

# Verification of solar irradiance probabilistic forecasts

Philippe Lauret<sup>a,\*</sup>, Mathieu David<sup>a</sup>, Pierre Pinson<sup>b</sup>

<sup>a</sup>*University of La Réunion - PIMENT laboratory, 15, avenue René Cassin, 97715 Saint-Denis*

<sup>b</sup>*Technical University of Denmark, Centre for Electric Power and Energy, 2800 Kgs. Lyngby, Denmark*

---

## Abstract

We propose a framework for evaluating the quality of solar irradiance probabilistic forecasts. The verification framework is based on visual diagnostic tools and a set of scoring rules mostly originating from the weather forecast verification community. Two types of probabilistic forecasts are used as a basis to illustrate the application of these verification approaches. The first one consists in ensemble forecasts commonly provided by national or international meteorological centres. The second one originates from statistical methods and produces a set of discrete quantile forecasts, the nominal proportions of which span the unit interval. These probabilistic forecasts are evaluated for two selected sites that experience very different climatic conditions. The first site is located in the continental US while the second one is situated on La Réunion Island. Although visual diagnostic tools can help identify deficiencies in generated forecasts, it is recommended that a set of numerical scores be used to assess the quality of probabilistic forecasts. In particular, the Continuous Ranked Probability Score (CRPS) seems to have all the features needed to evaluate a probabilistic forecasting system and, as such, may become a standard for verifying solar irradiance probabilistic forecasts and by extension probabilistic forecasts of solar power generation.

*Keywords:* probabilistic solar forecasting, evaluation framework, diagnostic tools, scoring rules, CRPS, Ignorance Score

---

## 1. Introduction

Forecasts of solar energy generation are of utmost importance for efficiently integrating solar power generation into existing power grids and to decrease associated costs. Indeed, power production from photovoltaic (PV) or solar thermal plants is highly variable since weather dependent. Therefore, accurate knowledge of the future production from solar power generation capacities is necessary to limit the needs for additional balancing services and potentially storage. Therefore, increasing the value of solar power generation through the improvement of solar irradiance or PV power forecasting models (both usually referred to

---

\*corresponding author

*Email addresses:* philippe.lauret@univ-reunion.fr (Philippe Lauret),  
mathieu.david@univ-reunion.fr (Mathieu David), ppin@elektro.dtu.dk (Pierre Pinson)

9 as “solar forecasting models”) is of paramount importance. In the realm of solar irradiance  
10 forecasting, Global Horizontal Irradiance (GHI) is a prominent key variable. Therefore, this  
11 work will use this variable to illustrate the application of the proposed evaluation framework.

12 Numerous works have been devoted to the development of models that generate point  
13 forecasts of solar power generation, commonly referred to as deterministic forecasts. Some  
14 of these models can be found in (Reikard, 2009; Dambreville et al., 2014; Marquez and  
15 Coimbra, 2011; Coimbra et al., 2013; Huang et al., 2013; Lauret et al., 2015; Voyant et al.,  
16 2017; Pedro and Coimbra, 2015; Lorenz and Heinemann, 2012). Furthermore, error metrics  
17 dedicated to evaluating the accuracy of these deterministic forecasts, like Mean Bias Error  
18 (MBE), Root Mean Square Error (RMSE) and Mean Absolute Error (MAE) together with  
19 skill-score measures (Hoff et al., 2013; Coimbra et al., 2013), are now quite standard and  
20 well accepted by the solar forecasting community.

21 However, a forecast is inherently uncertain and in a context of decision-making faced by  
22 the grid operator, a point forecast plus an uncertainty (or, better say, prediction) interval is  
23 of genuine added value. Put differently, reliable probabilistic predictions may contribute to a  
24 more efficient integration of intermittent sources in the energy network (Morales et al., 2014).  
25 Contrary to the wind power forecasting community where probabilistic forecasting appears  
26 to be a mature subject (Morales et al., 2014; Iversen et al., 2016; Jung and Broadwater,  
27 2014; Pinson et al., 2007), probabilistic solar forecasting is still in its infancy (Hong et al.,  
28 2016) albeit some recent works (Zamo et al., 2014; Sperati et al., 2016; Alessandrini et al.,  
29 2015; Grantham et al., 2016; Ben Bouallègue, 2015; David et al., 2016; Golestaneh et al.,  
30 2016b) tend to moderate this statement.

31 As mentioned by Pinson et al. (2007), the assessment of probabilistic forecasts is more  
32 complicated than for deterministic ones. Figure 1 shows an example of GHI probabilistic  
33 forecasts where point forecasts are enriched with prediction intervals (PIs). From the visual  
34 inspection of Figure 1, the deviation between observed GHI (black line) and deterministic  
35 forecasts (blue dashed line) is easily assessed. But it is quite difficult to state whether the  
36 prediction intervals are good or not. To objectively assess the performance of probabilistic  
37 forecasts and the methods used to generate those, it is necessary to employ appropriate  
38 diagnostic tools and quantitative scores.

39 According to Murphy (1993), goodness of weather forecasts can be characterized by three  
40 types namely consistency, quality and value. Consistency is related to the correspondence  
41 between forecasters’ judgment and their forecasts. Quality refers to the correspondence  
42 between forecasts and the observations and value is linked to the benefit (economical or  
43 others) gained from the use of these probabilistic forecasts in an operational context. In this  
44 work, we concentrate on the assessment of the quality of the models.

45 Several attributes characterize the quality of probabilistic forecasts (Wilks, 2014; Jolliffe  
46 and Stephenson, 2003) but two main properties, i.e. reliability and resolution are used  
47 to measure the quality of the forecasts (Jolliffe and Stephenson, 2003). A third attribute  
48 namely sharpness can be used to evaluate how informative the forecasts are. In the weather  
49 forecasting verification community, several diagnostic tools are used to characterize these  
50 required properties of reliability, resolution and sharpness. One can cite among others the  
51 reliability diagram (Pinson et al., 2010; Wilks, 2014) and rank histogram (Hamill, 2001;

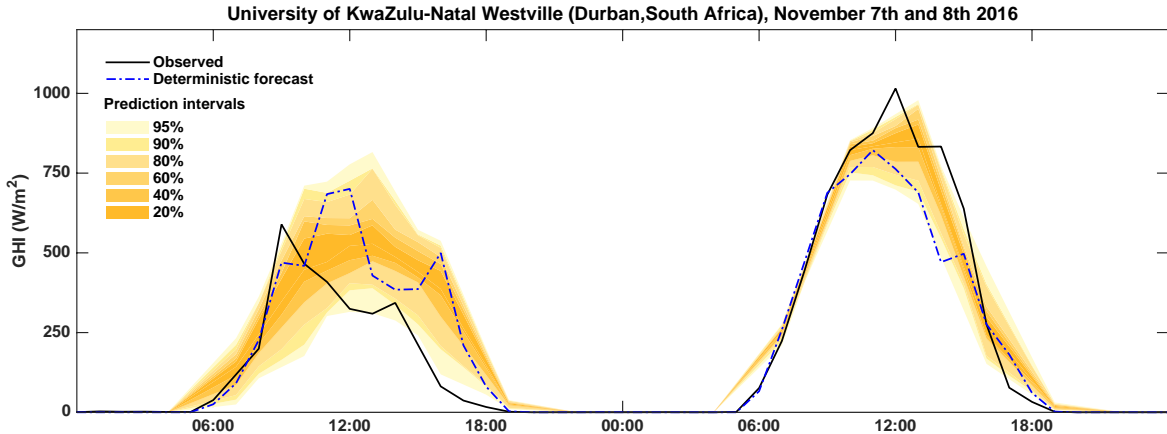


Figure 1: Example of probabilistic solar irradiance forecasts: 2 days of measured GHI at the University of KuwaZulu-Natal Westville (South Africa) (Brooks et al., 2015) and associated forecasts with prediction intervals provided by ECMWF-EPS.

52 Wilks, 2014) for assessing reliability. Regarding forecasts of continuous variable, there is  
 53 currently no visual tool to assess resolution. The sharpness property can be evaluated  
 54 through the use of sharpness diagrams (Pinson et al., 2007; Gneiting et al., 2007).

55 In addition to these tools that permit to visually assess the attributes of a forecasting  
 56 system, a metric called continuous ranked probability score (CRPS) (Hersbach, 2000) is  
 57 commonly used by the weather forecasting community to objectively quantify the overall skill  
 58 of the probabilistic forecasts. The CRPS is a metric capable of addressing both reliability  
 59 and resolution simultaneously. Indeed, the CRPS can be decomposed into three components  
 60 namely reliability, resolution and uncertainty. This decomposition provides a detailed picture  
 61 of the performance of the forecasting methods (Hersbach, 2000) and consequently may help  
 62 in the ranking of the probabilistic forecasts. A scoring rule originated from the information  
 63 theory called the logarithm or ignorance score metric has also been proposed for assessing  
 64 the quality of weather probabilistic forecasts (Roulston and Smith, 2002; Pinson et al., 2012).

65 Although solar probabilistic forecasting is not as mature as wind probabilistic forecasting  
 66 (Hong et al., 2016), some recent works (Alessandrini et al., 2015; Sperati et al., 2016; Zamo  
 67 et al., 2014; Grantham et al., 2016; David et al., 2016, 2018; Chu and Coimbra, 2017;  
 68 Golestaneh et al., 2016b; Verbois et al., 2018) proposed to assess the quality of the models  
 69 with some classical diagnostic tools originated from the weather verification community  
 70 like rank histogram and reliability diagram. This literature review also revealed that the  
 71 CRPS is a commonly used scoring rule. However, in our opinion, most of these works  
 72 did not conduct a detailed analysis of how to use and interpret the verification tools. For  
 73 instance, the CRPS formula proposed by (Hersbach, 2000) is restricted to ensemble forecasts  
 74 but David et al. (2018) and Lauret et al. (2017) used it to compute the CRPS of discrete  
 75 quantile forecasts. Moreover, most of the previous works that evaluated the overall skill  
 76 of competing methods through the use of the CRPS did not attempt to have a detailed  
 77 performance of the methods which is possible from the decomposition of the CRPS into

78 reliability, resolution and uncertainty. Besides, to our best knowledge, the ignorance score  
79 is not currently used by the solar forecasting community.

80 In addition, other metrics are proposed to assess the properties of prediction intervals  
81 such as Prediction Interval Coverage Probability (PICP), Prediction Interval Normalized  
82 Averaged Width (PINAW) (Khosravi et al., 2013; Chu and Coimbra, 2017; Lauret et al.,  
83 2017). PICP is related to the reliability of the probabilistic forecasts while PINAW gives  
84 a measure of the sharpness of the predictive distributions. However, as discussed below,  
85 these two metrics (PICP and PINAW) are not the most appropriate for measuring the  
86 quality of interval forecasts. It is also worth noting that a metric called coverage width-  
87 based criterion (CWC), which assesses the quality of the prediction intervals by combining  
88 PICP and PINAW has been proposed by (Khosravi et al., 2013). But as demonstrated by  
89 (Pinson and Tastu, 2014), this score can lead to possible misinterpretations of the results.  
90 Unfortunately, some researchers in the solar community (Scolari et al., 2016; Chu et al.,  
91 2015; Li et al., 2018) recently used this metric to assess the quality of their forecasting  
92 models. Furthermore, the CWC score has been recently cited in a reference paper (Yang  
93 et al., 2018) and a review paper (van der Meer et al., 2018).

94 This is why, we think that now is the time to take stock on the evaluation metrics of  
95 solar probabilistic forecasts. The objective of this work is therefore to provide the forecasting  
96 solar community a comprehensive overview of diagnostic tools and scoring rules that can  
97 be used to assess the performance of probabilistic forecasting methods. In particular, we  
98 propose an evaluation framework that may help the user to consistently evaluate the quality  
99 of the models. In others words, this paper aims at explaining how one should assess the  
100 quality of the probabilistic forecasts and how diagnostic tools and scores should be used and  
101 interpreted. In addition, we will propose a measure of resolution (through the decomposition  
102 of the CRPS) as this attribute is not currently assessed in the literature.

103 In this paper, two types of GHI probabilistic forecasts are used to illustrate the pro-  
104 posed verification framework. The first one is the ensemble forecast commonly provided  
105 by Ensemble Prediction Systems (EPS) of the Numerical Weather Predictions (NWP) of  
106 meteorological utilities such as ECMWF. The second one, denoted by quantile forecasts,  
107 is based on statistical methods and produces a set of quantiles spanning the unit interval.  
108 Both types generate forecasts represented by predictive distributions that can be modelled  
109 either by a Cumulative distribution function (CDF) or a Probability distribution Function  
110 (PDF).

111 Finally, note that in this paper, we restrict ourselves to the univariate context that corre-  
112 sponds to probabilistic forecasts that do not take into account spatio-temporal dependencies  
113 that are generated by stochastic processes like for instance cloud passing. The interested  
114 reader is referred to (Golestaneh et al., 2016a) who proposed a method to capture the  
115 spatio-temporal correlations in PV forecasts.

116 The remainder of this paper is organized as follows. Section 2 defines the probabilistic  
117 forecast as the estimation of a predictive distribution of the variable of interest (GHI in our  
118 case). Section 3 lists the properties required for skillful probabilistic forecasts while section  
119 4 proposes a verification framework and presents in details the verification tools. Section 5  
120 illustrates the application of these diagnostic tools on day-ahead forecasts provided by an

121 EPS and intra-day forecasts provided by quantile regression techniques. Finally, section 6  
122 gives some concluding remarks.

## 123 2. Nature of probabilistic forecasts of continuous variables

124 Probabilistic forecasts correspond to the estimation of the statistical distribution of a  
125 future event. Thus, a probabilistic forecast may be defined as a cumulative distribution  
126 function (CDF)  $F$  of a random variable  $X$ , such that  $F(x) = Pr(X \leq x)$ . This CDF can be  
127 summarized by a set of quantiles. The quantile  $q_\tau$ , at probability level  $\tau \in [0, 1]$ , is defined  
128 as follow

$$q_\tau = F^{-1}(\tau) = \inf\{x : F(x) \geq \tau\}. \quad (1)$$

129 A quantile  $q_\tau$  informs there is a probability  $\tau$  that the event  $x$  materializes below that  
130 quantile  $q_\tau$ . From a set of quantiles, prediction intervals (PIs) can be deduced. PIs define the  
131 range of values within which the observation is expected to be with a certain probability i.e.  
132 its nominal coverage rate (Pinson et al., 2007). To completely determine a PI, it is necessary  
133 to choose the way it should be centered on the probability density function (Pinson et al.,  
134 2007). The most common way is to center the PI on the median. Consequently, there is  
135 the same probability of risk below and above the median. Therefore, a central PI with a  
136 coverage rate of  $(1 - \alpha)100\%$  is estimated by using the  $\alpha/2$  quantile ( $\hat{q}_{\tau=\alpha/2}$ ) as the lower  
137 bound and the  $(1 - \alpha/2)$  quantile ( $\hat{q}_{\tau=1-\alpha/2}$ ) as the upper bound. More precisely, a PI with  
138  $(1 - \alpha)100\%$  nominal coverage rate is given by

$$\widehat{PI}_{(1-\alpha)100\%} = [\hat{q}_{\tau=\alpha/2}, \hat{q}_{\tau=1-\alpha/2}]. \quad (2)$$

139 In the realm of weather predictions, three ways to define this cumulative distribution  
140 are available: parametric CDFs, discrete estimates of a CDF via a non-parametric method  
141 and ensemble forecasts. Parametric CDFs are easy to set up and to assess. Nevertheless,  
142 regarding solar forecasts, they are seldom proposed in the literature because they suffer  
143 from a lack of calibration. Indeed, the distribution of future observations of the solar power  
144 can not be accurately reproduced by a single probabilistic law. David et al. (2016) gave an  
145 example with the GARCH model that assumes a Gaussian distribution.

146 An alternative to the parametric approach is the generation of discrete estimates of a  
147 CDF. This non-parametric method allows defining a predictive CDF without any assumption  
148 on the distribution of the future event. The forecast is provided as a set of quantiles spanning  
149 the unit interval. This kind of probabilistic forecast is also called quantile forecasts (Pinson  
150 et al., 2007). The Global Energy Forecasting Competition 2014 (GEFCom 2014) (Hong  
151 et al., 2016) is a good example of this approach. Indeed, the solar forecasts were to be  
152 expressed in the form of 99 quantiles with various nominal proportions between zero and  
153 one. Widely used statistical models, like Quantile Regressions (QR) or Gradient Boosting  
154 Decision Trees (GBDT) can estimate these predictive distributions.

155 The last type corresponds to ensemble forecasts classically generated by Numerical  
156 Weather Predictions (NWP) models. The distribution of the future event is given by an

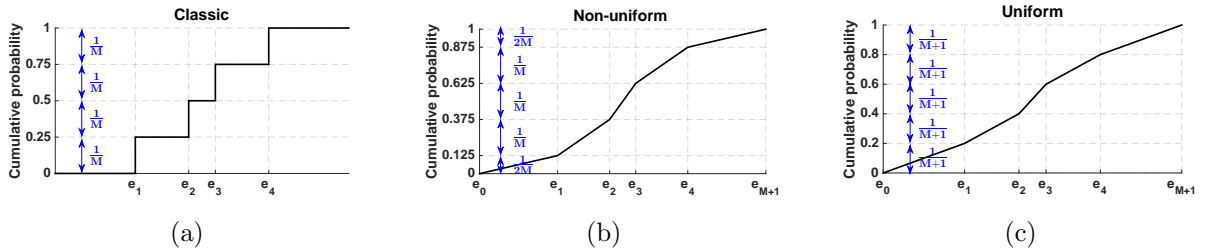


Figure 2: Different definitions of the CDF derived from an ensemble forecast ( $M = 4$ ): (a) classical; (b) non-uniform spacing of the cumulative probabilities and a linear interpolation between the members; (c) uniform spacing and a linear interpolation between the members.

157 ensemble of members that are not directly linked to the notion of quantiles. For example, in  
 158 the case of a NWP model, an ensemble forecast corresponds to a perturbed set of forecasts  
 159 computed by slightly changing the initial conditions of the control run and of the modeling  
 160 of unresolved phenomena (Leutbecher and Palmer, 2008). This ensemble prediction system  
 161 (EPS) allows representing the uncertainties of the prediction scheme. Nevertheless, ensem-  
 162 ble forecasts can be seen as discrete estimates of a CDF when they are sorted in ascending  
 163 order. In the literature, different ways to associate these sorted members to cumulative  
 164 probabilities are proposed. Considering  $M$  sorted members of an ensemble  $E = (e_1, \dots, e_M)$ ,  
 165 the most common definition in the domain of weather forecast assessment states that there  
 166 is a probability of  $1/M$  that the observation falls between two consecutive members  $e_j$  and  
 167  $e_{j+1}$  (Anderson, 1996; Hersbach, 2000). If we assign a null probability for future events that  
 168 fall outside the ensemble (i.e.  $x_{obs} < e_1$  or  $x_{obs} > e_M$ ), the predictive distribution can be  
 169 seen as a piecewise constant function

$$\widehat{F}(x) = \sum_{k=1}^M \alpha_k H(x - e_k). \quad (3)$$

170  $H$  is the Heaviside function which is 1 if the argument is positive and zero otherwise.  
 171 The weight  $\alpha_k = 1/M$  corresponds to the jump of probability that happens when  $x = e_k$ .  
 172 Figure 2(a) gives a visual representation of this classical definition of a CDF derived from  
 173 an ensemble with 4 members ( $M = 4$ ).

174 In the case of continuous variable, as the solar irradiance (GHI), the shape of the CDF  
 175 resulting from the preceding definition is obviously not realistic. Several works (Bröcker,  
 176 2012; Roulston and Smith, 2002; Pinson et al., 2010) proposed alternative approaches to  
 177 face this issue. Among others, these alternatives allow defining a continuous predictive  
 178 distribution and non-null probabilities outside the ensemble. We briefly present two other  
 179 ways to build a CDF from an ensemble forecast.

180 First, Bröcker (2012) proposes to preserve a jump of  $1/M$  between two members but to  
 181 assign a probability mass of  $1/2M$  for the events that fall outside of the ensemble. It results  
 182 in a non-uniform partition of the probability space  $[0; 1]$ . Figure 2(b) gives an example of  
 183 this definition for an ensemble with 4 members ( $M = 4$ ) and a linear interpolation between  
 184 the members. The tails of the distributions are bounded by  $e_0$  and  $e_{M+1}$ . The choice of

185 these limits are arbitrary. For continuous variables, Roulston and Smith (2002) proposed  
186 to use the minimum and the maximum of the climatology. Notice that this non-uniform  
187 definition amounts to consider each ensemble member  $i$  as a quantile with probability level  
188  $\tau(i) = \frac{i-0.5}{M}$ .

189 The second approach, described by (Pinson et al., 2010; Bröcker, 2012), assigns a prob-  
190 ability mass of  $1/(M + 1)$  between two members and for the events that fall outside of the  
191 ensemble. Note that using this definition that an ensemble member can be interpreted as  
192 a quantile forecast by considering its rank within the ensemble. The probability level  $\tau(i)$   
193 associated with the member of rank  $i$  is defined as:  $\tau(i) = \frac{i}{M+1}$ . This approach leads to  
194 an uniform spacing of the cumulative probabilities. Figure 2(c) presents graphically the  
195 shape of the CDF when considering this last definition and a linear interpolation between  
196 the members. As for the non-uniform definition, the boundaries of the CDF,  $e_0$  and  $e_{M+1}$ ,  
197 are arbitrarily chosen (see appendix A for more details).

198 Thus, when dealing with ensemble forecasts, three ways to build the CDF from the mem-  
199 bers are available. Unfortunately no definition can be favoured and each CDF construction  
200 has its pros and cons. The classic definition is the most used, specifically to compute the  
201 Continuous Rank Probability Score (CRPS, see section 4.2.3) with the methodology pro-  
202 posed by (Hersbach, 2000). As this commonly used definition assigns null probabilities to  
203 the events that fall outside of the ensemble, it can not be used to derive scores like ignorance  
204 (see section 4.2.4). The uniform and the non-uniform definitions requires to arbitrarily fix  
205 the boundaries of the CDF. Therefore, they are user dependent. Nevertheless, they allow  
206 designing continuous CDF that contains all the possible events. Thus, the procedure used to  
207 verify the quality of ensemble forecasts can be exactly the same as for the parametric CDFs  
208 or for the predictive distributions summarized by discrete quantiles estimated by some kind  
209 of statistical method. Bröcker (2012) showed that the non-uniform definition corresponds  
210 to a minimization of the CRPS. But, considering this definition, the optimal shape of the  
211 corresponding rank histogram (see section 4.1.2) is not flat. Indeed for this visual verifica-  
212 tion tool, the height of the first and last ranks should be the half of the other ones. Finally,  
213 if the aim is to compare different forecasting models, whatever the chosen definition, the  
214 ranking will remain the same. Nevertheless, a unique framework has to be defined to allow  
215 the comparison of different works.

### 216 3. Required properties for a skillful probabilistic system

217 As mentioned in the introduction, two main attributes (reliability and resolution) char-  
218 acterize the quality of probabilistic forecasts (Pinson et al., 2007). The evaluation of these  
219 two attributes can be complemented by a sharpness assessment.

#### 220 3.1. Reliability

221 Reliability or calibration refers to the statistical consistency between the forecasts and  
222 the observations. In other terms, the nominal coverage rate of the prediction intervals should  
223 be equal to the empirical one (e.g. a 90% PI should cover 90% of the observations). The  
224 reliability property is an important prerequisite as non reliable forecasts would lead to a  
225 systematic bias in subsequent decision-making processes (Pinson et al., 2007).

### 226 3.2. Resolution and sharpness

227 Resolution measures the capacity of a forecasting model to issue forecasts that are case-  
228 dependent. This important property, which is not easy to catch, is commonly not considered  
229 by the solar forecasting community. To understand concretely what resolution is, we will  
230 first define the climatological forecast (i.e. climatology). Imagine a distribution built from  
231 all the available past data of the parameter to forecast. The climatological forecast uses  
232 this unique distribution to forecast any future events. A high resolution forecasting system  
233 generates forecasts that differ from the climatology and, as a consequence, forecasts that are  
234 significantly different from each other. Climatological forecasts are perfectly reliable though  
235 having no resolution. Consequently, a skillful probabilistic forecasting system should issue  
236 reliable forecasts and with high resolution.

237 Sharpness evaluates how informative the forecasts are. Practically, sharpness refers to  
238 the concentration of the predictive distributions (Pinson et al., 2007; Gneiting et al., 2007)  
239 and can be measured by the average width of the prediction intervals. Unlike the two  
240 previous attributes, sharpness is a function of the forecasts only and does not depend on  
241 the observations. Consequently, a forecasting system can produce sharp forecasts yet being  
242 useless if those probabilistic forecasts are not reliable.

243 It must be emphasized here that these two components (sharpness and resolution) have  
244 different interpretations according a meteorologist’s point of view or a statistician’s point  
245 of view. In the meteorological literature (Wilks, 2014; Jolliffe and Stephenson, 2003), the  
246 sharpness property refers to the ability of a forecasting system to generate forecasts that are  
247 able to deviate from the climatological value of the variable to predict (also called predictand)  
248 whereas from a statistical point of view the sharpness property relates to the concentration  
249 of the predictive distributions (Pinson et al., 2007; Gneiting et al., 2007).

250 Similarly, from a meteorological point of view, resolution measures the ability of a fore-  
251 casting system to produce predictive distributions conditioned by the value of the predictand  
252 (i.e. forecasts that are case-dependent) (Pinson et al., 2007). From a statistical point of  
253 view, resolution amounts to evaluate the capacity of the forecast system to produce different  
254 density forecasts depending on the forecast conditions (i.e. the predictive distributions are  
255 not only conditioned by the value of the predictand) (Pinson et al., 2007). For instance, the  
256 prediction intervals may exhibit increasing widths with increasing forecast horizon. Also,  
257 regarding the solar irradiance (GHI), the width of the PIs may vary according the sun’s  
258 position in the sky - see for the instance the work of (Grantham et al., 2016). In this work,  
259 we will not provide such a conditional assessment. Instead, we will propose a measure of  
260 resolution through the decomposition of the CRPS. From a meteorological perspective, it is  
261 also worth noting that, for perfectly reliable forecasts, sharpness is identical to resolution.  
262 In this work, we will clearly distinguish the definition of sharpness and resolution. That is to  
263 say, sharpness will refer to the concentration of the prediction intervals while resolution will  
264 quantify the ability of the forecasting system to generate conditional predictive distributions.  
265 Finally, it must be noted that reliability can be improved by means of statistical techniques  
266 also called calibration techniques (Gneiting et al., 2005), whereas this is not possible for  
267 resolution.



## 268 4. Verification tools

269 Diagnostic tools are used to visually assess the quality of probabilistic forecasts, while  
270 numerical scores are used to quantify the skills of a forecasting system and to rank competing  
271 prediction methods.

### 272 4.1. Diagnostic tools

273 The first requirement of reliability can be assessed with the help of reliability diagrams  
274 (see sub-section 4.1.1) when considering quantile forecasts. Rank histograms (see sub-section  
275 4.1.2) can be used for ensemble forecasts while PIT histograms (see sub-section 4.1.3) are  
276 best suited for the evaluation of quantile forecasts. These two graphical tools also permit to  
277 have an idea of the ensemble dispersion and bias.

278 In this work, and similarly to (Pinson et al., 2007; Gneiting et al., 2007), the assessment  
279 of sharpness derives from a more statistical point of view with focus on the shape of the  
280 predictive distributions. For that purpose, the average width of the prediction intervals (see  
281 sub-section 4.1.4) is used to evaluate the sharpness of the predictive distributions (Pinson  
282 et al., 2007).

#### 283 4.1.1. Reliability diagram

284 The reliability diagram is a graphical verification display used to evaluate the reliability  
285 of the probabilistic forecasts. In this paper, we follow the methodology defined by (Pinson  
286 et al., 2010) that is especially designed for predictive distributions summarized by quantile  
287 forecasts. More precisely, quantile forecasts are reliable if their nominal proportions are equal  
288 to the proportions of the observed value. It means that, over an evaluation set of significant  
289 size, (statistically) the difference between observed and nominal probabilities should be as  
290 small as possible (Pinson et al., 2010). Notice that for ensemble forecasts, the uniform CDF  
291 or non uniform CDF (see section 2) must be chosen before applying this methodology.

292 This representation is attractive since the deviations from perfect reliability (i.e. the  
293 diagonal line) can be easily visualized (Pinson et al., 2010). Nonetheless, due to the finite  
294 sample of pairs of observation/forecast and also due to possibly serial correlation in the  
295 sequence of forecasts and observations, it is possible that observed proportions are not  
296 exactly along the diagonal, even if the forecasts are perfectly reliable. (Pinson et al., 2010).  
297 In other words, reliability diagrams can be misinterpreted since even for perfectly reliable  
298 forecasts, deviations from the ideal diagonal case can be observed.

299 To deal with the issue of limited number of pairs of observation/forecast, Bröcker and  
300 Smith (2007a) built reliability diagrams with consistency bars. In addition, Pinson et al.  
301 (2010) have proposed consistency bars taking into account the combined effect of serial cor-  
302 relation and limited data. Interpretation of reliability diagrams with consistency bars is  
303 that one cannot reject the hypothesis of the quantile forecasts being reliable if the observed  
304 proportions lie within the consistency bars. In practice, adding consistency bars to the rela-  
305 bility diagrams may reinforce the user's (possibly subjective) judgment about the reliability  
306 of the different models.

307 Finally, some preceding works (Chu and Coimbra, 2017; Lauret et al., 2017) proposed  
308 to evaluate the reliability component of a probabilistic system by calculating the prediction

309 interval coverage probability (PICP) (Khosravi et al., 2013). PICP permits one to assess  
310 the empirical coverage probability of the central prediction intervals. However, this metric  
311 is not suitable to assess the reliability of probabilistic forecasts because as noted by Pinson  
312 et al. (2007), both quantiles that define the prediction interval may be biased. In other  
313 words, PICP it is not sufficient to check if the nominal coverage of the intervals is respected.  
314 It is also necessary to verify that both quantiles defining the PI are unbiased.

#### 315 *4.1.2. Rank histogram*

316 The rank histogram is a graphical display initially designed for assessing ensemble fore-  
317 casts (Wilks, 2014). But, it can be extended to quantile forecasts by assuming that all  
318 evenly spaced forecasted quantiles form an ensemble. Rank histograms permit to assess the  
319 statistical consistency of the ensemble, that is, if the observation can be seen statistically  
320 just like another member of the ensemble (Wilks, 2014). A flat rank histogram is a neces-  
321 sary condition for ensemble consistency and shows an appropriate degree of dispersion of  
322 the ensemble. Put differently, the flatness of the rank histogram indicates that the ensemble  
323 members are statistically indistinguishable from the observations (Wilks, 2014). An under-  
324 dispersed ensemble (i.e. ensemble dispersion consistently too small) leads to a U-shape rank  
325 histogram and shows that the observation will often be an outlier in the distribution of  
326 ensemble members. Conversely, an over-dispersed ensemble (i.e. ensemble dispersion con-  
327 sistenty too large) gives a hump shape rank histogram and indicates that the observation  
328 may too often be in the middle of the ensemble distribution.

329 In addition, rank histograms can also detect deficiencies in ensemble calibration or relia-  
330 bility (Wilks, 2014). For instance, some unconditional biases can be revealed by asymmetric  
331 (triangle shape) rank histograms. Furthermore, overpopulation of the smallest (resp. high-  
332 est) ranks will correspond to an overforecasting (resp. underforecasting) bias. It must be  
333 stressed that one should be cautious when analyzing rank histograms. Indeed, as shown by  
334 (Hamill, 2001), a perfect flat rank histogram does not state that the corresponding forecast  
335 is reliable. Further, when the number of observations is limited, consistency bars can also  
336 be calculated with the procedure proposed by (Bröcker and Smith, 2007a). To build a rank  
337 histogram, it is necessary to find the rank of the observations when pooled within the or-  
338 dered ensemble and then plot the histogram of the ranks. For an ensemble of  $M$  members,  
339 the number of ranks of the histogram is  $M + 1$ . The histogram of verification ranks will be  
340 uniform with theoretical relative frequency of  $\frac{1}{M+1}$  if the consistency condition is met.

#### 341 *4.1.3. PIT histogram*

342 PIT histograms may help to assess the calibration property by verifying whether the  
343 observations can be seen as random samples of the predictive distributions (Gneiting et al.,  
344 2007). PIT histograms assess calibration of cumulative predictive distributions checking  
345 whether the observations can be considered as random samples of these distributions. Con-  
346 trary to rank histograms, PIT histograms require the computation of the predictive CDF.  
347 The PIT is the value that the predictive CDF has for a particular observation. PIT values  
348 can be calculated over a testing set of observations and one can then plot the histogram of  
349 the PIT values. Similarly to rank histograms, a flat PIT histogram is a necessary but not

350 sufficient condition to state that a forecast is reliable. As for rank histograms, departures  
 351 from flatness is a sign of conditional biases in the forecasts or over/under-dispersion.

#### 352 4.1.4. Sharpness diagram

353 A probabilistic forecast is sharp if prediction intervals are shorter on average than pre-  
 354 diction intervals derived from naïve methods, such as climatology or persistence.

355 Similarly to Pinson et al. (2007), we propose to assess the sharpness of the predictive  
 356 distributions by calculating the mean size of the central prediction interval denoted by  $\bar{\delta}^\alpha$   
 357 for different nominal coverage rates  $(1 - \alpha)\%$ . This leads to a graphical verification display  
 358 called  $\delta$ -diagrams. For an evaluation set of  $N$  forecasts,  $\bar{\delta}^\alpha$  is given by

$$\bar{\delta}^\alpha = \frac{1}{N} \sum_{i=1}^N (\hat{q}_{\tau=1-\alpha/2} - \hat{q}_{\tau=\alpha/2}). \quad (4)$$

359 Notice that Gneiting et al. (2007) proposed a diagnostic approach to evaluating proba-  
 360 bilistic forecasts that is based on the paradigm of maximizing the sharpness of the predictive  
 361 distributions subject to calibration. In the proposed evaluation framework, sharpness dia-  
 362 grams take the form of box-plots of the width of the prediction intervals.

363 As mentioned above, some researchers in the solar forecasting community used the  
 364 PINAW metric to measure sharpness. This metric is the average width of the  $(1 - \alpha)100\%$   
 365 prediction interval normalized by the mean of variable  $x$  to predict (e.g. here GHI) for  
 366 a testing set of  $N$  pairs of forecasts/observations. For a specific nominal coverage rate  
 367  $(1 - \alpha)100\%$ , PINAW reads as

$$\text{PINAW}(\alpha) = \frac{\sum_{i=1}^N (\hat{q}_{\tau=1-\alpha/2} - \hat{q}_{\tau=\alpha/2})}{\sum_{i=1}^N x}. \quad (5)$$

368 However, even if it can be interesting to compare the performance of forecasting methods  
 369 at different locations, it must stressed that the sharpness is a property of the forecasts only  
 370 and as such can not depend on the mean of the observations.

#### 371 4.2. Scores

372 Numerical scores provide summary measures for the evaluation of the quality of proba-  
 373 bilistic forecasts (Gneiting and Raftery, 2007). Scoring rules are based on the predictive  
 374 distribution of the forecast and on the observed value of the variable of interest. Scores  
 375 may help to rank competing probabilistic models. Scores are required to be proper (Bröcker  
 376 and Smith, 2007b; Gneiting and Raftery, 2007). A score is said to be proper if it insures  
 377 that the perfect forecasts should be given the best score value. If it is not the case, one  
 378 could then hedge the score, by finding tricks that permit to get better score values without  
 379 attempting to issue better forecasts. More generally, employing a score that is not proper  
 380 makes that one can never be sure of the validity of the results from an empirical comparison  
 381 or benchmarking of rival approaches (Pinson and Tastu, 2014). The scoring rules proposed  
 382 in this work (CRPS, Ignorance score, Interval score, quantile score) are proper. However,

383 this is not the case of the CWC score discussed in section 1 as demonstrated by (Pinson and  
 384 Tastu, 2014).

385 In addition to the property of propriety, a score can be local or non-local. A score is said  
 386 to be local if it depends only on the value of the predictive distribution at the observation,  
 387 not on other features of the functional form of the predictive PDF.

388 While different proper scores have been proposed in the literature (Bröcker and Smith,  
 389 2007b; Gneiting and Raftery, 2007), we focus here on proper scoring rules for probabilistic  
 390 forecasts of continuous variables and particularly on the following scores: Interval score,  
 391 quantile score, CRPS, and Ignorance Score.

#### 392 4.2.1. Interval Score (IS)

393 Following Winkler (1972), Gneiting and Raftery (2007) proposed a proper score to specif-  
 394 ically assess the quality of central  $(1 - \alpha)100\%$  prediction interval forecasts. This scoring  
 395 rule called Interval Score (IS), averaged over the  $N$  pairs of forecasts and observations, is  
 396 defined by

$$IS_\alpha = \frac{1}{N} \sum_{i=1}^N \left( U^i - L^i \right) + \frac{2}{\alpha} \left( L^i - x_{obs}^i \right) 1_{x_{obs}^i < L^i} + \frac{2}{\alpha} \left( x_{obs}^i - U^i \right) 1_{x_{obs}^i > U^i}, \quad (6)$$

397 where  $L^i$  and  $U^i$  represent respectively the  $\alpha/2$  lower quantile  $\hat{q}_{\tau=\alpha/2}$  and the  $1 - \alpha/2$  upper  
 398 quantile  $\hat{q}_{\tau=1-\alpha/2}$ . As shown by Equation 6, the IS rewards narrow prediction intervals but  
 399 penalizes (with the penalty term that depends on  $\alpha$ ) the forecasts for which the observation  
 400  $x_{obs}$  is outside the interval.

#### 401 4.2.2. Quantile Score (QS)

402 Some users may be interested by the performance of some specific quantiles ( e.g. over-  
 403 forecasting or underforecasting) and particularly those related to the tails of the predictive  
 404 distribution. Quantile Score (QS) permits to obtain detailed information about the fore-  
 405 cast quality at specific probability levels. As noted by (Bentzien and Friederichs, 2014), the  
 406 CRPS averages over the complete range of forecast thresholds through integration of the  
 407 Brier Score (see Appendix C). As a consequence, deficiencies in different parts of the distri-  
 408 bution, e.g. the tails of the distribution, might be hidden. Bentzien and Friederichs (2014)  
 409 recommend to extend the verification framework by calculating QS for different probability  
 410 levels. Notice also that, Bentzien and Friederichs (2014) proposed a decomposition of the  
 411 QS into its reliability and resolution components.

412 QS is based on an asymmetric piecewise linear function  $\psi_\tau$  called the check or pinball loss  
 413 function. The check function was first defined in the context of quantile regression (Koenker  
 414 and Bassett, 1978) and is given by

$$\psi_\tau(u) = \begin{cases} \tau u & \text{if } u \geq 0 \\ (\tau - 1)u & \text{if } u < 0, \end{cases} \quad (7)$$

415 with  $\tau$  representing the quantile probability level.

416 QS is given by the mean of the check function applied to the  $N$  pairs of observations  $x_{obs}^i$   
 417 and quantile forecasts for a specific probability level  $\tau$ ,  $\hat{q}_\tau^i$ . QS reads as

$$QS = \frac{1}{N} \sum_{i=1}^N \psi_\tau(x_{obs}^i - \hat{q}_\tau^i). \quad (8)$$

418 Finally, notice that Bröcker (2012) showed that the CRPS can be seen as a weighted sum  
 419 of quantiles scores applied to the quantiles derived from the non-uniform CDF.

#### 420 4.2.3. Continuous Rank Probability Score (CRPS) and its decomposition

421 The CRPS measures the difference between the predicted and observed cumulative dis-  
 422 tributions functions (CDF) (Hersbach, 2000). The formulation of the CRPS is

$$CRPS = \frac{1}{N} \sum_{i=1}^N \int_{-\infty}^{+\infty} [\hat{F}_{fcst}^i(x) - F_{x_{obs}}^i(x)]^2 dx, \quad (9)$$

423 where  $\hat{F}_{fcst}(x)$  is the predictive CDF of the variable of interest  $x$  (e.g. GHI) and  $F_{x_{obs}}(x)$  is a  
 424 cumulative-probability step function that jumps from 0 to 1 at the point where the forecast  
 425 variable  $x$  equals the observation  $x_{obs}$  (i.e.  $F_{x_{obs}}(x) = 1_{\{x \geq x_{obs}\}}$ ). The squared difference  
 426 between the two CDFs is averaged over the  $N$  forecast/observation pairs. The CRPS score  
 427 rewards concentration of probability around the step function located at the observed value  
 428 (Wilks, 2014). In other words, the CRPS penalizes lack of resolution of the predictive  
 429 distributions as well as biased forecasts. In addition, for deterministic forecasts, the CRPS  
 430 turns to be the MAE (Mean Absolute Error). This fact permits to compare directly the  
 431 performance of a probabilistic model against a deterministic one or equivalently evaluate  
 432 the added value brought by a probabilistic approach (Ben Bouallègue, 2015). Notice that  
 433 the CRPS is negatively oriented (smaller values are better) and the same dimension as the  
 434 forecasted variable.

435 For ensemble forecasts, Hersbach (2000) proposed a method to compute the CRPS using  
 436 the classical definition of the CDF (see section 2 and figure 2(a)). In the realm of weather  
 437 predictions, his method is widely used and at least embedded in one R-package (NCAR-  
 438 Research applications laboratory, 2015). Appendix B summarizes the Hersbach's method  
 439 to compute the CRPS for ensemble forecasts.

440 As mentioned above and as a proper score (Gneiting and Raftery, 2007), CRPS can be  
 441 further partitioned into the two main attributes of probabilistic forecasts namely reliability  
 442 and resolution. The decomposition of the CRPS leads to

$$CRPS = RELIABILITY + UNCERTAINTY - RESOLUTION. \quad (10)$$

443 The reliability term provides an estimation of the forecast biases while the resolution  
 444 term quantifies the improvement that results from issuing probability forecasts that are case  
 445 dependent. The uncertainty term cannot be modified by the forecast system and depends  
 446 only on the observations variability (Wilks, 2014). As the CRPS is negatively oriented, the

447 goal of a forecast system is to minimize (resp. maximize) as much as possible the reliability  
 448 term (resp. the resolution term). This decomposition of the CRPS may lead to a detailed  
 449 picture of the performance of the forecasting methods.

450 Regarding the calculation of these different terms, two possibilities exist. The first one  
 451 is based on the work of (Hersbach, 2000) and as such best suited for ensemble forecasts rep-  
 452 resented by the classical definition of the CDF. Appendix B gives the formulaes to calculate  
 453 the three terms. The second possibility makes use of the fact that CRPS is the integral of  
 454 the Brier Score over all the predictand thresholds. The Brier score is a proper score used  
 455 to evaluate probabilistic forecasts of binary predictands (Wilks, 2014). Appendix C gives  
 456 all the details regarding this second method. As the CRPS has the same unit as the vari-  
 457 able to predict, it can be normalized by the mean (e.g. mean GHI) or the maximum (e.g.  
 458 installed capacity) of the variable to forecast. The normalized CRPS permits to carry out  
 459 comparisons between different datasets (e.g. different locations).

460 We close this subsection related to the CRPS with the CRPS skill score (CRPSS). In a  
 461 similar manner that scores have been proposed to evaluate the skill of deterministic forecasts  
 462 (Coimbra et al., 2013), (Pedro et al., 2018) used the CRPSS to gauge the performance of  
 463 their probabilistic forecasting models against a reference easy-to-implement method i.e. the  
 464 persistence ensemble (PeEn). In that case, the CRPSS reads as  $CRPSS = 1 - \frac{CRPS_{new\ method}}{CRPS_{PeEn}}$ .

465 In this study, as our primary goal is to verify solar irradiance probabilistic forecasts and  
 466 not to compare and rank forecasting models, we do not detail the implementation of the  
 467 PeEn model. The interested reader should refer to (Pedro et al., 2018). However, as noted  
 468 by (Yang, 2019), the previous definition of the CRPSS may lead to some misinterpretations  
 469 of the skill score as the CRPS of the PeEn model varies according to certain parameters  
 470 (e.g. number of members of the ensemble, forecast lead time, etc.). To address this issue,  
 471 Yang (2019) proposed, instead of PeEn, a new baseline model called the complete-history  
 472 PeEn (CHPeEn) model that gives a nearly constant CRPS.

473 Another way to avoid a CRPSS that depends on the implementation of the reference  
 474 model, and to benefit from the decomposition of the CRPS mentioned above, is to use the  
 475 uncertainty part of the CRPS as the baseline value. The uncertainty component corresponds  
 476 to the CRPS of the climatology and is only sensitive to the observations variability and  
 477 therefore, for a given location and temporal resolution of the data, does not depend on any  
 478 other kind of parameters. Notice that, for meteorologists, when computing skill scores, the  
 479 baseline model is commonly climatology.

#### 480 4.2.4. Ignorance Score (IGN)

481 Initially proposed by (Good, 1952), this score is cited under various names: log score  
 482 (Gneiting and Raftery, 2007), divergence (Weijs et al., 2010) or ignorance score (Roulston  
 483 and Smith, 2002). Considering  $N$  verification pairs of probabilistic forecasts given by their  
 484 PDF  $\hat{f}^i(x)$  and outcomes  $x_{obs}^i$ , the ignorance (IGN) is defined as follow

$$IGN = -\frac{1}{N} \sum_{i=1}^N \log(\hat{f}^i(x_{obs}^i)). \quad (11)$$

485 This strictly proper score is appealing because it gathers interesting properties like ad-  
486 ditivity and locality (i.e. the score depends “only on the value of the probabilistic forecast  
487 at the verification” (Bröcker and Smith, 2007b)). Like the CRPS, the IGN is a negatively  
488 oriented score (smaller values are better). Based on the log function, this score is strongly  
489 affected by the large errors, when the observations fall far away from the highest forecasted  
490 probabilities. Equation 11 provides a simple way to compute the ignorance score from con-  
491 tinuous PDFs of parametric distributions or from predictive distributions (i.e. derived from  
492 discrete estimates, see section 2).

493 Considering ensemble forecasts, Roulston and Smith (2002) proposed a simple approach  
494 to compute the IGN. They used the “uniform” definition of the CDF derived from an  
495 ensemble forecast (see section 2 and figure 2(c)) combined with a linear interpolation of  
496 the probabilities between two consecutive members. Then, they applied Equation 11 to the  
497 corresponding PDF that is the first derivative of the CDF (see appendix A for more details).  
498 Thus, the ignorance score of an outcome  $x_{obs}$  that lies between two consecutive members  
499  $[e_k; e_{k+1}]$  of an ensemble forecast with  $M$  members is given by Equation 12. We propose here  
500 a slightly different formulation of the IGN defined in the article of (Roulston and Smith,  
501 2002). They defined the IGN using the binary logarithm (or log base 2) classically proposed  
502 by the field of information theory. We prefer here to use the common logarithm function (or  
503 log base 10) to coincide with the general framework of the IGN (see Equation 11) mainly  
504 used in the literature. For ensemble forecasts, IGN is given by

$$IGN = \log(M + 1) + \log \Delta X_k, \quad (12)$$

505 where

$$\begin{aligned} \Delta X_k &= e_{k+1} - e_k \text{ if } 1 < k < M \\ \Delta X_0 &= e_1 - e_0 \\ \Delta X_M &= e_{M+1} - e_M. \end{aligned} \quad (13)$$

506  $[e_0; e_{M+1}]$  is the a priori interval on which the outcome  $x_{obs}$  is expected to be. Roulston  
507 and Smith (2002) proposed to use the minimum and the maximum of the climatology as  
508 boundaries of this interval. One can notice that this formulation of the IGN assigns the  
509 highest probabilities to the smallest differences between consecutive members. For a verifi-  
510 cation dataset of  $N$  forecast-realization pairs, the ignorance score corresponds obviously to  
511 the arithmetical mean as in Equation 11. Notice that, unlike the CRPS, the ignorance score  
512 cannot be decomposed into reliability, resolution and uncertainty.

513 Finally, Tödter and Ahrens (2012) proposed a generalization of the IGN with an approach  
514 similar to Hersbach’s work (Hersbach, 2000) about the CRPS. They introduced a non-local  
515 version of the IGN for binary events and a new score called the Continuous Ranked Ignorance  
516 score (CRIGN) by analogy to the CRPS. For ensemble forecasts, no clear definition of the  
517 CDF to use to compute these non-local scores is provided. Thus, the CRIGN will not be  
518 addressed in this work.

Table 1: Diagnostic tools for reliability assessment

Diagnostic tool	Designed for	Pros	Cons	Remarks
Reliability diagram	Quantile forecasts	Departure from perfect reliability easily assessed Easy to build	Sensitivity to the finiteness of the data and serial correlation	Can be used for Ensemble if members are assigned specific probability levels (uniform/non uniform CDF)
Rank histogram	Ensemble forecasts	Lack of reliability and ensemble dispersion easily assessed Easy to build	Sensitivity to the finiteness of the data	Can be extended to quantile forecasts if quantiles are evenly spaced
PIT histogram	Quantile forecasts	Departure from perfect reliability easily assessed	- Need to specify the number of histograms bins. - Interpolation needed between the discrete quantiles to estimate the value the CDF attains at the observation.	Can be used for Ensemble (uniform CDF)

519 *4.3. Proposed evaluation framework*

520 Tables 1 and 2 summarize the diagnostic tools and scoring rules used to evaluate prob-  
521 abilistic forecasts generated either by ensemble methods or quantile techniques. Regarding  
522 pros and cons, and also the most common approaches already used in other fields (i.e.  
523 weather forecast verification and wind power forecasting), we propose to differentiate the  
524 methodologies and the tools to assess the quality of quantile forecasts and ensemble predic-  
525 tion systems (EPS).

526 Considering quantile forecasts, we advise to visually assess the quality of the forecasts  
527 using reliability diagrams with consistency bars. Then, to use the CRPS and its related  
528 decomposition as described in appendix C to quantify the overall performance of the methods  
529 and to measure the reliability and the resolution components.

530 For ensemble forecasts, we propose to use the rank histogram including consistency bars  
531 and the CRPS as defined by (Hersbach, 2000) (see appendix B) to respectively qualify and  
532 quantify the performances of the EPS. Indeed, these two tools does not require additional  
533 assumptions (i.e. to define the nature of the distribution and its boundaries) and they are  
534 already widely used.

535 For both type of forecasts, ignorance score, interval score, quantile score and sharpness  
536 diagrams can complement the characterization of the forecasting methods. However, sharp-  
537 ness diagrams must be interpreted with care because they are only relevant if the associated  
538 forecasts are reliable.

539 Finally, if interval score, quantile score and sharpness diagrams are computed for ensem-  
540 ble forecasts, it is important to clearly indicate the assumption done to obtain the quantiles  
541 (e.g. uniform or non-uniform spacing).

542 **5. Applications**

543 To illustrate the use of the different diagnostic tools and scores presented above, we will  
544 present in this section examples relative to two locations with different sky conditions. The  
545 first site, Desert Rock (USA), has an arid climate with a very sunny and stable sky. The  
546 second site, Tampon (Réunion island), is located in a tropical island and experiences a very  
547 variable sky. The experimental dataset corresponds to two consecutive years of recorded



Table 2: Scoring Rules

Scores	Pros	Cons	Remarks
CRPS	Same dimension as the variable to predict Can be normalized. Can be compared with MAE. Can be decomposed in reliability and resolution	No analytic formulae except for specific PDF (Gaussians, Student's t,...)- see R package scoringRules for details.	Specific formulae for Ensemble forecasts proposed by Herbasch (see Appendix A). Can be calculated through numerical integration (see Equation 9). Needs interpolation of uniform/non uniform CDF. Can be also computed through integration of the Brier Score (see Appendix C)
Ignorance Score	Easy to compute	Sensitive to the form of the PDF Cannot be decomposed into reliability and Resolution Cannot be normalized	Specific formula for Ensemble forecasts proposed by Roulston (see Equation 12). Otherwise, estimation of the value the PDF attains at the observation. Needs interpolation of uniform/non uniform CDF. Cannot be applied to predictive distributions with null probabilities
Quantile Score	Can be decomposed into reliability and resolution		
Interval Score	Very easy to compute Same Dimension as the variable to predict Can be normalized	Cannot be decomposed into reliability and Resolution	

Table 3: Main characteristic of the solar measurements

	<b>Desert Rock (USA)</b>	<b>Tampon (Réunion)</b>
Provider	SURFRAD	PIMENT
Position	36.6N, 116.0W	21.3S, 55.5E
Elevation	1007m	550m
Cimate type	Arid	Insular tropic
Period of record	2012-2013	2012-2013
Annual solar irradiation	2.105 MWh/m <sup>2</sup>	1.712 MWh/m <sup>2</sup>
Solar variability 1-h ( $\sigma\Delta kt^*$ )	0.146	0.241
Mean GHI (Testing set)	548 W/m <sup>2</sup>	458 W/m <sup>2</sup>
Uncertainty component of the CRPS	29.1%	33.1%

548 data of global horizontal irradiance (GHI). Table 3 gives detailed information about the  
549 data. The solar variability, quantified by the standard deviation of the changes in the clear  
550 sky index  $\sigma\Delta kt^*$  (Hoff and Perez, 2012), is the main difference between the two considered  
551 locations. We intentionally chose these two sites. Indeed, the solar variability is a key factor  
552 in the accuracy of deterministic forecasts. The higher the variability, the less accurate the  
553 forecasts are (Lauret et al., 2015). Finally, to build some of the models used in this work,  
554 we used the first year of data (2012) as training set and the second year of data (2013) as  
555 testing set. Therefore, all the metrics and visual tools presented hereafter are derived from  
556 the testing set.

557 Two forecasting time horizons will be addressed in this work. First, intra-day forecasts  
558 with lead times ranging from 1 to 6 hours will be appraised. These forecast are provided  
559 by state of the art forecasting models that generate predictive distributions from a set  
560 of quantiles spanning the unit interval. Second, day-ahead probabilistic forecasts will be  
561 studied. Generated by Numerical Weather Predictions (NWP) models, they are provided  
562 as ensemble forecasts.

563 *5.1. Intraday probabilistic forecasts*

564 Regarding intraday forecasts, the quality of four state-of-the-art probabilistic models  
 565 will be appraised. In this paper, we will not give the details of the implementation of these  
 566 models as they have already been described in previous works (David et al., 2018; Pedro  
 567 et al., 2018). In addition, we recall that the goal here is to illustrate the application of the  
 568 proposed evaluation framework and not to have a detailed evaluation of these models.

569 The selected models are based on two quantile regression techniques namely the quantile  
 570 regression forest (QRF) and the Gradient Boosting (GB) techniques.

571 Briefly, the proposed techniques estimate directly the set of quantiles from a regression  
 572 model  $Y = f(X)$  that relates the response variable  $Y$  (here GHI for lead time  $h = 1, 2, \dots, 6$   
 573 hours) to a set of predictor variables ( $X$ ). Two variants of regression models with different  
 574 sets of predictor variables are built. For the first variant described in (Lauret et al., 2017),  
 575 the vector of explanatory variables  $X$  consists of the actual measurement plus five past  
 576 ground measurements while the second one takes as additional inputs two geometrical solar  
 577 features related to the course of the sun in the sky namely the cosine of the zenith angle  
 578 ( $\cos(SZA)$ ) and the cosine of the hour angle ( $\cos(HA)$ ). The adding of the two variables  
 579 originates from the following reasons. First, some authors (Grantham et al., 2016; Lorenz  
 580 and Heinemann, 2012) showed a clear dependency of the forecasting error in relation to SZA.  
 581 Second, we expect that the hour angle will bring some information regarding the asymmetry  
 582 of the sky conditions between mornings and afternoons. This may be hold particularly for  
 583 site like Le Tampon that experiences such a dichotomy between mornings and afternoons.  
 584 Table 4 lists the acronyms of the resulting four quantile regression models.

Table 4: Acronyms related to the four quantile regression models

<b>Quantile regression techniques</b>	<b>Variant 1</b>	<b>Variant 2</b>
Quantile Regression Forest	QRF1	QRF2
Gradient Boosting	GB1	GB2

585 *5.2. Application of the diagnostic tools and scores to quantile forecasts*

586 We follow the evaluation framework proposed by (Pinson et al., 2007). More precisely,  
 587 we will first evaluate the reliability attribute by analysing the reliability diagrams. Then,  
 588 the sharpness property will be assessed with  $\bar{\delta}^\alpha$  diagrams. Notice that, although being at  
 589 this stage redundant with the reliability diagram, we also plot the PIT diagrams in order to  
 590 discuss possible issues related with the use of this graphical tool.

591 The CRPS and its related decomposition described in Appendix C are proposed to assess  
 592 the skills of the probabilistic models. We complement this analysis with the ignorance score  
 593 and the interval score. Detailed information about the performance of the probabilistic  
 594 models at specific probability levels will be given by the quantile score (QS). It must noted  
 595 that, in the following, the different figures plot the relative counterparts of the CRPS,  
 596 Interval Score and QS metrics. These relative metrics are normalized by dividing the absolute  
 597 values by the mean of the GHI for the considered testing period (see Table 3).

598 *5.2.1. Reliability assessment*

599 As mentioned above, reliability diagrams allow to graphically assessing the reliability of  
600 a set of quantile forecasts. Figures 3(a) and 3(b) plot the reliability diagrams (averaged over  
601 all the forecasting horizons) for the two selected sites. Consistency bars for a 90% confidence  
602 level are individually computed for each nominal proportion.

603 From the visual inspection of the reliability diagrams of Desert Rock, one can possibly  
604 state that the GB2 model is reliable as the observed proportions of all quantiles lie within the  
605 consistency bars. Conversely, for the others models, observed proportions of some quantiles  
606 lie outside the consistency bars. In particular, quantile forecasts generated by the QRF2  
607 model should not be considered reliable. In addition, notice the particular signature of  
608 the QRF2 model that corresponds to an over dispersed predictive distribution (i.e. an  
609 underconfident model).

610 For the site of Le Tampon, it seems that, except the GB2 model, all the other models  
611 lead to possible reliable quantile forecasts since all of their observed proportions lie within  
612 the consistency bars. Moreover, notice that the QRF1 model exhibits a reliability curve very  
613 close to the ideal diagonal case.

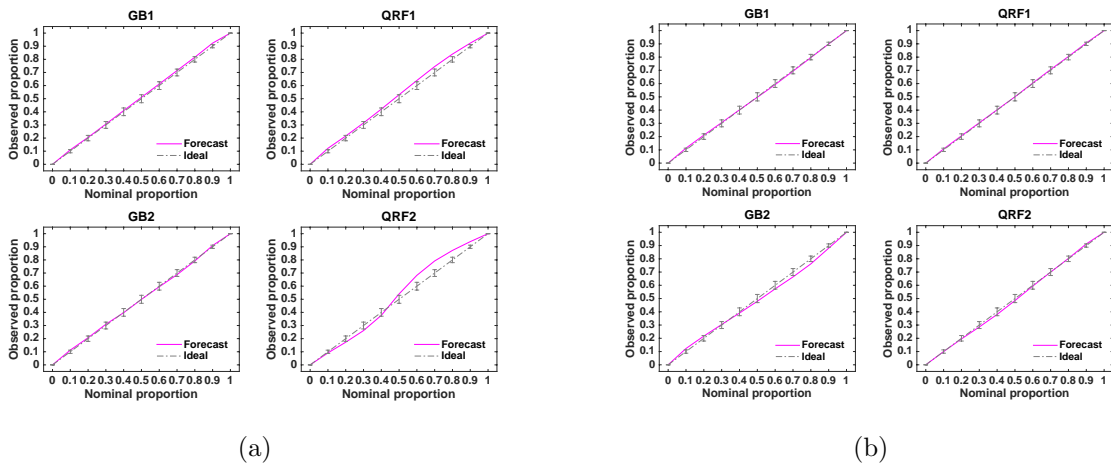


Figure 3: Reliability diagrams (a) Site of Desert Rock (b) Site of Le Tampon.

614 Figure 4 shows the PIT diagrams (averaged over all the lead times) related to the two  
615 sites. Following the preceding reliability analysis which possibly stated that the QRF1 model  
616 was reliable for the site of Le Tampon (see Figure 3(b)), one may expect a corresponding  
617 flat PIT histogram (Figure 4(b)). However, this is not the case. We suspect that this may  
618 come from the fact that one needs to specify the number of histograms bins to plot the PIT  
619 histogram. In addition, interpolation is needed between the discrete quantiles to estimate  
620 the value the CDF attains at the observation. This may motivate the choice of reliability  
621 diagrams against PIT histograms for assessing calibration. However, it is worth noting that,  
622 in accordance with the reliability diagram, the PIT histogram of the QRF2 method for  
623 Desert Rock confirms that this model corresponds to an over-dispersed forecasting system  
624 (i.e. too wide predictive distributions).

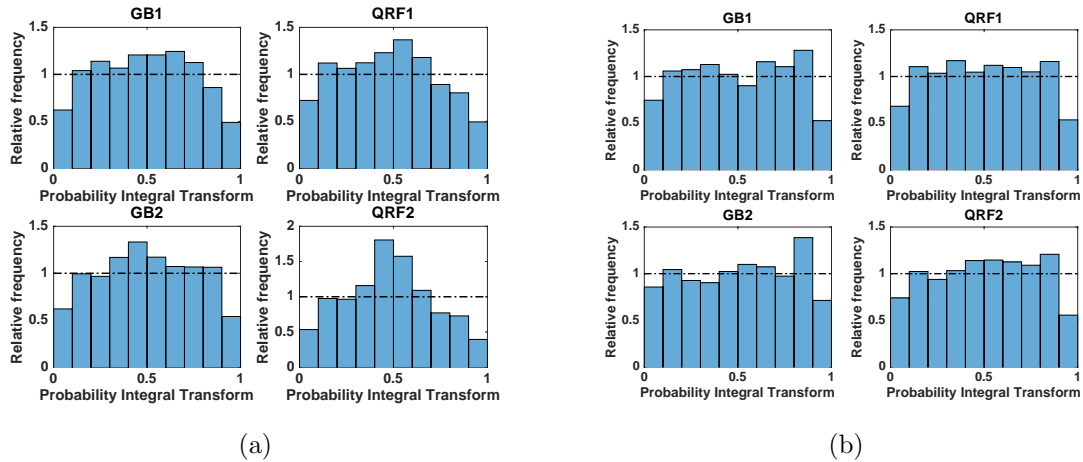


Figure 4: PIT diagrams (a) Site of Desert Rock (b) Site of Le Tampon.

### 5.2.2. Sharpness assessment

Figures 5(a) and 5(b) plot the  $\bar{\delta}^\alpha$  diagrams of the four models for different coverage rates. In this work, it must be noted that the  $\bar{\delta}^\alpha$  values have been averaged over all the lead times. One may first notice that prediction intervals are wider for the site of Le Tampon than for Desert Rock. As discussed in (Lauret et al., 2017), the variable sky conditions experienced by the site of Le Tampon have an impact on the shape of the predictive distributions. Conversely, the site of Desert Rock that experiences higher occurrences of clear and stable skies exhibits narrower prediction intervals.

For both sites, it appears the GB2 model leads to the lowest  $\bar{\delta}^\alpha$  values for all the forecasting horizons albeit the difference with the other models is less pronounced for the site of Desert Rock. At this point, the sharpness evaluation may favor the GB2 model for both sites. However, while the GB2 model may possibly generate reliable forecasts for the Desert Rock site, this may not be the case for Le Tampon site. If one attempts to select the best approach for both sites by combining the two previous separate reliability and sharpness assessments, the picture is less clear. Hence evaluating separately reliability and sharpness and drawing conclusions on the sole examination of either one of these diagnostic tools may be misleading. In the next section, we will use the CRPS and its related decomposition into reliability and resolution in an attempt to assess objectively and quantitatively the properties required for a skillfull probabilistic system.

### 5.2.3. CRPS and its decomposition into resolution and reliability

This section proposes a detailed picture of the performance of the probabilistic models by plotting the CRPS and its associated decomposition into reliability and resolution. Notice that the uncertainty part is given in Table 3. Figures 6(a) and 6(b) plot the relative CRPS in relation with the forecast horizon for the two considered sites.

As expected, the performance of the models decreases as the lead-time increases (i.e. the lower the CRPS, the better the model). One also may note that the site of Le Tampon, which experiences variable sky conditions compared to Desert Rock, yields higher CRPS

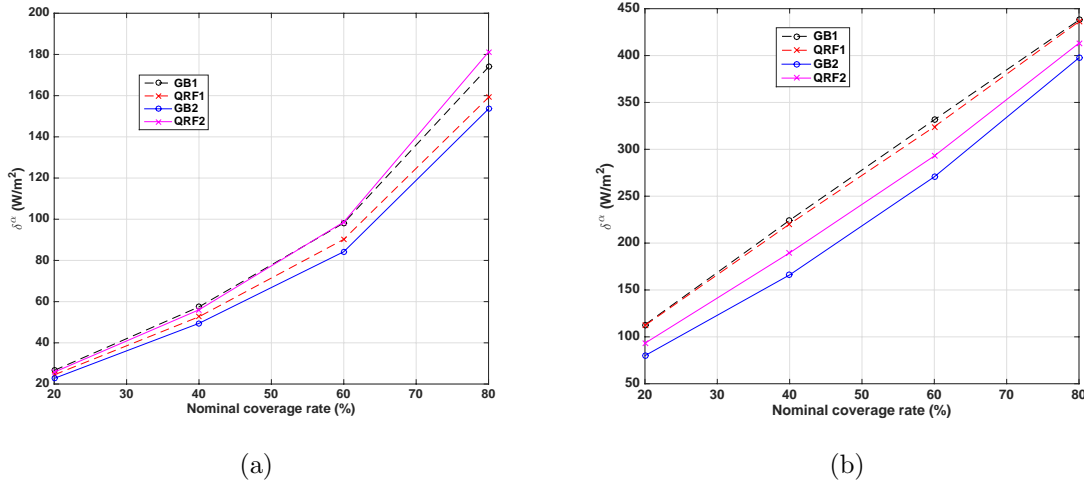


Figure 5: Sharpness diagrams (a) Site of Desert Rock (b) Site of Le Tampon

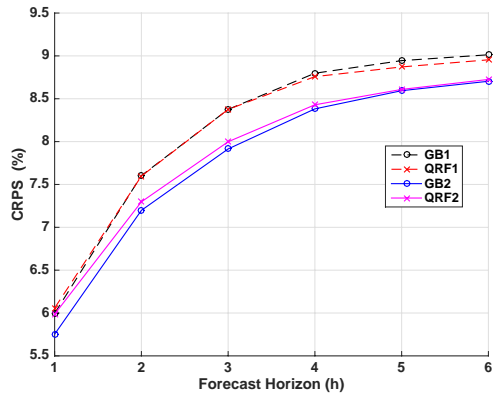
652 values. The interested reader is referred to (Lauret et al., 2017) where more details are  
653 given regarding the impact of the sky conditions on the quality of the probabilistic forecasts.  
654 As shown by Figures 6(a) and 6(b), the two non linear models that include the two solar  
655 geometric predictors namely zenith angle and hour angle (i.e. GB2 and QRF2 models)  
656 perform clearly better than the variant 1 models regardless the site. Thus, it appears that  
657 adding the two solar geometric variables brings a clear improvement and especially for a site  
658 like Le Tampon which is known to experience a morning/afternoon sky asymmetry. Unlike  
659 the previous separate analysis of reliability and sharpness, CRPS establishes a clear-cut  
660 ranking of the models. However, some inconsistencies appear with the reliability analysis  
661 which showed that the the QRF2 model (resp. the GB2 model) was non reliable for Desert  
662 Rock (resp. for Le Tampon). Therefore, in order to gain a better understanding of the CRPS  
663 results, we use the decomposition of the CRPS depicted in Appendix C. This decomposition,  
664 detailed in Appendix D, shows that the reliability component makes a small contribution  
665 to the CRPS and that the higher quality of the variant 2 models comes from the resolution  
666 attribute.

#### 667 5.2.4. Ignorance score (IGN)

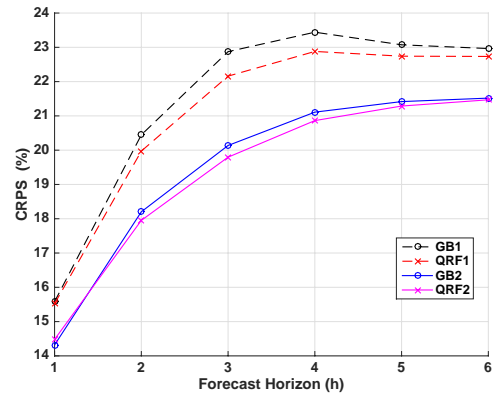
668 Figure 7 plots the ignorance score of the four models. This scoring rule confirms the su-  
669 periority of the variant 2 models although the QRF2 model appears to be the best performer.  
670 For this particular application, the ignorance score can complement the CRPS analysis and  
671 may increase the user’s confidence to select the QRF2 method.

#### 672 5.2.5. Interval score (IS)

673 Figure 8 shows the IS score for the 80% central prediction interval. Again, variant 2  
674 models perform better than the other models. In our opinion, this easy-to-calculate score  
675 can advantageously complete the set of proper scores available to the user.

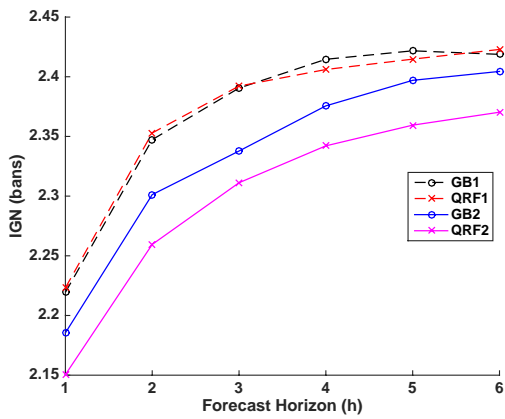


(a)

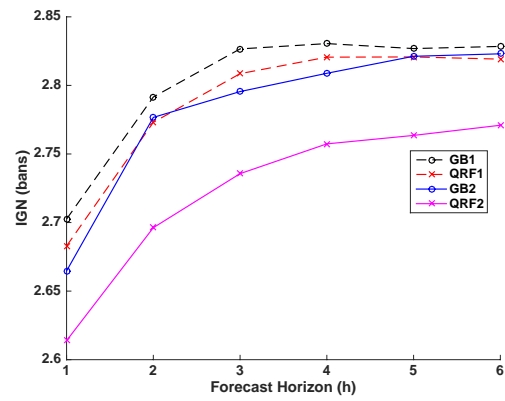


(b)

Figure 6: CRPS of the different intraday methods (a) Site of Desert Rock (b) Site of Le Tampon



(a)



(b)

Figure 7: Ignorance Score of the different intraday methods (a) Site of Desert Rock (b) Site of Le Tampon

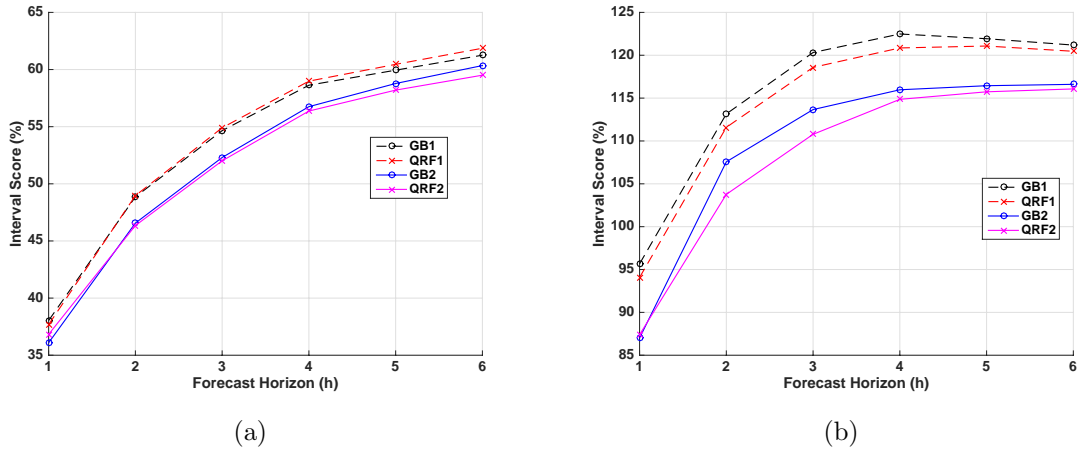


Figure 8: Interval Score ( $IS_{0.2}$ ) of the different intraday methods (a) Site of Desert Rock (b) Site of Le Tampon

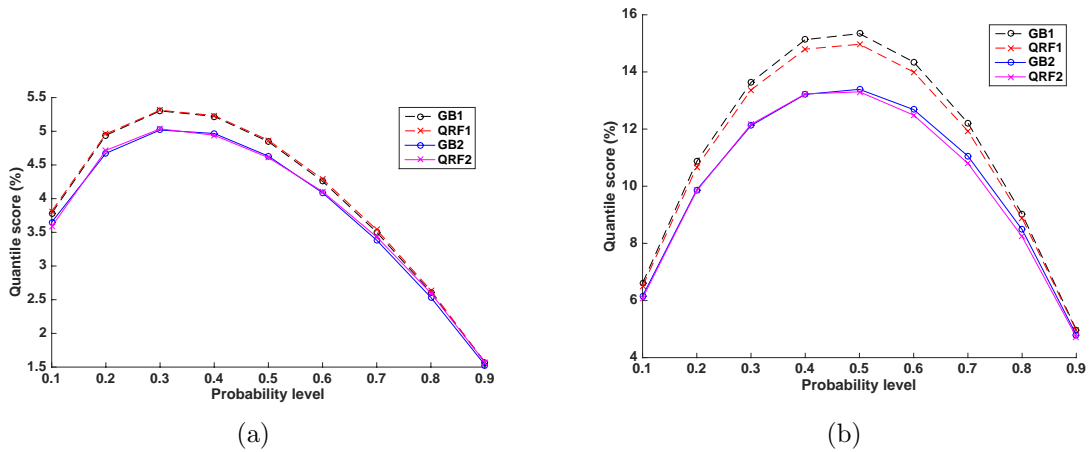


Figure 9: Quantile Score of the different intraday methods (a) Site of Desert Rock (b) Site of Le Tampon

### 676 5.2.6. Quantile score ( $QS$ )

677 Figure 9 plots the quantile score in relation with the probability levels ranging from 0.1  
 678 to 0.9. Again, this detailed analysis of the performance of the models favors the variant  
 679 2 models (and particularly for Le Tampon site). Figure 9(b) reveals a symmetric pattern  
 680 and shows that the highest quantiles and lowest quantiles are rather well estimated for Le  
 681 Tampon. Conversely, regarding the site of Desert Rock, an asymmetric pattern is observed  
 682 as the lowest quantiles are more penalized. This is possibly due to the high occurrences of  
 683 clear skies experienced by Desert Rock.

### 684 5.3. Day-ahead ensemble forecasts

685 The day-ahead ensemble predictions are provided by the Integrated Forecasting System  
 686 (IFS) of the European Centre of Medium-Range Weather Forecasts (ECMWF). We will

687 denote these ensemble forecasts as “ECMWF-EPS”. They consist in 50 perturbed members.  
688 The temporal resolution is of 3 hours and the spatial resolution is of  $0.2^\circ$  in both longitude  
689 and latitude. Consequently, 3h GHI (in  $Wh/m^2$ ) times series recorded on-site are compared  
690 with the nearest ECWMF pixel. In addition, we also propose a post-processed version of  
691 the original ECMWF-EPS forecasts. Indeed, the ensemble prediction systems of the NWP  
692 models commonly suffer from a lack of spread (Leutbecher and Palmer, 2008). To face  
693 this issue, Sperati et al. (2016) proposed a simple approach, named Variance Deficit (VD),  
694 to calibrate the ensemble forecasts. Their method spreads the initial ensemble forecasts  
695 by correcting their variance. The correction factor is evaluated from a training set. The  
696 calibrated ensemble forecasts will be denoted by “ECMWF-EPS + VD”.

697 To assess the quality and to finely characterize these two ensemble forecasting systems,  
698 we follow the approach used for the quantile forecasts (section 5.1) and proposed by (Pinson  
699 et al., 2007). But some of the tools used are different. Therefore, we will first assess the  
700 reliability of the forecast with rank histograms. Even if it is redundant, reliability diagrams  
701 will also be displayed to provide a comparison and to highlight the difficulty to read them in  
702 comparison with rank histograms. Then, we will complement the analysis with a sharpness  
703 evaluation. Lastly, the CRPS proposed by (Hersbach, 2000) and the IGN proposed by  
704 (Roulston and Smith, 2002) will be computed to rank the two ensemble prediction systems.  
705 As for the previous section, the metrics are normalized by the GHI of the test period (see  
706 Table 3).

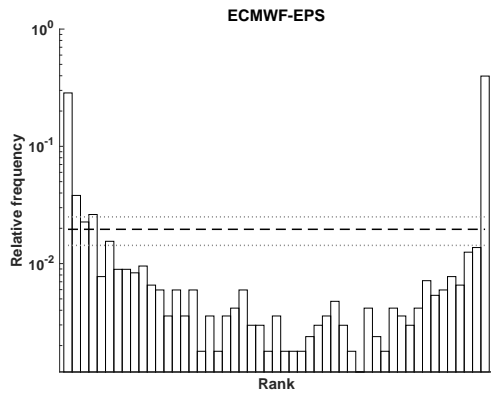
### 707 5.3.1. Reliability assessment

708 As exposed in section 3, the first property to check is the reliability of the forecasts. The  
709 most relevant tool for ensemble forecasts is the rank histogram because it graphically shows  
710 if the members of the ensemble are indistinguishable from the observations. Figures 10 and  
711 11 give the rank histograms before and after the calibration of the ensemble forecasts. As  
712 expected, the rank histograms of the initial ECMWF-EPS forecasts have a U-shape that  
713 corresponds to a lack of spread. The calibration reduces the under-dispersion but a bias  
714 appears for both sites and a large number of ranks remain out of the consistency band.  
715 Reliability diagrams plotted using the uniform method (see section 2 and Figures 12 and  
716 13), give the same conclusion. The calibrated forecasts are more reliable than the original  
717 ones but a large part of the curve falls outside of the consistency bars for both sites. With  
718 the reliability diagram, it is easy to see if the forecasts are reliable or not. But, in comparison  
719 with the rank histogram, it is more difficult to highlight the underlying issues due to a lack  
720 of reliability.

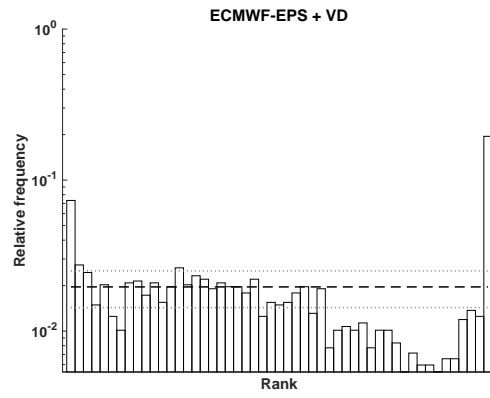
### 721 5.3.2. Sharpness assessment

722 At this point of our analysis, there is normally no need to lead further investigations about  
723 the sharpness of the prediction intervals. Indeed, none of the forecasts are reliable and a  
724 comparison of the sharpness of the forecasts could lead to a misunderstanding. Nevertheless,  
725 we do it for this study case to illustrate this issue. Figure 14 shows sharpness diagrams  
726 for coverage rates ranging from 0% to 100%, for the two sites and for the two considered  
727 ensemble forecasts. To compute the mean size of the central prediction interval  $\delta^{\bar{\alpha}}$ , we assume



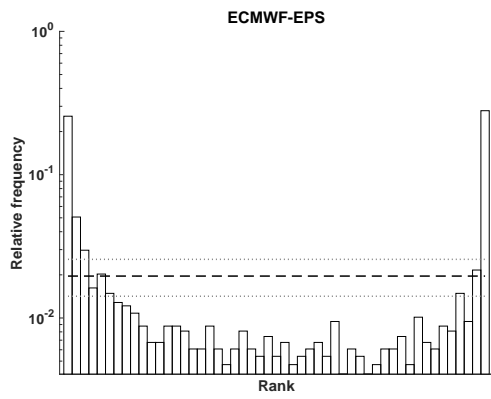


(a)

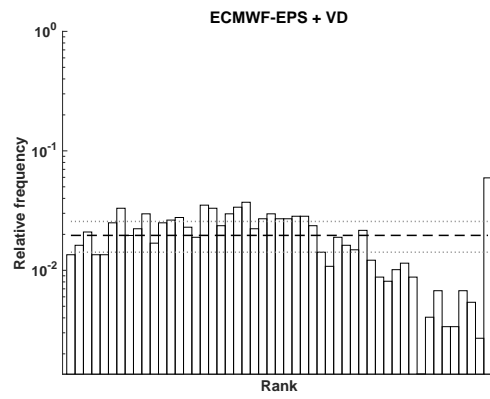


(b)

Figure 10: Rank histogram of raw ECMWF-EPS (a) and ECMWF-EPS calibrated with variance deficit (b) for Desert Rock

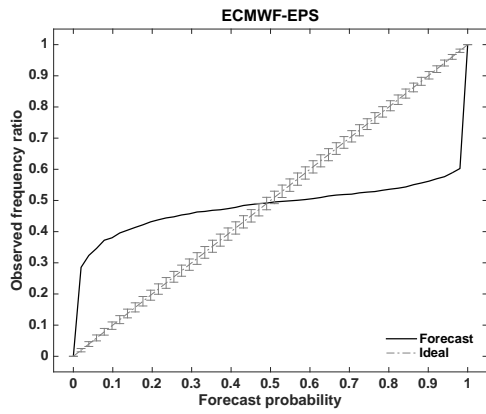


(a)

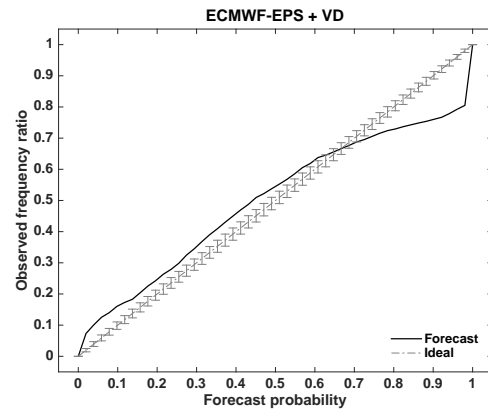


(b)

Figure 11: Rank histogram of raw ECMWF-EPS (a) and ECMWF-EPS calibrated with variance deficit (b) for Tampon

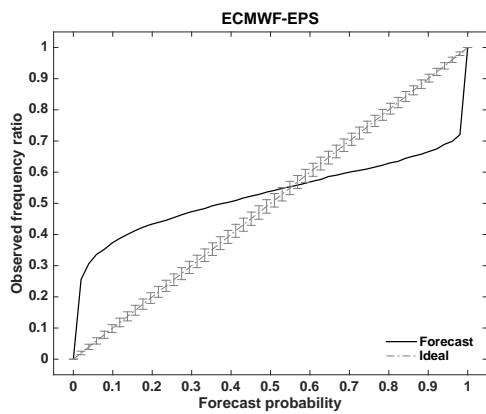


(a)

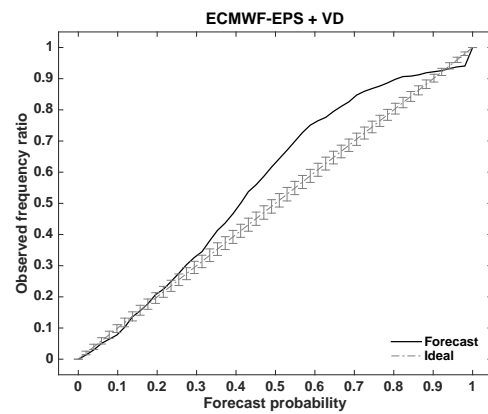


(b)

Figure 12: Reliability diagrams of raw ECMWF-EPS (a) and ECMWF-EPS calibrated with variance deficit (b) for Desert Rock



(a)



(b)

Figure 13: Reliability diagrams of raw ECMWF-EPS (a) and ECMWF-EPS calibrated with variance deficit (b) for Tampon

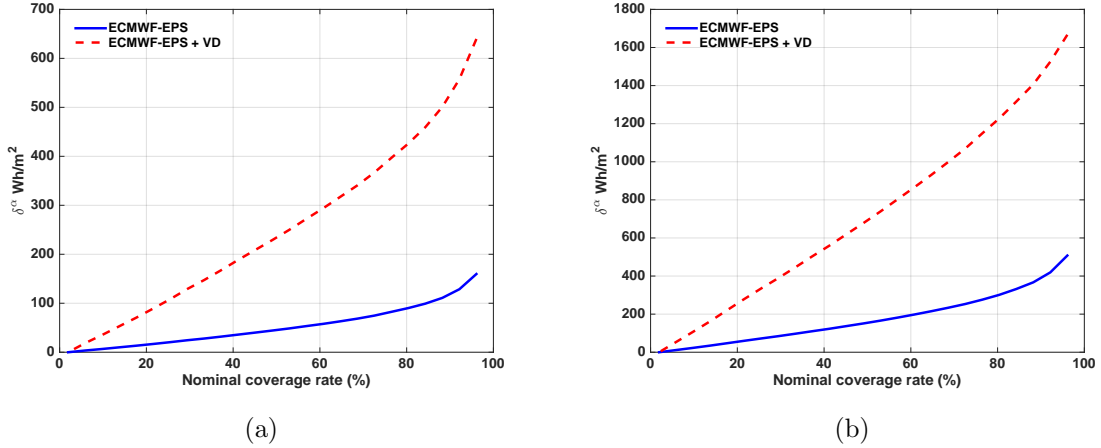


Figure 14: Sharpness of ECMWF-EPS and ECMWF-EPS + VD for Desert Rock (a) and for Tampon (b)

728 an uniform spacing of the quantiles derived from the ensemble (see section 2). As shown by  
 729 Figure 14, predictions intervals (PIs) of original ECMWF-EPS forecasts are narrower than  
 730 the calibrated ones. This is the consequence of the under-dispersion and therefore of the low  
 731 reliability of the ECMWF-EPS forecasts. So, in this case, even if narrow PIs are preferred,  
 732 sharpness diagrams should not be used as criteria to assess the quality of the forecasts.

### 733 5.3.3. CRPS and Ignorance score

734 To complete this analysis, Tables 5 and 6 give the IGN, the CRPS and its decompo-  
 735 sition. For Le Tampon and regarding both scores, the calibration brings an improvement.  
 736 The decomposition of the CRPS highlights that the calibration increases the reliability but  
 737 reduces the resolution. Regarding the site of Desert Rock, the two scores give an opposite  
 738 ranking. The IGN assigns a better score to the calibrated ensemble. Conversely, the CRPS  
 739 better rates the initial ECMWF forecasts. The decomposition of the CRPS shows that the  
 740 increase in reliability, resulting from the calibration, does not counter-balance the reduction  
 741 in resolution. Figure 15 illustrates this difference of scoring for a clear sky that has been  
 742 forecasted and occurred. The original ECMWF forecast (blue line) already contains the  
 743 observation (black line) and the associated CDF is very sharp. So, the IGN and the CRPS  
 744 are already relatively low. The VD method (red dashed line) spreads the CDF and the  
 745 observation falls close to the median of the calibrated CDF where the probability mass is  
 746 the highest. As it is a local score that depends only on the probability at the observation,  
 747 the IGN is slightly improved. Conversely, the CRPS, which takes into account the spread  
 748 of the CDF, increases significantly. Considering the large number of clear sky conditions  
 749 that are forecasted and observed at Desert Rock, the results obtained for this specific case  
 750 can be extended to a whole year. We can conclude that the VD calibration method spreads  
 751 blindly the ECMWF forecasts, even when it is not necessary. As it is a local score, the IGN  
 752 is not able to catch and to quantify such a behavior of forecasting models. Consequently, it  
 753 seems less robust than the CRPS.

Table 5: Scores for Desert Rock

	CRPS (%)	CRPS decomposition (%)			IGN
		Reliability	Resolution	Uncertainty	
ECMWF-EPS	6.97	1.77	37.9	43.1	9.67
ECMWF-EPS + VD	7.37	0.97	36.7	43.1	7.84

Table 6: Scores for Le Tampon

	CRPS (%)	CRPS decomposition (%)			IGN
		Reliability	Resolution	Uncertainty	
ECMWF-EPS	25.1	6.03	23.5	42.6	9.13
ECMWF-EPS + VD	23.1	2.41	21.9	42.6	7.89

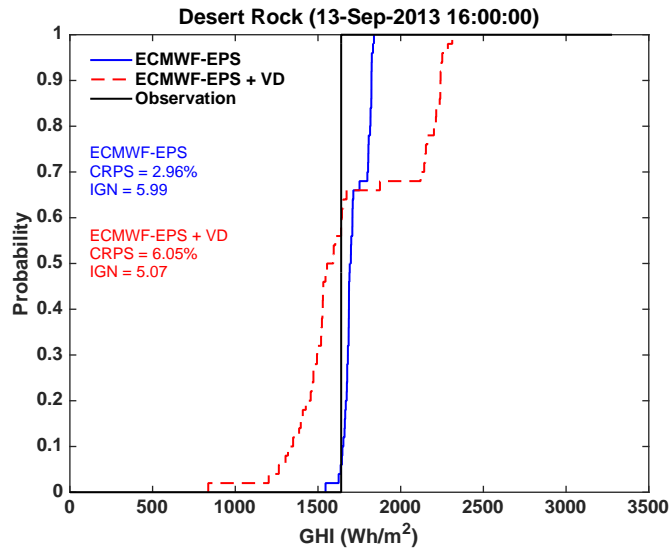


Figure 15: Illustration of the evolution of the CRPS and of the IGN between original and calibrated forecasts: case where these two scores give contradictory information. The CDFs are plotted using the classical definition for ensemble forecasts (see section 2).

## 754 6. Conclusions

755 In this work, we proposed a framework for evaluating solar probabilistic forecasts. Two  
756 types of solar probabilistic forecasts namely ensemble forecasts and quantile forecasts were  
757 used to illustrate the evaluation framework. This latter is based on visual diagnostic tools  
758 and scoring rules originally designed by the weather forecast verification community. For  
759 both types of probabilistic forecasts (quantile and ensemble forecasts), we proposed to follow  
760 the same approach to assess the quality of the models albeit some diagnostic tools are more  
761 appropriate depending on the type of forecast.

762 The proposed approach consists in first evaluating the reliability attribute. Graphical  
763 displays such as reliability diagrams and rank histograms with consistency bars, respectively  
764 for quantile forecasts and ensemble forecasts, are efficient, easy-to-build graphical tools ded-  
765 icated to this purpose. Once the reliability attribute checked, a sharpness analysis can be  
766 conducted. However, in our opinion, even if sharpness is an intuitive property that can be  
767 visually assessed with diagrams, it can only contribute to a qualitative evaluation of the  
768 forecasting methods. More generally, visual diagnostic tools cannot allow one to objectively  
769 conclude on a higher quality of a given model. Therefore, we recommend to systematically  
770 compute an overall score i.e. the CRPS which, in our opinion, might be a standard in assess-  
771 ing probabilistic forecasts of continuous variable. This proper score allows ranking  
772 models and its relative counterpart (i.e. CRPS normalized by the mean irradiance) permit  
773 to carry out sites' comparisons. Furthermore, the decomposition of the CRPS into reliability  
774 and resolution may provide additional insight into the performance of a forecasting system.

775 Also, we recommend to complement the CRPS scoring rule with a set of proper scores like  
776 interval score, ignorance score and quantile score. For instance, quantile score may provide  
777 detailed performance of the models at specific parts of the predictive distributions. Re-  
778 garding the ignorance score, although it can advantageously complement the CRPS results,  
779 attention should be paid to its use, as its locality makes it less robust than the CRPS.

780 Finally, when dealing with ensemble forecasts, dedicated verification tools, such as rank  
781 histograms and the CRPS proposed by (Hersbach, 2000), can be used without any additional  
782 assumptions. Indeed, they assume a classical definition of the underlying CDF and it is  
783 not necessary to define the CDF boundaries. However, care must be taken while deriving  
784 quantiles, prediction intervals and associated metrics from ensembles. As several possibilities  
785 are available, it is important to clearly state which one is used (e.g. uniform or non-uniform  
786 spacing). The authors of this paper have a preference for the uniform spacing because it  
787 defines the quantiles such that the members of the ensemble can be seen as a predictive  
788 distribution.

789 This work focused on the forecasting of the solar irradiance. However, the proposed  
790 methodology and associated tools can be extended to the evaluation of probabilistic forecasts  
791 of solar power generation.

792 **7. Appendices**

793 **Appendix A “Uniform” definition of the CDF and PDF derived from an en-**  
 794 **semble forecast**

795 Let  $E = (e_1, \dots, e_M)$  be an ensemble forecast with  $M$  members  $e_k$ ,  $k = 1, \dots, M$ . The  
 796 uniform definition of the resulting Cumulative Distribution Function (CDF) assigns a prob-  
 797 ability mass of  $1/(M + 1)$  between two consecutive members and for the events that fall  
 798 outside of the ensemble range. The tails of the CDF are bounded by  $e_0$  and  $e_{M+1}$  (see  
 799 figure 2(c)). Considering a linear interpolation between the consecutive members and the  
 800 two limits defined above, the analytic formulation of the CDF  $\hat{F}_k(x)$  corresponding to the  
 801 “uniform” definition is

$$\hat{F}_k(x) = \frac{x + (k\Delta X_k - e_k)}{(M + 1)\Delta X_k}, \quad (14)$$

802 where

$$\Delta X_k = e_{k+1} - e_k \text{ with } k = 0, \dots, M. \quad (15)$$

803 The corresponding Probability Density Function (PDF)  $\hat{f}_k(x)$  is the first derivative of  
 804 the CDF defined above i.e.

$$\hat{f}_k(x) = \frac{d\hat{F}_k(x)}{dx} = \frac{1}{(M + 1)\Delta X_k}. \quad (16)$$

805 **Appendix B Hersbach’s method to compute the CRPS from ensemble fore-**  
 806 **casts**

807 Here, we reproduce the methodology proposed by (Hersbach, 2000) to compute the CRPS  
 808 and its decomposition. Let  $E = (e_1, \dots, e_M)$  be an ensemble forecast with  $M$  members  $e_k$ ,  
 809  $k = 1, \dots, M$  and  $x_{obs}$  the observation. It is important to notice that Hersbach assumes a  
 810 classical definition of the CDF obtained from the ensemble (see figure 2(a)). Thus, the CRPS  
 811 could be seen as the sum of areas defined by the members  $E$ , the square of their associated  
 812 cumulative probability  $p_k$  and the position of the observation  $x_{obs}$ . One then have

$$CRPS = \sum_{k=0}^M \alpha_k p_k^2 + \beta_k (1 - p_k)^2, \quad (17)$$

813 with

$$p_k = \frac{k}{M}. \quad (18)$$

814 The values of  $\alpha$  and  $\beta$  are determined with the position of the observation  $x_{obs}$  when  
 815 pooled within the sorted members. Table 7 gives the values of  $\alpha$  and  $\beta$  for all the possible  
 816 cases. Some care must be taken for  $k = 0$  and  $k = M$ . Indeed, the corresponding intervals  
 817 (i.e.  $(-\infty, e_1]$  and  $[e_M, +\infty)$ ) contribute to the CRPS only if the observation falls outside

Table 7: Determination of  $\alpha$  and  $\beta$ 

$0 < k < M$	$\alpha_k$	$\beta_k$
$x_{obs} > e_{k+1}$	$e_{k+1} - e_k$	0
$e_{k+1} > x_{obs} > e_k$	$x_{obs} - e_k$	$e_{k+1} - x_{obs}$
$x_{obs} < e_k$	0	$e_{k+1} - e_k$
$k = 1, M$ (Outliers)	$\alpha_k$	$\beta_k$
$x_{obs} < e_1$	0	$e_1 - x_{obs}$
$x_{obs} > e_M$	$x_{obs} - e_M$	0

818 the range of the ensemble (see second part of table 7 about the outliers). Finally, considering  
819 a verification dataset of  $N$  forecast-realization pairs, the overall  $\overline{CRPS}$  corresponds to the  
820 mean of the CRPS obtained for each individual forecast i.e.  $\overline{CRPS} = \frac{1}{N} \sum_{i=1}^N CRPS_i$ .

821 Considering ensemble forecasts, the decomposition of the CRPS has no sense for a single  
822 forecast-realization pair. Indeed, such case has null uncertainty and resolution. Therefore,  
823 the decomposition of the  $\overline{CRPS}$  proposed by Hersbach is based on the mean values  $\bar{\alpha}_k =$   
824  $\frac{1}{N} \sum_{i=1}^N \alpha_k^i$  and  $\bar{\beta}_k = \frac{1}{N} \sum_{i=1}^N \beta_k^i$ . The components of the CRPS are

$$\overline{REL} = \sum_{k=0}^M \bar{g}_k [\bar{o}_k - p_k]^2, \quad (19)$$

$$\overline{UNC} = \frac{\sum_{i=1}^N \sum_{j=1}^i |x_{obs}^i - x_{obs}^j|}{N^2}, \quad (20)$$

$$\overline{CRPS}_{pot} = \sum_{k=0}^M \bar{g}_k \bar{o}_k (1 - \bar{o}_k), \quad (21)$$

$$\overline{RES} = \overline{UNC} - \overline{CRPS}_{pot}, \quad (22)$$

828 with

$$\bar{g}_k = \bar{\alpha}_k + \bar{\beta}_k, \quad (23)$$

$$\bar{o}_k = \frac{\bar{\beta}_k}{\bar{\alpha}_k + \bar{\beta}_k}. \quad (24)$$

## 830 Appendix C Decomposition of the CRPS through decomposition of the Brier 831 score

832 Hersbach (2000) showed that the CRPS can be calculated through the integration of the  
833 Brier Score over all possible values of the predictand. The Brier Score (BS) is a scoring  
834 rule used for the prediction of the occurrence of a specific event. Usually, such an event is  
835 characterized by a threshold value  $x$ . The event happened if  $x_{obs} \leq x$  and not happened if  
836  $x_{obs} > x$ . One can then have

$$CRPS = \int BS(x) dx = \int REL(x) dx - \int RES(x) dx + \int UNC(x) dx, \quad (25)$$

Table 8: Contingency Table for threshold  $x$ 

Probability $p_k$	Event occurred		Event not occurred	
	$x_{obs} \leq x$		$x_{obs} > x$	
0	$n_0$		$\hat{n}_0$	
...	...		...	
$i$	$n_k$		$\hat{n}_i$	
...	...		...	
1	$n_M$		$\hat{n}_M$	

837 with

$$REL(x) = \sum_{k=0}^M g_k(x) [o_k(x) - p_k]^2, \quad (26)$$

$$RES(x) = \sum_{k=0}^M g_k(x) [o_k(x) - o(x)]^2, \quad (27)$$

$$UNC(x) = o(x) [1 - o(x)]. \quad (28)$$

840 In our case, the integration over  $x$  of the different components ranges for values of GHI from  
841 0 to the maximum of the climatology.

842 For each value of the predictand  $x$ , terms necessary to compute the Brier Score compo-  
843 nents can be calculated from a 2x2 contingency table (see Table 8). In other words, the joint  
844 distribution of forecasts and observations for  $M + 1$  forecast probabilities can be summarized  
845 in a  $(M + 1) \times 2$  contingency table.

846 The total number of pairs of forecasts/observations  $N$  (i.e. the sample size) is given by  
847  $N = \sum_{k=0}^M n_k + \sum_{k=0}^M \hat{n}_k$ .

$$g_k(x) = \frac{l_k}{N}, \quad (29)$$

848 with  $l_k = n_i + \hat{n}_k$

$$o_k(x) = \frac{n_k}{l_k} \quad (30)$$

$$o(x) = \sum_{k=0}^M g_k(x) o_k(x). \quad (31)$$

850 Figure 16 shows the components of the CPRS through the decomposition of the Brier  
851 Score.

## 852 Appendix D Results of the CRPS decomposition for the intraday models

853 Figure 17 shows the resolution part of the CRPS which confirms the lack of resolution of  
854 the different models as the forecast horizon increases. Regarding resolution, the statements  
855 made regarding the CRPS still hold i.e. the two non-linear models (GB2 and QRF2) that  
856 include the solar geometric predictors lead to better resolution.



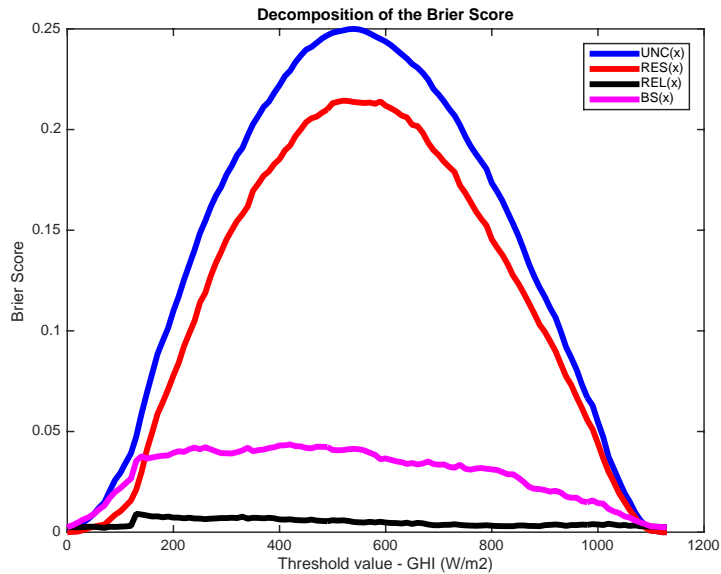
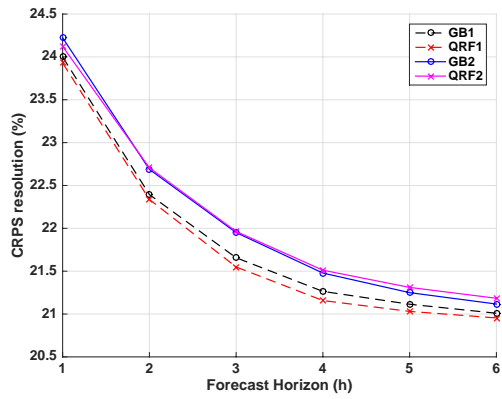


Figure 16: CRPS components through decomposition of the Brier Score (BS) - The area under each curve corresponds to the related CRPS component. Integration of  $BS(x)$  for all threshold values  $x$  gives the CRPS

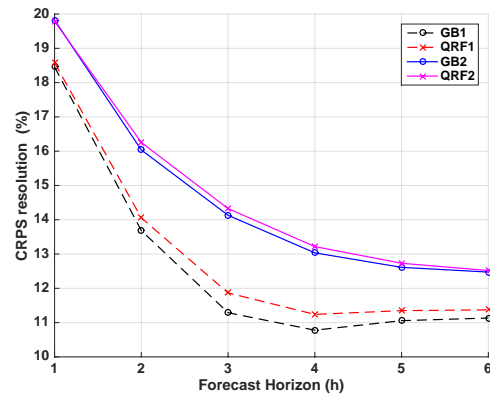
857 Figure 18 plots the reliability component of the CRPS. Surprisingly, the reliability do  
 858 not show a tendency to increase with the lead time. Indeed, we expect the reliability term  
 859 to increase with increasing forecast horizon (we recall that the reliability term is negatively  
 860 oriented i.e. a lower reliability value corresponds to a more reliable forecasts). However,  
 861 in agreement with the reliability assessment, the GB2 model exhibits the lowest reliability  
 862 for the site of Desert Rock while for Le Tampon, low reliability values are obtained with  
 863 the QRF1 model. Nonetheless, it must be noted that the reliability component weakly  
 864 contributes to the CRPS and that the higher quality of the probabilistic forecasts generated  
 865 by the variant 2 models originates from the resolution attribute.

## 866 References

867 Alessandrini, S., Delle Monache, L., Sperati, S., Cervone, G., 2015. An analog ensemble for short-term  
 868 probabilistic solar power forecast. *Applied Energy* 157, 95–110.  
 869 Anderson, J.L., 1996. A Method for Producing and Evaluating Probabilistic Forecasts from Ensemble Model  
 870 Integrations. *Journal of Climate* 9, 1518–1530.  
 871 Ben Bouallègue, Z., 2015. Assessment and added value estimation of an ensemble approach with a focus on  
 872 global radiation forecasts. *MAUSAN*, 541–550.  
 873 Bentzien, S., Friederichs, P., 2014. Decomposition and graphical portrayal of the quantile score. *Quart. J.*  
 874 *Roy. Meteor. Soc.* 140, 1924–1934.  
 875 Bröcker, J., 2012. Evaluating raw ensembles with the continuous ranked probability score. *Quarterly Journal*  
 876 *of the Royal Meteorological Society* 138, 1611–1617.  
 877 Bröcker, J., Smith, L.A., 2007a. Increasing the Reliability of Reliability Diagrams. *Weather and Forecasting*  
 878 22, 651–661.  
 879 Bröcker, J., Smith, L.A., 2007b. Scoring Probabilistic Forecasts: The Importance of Being Proper. *Weather*  
 880 *and Forecasting* 22, 382–388.

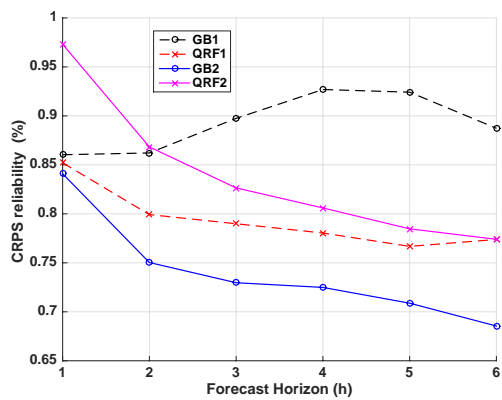


(a)

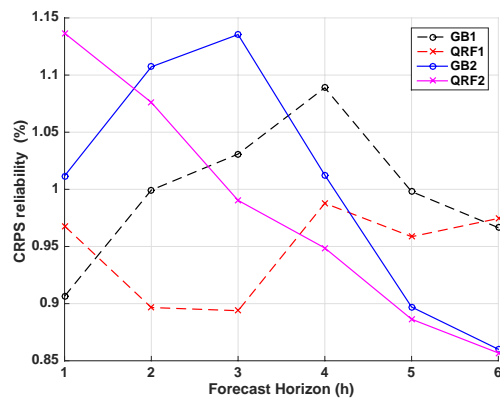


(b)

Figure 17: Resolution Component of the CRPS of the different intraday methods (a) Site of Desert Rock (b) Site of Le Tampon



(a)



(b)

Figure 18: Reliability Component of the CRPS of the different intraday methods (a) Site of Desert Rock (b) Site of Le Tampon

- 881 Brooks, M.J., Du Clou, S., Niekerk, V., L, W., Gauche, P., Leonard, C., Mouzouris, M.J., Meyer, R., Van der  
882 Westhuizen, N., Dyk, V., E, E., Vorster, F.J., 2015. SAURAN : a new resource for solar radiometric data  
883 in Southern Africa. *Journal of Energy in Southern Africa* 26.
- 884 Chu, Y., Coimbra, C.F., 2017. Short-term probabilistic forecasts for Direct Normal Irradiance. *Renewable*  
885 *Energy* 101, 526–536.
- 886 Chu, Y., Li, M., Pedro, H.T., Coimbra, C.F., 2015. Real-time prediction intervals for intra-hour DNI  
887 forecasts. *Renewable Energy* 83, 234–244.
- 888 Coimbra, C.F., Kleissl, J., Marquez, R., 2013. Overview of Solar-Forecasting Methods and a Metric for  
889 Accuracy Evaluation, in: *Solar Energy Forecasting and Resource Assessment*. Elsevier, pp. 171–194.
- 890 Dambreville, R., Blanc, P., Chanussot, J., Boldo, D., 2014. Very short term forecasting of the Global  
891 Horizontal Irradiance using a spatio-temporal autoregressive model. *Renewable Energy* 72, 291–300.
- 892 David, M., Mazorra Aguiar, L., Lauret, P., 2018. Comparison of intraday probabilistic forecasting of solar  
893 irradiance using only endogenous data. *International Journal of Forecasting* 34, 529–547.
- 894 David, M., Ramahatana, F., Trombe, P., Lauret, P., 2016. Probabilistic forecasting of the solar irradiance  
895 with recursive ARMA and GARCH models. *Solar Energy* 133, 55–72.
- 896 Gneiting, T., Balabdaoui, F., Raftery, A.E., 2007. Probabilistic forecasts, calibration and sharpness. *Journal*  
897 *of the Royal Statistical Society: Series B (Statistical Methodology)* 69, 243–268.
- 898 Gneiting, T., Raftery, A.E., 2007. Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the*  
899 *American Statistical Association* 102, 359–378.
- 900 Gneiting, T., Raftery, A.E., Westveld, A.H., Goldman, T., 2005. Calibrated Probabilistic Forecasting Using  
901 Ensemble Model Output Statistics and Minimum CRPS Estimation. *Monthly Weather Review* 133,  
902 1098–1118.
- 903 Golestaneh, F., Gooi, H.B., Pinson, P., 2016a. Generation and evaluation of spacetime trajectories of  
904 photovoltaic power. *Applied Energy* 176, 80 – 91.
- 905 Golestaneh, F., Pinson, P., Gooi, H.B., 2016b. Very Short-Term Nonparametric Probabilistic Forecasting of  
906 Renewable Energy Generation With Application to Solar Energy. *IEEE Transactions on Power Systems*  
907 31, 3850–3863.
- 908 Good, I.J., 1952. Rational Decisions. *Journal of the Royal Statistical Society. Series B (Methodological)* 14,  
909 107–114.
- 910 Grantham, A., Gel, Y.R., Boland, J., 2016. Nonparametric short-term probabilistic forecasting for solar  
911 radiation. *Solar Energy* 133, 465–475.
- 912 Hamill, T.M., 2001. Interpretation of Rank Histograms for Verifying Ensemble Forecasts. *Monthly Weather*  
913 *Review* 129, 550–560.
- 914 Hersbach, H., 2000. Decomposition of the Continuous Ranked Probability Score for Ensemble Prediction  
915 Systems. *Weather and Forecasting* 15, 559–570.
- 916 Hoff, T.E., Perez, R., 2012. Modeling PV fleet output variability. *Solar Energy* 86, 2177–2189.
- 917 Hoff, T.E., Perez, R., Kleissl, J., Renne, D., Stein, J., 2013. Reporting of irradiance modeling relative  
918 prediction errors: Reporting of irradiance modeling relative prediction errors. *Progress in Photovoltaics:*  
919 *Research and Applications* 21, 1514–1519.
- 920 Hong, T., Pinson, P., Fan, S., Zareipour, H., Troccoli, A., Hyndman, R.J., 2016. Probabilistic energy fore-  
921 casting: Global Energy Forecasting Competition 2014 and beyond. *International Journal of Forecasting*  
922 32, 896–913.
- 923 Huang, J., Korolkiewicz, M., Agrawal, M., Boland, J., 2013. Forecasting solar radiation on an hourly time  
924 scale using a Coupled AutoRegressive and Dynamical System (CARDS) model. *Solar Energy* 87, 136–149.
- 925 Iversen, E.B., Morales, J.M., Møller, J.K., Madsen, H., 2016. Short-term probabilistic forecasting of wind  
926 speed using stochastic differential equations. *International Journal of Forecasting* 32, 981–990.
- 927 Jolliffe, I., Stephenson, D., 2003. *Forecast Verification. A practitioner’s guide in atmospheric science*. Wiley,  
928 Chichester, England.
- 929 Jung, J., Broadwater, R.P., 2014. Current status and future advances for wind speed and power forecasting.  
930 *Renewable and Sustainable Energy Reviews* 31, 762–777.
- 931 Khosravi, A., Nahavandi, S., Creighton, D., 2013. Prediction Intervals for Short-Term Wind Farm Power

- 932 Generation Forecasts. *IEEE Transactions on Sustainable Energy* 4, 602–610.
- 933 Koenker, R., Bassett, G., 1978. Regression Quantiles. *Econometrica* 46, 33–50.
- 934 Lauret, P., David, M., Pedro, H., 2017. Probabilistic Solar Forecasting Using Quantile Regression Models. *Energies* 10, 1591.
- 935
- 936 Lauret, P., Voyant, C., Soubdhan, T., David, M., Poggi, P., 2015. A benchmarking of machine learning  
937 techniques for solar radiation forecasting in an insular context. *Solar Energy* 112, 446–457.
- 938 Leutbecher, M., Palmer, T.N., 2008. Ensemble forecasting. *Journal of Computational Physics* 227, 3515–  
939 3539.
- 940 Li, K., Wang, R., Lei, H., Zhang, T., Liu, Y., Zheng, X., 2018. Interval prediction of solar power using an  
941 Improved Bootstrap method. *Solar Energy* 159, 97–112.
- 942 Lorenz, E., Heinemann, D., 2012. Prediction of solar irradiance and photovoltaic power., in: *Comprehensive  
943 Renewable Energy*. Elsevier, Oxford, UK, pp. 239–292.
- 944 Marquez, R., Coimbra, C.F., 2011. Forecasting of global and direct solar irradiance using stochastic learning  
945 methods, ground experiments and the NWS database. *Solar Energy* 85, 746–756.
- 946 van der Meer, D., Widn, J., Munkhammar, J., 2018. Review on probabilistic forecasting of photovoltaic  
947 power production and electricity consumption. *Renewable and Sustainable Energy Reviews* 81, 1484 –  
948 1512.
- 949 Morales, J.M., Conejo, A.J., Madsen, H., Pinson, P., Zugno, M., 2014. Integrating Renewables in Electricity  
950 Markets. volume 205 of *International Series in Operations Research & Management Science*. Springer  
951 US, Boston, MA.
- 952 Murphy, A.H., 1993. What Is a Good Forecast? An Essay on the Nature of Goodness in Weather Forecasting.  
953 *Weather and Forecasting* 8, 281–293.
- 954 NCAR-Research applications laboratory, 2015. verification: Weather Forecast Verification Utilities. R  
955 package version 1.42.
- 956 Pedro, H.T., Coimbra, C.F., 2015. Nearest-neighbor methodology for prediction of intra-hour global hori-  
957 zontal and direct normal irradiances. *Renewable Energy* 80, 770–782.
- 958 Pedro, H.T., Coimbra, C.F., David, M., Lauret, P., 2018. Assessment of machine learning techniques for  
959 deterministic and probabilistic intra-hour solar forecasts. *Renewable Energy* 123, 191–203.
- 960 Pinson, P., McSharry, P., Madsen, H., 2010. Reliability diagrams for non-parametric density forecasts of  
961 continuous variables: Accounting for serial correlation. *Quarterly Journal of the Royal Meteorological  
962 Society* 136, 77–90.
- 963 Pinson, P., Nielsen, H.A., Møller, J.K., Madsen, H., Kariniotakis, G.N., 2007. Non-parametric probabilistic  
964 forecasts of wind power: required properties and evaluation. *Wind Energy* 10, 497–516.
- 965 Pinson, P., Reikard, G., Bidlot, J.R., 2012. Probabilistic forecasting of the wave energy flux. *Applied Energy*  
966 93, 364–370.
- 967 Pinson, P., Tastu, J., 2014. Discussion of “Prediction Intervals for Short-Term Wind Farm Generation Fore-  
968 casts” and “Combined Nonparametric Prediction Intervals for Wind Power Generation”. *IEEE Transactions  
969 on Sustainable Energy* 5, 1019–1020.
- 970 Reikard, G., 2009. Predicting solar radiation at high resolutions: A comparison of time series forecasts.  
971 *Solar Energy* 83, 342–349.
- 972 Roulston, M., Smith, L., 2002. Evaluating Probabilistic Forecasts Using Information Theory. *Monthly  
973 Weather Review* 130, 1653–1660.
- 974 Scolari, E., Sossan, F., Paolone, M., 2016. Irradiance prediction intervals for PV stochastic generation in  
975 microgrid applications. *Solar Energy* 139, 116–129.
- 976 Sperati, S., Alessandrini, S., Delle Monache, L., 2016. An application of the ECMWF Ensemble Prediction  
977 System for short-term solar power forecasting. *Solar Energy* 133, 437–450.
- 978 Tödter, J., Ahrens, B., 2012. Generalization of the Ignorance Score: Continuous Ranked Version and Its  
979 Decomposition. *Monthly Weather Review* 140, 2005–2017.
- 980 Verbois, H., Rusydi, A., Thiery, A., 2018. Probabilistic forecasting of day-ahead solar irradiance using  
981 quantile gradient boosting. *Solar Energy* 173, 313–327.
- 982 Voyant, C., Motte, F., Fouilloy, A., Notton, G., Paoli, C., Nivet, M.L., 2017. Forecasting method for global

983 radiation time series without training phase: Comparison with other well-known prediction methodologies.  
984 Energy 120, 199–208.

985 Weijs, S.V., van Nooijen, R., van de Giesen, N., 2010. Kullback–leibler divergence as a forecast skill score  
986 with classic reliability–resolution–uncertainty decomposition. Monthly Weather Review 138, 3387–3399.

987 Wilks, D.S., 2014. Statistical Methods in the Atmospheric Sciences. An Introduction. Elsevier Science,  
988 Burlington.

989 Winkler, R.L., 1972. A Decision-Theoretic Approach to Interval Estimation. Journal of the American  
990 Statistical Association 67, 187–191.

991 Yang, D., 2019. A universal benchmarking method for probabilistic solar irradiance forecasting. Solar  
992 Energy 184, 410–416.

993 Yang, D., Kleissl, J., Gueymard, C.A., Pedro, H.T., Coimbra, C.F., 2018. History and trends in solar  
994 irradiance and PV power forecasting: A preliminary assessment and review using text mining. Solar  
995 Energy 168, 60–101.

996 Zamo, M., Mestre, O., Arbogast, P., Pannekoucke, O., 2014. A benchmark of statistical regression methods  
997 for short-term forecasting of photovoltaic electricity production. Part II: Probabilistic forecast of daily  
998 production. Solar Energy 105, 804–816.