

1

2

3

4

A deep learning model for fish classification base on DNA barcode

5

6

7

8 Lina Jin^{1*}, Jiong Yu, Xiaoqian Yuan², Xusheng Du

9

10

11

12 ¹ School of Information Science and Engineering, Xinjiang University,

13 Urumqi, China.

14 ² School of Life Science and Technology, Xinjiang University, Urumqi,

15 China.

16

17 *Corresponding author

18 Email: jinlina@stu.xju.edu.cn (LJ)

19

20

21

22 **Abstract**

23 Fish is one of the most extensive distributed organisms in the world, fish
24 taxonomy is an important part of biodiversity and is also the basis of fishery resources
25 management. However, the morphological characters are so subtle to identify and
26 intact specimens are not available sometimes, making the research and application of
27 morphological method laborious and time-consuming. DNA barcoding based on a
28 fragment of the cytochrome c oxidase subunit I (COI) gene is a valuable molecular
29 tool for species identification and biodiversity studies. In this paper, a novel deep
30 learning classification approach that fuses Elastic Net-Stacked Autoencoder
31 (EN-SAE) with Kernel Density Estimation (KDE), named ESK-model, is proposed
32 bases on DNA barcode. In stage one, ESK-model preprocesses the original data from
33 COI fragments. In stage two, EN-SAE is used to learn the deep features and obtain the
34 outgroup score of each fish. In stage three, KDE is used to select the threshold base on
35 the outgroup scores and classify fish from different families. The effectiveness and
36 superiority of ESK-model have been validated by experiment on three dominant fish
37 families and comparisons with state-of-the-art methods. Those findings confirm that
38 the ESK-model can accurately classify fish from different family base on DNA
39 barcode.

40 **Introduction**

41 Fish is one of the most widely study group of aquatic organisms, about 27,683
42 fish species have most recently been catalogued into six classes, 62 orders and 540
43 families worldwide [1, 2]. Fish taxonomy and rapid species identification are the

44 fundamental premise of fishery biodiversity and fishery resources management, and
45 also an important part of marine biodiversity. As a traditional classification method,
46 morphological identification has successfully described nearly one million species on
47 the earth, which has laid a good foundation for species classification and identification
48 [3, 4]. However, routine species classification poses a challenge for fish classification
49 owing to four limitations. First, due to the differences of individual, gender and
50 geographical, phenotypic plasticity and genetic variability used for fish discrimination
51 can result in incorrect classification [5]. Second, with the deterioration of ecological
52 environment and disturbance of human activities, many fishery resources have been
53 seriously damaged, making it more difficult to collect fish specimens, especially for
54 those with less natural resources [6, 7]. Third, some fishes show subtle dissimilarity in
55 body shape, colors pattern, scale size and other external visible morphological
56 features, which cause confusion of the same species. Finally, the use of key not only
57 demands professional taxonomic knowledge, but also requires extensive experience
58 that misdiagnoses are common [8]. The limitations of morphology-based method, a
59 new technology to fish classification is needed.

60 Genomic approach is a new taxonomic technique combining molecular biology
61 with bioinformatics that uses DNA sequences as ‘barcodes’ to differentiate organisms
62 [5]. The DNA-based barcoding method is attainable to non-specialists. Many studies
63 have shown the effectiveness of DNA barcode technology for more than 15 years, it
64 has been extensive used in various fields such as species identification [9], discovery
65 of new species or cryptic species [10, 11], phylogeny and molecular evolution [12],

66 biodiversity survey and assessment [13, 14], customs inspection and quarantine [15],
67 conservation biology [16].

68 In the field of species classification, a short gene segment is used in DNA
69 barcoding, called the COI sequence, to build global standard dataset platforms,
70 universal technical rules and identification systems for animals' taxonomy [1]. COI
71 gene has the characteristics of high evolution rate, obvious interspecific variation,
72 relatively conservative within species, good universality of primers and easy
73 amplification [17]. Therefore, COI gene has been widespread employed as an
74 effective DNA barcode for species classification of varied animal lineages, including
75 bird [18, 19], Mosquito [20, 21], marine fish [22-24], freshwater fish [25-27]. DNA
76 barcode based on COI gene can be used to identify marine fish up to 98%, while
77 freshwater fish can be identified with 93% accuracy [28]. The approach base on DNA
78 barcode has been proven to be a valuable molecular tool for fish classification.

79 However, the complexity and high-dimensional characteristics in COI gene
80 sequences, analyzing these sequences reasonably and obtaining accessible information
81 that humans can classify fishes correctly are a major challenge. This issue requires a
82 multidisciplinary approach to deal with DNA sequences and to analyze the
83 information contained from data. Deep learning, a method of learning and extracting
84 useful representations from raw data, trains model, and then, uses the model to make
85 predictions, has made great progress in recent years [29]. Therefore, in this paper, we
86 propose a novel approach based on DNA barcode, use the deep learning model to
87 classify fish from different families and determine which fishes are regarded as

88 outgroup, called ESK-model. To verify the effectiveness of the model, three families
89 with many species and obvious interspecific variation were selected as the datasets.
90 First, the model preprocesses the original data that makes the COI gene sequences
91 into a matrix representation, then, converts them into numerical data. Second, the
92 model learns these data using EN-SAE model and obtains an outgroup score of each
93 fish. Finally, the KDE model is used to generate a threshold and to predict which fish
94 is outgroup base on threshold. The main contributions of our paper are as follows:

95 • We introduce a deep learning model to classify fish from different families and
96 determine which fish is outgroup based on DNA barcode, which is effective and
97 robust.

98 • To solve the model overfitting caused by COI gene sample of species in the
99 same family is limited, an Elastic Net is used for the model to increase the
100 generalization ability.

101 • We employ EN-SAE model to receive outgroup scores. The decision threshold
102 is automatically learned from organisms in same family by KDE model. An original
103 predictor is proposed based on the anomaly scores, while other classification works
104 often omit the importance of automatic learning threshold.

105 • We quantitatively evaluate the performance of our approach, and the results
106 demonstrate that our ESK-model outperforms state-of-the-art methods.

107 **Materials and Methods**

108 **Data description**

109 The COI sequences from three dominant families of fish in this study were

110 obtained from GenBank(www.ncbi.nlm.nih.gov), including Sciaenidae, Barbinae and
111 Mugilidae. Among them, Sciaenidae and Mugilidae belong to marine fish, Barbinae
112 belongs to freshwater fish. The genetic relationship and molecular divergence are
113 considered for selecting outgroups. The relevant information concerning the features,
114 specimen size and outgroup ratio of three families were summarized in Table 1.

115 **Table 1. Summary of datasets.**

Family	Feature	Instance	Outgroup ratio (%)
Sciaenidae	596	325	5.54
Barbinae	544	1022	2.35
Mugilidae	565	796	2.51

116 • Sciaenidae. The COI fragments contained 307 individuals of 21 species, 13
117 genera in Sciaenidae family. 18 homologous sequences in *Nemipterus virgatus*,
118 *Epinephelus awoara*, *Leiognathus equulus* and *Leiognathus ruconius* were selected
119 from different families, which were under the same order as Sciaenidae. After
120 processing, the length of COI gene fragment was 596 bp. Species of experimental
121 samples on Sciaenidae is shown in S1 Table.

122 • Barbinae. A total of 998 individuals from 103 species pertaining to 9 genera of
123 Barbinae were barcoded, which were 544 bp of COI gene sequence length. In
124 addition, 24 homologous sequences from 6 genera including *Foa brachygramma* and
125 *Cheilodipterus macrodon* belong to Apogonidae were used as outgroup. Species of
126 experimental samples on Barbinae is shown in S2 Table.

127 • Mugilidae. In this dataset, 776 Mugilidae sequences from 23 species belong to

128 7 genera were collected, which the length of COI gene was 565 bp. 20 homologous
 129 sequences in *Sphyraena pinguis* and *Sphyraena jello* from Mugiliformes were
 130 designated as outgroup. Species of experimental samples on Mugilidae is shown in S3
 131 Table.

132 Data preprocessing

133 Data definition

134 To facilitate the subsequent processing, DNA sequences can be represented by a
 135 matrix. The COI sequences for each family were formulated as follows:

$$136 \quad X = \begin{matrix} x_{11} & x_{12} & L & x_{1m} \\ x_{21} & x_{22} & L & x_{2m} \\ M & M & O & M \\ x_{n1} & x_{n2} & L & x_{nm} \end{matrix} \quad (2)$$

137 where n denotes the size of samples, and m denotes the number of features in
 138 each species.

139 One-hot code

140 One-hot encoding is the process of converting categorical variables into a form
 141 that is easy to use by machine learning algorithms, which are a combination of 0 and 1
 142 [30]. Therefore, the model encodes matrix into a numeric type of data using one-hot
 143 code. COI gene is composed of four bases, A, T, C, G. Each coded base was a 1×4
 144 vector [0, 0, a_i, 0], where a_i=1. Therefore, four bases were formulated as follows:

$$145 \quad \begin{matrix} A \\ T \\ C \\ G \end{matrix} = \begin{matrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \end{matrix} \quad (3)$$

146 Method introduction

147 **An overview of the ESK-model**

148 An overview of the proposed model is shown in Fig 1, ESK-model, which
149 consists of three stages: (1) the data preprocessing stage, (2) learning deep features
150 and computing each species outgroup score stage, and (3) deciding threshold base on
151 outgroup scores and classifying fishes from different family stage.

152 **Fig 1. An overview of ESK-model.** Three-dimensional visualization of data is
153 shown in Stage1, the distribution of the anomaly scores is shown in Stage 3.

154 In stage one, there are two main tasks: (1) preprocessing raw data by
155 representing the COI gene sequence in a matrix and (2) the one-hot code is performed
156 on the matrix because the features of each fish species need to be transformed into
157 numerical data. Finally, the preprocessed data are used as inputs for stage two.

158 In stage two, a deep learning network, EN-SAE, is used to learn deep features
159 from the data preprocessed in stage one. The model utilizes the EN-SAE model to
160 compress the digitalized data into a representation of the potential data to reconstruct
161 input, then, calculates the difference between input and output, and obtains an
162 outgroup score of each fish. Finally, the outgroup scores are used as inputs for stage
163 three.

164 In stage three, the KDE technique is used to learn the relationship between each
165 score from stage two, and then, fits the data distribution according to properties of the
166 outgroup scores. After that, the KDE model determines which fish is inner group and
167 which fish is outer group base on the threshold.

168 **Learning deep features and computing outgroup scores by EN-SAE**

169 Traditional AE is a three-layer neural network, including an input layer, an
170 output layer and a hidden layer. The structure of AE is symmetric, that is, the input
171 layer and output layer have the same number of nodes and the dimensions of each
172 node are the same too [31]. The purpose of AE is to compress input data and save
173 useful information to reconstruct input, and use the back propagation algorithm to
174 update the weights so that the output data is as similar to the input data as possible
175 [32]. However, the output data are not sufficient to yield a rewarding representation of
176 input. The reconstruction criterion with three-layer structure is unable to guarantee the
177 extraction of useful features as it can lead to the obvious solution “simply copy the
178 input” [33]. The SAE can greatly solve this problem.

179 The SAE model builds a deep neural networks base on AE by stacking several
180 AEs, puts the hidden representation of the upper layer as the input of the next AE. In
181 other word, extracting the compressed features of hidden layer into next AE to
182 training. In this way, training layer-by-layer can achieve input features compressed.
183 At the same time, more meaningful features of COI sequences are obtained. The
184 decoder can be reconstructed back into the input with a sufficiently small differences,
185 the structure of SAE is expressed in Fig 2.

186 **Fig 2. The structure of SAE.**

187 There are two basic steps in SAE training: encoder and decoder.

188 (1) Encoder: in this step, the activation function σ_e maps input data vector x to
189 hidden representation h that can compress the input data and retain more useful
190 representation, the typical form followed by a nonlinear representation:

191
$$h = \sigma_e(wx + b) \quad (4)$$

192 where x denotes input data vector, w is a weight matrix connecting the input
193 layer to hidden layer, b is bias vector belongs to nodes of latent layer, σ_e represents
194 activation function, such as Sigmoid, Relu, Tanh, etc.

195 (2) Decoder: in this step, the hidden representation h is mapped into
196 reconstruction vector y , the typical form as follows:

197
$$y = \sigma_d(w'h + b') \quad (5)$$

198 where w' is weight matrix connecting the latent layer to output layer, b' is bias
199 vector, σ_d represents activation function.

200 Loss function is defined to measure the reliability of SAE. SAE is trained to
201 reconstruct the features of input, and the weight of encoder and decoder are adjusted
202 to minimize the error between output and input. Thus, loss function is introduced, it is
203 represented by mean square error as follows:

204
$$L(w, b) = \frac{1}{n} \|y - x\|^2 \quad (6)$$

205 However, a COI gene fragment has too many features, which leads to the high
206 dimensionality of the training data. At the same time, fish species contained in each
207 family are limited, resulting in a relatively small dataset. Therefore, the model cannot
208 fully learn the characteristics of each fish species. In order to improve the
209 generalization ability of the proposed model, make the structure risky minimize, add
210 some kinds of constraint, reduce the weight of useless features. Base on this point,
211 Elastic Net composed of L1-norm and L2-norm is proposed in this method. The
212 structure of EN-SAE model is shown in Fig 3. It can also treat L1-norm and L2-norm

213 as penalty for loss function to restrict some parameters in the process of training.

214 **Fig 3. The structure of EN-SAE model.**

215 L1-norm also called Lasso regression, which contributes to generating a sparse
216 matrix. And it is defined as: $L_1(w) = \|w\| = \sum_i |w_i|$, where is the sum of the absolute
217 value of each element in weight vector w . Thus, it can be used to choose more
218 meaningful representations. When training model, the features are too many to select
219 what are contribute more for this model. So we dropped the connections that the
220 contribution of this model is so tiny, even if drop its have no impact on the model
221 [34]. It can reduce time consuming and study more useful features.

222 L2-norm also called Ridge regression, which is defined
223 as: $L_2(w) = \|w\|^2 = \sum_i |w_i|^2$, where is the sum of the squares of each element in weight
224 vector w . In the process of training, we usually tend to make the weight as small as
225 possible, because it is generally believed that the model with small parameters is
226 simpler and can fit different data effectively. Thus, L2-norm can void overfitting to
227 some extent and improve the generalization of model to adapt different fish families.

228 On the basis of proposed EN-SAE model, the outgroup score of each species
229 can be defined to measure whether fish is outgroup. The higher outgroup scores are,
230 the more likely they are to be treated as outgroup.

231 Therefore, the outgroup scores can be calculated by the following formula:

232
$$S(w, b) = \sum \|y - x\|^2 + \lambda_1 (\sum \|w\|^2) + \lambda_2 (\sum \|w\|) \quad (7)$$

233 where λ_1 is a parameter to adjust the L2-norm, λ_2 is a parameter to adjust the
234 L1-norm.

235 The EN-SAE model rejects high-dimensional features into low-dimensional
236 features step by step to obtain higher representation of COI sequences, which is
237 significantly more suitable for extract features and express data from original data.

238 **Analyzing the outgroup scores by using KDE**

239 KDE borrows its intuitive approach from the familiar histogram, which is among
240 the most common nonparametric density estimation techniques. KDE provides a
241 method of smoothing data points, and then, the distribution is fitted by the properties
242 of data itself. The decision threshold is ascertained by using KDE model base on the
243 outgroup scores. After that, the correct classification results of fish will be found.
244 Given the outgroup scores vector s , which obtained from EN-SAE model, KDE
245 estimates the probability density function (PDF) $p(s)$ in a nonparametric way:

$$246 \quad p(s) \approx \frac{1}{nh} \sum_{i=1}^n K\left(\frac{s - s_i}{h}\right) \quad (8)$$

247 where n is the size of the training dataset, $\{s_i\}$, $i = 1, 2, \dots, n$, is the training
248 dataset's outgroup scores vector, $K(\cdot)$ is the kernel function, and h is the bandwidth.

249 There are many kinds of kernel function, epanechnikov function is the most
250 common function in density estimation and also has a good effect. Therefore, the
251 epanechnikov is used to estimate the PDF:

$$252 \quad K_e(s) = \frac{3}{4}(1 - s^2) \quad (9)$$

253 After obtaining $p(s)$ of training the outgroup scores vector s by KDE, the
254 cumulative distribution function (CDF) $F(s)$ can be defined as follow:

$$255 \quad F(s) = \int_{-\infty}^s p(s) ds \quad (10)$$

256 Given a significance level parameter $\alpha \in [0,1]$ and combine with CDF, a decision
257 threshold s_α can be found, s_α satisfies following formula:

258
$$F(s_\alpha) = 1 - \alpha \quad (11)$$

259 If the outgroup scores of each species meet the condition $s \geq s_\alpha$, this species will
260 be considered as outgroup. On the contrary, they are ingroup. Confirmed by repeated
261 experiments that significance level parameter α is recommended to be set to 0.05.
262 ESK-model algorithm is summarized as shown in Algorithm 1.

Algorithm 1 ESK-model

Input: the COI sequences of each family

Output: the outgroup in matrix x

Step 1: Preprocessing data

Represent the DNA sequences by a matrix

Encode the matrix into a numeric type as matrix x

Step 2: Training EN-SAE model

Set the number of stacked AEs L .

Encoder process:

$$h_1 = \mathcal{S}_e(w_1 x + b_1)$$

for $i = 2$ to L do

$$h_i = \mathcal{S}_e(w_i x + b_i)$$

end for

Decoder process:

$$y_L = \mathcal{S}_d(w'_L x + b'_L)$$

for $j = L-1$ to 1 do

$$h_j = \mathcal{S}_d(w'_j x + b'_j)$$

end for

Step 3: Training KDE model

Calculate the outgroup score: $s = (y-x)^2 + I_1 w + I_2 w^2$

If $s < s_a$

the fish is ingroup

else

the fish is outgroup

end if

263 **Evaluation method**

264 To test performance of the proposed model, divide the sample into four situations
265 based on the actual classification and the ESK-model predicted classification. In
266 Table 2, four situations are illustrated with a confusion matrix. True positive (TP) is
267 the number of outgroups that are correctly classified as outgroup. True negative (TN)
268 is the number of ingroups that are correctly classified as ingroup. False positive (FP)
269 is the number of ingroups that are wrongly classified as outgroup. False negative (FN)
270 is the number of outgroups that are wrongly classified as ingroup.

271 **Table 2. Confusion matrix.**

	Actual	Positive	Negative
Forecast			
	Positive	TP	FP

Negative	FN	TN
----------	----	----

272 With the confusion matrix, the classification performance of all experiments was
273 measured by three criterions for Accuracy, Recall, and F-Measure. Those evaluation
274 equations are formulated as follows:

$$275 \quad Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (12)$$

$$276 \quad Recall = \frac{TP}{TP + FN} \quad (13)$$

$$277 \quad F - Measure = \frac{2TP}{2TP + FP + FN} \quad (14)$$

278 **Result**

279 **Impact of the number of stacked AEs on the classification** 280 **performance**

281 In the field of deep learning, the number of layers in the model is a critical factor,
282 because it directly affects the performance of the model. After all of the COI
283 sequences were prepared, the impact of the number of stacked AEs in our model on
284 the classification performance was also assessed. The outgroup scores trend with
285 various stacked AEs from 3 to 8 on Sciaenidae, Barbinae, Mugilidae is shown in Figs
286 4-6. The experimental results showed in Fig 4 demonstrate that, as the number of AEs
287 increased, the outgroup scores decreased rapidly on Sciaenidae when the number of
288 AEs was fewer than five. The outgroup scores gradually stabilized when the number
289 of AEs was greater than five. The outgroup scores on other two datasets showed the
290 same trend as those on Sciaenidae. These results reach the best classification
291 performance when the number of AEs was stacked to five.

292 **Fig 4. The outgroup scores trend on Sciaenidae.**

293 **Fig 5. The outgroup scores trend on Barbinae.**

294 **Fig 6. The outgroup scores trend on Mugilidae.**

295 Additionally, Table 3 illustrates the detailed data corresponding to Figs 4-6. The
296 results of Table 3 show that the outgroup scores of proposed model with five layers on
297 different datasets were 0.0193, 0.0197 and 0.01, respectively. Moreover, after the
298 number of AEs increased from 3 to 5, the outgroup scores on three datasets decreased
299 by approximately 29.04%, 41.02% and 16.90%, respectively. Those results indicate
300 that the proposed method can achieve low scores on identifying fish from different
301 families and the outgroup scores tend to be stable gradually.

302 **Table 3. The outgroup scores with different numbers of AEs on three datasets.**

Layers	Sciaenidae	Barbinae	Mugilidae
3	0.0272	0.0334	0.0213
4	0.0240	0.0218	0.0190
5	0.0193	0.0197	0.0177
6	0.0193	0.0196	0.0173
7	0.0185	0.0196	0.0173
8	0.0183	0.0196	0.0173

303 Note that the bold values denote the outgroup scores with five stacked AEs.

304 **Impact of Elastic Net on classification performance**

305 To evaluate effect of Elastic Net on the model performance, Stack Autoencoder-
306 Kernel Density Estimation (SK) and ESK-model were compared in Figs 7-9.

307 Evaluation method has been defined in previous section. As shown in Figs 7-9, all
308 evaluation indicators of ESK-model were higher than SK-model that without adding
309 Elastic Net.

310 **Fig 7. Accuracy on SK and ESK.**

311 **Fig 8. Recall on SK and ESK.**

312 **Fig 9. F-Measure on SK and ESK.**

313 In addition, Table 4 illustrates the detailed data corresponding to Figs 7-9. The
314 evaluation matrix (Accuracy, Recall, F-Measure) on Sciaenidae dataset increased by
315 approximately 0.0095, 0.0100 and 0.0052, respectively. Similarly, under the same
316 conditions, the evaluation matrix also increased in other two datasets. Those results
317 indicate that add Elastic Net can improve the performance of the ESK-model.

318 **Table 4. The evaluation matrix on SK and ESK models.**

	Sciaenidae	Barbinae	Mugilidae
SK	0.9528 0.9500 0.9744	0.9691 0.9900 0.9835	0.9212 0.9170 0.9567
ESK	0.9623 0.9600 0.9796	0.9938 0.9934 0.9967	0.9710 0.9694 0.9845

319 Note that the order of evaluation matrix is as follows Accuracy, Recall, F-measure.

320 **Performance evaluation with different methods**

321 We compared our method, ESK-model, with four state-of-art algorithms, one
322 class-support vector machine(OC-SVM) [35], K-nearest neighbor(KNN) [36],
323 isolation Forest(iForest) [37], autoencoder(AE) [38], to evaluate performance on the
324 task of sorting fishes from different families base on DNA barcode. Cross validation
325 was used for model training, and confusion matrix of different models on three

326 datasets is shown in Fig 10.

327 **Fig 10. Confusion matrix of five models on three datasets.**

328 In order to show the specific relationship between our method and other four
329 methods, we utilize histograms to compare the performance of three matrices.
330 Additionally, Table 5 exhibits the detailed data corresponding to Figs 11-13. As we
331 can see in Figs 11-13, ESK-model provides stable and efficient effects on three
332 datasets and generates the highest Accuracy, Recall and F-measure. Those results
333 show that ESK-model is superior to other methods.

334 **Fig 11. The evaluation matrix on Sciaenidae.**

335 **Fig 12. The evaluation matrix on Barbinae.**

336 **Fig 13. The evaluation matrix on Mugilidae.**

337 **Table 5. The evaluation matrix of three datasets.**

	OC-SVM	KNN	iForest	AE	ESK
Sciaenidae	0.7453 0.7300 0.8439	0.8491 0.8600 0.9149	0.9340 0.9500 0.9645	0.8302 0.8200 0.9011	0.9623 0.9600 0.9796
Barbinae	0.9599 0.9568 0.9779	0.9228 0.9402 0.9577	0.9722 0.9701 0.9848	0.9321 0.9269 0.9621	0.9938 0.9934 0.9967
Mugilidae	0.9544 0.9520 0.9754	0.9627 0.9607 0.9800	0.9378 0.9345 0.9661	0.9170 0.9127 0.9543	0.9710 0.9694 0.9845

338 Note that the best result is typeset in bold. The order of evaluation matrix is as
339 follows Accuracy, Recall, F-measure.

340 Discussion

341 This study set out with aim of constructing a novel deep learning model base on
342 DNA barcode with the employ of representative data to classify fishes from different
343 families and distinguish the outgroup. In this section, we discuss and analyze the

344 experimental results and findings.

345 A significant experimental result was that ESK-model achieved the best
346 discrimination performance when the number of stacked AEs was set to five. There
347 are several possible reasons for this result. The features of COI fragment can't be fully
348 learned when the number of stacked AEs is few. With the increase of the number of
349 AEs, the proposed model can learn the deeper hidden features of DNA sequences.
350 Obviously, when the number of AEs increased to five, the outgroup scores decreased
351 sharply. Experiments showed that increased the number of AEs did not improve
352 performance. The performance tended to be stable when the number of AEs was more
353 than five because the deep features had already fully learned. Hence, the prime
354 number of stacked AEs in the ESK-model was five.

355 Another considerable experimental result was that Elastic Net can improve the
356 performance of proposed model. A good model of deep learning usually requires
357 abundant data to training, while the limitation of obtaining the COI sequences of
358 fishes from different families, the problem of overfitting in small datasets is more and
359 more serious. To solve the overfitting problem in training process on small datasets is
360 of great importance. This model puts forward by using Elastic Net to solve overfitting
361 problem and improve the generalization ability of the model. Moreover, genetic
362 characteristics of fish belong to high-dimensional data, which is time-consuming
363 during training. However, directly combining a set of fully connected EN-SAE is
364 often useless to extract useful information. Elastic Net provides sparse connection
365 also can save training time. Therefore, Elastic Net can improve the performance of

366 proposed model.

367 The most surprising finding was that the proposed model could accurately
368 classify fish from different families. EN-SAE is used to calculate the outgroup scores,
369 when the outgroup scores are high, the probability of being identified as other families
370 is increased. The size of fish belonging to the same family is far more than that from
371 other families, EN-SAE can well fit and learn the characteristics of intraspecific fish
372 in the process of training. On the contrary, the number of fishes in different families is
373 relatively small, we can't get a good fitting effect, resulting in higher outgroup scores.
374 Therefore, they are more likely to be treated as outgroup in KDE-model. At the same
375 time, compared with other algorithms, it further confirms that the proposed model has
376 better performance in fish classification.

377 These positive results and findings suggest that the ESK-model based on deep
378 learning, with the utilization of DNA barcode technology, can effectively classify the
379 fish from different families.

380 **Conclusion**

381 In this study, we proposed the ESK-model that fuses EN-SAE model and KDE
382 technology for fish classification in different families through DNA barcode. The
383 experimental results and findings demonstrate the effectiveness of proposed model.

384 The main results and findings of this paper are as follows:

385 (1) The outgroup scores have leveled off when the number of stacked AEs was
386 set to five.

387 (2) Adding Elastic Net can prevent overfitting more effectively and improve the

388 generalization ability of the model.

389 (3) Compared with the current popular methods, our proposed model had better
390 performance in fish classification from different families by using COI sequences.

391 **References**

392 1. Xu L, Wang X, Van Damme K, Huang D, Li Y, Wang L, et al. Assessment
393 of fish diversity in the South China Sea using DNA taxonomy. *Fisheries Research*.
394 2021;233:105771. doi: 10.1016/j.fishres.2020.105771.

395 2. Fautin D, Dalton P, Incze LS, Leong J-AC, Pautzke C, Rosenberg A, et al.
396 An Overview of Marine Biodiversity in United States Waters. *PloS one*.
397 2010;5(8):e11914. doi: 10.1371/journal.pone.0011914.

398 3. Knowlton N, Weigt LA. New dates and new rates for divergence across the
399 Isthmus of Panama. *Proceedings of the Royal Society B: Biological Sciences*.
400 1998;265(1412):2257-63. doi: 10.1098/rspb.1998.0568.

401 4. Thu PT, Huang WC, Chou TK, Van NQ, Liao TY. DNA barcoding of coastal
402 ray-finned fishes in Vietnam. *PloS one*. 2019;14(9):e0222631. doi:
403 10.1371/journal.pone.0222631.

404 5. Hebert PD, Cywinska A, Ball SL, deWaard JR. Biological identifications
405 through DNA barcodes. *Proceedings Biological sciences*. 2003;270(1512):313-21.
406 doi: 10.1098/rspb.2002.2218. PubMed PMID: 12614582; PubMed Central PMCID:
407 PMC1691236.

408 6. Ramirez JL, Rosas-Puchuri U, Canedo RM, Alfaro-Shigueto J, Ayon P,
409 Zelada-Mazmela E, et al. DNA barcoding in the Southeast Pacific marine realm: Low

410 coverage and geographic representation despite high diversity. PloS one.
411 2020;15(12):e0244323. doi: 10.1371/journal.pone.0244323. PubMed PMID:
412 33370342; PubMed Central PMCID: PMC7769448.

413 7. Liang H, Meng Y, Luo X, Li Z, Zou G. Species identification of DNA
414 barcoding based on COI gene sequences in Bagridae catfishes. Journal of Fishery
415 Sciences of China. 2018;25(4):772. doi: 10.3724/sp.j.1118.2018.18036.

416 8. Xu L, Van Damme K, Li H, Ji Y, Wang X, Du F. A molecular approach to
417 the identification of marine fish of the Dongsha Islands (South China Sea). Fisheries
418 Research. 2019;213:105-12. doi: 10.1016/j.fishres.2019.01.011.

419 9. Ren BQ, Xiang XG, Chen ZD. Species identification of *Alnus* (Betulaceae)
420 using nrDNA and cpDNA genetic markers. Mol Ecol Resour. 2010;10(4):594-605.
421 doi: 10.1111/j.1755-0998.2009.02815.x. PubMed PMID: 21565064.

422 10. Newmaster SG, Fazekas AJ, Steeves RAD, Janovec J. Testing candidate
423 plant barcode regions in the Myristicaceae. Molecular Ecology Resources.
424 2008;8(3):480-90. doi: 10.1111/j.1471-8286.2007.02002.x.

425 11. Liu J, Moller M, Gao LM, Zhang DQ, Li DZ. DNA barcoding for the
426 discrimination of Eurasian yews (*Taxus* L., Taxaceae) and the discovery of cryptic
427 species. Mol Ecol Resour. 2011;11(1):89-100. doi:
428 10.1111/j.1755-0998.2010.02907.x. PubMed PMID: 21429104.

429 12. Necchi O, West JA, Rai SK, Ganesan EK, Rossignolo NL, de Goër SL.
430 Phylogeny and morphology of the freshwater red alga *Nemalionopsis*
431 *shawii* (Rhodophyta, Thoreales) from Nepal. Phycological Research. 2016;64(1):11-8.

432 doi: 10.1111/pre.12116.

433 13. Valentini A, Pompanon F, Taberlet P. DNA barcoding for ecologists. Trends
434 in ecology & evolution. 2009;24(2):110-7. doi: 10.1016/j.tree.2008.09.011. PubMed
435 PMID: 19100655.

436 14. Ji Y, Ashton L, Pedley SM, Edwards DP, Tang Y, Nakamura A, et al.
437 Reliable, verifiable and efficient monitoring of biodiversity via metabarcoding.
438 Ecology letters. 2013;16(10):1245-57. doi: 10.1111/ele.12162. PubMed PMID:
439 23910579.

440 15. Gathier G, van der Niet T, Peelen T, van Vugt RR, Eurlings MC, Gravendeel
441 B. Forensic identification of CITES protected slimming cactus (*Hoodia*) using DNA
442 barcoding. Journal of forensic sciences. 2013;58(6):1467-71. doi:
443 10.1111/1556-4029.12184. PubMed PMID: 23865560.

444 16. Liu J, Yan H-F, Newmaster SG, Pei N, Ragupathy S, Ge X-J, et al. The use
445 of DNA barcoding as a tool for the conservation biogeography of subtropical forests
446 in China. Diversity and Distributions. 2015;21(2):188-99. doi: 10.1111/ddi.12276.

447 17. Wang T, Qi D, Sun S, Liu Z, Du Y, Guo S, et al. DNA barcodes and their
448 characteristic diagnostic sites analysis of Schizothoracinae fishes in Qinghai province.
449 Mitochondrial DNA Part A. 2019; 30(4):592-601. doi:
450 10.1080/24701394.2019.1580273. Epub 2019 Apr 5. PMID: 30952197.

451 18. Hebert PDN, Stoeckle MY, Zemplak TS, Francis CM. Identification of Birds
452 through DNA Barcodes. PLOS Biology. 2004;2(10):e312. doi:
453 10.1371/journal.pbio.0020312.

-
- 454 19. Kerr KCR, Stoeckle MY, Dove CJ, Weigt LA, Francis CM, Hebert PDN.
455 Comprehensive DNA barcode coverage of North American birds. *Molecular Ecology*
456 *Notes*. 2007; 7(4):535-543. doi: 10.1111/j.1471-8286.2007.01670.x. PMID:
457 18784793; PMCID: PMC2259444.
- 458 20. Gang, Wang, Chunxiao, Li, Xiaoxia, Guo, et al. Identifying the Main
459 Mosquito Species in China Based on DNA Barcoding. *PloS one*. 2012; 7(10): e47051.
460 <https://doi.org/10.1371/journal.pone.0047051>.
- 461 21. Zhang J. Species identification of marine fishes in china with DNA
462 barcoding. *Evidence-based complementary and alternative medicine : eCAM*. 2011;
463 8(1):1-10. doi: 10.1155/2011/978253. PubMed PMID: 21687792; PubMed Central
464 PMCID: PMC3108176.
- 465 22. Steinke D, Zemplak TS, Boutillier JA, Hebert PDN. DNA barcoding of Pacific
466 Canada's fishes. *Marine Biology*. 2009;156(12):2641-7. doi:
467 10.1007/s00227-009-1284-0.
- 468 23. Thu PT, Huang WC, Chou TK, Van Quan N, Van Chien P, Li F, et al. DNA
469 barcoding of coastal ray-finned fishes in Vietnam. *PloS one*. 2019;14(9):e0222631.
470 doi: 10.1371/journal.pone.0222631. PubMed PMID: 31536551; PubMed Central
471 PMCID: PMC6752846.
- 472 24. Talaga S, Leroy C, Guidez A, Dusfour I, Girod R, Dejean A, et al. DNA
473 reference libraries of French Guianese mosquitoes for barcoding and metabarcoding.
474 *PloS one*. 2017;12(6):e0176993. doi: 10.1371/journal.pone.0176993. PubMed PMID:
475 28575090; PubMed Central PMCID: PMC5456030.

- 476 25. Decru E, Moelants T, De Gelas K, Vreven E, Verheyen E, Snoeks J.
477 Taxonomic challenges in freshwater fishes: a mismatch between morphology and
478 DNA barcoding in fish of the north-eastern part of the Congo basin. *Mol Ecol Resour.*
479 2016;16(1):342-52. doi: 10.1111/1755-0998.12445. PubMed PMID: 26186077.
- 480 26. Iyiola OA, Nneji LM, Mustapha MK, Nzeh CG, Oladipo SO, Nneji IC, et al.
481 DNA barcoding of economically important freshwater fish species from north-central
482 Nigeria uncovers cryptic diversity. *Ecology and evolution.* 2018;8(14):6932-51. doi:
483 10.1002/ece3.4210. PubMed PMID: 30073057; PubMed Central PMCID:
484 PMC6065348.
- 485 27. Wang T, Qi D, Sun S, Liu Z, Du Y, Guo S, et al. DNA barcodes and their
486 characteristic diagnostic sites analysis of Schizothoracinae fishes in Qinghai province.
487 Mitochondrial DNA Part A, DNA mapping, sequencing, and analysis.
488 2019;30(4):592-601. doi: 10.1080/24701394.2019.1580273. PubMed PMID:
489 30952197.
- 490 28. Ward RD, Hanner R, Hebert PD. The campaign to DNA barcode all fishes,
491 FISH-BOL. *Journal of fish biology.* 2009;74(2):329-56. doi:
492 10.1111/j.1095-8649.2008.02080.x. PubMed PMID: 20735564.
- 493 29. Jin S, Zeng X, Xia F, Huang W, Liu X. Application of deep learning methods
494 in biological networks. *Briefings in bioinformatics.* 2020;bbaa043. doi:
495 10.1093/bib/bbaa043. PubMed PMID: 32363401.
- 496 30. Chu Z, Yu J. An end-to-end model for rice yield prediction using deep
497 learning fusion. *Computers and Electronics in Agriculture.* 2020;174:105471. doi:

498 10.1016/j.compag.2020.105471.

499 31. Chen J, Sathe S, Aggarwal C, Turaga D. Outlier Detection with Autoencoder
500 Ensembles. Proceedings of the 2017 SIAM International Conference on Data Mining
501 (SDM). 2017;pp. 90-98. doi: 10.1137/1.9781611974973.11 .

502 32. Homoliak I. Convergence Optimization of Backpropagation Artificial Neural
503 Network Used for Dichotomous Classification of Intrusion Detection Dataset. Journal
504 of Computers. 2017;143-55. doi: 10.17706/jcp.12.2.143-155.

505 33. Vincent P, Larochelle H, Lajoie I, Bengio Y, Manzagol PA. Stacked
506 Denoising Autoencoders: Learning Useful Representations in a Deep Network with a
507 Local Denoising Criterion. Journal of Machine Learning Research.
508 2010;11(12):3371-408.

509 34. Taaffe K, Pearce B, Ritchie G. Using kernel density estimation to model
510 surgical procedure duration. International Transactions in Operational Research.
511 2018;28(1):401-18. doi: 10.1111/itor.12561.

512 35. Erfani SM, Rajasegarar S, Karunasekera S, Leckie C. High-dimensional and
513 large-scale anomaly detection using a linear one-class SVM with deep learning.
514 Pattern Recognition. 2016;58:121-34. doi: 10.1016/j.patcog.2016.03.028.

515 36. Hastie T, Tibshirani R. Discriminant adaptive nearest neighbor classification.
516 IEEE Transactions on Pattern Analysis & Machine Intelligence. 1996;18(6):607-16.

517 37. Liu FT, Ting KM, Zhou ZH. Isolation-Based Anomaly Detection. Acm
518 Transactions on Knowledge Discovery from Data. 2012;6(1):1-39.

519 38. Guo J, Liu G, Zuo Y, Wu J, editors. An Anomaly Detection Framework

520 Based on Autoencoder and Nearest Neighbor. 2018 15th International Conference on
521 Service Systems and Service Management (ICSSSM); 2018;pp. 1-6, doi:
522 10.1109/ICSSSM.2018.8464983.

523 **Supporting information**

524 **S1 Table. Species of experimental samples on Sciaenidae.**

525 **S2 Table. Species of experimental samples on Barbinae.**

526 **S3 Table. Species of experimental samples on Mugilidae.**

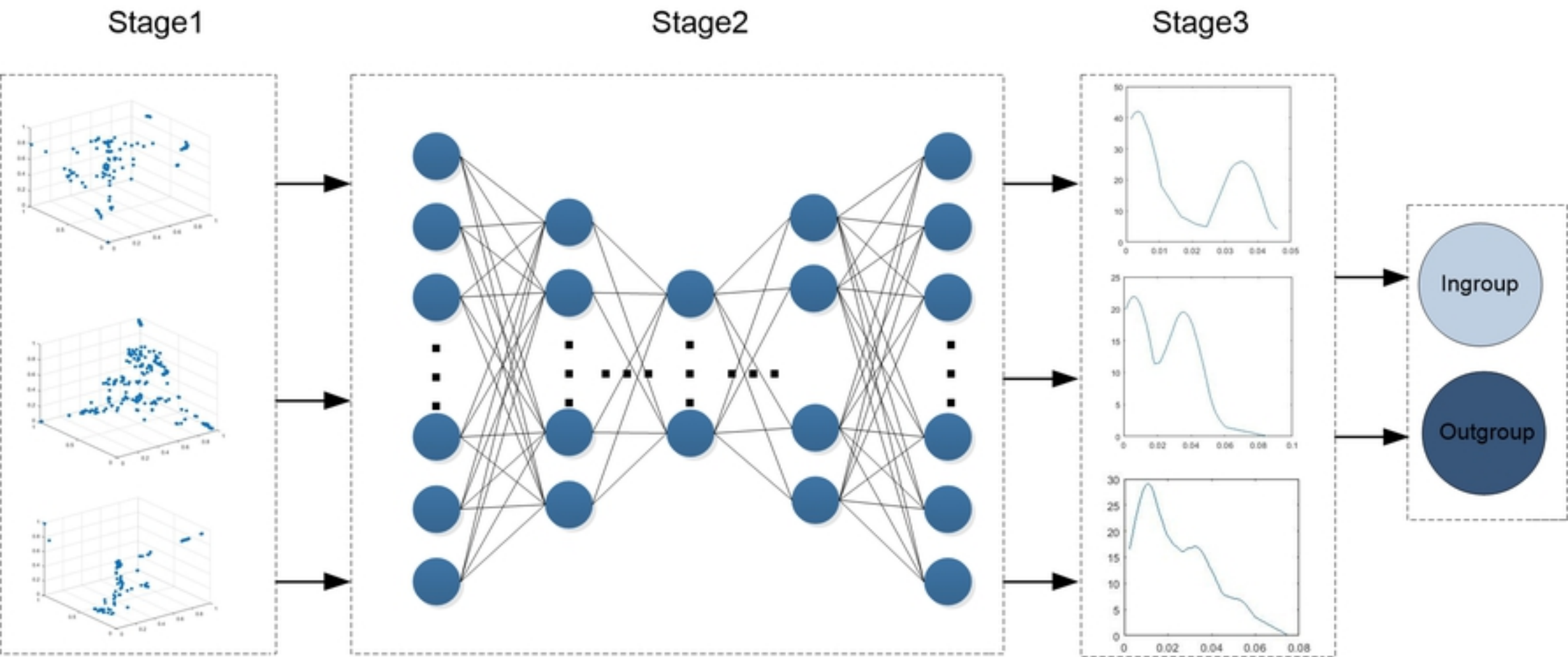
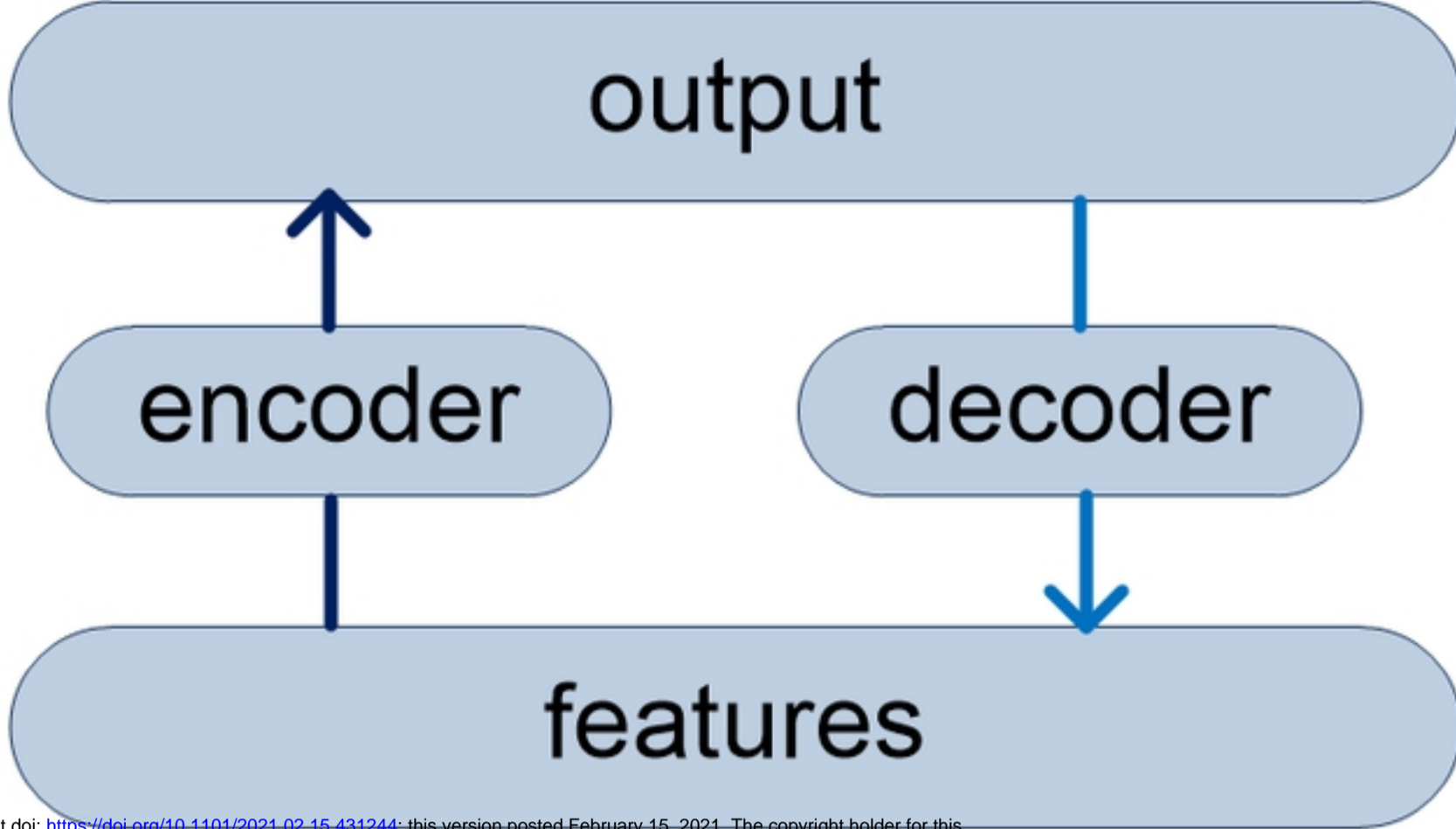


Fig1



bioRxiv preprint doi: <https://doi.org/10.1101/2021.02.15.431244>; this version posted February 15, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

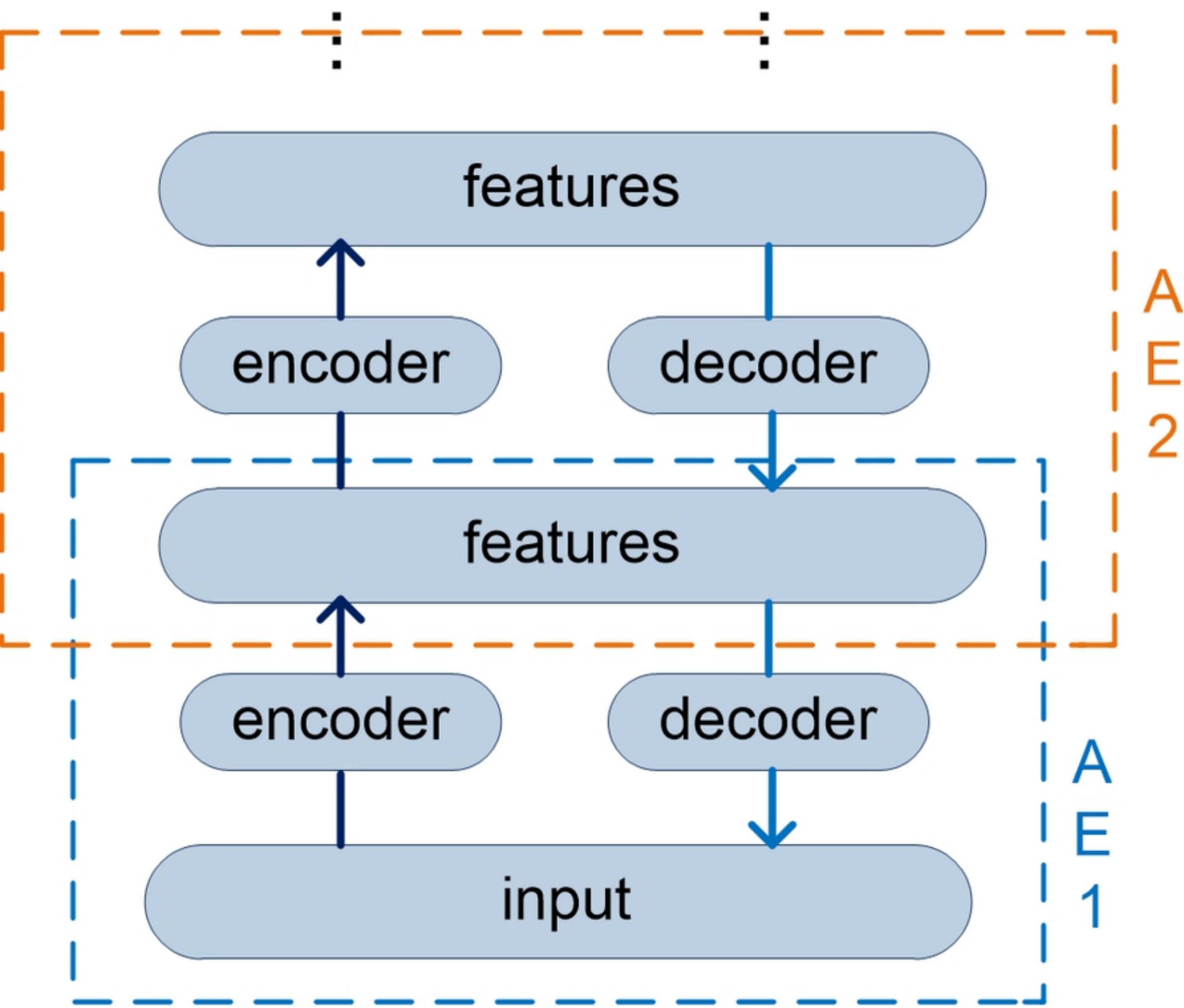


Fig2

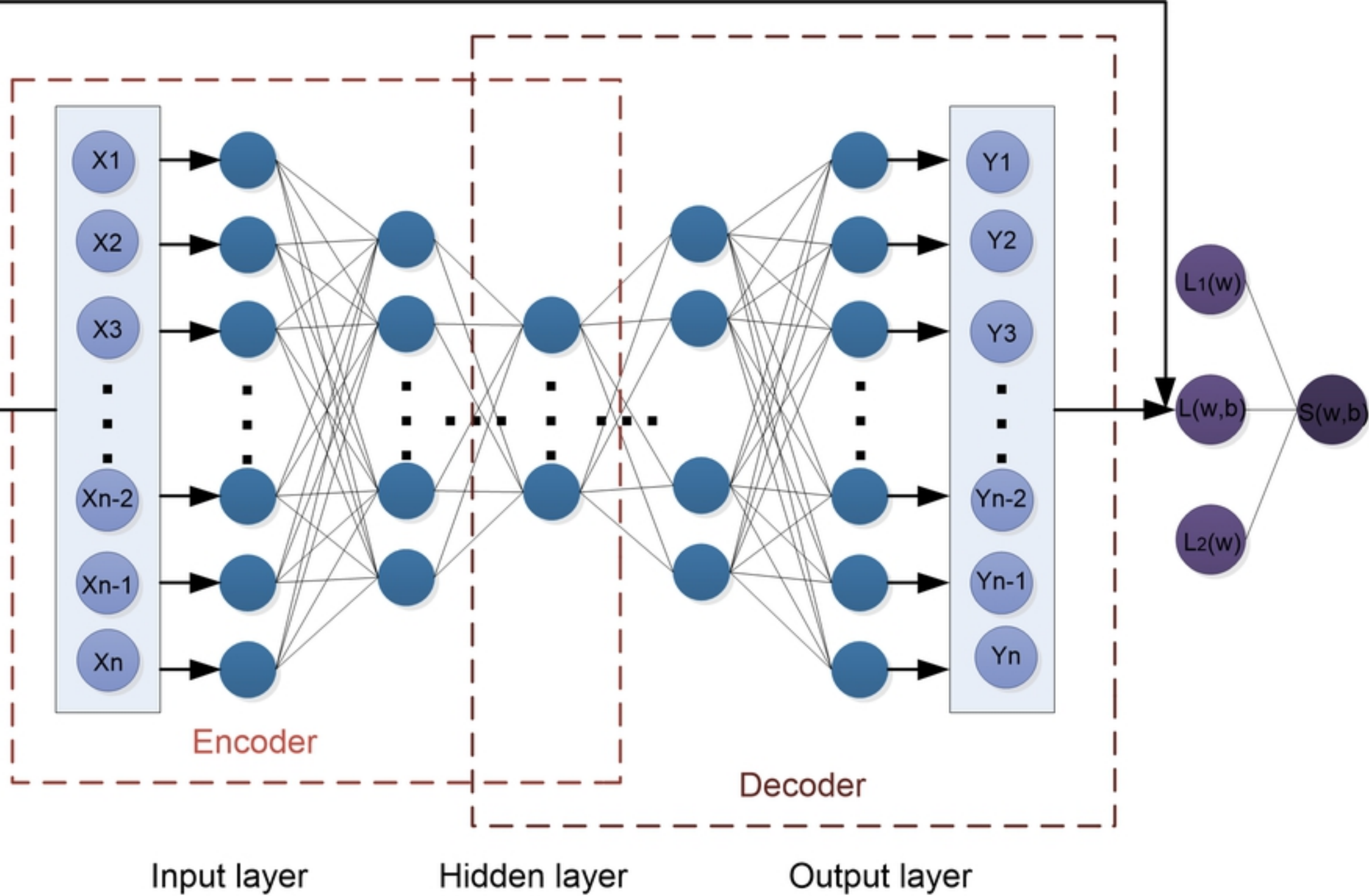


Fig3

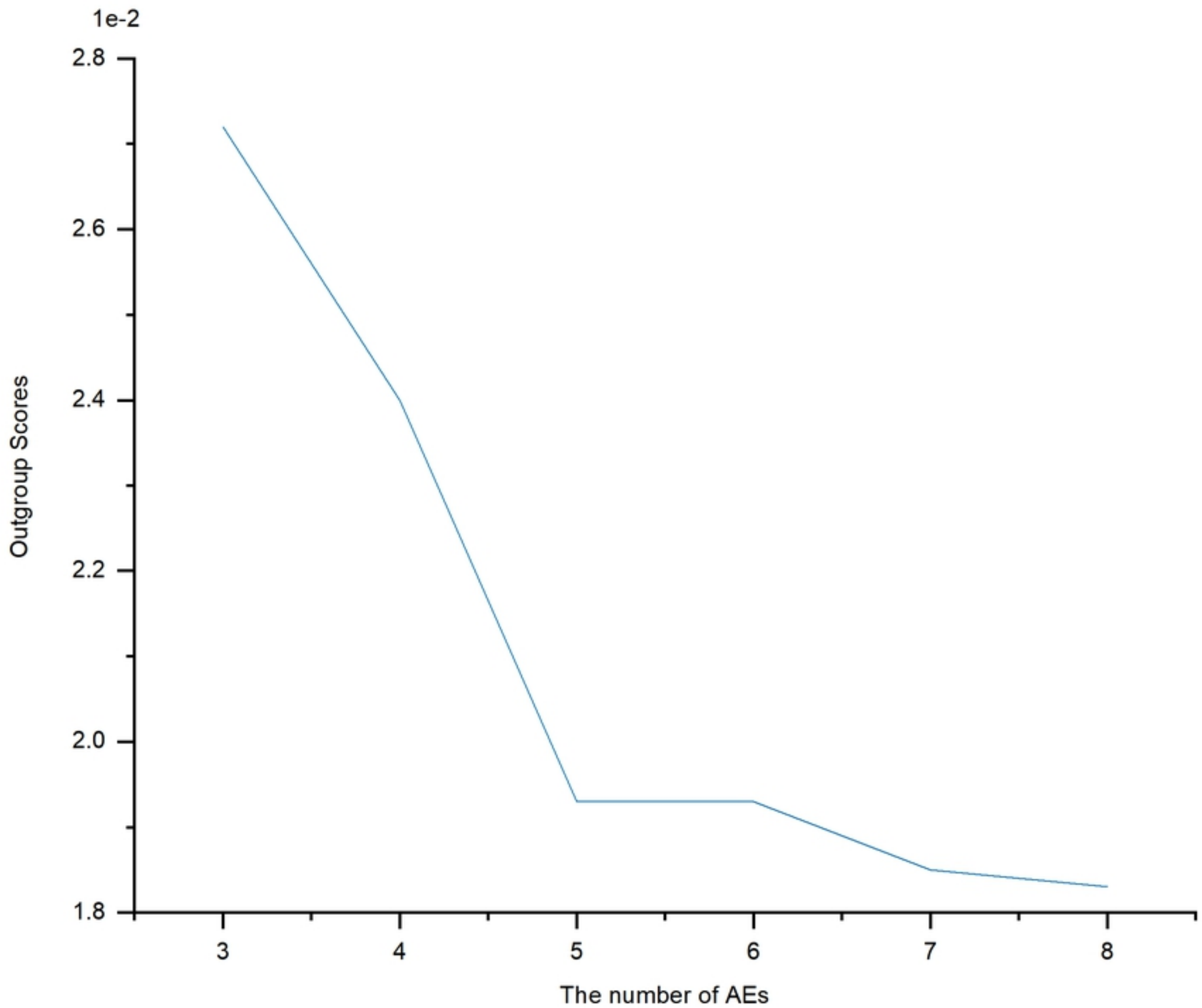


Fig4

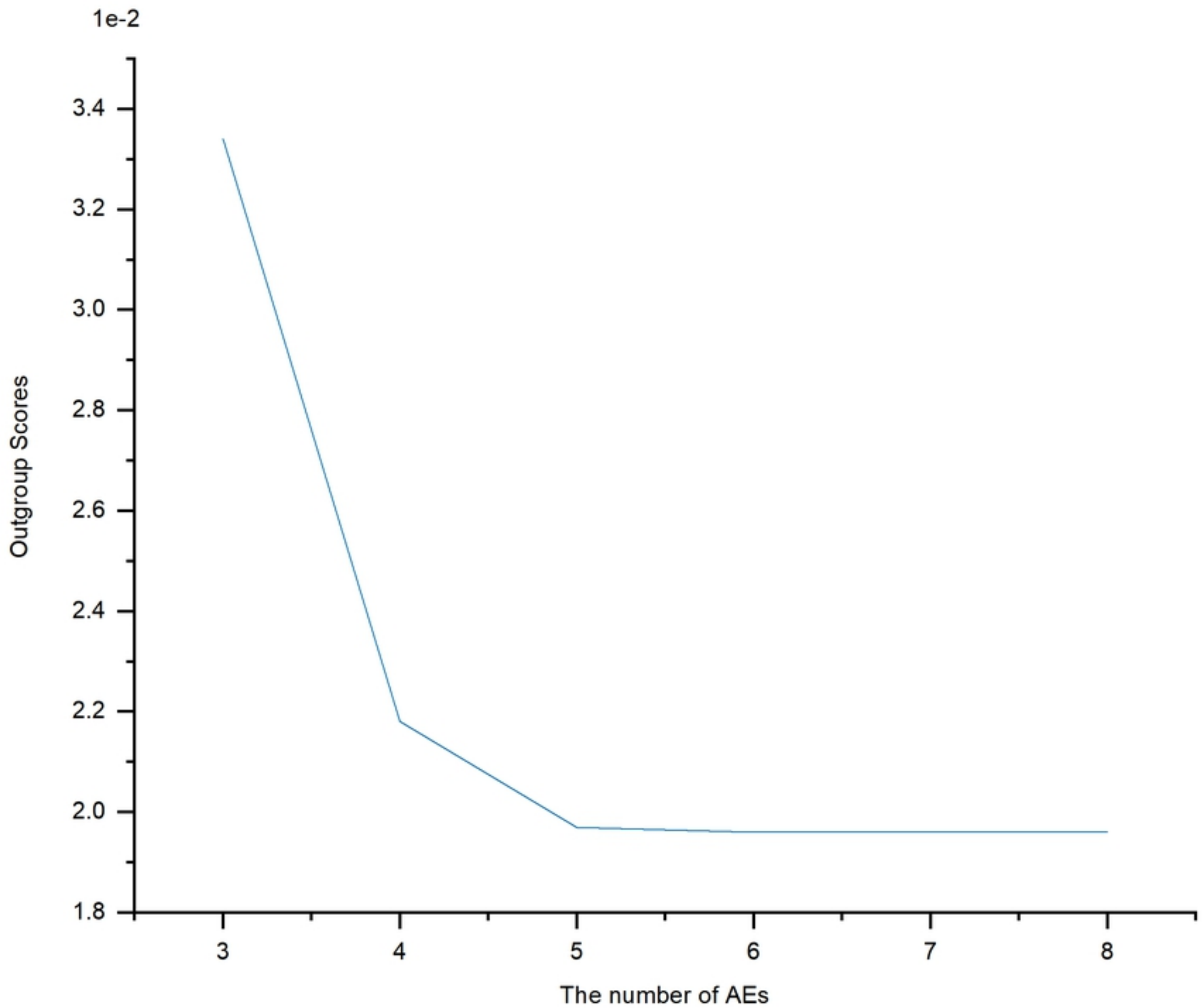


Fig5

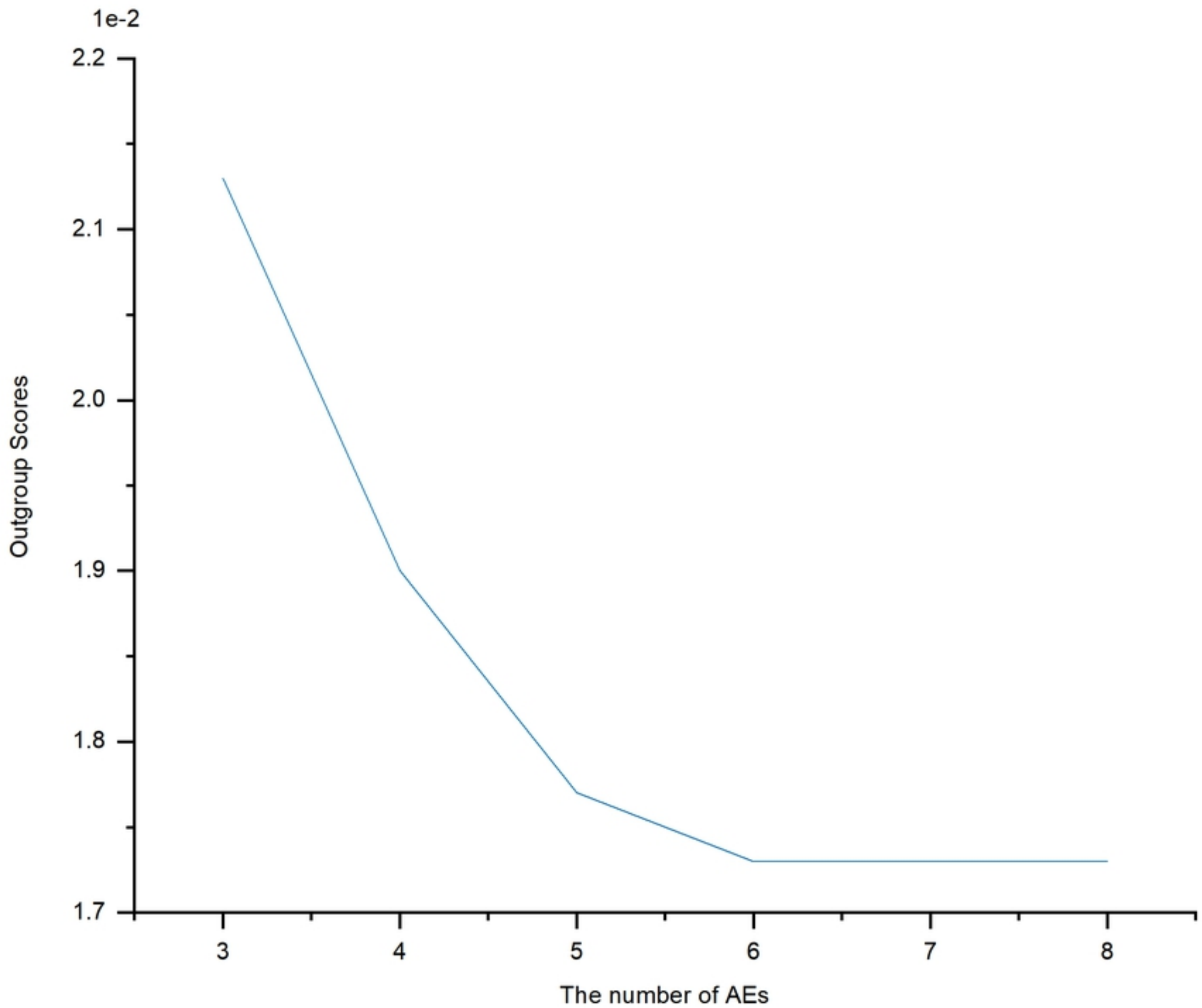


Fig6

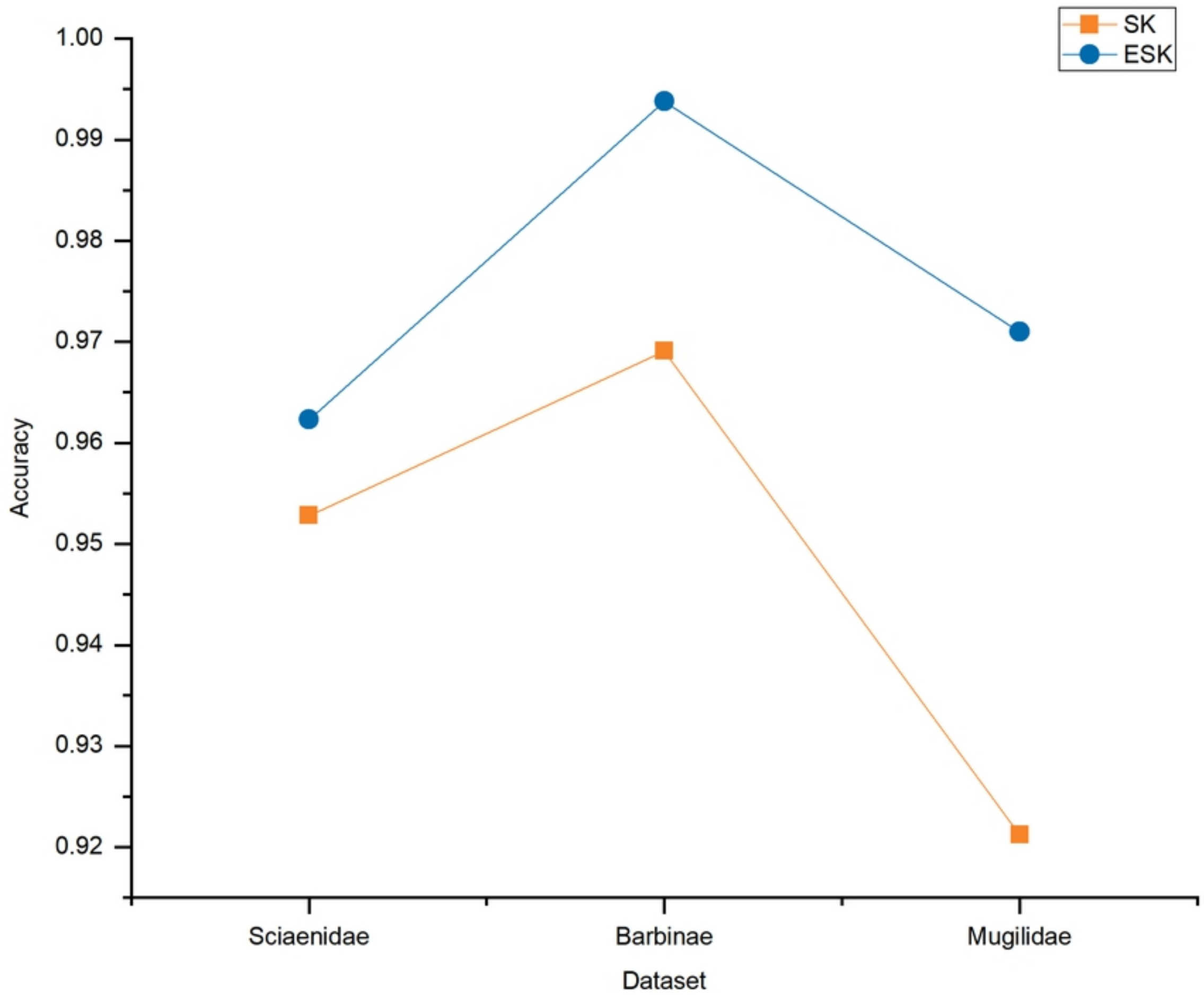


Fig7

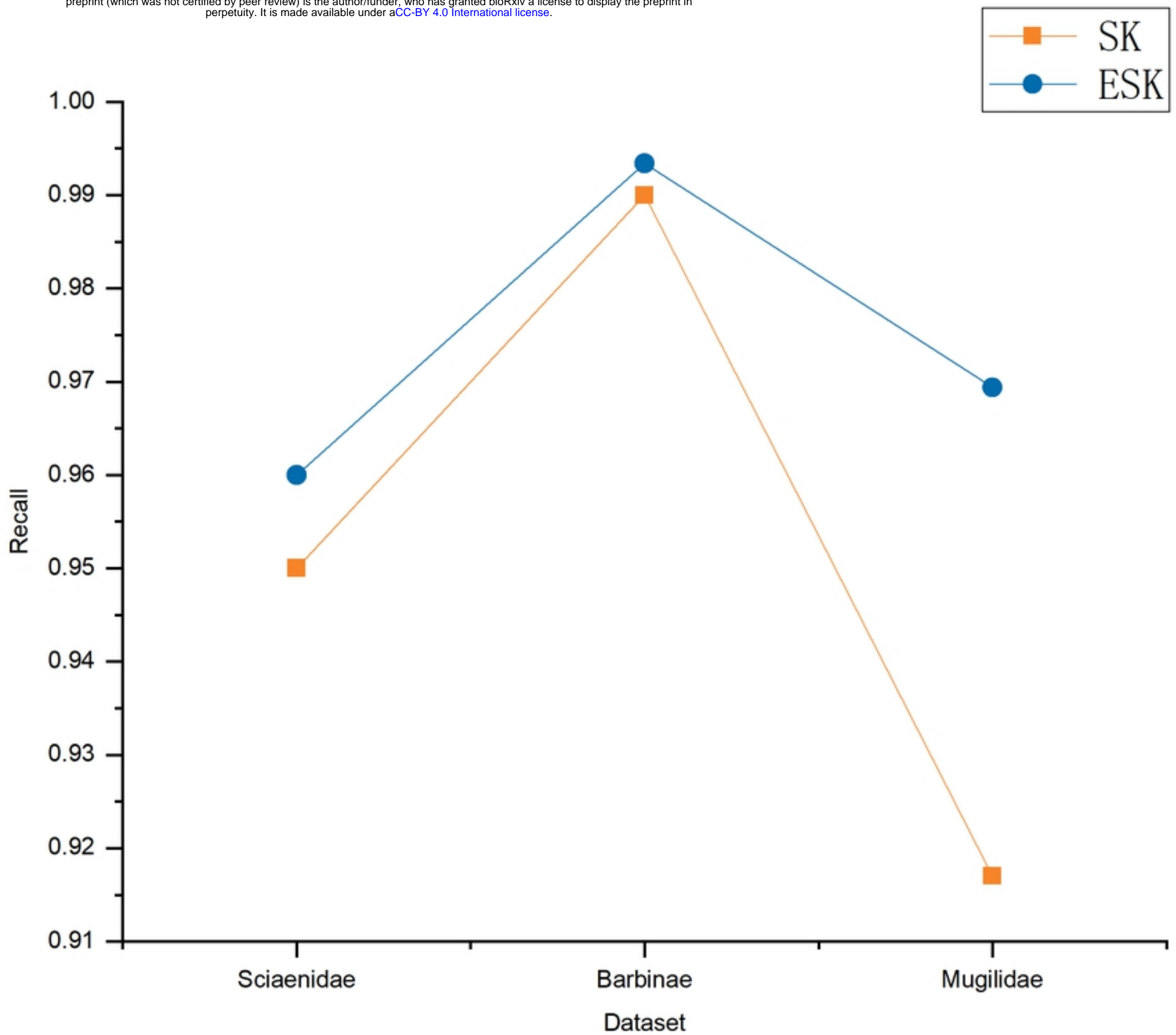


Fig8

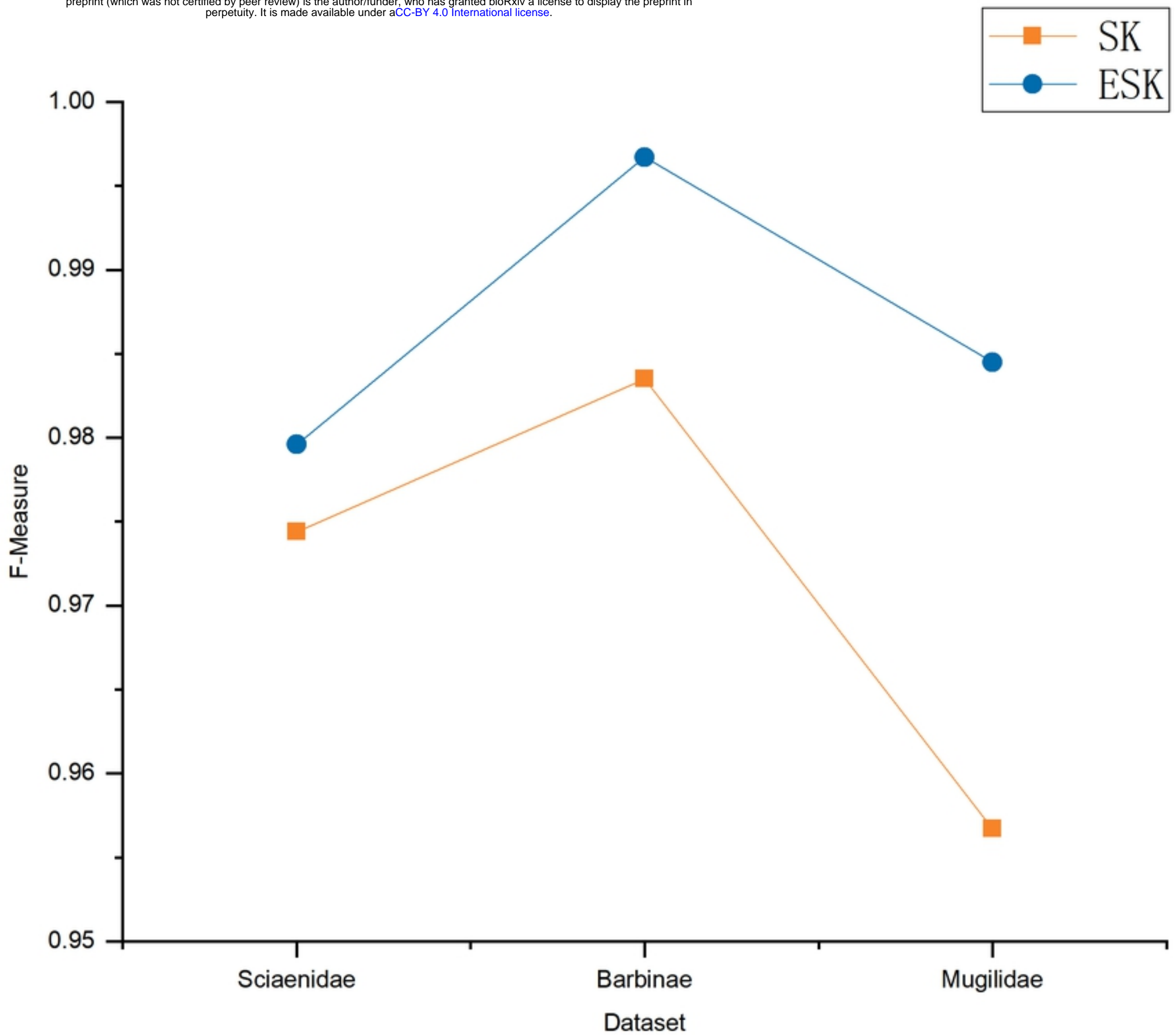


Fig9

		True Label										
		1	0	1	0	1	0	1	0	1	0	
Sciaenidae	1	73	27	86	14	95	5	82	18	96	4	1
	0	0	6	2	4	2	4	0	6	0	6	0
Barbinae	1	288	13	283	18	292	9	279	22	299	2	1
	0	0	23	6	17	0	23	0	23	0	23	0
Mugilidae	1	218	11	220	9	214	15	209	20	222	7	1
	0	0	12	0	12	0	12	0	12	0	12	0
		OC-SVM	KNN	iForest	AE	ESK						

Fig10

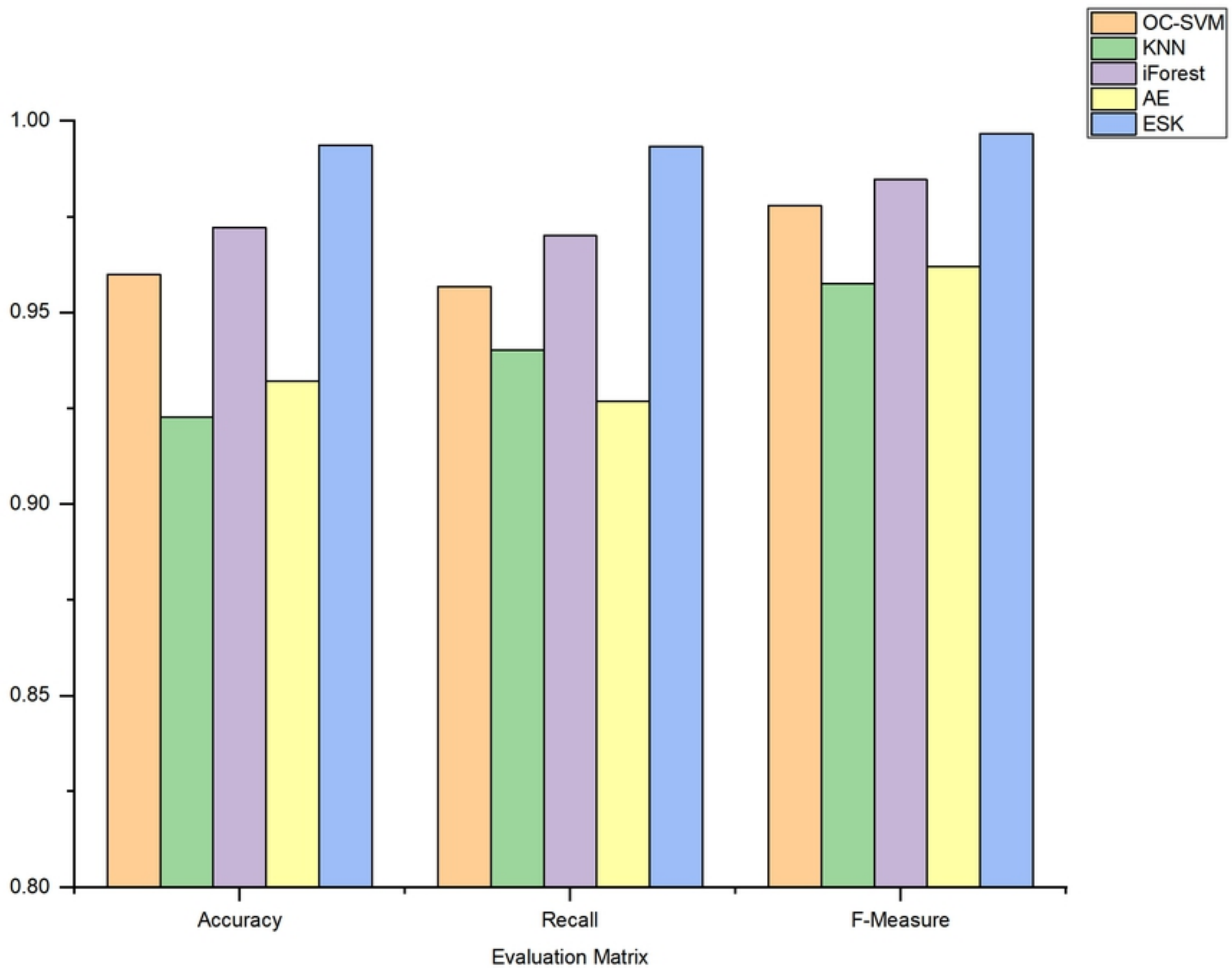


Fig11

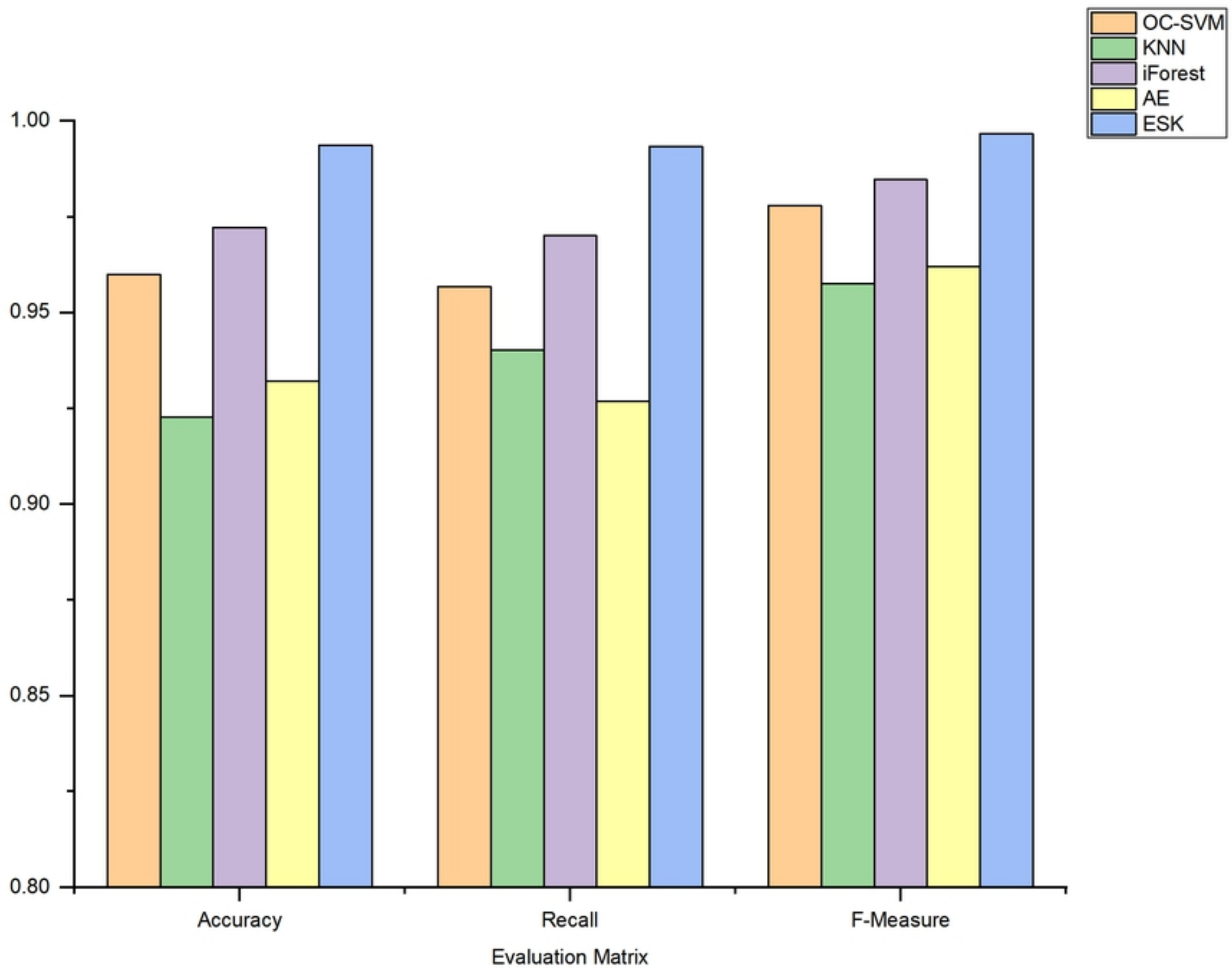


Fig12

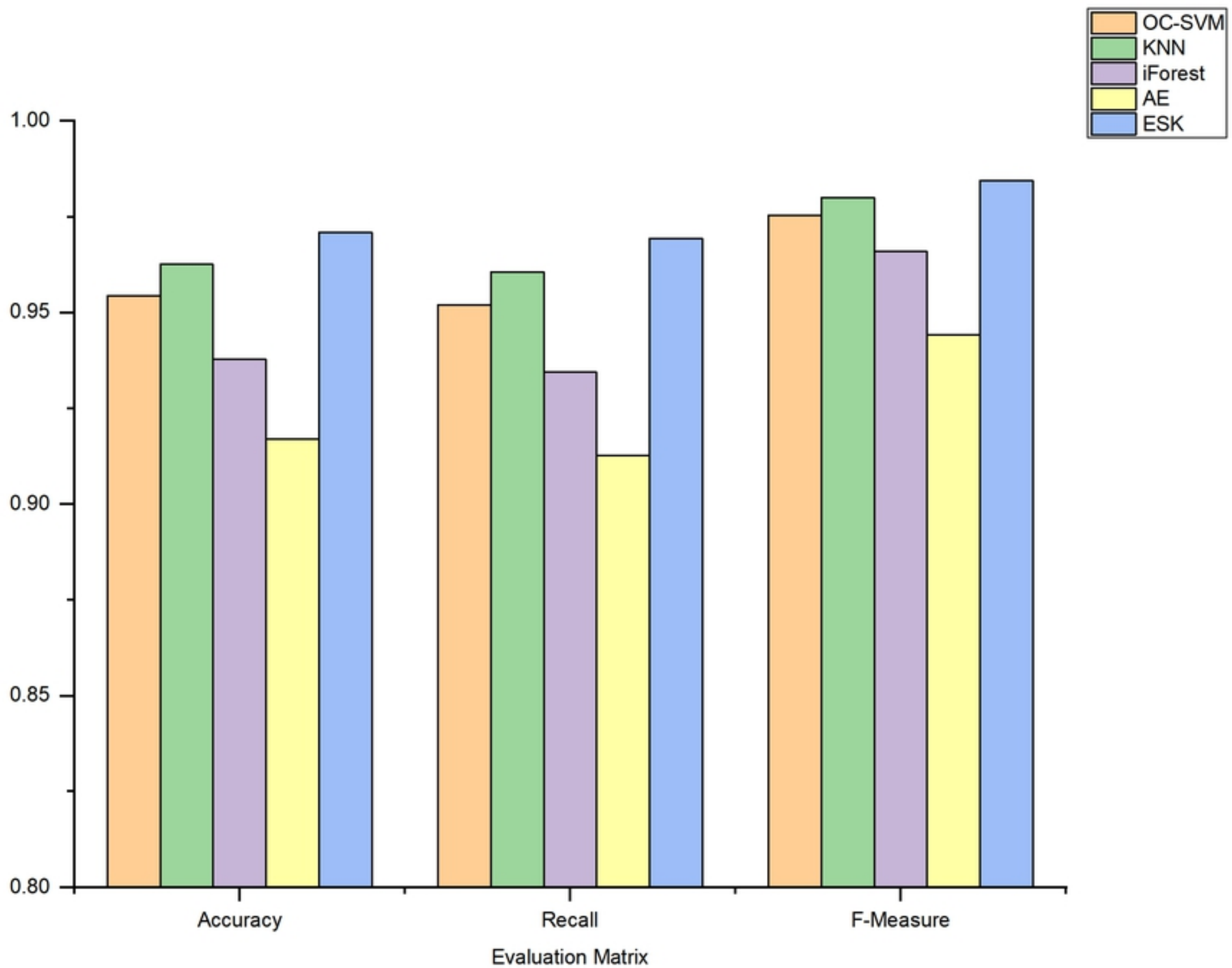


Fig13