

---

# Supplementary Material

## “Fast Bayesian Inference for Gaussian Cox Processes via Path Integral Formulation”

---

**Hideaki Kim**  
 NTT Human Informatics Laboratories  
 NTT Corporation  
 hideaki.kin.cn@hco.ntt.co.jp

### S1 Functional Determinant

Here we show that the functional determinant of integral operator  $\mathcal{K}$ , defined by  $|\mathcal{K}| \triangleq \lim_{J \rightarrow \infty} |\mathbf{K}\Delta|$  in (8), can be calculated by the product of its eigenvalues.

Under the discrete scenario (7), the  $l$ -th eigenvalue of the  $J \times J$  matrix,  $\mathbf{K}\Delta$ , should satisfy the following eigenvalue equation,

$$\sum_{j'=1}^J k(t_j, t_{j'})(t_{j'} - t_{j'-1})\tilde{v}_{j'}^l = \tilde{\lambda}_l \tilde{v}_j^l, \quad \forall j \in \{1, \dots, J\}, \quad (\text{S1})$$

where  $\tilde{\lambda}_l$  denotes the  $l$ -th eigenvalue, and  $\tilde{v}^l \triangleq (\tilde{v}_1^l, \dots, \tilde{v}_J^l)^\top$  is the  $J$ -element eigenvector for the  $l$ -th eigenvalue. The determinant of  $\mathbf{K}\Delta$  is thus given by the product of the finite number of eigenvalues,  $|\mathbf{K}\Delta| = \prod_{l=1}^L \tilde{\lambda}_l$ . Given the limit of the division number  $J \rightarrow \infty$  ( $\Delta \rightarrow \mathbf{0}$ ), the Nyström method [8] states that the eigenvalue equation (S1) converges to a homogeneous Fredholm integral equation, i.e., the eigenvalue equation of operator  $\mathcal{K}$ ,

$$\mathcal{K}v^l = \lambda_l v^l, \quad \text{or} \quad \int_{\mathcal{T}} k(t, s)v^l(s)ds = \lambda_l v^l(t), \quad t \in \mathcal{T}, \quad (\text{S2})$$

where  $v^l(t_j) = \tilde{v}_j^l$ , and  $\lambda_l = \lim_{J \rightarrow \infty} \tilde{\lambda}_l$ . The set of eigenvalues for the integral equation (S2) is countably infinite [9]. Therefore, the functional determinant of operator  $\mathcal{K}$ , that is,  $|\mathcal{K}| \triangleq \lim_{J \rightarrow \infty} |\mathbf{K}\Delta|$ , can be described as the product of the infinite number of eigenvalues of (S2):

$$|\mathcal{K}| \triangleq \lim_{J \rightarrow \infty} |\mathbf{K}\Delta| = \prod_{l=1}^{\infty} \lambda_l. \quad (\text{S3})$$

It should be noted that integral operator  $\mathcal{K}$  with a positive definite kernel function often has infinitely small positive eigenvalues, where the marginalization of Gaussian process prior by path integral (9) leads to the trivial value of zero,  $\sqrt{|\mathcal{K}|} = 0$ . Actually, the path integral provides a rational result when applied to the marginalization of the posterior probability (16), where the ratio of two functional determinants are considered.

## S2 Derivation of MAP Estimator Equation

We detail here the derivation of the integral equation (12) that yields the MAP estimator. The functional derivative of  $S(x(\mathbf{t}), \underline{x}(\mathbf{t}))$  should be zero on the MAP estimator  $\hat{x}(\mathbf{t})$ :

$$\begin{aligned} \delta S(\hat{x}(\mathbf{t}), \hat{\underline{x}}(\mathbf{t})) &= \int_{\mathcal{T}} \left[ \frac{\delta S}{\delta \hat{x}(\mathbf{t})} \delta x(\mathbf{t}) + \frac{\delta S}{\delta \hat{\underline{x}}(\mathbf{t})} \delta \underline{x}(\mathbf{t}) \right] dt + O((\delta x)^2) \\ &\simeq \int_{\mathcal{T}} \left[ \dot{\kappa}(\hat{x}(\mathbf{t})) - \sum_{n=1}^N \frac{\dot{\kappa}(\hat{x}(\mathbf{t}_n))}{\kappa(\hat{x}(\mathbf{t}_n))} \delta(\mathbf{t} - \mathbf{t}_n) + \frac{1}{2} \hat{\underline{x}}(\mathbf{t}) \right] \delta x dt + \int_{\mathcal{T}} \frac{1}{2} (\hat{x}(\mathbf{t}) - \mu) \delta \underline{x} dt \\ &= \int_{\mathcal{T}} \left[ \dot{\kappa}(\hat{x}(\mathbf{t})) - \sum_{n=1}^N \frac{\dot{\kappa}(\hat{x}(\mathbf{t}_n))}{\kappa(\hat{x}(\mathbf{t}_n))} \delta(\mathbf{t} - \mathbf{t}_n) + \hat{\underline{x}}(\mathbf{t}) \right] \delta x dt = 0, \end{aligned}$$

where the following relation was used,

$$\begin{aligned} \int_{\mathcal{T}} (\hat{x}(\mathbf{t}) - \mu) \delta \underline{x} dt &= \int_{\mathcal{T}} (\hat{x}(\mathbf{t}) - \mu) \int_{\mathcal{T}} k^*(\mathbf{t}, \mathbf{t}') \delta x(\mathbf{t}') dt' dt \\ &= \int_{\mathcal{T}} dt' \delta x(\mathbf{t}') \int_{\mathcal{T}} k^*(\mathbf{t}, \mathbf{t}') (\hat{x}(\mathbf{t}) - \mu) dt \\ &= \int_{\mathcal{T}} \hat{\underline{x}}(\mathbf{t}') \delta x dt'. \quad \because k^*(\mathbf{t}, \mathbf{t}') = k^*(\mathbf{t}', \mathbf{t}) \end{aligned}$$

Thus the following equation is derived,

$$\hat{\underline{x}}(\mathbf{t}) + \dot{\kappa}(\hat{x}(\mathbf{t})) = \sum_{n=1}^N \frac{\dot{\kappa}(\hat{x}(\mathbf{t}_n))}{\kappa(\hat{x}(\mathbf{t}_n))} \delta(\mathbf{t} - \mathbf{t}_n), \quad \mathbf{t} \in \mathcal{T}. \quad (\text{S4})$$

By applying operator  $\mathcal{K}$  to (S4), we realize an equation that derives the MAP estimator  $\hat{x}(\mathbf{t})$  as follows,

$$\hat{x}(\mathbf{t}) + \int_{\mathcal{T}} k(\mathbf{t}, \mathbf{t}') \dot{\kappa}(\hat{x}(\mathbf{t}')) dt' = \mu + \sum_{n=1}^N k(\mathbf{t}, \mathbf{t}_n) \frac{\dot{\kappa}(\hat{x}(\mathbf{t}_n))}{\kappa(\hat{x}(\mathbf{t}_n))}, \quad \mathbf{t} \in \mathcal{T}. \quad (\text{S5})$$

## S3 Derivation of Predictive Covariance

We detail the derivation of the predictive covariance shown in (13-15). The predictive inverse covariance (precision), denoted by  $\sigma^*(\mathbf{t}, \mathbf{t}')$ , is given by the second functional derivative of  $S$ , which is written as

$$\sigma^*(\mathbf{t}, \mathbf{t}') \triangleq \left. \frac{\delta^2 S(x, \underline{x})}{\delta x(\mathbf{t}) \delta x(\mathbf{t}')} \right|_{x=\hat{x}} = z^*(\mathbf{t}, \mathbf{t}') + h^*(\mathbf{t}, \mathbf{t}'),$$

where

$$\begin{aligned} z^*(\mathbf{t}, \mathbf{t}') &\triangleq - \sum_{n=1}^N \left. \frac{d^2 \log(\kappa(x))}{dx^2} \right|_{x=\hat{x}(\mathbf{t})} \delta(\mathbf{t} - \mathbf{t}_n) \delta(\mathbf{t}' - \mathbf{t}_n), \\ h^*(\mathbf{t}, \mathbf{t}') &\triangleq a(\mathbf{t}, \mathbf{t}') + k^*(\mathbf{t}, \mathbf{t}'), \\ a(\mathbf{t}, \mathbf{t}') &\triangleq \ddot{\kappa}(\hat{x}(\mathbf{t})) \delta(\mathbf{t} - \mathbf{t}'). \end{aligned} \quad (\text{S6})$$

Let the integral operators corresponding to  $\sigma(\mathbf{t}, \mathbf{t}')$ ,  $z(\mathbf{t}, \mathbf{t}')$ ,  $h(\mathbf{t}, \mathbf{t}')$ , and  $a(\mathbf{t}, \mathbf{t}')$  be denoted by  $\Sigma$ ,  $\mathcal{Z}$ ,  $\mathcal{H}$ , and  $\mathcal{A}$ , respectively, and their inverse counterparts by  $*$ . Using the fact that operator  $\mathcal{Z}^*$  is factorized as,

$$\begin{aligned} \mathcal{Z}^* &= \int_{\mathcal{T}} \cdot z^*(\mathbf{t}, \mathbf{t}') dt' = \mathcal{U}^\top \mathcal{Z}^{-1} \mathcal{U}, \\ \mathcal{Z}_{nn'} &\triangleq - \left. \frac{d^2 \log(\kappa(x))}{dx^2} \right|_{x=\hat{x}(\mathbf{t})} \delta_{nn'}, \quad \mathcal{U}_n \triangleq \int_{\mathcal{T}} \cdot \delta(\mathbf{t}' - \mathbf{t}_n) dt', \end{aligned}$$

we obtain the predictive covariance  $\sigma(\mathbf{t}, \mathbf{t}')$  with a finite (thus tractable)  $N$ -dimensional matrix representation as follows:

$$\begin{aligned}
\Sigma &= \int_{\mathcal{T}} \cdot \sigma(\mathbf{t}, \mathbf{t}') dt' = (\mathcal{Z}^* + \mathcal{H}^*)^* = (\mathcal{U}^\top \mathcal{Z}^{-1} \mathcal{U} + \mathcal{H}^*)^* \\
&= \mathcal{H} - \mathcal{H} \mathcal{U}^\top (\mathcal{Z} + \mathcal{U} \mathcal{H} \mathcal{U}^\top)^* \mathcal{U} \mathcal{H} \\
&= \int_{\mathcal{T}} \cdot \left[ h(\mathbf{t}, \mathbf{t}') - \mathbf{h}(\mathbf{t})^\top (\mathcal{Z} + \mathbf{H})^{-1} \mathbf{h}(\mathbf{t}') \right] dt', \\
\therefore \sigma(\mathbf{t}, \mathbf{t}') &= h(\mathbf{t}, \mathbf{t}') - \mathbf{h}(\mathbf{t})^\top (\mathcal{Z} + \mathbf{H})^{-1} \mathbf{h}(\mathbf{t}'),
\end{aligned}$$

where  $\mathbf{H}_{nn'} \triangleq h(\mathbf{t}_n, \mathbf{t}_{n'})$ ,  $\mathbf{h}(\mathbf{t}) \triangleq (h(\mathbf{t}, \mathbf{t}_1), \dots, h(\mathbf{t}, \mathbf{t}_N))^\top$ ; we used the Woodbury matrix identity in this derivation.

The remaining problem is how to obtain  $h(\mathbf{t}, \mathbf{t}')$ . Equation (S6) states that the operators  $\mathcal{H}$ ,  $\mathcal{A}$ , and  $\mathcal{K}$  hold the relation,  $\mathcal{H}^* = \mathcal{A} + \mathcal{K}^* \Leftrightarrow (\mathcal{I} + \mathcal{K} \mathcal{A}) \mathcal{H} = \mathcal{K}$ , which leads to the integral equation that  $h(\mathbf{t}, \mathbf{t}')$  should satisfy,

$$h(\mathbf{t}, \mathbf{s}) + \int_{\mathcal{T}} k(\mathbf{t}, \mathbf{t}') \ddot{\kappa}(\hat{\mathbf{x}}(\mathbf{t}')) h(\mathbf{t}', \mathbf{s}) dt' = k(\mathbf{t}, \mathbf{s}). \quad (\text{S7})$$

## S4 Derivation of Marginal Likelihood

We can obtain the marginal likelihood or evidence of GP models,  $p(\mathcal{D})$ , under Laplace approximation (13) by performing path integral (9) as follows:

$$\begin{aligned}
\log p(\mathcal{D}) &= \log \int \exp[-S(x(\mathbf{t}), \underline{\mathbf{x}}(\mathbf{t}))] \mathcal{D}x \\
&\simeq -S(\hat{\mathbf{x}}(\mathbf{t}), \hat{\underline{\mathbf{x}}}(\mathbf{t})) + \log \int e^{-\frac{1}{2} \iint_{\mathcal{T} \times \mathcal{T}} \sigma^*(\mathbf{t}, \mathbf{t}') (x(\mathbf{t}) - \hat{\mathbf{x}}(\mathbf{t})) (x(\mathbf{t}') - \hat{\mathbf{x}}(\mathbf{t}')) dt dt'} \mathcal{D}x \\
&= -S(\hat{\mathbf{x}}(\mathbf{t}), \hat{\underline{\mathbf{x}}}(\mathbf{t})) + \frac{1}{2} \log |\Sigma|.
\end{aligned}$$

Substituting (S4) into (11), we can write down  $S(\hat{\mathbf{x}}(\mathbf{t}), \hat{\underline{\mathbf{x}}}(\mathbf{t}))$  as

$$\begin{aligned}
S(\hat{\mathbf{x}}(\mathbf{t}), \hat{\underline{\mathbf{x}}}(\mathbf{t})) &= \frac{1}{2} \log |\mathcal{K}| - \sum_{n=1}^N \log \kappa(\hat{\mathbf{x}}(\mathbf{t}_n)) + \int_{\mathcal{T}} \kappa(\hat{\mathbf{x}}(\mathbf{t})) dt + \frac{1}{2} \int_{\mathcal{T}} (\hat{\mathbf{x}}(\mathbf{t}) - \mu) \hat{\underline{\mathbf{x}}}(\mathbf{t}) dt \\
&= \frac{1}{2} \log |\mathcal{K}| - \sum_{n=1}^N \log \kappa(\hat{\mathbf{x}}(\mathbf{t}_n)) + \int_{\mathcal{T}} \kappa(\hat{\mathbf{x}}(\mathbf{t})) dt \\
&\quad + \frac{1}{2} \int_{\mathcal{T}} (\hat{\mathbf{x}}(\mathbf{t}) - \mu) \left[ \sum_{n=1}^N \frac{\dot{\kappa}(\hat{\mathbf{x}}(\mathbf{t}_n))}{\kappa(\hat{\mathbf{x}}(\mathbf{t}_n))} \delta(\mathbf{t} - \mathbf{t}_n) - \dot{\kappa}(\hat{\mathbf{x}}(\mathbf{t})) \right] dt \\
&= \frac{1}{2} \log |\mathcal{K}| + \int_{\mathcal{T}} \kappa(\hat{\mathbf{x}}(\mathbf{t})) dt - \frac{1}{2} \int_{\mathcal{T}} (\hat{\mathbf{x}}(\mathbf{t}) - \mu) \dot{\kappa}(\hat{\mathbf{x}}(\mathbf{t})) dt \\
&\quad + \sum_{n=1}^N \left[ \frac{1}{2} (\hat{\mathbf{x}}(\mathbf{t}_n) - \mu) \frac{\dot{\kappa}(\hat{\mathbf{x}}(\mathbf{t}_n))}{\kappa(\hat{\mathbf{x}}(\mathbf{t}_n))} - \log \kappa(\hat{\mathbf{x}}(\mathbf{t}_n)) \right].
\end{aligned}$$

Furthermore, by using the matrix determinant lemma, we can rewrite  $\log |\Sigma|$  as

$$\begin{aligned}
\log |\Sigma| &= \log |\mathcal{H} - \mathcal{H} \mathcal{U}^\top (\mathcal{Z} + \mathcal{U} \mathcal{H} \mathcal{U}^\top)^* \mathcal{U} \mathcal{H}| \\
&= \log |\mathcal{H}| - \log |\mathcal{Z} + \mathcal{U} \mathcal{H} \mathcal{U}^\top| + \log |(\mathcal{Z} + \mathcal{U} \mathcal{H} \mathcal{U}^\top) - (\mathcal{U} \mathcal{H}) \mathcal{H}^* (\mathcal{H} \mathcal{U}^\top)| \\
&= \log |\mathcal{H}| - \log |\mathcal{Z} + \mathcal{U} \mathcal{H} \mathcal{U}^\top| + \log |\mathcal{Z}| \\
&= \log |\mathcal{H}| - \log |\mathbf{I}_N + \mathbf{Z}^{-1} \mathbf{H}|.
\end{aligned}$$

## S5 Calculation of $h(\mathbf{t}, \mathbf{s})$

We apply the *Gelerkin method* [2], a variant of the projection method, to solve equation (S7) with regard to  $h(\mathbf{t}, \mathbf{s})$ . Let

$$h(\mathbf{t}, \mathbf{s}) \simeq \sum_{l=1}^L \omega_l \varphi_l(\mathbf{t}) \varphi_l(\mathbf{s}),$$

and solve for the coefficients,  $\{\omega_l\}_{l=1}^L$ , using the set of residual equations derived from (S7),

$$\iint_{\mathcal{T} \times \mathcal{T}} r(\mathbf{t}, \mathbf{s}) \varphi_l(\mathbf{t}) \varphi_l(\mathbf{s}) dt ds = 0, \quad l = 1, \dots, L,$$

$$r(\mathbf{t}, \mathbf{s}) \triangleq h(\mathbf{t}, \mathbf{s}) + \int_{\mathcal{T}} k(\mathbf{t}, \mathbf{t}') \ddot{\kappa}(\hat{x}(\mathbf{t}')) h(\mathbf{t}', \mathbf{s}) dt' - k(\mathbf{t}, \mathbf{s}).$$

Using the relations of the eigenfunctions,

$$\int_{\mathcal{T}} k(\mathbf{t}, \mathbf{s}) \varphi_l(\mathbf{s}) d\mathbf{s} = \lambda_l \varphi_l(\mathbf{t}), \quad \int_{\mathcal{T}} h(\mathbf{t}, \mathbf{s}) \varphi_l(\mathbf{s}) d\mathbf{s} = \omega_l \varphi_l(\mathbf{t}),$$

we can solve the residual equations, which results in

$$\omega_l = \frac{\lambda_l}{1 + \lambda_l \Xi_l}, \quad \Xi_l = \int_{\mathcal{T}} \ddot{\kappa}(\hat{x}(\mathbf{t})) \varphi_l^2(\mathbf{t}) dt.$$

## S6 Representer Theorem

Formula (S5) provides the MAP estimator  $\hat{x}(\mathbf{t})$  for general Gaussian Cox processes, but an interesting representation is obtained under a quadratic link function.

If link function  $\kappa(x)$  is given in quadratic form,

$$\kappa(x) = x^2, \quad \dot{\kappa}(x) = 2x, \quad \ddot{\kappa}(x) = 2,$$

then (S5) reduces to a Fredholm integral equation of the second kind,

$$\hat{x}(\mathbf{t}) + 2 \int_{\mathcal{T}} k(\mathbf{t}, \mathbf{t}') \hat{x}(\mathbf{t}') dt' = \mu + 2 \sum_{n=1}^N k(\mathbf{t}, \mathbf{t}_n) \hat{x}(\mathbf{t}_n)^{-1}.$$

The linearity of the integral equation permits a representation of the form

$$\hat{x}(\mathbf{t}) - \mu(1 - 2\underline{h}(\mathbf{t})) = 2 \sum_{n=1}^N h(\mathbf{t}, \mathbf{t}_n) \hat{x}(\mathbf{t}_n)^{-1}, \quad (\text{S8})$$

where  $h(\mathbf{t}, \mathbf{s})$  is the positive semi-definite kernel defined in (S7), and  $\underline{h}(\mathbf{t}) \triangleq \int_{\mathcal{T}} h(\mathbf{t}, \mathbf{s}) d\mathbf{s}$ . (S8) states that MAP estimator  $\hat{x}(\mathbf{t})$  can be written as expansions in terms of the training examples, where the MAP estimation reduces to a finite-dimensional optimization problem corresponding to new kernel function  $h(\mathbf{t}, \mathbf{s})$ . Therefore, the MAP estimator of Gaussian Cox process involves a representer theorem under a quadratic link function, which is a generalization of Wahba's classical representer theorem [12]. Function  $h(\mathbf{t}, \mathbf{s})$  has been studied by [11, 3] as the *equivalent kernel*.

## S7 Kronecker Structure in Product of Eigenfunctions

For a set of  $L = \prod_d L_d$  points on a Cartesian product grid (the grid need not be regular),  $\mathbf{p} \in \mathcal{T}_1 \times \dots \times \mathcal{T}_D$ , matrix  $\Phi$ , defined by

$$\Phi_{\nu l} \triangleq \varphi_{\nu}(\mathbf{p}_l) = \prod_{d=1}^D \varphi_{\nu'_d}^{(d)}(p_{l_d}^{(d)}), \quad 1 \leq \nu'_d, l_d \leq L_d, \quad (\text{S9})$$

has a Kronecker structure as indicated by

$$\Phi = \Phi^{(1)} \otimes \dots \otimes \Phi^{(D)}, \quad \Phi_{\nu'_d l_d}^{(d)} \triangleq \varphi_{\nu'_d}^{(d)}(p_{l_d}^{(d)}). \quad (\text{S10})$$

Therefore, in the multi-dimensional input setting, exploiting the Kronecker structure can reduce the  $\mathcal{O}(L^2)$  computation in MAP estimation, or  $\{\sum_l \lambda_l \beta_l \varphi_l(\mathbf{p}_{\nu'})\}_{\nu'=1}^L$  and  $\{\sum_l \beta_l \varphi_l(\mathbf{p}_{\nu'})\}_{\nu'=1}^L$ , to  $\mathcal{O}(L)$  computation [10], although we did not employ it in this paper.

Table S1: Symbols and Definitions.

Symbol	Definition
$N$	number of data points
$\mathcal{T}$	observation region
$\mathcal{D} = \{\mathbf{t}_n \in \mathcal{T}\}_{n=1}^N$	observed data points
$k^{(*)}(\mathbf{t}, \mathbf{t}')$	(inverse) kernel function
$\mathcal{K}^{(*)} = \int_{\mathcal{T}} k^{(*)}(\mathbf{t}, \mathbf{t}') \cdot d\mathbf{t}'$	integral operator corresponding to $k^{(*)}(\mathbf{t}, \mathbf{t}')$
$ \mathcal{K} $	functional determinant of operator $\mathcal{K}$
$\int \cdot \mathcal{D}x$	path integral computation with respect to continuous path $x(\mathbf{t})$

Table S2: Summary of related works.  $\mathcal{O}_M$ ,  $\mathcal{O}_V$ ,  $\mathcal{O}_E$  represent the computational costs of the MAP/predictive mean, the predictive covariance, and the marginal likelihood/evidence lower bound, respectively.  $Q$  and  $P$  represent the number of gradient descent iterations and the number of MCMC samplings, respectively. For the other symbols, see the main text.

	Proposed	[13]	[1, 6]	[4]
$\mathcal{O}_M$	$(NL+L^2)Q$	$(N \cdot \min(N, L))Q$	$(NL^2+L^3)Q$	$(NL+L^2)P$
$\mathcal{O}_V$	$\frac{\mathcal{O}_M + N_{\text{mc}}L + (N+L) \cdot \min(N^2, L^2)}{(N+L) \cdot \min(N^2, L^2)}$	$\frac{\mathcal{O}_M + (N+L) \cdot \min(N^2, L^2)}{(N+L) \cdot \min(N^2, L^2)}$	$(NL^2+L^3)Q$	$(NL+L^2)P$
$\mathcal{O}_E$	$\frac{\mathcal{O}_M + N_{\text{mc}}L + (N+L) \cdot \min(N^2, L^2)}{(N+L) \cdot \min(N^2, L^2)}$	$\frac{\mathcal{O}_M + (N+L) \cdot \min(N^2, L^2)}{(N+L) \cdot \min(N^2, L^2)}$	$(NL^2+L^3)$	$(NL+L^2)$

## S8 Summary of Key Expression and Algorithm

We summarize the key expressions in Section 2.1-2.2 (see Table S1) and the proposed algorithm (see Algorithm 1).

## S9 Summary of Computational Complexity

We provide a summary of the related works about the computational complexity (see Table S2). Here we focus on the algorithms that could scale linearly with the number of data points.  $\mathcal{O}_V$  and  $\mathcal{O}_E$  of our proposed method reduces to those of [13] under the quadratic link function because the integral in Eq. (19) can be performed analytically. Note that the complexity of each algorithm could be reduced by utilizing a stochastic gradient descent algorithm or by exploiting the Kronecker structure.

## S10 Experimental Settings and Additional Experiments

### S10.1 Experimental Settings

For all the experiments in Section 5, we used a multiplicative Gaussian kernel,  $k(\mathbf{t}, \mathbf{s}) = \prod_d e^{-(\theta(t_d - s_d))^2}$ , where hyper-parameter  $\theta$  was optimized for each trial by grid search. We calculated the marginal likelihoods or the evidence lower bounds on the following sets of hyper-parameters:

$$\theta_{\lambda_1(t)} = \{.01, .02, \dots, 0.1\}, \quad \theta_{\lambda_2(t)} = \{0.5, 1.0, \dots, 5.0\}, \quad \theta_{\lambda_3(t)} = \{.01, .02, \dots, 0.1\},$$

for synthetic data, and

$$\theta_{2\text{D neuronal data}} = \{.01, .02, .03, \dots, 0.1\}, \quad \theta_{3\text{D taxi data}} = \{0.5, 1.0, 1.5, \dots, 5.0\},$$

for open real-world data sets. Then we adopted the hyper-parameter that maximized the marginal likelihood or the evidence lower bounds. Note that each of the parameter sets contains values close to those used in [1].

---

**Algorithm 1** Bayesian Inference via Path Integral Formulation
 

---

```

1: procedure INFERENCE( $\mathbf{t}, \mathcal{D}, \mathcal{T}, J, L, N_{\text{mc}}, N_{\text{gd}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \mu, \{\mathbf{p}_l\}_{l=1}^L$ )
2:    $\{\hat{\beta}_l, \hat{\omega}_l, \lambda_l, \varphi_l(\cdot)\}_{l=1}^L, \log p(\mathcal{D} | \boldsymbol{\theta}) = \text{TRAINING}(\mathcal{D}, \mathcal{T}, J, L, N_{\text{mc}}, N_{\text{gd}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \mu, \{\mathbf{p}_l\}_{l=1}^L)$ 
3:   Predictive mean / MAP:  $\hat{x}(\mathbf{t}) = \mu + \sum_n k(\mathbf{t}, \mathbf{t}_n) \gamma(\sum_l \hat{\beta}_l \varphi_l(\mathbf{t}_n)) - \sum_l \lambda_l \hat{\beta}_l \varphi_l(\mathbf{t})$ 
4:   Predictive covariance is evaluated by Eq. (14)
5: procedure TRAINING( $\mathcal{D}, \mathcal{T}, J, L, N_{\text{mc}}, N_{\text{gd}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \mu, \{\mathbf{p}_l\}_{l=1}^L$ )
6:    $\{\lambda_l, \varphi_l(\cdot)\}_{l=1}^L = \text{EIGENFUNCTION}(\mathcal{T}, J, L, k(\cdot, \cdot | \boldsymbol{\theta}))$ 
7:    $\{\hat{\beta}_l\}_{l=1}^L = \text{MAP}(\mathcal{D}, L, N_{\text{gd}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \mu, \{\mathbf{p}_l\}_{l=1}^L, \{\lambda_l, \varphi_l(\cdot)\}_{l=1}^L)$ 
8:    $\{\hat{\omega}_l\}_{l=1}^L = \text{HFUNCTION}(\mathcal{T}, L, N_{\text{mc}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \{\hat{\beta}_l, \lambda_l, \varphi_l(\cdot)\}_{l=1}^L)$ 
9:    $\log p(\mathcal{D} | \boldsymbol{\theta}) = \text{MARGINALLIKELIHOOD}(\mathcal{D}, \mathcal{T}, L, N_{\text{mc}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \mu, \{\hat{\beta}_l, \hat{\omega}_l, \lambda_l, \varphi_l(\cdot)\}_{l=1}^L)$ 
10:  return  $\{\hat{\beta}_l, \hat{\omega}_l, \lambda_l, \varphi_l(\cdot)\}_{l=1}^L, \log p(\mathcal{D} | \boldsymbol{\theta})$ 
11: procedure EIGENFUNCTION( $\mathcal{T}, J, L, k(\cdot, \cdot | \boldsymbol{\theta})$ )
12:  for  $d = 1, \dots, D$  do
13:     $w = T^{(d)} / J$ 
14:    for  $j, j' = 1, \dots, J$  do
15:       $\mathbf{K}[j, j'] = k^{(d)}(jw, j'w)$ 
16:      Solve  $\mathbf{K} \mathbf{v}_j = e_j \mathbf{v}_j \quad : e_1 > e_2 > \dots > e_J$ 
17:       $\mathbf{k}(\cdot) = (k^{(d)}(\cdot, w), k^{(d)}(\cdot, 2w), \dots, k^{(d)}(\cdot, Jw))^\top$ 
18:      for  $j = 1, \dots, L_d$  do
19:         $\lambda_j^{(d)}, \varphi_j^{(d)}(\cdot) = e_j w, \mathbf{k}(\cdot)^\top \mathbf{v}_j / (e_j \sqrt{w})$ 
20:     $U = \emptyset$ 
21:    for  $j_1 = 1, \dots, L_1$  do
22:      ...
23:      for  $j_D = 1, \dots, L_D$  do
24:         $U = U \cup \{(\prod_{d=1}^D \lambda_{j_d}^{(d)}, \prod_{d=1}^D \varphi_{j_d}^{(d)}(\cdot))\}$ 
25:    return  $U$ 
26: procedure MAP( $\mathcal{D}, L, N_{\text{gd}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \mu, \{\mathbf{p}_l\}_{l=1}^L, \{\lambda_l, \varphi_l(\cdot)\}_{l=1}^L$ )
27:  Initialize  $\boldsymbol{\beta} \equiv (\beta_1, \dots, \beta_L)$ 
28:  for  $i = 1, \dots, N_{\text{gd}}$  do
29:     $\boldsymbol{\delta} = \nabla_{\boldsymbol{\beta}} \sum_{l=1}^L [r(\mathbf{p}_l)]^2 \quad : r(\mathbf{p}_l)$  is defined in Eq. (18)
30:    Update  $\boldsymbol{\beta}$  by Adam( $\boldsymbol{\delta}$ )
31:  return  $\boldsymbol{\beta}$ 
32: procedure HFUNCTION( $\mathcal{T}, L, N_{\text{mc}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \{\hat{\beta}_l, \lambda_l, \varphi_l(\cdot)\}_{l=1}^L$ )
33:  Sample  $N_{\text{mc}}$  points on  $\mathcal{T}, \{\mathbf{u}_i\}_{i=1}^{N_{\text{mc}}}$ 
34:  for  $l = 1, \dots, L$  do
35:    Compute  $\Xi_l$  by Monte Carlo integration with  $\{\mathbf{u}_i\}_{i=1}^{N_{\text{mc}}}$  : See Eq. (19)
36:     $\omega_l = \lambda_l / (1 + \lambda_l \Xi_l)$ 
37:  return  $\{\omega_l\}_{l=1}^L$ 
38: procedure MARGINALLIKELIHOOD( $\mathcal{D}, \mathcal{T}, L, N_{\text{mc}}, \kappa(\cdot), k(\cdot, \cdot | \boldsymbol{\theta}), \mu, \{\hat{\beta}_l, \hat{\omega}_l, \lambda_l, \varphi_l(\cdot)\}_{l=1}^L$ )
39:  Sample  $N_{\text{mc}}$  points on  $\mathcal{T}, \{\mathbf{u}_i\}_{i=1}^{N_{\text{mc}}}$ 
40:  Compute  $\log p(\mathcal{D} | \boldsymbol{\theta})$  by Monte Carlo integration with  $\{\mathbf{u}_i\}_{i=1}^{N_{\text{mc}}}$  : See Eq. (20)
41:  return  $\log p(\mathcal{D} | \boldsymbol{\theta})$ 

```

---

Table S3: Results on three types of synthetic data with standard errors in brackets.

$\lambda_1(t)$	$L = 3$			$L = 5$			$L = 10$			$L = 20$		
	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time
PIF <sub>s</sub>	13.00 (2.24)	9.68 (3.48)	10.33 (0.37)	11.33 (3.10)	7.64 (3.44)	10.19 (0.39)	11.56 (3.43)	8.02 (3.38)	10.49 (0.47)	11.58 (3.26)	7.62 (3.22)	10.96 (0.75)
PIF <sub>e</sub>	12.65 (1.76)	9.56 (2.06)	10.09 (0.37)	12.04 (3.38)	9.26 (2.95)	10.01 (0.31)	11.89 (4.43)	9.50 (3.70)	10.11 (0.30)	12.29 (4.33)	9.30 (3.74)	10.45 (0.40)
PIF <sub>q</sub>	11.80 (2.19)	8.65 (2.68)	9.91 (0.24)	11.88 (3.15)	8.17 (2.93)	9.86 (0.30)	12.73 (4.30)	9.02 (3.92)	10.05 (0.31)	12.50 (4.40)	9.00 (4.29)	10.49 (0.58)
STVB	12.59 (2.17)	8.70 (2.91)	34.86 (1.15)	11.81 (2.45)	8.20 (2.97)	34.88 (0.99)	11.86 (2.82)	8.39 (3.31)	35.25 (1.13)	12.04 (2.99)	8.54 (3.58)	35.90 (1.21)
VBPP	12.12 (1.99)	8.65 (2.35)	24.69 (0.92)	12.18 (3.61)	8.44 (3.39)	25.12 (1.00)	11.69 (4.94)	7.68 (4.06)	26.31 (1.16)	15.99 (3.17)	10.92 (4.77)	27.38 (0.79)

$\lambda_2(t)$	$L = 3$			$L = 5$			$L = 10$			$L = 20$		
	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time
PIF <sub>s</sub>	15.57 (1.64)	15.38 (4.92)	10.02 (0.22)	15.05 (1.06)	10.66 (3.89)	9.86 (0.15)	14.57 (0.93)	11.10 (3.53)	10.44 (0.73)	14.46 (0.83)	10.96 (3.20)	10.51 (0.50)
PIF <sub>e</sub>	34.73 (13.79)	30.96 (8.88)	9.81 (0.28)	19.41 (8.76)	13.83 (5.81)	9.62 (0.21)	15.95 (2.24)	12.56 (2.71)	9.99 (0.51)	15.68 (1.67)	12.35 (1.99)	10.12 (0.40)
PIF <sub>q</sub>	21.83 (3.71)	29.12 (8.91)	9.64 (0.17)	15.23 (2.98)	11.31 (3.06)	9.57 (0.15)	14.46 (1.96)	10.02 (2.61)	10.04 (0.64)	13.05 (1.88)	8.65 (1.70)	10.15 (0.38)
STVB	15.37 (0.96)	9.51 (2.62)	34.47 (0.35)	14.71 (1.41)	10.24 (2.61)	34.45 (0.22)	14.48 (1.48)	10.03 (2.72)	36.10 (1.90)	13.38 (1.91)	9.10 (2.69)	36.18 (1.17)
VBPP	29.54 (4.17)	15.92 (0.83)	23.80 (0.36)	15.18 (1.67)	10.20 (1.43)	23.90 (0.29)	14.79 (1.77)	10.14 (2.07)	26.66 (1.67)	13.13 (1.69)	8.79 (1.45)	26.66 (0.79)

$\lambda_3(t)$	$L = 3$			$L = 5$			$L = 10$			$L = 20$		
	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time	IQL <sub>.5</sub>	IQL <sub>.85</sub>	Time
PIF <sub>s</sub>	45.17 (4.26)	35.01 (9.64)	11.05 (0.33)	32.59 (7.24)	20.21 (6.94)	10.80 (0.29)	25.60 (6.17)	14.75 (3.80)	10.87 (0.27)	27.10 (6.12)	15.81 (3.06)	10.90 (0.32)
PIF <sub>e</sub>	39.69 (4.55)	26.70 (6.15)	10.24 (0.27)	32.24 (7.14)	21.05 (4.79)	10.17 (0.39)	31.46 (6.17)	21.89 (3.92)	10.53 (0.44)	31.23 (5.83)	21.15 (4.20)	10.48 (0.32)
PIF <sub>q</sub>	39.82 (4.89)	25.42 (5.88)	10.47 (0.21)	30.81 (5.89)	20.03 (4.12)	10.42 (0.29)	29.89 (4.47)	18.73 (3.53)	10.41 (0.27)	31.36 (5.91)	19.96 (3.56)	10.46 (0.32)
STVB	42.41 (5.71)	23.01 (5.95)	40.02 (1.20)	30.89 (7.04)	16.75 (3.17)	48.83 (37.64)	28.87 (6.76)	17.37 (3.31)	38.98 (0.65)	28.59 (6.46)	16.75 (3.46)	38.35 (1.04)
VBPP	38.73 (4.35)	22.79 (4.14)	26.60 (0.74)	30.97 (4.88)	19.41 (3.23)	25.88 (0.70)	30.50 (5.37)	19.31 (3.67)	26.86 (0.84)	37.09 (6.90)	25.06 (3.57)	28.14 (0.90)

For fair comparison, we employed a popular (batch) gradient descent algorithm, *Adam* [5], to perform estimations for all compared methods. We equally set the number of iteration as 5,000, but used different learning parameters for the models: 0.5 for PIF<sub>e</sub>; 0.05 for PIF<sub>q</sub>; 0.05 for PIF<sub>s</sub>; 0.005 for STVB; 0.05 for VBPP. We implemented all the models by using TensorFlow-2.2, where for STVB we used the python code provided by Aglietti et al. [1]. Each of the CPU times reported is the amount of time required to calculate the MAP/predictive mean, the predictive covariance, and the marginal likelihood given a hyper-parameter, where computing the eigenfunctions of the kernel is of course included in the CPU time.

### S10.2 Figure 1A in tabular form

We provide the results in Figure 1A in tabular form (see Table S3) to make the results easy to review.

### S10.3 Experiments of How Stably the Proposed Scheme Works on Real-world Data

To demonstrate that our scheme is stable on real-world data, we ran additional experiments based on the real-world data used in Figure 2. For the 2D neuronal data ( $N_{\text{train}} = 583$ ,  $N_{\text{test}} = 29127$ ), we extracted two training datasets of 583 data points from the original test data randomly, which resulted in three training ( $N_{\text{train}} = 583$ ) and a test ( $N_{\text{test}} = 27961$ ) datasets. For the 3D taxi data ( $N_{\text{train}}$

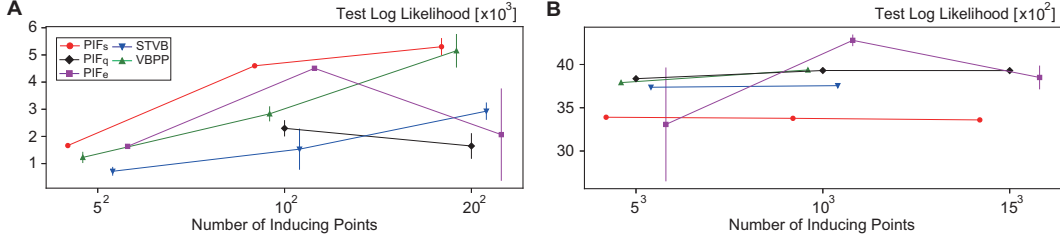


Figure S1: The predictive performances on open real-world data. The error bars represent the standard deviations. (A) 2D neuronal data. (B) 3D spatio-temporal taxi data.

= 1000,  $N_{\text{test}} = 3401$ ), we extracted two training datasets of 1000 data points from the original test data randomly, which resulted in three training ( $N_{\text{train}} = 1000$ ) and a test ( $N_{\text{test}} = 1401$ ) datasets. We evaluated three times the predictive performances of the compared models in terms of the test log likelihood, and calculated the means and the standard deviations (error bars) of the performances. Figure S1 shows the results.

#### S10.4 Experiments on Larger Taxi Dataset

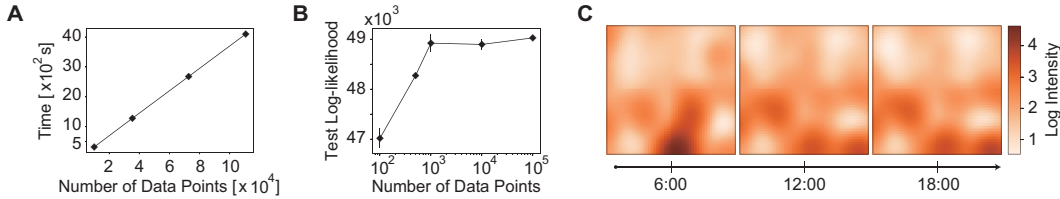


Figure S2: Results on large 3D taxi datasets. (A) The CPU times. (B) The test log-likelihoods. The error bars represent the standard deviations. (C) The estimated intensity function on the dataset with 110,705 data points.

To demonstrate that our scheme is scalable to large data size, we ran an additional experiment on a larger 3D taxi dataset [7]. The area considered was the same as that used in Experiments, and we took 10, 50, 100, and 150 weekdays from 1 July 2013, resulting in the training data sets containing from 9,971 to 110,705 data points. It should be noted that the maximum data size is comparable to that considered in [4] (113,020). We employed  $\text{PIF}_q$  with  $L = 10^3$ , and plot the execution times with respect to the number of data points (see Figure S2A). Also, we randomly divided the datasets ( $N = 110,705$ ) into the training ( $N = 100,000$ ) and the test ( $N = 10,705$ ) data, ran  $\text{PIF}_q$  on various sizes of subsets of the training data, and evaluated the test log-likelihoods. Figure S2B plots the test log-likelihoods as functions of the number of data points used for training; it shows that our approach ( $\text{PIF}_q$ ) can process large datasets effectively and recover the underlying intensity function more accurately with larger training datasets.

#### S10.5 Application of Stochastic Optimization Algorithm

Because the objective function to be minimized in the MAP estimation,  $\sum_{l=1}^L [r(\mathbf{p}_l)]^2$ , is given as a sum of residuals, we can apply a mini-batch gradient descent (MGD) algorithm to the optimization problem. With mini-batch size  $\tilde{L} \ll L$ , MGD reduces the computational complexity for each iteration to  $\mathcal{O}(NL + N\tilde{L} + \tilde{L}^2) \simeq \mathcal{O}(NL)$ , where the dominant cost stems from  $\{\gamma [\sum_l \beta_l \varphi_l(\mathbf{t}_n)]\}_{n=1}^N$ . Figure S3 shows the performances of MGD with  $\tilde{L} = 258$  and the batch gradient descent (BGD) on the 3D spatio-temporal taxi data used in Section 5 ( $N_{\text{train}} = 1000$ ), where the number of inducing points,  $L$ , was set as  $20^3$  for all models. In  $\text{PIF}_q$  and  $\text{PIF}_e$ , MGD (red line) converged much faster than BGD (blue line), which suggests that MGD would enhance the practical utility of our scheme under a large number of inducing points. In  $\text{PIF}_s$ , however, BGD achieved the convergence so rapidly and MGD worked poorly compared to BGD. Here we used the learning parameters of



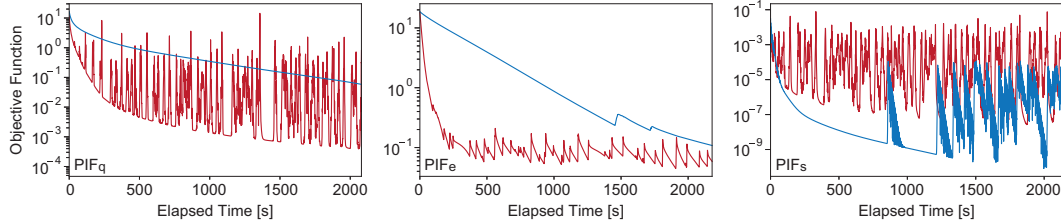


Figure S3: Training of MAP estimator on 3D spatio-temporal taxi data.  $L$  was set to  $20^3$  for each model. The blue and red lines represent the results yielded by batch gradient descent and mini-batch gradient descent with mini-batch size of 258, respectively. The maximum of the x-axis (elapsed time) in each figure equals the elapsed time that the batch gradient descent needed to execute the 5,000 epoch updates.

$10^{-3}$ ,  $10^{-2}$ , and  $10^{-5}$  for  $\text{PIF}_q$ ,  $\text{PIF}_e$ , and  $\text{PIF}_s$ , respectively, but examining the parameters of Adam [5] more carefully would speed up the convergence.

## References

- [1] Virginia Aglietti, Edwin V. Bonilla, Theodoros Damoulas, and Sally Cripps. Structured variational inference in continuous Cox process models. In *Advances in Neural Information Processing Systems* 32, 2019. <https://github.com/VirgiAgl/STVB>.
- [2] Kendall E. Atkinson. A survey of numerical methods for solving nonlinear integral equations. *The Journal of Integral Equations and Applications*, 4(1):15–46, 1992.
- [3] Seth Flaxman, Yee Whye Teh, and Dino Sejdinovic. Poisson intensity estimation with reproducing kernels. In *Artificial Intelligence and Statistics*, pages 270–279. PMLR, 2017.
- [4] S. T. John and James Hensman. Large-scale Cox process inference using variational Fourier features. In *International Conference on Machine Learning*, volume 80, pages 2362–2370. PMLR, 2018.
- [5] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [6] Chris Lloyd, Tom Gunter, Michael Osborne, and Stephen Roberts. Variational inference for Gaussian process modulated Poisson processes. In *International Conference on Machine Learning*, volume 37, pages 1814–1822. PMLR, 2015.
- [7] Luis Moreira-Matias, João Gama, Michel Ferreira, João Mendes-Moreira, and Luis Damas. Predicting taxi-passenger demand using streaming data. *IEEE Transactions on Intelligent Transportation Systems*, 14(3):1393–1402, 2013.
- [8] E. J. Nyström. Über die praktische Auflösung von Integralgleichungen mit Anwendungen auf Randwertaufgaben. *Acta Mathematica*, 54:185–204, 1930.
- [9] Andrei D. Polyanin and Alexander V. Manzhirov. *Handbook of Integral Equations*. CRC press, 1998.
- [10] Yunus Saatçi. *Scalable Inference for Structured Gaussian Process Models*. PhD thesis, University of Cambridge, 2011.
- [11] Peter Sollich and Christopher K. I. Williams. Understanding Gaussian process regression using the equivalent kernel. In *International Workshop on Deterministic and Statistical Methods in Machine Learning*, pages 211–228. Springer, 2004.
- [12] Grace Wahba. *Spline models for observational data*, volume 59. SIAM, 1990.
- [13] Christian J. Walder and Adrian N. Bishop. Fast Bayesian intensity estimation for the permanent process. In *International Conference on Machine Learning*, volume 70, pages 3579–3588. PMLR, 2017.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [\[Yes\]](#)
  - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 6
  - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See Section 6
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See Section 2
  - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) Also see supplementary materials
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[No\]](#) The data and instructions are provided in Section 5, but the code is proprietary.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Section 5
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) See Section 5
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Section 5
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) See Section 5 and References
  - (b) Did you mention the license of the assets? [\[Yes\]](#) See Section 5
  - (c) Did you include any new assets either in the supplemental material or as a URL? [\[No\]](#)
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[No\]](#)
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[Yes\]](#) See Section 5
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)