

Hourly Solar Irradiation Forecast Using Hybrid Local Gravitational Clustering and Group Method of Data Handling Methods

Khalil BENMOUIZA (✉ k.benmouiza@lagh-univ.dz)

Amar Telidji University of Laghouat <https://orcid.org/0000-0001-5222-4055>

Research Article

Keywords: Solar irradiation forecasting, local gravitational clustering, group method of data handling, clustering, phase space reconstruction.

Posted Date: December 3rd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1050483/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Environmental Science and Pollution Research on April 15th, 2022. See the published version at <https://doi.org/10.1007/s11356-022-20114-3>.

1 **Hourly Solar Irradiation Forecast Using Hybrid Local Gravitational Clustering and Group Method of Data**

2 **Handling Methods**

3
4 Khalil BENMOUIZA

5
6 Semiconductors and Functional Materials Laboratory, Amar Telidji University of Laghouat, Bp 37 C, Ghardaia
7 Road, Laghouat 03000, Algeria

8
9 Tel: +213 778 73 85 56

10
11 Corresponding author

12 **Khalil BENMOUIZA**

13 E-mail address: k.benmouiza@lagh-univ.dz

23 **Hourly Solar Irradiation Forecast Using Hybrid Local Gravitational Clustering and Group Method of Data**
24 **Handling Methods.**

25

26 **Khalil BENMOUIZA**

27 Semiconductors and Functional Materials Laboratory, Amar Telidji University of Laghouat, Bp 37 C, Ghardaia
28 Road, Laghouat 03000, Algeria

29

30 **Abstract**

31 The foundation for many solar energies uses as well as economic and environmental concerns is global solar
32 irradiation information. However, due to solar irradiation and measurements variations, reliable worldwide statistics
33 on solar irradiation are frequently impossible or difficult to acquire. In addition, more precise forecast of solar
34 irradiation plays an increasingly important role in electric energy planning and management due to integrating
35 photovoltaic solar systems into power networks. Hence, this paper proposes a new hybrid model for 1-hour ahead
36 solar irradiation forecasting called LGC- GMDH (local gravitational clustering- Group method of data handling). The
37 novel LGC- GMDH model is based on the local clustering that adequately captures the underlying features of the solar
38 irradiation time series. Each cluster is then forecasted using the GMDH method, which is a self-organized system that
39 is capable of handling very complicated nonlinear problems. Finally, these local forecasts are reconstructed in order
40 to obtain the global forecast. Comparative study between the proposed model and the traditional individual models
41 such as; backpropagation neural network (BP), supporting vector machines (SVM), long short-term memory (LSTM),
42 hybrid models such; BP-MLP, RNN-MLP, LSTM-MLP hybrid wavelet packet decomposition (WPD), convolutional
43 neural network (CNN) with LSTM-MLP, and ANFIS clustering shows that the proposed model overcomes
44 conventional model deficiencies and achieves more precise predicting outcome.

45

46 **Keywords:** Solar irradiation forecasting, local gravitational clustering, group method of data handling, clustering,
47 phase space reconstruction.

48

49 1. Introduction

50 Nowadays, a natural desire to develop renewable energies exists. Indeed, we are talking about energy transition,
51 which refers to the transition from the current energy system using non-renewable resources to an energy mix based
52 mainly on renewable resources. Thus, the inclusion of renewable energies in the electricity grid is one of the current
53 issues. Maintaining the balance of the network implies that at all times, electricity production is equal to electricity
54 consumption. However, photovoltaic energy, like wind energy, is so-called intermittent energy. Indeed, it strongly
55 depends on the weather. In order to be able to use renewable energies, solutions must therefore be found intermittently.
56 The development of storage technologies, consumption control, and the development of the electricity network are
57 thus envisaged. Likewise, the forecasting of photovoltaic production will have a primordial role in the future (Ahmed
58 et al. 2020; Guermoui et al. 2020; Lai et al. 2021).

59 A critical step in predicting the amount of electricity produced by a solar panel is predicting the amount of
60 energy received by the panel. Assuming that the inclination of the panels is low, we speak of a forecast of GHI (Global
61 Horizontal Irradiance) which is the total irradiation received by a horizontal surface on the ground. The methodology
62 used to predict irradiation depends on the desired forecast duration. Indeed, it is possible to use:

- 63 • Statistical models considering GHI measurements as time series. These are, for example, autoregressive
64 moving average models (ARMA) (Voyant et al. 2012), autoregressive integrated moving average models
65 (ARIMA) (Shadab et al. 2019), Kalman filter (Soubdhan et al. 2016), and Markov chain (Vindel and Polo
66 2014) are commonly used statistical techniques. These approaches are generally used for reasonably short
67 forecast horizons. However, statistical models are unable to cope with non-stationary data that is abruptly
68 changing. Solar irradiation data is non-stationary, resulting in poor forecast accuracy for sudden fluctuations.
- 69 • As more research is done, artificial intelligence (AI) and machine learning have gained much interest
70 regarding solar irradiation forecast models. Artificial neural network (ANN) (Premalatha and Valan Arasu
71 2016; Benali et al. 2019) , support vector machines (SVM) (Wang et al. 2019) , extreme learning machine
72 (ELM) (Burianek and Misak 2016), wavelet transform (Yadav and Behera 2014), deep learning (Lai et al.
73 2021), and ensemble learning are a few examples of typical AI techniques (Voyant et al. 2017; Alkhatat and
74 Mehmood 2021). These methods can handle the non-linearity presented in the solar irradiation time series,
75 and usually have high accuracy compared to statistical models.

- 76 • Numerical weather prediction (NWP) involves complex mathematical models used to simulate atmospheric
77 change and forecast weather patterns using current weather circumstances. In meteorology, NWP models
78 have met with a great deal of success, but because of their complexity, they cannot be easily implemented
79 (Perez et al. 2013; Verbois et al. 2018).
- 80 • Physical models by observing clouds and their movement are used to forecast the amount of energy received
81 on Earth (Chu et al. 2016; Blanc et al. 2017; Caldas and Alonso-Suárez 2019). This can be done by analysing
82 images taken by ground cameras (forecast up to 30 minutes), or satellite imagery (forecast up to 5 hours).
- 83 • Combinations of these different models to take advantage of their complementarity create a category named
84 hybrid forecasting models. Despite being very complicated, these models have become popular because they
85 can incorporate a number of models and provide more accurate results for forecasting solar irradiation.
86 Several models are proposed in different researches, a review of this category for forecasting solar irradiation
87 is given in Refs. (Guermoui et al. 2020; Álvarez-Alvarado et al. 2021).

88 However, the solar irradiation is a dynamic time series representing high non-linearity (Gan et al. 2012),
89 making it hard for the old forecasting models to obtain good results. Hence, as a solution clustering was introduced
90 recently as a category of hybrid models for forecasting purpose (Benmouiza and Cheknane 2013, 2018) . Clustering
91 (or data partitioning) is an unsupervised classification method that brings together a set of learning algorithms whose
92 goal is to group unlabelled data with similar properties. In general, data are assigned to groups (clusters) so that
93 observations within each group are similar to one another in terms of variables or attributes of interest, and the groups
94 themselves are distinct from one another in terms of variables or attributes (Ghayekhloo et al. 2015; Benmouiza and
95 Cheknane 2018).

96 Hence, based on this idea, we propose in this paper a hybrid model named LGC-GMDH in order to forecast the
97 measured hourly solar irradiation data for one hour ahead. It consists of using a hybrid local gravitational clustering
98 (LGC) (Wang et al. 2018) method and group method of data handling (GMDH) (Farlow 1981; Onwubolu 2016).
99 The main idea consists of dividing the data set into training and testing set using the k-fold cross-validation technique
100 (Kohavi 1995; Klipp et al. 2005). Then, the selected best set will be clustered into groups with the same proprieties
101 using the LGC clustering. However, and due to the dynamic behaviour of the solar irradiation time series, the phase
102 space reconstruction is used before the clustering phase (MacQueen 1967; Benmouiza and Cheknane 2013, 2018). It
103 is a method used to reconstruct full system dynamics using a single time series. Takens (Takens 1981) explained that

104 a single vector of observations describing a chaotic system may be regenerated into a series of multidimensional
105 vectors using this theorem. As long as the embedding size is sufficient enough, the regenerated vectors may exhibit a
106 wide range of essential characteristics of the real-time series.

107 After that, each obtained cluster from the previous step will be forecasted as a single system using GMDH. It is
108 a multi-layered network with a predetermined structure. It offers the advantage of expressing nonlinear dynamics
109 mathematically and allowing higher-order terms without instabilities (Farlow 1981; Onwubolu 2016). GMDH is also
110 used to find input–output connections in models. It looks for connections between one output and a vast number of
111 potential inputs. The network determines which inputs are relevant to the specified system. During training, the
112 network is constructed layer by layer. Each layer contains neurons with just two inputs; their output is a quadratic
113 function of their two inputs.

114 Next, a global forecast is reconstructed based on each local forecast. Then, error metrics such as
115 Mean absolute error (MAE), root mean square error (RMSE), coefficient of determination (R^2) and forecasting skill
116 (FS) are used in order to compare the obtained forecast versus the measured test solar irradiation time series.

117 The main contributions of this paper are:

- 118 - A new multi-branch hybrid structure with clustered inputs is developed, which considers the highly
119 complicated dynamics of the hourly solar irradiation time series and exhibits outstanding prediction
120 performance.
- 121 - The used LGC algorithm in this paper introduces the data mining and clustering notions in order to forecast
122 locally the solar irradiation time clustered series, which has similar characteristics.
- 123 - The proposed LGC- GMDH framework enhance greatly the hourly-ahead solar irradiation forecasting taking
124 the advantages of feature selection using LGC algorithm and the power of the multilayer neural network
125 using GMDH with its predetermined structure.
- 126 - A comparison between the results obtained from our proposed LGC- GMDH model and several forecasting
127 models proposed in literature namely, single models such as; the Back Propagation (BP) neural network,
128 recurrent neural network (RNN), long short-term memory (LSTM), Support vector machines (SVM), and
129 hybrid models such as BP-MLP, RNN-MLP, LSTM-MLP and the hybrid wavelet packet decomposition

130 (WPD), convolutional neural network (CNN) with LSTM-MLP. Also, the ANFIS clustering model is used
131 in this comparison in order to judge the goodness of our work.

132 The rest of this paper is organized as follows: Section 2 describes the principle of the proposed LGC- GMDH
133 model. Sections 3 and 4 introduce the used solar irradiation data as well as the error metrics, respectively. Section
134 5 presents the simulation results and discussions. The comparison between the proposed model and the other
135 models presented in the literature is shown in Section 6. The last section has devoted to the conclusion.

136 2. Methods

137 Our objective is to deliver the proper one-hour ahead forecast model in order to forecast the hourly solar
138 irradiation time series. Hence, the LGC-GMDH model is proposed. This part will discuss the model framework as
139 shown in Fig. 1, which summarises the adopted methodology of how this could work. The details of the whole process
140 are given in the following steps:

- 141 • *Step A: data pre-processing*

142 The first stage for the proposed methodology is the data. It is necessary to collect and analyse the data before
143 any prior use in the forecasting phase. The raw data should be divided into training and testing data sets. However,
144 this data should first pre-process since it contains missing data due to the lack of measurement and suffers from the
145 overfitting problem. To solve this, it was necessary to apply the k-fold cross-validation technique in order to solve the
146 overfitting issue. Part (A) from Fig.1 depicts a conceptual model of cross-validation in its simplest form. First, the
147 data were divided into k random subsets, S_1 , S_2 , and S_k , each of which was referred to as a fold. Each fold is about the
148 same size as the previous one. The model was then run k times to see how well it worked. On each occasion, the k
149 subset was utilized as a test set, while the other subsets were used as the training and validation sets, respectively.
150 After that, the model with the most outstanding performance was selected from among the conducted k tests (Kohavi
151 1995; Refaailzadeh et al. 2009; Wong 2015) .

- 152 • *Step B: data clustering*

153 The nonlinearity presented in the solar irradiation time series could affect the forecast model (Gan et al. 2012).
154 In this step, the clustering was proposed as a solution to discover similar patterns from the solar irradiation time series.
155 Clustering is an effective unsupervised learning method that has been used in various areas like pattern recognition,

156 image processing, bioinformatics, information retrieval, etc. Clustering is the process of separating a data collection
157 using a similarity metric into uncoupled groups. The items are identical in the same cluster and the components are
158 distinct in separate clusters (Benmouiza and Cheknane 2013, 2018). Several clustering techniques have been
159 developed and categorized into categories such as partition-based clustering, hierarchical clustering, grid-based
160 clustering, density-based clustering, model-based clustering, and so on (Wang et al. 2018).

161 However, before the clustering phase, the phase space reconstruction method is proposed in order to rebuild
162 the global behaviour of a dynamic system from a single variable of the system in question (MacQueen 1967; Khalil
163 and Ali 2016; Benmouiza and Cheknane 2018). Our goal is;

- 164 - Determine the minimum, appropriate, embedding dimension for phase space reconstruction for the time
165 series (Kennel et al. 1992; Benmouiza and Cheknane 2018);
- 166 - Identify regions of the reconstructed phase-space which has similar characteristics using local gravitational
167 clustering algorithm;

168 The theoretical background of both phase space reconstruction and LGC algorithm is given in detail in the
169 following sections.

- 170 • *Step C: data forecasting*

171 In this step, the obtained clusters from the previous step will be forecasted using the GMDH (Farlow 1981;
172 Onwubolu 2016). The main idea is to create local forecasts then a global forecast is reconstructed. The GMDH model
173 is a self-organized system in which the structure optimizes itself based on the information provided by the user. This
174 neural network attempts to build a function in a feed-forward neural network based on a second-degree transfer
175 function using the GMDH neural network as the starting point. A computer program determines the number of hidden
176 layers and neurons with the best deterministic transfer function in the generalized Markov decision-making network
177 (GMDH). As a result, the best possible model structure is discovered. Several nonlinear functions, such as the Volterra
178 series and the Kolmogorov-Gabor polynomial, are used to establish a connection between the input and output
179 variables (Farlow 1981; Onwubolu 2016; Vaishnav and Vajpai 2018).

180 The error metrics will be used in order to judge the goodness of the forecast by comparing the obtained data
181 with the measured one. Once the criteria are reached, the algorithm will end.

182

183

184

185

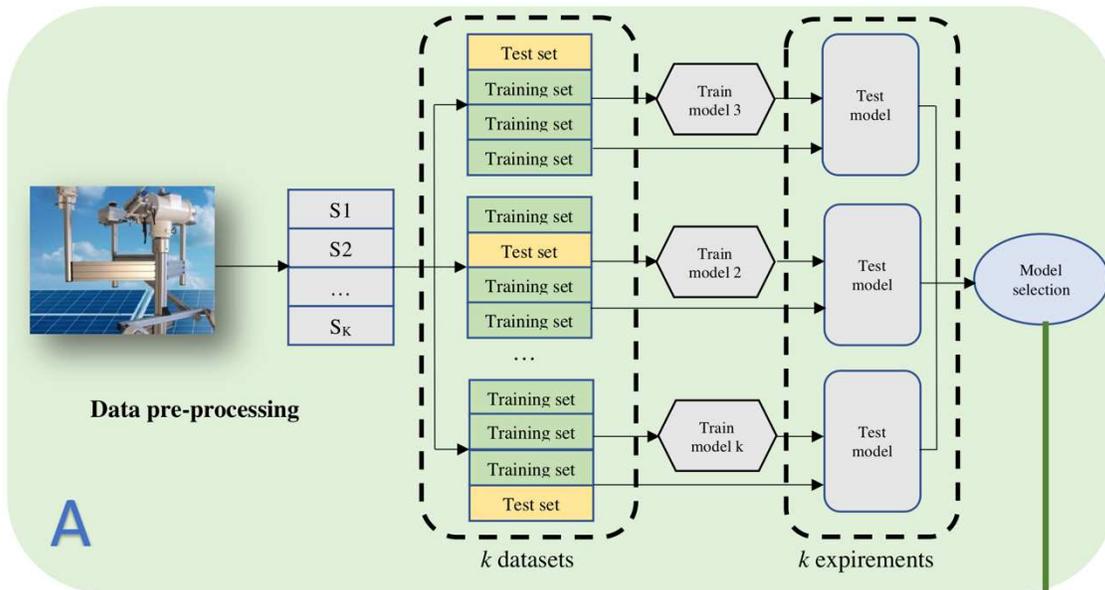
186

187

188

189

190



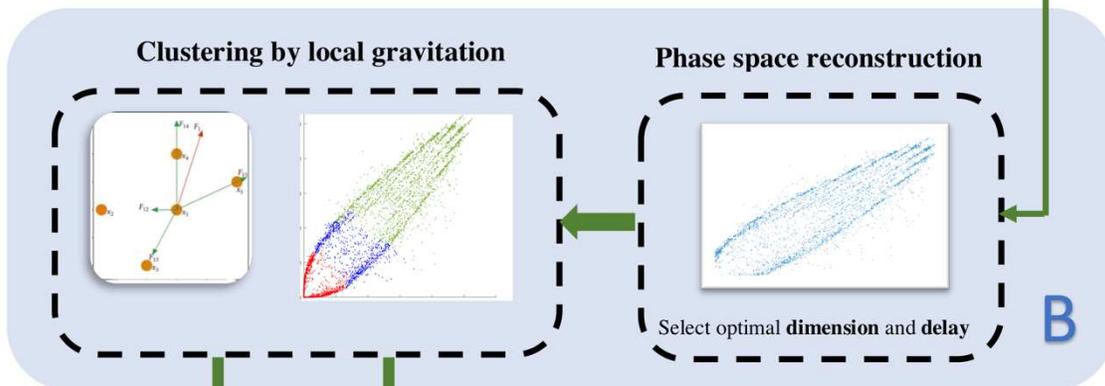
191

192

193

194

195



196

197

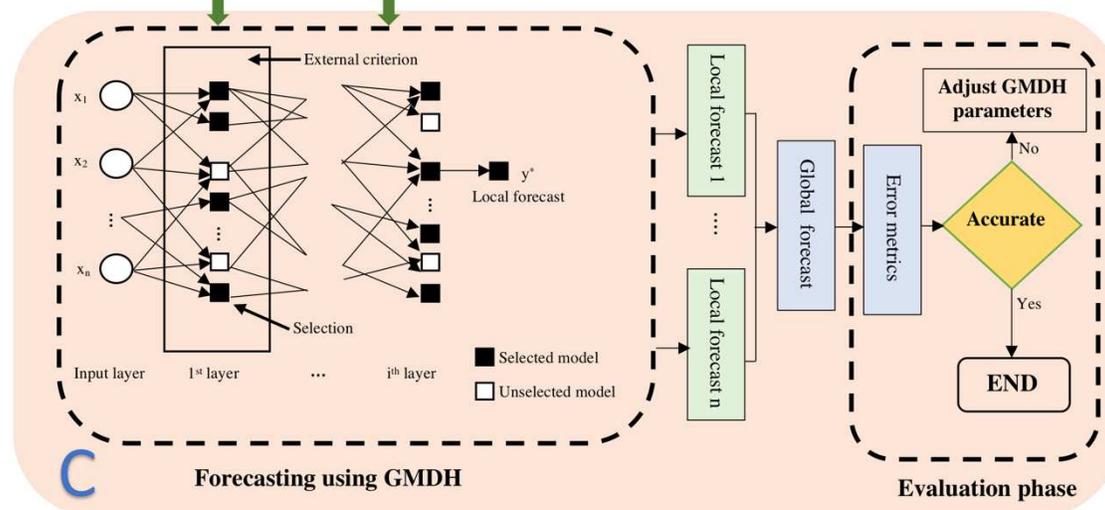
198

199

200

201

202



203

Fig.1 Flowchart of the proposed LGC- GMDH model.

204 In what follows, the mathematical modelling and the best parameters choosing for all the above-mentioned
 205 method are given in detail.

206 2.1 Phase space reconstruction

207 The phase space reconstruction introduced by Takens in 1981 (Takens 1981), is a technique that from a variable
 208 of a dynamic system, allows to reconstruct the global behaviour of its original system. It consists of embedding time
 209 series into high dimension space in order to understand its dynamics, which helps to provide a simple multidimensional
 210 representation of nonlinear time series. The principle of phase space reconstruction is based on the delay of the time
 211 series studied. We delay the time series in order to generate vectors of chosen dimensions which will constitute the
 212 reconstructed phase space (Benmouiza and Cheknane 2013, 2018).

213 The state vectors are defined as the set of d components corresponding to the d system's degrees of freedom.
 214 It is therefore theoretically necessary to have as many independent measurements as the system has degrees of
 215 freedom. In general, we will consider that we have only a sampled series resulting from the experimental observation
 216 of a fluctuating quantity and from which the system's phase space have to be estimated. A possible estimate consists
 217 in defining the set of vectors X_n ;

$$218 X_n = [x_n \ x'_n \ x''_n \ \dots]^T \quad (1)$$

219 Using finite differences, the downside of this method is its high sensitivity to noise. H. Whitney (Whitney
 220 1936) proposed an alternative considering the vectors:

$$221 X_n = [x_n \ x_{n-\tau_1} \ \dots \ x_{n-\tau_{d-1}}]^T \quad (2)$$

222 Where; d is the dimension of the estimated space and $\{\tau_i\}_{i=1,\dots,d-1}$ is a set of delays. d must verify $d \geq 2D +$
 223 1 , where D is the dimension of the state space (Whitney 1936). Authors in (Packard et al. 1980) have applied this
 224 method in the context of chaotic systems. A dynamic system will be said to be chaotic if it is a deterministic dynamic
 225 system whose behaviour is very dependent on the initial conditions in the sense that close initial trajectories diverge
 226 exponentially. The validity of the approach in the context of chaotic systems has been justified by F. Takens (Takens
 227 1981) ; for an infinite series , for any values of τ_i and with the sufficient condition on the dimension d .

$$228 d \geq 2D_2 + 1 \quad (3)$$

229 Where D_2 is the correlation dimension (Grassberger and Procaccia 1983), the space formed by the set of vectors
 230 given in Eq.(2) is topologically and dynamically equivalent to the real phase space in the sense that the geometric
 231 properties of the attractor are preserved.

232 The phase space reconstruction technique involves two parameters to be adjusted. The first is the dimension of
 233 the phase space we want to generate, and the second is the delay we want to apply to the time series.

234 *a) The choice of delay τ*

235 In the case of non-noisy observations over infinite times, the choice of the delay τ is irrelevant (Takens 1981),
 236 but it becomes critical (Michel and Flandrin 1996; Benmouiza and Cheknane 2013, 2018) when the observation time
 237 is limited.

- 238 • If τ is too small, all the coordinates are strongly correlated: $x_k \approx x_{k+1}$, the vectors defined by Eq. (2) with
 239 $\tau_i = i \tau$ are almost collinear, the phase space tends to be a straight line, and the estimated dimension tends
 240 to be 1.
- 241 • If, on the contrary, τ is too large, the coordinates are almost independent. The set of vectors $\{X_n\}_{n=1,\dots,N}$
 242 traverses almost all of the phase space: the system is close to a vector noise with independent coordinates,
 243 and the estimated dimension tends towards the value of the reconstruction dimension.

244 The most commonly used method is to choose τ at the value of the first zero of the estimated autocorrelation
 245 functions given by Eq. (4). It has also been proposed to take the first minimum of mutual information (Fraser and
 246 Swinney 1986) to determine τ .

247
$$C(\tau) = \frac{1}{N} \sum_{n=0}^{N-1} x_n x_{n+\tau} \quad (4)$$

248 *b) Choice of dimension d*

249 Let do $d_0 \in \mathbb{N}$ be the smallest dimension satisfying Eq. (3). If $d > d_0$, according to the theorem of F. Takens
 250 (Takens 1981), all the system dynamics are taken into account. Let $L = \max_n x_n - \min_n x_n$ and $N_{d\tau} = N - (d - 1)\tau$,
 251 the density of observations per unit of hypervolume of dimension d is $\delta = \frac{N_{d\tau}}{L^d}$. Consequently, for a series of N
 252 observations, the probability density for an observation vector to be found in a given neighborhood of characteristic
 253 size l_0 typically varies as $\left(\frac{l_0}{L}\right)^d$. It follows that the estimation variance of the probability density on a partition of the

254 neighbourhood varying as the inverse of the number of points is expressed as a function of $\left(\frac{L}{l_0}\right)^d$ with the
255 reconstruction dimension. It is therefore experimentally recommended to choose d as low as possible, i.e., $d = d_0$.
256 Several methods have been proposed in the literature to determine this minimum value, among them;

- 257 • The estimation of the dimension of the attractor by the algorithm of P. Grassberger and I. Procaccia
258 (Grassberger and Procaccia 1983) using the condition of Eq.(3),
- 259 • The study of the dynamics estimated by the reconstruction of different techniques such as singular value
260 analysis (Auvergne and M. 1988), the estimation of the intrinsic local dimension (Michel and Flandrin
261 1996), the false nearest neighbours method (Kennel et al. 1992) and the vector field method (DT and L
262 1992). The dimension is then chosen as the smallest value d for which two points close in dimension d
263 remain in dimension $d + 1$.

264 2.2 Clustering by local gravitation

265 Gravity theory has been used in clustering for a long time. Numerous gravity-based clustering algorithms have
266 been developed, which replicate the process of objects being attracted and merging by their gravitational force. Each
267 data point is considered an object by these methods and is assigned a mass in feature space. Local gravitation
268 clustering LGC (Wang et al. 2018) is a method that uses local gravitational attraction to cluster data points, in which
269 each data point is treated as a massed object that is attracted to its neighbours. According to (Wang et al. 2018), two
270 key stages are included in the LGC Algorithm. First, LGC differentiated interior points, border points, and unlabelled
271 points; and Second, LGC uses soft-connecting technologies to connect interior points. More details about this
272 algorithm are given in what follows,

273 2.2.1 LGC Algorithm

274 Inspired by Newton's theory of gravity, the local gravitation in the data clustering technique reflects the
275 relationship between a data point and its immediate neighbours. According to the theory of gravity, the attraction force
276 between two-point masses may be calculated using the following formula:

$$277 \quad \vec{F}_{ij} = G \frac{m_i m_j}{D_{ij}^2} \widehat{D}_{ij} \quad (5)$$

278 Where, \vec{F}_{ij} is the force between two points i and j . D_{ij} is the distance between two mass points m_i and m_j . G
 279 is the gravitational constant and \widehat{D}_{ij} is the line direction connecting two points along which the power acts. Because
 280 distances between points within the same local region are unlikely to change considerably, Eq. (5) can be simplified
 281 as follows:

$$282 \quad \vec{F}_{ij} = Gm_i m_j \widehat{D}_{ij} \quad (6)$$

283 The resultant force of point i with its k -nearest neighbors (LRF) is :

$$284 \quad \overrightarrow{LRF}(i, k) = \sum_{j=1}^k \vec{F}_{ij} = Gm_i \sum_{j=1}^k m_j \widehat{D}_{ij} \quad (7)$$

285 \widehat{D}_{ij} contains the directional information between point i and its neighbours, and m_j values are considered as
 286 weighting factors in composing the forces in the neighbourhood

287 Eq. (7) shows that the neighbours of points with greater masses affect over those with smaller masses, and
 288 vice versa. Therefore, authors in (Wang et al. 2018) suggest that the new concept of the LRF replace Newton's theory
 289 of gravity as expressed in Eq. (8).

$$290 \quad \overrightarrow{LRF}(i, k) = \frac{1}{m_i} \sum_{j=1}^k \widehat{D}_{ij} \quad (8)$$

291 Where the mas m_i is defined as ;

$$292 \quad m_i = \frac{1}{\sum_{j=1}^k \widehat{D}_{ij}} \quad (9)$$

293 The Centrality (CE) based local clustering methods are proposed to take advantage of the LRF's information. Eq. (10)
 294 is the definition of the CE of data point i ; k is the number of the neighbours.

$$295 \quad CE_i = \sum_{j=1}^k \cos(\vec{F}_j, \vec{D}_{ij}) / k \quad (10)$$

296

297 An LRF's approximate direction is indicated by a CE value larger than zero. Since

$$298 \quad -1 \leq \cos(\vec{F}_j, \vec{D}_{ij}) \leq 1 \quad (11)$$

299 And

300
$$-k \leq \sum_{j=1}^k \cos(\vec{F}_j, \vec{D}_{ij}) \leq k \quad (12)$$

301 The property of CE is:

302
$$-1 \leq CE_i \leq 1 \quad (13)$$

303 According to Eq. (10), the neighbour's LRF direction is nearly opposite to D_{ij} , and CE is small if a point is at
 304 the border of a cluster. The opposite is true if a point is at the centre of a cluster.

305 In summary, data points in the centre area are differentiated from those in the border region using the physical
 306 principles based LGC technique, which calculates the LRF for each data point. In the first phase of the algorithm,
 307 LGC calculates the LRF field for every data point. After classifying the data points into the interior, border, and
 308 unmarked points, LGC performs the clustering process in the second stage based on LRF information and connects
 309 the internal data points. The data points that have been selected as limit points are used to form a cluster.

310 2.3 Group method of data handling (GMDH)

311 GMDH modelling technique is an effective method for the identification of higher-order nonlinear systems.
 312 A.G. Ivakhnenko was the first person to use this term (Farlow 1981). Furthermore, GMDH is an inductive self-
 313 organizing algebraic model, so it is not essential to know the precise physical model in advance to use it successfully.
 314 Instead, during the training phase, GMDH automatically discovers the dominant relationships among the system
 315 variables. In other words, the optimum neuron's structure is automatically chosen to minimize the values of the
 316 prediction error criterion, and any superfluous neurons are removed from the network. As a result, the GMDH has
 317 excellent generalization capabilities and can accommodate the complexity of nonlinear systems (Kondo 1998).

318 2.3.1 GMDH Model

319 The GMDH network is a data-driven modelling method. It uses mathematical functions to describe the
 320 complicated nonlinear connections between the provided input and output data sets. The GMDH network is comprised
 321 several layers, each of which contains neurons. Each neuron has two inputs and a single output (Nazerfard et al. 2006)
 322 . The output of each neuron is computed using the Ivakhnenko (Farlow 1981) polynomial stated in Eq. (14):

323
$$g(x_1, x_2) = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2 + w_4x_1^2 + w_5x_2^2 \quad (14)$$

324 The GMDH networks are constructed by combining the two input variables in each layer. The coefficients of
325 the polynomials are then calculated for each combination and its associated output using the least square fitting method
326 or regression analysis. To assess and verify the output of the polynomials once they have been computed, external
327 accuracy requirements are used.

328 • **External criteria of accuracy**

329 The model adequacy is evaluated using the external criterion, which is known as the regularity criteria. It
330 evaluates and measures the GMDH network's mean square error for each model neuron to assess the model output.
331 When applied to the network, regularity also tells us which of the input combinations are more important. This method
332 also evaluates the neuron polynomial's capacity and fitness to produce the intended system output. The closer the
333 neuron polynomial fits the data, the lower the regularity requirement should be. The regularity criteria are given by
334 Eq. (15).

335
$$R^2 = \frac{\sum_{i=1}^N (y_i - g_i)^2}{\sum_{i=1}^N (y_i)^2} \quad (15)$$

336 Where; R is the regularity criteria measure, N is the number of samples, y is the desired output, g is the
337 GMDH neuron output.

338 Each neuron output is subjected to the regularity criteria in the GMDH network. In the next part, we will discuss
339 the sorting process used to decide which neurons will survive to go on to the next layer.

340 • **Sorting Out Procedure Description**

341 Since the number of neurons and the number of layers is not specified, GMDH modeling is self-organizing. In
342 order to ensure optimal performance, the best performing neurons in each layer of the GMDH are chosen using the
343 external criterion defined in the preceding part. Neurons that have less than a pre-set number of periodicity criterion
344 but which meet or exceed their established criteria on other neurons will be chosen for use. However, those that fall
345 below the requirements set on other neurons will be removed and discarded from the network. Moreover, each of the
346 layers keeps the lowest regularity criteria. When the following layer's lowest regularity criteria are greater than the
347 preceding layer's smallest regularity criterion, adding layers will no longer produce new regularities. The last layer's
348 neuron will provide the network's final output (Water et al. 2000). The main procedure for GMDH algorithm

349 implementation is described in the following steps (Nazerfard et al. 2006; Onwubolu 2016; Vaishnav and Vajpai
 350 2018).

351 • *Step (1)*

352 The data set is separated into a training set and a checking set. The GMDH neuron's weights are estimated by
 353 using training data, while checking data is utilized to organize the network structures. A heuristic method is used to
 354 divide the data, either by choosing random locations for each group or by analysing the variance of the data (C. 2008).

355 • *Step (2)*

356 Generate every combination of two inputs conceivable among all of the input variables. The number of possible
 357 combinations is represented by $n = \frac{m(m-1)}{2}$, where m and n are the numbers of input variables and the number of
 358 combinations, respectively. After that, expand the inputs to a quadratic polynomial Z for each combination.

359
$$Z = \begin{bmatrix} 1 & x_{11} & x_{21} & x_{11}x_{21} & x_{11}^2 & x_{21}^2 \\ 1 & x_{12} & x_{22} & x_{12}x_{22} & x_{12}^2 & x_{22}^2 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & x_{1N} & x_{2N} & x_{1N}x_{2N} & x_{1N}^2 & x_{2N}^2 \end{bmatrix} \quad (15)$$

360 In this way, the coefficients of the polynomial (g), which were previously stated in Eq.(14), are calculated for
 361 each combination in the training set. To get the polynomial weights, the least- squares fitting method is used, and the
 362 value obtained is given by Eq.(16).

363
$$W = (Z^T Z)^{-1} Z^T y \quad (16)$$

364 • *Step (3)*

365 The third step is to do a detailed evaluation and statistical testing of each polynomial using the testing data
 366 points. Eq. (17) is used to determine the outputs of each polynomial.

367
$$g = ZW \quad (17)$$

368 The criteria used to determine the regularity of each neuron in the first layer are determined after the
 369 polynomials have been computed. Neurons that meet the requirement for being "less than" a certain value are permitted
 370 to continue and become inputs for the neurons in the following layer. The remaining neurons are eliminated from the
 371 network.

372 • Step (4)

373 Finally, the whole process starting with the second step is repeated until the requirement for ending the GMDH
374 network is fulfilled, at which point the procedure ends. When the lowest regularity criterion in the current layer is no
375 longer smaller than the lowest regularity criterion in the preceding layer, the GMDH network will end. In order to get
376 the final GMDH model, the route of the neurons in each layer that corresponds to the lowest regularity criterion is
377 traced backward in time. The flowchart of the GMDH algorithm is given in Fig.2

378

379

380

381

382

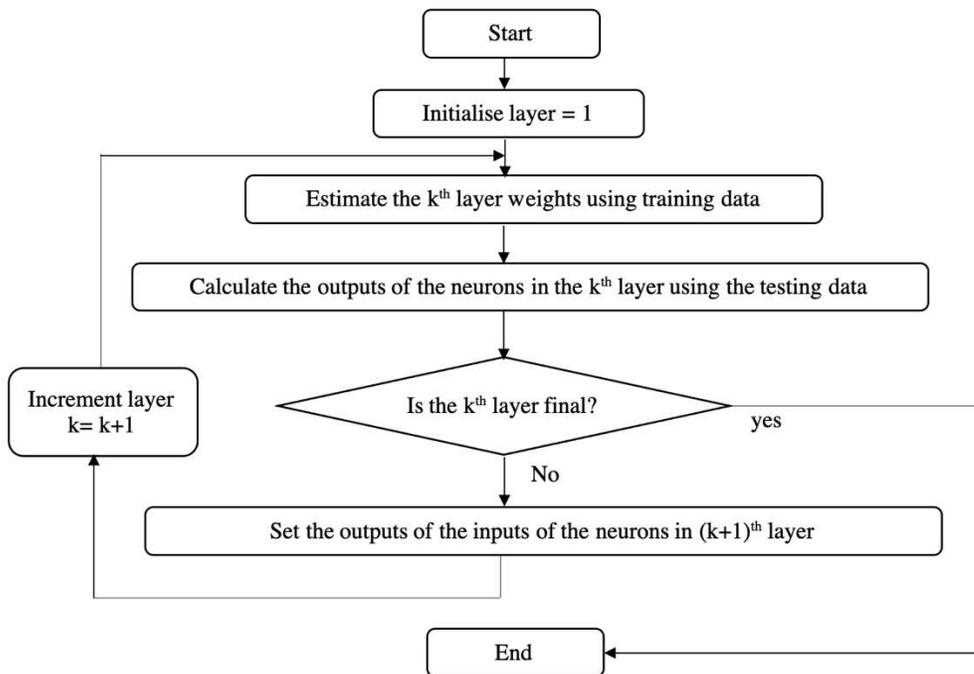
383

384

385

386

387



388

Fig.2 The GMDh network flowchart.

389 3. Data

390 The suggested approach was evaluated in Ghardaia, Algeria . It is situated in the heart of North Africa, along
391 the Mediterranean coast, between the latitudes of 19° and 38° North and the longitudes of 8° and 12° East. A significant
392 portion of the Sahara is located in its southern section (nearly 86 percent of the total area of the country). It is a key
393 player in the area of solar technology due to its location in the sun belt and excellent climatic conditions. Algeria has
394 the most significant solar energy potential in the MENA area. It has the highest effective solar potential where sunlight
395 length is approximately 7.3 hours in the North, 8.3 hours in the highlands, and 10+ hours in the southern regions. As

396 a whole, the entire country's sunlight duration surpasses 3000 hours/year, while in the Sahara, it may approach 3500
397 hours/year (BENMOUIZA 2015).

398 This paper is especially interested in a particular area in the southwest of Algeria, namely the City of Ghardaia
399 (positioned at 33.46° N, 3.78° E, at an altitude of 463 m). The accessible data is gathered by the enerMENA
400 meteorological network project's meteostation (Schüler et al. 2016) . Ghardaïa meteorological station measures the
401 global, DNI, and diffuse solar irradiation; furthermore, it monitors meteorological factors (air temperature, humidity,
402 wind speed, atmospheric pressure). Hence, as a result, a database that's been around for over nine years enabled us to
403 collect an in-depth knowledge of the area's solar potential. For DNI measurements, the meteostation is equipped with
404 CHP1 pyrheliometers, and for diffuse (DHI) and global horizontal irradiation measurements, it is equipped with
405 CMP11 thermal pyranometers (GHI). These instruments are placed on a Solys2 solar tracker with a sun sensor to
406 provide continuous power. Collecting the data is accomplished via the use of a Campbell Scientific CR1000 data
407 recorder.

408 **4. Error metrics**

409 A forecasting model's prediction values must be tested to determine its performance. To find the optimal
410 solution, calculate appropriate error metrics. A model's error metric is a method for measuring the quality of a model
411 and may be used by the forecaster to compare several models. It allows us to assess better how effectively the model
412 completes its various duties. To this end, we have chosen four metrics in order to evaluate the forecasted data noted
413 as \hat{y} versus measured one noted y and a number of observations noted N . These statical metrics are summarized as
414 follows (Benmouiza and Cheknane 2018; Botchkarev 2019; Huang et al. 2021) ;

415 **4.1 RMSE**

416 In mathematics, root means square error (RMSE) is defined as the square of the average of the errors' squares.
417 It is also described as a measure used to evaluate the quality of a forecasting model or predictor, among other things.
418 RMSE takes into account both variance (the dispersion of expected values from one another) and bias (the difference
419 between two anticipated values) (the distance of predicted value from its true value). This error measure is likewise
420 always positive, with lower values being preferable. The RMSE number is in the same unit as the predicted value,
421 which is a benefit of this approach. In comparison to Mean square error MSE, this makes it simpler to comprehend.
422 The RMSE may also be compared to the Mean Absolute Error MAE to see whether the prediction includes significant

423 but occasional mistakes. The greater the gap between RMSE and MAE, the greater the inconsistency of error
424 magnitude. This measure may gloss over issues related to poor data volume. Eq. (18) gives the calculation of RMSE.

425
$$RMSE = \sqrt{\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}} \quad (18)$$

426 **4.2 Mean Absolute Error MAE**

427 MAE is the average of the absolute difference between forecasted and actual values. The MAE informs us how
428 much inaccuracy we may anticipate from the prediction on average. The error numbers are in the predicted values
429 original units, and MAE = 0 implies that there is no mistake in the forecasted values. The lower the MAE number, the
430 better the model; a value of 0 indicates that the prediction is error-free. In other words, when comparing several
431 models, the model with the lowest MAE is seen to be superior. However, MAE does not show the proportional
432 magnitude of the mistake, making it impossible to distinguish between large and minor errors. It may be used with
433 other measures (Root Mean Square Error) to assess whether the mistakes are greater. Furthermore, MAE may obscure
434 issues related to low data volume. The MAE is calculated using the Eq.(19);

435
$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i| \quad (19)$$

436 **4.3 R-squared value**

437 R-squared or R^2 is a statistical measurement that illustrates how near the data are to the fitted regression line.
438 It is sometimes known as the coefficient of determination, or the coefficient of multiple determination for multiple
439 regression or the correlation coefficient. R-squared represents the proportion of the response variable variance
440 described by a linear model. R^2 values are between 0 and 100%. A correlation coefficient between 0% and 100%
441 implies that the model has little explanatory power in explaining the variability of the response data around their mean.
442 The 0% indicated that there is no fluctuation in the response data around its mean; therefore, the model has explained
443 nothing. When 100% appears, the model explains all of the response data's variability around its mean. Hence higher
444 R^2 values, the better the model fits the data. R^2 is calculated using Eq. (20);

445
$$R^2 = 1 - \left(\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2} \right) \quad (20)$$

446

447

448 4.4 Forecasting skill

449 The forecasting skill (FS) is defined as the degree of accuracy with which a prediction is associated with
450 estimating of actual values. It is calculated using Eq.(21) (Mazorra Aguiar et al. 2015; Schmidt et al. 2016) ;

$$451 \quad FS = 1 - \left(\frac{RMSE}{RMSE_{smart}} \right) \quad (21)$$

452 $RMSE_{smart}$ is the smart persistence model, which consists of forecasting clear sky index for each time horizon
453 h persists for the next time step given by;

$$454 \quad k_t^*(t+h) = k_t^*(t)GHI_{clear}(t+\Delta t) \quad (22)$$

455
456
457 k_t^* is the clear sky index given by;

$$458 \quad k_t^* = \frac{GHI}{GHI_{clear}} \quad (23)$$

459 Where ; GHI is the measured horizontal hourly solar radiation data at ground level and GHIclear is the
460 calculated clear sky hourly solar radiation data (W.M.O. 1981; Tadj et al. 2014).

461 A perfect forecast will have an FS equals 1, and an FS equals 0 means there is no improvement over the
462 reference model.

463 5 Simulation results

464 The main objective is to provide the best model for 1-hour ahead GHI forecasting. Hence, we have proposed
465 the clustering using by local gravitation method and the forecasting using GMDH. To this end, we will first train the
466 model from 1st January 2020 to 30th November 2020, and test it for 744 hours (total predicted hours) from 1st December
467 2020 to 31st December 2020. Then, two years have been chosen for training and 6 months for testing.

468 First, and using the time delay embedding technique, the phase space reconstruction of hourly solar irradiation
469 was obtained. We determined the best time delay for the phase space reconstruction to be one, with dimensions equal
470 to two. Next, the LGC algorithm was applied to the results obtained from the phase space reconstruction. The
471 clustering results from the LGC show that the number of clusters is 3.

472

473

474

475

476

477

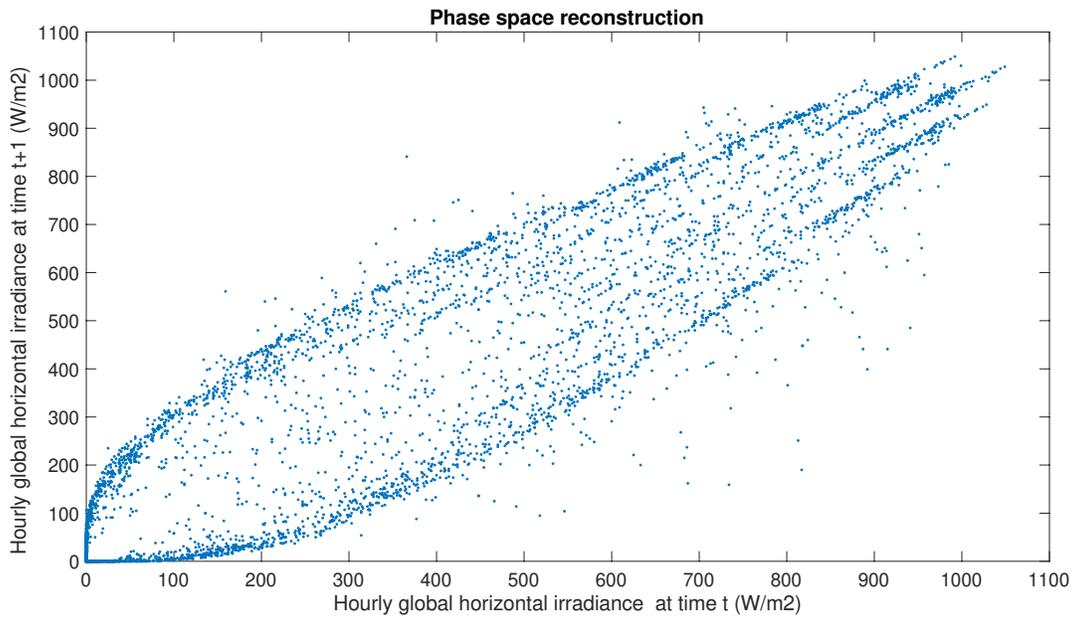
478

479

480

481

482



483

Fig.3 Phase space reconstruction of the GHI at time t and t+1 (1st January to 30th November 2020, Ghardaia).

484

The simulation results of the reconstructed phase space, the LGC clustering, and the three subsets obtained

485

from the LGC clustering algorithm of the solar irradiation are shown in Figs. 3 ,4, and 5, respectively.

486

487

488

489

490

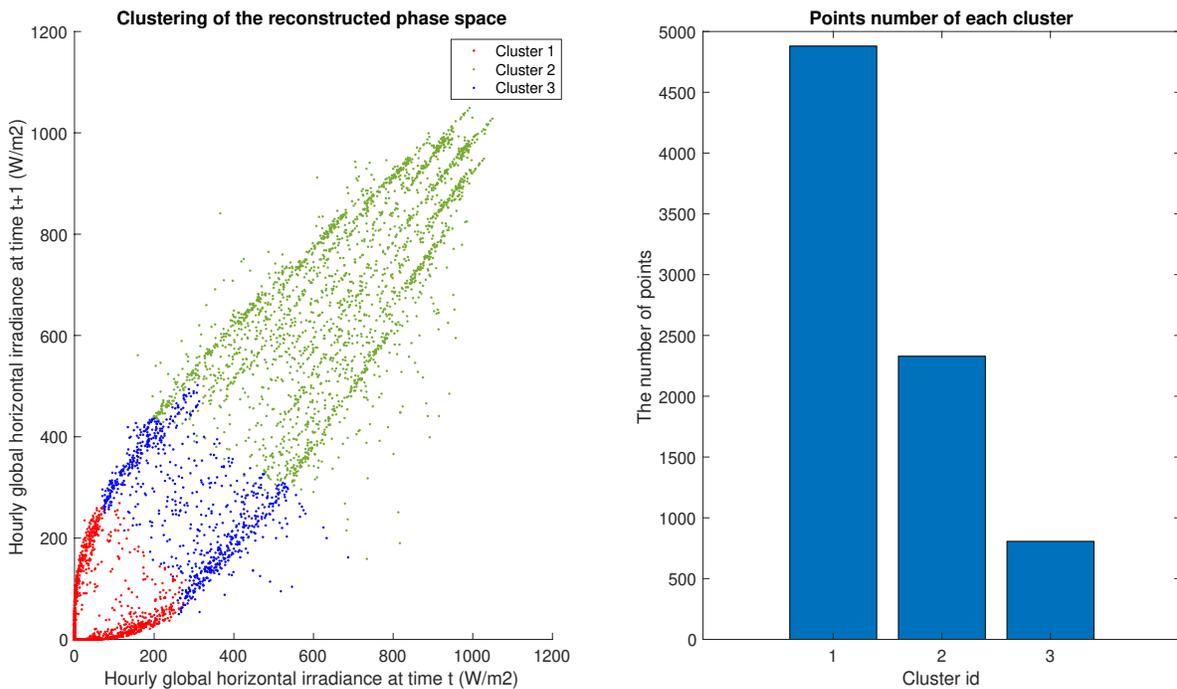
491

492

493

494

495



496

Fig.4 LGC for the reconstructed phase space of the training data (1st January to 30th November 2020, Ghardaia).

497 These figures show that the solar irradiation time series had been grouped into three clusters. High values solar
498 irradiation (cluster N°2) that describe the number of sunny hours throughout the year (for example, the number of
499 hours in the middle of the day when there are no clouds) are contained in this cluster. These values are contrasted with
500 medium values solar irradiation (cluster N°3) that depict hours when the sky is partly cloudy (for example, hours from
501 9 AM to 11 AM and from 2 PM to 4 PM). Finally, values of low solar irradiation (cluster N°1) describe hours when
502 the sky is entirely cloudy (for example, when it is raining or snowing) (or the hours of sunshine and sunset).

503

504

505

506

507

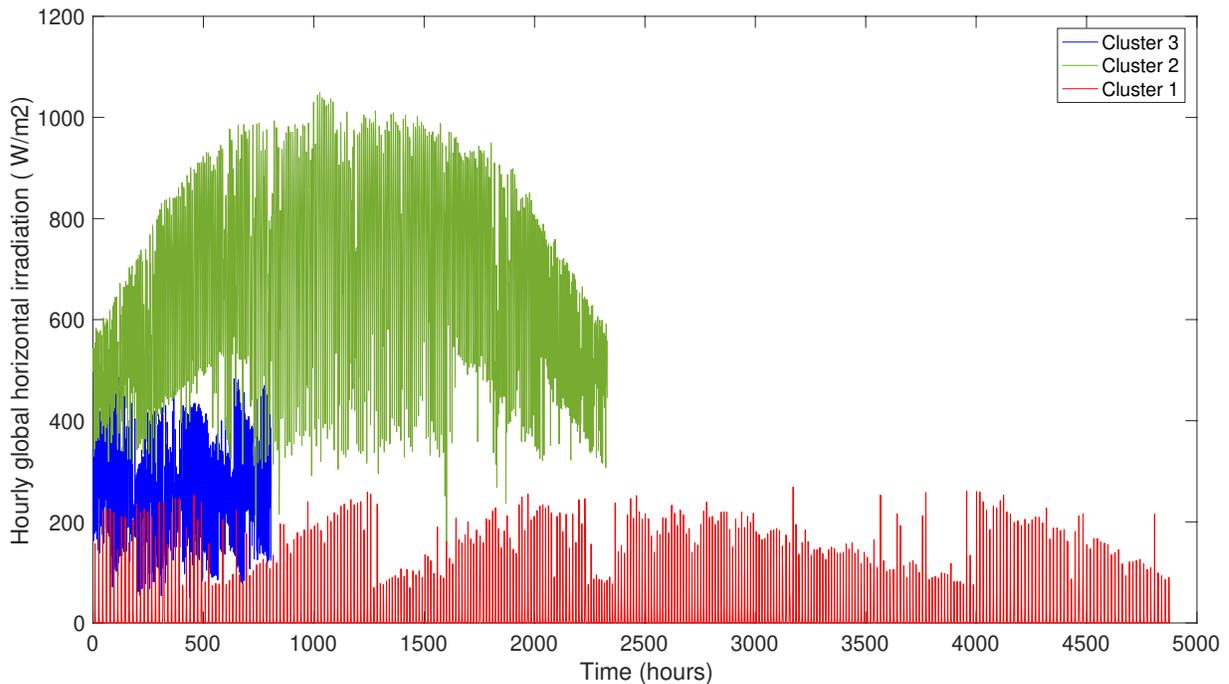
508

509

510

511

512



513

Fig.5 Obtained clusters from the LGC.

514

515

516

517

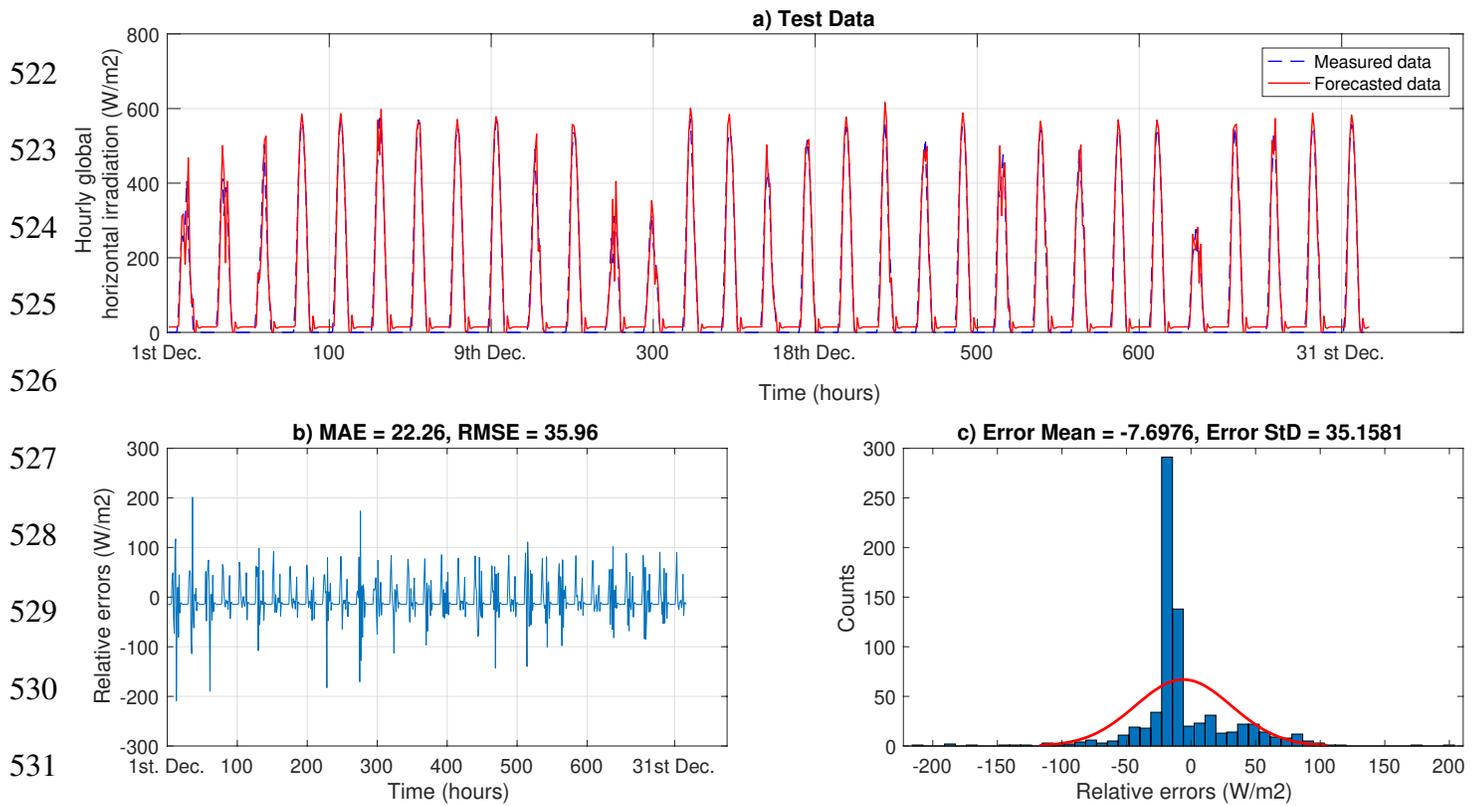
518

519

520

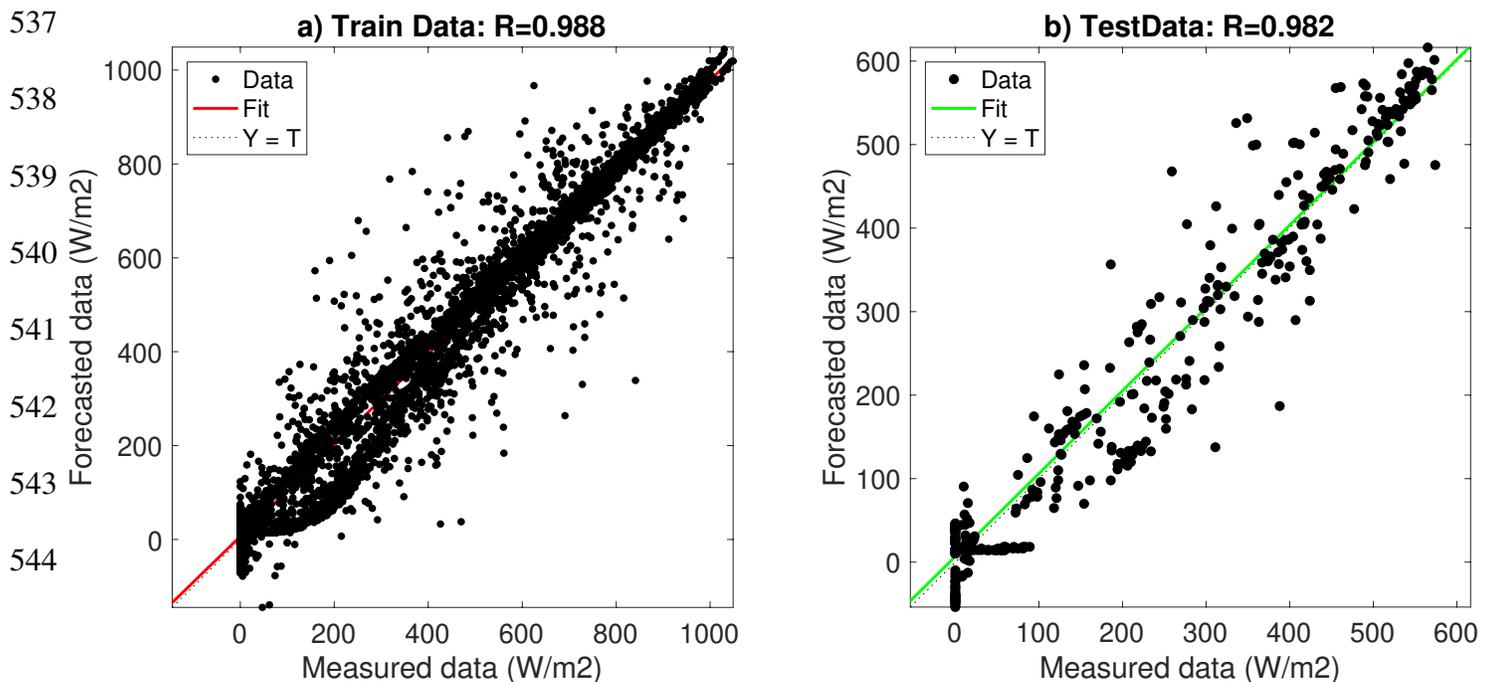
521

The next step consists of forecasting each cluster alone, then, obtain the global forecast. In the simulation, various numbers of parameters are chosen in order to achieve the best forecast. The calculation of errors has been used to assess the performance of the system. The simulation results for the best configuration of the forecasted and measured test data are shown in Fig.6 for December 2020; the red line represents the forecasted data, and the blue dot line shows the measured hourly global solar irradiation data. Furthermore, the relative error and its histogram are also shown in this figure.



532 Fig.6 a) Forecasted and measured testing data for December 2020 using hybrid LGC with GMDH model, b) the
 533 relative error, and c) the histogram.

534 In addition, the training and testing data versus its fit with their R-squared values are shown in Fig.7. The
 535 majority of the locations in both the predicted and measured series are within a few degrees of each other. Because of
 536 the large number of cloudy days, there are some delays in the system's response time.



545 Fig.7 The measured data versus estimated training (a) and testing (b) for Ghardaia 2020 using hybrid LGC with
 546 GMDH model.

547
 548 A further evaluation of the forecasted hourly global horizontal time series was carried out by computing the
 549 root mean squared errors (RMSE), mean absolute error (MAE), the R^2 value, and forecasting skill for training and
 550 testing data set between the actual data and the forecasted data for the period from the first of December 2020 to the
 551 last day of December 2020. Several tests were performed by changing the number of neurons and layers for each
 552 forecast. The results are expressed in Table 1.

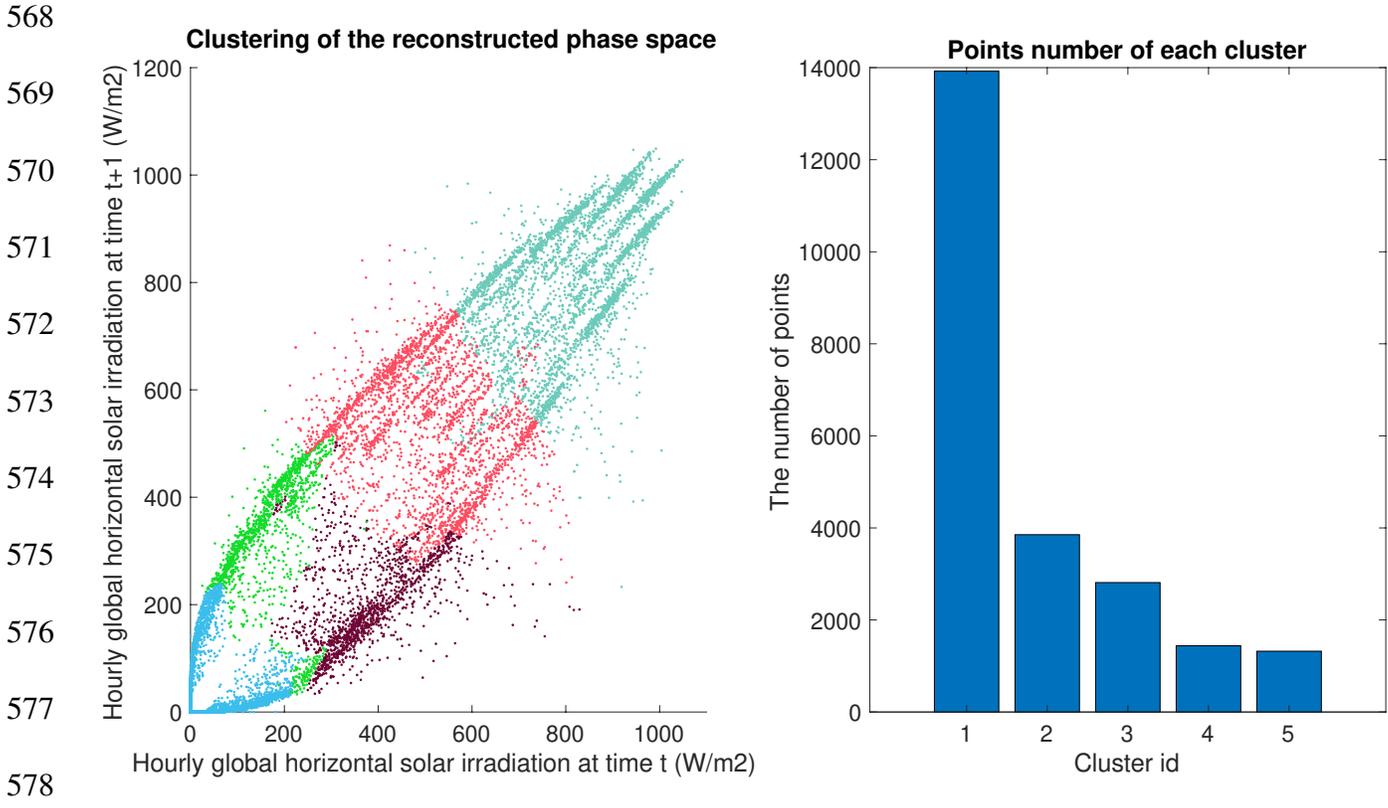
553 **Table 1** RMSE, MAE, R^2 , and FS results between measured and forecasted data for different neurones and layers.

layers	neurons	Training data			Testing data			
		RMSE (W/m ²)	MAE (W/m ²)	R^2	RMSE (W/m ²)	MAE (W/m ²)	R^2	FS (%)
5	25	51.61	31.46	0.986	37.59	26.43	0.978	47.85
3	10	53.46	33.96	0.985	37.90	27.48	0.978	47.79
20	50	49.65	29.12	0.987	36.71	25.18	0.980	47.98
60	100	48.61	27.26	0.988	35.96	22.26	0.982	48.12

554
 555 Following the results of Table 1, we can see clearly that a high number of layers and neurons gives the lowest
 556 errors and highest R^2 and forecasting skills. However, it takes more computational time. Hence, it is not necessary to
 557 increase their number since the results are close to each other. Moreover, from the results of Table 1 and Fig.6 which
 558 represent the plots of the best configuration we can see an overall RMSE was equal 35.96 W/m² and the MAE equals
 559 22.26 W/m², which may be regarded as excellent predicted values.

560 Furthermore, based on Fig. 7, the R-squared value computed by Eq. (20) is equal 0.982, which is a positive
 561 value with a good forecasting skill of 48.12%. Finally, based on the simulation findings, it was determined that this
 562 approach would be an excellent way to forecast solar radiation data.

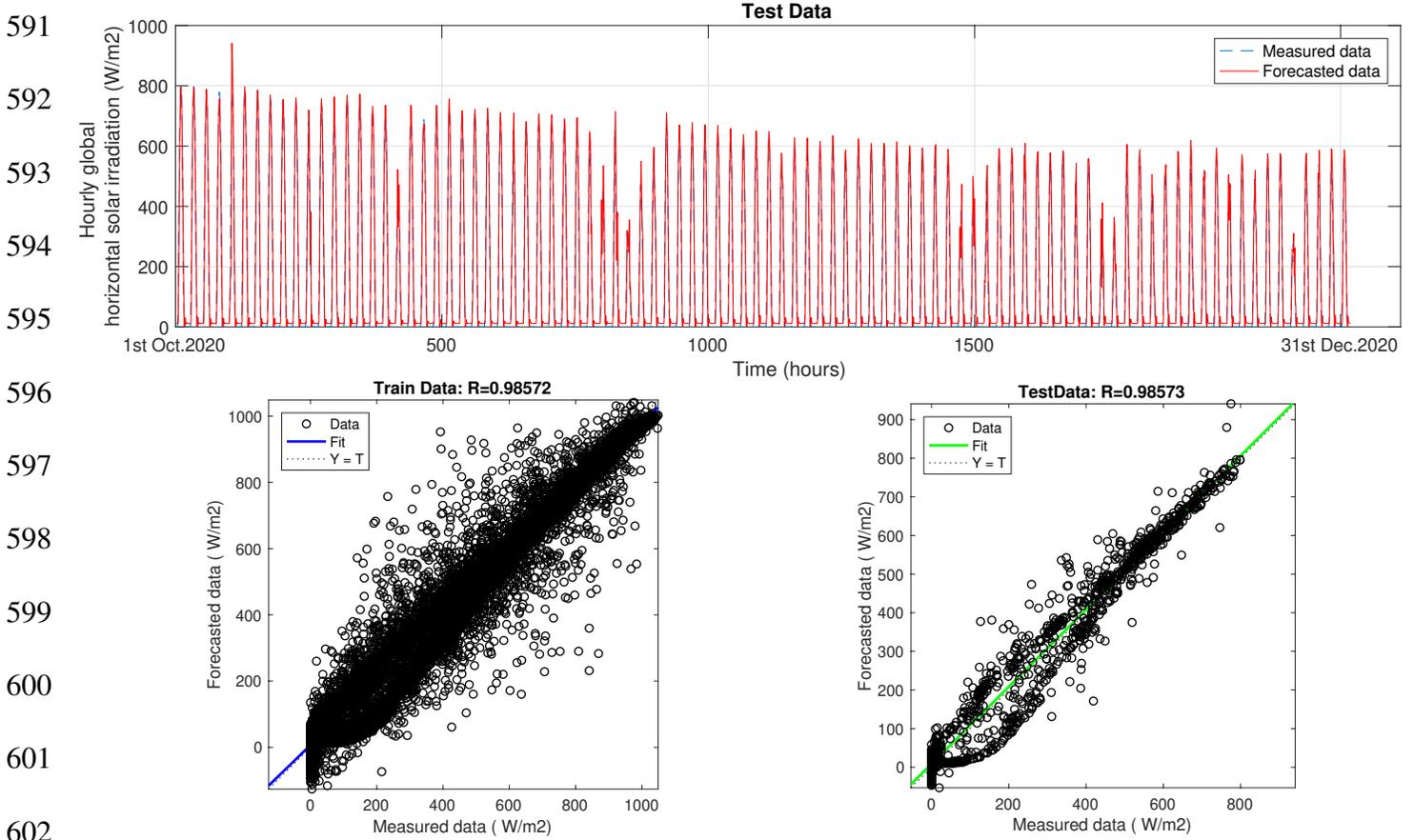
563 In the same way, we have tried larger data sets in order to forecast more than one month. Hence, we have
 564 chosen a solar irradiation time series measurement for three years from the 1st January 2018 to the 31st December 2020
 565 for the site of Ghardaia, Algeria. Thirty-three months were used for training and 3 months were used as forecast
 566 horizon. The same methodology applied for 1 month forecast has been applied for the case of three months. The results
 567 of the clustering of the reconstructed phase space are shown in Fig. 8.



579 Fig.8 LGC for the reconstructed phase space of the training data (1st January 2018 to 30th September 2020,
 580 Ghardaia).

581 In addition, the comparison between measured and forecasted hourly solar irradiation data for three months as
 582 well as the measured data versus forecasted data is shown in Fig.9. From these figures, we can clearly see the goodness
 583 of the LGC- GMDH model. It forecasts in a good manner even on cloudy days with an R^2 value is equals to 98.57%.
 584 These results clearly prove the robustness of the proposed model to forecast solar radiation time series.

585
 586
 587
 588
 589
 590



603 Fig.9 The measured versus forecasted test data LGC with GMDH model (1st October 2020 to 31th December 2020,
 604 Ghardaia).

605 **6. Comparison**

606 Various comparison tests were conducted with other models to check the accuracy of the proposed LGC-
 607 GMDH forecast model. Single forecasting methods such as Artificial Neural Networks (ANN) (Faceira et al. 2015;
 608 Voyant et al. 2017), Support Vector Machines (SVM) (Wang et al. 2019), Recurrent Neural Networks (RNN) (Zhang
 609 and Behera 2012), and their hybrid models have been extensively utilized in solar irradiation forecasting. As a result,
 610 the use of single benchmark models, such as those presented by back propagation BP (Laopaiboon et al. 2019), SVM,
 611 generic RNN, and LSTM (Malakar et al. 2021; Kumari and Toshniwal 2021), is studied extensively in this paper.
 612 New multivariate time series models have also been tested alongside the three models mentioned above, which use
 613 multivariate time series data: the BP-MLP model, the RNN-MLP model, the LSTM-MLP model , and WPD-CNN-
 614 LSTM-MLP (Huang et al. 2021) . Moreover, clustered ANFIS network using fuzzy c-means, subtractive clustering,

615 and grid partitioning model presented in (Benmouiza and Cheknane 2018) , which provides good results compared to
 616 other models is also considered in this comparison.

617 For the comparison purpose, solar irradiation data and other related climatic parameters (including air
 618 temperature, relative humidity, wind speed, and other data) for the location of Denver (39440 N, 105 and 110 W,
 619 Colorado, USA) in the United States are utilized to assess the performance of the proposed model. The data was taken
 620 from the Measurement and Instrumentation Data Center (MIDC), National Solar Irradiation Database supplied by the
 621 National Renewable Energy Laboratory, <https://midcdmz.nrel.gov/>. The dataset from 1st January 2012 to 31st
 622 December 2015, was utilized as the training dataset, and the dataset from 1st January 2016 to 31st December 2016 was
 623 used for testing the models. The original data is captured and recorded on a continuous minute basis ,then this data
 624 was transformed to hourly data.

625 The RMSE , MAE , R^2 and FS have been chosen as error metrics in order to judge the goodness of our proposed
 626 model. The results are expressed in Table 2 .

627 **Table 2** Comparison of RMSE , MAE , R^2 and FS for Denver (January 1, 2016 to December 31, 2016).

	Model	RMSE	MAE	R^2	FS
Single models	BP	76.9272	31.5464	0.9597	0.3167
	SVM	77.9152	52.8904	0.9599	0.3086
	RNN	70.5922	35.6135	0.9666	0.3730
	LTSM	68.7213	30.6064	0.9680	0.3896
Hybrid models	BP-MLP	63.1678	25.4427	0.9730	0.4389
	RNN-MLP	59.8450	23.1430	0.9766	0.4684
	LSTM-MLP	50.5665	19.5735	0.9832	0.5513
	WPD-CNN-LSTM-MLP	46.1336	20.3853	0.9858	0.590
Clustering models	ANFIS-clustering	45.5463	19.2985	0.9860	0.598
	LGC-GMDH	43.7268	19.1856	0.9887	0.613

628
 629 Moreover, the errors bars for these results are shown in Fig.10.

630
 631
 632
 633
 634

635

636

637

638

639

640

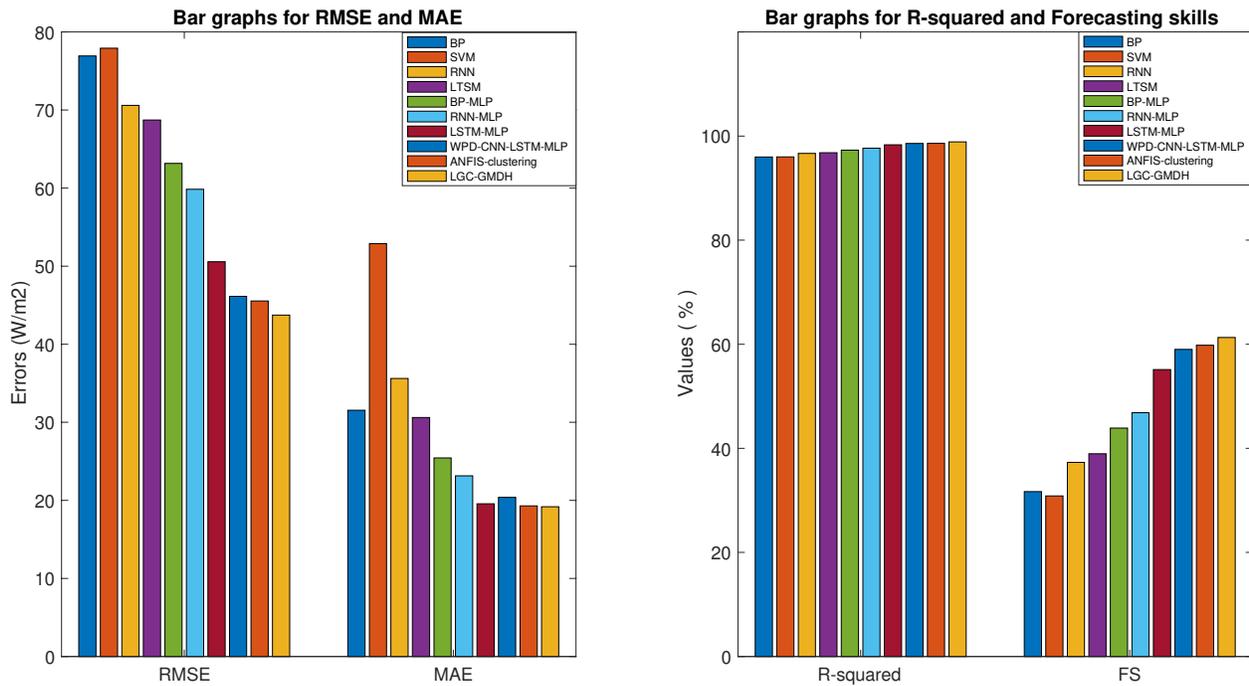
641

642

643

644

645



646 Fig. 10 Performance evaluation of different models for the region of Denver ,1st January 2016 to 31st December
 647 2016 using BP, SVM , RNN , LSTM , BP-MLP, RNN-MLP- LSTM-MLP, WPD-CNN-LTSM-MLP, ANFIS
 648 clustering and LGC-GMDH.

649 As can be seen, the prediction accuracy of single models, such as BP, SVM, RNN, and LSTM, is not
 650 favourable. However, the prediction error is significantly reduced once the models are integrated with the other
 651 network structure. The lowest RMSE values for single models equal 68.7213 W/m² for LSTM, and a FS equals to
 652 0.3167 for the BP model. In comparison, these errors are much lower in hybrid models with an RMSE equal to 46.1336
 653 W/m² for WPD-CNN-LSTM-MLP with a FS equal to 0.590. Furthermore, the category of clustering hybrid models
 654 gives more accuracy. The errors have been improved in the ANFIS-clustering model, which, RMSE, MAE, R², and
 655 FS are 45.5463 W/m², 19.2985 W/m², 98.60% and 59.8%, respectively. However, our proposed LGC-GMDH model
 656 has better performance than all mentioned ones, with an RMSE, MAE, R², and FS of 43.7268 W/m², 19.1856 W/m²,
 657 98.87% and 61.3%, respectively. This comparison concludes that the LGC-GMDH can be used to forecast the hourly
 658 global solar irradiation data, with condition to choose the best configuration for better forecasting performance.

659

660 **7. Conclusion**

661 This paper suggested a new hybrid LGC-GMDH machine learning model for 1-hour horizon global horizontal
662 irradiation forecasting. In this model, we have developed a sophisticated three-phase structure. The data were divided
663 into training and testing data sets using k-cross validation. The clustering using the local gravity algorithm was used
664 in order to get groups with similar dynamic characteristics. To this end, a phase space reconstruction was used in order
665 to build a 2-dimension representation of the global horizontal irradiation data. The obtained clusters were used as
666 inputs to the GMDH network in order to obtain the local forecasts. Then, a global forecast is then reconstructing. The
667 proposed model achieved an excellent forecasting result because it takes into consideration the frequency
668 characteristics of the irradiation time series and the power of the neural networks.

669 Comparing it with other models, the proposed model shows the lowest RMSE error equals 43.7268 W/m², and an R²
670 equals 0.9887, and the highest FS equals 0.613. It also illustrates that network structure has a significant impact on
671 the model's overall accuracy. The LGC-GMDH deep learning model has shown that they provide accurate hour-ahead
672 irradiation forecasts. There will be enough information to find the correct references for practical applications. Future
673 work can test optimisation methods for optimal network configuration choosing to refine results further.

674 **8. Declarations**

675
676 **Data availability** Data used within this research are available upon request from the corresponding author.

677 **Ethics approval** No ethical approval is foreseen since the research described by this paper only concerns the
678 analysis of meteorological data that have been made public by the laboratory that collected them and that no participant
679 has been involved not directly nor indirectly as the subject of the research itself.

680 **Consent to participate** No participant has been involved not directly nor indirectly in the research itself.

681 **Consent for publication** All the authors mentioned in the manuscript have agreed to authorship, read, and
682 approved the manuscript, and given consent for submission and subsequent publication of the manuscript. All named
683 authors have agreed with the order of authorship before submission, and all authors have agreed to the name of the
684 corresponding author. The manuscript has not been published anywhere else.

685 **Competing interests** The authors declare no competing interests.

686 **Funding** No funds, grants, or other support was received to assist with the preparation of this manuscript.

687

688 **Author contributions** The corresponding author [BENMOUIZA Khalil] contributed to the study
689 conception and design. In addition, he performed the material preparation, data collection and analysis. The first
690 draft of the manuscript was written by [BENMOUIZA Khalil]. Also, he read and approved the final manuscript.

691

692 **9. References**

693 Ahmed R, Sreeram V, Mishra Y, Arif MD (2020) A review and evaluation of the state-of-the-art in PV solar power
694 forecasting: Techniques and optimization. *Renew Sustain Energy Rev* 124:109792.

695 <https://doi.org/10.1016/J.RSER.2020.109792>

696 Alkhayat G, Mehmood R (2021) A review and taxonomy of wind and solar energy forecasting methods based on
697 deep learning. *Energy AI* 4:100060. <https://doi.org/10.1016/J.EGYAI.2021.100060>

698 Álvarez-Alvarado JM, Ríos-Moreno JG, Obregón-Biosca SA, et al (2021) Hybrid Techniques to Predict Solar
699 Radiation Using Support Vector Machine and Search Optimization Algorithms: A Review. *Appl Sci* 2021,
700 Vol 11, Page 1044 11:1044. <https://doi.org/10.3390/APP11031044>

701 Auvergne M, M. (1988) Singular value analysis applied to phase space reconstruction of pulsating stars. *A&A*
702 204:341–348

703 Benali L, Notton G, Fouilloy A, et al (2019) Solar radiation forecasting using artificial neural network and random
704 forest methods: Application to normal beam, horizontal diffuse and global components. *Renew Energy*
705 132:871–884. <https://doi.org/10.1016/J.RENENE.2018.08.044>

706 Benmouiza Khalil (2015) Quantification of solar radiation in Algeria, application to the sizing of photovoltaic
707 systems. University of Tlemcen

708 Benmouiza K, Cheknane A (2013) Forecasting hourly global solar radiation using hybrid k-means and nonlinear
709 autoregressive neural network models. *Energy Convers Manag* 75:.
710 <https://doi.org/10.1016/j.enconman.2013.07.003>

711 Benmouiza K, Cheknane A (2018) Clustered ANFIS network using fuzzy c-means, subtractive clustering, and grid
712 partitioning for hourly solar radiation forecasting. *Theor Appl Climatol* 1–13. [https://doi.org/10.1007/s00704-](https://doi.org/10.1007/s00704-018-2576-4)
713 018-2576-4

714 Blanc P, Remund J, Vallance L (2017) Short-term solar power forecasting based on satellite images. *Renew Energy*
715 Forecast From Model to Appl 179–198. <https://doi.org/10.1016/B978-0-08-100504-0.00006-8>

716 Botchkarev A (2019) A new typology design of performance metrics to measure errors in machine learning
717 regression algorithms. *Interdiscip J Information, Knowledge, Manag* 14:45–76. <https://doi.org/10.28945/4184>

718 Burianek T, Misak S (2016) Solar irradiance forecasting model based on extreme learning machine. *EEEIC 2016 -*
719 *Int Conf Environ Electr Eng*. <https://doi.org/10.1109/EEEIC.2016.7555445>

720 C. O (2008) Design of hybrid differential evolution and group method of data handling networks for modeling and
721 prediction. *Inf Sci Comput Sci Intell Syst Appl An Int J* 178:3616–3634.
722 <https://doi.org/10.1016/J.INS.2008.05.013>

723 Caldas M, Alonso-Suárez R (2019) Very short-term solar irradiance forecast using all-sky imaging and real-time
724 irradiance measurements. *Renew Energy* 143:1643–1658. <https://doi.org/10.1016/J.RENENE.2019.05.069>

725 Chu Y, Li M, Coimbra CFM (2016) Sun-tracking imaging system for intra-hour DNI forecasts. *Renew Energy*
726 96:792–799. <https://doi.org/10.1016/J.RENENE.2016.05.041>

727 DT K, L G (1992) Direct test for determinism in a time series. *Phys Rev Lett* 68:427–430.
728 <https://doi.org/10.1103/PHYSREVLETT.68.427>

729 Faceira J, Afonso P, Salgado P (2015) Prediction of Solar Radiation Using Artificial Neural Networks. *Lect Notes*
730 *Electr Eng* 321 LNEE:397–406. https://doi.org/10.1007/978-3-319-10380-8_38

731 Farlow SJ (1981) The GMDH Algorithm of Ivakhnenko. *Am Stat* 35:210. <https://doi.org/10.2307/2683292>

732 Fraser A, Swinney H (1986) Independent coordinates for strange attractors from mutual information. *Phys Rev A,*
733 *Gen Phys* 33:1134–1140

734 Gan M, Huang Y, Ding M, et al (2012) Testing for nonlinearity in solar radiation time series by a fast surrogate data
735 test method. *Sol Energy* 86:2893–2896. <https://doi.org/10.1016/j.solener.2012.04.021>

736 Ghayekhloo M, Ghofrani M, Menhaj MB, Azimi R (2015) A novel clustering approach for short-term solar
737 radiation forecasting. *Sol Energy* 122:1371–1383. <https://doi.org/10.1016/J.SOLENER.2015.10.053>

738 Grassberger P, Procaccia I (1983) Characterization of Strange Attractors. *Phys Rev Lett* 50:346.

739 <https://doi.org/10.1103/PhysRevLett.50.346>

740 Guermoui M, Melgani F, Gairaa K, Mekhalfi ML (2020) A comprehensive review of hybrid models for solar
741 radiation forecasting. *J Clean Prod* 258:120357. <https://doi.org/10.1016/J.JCLEPRO.2020.120357>

742 Huang X, Li Q, Tai Y, et al (2021) Hybrid deep neural model for hourly solar irradiance forecasting. *Renew Energy*
743 171:1041–1060. <https://doi.org/10.1016/J.RENENE.2021.02.161>

744 Kennel MB, Brown R, Abarbanel HDI (1992) Determining embedding dimension for phase-space reconstruction
745 using a geometrical construction. *Phys Rev A* 45:3403–3411. <https://doi.org/10.1103/PhysRevA.45.3403>

746 Khalil B, Ali C (2016) Density-based spatial clustering of application with noise algorithm for the classification of
747 solar radiation time series. In: 2016 8th International Conference on Modelling, Identification and Control
748 (ICMIC). IEEE, pp 279–283

749 Klipp E, Herwig R, Kowald A, et al (2005) *Systems Biology in Practice*. Wiley-VCH Verlag GmbH & Co. KGaA,
750 Weinheim, FRG

751 Kohavi R (1995) A study of cross-validation and bootstrap for accuracy estimation and model selection. 1137–1143

752 Kondo T (1998) GMDH neural network algorithm using the heuristic self-organization method and its application to
753 the pattern identification problem. *Proc SICE Annu Conf* 1143–1148.
754 <https://doi.org/10.1109/SICE.1998.742993>

755 Kumari P, Toshniwal D (2021) Long short term memory–convolutional neural network based deep hybrid approach
756 for solar irradiance forecasting. *Appl Energy* 295:117061. <https://doi.org/10.1016/J.APENERGY.2021.117061>

757 Lai CS, Zhong C, Pan K, et al (2021) A deep learning based hybrid method for hourly solar radiation forecasting.
758 *Expert Syst Appl* 177:114941. <https://doi.org/10.1016/J.ESWA.2021.114941>

759 Laopaiboon T, Ongsakul W, Panyainkaew P, Sasidharan N (2019) Hour-Ahead Solar Forecasting Program Using
760 Back Propagation Artificial Neural Network. *Proc Conf Ind Commer Use Energy, ICUE 2018-October*:
761 <https://doi.org/10.23919/ICUE-GESD.2018.8635756>

762 MacQueen J (1967) Some methods for classification and analysis of multivariate observations. In: *Proceedings of*
763 *the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*. The Regents

764 of the University of California

765 Malakar S, Goswami S, Ganguli B, et al (2021) Designing a long short-term network for short-term forecasting of
766 global horizontal irradiance. *SN Appl Sci* 2021 34 3:1–15. <https://doi.org/10.1007/S42452-021-04421-X>

767 Mazorra Aguiar L, Pereira B, David M, et al (2015) Use of satellite data to improve solar radiation forecasting with
768 Bayesian Artificial Neural Networks. *Sol Energy* 122:1309–1324.
769 <https://doi.org/10.1016/j.solener.2015.10.041>

770 Michel O, Flandrin P (1996) Application of methods based on higher-order statistics for chaotic time series analysis.
771 *Signal Processing* 53:133–148. [https://doi.org/10.1016/0165-1684\(96\)00082-5](https://doi.org/10.1016/0165-1684(96)00082-5)

772 Nazerfard E, Shouraki SB, Hakami V (2006) Evolutionary GMDH-based Identification of Building Blocks for
773 Binary-Coded Systems. 1900–1904. <https://doi.org/10.1109/ICTTA.2006.1684679>

774 Onwubolu G (2016) GMDH-Methodology and Implementation in MATLAB. *GMDH-Methodology Implement*
775 *MATLAB*. <https://doi.org/10.1142/P982>

776 Packard NH, Crutchfield JP, Farmer JD, Shaw RS (1980) Geometry from a Time Series. *Phys Rev Lett* 45:712.
777 <https://doi.org/10.1103/PhysRevLett.45.712>

778 Perez R, Lorenz E, Pelland S, et al (2013) Comparison of numerical weather prediction solar irradiance forecasts in
779 the US, Canada and Europe. *Sol Energy* 94:305–326. <https://doi.org/10.1016/J.SOLENER.2013.05.005>

780 Premalatha N, Valan Arasu A (2016) Prediction of solar radiation for solar systems by using ANN models with
781 different back propagation algorithms. *J Appl Res Technol* 14:206–214.
782 <https://doi.org/10.1016/J.JART.2016.05.001>

783 Refaeilzadeh P, Tang L, Liu H (2009) Cross-Validation. *Encycl Database Syst* 532–538.
784 https://doi.org/10.1007/978-0-387-39940-9_565

785 Schmidt T, Kalisch J, Lorenz E, Heinemann D (2016) Evaluating the spatio-temporal performance of sky-imager-
786 based solar irradiance analysis and forecasts. *Atmos Chem Phys* 16:3399–3412. [https://doi.org/10.5194/acp-
787 16-3399-2016](https://doi.org/10.5194/acp-16-3399-2016)

788 Schüler D, Wilbert S, Geuder N, et al (2016) The enerMENA meteorological network – Solar radiation

789 measurements in the MENA region. AIP Conf Proc 1734:150008. <https://doi.org/10.1063/1.4949240>

790 Shadab A, Said S, Ahmad S (2019) Box–Jenkins multiplicative ARIMA modeling for prediction of solar radiation: a
791 case study. Int J Energy Water Resour 2019 34 3:305–318. <https://doi.org/10.1007/S42108-019-00037-5>

792 Soubdhan T, Ndong J, Ould-Baba H, Do MT (2016) A robust forecasting framework based on the Kalman filtering
793 approach with a twofold parameter tuning procedure: Application to solar and photovoltaic prediction. Sol
794 Energy 131:246–259. <https://doi.org/10.1016/J.SOLENER.2016.02.036>

795 Tadj M, Benmouiza K, Cheknane A, Silvestre S (2014) Improving the performance of PV systems by faults
796 detection using GISTEL approach. Energy Convers Manag 80:298–304.
797 <https://doi.org/10.1016/J.ENCONMAN.2014.01.030>

798 Takens F (1981) Detecting strange attractors in turbulence. 366–381. <https://doi.org/10.1007/BFB0091924>

799 Vaishnav V, Vajpai J (2018) Seasonal Time Series Forecasting by Group Method of Data Handling. 2018 IEEE Int
800 Students' Conf Electr Electron Comput Sci SCEECS 2018. <https://doi.org/10.1109/SCEECS.2018.8546886>

801 Verbois H, Huva R, Rusydi A, Walsh W (2018) Solar irradiance forecasting in the tropics using numerical weather
802 prediction and statistical learning. Sol Energy 162:265–277. <https://doi.org/10.1016/J.SOLENER.2018.01.007>

803 Vindel JM, Polo J (2014) Markov processes and Zipf's law in daily solar irradiation at earth's surface. J Atmos
804 Solar-Terrestrial Phys 107:42–47. <https://doi.org/10.1016/J.JASTP.2013.10.017>

805 Voyant C, Muselli M, Paoli C, Nivet ML (2012) Numerical weather prediction (NWP) and hybrid ARMA/ANN
806 model to predict global radiation. Energy 39:341–355. <https://doi.org/10.1016/j.energy.2012.01.006>

807 Voyant C, Notton G, Kalogirou S, et al (2017) Machine learning methods for solar radiation forecasting: A review.
808 Renew Energy 105:569–582. <https://doi.org/10.1016/J.RENENE.2016.12.095>

809 W.M.O. (1981) Meteorological aspects of the utilization of solar radiation as an energy source, illustrate. Secretariat
810 of the World Meteorological Organization

811 Wang B, Che J, Wang B, Feng S (2019) A Solar Power Prediction Using Support Vector Machines Based on Multi-
812 source Data Fusion. 2018 Int Conf Power Syst Technol POWERCON 2018 - Proc 4573–4577.
813 <https://doi.org/10.1109/POWERCON.2018.8601672>

- 814 Wang Z, Yu Z, Philip Chen CL, et al (2018) Clustering by Local Gravitation. *IEEE Trans Cybern* 48:1383–1396.
815 <https://doi.org/10.1109/TCYB.2017.2695218>
- 816 Water PR, Kerckhoffs EJH, Van Welden D (2000) GMDH-based dependency modeling in the identification of
817 dynamic systems. In: *Proceedings of the 14th European Simulation Multiconference on Simulation and*
818 *Modelling: Enablers for a Better Quality of Life*. Society for Computer Simulation International, San Diego,
819 pp 211–218
- 820 Whitney H (1936) Differentiable manifold. *Ann Math* 37:645–680. <https://doi.org/10.2307/1968482>
- 821 Wong TT (2015) Performance evaluation of classification algorithms by k-fold and leave-one-out cross validation.
822 *Pattern Recognit* 48:2839–2846. <https://doi.org/10.1016/J.PATCOG.2015.03.009>
- 823 Yadav AP, Behera L (2014) Solar Radiation forecasting using neural networks and Wavelet Transform. *IFAC Proc*
824 *Vol 47:890–896*. <https://doi.org/10.3182/20140313-3-IN-3024.00218>
- 825 Zhang N, Behera PK (2012) Solar radiation prediction based on recurrent neural networks trained by Levenberg-
826 Marquardt backpropagation learning algorithm. *2012 IEEE PES Innov Smart Grid Technol ISGT 2012*.
827 <https://doi.org/10.1109/ISGT.2012.6175757>
- 828