1

# A universal DNA signature for the Tree of Life

3

4

5

Bruno A. S. de Medeiros[1,2,3], Liming Cai[4,5], Peter J. Flynn[4], Yujing Yan[4], Xiaoshan Duan[4,6],

Lucas C. Marinho[4,7], Christiane Anderson[8], and Charles C. Davis[4]

8

9

[1]Field Museum of Natural History, Chicago, Illinois, 60605, USA

[2]Department of Organismic and Evolutionary Biology, Museum of Comparative Zoology,
Harvard University, Cambridge, Massachusetts, 02138 USA

[3]Smithsonian Tropical Research Institute, Panama City, Panama

[4]Department of Organismic and Evolutionary Biology, Harvard University Herbaria,
Harvard University, Cambridge, Massachusetts, 02138 USA

[5]Department of Integrative Biology, The University of Texas at Austin, Austin, Texas, 78712
USA

[6]College of Forestry, Northwest Agriculture & Forestry University, Yangling 712100,
Shaanxi, China

[7]Departamento de Biologia, Universidade Federal do Maranhão, Av. dos Portugueses 1966,
Bacanga 65080-805, São Luís, Maranhão, Brazil

[8]University of Michigan Herbarium, 3600 Varsity Drive, Ann Arbor, Michigan 48108, USA

23

**Corresponding authors:**

Bruno A. S. de Medeiros, Field Museum of Natural History, Chicago, IL, 60605; E-mail:

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

29

30

# Abstract

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

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

53 # Introduction

54    For two decades, conventional DNA barcoding, which relies on standardized short

55    sequences (400–800 bp) for species identification[1–5], has enabled novel and massively

56    scalable science spanning evolution[4,6–9]; ecology[10–14] and paleontology[15–19]. Practical

57    applications of barcoding have also made major contributions to environmental health,

58    including the ability to authenticate medicinal plants[20], detect agricultural pests[21], and

59    monitor poaching and the trade of endangered species[22–27]. Despite these remarkable

60    achievements, conventional DNA barcoding suffers from at least four limitations. First,

61    barcodes are customized specifically for a taxon (e.g., plants, animals, and fungi), and

62    therefore are not universal. For example, commonly used plant barcodes from chloroplast

63    genes such as *mat*K and *rbc*L cannot be applied as barcodes for all plants[28,29], or for animals

64    and fungi. Second, conventional barcode loci may fail to distinguish closely related taxa, a

65    pervasive shortcoming in plants[2,30]. Third, reliance on a single locus may lead to spurious

66    results in the case of complex evolutionary scenarios such as hybridization in deep or

67    shallow time[31–34]. And fourth, the necessary comparison of homologous genes may fail

68    when PCR primers are not universal[35], the source DNA is fragmented[27], or paralogy and the

69    presence of pseudogenes confounds accurate orthology assessments[36,37].
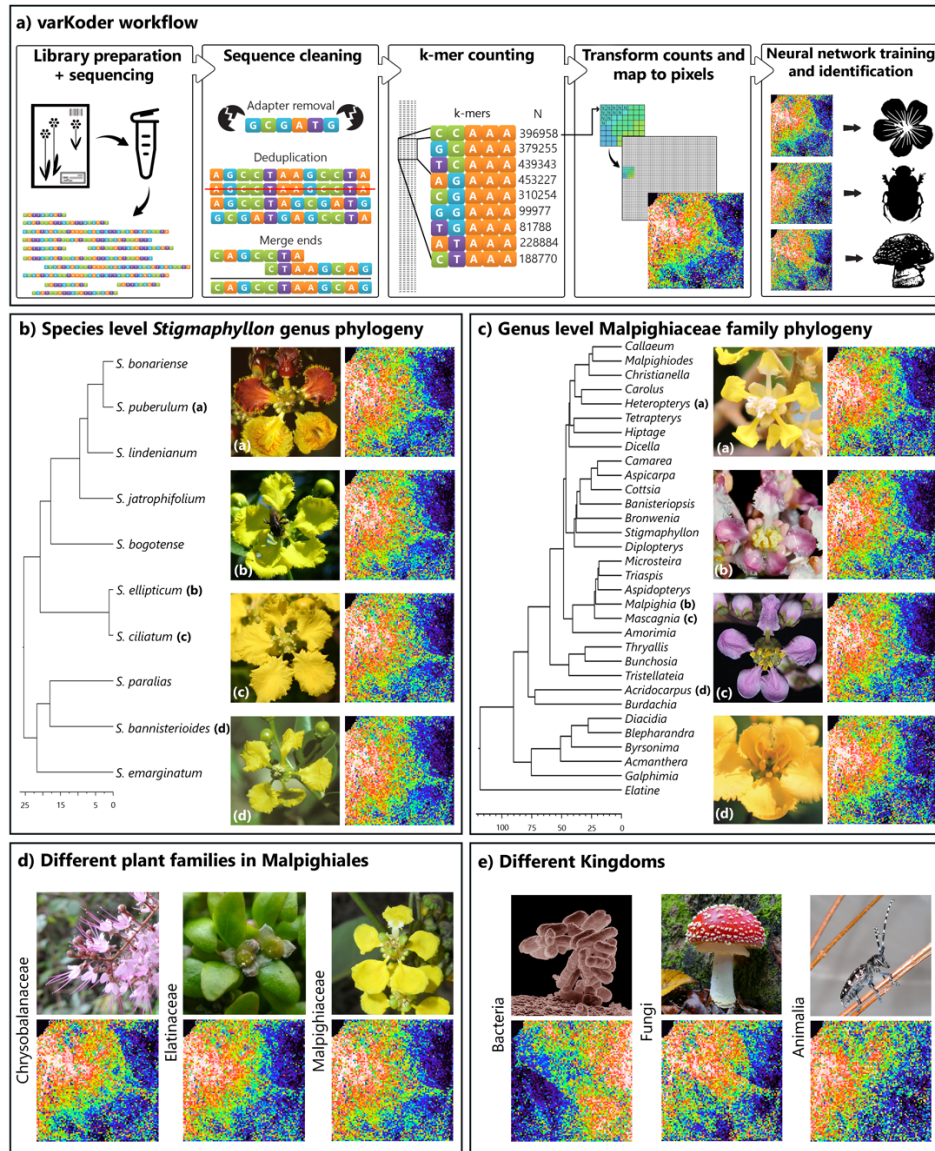
70

71    Newer alternatives to conventional barcoding have begun to address these challenges by

72    leveraging high-throughput sequencing and machine-learning powered by deep neural

73    networks. High-throughput sequencing facilitates more comprehensive assessments of

74    total genomic space[38,39]. For example, presence and absence patterns among short DNA

75    sequences (k-mers) from low-coverage reads (i.e., genome skims) can estimate overall

76    sequence distances, bypassing genome alignments entirely as implemented in *Skmer*[40].

77    Machine learning enables more complex sequence comparisons than conventional methods

78    that rely on homology and simple metrics[41]. Machine-learning models can cluster DNA

79    sequences without supervision[42,43] or classify sequences based on reference datasets[44–49].

80    In particular, neural networks are exceptionally powerful for sophisticated computer-

81    vision tasks, such as image classification[50]. Thus, the combination of low-coverage genome

82    skimming data and neural networks holds enormous promise for accurate and scalable

83    DNA barcoding, but its potential has yet to be fully realized[39].

84

85    Genomes differ substantially in many features beyond the simple nucleotide divergence

86    commonly used in conventional barcoding, but these genomic features have been

87    overlooked in species identification[31,51–55]. We propose that (1) relevant genomic features

88    can be captured by nucleotide composition with short k-mer counts and very small

89    sequence coverage; and (2) these counts can be used to distinguish species and higher taxa

90    efficiently and accurately using machine learning. Prior work on k-mer-based

91    representations of genome composition (i.e., genomic signatures) has shown high accuracy

92    can be achieved with high-coverage data or a large number of replicates per taxon,

93    particularly for identification at higher taxonomic ranks[42–47,56–63]. However, given the

94    millions of existing species and the sparse genetic data available, a practical scalable

95    method would require: (1) consistently high accuracy despite limited evolutionary

96    divergence; (2) fast computations; and (3) high accuracy with small training datasets (both

97    in number of samples and DNA data per sample). Here we developed a novel genomic

98    signature method, which we call **varKoding**, that integrates very low-coverage genome

99    skim data with optimized training of machine-learning models using two-dimensional

100   images representing genome composition (**Figure 1A**). We focus on images as forms of

101   genomic representation since they can be easily stored and accessed across computing

102   platforms, annotated with metadata and readily employed as input data in popular

103   machine learning frameworks such as pytorch[64]. Specifically, our method relies on raw

104   unassembled genomic reads sampling a very small fraction of a genome, since sequence

105   assembly is costly both in terms of DNA sequencing and computation[40,58] and sparse

106   sampling of genomic regions may be sufficient to summarize its features[39]. To develop and

107   optimize varKoding for accurate species identification, we generated a *de novo* genome

108   skim dataset including hundreds of samples derived primarily from historical herbarium

109   specimens for the diverse plant genus *Stigmaphyllon* (Malpighiaceae), which has received

110   extensive phylogenetic and taxonomic treatment[65–69]. Next, we explored the utility of

111   varKoding and compared it to alternatives at different phylogenetic depths from families to

112    species within the flowering plant order Malpighiales (Malpighiaceae, Chrysobalanaceae,

113    and Elatinaceae). Finally, we demonstrate the scalability of varKoding and its potential

114    application in forensics and related fields by testing it on (1) species-level datasets from

115    fungi, plants, animals, and bacteria; (2) massive datasets retrieved from the NCBI sequence

116    read archive (SRA); and (3) a previously published environmental DNA (eDNA) dataset.
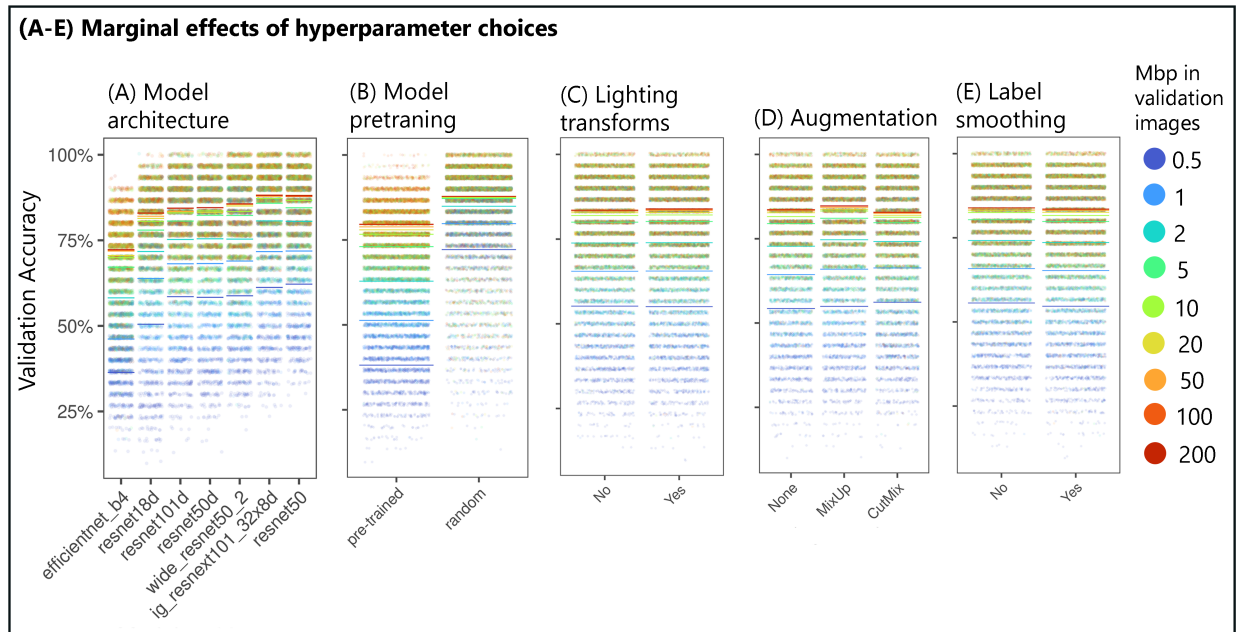
117



118

119    **Figure 1.** Overview of varKoding.  (**A**) Image generation workflow, depicting varKodes.

120    Images are natively grayscale, but here they are mapped to a rainbow color scale for

121    increased contrast. (**B**) Phylogeny and example varKodes of *Stigmaphyllon* species. (**C**)

122    Phylogeny and example varKodes of Malpighiaceae genera including their closest outgroup

123    (*Elatine*, Elatinaceae). Time trees in 1B and 1C were derived from an ongoing family-wide

124    phylogenomic investigation of the family Malpighiaceae (C. C. Davis personal

125    communication) using methods and fossil constraints described in Cai et al.[66]. (**D**)

126    Examples of varKodes from across plant families of Malpighiales, and (**E**) across kingdoms.

127    Chronograms depicted for each representative set with timelines in millions of years (Myr)

128    at the bottom of **B** and **C**.

# Results and Discussion

129

130    **Genomic signature images can be classified with generalized neural networks**

131    We first generated a novel kind of image representation of a genomic signature based on

132    raw reads, which we termed a **varKode**. varKodes map k-mers onto pixels of a 2-D image

133    based on their similarity and represent ranked k-mer frequencies as pixel brightness.

134    Variation in varKodes can be small but remain visually perceptible among species (**Figure**

135    **1B**) and genera (**Figure 1C**). Variation is more striking among higher levels of phylogenetic

136    divergence, such as between families in the order Malpighiales (**Figure 1D**) or different

137    kingdoms of eukaryotes and prokaryotes (**Figure 1E**). We expected, therefore, that neural

138    network architectures developed for image classification, (e.g., deep residual networks,

139    resnets[70] or vision transformers, ViT[71,72]) would be able to differentiate varKodes.
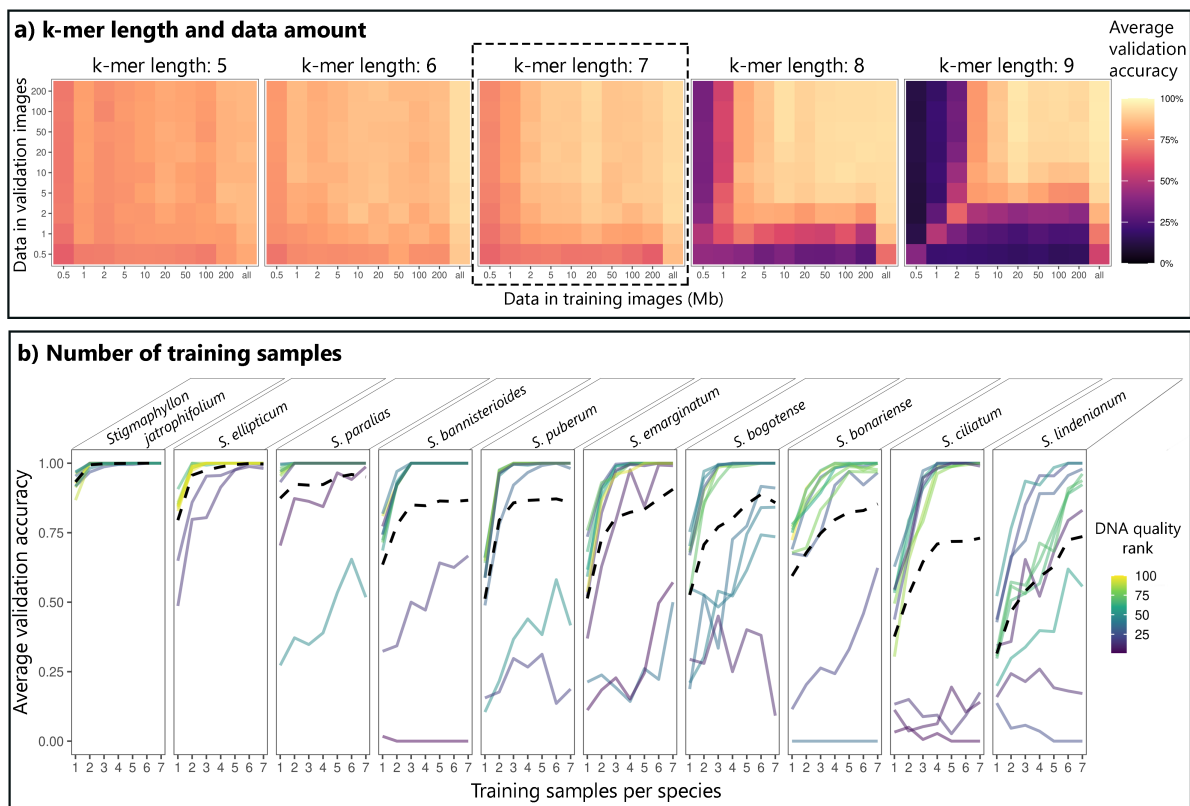
140

**(A-E) Marginal effects of hyperparameter choices**



**Figure 2**. 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.

We first optimized hyperparameters and training conditions to maximize accuracy for species-level identification of *Stigmaphyllon*. We identified that varKodes depicting k-mer length = 7 struck a good balance between accuracy and the amount of input sequence data (**Figure 3A**). Furthermore, models trained with augmented data from several subsampled sequences drawn from each individual exhibited substantially better performance (**Figure 3A**). A linear model demonstrated that neural network architectures and training methods designed for image classification of photographs[70,73–76] are extremely useful for varKode-based identification. Specifically, we observed increased accuracy with more parameter-rich neural network architectures (*ResNeXt101*[77], among those tested), augmentation with lighting transformations, *CutMix*[76] and *MixUp*[75]. Label smoothing[78] and pretraining models on generalized photographs decreased accuracy (**Figure 2**). Contrary to the widely held idea that deep neural networks require very large training datasets[60,79], the aforementioned approaches enabled training with very modest data amounts: four

161    biological replicates per taxon was sufficient for 100% median accuracy (**Figure 3B**).

162    Errors in species-level identification were concentrated among sequences derived from

163    herbarium samples that demonstrated evidence of DNA damage, as is sometimes reported

164    for ancient DNA[80] (**Figure** 3**B**). However, including low-quality training samples slightly

165    decreased mean validation accuracy—from 73% to 71%—for low-quality validation

166    samples, but had no effect on high-quality validation samples (89–90% mean accuracy,
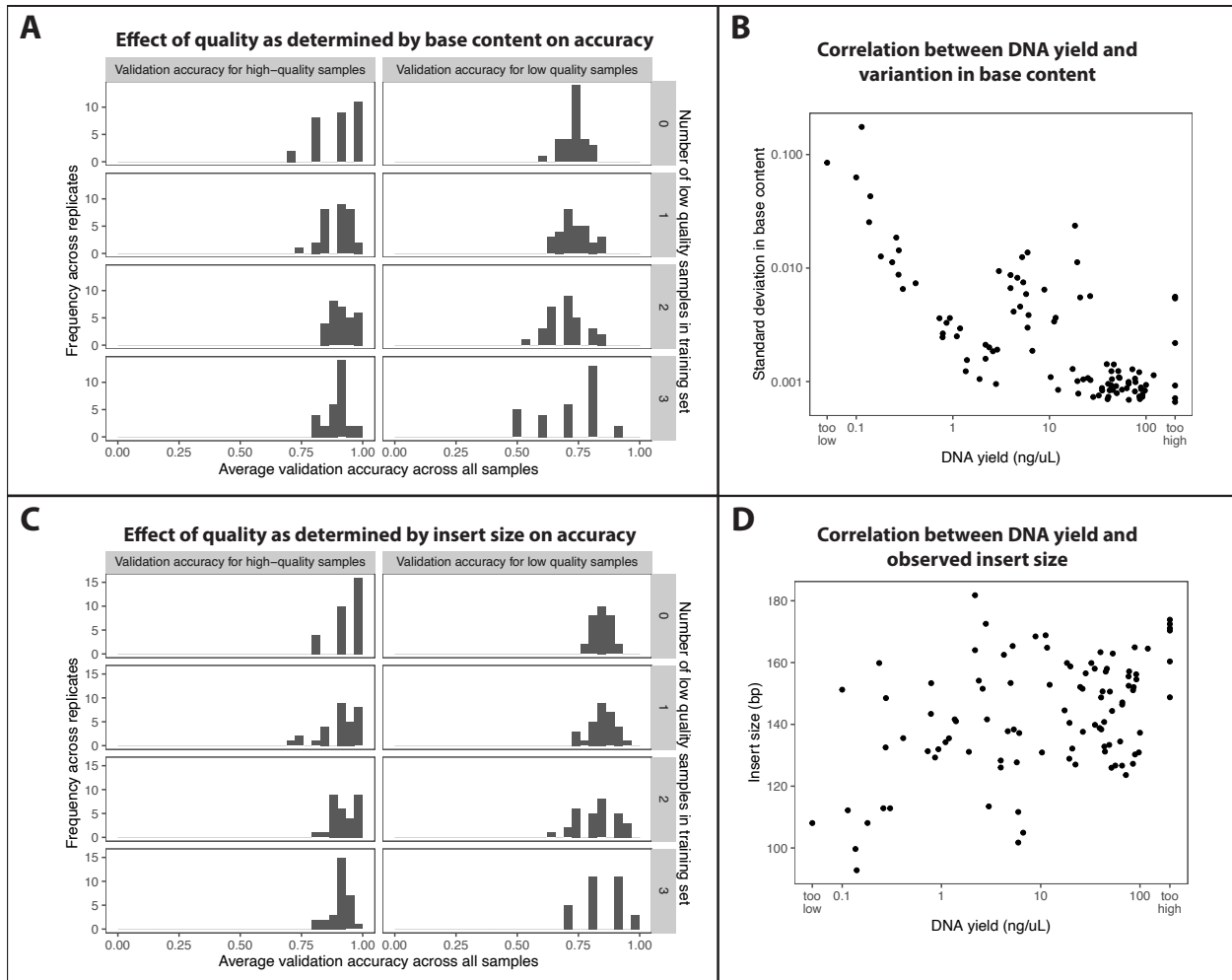
167    **Figure 4**).

168



169

170    **Figure 3. Neural network training of varKodes for species identification.** (**A**) Effect of

171    k-mer length and input data amount used to produce varKodes on validation accuracy.

172    Longer k-mers increase accuracy when more data are used. Mixing varKodes subsampled

173    from different amounts of data improves accuracy. Box with dashed line (k-mer length = 7)

174    strikes a good balance between model accuracy and amount of required data. (**B**)

175    Validation accuracy improves with increased number of training samples per species, but

176    even 3–4 samples are sufficient in most cases for achieving high accuracy. Each solid line

177     represents one sample, colored by DNA quality (i.e., variation in base pair frequencies).

178     Higher rank indicates better quality. Dashed lines represent averages across all samples.

179



180

181     **Figure 4.** Effect of the inclusion of low-quality training samples, inferred from variation in

182     base pair content (A, B) or insert size (C, D). Increasing the fraction of samples in the

183     training set that were low-quality did not strongly affect the average validation accuracy,

184     but it increased dispersion. Low-quality samples are the four samples with highest

185     variation in base-pair content or shortest insert size in raw reads for each species. Panels **B**

186     and **D** show the correlation of each quality metric with DNA extraction yield.

187

188

189    We hypothesized that lower-quality samples shared similar sequences resulting from

190    common patterns of DNA damage and greater levels of microbial or human contaminants,

191    resulting in spurious similarities in varKodes (**Figure 5**). Contaminants also are thought to

192    increase errors in other genome skim methods[81]. To mitigate this problem, we applied

193    multi-label classification[82] to our neural network models. Although single-label

194    classification models always return a single prediction (that is, an inferred label), multi-

195    label models can return zero or more predictions, avoiding spurious results when there is

196    uncertainty. For a set of samples with known labels used for validation, a prediction is a

197    true positive if the predicted label matches the actual label, and a false positive if not.

198    Failure to predict an actual label is deemed a false negative. For each validation sample, we

199    summarized predictions as (1) correct (true positives only); (2) incorrect (false positives

200    only); (3) ambiguous (multiple predictions, including true and false positives); or (4)

201    inconclusive (i. e. no prediction above the confidence threshold). For each test, we

202    summarized results across all validation samples using two metrics: precision (the sum of

203    all true positives divided by the sum of all true and false positives) and recall (the sum of all

204    true positives divided by the sum of all true positives and negatives).

205

varKodes for species of *Stigmaphyllon*

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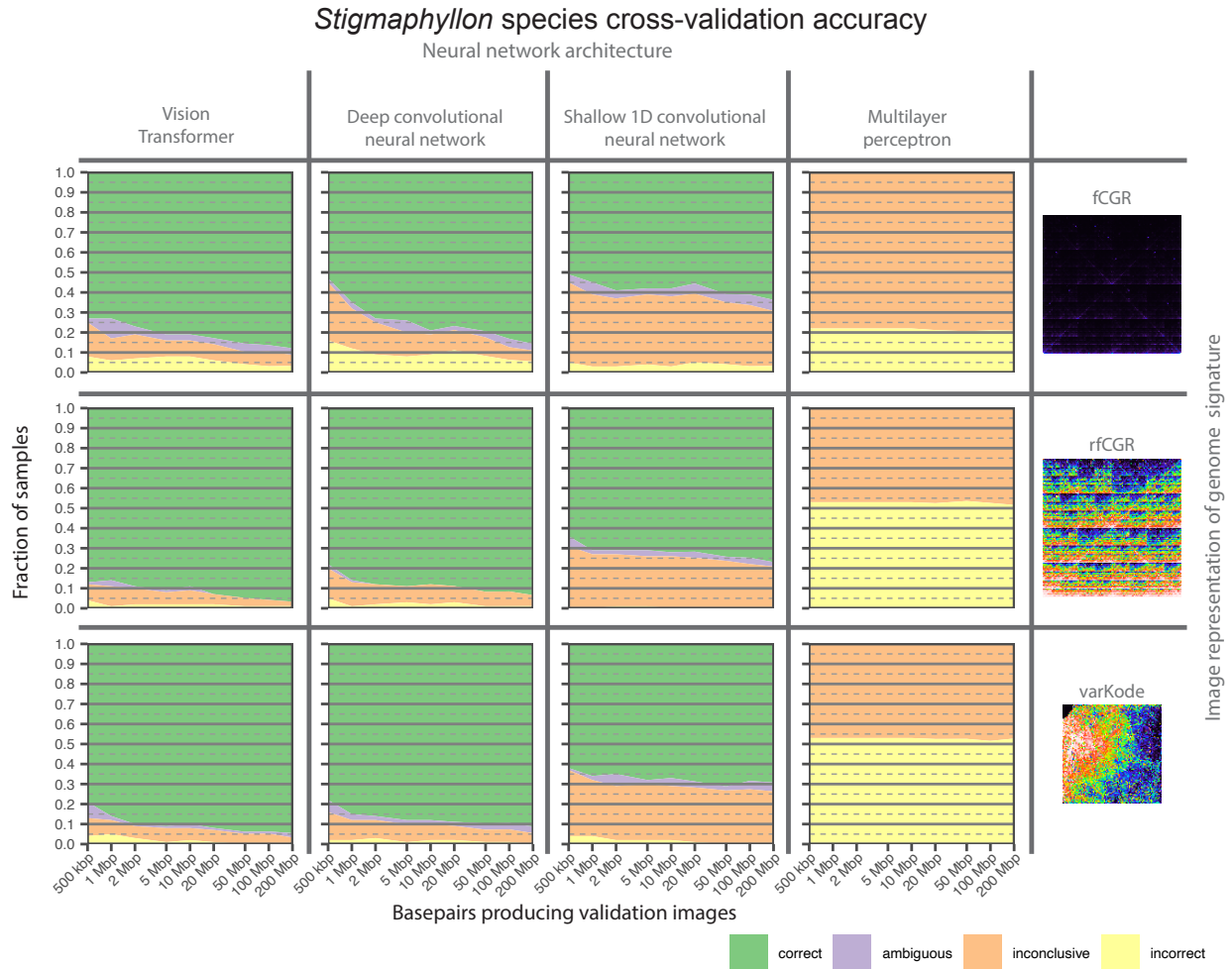ples are organized from lowest (left) to highest (right) DNA quality. Bounding boxes around each sample indicate the average validation accuracy across 30 random replicates with 7 training samples per species.

215

216    After optimizing these training conditions, we directly compared varKodes to an existing

217    method of genomic signature representation: the frequency chaos game representation

218    (*fCGR*)[56,59]. In *fCGR*s, k-mers are mapped to pixels based on their oriented sequence and

219    pixel brightness represents the rescaled k-mer frequency. To isolate the effects of pixel

220    mapping and brightness, we created a new representation combining *fCGR* mapping with

221    *varKode* ranked frequency transformation (*rfCGR*). Because raw sequence reads often

222    contain artifactual k-mers at very high frequencies, especially when low-quality DNA is

223    used to construct libraries, we hypothesized that *rfCGR*s would perform better than *fCGR*s,

224    where pixel brightness is linearly scaled to k-mer counts. By directly comparing these 3

225    kinds of representation combined with four neural network architectures, including (1)

226    two previously employed with *fCGR*s[42,44,60], (2) the optimal architecture in our initial tests

227    (ResNeXt101[77]), and (3) a Vision Transformer (ViT[71,72]), we found that ViT combined with

228    *rfCGR* representation maximizes performance (**Figure 6**). While fCGRs have been initially

229    proposed as tools to study single sequences[56], here we focus all of our comparisons on the

230    task of supervised classification-based genomic composition from very low coverage

231    sequencing. A multilayer perceptron, as employed in previous work[42,60], could not identify

232    any species correctly here (**Figure 6**). Similarly, a previously employed shallow 1D

233    convolutional neural network[44] underperformed more complex architectures (**Figure 6**).

234    *fCGR* showed much higher error rates than either *rfCGR* or *varKodes*, which yielded similar

235    results but with slightly higher accuracy for *rfCGR* (**Figure 6**). These results indicate that

236    deep complex neural networks, while not explicitly developed for genomic signature, are

237    necessary to extract features from very low-coverage data and distinguish closely related

238    species. Moreover, the method of k-mer frequency data transformation seems more

239    consequential than the mapping of k-mers to pixels for the performance of different image

240    representations. Due to its higher performance, we adopt the combination of *ViT* and

241    *rfCGR*s for subsequent tests.

242

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

In summary, we developed and tested a robust and scalable method of DNA barcoding capable of training with small amounts of data, and implemented it in the ***varKoder*** software, which can process sequence data, train an image-classification neural network using varKodes or rfCGRs, query new data with a trained neural network, and convert between the alternative k-mer mappings. These tasks are accomplished with widely used tools for sequence processing[83–87] and for neural network training[64,88–90].

256



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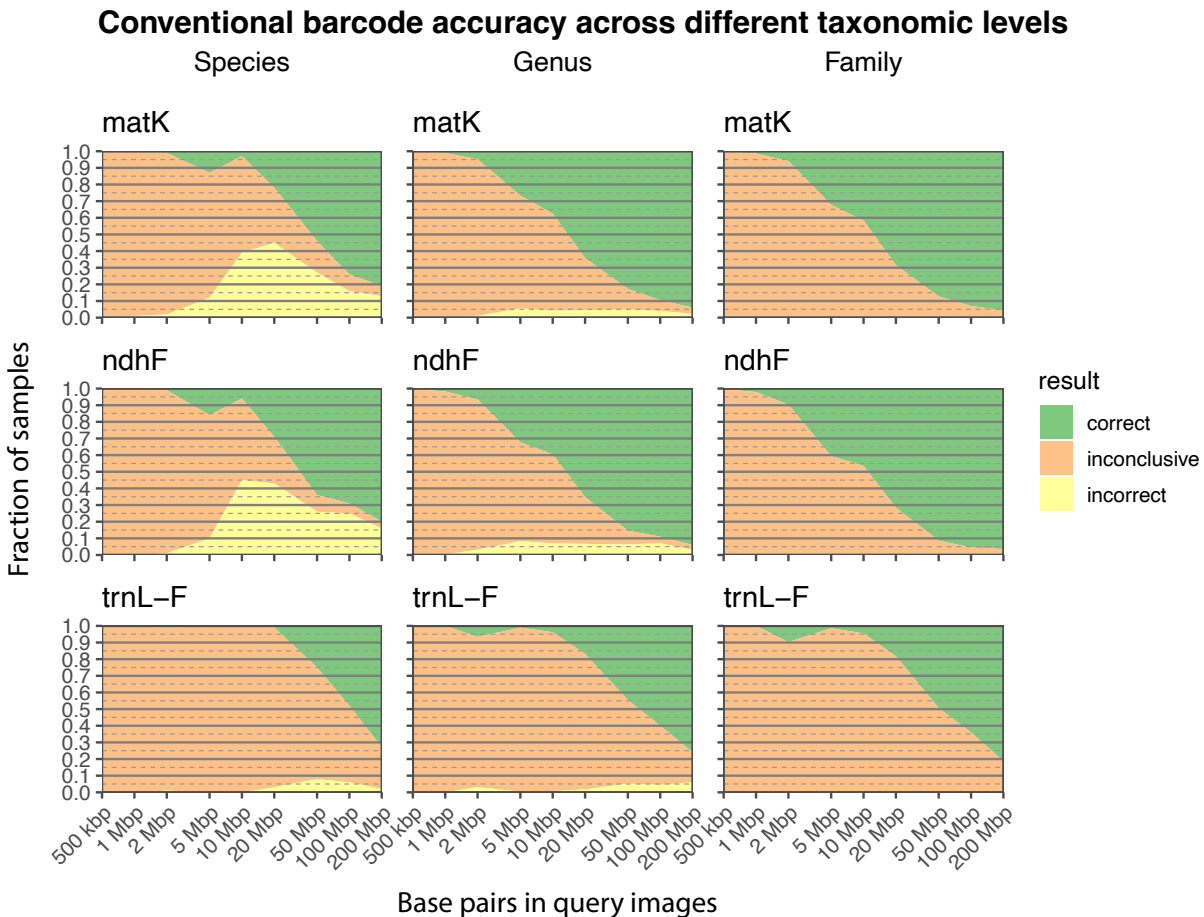
268 **varKodes are highly accurate for identification of species, genera, and families**

269 To test *varKoder* under a real-world scenario with heterogeneous data (e.g., large numbers

270 of taxa, multiple replicates per taxon, varying sequence depth and sample quality), our *de*

271 *novo* genomic data set included 287 accessions: 100 samples of *Stigmaphyllon* from our

272 initial development outlined above, plus additional genera in the families Malpighiaceae

273 (31 genera; 151 samples), Chrysobalanaceae (8 genera; 30 samples), and Elatinaceae (1

274 genus; 6 samples) in the order Malpighiales. We found high cross-validation accuracies for

275 species identity of *Stigmaphyllon* (87.0–96.7% correct, 94.6%–98.9% precision, 88.0%–

276 96.7% recall depending on data input amount; **Figure 7A**). Most errors were inconclusive

277 predictions (2.2–10%), instead of ambiguous (0–3%) or incorrect (1–4%) predictions.

278 *varKoder* is robust to the amount of input sequence data necessary for model training,

279 performing well even at the lower range of input data (**Figure 7A**). Assuming an average

280 genome size of about 2 Gbp for the average species of Malpighiaceae[91], the 500Kbp–

281 200Mbp of data used here represented exceptionally low coverages of about ~0.0002× –

282 0.107×. Such low coverages imply that we are likely not comparing homologous regions

283 across taxa, but rather more general genomic properties that can be inferred from

284 extremely sparse sampling. Moreover, when compared to cross-validation accuracies of

285 alternative barcoding methods, *varKoder* accuracy is higher than *Skmer*, which showed

286 46% correct predictions (57.5% precision, 46% recall) with minimal data amounts and

287 peaked at 79.1% for the larger data amounts (80% precision, 79.1% recall, **Figure 7A**). On

288 the other hand, conventional barcodes including individual plastid genes and nuclear

289 ribosomal ITS regions performed well for both BLAST-based (25–97% correct, 66.6–97.3%

290 precision, 25–97% recall depending on the gene) and phylogenetic-based (94–95% correct,

291 >99% precision, 97.2–98.4% recall for concatenated matrices) approaches when at least 50

292 Mbp of data was provided (**Figure 7A, Figure 8**). However, these results were much worse

293 when <50 Mbp of data were available (down to zero correct for BLAST). In this case,

294 unsuccessful locus assembly leading to inconclusive predictions as the primary reason for

295 the failure (**Figure 7A, Figure 8**), so we expect that alternative methods to BLAST(e.g. [48,92])

296 would not perform substantially better. Finally, an unsupervised clustering method based

297 on neural networks applied to *fCGR*s (*iDeLUCS*[93]) reached 24–60% clustering accuracy

298    depending on input data amount when prompted to cluster *Stigmaphyllon* sequences into

299    10 groups. In summary, *varKoder* reaches much higher accuracy for species determination

300    than existing methods for unprecedentedly small amounts of data and demonstrates

301    similar accuracies when greater amounts of sequence data are available.
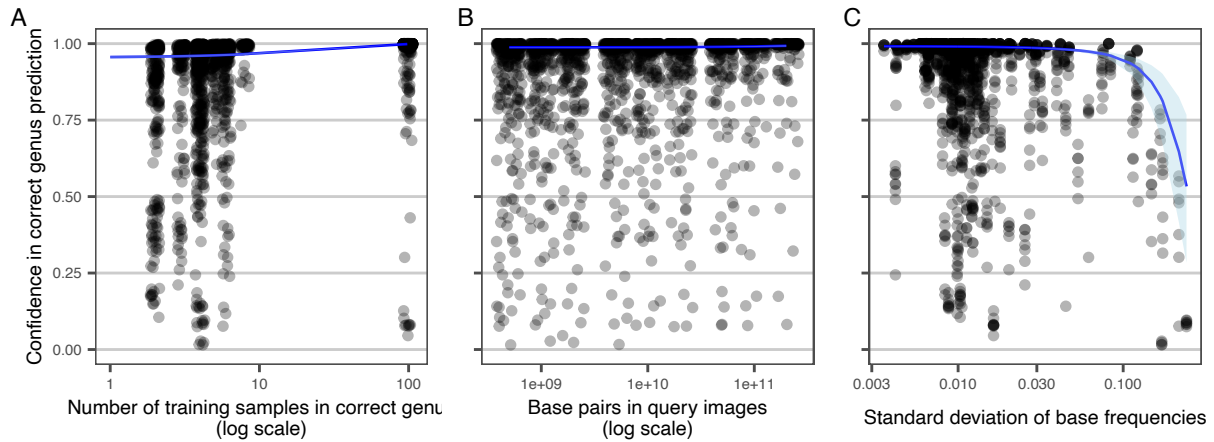
302

303



304

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

306    Malpighiales.

307

308    Genus-level identification yielded similar high accuracies with *varKoder* (86.1–93.3%

309    correct, 97.2%–97.7% precision, 86.4%–94.7% recall depending on input amount, **Figure**

310    **7B**), but with a higher rate of inconclusive predictions (4.5–11.5%). A linear model

311    demonstrated that this higher uncertainty can be attributed to two factors: (1) samples

312    exhibiting higher levels of DNA damage in genera other than *Stigmaphyllon*; and (2) genera

313    trained with fewer replicates (e.g., down to 3 samples for some genera; **Figures 9-10**).

314    Despite this trend, the vast majority of genera with fewer replicates and lower DNA quality

315    can still be correctly predicted, resulting in the >97% prediction and >86% recall across

316    the whole dataset. Additionally, samples within genera share fewer genetic similarities

317    than samples within species, which likely poses a more challenging classification problem.

318    However, the incorrect rate was very small in all cases (0.7–2.1%), with most errors being

319    inconclusive or ambiguous predictions. In contrast, *Skmer* exhibited better performance

320    when larger amounts of data were used (99.2% correct, 99.2% precision, 99.2% recall for

321    200 Mbp), but performed poorly for lower amounts of data like those commonly generated

322    from genome skim experiments (58.2% correct, 58.2% precision, 58.2% recall for 500

323    Kbp) (**Figure 7B**). Genus-level identifications using conventional barcodes in a

324    concatenated phylogeny were up to 98.1% correct (99.2% precision, 97.2%% recall) when

325    a large amount of data (200 Mbp) was available (**Figure 7B**). But like its application at

326    species-level identification, most predictions were inconclusive when less than 20 Mbp

327    reads were used (**Figure 7B**). Although genome skimming can be used to sequence

328    conventional barcodes, they are more often obtained with amplicon sequencing, which has

329    failure rates ranging from 15–75% even with highly optimized protocols[94], leading to an

330    even higher number of inconclusive predictions. At the family level, *Skmer* and *varKoder*

331    had near-perfect accuracy across all data amounts (>97% correct), while conventional

332    barcodes performed well when there were sufficiently large amounts of data (**Figures 8,**

333    **11**).
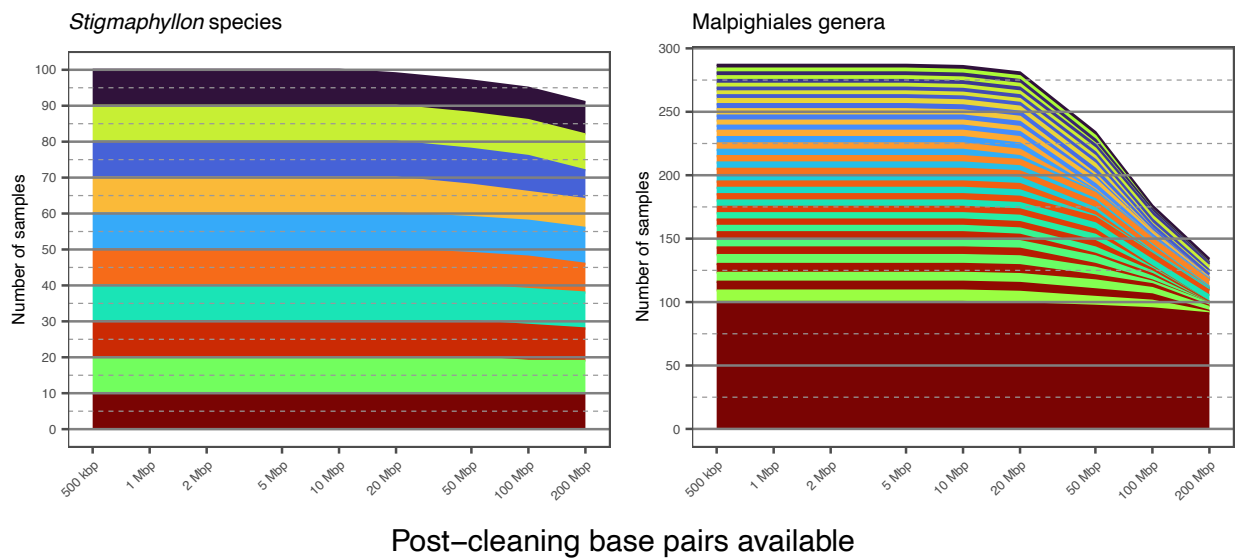
334

Factors affecting varKode prediction accuracy

**Figure 9.** Predictors of confidence in correct genus. A) Confidence increases with more training samples per genus. B) Amount of data per validation image has little effect. C) Validation samples with low quality have lower confidence.
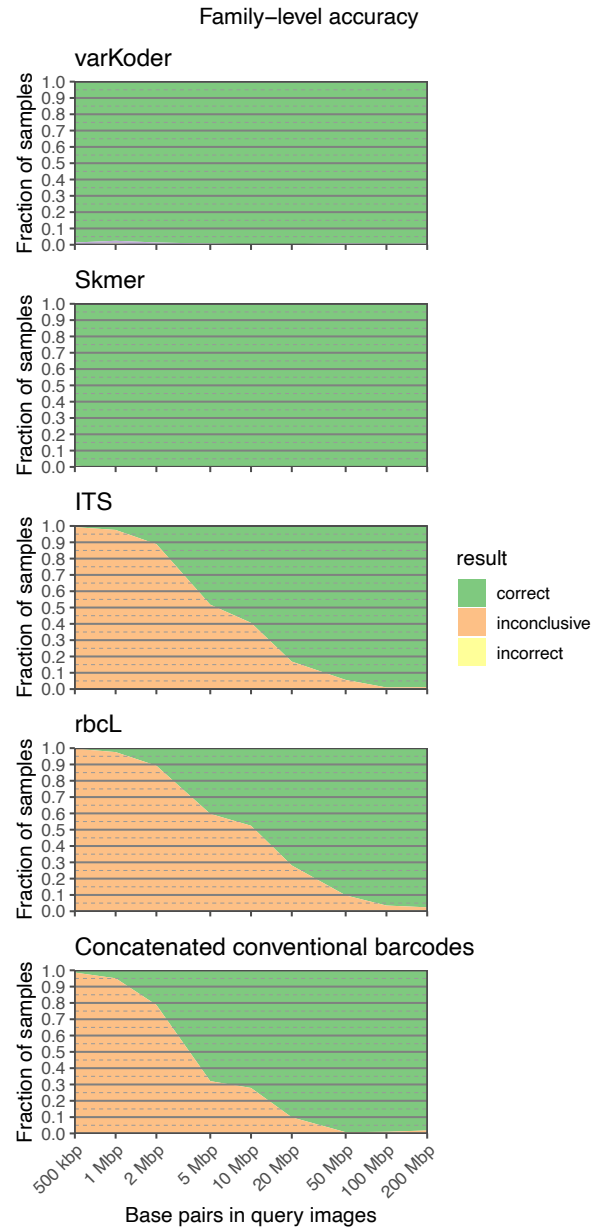


Number of samples available for different data amounts

**Figure 10.** Number of samples available for different data amounts in the Malpighiales and Eukaryote families datasets. Arbitrary colors are assigned to individual taxa.

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

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

To further test the universality of varKodes, we expanded to sequencing data from diverse clades of plants, fungi, animals, and bacteria (**Figure 7C**). These tests included species-level identification in insects (*Bembidion* beetles[54,95]) and lichen-forming fungi (*Xanthoparmelia*[96]), species and infra–specific taxon identification in coralroot orchids
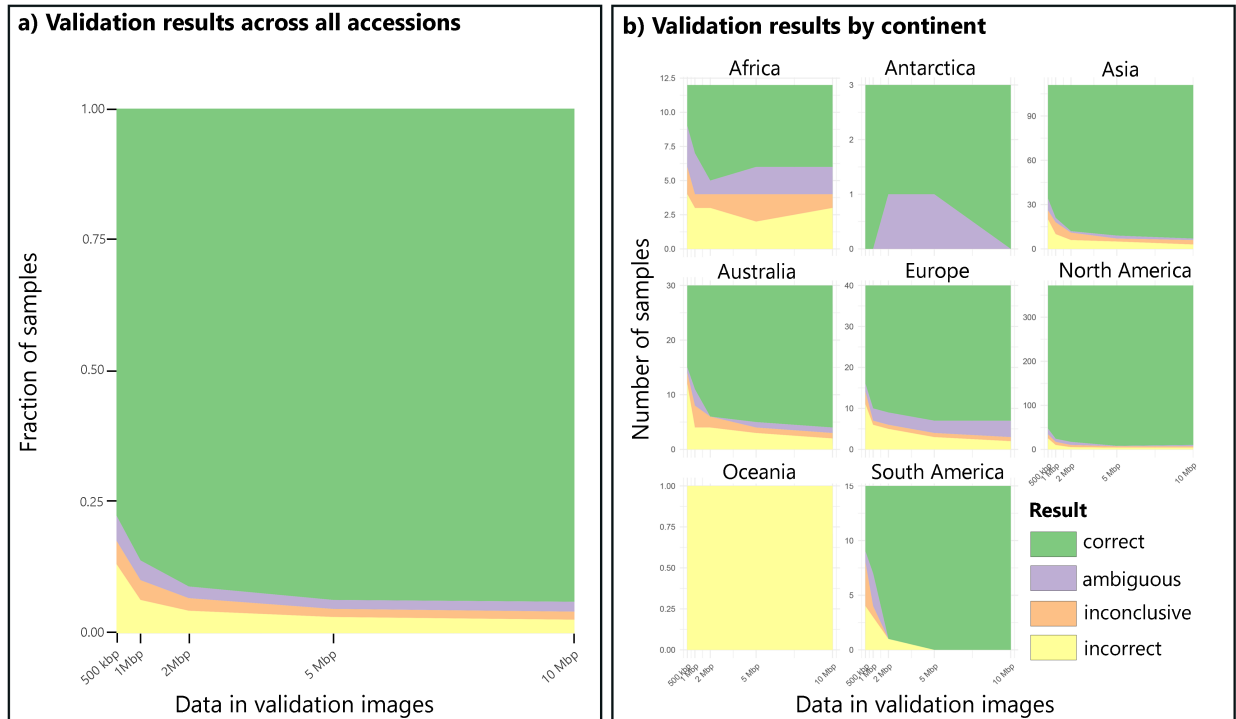
354    (*Corallorhiza*[97]), and clinical isolate identification of strains of human pathogenic bacteria

355    (*Mycobacterium tuberculosis*[98]). In all cases, we tested the performance of *varKoder* on taxa

356    included in the training set and on taxa not included in the training set. We identified

357    perfect species identification (100% correct, 100% precision, 100% recall) for beetles and

358    coralroot orchids included in the training set. For bacteria, 5.6% of the validation set

359    returned ambiguous predictions; the remaining samples were correctly identified (94.7%

360    precision, 100% recall). In lichen-forming fungi, which include DNA from both the fungal

361    and algal partners, and thus are more challenging, 10% of the test samples returned

362    incorrect predictions and another 10% were inclusive; the remainder were correct (89%

363    precision, 80% recall). For all cases, species or varieties not included in the training set

364    generally resulted in inconclusive results, with a minority yielding incorrect predictions

365    **(Figure 7C)**. Precision and recall using varKodes instead of *rfCGRs* were very similar for all

366    four datasets.

367



368

369    **Figure 12.** varKoder performance in predicting taxonomy for all data on SRA. Sample sizes

370    refer to the number of validation accessions available for each combination of platform,

371    sequencing strategy and taxonomic rank.

372

373

374    Finally, we tested the scalability of varKodes in three large-scale datasets: (1) all 861

375    eukaryotic families with Illumina data on NCBI SRA, (2) all taxa with multiple accessions on

376    NCBI SRA, including different sequencing platforms and library strategies (254,819

377    accessions and 14,151 taxa across all taxonomic ranks), and (3) a previously published

378    dataset of 2916 soil eDNA samples from all seven continents[99]. Owing to NCBI download

379    speed bottlenecks, we restricted varKode construction to a very limited maximum of 10

380    Mbp of DNA data in the former 2 cases. The family-level eukaryote data achieved a rate of

381    correct predictions of 65.2–81.3% across all kingdoms when families were included in the

382    training set (**Figure** 7**D**), with most errors being inconclusive predictions (17.5–33.1%).

383    Precision varied from 95.3% to 97.3% and recall from 67.9% to 78.3%. Similarly to the

384    species- and variety-level exercise, families not included in the training set often yielded

385    inconclusive predictions **(Figure 7D)**, suggesting a potential for varKoding to be used as a

386    discovery tool when reasonably well-sampled training data sets are available. The

387    expanded data with all taxa from NCBI SRA revealed that varKoding is robust to sequencing

388    platform and library preparation method (**Figure 12**). Predictions at the family level or

389    pooled for all the taxonomic hierarchy are accurate regardless of sequencing details (>94%

390    precision, >86% recall). The much higher accuracy when compared to the dataset based on

391    Eukaryotic families alone may be an effect of a completely random validation set instead of

392    stratified by family, resulting in higher representation of commonly sampled families. At

393    the genus and species level, results are more dependent on the sequencing method (**Figure**

394    **12**). For genera, precision/recall using 10Mbp of data varies from 90.8%/90.8% with

395    whole genome shotgun libraries in PacBio to 97.9%/97.6% with genotype-by-sequencing

396    in Illumina. Finally, the eDNA data shows promise in using varKoding to identify the

397    geographical origin of an environmental sample (**Figure 13**): in the validation set, at

398    10Mbp of DNA data, 94.0% of the samples had continent correctly identified, with 2.6%

399    being incorrect, 1.9% being ambiguous, and 1.5% being inconclusive (84.7% prediction,

400    84.5% recall). Precision and recall using varKodes instead of *rfCGRs* were very similar for

401    both datasets.

402

**Figure 13.** Varkoder performance in identifying the geographical origin of a soil metabarcoding sample. A) Performance across the whole dataset. B) Performance for each continent.

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

421    NCBI SRA was the bottleneck, taking over 70 days. By parallelizing processing over 40

422    cores, processing this data into varKodes was about 10 times faster, resulting in

423    approximately 18GB of data for all of these accessions. Training a model on more than 1.3

424    million images took about 45 hours using only 2 GPUs. Therefore, a model with the millions

425    of species on Earth could be trained in just a few days in a dedicated server, provided that

426    sequence data to generate varKodes can be transferred at high speeds. Although training

427    on large datasets requires powerful GPUs and large memory, training on small datasets and

428    querying is possible on personal computers in a few seconds to minutes. To reduce the

429    computational resources required for training new datasets, we provide a pre-trained

430    model from both varKodes and rfCGRs from all taxa on SRA using the huggingface hub

431    (https://huggingface.co/brunoasm/vit_large_patch32_224.NCBI_SRA). See Asprino et al.[100]

432    for details on the data used for this model. Whenever the data become available, a model

433    potentially trained on millions of species can easily be ported to devices without

434    continuous internet access. Moreover, the minimal data amounts needed for identification

435    could be generated in seconds in a portable Nanopore device. Finally, the library

436    preparation method based on shotgun sequencing is very simple and can be automated

437    with portable consumer devices, such as the Nanopore Voltrax. Together, these properties

438    allow for more widely distributed applications of varKoding, such as field-laboratory

439    environments[101] or proposed distributed genetic databases[102].

440

441    **Conclusions**

442    varKoding is universal, accurate, efficient, and holds tremendous promise for documenting

443    and discovering Earth's biodiversity. It achieves accurate identification with minimal data

444    compared to existing next-generation sequencing methods, while maintaining universal

445    applicability across the Tree of Life. Its modular framework relying on widely used image

446    formats and machine learning frameworks can evolve alongside advances in sequencing

447    technologies, bioinformatics, and machine learning, as exemplified here by the update in

448    image representation (*varKodes* to *rfCGR*s) and neural network architecture (resnext to

449    ViT) after initial testing. For these reasons, we expect it will contribute for the wider

450    adoption of genomic signatures on biodiversity assessments and ecological research,

451    overcoming current challenges[39]. Reference data for varKoding will be increasingly

452    available from ambitious efforts in genome sequencing[103–107]. However, we note that

453    reference data for varKoding is much easier and cost-effective to obtain from low-coverage

454    genome skims than high-quality contiguous genomes: the robustness to minimal levels of

455    coverage a central advantage of our method. For example, our cost for a 3× skim of

456    herbarium samples is about $34 per sample, versus a high-quality genome which may cost

457    tens-of-thousands of dollars each. Thus, varKoding shows tremendous promise for further

458    automating species identification from natural history collections[108–110].

459

460    We expect that varKoding will be invaluable to the biodiversity science community in

461    numerous ways, with many avenues remaining to be explored. One of them is the

462    identification of samples with poor-quality and degraded DNA, such as unidentified

463    fragmentary fossil and subfossil remains in natural history collections[108,111]. For example,

464    Malpighiales samples with signs of DNA damage could be correctly identified using

465    *varKoder* to species or genus in many cases and to family in almost every case. Future

466    research could explore the lower limits of sample quality and sequence coverage to achieve

467    accurate identification at different divergence levels. Moreover, a promising avenue of

468    research is to identify the genomic features driving the success of identification based on

469    such low sequence coverage. It is possible that the changes in repeat patterns are more

470    important drivers of genomic evolution than currently appreciated[31,51–55]. Finally, we

471    expect that new neural network architectures and forms of DNA representation will

472    continue to be explored. One limitation of varKoding, as applied here, is the challenging

473    identification of samples within mixed components such as lichens or environmental DNA.

474    However, with long-read sequencing, *varKodes* and *rfCGR*s from single reads could

475    potentially include sufficient data for that end. Moreover, mixed samples could be useful for

476    other ends: varKodes could be used to classify a set of sequences based on any kind of

477    metadata, beyond taxonomy as demonstrated by our test on the geographical origin of a

478    soil sample.

# Author contributions

479

480    BASM conceived varKodes and wrote the program *varKoder*. BASM and CCD designed the
481    research. CCD, CA and XD designed sampling and lab methodology for the new sequence
482    data. CCD, XD, YY, LCM, and CA collected the new sequence data. BASM and PJF collated
483    datasets from published data. BASM, CCD, LC, YY and PJF analyzed and interpreted the data.
484    BASM, CCD, LCM and PF prepared the figures. BASM and CCD wrote the manuscript with
485    key contributions from LC, YY, CA and PJF. All authors approved the manuscript.

# Acknowledgments

486

# Online Methods

500

## Sequence data

501

502    *Taxon sampling, DNA sequencing, assembly, and annotation for newly acquired genetic*
503    *data*—The newly generated plant data used here and the methods to obtain these data are

504     described in detail in a data descriptor article[100]. Briefly, they included members of the

505     large and diverse order Malpighiales[34]: Malpighiaceae (251 accessions representing 31

506     genera), Elatinaceae (6 accession for 1 genus), and Chrysobalanaceae (30 accessions for 8

507     genera). Malpighiaceae includes *Stigmaphyllon* with the most comprehensive species

508     sampling: 10 species and 10 accessions sampled per species. All 100 *Stigmaphyllon*

509     samples were sequenced specifically to build, validate, and test our identification models at

510     shallower phylogenetic depths, since their taxonomy has been extensively revised by

511     coauthor C. Anderson[68,69]. Each of these samples was labeled with species, genus, and

512     family names. The focus for the remainder of the Malpighiaceae, Chrysobalanaceae, and

513     Elatinaceae sampling was to identify a given sample to genus. In this case, among the non-

514     *Stigmaphyllon* samples we included 3–9 species per genus. Each accession in this case was

515     labeled with its corresponding genus and family identification. Unlike *Stigmaphyllon,* where

516     we included multiple accessions per species, there were no additional replicates per

517     species for our genus-level sampling. For this dataset, we used leave-one-out cross

518     validation in all assessments, and therefore there are no train and validation sets. For

519     additional information see Asprino et al.[100] .

520     *Public genomic data compilation*—To further understand the versatility of varKodes more

521     broadly across the Tree of Life, we tested species identification using genome skim data

522     sets from four genera of plants, animals, fungi, and a bacterial species. For each of the four

523     organismal clades, we trained a multi-label model that included five species with at least

524     three samples per species. This involved a plant data set from coralroot orchids (genus

525     *Corallorhiza*)[97], with five species (or varieties) with at least five samples per species, except

526     for *C. striata* var. *vreelandii* and *C. striata* var. *striata,* for which we included six and seven

527     samples each, respectively. The animal data consisted of a beetle data set in the genus

528     *Bembidion*[54,95], which included five species with five samples per species. The fungal

529     dataset focused on a lichen-forming fungus in the genus *Xanthoparmelia*[96]. Since the

530     *Xanthoparmelia* species were paraphyletic, we subsampled only monophyletic groups for

531     model training. In this case, four species included three samples per species (*X.*

532     *camtschadalis*, *X. mexicana*, *X. neocumberlandia*, and *X. coloradoensis*) and one species

533     included five samples per species (*X. chlorochroa*). One potential confounding factor for the

534   *Xanthoparmelia* model is that *Xanthoparmelia* is a lichen-forming fungus and thus genome

535   skim data represents a chimera of fungal and algal genomes representing both partners in

536   this unique symbiosis. Species of the algal symbiont *Trebouxia* are flexible generalists

537   across fungal species *Xanthoparmelia*. Since these genome skims are a mix of both algal

538   photobiont and fungus, we hypothesize that the accuracy of our model decreased because

539   of the more generalist nature of *Trebouxia*[112]. Finally, the bacterial data set included clinical

540   isolates from *Mycobacterium tuberculosis*, the species of pathogenic bacteria that causes

541   tuberculosis[98]. We included representatives of five monophyletic *M. tuberculosis* lineages

542   (L1, L2, L3, L4.1.i1.2.1, and L4.3.i2) with seven clinical isolates per lineage. In all these

543   cases, we labeled samples with the lowest-level taxonomic identification available (species,

544   subspecies or isolates). For taxa with two or more samples available, 20% (with a

545   minimum of 1) were randomly selected for the validation set, which also included all taxa

546   represented by a single sample (therefore, absent from the training set). The remaining

547   accessions were used in the training set. See Asprino et al.[100] for further information.

548   We also compiled two broad datasets from the NCBI SRA. The first consists of all 861

549   eukaryotic families with sequenced under the Illumina platform from whole genome

550   shotgun (WGS) libraries and up to 10 Mbp of data (download date March 7, 2023). This

551   comprised 8,222 accessions, including families of animals (5,642 accessions, 1,426

552   families), plants (2,705 accessions, 401 families) and fungi (1,572 accessions, 363 families).

553   We labeled samples with family name only and included taxa with at least two associated

554   accessions in the training set. Our validation set consisted of 20% randomly selected

555   accessions from each family (with a minimum of one), plus all accessions in families with a

556   single accession available (therefore not part of the training set). Only eight of the 8,222

557   samples included yielded less than 10Mbp after sequence cleanup for varKode preparation,

558   and all at least 100 Kbp. The second broad-scale dataset includes all taxa on NCBI SRA that

559   could be represented by at least 3 independent accessions. In this case, we included data

560   amounts of up to 20 Mbp, different sequencing platforms (Illumina, PacBio, Nanopore,

561   BGIseq) and library preparation methods (whole genome shotgun, RADseq, GBS)

562   downloaded on January 9, 2024. For taxa with too many sequences available (such as

563   humans, crops, disease agents, etc.), we randomly chose up to 20 accessions for each

564    combination of sequencing platform and library preparation method. The resulting dataset

565    includes 253,820 accessions associated 28,636 taxonomic labels. In the training set,

566    97.52% of the accessions included 10Mbp of cleaned data, with the remainder having at

567    least 500Kbp. Accessions were labeled with all NCBI taxonomy ranks available (from infra-

568    specific taxa to domain), the library preparation method, and the sequencing platform. The

569    validation set, in this case, consisted of a random selection of 10% of all samples, not

570    stratified by taxon. For additional information, see Asprino et al.[100] .

571

572    Our final dataset was assembled with the aim to extend varKoder beyond taxonomic

573    identification.  We compiled a global soil metagenome eDNA dataset labeled with continent

574    of origin from Ma et al.[99] We filtered out any metagenomic sample which lacked

575    information on continent in the Ma et al. This yielded 2916 soil metagenome samples

576    across all seven continents.  We downloaded 10Mbp DNA data for each sample directly

577    from NCBI. All code used to download and analyze these data can be found in the GitHub

578    repository for our study ([https://www.github.com/brunoasm/varkoder_development](https://www.github.com/brunoasm/varkoder_development)).

579

# 580    varKode design and testing

581    *Sequence data preprocessing*—Prior to the construction of images, raw reads were lightly

582    cleaned using the following steps: identical reads were de-duplicated using *clumpify.sh* as

583    implemented in *BBtools*[84,113], adapters were removed, low-quality tails trimmed, and

584    overlapping read pairs merged using *fastp*[86] with options "--detect_adapter_for_pe", "--

585    dedup", "--dup_calc_accuracy 1", "--disable_quality_filtering", "--disable_length_filtering", "--

586    trim_poly_g", "--merge", "--include_unmerged", . Next, we randomly selected subsets of

587    cleaned reads with predefined data amounts, ranging from 500 kbp to 200 Mbp, with

588    *BBtools*. These data subsets were used to generate a variety of input varKodes for a single

589    sample and all such images were used for training (see main text and Figure 2A). Finally,

590    we applied *dsk*[85] to count k-mers of a given length based on clean raw reads (i. e. k-mers

591    are counted for each read and their frequencies are pooled across reads). *dsk* exhibits good

592    performance with low memory requirements, which is ideal for potential applications

593    using varKodes on low-memory devices. We note that analyses for species-level public

594    datasets have low computational requirements and were performed on an Apple MacBook

595    with ARM processor architecture.

596    *varKode and rfCGR construction*— We designed novel images—**varKodes**—that portray

597    relative frequencies of k-mers from low-coverage raw Illumina reads. These are similar to a

598    frequency chaos game representation (*fCGR*) *sensu* Jeffrey[53], but optimized for raw reads in

599    which sequence orientation is unknown, and therefore canonical k-mers and their reverse

600    complement are indistinguishable. This averaging of canonical k-mer frequencies and their

601    reverse complements is widely used in the context of raw reads[40,61,62,114,115]. We call these

602    images varKodes because they en*CODE* the *VAR*iation in k-mer frequencies in a sample. We

603    name our method **varKoding** after varKodes, but notice that it is modular and can use

604    other kinds of DNA image representation. They are meant to represent a genomic signature

605    by mapping k-mer identity to pixel position in an image, such that k-mers with more

606    similar composition are closer together. Additionally, the brightness of these pixels

607    represents the abundance of the associated k-mer, but we use ranks instead of raw

608    frequencies to decrease the effect of overabundant and artifactual k-mers. In summary,

609    varKodes are produced by mapping k-mer counts onto a pre-computed map of k-mers to

610    pixels, and transforming frequency data to pixel brightness. varKode design employed t-

611    SNE[116] and the python libraries *numpy*[88] and *pillow*[117]. In addition to varKodes, here we

612    also developed a new image representation that uses the same pixel mapping as *fCGRs* but

613    represents k-mer abundance as ranks instead of raw frequencies. We named these ranked

614    frequency chaos game representation (*rfCGR*). Both varKodes and *fCGRs* are saved as 8-bit

615    PNG images including labels as exif metadata.

616    *Testing k-mer length and data amount*—We employed *fastai*[89] for, a high-level

617    implementation of neural networks based on *pytorch*[64] for training and prediction. All the

618    model architectures we applied are image classification models available from the *timm*

619    library[90], which have been widely tested using a variety of image types. To identify the

620    optimal training hyperparameters for our neural network, we conducted a series of tests

621    using the species-level data set for the genus *Stigmaphyllon*. We generated varKodes for

622    each of the *Stigmaphyllon* samples. We first tested the joint effect of k-mer length and input

623    data amount for neural network classification accuracy by selecting three samples per

624    species as a validation set; the remaining samples were used to train neural networks using

625    different amounts of input data across 10 randomly generated training sets. As input data

626    for both the validation and training sets, we randomly subsampled the original sequences

627    into fastq files containing from 500 Kb to 200 Mb (equivalent to about 1,700 to 670,000

628    2x150bp Illumina reads). In this test, we only included samples that yielded at least 200

629    million base pairs after cleaning. We also tested the effect of including images for all data

630    amounts during training. For each replicate, we applied the widely used image

631    classification neural network *resnet50* architecture[118] to classify varKodes and trained

632    models for 30 epochs. We visualized the distribution of validation accuracy for each

633    combination of input data amount and k-mer lengths to find a good balance between both.

634    Visualizations and code applied for training and evaluation is available in our GitHub

635    repository (https://www.github.com/brunoasm/varkoder_development).

636    *Neural network optimization*—After identifying an appropriate k-mer length and input data

637    used to produce varKodes (**Figure** 3), we next tested a series of neural network training

638    conditions. We varied the neural network model complexity, choosing from seven

639    commonly used architectures: *resnet50*[118], *resnet-D*[70] with different depths (18, 50, 101), a

640    wide *resnet50*[70], *efficientnet-B4*[119], and ResNeXt101[77]. We also tested the effect of the

641    following: random initial weights vs. pretrained weights from the *timm* library[90], presence

642    or absence of  lighting transforms, presence or absence of label smoothing, and presence or

643    absence of augmentation strategies (i.e., *CutMix*[76] or *MixUp*[75]). Because these parameters

644    may have complex interactions, we tested all combinations of architecture, pretraining,

645    transforms, label smoothing, and augmentation, with 20 replicates for each combination of

646    conditions. In each replicate, we randomly chose 20% of the samples for each species of

647    *Stigmaphyllon* as validation and trained the model using the remainder for 30 epochs.

648    Training was performed using all varKodes available for each sample (from 500kbp to

649    200Mbp). For validation, we separately evaluated whether each varKode with a different

650    amount of data was correctly identified. For each replicate and amount of data used to

651    validate varKodes, we recorded the average validation accuracy across the validation set.

652    We then applied a linear model to predict the effect of all training parameters and amount

653    of data in varKodes in the validation set on validation accuracy. Validation accuracy in this

654    case was arc-sin transformed for linear modling due to its bounded range of 0–1. We

655    started from the full model containing all parameters and their interactions and reduced

656    the model step-wise based on AIC scores (i. e. Akaike Information Criteria), as implemented

657    in the R function step. Visualizations and code applied for training and evaluation is

658    available in our GitHub repository

659    (https://www.github.com/brunoasm/varkoder_development).

660

661    *Testing sample number requirements*—A legitimate concern with complex neural networks

662    is that they may require vast amounts of training data and that typical skimming data sets

663    might be insufficient for them to be useful. We tested the robustness of our models to the

664    effect of the number of samples per species included in training by using from one to seven

665    samples per species as training set and the remaining as validation, with 50 replicates per

666    number of training samples. The batch size used in training was adjusted for the cases with

667    very few samples included, so that each training epoch included about 10 batches. We

668    included varKodes from 1Mbp to 200Mbp in both training and validation sets. In this case,

669    we applied the training parameters informed by our previous analyses: a *resnext101*

670    architecture, random initial weights, *CutMix* augmentation, and label smoothing for 30

671    epochs. We visualized the effect of the number of samples by plotting the average

672    validation accuracy of each sample against the number of training samples used in each

673    case. Visualizations and code applied for training and evaluation is available in our GitHub

674    repository (https://www.github.com/brunoasm/varkoder_development).

675

676    *Testing the effect of data quality*—Most of the cases with low accuracy corresponded to

677    samples with low DNA yield (**Figure** 3**B**). We identified that DNA extraction yield was

678    significantly correlated with two metrics of DNA quality: average insert size and variation

679    in nucleotide composition along reads[80] (**Figure 4**). *varKodes* produced from these samples

680   may be visually distinct from other samples of the same species (**Figure 5**). For this reason,

681   we further tested whether sample quality in training or validation impacted accuracy.

682   Using both quality metrics, we identified the five lowest quality samples for each species.

683   We next produced training sets using six randomly chosen samples per species, varying the

684   number of low-quality samples included in training from zero to four. We included

685   varKodes from 1Mbp to 200Mbp in both training and validation sets. We repeated this for

686   30 replicates for each number of low-quality samples. Like our tests with varying sample

687   numbers, we applied the following training parameters: a *resnext101* architecture, random

688   initial weights, *CutMix* augmentation, label smoothing for 30 epochs. For the validation set,

689   we separately recorded the accuracy for high- and low-quality samples. We then visualized

690   the effect of inclusion of low-quality samples in the training set by observing the

691   distribution of validation accuracies for high-quality and low-quality samples across the

692   range of number of low-quality samples included in the training set. Visualizations and

693   code applied for training and evaluation is available in our GitHub repository

694   ([https://www.github.com/brunoasm/varkoder_development](https://www.github.com/brunoasm/varkoder_development)).

695

696   *Implementation of varKoder*—Following all the tests described above, we implemented the

697   optimal neural network training strategies in a python program named ***varKoder***.

698   *varKoder* can process, train and query varKodes and is freely available on our GitHub:

699   https://github.com/brunoasm/varKoder. Because it employs standard neural network

700   frameworks (namely, *pytorch*[64], *fastai*[89], and *timm*[90]), any of the image classification models

701   and training hyperparamenters available now or in the future via these libraries can be

702   easily adapted and applied to varKode classification. Moreover, we have implemented a

703   multi-label model as the default to increase robustness to low-quality varKodes with little

704   diagnostic information in the training set. This was done by using an asymmetric multi-

705   label loss function[82] instead of the standard cross-entropy loss function used in single-label

706   classification. Analyses used development versions of *varKoder* starting with v.0.8.0.

707   Improvements suggested during the peer-review process are now implemented in

708   *varKoder* v.1.1.0.

# *varKoder* evaluation and comparison to alternatives

709  *varKoder*—To test *varKoder* performance on a complex dataset spanning multiple

710  taxonomic levels and varying phylogenetic depths, we used the Malpighiales dataset

711  including genera in Elatinaceae, Chrysobalanaceae and Malpighiaceae. Species of

712  *Stigmaphyllon* (Malpighiaceae) were labeled with species, genus, and family names; all

713  other samples were labeled with genus and family names. We tested the performance of

714  *varKoder* in each sample with leave-one-out cross-validation. For each sample, we retained

715  it as validation and trained a neural network using all the other samples. In preliminary

716  assessments, we found that a ViT[72] architecture combined with a multi-label model

717  sometimes led to instability in training for some datasets. For that reason, we used a two-

718  step approach. Models first were pre-trained for 20 epochs as single-label, using the least

719  inclusive taxonomic assignment available for each sample and a base learning rate of 0.05.

720  Next, we trained for an additional 10 epochs using the pre-trained weights but with a much

721  smaller learning rate (0.005) and a multi-label output. Training samples included varKodes

722  from 500 Kbp to 200 Mbp, and we recorded validation accuracy separately for varKodes

723  produced from each amount of data. We used an arbitrary confidence threshold of 0.7 to

724  make predictions in the multilabel models. For validation samples, we deemed a prediction

725  correct if only the correct taxon was predicted for each taxonomic rank (i.e., species, genus,

726  family). We deemed a prediction incorrect if one or more predictions passed the threshold

727  for a taxonomic rank, but none match the actual label. When predicted labels included both

728  the correct and incorrect taxa, we deemed it ambiguous. If the output prediction included

729  no taxon with confidence above the threshold, we considered it as inconclusive. As metrics

730  across all samples, we used prediction and recall, averaged across all predictions. We

731  visualized the fraction of correct, incorrect, ambiguous, and inconclusive samples for each

732  taxonomic rank and each amount of data used to produce varKodes. The code to reproduce

733  training conditions and evaluation tests is available on GitHub

734  (https://www.github.com/brunoasm/varkoder_development).

735  To test the joint effect of neural network architecture and image representation method,

736  we applied this cross-validation approach to all combinations of three image

738    representations and four neural network architectures. The architectures tested included:

739    (1) *ResNeXt101*[77], the optimal convolutional neural network architecture in our initial tests,

740    (2) *ViT*[72], a transformer-based architecture that became available after our initial testing,

741    (3) a neural network with two convolutional layers processing vectorized k-mer counts,

742    following Fiannaca et al[44] and (4) a multi-layer perceptron formed by a series of fully

743    connected layers as specified in Millán Arias et al[42]. The two latter have been previously

744    employed for *fCGR* data. The three representations tested include *varKodes* and *rfCGRs* as

745    developed here, and *fCGRs* as estimated by iDeLUCS[93]. In the latter case, we used iDeLUCS

746    functions to produce *fCGRs* as 2D python arrays of k-mer counts. Next, we rescaled these

747    counts to the range of 0–255 and rounded them to the nearest integer. These arrays were

748    then saved as 8-bit png images. In all tests, we employed the same data augmentation

749    methods and loss function as for *varKodes* and *rfCGRs*. All code used in *varKoder* analyses is

750    available on GitHub (https://www.github.com/brunoasm/varkoder_development)**.**

751    *Skmer*—To compare *varKoder* with alternative methods, we used fastq files cleaned and

752    subsampled by *varKoder* as input files to *Skmer*. In this case, we also used leave-one-out

753    cross-validation to evaluate performance. For each amount of input data (500Kbp to

754    200Mbp), we cycled through all samples, constructing a *Skmer* database with the "*skmer*

755    *reference*" command and including all samples but one and default settings. We then used

756    the "*skmer query*" command with default settings on the sample left out and deemed the

757    identification as correct if the sample in the reference database with closest estimated

758    genetic distance had the correct taxon label. Because *Skmer* could always query a sample

759    and there is no objective criterion to consider matches beyond the best match, the output

760    predictions can only be correct or incorrect, but not inconclusive or ambiguous. We

761    visualized the results similarly as we did with *varKoder*. The code to reproduce *Skmer*

762    analyses is available on GitHub

763    (https://www.github.com/brunoasm/varkoder_development)**.**

764    *Conventional plant barcodes* —For conventional barcodes, we applied standard BLAST- and

765    phylogeny-based methods, which do not involve machine learning. To infer phylogenies

766    from our genome skim data (Figure 1), we applied the *PhyloHerb* bioinformatic pipeline[120],

767    which has been applied recently to a taxa ranging from algae to flowering plants[121–123].

768    Briefly, this pipeline works as follows: for plastid loci, *PhyloHerb* maps raw short reads to a

769    database of land plant plastid genomes. Mapped reads are then assembled into scaffolds

770    using *SPAdes*[124] and plastid loci are identified using nucleotide BLAST searches with a

771    default e-value threshold of 1e-40. *PhyloHerb* then outputs orthologous plastid genes into

772    individual FASTA files, which are fed directly into MAFFT v7.407[125] for alignment.

773    Alignments are then concatenated into a super matrix using the 'conc' function within the

774    *PhyloHerb* package. Phylogenies for both individual locus and the concatenated alignment

775    were inferred with IQTREE v2.0.6 using the GTR+GAMMA model with 1000 ultrafast

776    bootstrap replicates[126].

777    To recover the conventional plant barcodes, *rbc*L, *mat*K, *trn*L-F, *ndh*F, and ITS, from our

778    Malpighiales genome skim data, we applied GetOrganelle v1.7.7.0[127] and *PhyloHerb*

779    v1.1.1[120] to automatically assemble and extract these DNA markers, respectively. Briefly,

780    the complete or subsampled genome skim data were first assembled into plastid genomes

781    or nuclear ribosomal regions using *GetOrganelle* with its default settings. Next, *PhyloHerb*

782    was applied to extract the relevant barcode genes using its built-in BLAST database. To test

783    whether these traditional barcodes provided accurate identification to species, genus, and

784    family, we ran an all-by-all BLASTn analysis for each individual gene across the same data

785    subsampling schemes as *Skmer* and *varKoder*. BLAST targets were always drawn from

786    assemblies using all the data available for each specimen, whereas queries included

787    assemblies from input data amounts varying from 500 Kbp to 200 Mbp. Within each BLAST

788    analysis for each one of the Malpighiales accessions, we deemed an identification to be

789    correct if the best non-self BLAST hit came from the same taxon, and incorrect otherwise.

790    We deemed it inconclusive if the locus could not be assembled for that amount of data. For

791    concatenated barcodes, we produced a phylogenetic tree for each amount of data and

792    deemed an identification to be correct if the sample with lowest patristic distance came

793    from the same taxon. We deemed it to be inconclusive when none of the genes in the

794    concatenated dataset could be assembled for a sample. We visualized results similarly to

795    *varKoder*, separately for each conventional barcoding gene and for the concatenated

796     dataset. The code to reproduce conventional barcode analyses is available on GitHub

797     (https://www.github.com/brunoasm/varkoder_development).

798     *iDeLUCS*—To evaluate the performance of *varKoder* with another deep learning based

799     sequence classifier, we applied the sequences assembled from the *PhyloHerb* pipeline to

800     *iDeLUCS*[93]. We first used concatenated sequences of five traditional plant barcodes (*rbc*L,

801     *mat*K, *trn*L-F, *ndh*F, and ITS) assembled from input reads varying from 500 Kbp to 200

802     Mbp. *iDeLUCS* was run with k-mer length of 6, 100 training epochs, 100 data augmentations

803     per sequence, and the SGD algorithm for neural network optimization. Unlike varKoder,

804     iDeLUCS does unsupervised clustering and therefore does not use labels during training.

805     Instead, all input accessions were set to be clustered into 10 groups (equal to the total

806     number of species) and the accuracy was evaluated with the *cluster_acc* function

807     implemented in *iDeLUCS.* We also applied the entire plastid genome and the nuclear

808     ribosomal sequence assemblies (ETS+18S+ITS1+5.8S+ITS2+28S) in *iDeLUCS* with the same

809     parameters to evaluate the impact of input data quality.

## 810    Application in diverse taxa

811     *Species-level identification in plants, animals, fungi, and bacteria*— For all four test cases

812     (*Corallorhiza*, *Bembidion*, *Xanthoparmelia*, and *Mycobacterium tubercolosum*), we applied

813     default *varKoder* v.0.8.0 parameters for generating *rfCGR* images, training each model, and

814     testing the accuracy of the trained model using the 'query' function. In all cases, we

815     included all the available data for each training or validation sample. To test if trained

816     models accurately predicted species identity, we queried them using extra genome skim

817     samples not used for training but from the same species included in the model. We also

818     tested genome skim test samples of species within the same genus *not* used in model

819     training. As in the case of Malpighiales, we set the threshold to make a prediction equal to

820     0.7 and used the same criteria to consider a prediction correct, incorrect, inconclusive, or

821     ambiguous. We separately evaluated results for taxa with representatives included in the

822     training set and taxa used only as queries, without conspecific samples in the training set.

823     The code to reproduce these analyses is available on GitHub

824     (https://www.github.com/brunoasm/varkoder_development).

825

826    *All eukaryotic families data set from SRA*—Each accession was labeled with its family

827    identification obtained from NCBI. Because of the larger size of this dataset, a leave-one-out

828    cross-validation approach would have been intractable. Therefore, we randomly selected

829    80% of the samples in each family as the training set and used the remainder for validation.

830    Similarly to Malpighiales, we used a two-step training method by pre-training as a single-

831    label model and finalizing with a multi-label model. Pre-training was done with a learning

832    rate of 0.1 and a batch size of 300 for 30 epochs. Final training was done with the same

833    batch size but a smaller base learning rate of 0.01 in 5 epochs with frozen body weights and

834    three epochs with unfrozen weights. The code to reproduce these analyses is available on

835    GitHub (https://www.github.com/brunoasm/varkoder_development).

836

837    *All taxa from SRA*— For each accession, we created *rfCGRs* from 500Kbp to 10Mbp of data.

838    Each accession was labeled with all the taxa in its taxonomic tree (that is, from infra-

839    specific taxa to domains of life), as well as library strategy (RAD, GBS or WGS) and

840    sequencing platform (Illumina, PACBIO, Nanopore or BGISEQ). We randomly selected 10%

841    of the samples as validation set, and eliminated from validation samples all labels absent

842    from the training set. We used a two-step training method. First, we pre-trained using a

843    single-label strategy, using as labels the concatenation of library strategy, sequencing

844    platform, kingdom, family and genus. For pretraining, we used a learning rate of 0.1, a

845    batch size of 500 and 30 epochs. We then used the weights of this pre-trained model as

846    starting weights for a multi-label model including all labels. We trained the model for

847    additional 50 epochs with unfrozen body weights and 10 epochs with frozen weights,

848    learning rate of 0.05 and batch size of 600. The code to reproduce these analyses is

849    available on GitHub (https://www.github.com/brunoasm/varkoder_development).

850

851    *Environmental metagenome global identification*—The downloaded soil metagenomes from

852    Ma et al.[99] were labeled by source continent. Similarly to the eukaryotic family data set

853    from SRA, we randomly selected 80% of the samples as the training set and used the

854    remaining 20% as the validation set. We used a two-step training method by pre-training

855    as a single-label model and finalizing with a multi-label model. Pre-training was done with

856    a learning rate of 0.1 and a batch size of 64 for 30 epochs. Final training was done with the

857    same batch size but a smaller base learning rate of 0.01 in 5 epochs with frozen body

858    weights and three epochs with unfrozen weights. The code to reproduce all these analyses

859    is available on GitHub (https://www.github.com/brunoasm/varkoder_development).

# References

861    1. Hebert, P. D. N., Ratnasingham, S. & de Waard, J. R. Barcoding animal life: cytochrome c
862    oxidase subunit 1 divergences among closely related species. Proc. R. Soc. Lond. B Biol. Sci.
863    270, S96–S99 (2003).
864
865    2. Kress, W. J. Plant DNA barcodes: Applications today and in the future. J. Syst. Evol. 55,
866    291–307 (2017).
867
868    3. Ratnasingham, S. & Hebert, P. D. N. BOLD: The Barcode of Life Data System
869    (www.barcodinglife.org). Mol. Ecol. Notes 7, 355–364 (2007).
870
871    4. Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C. & Willerslev, E. Towards next-
872    generation biodiversity assessment using DNA metabarcoding. Mol. Ecol. 21, 2045–2050
873    (2012).
874
875    5. Seifert, K. A. Progress towards DNA barcoding of fungi. Mol. Ecol. Resour. 9 Suppl s1, 83–
876    89 (2009).
877
878    6. Sharkey, M. J. et al. Minimalist revision and description of 403 new species in 11
879    subfamilies of Costa Rican braconid parasitoid wasps, including host records for 219
880    species. ZooKeys 1013, 1–665 (2021).
881
882    7. Lahaye, R. et al. DNA barcoding the floras of biodiversity hotspots. Proc. Natl. Acad. Sci. U.
883    S. A. 0709936105 (2008) doi:10.1073/pnas.0709936105.
884
885    8. Kuzmina, M. L. et al. Using herbarium-derived DNAs to assemble a large-scale DNA
886    barcode library for the vascular plants of Canada. Appl. Plant Sci. 5, apps.1700079 (2017).
887
888    9. Muñoz-Rodríguez, P. et al. A taxonomic monograph of Ipomoea integrated across
889    phylogenetic scales. Nat. Plants 5, 1136–1144 (2019).
890
891    10. Hebert, P. D., Penton, E. H., Burns, J. M., Janzen, D. H. & Hallwachs, W. Ten species in one:
892    DNA barcoding reveals cryptic species in the neotropical skipper butterfly Astraptes
893    fulgerator. Proc. Natl. Acad. Sci. U. S. A. 101, 14812–14817 (2004).
894

895    11. Zeale, M. R., Butlin, R. K., Barker, G. L., Lees, D. C. & Jones, G. Taxon-specific PCR for DNA
896    barcoding arthropod prey in bat faeces. Mol. Ecol. Resour. 11, 236–244 (2011).
897
898    12. Nitta, J. H., Meyer, J., Taputuarai, R. & Davis, C. C. Life cycle matters: DNA barcoding
899    reveals contrasting community structure between fern sporophytes and gametophytes.
900    Ecol. Monogr. 87, 278–296 (2016).
901
902    13. Kress, W. J. et al. Plant DNA barcodes and a community phylogeny of a tropical forest
903    dynamics plot in Panama. Proc. Natl. Acad. Sci. U. S. A. 106, 18621–18626 (2009).
904
905    14. Willis, C. G., Franzone, B. F., Xi, Z. & Davis, C. C. The establishment of Central American
906    migratory corridors and the biogeographic origins of seasonally dry tropical forests in
907    Mexico. Front. Genet. 5, 433 (2014).
908
909    15. Willerslev, E. et al. Ancient biomolecules from deep ice cores reveal a forested Southern
910    Greenland. Science 317, 111–114 (2007).
911
912    16. Crump, S. E. et al. Ancient plant DNA reveals High Arctic greening during the Last
913    Interglacial. Proc. Natl. Acad. Sci. U. S. A. 118, e2019069118 (2021).
914
915    17. Kjær, K. H. et al. A 2-million-year-old ecosystem in Greenland uncovered by
916    environmental DNA. Nature 612, 283–291 (2022).
917
918    18. Fierer, N. et al. Forensic identification using skin bacterial communities. Proc. Natl.
919    Acad. Sci. 107, 6477–81 (2010).
920
921    19. Rollo, F., Ubaldi, M., Ermini, L. & Marota, I. Ötzi's last meals: DNA analysis of the
922    intestinal content of the Neolithic glacier mummy from the Alps. Proc. Natl. Acad. Sci. U. S.
923    A. 99, 12594–12599 (2002).
924
925    20. Yu, J. et al. Progress in the use of DNA barcodes in the identification and classification of
926    medicinal plants. Ecotoxicol. Environ. Saf. 208, 111691 (2021).
927
928    21. Ashfaq, M. & Hebert, P. D. N. DNA barcodes for bio-surveillance: regulated and
929    economically important arthropod plant pests. Genome 59, 933–945 (2016).
930
931    22. Eaton, M. J. et al. Barcoding bushmeat: molecular identification of Central African and
932    South American harvested vertebrates. Conserv. Genet. 11, 1389–1404 (2010).
933
934    23. Liu, J. et al. Integrating a comprehensive DNA barcode reference library with a global
935    map of yews (Taxus L.) for forensic identification. Mol. Ecol. Resour. 18, 1115–1131 (2018).
936
937    24. Ogden, R., Dawnay, N. & McEwing, R. Wildlife DNA forensics—bridging the gap between
938    conservation genetics and law enforcement. Endanger. Species Res. 9, 179–195 (2009).
939

940   25. Williamson, J. et al. Exposing the illegal trade in cycad species (Cycadophyta:
941   Encephalartos) at two traditional medicine markets in South Africa using DNA barcoding.
942   Genome 59, 771–781 (2016).
943
944   26. Costa, F. O. & Carvalho, G. R. The Barcode of Life Initiative: synopsis and prospective
945   societal impacts of DNA barcoding of Fish. Genomics Soc. Policy 3, 29 (2007).
946
947   27. Gao, Z., Liu, Y., Wang, X., Wei, X. & Han, J. DNA mini-barcoding: a derived barcoding
948   method for herbal molecular identification. Front. Plant Sci. 10, (2019).
949
950   28. Molina, J. et al. Possible loss of the chloroplast genome in the parasitic flowering plant
951   Rafflesia lagascae (Rafflesiaceae). Mol. Biol. Evol. 31, 793–803 (2014).
952
953   29. Cai, L. et al. Deeply altered genome architecture in the endoparasitic flowering plant
954   Sapria himalayana Griff. (Rafflesiaceae). Curr. Biol. 31, 1002-1011.e9 (2021).
955
956   30. Richardson, J. E., Pennington, R. T., Pennington, T. D. & Hollingsworth, P. M. Rapid
957   diversification of a species-rich genus of neotropical rain forest trees. Science 293, 2242–
958   2245 (2001).
959
960   31. Wang, J., Luo, J., Ma, Y.-Z., Mao, X.-X. & Liu, J.-Q. Nuclear simple sequence repeat markers
961   are superior to DNA barcodes for identification of closely related Rhododendron species on
962   the same mountain. J. Syst. Evol. 57, 278–286 (2019).
963
964   32. Su, X., Wu, G., Li, L. & Liu, J. Species delimitation in plants using the Qinghai–Tibet
965   Plateau endemic Orinus (Poaceae: Tridentinae) as an example. Ann. Bot. 116, 35–48
966   (2015).
967
968   33. Lu, Z. et al. Species delimitation and hybridization history of a hazel species complex.
969   Ann. Bot. 127, 875–886 (2021).
970
971   34. Cai, L. et al. The perfect storm: gene tree estimation error, incomplete lineage sorting,
972   and ancient gene flow explain the most recalcitrant ancient angiosperm clade, Malpighiales.
973   Syst. Biol. 70, 491–507 (2021).
974
975   35. Clarke, L. J., Soubrier, J., Weyrich, L. S. & Cooper, A. Environmental metabarcodes for
976   insects: in silico PCR reveals potential for taxonomic bias. Mol. Ecol. Resour. 14, 1160–1170
977   (2014).
978
979   36. Song, H., Buhay, J. E., Whiting, M. F. & Crandall, K. A. Many species in one: DNA
980   barcoding overestimates the number of species when nuclear mitochondrial pseudogenes
981   are coamplified. Proc. Natl. Acad. Sci. U. S. A. 105, 13486–13491 (2008).
982
983   37. Xiong, H. et al. Species tree estimation and the impact of gene loss following whole-
984   genome duplication. Syst. Biol. 71, 1348–1361 (2022).
985

986   38. Straub, S. C. K. et al. Navigating the tip of the genomic iceberg: Next-generation
987   sequencing for plant systematics. Am. J. Bot. 99, 349–364 (2012).
988
989   39. Bohmann, K., Mirarab, S., Bafna, V. & Gilbert, M. T. P. Beyond DNA barcoding: The
990   unrealised potential of genome skim data in sample identification. Mol. Ecol. 1–14 (2020)
991   doi:10.1111/mec.15507.
992
993   40. Sarmashghi, S., Bohmann, K., P. Gilbert, M. T., Bafna, V. & Mirarab, S. Skmer: assembly-
994   free and alignment-free sample identification using genome skims. Genome Biol. 20, 34
995   (2019).
996
997   41. Borowiec, M. L. et al. Deep learning as a tool for ecology and evolution. Methods Ecol.
998   Evol. 13, 1640–1660 (2022).
999
1000  42. Arias, P. M., Alipour, F., Hill, K. A. & Kari, L. DeLUCS: Deep learning for unsupervised
1001  clustering of DNA sequences. PLOS ONE 17, e0261531 (2022).
1002
1003  43. Kari, L. et al. Mapping the space of genomic signatures. PLOS ONE 10, e0119815 (2015).
1004
1005  44. Fiannaca, A. et al. Deep learning models for bacteria taxonomic classification of
1006  metagenomic data. BMC Bioinformatics 19, (2018).
1007
1008  45. Linard, B., Swenson, K. & Pardi, F. Rapid alignment-free phylogenetic identification of
1009  metagenomic sequences. Bioinformatics (2019) doi:10.1093/bioinformatics/btz068.
1010
1011  46. Desai, H. P., Parameshwaran, A. P., Sunderraman, R. & Weeks, M. Comparative Study
1012  Using Neural Networks for 16S Ribosomal Gene Classification. J. Comput. Biol. 27, 248–258
1013  (2020).
1014
1015  47. Shang, J. & Sun, Y. CHEER: HierarCHical taxonomic classification for viral mEtagEnomic
1016  data via deep leaRning. Methods 189, 95–103 (2021).
1017
1018  48. Arias, P. M. et al. BarcodeBERT: Transformers for Biodiversity Analysis. Preprint at
1019  http://arxiv.org/abs/2311.02401 (2023).
1020
1021  49. Badirli, S., Akata, Z., Mohler, G., Picard, C. & Dundar, M. Fine-Grained Zero-Shot Learning
1022  with DNA as Side Information. Preprint at http://arxiv.org/abs/2109.14133 (2021).
1023
1024  50. Lecun, Y., Bengio, Y. & Hinton, G. Deep learning. Nature 521, 436–444 (2015).
1025
1026  51. Cong, Y., Ye, X., Mei, Y., He, K. & Li, F. Transposons and non-coding regions drive the
1027  intrafamily differences of genome size in insects. iScience 25, 104873 (2022).
1028
1029  52. Heckenhauer, J. et al. Genome size evolution in the diverse insect order Trichoptera.
1030  GigaScience 11, 1–19 (2022).
1031

1032    53. Schley, R. J. et al. The ecology of palm genomes: repeat-associated genome size
1033    expansion is constrained by aridity. New Phytol. 433–446 (2022) doi:10.1111/nph.18323.
1034
1035    54. Sproul, J. S., Barton, L. M. & Maddison, D. R. Repetitive DNA profiles Reveal Evidence of
1036    Rapid Genome Evolution and Reflect Species Boundaries in Ground Beetles. Syst. Biol. 0, 1–
1037    12 (2020).
1038
1039    55. de Medeiros, B. A. S. & Farrell, B. D. Whole-genome amplification in double-digest
1040    RADseq results in adequate libraries but fewer sequenced loci. PeerJ 6, e5089 (2018).
1041
1042    56. Jeffrey, H. J. Chaos game representation of gene structure. Nucleic Acids Res. 18, 2163–
1043    2170 (1990).
1044
1045    57. Deschavanne, P. J., Giron, A., Vilain, J., Fagot, G. & Fertil, B. Genomic signature:
1046    characterization and classification of species assessed by chaos game representation of
1047    sequences. Mol. Biol. Evol. 16, 1391–1399 (1999).
1048
1049    58. de la Fuente, R., Díaz-Villanueva, W., Arnau, V. & Moya, A. Genomic Signature in
1050    Evolutionary Biology: A Review. Biology 12, 322 (2023).
1051
1052    59. Avila Cartes, J., Anand, S., Ciccolella, S., Bonizzoni, P. & Della Vedova, G. Accurate and fast
1053    clade assignment via deep learning and frequency chaos game representation. GigaScience
1054    12, giac119 (2023).
1055
1056    60. Solis-Reyes, S., Avino, M., Poon, A. & Kari, L. An open-source k-mer based machine
1057    learning tool for fast and accurate subtyping of HIV-1 genomes. PLOS ONE 13, e0206409
1058    (2018).
1059
1060    61. Kislyuk, A., Bhatnagar, S., Dushoff, J. & Weitz, J. S. Unsupervised statistical clustering of
1061    environmental shotgun sequences. BMC Bioinformatics 10, 316 (2009).
1062
1063    62. Arias, P. M. et al. Environment and taxonomy shape the genomic signature of
1064    prokaryotic extremophiles. Sci. Rep. 13, 16105 (2023).
1065
1066    63. Murad, T., Ali, S., Khan, I. & Patterson, M. Spike2CGR: an efficient method for spike
1067    sequence classification using chaos game representation. Mach. Learn. 112, 3633–3658
1068    (2023).
1069
1070    64. Paszke, A. et al. PyTorch: An Imperative Style, High-Performance Deep Learning
1071    Library. in Advances in Neural Information Processing Systems 32 8024–8035 (Curran
1072    Associates, Inc., 2019).
1073
1074    65. Davis, C. C. & Anderson, W. R. A complete generic phylogeny of Malpighiaceae inferred
1075    from nucleotide sequence data and morphology. Am. J. Bot. 97, 2031–2048 (2010).
1076

1077    66. Cai, L. et al. Phylogeny of Elatinaceae and the tropical Gondwanan origin of the
1078    Centroplacaceae (Malpighiaceae, Elatinaceae) clade. PLOS ONE 11, e0161881 (2016).
1079
1080    67. Davis, C. C., Anderson, W. R. & Donoghue, M. J. Phylogeny of Malpighiaceae: evidence
1081    from chloroplast ndhF and trnL-F nucleotide sequences. Am. J. Bot. 88, 1830–1846 (2001).
1082
1083    68. Anderson, C. Revision of Ryssopterys and transfer to Stigmaphyllon (Malpighiaceae).
1084    Blumea 56, 73–104 (2011).
1085
1086    69. Anderson, C. Monograph of Stigmaphyllon (Malpighiaceae). Syst. Bot. Monogr. 51, 1–
1087    313 (1997).
1088
1089    70. He, T. et al. Bag of Tricks for Image Classification with Convolutional Neural Networks.
1090    Preprint at http://arxiv.org/abs/1812.01187 (2018).
1091
1092    71. Vaswani, A. et al. Attention Is All You Need. Preprint at
1093    https://doi.org/10.48550/arXiv.1706.03762 (2017).
1094
1095    72. Dosovitskiy, A. et al. An Image is Worth 16x16 Words: Transformers for Image
1096    Recognition at Scale. (2021).
1097
1098    73. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the Inception
1099    Architecture for Computer Vision. in 2016 IEEE Conference on Computer Vision and
1100    Pattern Recognition (CVPR) 2818–2826 (2016). doi:10.1109/CVPR.2016.308.
1101
1102    74. Smith, L. N. A disciplined approach to neural network hyper-parameters: Part 1 --
1103    learning rate, batch size, momentum, and weight decay. Preprint at
1104    http://arxiv.org/abs/1803.09820 (2018).
1105
1106    75. Zhang, H., Cisse, M., Dauphin, Y. N. & Lopez-Paz, D. mixup: Beyond Empirical Risk
1107    Minimization. Preprint at http://arxiv.org/abs/1710.09412 (2018).
1108
1109    76. Yun, S. et al. CutMix: Regularization Strategy to Train Strong Classifiers with Localizable
1110    Features. Preprint at https://doi.org/10.48550/arXiv.1905.04899 (2019).
1111
1112    77. Xie, S., Girshick, R., Dollár, P., Tu, Z. & He, K. Aggregated Residual Transformations for
1113    Deep Neural Networks. Preprint at https://doi.org/10.48550/arXiv.1611.05431 (2017).
1114
1115    78. Goodfellow, I., Bengio, Y. & Courville, A. Deep Learning. (MIT Press, 2016).
1116
1117    79. Christin, S., Hervet, É. & Lecomte, N. Applications for deep learning in ecology. Methods
1118    Ecol. Evol. 10, 1632–1644 (2019).
1119
1120    80. Weiß, C. L. et al. Temporal patterns of damage and decay kinetics of DNA retrieved from
1121    plant herbarium specimens. R. Soc. Open Sci. 3, 160239 (2016).
1122

1123    81. Rachtman, E., Balaban, M., Bafna, V. & Mirarab, S. The impact of contaminants on the
1124    accuracy of genome skimming and the effectiveness of exclusion read filters. Mol. Ecol.
1125    Resour. 20, 649–661 (2020).
1126
1127    82. Ben-Baruch, E. et al. Asymmetric Loss For Multi-Label Classification. Preprint at
1128    http://arxiv.org/abs/2009.14119 (2021).
1129
1130    83. Bushnell, B. BBMap. (2022).
1131
1132    84. Bushnell, B., Rood, J. & Singer, E. BBMerge – Accurate paired shotgun read merging via
1133    overlap. PLOS ONE 12, e0185056 (2017).
1134
1135    85. Rizk, G., Lavenier, D. & Chikhi, R. DSK: k-mer counting with very low memory usage.
1136    Bioinformatics 29, 652–653 (2013).
1137
1138    86. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor.
1139    Bioinformatics 34, i884–i890 (2018).
1140
1141    87. Tange, O. GNU Parallel 2018. (Ole Tange, 2018). doi:10.5281/zenodo.1146014.
1142
1143    88. Harris, C. R. et al. Array programming with NumPy. Nature 585, 357–362 (2020).
1144
1145    89. Howard, J. & Gugger, S. Fastai: A Layered API for Deep Learning. Information 11, 108
1146    (2020).
1147
1148    90. Wightman, R. PyTorch Image Models. GitHub repository (2019)
1149    doi:10.5281/zenodo.4414861.
1150
1151    91. Pellicer, J. & Leitch, I. J. The Plant DNA C-values database (release 7.1): an updated
1152    online repository of plant genome size data for comparative studies. New Phytol. 226, 301–
1153    305 (2020).
1154
1155    92. Fiannaca, A., La Rosa, M., Rizzo, R. & Urso, A. Analysis of DNA Barcode Sequences Using
1156    Neural Gas and Spectral Representation. in Engineering Applications of Neural Networks
1157    (eds. Iliadis, L., Pappadopoulos, H. & Jayne, C.) 212–221 (Springer, Heidelberg, 2013).
1158
1159    93. Millan Arias, P., Hill, K. A. & Kari, L. i DeLUCS: a deep learning interactive tool for
1160    alignment-free clustering of DNA sequences. Bioinformatics 39, btad508 (2023).
1161
1162    94. D'Ercole, J., Prosser, S. W. J. & Hebert, P. D. N. A SMRT approach for targeted amplicon
1163    sequencing of museum specimens (Lepidoptera)—patterns of nucleotide misincorporation.
1164    PeerJ 9, e10420 (2021).
1165
1166    95. Sproul, J. S. & Maddison, D. R. Cryptic species in the mountaintops: species delimitation
1167    and taxonomy of the Bembidion breve species group (Coleoptera: Carabidae) aided by

1168 genomic architecture of a century-old type specimen. Zool. J. Linn. Soc. 183, 556–583
1169 (2018).
1170
1171 96. Keuler, R. et al. Interpreting phylogenetic conflict: hybridization in the most speciose
1172 genus of lichen-forming fungi. Mol. Phylogenet. Evol. 174, 107543 (2022).
1173
1174 97. Barrett, C. F., Wicke, S. & Sass, C. Dense infraspecific sampling reveals rapid and
1175 independent trajectories of plastome degradation in a heterotrophic orchid complex. New
1176 Phytol. 218, 1192–1204 (2018).
1177
1178 98. Freschi, L. et al. Population structure, biogeography and transmissibility of
1179 Mycobacterium tuberculosis. Nat. Commun. 12, 6099 (2021).
1180
1181 99. Ma, B. et al. A genomic catalogue of soil microbiomes boosts mining of biodiversity and
1182 genetic resources. Nat. Commun. 14, 7318 (2023).
1183
1184 100. Asprino, R. et al. A dataset for benchmarking molecular identification tools based on
1185 genome skimming. Preprint at https://doi.org/10.32942/X2DW6K (2024).
1186
1187 101. Pomerantz, A. et al. Rapid in situ identification of biological specimens via DNA
1188 amplicon sequencing using miniaturized laboratory equipment. Nat. Protoc. 17, 1415–1443
1189 (2022).
1190
1191 102. Kimura, L. T. et al. Amazon Biobank: a collaborative genetic database for bioeconomy
1192 development. Funct. Integr. Genomics 23, 101 (2023).
1193
1194 103. Lewin, H. A. et al. The Earth BioGenome Project 2020: starting the clock. Proc. Natl.
1195 Acad. Sci. U. S. A. 119, e2115635118 (2022).
1196
1197 104. Ebenezer, T. E. et al. Africa: sequence 100,000 species to safeguard biodiversity.
1198 Nature 603, 388–392 (2022).
1199
1200 105. Cheng, S. et al. 10KP: A phylodiverse genome sequencing plan. GigaScience 7, giy013
1201 (2018).
1202
1203 106. Staff, E. A reference standard for genome biology. Nat. Biotechnol. 36, 1121–1121
1204 (2018).
1205
1206 107. i5K Consortium. The i5K Initiative: Advancing Arthropod Genomics for Knowledge,
1207 Human Health, Agriculture, and the Environment. J. Hered. 104, 595–600 (2013).
1208
1209 108. Davis, C. C. The herbarium of the future. Trends Ecol. Evol. (2022)
1210 doi:10.1016/j.tree.2022.11.015.
1211
1212 109. Davis, C. C. Collections are truly priceless. Science 383, 1035–1035 (2024).
1213

1214    110. Davis, C. C., Sessa, E. B., Paton, A., Antonelli, A. & Teisher, J. The destructive sampling
1215    conundrum and guidelines for effective and ethical sampling of herbaria. EcoEvoRxiv
1216    (2024) doi:10.32942/X2C603.
1217
1218    111. Card, D. C., Shapiro, B., Giribet, G., Moritz, C. & Edwards, S. V. Museum genomics. Annu.
1219    Rev. Genet. 55, 633–659 (2021).
1220
1221    112. Leavitt, S. D. et al. Fungal specificity and selectivity for algae play a major role in
1222    determining lichen partnerships across diverse ecogeographic regions in the lichen-
1223    forming family Parmeliaceae (Ascomycota). Mol. Ecol. 24, 3779–3797 (2015).
1224
1225    113. Bushnell, B. BBtools v.37.61. (2017).
1226
1227    114. Nissen, J. N. et al. Improved metagenome binning and assembly using deep variational
1228    autoencoders. Nat. Biotechnol. 39, 555–560 (2021).
1229
1230    115. Vurture, G. W. et al. GenomeScope: fast reference-free genome profiling from short
1231    reads. Bioinformatics 33, 2202–2204 (2017).
1232
1233    116. Maaten, L. van der & Hinton, G. Visualizing Data using t-SNE. J. Mach. Learn. Res. 9,
1234    2579–2605 (2008).
1235
1236    117. Clark, A. Pillow, Version 9.4.0. Software. https://pypi.org/project/Pillow/. (2023).
1237
1238    118. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition. arXiv
1239    1512.03385 (2015) doi:10.1109/CVPR.2016.90.
1240
1241    119. Tan, M. & Le, Q. V. EfficientNet: Rethinking Model Scaling for Convolutional Neural
1242    Networks. ArXiv abs/1905.11946, (2019).
1243
1244    120. Cai, L., Zhang, H. & Davis, C. C. PhyloHerb: A high-throughput phylogenomic pipeline
1245    for processing genome skimming data. Appl. Plant Sci. 10, e11475 (2022).
1246
1247    121. Marinho, L. C. et al. Plastomes resolve generic limits within tribe Clusieae (Clusiaceae)
1248    and reveal the new genus Arawakia. Mol. Phylogenet. Evol. 134, 142–151 (2019).
1249
1250    122. Lyra, G. de M. et al. Phylogenomics, divergence time estimation and trait evolution
1251    provide a new look into the Gracilariales (Rhodophyta). Mol. Phylogenet. Evol. 165, 107294
1252    (2021).
1253
1254    123. Marinho, L. C. et al. Phylogenetic Relationships of Tovomita (Clusiaceae): Carpel
1255    Number and Geographic Distribution Speak Louder than Venation Pattern. Syst. Bot. 46,
1256    102–108 (2021).
1257
1258    124. Bankevich, A. et al. SPAdes: a new genome assembly algorithm and Its applications to
1259    single-cell sequencing. J. Comput. Biol. 19, 455–477 (2012).

1260
1261  125. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
1262  Improvements in performance and usability. Mol. Biol. Evol. 30, 772–780 (2013).
1263
1264  126. Minh, B. Q. et al. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic
1265  Inference in the Genomic Era. Mol. Biol. Evol. 37, 1530–1534 (2020).
1266
1267  127. Jin, J.-J. et al. GetOrganelle: a fast and versatile toolkit for accurate de novo assembly of
1268  organelle genomes. Genome Biol. 21, 241 (2020).
1269

# Data Availability

1271  New data generated *de novo* genomic for this study is available on NCBI SRA under

1272  Bioproject PRJNA1052627. All datasets and metadata are thoroughly described in a

1273  companion Data Descriptor article[100] and deposited at Harvard dataverse

1274  (https://doi.org/10.7910/DVN/IMOX0S), which will be made public upon manuscript

1275  acceptance. A pretrained model on rfCGRs and varKodes for the all-SRA-taxa dataset is

1276  available at Huggingface hub

1277  (https://huggingface.co/brunoasm/vit_large_patch32_224.NCBI_SRA)

# Code Availability

1279  Code used in the initial development and test of varKoder is available on Github

1280  (https://www.github.com/brunoasm/varkoder_development). All code used to produce

1281  images is available in the development GitHub repository. The current version of varKoder

1282  is available at https://github.com/brunoasm/varKoder. Both repositories have been

1283  archived upon manuscript submission at the Figshare repository

1284  10.6084/m9.figshare.8304017, and will be made public upon acceptance.