

1 **Perspective**

2

**Sequencing Disparity in the Genomic Era**

3

Kyle T. David<sup>1\*</sup>, Alan E. Wilson<sup>2</sup>, and Kenneth M. Halanych<sup>1</sup>

4

5

*<sup>1</sup>Molette Biology Laboratory for Environmental and Climate Change Studies,*

6

*Department of Biological Sciences, Auburn University, Auburn, AL 36849*

7

*<sup>2</sup>School of Fisheries, Aquaculture, and Aquatic Sciences, Auburn University, Auburn AL*

8

*36849*

9

10

\*Corresponding Author: [kzd0038@auburn.edu](mailto:kzd0038@auburn.edu)

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32 Abstract:

33 Advances in sequencing technology have resulted in the expectation that genomic studies  
34 will become more representative of organismal diversity. To test this expectation, we  
35 explored species representation of nonhuman eukaryotes in the Sequence Read Archive.  
36 Though species richness has been increasing steadily, species evenness is decreasing over  
37 time. Moreover, the top 1% most-studied organisms increasingly represent a larger  
38 proportion of total experiments, demonstrating increasing bias in favor of a small  
39 minority of species. To better understand molecular processes and patterns, genomic  
40 studies should reverse current trends and adopt more comparative approaches.

41

## 42 **Who to Sequence?**

43 The use of model organisms, such as maize (*Zea mays*), the mouse (*Mus*  
44 *musculus*), and the fruit fly (*Drosophila melanogaster*), have contributed greatly to our  
45 understanding of biology via their tractability and large research communities (Müller  
46 and Grossniklaus 2010). Thus, when whole genome sequencing came of age, focusing on  
47 organisms that have been widely used as models was sensible. With cost-effective, high-  
48 throughput sequencing, however, many barriers that limited the use of non-model  
49 organisms have been removed. Advances in non-model and reduced representation  
50 genome sampling approaches have enabled researchers to sequence virtually any  
51 organism more cheaply and easily than ever before (Goldstein and King 2016). These  
52 advances enable comparative studies with broad sampling from across the tree of life that  
53 can elucidate the origins and variation in cellular mechanisms thereby ushering in a new  
54 era of discovery (Dunn and Ryan 2015). In light of this new approach, many researchers  
55 have predicted that the lines between model and non-model organisms will blur or  
56 disappear entirely (Davis 2004; Müller and Grossniklaus 2010; Goldstein and King 2016;  
57 Bolker 2017).

## 58 **Trends in High-Throughput Sequencing Biodiversity**

59 To explore how high-throughput sequencing efforts are distributed across the  
60 diversity of eukaryotic life, we accessed all nonhuman eukaryotic sequencing  
61 experiments in the Sequence Read Archive (SRA) using the Entrez Direct suite of UNIX  
62 commands. The SRA is a high-throughput sequence database administered by the DNA

63 Data Bank of Japan, European Nucleotide Archive, and the National Center for  
64 Biotechnology Information (NCBI). Note that experiments are defined in the SRA as “a  
65 unique sequencing result for a specific sample” and can be from experimental or  
66 descriptive research. Experiments may use one of many different sequencing strategies,  
67 though RNA-seq (37.6% of experiments) and whole genome sequencing (22.6% of  
68 experiments) are the most common. The search was executed on January 20<sup>th</sup> 2019 and  
69 returned 1,874,638 experiments. Of those, 29,578 (1.6%) experiments were removed  
70 either because they had pooled samples from multiple species or missing data.  
71 Experiments were restricted from those published between 2010, the first year with a  
72 sufficient (>100) number of species for our analysis, and December 2018, the most recent  
73 month with complete data at the time of the search. The final dataset includes 1,808,136  
74 high-throughput sequencing experiments from 24,288 unique species.

75 The top 1% of species with the most experiments represent 85.3% of all  
76 experiments accessed. The top 1% includes 120 animals, 99 plants, 14 fungi, and 9  
77 protists. Fifteen phyla are represented, although 84.1% of species are either streptophytes  
78 (green plants; n=98), chordates (n=73), or arthropods (n=30). At the species level, the  
79 mouse *Mus musculus* is the most represented by a wide margin with 523,192  
80 experiments, 4.4x more than any other species, and representing 28.9% of all  
81 experiments. All 13 model organisms officially recognized by the National Institute of  
82 Health (NIH)  
83 (<https://web.archive.org/web/20161123070020/http://modelorganisms.nih.gov/>) are  
84 featured in the top 1%, which together represent 44.7% of accessed experiments.

85 As expected, species richness has increased over time from 453 unique species  
86 sequenced in 2010 to 9,696 in 2018. However, despite an increase in the number of  
87 unique species sequenced over time, species evenness has decreased significantly (Fig.  
88 1A). This trend appears to be mediated at least in part by an increasing preference toward  
89 relatively fewer study species. The top 1% represented 49.7% of experiments in 2010;  
90 however, in 2017 the top 1% represented 80.4%. The top 1% of species for each month  
91 has been increasing at a rate of about 5.0% year<sup>-1</sup> (Fig. 1B).

92           Of the 24,288 species we accessed, 1,146 have significant increases in the number  
93 of experiments over time. Given that high-throughput sequencing has increased  
94 exponentially, seeing large increases in number of experiments over time for highly  
95 studied organisms is not surprising. Indeed, with the exception of *Plasmodium*  
96 *falciparum*, the top five species with the greatest increases were all NIH recognized  
97 models (in descending order: *M. musculus*, *Saccharomyces cerevisiae*, *P. falciparum*,  
98 *Arabidopsis thaliana*, and *D. melanogaster*).

99           We also explored change in relative frequency of experiments for each organism  
100 over time. In terms of relative frequency, the only NIH model that maintained a  
101 significant increase over time was the mouse. The two species with the largest decreases  
102 in relative frequency were *D. melanogaster* and *Caenorhabditis elegans*. While  
103 longstanding models are becoming less dominant, other organisms, such as the olive  
104 baboon (*Papio anubis*) and mummichog (*Fundulus heteroclitus*), appear to be receiving  
105 more attention (Data S1).

106           As anticipated, species richness in high-throughput sequencing experiments  
107 increased significantly over time (Fig. 1A). This finding demonstrates a trend toward  
108 sequencing greater taxonomic diversity, driven by recent initiatives such as those  
109 undertaken by Genome 10K and the Global Invertebrate Genomics Alliance (GIGA).  
110 Notably, in November 2018, the Earth BioGenome Project was launched, which aims to  
111 sequence all eukaryotic species genomes in 10 years.

112           Although more unique species are being sequenced over time, the disparity of  
113 sequencing efforts is widening, suggesting more focus is being put on relatively fewer  
114 species. In particular, the mouse has 1.8x more experiments than all other NIH models  
115 combined, and the number of experiments is growing linearly at a rate of 150.8 per  
116 month, 4.3x faster than any other organism. Mouse is also the only NIH model whose  
117 relative representation has increased, taking up a larger proportion of total sequencing  
118 experiments over time (increasing at a rate 3.5x that of any other organism).

### 119 **“All Models are Wrong”**

120           Model organisms have been a fundamental aspect of biology for at least 150 years  
121 (Müller and Grossniklaus 2010); however, they are not without problems. As statistician  
122 George Box famously articulated, “All models are wrong, but some are useful.” Many

123 researchers have previously commented on the pitfalls of model-centric research (Davis  
124 2004; Bolker 2017), which can result in disastrous and even lethal consequences as with  
125 the infamous fialuridine and thalidomide drug trials during the 20<sup>th</sup> century (Warkany  
126 1988; Xu, et al. 2014).

127 In particular, there are many questions for which mouse models are not well-  
128 suited, in spite of their popularity. In recent years, the mouse has been shown to be an  
129 imperfect representative of human disorders, particularly with regard to neurological and  
130 immune disease as well as cancer, resulting in frustratingly few real-world applications  
131 relative to the investment (Schnabel 2008; Geerts 2009; Seok, et al. 2013; Baker and  
132 Amor 2015; Perlman 2016). Additionally, early approximations of human gene count  
133 based on homology with mouse were overestimated by ~60,000, a number which was  
134 reduced to ~10,000 by a study that used the more compact genome of the pufferfish  
135 *Tetraodon nigroviridis* (Crollius, et al. 2000). Current estimates of human gene count rely  
136 on comparisons across many different species.

137 **While the model organism philosophy is instrumental to understanding certain**  
138 **intraspecific mechanisms, it is singularly inappropriate to address these questions in an**  
139 **evolutionary context and therefore not capable of answering questions on the origins and**  
140 **variation of such mechanisms.** Broad taxonomic sampling is a prerequisite to assess  
141 evolutionary processes and patterns. For example, a centralized nervous system was  
142 thought to be an ancestral character of bilaterians with a single evolutionary origin based  
143 on similar expression patterns in *D. melanogaster*, the annelid worm *Platynereis*  
144 *dumerilii*, and vertebrates. However a study in 2018, based on novel sequence data from  
145 additional groups within *Bilateria*, found different nervous-system architectures even  
146 between closely related taxa, suggesting nerve cords evolved within *Bilateria* multiple  
147 times (Martín-Durán, et al. 2018). Similar patterns remain to be explored for many other  
148 biological characters, such as gastrulation and segmentation (Dunn, et al. 2014).

### 149 **The Sangerian Shortfall**

150 Just as the Linnaean Shortfall describes how few species have been formally  
151 described, we define the Sangerian Shortfall as the lack of knowledge regarding most  
152 species' genomes. The 24,288 species represented in the SRA represent 2.0% of the

153 1,186,221 described eukaryotic species in NCBI's taxonomy database and only 0.0027%  
154 of the 8.7 million eukaryotic species thought to exist on earth. For species with whole  
155 genome sequences available, representation is even lower, 9,613 species as of December  
156 2018. Of the 14,927 species currently listed as endangered or critically endangered by the  
157 IUCN, only 2.6% had high-throughput sequence data.

## 158 **Concluding Remarks**

159 Advances in high-throughput sequencing technology have enabled researchers to  
160 sample from more species than ever before. In spite of this, genomic sampling is  
161 becoming more model-focused as relatively more attention is being paid to fewer species.  
162 Negative effects of this trend extend beyond biodiversity and evolutionary studies to  
163 many fields including pharmacology, development, genetics, and neurobiology. To  
164 improve taxonomic diversity in genomic studies we recommend 1) developing and  
165 improving funding resources to increase our understanding of biodiversity (such as the  
166 Dimensions of Biodiversity and PurSUIt programs offered by the USA National Science  
167 Foundation) and 2) taking steps to ensure reviewer panels are represented by researchers  
168 working on a variety of study systems to reduce bias in funding decisions. We predict  
169 that these recommendations, if acted upon, will not only improve our understanding of  
170 variation and diversity in nature but also foster collaborations across different research  
171 fields and systems. We conclude that molecular researchers should attempt to select  
172 models based on their relevance to biological questions over ease of use and use  
173 comparative approaches to address questions in an evolutionary framework whenever  
174 possible.

175 All code required to update experiments and reproduce results/figures are  
176 available at <https://github.com/KyleTDavid/SRA2019>. Original data files are available at  
177 <https://figshare.com/projects/SRA2019/39296>.

178

179 **Acknowledgments:** We would like to thank Doug Levey, three anonymous reviewers,  
180 and all members of the Molette and Phyletica labs, in particular Damien Waits, Caitlin  
181 Redak, Michael Tassia, and Drs. Yuanning Li and Jamie Oaks for their feedback and  
182 comments. This project was supported by the Alabama Agricultural Experiment Station,

183 the Hatch Program of the National Institute of Food and Agriculture, U.S. Department of  
184 Agriculture, Cell and Molecular Biosciences Fellowship, and the Schneller Endowed  
185 Chair Fund. This is Molette Biology Laboratory contribution XX and Auburn University  
186 Marine Biology Program contribution XX.

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215 **References**

- 216 Baker D, Amor S. 2015. Mouse models of multiple sclerosis: lost in translation?  
217 Current pharmaceutical design 21:2440-2452.
- 218 Bolker JA. 2017. Animal Models in Translational Research: Rosetta Stone or  
219 Stumbling Block? BioEssays 39:1700089.
- 220 Crollius HR, Jaillon O, Bernot A, Dasilva C, Bouneau L, Fischer C, Fizames C, Wincker  
221 P, Brottier P, Quétier F. 2000. Estimate of human gene number provided by  
222 genome-wide analysis using Tetraodon nigroviridis DNA sequence. Nature  
223 genetics 25:235.
- 224 Davis RH. 2004. The age of model organisms. Nature Reviews Genetics 5:69.
- 225 Dunn CW, Giribet G, Edgecombe GD, Hejnol A. 2014. Animal phylogeny and its  
226 evolutionary implications. Annual review of ecology, evolution, and systematics  
227 45:371-395.
- 228 Dunn CW, Ryan JF. 2015. The evolution of animal genomes. Current opinion in  
229 genetics & development 35:25-32.
- 230 Geerts H. 2009. Of mice and men. CNS drugs 23:915-926.
- 231 Goldstein B, King N. 2016. The future of cell biology: emerging model organisms.  
232 Trends in cell biology 26:818-824.
- 233 Martín-Durán JM, Pang K, Børve A, Lê HS, Furu A, Cannon JT, Jondelius U, Hejnol A.  
234 2018. Convergent evolution of bilaterian nerve cords. Nature 553:45.
- 235 Müller B, Grossniklaus U. 2010. Model organisms—a historical perspective. Journal  
236 of proteomics 73:2054-2063.
- 237 Perlman RL. 2016. Mouse models of human diseaseAn evolutionary perspective.  
238 Evolution, medicine, and public health 2016:170-176.
- 239 Schnabel J. 2008. Neuroscience: standard model. Nature News 454:682-685.
- 240 Seok J, Warren HS, Cuenca AG, Mindrinos MN, Baker HV, Xu W, Richards DR,  
241 McDonald-Smith GP, Gao H, Hennessy L. 2013. Genomic responses in mouse  
242 models poorly mimic human inflammatory diseases. Proceedings of the National  
243 Academy of Sciences 110:3507-3512.
- 244 Warkany J. 1988. Why I doubted that thalidomide was the cause of the epidemic of  
245 limb defects of 1959 to 1961. Teratology 38:217-219.
- 246 Xu D, Nishimura T, Nishimura S, Zhang H, Zheng M, Guo Y-Y, Masek M, Michie SA,  
247 Glenn J, Peltz G. 2014. Fialuridine induces acute liver failure in chimeric TK-NOG  
248 mice: a model for detecting hepatic drug toxicity prior to human testing. PLoS  
249 medicine 11:e1001628.

250

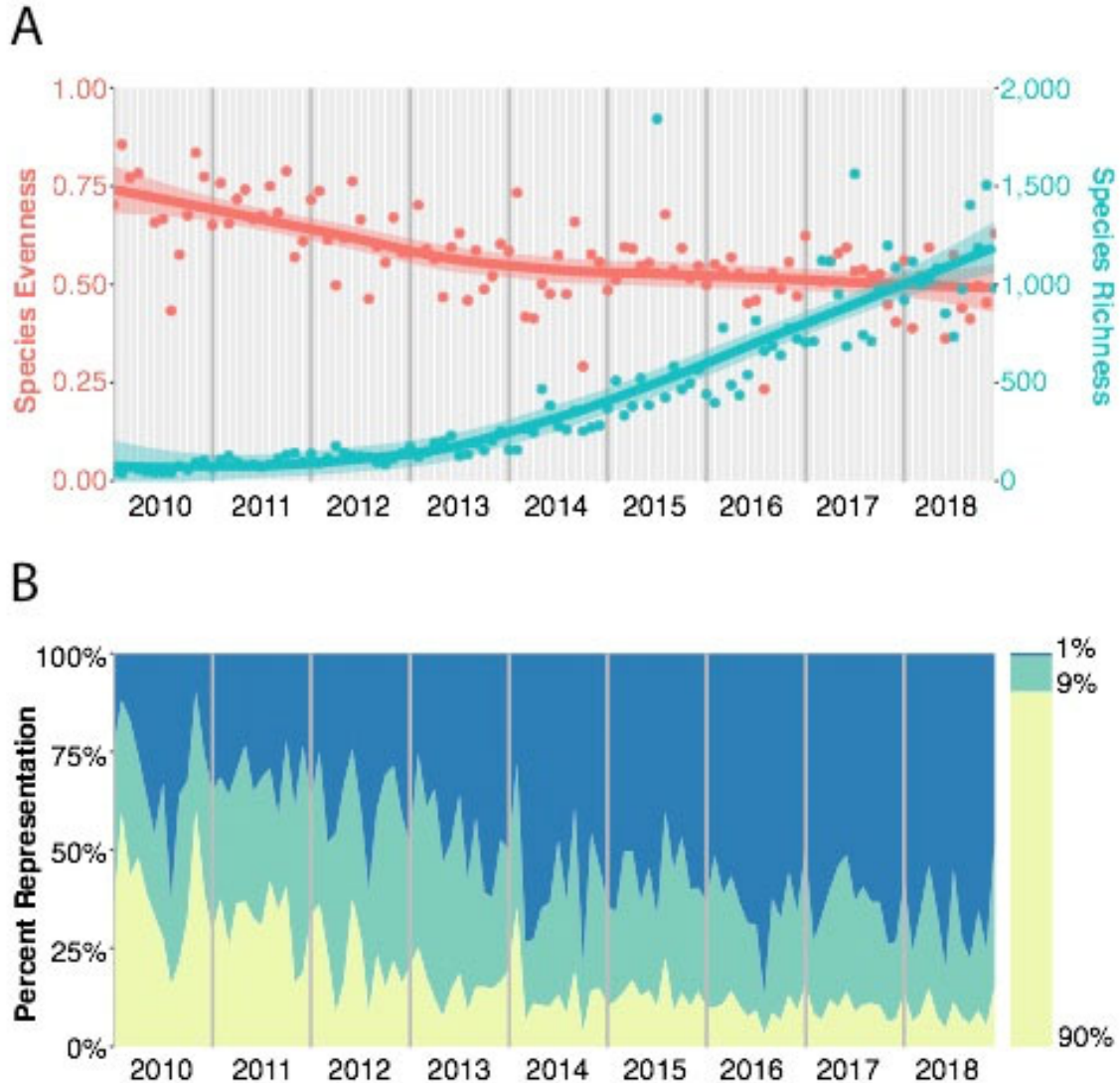
251

252

253

254





255

256 **Figure 1: Biodiversity in the Sequence Read Archive.** A) Pielou's species evenness  
 257 ( $H/\ln(S)$  where  $H$  is Shannon's diversity index and  $S$  is species richness) and species  
 258 richness of all nonhuman eukaryotic sequencing experiments in the SRA calculated for  
 259 each month between January 2010 and December 2018. Richness has increased over time  
 260 ( $p < 1.1E-32$ ) while evenness has decreased ( $p < 2.5E-13$ ). B) Relative representation of  
 261 the top 1%, top 2-10%, and bottom 90% of species in the SRA by number of  
 262 experiments for each month between January 2010 and July 2018.

263

264 **Data S1:** Summary statistics for each species represented in the SRA between January  
265 2010 and December 2018