

1 **An efficient model construction strategy to simulate microalgal lutein photo-production**  
2 **dynamic process**

3

4 Ehecatl Antonio del Rio-Chanona<sup>1,2</sup>, Fabio Fiorelli<sup>1</sup>, Dongda Zhang<sup>2, \*</sup>, Nur rashid Ahmed<sup>3</sup>, Keju  
5 Jing<sup>3</sup>, Nilay Shah<sup>2</sup>

6

7 1: Department of Chemical Engineering and Biotechnology, University of Cambridge, Pembroke  
8 Street, Cambridge CB2 3RA, UK.

9 2: Centre for Process Systems Engineering, Imperial College London, South Kensington Campus,  
10 London SW7 2AZ, UK.

11 3. Department of Chemical and Biochemical Engineering, College of Chemistry and Chemical  
12 Engineering, Xiamen University, Xiamen 361005, China.

13

14 \*: corresponding authors, email: [dongda.zhang11@imperial.ac.uk](mailto:dongda.zhang11@imperial.ac.uk), tel: 44 (0)7543785283.

15

16 **Running title: Dynamic Modelling of Algal Lutein Production**

17

18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40

## Abstract

Lutein is a high-value bioproduct synthesised by microalga *Desmodesmus* sp.. It has great potential for the food, cosmetics, and pharmaceutical industries. However, in order to enhance its productivity and to fulfil its ever-increasing global market demand, it is vital to construct accurate models capable of simulating the entire behaviour of the complicated dynamics of the underlying biosystem. To this aim, in this study two highly robust artificial neural networks are designed for the first time. Contrary to conventional artificial neural networks, these networks model the rate of change of the dynamic system, which makes them highly relevant in practice. Different strategies are incorporated into the current research to guarantee the accuracy of the constructed models, which include determining the optimal network structure through a hyperparameter selection framework, generating significant amounts of artificial data sets by embedding random noise of appropriate size, and rescaling model inputs through standardisation. Based on experimental verification, the high accuracy and great predictive power of the current models for long-term dynamic bioprocess simulation in both real-time and offline frameworks are thoroughly demonstrated. This research, therefore, paves the way to significantly facilitate the future investigation of lutein bioproduction process control and optimisation. In addition, the model construction strategy developed in this research has great potential to be directly applied to other bioprocesses.

**Keywords:** artificial neural network; dynamic simulation; lutein production; real-time framework; fed-batch operation; bioprocess modelling.

## 41 **1. Introduction**

42 The synthesis of valuable bioproducts from microalgae through photosynthetic related metabolic  
43 pathways is a promising sector (Mata et al. 2010) that can deliver a wide variety of commodities.  
44 One of the most well-known applications is the production of eco-friendly biofuels such as  
45 biodiesel and biohydrogen (Doshi et al. 2016; Zhang, Dechatiwongse, del Rio-Chanona, et al.  
46 2015), which are being developed to replace traditional transportation fuels. Another use is that  
47 as food supplements (Río et al. 2005; Xie et al. 2014), a capability in which they have been long  
48 employed and for which considerable growth of market demand is forecasted. Most importantly,  
49 there is a thriving international research interest, development and deployment of microalgae  
50 based sustainable and environmentally friendly technologies, by the health sector. This is with  
51 the aim of producing specialist high-value products such as lutein (Xie et al. 2014), C-  
52 phycocyanin (del Rio-Chanona, Zhang, et al. 2015), and astaxanthin (Zhang et al. 2016), for  
53 which traditional synthesis routes and refinery methods from existing non-renewable sources are  
54 expensive, energy intensive, and of low efficiency (Yen et al. 2011; Capelli et al. 2013).

55  
56 One high-value bioproduct that has found particular attention is lutein, a carotenoid of great  
57 interest to the health, pharmaceutical, cosmetics, and food industries (Yaakob et al. 2014;  
58 Fernández-Sevilla et al. 2010). Lutein has been widely used for the treatment of ophthalmic  
59 conditions and cancer, and applied as a natural colorant in cellular pigmentation and in the food  
60 industry (Ho et al. 2015; Xie et al. 2013). Because of its wide applications, its market demand in  
61 the US was estimated to increase significantly from \$150 million in 2000 to \$309 million in 2018  
62 (Marz 2011). Nonetheless, further growth in production is hampered by the fact that its current  
63 industrial feedstock is marigold, a plant which has an extremely low intracellular content of

64 lutein (0.02~0.1% wt (fresh flowers)) with low growth rate (Lin et al. 2015), explaining the  
65 current high process separation costs and low productivity.

66  
67 Therefore, microalgae with their rapid growth rate and ability to utilise plenty of low cost and  
68 abundant resources including solar energy, atmospheric CO<sub>2</sub> and wastewater, can provide  
69 significant improvements over existing industrial practice. Research into microalgal lutein  
70 production has already made considerable progress. It was found that nitrogen intake into the  
71 system is vital for the synthesis and accumulation of lutein (Xie et al. 2014; Ho et al. 2015). It  
72 was also discovered that low illumination is a decisive factor to ensure the suitable conditions for  
73 lutein production (Ho et al. 2012; Xie et al. 2013). Furthermore, the influence of different  
74 photobioreactor types and operating modes on microalgae growth and lutein synthesis were  
75 studied (Del Campo 2000; del Rio-Chanona, Zhang, et al. 2016). Different microalgae species  
76 such as *Scenedesmus obliquus*, *Chlorella sorokiniana*, *Chlorella zofingiensis* and *Desmodesmus*  
77 sp. were found to produce lutein with an intercellular lutein content that is between 6 and 15  
78 times higher than marigold (Shi et al. 2002; Sánchez et al. 2008; Del Campo 2000).

79  
80 Despite these achievements, one of the challenges that most severely prevents the  
81 industrialisation of the microalgae based lutein production process is to efficiently conduct a  
82 dynamic optimisation of the process. Precise control over a long-term bioprocess is  
83 indispensable to guarantee its safety and efficiency, as bioprocesses are very sensitive to the  
84 process operating conditions such as pH, temperature and nutrient supply. Hence, executing  
85 process optimisation within a control scheme can remarkably improve the process profitability  
86 (del Rio-Chanona, Zhang, et al. 2016). In order to resolve this challenge, it is essential to initially

87 construct a robust mathematical model, which is able to accurately simulate and predict the  
88 dynamic performance of the long-term photo-production system, for both microalgae growth and  
89 lutein production throughout the entire process (Zhang, Dechatiwongse, Del-Rio-Chanona, et al.  
90 2015).

91  
92 However, so far little attention has been paid on this aspect. At present, two methodologies have  
93 been predominantly developed for bioprocess simulation, namely kinetic modelling and artificial  
94 neural networks (ANNs). Specific to microalgal lutein production, a kinetic model including  
95 effects of both nutrient concentration and light intensity on biomass growth and lutein synthesis  
96 has not been proposed in the literature. In addition, although a kinetic model can display good  
97 accuracy and predictability, it is noticed that constructing such a complex model is always a  
98 difficult mathematical task (*e.g.* parameter estimation of highly nonlinear ordinary differential  
99 equation systems) even if the bioproduct synthesis mechanisms have been identified. Hence, its  
100 application in bioprocess control and real-time optimization has not yet been well conducted.

101  
102 On the contrary, ANNs have been widely used in traditional chemical engineering processes for  
103 process control (Mjalli 2005; Fissore et al. 2004), and their applications have been recently  
104 extended to biochemical systems (Witek-Krowiak et al. 2014; Rosales-Colunga et al. 2010).  
105 Furthermore, they have been successfully employed for the purpose of reproducing, controlling,  
106 and optimising the dynamical behaviour of microalgae based bioprocesses (García-Camacho et  
107 al. 2016; del Rio-Chanona, Manirafasha, et al. 2016). The key advantage of ANN over kinetic  
108 modelling is that the investigated system can be treated as a black-box, without the necessity to  
109 develop any empirical or analytical correlation. This significantly reduces the difficulty in model

110 structure design and model parameter estimation. The challenge in constructing ANNs, however,  
111 is the requirement of large and varied raw data sets to express good predictions in supervised  
112 learning, which is particularly time consuming for long-term bioprocesses (Witek-Krowiak et al.  
113 2014). As a result, there is not much development in the use of ANN for longer lasting biological  
114 systems (del Rio-Chanona, Manirafasha, et al. 2016).

115  
116 Therefore, to facilitate the industrialisation of lutein production and investigate the capability of  
117 ANNs in long-term bioprocess simulation, the work presented here aims to construct robust  
118 ANNs capable of accurately modelling and predicting the dynamic behaviour of microalgae  
119 biomass growth and lutein synthesis. Different strategies will be adopted to resolve the challenge  
120 arising from the limited amount of experimental data sets. In particular, *Desmodesmus* sp. is  
121 selected in the current study due to its highest intracellular lutein content (up to 5.0 mg/g) and  
122 great thermo-tolerant properties (highest growth rate at 35°C and can survive up to 46°C)  
123 compared to other algae species (Xie et al. 2014; Xie et al. 2013).

124  
125 The work is divided in the following sections. In this Section 1, a background introduction has  
126 been provided to both the object of investigation and the main tools that are used later in this  
127 study. In Section 2, the specific details of the implementation of ANNs in the current work are  
128 laid out, including a schematic of the structure of the ANNs and use of the “elbow rule” in  
129 deciding the key parameters. In Section 3, the results are presented and discussed, showing how  
130 with a small amount of experimental data points the designed ANNs can give robust predictions  
131 of the dynamic behaviours of the current investigated bioprocesses. This demonstrates their  
132 suitability for use in both real-time and offline optimisation frameworks. Finally, in the

133 conclusion section, the original findings are summarised and an overview of future avenues of  
134 research is provided.

135

## 136 **2. Theory and methods**

### 137 **2.1 Experiment setup**

138 Six fed-batch experimental processes have been carried out in our work. In these experiments,  
139 microalga *Desmodesmus* sp. F51 was used for lutein production, and the operating temperature  
140 was controlled at 35°C for biomass growth and lutein synthesis. This is because 35°C has been  
141 reported to be the optimal temperature for *Desmodesmus* sp. F51 growth (Xie et al. 2013). A 1 L  
142 tubular photobioreactor (15.5 cm in length and 9.5 cm in diameter) was used in these  
143 experiments with an external light source mounted on both sides. Initial biomass concentrations  
144 in these experiments were kept constant, and incident light intensities were set from 150  $\mu\text{mol m}^{-2}$   
145  $\text{s}^{-1}$  to 600  $\mu\text{mol m}^{-2} \text{s}^{-1}$ . Nitrate influent was supplied to the reactor from the 60<sup>th</sup> hour until the  
146 end of the process due to the consumption of initial culture nitrate, and its inflow rate was fixed  
147 at 3 mL  $\text{hr}^{-1}$ . Influent nitrate concentration was chosen as 0.1 M or 0.5 M in different  
148 experiments, and all the processes lasted for six days.

149

150 Biomass concentration, nitrate concentration, and lutein production in all the processes were  
151 measured once per 12 hours, thus in total for each experimental data set there are 12 data points.  
152 Amongst these data sets, four of them were used in the current study for ANN training (model  
153 construction), and the remaining two (Test 1 and Test 2, under different light intensity and  
154 inflow rate conditions) were used for the predictability verification of the ANNs. All of the  
155 experiments were in duplicate, and the detailed presentation of the experimental measurement

156 techniques can be found in the previous research (Xie et al. 2013). The detailed operating  
157 conditions of these experiments are listed in Table I.

158

## 159 **2.2 Selection of number of hidden layers**

160 In this study, two different ANN structures were explored and their performance tested. A first  
161 ANN with one hidden layer, and a second ANN with two hidden layers. This enabled us to  
162 determine which structure is the most suitable one for future process optimization and control,  
163 given that few theoretical bases would highlight one over the other. Both ANN structures were  
164 fully linked and implemented in Python 2.7 using *pybrain* as a library (Schaul et al. 2010). The  
165 function implemented in the hidden layer nodes was chosen to be a sigmoid function, with a  
166 linear function in the output layer. This choice comes from recommendations in the literature  
167 that suggest such a combination is robust and flexible (Nielsen 2014).

168

## 169 **2.3 Inputs and outputs of ANNs**

170 In all cases the ANNs received 5 inputs which are the key operating factors and state variables in  
171 the current experiments:

- 172 • Incident light intensity;
- 173 • Biomass concentration;
- 174 • Nitrate concentration;
- 175 • Lutein production;
- 176 • Nitrate inflow concentration.

177 Of these inputs, incident light intensity and nitrate inflow concentration are experimentally  
178 controllable and known at all times in the experimental settings. The other three inputs are



179 measured once per 12 hours throughout the entire experimental operating time course. All inputs  
180 fed to the ANN were rescaled, as the variables in the system under analysis are of very different  
181 magnitude and ANNs require similarly scaled inputs. This can benefit the network training in  
182 terms of speed and numerical performance, therefore producing a smooth reduction in error  
183 during training and promoting a more accurate gradient evaluation.

184

185 Furthermore, the outputs of the ANNs are the change of rate of the state variables including  
186 biomass concentration, nitrate concentration, and lutein production. These changes in the rate  
187 (either accumulation or consumption) are added into the current ANN inputs to predict the  
188 process states at the next time, which is schematically presented in Fig. 1. This strategy was  
189 adopted from our recent research in which using the change of states by giving their past  
190 information is found to give the network higher accuracy compared to directly predicting future  
191 states based on past ones (del Rio-Chanona, Manirafasha, et al. 2016).

192

#### 193 **2.4 Selection of number of neurons in hidden layers**

194 Through literature correlations between the number of inputs and outputs desired, an initial  
195 number of nodes in each hidden layer was estimated to be around 20 (Lawrence et al. 1996;  
196 Elisseff & Paugam-Moisy 1996). The exact number of neurons in the current ANNs, however,  
197 was tested through a hyper-parameter selection procedure, together with the number of training  
198 epochs (times that the ANN is trained). Four experimental data sets, each with 12 points, were  
199 employed in this hyper-parameter selection step.

200

201 This hyper-parameter selection was conducted in  $k$ -fold fashion, where a selection of 4 of the 5

202 data sets is used for training and the remaining data set is used to obtain an estimate of the  
203 maximum error along the trajectory and across all involved variables. Then another subset of 4 is  
204 selected and the resulting network is evaluated against the remaining data set. This procedure is  
205 repeated until the specific network parameter configuration has been tested against all the sets,  
206 and it is then possible to obtain an average of maximum error for each parameter configuration.  
207 These errors are compared, and the network configuration with the lowest error is chosen. The  
208 choices for the number of nodes in the hidden layers were 5, 10, 15, 20, and 30, and those for the  
209 number of training epochs were 15, 30, 50, 100, 150, 200, and 300. The following figures, *i.e.*  
210 Fig. 2 and Fig. 3, show the 3D landscapes representing the intersection of these two choices and  
211 the error of cross-validation of the constructed ANNs.

212

213 To better select the parameters for the ANNs, a framework known as “elbow rule” is used. The  
214 objective of this framework is to select the optimal number of layers and training epochs that  
215 would give the best trade-off between accuracy of the network, training time and potential of  
216 model over-fitting. While often increasing the number of epochs or layers brings further  
217 improvement, there are diminishing returns and the improvement is not worth the increasing  
218 training time. In addition, noises from an excessive number of internal parameters can create  
219 problems of over-fitting, thus deteriorating the ANN’s predictive capabilities.

220

221 From Fig. 2 and Fig. 3, it can be seen that for the number of epochs, the error decreases rapidly  
222 at the beginning and flattens long before the maximum of 600 training epochs. As a result, 400  
223 epochs were chosen in the current study for both ANNs. In terms of the number of nodes, in the  
224 single hidden layer case (Fig. 2) adding more than 20 nodes does not result in an improvement in

225 performance, while in the ANN with two hidden layers (Fig. 3) a significant increase in the  
226 number of nodes beyond 10 seems to even contribute to a decrease in performance. Therefore,  
227 20 nodes were picked for the first case, and 15 nodes per hidden layer were chosen in the second  
228 case.

229

## 230 **2.5 Training of artificial neural networks**

231 For ANN training, 4 sets of experimental data points, each containing 12 data points, were used.

232 To enhance the model accuracy, 50 replications of artificial data sets were produced based on the

233 original data sets with random noise added of 3 % of the variable size, and a further 50

234 replications with a 5 % noise. These proportions were selected based on the realistic assessment

235 of the accuracy of the current original experimental measurements. The strategy of embedding

236 adequate random noise into original data sets to generate significant amount of artificial data sets

237 has been found to improve ANNs modelling and predictive power when simulating other

238 biological systems, even when relatively little experimental data is available (del Rio-Chanona,

239 Manirafasha, et al. 2016).

240

## 241 **3. Results and discussion**

### 242 **3.1 Training results of the artificial neural networks**

243 Once trained and cross-validated, both ANNs are constructed. The current proposed ANN

244 construction strategy was implemented on a personal computer with low specifications in a

245 realistic time (3 hours for the ANN training). This is more time-efficient than constructing a

246 kinetic model for future process design, since it could be a time-consuming task, in particular at

247 the early research stage, to fully identify the biochemical kinetic mechanisms. It is useful to

248 mention that modern PCs (*e.g.* Core i7) should have more than acceptable performance (30  
249 minutes at most of ANN training). This is further enhanced by the fact that more advanced ANN  
250 libraries allow computation to take place in Graphics Processing Units (GPU) in place of  
251 traditional CPU usage. The peculiar architecture of GPUs allows for 10-100 times faster training  
252 for equivalent amounts of loss, when CPU computation is used for the same cases.

253

254 Fig. 4 and Fig. 5 show the comparison between the training data and the simulation results. In  
255 order to assess the effectiveness of the current ANN training procedure, in this case only initial  
256 operating conditions are provided, and the ANNs simulate the bioprocess behaviors throughout  
257 the 132 hours of operations. For illustrative purpose, the maximum absolute percentage error  
258 (MAPE) of the one hidden layer ANN is presented in Fig. 4.

259

260 With the above results, it can be clearly seen that the ANNs are capable of accurately modelling  
261 the process dynamics to which they have already been exposed, even if only the initial operating  
262 conditions are given. Furthermore, if these ANNs are used in process optimization and control,  
263 they should be able to predict the performance of the investigated system under the operating  
264 conditions which they have never encountered. It is for this reason that two additional  
265 experimental data sets obtained from different experiments, namely Test 1 and Test 2, were used  
266 to test the predictive power of the designed ANNs.

267

### 268 **3.2 Predictability of the ANNs on a real-time framework**

269 During an ongoing experiment, the ANN should be capable of predicting the dynamic behavior  
270 of the bioprocess several time steps ahead. For this, a real-time framework can be put into place,

271 where in future work the system can be either controlled or optimized. After every set number of  
272 hours (12 hours in the current work) new experimental measurements would become available  
273 and the exact state of the system at that time would be known. Therefore, the ANN would only  
274 need to be able to accurately predict the performance of the current lutein production process for  
275 12 hours in advance, as new system data would be available after this time interval.

276

277 Given that in the current experimental setting only 12 hours would be needed, both ANNs are  
278 used to predict the process behaviors of the two additional experimental tests (Test 1 and Test 2)  
279 after 12 hours once a measurement is given, and such prediction is repeated throughout the entire  
280 experiment operating time. Fig. 6 and Fig. 7 show the prediction results of the two ANNs for  
281 both experimental tests. The MAPE of both models are below 10%, expect for the one hidden  
282 layer ANN when predicting nitrate concentration at Experiment Test 2. The results shown in the  
283 two figures consist in strong proof that both models can be effectively used for the real-time  
284 control and optimization of the investigated bioprocess given their high predictive power.

285

### 286 **3.3 Predictability of the ANNs on offline framework**

287 Moreover, to thoroughly explore the feasibility of applying ANNs into an offline optimal control  
288 framework where the entire process behavior of an unknown experiment is predicted before its  
289 implementation (del Rio-Chanona, Dechatiwongse, et al. 2015), the current ANNs are used to  
290 simulate the processes Test 1 and Test 2. It follows the procedure that a single initial  
291 experimental point is supplied to the network, then the ANN computes the next state.  
292 Nonetheless, at the subsequent time step, instead of using the next experimental measurement as  
293 input to the ANN, the last computed simulated point is used. This means that, with the exception

294 of the first point at  $t=0$  (initial condition), the ANN is not supplied with any other experimental  
295 point. Thus, throughout the entire process, the ANN simulation errors are accumulated and the  
296 model to system mismatches are magnified. A competent working model should keep the growth  
297 of simulation errors contained within the time horizon investigated.

298

299 As an example, Fig. 8 compares the prediction results of the ANN with two hidden layers to the  
300 real experimental results. The MAPE of the two hidden layers ANN in both cases are mostly  
301 around 7% to 12%. From the figure, it can be appreciated that the prediction of the ANN  
302 matches the experimental results of both test sets well. This, as mentioned earlier, shows that the  
303 current constructed ANN can predict not only 12 hours in advance, but up to 132 hours with high  
304 accuracy. This result strongly indicates the great competence of the current ANN for long-term  
305 bioprocess modelling and offline optimization.

306

307 The small tendency for a slight consistent error seen in Test 2 might indicate a bias in the  
308 training sets, mostly in the concentration of nitrogen, where the increase of concentration  
309 influenced by the nitrogen injection is not well represented. This can be attributed to a relative  
310 lack of variety in the training sets, as the nitrogen inflow concentration has only two different  
311 quantities used as input. This seems to be further compounded by the fact that most differences  
312 commence around  $t=60^{th}$  hr, when nitrogen flow is switched on. Adding a relevant feature in a  
313 model usually decreases the modelling error. However, if the feature describes a variable or  
314 quantity in the system that does not vary much, this might actually cause more error than the  
315 addition of this feature would correct (Hagan et al. 2014). This could be solved by having more  
316 data sets with different nitrogen inflow concentrations as the main point of change. This analysis

317 can therefore support further experimental design.

318

### 319 **3.4 Comparison of the artificial neural networks**

320 In the current study, the two ANNs are created, one with one hidden layer and the other with two  
321 hidden layers, both of which are tested against Test 1 and Test 2. Some comparison results which  
322 highlight the difference in performance between the two ANNs are presented in Fig. 9. From the  
323 figure, it is observed that the performance of the ANN with two hidden layers is much better than  
324 the one with one hidden layer, when simulating the trajectory of the entire process (offline  
325 framework). Thus, it is concluded that although both ANNs can provide accurate  
326 implementations for a real-time framework, if the aim is to execute offline optimization, only the  
327 ANN with two hidden layers should be selected due to its higher accuracy and predictive  
328 capacities.

329

330 However, it is important to emphasize that the conclusion of a two hidden layers ANN being  
331 superior to a one hidden layer ANN cannot be considered as a general rule for bioprocess  
332 modelling and optimization. This is because there exists a trade-off between model accuracy and  
333 risk of model over-fitting. In other words, although it is possible to enhance the model fitting  
334 results by increasing the number of ANN layers (hence increasing the amount of model  
335 parameters), the addition of extra layers can result in an over-fitting to the constructed ANN and  
336 severely aggravate the ANN predictive capabilities. Moreover, attention should be also paid on  
337 the counter-balance between increasing ANN training times by adding more layers and  
338 decreasing returns in improvements of predictive power. Thus, when employing ANNs to  
339 simulate an unknown biosystem, it is necessary to adopt the current proposed hyper-parameter

340 selection framework to determine the optimal ANN structure.

341

### 342 **3.5 Strategy of training data rescaling**

343 Furthermore, the current study concluded that when rescaling experimental data points for ANN  
344 training, it is vital to choose a suitable rescaling method to guarantee the quality of the network.

345 For example, in the current study, it is found that using the standardization method to center  
346 training data on the mean seems to have worked best in scaling the points and obtaining a  
347 network with better predictive power. This can indicate that the ANN, when trained, becomes  
348 more flexible when data points that maximize the variance are present. A min-max normalization  
349 method would have not captured this characteristic as well (Hastie et al. 2009). In addition, this  
350 type of scaling displays a problem with the presence of larger outliers, as in such data sets it  
351 would force many of the scaled points to lie very close to each other in an attempt to include the  
352 outliers, which will inevitably decrease the accuracy of the ANN.

353

354 Another advantage of choosing a standardization method is connected to the fact that, in the  
355 current study, a sigmoid function was chosen for the hidden layers. Sigmoid functions have the  
356 tendency to saturate, meaning that they make little distinction between inputs that are on the  
357 extremes of the range and produce the same outputs. These functions have worse learning if the  
358 inputs are large for either sign, as the gradients inside the ANN are flattened close to values of 0,  
359 which can lead to serious problems in the ANN learning process (Nielsen 2014). A  
360 standardization method, however, solves this issue by concentrating inputs in a limited range  
361 around a mean.

362



363 **Conclusions**

364 In the current study, two ANNs were constructed to simulate a long-term dynamic biosystem for  
365 microalgae biomass growth and lutein production. To guarantee the high accuracy of the current  
366 models, different strategies were implemented during the model construction. These include  
367 using a standardization method to rescale inputs into the ANN, adding adequate random noise to  
368 efficiently generate sufficient amounts of artificial data sets which are infeasible to obtain in  
369 practice, and identifying the optimal ANN structure and parameter values through a hyper-  
370 parameter selection framework.

371  
372 By comparing the ANN simulation results with the training data sets, it is found that the current  
373 ANNs can accurately represent the dynamic behavior of the current investigated biosystem. By  
374 comparing the ANN prediction results against the two sets of test experimental results, it is  
375 concluded that both of the current designed ANNs can be effectively applied into real-time  
376 process optimization and control frameworks, which strongly indicates their high potential for  
377 future process design and optimization. Furthermore, the current research demonstrated that the  
378 ANN with two hidden layers is capable of predicting accurately the entire dynamic trajectory of  
379 an unknown process before its implementation, further suggesting its adequateness for its use  
380 even in the offline optimization framework and extended process prediction durations.

381  
382 Moreover, due to the necessity of process modelling and optimization for the design of industrial  
383 scale sustainable biochemicals production processes, laboratory scale experimental data provides  
384 an important test set and can be used as a starting point for pilot plant tests. Inaccuracies of  
385 ANNs can be rapidly and easily corrected by further phases of learning. Thus, it is notable that

386 the procedure and strategies presented in the current study for ANN construction can be directly  
387 transferred to other bioprocesses and make significant contributions to their further  
388 industrialization. In addition, effective optimization algorithms (such as stochastic and  
389 evolutionary algorithms) can be further developed and embedded into the real-time framework to  
390 facilitate the optimization of ANNs for future bioprocess design and optimal control.

391

### 392 **Acknowledgments**

393 Author E. A. del Rio-Chanona would like to acknowledge CONACyT scholarship No. 522530  
394 This work was also supported by the National Natural Science Foundation of China (No.  
395 31071488), the National High Technology Research and Development Program 863, China  
396 (No.2014AA021701).

397

### 398 **References**

- 399 Del Campo, J., 2000. Carotenoid content of chlorophycean microalgae: factors determining  
400 lutein accumulation in *Muriellopsis* sp. (Chlorophyta). *Journal of Biotechnology*, 76(1),  
401 pp.51–59.
- 402 Capelli, B., Bagchi, D. & Cysewski, G.R., 2013. Synthetic astaxanthin is significantly inferior to  
403 algal-based astaxanthin as an antioxidant and may not be suitable as a human nutraceutical  
404 supplement. *Nutrafoods*, 12(4), pp.145–152.
- 405 Doshi, A. et al., 2016. Economic and policy issues in the production of algae-based biofuels: A  
406 review. *Renewable and Sustainable Energy Reviews*, 64, pp.329–337.
- 407 Elisseeff, A. & Paugam-Moisy, H., 1996. Size of multilayer networks for exact learning: analytic  
408 approach. In *NIPS'96 Proceedings of the 9th International Conference on Neural*

409 *Information Processing Systems*. MIT Press Cambridge, pp. 162–168.

410 Fernández-Sevilla, J.M., Acién Fernández, F.G. & Molina Grima, E., 2010. Biotechnological  
411 production of lutein and its applications. *Applied Microbiology and Biotechnology*, 86(1),  
412 pp.27–40.

413 Fissore, D., Barresi, A.A. & Manca, D., 2004. Modelling of methanol synthesis in a network of  
414 forced unsteady-state ring reactors by artificial neural networks for control purposes.  
415 *Chemical Engineering Science*, 59(19), pp.4033–4041.

416 García-Camacho, F. et al., 2016. Artificial neural network modeling for predicting the growth of  
417 the microalga *Karlodinium veneficum*. *Algal Research*, 14, pp.58–64.

418 Hagan, M.T. et al., 2014. *Neural Network Design* 2nd ed., Martin Hagan.

419 Hastie, T., Tibshirani, R. & Friedman, J., 2009. *The Elements of Statistical Learning*, New York,  
420 NY: Springer New York. Available at: <http://link.springer.com/10.1007/978-0-387-84858-7>.

421 Ho, S.-H. et al., 2015. Effects of nitrogen source availability and bioreactor operating strategies  
422 on lutein production with *Scenedesmus obliquus* FSP-3. *Bioresource Technology*, 184,  
423 pp.131–138.

424 Ho, S.-H., Lu, W.-B. & Chang, J.-S., 2012. Photobioreactor strategies for improving the CO<sub>2</sub>  
425 fixation efficiency of indigenous *Scenedesmus obliquus* CNW-N: Statistical optimization of  
426 CO<sub>2</sub> feeding, illumination, and operation mode. *Bioresource Technology*, 105, pp.106–113.

427 Lawrence, S., Giles, C.L. & Tsoi, A.C., 1996. *What size neural network gives optimal*  
428 *generalization? Convergence properties of backpropagation*,

429 Lin, J.-H., Lee, D.-J. & Chang, J.-S., 2015. Lutein production from biomass: Marigold flowers  
430 versus microalgae. *Bioresource Technology*, 184, pp.421–428.

431 Marz, U., 2011. *The Global Market for Carotenoids*.

432 Mata, T.M., Martins, A.A. & Caetano, N.S., 2010. Microalgae for biodiesel production and other  
433 applications: A review. *Renewable and Sustainable Energy Reviews*, 14(1), pp.217–232.

434 Mjalli, F.S., 2005. Neural network model-based predictive control of liquid–liquid extraction  
435 contactors. *Chemical Engineering Science*, 60(1), pp.239–253.

436 Nielsen, M., 2014. Neural Networks and Deep Learning..

437 del Rio-Chanona, E.A., Manirafasha, E., et al., 2016. Dynamic modeling and optimization of  
438 cyanobacterial C-phycoyanin production process by artificial neural network. *Algal*  
439 *Research*, 13, pp.7–15.

440 del Rio-Chanona, E.A., Zhang, D., et al., 2015. Dynamic Simulation and Optimization for  
441 *Arthrospira platensis* Growth and C-Phycocyanin Production. *Industrial & Engineering*  
442 *Chemistry Research*, 54(43), pp.10606–10614.

443 del Rio-Chanona, E.A., Dechatiwongse, P., et al., 2015. Optimal Operation Strategy for  
444 Biohydrogen Production. *Industrial & Engineering Chemistry Research*, 54(24), pp.6334–  
445 6343.

446 del Rio-Chanona, E.A., Zhang, D. & Vassiliadis, V.S., 2016. Model-based real-time optimisation  
447 of a fed-batch cyanobacterial hydrogen production process using economic model  
448 predictive control strategy. *Chemical Engineering Science*, 142, pp.289–298.

449 Río, E. Del et al., 2005. Efficient one-step production of astaxanthin by the  
450 microalga *Haematococcus pluvialis* in continuous culture. *Biotechnology and*  
451 *Bioengineering*, 91(7), pp.808–815.

452 Rosales-Colunga, L.M., García, R.G. & De León Rodríguez, A., 2010. Estimation of hydrogen  
453 production in genetically modified *E. coli* fermentations using an artificial neural network.  
454 *International Journal of Hydrogen Energy*, 35(24), pp.13186–13192.

455 Sánchez, J.F. et al., 2008. Biomass and lutein productivity of *Scenedesmus almeriensis*:  
456 influence of irradiance, dilution rate and temperature. *Applied Microbiology and*  
457 *Biotechnology*, 79(5), pp.719–729.

458 Schaul, T. et al., 2010. PyBrain. *Journal of Machine Learning Research*, 11, pp.743–746.

459 Shi, X.-M., Jiang, Y. & Chen, F., 2002. High-Yield Production of Lutein by the Green Microalga  
460 *Chlorella protothecoides* in Heterotrophic Fed-Batch Culture. *Biotechnology Progress*,  
461 18(4), pp.723–727.

462 Witek-Krowiak, A. et al., 2014. Application of response surface methodology and artificial  
463 neural network methods in modelling and optimization of biosorption process. *Bioresource*  
464 *Technology*, 160, pp.150–160.

465 Xie, Y. et al., 2013. Phototrophic cultivation of a thermo-tolerant *Desmodesmus* sp. for lutein  
466 production: Effects of nitrate concentration, light intensity and fed-batch operation.  
467 *Bioresource Technology*, 144, pp.435–444.

468 Xie, Y.-P. et al., 2014. Simultaneous enhancement of CO<sub>2</sub> fixation and lutein production with  
469 thermo-tolerant *Desmodesmus* sp. F51 using a repeated fed-batch cultivation strategy.  
470 *Biochemical Engineering Journal*, 86(7), pp.33–40.

471 Yaakob, Z. et al., 2014. An overview: biomolecules from microalgae for animal feed and  
472 aquaculture. *Journal of Biological Research-Thessaloniki*, 21(1), p.6.

473 Yen, H.-W., Sun, C.-H. & Ma, T.-W., 2011. The Comparison of Lutein Production by  
474 *Scenedesmus* sp. in the Autotrophic and the Mixotrophic Cultivation. *Applied*  
475 *Biochemistry and Biotechnology*, 164(3), pp.353–361.

476 Zhang, D., Dechatiwongse, P., Del-Rio-Chanona, E.A., et al., 2015. Analysis of the  
477 cyanobacterial hydrogen photoproduction process via model identification and process

478 simulation. *Chemical Engineering Science*, 128, pp.130–146.

479 Zhang, D. et al., 2016. Dynamic modelling of *Haematococcus pluvialis* photoinduction for  
480 astaxanthin production in both attached and suspended photobioreactors. *Algal Research*, 13,  
481 pp.69–78.

482 Zhang, D., Dechatiwongse, P., del Rio-Chanona, E.A., et al., 2015. Dynamic modelling of high  
483 biomass density cultivation and biohydrogen production in different scales of flat plate  
484 photobioreactors. *Biotechnology and Bioengineering*, 112(12), pp.2429–2438.

485

486

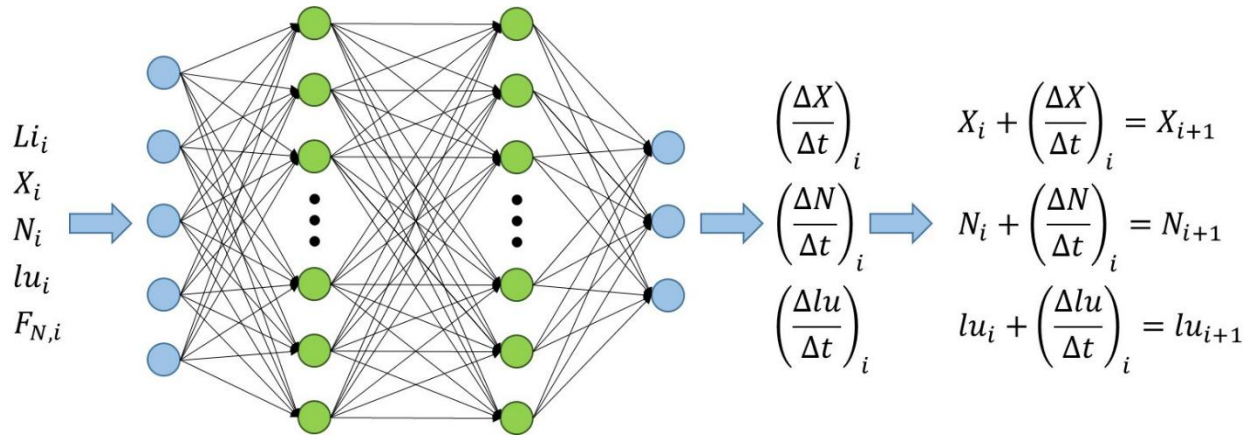
487 Table I: Operation conditions of the five experiments. Exp 1~Exp 4 are used for ANN training,  
488 and Test 1~Test 2 are used for ANN predictability verification.

Operation conditions	Exp 1	Exp2	Exp3	Exp4	Test 1	Test 2
Initial Biomass g L <sup>-1</sup>	0.07	0.07	0.07	0.07	0.07	0.07
Initial nitrate con. mM	8.8	30	8.8	8.8	30	8.8
Inflow rate mL h <sup>-1</sup>	3.0	3.0	3.0	3.0	3.0	3.0
Influent nitrate con. M	0.5	0.5	0.1	0.1	0.5	0.1
Light intensity $\mu\text{mol m}^{-2} \text{s}^{-1}$	300	600	150	600	480	300

489

490

491



492

493 Figure 1: Schematic of the current two hidden layer ANN structure. Time interval is chosen as

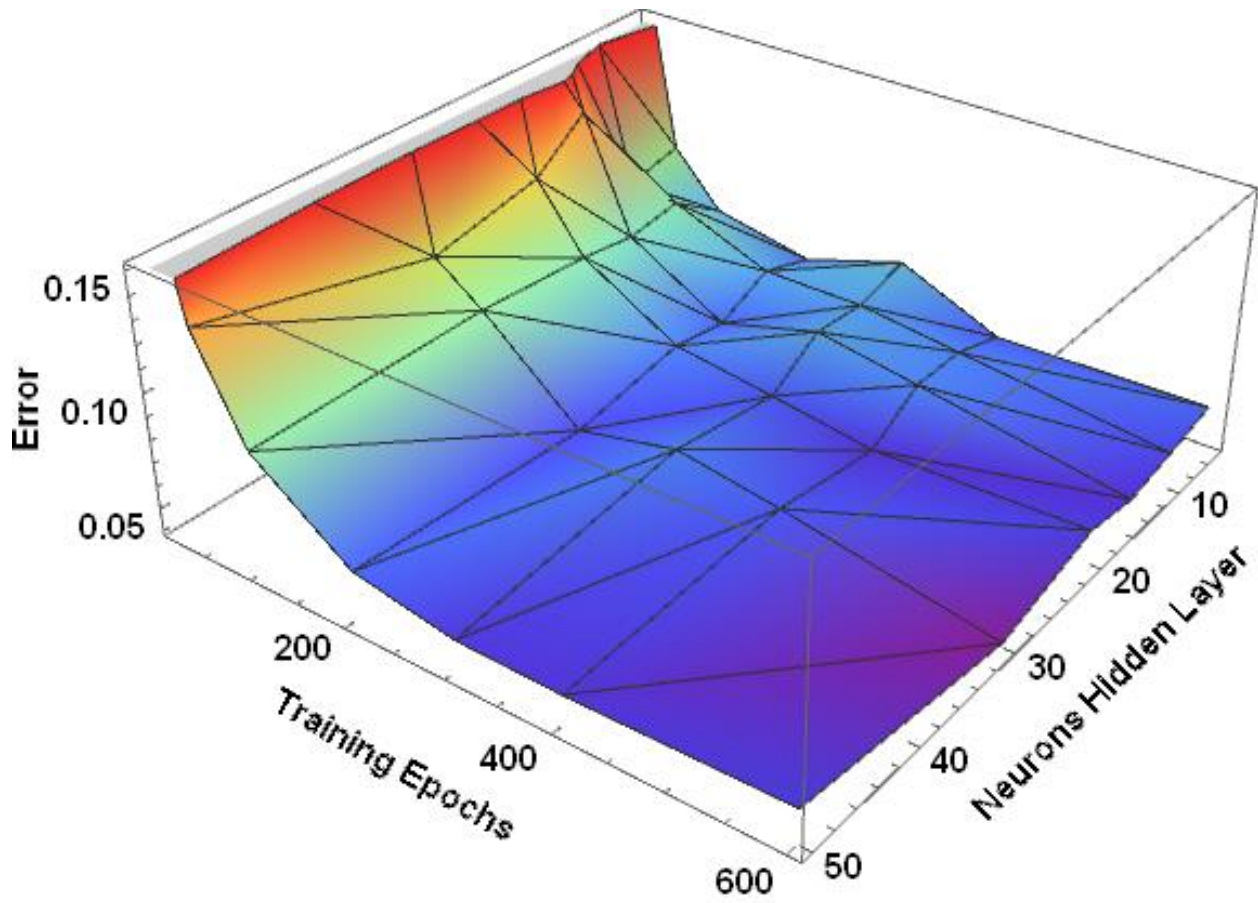
494 12 hours based on the real experimental implementation.

495

496

497





498

499

500

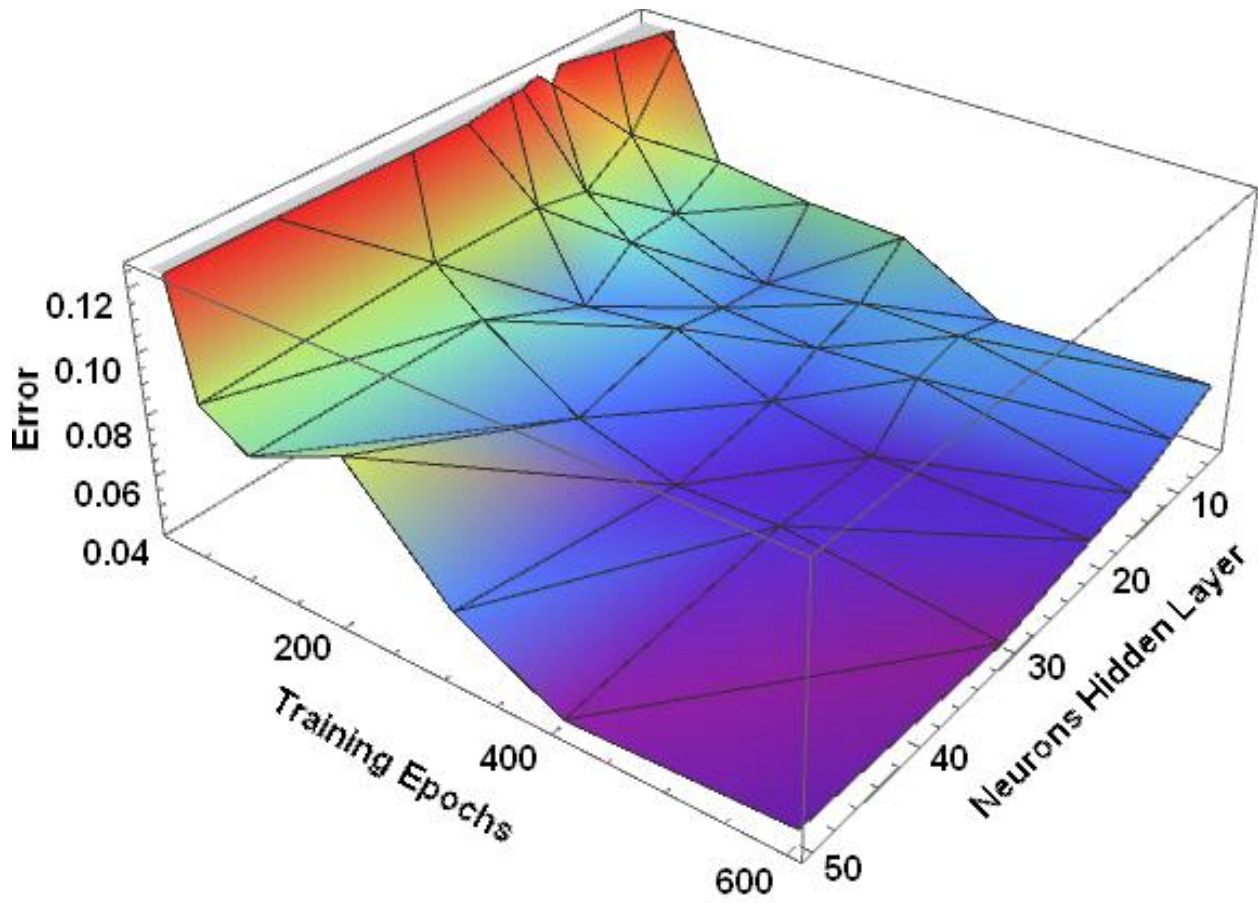
501

502

503

504

Figure 2: One hidden Layer elbow rule.

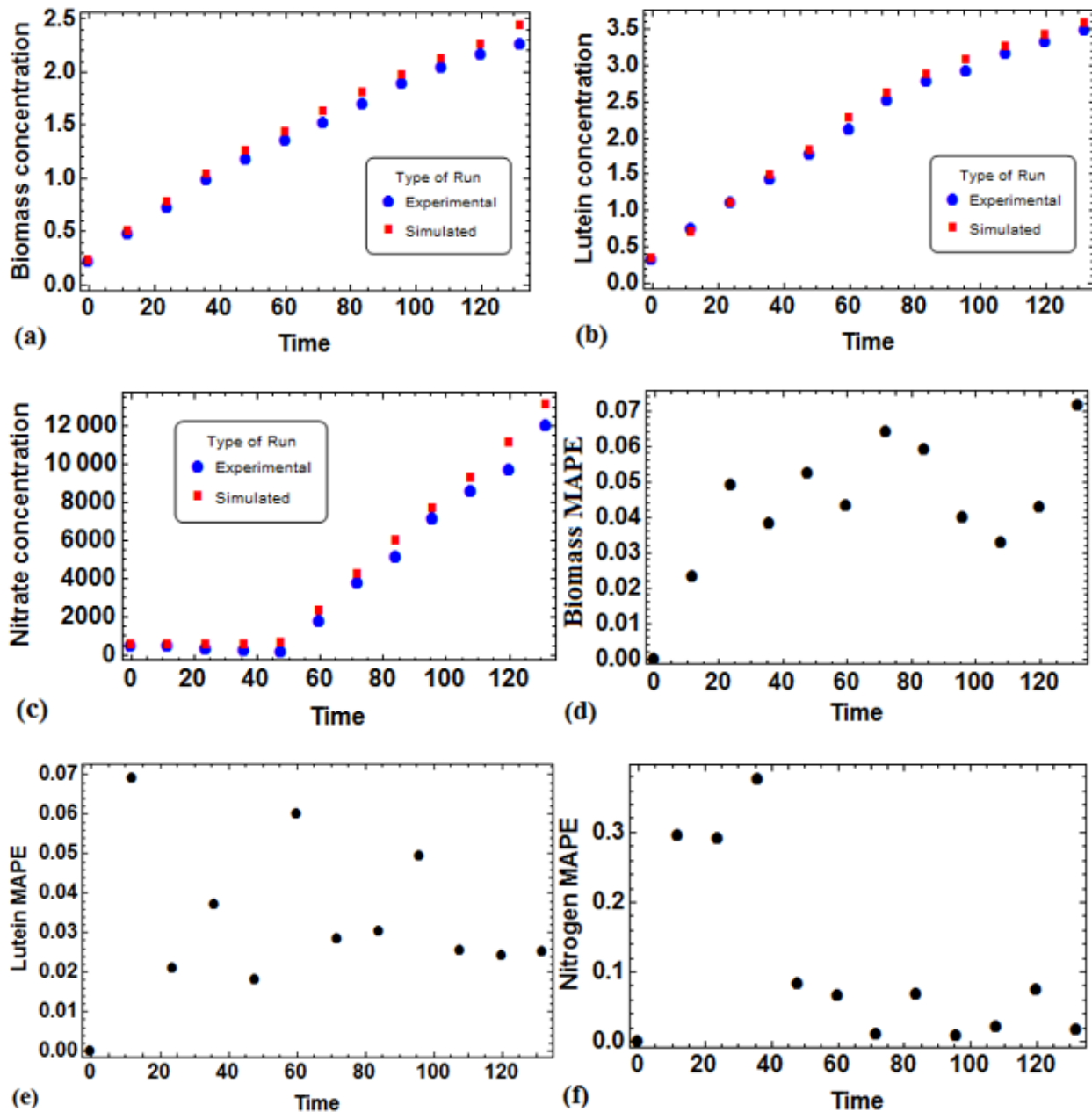


505  
506

Figure 3: Two hidden Layers elbow rule.

507

508



509

510 Figure 4: One hidden layer ANN process simulation results (training data set) when only an

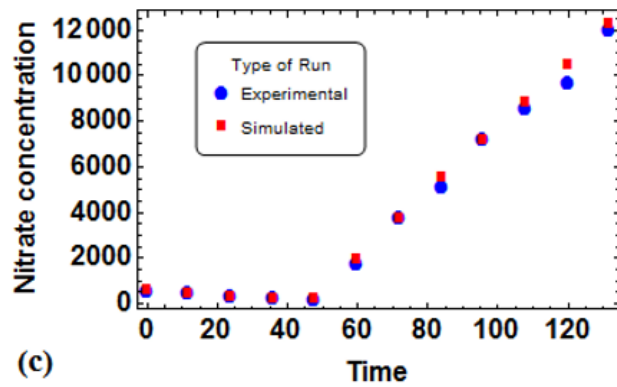
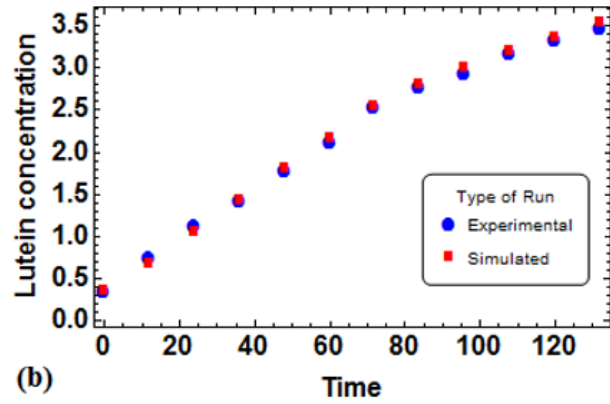
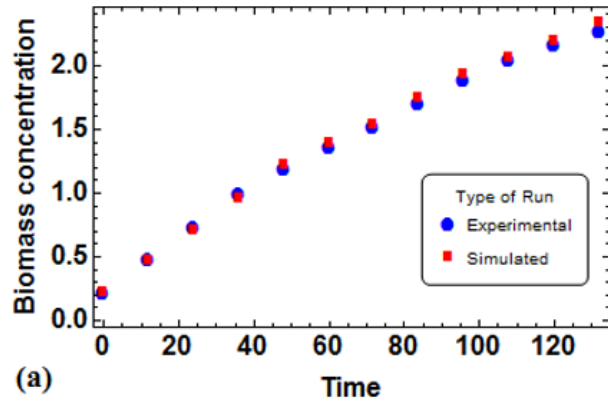
511 initial state is supplied and 132 hours of process operation time are simulated. Experimental data

512 points are averaged for convenience, and the measurement errors (error bars) are not presented in

513 the figures. The MAPE (mostly below 5%) indicates the high accuracy of the current ANN. The

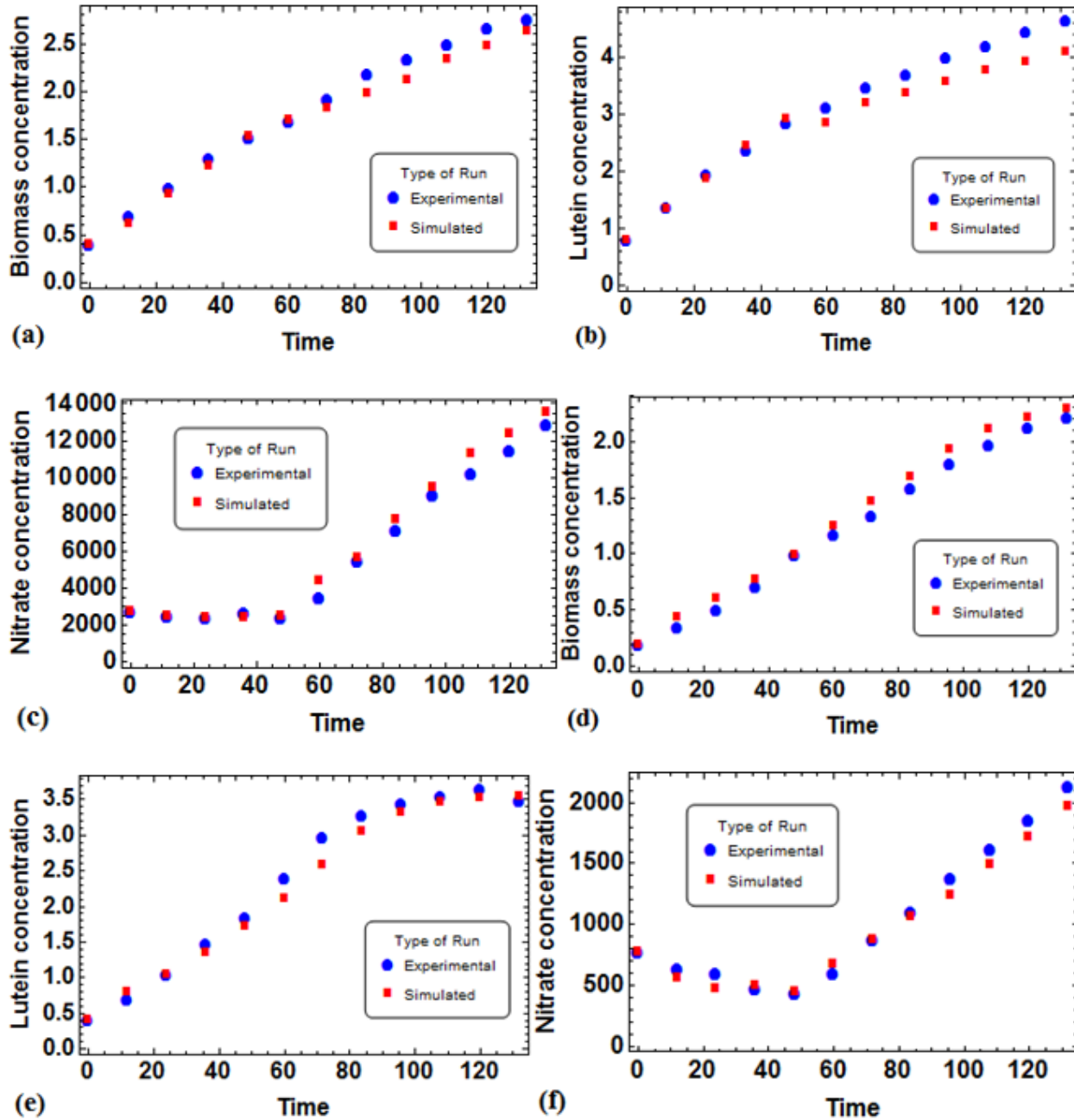
514 slightly higher MAPE in nitrogen at the beginning (12 hours to 36 hours) is explained at Section

515 3.3. Biomass concentration:  $\text{g L}^{-1}$ , lutein concentration and nitrate concentration:  $\text{mg L}^{-1}$ .



516

517 Figure 5: Two hidden layers ANN process simulation results (training data set) when only an  
 518 initial state is supplied and 132 hours of process operation time are simulated. The MAPE of this  
 519 ANN is mainly below 3%. Biomass concentration:  $\text{g L}^{-1}$ , lutein concentration and nitrate  
 520 concentration:  $\text{mg L}^{-1}$ .

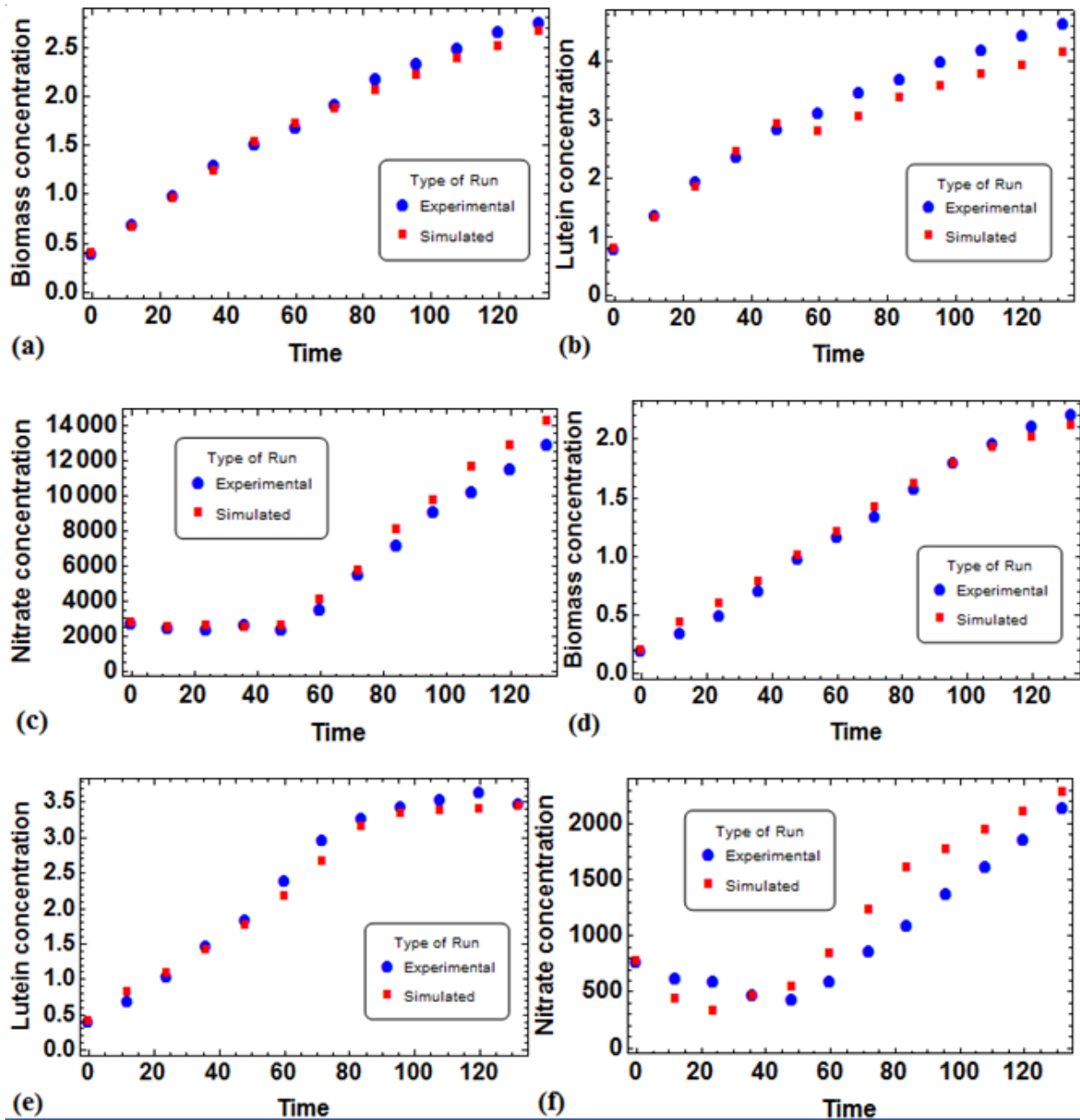


521 (a) (b) (c) (d) (e) (f)

522 Figure 6: Comparison of the two hidden layers ANN real-time prediction results with real

523 experimental data. (a), (b), and (c): Experiment Test 1. (d), (e), and (f): Experiment Test 2.

524 Biomass concentration: g L<sup>-1</sup>, lutein concentration and nitrate concentration: mg L<sup>-1</sup>.



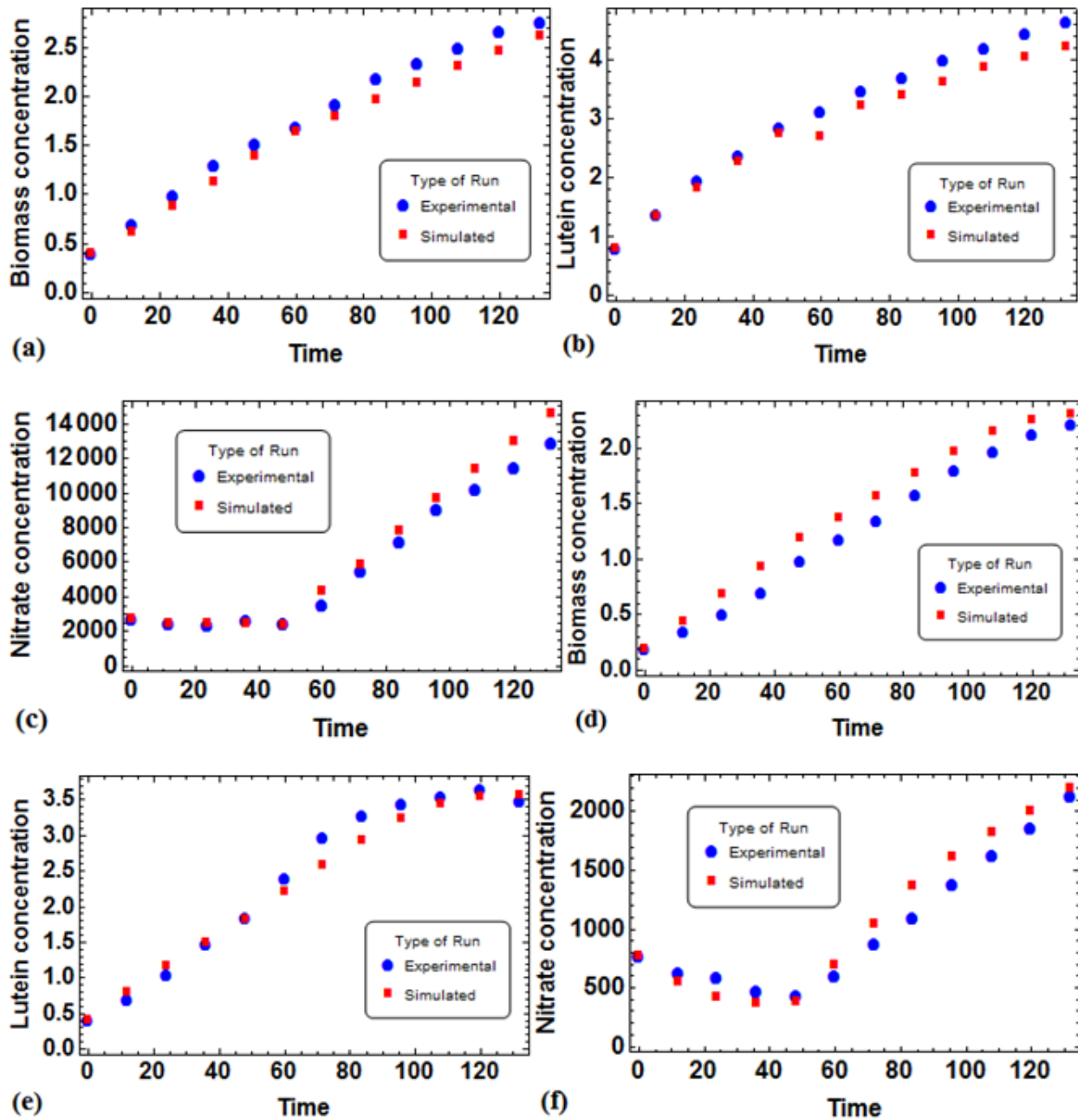
525

526 Figure 7: Comparison of the one hidden layer ANN real-time prediction results with real

527 experimental data. (a), (b), and (c): Experiment Test 1. (d), (e), and (f): Experiment Test 2.

528 Biomass concentration: g L<sup>-1</sup>, lutein concentration and nitrate concentration: mg L<sup>-1</sup>.



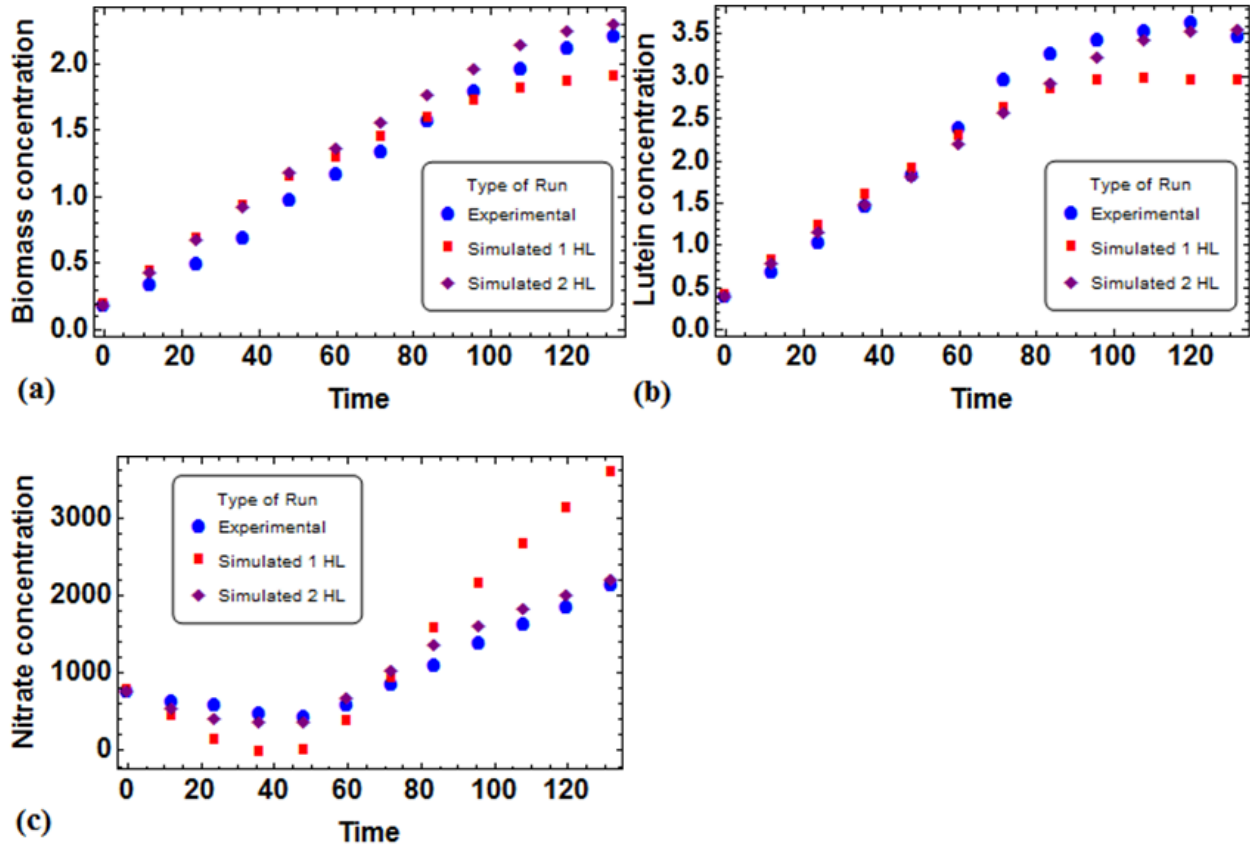


529

530 Figure 8: Comparison of the two hidden layers ANN offline prediction results with real

531 experimental data. (a), (b), and (c): Experiment Test 1. (d), (e), and (f): Experiment Test 2.

532 Biomass concentration: g L<sup>-1</sup>, lutein concentration and nitrate concentration: mg L<sup>-1</sup>.



533

534 Figure 9: Comparison of prediction results between one hidden layer ANN and two hidden layers

535 ANN in the offline framework (Experiment Test 2). Biomass concentration:  $g L^{-1}$ , lutein

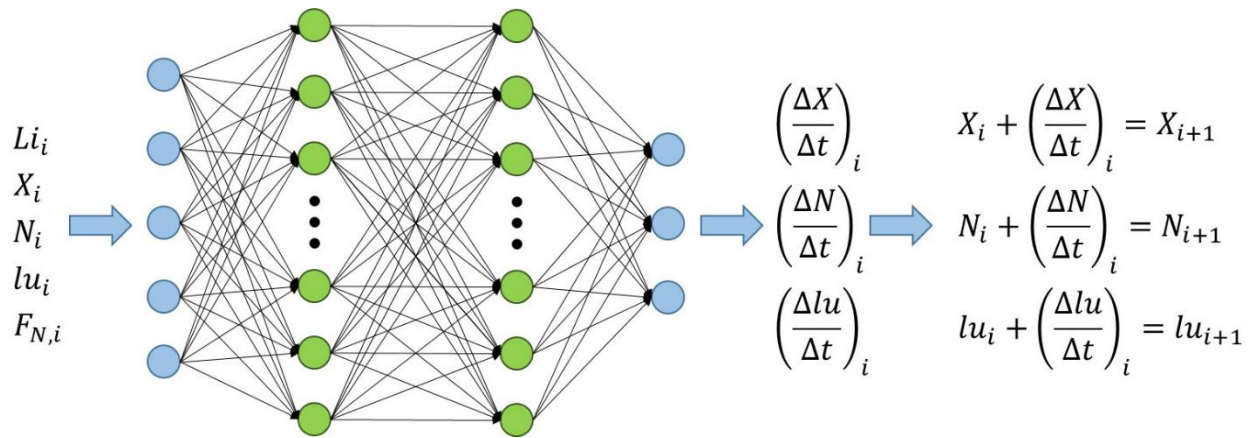
536 concentration and nitrate concentration:  $mg L^{-1}$ .

537

538

539





540

541 Graphical Table of Contents: Two robust artificial neural networks were constructed to simulate  
 542 the dynamic behaviour of microalgae growth and lutein production; different advanced strategies  
 543 were incorporated to guarantee the accuracy of the constructed models, including determining  
 544 the optimal network structure through a hyper-parameter selection framework, generating  
 545 artificial data sets by embedding appropriate random noise, and rescaling model inputs through  
 546 standardisation; the accuracy and predictive power of the models for long-term dynamic  
 547 bioprocess simulation in real-time and offline frameworks were demonstrated and verified  
 548 experimentally.

549

550

551

552

553

554

555