1 **Title:**

2 Performance and limitations of out-of-distribution detection for insect DNA (meta)barcoding

3

4 **Authors:**

5 Tomochika Fujisawa[1,*]

6 Takashi Imai[2]

7

8 **Author affiliations:**

9 [1] Center for Data Science and AI Innovation Research Promotion, Shiga University

10 1-1-1 Bamba, Hikone, Shiga, Japan

11 [2] Department of Data Science, Shiga University

12 1-1-1 Bamba, Hikone, Shiga, Japan

13

14 *Corresponding author

15 Tomochika Fujisawa t.fujisawa05@gmail.com

16

28

## ABSTRACT

Successful applications of DNA barcoding/metabarcoding rely on the accurate taxonomic identification of sequence fragments. When biological surveys with DNA (meta)barcoding target underexplored biological communities, sequence-based identification is often conducted using incomplete databases that do not fully cover the regional species pool. Consequently, specimens to be identified may include species not present in reference databases. Such unknown or "out-of-distribution" samples can cause misidentification if left undetected. A similarity cutoff is commonly used to detect out-of-distribution samples before taxonomic assignment, but its effectiveness has not been carefully studied. In this study, we evaluated the performance of out-of-distribution detection for DNA barcoding with genetic distance and deep learning metrics. Using extensively sampled datasets of multiple insect taxa, we measured the performance of identification and out-of-distribution detection under conditions in which genetic variations in species were sufficiently sampled. Although identification with DNA barcoding is a highly accurate process, even with short noisy fragments, out-of-distribution detection was more susceptible to a reduction in performance due to sequence noise and a lack of diagnosable characters. Our results provide guidelines for designing unknown-proof identification procedures by determining factors affecting out-of-distribution detection performance.

47

**INTRODUCTION**

The reliable identification of specimens to known taxonomic groups is the foundation of biological studies. Without accurate identification, subsequent practices, including conservation, biosecurity, and ecological monitoring may not be conducted reliably. Despite its importance, taxonomic expertise is a scarce resource, and the identification of specimens is a major bottleneck in large-scale ecological surveys (Van Klink et al. 2024).  Driven by the general lack of taxonomic expertise and the need for rapid and broad characterization of threatened biodiversity, replacing or complementing human identification with computational methods has attracted attention in recent decades (MacLeod et al. 2010; Gaston & O'Neil 2004).  In particular, methods based on DNA sequences have been considered promising because they can enable high-resolution, species-level identification that is accessible to non-experts. DNA barcoding (Hebert et al. 2003), which is the process of identification based on standardized short DNA fragments (e.g., CO1 fragments for animal identification), is the most successful project of such attempts. Currently, the Barcode of Life Data system has 16 million registered sequences, and identification using barcoding markers is routinely performed. The recent introduction of high-throughput sequencing technologies has broadened the scope of DNA barcoding applications.  A notable example is DNA metabarcoding (Taberlet et al. 2012), parallel sequencing and identification of barcoding markers, either from a bulk sample of organisms or environmental DNA. Metabarcoding has significantly expanded the scale of biological monitoring by increasing throughput (Srivathsan et al. 2019) and has widened research targets to previously neglected communities, such as meiofauna or soil arthropods (Macher et al. 2024; Depheide et al. 2019).

Although promising, modern applications of DNA metabarcoding are practiced under conditions that are substantially different from those for which DNA barcoding was initially

73  designed, posing new methodological challenges. For example, identification is conducted

74  using fragments shorter than the original barcoding markers (originally up to 1000 bp, but

75  now shorter than 400bp) (Leese et al. 2021; Miya et al. 2015; Leray et al. 2013) because of the

76  limitations associated with efficient PCR amplification and high-throughput sequencing.

77  Taxonomic identities are often recovered under noisy conditions involving more sequencing

78  errors and artifacts. In addition, metabarcoding surveys target largely unexplored biota whose

79  members are undescribed and there is a lack of representative sequences in reference

80  databases. Hence, the re-optimization of identification procedures has ensued since the

81  introduction of high-throughput technologies, including molecular protocols and bioinformatic

82  pipelines (Creedy et al. 2021; Alberdi et al. 2017).

83

84  One aspect of metabarcoding applications requiring close attention is the effect of samples of

85  the class not present in the reference database. Samples of classes that are not present in the

86  reference are called by different names in different application fields, reflecting the nature of

87  such samples: unknown, novelty, anomaly, outliers and out-of-distribution. We use out-of-

88  distribution (hereafter, OOD) samples in this study because it is a general term that sufficiently

89  encompasses our task.  Also, in the context of sequence-based taxonomic identification,

90  absence from a reference dataset does not immediately define the exact nature of a sample

91  such as "novelty".

92

93   It has been shown that the current sequence databases do not fully represent the diversity of

94  life, and this trend is especially prominent for highly diverse groups such as arthropods.

95  According to previous studies, less than 20% of described invertebrate species have sequence

96  records in public repositories (Keck et al. 2023), and only approximately 20% of terrestrial

97  arthropod species have been formally described (Stork 2018). The use of underrepresented

98   databases for metabarcoding surveys inevitably results in OOD samples. Indeed, large-scale

99   metabarcoding applications have reported many unknown species, even from among

100  supposedly well-studied fauna (Buchner et al. 2024). Because current reference databases are

101  incomplete and undetected OOD samples certainly cause misidentification, molecular

102  identification methods under metabarcoding projects require the appropriate handling of OOD

103  samples. In short, any identification algorithm should be able to say, "I do not know."

104

105  The treatment of OOD samples has been an important issue for practical applications involving

106  DNA barcoding because encountering them is the norm rather than an exception in most

107  conditions. Protocols for detecting unknowns with distance thresholds have been considered

108  even in the very early stages of DNA barcoding studies (Meier et al. 2006). Empirical

109  thresholds with sequence similarity are still commonly used to remove putative unknown

110  samples or retain them for higher-class assignments (e.g. a 97% similarity threshold). This

111  approach also includes thresholds using reliability scores instead of distance, such as bootstrap

112  uncertainty scores (Murali et al. 2018; Porter et al. 2014). More recently, methods that

113  explicitly model the encounters of unknown or new species have been introduced. Methods

114  such as PROTAX and BayesANT (Zito et al. 2023; Somervuo et al. 2016) explicitly model the

115  probability of finding unknowns under species sampling models and infer the sample's

116  posterior probability of being an unknown species.

117

118  Nevertheless, apart from these studies, performance evaluations of OOD detection procedures

119  have not often been conducted systematically. Performance evaluation often does not enable

120  us to distinguish between the failure of OOD detection and misidentification of in-distribution

121  samples, even if these two types of errors represent failures at different steps in the

122  identification process. Critical issues, such as the methods that are favorable for use under

123 certain conditions, have not been clearly addressed, and the major determinants of detection

124 performance are unknown.

125

126 In this study, we evaluated molecular taxonomic assignment methods for DNA barcoding, with

127 a specific focus on the effects of incomplete databases and the presence of OOD samples. We

128 used insect DNA barcoding data to conduct the performance evaluation. Insects are major

129 targets of DNA barcoding projects because of their extreme diversity and the difficulties

130 associated with manual identification. Recent large-scale inventorying efforts (Roslin et al.

131 2022; Hebert et al. 2016) have enabled performance testing under ideal conditions, where a

132 sufficient number of samples are available to characterize species genetic variations across

133 clades. The large sample size also enables the training of parameter-rich machine learning

134 models, including deep learning models. Recently, deep learning has been successfully applied

135 to various biological sequence analyses, including taxonomic classification (Romeijn et al.

136 2024; Ziemski et al. 2021; Busia et al. 2018), but its performance in incomplete databases is

137 unexplored.

138

139 We tested the performance of conventional and deep learning algorithms for taxonomic

140 assignment and OOD detection using extensively sampled insect taxa. We explored the

141 performance limits of identification and OOD detection methods by focusing on insect taxa

142 with sufficient within- and between-species sampling. We showed that both conventional and

143 deep learning identification methods are highly accurate for taxonomic assignment and robust

144 for short and noisy sequences, whereas OOD detection is more prone to performance reduction

145 due to noise and limited availability of information in short fragments.

146

**MATERIALS AND METHODS**

*Data acquisition*

We tested the performance of the classification and OOD detection models using datasets downloaded from the Barcode of Life Data (BOLD) System database (Ratnasingham & Hebert 2007). We first conditioned database entries by geographic regions where comprehensive inventorying of regional insect fauna was underway (mainly North American and EU countries), and then selected insect genera with sufficient sample size, taxonomic coverage, and geographic extent. Genera were selected as targets when at least 15 species were represented by 15 or more individuals (Fifteen individuals within species is a sample size large enough for correct estimation of within-species genetic variations upon training) (Zhang et al. 2010; Matz & Nielsen 2005). We applied this criterion to four major insect groups, i.e., Hymenoptera (bees and wasps), Diptera (flies), Lepidoptera (moths and butterflies), and Coleoptera (beetles), as these groups had the densest and broadest samples in the BOLD database. Then, twenty candidate genera were randomly selected. Classification models were trained to conduct species-level identification within the genus. A full list of BOLD accession numbers and sequence alignments are available in the Supplementary Data.

*Data preparation*

The downloaded CO1 sequences of the 20 genera were filtered according to length and sequenced regions. Only fragments with length > 400 bp and < 1000 bp and fragments with the "5BP" DNA barcoding region were retained. Sequences with missing bases in the latter half of the 5BP region were discarded because they did not contain the short barcoding regions used for performance evaluation. Samples identified only at the genus level were removed, but samples with unconventional labels, such as "sp. 10" or "sp. DNAS-***" were included in OOD datasets because these names were consistently applied to multiple samples and likely

172   represented true unknown species. To reduce the adverse effects of overrepresented species,

173   species with >125 samples were randomly resampled to reduce the number of samples to 125,

174   which was sufficient to characterize the genetic and haplotype diversity of the focal species

175   (Zhang et al. 2010).

176

177   The filtered sequences were then aligned using MAFFT (v.7.453, Katoh et al. 2013) with default

178   parameters. Aligned sequences were split into in-distribution (ID) and OOD samples based on

179   the number of individuals in the species (based on the logic that rare species are more likely to

180   be OOD). Species with ≥ 15 samples were assigned to ID and the other species were assigned to

181   OOD. Classification models were first trained to classify ID samples into their taxonomic groups

182   and were subsequently exposed to OOD samples to test whether the models could correctly

183   detect them.

184

185   We compiled multiple datasets covering various parameters, including the total number of

186   samples, fragment lengths, and sequencing errors. We prepared two short alignments by

187   selecting the 350-650 region (300 bp) and 350-500 (150 bp) within full alignments. These 300

188   bp and 150 bp regions largely overlap with the short barcoding regions proposed by Lelay et

189   al. (2013) and Leese et al. (2021), respectively.  We also randomly halved the number of

190   samples to create smaller datasets in which at least five samples per species were retained in

191   the training process. We named these halved datasets "*Small*" and original datasets "*Sufficient*."

192   To simulate sequencing errors, bases were randomly swapped with a noise rate of 0.02, where

193   randomly selected 2% of bases in a fragment were replaced with one of the alternative bases

194   with an equal probability of 1/3. Errors were introduced only into test datasets because only

195   clean reference databases were available for model training in realistic applications. The final

196   datasets covered two database size categories {"Sufficient," "Small"}, two noise levels {0.0%,

197     2%} and three fragment lengths {650 bp, 300 bp, 150 bp}.

198

199     *Deep learning model for taxonomic classification*

200     ---CNN model

201     The convolutional neural network (CNN) classification model employed a typical convolutional

202     architecture used in multiple studies (Jiang et al. 2023; Zheng et al. 2019; Busai et al. 2019),

203     consisting of convolutional blocks for feature extraction and subsequent fully connected (FC)

204     classification layers. The convolutional part has three consecutive convolutional blocks with

205     each sequentially consisting of the 1D convolution, batch normalization, 1D max pooling,

206     rectified linear unit (ReLU, $ReLU(x) = x$ if $x > 0$ otherwise 0), and dropout (rate = 0.15 for CNN

207     layers and 0.25 for FC layers) layers.  There were 64, 128, and 128 channels in the first, second,

208     and third convolutional blocks, respectively. Hyperparameters, including the number of

209     channels and dropout rates, were determined using cross-validation runs on a partial dataset.

210

211     An input DNA sequence with length L was encoded in an L × 4 matrix whose rows were four-

212     dimensional one-hot vectors.  For instance, a base letter "A" was represented as a row [1, 0, 0,

213     0], and "T" as [0, 1, 0, 0]. Noncanonical base letters (N, R, Y, etc.) were represented as [0, 0, 0,

214     0]. For each convolution process, the length was halved by one-dimensional (1D) max pooling

215     with a size of two, resulting in an L/8 ×128-dimensional output. One-dimensional global

216     average pooling was then applied to the outputs to obtain a 128-dimensional feature vector.

217     Subsequently, classification was performed with three FC layers to classify the input sequences

218     into known taxa. The softmax function was applied to the final output of the FC layer to obtain

219     prediction probabilities. Details of the neural network architecture are presented in

220     Supplementary figure S1.

221

222    Throughout the performance evaluation process, the models were trained using the Adam

223    algorithm with a cross-entropy loss. Default hyperparameter settings provided by Keras were

224    used (batch size=16, learning rate=0.001). The convergence of loss was visually assessed.

225

226    ---Deep learning methods for OOD detection

227    In addition to the taxonomic classification model described above, we implemented deep

228    learning methods for out-of-distribution (OOD) detection. OOD detection is the task of

229    separating samples into two categories: *IN-DISTRIBUTION,* hereafter, *ID*, which includes

230    samples from classes present in the training data, and *OUT-OF-DISTRIBUTION* or *OOD*, which

231    includes samples from classes NOT present in the training data (Zhang et al. 2024). The

232    accepted ID samples were subsequently classified into ID classes. We employed three methods

233    based on the prediction uncertainty scores. We selected these methods based on their reported

234    performances (Zhang et al. 2023) and implementation complexities. Methods designed to work

235    without explicit OOD sample exposure during the training phase were selected, because OOD

236    exposure is not feasible for real barcoding applications. We also excluded methods that

237    required complex optimization of hyperparameters.

238

239    Three OOD scores were calculated from the output obtained from intermediate FC layers.

240    When the output of the penultimate FC layers was $g(x)$, the following transformation to $g(x)$

241    was applied in the final FC layer:

242

243
$$f(x) = mg(x) + a$$

244

245    Here, $f(x)$ is the output of the final FC layer, which is a vector of length equal to the number of

246    classes; $m$ is a weight matrix; and $a$ is an offset vector. Each of these parameters were

247    optimized in the training process. These intermediate outputs, *f(x)* and *g(x)*, contain useful

248    information for discriminating OODs from ID samples (Supplementary figure S2).

249

250    --Maximum softmax probability

251    The maximum softmax probability (MSP) is commonly used as the prediction probability for

252    neural network classification. The MSP score is defined as a function of the processes of

253    exponentiation and scaling of *f(x)*, the output of the final FC layer, and its maximum value:

254

255
$$MSP(x) = \max_k \left( \frac{\exp(f_k(x))}{\sum_{k=1}^{K} \exp(f_k(x))} \right)$$

256

257     Here, $f_k(x)$ is the k-th component of the vector *f(x)*. The kth class that yields the MSP (i.e.,

258    $\underset{k}{\arg\max}$) is a predicted assignment of sample *x*. Importantly, this predicted class is chosen only

259    from the classes present in the training dataset, regardless of whether the sample is of the OOD

260    type.  Hence, OOD detection is required to avoid the erroneous assignment of an OOD sample to

261    a known class.  Hendrycks and Gimpel (2016) proposed MSP as a metric for prediction

262    uncertainty and showed that MSP scores of OOD samples were consistently lower than those of

263    the ID sample, and a cutoff by a threshold of prediction probability helped to successfully

264    detect OOD samples.

265

266    --Energy score

267    Liu et al. (2020) introduced the "energy score" of a neural network model for OOD detection.

268    The log energy score of a neural network is defined as

269

270
$$E(x) = -\log\left(\sum_{k=1}^{K} \exp(f_k(x))\right)$$

271

272 The log energy score is the logarithm of the softmax denominator in *MSP(x)*. The energy score

273 is interpreted as the relative log-likelihood score of a model given a sample *x, Pr(x|model)*. The

274 energy score of OOD samples was consistently lower than that of ID samples because the

275 likelihood of obtaining such samples is less for models trained only with ID samples. Liu et al.

276 (2020) reported that the threshold of sample energy values outperformed the softmax

277 probability for multiple OOD detection tasks.

278

279 --Mahalanobis distance

280 Lee et al. (2018) developed a distance-based OOD detection method. The Mahalanobis distance

281 of a sample from a class center is defined as

282

283
$$d_k(x) = (g(x) - \widehat{\mu_k})^T \hat{\Sigma}^{-1} (g(x) - \widehat{\mu_k})$$

284

285 where $\mu_k$ is the k-th class center value, and $\Sigma$ is a variance-covariance matrix of *g(x)*.

286 Mahalanobis distance measures the distance from the k-th class center, assuming that the

287 distribution of *g(x)* follows a multivariate normal distribution with a mean $\mu_k$ and a single

288 variance-covariance matrix, $\hat{\Sigma}^{-1}$, which are empirically estimated from a distribution of *g(x)* in

289 a training data set. Lee et al. (2018) proposed the following negative Mahalanobis distance to

290 the closest distribution center as an uncertainty metric for OOD detection:

291

292
$$M(x) = \max_k (-d_k(x))$$

293

294 --Majority voting for OOD detection

295 In addition to independent OOD detection procedures with the above metrics, we devised a

296 process for OOD detection with majority voting for the above three detectors. With this

297 approach, a sample was treated as an OOD sample if two of the three methods "vote" for the

298 presence of OOD.

299 All deep-learning models were implemented in Python using the Keras library. The code is

300 available at https://github.com/tfujisawa/barcoding_cnn.

301

302 *Model training and performance test*

303 The CNN models were trained with 70% of the ID data and their baseline prediction accuracy,

304 the proportion of correct identifications to the total identification trials of the test samples, was

305 measured. We then calculated the OOD scores (softmax probability, energy score, and

306 Mahalanobis distance) for all ID test samples and obtained class thresholds by accounting for

307 the 95% quantiles of all classes (Supplementary figure S3). After setting the thresholds, the

308 model was exposed to OOD samples, and their scores were calculated. Samples with more

309 extreme values than class-wise threshold values were classified as OODs. The proportion of

310 OOD samples falsely classified as ID samples was measured as the false negative rate at a 95%

311 threshold (FNR@95%). The training and evaluation processes were repeated 20 times for each

312 dataset. The effects of fragment length, dataset size, noise level, and methods to determine the

313 identification performance were assessed using multivariate linear regression. To assess the

314 difficulty of the classification tasks, we calculated the proportion of misidentifications with

315 zero genetic distances, i.e., the proportion of cases in which heterospecific specimens had

316 identical sequences that led to misidentification of ID samples. The proportion of OOD

317 specimens with zero genetic distances from any ID sample was also calculated. These zero-

318 genetic-distance proportions determine the upper limits of classification accuracy and OOD

319 detection (i.e., "perfect classifier". Ziemski et al. 2021).

320

321 *Classification and OOD detection methods with distance*

322 Conventional classification methods based on sequence distances were used as performance

323 baselines. We measured the pairwise K2P genetic distance for the aligned sequences and the

324 BLAST percentage similarity for the unaligned matrices (Altschul et al. 1990). Distance-based

325 classification was then performed using the 1-nearest neighbor criterion (1NN), where a new

326 sample was assigned to the taxon of a sample with the smallest distance from it. Although 1NN

327 based on the K2P or BLAST distance is the simplest distance-based classification algorithm, it is

328 still widely used and often outperforms more sophisticated algorithms (Leray et al. 2022;

329 Hleap et al. 2021). For OOD detection tasks, we calculated the minimum distances from

330 samples within their own class/species and set class OOD thresholds by taking the 95%

331 quantiles of their minimum distances. When the distance between a test OOD sample and its

332 nearest neighbor is greater than the class OOD threshold, the sample is classified as an OOD

333 sample. The above procedure is similar to the "best close match" procedure proposed in Meier

334 et al. (2006) although the quantile calculation process is different.

335

336 *Gradient-based attribution*

337 We visualized the region responsible for classification decisions using a one-dimensional

338 gradient-based class activation map (GradCAM, Selvaraju et al. 2016). GradCAM localizes the

339 region of importance by measuring the effects of CNN features on the classification

340 probabilities. Specifically, the GradCAM score on window $w$, is defined as

341

342
$$L_{w,GradCAM} = \text{ReLU}\left(\sum_{n=1}^{N} \alpha_n A_{n,w}\right)$$

343

344 $L_{w,GradCAM}$ is a weighted average of $N$ CNN features, $A_{n,w}$, calculated on the window $w$ with the

345    weight $\alpha_n$. The weight is a feature importance, measured as an averaged partial derivative of

346    MSP(x) with respect to $A_{n,w}$, $\alpha_n = \frac{1}{W}\sum_{w=1}^{W}\frac{\partial MSP(x)}{\partial A_{n,w}}$, where $W$ is the total number of windows on a

347    sequence. An interpretation of importance is that when a unit change in a CNN feature ($A_{n,w}$)

348    results in a significant change in the prediction probability (MSP(x)), $A_{n,w}$ is considered to be

349    important in the prediction process. We also implemented an activation map of the energy

350    score to visualize the region responsible for OOD detection decisions by replacing the gradient

351    of the MSP(x) in the weight calculation with the gradient of the energy score.

352

353
$$L_{w,Grad-Energy} = \text{ReLU}\left(\sum_{n=1}^{N} \beta_n A_{n,w}\right)$$

354

355    Here, the weight $\beta_n$ is defined as $\beta_n = \frac{1}{W}\sum_{w=1}^{W}\frac{\partial E(x)}{\partial A_{n,w}}$. In this case, the effect of the CNN features on

356    the energy score was measured. In the current study, the window size was set to 8 bp, resulting

357    in 85 windows in a 680 bp fragment. We compared the GradCAM and Grad-Energy scores with

358    genetic variations measured in 16 bp windows.

359

360    *Regression by population genetic metrics*

361    A set of population genetic metrics was calculated for each genus dataset to identify the

362    determinants of the identification performance. Genetic distance-related metrics, including

363    average within-species and average and minimum between-species genetic distances, were

364    calculated from alignments. Dataset completeness was measured by the average number of

365    samples per species, the total number of species, and "taxonomic completeness," defined by the

366    number of ID species divided by the total number of all species. To identify factors affecting

367    model performance, multivariate regression modeling was conducted using the above metrics

368    as explanatory variables, and identification accuracy and false negative rate as responses. Least

369     absolute shrinkage and selection operator (LASSO) procedures were used to select important

370     explanatory variables.

371 **RESULTS**

372 *BOLD dataset profiles*

373 We collected 34,408 COI sequences from 20 genera in the BOLD database. Of the 13,078

374 examined genera in the four target orders, only 82 (0.6%) met the sample size criteria. The

375 number of in-distribution (ID) and out-of-distribution (OOD) samples were 28,422 and 5,986,

376 respectively. The number of species within the selected genera ranged from 15 to 68, and the

377 average number of samples per species was 45. The number of OOD samples per genus ranged

378 from 23 to 910, and the proportion of OOD samples to total samples was 0.17.

379

380 *Accuracy of identification and OOD detection*

381 Both deep learning and distance-based methods were highly accurate in ID identification tasks,

382 especially when trained with sufficiently large datasets. For whole 650 bp fragments, the

383 average baseline prediction accuracy of the CNN model was 0.97, and the two conventional

384 methods were as accurate as the CNN model (0.971 for k2p distance, 0.973 for BLAST; Fig. 2

385 and Table. 1). The accuracy decreased with reduced fragment sizes for all methods, and the

386 CNN slightly outperformed the conventional methods when the fragment length was 150 bp

387 (0.960 for CNN, 0.945 for k2p distance, and 0.946 for BLAST). Multiple linear regression

388 analyses showed that shorter fragments were significantly associated with lower accuracy, and

389 the CNN model exhibited slightly higher accuracy, but the difference was not significant. When

390 the training datasets were smaller, the CNN performance decreased, whereas the distance

391 methods were less affected. The introduction of 2% noise to the sequence reduced the

392 identification performance; however, the reduction in accuracy was within 2% under most

393 conditions (average accuracy decrease = 0.011, p<<0.001).

394

395 The performance of OOD detection tasks generally exhibited patterns similar to those of

396 identification tasks. The CNN model underperformed conventional methods with long

397 fragments (FNR@95% 0.128 for CNN, 0.103 for k2p distance, and 0.101 for BLAST for 650 bp

398 fragments; Figure 3 and Table 2) but outperformed with shorter fragments. However, the

399 effect of reduced fragment size was more pronounced (FNR@95%: 0.156 for CNN, 0.176 for

400 k2p distance, and 0.170 for BLAST for 150bp fragments). There was no significant difference in

401 FNR between the detection methods. Among the deep learning methods for OOD detection, the

402 performances of the energy score and Mahalanobis distance were closely matched, and these

403 methods significantly outperformed MSP (Supplementary figure S4 and Table S2). Consensus

404 across the three methods generally resulted in better performance, but the improvement was

405 not significant, and the best-performing methods depended on the datasets. The proportion of

406 OOD samples with zero distance from ID samples was 0.061 for 650 bp and reached 0.161 for

407 150 bp. The performance of the OOD detection method was close to these optimal values,

408 although the gaps were greater for longer fragments ($FNR_{perfect}$=0.061 vs. $FNR_{CNN}$=0.128 for

409 650 bp and 0.156 vs. 0.161 for 150 bp).

410

411 Sequences with noise significantly compromised the OOD detection performance for each

412 method (Average FNR increase=0.064, p<<0.001). In reduced-size datasets, FNRs of the

413 distance-based methods were slightly improved. The average *decrease* in the FNRs for small

414 datasets over sufficient datasets was 0.0079, and linear regression analysis showed that the

415 effect was significant (p=0.0026).  These counterintuitive results are attributable to the

416 reduced number of identical sequences shared between OOD samples and their nearest ID

417 counterparts.

418

419 *Regression modeling*

420 The results of the multiple regression analysis with LASSO variable selection are summarized

421    in Table 3. Regression modeling showed that identification accuracy was positively correlated

422    with the number of samples per species and the minimum between-species distance and

423    negatively correlated with the number of classes and taxonomic completeness. FNR@95% was

424    negatively correlated with the minimum between-species distance and positively correlated

425    with the average within-species distance. A minor negative effect of the number of classes was

426    also observed, while other variables were excluded. In addition, the identification accuracy and

427    FNR were significantly correlated (Pearson's r = -0.48, Figure 4), indicating that when the

428    model correctly identified the ID classes, its OOD detection ability was accurate.

429

430    *Gradient-based attribution*

431    The sequence regions important for classification localized by GradCAM largely corresponded

432    to regions with high genetic variation in the alignment. The genetic variations in 16 bp

433    windows were strongly correlated with average GradCAM scores on the same windows

434    (Pearson's $r$ = 0.41-0.76), and peaks were often aligned with the highest genetic variations.

435    This trend was consistently observed across fragment lengths (Figure 5). In contrast, a weaker

436    correspondence was observed between the regions of importance in energy-based OOD

437    detection and regions with high genetic variation (Figure 6), and the grad-energy score was

438    less correlated with genetic variation (Pearson's $r$ = 0.06 - 0.45) that were always lower than

439    those of GradCAM.

440

**DISCUSSION**

Under the conditions considered in this study, sequence-based identification methods were highly accurate and robust against sequence noise when sufficient samples were available. The models slightly underperformed with short fragment lengths or smaller training datasets, as reported in previous studies (Porter & Hajibabaei 2018), but the best-performing models retained ~95% accuracy. Regardless of minor performance differences, DNA barcoding identification methods appear to have already been optimized, and their performance is very close to that of the ideal classifier. The reduction in accuracy was largely due to the reduced number of diagnostic variations, as reported by Ziemski et al. (2021).

By contrast, sequence-based out-of-distribution (OOD) detection was a more refractory task, with higher error rates for short and noisy fragments. In addition, the performance was more counterintuitively dependent on the database size (e.g., improved FNR@95% with *smaller* databases). More importantly, the accuracy of OOD detection was more strongly limited by samples that were undiagnosable by sequencing alone, comprising ~16% of OOD samples under some extreme conditions. Although there is room for improvement in OOD detection using long fragments and smaller databases, a strong limiting factor is the lack of diagnostic characters for short fragments.

The risk of overlooking OOD samples has been recognized but considered difficult to quantify (Virgilio et al. 2010). This study provides a coarse estimate of these risks. Assuming that the current proportion (17%) of bulk specimens is of the OOD type, up to ~3% of the total specimens may be misidentified as referenced species, with errors increasing when targets containing more unknown samples or noisy sequences are used. Under such conditions, sequence-based surveys may significantly underestimate unknown biodiversity even with the

466  best identification methods. Because most insect species have not been barcoded and are

467  highly likely to have variations that will be missed during fragment truncation, it is prudent to

468  use as long fragments as possible to minimize the risk of overlooking them.  Because the

469  performance of OOD detection was correlated with within-species distance and minimum

470  interspecific distance, these metrics may be used to help determine the appropriate marker

471  length. In addition, the correlation between ID classification and OOD detection performance

472  can be used to assess potential risks and improve detection performance during the training

473  process (Vaze et al. 2021).

474

475  In this study, we compared distance-based methods with deep learning. Although their overall

476  performances were similar, the general tendency was that short fragments favored the CNN

477  model, and longer fragments favored the distance methods. Deep learning classification

478  performed poorly with small database sizes and long fragment lengths. This performance

479  reduction may reflect difficulties in optimizing highly parameter-rich models (~110k

480  parameters).  Our deep learning OOD detection methods exhibited a performance comparable

481  to that of the distance methods. Therefore, they may be used to mitigate the reported

482  performance reduction of deep learning models due to incomplete training databases (Romeijn

483  et al. 2024).

484

485  Deep learning models also provide useful information for interpreting results. Gradient-based

486  explanations of CNN classification showed that highly variable regions in the alignment were

487  informative for the classification tasks of in-distribution (ID) samples but not necessarily for

488  OOD detection tasks.  This result may reflect the different natures of the two tasks. For

489  example, in an extreme case, a site might be completely invariable among ID samples, while an

490  OOD species harbors a diagnosable difference at the same site. Under these conditions, the site

491    is uninformative for classification but highly informative for OOD detection. Such site

492    importance score may be useful for distinguishing favorable sites for different tasks and for

493    selecting informative markers.

494

495    We applied only a limited number of OOD detection methods in this study, and methodological

496    improvements may exist. For example, training models with ID and synthetic OOD samples

497    may potentially improve detection performance. Nevertheless, as taxonomic coverage and

498    within-species sample sizes improve, identification success will depend more on sequence

499    variation than on identification procedures. Except for using longer fragments, a

500    straightforward improvement is sequencing multilocus markers to increase the available

501    diagnosable variations. Although sequencing short multilocus markers can be cost-effective

502    (Wang et al. 2023), they may result in incongruent species compositions even in a single

503    community because of different PCR affinities. An alternative approach to OOD tolerant

504    identification is to use additional information, such as geographic locations, environmental

505    information, and morphological features. For example, fine-grained geographic information

506    can not only be used to identify species, but also to detect possible OOD samples because insect

507    communities can have extremely high geographic turnover (Srivathsan et al. 2023; Arribas et

508    al. 2021). When closely related species occupy different niches, environmental niche modeling

509    can provide additional information for identification (Yang et al. 2024). A similar approach

510    may be used for OOD detection. High-throughput imaging is another potentially useful tool to

511    supplement metabarcoding (Fujisawa et al. 2023; Wührl et al. 2022) and has been successfully

512    used to verify the metabarcoding results (Panel et al. 2025).  Machine learning algorithms may

513    help integrate multiple information sources because designing "multimodal" models

514    combining different types of data is easier than using conventional statistical methods.

515

516    In summary, sequence-based identification with DNA barcoding is highly accurate,

517    owing to collective efforts for performance improvement.  However, incomplete databases and

518    the presence of OOD samples still pose methodological challenges, and a careful experimental

519    design is required to avoid overlooking these unknowns. In the future, machine learning

520    models should integrate multiple sources of information for more robust and unknown-proof

521    identification. The rapid accumulation of DNA sequence databases and additional ecological

522    information, along with advanced machine learning algorithms, may enable the deployment of

523    integrated biodiversity monitoring systems.

524

525

526    Acknowledgements

528

529

530

**References**

Alberdi, A., Aizpurua, O., Gilbert, M. T. P., & Bohmann, K. (2018). Scrutinizing key steps for reliable metabarcoding of environmental samples. *Methods in Ecology and Evolution*, *9*(1), 134–147. https://doi.org/10.1111/2041-210X.12849

Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment search tool. *Journal of Molecular Biology*, *215*(3), 403–410. https://doi.org/10.1016/S0022-2836(05)80360-2

Arribas, P., Andújar, C., Salces-Castellano, A., Emerson, B. C., & Vogler, A. P. (2021). The limited spatial scale of dispersal in soil arthropods revealed with whole-community haplotype-level metabarcoding. *Molecular Ecology*, *30*(1), 48–61. https://doi.org/10.1111/mec.15591

Buchner, D., Sinclair, J. S., Ayasse, M., Beermann, A. J., Buse, J., Dziock, F., Enss, J., Frenzel, M., Hörren, T., Li, Y., Monaghan, M. T., Morkel, C., Müller, J., Pauls, S. U., Richter, R., Scharnweber, T., Sorg, M., Stoll, S., Twietmeyer, S., … Leese, F. (2024). Upscaling biodiversity monitoring: Metabarcoding estimates 31,846 insect species from Malaise traps across Germany. *Molecular Ecology Resources*, *December 2023*, 1–26. https://doi.org/10.1111/1755-0998.14023

Busia, A., Dahl, G. E., Fannjiang, C., Alexander, D. H., Dorfman, E., Poplin, R., McLean, C. Y., Chang, P.-C., & DePristo, M. (2019). A deep learning approach to pattern recognition for short DNA sequences. *BioRxiv*, 353474. https://www.biorxiv.org/content/10.1101/353474v3%0Ahttps://www.biorxiv.org/content/10.1101/353474v3.abstract

Dopheide, A., Tooman, L. K., Grosser, S., Agabiti, B., Rhode, B., Xie, D., Stevens, M. I., Nelson, N., Buckley, T. R., Drummond, A. J., & Newcomb, R. D. (2019). Estimating the biodiversity of terrestrial invertebrates on a forested island using DNA barcodes and metabarcoding data. *Ecological*

554    *Applications*, *29*(4), 0–14. https://doi.org/10.1002/eap.1877

555    Fujisawa, T., Noguerales, V., Meramveliotakis, E., Papadopoulou, A., & Vogler, A. P. (2023).

556    Image-based taxonomic classification of bulk insect biodiversity samples using deep learning and

557    domain adaptation. *Systematic Entomology*, *48*(3), 387–401.

558    https://doi.org/https://doi.org/10.1111/syen.12583

559    Gaston, K. J., & O'Neill, M. A. (2004). Automated species identification: Why not? In *Philosophical*

560    *Transactions of the Royal Society B: Biological Sciences* (Vol. 359, Issue 1444, pp. 655–667). Royal

561    Society. https://doi.org/10.1098/rstb.2003.1442

562    Hebert, P. D. N., Cywinska, A., Ball, S. L., & DeWaard, J. R. (2003). Biological identifications

563    through DNA barcodes. *Proceedings of the Royal Society B: Biological Sciences*, *270*(1512), 313–

564    321. http://dx.doi.org/10.1098/rspb.2002.2218

565    Hebert, P. D. N., Ratnasingham, S., Zakharov, E. v., Telfer, A. C., Levesque-Beaudin, V., Milton, M.

566    A., Pedersen, S., Jannetta, P., & Dewaard, J. R. (2016). Counting animal species with DNA

567    barcodes: Canadian insects. *Philosophical Transactions of the Royal Society B: Biological Sciences*,

568    *371*(1702). https://doi.org/10.1098/rstb.2015.0333

569    Hendrycks, D., & Gimpel, K. (2016). A Baseline for Detecting Misclassified and Out-of-Distribution

570    Examples in Neural Networks. *5th International Conference on Learning Representations, ICLR*

571    *2017 - Conference Track Proceedings*, 1–12. http://arxiv.org/abs/1610.02136

572    Hleap, J. S., Littlefair, J. E., Steinke, D., Hebert, P. D. N., & Cristescu, M. E. (2021). Assessment of

573    current taxonomic assignment strategies for metabarcoding eukaryotes. *Molecular Ecology*

574    *Resources*, *21*(7), 2190–2203. https://doi.org/https://doi.org/10.1111/1755-0998.13407

575    Jiang, Y., Balaban, M., Zhu, Q., & Mirarab, S. (2023). DEPP: Deep Learning Enables Extending

576      Species Trees using Single Genes. *Systematic Biology*, *72*(1), 17–34.

577      https://doi.org/10.1093/sysbio/syac031

578      Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7:

579      improvements in performance and usability. *Molecular Biology and Evolution*, *30*(4), 772–780.

580      https://doi.org/10.1093/molbev/mst010

581      Keck, F., Couton, M., & Altermatt, F. (2023). Navigating the seven challenges of taxonomic

582      reference databases in metabarcoding analyses. *Molecular Ecology Resources*, *23*(4), 742–755.

583      https://doi.org/10.1111/1755-0998.13746

584      Lee, K., Lee, K., Lee, H., & Shin, J. (2018). A simple unified framework for detecting out-of-

585      distribution samples and adversarial attacks. *Advances in Neural Information Processing Systems*,

586      *2018-Decem*(LID), 7167–7177.

587      Leese, F., Sander, M., Buchner, D., Elbrecht, V., Haase, P., & Zizka, V. M. A. (2021). Improved

588      freshwater macroinvertebrate detection from environmental DNA through minimized nontarget

589      amplification. *Environmental DNA*, *3*(1), 261–276. https://doi.org/10.1002/edn3.177

590      Leray, M., Yang, J. Y., Meyer, C. P., Mills, S. C., Agudelo, N., Ranwez, V., Boehm, J. T., &

591      Machida, R. J. (2013). A new versatile primer set targeting a short fragment of the mitochondrial

592      COI region for metabarcoding metazoan diversity: application for characterizing coral reef fish gut

593      contents. *Frontiers in Zoology*, *10*(1), 34. https://doi.org/10.1186/1742-9994-10-34

594      Leray, M., Knowlton, N., & Machida, R. J. (2022). MIDORI2: A collection of quality controlled,

595      preformatted, and regularly updated reference databases for taxonomic assignment of eukaryotic

596      mitochondrial sequences. *Environmental DNA*, *4*(4), 894–907. https://doi.org/10.1002/edn3.303

597      Liu, W., Wang, X., Owens, J. D., & Li, Y. (2020). Energy-based Out-of-distribution Detection.

598 *Advances in Neural Information Processing Systems*, *2020-Decem*(NeurIPS).

599 http://arxiv.org/abs/2010.03759

600 Macher, J.-N., Martínez, A., Çakir, S., Cholley, P.-E., Christoforou, E., Curini Galletti, M., van

601 Galen, L., García-Cobo, M., Jondelius, U., de Jong, D., Leasi, F., Lemke, M., Rubio Lopez, I.,

602 Sánchez, N., Sørensen, M. V., Todaro, M. A., Renema, W., & Fontaneto, D. (2024). Enhancing

603 metabarcoding efficiency and ecological insights through integrated taxonomy and DNA reference

604 barcoding: A case study on beach meiofauna. *Molecular Ecology Resources*, *24*(7), e13997.

605 https://doi.org/https://doi.org/10.1111/1755-0998.13997

606 MacLeod, N., Benfield, M., & Culverhouse, P. (2010). Time to automate identification. In *Nature*

607 (Vol. 467, Issue 7312, pp. 154–155). Nature Publishing Group. https://doi.org/10.1038/467154a

608 Matz, M. v, & Nielsen, R. (2005). A likelihood ratio test for species membership based on DNA

609 sequence data. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *360*(1462),

610 1969–1974. https://doi.org/10.1098/rstb.2005.1728

611 Meier, R., Shiyang, K., Vaidya, G., & Ng, P. K. L. (2006). DNA barcoding and taxonomy in diptera:

612 A tale of high intraspecific variability and low identification success. *Systematic Biology*, *55*(5),

613 715–728. https://doi.org/10.1080/10635150600969864

614 Miya, M., Sato, Y., Fukunaga, T., Sado, T., Poulsen, J. Y., Sato, K., Minamoto, T., Yamamoto, S.,

615 Yamanaka, H., Araki, H., Kondoh, M., & Iwasaki, W. (2015). MiFish, a set of universal PCR

616 primers for metabarcoding environmental DNA from fishes: Detection of more than 230 subtropical

617 marine species. *Royal Society Open Science*, *2*(7). https://doi.org/10.1098/rsos.150088

618 Murali, A., Bhargava, A., & Wright, E. S. (2018). IDTAXA: A novel approach for accurate

619 taxonomic classification of microbiome sequences. *Microbiome*, *6*(1).

620 https://doi.org/10.1186/s40168-018-0521-5

621    Penel, B., Meynard, C. N., Benoit, L., Boudonne, A., Clamens, A.-L., Soldati, L., Migeon, A.,

622    Chapuis, M.-P., Piry, S., Kergoat, G., & Haran, J. (2025). The best of two worlds: toward large-scale

623    monitoring of biodiversity combining COI metabarcoding and optimized parataxonomic validation.

624    *Ecography*, *n/a*(n/a), e07699. https://doi.org/https://doi.org/10.1111/ecog.07699

625    Porter, T. M., Gibson, J. F., Shokralla, S., Baird, D. J., Golding, G. B., & Hajibabaei, M. (2014).

626    Rapid and accurate taxonomic classification of insect (class Insecta) cytochrome c oxidase subunit 1

627    (COI) DNA barcode sequences using a naïve Bayesian classifier. *Molecular Ecology Resources*,

628    *14*(5), 929–942. https://doi.org/https://doi.org/10.1111/1755-0998.12240

629    Porter, T. M., & Hajibabaei, M. (2018). Automated high throughput animal CO1 metabarcode

630    classification. *Scientific Reports*, *8*(1), 4226. https://doi.org/10.1038/s41598-018-22505-4

631    Ratnasingham, S., & Hebert, P. D. N. (2007). BOLD: The Barcode of Life Data System

632    (http://www.barcodinglife.org). *Molecular Ecology Notes*, *7*(3), 355–364.

633    https://doi.org/10.1111/j.1471-8286.2007.01678.x

634    Roslin, T., Somervuo, P., Pentinsaari, M., Hebert, P. D. N., Agda, J., Ahlroth, P., Anttonen, P., Aspi,

635    J., Blagoev, G., Blanco, S., Chan, D., Clayhills, T., DeWaard, J., DeWaard, S., Elliot, T., Elo, R.,

636    Haapala, S., Helve, E., Ilmonen, J., … Mutanen, M. (2022). A molecular-based identification

637    resource for the arthropods of Finland. *Molecular Ecology Resources*, *22*(2), 803–822.

638    https://doi.org/10.1111/1755-0998.13510

639    Romeijn, L., Bernatavicius, A., & Vu, D. (2024). MycoAI: Fast and accurate taxonomic

640    classification for fungal ITS sequences. *Molecular Ecology Resources*, *August*, 1–22.

641    https://doi.org/10.1111/1755-0998.14006

642    Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2016). *Grad-CAM:*

643    *Visual Explanations from Deep Networks via Gradient-based Localization*.

644    https://doi.org/10.1007/s11263-019-01228-7

645    Somervuo, P., Koskela, S., Pennanen, J., Henrik Nilsson, R., & Ovaskainen, O. (2016). Unbiased

646    probabilistic taxonomic classification for DNA barcoding. *Bioinformatics*, *32*(19), 2920–2927.

647    https://doi.org/10.1093/bioinformatics/btw346

648    Srivathsan, A., Ang, Y., Heraty, J. M., Hwang, W. S., Jusoh, W. F. A., Kutty, S. N., Puniamoorthy,

649    J., Yeo, D., Roslin, T., & Meier, R. (2023). Convergence of dominance and neglect in flying insect

650    diversity. *Nature Ecology & Evolution*, *7*(7), 1012–1021. https://doi.org/10.1038/s41559-023-02066-

651    0

652    Srivathsan, A., Hartop, E., Puniamoorthy, J., Lee, W. T., Kutty, S. N., Kurina, O., & Meier, R.

653    (2019). Rapid, large-scale species discovery in hyperdiverse taxa using 1D MinION sequencing.

654    *BMC Biology*, *17*(1), 96. https://doi.org/10.1186/s12915-019-0706-9

655    Stork, N. E. (2018). How Many Species of Insects and Other Terrestrial Arthropods Are There on

656    Earth? *Annual Review of Entomology*, *63*(Volume 63, 2018), 31–45.

657    https://doi.org/https://doi.org/10.1146/annurev-ento-020117-043348

658    Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C., & Willerslev, E. (2012). Towards next-

659    generation biodiversity assessment using DNA metabarcoding. *Molecular Ecology*, *21*(8), 2045–

660    2050. https://doi.org/10.1111/j.1365-294X.2012.05470.x

661    van Klink, R., Sheard, J. K., Høye, T. T., Roslin, T., do Nascimento, L. A., & Bauer, S. (2024).

662    Towards a toolkit for global insect biodiversity monitoring. In *Philosophical Transactions of the*

663    *Royal Society B: Biological Sciences* (Vol. 379, Issue 1904). Royal Society Publishing.

664    https://doi.org/10.1098/rstb.2023.0101

665    Vaze, S., Han, K., Vedaldi, A., & Zisserman, A. (2021). Open-Set Recognition: a Good Closed-Set

666    Classifier is All You Need? *Arxiv.Org*, 1–27. http://arxiv.org/abs/2110.06207

667    Virgilio, M., Backeljau, T., Nevado, B., & de Meyer, M. (2010). Comparative performances of DNA

668    barcoding across insect orders. In *BMC Bioinformatics* (Vol. 11).

669    http://www.biomedcentral.com/1471-2105/11/206

670    Wang, Z., Liu, X., Liang, D., Wang, Q., Zhang, L., & Zhang, P. (2023). VertU: universal multilocus

671    primer sets for eDNA metabarcoding of vertebrate diversity, evaluated by both artificial and natural

672    cases. *Frontiers in Ecology and Evolution*, *11*(June), 1–18.

673    https://doi.org/10.3389/fevo.2023.1164206

674    Wührl, L., Pylatiuk, C., Giersch, M., Lapp, F., von Rintelen, T., Balke, M., Schmidt, S., Cerretti, P.,

675    & Meier, R. (2022). DiversityScanner: Robotic handling of small invertebrates with machine

676    learning methods. *Molecular Ecology Resources*, *22*(4), 1626–1638. https://doi.org/10.1111/1755-

677    0998.13567

678    Yang, C., Wang, Y., Li, X., Li, J., Yang, B., Orr, M. C., & Zhang, A. (2024). Environmental niche

679    models improve species identification in DNA barcoding. *Methods in Ecology and Evolution*,

680    *15*(12), 2343–2358. https://doi.org/10.1111/2041-210X.14440

681    Zhang, A. B., He, L. J., Crozier, R. H., Muster, C., & Zhu, C.-D. (2010). Estimating sample sizes for

682    DNA barcoding. *Molecular Phylogenetics and Evolution*, *54*(3), 1035–1039.

683    http://www.sciencedirect.com/science/article/B6WNH-4X7GMG7-

684    1/2/7fb4176f978148802690f8769023007e

685    Zhang, J., Yang, J., Wang, P., Wang, H., Lin, Y., Zhang, H., Sun, Y., Du, X., Li, Y., Liu, Z., Chen,

686    Y., & Li, H. (2023). OpenOOD v1.5: Enhanced Benchmark for Out-of-Distribution Detection.

687    *Arxiv.Org*. http://arxiv.org/abs/2306.09301

688    Zheng, W., Yang, L., Genco, R. J., Wactawski-Wende, J., Buck, M., & Sun, Y. (2019). SENSE:

689    Siamese neural network for sequence embedding and alignment-free comparison. *Bioinformatics*,

690    *35*(11), 1820–1828. https://doi.org/10.1093/bioinformatics/bty887

691    Ziemski, M., Wisanwanichthan, T., Bokulich, N. A., & Kaehler, B. D. (2021). Beating Naive Bayes

692    at Taxonomic Classification of 16S rRNA Gene Sequences. *Frontiers in Microbiology*, *12*(June), 1–

693    9. https://doi.org/10.3389/fmicb.2021.644487

694    Zito, A., Rigon, T., & Dunson, D. B. (2023). Inferring taxonomic placement from <scp>DNA</scp>

695    barcoding aiding in discovery of new taxa. *Methods in Ecology and Evolution*, *14*(2), 529–542.

696    https://doi.org/10.1111/2041-210X.14009

697

698

699 **Tables**

700

701 Table 1.

702 Baseline prediction accuracy of three identification algorithms and the perfect classifier under

703 different database sizes, noise levels, and fragment lengths. The best performing method in a

704 set of parameters is indicated by values in boldface.

| | Database Size | Sufficient | | | | Small | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Method | CNN | Distance | BLAST | Perfect | CNN | Distance | BLAST | Perfect |
| Noise Level | Fragment Length | | | | | | | | |
| 0 | 650 | 0.97 | 0.972 | **0.973** | 0.98 | 0.95 | 0.968 | **0.971** | 0.982 |
| | 300 | **0.967** | 0.957 | 0.958 | 0.963 | **0.962** | 0.959 | 0.955 | 0.968 |
| | 150 | **0.96** | 0.945 | 0.946 | 0.951 | **0.958** | 0.945 | 0.942 | 0.952 |
| 0.02 | 650 | 0.953 | 0.965 | **0.968** | | 0.929 | 0.966 | **0.967** | |
| | 300 | **0.963** | 0.937 | 0.953 | | **0.956** | 0.944 | 0.95 | |
| | 150 | **0.949** | 0.923 | 0.933 | | **0.945** | 0.93 | 0.931 | |

705
706

707

708

709

710

711

712    Table 2.

713    False negative rates at a 95% threshold for three OOD detection methods and a perfect

714    classifier under various parameter settings. A lower FNR indicates better performance.  the

715    best performing method in a set of parameters is indicated by values in boldface.  For the CNN

716    model, the results of majority voting (MV) with multiple methods are shown.

| | Database Size | Sufficient | | | | Small | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Method | CNN(MV) | Distance | BLAST | Perfect | CNN(MV) | Distance | BLAST | Perfect |
| Noise Level | Fragment Length | | | | | | | | |
| 0 | 650 | 0.128 | 0.105 | **0.101** | 0.061 | 0.13 | **0.099** | 0.101 | 0.046 |
| | 300 | **0.126** | 0.14 | 0.142 | 0.116 | 0.137 | **0.125** | 0.133 | 0.09 |
| | 150 | **0.156** | 0.176 | 0.17 | 0.161 | **0.156** | 0.165 | 0.162 | 0.132 |
| 0.02 | 650 | 0.171 | 0.147 | **0.142** | | 0.175 | **0.136** | **0.136** | |
| | 300 | **0.186** | 0.205 | 0.198 | | 0.193 | **0.181** | 0.183 | |
| | 150 | **0.242** | 0.272 | 0.27 | | **0.243** | 0.244 | 0.256 | |

717
718

719

720

721

722

723

724

725     Table 3.

726     Regression coefficients estimated by multivariate regression modeling for baseline accuracy

727     and false negative rate. Coefficients dropped by the LASSO variable selection are indicated by

728     "." signs. *Dw*: Within-species genetic distance. *Dbt* : Between-species genetic distance

| | Explanatory | variables | | | | |
|---|---|---|---|---|---|---|
| Response | No. classes | Average *Dw* | Average *Dbt* | Minimum *Dbt* | completeness | No. samples per species |
| Baseline Accuracy | -1.35e-05 | . | . | 0.040 | -0.051 | 6.39e-04 |
| False negative rate | -0.0007 | 1.031 | . | -1.16 | . | . |

729

730

**Figures**

Figure 1. A schematic diagram of the classification model, data acquisition and analysis

procedures.

747    Figure 2.

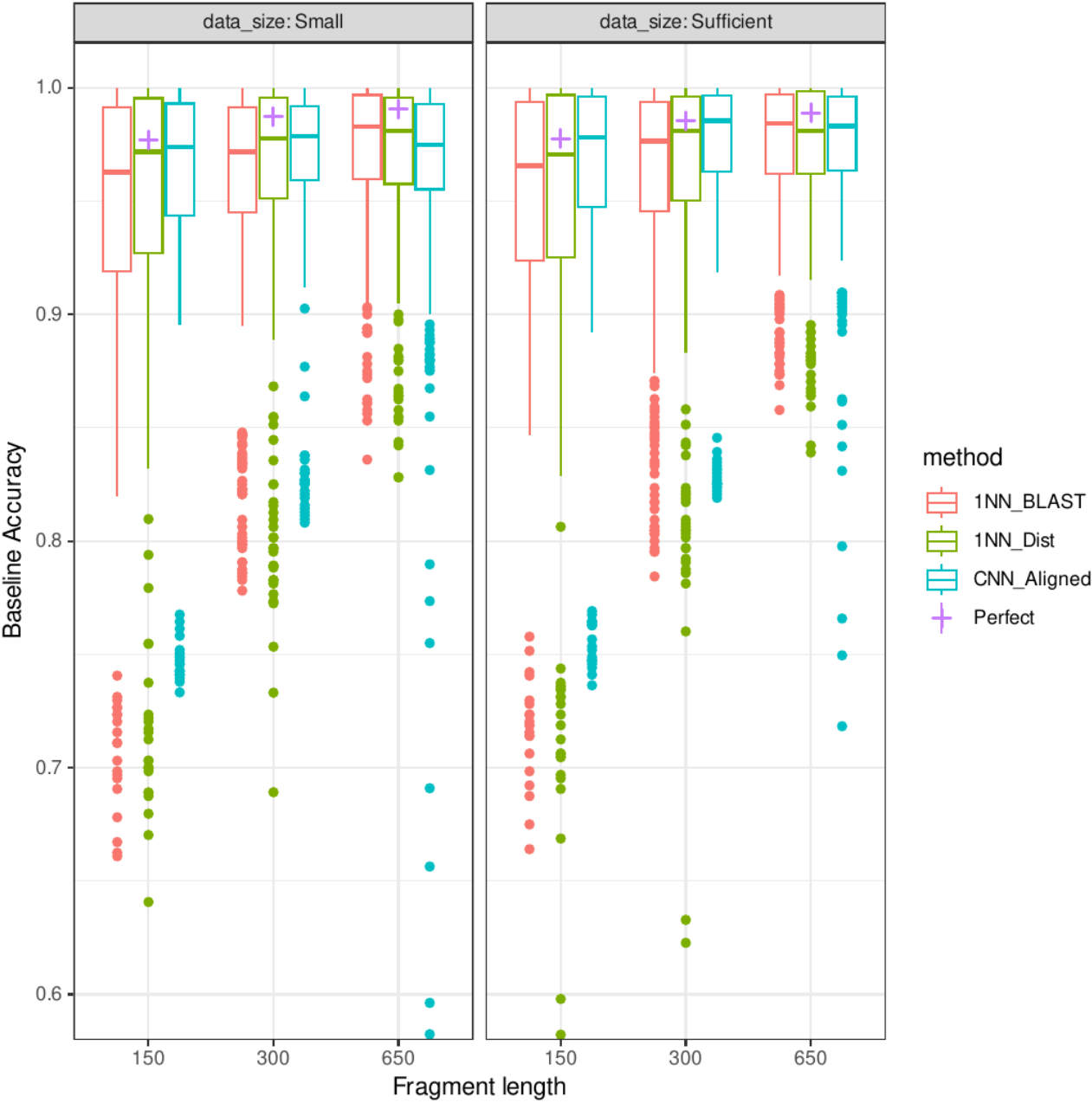748    Effects of fragment lengths and database sizes on baseline prediction accuracy in the noiseless

749    dataset.



750

751

752

753

754

755 Figure 3.

756 Effects of fragment lengths and database sizes on false negative rates of OOD detection in the
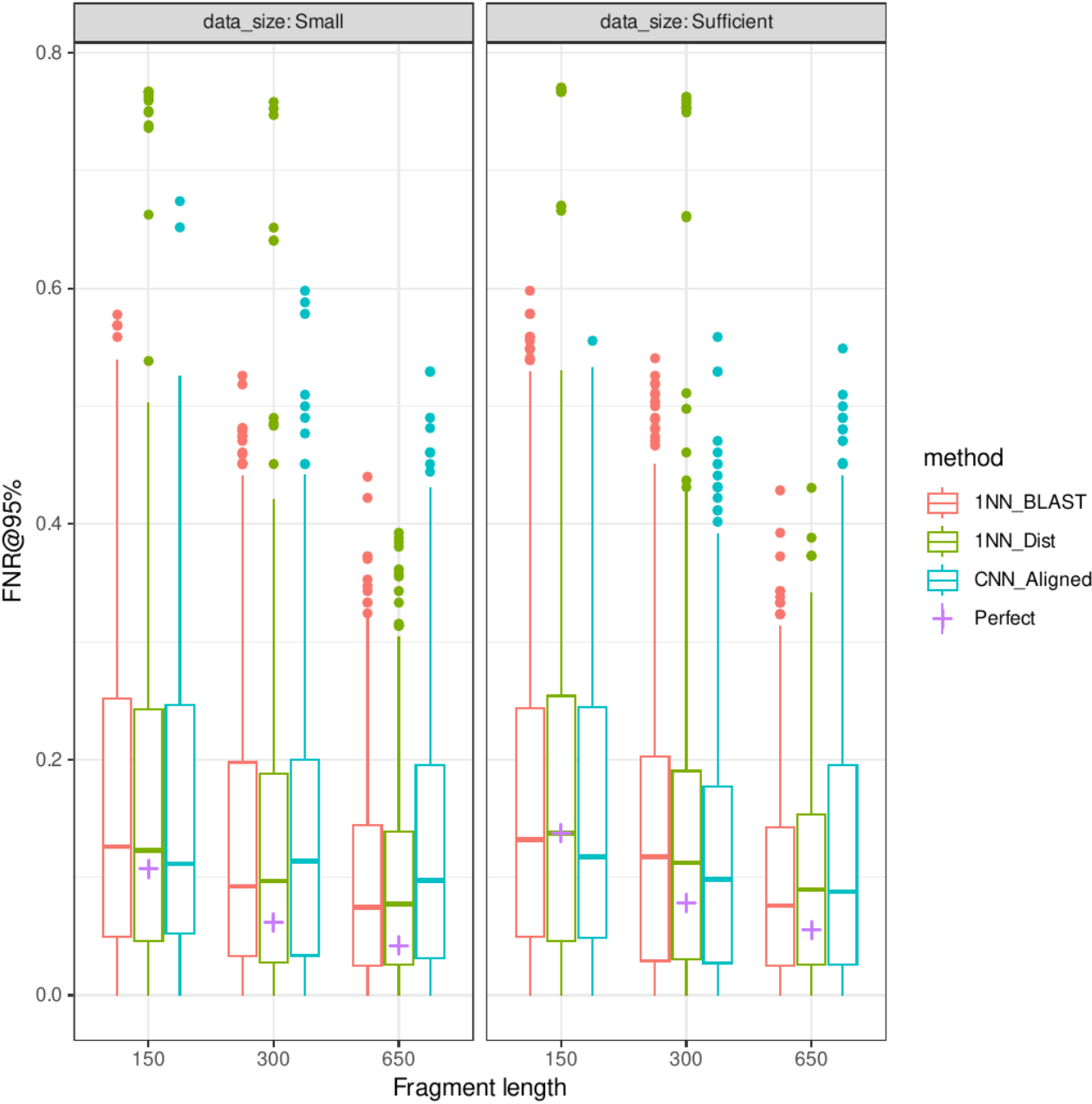
757 noiseless dataset.



758

759

760

761

762

763  Figure 4.

764  Relationship between the accuracy of the CNN classifier and the false negative rate of the CNN
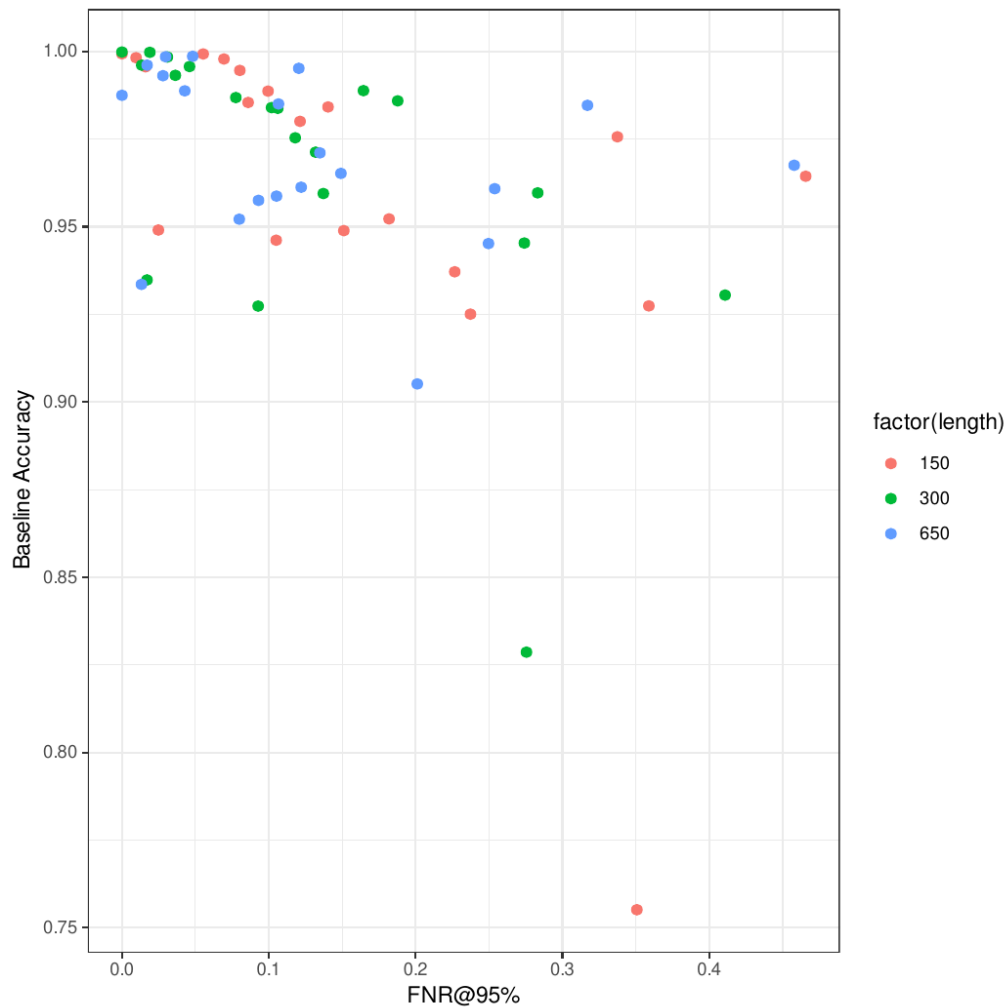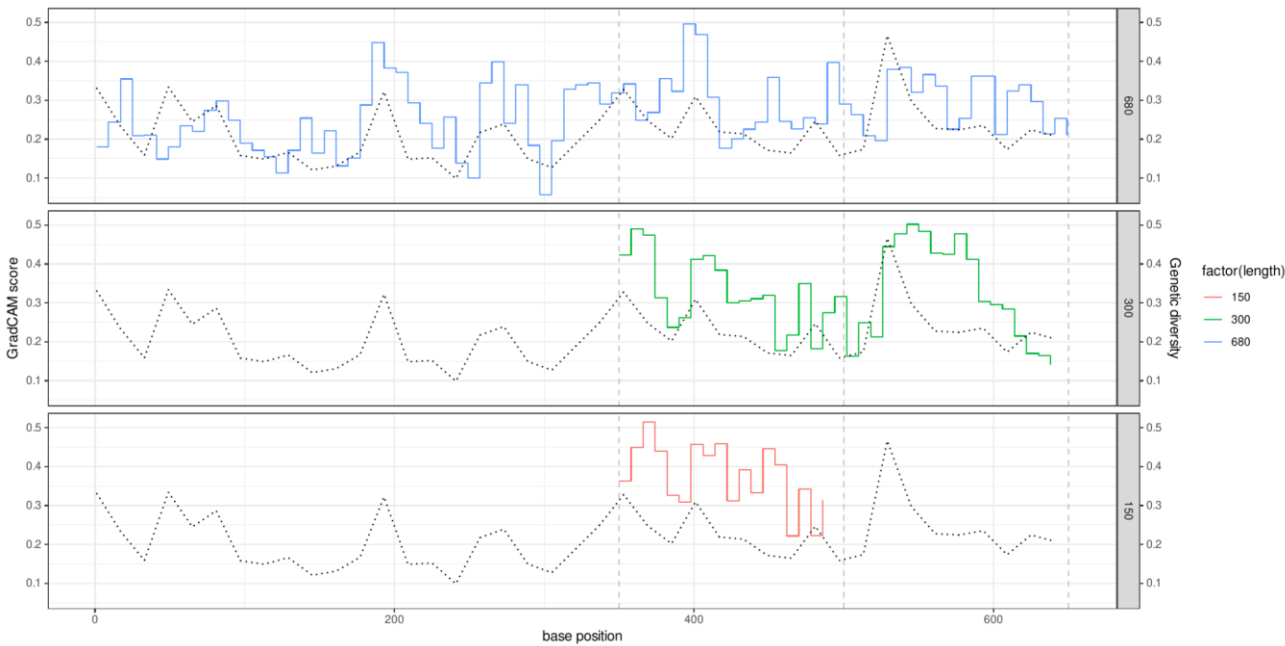
765  OOD detector.



766

767

768

769

770

771

772

773

774

775   Figure 5.

776   Spatial distribution of the average GradCAM score for in-distribution samples of the

777   *Cryptocephalus* (leaf beetle) dataset. Solid step lines represent the GradCAM score for the 8-bp

778   windows, and black dotted lines represent genetic variations in the 16-bp windows.



779

780

781

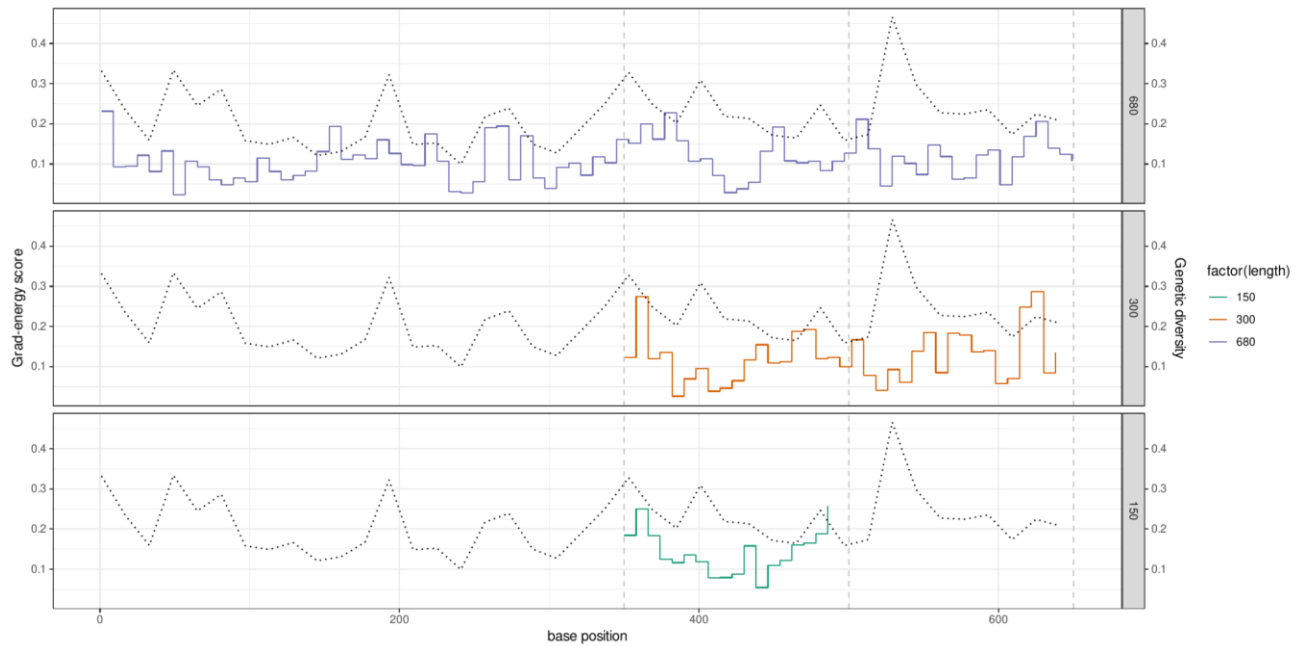782

783

784

785

786

787

788

789

790

791

792    Figure 6.

793    Spatial distribution of the average grad-energy scores for OOD samples from the

794    *Cryptocephalus* data set. Solid step lines represent the grad-energy score for 8-bp windows,

795    and black dotted lines represent genetic variations in 16-bp windows.
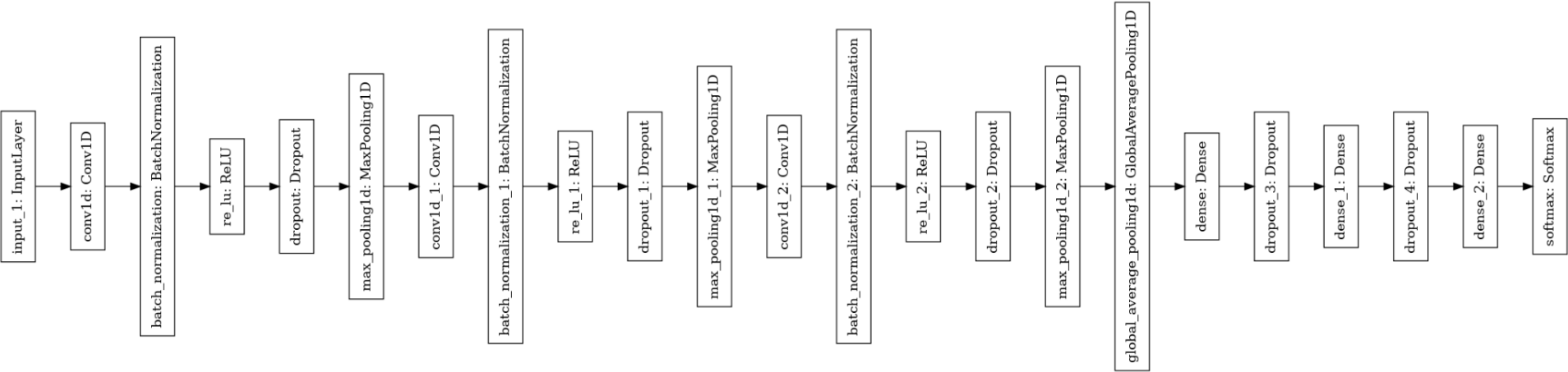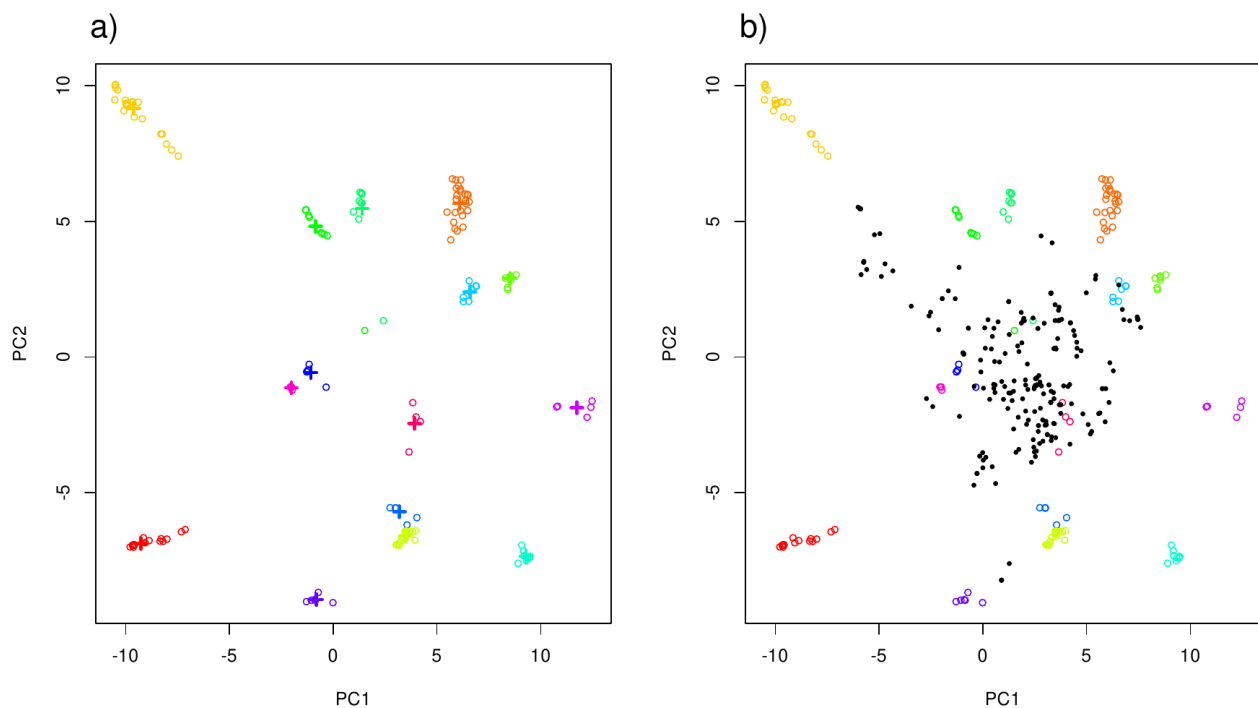


796

797
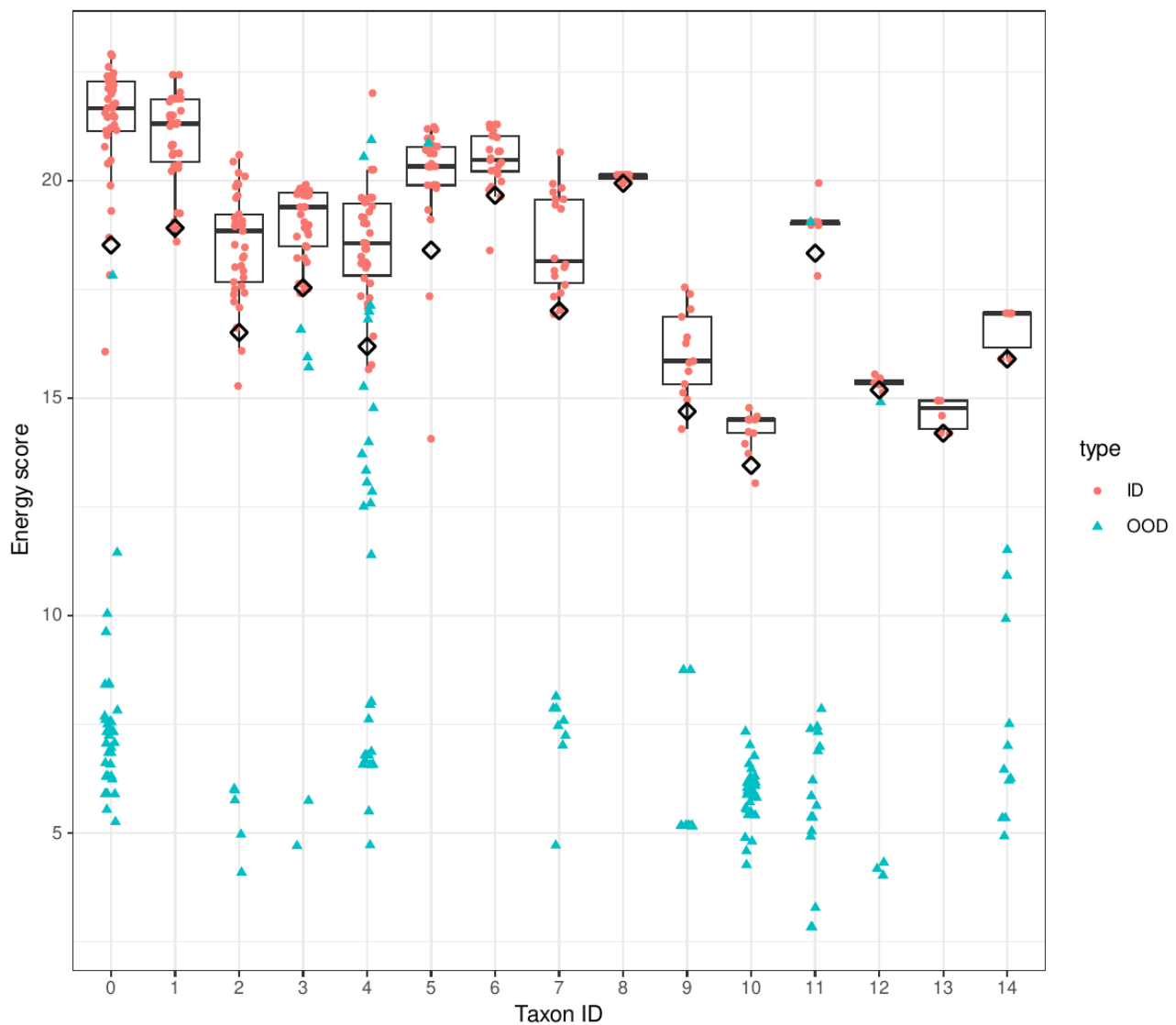
798  **Supplementary figures**

799  Supplementary figure S1. A diagram showing the detailed architecture of the CNN model
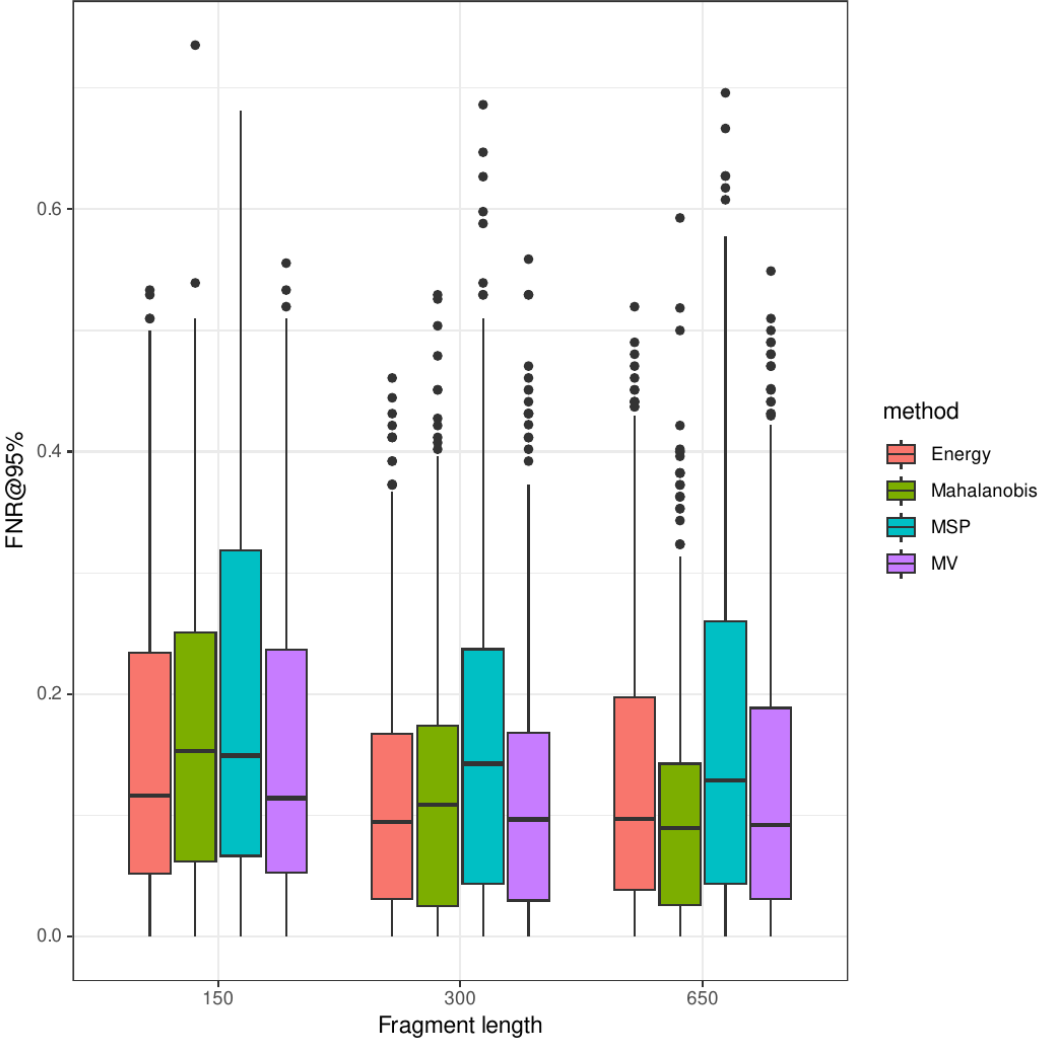


800

801

802

803

804

805

806   Supplementary figure S2. Exemplary distribution of *g(x)*, showing the outputs of the

807   penultimate FC layer. PCA was applied to reduce the dimensionality for visualization. a) Dots in

808   colors represent in-distribution (ID) samples of different species, while crosses in

809   corresponding colors are class centers; $\mu_k$. b) the same plots with OOD samples are shown in

810   black dots. ID samples were frequently clustered in linearly separable groups in the

811   intermediate output space, while OOD samples were placed between such groups. Hence,

812   distances from class centers to samples can be used to measure the OOD status of samples.

813



814

815

816

817

818

819

820

821    Supplementary figure S3. Distribution of energy scores for ID and OOD samples from the

822    Drosophila dataset. Taxon IDs of the OOD samples were assigned based on the predictions of

823    the CNN classifier. Open squares indicate the 95% quantiles of the energy scores of ID samples.

824    Samples with lower energy scores with these thresholds were detected as OODs. OOD samples

825    missed by these procedures, such as those with high energy scores in Taxon ID 4, were

826    considered false negatives.



827

828

829

830

831    Supplementary figure S4. False negative rates (FNR@95%) of four OOD detection methods and

832    their relationships with fragment lengths. Results on the noiseless sufficient-size dataset are

833    shown. MSP: Maximum Softmax Probability, MV: Majority Voting.



834

835

836

837

838     Supplementary table 1. Dataset summary.

| Genus | Dataset Code | Order | Common Name | No.ID.samples | No.ID-species | No.OOD samples |
|---|---|---|---|---|---|---|
| *Drosophila* | dro15 | Diptera | fruit fly | 1080 | 15 | 162 |
| *Megaselia* | meg42 | Diptera | scuttle fly | 2296 | 42 | 348 |
| *Aedes* | aed19 | Diptera | tiger mosquito | 1169 | 19 | 102 |
| *Atheta* | ath18 | Coleoptera | rove beetle | 707 | 18 | 228 |
| *Cryptocephalus* | cry22 | Coleoptera | leaf beetle | 828 | 22 | 304 |
| *Pterostichus* | pte45 | Coleoptera | ground beetle | 2383 | 44 | 169 |
| *Dolerus* | dol24 | Hymenoptera | sawfly | 907 | 24 | 135 |
| *Megachile* | mch15 | Hymenoptera | leafcutter bee | 649 | 15 | 240 |
| *Lassioglossum* | las68 | Hymenoptera | sweat bee | 2960 | 68 | 910 |
| *Euxoa* | eux40 | Lepidoptera | owlet moth | 1623 | 40 | 778 |
| *Phyllonorycter* | phy53 | Lepidoptera | leaf mining moth | 2086 | 53 | 439 |
| *Acleris* | acl32 | Lepidoptera | leaf roller moth | 1645 | 32 | 138 |
| *Culicoides* | cul26 | Diptera | biting midge | 992 | 26 | 248 |
| *Amara* | ama25 | Coleoptera | sun beetle | 1142 | 25 | 260 |
| *Catocala* | cat62 | Lepidoptera | underwing moth | 2135 | 62 | 355 |
| *Andrena* | and28 | Hymenoptera | mining bee | 873 | 28 | 565 |
| *Bombus* | bom24 | Hymenoptera | bumble bee | 1187 | 24 | 167 |
| *Corynoptera* | cor19 | Diptera | fungus gnat | 1927 | 19 | 23 |
| *Bembidion* | bem39 | Coleoptera | ground beetle | 1053 | 39 | 226 |
| *Caloptilia* | cal17 | Lepidoptera | leaf mining moth | 780 | 17 | 189 |

839

840    Supplementary table 2

841    False negative rates at the 95% threshold for the four OOD detection methods of the deep

842    learning model. The best performing methods are indicated in boldface. The results for a

843    noiseless, sufficiently sized dataset are shown.  MSP: Maximum Softmax Probability, MV:

844    Majority Voting.

845

| | Database | Sufficient | | | |
|---|---|---|---|---|---|
| | Method | MSP | Energy | Mahalanobis | MV |
| Noise level | Fragment length | | | | |
| 0.0 | 650 | 0.169 | 0.13 | **0.11** | 0.128 |
| | 300 | 0.171 | **0.124** | **0.124** | 0.126 |
| | 150 | 0.202 | 0.157 | 0.167 | **0.156** |

846

847