

Evaluating the potential and challenges of an uncertainty quantification method for long short-term memory models for soil moisture predictions

Kuai Fang¹, Daniel Kifer², Kathryn Lawson², and Chaopeng Shen²

¹Stanford University

²Pennsylvania State University

November 22, 2022

Abstract

Recently, recurrent deep networks have shown promise to harness newly available satellite-sensed data for long-term soil moisture projections. However, to be useful in forecasting, deep networks must also provide uncertainty estimates. Here we evaluated Monte Carlo dropout with an input-dependent data noise term (MCD+N), an efficient uncertainty estimation framework originally developed in computer vision, for hydrologic time series predictions. MCD+N simultaneously estimates a heteroscedastic input-dependent data noise term (a trained error model attributable to observational noise) and a network weight uncertainty term (attributable to insufficiently-constrained model parameters). Although MCD+N has appealing features, many heuristic approximations were employed during its derivation, and rigorous evaluations and evidence of its asserted capability to detect dissimilarity were lacking. To address this, we provided an in-depth evaluation of the scheme's potential and limitations. We showed that for reproducing soil moisture dynamics recorded by the Soil Moisture Active Passive (SMAP) mission, MCD+N indeed gave a good estimate of predictive error, provided that we tuned a hyperparameter and used a representative training dataset. The input-dependent term responded strongly to observational noise, while the model term clearly acted as a detector for physiographic dissimilarity from the training data, behaving as intended. However, when the training and test data were characteristically different, the input-dependent term could be misled, undermining its reliability. Additionally, due to the data-driven nature of the model, the two uncertainty terms are correlated. This approach has promise, but care is needed to interpret the results.

Abstract

Recently, recurrent deep networks have shown promise to harness newly available satellite-sensed data for long-term soil moisture projections. However, to be useful in forecasting, deep networks must also provide uncertainty estimates. Here we evaluated Monte Carlo dropout with an input-dependent data noise term (MCD+N), an efficient uncertainty estimation framework originally developed in computer vision, for hydrologic time series predictions. MCD+N simultaneously estimates a heteroscedastic input-dependent data noise term (a trained error model attributable to observational noise) and a network weight uncertainty term (attributable to insufficiently-constrained model parameters). Although MCD+N has appealing features, many heuristic approximations were employed during its derivation, and rigorous evaluations and evidence of its asserted capability to detect dissimilarity were lacking. To address this, we provided an in-depth evaluation of the scheme’s potential and limitations. We showed that for reproducing soil moisture dynamics recorded by the Soil Moisture Active Passive (SMAP) mission, MCD+N indeed gave a good estimate of predictive error, provided that we tuned a hyperparameter and used a representative training dataset. The input-dependent term responded strongly to observational noise, while the model term clearly acted as a detector for physiographic dissimilarity from the training data, behaving as intended. However, when the training and test data were characteristically different, the input-dependent term could be misled, undermining its reliability. Additionally, due to the data-driven nature of the model, the two uncertainty terms are correlated. This approach has promise, but care is needed to interpret the results.

1 Introduction**1.1 Time series deep learning for hydrologic predictions**

Recently, we have witnessed the rise of data-driven models, including those based on deep learning (DL), across various scientific disciplines (Shen et al., 2018; Schmidhuber, 2015; LeCun et al., 2015; Goodfellow et al., 2016). In hydrology, time series DL has been employed in predictions of soil moisture (Fang et al., 2017, 2018; Fang & Shen, 2020), water level in urban water networks (D. Zhang et al., 2018), streamflow (Kratzert et al., 2018; Feng et al., 2019), water table depth (J. Zhang et al., 2018), and weather (Wilson et al., 2018), among other applications. A defining characteristic of DL is the depth of the neural network which enables intermediate layers to perform *representa-*

tion learning – automatically deriving problem-relevant features which are then used to predict the outputs (Bengio, 2009). Provided that there is enough training data, this characteristic implies that few pre-processing steps and human-defined features are needed. In some tasks, the networks can engineer better features than human experts (Schmidhuber, 2015).

In our previous work, we showed that a recurrent DL approach, called long short-term memory (LSTM), could learn from the soil moisture dynamics measured by the Soil Moisture Active Passive (SMAP) mission (Fang et al., 2017). A model trained on only one year of data can make strong predictions for another year. Despite the large number of parameters, the DL model did not overfit and was more robust than regularized linear regression and autoregressive models. With 3 years of training data, LSTM could successfully predict multi-year trends in soil moisture for years not included in the training data (Fang et al., 2018). Despite SMAP’s own limitations, this flexible model can be beneficial in a data fusion setting for long-term projections. There remains a substantial potential to utilize DL to improve accuracies for various hydrologic modeling applications with other variables of interest.

1.2 Uncertainties for data-driven models

Despite significant progress with DL models for hydrology, none of the above-mentioned studies addressed model uncertainties, here referring to the estimation of prediction errors. For many practical and scientific purposes, e.g. ensemble data assimilation (De Lannoy et al., 2007) and decision support (Lamontagne et al., 2018), it is as important to obtain the confidence of a prediction as to obtain the prediction itself (Beven, 1989; Pappenberger & Beven, 2006; Ajami et al., 2008). This is even more critical for hydrologic DL models, considering the alien nature of DL models to most hydrologic users. However, no big-data work so far in hydrology has reported uncertainty estimation methods for time series DL models.

Multiple classes of methods have arisen from Bayesian probability theory to estimate uncertainties, with different advantages and disadvantages. For example, the Markov Chain Monte Carlo (MCMC) method adaptively generates new samples that gradually approach the posterior distribution of model parameters (Vrugt et al., 2008). In the context of hydrologic modeling, these models are typically process-based ones with a low-

80 dimensional (10) parameter set. The uncertainty estimate is obtained from sampling
 81 parameter sets from this posterior distribution which is incrementally improved. Unfor-
 82 tunately, MCMC is intractable for DL models that have orders-of-magnitude more pa-
 83 rameters. Aside from the computational cost, another difficulty of this approach is struc-
 84 tural errors from the forward model, as such an approach assumes that the error comes
 85 from uncertainty in the model parameters only (and not from the structure of the model),
 86 but model structure is known to strongly control the errors (Butts et al., 2004).

87 Uncertainty for data-driven models is not a monolithic quantity. It consists of sev-
 88 eral distinct components that can be mathematically modeled as follows. Consistent with
 89 the machine learning literature, the target variable Y (e.g. soil moisture) is a function
 90 of the input X and some random noise whose distribution has dependence on X . In other
 91 words, $Y = f(X) + \epsilon_X$. This function f is unknown and furthermore, due to measure-
 92 ment error, we may have a noisy version \tilde{X} of the inputs (instead of the true X) (Kavetski
 93 et al., 2006). There exists some unknown function f^* that serves as the best predictor
 94 of Y given noisy input \tilde{X} , i.e. $f^*(\tilde{X}) \approx Y$. Now, since f^* is unknown, the goal of ma-
 95 chine learning is to approximate it using a function g with parameters W (hence we write
 96 g_W). Neural networks are known as *universal approximators* (Hornik, 1991) which means
 97 that, under mild regularity conditions that depend on a chosen error metric, any func-
 98 tion can be approximated to any desired level of accuracy by a sufficiently large neural
 99 network with the right choice of weight parameters W^* . However, since W^* is also un-
 100 known, it must be estimated from the data, leading to network weight uncertainty. The
 101 network g_W learned from the data has weights W that are different from W^* (network
 102 weight uncertainty). To summarize, we have 3 sources of error/uncertainty: data noise
 103 (predicting Y using f^*), model mis-specification error (approximating f^* with g_{W^*}), and
 104 network weight uncertainty (approximating g_{W^*} with g_W).

105 Of the three uncertainty terms mentioned above, without improvement in data qual-
 106 ity, only the data noise cannot be reduced by collecting more data. However, data noise
 107 is often related to certain attributes that are known and is thus also input-dependent.
 108 For example, in our case of learning SMAP observations (Fang et al., 2017), SMAP ob-
 109 servations are highly uncertain in regions with large vegetation water content (VWC).
 110 Hence, the magnitude of SMAP data noise could potentially be estimated based on pre-
 111 cipitation and land cover types. The network weight uncertainty, on the other hand, re-
 112 sults from insufficient training data and can be reduced by more data collection (and more

113 effort). As the amount of training data increases, the parameters are better constrained
114 and the prediction uncertainty decreases. The mis-specification error is more pronounced
115 with process-based models, which impose strong constraints on the function space. If these
116 constraints differ from the actual physics, they could be inadequate or inappropriate for
117 the modeling task, under which condition it could be said the model is *mis-specified*. For
118 DL models, as long as the appropriate basic architecture is selected, the effect of mis-
119 specified structure is minor as the constraints are universal approximators. The basic
120 architecture of deep networks such as LSTM is so versatile that these networks can ap-
121 proximate a large range of problems, from speech recognition (Graves et al., 2013), to
122 handwriting synthesis (Graves, 2013), to brain wave interpretation (Kumar et al., 2019),
123 to improving health care (Miotto et al., 2017). Hence in practice the approximation er-
124 ror is dominated by data noise and network weight uncertainty.

125 Some may recognize that the data noise and network weight uncertainty terms are
126 sometimes referred to as the *aleatoric* and *epistemic* uncertainties in the literature of ma-
127 chine learning and some other domains. For example, Kiureghian and Ditlevsen (2009)
128 asserted that “Uncertainties are characterized as epistemic, if the modeler sees a pos-
129 sibility to reduce them by gathering more data or by refining models. Uncertainties are
130 categorized as aleatory if the modeler does not foresee the possibility of reducing them”.
131 This categorization is simple to grasp and is in general agreement with the machine learn-
132 ing literature (Kendall & Gal, 2017; Senge et al., 2014; Depeweg et al., 2017), as well as
133 some hydrology papers (Nearing, Mocko, et al., 2016; Gong et al., 2013; Behrouz & Al-
134 imohammadi, 2018). Data-driven modelers have become accustomed to highly noisy data
135 and have regarded such noise (after due effort in data curation) as irreducible. On the
136 other hand, their knowledge comes from the training data and hence they regard the pa-
137 rameter uncertainty (of a data-driven model) as *epistemic*. However, these definitions
138 clash with some other definitions known to hydrology. On a philosophical level, it is quite
139 difficult to clearly define the limit of what is knowable and what is unknowable, which
140 can be witnessed by a series of historical debates (Beven, 2016; Nearing, Tian, et al., 2016).
141 For example, some would regard noise with data (e.g. precipitation), and observations
142 (e.g. soil moisture readings from SMAP), as epistemic (Beven, 2016), while to a machine
143 learning scientist they would most likely be considered aleatoric. Because the purpose
144 of this paper is largely to evaluate the methods that estimate errors with LSTM mod-
145 els, we avoided the controversial terms.

1.3 Background on Monte-Carlo dropout

Here we examine Monte Carlo dropout with a data noise term (MCD+N). The first part of MCD+N, proposed by Gal and Ghahramani (2016) (hereafter called GG16), can be interpreted as measuring the disagreement among ensemble members generated by applying dropout. The second part of MCD+N is a heteroscedastic input-dependent model for observational noise, proposed by Kendall and Gal (2017) (hereafter called KG17).

The foundational ideas are:

- Dropout (Srivastava et al., 2014) is a training technique that is used to prevent overfitting in deep networks - during each iteration of back-propagation, randomly selected units are ignored. It was originally interpreted as an efficient way of simulating an ensemble of deep networks. GG16 provided another interpretation, that dropout training of deep networks was an approximation of training Gaussian process (GP) models (Rasmussen & Williams, 2005). GG16 proposed the use of dropout during prediction to create random predictions and postulated that the variability of these predictions was a good measure of network weight uncertainty. This use of dropout is called Monte Carlo Dropout (MCD). It is worth noting that this term does not seek to approximate the bias of the network.
- An second output unit can be added to the deep network to be implicitly supervised. With a proper scoring function during training, this unit can be interpreted as an estimate of the variance of the network’s prediction from its original output unit. The goal of the secondary unit is to measure data noise and model it as a function of the inputs.

GG16 revealed a new and surprisingly convenient path toward estimating uncertainty for DL models. A GP models data as multi-variate Gaussian distributions with covariance functions. Without the need for sampling, a GP model could directly prescribe the predictive distribution at a new point. Earlier work showed that with the right activation functions, a neural network with one or more hidden layers and a Gaussian prior on the weights would converge in distribution to a GP as the size of the hidden layers grows to infinity (Neal, 1996; Lee et al., 2018; Matthews et al., 2018). Extending along this avenue, GG16 developed a theoretical framework casting dropout (Srivastava et al.,

176 2014) as an approximate GP, where the sampling of the distribution could be achieved
177 by applying dropout during model testing.

178 GG16’s GP interpretation of dropout training is heuristic in the sense that it in-
179 volves approximations whose accuracies were not quantified (and is a subject for debate
180 (Osband et al., 2016)). Moreover, with respect to the GP argument, it has never been
181 systematically shown in previous studies (Gal & Ghahramani, 2016; Kendall & Gal, 2017;
182 Vandal et al., 2018) that the MCD estimate would predict a smaller error for an instance
183 more similar to the training dataset, and a larger error for instances that are unlike the
184 training data. One barrier was that for the tasks examined in many DL applications, it
185 was difficult to define and visualize proximity. Hence, the effectiveness of the MCD en-
186 semble to quantify similarity has yet to be evidenced.

187 The MCD+N method is appealing due to its simplicity and its support for arbi-
188 trary network architectures. The resulting uncertainty estimates also proved useful in
189 an image segmentation task (Kendall & Gal, 2017). Consequently, the scheme has gar-
190 nered an enormous amount of popularity, which can be witnessed by the high citation
191 count of GG16 (cited 1620 times at the time of writing this article) and KG17. However,
192 the limitations and properties of this method have not been adequately examined. Since
193 the input-dependent uncertainty is estimated by the trained network, it is natural to ques-
194 tion its accuracy in the event that the test data comes from a fundamentally different
195 distribution than the training data the network is based on, i.e., the test data is *out of*
196 *distribution*. Another question is whether the combined uncertainty estimate is of high
197 quality given representative or unrepresentative training data. This work constitutes the
198 first report on MCD+N in hydrology and perhaps also one of the most thorough eval-
199 uations of this scheme in DL, revealing both its potential and limitations.

200 **1.4 Research questions**

201 The goal of this paper is not to promote the MCD+N scheme but to use experi-
202 ments to evaluate the quality and limitations of the scheme for the case of soil moisture
203 predictions, which is the first hydrologic dataset encountered by this method. While satel-
204 lites provide global-scale coverage of surface soil moisture, many other hydrologic data,
205 e.g. streamflow and groundwater levels, are available only locally. Even with satellites,
206 there are regions beyond the scope of satellite, e.g. high latitudes and areas covered with

207 dense vegetation canopy. Therefore, we are concerned with the quality of MCD+N es-
 208 timates when the training data is biased in only part of the domain. We ask the follow-
 209 ing questions:

210 (1) When the training data is representative of the spatial domain, can the MCD+N
 211 uncertainty terms help us anticipate predictive error as measured by unbiased RMSE?

212 (2) Do the two uncertainty estimates behave as asserted, i.e., does the data noise
 213 term respond to stochasticity in the data and does the network weight uncertainty term
 214 respond to dissimilar cases?

215 (3) When a network directly predicts input-dependent uncertainty via a secondary
 216 output unit, is this estimate reliable for time series that are out of the training data dis-
 217 tribution?

218 (4) How are these results affected by hyperparameters such as the dropout rate and
 219 priors on the input-dependent uncertainty output units?

220 It is worth mentioning that the goal of this paper is not to promote the MCD+N
 221 scheme but to use carefully-designed experiments to evaluate its quality.

222 **2 Methods and datasets**

223 As an overview, we trained a probabilistic time series DL model to learn the level-
 224 3 SMAP surface soil moisture product. The input to this DL model included climatic
 225 forcing data and constant geophysical attributes. In addition to the SMAP product, the
 226 network also estimates the input-dependent data noise. The network weight uncertainty
 227 is then estimated via the MCD procedure, which runs many forward realizations of the
 228 stochastic dropout masks during inference (making soil moisture predictions about a new
 229 instance).

230 **2.1 SMAP and input data**

231 The SMAP level 3 radiometer product (L3_SM_P, version 4) measures the global
 232 surface soil moisture since April 2015, with a moisture-dependent sensing depth that is
 233 less than 5 cm. The spatial resolution of L3_SM_P is 36 km, with a revisit time of 2 to
 234 3 days. The DL model was trained with seven climatic forcing inputs: precipitation, tem-
 235 perature, radiation, humidity, pressure, and wind speed (two directions). We obtained

236 the forcing data from North American Land Data Assimilation System phase II (NL-
237 DAS2) (Xia et al., 2015). In addition, the DL model also used static geographic attributes,
238 e.g. soil texture and attributes, from the World Soil Information (ISRICWISE) database
239 (Batjes, 1995), and land surface characteristics from SMAP flags.

240 **2.2 Time series deep learning**

241 The LSTM model used the atmospheric forcing time series and static land surface
242 characteristics described above as inputs. Each valid SMAP pixel over the continental
243 United States (CONUS) was treated as a training instance. Spatial autocorrelation was
244 not explicitly modeled but could be implicitly considered due to the spatial autocorre-
245 lation in the inputs. During training, we used a mini-batch size of 100. A mini-batch bun-
246 dles a small number of training instances together to perform weight updates via vari-
247 ations of stochastic gradient descent (typical deep learning training algorithms cycle over
248 mini-batches while performing updates). The loss function is summed over the mini-batch.
249 This procedure allows for more effective use of the memory of the Graphical Processor
250 Units (GPUs).

251 Because surface soil moisture has short memory, each instance in the mini-batch
252 is 30 days of data randomly taken from the available training data of a randomly selected
253 SMAP pixel. 500 epochs were performed for a training job for our CONUS-scale exper-
254 iment. An epoch has approximately the same number of forward runs as the number of
255 instances. In our case, each epoch contains around 888 mini-batches.

256 Recurrent Neural Networks make use of sequential information by updating hid-
257 den states based on both inputs of the current time step and network states of previous
258 time steps. By implementing a *memory cell* and *gates*, LSTM addressed the *vanishing*
259 *gradient* issue that has prevented effective training for vanilla recurrent networks (Hochreiter
260 & Schmidhuber, 1997). While there are several versions of LSTM units, we use the one
261 specified by the following equations:

$$\text{(input transformation)} \quad x^{(t)} = \text{ReLU}(W_{xx}x_0^{(t)} + b_{xx}) \quad (1)$$

$$\text{(input node)} \quad g^{(t)} = \tanh(\mathcal{D}(W_{gx})x^{(t)} + \mathcal{D}(W_{gh})h^{(t-1)}) + b_g \quad (2)$$

$$\text{(input gate)} \quad i^{(t)} = \sigma(\mathcal{D}(W_{ix})x^{(t)} + \mathcal{D}(W_{ih})h^{(t-1)}) + b_i \quad (3)$$

$$\text{(forget gate)} \quad f^{(t)} = \sigma(\mathcal{D}(W_{fx})x^{(t)} + \mathcal{D}(W_{fh})h^{(t-1)}) + b_f \quad (4)$$

$$\text{(output gate)} \quad o^{(t)} = \sigma(\mathcal{D}(W_{ox})x^{(t)} + \mathcal{D}(W_{oh})h^{(t-1)}) + b_o \quad (5)$$

$$\text{(cell state)} \quad s^{(t)} = \mathcal{D}(g^{(t)}) \odot i^{(t)} + s^{(t-1)} \odot f^{(t)} \quad (6)$$

$$\text{(hidden gate)} \quad h^{(t)} = \tanh(s^{(t)}) \odot o^{(t)} \quad (7)$$

$$\text{(output layer)} \quad f^{(t)} = W_{hy}h^{(t)} + b_y \quad (8)$$

262 The superscript t refers to the time step. For a time step t , the vector of raw in-
 263 puts is $x_0^{(t)}$, the state of the hidden cells is denoted by $h^{(t)}$, the state of memory cells is
 264 denoted by $s^{(t)}$, and the output of the network by $f^{(t)}$. *ReLU* refers to Rectified Lin-
 265 ear units (Glorot et al., 2011). In this equation, σ and *tanh* refer to sigmoid and hyper-
 266 bolic tangent functions, respectively, and they are used as the activation function in the
 267 network. \odot represents point-wise multiplication. The W 's and b 's are the trainable con-
 268 nection weights and constant bias parameters in the network, which are shared by all
 269 time steps. \mathcal{D} is the Dropout operator (Srivastava et al., 2014), which randomly sets some
 270 of the network connections to zero in order to reduce overfitting. During each iteration,
 271 the dropout mask is randomly initialized and remains the same for all time steps. More
 272 details of dropout are provided in Section 2.3.2.

273 **2.3 Probabilistic LSTM Model**

274 Overall, the uncertainty of the model is comprised of an input-dependent data noise
 275 term (Section 2.3.1) and a network weight uncertainty term (Section 2.3.2), following Kendall
 276 and Gal (2017). We let the DL network learn and predict the variance of the input-dependent
 277 uncertainty based on inputs to LSTM. Network weight uncertainty results from insuf-
 278 ficient training data, and according to GG16, is estimated by Monte Carlo Dropout.

279 **2.3.1 Input-dependent data noise**

280 It is well known that SMAP observations are highly uncertain in regions with high
 281 vegetation water content (VWC) due to instrumental limitations. This kind of uncer-

282 tainty can be captured based on many input variables such as vegetation cover and tem-
 283 perature. However, instead of manually prescribing a model for the error, we let the net-
 284 work estimate it and provide it as an output, following KG17. For a model prediction
 285 f , the corresponding observation and error vectors are y and $\epsilon = y - f$, respectively.
 286 We assume the errors come from a Gaussian distribution, with a variance σ_x^2 that is de-
 287 pendent on the input data x : $\epsilon \sim \mathcal{N}(0, \sigma_x^2)$ and $y \sim \mathcal{N}(f, \sigma_x^2)$. Given n data points
 288 (regardless of space or time) $\mathbf{y} = \{y_1, \dots, y_n\}$ and corresponding model predictions $\mathbf{f} =$
 289 $\{f_1, \dots, f_n\}$ and standard deviations $\sigma_{\mathbf{x}} = \{\sigma_{x,1}, \dots, \sigma_{x,n}\}$, the likelihood function is

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_{x,i}^2}} \exp\left[-\frac{(y_i - f_i)^2}{2\sigma_{x,i}^2}\right] \quad (9)$$

290 We ask the LSTM model to output an estimate variance, $\hat{\sigma}_x^2$, for σ_x^2 . For numer-
 291 ical stability, the network will predict $s = \log(\hat{\sigma}_x^2)$. Hence, the LSTM model will have
 292 two nodes at the output layer: $(\mathbf{f}, \mathbf{s}) = F^W(x)$, where F^W is the trained LSTM model
 293 and W is the weight in the network. There is no directly supervising data for s . Rather,
 294 it is implicitly supervised by the regression task. As the network cannot reduce random
 295 errors that cannot be predicted based on the inputs, it is forced to learn the error mag-
 296 nitude. For N SMAP pixels (N is the mini-batch size during training), each with T time
 297 steps, the loss function \mathcal{L} to be minimized is the negative logarithm of Equation 9 across
 298 the data points:

$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^N \sum_{t=1}^T \mathbf{1}_{i,t} [(y_{i,t} - f_{i,t})^2 \exp(-s_{i,t}) + s_{i,t}] \quad (10)$$

299 where i and t are the spatial and temporal indices, respectively, and $\mathbf{1}_{i,t}$ is 1 when there
 300 is a valid SMAP observation and 0 when there is not. Naturally, the $s_{i,t}$ term also serves
 301 as a regularization term to prevent the training from unreservedly decreasing the $\exp(-s_{i,t})$
 302 term to minimize the loss function.

303 **2.3.2 MCD for network weight uncertainty**

304 Each weight update step consists of a forward pass (in which the prediction of the
 305 network is computed) and a back-propagation pass (in which this information is used
 306 to compute an approximate gradient for updating the weights). In the dropout method,
 307 a randomly chosen set of nodes is ignored for each weight update step (the ignored nodes
 308 do not affect the prediction in the forward pass). The choice of which nodes to keep and
 309 which to (temporarily) drop is implemented via the *dropout mask*.

310 GG16 proposed the use of dropout during the test step (inference) to generate ran-
 311 dom predictions. MCD runs M forward realizations, $f^{\widehat{W}_j}, j \in 1, \dots, M$, with each set
 312 of weights \widehat{W}_j obtained by randomly sampling dropout masks at the same locations where
 313 dropout is applied during training. In contrast, the normal use of dropout during infer-
 314 ence would turn the dropout operators into a multiplication operation with constant scalars
 315 related to the dropout rate, with all connections enabled. The average of the MCD re-
 316 alizations becomes the overall prediction, and their variance is interpreted as a measure
 317 of uncertainty. GG16 recommended that MCD only be used for networks that are also
 318 trained using dropout. The mean and variance of the MCD ensemble for a prediction
 319 f are:

$$E[\mathbf{f}] \approx \frac{1}{M} \sum_{m=1}^M f^{\widehat{W}_m}(x) \quad (11)$$

$$\sigma_{mc}^2[\mathbf{f}] \approx \frac{1}{M} \sum_{m=1}^M f^{\widehat{W}_m}(x)^2 - E[\mathbf{f}]^2 \quad (12)$$

320 MCD can be interpreted intuitively from an ensemble simulation perspective, just
 321 like dropout training (Srivastava et al., 2014). Each realization of the dropout mask forms
 322 a sub-network. The random predictions arising from multiple randomly chosen masks
 323 can then be viewed as predictions coming from an ensemble of related sub-networks. These
 324 sub-networks would be in stronger agreement (hence smaller variance) in regions where
 325 the input space is well conditioned by known data points. Further away from the train-
 326 ing data, the sub-networks may diverge more significantly. Nevertheless, it is very chal-
 327 lenging to formally prove this intuition.

328 The primary contribution of GG16 was that they noted connections between dropout
 329 training and variational Bayesian inference of GP (an overview of their arguments and
 330 a discussion of issues can be found in Appendix A). Their main argument was that if
 331 variational inference was conducted with respect to network weights, with a special set
 332 of variational distributions, it would approximately lead to the same loss function as dropout
 333 training with mini-batching, as described in Equation 10. In this way, each realization
 334 with a set of randomly sampled dropout masks is equivalent to sampling from the pos-
 335 terior variational distribution. Although the approximation error was generally not quan-
 336 tified, this connection inspired their proposal of using MCD as an estimate of model un-
 337 certainty (since this is what the posterior distribution of a GP corresponds to). In com-
 338 puter vision tasks, GG16 and KG17 found that MCD was useful as an uncertainty mea-

339 sure – the estimated uncertainty tended to be large when the prediction of the network
 340 was inaccurate.

341 2.3.3 Combining uncertainties

In their analysis of the connections between GP and MCD in deep networks, GG16 noted that the variance of the posterior distribution depends on the variance of the prior as well as the dropout retention rate β . These factors suggest that the network weight uncertainty term needs to be calibrated. GG16 suggested linearly scaling the model uncertainty term to match the predictive error magnitude, i.e.,

$$\sigma_{mc}^2(f_{i,t}) \approx \alpha \left\{ \frac{1}{M} \sum_{m=1}^M f_{i,t}^{\widehat{W}_m}(x)^2 - \left[\frac{1}{M} \sum_{m=1}^M f_{i,t}^{\widehat{W}_m}(x) \right]^2 \right\} \quad (13)$$

Another option is to find β^* , the optimum β value, to best capture the correct uncertainty magnitude, i.e.,

$$\sigma_{mc}^2(f_{i,t}) \approx \frac{1}{M} \sum_{m=1}^M f_{i,t}^{\widehat{W}_m(\beta^*)}(x)^2 - \left[\frac{1}{M} \sum_{m=1}^M f_{i,t}^{\widehat{W}_m(\beta^*)}(x) \right]^2 \quad (14)$$

342 Here $f_{i,t}^{\widehat{W}_m(\beta^*)}(x)$ is the prediction for input x when the network uses the weight param-
 343 eters \widehat{W}_m obtained by applying dropout with rate β to the trained network.

344 Given $y \sim \mathcal{N}(f, \sigma_x^2)$ and the model uncertainty as calculated in Equation 12, the
 345 total uncertainty variance is σ_{comb}^2 :

$$\sigma_{comb}^2 = \sigma_{mc}^2 + \sigma_x^2 \quad (15)$$

346 where (i, t) are dropped for brevity.

347 The hyperparameter β^* or α , depending on which calibration method was chosen,
 348 needs to be tuned. For the scope of this work, we chose to tune β^* as it is a simpler pro-
 349 cedure, and we found a constant β^* to be sufficient for improving the quality of the un-
 350 certainty. We used the first year of the SMAP data as training data, and the second year
 351 as the validation data for hyperparameter tuning. Hyperparameters were adjusted so that
 352 the estimated combined error σ_{comb}^2 matched the predictive error in the spatial regions
 353 where the model was trained. To avoid over-tuning, we did a lazy search (meaning with-
 354 out sophisticated searching) for a uniform β^* value in all layers and locations, although
 355 we recognize that β^* could, in theory, be different from location to location. The third
 356 year of SMAP data was used as a test dataset entirely for the purpose of evaluation.

2.4 Evaluation of the uncertainty quality

In all of our experiments, we used the level-3 SMAP surface soil moisture product over the CONUS as the training target. As mentioned earlier, we used the first year of data (2015/04 - 2016/03) as the training data, the second (2016/04 - 2017/03) for validation and hyperparameter tuning, and the third (2017/04 - 2018/03) as the test data for the evaluation of metrics. The quality of uncertainty was evaluated by both the predictive errors and the cumulative distribution of the likelihood function. For the predictive errors, we compared the magnitude of σ_{comb} , the standard deviation of the combined errors, to that of the unbiased root-mean-square error (ubRMSE) when predicting SMAP surface soil moisture in the test period. We also calculated the Pearson's correlation coefficient (R) between ubRMSE and σ_{comb} .

Similar to KG17 and Vandal et al. (2018), we calculated an error exceedance likelihood, $p_{ee}(|e| > |y - f|; \sigma^2) = 1 - \frac{\text{erf}(-|y-f|)}{2\sigma}$, $e \sim \mathcal{N}(0, \sigma^2)$, which is the self-assessed chance that an error of this magnitude ($|y - f|$) or worse could happen, given an uncertainty estimate σ^2 . By this definition, if the uncertainty estimate is perfect, for a large error marked with a 0.01 exceedance likelihood, we expect to see that it is exceeded roughly 1% of the time. Similarly, for an error estimate exceeded 40% of the time, we expect to see a calculated error exceedance likelihood of 0.4. As a result, when the cumulative distribution function (CDF) of p_{ee} is plotted (called the calibration plot in KG17), we would like to see it being close to a one-to-one line. We further calculated d , the maximum distance of the CDF from the 1:1 line, also called the Kolmogorov-Smirnov distance between two empirical CDFs. d thus serves as a succinct measure of the quality of the uncertainty estimate. A d value of 0 would mean a perfect uncertainty quality, while a d value close to 0.5 would suggest very poor quality. The error exceedance likelihoods calculated using σ_x , σ_{mc} , and σ_{comb} as σ^2 are referred to as p_x , p_{mc} , and p_{comb} , respectively. Evaluating p_{ee} separately with these variances helps us to understand how each component of the uncertainty estimate works.

2.5 Training experiments and evaluations

2.5.1 CONUS-scale generalization test

We trained a LSTM model over the entire CONUS from 2015/04 to 2016/03, with spatial downsampling done by picking 1 pixel from every patch of 2 x 2 pixels. To eval-

388 uate the overall quality of the uncertainty estimation, we ran both a temporal test and
 389 a regular spatial test. In the temporal generalization test, the model was tested on the
 390 same pixels as the training set but with the third year of data (2017/04 to 2018/03). In
 391 the regular spatial generalization test, the model was tested on the same period as the
 392 training set, but with the neighboring pixel in the diagonal direction, which was not part
 393 of the model’s training data.

394 **2.5.2 Noise perturbation experiments**

395 According to the theory discussed by KG17, the input-dependent data noise term
 396 could directly detect observation error, while the model parameter uncertainty could not.
 397 To test this theory, we examined how the input-dependent data noise (σ_x) and network
 398 weight uncertainty (σ_{mc}) each responded to noise introduced to the learning target. Here
 399 we prescribed an independent zero-mean Gaussian relative noise value with variance σ_{noise}^2 ,
 400 which was added to the observation data as

$$y_{noise} = y + \mathcal{N}(0, \sigma_{noise}^2) \quad (16)$$

401 Ten independent models were trained by adding different levels of noise as $\sigma_{noise} \in \{0.1, 0.2, \dots, 1.0\}$.
 402 The results of the noise perturbation experiments are presented in Section 3.2.

403 **2.5.3 Spatial extrapolation experiments**

404 As discussed earlier, a primary objective of uncertainty analysis is to measure the
 405 model confidence when making predictions for new and potentially unfamiliar instances.
 406 For example, a GP assigns high posterior uncertainty to instances that are dissimilar from
 407 the training data and low posterior variance to instances that are similar. Ideally, a neu-
 408 ral network trained with dropout would exhibit similar behavior.

409 Thus we tested how the proposed uncertainty estimates respond to instances sim-
 410 ilar to (or dissimilar from) the training dataset with two sets of experiments. Similar-
 411 ity, defined as the proximity between instances in a space spanned by inputs that are rel-
 412 evant to the prediction target, can be difficult to judge, so here we use geographic prox-
 413 imity and ecoregion hierarchy as proxies. Based on US Environmental Protection Agency
 414 (EPA) Ecoregions, which are areas where ecosystems are generally similar (McMahon
 415 et al., 2001), we divided the entire CONUS into 17 sub-regions of relative similar sizes.
 416 To achieve this, we broke the largest ecoregion into several smaller ones and merged the

417 smallest ecoregions into bigger ones. The ecoregions are hierarchical, i.e., ecoregions under
 418 the same level-1 or level-2 codes will be more similar to each other than the ones with
 419 different level-1 or level-2 codes. These ecoregions represent a wide diversity of landscapes,
 420 land covers, soils, and climates over the CONUS.

421 In the first set of experiments, we trained a LSTM model on each of the ecoregions
 422 using year one data, adjusted hyperparameters on these training ecoregions using year
 423 two data, and examined standard deviations for data noise (σ_x), *networkweightuncertainty*(σ_{mc}),
 424 and combined uncertainty σ_{comb} when the model was tested in other regions with year
 425 three data. Our hypothesis was that if MCD indeed captures the network weight uncer-
 426 tainty, then σ_{mc} should be small in regions similar to the training region and large in
 427 dissimilar regions. For comparison, we also attempted a different division strategy, 18
 428 level-2 hydrologic cataloging units (HUC2), and show the results in the Appendix.

429 In the second set of experiments, we trained the models on several combinations
 430 of ecoregions. Some of these ecoregion combinations are dispersed throughout different
 431 parts of the CONUS (hence were more likely to be representative of the background test-
 432 ing data), while three of the combinations were clustered towards only part of the CONUS
 433 (hence were more likely to be biased). These tests allowed us to examine whether use-
 434 ful uncertainty measures could be produced using a small subset of available data.

435 **3 Results and Discussion**

436 **3.1 Uncertainty quality**

437 We first examined the impacts of the dropout retention rate β on uncertainty es-
 438 timates and predictive error. The network weight uncertainty was clearly a function of
 439 β , and we found $\beta \approx 0.4$ to be an approximate value that enabled both accurate pre-
 440 dictions and high-quality uncertainty estimates during the validation period (Appendix
 441 B, Figure B.1). This was the case for either CONUS-scale models or regional-scale mod-
 442 els. To avoid fine tuning, we used $\beta = 0.4$ for all of our evaluations. This result also
 443 suggests that it is useful to calibrate the network weight uncertainty before using it to
 444 anticipate errors.

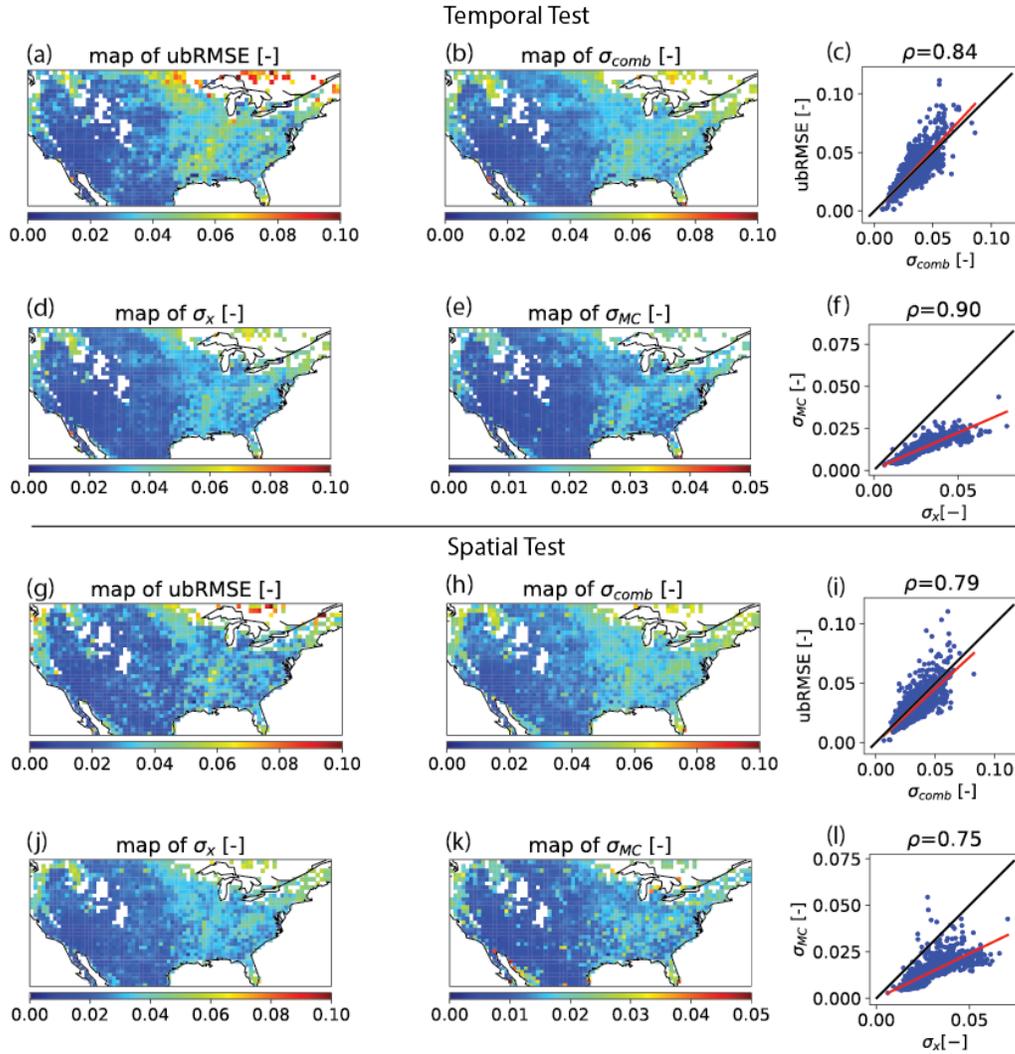
445 The spatial patterns of both data noise (σ_x) and model uncertainty (σ_{mc}) agreed
 446 more or less with the predictive metric of unbiased root-mean-square error (ubRMSE),
 447 and were larger in the eastern CONUS than in the western CONUS (Figure 1 maps).

448 In particular, the northern central CONUS and northeast and northwest coastal regions
 449 had large ubRMSE along with large σ_x . The eastern half of the CONUS, in general, had
 450 larger annual precipitation than the western half. The magnitudes of soil moisture fluc-
 451 tuations, and consequently the magnitudes of measurement errors, were larger. In the
 452 northern CONUS, forest land cover is prominent and a larger fraction of precipitation
 453 falls as snow, so the SMAP signal is adversely impacted by large vegetation water con-
 454 tent (VWC) (O'Neill et al., 2016). Soil moisture cannot be accurately sensed below freez-
 455 ing conditions, which further reduces the amount of training data available (Fang et al.,
 456 2018). As a result, the northeastern and northwestern (along the Rocky mountains) forests
 457 had the highest ubRMSE. The lowest errors were found on the Great Plains and in the
 458 southeastern CONUS, due to arid conditions and reduced forest cover, with associated
 459 low VWC. The predicted σ_x automatically captured these spatial patterns. A belt-like
 460 region with large errors was found along the Mississippi River, which descends along curved
 461 state boundaries into the Gulf of Mexico in the south. This large noise may be associ-
 462 ated with (i) signal leakage from the Mississippi River; or (ii) extensive irrigation due
 463 to cultivated crops along the Mississippi, but, interestingly, σ_x captured it nonetheless.

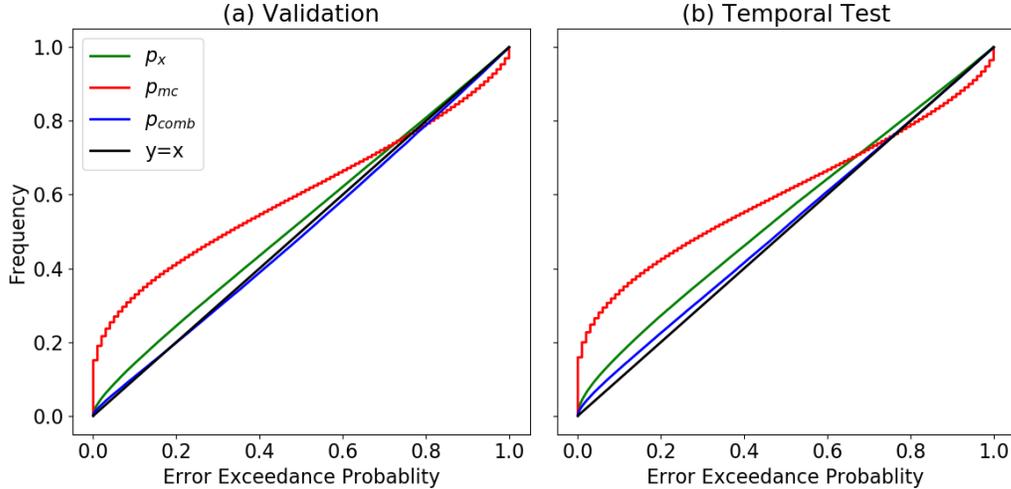
464 On scatter plots of these results, we note a high Pearson's correlation coefficient
 465 value ($R=0.84$) between ubRMSE and σ_{comb} with a small under-estimation bias (Fig-
 466 ure 1c). For the regular spatial generalization test, the correlation was still around 0.79
 467 (Figure 1i). The relationship between σ_{comb} and ubRMSE was heteroscedastic, with more
 468 spread toward the wetter range. In addition, we found that σ_x was larger than σ_{mc} in
 469 both cases, but the two terms were correlated (Figure 1f, Figure 1l).

470 These results suggest that for cases of temporal prolongation or mild spatial ex-
 471 trapolation, it is possible to anticipate model predictive errors using σ_{comb} , while using
 472 either σ_x or σ_{mc} alone would result in under-estimation of the error. In particular, we
 473 can anticipate that if the predicted σ_{comb} is below 0.03, the actual model error will be
 474 closely bounded to the range of 0–0.03. When σ_{comb} is larger than 0.05, however, we
 475 should anticipate large errors, even though ubRMSE may be coincidentally small. The
 476 results suggest that we can use the σ_{comb} map to identify regions where SMAP does not
 477 function properly. In addition, as observed by Pan, Cai, Chaney, Entekhabi, and Wood
 478 (2016), the low uncertainty in the southeast coastal plains is noteworthy. The small er-
 479 ror indicates that SMAP has a reasonable value in this region.

480 The calibration plots of error exceedance likelihoods (Figure 2) show the quality
481 of each uncertainty-estimating component. p_{mc} in both panels lies above the 1:1 line to-
482 ward the left end (e.g. for a p_{mc} of 0.2, a cumulative frequency of $\tilde{0}.39$ is obtained), which
483 means that large predictive errors occurred more frequently than anticipated. Hence, the
484 pattern means that σ_{mc} alone under-estimated the uncertainty toward the large-error
485 range. On the other hand, if we had only considered σ_x , the uncertainty would be slightly
486 under-estimated. In both validation and temporal tests, σ_{comb} was closer to the one-to-
487 one line than either individual component. Since the validation period was employed to
488 identify the optimal β , p_{comb} was almost perfect. In the test period, there was a slightly
489 bigger gap between p_{comb} and the 1:1 line, but the difference still remained small, with
490 a KolmogorovSmirnov distance of 0.027.



491 **Figure 1.** Model error and uncertainty estimates of temporal and spatial generalization tests
 492 over the CONUS. The top two rows (a-f) show temporal test results, and the bottom two rows
 493 ((g)-(l)) show spatial test results. For each of these tests, the left two columns show maps of
 494 model test error (unbiased root-mean-square error, $ubRMSE$) and three uncertainty estimates:
 495 data noise (σ_x), network weight uncertainty (σ_{mc}), and combined uncertainty (σ_{comb}). Note
 496 that the plots of σ_{mc} ((e), (k)) have a narrower numeric range for the same color range as the
 497 other uncertainty estimates, as the range of σ_{mc} is smaller than those of the others. For the
 498 two maps in each row, the one-to-one comparison is shown on the right column, with each point
 499 corresponding to one pixel on the maps, red lines representing lines of best fit, and black lines
 500 representing $y = x$.



501 **Figure 2.** Calibration plots of error exceedance likelihoods computed using network weight
 502 uncertainty (p_{mc}), data noise (p_x), and combined error (p_{comb}) for the (a) validation set
 503 (2016/04-2017/03) and (b) test set (2017/04-2018/03) of the CONUS-scale temporal gener-
 504 alization test. x-axes are estimated error exceedance likelihoods (p_{ee}) based on the different
 505 variances given, and y-axes are the cumulative frequencies, so these curves are the cumulative
 506 distribution functions (CDFs) of p_{ee} , given an uncertainty estimate. The left end of the x-axis
 507 represents large errors, and the right end represents smaller errors. An ideal uncertainty estimate
 508 would produce a CDF that is identical to a 1:1 plot (black lines). The uncertainty qualities, d
 509 values (maximum distance of the CDF from the 1:1 line, section 2.4), of p_x , p_{mc} , and p_{comb} were
 510 0.045, 0.230, and 0.015 for the validation set, and 0.072, 0.241, and 0.027 for the temporal test,
 511 respectively.

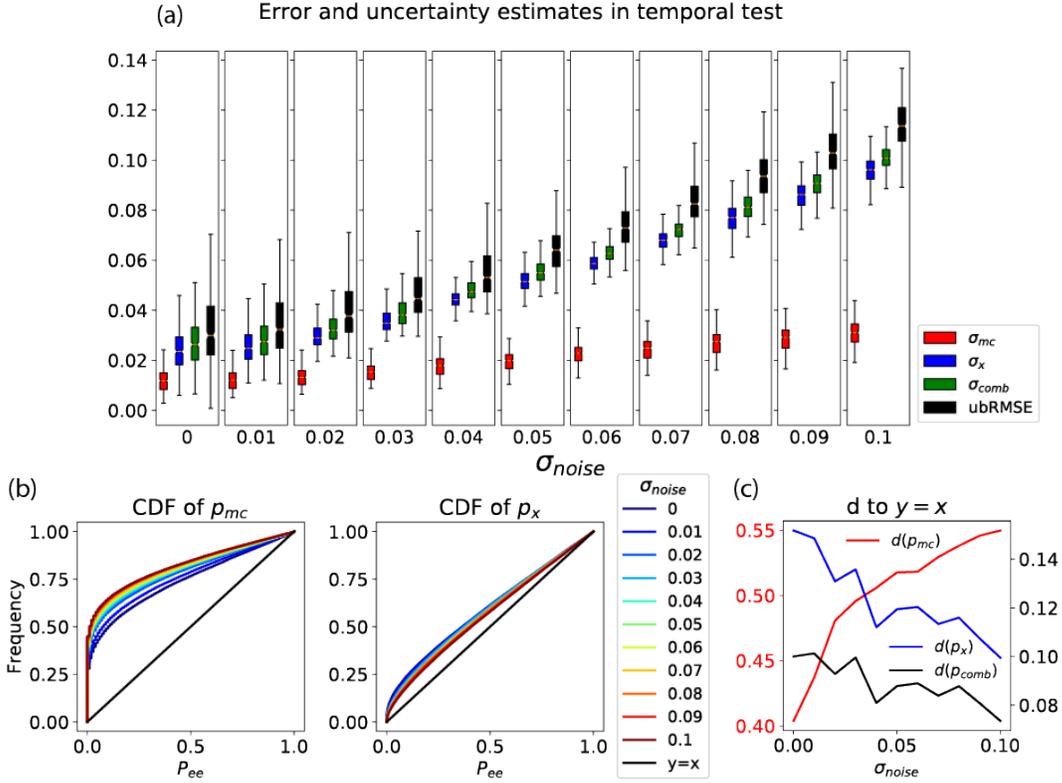
512 3.2 Responses of uncertainty estimates to noisy targets

513 The observation that the two uncertainty estimates were correlated needed further
 514 investigation. Were they correlated because they partially measured the same type of
 515 uncertainty, or because the presence of different uncertainties themselves were correlated
 516 in the SMAP prediction task? In other words, were they correlated because regions with
 517 smaller amounts of training data (leading to larger network weight uncertainties) also
 518 tended to have higher data uncertainties? We thus added noise into the observations to
 519 increase the apparent data uncertainty. In the ideal case, this would cause σ_{mc} to remain
 520 unchanged and σ_x to increase by the same amount as the noise.

521 When the model was trained on the whole CONUS without added noise, the me-
522 dian ubRMSE was around 0.03, smaller than the design accuracy of SMAP. When we
523 added Gaussian random noise, test error and estimated uncertainties all increased. σ_{comb}
524 maintained roughly the same magnitude as $ubRMSE$, with a slight under-estimation
525 (Figure 3a). σ_x responded much more strongly to noise than σ_{mc} , which shows that the
526 proposed data noise scheme is effective at estimating random noise with the target. LSTM
527 could not predict the random noise, and the part that was uncapturable was correctly
528 attributed to the data noise term, especially toward the high noise levels. This result shows
529 that this decomposition of uncertainty could be reasonable at least when the training
530 data are representative.

531 We note in Figure 3a that σ_{mc} also increased with noise, albeit gradually. This ob-
532 servation is consistent with the spatial patterns shown in Figure 1 and the correlation
533 between the two uncertainty terms, and is not in conflict with the meaning of the two
534 terms. Unsurprisingly, significant observational noise led to reduced useful supervising
535 data and thus more ambiguous network weights. Even though σ_{mc} can, in theory, be re-
536 duced by the addition of more data, when noise is significant, the demand for data is am-
537 plified. As a result, the resulting training data is not sufficient at high noise levels.

538 We wanted to see how the quality of two uncertainty estimates changed with the
539 noise in observational data. As Figure 3b and c show, the quality of σ_x increased with
540 noise, as the data noise component could explain more of the total uncertainty. The net-
541 work weight component, on the contrary, was less and less important with respect to the
542 total error. This observation agrees with the naming of the data noise term.



543 **Figure 3.** Performance of model trained by noise-added observations. (a) shows matrices of
 544 uncertainty estimates for network weight uncertainty (σ_{mc}), input-dependent data noise (σ_x),
 545 and combined uncertainty (σ_{comb}), as well as test error ($ubRMSE$). (b) shows calibration plots
 546 of error exceedance likelihoods for different noise levels (p_{mc} , p_x). (c) shows the uncertainty
 547 quality, d (the maximum distance between each CDF and the one-to-one line), varied with noise
 548 added to observations. $d(p_{mc})$ is plotted using the left y-axis while $d(p_x)$ and $d(p_{comb})$ are plotted
 549 using the right y-axis.

550 3.3 Response of uncertainty estimates to dissimilarity

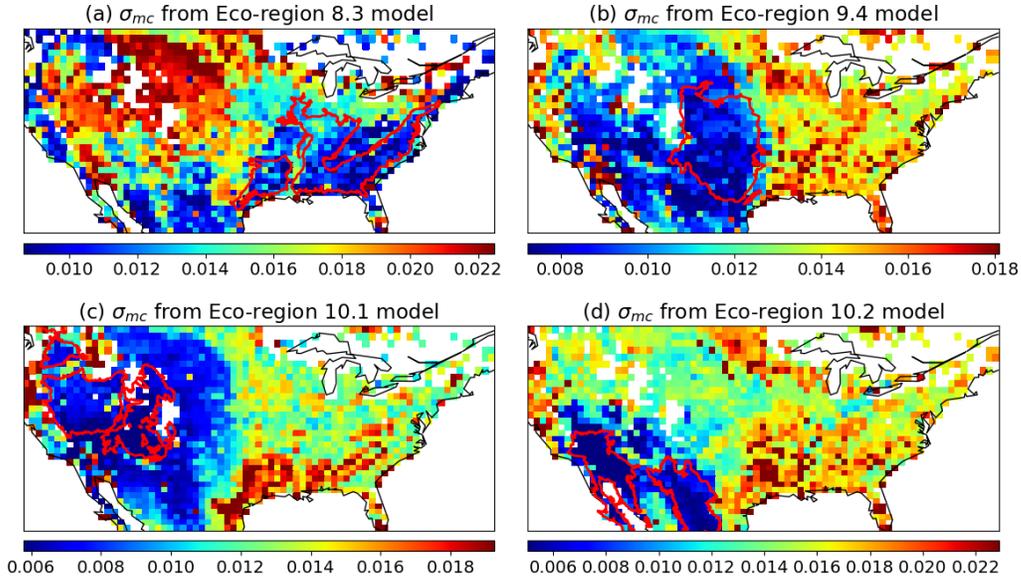
551 The results in Sections 3.1- 3.2 were obtained from models trained on the entire
 552 CONUS. In the following sections we show results from models trained over parts of the
 553 CONUS, which explore how the uncertainty terms respond to out-of-training instances.
 554 We questioned whether the network parameter uncertainty adequately captured dissim-
 555 ilarity.

556 Overall, we see a clear influence of geographic proximity on network weight uncer-
 557 tainty, σ_{mc} , as a result of spatial autocorrelation in the attributes. When we tested mod-

558 els that were only trained on a single level-2 ecoregion, σ_{mc} was smallest inside the train-
559 ing region, somewhat larger in neighboring regions, and much larger further away (Fig-
560 ure 4). We only show models trained on four of the level-2 ecoregions here, but other cases
561 behaved similarly. We show several results in Figure C.1 in Appendix C with similar re-
562 sults when using HUC2 as training regions. These results provided the clearest visual
563 evidence so far that MCD does detect dissimilarity.

564 However, spatial distance itself was not the causal factor for autocorrelation. There
565 is a visible contrast along the eastern edge of the training ecoregion in Figure 4b. This
566 gradient shows where the Great Plains descends to the central plains, and also the di-
567 vide between the drier western half and the wetter eastern half. Some pixels immediately
568 adjacent to the east of the training ecoregion had much larger σ_{mc} than the western neigh-
569 boring pixels, which suggests the model used precipitation and temperature as impor-
570 tant factors in deciding similarity in terms of soil moisture dynamics.

571 It is important to remember that σ_{mc} also depends on the training data, so while
572 it tends to be reciprocal, it may not always be. For example, when the model was trained
573 on ecoregion 8.3 (Southeastern Plains, 4a), it regarded the the western coastal regions
574 and some parts of the southwestern hot desert (parts of ecoregion 10.2, which is the red-
575 highlighted training region selected in Figure 4d) as being similar, and regarded the north-
576 ern high plains (including ecoregion 9.4 and 10.1, which are training regions highlighted
577 in Figure 4c and d, respectively) as being dissimilar. As expected, models trained on ecore-
578 gion 9.4 and 10.1 (results shown in Figure 4c and d) also identified ecoregion 8.3 (train-
579 ing region in 4a) as being dissimilar. However, the model trained on ecoregion 10.2, most
580 of which was found to be similar to ecoregion 8.3 by the model in Figure 4a, regarded
581 the ecoregion 8.3 as dissimilar. This might be due to the more homogeneous environ-
582 ment of ecoregion 10.2 (hot desert). When a model is trained here, it has limited knowl-
583 edge of what soil moisture may do in a wetter environment. When the model was trained
584 in ecoregion 8.3 (wetter and relatively more diverse), it was trained on data with larger
585 gradients in rainfall and appeared to be more confident to predict in ecoregion 10.2.

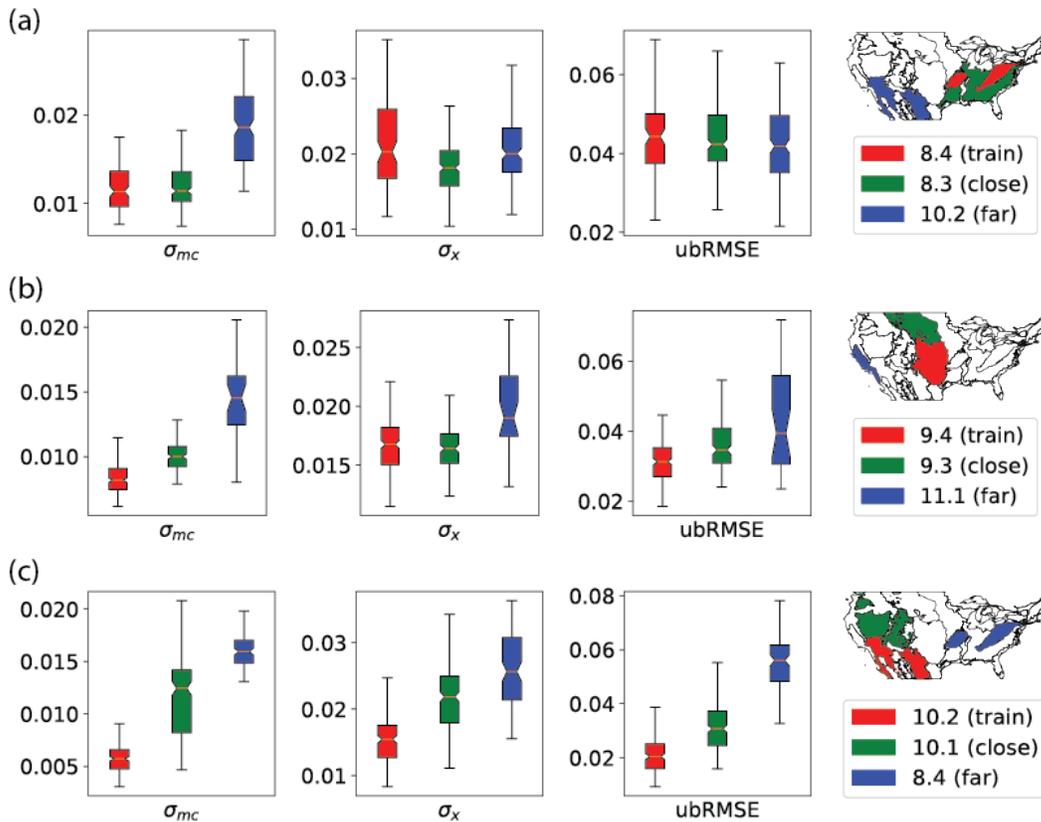


586 **Figure 4.** Maps of network weight uncertainty (σ_{mc}) when the LSTM model was trained
 587 on single level-2 ecoregions. The training region for each model instance is highlighted by the
 588 red polygon. The four selected ecoregions are a) 8.3 Southeastern Plains; b) 9.4 South-central
 589 Semiarid Prairies; c) 10.1 Cold Deserts; d) 10.2 Warm Deserts

590 The responses to similarity can become more clear via bar plots based on the ecore-
 591 gion hierarchy (Figure 5), where the model was trained on one level-2 ecoregion and tested
 592 on another one belonging to the same level-1 ecoregion (the close ecoregion), and another
 593 one belonging to a different level-1 (the far ecoregion). In all three cases, σ_{mc} was much
 594 larger for the far ecoregions as compared to the close ones. Similar to what was suggested
 595 in Figure 4, σ_{mc} correctly provided warnings for instances that were dissimilar to the train-
 596 ing region, and could discern that one region was more dissimilar than another.

597 In contrast, σ_x was not controlled by ecoregion similarity, but represented a pre-
 598 diction of the error based on the inputs, especially precipitation. The predictions seemed
 599 to be largely correct when we qualitatively examined Figure 5, although they may not
 600 be quantitatively perfect. In case (a), σ_x was smaller for both close and far ecoregions
 601 than for the training ecoregion (Figure 5a). Here the model was trained in the north-
 602 eastern region, which has heavy forest cover and more months in a year with frozen soil,
 603 and thus large measurement error. It was tested in ecoregion 10.2, which has much drier
 604 conditions, and should therefore have smaller errors. This was reflected in the smaller

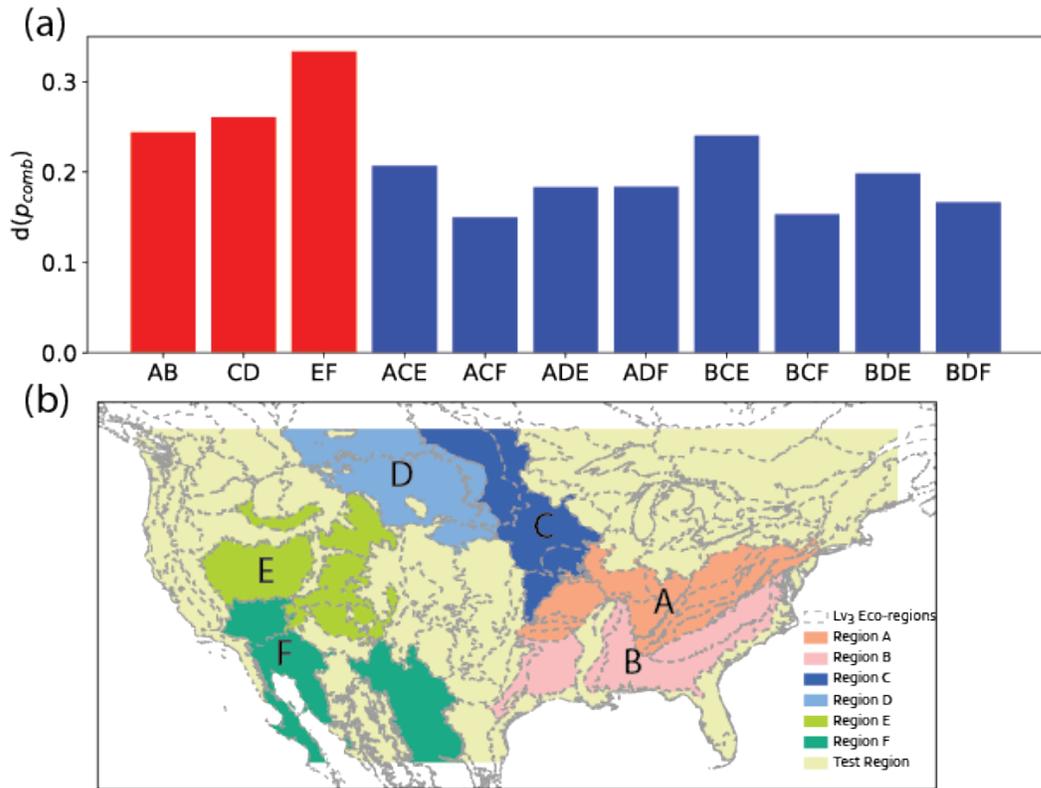
605 σ_x for ecoregion 10.2, but we would have expected the σ_x to be even smaller than the
 606 actual estimate. In case (b), σ_x was similar for the training and the close ecoregions, and
 607 larger for the coastal ecoregion of 11.1 (Figure 5b). Ecoregion 11.1 has larger rainfall than
 608 the inland regions and thus larger error, which was correctly captured by σ_x . In case (c),
 609 the model was trained in a drier region and tested in ecoregion 8.4, which is both dif-
 610 ferent (higher σ_{mc} expected) and much wetter (higher σ_x expected). Therefore, σ_x and
 611 σ_{mc} seemed to indeed reflect different parts of the uncertainty and agreed with our ex-
 612 pectation in terms of the general patterns, but quantitatively the quality could be lim-
 613 ited by the training data (Figure 5c). We show similar results from HUC2 training re-
 614 gions in Figure C.2 in Appendix C.



615 **Figure 5.** Metrics of performance when we trained the model in one level-2 ecoregion, and
 616 tested in two other level-2 ecoregions: one similar to the training region (from the same level-1
 617 ecoregion), one farther away (from a different level-1 ecoregion). Performance metrics are network
 618 weight uncertainty (σ_{mc}), input-dependent data noise (σ_x), and test error ($ubRMSE$).

619 As σ_x is dependent on the training region, to further explore its limitations we in-
620 vestigated the performance of models when they were trained on several ecoregion com-
621 binations and tested on the rest of the CONUS. When the ecoregion combinations spanned
622 across the CONUS, occupying a variety of landscapes in the CONUS domain (blue bars
623 on Figure 6a), the estimated uncertainties were of higher quality. When the chosen ecore-
624 gions were clustered in only part of the CONUS domain (grouped as AB, CD, or EF,
625 shown in Figure 6b), the estimated uncertainties were of much lower quality (higher d
626 values). The combination EF had the lowest uncertainty quality, as these two regions
627 are clustered together in the western arid landscape. Due to this aridity, the model trained
628 there predicts small soil moisture fluctuations and also small σ_x when tested on other
629 regions, resulting in significant under-estimation of the data noise term. We also noticed
630 that whenever region F (warm deserts) was included in a combination in place of region
631 E (cold deserts), the quality tended to be lower. This is presumably because the arid-
632 ity of E is less extreme than F. As a result, including F instead of E expands the cov-
633 erage of the training data in terms of the aridity scenarios.

644 This result can be explained by the fact that the data noise term was a trained out-
645 put from the network, and was thus also conditioned by the training data. It provides
646 direct evidence that σ_x could be misled by a strongly biased or unrepresentative train-
647 ing set. It is worth noting that the more representative sets (first three combinations)
648 only sampled a fraction of the domain and are still far from representing the wide di-
649 versity of soil, land cover, and terrain combinations over the CONUS. However, they did
650 provide more variety in the training data, and so it follows that σ_x reported by a model
651 trained on one of these more varied datasets was more representative than σ_x reported
652 by a model trained on a more biased training dataset.



634 **Figure 6.** Evaluation of uncertainty quality (smaller d for higher quality) when models were
 635 trained on different combinations of ecoregions. The metrics were calculated in common regions
 636 of the CONUS that were outside of the training set. (a) Quality metric for combined error ex-
 637 ceedance likelihoods ($d(p_{comb})$), the lower the better) of 11 combinations of regions, where 3 red
 638 bars show region combinations that are spatially clustered (AB, CD, EF) and 8 blue bars show
 639 region combinations that are spatially dispersed. Letters denote which regions are combined (e.g.
 640 ACE refers to a combination of regions A, C, and E). (b) Map of regions, some of which are com-
 641 posed of multiple level-3 ecoregions. A: ecoregions 8.3.1, 8.3.2, 8.3.3, 8.3.4, and 8.4; B: ecoregions
 642 8.3.4, 8.3.5, 8.3.6, 8.3.7, and 8.3.8; C: ecoregion 9.2; D: ecoregion 9.3; E: ecoregions 10.1.4, 10.1.5,
 643 10.1.6, 10.1.7, and 10.1.8; F: ecoregion 10.2.

653 3.4 Further discussion, limitations, and future work

654 The data noise term σ_x , which is essentially a trained, network-predicted error model,
 655 is shown to be a powerful technique with important implications for hydrology to sim-
 656 plify our workflow. Its quality and clear response to data noise suggest the plausibility
 657 of training such error models with very loose specifications of data noise. In the past,

658 a wealth of research has been dedicated to modeling error, e.g., specify error structures
659 and adjustments for heteroscedasticity and autocorrelation (Evin et al., 2013; Göttinger
660 & Bárdossy, 2008; Smith et al., 2015). The proposed procedure greatly relaxes the as-
661 sumptions we need to make to obtain error models. The complex, possibly nonlinear,
662 and potentially time-varying dependencies of the error on input terms can hardly be pre-
663 scribed by experts. We can conveniently delegate such estimation to the deep learning
664 algorithm itself, with the requirement that the training data must be representative.

665 The uncertainty with respect to climate or weather projections, a large and chal-
666 lenging research topic, has not been quantified here. For short-term forecast problems,
667 the impacts of weather prediction error could potentially be assessed using weather fore-
668 casts from the past as atmospheric forcing data inputs to the model. As with other DL
669 models, however, this work does not assume the forcings or the target observations to
670 be perfect. The Artificial Intelligence community has worked extensively with data “*in*
671 *the wild*”, i.e. large but low-quality datasets, and DL models appear to deliver good per-
672 formance even if there is significant noise (Izadinia et al., 2015; Stadelmann et al., 2018;
673 Huang et al., 2016). What *will* mislead models are systematic errors.

674 The MCD+N method is simple to implement, but a lot remains to be understood.
675 Although the two uncertainty terms were computed using very different methods and
676 our experiments show they measure different uncertainty sources, their high level of cor-
677 relation shows that they are not orthogonal, i.e. independent, quantities. Although per-
678 haps unsatisfying, the correlation is consistent with their definitions and the proposed
679 GP interpretation of network weight uncertainty (which was called the epistemic uncer-
680 tainty in KG17). For data-driven models, knowledge comes from training data. When
681 the training data has large amounts of noise, the knowledge of the model is negatively
682 impacted, as reflected by the network weight uncertainty. In other words, noise in train-
683 ing data makes the model less certain of its own predictions. To further complicate our
684 understanding, the correlation between network weight and data noise uncertainties also
685 reflects the overall pattern of moisture variation and SMAP accuracy as functions of an-
686 nual precipitation over the CONUS. Regions with high annual precipitation and high per-
687 centages of precipitation as snow also have high percentages of forest cover, and there-
688 fore high vegetation water content, which is known to lead to large uncertainty in SMAP
689 measurements. Other datasets without these associations could help to disentangle the

690 effects of these factors. Even entangled, however, these factors are good estimators of
691 prediction error and are thus still useful.

692 It could be hypothesized that the correlation between network weight and data noise
693 uncertainties will be lower if we have a much larger dataset, as the data quantity could
694 compensate for the quality, as shown in studies using noisy data “*in the wild*”. However,
695 as this is merely the first paper in hydrology to examine the MCD+N scheme, we leave
696 the testing of this hypothesis to future work with more data quantity and diversity.

697 Due to its data-driven nature, the data noise uncertainty estimate is still conditioned
698 by data, making it vulnerable to biased training data. This observation exposes an in-
699 herent limitation with any purely data-driven method, which is that it is difficult to as-
700 sess the quality of data based only on the data itself. Future integration of knowledge
701 or process-based models could potentially reduce this barrier. For example, process-based
702 models could be constructed to introduce physics relationships that were not adequately
703 represented in the training data. How to properly combine two classes of models is an
704 active area of research (Karpatne et al., 2017; Shen et al., 2018), and other methods such
705 as Stein variational gradient descent training (Liu & Wang, 2016; Mo et al., 2018) could
706 also be considered.

707 MCD seemed to have automatically identified similarities in the inputs (atmospheric
708 forcing data, soils, slope, land cover), which manifested as smaller network weight un-
709 certainties for neighboring regions. These similarities are not entirely based on geographic
710 proximity. Compared to geostatistical methods such as Kriging (a GP that parameter-
711 izes covariance functions over geographic distance), input-parameterized similarity fa-
712 cilitates physical interpretation and relieves us from the burden of identifying and tun-
713 ing appropriate forms and parameters of covariance functions. An immediate next step
714 could be to examine the most important physical input parameters that were employed
715 by the MCD dissimilarity detector, to determine whether the network has made a physically-
716 meaningful selection of attributes.

717 The theory behind the success of MCD needs further development, but this is one
718 intuitive explanation for how it works: A deep network is composed of neurons. Each
719 neuronal unit has inputs x_1, \dots, x_k , corresponding weights w_1, \dots, w_k , a bias term b , and
720 an activation function \mathbf{g} . The output of the unit is $\mathbf{g}(b + \sum_i x_i w_i)$. During training with
721 dropout, the neuron only uses a Bernoulli random sample of its inputs to create an out-

722 put, such that a random subset of the terms in the summation are removed. Thus the
723 unit is conditioned to produce approximately the same output from different subsets of
724 its input; otherwise training would not be stable. In other words, the neuronal unit learns
725 about redundancies in its inputs that occur during training, and takes advantage of them
726 so that different subsets of its inputs can produce approximately the same output. When
727 the testing data are not represented by the training data, the characteristics of the in-
728 puts to the neuronal unit change. The same types of redundancies that held in the train-
729 ing data would not be expected to hold in the testing data. Hence, the random summa-
730 tions would no longer result in similar outputs, causing an observable increase in vari-
731 ability. Future work could test this intuition and further improve the MCD formulation.
732 As a side note, this redundancy requirement would be a very powerful constraint, which
733 could ensure that a trained neural system produces robust outcomes.

734 Uncertainty estimation has long been a focus in hydrology and other domains. How-
735 ever, very often the quality of the uncertainty estimate has not been thoroughly eval-
736 uated. Our results show that there could be many subtleties and limitations with state-
737 of-the-art uncertainty estimates. For example, one could employ the MCD+N method
738 for a model to produce an uncertainty estimate for a new instance, without realizing the
739 limitations of the data noise term when this new instance is outside of the training data
740 distribution. More importantly, an improper uncertainty estimate could provide a false
741 sense of reliability. Therefore, we recommend carefully evaluating the uncertainty esti-
742 mate before applying it in a production setting.

743 **4 Conclusions**

744 Uncertainty estimation is an essential task for hydrology, but it is new for hydro-
745 logic time series deep learning. Our evaluation with soil moisture predictions shows that
746 MCD+N can indeed help to estimate model error. MCD+N proposed an input-dependent
747 data noise term and a network weight uncertainty term, which are new concepts for hy-
748 drology. While the two terms were correlated for a CONUS-scale model, our experiments
749 showed they indeed primarily targeted different uncertainty sources. The proposed data
750 noise term is essentially a data-driven error model that greatly simplifies error quanti-
751 fication, without the need for explicit assumptions. Most observational noise was correctly
752 attributed to the data noise term in our experiments. Additionally, our results provided
753 the first strong supporting evidence that Monte Carlo dropout does act as a dissimilar-

754 ity detector, while the data noise term does not. These *work-as-intended* behaviors gives
 755 us some confidence that MCD+N is a useful tool. However, uncertainty estimation is
 756 not a replacement for data acquisition. We showed that both terms are dependent on
 757 the training data. If the training data are not representative, not only will the error in-
 758 crease noticeably, but the quality of the data noise estimate may also deteriorate. For-
 759 tunately, we only need a small set of data covering the input space to serve as a repre-
 760 sentative training set. To improve the uncertainty quality, we should strive to include
 761 extreme cases in the training set. The MCD+N scheme had promise, but should not be
 762 used with blind trust.

763 Acknowledgments

764 All data used in this study, including forcing data from NLDAS-2 ¹, land surface char-
 765 acteristics (including soil texture from ISRIC-WISE ², land cover from NLCD ³, and NDVI
 766 from ⁴), and SMAP measurements, are available from public sources. The LSTM code
 767 can be openly downloaded from the open-source repository ⁵. KF was sponsored par-
 768 tially by the Biological and Environmental Research program from the U.S. Department
 769 of Energy under contract DE-SC0016605. CS was supported by the National Science Foun-
 770 dation under grant EAR #1832294, and a seed grant from the Penn State Institutes of
 771 Energy and the Environment.

772 References

- 773 Ajami, N. K., Hornberger, G. M., & Sunding, D. L. (2008, nov). Sustainable
 774 water resource management under hydrological uncertainty. *Water Re-*
 775 *sources Research*, 44(11). Retrieved from [http://doi.wiley.com/10.1029/](http://doi.wiley.com/10.1029/2007WR006736)
 776 [2007WR006736](http://doi.wiley.com/10.1029/2007WR006736) doi: 10.1029/2007WR006736
- 777 Batjes, N. H. (1995). *A Homogenized Soil Data File for Global Environmental Re-*
 778 *search: A Subset of FAO, ISRIC, and NRCS profiles (Version 1.0)* (Tech. Rep.
 779 No. No. 95/10b). Wageningen: ISRIC.
- 780 Behrouz, M., & Alimohammadi, S. (2018, aug). Uncertainty Analysis of

¹ https://hydro1.gesdisc.eosdis.nasa.gov/data/NLDAS/NLDAS_FORA0125_H.002/

² <https://www.isric.org/projects/world-inventory-soil-emission-potentials-wise>

³ <https://www.mrlc.gov/data/nlcd-2016-land-cover-conus>

⁴ <https://ecocast.arc.nasa.gov/data/pub/gimms/3g.v1/>

⁵ <https://github.com/mhpi/hydroDL>

- 781 Flood Control Measures Including Epistemic and Aleatory Uncertainties:
 782 Probability Theory and Evidence Theory. *Journal of Hydrologic Engi-*
 783 *neering*, 23(8), 04018033. Retrieved from [http://ascelibrary.org/](http://ascelibrary.org/doi/10.1061/(ASCE)HE.1943-5584.0001675)
 784 [doi/10.1061/](http://doi/10.1061/(ASCE)HE.1943-5584.0001675)
 785 [doi: 10.1061/](http://doi/10.1061/(ASCE)HE.1943-5584.0001675)
 (ASCE)HE.1943-5584.0001675
- 786 Bengio, Y. (2009). Learning Deep Architectures for AI. *Foundations and*
 787 *Trends® in Machine Learning*, 2(1), 1–127. Retrieved from [http://www](http://www.nowpublishers.com/article/Details/MAL-006)
 788 [.nowpublishers.com/article/Details/MAL-006](http://www.nowpublishers.com/article/Details/MAL-006) doi: 10.1561/2200000006
- 789 Beven, K. (1989, jan). Changing ideas in hydrology The case of physically-based
 790 models. *Journal of Hydrology*, 105(1-2), 157–172. Retrieved from [http://](http://linkinghub.elsevier.com/retrieve/pii/0022169489901017)
 791 linkinghub.elsevier.com/retrieve/pii/0022169489901017 doi: 10.1016/
 792 [0022-1694\(89\)90101-7](http://doi/10.1016/0022-1694(89)90101-7)
- 793 Beven, K. (2016, jul). Facets of uncertainty: epistemic uncertainty, non-
 794 stationarity, likelihood, hypothesis testing, and communication. *Hydro-*
 795 *logical Sciences Journal*, 61(9), 1652–1665. Retrieved from [http://](http://www.tandfonline.com/doi/full/10.1080/02626667.2015.1031761)
 796 www.tandfonline.com/doi/full/10.1080/02626667.2015.1031761 doi:
 797 [10.1080/02626667.2015.1031761](http://doi/10.1080/02626667.2015.1031761)
- 798 Butts, M. B., Payne, J. T., Kristensen, M., & Madsen, H. (2004, oct). An
 799 evaluation of the impact of model structure on hydrological modelling
 800 uncertainty for streamflow simulation. *Journal of Hydrology*, 298(1-4),
 801 242–266. Retrieved from [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S0022169404002471?casa={_}token=XRb1PqQkAFMAAAA:1tPq9pDTvI7kRav1465t0bC86fnWDkhYp53VC2kB-QbERZ5kAPpPaB6Uo3bM8J-X{_}awVEHd1AQ)
 802 [article/pii/S0022169404002471?casa={_}token=XRb1PqQkAFMAAAA:](https://www.sciencedirect.com/science/article/pii/S0022169404002471?casa={_}token=XRb1PqQkAFMAAAA:1tPq9pDTvI7kRav1465t0bC86fnWDkhYp53VC2kB-QbERZ5kAPpPaB6Uo3bM8J-X{_}awVEHd1AQ)
 803 [1tPq9pDTvI7kRav1465t0bC86fnWDkhYp53VC2kB-QbERZ5kAPpPaB6Uo3bM8J-X{\](https://www.sciencedirect.com/science/article/pii/S0022169404002471?casa={_}token=XRb1PqQkAFMAAAA:1tPq9pDTvI7kRav1465t0bC86fnWDkhYp53VC2kB-QbERZ5kAPpPaB6Uo3bM8J-X{_}awVEHd1AQ)
 804 [_}awVEHd1AQ](https://www.sciencedirect.com/science/article/pii/S0022169404002471?casa={_}token=XRb1PqQkAFMAAAA:1tPq9pDTvI7kRav1465t0bC86fnWDkhYp53VC2kB-QbERZ5kAPpPaB6Uo3bM8J-X{_}awVEHd1AQ) doi: 10.1016/J.JHYDROL.2004.03.042
- 805 De Lannoy, G. J. M., Reichle, R. H., Houser, P. R., Pauwels, V. R. N., & Ver-
 806 hoest, N. E. C. (2007, sep). Correcting for forecast bias in soil moisture
 807 assimilation with the ensemble Kalman filter. *Water Resources Research*,
 808 43(9). Retrieved from <http://doi.wiley.com/10.1029/2006WR005449> doi:
 809 [10.1029/2006WR005449](http://doi/10.1029/2006WR005449)
- 810 Depeweg, S., Hernández-Lobato, J. M., Doshi-Velez, F., & Udluft, S. (2017, oct).
 811 Decomposition of Uncertainty in Bayesian Deep Learning for Efficient and
 812 Risk-sensitive Learning. In *Proceedings of the 35 th international confer-*
 813 *ence on machine learning, stockholm, sweden, pmlr 80, 2018*. Retrieved from

- 814 <http://arxiv.org/abs/1710.07283>
- 815 Evin, G., Kavetski, D., Thyer, M., & Kuczera, G. (2013, jul). Pitfalls and im-
816 provements in the joint inference of heteroscedasticity and autocorrelation
817 in hydrological model calibration. *Water Resources Research*, *49*(7), 4518–
818 4524. Retrieved from <http://doi.wiley.com/10.1002/wrcr.20284> doi:
819 10.1002/wrcr.20284
- 820 Fang, K., Pan, M., & Shen, C. (2018). The Value of SMAP for Long-Term
821 Soil Moisture Estimation With the Help of Deep Learning. *IEEE Trans-*
822 *actions on Geoscience and Remote Sensing*, *PP*(DI), 1–13. Retrieved
823 from <https://ieeexplore.ieee.org/document/8497052/> doi: 10.1109/
824 TGRS.2018.2872131
- 825 Fang, K., & Shen, C. (2020, jan). Near-real-time forecast of satellite-based soil
826 moisture using long short-term memory with an adaptive data integration
827 kernel. *Journal of Hydrometeorology*, *JHM-D-19-0169.1*. Retrieved from
828 <http://journals.ametsoc.org/doi/10.1175/JHM-D-19-0169.1> doi:
829 10.1175/JHM-D-19-0169.1
- 830 Fang, K., Shen, C., Kifer, D., & Yang, X. (2017, nov). Prolongation of SMAP
831 to Spatiotemporally Seamless Coverage of Continental U.S. Using a Deep
832 Learning Neural Network. *Geophysical Research Letters*, *44*(21), 11,030–
833 11,039. Retrieved from <http://doi.wiley.com/10.1002/2017GL075619> doi:
834 10.1002/2017GL075619
- 835 Feng, D., Fang, K., & Shen, C. (2019). Enhancing streamflow forecast and extract-
836 ing insights using long short term memory networks that assimilate recent
837 observations. <https://arxiv.org/abs/1912.08949>.
- 838 Gal, Y., & Ghahramani, Z. (2016, jun). Dropout as a Bayesian Approximation:
839 Representing Model Uncertainty in Deep Learning. *Proceedings of The*
840 *33rd International Conference on Machine Learning*, *48*. Retrieved from
841 <http://arxiv.org/abs/1506.02142>
- 842 Glorot, X., Bordes, A., & Bengio, Y. (2011). Deep Sparse Rectifier Neural Net-
843 works. In G. Gordon, D. Dunson, & M. Dudík (Eds.), *Proceedings of the*
844 *fourteenth international conference on artificial intelligence and statistics* (pp.
845 315–323). PMLR. Retrieved from [http://proceedings.mlr.press/v15/
846 glorot11a.html](http://proceedings.mlr.press/v15/glorot11a.html)

- 847 Gong, W., Gupta, H. V., Yang, D., Sricharan, K., & Hero, A. O. (2013). Esti-
 848 mating epistemic and aleatory uncertainties during hydrologic modeling: An
 849 information theoretic approach. *Water Resources Research*, *49*(4), 2253–2273.
 850 Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/full/](https://agupubs.onlinelibrary.wiley.com/doi/full/10.1002/wrcr.20161)
 851 [10.1002/wrcr.20161](https://doi.org/10.1002/wrcr.20161) doi: 10.1002/wrcr.20161
- 852 Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. The MIT Press.
 853 Retrieved from <https://www.deeplearningbook.org/>
- 854 Götzing, J., & Bárdossy, A. (2008, dec). Generic error model for calibration and
 855 uncertainty estimation of hydrological models. *Water Resources Research*,
 856 *44*(12). Retrieved from <http://doi.wiley.com/10.1029/2007WR006691> doi:
 857 [10.1029/2007WR006691](https://doi.org/10.1029/2007WR006691)
- 858 Graves, A. (2013, aug). Generating Sequences With Recurrent Neural Networks.
 859 *arXiv:1308.0850*. Retrieved from <http://arxiv.org/abs/1308.0850>
- 860 Graves, A., Mohamed, A.-r., & Hinton, G. (2013, may). Speech recognition with
 861 deep recurrent neural networks. In *2013 IEEE International Conference on Acous-*
 862 *tics, Speech and Signal Processing* (pp. 6645–6649). Vancouver, Canada: IEEE.
 863 Retrieved from <http://ieeexplore.ieee.org/document/6638947/> doi: 10
 864 [.1109/ICASSP.2013.6638947](https://doi.org/10.1109/ICASSP.2013.6638947)
- 865 Hochreiter, S., & Schmidhuber, J. (1997, nov). Long Short-Term Mem-
 866 ory. *Neural Computation*, *9*(8), 1735–1780. Retrieved from [http://](http://www.mitpressjournals.org/doi/10.1162/neco.1997.9.8.1735)
 867 www.mitpressjournals.org/doi/10.1162/neco.1997.9.8.1735 doi:
 868 [10.1162/neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735)
- 869 Hornik, K. (1991). Approximation capabilities of multilayer feedforward
 870 networks. *Neural Networks*, *4*(2), 251–257. Retrieved from [http://](http://linkinghub.elsevier.com/retrieve/pii/089360809190009T)
 871 linkinghub.elsevier.com/retrieve/pii/089360809190009T doi:
 872 [10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T)
- 873 Huang, W., He, D., Yang, X., Zhou, Z., Kifer, D., & Giles, C. L. (2016). De-
 874 tecting Arbitrary Oriented Text in the Wild with a Visual Attention
 875 Model. In *Proceedings of the 2016 ACM on multimedia conference - mm*
 876 *'16* (pp. 551–555). New York, New York, USA: ACM Press. Retrieved
 877 from <http://dl.acm.org/citation.cfm?doid=2964284.2967282> doi:
 878 [10.1145/2964284.2967282](https://doi.org/10.1145/2964284.2967282)
- 879 Izadinia, H., Russell, B. C., Farhadi, A., Hoffman, M. D., & Hertzmann, A. (2015).

- 880 Deep Classifiers from Image Tags in the Wild. In *Proceedings of the 2015*
 881 *workshop on community-organized multimodal mining: Opportunities for novel*
 882 *solutions* (pp. 13–18). ACM. doi: 10.1145/2814815.2814821
- 883 Karpatne, A., Atluri, G., Faghmous, J. H., Steinbach, M., Banerjee, A., Ganguly, A.,
 884 ... Kumar, V. (2017, oct). Theory-Guided Data Science: A New Paradigm
 885 for Scientific Discovery from Data. *IEEE Transactions on Knowledge and Data*
 886 *Engineering*, 29(10), 2318–2331. Retrieved from [http://ieeexplore.ieee](http://ieeexplore.ieee.org/document/7959606/)
 887 [.org/document/7959606/](http://ieeexplore.ieee.org/document/7959606/) doi: 10.1109/TKDE.2017.2720168
- 888 Kavetski, D., Kuczera, G., & Franks, S. W. (2006, mar). Bayesian analysis of
 889 input uncertainty in hydrological modeling: 2. Application. *Water Re-*
 890 *sources Research*, 42(3). Retrieved from [http://doi.wiley.com/10.1029/](http://doi.wiley.com/10.1029/2005WR004376)
 891 [2005WR004376](http://doi.wiley.com/10.1029/2005WR004376) doi: 10.1029/2005WR004376
- 892 Kendall, A., & Gal, Y. (2017, mar). What Uncertainties Do We Need in Bayesian
 893 Deep Learning for Computer Vision? *Advances in Neural Information Process-*
 894 *ing Systems* 30, 16(4), 5574—5584. Retrieved from [http://arxiv.org/abs/](http://arxiv.org/abs/1703.04977)
 895 [1703.04977](http://arxiv.org/abs/1703.04977)
- 896 Kiureghian, A. D., & Ditlevsen, O. (2009, mar). Aleatory or epistemic? Does
 897 it matter? *Structural Safety*, 31(2), 105–112. Retrieved from [https://](https://www.sciencedirect.com/science/article/pii/S0167473008000556)
 898 www.sciencedirect.com/science/article/pii/S0167473008000556 doi:
 899 [10.1016/J.STRUSAFE.2008.06.020](https://www.sciencedirect.com/science/article/pii/S0167473008000556)
- 900 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall-
 901 runoff modelling using Long Short-Term Memory (LSTM) networks. *Hydrolog-*
 902 *ogy and Earth System Sciences*, 22(11), 6005–6022. Retrieved from [https://](https://www.hydrol-earth-syst-sci.net/22/6005/2018/)
 903 www.hydrol-earth-syst-sci.net/22/6005/2018/ doi: 10.5194/hess-22-6005
 904 -2018
- 905 Kumar, S., Sharma, A., & Tsunoda, T. (2019, dec). Brain wave classification us-
 906 ing long short-term memory network based OPTICAL predictor. *Scientific Re-*
 907 *ports*, 9(1). doi: 10.1038/s41598-019-45605-1
- 908 Lamontagne, J. R., Reed, P. M., Link, R., Calvin, K. V., Clarke, L. E., & Edmonds,
 909 J. A. (2018, mar). Large Ensemble Analytic Framework for Consequence-
 910 Driven Discovery of Climate Change Scenarios. *Earth's Future*, 6(3), 488–
 911 504. Retrieved from <http://doi.wiley.com/10.1002/2017EF000701> doi:
 912 [10.1002/2017EF000701](http://doi.wiley.com/10.1002/2017EF000701)

- 913 LeCun, Y., Bengio, Y., & Hinton, G. (2015, may). Deep learning. *Nature*,
 914 521(7553), 436–444. Retrieved from <http://www.nature.com/articles/nature14539> doi: 10.1038/nature14539
- 916 Lee, J., Bahri, Y., Novak, R., Schoenholz, S. S., Pennington, J., & Sohl-Dickstein,
 917 J. (2018). Deep Neural Networks as Gaussian Processes. In *International
 918 conference on learning representations*. Retrieved from <http://arxiv.org/abs/1711.00165>
- 920 Liu, Q., & Wang, D. (2016). Stein variational gradient descent: A general pur-
 921 pose bayesian inference algorithm. In *Advances in neural information process-
 922 ing systems* (pp. 2378–2386).
- 923 Matthews, A. G. d. G., Rowland, M., Hron, J., Turner, R. E., & Ghahramani, Z.
 924 (2018). Gaussian Process Behaviour in Wide Deep Neural Networks. In *In-
 925 ternational conference on learning representations* (pp. 1–36). Retrieved from
 926 <http://arxiv.org/abs/1804.11271>
- 927 McMahon, G., Gregonis, S. M., Waltman, S. W., Omernik, J. M., Thorson, T. D.,
 928 Ffreehouf, J. A., . . . Keys, J. E. (2001, apr). Developing a Spatial Frame-
 929 work of Common Ecological Regions for the Conterminous United States.
 930 *Environmental Management*, 28(3), 293–316. Retrieved from [http://](http://link.springer.com/10.1007/s0026702429)
 931 link.springer.com/10.1007/s0026702429 doi: 10.1007/s0026702429
- 932 Miotto, R., Wang, F., Wang, S., Jiang, X., & Dudley, J. T. (2017, may). Deep learn-
 933 ing for healthcare: Review, opportunities and challenges. *Briefings in Bioinfor-
 934 matics*, 19(6), 1236–1246. doi: 10.1093/bib/bbx044
- 935 Mo, S., Zhu, Y., Zabararas, J., Nicholas, Shi, X., & Wu, J. (2018). Deep convolutional
 936 encoder-decoder networks for uncertainty quantification of dynamic multiphase
 937 flow in heterogeneous media. *Water Resources Research*.
- 938 Neal, R. M. (1996). Priors for Infinite Networks. In *Bayesian learning for neu-
 939 ral networks* (pp. 29–53). New York, NY: Springer New York. Retrieved from
 940 http://link.springer.com/10.1007/978-1-4612-0745-0_2 doi: 10.1007/
 941 978-1-4612-0745-0_2
- 942 Nearing, G. S., Mocko, D. M., Peters-Lidard, C. D., Kumar, S. V., & Xia, Y. (2016,
 943 mar). Benchmarking NLDAS-2 Soil Moisture and Evapotranspiration to
 944 Separate Uncertainty Contributions. *Journal of Hydrometeorology*, 17(3),
 945 745–759. Retrieved from <http://journals.ametsoc.org/doi/10.1175/>

- 946 JHM-D-15-0063.1 doi: 10.1175/JHM-D-15-0063.1
- 947 Nearing, G. S., Tian, Y., Gupta, H. V., Clark, M. P., Harrison, K. W., & Weijs,
948 S. V. (2016, jul). A philosophical basis for hydrological uncertainty. *Hydrolog-
949 ical Sciences Journal*, 61(9), 1666–1678. Retrieved from [http://dx.doi.org/
950 10.1080/02626667.2016.1183009](http://dx.doi.org/10.1080/02626667.2016.1183009) doi: 10.1080/02626667.2016.1183009
- 951 O’Neill, P., Chan, S., Njoku, E., Jackson, T., & Bindlish, R. (2016). *SMAP L3
952 Radiometer Global Daily 36 km EASE-Grid Soil Moisture, Version 4*. Boul-
953 der, Colorado USA: NASA National Snow and Ice Data Center Distributed
954 Active Archive Center. Retrieved from [https://nsidc.org/data/SPL3SMP/
955 versions/4](https://nsidc.org/data/SPL3SMP/versions/4) doi: 10.5067/OBBHQ5W22HME
- 956 Osband, I., Blundell, C., Pritzel, A., & Van Roy, B. (2016, feb). Deep Exploration
957 via Bootstrapped DQN. *NIPS 2016 Bayesian Deep Learning Workshop*, 26–28.
958 Retrieved from <http://arxiv.org/abs/1602.04621>
- 959 Pan, M., Cai, X., Chaney, N. W., Entekhabi, D., & Wood, E. F. (2016, sep). An
960 initial assessment of SMAP soil moisture retrievals using high-resolution model
961 simulations and in situ observations. *Geophysical Research Letters*, 43(18),
962 9662–9668. Retrieved from [http://doi.wiley.com/10.1002/2016GL069964
963 doi: 10.1002/2016GL069964](http://doi.wiley.com/10.1002/2016GL069964)
- 964 Pappenberger, F., & Beven, K. J. (2006, may). Ignorance is bliss: Or seven reasons
965 not to use uncertainty analysis. *Water Resources Research*, 42(5), 1–8. Re-
966 trieved from <http://doi.wiley.com/10.1029/2005WR004820> doi: 10.1029/
967 2005WR004820
- 968 Rasmussen, C. E., & Williams, C. K. I. (2005). *Gaussian Processes for Machine
969 Learning (Adaptive Computation and Machine Learning Series)*. The MIT
970 Press.
- 971 Schmidhuber, J. (2015, jan). Deep learning in neural networks: An overview. *Neu-
972 ral Networks*, 61(10), 85–117. Retrieved from [https://linkinghub.elsevier
973 .com/retrieve/pii/S0893608014002135](https://linkinghub.elsevier.com/retrieve/pii/S0893608014002135) doi: 10.1016/j.neunet.2014.09.003
- 974 Senge, R., Bösner, S., Dembczyński, K., Haasenritter, J., Hirsch, O., Donner-
975 Banzhoff, N., & Hüllermeier, E. (2014, jan). Reliable classification: Learning
976 classifiers that distinguish aleatoric and epistemic uncertainty. *Information
977 Sciences*, 255, 16–29. Retrieved from [https://www.sciencedirect.com/
978 science/article/pii/S0020025513005410](https://www.sciencedirect.com/science/article/pii/S0020025513005410) doi: 10.1016/J.INS.2013.07.030

- 979 Shen, C., Laloy, E., Elshorbagy, A., Albert, A., Bales, J., Chang, F.-j., ... Tsai,
 980 W.-P. (2018, nov). HESS Opinions: Incubating deep-learning-powered hydro-
 981 logic science advances as a community. *Hydrology and Earth System Sciences*,
 982 22(11), 5639–5656. Retrieved from [https://www.hydrol-earth-syst-sci](https://www.hydrol-earth-syst-sci.net/22/5639/2018/)
 983 [.net/22/5639/2018/](https://www.hydrol-earth-syst-sci.net/22/5639/2018/) doi: 10.5194/hess-22-5639-2018
- 984 Smith, T., Marshall, L., & Sharma, A. (2015, sep). Modeling
 985 residual hydrologic errors with Bayesian inference. *Journal of Hydrology*, 528,
 986 29–37. Retrieved from [https://www](https://www.sciencedirect.com/science/article/pii/S0022169415004011?casa={_}token=qzNGTPKWZooAAAAA:jyFlvInezUh780kFBSn-7tVwF9PdX2EERKiKT0lpjEKrKyfBIEhPqcQdge1Tzb3gC7tJz0uZJQ)
 987 [.sciencedirect.com/science/article/pii/S0022169415004011](https://www.sciencedirect.com/science/article/pii/S0022169415004011?casa={_}token=qzNGTPKWZooAAAAA:jyFlvInezUh780kFBSn-7tVwF9PdX2EERKiKT0lpjEKrKyfBIEhPqcQdge1Tzb3gC7tJz0uZJQ)
 988 [?casa={_}token=qzNGTPKWZooAAAAA:jyFlvInezUh780kFBSn](https://www.sciencedirect.com/science/article/pii/S0022169415004011?casa={_}token=qzNGTPKWZooAAAAA:jyFlvInezUh780kFBSn-7tVwF9PdX2EERKiKT0lpjEKrKyfBIEhPqcQdge1Tzb3gC7tJz0uZJQ)
 989 [-7tVwF9PdX2EERKiKT0lpjEKrKyfBIEhPqcQdge1Tzb3gC7tJz0uZJQ](https://www.sciencedirect.com/science/article/pii/S0022169415004011?casa={_}token=qzNGTPKWZooAAAAA:jyFlvInezUh780kFBSn-7tVwF9PdX2EERKiKT0lpjEKrKyfBIEhPqcQdge1Tzb3gC7tJz0uZJQ) doi:
 990 10.1016/J.JHYDROL.2015.05.051
- 991 Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., & Salakhutdinov,
 992 R. (2014). Dropout: A Simple Way to Prevent Neural Networks from
 993 Overfitting. *Journal of Machine Learning Research*, 15, 1929–1958. Re-
 994 trieved from <http://jmlr.org/papers/v15/srivastava14a.html> doi:
 995 10.1214/12-AOS1000
- 996 Stadelmann, T., Amirian, M., Arabaci, I., Arnold, M., Duivesteijn, G. F., Elezi, I.,
 997 ... Tuggener, L. (2018, sep). Deep Learning in the Wild. In *Annpr 2018*
 998 (pp. 17–38). Springer, Cham. Retrieved from [http://link.springer.com/](http://link.springer.com/10.1007/978-3-319-99978-4{_}2)
 999 [10.1007/978-3-319-99978-4{_}2](http://link.springer.com/10.1007/978-3-319-99978-4{_}2) doi: 10.1007/978-3-319-99978-4_2
- 1000 Vandal, T., Kodra, E., Dy, J., Ganguly, S., Nemani, R., & Ganguly, A. R. (2018).
 1001 Quantifying Uncertainty in Discrete-Continuous and Skewed Data with
 1002 Bayesian Deep Learning. In *Proceedings of the 24th acm sigkdd international*
 1003 *conference on knowledge discovery & data mining - kdd '18* (pp. 2377–2386).
 1004 New York, New York, USA: ACM Press. Retrieved from [http://dx.doi.org/](http://dx.doi.org/10.1145/3219819.3219996)
 1005 [10.1145/3219819.3219996](http://dx.doi.org/10.1145/3219819.3219996) doi: 10.1145/3219819.3219996
- 1006 Vrugt, J. A., ter Braak, C. J. F., Clark, M. P., Hyman, J. M., & Robinson, B. A.
 1007 (2008, dec). Treatment of input uncertainty in hydrologic modeling: Doing
 1008 hydrology backward with Markov chain Monte Carlo simulation. *Water Re-*
 1009 *sources Research*, 44(12), 1–15. Retrieved from [http://doi.wiley.com/](http://doi.wiley.com/10.1029/2007WR006720)
 1010 [10.1029/2007WR006720](http://doi.wiley.com/10.1029/2007WR006720) doi: 10.1029/2007WR006720
- 1011 Wilson, T., Tan, P.-n., & Luo, L. (2018, nov). A Low Rank Weighted Graph Con-

- 1012 volutional Approach to Weather Prediction. In *2018 IEEE International Confer-*
1013 *ence on Data Mining (ICDM)* (pp. 627–636). IEEE. Retrieved from [https://](https://ieeexplore.ieee.org/document/8594887/)
1014 ieeexplore.ieee.org/document/8594887/ doi: 10.1109/ICDM.2018.00078
- 1015 Xia, Y., Ek, M. B., Wu, Y., Ford, T., & Quiring, S. M. (2015, oct). Compari-
1016 son of NLDAS-2 Simulated and NASMD Observed Daily Soil Moisture. Part I:
1017 Comparison and Analysis. *Journal of Hydrometeorology*, *16*(5), 1962–1980. Re-
1018 trieved from <http://journals.ametsoc.org/doi/10.1175/JHM-D-14-0096.1>
1019 doi: 10.1175/JHM-D-14-0096.1
- 1020 Zhang, D., Lindholm, G., & Ratnaweera, H. (2018, jan). Use long short-
1021 term memory to enhance Internet of Things for combined sewer over-
1022 flow monitoring. *Journal of Hydrology*, *556*, 409–418. Retrieved from
1023 <https://linkinghub.elsevier.com/retrieve/pii/S0022169417307722>
1024 doi: 10.1016/j.jhydrol.2017.11.018
- 1025 Zhang, J., Zhu, Y., Zhang, X., Ye, M., & Yang, J. (2018, jun). Developing a Long
1026 Short-Term Memory (LSTM) based model for predicting water table depth in
1027 agricultural areas. *Journal of Hydrology*, *561* (April), 918–929. Retrieved from
1028 <https://linkinghub.elsevier.com/retrieve/pii/S0022169418303184>
1029 doi: 10.1016/j.jhydrol.2018.04.065

1030 **A The MCD theory and its potential issues**

1031 The derivations from GG16 (Gal & Ghahramani, 2016) are quite lengthy, so here
1032 we only highlight a few main steps. The prototype network analyzed is a two-layer net-
1033 work written as $\mathbf{f} = \sigma(\mathbf{x}\mathbf{W}^{(1)} + \mathbf{b})\mathbf{W}^{(2)}$, where σ is a nonlinear activation function such
1034 as TanH or ReLU and $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$ are the weights for the first and second layers,
1035 respectively. Adding in dropout operators, we obtain $\mathbf{f} = (\sigma(\mathbf{x}(\mathbf{z}^{(1)}\mathbf{W}^{(1)} + \mathbf{b}))(\mathbf{z}^{(2)}\mathbf{W}^{(2)}))$,
1036 where $z^{(1)} \sim \text{Bernoulli}(\beta^{(1)})$ and $z^{(2)} \sim \text{Bernoulli}(\beta^{(2)})$ are dropout masks of the same
1037 sizes as $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$, respectively. $\beta^{(k)}$ is the probability that a connection on
1038 the k -th layer is retained during dropout, or one minus the “dropout rate” in many DL
1039 packages. Hence we refer to it as the dropout retention rate.

1040 In a standard Bayesian inference framework, we (i) start with a prior distribution
1041 of model parameters, e.g. $p(\mathbf{W}) = \mathcal{N}(0, I)$; (ii) confront the model with the data (eval-
1042 uating the likelihood function) and calculate the posterior distribution of the parame-
1043 ter sets using Bayes law (i.e. given the training dataset (\mathbf{X}, \mathbf{Y}) , $p(\mathbf{W}|\mathbf{X}, \mathbf{Y}) =$

1044 $p(\mathbf{Y}|\mathbf{W}, \mathbf{X})p(\mathbf{W})/p(\mathbf{Y}|\mathbf{X})$); and (iii) use the posterior distribution to make predictions
 1045 as well as estimate predictive uncertainty for new test instances X^* :

$$p(\mathbf{Y}^*|\mathbf{X}^*) = \int \mathbf{p}(\mathbf{Y}^*|\mathbf{X}^*, \mathbf{W})\mathbf{p}(\mathbf{W}|\mathbf{X}, \mathbf{Y})dW \quad (\text{A.1})$$

1046 The posterior distribution $p(\mathbf{W}|\mathbf{X}, \mathbf{Y})$ is the distribution that most likely gener-
 1047 ated the observed data. However, this distribution cannot be easily estimated as the marginal
 1048 distribution $p(\mathbf{Y}|\mathbf{X})$ cannot be evaluated analytically, and is intractable for very high-
 1049 dimensional deep networks. A viable approach is to replace this distribution with a *vari-*
 1050 *ational* distribution $q(W)$, whose structure is easier to work with in the integral. Vari-
 1051 ational inference turns the inference problem into an optimization problem, where we
 1052 minimize the Kullback-Leibler divergence between the variational distribution and the
 1053 posterior distribution, $\mathbf{KL}(q(\mathbf{W})||p(\mathbf{W}|\mathbf{X}, \mathbf{Y}))$, which measures the dissimilarity between
 1054 distributions. Typically, this task is further turned into the problem of maximizing the
 1055 *log evidence lower bound* (LELB)

$$\mathcal{L} = \int q(W) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{W})d\omega - \mathbf{KL}(q(\mathbf{W})||p(\mathbf{W}|\mathbf{X}, \mathbf{Y})) \quad (\text{A.2})$$

1056 This procedure optimizes both the weights of the neural network and the varia-
 1057 tional parameters. As a result, after we solve this minimization problem we will have ob-
 1058 tained both a functional neural network and a variational distribution that can be eas-
 1059 ily sampled from. In the case of GG16, the authors would like to prove that dropout train-
 1060 ing corresponds to *some* form of variational distribution. They defined their variational
 1061 distributions for the weights of layer 1, $\mathbf{W}^{(1)}$, as a Gaussian mixture which can be fac-
 1062 torized over each row vector:

$$q(\mathbf{W}^{(1)}) = \prod_{q=1}^Q q(\mathbf{w}_q) \quad (\text{A.3})$$

$$q(\mathbf{w}_q) = \beta^{(1)}\mathcal{N}(\mathbf{m}_q, \sigma^2\mathbf{I}_K) + (1 - \beta^{(1)})\mathcal{N}(0, \sigma^2\mathbf{I}_K) \quad (\text{A.4})$$

1063 where $\mathbf{W}^{(1)}$ is of the size $Q \times K$ and \mathbf{w}_q is a row vector in $\mathbf{W}^{(1)}$. Similar distri-
 1064 butions were put on $\mathbf{W}^{(2)}$. This variational distribution can further be re-parameterized
 1065 as the following

$$\mathbf{W}^{(1)} = z^{(1)}(M^{(1)} + \sigma\epsilon^{(1)}) + (1 - z^{(1)})\sigma\epsilon^{(1)} \quad (\text{A.5})$$

$$\mathbf{W}^{(2)} = z^{(2)}(M^{(2)} + \sigma\epsilon^{(2)}) + (1 - z^{(2)})\sigma\epsilon^{(2)} \quad (\text{A.6})$$

$$\mathbf{b} = \mathbf{m} + \sigma\epsilon \quad (\text{A.7})$$

1066 The parameterization allows the integral in Eq. A.2 to be estimated using Monte
1067 Carlo integration, i.e.,

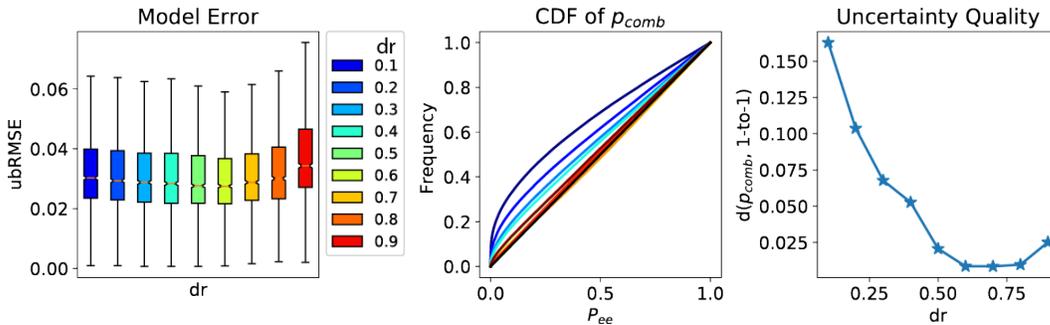
$$\mathcal{L}_{GP-MC} = \sum_{m=1}^M \log p(\mathbf{y}_m | \mathbf{x}_m, \widehat{\mathbf{W}}_m^{(1)}, \widehat{\mathbf{W}}_m^{(2)}, \widehat{\mathbf{b}}_m) - KL(q(\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}) || p(\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b})) \quad (\text{A.8})$$

1068 where $\widehat{W}_n^{(1)}$, $\widehat{W}_n^{(2)}$, and \widehat{b}_n are the weights for the n -th realization. GG16 argued
1069 that when σ is small, we simply have $\widehat{\mathbf{W}}^{(1)} \approx \widehat{\mathbf{z}}_n^{(1)} \mathbf{M}^{(1)}$, $\widehat{\mathbf{W}}^{(2)} \approx \widehat{\mathbf{z}}_n^{(2)} \mathbf{M}^{(2)}$, $\widehat{\mathbf{b}}^n \approx \mathbf{m}$. In
1070 other words, applying a stochastic dropout mask on the weights is approximately draw-
1071 ing a sample from the variational distribution in Eq. A.7, and the summation term sim-
1072 ply amounts to the sum of squared loss for training with dropout and mini-batching. Some
1073 other approximations that take advantage of the large size of deep networks were fur-
1074 ther employed to handle the KL term. Furthermore, by stacking more layers, the same
1075 derivation was extended to multi-layer networks.

1076 While it is fortunate that such an interpretation for dropout could exist, there were
1077 many approximate steps in this derivation. In particular, we have the following concerns:
1078 (i) the Bernoulli distribution and the Gaussian mixture that it approximates might not
1079 be competent enough as a variational distribution. The Gaussian mixture itself, as shown
1080 in the derivation, must have small variances, and it is uncertain if such strong limita-
1081 tions are valid for Bayesian inference; (ii) the Gaussian prior over the parameters $W \sim$
1082 $\mathcal{N}(1, I)$ is coincidental but not necessarily optimal; (iii) with many approximations stacked
1083 up in the derivation, it is dubious if the conclusion still converges to the declared final
1084 outcome; and (iv) the derivation was only demonstrated for simple multi-layer neural
1085 networks. This derivation has yet to be shown to work for complex recurrent networks
1086 like LSTM. It is not certain if LSTM with dropout training is a deep GP. While these
1087 concerns are difficult to address analytically at the moment, we can experimentally ver-
1088 ify the effectiveness of MCD and answer the research questions presented at the end of
1089 the Introduction section.

1090 B Calibration of dropout rate

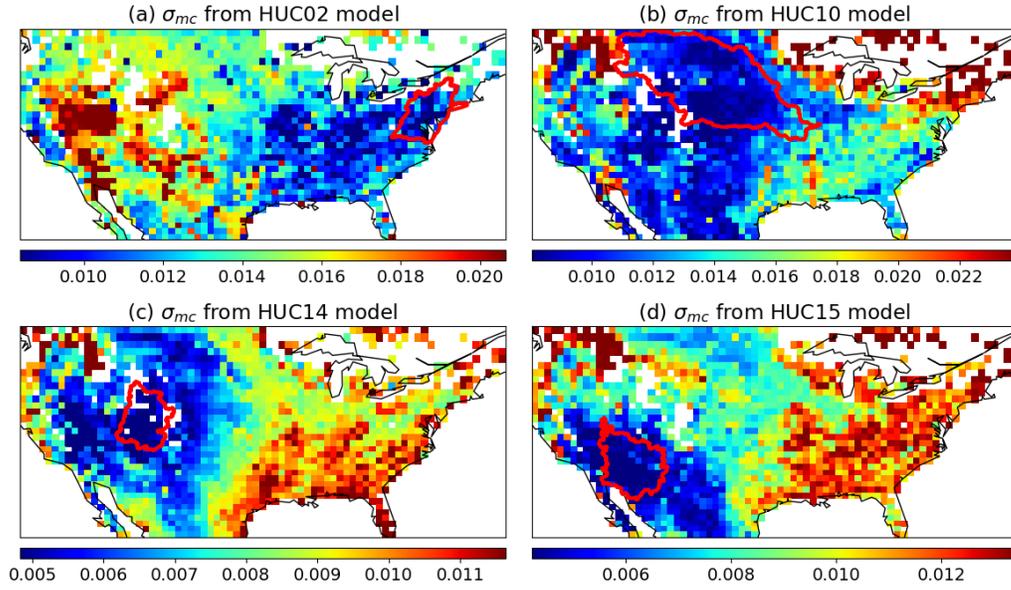
1091 Here we examine the role that dropout retention rate (β) plays in the uncertainty
 1092 estimate terms and the predictive error. In the MCD theory, the variational distribu-
 1093 tion for the parameters are Gaussian mixtures with very small variances, and the weights
 1094 before them are from a Bernoulli distribution (Appendix A). The dropout rate ($dr =$
 1095 $1 - \beta$) should be carefully calibrated. We trained the model from 2015/04 to 2016/03
 1096 using $\beta \in \{0.1, 0.2, \dots, 0.9\}$. The best β was chosen based on both the error and qual-
 1097 ity of the uncertainty estimate in the validation set (2016/04 - 2017/03). As figure B.1
 1098 shows, both $ubRMSE$ and σ_{comb} are affected by the dropout rate. We chose the model
 1099 trained with $dr = 0.6$, or $\beta = 0.4$, as it simultaneously gave the smallest $ubRMSE$
 1100 and the best uncertainty quality, as measured by d , the Kolmogorov-Smirnov statistic
 1101 (maximum distance) between the CDF of the error exceedance likelihoods and the one-
 1102 to-one line.



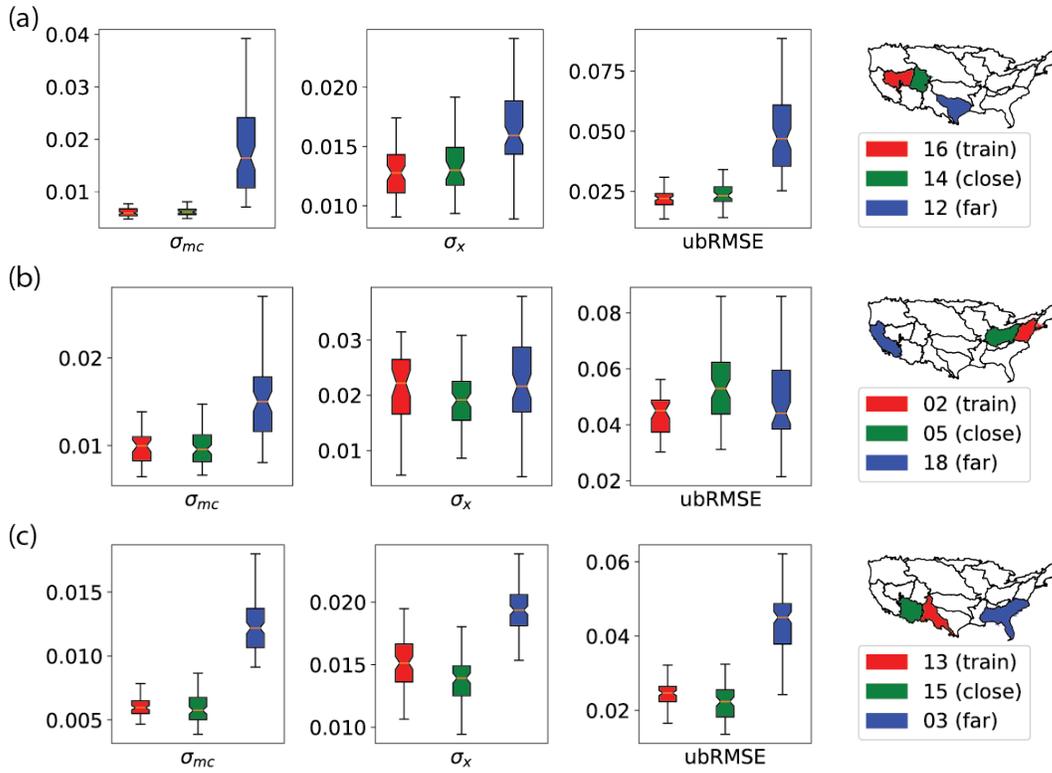
1103 **Figure B.1.** Performance of uncertainty models with different dropout rates ($dr = 1 - \beta$). (a)
 1104 $ubRMSE$ as a function of dr . (b) The CDF curves of the error exceedance likelihoods. (c) The
 1105 Kolmogorov-Smirnov statistic as a function of dr . We found that $dr = 0.6$ offers a balance of
 1106 small d as well as small $ubRMSE$.

1107 C Test on hydrologic basins instead of ecoregions.

1108 In practice, hydrologic models are commonly developed based on basins instead of
 1109 ecoregions. Hence, to provide more insights, we trained models on each of the 18 2-digits
 1110 hydrologic cataloging unit (HUC02) basins dividing CONUS. Similar to the ecoregion
 1111 experiments, the models were trained over year 2015, validated over 2016 and tested over
 1112 2017. We reproduced the figure 4 and 5 as C.1 and C.2 correspondingly, and they re-
 1113 vealed similar pattern as we discussed in section 3.3.



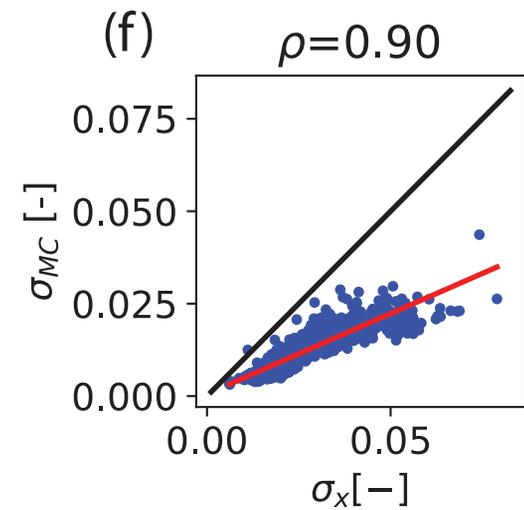
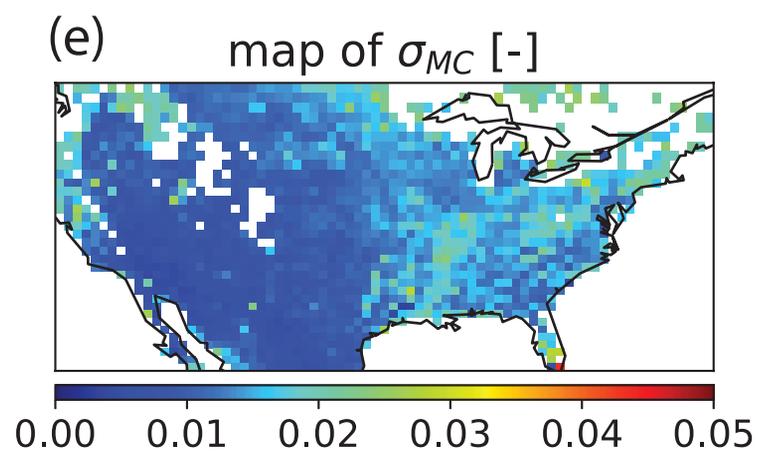
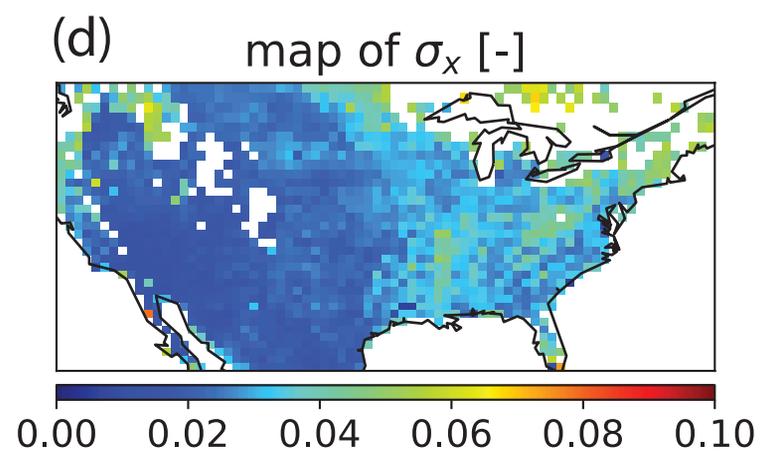
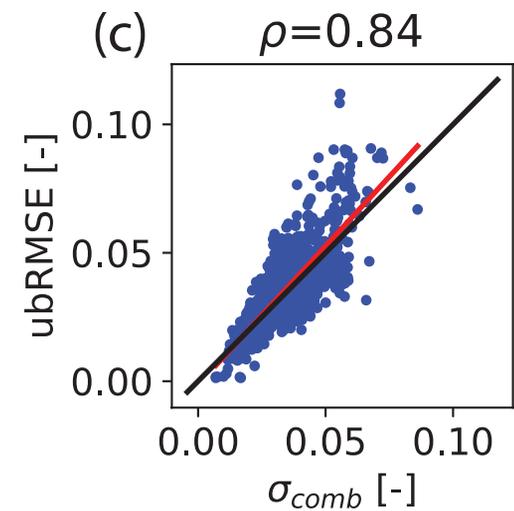
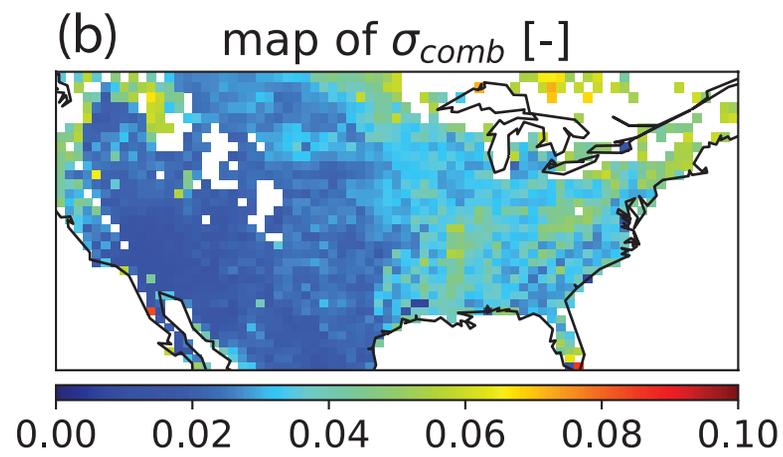
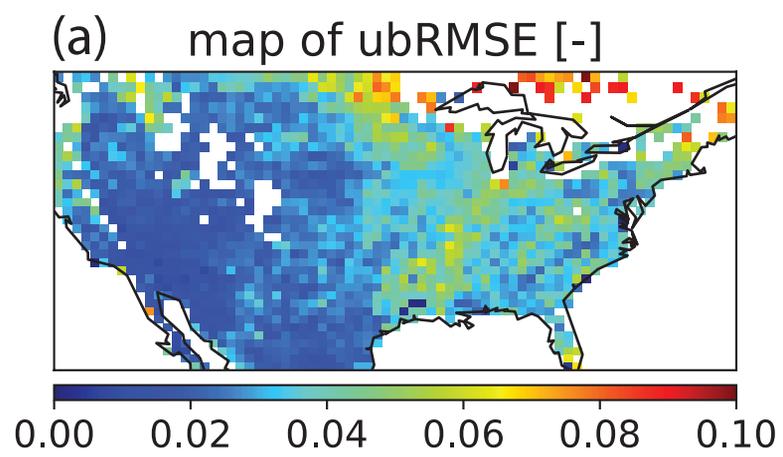
1114 **Figure C.1.** Maps of σ_{mc} when the LSTM model is trained in one of the HUC2 basins. The
 1115 training region is highlighted by the red polygon.



1116 **Figure C.2.** Metrics of performance when we trained the model in a HUC2, and tested in two
 1117 other HUC2s: one similar to the training region, one farther away, in a different physiographic
 1118 region.

Figure 1.

Temporal Test



Spatial Test

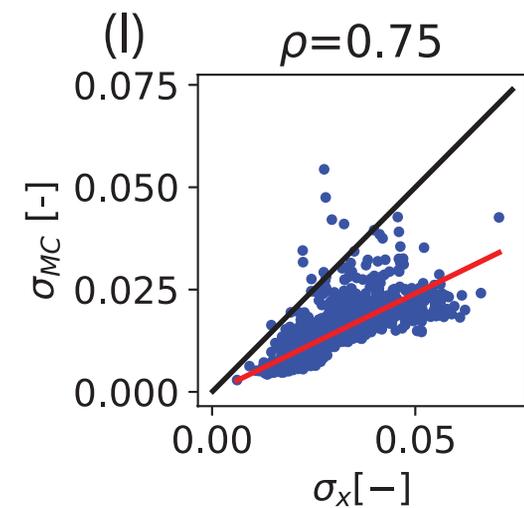
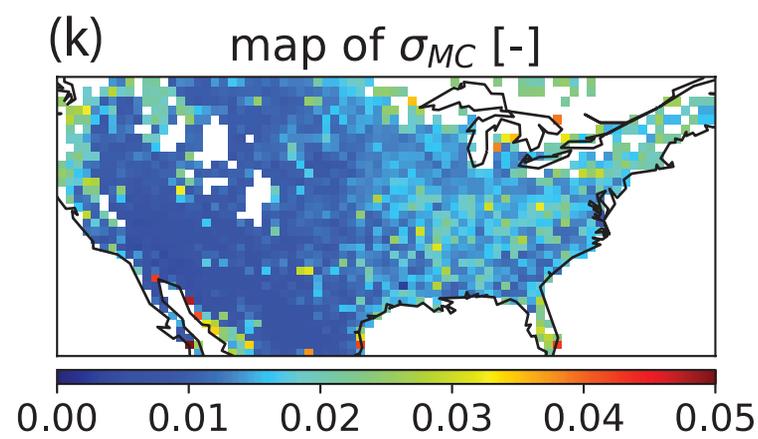
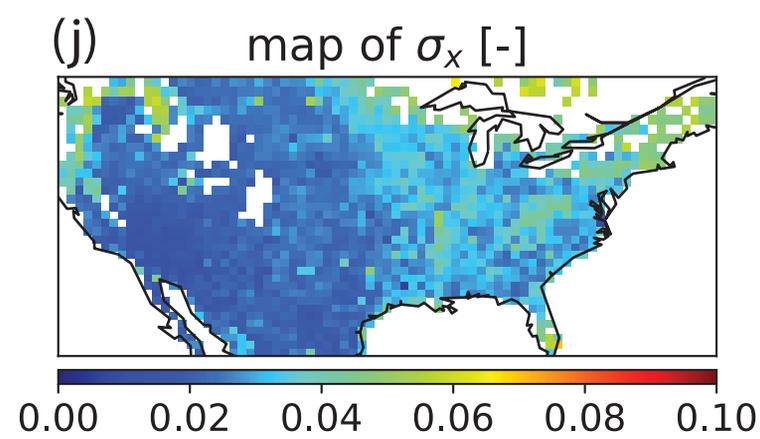
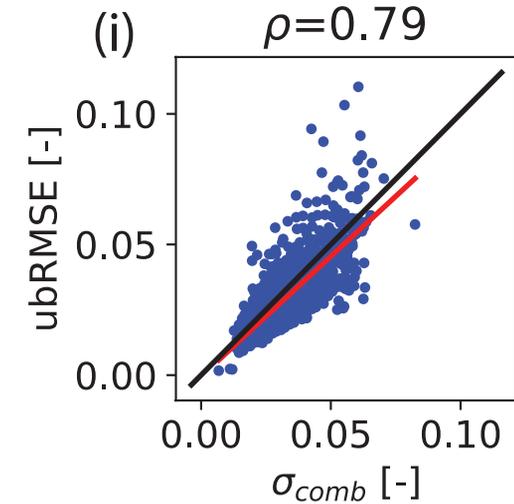
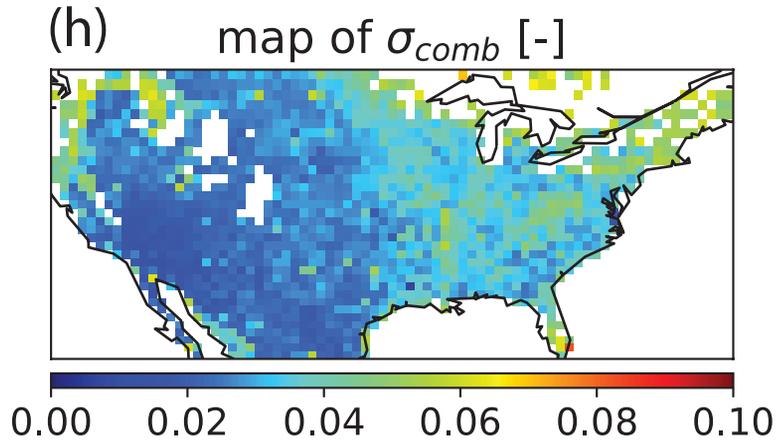
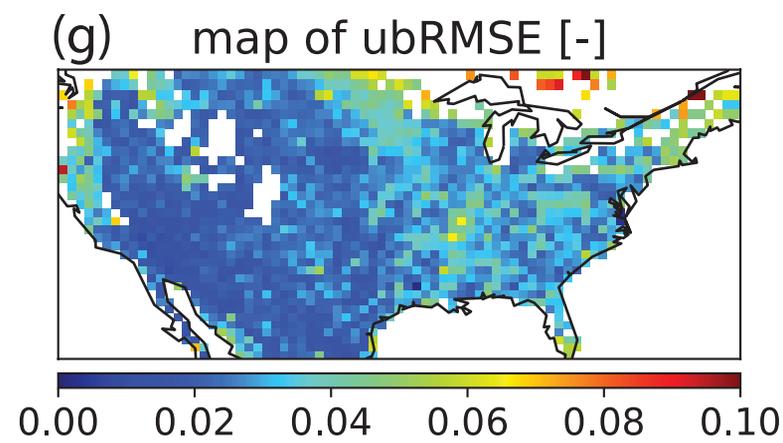
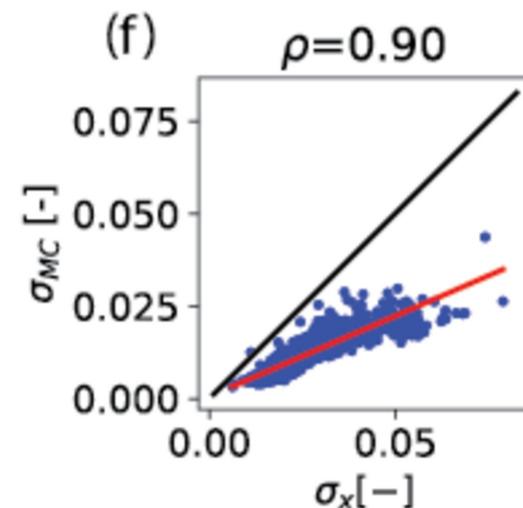
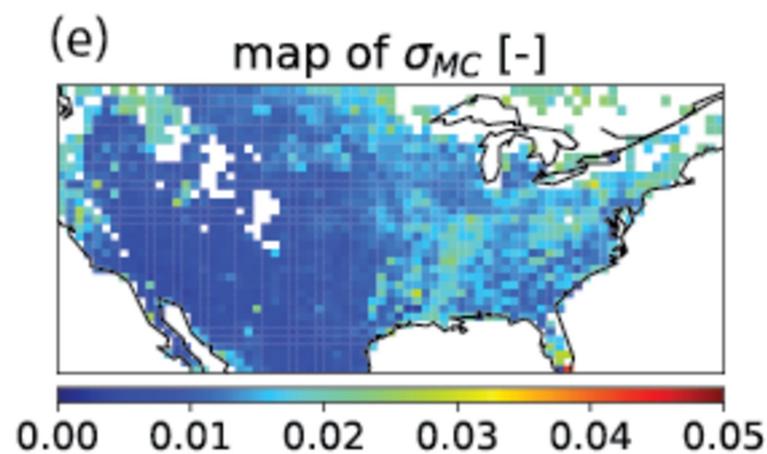
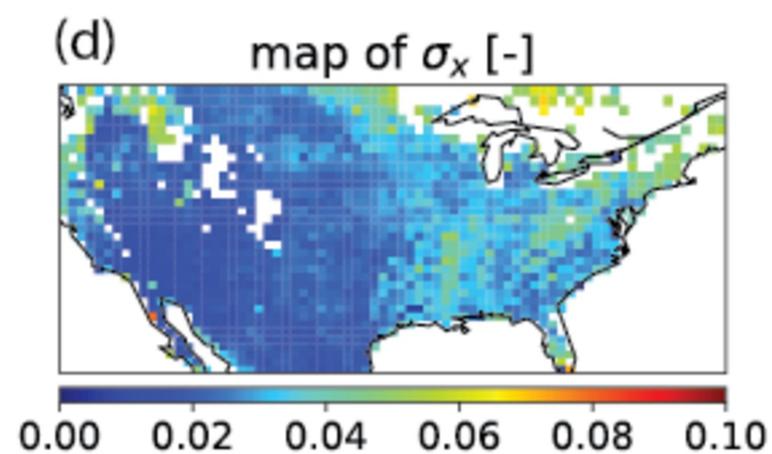
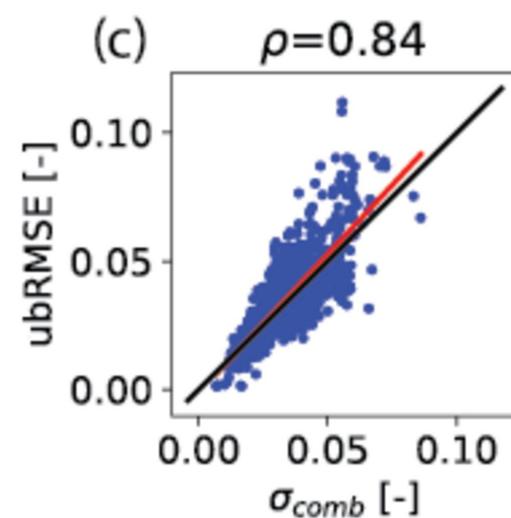
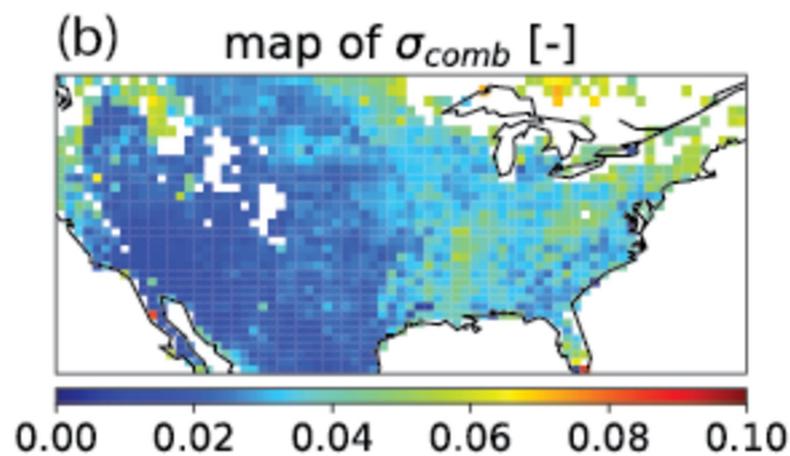
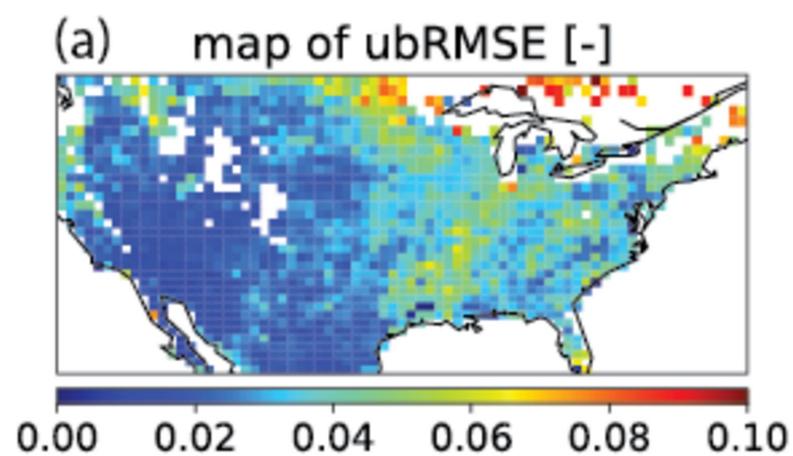


Figure 1 png ver.

Temporal Test



Spatial Test

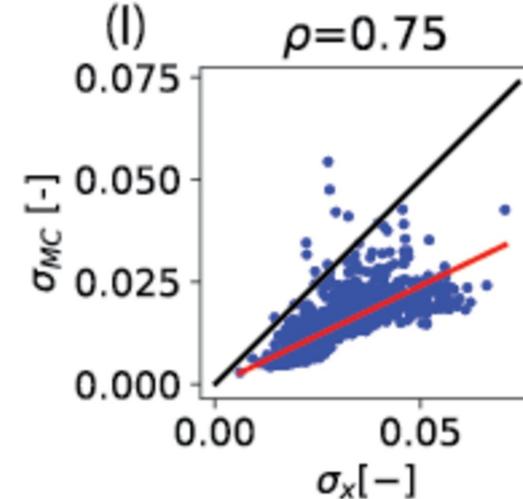
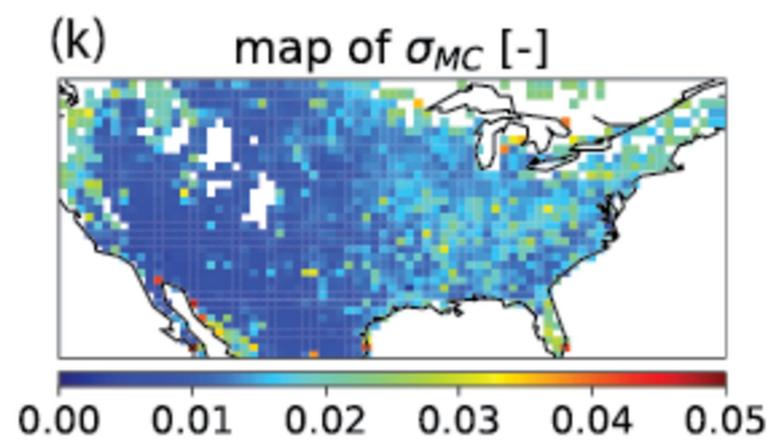
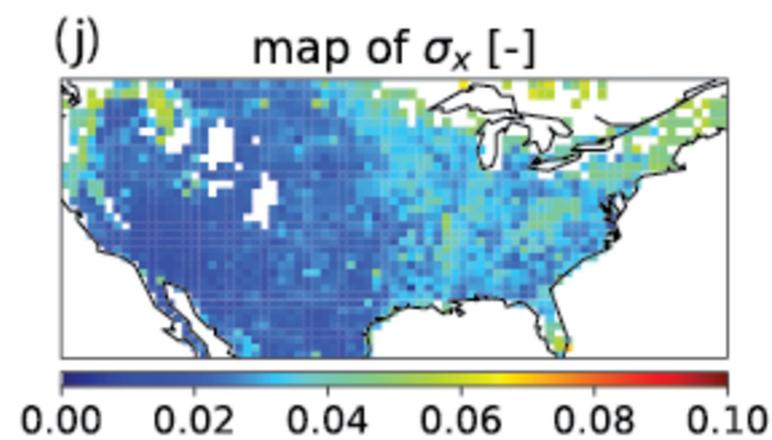
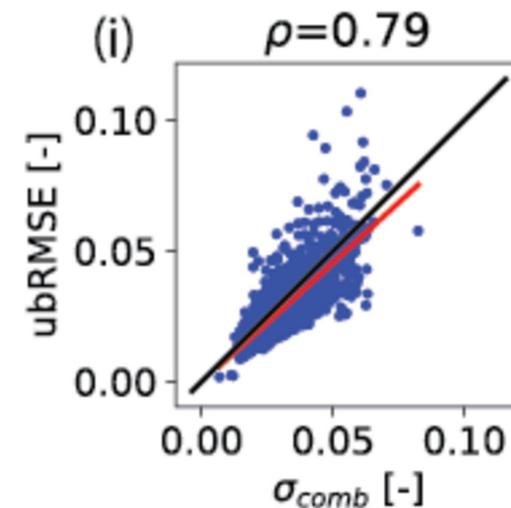
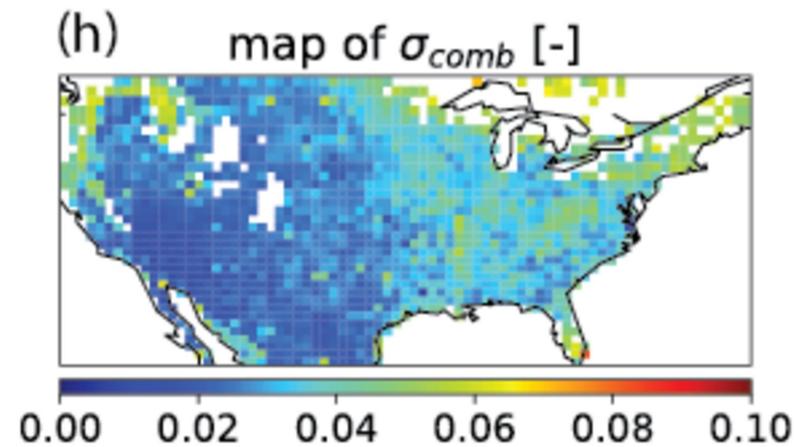
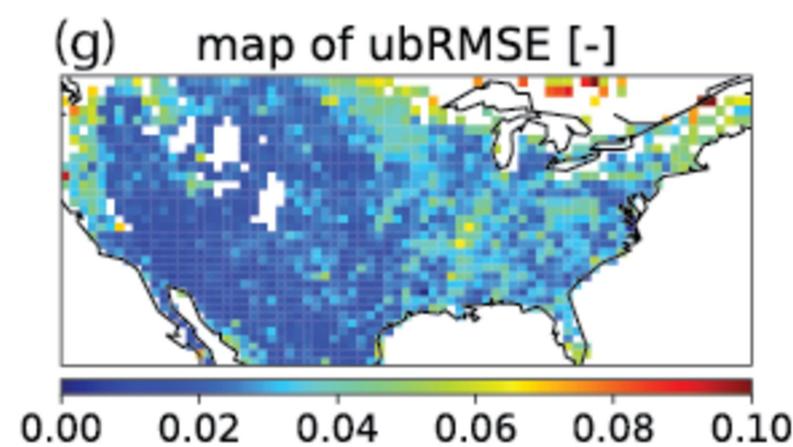
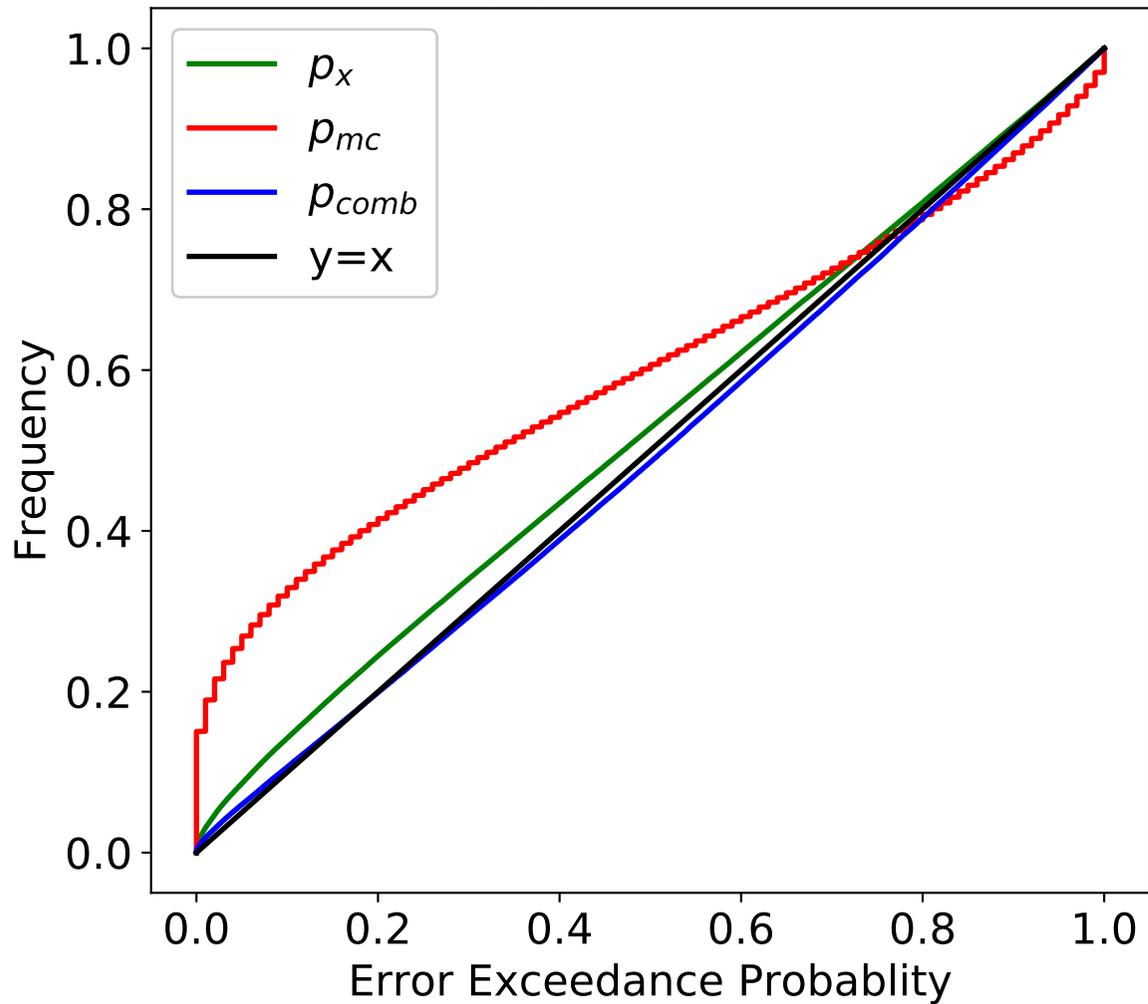


Figure 2.

(a) Validation



(b) Temporal Test

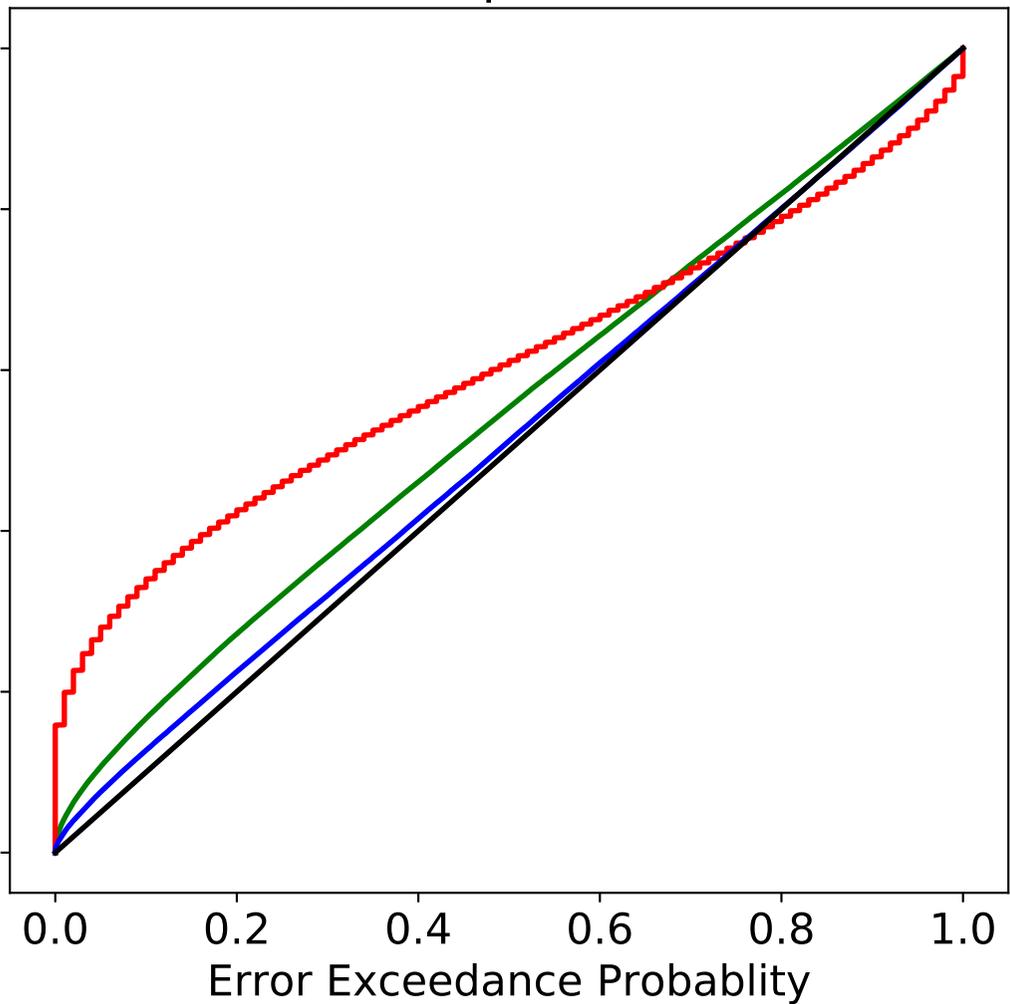
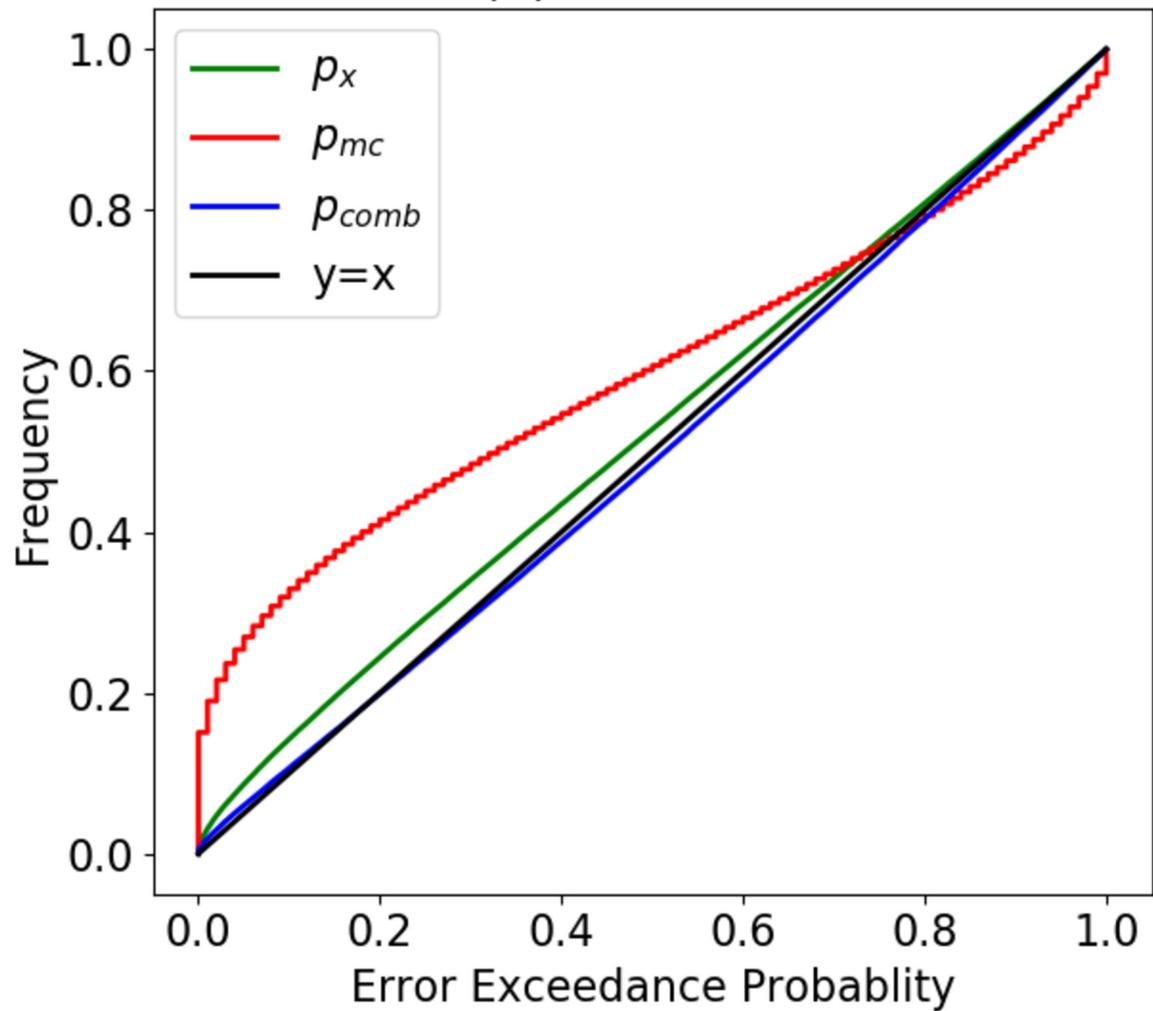


Figure 2 png ver.

(a) Validation



(b) Temporal Test

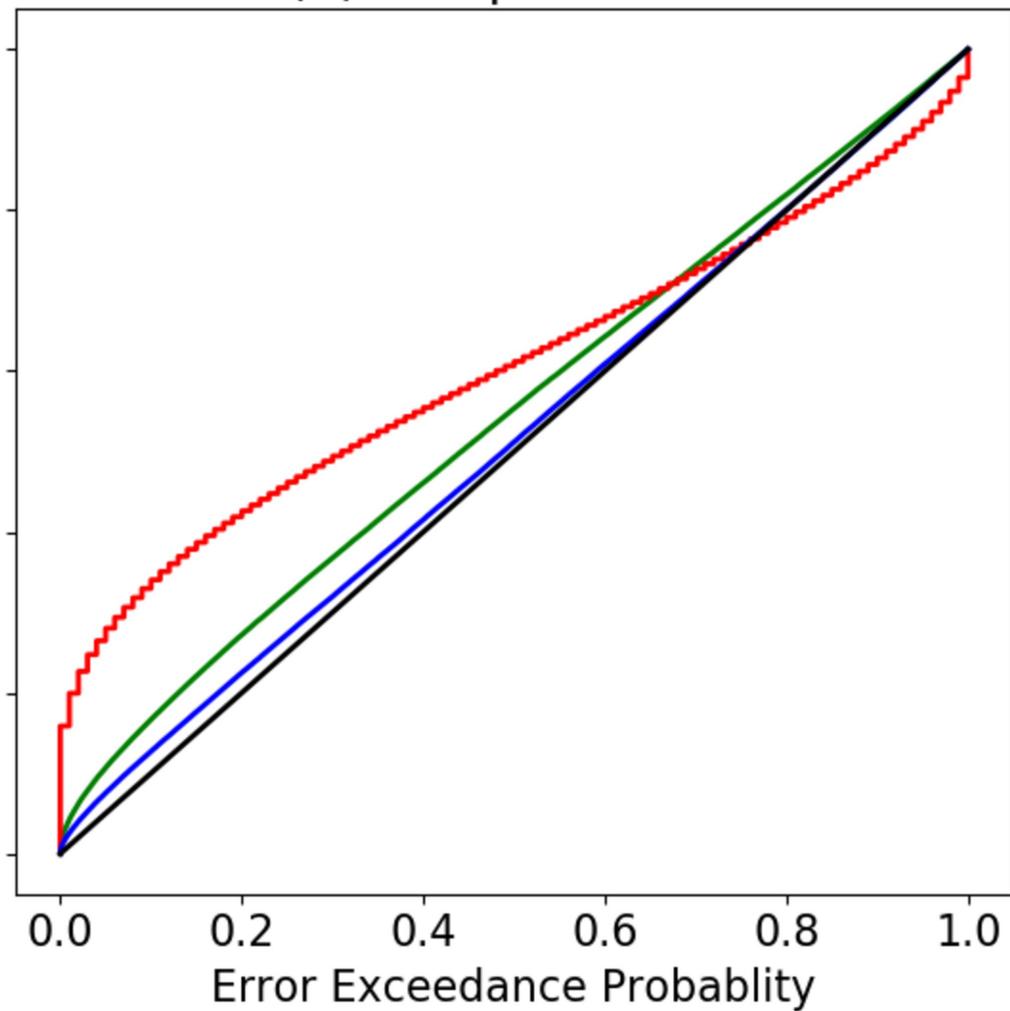
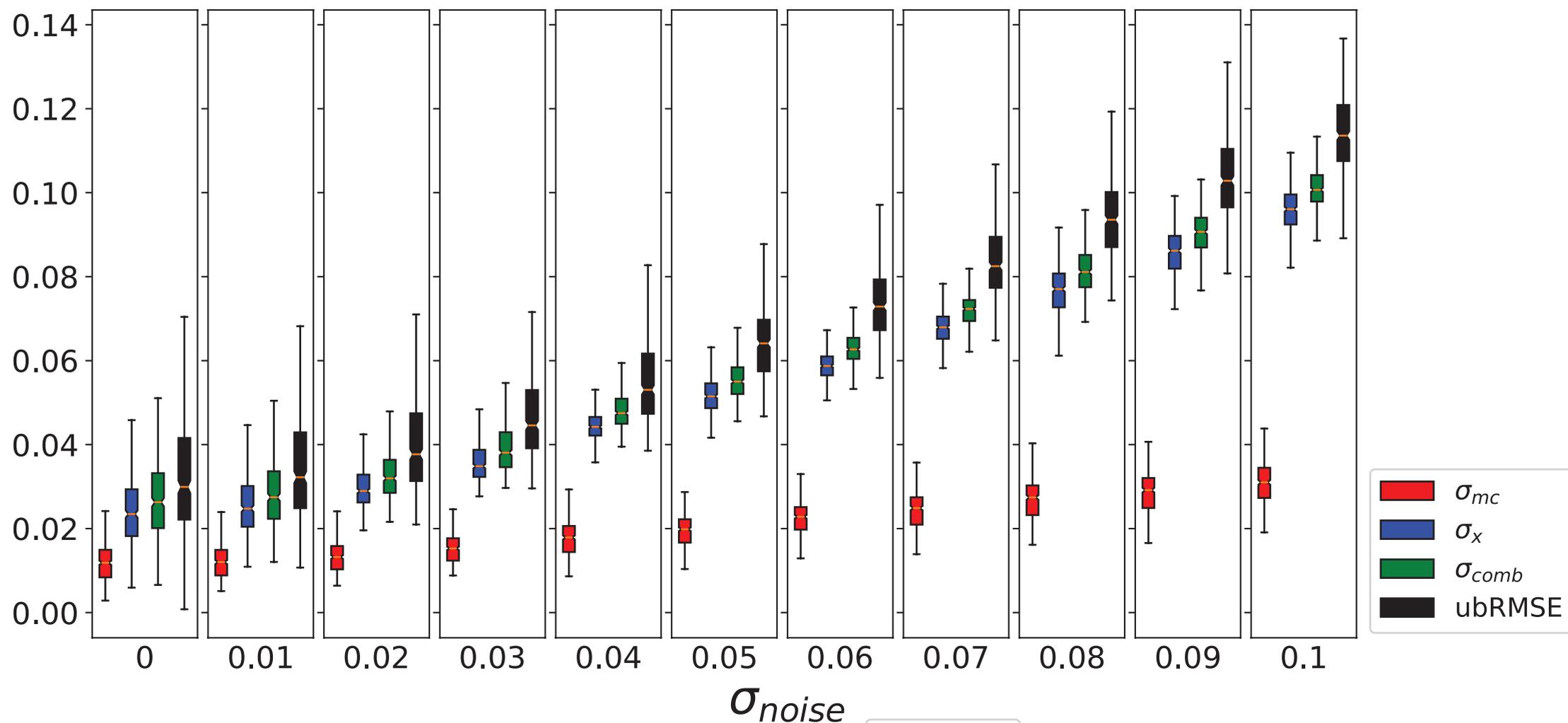
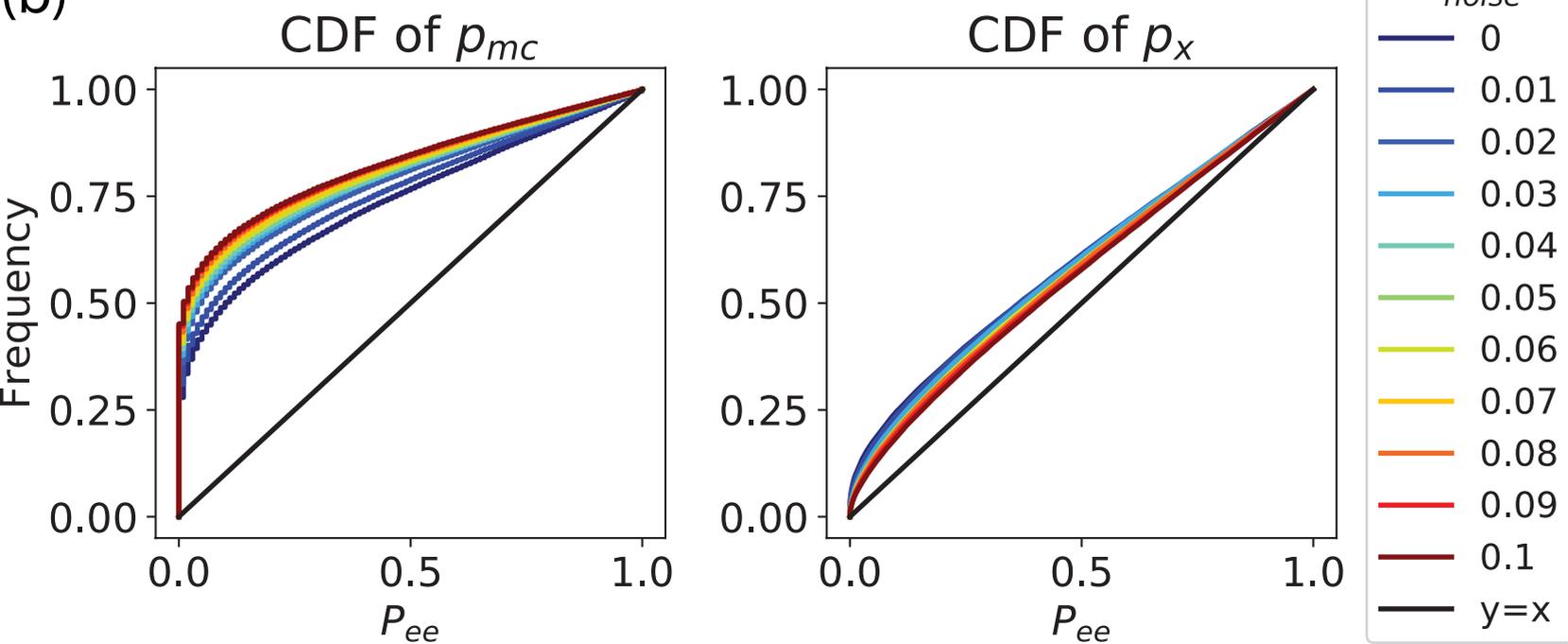


Figure 3.

(a) Error and uncertainty estimates in temporal test



(b)



(c)

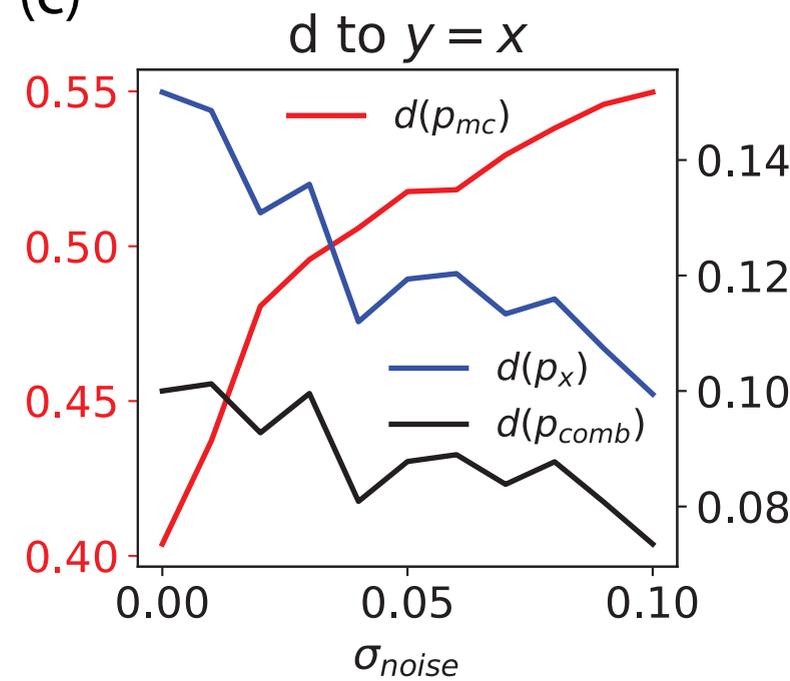
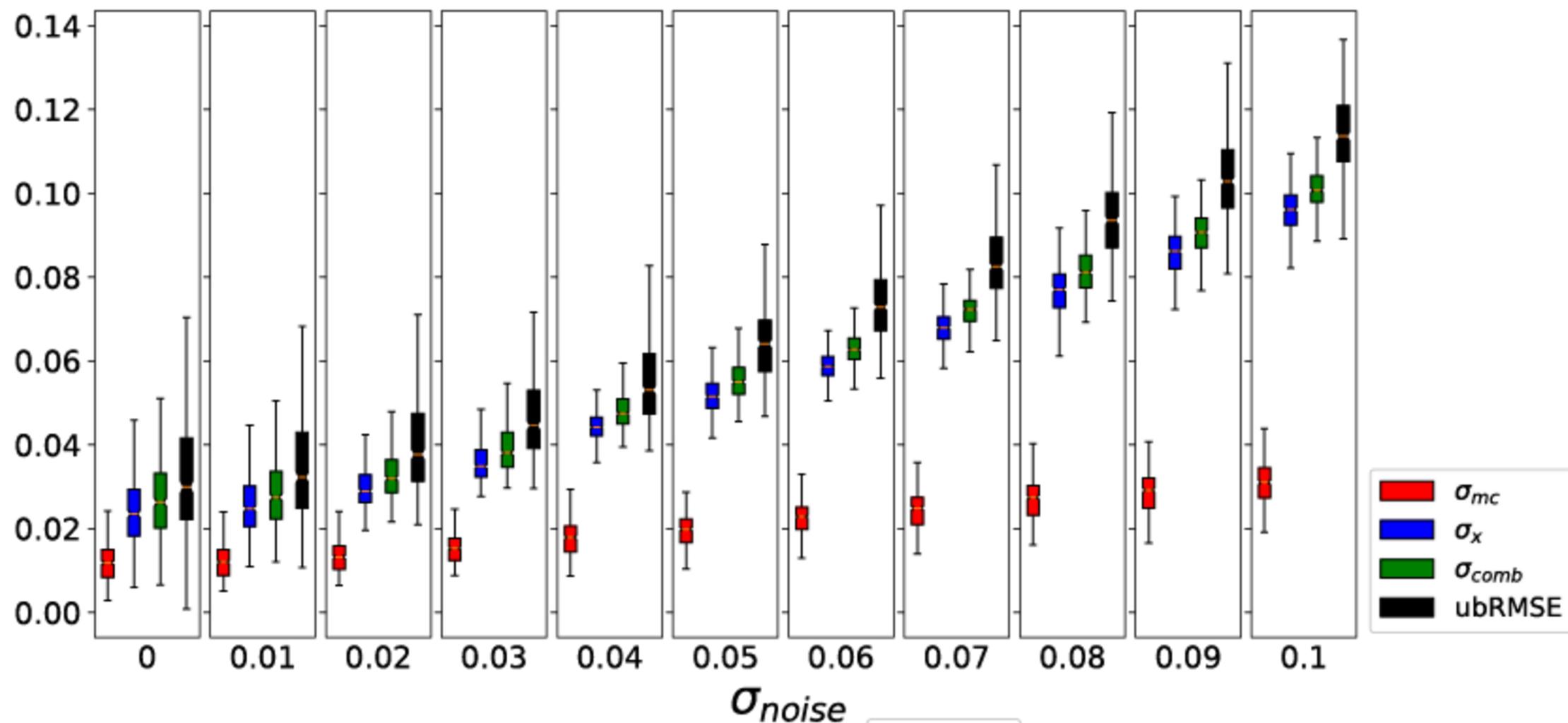
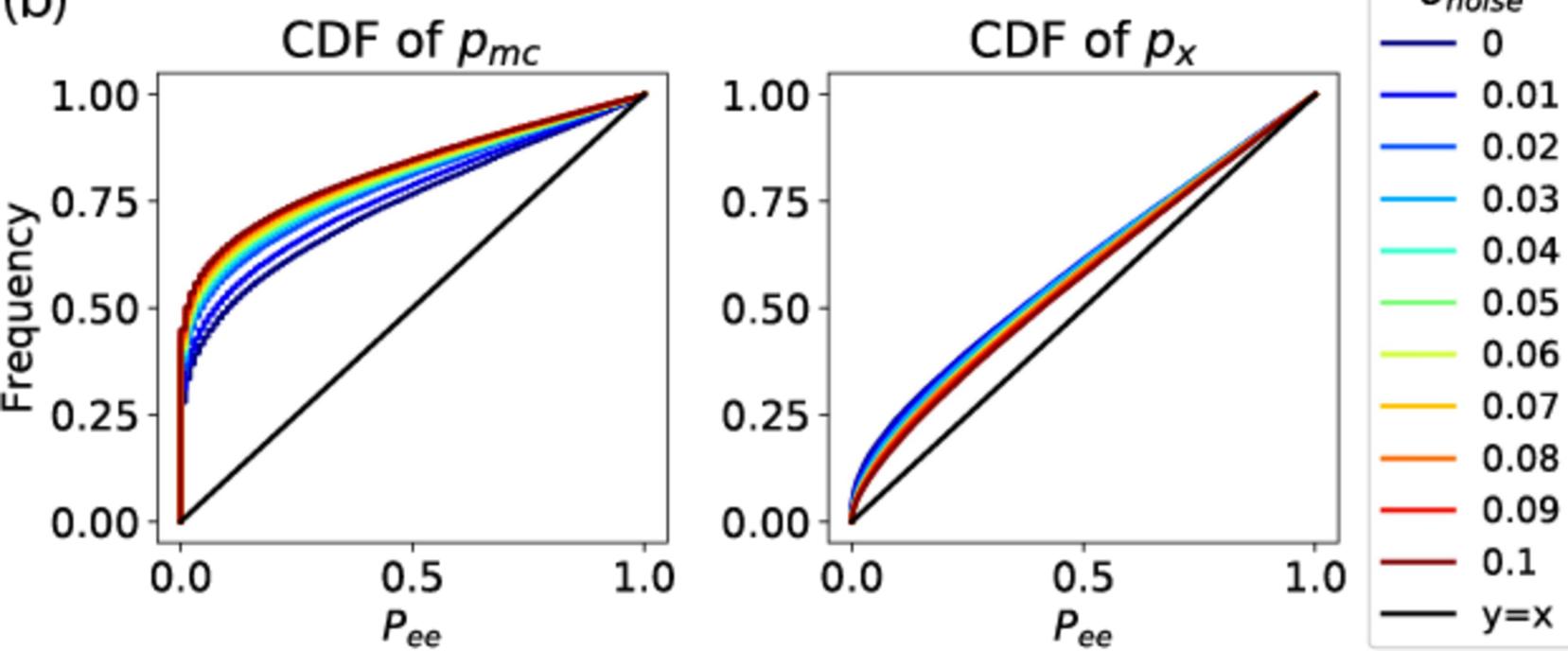


Figure 3 png ver.

(a) Error and uncertainty estimates in temporal test



(b)



(c)

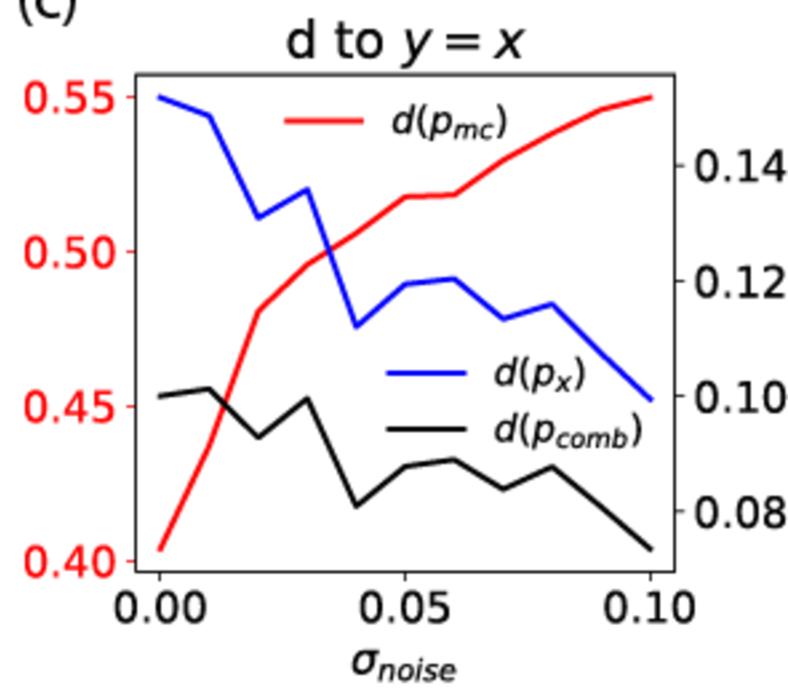
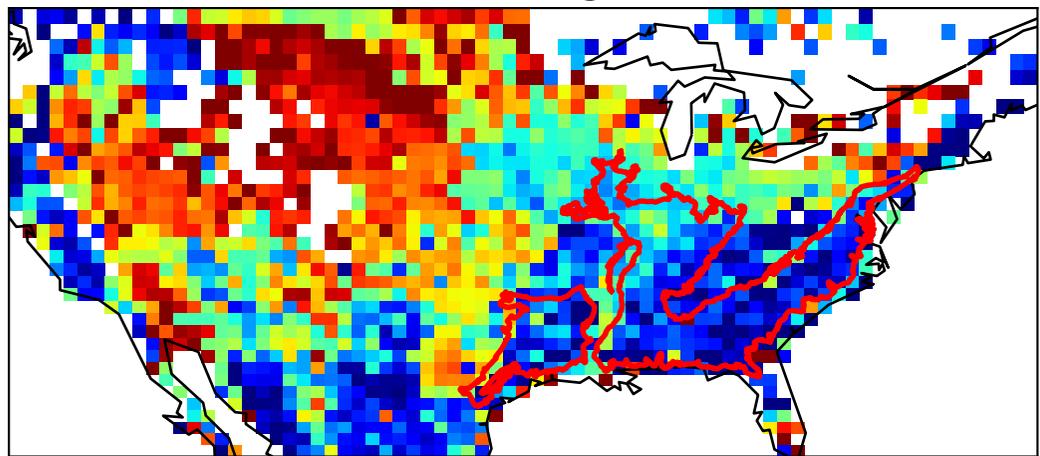
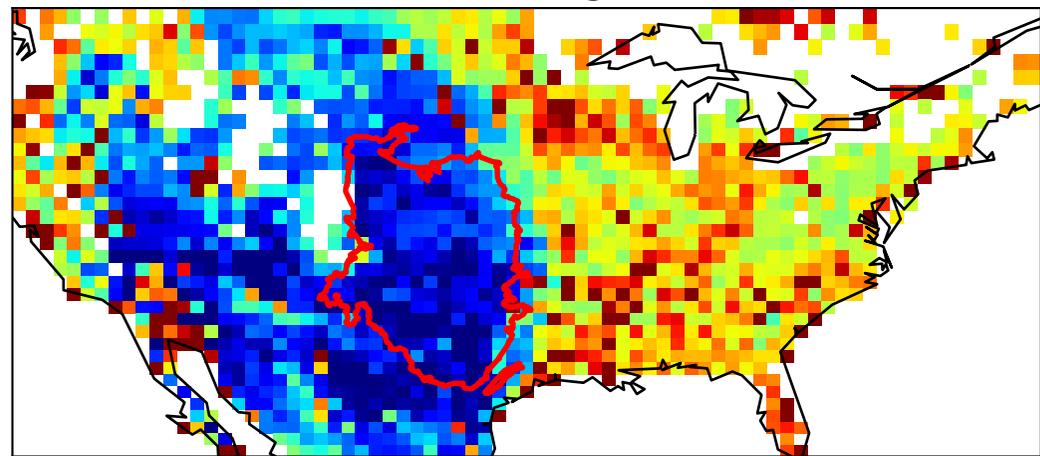


Figure 4.

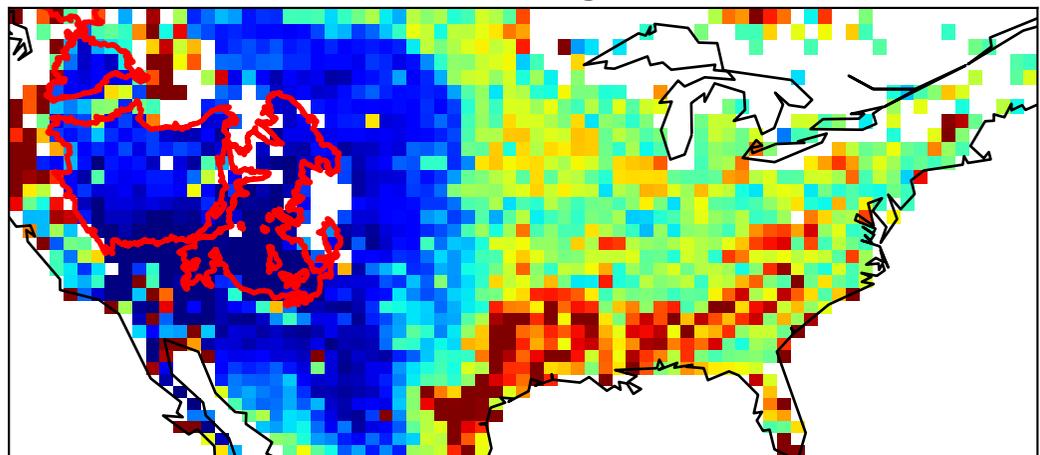
(a) σ_{mc} from Eco-region 05 model



(b) σ_{mc} from Eco-region 10 model



(c) σ_{mc} from Eco-region 12 model



(d) σ_{mc} from Eco-region 13 model

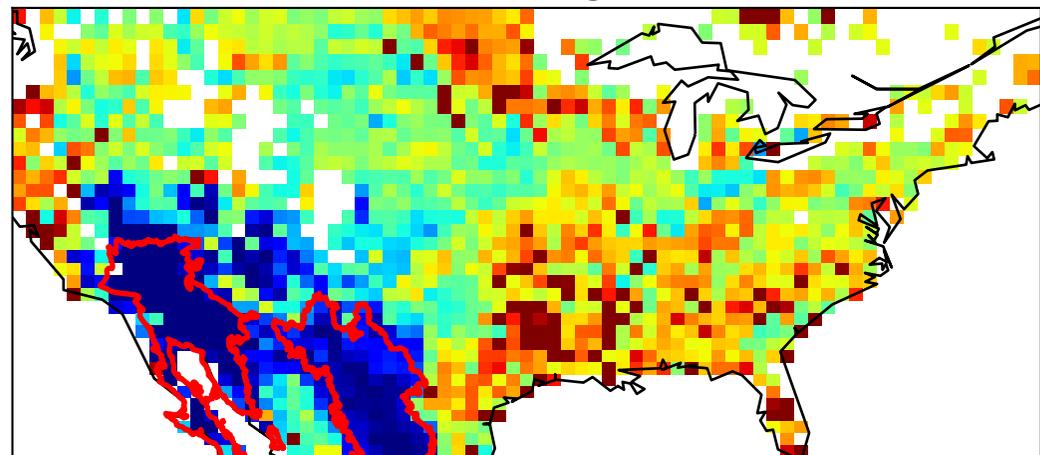
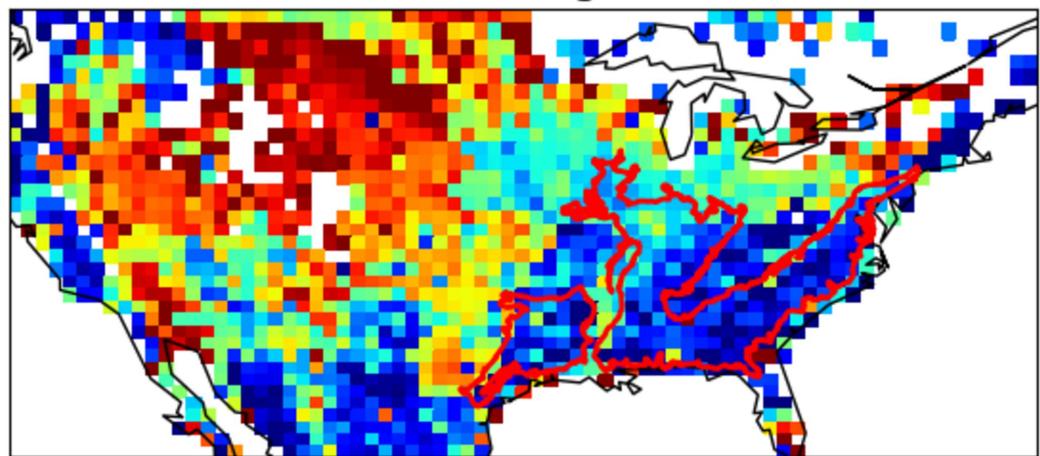
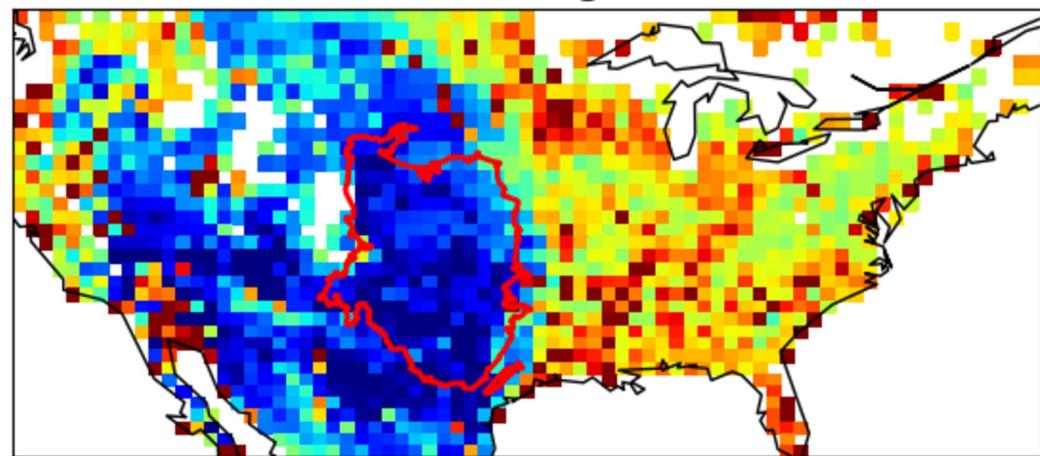


Figure 4 png ver.

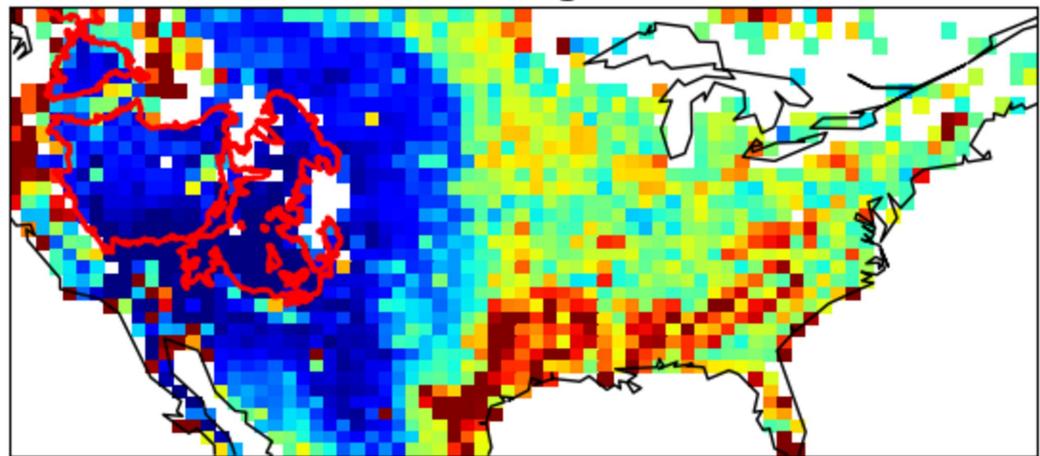
(a) σ_{mc} from Eco-region 8.3 model



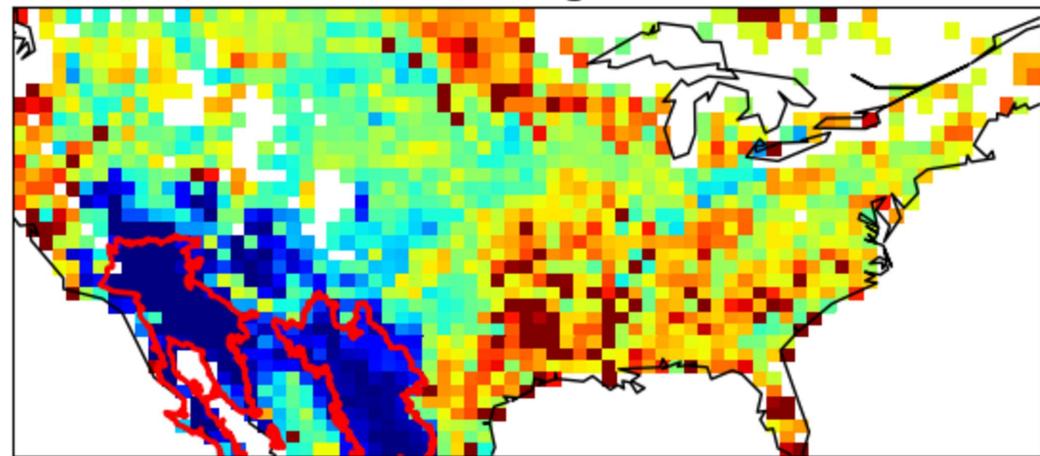
(b) σ_{mc} from Eco-region 9.4 model



(c) σ_{mc} from Eco-region 10.1 model



(d) σ_{mc} from Eco-region 10.2 model

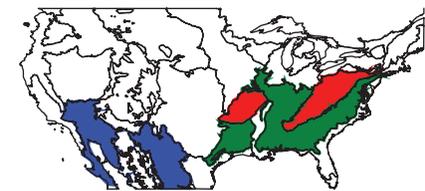
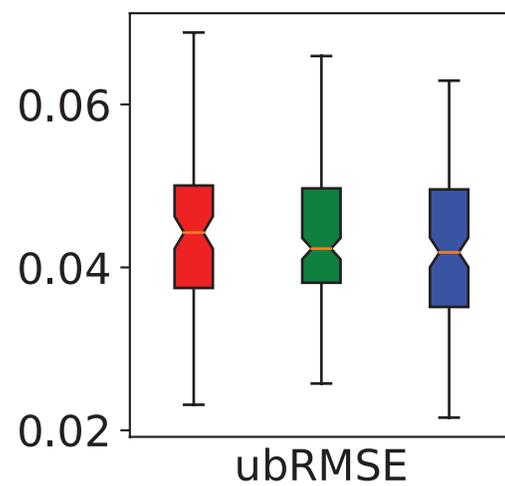
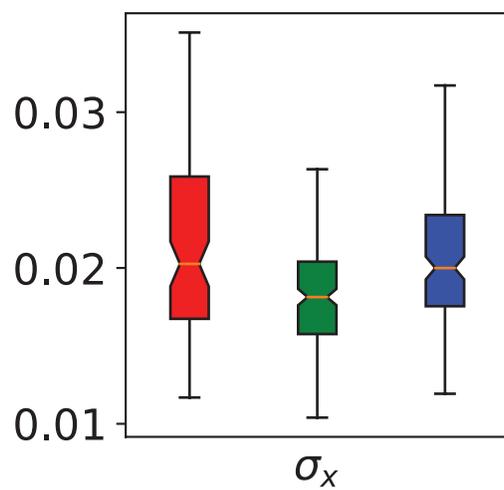
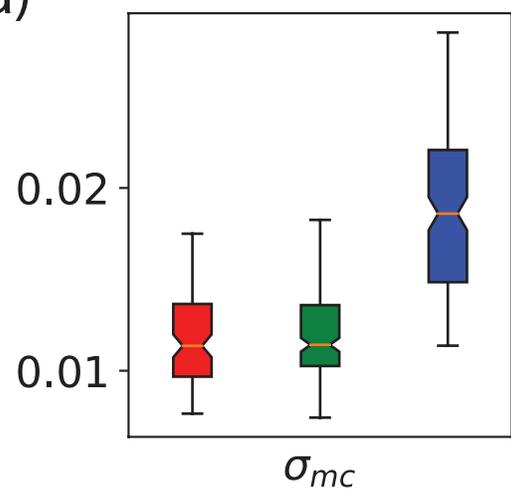


0.006 0.008 0.010 0.012 0.014 0.016 0.018

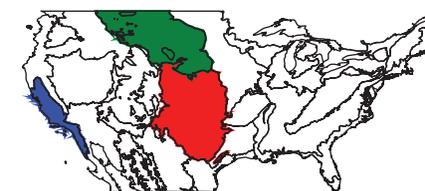
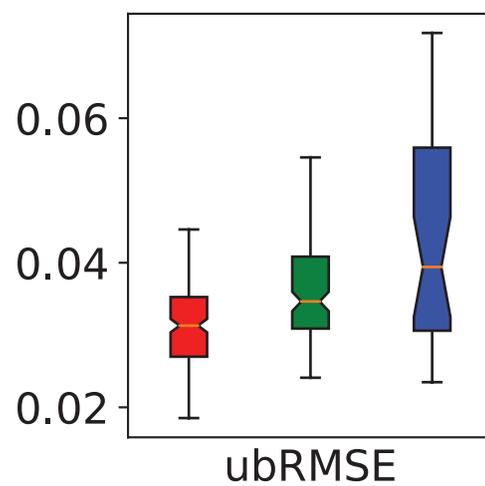
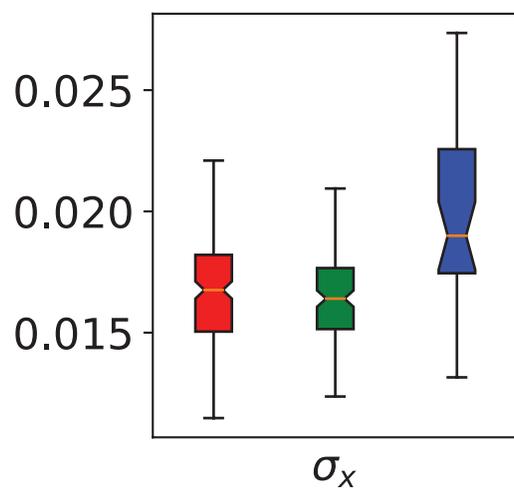
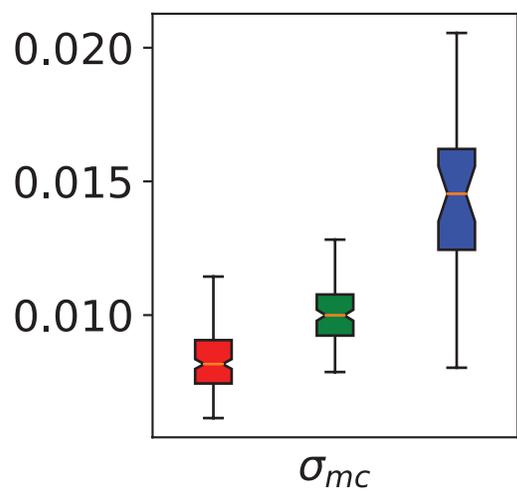
0.006 0.008 0.010 0.012 0.014 0.016 0.018 0.020 0.022

Figure 5.

(a)



(b)



(c)

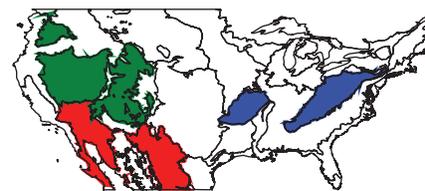
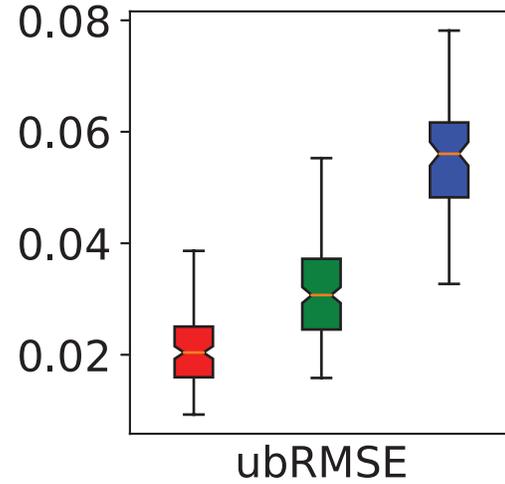
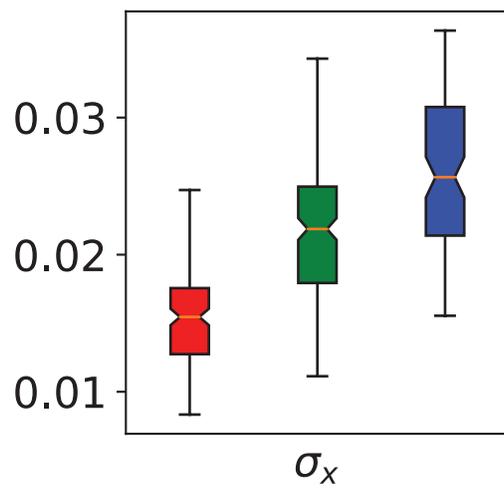
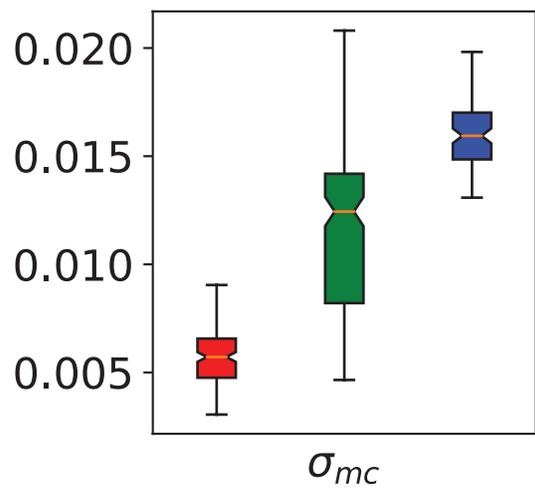
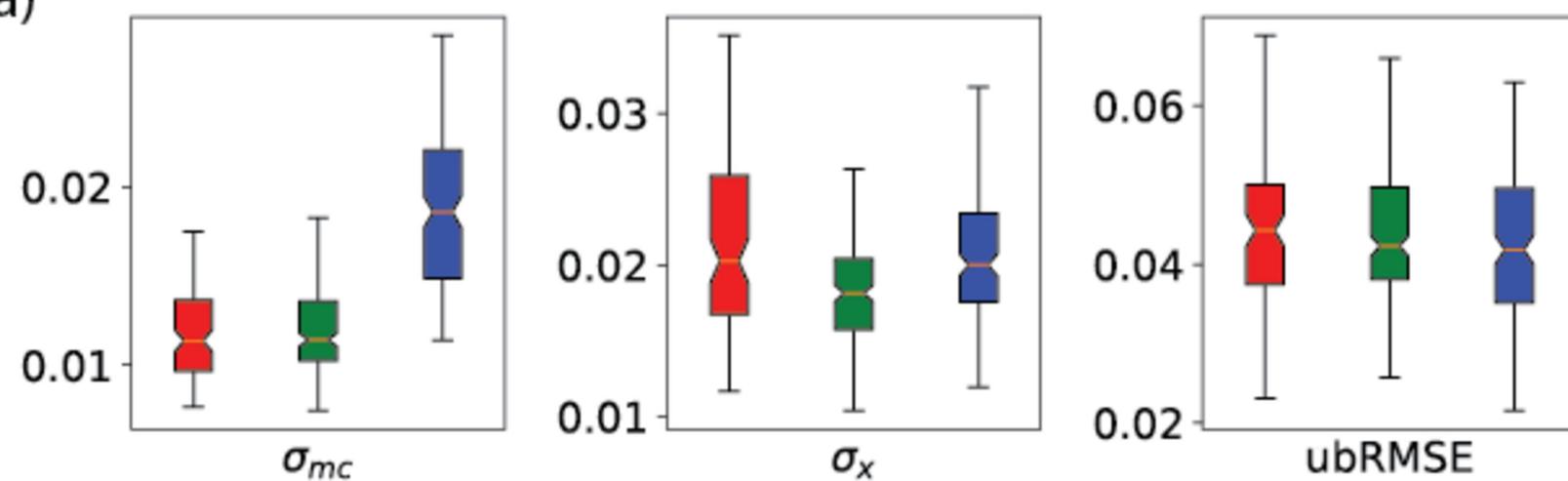
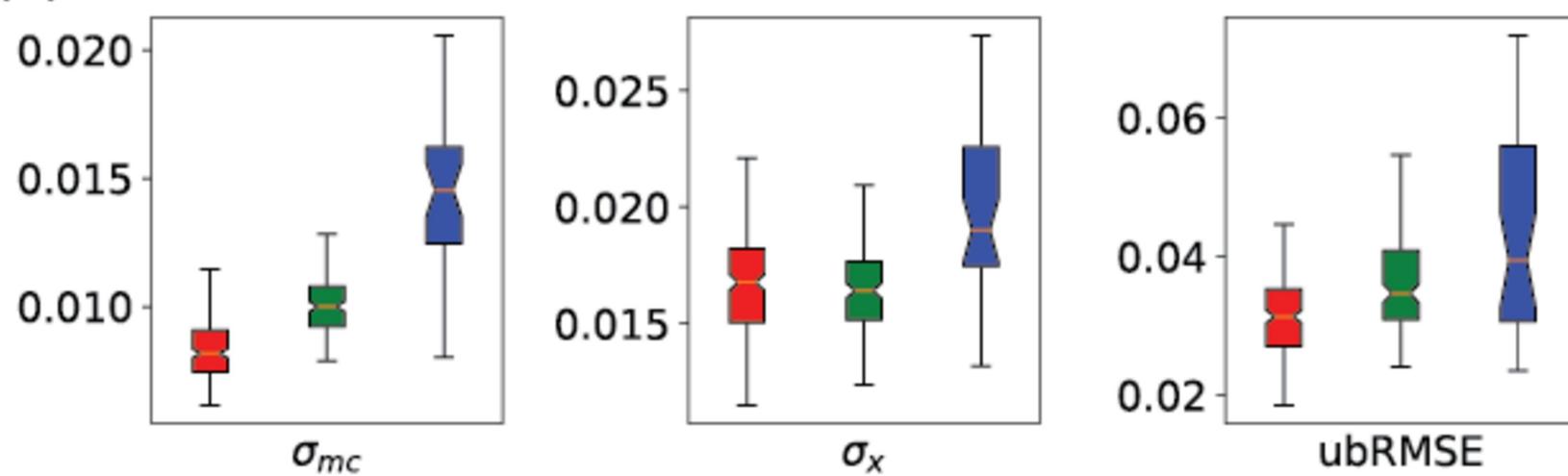


Figure 5 png ver.

(a)



(b)



(c)

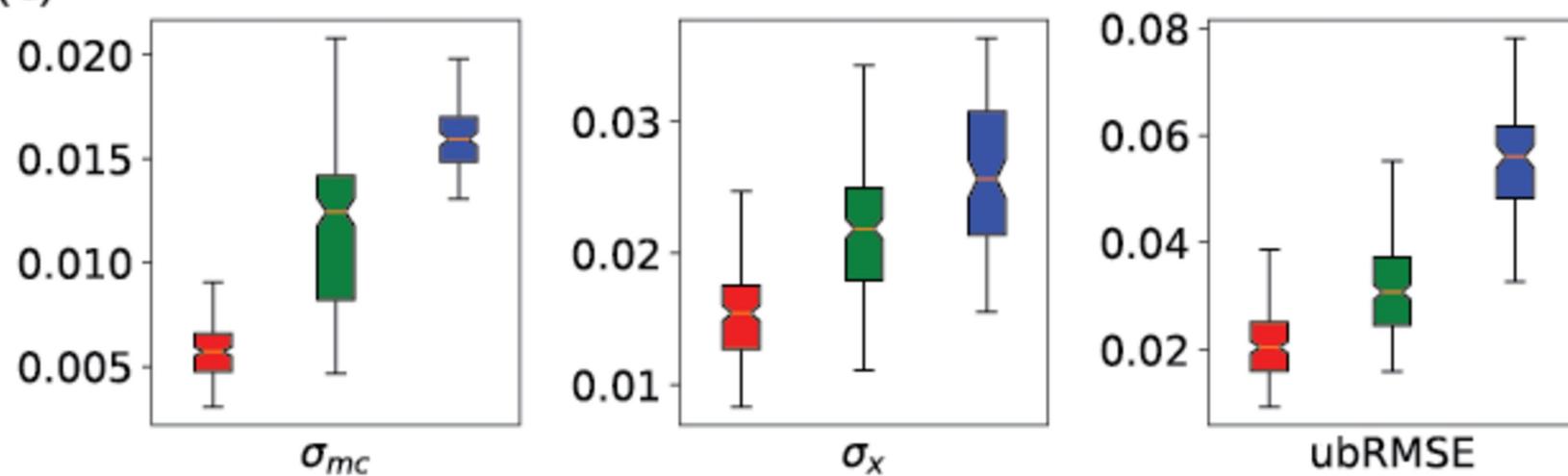


Figure 6.

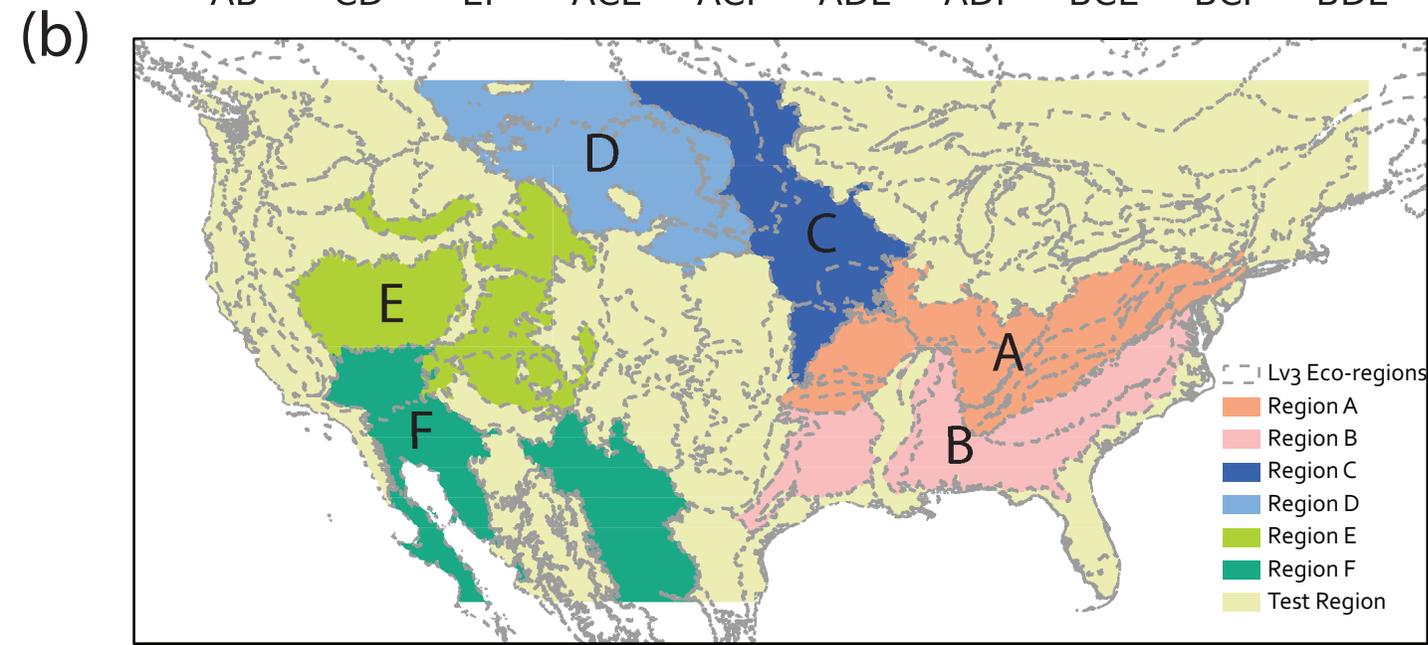
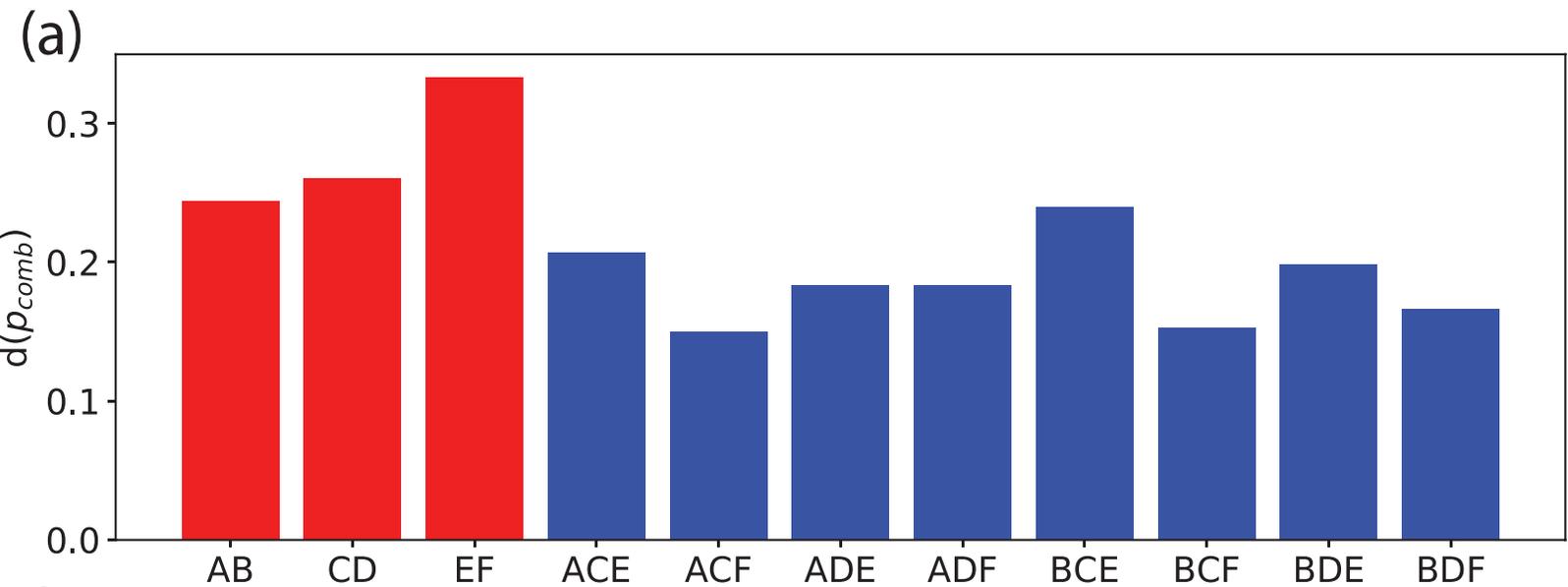


Figure 6 png ver.

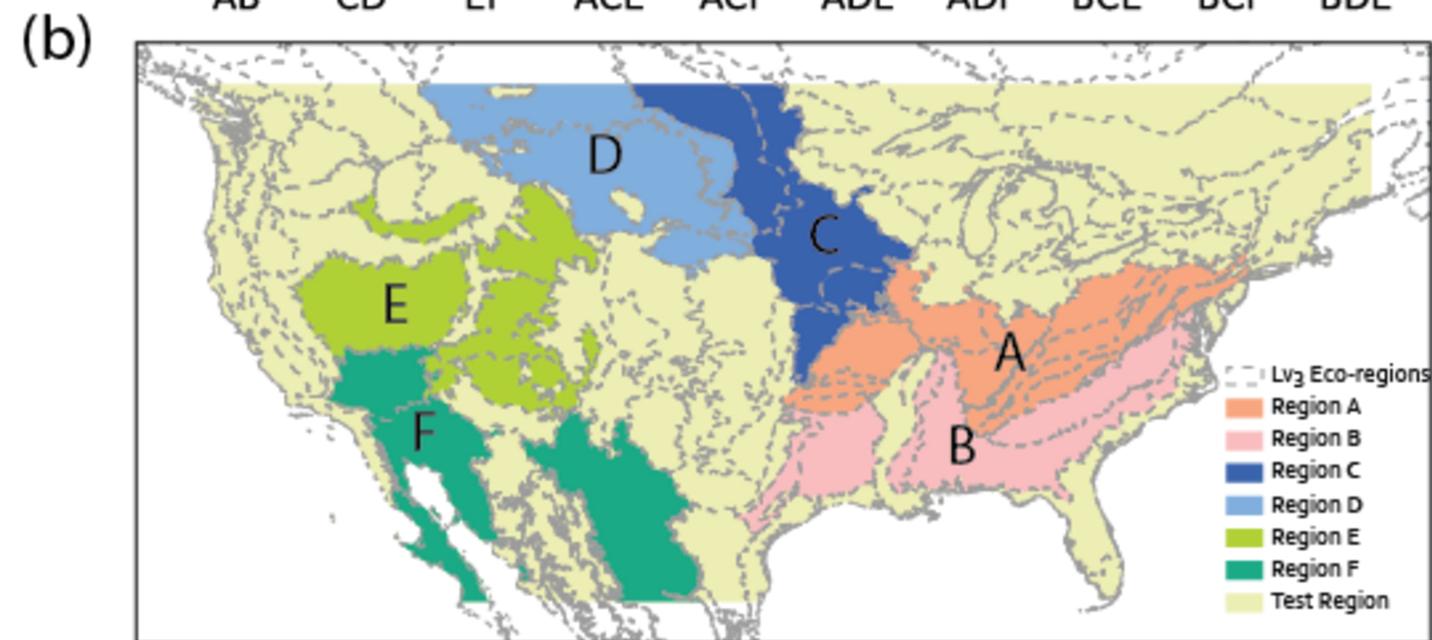
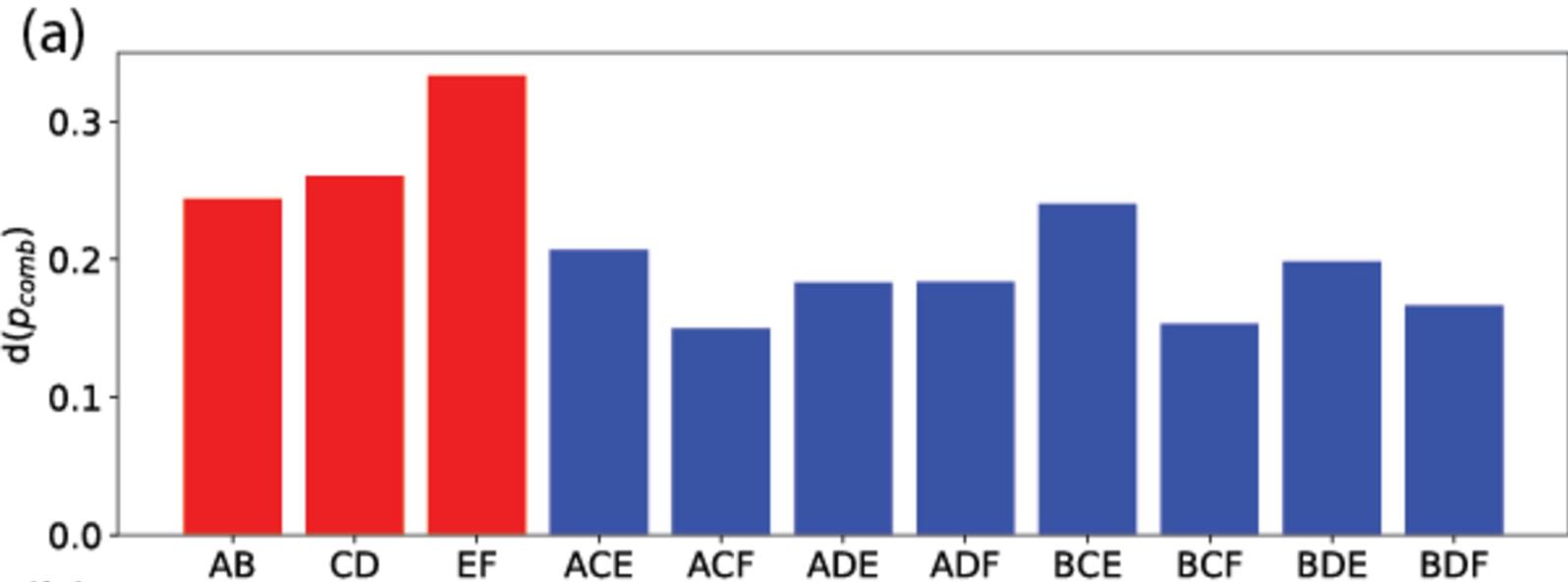


Figure B1.

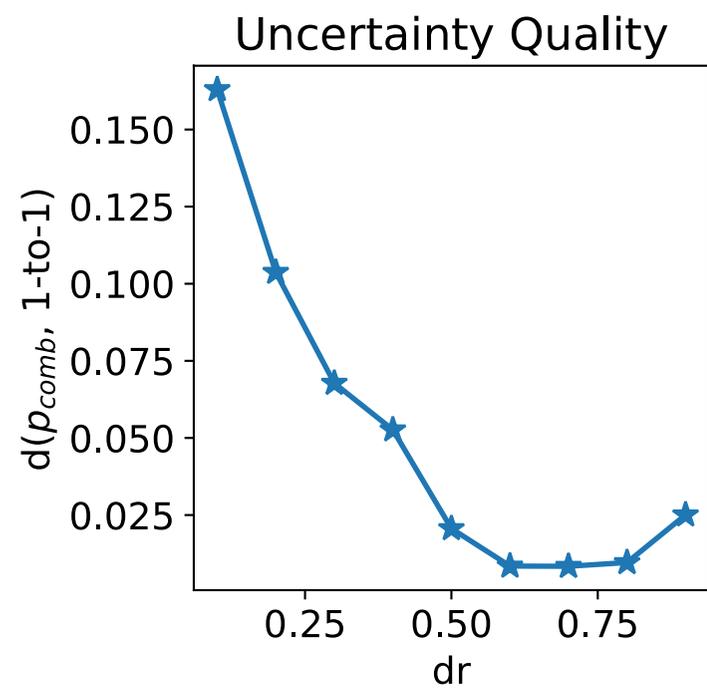
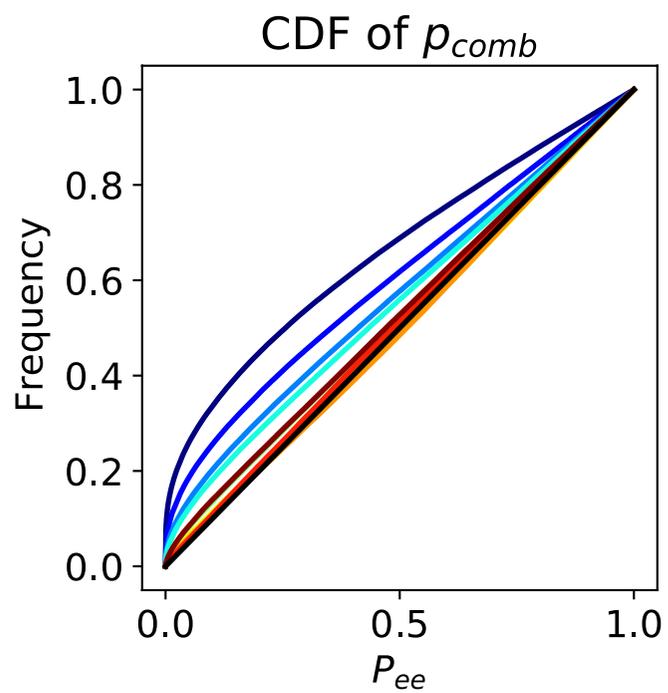
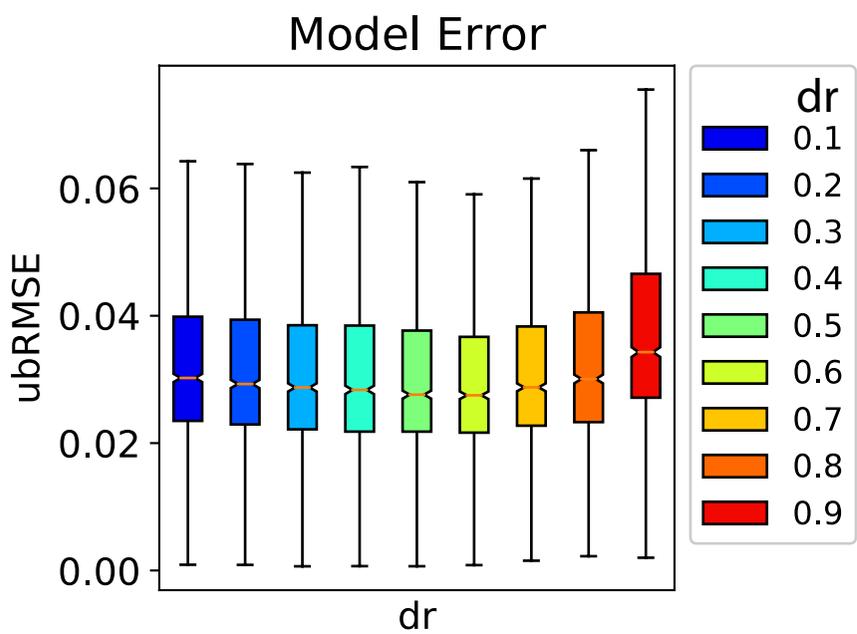
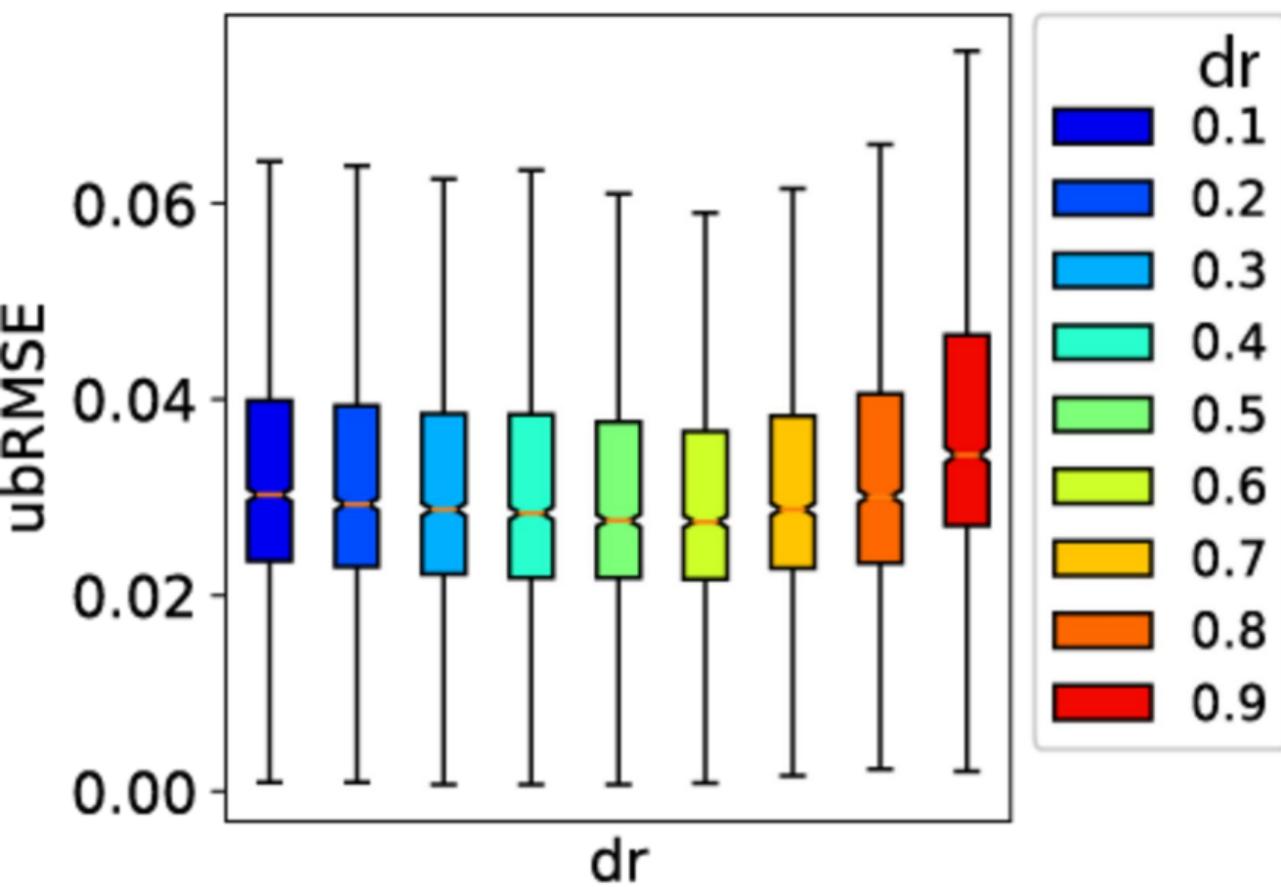
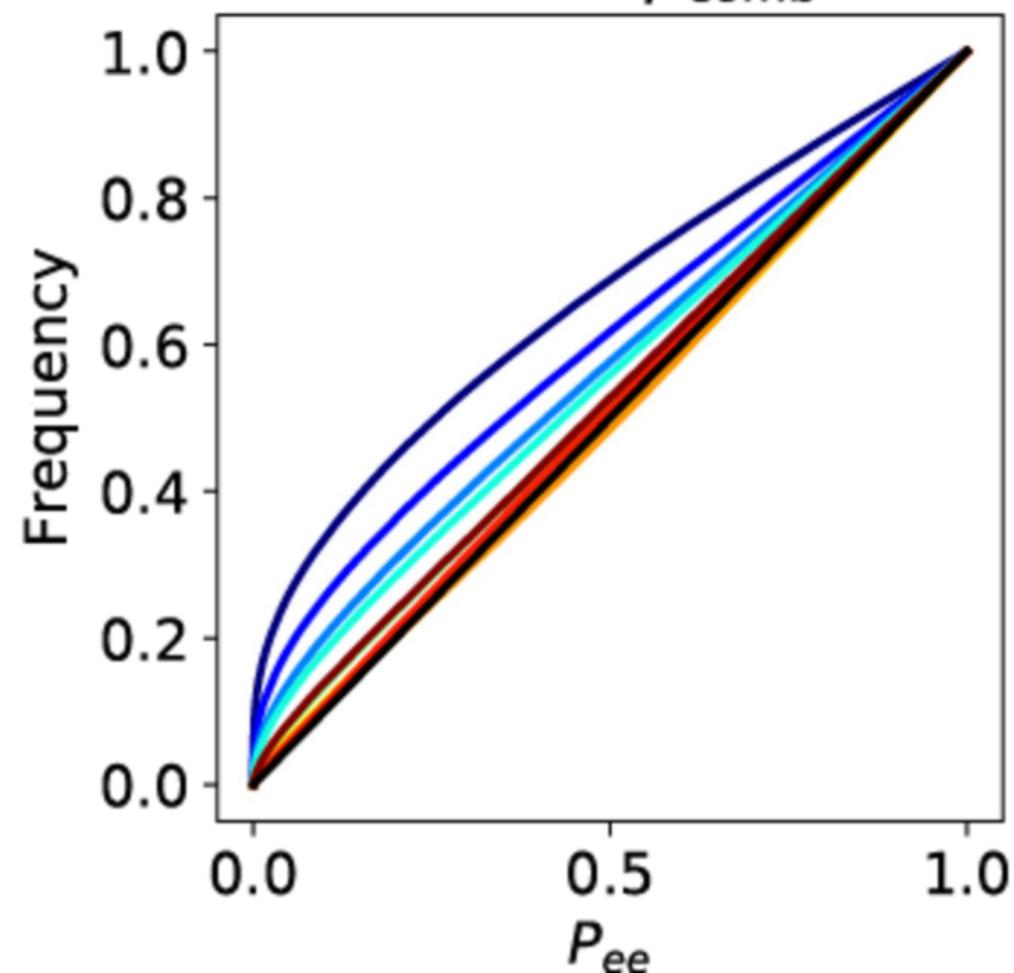


Figure B1 png ver.

Model Error

CDF of p_{comb} 

Uncertainty Quality

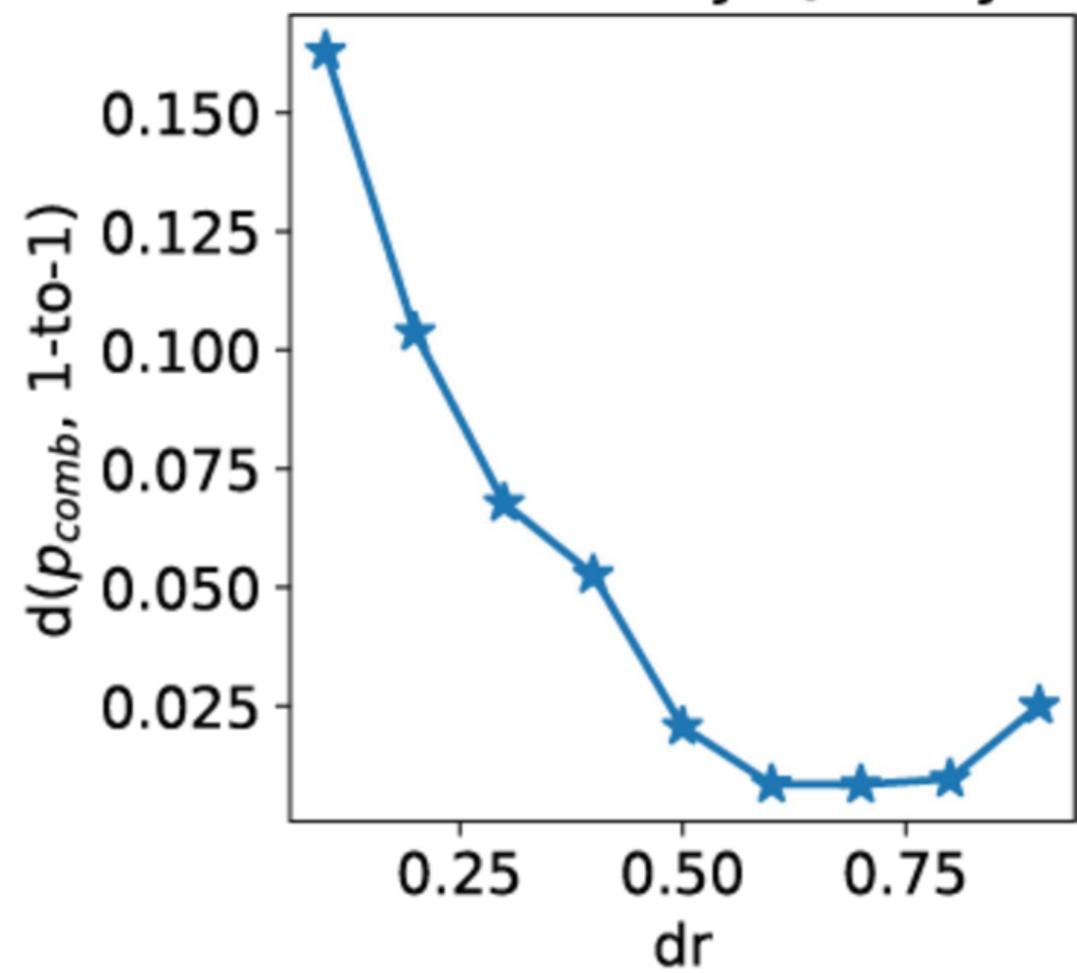
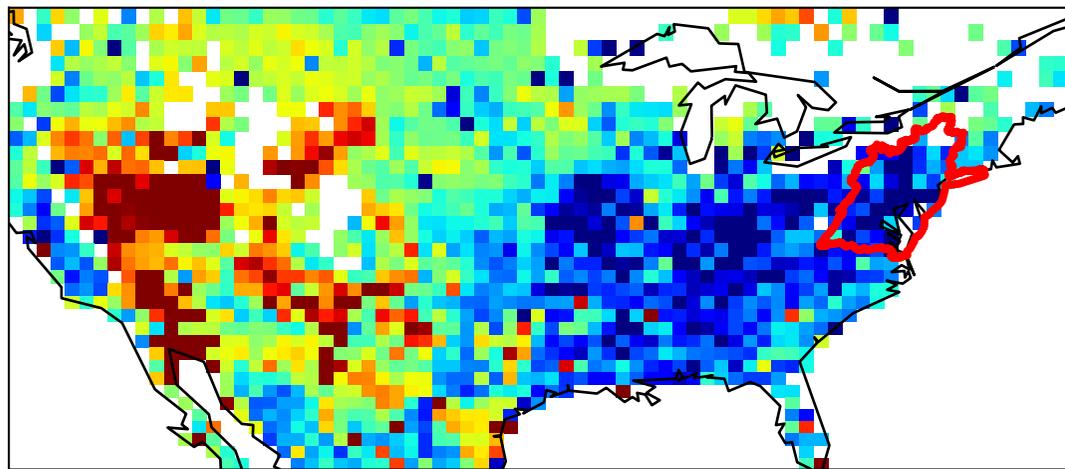


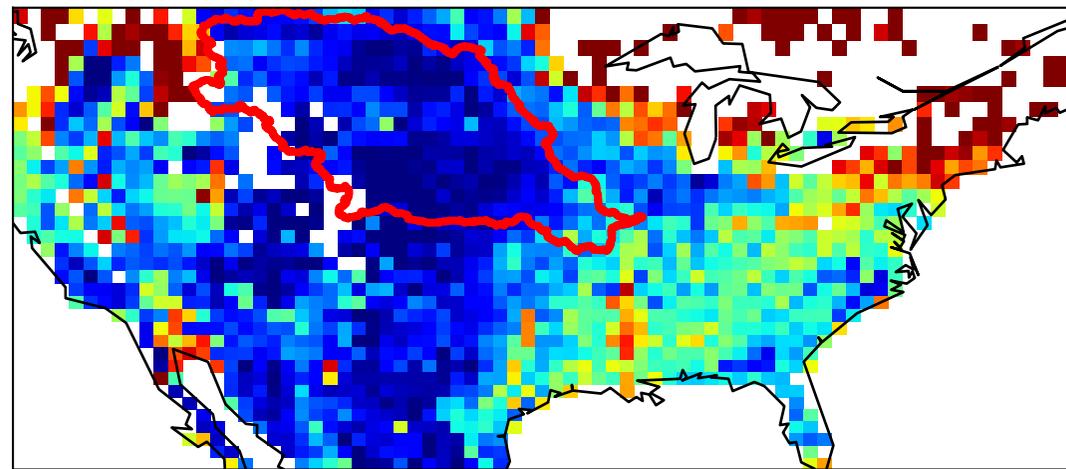
Figure C1.

(a) σ_{mc} from HUC02 model



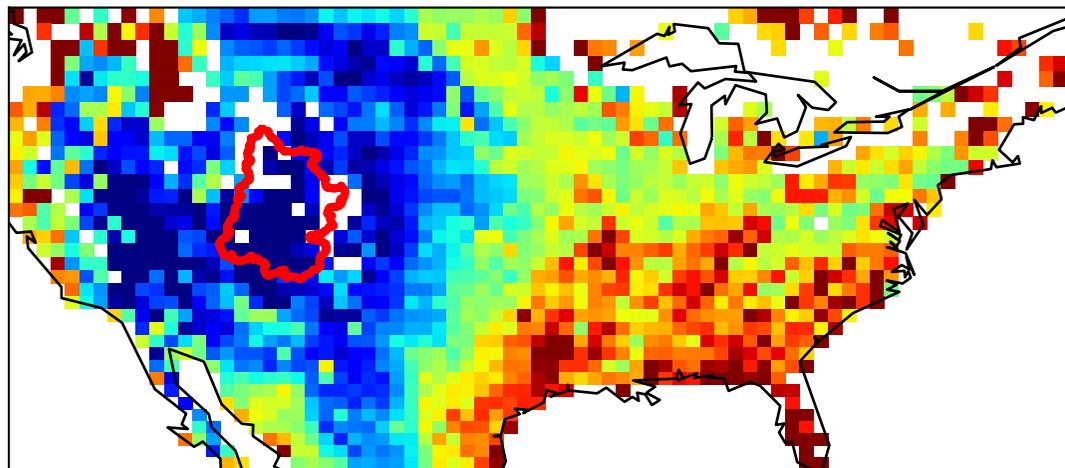
0.010 0.012 0.014 0.016 0.018 0.020

(b) σ_{mc} from HUC10 model



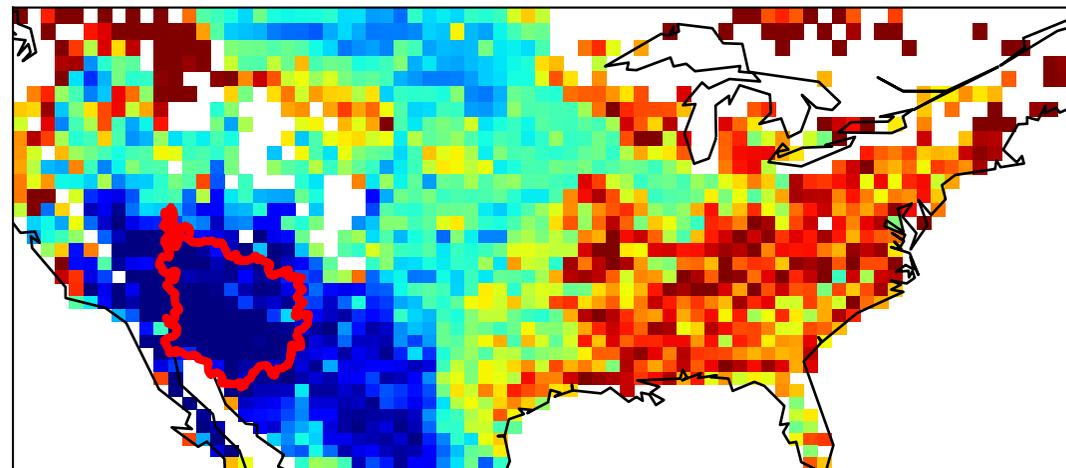
0.010 0.012 0.014 0.016 0.018 0.020 0.022

(c) σ_{mc} from HUC14 model



0.005 0.006 0.007 0.008 0.009 0.010 0.011

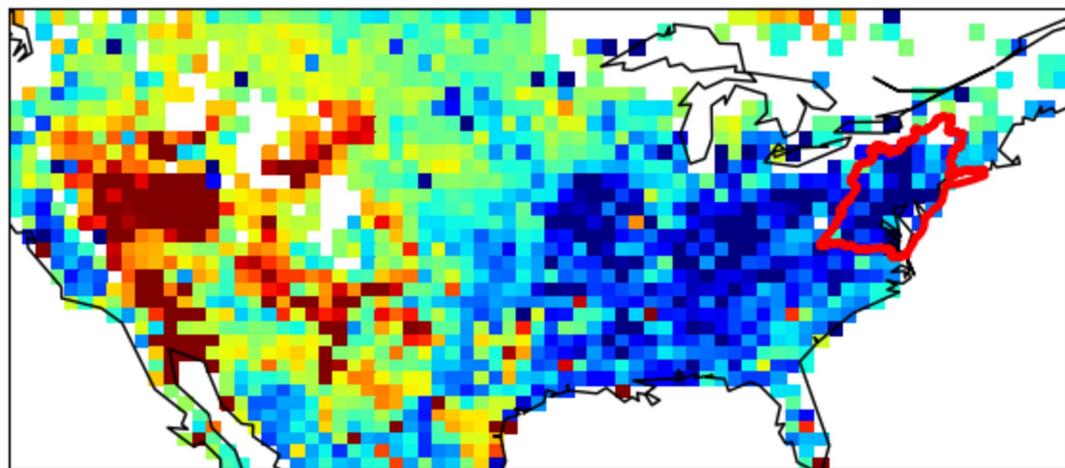
(d) σ_{mc} from HUC15 model



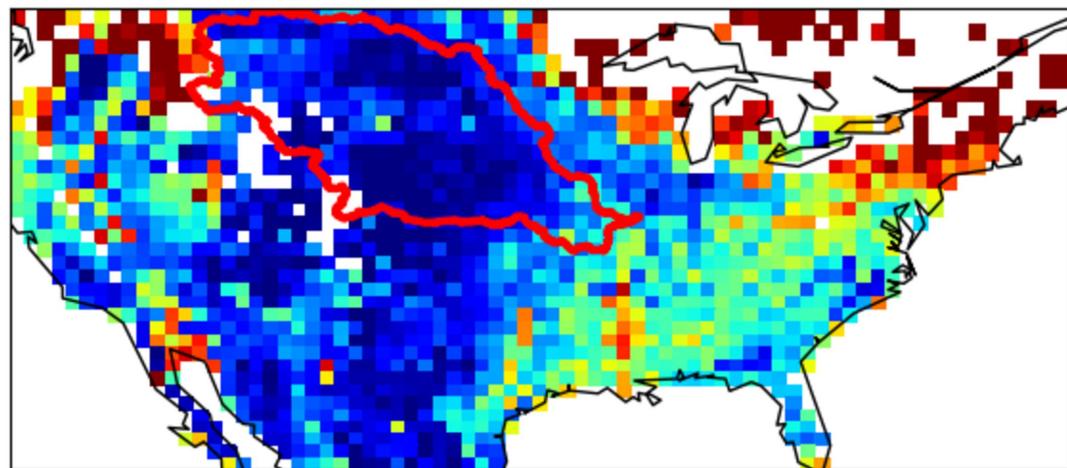
0.006 0.008 0.010 0.012

Figure C1 png ver.

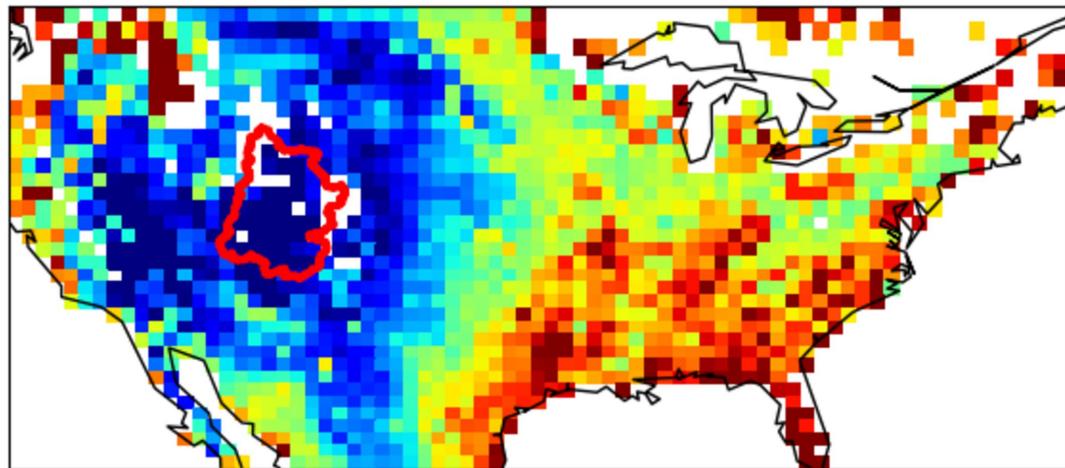
(a) σ_{mc} from HUC02 model



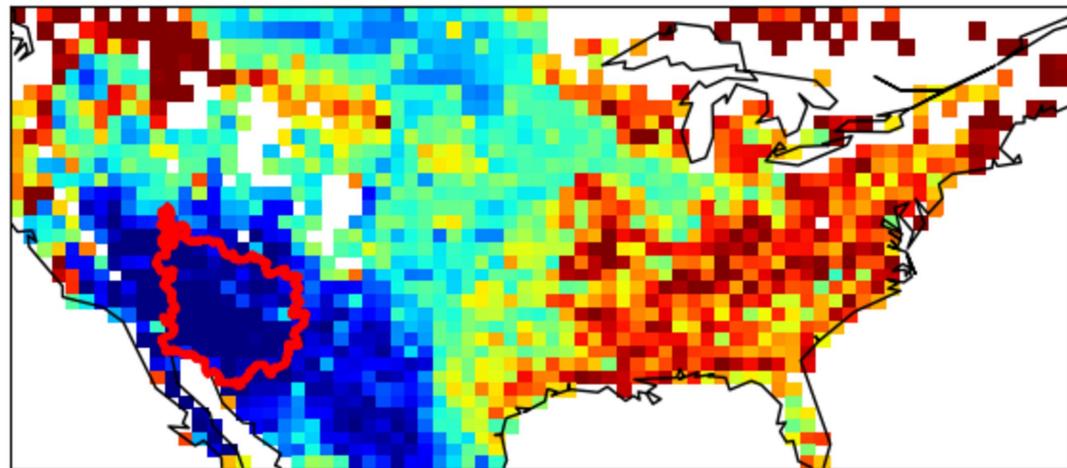
(b) σ_{mc} from HUC10 model



(c) σ_{mc} from HUC14 model



(d) σ_{mc} from HUC15 model



0.005 0.006 0.007 0.008 0.009 0.010 0.011

0.006 0.008 0.010 0.012

Figure C2.

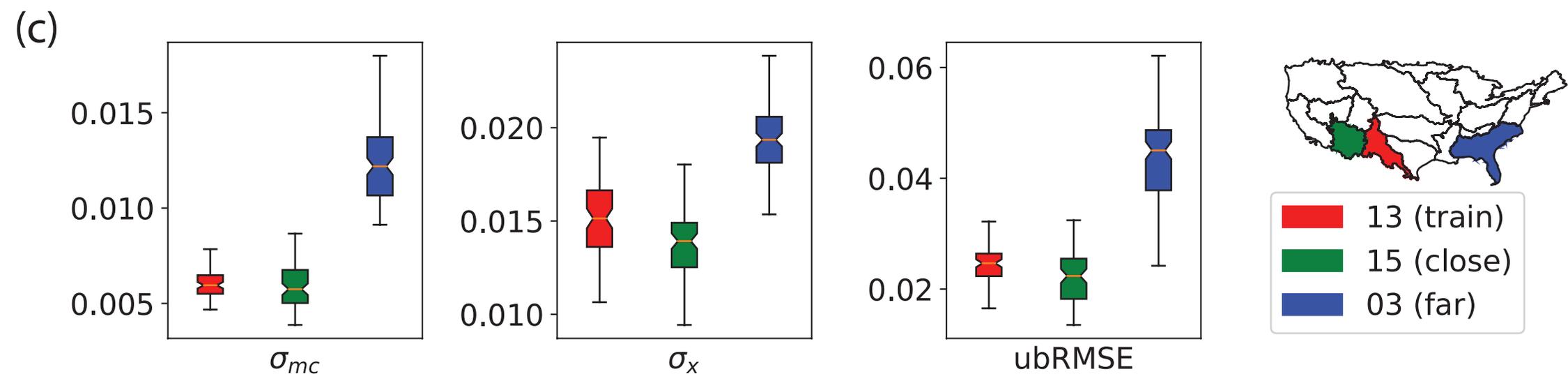
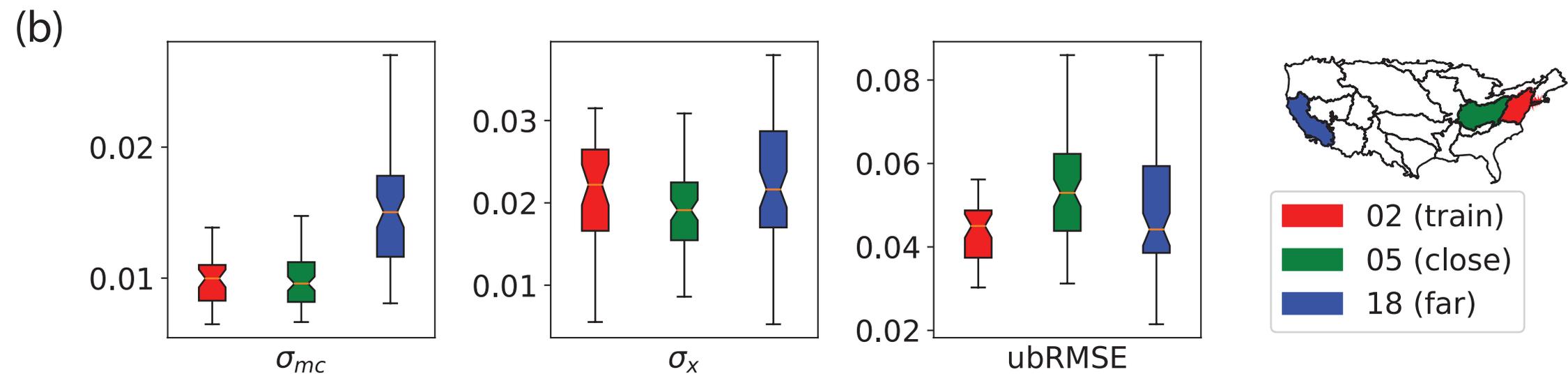
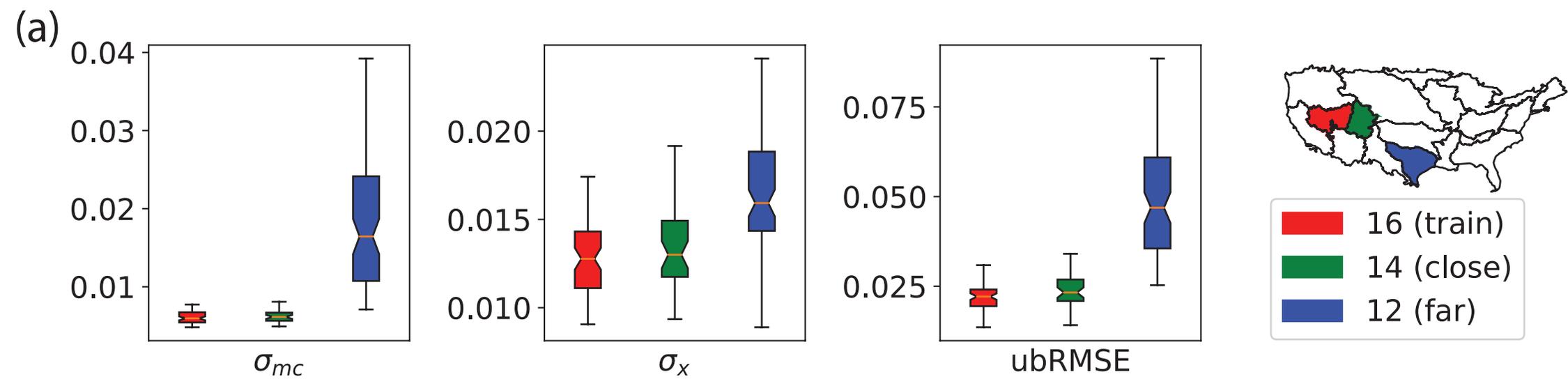
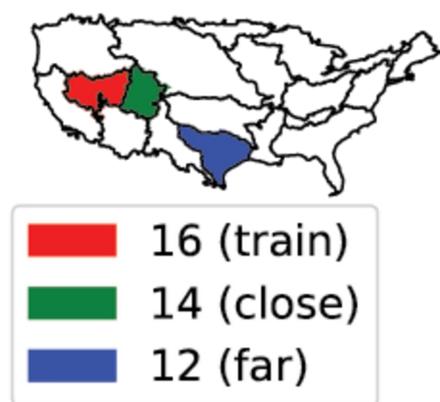
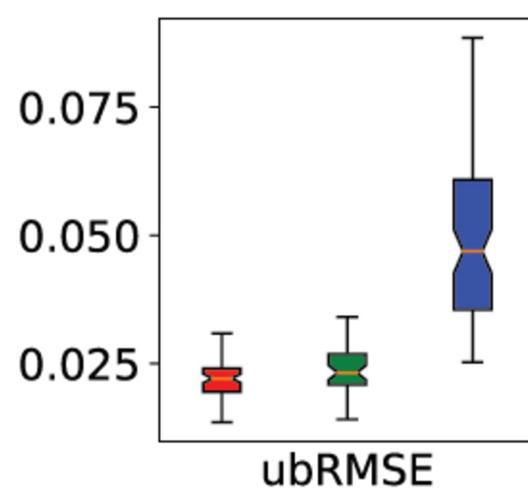
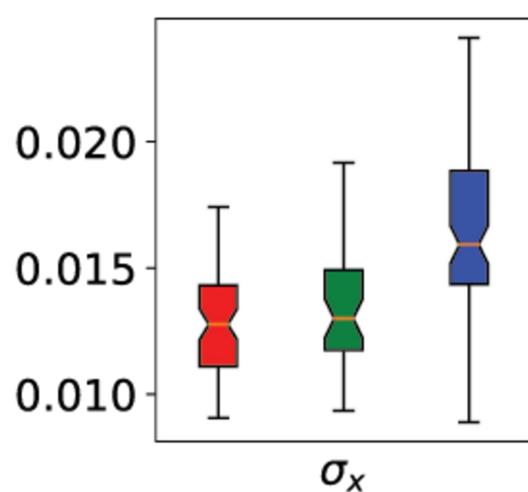
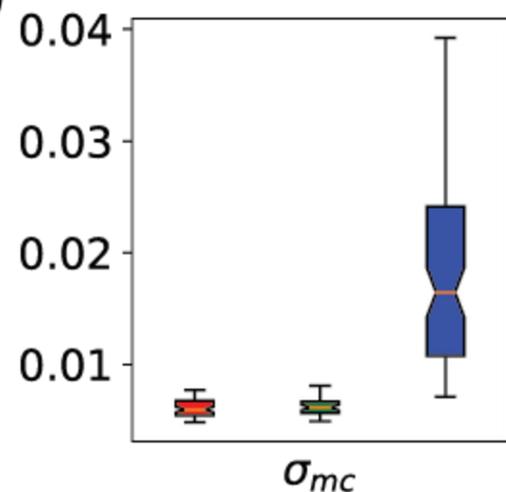
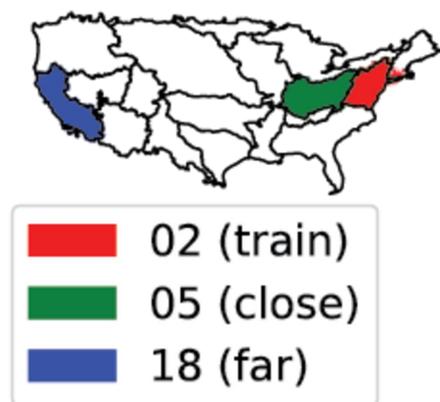
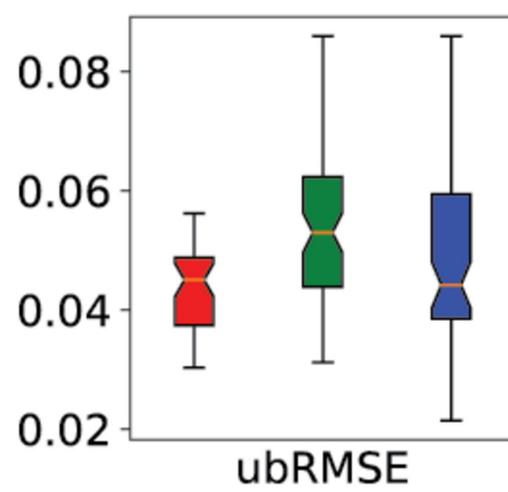
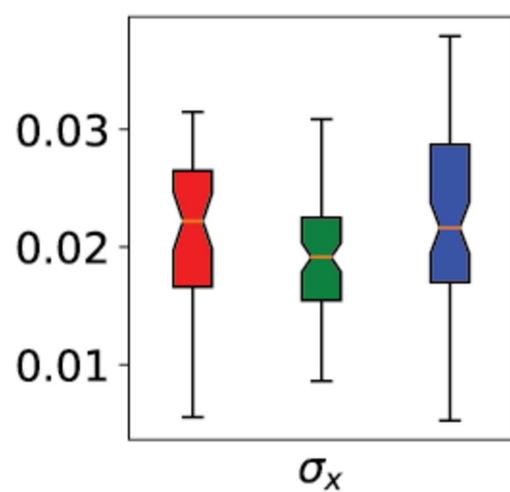
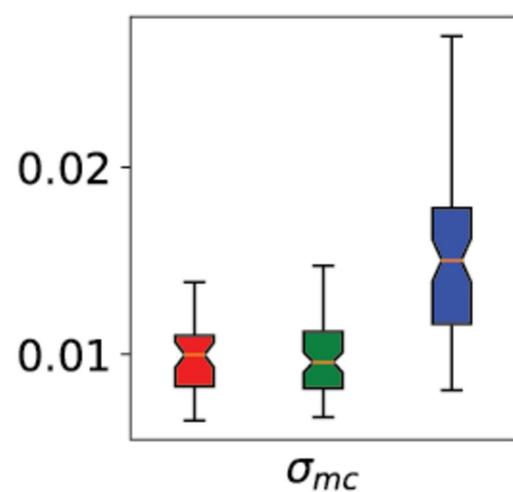


Figure C2 png ver.

(a)



(b)



(c)

