



**University of  
Zurich<sup>UZH</sup>**

# Generalised Spatial Fusion Model Framework for Multivariate Analysis of Point and Areal Data

Craig Wang<sup>a</sup>, Prof. Milo Puhan<sup>b</sup>, Prof. Reinhard Furrer<sup>a,c</sup>

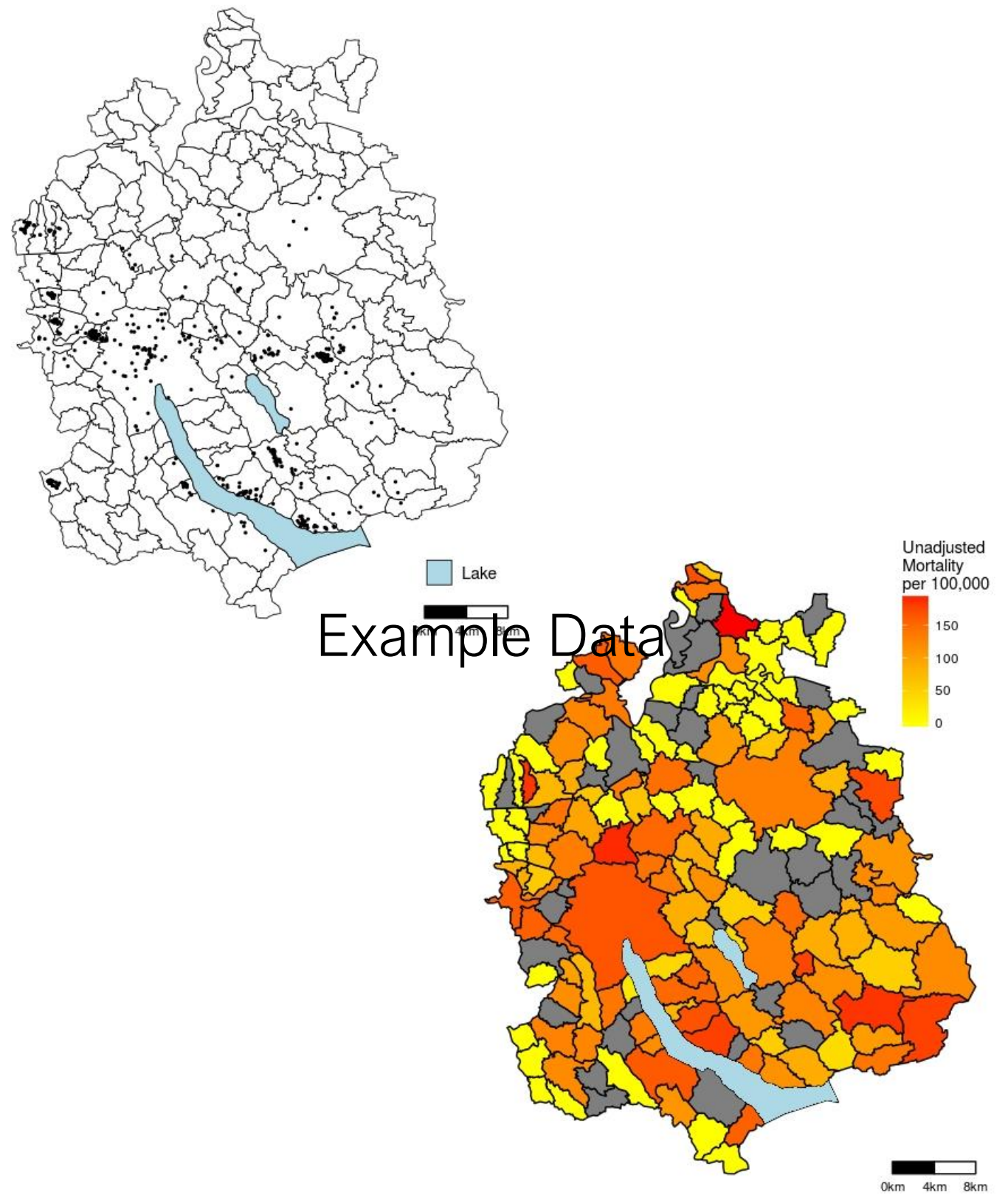
<sup>a</sup> Applied Statistics Group, Department of Mathematics, University of Zurich, Switzerland

<sup>b</sup> Epidemiology, Biostatistics and Prevention Institute, University of Zurich, Switzerland

<sup>c</sup> Department of Computational Science, University of Zurich, Switzerland

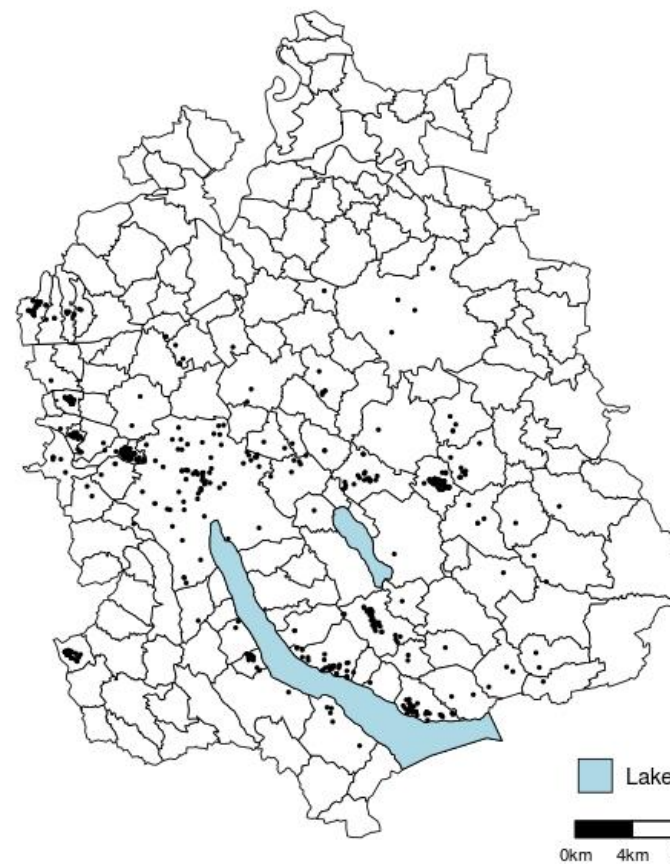


# Motivations



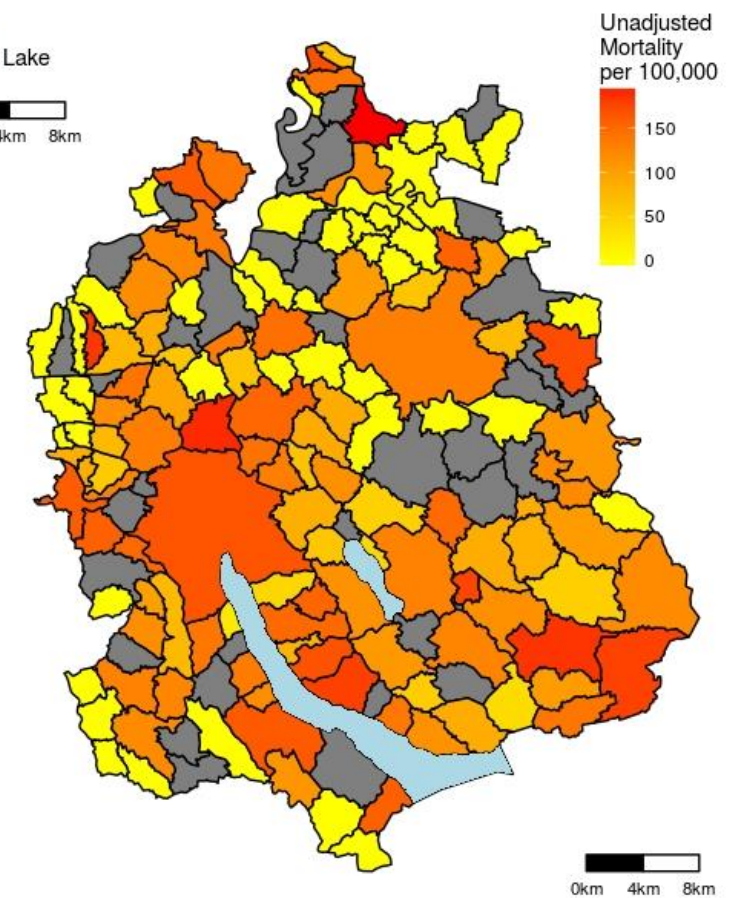
# Motivations

- Scarce observations



Few data points

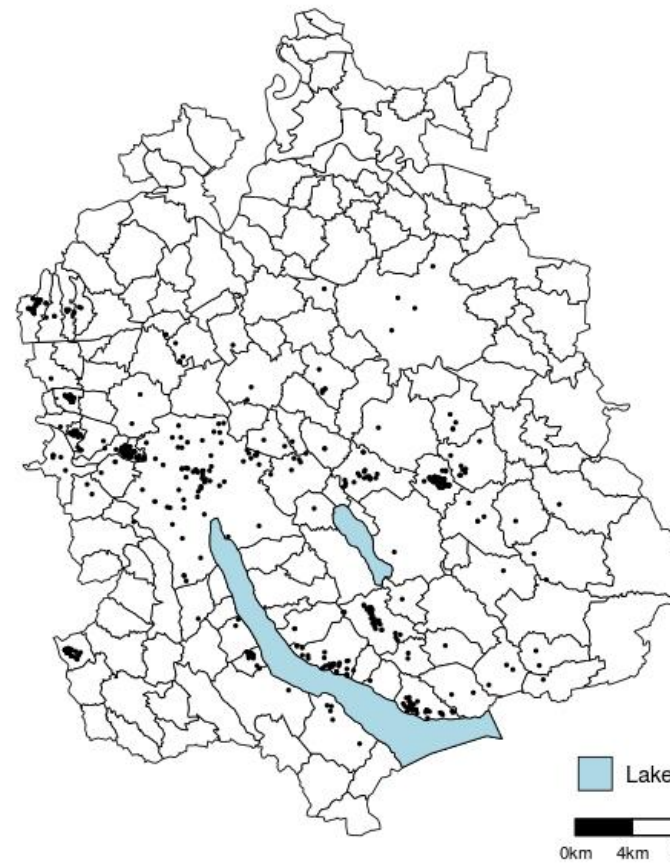
Missing areas





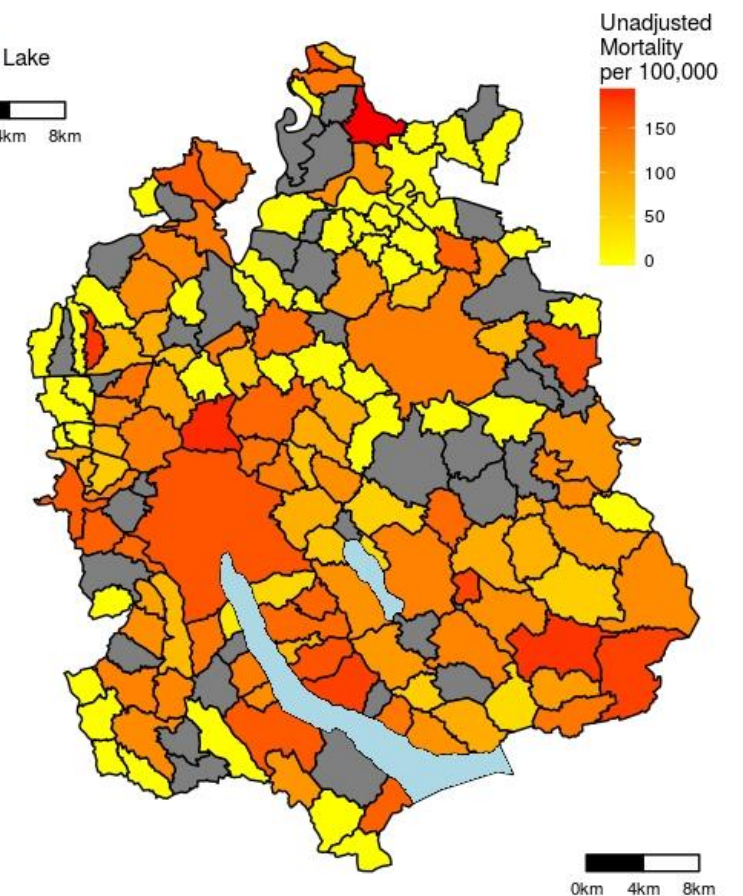
# Motivations

- Scarce observations
- Multiple data sources
- Different resolutions



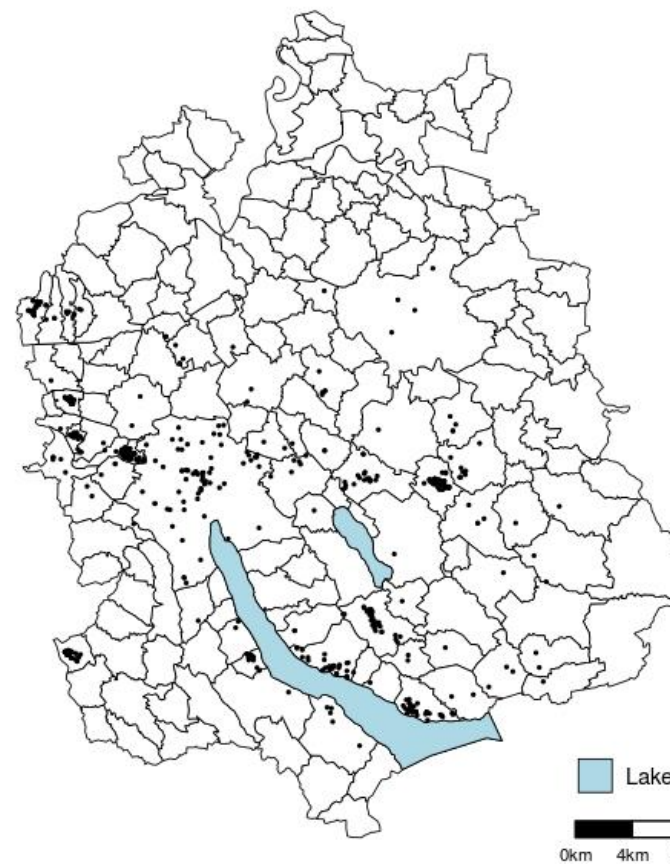
Source 1: Point data

Source 2: Areal data



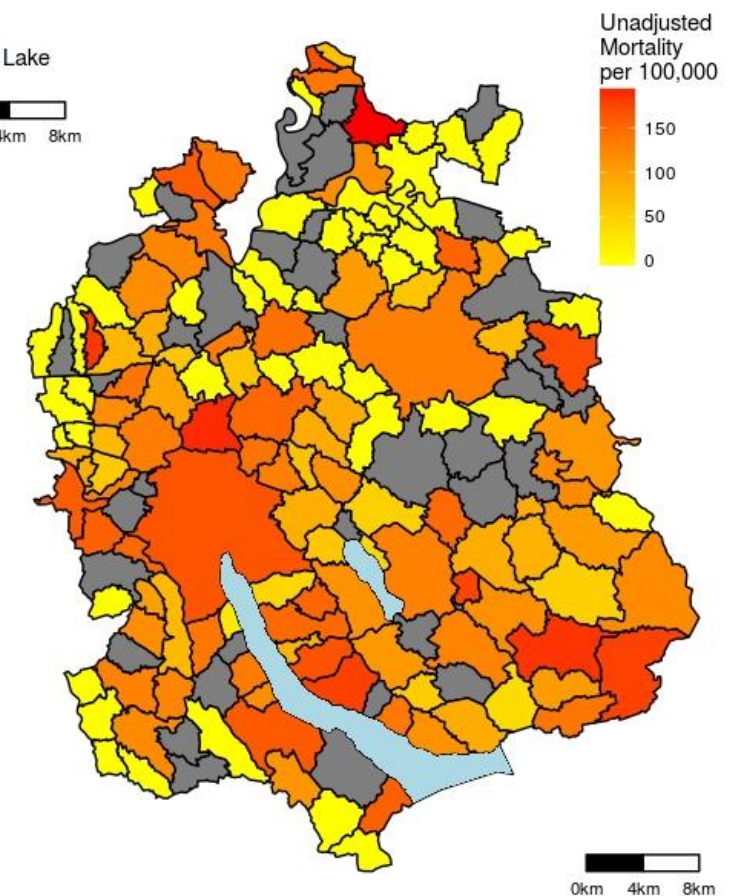
# Motivations

- Scarce observations
- Multiple data sources
- Different resolutions
- Different distributions



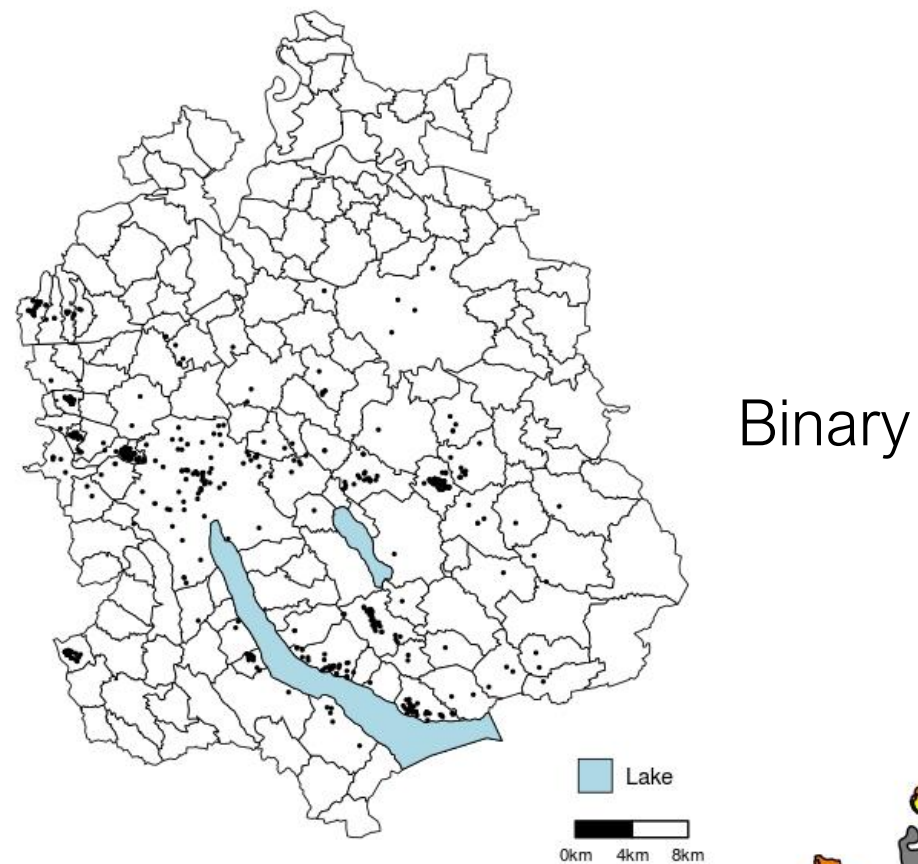
Lung function ~ Normal

Mortality ~ Poisson

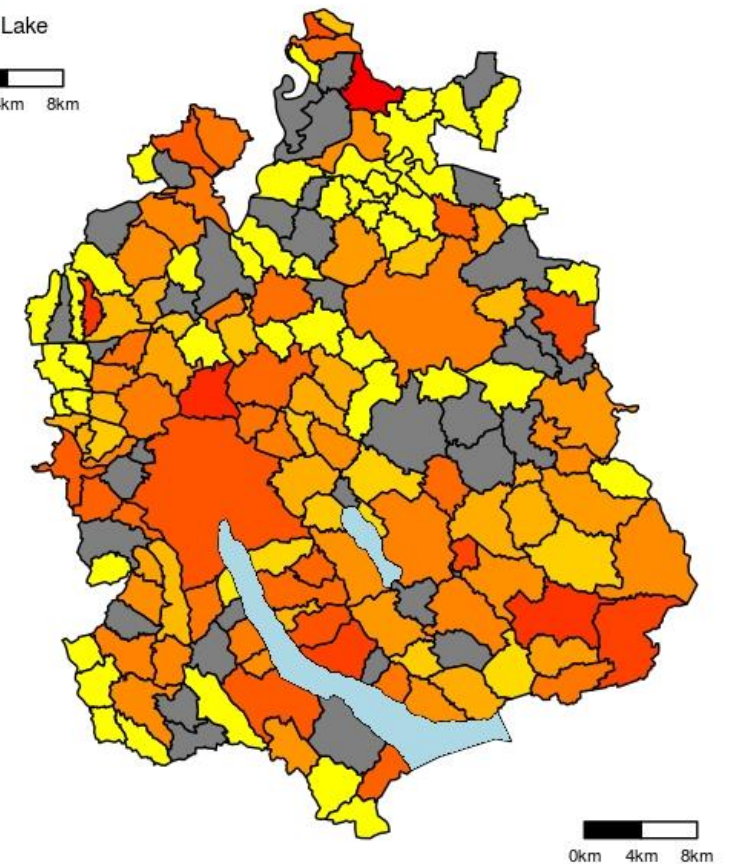


# Motivations

- Scarce observations
- Multiple data sources
- Different resolutions
- Different distributions



Counts





# Existing Spatial Fusion Models

- Also known as data assimilation (Banerjee et al. 2014) or Bayesian melding (Fuentes and Raftery, 2005)
- Fusion models include Sahu et al. (2010), Berrocal et al. (2010), Goovaerts (2010), Liu et al. (2011), Moraga et al. (2017), and Shi & Kang (2017)
- Some drawbacks
  - Constraint on distribution assumption
  - Model specification for single application
  - Designing of custom Monte Carlo sampler

# A Framework for Spatial Fusion Models

- Point data  $\mathbf{s} = \{s_1, \dots, s_n\}$

$$f(\mathbb{E}[Y(\mathbf{s})|w(\mathbf{s})]) = \mathbf{X}_s^T(\mathbf{s})\boldsymbol{\beta}_s + \underline{w(\mathbf{s})}$$

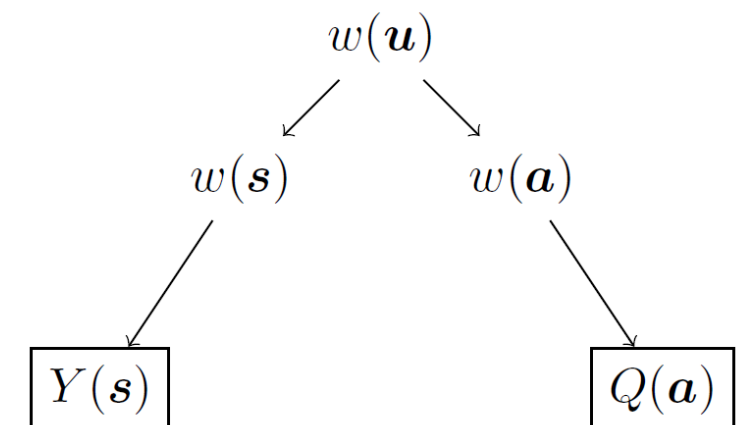
- Areal data  $\mathbf{a} = \{a_1, \dots, a_m\}$

$$g(\mathbb{E}[Q(\mathbf{a})|w(\mathbf{a})]) = \mathbf{X}_a^T(\mathbf{a})\boldsymbol{\beta}_a + \underline{w(\mathbf{a})}$$

- Data fusion via sampling points  $\mathbf{s}' = \{s'_1, \dots, s'_l\}$

$$\underline{w(a_i)} = \int_{\mathbf{s} \in a_i} w(\mathbf{s}) d\mathbf{s} \approx \frac{1}{L} \sum_{j=1, s'_j \in a_i}^L \underline{w(s'_j)}$$

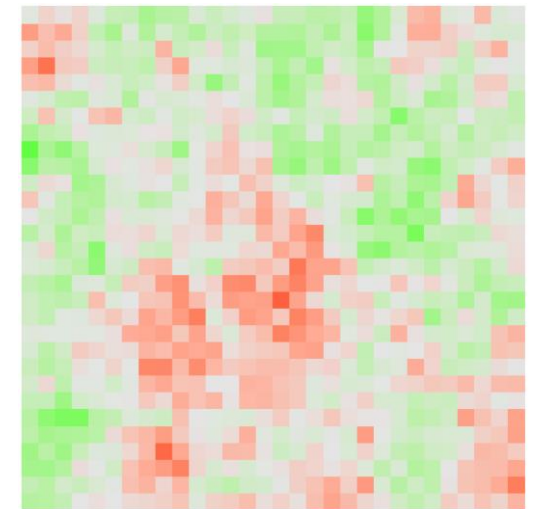
$$w(\mathbf{u}) = [\underline{w(\mathbf{s})} \quad \underline{w(\mathbf{s}')} ] \sim GP(0, C(\cdot, \cdot; \theta))$$



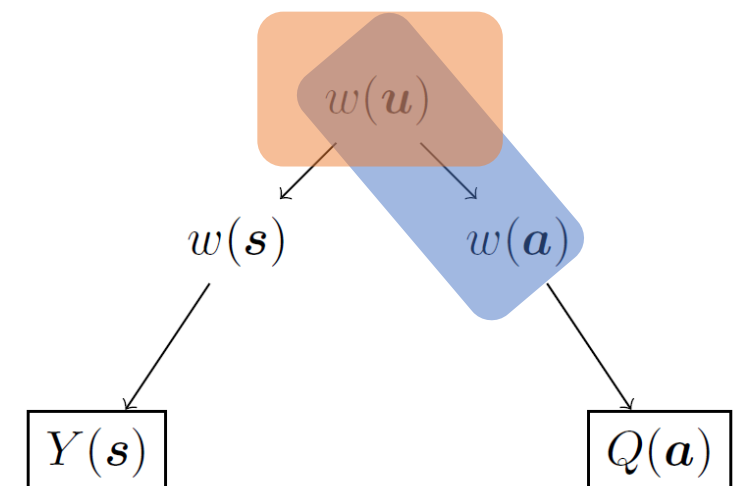


# Modelling Approach

- Bayesian hierarchical model
- Common latent spatial process
- Approximation to stochastic integral
- Efficient Computation
  - Nearest neighbour Gaussian Process
  - Stan modelling language



A simulated Gaussian spatial process.



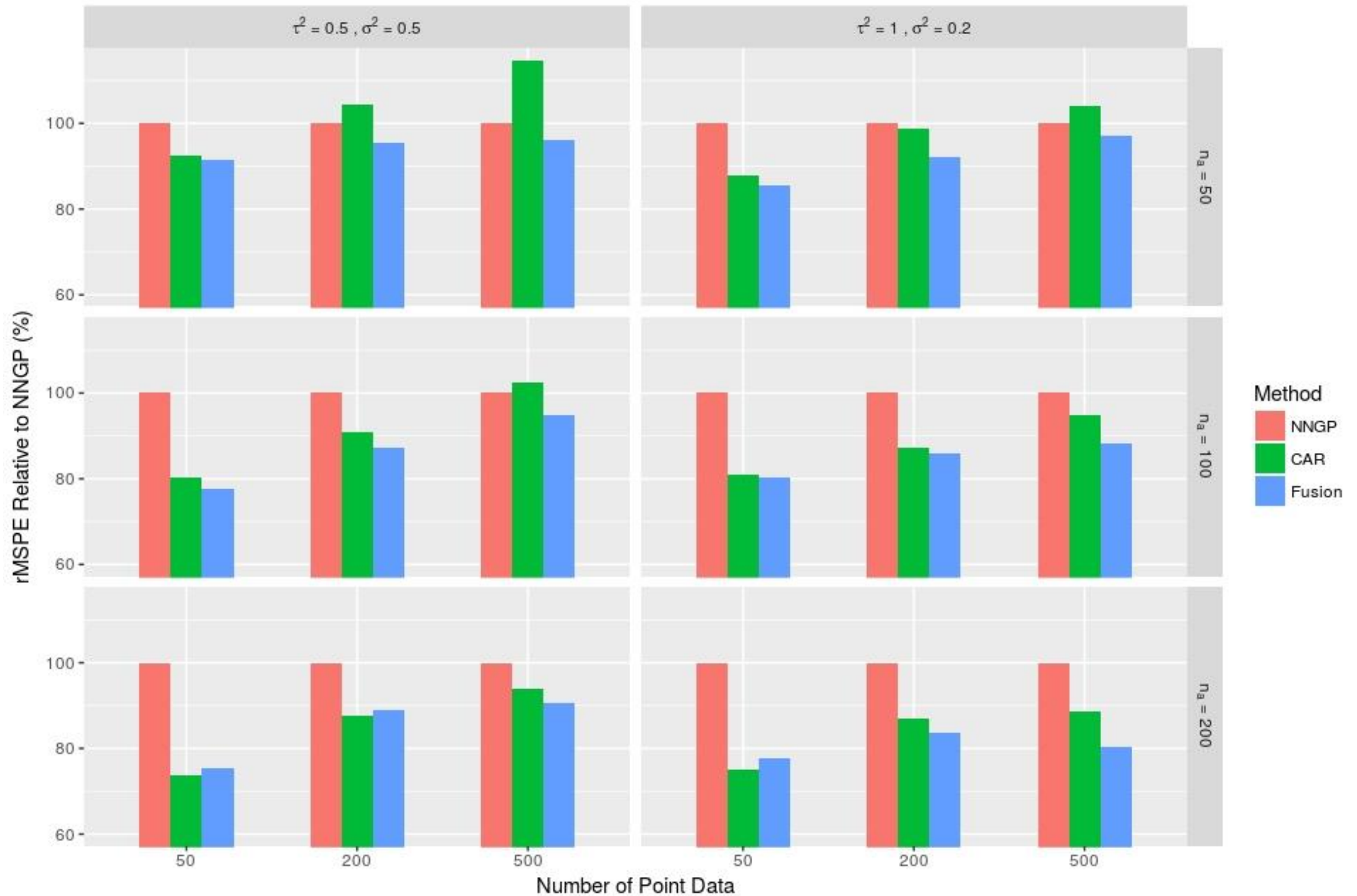
# Simulation Study

$$Y(\mathbf{s})|\boldsymbol{\beta}_s, w(\mathbf{s}) \sim \text{Normal}(X_s^\top \boldsymbol{\beta}_s + w(\mathbf{s}), \tau^2)$$

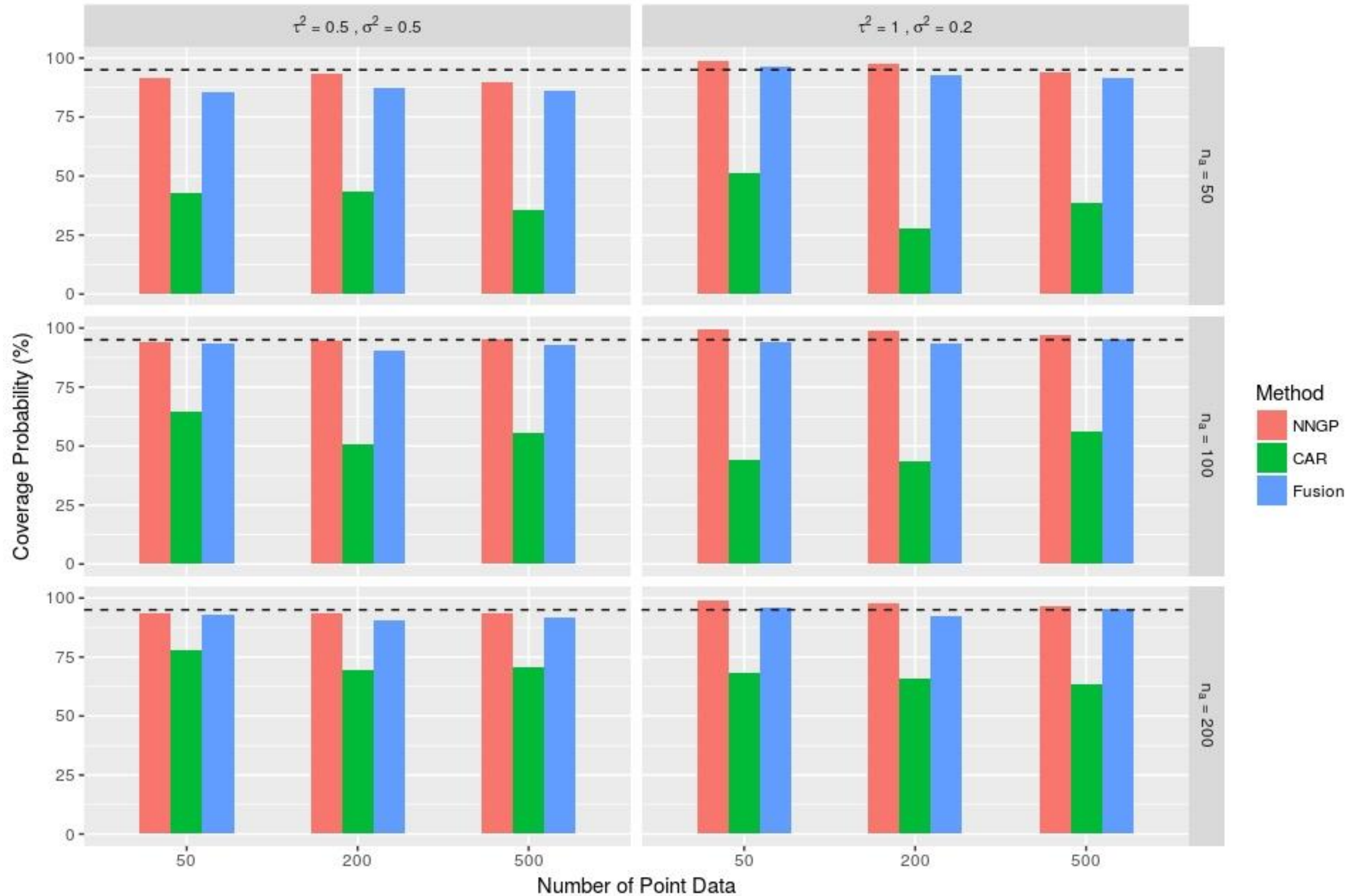
$$Q(\mathbf{a})|\boldsymbol{\beta}_a, w(\mathbf{a}) \sim \text{Poisson}(X_a^\top \boldsymbol{\beta}_a + w(\mathbf{a}))$$

- 18 different scenarios
  - Different sample size
  - Different spatial signal-to-noise ratio
- Methods
  - nearest neighbor Gaussian process (NNGP) – point data only
  - conditional autoregressive model (CAR) + kriging – areal data only
  - fusion model – both point and areal data
- Evaluation criteria on latent process prediction  $w(\mathbf{s})$ 
  - root mean squared prediction error (rMSPE)
  - coverage probability of 95% posterior credible intervals

# Simulation Performance - rMSPE



# Simulation Performance - Coverage





# Case Study: LuftiBus - SNC

- LuftiBus dataset
  - Time: 2003 to 2012
  - Location: Switzerland
  - Variables: spirometry, demographics
- Swiss National Cohort (SNC)
  - Long-term population-based cohort
  - Variables: mortality, demographics
- Variables of interest
  - Point: Forced Expiratory Volume in one second (FEV1)
  - Areal: Mortality caused by respiratory disease and lung cancer

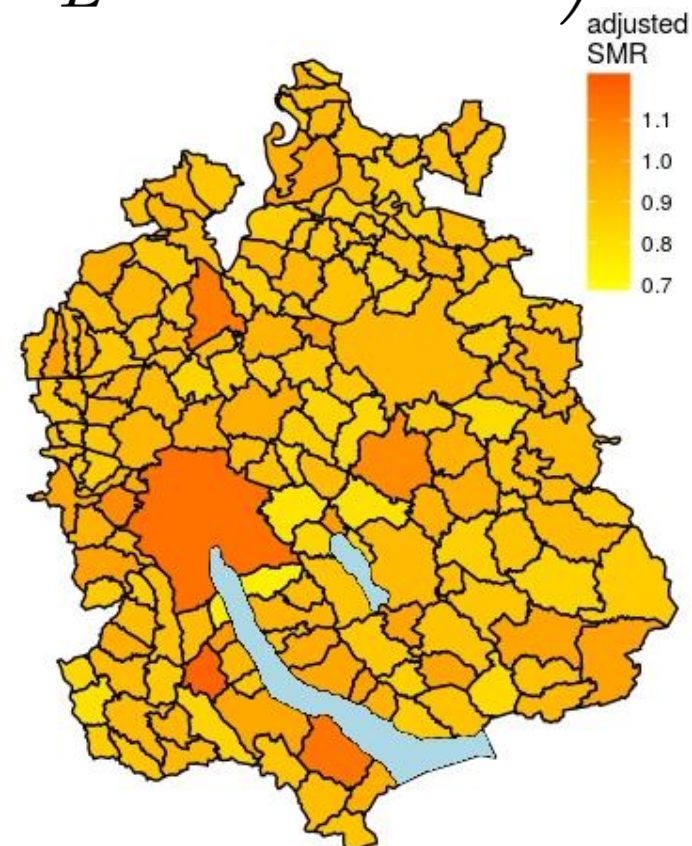
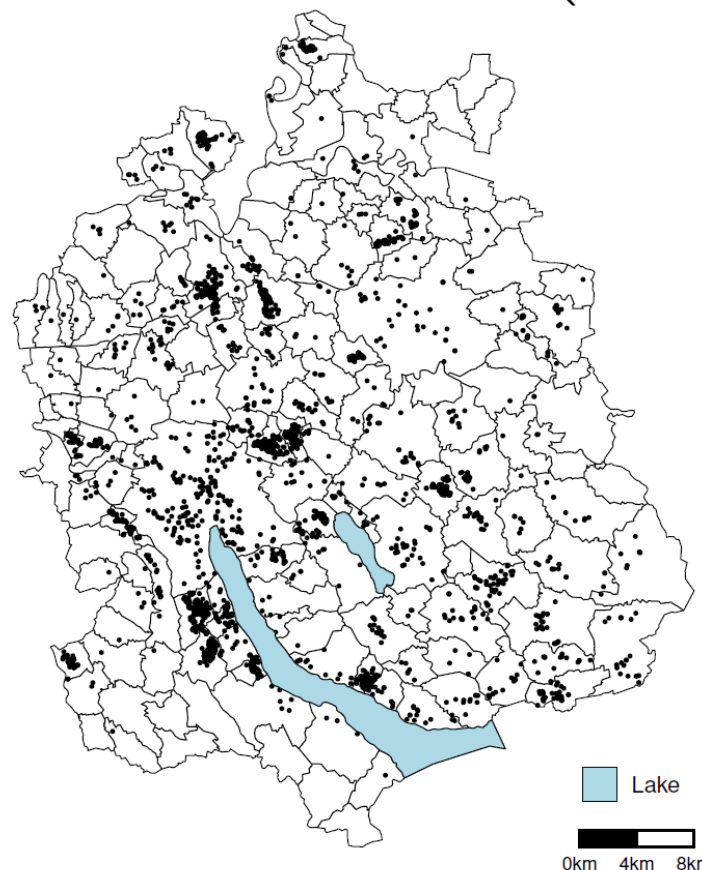


# Case Study: Fusion Modelling

What is the spatial distribution of underlying risk in respiratory disease and lung cancer?

$$\text{FEV1} \sim \text{Normal}(\beta_0 + \beta_1 \times \text{gender} + \beta_2 \times \text{age} - w(\mathbf{s}), \tau^2)$$

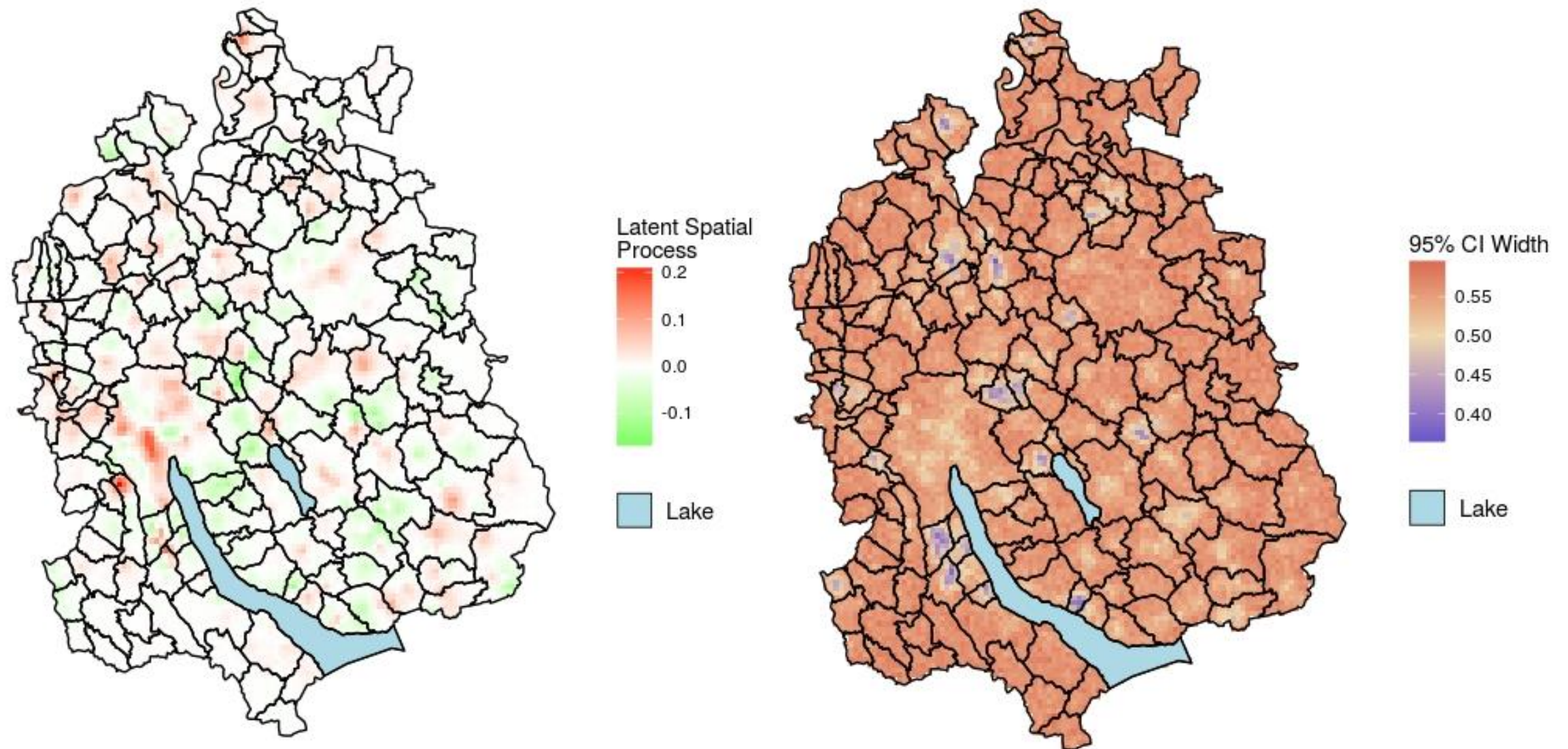
$$\text{Mortality} \sim \text{Poisson}\left(\text{Expected} \times \frac{1}{L} \sum \exp(w(\mathbf{s}'))\right)$$





# Case Study: Results

Posterior Median of the Latent Spatial Process and its 95% CI Width



# Summary

- Introduced generalised spatial fusion framework
  - Multivariate analysis of spatial data with different resolution and distribution
  - Mimicking data generating process, account for ecological bias
  - Can be easily adapted to different problems
- Utilized Stan modelling language and NNGP
- Simulation showed fusion models perform better in latent process prediction
- If you have “fusible” data, try it out!



# References

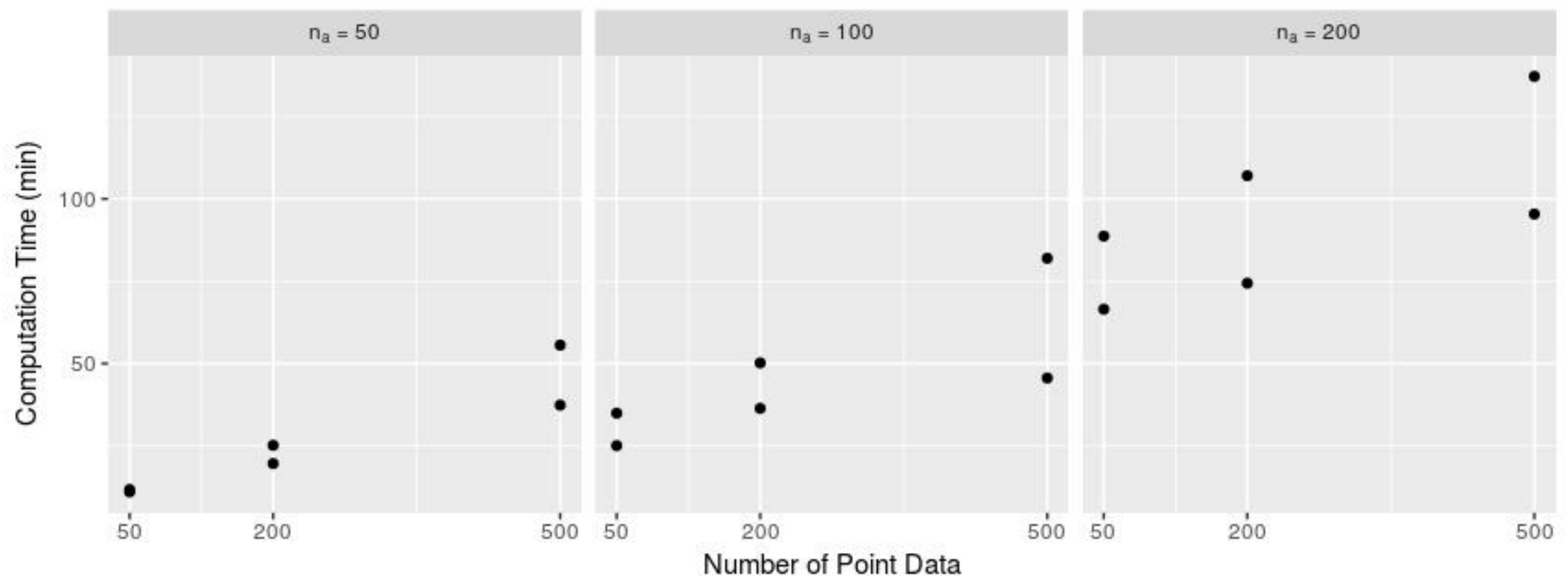
- Banerjee, Carlin & Gelfand, (2014). Hierarchical Modeling and Analysis for Spatial Data. CRC.
- Berrocal, Gelfand & Holland, (2010), *A spatio-temporal downscaler for output from numerical models*, Journal of Agricultural, Biological, and Environmental Statistics.
- Datta, Banerjee, Finley & Gelfand, (2016), *Hierarchical nearest-neighbor Gaussian process models for large geostatistical datasets*, Journal of the American Statistical Association.
- Fuentes & Raftery, (2005), *Model evaluation and spatial interpolation by Bayesian combination of observations with outputs from numerical models*, Biometrics.
- Goovaerts, (2010). *Combining areal and point data in geostatistical interpolation: Applications to soil science and medical geography*. Mathematical Geosciences.
- Liu, Le & Zidek, (2011), *An empirical assessment of Bayesian melding for mapping ozone pollution*, Environmetrics.
- Moraga, Cramb, Mengersen & Pagano, (2017), *A geostatistical model for combined analysis of point-level and area-level data using INLA and SPDE*, Spatial Statistics.
- Sahu, Gelfand & Holland, (2010), *Fusing point and areal level space-time data with application to wet deposition*, Journal of the Royal Statistical Society Series C.
- Shi & Kang, (2017), *Spatial data fusion for large non-Gaussian remote sensing datasets*, Stat.
- Wang, Puhan & Furrer (2018), *Generalized spatial fusion model framework for joint analysis of point and areal data*, Spatial Statistics.

# Backup Slides

- Accounting for ecological bias

$$\exp(w(a_i)) = |a_i|^{-1} \int_{\mathbf{s} \in a_i} \exp(w(\mathbf{s})) d\mathbf{s} \approx \frac{1}{L} \sum_{j=1, s'_{ij} \in a_i}^L \exp(w(s'_{ij}))$$

- Computation time for fusion model



# Stan Model Code

```
parameters{
  vector[pn] beta; // coefficients
  real<lower = 0> sigma_sq;
  real<lower = 0> tau_sq;
  real<lower = 0> phi;
  vector[n+aL] w;
}

transformed parameters{
  vector[n] A0w;
  vector[a] A1w;
  A0w = w[1:n];
  A1w = log(A1 * exp(w));
}

model{
  beta ~ normal(0, 5);
  sigma_sq ~ inv_gamma(2, 0.1);
  tau_sq ~ inv_gamma(2, 1);
  phi ~ normal(700, 100);
  w ~ nngp_w(sigma_sq, phi, neardist, neardistM, nearind, n+aL, M);
  Y ~ normal(Xn * beta - A0w, sqrt(tau_sq));
  Q ~ poisson_log(offset + A1w);
}
```

# Nearest neighbor Gaussian process

- Datta et al. (2016) proposed a low rank process that achieves  $O(nk^3)$  where  $k$  is the number of nearest neighbors

- Instead of specifying  $w(\mathbf{s})$  with a  $n \times n$  covariance matrix,

$$w(\mathbf{s}) \sim N \left( C_{\mathbf{s}, N(\mathbf{s})} C_{N(\mathbf{s})}^{-1} w_{N(\mathbf{s})}, C(\mathbf{s}, \mathbf{s}) - C_{\mathbf{s}, N(\mathbf{s})} C_{N(\mathbf{s})}^{-1} C_{N(\mathbf{s}), \mathbf{s}} \right)$$

where  $C_{\mathbf{s}, N(\mathbf{s})}$  is the  $1 \times k$  cross-covariance matrix between location  $\mathbf{s}$  and each neighbors, and  $C_{N(\mathbf{s})}$  is the  $k \times k$  covariance matrix of  $w_{N(\mathbf{s})}$

- Computational complexity comparison with  $n = 10,000$ 
  - Spatial regression with full Gaussian process  $O(10^{12})$
  - Gaussian predictive process with 500 knots  $O(2.5 \times 10^9)$
  - Nearest neighbor Gaussian process with 10 neighbors  $O(10^7)$