

1 **Title:** Long-read HiFi Sequencing Correctly Assembles Repetitive *heavy fibroin* Silk Genes in

2 New Moth and Caddisfly Genomes

3 **Authors:**

4 <sup>1,\*</sup>, † Akito Y. Kawahara, 0000-0002-3724-4610

5 <sup>1,2</sup>, † Caroline G. Storer, 0000-0002-0349-0653

6 <sup>1,3</sup> Amanda Markee

7 <sup>4,5</sup> Jacqueline Heckenhauer, 0000-0001-8771-9154

8 <sup>6</sup> Ashlyn Powell

9 <sup>1</sup> David Plotkin, 0000-0002-2339-655X

10 <sup>7</sup> Scott Hotaling, 0000-0002-5965-0986

11 <sup>8</sup> Timothy P. Cleland, 0000-0001-9198-2828

12 <sup>9</sup> Rebecca B. Dikow

13 <sup>10</sup> Torsten Dikow, 0000-0003-4816-2909

14 <sup>11,12</sup> Ryoichi B. Kuranishi, 0000-0002-6353-0450

15 <sup>1</sup> Rebeccah Messcher

16 <sup>4,5,13</sup> Steffen U. Pauls, 0000-0002-6451-3425

17 <sup>14</sup> Russell J. Stewart, 0000-0002-8389-8877

18 <sup>15</sup> Koji Tojo

19 <sup>6,9,\*</sup> Paul B. Frandsen, 0000-0002-4801-7579 (Last & corresponding author)

20

21

22

23

24 **Affiliations:**

25 <sup>1</sup>McGuire Center for Lepidoptera and Biodiversity, Florida Museum of Natural History,

26 University of Florida, Gainesville, FL 32611, USA

27 <sup>2</sup>Pacific Biosciences, 1305 O'Brien Dr., Menlo Park, CA 94025 USA

28 <sup>3</sup>School of Natural Resources and the Environment, University of Florida, Gainesville, FL

29 32611, USA

30 <sup>4</sup>LOEWE Centre for Translational Biodiversity Genomics (LOEWE-TBG), Frankfurt 60325,

31 Germany

32 <sup>5</sup>Department of Terrestrial Zoology, Senckenberg Research Institute and Natural History

33 Museum Frankfurt, Frankfurt 60325, Germany

34 <sup>6</sup>Department of Plant and Wildlife Sciences, Brigham Young University, Provo, UT 84602, USA

35 <sup>7</sup>School of Biological Sciences, Washington State University, Pullman, WA, USA

36 <sup>8</sup>Museum Conservation Institute, Smithsonian Institution, Suitland, MD 20746, USA

37 <sup>9</sup>Data Science Lab, Office of the Chief Information Officer, Smithsonian Institution,

38 Washington, DC 20002, USA

39 <sup>10</sup>Department of Entomology, National Museum of Natural History, Smithsonian Institution,

40 Washington, DC, USA

41 <sup>11</sup>Graduate School of Science, Chiba University, Chiba 263-8522, Japan

42 <sup>12</sup>Kanagawa Institute of Technology, Kanagawa 243-0292, Japan

43 <sup>13</sup>Institute for Insect Biotechnology, Justus-Liebig-University, Gießen 35390, Germany

44 <sup>14</sup>Department of Biomedical Engineering, University of Utah, Salt Lake City, UT 84112, USA

45 <sup>15</sup>Department of Biology, Shinshu University, Nagano, Japan

46 †Co-first authors

47 **Abstract** (250 words maximum)

48 Insect silk is an incredibly versatile biomaterial. Lepidoptera and their sister lineage, Trichoptera,  
49 display some of the most diverse uses of silk with varying strength, adhesive qualities and elastic  
50 properties. It is well known that silk fibroin genes are long (> 20 kb) and have many repetitive  
51 motifs. These features make these genes challenging to sequence. Most research thus far has  
52 focused on conserved N- and C-terminal regions of fibroin genes because a full comparison of  
53 repetitive regions across taxa has not been possible. Using the PacBio Sequel II system and  
54 SMRT sequencing, we generated high fidelity (HiFi) long-read genomic and transcriptomic  
55 sequences for the Indianmeal moth (*Plodia interpunctella*) and genomic sequences for the  
56 caddisfly, *Eubasilissa regina*. Both genomes were highly contiguous (N50 = 9.7 Mbp/32.4 Mbp,  
57 L50 = 13/11) and complete (BUSCO Complete = 99.3%/95.2%), with complete and contiguous  
58 recovery of silk *heavy fibroin* gene sequences. This study demonstrates that HiFi long-read  
59 sequencing can significantly help our understanding of genes with highly contiguous, repetitive  
60 regions.

61

62 **Keywords:**

63 Caddisfly, data description, Indianmeal moth, Lepidoptera, moth, PacBio, transcriptome,  
64 Trichoptera

65

66 **Main Content**

67 **Data Description**

68 **Background**

69 Many phenotypic traits across the tree of life are controlled by repeat-rich genes [1].  
70 There are many examples, such as antifreeze proteins in fish [2], keratin in mammals, and resilin  
71 in insects [1]. Silk is a fundamental biomaterial that is produced by many arthropods, and silk  
72 genes are often long (> 20 kb) and contain repetitive motifs [3]. Accurately sequencing through  
73 repeat-rich genomic regions is critical to understand how functional genes dictate phenotypes,  
74 but research thus far has been unable to quantify these regions. In the case of silk genes, this is  
75 essential as these regions control the strength and elasticity properties of silk fibers [4-6].

76 Lepidoptera (moths and butterflies) and their sister lineage Trichoptera (caddisflies)  
77 display some of the most diverse uses of silk from spinning cocoons to prey capture nets and  
78 protective armorment [7]. A complete heavy fibroin for the model silkworm moth, *Bombyx mori*,  
79 was assembled over two decades ago using BAC libraries [8]. Recently, a combination of  
80 nanopore and Illumina sequencing technologies helped generate a full heavy fibroin sequence of  
81 *B. mori*, but large regions of the genome remain unassembled [3]. Similarly, we have had similar  
82 problems from the Nanopore and Illumina hybrid assemblies in caddisfly genomes [e.g., 9],  
83 where we were unable to assemble the complete H-Fibroin genes despite intensive efforts for  
84 ~20 species. In these assemblies, the biggest hindrance was sequencing single strands across  
85 large repeat regions and limited illumina polishing due to higher error rates in Nanopore data.  
86 The lack of full coverage was largely due to the fact that Nanopore and Illumina sequencing  
87 approaches introduce uncertainty for direct inference of function. Therefore, most research thus  
88 far has been limited, and focused only on conserved N- and C-terminal regions [e.g., 10].  
89 Complete high-fidelity fully phased fibroin sequences are critical for advancing biomaterials  
90 discovery for insect silks.

91

## 92 **Context**

93 We generated HiFi long-read genomic sequences for the Indianmeal moth (*Plodia*  
94 *interpunctella*), and the caddisfly species *Eubasilissa regina*, with the PacBio Sequel II system.  
95 Our goal was to recover the area of the genome that has been nearly impossible to sequence due  
96 to its repeated regions. We chose these two taxa as they represent two species with very different  
97 life histories – *Plodia interpunctella* is an important model organism in Lepidoptera whose  
98 larvae feed on a wide variety of grains and stored food products, and secrete large amounts of  
99 thin silken webbing at their feeding sites; they also use silk to create a cocoon during pupation  
100 [11-12]. *Eubasilissa regina*, on the other hand, is a member of Trichoptera, whose larvae secrete  
101 silk in aquatic environments in order to produce protective silk cases made of broader leaf pieces  
102 from deciduous trees, cut to size [13]. These new resources not only expand our knowledge of a  
103 primary silk gene in Lepidoptera and Trichoptera, but also contribute new, high-quality genomic  
104 resources for aquatic insects and arthropods which have thus far been underrepresented in  
105 genome biology [14-16].

106

## 107 **Methods**

### 108 **Sample information and sequencing**

109 A single adult specimen of each species was sampled for inclusion in the present study.  
110 For *P. interpunctella*, we used a specimen from the PiW3 colony line at the USDA lab (1600 SW  
111 23 Dr. Gainesville, FL, USA), and its entire body was used for extraction, given its small size.  
112 For *E. regina*, a wild-caught female adult specimen (#AK0WP01) from Enzan, Yamanashi,  
113 Japan (N35°43'24" E138°50'33", elevation ~4,840 ft), originally deposited in the Smithsonian  
114 Institution, National Museum of Natural History (USNMENT01414923), was used. The head

115 and thorax were macerated and DNA was extracted. The remainder of the body is preserved as a  
116 frozen tissue sample in the lab of PRB at BYU. Both specimens were flash frozen in LN2 and  
117 DNA was extracted using Quick-DNA HMW MagBead Kit (Zymo Research). Extractions with  
118 at least 1  $\mu\text{g}$  of high-molecular weight ( $> 40\text{kb}$ ) were sheared and the BluePippin system (Sage  
119 Science, Beverly, MA, USA) was used to collect fractions containing 15 kb fragments for library  
120 preparation. Sequencing libraries were prepared for each species using the SMRTbell Express  
121 Template Prep Kit 2.0 (PacBio, Menlo Park, CA, USA) and following the ultra-low protocol. All  
122 sequencing was performed using the PacBio Sequel II system. For *P. interpunctella*, the genomic  
123 library was sequenced on a single 8M SMRTcell and *E. regina* was sequenced on three 8M  
124 SMRTcells, all with 30 hour movie times. For the *P. interpunctella* Iso-seq transcriptome, RNA  
125 was extracted using TRIzol (Invitrogen) from freshly dissected silk glands of caterpillars and  
126 following manufacturer's protocol. This species has relatively small body size compared to other  
127 Lepidoptera, so we waited until caterpillars reached their maximum size (during the fifth instar)  
128 before dissection, in order to maximize yield. Sequencing libraries were prepared following the  
129 PacBio IsoSeq Express 2.0 Workflow and using the NEBNext Single Cell/Low Input cDNA  
130 Synthesis & Amplification Module for the SMRTbell Express Template Prep Kit 2.0. The  
131 resulting library was multiplexed and sequenced on a single Sequel II PacBio SMRT cell for 30  
132 hrs. Library preparation and sequencing was carried out at DNA Sequencing Center at Brigham  
133 Young University (Provo, UT, USA).

134 Genomic HiFi reads were generated by circular consensus sequencing (CCS) where  
135 consensus sequences have three or more passes with quality values equal to or greater than 20,  
136 from the subreads.bam files and using pbccs tool (v.6.0.0) in the *pbccs* package  
137 (<https://github.com/PacificBiosciences/pbccs>). Using the same *pbccs* package and

138 the Iso-seq v3 tools, high quality (> Q30) transcripts were generated from HiFi read clustering  
139 without polishing.

140

#### 141 **Genome size estimations and genome profiling**

142 Estimation of genome characteristics such as size, heterozygosity, and repetitiveness were  
143 conducted using a *k*-mer distribution-based approach. After counting *k*-mers with K-Mer  
144 Counter (KMC) v.3.1.1 and a *k*-mer length of 21 (-m 21), we generated a histogram of *k*-mer  
145 frequencies with KMC transform histogram [17]. We then generated genome *k*-mer profiles on  
146 the *k*-mer count histogram using the GenomeScope 2.0 web tool [18] with the *k*-mer length set to  
147 21 and the ploidy set to 2.

148

#### 149 **Sequence assembly and analysis**

150 For both genomes, reads were then assembled into contigs using the assembler Hifiasm  
151 v0.13-r307 with aggressive duplicate purging enabled (option -l 2) [19]. The primary contig  
152 assembly was used for all downstream analyses. Genome contiguity was measured using  
153 *assembly\_stats.py* [20] and genome completeness was determined using Busco v.5.2.2 [21] and  
154 the obd10 reference Endopterygota. Contamination in the genome was assessed by creating  
155 Taxon-annotated GC-coverage (TAGC) plots using BlobTools v1.0 [22]. First, assemblies were  
156 indexed using *samtools faidx* then HiFi reads were mapped back to the indexed assemblies using  
157 *minimap2* [23] with *-ax asm20*. The resulting bam files were sorted with *samtools sort*.  
158 Taxonomic assignment was performed via Megablast and using the NCBI nucleotide database  
159 with parameters *-outfmt 6 qseqid staxids bitscore std' -max\_target\_seqs 1 -max\_hsps 1 -e value*  
160 *1e-25*. BlobPlots were created by making a blobtools database from the assembly file, blast  
161 results, and mapping results using *blobtools create* and plots were created using *blobtools plot*.

162

### 163 **Genome statistics**

164 All samples, raw sequence reads, and assemblies were deposited to GenBank (Table 1).

165 We generated 35.7 Gbp (41x coverage) and 15.7 Gbp (44x coverage) of PacBio HiFi sequence

166 for *E. regina* and *P. interpunctella*, respectively. We assembled those reads into two contiguous

167 genome assemblies. The assembly for *E. regina* has the highest contig N50 of any Trichoptera

168 genome assembly to date. It contains 123 contigs, a contig N50 of 32.4 Mbp, GC content of

169 32.68%, and a total length of 917,780,411 bp. GenomeScope 2.0 estimated a genome size of

170 854,331,742 bp with 75.3% unique sequence

171 (<http://genomescope.org/genomescope2.0/analysis.php?code=ghDHLpAQUkIKK4e5yH88>).

172 Despite recent analyses showing no evidence of whole-genome duplication in caddisflies

173 (Heckenhauer et al. 2022), the findings in this study may be an indication of tetraploidy. Future

174 research should be done to further examine these patterns.

175 The *P. interpunctella* assembly represents a substantial improvement to existing, publicly

176 available genome assemblies (Table 2). After contaminated contigs were removed (three contigs

177 contaminated with *Wolbachia* were identified), the resulting assembly comprises 118 contigs

178 with a cumulative length of 300,731,903 bp. It exhibits a contig N50 of 9.7Mbp and a GC

179 content of 35.41%. The genome size estimated by GenomeScope 2.0 was 275,458,564 bp with

180 87.1% unique sequence (<http://genomescope.org/genomescope2.0/>

181 [analysis.php?code=96nVnnk42W5nlBWlFHFj](http://genomescope.org/genomescope2.0/analysis.php?code=96nVnnk42W5nlBWlFHFj)).

182

### 183 **Heavy fibroin gene annotation**



184 We extracted *heavy fibroin (H-fibroin)* silk genes from both the *P. interpunctella* and *E.*  
185 *regina* assemblies. For *P. interpunctella*, we also searched existing, short-read based assemblies.  
186 We downloaded two short-read based genome assemblies for *P. interpunctella*,  
187 GCA\_001368715.1 and GCA\_900182495.1 from NCBI (<https://www.ncbi.nlm.nih.gov/>). Since  
188 the internal region of *H-fibroin* is known to be repetitive, the more conserved N- and C-termini  
189 amino acids were blasted against the genomes with tblastn. For *P. interpunctella*, we used the  
190 terminal sequences published in [24] and for *E. regina*, we used the terminal sequences  
191 published in [5]. We then extracted the sequences and 500 bps of flanking regions from the  
192 assembly and annotated them using Augustus v.3.3.2 [25]. Spurious introns (those that did not  
193 affect reading frames and were not supported by transcript evidence) were manually removed.  
194 Annotated sequences are provided in the *Gigascience* GigaDB repository [26].

195 We recovered full-length *H-fibroin* sequences in both genomes. To our knowledge, the  
196 only other previously published full-length lepidopteran *H-fibroin* sequence was from a BAC  
197 library-based sequence of the model organism, *B. mori*. We compared our assembly of the *P.*  
198 *interpunctella* *H-fibroin* sequence with that from a previously published Illumina-based genome  
199 assembly of the same species (Table 2). Where the Illumina-based assembly only recovered the  
200 conserved terminal regions and a small number of repetitive elements, our assembly recovered  
201 the full-length gene, including the full complement of repetitive motifs (Figures 1, 2).  
202 Specifically, the *P. interpunctella* genome had a *H-fibroin* sequence that was 14,866 bp (whole  
203 gene with introns; 4,714 AA), and a molecular weight of 413,334.41 Da. For *E. regina*, we  
204 recovered the full-length sequence of *H-fibroin*, which was 25,250 bp (whole gene with introns;  
205 8,386 AA), and a molecular weight of 815,864.95 Da, with repeated regions (Figure 3). The  
206 recovery of this *H-fibroin* sequence marks the third complete, published *H-fibroin* sequence in

207 Trichoptera [27-28]. Our work shows that high quality, long-read sequencing can be used to  
208 successfully assemble difficult regions of non-model organisms without the use of expensive and  
209 tedious BAC methods. While our study is focused on the repetitive silk gene, *H-fibroin*, these  
210 results likely extend to other long, repetitive proteins that have previously proven difficult to  
211 assemble.

212

### 213 **Genome annotation**

214 For the structural annotations of the genomes, we masked and annotated repetitive  
215 elements using RepeatMasker [29] after identifying and classifying them *de novo* with  
216 RepeatModeler2 [30] following [31]. For species specific gene model training, we used BUSCO  
217 v.4.1.4 [21] with the Endopterygota odb10 core ortholog sets [32] with the `-long` option in  
218 genome mode. In addition, we predicted genes with the homology-based gene prediction  
219 GeMoMaPipeline of GeMoMa v1.6.4 [33-34] using previously published genomes. For *E.*  
220 *regina* we used the genome of *Agypnia vestita* (JADDOH000000000.1) [35] and for *P.*  
221 *interpunctella* we used the genome of *Bombyx mori* (GCF\_014905235) as reference. We then  
222 used the MAKER v3.01.03 pipeline [36] to generate additional *ab initio* gene predictions with  
223 the proteins predicted from GeMoMa for protein homology evidence and the augustus-generated  
224 gene prediction models from BUSCO for gene prediction. For EST evidence, we used the  
225 transcriptome of *Ptilostomis semifasciata* (111015\_I297\_FCD05HRACXX\_  
226 L1\_INSbttTHRAAPEI-17, 1kite.org) for *E. regina* and Iso-seq data for *P. interpunctella*.  
227 Evidence used in Maker and the Maker config files can be found in the *Gigascience* GigaDB  
228 repository [26].

229 To add functional annotations to the predicted proteins, we blasted the predicted proteins  
230 against the ncbi-blast protein database using BlastP in blast.2.9 with an e-value cutoff of  $10^{-4}$   
231 and `-max_target_seqs` set to 10 (see repository). We then used the command line version of  
232 Blast2GO v.1.4.4 [37] to assign functional annotation and GO terms.

233

### 234 **Validation and quality control**

235 In addition to full-length *H-fibroin* sequences, we recovered a high number of single copy  
236 orthologs in each genome with BUSCO. The *E. regina* genome contained 95.2% of an  
237 Endopterygota core gene collection (comprised of 2124 genes) indicating an almost complete  
238 coverage of known single copy orthologs in the coding fraction. While the number of single-  
239 copy orthologs recovered in the new *P. interpunctella* genome was similar to earlier published  
240 genomes (99.3% of the Endopterygota core gene collection, 99.1% of the Lepidoptera core gene  
241 collection), the full-length sequence of *H-fibroin* only recovered in the HiFi based genome gives  
242 some indication of how other portions of the genome may have assembled. Following  
243 contamination screening by NCBI, we filtered out three instances of *Wolbachia* contamination in  
244 the *P. interpunctella* genome. BlobPlots for both genomes revealed low levels of contamination  
245 (Supplementary Figures 1, 2).

246

### 247 **Structural and functional annotation**

248 A total of 56.26% of the *E. regina* genome was classified as repetitive (54.2% interspersed  
249 repeats). More than half of the interspersed repeats, 29.87%, could not be classified by  
250 comparison with known repeat databases, and therefore may be specific for Trichoptera. Of the  
251 repeats that were classified, retroelements were the most abundant, comprising 15.35% (of which

252 14.55% are LINEs) of the genome. The relatively high proportion of repetitive sequence  
253 supports previous studies which suggest that repetitive element expansion occurred in lineages of  
254 tube case-making caddisflies, such as the closely related genera *Agrypnia* and *Hesperophylax* [9,  
255 35]. In contrast, a total of 31.94% of the *P. interpunctella* genome assembly was masked as  
256 repeats. A total of 23.04% of the annotated repeats were interspersed repeats. Details on the  
257 repeat classes are given in the *Gigascience* GigaDB repository [26].

258 The genome annotations resulted in the prediction of 16,937 and 60,686 proteins in *P.*  
259 *interpunctella* and *E. regina*, respectively. Of the annotated proteins, for *E. regina*, 28,358  
260 showed significant sequence similarity to entries in the NCBI nr database, of those 12,550 were  
261 mapped to GO terms, and 5,652 were functionally annotated with Blast2GO. For *P.*  
262 *interpunctella*, 16,349 were verified by BLAST, 12,410 were mapped to GO terms, and 9,711  
263 were functionally annotated in Blast2GO.

264 The major biological process found in the two genomes were cellular (*E. regina*: 2,326  
265 genes; *P. interpunctella*: 4,725 genes) and metabolic (*E. regina*: 2,454 genes; *P. interpunctella*:  
266 3,699 genes) processes. Binding (*E. regina*: 2,382 genes; *P. interpunctella*: 4,405 genes) and  
267 catalytic activity (*E. regina*: 2,778 genes; *P. interpunctella*: 3,893 genes) were the largest  
268 subcategories in molecular function. Regarding the cellular component category, most genes  
269 were assigned to the cell (1,553 genes) and membrane (1,491 genes) subcategory in *E. regina*  
270 and to the cellular anatomical entity subcategory in *P. interpunctella* (5,602 genes). The major  
271 biological process found in both genomes were cellular and metabolic processes.

## 272 **Re-use potential**

273 We provide a complete genome of two species of silk-producing insects in the superorder  
274 Amphiesmenoptera, the moth *P. interpunctella* and the caddisfly *E. regina*, and recover the

275 difficult-to-sequence repetitive regions of both genomes with HiFi sequencing. *P. interpunctella*  
276 is currently being developed in multiple labs as a model organism and this genome assembly will  
277 facilitate molecular genetics research on this species. We show that PacBio HiFi sequencing  
278 allows for accurate generation of repetitive protein-coding regions of the genome (silk *fibroins*),  
279 and this likely applies to other similarly repetitive regions of the genome. For Trichoptera, there  
280 are only four other HiFi genome assemblies available on Genbank, only one of which has been  
281 published [38] and insects have generally been neglected (relative to their total species diversity)  
282 with respect to genome sequencing efforts [15-16], which is especially true for aquatic insects  
283 [14]. These data serve as the first step to study the evolution of adhesive silk in  
284 Amphiesmenoptera, which is an innovation that is beneficial for survival in aquatic and  
285 terrestrial environments. Finally, the Iso-seq data that we provide serve as useful resources for  
286 the translational aspects of silk – these data provide information on how Amphiesmenoptera  
287 genetically modulate and regulate different silk properties, that allow them using silk for  
288 different purposes such as for nets, cases, and cocoons in both terrestrial and aquatic  
289 environments.

#### 290 **Availability of source code and requirements**

291 All custom-made scripts used in this study are available on GitHub  
292 (<https://github.com/AshlynPowell/silk-gene-visualization/tree/main>).

#### 293 **Availability of supporting data**

294 Raw sequence data, genome assemblies, and sample information are all available from NCBI  
295 (accession can be found in Table 1). All supporting data and materials are available on GigaDB.

#### 296 **Declarations**

297 All authors have nothing to declare.

298 **Competing interests**

299 The authors declare that they have no competing interests.

300

301

302 **Funding**

303 Smithsonian National Museum of Natural History Global Genome Initiative (GGI-Peer-2018-

304 182) to T.P.C., R.D., T.D., A.Y.K., Smithsonian Museum Conservation Institute Federal and

305 Trust funds to T.P.C.. and P.B.F. University of Florida Research Opportunity Seed Fund internal

306 award number AWD06265 to PIs AYK and CGS.

307 **Authors' contributions**

308 AYK: Designed project, collected samples, provided computational resources, manuscript  
309 writing.

310 CGS: Designed project, data analysis, manuscript writing.

311 AM: Sample preparation, manage colonies, manuscript writing.

312 JH: Data analysis, manuscript writing.

313 AP: Data analysis, manuscript writing.

314 DP: Data file management, manuscript writing.

315 SH: Visualization, manuscript writing.

316 TPC: Grant writing, manuscript writing.

317 RBD: Grant writing, manuscript writing.

318 TD: Grant writing, manuscript writing.

319 RBK: Collected samples, manuscript writing.

320 RM: Helped with sample preparation, manage colonies, manuscript writing.

321 SUP: provided computational resources, manuscript writing.

322 RJS: Grant writing, manuscript writing.

323 KT: Collected samples, manuscript writing.

324 PBF: Designed project, collected samples, conducted analyses, provided computational  
325 resources, manuscript writing.

### 326 **Acknowledgements**

327 We thank the BYU, UF and LOEWE-Centre for Translational Biodiversity Genomics (TBG)

328 High Performance Computing clusters for providing the computational resources needed to

329 complete this study. TBG is funded by the Hessen State Ministry of Higher Education, Research

330 and the Arts (HMWK), which financially supported J.H. and S.U.P.. S.H. was supported by NSF

331 award #OPP-1906015.

### 332 **Authors' information**

333 AM: Graduate Student at the University of Florida, School of Natural Resources and

334 Environment. Studies variation in lepidopteran silk production.

335 AP: Undergraduate Student at Brigham Young University. Studies the genomics of caddisfly

336 silk.

337 AYK: Associate Curator at the University of Florida. Works on Lepidoptera, genomics, and

338 evolution.

339 CGS: formally Assistant Scientist at the University of Florida. Expert on silks and Lepidoptera.

340 Currently works at Pacific Biosciences.

341 DP: Project Manager at the Florida Museum of Natural History, University of Florida. Works on  
342 Lepidoptera systematics and evolution.

343 JH: Postdoctoral researcher at LOEWE TBG, interested in biodiversity genomics, comparative  
344 genomics, evolution and phylogenomics.

345 KT: Researcher at Shinshu University, Japan, specialist of Trichoptera.

346 PBF: Assistant Professor of Genetics, Genomics, and Biotechnology at Brigham Young  
347 University, specialist on the genomics of caddisflies and their silk.

348 RBD: Data Scientist at Smithsonian National Museum, Washington, D.C.

349 RJS: Professor of Bioengineering, University of Utah, specialist on the biomechanics of  
350 caddisfly silk.

351 RBK: Guest Professor of Kanagawa Institute of Technology, specialist of Trichoptera.

352 RM: Biological Scientist at the University of Florida researching molecular biology and genetics  
353 of organisms.

354 SH: Postdoctoral Research Associate at Washington State University and a specialist in insect  
355 genomics.

356 SUP: Entomologist at Senckenberg Research Institute and Natural History Museum Frankfurt  
357 and Justus-Liebig-University, interested in the evolution and ecology of freshwater insects,  
358 especially Trichoptera.

359 TD: Researcher and curator at Smithsonian National Museum of Natural History, Washington,  
360 D.C., Works on biodiversity of flies and is curator for aquatic insects.

361 TPC: Physical Scientist at the Smithsonian Museum Conservation Institute.

362



363 **References**

- 
- 364 1 Numata K. How to define and study structural proteins as biopolymer materials. *Polym. J.*,  
365 2020; 52: 1043–1056. doi:10.1038/s41428-020-0362-5.
- 366 2 Davies PL, Hew CL. Biochemistry of fish antifreeze proteins. *FASEB J*, 1990; 4: 2460–  
367 2468. doi:10.1096/fasebj.4.8.2185972
- 368 3 Kono N, Nakamura H, Ohtoshi R et al. The bagworm genome reveals a unique fibroin  
369 gene that provides high tensile strength. *Commun. Biol.*, 2019; 2: 1–9. doi:10.1038/s42003-  
370 019-0412-8.
- 371 4 Stewart RJ, Wang CS. Adaptation of caddisfly larval silks to aquatic habitats by  
372 phosphorylation of H-fibroin serines. *Biomacromolecules*, 2010; 11: 969–974.  
373 doi:10.1021/bm901426d.
- 374 5 Ashton NN, Roe DR, Weiss RB et al. Self-tensioning aquatic caddisfly silk: Ca<sup>2+</sup>-  
375 dependent structure, strength, and load cycle hysteresis. *Biomacromolecules*, 2013; 14: 3668–  
376 3681. doi:10.1021/bm401036z.
- 377 6 You Z, Ye X, Ye L et al. Extraordinary mechanical properties of composite silk through  
378 heritable transgenic silkworm expressing recombinant major ampullate spidroin. *Sci. Rep.*,  
379 2018; 8: 1–4. doi:10.1038/s41598-018-34150-y.
- 380 7 Sutherland TD, Young JH, Weisman S et al. Insect silk: one name, many materials. *Annu.*  
381 *Rev. Entomol.*, 2010; 55: 171–188. doi:10.1146/annurev-ento-112408-085401.
- 382 8 Zhou CZ, Confalonieri F, Medina N et al. Fine organization of *Bombyx mori* fibroin heavy  
383 chain gene. *Nucleic Acids Res.*, 2000; 28: 2413–2419. doi:10.1093/nar/28.12.2413.
- 384 9 Heckenhauer J, Frandsen PB, Sproul JS et al. Genome size evolution in the diverse insect  
385 order Trichoptera. *GigaScience*, 2022; 11: giac011. doi:10.1093/gigascience/giac011.
-

- 386 10 Yonemura N, Mita K, Tamura T et al. Conservation of silk genes in Trichoptera and  
387 Lepidoptera. J. Mol. Evol., 2009; 68: 641–653. doi:10.1007/s00239-009-9234-5.
- 388 11 Rutschky CW, Calvin D. Indian meal moth. 1990; [https://extension.psu.edu/indian-meal-](https://extension.psu.edu/indian-meal-moth)  
389 moth. Accessed March 2022.
- 390 12 Fasulo TR, Knox MA. Indianmeal moth - *Plodia interpunctella* (Hübner). 1998;  
391 [https://entnemdept.ufl.edu/creatures/urban/stored/indianmeal\\_moth.HTM](https://entnemdept.ufl.edu/creatures/urban/stored/indianmeal_moth.HTM). Accessed March  
392 2022.
- 393 13 Wiggins, G. B. The caddisfly family Phryganeidae (Trichoptera)., 1998; University of  
394 Toronto Press, Toronto, Buffalo, London.
- 395 14 Hotaling S, Kelley JL, Frandsen PB. Aquatic insects are dramatically underrepresented in  
396 genomic research. Insects, 2020; 11: 601. doi:10.3390/insects11090601.
- 397 15 Hotaling S, Sproul JS, Heckenhauer J et al. Long reads are revolutionizing 20 years of  
398 insect genome sequencing. Genome Biol. Evol., 2021; 13: evab138. doi:10.1093/gbe/evab138.
- 399 16 Hotaling S, Kelley JL, Frandsen PB. Toward a genome sequence for every animal: Where  
400 are we now? P. Natl. Acad. Sci. USA, 2021; 118: e2109019118.  
401 doi:10.1073/pnas.2109019118.
- 402 17 Kokot M, Długosz M, Deorowicz S. KMC 3: counting and manipulating k-mer statistics.  
403 Bioinformatics, 2017; 33: 2759–2761. doi:10.1093/bioinformatics/btx304.
- 404 18 Ranallo-Benavidez TR, Jaron KS, Schatz MC. GenomeScope 2.0 and Smudgeplot for  
405 reference-free profiling of polyploid genomes. Nat. Commun., 2020; 11: 1432.  
406 doi:10.1038/s41467-020-14998-3.
-

- 407 19 Cheng H, Concepcion GT, Feng X et al. Haplotype-resolved de novo assembly using  
408 phased assembly graphs with hifiasm. *Nat. Methods*, 2021; 18: 170–175. doi:10.1038/s41592-  
409 020-01056-5.
- 410 20 Trizna M. assembly\_stats 0.1.4 | Zenodo. 2020;  
411 [https://zenodo.org/record/3968775#.Yi\\_b4C1h2FU](https://zenodo.org/record/3968775#.Yi_b4C1h2FU). Accessed March 2022.  
412 doi:10.5281/zenodo.3968775.svg.
- 413 21 Manni M, Berkeley MR, Seppey M et al. BUSCO update: novel and streamlined  
414 workflows along with broader and deeper phylogenetic coverage for scoring of eukaryotic,  
415 prokaryotic, and viral genomes. *Mol. Biol. Evol.*, 2021; 38: 4647–4654.  
416 doi:10.1093/molbev/msab199.
- 417 22 Laetsch DR, Blaxter ML. BlobTools: Interrogation of genome assemblies [version 1; peer  
418 review: 2 approved with reservations]. *F1000Research*, 2017; 6: 1287.  
419 doi:10.12688/f1000research.12232.1.
- 420 23 Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 2018; 34:  
421 3094–3100. doi:10.1093/bioinformatics/bty191.
- 422 24 Kludkiewicz B, Kucerova L, Konikova T et al. The expansion of genes encoding soluble  
423 silk components in the greater wax moth, *Galleria mellonella*. *Insect Biochem. Molec.*, 2019;  
424 106: 28–38. doi:10.1016/j.ibmb.2018.11.003.
- 425 25 Stanke M, Diekhans M, Baertsch R et al. Using native and syntenically mapped cDNA  
426 alignments to improve de novo gene finding. *Bioinformatics*, 2008; 24: 637–644.  
427 doi:10.1093/bioinformatics/btn013.
- 428 26 Kawahara AY, Storer CG, Markee A et al. Supporting data for “Long-read HiFi  
429 sequencing correctly assembles repetitive *heavy fibroin* silk genes in new moth and caddisfly

430 genomes". GigaScience Database. 2022. doi:XXXX. [Repository for supplemental files for  
431 this paper]

432 27 Luo S, Tang M, Frandsen PB et al. The genome of an underwater architect, the caddisfly  
433 *Stenopsyche tienmushanensis* Hwang (Insecta: Trichoptera). GigaScience, 2018; 7: giy143.  
434 doi:10.1093/gigascience/giy143.

435 28 Frandsen PB, Bursell MG, Taylor AM et al. Exploring the underwater silken architectures  
436 of caddisworms: comparative silkomics across two caddisfly suborders. Philos. T. R. Soc. B,  
437 2019; 374: 20190206. doi:10.1098/rstb.2019.0206.

438 29 Smit AFA, Hubley R, Green P. RepeatMasker Open-4.0. 2013–2015;  
439 <http://www.repeatmasker.org>. Accessed January 2022.

440 30 Flynn JM, Hubley R, Goubert C et al. RepeatModeler2 for automated genomic discovery  
441 of transposable element families. P. Natl. Acad. Sci. USA, 2020; 117: 9451–9457.  
442 doi:10.1073/pnas.1921046117.

443 31 Heckenhauer J, Frandsen PB, Gupta DK et al. Annotated draft genomes of two caddisfly  
444 species *Plectrocnemia conspersa* CURTIS and *Hydropsyche tenuis* NAVAS (Insecta:  
445 Trichoptera). Genome Biol. Evol., 2019; 11: 3445–3451. doi:10.1093/gbe/evz264.

446 32 Kriventseva EV, Kuznetsov D, Tegenfeldt F et al. OrthoDB v10: sampling the diversity of  
447 animal, plant, fungal, protist, bacterial and viral genomes for evolutionary and functional  
448 annotations of orthologs. Nucleic Acids Res., 2019; 47: D807–811. doi:10.1093/nar/gky1053

449 33 Keilwagen J, Wenk M, Erickson JL et al. Using intron position conservation for homology-  
450 based gene prediction. Nucleic Acids Res., 2016; 44: e89. doi:10.1093/nar/gkw092.

---

- 451 34 Keilwagen J, Hartung F, Paulini M et al. Combining RNA-seq data and homology-based  
452 gene prediction for plants, animals and fungi. *BMC Bioinform.*, 2018; 19: 1–2.  
453 doi:10.1186/s12859-018-2203-5.
- 454 35 Olsen LK, Heckenhauer J, Sproul JS et al. Draft genome assemblies and annotations of  
455 *Agrypnia vestita* Walker, and *Hesperophylax magnus* Banks reveal substantial repetitive  
456 element expansion in tube case-making caddisflies (Insecta: Trichoptera). *Genome Biol.*  
457 *Evol.*, 2021; 13: evab013. doi:10.1093/gbe/evab013.
- 458 36 Campbell MS, Holt C, Moore B et al. Genome annotation and curation using MAKER and  
459 MAKER-P. *Curr. Protocols Bioinform.* 2014; 48: 4–11.  
460 doi:10.1002/0471250953.bi0411s48.
- 461 37 Conesa A, Götz S. Blast2GO: a comprehensive suite for functional analysis in plant  
462 genomics. *Int. J. Plant Genomics.* 2008; 2008: 1–12. doi:10.1155/2008/619832.
- 463 38 Ríos-Touma B, Holzenthal RW, Rázuri-Gonzales E et al. De novo genome assembly and  
464 annotation of an Andean caddisfly, *Atopsyche davidsoni* Sykora, 1991, a model for genome  
465 research of high-elevation adaptations. *Genome Biol. Evol.*, 2022; 14: evab286.  
466 doi:10.1093/gbe/evab286.

467

468

469

470

471

472

473 **Tables and Figures**

---

474 Table 1. Specimen accession and data type information.

	<b>BioProject</b>	<b>BioSample</b>	<b>Assembly</b>	<b>SRA</b>	<b>Sequence type</b>
<i>P. interpunctella</i>	PRJNA741212	SAMN20990134	NA	SRR15699974	Transcriptome
<i>P. interpunctella</i>	PRJNA741212	SAMN19857939	JAJAFS000000000	SRR15658214	Genome
<i>E. regina</i>	PRJNA741212	SAMN20522324	JAINEB000000000	SRR15651978	Genome

475

476 Table 2. Assembly genome stats for the species sampled in this study.

	<i>P. interpunctella</i>	<i>E. regina</i>	<i>P. interpunctella</i>	<i>P. interpunctella</i>
Reference	This study	This study	GCA_001368715.1	GCA_900182495.1
Platform	PacBio Sequel II	PacBio Sequel II	Illumina MiSeq/HiSeq	Illumina MiSeq/HiSeq
Coverage	44x	41x	100x	50x
Total ungapped length	300,731,903	917,780,411	364,621,386	364,623,808
Total gapped length	NA	NA	382,235,502	381,952,380
Number of scaffolds	NA	NA	7,743	10,542
Scaffold N50	NA	NA	5,094,612	1,270,674
Scaffold L50	NA	NA	23	75
Number of contigs	118	123	17,717	17,725
Contig N50	9,707,027	32,427,664	302,097	298,497
Contig L50	13	11	314	319
GC content	35.41%	32.68%	35.1%	35.1%
Shortest Contig	452	15,452	258	258
Longest Contig	13,555,736	57,864,696	2,314,344	2,314,344

Median Contig	161,724	36,760	1,714	1,719
Mean Contig	2,548,575	7,401,455	20,580	20,571

477

478

479

480 Table 3. Genome completeness by sample studied. Values shown are BUSCO scores for the

481 Endopterygota ODB10 data set.

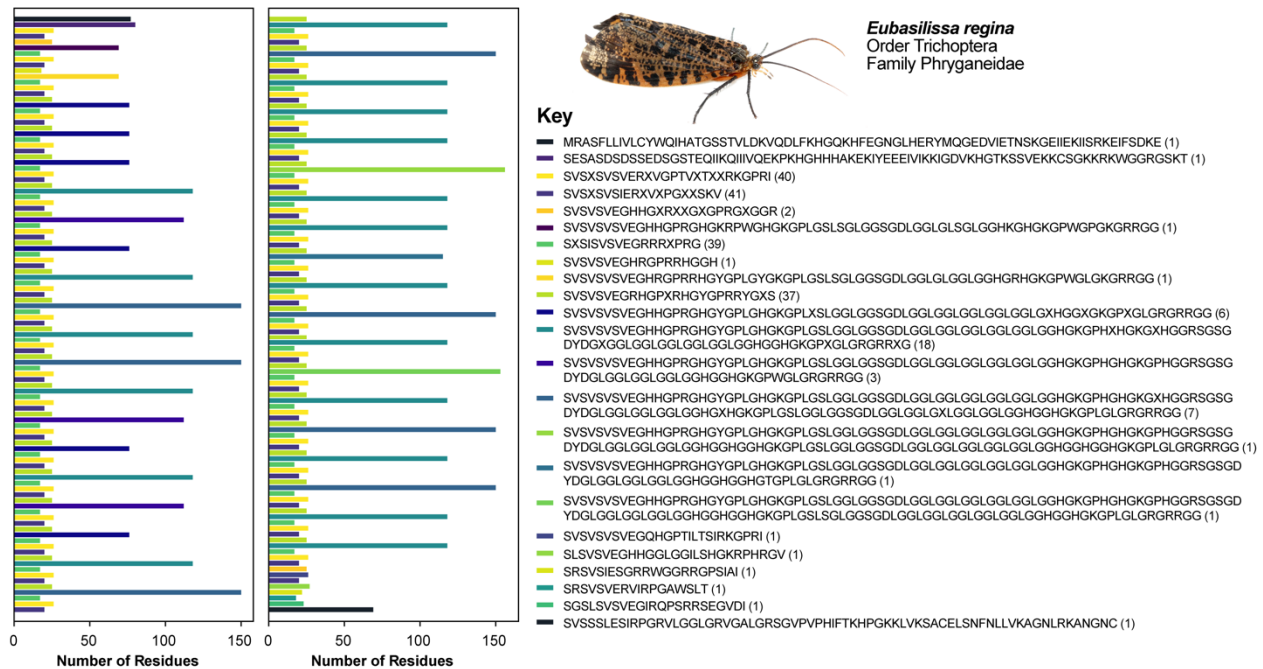
	<i>P. interpunctella</i>	<i>E. regina</i>	<i>P. interpunctella</i>	<i>P. interpunctella</i>
	This study	This study	GCA_001368715	GCA_900182495
Complete BUSCOs	2110	2021	2103	2105
Complete and single-copy	2097	2013	2074	2077
Complete and duplicated	13	14	29	28
Fragmented	5	63	10	8
Missing	9	34	11	11
Total groups searched	2124	2124	2124	2124
% complete	99.3	95.2	99.0	99.1

482





491 **Figure 2.** Schematic of the identity and ordering of repeat motifs in *P. interpunctella*. On the  
492 right panel are the repetitive units with the *N*-terminus at the beginning and the *C*-terminus at the  
493 end. The number in parenthesis refers to the number of times that particular motif is repeated  
494 across the gene. The color corresponds with the ordering of the repeats shown on the left. The  
495 gene is split into two panels, starting in the left panel and continuing in the right panel. “X”  
496 implies a site that is variable.



502  
503  
504  
505

506

507

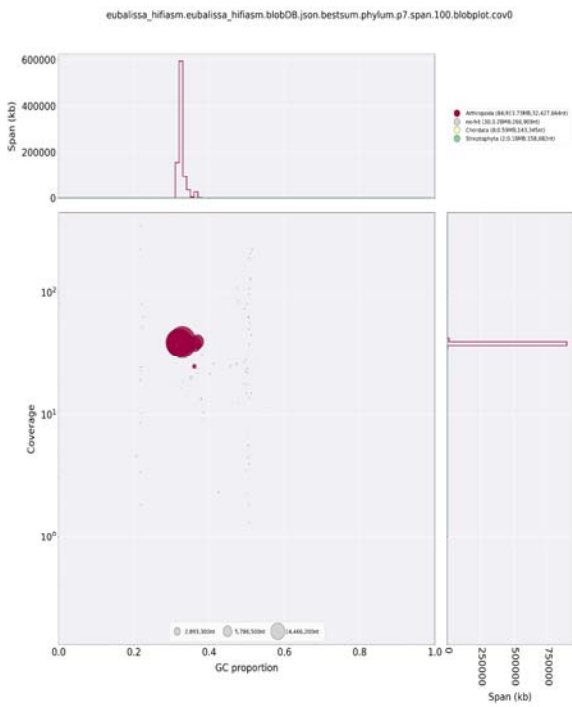
508

509

510

511

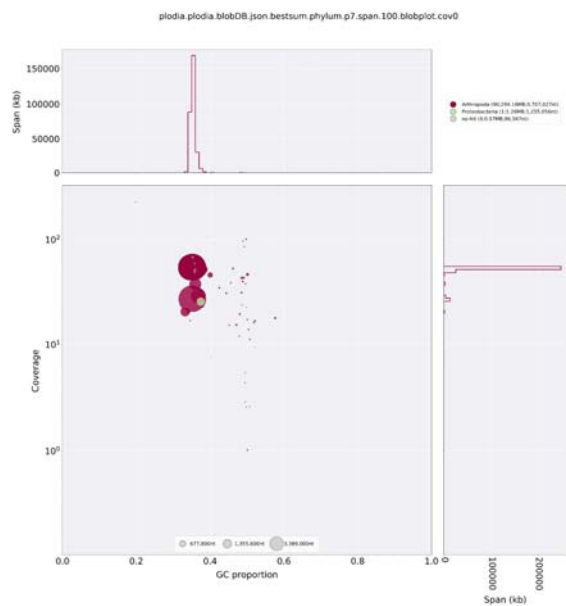
512 **Supplementary Files**



513

514 **Supplementary Figure 1 (in GigaDB). BlobPlot for *E. regina*.**

515



516

517 **Supplementary Figure 2 (in GigaDB).** BlobPlot for *P. interpunctella*.