# Automatic identification and morphological comparison of bivalve and brachiopod fossils based on deep learning

**Jiarui Sun** [1], **Xiaokang Liu** [1,2], **Yunfei Huang** [3], **Fengyu Wang** [1], **Yongfang Sun** [1], **Jing Chen** [4], **Daoliang Chu** [1], **Haijun Song** [Corresp. 1]

[1] State Key Laboratory of Biogeology and Environmental Geology, School of Earth Sciences, China University of Geosciences, Wuhan, Hubei, China

[2] Department of Biology, University of Fribourg, Fribourg, Switzerland

[3] School of Geosciences, Yangtze University, Wuhan, Hubei, China

[4] Yifu Museum, China University of Geosciences, Wuhan, Hubei, China

Corresponding Author: Haijun Song
Email address: haijunsong@cug.edu.cn

Fossil identification is an essential and fundamental task for conducting palaeontological research. Because the manual identification of fossils requires extensive experience and is time-consuming, automatic identification methods are proposed. However, these studies are limited to a few or dozens of species, which is hardly adequate for the needs of research. This study enabled the automatic identification of hundreds of species based on a newly established fossil dataset. An available "bivalve and brachiopod fossil image dataset" (BBFID, containing > 16,000 "image-label" data pairs, taxonomic determination completed). The bivalves and brachiopods contained in BBFID are closely related in morphology, ecology and evolution that have long attracted the interest of researchers. We achieved > 80% identification accuracy at 22 genera and ~64% accuracy at 343 species using EfficientNetV2s architecture. The intermediate output of the model was extracted and downscaled to obtain the morphological feature space of fossils using t-distributed stochastic neighbor embedding (t-SNE). We found a distinctive boundary between the morphological feature points of bivalves and brachiopods in fossil morphological feature distribution maps. This study provides a possible method for studying the morphological evolution of fossil clades using computer vision in the future.

# Automatic identification and morphological comparison of bivalve and brachiopod fossils based on deep learning

4   Jiarui Sun[1], Xiaokang Liu[1,2], Yunfei Huang[3], Fengyu Wang[1], Yongfang Sun[1], Jing Chen[4],

5   Daoliang Chu[1], Haijun Song[1]*

6

7   [1] State Key Laboratory of Biogeology and Environmental Geology, School of Earth Sciences,

8   China University of Geosciences, Wuhan, 430074, China

9   [2] Department of Biology, University of Fribourg, Fribourg, Switzerland

10  [3] School of Geosciences, Yangtze University, Wuhan 430100, China

11  [4] Yifu Museum, China University of Geosciences, Wuhan, 430074, China

12

13  *Corresponding author:

14  Haijun Song[1]

15  State Key Laboratory of Biogeology and Environmental Geology, School of Earth Sciences, China

16  University of Geosciences, Wuhan, 430074, China

17  Email address: haijunsong@cug.edu.cn

18

19

## Abstract

Fossil identification is an essential and fundamental task for conducting palaeontological research. Because the manual identification of fossils requires extensive experience and is time-consuming, automatic identification methods are proposed. However, these studies are limited to a few or dozens of species, which is hardly adequate for the needs of research. This study enabled the automatic identification of hundreds of species based on a newly established fossil dataset. An available "bivalve and brachiopod fossil image dataset" (BBFID, containing > 16,000 "image-label" data pairs, taxonomic determination completed). The bivalves and brachiopods contained in BBFID are closely related in morphology, ecology and evolution that have long attracted the interest of researchers. We achieved > 80% identification accuracy at 22 genera and ~64% accuracy at 343 species using EfficientNetV2s architecture. The intermediate output of the model was extracted and downscaled to obtain the morphological feature space of fossils using t-distributed stochastic neighbor embedding (t-SNE). We found a distinctive boundary between the morphological feature points of bivalves and brachiopods in fossil morphological feature distribution maps. This study provides a possible method for studying the morphological evolution of fossil clades using computer vision in the future.

**Key words:**

Fossil identification; Machine learning; Invertebrate; Morphology; Convolutional neural network.

## Introduction

40

41    Fossil identification is a fundamental task in palaeontological research and has a wide range

42    of applications, including biostratigraphic dating (Yin *et al.* 2001; Gradstein *et al.* 2012), biological

43    evolution (Alroy *et al.* 2008; Fan *et al.* 2020; Song *et al.* 2021), palaeoenvironmental

44    reconstruction (Flügel and Munnecke 2010; Scotese *et al.* 2021), and palaeoelevational estimation

45    (Su *et al.* 2019). Because taxonomic identification requires a large amount of prior knowledge as

46    a foundation, researchers need several years of training to accumulate enough experience to ensure

47    the reliability of identification. However, the actual identification process still takes considerable

48    time and is susceptible to subjective factors. The identification accuracy of some genera is even

49    lower than 80% (Hsiang *et al.* 2019). In many fields of palaeontology, deep convolutional neural

50    network (DCNN) has a significant advantage over humans, such as the identification of cut and

51    trampling marks on bones (Byeon *et al.* 2019), the discrimination of dinosaur tracks (Lallensack

52    *et al.* 2022), and the quantification of plant mimesis (Fan *et al.* 2022). To reduce the workload and

53    work difficulty for researchers, automatic fossil identification methods relying on machine

54    learning have been proposed extensively in recent years, among which models using convolutional

55    neural networks (CNNs) [e.g., VGG-16 (Simonyan and Zisserman 2014), Inception-ResNet

56    (Szegedy *et al.* 2017), GoogLeNet (Szegedy *et al.* 2015), etc.] have achieved good results

57    (Dionisio *et al.* 2020; Liu and Song 2020; Liu *et al.* 2022; Niu and Xu 2022; Wang *et al.* 2022; Ho

58    *et al.* 2023). Other supervised (e.g., Naïve Bayes) algorithms also achieved ≥70% accuracy in

59    ammonoid species identification (Foxon *et al.* 2021). This method can assist researchers in fossil

60    identification, reduce the work stress of non-palaeontologists, and enable better identification and

61    application of fossil materials in research. Furthermore, for identifying poorly preserved fossils,

62    neural networks still maintain high identification accuracy (Bourel *et al.* 2020). Neural network in

63    fossil identification is still at an early stage of development, beacuse professional palaeontologists

64    have advantages that such models do not. The ability to take into account complex contextual

65    information is one of them. But in the face of the reality that there are less and less experts on

66    taxonomy, neural network can provide a useful aid to manual identification rather than replace it

67    (De Baets *et al.* 2021). It is still worth studying, and asmore training data are available, the

68    reliability and applicability of models will become better.

69        The training of automatic taxonomy identification models (ATIM) requires a large dataset of

70    labelled fossil images (Liu *et al.* 2022). In general identification field, machine-learning datasets

71    contain millions of images (e.g., ImageNet dataset), which far exceed fossil datasets. The lack of

72    high-resolution (genus-level) fossil labels in the field of palaeontology is mainly due to the tedious

73    and time-consuming process of dataset building. Machine learning has now achieved good results

74    in fossil identification (above the genus level). Liu and Song (2020) achieved 95% accuracy for

75    22 fossil and abiotic grain groups during carbonate microfacies analysis. While 90% accuracy was

76    achieved in the automatic identification of 50 fossil clades relying on web crawlers (Liu *et al.*

77    2022), genus- and species-level automatic identification focused mainly on a few taxa (mostly <

78    10). Dionisio *et al.* (2020) performed automatic identification of 9 radiolarian genera, obtaining

79    91.85% accuracy. Niu and Xu (2022) performed automatic identification of fossils covering 113

80    graptolite species or subspecies. However, similar studies targeting a large number of taxa are less

81  common (Fig. 1, details of the relevant studies can be found in the Appendix S1). In practice, it is

82  common to identify a large number of fossil categories. However, current automatic identification

83  studies are limited to a few or dozens of taxa, which is hardly adequate for the needs of research.

84  There is a gap in automatic fossil identification studies for hundreds of taxa. This study provides

85  new practice in this field based on a newly established fossil dataset. In addition, previous studies

86  all focus on the same fossil clade (e.g., radiolarians, brachiopods, etc.). It is unclear whether fossils

87  in different phyla can achieve automatic fossil identification.

88      Brachiopods and bivalves are the two most common invertebrate clades in the Phanerozoic

89  (Sepkoski 1981; Clapham *et al.* 2006; Benton and Harper 2020). Brachiopods are the dominant

90  fossil animals of the Paleozoic, but their diversity is now far less than that of bivalves (Thayer

91  1986). The start of this transition occurred at the Permian-Triassic mass extinction (PTME), when

92  the marine benthic faunas changed from brachiopod-dominated Paleozoic evolutionary fauna to

93  mollusk-dominated modern evolutionary fauna (Fraiser and Bottjer 2007; Dai *et al.* 2023). The

94  reasons for the dominance of bivalves over brachiopods have long attracted the attention of

95  palaeontologists (Ballanti *et al.* 2012; Payne *et al.* 2014). The similarities and differences between

96  them in morphology and physiological mechanisms may be an important perspective. Whether

97  bivalves and brachiopods influenced each other evolutionarily is a controversial issue, also known

98  as "ships that pass in the night" (Gould and Calloway 1980; Fraiser and Bottjer 2007; Liow *et al.*

99  2015). Bivalves feed more efficiently at high algal concentrations than articulate brachiopods,

100  which is thought to be the reason for the physiological perspective (Rhodes and Thompson 1993).

101  Morphologically, the prosperity of the bivalves after PTME cannot be attributed to their

102    morphological innovations (Fraiser and Bottjer 2007), while bivalves suppressed brachiopod

103    evolution (Liow *et al.* 2015). However, they still have certain similarities, for instance they both

104    tend to become smaller under heat stress (Piazza *et al.* 2020). The close relationship and significant

105    differences between them have attracted researchers' interest, and conducting morphological

106    studies is the first step. However, the similar morphological features between them have caused

107    problems for researchers to identify them accurately. Automatic identification of brachiopods has

108    been carried out previously. Wang *et al.* (2022) used the transposed convolutional neural network

109    to realize the automatic identification of fossils with a relatively small dataset and they achieved

110    97% identification accuracy for five brachiopod species based on 630 training images. In this

111    study, we enabled the automatic identification of hundreds of taxa (bivalve and brachiopod) based

112    on a newly established fossil dataset. We built a "bivalve and brachiopod fossil image dataset"

113    (BBFID) (16,596 labelled fossil images covering 870 genera and 2033 species) for the first time

114    by collecting and sorting a large amount of published literature. We built ATIMs using transfer

115    learning in VGG-16 (Simonyan and Zisserman 2014), Inception-ResNet-v2 (Szegedy *et al.* 2017),

116    and EfficientNetV2s (Tan and Le 2021) architectures, which have performed well in general

117    identifications. Furthermore, we extracted the process outputs of the model as fossil features and

118    downscaled them to two-dimensional data using t-SNE (Van der Maaten and Hinton 2008).

119    Plotting them in a two-dimensional space is an effective way to compare morphological

120    differences between bivalves and brachiopods.

## Materials and Data

121

122　The BBFID used for training ATIMs contains bivalve-part (BBFID-1) and brachiopod-part

123　(BBFID-2), all collected from published literature and monographs (see Appendix). Detailed data

124　on the number of each taxon are given in the Appendix (Table S1, S2). This study collected fossil

125　images from publications that were of diverse origin. This makes use of the large amount of data

126　that already exists and allows for better use of data from previous studies.

127　We used Adobe Acrobat Pro DC to capture accurately named bivalve and brachiopod fossil

128　images from the collected literature and saved them as bmp, jpg, or png images to minimize the

129　quality loss of the images. These fossils are Carboniferous ($< 1.5‰$), Permian (majority), Triassic

130　(majority), Jurassic ($< 1.5\%$), and Quaternary ($< 0.5‰$) in age. Permian and Triassic fossils make

131　up the vast majority of the dataset. Their overlapping occurrences, having undergone the same

132　geological events, are of great importance in fossil identification and in studies of morphological

133　evolution. Those that could not be saved due to the encryption of PDF files in the literature were

134　screenshotted as png files using Snipaste. The majority of images collected from plates are single

135　animal images, and the effect of plate numbering was avoided as much as possible.

136　We obtained more than 16,000 fossil images from 188 publications and performed data

137　cleaning. The contribution of each publication to the dataset is given in the Appendix (Appendix

138　S2). The contribution of publications is sufficiently balanced for the training model. During the

139　data collection stage, we collected as many fossil images as possible. These images were taken at

140　any viewpoint and in any orientation. Different views of the same specimen were treated as

141    different instances and labelled separately. To ensure the reliability of the dataset, we checked the

142    bivalve and brachiopod images and corresponding labels. Because the taxonomic system of

143    bivalves and brachiopods is continuously improving (Konopleva *et al.* 2017; Sulser *et al.* 2010),

144    we categorized the genera whose taxonomic names and positions had been changed in the

145    literature. Additionally, we removed poorly preserved fossil images, which contain two cases. The

146    first case is images with uncertain taxonomic names. The other discarded images are obtained from

147    scanned published documents (mostly monographs published in the last century) that are poorly

148    pixelated and difficult to identify even for palaeontologists. In both cases, the ambiguous images

149    are discarded based on whether the experts can distinguish the fossils or not. There is no filtering

150    based on deep learning preference, so this operation does not affect the utility of the deep learning

151    method.

152         Our dataset was randomly divided into the training set (60%), validation set (20%), and test

153    set (20%) to train, tune, and test the model. Such a distribution is intended for the test set to cover

154    a sufficient number of taxa to make the accuracy more reliable. Because the validation set is used

155    as a reference for the tuning process, the identification accuracy of this part may have artificial

156    bias and is not universally meaningful. Thus, the final accuracy was calculated using a separate

157    test set to evaluate model performance.

158         The final BBFID contains 870 genera, with 16,596 sets of "image-label" data pairs. All images

159    have genus labels, with 14,185 items having higher-resolution species labels. BBFID-1 contains

160    379 genera and 889 species, with 8,144 sets of image-label data pairs. BBFID-2 contains 491

161    genera and 1,144 species, with 8,452 sets of data pairs. A total of about 2,300 genera of bivalves

162   have been described to date, and 1,700 genera of brachiopods have been described (Pitrat and

163   Moore 1965; Nevesskaja 2003). BBFID covers 16.4% and 28.8% of the described bivalve and

164   brachiopod genus-level classifications, respectively. Genus distributions of BBFID are shown in

165   Figure 2. The BBFID-1 dataset consists of ~85% black and white images, and the rest are in colour.

166   In situ photographs of fossils (images with rocks in the background) occupy ~25% of BBFID-1,

167   while all other fossil photographs have plain white/black backgrounds. The BBFID-2 dataset

168   consists of ~95% black and white images, and the rest are in colour. In situ photographs of fossils

169   (images with rocks in the background) occupy ~1% (61 images) of BBFID-2, while all other fossil

170   photographs have plain white/black backgrounds. The BBFID-2 (brachiopod dataset) has more

171   plain white/black background photographs, because brachiopod fossils are more robust and easier

172   to preserve intact than bivalve fossils. Therefore the former can be imaged to obtain more complete

173   pictures of the fossils without the rocky background.

174        To meet the requirements of machine learning, each taxon should have at least three items.

175   Therefore, we chose the categories with > 2 items of BBFID to perform the model training, which

176   contains 16,389 sets of "image-label" data pairs. BBFID contains images of the whole shells and

177   detailed images. Detailed images refer to all non-full shell face images as well as photographs not

178   in front view, such as structures of fossils. Of the selected images, detailed images occupy ~18%

179   in BBFID-1, ~40% in BBFID-2, and ~29% in the overall BBFID. The number of detailed images

180   (the categories with > 2 items) and the exact number of detailed images in each dataset (training

181   set, validation set, and test set) of the common genera are available in the Appendix (Table S1,

182   Table S2).

## Methods

*Convolutional Neural Network*

Convolutional neural networks (CNNs) perform well in general recognition and have been used in the automatic identification of palaeontological fossils (Dionisio *et al.* 2020; Liu and Song 2020; Kiel 2021; Liu *et al.* 2022; Niu and Xu 2022; Wang *et al.* 2022; Ho *et al.* 2023). In this study, three pre-trained models of convolutional neural networks with good classification performance on the ImageNet dataset (Deng *et al.* 2009) namely VGG-16 (Simonyan and Zisserman 2014), Inception-ResNet-v2 (Szegedy *et al.* 2017), and EfficientNetV2s (Tan and Le 2021) were selected and suitably modified (Fig. 3). VGG-16 and Inception-ResNet-v2 have been proven to automatically identify fossils and perform well (Hsiang *et al.* 2019; Liu *et al.* 2022). We retained their main architecture, removed the top softmax layer and/or fully connected layer depending on fossil categories, and added a fully connected layer (with 256 output and Relu activation function), batch normalization layer (Ioffe and Szegedy 2015), dropout layer (with rate = 0.2), and fully connected layer (with output as fossil categories) (Fig. 3).

In fossil identification, CNNs first decode the fossil images to obtain the tensor that can be operated, and the model operates on these values to establish the correspondence between the image data and the fossil name. CNNs use convolutional operations to process image data and gradient descent to minimize the loss function to train the model (LeCun *et al.* 1998). The neural network can be divided into multiple network layers. More specifically, the convolutional, pooling, and fully connected layers play a crucial role in the automatic identification process. The

203    convolutional layers transform an image by sweeping a kernel over each pixel and performing a

204    mathematical operation. The pooling layer reduces the amount of computation, making the model

205    easier to train (Giusti *et al.* 2013). The fully connected layer and activation function (Relu) fit the

206    correspondence between fossils and labels (Nair and Hinton 2010) and output the predicted labels

207    and probabilities we need at the top layer.

208         VGG-16 is a classic DCNN proposed by Simonyan and Zisserman (2014), which uses 16

209    layers and 3 × 3 convolutional kernels (convolution filters) to achieve good performance. Then,

210    He *et al.* (2016) proposed a new residual connectivity method and applied it to Inception-ResNet-

211    v2, which makes the network easier to optimize and allows the use of a deeper network to improve

212    performance. EfficientNetV2 is currently a more advanced open-source image classification model

213    using the training-aware neural architecture search and scaling method to improve training speed

214    and parameter efficiency (Tan and Le 2021).

215    *Data preprocess*

216         Deep learning models have requirements for input data size. However, the images in our

217    dataset were of different sizes and the labels were also inappropriate for model training. Thus, data

218    needed to be preprocessed. To match the model's requirement, all images were resized to a uniform

219    size (slightly different depending on the model: VGG-16, 224*224; Inception-ResNet-v2,

220    299*299; EfficientNetV2s, 384*384.). Further details are shown in Fig. 3. To improve their

221    generalization ability, we randomly adjusted the image (training set and validation set) brightness

222    (within ± 0.5) and contrast (within 0 to + 10) (Simonyan and Zisserman 2014; Szegedy et al. 2015;

223    He et al. 2016; Liu and Song, 2020). In addition, the images were normalized and standardized [all

224    images were processed using the following equation: x_new=(x-mean)/std, mean=(0.485, 0.456,

225    0.406), std=(0.229, 0.224, 0.225). Mean and std are empirical values, which are calculated from a

226    large number of images]. We conducted the discrete one-hot coding for image labels. Finally, a

227    one-to-one correspondence between the images and the labels was established, and we obtained

228    the processed machine-learning dataset.

229    *Training Methodology*

230        Achieving high accuracy in multiclass fossil identification using neural networks requires a

231    large dataset as a basis. Although we built the bivalve and brachiopod dataset manually, it was still

232    insufficient to train a model with random initialization of parameters to converge and achieve the

233    best results. Therefore, we applied transfer learning in the model training process, an effective way

234    to train a model on a small dataset (Tan *et al.* 2018; Brodzicki *et al.* 2020; Koeshidayatullah *et al.*

235    2020). Transfer learning uses parameters trained by general identification tasks for initialization

236    to accelerate the convergence of the new model. It is feasible to use this to reuse the general

237    identification model parameters for palaeontological fossil identifications (Pires de Lima *et al.*

238    2020). This is why we only envision applying this method to the automatic identification of

239    common fossils, while fossils with too few specimens will still need to rely on palaeontologists.

240        In this study, each model was loaded with pre-trained parameters that were originally trained

241    on ImageNet. This method greatly reduces the amount of data required for automatic identification,

242    greatly expanding their application scenarios.

243    We coded in Python and relied on the Tensorflow scientific computing library (Abadi *et al.*

244    2016) to train the model. The training process was performed using the Adam optimizer (Kingma

245    and Ba 2014). The loss function uses the categorical cross-entropy loss function (Botev *et al.*

246    2013), and the accuracy is used as an evaluation metric during training. The final model

247    performance is presented using a confusion matrix and F1 score (Figs. 4-6). The confusion matrix

248    contains recall and precision, which represent two perspectives of identification performance.

249    Recall represents the proportion of "items correctly identified as a specific taxon" to "the total

250    items belong to that taxon". Precision represents the proportion of "items correctly identified as a

251    specific taxon" to "the total items identified as that taxon". The F1 score is the harmonic mean of

252    the recall and precision, which can represent both false positives and false negatives (Sarkar *et al.*

253    2018).

254    To facilitate training, the learning rate is adjusted with validation loss in training. If "validation

255    loss" down less than 0.0001 lasts 5 epochs, the learning rate will be halved using

256    "callbacks.ReduceLROnPlateau()" function. Additionally, to prevent overfitting, EarlyStopping

257    (a method to stop training when the model performs optimally) was set to ensure the good

258    performance of the model in the test set (Ying 2019). During the training process, the model saves

259    architecture and parameters with the highest accuracy in the validation set in real-time for rapid

260    deployment in subsequent applications. Because BBFID contains both the genus tags and species

261    tags, we set the model to the genus mode (only read the genus tag) and species mode (read both

262    genus tag and species tag) during model training and testing. Because of dataset size, the model's

263    architecture and hyperparameters significantly affect its performance; thus, we trained models and

264 compared their performance under different scenarios (Table 1).

265 We chose the different sizes of the datasets to train models according to the taxonomic levels.

266 At the genus-level, we set three scales to explore model performance using different volumes of

267 datasets. These three scales are the number of each genus > 100 images (scale A), > 50 images

268 (scale B), and > 10 images (scale C) (Table 2). Among them, scale B/C contains all genera with

269 more than 50/10 pictures, the same for other scales. The numbers of taxa in BBFID-1 are 13 (scale

270 A), 34 (scale B), and 156 (scale C), whereas the numbers of items in BBFID-2 are 9 (scale A), 32

271 (scale B), and 223 (scale C). They display a clear gradient to match our research needs. For the

272 selection of data adequacy (i.e., data scale) at the species level, we selected scale B (number of

273 each species > 50 images) and scale C (number of each species > 10 images) for training and

274 testing, according to the performance of the genus mode. Furthermore, we also tried two larger

275 scales: scale D (number of each species > 8 images) and scale E (number of each species > 6

276 images). There are four gradients in total to find the range that covers more genera with guaranteed

277 accuracy. In addition, for BBFID, we added two larger scales (the number of each taxon > 4 images

278 and > 2 images) to explore the model performance in small datasets. As mentioned earlier, all data

279 (scales A, B, C, D, and E) were randomly divided into the training set, validation set, and test set

280 in the ratios of 60%, 20%, and 20%, which is the ideal situation. In order to try a larger data scale,

281 we discarded the requirement that the validation set cover all species. Therefore, the number of

282 single-taxon images > 2 was the maximum data size we could try, because all taxonomic units

283 shall be covered in the training set and test set.

284 Model architecture plays a pivotal role in models. Thus, we used BBFID-1 (scale A) to test

285   model identification accuracy at the genus level under three different model architectures (i.e.,

286   VGG-16, Inception-ResNet-v2, and EfficientNetV2s). Subsequently, the best architecture was

287   selected to build ATIM, trained and tested using different scales of BBFID-1, BBFID-2, and

288   BBFID, respectively, to obtain the corresponding model performance (Table 2).

289       Considering that a particular identification model cannot identify arbitrary fossil taxa, it is

290   necessary to establish a method for measuring the applicability of the model. We divide the entire

291   BBFID into "applicable" and "inapplicable". Anything in the training set is considered "applicable"

292   and anything not in training set is considered "unapplicable". Binary classification training based

293   on Inception-ResNet-v2 was performed and the "Applicability Model" (AM) was obtained. Users

294   can use the AM to determine the applicability.

295   *Dimensionality reduction method*

296       In this study, we employed a downscaling method of t-SNE that uses a probability measure of

297   similarity and expresses probabilities as spatial distances (Van der Maaten and Hinton 2008). To

298   compare fossil morphology, we extracted the output of the last maximum pooling layer as fossil

299   features and downscaled the high-dimensional data of fossil features to a two-dimensional plane

300   using t-SNE. Next, we visualized that to analyze easily the morphological differences and

301   similarities between bivalves and brachiopods. The model training and downscaled visualization

302   codes were referenced from some open-source projects (Liu and Song 2020; Liu *et al.* 2022).

303 **Results**

304 *Model performance between different architectures and hyperparameters*

305      Different architectures perform differently using BBFID-1 (scale A, genus level), with the

306 best performance of 83.02% obtained with the EfficientNetV2s architecture and the

307 corresponding hyperparameters (Table 1). The results of confusion matrix for this identification

308 task are shown in Figure 4. The identification recalls were > 79% for all categories except the

309 genera *Pteria* (0.71, test set: 28 items), *Bakevellia* (0.72, test set: 47 items), and *Halobia* (0.72,

310 test set: 65 items), whereas the accuracies of *Quemocuomegalodon* (1.00, test set: 21 items)*,* and

311 *Monotis* (0.91, test set: 67 items) exceeded 90%.

312 *Model performance using different data scales*

313     The three model architectures (VGG-16, Inception-ResNet-v2, and EfficientNetV2s) were tested

314 in  BBFID-1  and  the  EfficientNetV2s  architecture  was  found  to  perform  best.  We  used

315 EfficientNetV2s architecture that performed well on BBFID-1 and corresponding hyperparameters

316 to build other models (genus mode) using different data scales, which performed as expected under

317 different datasets (Table 2). The accuracy of BBFID-1 (scale A) was 82.10%, whereas those of

318 scales B and C were 71.73% and 58.34% respectively, with the loss increasing by decreasing

319 accuracy for all three. The accuracy of BBFID-2 was 85.43%, 71.35%, and 50.04% for the three

320 dataset scales, whereas the identification accuracy of scale A exceeded 85%. Furthermore, in four

321 categories, more than 90% of images were identified correctly (Fig. 5). The accuracy of model

322    training by BBFID was 81.45%, 70.66%, and 53.71% at the three scales, and the performance of

323    each scale was similar to the performance of the corresponding bivalve and brachiopod individual

324    identifications. In species mode, the models also performed similarly (Table 2), with the accuracy

325    of BBFID at scale C (148 categories for bivalves, 195 categories for brachiopods) of more than

326    60% (see Appendix S3 for confusion matrix and evaluation metrics). The accuracies of Scale D

327    (bivalve 179 categories, brachiopod 265 categories) and scale E (bivalve 241 categories,

328    brachiopod 396 categories) ranged from 51% to 59%. All these models in EfficientNetV2s

329    architecture met the early stopping condition and terminated training before 50 epochs, and the

330    training set accuracy was close to 100% at this point. This indicates that the models completed

331    fitting to the training set. The training process of BBFID (scale A) shows that the model basically

332    converged about 20 epochs (Fig. 7), and its training set accuracy finally reached ~100%, while the

333    maximum validation accuracy was over 80% (Table 2).

334       We extracted the process output from the ATIM (Order 22) and summed the same point data

335    in each dimension to draw a feature map (Fig. 8). We also used the output of the top maximum

336    pooling layer as fossil features and then used t-SNE (Van der Maaten and Hinton 2008) for

337    dimension reduction, which achieved good results of morphology clustering and comparison (Fig.

338    9). The classification of each taxon in Figure 9 is clear, and the t-SNE results are similar between

339    the training set (Fig. 9A) and the validation set and test set (Fig. 9B). However, the individual

340    clusters obtained from the training set are more concentrated and the boundaries between different

341    categories are clearer than the latter due to the training process (Fig. 9). Additional t-SNE

342    calculation for more categories (444 categories, based on Order 34) was also performed (see

343    Appedix).

344    **Discussion**

345    *Identification accuracy*

346        The ranking of automatic identification performance among three architectures trained by

347    BBFID-1 (Table 1) is comparable to general task results (Simonyan and Zisserman 2014; Szegedy

348    *et al.* 2017; Tan and Le 2021), indicating that transfer learning is useful. It is feasible to apply pre-

349    trained parameters of the general model to the ATIM in the field of palaeontology using transfer

350    learning. The identification accuracy (> 80%) on genus mode is similar to some previous studies

351    that built upon ResNet architecture (Romero *et al.* 2020). Romero *et al.* (2020) achieved an

352    accuracy of 83.59% using the external morphology of pollen grains, increasing to 90% with the

353    addition of an internal structure using Airyscan confocal superresolution microscopy. Adding the

354    sequential internal structure of bivalves and brachiopods may be a way to improve identification

355    accuracy.

356        The fossil images used in this study contain pictures of the whole shells and detailed pictures,

357    such as structures of fossils. The detailed images contain different information than the whole shell

358    images. Since no specific labels have been added to the detail images, the identification accuracy

359    was adversely affected by this factor. For the accuracy of different parts of the dataset, the accuracy

360    of the validation set was comparable to that of the test set, but lower than that of the training set.

361    Because the model was trained using the training set, the identification performance was better in

362  this part. However, the data from the validation and test sets were not used to train the models.

363  Accordingly, the results were slightly worse compared with the training set. Furthermore, the

364  validation set was purposefully optimized in the conditioning.

365       The accuracy of the model using selected architecture and parameters (Table 1, Order 11) on

366  genus mode exceeded 80% using BBFID-1 (scale A). In contrast, the accuracy decreases between

367  scale B and scale C stems from the single taxon images decrease and confusion caused by the

368  categories increase. Nevertheless, the identification accuracy of scale C (156 categories) was still

369  close to 60%. In addition, the model based on BBFID-2 achieved similar accuracy to the model

370  based on BBFID-1 at all scales. The identification accuracy at scale A exceeded 80%, which is

371  close to or even exceeds the identification level of palaeontologists (Hsiang *et al.* 2019). Hsiang

372  *et al.* (2019) collected the accuracy of foraminiferal identification by palaeontologists and found

373  that human accuracy is only 71.4%, which is lower than automatic identification (87.4%). Another

374  study of planktonic foraminifera covering 300 specimens reported an average identification

375  accuracy of <78% for 21 experts (Al-Sabouni *et al.* 2018). In an automatic identification of modern

376  dinoflagellates, the expert's accuracy was also only 72% (Culverhouse *et al.* 2003). Austen *et al.*

377  (2016) found that the accuracy of experts in bumblebees was even lower than 60%. It must be

378  noted, however, that the above-mentioned studies differ from this study in terms of the taxa and

379  there may be differences in the difficulty of identification.

380       As mentioned previously, this study achieved automatic identification of fossils including 22

381  genera of bivalves and brachiopods, with a test set accuracy > 80%. The obtained model performed

382  relatively well considering the volume of categories and datasets in this task. Dionisio *et al.* (2020)

383    also trained a model for identifying radiolarian fossils (containing only nine genera with 929

384    photographs) automatically. The accuracy of the CNN model is 91.85%, higher than ours. The

385    average number of images per genus used in this study was comparable to ours; however, they

386    used SEM photographs from the same source. Fewer extraneous factors and fewer categories

387    might have contributed to slightly higher accuracy. Models for the automatic identification of

388    pollen from 16 genera were also proposed with accuracies between 83% and 90%, also using

389    microscopic images (Romero *et al.* 2020).

390        Moreover, models based on BBFID performed similarly to the models based on the

391    corresponding scale of BBFID-1 or BBFID-2, which indicates that the ATIM is not easily affected

392    by the similar morphology between bivalves and brachiopods with sufficient data volume (as

393    further demonstrated by the confusion matrix). The models are highly reliable in bivalves and

394    brachiopods identification at the genus level, which provides a basis for our subsequent

395    comparison of their morphology. Moreover, the identification accuracy of BBFID (scale C,

396    including 379 taxa) was 53.71%, which is understandable considering the large taxonomic unit

397    number with the relatively limited training set. Large-scale automatic fossil identification based

398    on a small dataset is feasible. However, it must be noted that the categories with fewer figures are

399    more concentrated in the literature, which might have led to the similarity between the test set and

400    the training set. Thus, these accuracies cannot objectively generalize the performance and ability

401    of models.

402        Regarding species-level automatic identification performance, we achieved an accuracy of

403    82.83% for 16 species identification, with several species attributed to the same genus with

404    relatively similar morphology. Although Kong *et al.* (2016) automatically identified three pollen

405    species of the same genus in a confusing species classification task with 86.13% accuracy, it must

406    be noted that their pollen task relied more on confusing information such as a texture for

407    identification. Importantly, the identification accuracy of mixed data scale C in the species mode

408    is similar to, or even slightly higher than, that in the genus mode. This implies that the number of

409    taxonomic categories can have a greater impact on automatic identification performance relative

410    to the differences between taxonomic units. The relationship between the number of categories

411    and the accuracy corroborates this (Appendix S4). The two correspond well to the logarithmic

412    relationship ($R^2$=0.8975).

413        Although we independently built a dataset containing >16,000 images, it is still small for

414    machine learning. Most studies in automatic fossil identification have focused on a few categories

415    and large sample sizes (Liu and Song 2020; Liu *et al.* 2022; Niu and Xu 2022; Wang *et al.* 2022;),

416    which undoubtedly helps improve performance. Niu and Xu (2022) used a dataset of 34,000

417    graptolites to perform an automatic identification study of 41 genera, which resulted in 86%

418    accuracy. In contrast, the identification accuracy of 47 genera in this study was 76.26%, which

419    demonstrates the importance of larger data sets.

420    *Analysis of identification results*

421        We tested models in genus mode using BBFID-1, BBFID-2, and BBFID (scale A) and

422    obtained a confusion matrix (Figs. 4, 5, 6), which truly reflects the model performance and

423    misidentification. Example images of all 22 genera in this scenario are shown in the Appendix S5

424    for a better comparison of morphological differences. In the confusion matrix, the vertical axis

425    represents the "true" genus name, whereas the horizontal axis represents the "predicted" genus

426    name. The numbers in the matrix represent the proportion of "true" genera identified as "predicted"

427    genera, and the larger the proportion, the darker the squares. The model performed well in the

428    automatic identification of bivalves and brachiopods respectively, and misidentification was

429    maintained at a low level.

430        In the hybrid auto-identification model (i.e., model based on BBFID), the overall performance

431    was good although the accuracy (81.90%) decreased slightly compared to the separate auto-

432    identification accuracies of bivalves and brachiopods (i.e., accuracy testing by BBFID-1 or

433    BBFID-2). The genus *Quemocuomegalodon* maintained a high identification recall (1.00) in the

434    bivalve categories, whereas the recall of *Proyalina* increased from 0.88 to 0.92. Other categories

435    decreased slightly. Most of the brachiopod categories showed significant or stable increases,

436    whereas only two genera exhibited recall decreases (*Araxathris* from 0.76 to 0.68 and *Paryphella*

437    from 0.77 to 0.72). The change in the recall may be related to the change in the distribution of the

438    training set. Among these misidentified categories, two cases were distinctive, each exceeding

439    0.20 of their respective categories in the test set. The bivalve *Pteria* was misidentified as *Bakevellia*

440    (0.25), and the brachiopod *Paryphella* was misidentified as *Fusichonetes* (0.24), with

441    morphological similarity being the main reason for misidentification. For example, the shells of

442    both *Pteria* and *Bakevellia* have similar outlines and are anteriorly oblique. The posterior ear is

443    larger than the anterior ear. Distinctive concentric rings are visible on the shell surface. All these

444    features are very similar.

445       Importantly, the vast majority of misidentifications in the hybrid auto-identification model

446   occurred within categories (i.e., bivalves were misidentified as other bivalves and brachiopods

447   were misidentified as other brachiopods), whereas misidentifications between broad categories

448   were relatively rare. For example, only 0.04 of the brachiopod *Araxathris* were misidentified as

449   bivalve *Daonella* and 0.04 as bivalve *Eumorphotis*, which indicates that bivalves and brachiopods

450   have considerable morphological differences.

451       The above are all cases where the input fossil taxon is included in the training set, but in

452   reality, there are many fossil taxa that are not included in the training set. To deal with this

453   exception, we propose an AM to identify such cases. The accuracy of AM (suitable for Order 22)

454   is 85.54%. When the training is completed, the user can use the AM to verify whether the taxon

455   of the input images is included in the training set and the usability of the genus/species

456   identification model. If the result is "applicable", the fossil will be identified automatically. If the

457   result is "inapplicable", the identification model will give the name of the fossil taxon that is most

458   similar to it, and the user can continue the manual identification based on that taxon.

459   *Morphological analysis of fossils*

460       Fossils have complex and variable high-dimensional morphological features, which are

461   difficult to visualize and analyze. Deep learning can extract features, downscale dimensions of

462   data, and exclude the influence of human bias to fully reflect the fossil features. Neural networks

463   can extract features more efficiently than manually selected features, although the majority of the

464   data extracted by models are too abstract for the human eye (Keceli *et al.* 2017). The accuracy of

465     supervised classification of ammonoids using human-selected geometric features was similar to

466     the accuracy in this study (Foxon 2021).

467       Machine learning can quantify morphological features and compare differences. In the feature

468     map (Fig. 8), we can observe the identification features used by the convolutional neural network.

469     However, the supervised deep learning used in this paper is a "result reason" approach that cannot

470     verify the correctness of the taxonomic practice. Models may use some features not used by experts

471     to identify, which does not mean that the taxonomic practice is wrong. A possible scenario is that

472     there are multiple differences between the two taxa, with experts and models choosing different

473     perspectives. The model establishes a relationship between the input (fossil image, i.e.,

474     morphological features) and the output (taxon), and its ability to accurately identify fossil taxa

475     indicates that taxonomic practice is well correlated with fossil morphology. However, the features

476     used by the models sometimes differ from those used by humans (Liu *et al.* 2022). Input-output

477     relationships are established by feature extraction through convolutional neural networks.

478     Automatic identification relies on these features that are similar with the working process of

479     experts. The features extracted by the model are diverse, such as the umbilicus, ribs, and inner

480     whorl of the ammonoid, spires and apices of gastropod, and growth lines and radial ribs of bivalve

481     and brachiopod (Liu *et al*. 2022). For the identification results, there is no difference between the

482     model's identification using images (actually fossil morphology) and the expert's identification

483     using characterization. This is essentially determined by the prior knowledge, which is obtained

484     by taxonomic practice. In the future, unsupervised learning may be able to provide unique insights

485     to evaluate taxonomic practice.

486     In the downscaled visualization of this model for the validation and test sets, the brachiopods

487     and bivalves are clearly demarcated, but a few points are still mixed (Fig. 9B). A clear boundary

488     means that the brachiopod and bivalve fossils are sufficiently morphologically distinct, so that the

489     model can extract the differences well and represent them quantitatively. This demonstrates the

490     unique potential of deep learning models for fossil feature extraction. Without inputting any prior

491     knowledge other than the genus name (e.g., the model does not know which genus belongs to

492     bivalve or brachiopod), the model computationally obtains information on the morphological

493     differences between bivalve and brachiopod, which is compatible with the expert's classification.

494     In addition to the distinction between bivalve and brachiopod, the t-SNE gives an indication of the

495     similarity of fossil morphology. For example, in Fig. 9B numbers 8(*Pteria*) and 17(*Bakevellia*)

496     overlap more, which demonstrates their more similar morphology. This is the same as the

497     traditional morphological view. In the future, it may be possible to use this feature to find similar

498     classification boundaries relying on models to perceive more detailed information about fossils

499     (e.g., ornamental features and 3D-morphology), which in turn could allow for quantitative

500     differentiation of gradual features (Klinkenbuß *et al.* 2020; Edie *et al.* 2023). That could not only

501     provide new possible perspectives for exploring fossil classification and biomorphological

502     evolution, but also try to explore whether there are important features that have been overlooked

503     by experts. In terms of the distribution area, the distribution of bivalve points is more extensive

504     than that of brachiopods, indicating that bivalves have greater morphological variability than

505     brachiopods in our dataset (but the effect of image context is not excluded here). Overall, the fossil

506     features extracted by CNNs can reflect the morphological characteristics of organisms to some

507    extent.

508        CNNs can complement existing methods for morphological studies such as morphological

509    matrix (Dai *et al.* 2021), landmark (Bazzi *et al.* 2018), fractal dimensions (Wiese *et al.* 2022),

510    ornamentation index (Miao *et al.* 2022), conch properties (De Baets 2021), and 3D morphological

511    methods (Klinkenbuß *et al.* 2020) and provide new perspectives for studying the morphological

512    evolution of fossils in the future. Geometric morphometry requires the extraction of fossil features

513    by labelling manually and performing descending operations (e.g., principal component analysis),

514    which has proven to be very effective (Aguirre *et al.* 2016; Topper *et al.* 2017;). In this method,

515    fossil features are selected by experts, with biological significance and better interpretation.

516    However, it is also influenced by human factors, and some features may be missed (Villier and

517    Korn 2004, Dai *et al.* 2021). Artificial intelligence differs in that it can obtain all information

518    displayed in fossil images (not just a few dozen points). These obtained features are then

519    downscaled (e.g., t-SNE used in this paper) to get the final fossil features. However, due to the

520    black-box character of deep learning, the features obtained are poorly interpretable, and whether

521    they are biologically meaningful needs further study in the future. Therefore, the advantage of

522    artificial intelligence mainly lies in the feature extraction, which reduces the subjective influence

523    and the time cost of manual marking. On the other hand, manual feature extraction is difficult to

524    orient to a large number of specimens and is based only on some specific species. However, deep

525    learning is capable of obtaining information from more specimens at the scale of big data, such as

526    intraspecific differences, spatial and temporal differences, etc., due to its ability to automate the

527    extraction of fossil features. Moreover, combining 3D information of fossils for palaeontological

528    studies is also promising (Hou *et al.* 2020).

## Conclusions

530    In this study, we used machine learning to automate fossil identification based on the practical

531    needs of palaeontological research. We built a bivalve and brachiopod fossil dataset by collecting

532    open literature, with > 16,000 "image-label" data pairs. Using these data, we compared the

533    performance of several convolutional neural network models based on VGG-16, Inception-

534    ResNet-v2, and EfficientNetV2s, which are commonly used in the field of image classification

535    and fossil identification. For this identification task, we found that EfficientNetV2s has the best

536    performance.

537    We finally achieved automatic fossil identification including 22 fossil genera (genus mode,

538    based on BBFID, including 13 bivalve genera and 9 brachiopod genera) and 16 fossil species

539    (species mode, based on BBFID, including 8 bivalve species and 8 brachiopod species), both with

540    > 80% accuracy. Furthermore, we conducted a study on the multiple categories' automatic fossil

541    identification at the species level, and the test accuracy was ~64% based on BBFID (scale C,

542    containing 343 bivalves and brachiopods). Models performed well in the automatic identification

543    of multiple categories with a small dataset. These models can be deployed to a web platform

544    [www.ai-fossil.com, (Liu *et al.* 2022)] in the future to make them accessible more easily and usable

545    by researchers. At present, automatic fossil identification must be based on expert consensus,

546    which is precisely why we emphasize the use of this model primarily for common fossil categories

547    to aid in identification. With more taxa be included, we can use the output from deep learning

548  models to accelerate the systematic palaeontology work during research rather than replace it and

549  contribute to quantitative assessment of morphology (De Baets 2021). Therefore, the researchers

550  can focus on most challenging and ambiguous identification cases. When a new taxon is found,

551  the AM output is "unapplicable", and experts can perform further taxonomic studies on it. When

552  experts decide to establish a new species, the fossil differences given by the algorithm can assist

553  them in making determinations, which is what the model excels at. But ultimately the

554  establishment of new species still depends on how taxonomists apply the results of deep learning.

555  We believe that there will be many palaeontologists working on fossil taxonomy and creating a

556  steady stream of a priori knowledge to promote the interdisciplinary relationship between

557  palaeontology and computer science together with AI researchers.

558      However, it must be noted that the model is an exploratory experiment and can currently only

559  serve as a useful assist to manual identification, not a complete replacement for it, at least for now.

560  The current model still relies on a manually created taxonomy and uses it as a priori knowledge

561  for model training. Current models are not able to combine all biological features (now only use

562  morphological data) to build the taxonomy by themselves. However, when experts have completed

563  the taxonomic criteria, researchers can use AI to identify fossils based on those criteria, reducing

564  repetitive identification work and allowing palaeontologists to have more time and energy for other

565  more creative research work.

566      We also used machine learning to extract high-dimensional data of fossil morphology and

567  downscaled them to obtain fossil morphological feature distribution maps, which present the

568  similarity of fossil morphology in a visual way. It was found that the bivalve and brachiopod

569    distribution regions have distinctive boundaries, and the morphological differences between the

570    two are obvious enough from the neural network perspective. In this process, models based on

571    deep learning are not absolutely objective. In contrast, palaeontologists play a crucial role. This is

572    precisely why we chose researcher consensus as a priori knowledge. Furthermore, we downscaled

573    the fossil features to cast the map and observe their morphological distribution. Compared with

574    the manually selected features, features based on the models are more objective and can better

575    reflect the morphological characteristics of fossils, which are still derived based on the consensus

576    of researchers on fossil taxonomy to a certain extent. In the future, this can be used as a basis to

577    quantify morphological information, analyze their morphological spatial distribution, and provide

578    a new perspective for exploring biological evolution.

579

580    **Data Availability Statement**

581    BBFID is available from the Zenodo digital repository: https://doi.org/10.5281/zenodo.7248779.

582    The main code and models of this study can be found at: https://doi.org/10.5281/zenodo.8126697.

583

584    **References**

585    Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G. and Isard, M.
586        2016. TensorFlow: a system for Large-Scale machine learning. 265–283. *12th USENIX symposium on*
587        *operating systems design and implementation (OSDI 16)*.
588    Aguirre, M. L., Richiano, S., Alvarez, A. and Farinati, E. A. 2016. Reading shell shape: implications for
589        palaeoenvironmental reconstructions. A case study for bivalves from the marine Quaternary of Argentina
590        (south-western Atlantic). *Historical Biology*, **28**, 753-773.
591    Al-Sabouni, N., Fenton, I. S., Telford, R. J. and Kučera, M. 2018. Reproducibility of species recognition in modern
592        planktonic foraminifera and its implications for analyses of community structure. *Journal of*

593      *Micropalaeontology*, **37**, 519-534.

594    Alroy, J., Aberhan, M., Bottjer, D. J., Foote, M., Fursich, F. T., Harries, P. J., Hendy, A. J. W., Holland, S. M., Ivany,

595      L. C., Kiessling, W., Kosnik, M. A., Marshall, C. R., Mcgowan, A. J., Miller, A. I., Olszewski, T. D.,

596      Patzkowsky, M. E., Peters, S. E., Villier, L., Wagner, P. J., Bonuso, N., Borkow, P. S., Brenneis, B., Clapham,

597      M. E., Fall, L. M., Ferguson, C. A., Hanson, V. L., Krug, A. Z., Layou, K. M., Leckey, E. H., Nurnberg, S.,

598      Powers, C. M., Sessa, J. A., Simpson, C., Tomasovych, A. and Visaggi, C. C. 2008. Phanerozoic trends in

599      the global diversity of marine invertebrates. *Science*, **321**, 97–100.

600    Austen, G. E., Bindemann, M., Griffiths, R. A. and Roberts, D. L. 2016. Species identification by experts and non-

601      experts: comparing images from field guides. *Scientific Reports*, **6**, 1-7.

602    Baets, K. D. 2021. Performance of machine-learning approaches in identifying ammonoid species based on conch

603      properties. *Peer Community in Paleontology*, **1**, 100010.

604    Ballanti, L. A., Tullis, A. and Ward, P. D. 2012. Comparison of oxygen consumption by Terebratalia transversa

605      (Brachiopoda) and two species of pteriomorph bivalve molluscs: implications for surviving mass extinctions.

606      *Paleobiology*, **38**, 525–537.

607    Bazzi, M., Kear, B. P., Blom, H., Ahlberg, P. E. and Campione, N. E. 2018. Static dental disparity and morphological

608      turnover in sharks across the end-Cretaceous mass extinction. *Current Biology*, **28**, 2607–2615.

609    Benton, M. J. and Harper, D. A. 2020. *Introduction to paleobiology and the fossil record*. John Wiley & Sons, 642

610      pp.

611    Botev, Z. I., Kroese, D. P., Rubinstein, R. Y. and L'Ecuyer, P. 2013. The cross-entropy method for optimization.

612      *Handbook of statistics*, **31**, 35–59.

613    Bourel, B., Marchant, R., De Garidel-Thoron, T., Tetard, M., Barboni, D., Gally, Y. and Beaufort, L. 2020. Automated

614      recognition by multiple convolutional neural networks of modern, fossil, intact and damaged pollen grains.

615      *Computers & Geosciences*, **140**, 104498.

616    Brodzicki, A., Piekarski, M., Kucharski, D., Jaworek-Korjakowska, J. and Gorgon, M. 2020. Transfer Learning

617      Methods as a New Approach in Computer Vision Tasks with Small Datasets. *Foundations of Computing and*

618      *Decision Sciences*, **45**, 179–193.

619    Byeon, W., Dominguez-Rodrigo, M., Arampatzis, G., Baquedano, E., Yravedra, J., Mate-Gonzalez, M. A. and

620      Koumoutsakos, P. 2019. Automated identification and deep classification of cut marks on bones and its

621      paleoanthropological implications. *Journal of Computational Science*, **32**, 36–43.

622    Clapham, M. E., Bottjer, D. J., Powers, C. M., Bonuso, N., Fraiser, M. L., Marenco, P. J., Dornbos, S. Q. and Pruss,

623      S. B. 2006. Assessing the ecological dominance of Phanerozoic marine invertebrates. *Palaios*, **21**, 431–441.

624    Culverhouse, P. F., Williams, R., Reguera, B., Herry, V. and González-Gil, S. 2003. Do experts make mistakes? A

625      comparison of human and machine identification of dinoflagellates. *Marine Ecology Progress Series*, **247**,

626      17-25.

627    Dai, X., Davies, J. H. F. L., Yuan, Z., Brayard, A., Ovtcharova, M., Xu, G., Liu, X., Smith, C. P. A., Schweitzer, C.

628      E., Li, M., Perrot, M. G., Jiang, S., Miao, L., Cao, Y., Yan, J., Bai, R., Wang, F., Guo, W., Song, H., Tian,

629      L., Dal Corso, J., Liu, Y., Chu, D. and Song, H. 2023. A mesozoic fossil lagerstätte from 250.8 million years

630      ago shows a modern-type marine ecosystem. *Science*, **379**, 567-572.Dai, X., Korn, D. and Song, H. 2021.

631      Morphological selectivity of the Permian-Triassic ammonoid mass extinction. *Geology*, **49**, 1112–1116.

632    Deng, J., Dong, W., Socher, R., Li, L. J., Li, K. and Li, F. F. 2009. ImageNet: A Large-Scale Hierarchical Image

633      Database. 248–255. *IEEE-Computer-Society Conference on Computer Vision and Pattern Recognition*

634        *Workshops*. IEEE, Miami Beach, FL.

635 Dionisio, A., Solano, G., Quisote, M. and Marquez, E. 2020. A Radiolarian Classifier using Convolutional Neural
636         Networks. 1–5. *International Conference on Artificial Intelligence and Signal Processing (AISP)*. VIT AP
637         Univ, Amaravati, India.

638 Edie, S. M., Collins, K. S. and Jablonski, D. 2023. High-throughput micro-ct scanning and deep learning segmentation
639         workflow for analyses of shelly invertebrates and their fossils: Examples from marine bivalvia. *Frontiers in
640         Ecology and Evolution*, **11**, 1127756.

641 Fan, J. X., Shen, S. Z., Erwin, D. H., Sadler, P. M., Macleod, N., Cheng, Q. M., Hou, X. D., Yang, J., Wang, X. D.,
642         Wang, Y., Zhang, H., Chen, X., Li, G. X., Zhang, Y. C., Shi, Y. K., Yuan, D. X., Chen, Q., Zhang, L. N., Li,
643         C. and Zhao, Y. Y. 2020. A high-resolution summary of Cambrian to Early Triassic marine invertebrate
644         biodiversity. *Science*, **367**, 272–277.

645 Fan, L., Xu, C., Jarzembowski, E. A. and Cui, X. 2022. Quantifying plant mimesis in fossil insects using deep learning.
646         *Historical Biology*, **34**, 907-916.

647 FlÜgel, E. and Munnecke, A. 2010. *Microfacies of carbonate rocks: analysis, interpretation and application*.
648         Springer, Berlin, 924 pp.

649 Foxon, F. 2021. Ammonoid taxonomy with supervised and unsupervised machine learning algorithms. *PaleorXiv
650         ewkx9, ver. 3*.https://doi.org/10.31233/osf.io/ewkx9.

651 Fraiser, M. L. and Bottjer, D. J. 2007. When bivalves took over the world. *Paleobiology*, **33**, 397-413.

652 Giusti, A., Ciresan, D. C., Masci, J., Gambardella, L. M. and Schmidhuber, J. 2013. Fast Image Scanning with Deep
653         Max-Pooling Convolutional Neural Networks. 4034–4038. *20th IEEE International Conference on Image
654         Processing (ICIP)*. IEEE, Melbourne, Australia.

655 Gould, S. J. and Calloway, C. B. 1980. Clams and brachiopods—ships that pass in the night. *Paleobiology*, **6**, 383-
656         396.

657 Gradstein, F. M., Ogg, J. G., Schmitz, M. and Ogg, G. 2012. *The geologic time scale 2012*. Elsevier, 1144 pp.

658 He, K. M., Zhang, X. Y., Ren, S. Q. and Sun, J. 2016. Deep Residual Learning for Image Recognition. 770–778. *2016
659         IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Seattle, WA.

660 Ho, M., Idgunji, S., Payne, J. L. and Koeshidayatullah, A. 2023. Hierarchical multi-label taxonomic classification of
661         carbonate skeletal grains with deep learning. *Sedimentary Geology*, **443**, 106298.

662 Hou, Y. M., Cui, X. D., Canul-Ku, M., Jin, S. C., Hasimoto-Beltran, R., Guo, Q. H. and Zhu, M. 2020. ADMorph: A
663         3D Digital Microfossil Morphology Dataset for Deep Learning. *IEEE Access*, **8**, 148744–148756.

664 Hou, C., Lin, X., Huang, H., Xu, S., Fan, J., Shi, Y. and Lv, H. 2023. Fossil image identification using deep learning
665         ensembles     of     data     augmented     multiviews.     *arXiv*     *preprint
666         arXiv*.https://doi.org/10.48550/arXiv.2302.080622302.08062.

667 Hsiang, A. Y., Brombacher, A., Rillo, M. C., Mleneck-Vautravers, M. J., Conn, S., Lordsmith, S., Jentzen, A.,
668         Henehan, M. J., Metcalfe, B., Fenton, I. S., Wade, B. S., Fox, L., Meilland, J., Davis, C. V., Baranowskils,
669         U., Groeneveld, J., Edgar, K. M., Movellan, A., Aze, T., Dowsett, H. J., Miller, C. G., Rios, N. and Hull, P.
670         M. 2019. Endless Forams: > 34,000 Modern Planktonic Foraminiferal Images for Taxonomic Training and
671         Automated Species Recognition Using Convolutional Neural Networks. *Paleoceanography and
672         Paleoclimatology*, **34**, 1157–1177.

673 Huang, Y., Tong, J. and Fraiser, M. L. 2018. A Griesbachian (Early Triassic) mollusc fauna from the Sidazhai Section,
674         Southwest China, with paleoecological insights on the proliferation of genus *Claraia* (bivalvia). *Journal of

675        *Earth Science*, **29**, 794-805.

676    Ioffe, S. and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal

677        covariate shift. 448–456. *International conference on machine learning*. PMLR,

678    Keceli, A. S., Kaya, A. and Keceli, S. U. 2017. Classification of radiolarian images with hand-crafted and deep

679        features. *Computers & Geosciences*, **109**, 67–74.

680    Kiel, S. 2021. Assessing bivalve phylogeny using Deep Learning and computer vision approaches.

681        *bioRxiv*.https://doi.org/10.1101/2021.04.08.438943.

682    Kingma, D. P. and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint*

683        *arXiv:1412.6980*.https://doi.org/10.48550/arXiv.1412.6980.

684    Klinkenbuß, D., Metz, O., Reichert, J., Hauffe, T., Neubauer, T. A., Wesselingh, F. P. and Wilke, T. 2020.

685        Performance of 3D morphological methods in the machine learning assisted classification of closely related

686        fossil bivalve species of the genus dreissena. *Malacologia*, **63**, 95-105.

687    Koeshidayatullah, A., Morsilli, M., Lehrmann, D. J., Al-Ramadan, K. and Payne, J. L. 2020. Fully automated

688        carbonate petrography using deep convolutional neural networks. *Marine and Petroleum Geology*, **122**,

689        104687.

690    Kong, S., Punyasena, S. and Fowlkes, C. 2016. Spatially Aware Dictionary Learning and Coding for Fossil Pollen

691        Identification. 1305–1314. *29th IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Las

692        Vegas, NV.

693    Konopleva, E. S., Bolotov, I. N., Vikhrev, I. V., Gofarov, M. Y. and Kondakov, A. V. 2017. An integrative approach

694        underscores the taxonomic status of Lamellidens exolescens, a freshwater mussel from the Oriental tropics

695        (Bivalvia: Unionidae). *Systematics and Biodiversity*, **15**, 204–217.

696    Lallensack, J. N., Romilio, A. and Falkingham, P. L. 2022. A machine learning approach for the discrimination of

697        theropod and ornithischian dinosaur tracks. *Journal of The Royal Society Interface*, **19**, 20220588.

698    Lecun, Y., Bottou, L., Bengio, Y. and Haffner, P. 1998. Gradient-based learning applied to document recognition.

699        *Proceedings of the IEEE*, **86**, 2278–2324.

700    Liow, L. H., Reitan, T. and Harnik, P. G. 2015. Ecological interactions on macroevolutionary time scales: Clams and

701        brachiopods are more than ships that pass in the night. *Ecology Letters*, **18**, 1030-1039.

702    Liu, X., Jiang, S., Wu, R., Shu, W., Hou, J., Sun, Y., Sun, J., Chu, D., Wu, Y. and Song, H. 2022. Automatic taxonomic

703        identification based on the Fossil Image Dataset (> 415,000 images) and deep convolutional neural networks.

704        *Paleobiology*.1–22.

705    Liu, X. K. and Song, H. J. 2020. Automatic identification of fossils and abiotic grains during carbonate microfacies

706        analysis using deep convolutional neural networks. *Sedimentary Geology*, **410**, 105790.

707    Miao, L. Y., Dai, X., Song, H. C., Backes, A. R. and Song, H. J. 2022. A new index for quantifying the ornamentational

708        complexity of animals with shells. *Ecology and Evolution*, **12**, e9247.

709    Nair, V. and Hinton, G. E. 2010. Rectified linear units improve restricted boltzmann machines. *Icml*.

710    Nevesskaja, L. A. 2003. Morphogenesis and ecogenesis of bivalves in the Phanerozoic. *Paleontological Journal*, **37**,

711        S591-S715.

712    Niu, Z. B. and Xu, H.-H. 2022. AI-based graptolite identification improves shale gas exploration.

713        *bioRxiv*.https://doi.org/10.1101/2022.01.17.476477.

714    Payne, J. L., Heim, N. A., Knope, M. L. and Mcclain, C. R. 2014. Metabolic dominance of bivalves predates

715        brachiopod diversity decline by more than 150 million years. *Proceedings of the Royal Society B-Biological*

716         *Sciences*, **281**, 20133122.

717    Piazza, V., Ullmann, C. V. and Aberhan, M. 2020. Temperature-related body size change of marine benthic
718         macroinvertebrates across the early toarcian anoxic event. *Scientific Reports*, **10**, 4675.

719    Pires De Lima, R., Welch, K. F., Barrick, J. E., Marfurt, K. J., Burkhalter, R., Cassel, M. and Soreghan, G. S. 2020.
720         Convolutional Neural Networks as an aid to 131 biostratigraphy and micropaleontology: a test on late
721         paleozoic microfossils. *Palaios*, **35**, 391–402.

722    Pitrat, C. W. and Moore, R. C. 1965. *Treatise on invertebrate paleontology part h, brachiopoda*. Ed. RC Moore. Univ.
723         of Kansas Press and Geol. Soc. America, 522 pp.

724    Punyasena, S. W., Tcheng, D. K., Wesseln, C. and Mueller, P. G. 2012. Classifying black and white spruce pollen
725         using layered machine learning. *New Phytologist*, **196**, 937-944.

726    Romero, I. C., Kong, S., Fowlkes, C. C., Jaramillo, C., Urban, M. A., Oboh-Ikuenobe, F., D'apolito, C. and Punyasena,
727         S. W. 2020. Improving the taxonomy of fossil pollen using convolutional neural networks and
728         superresolution microscopy. *Proceedings of the National Academy of Sciences of the United States of*
729         *America*, **117**, 28496–28505.

730    Rhodes, M. C. and Thompson, R. J. 1993. Comparative physiology of suspension-feeding in living brachiopods and
731         bivalves: Evolutionary implications. *Paleobiology*, **19**, 322-334.

732    Sarkar, D., Bali, R. and Ghosh, T. 2018. *Hands-on transfer learning with python: Implement advanced deep learning*
733         *and neural network models using tensorflow and keras*. Packt Publishing Ltd, 430 pp.

734    Scotese, C. R., Song, H. J., Mills, B. J. W. and Van Der Meer, D. G. 2021. Phanerozoic paleotemperatures: The earth's
735         changing climate during the last 540 million years. *Earth-Science Reviews*, **215**, 103503.

736    Sepkoski, J. J. 1981. A factor analytic description of the Phanerozoic marine fossil record. *Paleobiology*, **7**, 36–53.

737    Simonyan, K. and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv*
738         *preprint arXiv:1409.1556*.https://doi.org/10.48550/arXiv.1409.1556.

739    Solano, G. A., Gasmen, P. and Marquez, E. J. 2018. Radiolarian classification decision support using supervised and
740         unsupervised learning approaches. 1-6. *2018 9th International Conference on Information, Intelligence,*
741         *Systems and Applications (IISA)*.

742    Song, H. J., Kemp, D. B., Tian, L., Chu, D. L., Song, H. Y. and Dai, X. 2021. Thresholds of temperature change for
743         mass extinctions. *Nature Communications*, **12**, 4694.

744    Su, T., Farnsworth, A., Spicer, R. A., Huang, J., Wus, F. X., Liu, J., Li, S. F., Xing, Y. W., Huang, Y. J., Deng, W. Y.
745         D., Tang, H., Xu, C. L., Zhao, F., Srivastava, G., Valdes, P. J., Deng, T. and Zhou, Z. K. 2019. No high
746         Tibetan Plateau until the Neogene. *Science Advances*, **5**, eaav2189.

747    Sulser, H., García-Ramos, D., Kürsteiner, P. and Menkveld-Gfeller, U. 2010. Taxonomy and palaeoecology of
748         brachiopods from the South-Helvetic zone of the Fäneren region (Lutetian, Eocene, NE Switzerland). *Swiss*
749         *Journal of Geosciences*, **103**, 257–272.

750    Szegedy, C., Ioffe, S., Vanhoucke, V. and Alemi, A. A. 2017. Inception-v4, inception-resnet and the impact of residual
751         connections on learning. *Thirty-first AAAI conference on artificial intelligence*.

752    Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V. and Rabinovich, A.
753         2015. Going deeper with convolutions. 1–9. *Proceedings of the IEEE conference on computer vision and*
754         *pattern recognition*.

755    Tan, C. Q., Sun, F. C., Kong, T., Zhang, W. C., Yang, C. and Liu, C. F. 2018. A Survey on Deep Transfer Learning.
756         270–279. *27th International Conference on Artificial Neural Networks (ICANN)*. Springer International

757        Publishing Ag, Rhodes, Greece.

758 Tan, M. and Le, Q. 2021. Efficientnetv2: Smaller models and faster training. 10096–10106. *International Conference*
759        *on Machine Learning*. PMLR.

760 Thayer, C. W. 1986. Are brachiopods better than bivalves? Mechanisms of turbidity tolerance and their interaction
761        with feeding in articulates. *Paleobiology*, **12**, 161-174.

762 Topper, T. P., Strotz, L. C., Skovsted, C. B. and Holmer, L. E. 2017. Do brachiopods show substrate‑related
763        phenotypic variation? A case study from the Burgess Shale. *Palaeontology*, **60**, 269-279.

764 Van Der Maaten, L. and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of Machine Learning Research*, **9**,
765        2579–2605.

766 Villier, L. and Korn, D. 2004. Morphological disparity of ammonoids and the mark of Permian mass extinctions.
767        *Science*, **306**, 264-266.

768 Wang, F., Chen, J., Dai, X. and Song, H. 2017. A new Dienerian (Early Triassic) brachiopod fauna from South China
769        and implications for biotic recovery after the Permian–Triassic extinction. *Papers in Palaeontology*, **3**, 425-
770        439.

771 Wang, H. Z., Li, C. F., Zhang, Z. F., Kershaw, S., Holmer, L. E., Zhang, Y., Wei, K. Y. and Liu, P. 2022. Fossil
772        brachiopod identification using a new deep convolutional neural network. *Gondwana Research*, **105**, 290–
773        298.

774 Wang, B., Sun, R., Yang, X., Niu, B., Zhang, T., Zhao, Y., Zhang, Y., Zhang, Y. and Han, J. 2023. Recognition of
775        rare microfossils using transfer learning and deep residual networks. *Biology*, **12**, 16.

776 Wiese, R., Harrington, K., Hartmann, K., Hethke, M., Von Rintelen, T., Zhang, H., Zhang, L.-J. and Riedel, F. 2022.
777        Can fractal dimensions objectivize gastropod shell morphometrics? A case study from lake lugu (sw China).
778        *Ecology and Evolution*, **12**, e8622.

779 Yin, H. F., Zhang, K. X., Tong, J. N., Yang, Z. Y. and Wu, S. B. 2001. The Global Stratotype Section and Point
780        (GSSP) of the Permian-Triassic Boundary. *Episodes*, **24**, 102–114.

781 Ying, X. 2019. An overview of overfitting and its solutions. *Journal of Physics: Conference Series*, **1168**, 022022.

782 Zhang, T., Wang, B., Li, D., Niu, B., Sun, J., Sun, Y., Yang, X., Luo, J. and Han, J. 2020. Artificial intelligence
783        identification of multiple microfossils from the Cambrian kuanchuanpu formation in southern shaanxi, China.
784        *Acta Geologica Sinica - English Edition*, **94**, 189-197.

785

786 **Figures and Tables:**

787

788 **FIG. 1.** Number of taxa and accuracy for automatic fossil identification studies based on deep

789 learning (Punyasena et al. 2012; Kong et al. 2016; Solano et al. 2018; Dionisio et al. 2020; Hou et

790 al. 2020; Liu and Song 2020; Pires De Lima et al. 2020; Zhang et al. 2020; Foxon 2021; Lallensack

791 et al. 2022; Niu and Xu 2022; Wang et al. 2022; Ho et al. 2023; Hou et al. 2023; Liu et al. 2023;

792    Wang et al. 2023)**.**

793    **FIG. 2.** Number of samples for each taxon at the genus level in (A) BBFID-1 and (B) BBFID-2

794    (scale B) and the distribution in subsets.

795

796    **FIG. 3.** DCNN architectures used in this study. Automatic identification model architectures of

797    A, B, C are modified from VGG-16 (Simonyan and Zisserman 2014), Inception-ResNet-v2

798    (Szegedy *et al.* 2017), and EfficientNetV2s (Tan and Le 2021) respectively.

799

800    **FIG. 4.** Confusion matrix and evaluation metrics of models trained by BBFID-1 (scale A) on genus

801    mode. The horizontal axis is the predicted label, and the vertical axis is the true label. Colors and

802    values represent the proportion of the corresponding taxon identified as the predicted label taxon.

803

804    **FIG. 5.** Confusion matrix and evaluation metrics of models trained by BBFID-2 (scale A) on genus

805    mode. Colors and values represent the proportion of the corresponding taxon identified as the

806    predicted label taxon.

807

808    **FIG. 6.** Confusion matrix and evaluation metrics of models trained by BBFID (scale A) on genus

809    mode. Colors and values represent the proportion of the corresponding taxon identified as the

810    predicted label taxon. The categories marked in red are brachiopods, and the others are bivalves.

811

812    **FIG. 7**. The training process of ATIM on genus mode using BBFID (scale A) (Order 22).

813

814 **FIG. 8.** Feature maps of the bivalve (*Claraia*) and brachiopod (*Lichuanorelloides*) fossils in

815 BBFID, plotted by extracting model (Order 22) intermediate output. Fossil images are from

816 Huang *et al.* (2018), and Wang *et al.* (2017).

817

818 **FIG. 9.** Fossil morphological feature distribution maps. (A) Training set data and (B) validation

819 set and test set data were fitted simultaneously using t-SNE. The accuracy of the original

820 identification model is 81.01%. The horizontal and vertical coordinates in the figure are the two

821 dimensions obtained by t-SNE (n_components=2, perplexity=10, init='pca', learning_rate=1,

822 n_iter= 6000, n_iter_without_progress=6000). The numbers represent different genera, where the

823 black numbers represent the bivalves and the red numbers represent the brachiopods. The detailed

824 correspondence is 0: *Pseudospiriferina*, 1: *Quemocuomegalodon*, 2: *Burmirhynchia*, 3:

825 *Promyalina*, 4: *Araxathyris*, 5: *Spiriferina*, 6: *Costatoria*, 7: *Fusichonetes*, 8: *Pteria*, 9: *Paryphella*,

826 10: *Neoschizodus*, 11: *Prelissorhynchia*, 12: *Juxathyris*, 13: *Piarorhynchella*, 14: *Leptochondria*,

827 15: *Daonella*, 16: *Unionites*, 17: *Bakevellia*, 18: *Halobia*, 19: *Eumorphotis*, 20: *Monotis*, 21:

828 *Claraia*.

829

830 **TABLE 1.** Identification accuracy training on BBFID-1 (scale A) at the genus level with different

831 model architectures and hyperparameters. Architectures in this table are shown in Fig. 3.

832 "Trainable layers of functional layers" represents the size of the parameters that can be trained.

833 "None" means that all layers of the backbone are frozen and the parameters involved in these

834    layers cannot be trained. These parameters maintain the values at the time of model initialization.

835    "Half layers" means that half of the backbone layer parameters are frozen, while "All layers"

836    means that all parameters of this model are not frozen and can be updated during the training

837    process. This setting has an impact on both the model training process and the model performance.
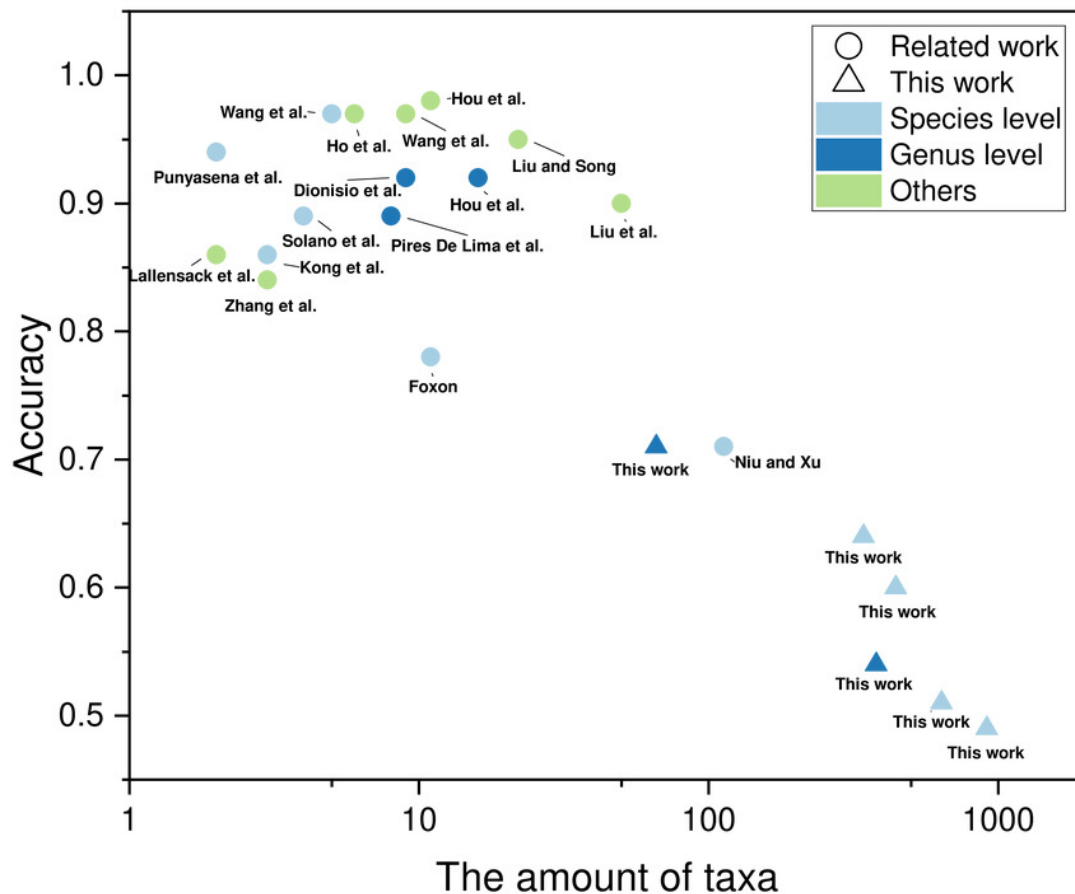
838

839    **TABLE 2.** Model performance using BBFID-1, BBFID-2 and BBFID in EfficientNetV2s

840    architecture. Learning rate starts from 1e-4 and the epoch is limited to less than 51. Test accuracy

841    / Test loss means the accuracy / loss of the saved model.

842

# Figure 1

Number of taxa and accuracy for automatic fossil identification studies based on deep learning.
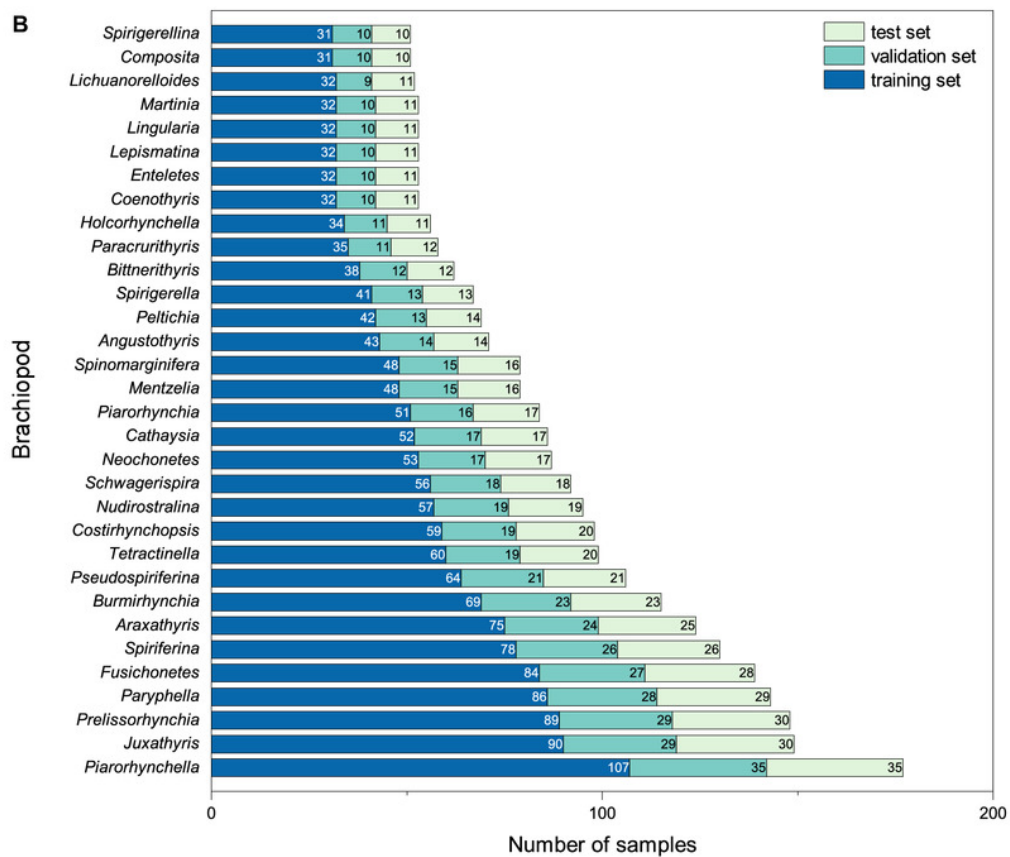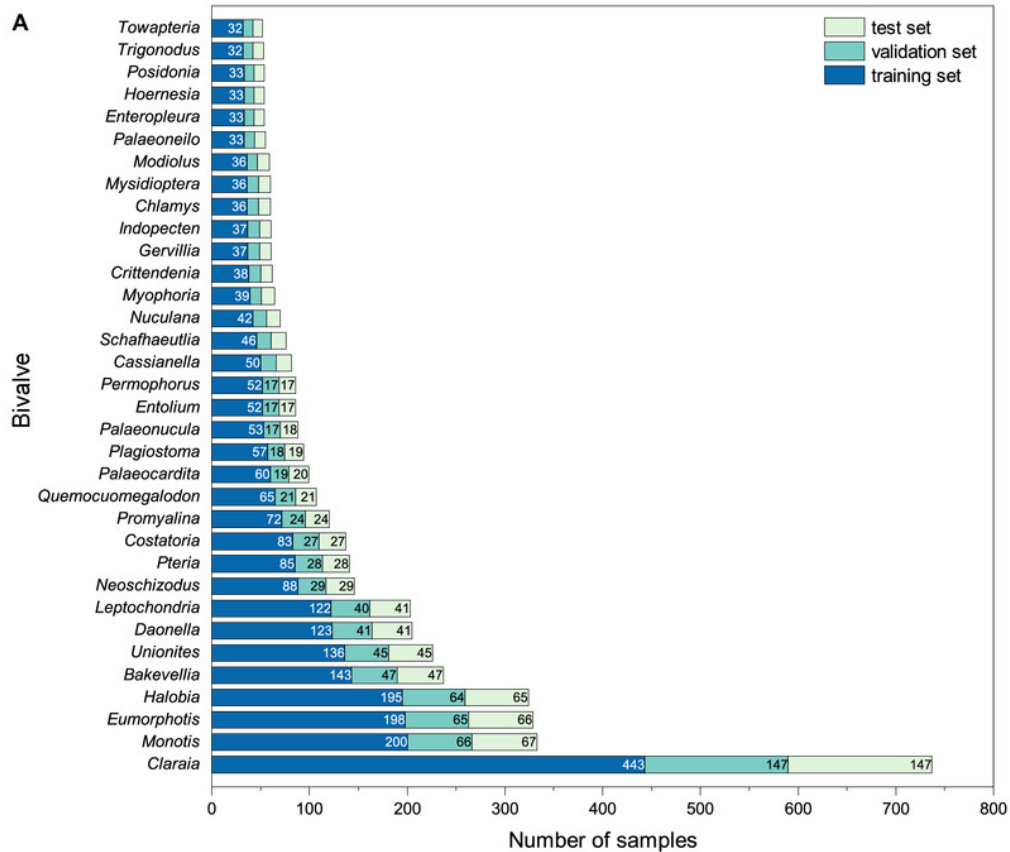
Number of taxa and accuracy for automatic fossil identification studies based on deep learning (Punyasena et al. 2012; Kong et al. 2016; Solano et al. 2018; Dionisio et al. 2020; Hou et al. 2020; Liu and Song 2020; Pires De Lima et al. 2020; Zhang et al. 2020; Foxon 2021; Lallensack et al. 2022; Niu and Xu 2022; Wang et al. 2022; Ho et al. 2023; Hou et al. 2023; Liu et al. 2023; Wang et al. 2023) **.**

# Figure 2

Number of samples for each taxon at the genus level in (A) BBFID-1 and (B) BBFID-2 (scale B) and the distribution in subsets.
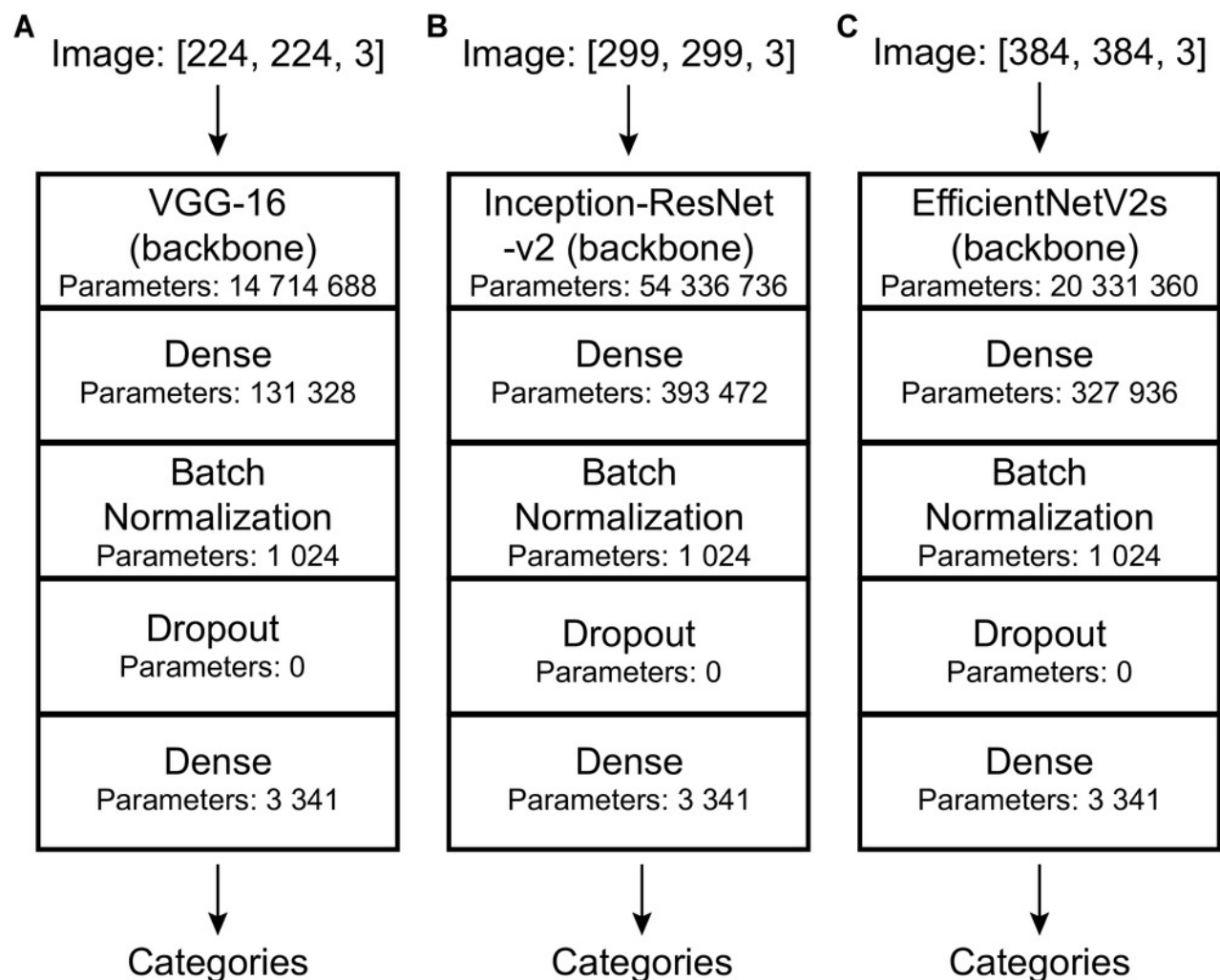
# Figure 3

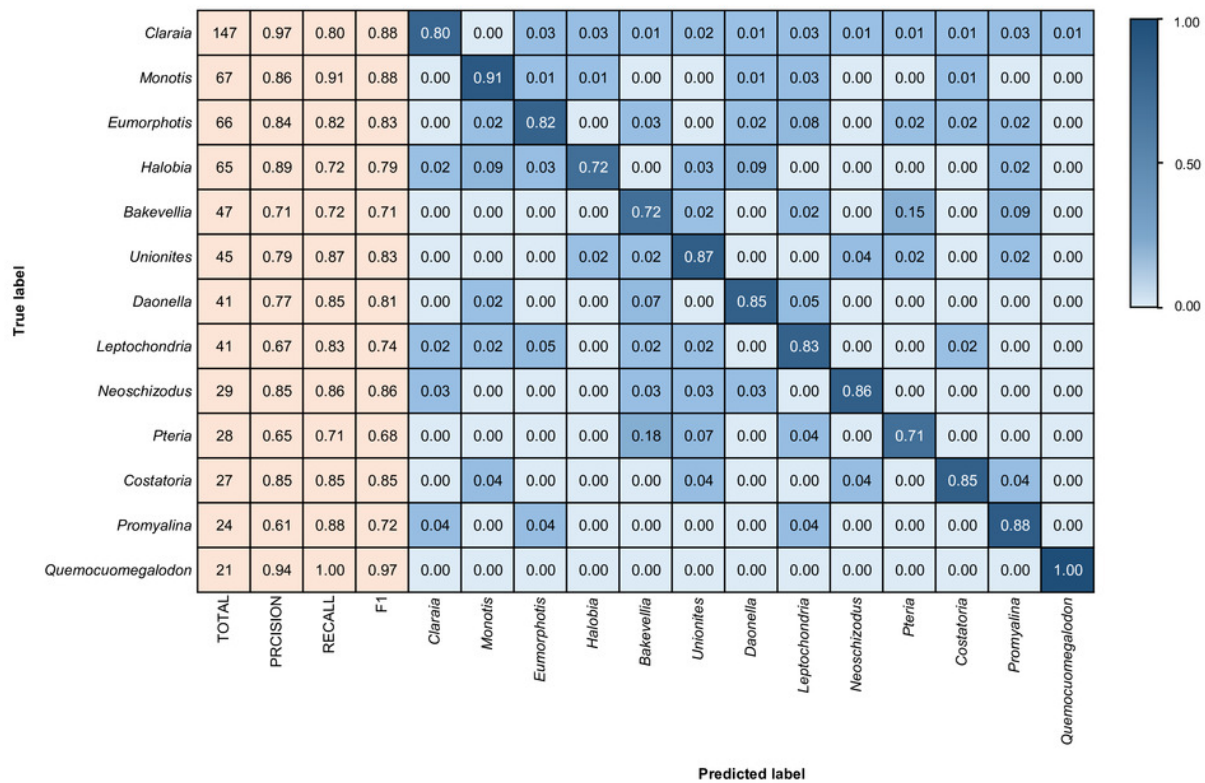DCNN architectures used in this study.

Automatic identification model architectures of A, B, C are modified from VGG-16 (Simonyan and Zisserman 2014), Inception-ResNet-v2 (Szegedy, et al. 2017), and EfficientNetV2s (Tan and Le 2021) respectively.

# Figure 4

Confusion matrix and evaluation metrics of models trained by BBFID-1 (scale A) on genus mode.

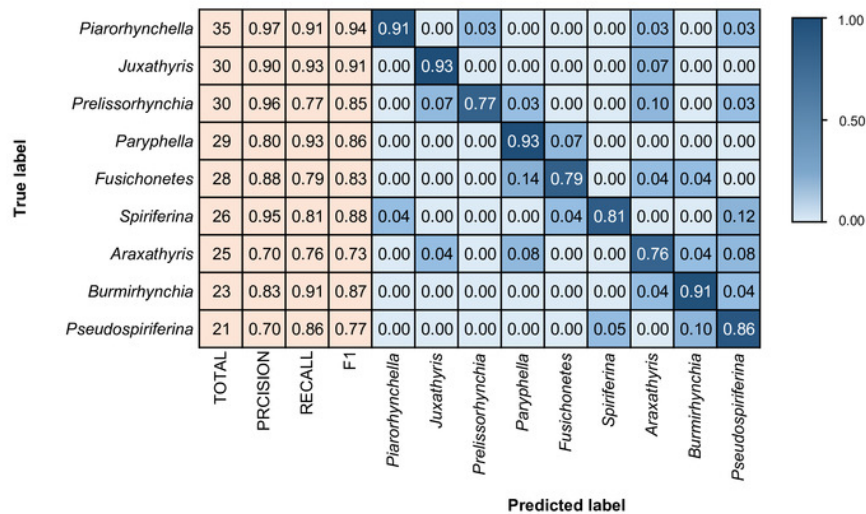The horizontal axis is the predicted label, and the vertical axis is the true label. Colors and values represent the proportion of the corresponding taxon identified as the predicted label taxon.

# Figure 5

Confusion matrix and evaluation metrics of models trained by BBFID-2 (scale A) on genus mode.

Colors and values represent the proportion of the corresponding taxon identified as the predicted label taxon.

# Figure 6

Confusion matrix and evaluation metrics of models trained by BBFID (scale A) on genus mode.

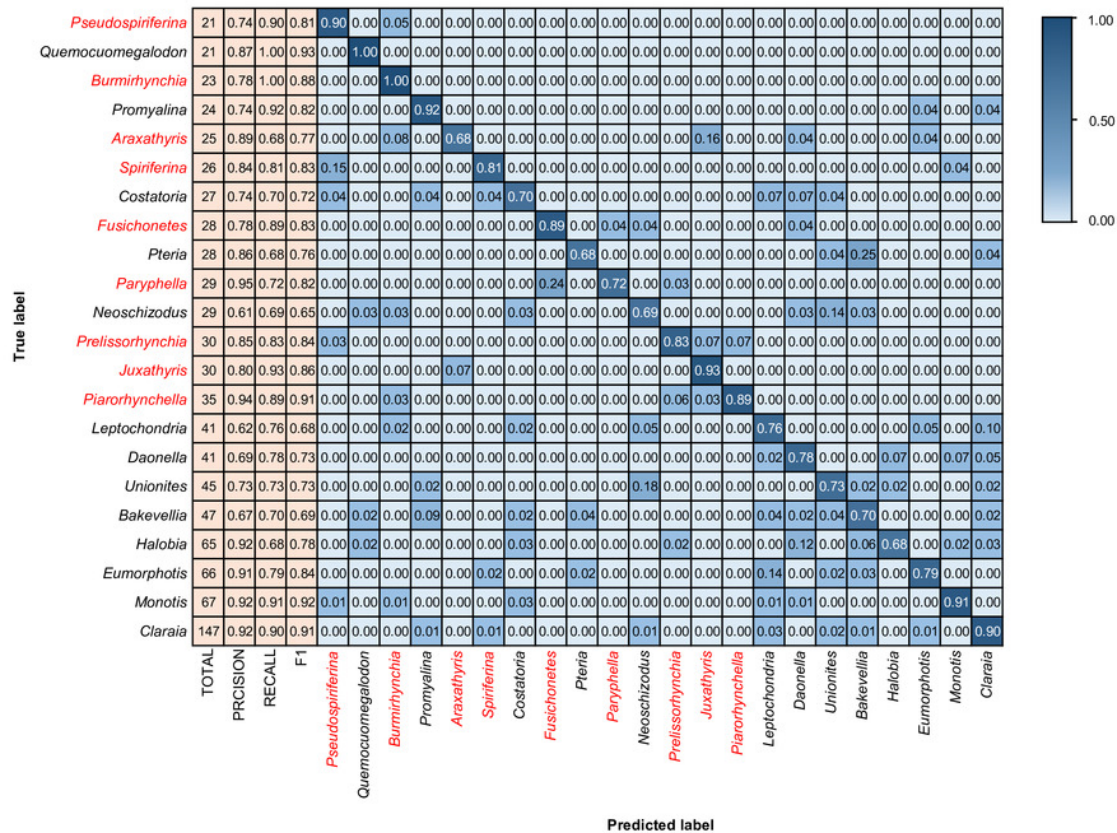Colors and values represent the proportion of the corresponding taxon identified as the predicted label taxon. The categories marked in red are brachiopods, and the others are bivalves.

# Figure 7

The training process of ATIM on genus mode using BBFID (scale A) (Order 22).

# Figure 8

Feature maps of the bivalve (*Claraia*) and brachiopod (*Lichuanorelloides*) fossils in BBFID, plotted by extracting model (Order 22) intermediate output.

Fossil images are from Huang *et al.* (2018), and Wang *et al.* (2017).

# Figure 9

Fossil morphological feature distribution maps.

(A) Training set data and (B) validation set and test set data were fitted simultaneously using t-SNE. The accuracy of the original identification model is 81.01%. The horiz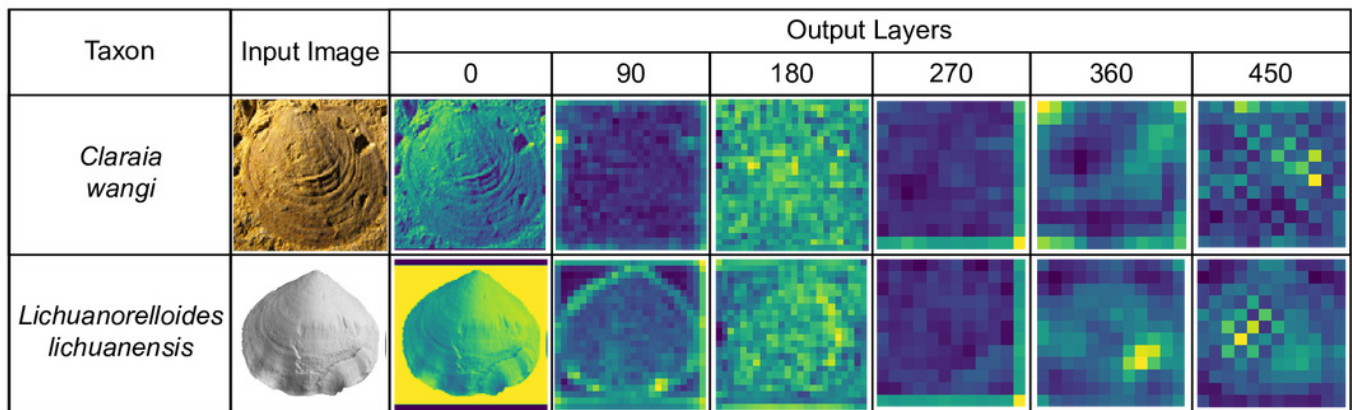ontal and vertical coordinates in the figure are the two dimensions obtained by t-SNE (n_components=2, perplexity=10, init='pca', learning_rate=1, n_iter= 6000, n_iter_without_progress=6000). The numbers represent different genera, where the black numbers represent the bivalves and the red numbers represent the brachiopods. The detailed correspondence is 0: *Pseudospiriferina*, 1: *Quemocuomegalodon*, 2: *Burmirhynchia*, 3: *Promyalina*, 4: *Araxathyris*, 5: *Spiriferina*, 6: *Costatoria*, 7: *Fusichonetes*, 8: *Pteria*, 9: *Paryphella*, 10: *Neoschizodus*, 11: *Prelissorhynchia*, 12: *Juxathyris*, 13: *Piarorhynchella*, 14: *Leptochondria*, 15: *Daonella*, 16: *Unionites*, 17: *Bakevellia*, 18: *Halobia*, 19: *Eumorphotis*, 20: *Monotis*, 21: *Claraia*.

**Table 1**(on next page)

Identification accuracy training on BBFID-1 (scale A) at the genus level with different model architectures and hyperparameters.

Architectures in this table are shown in Fig. 3. "Trainable layers of functional layers" represents the size of the parameters that can be trained. "None" means that all layers of the backbone are frozen and the parameters involved in these layers cannot be trained. These parameters maintain the values at the time of model initialization. "Half layers" means that half of the backbone layer parameters are frozen, while "All layers" means that all parameters of this model are not frozen and can be updated during the training process. This setting has an impact on both the model training process and the model performance.

1   **TABLE 1.** Identification accuracy training on BBFID-1 (scale A) at the genus level with different model architectures and hyperparameters. Architectures

2   in this table are shown in Fig. 3. "Trainable layers of functional layers" represents the size of the parameters that can be trained. "None" means that all

3   layers of the backbone are frozen and the parameters involved in these layers cannot be trained. These parameters maintain the values at the time of model

4   initialization. "Half layers" means that half of the backbone layer parameters are frozen, while "All layers" means that all parameters of this model are not

5   frozen and can be updated during the training process. This setting has an impact on both the model training process and the model performance.

| Order | Backbone | Batch size | Trainable layers of functional layers | Reduce LR on plateau | Epochs | Max. training accuracy | Min. training loss | Max. validation accuracy | Min. validation loss | Test accuracy | Test loss |
|-------|----------|-----------|----------------------------------------|---------------------|--------|------------------------|--------------------|--------------------------|----------------------|---------------|-----------|
| 1 | VGG-16 | 32 | None | Yes | 50 | 0.8648 | 0.4212 | 0.6444 | 1.1440 | 0.6281 | 1.2512 |
| 2 | VGG-16 | 32 | Half layers | Yes | 40 | 0.9959 | 0.0181 | 0.7515 | 0.9126 | 0.7330 | 0.8444 |
| 3 | VGG-16 | 32 | All layers | Yes | 50 | 0.7670 | 0.6080 | 0.5698 | 1.3465 | 0.5386 | 1.4802 |
| 4 | VGG-16 | 32 | All layers | No | 36 | 0.3609 | 1.8002 | 0.3338 | 2.0523 | 0.0957 | 3.0871 |
| 5 | Inception-ResNet-v2 | 8 | None | Yes | 50 | 0.3236 | 1.9945 | 0.3385 | 2.0345 | 0.3225 | 2.1000 |
| 6 | Inception-ResNet-v2 | 8 | Half layers | Yes | 50 | 0.7363 | 0.7163 | 0.5263 | 1.4931 | 0.4877 | 1.5584 |
| 7 | Inception-ResNet-v2 | 8 | All layers | Yes | 46 | 0.9959 | 0.0216 | 0.7934 | 1.2041 | 0.7778 | 2.5044 |
| 8 | Inception-ResNet-v2 | 8 | All layers | No | 46 | 0.9805 | 0.0602 | 0.7981 | 0.8178 | 0.6590 | 1.2590 |
| 9 | EfficientNetV2s | 8 | None | Yes | 50 | 0.5693 | 1.2799 | 0.5419 | 1.4210 | 0.4923 | 1.5424 |
| 10 | EfficientNetV2s | 8 | Half layers | Yes | 50 | 0.9708 | 0.1013 | 0.7624 | 0.8314 | 0.7515 | 0.8633 |
| 11 | EfficientNetV2s | 8 | All layers | Yes | 44 | 0.9959 | 0.0139 | 0.8338 | 0.6130 | 0.8302 | 0.6807 |
| 12 | EfficientNetV2s | 8 | All layers | No | 37 | 0.9825 | 0.0578 | 0.8136 | 0.7905 | 0.7886 | 0.8122 |

6

**Table 2**(on next page)

Model performance using BBFID-1, BBFID-2 and BBFID in EfficientNetV2s architecture.

Learning rate starts from 1e-4 and the epoch is limited to less than 51. Test accuracy / Test loss means the accuracy / loss of the saved model.

1  **TABLE 2.** Model performance using BBFID-1, BBFID-2 and BBFID in EfficientNetV2s architecture. Learning rate starts from 1e-4 and the epoch is
2  limited to less than 51. Test accuracy / Test loss means the accuracy / loss of the saved model.
3

| Order | MODE | Dataset | Scale | > x images each taxon | Number of categories | Learning rate in the end | Epochs | Max. training accuracy | Min. training loss | Max. validation accuracy | Min. validation loss | Last epoch test accuracy | Last epoch test loss | Test accuracy | Test loss |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 13 | Genus | BBFID-1 | C | 10 | 156 | 1.25E-05 | 49 | 0.9972 | 0.0080 | 0.5990 | 1.8758 | 0.5848 | 1.9234 | 0.5834 | 1.9320 |
| 14 | Genus | BBFID-1 | B | 50 | 34 | 1.25E-05 | 34 | 0.9939 | 0.0281 | 0.7185 | 1.1308 | 0.6916 | 1.1420 | 0.7173 | 1.1142 |
| 15 | Genus | BBFID-1 | A | 100 | 13 | 5.00E-05 | 29 | 0.9866 | 0.0446 | 0.8090 | 0.6661 | 0.8256 | 0.6719 | 0.8210 | 0.6650 |
| 16 | Genus | BBFID-2 | C | 10 | 223 | 1.00E-04 | 22 | 0.9908 | 0.0848 | 0.5320 | 2.1067 | 0.4919 | 2.2493 | 0.5004 | 2.2021 |
| 17 | Genus | BBFID-2 | B | 50 | 32 | 5.00E-05 | 21 | 0.9929 | 0.0483 | 0.7370 | 0.9765 | 0.7170 | 1.0273 | 0.7135 | 1.0625 |
| 18 | Genus | BBFID-2 | A | 100 | 9 | 5.00E-05 | 25 | 0.9878 | 0.0486 | 0.8636 | 0.5007 | 0.8259 | 0.5409 | 0.8543 | 0.4904 |
| 19 | Genus | BBFID | C | 10 | 379 | 2.50E-05 | 35 | 0.9974 | 0.0134 | 0.5567 | 2.0772 | 0.5353 | 2.2279 | 0.5371 | 2.2333 |
| 20 | Genus | BBFID | B | 50 | 66 | 2.50E-05 | 27 | 0.9933 | 0.0299 | 0.7335 | 1.1080 | 0.7192 | 1.1866 | 0.7066 | 1.2000 |
| 21 | Genus | BBFID | / | 60 | 47 | 1.25E-05 | 34 | 0.9961 | 0.0177 | 0.7538 | 1.0721 | 0.7742 | 0.8506 | 0.7626 | 0.8921 |
| 22 | Genus | BBFID | A | 100 | 22 | 5.00E-05 | 26 | 0.9907 | 0.0335 | 0.8261 | 0.6590 | 0.8190 | 0.6615 | 0.8145 | 0.6759 |
| 23 | Species | BBFID-1 | E | 6 | 241 | 5.00E-05 | 31 | 0.9949 | 0.0345 | 0.6117 | 1.8168 | 0.5971 | 1.9054 | 0.6080 | 1.9233 |
| 24 | Species | BBFID-1 | D | 8 | 179 | 1.00E-04 | 28 | 0.9938 | 0.0645 | 0.6251 | 1.6484 | 0.5810 | 1.8759 | 0.6299 | 1.6987 |
| 25 | Species | BBFID-1 | C | 10 | 148 | 2.50E-05 | 32 | 0.9975 | 0.0289 | 0.6629 | 1.4035 | 0.6642 | 1.4147 | 0.6790 | 1.4161 |
| 26 | Species | BBFID-1 | B | 50 | 8 | 5.00E-05 | 27 | 0.9871 | 0.0789 | 0.7460 | 0.7560 | 0.7984 | 0.7489 | 0.8140 | 0.6747 |
| 27 | Species | BBFID-2 | E | 6 | 396 | 1.00E-04 | 23 | 0.9950 | 0.0726 | 0.5128 | 2.3015 | 0.4677 | 2.5728 | 0.4813 | 2.5160 |
| 28 | Species | BBFID-2 | D | 8 | 265 | 1.00E-04 | 28 | 0.9983 | 0.0492 | 0.5590 | 1.9957 | 0.5411 | 2.0768 | 0.5349 | 2.1075 |
| 29 | Species | BBFID-2 | C | 10 | 195 | 1.00E-04 | 25 | 0.9969 | 0.0647 | 0.6162 | 1.6714 | 0.5540 | 1.9768 | 0.5791 | 1.8711 |
| 30 | Species | BBFID-2 | B | 50 | 8 | 5.00E-05 | 24 | 0.9968 | 0.0472 | 0.9494 | 0.1308 | 0.9615 | 0.1806 | 0.9519 | 0.1610 |
| 31 | Species | BBFID | / | 2 | 1436 | 5.00E-05 | 41 | 0.9956 | 0.0271 | 0.4975 | 2.4540 | 0.4274 | 2.8980 | 0.4330 | 2.9233 |
| 32 | Species | BBFID | / | 4 | 914 | 1.00E-04 | 28 | 0.9920 | 0.0758 | 0.4958 | 2.4228 | 0.4707 | 2.5650 | 0.4899 | 2.5005 |
| 33 | Species | BBFID | E | 6 | 637 | 1.00E-04 | 25 | 0.9934 | 0.0677 | 0.5521 | 2.0340 | 0.5067 | 2.3438 | 0.5142 | 2.2276 |
| 34 | Species | BBFID | D | 8 | 444 | 5.00E-05 | 26 | 0.9975 | 0.0291 | 0.6148 | 1.6785 | 0.5752 | 1.8458 | 0.5957 | 1.8470 |
| 35 | Species | BBFID | C | 10 | 343 | 2.50E-05 | 34 | 0.9991 | 0.0143 | 0.6472 | 1.5119 | 0.6476 | 1.4888 | 0.6397 | 1.5602 |
| 36 | Species | BBFID | B | 50 | 16 | 1.00E-04 | 23 | 0.9787 | 0.1037 | 0.8399 | 0.5760 | 0.8283 | 0.5472 | 0.8283 | 0.5487 |

4