

---

# Quantifying Causal Contribution in Rare Event Data

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We introduce a framework for causal discovery and attribution of causal influence  
2 for rare events in time series data—where the interest is in identifying causal links  
3 and root causes of individual discrete events rather than the types of these events.  
4 Specifically, we build on the theory of temporal point processes, and describe  
5 a discrete-time analogue of Hawkes processes to model the occurrence of self-  
6 exciting rare events with instantaneous effects. We then introduce several scores  
7 to measure causal influence among individual events. These statistics are drawn  
8 from causal inference and temporal point process theories, describe complementary  
9 aspects of causality in temporal event data, and obey commonly used axioms for  
10 feature attribution. We demonstrate the efficacy of our model and the proposed  
11 influence scores on real and synthetic data.

## 12 1 Introduction

13 The field of causal inference studies causal links among random variables of interest, disentangling  
14 causal effects from simple statistical associations [25, 27]. For example, quantifying the causal effects  
15 of a medical treatment on patient outcomes concerns two primary random variables—treatment and  
16 outcome—potentially along with other covariates to consider. In causal discovery, the aim is to  
17 recover causal links among finitely many well-defined random variables from which a finite sample  
18 is observed. However, many causal questions in real world applications take on a different form  
19 that do not appeal to these descriptions. Many applications in *root cause analysis* comprise singular  
20 discrete events that unfold in time, and the objective is to recover causal links and chains among these  
21 individual events [40]. For example, in system administration and operations (recently, *AIOps*) it is  
22 often required to establish root causes of some adverse events such as failures and outages to other  
23 events in the data such as deployments and failures in dependencies. In the study of electronic health  
24 records, one may be interested in causally tracing changes in a patient’s trajectory to treatments. These  
25 examples can be viewed as establishing causal links among individual rare events unfolding in time.

26 In multivariate event streams, where events can be identified as members of finitely many types, these  
27 questions extend to whether one type of event Granger-causes another [1, 8]. However, there exists  
28 no framework for defining this problem in the language of causal discovery and for attribution of  
29 causal effects among *individual events* as opposed to *types of events*. We aim to address this problem  
30 in this work, paving the way to a unified and consistent methodology for identifying root causes in  
31 event streams.

32 In this paper, our objectives are twofold. We will first introduce a novel time series model for rare  
33 “event” data where occurrences of events will be represented as binary random variables in discrete  
34 time. Our model is inspired by the rich theory on temporal point processes (TPP) [7] and self-exciting  
35 point processes [12, 13]; and will serve to represent event data in a statistical framework that is  
36 amenable to causal analysis. We will then use this model to utilize the tools of time-series causal

37 inference and discovery, introducing two measures of causal contribution among events that are  
 38 analogous to causal effects as defined by Pearl[25]. Finally, we will link these quantities to existing  
 39 results on Hawkes processes and Granger causality. Together, our model and causal attribution  
 40 scores constitute a framework for fitting rare event processes and attributing causal influence among  
 41 individual events.

## 42 2 Preliminaries

43 **Causal Inference** Causal inference focuses on drawing causal conclusions from data, establishing  
 44 some variables as the *causes* of others as opposed to simply recovering associational relationships  
 45 (*i.e.*, correlations) among them. One formalism often used in causal inference is the *causal Bayesian*  
 46 *network* (CBN) [25, 30], which is a Bayesian network that obeys the causal Markov condition, *i.e.*,  
 47 that the joint distribution of random variables  $X_1, X_2, \dots, X_N$  decomposes as

$$P(x_1, x_2, \dots, x_N) = \prod_i P(x_i | PA_i = pa_i),$$

48 where  $PA_i$  denotes the set of variables  $X_j$  that are the parents of  $X_i$  in the CBN, and  
 49  $pa_i$  the corresponding random variates. Moreover, each conditional  $P(x_i | PA_i = pa_i)$  rep-  
 50 represents an independent *causal* mechanism. That is, to obtain the *interventional* distribution  
 51  $P(X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_{N-1} | do(X_j = x_j))$  it suffices to replace the term  $P(X_j | PA_j)$  with  
 52  $\delta_{X_j, x_j}$  where  $\delta_{a,b}$  denotes Kronecker’s delta. Note that this distribution is different from the condi-  
 53 tional that would result from simply observing that  $X_j = x_j$ , and corresponds to the distributions of  
 54  $\{X_i | i \neq j\}$  when  $X_j$  is actively determined (*i.e.*, intervened on).

55 **Structural Causal Models** While CBNs suffice to completely specify all possible interventional  
 56 distributions of a set of variables, a stricter formalism is needed to answer so called *counterfactual*  
 57 queries that allow answering “what if?” questions for individual observations. Under *structural causal*  
 58 *models* (SCMs, also referred to as functional causal or structural equation models), every variable  
 59 is written as a function of its parents and an unobserved noise variable,  $X_i = f_i(PA_i, U_i)$ , where  
 60  $U_i$  are statistically mutually independent.  $f_i(PA_i, U_i)$  is an SCM for the conditional  $P(X_i | PA_i)$  if  
 61  $f_i(pa_i, U_i)$  is distributed according to  $P(X_i | PA_i = pa_i)$  for almost all  $pa_i$ .

62 **Granger causality** Time ordering of data significantly facilitates reasoning about causal relations,  
 63 as causal effects can only act forward in time. However, variables measured *simultaneously* in time,  
 64 up to the temporal granularity available, still present an issue as the causal ordering among these  
 65 variables are not determined [27, Ch. 10]. However, in the absence of causal influence among  
 66 simultaneous measurements of variables, or so called *instantaneous effects*, causal influence can  
 67 be captured using the formalism of *Granger causality* [11]. Let  $\{\mathbf{X}_t\}_{t=1}^T$  denote a discrete-time  
 68 vector-valued stochastic process where  $\mathbf{X}_t = [X_t^{(1)}, X_t^{(2)}, \dots, X_t^{(d)}]$ . A time-series  $X^{(i)}$  is said to  
 69 Granger-cause another time series  $X^{(j)}$  if the past of  $X^{(i)}$  improves the predictions of  $X^{(j)}$  given all  
 70 past information about  $\{X^{(j')} | j' \neq i\}$ .

71 **Temporal Point Processes** A TPP specifies the full generative model for random sequences of  
 72 points  $(t_1, t_2, \dots, t_n)$  on a bounded subset of the real line, where  $0 < t_1 < t_2 < \dots \leq T$  and the  
 73 variable  $n$  is also random [7]. The *conditional intensity* function of the TPP

$$\lambda^*(t)dt = \mathbb{P}\{\text{next event is in } [t, t + dt) | \mathcal{H}_t\},$$

74 completely determines the process and is often used to characterize TPP models. Here  $\mathcal{H}_t$  denotes  
 75 the history (filtration) up to time  $t$ —specified by the set of points up to time  $t$ ,  $\{t_i | t_i < t\}$ . Intuitively,  
 76 the conditional intensity function specifies the arrival rate of events per unit time, in the infinitesimal  
 77 interval after  $t$ .

78 **Hawkes process** [13, 12]. A (univariate) Hawkes process is given by the conditional intensity<sup>1</sup>

$$\lambda^*(t) = \mu + \sum_{t_i < t} \alpha \varphi(t - t_i). \quad (1)$$

79 Here  $\mu > 0$  is a *background intensity*—the arrival rate of events if previous events had no effect on the  
 80 present. The *delay density*  $\varphi(s)$  determines the temporal profile of interactions between points—with  
 81  $\int \varphi(s)ds = 1$  w.l.o.g. Moreover, the function  $\varphi$  is always “causal,”  $\varphi(s) = 0, \forall s < 0$ , nonnegative  
 82  $\varphi(s) \geq 0$ , and is often monotonically decreasing over all  $s > 0$ . The parameter  $\alpha > 0$  is the so-called  
 83 *infectivity* or *branching* parameter.

84 Multivariate TPPs model *marks*  $y_i \in \{1, \dots, d\}$  for each event  $t_i$ . That is, *events* are now observed  
 85 as ordered pairs  $(t_i, y_i)$ . Practically,  $y_i$  often represent membership to an entity, such as a user on a  
 86 social network, host on a computer network, etc. The conditional intensity of the multivariate Hawkes  
 87 process (MHP) is written separately for each mark  $k$  as

$$\lambda_k^*(t) = \mu_k + \sum_m \sum_{t_i < t | y_i = m} \alpha_{km} \varphi(t - t_i).$$

88 Note that the background intensity of each mark is now different, and the infectivity parameters  
 89 can now be arranged along a matrix  $\mathbf{A}_{km} = \alpha_{km}$ , where each element describes the directional  
 90 *infectivity* of one mark over the other. Eichler *et al.* [8] show that  $\alpha_{km} > 0$  implies that the process  $m$   
 91 Granger-causes  $k$ ; while Achab *et al.* have shown how to recover  $\mathbf{A}$  via moment-matching estimators  
 92 [1]. While Hawkes processes are defined in continuous time ( $t_i \in \mathbb{R}$ ), in this paper we will explore  
 93 their discrete time analogues ( $t_i \in \mathbb{N}$ ) starting from the next section.

### 94 3 Quantifying Causal Contribution with Discrete-Time Hawkes 95 Processes

96 **Basic Observations** General TPPs model a wide range of occurrence patterns such as self-excitation  
 97 [12], self-inhibition [16], quasi-periodicity [6], etc. for discrete events in continuous time. However,  
 98 much of the established literature in causal inference deals with a finite set of random variables as  
 99 opposed to continuous time processes, leading to conceptual difficulties in analyzing cause-effect  
 100 relationships in continuous-time stochastic processes. Similarly, in many application domains time is  
 101 inherently *quantized*, *i.e.*, the data is sampled in discrete time—events can often “co-occur” with no  
 102 temporal ordering implied among them—and a continuous-time process serves as an approximation.  
 103 For example, neural spike trains are recorded with finite sampling rates, or many rare events in  
 104 computer systems logs are recorded in a predetermined time resolution. Therefore, in this section,  
 105 we start with the introduction of a discrete-time analogue of self-exciting temporal point processes  
 106 which will serve primarily to reconcile notation between causal inference and TPPs, as well as having  
 107 the added benefit of removing any statistical bias that results from using continuous-time models for  
 108 discrete data.

109 **Discrete-Time Hawkes Processes** In our formalism, the occurrences of “events”<sup>2</sup> or “points” are  
 110 interpreted as those times  $t \in \mathbb{Z}_{>0}$  of  $X_t$  where  $X_t = 1$ . Such models have been called discrete-time  
 111 point processes[38], such as in determinantal point processes [19] or discrete-time renewal processes  
 112 [9]. In addition to modeling discretely sampled events, our model builds on Hawkes processes to  
 113 model excitation patterns among them. We introduce the discrete-time Hawkes process (DTHP)  
 114 below.

115 **Definition 1.** (*Discrete-time Hawkes Process (DTHP)*) A binary-valued stochastic process  $\{X_t \in$   
 116  $\{0, 1\}\}_{t \in \mathbb{Z}_{>0}}$  is a discrete-time Hawkes process if, for all  $t$ ,

$$p_t := \mathbb{P}\{X_t = 1 | X_{1:t-1}\} = 1 - \exp\left(-\mu - \sum_{s=1}^{t-1} X_s g(t-s)\right).$$

117 where  $g(\tau) : \mathbb{Z} \rightarrow \mathbb{R}_{\geq 0}$  is a nonnegative function that satisfies  $g(\tau) = 0, \forall \tau < 0$ .

118 We observe that for all  $s < t$ ,  $\mathbb{E}[X_t | X_s = 1] > \mathbb{E}[X_t | X_s = 0]$ , therefore the process preserves the  
 119 self-excitation property of Hawkes processes, *i.e.*, that events only increase the probability of future  
 120 event occurrences.

<sup>2</sup>Not to be confused with the events of the underlying probability space, we reserve this term exclusively to refer to occurrences of 1 in a discrete-time binary process.

121 Our construction of DTHP admits the continuous-time Hawkes process as a limit case, *i.e.*, it tends  
 122 to a continuous-time Hawkes process as events become “infinitely” rare. In the same light, we can  
 123 examine a *rare event limit* or how the probability of events behaves as events become increasingly  
 124 rare. We will use these limits to derive approximations to the true causal effects that are expressed  
 125 simply in terms of the learned parameters of our model. Concretely, bounding the probability of  
 126 occurrence of points such that  $\forall t, p_t \leq \bar{p}$  we observe as  $\bar{p} \rightarrow 0$ , these probabilities also admit a linear  
 127 approximation in the effects of past points.

128 **Proposition 1.**  $p_t = \mu + \sum_{s=1}^{t-1} X_s g(t-s) + O(\bar{p}^2)$  as  $\bar{p} \rightarrow 0$ .

129 Another benefit of casting event occurrences in discrete time is that it enables the use of concepts from  
 130 traditional time-series analysis and the well-established literature of causal inference; specifically  
 131 causal Bayesian networks [25], and causality in time series [26, 27]. Moreover, in order to set our  
 132 new model in this framework, we can write an SCM that results in joint distributions equivalent to  
 133 the DTHP, which follows from observing that each  $X_t$  can be written as a function of an independent  
 134 source of noise and parent variables  $X_{1:t-1}$ .

135 **Definition 2. (DTHP SCM)** Let  $\{X_t\}$  are determined by the structural equations

$$X_t = \llbracket U_t \leq \lambda(X_{1:t-1}) \rrbracket \quad \text{where } \lambda(X_{1:t-1}) = \mu + \alpha \sum_{s<t} X_s g(t-s),$$

136  $U_t$  are independent standard exponential random variables and  $\llbracket \cdot \rrbracket$  denotes the indicator function.

137 Apart from rendering the mathematical objects conceptually simpler, DTHP enables using the  
 138 language of causal graphical models. Note that our choice of  $1 - \exp(-x)$  as a link function in  
 139 Definition 1 is one of many possible that would yield similar and tighter approximations. However,  
 140 for the purposes, this function suffices to demonstrate the key links between self-exciting point  
 141 processes and measures of causal contribution.

142 **Multivariate DTHP** We can now extend the DTHP to multivariate processes, where the interest is  
 143 in multiple related types of events.

144 **Definition 3. (Multivariate DTHP)** A binary vector valued process  $\mathbf{X}_t = (X_t^{(1)}, \dots, X_t^{(d)}) \in$   
 145  $\{0, 1\}^d$  is a multivariate DTHP if for all  $k, t$

$$p_t^{(k)} := \mathbb{P}\{X_t^{(k)} = 1 | \mathbf{X}_{1:t-1}\} = 1 - \exp\left(-\mu^{(k)} - \sum_{m=1}^d \sum_{s=1}^{t-1} X_s^{(m)} g_{m \rightarrow k}(t-s)\right).$$

146 Here,  $g_{m \rightarrow k}$  determine the decay profile of effects of events in type  $m$  on events of type  $k$ . In the  
 147 remainder of this paper, we will assume a more specific form for this quantity,  $g_{m \rightarrow k}(t-s) =$   
 148  $\mathbf{A}_{km} g(t-s)$  where  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and we assume  $\sum_{\tau=1}^{\infty} g(\tau) = 1$  without loss of generality.

149 We can also rely on previous results in time series causal discovery [27] to repeat a result similar to  
 150 those of [8] and [1] for DTHP.

151 **Proposition 2.** Events  $\{X_t^{(m)}\}$  Granger-cause events  $\{X_t^{(k)}\}$  if and only if  $\mathbf{A}_{km} > 0$ .

152 In the discrete-time world, however, we encounter another conceptual difficulty: continuous-time  
 153 TPPs are built on the *simplicity* assumption [7], that specify that no two points can co-occur on the  
 154 same point  $t' \in \mathbb{R}$  almost surely. This is a somewhat restrictive requirement in discrete time where  
 155 one may be interested in multiple types of points occurring together while being causally related, *i.e.*,  
 156 via *instantaneous* effects. Our formulation in Definition 3 disallows any such interactions between  
 157 ‘simultaneous’ variables  $X_t^{(m)}$  and  $X_t^{(k)}$ . In order to incorporate such effects for more realistic  
 158 modeling, we can extend the model as follows. For brevity, we denote

$$\lambda_t^{(k)} = \lambda^{(k)}(\mathbf{X}_{1:t-1}) := \sum_{m=1}^d \sum_{s=1}^{t-1} X_s^{(m)} \mathbf{A}_{km} g(t-s),$$

159 and define

$$p_t^{(k)} = 1 - \exp\left(-\lambda_t^{(k)} - \sum_{X^{(m)} \in PA_k^{(B)}} \mathbf{B}_{km} X_t^{(m)}\right),$$

160 where we define  $\mathbf{B} \in \mathbb{R}^{d \times d}$  as the weighted adjacency matrix of a graph that specifies the instanta-  
 161 neous causal effects among types of events, and  $PA_k^{(B)}$  to denote the set of parents of  $X^{(k)}$  along  
 162 this graph.

163 **Quantifying Causal Contribution** We can now build on the DTHP to introduce our method for  
 164 quantifying causal influence among observed events themselves. Specifically, we will focus on  
 165 quantifying causal contributions given a fitted multivariate DTHP model (Definition 3), where we will  
 166 currently ignore instantaneous effects for notational brevity. However, extensions of our arguments to  
 167 the case with instantaneous effects and implications for continuous-time Hawkes processes can be  
 168 derived from our framework.

169 Our problem can be formulated as follows. Given a finite realization of the process  $\{\mathbf{X}_t = \mathbf{x}_t\}_{t=1}^T$ ,  
 170 we seek to quantify the causal contribution of  $X_s^{(m)}$  on  $X_t^{(k)}$ , where  $s < t$ , and when such events are  
 171 “rare.” In practice, such quantification is only relevant when  $x_s^{(m)} = x_t^{(k)} = 1$  as, under our model,  
 172 events are assumed to be mutually exciting and we do not intuitively expect that the absence of an  
 173 event is the cause of another.

174 Such a notion of causal influence, of  $X_s^{(m)}$  on  $X_t^{(k)}$ , can be built on several familiar quantities in  
 175 causal inference. For example, one could consider the *average causal effect*  $\text{ACE}(X_s^{(m)} \rightarrow X_t^{(k)}) =$   
 176  $\mathbb{E} \left[ X_t^{(k)} \mid \text{do}(X_s^{(m)} = 1) \right] - \mathbb{E} \left[ X_t^{(k)} \mid \text{do}(X_s^{(m)} = 0) \right]$ , measuring the added probability of an event  
 177 on  $X_t^{(k)}$  when  $X_s^{(m)}$  is intervened on [15]. However, this quantity disregards the fact that the entire  
 178 history  $\mathbf{X}_{1:T}$  is observed. Moreover, ACE also does not take into account how (marginally) rare  
 179 the target event  $\{X_t^{(k)} = 1\}$  is. In this light, we define our first measure of causal influence on  
 180 a different quantity, the *direct effect* [25, Sec 4.5], which refers to the isolated effect of changing  
 181 only a single parent  $X_s^{(m)}$  having observed all other parents of  $X_t^{(k)}$ . We will denote this quantity  
 182  $\text{DE}(X_s^{(m)} \rightarrow X_t^{(k)})$ , defined

$$\mathbb{E} \left[ X_t^{(k)} \mid \text{do}(X_s^{(m)} = 1, \mathbf{X}_{-(m,s)} = \mathbf{x}_{-(m,s)}) \right] - \mathbb{E} \left[ X_t^{(k)} \mid \text{do}(X_s^{(m)} = 0, \mathbf{X}_{-(m,s)} = \mathbf{x}_{-(m,s)}) \right],$$

183 where the notation  $\mathbf{X}_{-(m,s)}$  is used to refer to all variables in the history except  $X_s^{(m)}$ .

184 We can now show that under the DTHP SCM and in the rare event regime, the direct effect yields a  
 185 convenient approximation. Namely,

186 **Proposition 3.**  $\text{DE}(X_s^{(m)} \rightarrow X_t^{(k)}) = \mathbf{A}_{km} g(t-s) + O(\bar{p}^2)$ .

187 This result links the proposed contribution measure to a well known quantity in the analysis of  
 188 Hawkes processes, namely the incremental intensity due to a previous event in the Hawkes process,  
 189 *i.e.*, the summand in  $\lambda^{(k)}(\mathbf{X}_{1:t-1})$  due to  $X_s^{(m)}$ .

190 The direct effect is based on a “total” intervention on all of the parents of  $X_t^{(k)}$ , comparing the  
 191 intervention where there is a source event at  $X_s^{(m)}$  to one where there is not. In this sense, it already  
 192 takes into account the full information available ( $\mathbf{x}_{1:t}$ ). However, it is still scaled in terms of the  
 193 marginal probability of  $\{X_t^{(k)} = 1\}$ . In order to quantify the proportion of influence of each past  
 194 event on a given target event, we can define a normalized quantity.

195 **Definition 4** (Normalized Direct Effect). *The normalized effect is defined*

$$\widetilde{\text{DE}}(X_s^{(m)} \rightarrow X_t^{(k)}) = \frac{\text{DE}(X_s^{(m)} \rightarrow X_t^{(k)})}{\lambda_k^t}.$$

196 Note that,  $\widetilde{\text{DE}}(X_s^{(m)} \rightarrow X_t^{(k)})$  is exactly equivalent to the posterior “parent” distribution in the  
 197 immigration-birth representation of Hawkes processes [14, 3]. Indeed, this representation of Hawkes  
 198 processes captures an intuitive notion of a causal chain of events. As previously indicated, we expect  
 199 that in an unconfounded system, the causes of events can only be (a combination of) other events,  
 200 but not the lack thereof. Similarly, we are more rarely interested in causal questions such as “what  
 201 previous event caused the lack of an event at time  $t$ ?” In this sense, the immigration-birth process

202 naturally captures an intuitive notion of causality among events. Our results show that the influence  
 203 of a direct cause, in the direct parenthood sense of a Hawkes process, is analogous to the direct effect  
 204 in causal inference. We describe this link in detail in Appendix B. Finally, we observe that for the  
 205 normalized direct effects to add to one, a summand  $\mu^{(k)}/\lambda_t^{(k)}$  is also required. This quantity can be  
 206 thought of as the probability that an event has no observed causal parent. In the immigration-birth  
 207 interpretation, the same quantity can be understood as the probability that an event is an “immigrant,”  
 208 and not a descendant of any previous events. Finally, the following result links Granger causality to  
 209 our approximate contribution measure  $\text{DE}(X_s^{(m)} \rightarrow X_t^{(k)})$ .

210 **Proposition 4.** *Assume  $\forall \tau, g(\tau) > 0$ . Then,  $\text{DE}(X_s^{(m)} \rightarrow X_t^{(k)}) > 0$  if and only if  $X^{(m)}$  Granger-*  
 211 *causes  $X^{(k)}$ .*

212 We can build on these observations to define a *total effect* of a single event at  $X_s^{(m)}$  on an event at  
 213  $X_t^{(k)}$ , by summing over all indirect paths of influence, weighted by their normalized direct effects.  
 214 In the following, let  $\mathcal{B}_{s,t}$  define the set of all points  $\{X_{s'}^{(k')} = 1 | s < s' < t, k' \in \{1, \dots, d\}\}$ , and  
 215  $\mathcal{P}_o(\mathcal{B}_{s,t})$  all ordered sets in the power set of  $\mathcal{B}_{s,t}$  such that temporal ordering is preserved. In other  
 216 words,

$$\mathcal{P}_o(\mathcal{B}_{s,t}) = \{(X_{s_1}^{k_1}, \dots, X_{s_n}^{k_n}) | \forall n \in [|\mathcal{B}_{s,t}|], s_i < s_{i+1} \forall i\}.$$

217 Note that the empty ordered set  $\emptyset \in \mathcal{P}_o(\mathcal{B}_{s,t})$ . For brevity, let us also define the path effect

$$\widetilde{\text{DE}}((X_{s_1}^{k_1}, \dots, X_{s_n}^{k_n})) := \prod_{i=1}^{n-1} \widetilde{\text{DE}}(X_{s_i}^{k_i} \rightarrow X_{s_{i+1}}^{k_{i+1}}).$$

218 We heuristically define the total effect as,

219 **Definition 5 (Total Effect).** *The total effect  $\text{TE}(X_s^{(m)} \rightarrow X_t^{(k)})$  is defined*

$$\sum_{\mathbf{Z} \in \mathcal{P}_o(\mathcal{B}_{s,t})} \widetilde{\text{DE}}((X_s^{(m)}, \mathbf{Z}, X_t^{(k)})), \quad (2)$$

220 *where we use the notation  $(X_s^{(m)}, \mathbf{Z}, X_t^{(k)})$  to denote the sequence generated by prepending (resp.*  
 221 *appending)  $X_s^{(m)}$  (resp.  $X_t^{(k)})$  to the sequence  $\mathbf{Z}$ .*

222 The simple intuition behind our definition is hidden away by the cumbersome notation required.  
 223 In other terms, the total effect captures the total influence an event has on a descendant, summing  
 224 over all paths of descendance—direct or indirect. While (2) seemingly requires summing over  
 225 exponentially many paths, its computation can be greatly accelerated via simple heuristics such as  
 226 dropping connections below a certain NDE.

227 Proposition 4 highlights that one event can be the cause of another in our sense of DE only if there  
 228 is a Granger-causality relationship between their marks. Note, however, that the same is not true  
 229 for our definition of total effects, where one mark can indirectly cause events in another mark. To  
 230 understand this relationship, and to contrast the two measures, assume a multivariate Hawkes process  
 231 of three marks is used to represent “delay” events of three consecutive trains, where the delay of the  
 232 first train directly causes a delay in the second, and a delay in the second causes a delay of the third  
 233 train. Using our measures of influence, and with perfect information, we will always attribute direct  
 234 causation to the previous train only. However, through total effects, we can attribute the third train’s  
 235 delay to that of the first. Finally note that, in the sense of Granger causality, the first train cannot be  
 236 said to cause the delay of the third train as, given knowledge of the second train, we cannot better  
 237 predict the delay of the third train. Although we will not make a rigorous argument in this work, the  
 238 DE measure, when viewed as an attribution method, readily satisfies the axioms of [33].

239 Finally, let us highlight that methods proposed in this section can be viewed as the parts of a single  
 240 framework. Given sparse event data that are sampled in discrete time and can be identified as one of  
 241 finitely many types, our framework only makes the additional assumption that past events will have  
 242 linear and additive (self-exciting) effects on future events. Under these assumptions, to identify the  
 243 causal effects among individual events we (i) fit a DTHP model to the observed sequence and (ii) use  
 244 DE, NDE and TE as measures of causal contribution to trace individual events to their causal parents.

## 245 4 Related Work

246 The Hawkes process has been studied commonly to establish Granger causality—i.e., causal links  
247 among different types of events as opposed to individual events. Eichler *et al.* explore the link  
248 between Hawkes process infectivity kernels and Granger causality [8]. Achab *et al.* use this link and  
249 previous results on moment-matching methods for Hawkes process estimation to introduce a fast  
250 algorithm for uncovering Granger causality [1]. Xu *et al.* consider group sparsity regularization for  
251 a more precise recovery of the Granger causal graph [39]. Notably, Prabhakar *et al.* introduced an  
252 algorithm for Granger-causal discovery directly from a cross-spectral estimate of multivariate TPPs,  
253 without making any parametric assumptions on the form of the conditional intensity [28]. We also  
254 refer the reader to [34] for a discussion of causal discovery in multitype event sequences. Many other  
255 works, which focus on more effective methods for recovering the infectivity matrix of a Hawkes  
256 process, can be seen as causal discovery algorithms in the context of multitype event sequences.  
257 Among these we can cite [21] who use an EM algorithm for better stability, [10] who work with more  
258 general transmission models, [35] who employ Bayesian inference for more accurate recovery of the  
259 graph, and [37] who employ low rank factorizations for improved scalability.

260 To our knowledge, a “discretized” Hawkes process appears only in [22], who allow each time step  
261 to have more than one points—i.e., work with time series of positive integers instead of binary  
262 sequences. Other discrete time point processes, towards recovering Granger-causal structure, have  
263 also been introduced in the context of neural structure learning [17].

264 Sun and Janzing study a similar form for causal discovery in arbitrary causal graphs of binary  
265 variables, although their setting is more general and their methods do not address temporal data [32].  
266 Similar to our framework, their probabilities of occurrence also admit linear approximations around  
267 0, although the authors do not explore this direction. In [5], Budhatoki *et al.* discuss methods for root  
268 cause analysis of outlier values, which could be regarded as rare events.

269 Recently, Tran *et al.* introduced QTree [36], a method that draws from extreme value theory and  
270 causal inference to infer graphs (more specifically, root-directed trees) of causal influence among  
271 nodes where *simultaneous* outlier events occur jointly. The max-linear Bayesian network model  
272 used is able to handle missing values as well as infer graphs of influence among network nodes  
273 in a robust fashion. Moreover, the authors employ the Chu-Liu-Edmonds algorithm for minimum  
274 cost arborescence to heuristically recover root-directed trees, as required by their application in  
275 uncovering hidden river networks. CAUSE, by Zhang *et al.*, is the closest to our work [40]. Here,  
276 the authors consider an axiomatic causal attribution method that obeys the axioms of [33]. Notably,  
277 the method considered attributes causal influence among events, using an “explainable” recurrent  
278 point process—a neural TPP model. The authors then show that an aggregation of these influence  
279 scores can be interpreted as a measure for Granger causality among event marks. However, the  
280 neural network-based model used and the attribution methods make computation under this method  
281 prohibitively costly. Finally, in “counterfactual” TPPs [23], Noorbakhsh and Gomez-Rodriguez,  
282 describe a structural causal model analogue of Lewis’ thinning algorithm which they then use to  
283 answer counterfactual queries in observed point sequences.

## 284 5 EXPERIMENTS

285 **Model Performance** We start by validating the performance of DTHP on three data sets for the  
286 task of inferring the latent network of influence among event types—the first step of the framework  
287 we propose. We compare the performance of DTHP with two baselines: QTree [36] and CAUSE  
288 [40]. To contrast these baselines with ours, QTree is able to handle missing values and can work with  
289 general real-valued variables to infer both general graphs and trees of influence. However, QTree  
290 only works with instantaneous effects, i.e., assumes that each time step is i.i.d. CAUSE [40] is based  
291 on neural TPPs and does not consider instantaneous effects. Both algorithms are developed for causal  
292 discovery in sequences of rare events. For both baselines, we use repositories made available by the  
293 authors and keep the original hyperparameters included in the libraries.<sup>3</sup>

<sup>3</sup>see <https://github.com/razhangwei/CAUSE>,  
<https://github.com/princengoc/QTree>.

Metric Model	AUC			F1		
	QTree	CAUSE	DTHP (ours)	QTree	CAUSE	DTHP (ours)
hawkes-1	0.248	0.431	<b>0.830</b>	0.114	0.286	<b>0.471</b>
hawkes-2	0.563	0.610	<b>0.765</b>	0.265	0.254	<b>0.467</b>
hawkes-3	0.563	0.536	<b>0.736</b>	0.255	0.242	<b>0.424</b>
danube	<b>0.897</b>	0.628	0.841	<b>0.800</b>	0.118	0.308
lower-colorado	<b>0.712</b>	0.639	0.701	<b>0.450</b>	0.200	0.214
middle-colorado	<b>0.951</b>	0.734	0.563	<b>0.909</b>	0.286	0.235
upper-colorado	<b>0.931</b>	0.570	0.660	<b>0.875</b>	0.267	0.333
Connectomics-1	0.499	0.525	<b>0.623</b>	0.185	0.186	<b>0.234</b>
Connectomics-2	0.519	0.514	<b>0.639</b>	0.179	0.178	<b>0.243</b>
Connectomics-3	0.590	0.508	<b>0.670</b>	0.206	0.175	<b>0.267</b>
Connectomics-4	0.702	0.527	<b>0.738</b>	0.301	0.185	<b>0.348</b>
Connectomics-5	0.730	0.515	<b>0.745</b>	0.320	0.186	<b>0.357</b>
Connectomics-6	0.859	0.715	<b>0.880</b>	0.487	0.307	<b>0.545</b>

Table 1: Experiment results comparing QTree, CAUSE, and DTHP algorithms given in AUC and maximum F1 scores (higher better). Top scores in each row are given in bold.

294 The objective of all experiments is the recovery of an underlying causal graph from observed time  
295 series. We use three groups of data sets, the first of which is simulated, and the others taken from  
296 real applications. Further details on synthetic data generation and benchmark data sets are given in  
297 Appendix C.

- 298 • We simulate data from **continuous-time multivariate Hawkes processes** using tick [2].
- 299 • The **River Basin Data Sets** include data collected from two river basins in Europe and the  
300 US [36], for the so-called *hidden river* discovery task. We experiment with four data sets,  
301 belonging to the Danube, as well as lower, middle and upper sections of the Lower Colorado  
302 river basin.
- 303 • We use the neural connectome data set from the **Chalearn Connectomics** challenge [4]. The  
304 data set includes realistically generated spike trains from neuronal networks [31] Specifically,  
305 we perform experiments on the **small** data sets which are numbered in increasing order from  
306 the most challenging setting to the least.

307 All data sets have ground truth causal networks available. For the river basin data sets, we use raw  
308 measurements in the QTree algorithm, however threshold the data to convert it into binary time series  
309 for use in CAUSE and DTHP models. For the neural connectome data sets, we threshold each data  
310 set at the 99th percentile, obtaining binary time series used in all of the algorithms. We use both  
311 versions of the QTree algorithm, with and without the minimum cost arborescence step, and report  
312 the best results. As CAUSE is designed for continuous time data sets, we “dequantize” binary time  
313 series by adding random noise drawn from a uniform distribution between 0 and 1 to each timestamp.

314 Our results are presented in Table 1. We report the area under the ROC curve (AUC) and maximum  
315 attained F1 score for edge classification. As expected, our model is significantly superior in the  
316 Hawkes process data sets, and the QTree algorithm dominates in the river basin data sets where it was  
317 designed to perform well. Both CAUSE and DTHP perform significantly below the QTree baseline  
318 in the river data sets. We believe this is primarily due to two reasons. First, neither model performs  
319 Bayesian treatment of missing values which is especially important in the Lower Colorado river basin  
320 data sets. Second, these algorithms do not search for the best tree with minimum cost arborescence.  
321 Let us note, however, that running the Chu-Liu-Edmonds algorithm alone on graphs recovered by  
322 DTHP and CAUSE also did not yield significantly better results. Still, DTHP appears to perform  
323 slightly more favorably than CAUSE, which does not address instantaneous effects.

324 In the Connectomics experiments, we find that our algorithm significantly outperforms baselines.  
325 This matches our expectation as the Connectomics data set is both high-dimensional (100 marks),  
326 and features both delayed and instantaneous effects. Our model is the only one designed to capture  
327 all such patterns simultaneously. Overall, we can conclude that DTHP generally yields favorable  
328 performance in modeling sparse binary time series where instantaneous effects occur.

329 **Causal Influence Scores** For a demonstration of our causal influence scores, we present a set of  
330 experiments on synthetic data. Here, our aim is to first exhibit the general difficulty of attributing



$d$	$p_E$	True	Fitted
5	0.05	$0.735 \pm 0.156$	$0.604 \pm 0.156$
	0.1	$0.717 \pm 0.062$	$0.641 \pm 0.081$
	0.2	$0.672 \pm 0.069$	$0.607 \pm 0.070$
	0.5	$0.503 \pm 0.063$	$0.473 \pm 0.058$
10	0.05	$0.826 \pm 0.097$	$0.600 \pm 0.135$
	0.1	$0.793 \pm 0.046$	$0.600 \pm 0.077$
	0.2	$0.722 \pm 0.045$	$0.569 \pm 0.047$
	0.5	$0.528 \pm 0.038$	$0.449 \pm 0.033$

Table 2: Comparison of results when retrieving the Hawkes process parent event using normalized direct effect. Numbers reported are means and standard deviations of recall at top 1—*i.e.*, among events with known parents, the ratio of those with the top NDE score assigned to the correct parent.  $d$  denotes the dimensionality of the Hawkes process, and  $p_E$  is the prior for sparsity. Higher  $p_E$  implies lower sparsity.

331 causal influence among rare events, even with perfect information. To this end, we draw from  
 332 multivariate Hawkes processes while keeping record of the parents of each event. We regard these  
 333 parenthood relationships as the ground truth causes of events, and measure if the direct effects  
 334 computed as per Definition 4 correctly recover the causes. We consider only those events that have a  
 335 parent in the branching process, and compute recall (at top 1).

336 Results, for varying dimensionality and degrees of sparsity in the infectivity matrix, are presented  
 337 in Table 2. Here, we observe that direct effects computed with known parameters already fall to an  
 338 accuracy of around 50% when 50% of the edges in the ground truth graph are active—highlighting  
 339 a general ambiguity with assigning causes among events when many such causes are possible.  
 340 Moreover, we find that causal attribution with fitted parameters (“Fitted”) performs slightly worse  
 341 than when ground truth parameters are known (“True”), but also that it is relatively robust. However,  
 342 as expected, robustness decreases when dimensionality is increased. Further details are available in  
 343 Appendix C.

## 344 6 CONCLUSION

345 In this paper we introduced a framework for attributing causal influence among individual events  
 346 observed in time. Assuming only that events are sparse and there is a quasi-linear and monotonic  
 347 relationship among their probabilities of occurrence our method proceeds by fitting a newly introduced  
 348 discrete-time process model, and performing causal attribution via simple quantities based on the  
 349 fitted parameters of this model Our analysis was cast in a discrete-time framework, enabling unbiased  
 350 estimation in many real-world scenarios where data is sampled with finite rates and instantaneous  
 351 effects are also present. Finally, our numerical experiments validate the efficacy of our model for  
 352 the unique scenarios it addresses, as well as the intuition behind the causal contribution metrics  
 353 we proposed in this work. While our method can address many discrete event scenarios, its main  
 354 limitation is that it only allows excitation relationships among events. Several directions remain as  
 355 next steps to our work, such as extending the model with real-valued marks and inhibitory effects to  
 356 address more general sparse discrete event sequences.

## References

- 358 [1] Massil Achab, Emmanuel Bacry, Stéphane Gaïffas, Iacopo Mastromatteo, and Jean-François  
359 Muzy. Uncovering causality from multivariate Hawkes integrated cumulants. *The Journal of*  
360 *Machine Learning Research*, 18(1):6998–7025, 2017.
- 361 [2] Emmanuel Bacry, Martin Bompaire, Philip Deegan, Stéphane Gaïffas, and Søren Poulsen.  
362 tick: a python library for statistical learning, with an emphasis on hawkes processes and  
363 time-dependent models. *J. Mach. Learn. Res.*, 18(1):7937–7941, 2017.
- 364 [3] Emmanuel Bacry, Iacopo Mastromatteo, and Jean-François Muzy. Hawkes processes in finance.  
365 *Market Microstructure and Liquidity*, 1(01):1550005, 2015.
- 366 [4] Demian Battaglia, Isabelle Guyon, Vincent Lemaire, Javier Orlandi, Bisakha Ray, and Jordi  
367 Soriano. *Neural connectomics challenge*. Springer, 2017.
- 368 [5] Kailash Budhathoki, Lenon Minorics, Patrick Blöbaum, and Dominik Janzing. Causal structure-  
369 based root cause analysis of outliers. In *International Conference on Machine Learning*, pages  
370 2357–2369. PMLR, 2022.
- 371 [6] David Roxbee Cox. *Renewal theory*. Methuen, 1962.
- 372 [7] Daryl J. Daley and David Vere-Jones. *An introduction to the theory of point processes: Volume*  
373 *I: elementary theory and methods*. Springer Science & Business Media, 2007.
- 374 [8] Michael Eichler, Rainer Dahlhaus, and Johannes Dueck. Graphical modeling for multivariate  
375 hawkes processes with nonparametric link functions. *Journal of Time Series Analysis*, 38(2):225–  
376 242, 2017.
- 377 [9] Willliam Feller. *An introduction to probability theory and its applications*. John Wiley & Sons,  
378 1957.
- 379 [10] Manuel Gomez Rodriguez, David Balduzzi, and Bernhard Schölkopf. Uncovering the temporal  
380 dynamics of diffusion networks. In *Proceedings of the 28th International Conference on*  
381 *Machine Learning*, 2011.
- 382 [11] Clive WJ Granger. Investigating causal relations by econometric models and cross-spectral  
383 methods. *Econometrica: journal of the Econometric Society*, pages 424–438, 1969.
- 384 [12] Alan G. Hawkes. Point spectra of some mutually exciting point processes. *Journal of the Royal*  
385 *Statistical Society. Series B (Methodological)*, pages 438–443, 1971.
- 386 [13] Alan G. Hawkes. Spectra of some self-exciting and mutually exciting point processes.  
387 *Biometrika*, 58(1):83–90, 1971.
- 388 [14] Alan G. Hawkes and David Oakes. A cluster process representation of a self-exciting process.  
389 *Journal of Applied Probability*, 11(3):493–503, 1974.
- 390 [15] Paul W Holland. Statistics and causal inference. *Journal of the American statistical Association*,  
391 81(396):945–960, 1986.
- 392 [16] Valerie Isham and Mark Westcott. A self-correcting point process. *Stochastic processes and*  
393 *their applications*, 8(3):335–347, 1979.
- 394 [17] Sanggyun Kim, David Putrino, Soumya Ghosh, and Emery N Brown. A granger causality  
395 measure for point process models of ensemble neural spiking activity. *PLoS computational*  
396 *biology*, 7(3):e1001110, 2011.
- 397 [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
398 *arXiv:1412.6980*, 2014.
- 399 [19] Alex Kulesza, Ben Taskar, et al. Determinantal point processes for machine learning. *Founda-*  
400 *tions and Trends in Machine Learning*, 5(2–3):123–286, 2012.
- 401 [20] Patrick J. Laub, Thomas Taimre, and Philip K. Pollett. Hawkes Processes. *arXiv:1507.02822*  
402 *[math, q-fin, stat]*, July 2015. arXiv: 1507.02822.
- 403 [21] Scott Linderman and Ryan Adams. Discovering latent network structure in point process data.  
404 In *International Conference on Machine Learning*, pages 1413–1421, 2014.
- 405 [22] Scott W. Linderman and Ryan P. Adams. Scalable bayesian inference for excitatory point  
406 process networks. *arXiv preprint arXiv:1507.03228*, 2015.

- 407 [23] Kimia Noorbakhsh and Manuel Gomez Rodriguez. Counterfactual temporal point processes. In  
408 *Neural Information Processing Systems*, 2019.
- 409 [24] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
410 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative  
411 style, high-performance deep learning library. *Advances in neural information processing*  
412 *systems*, 32:8026–8037, 2019.
- 413 [25] Judea Pearl. *Causality: Models, reasoning, and inference*. Cambridge University Press, 2000.
- 414 [26] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Causal inference on time series using  
415 restricted structural equation models. In *Advances in Neural Information Processing Systems*,  
416 pages 154–162, 2013.
- 417 [27] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: founda-*  
418 *tions and learning algorithms*. The MIT Press, 2017.
- 419 [28] Karthir Prabhakar, Sangmin Oh, Ping Wang, Gregory D Abowd, and James M Rehg. Temporal  
420 causality for the analysis of visual events. In *2010 IEEE Computer Society Conference on*  
421 *Computer Vision and Pattern Recognition*, pages 1967–1974. IEEE, 2010.
- 422 [29] Aleksandr Simma and Michael I. Jordan. Modeling events with cascades of Poisson processes.  
423 *arXiv preprint arXiv:1203.3516*, 2012.
- 424 [30] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction,*  
425 *and search*. MIT press, 2000.
- 426 [31] Olav Stetter, Demian Battaglia, Jordi Soriano, and Theo Geisel. Model-free reconstruction of  
427 excitatory neuronal connectivity from calcium imaging signals. *PLoS Computational Biology*,  
428 8(8), 2012.
- 429 [32] Xiaohai Sun and Dominik Janzing. Exploring the causal order of binary variables via exponential  
430 hierarchies of markov kernels. In *15th European Symposium on Artificial Neural Networks*  
431 *(ESANN 2007)*, pages 465–470. D-Side, 2007.
- 432 [33] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
433 *International Conference on Machine Learning*, pages 3319–3328. PMLR, 2017.
- 434 [34] Nikolaj Theodor Thams. Causal structure learning in multivariate point processes. Master’s  
435 thesis, University of Copenhagen, 2019.
- 436 [35] Long Tran, Mehrdad Farajtabar, Le Song, and Hongyuan Zha. Netcodec: Community detection  
437 from individual activities. In *Proceedings of the 2015 SIAM International Conference on Data*  
438 *Mining*, pages 91–99. SIAM, 2015.
- 439 [36] Ngoc Mai Tran, Johannes Buck, and Claudia Klüppelberg. Causal discovery of a river network  
440 from its extremes. *arXiv preprint arXiv:2102.06197*, 2021.
- 441 [37] Ali Caner Türkmen, Gökhan Çapan, and Ali Taylan Cemgil. Clustering event streams with low  
442 rank hawkes processes. *IEEE Signal Processing Letters*, 27:1575–1579, 2020.
- 443 [38] Ali Caner Türkmen, Tim Januschowski, Yuyang Wang, and Ali Taylan Cemgil. Forecasting in-  
444 termittent and sparse time series: A unified probabilistic framework via deep renewal processes.  
445 *Plos one*, 16(11):e0259764, 2021.
- 446 [39] Hongteng Xu, Mehrdad Farajtabar, and Hongyuan Zha. Learning granger causality for hawkes  
447 processes. In *International Conference on Machine Learning*, pages 1717–1726. PMLR, 2016.
- 448 [40] Wei Zhang, Thomas Panum, Somesh Jha, Prasad Chalasani, and David Page. Cause: Learning  
449 granger causality from event sequences using attribution methods. In *International Conference*  
450 *on Machine Learning*, pages 11235–11245. PMLR, 2020.

## 451 **A Proofs of Propositions**

### 452 **Proposition 1**

453 *Proof.* Let  $\lambda_t := \mu + \sum_{s=1}^{t-1} X_s g(t-s)$ . Note the Taylor series approximation of  $p_t$  around 0 is,

$$\bar{p} \geq p_t = \lambda_t - \frac{\lambda_t^2}{2} + O(\lambda_t^3).$$

454 and also note that  $\bar{p}^2 \sim \frac{\lambda_t^2}{2}$  as  $\bar{p} \rightarrow 0$ . Therefore  $p_t = \lambda_t + O(\bar{p}^2)$ .  $\square$

### 455 **Proposition 2**

456 *Proof.* We will follow the arguments of [27, Theorem 10.3], assuming causal sufficiency (as required  
457 by Granger causality in general). By definition, there exists a link between  $X^{(m)}$  and  $X^{(k)}$  in the  
458 *summary graph* only when there exists a link from  $X_s^{(m)}$  to  $X_t^{(k)}$  for some  $s < t$ . However, by  
459 definition of Hawkes SCM (Definition 2, extended analogously to the multivariate case), such a link  
460 only exists if  $\mathbf{A}_{km} > 0$ .  $\square$

### 461 **Proposition 3**

462 *Proof.* Let

$$\lambda_t^{k,0} := \mu^{(k)} + \sum_{m', s' < t | (m', s') \neq (m, s)} x_{s'}^{(m')} \mathbf{A}_{km'} g(t-s'). \quad (3)$$

463 It follows from Proposition 1 that

$$\begin{aligned} \text{DE}(X_s^{(m)} \rightarrow X_t^{(k)}) &= \mathbb{E} \left[ X_t^{(k)} \mid \text{do}(X_s^{(m)} = 1, \mathbf{X}_{-(m,s)} = \mathbf{x}_{-(m,s)}) \right] \\ &\quad - \mathbb{E} \left[ X_t^{(k)} \mid \text{do}(X_s^{(m)} = 0, \mathbf{X}_{-(m,s)} = \mathbf{x}_{-(m,s)}) \right], \\ &= \lambda_t^{k,0} + \mathbf{A}_{km} g(t-s) + O(\bar{p}^2) - (\lambda_t^{k,0} + O(\bar{p}^2)) \\ &= \mathbf{A}_{km} g(t-s) + O(\bar{p}^2). \end{aligned}$$

464  $\square$

### 465 **Proposition 4**

466 *Proof.* From Proposition 3 and 2, this immediately holds for an approximation of direct effects  
467  $\text{DE}(X_s^{(m)} \rightarrow X_t^{(k)}) \approx \mathbf{A}_{km} g(t-s)$ . To see that it also holds exactly, let  $\lambda_t^{k,0}$  be defined as in (3)  
468 and note that

$$\begin{aligned} \text{DE}(X_s^{(m)} \rightarrow X_t^{(k)}) &= \mathbb{E} \left[ X_t^{(k)} \mid \text{do}(X_s^{(m)} = 1, \mathbf{X}_{-(m,s)} = \mathbf{x}_{-(m,s)}) \right] \\ &\quad - \mathbb{E} \left[ X_t^{(k)} \mid \text{do}(X_s^{(m)} = 0, \mathbf{X}_{-(m,s)} = \mathbf{x}_{-(m,s)}) \right], \\ &= \exp(-\lambda_t^{k,0}) - \exp(-\lambda_t^{k,0} - \mathbf{A}_{km} g(t-s)), \end{aligned}$$

469 from where it is apparent that  $A_{km} = 0$  implies  $\text{DE}(X_s^{(m)} \rightarrow X_t^{(k)}) = 0, \forall s, t$ . Conversely, assuming  
470  $g(t-s) > 0, A_{km} = 0$  implies  $\text{DE}(X_s^{(m)} \rightarrow X_t^{(k)}) = 0$  completing the proof.  $\square$

## 471 B Equivalence to Hawkes’ Branching Process Interpretation

472 Owing to the convenient additive form of its intensity, the Hawkes process lends itself to interpretation  
 473 as a Poisson-cluster process, or an infinite cascade of Poisson processes. This description of the  
 474 process is sometimes intuitively called an *immigration-birth* or *branching* representation[14, 7].  
 475 Below, we describe a *new* generative process, one which does not rely on the conditional intensity  
 476 as in (1). Here, individual points will be denoted as ordered pairs  $(s_n, z_n)$  where  $s_n$  denotes the  
 477 timestamp, and  $z_n$  the timestamp of the parent event which gave *birth* to the point at  $s_n$ .

- 478 1. Draw  $N_0 \sim \text{Poisson}(\mu \times T)$ . Let  $\mathcal{D}_0 = \{(s_i, 0)\}_{i=1}^{N_0}$  where  $s_i$  are drawn uniformly at  
 479 random in  $(0, T]$ . These points are the so-called *immigrants*.
- 480 2. For each generation  $j$ , starting from  $j = 1$  we draw the children of each point in the previous  
 481 generation.
  - 482 • Letting  $\mathcal{D}_{j-1} = \{(s_i^{(j-1)}, z_i^{(j-1)})\}$ , draw  $N_{s_i^{(j-1)}} \sim \mathcal{PO}(\alpha)$  for each  $s_i^{(j-1)}$
  - 483 • Let
 
$$484 \mathcal{D}_j^{s_i^{(j-1)}} = \{(\tau_k + s_i^{(j-1)}, s_i^{(j-1)})\}_{k=1}^{N_{s_i^{(j-1)}}},$$
 where we draw  $\tau_k \sim g$  i.i.d.
    - 485 • Let  $\mathcal{D}_j = \bigcup_{s_i^{(j-1)}} \mathcal{D}_j^{s_i^{(j-1)}}$ .
    - 486 • Stop if there exist no  $(s_i^{(j)}, z_i^{(j)}) \in \mathcal{D}_j$  such that  $s_i^{(j)} \leq T$ .
- 487 3. Return  $\mathcal{D} = \{(s_i, z_i) \in \bigcup_j \mathcal{D}_j | s_i \leq T\}$ .

488 Somewhat surprisingly, due to the Poisson superposition property, this process is equivalent to the  
 489 process determined by the conditional intensity function of (1). Moreover, if one uses this method  
 490 of generating a Hawkes draw, an auxiliary *parenthood* variable,  $z_i$  which refers to the (timestamp  
 491 of) point which “gave birth” to it, s.t.  $z_i < s_i$  always holds. Moreover, if these parenthood variables  
 492 were known from the beginning, optimal parameters  $\{\mu, \alpha, g\}$  could be recovered in a closed-form  
 493 maximization step since they would just be parameters of iid Poisson process observations.

494 The discarded parenthood variables  $z_i$  define a *forest* of immigrants (root nodes) and their descendants.  
 495 It is this observation that underlies the EM algorithm for Hawkes processes [14, 3, 29, 20], which  
 496 proceeds by (E) inferring the parent of each variable (computing  $\mathbb{P}\{z_i = s_j\}$  where  $s_j < z_i$ ), and (M)  
 497 maximizing  $\{\mu, \alpha, g\}$  under the expected complete data likelihood. By consulting [29], for example,  
 498 one can see our approximate normalized direct effect (for the univariate case)  $\alpha g(t_i - t_j)/\lambda_t$  appears  
 499 as the “posterior” probability  $\mathbb{P}\{z_i = t_j\}$ . While our exposition here is concerned only with the  
 500 univariate Hawkes process, its extensions to multivariate processes follow easily.

501 Using the same statistical foundation as above we can now argue that our approximated normalized  
 502 direct effects coherently describe a graph where each node is a point and each edge is weighted by the  
 503 probability of parenthood. In this formalism, our definition of the total effect also appears as the total  
 504 path weight where a path weight is defined as the product of the weights of edges it is composed of.

## 505 C Further Details on Experiments

### 506 C.1 Model Performance

507 **Generated Hawkes processes** Data sets are generated with the `SimuHawkesExpKernels` class  
 508 provided in tick [2]. Namely, we generate infectivity matrices  $\mathbf{A} = \mathbf{W} \odot \mathbf{Y}$  where  $\mathbf{A} \in \mathbb{R}^{d \times d} \odot$   
 509 denotes the Hadamard product,  $\mathbf{W}_{km} \stackrel{iid}{\sim} \text{Exp}(1)$ , and  $\mathbf{Y} \stackrel{iid}{\sim} \text{Bernoulli}(0.1)$ . We then adjust the  
 510 spectral radius of the matrix to  $\rho$ . We set the baseline intensities  $\mu_k = 0.05$ , and the maximum  
 511 number of jumps to 5000. The three data sets `hawkes-1`, `hawkes-2`, and `hawkes-3` are sampled with  
 512 parameters  $(\rho, d) = (0.5, 10), (0.4, 20), (0.3, 30)$  ranked from least to most challenging respectively.  
 513 We then binarize these data sets by quantizing time along the unit grid and setting a time interval to 1

514 if the interval contains a sampled point. The resulting data sets have points in 5.5%, 8.2%, and 6.8%  
515 of intervals respectively.

516 **Lower Colorado River Basin Data Sets** Except for use in the QTree algorithm, the data sets are  
517 preprocessed by binarizing at the 0.99-quantile and filling missing values with 0.

518 **Connectomics Data Set** We first preprocess the data set by taking the first difference of the raw  
519 action potentials. Except for QTree, we binarize the data by setting a cutoff at at the 99th percentile.  
520 In practice, this percentile is also close to the recommended binarization cutoff, 0.12.

521 **Baselines and Hyperparameters** We use the `ExplainableRecurrentPointProcess` class from  
522 the CAUSE library, and use the default hyperparameters as defined in the training script. By default,  
523 the model uses a hidden layer size of 64, embedding dimension of 64, batch size of 64, no dropout or  
524 L2 regularization, learning rate of 0.001, 200 epochs and the Adam optimizer. We use the `QTree`  
525 class of the QTree library, leaving default hyperparameters `smallR = 0.05`, `q = 0.8`.

526 **Implementation of DTHP** We use our own implementation for the DTHP model, using PyTorch  
527 [24]. We implement maximum likelihood optimization for the proposed discrete time model, with  
528 added regularization for the graph such that the total loss function is

$$\ell(\mu, \theta, \mathbf{A}) = \log \sum_{k,t} p(X_t^{(k)} | \mathcal{H}_t, \mu, \theta, \mathbf{A}) + \gamma \|\mathbf{A}\|_F.$$

529 In our experiments, we heuristically set  $\gamma = 10$ . We use the implementation of the Adam optimizer  
530 [18] implemented in PyTorch for optimization, setting the learning rate to 0.01. We train for 10K  
531 epochs on the Connectomics data set, and 5K epochs on the other data sets. In practice, we truncate  
532 the history of each point where influences can flow to a certain maximum history, and set this value  
533 to 1 in the river data sets and 5 in the synthetic and connectome data sets.

## 534 C.2 Causal Influence

535 For the causal influence estimation experiments, we generate infectivity matrices  $\mathbf{A} = \mathbf{U} \odot \mathbf{Y}$  where  
536  $\mathbf{A} \in \mathbb{R}^{d \times d}$ ,  $\mathbf{U}_{km} \stackrel{iid}{\sim} \text{Uniform}[0, 1]$ , and  $\mathbf{Y}_{km} \stackrel{iid}{\sim} \text{Bernoulli}(p_E)$ , set  $\mu = (2d)^{-1}$ , and  $\theta = 0.33$ . We  
537 use our own implementation of a Hawkes process branching sampler to draw from a Hawkes process  
538 while retaining the parent identifiers  $z_i$  as explained in Appendix B.

539 For experiments where the infectivity matrix  $\mathbf{A}$  is estimated (denoted “Fitted” in the results), we run  
540 DTHP setting maximum lag to 5 and the number of epochs to 3K.