



# A method for quantifying uncertainty in spatially interpolated meteorological data with application to daily maximum air temperature

Conor T. Doherty<sup>1,2</sup>, Weile Wang<sup>1</sup>, Hirofumi Hashimoto<sup>1,2</sup>, and Ian G. Brosnan<sup>1</sup>

<sup>1</sup>NASA Ames Research Center, Moffett Field, CA 94035, USA

<sup>2</sup>Applied Environmental Science, California State University Monterey Bay, Seaside, CA 93955, USA

**Correspondence:** Conor T. Doherty (conor.t.doherty@nasa.gov)

Received: 19 June 2024 – Discussion started: 17 July 2024

Revised: 6 January 2025 – Accepted: 29 January 2025 – Published: 26 May 2025

**Abstract.** Uncertainty is inherent in gridded meteorological data, but this fact is often overlooked when data products do not provide a quantitative description of prediction uncertainty. This paper describes, applies, and evaluates a method for quantifying prediction uncertainty in spatially interpolated estimates of meteorological variables. The approach presented here, which we will refer to as DNK for “detrend, normal score, krige”, uses established methods from geostatistics to produce not only point estimates (i.e., a single number) but also predictive distributions for each location. Predictive distributions quantitatively describe uncertainty in a manner suitable for propagation into physical models that take meteorological variables as inputs. We apply the method to interpolate daily maximum near-surface air temperature ( $T_{\max}$ ) and then validate the uncertainty quantification by comparing theoretical versus actual coverage of prediction intervals computed at locations where measurement data were held out from the estimation procedure. We find that, for most days, the predictive distributions accurately quantify uncertainty and that theoretical versus actual coverage levels of prediction intervals closely match one another. Even for days with the worst agreement, the predictive distributions meaningfully convey the relative certainty of predictions for different locations in space. After validating the methodology, we demonstrate how the magnitude of prediction uncertainty varies significantly in both space and time. Finally, we examine spatial correlation in predictions and errors using conditional Gaussian simulation to sample from the joint spatial predictive distribution. In summary, this work demonstrates the efficacy and value of describing un-

certainty in gridded meteorological data products using predictive distributions.

## 1 Introduction

Interpolated meteorological data products are widely used in the geosciences, but relatively little attention is paid to the errors they contain. For example, when studying terrestrial fluxes of carbon, water, and energy over a large spatial domain (e.g.,  $\geq 100 \text{ km}^2$ ), it is necessary to work with gridded meteorological data. Ground-based weather stations may be sparse or only cover a small fraction of the study area, so gridded estimates, rather than station measurements, of meteorological variables are used by models of land surface processes (Zeng et al., 2020; Volk et al., 2024). In many gridded data products, the values are point estimates (i.e., a single number rather than a range or distribution). When given only point estimates, data users do not know and cannot propagate the uncertainty in the meteorological inputs to their model. While users may refer to point estimate accuracy statistics for the data product, these statistics only capture errors at locations where measurements are available. For applications that are particularly sensitive to meteorological inputs, such as evapotranspiration modeling, uncertainty in gridded data can contribute significantly to downstream model errors (Doherty et al., 2022). While geostatistical uncertainty quantification is standard practice in other domains like mining (Rossi and Deutsch, 2014), oil and gas exploration (Pyrzcz and Deutsch, 2014), and hydrogeology

(Kitanidis, 1997), these methods are not as widely used in popular near-surface meteorological data products. Understanding uncertainty in gridded meteorological data is necessary to evaluate the robustness of scientific findings, especially when designing and implementing public policy based on those findings (Morgan and Henrion, 1990).

One approach to producing gridded near-surface meteorological data is statistical interpolation, where gridded values are estimated by interpolating between measurements at ground-based weather stations. For North America, Daymet (Thornton et al., 1997; Thornton et al., 2021) and PRISM (Daly et al., 2008), which produce estimates of several meteorological variables on fine spatial grids ( $\sim 1 \text{ km}^2$ ), are widely used statistical interpolation products. A related product, NEX-GDM (Hashimoto et al., 2019), uses machine learning and a wide range of inputs to produce high-resolution gridded meteorological values. For Europe, E-OBS (Cornes et al., 2018) is an interpolated data product with 12 km spatial resolution. Regarding uncertainty, Daly et al. (2008) describe a method for creating prediction intervals, but the resulting maps are not publicly distributed. Thornton et al. (2021) include an extensive accuracy assessment using cross-validation, but the methodology does not produce spatially resolved uncertainty estimates. The E-OBS methodology is the most similar to this work in terms of modeling the data as a Gaussian random field and producing predictive distributions for each grid cell. However, there are important differences in the data processing, the effects of which may explain some of the results in Cornes et al. (2018). We address this topic in the discussion. Spatial machine learning methods including quantile random forest (QRF) can also be used to produce prediction intervals for uncertainty quantification (Hengl et al., 2018; Milà et al., 2023). However, this approach does not take into account spatial correlation and is sometimes combined with other geostatistical methods that do (Milà et al., 2023). Finally, Bayesian methods using a Gaussian Markov random field (GMRF) model have also been used to perform probabilistic interpolation of meteorological data. For example, Fioravanti et al. (2023) applied these methods to air temperature (but did not quantitatively validate the uncertainty quantification) and Ingebrigtsen et al. (2014) applied them to precipitation data. In future work, it would be instructive to compare predictive distributions produced using QRF and Bayesian GMRF-based methods with those produced using the more classical geostatistical methods described in this work.

Another class of approaches to producing gridded data products is data assimilation, which combines dynamic physical models with data-driven adjustments. This work is focused on statistical interpolation rather than data assimilation, but we give a brief overview of the latter for context. A wide range of data assimilation products are available including regional products like RTMA (De Pondeca et al., 2011) and CONUS404 (Rasmussen et al., 2023) and global ones like MERRA-2 (Gelaro et al., 2017) and ERA5 (Hersbach et

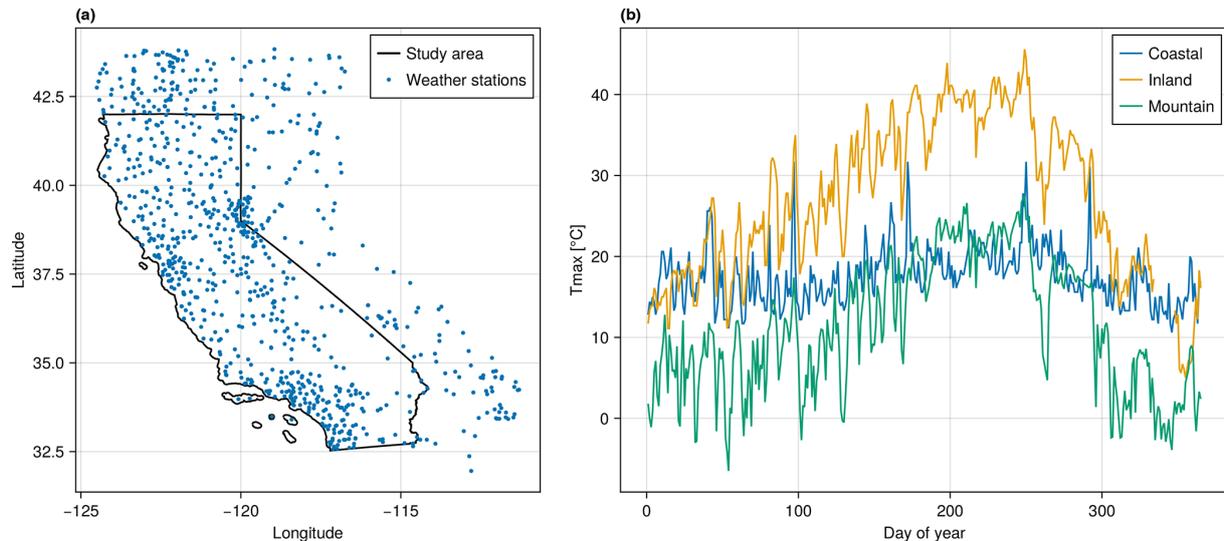
al., 2020; Bell et al., 2021). Some assimilation products, like ERA5, express uncertainty using an ensemble of model runs, where a greater magnitude of spread in the ensemble is taken to indicate greater uncertainty. However, the computational expense of large-scale climate simulations generally means that the resulting data products have relatively coarse spatial resolution (31 km horizontal resolution for ERA5) and that ensembles are not large enough (tens of ensemble members) to characterize stable empirical distributions. In contrast, the approach described in this work is computationally efficient enough to be run at fine spatial resolution over large areas while also giving a robust description of the predictive distribution.

In this paper we present and analyze a statistical method to produce a spatially and temporally resolved uncertainty quantification and apply it to the interpolation of daily maximum near-surface air temperature ( $T_{\text{max}}$ ). We will refer to the approach for estimation and uncertainty quantification as DNK for “detrend, normal score, krigé”. The basic approach of DNK is well-established in geostatistics, appearing in textbooks such as Olea (1999) and Goovaerts (1997). While kriging and related spatial regression methods have previously been used for meteorological data interpolation, they are most commonly used to produce gridded point estimates. A central component of this work is to test the validity of predictive distributions generated using DNK and, as such, their utility for uncertainty quantification. Uncertainty is not intrinsic to macroscale physical phenomenon but rather is a property of the combination of data and a model (Goovaerts, 1997), which means that there is not an objective “correct” predictive distribution for a given unknown value. However, we can assess the validity of a collection of predictive distributions, in aggregate, by testing the rate at which true measurement data fall within prediction intervals relative to those intervals’ theoretical coverage. If the validity of predictive distributions can be established, then DNK can accurately quantify uncertainty in gridded meteorological data.

## 2 Methods

### 2.1 Input data

This study uses two sets of input data: daily maximum air temperature at 2 m ( $T_{\text{max}}$ ) and elevation.  $T_{\text{max}}$  data are provided by Thornton et al. (2022) for stations in the Global Historical Climatology Network (GHCN) (Menne et al., 2012), a database of measurement data from ground-based weather stations across the world. The GHCNd (daily) data are processed as described in Thornton et al. (2021) to correct for temperature sensor biases and inconsistencies in time of observation. Figure 1a shows the spatial distribution of weather stations across the study area. We use data from 2022 to conduct the validation study. In 2022, the number of weather stations within the California state boundary ranges between



**Figure 1.** Study area, locations of weather stations, and sample  $T_{\max}$  time series. Panel (a) shows the study area and weather station locations. Black lines mark the bounds of the study area (the state of California). Blue dots mark the locations of GHCN weather stations that were active in 2022. Panel (b) shows sample time series for three weather stations, each from a different type of location.

524 and 542 stations depending on the day of the year (DOY). Data from stations within the state boundaries are used as ground truth for validation. Data from stations outside the study area contribute to predictions at locations near the boundary, but these stations are not, themselves, used as validation locations. Elevation data are sourced from a digital elevation model (DEM) with 90 m resolution for the western United States (Hanser, 2008). The DEM is clipped to the boundaries of the study area and then resampled using mean resampling to a grid with 1 km<sup>2</sup> grid cells. Figure 1b shows example  $T_{\max}$  time series for three stations from different types of regions: “USC00045795” (blue, coastal), “USW00053119” (orange, inland), and “USS0019L45S” (green, mountain).

## 2.2 Probabilistic interpolation procedure

The interpolation procedure, which we summarize as detrend, normal score, krigé (DNK), combines established methods from geostatistics to produce predictive distributions. We briefly summarize the procedure here, before providing greater detail in the following sub-sections. At a high level, the data are modeled as being observations of a stationary Gaussian random field. The “detrend” and “normal score” steps are transformations applied to make the data better satisfy the assumptions of this mathematical model. The procedure consists of the following steps:

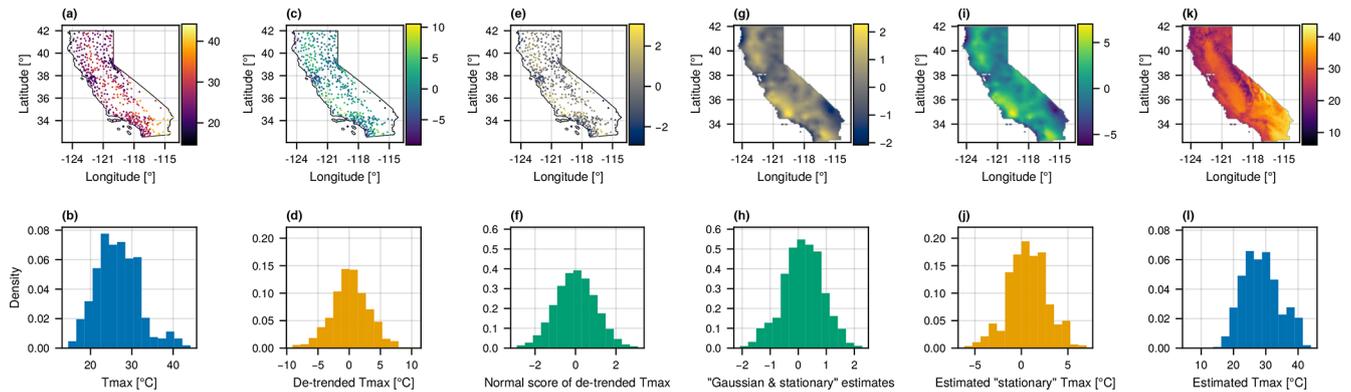
1. detrend – estimate and subtract spatial trends from measurement data
2. normal score – apply quantile-to-quantile mapping transforming observed empirical distribution to a standard normal distribution

3. krigé – produce marginal (or joint) predictive distributions at prediction locations using ordinary kriging (or conditional Gaussian simulation)
4. inverse normal score transformation – apply the inverse of the mapping from step (2) to the output of step (3), transforming predicted values back to the “original” (detrended) distribution
5. add back the trend – estimate spatial trends at prediction locations and add them to the output of step (4) to produce final predictive distributions.

The procedure is also represented graphically in Fig. 2, where each step in the previous list corresponds to a transition from one column to the next (moving from left to right). We will refer to Fig. 2 in the following sub-sections. All calculations for this study were performed using the Julia programming language (Bezanson et al., 2017), and we make significant use of the GeoStats.jl package (Hoffmann, 2018) in particular. All package names and versions that were used can be found in the Manifest.toml file provided (see “Code availability”).

### 2.2.1 Detrending

The first transformation models and then subtracts local spatial trends. The purpose of trend modeling is to identify and remove variation at coarse spatial scales that would otherwise make the data nonstationary. We use the term “trend” to refer to both large-scale variation in longitude–latitude space and variation due to change in elevation (lapse rate). While some prior approaches assume a fixed lapse rate (Hart et al., 2009), we allow the lapse rate to vary in space (Thornton et



**Figure 2.** Illustration of data-processing steps. Each column corresponds to a step in the processing and estimation pipeline. The top row shows point (a, c, e) and gridded (g, i, k) data values. The bottom row shows histograms of the data at each processing step. Panels (a) and (b) show the measured  $T_{\max}$  values. Panels (c) and (d) show the  $T_{\max}$  values after spatial trends have been estimated and subtracted. Panels (e) and (f) show the detrended  $T_{\max}$  data after the normal score transformation has been applied. Panels (g) and (h) show gridded estimates produced by OK in the detrended and Gaussian space. Panels (i) and (j) show gridded estimates after reverting the normal score transformation using quantile information from the distributions in panels (c) and (d). Panels (k) and (l) show the final gridded estimates after reading spatial trends.

al., 1997, 2021; Daly et al., 2008; Cornes et al., 2018). Local trends are estimated by regressing  $T_{\max}$  values on spatial coordinates and elevation:

$$T_{\max_i} = x_i + y_i + z_i + \epsilon_i, \quad (1)$$

where  $x_i$  and  $y_i$  are the spatial coordinates and  $z_i$  is the elevation at weather station  $i$ . The residual  $\epsilon_i$ , calculated as the observation minus the fitted value from the trend model, is taken as the detrended  $T_{\max}$  value. Detrending is performed locally using a 100 km search radius around each weather station. Detrended station-level  $T_{\max}$  values are shown in Fig. 2b and c. A similar procedure is applied to “add back” the trend after estimation: the trend parameters are estimated at locations centered on each grid cell, again using data from weather stations within 100 km. The cell-wise trend value is calculated using the regression parameter estimates and the corresponding coordinates of the cell and elevation from a digital elevation model. The gridded estimates with the trend added back is shown in Fig. 2k and l. The purpose of detrending is not to explain all variation in  $T_{\max}$  values. Rather, the purpose is to control for large-scale variation such that the residuals, after detrending, can be modeled as a random field whose mean does not vary deterministically in space. While it could be useful to incorporate other covariates (e.g., distance to the ocean), doing so would likely require a non-linear trend model.

## 2.2.2 Normal score transformation

The second transformation, a “normal score transformation”, transforms the data to be approximately Gaussian. The transformation is done by mapping quantiles of the empirical distribution of detrended  $T_{\max}$  values to the corresponding quantiles of a standard normal distribution. The transformed

station-level data are shown in Fig. 2c and d. The quantile information from the original empirical distribution is saved so that the inverse transformation can be applied after estimation. The gridded estimates, after reverting the normal score transformation, are shown in Fig. 2i and j. This transformation is implemented in `TableTransforms.jl`. Inclusion of this step is one place where our approach differs from the E-OBS methodology (Cornes et al., 2018) and could explain some of the differences in results. We address this difference in the discussion.

## 2.2.3 Estimation of marginal predictive distributions using ordinary kriging

The primary estimation method we consider is ordinary kriging (OK) (see, e.g., Goovaerts, 1997; Olea, 1999; Anderes, 2012), which gives analytical solutions for both a point estimate and the variance of a predictive distribution. For a random field that is covariance stationary and Gaussian, the OK prediction mean and variance completely characterize the predictive distribution. In general, spatially distributed  $T_{\max}$  data satisfy neither assumption, which is why we first apply detrending and normal score transformations. We apply OK locally at each prediction location using measurement data within a 100 km radius of the location in the estimation. Figure 2g and h show an example of the gridded estimates produced by OK from the point measurement data shown in Fig. 2e and f. Samples (or quantiles) of a given predictive distribution are generated by drawing samples from (or calculating quantiles for) the Gaussian predictive distribution described by the OK prediction mean and variance and then applying the inverse normal score and detrending transformations as described in the prior two sections. The validation scheme for local (marginal) predictive distribu-

tions is described in Sect. 2.3, and the results of the validation study are presented in Sect. 3.1.

#### 2.2.4 Sampling from the joint predictive distribution using conditional Gaussian simulation

In addition to OK, we also demonstrate spatial uncertainty quantification using conditional Gaussian simulation (CGS). Samples generated by CGS are equally probable “realizations” of the underlying random field that produced the measurement data. For a stationary Gaussian random field, local predictive distributions (i.e., the marginal predictive distribution at a given grid cell) generated by CGS are equivalent to the distribution defined by the kriging variance (see, e.g., Goovaerts, 2001). However, CGS produces realizations that are spatially coherent with respect to the model of spatial covariance (see Sect. 2.2.5, Variography), whereas gridded estimates produced by OK do not. As such, CGS can be used to express “spatial uncertainty”, or spatial correlation in errors, as described by the joint predictive distribution over multiple grid cells.

We apply CGS using a method based on a decomposition of the conditional distribution covariance matrix, commonly referred to in geostatistics literature as the “LU method” of CGS (Alabert, 1987). The LU method samples from the exact multivariate Gaussian predictive distribution. For larger-scale simulations, the LU method becomes computationally infeasible and other algorithms and approximations may be required. However, in this example, we are drawing samples for a small number of locations, so the LU method works well. As with OK, inverse normal score and detrending transformations are applied to produce the final predictive distributions.

This study does not include a comprehensive assessment of spatial uncertainty quantification. Instead, we present a case study (Sect. 3.2) using data from two weather stations close to one another where we expect the  $T_{\max}$  values and prediction errors to be correlated. Quantitative results using  $10^4$  samples from the joint and marginal predictive distributions are shown that illustrate how CGS describes spatial uncertainty.

#### 2.2.5 Variography

Both OK and CGS rely on a model of spatial covariance to produce gridded  $T_{\max}$  estimates and prediction uncertainty. Traditionally in geostatistics, spatial variation is represented using a semi-variogram, which is a function that describes the decrease in correlation between two locations as the distance between them increases (see, e.g., Olea, 1999; Goovaerts, 1997). In this study, variography is performed using functions from the GeoStats.jl library (Hoffmann, 2018). We model the theoretical semi-variogram using the pentaspherical function:

$$\gamma(h) = (s - n) \left[ \left( \frac{15}{8} \left( \frac{h}{r} \right) - \frac{5}{4} \left( \frac{h}{r} \right)^3 + \frac{3}{8} \left( \frac{h}{r} \right)^5 \right) \cdot 1_{(0,r)}(h) + 1_{[r,\infty)}(h) \right] + n \cdot 1_{(0,\infty)}(h), \quad (2)$$

where the variable  $h$  is the lag distance;  $s$ ,  $n$ , and  $r$  are parameters estimated to fit the empirical semi-variogram representing the sill (value of  $\gamma$  as  $h \rightarrow \infty$ ), nugget (value of  $\gamma$  as  $h \rightarrow 0$ ), and range (roughly the value of  $h$  where  $\gamma$  “levels off”), respectively; and  $1_{(l,u)}(h)$  is an indicator function that is 1 if  $l < h < u$  and 0 otherwise. The theoretical model is fit to the empirical semi-variogram by minimizing the sum of squared errors with equal weight given to each lag bin. Model fitting is performed using the detrended and normal score data (the data in Fig. 2e and f) for each day independently (i.e., there is a different semi-variogram model for each day). At shorter lag distances, the pentaspherical model produces semi-variances between that of exponential and spherical models, which have been used previously to interpolate near-surface air temperature (Cornes et al., 2018; Hudson and Wackernagel, 1994). Running the analysis using these different semi-variogram models produces, in aggregate, very similar validation statistics. The pentaspherical model produces better accuracy in the CGS case study, and we use it throughout the analysis for consistency. Interested readers can reproduce all results and figures in this paper using different semi-variogram models by replacing the model type passed to the semi-variogram fitting function (see “Code availability”). Examples of fitted semi-variograms are shown in Fig. S1 in the Supplement.

At present, we do not model or try to exploit temporal correlation in  $T_{\max}$  or prediction errors. This topic will be addressed in future work.

### 2.3 Validation of local predictive distributions

We implement a validation scheme that evaluates the accuracy of the OK predictive distributions in quantifying local prediction uncertainty. Uncertainty is not intrinsic to the physical phenomenon, and, as such, there is not an objective correct predictive distribution. However, a collection of valid predictive distributions should produce statistics that reflect appropriate levels of confidence in aggregate. The main quantitative validation results in this study are, for local (marginal) predictive distributions, computed using OK. We present some quantitative results for a CGS case study, but we do not conduct a comprehensive evaluation of spatial uncertainty quantification.

To assess the validity of local predictive distributions, we use a strategy based on prediction intervals described by Deutsch (1997). For each day of the year, we perform leave-one-out (LOO) cross-validation with each weather station. The measurement at the left-out station is not used in trend modeling, variography, or estimation. Predictions are made for the point at the center of the grid cell containing the

weather station that will be used for validation. Estimates of the multiples of 0.005 quantiles are produced for each predictive distribution, from which centered prediction intervals are calculated with  $q$  coverage for  $q \in Q = \{0.01, 0.02, \dots, 0.99\}$ . For each  $q$  prediction interval, let  $q_{\text{low}} = (1 - q)/2$  and  $q_{\text{upp}} = (1 + q)/2$  be the lower and upper quantiles bounding the theoretical prediction interval. Then for  $T_{\text{max}_i}$ , the measured  $T_{\text{max}}$  at weather station  $i$ , define an indicator function:

$$\xi(T_{\text{max}_i}; q) = \begin{cases} 1, & \text{if } q^*(T_{\text{max}_i}) \in (q_{\text{low}}, q_{\text{upp}}) \\ 0, & \text{otherwise} \end{cases}, \quad (3)$$

where  $q^*(T_{\text{max}_i})$  denotes the quantile of the true  $T_{\text{max}_i}$  relative to the predictive distribution for location  $i$ . Then for each  $q$  prediction interval, we compute an average over all

$n$  stations as  $\bar{\xi}(q) = \frac{1}{n} \sum_{i=1}^n \xi(T_{\text{max}_i}; q)$ . For exactly valid predictive distributions,  $\bar{\xi}(q) = q$  for any  $q$ . To summarize the errors, we calculate the mean bias error as  $\frac{1}{99} \sum_{q \in Q} \bar{\xi}(q) - q$

and the mean absolute error (MAE) as  $\frac{1}{99} \sum_{q \in Q} |\bar{\xi}(q) - q|$ . The MAE weights errors being too confident and too conservative equally. The bias indicates whether the predictive distributions are too confident (negative) or too conservative (positive) on average.

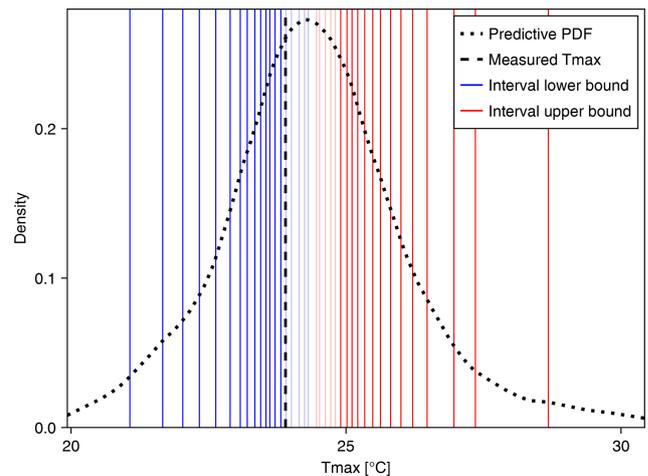
Figure 3 shows a representation of how the  $\xi$  function is evaluated. The vertical blue and red lines denote the lower and upper bounds, respectively, of nested prediction intervals. The translucent lines indicate intervals that do not contain the true measured value, which is denoted by the vertical dashed line. The solid lines indicate that the interval does contain the true value. The prediction intervals are determined by the predictive distribution, the probability density function (PDF) of which is shown as a dotted curve.

Figure 3 shows a representation of how the  $\xi$  function is evaluated. The vertical blue and red lines denote the lower and upper bounds, respectively, of nested prediction intervals. The translucent lines indicate intervals that do not contain the true measured value, which is denoted by the vertical dashed line. The solid lines indicate that the interval does contain the true value. The prediction intervals are determined by the predictive distribution, the probability density function (PDF) of which is shown as a dotted curve.

### 3 Results

#### 3.1 Local uncertainty quantification using ordinary kriging

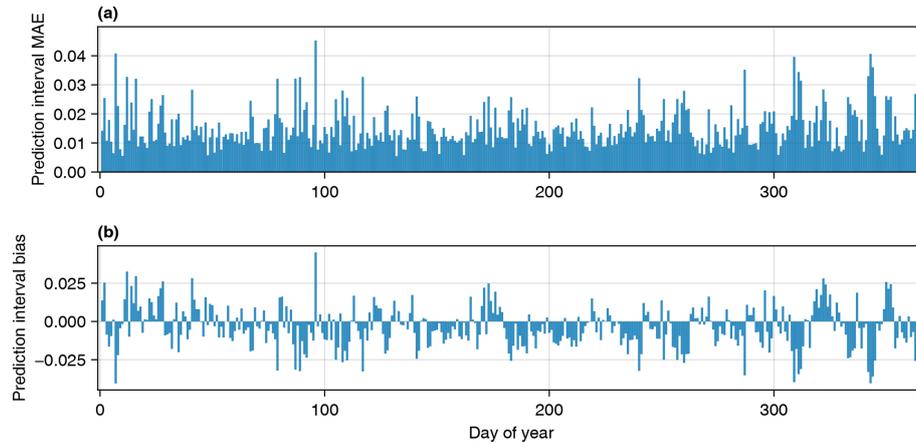
We first present the validation of local uncertainty quantification using the predictive distributions from OK. Figure 4 shows the average accuracy of the predictive distributions using the LOO validation scheme described in Sect. 2.5. Figure 4a shows the mean absolute error (MAE) in terms of the predicted versus actual proportion of  $T_{\text{max}}$  values that fall within a given prediction interval. For example, an MAE of 0.01 indicates that for a  $p\%$  prediction interval (where  $p = q \times 100$ ), the true value fell within that interval  $(p \pm 1)\%$  of the time. The largest MAE of 0.045 occurs on DOY 96. The median MAE is 0.013, and 82% of days had an MAE less than 0.02. Figure 4b shows the average bias for each day. For example, a bias of  $-0.01$  indicates that for a  $p\%$  prediction interval, the true value fell within that prediction interval



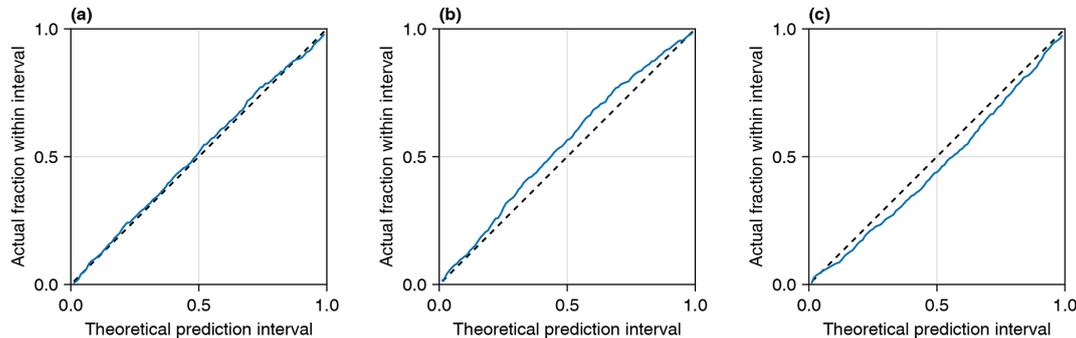
**Figure 3.** Prediction intervals versus the measured  $T_{\text{max}}$ . The lower bounds (blue) and upper bounds (red) of prediction intervals with coverage  $q \in \{0.05, 0.10, \dots, 0.95\}$  for a single prediction location are shown. (Validation statistics are computed using  $q \in \{0.01, 0.02, \dots, 0.99\}$ , but we show fewer intervals here for legibility.) Intervals that contain the measured  $T_{\text{max}}$  (dashed black line) are drawn as solid lines, and the intervals that do not are drawn as partially transparent. The dotted black curve shows a kernel density estimate of the predictive distribution.

$(p - 1)\%$  of the time. Positive bias indicates that the prediction intervals are too conservative on average, and negative bias indicates that the prediction intervals are too confident on average. The largest positive bias of  $+0.045$  occurred on DOY 96, and the largest negative bias of  $-0.040$  occurred on DOY 7. The median bias was  $-0.005$ . There are temporal patterns with the bias metric errors of a similar sign and magnitude persisting for periods ranging from a few days to multiple weeks.

Figure 5 shows validation results for 3 individual days. In addition to showing the day with the median MAE, we highlight the two days with the most significant errors in the sample to give a sense of the “worst-case” accuracy in uncertainty quantification. Values in the  $x$  direction correspond to theoretical prediction intervals centered on the median. Values in the  $y$  direction are the actual proportion of true  $T_{\text{max}}$  values that fall within the given theoretical interval. For example, a point at 0.25, 0.3 would mean that 30% of true values fell within the corresponding 25% prediction intervals. Figure 5a shows results for DOY 304, which had the median MAE of 0.013. For DOY 304, the actual proportions closely track the theoretical prediction intervals, with the largest error occurring at the 74% intervals, which contained 77.1% of the true values. Figure 5b and c show the same information for DOY 96 and DOY 7, which are the days with the largest positive and negative biases of  $+0.045$  and  $-0.040$ , respectively. The largest error for DOY 96 occurs at the 60% prediction intervals, which contained 68.1% of true  $T_{\text{max}}$  values.



**Figure 4.** Predictive-distribution quantitative validation statistics. Validity is assessed by comparing the theoretical coverage versus the rate that measurement data actually fell within the prediction intervals. The calculation uses an indicator function  $\xi$  described in Eq. (1). Panel (a) shows the mean absolute error (MAE) for each day of the year, calculated as  $\frac{1}{99} \sum_{p=1}^{99} |\overline{\xi(q)} - q|$ . Panel (b) shows the mean bias error for each day of the year:  $\frac{1}{99} \sum_{p=1}^{99} \overline{\xi(q)} - q$ . Positive bias indicates that predictive distributions were too conservative on average, and negative bias indicates they were too confident on average.



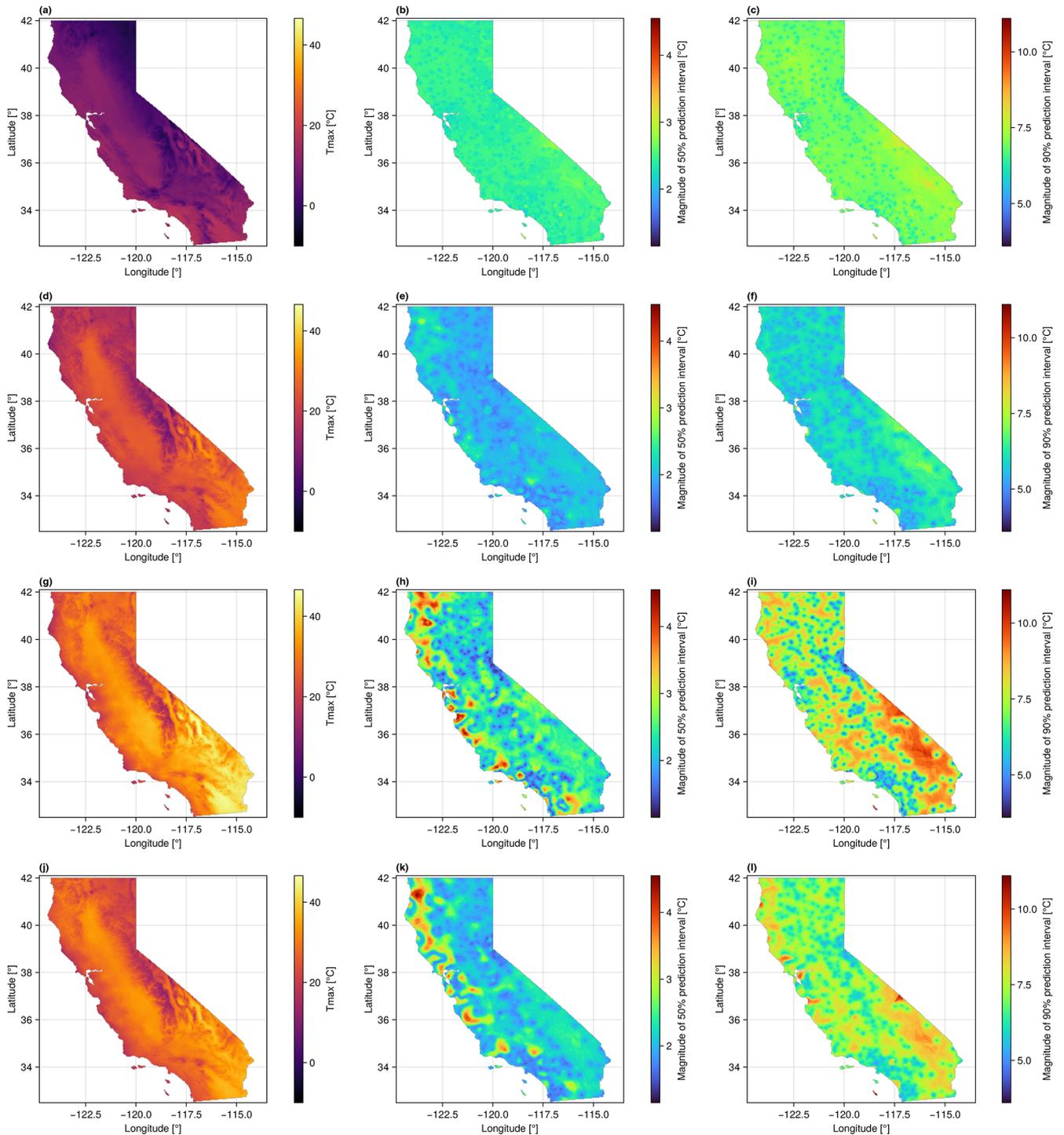
**Figure 5.** Theoretical versus actual coverage of prediction intervals.  $x$  axes correspond to the coverage of the theoretical prediction interval, and  $y$  axes correspond to the proportion of values that actually fell within the intervals. Panel (a) shows results for DOY 304, which had the median MAE of all the days in the sample. Panels (b) and (c) show DOY 96 and 7, which had the largest positive and negative biases, respectively.

The largest error for DOY 7 occurs at the 62 % prediction intervals, which contained 54.9 % of true  $T_{\max}$  values.

When applying our method across a full spatial domain, we observe that local predictive distributions vary in space and time. Figure 6 shows predicted  $T_{\max}$  and uncertainty statistics (spread of predictive distributions) for 4 different days, where each row corresponds to a single day of the year. The days are approximately equally spaced and cover different seasons, including the first days of January (Fig. 6a–c), April (Fig. 6d–f), July (Fig. 6g–i), and October (Fig. 6j–l). The first column (Fig. 6a, d, g, and j) shows the median of the predictive distribution for each 1 km<sup>2</sup> grid cell. The second column (Fig. 6b, e, h, and k) shows the magnitude of the 50 % prediction interval of each predictive distribution (the 75th percentile minus the 25th percentile). The third column

(Fig. 6c, f, i, and l) shows the magnitude of the 90 % prediction interval of each predictive distribution (the 95th percentile minus the 5th percentile). Each column uses a single common color gradient (the gradients do not vary between rows).

The  $T_{\max}$  prediction uncertainty varies in both space and time. Comparing the uncertainty maps in time (between rows), the magnitude, spatial patterns, and the magnitude of variation in prediction uncertainty all vary. Spatial patterns also differ based on the magnitude of the prediction interval. OK prediction variance depends only on the distances to measurement data locations, not the measurement values themselves (see, e.g., Olea, 1999). However, the forward and inverse normal score transformations are nonlinear, so the shape of the final prediction distribution is influenced by the



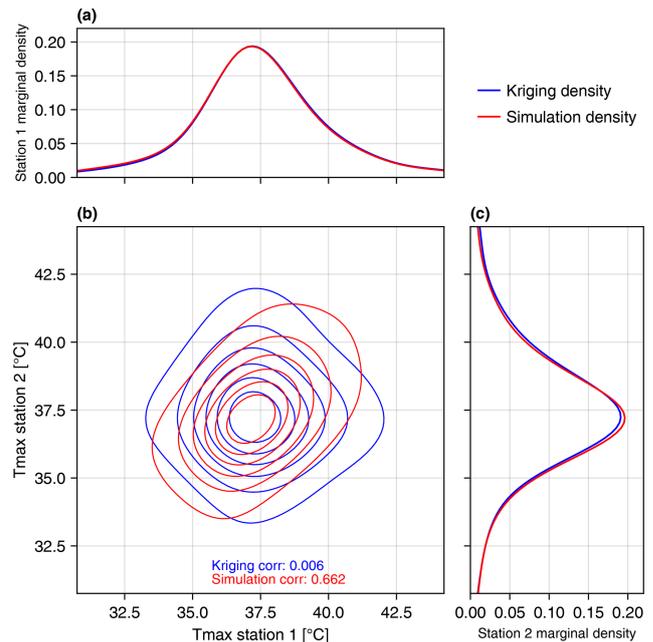
**Figure 6.**  $T_{\max}$  and prediction uncertainty maps for 4 different days. Each row shows results for a different day. The days are DOY 1 (a–c), DOY 91 (d–f), DOY 182 (g–i), and DOY 274 (j–l) from 2022. The first column (a, d, g, j) shows the median of the  $T_{\max}$  predictive distribution at each grid cell. The second column (b, e, h, k) shows the magnitude of the 50% prediction interval for each grid cell. The third column (c, f, i, l) shows the magnitude of the 90% prediction interval for each grid cell.

actual measurement and prediction values. The 90 % prediction intervals (rightmost column) appear to be controlled primarily by the local spatial density of measurements, with the locations of weather stations standing out as local minima. The patterns in the 50 % prediction intervals (center column) are more complex. The fact that many of the locations with the largest uncertainties are near the Pacific coast may be caused by the trend model failing to capture local  $T_{\max}$  trends, where cooling from the ocean confounds the usual negative correlation between temperature and elevation. This could cause the detrended measurements near the coast to be outliers on the low end of the overall distribution. As a result, prediction intervals in this part of the distribution get “stretched out” by the inverse normal score transformation. This issue would not occur in the mountains where the expected relationship between temperature and elevation holds, and the resulting detrended data are not outliers. A more sophisticated approach to trend modeling may improve accuracy and reduce uncertainty near the coast.

### 3.2 Case study: spatial uncertainty quantification using conditional Gaussian simulation

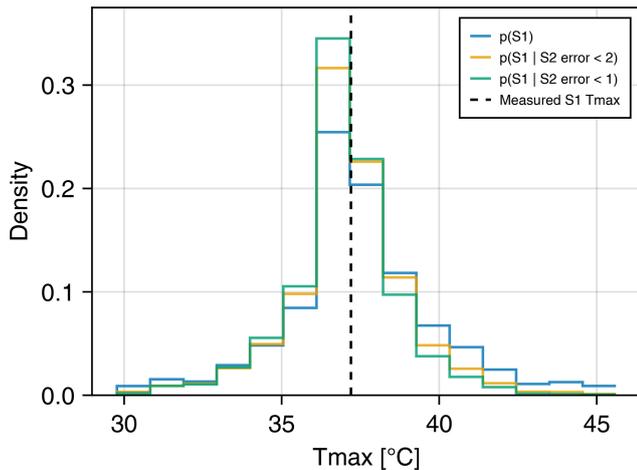
In addition to local (cell-wise marginal) uncertainty, we can represent spatial uncertainty using CGS to sample from the joint predictive distribution over multiple grid cells. Figure 7 shows the predictive distributions for grid cells containing two weather stations, USC00043747 (“Station 1”) and USW00053119 (“Station 2”), near Hanford, California, for DOY 182. The centers of the grid cells containing the stations are 1.4 km apart. The red lines show kernel density estimates of the joint and marginal predictive distributions generated using CGS. The blue lines show the same predictive distributions generated using OK, where the “joint distribution” is generated by sampling independently from the marginal distributions. The empirical distributions were generated using  $10^4$  samples each. While the OK samples show no correlation between locations (by construction), the samples produced by CGS show a correlation of approximately 0.66 between predictions at the two locations. This shows how errors in predictions at nearby locations are likely to be correlated with one another. The marginal distributions generated using OK and CGS are virtually identical at each of the two locations (Fig. 7a and c), as expected.

Figure 8 shows histograms of predictive distributions for  $T_{\max}$  estimates at Station 1 when conditioned on different information about  $T_{\max}$  at Station 2. Measurements from both stations are held out from the estimation procedure. The dashed black line is the  $T_{\max}$  measured at Station 1. The blue histogram is the unconditional marginal distribution for the grid cell containing Station 1. The orange and green histograms show empirical conditional distributions at Station 1 given that the error at Station 2 is less than  $2^\circ$  (“conditional 2”) and less than  $1^\circ$  (“conditional 1”), respectively. These empirical distributions are generated by filtering the



**Figure 7.** Joint and marginal predictive distributions using OK versus CGS. Predictive distributions were generated for two nearby grid cells that contain weather stations. Panels (a) and (c) show the marginal (local) predictive distributions at the two locations, and panel (b) shows the joint (spatial) predictive distribution. The predictive distributions from OK are in blue, and the distributions generated by CGS are in red. The two estimation methods produce the same marginal distributions, but the joint distributions are different. The distribution from CGS reflects the spatial correlation in predictions between the two locations.

simulation ensemble. Conditioning the predictive distribution reduces the standard deviation from  $2.4^\circ\text{C}$  in the unconditional distribution to  $1.9$  and  $1.7^\circ\text{C}$  for conditional 2 and conditional 1, respectively. If the mean of the predictive distribution is taken as a point estimate, the error for the unconditional distribution is  $+0.25^\circ\text{C}$ . The error is reduced to  $-0.08$  for conditional 2 but then increases (in magnitude) to  $-0.25$  for conditional 1. The degree to which the point estimates and prediction variance change is controlled by the semi-variogram model. The semi-variogram for DOY 182 has the smallest nugget (i.e., the strongest correlation at short distances) among the days shown in Fig. 6. This means the model-implied correlation between the two stations will be lower on the other 3 d. Plots of empirical and fitted theoretical semi-variograms for the 4 d can be found in the Supplement. In general, we observe that the fitted semi-variograms have relatively large nugget values. This may be partly an artifact of the semi-variogram estimation process, given that there are few pairs of weather stations within 1–10 km of one another. However, there does appear to be significant “real” variance in measured  $T_{\max}$  values even at short spatial scales. For example, the two stations used in this example are 1.4 km apart and at nearly identical elevation but have an average



**Figure 8.** Marginal predictive distributions without and with additional spatial information. Three empirical predictive distributions at Station 1 (S1) are shown including the unconditional distribution (blue) and two distributions conditioned on predictive accuracy at nearby Station 2 (S2). Conditioning is performed by filtering the ensemble and leaving only realizations, where the prediction at S2 was within 2 °C (orange) or within 1 °C (green) of the measured value. Conditioning on additional spatial information reduces the spread of the marginal predictive distribution.

difference in  $T_{\max}$  of nearly 1.5 °C in 2022. This difference may be explained by site-specific effects, such as land cover or other features near the weather stations, rather than physical variation that would persist under idealized conditions.

#### 4 Discussion and conclusion

This work presents and validates an approach to quantifying spatially and temporally resolved prediction uncertainty in interpolated meteorological data products. In the quantitative validation study, the DNK method produced highly accurate uncertainty quantification. Even in the least accurate cases, the predictive distributions are qualitatively informative and are still sufficiently accurate to be useful for many applications. Accuracy assessment of point estimates is useful but fundamentally cannot describe prediction uncertainty at times and locations where measurement data are unavailable. In this application, we have measurements at hundreds of weather stations and predict  $T_{\max}$  values at hundreds of thousands of grid cells, some of which are nearly 100 km from the nearest measurement. This means that the number of locations where we can assess prediction uncertainty using actual measurement data is vanishingly small compared to the number of locations where we do not have measurements. We show that the magnitude of prediction uncertainty varies significantly in space and in time and that using average error statistics will overestimate prediction uncertainty in some cases and dramatically underestimate it in others.

The DNK methodology described in this paper could be applied to a wide range of applications due to its generality and relatively low computational cost. The main application area motivating this work was modeling land surface fluxes of water and carbon, given that rates of evapotranspiration (Volk et al., 2024) and primary production (Zeng et al., 2020) are particularly sensitive to near-surface meteorological conditions. However, gridded meteorological data are used to make predictions and draw conclusions about many other phenomena including crop yield (Lobell et al., 2015), vegetation phenology (White et al., 1997), economic productivity (Burke et al., 2015), human conflict (Hsiang et al., 2013), and others. Using gridded predictive distributions rather than point estimates can help ensure the robustness of scientific conclusions given uncertainty in model inputs. The decision to use cell-wise marginal (OK) or full joint (CGS) predictive distributions depends on various factors including the size of grid cells, the distances between locations being compared, and the sensitivity of the analysis to spatial correlation in prediction errors (e.g., for causal inference). Users of gridded meteorological data products can propagate uncertainty through their analysis by running models multiple times using either random samples or a preset collection of quantiles from the distribution. This approach does not require any additional modeling choices or assumptions because the relevant information about uncertainty in the model inputs is expressed by the predictive distributions. Also, computationally expensive models may use the gridded variable uncertainty to prioritize sensitivity analyses and reduce the total number of model runs required.

Producing a predictive distribution using DNK, rather than a single point estimate, requires only marginally more computation. The main computation required in OK is solving for the kriging weights  $\lambda_s$ , which requires solving a  $(S + 1) \times (S + 1)$  system of linear questions for  $S$  stations. The prediction mean  $\mu_{\text{OK}}$  and variance  $\sigma_{\text{OK}}^2$  are both functions of the  $S$  station data and the weights  $\lambda_s$ . Drawing samples from the predictive distribution only requires drawing samples from  $\text{Normal}(\mu_{\text{OK}}, \sigma_{\text{OK}}^2)$  and then, for each sample, applying the reverse normal score transformation and adding back the trend. Computing the joint conditional distribution and sampling using CGS is more computationally expensive, as it requires solving the system of kriging equations as well as factorizing a  $C \times C$  matrix, where  $C$  is the number of prediction grid cells and often  $C \gg S$ . There exist more computationally scalable methods for conditional Gaussian simulation (Gómez-Hernández and Journel, 1993; Gutjahr et al., 1997; Gómez-Hernández and Srivastava, 2021) that we do not discuss here.

The normal score transformation used in this study is an important step that may explain differences in our validation results compared to prior studies. Cornes et al. (2018) evaluated the predictive distributions of the E-OBS data product using an approach similar to the one used in this work. They found that, on average, the predictive distributions were

overconfident and that measured values fell in the tails of the distributions at a higher-than-expected rate, although this trend varied regionally. Our study is conducted over a much smaller area and with different prediction spatial scales (1 km versus 12 km), so we cannot definitively explain the differences in findings. However, omission of the normal score transformation could explain underdispersion in predictive distributions. We observe that the distributions of the detrended observation data generally have positive excess kurtosis, meaning they have heavier tails than a Gaussian distribution. As such, approximating predictive distributions with a Gaussian could fail to account for additional probability mass in the tails. Surrounding the calculation of the predictive distributions with forward and inverse normal score transformations does a better job of capturing the shape of the original data distribution.

Our analysis also provides perspective on when the additional computational cost of CGS is justified. Cornes et al. (2018) note the high cost of CGS, including the fact that it was prohibitive for earlier versions of the data product. In general, CGS is required when spatial correlation in errors matters for the application, like prediction uncertainty for spatial averages. However, when examining a single location in isolation, the predictive distribution given by OK is valid and CGS does not add value. The semi-variogram model, with the nugget in particular, determines the degree of correlation at shorter spatial scales. For distances beyond some correlation length, the joint distribution can be reasonably approximated as the product of the marginals. While air temperature is traditionally modeled as varying smoothly in space with semi-variogram nuggets close (Cornes et al., 2018) or equal to 0 (Hudson and Wackernagel, 1994), observational  $T_{\max}$  data do not necessarily support this assumption (see Sect. 3.2). This topic merits further inquiry.

Users of the DNK methodology should be aware of some practical limitations. First, we reiterate that there is no objectively correct  $T_{\max}$  predictive distribution for a given location. Uncertainty is a property of the measurement data and modeling decisions, and making different modeling decisions will produce different predictive distributions. Different modeling decisions could lead to larger or smaller errors in point estimates on average but still produce predictive distributions that are valid (i.e., where true values fall in prediction intervals at the prescribed rate). In addition, even large samples from predictive distributions will necessarily suffer from deficiencies inherent in the estimation and sampling process. Covariance stationarity and multi-Gaussianity are strong assumptions that are relied upon for the validity of the predictive distributions, and the transformations made to satisfy these assumptions are imperfect. The normal score transformation requires estimating an empirical distribution from the weather station data. Given that we generally wish to draw on samples larger than the size of the measurement data, reverting the normal score transformation necessarily requires interpolation and extrapolation of that distribution

(see Goovaerts, 1997, for a detailed discussion). In practice, this can produce artifacts like clusters of similar values, particularly near undersampled portions of the distribution. Regarding stationarity, trend modeling can also strongly influence predictive distributions. Like uncertainty itself, a “spatial trend” is not an objectively observable phenomenon. Reliably estimating spatial trends can be difficult when measurement data are sparse or when other physical phenomena, like cooling or warming due to coastal proximity, confound the basic estimation procedure. Limitations in trend modeling contributed to our use of OK (locally constant unknown mean) rather than simple kriging (SK) (locally constant known mean), despite the latter being theoretically justifiable for detrended data. Results using OK versus SK were very similar, with OK performing marginally better likely due to coarse-scale variation that was still present after modeling and subtracting spatial trends. Finally, note that the validation results for this study were produced using high-quality  $T_{\max}$  data from GHCN. If the data contain non-random errors, the assumptions of the DNK procedure will no longer be satisfied and the quality of uncertainty quantification may decline. For less densely sampled study areas, prediction uncertainty will grow for many locations but the predictive distributions should still be valid. It would be valuable to evaluate the DNK methodology using poor-quality and sparse data, but doing so would require a different experimental design and is beyond the scope of this study.

There are many potential avenues for future work building on the methods and results described in this paper. One important area for further study is the analyzing of the effects of trend estimation on characteristics and robustness of predictive distributions. For a covariance stationary and multi-Gaussian random field, predictive distributions will be valid over a sufficiently large sample. This indicates that invalid predictive distributions are driven primarily by the transformations we apply (and revert) to make the data stationary and Gaussian. Relatedly, it would be useful to find ways of incorporating additional physical information not explained by a large-scale spatial trend. The strength of data products like NEX-GDM and PRISM come from the use of additional physical information (e.g., coastal proximity, slope, aspect) in predictions. It is not immediately clear how this information could be incorporated into the underlying mathematical model from which our predictive distributions are derived, but doing so could produce more precise estimates accompanied by valid predictive distributions. Lastly, the DNK method should be tested for the interpolation of meteorological variables other than  $T_{\max}$ . The methodology seems likely to transfer to certain variables like daily minimum air temperature and humidity, but it is not a given that all variables can be approached the same way.

Going forward, it will be valuable not only to produce the “best available” gridded meteorological data products but also to produce spatially and temporally resolved uncertainty quantification. Accounting for uncertainty in model

inputs is important considering the robustness of scientific conclusions. It is also important for guiding the design and implementation of science-informed policies. Given uncertain information, conclusions about the “best” policy option may differ when using a deterministic versus probabilistic cost–benefit analysis (Morgan and Henrion, 1990). Similarly, there may be asymmetric consequences for over- or underestimation of a given model input. In this scenario, using a predictive distribution rather than a point estimate allows policy-makers to quantitatively assess tradeoffs between maximizing expected outcomes and minimizing risk. More broadly, accounting for uncertainty in scientific models is necessary not only for designing informed policies but also for building and maintaining trust in science-informed policymaking.

*Code availability.* All code to perform analysis and generate figures is available on Zenodo at <https://doi.org/10.5281/zenodo.12171025> (Doherty, 2024) and on GitHub at <https://github.com/conordoherty/met-uncertainty-paper>.

*Data availability.*  $T_{\max}$  data from Thornton et al. (2022) are available for download from the ORNL DAAC (Oak Ridge National Laboratory Distributed Active Archive Center) at <https://doi.org/10.3334/ORNLDAAC/2132> (Thornton et al., 2022). The DEM used in this study is available on Zenodo (<https://doi.org/10.5281/zenodo.14602669>, Doherty, 2024). It was created by resampling the dataset from Hanser (2008), which can be found at <https://www.sciencebase.gov/catalog/item/542aebf9e4b057766eed286a>.

*Supplement.* The supplement related to this article is available online at <https://doi.org/10.5194/gmd-18-3003-2025-supplement>.

*Author contributions.* CD and IB developed the concept and acquired funding. CD developed the methodology, developed the software, performed the experiments and formal analysis, and wrote the initial draft of the manuscript. WW contributed to methodology development and formal analysis. HH contributed to methodology development, validation, and software testing. All authors contributed to reviewing and editing the manuscript.

*Competing interests.* The contact author has declared that none of the authors has any competing interests.

*Disclaimer.* Publisher’s note: Copernicus Publications remains neutral with regard to jurisdictional claims made in the text, published maps, institutional affiliations, or any other geographical representation in this paper. While Copernicus Publications makes every effort to include appropriate place names, the final responsibility lies with the authors.

*Acknowledgements.* Conor T. Doherty is supported by an appointment to the NASA Postdoctoral Program at NASA Ames Research Center, administered by Oak Ridge Associated Universities under contract with NASA. Conor T. Doherty thanks Forrest Melton for support on the initial conceptualization and proposal for this work, A. J. Purdy and Lee Johnson for helpful feedback on the manuscript, and Kyle Kabasares for performing additional software testing for reproducibility. Conor T. Doherty also thanks Júlio Hoffmann for his generous assistance in using the GeoStats.jl library, which he created.

*Financial support.* This research was supported by an appointment to the NASA Postdoctoral Program at Ames Research Center, administered by Oak Ridge Associated Universities under contract with NASA.

*Review statement.* This paper was edited by Rohitash Chandra and reviewed by Carles Milà and two anonymous referees.

## References

- Alabert, F.: The practice of fast conditional simulations through the LU decomposition of the covariance matrix, *Math. Geol.*, 19, 369–386, <https://doi.org/10.1007/BF00897191>, 1987.
- Anderes, E. B.: Kriging, in: *Encyclopedia of Environmetrics*, edited by: El-Shaarawi, A. H. and Piegorsch, W. W., Wiley, <https://doi.org/10.1002/9780470057339.vak003.pub2>, 2012.
- Bell, B., Hersbach, H., Simmons, A., Berrisford, P., Dahlgren, P., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Radu, R., Schepers, D., Soci, C., Villaume, S., Bidlot, J., Haimberger, L., Woollen, J., Buontempo, C., and Thépaut, J.: The ERA5 global reanalysis: Preliminary extension to 1950, *Q. J. Roy. Meteorol. Soc.*, 147, 4186–4227, <https://doi.org/10.1002/qj.4174>, 2021.
- Bezanson, J., Edelman, A., Karpinski, S., and Shah, V. B.: Julia: A fresh approach to numerical computing, *SIAM Rev.*, 59, 65–98, 2017.
- Burke, M., Hsiang, S. M., and Miguel, E.: Global non-linear effect of temperature on economic production, *Nature*, 527, 235–239, <https://doi.org/10.1038/nature15725>, 2015.
- Cornes, R. C., Van Der Schrier, G., Van Den Besselaar, E. J. M., and Jones, P. D.: An Ensemble Version of the E-OBS Temperature and Precipitation Data Sets, *J. Geophys. Res.-Atmos.*, 123, 9391–9409, <https://doi.org/10.1029/2017JD028200>, 2018.
- Daly, C., Halbleib, M., Smith, J. L., Gibson, W. P., Doggett, M. K., Taylor, G. H., Curtis, J., and Pasteris, P. P.: Physiographically sensitive mapping of climatological temperature and precipitation across the conterminous United States, *Int. J. Climatol.*, 28, 2031–2064, <https://doi.org/10.1002/joc.1688>, 2008.
- De Pondeca, M. S. F. V., Manikin, G. S., DiMego, G., Benjamin, S. G., Parrish, D. F., Purser, R. J., Wu, W.-S., Horel, J. D., Myrick, D. T., Lin, Y., Aune, R. M., Keyser, D., Colman, B., Mann, G., and Vavra, J.: The Real-Time Mesoscale Analysis at NOAA’s National Centers for Environmental Prediction: Current Status and Development, *Weather Forecast.*, 26, 593–612, <https://doi.org/10.1175/WAF-D-10-05037.1>, 2011.

- Deutsch, C.: Direct assessment of local accuracy and precision, *Geostatistics Wollongong*, 96, 115–125, 1997.
- Doherty, C.: Code for revision of “A Method for Quantifying Uncertainty in Spatially Interpolated Meteorological Data with Application to Daily Maximum Air Temperature”, Zenodo [code], <https://doi.org/10.5281/zenodo.14602669>, 2024 (code also available at: <https://github.com/conordoherty/met-uncertainty-paper>, last access: 5 May 2025).
- Doherty, C. T., Johnson, L. F., Volk, J., Mauter, M. S., Bambach, N., McElrone, A. J., Alfieri, J. G., Hipps, L. E., Prueger, J. H., Castro, S. J., Alsina, M. M., Kustas, W. P., and Melton, F. S.: Effects of meteorological and land surface modeling uncertainty on errors in winegrape ET calculated with SIMS, *Irrig. Sci.*, 40, 515–530, <https://doi.org/10.1007/s00271-022-00808-9>, 2022.
- Fioravanti, G., Martino, S., Cameletti, M., and Toreti, A.: Interpolating climate variables by using INLA and the SPDE approach, *Int. J. Climatol.*, 43, 6866–6886, <https://doi.org/10.1002/joc.8240>, 2023.
- Gelaro, R., McCarty, W., Suárez, M. J., Todling, R., Molod, A., Takacs, L., Randles, C. A., Darmenov, A., Bosilovich, M. G., Reichle, R., Wargan, K., Coy, L., Cullather, R., Draper, C., Akella, S., Buchard, V., Conaty, A., Da Silva, A. M., Gu, W., Kim, G.-K., Koster, R., Lucchesi, R., Merkova, D., Nielsen, J. E., Parityka, G., Pawson, S., Putman, W., Rienecker, M., Schubert, S. D., Sienkiewicz, M., and Zhao, B.: The Modern-Era Retrospective Analysis for Research and Applications, Version 2 (MERRA-2), *J. Climate*, 30, 5419–5454, <https://doi.org/10.1175/JCLI-D-16-0758.1>, 2017.
- Gómez-Hernández, J. J. and Journel, A. G.: Joint Sequential Simulation of MultiGaussian Fields, in: *Geostatistics Tróia '92*, vol. 5, edited by: Soares, A., Springer Netherlands, Dordrecht, 85–94, [https://doi.org/10.1007/978-94-011-1739-5\\_8](https://doi.org/10.1007/978-94-011-1739-5_8), 1993.
- Gómez-Hernández, J. J. and Srivastava, R. M.: One Step at a Time: The Origins of Sequential Simulation and Beyond, *Math. Geosci.*, 53, 193–209, <https://doi.org/10.1007/s11004-021-09926-0>, 2021.
- Goovaerts, P.: *Geostatistics for Natural Resources Evaluation*, Oxford University Press, <https://doi.org/10.1093/oso/9780195115383.001.0001>, 1997.
- Goovaerts, P.: Geostatistical modelling of uncertainty in soil science, *Geoderma*, 103, 3–26, [https://doi.org/10.1016/S0016-7061\(01\)00067-2](https://doi.org/10.1016/S0016-7061(01)00067-2), 2001.
- Gutjahr, A., Bullard, B., and Hatch, S.: General joint conditional simulations using a fast fourier transform method, *Math. Geol.*, 29, 361–389, <https://doi.org/10.1007/BF02769641>, 1997.
- Hanser, S. E.: Elevation in the Western United States (90 meter DEM), USGS [data set], <https://www.sciencebase.gov/catalog/item/542aebf9e4b057766eed286a> (last access: 19 September 2023), 2008.
- Hart, Q. J., Brugnach, M., Temesgen, B., Rueda, C., Ustin, S. L., and Frame, K.: Daily reference evapotranspiration for California using satellite imagery and weather station measurement interpolation, *Civ. Eng. Environ. Syst.*, 26, 19–33, <https://doi.org/10.1080/10286600802003500>, 2009.
- Hashimoto, H., Wang, W., Melton, F. S., Moreno, A. L., Ganguly, S., Michaelis, A. R., and Nemani, R. R.: High-resolution mapping of daily climate variables by aggregating multiple spatial data sets with the random forest algorithm over the conterminous United States, *Int. J. Climatol.*, 39, 2964–2983, <https://doi.org/10.1002/joc.5995>, 2019.
- Hengl, T., Nussbaum, M., Wright, M. N., Heuvelink, G. B. M., and Gräler, B.: Random forest as a generic framework for predictive modeling of spatial and spatio-temporal variables, *Peer J.*, 6, e5518, <https://doi.org/10.7717/peerj.5518>, 2018.
- Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., De Chiara, G., Dahlgren, P., Dee, D., Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer, A., Haimberger, L., Healy, S., Hogan, R. J., Hólm, E., Janisková, M., Keeley, S., Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., De Rosnay, P., Rozum, I., Vamborg, F., Villaume, S., and Thépaut, J.: The ERA5 global reanalysis, *Q. J. Roy. Meteorol. Soc.*, 146, 1999–2049, <https://doi.org/10.1002/qj.3803>, 2020.
- Hoffmann, J.: GeoStats.jl – High-performance geostatistics in Julia, *J. Open Sour. Softw.*, 3, 692, <https://doi.org/10.21105/joss.00692>, 2018.
- Hsiang, S. M., Burke, M., and Miguel, E.: Quantifying the Influence of Climate on Human Conflict, *Science*, 341, 1235367, <https://doi.org/10.1126/science.1235367>, 2013.
- Hudson, G. and Wackernagel, H.: Mapping temperature using kriging with external drift: Theory and an example from scotland, *Int. J. Climatol.*, 14, 77–91, <https://doi.org/10.1002/joc.3370140107>, 1994.
- Ingebrigtsen, R., Lindgren, F., and Steinsland, I.: Spatial models with explanatory variables in the dependence structure, *Spat. Statist.*, 8, 20–38, 2014.
- Kitanidis, P. K.: *Introduction to Geostatistics: Applications in Hydrogeology*, Cambridge University Press, ISBN 9780521587471, 1997.
- Lobell, D. B., Thau, D., Seifert, C., Engle, E., and Little, B.: A scalable satellite-based crop yield mapper, *Remote Sens. Environ.*, 164, 324–333, <https://doi.org/10.1016/j.rse.2015.04.021>, 2015.
- Menne, M. J., Durre, I., Vose, R. S., Gleason, B. E., and Houston, T. G.: An Overview of the Global Historical Climatology Network-Daily Database, *J. Atmos. Ocean. Tech.*, 29, 897–910, <https://doi.org/10.1175/JTECH-D-11-00103.1>, 2012.
- Milà, C., Ballester, J., Basagaña, X., Nieuwenhuijsen, M. J., and Tonne, C.: Estimating daily air temperature and pollution in Catalonia: A comprehensive spatiotemporal modelling of multiple exposures, *Environ. Pollut.*, 337, 122501, <https://doi.org/10.1016/j.envpol.2023.122501>, 2023.
- Morgan, M. G. and Henrion, M.: *Uncertainty: A Guide to Dealing with Uncertainty in Quantitative Risk and Policy Analysis*, Cambridge University Press, <https://doi.org/10.1017/CBO9780511840609>, 1990.
- Olea, R. A.: *Geostatistics for Engineers and Earth Scientists*, Springer US, Boston, MA, <https://doi.org/10.1007/978-1-4615-5001-3>, 1999.
- Pyrz, M. and Deutsch, C. V.: *Geostatistical Reservoir Modelling*, in: 2nd Edn., Oxford University Press, New York, p. 448, ISBN 9780199731442, 2014.
- Rasmussen, R. M., Chen, F., Liu, C. H., Ikeda, K., Prein, A., Kim, J., Schneider, T., Dai, A., Gochis, D., Dugger, A., Zhang, Y., Jaye, A., Dudhia, J., He, C., Harrold, M., Xue, L., Chen, S., Newman, A., Dougherty, E., Abolafia-Rosenzweig, R., Lybarger, N.

- D., Viger, R., Lesmes, D., Skalak, K., Brakebill, J., Cline, D., Dunne, K., Rasmussen, K., and Miguez-Macho, G.: CONUS404: The NCAR–USGS 4-km Long-Term Regional Hydroclimate Reanalysis over the CONUS, *B. Am. Meteorol. Soc.*, 104, E1382–E1408, <https://doi.org/10.1175/BAMS-D-21-0326.1>, 2023.
- Rossi, M. E. and Deutsch, C. V.: *Mineral Resource Estimation*, Springer Netherlands, Dordrecht, <https://doi.org/10.1007/978-1-4020-5717-5>, 2014.
- Thornton, M. M., Shrestha, R., Wei, Y., Thornton, P. E., Kao, S.-C., and Wilson, B. E.: Daymet: Station-Level Inputs and Cross-Validation for North America, Version 4 R1, ORNL DAAC [data set], <https://doi.org/10.3334/ORNLDAAC/2132>, 2022.
- Thornton, P. E., Running, S. W., and White, M. A.: Generating surfaces of daily meteorological variables over large regions of complex terrain, *J. Hydrol.*, 190, 214–251, [https://doi.org/10.1016/S0022-1694\(96\)03128-9](https://doi.org/10.1016/S0022-1694(96)03128-9), 1997.
- Thornton, P. E., Shrestha, R., Thornton, M., Kao, S.-C., Wei, Y., and Wilson, B. E.: Gridded daily weather data for North America with comprehensive uncertainty quantification, *Sci. Data*, 8, 190, <https://doi.org/10.1038/s41597-021-00973-0>, 2021.
- Volk, J. M., Huntington, J. L., Melton, F. S., Allen, R., Anderson, M., Fisher, J. B., Kilic, A., Ruhoff, A., Senay, G. B., Minor, B., Morton, C., Ott, T., Johnson, L., Comini De Andrade, B., Carrara, W., Doherty, C. T., Dunkerly, C., Friedrichs, M., Guzman, A., Hain, C., Halverson, G., Kang, Y., Knipper, K., Laipelt, L., Ortega-Salazar, S., Pearson, C., Parrish, G. E. L., Purdy, A., ReVelle, P., Wang, T., and Yang, Y.: Assessing the accuracy of OpenET satellite-based evapotranspiration data to support water resource and land management applications, *Nat. Water*, 2, 193–205, <https://doi.org/10.1038/s44221-023-00181-7>, 2024.
- White, M. A., Thornton, P. E., and Running, S. W.: A continental phenology model for monitoring vegetation responses to interannual climatic variability, *Global Biogeochem. Cy.*, 11, 217–234, <https://doi.org/10.1029/97GB00330>, 1997.
- Zeng, J., Matsunaga, T., Tan, Z.-H., Saigusa, N., Shirai, T., Tang, Y., Peng, S., and Fukuda, Y.: Global terrestrial carbon fluxes of 1999–2019 estimated by upscaling eddy covariance data with a random forest, *Sci. Data*, 7, 313, <https://doi.org/10.1038/s41597-020-00653-5>, 2020.