

1 **Title:** iNaturalist is an open science resource for ecological genomics by enabling rapid and tractable
2 records of initial observations of sequenced specimens

3 **Running Title:** iNaturalist for genomics

4 **Author:** Jay Keche Goldberg, Department of Ecology and Evolutionary Biology, University of Arizona,
5 Tucson, AZ, USA

6 **For correspondence:** jaykgold@arizona.edu

8 **Abstract**

9 The rapidly growing body of publicly available sequencing data for rare species and/or wild-
10 caught samples is accelerating the need for detailed records of the specimens used to generate
11 datasets. Many already published datasets are unlikely to ever be reused, not due to problems with the
12 data themselves, but due to their questionable or unverifiable origins. In this paper, I present iNaturalist
13 – a pre-existing citizen science platform that allows people to post photo observations of organisms in
14 nature – as a tool that allows genomics researchers to rapidly publish observations of specimens used to
15 generate sequencing datasets. This practice aligns with the values of the open science movement; and I
16 also discuss how iNaturalist, along with other online resources, can be used to create an open genomics
17 pipeline that enables future replication studies and ensures the value of genomics datasets to future
18 research.

20 **Introduction**

21 The number of high-quality published genomes has increased rapidly in recent years (Kress et al.
22 2022) and the feasibility of sequencing multiple individuals of species with large heterozygous genomes
23 has enabled pan-genomics with eukaryotic organisms (Golicz et al. 2020). Once restricted to prokaryotes
24 with small genomes (Rasko et al. 2008), there are now several plant and animal species with publicly
25 available pangenome databases (Gao et al. 2019; Tong et al. 2022). Evolutionary biologists are routinely
26 using whole genome sequencing to observe responses to climate change (Waldvogel et al. 2020) and
27 experimental manipulation (Kovács and Dragoš 2019) in real time. Many labs and consortia are
28 publishing genomes as fast as possible to make them available to the broader scientific community
29 (Mathers et al. 2022), but often publish their data in minimalist reports (Smith et al. 2017) that
30 sometimes lack even basic descriptions of the data itself (Hains et al. 2020). The explosion of genomic
31 data, while scientifically exciting, presents a dilemma if details regarding the collection of source
32 specimen(s) are not properly recorded and made available to the broader scientific community.
33 Datasets originating from wild specimens require more rigorous documentation of the originating
34 samples to ensure their long-term value – especially when they represent rare, cryptic, or species that
35 are likely to see their taxonomic identity altered (such as members of poorly resolved clades). This is a
36 significant gap for open genomics research, but iNaturalist – a platform where users post observations
37 of wildlife and experts identify them – could be a valuable tool for researchers who wish to improve the
38 reusability of their data and help create a fully open genomics data pipeline (Figure 1). The publicly
39 accessible nature of iNaturalist observations (one can access them without an account on the platform)
40 makes it ideal for tackling the lack of robust, easily accessible, information regarding the originating
41 samples used to generate publicly available sequencing datasets.

43 **What is iNaturalist?**

44 iNaturalist is a citizen science platform that allows users to upload photos from an internet
45 connected device (smartphone, computer, etc.). Knowledgeable identifiers – often actively publishing
46 researchers or museum curators – then identify these organisms. These photo observations are also
47 accompanied by metadata – the date/time and location at which the photo was taken – and sometimes
48 include specific notes regarding the sex/life stage/etc. of the observed organism (these are often filled in

49 by identifiers). Any discussion of the observations by the observer and identifiers is also recorded and
50 associated with it. iNaturalist has already proven its value to ecologists and is cited in several papers
51 regarding invasion dynamics (Serniak et al. 2022) and animal behavior (Vardi et al. 2021).

52

53 **An open genomics pipeline**

54 Open access journals lacking paywalls have become commonplace and many funding agencies
55 mandate that results be published in them. Public repositories for various forms of data (GenBank,
56 Dryad, etc.) – and the code needed to analyze them (Github) – exist and are often free to contribute to.
57 Some model species and popular study clades even have their own dedicated repositories (e.g. Flybase,
58 Sol Genomics Network). Resources for publishing step-by-step methodologies (protocols.io) also exist.
59 Yet, until the advent of iNaturalist there was no way to freely publish open access natural history
60 observations other than within peer-reviewed publications. Now, however, it is possible to instantly
61 upload photos from the field, have them automatically associated with key metadata (time and
62 location), and make them freely available to both the scientific community and broader public using
63 iNaturalist. This makes it a valuable tool for ecological and evolutionary geneticists to improve their data
64 pipelines and better align with open science practices.

65 iNaturalist's utility lies in how it allows researchers to associate publications with field
66 observations via their unique URLs (example user profile and observation can be found in Web
67 Resources) that provide an easy-to-follow paper trail. This allows future researchers to verify the
68 identity of the initial sample/specimen and collection details. This is critical for species that are likely to
69 have their taxonomy revised as their identity can be followed through disagreements between
70 systematists based on their observable traits. It also allows for minor identification errors (e.g. those
71 that do not meaningfully alter the outcome of a study) to be easily resolved through the community
72 identification process and for these resolutions to be linked to the final dataset and/or publication
73 without the need for formal corrections.

74 When accessed in bulk through the Global Biodiversity Information Facility (GBIF), sets of
75 iNaturalist observations can be given digital object identifiers (DOIs) that enable replication studies
76 (Forti et al. 2022a/b); and, within the iNaturalist platform, observations can be collected into projects.
77 This makes it easy for researchers to keep track of the samples used in a study that they may be
78 planning, carrying out, or have already published. Any projects that an observation is a part of are
79 shown underneath the observation, thus making it easy to track how researchers have used a
80 sample/dataset over time. In addition to tracking important metadata regarding the use of scientific
81 samples for open and repeatable science, this gives the public deeper insight into the science of the
82 species they see in daily life and a direct line to the researchers conducting it.

83

84 **Conclusions**

85 As the genomics revolution continues to open doors to research on the ecology and evolution of
86 previously impossible-to-study species, the need for better documentation of data origins will increase
87 dramatically. iNaturalist is a platform that researchers can use to rapidly publish field observations of
88 samples/specimens that are eventually used in sequencing projects. When combined with other open
89 science resources, it creates an open genomics data pipeline that allows both the scientific community
90 and public-at-large to have better insight into the process behind genomics research.

91

92 **Acknowledgements**

93 I would like to acknowledge the California Academy of Sciences and National Geographic Society
94 for enabling the iNaturalist initiative and all developers who have worked on the project in any way. I
95 would also like to thank Margaret Wilch for introducing me to iNaturalist and Judith Bronstein for
96 encouraging me to write this manuscript, as well as helpful comments on an early draft for this

97 manuscript. My funding is currently provided by a National Science Foundation postdoctoral research
98 fellowship in biology (PRFB #2010772) and the University of Arizona.

99

100 **Data Availability Statement**

101 No new data or code was generated during the preparation of this manuscript. Links to all
102 mentioned platforms can be found in the 'Web Resources' section.

103

104 **Conflict of Interest Statement**

105 I declare that I have no conflict of interest associated with the contents of this manuscript; and
106 that I am not affiliated with iNaturalist (or its parent organizations) in any way beyond that of other
107 enthusiastic users.

108

109 **Web Resources**

110 *iNaturalist*

111 Homepage: <https://www.inaturalist.org/>

112 GBIF Homepage: <https://www.gbif.org/>

113 iNaturalist User Profile: <https://www.inaturalist.org/people/6089000>

114 Example Observation: <https://www.inaturalist.org/observations/134334492>

115

116 *Public Information Repositories*

117 Protocols.io: <https://www.protocols.io/>

118 Dryad: <https://datadryad.org/stash>

119 Github: <https://github.com/>

120 GenBank: <https://www.ncbi.nlm.nih.gov/genbank/>

121 European Nucleotide Archive (ENA): <https://www.ebi.ac.uk/ena/browser/home>

122 FlyBase: <https://flybase.org/>

123 WormBase: <https://wormbase.org/>

124 The Arabidopsis Information Resource (TAIR): <https://www.arabidopsis.org/>

125 Sol Genomics Network: <https://solgenomics.net/>

126 Saccharomyces Genome Database: <https://www.yeastgenome.org/>

127

128 *Biology focused pre-print servers*

129 BioRxiv: <https://www.biorxiv.org/>

130 EcoEvoRxiv: <https://ecoevorxiv.org/>

131 MedRxiv: <https://www.medrxiv.org/>

132 Zenodo: <https://zenodo.org/>

133

134

135 **Literature Cited**

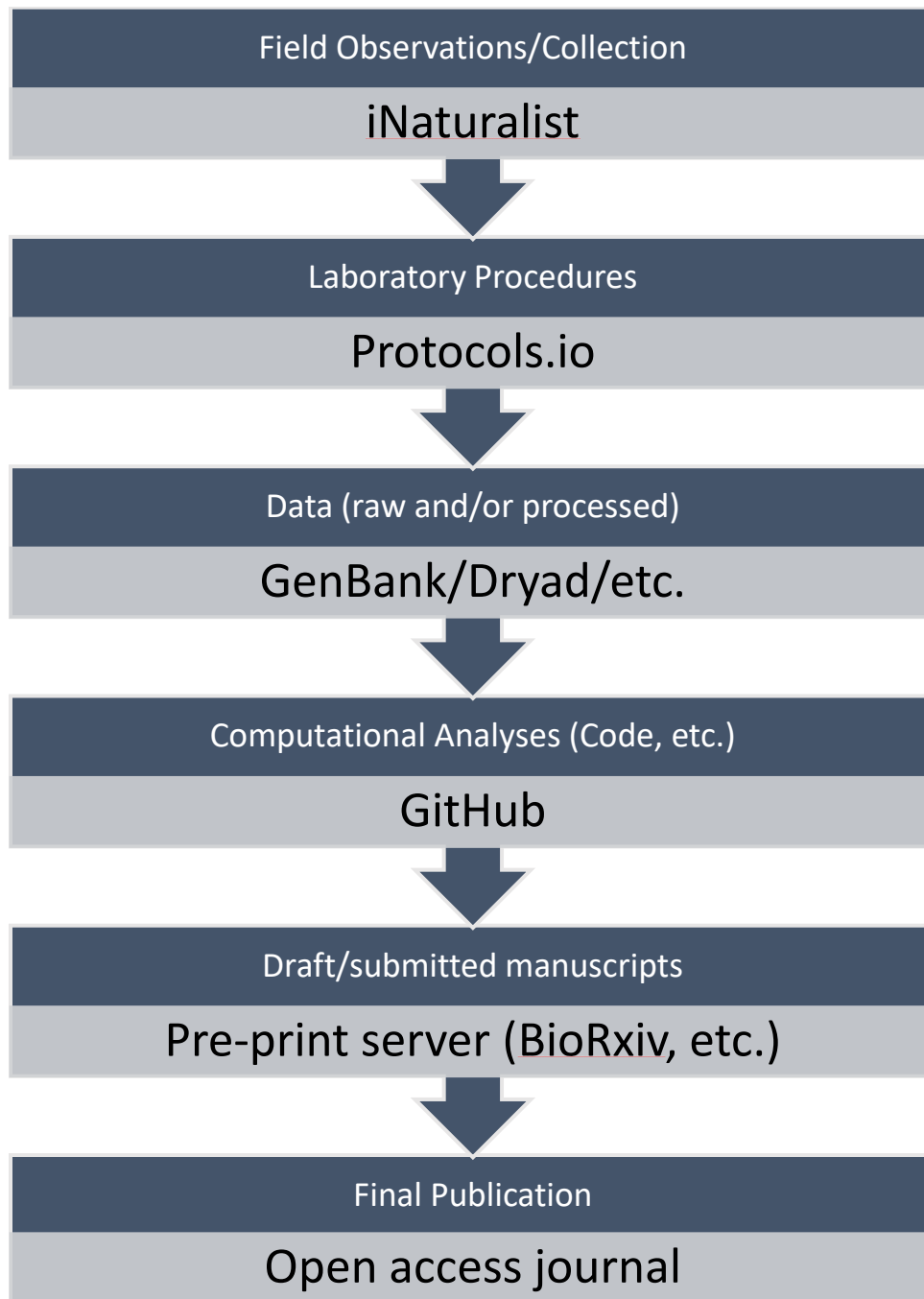
- 136 Cheng, H., Concepcion, G.T., Feng, X., Zhang, H., Li, H., 2021. Haplotype-resolved de novo assembly using
137 phased assembly graphs with hifiasm. *Nat Methods* 18, 170–175.
138 <https://doi.org/10.1038/s41592-020-01056-5>
- 139 Forti, L.R., Hepp, F., de Souza, J.M., Protazio, A., Szabo, J.K., 2022a. Climate drives anuran breeding
140 phenology in a continental perspective as revealed by citizen-collected data. *Diversity and*
141 *Distributions* 28, 2094–2109. <https://doi.org/10.1111/ddi.13610>
- 142 Forti, L.R., Hepp, F., Souza, J.M. de, Protazio, A., Szabo, J.K., 2022b. Climate drives anuran breeding
143 phenology in a continental perspective as revealed by citizen-collected data. Zenodo
144 <https://doi.org/10.5281/zenodo.6811407>
- 145 Gao, L., Gonda, I., Sun, H., Ma, Q., Bao, K., Tieman, D.M., Burzynski-Chang, E.A., Fish, T.L., Stromberg,
146 K.A., Sacks, G.L., Thannhauser, T.W., Foolad, M.R., Diez, M.J., Blanca, J., Canizares, J., Xu, Y., van
147 der Knaap, E., Huang, S., Klee, H.J., Giovannoni, J.J., Fei, Z., 2019. The tomato pan-genome
148 uncovers new genes and a rare allele regulating fruit flavor. *Nat Genet* 51, 1044–1051.
149 <https://doi.org/10.1038/s41588-019-0410-2>
- 150 Golicz, A.A., Bayer, P.E., Bhalla, P.L., Batley, J., Edwards, D., 2020. Pangenomics Comes of Age: From
151 Bacteria to Plant and Animal Applications. *Trends in Genetics* 36, 132–145.
152 <https://doi.org/10.1016/j.tig.2019.11.006>
- 153 Hains, T., O’Neill, K., Velez, J., Speed, N., Clubb, S., Oleksyk, T., Pirro, S., 2020. The complete genome
154 sequences of 22 parrot species (Psittaciformes, Aves). *F1000 Research*
155 <https://doi.org/10.12688/f1000research.25560.1>
- 156 Hon, T., Mars, K., Young, G., Tsai, Y.-C., Karalius, J.W., Landolin, J.M., Maurer, N., Kudrna, D., Hardigan,
157 M.A., Steiner, C.C., Knapp, S.J., Ware, D., Shapiro, B., Peluso, P., Rank, D.R., 2020. Highly
158 accurate long-read HiFi sequencing data for five complex genomes. *Sci Data* 7, 399.
159 <https://doi.org/10.1038/s41597-020-00743-4>
- 160 Kim, M., Jung, J.-K., Shim, E.-J., Chung, S.-M., Park, Y., Lee, G.P., Sim, S.-C., 2021. Genome-wide SNP
161 discovery and core marker sets for DNA barcoding and variety identification in commercial
162 tomato cultivars. *Scientia Horticulturae* 276, 109734.
163 <https://doi.org/10.1016/j.scienta.2020.109734>
- 164 Kovács, Á.T., Dragoš, A., 2019. Evolved Biofilm: Review on the Experimental Evolution Studies of *Bacillus*
165 *subtilis* Pellicles. *Journal of Molecular Biology, Underlying Mechanisms of Bacterial Phenotypic*
166 *Heterogeneity and Sociobiology* 431, 4749–4759. <https://doi.org/10.1016/j.jmb.2019.02.005>
- 167 Kress, W.J., Soltis, D.E., Kersey, P.J., Wegrzyn, J.L., Leebens-Mack, J.H., Gostel, M.R., Liu, X., Soltis, P.S.,
168 2022. Green plant genomes: What we know in an era of rapidly expanding opportunities.
169 *Proceedings of the National Academy of Sciences* 119, e2115640118.
170 <https://doi.org/10.1073/pnas.2115640118>
- 171 Mathers, T.C., Mugford, S.T., Wouters, R.H.M., Heavens, D., Botha, A.-M., Swarbreck, D., Van
172 Oosterhout, C., Hogenhout, S.A., 2022. Aphidinae comparative genomics resource. Zenodo
173 <https://doi.org/10.5281/zenodo.5908005>
- 174 Rasko, D.A., Rosovitz, M.J., Myers, G.S.A., Mongodin, E.F., Fricke, W.F., Gajer, P., Crabtree, J., Sebaihia,
175 M., Thomson, N.R., Chaudhuri, R., Henderson, I.R., Sperandio, V., Ravel, J., 2008. The
176 Pangenome Structure of *Escherichia coli*: Comparative Genomic Analysis of *E. coli* Commensal
177 and Pathogenic Isolates. *Journal of Bacteriology* 190, 6881–6893.
178 <https://doi.org/10.1128/JB.00619-08>
- 179 Serniak, L.T., Chan, S.S., Lajtha, K., 2023. Predicting habitat suitability for *Amyntas* spp. in the United
180 States: a retrospective analysis using citizen science data from iNaturalist. *Biol Invasions* 25,
181 817–825. <https://doi.org/10.1007/s10530-022-02947-8>
- 182 Smith, D.R., 2017. Goodbye genome paper, hello genome report: the increasing popularity of ‘genome

183 announcements' and their impact on science. *Brief Funct Genomics* 16, 156–162.
184 <https://doi.org/10.1093/bfgp/elw026>

185 Tong, X., Han, M.-J., Lu, K., Tai, S., Liang, S., Liu, Yucheng, Hu, H., Shen, J., Long, A., Zhan, C., Ding, X., Liu,
186 S., Gao, Q., Zhang, B., Zhou, Linli, Tan, D., Yuan, Y., Guo, N., Li, Y.-H., Wu, Z., Liu, L., Li, C., Lu, Y.,
187 Gai, T., Zhang, Y., Yang, R., Qian, H., Liu, Yanqun, Luo, J., Zheng, L., Lou, J., Peng, Y., Zuo, W.,
188 Song, J., He, S., Wu, S., Zou, Y., Zhou, Lei, Cheng, L., Tang, Y., Cheng, G., Yuan, L., He, W., Xu, J.,
189 Fu, T., Xiao, Y., Lei, T., Xu, A., Yin, Y., Wang, J., Monteiro, A., Westhof, E., Lu, C., Tian, Z., Wang,
190 W., Xiang, Z., Dai, F., 2022. High-resolution silkworm pan-genome provides genetic insights into
191 artificial selection and ecological adaptation. *Nat Commun* 13, 5619.
192 <https://doi.org/10.1038/s41467-022-33366-x>

193 Vardi, R., Berger-Tal, O., Roll, U., 2021. iNaturalist insights illuminate COVID-19 effects on large
194 mammals in urban centers. *Biological Conservation* 254, 108953.
195 <https://doi.org/10.1016/j.biocon.2021.108953>

196 Waldvogel, A.-M., Feldmeyer, B., Rolshausen, G., Exposito-Alonso, M., Rellstab, C., Kofler, R., Mock, T.,
197 Schmid, K., Schmitt, I., Bataillon, T., Savolainen, O., Bergland, A., Flatt, T., Guillaume, F.,
198 Pfenninger, M., 2020. Evolutionary genomics can improve prediction of species' responses to
199 climate change. *Evolution Letters* 4, 4–18. <https://doi.org/10.1002/evl3.154>
200



201
202
203
204
205
206

Figure 1. A flowchart outlining an example “open genomics pipeline” with five key steps and their corresponding open science platform. The first step in this pipeline, publicly recording the initial field observations/collection associated with a study, is the aspect that iNaturalist fulfills. The precise steps, and platforms used to carry them out, necessary for the best open science practices will vary, given the wealth of system-specific databases such as FlyBase or the Sol Genomics Network.