

Joule

Inferences on Long Duration Storage Requirements for Wind and Solar Energy from a Novel Space-Time Simulator --Manuscript Draft--

Manuscript Number:	JOULE-D-21-00787
Full Title:	Inferences on Long Duration Storage Requirements for Wind and Solar Energy from a Novel Space-Time Simulator
Article Type:	Research Article
Keywords:	Long Duration Storage; Wind and Solar Generation; Texas Interconnection; Stochastic Simulators; Macro-Energy Systems
Corresponding Author:	Yash Vijay Amonkar Columbia University New York, NY UNITED STATES
First Author:	Yash Amonkar
Order of Authors:	Yash Amonkar David J Farnham Upmanu Lall
Abstract:	The space-time variability of potential electricity generation is a critical factor for the design of wind and solar dominated power systems. Limited historical data is available on incoming solar radiation and wind. Consequently, a stochastic model for simulating the co-variation of wind and solar fields is needed 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 a novel k-Nearest Neighbor Space-Time Simulator that accounts for the space-time dependence in high dimensional wind and solar fields and can simulate synthetic wind and solar fields of arbitrary length. Probability distributions of the severity-duration-frequency of regional energy droughts relative to a target are produced. As expected, the severity-duration-frequency of shortfalls is much greater when space-time dependence is properly considered than when it is not.
Opposed Reviewers:	
Suggested Reviewers:	Ashish Sharma, Ph.D. Professor, University of New South Wales Faculty of Engineering ashish.sharma@unsw.edu.au The proposed algorithm builds up on the paper by Lall and Sharma (1996). Dr. Sharma will have a good understanding of the algorithm and be able to provide critiques too. Jesse Jenkins, Ph.D. Assistant Professor, Princeton University jessejenkins@princeton.edu Lead Author of "Getting to Zero Carbon Emissions in the Electric Power Sector" - Joule (2018) Stefan Pfenninger, Ph.D. Assistant Professor, TU Delft: Technische Universiteit Delft S.Pfenninger@tudelft.nl Author of the paper - "Impacts of Inter-annual Wind and Solar Variations on the European Power System." Joule (2018)
Additional Information:	
Question	Response
Standardized datasets A list of datatypes considered standardized under Cell Press policy is available here . Does this manuscript	No

report new standardized datasets?	
Original Code Does this manuscript report original code?	Yes
Reviewers must have anonymous access to these original code that is free-of-cost. Please provide code location and instructions for access here. Please consult this Author's guide for more information: " How standardized datasets and original code accompany Cell Press manuscripts from submission through publication " or email us at joule@cell.com . as follow-up to " Original Code Does this manuscript report original code?"	https://github.com/yashamonkar/LDS-Inferences

Department of Earth and Environmental Engineering
Columbia University
New York, NY 10027
June 2021

Dear Editor,

I am writing to submit the manuscript entitled, “Inferences on Long Duration Storage Requirements for Wind and Solar Energy from a Novel Space-Time Simulator” by Yash Amonkar, David J. Farnham, and Upmanu Lall for consideration for publication in *Joule*.

The need for long duration storage to help manage the variability of wind and solar power generation has been increasingly discussed, including in Dowling et al. (2020) in *Joule*. Estimating the necessary storage capacity for a given system requires relatively long data records to identify the statistics of wind and solar “droughts”. However, long records of wind and solar potential are often not available. We address this issue by proposing a novel simulator capable of producing synthetic, realistic, and long-record wind and solar spatio-temporal fields. These fields can then be used to estimate the severity and duration of events when wind and solar generation potential is low, and storage may be needed to bridge the gap in supply.

Our attached paper presents a novel k-nearest neighbors’ model capable of simulating joint wind and solar fields with hundreds of individual sites. The application of the simulator to ERCOT region with 40-yrs of gridded reanalysis data is demonstrated. The ability of the simulator to infer long-duration storage requirements and the corresponding sampling uncertainty is shown. Finally, the importance of the properly simulating the spatial dependence when modelling such large fields is demonstrated by comparing our novel simulator with another model that adequately fits individual sites but fails to capture the spatial dependence.

We have selected the journal *Joule* due to its focus on sustainable energy and interdisciplinary solutions for providing sustainable energy solutions. We think that the novel methodologies and ideas presented in our paper will be of interest and utility to many researchers, industry professionals and academics working at the intersection of sustainable energy, energy storage and its interaction with grid reliability.

This manuscript has not been submitted elsewhere. If you require additional information regarding this manuscript, please contact me at yva2000@columbia.edu.

Thank you for your time and consideration,
Yash Amonkar

Inferences on Long Duration Storage Requirements for Wind and Solar Energy from a Novel Space-Time Simulator

Yash Amonkar^{1,2,4,*}, David J Farnham³, 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

⁴Lead Author

*Correspondence: yva2000@columbia.edu

Abstract

Summary

The space-time variability of potential electricity generation is a critical factor for the design of wind and solar dominated power systems. Limited historical data is available on incoming solar radiation and wind. Consequently, a stochastic model for simulating the co-variation of wind and solar fields is needed 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 a novel *k*-Nearest Neighbor Space-Time Simulator that accounts for the space-time dependence in high dimensional wind and solar fields and can simulate synthetic wind and solar fields of arbitrary length. Probability distributions of the severity-duration-frequency of regional energy droughts relative to a target are produced. As expected, the severity-duration-frequency of shortfalls is much greater when space-time dependence is properly considered than when it is not.

Context and Scale

Many regions plan to integrate more wind and solar generation into the energy grid. Greater reliance on wind and solar generation introduces power supply variability that can pose risks of undersupply. Designing a reliable grid thus requires estimating the frequency, duration, and severity of periods of low wind and/or solar generation potential. We develop and apply a novel statistical simulation model that can produce realistic, synthetic realizations of wind and solar potential across a region, which

allows us to estimate the probability of extreme events that are not necessarily represented in the relatively short observational records available. This tool allows us to estimate the probability of regional wind and solar energy “droughts”, and hence allows us to estimate the long duration storage (LDS) needed to achieve desired grid reliability. An application to the Texas Interconnection demonstrates how better estimates of LDS requirements can be obtained through simulation.

Highlights

Space-time correlations in wind and solar fields can be spatially and seasonally non-homogeneous in a region like Texas.

The probability distribution of the aggregate long duration storage (LDS) needed to withstand renewable energy shortages in this setting needs to account for this variable dependence.

A machine learning approach that can capture the joint co-variation of wind and solar fields is presented.

An example application to the Texas Interconnection demonstrates the utility of the approach.

Key Words

Long Duration Storage, Wind and Solar Generation, Texas Interconnection, Stochastic Simulators, Macro-Energy Systems

1 Introduction

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

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

13 Long Duration Storage (LDS), defined as storage needed to meet deficits
14 for duration greater than 10 hours^{11 12}, is critical to economically meet grid
15 reliability targets while relying primarily on wind and solar generation⁹. Many
16 recent macro scale electricity studies focusing on renewable electric grids and
17 economy wide de-carbonization models commonly include LDS and expansion of
18 long-distance transmission capacity across the continent to smooth the variation
19 in renewable production⁹. Such an approach necessitates proper consideration
20 of the temporal and spatial dependence structure of available wind and solar
21 energy including their cross-dependence.

22 Given a candidate regional configuration of wind and solar generators, sizing
23 LDS economically for a regional grid requires estimates of the probability
24 of potential energy shortages for different durations. 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 inter-annual 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 USA¹⁷. 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. The reliability of reanalysis products is time-dependent with the post 1970s satellite measurement era considered to be the most reliable, effectively setting the data record length to 40 years²³. 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 paper²⁴ 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 spatial or temporal correlation across wind and solar sources, or the stochastic properties of the duration and severity of aggregate 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 long duration energy shortage across a spatial 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 while reproducing the space and time dependence structure of the wind and solar fields, so that more reliable estimates of the severity and duration of regional total wind and solar energy potential and their uncertainty can be estimated.

71 The wind and solar data from the ERA-5 reanalysis product¹⁹ are used for the
72 development and testing of a stochastic spatio-temporal model that can provide
73 insights as to the variation of the aggregate energy production from a set of
74 spatially distributed wind and solar generation facilities. We take the Electric
75 Reliability Council of Texas (ERCOT) - Texas Interconnection region²⁵ as a
76 target example to explore the historical record and to demonstrate the perfor-
77 mance of our algorithm. LDS considerations motivate the use of daily data on
78 potential wind and solar resource. The implicit assumptions is that chemical
79 batteries help smooth out the sub-daily time-scale shortages¹¹. The potential
80 wind power is derived from hourly wind speed data at 100m using a turbine
81 rating curve and then averaged for the day. The downward surface solar radi-
82 nation is taken as the mean value for the day (see the Experimental Procedure
83 section).

84 Over a large region (e.g., the Texas Interconnection), the wind and solar
85 generation assets are likely to be spatially distributed throughout the region²⁶.
86 Non-homogeneous and non-local space and time correlations in the potential
87 energy production across the assets utilized by a grid operator are possible.
88 The annual and seasonal variation of the daily wind and solar energy potential
89 across 216 grid points using daily averages of daily wind and solar reanalysis
90 data for our example application to Texas is illustrated in Figures 1 and S1.

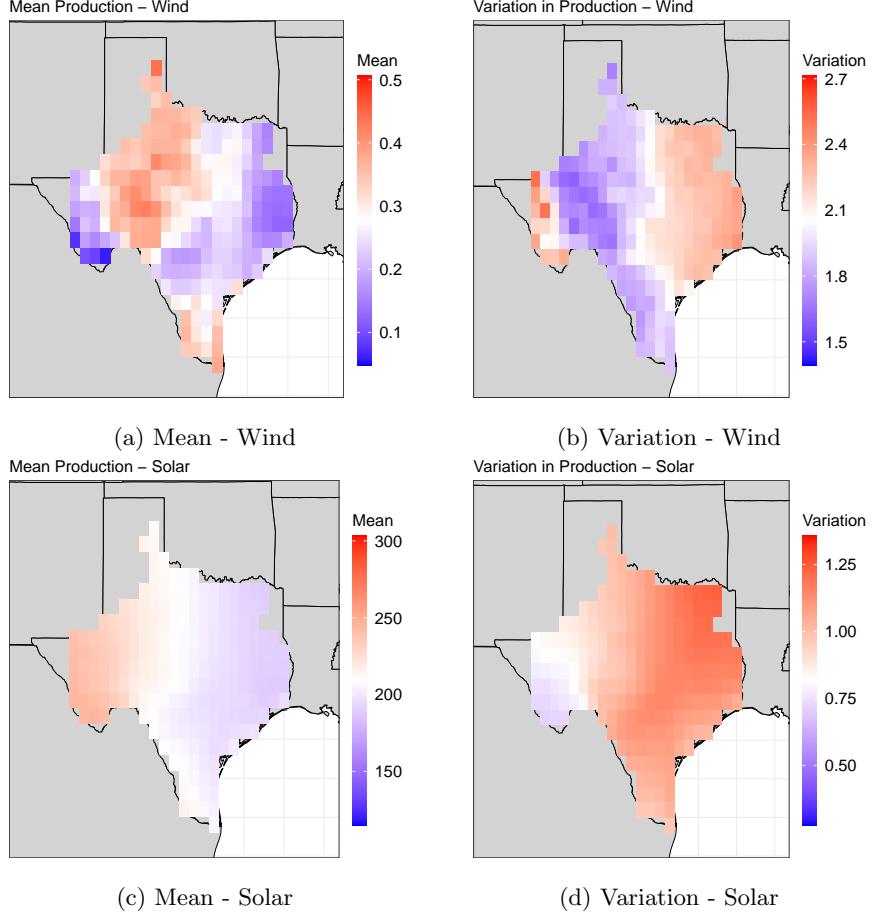


Figure 1: Mean and variation in daily wind capacity factors and downward surface solar radiation (W/m^2) across the Texas Interconnection. (a) Mean daily wind capacity factors. (b) Variation in daily wind capacity factors. (c) Mean daily downward surface solar radiation. (d) Variation in daily downward surface solar radiation. 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 radiation often exhibit significant spatial correlation
 92 that varies by location and time of year and needs to be accounted for in an
 93 analysis of potential renewable energy droughts or LDS system sizing²⁷. As
 94 seen in Figure 1 and S1, wind and solar along the Gulf of Mexico and the
 95 land-area adjoining Louisiana are regions with lower generation potential but
 96 with high variability. The mean wind capacity factor and its variability (Figure
 97 1) is non-homogeneous. The highest capacity factors are in the north-western
 98 and southern-most portions of the interconnection. The highest variability is
 99 in the eastern portion of the interconnection (Figure S1). Daily wind capacity

100 factors are generally highest during spring and lowest during summer while wind
101 variation is highest during fall and lowest in summer and spring. The mean daily
102 radiation and its variability (Figure 1) is more homogeneous and a function of
103 the season, with low mean radiation and high variability in winter (DJF) and
104 high mean radiation and low variability in the summer (JJA).

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

110 The new k-Nearest Neighbors Space-Time Simulator (KSTS) algorithm is
111 structured as follows. A model for temporal variability at each site and for
112 each variable (wind and solar) is considered first. This entails defining a state
113 space through an embedding of the time series. A time series simulation can
114 then be achieved by sequentially drawing from the k nearest neighbor successors
115 of the embedding at each time step, but this will not preserve spatial depen-
116 dence. Spatial dependence is then introduced by identifying the most likely
117 neighbors of the full spatial field by aggregating neighbor likelihoods for each
118 site/variable. If the state space evolution at two sites is similar (i.e., identified
119 by the same neighbors in time), then the evolution of those two sites would be
120 fully synchronous. Thus the similarity in the selection of neighbors reflects the
121 similarity in dynamics and provides a useful basis for space-time conditioning
122 of a random field’s dynamics. The k-nearest neighbors identified across all the
123 sites as the most similar at a given time, are then used to randomly draw a
124 full spatial field for the next time step, using a kernel function that accounts
125 for their degree of similarity through a probability measure. The process is
126 repeated sequentially to generate a time series simulation of the spatial field.

127 The target variables, wind power and downward surface solar radiation, have
128 non-Gaussian and skewed distributions. They are also potentially bounded,
129 since the maximum wind a generator can use is limited and the maximum in-
130 coming solar radiation on a calendar day is also limited. Consequently, the
131 probabilistic sampling using k-nearest neighbors provides an effective approach
132 to sampling from a nonparametric distribution applied to each target variable.
133 The seasonality in the variables is accounted for by restricting search of k-nearest
134 neighbors using a moving window around the Day of Year (DOY). This new al-
135 gorithm generalizes a k-Nearest neighbor algorithm²⁸ used for univariate or low
136 dimensional multivariate simulations of non-Gaussian and nonlinear dependence
137 that has been used extensively for other climate variables^{29 30 31 32}.

138 We apply our new KSTS algorithm to assess the severity-duration-frequency
139 of long duration storage needs associated with the aggregate regional energy pro-
140 duction. We show that the simulator captures the regional aggregate as well as
141 the site by site probabilities of wind and solar energy produced including the
142 spatial correlation within and across the two fields and the temporal autocor-
143 relation at each site. We also run a simulation (henceforth termed KNN) that
144 preserves the time series structure but not the spatial structure or the wind-solar
145 dependence. As one may expect, this demonstrates a significant underestima-

146 tion of the regional LDS probabilities.

147 For the application presented, we use the 40-yr gridded daily wind and so-
148 lar data from the ERA-5¹⁹ reanalysis dataset for 216 sites (grids/nodes) in the
149 Texas Interconnection. Using the KNN or KSTS algorithm one can generate a
150 large number (e.g., 100) of synthetic 40-year simulations (or equivalently a 4000
151 year simulation) of the daily wind and solar fields, without and with spatial de-
152 pendence preserved, respectively. From each simulation we extract the duration
153 and severity of each drought event as a shortage in aggregate energy produced
154 across the grid, relative to a target threshold. The probabilities of LDS severity
155 and duration can then be assessed from this derived set of events. If multiple
156 simulations of 40 years are generated, then one can get also get an estimate of
157 the uncertainty associated with the probability of a LDS severity-duration given
158 40 years of data. If a single long simulation is generated, then we can estimate
159 LDS severity-duration probabilities with reduced uncertainty using the longer
160 synthetic record.

161 Results

162 We present an evaluation of the severity-duration-frequency of the aggregate en-
163 ergy droughts for the Texas Interconnection with (KSTS) and without (KNN)
164 preserving the spatial structure and wind-solar dependence in simulations. For
165 illustrative purposes a uniform land allocation (and consequently capacity within
166 wind and solar separately) for generation is considered. For both types of simu-
167 lations, we generated 48 realizations of 40 years of daily wind and solar data at
168 each of the 216 sites. In an actual use case, an optimization model may be used
169 to allocate wind and solar capacity across the sites and estimate the regional
170 size of LDS capacity using the simulations developed. The results presented
171 here illustrate the importance of getting the space-time dependence right in the
172 simulations for a proper estimation of the regional LDS capacity given a candi-
173 date spatial configuration of wind and solar generation. Detailed performance
174 statistics of the simulator are presented in the supplement.

175 **Severity-Duration-Frequency of Energy Droughts**

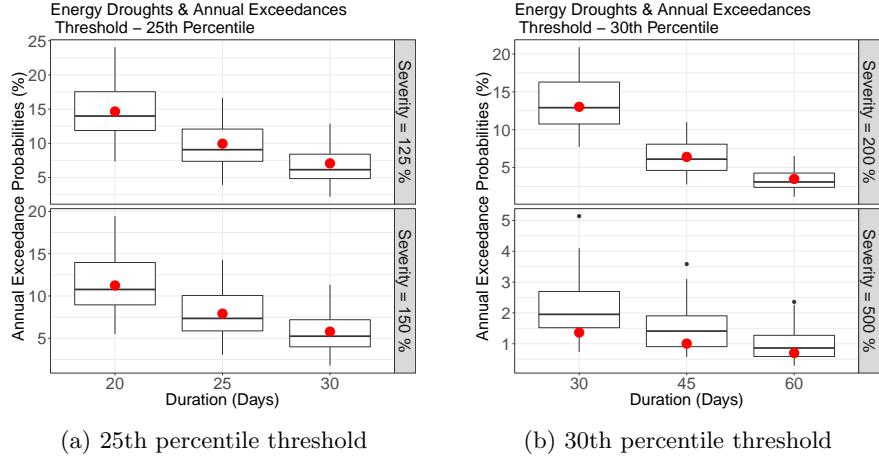


Figure 2: Probability of joint annual exceedances for energy droughts given a duration and severity with threshold values of (a) 25th percentile and (b) 30th percentile. The red dot denotes the exceedance probability calculated from the reanalysis data. The boxplots denote the uncertainty in the 48 generated 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.

176 Energy droughts are defined as continuous periods when the daily production
 177 falls below a target threshold. The threshold value, changing every calendar
 178 day in a year, can be thought of as a forward contract's daily obligation to be
 179 supplied based on the seasonality of the historical reanalysis data. The severity
 180 of the drought is the accumulated deficit in production over the duration of the
 181 event, i.e., the level of default on a potential contract covering the period, while
 182 the duration of the event is the duration during which the deficit exists. Figure
 183 2 (a) shows the annual exceedance probabilities for energy droughts of duration
 184 20, 25 and 30 days with severity of 100% and 150% when the target threshold
 185 is the 25th percentile of the distribution of energy that could be produced over
 186 that period based on the historical data. The severity of energy droughts was
 187 scaled by the mean daily historical production, with a severity of 100% denoting
 188 a shortfall equal to the mean daily historical value. The annual exceedance
 189 probabilities were computed using local regression (Locfit)³³ with the number
 190 of exceedances regressed against the duration and severity using a Poisson link
 191 function. (see Experimental Procedures Section and Supplementary Materials)

192 The KSTS simulations bracket the exceedance probabilities seen in the re-

193 analysis data (Figure 2). For example, an energy drought with duration over 30
194 days with a severity of 150% relative to a threshold guaranteeing delivery set at
195 the 25th percentile of daily regional generation, has an annual exceedance prob-
196 ability of $\sim 5\%$, based on the reanalysis data. This corresponds to an event that
197 may be expected to be exceeded once every 20 years. The median exceedance
198 probability from the simulations is quite close to this, but with considerable un-
199 certainty around that value. The 25th to 75th percentiles from the simulations
200 are around 4% to 7.5% with the 5th and 95th percentiles extending from 2% to
201 11%, demonstrating the limitations of using solely the original 40 year record
202 for such evaluations.

203 Results from increasing the target threshold to the 30th percentile of daily
204 regional energy production and looking at higher severity and longer duration
205 droughts are shown in Figure 2 (b). The KSTS simulations bracket the ex-
206 ceedance probabilities seen in the reanalysis data for the severity of 200%. The
207 simulations show higher exceedance probabilities than the data for the 500%
208 case, which is not surprising considering these are rare events with mean an-
209 nual exceedance probabilities of 0.5-1.5% and thus are difficult to identify given
210 relatively short data records. The severity/duration probabilities from the his-
211 torical record of 40 years have high uncertainty for events that are rarer than
212 perhaps once every 10 years (annual exceedance probability of 0.1) given this
213 record length^{34 35}. The simulations show that these extreme events could oc-
214 cur far more frequently than if we rely on just the short historical records. In
215 these illustrations, we considered specific thresholds for supply guarantee, spe-
216 cific drought durations and severity levels and present the range of probabilities
217 of exceedance from the simulations. In a system design optimization model,
218 for a candidate spatial configuration of generation, the simulator would provide
219 the probability distribution for a candidate LDS capacity that is considered to
220 meet the deficit over a specified duration (e.g., specified by a contract). Al-
221 ternately, one could also compute the probability distribution of the shortage
222 beyond the candidate LDS to assess potential penalties for non-delivery, if those
223 were considered in the optimization model.

224 Annual exceedance probabilities for different combinations of duration and
225 severity, and for multiple thresholds are provided in Figure S4. The entire joint
226 distribution of duration and severity for all energy droughts in the data and the
227 generated simulations relative to a threshold for thresholds at the 25th, 30th,
228 35th, 40th percentile are shown in Figure S5. We see that KSTS is effective for
229 representing a range of energy droughts. Similar boxplot estimates (Figure 2) for
230 the KNN algorithm generated simulations are not shown since the simulations
231 show no occurrence of energy droughts at these thresholds unlike the reanalysis
232 data.

233 KSTS Reproduces the Aggregate Generation

234 The simulations from both KSTS and KNN reproduce temporal dynamics and
 235 data characteristics across both wind and solar fields at individual sites. The
 236 moments (mean and standard deviation), minimum and maximum for individual
 237 sites in KSTS and KNN simulations are representative of the underlying data
 238 (Figure S6). Both simulators are able to reproduce the quantiles (Figure S7,
 239 Figure S8), underlying probability distribution (Figure S9), auto-correlation
 240 structure (Figure S10), and site-level seasonality (Figure S11).

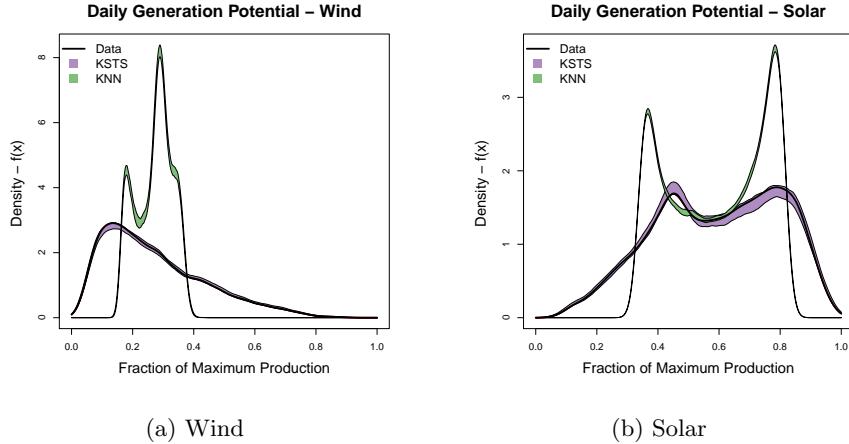


Figure 3: Kernel density estimate / Probability density function (PDF) of the daily aggregated energy production across the Texas Interconnection simulated using KSTS (purple) and KNN (green). The black line denotes the observed data pdf. The purple and green regions show the mid 90th (5th-95th) percentile interval regions from the individual pdfs computed from 48 simulations from each simulator. (a) Wind, (b) Solar. Equal land and capacity allocation is assumed across all the nodes/grids.

241 The kernel density estimate of aggregated daily energy generation potential
 242 across the Texas Interconnection is shown in Figure 3 for the historical reanalysis
 243 record (black), and for the KSTS (purple) and KNN (green) simulations.
 244 The degree to which adequate consideration of the spatial dependence and the
 245 wind-solar correlation leads to a proper representation of the potential for en-
 246 ergy production is illustrated through the fidelity of the KSTS simulations to
 247 the density function from the observations, and the marked departure of the
 248 KNN based simulations. It is clear that modeling spatial and cross field depen-
 249 dence is important to get the right frequency of the tail events (i.e., for LDS
 250 probabilities), even if the site-level production is adequately simulated without
 251 considering spatial dependence.

²⁵² **KSTS Reproduces Cross-Field Dependence**

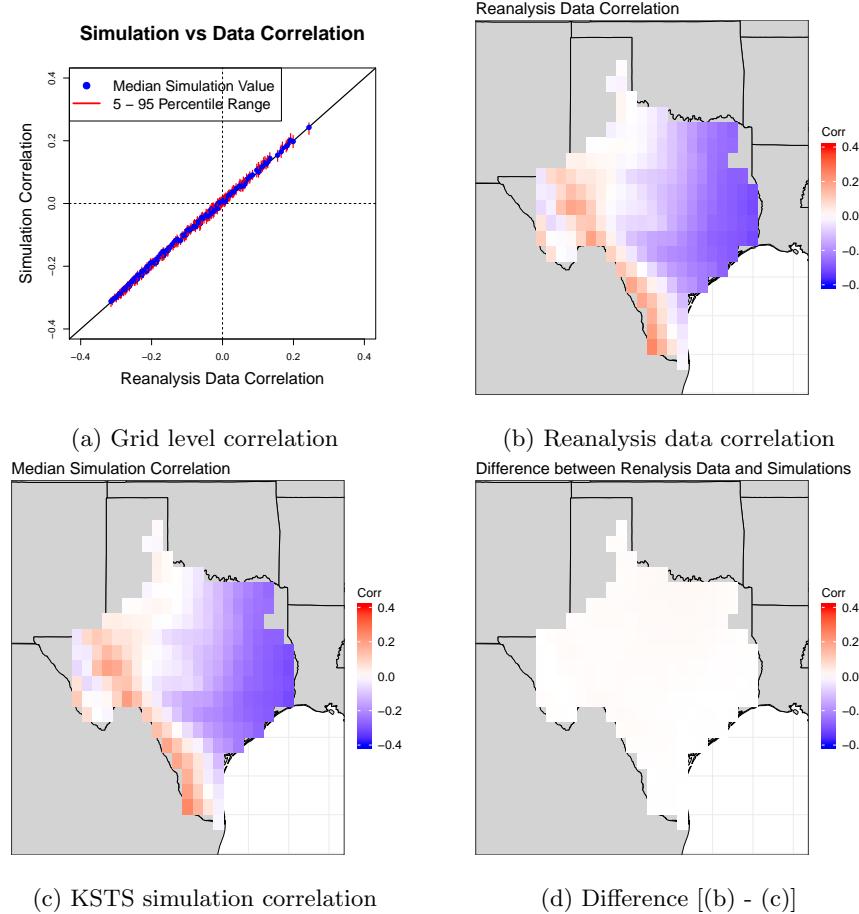


Figure 4: Pearson correlation between wind and solar at each grid point based on simultaneous simulations of wind and solar using KSTS. (a) Simulation correlation vs reanalysis data correlation between wind and solar where the red lines denote the mid 90th (5th-95th) percentile range and the blue dots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

²⁵³ From Figure 4 we note that the grid-wise correlation between daily wind and
²⁵⁴ solar across ERCOT is well reproduced by the KSTS simulations that are based
²⁵⁵ on simultaneous modeling of the wind and solar fields. By comparison, the KNN
²⁵⁶ simulations do not exhibit grid-wise wind-solar correlations consistent with the

257 reanalysis data (Figure S12). Furthermore, the spatial correlation structure
258 across all grids within a field for both wind and solar is also well reproduced by
259 the KSTS simulations unlike the KNN simulations (Figure S13). The seasonal
260 variation in the correlation between wind and solar is also well modeled by the
261 KSTS algorithm (Figure S14 and Figure S15).

262 Discussion

263 The primary contribution of this paper is the presentation of a novel k -Nearest
264 Neighbor Space Time Simulator (KSTS) and its application to the joint wind-
265 solar fields across the Texas Interconnection. We demonstrated the importance
266 of using a stochastic simulator that can properly reproduce the marginal prob-
267 ability densities of wind and solar at each site, as well as the cross-field spatial
268 dependence structure if good estimates of the severity-duration and frequency
269 of long duration renewable energy droughts, are of interest. So far, much of the
270 development of renewable electricity sources has focused on local microgrids,
271 but there has been growing interest in national and regional grids³⁶. As the
272 scale is increased, there is evidence that LDS is an effective and economic com-
273 ponent of the design of these regional systems^{10 11 37}. However, most of the
274 models developed and applied at these scales are deterministic and use rela-
275 tively short records³⁸. They do not consider the possible contracting structures
276 for guaranteed delivery and the associated default penalties. The probabilities
277 of the severity and duration of defaults as well as the penalties and LDS costs
278 would ultimately determine economically optimal resource allocations. We
279 anticipate and are planning to develop stochastic simulation-optimization models
280 to address a range of questions associated with such designs and contracts. The
281 KSTS simulator is motivated by this context, and it was important to under-
282 stand how important modeling spatial dependence to assess rare probabilities
283 of the duration and severity of shortages from a grid was for a real system.

284 From the application to the Texas Interconnection, we note that there is
285 substantial seasonal variability in the spatial expression of potential wind and
286 solar resource. This is not a surprise. The point by point wind-solar correlation
287 varies substantially by location and by season, as does the spatial correlation
288 structure for wind and solar and their cross-dependence. If these factors are
289 ignored, then the resulting regional LDS probability distributions are compro-
290 mised quite significantly. The simulations assuming a uniform land and capacity
291 allocation scenario show large uncertainties in the annual exceedance probabil-
292 ities for the severity-duration-threshold combinations considered, as well as po-
293 tentially higher exceedance probabilities than computed from the 40 year data
294 record for the more extreme threshold-severity-duration combinations.

295 The KSTS simulator is nonparametric and is appropriate for this setting
296 where the target variables are bounded with non-Gaussian distributions with
297 space and time dependence across variables changing by season. Since KSTS is
298 based on sampling the observed data, it can be thought of as a spatio-temporal
299 bootstrap procedure, where a spatio-temporal kernel is used at each time step

300 to sample a historical field with probabilities determined by the kernel and
301 a distance metric applied to the temporal state space for each variable. The
302 temporal sequences of potential energy produced at each site and across the
303 region are different even though the individual daily values are resampled from
304 the historical record. This allows the analysis of the range of drought severity-
305 duration-frequency using an extended sample. Extensions of the simulator to
306 hourly or other time scales are feasible. They would need to consider the diurnal
307 cycle as well as the seasonal cycle, and we are exploring computationally efficient
308 strategies for an algorithm that can address this while maintaining spatial and
309 cross field dependence.

310 The KSTS simulator exploits the similarity in the temporal evolution across
311 the fields and grid points. The potential next step would be developing an
312 algorithm which is capable of capturing the heterogeneity in dynamics across
313 even larger regions. This becomes important when the spatial scale of the
314 simulation is expanded from the Texas Interconnection to either the Western
315 or Eastern Interconnection or the entire North American continent. Such a
316 large scale makes it more likely that the wind and solar availability is driven by
317 different climate dynamics and consequently their temporal evolution structure
318 would be heterogeneous when compared to just Texas.

319 **Experimental Procedures**

320 **Resource Availability**

321 **Lead Contact**

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

324 **Materials Availability**

325 This study did not generate new unique materials.

326 **Data and Code Availability**

327 The KSTS and KNN generated simulations use wind and solar data spanning 40-
328 yrs (1979-2018) across the Texas Interconnection and are taken from the ERA-5
329 reanalysis dataset¹⁹, which can be accessed publicly. All code used in this study
330 is made publicly available on Github at <https://github.com/yashamonkar/LDS-Inferences>.

332 **Wind and Solar Data**

333 The variables used are wind speed at 100 meter altitude and downward surface
334 solar radiation from the ERA-5 reanalysis dataset¹⁹. The spatial grid size of
335 the data is set at 0.5° lat \times 0.5° lon and contains 216 grid points across the
336 Texas Interconnection domain (Figure S16).

Wind power is estimated by converting the 100 m wind speed to wind power using the wind turbine power curve from a V90-2.0MW Vestas turbine (as shown in Fig S17). We then estimate wind power availability per unit area (W/m^2) by assuming a required land area per turbine equal to four times the rotor diameter in one direction and seven times the rotor diameter in the other direction³⁹. This means that each turbine is assumed to take up 0.2268 km^2 given the rotor diameter of 90 m for the V90-2.0MW Vestas turbine. The data are converted to the daily time step by taking the mean of the hourly data for each day and the dataset spans the 40 years from January 1st, 1979 to December 31st, 2018.

The solar variable is the downward surface solar radiation (W/m^2) and is converted to the daily time step by taking the mean of the hourly data. The solar data are not converted to capacity factors. The conversion to capacity factors or solar power involves making assumptions about the PV technology, tilt and tracking capabilities of the array and are not made in this study.

Our example application assumes a uniform land allocation for both wind and solar generators.

Energy Deficits and Drought Spells

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

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

where, y_t is the aggregated daily energy deviation at day t ; $x_{i,t}$ is the normalized wind or solar value at site i and day t ; $\widetilde{x_{i,T}}$ is the normalized DOY percentile based on the selected threshold for site i and day DOY(t); n is the total number of grid points (216) times the fields (wind and solar). The aggregated deviation y_t can take a positive (surplus) or negative (deficit) value on any particular day, while the cumulative deficit, the variable of interest is computed as,

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

where, z_t and y_t are the cumulative deficit and daily deviation at day t respectively. While y_t can either be positive or negative, the cumulative deficit takes a lower value of 0 (surplus) and is restricted to positive values (periods of energy deficit). Energy Droughts for a selected threshold percentile are defined to occur during instances of consecutive days with positive values of cumulative deficit. Severity of a drought event is defined as the maximum cumulative deficit during the drought period, while the duration is the spell length in days.

373 **Annual Exceedance Probability**

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

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

379 where, $C(e_i)$ is the count of exceedances for drought event i with duration
380 d_i and severity s_i and n is the total number of drought events. The count
381 of exceedances $C(e_i)$ is regressed against the severity s_i and duration d_i using
382 Poisson regression. The methodology used is local regression using the locfit
383 package³³.

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

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

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

391 **Simulator Hyperparameters**

392 For both KSTS and KNN methods, the seasonality in the data is accounted
393 for by restricting the search of nearest neighbors to a ± 30 days moving window
394 across the years around the day of the year (DOY). The number of nearest
395 neighbors (k) selected is approximately \sqrt{n} . With 40 years and 61 days per
396 year, \sqrt{n} is ~ 50 , where n is the number of possible candidate neighbors after
397 accounting for the moving window²⁸. A lag-1 dependence structure for the
398 state space is assumed. 48 independent simulation realizations, each of the
399 same length as the reanalysis data (40 yrs = 14610 days), are generated using
400 the KSTS and KNN algorithms.

401 The KNN algorithm is fit to each grid individually using the hyper-parameters
402 specified above with the algorithm which is outlined in Lall and Sharma²⁸.

403 **k-Nearest Neighbor Space-Time Simulator (KSTS)**

404 The general structure and a cartoon example application of the KSTS algorithm
405 is illustrated in Figures S18 and 5 respectively. The algorithm leads to a space-
406 time simulation process that is Markovian (or corresponds to a state space
407 formed by the embedding) in time.

408 **KSTS Algorithm**

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

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

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

417 Case 1 represents simple dependence on the two previous values. Case 2
 418 represents dependence on the past two values and values 12 and 24 steps be-
 419 fore the current value allowing for monthly and inter-annual dependence for
 420 monthly data. Case 3 represents incorporation of a temporal dependence struc-
 421 ture unique to the data. The state space $D_{i,t}$ is defined for each time series at
 422 site i and time t , whereas $D_{i,T}$ are all the historic vectors which correspond to
 423 the selected embedding structure for site i .

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

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

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

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

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

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

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

438

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

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

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

446 where s is the total number of sites. The similarity matrix S_t has the same
 447 length as the number of unique time indices in the data.

448 **Step 5:- Curtail and scale the similarity matrix S_t .**

449 The similarity matrix S_t is ordered and curtailed to its highest k values.
 450 The time indices associated the k highest values of S_t are selected as the k-
 451 nearest neighbor candidates for the entire spatial field. The probabilities of the
 452 associated k neighbors are scaled to add up to 1.

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

453

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

455 Using the discrete probability mass function S_t , sample a single value and
 456 re-sample entire fields across all sites from the time index which corresponds to
 457 the selected value in S_t as data for the simulation at time $t + 1$. Return to Step
 458 2 if further time-steps are needed for the simulation..

459 Refer supplementary materials for further details on the algorithm and hyper-
 460 parameter selection.

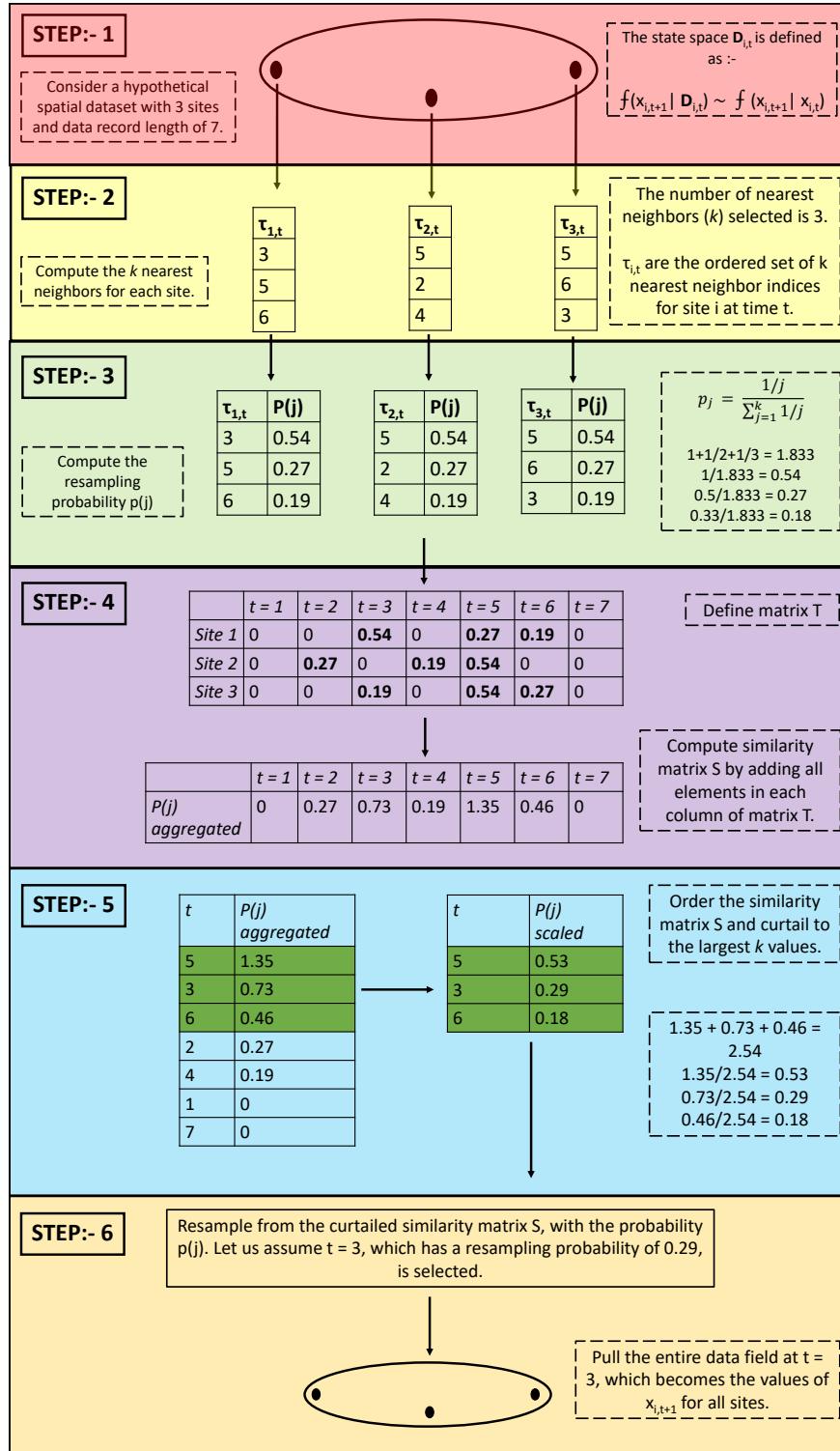


Figure 5: Cartoon of the KSTS algorithm for an example application with a spatial dataset consisting of 3 grids/sites and data record (time) length of 7. ¹⁸

461 **Acknowledgment**

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

465 **Author Contributions**

466 Y.A developed the code and performed the computations. Y.A and D.J.F de-
467 signed the analysis, conceived experiments and simulation checks with super-
468 vision from U.L. D.J.F provided the data. Y.A took the lead in writing the
469 manuscript with all authors discussing and contributing to the final manuscript.

470 **Declaration of Interests**

471 The authors declare no competing interests.

472 **References**

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

- 487 [7] Marc Beaudin, Hamidreza Zareipour, Anthony Schellenberglabe,
488 and William Rosehart. Energy storage for mitigating the vari-
489 ability of renewable electricity sources: An updated review.
490 *Energy for Sustainable Development*, 14(4):302–314, December
491 2010. ISSN 0973-0826. doi: 10.1016/j.esd.2010.09.007. URL
492 <https://www.sciencedirect.com/science/article/pii/S0973082610000566>.
- 493 [8] Jacques Després, Silvana Mima, Alban Kitous, Patrick Criqui, Noure-
494 dine Hadjsaid, and Isabelle Noirot. Storage as a flexibility option
495 in power systems with high shares of variable renewable energy
496 sources: a POLES-based analysis. *Energy Economics*, 64:638–650,
497 May 2017. ISSN 0140-9883. doi: 10.1016/j.eneco.2016.03.006. URL
498 <https://www.sciencedirect.com/science/article/pii/S0140988316300445>.
- 499 [9] Jesse D. Jenkins, Max Luke, and Samuel Thernstrom. Getting to Zero
500 Carbon Emissions in the Electric Power Sector. *Joule*, 2(12):2498–2510,
501 December 2018. ISSN 2542-4351. doi: 10.1016/j.joule.2018.11.013. URL
502 <https://www.sciencedirect.com/science/article/pii/S2542435118305622>.
- 503 [10] Matthew R. Shaner, Steven J. Davis, Nathan S. Lewis, and Ken
504 Caldeira. Geophysical constraints on the reliability of solar and wind
505 power in the United States. *Energy & Environmental Science*, 11(4):
506 914–925, April 2018. ISSN 1754-5706. doi: 10.1039/C7EE03029K. URL
507 <https://pubs.rsc.org/en/content/articlelanding/2018/ee/c7ee03029k>.
508 Publisher: The Royal Society of Chemistry.
- 509 [11] Jacqueline A. Dowling, Katherine Z. Rinaldi, Tyler H. Ruggles,
510 Steven J. Davis, Mengyao Yuan, Fan Tong, Nathan S. Lewis, and
511 Ken Caldeira. Role of Long-Duration Energy Storage in Variable
512 Renewable Electricity Systems. *Joule*, 4(9):1907–1928, Septem-
513 ber 2020. ISSN 2542-4351. doi: 10.1016/j.joule.2020.07.007. URL
514 <https://www.sciencedirect.com/science/article/pii/S2542435120303251>.
- 515 [12] ARPA-E. Duration Addition to electricitY Storage, 2018. URL
516 <https://arpa-e.energy.gov/technologies/programs/days>.
- 517 [13] Seán Collins, Paul Deane, Brian Ó Gallachóir, Stefan Pfenninger,
518 and Iain Staffell. Impacts of Inter-annual Wind and Solar Varia-
519 tions on the European Power System. *Joule*, 2(10):2076–2090, Octo-
520 ber 2018. ISSN 2542-4351. doi: 10.1016/j.joule.2018.06.020. URL
521 <https://www.sciencedirect.com/science/article/pii/S254243511830285X>.
- 522 [14] Luc Bonnafous, Upmanu Lall, and Jason Siegel. A water risk index for
523 portfolio exposure to climatic extremes: conceptualization and an appli-
524 cation to the mining industry. *Hydrology and Earth System Sciences*, 21
525 (4):2075–2106, April 2017. ISSN 1027-5606. doi: 10.5194/hess-21-2075-
526 2017. URL <https://hess.copernicus.org/articles/21/2075/2017/>.
527 Publisher: Copernicus GmbH.

- 528 [15] Shaleen Jain and Upmanu Lall. Floods in a changing climate: Does the
 529 past represent the future? *Water Resources Research*, 37(12):3193–3205,
 530 2001. ISSN 1944-7973. doi: <https://doi.org/10.1029/2001WR000495>. URL
 531 <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2001WR000495>.
 532 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2001WR000495>.
- 533 [16] James Doss-Gollin, David J. Farnham, Scott Steinschneider,
 534 and Upmanu Lall. Robust Adaptation to Multiscale Cli-
 535 mate Variability. *Earth's Future*, 7(7):734–747, 2019. ISSN
 536 2328-4277. doi: <https://doi.org/10.1029/2019EF001154>. URL
 537 <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019EF001154>.
 538 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2019EF001154>.
- 539 [17] David J. Farnham. *Identifying and Modeling Spatio-temporal Structures in*
 540 *High Dimensional Climate and Weather Datasets with Applications to Wa-*
 541 *ter and Energy Resource Management*. PhD thesis, Columbia University,
 542 2018. URL <https://doi.org/10.7916/D8321CTB>.
- 543 [18] Scott Chamberlain. 'NOAA' Weather Data from R
 544 [R package rnoaa version 1.3.2], February 2021. URL
 545 <https://CRAN.R-project.org/package=rnoaa>. Publisher: Compre-
 546 hensive R Archive Network (CRAN).
- 547 [19] Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András
 548 Horányi, Joaquín Muñoz-Sabater, Julien Nicolas, Carole Peubey, Raluca
 549 Radu, Dinand Schepers, Adrian Simmons, Cornel Soci, Saleh Abdalla,
 550 Xavier Abellan, Gianpaolo Balsamo, Peter Bechtold, Gionata Biavati,
 551 Jean Bidlot, Massimo Bonavita, Giovanna De Chiara, Per Dahlgren,
 552 Dick Dee, Michail Diamantakis, Rossana Dragani, Johannes Flemming,
 553 Richard Forbes, Manuel Fuentes, Alan Geer, Leo Haimberger,
 554 Sean Healy, Robin J. Hogan, Elías Hólm, Marta Janisková,
 555 Sarah Keeley, Patrick Laloyaux, Philippe Lopez, Cristina Lupu, Gabor
 556 Radnoti, Patricia de Rosnay, Iryna Rozum, Freja Vamborg, Sébastien
 557 Villaume, and Jean-Noël Thépaut. The ERA5 global reanalysis. *Quar-
 558 terly Journal of the Royal Meteorological Society*, 146(730):1999–2049,
 559 2020. ISSN 1477-870X. doi: <https://doi.org/10.1002/qj.3803>. URL
 560 <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3803>.
 561 _eprint: <https://rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/qj.3803>.
- 562 [20] Patrick Laloyaux, Eric de Boisseson, Magdalena Balmaseda, Jean-
 563 Raymond Bidlot, Stefan Broennimann, Roberto Buizza, Per Dalgren,
 564 Dick Dee, Leopold Haimberger, Hans Hersbach, Yuki Kosaka, Matthew
 565 Martin, Paul Poli, Nick Rayner, Elke Rustemeier, and Dinand Schep-
 566 pers. CERA-20C: A Coupled Reanalysis of the Twentieth Century.
 567 *Journal of Advances in Modeling Earth Systems*, 10(5):1172–1195,
 568 2018. ISSN 1942-2466. doi: <https://doi.org/10.1029/2018MS001273>. URL
 569 <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001273>.
 570 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2018MS001273>.

- 571 [21] Ronald Gelaro, Will McCarty, Max J. Suárez, Ricardo Todling, An-
 572 drea Molod, Lawrence Takacs, Cynthia A. Randles, Anton Darmenov,
 573 Michael G. Bosilovich, Rolf Reichle, Krzysztof Wargan, Lawrence Coy,
 574 Richard Cullather, Clara Draper, Santha Akella, Virginie Buchard,
 575 Austin Conaty, Arlindo M. da Silva, Wei Gu, Gi-Kong Kim, Ran-
 576 dal Koster, Robert Lucchesi, Dagmar Merkova, Jon Eric Nielsen,
 577 Gary Partyka, Steven Pawson, William Putman, Michele Rienecker,
 578 Siegfried D. Schubert, Meta Sienkiewicz, and Bin Zhao. The Modern-
 579 Era Retrospective Analysis for Research and Applications, Version
 580 2 (MERRA-2). *Journal of Climate*, 30(14):5419–5454, July 2017.
 581 ISSN 0894-8755, 1520-0442. doi: 10.1175/JCLI-D-16-0758.1. URL
 582 <https://journals.ametsoc.org/view/journals/clim/30/14/jcli-d-16-0758.1.xml>.
 583 Publisher: American Meteorological Society Section: Journal of Climate.
- 584 [22] D. P. Dee, S. M. Uppala, A. J. Simmons, P. Berrisford, P. Poli,
 585 S. Kobayashi, U. Andrae, M. A. Balmaseda, G. Balsamo, P. Bauer, P. Bech-
 586 told, A. C. M. Beljaars, L. van de Berg, J. Bidlot, N. Bormann, C. Delsol,
 587 R. Dragani, M. Fuentes, A. J. Geer, L. Haimberger, S. B. Healy, H. Hers-
 588 bach, E. V. Hólm, L. Isaksen, P. Kållberg, M. Köhler, M. Matricardi, A. P.
 589 McNally, B. M. Monge-Sanz, J.-J. Morcrette, B.-K. Park, C. Peubey, P. de
 590 Rosnay, C. Tavolato, J.-N. Thépaut, and F. Vitart. The ERA-Interim
 591 reanalysis: configuration and performance of the data assimilation system.
 592 *Quarterly Journal of the Royal Meteorological Society*, 137(656):553–597,
 593 2011. ISSN 1477-870X. doi: <https://doi.org/10.1002/qj.828>. URL
 594 <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.828>.
 595 _eprint: <https://rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/qj.828>.
- 596 [23] Andreas Sterl. On the (In)Homogeneity of Reanalysis Products. *Journal*
 597 *of Climate*, 17(19):3866–3873, October 2004. ISSN 0894-8755, 1520-
 598 0442. doi: 10.1175/1520-0442(2004)017j3866:OTIOP;2.0.CO;2. URL
 599 https://journals.ametsoc.org/view/journals/clim/17/19/1520-0442_2004_017j3866_0tiorp2.0.co2.xml.
 600 American Meteorological Society Section : Journal of Climate.
- 601 [24] Katherine Z. Rinaldi, Jacqueline A. Dowling, Tyler H. Ruggles, Ken
 602 Caldeira, and Nathan S. Lewis. Wind and Solar Resource Droughts in
 California Highlight the Benefits of Long-Term Storage and Integration
 603 with the Western Interconnect. *Environmental Science & Technology*, 55
 604 (9):6214–6226, May 2021. ISSN 0013-936X. doi: 10.1021/acs.est.0c07848.
 605 URL <https://doi.org/10.1021/acs.est.0c07848>. Publisher: American
 606 Chemical Society.
- 607 [25] Electric Reliability Council of Texas. About ERCOT, 2021. URL
 608 <http://www.ercot.com/about>.
- 609 [26] Electric Reliability Council of Texas. Impact of increased wind
 610 resrouces in the ERCOT region. Technical report, June 2020. URL
 611 <http://www.ercot.com/content/wcm/lists/200196/WindOnePagerJune2020.pdf>.

- 605 [27] Andrew Kumler, Ignacio Losada Carreño, Michael T. Craig, Bri-
 606 Mathias Hodge, Wesley Cole, and Carlo Brancucci. Inter-annual
 607 variability of wind and solar electricity generation and capacity val-
 608 ues in Texas. *Environmental Research Letters*, 14(4):044032, April
 609 2019. ISSN 1748-9326. doi: 10.1088/1748-9326/aaf935. URL
 610 <https://doi.org/10.1088/1748-9326/aaf935>. Publisher: IOP Publish-
 611 ing.
- 612 [28] Upmanu Lall and Ashish Sharma. A Nearest Neighbor
 613 Bootstrap For Resampling Hydrologic Time Series. *Wa-
 614 ter Resources Research*, 32(3):679–693, 1996. ISSN 1944-
 615 7973. doi: <https://doi.org/10.1029/95WR02966>. URL
 616 <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/95WR02966>.
 617 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/95WR02966>.
- 618 [29] Balaji Rajagopalan and Upmanu Lall. A k-nearest-neighbor sim-
 619 ulator for daily precipitation and other weather variables. *Wa-
 620 ter Resources Research*, 35(10):3089–3101, 1999. ISSN 1944-
 621 7973. doi: <https://doi.org/10.1029/1999WR900028>. URL
 622 <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/1999WR900028>.
 623 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/1999WR900028>.
- 624 [30] James Prairie, Balaji Rajagopalan, Upmanu Lall, and Terrance
 625 Fulp. A stochastic nonparametric technique for space-time disag-
 626 gregation of streamflows. *Water Resources Research*, 43(3), 2007.
 627 ISSN 1944-7973. doi: <https://doi.org/10.1029/2005WR004721>. URL
 628 <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2005WR004721>.
 629 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2005WR004721>.
- 630 [31] James Prairie, Kenneth Nowak, Balaji Rajagopalan, Upmanu
 631 Lall, and Terrance Fulp. A stochastic nonparametric approach
 632 for streamflow generation combining observational and paleore-
 633 constructed data. *Water Resources Research*, 44(6), 2008. ISSN
 634 1944-7973. doi: <https://doi.org/10.1029/2007WR006684>. URL
 635 <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2007WR006684>.
 636 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2007WR006684>.
- 637 [32] Francisco Assis Souza Filho and Upmanu Lall. Seasonal to interannual
 638 ensemble streamflow forecasts for Ceará, Brazil: Applications of a mul-
 639 tivariate, semiparametric algorithm. *Water Resources Research*, 39(11),
 640 2003. ISSN 1944-7973. doi: <https://doi.org/10.1029/2002WR001373>. URL
 641 <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2002WR001373>.
 642 _eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2002WR001373>.
- 643 [33] Catherine Loader. Local Regression, Likelihood and Density Es-
 644 timation [R package locfit version 1.5-9.4], March 2020. URL
 645 <https://CRAN.R-project.org/package=locfit>. Publisher: Comprehen-
 646 sive R Archive Network (CRAN).

- 647 [34] Gary D. Tasker. Effective record length for the T-year
 648 event. *Journal of Hydrology*, 64(1):39–47, July 1983. ISSN
 649 0022-1694. doi: 10.1016/0022-1694(83)90059-8. URL
 650 <https://www.sciencedirect.com/science/article/pii/0022169483900598>.
- 651 [35] Robert Link, Thomas B. Wild, Abigail C. Snyder, Mohamad I. Hejazi,
 652 and Chris R. Vernon. 100 years of data is not enough to estab-
 653 lish reliable drought thresholds. *Journal of Hydrology X*, 7:100052,
 654 April 2020. ISSN 2589-9155. doi: 10.1016/j.hydroa.2020.100052. URL
 655 <https://www.sciencedirect.com/science/article/pii/S2589915520300031>.
- 656 [36] Patricia J. Levi, Simon Davidsson Kurland, Michael Carbajales-Dale,
 657 John P. Weyant, Adam R. Brandt, and Sally M. Benson. Macro-
 658 Energy Systems: Toward a New Discipline. *Joule*, 3(10):2282–2286,
 659 October 2019. ISSN 2542-4351. doi: 10.1016/j.joule.2019.07.017. URL
 660 <https://www.sciencedirect.com/science/article/pii/S2542435119303617>.
- 661 [37] Micah S. Ziegler, Joshua M. Mueller, Gonçalo D. Pereira, Juhyun
 662 Song, Marco Ferrara, Yet-Ming Chiang, and Jessika E. Trancik.
 663 Storage Requirements and Costs of Shaping Renewable Energy
 664 Toward Grid Decarbonization. *Joule*, 3(9):2134–2153, Septem-
 665 ber 2019. ISSN 2542-4351. doi: 10.1016/j.joule.2019.06.012. URL
 666 <https://www.sciencedirect.com/science/article/pii/S2542435119303009>.
- 667 [38] Leonard Göke and Mario Kendziora. The adequacy of time-series reduc-
 668 tion for renewable energy systems. *arXiv:2101.06221 [econ, q-fin]*, January
 669 2021. URL <http://arxiv.org/abs/2101.06221>. arXiv: 2101.06221.
- 670 [39] Xi Lu and Michael B. McElroy. Chapter 4 - Global Potential for
 671 Wind-Generated Electricity. In Trevor M. Letcher, editor, *Wind*
 672 *Energy Engineering*, pages 51–73. Academic Press, January 2017.
 673 ISBN 978-0-12-809451-8. doi: 10.1016/B978-0-12-809451-8.00004-7. URL
 674 <https://www.sciencedirect.com/science/article/pii/B9780128094518000047>.

675 Supplementary Materials

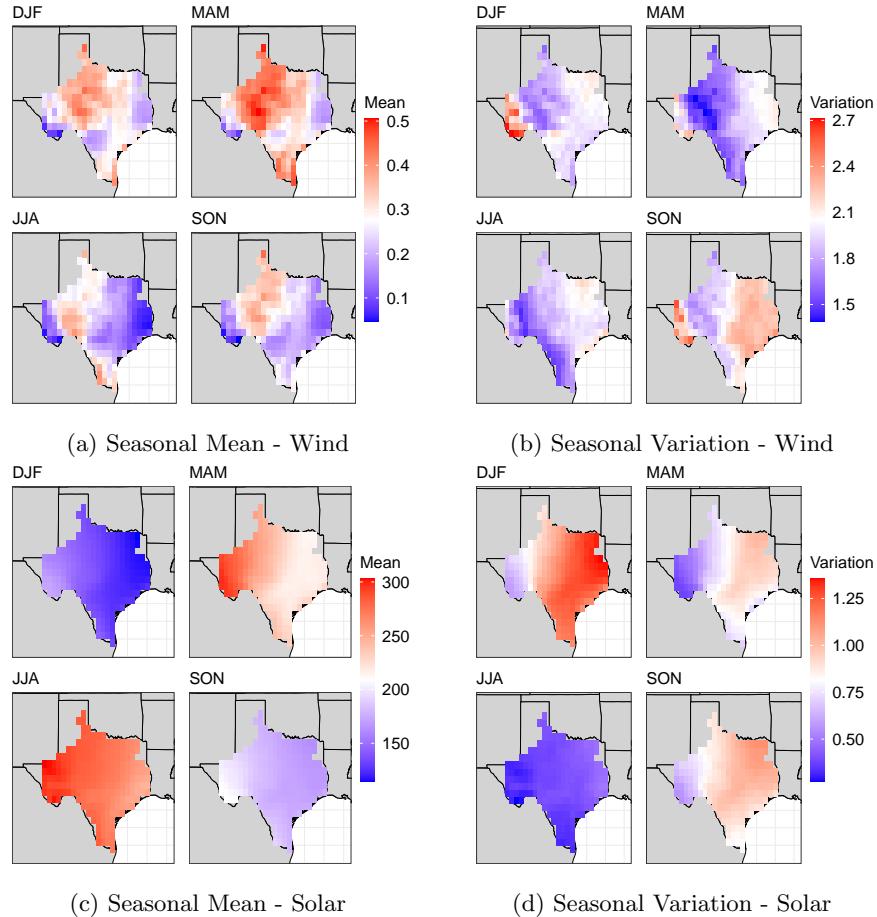


Figure S1: Seasonal mean and variation in daily wind capacity factors and downward surface solar radiation ($\text{W}/\text{m} \cdot \text{sq}$) across the Texas Interconnection. (a) Mean daily wind capacity factors by season. (b) Variation in daily wind capacity factors by season. (c) Mean downward surface solar radiation by season. (d) Variation in downward surface solar radiation 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).

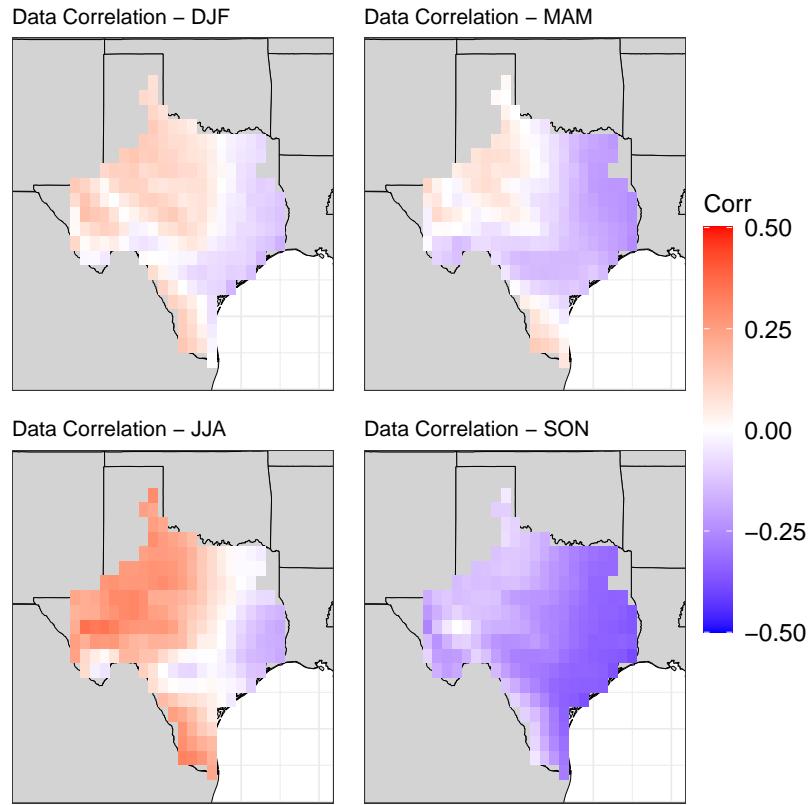


Figure S2: Seasonal correlation between daily wind capacity factors and downward surface solar radiation in the ERA-5 reanalysis dataset at each grid point. (top-left) Dec-Jan-Feb (DJF). (top-right) Mar-Apr-May (MAM). (bottom-left) Jun-Jul-Aug (JJA). (bottom-right) Sep-Oct-Nov (SON). The correlations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

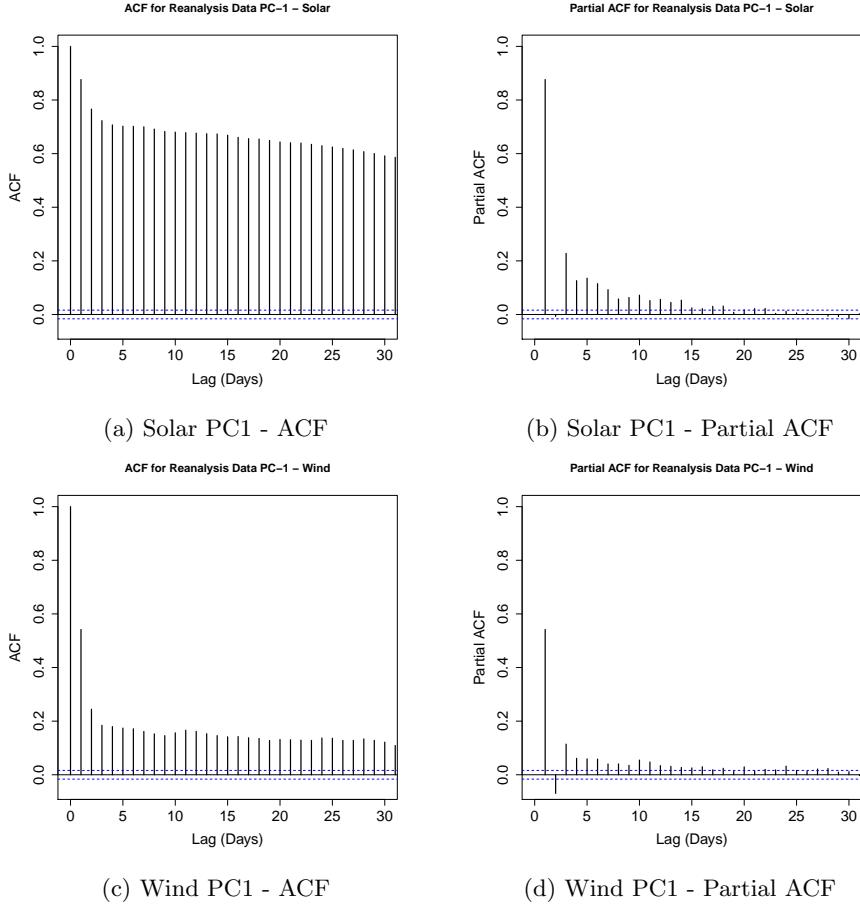


Figure S3: The auto-correlation (ACF) and partial auto-correlation (PACF) for PC-1 of the reanalysis data. (a) Solar PC-1 ACF (b) Solar PC-1 partial ACF (c) Wind PC-1 ACF (d) Wind PC-1 partial ACF. The fractional variance explained by PC-1 for solar and wind fields is 80 % and 63 % respectively. The blue dashed line denotes the significance level for the record length of 14610 days (40 yrs).

Annual Exceedance for Threshold-Duration-Severity

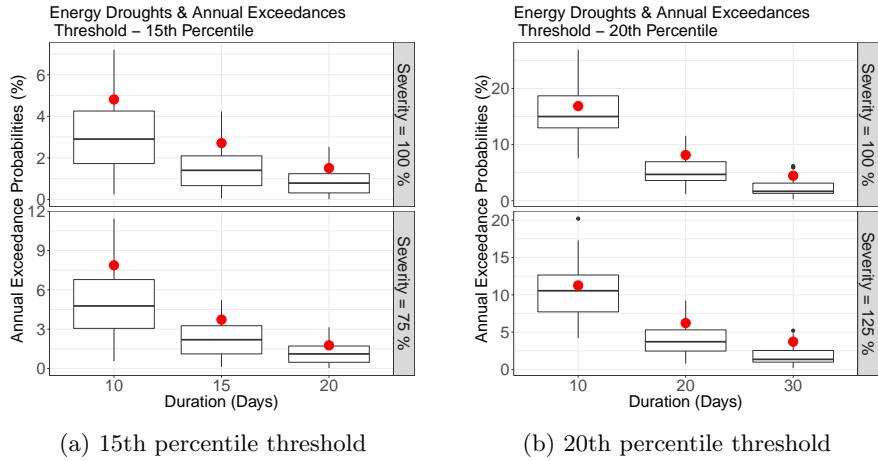


Figure S4: Probability of joint annual exceedances for energy droughts given a duration and severity with threshold values of (a) 15th percentile, and (b) 20th percentile. The red dot denotes the exceedance probability calculated from the reanalysis data. The boxplots denote the uncertainty in the 48 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.

677 Severity vs Duration plots for different thresholds

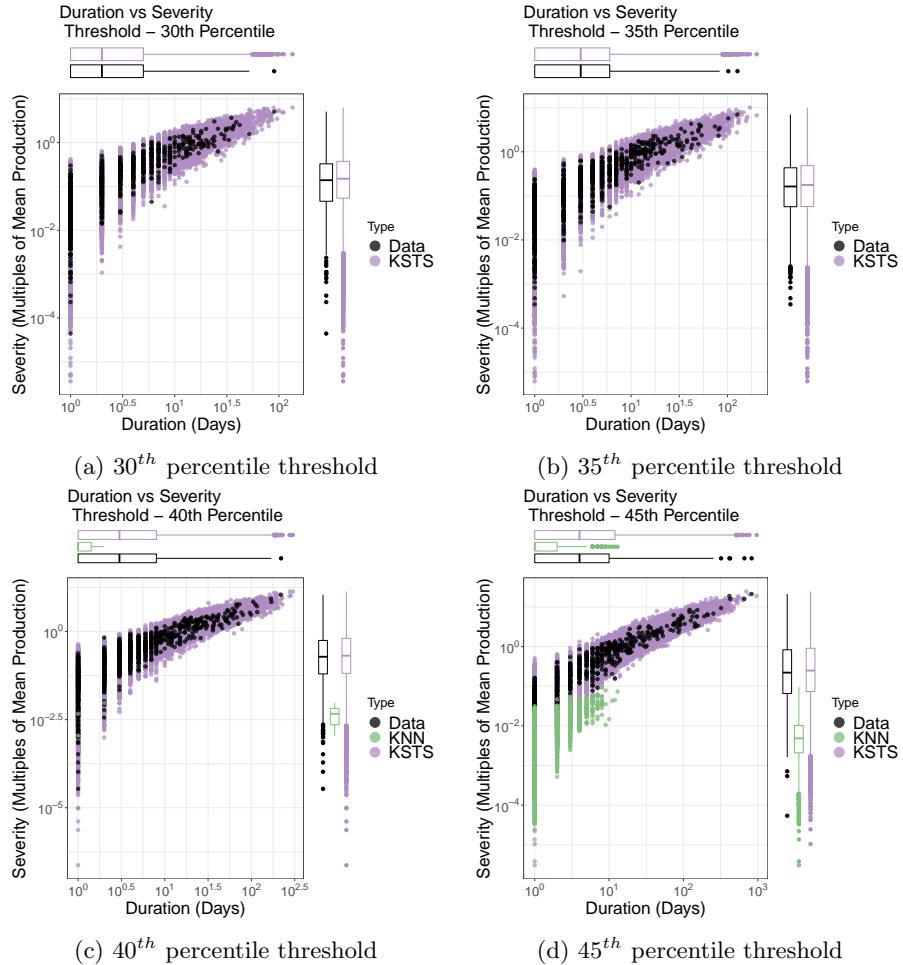


Figure S5: Severity versus duration for all energy droughts with marginal distributions (in boxplots for both variables) for the data (black), KSTS (purple) and KNN (green) simulations using threshold values of (a) 30th percentile, (b) 35th percentile, (c) 40th percentile, and (d) 45th percentile. For the lower thresholds, the KNN simulations do not have occurrence of energy droughts and hence do not appear in the plots.

678 KSTS and KNN Simulations - Individual Site 679 Characteristics

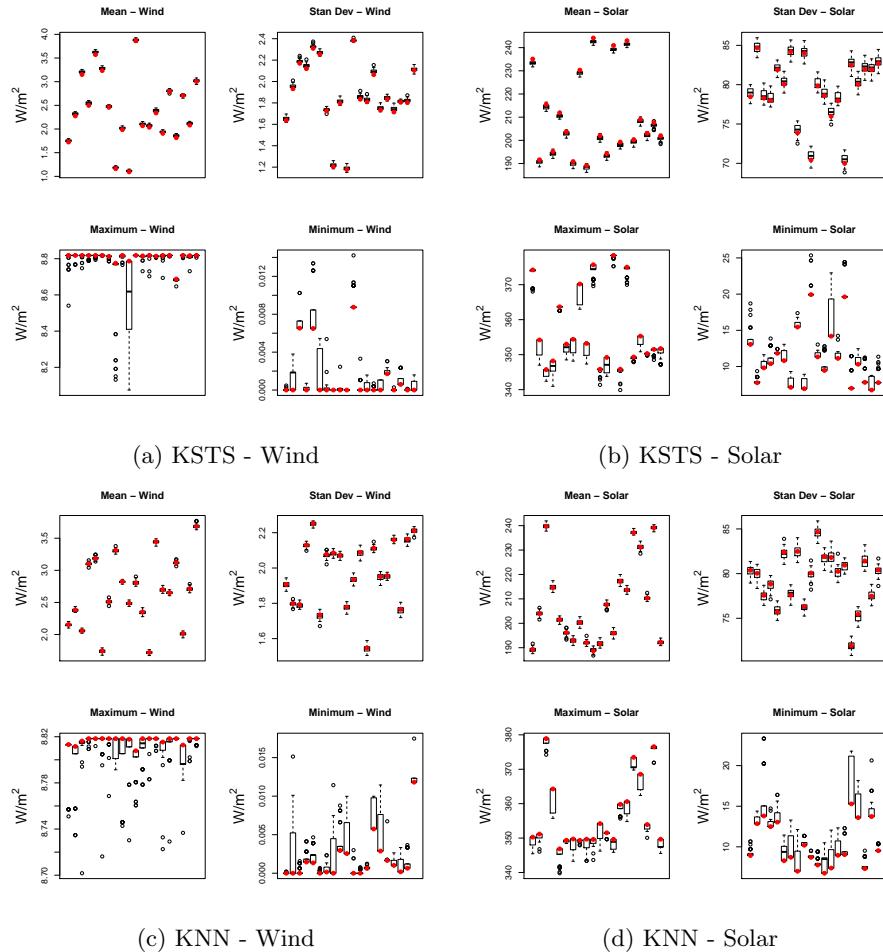


Figure S6: Simulation skill assessments for individual sites in the wind and solar fields for both KSTS and KNN simulations. (a) KSTS wind. (b) KSTS solar. (c) KNN wind. (d) KNN solar. For each sub-plot, we show the mean (top-left), the standard deviation (top-right), the minimum (bottom-left) and the 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.

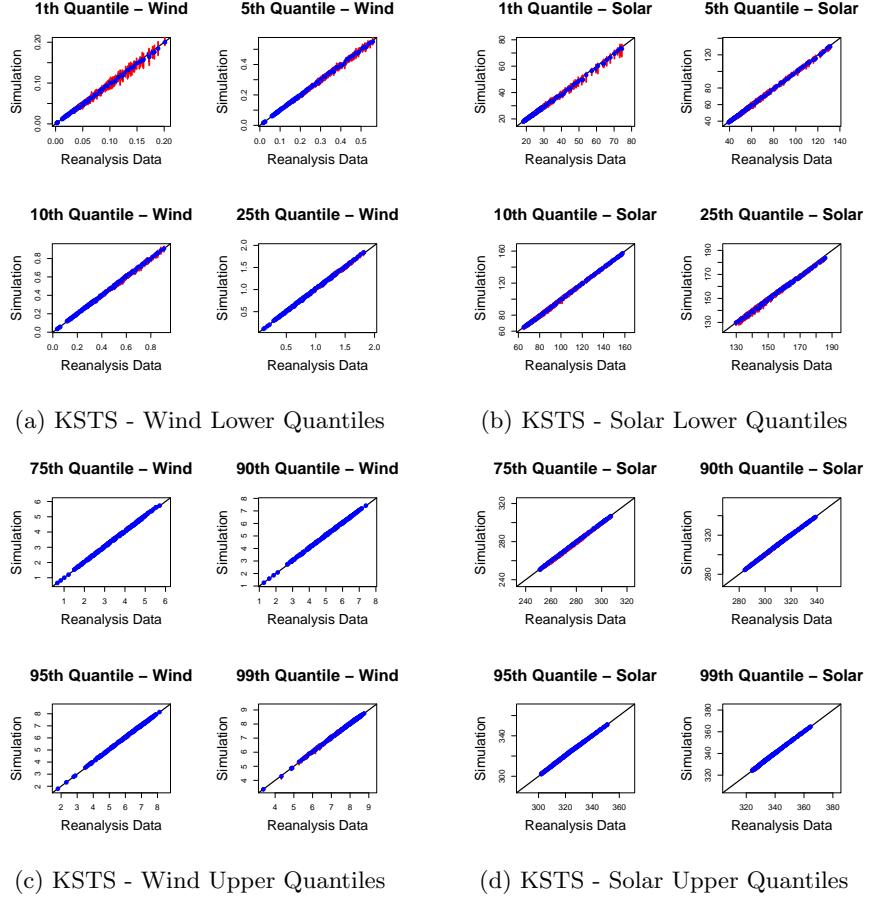


Figure S7: Simulation vs reanalysis data quantiles plots for the 1st, 5th, 10th, 25th, 75th, 90th, 95th and 99th percentiles for KSTS simulations. The plots denote the quantiles for all 216 grid points in the wind and solar fields. The red lines denote the mid 90th (5th-95th) percentile range and the blue dots denote the median value in the simulation spread.

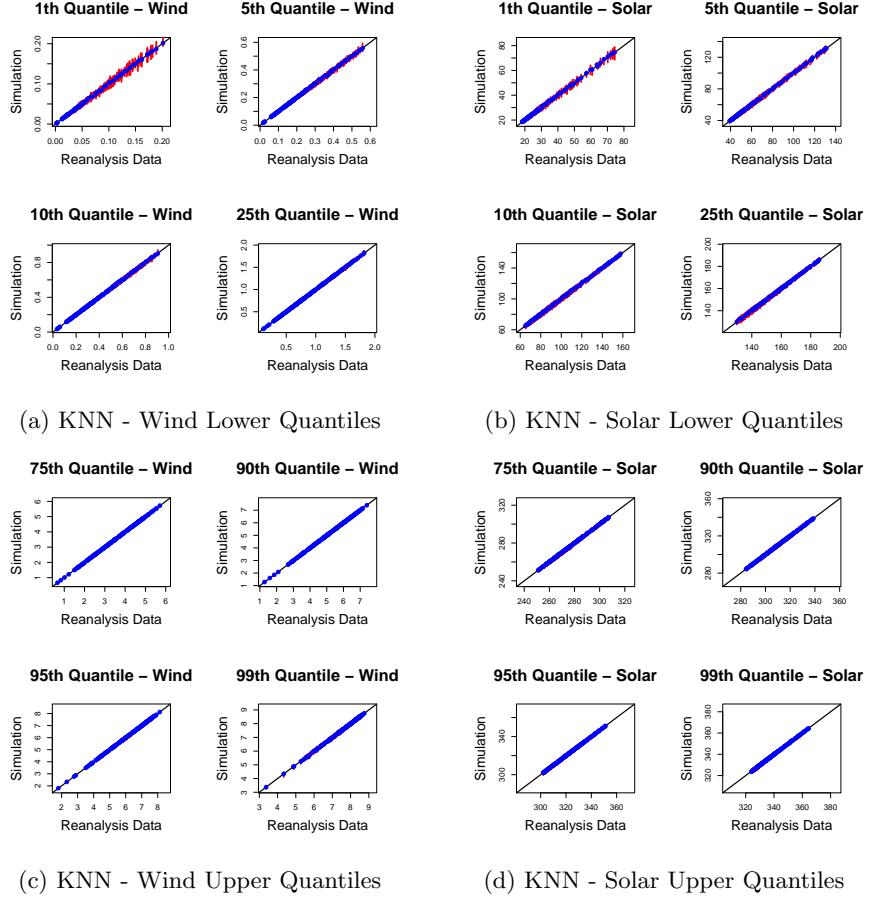


Figure S8: Simulation vs reanalysis data quantiles plots for the 1st, 5th, 10th, 25th, 75th, 90th, 95th and 99th percentiles for KNN simulations. The plots denote the quantiles for all 216 grid points in the wind and solar fields. The red lines denote the mid 90th (5th-95th) percentile range and the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

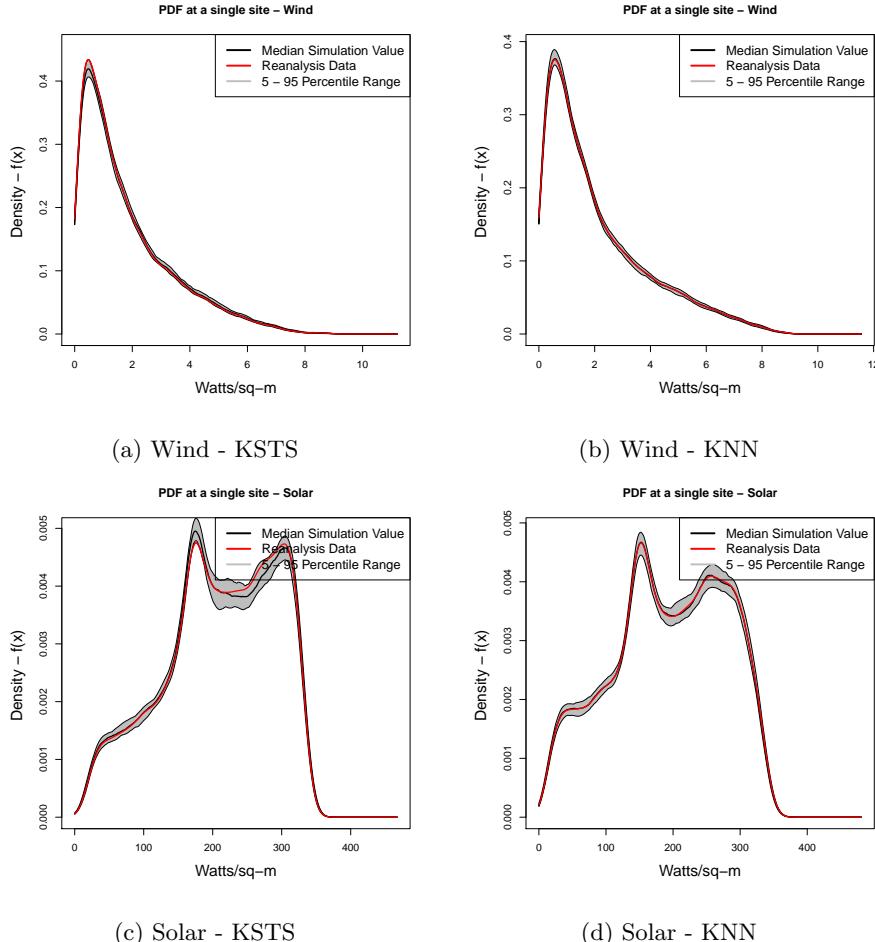


Figure S9: Kernel density estimate / Probability density function (PDF) plots for a single randomly selected grid for wind and solar. The red line denotes the reanalysis data probability density function for the selected site and the black line denotes the median simulation density. The grey region is the mid 90th (5th-95th) percentile range of the simulation spread. The grid that is shown in this figure 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.

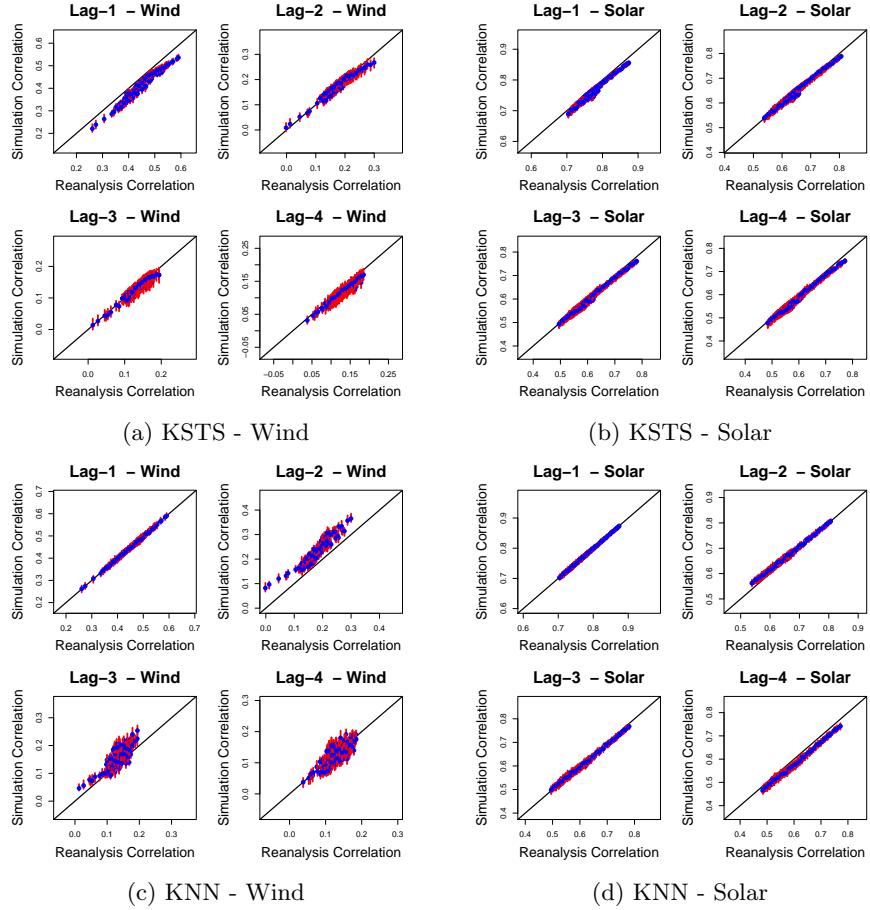


Figure S10: Simulation vs reanalysis data auto-correlation plots for lag 1,2,3 and 4 for all grid points. (a) Wind KSTS simulations. (b) Solar KSTS simulations. (c) Wind KNN simulations. (d) Solar KNN simulations. The red lines denote the mid 90th (5th-95th) percentile range and the blue dots denote the median value in the simulation spread.

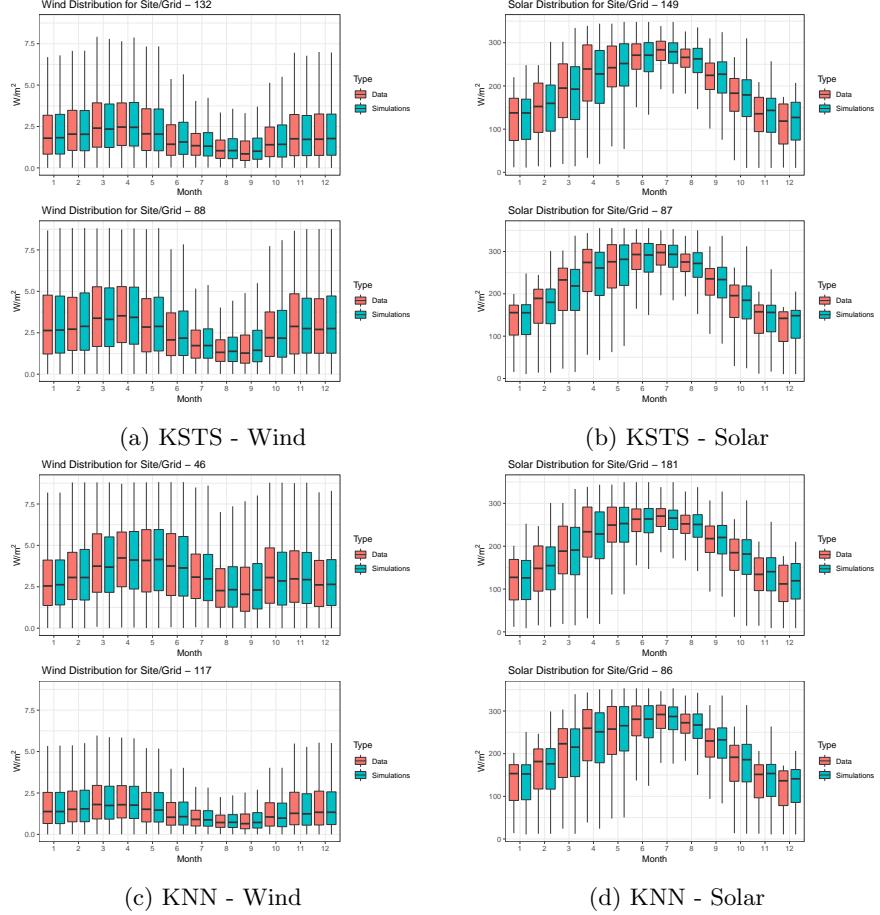


Figure S11: Seasonality / Monthly distribution of the reanalysis data and simulations. The red and green boxplots denote the reanalysis data and simulations respectively. (a) Wind KSTS simulations. (b) Solar KSTS simulations. (c) Wind KNN simulations. (d) Solar KNN simulations. Two grids are randomly 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

680 Cross-Field Dependence

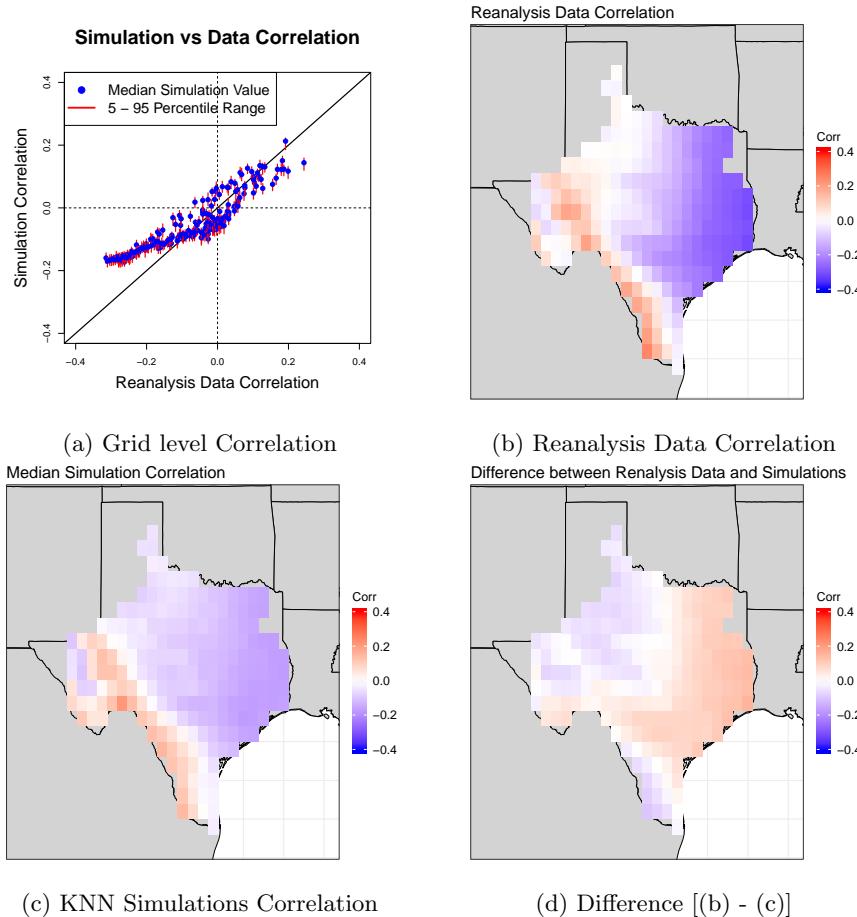


Figure S12: Pearson correlation between wind and solar at each grid point based on simultaneous simulations of wind and solar using KNN. (a) Simulation correlation vs reanalysis data correlation between wind and solar where the red lines denote the mid 90th (5th-95th) percentile range and the blue dots 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).

681 Individual Field (Wind and Solar) Spatial Correlations

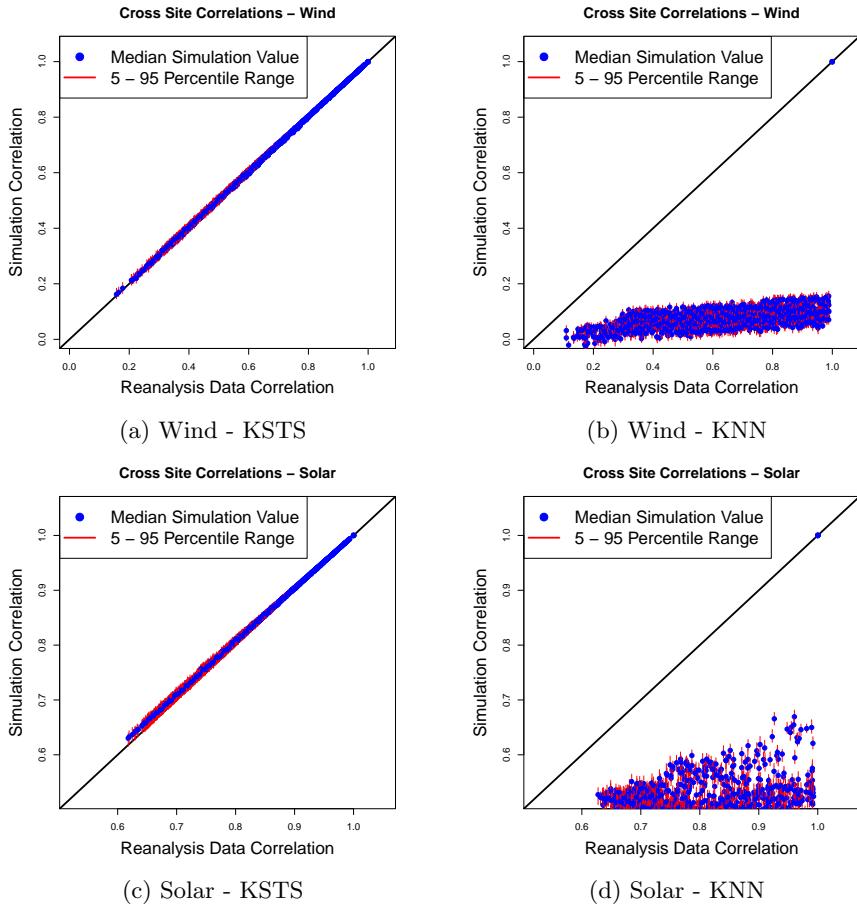


Figure S13: Simulation vs reanalysis data cross site correlation plots for individual fields. (a) Wind KSTS. (b) Wind KNN. (c) Solar KSTS (d) Solar KNN. 40 random grids out of 216 are selected and the 40x40 cross correlation values are computed and plotted instead of the entire 216 x 216 correlation values. The correlations are computed using Pearson's method.

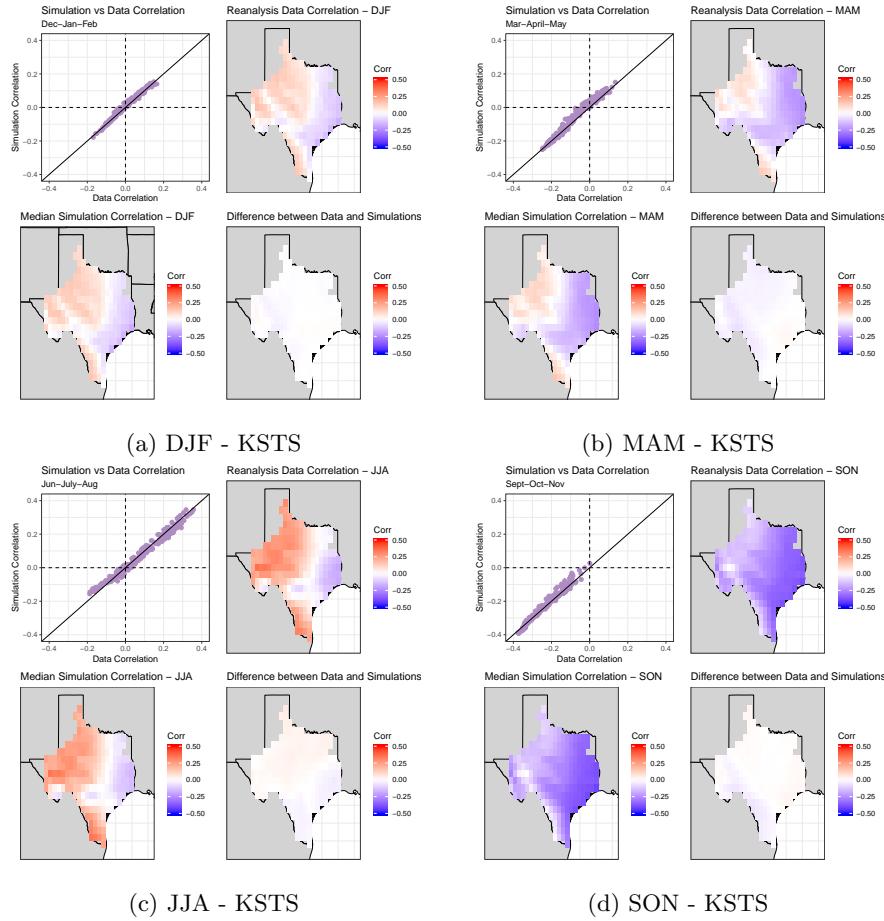


Figure S14: Seasonal correlation between wind and solar at each grid point for KSTS simulations. (a) Dec-Jan-Feb (DJF). (b) Mar-Apr-May (MAM). (c) Jun-Jul-Aug (JJA). (d) Sep-Oct-Nov (SON). For each subplot: (top-left) - median simulation vs reanalysis data correlation between wind and solar. (bottom-left) - Plot of the median simulation correlations. (top-right) - Plot of the reanalysis data correlations. (bottom-right) - Plot of the difference between data and median simulations correlations. The correlations are computed using Pearson's method.

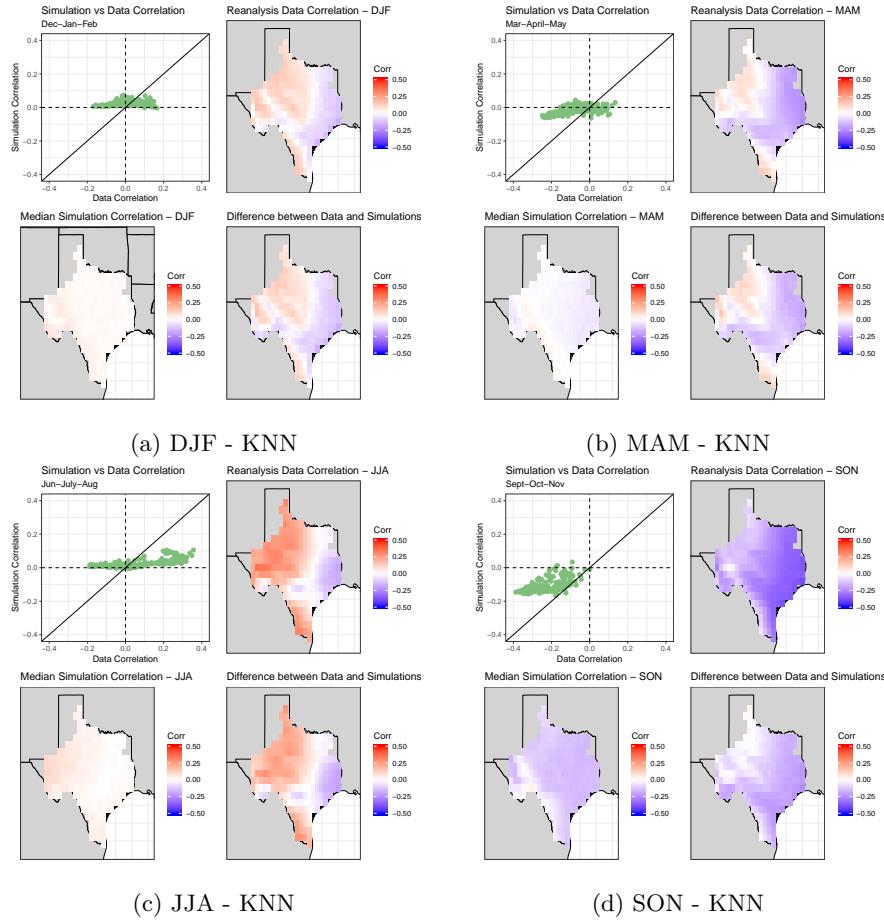


Figure S15: Seasonal correlation between wind and solar at each grid point for KSTS simulations. (a) Dec-Jan-Feb (DJF). (b) Mar-Apr-May (MAM). (c) Jun-Jul-Aug (JJA). (d) Sep-Oct-Nov (SON). For each subplot: (top-left) - median simulation vs reanalysis data correlation between wind and solar. (bottom-left)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

682 Data

683 The Electric Reliability Council of Texas (ERCOT - Fig. S16), functions as
 684 an Independent System Operator and the balancing authority for the Texas
 685 Interconnection and manages about 90% of state's electric load. ERCOT covers
 686 about 75% of the land area in Texas¹.

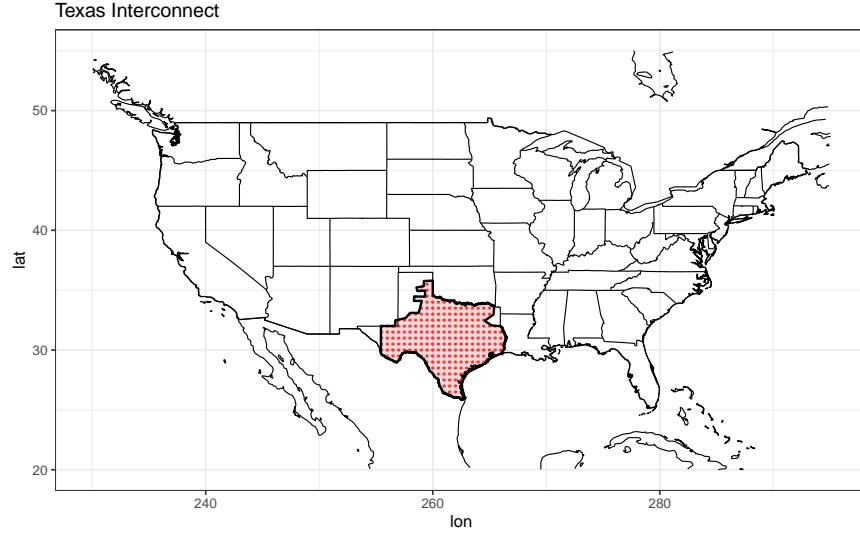


Figure S16: Texas Interconnection / ERCOT domain plot - The red shaded region denotes the area administered. The red dots (216) are the locations of the grid points (0.5° lat \times 0.5° lon) from the ERA-5 reanalysis dataset.

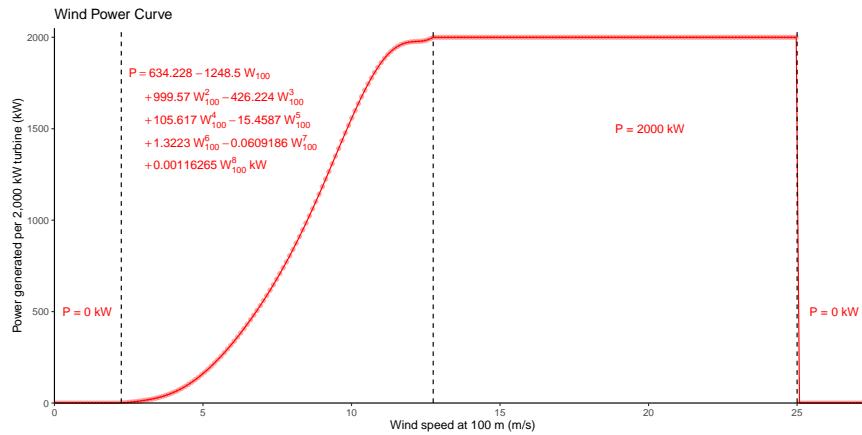


Figure S17: Wind power curve for a V90-2.0MW Vestas turbine.

687 **KSTS Description**

688 The algorithm first models the temporal variability at each site and for each
689 variable field. A state space $D_{i,t}$ is defined through an embedding of the time
690 series with a delay parameter ϕ and an embedding dimension m , where i is
691 the index for the site/variable combination. Following Lall and Sharma ², we
692 consider that the conditional density $f(x_{i,t+1}|D_{i,t})$ is defined through the $x_{i,t'+1}$
693 corresponding to the time indices t' associated with the k-nearest neighbors
694 (knn) of $D_{i,t}$ in the historical data set. The kernel function p_j associated with
695 the jth nearest neighbor is proportional to the rank j of the neighbor ². The
696 sequential drawing from the knn successors at each time step using the specified
697 kernel leads to a simulation of $x_{i,t}$ in a time series dependence structure. Since
698 the procedure leads to a resampling of the historical data, the algorithm can
699 be considered to be a bootstrap which preserves the time dependence in serial
700 data.

701 Given a state space $D_{i,t}$ at site i and time t , the k-nearest neighbor algorithm
702 is used to identify a set of time indices $\tau_{i,t}$ that correspond to the time instances
703 corresponding to the nearest neighbors of site i at time t . An example, for site,
704 with an embedding $D_{i,t}$, defined as (x_t, x_{t-1}) which takes values (10,9.8) at time
705 $t=100$ days, the nearest neighbors would be identified as the closest historical
706 vectors to (10,9.8) using the data for this site.

707 The time instances at which these neighbors occur in the historical time
708 series are then recorded in the order of nearest neighbors in $\tau_{i,t}$. The kernel
709 $p_{i,j}$ then associates a probability proportional to $1/j$ for the jth element (time
710 instance of the jth nearest neighbor of $D_{i,t}$) in $\tau_{i,t}$, for the first k neighbors and 0
711 else. For space-time neighbors across all sites, i.e. to address spatial dependence,
712 we now identify appropriate k-nearest neighbors by finding the time indices in
713 the historical data that have the highest likelihood of being selected across all
714 sites given their associated resampling probabilities.

715 Define $T_{i,t}$ as a matrix such that the rows and columns are pointers for sites
716 and unique time indices from the historical data respectively. The columns
717 record the resampling probabilities associated with the time indices of the k-
718 nearest neighbors for all sites at time t . Other measures of similarity of the
719 spatial neighbors of the temporal process could also be considered. The simi-
720 larity matrix S_t is the sum of the resampling probabilities associated with each
721 unique time index across all sites. The curtailment of the similarity matrix S_t
722 is carried out by selection of the time indices which correspond to the k highest
723 resampling probability values in S_t , now designated as the k-nearest neighbor
724 candidates for the entire spatial field. The full spatial field of the simulation
725 for the next time step is resampled after re-scaling probability values (such that
726 they add to 1) of the curtailed similarity matrix S_t .

727 **Hyper-parameters of the Algorithm**

728 **Resampling Kernel Weight Function**

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

by Lall and Sharma². 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 to Lall and Sharma² 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 ². 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 process³. 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 data².

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 setup⁴. 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 model⁵.

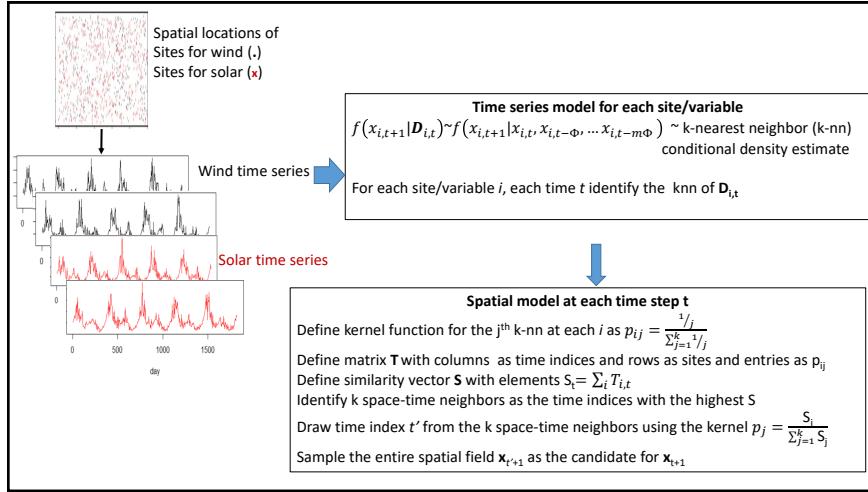


Figure S18: Schematic representation of the k-Nearest Neighbor Simulator (KSTS) with application to wind and solar fields.

766 Supplementary Materials - References

- 767 1. Electric Reliability Council of Texas. About ERCOT, 2021. URL <http://www.ercot.com/about>
- 768 2. Upmanu Lall and Ashish Sharma. A Nearest Neighbor Bootstrap for Re-
769 sampling Hydrologic Time Series. *Water Resources Research* 32.3 (1996):
770 679-693.
- 771 3. Kennel, Matthew B., Reggie Brown, and Henry DI Abarbanel. "Deter-
772 mining embedding dimension for phase-space reconstruction using a geo-
773 metrical construction." *Physical review A* 45.6 (1992): 3403.
- 774 4. Yakowitz, S., and M. Karlsson. "Nearest neighbor methods for time series,
775 with application to rainfall/runoff prediction." *Advances in the statistical
776 sciences: Stochastic hydrology*. Springer, Dordrecht, 1987. 149-160.
- 777 5. Souza Filho, Francisco Assis, and Upmanu Lall. "Seasonal to interan-
778 nual ensemble streamflow forecasts for Ceara, Brazil: Applications of a
779 multivariate, semiparametric algorithm." *Water Resources Research* 39.11
780 (2003).