

# Efficient Non-parametric Bayesian Hawkes Processes

Rui Zhang<sup>1,2</sup>, Christian Walder<sup>1,2</sup>, Marian-Andrei Rizoiu<sup>3</sup> and Lexing Xie<sup>1,2</sup>

<sup>1</sup>The Australian National University, Australia

<sup>2</sup>Data61 CSIRO, Australia

<sup>3</sup>University of Technology Sydney, Australia

{firstname}. {lastname}@[anu.edu.au<sup>1</sup>, data61.csiro.au<sup>2</sup>, uts.edu.au<sup>3</sup>]

## Abstract

In this paper, we develop an efficient non-parametric Bayesian estimation of the kernel function of Hawkes processes. The non-parametric Bayesian approach is important because it provides flexible Hawkes kernels and quantifies their uncertainty. Our method is based on the cluster representation of Hawkes processes. Utilizing the stationarity of the Hawkes process, we efficiently sample random branching structures and thus, we split the Hawkes process into clusters of Poisson processes. We derive two algorithms — a block Gibbs sampler and a maximum a posteriori estimator based on expectation maximization — and we show that our methods have a linear time complexity, both theoretically and empirically. On synthetic data, we show our methods to be able to infer flexible Hawkes triggering kernels. On two large-scale Twitter diffusion datasets, we show that our methods outperform the current state-of-the-art in goodness-of-fit and that the time complexity is linear in the size of the dataset. We also observe that on diffusions related to online videos, the learned kernels reflect the perceived longevity for different content types such as music or pets videos.

## 1 Introduction

The Hawkes process [Hawkes, 1971] is a useful model of self-exciting point data in which the occurrence of a point increases the likelihood of arrival of new points. More specifically, every point causes the conditional intensity function  $\lambda$  — which modulates the arrival rate of new points — to increase. An alternative representation of the Hawkes process is a cluster of Poisson processes [Hawkes and Oakes, 1974], which categorizes points into *immigrants* and *offspring*. Immigrant points are generated independently at a background rate  $\mu$ ; offspring points are triggered by existing points at a rate of  $\phi$ . Points can therefore be structured into clusters, where each cluster contains a point and the offspring it directly generated. This leads to a tree structure, also known as the branching structure (an example is shown in Fig. 1).

**Background & Motivations**  $\phi$  is important as it is shared and decides the class of the whole process and recently the

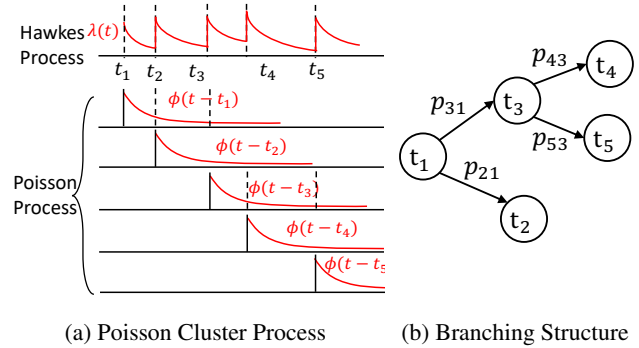


Figure 1: The cluster Representation of a Hawkes Process. **(a)** A Hawkes process with decaying triggering kernel  $\phi(\cdot)$  has intensity  $\lambda(t)$  which increases after each new point is generated. It can be represented as a cluster of Poisson processes  $PP(\phi(t - t_i))$  associated with each  $t_i$ . **(b)** The branching structure corresponding to the triggering relationships shown in (a), where an edge  $t_i \rightarrow t_j$  means that  $t_i$  triggered  $t_j$  with the probability  $p_{ji}$  (calculated as Eq. (9)).

Hawkes process with various  $\phi$  has been studied. Mishra *et al.* [2016] employ the branching factor of the Hawkes process with the power-law kernel to predict popularity of tweets; Kurashima *et al.* [2018] predict human actions using a Hawkes process equipped with exponential, Weibull and Gaussian mixture kernels; online popularity unpredictability is explained using the Hawkes process with a variant of the exponential kernel by Rizoiu *et al.* [2018]. However, most work employs Hawkes process with parametric kernels, which encodes strong assumptions, and limits the expressivity of the model. **Can we design a practical approach to learn flexible representations of the optimal Hawkes kernel function  $\phi$  from data?**

A typical solution is the non-parametric estimation of the kernel function [Lewis and Mohler, 2011; Zhou *et al.*, 2013b; Bacry and Muzy, 2014; Eichler *et al.*, 2017]. These are all frequentist methods which do not quantify the uncertainty of the learned kernels. There exists work [Rasmussen, 2013; Linderman and Adams, 2015; Rousseau *et al.*, 2018] on the Bayesian inference for the Hawkes process. To scale past toy-dataset sizes these methods require either parametric triggering kernels or discretization of the input domain, which in turn leads to poor scaling with the dimension of the domain

Table 1: Related Works in Non-parametric Hawkes Processes.

Methods	Time Complexity	Bayesian	Continuous	Non-parametric
Zhou <i>et al.</i> [2013a]	$O(n^3)$	×	✓	×
Xu <i>et al.</i> [2016]	$O(n^3)$	×	✓	×
Lewis and Mohler [2011]	$O(n^3)$	×	×	✓
Zhou <i>et al.</i> [2013b]	$O(n^3)$	×	×	✓
Rasmussen [2013]	$O(n)$	✓	✓	×
Linderman and Adams [2015]	$O(n)$	✓	interval-sensored	×
Rousseau <i>et al.</i> [2018]	unspecified	✓	✓	✓
Ours	$O(n)$	✓	✓	✓

and sensitivity to the choice of discretization. The work closest to our own is that of Rousseau *et al.* [2018], however their main contributions are theoretical; on the practical side they resort to an unscalable Markov chain Monte Carlo (MCMC) estimator. We comparatively summarize related works in Table 1. To the best of our knowledge, our work is the first work proposing a Bayesian non-parametric Hawkes process estimation procedure, with a linear time complexity allowing it to be applied to real-world datasets, and without requiring discretization of domains.

**Contributions** In this paper, we propose a general framework for the efficient non-parametric Bayesian inference of Hawkes processes.

(1) We exploit block Gibbs sampling [Ishwaran and James, 2001] to iteratively sample the latent branching structure, the background intensity  $\mu$  and the triggering kernel  $\phi$ . In each iteration, the point data are decomposed as a cluster of Poisson processes based on the sampled branching structure. This is exemplified in Fig. 1, in which a Hawkes process (shown on the top temporal axis of Fig. 1a) is decomposed into several Poisson processes (the following temporal axes); the corresponding branching structure is shown in Fig. 1b. The posterior  $\mu$  and  $\phi$  are estimated using the resulting cluster processes. Our framework is close to the stochastic Expectation-Maximization (EM) algorithm [Celeux and Diebolt, 1985] where posterior  $\mu$  and  $\phi$  are estimated [Lloyd *et al.*, 2015; Walder and Bishop, 2017; Gugushvili *et al.*, 2018] in the M-step and random samples of  $\mu$  and  $\phi$  are drawn. We adapt the approach of the recent non-parametric Bayesian estimation for Poisson process intensities, termed Laplace Bayesian Poisson process (LBPP) [Walder and Bishop, 2017], to estimate the posterior  $\phi$  given the sampled branching structure.

(2) We utilize the stationarity of the Hawkes Process to speed up sampling and computing the probability of the branching structure. We theoretically show our method to be of linear time complexity. Furthermore, we explore the connection with the EM algorithm [Dempster *et al.*, 1977] and develop a second variant of our method, as an approximate EM algorithm.

(3) We empirically show our method enjoys linear time complexity and can infer known analytical kernels, i.e., exponential and sinusoidal kernels. On two large-scale social media datasets, our method outperforms the current state-of-the-art non-parametric approaches and the learned kernels reflect the perceived longevity for different content types.

## 2 Preliminaries

In this section, we introduce the prerequisites of our work: the Hawkes process and LBPP.

**The Hawkes process [Hawkes, 1971].** Introduced in Section 1, the Hawkes process can be specified using the conditional intensity function  $\lambda$  which modulates the arrival rate of points. Mathematically, conditioned on a set of points  $\{t_i\}_{i=1}^N$ , the intensity  $\lambda$  is expressed as:

$$\lambda(t) = \mu + \sum_{t_i < t} \phi(t - t_i), \quad (1)$$

where  $\mu > 0$ , considered as a constant, and  $\phi(\cdot) : \mathbb{R} \rightarrow [0, \infty)$  are the background immigrant intensity and the triggering kernel. The log-likelihood of  $\{t_i\}_{i=1}^N$  given  $\mu$  and  $\phi(\cdot)$  is [Rubin, 1972]:

$$\log p(\{t_i\}_{i=1}^N | \mu, \phi(\cdot)) = \sum_{i=1}^N \log \lambda(t_i) - \int_{\Omega} \lambda(t) dt, \quad (2)$$

where  $\Omega$  is the sampling domain of  $\{t_i\}_{i=1}^N$ .

**Laplace Bayesian Poisson process (LBPP) [Walder and Bishop, 2017]** has been proposed for the non-parametric Bayesian estimation of the intensity of a Poisson process. To satisfy non-negativity of the intensity function, LBPP models the intensity function  $\lambda$  as a permanental process [Shirai and Takahashi, 2003], i.e.,  $\lambda = g \circ f$  where the link function  $g(z) = z^2/2$  and  $f(\cdot)$  obeys a Gaussian process (GP) prior. Alternative link functions include  $\exp(\cdot)$  [Møller *et al.*, 1998; Diggle *et al.*, 2013] and  $g(z) = \lambda^*(1 + \exp(-z))^{-1}$  [Adams *et al.*, 2009] where  $\lambda^*$  is constant.

The choice  $g(z) = z^2/2$  has the analytical advantages; for some covariances the log-likelihood can be computed in closed form [Lloyd *et al.*, 2015; Flaxman *et al.*, 2017]. LBPP exploits the Mercer expansion [Mercer, 1909] of the GP covariance function  $k(x, y) \equiv \text{Cov}(f(x), f(y))$ ,

$$k(x, y) = \sum_{i=1}^K \lambda_i e_i(x) e_i(y), \quad (3)$$

where for non-degenerate kernels,  $K = \infty$ . The eigenfunctions  $\{e_i(\cdot)\}_i$  are chosen to be orthonormal in  $L^2(\Omega, m)$  for some sample space  $\Omega$  with measure  $m$ .  $f(\cdot)$  can be represented as a linear combination of  $e_i(\cdot)$  [Rasmussen and Williams, 2005],  $f(\cdot) = \omega^T e(\cdot)$ , and  $\omega$  has a Gaussian prior, i.e.,  $\omega \sim \mathcal{N}(0, \Lambda)$  where  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_K)$  is a diagonal covariance matrix and  $e(\cdot) = [e_1(\cdot), \dots, e_K(\cdot)]^T$  is

a vector of basis functions. Computing the posterior distribution of the intensity function  $\lambda(\cdot)$  is equivalent to estimating the posterior distribution of  $\omega$  which, in LBPP, is approximated by a normal distribution (a.k.a Laplace approximation [Rasmussen and Williams, 2005]), i.e.,

$$\log p(\omega|X, \Omega, k) \approx \log \mathcal{N}(\omega|\hat{\omega}, Q), \quad (4)$$

where  $X \equiv \{t_i\}_{i=1}^N$  is a set of point data,  $\Omega$  the sample space and  $k$  the Gaussian process kernel function.  $\hat{\omega}$  is selected as the mode of the true posterior and  $Q$  the negative inverse Hessian of the true posterior at  $\hat{\omega}$ :

$$\hat{\omega} = \underset{\omega}{\operatorname{argmax}} \log p(\omega|X, \Omega, k), \quad (5)$$

$$Q^{-1} = - \left. \frac{\partial^2}{\partial \omega \partial \omega^T} \log p(\omega|X, \Omega, k) \right|_{\omega=\hat{\omega}}. \quad (6)$$

The approximate posterior distribution of  $f(t)$  is a normal distribution [Rasmussen and Williams, 2005]:

$$f(t) \sim \mathcal{N}(\hat{\omega}^T e(t), e(t)^T Q e(t)) \equiv \mathcal{N}(\nu, \sigma^2). \quad (7)$$

Furthermore, the posterior of  $\lambda(t)$  is a Gamma distribution:

$$\text{Gamma}(x|\alpha, \beta) \equiv \beta^\alpha x^{\alpha-1} e^{-\beta x} / \Gamma(\alpha), \quad (8)$$

where  $\alpha = (\nu^2 + \sigma^2)^2 / (4\nu^2\sigma^2 + 2\sigma^4)$  and  $\beta = (\nu^2 + \sigma^2) / (2\nu^2\sigma^2 + \sigma^4)$ .

### 3 Inference via Sampling

We now detail our efficient non-parametric Bayesian estimation algorithm, which samples the posterior of  $\mu$  (constant background intensity) and  $\phi(\cdot)$  (the triggering kernel). Our method starts with random  $\mu_0, \phi_0(\cdot)$  and iterates by cycling through the following four steps ( $k$  is the iteration index):

- i Calculate  $p(\mathcal{B}|X, \phi_{k-1}, \mu_{k-1})$ , the distribution of the branching structure  $\mathcal{B}$  given the data  $X$ , triggering kernel  $\phi_{k-1}$ , and background intensity  $\mu_{k-1}$  (see Section 3.1).
- ii Sample a  $\mathcal{B}_k$  as per  $p(\mathcal{B}|X, \phi_{k-1}, \mu_{k-1})$  (see Section 3.1).
- iii Estimate  $p(\phi|\mathcal{B}_k, X)$  (Section 3.3) and  $p(\mu|\mathcal{B}_k, X)$  (Section 3.2).
- iv Sample a  $\phi_k(\cdot)$  and  $\mu_k$  from  $p(\phi(\cdot)|\mathcal{B}_k, X)$  and  $p(\mu|\mathcal{B}_k, X)$ , respectively.

By standard Gibbs sampling arguments, the samples of  $\phi(\cdot)$  and  $\mu$  drawn in the step (iv) converge to the desired posterior, modulo the Laplace approximation in (iii). As the method is based on block Gibbs sampling [Ishwaran and James, 2001], we term it *Gibbs-Hawkes*.

#### 3.1 Conditional Distribution and Sampling of the Branching Structure

The branching structure  $\mathcal{B}$  has a data structure of tree (as Fig. 1(b)) and consists of independent triggering events. Thus, the probability of the branching structure  $\mathcal{B}$  is the product of probabilities of triggering events, i.e.,  $p(\mathcal{B}) = \prod_{i=1}^N p_{ij}$  where  $p_{ij}$  is the probability of  $t_j$  triggering  $t_i$ . Given  $\mu$  and  $\phi(\cdot)$ ,  $p_{ij}$  is the ratio between  $\phi(t_i - t_j)$  and  $\lambda(t_i)$  (see e.g. [Lewis and Mohler, 2011]):

$$p_{ij} \equiv \phi(t_i - t_j) / (\mu + \sum_{t_k < t_i} \phi(t_i - t_k)), \quad j \geq 1, \quad (9)$$

Similarly, the probability of point  $t_i$  being from  $\mu$ , say  $p_{i0}$ , is:

$$p_{i0} \equiv \mu / (\mu + \sum_{t_k < t_i} \phi(t_i - t_k)). \quad (10)$$

Given these probabilities we may sample a branching structure by sampling a parent for each  $t_i$  according to probabilities  $\{p_{ij}\}_{j \geq 0}$ . The sampled branching structure separates a set of points into immigrants and offspring (introduced in Section 1). Immigrants can be regarded as a sequence generated from  $\text{PP}(\mu)$ , where  $\text{PP}(\cdot)$  is a Poisson process which has an intensity  $\mu$ , and used to estimate the posterior  $\mu$ .

The key property which we exploit in the subsequent Section 3.2 and Section 3.3 is the following. Denote by  $\{t_k^{(i)}\}_{k=1}^{N_{t_i}}$  the  $N_{t_i}$  offspring generated by point  $t_i$ . If such a sequence is *aligned* to an origin at  $t_i$ , yielding  $S_{t_i} \equiv \{t_k^{(i)} - t_i\}_{k=1}^{N_{t_i}}$ , then the aligned sequence is drawn from  $\text{PP}(\phi(\cdot))$  over  $[0, T - t_i]$  where  $[0, T]$  is the sample domain of the Hawkes process. The posterior distribution of  $\phi(\cdot)$  is estimated on all such aligned sequences.

#### 3.2 Posterior Distribution of $\mu$

Continuing from the observations in Section 3.1, note that if we are given a set of points  $\{t_i\}_{i=1}^N$  generated by  $\text{PP}(\mu)$  over  $\Omega = [0, T]$ , the likelihood for  $\{t_i\}_{i=1}^N$  is the Poisson likelihood,  $p(\{t_i\}_{i=1}^N | \mu, \Omega) = e^{-\mu T} (\mu T)^N / N!$ . For simplicity, we place a conjugate (Gamma) prior on  $\mu T$ ,  $\mu T \sim \text{Gamma}(\alpha, \beta)$ ; the Gamma-Poisson conjugate family conveniently gives the posterior distribution of  $\mu T$ , i.e.,  $p(\mu T | \{t_i\}_{i=1}^N, \alpha, \beta) = \text{Gamma}(\alpha + N, \beta + 1)$ . We choose the scale  $\alpha$  and the rate  $\beta$  in the Gamma prior by making the mean of the Gamma posterior equal to  $N$  and the variance  $N/2$ , which is easily shown to correspond to  $\alpha = N$  and  $\beta = 1$ . Finally, due to conjugacy we obtain the posterior

$$p(\mu | \{t_i\}_{i=1}^N, \alpha, \beta) = \text{Gamma}(2N, 2T). \quad (11)$$

#### 3.3 Posterior Distribution of $\phi(\cdot)$

We handle the posterior distribution of the triggering kernel  $\phi(\cdot)$  given the branching structure in an analogous manner to the LBPP method of Walder and Bishop [2017]. That is, we assume that  $\phi(\cdot) = f^2(\cdot)/2$  where  $f(\cdot)$  is Gaussian process distributed as described in Section 2. In line with [Walder and Bishop, 2017], we consider the sample domain  $[0, \pi]$  and the so-called *cosine kernel*,

$$k(x, y) = \sum_{\gamma \geq 0} \lambda_\gamma e_\gamma(x) e_\gamma(y), \quad (12)$$

$$\lambda_\gamma \equiv 1 / (a(\gamma^2)^m + b), \quad (13)$$

$$e_\gamma(x) \equiv (2/\pi)^{1/2} \sqrt{1/2}^{[\gamma=0]} \cos(\gamma x). \quad (14)$$

Here,  $\gamma$  is a multi-index with non-negative (integral) values,  $[\cdot]$  is the indicator function,  $a$  and  $b$  are parameters controlling the prior smoothness, and we let  $m = 2$ . This basis is orthonormal w.r.t. the Lebesgue measure on  $\Omega = [0, \pi]$ . The expansion Eq. (12) is an explicit kernel construction based on the Mercer expansion as per Eq. (3), but other kernels may be used, for example by Nyström approximation [Flaxman *et al.*, 2017] of the Mercer decomposition.

As mentioned at the end of Section 3.1, by conditioning on the branching structure we may estimate  $\phi(\cdot)$  by considering

the *aligned* sequences. In particular, letting  $S_{t_i}$  denote the aligned sequence generated by  $t_i$ , the joint distribution of  $\omega$  and  $\{S_{t_i}\}_i$  is calculated as [Walder and Bishop, 2017]

$$\begin{aligned} & \log p(\omega, \{S_{t_i}\}_i | \Omega, k) \\ &= \sum_i \left\{ \sum_{\Delta t \in S_{t_i}} \log \frac{1}{2} (\omega^T e(\Delta t))^2 - \frac{1}{2} \int_0^{T-t_i} (\omega^T e(t))^2 dt \right\} \\ & \quad - \frac{1}{2} \log [(2\pi)^K |\Lambda|] - \frac{1}{2} \omega^T \Lambda^{-1} \omega. \end{aligned} \quad (15)$$

Note that there is a subtle but important difference between the integral term above and that of Walder and Bishop [2017], namely the limit of integration; closed-form expressions for the present case are provided in the online supplement [app, 2019, annex A]. Putting the above equation into Eq. (5) and Eq. (6), and we obtain the mean  $\hat{\omega}$  and the covariance  $Q$  of the (Laplace) approximate log-posterior in  $\omega$ :

$$\hat{\omega} = \underset{\omega}{\operatorname{argmax}} \log p(\omega, \{S_{t_i}\}_i | \Omega, k) \quad (16)$$

$$Q^{-1} = - \sum_i \left\{ \sum_{\Delta t \in S_{t_i}} \frac{2e(\Delta t)e(\Delta t)^T}{(\hat{\omega}^T e(\Delta t))^2} - \int_0^{T-t_i} e(t)e(t)^T dt \right\} + \Lambda^{-1}. \quad (17)$$

Then, the posterior  $\phi$  is achieved by Eqs. (7) and (8).

### 3.4 Computational Complexity

For LBPP, constructing Eq. (15) and Eq. (17) takes  $O(N_o K^2)$  where  $K$  is the number of basis functions and  $N_o$  is the number of offspring. Optimizing  $\omega$  (Eq. (16)) is a concave problem, which can be solved efficiently. If L-BFGS is used,  $O(CK)$  will be taken to calculate the gradient on each  $\omega$  where  $C$  is the number of steps stored in memory. Computing  $Q$  requires inverting a  $K \times K$  matrix, which is  $O(K^3)$ . As a result, the complexity of estimating  $\phi|B$  is  $O((N_o + K)K^2)$ . In terms of estimating  $\mu|B$  taking  $O(1)$ , the complexity of estimating  $\mu, \phi|B$  is linear to the number of data. While the naive complexity for  $p_{ij}$  is  $O(N^2)$ , Halpin [2012] provides an optimized approach to reduce it to  $O(N)$ , which relies on the stationary of Hawkes processes. In the step of sampling branching structures, points with negligible impacts on another point are not sampled as its parents. Interestingly, in comparison with LBPP, while our model is in some sense more complex, it enjoys a more favorable computational complexity. In summary, we have the following complexities per iteration and in Section 5, we validate the complexity on both synthetic and real data.

Operation	Complexity
$\mu B$	$O(1)$
$p_{ij}$	$O(N)$
$\phi B$	$O((N_o + K)K^2)$
overall	$O((N + K)K^2)$

## 4 Maximum-A-Posterior Estimation

We explore a connection between the sampler of section 3 and the EM algorithm, which allows us to introduce an analogous approximate *maximum a posteriori* (M.A.P.) scheme.

**Relationship to EM** In Section 1 we mentioned the connection between our method and the stochastic EM algorithm

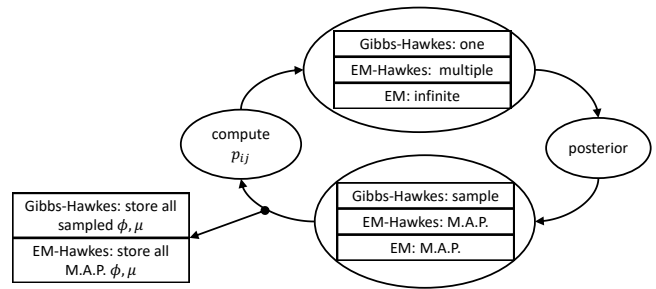


Figure 2: A visual summary of the Gibbs-Hawkes, EM-Hawkes and the EM algorithms. The differences between them are (1) the number of sampled branching structures and (2) selected  $\phi$  and  $\mu$  for  $p_{ij}$ . In contrast with with Gibbs-Hawkes, the EM-Hawkes method draws multiple branching structures at once and calculates  $p_{ij}$  using M.A.P.  $\phi$  and  $\mu$ . The EM algorithm is equivalent to sampling infinite branching structures and exploiting M.A.P. or constrained M.L.E.  $\phi$  and  $\mu$  to calculate  $p_{ij}$  (see Section 4).

[Celeux and Diebolt, 1985]. The difference is in the M-step; to perform EM [Dempster *et al.*, 1977] we need only modify our sampler by: (a) sampling infinite branching structures at each iteration, and (b) re-calculating the probability of the branching structure with the M.A.P.  $\mu$  and  $\phi(\cdot)$ , given the infinite set of branching structures. More specifically, maximizing the expected log posterior distribution to estimate M.A.P.  $\mu$  and  $\phi(\cdot)$  given infinite branching structures is equivalent to maximizing the EM objective in the M-step (see the online supplement [app, 2019, annex B] for the formal derivation). Finally, note that the above step (b) is identical to the E-step of the EM algorithm.

**EM-Hawkes** Following the discussion above, we propose *EM-Hawkes*, an approximate EM algorithm variant of Gibbs-Hawkes proposed in Section 3. Specifically, at each iteration EM-Hawkes (a) samples a finite number of cluster assignments (to approximate the expected log posterior distribution), and (b) finds the M.A.P. triggering kernels and background intensities rather than sampling them as per block Gibbs sampling (the M-step of the EM algorithm). An overview of the Gibbs-Hawkes, EM-Hawkes and EM algorithm is illustrated in Fig. 2.

Note that under our LBPP-like posterior, finding the most likely triggering kernel  $\phi(\cdot)$  is intractable (see the online supplement [app, 2019, annex C]). As an approximation we take the element-wise mode of the *marginals* of the  $\phi(t_i)$  to approximate the mode of the joint distribution of the  $\phi(t_i)$ .

## 5 Experiments

We now evaluate our proposed approaches — Gibbs-Hawkes and EM-Hawkes — and compare them to three baseline models, on synthetic data and on two large Twitter online diffusion datasets. The three baselines are:

- A **naive parametric Hawkes** equipped with a constant background intensity and an exponential (Exp) triggering kernel  $\phi = a_1 a_2 \exp(-a_2 t)$ ,  $a_1, a_2 > 0$ , estimated by maximum likelihood.

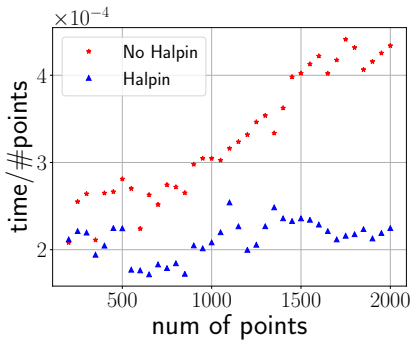


Figure 3: Computation time for calculating  $p_{ij}$  and sampling branching structures, with and without Halpin’s speed up.

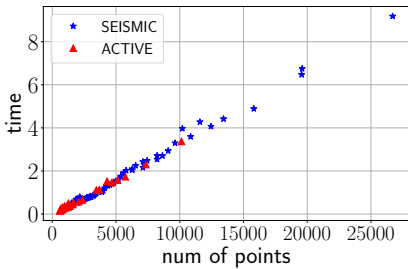


Figure 4: Running time per iteration on ACTIVE and SEISMIC.

- **Ordinary differential equation (ODE)**-based non-parametric non-bayesian Hawkes [Zhou *et al.*, 2013b]. The code is publicly available [Bacry *et al.*, 2017].
- **Wiener-Hopf (WH)** equation based non-parametric non-bayesian Hawkes [Bacry and Muzy, 2014]. The code is publicly available [Bacry *et al.*, 2017].

## 5.1 Synthetic Data

We employ two toy Hawkes processes to generate data, both having the same background intensity  $\mu = 10$ , and cosine (Eq. (18)) and exponential (Eq. (19)) triggering kernels respectively:

$$\phi_{\cos}(t) = \cos(3\pi t) + 1, \quad t \in [0, 1]; \quad 0, \quad \text{otherwise}; \quad (18)$$

$$\phi_{\exp}(t) = 5\exp(-5t), \quad t > 0. \quad (19)$$

**Prediction** For three baseline models and EM-Hawkes, the predictions  $\mu_{\text{pred}}$  and  $\phi_{\text{pred}}(\cdot)$  are taken to be the M.A.P. values, while for Gibbs-Hawkes we use the posterior mean.

**Evaluation** Each toy model generates 400 point sequences over  $\Omega = [0, \pi]$ , which are evenly split into 40 groups, 20 for training and 20 for test. Each of the three methods fit on each group, *i.e.*, summing log-likelihoods for 10 sequences (for the parametric Hawkes) or estimating the log posterior probability of the Hawkes process given 10 sequences (for Gibbs-Hawkes and EM-Hawkes) or fitting the superposition of 10 sequences [Xu *et al.*, 2018]. Since the true models are known, we evaluate fitting results using the L2 distance between predicted and true  $\mu$  and  $\phi(\cdot)$ :

$$d_{L2}(g_{\text{pred}}, g_{\text{true}}) = \left( \int_{\Omega} (g_{\text{pred}}(t) - g_{\text{true}}(t))^2 dt \right)^{1/2}. \quad (20)$$

Table 2: Empirical performance comparison between algorithms (columns) with different measures (rows). *Top*: L2 distance to known  $\phi$  and  $\mu$ , *bottom*: mean predictive log likelihood on real data.

Data	Exp	ODE	WH	Gibbs	EM
$\phi_{\cos}$	0.809	0.677	1.225	0.414	<b>0.390</b>
$\mu_{\cos}$	<b>1.229</b>	1.262	30.826	1.381	2.109
$\phi_{\exp}$	<b>0.189</b>	0.965	1.581	0.235	0.221
$\mu_{\exp}$	<b>1.516</b>	5.471	82.078	1.818	3.617
activeYT	2.369	2.370	1.315	2.580	<b>2.592</b>
SEISMIC	3.335	3.357	2.131	3.576	<b>3.578</b>

**Experimental details** For Gibbs-Hawkes and EM-Hawkes, we must select parameters of the GP kernel (Eqs. (12) to (14)). Having many basis functions leads to a high fitting accuracy, but low speed. We found that using 32 basis functions provides a suitable balance. For kernel parameters  $a, b$  of Eq. (13), we choose  $a, b = 0.002$ . 5000 iterations are run to fit each group and first 1000 are ignored (*i.e.* *burned-in*).

**Results** The top of Table 2 shows the mean L2 distance between the learned and the true  $\phi(\cdot)$  and  $\mu$  on toy data. Gibbs-Hawkes and EM-Hawkes get closest triggering kernels to  $\phi_{\cos}$ ; naturally, the parametric Hawkes – which uses an exponential kernel – fits  $\phi_{\exp}$  best. The parametric model retrieves  $\mu$  slightly better on both synthetic datasets. The learned triggering kernels for  $\phi_{\exp}$  and  $\phi_{\cos}$  are shown in Fig. 5a and in the online supplement [app, 2019]. The ODE-based method performs well on  $(\mu_{\cos}, \phi_{\cos})$  but badly on  $(\mu_{\exp}, \phi_{\exp})$ . Notably, tuning the hyper-parameters of the WH method is challenging, and Table 2 shows the best result obtained after a rather exhaustive experimentation. In summary, compared with state-of-the-art methods, our approaches achieve better performances for data generated by kernels from several parametric classes; as expected the parametric models are only effective for data generated from their own class.

**Effect of Halpin’s Procedure** In Section 3.4 we show that using Halpin’s procedure reduces the complexity of calculating  $p_{ij}$  from quadratic to linear. We now empirically validate this speed up. To distinguish between quadratic and linear complexity, we compute the ratio between running time and data size, as shown in Fig. 3. The ratio when using Halpin’s procedure remains roughly constant as data size increases (the ratio increases linearly without the optimization), which implies that Halpin’s procedure renders linear calculation of  $p_{ij}$  and of branching structures. Later, we will show the linear complexity of our method on real data.

## 5.2 Twitter Diffusion Data

We evaluate the performance of our two proposed approaches (Gibbs-Hawkes and EM-Hawkes) on two Twitter datasets, containing retweet cascades. A retweet cascade contains an original tweet, together with its direct and indirect retweets. Current state of the art diffusion modeling approaches [Zhao *et al.*, 2015; Mishra *et al.*, 2016; Rizoio *et al.*, 2018] are based on the self-exciting assumption: users get in contact with on-

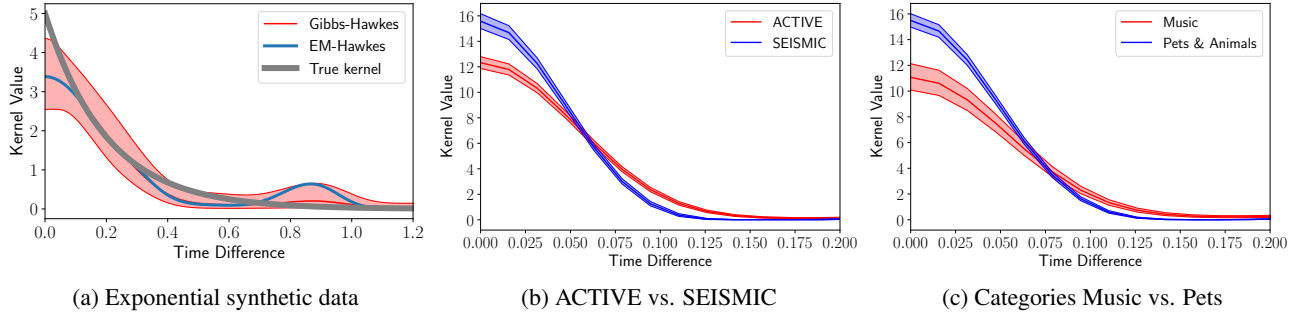


Figure 5: Learned Hawkes triggering kernels using our non-parametric Bayesian approaches. Each red or blue area shows the estimated posterior distributions of  $\phi$ , while the solid lines indicate the 10, 50 and 90 percentiles. (a) a synthetic dataset simulated using  $\phi_{\text{exp}}(t)$  (Eq. (19)), shown in gray, is fit using Gibbs-Hawkes (in red) and EM-Hawkes (in blue); (b) Twitter data in ACTIVE (in red) and SEISMIC (in blue); (c) Twitter data associated with two categories in the ACTIVE set: Music (in red) and Pets & Animals (in blue)

line content, and then diffuse it to their friends, therefore generating a cascading effect. The two datasets we use have been employed in prior work and they are publicly available:

- ACTIVE [Rizoio *et al.*, 2018] contains 41k retweet cascades, each containing at least 20 (re)tweets with links to Youtube videos. It was collected in 2014 and each Youtube video (and therefore each cascade) is associated with a Youtube category, e.g., *Music* or *News*.
- SEISMIC [Zhao *et al.*, 2015] contains 166k randomly sampled retweet cascades, collected in from Oct 7 to Nov 7, 2011. Each cascade contains at least 50 tweets.

**Setup** The temporal extent of each cascade is scaled to  $[0, \pi]$ , and assigned to either training or test data with equal probability. We bundle together groups of 30 cascades of similar size, and we estimate one Hawkes process for each bundle. Unlike for the synthetic dataset, for the retweet cascades dataset there is no *true* Hawkes process to evaluate against. Instead, we measure using log-likelihood how well the learned model generalizes to the test set. We use the same hyper-parameters values as for the synthetic data.

**Fitting Performance** For each dataset, we calculate the log-likelihood per event for each tweet cascade obtained by three baselines and our approaches (Table 2). Visibly, our proposed methods consistently outperform baselines, with EM-Hawkes performing slightly better than Gibbs-Hawkes (by 0.6% for ACTIVE and 0.4% for SEISMIC). This seems to indicate that online diffusion is influenced by factors not captured by the parametric kernel, therefore justifying the need to learn the Hawkes kernels non-parametrically. As mentioned in the synthetic data part, the WH-based method has a disadvantage of hard-to-tune hyper-parameters, which leads to the worst performance among all methods.

**Scalability** To validate the linear complexity of our method, we record running time per iteration of Gibbs-Hawkes on ACTIVE and SEISMIC in Fig. 4. The running time rises linearly with the number of points increasing, in line with the theoretical analysis. Linear complexity makes our method scalable and applicable on large datasets.

**Interpretation** We show in Fig. 5a the learned kernels for information diffusions. We notice that the learned kernels appear to be decaying and long-tailed, in accordance with the prior literature. Fig. 5b shows that the kernel learned on SEISMIC is decaying faster than the kernel learned on ACTIVE. This indicates that non-specific (i.e. random) cascades have a faster decay than video-related cascades, presumably due to the fact that Youtube videos stay longer in the human attention. This connection between the type of content and the speed of the decay seems further confirmed in Fig. 5c, where we show the learned kernels for two categories in ACTIVE: *Music* and *Pets & Animals*. Cascades relating to *Pets & Animals* have a faster decaying kernel than *Music*, most likely because Music is an ever-green content.

## 6 Conclusions

In this paper, we provided the first non-parametric Bayesian inference procedure for the Hawkes process which requires no discretization of the input domain and enjoys a linear time complexity. Our method iterates between two steps. First, it samples the branching structure, effectively transforming the Hawkes process into a cluster of Poisson processes. Next, it estimates the Hawkes triggering kernel using a non-parametric Bayesian estimation of the intensity of the cluster Poisson processes. We provide both a full posterior sampler and an EM estimation algorithm based on our ideas. We demonstrated our approach can infer flexible triggering kernels on simulated data. On two large Twitter diffusion datasets, our method outperforms the state-of-the-art in held-out likelihood. Moreover, the learned non-parametric kernel reflects the intuitive longevity of different types of content. The linear complexity of our approach is corroborated on both the synthetic and real problems. The present framework is limited to the univariate unmarked Hawkes process and will be extended to marked multivariate Hawkes process.

## Acknowledgement

This research was supported in part by the Australian Government through the Australian Research Council’s Discovery Projects funding scheme (project DP180101985).

## References

- R. P. Adams, I. Murray, and D. JC MacKay. Tractable non-parametric bayesian inference in poisson processes with gaussian process intensities. In *ICML*, 2009.
- Appendix: Efficient non-parametric bayesian hawkes processes, 2019.
- E. Bacry and J.-F. Muzy. Second order statistics characterization of hawkes processes and non-parametric estimation. *arXiv preprint arXiv:1401.0903*, 2014.
- E. Bacry, M. Bompierre, S. Gaïffas, and S. Poulsen. tick: a Python library for statistical learning, with a particular emphasis on time-dependent modeling. *ArXiv e-prints*, 2017.
- G. Celeux and J. Diebolt. The sem algorithm: A probabilistic teacher algorithm derived from the EM algorithm for the mixture problem. *Comp. Stat. Quarterly*, 2:73–82, 1985.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *JOURNAL OF THE ROYAL STATISTICAL SOCIETY, SERIES B*, 39(1):1–38, 1977.
- P. J. Diggle, P. Moraga, B. Rowlingson, B. M. Taylor, et al. Spatial and spatio-temporal log-gaussian cox processes: extending the geostatistical paradigm. *Statistical Science*, 28(4):542–563, 2013.
- M. Eichler, R. Dahlhaus, and J. Dueck. Graphical modeling for multivariate hawkes processes with nonparametric link functions. *Journal of Time Series Analysis*, 38(2):225–242, 2017.
- S. Flaxman, Y. W. Teh, and D. Sejdinovic. Poisson intensity estimation with reproducing kernels. In *AISTATS*, pages 270–279, 2017.
- S. Gugushvili, F. van der Meulen, M. Schauer, and P. Spreij. Fast and scalable non-parametric bayesian inference for poisson point processes. *arXiv preprint arXiv:1804.03616*, 2018.
- P. F. Halpin. An em algorithm for hawkes process. *Psychometrika*, 2, 2012.
- A. G. Hawkes and D. Oakes. A cluster process representation of a self-exciting process. *Journal of Applied Probability*, 11(3):493–503, 1974.
- A. G. Hawkes. Spectra of some self-exciting and mutually exciting point processes. *Biometrika*, 58(1):83–90, 1971.
- H. Ishwaran and L. F. James. Gibbs sampling methods for stick-breaking priors. *Journal of the American Statistical Association*, 96(453):161–173, 2001.
- T. Kurashima, T. Althoff, and J. Leskovec. Modeling interdependent and periodic real-world action sequences. In *WWW*, page 803, 2018.
- E. Lewis and G. Mohler. A nonparametric em algorithm for multiscale hawkes processes. *Journal of Nonparametric Statistics*, 1(1):1–20, 2011.
- S. W. Linderman and R. P. Adams. Scalable bayesian inference for excitatory point process networks. *arXiv preprint arXiv:1507.03228*, 2015.
- C. Lloyd, T. Gunter, M. Osborne, and S. Roberts. Variational inference for gaussian process modulated poisson processes. In *ICML*, pages 1814–1822, 2015.
- J. Mercer. Functions of positive and negative type, and their connection with the theory of integral equations. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 209:415–446, 1909.
- S. Mishra, M.-A. Rizoïu, and L. Xie. Feature Driven and Point Process Approaches for Popularity Prediction. *CIKM*, 2016.
- J. Møller, A. R. Syversveen, and R. P. Waagepetersen. Log gaussian cox processes. *Scandinavian journal of statistics*, 25(3):451–482, 1998.
- C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005.
- J. G. Rasmussen. Bayesian inference for hawkes processes. *Methodology and Computing in Applied Probability*, 15(3):623–642, 2013.
- M.-A. Rizoïu, S. Mishra, Q. Kong, M. Carman, and L. Xie. Sir-hawkes: Linking epidemic models and hawkes processes to model diffusions in finite populations. In *WWW*, pages 419–428, 2018.
- J. Rousseau, S. Donnet, and V. Rivoirard. Nonparametric bayesian estimation of multivariate hawkes processes. *arXiv preprint arXiv:1802.05975*, 2018.
- I. Rubin. Regular point processes and their detection. *IEEE Transactions on Information Theory*, 18(5):547–557, 1972.
- T. Shirai and Y. Takahashi. Random point fields associated with certain fredholm determinants ii: Fermion shifts and their ergodic and gibbs properties. *The Annals of Probability*, 31(3):1533–1564, 2003.
- C. J. Walder and A. N. Bishop. Fast bayesian intensity estimation for the permanental process. In *ICML*, 2017.
- H. Xu, M. Farajtabar, and H. Zha. Learning granger causality for hawkes processes. In *ICML*, pages 1717–1726, 2016.
- H. Xu, D. Luo, X. Chen, and L. Carin. Benefits from superposed hawkes processes. *AISTATS*, 2018.
- Q. Zhao, M. A. Erdogdu, H. Y. He, A. Rajaraman, and J. Leskovec. Seismic: A self-exciting point process model for predicting tweet popularity. In *SIGKDD*. ACM, 2015.
- K. Zhou, H. Zha, and L. Song. Learning social infectivity in sparse low-rank networks using multi-dimensional hawkes processes. In *AISTATS*, pages 641–649, 2013.
- K. Zhou, H. Zha, and L. Song. Learning triggering kernels for multi-dimensional hawkes processes. In *ICML*, pages 1301–1309, 2013.

## A Computing the Integral Term of the Log-likelihood

We consider  $\Omega = [0, T]$ , the background intensity  $\mu$ , the triggering kernel  $\phi(\cdot) = 1/2f(\cdot)^2$ ,  $f(\cdot) = \boldsymbol{\omega}^T \mathbf{e}(\cdot)$ , and data  $\{t_i\}_{i=1}^N$ , and the integral term in the log-likelihood is calculated as below

$$\begin{aligned}
 & \text{Integral Term} \\
 &= -\frac{1}{2} \sum_{i=1}^N \int_0^T f^2(t - t_i) dt \\
 &= -\frac{1}{2} \sum_{i=1}^N \int_0^{T-t_i} \left[ \sum_{k=1}^K \omega_k e_k(t) \right]^2 dt \\
 &= -\frac{1}{2} \sum_{i=1}^N \sum_{k=1}^K \sum_{k'=1}^K \omega_k \omega_{k'} \underbrace{\int_0^{T-t_i} e_k(t) e_{k'}(t) dt}_{U_{kk'}^{(i)}} \\
 &= -\frac{1}{2} \sum_{i=1}^N \boldsymbol{\omega}^T U^{(i)} \boldsymbol{\omega}. \tag{21}
 \end{aligned}$$

In our case, Eq. (14) has  $d = 1$ , i.e.,  $\phi_k(x) = (2/\pi)^{1/2} \sqrt{1/2}^{[k-1=0]} \cos[(k-1)x]$ ,  $k = 1, 2, \dots$ . The matrix  $U^{(i)}$  is calculated as below:

$$\begin{aligned}
 U_{1,1}^{(i)} &= \int_0^{T-t_1} \frac{1}{\pi} dt = \frac{T-t_i}{\pi}, \\
 U_{k>1,1}^{(i)} &= U_{1,k>1}^{(i)} = \frac{\sqrt{2} \sin[(k-1)(T-t_i)]}{\pi(k-1)}, \\
 U_{k,k(k>1)} &= \frac{1}{\pi} \left\{ T-t_i + \frac{\sin[2(k-1)(T-t_i)]}{2(k-1)} \right\}, \\
 U_{k,k'(k \neq k')} &= \frac{1}{\pi} \left\{ \frac{\sin[(k-k')(T-t_i)]}{k-k'} + \frac{\sin[(k+k'-2)(T-t_i)]}{k+k'-2} \right\}.
 \end{aligned}$$

## B M.A.P. $\mu$ and $\phi$ Given Infinite Branching Structures

M.A.P.  $\mu$  and  $\phi$  given Infinite branching structures is written as:

$$\begin{aligned}
 & \operatorname{argmax}_{\boldsymbol{\omega}, \mu} \mathbb{E}_B [\log p(\boldsymbol{\omega}, \mu | B, \{t_i\}_{i=1}^N, \Omega, k)] \\
 &= \operatorname{argmax}_{\boldsymbol{\omega}, \mu} \underbrace{\mathbb{E}_B [\log p(\{t_i\}_{i=1}^N | \boldsymbol{\omega}, \mu, B, \Omega, k)]}_{\text{Expected Log-likelihood}} + \underbrace{\log p(\boldsymbol{\omega}) + \log p(\mu)}_{\text{Constraints}} \\
 &= \operatorname{argmax}_{\boldsymbol{\omega}, \mu} \sum_{i=1}^N \left\{ \sum_{t_j < t_i} p_{ij} \log \frac{1}{2} [\boldsymbol{\omega}^T \mathbf{e}(t_i - t_j)]^2 - p_{i0} \log \mu - \frac{1}{2} \int_0^T [\boldsymbol{\omega}^T \mathbf{e}(t - t_i)]^2 dt \right\} - (\beta + 1) \mu T \\
 & \quad - \frac{1}{2} \boldsymbol{\omega}^T \Lambda^{-1} \boldsymbol{\omega} - (\alpha - 1) \log \mu, \tag{22}
 \end{aligned}$$

where  $B$  represents the branching structure,  $p_{ij}$  the probabilities of triggering relationships shown as Eq. (9) and Eq. (10), and  $\alpha, \beta$  are parameters of the Gamma prior of  $\mu T$ . The second line is obtained using Bayes' rule, which shows M.A.P.  $\mu$  and  $\phi$  given Infinite branching structures is equivalent to maximizing the constrained expected log-likelihood, i.e., the objective function for the M-step of the EM algorithm and the third line is an explicit expression of the second line.

## C Mode-Finding the Triggering Kernel

Here we demonstrate in detail the computational challenges involved in finding the posterior mode with respect to the value of the triggering kernel at multiple point locations. Consider the triggering kernel  $\phi(\cdot) = \frac{1}{2} f^2(\cdot)$  where  $f(\cdot)$  is Gaussian process distributed. For a dataset  $\{t_i\}_{i=1}^N$ ,  $\mathbf{X} \equiv \{f(t_i)\}_{i=1}^N = \{X_i\}_{i=1}^N$  has a normal distribution, i.e.,  $\{f(t_i)\}_{i=1}^N \sim \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma})$  where  $\mathbf{m}$  and  $\boldsymbol{\Sigma}$  are the mean and the covariance matrix. The distribution of  $\mathbf{Y} \equiv \{\phi(t_i)\}_{i=1}^N = \{Y_i\}_{i=1}^N$  is derived as below where  $F$



is the cumulative density function and  $f$  the probabilistic density function.

$$\begin{aligned}
F_{\mathbf{Y}}(\mathbf{y}) &= P(-\sqrt{2y_i} < X_i < \sqrt{2y_i}, i = 1, \dots, N) \\
&= \int_{-\sqrt{2y_1}}^{\sqrt{2y_1}} \dots \int_{-\sqrt{2y_N}}^{\sqrt{2y_N}} \frac{1}{\sqrt{(2\pi)^N \Sigma^{-1}}} \exp\left[-\frac{(\mathbf{X} - \mathbf{m})^T \Sigma^{-1} (\mathbf{X} - \mathbf{m})}{2}\right] dX_1 \dots dX_N, \\
f_{\mathbf{Y}}(\mathbf{y}) &= \frac{\partial^N}{\partial y_1 \dots \partial y_N} F_{\mathbf{Y}}(\mathbf{y}) \\
&= \frac{1}{\sqrt{(2\pi)^N \Sigma^{-1}}} \left(\prod_{i=1}^N \frac{1}{2\sqrt{2y_i}}\right) \sum_{\mathbf{X} \in \times_{i=1}^N \{\sqrt{2y_i}, -\sqrt{2y_i}\}} \exp\left[-\frac{(\mathbf{X} - \mathbf{m})^T \Sigma^{-1} (\mathbf{X} - \mathbf{m})}{2}\right], \tag{23}
\end{aligned}$$

where  $\times$  is the Cartesian product. There are  $2^N$  summations of exponential functions, which is intractable.

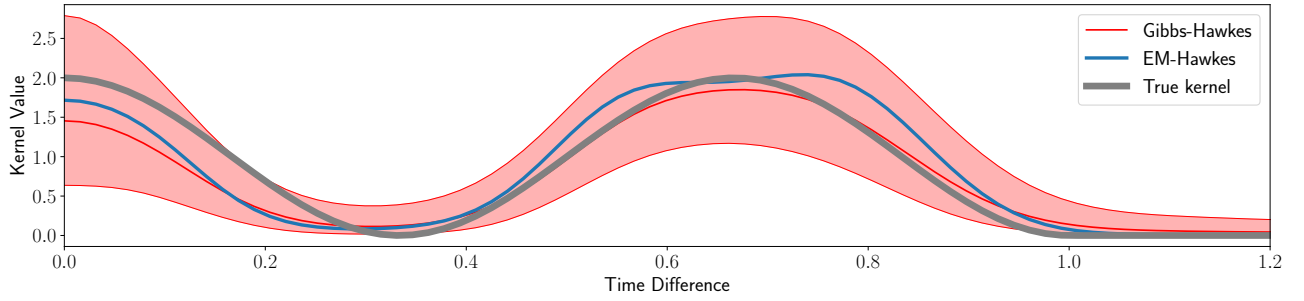


Figure 6: Triggering kernels estimated by the Gibbs-Hawkes method (Section 3) and the EM-Hawkes method (Section 4). The true kernel is plotted as the bold gray curve. We plot the median (red) and [0.1, 0.9] interval (filled red) of the approximate predictive distribution, along with the triggering kernel inferred by the EM Hawkes method (blue). The hyper-parameters  $a$  and  $b$  of the Gaussian process kernel (Eq. (13)) are set to 0.002.