

1 **Realizing the promise of biodiversity genomics with highly accurate long reads**

2

3 Scott Hotaling^{1*}, Edward R. Wilcox², Jacqueline Heckenhauer^{3,4}, Russell J. Stewart⁵, and Paul
4 B. Frandsen^{3,6,7*}

5

6 **Affiliations:**

7 ¹ Department of Watershed Sciences, Utah State University, Logan, UT, USA

8 ² DNA Sequencing Center, Department of Biology, Brigham Young University, Provo, UT, USA

9 ³ LOEWE Centre for Translational Biodiversity Genomics (LOEWE-TBG), Frankfurt, Germany

10 ⁴ Department of Terrestrial Zoology, Senckenberg Research Institute and Natural History
11 Museum Frankfurt, Frankfurt 60325, Germany

12 ⁵ Department of Biomedical Engineering, University of Utah, Salt Lake City, UT, USA

13 ⁶ Department of Plant and Wildlife Sciences, Brigham Young University, Provo, UT, USA

14 ⁷ Data Science Lab, Smithsonian Institution, Washington, DC, USA

15

16 ***Authors for Correspondence:**

17 Scott Hotaling, Department of Watershed Sciences, Utah State University, Logan, UT, USA;

18 Email: scott.hotaling1@gmail.com; Phone: (828) 507-9950

19

20 Paul B. Frandsen, Department of Plant and Wildlife Sciences, Brigham Young University, Provo,
21 UT, USA; Email: paul_frandsen@byu.edu; Phone: (804) 422-2283

22

23 **Abstract:**

24 Generating the most contiguous, accurate genome assemblies given available sequencing
25 technologies is a long-standing challenge in genome science. With the rise of long-read
26 sequencing, assembly challenges have shifted from merely increasing contiguity to correctly
27 assembling complex, repetitive regions of interest, ideally in a phased manner. At present,
28 researchers largely choose between two types of long read data: longer, but less accurate

29 sequences, often generated with Oxford Nanopore (ONT) technology, or highly accurate, but
30 shorter, reads typically generated with Pacific Biosciences HiFi. To understand how both
31 technologies influence genome assembly and to clarify how scale of data (i.e., mean length and
32 sequencing depth) influence outcomes, we compared genome assemblies for a caddisfly,
33 *Hesperophylax magnus*, generated with ONT and HiFi data. Despite shorter reads and less
34 coverage, HiFi reads outperformed ONT reads in all assembly metrics tested and allowed for
35 accurate assembly of the repetitive ~20-Kb *H-fibroin* gene. Next, we quantified the influence of
36 data type on genome assemblies across 6,750 plant and animal genomes. We show that HiFi
37 reads consistently outperform all other data types for both plants and animals and may
38 represent a particularly valuable tool for assembling complex plant genomes. To realize the
39 promise of biodiversity genomics, we call for greater uptake of highly accurate long-reads in
40 future studies.

41

42 **Keywords:** Insecta, Oxford Nanopore, PacBio, HiFi, caddisfly, genome biology

43

44 **Significance statement:**

45 Understanding how types of sequence data influence genome assembly is an important aspect
46 of genome science. In general, more data—i.e., longer reads, greater depth of coverage—often
47 yields better genome assemblies. However, it is unclear how highly accurate long-read
48 sequence data (e.g., PacBio HiFi) compare to noisier long-read data. We showed that HiFi
49 outperformed noisier long-read data for a caddisfly species in terms of assembly contiguity and
50 resolution of the highly repetitive ~20-Kb *H-fibroin* gene. We also showed that this
51 outperformance likely extends to all animals and plants via a field-wide meta-analysis. Thus,
52 long-read accuracy should be emphasized in future genome studies.

53

54 **Body:**

55 As genome sequencing has been revolutionized by high-throughput sequencing, a general rule
56 has emerged: more data—e.g., longer reads or greater depth of coverage—yields more

57 contiguous, accurate genome assemblies. This is particularly evident when read length is
58 considered; third-generation long reads, which are often tens or even hundreds of thousands of
59 base pairs in length, have dramatically improved genome assemblies across the Tree of Life
60 (Hotaling, et al. 2021a; Hotaling, et al. 2021b; Marks, et al. 2021; Rhie, et al. 2021). For
61 coverage, an increase can limit the impacts of erroneous read calls through more replication of
62 potential variants (Sims, et al. 2014). However, a shortcoming of second-generation, short-read
63 platforms (e.g., Illumina) for genome assembly is that no amount of data will allow for resolution
64 of repeat-driven gaps that exceed read lengths (Sims, et al. 2014). The power of long reads to
65 mitigate this issue is well-documented (e.g., Hotaling, et al. 2022) but not all long reads are
66 created equal; indeed, different platforms yield different length versus error profiles
67 (Amarasinghe, et al. 2020; De Coster, et al. 2021). This difference is particularly important since
68 some long, repeat-rich genomic regions—including many genes of phenotypic relevance—pose
69 assembly challenges even when long read data are used. Generally speaking, past evidence
70 would suggest that to resolve difficult genomic regions with long read data, longer reads at
71 greater depth of coverage will outperform shorter reads and/or less dense coverage. However, it
72 is unclear how this expectation jibes with read accuracy, particularly for two common types of
73 long-read data that are currently in use: longer but noisier long reads [e.g., Oxford Nanopore
74 (ONT)] and more accurate, but shorter, long reads (e.g., PacBio HiFi).

75
76 We inadvertently tested this prediction—that longer, higher coverage, and noisier reads would
77 outperform shorter, but more accurate, reads at lower sequencing depth—in a silk-producing
78 caddisfly, *Hesperophylax magnus*, and the highly repetitive *heavy fibroin chain* gene (*H-fibroin*).
79 Caddisflies share an evolutionary origin of silk with butterflies and moths, including the primary
80 protein component of silk, *H-fibroin* (Frandsen, et al. 2019). *H-fibroin* commonly spans ~20 Kb
81 and its amino acid sequence consists of conserved termini and a highly repetitive internal
82 region. Early efforts to assemble the complete *H-fibroin* gene with short reads were

83 unsuccessful due to its long repetitive region (Ashton, et al. 2013; Yonemura, et al. 2009).

84 However, long-read assemblies have since yielded full-length sequences (Frandsen, et al.

85 2019; Kawahara, et al. 2022; Luo, et al. 2018).

86

87 In this study, we compared two long-read genome assemblies for *H. magnus* produced with two

88 technologies: ONT and HiFi. For comparison, we considered genome-wide metrics as well as

89 accurate assembly of the *H-fibroin* gene. We chose to focus on *H-fibroin* as a surrogate for

90 complex but phenotypically important genes where we expected the benefits of highly accurate

91 long reads to be most obvious. To estimate whether assembly outperformance with HiFi data

92 was unique to our focal caddisfly or reflects a broader trend in genome biology, we performed a

93 meta-analysis of contig N50, assembly length, and sequencing technology for all publicly

94 available plant and animal genomes on GenBank. For our caddisfly case study and meta-

95 analysis, highly accurate HiFi sequence data dramatically outperform all other types of

96 sequence data.

97

98 For our caddisfly genome comparison, we sequenced two individuals of *H. magnus* from the

99 same population. For the first individual, we generated a combination of noisy ONT sequencing

100 (R.9.4.1, LSK-109 ligation library prep kit) and Illumina sequencing (NovaSeq). For the second

101 individual, we only generated HiFi reads via CCS sequencing on the PacBio Sequel II platform.

102 We used Guppy v.5.0.11 to base-call the ONT reads and removed all reads under 5 Kb for

103 further analysis. We used SMRTlink v.10 to generate HiFi reads (reads with quality >Q20). To

104 assemble genomes for the two individuals, we tested a range of assemblers. For ONT, we

105 tested Canu v.1.8 (Koren, et al. 2017), wtdbg2 v.2.4 (Ruan and Li 2020), and a hybrid approach

106 with MaSuRCA (Zimin, et al. 2013). For HiFi, we tested Hifiasm (Cheng, et al. 2021) and

107 HiCanu (Nurk, et al. 2020). We selected the best assembly based on contiguity and

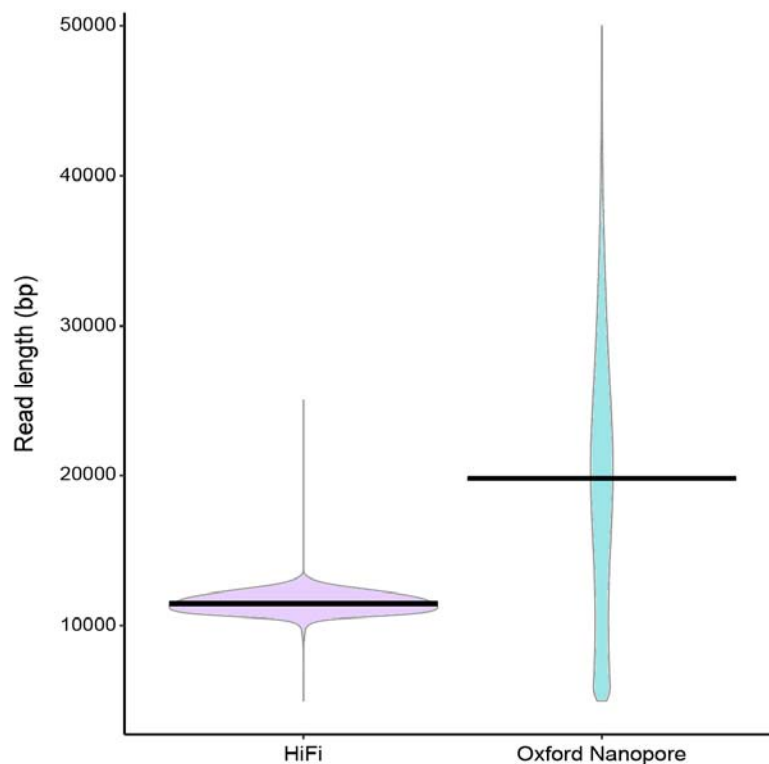
108 “Benchmarking Universal Single-Copy Orthologs” (BUSCO) scores. We ran BUSCO v.5.2.2

109 (Manni, et al. 2021) using the 1,367 reference genes in the OrthoDB v.10 Insecta gene set
110 (Kriventseva, et al. 2019). To evaluate recovery and assembly of *H-fibroin*, we used tblastn to
111 identify conserved terminal sequences from existing transcriptomes (Ashton, et al. 2013). We
112 then extracted *H-fibroin* from both assemblies with 1,000 additional bps from each terminus and
113 annotated it using Augustus v.3.3.2 (Stanke, et al. 2006). To visualize mismatches between
114 reads and the assembly, we mapped raw reads to the assembled *H-fibroin* gene using
115 Minimap2 (Li 2018) and visualized the results in Geneious 2022.0.2 (Kearse, et al. 2012).

116
117 To assess the influence of sequencing technology on genome assembly across the Tree of Life,
118 we extracted all available genome assemblies for plants (class Embryophyta) and animals
119 (class Metazoa) from GenBank using the “summary genome” function in v.10.9.0 of the NCBI
120 Datasets command-line tool on 13 November 2021. We then used the “lineage” function of
121 TaxonKit (Shen and Xiong 2019) to retrieve taxonomic information for all entries in our genome
122 assembly list. Next, we gathered additional metadata (e.g., sequencing technology) for each
123 entry using a custom web scraper script (modified from Hotaling, et al. 2021b). Next, we
124 removed duplicate assemblies and alternative haplotypes for a given assembly that were
125 identified through keyword searching in either the BioProject information or assembly title.

126
127 We binned assemblies into four sequencing technology categories: short-reads (e.g., Illumina),
128 long-read ONT (ONT long-reads with or without short-reads), long-read PacBio (non-HiFi
129 PacBio long-reads with or without short reads), or HiFi (any assembly where HiFi long-reads
130 were used). To assist in categorization, we considered any assembly that was generated before
131 2017—when long-read assemblies began to emerge (Hotaling, et al. 2021b)—to be a short-read
132 assembly. We also used self-reported information on the genome assembly algorithms used to
133 classify assemblies. For instance, if an assembly only reported PacBio sequence data but a
134 HiFi-specific assembler (e.g., Hifiasm) was used, we classified it as a HiFi assembly. After

135 binning, we removed any assembly for which the sequence data type used could not be
136 established. We tested for differences in the distributions of assembly size and contig N50
137 among plant and animals or our sequence type categories within the overarching plant or
138 animal grouping using Welch Two Sample T-tests or one-way ANOVAs followed by Tukey HSD
139 tests in R v3.6.3 (R Core Team 2021). While all statistical tests were performed on the
140 untransformed data, we visualized log-transformed comparisons using ggplot2 (Wickham 2011).
141



142
143 **Figure 1.** Violin plots of read lengths for the HiFi and Oxford Nanopore data sets used to assemble
144 *Hesperophylax magnus* genomes in this study. Width of the colored areas indicate numbers of reads according
145 to lengths on the y-axis. Dark lines represent the medians of each distribution.
146

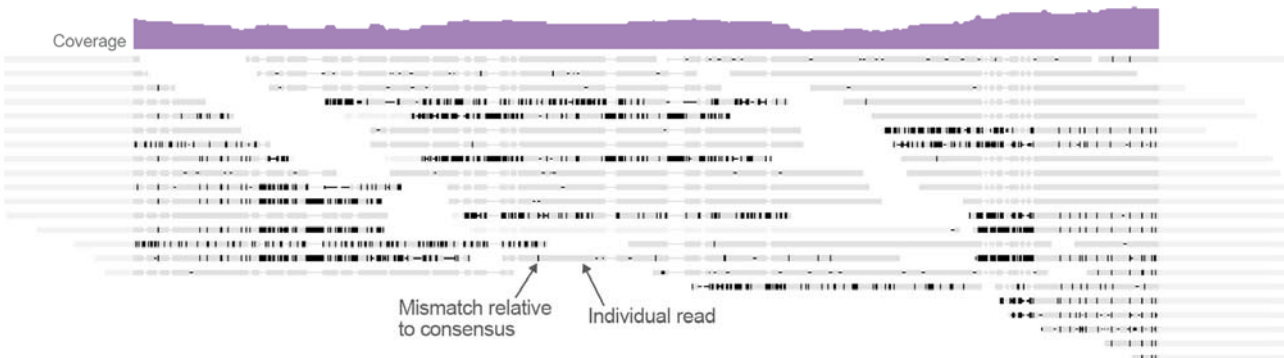
147 The ONT library had a wider distribution of read lengths with a median of 19.9 Kb for 33.5 Gb of
148 raw data. The HiFi dataset had a median read length of 11.3 kb bp for 28 Gb of raw data (Fig.
149 1). The best ONT assembly (Genbank #GCA_016648045.1) was generated with MaSuRCA and
150 spanned 1.23 Gb. For HiFi, the best assembly was produced with Hifiasm (Genbank

151 #JAIUSX000000000) and was nearly identical in length at 1.22 Gb. However, we observed a
152 dramatic difference in contiguity; the ONT assembly had a contig N50 of 0.7 Mb versus 11.2 Mb
153 for the HiFi assembly. The ONT assembly also contained fewer complete BUSCOs (93%)
154 versus the HiFi assembly (95.6%). In both assemblies, the full-length *H-fibroin* locus was
155 present but the quality of the annotations differed greatly (Fig. 2). For ONT, Augustus annotated
156 a dozen genes in the ~30Kb region, most of which did not include the characteristic repeats
157 known from previous data. For HiFi, the annotation included a single gene with a single intron in
158 the n-terminus region. The second exon was large (25.3 Kb) and fully in-frame including a well-
159 resolved repetitive structure, giving high confidence in the accuracy of the assembly.
160

a Oxford Nanopore reads mapped against the consensus Oxford Nanopore *H-fibroin* assembly

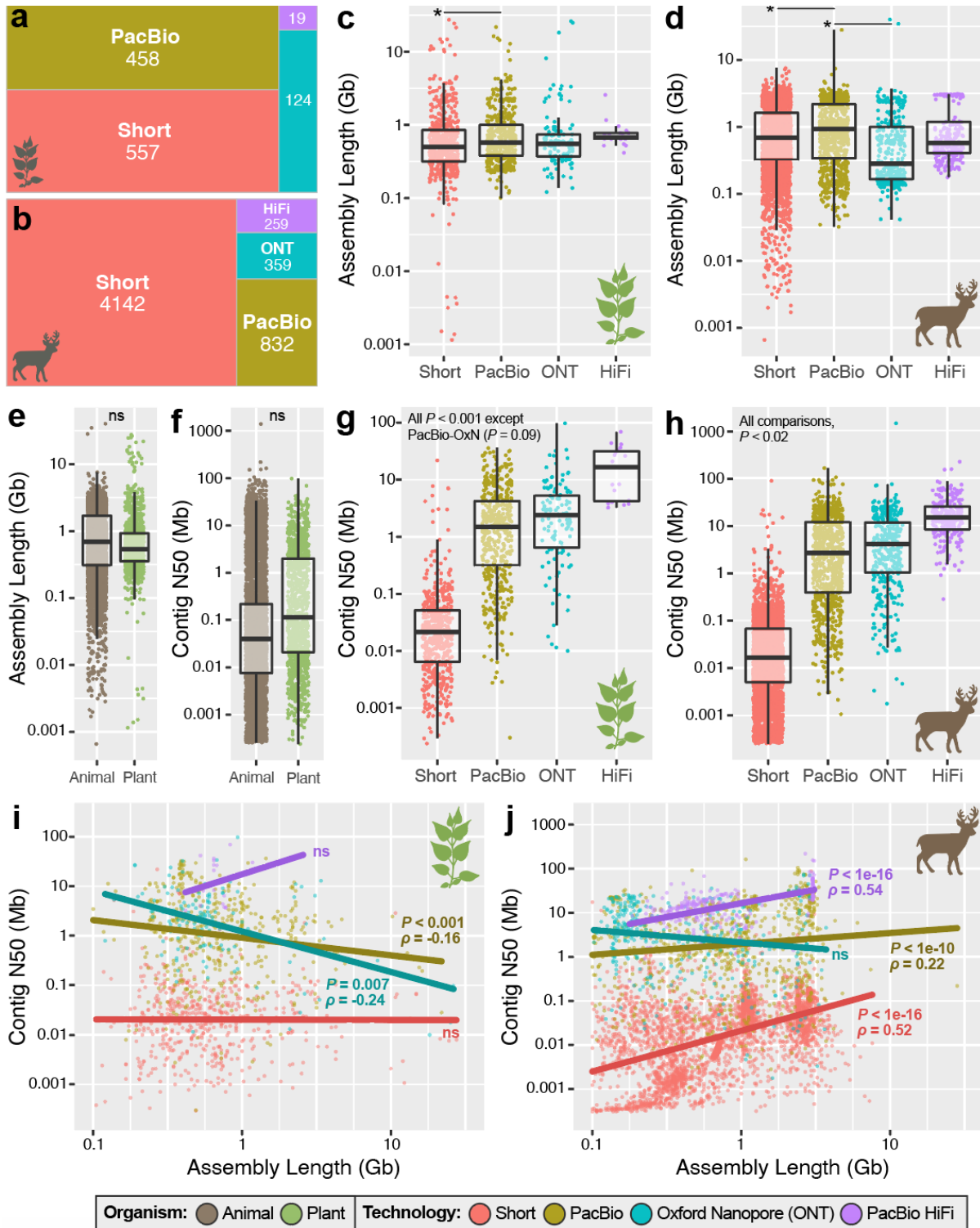


b PacBio HiFi reads mapped against the primary HiFi *H-fibroin* assembly



161
162 **Figure 2.** A case study comparing the capacity for two long-read sequencing technologies to assemble the
163 complex gene underlying silk production in caddisflies, *H-fibroin*. (a) Raw Oxford Nanopore (ONT) reads

164 mapped to the consensus *H-fibroin* sequence from the ONT assembly. (b) Raw HiFi reads mapped to the
 165 primary *H-fibroin* sequence from the phased HiFi assembly. Dark lines indicate mismatches relative to
 166 consensus. In (b), mismatches reflect an *H-fibroin* length polymorphism that can be resolved by subsampling
 167 reads based on their allele-specificity.
 168



169

170 **Figure 3.** Sequencing technology representation and genome assembly quality across all animal and plant
171 assemblies deposited in GenBank as of November 2021. A breakdown of the sequencing technology used for
172 genome assemblies in (a) plants and (b) animals. Total assembly length broken down by sequencing
173 technology for (c) plants and (d) animals. (e) Assembly length across all plant and animal genomes, regardless
174 of technology. (f) Contig N50 across all plant and animal genomes and broken down by technology for (g)
175 plants and (h) animals. Spearman's correlations between contig N50 and assembly length for (i) plants and (j)
176 animals. For (i) and (j), correlation statistics were generated for the full data sets but for visualization,
177 assemblies less than 0.1 Gb in length or with contig N50 > 1 Gb have been excluded. For (c-h), asterisks and
178 thin dark lines indicate significant differences at $P < 0.05$.

179
180 After filtering, our data set contained 6,750 genome assemblies (animals = 5,592; plants =
181 1,158; Table S1). For plants, short-read assemblies (48.1%; $N = 557$) and long-read assemblies
182 generated with non-HiFi PacBio data were similarly common (39.6%; $N = 458$; Fig. 3a). ONT
183 assemblies, however, were much less common (10.7%) and HiFi assemblies were exceptionally
184 rare, comprising just 1.6% of all assemblies ($N = 19$; Fig. 3a). For animals, the majority of
185 assemblies were generated with short-read data (74.1%; $N = 4,142$). Non-HiFi long-read
186 assemblies, generated with either PacBio or ONT reads, were also common, comprising 21.4%
187 of the data set. HiFi assemblies were again the least common with just 259 assemblies (4.6%;
188 Fig. 3b).

189
190 On average, animal genome assemblies were neither longer (P , Welch T-test = 0.80, Fig. 3e)
191 nor more contiguous than those of plants (P , Welch T-test = 0.10, Fig. 3f). When broken down
192 by technology, assembly lengths did differ for plants (P , one-way ANOVA = 0.006) and animals
193 (P , one-way ANOVA < 0.001). For plants, only one length comparison was significantly
194 different—non-HiFi PacBio assemblies were longer than those generated with short-reads (P ,
195 Tukey HSD = 0.005; Fig. 3c). For animals, non-HiFi PacBio assemblies were longer than both
196 short-read (P , Tukey HSD = 0.002) and ONT assemblies (P , Tukey HSD < 0.001; Fig. 3d). In
197 terms of assembly contiguity, contig N50 was significantly different for all comparisons in plants
198 (P , Tukey HSD < 0.001) and animals (P , Tukey HSD < 0.02) with one exception: in plants, non-

199 HiFi PacBio assemblies were not different from ONT assemblies (P , Tukey HSD = 0.09; Fig. 3g-
200 h). For both groups, HiFi reads dramatically outperformed all other long-read technologies
201 tested. In plants, the average HiFi assembly was 501% more contiguous (mean contig N50 =
202 20.5 Mb) than assemblies generated with other long-reads (mean contig N50 = 4.1 Mb; Fig. 3g).
203 For animals, HiFi assemblies were 226% more contiguous (mean contig N50 = 20.9 Mb) versus
204 other long-read assemblies (mean contig N50 = 9.3 Mb; Fig. 3h).

205
206 When assembly size was compared to contiguity, striking patterns emerged. For plants, when
207 non-HiFi long reads are used, contiguity declines with increasing assembly length ($P < 0.008$;
208 Spearman's ρ , PacBio = -0.16, Spearman's ρ , ONT = -0.24; Fig. 3i). The same trend isn't
209 present for short reads (P , Spearman's $\rho = 0.18$) nor HiFi (P , Spearman's $\rho = 0.37$). However,
210 for HiFi, this lack of significance is likely a product of small sample size ($N = 19$; Fig. 3i). For
211 animals, however, three of four read types (short, PacBio, HiFi) exhibit positive correlations
212 between contig N50 and assembly length (P , Spearman's $\rho < 1e-10$) with the steepest trends
213 for the two most accurate sequencing technologies: HiFi (Spearman's $\rho = 0.54$) and short-reads
214 (Spearman's $\rho = 0.52$; Fig. 3j).

215
216 As high-throughput sequencing technologies have matured, so has a common strategy: to
217 generate better assemblies, more data is better. While simplistic, this “more is better” approach
218 has been supported by empirical data and echoed by genome sequencing overviews (e.g.,
219 Ekblom and Wolf 2014). The practicality of this approach has also been empirically observed in
220 the long-read era. For instance, benchmarking of long-read assemblies in maize found that
221 lower sequencing depths ($< 30x$) with mean read lengths less than 11 Kb yielded highly
222 fragmented assemblies (Ou, et al. 2020). Stepping back, the outperformance of long-reads
223 relative to short-reads in terms of basic assembly metrics is dramatic and independent of
224 taxonomy (Hotaling, et al. 2021b; Marks, et al. 2021). Thus, support exists for the premise that

225 to generate the most high-quality assemblies, researchers should maximize depth of coverage
226 and mean read length regardless of the technology being used.

227

228 However, since we cannot maximize perfectly accurate deep sequencing with the longest
229 possible reads, achieving the best possible genome assemblies under the current landscape of
230 sequencing technologies appears to require more nuance. Our results highlight that in at least
231 one instance, a smaller amount of shorter, but more accurate, HiFi reads outperformed ONT
232 data for genome assembly in the same taxon. Since this outperformance likely stems from base
233 pair accuracy and the potential for sequence errors to confound assembly, we expected the
234 benefit of highly accurate reads to scale with genome and/or gene-region complexity. The
235 dramatic outperformance of HiFi reads when assembling the *H-fibroin* locus supports this
236 expectation. Similar results have also been obtained for other taxonomic groups. For the leaf
237 rust fungus, *Puccinia triticina*, HiFi reads outperformed noisier long-reads, particularly in areas
238 of the genome that were difficult to assemble (Duan, et al. 2022). Among rice genomes,
239 however, the results were less clear: ONT ultra-long reads yielded higher overall contiguity but
240 more errors than HiFi (Lang, et al. 2020).

241

242 Moving beyond our case study to a large-scale comparison of animal and plant genome
243 assemblies, key broader themes emerged. First, HiFi assemblies were significantly more
244 contiguous than all other types of sequence data tested. And, given the benefits it affords, highly
245 accurate long-read sequencing—e.g., HiFi—remains underrepresented, particularly in plant
246 genetics. Second, all significant relationships between genome size and contiguity for
247 sequencing technologies in animals were positive (Fig. 3j) whereas no correlation was positive
248 in plants (Fig. 3i). This suggests that as genome size increases in plants, so too does
249 complexity, likely at a rate that outpaces the capacity for modern assembly algorithms to
250 assemble it. Notable, however, was the sharply positive trend for HiFi sequencing in plants

251 where contig N50 appears to rapidly increase with assembly length (Fig. 3i). While not
252 statistically significant at $P < 0.05$ —likely due to a low sample size—this pattern suggests that HiFi
253 and similarly accurate long-read technologies may represent a valuable means to overcome
254 challenges of genome complexity in plants. Finally, we acknowledge that while we focused on
255 specific technologies in this study and found strong evidence for HiFi efficacy, our point is more
256 about long-read accuracy than specific technologies. Indeed, other technologies—including
257 ONT—will likely approach and may even surpass HiFi in terms of accuracy in the months and
258 years to come.

259
260 With the rise of highly accurate long-read sequencing, little room remains for another revolution
261 in genome assembly quality. It is now possible to generate reference-quality assemblies for
262 virtually any species with modest resources. The only exceptions, for now, are species that are
263 very small, difficult to obtain, and/or with exceptionally large or complex genomes. What is
264 lacking now are appropriate metrics for contemporary genome assembly benchmarking. For
265 instance, contig N50—the most common metric for assessing contiguity—scales with assembly
266 length in highly accurate long-read assemblies (Fig. 3i-j) and thus, its upper limit is tied to
267 chromosome length, making comparisons among groups difficult. For gene content
268 assessment—i.e., BUSCO scores—one challenge lies in how accurately BUSCO scores reflect
269 true gene content, particularly when more repeat-rich genes are considered. For instance, the
270 median gene length in the 1,467-gene “Insecta” gene set is ~1 Kb (longest = 9.1 Kb; Hotaling, et
271 al. 2021b) yet phenotypically relevant genes like *H-fibroin* can be *much* longer (i.e., >20 Kb).
272 Thus, while BUSCO scores for the ONT and HiFi assemblies of *H. magnus* reflect marginal
273 differences in their gene completeness (93% vs. 95.6%), the true gap is likely much greater. We
274 expect differences in assembly quality to scale with genomic region complexity—a result that is
275 not captured by BUSCO scores.

276

277 **Acknowledgements:**

278 J.H was supported by the LOEWE-Centre for Translational Biodiversity Genomics, which was
279 funded by the Hessen State Ministry of Higher Education, Research and the Arts.

280

281 **Data availability:**

282 Both genome assemblies generated for this study are available on Genbank
283 (#GCA_016648045.1, #JAIUSX000000000) as well as the raw reads used (ONT: SRX9290148;
284 Illumina: SRX9290147; HiFi: SRR15840267). The full data set used for the animal and plant
285 genome analysis is provided in Table S1.

286

287 **References:**

288 Amarasinghe SL, et al. 2020. Opportunities and challenges in long-read sequencing data
289 analysis. *Genome Biology* 21: 1-16.

290 Ashton NN, Roe DR, Weiss RB, Cheatham III TE, Stewart RJ 2013. Self-tensioning aquatic
291 caddisfly silk: Ca²⁺-dependent structure, strength, and load cycle hysteresis.
292 *Biomacromolecules* 14: 3668-3681.

293 Cheng H, Concepcion GT, Feng X, Zhang H, Li H 2021. Haplotype-resolved de novo assembly
294 using phased assembly graphs with hifiasm. *Nature Methods* 18: 170-175.

295 De Coster W, Weissensteiner MH, Sedlazeck FJ 2021. Towards population-scale long-read
296 sequencing. *Nature Reviews Genetics* 22: 572-587.

297 Duan H, et al. 2022. Physical separation of haplotypes in dikaryons allows benchmarking of
298 phasing accuracy in Nanopore and HiFi assemblies with Hi-C data. *Genome Biology* 23: 1-27.

- 299 Ekblom R, Wolf JB 2014. A field guide to whole-genome sequencing, assembly and annotation.
300 *Evol Appl* 7: 1026-1042. doi: 10.1111/eva.12178
- 301 Frandsen PB, et al. 2019. Exploring the underwater silken architectures of caddisworms:
302 comparative silkomics across two caddisfly suborders. *Philosophical Transactions of the Royal*
303 *Society B* 374: 20190206.
- 304 Hotaling S, Desvignes T, Sproul JS, Lins LS, Kelley JL 2022. Pathways to polar adaptation in
305 fishes revealed by long-read sequencing. *Mol Ecol*.
- 306 Hotaling S, Kelley JL, Frandsen PB 2021a. Toward a genome sequence for every animal:
307 Where are we now? *Proceedings of the National Academy of Sciences* 118: e2109019118.
- 308 Hotaling S, et al. 2021b. Long-reads are revolutionizing 20 years of insect genome sequencing.
309 *Genome Biology and Evolution* evab138.
- 310 Kawahara AY, et al. 2022. Long-read HiFi Sequencing Correctly Assembles Repetitive heavy
311 fibroin Silk Genes in New Moth and Caddisfly Genomes. *bioRxiv*.
- 312 Kearse M, et al. 2012. Geneious Basic: an integrated and extendable desktop software platform
313 for the organization and analysis of sequence data. *Bioinformatics* 28: 1647-1649.
- 314 Koren S, et al. 2017. Canu: scalable and accurate long-read assembly via adaptive k-mer
315 weighting and repeat separation. *Genome research* 27: 722-736.
- 316 Kriventseva EV, et al. 2019. OrthoDB v10: sampling the diversity of animal, plant, fungal, protist,
317 bacterial and viral genomes for evolutionary and functional annotations of orthologs. *Nucleic*
318 *Acids Res* 47: D807-D811.

- 319 Lang D, et al. 2020. Comparison of the two up-to-date sequencing technologies for genome
320 assembly: HiFi reads of Pacific Biosciences Sequel II system and ultralong reads of Oxford
321 Nanopore. *GigaScience* 9: giaa123.
- 322 Li H 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34: 3094-
323 3100.
- 324 Luo S, Tang M, Frandsen PB, Stewart RJ, Zhou X 2018. The genome of an underwater
325 architect, the caddisfly *Stenopsyche tienmushanensis* Hwang (Insecta: Trichoptera).
326 *GigaScience* 7: giy143.
- 327 Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM 2021. BUSCO update: novel and
328 streamlined workflows along with broader and deeper phylogenetic coverage for scoring of
329 eukaryotic, prokaryotic, and viral genomes. *Mol Biol Evol* 38: 4647-4654.
- 330 Marks RA, Hotaling S, Frandsen PB, VanBuren R 2021. Representation and participation
331 across 20 years of plant genome sequencing. *Nature Plants* 7: 1571-1578.
- 332 Nurk S, et al. 2020. HiCanu: accurate assembly of segmental duplications, satellites, and allelic
333 variants from high-fidelity long reads. *Genome research* 30: 1291-1305.
- 334 Ou S, et al. 2020. Effect of sequence depth and length in long-read assembly of the maize
335 inbred NC358. *Nat Commun* 11: 1-10.
- 336 R Core Team 2021. R: A language and environment for statistical computing.
- 337 Rhie A, et al. 2021. Towards complete and error-free genome assemblies of all vertebrate
338 species. *Nature* 592: 737-746.

- 339 Ruan J, Li H 2020. Fast and accurate long-read assembly with wtdbg2. *Nature Methods* 17:
340 155-158.
- 341 Shen W, Xiong J 2019. TaxonKit: a cross-platform and efficient NCBI taxonomy toolkit. *bioRxiv*:
342 513523.
- 343 Sims D, Sudbery I, Illott NE, Heger A, Ponting CP 2014. Sequencing depth and coverage: key
344 considerations in genomic analyses. *Nature Reviews Genetics* 15: 121-132.
- 345 Stanke M, et al. 2006. AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids*
346 *Res* 34: W435-W439.
- 347 Wickham H 2011. ggplot2. *Wiley Interdisciplinary Reviews: Computational Statistics* 3: 180-185.
- 348 Yonemura N, Mita K, Tamura T, Sehnael F 2009. Conservation of silk genes in Trichoptera and
349 Lepidoptera. *Journal of Molecular Evolution* 68: 641-653.
- 350 Zimin AV, et al. 2013. The MaSuRCA genome assembler. *Bioinformatics* 29: 2669-2677.
351