

Proceeding Paper

# Time Series Chlorophyll-A Concentration Data Analysis: A Novel Forecasting Model for Aquaculture Industry <sup>†</sup>

Elias Eze <sup>1,\*</sup> , Sam Kirby <sup>2</sup>, John Attridge <sup>2</sup> and Tahmina Ajmal <sup>1</sup> 

<sup>1</sup> Institute for Research in Applicable Computing (IRAC), School of Computer Science and Technology, University of Bedfordshire, Luton LU1 3JU, UK; tahmina.ajmal@beds.ac.uk

<sup>2</sup> Chelsea Technology Group, 55 Central Avenue, West Molesey, Surrey KT8 2QZ, UK; skirby@chelsea.co.uk (S.K.); jattridge@chelsea.co.uk (J.A.)

\* Correspondence: elias.eze1@beds.ac.uk

<sup>†</sup> Presented at the 7th International conference on Time Series and Forecasting, Gran Canaria, Spain, 19–21 July 2021.

**Abstract:** Eutrophication in fresh water has become a critical challenge worldwide and chlorophyll-a content is a key water quality parameter that indicates the extent of eutrophication and algae concentration in a body of water. In this paper, a forecasting model for the high accuracy prediction of chlorophyll-a content is proposed to enable aquafarm managers to take remediation actions against the occurrence of toxic algal blooms in the aquaculture industry. The proposed model combines the ensemble empirical mode decomposition (EEMD) technique and a deep learning (DL) long short-term memory (LSTM) neural network (NN). With this hybrid approach, the time-series data are firstly decomposed with the aid of the EEMD algorithm into manifold intrinsic mode functions (IMFs). Secondly, a multi-attribute selection process is employed to select the group of IMFs with strong correlations with the measured real chlorophyll-a dataset and integrate them as inputs for the DL LSTM NN. The model is built on water quality sensor data collected from the Loch Duart salmon aquafarm in Scotland. The performance of the proposed novel hybrid predictive model is validated by comparing the results against the dataset. To measure the overall accuracy of the proposed novel hybrid predictive model, the Mean Absolute Error (MAE), Mean Square Error (MSE), Root Mean Square Error (RMSE), and Mean Absolute Percentage Error (MAPE) were used.

**Keywords:** water quality; aquaculture; forecasting; chlorophyll-a time-series data; deep learning LSTM

check for  
updates

**Citation:** Eze, E.; Kirby, S.; Attridge, J.; Ajmal, T. Time Series Chlorophyll-A Concentration Data Analysis: A Novel Forecasting Model for Aquaculture Industry. *Eng. Proc.* **2021**, *2*, 5027. <https://doi.org/10.3390/engproc2021005027>

Academic Editors: Ignacio Rojas, Fernando Rojas, Luis Javier Herrera and Hector Pomare

Published: 29 June 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Eutrophication in freshwater bodies is an organic process usually caused by the increased enrichment of nutrients which can pollute water quality and adversely affect aquatic ecosystems. The extent of eutrophication in fresh water can be estimated through chlorophyll-a concentration monitoring. In the aquaculture industry, this natural process of nutrient enrichment also results in structural changes to the aquatic ecosystem through increased algae production, the depletion of fish species, and the prevalent degradation of overall water quality [1,2]. Chlorophyll-a concentration is representative of the state of freshwater quality and has generally been used as a key indicator for measuring algal blooms [3].

According to Gao and Zhang [4], eutrophication has become a ubiquitous fresh-water-quality pollutant in China. Similarly, a study conducted by Jules et al. [5] estimated the annual damage costs of the eutrophication of fresh water in England and Wales to be \$105–160 million (£75.0–114.3 m). Given the link between the adverse effect of eutrophication in freshwater and the stagnation of wild fishery populations, the aquaculture industry has emerged as a crucial means of providing protein to our constantly growing population. Therefore, the monitoring of water quality parameters (for instance, algal

biomass and cyanobacteria) through chlorophyll-a concentration is increasingly favoured over laboratory analysis and similar traditional methods because of the high cost and labour-intensive requirements associated with them [6]. The effective monitoring and prediction of chlorophyll-a concentrations is a promising approach for the routine estimation of phytoplankton biomass in the aquaculture ecosystems of the Nile tilapia (*Oreochromis niloticus*) [7]. Sensory monitoring of the chlorophyll-a concentration is an effective approach for reliably assessing the trophic state of freshwater bodies given its strong affinity to the abundance of phytoplankton, cyanobacteria, and biomass, which affect the turbidity and general colouration of fresh water [8].

Several studies have been conducted to establish a means of coping with water quality impairments caused by algal biomass using conventional numerical modelling methods, least squares support vector regression (LSSVR), neural networks methods such as Radial Basis Function neural network (RBFNN), Back Propagation neural network (BPNN) algorithms, and machine learning methods to predict chlorophyll-a concentrations as an indicator for future water quality changes [9–12]. However, the challenge with traditional numerical methods, LSSVR, and neural networks such as RBFNN and BPNN is the inherent weakness of the long-term dependency problem. Research has shown that deep learning long short-term memory (LSTM) neural networks can overcome the above-mentioned weakness and can provide efficient applicability and reliability for water quality parameter prediction [13,14]. Additionally, combining the ensemble empirical mode decomposition (EEMD) method with deep learning LSTM neural network has demonstrated clear advantages over traditional LSTM neural networks in terms of improved water quality parameter prediction accuracy in the aquaculture environment [13]. In this paper, a novel deep learning-based hybrid chlorophyll-a prediction model for the aquaculture industry is proposed.

## 2. Data Source

### 2.1. The Study Area Description and Datasets Analysis

Loch Duart is an independent Scottish salmon aquafarm industry, which has its headquarters in Scourie, Sutherland, in north-west Scotland. The salmon farming company owns and operates eight sea-sites and two hatcheries in Sutherland and the Outer Hebrides. In Loch Duart, salmon are hatched and grown in the cold, clear fresh water of north-west Scotland. The salmon farming company annually harvests approximately 5000 tons of fresh salmon. Chlorophyll-a ( $\mu\text{g/L}$ ) time-series data were collected via a TriLux multi-parameter sensor probe. The sensor deployment took place at one of their sheltered sites along the coast (see Figure 1a). The telemetry unit was secured to the metal walkway around the outside of the net pens and the sensor was situated on the outside of one of the outermost pens, nearest to the feed barge.

A TriLux multi-parameter fluorometer/sensor (see Figure 1b) developed by Chelsea Technology Group was used for measuring and collecting a total of 22,708 sets of a non-linear, non-stationary water-quality parameter time-series dataset at Loch Duart salmon aquafarm between May and October 2020. The water quality parameters include chlorophyll-a (470), turbidity, and chlorophyll-a (530).

Generally, the 470 channel measures chlorophyll fluorescence from direct excitation of chlorophyll-a that usually strongly correlates with phytoplankton biomass in freshwater. Table 1 shows the list of other sensors developed by Chelsea Technology Group and the corresponding parameters that each of them measures.



Table 1. Cont.

		Fluorometers						Active Fluorometers					Optical Sensors			
		UniLux	TriLux	UviLux	VLux AlgaePro	VLux TPro	VLux FuelPro	VLux OilPro	LabSTAF	FastOcean APD	FastOcean	Act2 Lab	FastBallast	PAR Sensor	GlowTracka	UniLux Turbidity
Active Fluorometers	Variable Fluorescence							✗	✗	✗	✗	✗				
	Fluorescence Light Curves (FLC)							✗		✗						
	Phytoplankton Primary Productivity							✗	✗	✗						
	Phytoplankton Cell Counting												✗			
Optical Sensors	PAR													✗		
	Bioluminescence														✗	
	Turbidity	✗	✗		✗	✗	✗	✗								✗
	Absorbance							✗	✗	✗						

2.2. Data Pre-Treatment, Filling and Correction

Water-quality parameter time-series dataset defects usually result in excessive deviation between the measured original water-quality parameter values and the prediction results. The basis of accurate time-series analysis and the development of effective and reliable predictive models is high-quality sample data. To provide a concise, accurate dataset for the prediction model and improve prediction accuracy, the measured water-quality parameter data was carefully pre-processed. Generally, the issue of missing data is often inevitable with automatic water quality monitoring systems. The water-quality parameters like turbidity, chlorophyll-a (470), and chlorophyll-a (530) were automatically measured for 10 months at 10 min intervals. To fill in any missing data, a filling-in approach called linear interpolation algorithm [16] is applied to achieve a better estimation effect that can accurately approximate the missing data values. In data analysis, a linear interpolation algorithm assumes the ratio of two separate known data and a single unknown datum to be a linear interrelation. Therefore, to obtain the missing, unknown water quality parameter value, the linear interpolation technique applies the slope of the presumed line to compute the time-series dataset increment.

**Definition 1.** Time series nature of the measured parameter (Chlorophyll-a (470)).

The automated water quality sensory system at Loch Duart salmon aquafarm measures the time series water quality parameters at a constant time interval everyday which can be denoted as  $\beta$ , so that  $n$  length time-series of the measured parameters' datasets is defined as (1)

$$S_{i,n} = \{(X_{i,1}, T_1), (X_{i,2}, T_2), \dots, (X_{i,n}, T_n)\} \tag{1}$$

where  $X_{i,l}$  represents the values of the measured  $i^{th}$  time-series water-quality parameters by the automatic sensory system at time  $T_l$  ( $1 \leq i \leq \beta, 1 \leq l \leq n$ ), so that for a given  $T_l$ , the sampling time interval is constant at  $\Delta T = (T_{l+1} - T_l) = 5$  min. Therefore, if the original value  $X_{i,l}$  is missing, its estimated value  $\hat{X}_{i,l}$  can be obtained with the problem of minimum, which is given as  $|\hat{X}_{i,l} - X_{i,l}|$ , changed into the missing value estimation

problem. Based on the measured data  $X_{i,x}$  and  $X_{i,y}$  at time  $T_{i,x}$  and  $T_{i,y}$ , respectively, the linear imputation function  $L(t)$  could be formulated for the time-series water-quality parameter monitoring systems as:

$$L(t) = X_{i,x} + \left( \frac{X_{i,x} - X_{i,y}}{T_{i,x} - T_{i,y}} \right) \cdot (t - T_{i,x}). \tag{2}$$

For any missing time-series water-quality parameter data at any given moment, the linear interpolation algorithm firstly finds the two closest moments  $T_{i,x}$  and  $T_{i,y}$  ( $T_{i,x} < t < T_{i,y}$ ), and estimates the lost data value at time  $t$  with the help of the known measured data  $X_{i,x}$  and  $X_{i,y}$  of  $T_{i,x}$  and  $T_{i,y}$  moments based on Equation (2), i.e.,  $\hat{X}_n = L(t)$ .

### 3. Proposed Model

The EEMD technique and deep learning LSTM NN were merged to form the chlorophyll-a hybrid prediction model. A detailed implementation processes of the applied EEMD technique is shown in full in [13]. The LSTM deep learning NN approach is described in full detail in Section 3.1. The original chlorophyll-a (470) dataset is decomposed effectively by the application of the EEMD technique into  $n$  disparate IMFs and a residual item. The IMF components that are contained within individual frequency bands are independently different and usually change with the variation of the chlorophyll-a (470) time-series data  $x(t)$ . Likewise, the trend of  $x(t)$  is generally demonstrated by the corresponding ensemble residual item as the output of the decomposition process implementation.

#### 3.1. Deep Learning LSTM Neural Networks

Deep learning LSTM NNs are a special type of recurrent NN (RNN) with significant improvement in the ability to learn long-term dependencies which gives it an advantage over other artificial neural networks such as BPNN and RBFNN. Figure 2a illustrates a typical schematic diagram of a traditional RNN node with the previous hidden state represented by  $h_{t-1}$ , activation tanh function, current input sample by  $X_t$ , current output by  $h_t$ , and the current hidden state by  $h_t$ . As depicted in Figure 2, all RNNs generally have the form of a chain of repeating modules of NNs. These repeating modules generally have a very basic structure in standard RNNs, like a single tanh layer only. However, a deep learning LSTM which stores information with the aid of purpose-built memory cells maintains similar chain-like structure, but with a differently structured repeating module (see Figure 2b).

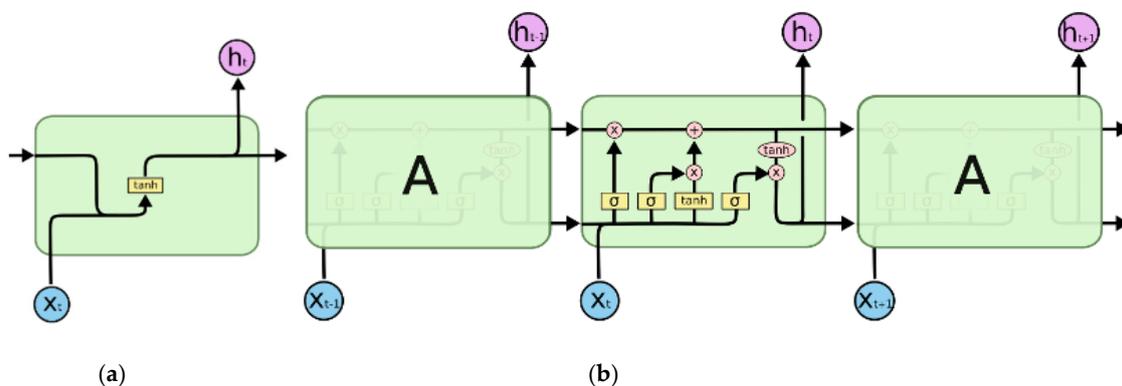


Figure 2. (a,b): Typical schematic diagram of (a) Traditional RNN node, and (b) Chained LSTM blocks.

The equations below illustrate the calculation processes involved in deep learning LSTM NNs.

(a) Forget gate equation:

$$F_t = \sigma(W_f \times [h_{t-1}, X_t] + b_f) \quad (3)$$

where  $F_t$  represents a vector that has a range from 0 to 1 as its values;  $W_f$ ,  $\sigma$ , and  $b_f$  represent the weight matrices, sigmoid function, and the bias of forget gate, respectively. The  $\sigma$  is used to find out whether the new information is unnecessary, in which case the information ignored and discarded, or necessary and used for updating. Finally, the tanh function is used to add weight to individual values that pass and determines their level of relevance, and ranges from  $-1$  to  $1$ . Inside the input gate and the output gate, same operations are repeated, which are shown in (4)–(7).

(b) Input gate equations:

$$I_t = \sigma(W_i \times [h_{t-1}, X_t] + b_i) \quad (4)$$

$$\hat{I}_t = \tanh(W_i \times [h_{t-1}, X_t] + b_i) \quad (5)$$

(c) Output gate equations:

$$O_t = \sigma(W_o \times [h_{t-1}, X_t] + b_o) \quad (6)$$

$$h_t = O_t \times \tanh(C_t) \quad (7)$$

(d) Cell state equation:

$$C_t = \{(F_t \times C_{t-1}) + (I_t \times \hat{I}_t)\} \quad (8)$$

where  $W_i$  and  $W_o$  denote the weight matrixes,  $b_i$  and  $b_o$  denote the bias vectors of the network of both input gate and output gate, and the hyperbolic tangent function is denoted by the tanh function.

### 3.2. Proposed Water Quality Prediction Model

The proposed hybrid EEMD-LSTM deep learning NN-based water-quality parameter prediction model is depicted in Figure 3. With the proposed novel water quality forecasting model, the measured real water-quality parameter content dataset undergoes decomposition processes into disparate components by applying the EEMD method for the purpose of improving the prediction accuracy of the proposed predictive model. The full procedures demonstrated in Figure 3 show the three important steps which were followed in developing the novel hybrid water quality parameters prediction solution. Firstly, the water quality parameters dataset  $x(t)$  generates multiple, distinct IMF components and a corresponding residue  $R_N(t)$  from the decomposition processes via the applied EEMD method in the input layer of Figure 3. The decomposition of  $x(t)$  is carried out by means of an iterative sifting procedure as given below:

$$x(t) = \sum_{i=1}^N IMF_i(t) + R_N(t) \quad (9)$$

Subsequently, the separate IMF components and their corresponding residue undergo a process of normalization in the second step and are then used for prediction by the DL LSTM in the hidden layer of Figure 3. Lastly, in step three, individual prediction results undergo a reverse normalization process before they are efficiently combined together with the aid of a summation operation by the summation function to get the final predicted values as shown in the output layer of Figure 3. As clearly illustrated using the extended forecasting model with multiple hidden DL LSTM layers ( $LSTM_{1,1}$ ,  $LSTM_{1,2}$ ,  $\dots$ ,  $LSTM_{m,1}$ ,

up to LSTM<sub>m,n</sub>) in Figure 3, individual hidden layers of the stacked DL LSTM are equipped with multiple memory cells which earn the proposed prediction model the name ‘deep learning’ NN [17].

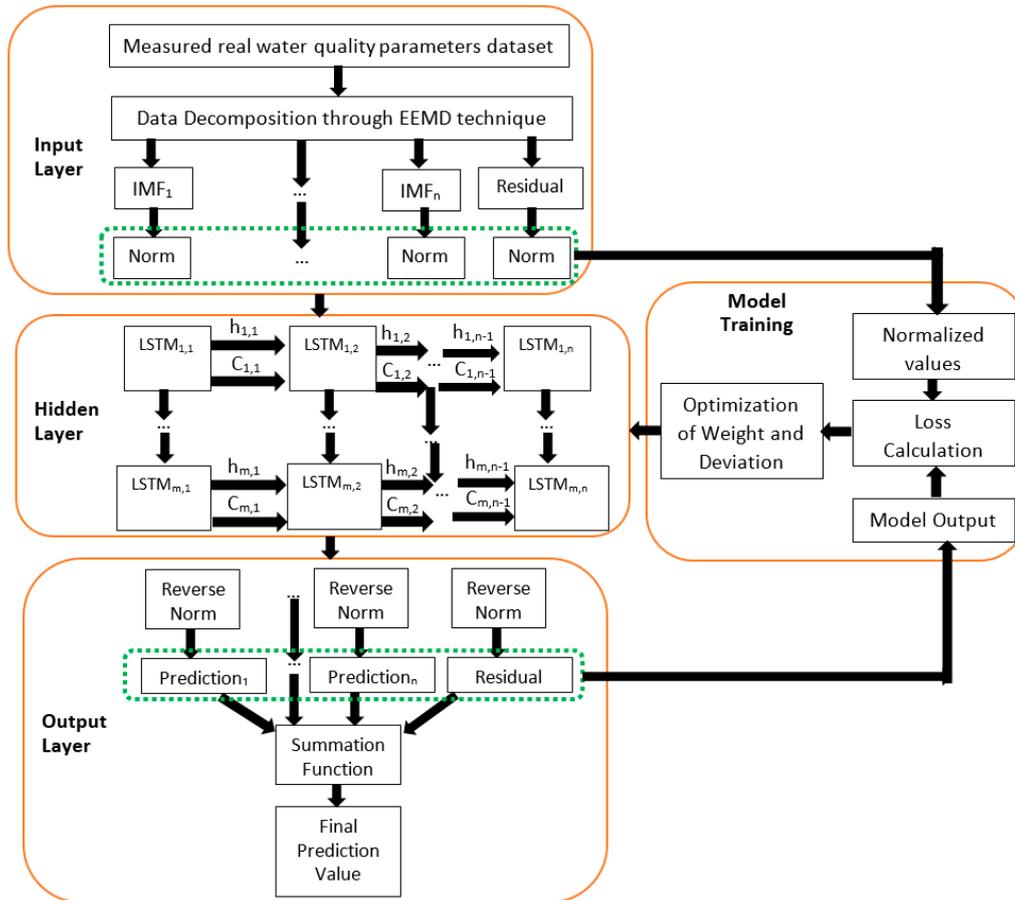


Figure 3. Proposed hybrid EEMD–LSTM deep learning water quality prediction Model.

#### 4. Performance Evaluation

For the evaluation of the proposed hybrid EEMD–LSTM deep learning water-quality prediction model, four performance evaluation metrics were introduced to evaluate its prediction accuracy. These metrics include MAE, MSE, RMSE, and MAPE. The mathematical formulae are expressed as follows:

$$MAE = \frac{1}{n} \sum_{i=1}^n |M_i - F_i| \tag{10}$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (M_i - F_i)^2 \tag{11}$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (M_i - F_i)^2} \tag{12}$$

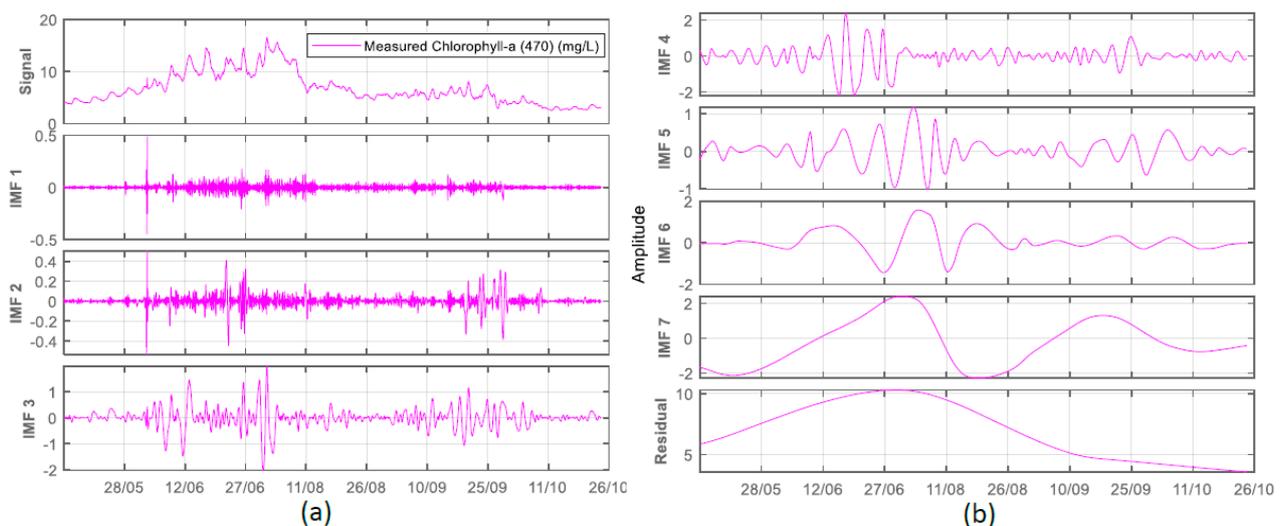
$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{M_i - F_i}{M_i} \right|. \tag{13}$$

In (10)–(13) above,  $n$  denotes the number of data points in the dataset, and  $V_i$  and  $F_i$  represent the measured real chlorophyll-a values and the forecasted values, respectively.

The closer these four performance evaluation metrics tend towards 0, the higher the overall forecasting and fitting accuracy of the proposed solution.

## 5. Results and Discussions

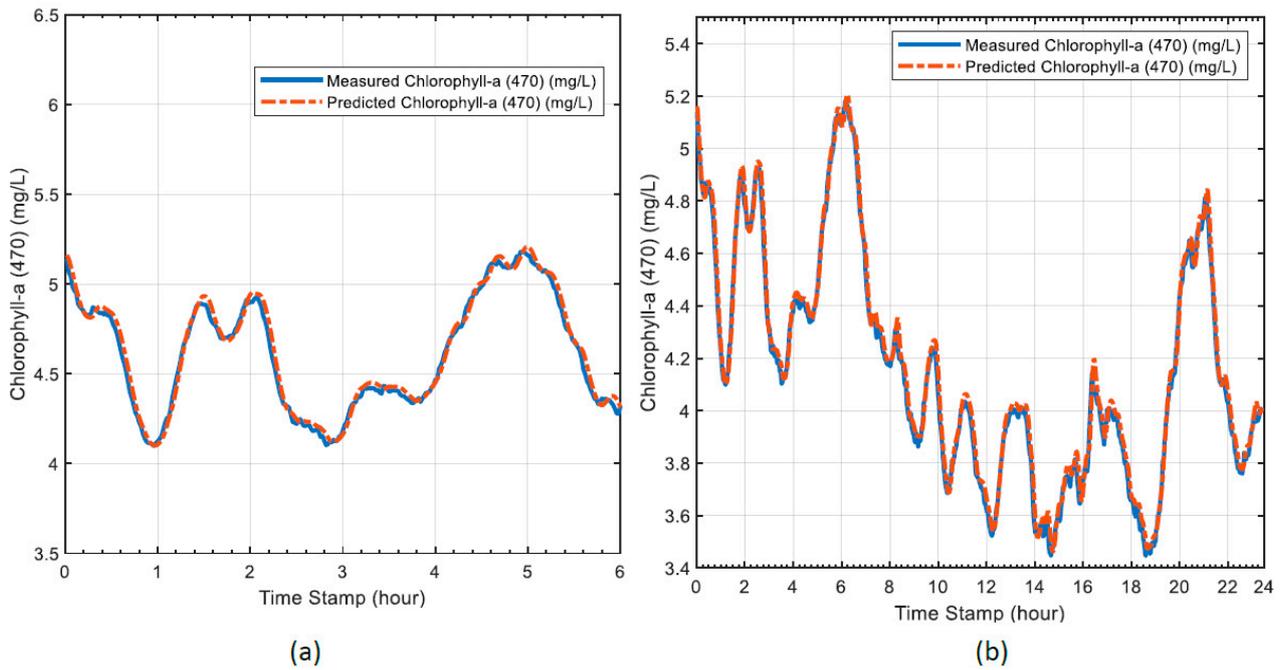
In this study, decomposing the Chelsea's TriLux multiparameter fluorometer measured chlorophyll-a content dataset is an intrinsic aspect of the novel prediction model for ensuring high short-term prediction accuracy. The EEMD method decomposes the real chlorophyll-a content dataset into seven individually stable IMF components (IMF 1–7) and one residual item as depicted in Figure 4a,b. The obtained IMFs from the original chlorophyll-a (470) dataset decomposition with the EEMD method is shown in Figure 4a,b.



**Figure 4.** (a,b). Chlorophyll-a (470) dataset decomposition through the EEMD method showing (a) 1 to 3 of the resultant 7 IMFs, and (b) 4 to 7 of the resultant 7 IMFs.

The graphs in Figure 5a,b clearly show that the novel hybrid forecasting model provided good results for short-term (6 h) and long-term (24 h) forecast scenarios. With chlorophyll-a (470) concentration data, the matching trends in both Figure 5a,b further show that the model can successfully predict, with a high level of accuracy, the presence of algal bacteria such as cyanobacteria, which is a harmful alga that produces odorous and toxic substances leading to severe problems for different species of fish in the aquaculture industry.

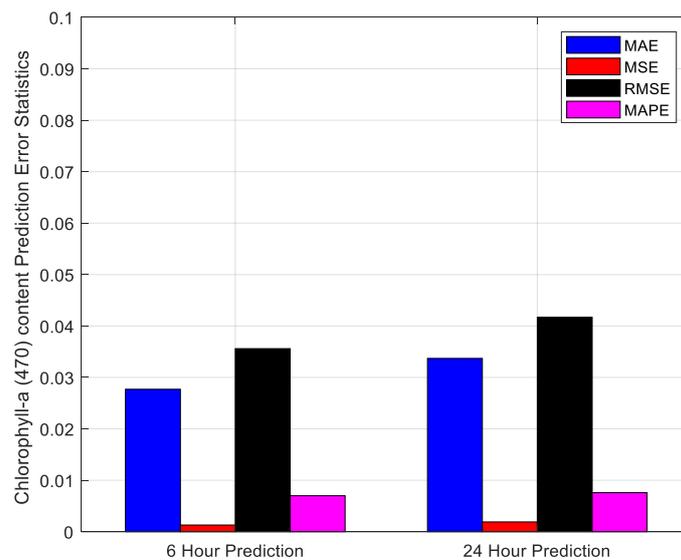
The proposed model improved the prediction accuracy due to the application of the EEMD method, which enabled the predictive model to manifest the temporal features of the chlorophyll-a (470) content time-series data. This was done through the multi-feature selection process of the EEMD method which allowed for the selection of certain groups of IMFs that strongly correlate with the Chelsea's TriLux multi-parameter fluorometer measured chlorophyll-a data and integrate them into inputs for the deep learning LSTM neural network. Table 2 and Figure 6 present the error statistics for both 6 h and 24 h forecast results. Although these are minimal errors, the overall prediction accuracy could be further improved with an increase in data availability because the deep learning LSTM chain structure tends to be more complex and performs better with big data.



**Figure 5.** (a,b). Performance comparison of real Chlorophyll-a (470) parameter values and the predicted values: (a) half-day (6 h), and (b) one day (24 h) prediction results.

**Table 2.** Error statistics for 6 h and 24 h chlorophyll-a (470) content prediction.

Error Statistics	6 Hour Prediction	24 Hour Prediction
MSE	0.0013	0.0019
MAE	0.0277	0.0337
RMSE	0.0356	0.0417
MAPE	0.0070	0.0076



**Figure 6.** Chlorophyll-a (470) content prediction error statistics for 6 h and 24 h.

### 6. Conclusions

Timely prediction of toxic algal blooms with the help of real chlorophyll-a (470) sensor time-series data in aquatic ecosystems can allow for the effective operation and management of the aquaculture industry by providing useful information that can facilitate

the decision-making process in aquafarming. In this study, we present a novel hybrid model to forecast chlorophyll-a content through the combination of the potential of the EEMD technique and a DL LSTM neural network approach. The actual experimental data from Loch Duart Salmon aquafarm show that the proposed model provides impressive results with high prediction accuracy. For future work, varieties of water quality parameter time-series datasets measured from different aquafarming sites will be considered to broaden the application horizon of the proposed forecasting model.

**Funding:** This research was funded by Innovate UK/BBSRC (ref: 86204028, BB/S020896/1).

**Acknowledgments:** The authors wish to thank the anonymous reviewers for their comments and suggestions, which have greatly helped in improving the quality of the paper.

**Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

## References

1. Chislock, M.F.; Doster, E.; Zitomer, R.A.; Wilson, A.E. Eutrophication: Causes, Consequences, and Controls in Aquatic Ecosystems. *Nat. Educ. Knowl.* **2013**, *4*, 1–10.
2. Howarth, R.; Chan, F.; Conley, D.J.; Garnier, J.; Doney, S.C.; Marino, R.; Billen, G. Coupled biogeochemical cycles: Eutrophication and hypoxia in temperate estuaries and coastal marine ecosystems. *Front. Ecol. Environ.* **2011**, *9*, 18–26. [CrossRef]
3. Kim, B.C.; Jung, S.M.; Jang, C.W.; Kim, J.K. Comparison of BOD, COD and TOC as the indicator of organic matter pollution in streams and reservoirs of Korea. *J. Korean Soc. Environ. Eng.* **2007**, *29*, 640–643.
4. Gao, C.; Zhang, T. Eutrophication in a Chinese context: Understanding various physical and socio-economic aspects. *Ambio* **2010**, *39*, 385–393. [CrossRef] [PubMed]
5. Pretty, J.N.; Mason, C.F.; Nedwell, D.B.; Hine, R.E.; Leaf, S.; Dils, R. Environmental Costs of Freshwater Eutrophication in England and Wales. *Environ. Sci. Technol.* **2003**, *37*, 201–208. [CrossRef] [PubMed]
6. Chelsea Technologies. Aquaculture. Available online: <https://chelsea.co.uk/application-category/aquaculture> (accessed on 13 April 2021).
7. El-Otify, A.M. Evaluation of the physicochemical and chlorophyll-a conditions of a subtropical aquaculture in Lake Nasser area, Egypt. *Beni-Suef Univ. J. Basic Appl. Sci.* **2015**, *4*, 327–337. [CrossRef]
8. Ha, N.T.; Koike, K.; Nhuan, M.T. Improved Accuracy of Chlorophyll-a Concentration Estimates from MODIS Imagery Using a Two-Band Ratio Algorithm and Geostatistics: As Applied to the Monitoring of Eutrophication Processes over Tien Yen Bay (Northern Vietnam). *Remote Sens.* **2013**, *6*, 421–442. [CrossRef]
9. Shumway, S.E. A review of the effects of algal blooms on shellfish and aquaculture. *J. World Aquac. Soc.* **1990**, *21*, 65–104. [CrossRef]
10. Shin, Y.; Kim, T.; Hong, S.; Lee, S.; Lee, E.; Hong, S.; Lee, C.; Kim, T.; Park, M.S.; Park, J.; et al. Prediction of chlorophyll-a concentrations in the Nakdong River using machine learning methods. *Water* **2020**, *12*, 1822. [CrossRef]
11. Wang, X.; Wang, G.; Zhang, X. Prediction of Chlorophyll-a content using hybrid model of least squares support vector regression and radial basis function neural networks. In Proceedings of the 2016 Sixth International Conference on Information Science and Technology (ICIST), Dalian, China, 6–8 May 2016; pp. 366–371.
12. Syariz, M.A.; Lin, C.H.; Nguyen, M.V.; Jaelani, L.M.; Blanco, A.C. WaterNet: A convolutional neural network for chlorophyll-a concentration retrieval. *Remote Sens.* **2020**, *12*, 1966. [CrossRef]
13. Eze, E.; Ajmal, T. Dissolved Oxygen Forecasting in Aquaculture: A Hybrid Model Approach. *Appl. Sci.* **2020**, *10*, 7079. [CrossRef]
14. Hu, Z.; Zhang, Y.; Zhao, Y.; Xie, M.; Zhong, J.; Tu, Z.; Liu, J. A water quality prediction method based on the deep LSTM network considering correlation in smart mariculture. *Sensors* **2019**, *19*, 1420. [CrossRef] [PubMed]
15. Chelsea Technologies. TriLux. Available online: <https://chelsea.co.uk/products/trilux/> (accessed on 13 April 2021).
16. Pan, L.; Li, J.; Luo, J. A temporal and spatial correction based missing values imputation algorithm in wireless sensor networks. *Chin. J. Comput.* **2010**, *33*, 1–10. [CrossRef]
17. Jason Brownlee, Stacked Long Short-Term Memory Networks Develop Sequence Prediction Models in Keras. 14 August 2019. Available online: <https://machinelearningmastery.com/stacked-long-short-term-memorynetworks/> (accessed on 19 February 2021).