

# Structured Variational Inference in Continuous Cox Process Models

Virginia Aglietti<sup>1,2</sup>, Edwin V. Bonilla<sup>3</sup>, Theodoros Damoulas<sup>1,2</sup>, Sally Cripps<sup>5</sup>

<sup>1</sup>University of Warwick <sup>2</sup>The Alan Turing Institute <sup>3</sup>CSIRO's Data61 <sup>5</sup>The University of Sydney

V.Aglietti@warwick.ac.uk, Edwin.Bonilla@data61.csiro.au, T.Damoulas@warwick.ac.uk, Sally.Cripps@sydney.edu.au



THE UNIVERSITY OF SYDNEY

OxWaSP

The Alan Turing Institute

WARWICK  
THE UNIVERSITY OF WARWICK

## CONTRIBUTIONS

We propose a scalable *structured* variational inference algorithm for *continuous* sigmoidal Cox processes. Contributions:

- **Scalable inference in continuous input spaces** via a process superposition.
- **Efficient structured posterior estimation** giving a posterior capturing the complex variable dependencies in the model
- **State-of-the-art performance** when compared to alternative inference schemes, link functions, augmentation schemes and representations of the input space.

## THE LIKELIHOOD FUNCTION

**Discrete likelihood:**

$$p(\mathbf{Y}|\mathbf{f}) = \prod_{n=1}^N \text{Poisson}(y_n; \lambda(\mathbf{x}))$$

**Continuous likelihood:**

$$\mathcal{L}(N, \{\mathbf{x}_1, \dots, \mathbf{x}_n\} | \lambda(\mathbf{x})) = \exp\left(-\int_{\tau} \lambda(\mathbf{x}) d\mathbf{x}\right) \prod_{n=1}^N \lambda(\mathbf{x}_n)$$

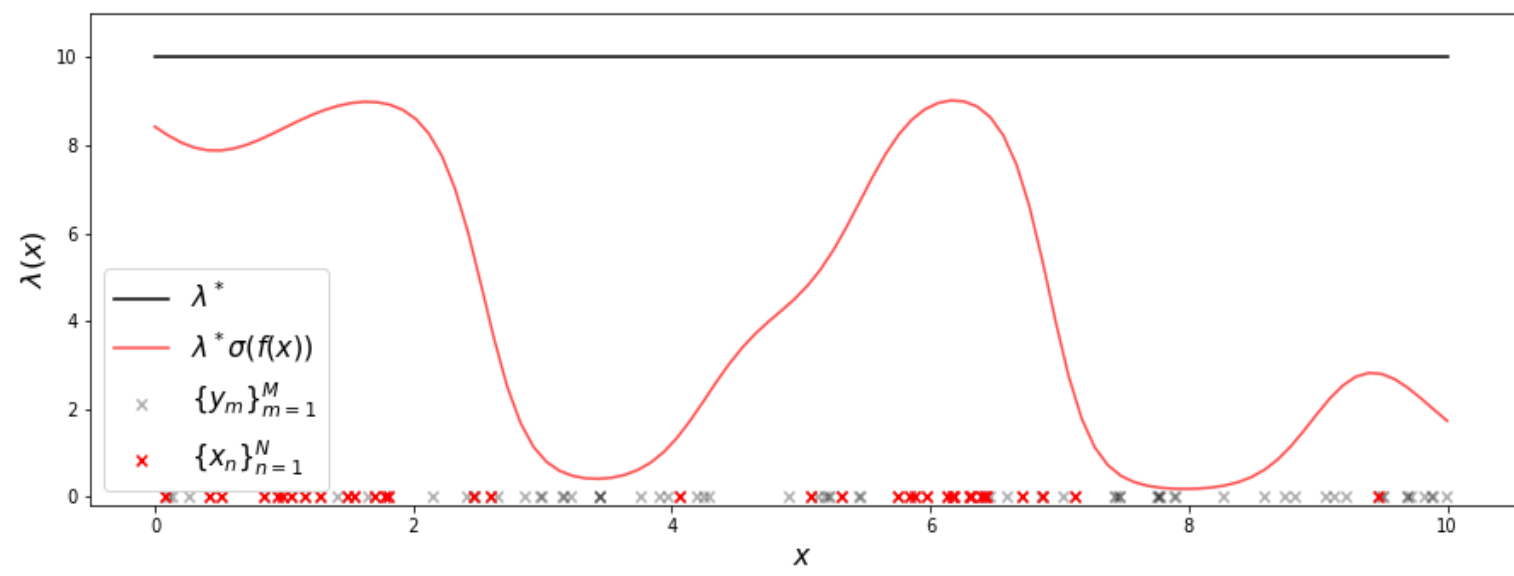


Figure 1: Superposition of two Poisson Point Processes with intensities  $\lambda^*\sigma(f(\mathbf{x}))$  and  $\lambda^*\sigma(-f(\mathbf{x}))$ .

	LGCP	SGCP [1]	Gunter et al.(2014)	VBPP [2]	Lian et al. (2015)	MFVB [3]
Inference	MCMC	MCMC	MCMC	VI-MF	VI-MF	VI-MF
$\mathcal{O}$	$N^3$	$(N+M)^3$	$(N+M)^3$	$NK^2$	$NK^2$	$NK^2$
$\lambda(x)$	$\exp(f(x))$	$\lambda^*\sigma(f(x))$	$\lambda^*\sigma(f(x))$	$(f(x))^2$	$(f(x))^2$	$\lambda^*\sigma(f(x))$
Tractability	$\sum$	Thinning	Adaptive Thinning	Functional form	$\sum$	Integral approximation

Table 1: Summary of related work.  $f$  is continuous,  $\sum$  is discrete.  $M$  represents the number of thinned points derived from the thinning algorithm.  $K$  are the number of inducing inputs.

## Augmentation via superposition

$$\text{Full joint distribution } \mathcal{L}(\{\mathbf{x}_n\}_{n=1}^N, \{\mathbf{y}_m\}_{m=1}^M, M, \mathbf{f}, \lambda^* | \mathcal{X}, \boldsymbol{\theta}):$$

$$\frac{(\lambda^*)^{N+M} \exp(-\lambda^* \int_{\mathcal{X}} d\mathbf{x})}{N!M!} \prod_{n=1}^N \sigma(f(\mathbf{x}_n)) \prod_{m=1}^M \sigma(-f(\mathbf{y}_m)) p(\mathbf{f}) p(\lambda^*)$$

## STRUCTURED VARIATIONAL INFERENCE

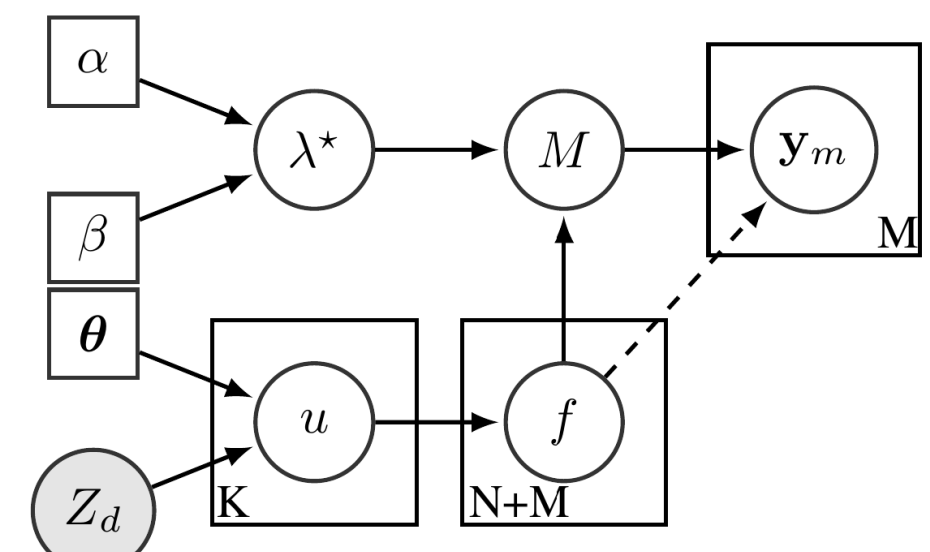


Figure 2: Posterior distribution accounting for all model dependencies. The dashed line represents the assumed factorization.

$$Q(\mathbf{f}, \mathbf{u}, M, \{\mathbf{y}_m\}_{m=1}^M, \lambda^*) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})q(\lambda^*) \times q(\{\mathbf{y}_m\}_{m=1}^M|M)q(M|\mathbf{f}, \lambda^*)$$

$$q(\mathbf{u}) = \mathcal{N}(\mathbf{m}, \mathbf{S}) \quad q(\lambda^*) = \text{Gamma}(\alpha, \beta)$$

$$q(\{\mathbf{y}_m\}_{m=1}^M|M) = \prod_{m=1}^M \sum_{s=1}^S \pi_s \mathcal{N}_T(\mu_s, \sigma_s^2; \mathcal{X})$$

$$q(M|\mathbf{f}, \lambda^*) = \text{Poisson}(\eta) \quad \eta = \lambda^* \int_{\mathcal{X}} \sigma(-f(\mathbf{x})) d\mathbf{x}$$

## THE EVIDENCE LOWER BOUND

$$\mathcal{L}_{\text{elbo}} = T_0 + \underbrace{\mathbb{E}_Q[M \log(\lambda^*)]}_{T_1} - \underbrace{\mathbb{E}_Q[\log(M!)]}_{T_2}$$

$$+ \sum_{n=1}^N \mathbb{E}_{q(\mathbf{f})}[\log(\sigma(f(\mathbf{x}_n)))] + \underbrace{\mathbb{E}_Q\left[\sum_{m=1}^M \log(\sigma(-f(\mathbf{y}_m)))\right]}_{T_3}$$

$$- \underbrace{\mathcal{L}_{\text{kl}}^{\mathbf{u}} - \mathcal{L}_{\text{kl}}^{\lambda^*}}_{T_4} - \underbrace{\mathcal{L}_{\text{ent}}^M - \mathcal{L}_{\text{ent}}^{\{\mathbf{y}_m\}_{m=1}^M}}_{T_5}$$

where  $T_0 = N(\psi(\alpha) - \log(\beta)) - V \frac{\alpha}{\beta} - \log(N!)$ ,  $V = \int_{\mathcal{X}} d\mathbf{x}$ ,  $\psi(\cdot)$  is the digamma function and  $q(\mathbf{f}) = \mathcal{N}(\mathbf{A}\mathbf{m}, \mathbf{K}_{xx} - \mathbf{A}\mathbf{K}_{zx} + \mathbf{A}\mathbf{S}\mathbf{A}^T)$ .

We derive expressions for  $T_i$ ,  $i = 1, \dots, 5$  that avoid sampling from the full joint posterior and computing the GP on the stochastic locations.

**Time complexity:**  $\mathcal{O}(K^3)$  **Space complexity:**  $\mathcal{O}(K^2)$

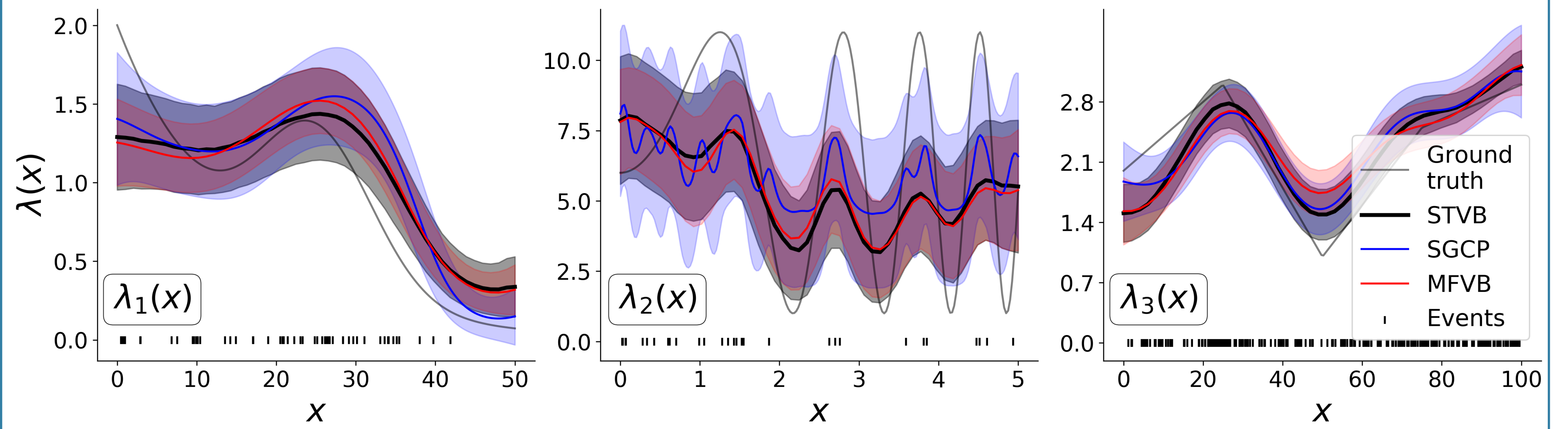
## KEY REFERENCES

[1] Adams, R. P., Murray, I., and MacKay, D. J. *Tractable nonparametric bayesian inference in poisson processes with gaussian process intensities* Proceedings of the 26th Annual International Conference on Machine Learning, pages 9–16 (2009).

[2] Lloyd, C., Gunter, T., Osborne, M. A., and Roberts, S. J. *Variational Inference for Gaussian Process Modulated Poisson Processes* International Conference on Machine Learning, pages 1814–1822 (2015).

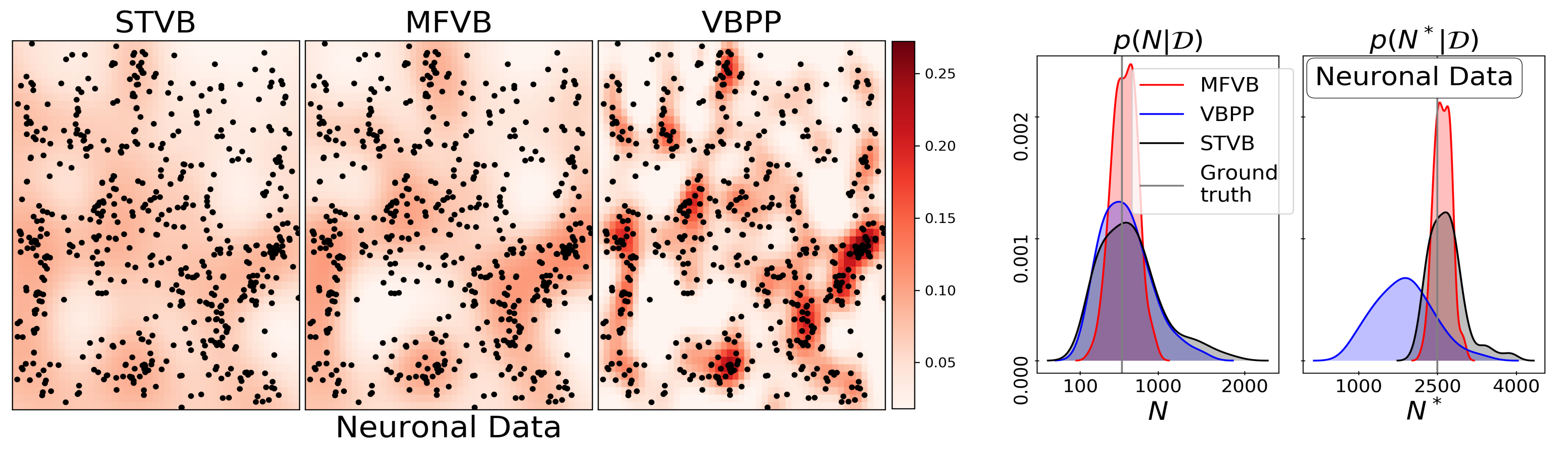
[3] Donner, C. and Oppor, M. *Efficient bayesian inference of sigmoidal gaussian cox processes* The Journal of Machine Learning Research, 9(1):2710–2743 (2018).

## SYNTHETIC DATA



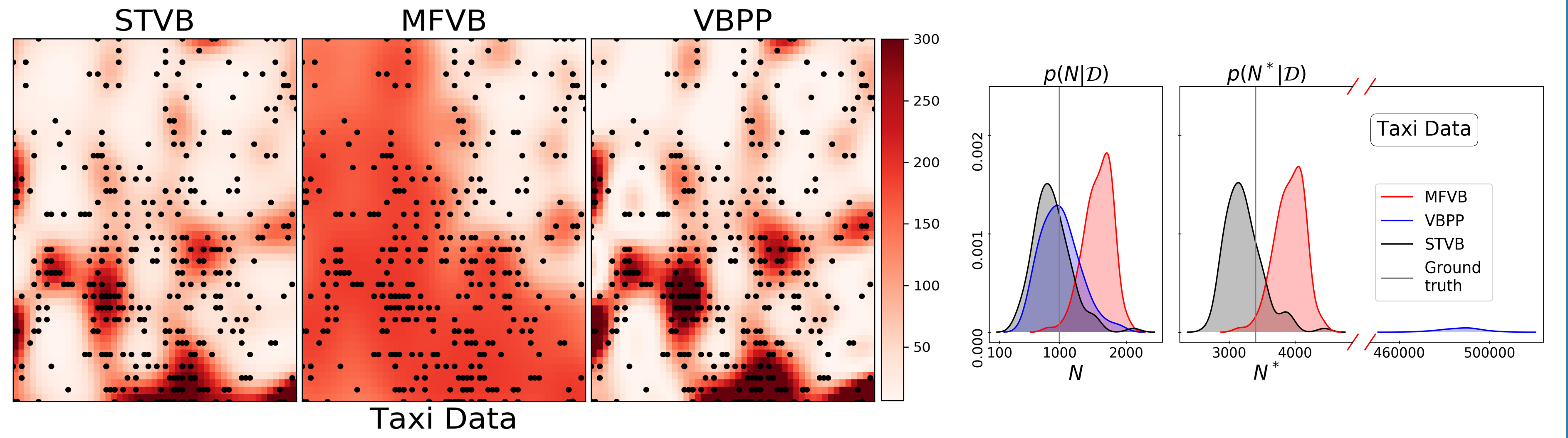
	$\lambda_1(x)$					$\lambda_2(x)$					$\lambda_3(x)$					CPU time (s)
	$l_2$	$\ell_{\text{test}}$	NLPL	EC 30% CI	EC 40% CI	$l_2$	$\ell_{\text{test}}$	NLPL	EC 30% CI	EC 40% CI	$l_2$	$\ell_{\text{test}}$	NLPL	EC 30% CI	EC 40% CI	
STVB	<b>3.44</b>	<b>-1.39</b>	4.71	<b>0.81</b>	<b>0.72</b>	46.28	56.04	5.62	<b>0.91</b>	<b>0.88</b>	<b>7.39</b>	153.98	6.41	<b>0.99</b>	0.97	315.59
	(1.43)	(1.05)	(0.51)	(0.27)	(0.27)		(9.95)	(4.47)	(0.72)	(0.24)		(2.76)	(11.91)	(0.64)	(0.03)	(0.09)
MFVB	4.56	-2.84	4.74	0.76	0.61	44.44	55.35	5.52	0.89	0.84	8.17	155.08	5.82	0.97	0.91	0.01
	(1.43)	(1.0)	(0.1)	(0.25)	(0.28)		(10.7)	(4.72)	(1.29)	(0.23)		(3.43)	(10.20)	(0.61)	(0.09)	(0.14)
VBPP	9.19	-7.71	8.91	0.75	0.41	48.15	<b>56.82</b>	5.20	0.76	0.45	20.54	152.82	8.35	0.83	0.43	0.44
	(2.32)	(3.31)	(1.19)	(0.21)	(0.25)		(13.16)	(4.42)	(1.33)	(0.26)		(6.53)	(11.43)	(2.28)	(0.19)	(0.14)
SGCP	4.22	<b>-1.39</b>	<b>4.21</b>	0.39	0.27	<b>43.50</b>	55.05	<b>3.77</b>	0.64	0.14	<b>14.44</b>	<b>165.66</b>	<b>4.78</b>	0.49	0.34	2764.88
	(1.88)	(1.28)	(1.04)	(0.28)	(0.22)		(8.69)	(1.35)	(0.54)	(0.09)		(2.97)	(2.12)	(0.33)	(0.03)	(0.07)
LGCP	67.76	-5.26	26.26	0.08	0.03	106.74	28.56	15.75	0.04	0.00	19.24	147.67	10.84	<b>0.99</b>	<b>0.99</b>	4.74
	(24.38)	(8.84)	(8.09)	(0.12)	(0.09)		(13.89)	(6.88)	(3.36)	(0.08)		(6.44)	(11.76)	(1.36)	(0.00)	(0.12)

## APPLICATION I: NEURONAL DATA



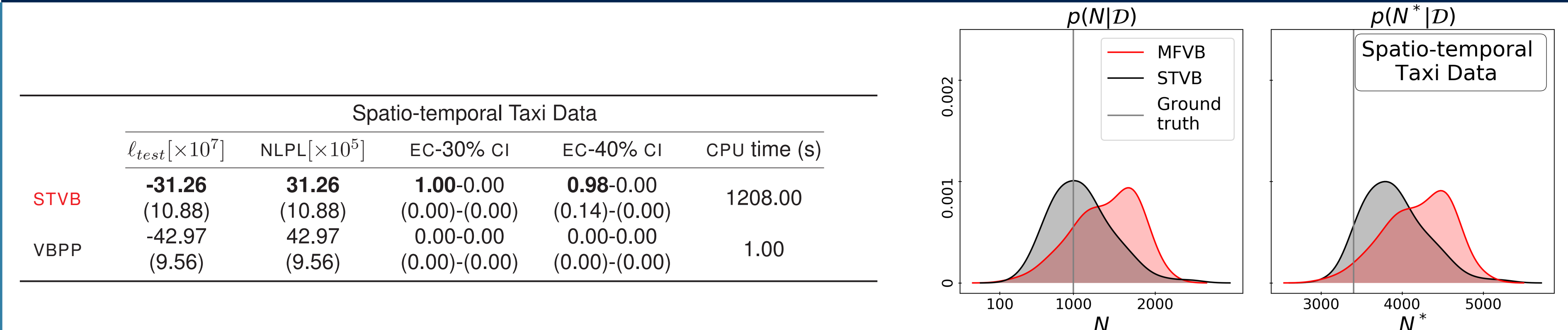
	Neuronal data				
	$\ell_{\text{test}} [\times 10^3]$	NLPL	EC-30% CI	EC-40% CI	CPU time (s)
STVB	-84.55 (16.05)	<b>10.10</b> (7.02)	<b>1.00-1.00</b> (0.00)-(0.00)	<b>0.99-0.56</b> (0.10)-(0.50)	193.07
MFVB	<b>-83.54</b> (4.60)	10.71 (3.39)	<b>1.00-0.03</b> (0.00)-(0.17)	0.78-0.00 (0.41)-(0.00)	0.35
VBPP	-83.89 (12.49)	11.39 (8.18)	<b>1.00-0.00</b> (0.00) - (0.00)	0.83-0.00 (0.38)-(0.00)	26.23

## APPLICATION II: TAXI DATA



	Taxi data				
	$\ell_{\text{test}} [\times 10^3]$	NLPL	EC-30% CI	EC-40% CI	CPU time (s)
STVB	<b>-27.96</b> (9.16)	<b>27.96</b> (9.16)	<b>0.81-0.37</b> (0.39)-(0.48)	<b>0.09-0.01</b> (0.29)-(0.10)	290.34
MFVB	-40.8 (6.41)	40.65 (6.41)	0.00-0.00 (0.00)-(0.00)	0.00-0.00 (0.00)-(0.00)	0.24
VBPP	-31.32 (8.18)	31.32 (8.18)	<b>0.98-0.00</b> (0.14)-(0.00)	<b>0.48-0.00</b> (0.50)-(0.00)	3.62

## APPLICATION III: SPATIO-TEMPORAL TAXI DATA



	Spatio-temporal Taxi Data				
	$\ell_{\text{test}} [\times 10^7]$	NLPL [ $\times 10^5$ ]	EC-30% CI	EC-40% CI	CPU time (s)
STVB	<b>-31.26</b> (10.88)	<b>31.26</b> (10.88)	<b>1.00-0.00</b> (0.00)-(0.00)	<b>0.98-0.00</b> (0.14)-(0.00)	1208.00
VBPP	-42.97 (9.56)	42.97 (9.56)	0.00-0.00 (0.00)-(0.00)	0.00-0.00 (0.00)-(0.00)	1.00

## FUTURE RESEARCH

- Test the algorithm in higher dimensional settings.
- Develop a scalable fully structured variational inference scheme by relaxing the factorization assumption in the posterior.

## MORE INFORMATION

**Paper:** <https://arxiv.org/pdf/1906.03161.pdf>

**Python Code:** <https://github.com/VirgiAgl/STVB>