

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

# A universal DNA signature for the Tree of Life

Bruno A. S. de Medeiros<sup>1,2,3</sup>, Liming Cai<sup>4,5</sup>, Peter J. Flynn<sup>4</sup>, Yujing Yan<sup>4</sup>, Xiaoshan Duan<sup>4,6</sup>,  
Lucas C. Marinho<sup>4,7</sup>, Christiane Anderson<sup>8</sup>, and Charles C. Davis<sup>4</sup>

<sup>1</sup>Field Museum of Natural History, Chicago, Illinois, 60605, USA

<sup>2</sup>Department of Organismic and Evolutionary Biology, Museum of Comparative Zoology,  
Harvard University, Cambridge, Massachusetts, 02138 USA

<sup>3</sup>Smithsonian Tropical Research Institute, Panama City, Panama

<sup>4</sup>Department of Organismic and Evolutionary Biology, Harvard University Herbaria,  
Harvard University, Cambridge, Massachusetts, 02138 USA

<sup>5</sup>Department of Integrative Biology, The University of Texas at Austin, Austin, Texas, 78712  
USA

<sup>6</sup>College of Forestry, Northwest Agriculture & Forestry University, Yangling 712100,  
Shaanxi, China

<sup>7</sup>Departamento de Biologia, Universidade Federal do Maranhão, Av. dos Portugueses 1966,  
Bacanga 65080-805, São Luís, Maranhão, Brazil

<sup>8</sup>University of Michigan Herbarium, 3600 Varsity Drive, Ann Arbor, Michigan 48108, USA

**Corresponding authors:**

Bruno A. S. de Medeiros, Field Museum of Natural History, Chicago, IL, 60605; E-mail:  
[bdemedeiros@fieldmuseum.org](mailto:bdemedeiros@fieldmuseum.org)

27 Charles C. Davis, Department of Organismic and Evolutionary Biology, Harvard University  
28 Herbaria, Cambridge, MA 02138, USA; E-mail: [cdavis@oeb.harvard.edu](mailto:cdavis@oeb.harvard.edu)

29

30

# 31 Abstract

32 Species identification using DNA barcodes has revolutionized biodiversity sciences and  
33 society at large. However, conventional barcoding methods may lack power and universal  
34 applicability across the Tree of Life. Alternative methods based on whole genome  
35 sequencing are hard to scale due to large data requirements. Here, we develop a novel  
36 DNA-based identification method, varKoding, using exceptionally low-coverage genome  
37 skim data to create two-dimensional images representing the genomic signature of a  
38 species. Using these representations, we train neural networks for taxonomic  
39 identification. Applying a taxonomically verified novel genomic dataset of Malpighiales  
40 plant accessions, we optimize training hyperparameters and find the highest performance  
41 by combining a transformer architecture with a new modified chaos game representation.  
42 Remarkably, >91% precision is achieved despite minimal input data, exceeding alternative  
43 methods tested. We illustrate the broad utility of varKoding across several focal clades of  
44 eukaryotes and prokaryotes. We also train a model capable of identifying all species in  
45 NCBI SRA using less than 10 Mbp sequencing data with 96% precision and 95% recall and  
46 robust to sequencing platforms. Enhanced computational efficiency and scalability,  
47 minimal data inputs robust to sequencing details, and modularity for further development  
48 make varKoding an ideal approach for biodiversity science.

49  
50 **Keywords:** biodiversity science, computer vision, DNA barcoding, DNA signature,  
51 Malpighiaceae, natural history collections, neural networks, species identification,  
52 taxonomy

53

## 54 **Introduction**

55 For two decades, conventional DNA barcoding, which relies on standardized short  
56 sequences (400–800 bp) for species identification<sup>1–5</sup>, has enabled novel and massively  
57 scalable science spanning evolution<sup>4,6–9</sup>; ecology<sup>10–14</sup> and paleontology<sup>15–19</sup>. Practical  
58 applications of barcoding have also made major contributions to environmental health,  
59 including the ability to authenticate medicinal plants<sup>20</sup>, detect agricultural pests<sup>21</sup>, and  
60 monitor poaching and the trade of endangered species<sup>22–27</sup>. Despite these remarkable  
61 achievements, conventional DNA barcoding suffers from at least four limitations. First,  
62 barcodes are customized specifically for a taxon (e.g., plants, animals, and fungi), and  
63 therefore are not universal. For example, commonly used plant barcodes from chloroplast  
64 genes such as *matK* and *rbcL* cannot be applied as barcodes for all plants<sup>28,29</sup>, or for animals  
65 and fungi. Second, conventional barcode loci may fail to distinguish closely related taxa, a  
66 pervasive shortcoming in plants<sup>2,30</sup>. Third, reliance on a single locus may lead to spurious  
67 results in the case of complex evolutionary scenarios such as hybridization in deep or  
68 shallow time<sup>31–34</sup>. And fourth, the necessary comparison of homologous genes may fail  
69 when PCR primers are not universal<sup>35</sup>, the source DNA is fragmented<sup>27</sup>, or paralogy and the  
70 presence of pseudogenes confounds accurate orthology assessments<sup>36,37</sup>.

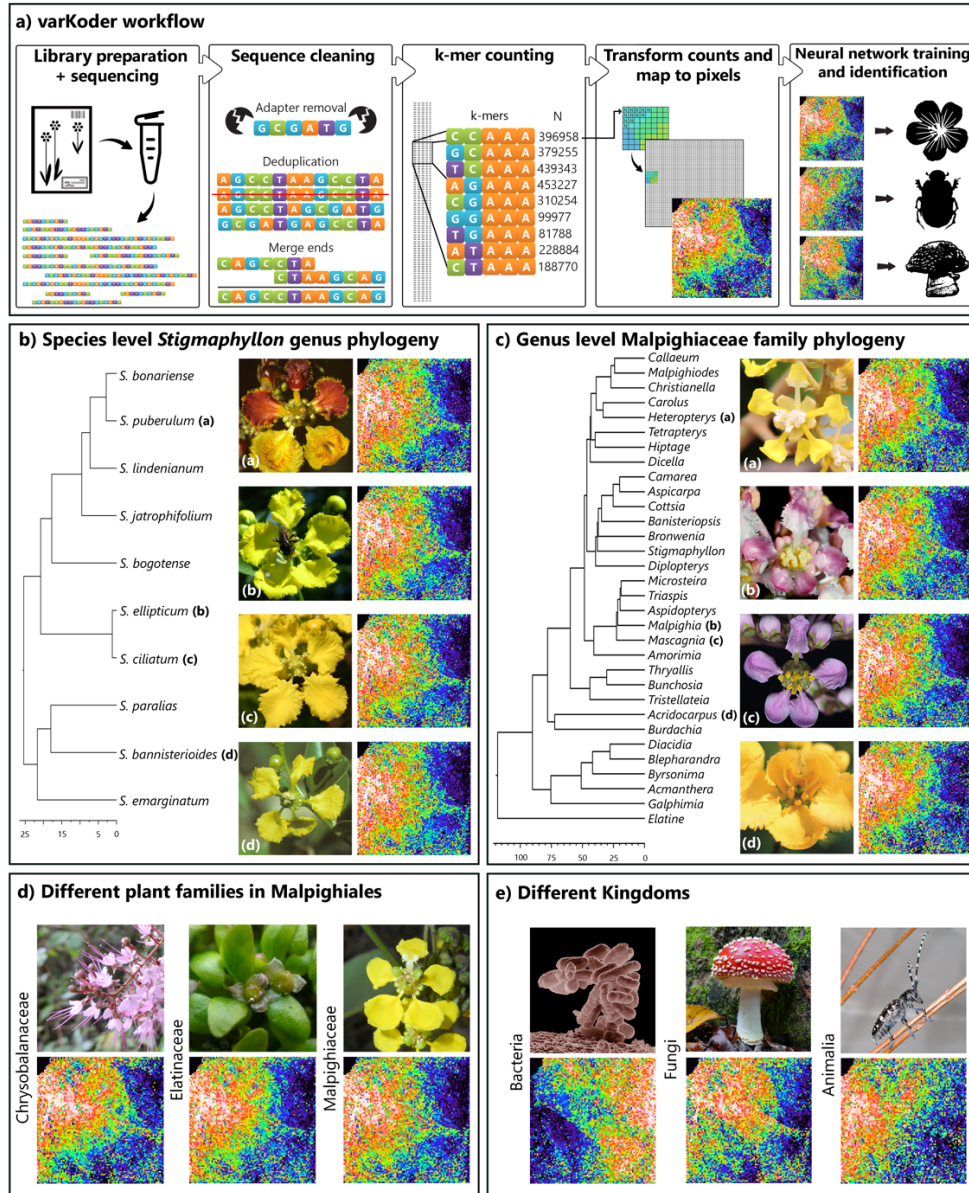
71

72 Newer alternatives to conventional barcoding have begun to address these challenges by  
73 leveraging high-throughput sequencing and machine-learning powered by deep neural  
74 networks. High-throughput sequencing facilitates more comprehensive assessments of  
75 total genomic space<sup>38,39</sup>. For example, presence and absence patterns among short DNA  
76 sequences (k-mers) from low-coverage reads (i.e., genome skims) can estimate overall  
77 sequence distances, bypassing genome alignments entirely as implemented in *Skmer*<sup>40</sup>.  
78 Machine learning enables more complex sequence comparisons than conventional methods  
79 that rely on homology and simple metrics<sup>41</sup>. Machine-learning models can cluster DNA  
80 sequences without supervision<sup>42,43</sup> or classify sequences based on reference datasets<sup>44–49</sup>.

81 In particular, neural networks are exceptionally powerful for sophisticated computer-  
82 vision tasks, such as image classification<sup>50</sup>. Thus, the combination of low-coverage genome  
83 skimming data and neural networks holds enormous promise for accurate and scalable  
84 DNA barcoding, but its potential has yet to be fully realized<sup>39</sup>.

85  
86 Genomes differ substantially in many features beyond the simple nucleotide divergence  
87 commonly used in conventional barcoding, but these genomic features have been  
88 overlooked in species identification<sup>51-55</sup>. We propose that (1) relevant genomic features  
89 can be captured by nucleotide composition with short k-mer counts and very small  
90 sequence coverage; and (2) these counts can be used to distinguish species and higher taxa  
91 efficiently and accurately using machine learning. Prior work on k-mer-based  
92 representations of genome composition (i. e. DNA signatures) has shown high accuracy can  
93 be achieved with high-coverage data or a large number of replicates per taxon, particularly  
94 for identification at higher taxonomic ranks<sup>42-47,56-63</sup>. However, given the millions of  
95 existing species and the sparse genetic data available, a practical scalable method would  
96 require: (1) consistently high accuracy despite limited evolutionary divergence; (2) fast  
97 computations; and (3) high accuracy with small training datasets (both in number of  
98 samples and DNA data per sample). Here we developed a novel DNA signature method,  
99 which we call **varKoding**, that integrates very low-coverage genome skim data with  
100 optimized training of machine-learning models using two-dimensional images  
101 representing genome composition (**Figure 1A**). To develop and optimize varKoding for  
102 accurate species identification, we generated a *de novo* genome skim dataset including  
103 hundreds of samples derived primarily from historical herbarium specimens for the  
104 diverse plant genus *Stigmaphyllon* (Malpighiaceae), which has received extensive  
105 phylogenetic and taxonomic treatment<sup>64-68</sup>. Next, we explored the utility of varKoding and  
106 compared it to alternatives at different phylogenetic depths from families to species within  
107 the flowering plant order Malpighiales (Malpighiaceae, Chrysobalanaceae, and  
108 Elatinaceae). Finally, we demonstrate the scalability of varKoding and its potential  
109 application in forensics and related fields by testing it on (1) species-level datasets from

110 fungi, plants, animals, and bacteria; (2) massive datasets retrieved from the NCBI sequence  
111 read archive (SRA); and (3) a previously published environmental DNA (eDNA) dataset.



112

113 **Figure 1.** Overview of varKoding. (A) Image generation workflow, depicting varKodes. Images are natively  
 114 grayscale, but here they are mapped to a rainbow color scale for increased contrast. (B) Phylogeny and  
 115 example varKodes of *Stigmaphyllon* species. (C) Phylogeny and example varKodes of Malpighiaceae genera  
 116 including their closest outgroup (*Elatine*, Elatinaceae). Time trees in 1B and 1C (D) Examples of varKodes  
 117 from across plant families of Malpighiales, and (E) across kingdoms. Chronograms depicted for each  
 118 representative set with timelines in millions of years (Myr) at the bottom of B and C. These were derived  
 119 from an ongoing family-wide phylogenomic investigation of the family Malpighiaceae (C. C. Davis personal  
 120 communication) using methods and fossil constraints described in Cai et al.<sup>65</sup>; dates inferred are consistent  
 121 with earlier findings<sup>65</sup> and were not applied in this study for quantitative analyses.

## 122 Results and Discussion

### 123 DNA signature images can be classified with generalized neural networks

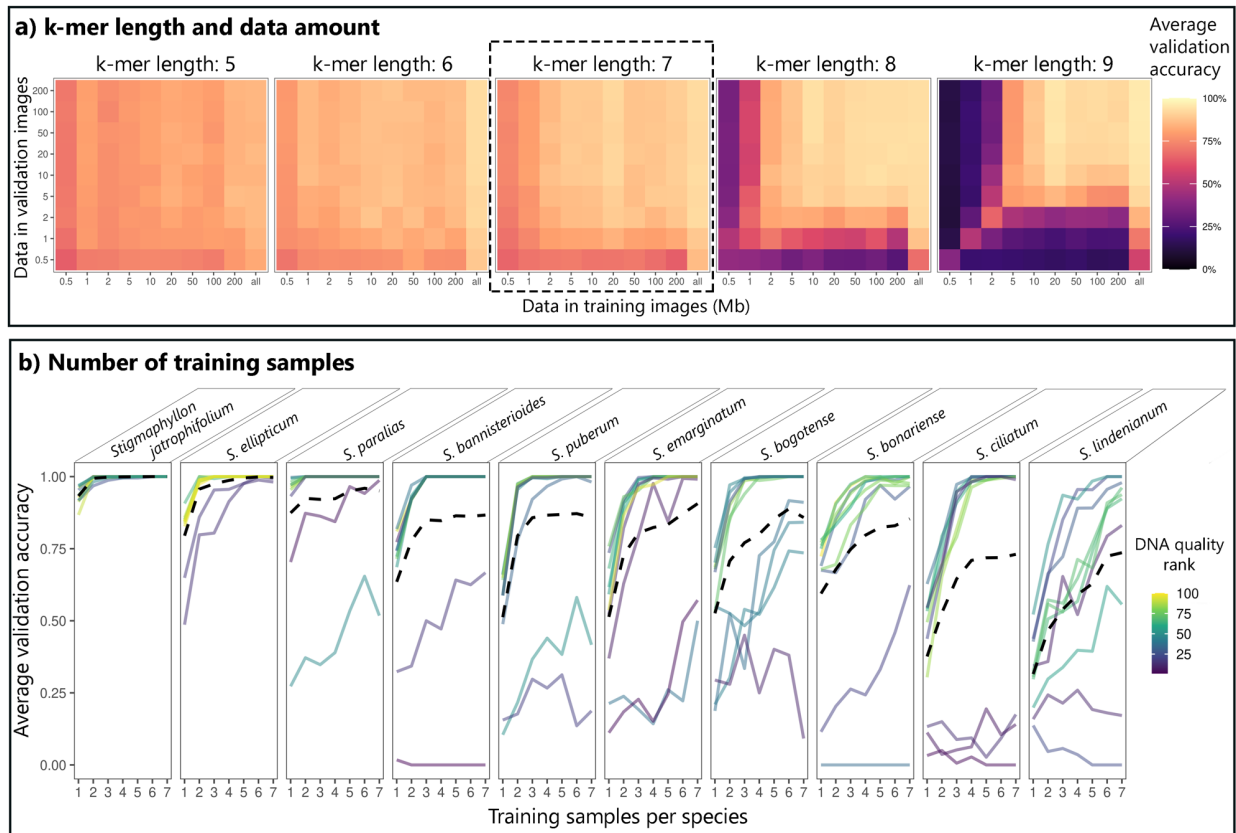
124 We first generated a novel kind of image representation of a DNA signature, which we  
125 termed a **varKode**. varKodes map k-mers onto pixels of a 2-D image based on their  
126 similarity and represent ranked k-mer frequencies as pixel brightness. Variation in  
127 varKodes can be small but remain visually perceptible among species (**Figure 1B**) and  
128 genera (**Figure 1C**). Variation is more striking among higher levels of phylogenetic  
129 divergence, such as between families in the order Malpighiales (**Figure 1D**) or different  
130 kingdoms of eukaryotes and prokaryotes (**Figure 1E**). We expected, therefore, that neural  
131 network architectures developed for image classification, (e.g., deep residual networks,  
132 resnets<sup>69</sup> or vision transformers, ViT<sup>70,71</sup>) would be able to differentiate varKodes.

133

134 We first optimized hyperparameters and training conditions to maximize accuracy for  
135 species-level identification of *Stigmaphyllon*. We identified that varKodes depicting k-mer  
136 length = 7 struck a good balance between accuracy and the amount of input sequence data  
137 (**Figure 2A**). Furthermore, models trained with augmented data from several subsampled  
138 sequences drawn from each individual exhibited substantially better performance (**Figure**  
139 **2A**). A linear model demonstrated that neural network architectures and training methods  
140 designed for image classification of photographs<sup>69,72-75</sup> are extremely useful for varKode-  
141 based identification. Specifically, we observed increased accuracy with more parameter-  
142 rich neural network architectures (*ResNeXt101*<sup>76</sup>, among those tested), augmentation with  
143 lighting transformations, *CutMix*<sup>75</sup> and *MixUp*<sup>74</sup>. Label smoothing<sup>77</sup> and pretraining models  
144 on generalized photographs decreased accuracy (**Figure 3**). Contrary to the widely held  
145 idea that deep neural networks require very large training datasets<sup>60,78</sup>, the  
146 aforementioned approaches enabled training with very modest data amounts: four  
147 biological replicates per taxon was sufficient for 100% median accuracy (**Figure 2B**).  
148 Errors in species-level identification were concentrated among sequences derived from  
149 herbarium samples that demonstrated evidence of DNA damage, as is sometimes reported  
150 for ancient DNA<sup>79</sup> (**Figure 2B**). However, including low-quality training samples slightly



151 decreased mean validation accuracy—from 73% to 71%—for low-quality validation  
 152 samples, but had no effect on high-quality validation samples (89–90% mean accuracy,  
 153 **Figure 4A**).

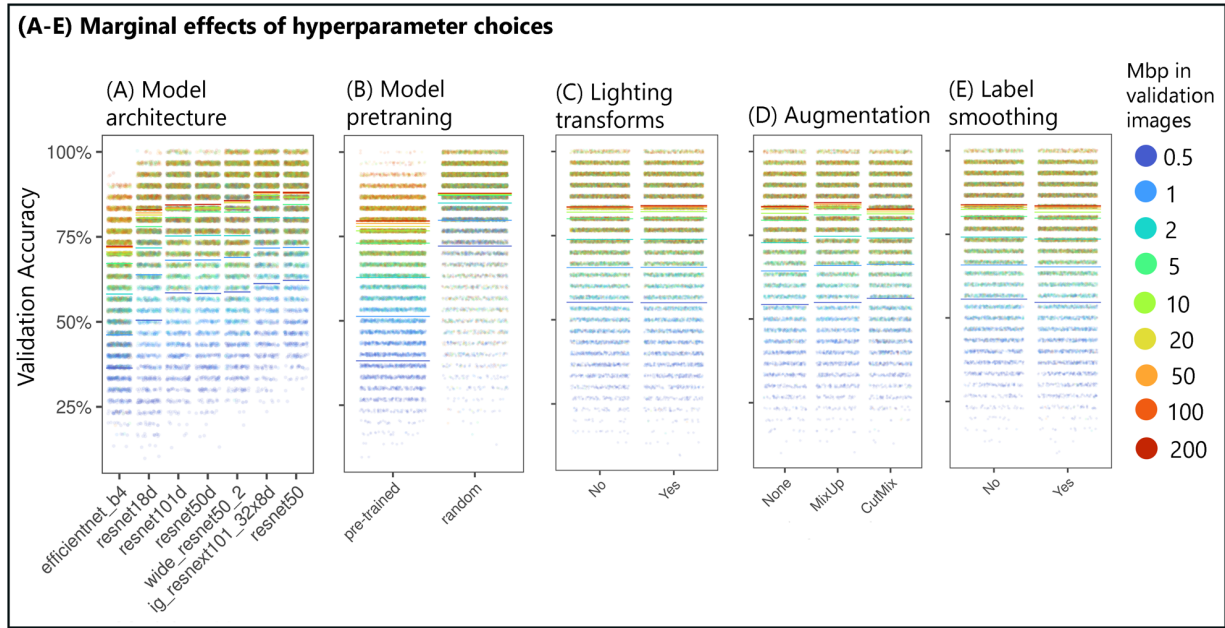


154

155 **Figure 2. Neural network training of varKodes for species identification. (A)** Effect of k-mer length and  
 156 input data amount used to produce varKodes on validation accuracy. Longer k-mers increase accuracy when  
 157 more data are used. Mixing varKodes subsampled from different amounts of data improves accuracy. Box  
 158 with dashed line (k-mer length = 7) strikes a good balance between model accuracy and amount of required  
 159 data. **(B)** Validation accuracy improves with increased number of training samples per species, but even 3–4  
 160 samples are sufficient in most cases for achieving high accuracy. Each solid line represents one sample,  
 161 colored by DNA quality (i.e., variation in base pair frequencies). Higher rank indicates better quality. Dashed  
 162 lines represent averages across all samples.

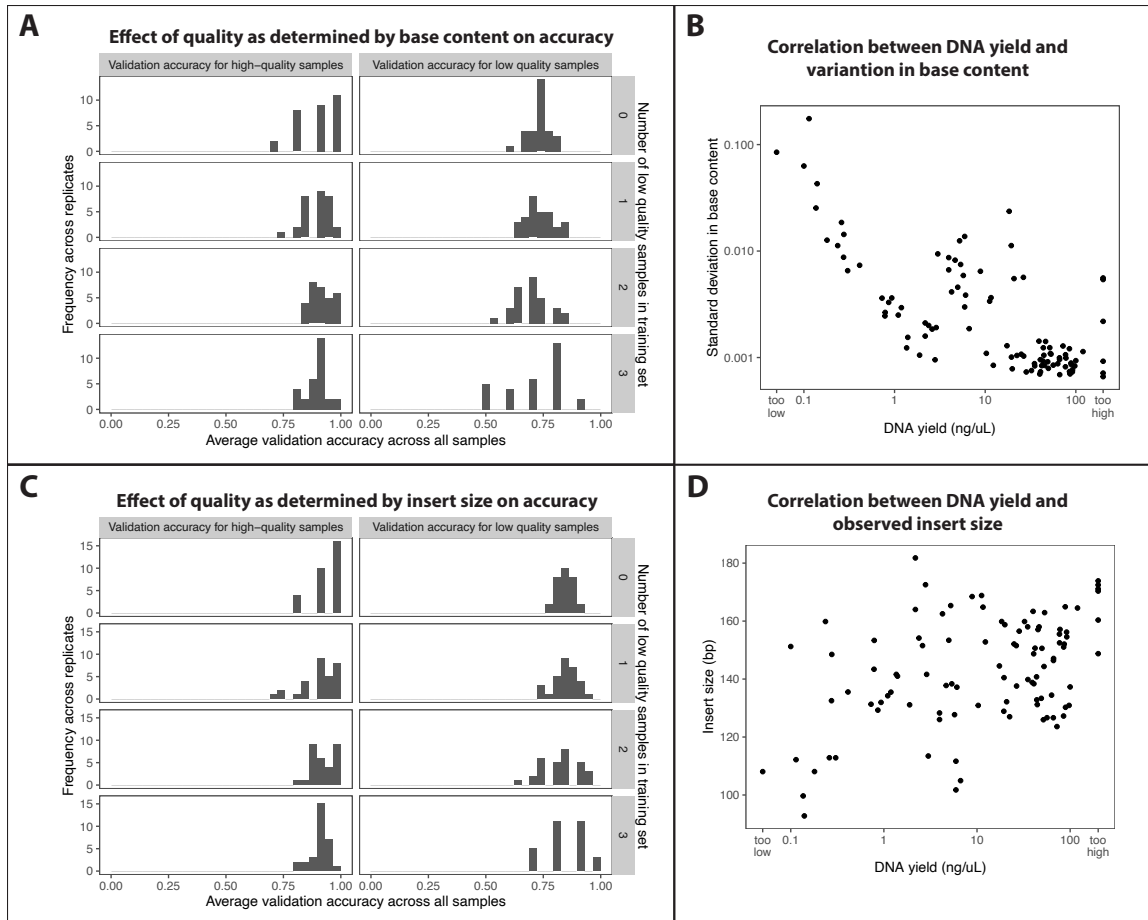
163

164



165  
166  
167  
168  
169

**Figure 3.** Marginal effects of neural network model and training options. Dots represent individual replicates, and bars depict averages. All parameters were identified to be significant in a linear model: more complex model architectures, lighting transformations, and augmentation methods *MixUp* and *CutMix* improved accuracy. However, pretraining with large image datasets and label smoothing decreased accuracy.



170

171

172

173

174

175

176

177

178

179

180

181

182

183

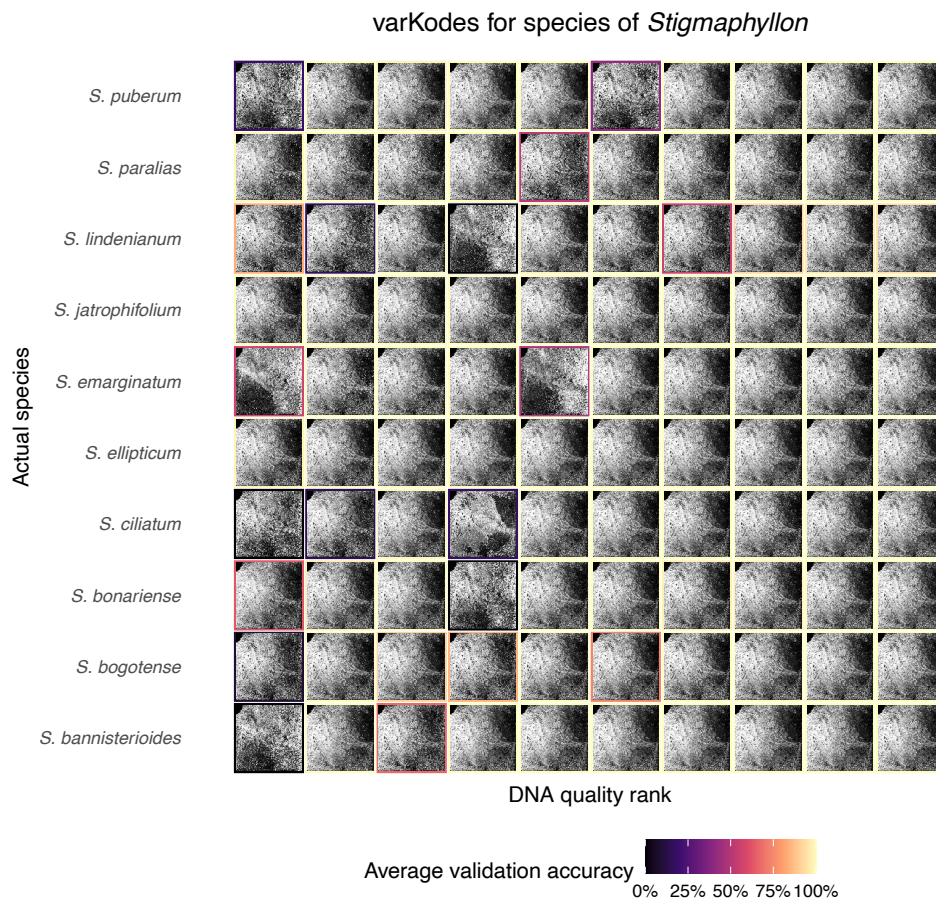
184

185

**Figure 4.** Effect of the inclusion of low-quality training samples, inferred from variation in base pair content (A, B) or insert size (C, D). Increasing the fraction of samples in the training set that were low-quality did not strongly affect the average validation accuracy, but it increased dispersion. Low-quality samples are the four samples with highest variation in base-pair content or shortest insert size in raw reads for each species. Panels B and D show the correlation of each quality metric with DNA extraction yield.

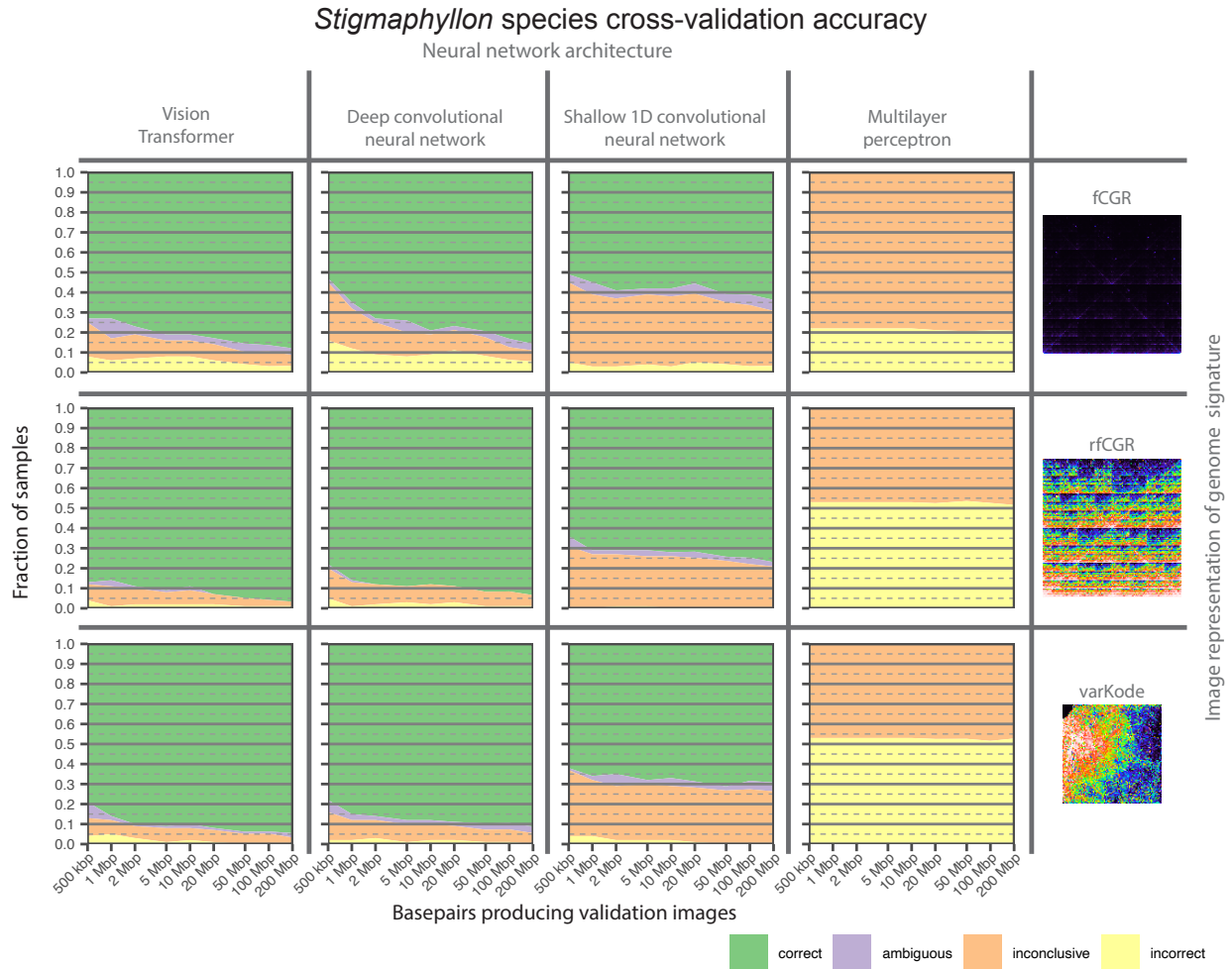
We hypothesized that lower-quality samples shared similar sequences resulting from common patterns of DNA damage and greater levels of microbial or human contaminants, resulting in spurious similarities in varKodes (Figure 5). Contaminants also are thought to increase errors in other genome skim methods<sup>80</sup>. To mitigate this problem, we applied multi-label classification<sup>81</sup> to our neural network models. Although single-label classification models always return a single prediction (that is, an inferred label), multi-label models can return zero or more predictions, avoiding spurious results when there is uncertainty. For a set of samples with known labels used for validation, a prediction is a true positive if the predicted label matches the actual label, and a false positive if not.

186 Failure to predict an actual label is deemed a false negative. For each validation sample, we  
 187 summarized predictions as (1) correct (true positives only); (2) incorrect (false positives  
 188 only); (3) ambiguous (multiple predictions, including true and false positives); or (4)  
 189 inconclusive (i. e. no prediction above the confidence threshold of 0.7). For each test, we  
 190 summarized results across all validation samples using two metrics: precision (the sum of  
 191 all true positives divided by the sum of all true and false positives) and recall (the sum of all  
 192 true positives divided by the sum of all true positives and negatives).



193  
 194 **Figure 5.** Low-quality DNA may lead to spurious patterns of similarity in varKodes. Samples with lower  
 195 quality show varKode patterns divergent from their species more often than high-quality ones. These  
 196 divergent patterns may be similar between low-quality samples across species. These samples also show  
 197 reduced validation accuracy in a single-label model. For each sample, we show the varKodes produced from  
 198 all DNA data available. Within each species, samples are organized from lowest (left) to highest (right) DNA  
 199 quality. Bounding boxes around each sample indicate the average validation accuracy across 30 random  
 200 replicates with 7 training samples per species.

201  
202 After optimizing these training conditions, we directly compared varKodes to an existing  
203 method of DNA signature representation: the frequency chaos game representation  
204 (*fCGR*)<sup>56,59</sup>. In *fCGR*s, k-mers are mapped to pixels based on their oriented sequence and  
205 pixel brightness represents the rescaled k-mer frequency. To isolate the effects of pixel  
206 mapping and brightness, we created a new representation combining *fCGR* mapping with  
207 *varKode* ranked frequency transformation (*rfCGR*). By directly comparing these 3 kinds of  
208 representation combined with four neural network architectures, including (1) two  
209 previously employed with *fCGR*s<sup>42,44,60</sup>, (2) the optimal architecture in our initial tests  
210 (ResNeXt101<sup>76</sup>), and (3) a Vision Transformer (ViT<sup>70,71</sup>), we found that ViT combined with  
211 *rfCGR* representation maximizes performance (**Figure 6**). A multilayer perceptron, as  
212 employed in previous work<sup>42,60</sup>, could not identify any species correctly here (**Figure 6**).  
213 Similarly, a previously employed shallow 1D convolutional neural network<sup>44</sup>  
214 underperformed more complex architectures (**Figure 6**). *fCGR* showed much higher error  
215 rates than either *rfCGR* or *varKodes*, which yielded similar results but with slightly higher  
216 accuracy for *rfCGR* (**Figure 6**). These results indicate that deep complex neural networks,  
217 while not explicitly developed for DNA signature, are necessary to extract features from  
218 very low-coverage data and distinguish closely related species. Moreover, the method of k-  
219 mer frequency data transformation seems more consequential than the mapping of k-mers  
220 to pixels for the performance of different image representations. Due to its higher  
221 performance, we adopt the combination of *ViT* and *rfCGR*s for subsequent tests.



222

223 **Figure 6.** Effect of image representation and neural network architecture on cross-validation accuracy of  
 224 species identification in *Stigmaphyllon*. One example for each image representation is shown, drawn from the  
 225 same DNA data (SRA accession XXXX) and mapped to a rainbow color scale for increased contrast. See text for  
 226 details on architectures.

227

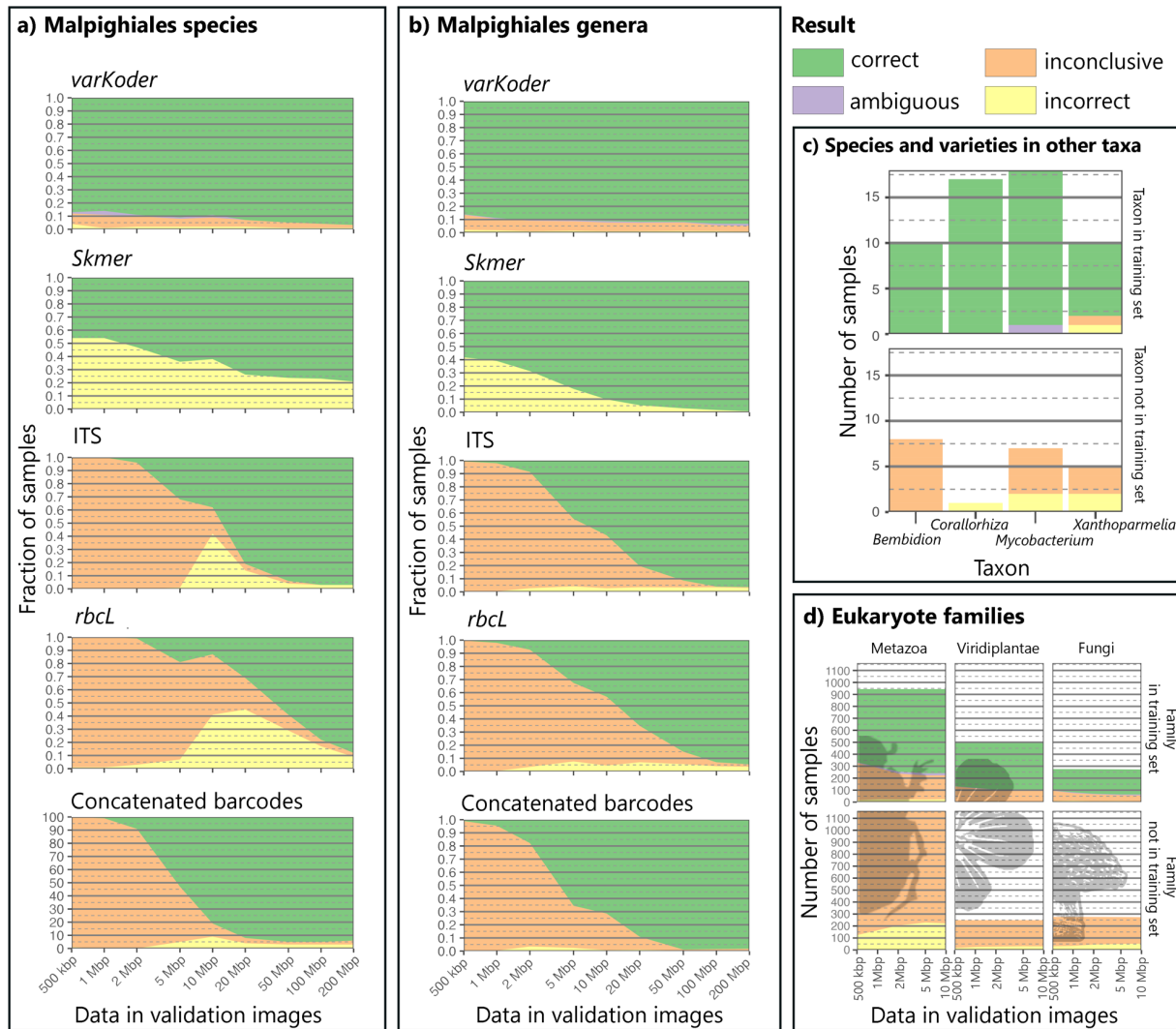
228 In summary, we developed and tested a robust and scalable method of DNA barcoding  
 229 capable of training with small amounts of data, and implemented it in the **varKoder**  
 230 software, which can process sequence data, train an image-classification neural network  
 231 using varKodes or rfCGRs, query new data with a trained neural network, and convert  
 232 between the alternative k-mer mappings. These tasks are accomplished with widely used  
 233 tools for sequence processing<sup>82–86</sup> and for neural network training<sup>87–90</sup>.

234

**235 varKodes are highly accurate for identification of species, genera, and families**

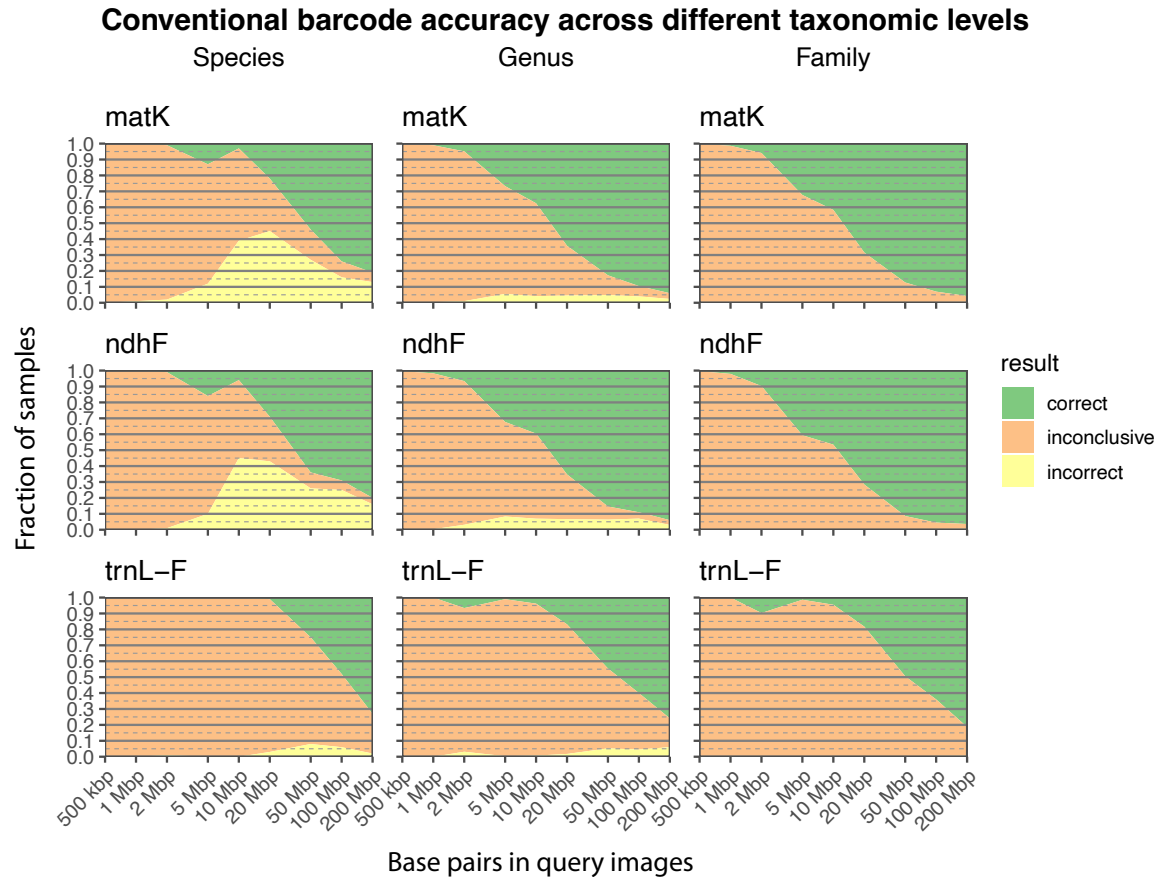
236 To test varKoder under a real-world scenario with heterogeneous data (e.g., large numbers  
237 of taxa, multiple replicates per taxon, varying sequence depth and sample quality), our *de*  
238 *novo* genomic data set included 287 accessions: 100 samples of *Stigmaphyllon* from our  
239 initial development outlined above, plus additional genera in the families Malpighiaceae  
240 (31 genera; 151 samples), Chrysobalanaceae (8 genera; 30 samples), and Elatinaceae (1  
241 genus; 6 samples) in the order Malpighiales. We found high cross-validation accuracies for  
242 species identity of *Stigmaphyllon* (87.0–96.7% correct, 94.6%–98.9% precision, 88.0%–  
243 96.7% recall depending on data input amount; **Figure 7A**). Most errors were inconclusive  
244 predictions (2.2–10%), instead of ambiguous (0–3%) or incorrect (1–4%) predictions.  
245 *varKoder* is robust to the amount of input sequence data necessary for model training,  
246 performing well even at the lower range of input data (**Figure 7A**). Assuming an average  
247 genome size of about 2 Gbp for the average species of Malpighiaceae<sup>91</sup>, the 500Kbp–  
248 200Mbp of data used here represented exceptionally low coverages of about  $\sim 0.0002\times$  –  
249  $0.107\times$ . Moreover, when compared to cross-validation accuracies of alternative barcoding  
250 methods, *varKoder* accuracy is higher than *Skmer*, which showed 46% correct predictions  
251 (57.5% precision, 46% recall) with minimal data amounts and peaked at 79.1% for the  
252 larger data amounts (80% precision, 79.1% recall, **Figure 7A**). On the other hand,  
253 conventional barcodes including individual plastid genes and nuclear ribosomal ITS  
254 regions performed well for both BLAST-based (25–97% correct, 66.6–97.3% precision, 25–  
255 97% recall depending on the gene) and phylogenetic-based (94–95% correct, >99%  
256 precision, 97.2–98.4% recall for concatenated matrices) approaches when at least 50 Mbp  
257 of data was provided (**Figure 7A, Figure 8**). However, these results were much worse  
258 when <50 Mbp of data were available (down to zero correct for BLAST), with unsuccessful  
259 locus assembly leading to inconclusive predictions as the primary reason for the failure  
260 (**Figure 7A, Figure 8**). Finally, an unsupervised clustering method based on neural  
261 networks applied to *fCGRs* (*iDeLUCS*<sup>92</sup>) reached 24–60% clustering accuracy depending on  
262 input data amount when prompted to cluster *Stigmaphyllon* sequences into 10 groups  
263 (**Table 1**). In summary, *varKoder* reaches much higher accuracy for species determination

264 than existing methods for unprecedentedly small amounts of data and demonstrates  
 265 similar accuracies when greater amounts of sequence data are available.  
 266



267 **Figure 7.** Performance of *varKoder* and alternative barcoding methodologies across different data sets. (A)  
 268 Leave-one-out cross-validation to identify species of Malpighiales using different approaches and amounts of  
 269 data to assemble query samples. (B) Same as (A), but for genera. (C) Performance for species-level  
 270 identification across different publicly-available datasets: *Bembidion* beetles, *Corallorhiza* orchids,  
 271 *Mycobacterium tuberculosis* bacteria, and *Xanthoparmelia* fungi. All query samples used as much data as were  
 272 available. (D) Performance for Eukaryote family-level identification for different amounts of input data.  
 273  
 274





275

276

277

**Figure 8.** Accuracy of conventional barcode loci for species, genera and families within the Malpighiales.

Input	rbcl+matK+ndhF+ITS	plastid+ITS full assembly	278
<b>200 mb</b>	0.59	<del>0.24</del>	<del>279</del>
<b>100 mb</b>	0.6	0.25	280
<b>50 mb</b>	0.29	0.26	281
<b>20 mb</b>	0.27	<del>0.23</del>	<del>282</del>
<b>10 mb</b>	0.29	<del>0.27</del>	<del>283</del>
<b>5 mb</b>	0.24	0.28	284
<b>2 mb</b>	0.27	0.53	284

285

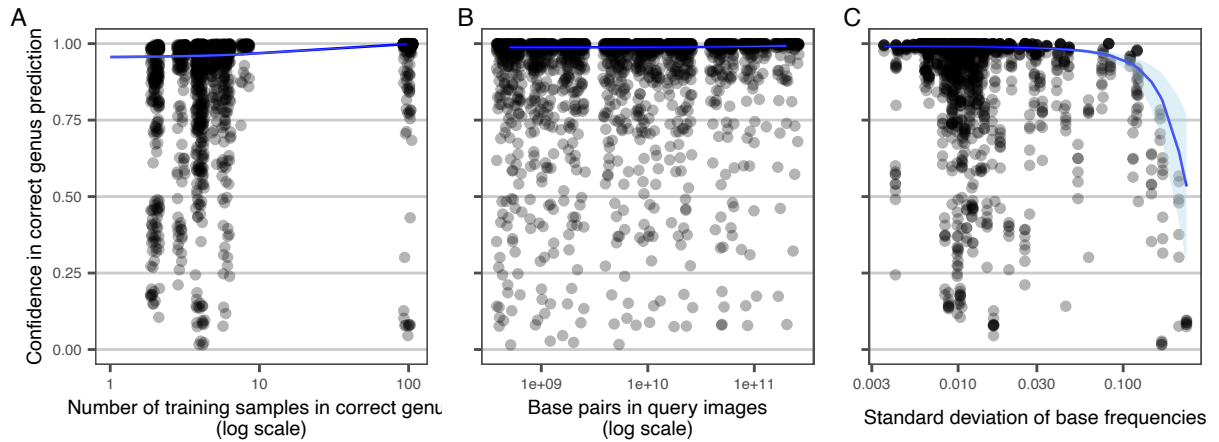
**Table 1.** Accuracy in deuces classification by data amount and plastid genes included.

286

287 Genus-level identification yielded similar high accuracies with *varKoder* (86.1–93.3%  
 288 correct, 97.2%–97.7% precision, 86.4%–94.7% recall depending on input amount, **Figure**  
 289 **7B**), but with a higher rate of inconclusive predictions (4.5–11.5%). A linear model  
 290 demonstrated that this higher uncertainty can be attributed to two factors: (1) samples

291 exhibiting higher levels of DNA damage in genera other than *Stigmaphyllon*; and (2) genera  
292 trained with fewer replicates (e.g., down to 3 samples for some genera; **Figures 9–10**).  
293 Despite this trend, the vast majority of genera with fewer replicates and lower DNA quality  
294 can still be correctly predicted, resulting in the >97% prediction and >86% recall across  
295 the whole dataset. Additionally, samples within genera share fewer genetic similarities  
296 than samples within species, which likely poses a more challenging classification problem.  
297 However, the incorrect rate was very small in all cases (0.7–2.1%), with most errors being  
298 inconclusive or ambiguous predictions. In contrast, *Skmer* exhibited better performance  
299 when larger amounts of data were used (99.2% correct, 99.2% precision, 99.2% recall for  
300 200 Mbp), but performed poorly for lower amounts of data like those commonly generated  
301 from genome skim experiments (58.2% correct, 58.2% precision, 58.2% recall for 500  
302 Kbp) (**Figure 7B**). Genus-level identifications using conventional barcodes in a  
303 concatenated phylogeny were up to 98.1% correct (99.2% precision, 97.2% recall) when  
304 a large amount of data (200 Mbp) was available (**Figure 7B**). But like its application at  
305 species-level identification, most predictions were inconclusive when less than 20 Mbp  
306 reads were used (**Figure 7B**). Although genome skimming can be used to sequence  
307 conventional barcodes, they are more often obtained with amplicon sequencing, which has  
308 failure rates ranging from 15–75% even with highly optimized protocols<sup>93</sup>, leading to an  
309 even higher number of inconclusive predictions. At the family level, *Skmer* and *varKoder*  
310 had near-perfect accuracy across all data amounts (>97% correct), while conventional  
311 barcodes performed well when there were sufficiently large amounts of data (**Figures 8,**  
312 **11**).

## Factors affecting varKode prediction accuracy



313

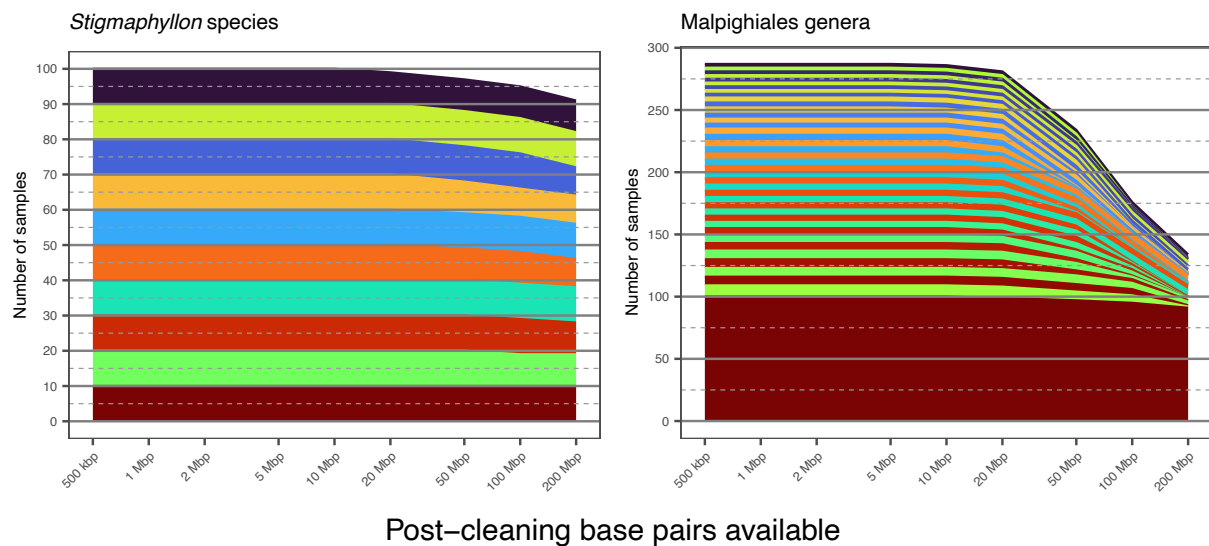
314 **Figure 9.** Predictors of confidence in correct genus. A) Confidence increases with more training samples per

315 genus. B) Amount of data per validation image has little effect. C) Validation samples with low quality have

316 lower confidence.

317

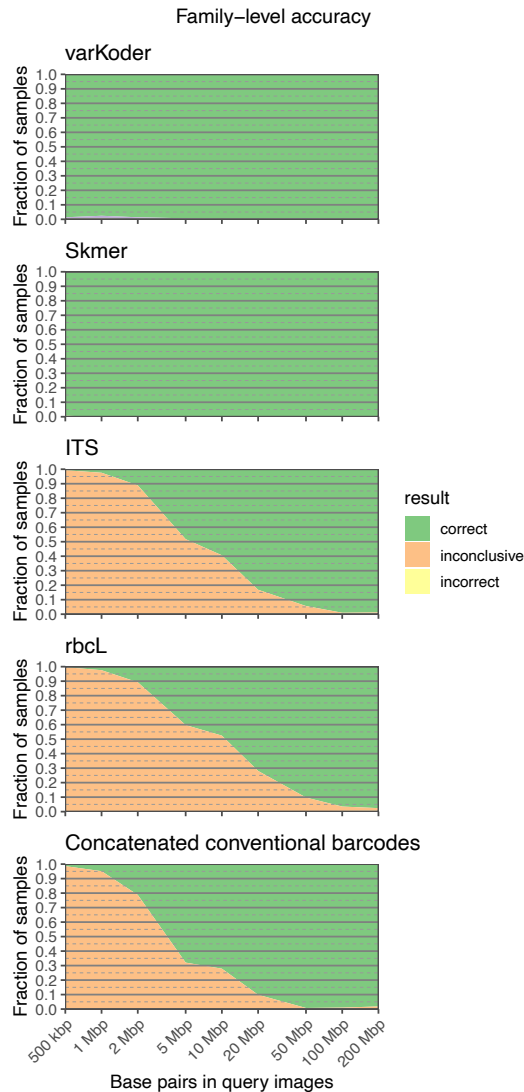
## Number of samples available for different data amounts



318

319 **Figure 10.** Number of samples available for different data amounts in the Malpighiales and Eukaryote

320 families datasets. Arbitrary colors are assigned to individual taxa.



321  
 322 **Figure 11.** Comparison of *varKoder*, *Skmer*, and conventional barcode accuracy for identifying families of  
 323 Malpighiales.

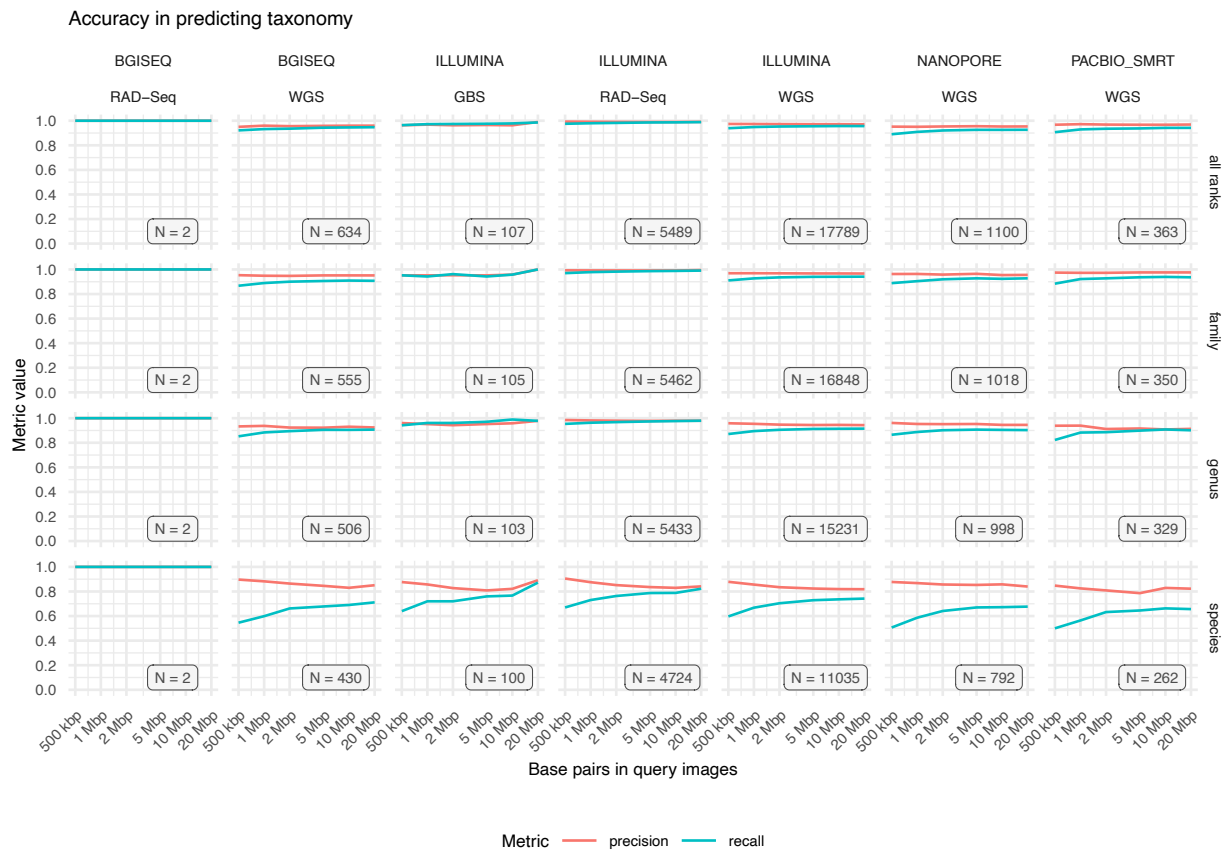
324  
 325 **varKodes are universal and scalable across the Tree of Life**

326 To further test the universality of varKodes, we expanded to sequencing data from diverse  
 327 clades of plants, fungi, animals, and bacteria (**Figure 7C**). These tests included species-level  
 328 identification in insects (*Bembidion* beetles<sup>54,94</sup>) and lichen-forming fungi  
 329 (*Xanthoparmelia*<sup>95</sup>), species and infra-specific taxon identification in coralroot orchids  
 330 (*Corallorhiza*<sup>96</sup>), and clinical isolate identification of strains of human pathogenic bacteria  
 331 (*Mycobacterium tuberculosis*<sup>97</sup>). In all cases, we tested the performance of *varKoder* on taxa  
 332 included in the training set and on taxa not included in the training set. We identified

333 perfect species identification (100% correct, 100% precision, 100% recall) for beetles and  
334 coralroot orchids included in the training set. For bacteria, 5.6% of the validation set  
335 returned ambiguous predictions; the remaining samples were correctly identified (94.7%  
336 precision, 100% recall). In lichen-forming fungi, which include DNA from both the fungal  
337 and algal partners, and thus are more challenging, 10% of the test samples returned  
338 incorrect predictions and another 10% were inclusive; the remainder were correct (89%  
339 precision, 80% recall). For all cases, species or varieties not included in the training set  
340 generally resulted in inconclusive results, with a minority yielding incorrect predictions  
341 **(Figure 7C)**.

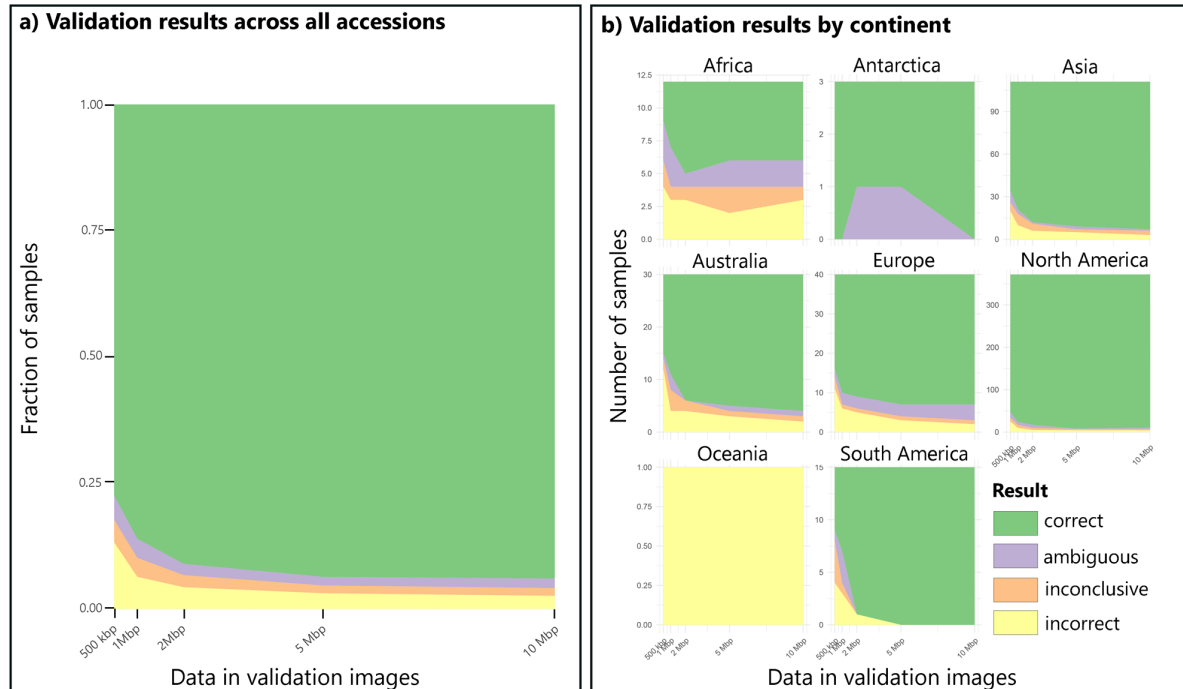
342  
343 Finally, we tested the scalability of varKodes in three large-scale datasets: (1) all 861  
344 eukaryotic families with Illumina data on NCBI SRA, (2) all taxa with multiple accessions on  
345 NCBI SRA, including different sequencing platforms and library strategies (254,819  
346 accessions and 14,151 taxa across all taxonomic ranks), and (3) a previously published  
347 dataset of 2916 soil eDNA samples from all seven continents<sup>98</sup>. Owing to NCBI download  
348 speed bottlenecks, we restricted varKode construction to a very limited maximum of 10  
349 Mbp of DNA data in the former 2 cases. The family-level eukaryote data achieved a rate of  
350 correct predictions of 65.2–81.3% across all kingdoms when families were included in the  
351 training set **(Figure 7D)**, with most errors being inconclusive predictions (17.5–33.1%).  
352 Precision varied from 95.3% to 97.3% and recall from 67.9% to 78.3%. Similarly to the  
353 species- and variety-level exercise, families not included in the training set often yielded  
354 inconclusive predictions **(Figure 7D)**, suggesting a potential for varKoding to be used as a  
355 discovery tool when reasonably well-sampled training data sets are available. The  
356 expanded data with all taxa from NCBI SRA revealed that varKoding is robust to sequencing  
357 platform and library preparation method **(Figure 12)**. Predictions at the family level or  
358 pooled for all the taxonomic hierarchy are accurate regardless of sequencing details (>94%  
359 precision, >86% recall). The much higher accuracy when compared to the dataset based on  
360 Eukaryotic families alone may be an effect of a completely random validation set instead of  
361 stratified by family, resulting in higher representation of commonly sampled families. At  
362 the genus and species level, results are more dependent on the sequencing method **(Figure**

363 12). For genera, precision/recall using 10Mbp of data varies from 90.8%/90.8% with  
 364 whole genome shotgun libraries in PacBio to 97.9%/97.6% with genotype-by-sequencing  
 365 in Illumina. Finally, the eDNA data shows promise in using varKoding to identify the  
 366 geographical origin of an environmental sample: in the validation set, at 10Mbp of DNA  
 367 data, 94.0% of the samples had continent correctly identified, with 2.6% being incorrect,  
 368 1.9% being ambiguous, and 1.5% being inconclusive (84.7% prediction, 84.5% recall)  
 369 (Figure 13).



370  
 371 **Figure 12.** varKoder performance in predicting taxonomy for all data on SRA. Sample sizes refer to the  
 372 number of validation accessions available for each combination of platform, sequencing strategy and  
 373 taxonomic rank.

374  
 375



376

377 **Figure 13.** Varkoder performance in identifying the geographical origin of a soil metabarcoding sample. A)

378 Performance across the whole dataset. B) Performance for each continent.

379

380 A single model classifying all of life is not possible with conventional barcodes. *Skmer*, the  
 381 state-of-the-art genome skimming alternative, cannot be scaled to a dataset of this size: our  
 382 attempt to apply it to Eukaryote families could not be finished after more than 40 days  
 383 using 32 high-performance computing cores. In general, conventional barcodes, when  
 384 derived from genome skimming data, require memory- and processor-intensive sequence  
 385 assembly, and *Skmer* relies on pairwise all-by-all sample comparisons; its computing time  
 386 and required storage both increase quadratically with the number of samples. Neural  
 387 network models, on the other hand, have a fixed size, independent of the number of  
 388 samples used in training, and training time scales linearly with the number of input  
 389 samples. Our most complex model, trained on all taxa available from the NCBI SRA, has  
 390 about 1.3GB of disk size. varKode images also are tiny replacements (8.2 KB on average for  
 391 k-mer length of 7) for much larger genomic data sets (on average, 144 MB per sample  
 392 here). Downloading up to 20Mb of sequence data for over 250,000 accessions from the  
 393 NCBI SRA was the bottleneck, taking over 70 days. By parallelizing processing over 40  
 394 cores, processing this data into varKodes was about 10 times faster, resulting in

395 approximately 18GB of data including all of these accessions. Training a model on more  
396 than 1.3 million images took about 45 hours using only 2 GPUs. Although training on large  
397 datasets requires powerful GPUs and large memory, training on small datasets and  
398 querying is possible on personal computers in a few seconds to minutes. To reduce the  
399 computational resources required for training new datasets, we provide a pre-trained  
400 model from both varKodes and rfCGRs from all taxa on SRA using the huggingface hub  
401 ([https://huggingface.co/brunoasm/vit\\_large\\_patch32\\_224.NCBI\\_SRA](https://huggingface.co/brunoasm/vit_large_patch32_224.NCBI_SRA)). See Asprino et al.<sup>99</sup>  
402 for details on the data used for this model. Whenever the data become available, a model  
403 potentially trained on millions of species easily can be ported to devices without  
404 continuous internet access. Moreover, the minimal data amounts needed for identification  
405 could be generated in seconds in a portable Nanopore device. Finally, the library  
406 preparation method based on shotgun sequencing is very simple and can be automated  
407 with portable consumer devices, such as the Nanopore Voltrax. Together, these properties  
408 allow for more widely distributed applications of varKoding, such as field-laboratory  
409 environments<sup>100</sup> or proposed distributed genetic databases<sup>101</sup>.

410

## 411 **Conclusions**

412 varKoding is universal, accurate, efficient, and holds tremendous promise for documenting  
413 and discovering Earth's biodiversity. It achieves accurate identification with minimal data  
414 compared to existing next-generation sequencing methods, while maintaining universal  
415 applicability across the Tree of Life. Its modular framework can evolve alongside advances  
416 in sequencing technologies, bioinformatics, and machine learning, as exemplified here by  
417 the update in image representation (*varKodes* to *rfCGRs*) and neural network architecture  
418 (resnext to ViT) after initial testing. For these reasons, we expect it will contribute for the  
419 wider adoption of DNA signatures on biodiversity assessments and ecological research,  
420 overcoming current challenges<sup>39</sup>. Reference data for varKoding will be increasingly  
421 available from ambitious efforts in genome sequencing<sup>102–106</sup>. However, we note that  
422 reference data for varKoding is much easier and cost-effective to obtain from low-coverage  
423 genome skims than high-quality contiguous genomes: the robustness to minimal levels of  
424 coverage a central advantage of our method. For example, our cost for a 3× skim of



425 herbarium samples is about \$34 per sample, versus a high-quality genome which may cost  
426 tens-of-thousands of dollars each. Thus, varKoding shows tremendous promise for further  
427 automating species identification from natural history collections<sup>107–109</sup>.

428  
429 We expect that varKoding will be invaluable to the biodiversity science community in  
430 numerous ways, with many avenues remaining to be explored. One of them is the  
431 identification of samples with poor-quality and degraded DNA, such as unidentified  
432 fragmentary fossil and subfossil remains in natural history collections<sup>107,110</sup>. For example,  
433 Malpighiales samples with signs of DNA damage could be correctly identified using  
434 *varKoder* to species or genus in many cases and to family in almost every case. Future  
435 research could explore the lower limits of sample quality and sequence coverage to achieve  
436 accurate identification at different divergence levels. Finally, we expect that new neural  
437 network architectures and forms of DNA representation will continue to be explored. One  
438 limitation of varKoding, as applied here, is the challenging identification of samples within  
439 mixed components such as lichens or environmental DNA. However, with long-read  
440 sequencing, *varKodes* and *rfCGRs* from single reads could potentially include sufficient data  
441 for that end. Moreover, mixed samples could be useful for other ends: *varKodes* could be  
442 used to classify a set of sequences based on any kind of metadata, beyond taxonomy as  
443 demonstrated by our test on the geographical origin of a soil sample.

444

## 445 **Author contributions**

446 BdM conceived *varKodes* and wrote the program *varKoder*. BdM and CCD designed the  
447 research. CCD, XD, YY, LCM, and CA collected the new sequence data. BdM, CCD, LC, YY, PJF  
448 analyzed and interpreted the data. LCM prepared the figures. BdM and CCD wrote the  
449 manuscript with key contributions from LC, YY and PJF. All authors approved the  
450 manuscript.

## 451 Acknowledgments

452 BdM received postdoctoral fellowships from the Harvard University Museum of  
453 Comparative Zoology and the Smithsonian Tropical Research Institute during the early  
454 stages of this study. LC was supported by Harvard University and by a Stengl Wyer  
455 scholarship from the University of Texas at Austin. PF was supported by LVMH Research,  
456 and Dior Science. YY was supported by a postdoctoral fellowship from Harvard University  
457 Herbaria. CCD was supported by Harvard University, LVMH Research, Dior Science, and  
458 National Science Foundation grants DEB-1355064 and DEB-0544039. Computations were  
459 performed at the Harvard Cannon Cluster and the Field Museum Grainger Bioinformatics  
460 Center. We thank the Bauer Core Facility, and especially Claire Reardon, at Harvard  
461 University for providing technical support during the laboratory process. We thank Renata  
462 Asprino and Kylee Peterson for their assistance in obtaining the newly sequenced data  
463 under Harvard's Binding Participation Agreement. The team at Sound Solutions for  
464 Sustainable Science carefully edited early versions of our manuscript.

## 465 References

- 466 1. Hebert, P. D. N., Ratnasingham, S. & de Waard, J. R. Barcoding animal life: cytochrome c  
467 oxidase subunit 1 divergences among closely related species. *Proc. R. Soc. Lond. B* 270,  
468 S96–S99 (2003).  
469
- 470 2. Kress, W. J. Plant DNA barcodes: Applications today and in the future. *J. Syst. Evol.* 55,  
471 291–307 (2017).  
472
- 473 3. Ratnasingham, S. & Hebert, P. D. N. BOLD: The Barcode of Life Data System. *Mol. Ecol.*  
474 *Notes* 7, 355–364 (2007).  
475
- 476 4. Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C. & Willerslev, E. Towards next-  
477 generation biodiversity assessment using DNA metabarcoding. *Mol. Ecol.* 21, 2045–2050  
478 (2012).  
479
- 480 5. Seifert, K. A. Progress towards DNA barcoding of fungi. *Mol. Ecol. Resour.* 9, 83–89  
481 (2009).  
482

- 483 6. Sharkey, M. J. et al. Minimalist revision and description of 403 new species in 11  
484 subfamilies of Costa Rican braconid parasitoid wasps, including host records for 219  
485 species. *ZooKeys* 1013, 1–665 (2021).  
486
- 487 7. Lahaye, R. et al. DNA barcoding the floras of biodiversity hotspots. *Proc. Natl. Acad. Sci.*  
488 *USA* 105, 2923–2928 (2008).  
489
- 490 8. Kuzmina, M. L. et al. Using herbarium-derived DNAs to assemble a large-scale DNA  
491 barcode library for the vascular plants of Canada. *Appl. Plant Sci.* 5, apps.1700079 (2017).  
492
- 493 9. Muñoz-Rodríguez, P. et al. A taxonomic monograph of *Ipomoea* integrated across  
494 phylogenetic scales. *Nat. Plants* 5, 1136–1144 (2019).  
495
- 496 10. Hebert, P. D. N., Penton, E. H., Burns, J. M., Janzen, D. H. & Hallwachs, W. Ten species in  
497 one: DNA barcoding reveals cryptic species in the neotropical skipper butterfly *Astrartes*  
498 *fulgerator*. *Proc. Natl. Acad. Sci. USA* 101, 14812–14817 (2004).  
499
- 500 11. Zeale, M. R., Butlin, R. K., Barker, G. L., Lees, D. C. & Jones, G. Taxon-specific PCR for DNA  
501 barcoding arthropod prey in bat faeces. *Mol. Ecol. Resour.* 11, 236–244 (2011).  
502
- 503 12. Nitta, J. H., Meyer, J., Taputuarai, R. & Davis, C. C. Life cycle matters: DNA barcoding  
504 reveals contrasting community structure between fern sporophytes and gametophytes.  
505 *Ecol. Monogr.* 87, 278–296 (2016).  
506
- 507 13. Kress, W. J. et al. Plant DNA barcodes and a community phylogeny of a tropical forest  
508 dynamics plot in Panama. *Proc. Natl. Acad. Sci. USA* 106, 18621–18626 (2009).  
509
- 510 14. Willis, C. G., Franzone, B. F., Xi, Z. & Davis, C. C. The establishment of Central American  
511 migratory corridors and the biogeographic origins of seasonally dry tropical forests in  
512 Mexico. *Front. Genet.* 5, 433 (2014).  
513
- 514 15. Willerslev, E. et al. Ancient biomolecules from deep ice cores reveal a forested Southern  
515 Greenland. *Science* 317, 111–114 (2007).  
516
- 517 16. Crump, S. E. et al. Ancient plant DNA reveals High Arctic greening during the Last  
518 Interglacial. *Proc. Natl. Acad. Sci. USA* 118, e2019069118 (2021).  
519
- 520 17. Kjær, K. H. et al. A 2-million-year-old ecosystem in Greenland uncovered by  
521 environmental DNA. *Nature* 612, 283–291 (2022).  
522
- 523 18. Fierer, N. et al. Forensic identification using skin bacterial communities. *Proc. Natl.*  
524 *Acad. Sci. USA* 107, 6477–6481 (2010).  
525
- 526 19. Rollo, F., Ubaldi, M., Ermini, L. & Marota, I. Ötzi's last meals: DNA analysis of the  
527 intestinal content of the Neolithic glacier mummy from the Alps. *Proc. Natl. Acad. Sci. USA*  
528 99, 12594–12599 (2002).

- 529  
530 20. Yu, J. et al. Progress in the use of DNA barcodes in the identification and classification of  
531 medicinal plants. *Ecotoxicol. Environ. Saf.* 208, 111691 (2021).  
532
- 533 21. Ashfaq, M. & Hebert, P. D. N. DNA barcodes for bio-surveillance: regulated and  
534 economically important arthropod plant pests. *Genome* 59, 933–945 (2016).  
535
- 536 22. Eaton, M. J. et al. Barcoding bushmeat: molecular identification of Central African and  
537 South American harvested vertebrates. *Conserv. Genet.* 11, 1389–1404 (2010).  
538
- 539 23. Liu, J. et al. Integrating a comprehensive DNA barcode reference library with a global  
540 map of yews (*Taxus* L.) for forensic identification. *Mol. Ecol. Resour.* 18, 1115–1131 (2018).  
541
- 542 24. Ogden, R., Dawnay, N. & McEwing, R. Wildlife DNA forensics—bridging the gap between  
543 conservation genetics and law enforcement. *Endanger. Species Res.* 9, 179–195 (2009).  
544
- 545 25. Williamson, J. et al. Exposing the illegal trade in cycad species (Cycadophyta:  
546 *Encephalartos*) at two traditional medicine markets in South Africa using DNA barcoding.  
547 *Genome* 59, 771–781 (2016).  
548
- 549 26. Costa, F. O. & Carvalho, G. R. The Barcode of Life Initiative: synopsis and prospective  
550 societal impacts of DNA barcoding of Fish. *Genomics Soc. Policy* 3, 29–40 (2007).  
551
- 552 27. Gao, Z., Liu, Y., Wang, X., Wei, X. & Han, J. DNA mini-barcoding: a derived barcoding  
553 method for herbal molecular identification. *Front. Plant Sci.* 10, 987 (2019).  
554
- 555 28. Molina, J. et al. Possible loss of the chloroplast genome in the parasitic flowering plant  
556 *Rafflesia lagascae* (Rafflesiaceae). *Mol. Biol. Evol.* 31, 793–803 (2014).  
557
- 558 29. Cai, L. et al. Deeply altered genome architecture in the endoparasitic flowering plant  
559 *Sapria himalayana* Griff. (Rafflesiaceae). *Curr. Biol.* 31, 1002–1011 (2021).  
560
- 561 30. Richardson, J. E., Pennington, R. T., Pennington, T. D. & Hollingsworth, P. M. Rapid  
562 diversification of a species-rich genus of neotropical rain forest trees. *Science* 293, 2242–  
563 2245 (2001).  
564
- 565 31. Wang, J., Luo, J., Ma, Y.-Z., Mao, X.-X. & Liu, J.-Q. Nuclear simple sequence repeat markers  
566 are superior to DNA barcodes for identification of closely related *Rhododendron* species on  
567 the same mountain. *J. Syst. Evol.* 57, 278–286 (2019).  
568
- 569 32. Su, X., Wu, G., Li, L. & Liu, J. Species delimitation in plants using the Qinghai–Tibet  
570 Plateau endemic *Orinus* (Poaceae: Tridentinae) as an example. *Ann. Bot.* 116, 35–48  
571 (2015).  
572
- 573 33. Lu, Z. et al. Species delimitation and hybridization history of a hazel species complex.  
574 *Ann. Bot.* 127, 875–886 (2021).

- 575  
576 34. Cai, L. et al. The perfect storm: gene tree estimation error, incomplete lineage sorting,  
577 and ancient gene flow explain the most recalcitrant ancient angiosperm clade, Malpighiales.  
578 *Syst. Biol.* 70, 491–507 (2021).  
579
- 580 35. Clarke, L. J., Soubrier, J., Weyrich, L. S. & Cooper, A. Environmental metabarcodes for  
581 insects: in silico PCR reveals potential for taxonomic bias. *Mol. Ecol. Resour.* 14, 1160–1170  
582 (2014).  
583
- 584 36. Song, H., Buhay, J. E., Whiting, M. F. & Crandall, K. A. Many species in one: DNA  
585 barcoding overestimates the number of species when nuclear mitochondrial pseudogenes  
586 are coamplified. *Proc. Natl. Acad. Sci. USA* 105, 13486–13491 (2008).  
587
- 588 37. Xiong, H. et al. Species tree estimation and the impact of gene loss following whole-  
589 genome duplication. *Syst. Biol.* 71, 1348–1361 (2022).  
590
- 591 38. Straub, S. C. K. et al. Navigating the tip of the genomic iceberg: Next-generation  
592 sequencing for plant systematics. *Am. J. Bot.* 99, 349–364 (2012).  
593
- 594 39. Bohmann, K., Mirarab, S., Bafna, V. & Gilbert, M. T. P. Beyond DNA barcoding: The  
595 unrealised potential of genome skim data in sample identification. *Mol. Ecol.* 29, 2882–  
596 2895 (2020).  
597
- 598 40. Sarmashghi, S., Bohmann, K., P. Gilbert, M. T., Bafna, V. & Mirarab, S. Skmer: assembly-  
599 free and alignment-free sample identification using genome skims. *Genome Biol.* 20, 34  
600 (2019).  
601
- 602 41. Borowiec, M. L. et al. Deep learning as a tool for ecology and evolution. *Methods Ecol.*  
603 *Evol.* 13, 1640–1660 (2022).  
604
- 605 42. Arias, P. M., Alipour, F., Hill, K. A. & Kari, L. DeLUCS: Deep learning for unsupervised  
606 clustering of DNA sequences. *PLoS One* 17, e0261531 (2022).  
607
- 608 43. Kari, L. et al. Mapping the space of genomic signatures. *PLoS One* 10, e0119815 (2015).  
609
- 610 44. Fiannaca, A. et al. Deep learning models for bacteria taxonomic classification of  
611 metagenomic data. *BMC Bioinformatics* 19, 198 (2018).  
612
- 613 45. Linard, B., Swenson, K. & Pardi, F. Rapid alignment-free phylogenetic identification of  
614 metagenomic sequences. *Bioinformatics* 35, 3303–3312 (2019).  
615
- 616 46. Desai, H. P., Parameshwaran, A. P., Sunderraman, R. & Weeks, M. Comparative Study  
617 Using Neural Networks for 16S Ribosomal Gene Classification. *J. Comput. Biol.* 27, 248–258  
618 (2020).  
619

- 620 47. Shang, J. & Sun, Y. CHEER: HierarCHical taxonomic classification for viral mEtagEnomic  
621 data via deep leaRning. *Methods* 189, 95–103 (2021).  
622
- 623 48. Arias, P. M. et al. BarcodeBERT: Transformers for Biodiversity Analysis. Preprint at  
624 <https://arXiv.org/abs/2311.02401> (2023).  
625
- 626 49. Badirli, S., Akata, Z., Mohler, G., Picard, C. & Dundar, M. Fine-Grained Zero-Shot Learning  
627 with DNA as Side Information. Preprint at <https://arXiv.org/abs/2109.14133> (2021).  
628
- 629 50. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* 521, 436–444 (2015).
- 630
- 631 51. Cong, Y., Ye, X., Mei, Y., He, K. & Li, F. Transposons and non-coding regions drive the  
632 intrafamily differences of genome size in insects. *iScience* 25, 104873 (2022).  
633
- 634 52. Heckenhauer, J. et al. Genome size evolution in the diverse insect order Trichoptera.  
635 *GigaScience* 11, giac015 (2022).  
636
- 637 53. Schley, R. J. et al. The ecology of palm genomes: repeat-associated genome size  
638 expansion is constrained by aridity. *New Phytol.* 236, 433–446 (2022).  
639
- 640 54. Sproul, J. S., Barton, L. M. & Maddison, D. R. Repetitive DNA profiles Reveal Evidence of  
641 Rapid Genome Evolution and Reflect Species Boundaries in Ground Beetles. *Syst. Biol.* 70,  
642 1111–1122 (2020).  
643
- 644 55. de Medeiros, B. A. S. & Farrell, B. D. Whole-genome amplification in double-digest  
645 RADseq results in adequate libraries but fewer sequenced loci. *PeerJ* 6, e5089 (2018).  
646
- 647 56. Jeffrey, H. J. Chaos game representation of gene structure. *Nucleic Acids Res.* 18, 2163–  
648 2170 (1990).  
649
- 650 57. Deschavanne, P. J., Giron, A., Vilain, J., Fagot, G. & Fertil, B. Genomic signature:  
651 characterization and classification of species assessed by chaos game representation of  
652 sequences. *Mol. Biol. Evol.* 16, 1391–1399 (1999).  
653
- 654 58. de la Fuente, R., Díaz-Villanueva, W., Arnau, V. & Moya, A. Genomic Signature in  
655 *Evolutionary Biology: A Review.* *Biology* 12, 322 (2023).  
656
- 657 59. Avila Cartes, J., Anand, S., Ciccolella, S., Bonizzoni, P. & Della Vedova, G. Accurate and fast  
658 clade assignment via deep learning and frequency chaos game representation. *GigaScience*  
659 12, giac119 (2023).  
660
- 661 60. Solis-Reyes, S., Avino, M., Poon, A. & Kari, L. An open-source k-mer based machine  
662 learning tool for fast and accurate subtyping of HIV-1 genomes. *PLoS One* 13, e0206409  
663 (2018).

- 664  
665 61. Kislyuk, A., Bhatnagar, S., Dushoff, J. & Weitz, J. S. Unsupervised statistical clustering of  
666 environmental shotgun sequences. *BMC Bioinformatics* 10, 316 (2009).  
667
- 668 62. Arias, P. M. et al. Environment and taxonomy shape the genomic signature of  
669 prokaryotic extremophiles. *Sci. Rep.* 13, 16105 (2023).  
670
- 671 63. Murad, T., Ali, S., Khan, I. & Patterson, M. Spike2CGR: an efficient method for spike  
672 sequence classification using chaos game representation. *Mach. Learn.* 112, 3633–3658  
673 (2023).  
674
- 675 64. Davis, C. C. & Anderson, W. R. A complete generic phylogeny of Malpighiaceae inferred  
676 from nucleotide sequence data and morphology. *Am. J. Bot.* 97, 2031–2048 (2010).  
677
- 678 65. Cai, L. et al. Phylogeny of Elatinaceae and the tropical Gondwanan origin of the  
679 Centropalacaceae (Malpighiaceae, Elatinaceae) clade. *PLoS One* 11, e0161881 (2016).  
680
- 681 66. Davis, C. C., Anderson, W. R. & Donoghue, M. J. Phylogeny of Malpighiaceae: evidence  
682 from chloroplast *ndhF* and *trnL-F* nucleotide sequences. *Am. J. Bot.* 88, 1830–1846 (2001).  
683
- 684 67. Anderson, C. Revision of *Ryssopterys* and transfer to *Stigmaphyllon* (Malpighiaceae).  
685 *Blumea* 56, 73–104 (2011).  
686
- 687 68. Anderson, C. Monograph of *Stigmaphyllon* (Malpighiaceae). *Syst. Bot. Monogr.* 51, 1–  
688 313 (1997).  
689
- 690 69. He, T. et al. Bag of Tricks for Image Classification with Convolutional Neural Networks.  
691 Preprint at <https://doi.org/10.48550/arXiv.1812.01187> (2018).  
692
- 693 70. Vaswani, A. et al. Attention Is All You Need. Preprint at  
694 <https://doi.org/10.48550/arXiv.1706.03762> (2017).  
695
- 696 71. Dosovitskiy, A. et al. An Image is Worth 16x16 Words: Transformers for Image  
697 Recognition at Scale. Preprint at <https://doi.org/10.48550/arXiv.2010.11929> (2021).  
698
- 699 72. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the Inception  
700 Architecture for Computer Vision. in 2016 IEEE Conference on Computer Vision and  
701 Pattern Recognition (CVPR) 2818–2826 (IEEE, 2016).  
702
- 703 73. Smith, L. N. A disciplined approach to neural network hyper-parameters: Part 1 -  
704 learning rate, batch size, momentum, and weight decay. Preprint at  
705 <https://arXiv.org/abs/1803.09820> (2018).  
706
- 707 74. Zhang, H., Cisse, M., Dauphin, Y. N. & Lopez-Paz, D. mixup: Beyond Empirical Risk  
708 Minimization. Preprint at <https://doi.org/10.48550/arXiv.1710.09412> (2018).  
709

- 710 75. Yun, S. et al. CutMix: Regularization Strategy to Train Strong Classifiers with Localizable  
711 Features. Preprint at <https://doi.org/10.48550/arXiv.1905.04899> (2019).  
712
- 713 76. Xie, S., Girshick, R., Dollár, P., Tu, Z. & He, K. Aggregated Residual Transformations for  
714 Deep Neural Networks. Preprint at <https://doi.org/10.48550/arXiv.1611.05431> (2017).  
715
- 716 77. Goodfellow, I., Bengio, Y. & Courville, A. Deep Learning (MIT Press, 2016).  
717
- 718 78. Christin, S., Hervet, É. & Lecomte, N. Applications for deep learning in ecology. *Methods*  
719 *Ecol. Evol.* 10, 1632–1644 (2019).  
720
- 721 79. Weiß, C. L. et al. Temporal patterns of damage and decay kinetics of DNA retrieved from  
722 plant herbarium specimens. *R. Soc. Open Sci.* 3, 160239 (2016).  
723
- 724 80. Rachtman, E., Balaban, M., Bafna, V. & Mirarab, S. The impact of contaminants on the  
725 accuracy of genome skimming and the effectiveness of exclusion read filters. *Mol. Ecol.*  
726 *Resour.* 20, 649–661 (2020).  
727
- 728 81. Ben-Baruch, E. et al. Asymmetric Loss For Multi-Label Classification. Preprint at  
729 <https://doi.org/10.48550/arXiv.2009.14119> (2021).  
730
- 731 82. Bushnell, B. BBMap. v.37.61 (2022). Available at:  
732 <http://sourceforge.net/projects/bbmap/>  
733
- 734 83. Bushnell, B., Rood, J. & Singer, E. BBMerge – Accurate paired shotgun read merging via  
735 overlap. *PLoS One* 12, e0185056 (2017).  
736
- 737 84. Rizk, G., Lavenier, D. & Chikhi, R. DSK: k-mer counting with very low memory usage.  
738 *Bioinformatics* 29, 652–653 (2013).  
739
- 740 85. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor.  
741 *Bioinformatics* 34, i884–i890 (2018).  
742
- 743 86. Tange, O. GNU Parallel 2018 (Ole Tange, 2018). Available at  
744 <https://doi.org/10.5281/zenodo.1146014>  
745
- 746 87. Harris, C. R. et al. Array programming with NumPy. *Nature* 585, 357–362 (2020).  
747
- 748 88. Paszke, A. et al. PyTorch: An Imperative Style, High-Performance Deep Learning  
749 Library. in *Advances in Neural Information Processing Systems* 32 8024–8035 (Curran  
750 Associates, Inc., 2019).  
751
- 752 89. Howard, J. & Gugger, S. Fastai: A Layered API for Deep Learning. *Information* 11, 108  
753 (2020).  
754
- 755 90. Wightman, R. PyTorch Image Models (2019). Available at [10.5281/zenodo.4414861](https://doi.org/10.5281/zenodo.4414861).



- 756  
757 91. Pellicer, J. & Leitch, I. J. The Plant DNA C-values database (release 7.1): an updated  
758 online repository of plant genome size data for comparative studies. *New Phytol.* 226, 301–  
759 305 (2020).  
760
- 761 92. Millan Arias, P., Hill, K. A. & Kari, L. iDeLUCS: a deep learning interactive tool for  
762 alignment-free clustering of DNA sequences. *Bioinformatics* 39, btad508 (2023).  
763
- 764 93. D'Ercole, J., Prosser, S. W. J. & Hebert, P. D. N. A SMRT approach for targeted amplicon  
765 sequencing of museum specimens (Lepidoptera)—patterns of nucleotide misincorporation.  
766 *PeerJ* 9, e10420 (2021).  
767
- 768 94. Sproul, J. S. & Maddison, D. R. Cryptic species in the mountaintops: species delimitation  
769 and taxonomy of the *Bembidion breve* species group (Coleoptera: Carabidae) aided by  
770 genomic architecture of a century-old type specimen. *Zool. J. Linn. Soc.* 183, 556–583  
771 (2018).  
772
- 773 95. Keuler, R. et al. Interpreting phylogenetic conflict: hybridization in the most speciose  
774 genus of lichen-forming fungi. *Mol. Phylogenet. Evol.* 174, 107543 (2022).  
775
- 776 96. Barrett, C. F., Wicke, S. & Sass, C. Dense infraspecific sampling reveals rapid and  
777 independent trajectories of plastome degradation in a heterotrophic orchid complex. *New*  
778 *Phytol.* 218, 1192–1204 (2018).  
779
- 780 97. Freschi, L. et al. Population structure, biogeography and transmissibility of  
781 *Mycobacterium tuberculosis*. *Nat. Commun.* 12, 6099 (2021).  
782
- 783 98. Ma, B. et al. A genomic catalogue of soil microbiomes boosts mining of biodiversity and  
784 genetic resources. *Nat. Commun.* 14, 7318 (2023).  
785
- 786 99. Asprino, R. et al. A dataset for benchmarking molecular identification tools based on  
787 genome skimming. (in preparation for Scientific Data).  
788
- 789 100. Pomerantz, A. et al. Rapid in situ identification of biological specimens via DNA  
790 amplicon sequencing using miniaturized laboratory equipment. *Nat. Protoc.* 17, 1415–1443  
791 (2022).  
792
- 793 101. Kimura, L. T. et al. Amazon Biobank: a collaborative genetic database for bioeconomy  
794 development. *Funct. Integr. Genomics* 23, 101 (2023).  
795
- 796 102. Lewin, H. A. et al. The Earth BioGenome Project 2020: starting the clock. *Proc. Natl.*  
797 *Acad. Sci. USA* 119, e2115635118 (2022).  
798
- 799 103. Ebenezer, T. E. et al. Africa: sequence 100,000 species to safeguard biodiversity.  
800 *Nature* 603, 388–392 (2022).  
801

- 802 104. Cheng, S. et al. 10KP: A phylodiverse genome sequencing plan. *GigaScience* 7, giy013  
 803 (2018).  
 804  
 805 105. Nature Biotechnology Editorial. A reference standard for genome biology. *Nat.*  
 806 *Biotechnol.* 36, 1121 (2018).  
 807  
 808 106. i5K Consortium. The i5K Initiative: Advancing Arthropod Genomics for Knowledge,  
 809 Human Health, Agriculture, and the Environment. *J. Hered.* 104, 595–600 (2013).  
 810  
 811 107. Davis, C. C. The herbarium of the future. *Trends Ecol. Evol.* 38, 229–240 (2022).  
 812  
 813 108. Davis, C. C. Collections are truly priceless. *Science* 383, 1035 (2024).  
 814  
 815 109. Davis, C. C., Sessa, E. B., Paton, A., Antonelli, A. & Teisher, J. The destructive sampling  
 816 conundrum and guidelines for effective and ethical sampling of herbaria. Preprint at  
 817 <https://doi.org/10.32942/X2C603> (2024).  
 818  
 819 110. Card, D. C., Shapiro, B., Giribet, G., Moritz, C. & Edwards, S. V. Museum genomics. *Annu.*  
 820 *Rev. Genet.* 55, 633–659 (2021).  
 821

## 822 **Online Methods**

### 823 **Sequence data**

824 *Taxon sampling, DNA sequencing, assembly, and annotation for newly acquired genetic*  
 825 *data*—The newly generated plant data used here and the methods to obtain these data are  
 826 described in detail in a data descriptor article<sup>99</sup>. Briefly, they included members of the large  
 827 and diverse order Malpighiales<sup>34</sup>: Malpighiaceae (251 accessions representing 31 genera),  
 828 Elatinaceae (6 accession for 1 genus), and Chrysobalanaceae (30 accessions for 8 genera).  
 829 Malpighiaceae includes *Stigmaphyllon* with the most comprehensive species sampling: 10  
 830 species and 10 accessions sampled per species. All 100 *Stigmaphyllon* samples were  
 831 sequenced specifically to build, validate, and test our identification models at shallower  
 832 phylogenetic depths, since their taxonomy has been extensively revised by coauthor C.  
 833 Anderson<sup>67,68</sup>. Each of these samples was labeled with species, genus, and family names.  
 834 The focus for the remainder of the Malpighiaceae, Chrysobalanaceae, and Elatinaceae  
 835 sampling was to identify a given sample to genus. In this case, among the non-

836 *Stigmaphyllon* samples we included 3–9 species per genus. Each accession in this case was  
837 labeled with its corresponding genus and family identification. Unlike *Stigmaphyllon*, where  
838 we included multiple accessions per species, there were no additional replicates per  
839 species for our genus-level sampling. For this dataset, we used leave-one-out cross  
840 validation in all assessments, and therefore there are no train and validation sets. For  
841 additional information see Asprino et al.<sup>99</sup>.

842 *Public genomic data compilation*—To further understand the versatility of varKodes more  
843 broadly across the Tree of Life, we tested species identification using genome skim data  
844 sets from four genera of plants, animals, fungi, and a bacterial species. This involved a plant  
845 data set from coralroot orchids (genus *Corallorhiza*)<sup>96</sup>, a beetle data set in the genus  
846 *Bembidion*<sup>54,94</sup>, a lichen-forming fungus in the genus *Xanthoparmelia*<sup>95</sup>, and a bacterial data  
847 set of clinical isolates from *Mycobacterium tuberculosis*, the species of pathogenic bacteria  
848 that causes tuberculosis<sup>97</sup>. In all these cases, we labeled samples with the lowest-level  
849 taxonomic identification available (species, subspecies or isolates). For taxa with two or  
850 more samples available, 20% (with a minimum of 1) were randomly selected for the  
851 validation set, which also included all taxa represented by a single sample (therefore,  
852 absent from the training set). The remaining accessions were used in the training set. See  
853 Asprino et al.<sup>99</sup> for further information.

854

855 We also compiled two broad datasets from the NCBI SRA. The first consists of all 861  
856 eukaryotic families with sufficient sequence read data using the Illumina platform and  
857 whole genome shotgun sequencing. We labeled samples with family name only and  
858 included taxa with at least two associated accessions in the training set. Our validation set  
859 consisted of 20% randomly selected accessions from each family (with a minimum of one),  
860 plus all accessions in families with a single accession available (therefore not part of the  
861 training set). The second broad-scale dataset includes all taxa on NCBI SRA that could be  
862 represented by at least 3 independent accessions. In this case, we included different  
863 sequencing platforms (Illumina, PacBio, Nanopore, BGISEQ) and library preparation  
864 methods (whole genome shotgun, RADseq, GBS). For taxa with too many sequences (such  
865 as humans, crops, etc.), we randomly chose up to 20 accessions for each combination of

866 sequencing platform and library preparation method. Accessions were labeled with all  
867 NCBI taxonomy ranks available for a sample, the library preparation method, and the  
868 sequencing platform. The validation set, in this case, consisted of a random selection of  
869 10% of all samples, not stratified by taxon.

870  
871 Our final dataset was assembled with the aim to extend varKoder beyond taxonomic  
872 identification. We compiled a global soil metagenome eDNA dataset labeled with continent  
873 of origin from Ma et al.<sup>98</sup> We filtered out any metagenomic sample which lacked  
874 information on continent in the Ma et al. metadata. This yielded 2916 soil metagenome  
875 samples across all seven continents. We downloaded 10Mbp DNA data for each sample  
876 directly from NCBI. All metadata for the samples and code used to download and analyze  
877 these data can be found in the GitHub repository for our study.

878

## 879 **varKode design and testing**

880 *Sequence data preprocessing*—Prior to the construction of images, raw reads were lightly  
881 cleaned using the following steps: identical reads were de-duplicated using *clumpify.sh* as  
882 implemented in *BBtools*<sup>83,111</sup>, adapters were removed, low-quality tails trimmed, and  
883 overlapping read pairs merged using *fastp*<sup>85</sup> with options "--detect\_adapter\_for\_pe", "--  
884 dedup", "--dup\_calc\_accuracy 1", "--disable\_quality\_filtering", "--disable\_length\_filtering", "--  
885 trim\_poly\_g", "--merge", "--include\_unmerged", . Next, we randomly selected subsets of  
886 cleaned reads with predefined data amounts, ranging from 500 kbp to 200 Mbp, with  
887 *BBtools*. These data subsets were used to generate a variety of input varKodes for a single  
888 sample and all such images were used for training (see main text and Figure 2A). Finally,  
889 we applied *dsk*<sup>84</sup> to count k-mers of a given length based on clean raw reads (i. e. k-mers  
890 are counted for each read and their frequencies are pooled across reads). *dsk* exhibits good  
891 performance with low memory requirements, which is ideal for potential applications  
892 using varKodes on low-memory devices. We note that analyses for species-level public  
893 datasets have low computational requirements and were performed on an Apple MacBook  
894 with ARM processor architecture.

895 *varKode and rfCGR construction*— We designed novel images—**varKodes**—that portray  
896 relative frequencies of k-mers from low-coverage raw Illumina reads. These are similar to a  
897 frequency chaos game representation (*fCGR*) *sensu* Jeffrey<sup>53</sup>, but optimized for raw reads in  
898 which sequence orientation is unknown, and therefore canonical k-mers and their reverse  
899 complement are indistinguishable. This averaging of canonical k-mer frequencies and their  
900 reverse complements is widely used in the context of raw reads<sup>40,61,62,112,113</sup>. We call these  
901 images varKodes because they enCODE the VARIation in k-mer frequencies in a sample. We  
902 name our method **varKoding** after varKodes, but notice that it is modular and can use  
903 other kinds of DNA image representation. They are meant to represent a DNA signature by  
904 mapping k-mer identity to pixel position in an image, such that k-mers with more similar  
905 composition are closer together. Additionally, the brightness of these pixels represents the  
906 abundance of the associated k-mer, but we use ranks instead of raw frequencies to  
907 decrease the effect of overabundant and artifactual k-mers. In summary, varKodes are  
908 produced by mapping k-mer counts onto a pre-computed map of k-mers to pixels, and  
909 transforming frequency data to pixel brightness. varKode design employed t-SNE<sup>114</sup> and  
910 the python libraries *numpy*<sup>87</sup> and *pillow*<sup>115</sup>. In addition to varKodes, here we also developed  
911 a new image representation that uses the same pixel mapping as *fCGRs* but represents k-  
912 mer abundance as ranks instead of raw frequencies. We named these ranked frequency  
913 chaos game representation (*rfCGR*). Both varKodes and *fCGRs* are saved as 8-bit PNG  
914 images including labels as exif metadata.

915 *Testing k-mer length and data amount*—We employed *fastai*<sup>89</sup> for, a high-level  
916 implementation of neural networks based on *pytorch*<sup>88</sup> for training and prediction. All the  
917 model architectures we applied are image classification models available from the *timm*  
918 library<sup>90</sup>, which have been widely tested using a variety of image types. To identify the  
919 optimal training hyperparameters for our neural network, we conducted a series of tests  
920 using the species-level data set for the genus *Stigmaphyllon*. We generated varKodes for  
921 each of the *Stigmaphyllon* samples. We first tested the joint effect of k-mer length and input  
922 data amount for neural network classification accuracy by selecting three samples per  
923 species as a validation set; the remaining samples were used to train neural networks using  
924 different amounts of input data across 10 randomly generated training sets. As input data

925 for both the validation and training sets, we randomly subsampled the original sequences  
926 into fastq files containing from 500 Kb to 200 Mb (equivalent to about 1,700 to 670,000  
927 2x150bp Illumina reads). In this test, we only included samples that yielded at least 200  
928 million base pairs after cleaning. We also tested the effect of including images for all data  
929 amounts during training. For each replicate, we applied the widely used image  
930 classification neural network *resnet50* architecture<sup>116</sup> to classify varKodes and trained  
931 models for 30 epochs. We visualized the distribution of validation accuracy for each  
932 combination of input data amount and k-mer lengths to find a good balance between both.  
933 Visualizations and code applied for training and evaluation is available in our GitHub  
934 repository.

935 *Neural network optimization*—After identifying an appropriate k-mer length and input data  
936 used to produce varKodes (**Figure 2**), we next tested a series of neural network training  
937 conditions. We varied the neural network model complexity, choosing from seven  
938 commonly used architectures: *resnet50*<sup>116</sup>, *resnet-D*<sup>69</sup> with different depths (18, 50, 101), a  
939 wide *resnet50*<sup>69</sup>, *efficientnet-B4*<sup>117</sup>, and ResNeXt101<sup>76</sup>. We also tested the effect of the  
940 following: random initial weights vs. pretrained weights from the *timm* library<sup>90</sup>, presence  
941 or absence of lighting transforms, presence or absence of label smoothing, and presence or  
942 absence of augmentation strategies (i.e., *CutMix*<sup>75</sup> or *MixUp*<sup>74</sup>). Because these parameters  
943 may have complex interactions, we tested all combinations of architecture, pretraining,  
944 transforms, label smoothing, and augmentation, with 20 replicates for each combination of  
945 conditions. In each replicate, we randomly chose 20% of the samples for each species of  
946 *Stigmaphyllon* as validation and trained the model using the remainder for 30 epochs.  
947 Training was performed using all varKodes available for each sample (from 500kbp to  
948 200Mbp). For validation, we separately evaluated whether each varKode with a different  
949 amount of data was correctly identified. For each replicate and amount of data used to  
950 validate varKodes, we recorded the average validation accuracy across the validation set.  
951 We then applied a linear model to predict the effect of all training parameters and amount  
952 of data in varKodes in the validation set on validation accuracy. Validation accuracy in this  
953 case was arc-sin transformed for linear modeling due to its bounded range of 0–1. We  
954 started from the full model containing all parameters and their interactions and reduced

955 the model step-wise based on AIC scores (i. e. Akaike Information Criteria), as implemented  
956 in the R function step. Visualizations and code applied for training and evaluation is  
957 available in our GitHub repository.

958

959 *Testing sample number requirements*—A legitimate concern with complex neural networks  
960 is that they may require vast amounts of training data and that typical skimming data sets  
961 might be insufficient for them to be useful. We tested the robustness of our models to the  
962 effect of the number of samples per species included in training by using from one to seven  
963 samples per species as training set and the remaining as validation, with 50 replicates per  
964 number of training samples. The batch size used in training was adjusted for the cases with  
965 very few samples included, so that each training epoch included about 10 batches. We  
966 included varKodes from 1Mbp to 200Mbp in both training and validation sets. In this case,  
967 we applied the training parameters informed by our previous analyses: a *resnext101*  
968 architecture, random initial weights, *CutMix* augmentation, and label smoothing for 30  
969 epochs. We visualized the effect of the number of samples by plotting the average  
970 validation accuracy of each sample against the number of training samples used in each  
971 case. Visualizations and code applied for training and evaluation is available in our GitHub  
972 repository.

973

974 *Testing the effect of data quality*—Most of the cases with low accuracy corresponded to  
975 samples with low DNA yield (**Figure 2B**). We identified that DNA extraction yield was  
976 significantly correlated with two metrics of DNA quality: average insert size and variation  
977 in nucleotide composition along reads<sup>79</sup> (**Figure 4**). *varKodes* produced from these samples  
978 may be visually distinct from other samples of the same species (**Figure 5**). For this reason,  
979 we further tested whether sample quality in training or validation impacted accuracy.  
980 Using both quality metrics, we identified the five lowest quality samples for each species.  
981 We next produced training sets using six randomly chosen samples per species, varying the  
982 number of low-quality samples included in training from zero to four. We included

983 varKodes from 1Mbp to 200Mbp in both training and validation sets. We repeated this for  
984 30 replicates for each number of low-quality samples. Like our tests with varying sample  
985 numbers, we applied the following training parameters: a *resnext101* architecture, random  
986 initial weights, *CutMix* augmentation, label smoothing for 30 epochs. For the validation set,  
987 we separately recorded the accuracy for high- and low-quality samples. We then visualized  
988 the effect of inclusion of low-quality samples in the training set by observing the  
989 distribution of validation accuracies for high-quality and low-quality samples across the  
990 range of number of low-quality samples included in the training set. Visualizations and  
991 code applied for training and evaluation is available in our GitHub repository.

992

993 *Implementation of varKoder*—Following all the tests described above, we implemented the  
994 optimal neural network training strategies in a python program named ***varKoder***.  
995 *varKoder* can process, train and query varKodes and is freely available on our GitHub:  
996 <https://github.com/brunoasm/varKoder>. Because it employs standard neural network  
997 frameworks (namely, *pytorch*<sup>88</sup>, *fastai*<sup>89</sup>, and *timm*<sup>90</sup>), any of the image classification models  
998 and training hyperparameters available now or in the future via these libraries can be  
999 easily adapted and applied to varKode classification. Moreover, we have implemented a  
1000 multi-label model as the default to increase robustness to low-quality varKodes with little  
1001 diagnostic information in the training set. This was done by using an asymmetric multi-  
1002 label loss function<sup>81</sup> instead of the standard cross-entropy loss function used in single-label  
1003 classification. Analyses used development versions of *varKoder* starting with v.0.8.0.  
1004 Improvements suggested during the peer-review process are now implemented in  
1005 *varKoder* v.1.1.0.

## 1006 ***varKoder* evaluation and comparison to alternatives**

1007 *varKoder*—To test *varKoder* performance on a complex dataset spanning multiple  
1008 taxonomic levels and varying phylogenetic depths, we used the Malpighiales dataset  
1009 including genera in Elatinaceae, Chrysobalanaceae and Malpighiaceae. Species of  
1010 *Stigmaphyllon* (Malpighiaceae) were labeled with species, genus, and family names; all



1011 other samples were labeled with genus and family names. We tested the performance of  
1012 *varKoder* in each sample with leave-one-out cross-validation. For each sample, we retained  
1013 it as validation and trained a neural network using all the other samples. In preliminary  
1014 assessments, we found that a ViT<sup>71</sup> architecture combined with a multi-label model  
1015 sometimes led to instability in training for some datasets. For that reason, we used a two-  
1016 step approach. Models first were pre-trained for 20 epochs as single-label, using the least  
1017 inclusive taxonomic assignment available for each sample and a base learning rate of 0.05.  
1018 Next, we trained for an additional 10 epochs using the pre-trained weights but with a much  
1019 smaller learning rate (0.005) and a multi-label output. Training samples included varKodes  
1020 from 500 Kbp to 200 Mbp, and we recorded validation accuracy separately for varKodes  
1021 produced from each amount of data. We used an arbitrary confidence threshold of 0.7 to  
1022 make predictions in the multilabel models. For validation samples, we deemed a prediction  
1023 correct if only the correct taxon was predicted for each taxonomic rank (i.e., species, genus,  
1024 family). We deemed a prediction incorrect if one or more predictions passed the threshold  
1025 for a taxonomic rank, but none match the actual label. When predicted labels included both  
1026 the correct and incorrect taxa, we deemed it ambiguous. If the output prediction included  
1027 no taxon with confidence above the threshold, we considered it as inconclusive. As metrics  
1028 across all samples, we used prediction and recall, averaged across all predictions. We  
1029 visualized the fraction of correct, incorrect, ambiguous, and inconclusive samples for each  
1030 taxonomic rank and each amount of data used to produce varKodes. The code to reproduce  
1031 training conditions and evaluation tests is available on GitHub.

1032 To test the joint effect of neural network architecture and image representation method,  
1033 we applied this cross-validation approach to all combinations of three image  
1034 representations and four neural network architectures. The architectures tested included:  
1035 (1) *ResNeXt101*<sup>76</sup>, the optimal convolutional neural network architecture in our initial tests,  
1036 (2) *ViT*<sup>71</sup>, a transformer-based architecture that became available after our initial testing,  
1037 (3) a neural network with two convolutional layers processing vectorized k-mer counts,  
1038 following Fiannaca et al<sup>44</sup> and (4) a multi-layer perceptron formed by a series of fully  
1039 connected layers as specified in Millán Arias et al<sup>42</sup>. The two latter have been previously  
1040 employed for *fCGR* data. The three representations tested include *varKodes* and *rfCGRs* as

1041 developed here, and *fCGRs* as estimated by iDeLUCS<sup>92</sup>. In the latter case, we used iDeLUCS  
1042 functions to produce *fCGRs* as 2D python arrays of k-mer counts. Next, we rescaled these  
1043 counts to the range of 0–255 and rounded them to the nearest integer. These arrays were  
1044 then saved as 8-bit png images. All code used in *varKoder* analyses is available on GitHub.

1045 *Skmer*—To compare *varKoder* with alternative methods, we used fastq files cleaned and  
1046 subsampled by *varKoder* as input files to *Skmer*. In this case, we also used leave-one-out  
1047 cross-validation to evaluate performance. For each amount of input data (500Kbp to  
1048 200Mbp), we cycled through all samples, constructing a *Skmer* database with the "*skmer*  
1049 *reference*" command and including all samples but one and default settings. We then used  
1050 the "*skmer query*" command with default settings on the sample left out and deemed the  
1051 identification as correct if the sample in the reference database with closest estimated  
1052 genetic distance had the correct taxon label. Because *Skmer* could always query a sample  
1053 and there is no objective criterion to consider matches beyond the best match, the output  
1054 predictions can only be correct or incorrect, but not inconclusive or ambiguous. We  
1055 visualized the results similarly as we did with *varKoder*. The code to reproduce *Skmer*  
1056 analyses is available on GitHub.

1057 *Conventional plant barcodes* —To infer phylogenies from our genome skim data (Figure 1),  
1058 we applied the *PhyloHerb* bioinformatic pipeline<sup>118</sup>, which has been applied recently to a  
1059 taxa ranging from algae to flowering plants<sup>119–121</sup>. Briefly, this pipeline works as follows: for  
1060 plastid loci, *PhyloHerb* maps raw short reads to a database of land plant plastid genomes.  
1061 Mapped reads are then assembled into scaffolds using *SPAdes*<sup>122</sup> and plastid loci are  
1062 identified using nucleotide BLAST searches with a default e-value threshold of 1e-40.  
1063 *PhyloHerb* then outputs orthologous plastid genes into individual FASTA files, which are fed  
1064 directly into MAFFT v7.407<sup>123</sup> for alignment. Alignments are then concatenated into a  
1065 super matrix using the 'conc' function within the *PhyloHerb* package. Phylogenies for both  
1066 individual locus and the concatenated alignment were inferred with IQTREE v2.0.6 using  
1067 the GTR+GAMMA model with 1000 ultrafast bootstrap replicates<sup>124</sup>.

1068 To recover the traditional plant barcodes, *rbcl*, *matK*, *trnL-F*, *ndhF*, and ITS, from our  
1069 Malpighiales genome skim data, we applied GetOrganelle v1.7.7.0<sup>125</sup> and *PhyloHerb*

1070 v1.1.1<sup>118</sup> to automatically assemble and extract these DNA markers, respectively. Briefly,  
1071 the complete or subsampled genome skim data were first assembled into plastid genomes  
1072 or nuclear ribosomal regions using *GetOrganelle* with its default settings. Next, *PhyloHerb*  
1073 was applied to extract the relevant barcode genes using its built-in BLAST database. To test  
1074 whether these traditional barcodes provided accurate identification to species, genus, and  
1075 family, we ran an all-by-all BLASTn analysis for each individual gene across the same data  
1076 subsampling schemes as *Skmer* and *varKoder*. BLAST targets were always drawn from  
1077 assemblies using all the data available for each specimen, whereas queries included  
1078 assemblies from input data amounts varying from 500 Kbp to 200 Mbp. Within each BLAST  
1079 analysis for each one of the Malpighiales accessions, we deemed an identification to be  
1080 correct if the best non-self BLAST hit came from the same taxon, and incorrect otherwise.  
1081 We deemed it inconclusive if the locus could not be assembled for that amount of data. For  
1082 concatenated barcodes, we produced a phylogenetic tree for each amount of data and  
1083 deemed an identification to be correct if the sample with lowest patristic distance came  
1084 from the same taxon. We deemed it to be inconclusive when none of the genes in the  
1085 concatenated dataset could be assembled for a sample. We visualized results similarly to  
1086 *varKoder*, separately for each conventional barcoding gene and for the concatenated  
1087 dataset. The code to reproduce conventional barcode analyses is available on GitHub.

1088 *iDeLUCS*—To evaluate the performance of *varKoder* with another deep learning based  
1089 sequence classifier, we applied the sequences assembled from the *PhyloHerb* pipeline to  
1090 *iDeLUCS*<sup>92</sup>. We first used concatenated sequences of five traditional plant barcodes (*rbcL*,  
1091 *matK*, *trnL-F*, *ndhF*, and ITS) assembled from input reads varying from 500 Kbp to 200  
1092 Mbp. *iDeLUCS* was run with k-mer length of 6, 100 training epochs, 100 data augmentations  
1093 per sequence, and the SGD algorithm for neural network optimization. All input sequences  
1094 were set to be clustered into 10 groups (equal to the total number of species) and the  
1095 accuracy was evaluated with the *cluster\_acc* function implemented in *iDeLUCS*. We also  
1096 applied the entire plastid genome and the nuclear ribosomal sequence assemblies  
1097 (ETS+18S+ITS1+5.8S+ITS2+28S) in *iDeLUCS* with the same parameters to evaluate the  
1098 impact of input data quality.

## 1099 **Application in diverse taxa**

1100 *Species-level identification in plants, animals, fungi, and bacteria*—For each of the four  
1101 organismal clades, we trained a multi-label model that included five species with at least  
1102 three samples per species. For *Bembidion*, we included five species with five samples per  
1103 species. For *Corallorhiza*, we included five species (or varieties) with at least five samples  
1104 per species, except for *C. striata* var. *vreelandii* and *C. striata* var. *striata*, for which we  
1105 included six and seven samples each, respectively. For *Mycobacterium tuberculosis*, we  
1106 included representatives of five monophyletic *M. tuberculosis* lineages (L1, L2, L3,  
1107 L4.1.i1.2.1, and L4.3.i2) with seven clinical isolates per lineage. Samples for *Bembidion*,  
1108 *Corallorhiza*, and *M. tuberculosis* isolates all formed monophyletic groups, whereas  
1109 *Xanthoparmelia* species did not. Since the *Xanthoparmelia* species were paraphyletic, we  
1110 subsampled only monophyletic groups for model training. In this case, four species  
1111 included three samples per species (*X. camtschadalis*, *X. mexicana*, *X. neocumberlandia*, and  
1112 *X. coloradoensis*) and one species included five samples per species (*X. chlorochroa*). One  
1113 potential confounding factor for the *Xanthoparmelia* model is that *Xanthoparmelia* is a  
1114 lichen-forming fungus and thus genome skim data represents a chimera of fungal and algal  
1115 genomes representing both partners in this unique symbiosis. Species of the algal symbiont  
1116 *Trebouxia* are flexible generalists across fungal species *Xanthoparmelia*. Since these  
1117 genome skims are a mix of both algal photobiont and fungus, we hypothesize that the  
1118 accuracy of our model decreased because of the more generalist nature of *Trebouxia*<sup>126</sup>.

1119  
1120 For all four test cases, we applied default *varKoder* v.0.8.0 parameters for generating *rfCGR*  
1121 images, training each model, and testing the accuracy of the trained model using the ‘query’  
1122 function. In all cases, we included all the available data for each training or validation  
1123 sample. To test if trained models accurately predicted species identity, we queried them  
1124 using extra genome skim samples not used for training but from the same species included  
1125 in the model. We also tested genome skim test samples of species within the same genus  
1126 *not* used in model training. As in the case of Malpighiales, we set the threshold to make a  
1127 prediction equal to 0.7 and used the same criteria to consider a prediction correct,  
1128 incorrect, inconclusive, or ambiguous. We separately evaluated results for taxa with

1129 representatives included in the training set and taxa used only as queries, without  
1130 conspecific samples in the training set. The code to reproduce these analyses is available on  
1131 GitHub.

1132  
1133 *All eukaryotic families data set from SRA*—Each accession was labeled with its family  
1134 identification obtained from NCBI. Because of the larger size of this dataset, a leave-one-out  
1135 cross-validation approach would have been intractable. Therefore, we randomly selected  
1136 80% of the samples in each family as the training set and used the remainder for validation.  
1137 Similarly to Malpighiales, we used a two-step training method by pre-training as a single-  
1138 label model and finalizing with a multi-label model. Pre-training was done with a learning  
1139 rate of 0.1 and a batch size of 300 for 30 epochs. Final training was done with the same  
1140 batch size but a smaller base learning rate of 0.01 in 5 epochs with frozen body weights and  
1141 three epochs with unfrozen weights. The code to reproduce these analyses is available on  
1142 GitHub.

1143  
1144 *All taxa from SRA*—For each accession, we created *rfCGRs* from 500Kbp to 10Mbp of data.  
1145 Each accession was labeled with all the taxa in its taxonomic tree, as well as library strategy  
1146 (RAD, GBS or WGS) and sequencing platform (Illumina, PACBIO, Nanopore or BGISEQ). We  
1147 randomly selected 10% of the samples as validation set, and eliminated from validation  
1148 samples all labels absent from the training set. We used a two-step training method. First,  
1149 we pre-trained using a single-label strategy, using as labels the concatenation of library  
1150 strategy, sequencing platform, kingdom, family and genus. For pretraining, we used a  
1151 learning rate of 0.1, a batch size of 500 and 30 epochs. We then used the weights of this  
1152 pre-trained model as starting weights for a multi-label model including all labels. We  
1153 trained the model for additional 50 epochs with unfrozen body weights and 10 epochs with  
1154 frozen weights, learning rate of 0.05 and batch size of 600. The code to reproduce these  
1155 analyses is available on GitHub.

1156  
1157 *Environmental metagenome global identification*—The downloaded soil metagenomes from  
1158 Ma et al.<sup>98</sup> were labeled by source continent. Similarly to the eukaryotic family data set  
1159 from SRA, we randomly selected 80% of the samples as the training set and used the

1160 remaining 20% as the validation set. We used a two-step training method by pre-training  
1161 as a single-label model and finalizing with a multi-label model. Pre-training was done with  
1162 a learning rate of 0.1 and a batch size of 64 for 30 epochs. Final training was done with the  
1163 same batch size but a smaller base learning rate of 0.01 in 5 epochs with frozen body  
1164 weights and three epochs with unfrozen weights. The code to reproduce all these analyses  
1165 is available on GitHub.

## 1166 **Methods-only references**

- 1167 111. Bushnell, B. BBtools v.37.61 (2017).  
1168
- 1169 112. Nissen, J. N. et al. Improved metagenome binning and assembly using deep variational  
1170 autoencoders. *Nat. Biotechnol.* 39, 555–560 (2021).  
1171
- 1172 113. Vurture, G. W. et al. GenomeScope: fast reference-free genome profiling from short  
1173 reads. *Bioinformatics* 33, 2202–2204 (2017).  
1174
- 1175 114. Maaten, L. van der & Hinton, G. Visualizing Data using t-SNE. *J. Mach. Learn. Res.* 9,  
1176 2579–2605 (2008).  
1177
- 1178 115. Clark, A. Pillow, Version 9.4.0. Software. (2023). Available at:  
1179 <https://doi.org/10.5281/zenodo.7498081>  
1180
- 1181 116. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition.  
1182 Preprint at <https://doi.org/10.48550/arXiv.1512.03385> (2015).  
1183
- 1184 117. Tan, M. & Le, Q. V. EfficientNet: Rethinking Model Scaling for Convolutional Neural  
1185 Networks. Preprint at <https://doi.org/10.48550/arXiv.1905.11946> (2019).  
1186
- 1187 118. Cai, L., Zhang, H. & Davis, C. C. PhyloHerb: A high-throughput phylogenomic pipeline  
1188 for processing genome skimming data. *Appl. Plant Sci.* 10, e11475 (2022).  
1189
- 1190 119. Marinho, L. C. et al. Plastomes resolve generic limits within tribe Clusieae (Clusiaceae)  
1191 and reveal the new genus *Arawakia*. *Mol. Phylogenet. Evol.* 134, 142–151 (2019).  
1192
- 1193 120. Lyra, G. de M. et al. Phylogenomics, divergence time estimation and trait evolution  
1194 provide a new look into the Gracilariales (Rhodophyta). *Mol. Phylogenet. Evol.* 165, 107294  
1195 (2021).  
1196

- 1197 121. Marinho, L. C. et al. Phylogenetic Relationships of Tovomita (Clusiaceae): Carpel  
1198 Number and Geographic Distribution Speak Louder than Venation Pattern. *Syst. Bot.* 46,  
1199 102–108 (2021).  
1200
- 1201 122. Bankevich, A. et al. SPAdes: a new genome assembly algorithm and Its applications to  
1202 single-cell sequencing. *J. Comput. Biol.* 19, 455–477 (2012).  
1203
- 1204 123. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:  
1205 Improvements in performance and usability. *Mol. Biol. Evol.* 30, 772–780 (2013).  
1206
- 1207 124. Minh, B. Q. et al. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic  
1208 Inference in the Genomic Era. *Mol. Biol. Evol.* 37, 1530–1534 (2020).  
1209
- 1210 125. Jin, J.-J. et al. GetOrganelle: a fast and versatile toolkit for accurate de novo assembly of  
1211 organelle genomes. *Genome Biol.* 21, 241 (2020).  
1212
- 1213 126. Leavitt, S. D. et al. Fungal specificity and selectivity for algae play a major role in  
1214 determining lichen partnerships across diverse ecogeographic regions in the lichen-  
1215 forming family Parmeliaceae (Ascomycota). *Mol. Ecol.* 24, 3779–3797 (2015).  
1216

## 1217 **Data Availability**

1218 New data generated for this study is described in a Data Descriptor article containing  
1219 accession numbers (doi to be updated upon acceptance).

## 1220 **Code Availability**

1221 Code used in the initial development and test of *varKoder* is available on Github, including  
1222 all code used to produce figures for this manuscript. *varKoder* releases and source code are  
1223 available at <https://github.com/brunoasm/varKoder>.