



MOTIVATION

- Gaussian random fields (GRFs) are valuable mathematical tools to fit spatial datasets for applications dealing with measurements regularly or irregularly located across geographical regions.
- The Exascale GeoStatistics project (ExaGeoStat) aims at optimizing the likelihood function for spatial data on massively parallel systems.
- The likelihood evaluation may require generating a covariance matrix Σ and performing the Cholesky factorization of Σ , which necessitates $O(n^3)$ operations and $O(n^2)$ memory. This implies that some of these standard methods for GRFs are not capable to cope with large datasets.
- The objective of ExaGeoStat is to enable statisticians to tackle computationally challenging scientific problems at large scale, while abstracting the hardware complexity.
- ExaGeoStatR further raises the game by providing wrappers in R, i.e., the *de facto* programming environment for statisticians, to the ExaGeoStat main functions. This renders large scale simulations feasible on massively parallel systems, while keeping the simplicity of the R environment.

PROBLEM STATEMENT

- Suppose Z is a stationary GRF and the main domain $D \subset \mathbb{R}^d$ at n locations, s_1, \dots, s_n . The random vector $\{Z(s_1), \dots, Z(s_n)\}^T$ follows a multivariate Gaussian distribution:

$$\forall \{s_1, \dots, s_n\} \subset D, \quad \{Z(s_1), \dots, Z(s_n)\}^T \sim \mathcal{N}_n(\mu, \Sigma), \quad (1)$$

where μ and Σ are the mean vector and the covariance matrix of the n -dimensional multivariate normal distribution.

- Given μ and Σ , the likelihood of observing $\mathbf{z} = \{z(s_1), \dots, z(s_n)\}^T$ at the n locations is

$$L(\mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mu)^T \Sigma^{-1} (\mathbf{z} - \mu) \right\} \quad (2)$$

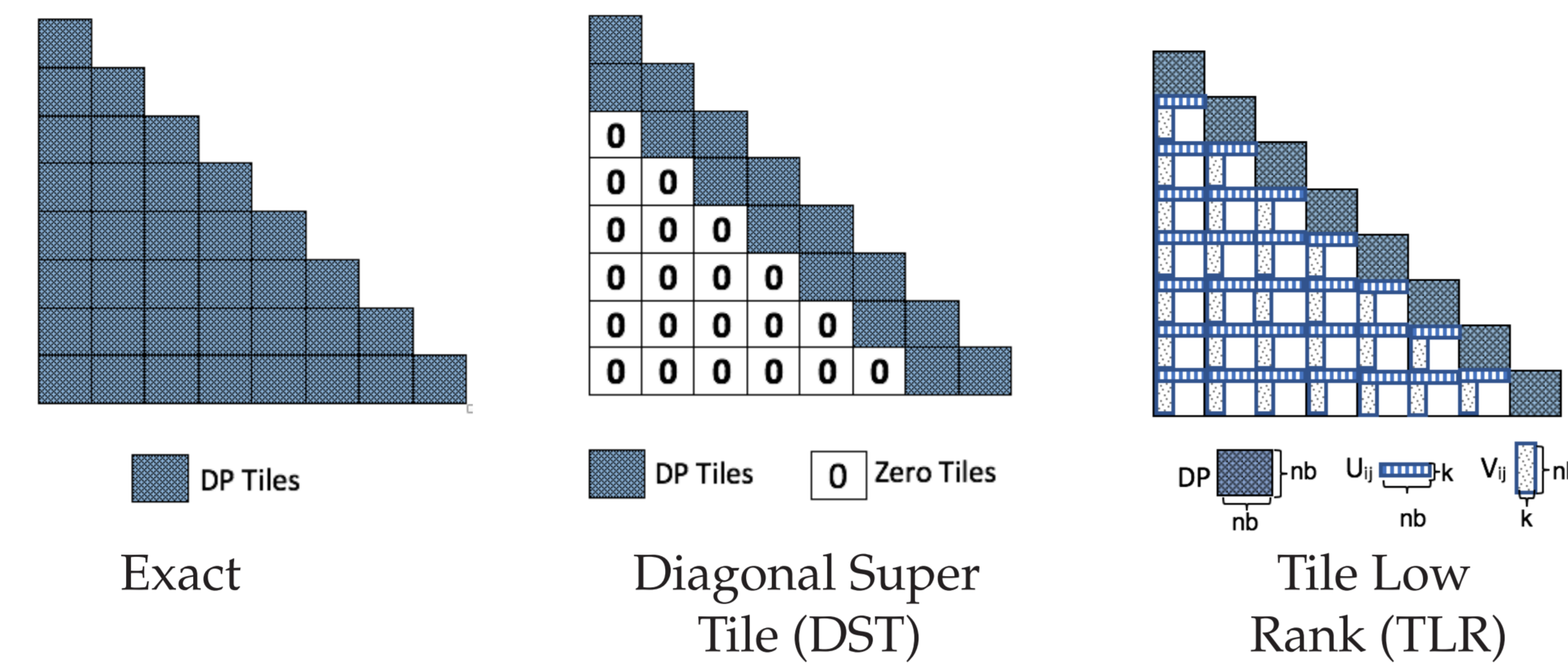
- The (i, j) -th element of Σ is $\Sigma_{ij} = C(s_i, s_j)$, where the covariance function $C(s_i, s_j)$ is assumed to have a parametric form with unknown vector of parameters θ .
- We choose the isotropic Matérn covariance kernel,

$$\Sigma_{ij} = \frac{\sigma^2}{2\nu-1\Gamma(\nu)} \left(\frac{\|s_i - s_j\|}{\beta} \right)^\nu \mathcal{K}_\nu \left(\frac{\|s_i - s_j\|}{\beta} \right), \quad (3)$$

where $\|s_i - s_j\|$ is the distance between s_i and s_j , and the variance, spatial range, and smoothness parameters, i.e., σ^2 , β , and ν , respectively, determine the properties of the GRF.

EXAGEOSTATR FEATURES

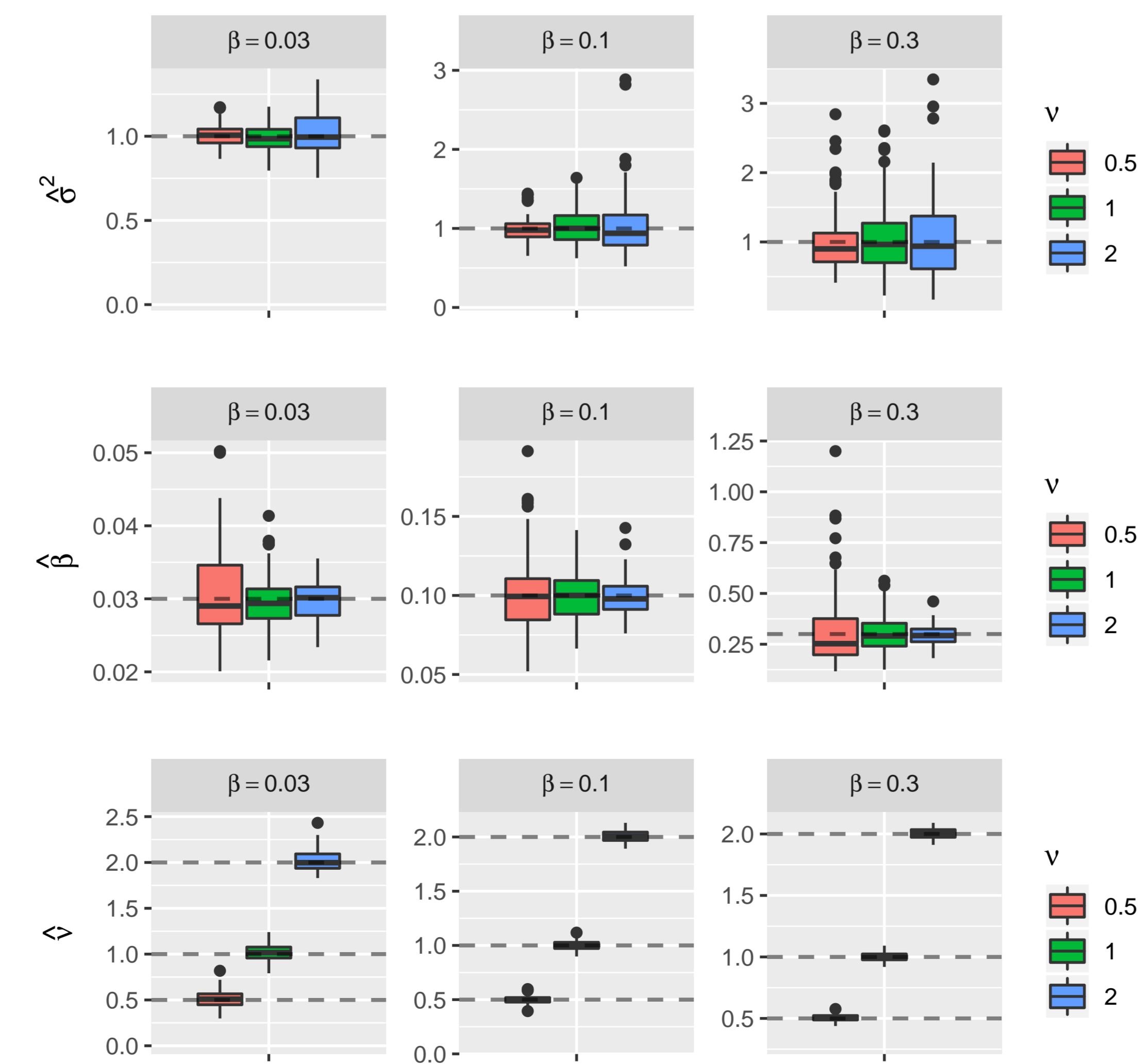
- ExaGeoStatR variant computation techniques:



- Overview of ExaGeoStatR functions:

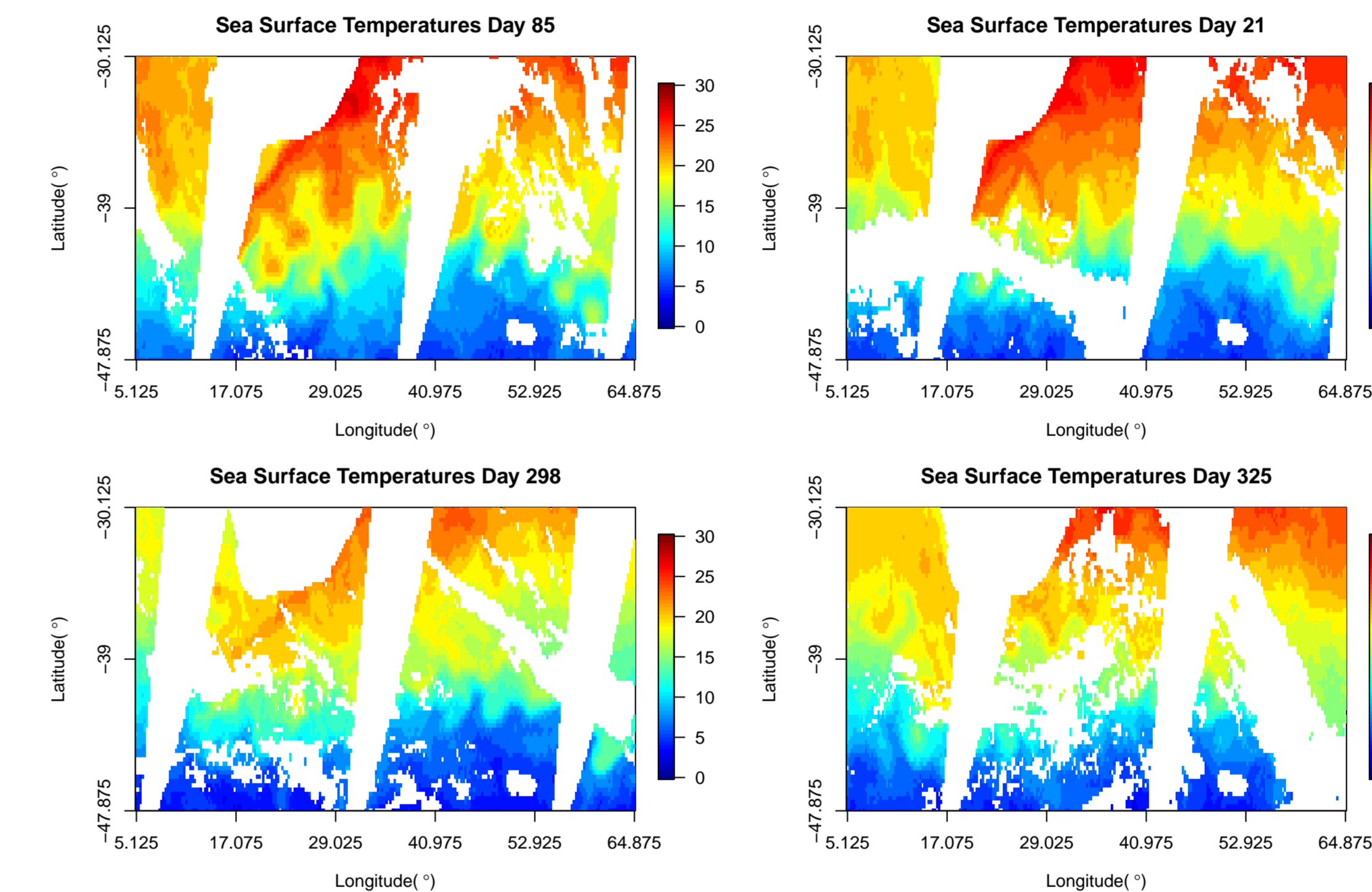
| Function Name | Description |
|---------------------|--|
| exageostat_init | Initiate ExaGeoStat instance. |
| simulate_data_exact | Generate \mathbf{Z} measurements vector. |
| simulate_obs_exact | Generate \mathbf{Z} measurements vector on n given 2D locations. |
| exact_mle | Exact parameter vector evaluation. |
| dst_mle | DST approximation parameter vector evaluation. |
| tlr_mle | TLR approximation parameter vector evaluation. |
| exageostat_finalize | Finalize active ExaGeoStat instance. |

SYNTHETIC DATA ACCURACY ESTIMATION

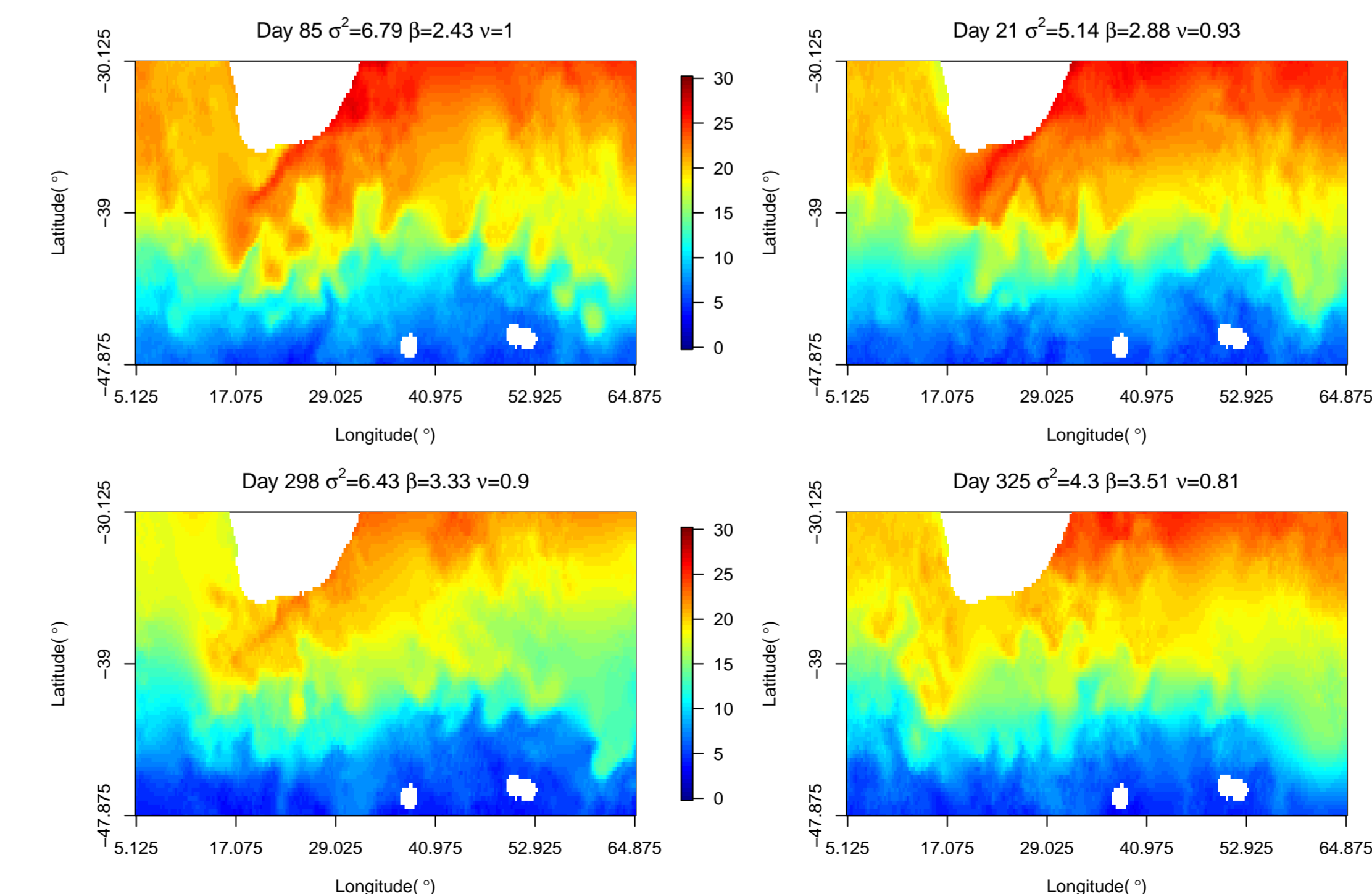


SEA SURFACE TEMPERATURE DATA

- The sea surface temperature collected by satellite for the Agulhas and surrounding areas off the shore of South Africa.
- The region is abstracted into 72×240 regular grid. The spatial resolution is approximately 25 kilometers.
- The data covers 331 days, from Jan. 1st to Nov. 26th, 2004. We selected four days, day 21, day 85, day 298, and day 325.
- The original sea surface temperatures in Celsius ($^{\circ}\text{C}$) where the locations with unavailable values are colored in white.

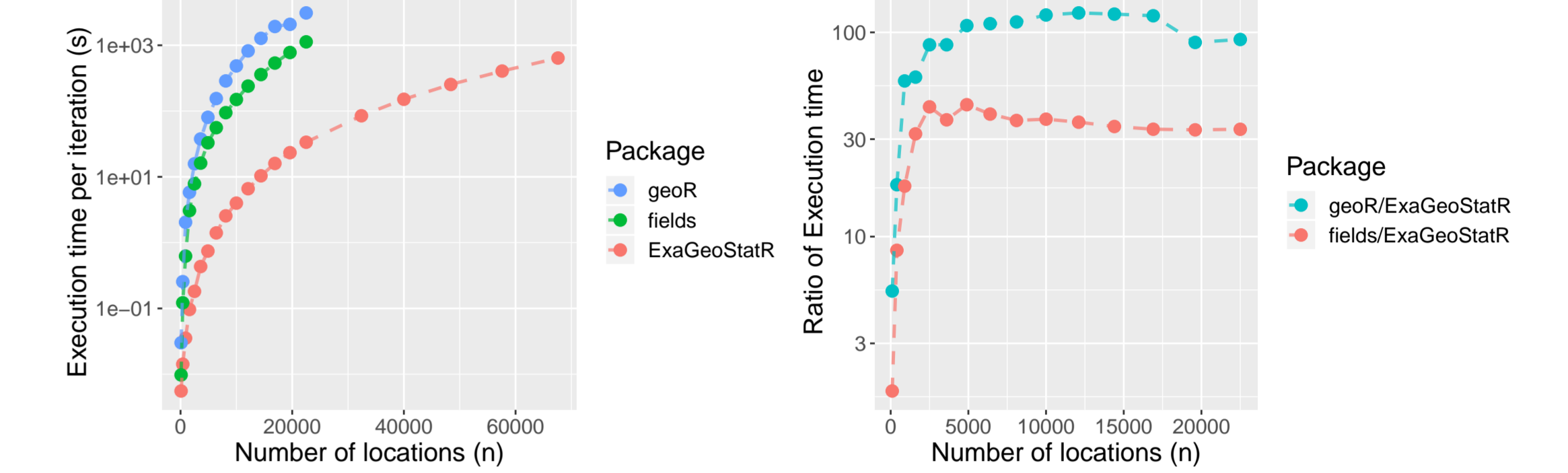


- The predicted sea surface temperatures based on the linear mean structure and the kriging results where the land area is not predicted and colored in white.

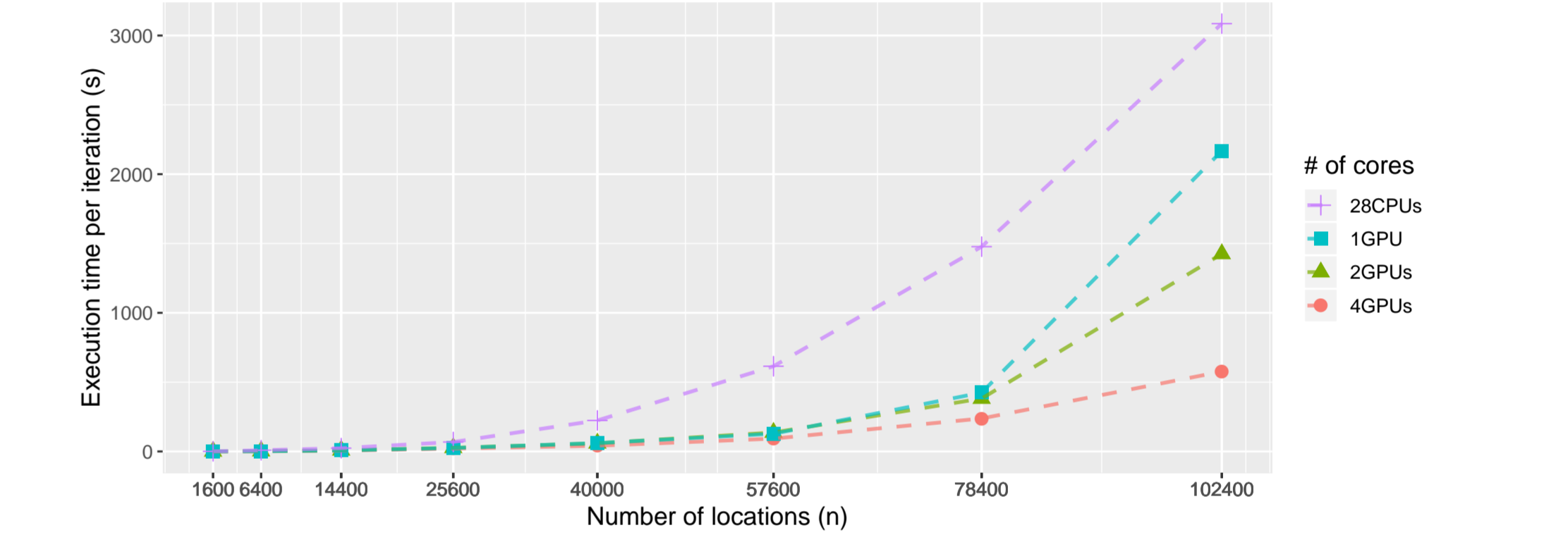


PERFORMANCE RESULTS

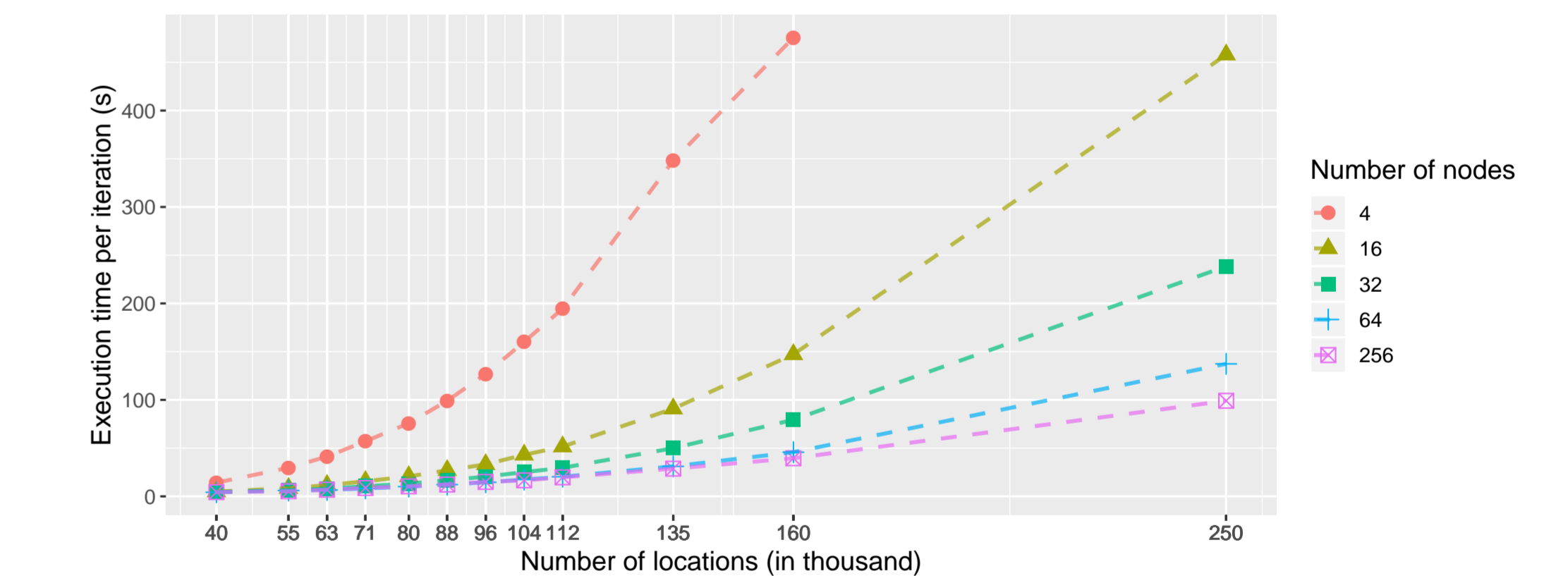
- Performance on a dual-socket 8-core Intel Sandy Bridge processor running at 2.0 GHz:



- Performance on a dual-socket 14-core Intel Broadwell processor running at 2.4 GHz equipped with 8x NVIDIA K80s:



- Performance on Cray XC40 system, each node is a dual-socket 16-core Intel Haswell processor running at 2.3 GHz:



REFERENCES

- S. Abdulah, H. Ltaief, Y. Sun, M. G. Genton, and D. E. Keyes. "ExaGeoStat: A high performance unified software for geostatistics on manycore systems." IEEE Transactions on Parallel and Distributed Systems 29, no. 12 (2018): 2771-2784.
- D. Nychka, R. Furrer, and S. Sain. "fields: Tools for spatial data." National Center for Atmospheric Research (2005).
- P. J. Ribeiro and P. J. Diggle. "geoR: Analysis of geostatistical data." R package version 1.5.2 (2016).