Section 5). In contrast, it is usually trivial to obtain ϕ_i^* from any parametrization of Φ_i^* , since differentiation is easier than integration [Omi *et al.*, 2019]. More generally, there are three important aspects that one has to keep in mind when specifying the distribution $P_i^*(\tau_i)$:

¹We again use * to denote conditioning on the history \mathcal{H}_{t_i} .

- *Flexibility*: Does the given parametrization of $P_i^*(\tau_i)$ allow us to approximate any distribution, e.g., a multimodal one?
- *Closed-form likelihood*: Can we compute either the CDF F_i^* , SF S_i^* or CHF Φ_i^* analytically? These functions are involved in the log-likelihood computation (Section 5), and therefore should be computed in closed form for efficient model training. Approximating these functions with Monte Carlo or numerical quadrature is slower and less accurate.
- *Closed-form sampling*: Can we draw samples from $P_i^*(\tau_i)$ analytically? In the best case, this should be done with inversion sampling [Rasmussen, 2011], which requires analytically inverting either F_i^* , S_i^* or Φ_i^* . Inverting these functions via numerical root-finding is again slower and less accurate. Approaches based on thinning [Ogata, 1981] are also not ideal, since they do not benefit from parallel hardware like GPUs. Moreover, closed-form inversion sampling enables the reparametrization trick [Mohamed *et al.*, 2020], which allows us to train TPPs using sampling-based losses (Section 5.2).

Existing approaches offer different trade-offs between the above criteria. For example, a simple unimodal distribution offers closed-form sampling and likelihood computation, but lacks flexibility [Du *et al.*, 2016]. One can construct more expressive distributions by parametrizing the cumulative hazard Φ_i^* either with a mixture of kernels [Okawa *et al.*, 2019; Zhang *et al.*, 2020b] or a neural network [Omi *et al.*, 2019], but this will prevent closed-form sampling. Specifying the PDF f_i^* with a mixture distribution [Shchur *et al.*, 2020a] or Φ_i^* with invertible splines [Shchur *et al.*, 2020b] allows to define a flexible model where both sampling and likelihood computation can be done analytically. Finally, parametrizations that require approximating Φ_i^* via Monte Carlo integration are less efficient and accurate than all of the above-mentioned approaches [Omi *et al.*, 2019].

As a side note, a more flexible parametrization might also be more difficult to train or more prone to overfitting. Therefore, the choice of the parametrization is an important modeling decision that depends on the application.

Lastly, we would like to point out that the view of a TPP as an autoregressive model naturally connects to the traditional conditional intensity characterization (Equation 1). The conditional intensity $\lambda^*(t)$ can be defined by stitching together the hazard functions ϕ_i^*

$$\lambda^*(t) = \begin{cases} \phi_1^*(t) & \text{if } 0 \leq t \leq t_1 \\ \phi_2^*(t - t_1) & \text{if } t_1 < t \leq t_2 \\ \vdots & \\ \phi_{N+1}^*(t - t_N) & \text{if } t_N < t \leq T \end{cases} \quad (3)$$

In the TPP literature, the hazard function ϕ_i^* is often called “intensity,” even though the two are, technically, different mathematical objects.

3.4 Modeling the Marks

In a *marked* autoregressive TPP, one has to parametrize the conditional distribution $P_i^*(\tau_i, m_i)$ using the history embed-

ding \mathbf{h}_i . We first consider categorical marks, as they are most often used in practice.

Conditionally independent models factorize the distribution $P_i^*(\tau_i, m_i)$ into a product of two independent distributions $P_i^*(\tau_i)$ and $P_i^*(m_i)$ that are both parametrized using \mathbf{h}_i [Du *et al.*, 2016; Xiao *et al.*, 2017b]. The time distribution $P_i^*(\tau_i)$ can be parametrized using any of the choices described in Section 3.3, such as the hazard function $\phi_i^*(\tau)$. The mark distribution $P_i^*(m_i)$ is a categorical distribution with probability mass function $p_i^*(m_i = k)$. In this case, the conditional intensity for mark k is computed as

$$\lambda_k^*(t) = p_i^*(m_i = k) \cdot \phi_i^*(t - t_{i-1}), \quad (4)$$

where $(i-1)$ is the index of the most recent event before time t . Note that if we set $K = 1$, we recover the definition of the intensity function from Equation 3.

While conditionally independent models require fewer parameters to specify $P_i^*(\tau_i, m_i)$, recent works suggest that this simplifying assumption may hurt predictive performance [Enguehard *et al.*, 2020]. There are two ways to model dependencies between τ_i and m_i that we consider below.

Time conditioned on marks [Zuo *et al.*, 2020; Enguehard *et al.*, 2020]. In this case, we must specify a separate distribution $P_i^*(\tau_i | m_i = k)$ for each mark $k \in \{1, \dots, K\}$. Suppose that for each k we represent $P_i^*(\tau_i | m_i = k)$ with a hazard function $\phi_{ik}^*(\tau)$. Then the conditional intensity $\lambda_k^*(t)$ for mark k is computed simply as

$$\lambda_k^*(t) = \phi_{ik}^*(t - t_{i-1}). \quad (5)$$

It is possible to model the dependencies across marks on a coarser grid in time, which significantly improves the scalability in the number of marks [Türkmen *et al.*, 2019b].

Marks conditioned on time. Here, the inter-event time is distributed according to $P_i^*(\tau_i)$, and for each τ we need to specify a distribution $P_i^*(m_i | \tau_i = \tau)$. We again assume that the time distribution $P_i^*(\tau_i)$ is described by a hazard function $\phi_i^*(\tau)$, and $P_i^*(m_i | \tau_i = \tau)$ can be parametrized, *e.g.*, using a Gaussian process [Biloš *et al.*, 2019]. In this case the conditional intensity $\lambda_k^*(t)$ is computed as

$$\lambda_k^*(t) = p_i^*(m_i = k | \tau_i = t - t_{i-1}) \cdot \phi_i^*(t - t_{i-1}), \quad (6)$$

where we used notation analogous to Equation 4. The term $\phi_i^*(t - t_{i-1})$ is often referred to as ‘‘ground intensity.’’

Other mark types. A conditionally independent model can easily handle any type of marks by specifying an appropriate distribution $P_i^*(m_i)$. Dependencies between continuous marks and the inter-event time can be incorporated by modeling the joint density $f_i^*(\tau_i, m_i)$ [Zhu *et al.*, 2020].

4 Continuous-time State Evolution

Another line of research has studied neural TPPs that operate completely in continuous time. Such models define a left-continuous state $\mathbf{h}(t)$ at all times $t \in [0, T]$. The state is initialized to some value $\mathbf{h}(0)$. Then, for each event i the state is updated as

$$\begin{aligned} \mathbf{h}(t_i) &= \text{Evolve}(\mathbf{h}(t_{i-1}), t_{i-1}, t_i) \\ \lim_{\varepsilon \rightarrow 0} \mathbf{h}(t_i + \varepsilon) &= \text{Update}(\mathbf{h}(t_i), \mathbf{y}_i) \end{aligned} \quad (7)$$

The $\text{Evolve}(\mathbf{h}(t_{i-1}), t_{i-1}, t_i)$ procedure evolves the the state continuously over the time interval $(t_{i-1}, t_i]$ between the events. Such state evolution can be implemented either via exponential decay [Mei and Eisner, 2017] or, more generally, be governed by an ordinary differential equation [Rubanova *et al.*, 2019; Jia and Benson, 2019]. The $\text{Update}(\mathbf{h}(t_i), \mathbf{y}_i)$ operation performs an instantaneous update to the state, similarly to the recurrent encoder from last section.

While the above procedure might seem similar to a recurrent model from Section 3.2, the continuous-time model uses the state $\mathbf{h}(t)$ differently from the autoregressive model. In a continuous-time model, the state $\mathbf{h}(t)$ is used to directly define the intensity λ_k^* for each mark k as

$$\lambda_k^*(t) = g_k(\mathbf{h}(t)), \quad (8)$$

where $g_k : \mathbb{R}^H \rightarrow \mathbb{R}_{>0}$ is a non-linear function that maps the hidden state $\mathbf{h}(t)$ to the value of the conditional intensity for mark k at time t . Such function, for example, can be implemented as $g_k(\mathbf{h}(t)) = \text{softplus}(\mathbf{w}_k^T \mathbf{h}(t))$ or a multi-layer perceptron [Chen *et al.*, 2021b].

To summarize, in a continuous-time state model, the state $\mathbf{h}(t)$ is defined for all $t \in [0, T]$, and the intensity $\lambda_k^*(t)$ at time t depends only on the current state $\mathbf{h}(t)$. In contrast, in an autoregressive model, the discrete-time state \mathbf{h}_i is updated only after an event occurs.. Hence, the state \mathbf{h}_i defines the entire conditional distribution $P_i^*(\tau_i, m_i)$, and therefore the intensity $\lambda_k^*(t)$ in the interval $(t_i, t_{i+1}]$.

Discussion. Continuous-time models have several advantages compared to autoregressive ones. For example, they provide a natural framework for modeling irregularly-sampled time series, which is valuable in medical applications [Enguehard *et al.*, 2020]. By modeling the unobserved attributes as a function of the state $\mathbf{h}(t)$, it is easy to estimate each attribute at any time t [Rubanova *et al.*, 2019]. These models are also well-suited for modeling spatio-temporal point processes (*i.e.*, with marks in $\mathcal{M} = \mathbb{R}^D$) [Jia and Benson, 2019; Chen *et al.*, 2021b].

However, this flexibility comes at a cost: evaluating both the state evolution (Equation 7) and the model likelihood (Equation 9) requires numerically approximating intractable integrals. This makes training in continuous-time models slower than for autoregressive ones. Sampling similarly requires numerical approximations.

5 Parameter Estimation

5.1 Maximum Likelihood Estimation

Negative log-likelihood (NLL) is the default training objective for both neural and classical TPP models. NLL for a single sequence X with categorical marks is computed as

$$\begin{aligned} -\log p(X) &= -\sum_{i=1}^N \sum_{k=1}^K \mathbb{1}(m_i = k) \log \lambda_k^*(t_i) \\ &\quad + \sum_{k=1}^K \left(\int_0^T \lambda_k^*(u) du \right). \end{aligned} \quad (9)$$

The log-likelihood can be understood using the following two facts. First, the quantity $\lambda_k^*(t_i) dt$ corresponds to the

probability of observing an event of type k in the infinitesimal interval $[t_i, t_i + dt)$ conditioned on the past events \mathcal{H}_{t_i} . Second, we can compute the probability of not observing any events of type k in the rest of the interval $[0, T]$ as $\exp\left(-\int_0^T \lambda_k^*(u) du\right)$. By taking the logarithm, summing these expressions for all events (t_i, m_i) and event types k , and finally negating, we obtain Equation 9.

Computing the NLL for TPPs can sometimes be challenging due to the presence of the integral in the second line of Equation 9. One possible solution is to approximate this integral using Monte Carlo integration [Mei and Eisner, 2017] or numerical quadrature [Rubanova *et al.*, 2019; Zuo *et al.*, 2020]. However, some autoregressive neural TPP models allow us to compute the NLL analytically, which is more accurate and computationally efficient. We demonstrate this using the following example.

Suppose we model $P_i^*(\tau_i, m_i)$ using the “time conditioned on marks” approach from Section 3.4. That is, we specify the distribution $P_i^*(\tau_i | m_i = k)$ for each mark k with a cumulative hazard function $\Phi_{ik}^*(\tau)$. By combining Equation 5 and Equation 9, we can rewrite the expression for the NLL as

$$-\log p(X) = -\sum_{i=1}^N \sum_{k=1}^K \mathbb{1}(m_i = k) \log \phi_{ik}^*(\tau_i) + \sum_{i=1}^{N+1} \sum_{k=1}^K \Phi_{ik}^*(\tau_i). \quad (10)$$

Assuming that our parametrization allows us to compute $\Phi_{ik}^*(\tau)$ analytically, we are now able to compute the NLL in closed form (without numerical integration). Remember that the hazard function ϕ_{ik}^* can be easily obtained by differentiation as $\phi_{ik}^*(\tau) = \frac{\partial}{\partial \tau} \Phi_{ik}^*(\tau)$. Finally, note that the NLL can also be expressed in terms of, *e.g.*, the conditional PDFs f_{ik}^* or survival functions S_{ik}^* (Section 3.3).

Evaluating the NLL in Equation 10 can be still computationally expensive when K , the number of marks, is extremely large. Several works propose approximations based on noise-contrastive-estimation that can be used in this regime [Guo *et al.*, 2018; Mei *et al.*, 2020].

Training. Usually, we are given a training set $\mathcal{D}_{\text{train}}$ of sequences that are assumed to be sampled i.i.d. from some unknown data-generating process. The TPP parameters (*e.g.*, weights of the encoder in Section 3) are learned by minimizing the NLL of the sequences in $\mathcal{D}_{\text{train}}$. This is typically done with some variant of (stochastic) gradient descent. In practice, the NLL loss is often normalized per sequence, *e.g.*, by the interval length T [Enguehard *et al.*, 2020]. Importantly, this normalization constant cannot depend on X , so it would be incorrect to, for example, normalize the NLL by the number of events N in each sequence. Finally, the i.i.d. assumption is not appropriate for all TPP datasets; [Boyd *et al.*, 2020] show how to overcome this challenge by learning sequence embeddings.

5.2 Alternatives to MLE

TPPs can be trained using objective functions other than the NLL. Often, these objectives can be expressed as

$$\mathbb{E}_{X \sim p(X)} [f(X)]. \quad (11)$$

Such sampling-based losses have been used by several approaches for learning generative models from the training sequences. These approaches aim to maximize the similarity between the training sequences in $\mathcal{D}_{\text{train}}$ and sequences X generated by the TPP model $p(X)$ using a scoring function $f(X)$. Examples include procedures based on Wasserstein distance [Xiao *et al.*, 2017a], adversarial losses [Yan *et al.*, 2018; Wu *et al.*, 2018] and inverse reinforcement learning [Li *et al.*, 2018].

Sometimes, the objective function of the form (11) arises naturally based on the application. For instance, in reinforcement learning, a TPP $p(X)$ defines a stochastic policy and $f(X)$ is the reward function [Upadhyay *et al.*, 2018]. When learning with missing data, the missing events X are sampled from the TPP $p(X)$, and $f(X)$ corresponds to the NLL of the observed events [Gupta *et al.*, 2021]. Finally, in variational inference, the TPP $p(X)$ defines an approximate posterior and $f(X)$ is the evidence lower bound (ELBO) [Shchur *et al.*, 2020b; Chen *et al.*, 2021a].

In practice, the gradients of the loss (Equation 11) w.r.t. the model parameters usually cannot be computed analytically and therefore are estimated with Monte Carlo. Earlier works used the score function estimator [Upadhyay *et al.*, 2018], but modern approaches rely on the more accurate pathwise gradient estimator (also known as the “reparametrization trick”) [Mohamed *et al.*, 2020]. The latter relies on our ability to sample with reparametrization from $P_i^*(\tau_i, m_i)$, which again highlights the importance of the chosen parametrization for the conditional distribution, as described in Section 3.3.

On a related note, sampling-based losses for TPPs (Equation 11) can be non-differentiable, since N , the number of events in a sequence, is a discrete random variable. This problem can be solved by deriving a differentiable relaxation to the loss [Shchur *et al.*, 2020b].

6 Applications

The literature on neural TPPs mostly considers their applications in web-related domains, *e.g.*, for modeling user activity on social media. Most existing applications of neural TPPs fall into one of two categories:

- Prediction tasks, where the goal is to predict the time and / or type of future events;
- Structure discovery, where the tasks is to learn dependencies between different event types.

We now discuss these in more detail.

6.1 Prediction

Prediction is among the key tasks associated with temporal models. In case of a TPP, the goal is usually to predict the times and marks of future events given the history of past events. Such queries can be answered using the conditional distribution $P_i^*(\tau_i, m_i)$ defined by the neural TPP

model. Nearly all papers mentioned in previous sections feature numerical experiments on such prediction tasks. Some works combine elements of neural TPP models (Section 3) with other neural network architectures to solve specific real-world prediction tasks related to event data. We give some examples below.

Recommender systems are a recent application area for TPPs. Here, the goal is to predict the next most likely purchase event, in terms of both the time of purchase and the type (*i.e.*, item), given a sequence of customer interactions or purchases in the past. Neural TPPs are especially well-suited for this task, as they can learn embeddings for large item sets using neural networks, similar to other neural recommendation models. Moreover, representing the temporal dimension of purchase behavior enables *time-sensitive* recommendations, *e.g.*, used to time promotions. For example, [Kumar *et al.*, 2019] address the next item prediction problem by embedding the event history to a vector, but are concerned only with predicting the next item type (mark). This approach is extended to a full model of times and events, with a hierarchical RNN model for intra- and inter-session activity represented on different levels, in [Vassøy *et al.*, 2019].

Another common application of event sequence prediction is within the *human mobility* domain. Here, events are *spatio-temporal*, featuring coordinates in both time and space, potentially along with other marks. Examples include predicting taxi trajectories in a city (*e.g.*, the ECML/PKDD 2015 challenge²), or user check-ins on location-based social media. [Yang *et al.*, 2018] address this task with a full neural TPP, on four different mobility data sets. This extends the approach in DeepMove [Feng *et al.*, 2018], where the authors similarly use an RNNs to compute embeddings of timestamped sequences, but limit the predictions to the next location alone.

Other applications include clinical event prediction [Enguehard *et al.*, 2020], predicting timestamped sequences of interactions of patients with the health system; human activity prediction for assisted living [Shen *et al.*, 2018], and demand forecasting in sparse time series [Türkmen *et al.*, 2019a].

6.2 Structure Discovery & Modeling Networks

In prediction tasks we are interested in the conditional distributions learned by the model. In contrast, in structure discovery tasks the parameters learned by the model are of interest.

For example, in *latent network discovery* applications we observe event activity generated by K users, each represented by a categorical mark. The goal is to infer an influence matrix $A \in \mathbb{R}^{K \times K}$ that encodes the dependencies between different marks [Linderman and Adams, 2014]. Here, the entries of A can be interpreted as edge weights in the network. Historically, this task has been addressed using non-neural models such as the Hawkes process [Hawkes, 1971]. The main advantage of network discovery approaches based on neural TPPs [Zhang *et al.*, 2021] is their ability to handle more general interaction types.

Learning *Granger causality* is another task that is closely related to the network discovery problem. [Eichler *et al.*,

2017] and [Achab *et al.*, 2017] have shown that the influence matrix A of a Hawkes process completely captures the notion of Granger causality in multivariate event sequences. Recently, [Zhang *et al.*, 2020b] generalized this understanding to neural TPP models, where they used the method of integrated gradients to estimate dependencies between event types, with applications to information diffusion on the web.

Neural TPPs have also been used to model *information diffusion and network evolution* in social networks [Trivedi *et al.*, 2019]. The DyRep approach by Trivedi *et al.* generalizes an earlier non-neural framework, COEVOLVE, by [Farajtabar *et al.*, 2017]. Similarly, Know-Evolve [Trivedi *et al.*, 2017] models dynamically evolving knowledge graphs with a neural TPP. A related method, DeepHawkes, was developed specifically for modeling item popularity in information cascades [Cao *et al.*, 2017].

Other applications. Neural TPPs have been featured in works in other research fields and application domains. For instance, [Huang *et al.*, 2019] proposed an RNN-based Poisson process model for speech recognition. [Sharma *et al.*, 2018] developed a latent-variable neural TPP method for modeling the behavior of larval zebrafish. Performing inference in the model allowed the authors to detect distinct behavioral patterns in zebrafish activity. Lastly, [Upadhyay *et al.*, 2018] showed how to automatically choose timing for interventions in an interactive environment by combining a TPP model with the framework of reinforcement learning.

7 Open Challenges

We conclude with a discussion of what, in our opinion, are the main challenges that the field of neural TPPs currently faces.

7.1 Experimental Setup

Lack of standardized experimental setups and high-quality benchmark datasets makes a fair comparison of different neural TPP architectures problematic.

Each neural TPP model consists of multiple components, such as the history encoder and the parametrization of the conditional distributions (Sections 3 and 4). New architectures often change all these components at once, which makes it hard to pinpoint the source of empirical gains. Carefully-designed ablation studies are necessary to identify the important design choices and guide the search for better models.

On a related note, the choice of baselines varies greatly across papers. For example, papers proposing autoregressive models (Section 3) rarely compare to continuous-time state models (Section 4), and vice versa. Considering a wider range of baselines is necessary to fairly assess the strengths and weaknesses of different families of approaches.

Finally, it is not clear whether the datasets commonly used in TPP literature actually allow us to find models that will perform better on real-world prediction tasks. In particular, Enguehard *et al.* point out that two popular datasets (MIMIC-II and StackOverflow) can be “solved” by a simple history-independent baseline. Also, common implicit assumptions, such as treating the training sequences as i.i.d. (Section 5), might not be appropriate for existing datasets.

²<http://www.geolink.pt/ecmlpkdd2015-challenge/>

To conclude, open-sourcing libraries with reference implementations of various baseline methods and collecting large high-quality benchmark datasets are both critical next steps for neural TPP researchers. A recently released library with implementations of some autoregressive models by [Enguehard *et al.*, 2020] takes a step in this direction, but, for instance, doesn't include continuous-time state models.

7.2 Evaluation Metrics

Many of the metrics commonly used to evaluate TPP models are not well suited for quantifying their predictive performance and have subtle failure modes.

For instance, consider the NLL score (Section 5), one of the most popular metrics for quantifying predictive performance in TPPs. As mentioned in Section 3.4, the NLL consists of a continuous component for the time t_i and a discrete component for the mark m_i . Therefore, reporting a single NLL score obscures information regarding the model's performance on predicting marks and times separately. For example, the NLL score is affected disproportionately by errors in marks as the number of marks increase. Moreover, the NLL can be "fooled" on datasets where the arrival times t_i are measured in discrete units, such as seconds—a flexible model can produce an arbitrarily low NLL by placing narrow density "spikes" on these discrete values [Uria *et al.*, 2013].

More importantly, the NLL is mostly irrelevant as a measure of error in real-world applications—it yields very little insight into model performance from a domain expert's viewpoint. However, other metrics, such as accuracy for the mark prediction, and mean absolute error (MAE) or mean squared error (MSE) for inter-event time prediction are even less suited for evaluating neural TPPs. These metrics measure the quality of a *single* future event prediction, and MAE/MSE only take a *point estimate* of τ_i into account. One doesn't need to model the entire distribution over τ_i (as done by a TPP) to perform well w.r.t. such metrics. If only single-event predictions are of interest, one could instead use a simple baseline that only models $P_i^*(m_i)$ or produces a point estimate τ_i^{pred} . This baseline can be trained by directly minimizing absolute or squared error for inter-event times or cross-entropy for marks. TPPs are, in contrast, probabilistic models trained with the NLL loss; so comparing them to point-estimate baselines using above metrics is unfair.

The main advantage of neural TPPs compared to simple "point-estimate" methods is their ability to sample entire trajectories of future events. Such probabilistic forecasts capture uncertainty in predictions and are able to answer more complex prediction queries (*e.g.*, "How many events of type k_1 will happen immediately after an event of type k_2 ?"). Probabilistic forecasts are universally preferred to point estimates in the neighboring field of time series modeling [Gneiting and Katzfuss, 2014; Alexandrov *et al.*, 2020]. A variety of metrics for evaluating the quality of such probabilistic forecasts for (regularly-spaced) time series have been proposed [Gneiting *et al.*, 2008], but they haven't been generalized to marked continuous-time event sequences. Developing metrics that are based on entire sampled event sequences can help us unlock the full potential of neural TPPs and allow us to better compare different models. More generally, we should take

advantage of the fact that TPP models learn an entire distribution over trajectories and rethink how these models are applied to prediction tasks.

7.3 Applications

While most recent works have focused on developing new architectures for better prediction and structure discovery, other applications of neural TPPs remain largely understudied.

Applications of classical TPPs go beyond the above two tasks. For example, latent-variable TPP models have been used for event clustering [Mavroforakis *et al.*, 2017; Xu and Zha, 2017] and change point detection [Altieri *et al.*, 2015]. Other applications include anomaly detection [Li *et al.*, 2017], optimal control [Zarezade *et al.*, 2017; Tabibian *et al.*, 2019] and fighting the spread of misinformation online [Kim *et al.*, 2018].

Moreover, most papers on neural TPPs consider datasets originating from the web and related domains, such as recommender systems and knowledge graphs. Meanwhile, traditional application areas for TPPs, like neuroscience, seismology and finance, have not received as much attention. Adapting neural TPPs to these traditional domains requires answering exciting research questions. To name a few, spike trains in neuroscience [Aljadeff *et al.*, 2016] are characterized by both high numbers of marks (*i.e.*, neurons that are being modeled) and high rates of event occurrence (*i.e.*, firing rates), which requires efficient and scalable models. Applications in seismology usually require interpretable models [Bray and Schoenberg, 2013], and financial datasets often contain dependencies between various types of assets that are far more complex than self-excitation, commonly encountered in social networks [Bacry *et al.*, 2015].

To summarize, considering new tasks and new application domains for neural TPP models is an important and fruitful direction for future work.

References

- [Achab *et al.*, 2017] Massil Achab, Emmanuel Bacry, Stéphane Gaïffas, Iacopo Mastromatteo, and Jean-François Muzy. Uncovering causality from multivariate Hawkes integrated cumulants. In *ICML*, 2017.
- [Alexandrov *et al.*, 2020] Alexander Alexandrov, Konstantinos Benidis, Michael Bohlke-Schneider, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Danielle C. Maddix, Syama Rangapuram, David Salinas, Jasper Schulz, Lorenzo Stella, Ali Caner Türkmen, and Yuyang Wang. GluonTS: Probabilistic and neural time series modeling in Python. *JMLR*, 21(116), 2020.
- [Aljadeff *et al.*, 2016] Johnatan Aljadeff, Benjamin J Lansdell, Adrienne L Fairhall, and David Kleinfeld. Analysis of neuronal spike trains, deconstructed. *Neuron*, 91(2), 2016.
- [Altieri *et al.*, 2015] Linda Altieri, E Marian Scott, Daniela Cocchi, and Janine B Illian. A changepoint analysis of spatio-temporal point processes. *Spatial Statistics*, 2015.
- [Bacry *et al.*, 2015] Emmanuel Bacry, Iacopo Mastromatteo, and Jean-François Muzy. Hawkes processes in finance. *Market Microstructure and Liquidity*, 1, 2015.

- [Biloš *et al.*, 2019] Marin Biloš, Bertrand Charpentier, and Stephan Günnemann. Uncertainty on asynchronous time event prediction. *NeurIPS*, 2019.
- [Boyd *et al.*, 2020] Alex Boyd, Robert Bamler, Stephan Mandt, and Padhraic Smyth. User-dependent neural sequence models for continuous-time event data. *NeurIPS*, 2020.
- [Bray and Schoenberg, 2013] Andrew Bray and Frederic Paik Schoenberg. Assessment of point process models for earthquake forecasting. *Statistical science*, 2013.
- [Brockmeyer *et al.*, 1948] E Brockmeyer, HL Halstrøm, and A Jensen. The life and works of AK Erlang. *Transactions of the Danish Academy of Technical Sciences*, 2, 1948.
- [Cao *et al.*, 2017] Qi Cao, Huawei Shen, Keting Cen, Wentao Ouyang, and Xueqi Cheng. Deephawkes: Bridging the gap between prediction and understanding of information cascades. In *CIKM*, 2017.
- [Chen *et al.*, 2021a] Ricky TQ Chen, Brandon Amos, and Maximilian Nickel. Learning neural event functions for ordinary differential equations. *ICLR*, 2021.
- [Chen *et al.*, 2021b] Ricky TQ Chen, Brandon Amos, and Maximilian Nickel. Neural spatio-temporal point processes. *ICLR*, 2021.
- [Cramér, 1969] Harald Cramér. Historical review of Filip Lundberg’s works on risk theory. *Scandinavian Actuarial Journal*, 1969(sup3), 1969.
- [Daley and Vere-Jones, 2007] Daryl J Daley and David Vere-Jones. *An introduction to the theory of point processes: Volume II: general theory and structure*. 2007.
- [Du *et al.*, 2016] 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 *KDD*, 2016.
- [Eichler *et al.*, 2017] Michael Eichler, Rainer Dahlhaus, and Johannes Dueck. Graphical modeling for multivariate hawkes processes with nonparametric link functions. *Journal of Time Series Analysis*, 38(2):225–242, 2017.
- [Enguehard *et al.*, 2020] 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*, 2020.
- [Farajtabar *et al.*, 2017] Mehrdad Farajtabar, Yichen Wang, Manuel Gomez-Rodriguez, Shuang Li, Hongyuan Zha, and Le Song. COEVOLVE: A joint point process model for information diffusion and network evolution. *JMLR*, 18(1), 2017.
- [Feng *et al.*, 2018] Jie Feng, Yong Li, Chao Zhang, Funing Sun, Fanchao Meng, Ang Guo, and Depeng Jin. Deepmove: Predicting human mobility with attentional recurrent networks. In *WWW*, 2018.
- [Gneiting and Katzfuss, 2014] Tilmann Gneiting and Matthias Katzfuss. Probabilistic forecasting. *Annual Review of Statistics and Its Application*, 2014.
- [Gneiting *et al.*, 2008] Tilmann Gneiting, Larissa I Stanberry, Eric P Gritmit, Leonhard Held, and Nicholas A Johnson. Assessing probabilistic forecasts of multivariate quantities, with an application to ensemble predictions of surface winds. 2008.
- [Gomez-Rodriguez and Valera, 2018] Manuel Gomez-Rodriguez and Isabel Valera. Learning with temporal point processes. *Tutorial at ICML*, 2018.
- [Guo *et al.*, 2018] Ruocheng Guo, Jundong Li, and Huan Liu. INITIATOR: Noise-contrastive estimation for marked temporal point process. In *IJCAI*, 2018.
- [Gupta *et al.*, 2021] Vinayak Gupta, Srikanta Bedathur, Sourangshu Bhattacharya, and Abir De. Learning temporal point processes with intermittent observations. In *AISTATS*, 2021.
- [Hawkes, 1971] Alan G Hawkes. Point spectra of some mutually exciting point processes. *Journal of the Royal Statistical Society: Series B*, 33(3), 1971.
- [Huang *et al.*, 2019] Hengguan Huang, Hao Wang, and Brian Mak. Recurrent Poisson process unit for speech recognition. In *AAAI*, 2019.
- [Jia and Benson, 2019] Junteng Jia and Austin R Benson. Neural jump stochastic differential equations. *NeurIPS*, 2019.
- [Kim *et al.*, 2018] Jooyeon Kim, Behzad Tabibian, Alice Oh, Bernhard Schölkopf, and Manuel Gomez-Rodriguez. Leveraging the crowd to detect and reduce the spread of fake news and misinformation. In *WSDM*, 2018.
- [Kumar *et al.*, 2019] Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in temporal interaction networks. In *KDD*, 2019.
- [Li *et al.*, 2017] Shuang Li, Yao Xie, Mehrdad Farajtabar, Apurv Verma, and Le Song. Detecting changes in dynamic events over networks. *IEEE Transactions on Signal and Information Processing over Networks*, 3(2), 2017.
- [Li *et al.*, 2018] Shuang Li, Shuai Xiao, Shixiang Zhu, Nan Du, Yao Xie, and Le Song. Learning temporal point processes via reinforcement learning. *NeurIPS*, 2018.
- [Linderman and Adams, 2014] Scott Linderman and Ryan Adams. Discovering latent network structure in point process data. In *ICML*, 2014.
- [Mavroforakis *et al.*, 2017] Charalampos Mavroforakis, Isabel Valera, and Manuel Gomez-Rodriguez. Modeling the dynamics of learning activity on the web. In *WWW*, 2017.
- [Mei and Eisner, 2017] Hongyuan Mei and Jason Eisner. The neural Hawkes process: A neurally self-modulating multivariate point process. *NeurIPS*, 2017.
- [Mei *et al.*, 2020] Hongyuan Mei, Tom Wan, and Jason Eisner. Noise-contrastive estimation for multivariate point processes. *NeurIPS*, 2020.
- [Mohamed *et al.*, 2020] Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte carlo gradient estimation in machine learning. *JMLR*, 21(132), 2020.

- [Ogata, 1981] Yoshihiko Ogata. On Lewis' simulation method for point processes. *Transactions on Information Theory*, 27(1), 1981.
- [Okawa *et al.*, 2019] 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 *KDD*, 2019.
- [Omi *et al.*, 2019] Takahiro Omi, Naonori Ueda, and Kazuyuki Aihara. Fully neural network based model for general temporal point processes. *NeurIPS*, 2019.
- [Rasmussen, 2011] Jakob Gulddahl Rasmussen. Temporal point processes. *Lecture Notes*, 2011.
- [Rubanova *et al.*, 2019] Yulia Rubanova, Ricky TQ Chen, and David Duvenaud. Latent ordinary differential equations for irregularly-sampled time series. *NeurIPS*, 2019.
- [Sharma *et al.*, 2018] Anuj Sharma, Robert Johnson, Florian Engert, and Scott Linderman. Point process latent variable models of larval zebrafish behavior. *NeurIPS*, 2018.
- [Shchur *et al.*, 2020a] Oleksandr Shchur, Marin Biloš, and Stephan Günnemann. Intensity-free learning of temporal point processes. *ICLR*, 2020.
- [Shchur *et al.*, 2020b] Oleksandr Shchur, Nicholas Gao, Marin Biloš, and Stephan Günnemann. Fast and flexible temporal point processes with triangular maps. *NeurIPS*, 2020.
- [Shen *et al.*, 2018] Yang Shen, Bingbing Ni, Zefan Li, and Ning Zhuang. Egocentric activity prediction via event modulated attention. In *ECCV*, 2018.
- [Sutskever, 2013] Ilya Sutskever. *Training recurrent neural networks*. PhD thesis, University of Toronto, 2013.
- [Tabibian *et al.*, 2019] Behzad Tabibian, Utkarsh Upadhyay, Abir De, Ali Zaregade, Bernhard Schölkopf, and Manuel Gomez-Rodriguez. Enhancing human learning via spaced repetition optimization. *PNAS*, 116(10), 2019.
- [Trivedi *et al.*, 2017] Rakshit Trivedi, Hanjun Dai, Yichen Wang, and Le Song. Know-Evolve: Deep temporal reasoning for dynamic knowledge graphs. In *ICML*, 2017.
- [Trivedi *et al.*, 2019] Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, and Hongyuan Zha. DyRep: Learning representations over dynamic graphs. In *ICLR*, 2019.
- [Türkmen *et al.*, 2019a] Ali Caner Türkmen, Yuyang Wang, and Tim Januschowski. Intermittent demand forecasting with deep renewal processes. *TPP @ NeurIPS*, 2019.
- [Türkmen *et al.*, 2019b] Ali Caner Türkmen, Yuyang Wang, and Alexander J Smola. FastPoint: Scalable deep point processes. In *ECML PKDD*, 2019.
- [Upadhyay *et al.*, 2018] Utkarsh Upadhyay, Abir De, and Manuel Gomez-Rodriguez. Deep reinforcement learning of marked temporal point processes. *NeurIPS*, 2018.
- [Uria *et al.*, 2013] Benigno Uria, Iain Murray, and Hugo Larochelle. Rnade: The real-valued neural autoregressive density-estimator. *NeurIPS*, 2013.
- [Vassøy *et al.*, 2019] Bjørnar Vassøy, Massimiliano Ruocco, Eliezer de Souza da Silva, and Erlend Aune. Time is of the essence: a joint hierarchical rnn and point process model for time and item predictions. In *WSDM*, 2019.
- [Vaswani *et al.*, 2017] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017.
- [Wu *et al.*, 2018] Qitian Wu, Chaoqi Yang, Hengrui Zhang, Xiaofeng Gao, Paul Weng, and Guihai Chen. Adversarial training model unifying feature driven and point process perspectives for event popularity prediction. In *CIKM*, 2018.
- [Xiao *et al.*, 2017a] Shuai Xiao, Mehrdad Farajtabar, Xiaojing Ye, Junchi Yan, Le Song, and Hongyuan Zha. Wasserstein learning of deep generative point process models. *NeurIPS*, 2017.
- [Xiao *et al.*, 2017b] Shuai Xiao, Junchi Yan, Xiaokang Yang, Hongyuan Zha, and Stephen Chu. Modeling the intensity function of point process via recurrent neural networks. In *AAAI*, 2017.
- [Xu and Zha, 2017] Hongteng Xu and Hongyuan Zha. A Dirichlet mixture model of Hawkes processes for event sequence clustering. *NeurIPS*, 2017.
- [Yan *et al.*, 2018] Junchi Yan, Xin Liu, Liangliang Shi, Changsheng Li, and Hongyuan Zha. Improving maximum likelihood estimation of temporal point process via discriminative and adversarial learning. In *IJCAI*, 2018.
- [Yan, 2019] Junchi Yan. Recent advance in temporal point process: from machine learning perspective. *SJTU Technical Report*, 2019.
- [Yang *et al.*, 2018] Guolei Yang, Ying Cai, and Chandan K Reddy. Recurrent spatio-temporal point process for check-in time prediction. In *CIKM*, 2018.
- [Zaregade *et al.*, 2017] Ali Zaregade, Abir De, Utkarsh Upadhyay, Hamid R Rabiee, and Manuel Gomez-Rodriguez. Steering social activity: A stochastic optimal control point of view. *JMLR*, 18, 2017.
- [Zhang *et al.*, 2020a] Qiang Zhang, Aldo Lipani, Omer Kirnap, and Emine Yilmaz. Self-attentive Hawkes process. In *ICML*, 2020.
- [Zhang *et al.*, 2020b] Wei Zhang, Thomas Panum, Somesh Jha, Prasad Chalasani, and David Page. CAUSE: Learning granger causality from event sequences using attribution methods. In *ICML*, 2020.
- [Zhang *et al.*, 2021] Qiang Zhang, Aldo Lipani, and Emine Yilmaz. Learning neural point processes with latent graphs. In *WWW*, 2021.
- [Zhu *et al.*, 2020] Shixiang Zhu, Henry Shaowu Yuchi, and Yao Xie. Adversarial anomaly detection for marked spatio-temporal streaming data. In *ICASSP*, 2020.
- [Zuo *et al.*, 2020] Simiao Zuo, Haoming Jiang, Zichong Li, Tuo Zhao, and Hongyuan Zha. Transformer Hawkes process. In *ICML*, 2020.