



<http://www.diva-portal.org>

Postprint

This is the accepted version of a paper published in *Journal of Renewable and Sustainable Energy*. This paper has been peer-reviewed but does not include the final publisher proof-corrections or journal pagination.

Citation for the original published paper (version of record):

Shepero, M., Munkhammar, J., Widén, J. (2019)

A generative hidden Markov model of the clear-sky index

Journal of Renewable and Sustainable Energy, 11: 043703

<https://doi.org/10.1063/1.5110785>

Access to the published version may require subscription.

N.B. When citing this work, cite the original published paper.

Permanent link to this version:

<http://urn.kb.se/resolve?urn=urn:nbn:se:uu:diva-389945>

A generative hidden Markov model of the clear-sky index2 Mahmoud Shepero,^{1, a)} Joakim Munkhammar,^{1, b)} and Joakim Widén^{1, c)}3 *Department of Engineering Sciences, Uppsala University, P.O. Box 534,*
4 *SE-751 21 Uppsala, Sweden*

5 (Dated: 3 July 2019)

Clear-sky index (CSI) generative models are of paramount importance in, e.g., studying the integration of solar power in the electricity grid. Several models have recently been proposed with methodologies that relate to hidden Markov models (HMMs). In this paper, we formally employ HMMs, with Gaussian distributions, to generate CSI time-series. The authors propose two different methodologies. The first is a completely data-driven approach, where an HMM with Gaussian observation distributions is proposed. In the second, the means of these Gaussian observation distributions were pre-defined based on the fraction of time of bright sunshine from the site. Finally, the authors also propose a novel method to improve the autocorrelation function (ACF) of HMMs in general. The two methods were tested on two data sets representing two different climate regions. The performance of the two methodologies varied between the two data sets and among the compared performance metrics. Moreover, both proposed methods underperformed in reproducing the ACF as compared to state-of-the-art models. However, the method proposed to improve the ACF was able to reduce the mean absolute error (MAE) of the ACF by up to 19%. In summary, the proposed models were able to achieve a Kolmogorov-Smirnov (K-S) test score as low as 0.042, and MAE of the ACF as low as 0.012. These results are comparable with the state-of-the-art models. Moreover, the proposed models were fast to train. HMMs are shown to be viable CSI generative models.

6 Keywords: Hidden Markov models, HMM, clear-sky index, CSI, generative clear-sky
7 index model, autocorrelation functiona)Electronic mail: mahmoud.shepero@angstrom.uu.seb)Electronic mail: joakim.munkhammar@angstrom.uu.sec)Electronic mail: joakim.widen@angstrom.uu.se

INTRODUCTION

9 A. Literature overview

10 The variability of solar irradiance on the Earth's surface affects many solar engineer-
11 ing applications such as photovoltaic (PV) power generation¹⁻³. By quantifying the solar
12 irradiance variability, the design and operation of power systems with large amounts of
13 grid-connected distributed PV power generation can be improved^{2,4}. However, quantify-
14 ing and reproducing the complexity of solar irradiance variability over time is a challenge
15 for statistical and machine learning models⁵. The value of high resolution irradiance data
16 might increase in the future^{4,6}. Consequently, several recent models have been proposed to
17 generate high resolution irradiance data using low resolution data⁷⁻¹⁰.

18 The solar irradiance is often normalized using the clear-sky irradiance, producing the
19 clear-sky index (CSI). The CSI has interesting features of statistical complexity, particularly
20 on minute to instantaneous scale^{11,12}. High resolution is favorable in order to take short
21 bursts of overirradiance into account¹³, but resolutions above 1 second are arguably less
22 informative and only increases the complexity of data management¹⁴.

23 There exist irradiation estimates from ground measurements or from satellite data, see,
24 e.g., Ref. 15, and from software, such as Meteonorm¹⁶. Despite this, more high resolution
25 data are needed for many locations¹¹, and when data are insufficient, realistic synthetic data
26 generated by generative models are useful⁵. Also, when modeling CSI variability, it is of
27 importance to quantify the model output accuracy so that it generates realistic data^{5,17-20}.
28 In terms of modeling, the distribution of minute to instantaneous CSI can be modeled by
29 two^{20,21}, three^{22,23}, n -single peak distributions⁵, or n -multi peak distributions²⁴.

30 In order to ensure that the model generates CSI samples that resemble that of the training
31 CSI data, models are evaluated based on firstly distribution goodness-of-fit, and secondly
32 autocorrelation function (ACF) likeness. Model evaluation and selection, in terms of out-
33 put probability distribution goodness-of-fit, are typically measured with, e.g., Kolmogorov-
34 Smirnov (K-S) statistics^{12,17,25-27}. ACF likeness compares the temporal variability between
35 models and data, and the ACF measured over a set of lags has commonly been used^{17,19,20}.
36 A metric for model selection is for example a mean absolute error (MAE) over lags in the
37 ACF⁵.

Models aiming to quantify the CSI and generate realistic synthetic CSI data, with proper ACF likeness, include Gaussian-Markov¹⁹, auto-regressive Gaussian²⁸, neural networks²⁹, copula modeling^{17,26}, fractal cloud modeling³⁰, Markov-chains^{5,12,18,31-34} and Dirichlet process Gaussian mixture (DPGM)²⁴. These CSI generators are practical, since they utilize some existing data set of lower resolution or averaged CSI data to estimate higher resolution data (temporal or spatial), see, e.g., Refs. 12, 25, 31, 34-36, or smaller amounts of CSI data to generate unlimited amounts of data^{5,17,18,26}.

In particular Markov-chains, or generally Markov models, have been useful as CSI generators, where Ref. 34 generated minute resolution, while in Ref. 32 daily resolution was modeled. In Ref. 18 a two-state Markov-chain mixture probability distribution was used, similar to the one in the study of Morf³¹, where a general model for generating clear and cloudy periods was constructed. As a generalization of the two-state models an n -state piecewise uniform Markov-chain mixture distribution model was developed in Refs. 5 and 10, which had low model complexity, yet high accuracy. Generally these models have presented high accuracy in both distribution goodness-of-fit and ACF likeness.

A novel approach was recently proposed by Frimane *et al.*²⁴ where a DPGM model was employed to generate CSI samples. Unlike previous models, the DPGM model is a nonparametric model that selects the number of CSI clusters based on the training data. The daily distribution of CSI is considered to be represented by an infinite mixture of Gaussian distributions. A Dirichlet process was used to cluster the daily distribution of CSI into a number of clusters, with some probability of creating a new cluster for the observed daily CSI. Thus the number of clusters was inferred from the data. A Markov chain was then learned from the data, which was used to generate the CSI time-series.

One conclusion from the literature is that mixture models based on stochastic processes (Markov-chain or otherwise) reproduce the CSI variability comparatively well, and that models which can utilize meteorological observables and produce realistic output with high accuracy are particularly useful.

An improved understanding regarding three-state modeling connected to meteorological variables²², in particular in combination with hidden Markov models (HMMs), have been proposed but not yet investigated in the literature¹⁸. Also, it can be concluded that HMMs with Gaussian distributions have not been investigated in the literature either. In this paper both a three-state Gaussian HMM model, with connections to measured fraction of time of

71 light sunshine, and a general machine learning-based HMM for n states are developed and
72 investigated.

73 The developed models are beneficial to generate CSI time-series in locations of interest.
74 The generated time-series can be then converted to the global irradiance, which can be used
75 in applications such as PV potential studies.

75 B. Contribution

76 This study aims to develop an HMM for the CSI based on n Gaussian distributions. The
77 model extends previous Markov-chain mixture distribution models^{5,18} and previous literature
78 through

- 79 1. Developing an HMM generative model for the CSI. Previous contributions did not
80 explicitly employ HMMs.
- 81 2. Connecting the HMM model to previous meteorological CSI models.
- 82 3. Proposing a methodology to improve the ACF of the fitted HMMs; to the best of the
83 authors knowledge this method has not been proposed before in the literature.

84 The model is trained and tested on solar irradiance data sets for Hawaii and Norrköping,
85 which were also used in Refs. 5, 17, and 18.

86 This paper is organized as follows. The proposed HMMs are described in Section II. In
87 Section III, the results obtained from implementing the model are presented. A discussion
88 is provided in Section IV. Finally, conclusions are drawn in Section V.

89 II. METHODS

90 A general introduction to HMMs is provided in Section II A. In Section II B, the HMM
91 is applied to CSI modeling. This section connects HMMs to CSI observations in a certain
92 location. Section II C presents the two proposed methods to employ HMMs in CSI time-
93 series generation. Performance metrics, used in comparing the various models, are described
94 in Section II D. Finally, data used in learning and comparing the models are introduced in
95 Section II E.

Hidden Markov model

97 Here, a brief introduction to HMMs is provided. For more detailed information regarding
98 HMMs Refs. 37 and 38 can be consulted.

99 A discrete-time HMM is a state space model that is characterized by a set of observa-
100 tions $\{X_t\}_{t=1}^T$. These observations were observed at discrete time-steps $t \in \{1, 2, \dots, T\}$.
101 Furthermore, the observations are dependent on a set of hidden states $\{S_t\}_{t=1}^T$ such that the
102 observations are conditionally independent given the states, see Fig. 1.

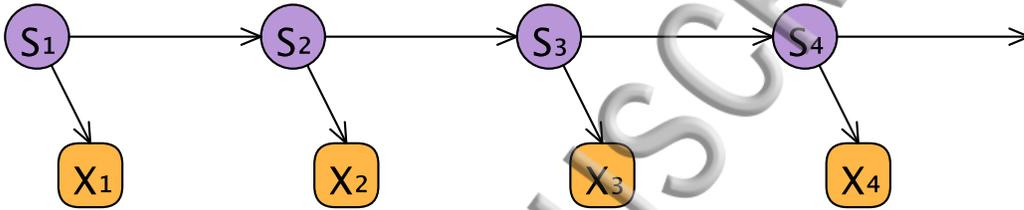


FIG. 1. A graphical representation of the HMM for four time-steps. The arrows represent conditional dependencies. The hidden states $\{S_1, S_2, \dots\}$ are represented by circles and the observations $\{X_1, X_2, \dots\}$ by squares.

The transition probabilities of the hidden states follow a Markov chain. A Markov chain is a memoryless process such that the hidden state S_t is only dependent on the previous hidden state S_{t-1} and independent of the previous state trajectory, i.e.,

$$\begin{aligned} P(S_t = j | S_1, \dots, S_{t-1} = i) &= P(S_t = j | S_{t-1} = i) \\ &= p_{ij}, \end{aligned} \quad (1)$$

103 where p_{ij} is the probability of transitioning from state i to state j . A transition matrix A
104 between the states can be formed:

$$A = \begin{bmatrix} p_{11} & \cdots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \cdots & p_{nn} \end{bmatrix}, \quad (2)$$

106 where n is the number of hidden states, i.e., the cardinality of $\{S_t\}_{t=1}^T$. The transition matrix
107 A has to satisfy

$$\sum_{j=1}^n p_{ij} = 1 \quad (3)$$

for every initial state i .

The observations X_t are only dependent on the hidden state at time t

$$\begin{aligned} P(X_t = x_t | S_1, \dots, S_t = i, X_1, \dots, X_{t-1}) = \\ P(X_t = x_t | S_t = i) = b_i(x_t), \end{aligned} \quad (4)$$

where $b_i(x_t)$ is the observation distribution of state i , e.g., a Gaussian distribution.

The HMM is characterized by the model parameters θ . The model parameters θ consist of the initial distribution of the hidden states, the transition matrix A , and the parameters of the observation distributions $b_i(x_t)$. The parameters $b_i(x_t)$ can be, for example, the mean μ_i and variance σ_i^2 for a univariate Gaussian observations.

Assuming that the hidden process is aperiodic and irreducible and assuming that the process is stationary, the initial distribution can be replaced by the stationary distribution $\boldsymbol{\pi}$ of the Markov chain, see, e.g., Ref. 39, p. 66. In other words, if the process has been going on for a long time, e.g., the CSI time-series started long ago, the initial distribution loses importance and the stationary distribution becomes more important. The stationary distribution $\boldsymbol{\pi}$ can be estimated by solving

$$\boldsymbol{\pi} A = \boldsymbol{\pi}. \quad (5)$$

In HMMs the model parameters θ are learned from sequences of observations using the Baum-Welch algorithm. The Baum-Welch algorithm is of computational complexity $\mathcal{O}(n^2T)$, see Ref. 37, where n is the number of hidden states and T is the length of the training time-series.

B. Clear-sky index as HMM observations

This section draws connections between HMMs and CSI observations in a certain location. In other words, this section shows that the CSI time-series can be viewed as the observation of an HMM.

Assume that the hidden states S_t of the HMM represent the status of the cloud cover in the sky at every time-step t in a certain location. Following Hollands and Suehrcke²², the cloud cover can be categorized into three categories, $S \in \{1, 2, 3\}$. These categories can be conceptualized as states of the atmosphere²². For clarity, we propose, here, to call these

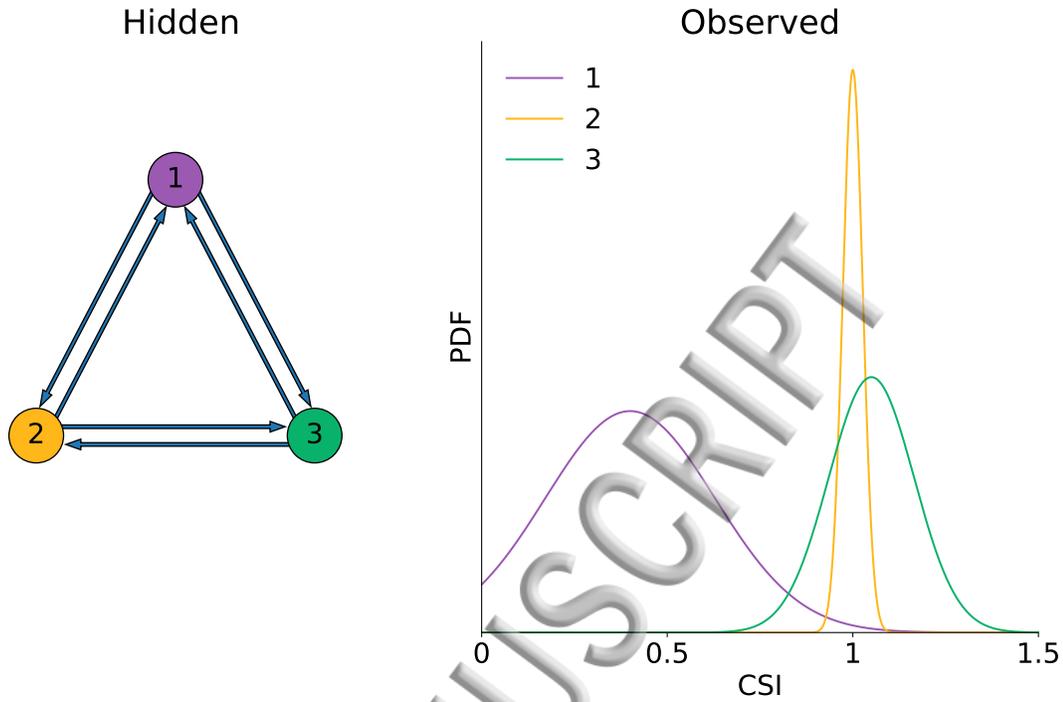


FIG. 2. A diagram connecting the hidden states of a three state HMM—representing the cloud cover states—and the CSI. The hidden states are represented by the circles and the between-state transitions are represented by arrows. The remain-in-state transitions are not presented. The observed distribution of every state is plotted on the right. State 1 represents the obscured state, state 2 represents the unobscured in clear-sky state, and state 3 represents the unobscured in partially cloudy sky.

134 hidden states obscured (1), unobscured in clear sky (2), and unobscured in partially cloudy
 135 sky (3). Consequently, the transition matrix A is a (3×3) -matrix.

136 As shown in Fig. 2, each hidden state is coupled with an observation distribution, defined
 137 by a probability density function (PDF). In such case, the CSI measurements in a certain
 138 location can be conceptualized as observations from the hidden—or unobserved—states in
 139 the HMM model. These hidden states represent the cloud cover in the sky.

140 Pursuing the same derivation as in Ref. 18, few model parameters can be extracted from
 141 meteorological measurements. The probability of having an obscured sun for δ_1 consecutive
 142 time steps—then switching to a different state—is⁴⁰ (p. 308)

$$P(\delta_1) = (1 - p_{11})p_{11}^{\delta_1 - 1}. \quad (6)$$

The expected value of obscured durations $\mathbb{E}[\delta_1]$ can be estimated as ^{18,40}

$$\mathbb{E}[\delta_1] = \frac{1}{1 - p_{11}}. \quad (7)$$

Applying Eq. 7 for all the three hidden states results in

$$\begin{aligned} p_{11} &= \frac{\mathbb{E}[\delta_1] - 1}{\mathbb{E}[\delta_1]}, \\ p_{22} &= \frac{\mathbb{E}[\delta_2] - 1}{\mathbb{E}[\delta_2]}, \\ p_{33} &= \frac{\mathbb{E}[\delta_3] - 1}{\mathbb{E}[\delta_3]}, \end{aligned} \quad (8)$$

144 where $\mathbb{E}[\delta_2]$ and $\mathbb{E}[\delta_3]$ are the mean durations of unobscured in clear-sky and unobscured in
 145 partially cloudy sky, respectively. These can be directly observed or calculated from other
 146 measured values in a certain location, see, e.g., Refs. 18, 22, 31, 41, and 42.

147 Hollands and Suehrcke²² proposed Gaussian distributions for each hidden state. The
 148 means of the Gaussian distributions were proposed to be $\mu_2 = 1$, $\mu_3 = 1.04$, and

$$149 \quad \mu_1 = 0.1205 + 0.3341\mathcal{K}, \quad (9)$$

150 where \mathcal{K} is the ratio between the mean CSI in the location and the CSI at the peak of the
 151 distribution. Here and as proposed in Refs. 22 and 41, we calculate $\mathcal{K} = \sqrt{\tau}$, where τ is the
 152 fraction of time of bright sunshine. Here, τ is defined as the fraction of time where the CSI
 153 is larger than 0.95, see Ref. 18.

154 The HMM has further analytic properties as regards the temporal variability, e.g., the
 155 ACF $\rho(k)$ for the k -th time-lag of the observed CSI time-series X_t is defined as⁴³:

$$156 \quad \rho(k) = \frac{\mathbb{E}[(X_t - \mu)(X_{t+k} - \mu)]}{\sigma^2} \quad (10)$$

where μ and σ^2 represent the mean and variance of the time-series. In HMMs, Eq. 10 can
 be shown to follow⁴⁰ (p. 310)

$$\begin{aligned} \rho(k) &= \frac{\sum_i \pi_i \mu_i \sum_j \mu_j (A^k)_{ij} - \mu^2}{\sigma^2}, \\ &= \frac{\boldsymbol{\pi} \text{diag}(\boldsymbol{\mu}) A^k \boldsymbol{\mu}' - \mu^2}{\sigma^2}, \end{aligned} \quad (11)$$

157 where $(A^k)_{ij}$ is the value in i th row and j th column of the A^k matrix, $\boldsymbol{\mu}$ is a vector containing
 158 the means of the observation distributions of the hidden states; and $\boldsymbol{\mu}$ and σ^2 are the mean

159 and variance of the mixture distribution of the observations, respectively. In case of Gaussian
 160 mixture distributions⁴⁰ (p. 11),

$$161 \quad \mu = \sum_{i=1}^n \pi_i \mu_i, \quad (12)$$

162 and

$$163 \quad \sigma^2 = \sum_{i=1}^n \pi_i (\mu_i^2 + \sigma_i^2) - \mu^2. \quad (13)$$

164 Equation 11 is applicable only for $k > 0$ ³⁹ (p. 70). For $k = 0$, $\rho(0) = 1$ by definition⁴⁴.
 165 One limitation of the HMMs is that they can only produce exponentially decaying ACFs⁴⁵,
 166 observe that $\rho(k) \propto A^k$ for all $k > 0$ in Eq. 11.

167 An HMM with n hidden states and univariate Gaussian observations, can be defined by
 168 $(n^2 + 2n)$ parameters; n^2 parameters in A , n parameters in $\boldsymbol{\mu}$, and n parameters in $\boldsymbol{\sigma}^2$ which
 169 is the vector containing the variances of the n observation distributions.

170 Fitting a Gaussian mixture model to the CSI distribution provides an estimate of $\boldsymbol{\mu}$,
 171 $\boldsymbol{\sigma}^2$, and $\boldsymbol{\pi}$, see, e.g., Refs. 22 and 23. The stationary distribution of the Markov process
 172 $\boldsymbol{\pi}$ controls the contribution of each Gaussian observation distribution to the CSI mixture
 173 distribution.

174 In order to estimate the transition matrix A , Eqs. 3, 5, 8 and 11 can, theoretically,
 175 be employed. However, Eq. 11 is nonlinear, because of the term A^k . As a result, the
 176 complexity grows as the number of hidden states grows. For example, for an HMM with
 177 $n = 2$, Munkhammar and Widén¹⁸ showed that solving Eqs. 3 and 8 was sufficient to
 178 estimate A .

179 For an HMM with $n = 3$, Eqs. 3, 5 and 8 provide 8 independent equations as the 9th
 180 equation can be derived using the remaining 8 equations. Equation 11 can provide the 9th
 181 independent equation, however, only if $k > 1$ since the equation for $\rho(1)$ can be derived from
 182 the previously mentioned 8 equations. Consequently, a nonlinear equation solving method
 183 is needed if this approach is to be pursued.

184 C. Model implementation

185 We assume that the hidden states of an HMM represent some hidden, or unobserved,
 186 cloud cover states. These cloud cover states in turn control the CSI observations. Using the
 187 CSI observations, an HMM with various number of states can be fitted using the Baum-Welch

algorithm. These models will be called the Gaussian hidden Markov models (GHMMs) in
189 further discussions, since the observation distributions for each hidden state are assumed
190 to be Gaussian. The Gaussian distribution was chosen since it was shown to adequately
191 represent the CSI time-series in Refs. 22 and 24. Selecting the optimal number of hidden
192 states n will be based on comparing the fitted models using some performance metrics.

The GHMM with $n = 2$ shares some similarities with the models proposed in Refs.
193 18, 31, and 46. Here, unlike Ref. 18, the observation distributions are Gaussian, and they
194 are not truncated at a middle point. In comparison, Munkhammar and Widén¹⁸ used
195 Gaussian, log-normal and polynomial distributions in their two-state model.
196

It is important to note that in all HMMs trained here the Baum-Welch algorithm was
197 employed to estimate the model parameters θ . The Baum-Welch algorithm updates the
198 model parameters θ in each iteration such that θ increases the log-likelihood of the model.
199 The final model is thus the model that locally maximizes the log-likelihood. As a result,
200 for the GHMM with $n = 3$, the estimated means of the observation distributions of the
201 hidden states might not represent the cloud cover assumptions proposed by Hollands and
202 Suehrcke²².
203

An alternative approach to the GHMM with $n = 3$ is to fix the means of the observa-
204 tion distributions, and thereby connect it to the cloud cover assumptions of Hollands and
205 Suehrcke²². In this alternative model, the Baum-Welch algorithm was employed to esti-
206 mate all the model parameters θ except the means of the observation distributions μ . The
207 first mean μ_1 was pre-defined according to Eq. 9, and μ_2 and μ_3 were set to 1.0 and 1.04,
208 respectively²². This alternative modeling approach will be called the fixed-means Gaussian
209 hidden Markov model (FMGHMM) in further discussions.
210

In similarity with Ref. 5 this paper uses HMMs with $n > 3$, but unlike Ref. 5 Gaussian,
211 instead of uniform, observation distributions are employed. A non-parametric model for the
212 number of hidden states was recently developed in Ref. 24. In this model, the number of
213 hidden states was estimated from the training data. Thus, it can be understood as a GHMM
214 with an infinite number of hidden states⁴⁷. Nevertheless, the authors used mixed Gaussian
215 observation distributions.
216

To summarize this section, two modeling strategies were proposed: the GHMMs and the
217 FMGHMMs. Using the GHMMs, the number of hidden states and the shapes of the obser-
218 vation distributions were learned completely from the data. Consequently, the hidden states
219

of the GHMMs are not guaranteed to reflect the assumptions of Hollands and Suehrcke²².
On the other hand, the FMGHMM, used the assumptions in Ref. 22 to fix the means of
the observation distributions. This makes the FMGHMM more connected to the theoretical
properties of the CSI. The FMGHMM is, thus, only meaningful if the number of hidden
states $n = 3$.

In this paper, the Python HMM package `hmmlearn v0.2.1`⁴⁸ was used to fit the proposed
models.

D. Model evaluation metrics

Theis, Oord, and Bethge⁴⁹ compared the various metrics used in evaluating generative
models, e.g., HMMs. The authors concluded that generative models need to be evaluated
with respect to their applications.

Consequently, we propose the following performance metrics to compare the proposed
models:

1. Log-likelihood \mathcal{L} of the training and test data sets.
2. K-S test score K between both the training and test data sets and samples from the
model.
3. MAE of the ACF ϵ between the test data and samples from the model.
4. Training time.

The log-likelihood, $\mathcal{L} = \log(P(X_1, \dots, X_T|\theta))$, represents the natural logarithm of the prob-
ability of observing the CSI time-series given the model parameters θ . This metric was
measured for both the training and the test data sets.

Following the conclusions of Theis, Oord, and Bethge⁴⁹ and to evaluate the models with
regard to their applications, two metrics were added to measure firstly the goodness-of-fit
of the model distribution and the distribution of the CSI time-series; and secondly the ACF
likeness between the generated samples and the ACF of the CSI time-series. As stated
before, these two metrics were commonly used in the literature. The K-S test score was
employed before in Refs. 5, 12, 17, 18, 25–27. The MAE of the ACF was employed before in
Refs. 5 and 18. Consequently, the chosen metrics enable comparing the proposed methods to

249 as regards different error metrics.

250 To evaluate the goodness-of-fit of the distribution, the K-S test was used. The K-S test
251 score K evaluates the similarity between two distributions

$$252 \quad K = \max |F_1(x) - F_2(x)|, \quad (14)$$

253 where $F_1(x)$ and $F_2(x)$ are the empirical cumulative distribution functions of the two com-
254 pared distributions⁵⁰.

255 To evaluate the ACF likeness, the MAE of the ACF ϵ was evaluated, which is estimated
256 as

$$257 \quad \epsilon = 1/60 \sum_{k=1}^{60} |\hat{\rho}(k) - \rho(k)|, \quad (15)$$

258 where $\hat{\rho}(k)$ is ACF of the test data, and $\rho(k)$ is the ACF of the samples from the model.
259 The ACF error was calculated for one hour, i.e., 60 time-lags. This makes our ϵ results
260 comparable with that of Refs. 5 and 18.

261 The reason for only recording the MAE of the ACF ϵ on the test data is that by visual
262 inspection, as will be shown later in Figs. 5, 8 and 11, there was no difference between the
263 ACF of the training and test data for both locations.

264 Finally, the training time was measured in seconds.

265 E. Data

266 The data used to train and evaluate the models were based on radiometer measurements
267 of global horizontal irradiance (GHI) for one year obtained from the Swedish Meteorological
268 and Hydrological Institute (SMHI) for Norrköping, Sweden (59°35'31" N 17°11'8" E)⁵¹, and
269 from the National Renewable Energy Laboratory (NREL) radiometer array in Oahu, Hawaii,
270 USA (21°31' N, 158°09' W)⁵². The data were recorded during the years 2008 and 2010 for
271 Norrköping and Hawaii, respectively. The Kipp & Zonen CM 21⁵³ and LI-COR LI-200⁵⁴
272 pyranometers were used in Norrköping and Hawaii, respectively.

273 Data from both locations represent instantaneous irradiance with one minute resolution,
274 averaged from raw data recorded with higher resolutions. From these data sets 120 data
275 points (minutes) per day (centered around noon each day) for both locations were used in
276 order to avoid low solar angles. Thus, the selected data sets totaled 43,800 (120×365) data

278 points from each location. These particular data sets were previously selected for compar-
279 ative reasons since they were also used in Refs. 17 and 26 to generate an autocorrelation
280 model of the CSI using copulas, and in Refs. 5 and 18 for the Markov-chain mixture dis-
281 tribution models. Testing the performance of the proposed methods on data from different
282 locations, and representing various solar altitude angles, climate conditions, air masses, etc.,
is encouraged for future works.

283 The CSI normalization for the Norrköping data was made with the Ineichen-Perez clear-
284 sky irradiance model⁵⁵ and for the Hawaii data the McClear clear-sky irradiance model⁵⁶ was
285 used. The use of these clear-sky irradiance models for obtaining the CSI for both data sets is
286 described, tested, and was found optimal in Refs. 17 and 26. Thus, these particular clear-sky
287 irradiance models were chosen in order to enable comparisons with previous papers that used
288 the same data sets. Moreover, the two clear-sky irradiance models had similar performances
289 on the data sets, which further motivated using both models to enable comparison of results
290 with previous studies. For a detailed review on the recommended clear-sky irradiance models
291 for each climate zone, Ref. 57 can be consulted.

292 The CSI data from both locations were divided into training data and test data. The
293 division was made on every-other-day basis in order to minimize the seasonality effect on
294 training of the models. In other words, the training data were the data recorded on the days
295 1,3,...,365; and the test data were the data from the days 2,4,...,364.

296 For each scenario, the model was used to generate a synthetic time-series with length
297 1 million, i.e., $T = 10^6$. This sample is large enough to ensure stable estimates of the
298 performance metrics of the models.

299 III. RESULTS

300 This section provides the results of fitting the proposed HMMs to the CSI data. Sec-
301 tion III A provides the results of the GHMMs. In Section III B, the results of the FMGHMM
302 are presented. Section III C improves the ACF of the fitted HMMs, with the FMGHMM
303 taken as an example. In Section III D our results are compared to the results of some
304 previously proposed models.

Gaussian hidden Markov models

306 In this section, the performance results of the GHMMs are presented. The Baum-Welch
307 algorithm was employed to learn all the model parameters θ including the means of the
308 observation distributions μ —unlike the FMGHMMs results to be presented in Section III B.

309 As shown in Figs. 3a and 3b, the log-likelihood \mathcal{L} of the training and the test data
310 sets increased as the number of hidden states increased. The rate of increase in the log-
311 likelihood, however, decreased as the number of hidden states increased. This indicated
312 that the improvement increments in model fitting to the data was high for the models with
313 a lower number of hidden states.

314 The K-S test scores K from the models were lower than 0.1 for both locations and for
315 both the training and the test data sets, see Figs. 3c and 3d. The K-S test score was 0.050
316 for both the training and the test data sets for $n = 3$ in the case of Hawaii. In the case of
317 Norrköping, the scores were 0.047 and 0.048 for the training and test data sets, respectively.
318 The minimum K-S test score on the test data set was 0.008 at $n = 9$, for the Hawaii case,
319 and 0.012 at $n = 9$, for the Norrköping case.

320 Figures 4a and 4b compare the PDFs of the samples from the fitted GHMMs to the test
321 data. The Python class `scipy.stats.gaussian_kde` was used, with the default parameters,
322 to smooth the histograms and perform a kernel density estimation (KDE). Still, the exact
323 histograms for both locations are presented in Figs. 4c and 4d.

324 Figures 4a and 4c show that for Hawaii the model with two hidden states, $n = 2$, under-
325 represented the CSI peak at $\text{CSI} \approx 1$. In addition, the first mode—peak—of the bimodal
326 CSI distribution was situated at a higher CSI compared with the test data, 0.45 and 0.333,
327 respectively. For $n = 3$, the model results seem to match the cloud cover assumptions of
328 Hollands and Suehrcke²². Moreover, the PDF of the model samples decently represented the
329 second mode of the test data. Still the model estimated the first mode of the distribution
330 at a higher CSI; 0.433 compared with 0.333 for the test data. The exact closeness scores
331 between the distributions were presented in Figs. 3c and 3d.

332 For the Norrköping data, the CSI distribution was more skewed towards the obscured
333 cloud cover states, see Figs. 4b and 4d. The first mode value of the PDF was at $\text{CSI} = 0.19$
334 for the test data. The GHMM with $n = 2$ under-represented the mode at $\text{CSI} \approx 1$. In
335 addition, the first mode of the PDF was estimated at $\text{CSI} = 0.33$. For $n = 3$, the model

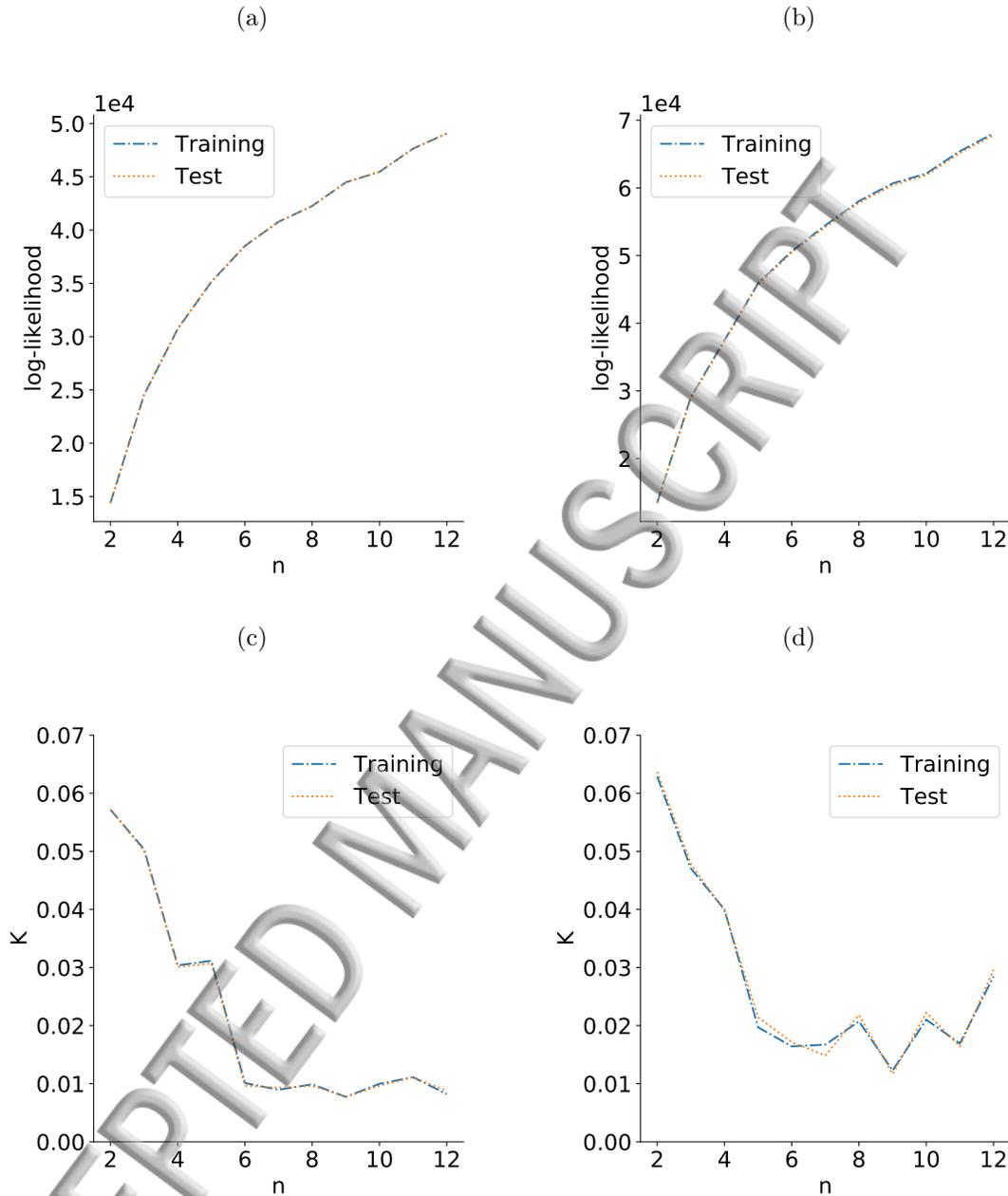


FIG. 3. The log-likelihood \mathcal{L} , in (a) and (b), and the K-S test scores K , in (c) and (d), results for the GHMMs with different number of hidden states n . (a) and (c) present the results of Hawaii. (b) and (d) present the results of Norrköping.

336 allocated two distributions for the first half of the PDF. This improved the fit to the obscured
 337 part of the PDF, but still under-represented the mode at $\text{CSI} \approx 1$. The model with $n = 5$
 338 adequately represented the test data. This model allocated two hidden states each to the
 339 two unobscured-sun states, and three hidden states to represent the obscured-sun state.

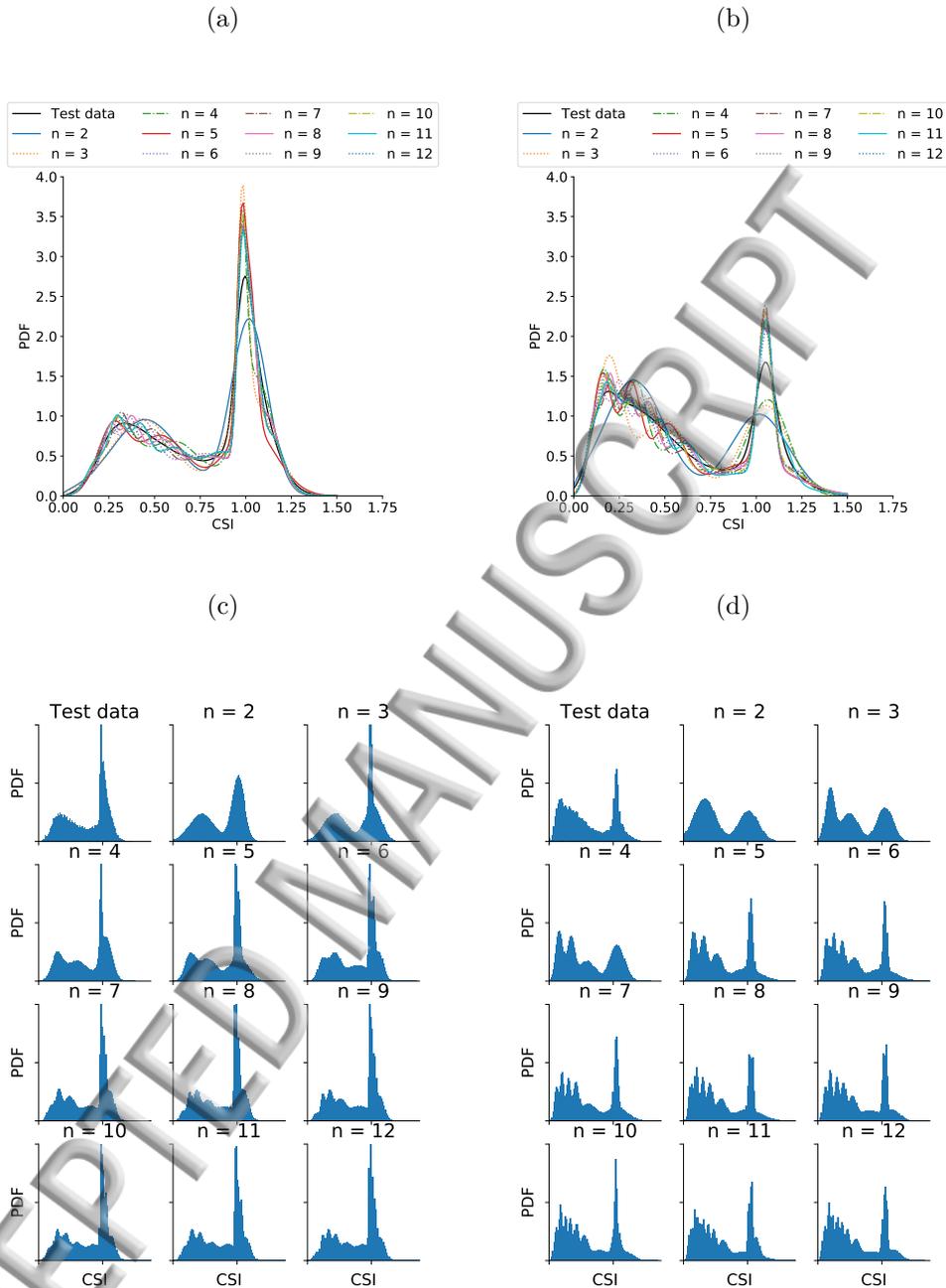


FIG. 4. Kernel density estimation (KDE), in (a) and (b), and Histograms in (c) and (d), of the CSI of the test data and samples from the GHMMs with different number of hidden states n . (a) and (c) present the results of Hawaii. (b) and (d) present the results of Norrköping.

TABLE I. The observation distributions for the GHMM model with $n = 3$ for both locations. All distributions are Gaussians.

| State | Hawaii | | Norrköping | |
|-------|---------|--------------|------------|--------------|
| | μ_i | σ_i^2 | μ_i | σ_i^2 |
| 1 | 0.433 | 0.0312 | 0.190 | 0.0048 |
| 2 | 0.982 | 0.00032 | 0.490 | 0.019 |
| 3 | 1.025 | 0.016 | 1.05 | 0.0154 |

340 In order to facilitate further comparisons with the FMGHMM, the results of the GHMM
 341 with $n = 3$ are summarized in Tables I and II. Table I presents the fitted observation
 342 distributions of the GHMM with three hidden states for both locations. For the Hawaii
 343 data set, the distributions resembled the three clusters from Ref. 22, presented here in
 344 Section II B. In Norrköping, however, the HMM combined the two unobscured-sun states—
 345 from Section II B—into one state, state 3 in Table I. The obscured-sun state was represented
 346 by the remaining two distributions. Table II presents the performance metrics of the GHMM
 347 with $n = 3$.

348 As regards the ACF, the GHMM did not accurately represent the ACF of neither the
 349 training nor the test data—both had similar ACFs—see Fig. 5. For Hawaii, the ACF from
 350 model samples did not improve as the number of hidden states increased. For Norrköping,
 351 the ACF of samples from the models with $n > 3$ closely resembled that of the training and
 352 test data.

353 The MAE of the ACF ϵ is presented in Fig. 6a. For Hawaii, ϵ stagnated at approximately
 354 0.2 or 20%. This was perhaps expected given that Rydén, Teräsvirta, and Åsbrink⁴⁵ also
 355 noticed that their HMM failed to properly represent the ACF of the data. The authors
 356 mentioned that models fitted with maximizing the likelihood cannot accurately capture the
 357 profile of the ACF. A method is proposed and tested in Section III C to improve the ACF
 358 of HMMs.

359 The training time of different models is presented in Fig. 6b. As expected from the
 360 Baum-Welch algorithm the training time increased with the square of the number of states.

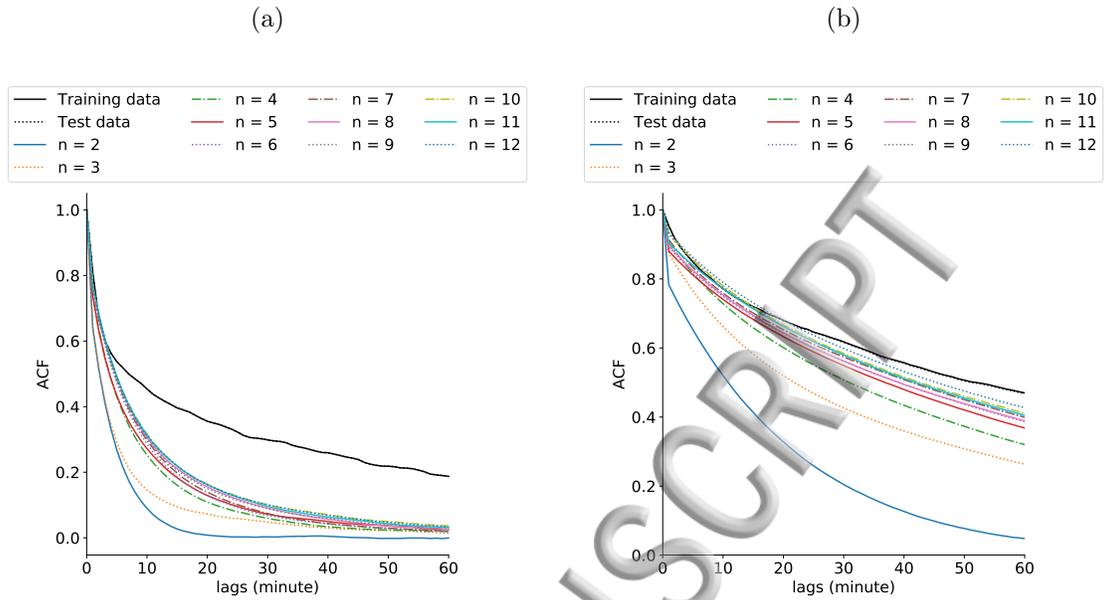


FIG. 5. The ACF of the training and test data and samples from the GHMMs with different number of hidden states n . (a) Hawaii, and (b) Norrköping. In both locations the training and test data had the same ACF.

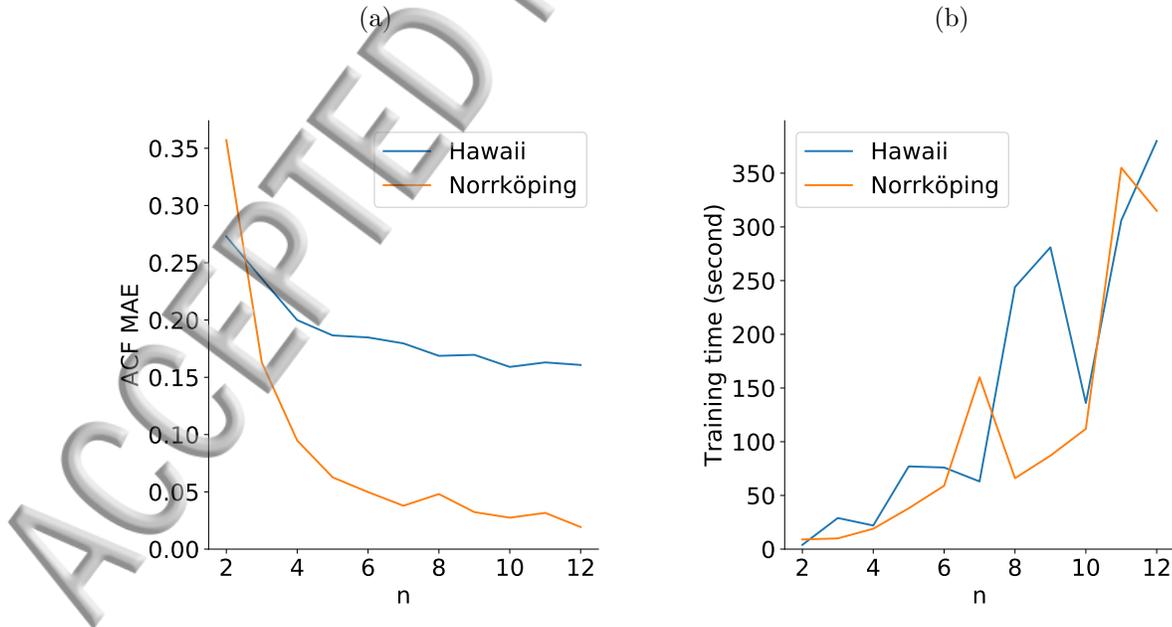


FIG. 6. (a) The MAE of the ACF, ϵ , see Section II D, and (b) the training time for the GHMMs with different number of hidden states n .

TABLE II. The performance metrics for the GHMM for both locations and data sets for $n = 3$. Note that the MAE of the ACF ϵ was only estimated on the test data set. On the other hand, the log-likelihood \mathcal{L} and the K-S test score K were estimated for both data sets.

| Location | Training | | Test | | |
|------------|---------------|-------|---------------|-------|------------|
| | \mathcal{L} | K | \mathcal{L} | K | ϵ |
| Hawaii | 24534.25 | 0.050 | 24557.19 | 0.050 | 0.24 |
| Norrköping | 29058.64 | 0.047 | 28893.89 | 0.048 | 0.16 |

TABLE III. The performance metrics for the FMGHMM for both locations and data sets. Note that the MAE of the ACF ϵ was only estimated on the test data set. On the other hand, the log-likelihood \mathcal{L} and the K-S test score K were estimated for both data sets.

| Location | Training | | Test | | |
|------------|---------------|-------|---------------|-------|------------|
| | \mathcal{L} | K | \mathcal{L} | K | ϵ |
| Hawaii | 18883.97 | 0.118 | 18897.70 | 0.118 | 0.22 |
| Norrköping | 24275.97 | 0.047 | 24222.60 | 0.045 | 0.22 |

361 B. Fixed-means Gaussian hidden Markov model

362 In this section, the results of applying the Baum-Welch algorithm to learn the FMGHMM
 363 parameters θ are presented. The model parameters here do not include the means of the
 364 observation distributions μ , since they were pre-defined based on the model proposed in
 365 Ref. 22 as described in Section II C.

366 Table III presents the results of the FMGHMM as regards the performance metrics. The
 367 K-S test score K was 0.118 for Hawaii for both the training and test data sets, respectively.
 368 For Norrköping, the K scores were 0.047 and 0.045 respectively for the training and test
 369 data sets. As regards the MAE of the ACF ϵ , it was high, 0.22 or 22%, for both locations.
 370 The training time was 6 and 5 seconds for Hawaii and Norrköping, respectively.

371 As regards the observed distribution, Fig. 7 presents a comparison between the test
 372 data and samples from the FMGHMM. The generated PDF seems—by visual inspection—
 373 to adequately represent the PDF of the CSI. Nonetheless, the accuracy of the model was
 374 higher in case of Norrköping in comparison with Hawaii as regards the observed distribution,

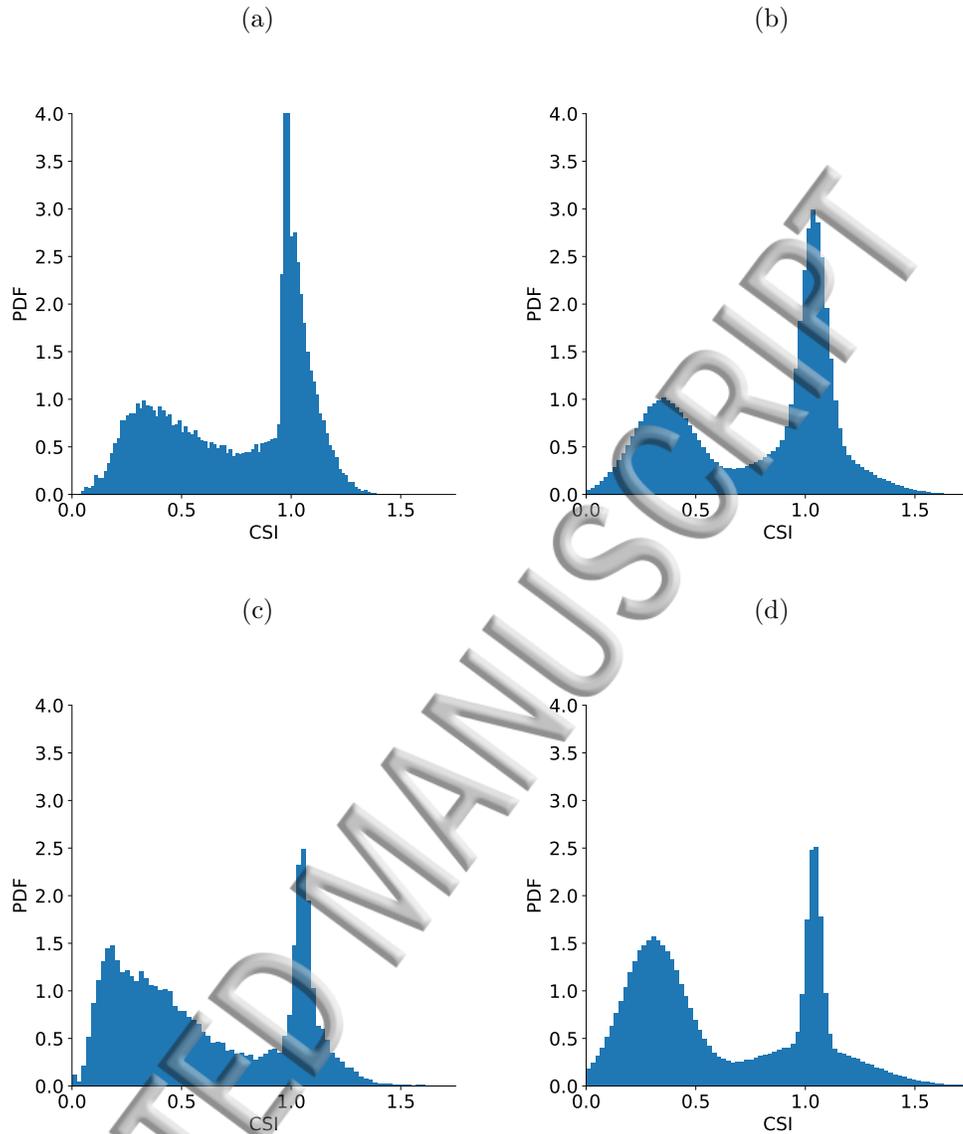


FIG. 7. Histograms of the CSI of the test data (a) and (c); and samples from the FMGHMMs (b) and (d). (a) and (b) present the results of Hawaii. (c) and (d) present the results of Norrköping.

375 which was expected from the results of the K-S test scores presented before.

376 Table IV depicts the parameters of the fitted distributions. Note that the mean of the
 377 obscured state μ_1 was estimated for both locations, using Eq. 9, to be 0.351 and 0.307 for
 378 Hawaii and Norrköping, respectively. However, the peaks of the first mode of the CSI PDF
 379 were at 0.333 and 0.19 for the Hawaii and Norrköping test data sets, respectively.

380 Another important observation is that the Baum-Welch algorithm did not fit the variances
 381 as perhaps expected. State 2, the state of unobscured with clear-sky, is expected to have a

TABLE IV. The observation distributions for the FMGHMM for both locations. All distributions are Gaussians.

| State | Hawaii | | Norrköping | |
|-------|---------|--------------|------------|--------------|
| | μ_i | σ_i^2 | μ_i | σ_i^2 |
| 1 | 0.351 | 0.0198 | 0.307 | 0.0204 |
| 2 | 1.0 | 0.0526 | 1.0 | 0.0692 |
| 3 | 1.04 | 0.0031 | 1.04 | 0.0010 |

382 narrow distribution, i.e., small variance, due to the concentration of the observed CSI around
 383 1.0 when the sky is clear. On the other hand, the distribution of state 3, unobscured with
 384 partially cloudy sky, is expected to be wider than that of state 2. As shown in Table IV, the
 385 Baum-Welch algorithm found that the optimal fitted models have a wider distribution in
 386 state 2 than in state 3. It should be noted that the Baum-Welch algorithm finds the model
 387 with highest log-likelihood \mathcal{L} .

388 Figure 8 shows the ACF of the training and test data along with samples from FMGHMMs
 389 for both locations for up to 60 minutes lags. The ACFs of the training and test data were
 390 similar. This was true for both locations. The models produced exponentially decaying
 391 ACFs as expected. However, there was a discrepancy between the ACFs of the model and
 392 that of the data in both locations. The MAE of the ACF ϵ was 0.22, or 22%, for both
 393 locations, see Table III.

394 Table III presents the performance metrics of the FMGHMM. In comparison Table II
 395 presented the same performance metrics of the GHMMs with $n = 3$. The log-likelihood \mathcal{L}
 396 was lower for both locations using the FMGHMM in comparison with the GHMM. The K-S
 397 test score K was higher for Hawaii using the FMGHMM compared with the GHMM, 0.118
 398 compared with 0.05. For Norrköping, the differences were negligibly small, 0.003, between
 399 the two models as regards the K of the test data set. The MAE of the ACF ϵ , when using
 400 the FMGHMM instead of the GHMM, was lower by 0.02 in the case of Hawaii, and higher
 401 by 0.06 in the case of Norrköping.

402 To summarize the comparison between the GHMM and the FMGHMM, the GHMM—
 403 when compared to the FMGHMM—was better for Hawaii in the K-S test score and worse
 404 in the MAE of the ACF, while for Norrköping, the GHMM was comparable in the K-S test

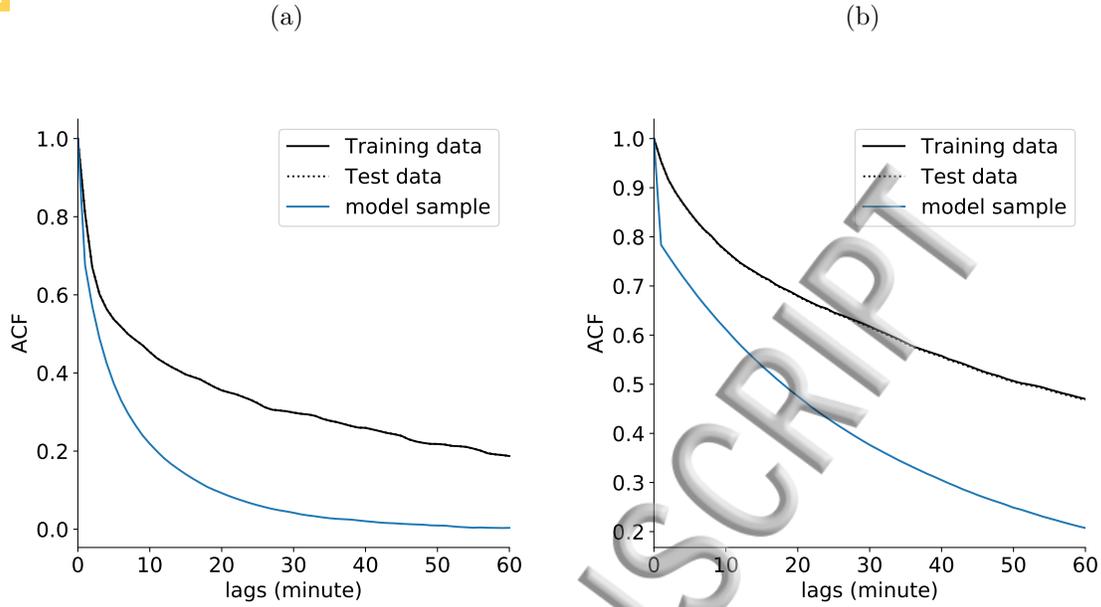


FIG. 8. The ACF for the fitted FMGHMMs, for Hawaii in (a) and Norrköping in (b). In both locations, the training and test data had the same ACF.

405 score and better in the MAE of the ACF.

406 C. Improving the autocorrelation function

407 In this section, we improve the ACF of the HMMs. The results of this improvement on
408 the FMGHMM are presented.

409 The ACF of the HMM can be calculated directly using Eq. 11. A careful inspection of
410 Eq. 11 shows that the ACF of the samples depends on both the hidden state trajectory and
411 the observation distributions. The transition matrix A controls how likely the process is to
412 persist in a certain hidden state, thus controlling the hidden state trajectory. The variables
413 $\boldsymbol{\pi}$, $\boldsymbol{\mu}$, and $\boldsymbol{\sigma}^2$ control the impacts of the dispersion of the observation distributions on the
414 ACF.

415 For example, we can take the distributions for Hawaii presented in Table I, and assume
416 an HMM with a hidden state trajectory $\{1, 1, \dots, 2, 2, \dots, 3, 3, \dots\}$, i.e., a highly persistent
417 hidden Markov process with $\boldsymbol{\pi} = [1/3, 1/3, 1/3]$. The ACF of such a process will highly
418 depend on the dispersion of the observation distributions, i.e., $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}^2$. This is a result
419 of the fact that the consecutive samples from the process are conditionally independent

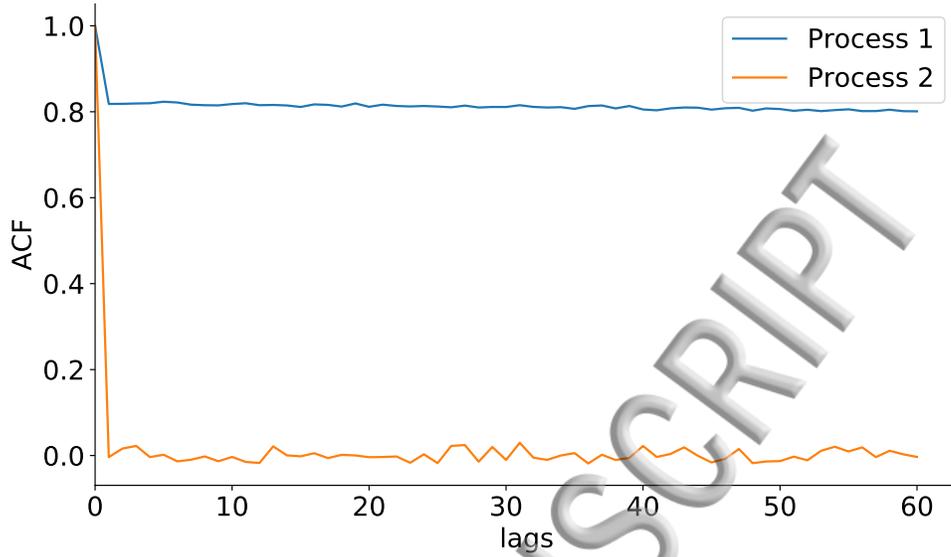


FIG. 9. An example of the ACF of samples from two HMM processes. Both processes have the same observation distributions as the ones presented for Hawaii in Table I, and both have $\pi = [1/3, 1/3, 1/3]$. Process 1 has a hidden state trajectory $\{1, 1, \dots, 2, 2, \dots, 3, 3, \dots\}$, process 2 randomly shuffles this hidden state trajectory. In other words, process 2 has an uncorrelated random hidden state trajectory.

420 given the states, more formally ($X_t \perp X_{t-1} | S_t$). This is to say that consecutive samples are
 421 independent even if they were drawn from the same distribution—or state, see Fig. 1.

422 On the other hand, an HMM with the same stationary distribution π and observation
 423 distributions as in the previous example but with a randomly shuffled hidden state trajectory
 424 will have a lower ACF for the same lags compared with the previous example. This is due
 425 to both the independent sampling of the observations and the rapidly switching hidden
 426 state trajectory. A comparison between the ACFs of the two example processes with 9000
 427 time-steps is provided in Fig. 9.

428 Any adjustments to the variables π , μ , and σ^2 of the model changes the distribution of
 429 the CSI generated by the model. Consequently, the only method to improve the ACF of the
 430 HMM is by adjusting the transition matrix A while maintaining the stationary distribution
 431 π constant.

432 As proven in the Appendix, a scaling constant $\phi \in (0, 1)$ can be multiplied by the off-
 433 diagonal elements of A followed by re-scaling of the diagonal elements using $p_{ii} = 1 -$

TABLE V. The performance metrics of the proposed scaling method tested on the test data and using the FMGHMM. Note that the case $\phi = 1$ is the original fitted models, presented before in Section III B.

| ϕ | Hawaii | | | Norrköping | | |
|--------|---------------|-------|------------|---------------|-------|------------|
| | \mathcal{L} | K | ϵ | \mathcal{L} | K | ϵ |
| 1 | 18897.70 | 0.118 | 0.22 | 24222.60 | 0.045 | 0.22 |
| 0.8 | 18820.36 | 0.117 | 0.19 | 24205.54 | 0.045 | 0.17 |
| 0.6 | 18541.54 | 0.116 | 0.15 | 24142.82 | 0.047 | 0.10 |
| 0.5 | 18288.97 | 0.114 | 0.12 | 24084.92 | 0.044 | 0.08 |
| 0.4 | 17922.33 | 0.115 | 0.09 | 23999.54 | 0.048 | 0.03 |
| 0.3 | 17383.42 | 0.117 | 0.05 | 23871.73 | 0.045 | 0.03 |
| 0.2 | 16547.13 | 0.112 | 0.06 | 23668.90 | 0.048 | 0.07 |

434 $\phi \sum_{j \neq i} p_{ij}$. Such a method, firstly, will ensure that A satisfies Eq. 3, secondly it will not
 435 change the stationary distribution π of the process, and thirdly it will alter the ACF by
 436 making the hidden state trajectory more persistent.

437 Here, the authors scaled the previously fitted FMGHMMs in Section III B to im-
 438 prove the ACF. The values of the scaling constant ϕ were chosen arbitrarily in the set
 439 $\{1, 0.8, 0.6, 0.5, 0.4, 0.3, 0.2\}$. Note that $\phi = 1$ represents the same models fitted in Sec-
 440 tion III B, i.e., no scaling takes place.

441 Table V compares the scaled models with regard to the log-likelihood \mathcal{L} , the K-S test score
 442 K , and the MAE of the ACF ϵ . As seen in the table, \mathcal{L} was highest for the previously fitted
 443 models with no scaling, i.e., $\phi = 1$. This explains why the Baum-Welch algorithm previously
 444 selected these models. The Baum-Welch algorithm blindly maximizes the log-likelihood with
 445 little regard to the remaining performance metrics⁴⁵.

446 Further research is needed to explain why improving the ACF deteriorated the log-
 447 likelihood in our case. The authors suspect that the intersecting distributions of states
 448 2 and 3 might explain this phenomenon. The distribution of state 3 is narrower than that
 449 of state 2, see Table IV, which might caused the model to have higher likelihood by switch-
 450 ing to this state in comparison with persisting in state 3. Nonetheless, further research
 451 investigating this behavior is needed, as stated before.

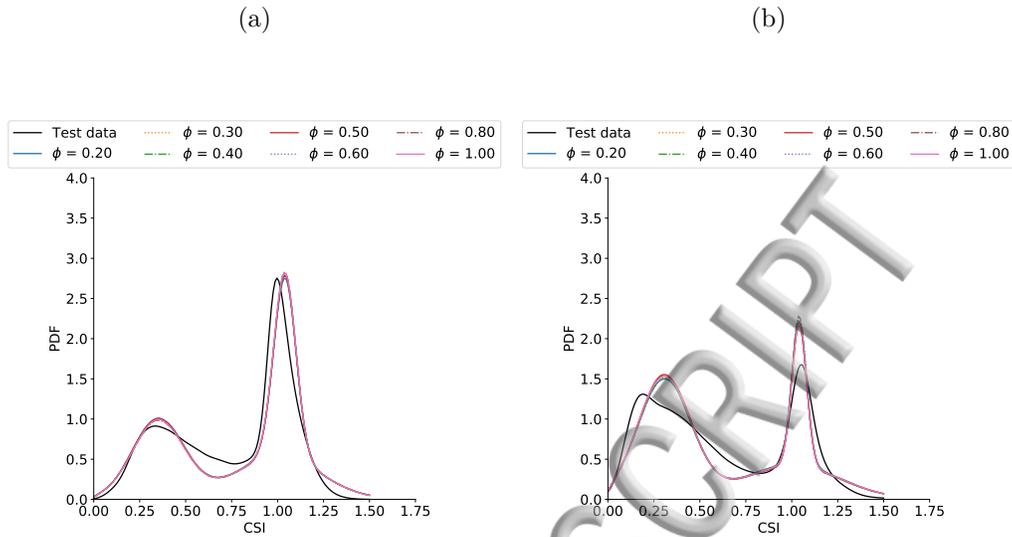


FIG. 10. Distribution of the test data and samples from the FMGHMMs with different scaling constants ϕ . (a) Hawaii, and (b) Norrköping.

452 Table V also shows that the differences in the K-S test score were negligibly small. Thus,
 453 the proposed scaling did not impact the generated CSI distribution, see also Fig. 10. On
 454 the other hand, the proposed scaling improved the MAE of the ACF ϵ compared with
 455 the original—unscaled—models. Figure 11 presents the ACFs of samples from the scaled
 456 models. A significant improvement in the ACFs can be observed.

457 D. Comparison with previously proposed models

458 In this section, a comparison between our proposed models and previous models is pro-
 459 vided. We compare our results to that of previous models which used the same data sets
 460 and performance metrics.

461 Table VI compares the models as regards the K-S test score and the MAE of the ACF.
 462 The GHMM did perform comparably with other models when it comes to the K-S test score
 463 K . However, when it comes to the MAE of the ACF ϵ it lagged behind other models. It
 464 is important to note that we did not scale the GHMM. If we scale the GHMMs using the
 465 method proposed in Section III C, the ACF likeness improved. For Hawaii scaling with
 466 $\phi = 0.2$ resulted in the lowest ϵ among the proposed values of ϕ . In this case, we achieved
 467 $\epsilon = 0.051$ and $K = 0.050$. Similarly for Norrköping, scaling with $\phi = 0.4$ achieved $\epsilon = 0.012$

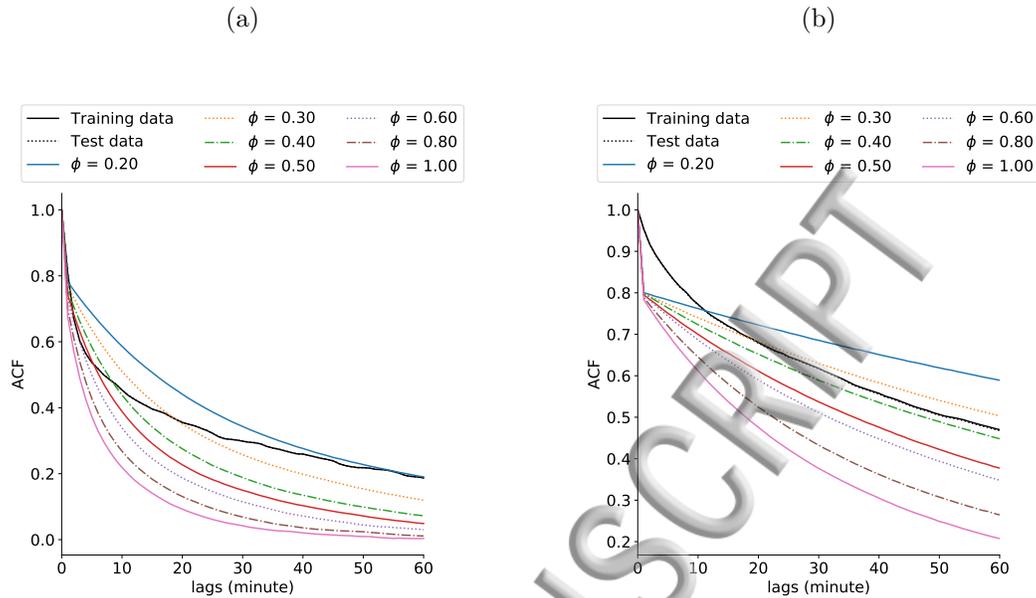


FIG. 11. The ACF of the test data and samples from the FMGHMMs with different scaling constants ϕ . (a) Hawaii, and (b) Norrköping. In both locations the training and test data had the same ACF.

468 and $K = 0.042$.

469 Compared with other proposed models, the FMGHMM scaled with $\phi = 0.3$ for Hawaii
 470 did not perform well as regards the K-S test score K . Neither did any FMGHMM presented
 471 in Table V. For Norrköping, this problem was not observed. Moreover, the MAEs of the
 472 ACF of the FMGHMM in both locations were close to that of previous models.

473 IV. DISCUSSION

474 This paper extends the previous models in Refs. 5 and 18 by considering a Gaussian
 475 observation distribution. In addition, it connects the new model to previous meteorological
 476 CSI models and proposes a method to improve the ACF likeness between samples from the
 477 fitted models and CSI data.

478 Two HMMs were developed: GHMM and FMGHMM. In the GHMM, the model pa-
 479 rameters were fully learned from the data. The FMGHMM pre-defined the means of the
 480 observation distribution of each hidden state. Still, the remaining model parameters were
 481 learned from the data.

TABLE VI. Comparison of the performance of the proposed models with the previously developed models that used the same data sets—the results of the previously developed models were presented before in Ref. 5.

| Model | Hawaii | | Norrköping | |
|-----------------------------|--------|------------|------------|------------|
| | K | ϵ | K | ϵ |
| Refs. 17 and 26 | 0.004 | 0.028 | 0.009 | 0.053 |
| Ref. 18 | 0.053 | 0.077 | 0.077 | 0.087 |
| Ref. 5, $n = 30$ | 0.047 | 0.038 | 0.044 | 0.023 |
| GHMM, $n = 3$ | 0.050 | 0.240 | 0.048 | 0.160 |
| FMGHMM, scaled $\phi = 0.3$ | 0.117 | 0.050 | 0.045 | 0.030 |

482 In the GHMM, we varied the number of hidden states from 2 to 12. Since the model
 483 parameters were freely learned from the data, the fitted distributions for $n = 3$ were not
 484 guaranteed to resemble the cloud cover categories proposed by Ref. 22, see Table I.

485 The K-S test score K of the GHMMs improved in both locations as the number of hidden
 486 states increased, see Figs. 3c and 3d. The MAE of the ACF ϵ also improved as the number
 487 of hidden states increased, however, it quickly stagnated for Hawaii, see Fig. 6a.

488 As regards the FMGHMM, three hidden states were defined: obscured, unobscured in
 489 clear sky, and unobscured in partially cloudy sky. These states were defined following
 490 Hollands and Suehrcke²², and they represent the cloud cover above a certain location. The
 491 means of the Gaussian observation distributions of the three hidden states were pre-defined
 492 using the fraction of bright sunshine. The Baum-Welch algorithm was then used to learn
 493 the remaining model parameters.

494 FMGHMMs produced K-S test score, K , of 0.118 for Hawaii and 0.045 for Norrköping.
 495 Unexpectedly, the variances of the fitted observation distributions were wider for the unob-
 496 scured in clear-sky state than for the unobscured in partially cloudy state, see Table IV.

497 Comparing the GHMM with $n = 3$ and FMGHMM showed that for Hawaii, the K-S test
 498 score K increased by 0.068 for Hawaii and decreased by 0.003 for Norrköping (negligible
 499 decrease). The MAE of the ACF ϵ decreased by 0.02 for Hawaii and increased by 0.06 for
 500 Norrköping.

501 Increasing the number of hidden states improved the GHMMs K-S test score. This,

502 Nonetheless, comes at the cost of the training time. Moreover, the K-S test scores of the
503 simpler GHMMs, the models with fewer hidden states, were still comparable with the pre-
504 viously developed models in the literature, e.g., Refs. 5, 17, 18, and 26.

505 Fitting the HMM using the Baum-Welch algorithm might cause the ACF of the HMM
506 to underperform when compared to the training data. A model with slightly lower log-
507 likelihood, might have higher ACF likeness compared with the Baum-Welch optimal model.
508 In this paper, a method to improve the ACF likeness of the HMM models is proposed. This
509 method ensures that the generated distribution from the HMM does not change.

510 Generating CSI time-series from a fitted model is a two-step process. In the first step, the
511 hidden state trajectory, i.e., the Markov process $\{S_1, \dots, S_T\}$, is sampled using the method
512 detailed in Ref. 58. In the second step, the CSI observations can be randomly sampled from
513 the distribution associated with each hidden state generated in the first step.

514 V. CONCLUSIONS

515 This paper showed that hidden Markov models (HMMs) with Gaussian observation dis-
516 tributions can be employed as generative models for the clear-sky index (CSI). These models
517 can also be adapted to pre-define the means of the Gaussian distribution associated with
518 each hidden state. Fitting the models solely by maximizing the log-likelihood might cause
519 the models to underperform when it comes to other performance metrics, e.g., autocorre-
520 lation function (ACF) likeness. Care must be taken to ensure that the ACF of the fitted
521 HMM resembles the ACF of the CSI data. In this paper, a novel method to improve the
522 ACF likeness of HMMs was proposed. Increasing the number of hidden states can improve
523 the goodness-of-fit of the CSI distribution. This, however, comes at the cost of the training
524 time, which grows as a function of the square of the number of hidden states.

525 Future contributions might expand this work by employing other distributions for the
526 hidden states, e.g., log-normal and Weibull distributions. In addition, future works might
527 explore the impacts of using autoregressive HMMs on the ACF likeness of the CSI model.
528 Testing the performance of the method on various CSI time-series and in different condi-
529 tions, e.g., solar altitude angles and air masses, is encouraged for future works. Finally,
530 developing a method, similar to the one proposed here, to generate the direct component of
531 the irradiance is left for future works.

ACKNOWLEDGMENTS

533 This work was financially supported by SamspeL 2016–2020 in the project “Development
 534 and evaluation of forecasting models for solar power and electricity use over space and time”,
 535 financed primarily by the Swedish Energy Agency. This work forms part of the Swedish
 536 strategic research programme StandUp for Energy.

537 Appendix: Improving the ACF of the HMM

As stated in Eq. 5, the stationary distribution π of the Markov chain can be estimated by solving

$$\begin{aligned} \pi A &= \pi, \\ \pi(A - I) &= \mathbf{0}, \end{aligned} \quad (\text{A.1})$$

538 where I is the identity matrix. Notice that multiplying both sides of Eq. A.1 by any constant
 539 ϕ does not change the solution. Equation A.1 can be expanded by substituting the diagonals
 540 of A by $p_{ii} = 1 - \sum_{j \neq i} p_{ij}$ to

$$541 \quad [\pi_1, \dots, \pi_n] \begin{bmatrix} -\sum_{j \neq 1} p_{1j} & \cdots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \cdots & -\sum_{j \neq n} p_{nj} \end{bmatrix} = \mathbf{0}. \quad (\text{A.2})$$

542 Equation A.2 indicates that any constant ϕ will only scale the off-diagonal elements of the
 543 transition matrix A .

544 The values of the constant ϕ need to satisfy some conditions. Firstly, ϕ cannot be 0,
 545 otherwise no single solution to Eq. A.1 exists. Secondly, all the off-diagonal elements and
 546 the diagonal elements, re-scaled as $p_{ii} = 1 - \phi \sum_{j \neq i} p_{ij}$, of the transition matrix A have to
 547 be in $[0,1]$ as they are probabilities.

548 Choosing $\phi \in (0,1)$ will decrease the off-diagonal transition probabilities and increase the
 549 diagonal transition probabilities of the transition matrix A . Consequently, the process will
 550 be more persistent, and still has the same stationary distribution π .

551 To further explain the effect of the scaling method note that the largest eigenvalue of the
 552 transition matrix A is always 1. And the stationary distribution π is the left eigenvector
 553 corresponding to this eigenvalue, see Eq. A.1 and⁵⁹ (p. 15). Thus the previously described

scaling procedure thus changes the remaining eigenvalues of the transition matrix A . Consequently, the scaling procedure changes the speed of convergence of the Markov chain. The speed of convergence of a Markov chain can be shown to follow⁵⁹ (p. 17)

$$|A^k - A^\infty| \leq \alpha |\lambda_2|^k, \quad (\text{A.3})$$

where α is a constant which satisfies $\alpha > 0$, and λ_i is the an eigenvalue of the matrix A such that $1 = |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$.

REFERENCES

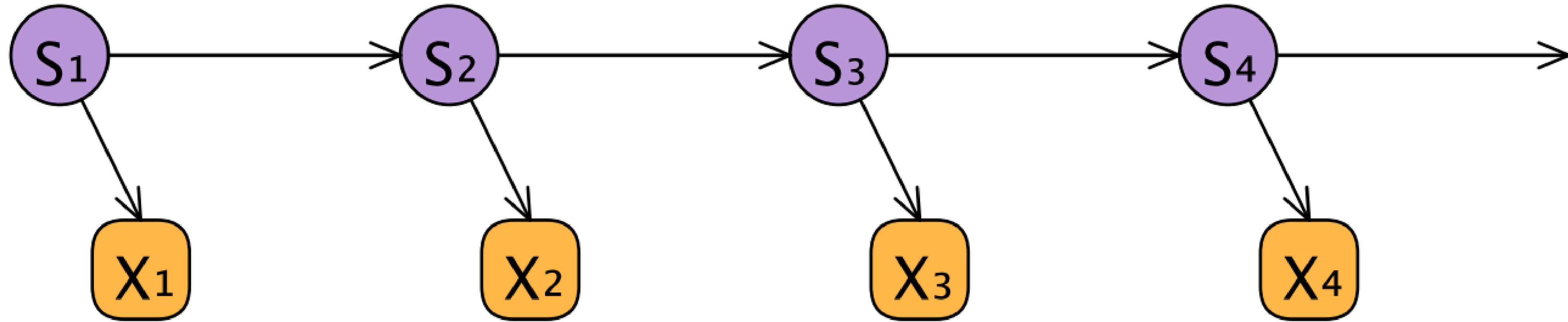
- ¹M. Lave, J. Kleissl, and E. Arias-Castro, “High-frequency irradiance fluctuations and geographic smoothing,” *Solar Energy* **86**, 2190–2199 (2012).
- ²M. H. Bollen and F. Hassan, *Integration of distributed generation in the power system* (John Wiley & Sons, 2011).
- ³J. Kleissl, *Solar energy forecasting and resource assessment* (Academic Press, 2013).
- ⁴M. Lave, M. J. Reno, and R. J. Broderick, “Characterizing local high-frequency solar variability and its impact to distribution studies,” *Solar Energy* **118**, 327–337 (2015).
- ⁵J. Munkhammar and J. Widén, “An N-state Markov-chain mixture distribution model of the clear-sky index,” *Solar Energy* **173**, 487–495 (2018).
- ⁶J. Deboever, S. Grijalva, M. J. Reno, and R. J. Broderick, “Fast quasi-static time-series (QSTS) for yearlong PV impact studies using vector quantization,” *Solar Energy* **159**, 538–547 (2018).
- ⁷C. F. Peruchena, M. Larrañeta, M. Blanco, and A. Bernardos, “High frequency generation of coupled GHI and DNI based on clustered Dynamic Paths,” *Solar Energy* **159**, 453–457 (2018).
- ⁸W. Zhang, W. Kleiber, A. R. Florita, B.-M. Hodge, and B. Mather, “A stochastic down-scaling approach for generating high-frequency solar irradiance scenarios,” *Solar Energy* **176**, 370–379 (2018).
- ⁹J. Widén and J. Munkhammar, “Spatio-temporal downscaling of hourly solar irradiance data using gaussian copulas,” in *46TH IEEE Photovoltaic Specialists Conference (PVSC)* (2019).
- ¹⁰J. M. Bright, “The impact of globally diverse GHI training data: evaluation through ap-

- 584 application of a simple Markov chain downscaling methodology,” [Renewable and Sustainable](#)
585 [Energy](#) **11**, 023703 (2019).
- 586 ¹¹G. Lohmann, “Irradiance variability quantification and small-scale averaging in space and
587 time: a short review,” [Atmosphere](#) **9**, 264 (2018).
- 588 ¹²J. M. Bright, C. J. Smith, P. G. Taylor, and R. Crook, “Stochastic generation of synthetic
589 minutely irradiance time series derived from mean hourly weather observation data,” [Solar](#)
590 [Energy](#) **115**, 229–242 (2015).
- 591 ¹³G. H. Yordanov, T. O. Saetre, and O.-M. Midtgård, “100-millisecond resolution for accu-
592 rate overirradiance measurements,” [IEEE Journal of Photovoltaics](#) **3**, 1354–1360 (2013).
- 593 ¹⁴G. M. Lohmann and A. H. Monahan, “Effects of temporal averaging on short-term irradi-
594 ance variability under mixed sky conditions,” [Atmospheric Measurement Techniques](#) **11**,
595 [3131–3144](#) (2018).
- 596 ¹⁵N. A. Engerer, J. M. Bright, and S. Killinger, “Himawari-8 enabled real-time distributed
597 PV simulations for distribution networks,” in [44th IEEE Photovoltaic Specialist Conference](#)
598 [\(PVSC\)](#) (2017) pp. 1405–1410.
- 599 ¹⁶Meteotest, “[Meteonorm irradiance software](#),” (2017), [Accessed 8 Aug 2017].
- 600 ¹⁷J. Munkhammar and J. Widén, “An autocorrelation-based copula model for generating
601 realistic clear-sky index time-series,” [Solar Energy](#) **158**, 9 – 19 (2017).
- 602 ¹⁸J. Munkhammar and J. Widén, “A Markov-chain probability distribution mixture ap-
603 proach to the clear-sky index,” [Solar Energy](#) **170**, 174–183 (2018).
- 604 ¹⁹B. J. Brinkworth, “Autocorrelation and stochastic modelling of insolation sequences,” [Solar](#)
605 [Energy](#) **19**, 343 – 347 (1977).
- 606 ²⁰A. Skartveit and J. A. Olseth, “The probability density and autocorrelation of short-term
607 global and beam irradiance,” [Solar Energy](#) **49**, 477–487 (1992).
- 608 ²¹J. Munkhammar, J. Rydén, J. Widén, and D. Lingfors, “Simulating dispersed photovoltaic
609 power generation using a bimodal mixture model of the clear-sky index,” in [31st European](#)
610 [Photovoltaic Solar Energy Conference and Exhibition \(EU PVSEC\)](#) (2015) pp. 1560 –
611 1567.
- 612 ²²K. T. Hollands and H. Suehrcke, “A three-state model for the probability distribution of
613 instantaneous solar radiation, with applications,” [Solar Energy](#) **96**, 103–112 (2013).
- 614 ²³J. Widén, M. Shepero, and J. Munkhammar, “On the properties of aggregate clear-sky
index distributions and an improved model for spatially correlated instantaneous solar

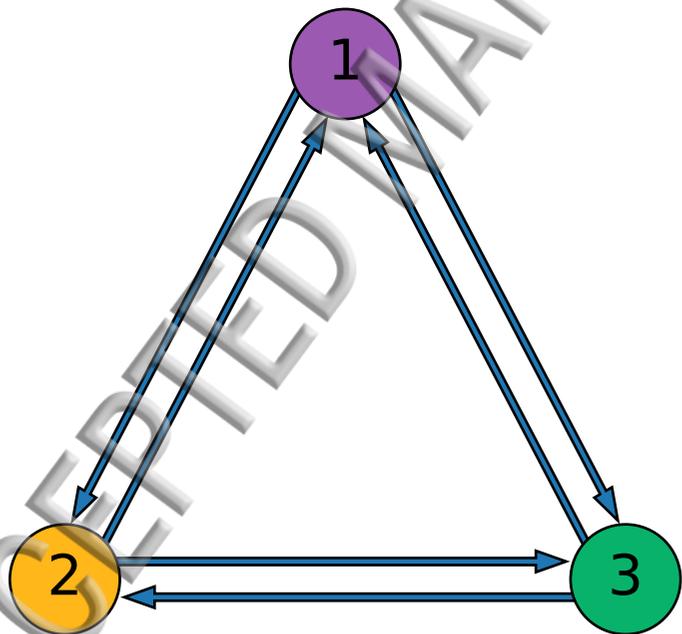
- 616 ²⁴Á. Frimane, T. Soubdhan, J. M. Bright, and M. Aggour, “Nonparametric Bayesian-based
617 recognition of solar irradiance conditions: Application to the generation of high temporal
618 resolution synthetic solar irradiance data,” *Solar Energy* **157**, 566–580 (2017).
- 619 ²⁵J. M. Bright, O. Babacan, J. Kleissl, P. G. Taylor, and R. Crook, “A synthetic, spatially
620 decorrelating solar irradiance generator and application to a LV grid model with high PV
621 penetration,” *Solar Energy* **147**, 83 – 98 (2017).
- 622 ²⁶J. Munkhammar and J. Widén, “An autocorrelation-based copula model for producing
623 realistic clear-sky index and photovoltaic power generation time-series,” in *44th IEEE
624 Photovoltaic Specialist Conference (PVSC)* (2017) pp. 3067–3072.
- 625 ²⁷A. P. Grantham, P. J. Pudney, and J. W. Boland, “Generating synthetic sequences of
626 global horizontal irradiation,” *Solar Energy* **162**, 500 – 509 (2018).
- 627 ²⁸R. Aguiar and M. Collares-Pereira, “TAG: A time-dependent, autoregressive, Gaussian
628 model for generating synthetic hourly radiation,” *Solar Energy* **49**, 167 – 174 (1992).
- 629 ²⁹C. Voyant, M. Muselli, C. Paoli, and M.-L. Nivet, “Optimization of an artificial neural
630 network dedicated to the multivariate forecasting of daily global radiation,” *Energy* **36**,
631 348 – 359 (2011).
- 632 ³⁰G. M. Lohmann, A. Hammer, A. H. Monahan, T. Schmidt, and D. Heinemann, “Simu-
633 lating clear-sky index increment correlations under mixed sky conditions using a fractal
634 cloud model,” *Solar Energy* **150**, 255 – 264 (2017).
- 635 ³¹H. Morf, “The stochastic two-state solar irradiance model (STSIM),” *Solar Energy* **62**,
636 101–112 (1998).
- 637 ³²R. J. Aguiar, M. Collares-Pereira, and J. P. Conde, “Simple procedure for generating
638 sequences of daily radiation values using a library of Markov transition matrices,” *Solar
639 Energy* **40**, 269–279 (1988).
- 640 ³³E. Palomo, “Hourly solar radiation time series as first-order Markov chains,” in *Actes du
641 International Solar Energy Society Solar World Congress* (1989) pp. 2146–2150.
- 642 ³⁴B. O. Ngoko, H. Sugihara, and T. Funaki, “Synthetic generation of high temporal resolu-
643 tion solar radiation data using Markov models,” *Solar Energy* **103**, 160–170 (2014).
- 644 ³⁵J. Wegener, M. Lave, J. Luoma, and J. Kleissl, “Temporal downscaling of irradiance
645 data via Hidden Markov Models on Wavelet coefficients: Application to California Solar
646 Initiative data,” UC San Diego 2012 (2012).

- 648 synthetic five-minute solar irradiance values from hourly observations,” *Solar Energy* **147**,
649 [209–221 \(2017\)](#).
- 650 ³⁷L. R. Rabiner, “A tutorial on hidden Markov models and selected applications in speech
651 recognition,” *Proceedings of the IEEE* **77**, 257–286 (1989).
- 652 ³⁸D. Jurafsky and J. H. Martin, *Speech and Language Processing: An Introduction to Nat-
653 ural Language Processing, Computational Linguistics, and Speech Recognition*, 2nd ed.
654 (Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 2009).
- 655 ³⁹I. L. MacDonald and W. Zucchini, *Hidden Markov and other models for discrete-valued
656 time series*, Vol. 110 (CRC Press, 1997).
- 657 ⁴⁰S. Frühwirth-Schnatter, *Finite mixture and Markov switching models* (Springer Science &
658 Business Media, 2006).
- 659 ⁴¹H. Suehrcke, “On the relationship between duration of sunshine and solar radiation on the
660 earth’s surface: Ångström’s equation revisited,” *Solar Energy* **68**, 417–425 (2000).
- 661 ⁴²H. Suehrcke, R. S. Bowden, and K. T. Hollands, “Relationship between sunshine duration
662 and solar radiation,” *Solar Energy* **92**, 160–171 (2013).
- 663 ⁴³P. M. Broersen, *Automatic autocorrelation and spectral analysis* (Springer Science & Busi-
664 ness Media, 2006).
- 665 ⁴⁴R. J. Hyndman and G. Athanasopoulos, *Forecasting: principles and practice* (OTexts:
666 Melbourne, Australia, 2018) [Accessed 25 Feb 2019].
- 667 ⁴⁵T. Rydén, T. Teräsvirta, and S. Åsbrink, “Stylized facts of daily return series and the
668 hidden Markov model,” *Journal of applied econometrics* **13**, 217–244 (1998).
- 669 ⁴⁶M. Jurado, J. M. Caridad, and V. Ruiz, “Statistical distribution of the clearness index with
670 radiation data integrated over five minute intervals,” *Solar Energy* **55**, 469–473 (1995).
- 671 ⁴⁷M. J. Beal, Z. Ghahramani, and C. E. Rasmussen, “The infinite hidden Markov model,”
672 in *Advances in neural information processing systems* (2002) pp. 577–584.
- 673 ⁴⁸<https://github.com/hmmlearn/hmmlearn/>.
- 674 ⁴⁹L. Theis, A. v. d. Oord, and M. Bethge, “A note on the evaluation of generative models,”
675 in *International Conference on Learning Representations (ICLR 2016)* (2016) pp. 1–10.
- 676 ⁵⁰F. J. Massey Jr, “The Kolmogorov-Smirnov test for goodness of fit,” *Journal of the Amer-
677 ican statistical Association* **46**, 68–78 (1951).
- 678 ⁵¹SMHI, “Norrköping, Sweden: SMHI,” (2008).

- 52 M. Sengupta and A. Andreas, “Oahu solar measurement grid (1-year archive): 1-second
680 solar irradiance; Oahu, Hawaii (data),” Tech. Rep. NREL/DA-5500-56506 (National Re-
681 newable Energy Lab.(NREL), Golden, CO (United States), 2010).
- 682 ⁵³Kipp & Zonen, “CM 21 precision pyranometer instruction manual,” (2019), [Accessed 28
683 June 2019].
- 684 ⁵⁴LI-COR, “LI-200R Pyranometer,” (2019), [Accessed 28 June 2019].
- 685 ⁵⁵P. Ineichen and R. Perez, “A new airmass independent formulation for the Linke turbidity
686 coefficient,” *Solar Energy* **73**, 151–157 (2002).
- 687 ⁵⁶Soda-service, “CAMS McClear Service for estimating irradiation under clear-sky,” (2016),
688 [Accessed 5 May 2016].
- 689 ⁵⁷X. Sun, J. M. Bright, C. A. Gueymard, B. Acord, P. Wang, and N. A. Engerer, “Worldwide
690 performance assessment of 75 global clear-sky irradiance models using Principal Compo-
691 nent Analysis,” *Renewable and Sustainable Energy Reviews* **111**, 550–570 (2019).
- 692 ⁵⁸J. Widén, A. M. Nilsson, and E. Wäckelgård, “A combined Markov-chain and bottom-up
693 approach to modelling of domestic lighting demand,” *Energy and Buildings* **41**, 1001–1012
694 (2009).
- 695 ⁵⁹P. Lorek, *Speed of convergence to stationarity for stochastically monotone Markov chains*,
696 Ph.D. thesis, University of Wrocław (2007).



Hidden



Observed

