

A k-Nearest Neighbor Space-Time Simulator with applications to large-scale wind and solar power modelling

Yash Amonkar^{1,2,5,*}, David J Farnham^{3,4}, and Upmanu Lall^{1,2}

¹Columbia Water Center, Columbia University, New York, New York, USA - 10027

²Department of Earth and Environmental Engineering, Columbia University, New York, New York, USA - 10027

³Department of Global Ecology, Carnegie Institution for Science, Stanford, CA, USA

⁴Now at ClimateAi, San Francisco, USA

⁵Lead Author

*Correspondence: yva2000@columbia.edu

Abstract

Summary

We develop and present a k-Nearest Neighbor Space-Time Simulator that accounts for the spatiotemporal dependence in high dimensional hydroclimatic fields (e.g. wind and solar) and can simulate synthetic realizations of arbitrary length. We illustrate how this statistical simulation tool can be used in the context of regional power system planning under a scenario of high reliance on wind and solar generation and when long historical records of wind and solar power generation potential are not available. We show how our simulation model can be used to assess the probability distribution of the severity and duration of energy “droughts” at the network scale that need to be managed by long duration storage or alternate energy sources. We present this estimation of supply side shortages for the Texas Interconnection.

Bigger Picture

A novel statistical simulation model that can produce realistic, synthetic realizations of hydroclimatic fields across a region is developed. This k-nearest neighbor based space time simulator can be applied to single or multiple hydroclimatic fields across a large domain. The algorithm facilitates the estimation of the probability of extreme events that are not necessarily represented in relatively short observational records. We apply this algorithm to wind and solar fields. Many regions plan to integrate more wind and solar generation into the energy grid, increasing

power supply variability that can pose risks of undersupply. This simulation tool facilitates the estimation of the probability of regional wind and solar energy “droughts”, and hence allows for the estimation of the storage needed to achieve desired supply-side reliability.

Maturity level categorization - 3 (Development/Pre-production)

In brief

A novel daily time step stochastic simulator capable of capturing the joint dynamics of multiple high dimensional fields across a large domain is presented. The proposed algorithm’s utility is demonstrated via application to joint wind and solar fields across the Texas Interconnection. The risks of undersupply from infrequent, persistent periods of suppressed wind and solar availability are analyzed.

Highlights

An algorithm capable of generating simulations of high dimensional fields is presented.

An application to wind-solar fields across the Texas Interconnection is demonstrated

The approach captures the joint co-variation between wind-solar fields.

Key words

k-Nearest Neighbors, High Dimensional Simulations, Long Duration Storage, Wind and Solar Generation, Texas Interconnection, Stochastic Simulators, Macro-Energy Systems

1 Introduction

2 Many countries and individual states within the United States are mandating
3 reductions in carbon emissions to mitigate anthropogenic climate change, es-
4 specially from the power sector^{1 2 3 4}. At the same time, the costs of wind and
5 solar electricity generation technologies have declined substantially over the last
6 decade⁵. These two factors are spurring increasing deployment of wind and solar
7 based electricity generation.

8 A target system reliability requirement of 99.97 %⁶ necessitates the ad-
9 dition of energy storage, fossil or hydro power sources or significant overca-
10 pacity to buffer supply variations if there is high penetration of variable solar
11 and wind generation^{7 8}. Studies show future scenarios with wind-heavy and/or
12 solar-heavy grid mixes would need long term and even seasonal storage to cost-
13 effectively meet current reliability standards^{9 10}.

14 Long Duration Storage (LDS), defined as storage needed to meet deficits
15 for duration greater than 10 hours^{11 12}, is one option to economically meet grid
16 reliability targets while relying primarily on wind and solar generation⁹. Many
17 recent macro scale electricity studies focusing on renewable electric grids and
18 economy wide de-carbonization models commonly include LDS and expansion
19 of long-distance transmission capacity to smooth the variation in renewable pro-

duction⁹. Such an approach necessitates proper consideration of the temporal and spatial dependence structure of available wind and solar energy including their cross-dependence.

Given a candidate regional configuration of wind and solar generators, sizing LDS economically for a regional grid requires estimates of the probability of potential energy shortages for different duration along with estimates of the demand profile. The estimation of these probabilities to assure high system reliability requires long data records, potentially over many decades. Collins et al¹³ show the pitfalls of modelling energy systems that rely on variable generation using short data records, and note the substantial impact on European power generation costs due to interannual climate variability. Dowling et al¹¹ analyzed LDS sizing and found that the estimated requirement increased as the record length was increased from 1 to 6 years, emphasizing that long data records are needed to properly estimate LDS requirements. This observation is unsurprising given the low frequency behavior of weather and climate, that is well known to have quasi-periodic modes at seasonal to interannual to decadal time scales^{14 15 16}.

The potential for persistent and long duration solar and wind “droughts” and their potential teleconnection to climate modes was illustrated using several long record stations in the United States¹⁷. The availability of long-historical wind and solar data records, however, is restricted to a few sites, for example, airports in the United States¹⁸. Decades long reanalysis datasets^{19 20 21 22} are consequently used to generate gridded wind and solar data records. This data can be used with deterministic optimization methods to compute reliability, capacity allocation, siting and least cost optimization solutions.

An analysis of 39 years of hourly historical (reanalysis) wind and solar data demonstrated the importance of LDS to reduce costs for a wind-solar based electricity system if high reliability is desired¹¹. A subsequent paper²³ focused on the Western Interconnection and derived the frequency of solar and wind droughts of different durations using a 39 year historical record. They define a drought when the production from a source drops below a specified threshold. However, they do not explicitly consider the stochastic properties of the duration and severity of wind and solar energy droughts. Further, their analysis is limited to what can be extracted from the historical record. A primary goal of our paper is to provide a stochastic analysis capability to assess the probability of the severity and duration of the aggregate supply side energy shortages across a region with both wind and solar generators. In other words, our goal is to develop a flexible methodology capable of estimating the exceedance probability (including its uncertainty) for any wind and/or solar generation shortage event of any given duration and severity and for any portfolio distribution of wind and solar collectors over a domain.

While the instrumental data themselves encode the space-time dependence structure which arises due to seasonality, geography and other climate variations, a finite record is basically a sample or realization from the underlying stochastic process. In this paper we address the challenge of developing a stochastic simulator that can synthetically extend these reanalysis data records

66 while reproducing the space and time dependence structure of the wind and
67 solar fields, so that more reliable estimates of the severity and duration of re-
68 gional wind and solar energy potential and their uncertainty can be computed.
69 The wind and solar data from the ERA-5 reanalysis product¹⁹ are used for the
70 development and testing of a stochastic spatio-temporal model that can provide
71 insights as to the variation of the aggregate energy production from a set of
72 spatially distributed wind and solar generation facilities. We take the Electric
73 Reliability Council of Texas (ERCOT) - Texas Interconnection region²⁴ as a
74 target example to explore the historical record and to demonstrate the perfor-
75 mance of our algorithm. While LDS considerations motivate the use of daily
76 data on potential wind and solar resource, the model can be used to simulate
77 any spatio-temporal data, including climate or environmental fields. An im-
78 plicit assumption in the choice of time scale for the energy application is that
79 chemical batteries help smooth out the sub-daily time-scale shortages¹¹. The
80 daily wind and solar capacity factors are computed at the hourly timescale and
81 then averaged over the entire day. (See the Experimental Procedure section for
82 additional methodological details).

83 Over a large region (e.g., the Texas Interconnection), the wind and solar
84 generation assets are likely to be spatially distributed throughout the region²⁵.
85 Non-homogeneous and non-local space and time correlations in the potential
86 energy production across the assets utilized by a grid operator are possible.
87 The annual and seasonal variation of the daily wind and solar energy potential
88 across 216 grid points using daily averages of wind and solar capacity factors
89 from reanalysis data for our example application to Texas are illustrated in
90 Figures 1 and S1.

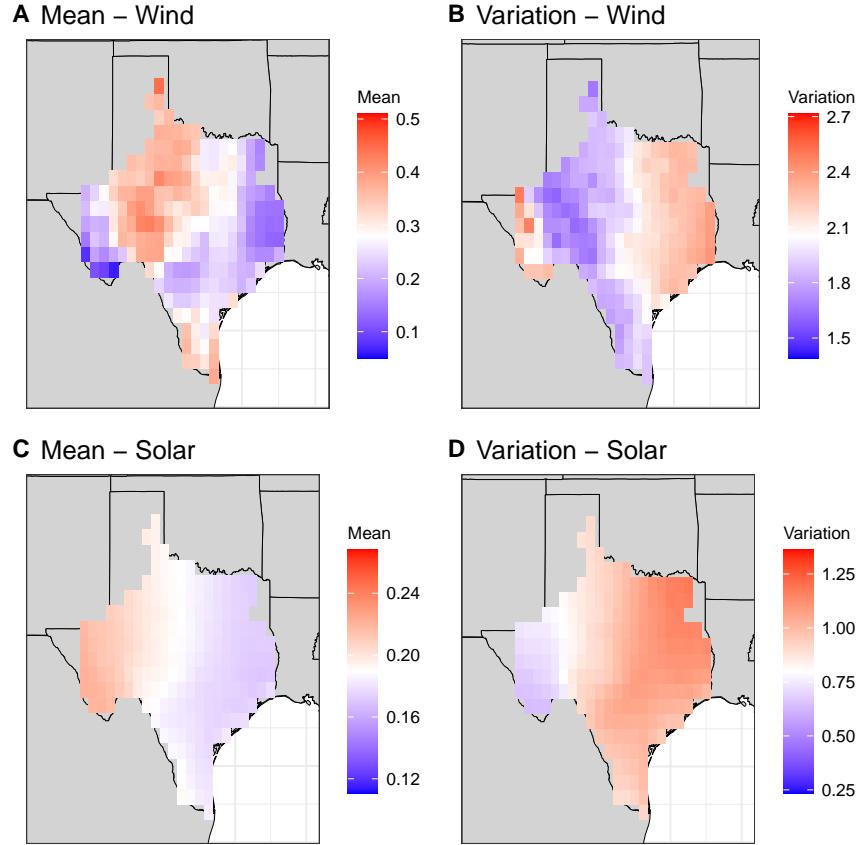


Figure 1: Mean and variation in daily wind and solar capacity factors across the Texas Interconnection. (A) Mean daily wind capacity factors. (B) Variation in daily wind capacity factors. (C) Mean daily solar capacity factors. (D) Variation in daily solar capacity factors. The variation is computed as the difference between the 90th and 10th percentile divided by the mean for each grid.

91 Daily wind and solar fields often exhibit variability that changes by location
 92 and time of year and needs to be accounted for in an analysis of potential
 93 renewable energy droughts or LDS system sizing²⁶. As seen in Figure 1 and S1,
 94 wind and solar along the Gulf of Mexico and the land-area adjoining Louisiana
 95 are regions with relatively low generation potential but with relatively high
 96 variability. The mean wind capacity factor and its variability (Figure 1) is
 97 non-homogeneous. The highest capacity factors are in the north-western and
 98 southern-most portions of the interconnection. The highest variability is in the
 99 eastern portion of the interconnection. Daily wind capacity factors are generally
 100 highest during spring while wind variation is highest during fall and lowest in
 101 summer and spring (Figure S1). The mean daily solar capacity factors and their

102 variability (Figure 1) are more homogeneous and are a function of the season,
103 with low mean radiation and high variability in winter (DJF) and high mean
104 radiation and low variability in the summer (JJA) (Figure S1).

105 The seasonal cross-field spatial correlation between wind and solar is illus-
106 trated in Figure S2 where significant local and non-local spatial correlation
107 structures are evident. The temporal dependence structure explored through
108 the dominant principal component of each field also shows heterogeneity be-
109 tween fields (Figure S3).

110 k-Nearest Neighbor Algorithm

111 We now discuss the historical development and associated literature of the k-
112 nearest neighbor algorithm. The k-nearest neighbor algorithm, a non-parametric
113 method, has been used in traditional problems of classification and regression
114 across fields²⁷. The algorithm serves as a simple first choice in most cases
115 where the underlying data distribution characteristics are not known a priori.
116 The algorithm has its origins in discriminant analysis²⁸. Yakowitz²⁹ and Karlson^{30 31}
117 first developed and utilized a nearest neighbor regression methodology
118 in a time series context for use in rainfall-runoff forecasting. They showed that
119 the method, when used in a time series context has attractive convergence prop-
120 erties, being asymptotically optimal for finite data sets.

121 Lall and Sharma³² developed a nearest neighbor algorithm based simula-
122 tor/resampling scheme for time series data, with applications for hydrological
123 time series. The resampling scheme, referred to as nearest neighbor bootstrap
124 in their work, preserves the dependence in a probabilistic sense, without making
125 any assumptions about the distributional form and marginal densities of the un-
126 derlying process. They also introduced a new resampling kernel to weigh the k
127 successors rather than having uniform weights. They make the assumption that
128 in the space of the nearest neighbors, the local density of the future resampled
129 value can be approximated as a Poisson process. The kernel has the attractive
130 properties of bandwidth and shape adapting to local sampling density changes
131 along with the dimension of the feature vector; and decreases monotonically
132 with distance of the neighbors.

133 Another study³³ introduced a k-nearest neighbors simulator for multivariate
134 time series data following on earlier work³², which was a univariate simulator.
135 The multivariate knn simulation model, a non-parametric approximation of a
136 multivariate lag-1 Markov process, was shown to simulate daily sequences of
137 solar radiation, wind speed, maximum and minimum temperature and precip-
138 itation at a single site. The model simulations preserve the marginal densities
139 of the variables along with the cross-correlations and spell lengths, crucial in-
140 dices for climatological variables. Nowak et al³⁴ developed a disaggregation
141 method which generates multi-site daily flows from a simulated annual value
142 via the knn resampling scheme. While the above described algorithms were
143 all non-parametric, Filho and Lall³⁵ developed a multivariate semi-parametric
144 approach for multi-site streamflow forecasting conditional on external climate
145 predictors using the knn resampling scheme. The key innovation in their work

146 included an adaptive strategy to compute scaling weights for the knn resampling
147 approach, which are the regression coefficients of the external predictors from
148 a parametric regression model. These scaling weights ensure that the relative
149 importance of the predictor vectors is accounted for in the resampling scheme.

150 The structure of the new k-Nearest Neighbors Space Time Simulator (KSTS)
151 algorithm that is presented here is as follows:- A model for temporal variability
152 at each site and for each variable (wind and solar) is considered first. This entails
153 defining a state space through an embedding of the time series. A time series
154 simulation can then be achieved by sequentially drawing from the successors
155 of the k-nearest neighbor of the embedding at each time step, but this will
156 not preserve spatial dependence. Spatial dependence is then introduced by
157 identifying the most likely neighbors of the full spatial field by aggregating
158 neighbor likelihoods for each site/variable. If the state space evolution at two
159 sites is similar (i.e., identified by the same neighbors in time), then the evolution
160 of those two sites would be fully synchronous. Thus the similarity in the selection
161 of neighbors reflects the similarity in dynamics and provides a useful basis for
162 space-time conditioning of a random field's dynamics. The k-nearest neighbors
163 identified across all the sites as the most similar at a given time, are then used to
164 randomly draw a full spatial field for the next time step, using a kernel function
165 that accounts for their degree of similarity through a probability measure. The
166 process is repeated sequentially to generate a time series simulation of the spatial
167 field.

168 The target variables, wind and solar capacity factors, have non-Gaussian
169 skewed distributions and are bounded. The probabilistic sampling using k-
170 nearest neighbors provides an effective approach to sampling from such a non-
171 parametric distribution applied to each target variable. The seasonality in the
172 variables is accounted for by restricting search of k-nearest neighbors using a
173 moving window around the Day of Year (DOY). This method generalizes to a
174 higher dimensional space, the k-Nearest neighbor algorithm³² used for univari-
175 ate or low dimensional multivariate simulations of non-Gaussian and nonlinear
176 dependence that has been used extensively for other climate variables^{33 34 35 36 37}

177 .

178 We apply our new KSTS algorithm to assess the severity, duration, and fre-
179 quency of long duration storage needs associated with aggregate regional energy
180 production. We show that the simulator captures the regional aggregate as well
181 as the site by site probabilities of wind and solar energy potential including the
182 spatial correlation within and across the two fields and the temporal autocor-
183 relation at each site. This study uses the issue of LDS sizing and requirement
184 from the supply-side perspective of renewable energy producers to illustrate the
185 utility of the proposed spatio-temporal field simulation algorithm. We recognize
186 that both supply and demand (load) are needed to assess energy storage needs
187 on a grid level, with the net load (demand – renewables) being of particular
188 interest. As such, our application of the simulation algorithm should be viewed
189 as illustrative but should not be seen as an estimation of the actual LDS needs
190 on the Texas Interconnection. To properly contextualize our application, we
191 consider a target firm energy contract from renewables across the domain, and

192 compute the drought statistics with reference to that contract. We also run a
193 simulation (henceforth termed KNN) that preserves the time series structure
194 but not the spatial structure or the wind-solar dependence. As one may expect,
195 this demonstrates a significant underestimation of the regional LDS probabili-
196 ties. The relative utility and performance of other statistical models relative to
197 the KSTS model and relevant literature review is discussed in the Experimental
198 Procedures section.

199 For the application presented, we use the 71-yr gridded daily wind and so-
200 lar data from the ERA-5¹⁹ reanalysis dataset for 216 sites (grids/nodes) in the
201 Texas Interconnection. Using the KNN or KSTS algorithm one can generate a
202 large number (e.g., 100) of synthetic 71-year simulations (or equivalently a 7100
203 year simulation) of the daily wind and solar fields, without and with spatial de-
204 pendence preserved, respectively. From each simulation we extract the duration
205 and severity of each drought event, which is defined as a shortage in aggregate
206 energy produced across the grid relative to a target threshold. The probabilities
207 of drought severity and duration can then be assessed from this derived set of
208 events. If multiple simulations of 71 years are generated, then one can get also
209 get an estimate of the uncertainty associated with the probability of severity-
210 duration given 71 years of data. If a single long simulation is generated, then
211 we can estimate LDS severity-duration probabilities with reduced uncertainty
212 using the longer synthetic record. While we make inferences on LDS statistics
213 from a purely supply-side perspective, the primary purpose of the example is
214 showing the application of the novel KSTS algorithm to a high dimensional
215 problem of interest.

216 Results

217 We present an evaluation of the severity, duration, and frequency of the aggre-
218 gate energy droughts for the Texas Interconnection with (KSTS) and without
219 (KNN) preserving the spatial structure and wind-solar dependence in simula-
220 tions. For illustrative purposes, a uniform installed capacity allocation of wind
221 and solar generation across all grid points is considered. For both types of sim-
222 ulations, we generated 48 realizations of 71 years of daily wind and solar data
223 at each of the 216 sites. In an actual use case, a stochastic optimization model
224 would use wind and solar capacities reflective of the Texas Interconnection along
225 with the demand to allocate resources and estimate size of LDS capacity using
226 the simulations developed. The results presented here illustrate the importance
227 of getting the space-time dependence right in the simulations for a proper esti-
228 mation of the regional LDS capacity given a candidate spatial configuration of
229 wind and solar generation. Detailed performance statistics of the simulator are
230 presented in the supplement.

²³¹ **Severity, Duration, and Frequency of Energy Droughts**

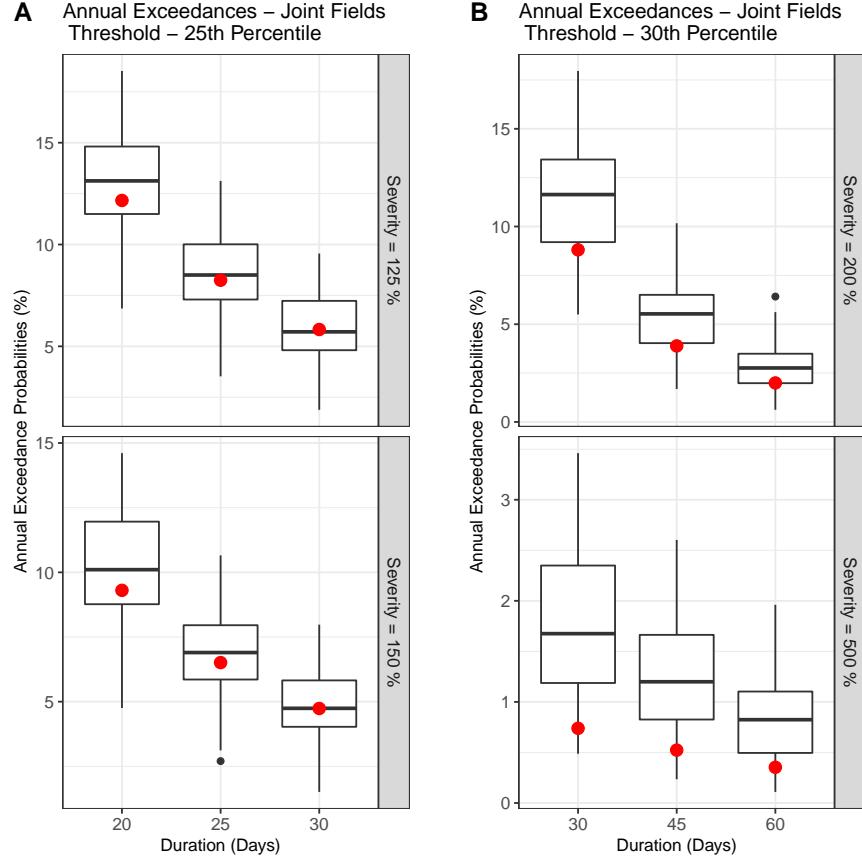


Figure 2: Probability of annual exceedances for energy droughts given a duration and severity with threshold values of (A) 25th percentile and (B) 30th percentile. The red dot denotes the exceedance probability calculated from the reanalysis data. The boxplots denote the uncertainty in the 48 generated simulations using KSTS. The duration is in days and the severity is denoted in terms of percentage of the mean historical reanalysis value. For each box-plot, the thick black horizontal line across the box denotes the median of the annual exceedance probabilities from the simulations and the edges of the box denote the 25th and 75th percentiles, and the lower and upper extents of the vertical lines outside the box denote the 5th and 95th percentiles.

²³² Energy droughts are defined as continuous periods when the daily production
²³³ falls below a target threshold. The threshold value, changing every calendar
²³⁴ day in a year, can be thought of as a forward contract's daily obligation to be
²³⁵ supplied based on the seasonality of the historical reanalysis data. Examples of

such contracts would be where renewable power producers bid in the day ahead market but also buy options from natural gas producers (reliable sources) to hedge their risks in case of lower than anticipated production^{38 39}. The forward contract example in our study is essentially a pre-bid power delivery promise (corresponding to the threshold) and the energy droughts are the periods where the producer will not be able to meet their obligations.

The severity of the drought is the accumulated deficit in production over the duration of the event, i.e., the level of default on a potential contract covering the period, while the duration of the event is the duration during which the deficit exists. Figure 2 (A) shows the annual exceedance probabilities for energy droughts of duration 20, 25 and 30 days with severity of 100% and 150% when the target threshold is the 25th percentile of the distribution of energy that could be produced over that period based on the historical data. The severity of energy droughts was scaled by the mean daily historical production, with a severity of 100% denoting a shortfall equal to the mean daily historical value. The annual exceedance probabilities were computed using local regression (Locfit)⁴⁰ with the number of exceedances regressed against the duration and severity using a Poisson link function. (see Experimental Procedures Section and Supplementary Materials)

The KSTS simulations bracket the exceedance probabilities seen in the reanalysis data (Figure 2). For example, an energy drought with duration over 30 days with a severity of 150% relative to a threshold guaranteeing delivery set at the 25th percentile of daily regional generation, has an annual exceedance probability of ~ 5%, based on the reanalysis data. This corresponds to an event that may be expected to be exceeded once every 20 years. The median exceedance probability from the simulations is quite close to this, but with considerable uncertainty around that value. The 25th to 75th percentiles from the simulations are around 4% to 6% with the 5th and 95th percentiles extending from 2% to 8%, demonstrating the limitations of using solely the original 71 year record for such evaluations.

Results from increasing the target threshold to the 30th percentile of daily regional energy production and looking at higher severity and longer duration droughts are shown in Figure 2 (B). The KSTS simulations bracket the exceedance probabilities seen in the reanalysis data for the severity of 200%. The simulations show higher exceedance probabilities than the data for the 500% case, which is not surprising considering these are rare events with mean annual exceedance probabilities of 0.5-1.5% and thus are difficult to identify given relatively short data records. The severity/duration probabilities from the historical record of 71 years have high uncertainty for events that are rarer than perhaps once every 10 years (annual exceedance probability of 0.1) given this record length^{41 42}. The simulations show that these extreme events could occur far more frequently than would be estimated from just the short historical records. In these illustrations, we consider specific thresholds for supply guarantees, specific drought durations and severity levels, and present the range of probabilities of exceedance from the simulations. In a system design optimization model, for a candidate spatial configuration of generation, the simulator would provide

282 the probability distribution for a candidate LDS capacity that is considered to
283 meet the deficit over a specified duration (e.g., specified by a contract). Alternately,
284 one could also compute the probability distribution of the shortage beyond the candidate LDS to assess potential penalties for non-delivery, if those
285 were considered in the optimization model.

286 Annual exceedance probabilities for different combinations of duration and
287 severity, and for multiple thresholds and wind-solar individual fields are provided
288 in Figures S4 and S6. The entire joint distribution of duration and severity for all
289 energy droughts in the data and the generated simulations relative to a threshold
290 for thresholds at the 25th, 30th, 35th, and 40th percentile are shown in Figure
291 S5. We see that KSTS is effective for representing the range of energy droughts.
292 Similar boxplot estimates for the KNN algorithm generated simulations are not
293 shown since the simulations show no occurrences of energy droughts at these
294 thresholds.
295

296 **KSTS Reproduces the Aggregate Generation**

297 The simulations from both KSTS and KNN reproduce temporal dynamics and
 298 data characteristics across both wind and solar fields at individual sites. The
 299 moments (mean and standard deviation), minimum and maximum for individual
 300 sites in KSTS and KNN simulations are representative of the underlying data
 301 (Figure S7). Both simulators are able to reproduce the quantiles (Figure S8,
 302 Figure S9), underlying probability distribution (Figure S10), auto-correlation
 303 structure (Figure S11), and site-level seasonality (Figure S12). The distribution
 304 of the aggregate generation over the full domain, however, is properly repro-
 305 duced by the KSTS simulator but not by the KNN simulator.

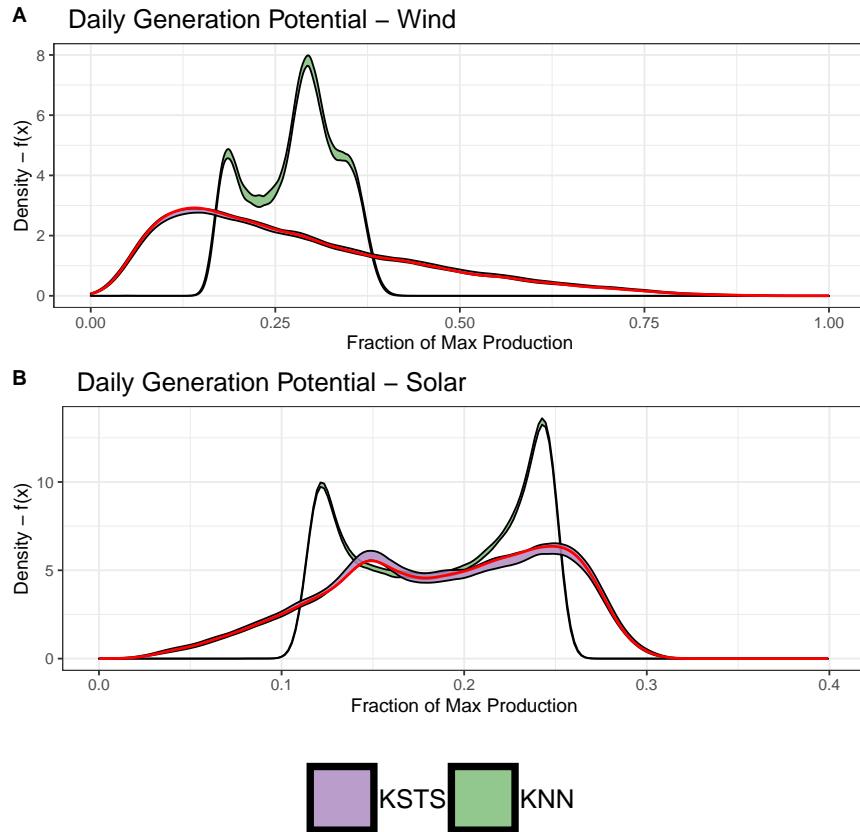


Figure 3: Kernel density estimate / Probability density function (PDF) of the daily aggregated energy production across the Texas Interconnection simulated using KSTS (purple) and KNN (green). The red line denotes the observed data pdf. The purple and green regions show the mid 90th (5th-95th) percentile interval regions from the individual pdfs computed from 48 simulations from each simulator. (A) Wind. (B) Solar.

306 The kernel density estimate of aggregated daily energy generation potential
307 across the Texas Interconnection is shown in Figure 3 for the historical reanalysis
308 record (black), and for the KSTS (purple) and KNN (green) simulations.
309 The degree to which adequate consideration of the spatial dependence and the
310 wind-solar correlation leads to a proper representation of the potential for en-
311 ergy production is illustrated through the fidelity of the KSTS simulations to
312 the density function from the observations, and the marked departure of the
313 KNN based simulations. It is clear that modeling spatial and cross field depen-
314 dence is important to get the right frequency of the tail events (i.e., for LDS
315 probabilities), even if the site-level production is adequately simulated without
316 considering spatial dependence.

³¹⁷ **KSTS Reproduces Cross-Field Dependence**

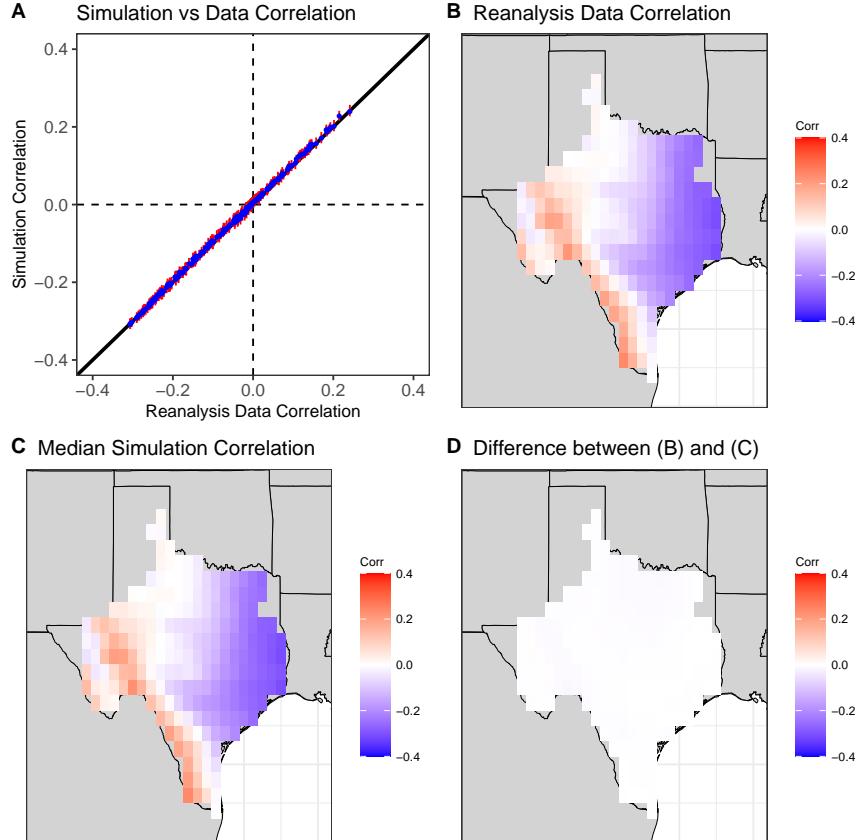


Figure 4: Pearson correlation between wind and solar at each grid point based on simultaneous simulations of wind and solar using KSTS. (A) Simulation correlation vs reanalysis data correlation where the red lines denote the mid 90th (5th-95th) percentile range and the blue dots denote the median value in the simulation spread generated using the KSTS method. (B) Map of the grid-wise correlations in the reanalysis data record. (C) Map of the grid-wise median simulation correlations using KSTS. (D) Map of the difference between (B) and (C).

³¹⁸ From Figure 4 we note that the grid-wise correlation between wind and solar
³¹⁹ across ERCOT is well reproduced by the KSTS simulations, which are based on
³²⁰ simultaneous modeling of the wind and solar fields. By comparison, the KNN
³²¹ simulations do not exhibit grid-wise wind-solar correlations consistent with the
³²² reanalysis data (Figure S13). Furthermore, the spatial correlation structure
³²³ across all grids within a field for both wind and solar is also well reproduced by

324 the KSTS simulations unlike the KNN simulations (Figure S14). The seasonal
325 variation in the correlation between wind and solar is also well modeled by the
326 KSTS algorithm (Figure S15 and Figure S16).

327 Discussion

328 The primary contribution of this paper is the presentation of a *k*-Nearest Neighbor
329 Space Time Simulator (KSTS) and its application to the joint wind-solar
330 fields across the Texas Interconnection. We demonstrate the importance of using
331 a stochastic simulator that can properly reproduce the marginal probability
332 densities of wind and solar at each site, as well as the cross-field spatial depen-
333 dence structure if good estimates of the severity-duration and frequency of long
334 duration renewable energy droughts, are of interest. These resource droughts
335 are analyzed from a purely supply side perspective in this study with demand
336 (load) and installed solar and wind capacity data needed for further detailed
337 analysis. This algorithm seeks to estimate the probability (and associated un-
338 certainty) of the duration and severity of resource droughts integrated over the
339 spatial domain, through simulation. So far, much of the development of re-
340 newable electricity sources has focused on local microgrids, but there has been
341 growing interest in national and regional grids⁴³. As the scale is increased, there
342 is evidence that LDS is an effective and economic component of the design of
343 these regional systems^{10 11 44}. However, most of the models developed and ap-
344 plied at these scales are deterministic and use relatively short records with a
345 potential to lead to biased results⁴⁵. They do not consider the possible con-
346 tracting structures for guaranteed delivery and the associated default penalties.
347 The probabilities of the severity and duration of defaults as well as the pen-
348 alties and LDS costs would ultimately determine economically optimal resource
349 allocations. We anticipate and are planning to develop stochastic simulation-
350 optimization models to address a range of questions associated with such designs
351 and contracts. The KSTS simulator is motivated by this context, and it was
352 important to understand how critical it is to model spatial dependence when
353 assessing the characteristics of energy shortages on a grid.

354 From the application to the Texas Interconnection, we note that there is
355 substantial seasonal variability in the spatial expression of potential wind and
356 solar resource. This is not a surprise. The point by point wind-solar correlation
357 varies substantially by location and by season, as does the spatial correlation
358 structure for wind and solar and their cross-dependence. If these factors are
359 ignored, then the resulting regional LDS probability distributions are compro-
360 mised quite significantly. These simulations show large uncertainties in the
361 annual exceedance probabilities for the severity, duration, and threshold combi-
362 nations considered, as well as potentially higher exceedance probabilities than
363 computed from the 71 year data record for the more extreme severity, duration,
364 and threshold combinations.

365 The KSTS simulator is nonparametric and is appropriate for this setting
366 where the target variables are bounded with non-Gaussian distributions with

367 space and time dependence across variables changing by season. Since KSTS is
368 based on sampling the observed data, it can be thought of as a spatio-temporal
369 bootstrap procedure, where a spatio-temporal kernel is used at each time step
370 to sample a historical field with probabilities determined by the kernel and
371 a distance metric applied to the temporal state space for each variable. The
372 temporal sequences of potential energy produced at each site and across the
373 region are different even though the individual daily values are resampled from
374 the historical record. This allows the analysis of the range of drought severity,
375 duration, and frequency using an extended sample.

376 **Applicability to Other Problems**

377 The KSTS algorithm could be used for any spatio-temporal simulation problem
378 where the preservation of spatial dependence is of interest, and the temporal
379 dynamics are modeled through a Markovian process or through a time domain
380 embedding, as illustrated in the methodology section. Typical examples would
381 be any weather or climate fields where maintaining the space and time consis-
382 tency across multiple variables is of interest. An example that is similar to the
383 current context is a copula based model that was developed to model risk of
384 national livestock losses in Mongolia using spatially distributed livestock loss
385 data over time⁴⁶. Many of the existing space-time simulators were developed
386 in a Markovian framework with random variables considered to be drawn from
387 the Exponential family of distributions.

388 Extension of the KSTS simulator to other time scales (e.g., hourly) is feasible.
389 An hourly simulator would need to consider the diurnal cycle, in addition to
390 the seasonal cycle, and we are exploring computationally efficient strategies for
391 an algorithm that can address this while maintaining spatial and cross field
392 dependence.

393 The KSTS simulator can also be applied to simultaneously modelling mul-
394 tiple streamflow or weather stations across a watershed while preserving the
395 internal dependence structure. Such streamflow data exhibit spatio-temporal
396 correlation patterns due to their position in the river network, altitude and a
397 host of other hydrological variables³⁴, making application of KSTS attractive.

398 **Limitations and Next Steps**

399 Since KSTS is a hybrid resampling (bootstrap) method, it cannot simulate
400 values not seen in the historical record. This is not a major issue for wind and
401 solar capacity factors, since the lower and upper extremities of the distribution
402 for both wind and solar are recorded in the reanalysis (historical) data enabling
403 KSTS to generate daily simulations that span the entire distribution of both
404 fields.

405 In the general case of other hydroclimatic variables, extrapolation to val-
406 ues not seen in the historical record is also possible⁴⁶. If a parametric or non-
407 parametric marginal probability distribution is fit to the time series of a variable,
408 with parameters that may vary by season, one could draw observations from that

409 distribution that are consistent with the k-nearest neighbor value selected for
410 simulation. If the rank (small to big) of the k-nearest neighbor value in the his-
411 torical data is j , then an estimate of its corresponding cumulative distribution
412 function $F(x)$ is $j/(n+1)$, where n is the sample size⁴⁷. Accounting for uncer-
413 tainty, one can consider that $F(x)$ lies between $(j-0.5)/(n+1)$ and $(j+0.5)/(n+1)$.
414 For the largest/smallest value on record the intervals would be $((n-0.5)/(n+1),$
415 1) and $(0,1.5/(n+1))$, respectively. Consequently, if sampling values not seen
416 in the historical record is of interest, one can first sample uniformly from this
417 interval and then sample the corresponding value from the marginal distribu-
418 tion of x . This does not change the basic structure of the KSTS algorithm, but
419 allows values to be simulated from an appropriate probability distribution for
420 each variable considered.

421 The KSTS simulator exploits the similarity in the temporal evolution across
422 the fields and grid points. The potential next step would be developing an
423 algorithm which is capable of capturing the heterogeneity in dynamics across
424 even larger regions. This becomes important when the spatial scale of the
425 simulation is expanded from the Texas Interconnection to either the Western
426 or Eastern Interconnection or the entire North American continent. Such a
427 large scale makes it more likely that the wind and solar availability in some
428 sub-regions is driven by disparate atmospheric dynamics and consequently their
429 temporal evolution structure would be heterogeneous when compared to just
430 Texas.

431 **Experimental Procedures**

432 **Resource Availability**

433 **Lead Contact**

434 Further information and requests for resources and materials should be directed
435 to Yash Amonkar yva2000@columbia.edu

436 **Materials Availability**

437 This study did not generate new unique materials.

438 **Data and Code Availability**

439 The KSTS and KNN generated simulations use wind and solar data spanning 71-
440 yrs (1950-2020) across the Texas Interconnection and are taken from the ERA-5
441 reanalysis dataset¹⁹, which can be accessed publicly. All code used in this study
442 is made publicly available on Github at <https://github.com/yashamonkar/LDS-Inferences>.

444 **Wind and Solar Data**

445 The ERA-5 reanalysis variables used are wind speeds at 100 meter altitude and
446 downward surface solar radiation¹⁹. The spatial grid size of the data is set at
447 0.5° lat $\times 0.5^\circ$ lon and contains 216 grid points across the Texas Interconnection
448 domain (Figure S17). The wind speed and solar radiation at each hour are
449 converted to hourly wind and solar power respectively, and the daily capacity
450 factors are computed as the average across 24 hours.

451 Wind power is estimated by converting the 100 m wind speed to wind power
452 using the wind turbine power curve from a V90-2.0MW Vestas turbine (as shown
453 in Figure S18). The data are converted to the daily time step by taking the
454 mean of the hourly capacity factors for each day and the dataset spans the 71
455 years from January 1st, 1950 to December 31st, 2020.

456 The solar variable is the downward surface solar radiation (W/m^2) and is
457 converted to capacity factor at the hourly level by accounting for the dependence
458 of photovoltaic performance on temperature⁴⁸ (Figure S19). We then compute
459 a capacity factor for each day by taking the mean of the 24 hourly values for
460 that day.

461 **Energy Deficits and Drought**

462 The daily energy deficit is defined as the daily deviation below a percentile
463 threshold for that day of year (DOY) for each site. The deviation could be
464 positive if that day's value is greater than the selected threshold percentile
465 value or negative if it is lower. The daily energy deviation across the field is
466 computed by aggregating the daily site deviation and is given by,

$$y_t = \sum_{i=1}^n (x_{i,t} - \widetilde{x}_{i,T})$$

467 where, y_t is the aggregated daily energy deviation at day t ; $x_{i,t}$ is the nor-
468 malized wind or solar value at site i and day t ; $\widetilde{x}_{i,T}$ is the normalized DOY
469 percentile based on the selected threshold for site i and day DOY(t); n is the
470 total number of grid points (216) times the fields (wind and solar). The ag-
471 gregated deviation y_t can take a positive (surplus) or negative (deficit) value
472 on any particular day, while the cumulative deficit, the variable of interest is
473 computed as,

$$z_1 = \max(0, -y_1)$$
$$z_t = \max(0, z_{t-1} - y_t)$$

474 where, z_t and y_t are the cumulative deficit and daily deviation at day t
475 respectively. While y_t can either be positive or negative, the cumulative deficit
476 takes a lower value of 0 (surplus) and is restricted to positive values (periods of
477 energy deficit). Energy Droughts for a selected threshold percentile are defined
478 to occur during instances of consecutive days with positive values of cumulative

479 deficit. Severity of a drought event is defined as the maximum cumulative deficit
480 during the drought period, while the duration is the spell length in days.

481 Annual Exceedance Probability

482 The previous section is used to compute the duration and severity for all energy
483 droughts in the data and the generated simulations. The number of exceedances
484 (e_i) for each drought i include all drought events in the data record (or individual
485 simulation realizations) having a greater severity and greater duration than
486 event i , which are computed as,

$$487 C(e_i) = \sum_{j=1}^n (d_i > d_j) \cap \sum_{j=1}^n (s_i > s_j) \quad (1)$$

487 where, $C(e_i)$ is the count of exceedances for drought event i with duration
488 d_i and severity s_i and n is the total number of drought events. The count
489 of exceedances $C(e_i)$ is regressed against the severity s_i and duration d_i using
490 Poisson regression. The methodology used is local regression using the locfit
491 package⁴⁰.

492 After the model fitting process, the count of exceedances $C(e_t)$ is estimated
493 using the fitted model for the required duration d_t and severity s_t for a desired
494 drought event t . The number of years of the record (yr) is then used to scale
495 the number of exceedances to get the annual exceedance probability (p) using
496 the formula:-

$$497 p_t = \frac{C(e_t) \times 100}{yr} \quad (2)$$

497 where, p_t is the annual exceedance percentage for a drought event t with
498 severity s_t and duration d_t .

499 Fitting Other Models

500 We considered and tested other strategies for spatio-temporal simulation with
501 the Texas Interconnection data prior to developing and testing the KSTS algo-
502 rithm. A brief review of those efforts is presented below. The Autoregressive
503 Integrated Moving Average (ARIMA) model was first fit to sites individually,
504 using the Akaike Information Criterion to select model order⁴⁹. The results,
505 not displayed here, failed to capture the underlying data generating process
506 with significant departures even from the base moments. ARIMA and other
507 similar parametric approaches assume normality of the underlying distributions
508 making them a bad fit to this joint modelling problem where both wind and
509 solar capacity factors are non-normal, bounded and multi-modal distributions
510 (Figure S10). The serial dependence structure of the wind and solar data is
511 also nonlinear and the use of these linear models contributes to biases in the

512 simulations. Generalizations such as Vector Autoregressive processes and Space-
513 Time autoregressive models suffer from the same problems, in addition to the
514 challenge of fitting a high dimensional covariance matrix.

515 Another potential class of non-parametric machine learning based models ap-
516 plicable to the current problem are Generative Adversarial Networks (GANs)⁵⁰.
517 GANs have been used for renewable simulations (scenario generation) and do
518 not assume normality of the underlying data^{51 52}. However, while GANs can
519 model complex spatial dependencies, they require a large amount of temporal
520 data to fit the model. Our initial efforts at fitting GANs to the current data did
521 not lead to a model that had skillful temporal evolution characteristics, making
522 a direct comparison infeasible.

523 Finally, different types of Hidden Markov Models (HMM)⁵³ were also ex-
524 plored to simulate the wind-solar fields. The application of a non-homogeneous
525 HMM with spatial covariance modeled using variograms led to unsatisfactory
526 results due to inadequate representation of the spatio-temporal correlation struc-
527 ture.

528 Simulator Hyperparameters

529 For both KSTS and KNN, the seasonality in the data is accounted for by re-
530 stricting the search of nearest neighbors to a ± 30 days moving window across
531 the years around the day of the year (DOY). The number of nearest neighbors
532 (k) selected is approximately \sqrt{n} . With 71 years and 61 days per year, \sqrt{n} is
533 ~ 65 , where n is the number of possible candidate neighbors after accounting
534 for the moving window³². A lag-1 dependence structure for the state space is
535 assumed. 48 independent simulation realizations, each of the same length as
536 the reanalysis data (71 yrs = 25933 days), are generated using the KSTS and
537 KNN algorithms. The KNN algorithm is fit to each grid individually using the
538 hyper-parameters specified above with the algorithm which is outlined in Lall
539 and Sharma³².

540 k-Nearest Neighbor Space-Time Simulator (KSTS)

541 The general structure and a cartoon example application of the KSTS algorithm
542 is illustrated in Figures S20 and 5 respectively. The algorithm leads to a space-
543 time simulation process that is Markovian (or corresponds to a state space
544 formed by the embedding) in time.

545 KSTS Algorithm

546 Step 1:- Define the composition of the state space $D_{i,t}$.

547 Define a state space $D_{i,t}$ of dimension m which is the number of embedding
548 delay lags. The state space can be a single lag, multiple lags and/or disjoint
549 lags allowing for custom time dependencies. The embedding selected for the
550 simulator application could be,

551 Case 1 $D_{i,t} := (x_{t-1}, x_{t-2}); m = 2$

552 Case 2 $D_{i,t} \text{ :- } (x_{t-\tau}, x_{t-2\tau}, x_{t-\phi}, x_{t-2\phi}); m = 4, \tau = 1, \phi = 12$
 553 Case 3 $D_{i,t} \text{ :- } (x_{t-1}, x_{t-4} x_{t-7}); m = 3$

554 Case 1 represents simple dependence on the two previous values. Case 2
 555 represents dependence on the past two values and values 12 and 24 steps be-
 556 fore the current value allowing for monthly and interannual dependence for
 557 monthly data. Case 3 represents incorporation of a temporal dependence struc-
 558 ture unique to the data. The state space $D_{i,t}$ is defined for each time series at
 559 site i and time t , whereas $D_{i,T}$ are all the historic vectors which correspond to
 560 the selected embedding structure for site i .

561 **Step 2:- Compute the k-nearest neighbors for all sites at time t .**

562 At time step t and site i using the current state space vector $D_{i,t}$, identify
 563 the k -nearest neighbors using the weighted Euclidean distance measure,

$$r_{i,t} = \left(\sum_{j=1}^m w_j ([D_{i,t}]_j - [D_{i,T}]_j)^2 \right)^{1/2}$$

564 where, $[D_{i,t}]_j$ and $[D_{i,T}]_j$ are the j^{th} components of $D_{i,t}$ and $D_{i,T}$ respectively
 565 and w_j are the weights assigned to each of the embedding lags j . This is repeated
 566 for all sites. The ordered set of time indices which correspond to the k nearest
 567 neighbors (as defined by the euclidean distances stored in $r_{i,t}$) of site i at time
 568 t are stored in $\tau_{i,t}$.

569 **Step 3:- Compute resampling probabilities for k nearest neighbor in-
 570 dices using a discrete kernel p_j .**

$$p_j = \frac{1/j}{\sum_{j=1}^k 1/j}$$

571 where p_j is the resampling probability for the j th element (time instance of
 572 the j th nearest neighbor of $D_{i,t}$) in $\tau_{i,t}$. The resampling kernel stays the same
 573 across all time t and across all sites, and is pre-computed and stored prior to
 574 simulation. It is a function of the number of neighbors k and not the distances.

575

576 **Step 4:- Define $T_{i,t}$ and similarity vector S_t for time t .**

577 Define $T_{i,t}$ as a matrix where the rows and columns correspond to the sites and
 578 unique time indices from the historical data respectively. The columns record
 579 the resampling probabilities associated with the time indices for the k -nearest
 580 neighbors in $\tau_{i,t}$ for each site i , with values being 0 for other time indices. The
 581 similarity vector S_t is then defined as the sum of all elements in each column in
 582 $T_{i,t}$.

$$S_t = \sum_{i=1}^s T_{i,t}$$

583 where s is the total number of sites. The similarity vector S_t has the same
584 length as the number of unique time indices in the data.

585 **Step 5:- Curtail and scale the similarity vector S_t .**

586 The similarity vector S_t is ordered and curtailed to its highest k values.
587 The time indices associated the k highest values of S_t are selected as the k-
588 nearest neighbor candidates for the entire spatial field. The probabilities of the
589 associated k neighbors are scaled to add up to 1.

$$[S_t]_j = \frac{[S_t]_j}{\sum_{j=1}^k [S_t]_j}$$

590

591 **Step 6:- Re-sample the full spatial field for time $t + 1$.**

592 Using the discrete probability mass function S_t , sample a single value and
593 re-sample entire fields across all sites from the time index which corresponds
594 to the next time step of selected value in S_t as data for the simulation at time
595 $t + 1$. Return to Step 2 if further time-steps are needed for the simulation.

596 Refer supplementary materials for further details on the algorithm and hyper-
597 parameter selection.

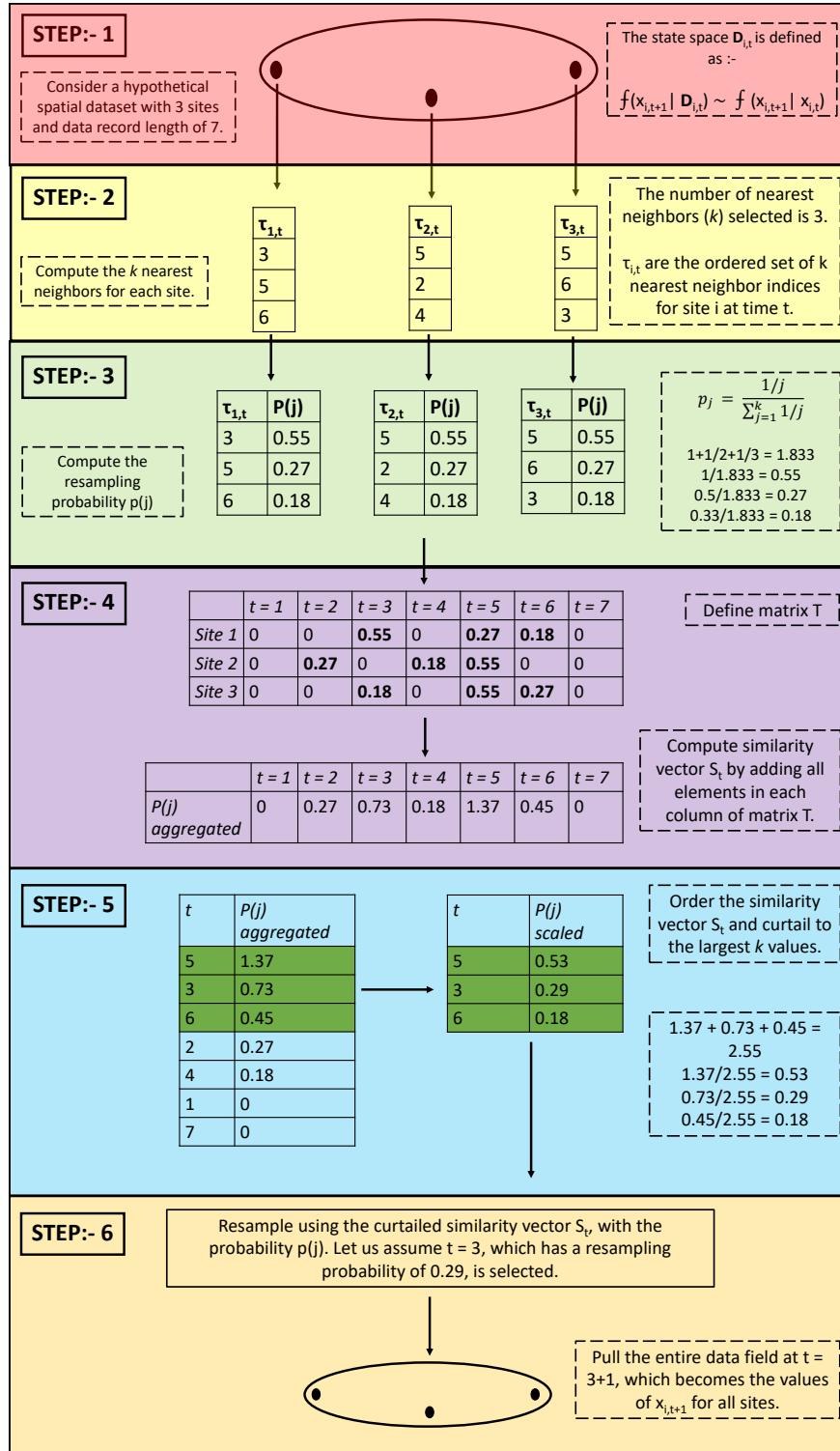


Figure 5: Cartoon example application²³ of the KSTS algorithm to a spatial dataset consisting of 3 grids/sites and data record (time) length of 7.

598 **Acknowledgment**

599 Y.A acknowledges support from the Cheung-Kong Innovation Doctoral Fellow-
600 ship. D.J.F. was supported by a gift from Gates Ventures LLC to the Carnegie
601 Institution for Science.

602 **Author Contributions**

603 Y.A developed the code and performed the computations. Y.A and D.J.F de-
604 signed the analysis, conceived experiments and simulation checks with supervi-
605 sion from U.L. who introduced the algorithm. D.J.F provided the data. Y.A
606 took the lead in writing the manuscript with all authors discussing and con-
607 tributing to the final manuscript.

608 **Declaration of Interests**

609 The authors declare no competing interests.

610 **References**

- 611 [1] California Legislative Information. Senate Bill -100 California
612 Renewables Portfolio Standard Program: emissions of greenhouse
613 gases., 2018. URL [https://leginfo.legislature.ca.gov/faces/
614 billNavClient.xhtml?bill_id=201720180SB100](https://leginfo.legislature.ca.gov/faces/billNavClient.xhtml?bill_id=201720180SB100).
- 615 [2] Jeff Deyette. States March toward 100% Clean Energy—Who's
616 Next?, August 2019. URL [https://blog.ucsusa.org/jeff-deyette/
617 states-march-toward-100-clean-energy-whos-next](https://blog.ucsusa.org/jeff-deyette/states-march-toward-100-clean-energy-whos-next). Section: Energy.
- 618 [3] European Comission. REGULATION OF THE EUROPEAN PARLIA-
619 MENT AND OF THE COUNCIL establishing the framework for achiev-
620 ing climate neutrality and amending Regulation (EU) 2018/1999 (European
621 Climate Law), 2020. URL [https://eur-lex.europa.eu/legal-content/
622 EN/TXT/?qid=1588581905912&uri=CELEX:52020PC0080](https://eur-lex.europa.eu/legal-content/EN/TXT/?qid=1588581905912&uri=CELEX:52020PC0080).
- 623 [4] New York State Legislator. New York's Climate Leadership and Commu-
624 nity Protection Act (CLCPA), 2019. URL <https://climate.ny.gov/>.
- 625 [5] International Renewable Energy Agency. Renewable Power Generation
626 Costs in 2019. Technical report, 2020. URL [/publications/2020/Jun/
627 Renewable-Power-Costs-in-2019](https://publications/2020/Jun/Renewable-Power-Costs-in-2019).
- 628 [6] North American Electric Reliability Corporation. (NERC). 2012 State of
629 Reliability. Technical report, 2012. URL https://www.nerc.com/files/2012_sor.pdf.
- 630

- 631 [7] Marc Beaudin, Hamidreza Zareipour, Anthony Schellenberglabe, and
632 William Rosehart. Energy storage for mitigating the variability
633 of renewable electricity sources: An updated review. *Energy for*
634 *Sustainable Development*, 14(4):302–314, December 2010. ISSN 0973-0826.
635 doi: 10.1016/j.esd.2010.09.007. URL <https://www.sciencedirect.com/science/article/pii/S0973082610000566>.
- 637 [8] Jacques Després, Silvana Mima, Alban Kitous, Patrick Criqui, Nouredine
638 Hadjsaid, and Isabelle Noirot. Storage as a flexibility option in power
639 systems with high shares of variable renewable energy sources: a POLES-
640 based analysis. *Energy Economics*, 64:638–650, May 2017. ISSN 0140-9883.
641 doi: 10.1016/j.eneco.2016.03.006. URL <https://www.sciencedirect.com/science/article/pii/S0140988316300445>.
- 643 [9] Jesse D. Jenkins, Max Luke, and Samuel Thernstrom. Getting to Zero Car-
644 bon Emissions in the Electric Power Sector. *Joule*, 2(12):2498–2510, Decem-
645 ber 2018. ISSN 2542-4351. doi: 10.1016/j.joule.2018.11.013. URL <https://www.sciencedirect.com/science/article/pii/S2542435118305622>.
- 647 [10] Matthew R. Shaner, Steven J. Davis, Nathan S. Lewis, and Ken Caldeira.
648 Geophysical constraints on the reliability of solar and wind power in the
649 United States. *Energy & Environmental Science*, 11(4):914–925, April 2018.
650 ISSN 1754-5706. doi: 10.1039/C7EE03029K. URL <https://pubs.rsc.org/en/content/articlelanding/2018/ee/c7ee03029k>. Publisher: The
651 Royal Society of Chemistry.
- 653 [11] Jacqueline A. Dowling, Katherine Z. Rinaldi, Tyler H. Ruggles, Steven J.
654 Davis, Mengyao Yuan, Fan Tong, Nathan S. Lewis, and Ken Caldeira.
655 Role of Long-Duration Energy Storage in Variable Renewable Electric-
656 ity Systems. *Joule*, 4(9):1907–1928, September 2020. ISSN 2542-4351.
657 doi: 10.1016/j.joule.2020.07.007. URL <https://www.sciencedirect.com/science/article/pii/S2542435120303251>.
- 659 [12] ARPA-E. Duration Addition to electricY Storage, 2018. URL <https://arpa-e.energy.gov/technologies/programs/days>.
- 661 [13] Seán Collins, Paul Deane, Brian Ó Gallachóir, Stefan Pfenninger, and
662 Iain Staffell. Impacts of Inter-annual Wind and Solar Variations on
663 the European Power System. *Joule*, 2(10):2076–2090, October 2018.
664 ISSN 2542-4351. doi: 10.1016/j.joule.2018.06.020. URL <https://www.sciencedirect.com/science/article/pii/S254243511830285X>.
- 666 [14] Luc Bonnafous, Upmanu Lall, and Jason Siegel. A water risk index for
667 portfolio exposure to climatic extremes: conceptualization and an applica-
668 tion to the mining industry. *Hydrology and Earth System Sciences*, 21(4):
669 2075–2106, April 2017. ISSN 1027-5606. doi: 10.5194/hess-21-2075-2017.
670 URL <https://hess.copernicus.org/articles/21/2075/2017/>. Pub-
671 lisher: Copernicus GmbH.

- 672 [15] Shaleen Jain and Upmanu Lall. Floods in a changing cli-
 673 mate: Does the past represent the future? *Water Resources*
 674 *Research*, 37(12):3193–3205, 2001. ISSN 1944-7973. doi:
 675 <https://doi.org/10.1029/2001WR000495>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2001WR000495>. eprint:
 676 <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2001WR000495>.
- 678 [16] James Doss-Gollin, David J. Farnham, Scott Steinschneider, and
 679 Upmanu Lall. Robust Adaptation to Multiscale Climate Vari-
 680 ability. *Earth's Future*, 7(7):734–747, 2019. ISSN 2328-4277.
 681 doi: <https://doi.org/10.1029/2019EF001154>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019EF001154>. eprint:
 682 <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2019EF001154>.
- 684 [17] David J. Farnham. *Identifying and Modeling Spatio-temporal Structures*
 685 *in High Dimensional Climate and Weather Datasets with Applications to*
 686 *Water and Energy Resource Management*. PhD thesis, Columbia Univer-
 687 sity, 2018. URL <https://doi.org/10.7916/D8321CTB>.
- 688 [18] Scott Chamberlain. 'NOAA' Weather Data from R [R package rnoaa ver-
 689 sion 1.3.2], February 2021. URL <https://CRAN.R-project.org/package=rnoaa>. Publisher: Comprehensive R Archive Network (CRAN).
- 691 [19] Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András
 692 Horányi, Joaquín Muñoz-Sabater, Julien Nicolas, Carole Peubey, Raluca
 693 Radu, Dinand Schepers, Adrian Simmons, Cornel Soci, Saleh Abdalla,
 694 Xavier Abellan, Giampaolo Balsamo, Peter Bechtold, Gionata Biavati,
 695 Jean Bidlot, Massimo Bonavita, Giovanna De Chiara, Per Dahlgren,
 696 Dick Dee, Michail Diamantakis, Rossana Dragani, Johannes Flemming,
 697 Richard Forbes, Manuel Fuentes, Alan Geer, Leo Haimberger, Sean
 698 Healy, Robin J. Hogan, Elías Hólm, Marta Janisková, Sarah Keeley,
 699 Patrick Laloyaux, Philippe Lopez, Cristina Lupu, Gabor Radnoti, Pa-
 700 tricia de Rosnay, Iryna Rozum, Freja Vamborg, Sébastien Villaume,
 701 and Jean-Noël Thépaut. The ERA5 global reanalysis. *Quarterly*
 702 *Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020.
 703 ISSN 1477-870X. doi: <https://doi.org/10.1002/qj.3803>. URL <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3803>. eprint:
 704 <https://rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/qj.3803>.
- 706 [20] Patrick Laloyaux, Eric de Boisseson, Magdalena Balmaseda, Jean-
 707 Raymond Bidlot, Stefan Broennimann, Roberto Buizza, Per Dal-
 708 hgren, Dick Dee, Leopold Haimberger, Hans Hersbach, Yuki
 709 Kosaka, Matthew Martin, Paul Poli, Nick Rayner, Elke Ruste-
 710 meier, and Dinand Schepers. CERA-20C: A Coupled Reanaly-
 711 sis of the Twentieth Century. *Journal of Advances in Modeling*
 712 *Earth Systems*, 10(5):1172–1195, 2018. ISSN 1942-2466. doi:
 713 <https://doi.org/10.1029/2018MS001273>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001273>. eprint:
<https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2018MS001273>.

- 714 onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001273. _eprint:
715 <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2018MS001273>.
- 716 [21] Ronald Gelaro, Will McCarty, Max J. Suárez, Ricardo Todling, An-
717 drea Molod, Lawrence Takacs, Cynthia A. Randles, Anton Darmenov,
718 Michael G. Bosilovich, Rolf Reichle, Krzysztof Wargan, Lawrence Coy,
719 Richard Cullather, Clara Draper, Santha Akella, Virginie Buchard, Austin
720 Conaty, Arlindo M. da Silva, Wei Gu, Gi-Kong Kim, Randal Koster,
721 Robert Lucchesi, Dagmar Merkova, Jon Eric Nielsen, Gary Partyka,
722 Steven Pawson, William Putman, Michele Rienecker, Siegfried D. Schu-
723 bert, Meta Sienkiewicz, and Bin Zhao. The Modern-Era Retrospective
724 Analysis for Research and Applications, Version 2 (MERRA-2). *Journal*
725 *of Climate*, 30(14):5419–5454, July 2017. ISSN 0894-8755, 1520-0442.
726 doi: 10.1175/JCLI-D-16-0758.1. URL [https://journals.ametsoc.org/](https://journals.ametsoc.org/view/journals/clim/30/14/jcli-d-16-0758.1.xml)
727 [view/journals/clim/30/14/jcli-d-16-0758.1.xml](https://journals.ametsoc.org/view/journals/clim/30/14/jcli-d-16-0758.1.xml). Publisher: Ameri-
728 can Meteorological Society Section: Journal of Climate.
- 729 [22] D. P. Dee, S. M. Uppala, A. J. Simmons, P. Berrisford, P. Poli,
730 S. Kobayashi, U. Andrae, M. A. Balmaseda, G. Balsamo, P. Bauer, P. Bech-
731 told, A. C. M. Beljaars, L. van de Berg, J. Bidlot, N. Bormann, C. Delsol,
732 R. Dragani, M. Fuentes, A. J. Geer, L. Haimberger, S. B. Healy, H. Hers-
733 bach, E. V. Hólm, L. Isaksen, P. Kållberg, M. Köhler, M. Matricardi, A. P.
734 McNally, B. M. Monge-Sanz, J.-J. Morcrette, B.-K. Park, C. Peubey, P. de
735 Rosnay, C. Tavolato, J.-N. Thépaut, and F. Vitart. The ERA-Interim re-
736 analysis: configuration and performance of the data assimilation system.
737 *Quarterly Journal of the Royal Meteorological Society*, 137(656):553–597,
738 2011. ISSN 1477-870X. doi: <https://doi.org/10.1002/qj.828>. URL <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.828>. _eprint:
739 <https://rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/qj.828>.
- 740 [23] Katherine Z. Rinaldi, Jacqueline A. Dowling, Tyler H. Ruggles, Ken
741 Caldeira, and Nathan S. Lewis. Wind and Solar Resource Droughts in
742 California Highlight the Benefits of Long-Term Storage and Integration
743 with the Western Interconnect. *Environmental Science & Technology*, 55
744 (9):6214–6226, May 2021. ISSN 0013-936X. doi: 10.1021/acs.est.0c07848.
745 URL <https://doi.org/10.1021/acs.est.0c07848>. Publisher: American
746 Chemical Society.
- 747 [24] Electric Reliability Council of Texas. About ERCOT, 2021. URL <http://www.ercot.com/about>.
- 748 [25] Electric Reliability Council of Texas. Impact of increased wind resrouces in
749 the ERCOT region. Technical report, June 2020. URL http://www.ercot.com/content/wcm/lists/200196/Wind_One_Pager_June_2020.pdf.
- 750 [26] Andrew Kumler, Ignacio Losada Carreño, Michael T. Craig, Bri-Mathias
751 Hodge, Wesley Cole, and Carlo Brancucci. Inter-annual variability
752 of wind and solar electricity generation and capacity values in Texas.

- 756 Environmental Research Letters, 14(4):044032, April 2019. ISSN 1748-
 757 9326. doi: 10.1088/1748-9326/aaf935. URL <https://doi.org/10.1088/1748-9326/aaf935>. Publisher: IOP Publishing.
- 759 [27] T. Cover and P. Hart. Nearest neighbor pattern classification. *IEEE
 760 Transactions on Information Theory*, 13(1):21–27, January 1967. ISSN
 761 1557-9654. doi: 10.1109/TIT.1967.1053964. Conference Name: IEEE
 762 Transactions on Information Theory.
- 763 [28] B. W. Silverman and M. C. Jones. E. Fix and J.L. Hodges (1951): An Im-
 764 portant Contribution to Nonparametric Discriminant Analysis and Den-
 765 sity Estimation: Commentary on Fix and Hodges (1951). *International
 766 Statistical Review / Revue Internationale de Statistique*, 57(3):233–238,
 767 1989. ISSN 0306-7734. doi: 10.2307/1403796. URL <https://www.jstor.org/stable/1403796>. Publisher: [Wiley, International Statistical Institute
 768 (ISI)].
- 770 [29] S. Yakowitz. Nearest-Neighbour Methods for Time Series Analysis. *Journal of Time Series Analysis*, 8(2):235–247, 1987. ISSN 1467-9892. doi:
 771 10.1111/j.1467-9892.1987.tb00435.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9892.1987.tb00435.x>.
 772 _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9892.1987.tb00435.x>.
- 776 [30] M. Karlsson and S. Yakowitz. Nearest-neighbor methods for
 777 nonparametric rainfall-runoff forecasting. *Water Resources
 778 Research*, 23(7):1300–1308, 1987. ISSN 1944-7973. doi:
 779 10.1029/WR023i007p01300. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/WR023i007p01300>.
 780 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/WR023i007p01300>.
- 782 [31] M. Karlsson and S. Yakowitz. Rainfall-runoff forecasting methods, old and
 783 new. *Stochastic Hydrology and Hydraulics*, 1(4):303–318, December 1987.
 784 ISSN 1435-151X. doi: 10.1007/BF01543102. URL <https://doi.org/10.1007/BF01543102>.
- 786 [32] Upmanu Lall and Ashish Sharma. A Nearest Neighbor
 787 Bootstrap For Resampling Hydrologic Time Series. *Water
 788 Resources Research*, 32(3):679–693, 1996. ISSN 1944-7973. doi:
 789 <https://doi.org/10.1029/95WR02966>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/95WR02966>.
 790 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/95WR02966>.
- 792 [33] Balaji Rajagopalan and Upmanu Lall. A k-nearest-neighbor sim-
 793 ulator for daily precipitation and other weather variables. *Water
 794 Resources Research*, 35(10):3089–3101, 1999. ISSN 1944-7973. doi:
 795 <https://doi.org/10.1029/1999WR900028>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/1999WR900028>.
 796 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/1999WR900028>.

- 798 [34] Kenneth Nowak, James Prairie, Balaji Rajagopalan, and Upmanu Lall.
 799 A nonparametric stochastic approach for multisite disaggregation of
 800 annual to daily streamflow. *Water Resources Research*, 46(8), 2010.
 801 ISSN 1944-7973. doi: 10.1029/2009WR008530. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2009WR008530>.
 802 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2009WR008530>.
- 804 [35] Francisco Assis Souza Filho and Upmanu Lall. Seasonal to in-
 805 terannual ensemble streamflow forecasts for Ceara, Brazil: Ap-
 806 plications of a multivariate, semiparametric algorithm. *Water*
 807 *Resources Research*, 39(11), 2003. ISSN 1944-7973. doi:
 808 <https://doi.org/10.1029/2002WR001373>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2002WR001373>.
 809 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2002WR001373>.
- 811 [36] James Prairie, Balaji Rajagopalan, Upmanu Lall, and Terrance Fulp.
 812 A stochastic nonparametric technique for space-time disaggregation of
 813 streamflows. *Water Resources Research*, 43(3), 2007. ISSN 1944-7973.
 814 doi: <https://doi.org/10.1029/2005WR004721>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2005WR004721>.
 815 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2005WR004721>.
- 817 [37] James Prairie, Kenneth Nowak, Balaji Rajagopalan, Upmanu Lall,
 818 and Terrance Fulp. A stochastic nonparametric approach for stream-
 819 flow generation combining observational and paleoreconstructed data.
 820 *Water Resources Research*, 44(6), 2008. ISSN 1944-7973. doi:
 821 <https://doi.org/10.1029/2007WR006684>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2007WR006684>.
 822 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2007WR006684>.
- 824 [38] Nayara Aguiar, Vijay Gupta, and Pramod P. Khargonekar. A Real Op-
 825 tions Market-Based Approach to Increase Penetration of Renewables. *IEEE*
 826 *Transactions on Smart Grid*, 11(2):1691–1701, March 2020. ISSN 1949-
 827 3061. doi: 10.1109/TSG.2019.2942258. Conference Name: IEEE Transac-
 828 tions on Smart Grid.
- 829 [39] David D’Achiardi, Nayara Aguiar, Stefanos Baros, Vijay Gupta, and
 830 Anuradha M. Annaswamy. Reliability Contracts Between Renewable
 831 and Natural Gas Power Producers. *IEEE Transactions on Control of*
 832 *Network Systems*, 6(3):1075–1085, September 2019. ISSN 2325-5870. doi:
 833 10.1109/TCNS.2019.2919857. Conference Name: IEEE Transactions on
 834 Control of Network Systems.
- 835 [40] Catherine Loader. Local Regression, Likelihood and Density Estima-
 836 tion [R package locfit version 1.5-9.4], March 2020. URL <https://CRAN.R-project.org/package=locfit>. Publisher: Comprehensive R Archive
 837 Network (CRAN).

- 839 [41] Gary D. Tasker. Effective record length for the T-year event. *Journal*
840 of *Hydrology*, 64(1):39–47, July 1983. ISSN 0022-1694. doi: 10.1016/
841 0022-1694(83)90059-8. URL <https://www.sciencedirect.com/science/article/pii/0022169483900598>.
- 843 [42] Robert Link, Thomas B. Wild, Abigail C. Snyder, Mohamad I. Hejazi,
844 and Chris R. Vernon. 100 years of data is not enough to establish
845 reliable drought thresholds. *Journal of Hydrology X*, 7:100052, April
846 2020. ISSN 2589-9155. doi: 10.1016/j.hydroa.2020.100052. URL <https://www.sciencedirect.com/science/article/pii/S2589915520300031>.
- 848 [43] Patricia J. Levi, Simon Davidsson Kurland, Michael Carabajales-Dale,
849 John P. Weyant, Adam R. Brandt, and Sally M. Benson. Macro-Energy
850 Systems: Toward a New Discipline. *Joule*, 3(10):2282–2286, October
851 2019. ISSN 2542-4351. doi: 10.1016/j.joule.2019.07.017. URL <https://www.sciencedirect.com/science/article/pii/S2542435119303617>.
- 853 [44] Micah S. Ziegler, Joshua M. Mueller, Gonçalo D. Pereira, Juhyun Song,
854 Marco Ferrara, Yet-Ming Chiang, and Jessika E. Trancik. Storage Re-
855 quirements and Costs of Shaping Renewable Energy Toward Grid De-
856 carbonization. *Joule*, 3(9):2134–2153, September 2019. ISSN 2542-4351.
857 doi: 10.1016/j.joule.2019.06.012. URL <https://www.sciencedirect.com/science/article/pii/S2542435119303009>.
- 859 [45] Leonard Göke and Mario Kendziora. The adequacy of time-series reduc-
860 tion for renewable energy systems. *arXiv:2101.06221 [econ, q-fin]*, January
861 2021. URL <http://arxiv.org/abs/2101.06221>. arXiv: 2101.06221.
- 862 [46] Upmanu Lall, Naresh Devineni, and Yasir Kaheil. An Empirical,
863 Nonparametric Simulator for Multivariate Random Variables
864 with Differing Marginal Densities and Nonlinear Dependence
865 with Hydroclimatic Applications. *Risk Analysis*, 36(1):
866 57–73, 2016. ISSN 1539-6924. doi: 10.1111/risa.12432. URL
867 <https://onlinelibrary.wiley.com/doi/abs/10.1111/risa.12432>.
868 _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/risa.12432>.
- 869 [47] Lasse Makkonen. Bringing Closure to the Plotting Position Controversy.
870 *Communications in Statistics - Theory and Methods*, 37(3):460–467, Jan-
871 uary 2008. ISSN 0361-0926. doi: 10.1080/03610920701653094. URL
872 <https://doi.org/10.1080/03610920701653094>. Publisher: Taylor &
873 Francis _eprint: <https://doi.org/10.1080/03610920701653094>.
- 874 [48] Philip E. Bett and Hazel E. Thornton. The climatological relationships
875 between wind and solar energy supply in Britain. *Renewable Energy*,
876 87:96–110, March 2016. ISSN 0960-1481. doi: 10.1016/j.renene.2015.
877 10.006. URL <https://www.sciencedirect.com/science/article/pii/S0960148115303591>.

- 879 [49] Rob J. Hyndman and George Athanasopoulos. Forecasting: principles and
880 practice. OTexts, May 2018. ISBN 978-0-9875071-1-2. Google-Books-ID:
881 _bBhDwAAQBAJ.
- 882 [50] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-
883 Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative
884 adversarial networks. Communications of the ACM, 63(11):139–144, Oc-
885 tober 2020. ISSN 0001-0782. doi: 10.1145/3422622. URL <https://doi.org/10.1145/3422622>.
- 887 [51] Yize Chen, Yishen Wang, Daniel Kirschen, and Baosen Zhang. Model-Free
888 Renewable Scenario Generation Using Generative Adversarial Networks.
889 IEEE Transactions on Power Systems, 33(3):3265–3275, May 2018. ISSN
890 1558-0679. doi: 10.1109/TPWRS.2018.2794541. Conference Name: IEEE
891 Transactions on Power Systems.
- 892 [52] Yize Chen, Xiyu Wang, and Baosen Zhang. An Unsupervised Deep Learn-
893 ing Approach for Scenario Forecasts. In 2018 Power Systems Computation
894 Conference (PSCC), pages 1–7, June 2018. doi: 10.23919/PSCC.2018.
895 8442500.
- 896 [53] L. Rabiner and B. Juang. An introduction to hidden Markov models.
897 IEEE ASSP Magazine, 3(1):4–16, January 1986. ISSN 1558-1284. doi:
898 10.1109/MASSP.1986.1165342. Conference Name: IEEE ASSP Magazine.

899 Supplementary Materials

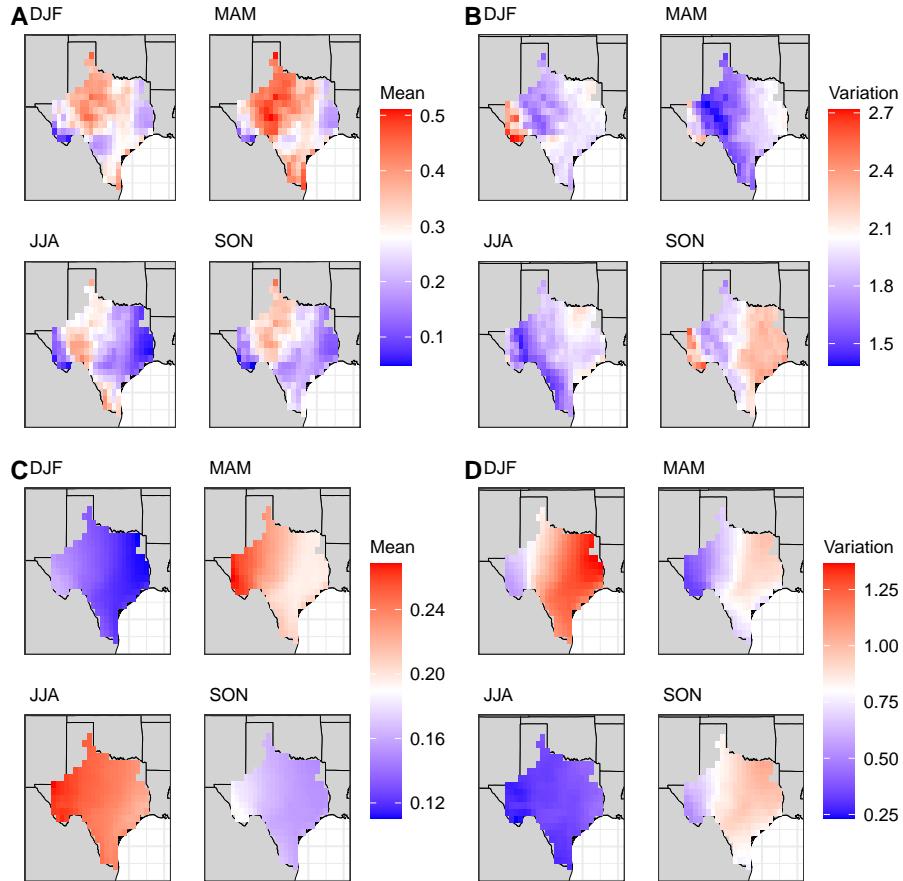


Figure S1: Seasonal mean and variation in daily wind and solar capacity factors across the Texas Interconnection. (A) Mean daily wind capacity factors by season. (B) Variation in daily wind capacity factors by season. (C) Mean daily solar capacity factors by season. (D) Variation in daily solar capacity factors by season. The seasonal variation is computed as the difference between the 90th and 10th percentile divided by the mean for each grid point for each season. The sub-plots are arranged as follows :- top left - Dec-Jan-Feb (DJF), top right - Mar-Apr-May (MAM), bottom left - Jun-Jul-Aug (JJA), bottom right - Sept-Oct-Nov (SON).

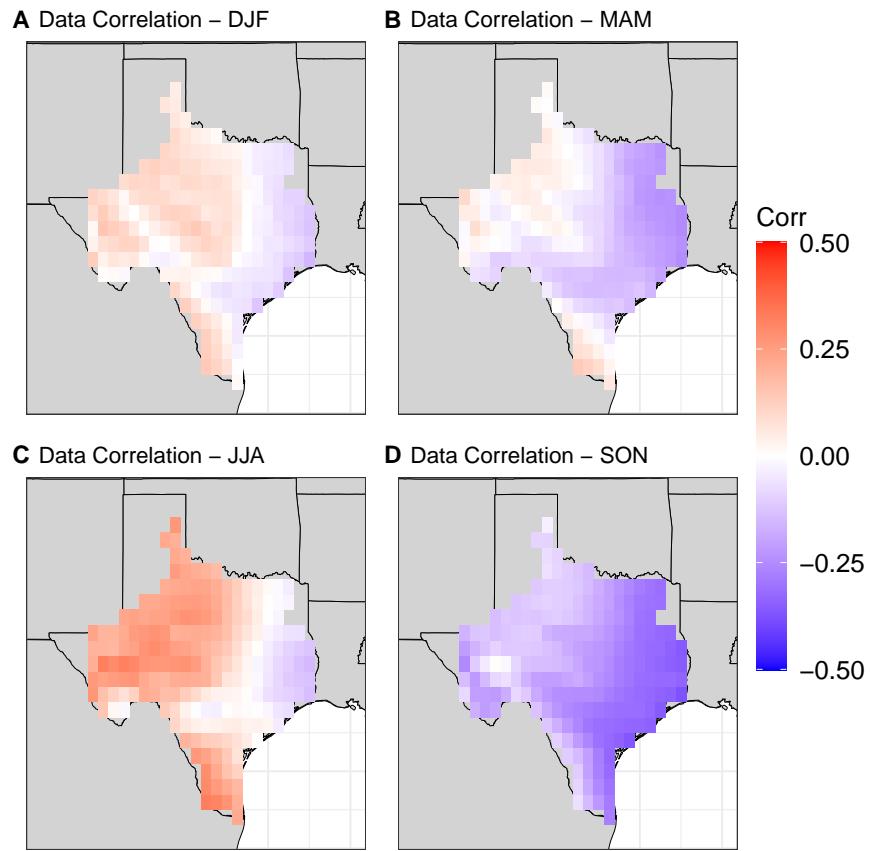


Figure S2: Seasonal correlation between daily wind and solar capacity factors in the ERA-5 reanalysis dataset at each grid point. (A) Dec-Jan-Feb (DJF). (B) Mar-Apr-May (MAM). (C) Jun-Jul-Aug (JJA). (D) Sep-Oct-Nov (SON). The correlations are computed using Pearson's method.

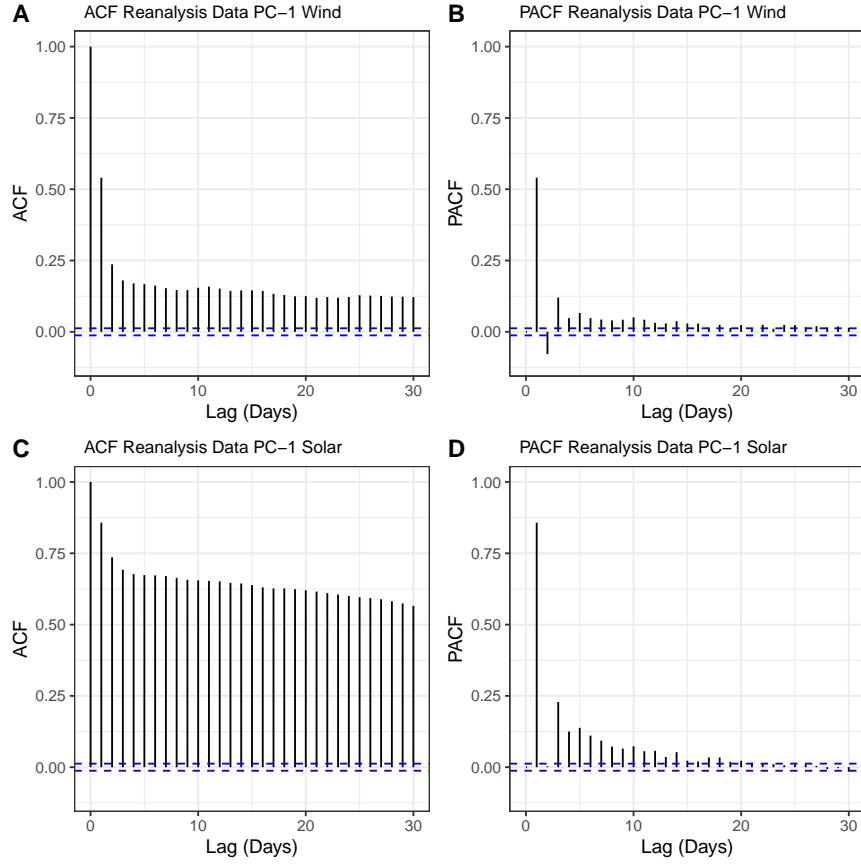


Figure S3: The auto-correlation (ACF) and partial auto-correlation (PACF) of the first Principal Component(PC) of the individual reanalysis data fields. (A) Wind PC-1 ACF. (B) Wind PC-1 PACF. (C) Solar PC-1 ACF. (D) Solar PC-1 PACF The fractional variance explained by PC-1 for solar and wind fields is 79 % and 63 % respectively. The blue dashed line denotes the significance level for the record length of 25933 days (71 yrs).

900 Annual Exceedance for Threshold-Duration-Severity

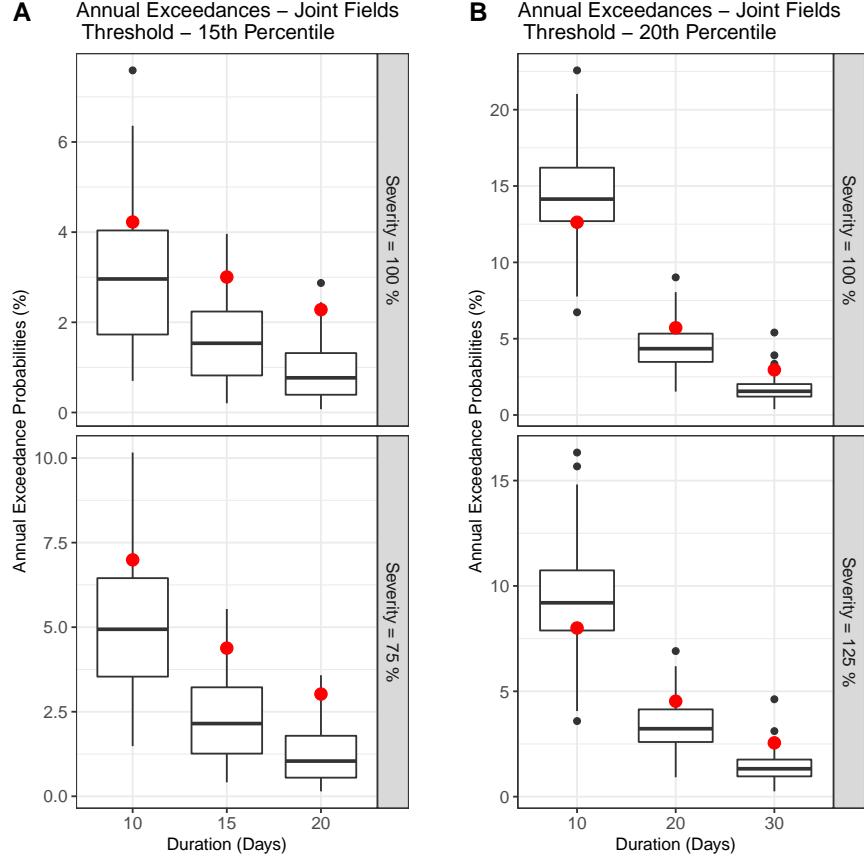


Figure S4: Probability of annual exceedances for energy droughts given a duration and severity with threshold values of (A) 15th percentile, and (B) 20th percentile. The red dot denotes the exceedance probability calculated from the reanalysis data. The boxplots denote the uncertainty in the 48 generated simulations using KSTS. The duration is in days and the severity is denoted in terms of percentage of the mean historical reanalysis value. For each box-plot, the thick black horizontal line across the box denotes the median of the annual exceedance probabilities from the simulations and the edges of the box denote the 25th and 75th percentiles, and the lower and upper extents of vertical lines outside the box denote the 5th and 95th percentiles.

₉₀₁ Severity vs Duration plots for different thresholds

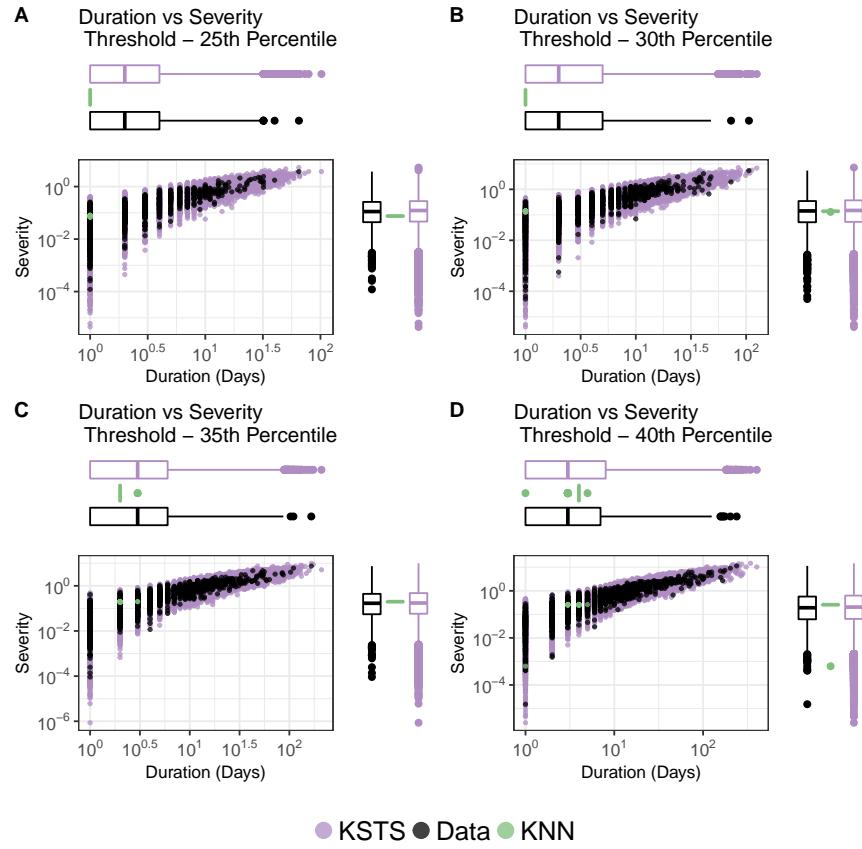


Figure S5: Duration versus Severity for all (wind and solar aggregated) energy droughts with marginal distributions (in boxplots for both variables) for the data (black), KSTS (purple) and KNN (green) simulations using threshold values of (A) 25th percentile, (B) 30th percentile, (C) 35th percentile, and (D) 40th percentile.

902 Energy Droughts for Wind and Solar Fields Individually

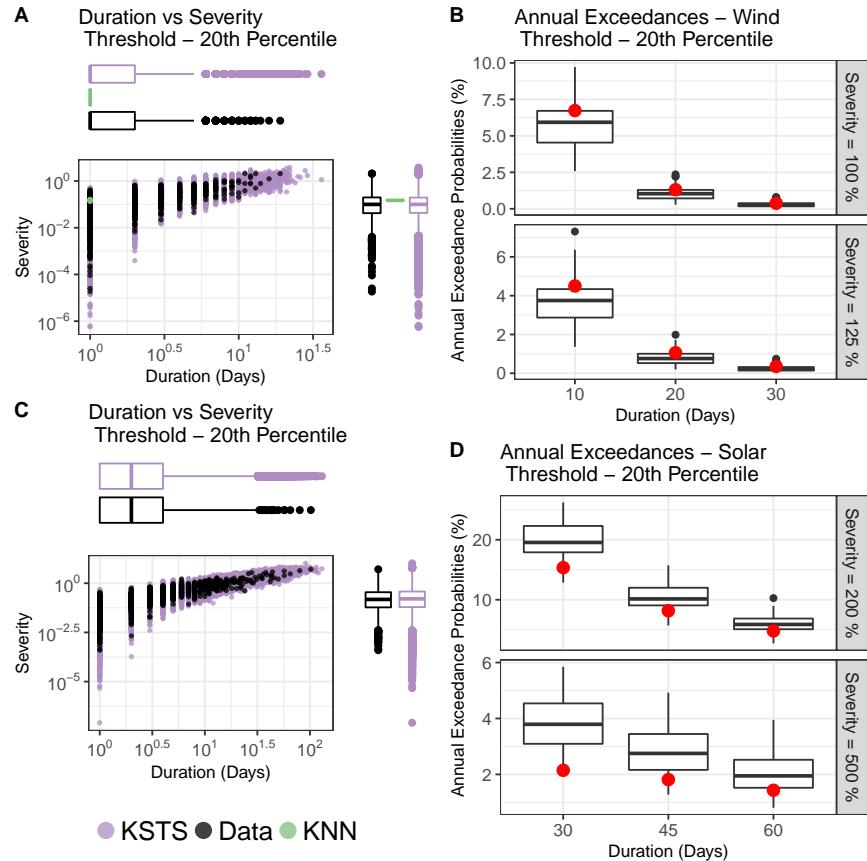


Figure S6: Duration versus Severity plots for all (A) wind and (C) solar energy droughts with marginal distributions (in boxplots for both variables) for the data (black), KSTS (purple) and KNN (green) simulations using 20th percentile as threshold. Probability of annual exceedances for (B) wind and (D) solar energy droughts given a duration and severity with threshold values of 20th percentile. The red dot denotes the exceedance probability calculated from the reanalysis data. The boxplots denote the uncertainty in the 48 generated simulations using KSTS. The duration is in days and the severity is denoted in terms of percentage of the mean historical reanalysis value. For each box-plot, the thick black horizontal line across the box denotes the median of the annual exceedance probabilities from the simulations and the edges of the box denote the 25th and 75th percentiles, and the lower and upper extents of vertical lines outside the box denote the 5th and 95th percentiles.

903 **KSTS and KNN Simulations - Individual Site**
 904 **Characteristics**

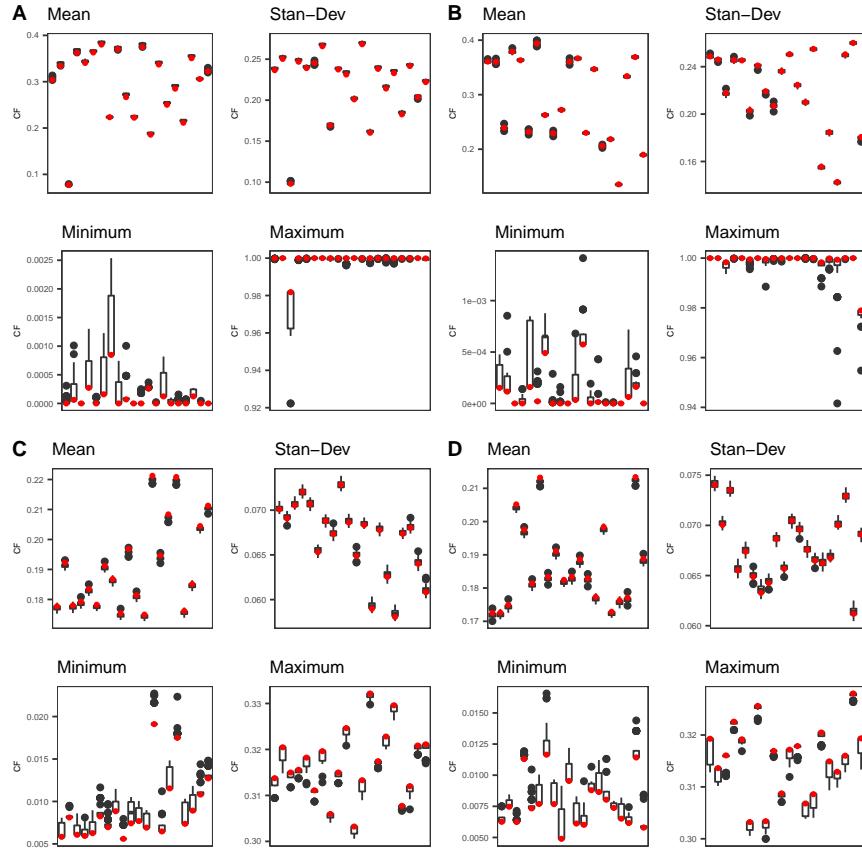


Figure S7: Simulation skill assessments for individual sites in the wind and solar fields for both KSTS and KNN simulations. (A) KSTS wind. (B) KNN wind. (C) KSTS solar. (D) KNN solar. For each sub-plot, we show the mean (top-left), the standard deviation (top-right), the minimum (bottom-left) and the maximum (bottom-right). Red dots denote the reanalysis data value and box-plots denote spread among the 48 simulations. Each subplot includes results for 20 randomly selected grid-points out of the 216 total grids.

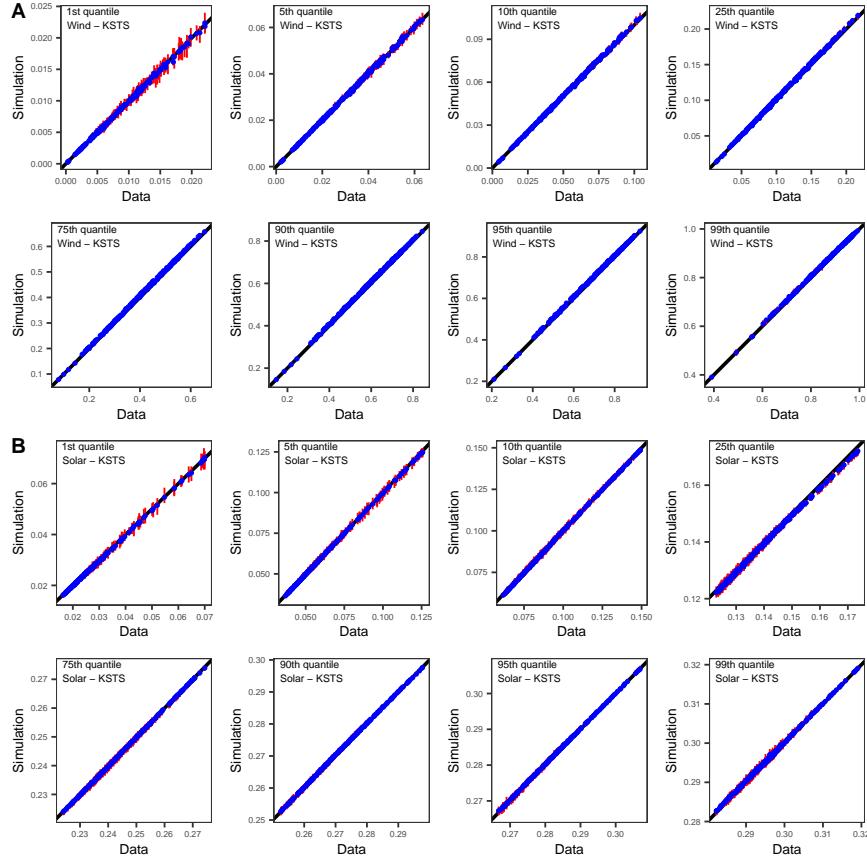


Figure S8: Simulation vs reanalysis data quantile plots for the 1st, 5th, 10th, 25th, 75th, 90th, 95th and 99th percentiles for KSTS simulations. (A) Wind KSTS - Top two rows. (B) Solar KSTS - Bottom two rows. The plots denote the quantiles for all 216 grid points in the wind and solar fields. The red lines denote the mid 90th (5th-95th) percentile range and the blue dots denote the median value in the simulation spread.

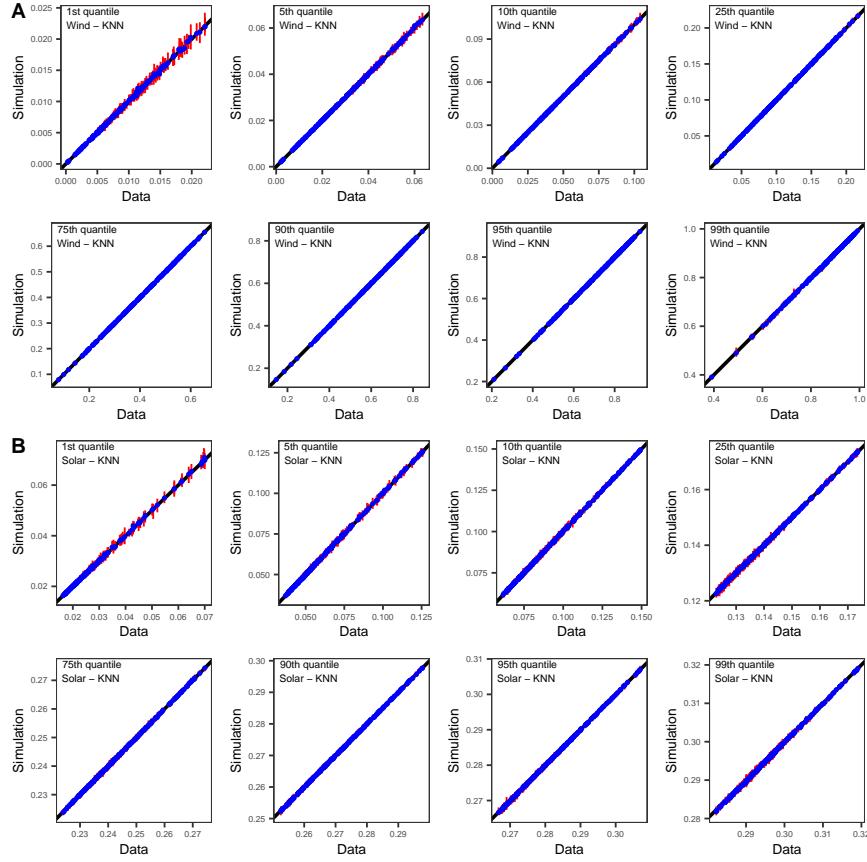


Figure S9: Simulation vs reanalysis data quantile plots for the 1st, 5th, 10th, 25th, 75th, 90th, 95th and 99th percentiles for KNN simulations. (A) Wind KSTS - Top two rows. (B) Solar KSTS - Bottom two rows. The plots denote the quantiles for all 216 grid points in the wind and solar fields. The red lines denote the mid 90th (5th-95th) percentile range and the blue dots denote the median value in the simulation spread.

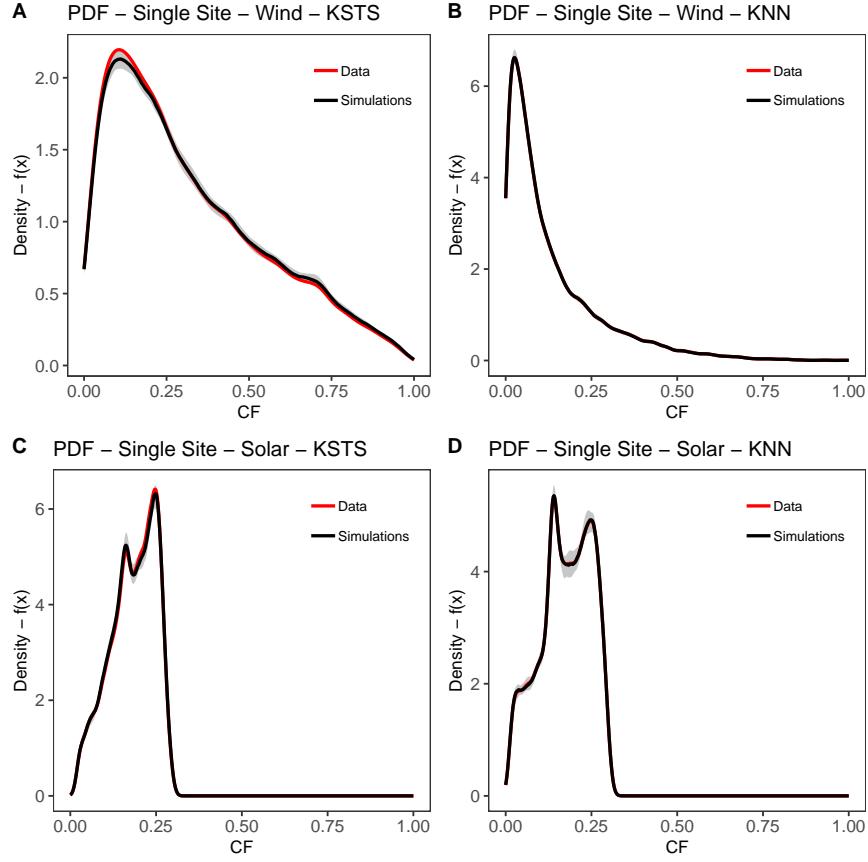


Figure S10: Kernel density estimate / Probability density function (PDF) plots for a single randomly selected grid for wind and solar. The red line denotes the reanalysis data probability density function for the selected site and the black line denotes the median simulation density. The grey region is the mid 90th (5th-95th) percentile range of the simulation spread. The grid point is selected at random separately for KSTS and KNN. (A) Wind KSTS simulation. (B) Wind KNN simulation. (C) Solar KSTS simulation. (D) Solar KNN simulation.

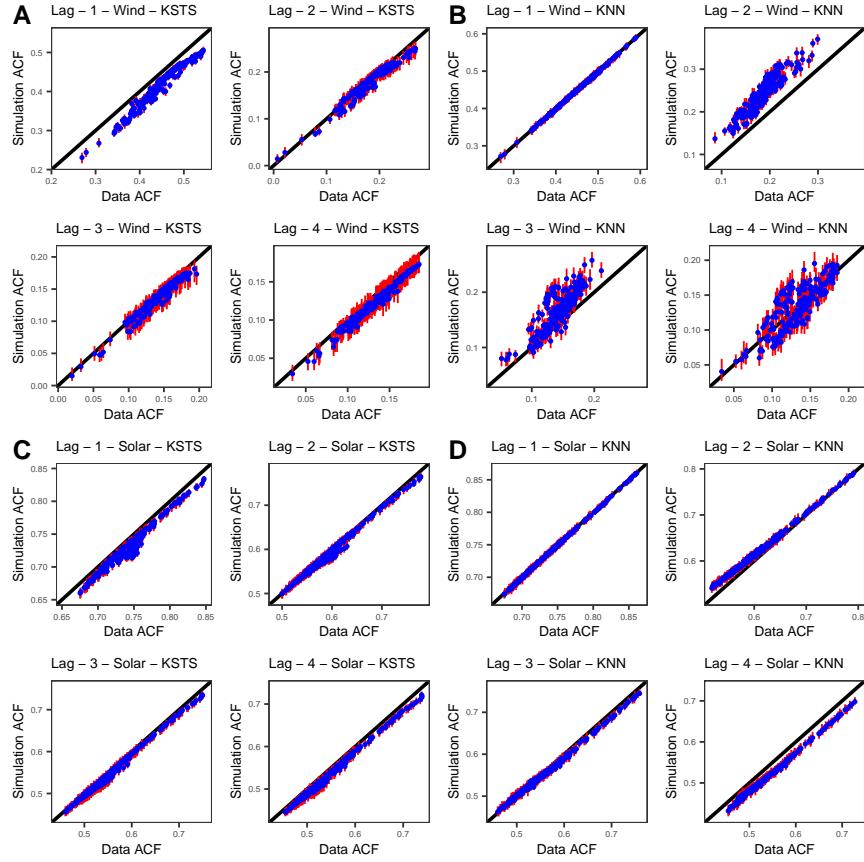


Figure S11: Simulation vs reanalysis data auto-correlation plots for lag 1,2,3 and 4 for all grid points. (A) Wind KSTS simulations. (B) Wind KNN simulations. (C) Solar KSTS simulations. (D) Solar KNN simulations. The red lines denote the mid 90th (5th-95th) percentile range and the blue dots denote the median value in the simulation spread.

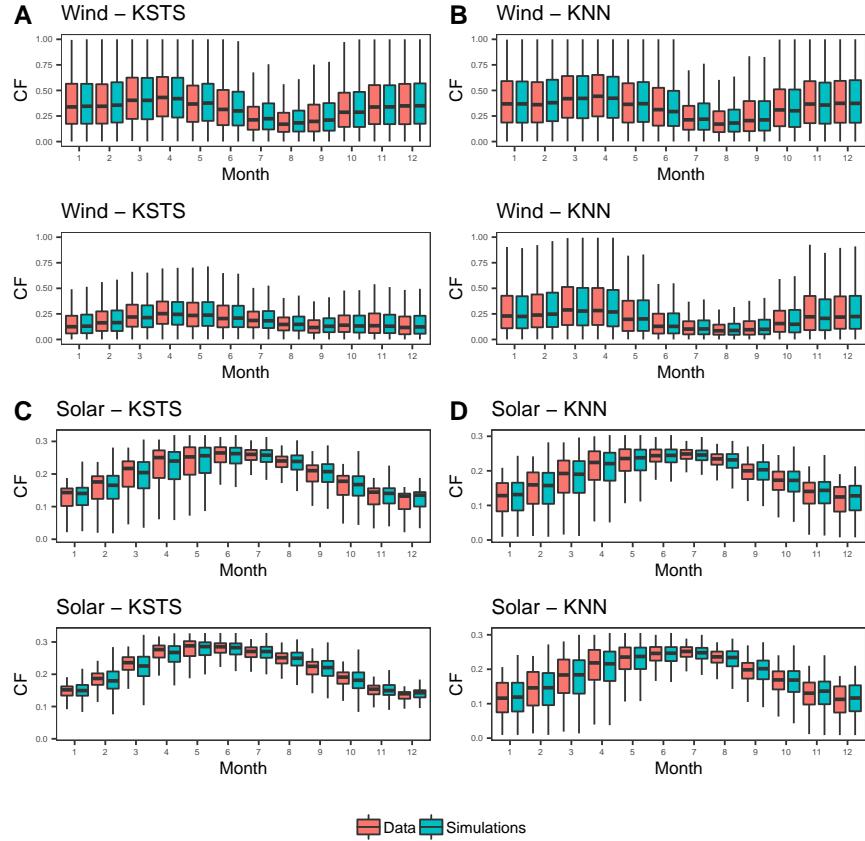


Figure S12: Seasonality / Monthly distribution of the reanalysis data and simulations. The red and green boxplots denote the reanalysis data and simulations respectively. (A) Wind KSTS simulations. (B) Wind KNN simulations. (C) Solar KSTS simulations. (D) Solar KNN simulations. Two grid points are randomly selected for wind and solar. The grids are selected at random separately for KSTS and KNN. Months are numbered in accordance with the Gregorian calendar

905 Cross-Field Dependence

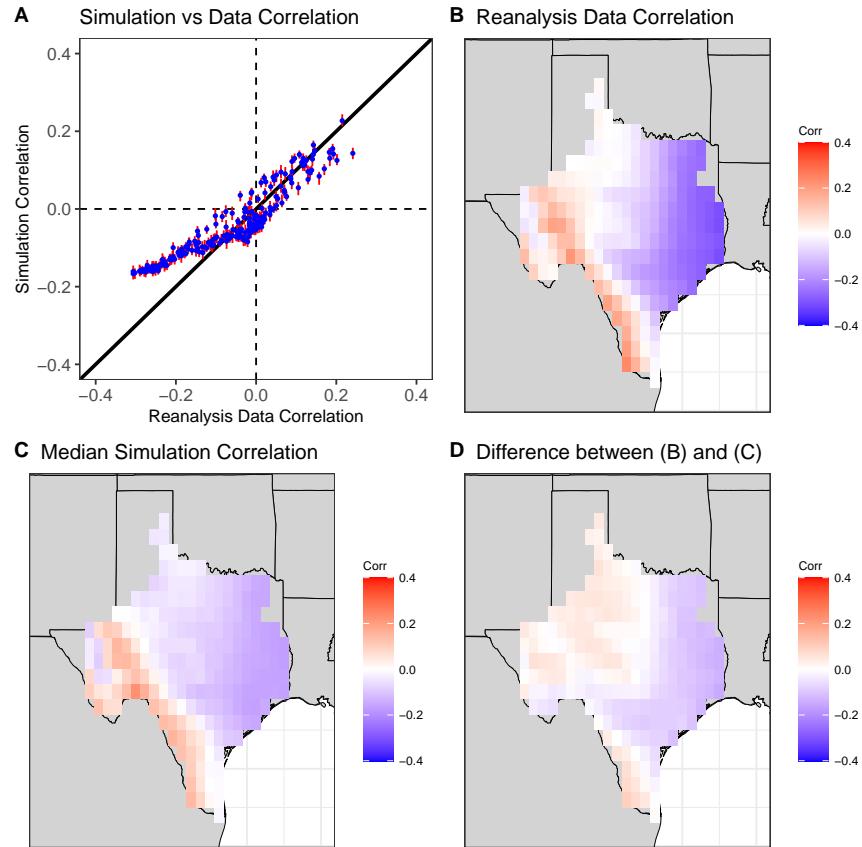


Figure S13: Pearson correlation between wind and solar at each grid point based on simultaneous simulations of wind and solar using KNN. (A) Simulation correlation vs reanalysis data correlation between wind and solar where the red lines denote the mid 90th (5th-95th) percentile range and the blue dots denote the median value in the simulation spread generated using the KNN Method. (B) Map of the grid-wise correlations in the reanalysis data record. (C) Map of the grid-wise median simulation correlations using KNN. (D) Map of the difference between (B) and (C).

906 Individual Field (Wind and Solar) Spatial Correlations

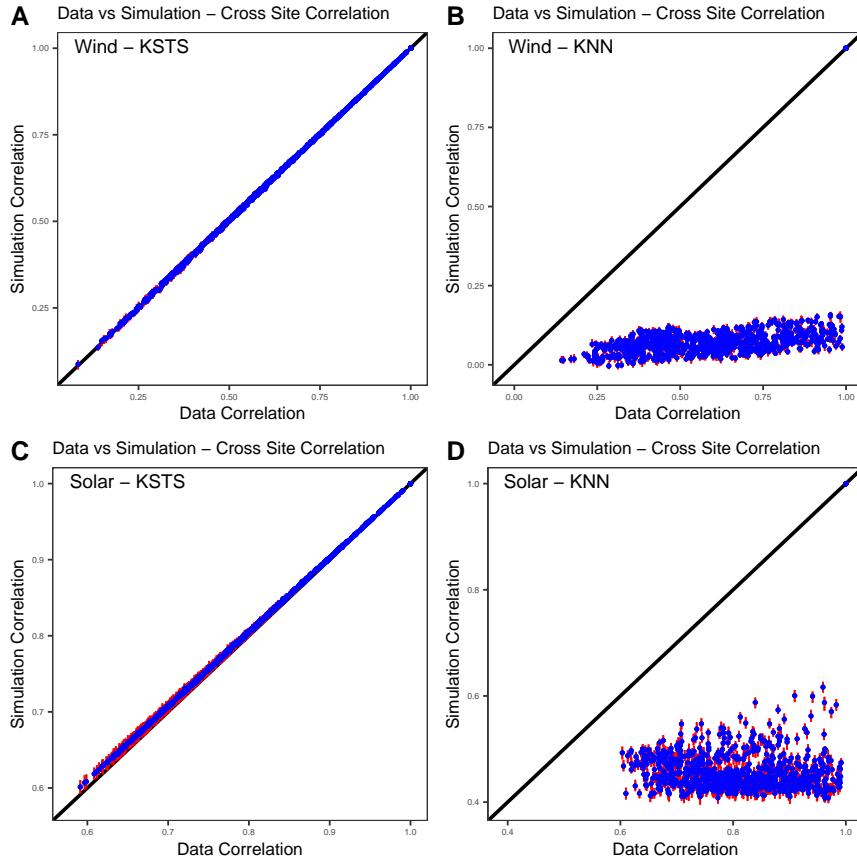


Figure S14: Simulation vs reanalysis data cross site correlation plots for individual fields. (a) Wind KSTS. (b) Wind KNN. (c) Solar KSTS (d) Solar KNN. 40 grids out of 216 are randomly selected and the 40x40 cross correlation values are computed and plotted instead of the entire 216 x 216 correlation values. The correlations are computed using Pearson's method. The red lines denote the mid 90th (5th-95th) percentile range and the blue dots denote the median value in the simulation spread.

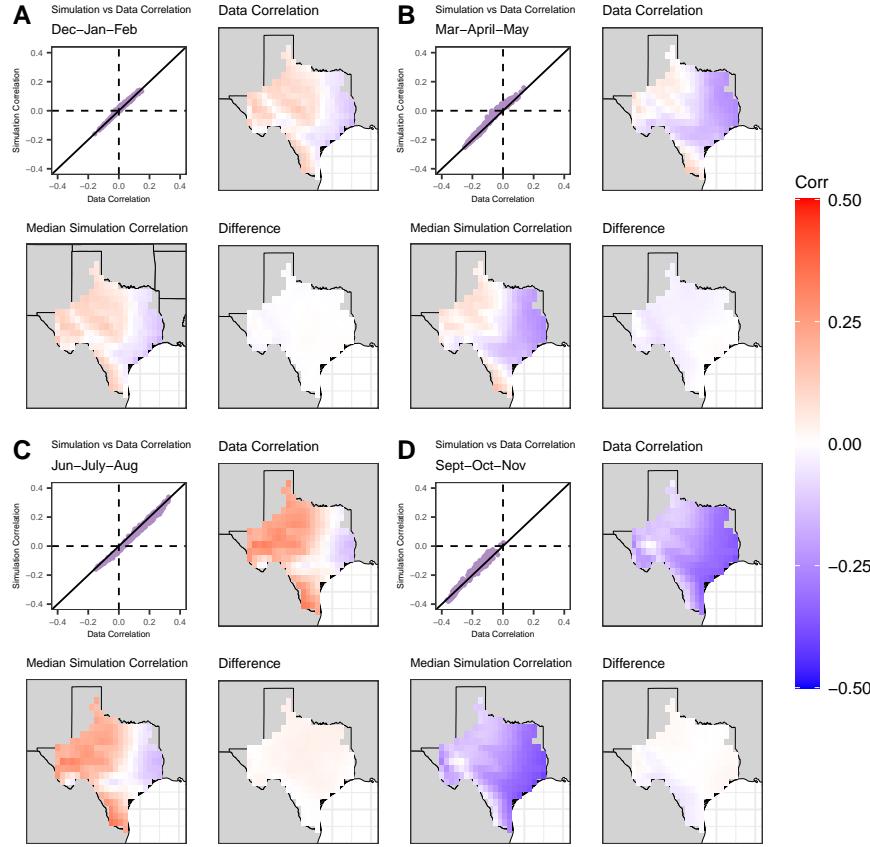


Figure S15: Seasonal correlation between wind and solar at each grid point for KSTS simulations. (A) Dec-Jan-Feb (DJF). (B) Mar-Apr-May (MAM). (C) Jun-Jul-Aug (JJA). (D) Sep-Oct-Nov (SON). For each subplot: (top-left) - median simulation vs reanalysis data correlation between wind and solar. (top-right) - Plot of the reanalysis data correlations. (bottom-left) - Plot of the median simulation correlations. (bottom-right) - Plot of the difference between data and median simulations correlations. The correlations are computed using Pearson's method.

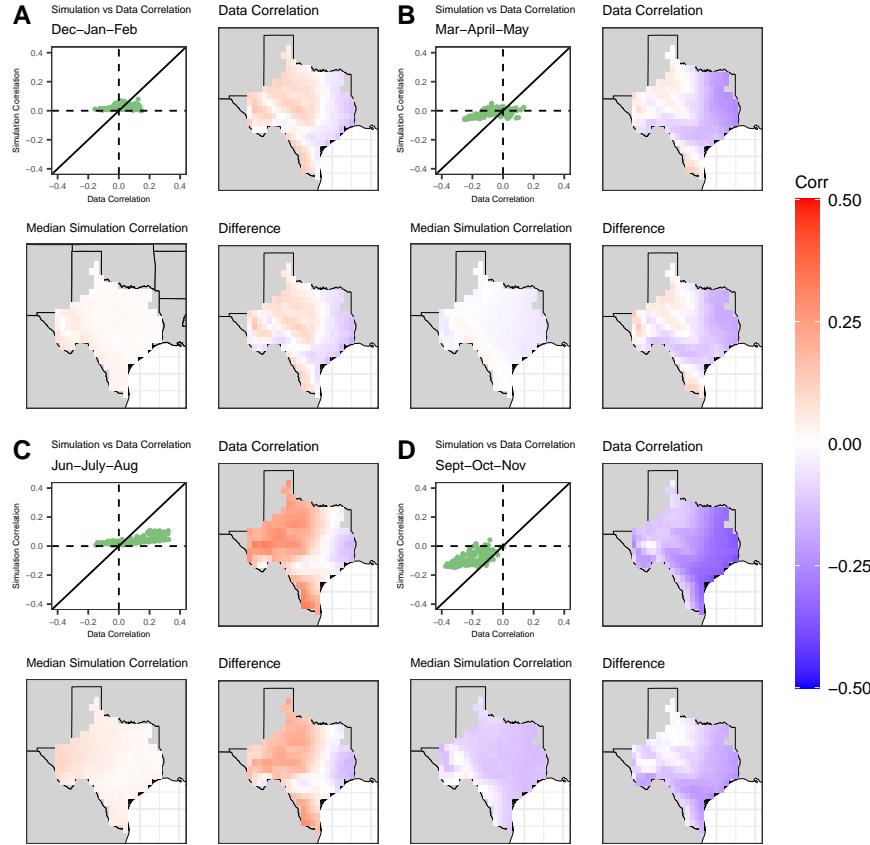


Figure S16: Seasonal correlation between wind and solar at each grid point for KNN simulations. (A) Dec-Jan-Feb (DJF). (B) Mar-Apr-May (MAM). (C) Jun-Jul-Aug (JJA). (D) Sep-Oct-Nov (SON). For each subplot: (top-left) - median simulation vs reanalysis data correlation between wind and solar. (top-right) - Plot of the reanalysis data correlations. (bottom-left) - Plot of the median simulation correlations. (bottom-right) - Plot of the difference between data and median simulations correlations. The correlations are computed using Pearson's method.

907 Data

908 The Electric Reliability Council of Texas (ERCOT - Figure S17), functions as
 909 an Independent System Operator and the balancing authority for the Texas
 910 Interconnection and manages about 90% of state's electric load. ERCOT covers
 911 about 75% of the land area in Texas¹.

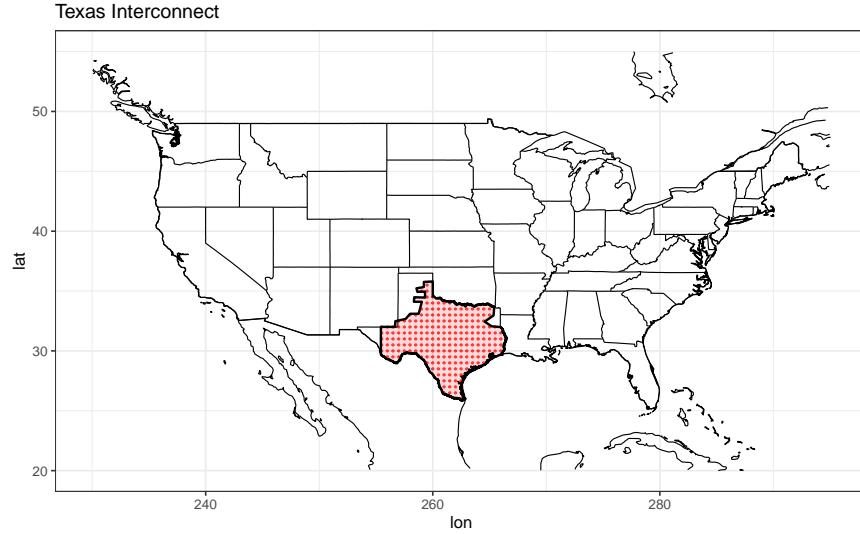


Figure S17: Texas Interconnection / ERCOT domain plot - The red shaded region denotes the area administered by ERCOT. The red dots (216) are the locations of the grid points (0.5° lat \times 0.5° lon) from the ERA-5 reanalysis dataset.

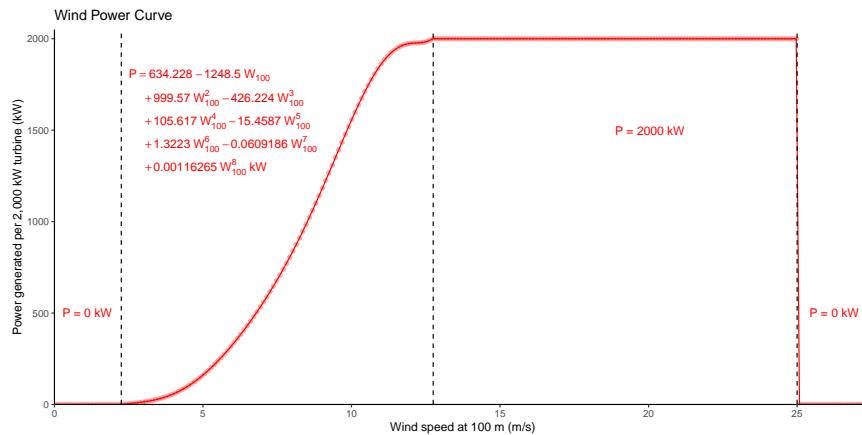


Figure S18: Wind Power Curve for a V90-2.0MW Vestas turbine.

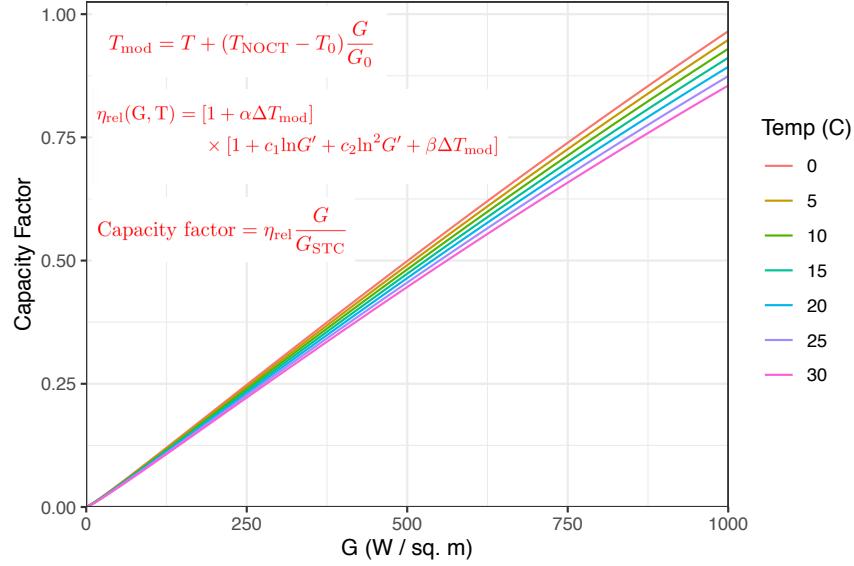


Figure S19: Relationship between solar capacity factor, hourly radiation data (G), and temperature as per Bett and Thornton ². Variables are defined in the following way:

$$G' = \frac{G}{G_{\text{STC}}}$$

$$\Delta T_{\text{mod}} = T_{\text{mod}} - T_{\text{STC}}$$

. The constant values are the following $\alpha = 1.2 \times 10^{-3} K^{-1}$, $\beta = -4.6 \times 10^{-3} K^{-1}$, $c_1 = 0.033$, $c_2 = -0.0092$, $T_{\text{NOCT}} = 48^\circ C$, $T_o = 20^\circ C$, $G_o = 800 \frac{\text{W}}{\text{sq-m}}$, $G_{\text{STC}} = 1000 \frac{\text{W}}{\text{sq-m}}$, $T_{\text{STC}} = 25^\circ C$.

912 **KSTS Description**

913 The algorithm first models the temporal variability at each site and for each
914 variable field. A state space $D_{i,t}$ is defined through an embedding of the time
915 series with a delay parameter ϕ and an embedding dimension m , where i is
916 the index for the site/variable combination. Following Lall and Sharma ³, we
917 consider that the conditional density $f(x_{i,t+1}|D_{i,t})$ is defined through the $x_{i,t'+1}$
918 corresponding to the time indices t' associated with the k-nearest neighbors
919 (knn) of $D_{i,t}$ in the historical data set. The kernel function p_j associated with
920 the jth nearest neighbor is proportional to the rank j of the neighbor ³. The
921 sequential drawing from the knn successors at each time step using the specified
922 kernel leads to a simulation of $x_{i,t}$ in a time series dependence structure. Since
923 the procedure leads to a resampling of the historical data, the algorithm can
924 be considered to be a bootstrap which preserves the time dependence in serial
925 data.

926 Given a state space $D_{i,t}$ at site i and time t , the k-nearest neighbor algorithm
927 is used to identify a set of time indices $\tau_{i,t}$ that correspond to the time instances
928 corresponding to the nearest neighbors of site i at time t . An example of the
929 nearest neighbors for a site with an embedding $D_{i,t}$, defined as (x_t, x_{t-1}) taking
930 value of (10,9.8) at time t=100 days would be the closest historical vectors to
931 (10,9.8) using the data for this site.

932 The time instances at which these neighbors occur in the historical time series
933 are then recorded in the order in $\tau_{i,t}$. The kernel $p_{i,j}$ associates a probability
934 proportional to $1/j$ for the j th element (time instance of the j th nearest neighbor
935 of $D_{i,t}$) in $\tau_{i,t}$, for the first k neighbors and 0 elsewhere. For space-time neighbors
936 across all sites, i.e. to address spatial dependence, we now identify appropriate
937 k-nearest neighbors by finding the time indices in the historical data that have
938 the highest likelihood of being selected across all sites given their associated
939 resampling probabilities.

940 Define $T_{i,t}$ as a matrix such that the rows and columns are pointers for sites
941 and unique time indices from the historical data respectively. The columns
942 record the resampling probabilities associated with the time indices of the k-
943 nearest neighbors for all sites at time t . The similarity vector S_t is the sum of
944 the resampling probabilities associated with each unique time index across all
945 sites. The curtailment of the similarity vector S_t is carried out by selection of
946 the time indices which correspond to the k highest resampling probability values
947 in S_t , now designated as the k-nearest neighbor candidates for the entire spatial
948 field. The full spatial field of the simulation for the next time step is resampled
949 after re-scaling probability values (such that they add to 1) of the curtailed
950 similarity vector S_t . Other measures of similarity of the spatial neighbors of the
951 temporal process could also be considered.

952 **Hyper-parameters of the Algorithm**

953 **Resampling Kernel Weight Function**

954 The resampling kernel p_j selected for the simulator is the one proposed

by Lall and Sharma³. This resampling kernel decreases monotonically with increase in distance, with the bandwidth and kernel shape varying with the local sampling density. Overall, the kernel is adaptive to the dimensionality of the state space, with implicit dependence through the distance calculations. Further, the resampling weights need to be computed only once and stored, which significantly reduces computation time.

$$p_j = \frac{1/j}{\sum_{j=1}^k 1/j}$$

Other options for the kernel include a uniform kernel ($p_j = 1/k$) or a power kernel based on the distances of the k neighbors. Refer Lall and Sharma³ for further details on the behavior of the kernel in the boundary region, for bounded data and comparison to a uniform kernel.

Number of Neighbors (k) and Model Order (m)

One method to choose model hyper-parameters involves criterion that minimize the mean squared error in forecast. The generalized cross validation (GCV) score was suggested to select k and m ³. The selected number of nearest neighbors k and the order of the feature vector m are the ones which minimize the GCV score, which is given by

$$GCV = \frac{\sum_{i=1}^n e_i^2/n}{\left(1 - \frac{1}{\sum_{j=1}^k 1/j}\right)^2}$$

where, e_i is the forecast error at point i for the model fit to all the data without it and n is the total number of points. The selection of these parameters by GCV is most appropriate if the model errors e_i are normally distributed or if the variables are transformed such that model errors are normally distributed. Non-normality of the errors may lead to sub-optimal choice of k and m with respect to its conditional mean and variance.

Another method to select the model lags in the feature vector is the false nearest neighbors algorithm which determines the embedding dimension for the process⁴. Finally, an ad-hoc choice of $k = n^{0.5}$ is suggested across the knn literature, with low sensitivity around this value. Further suggestions include trying various combinations of k and m followed by visual examinations of the simulation attributes and data³.

Scaling Weights (w)

The simplest selection choice for the weights w , which weigh the euclidean distance of the selected lags m is to be specified *a priori* with uniform values. The weights can also be selected such that they minimize the forecast error in least squares sense when used in a knn regression setup⁵. An alternate adaptive strategy is to compute scaling weights (w) for the knn resampling approach such that they are the regression coefficients of the selected external predictors from a parametric regression model⁶.

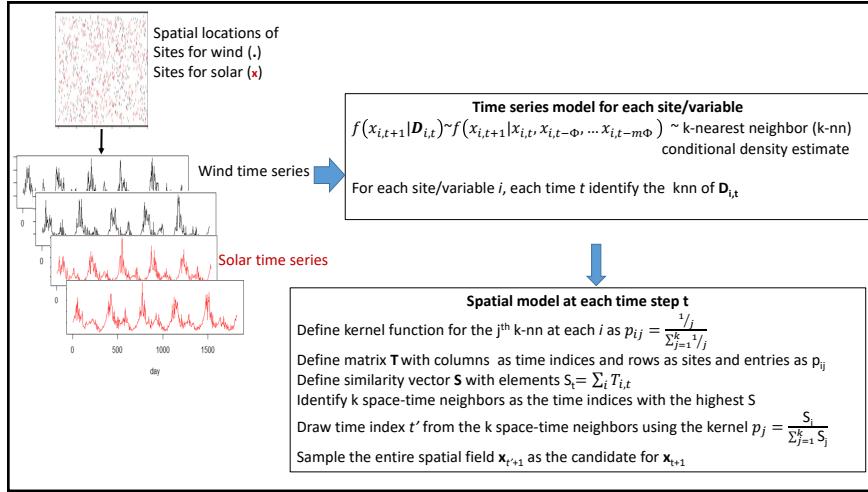


Figure S20: Schematic representation of the k-Nearest Neighbor Space Time Simulator (KSTS) with application to wind and solar fields.

991 Supplementary Materials - References

- 992 1. “About ERCOT.” Electric Reliability Council of Texas. <http://www.ercot.com/about>. Accessed 10 July 2021.
- 993 2. Bett, Philip E., and Hazel E. Thornton. “The climatological relationships
994 between wind and solar energy supply in Britain.” Renewable Energy 87
995 (2016): 96-110.
- 996 3. Lall, Upmanu, and Ashish Sharma. “A nearest neighbor bootstrap for
997 resampling hydrologic time series.” Water resources research 32.3 (1996):
998 679-693.
- 999 4. Kennel, Matthew B., Reggie Brown, and Henry DI Abarbanel. “Determining
1000 embedding dimension for phase-space reconstruction using a geo-
1001 metrical construction.” Physical review A 45.6 (1992): 3403.
- 1002 5. Yakowitz, S., and M. Karlsson. “Nearest neighbor methods for time series,
1003 with application to rainfall/runoff prediction.” Advances in the statistical
1004 sciences: Stochastic hydrology. Springer, Dordrecht, 1987. 149-160.
- 1005 6. Souza Filho, Francisco Assis, and Upmanu Lall. “Seasonal to interan-
1006 nual ensemble streamflow forecasts for Ceara, Brazil: Applications of a
1007 multivariate, semiparametric algorithm.” Water Resources Research 39.11
1008 (2003).
- 1009 7. Lall, Upmanu, Naresh Devineni, and Yasir Kaheil. “An empirical, non-
1010 parametric simulator for multivariate random variables with differing marginal
1011

¹⁰¹² densities and nonlinear dependence with hydroclimatic applications.” Risk
¹⁰¹³ Analysis 36.1 (2016): 57-73.