

1 Hybridization dynamics and extensive introgression in the *Daphnia*  
2 *longispina* species complex: new insights from a high-quality *Daphnia*  
3 *galeata* reference genome

4

5 Jana Nickel<sup>1</sup>, Tilman Schell<sup>2</sup>, Tania Holtzem<sup>3</sup>, Anne Thielsch<sup>4</sup>, Stuart R. Dennis<sup>5</sup>, Birgit C. Schlick-  
6 Steiner<sup>3</sup>, Florian M. Steiner<sup>3</sup>, Markus Möst<sup>3</sup>, Markus Pfenninger<sup>2,6,7</sup>, Klaus Schwenk<sup>4</sup>, Mathilde  
7 Cordellier<sup>1\*</sup>

8 <sup>1</sup> Institute of Zoology, Universität Hamburg, Germany

9 <sup>2</sup> LOEWE Centre for Translational Biodiversity Genomics (LOEWE-TBG), Frankfurt am Main, Germany

10 <sup>3</sup> Department of Ecology, University of Innsbruck, Innsbruck, Austria

11 <sup>4</sup> Molecular Ecology, Institute for Environmental Sciences, University Koblenz-Landau, Landau in der  
12 Pfalz, Germany

13 <sup>5</sup> Dept of Aquatic Ecology, EAWAG Dübendorf, Switzerland

14 <sup>6</sup> Molecular Ecology, Senckenberg Biodiversity and Climate Research Centre, Frankfurt, Germany

15 <sup>7</sup> IoME, Gutenberg University, Mainz, Germany

16

17 \*corresponding author: Mathilde Cordellier, Department of Biology, Institute of Zoology, University of  
18 Hamburg, Hamburg, Germany, [mathilde.cordellier@uni-hamburg.de](mailto:mathilde.cordellier@uni-hamburg.de)

19

20

## 21 Abstract

22 Hybridization and introgression are recognized as an important source of variation that influence adaptive  
23 processes; both phenomena are frequent in the genus *Daphnia*, a keystone zooplankton taxon in freshwater  
24 ecosystems that comprises several species complexes. To investigate genome-wide consequences of  
25 introgression between species, we provide here the first high-quality genome assembly for a member of the  
26 *Daphnia longispina* species complex, *Daphnia galeata*. We further re-sequenced 49 whole genomes of three  
27 species of the complex and their interspecific hybrids both from genotypes sampled in the water column  
28 and from single resting eggs extracted from sediment cores. Populations from habitats with diverse  
29 ecological conditions offered an opportunity to study the dynamics of hybridization linked to ecological  
30 changes and revealed a high prevalence of hybrids. Using phylogenetic and population genomic approaches,  
31 we provide first insights into the intra- and interspecific genome-wide variability in this species complex  
32 and identify regions of high divergence. Finally, we assess the length of ancestry tracts in hybrids to  
33 characterize introgression patterns across the genome. Our analyses uncover a complex history of  
34 hybridization and introgression reflecting multiple generations of hybridization and backcrossing in the  
35 *Daphnia longispina* species complex. Overall, this study and the new resources presented here pave the way  
36 for a better understanding of ancient and contemporary gene flow in the species complex and facilitate future  
37 studies on resting egg banks accumulating in lake sediment.

## 38 Keywords

39 introgression, hybridization, resting eggs, species complex, whole-genome amplification, genome  
40 assembly

41

## 42 Introduction

43 Gene flow between species can be pervasive and can affect substantial parts of the genome. Hybridization  
44 and introgression are recognized as an important source of variation that can influence adaptive processes  
45 in plants, animals, yeast, and fungi (reviewed in Abbott, *et al.* 2013; Arnold and Martin 2009). The amount  
46 of realized gene flow varies among taxa and along the genome; it is governed by intrinsic genomic features  
47 such as recombination rate, structural variation and intrinsic incompatibilities, as well as the species' biology  
48 and ecology including ecological and sexual selection, migration, and mode of reproduction.

49 How can species in diversifying clades frequently hybridize and show introgression but nevertheless  
50 maintain species boundaries? A growing body of literature provides examples for a high variety of systems  
51 where speciation occurs in the face of gene flow (e.g. Fraïsse, *et al.* 2014; Martin, *et al.* 2019; Meier, *et al.*  
52 2017). However, it is important to recognize that these systems are distributed along a wide spectrum. On  
53 one side of this spectrum, hybridization occurs but is not followed by introgression for several reasons such  
54 as reduced hybrid fertility or strong selection against hybrid phenotypes, leading to rapid hybrid breakdown.  
55 Barth, *et al.* (2020) found that species boundaries in tropical eels are stable despite millions of years of  
56 hybridization, and also observed very few admixed individuals beyond F1 and first-generation backcrosses.  
57 The hybrid breakdown observed in this system reduces the likelihood of introgression via backcrossing. On  
58 the other side of the spectrum, hybridization is followed by introgression, and ongoing exchange of genetic  
59 information between species (e.g. Butlin, *et al.* 2014; Doellman, *et al.* 2018; Kaiser, *et al.* 2021; Martin, *et*  
60 *al.* 2013). Several empirical studies (Canestrelli, *et al.* 2017; Schreiber and Pfenninger) as well as theoretical  
61 models (Flaxman, *et al.* 2014; Rafajlović, *et al.* 2016; Yeaman and Whitlock 2011) suggest the possibility  
62 of intermediate constant equilibrium states, meaning that certain parts of the genome remain diverged  
63 ('islands' or 'continents of divergence'), while others are freely exchanged among closely related species  
64 without ever reaching complete genomic isolation.

65 Recurrent hybridization and introgression are frequent in the genus *Daphnia* (Crustacea, Cladocera)  
66 Members of the genus have served as ecological model organisms for over a century (e.g. Miner, *et al.*

67 2012), and the first crustacean genome to be sequenced was that of a member of the *Daphnia pulex* species  
68 complex (Colbourne, *et al.* 2011). Since then, the genomes of 45 crustaceans have been sequenced with a  
69 focus on species of economic or medical interest (NCBI, last accessed January 2021). Despite their key role  
70 in marine and freshwater food webs around the globe, genomic resources for zooplanktonic species are still  
71 scarce. In many aquatic food webs, zooplanktonic crustaceans link primary production by phytoplankton  
72 and secondary consumers, such as planktivorous fish and larger invertebrate species (Lampert and Sommer  
73 2007). (e.g. Gannon and Stemberger 1978; Gliwicz 1990)

74 *Daphnia* are highly phenotypically plastic and a textbook example for inducible defense mechanisms  
75 (Tollrian and Harvell 1999), as they respond to variation in predation risk through spectacular changes in  
76 morphology. Further, *Daphnia* are cyclical parthenogens and hence able to alternate between asexual and  
77 sexual reproduction. They reproduce asexually through longer periods of time, and the product of sexual  
78 reproduction events (usually seasonal) are resting eggs able to withstand adverse conditions for decades and  
79 even centuries (Frisch, *et al.* 2014). Resting eggs extracted from sediment cores can be hatched, and ancient  
80 genotypes brought to life (reviewed in Orsini, *et al.* 2013). Moreover, the DNA preserved in those resting  
81 eggs can be directly analyzed with various molecular methods (e.g. Cousyn, *et al.* 2001; Dziuba, *et al.* 2020;  
82 Lack, *et al.* 2018). Thus, cyclical parthenogenesis, biological archives in lake sediments and high levels of  
83 phenotypic plasticity make *Daphnia* a particularly interesting model for evolutionary studies.

84 The genus *Daphnia* is composed of two subgenera, *Ctenodaphnia* and *Daphnia*, and two groups are  
85 delimited within the subgenus *Daphnia*: the *D. pulex* group *sensu lato* and the *D. longispina* group *sensu*  
86 *lato* (see Adamowicz, *et al.* 2009). The latter is sometimes also referred to as subgenus *Hyalodaphnia* and  
87 includes the *Daphnia longispina* species complex (DLSC) (Petrusek, *et al.* 2008a). The two *Daphnia* groups  
88 are highly differentiated and share their most recent common ancestor around 30 Mya (MRCA *D. longispina*  
89 – *D. pulex* group, MRCA *D. longispina* – *D. pulex* group, Cornetti, *et al.* 2019). Members of the genus  
90 *Daphnia* show little variation in chromosome number, with most species having 10 pairs of chromosomes,  
91 except for the *D. pulex* group with n=12 (Beaton and Hebert 1994; Trentini 1980). All sequenced and

92 assembled *Daphnia* genomes so far belong either to the *D. pulex* group or the subgenus *Ctenodaphnia*,  
93 however no high-quality reference genome of the third major group, the *D. longispina* group  
94 (*Hyalodaphnia*) is published.

95 The prevalence of hybridization in the genus *Daphnia* across taxa and ecosystems and its impact on their  
96 evolutionary history has intrigued researchers for decades (e.g. Schwenk 1993; Vergilino, *et al.* 2011; Wolf  
97 1987). In contrast to many other well-studied hybrid systems (Barton and Hewitt 1985) with clear defined  
98 hybrid zones where species' ranges overlap, the distribution of *Daphnia* species and their hybrids is more  
99 of a fragmented nature: they occupy lake and pond ecosystems that vary in their ecological characteristics  
100 and hence constitute a mosaic across the landscape. Ecologically differentiated taxa and their hybrids are  
101 thus distributed across habitat patches (Harrison 1986). Within these patches, the possibility to interrogate  
102 biological archives also revealed fluctuations in *Daphnia* community composition over time (e.g. Alric, *et*  
103 *al.* 2016; Brede, *et al.* 2009), associated with hybridization events among species in some cases. Variation  
104 in hybridization events across time and among habitats has often been observed in correlation with  
105 ecological changes, such as eutrophication or global change (Brede, *et al.* 2009; Cordellier, *et al.* 2021;  
106 Dziuba, *et al.* 2020; Keller, *et al.* 2008; Rellstab, *et al.* 2011; Spaak, *et al.* 2012).

107 Members of the *Daphnia longispina* species complex inhabit many large ponds and lakes in central and  
108 northern Europe, and three of them have been particularly well studied: *Daphnia galeata*, *Daphnia*  
109 *longispina* and *Daphnia cucullata* (Petrusek, *et al.* 2008a). These species can coexist, but earlier studies  
110 suggest gene flow among them is limited (Spaak 2004). Despite their obviously ancient divergence  
111 (Schwenk, *et al.* 2000), DLSC species are still able to form interspecific hybrids, although not all  
112 combinations are equally likely to lead to viable and fertile individuals (Schwenk, *et al.* 2001). A mechanism  
113 preventing gene flow among species might be their different ecological preferences, e.g., regarding trophic  
114 level (Spaak, *et al.* 2012), food quality (Seidendorf, *et al.* 2007), and predation pressure (Spaak and Hoekstra  
115 1997);(Petrusek, *et al.* 2008b).

116 Up to now, genetic markers available to study hybridization in the DLSC are limited to allozymes (Wolf  
117 and Mort 1986), a few mitochondrial regions (Schwenk 1993), a dozen microsatellite markers (Brede, *et al.*  
118 2006; Thielsch, *et al.* 2012) and a few further nuclear loci (Billiones, *et al.* 2004; Rusek, *et al.* 2015; Skage,  
119 *et al.* 2007). Seminal studies such as Brede, *et al.* (2009) and Limburg and Weider (2002) first made use of  
120 microsatellite markers to analyze environmentally driven shifts in allelic frequencies, species and hybrid  
121 composition of the DLSC communities in Lake Constance and Belauer See over time, respectively. Further,  
122 a number of studies addressed the spatial distribution of DLSC species/taxa with these markers (e.g. Griebel,  
123 *et al.* 2016; Ma, *et al.* 2019; Thielsch, *et al.* 2017). These low-resolution markers allowed to identify hybrid  
124 individuals and brought evidence for introgression but could not provide the resolution necessary to either  
125 assess how pervasive introgression is or how it varies across the genome. Further, it is not clear whether  
126 introgression occurs among all three species to the same extent. Given the ubiquitous hybridization among  
127 the DLSC taxa, the question also arises why they are still well distinguishable species. Whether the DLSC  
128 represents a case of incipient speciation, introgression after secondary contact, speciation reversal, or has  
129 reached an intermediate constant equilibrium state, among other possibilities, can only be answered with  
130 genome-wide analyses empowered by a high-quality genome assembly.

131 Here, we present a high-quality assembly for *Daphnia galeata*, thus filling an important gap for *Daphnia*  
132 whole-genome studies. Furthermore, to facilitate genome-wide assessments of divergence across species  
133 and of introgression between species, we conducted genome-wide resequencing studies in the DLSC. We  
134 analyzed whole-genome sequences of parental species and their interspecific hybrids, both from genotypes  
135 obtained in the wild and maintained in laboratories, and from single resting eggs extracted from sediment  
136 cores. We provide first insights into the intra- and interspecific genome-wide variability in this species  
137 complex and identify regions of high divergence. We reconstructed the phylogenetic relationships in the  
138 species complex using whole mitochondrial genomes. Finally, we assess the length of ancestry tracts in  
139 different classes of hybrids to characterize introgression patterns. Our study paves the way for long-awaited  
140 analyses on the dynamics of introgression in this complex and exploitation of the unique opportunity this  
141 group has to offer: a window of more than one hundred years of evolution in action.

## 142 Results

### 143 Genome assembly

144 The raw assembly was obtained by combining PacBio long reads (1,679,290, 11.52 Gb) and Illumina short  
145 reads (70,310,338, 9.79 Gb after trimming) and using the hybrid assembler RA ([https://github.com/lbcb-](https://github.com/lbcb-sci/ra)  
146 [sci/ra](https://github.com/lbcb-sci/ra)). It originally comprised 1,415 contig sequences covering a total length of 153.6 Megabases (Mb),  
147 with an N50 value of 172 kilobases (kb) and a slightly elevated GC content (40.02% Supplementary  
148 Methods Table 3) compared with the values expected for a *Daphnia* species (see Table 1). According to an  
149 analysis based on coverage and GC content of the contig sequences conducted with blobtools (Laetsch and  
150 Blaxter 2017), a portion of the assembly consisted of non-*Daphnia* contigs, which could then be removed  
151 (267 contigs, equaling 22.97 Mb). Consequently, GC content decreased to 38.75%, nearing the values  
152 obtained for other *Daphnia* assemblies (see Table 1 for an overview). The application of this filter as well  
153 as the exclusion of the mitochondrial genome led to a decrease in the number of sequences and the total  
154 length of the assembly. Iterative scaffolding led to a decrease in the total number of sequences. This together  
155 with a substantial increase in N50 resulted in a highly contiguous assembly, with a total length of 133,304,63  
156 basepairs (bp), an N50 of 756.7 kb and only 346 sequences, i.e., on average 30 sequences per chromosome.  
157 Contiguity statistics for the different assembling steps are given in Supplementary Methods Table 3.

158 Mapping the filtered Illumina reads with bwa mem (Li 2013) and PacBio reads with Minimap 2.17 (Li  
159 2018), resulted in a mapping rate of respectively 94.1% and 85.5%. The coverage distribution can be seen  
160 in Figure S1B.

161 According to blobtools results, no contamination could be identified in the final assembly (Figure S1A).  
162 Remaining scaffolds (12, amounting to a total length 1.79Mb) with taxonomic assignment other than  
163 Arthropoda were kept because coverage and GC are similar to *D. galeata* scaffolds and taxonomic  
164 assignment alone might be false positive. Further, the completeness assessment through BUSCO (Simão, *et*  
165 *al.* 2015, Arthropoda set, odb9) indicated 95.7% of complete single copy core orthologs and a very low  
166 duplication rate (C: 95.7% [S: 94.7%, D: 1.0%], F: 0.8%, M: 3.5%, n: 1066). The genome size was estimated

167 based on mapped nucleotides and mode of the coverage distribution by backmap 0.3  
168 (<https://github.com/schell/Backmap>), resulting in 156.86Mb and 178.03Mb for Illumina (52x) and PacBio  
169 (26x) respectively, and by k-mer based approach using GenomeScope resulting in a size of 150.6Mb.

170 When compared to other published full genomes for *Daphnia* species, the *D. galeata* final assembly is  
171 shorter than both *D. pulex* assemblies (Colbourne, *et al.* 2011; Ye, *et al.* 2017), and roughly the same size  
172 as *D. magna* (Lee, *et al.* 2019), which also has 10 chromosomes (Table 1). The GC content is lower, which  
173 can be attributed to the strict filtering for contamination applied pre- and post-assembly, a procedure not  
174 applied in the other *Daphnia* assemblies, to our knowledge. Even though Lee *et al.* (2019) and Ye *et al.*  
175 (2017) treated the animals with antibiotics before sequencing this suggests that these genome assemblies  
176 contain more contigs of bacterial origin than the *D. galeata* assembly. Thanks to the use of long-read data,  
177 iterative scaffolding and gap filling, the number and length of assembly gaps (Ns) is substantially lower and  
178 contiguity is high (but see Table 3 in Supplementary Methods).

## 179 Genome annotation

180 After applying RepeatMasker (Smit, *et al.* 2013-2015) with the custom repeat library described in the  
181 methods section, 21.9% of the assembly was masked. The distribution of masked fraction per repeat element  
182 can be found in Supplementary Methods Table 5.

183 The final annotation with MAKER (Holt and Yandell 2011) predicts 15,845 genes with a median length of  
184 2,097 base pairs. There is an average of 1.06 mRNAs per gene and 7 exons/mRNA (Table 1). The total  
185 number of predicted mRNA substantially differs from the number of transcripts previously published for  
186 this species (32,903, Huylmans, *et al.* 2016). This is not surprising, as this transcriptome assembly did not  
187 make use of protein evidence we included here, and might contain isoforms. Further, it was based on a pool  
188 of mRNA from different clonal lines, and the assembly process might have been impeded by allelic  
189 diversity. As further quality criterion, the Annotation Editing Distance (AED) was compared across the  
190 three MAKER rounds and is visualized in Supplementary Methods Figure 4. AED improved mostly  
191 between rounds 1 and 2 of the annotation but only marginally with a further round.



192 **Table 1:** Assembly metrics and annotation statistics for the present assembly and two previously published *Daphnia*  
 193 assemblies. Contiguity statistics of the annotation were calculated excluding tRNAscan results. BUSCO 3.0.2 was  
 194 executed in protein mode for the different MAKER rounds. Conserved Domain Arrangements (CDAs) were searched  
 195 with Pfam scan 1.6 and DOGMA 3.4. Results for BUSCO and DOGMA completeness statistics are given in percent.

Species	<i>D. galeata</i>	<i>D. pulex</i> (Ye <i>et al</i> )	<i>D. magna</i> (Lee <i>et al</i> )
Strain	M5	PA42	SK
<b>Assembly metrics</b>			
# scaffolds	346	493	4,192
Largest scaffold (bp)	2,950,711	7,584,612	16,359,456
Total length (bp)	133,304,630	189,550,516	122,937,721
N50 (bp)	756,671	1,160,003	10,124,675
L50 (bp)	48	36	5
GC (%)	38.75	40.39	40.54
# N's	120,845	4,006,006	82,97,703
# N's per 100 kbp	90.65	2113.42	6749.52
<b>Annotation</b>			
<b>Number</b>			
Gene	15,845	18,440	15,721
mRNA	16,774	18,440	15,721
Exon	117,364	128,688	95,203
CDS	119,402	118,916	94,047
<b>Mean</b>			
mRNAs/gene	1.06	1	1
Exons/mRNA	7.00	6.98	6.06
CDSs/mRNA	7.12	6.45	5.98
<b>Median length (bp)</b>			
Gene	2,097	1,919.5	1,586
mRNA	2,142	1,919.5	1,521
Exon	167	162	160
Intron	74		
CDS	152	144	159
<b>Total space (bp)</b>			
Gene	51,689,473	53,936,938	37,505,261
mRNA	51,689,329	53,936,938	36,178,687
Exon	29,314,592	30,208,483	22,336,755
CDS	25,132,876	23,586,918	21,881,778
<b>Single</b>			
Exon mRNA	663	144	1,775
CDS mRNA	710	0	0
<b>BUSCO N=1066</b>			
C	94.3	94.1	97.0
S	91.7	82.6	95.3

D	2.6	11.5	1.7
F	0.7	3.5	1.7
M	5.0	2.4	1.3

DOGMA N=4222 93.63 91.43 93.91

196  
197 A high percentage of protein sequences could be annotated: 15,898 (94.78%) with InterProScan (Jones, *et*  
198 *al.* 2014) and 15,960 (95.15%) with blast against Swiss-Prot. With this combination of searches, a hit within  
199 InterProScan and blast was found for 16,675 protein sequences (99.41%). GeneOntology annotation was  
200 possible for 9555 sequences (56.96%). A detailed overview of the functional annotated sequences per  
201 database or search algorithm is shown in Supplementary Methods Table 7.

## 202 Genotyping

203 Short-read sequence data were generated for 72 individuals: 17 unamplified DNA samples from isofemale  
204 clonal lines and 55 whole genome amplification (WGA) samples (conducted on single resting eggs) that  
205 passed PCR contamination checks. After screening for contamination and removing datasets with only very  
206 few reads mapping to the *D. galeata* reference, 49 single genotypes remained: 32 from resting eggs and 17  
207 from clonally propagated lines, established from individuals sampled in the water column or hatched from  
208 resting eggs (Figure 1A, Table S2). Data gained from clonal lines with a species attribution were used as  
209 “parental species” data: five samples for *D. galeata*, four for *D. longispina*, and three for *D. cucullata*. The  
210 parental clones are part of two larger clone panels representing the parental species and their diversity in  
211 several European lakes. Their identity was established prior to this study either based on mitochondrial and  
212 microsatellite markers (M5, LC3\_6, J2, Herrmann, *et al.* 2017) or morphological examination,  
213 mitochondrial markers and factorial correspondence analyses based on microsatellite markers (Alric, *et al.*  
214 2016; Möst 2013), including hybrids and historical resting eggs, which separates parental species and  
215 hybrids (e.g. Alric, *et al.* 2016; Dlouha, *et al.* 2010; Rellstab, *et al.* 2011; Yin, *et al.* 2014). In addition, data  
216 were available for four resting eggs from Arendsee (AR), 12 resting eggs from Dobersdorfer See (DOB),  
217 five clonal lines and eight resting eggs from Eichbaumsee (EIC), and eight resting eggs from Selenter See  
218 (SE) (Table S2). While the analysis of eggs from older sediment layers was attempted, biological material

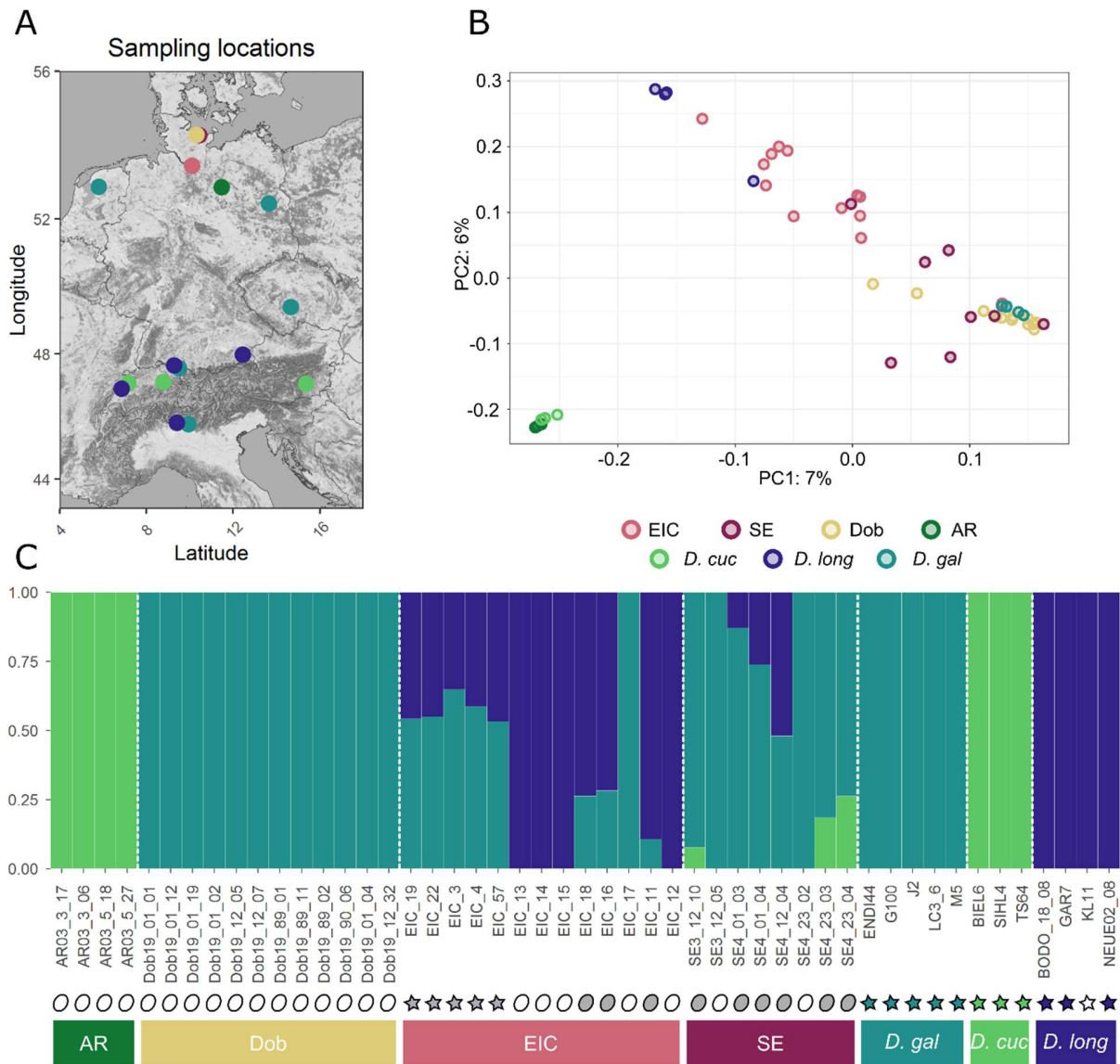
219 was either limited, of poor quality, or contaminated. Our isotope dating for DOB was inconclusive: either  
220 slides of the cored location or a high sedimentation rate meant the top 30 cm of the core didn't show the  
221 usual isotope peaks, thus preventing precise dating. EIC samples were recent since they were collected from  
222 surface bank sand. For SE, the oldest eggs analyzed here originated from the 2-3cm layer of the core, which  
223 corresponds to max. ~17 years (pers. comm Thorbjørn Andersen). For AR the oldest eggs for which results  
224 were obtained were isolated from the 4-5cm layer, corresponding to ~2005 (pers. comm. Miklos Balint).

225 An average of 89.9% (range: 31.7-98.6%) reads aligned to the reference genome with a mean coverage of  
226 10.26x (range: 0.34-52.30x) (Table S3). The final SNP data set for subsequent analyses after quality-  
227 filtering included 3,240,339 SNPs across the 49 samples. To rule out possible reference bias we compared  
228 mapping rates of reads with the reference allele and to the alternative allele at heterozygous sites. We found  
229 no preferential mapping of the reference allele, as all species categories and the hybrids had a median  
230 distribution close to 0.5 (Figure S2).

### 231 **Principal Component Analysis**

232 In a PCA including all genotypes and conducted with SNPRelate v1.20.1 (Zheng, *et al.* 2012), the parental  
233 species genotypes grouped in three very distinct clusters. *D. cucullata* separated from *D. galeata* along PC1,  
234 which explained 7% of the variation. *D. longispina* separated from *D. galeata* and *D. cucullata* along PC2  
235 which explained 6% of the variation (Figure 1B). Although sampled in different lakes, all parental species  
236 genotypes were grouped in tight clusters along the two axes with little evidence for population substructure.  
237 Population AR clustered with the *D. cucullata* reference individuals while population samples from DOB,  
238 EIC and SE were more spread out, mostly between the *D. galeata* and *D. longispina* clusters.

239



240

241 **Figure 1:** Parental species: *D. gal*: *D. galeata*, *D. long*: *D. longispina*, *D. cuc*: *D. cucullata*, populations: AR: Arendsee,  
 242 Dob: Dobersdorfer See, SE: Selenter See, EIC: Eichbaumsee. Color coding is consistent throughout panels A and B  
 243 **A.** Map of the sampling locations. **B.** PCA plot obtained with SNPrelate, including loci with linkage  $r^2 < 0.5$  within  
 244 500-kb sliding windows. **C.** Admixture plot obtained with K=3. Symbols indicate the sample type: oval for genotypes  
 245 sequenced directly from resting eggs, stars for genotypes sampled in the water column and propagated clonally in the  
 246 laboratory prior to sequencing. Symbol filling indicates how these genotypes were classified in subsequent analyses:  
 247 white for non-admixed genotypes, grey for admixed genotypes, green, blue and teal for genotypes used as  
 248 representatives for parental species. Bottom bars are color coded to match the color scheme used in panels A and B.

## 249 Admixture analyses uncover hybrids

250 The PCA results are confirmed by an admixture analysis conducted with ADMIXTURE (Alexander and  
 251 Lange 2011) with K=3, supported by the lowest cross-validation error of the tested K values. The known

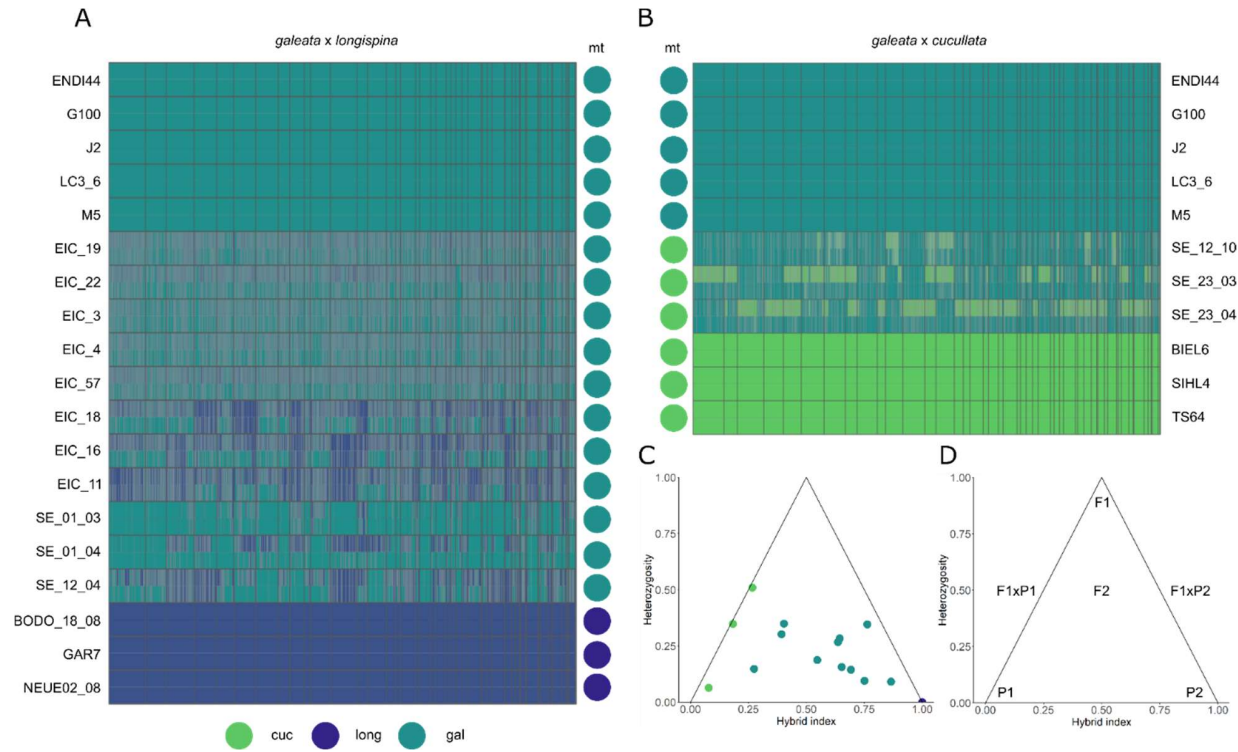
252 parental species genotypes were clearly separated into three clusters (Figure 1C). While we detect no  
253 evidence of admixture in the AR and DOB samples and, based on our parental species, consider them to  
254 belong to the species *D. cucullata* and *D. galeata*, respectively, the two other populations seem to consist  
255 mostly of admixed individuals. The five EIC samples sequenced after clonal propagation were all found to  
256 be admixed (*D. galeata* and *D. longispina*), while the EIC resting eggs were either admixed (3) or belonged  
257 to one of the parental species (5). SE resting eggs present all combinations of admixture except *D. cucullata*  
258  $\times$  *D. longispina*: *D. galeata* (2), admixed between *D. galeata* and *D. cucullata* (3), and admixed between *D.*  
259 *galeata* and *D. longispina* (3).

## 260 Ancestry painting

261 Based on results obtained in the ADMIXTURE analysis, two pairs of species and their putative hybrids  
262 were analyzed with an “ancestry painting” approach, outlined in Barth, *et al.* (2020) and Runemark, *et al.*  
263 (2018a): *D. galeata* and *D. longispina* parental genotypes and putative hybrids between them from  
264 populations EIC and SE, and *D. galeata* and *D. cucullata* parental genotypes and putative hybrids between  
265 them from population SE. Briefly, after identifying fixed sites for each of the species in the analyzed pair,  
266 heterozygosity was calculated for these sites and a hybrid index derived from the obtained results  
267 ([https://github.com/mmatschiner/tutorials/tree/master/analysis\\_of\\_introgression\\_with\\_snp\\_data](https://github.com/mmatschiner/tutorials/tree/master/analysis_of_introgression_with_snp_data)). Further,  
268 information on the maternal species is used to tentatively categorize the admixed individuals. For a first-  
269 generation hybrid (F1) the expectation would be 50% of the nuclear genome being derived from each  
270 parental species (hybrid index  $\approx$  0.5) and mostly heterozygous fixed sites (heterozygosity  $\approx$  1.0). Individuals  
271 originating from the backcrossing of F1 with one of the parental species are expected to have hybrid index  
272 values around 0.25 or 0.75. (Figure 2D). We consider individuals with intermediate hybrid indices ( $>0.25$   
273 and  $<0.75$ ) and lower heterozygosity ( $<0.5$ ) to be later-generation hybrids, meaning they have one or  
274 multiple hybrid ancestors we are not able to classify further (Slager, *et al.* 2020) We consider individuals  
275 with a hybrid index of  $\leq 0.25$  or  $\geq 0.75$  to be backcrossed with the respective parental species in at least one  
276 generation and the majority of the genome derives from one species. This definition is broad and will be  
277 refined with the addition of a greater number of parental genotypes.

278 The comparison of genotypes from parental species *D. galeata* and *D. longispina* (five and three individuals,  
279 respectively) allowed identifying a total of 335,052 fixed sites between the two species. Due to the quality  
280 filters applied to parental and hybrid genotypes, we could analyze 131,914 of these fixed sites in the putative  
281 hybrids, where read coverage was sufficient. The diploid genotypes were then plotted for all hybrids as  
282 homozygous for either of the parental species or as heterozygous (Figure 2A for the 50 longest scaffolds).  
283 The *D. longispina* reference clone KL11 was excluded from further analysis due to issues with missing data.  
284 All eleven genotypes from SE and EIC identified as likely *D. galeata* x *D. longispina* hybrids in the  
285 ADMIXTURE analysis possessed a *D. galeata* mitochondrial genome. The proportion of the maternal *D.*  
286 *galeata* genome in these hybrids, however, varied greatly, between 27.4% and 86.6%, and they all showed  
287 low heterozygosity, between 9.1% and 34.6% (Figure 2C, Table 2). These values are unlikely for F1 hybrids  
288 or backcrosses of F1 with one of the parental species (Table 2).

289 Comparing genotypes from the parental species *D. galeata* and *D. cucullata* (five and three individuals,  
290 respectively) led to identifying 715,438 fixed sites between the two species (due to quality filtering, 275,216  
291 of these sites were further analyzed). All three *D. galeata* x *D. cucullata* hybrids carried a *D. cucullata*  
292 mitochondrial genome, their hybrid index varied between 0.079 and 0.267 and their heterozygosity ranged  
293 from 6.4% to 50.9 % (Figure 2B&C, Table 2). The individual SE\_23\_04 is most likely the result of a  
294 backcrossing with *D. galeata*; however, it is difficult to determine what backcrossed with it: either an F1  
295 hybrid or a later generation hybrid i.e., that resulted from several generations of admixture. Haplotype  
296 information would be needed to gain certainty. The other two hybrids' lower hybrid index hints  
297 backcrossing with *D. galeata*, according to the criteria defined above.



298

299  
 300  
 301  
 302  
 303  
 304

**Figure 3: Panels A & B** Ancestry painting of the hybrid individuals identified through the admixture analysis. Each row represents an individual. Colored circles on the side indicate the mitochondrial identity of the individuals, based on the analysis of full mitochondrial genomes. Scaffolds are sorted by length and separated by thin grey lines. In panels **A** and **B**, the five upper rows represent individuals assigned to the parental species *D. galeata*. In **A**, the last three rows correspond to individuals assigned to the parental species *D. longispina*. In **B**, the last three rows correspond to individuals assigned to the parental species *D. cucullata*. Triangle plots summarizing **C**, the hybrid index and mitochondrial species identity for all individuals identified as admixed **D**, the hypothetical expected means of parental species (P1 and P2) and hybrid classes (F1xP1 and F1xP2: backcrosses with parental species P1 and P2, respectively).

305 **Table 2:** Data derived from ancestry painting analysis and based on the fixed sites inferred from analyzing parental  
306 species genotypes. Maternal species attribution is based on mitochondrial phylogeny, hybrid attribution is based on  
307 the ADMIXTURE plot.

Sample	Hybrid Index	Heterozygosity	Maternal species	Hybrid	Interpretation
SE_12_10	0.079	0.064	cuc	gal x cuc	Backcross gal
SE_23_03	0.183	0.348	cuc	gal x cuc	Backcross gal
SE_23_04	0.267	0.509	cuc	gal x cuc	Unclear
EIC_19	0.653	0.156	gal	gal x long	Later-generation
EIC_22	0.644	0.283	gal	gal x long	Later-generation
EIC_3	0.751	0.095	gal	gal x long	Backcross gal
EIC_4	0.693	0.144	gal	gal x long	Later-generation
EIC_57	0.637	0.267	gal	gal x long	Later-generation
EIC_18	0.393	0.302	gal	gal x long	Later-generation
EIC_16	0.403	0.349	gal	gal x long	Later-generation
EIC_11	0.274	0.148	gal	gal x long	Later-generation
SE_01_03	0.866	0.091	gal	gal x long	Backcross gal
SE_01_04	0.763	0.346	gal	gal x long	Backcross gal
SE_12_04	0.547	0.187	gal	gal x long	Later-generation

### 308 Population genomics parameters

309 To calculate genome-wide nucleotide diversity ( $\pi$ ), between-taxon differentiation ( $F_{ST}$ ), and between-taxon  
310 divergence ( $d_{xy}$ ) within 100-kb sliding windows, we took advantage of the inference made with  
311 ADMIXTURE and pooled all genotypes which were unambiguously assigned to either of the parental  
312 species clusters. Consequently, a total of seven genotypes from four populations were classified as *D.*  
313 *cucullata*, eight from five populations as *D. longispina*, and 20 genotypes from eight populations as *D.*  
314 *galeata* (Table S2). All values ( $d_{xy}$ ,  $\pi$  and  $F_{ST}$ ) were calculated with the script popgenWindows.py  
315 ([github.com/simonhmartin/genomics\\_general](https://github.com/simonhmartin/genomics_general) release 0.3) for each species pair and are plotted for the 50  
316 largest scaffolds in Figure 3.

317 The window-based  $F_{ST}$  values for all three possible pairs among the three species averaged 0.274 for *D.*  
318 *galeata* vs *D. longispina*, 0.343 for *D. cucullata* vs *D. longispina* and 0.364 for *D. galeata* vs *D. cucullata*.  
319 The mean sequence divergence  $d_{xy}$  for the three pairs was 0.018 for *D. galeata* vs *D. longispina* and 0.022  
320 for both *D. cucullata* vs *D. longispina* and *D. cucullata* vs *D. galeata*. Both parameters show similar



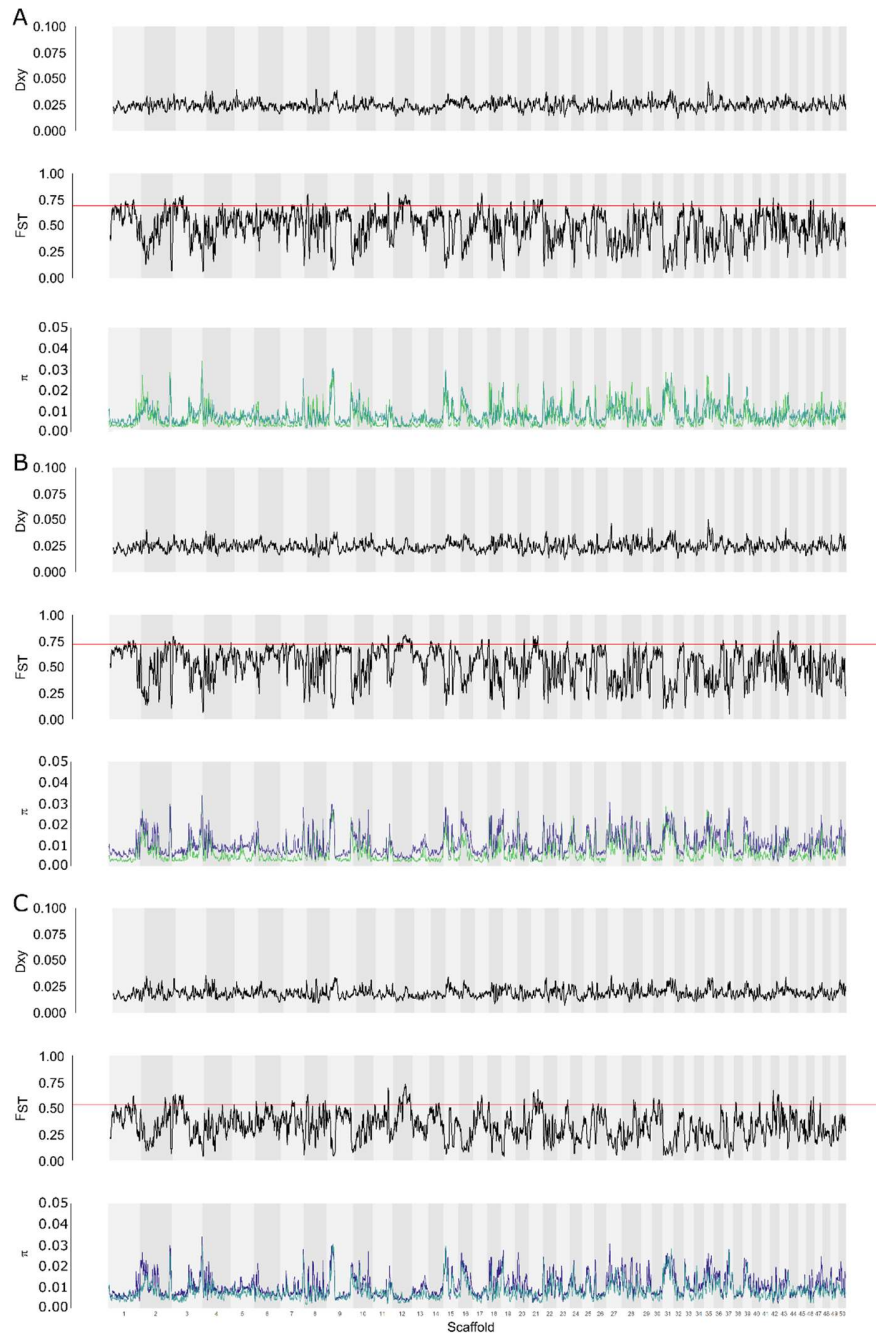
321 patterns, with lower values on average when comparing *D. galeata* to *D. longispina* than when comparing  
322 *cucullata* to either one of the other species. These patterns confirm the results obtained with other analyses,  
323 for example, the higher number of fixed sites between *D. galeata* and *D. cucullata* in the ancestry painting  
324 analysis.

325 The window-based estimates show high variability in levels of differentiation and divergence along the  
326 genome. Further, regions of high or low differentiation are mostly associated with depleted or high  
327 nucleotide diversity, respectively (see scaffolds 2 and 9 for example). However, the genome being  
328 represented by unordered scaffolds instead of chromosomes makes this difficult to interpret further.

329 Nucleotide diversity ( $\pi$ ) to quantify the level of genetic variation within each taxon was on average higher  
330 for *D. longispina* (1.18%) than for the other two species (0.95% and 0.85% for *D. galeata* and *D. cucullata*,  
331 respectively). This cannot be explained by the differences in group sample sizes, since *D. galeata* was the  
332 group with the largest sample size (and highest number of sampled populations). To ensure our window-  
333 based estimates were not biased because of the overrepresentation of some populations in a group (e.g. DOB  
334 in the *galeata* group), we also calculated these indices using only one individual from each population per  
335 species; if one population contained multiple individuals, we picked one individual at random to represent  
336 this population (see Table S2 for a listing of the used genotypes - results shown in Figure S4).

337 Many more highly differentiated windows and genes were shared among two or all species pairs than would  
338 be expected by a random intersection (Figure S5). For example, a total of 2575 10kb windows had an  $F_{ST}$   
339 value within the 95<sup>th</sup> percentile in the pair *D. galeata*/*D. longispina* and 2569 in the pair *D. galeata*/*D.*  
340 *cucullata*. The mean expected number of windows in common between these two pairwise comparisons  
341 was 113, but the number of windows in common observed in the data was 1601. A similar pattern was  
342 observed in all other intersections. This result suggests that the location of differentiated genome parts is  
343 not due to random processes but has biological significance. A GO-enrichment analysis of these isolated  
344 genes to shed light on the function of these species-specific genes, however, was not possible, because of  
345 the low number of genes with GO annotation. For the pair *D. galeata*/*cucullata*, only 12% of the genes in

346 the outlier windows were annotated with Gene Ontologies, for the pair *D. galeata/longispina* it was 11%  
347 and for the *D. cucullata/longispina* pair it was 10%.



348  
349 **Figure 4:** Window-based statistics for the pairs **A.** *D. galeata* / *D. cucullata*, **B.** *D. cucullata* / *D. longispina* and **C.** *D.*  
350 *galeata* / *D. longispina*, shown for the 50 largest scaffolds in 100kb windows with 10kb step size – calculations are for  
351 all non-admixed individuals unambiguously assigned to parental species according to the ADMIXTURE analysis. **In**  
352 **each panel from top to bottom:**  $d_{xy}$  values, pairwise  $F_{ST}$  values with a red horizontal line indicating the 95<sup>th</sup> percentile,  
353 nucleotide diversity ( $\pi$ ) for *D. galeata* (teal), *D. longispina* (dark blue), and *D. cucullata* (lime green).

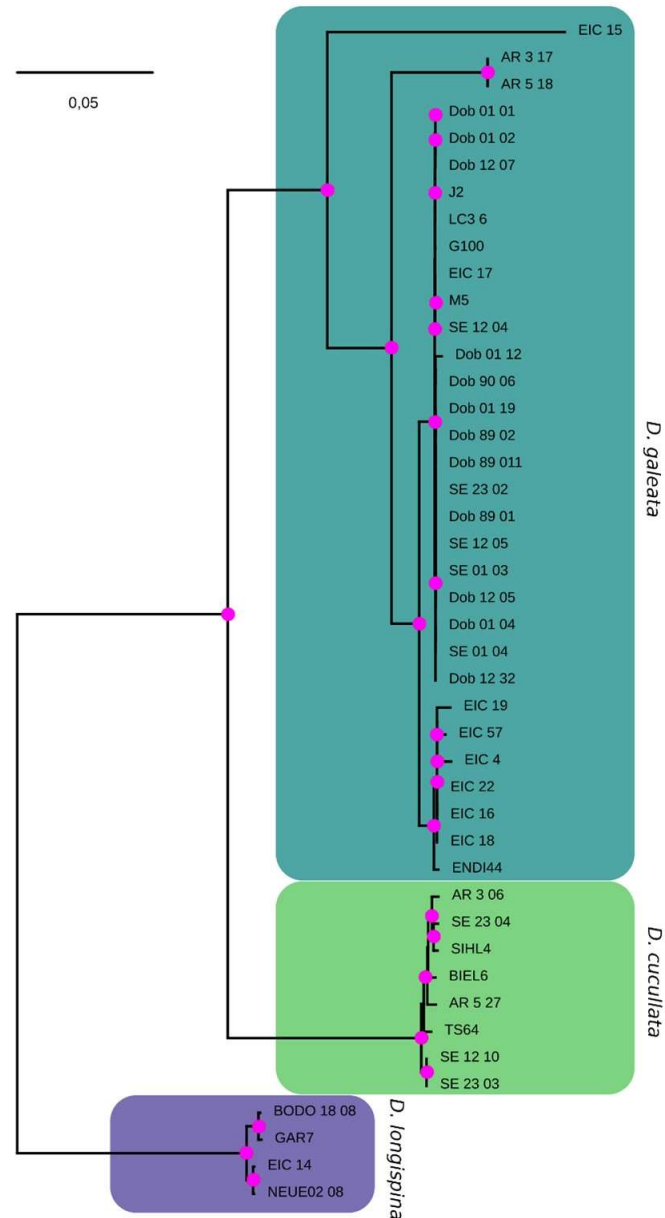
## 354 Phylogeny based on complete mitochondrial genomes

355 Phylogenetic reconstruction based on the mitochondrial protein-coding and ribosomal RNA genes were  
356 largely consistent with earlier mitochondrial phylogenies based on single or few mitochondrial genes (e.g.  
357 Adamowicz, *et al.* 2009; Petrusek, *et al.* 2012). We identified highly supported clades comprising the  
358 respective parental genotypes, hence representing *D. longispina*, *D. cucullata*, and *D. galeata* mitochondrial  
359 haplotypes (Figure 4). *D. cucullata* and *D. galeata* mitochondrial haplotypes clustered as sister groups.  
360 While the mitochondrial haplotypes in the *D. longispina* and *D. cucullata* clusters do not show much  
361 divergence, the *D. galeata* haplotype cluster also contains deeper branching events (haplotype EIC\_15 and  
362 AR3\_17 / AR5\_18). Further, although all samples from AR were unequivocally categorized as *D. cucullata*  
363 in the ADMIXTURE analysis and clustered with *D. cucullata* parental genotypes in the PCA, two of them  
364 have mitochondrial haplotypes falling into the *D. galeata* cluster (AR3\_17 and AR5\_18). A similar  
365 mismatch was also observed for EIC\_15, which clusters with *D. longispina* when considering nuclear SNP  
366 and with *D. galeata* when considering the mitochondrial genome. Within the species clusters, we observed  
367 a grouping by lake with many haplotypes being either identical or very similar when originating from the  
368 same location. The trees obtained with either only protein-coding genes (CODON model) or protein-coding  
369 and ribosomal RNA genes but with a mixed model (DNA for rRNAs and CODON for DNA) were all  
370 consistent with the tree shown here and are therefore not included.

## 371 Patterns of introgression

372 We tested all four northern Germany populations (EIC, DOB, AR and SE) for admixture between the three  
373 reference species with  $f_3$  statistics tests (Table 3) and considered a Z-score  $< -3$  as significant (following  
374 Patterson, *et al.* 2012; Reich, *et al.* 2009). Negative and significant values ( $f_3 = -0.19$ ) using EIC as the test  
375 population and *D. galeata* and *D. longispina* as the source populations indicated mixed ancestry from these  
376 two or closely related populations. For population SE, the  $f_3$  test supports both admixed ancestry from *D.*  
377 *galeata* and *D. longispina* ( $f_3 = -0.09$ ) and *D. galeata* and *D. cucullata* ( $f_3 = -0.15$ ). All tests for population  
378 DOB and AR were positive providing no evidence of admixture events. The supported introgression events

379 are consistent with the results in our previous analyses conducted with ADMIXTURE and the ancestry  
380 painting approach.



381  
382 **Figure 5:** Maximum-likelihood tree reconstructed from mitochondrial protein-coding and ribosomal RNA genes of  
383 parental species, clones sampled in the water column and resting eggs sequenced in this study. The tree reveals distinct  
384 and highly supported clusters corresponding to *D. galeata*, *D. cucullata* and *D. longispina* mitotypes (as defined by  
385 the respective parental species and a sister taxa relationship between *D. galeata* and *D. cucullata*). Here, the best tree  
386 ( $\log L = -47950.82$ ) rooted with outgroup *D. laevis* is depicted. Magenta dots indicate Shimodaira-Hasegawa  
387 approximate likelihood ratio test values  $\geq 80\%$  and ultrafast bootstrap support values  $\geq 95\%$  calculated from 10,000  
388 bootstrap replicates (SH-aLRT / UFboot). The scale bar corresponds to 0.05 nucleotide substitutions per nucleotide  
389 site. KL11 was excluded due to missing data.

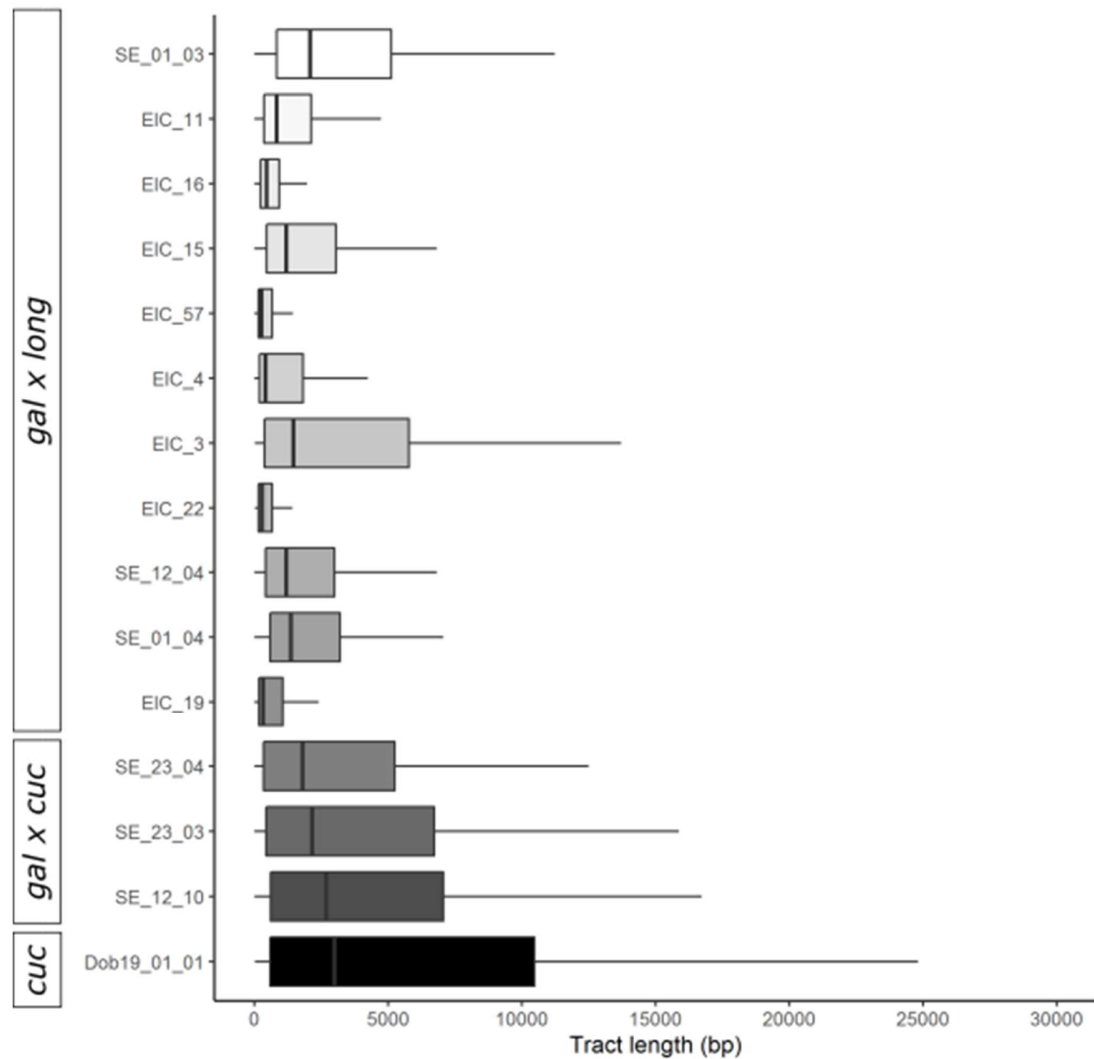
390 **Table 3:** Summary of the  $f_3$  statistic for admixture in the form (C; A, B). A significantly (Z-score < -3, in bold)  
 391 negative  $f_3$  value implies that the target population C is admixed. SE: standard error

Source population A	Source population B	Target population C	$f_3$	SE	Z-score	# Sites
gal	long	AR	5.48226	0.154392	35.509	1,072,106
gal	cuc	AR	0.245124	0.006208	39.487	791,570
long	cuc	AR	0.225393	0.005475	41.165	875,293
gal	long	Dob	0.152725	0.003233	47.242	752,265
gal	cuc	Dob	0.150027	0.003028	49.546	815,035
long	cuc	Dob	1.763147	0.049657	35.506	1,059,418
gal	long	EIC	-0.193114	0.003224	<b>-59.898</b>	864,328
gal	cuc	EIC	0.086614	0.002623	33.025	1,119,715
long	cuc	EIC	0.014322	0.001774	8.071	1,145,530
gal	long	SE	-0.093233	0.003222	<b>-28.939</b>	1,029,744
gal	cuc	SE	-0.154696	0.003924	<b>-39.426</b>	955,516
long	cuc	SE	0.288212	0.011203	25.727	1,092,186

392

393 We performed local ancestry inference with Loter (Dias-Alves, *et al.* 2018) to trace genome-wide  
 394 introgression among the hybrids and infer additional details about the parental species and backcross history  
 395 from haplotype information. The results were summarized genome-wide for the ancestry proportion,  
 396 heterozygosity of ancestry and the number of ancestry transitions where each ancestry tract is counted when  
 397 the state of an SNP changes to the other species or at the end of a scaffold (Figure 5). The three *D. galeata*  
 398 x *cucullata* hybrids were all found to have high *galeata* ancestry (73.9%-94.2%) and the individuals  
 399 SE\_23\_03 and SE\_23\_04 were confirmed as the offspring of a later-generation hybrid and a pure *galeata*  
 400 parent (Figure S6B).

401 Seven *D. galeata* x *longispina* hybrids had very high *galeata* ancestry (81.4%-97.9%), and visual inspection  
 402 of the ancestry tracts (Figure S6A) revealed very short *longispina* tracts and scaffolds with multiple  
 403 breakpoints indicating multiple generations of recombination. Four *D. galeata* x *longispina* hybrids had  
 404 lower *galeata* ancestry (27.7%-59.0%) and the presence of complete *longispina* scaffolds implying some  
 405 backcrossing with *longispina*. The haplotype phasing confirmed that the parents of all *D. galeata* x  
 406 *longispina* hybrids were also of hybrid origin. The average and maximum ancestry tract length for all *D.*  
 407 *galeata* x *longispina* hybrids is shorter than those for *D. galeata* x *cucullata* hybrids.



408

409 **Figure 6:** Distribution of the ancestry tract length where each ancestry tract represents the state of a SNP changing to  
410 the other species or the end of a scaffold in the local ancestry inference for each admixed individual and one non-  
411 admixed *D. cucullata* individual. The non-admixed individual displays the ancestry tract length distribution when all  
412 scaffolds derive from the same species. Hybrid type (according to ADMIXTURE analysis) is given on the left side.

## 413 Discussion

### 414 A reference genome for studying a species complex

415 *Daphnia* are a key species in freshwater habitats. Previous studies have established reference genomes for  
416 the model species *D. pulex* (Colbourne, *et al.* 2011; Ye, *et al.* 2017) and *D. magna* (Lee, *et al.* 2019). No  
417 high-quality reference genome for species belonging to the *Daphnia longispina* species complex was  
418 available so far. To date, it is unclear whether the ecological differentiation and/or intrinsic incompatibilities

419 drive and maintain divergence between DLSC species. Besides its utility for studies of hybridization events  
420 in the DLSC, the new assembly we present here will thus allow us to better understand the evolution of a  
421 key species in European freshwaters.

422 Even though the onset of the DLSC radiation was dated to 5-7 Mya based on nuclear and mitochondrial  
423 markers (Adamowicz, *et al.* 2009; Schwenk 1993; Taylor, *et al.* 1996), several factors confirm the suitability  
424 of this reference for all tested species. Mapping success and coverage of whole-genome data from *D.*  
425 *cucullata* and *D. longispina* to the reference genome were high, and we found no evidence of reference bias.  
426 This assembly clearly benefited from advances both in the sequencing technologies and assembly and post-  
427 processing algorithms since the first *Daphnia* genome (Colbourne, *et al.* 2011). The metrics used for  
428 assessing its quality reveal that in particular, the combination of long and short read technologies led to  
429 highly contiguous and accurate scaffolds. Although we likely did not recover the genome in its full length  
430 (133Mb out of an estimated 156Mb), and the N50 value is lower than those obtained for *D. pulex* (Ye, *et al.*  
431 2017) and *D. magna* (Lee, *et al.* 2019), iterative scaffolding allowed for a very efficient gap-closing, and an  
432 exceptionally low number of mismatches, compared to the other *Daphnia* assemblies.

### 433 Pervasive introgression in the *Daphnia longispina* species complex

434 We utilize a method that allows us to interrogate biological archives and analyze whole *Daphnia* genomes  
435 directly from the resting egg bank (Lack, *et al.* 2018) without hatching and culturing several clonal lineages.  
436 This provides a wide sweep of populations, past and present, with each egg being the product of local sexual  
437 recombination.

438 While no evidence of introgression was found in the DOB population, the three other locations host a variety  
439 of admixed genotypes. SE & EIC can even be considered hybridization hotspots with more than 60% of  
440 individuals having hybrid ancestry, as revealed in the ADMIXTURE analysis. However, Kong and Kubatko  
441 (2021) very recently showed that ADMIXTURE is sensitive to unequal contributions by the parental  
442 species, and we thus sought to support these inferences by  $f_3$  calculations and using an ancestry painting  
443 approach.

444 In DOB, ADMIXTURE delivered unequivocal results. Further, the  $f_3$  index indicated that no introgression  
445 was detectable in this population. However, the PCA plot shows that some of the DOB genotypes are near  
446 hybrid individuals. ANGSD results are similar but these genotypes nearer the parental species (Figure S3A).  
447 To address these slightly contradictory results, we therefore conducted an ancestry painting on two DOB  
448 genotypes, (12\_07 and 89\_02, Table S7). Both genotypes had a very low heterozygosity, thus confirming  
449 the ADMIXTURE and  $f_3$  outcomes. A possible explanation would be that these two genotypes carry  
450 variation that is not reflected in our limited sampling of the parental species. When comparing to  
451 microsatellite-based analysis including many more populations and data points (e.g. Thielsch, *et al.* 2009),  
452 the *D. galeata* cluster has often been larger and more diverse than the others. The seemingly two “stray”  
453 DOB genotypes are therefore likely well within the species variation boundaries. All mitochondrial  
454 haplotypes were clustered together in the phylogenetic reconstruction as well.

455 In AR, despite the high resting egg density found in the sediment, only very few could be successfully  
456 genotyped. While all inferences based on nuclear markers (PCA, ADMIXTURE,  $f_3$ ) indicated an absence  
457 of hybridization or introgression in this population, the mitochondrial phylogenetic reconstruction showed  
458 diverging results. From a nuclear point of view, all genotypes could be categorized as *D. cucullata*, but two  
459 out of four AR individuals presented the mitochondrial genome of another species, i.e., *D. galeata*.  
460 However, the phylogenetic reconstruction shows that the two AR mitochondrial haplotypes form a cluster  
461 separate from the main *D. galeata* cluster, which hints at different evolutionary history for these  
462 mitochondrial genomes. Such distinct lineages within a species and mito-nuclear discordances were also  
463 found by Thielsch, *et al.* (2017) in the DLSC and mitochondrial capture has been detected in other *Daphnia*  
464 species (Marková, *et al.* 2013). It is an interesting phenomenon in the DLSC that merits to be further  
465 investigated in the future with broader sampling.

466 In EIC, both pelagic samples and resting eggs were analyzed. Genotypes sampled alive from the water  
467 column were all inferred to be admixed to various degrees, three resting egg samples were also admixed,  
468 and the remaining five were assigned to either one of the 2 parental species. We conducted the ancestry



469 painting approach on all EIC individuals; the fixed sites heterozygosity of the individuals categorized as  
470 non-admixed in ADMIXTURE was indeed near zero (Table S7). Such high abundances of *D. galeata* x  
471 *longispina* hybrid resting eggs in periods of rapidly changing environmental conditions (i.e. eutrophication)  
472 have also been recorded in Lake Constance (Brede, *et al.* 2009). The high frequency of *D. galeata* x  
473 *longispina* hybrids observed here might be due to similar ecological history: the lake Eichbaumsee was  
474 created through sand excavation for construction work around ~40 years ago and is characterized by extreme  
475 eutrophication and even hypertrophy that could not be remediated. The presence of later-generation hybrids  
476 and backcrosses with *D. galeata* and *D. longispina* and short ancestry tract length suggest that both species  
477 have been present and hybridizing for most of the lake's short history, or even that it was colonized by  
478 individuals of hybrid origin. However, we only obtained contemporary samples for EIC and analysis of  
479 resting eggs from sediment cores are needed to distinguish between the two hypotheses.

480 In SE, diversity is high, both in terms of species combinations in admixed individuals and in terms of degrees  
481 of introgression. Although we analyzed eggs from sediments cores, they originated from the first centimeters  
482 and there is, therefore, no clear temporal pattern that separates the different hybrid combinations found here.  
483 Strikingly, while SE and DOB are geographically very close to each other (~10 km), and dispersal of resting  
484 eggs through e.g., waterfowl or storms would be possible (Figuerola, *et al.* 2005; Frisch, *et al.* 2007; Pietrzak  
485 and Slusarczyk 2006), the *Daphnia* communities are quite different. This might be due to their different  
486 eutrophication levels, reflect the fact that initial colonization was followed by the establishment of different  
487 species, or a combination of both. The observed diversity at such a small spatial scale underlines the mosaic  
488 nature of freshwater habitats and the usefulness of approaches including many populations to fully  
489 understand genetic diversity arising from colonization and hybridization events in the DLSC.

490 Previous studies using mitochondrial and few nuclear markers (e.g. Alric, *et al.* 2016; Thielsch, *et al.* 2012)  
491 were able to categorize hybrids into F1, F2 and backcrosses. However, due to the low resolution of the used  
492 markers, further categorizing and above all identification of genome-wide breaking points was not possible  
493 at the time. The *D. galeata* reference genome and resequencing data offer now a much higher resolution to

494 assess later generation hybrids and patterns across the genome. In general, hybrids identified in this study  
495 seem to have a complex history of multiple generations of hybridization and backcrossing with both parental  
496 species that we are not able to detangle using only ancestry paintings. The local ancestry inference revealed  
497 that the average ancestry tract length for *D. galeata x longispina* hybrids from EIC and SE is shorter than  
498 those for *D. galeata x cucullata* hybrids. There are several explanations for the observed pattern. One is that  
499 more generations of recombination led to shorter introgressed tracts, and the *D. galeata x longispina* hybrids  
500 are therefore the result of a greater number of sexual generations than the *D. galeata x cucullata* hybrids.  
501 The genomic mosaic of ancestry segments for all hybrid individuals is also characterized by multiple  
502 breakpoints within the same scaffolds, which is only possible after multiple generations of recombination.  
503 However, data on genome-wide recombination rates and selection are needed to reach solid conclusions  
504 about the correlation between tract length and age of the hybridization event in the individual's ancestors.  
505 Alternatively, reproductive isolation might be lower between *D. galeata* and *D. longispina* than between *D.*  
506 *galeata* and *D. cucullata*, thus leading to faster introgression in the former case.

507 As evidenced by the comparison of genomic windows of higher divergence between species pairs, the  
508 introgression pattern is not random: a given region exhibiting high  $F_{ST}$  values between the *D. galeata* and  
509 *D. longispina* genotypes is also likely to show similarly high  $F_{ST}$  values in the *D. galeata/ D. cucullata* pair.  
510 Further, some parts of the genome seem to be effectively shielded from introgression. About a quarter of all  
511 genes (4136) are in regions that are highly differentiated between at least two species and about 5% (859)  
512 in parts of the genome that are isolated among all three species of the complex. This is much more than  
513 expected by chance (Figure S3) and is thus likely due to selection against introgression. It seems plausible  
514 to search among these for genes that conserve the specific identity of the involved taxa, despite incomplete  
515 reproductive isolation. Genes responsible for the observed ecological divergence among the taxa (Schwenk,  
516 *et al.* 2000) or genetic incompatibilities are most likely candidates to be found in the observed divergent  
517 regions. Given the ancient divergence, the speciation process in the DLSC might have attained a selection-  
518 migration-drift equilibrium, for which there is growing empirical evidence in other species like stick insects  
519 (Riesch, *et al.* 2017), flycatchers (Burri, *et al.* 2015), and non-biting midges (Schreiber and Pfenninger

520 2020). However, the current snapshot could equally likely be a consequence of one or several pulses of  
521 hybridization. To assess the stability of the equilibrium, data showing that the introgression/selection  
522 process is ongoing and constant across an extended period of time would be required and *Daphnia* offers  
523 the unique opportunity to go back in time to test these alternative hypotheses.

#### 524 **New evidence for cytonuclear discordance**

525 The genome-wide perspective also elucidated discordance between nuclear and mitochondrial patterns. The  
526 phylogeny based on mitochondrial genomes conforms to previously inferred relationships in the DLSC and  
527 suggests *D. galeata* and *D. cucullata* are sister species, with *D. longispina* as an outgroup (Adamowicz, *et*  
528 *al.* 2009; Petrusek, *et al.* 2012). However, several of our analyses based on nuclear SNPs challenge this  
529 view and suggest different evolutionary histories for mitochondrial and nuclear genomes. The ancestry  
530 painting approach relies on the identification of fixed sites for species in a pairwise comparison. More sites  
531 were found to be fixed between *D. galeata* and *D. cucullata* (715,438) than between *D. galeata* and *D.*  
532 *longispina* (335,052), which implies a greater divergence between members of the former pair. Further,  $F_{ST}$   
533 values were on average higher between *D. galeata* and *D. longispina* (0.274) than between *D. galeata* and  
534 *D. cucullata* (0.364). Reports of cytonuclear discordance are common both in plants (e.g. Folk, *et al.* 2017;  
535 Huang, *et al.* 2014; Lee-Yaw, *et al.* 2019; Stephens, *et al.* 2015) and animals (e.g. Llopart, *et al.* 2014; Melo-  
536 Ferreira, *et al.* 2014; Sarver, *et al.* 2021). Several processes can lead to this discordance among closely  
537 related species: incomplete lineage sorting causing phylogenetic reconstructions based on mitochondrial  
538 markers to differ from the true phylogeny of the taxa, or selection causing the fixation of different  
539 mitochondrial genomes in different places from standing variation within species (e.g. Barrett and Schluter  
540 2008). Alternatively, cytonuclear discordance may reflect hybridization between species and cytoplasmic  
541 introgression, accompanied or not by selection (reviewed in Sloan, *et al.* 2017). The latter explanation would  
542 be quite conceivable in the DLSC.

## 543 Conclusion

544 We here provide the first high-quality resources to study genome-wide patterns of divergence in the *Daphnia*  
545 *longispina* species complex, an ecologically important taxon in European freshwater habitats. By  
546 quantifying intra- and interspecific diversity, we provide a first glimpse into introgressive hybridization and  
547 lay the ground for further studies aiming at understanding how species boundaries are maintained in the  
548 face of gene flow.

549 Unlike for *D. pulex* and *D. magna*, no linkage groups are known for any species of the DLSC. Hi-C  
550 sequencing data will be added in the future to order scaffolds into larger, potentially chromosome-scale  
551 scaffolds. Such an approach holds promise in a species complex where laboratory crossings for F2 panels  
552 and traditional mapping are nearly impossible. This will allow discovering structural variants, identifying  
553 recombination breakpoints along each chromosome and thus provide a deeper understanding of the  
554 introgression patterns observed here. The functional role of genes in the regions of high divergence  
555 uncovered through this first analysis is yet unclear and will be addressed in future studies.

556 Finally, wider sampling, with the inclusion of more populations as well as more members of the species  
557 complex, and the reconstruction of a nuclear based phylogeny are necessary to reach conclusions about the  
558 species relationships and eventually identify the causes of the pattern uncovered here.

559

## 560 Materials & Methods

### 561 Sampling

562 The clonal line used for genome sequencing and assembly, M5, was hatched from a resting egg isolated  
563 from the upper layers (first 5cm, corresponding to the years 2000-2010) of a sediment core taken in Lake  
564 Müggelsee in 2010. Further, single genotypes representing the parental species from various locations were  
565 used in this study, henceforth “parental species genotypes”. Most of them were established from individuals  
566 sampled from the water column and are still maintained through asexual reproduction as monoclonal

567 cultures in the laboratory. Thus, all individuals of a clonal line are the same genotype and can be pooled to  
568 achieve sufficient amounts of genomic DNA. The species identity for these genotypes was established  
569 through a combination of methods: morphology, mitochondrial sequences, and nuclear markers.

570 Sediment cores were collected from Dobersdorfer See (DOB), Selenter See (SE) and Arendsee (AR),  
571 Germany using a gravity corer (Uwitec, Mondsee, AT) (Table S4). Samples were taken from the deepest  
572 part of the lakes to minimize past disturbance of the sediment. Cores were cut horizontally into 1cm layers  
573 and the layers were stored at 4 °C in the dark to prevent hatching. Sediment rate of the three lakes was  
574 determined using radioisotope dating ( $^{137}\text{Cs}$  and  $^{210}\text{Pb}$ ).

575 In addition, lake sediment from the shoreline of Eichbaumsee (EIC) (Table S4), Germany was collected by  
576 hand and stored at 4 °C. The exact age of the sediment is unknown but the upper layers most likely contain  
577 recent eggs from the last few years. Zooplankton samples were taken from Eichbaumsee using a plankton  
578 net (mesh size 150  $\mu\text{m}$ ) from which six *Daphnia* clonal lines were established in a laboratory setting with  
579 artificial medium (Aachener Daphnien Medium, ADaM Klüttgen, *et al.* 1994).

580 All sampling locations are plotted in Figure 1A and information on all samples is provided in Table S2.

## 581 **Genome sequencing**

### 582 **DNA extraction for genome sequencing with Illumina & PacBio**

583 DNA was extracted from around 60 clonal M5 individuals collected from batch cultures maintained in  
584 ADaM, and fed with the algae *Acutodesmus obliquus*, cultivated in medium modified after (Zehnder and  
585 Gorham 1960). Extraction was conducted following a phenol chloroform-based protocol with an RNase  
586 step and subsequently sequenced on an Illumina HiSeq4000 at BGI China. Additionally, tissue samples with  
587 around 3000 individuals were sent to BGI for DNA extraction and PacBio sequencing.

## 588 Re-sequencing (Population genomics approach)

### 589 DNA extraction from batch cultures for re-sequencing

590 For clonal lines used as reference for the parental species, individuals were raised in batch cultures and  
591 treated with antibiotics prior to collection and storage at -20 or -80°C. DNA was extracted with either a  
592 phenol chloroform method, a (modified) CTAB protocol or a rapid desalting method (MasterPure™  
593 Complete DNA and RNA Purification Kit; Lucigen Corporation).

594 Total genomic DNA was isolated from 20 pooled adult *Daphnia* for each of the five EIC clonal lines using  
595 a CTAB extraction method (Doyle and Doyle 1987).

### 596 Whole Genome amplification on resting eggs for re-sequencing

597 To isolate *Daphnia* resting eggs from the sediment each sediment layer was sieved using a sieve with 125  
598 µm mesh size and small amounts of the remaining sediment were resuspended in distilled water. Ehippia  
599 were eye spotted under a stereomicroscope, counted and transferred to 1.5 mL tubes. The water was removed  
600 and ehippia stored at -20 °C in the dark until further analysis.

601 The ehippia were then opened under a binocular with insect needles and tweezers previously treated under  
602 a clean bench (UV sterilization) and with DNase away (Thermo Fisher). Eggs that were already damaged,  
603 had an uneven shape or were orange, which is evidence for degradation, were discarded. The resting egg  
604 separated from the ehippial casing washed in 15 µl sterile 1x PBS and then transferred in 1 µl 1x PBS to a  
605 new tube with 2 µl fresh 1x PBS. The isolated eggs were stored at -80 °C at least overnight.

606 For whole genome amplification of single eggs, the REPLI-g Mini Kit (Qiagen) was used. This kit is  
607 enabling unbiased amplification of genomic loci via Multiple Displacement Amplification (MDA). The  
608 isolated resting eggs were thawed on ice and the whole genome was amplified following the manufacturer's  
609 protocol for amplification of genomic DNA from blood or cells. Briefly, denaturation buffer was added to  
610 the prepared resting eggs in 3 µl 1x PBS and amplified by REPLI-g Mini DNA Polymerase under isothermal  
611 conditions for 16 hours.

612 The amplified product was quantified on a Nanodrop spectrophotometer (Thermo Fisher) and with a Qubit  
613 Fluorometer (Thermo Fisher). Successful amplifications were purified with 0.4 x Agencourt AMPure XP  
614 magnetic beads (Beckman Coulter) to remove small fragments and eluted in 60 µl 1x TE buffer.

615 Fragments of the mitochondrial gene 16S rRNA gene were amplified to check successful amplification of  
616 *Daphnia* DNA using the universal cladoceran primers S1 and S2 (Schwenk, *et al.* 1998) and a low presence  
617 of bacterial DNA using universal primers for the bacterial 16S rDNA gene (Nadkarni, *et al.* 2002). Only  
618 samples with a successful amplification of the *Daphnia* 16S fragment and low amplification of the bacterial  
619 16S fragment indicating low bacterial contamination were used for sequencing steps.

## 620 Library preparation and sequencing of re-sequencing samples

621 After quantification and quality control of the DNA using Nanodrop and Qubit instruments, libraries were  
622 prepared either directly in-house with the NEBNext® Ultra™ II DNA Library Prep Kit for Illumina® (New  
623 England Biolabs), or at the sequencing company Novogene (Cambridge, UK). Resequencing (paired-end  
624 150bp reads) was then performed either at Novogene (UK) Company Limited or the Functional Genomics  
625 Center (ETH Zurich and University of Zurich) on Illumina NovaSeq 6000 and HiSeq4000 instruments.

626 Details on the procedure used for each sample are provided in Table S2.

## 627 Genome assembly and annotation

628 We provide here a summarized version of the procedure used to assemble and annotate the genome.

629 Details can be found in Supplementary Methods.

## 630 Raw data QC

631 Illumina reads were trimmed and the adapter removed using a combination of Trimmomatic 0.38 (Bolger,  
632 *et al.* 2014), FastQC 0.11.7 (Andrews 2010) and MultiQC 1.6 (Ewels, *et al.* 2016) within autotrim 0.6.1  
633 (Waldvogel, *et al.* 2018). To filter out reads possibly originating from contamination from known sources  
634 (see below), a FastQ Screen like approach was chosen. In brief, the reads are separated by results of mapping  
635 behavior to different genomes. Positive controls consisted of genome data for other *Daphnia* species

636 (dmagna-v2.4 and *Daphnia\_pulex\_PA42\_v3.0*, see Supplementary Methods for accession numbers), and  
637 negative control i.e. sequences deemed undesirable for genome assembly consisted of genome data from  
638 human, bacteria, viruses and the algae used to feed the batch cultures. The resulting database comprised  
639 108,163 sequences (total sequence space 42.2 Gb). Both Illumina reads and PacBio subreads were mapped  
640 against the database with NextGenMap (Sedlazeck, *et al.* 2013) and minimap2 (Li 2018), respectively.

641 Reads did only pass the filtering if they either did not map to the database at all or had at least one hit against  
642 one of the two *Daphnia* genomes. Table 2 in Supplementary methods gives an overview of the effect of  
643 different filtering steps.

#### 644 **Assembly and contamination screening**

645 All paired and unpaired contamination filtered Illumina reads as well as the contamination filtered PacBio  
646 reads were used as input for RA 0.2.1 (<https://github.com/rvaser/ra>). Blobtools 1.0 (Laetsch and Blaxter  
647 2017) was used to screen the resulting assembly for possible unidentified contamination in the hybrid  
648 assembly. Briefly, bwa mem 0.7.17 (Li 2013) was used to map Illumina reads back to the assembly and  
649 taxonomic assignment was done by sequence similarity search with blastn 2.9.0+ (Camacho, *et al.* 2009).  
650 Contamination with different bacteria was clearly identifiable, and contigs with coverage below 10x and/or  
651 GC content above 50% were removed. Additionally, PacBio reads mapping to these contigs were removed  
652 to minimize false scaffolding in further steps. The contig corresponding to the mitochondrial genome was  
653 identified after a blast search against available mitochondrial genomes for this species and removed from  
654 the assembly.

#### 655 **Scaffolding and gap closing**

656 The blobtools filtered PacBio reads were used for scaffolding and gapclosing, which was conducted in three  
657 iterations. Each iteration consisted of a scaffolding step with SSPACE LongRead 1-1 (Boetzer and Pirovano  
658 2014), a gap closing step with LR Gapcloser ([https://github.com/CAFS-bioinformatics/LR\\_Gapcloser](https://github.com/CAFS-bioinformatics/LR_Gapcloser);  
659 commit 156381a), and a step to polish former gap parts with short reads using bwa mem 0.7.17-r1188 and



660 Pilon 1.23 (Walker, *et al.* 2014) in a pipeline developed to this effect, wtdbg2-racon-pilon.pl 0.4  
661 (<https://github.com/schell/wtdbg2-racon-pilon>).

## 662 Assembly quality assessment

663 Contiguity was analyzed with Quast 5.0.2 (Gurevich, *et al.* 2013) at different stages of the assembly process.  
664 Further, mapping rate, coverage and insert size distribution were assessed by mapping Illumina and PacBio  
665 reads with bwa mem and Minimap 2.17 respectively. To show absence of contamination in the assembly  
666 blobtools was ran as above. The genome size was estimated by dividing the mapped nucleotides by the  
667 mode of the coverage distribution of the Illumina reads by backmap 0.3  
668 (<https://github.com/schell/backmap>), resulting in 156.86Mb (with the obtained assembly length amounting  
669 to 85% of this estimated length). Additionally, the genome size was estimated using a k-mer based approach  
670 by creating a histogram from raw Illumina reads with Jellyfish 1.1.12 (Marçais and Kingsford 2011) and  
671 running the GenomeScope web application (<http://qb.cshl.edu/genomescope/>) resulting in a genome size  
672 estimate of 150.6Mb.

673 Completeness in terms of single copy core orthologs of the final scaffolds was assessed with BUSCO 3.0.2  
674 (Simão, *et al.* 2015), using the Arthropoda set (odb9).

## 675 Genome Annotation

676 RepeatModeler 2.0 (Smit and Hubley 2015) was run to identify *D. galeata* specific repeats. The 1,115  
677 obtained repeat families were combined with 237 *D. pulex* and 1 *D. pulicaria* repeat sequences from  
678 RepBase release 20181026 to create the final repeat library. The genome assembly was then soft masked  
679 with RepeatMasker 4.1.0 (Smit, *et al.* 2013-2015), resulting in 21.9% of the assembly being masked.

680 Gene prediction models were produced with Augustus 3.3.2 (Stanke, *et al.* 2008), GeneMark ET  
681 4.48\_3.60\_lic (Lomsadze, *et al.* 2005) and SNAP 2006-07-28 (Korf 2004). The Augustus model was based  
682 on the soft masked assembly and the *D. galeata* transcriptome (HAFN01.1, Huylmans, *et al.* 2016). The  
683 GeneMark model was obtained by first mapping trimmed RNAseq reads to the assembly with HISAT 2.1.0

684 (Kim, *et al.* 2019) and then processing the resulting bam file with bam2hints and filterIntronsFindStrand.pl  
685 from Augustus to create a gff file with possible introns, which was finally fed into GeneMark.

686 The structural annotation was conducted in MAKER 2.31.10 (Holt and Yandell 2011). Briefly, the  
687 unmasked genome assembly, the species own transcriptome assembly as ESTs, the complete Swiss-Prot  
688 2019\_10 (UniProt Consortium, 2019) and the protein sequences resulting from *D. magna* (Lee, *et al.* 2019),  
689 as well as *D. pulex* (Ye, *et al.* 2017) genome annotations as protein evidence, were used as input for  
690 MAKER. In total three iterations of MAKER with retraining of the Augustus and SNAP model in between  
691 the iterations were conducted.

692 The quality of the structural annotation was assessed by comparing values as number of genes, gene space,  
693 etc. to existing annotations for other *Daphnia* genomes. Furthermore, core orthologs from BUSCO's  
694 Arthropoda (odb9) set and conserved domain arrangements from the Arthropoda reference set of DOGMA  
695 3.4 (Dohmen, *et al.* 2016) were searched in the annotated protein set.

696 The functional annotation was conducted using InterProScan 5.39-77.0 (Jones, *et al.* 2014) as well as a blast  
697 against the Swiss-Prot 2019\_10.

## 698 Population samples

### 699 Raw data QC and contamination check

700 The quality of raw reads was checked using FastQC v0.11.5. Adapter trimming and quality filtering were  
701 performed using Trimmomatic v0.36 with the following parameters: ILLUMINACLIP:TruSeq3-  
702 PE.fa:2:30:10 TRAILING:20 SLIDINGWINDOW:4:15 MINLEN:70. For samples sequenced on a  
703 NovaSeq6000 instrument and presenting a typical polyG tail, the program fastp (Chen, *et al.* 2018) was used  
704 for trimming as well. To assess contamination in the WGA samples FastQ Screen v0.14.0 with the bwa  
705 mapping option was used (Wingett and Andrews 2018). A custom database was built to map trimmed reads  
706 against possible contaminants that included general common contaminants such as *Homo sapiens*, the  
707 UniVec reference database, a bacterial and a viral reference set as well as the *D. galeata* genome and

708 *Acutodesmus obliquus* draft genome (see Supplementary Methods for accession numbers). Samples with  
709 <25 % reads mapped to the *D. galeata* genome (and 25% contamination) were excluded from further  
710 analysis because the whole amplification of the resting egg most likely failed.

#### 711 Mapping to reference genome and variant calling

712 The variant calling was performed within the Genome Analysis Toolkit (GATK v4.1.4.0; McKenna, *et al.*  
713 2010) program according to GATK4 best practices (Van der Auwera, *et al.* 2013). The trimmed reads were  
714 mapped to the *D. galeata* genome using the BWA-MEM algorithm in BWA v0.7.17 with the -M parameter  
715 and adding read group identifiers for Picard compatibility (Li and Durbin 2009). PCR duplicates were  
716 marked and filtered out in the BAM file using Picard v2.21.1 (<http://broadinstitute.github.io/picard/>).

717 To call variants for each sample GATK HaplotypeCaller in GATK was used with the --emitRefConfidence  
718 GVCF option resulting in a genomic variant call format (gVCF) file with information on each position for  
719 each individual (Poplin, *et al.* 2018). All gVCF files were consolidated using CombineGVCFs and joint  
720 genotyping was performed with GenotypeGVCFs. The VCF file was filtered to include only SNPs and hard  
721 filtering was performed to remove variants with a QualByDepth <10, StrandOddsRatio >3, FisherStrand  
722 >60, mapping quality <40, MappingQualityRankSumTest <-8 and ReadPosRankSumTest <-5.

723 Subsequently, we removed sites with either very high coverage (>450) or for which genotypes were missing  
724 for more than 20% of the individuals using VCFtools v0.1.5 (Danecek, *et al.* 2011). The final SNP data set  
725 for downstream analyses included 3,240,339 SNPs across the 49 samples.

726 In addition, GenotypeGVCFs was run with the --include-non-variant-sites option to output all variant as  
727 well as invariant genotyped sites. The final invariant data set included 127,530,229 sites after removing  
728 indels and multi-allelic sites with BCFtools v1.9 (Li 2011) and is used for population genomic analysis to  
729 be able to calculate the total number of genotyped sites (variant and invariant) within a genomic window.

730 As we mapped all different species to the reference *D. galeata* genome we assessed possible reference bias  
731 by checking the distribution of reference and alternative alleles observed at heterozygous genotypes based

732 on Pinsky, *et al.* (2021). We pooled all genotypes which were unambiguously assigned to either of the  
733 parental species clusters *D. galeata*, *D. cucullata* and *D. longispina* as was done for the population genomic  
734 parameters (Table S2) or classified as hybrids using ADMIXTURE inference. Without reference bias, we  
735 would expect that in heterozygous genotypes the reference and the alternative allele are on average  
736 represented by 50% of the reads. An indication of reference bias would be that the *D. galeata* reference  
737 allele would be more frequent.

## 738 Phylogenetic and population genetics inferences

### 739 Mitochondrial genome assemblies and phylogenetic analyses

740 All reads were used to produce mitochondrial genome assemblies using the “de novo assembly” and “find  
741 mitochondrial scaffold” modules provided in MitoZ v2.4 with default settings (Meng, *et al.* 2019). For some  
742 samples, this was not sufficient and we used two approaches to recover a complete mitogenome: either the  
743 mitochondrial baiting and iterative mapping implemented in MITObim v1.9.1 (Hahn, *et al.* 2013) with the  
744 *D. galeata* mitochondrial reference genome or the modified baiting and iterative mapping in GetOrganelle  
745 v1.7.1 (Jin, *et al.* 2020) with the animal database and k-mer values set to 21, 45, 65, 85 and 105. The  
746 procedure used for each dataset is given in Table S5.

747 We annotated the mitochondrial genome assemblies with the mitochondrial annotation web server MITOS2  
748 (Bernt, *et al.* 2013) using the mitochondrial codon code 05 for invertebrates. Automated genome annotation  
749 identified thirteen protein-coding genes (PCGs), two ribosomal RNA genes (rRNAs), and twenty-two  
750 transfer RNA genes (tRNAs). Initially, the mitochondrial genes (PCGs and rRNAs, Table S6) were  
751 individually aligned with MUSCLE v3.8.1551 (Edgar 2004) and visually checked for their quality. The  
752 mitochondrial genome assemblies with discrepancies, i.e., a lot of missing data and/or split features were  
753 excluded from further analysis. The final data set (Table S6) included 44 mitochondrial genomes from this  
754 study and the previously published mitochondrial genome of *Daphnia laevis* (Martins Ribeiro, *et al.* 2019/  
755 accession number: NC\_045243.1/, accession number: NC\_045243.1). The mitochondrial genes of the final  
756 data set were individually realigned with MUSCLE v3.8.1551 (Edgar 2004) and MACSE v2.05 (Ranwez,

757 *et al.* 2018) and concatenated into a mitochondrial DNA matrix (Table S6) using SequenceMatrix v1.8.1  
758 (Vaidya, *et al.* 2011). During this step, we used MACSE v2.05 to realign PCG genes keeping the information  
759 about codon position (gene partitioning) and to remove STOP codons. The final dataset consisted of the  
760 concatenation matrix of the thirteen protein-coding (PCGs) and the two structural ribosomal RNA (rRNAs)  
761 genes. With this alignment, phylogenetic trees were reconstructed using IQ-TREE v1.6.12 (Nguyen, *et al.*  
762 2015). We initially partitioned the alignment into a full partition model, i.e., each gene and all three codon  
763 positions for PCGs, and then ran IQ-TREE with partition analyses (-spp, Chernomor, *et al.* 2016),  
764 ModelFinder (-m MFP+MERGE, Kalyaanamoorthy, *et al.* 2017) and 10,000 ultrafast bootstrap (-bb 10000,  
765 Hoang, *et al.* 2018) and SH-like approximate likelihood ratio test (-alrt 10000, Guindon, *et al.* 2010)  
766 replicates. The resulting trees were visualized in R (R Core Team 2017) using the multifunctional  
767 phylogenetics package phytools (Revell 2012).

#### 768 Ancestry and population structure

769 A principal component analysis was conducted in R v3.6.2 (R Core Team 2017) with the package  
770 SNPRelate v1.20.1 (Zheng, *et al.* 2012). Linkage disequilibrium (LD) was calculated within a 500-kb  
771 sliding window and LD-pruned for  $r^2$  values  $>0.5$  before conducting the PCA for all sites using the  
772 snpgdsPCA function with default settings. The relative large LD value was chosen because clonal  
773 reproduction and the overlap of generations due to diapause leads to increased linkage disequilibrium in  
774 *Daphnia* (Brede, *et al.* 2009).

775 Genetic admixture was estimated using ADMIXTURE v1.3.0 (Alexander and Lange 2011). The SNP set  
776 VCF file was converted to BED format using plink v1.90b6.13 (Chang, *et al.* 2015). The log-likelihood  
777 values were estimated for one to five genetic clusters (K) of ancestral populations and admixture analysis  
778 were run for the most appropriate K value with 10-fold cross-validation. We also conducted the PCA and  
779 Admixture analysis using PCAngsd implemented in ANGSD and NgsAdmix, respectively (Korneliussen,  
780 *et al.* 2014) to take genotype likelihoods into account (details in Supplementary methods). The results did  
781 not differ substantially and are shown in Figure S3.

782 However, using such a population genetic clustering approach to estimate ancestry coefficients is not  
783 directly equivalent to the proportion of hybrid ancestry in each individual and should be interpreted with  
784 caution (Kong and Kubatko 2021; Lawson, *et al.* 2018). The results of the ADMIXTURE analysis suggested  
785 that the dataset included hybrids between *D. longispina* and *D. galeata* as well as *D. cucullata* and *D.*  
786 *galeata*. We then followed the “ancestry painting” procedure outlined in Barth, *et al.* (2020) and Runemark,  
787 *et al.* (2018b), and classified sites according to their  $F_{ST}$  values when comparing parental species sets. Unlike  
788 the PCA and the admixture analysis, this approach requires the user to define parental genotypes; the  
789 individuals belonging to these sets are indicated with stars in Figure 1C. Fixed sites are those where a  
790 specific allele is fixed in all individuals belonging to one parental species and another allele fixed in the  
791 other parental species. To show the ancestry of the hybrid individuals each fixed site was plotted in an  
792 “ancestry painting” if at least 80% of genotypes were complete using available ruby scripts  
793 ([https://github.com/mmatschiner/tutorials/tree/master/analysis\\_of\\_introgession\\_with\\_snp\\_data](https://github.com/mmatschiner/tutorials/tree/master/analysis_of_introgession_with_snp_data)). These  
794 scripts calculate the heterozygosity of each individual and visualize regions that are possibly affected by  
795 introgression. The mitochondrial genome assembly from each individual was used to determine the maternal  
796 species and the proportion of the genome derived from the maternal species was then calculated for each  
797 hybrid. For gal x cuc hybrids the hybrid index scale ranges from 0 (gal) to 1 (cuc) and for gal x long hybrids  
798 from 0 (long) to 1 (gal).

### 799 Window-based population parameters

800 To assess genome-wide genetic differentiation between the clusters identified with admixture, we calculated  
801 nucleotide diversity ( $\pi$ ), between-taxon differentiation ( $F_{ST}$ ), and between-taxon divergence ( $d_{xy}$ ) using the  
802 Python script popgenWindows.py ([github.com/simonhmartin/genomics\\_general](https://github.com/simonhmartin/genomics_general) release 0.3, Martin, *et al.*  
803 2020) with a sliding 100-kb window, a step size of 10kb and at least 20kb genotyped sites within each  
804 window. To compare species pairs we only considered individuals assigned to parental species based on  
805 ADMIXTURE results (Table S2). In addition, we also calculated these parameters using one randomly  
806 chosen individual from each population per species to check if the estimates are biased because of the  
807 overrepresentation of some populations in a species group (Table S2).

808 Sets of outlier windows were defined as those with  $F_{ST}$  values in the upper 95<sup>th</sup> percentile of the distribution  
809 for each of the 3 pairwise comparisons. Further, the genes in these windows were extracted using the  
810 annotation file. We used a randomization approach to assess whether the observed intersections (i.e. outlier  
811  $F_{ST}$  windows occurring in both species) between all seven possible species comparisons are larger or smaller  
812 than expected by chance. For this, we randomly drew the observed number of windows, respectively genes  
813 from the total number of 10kb windows in the assembly (13,330), respectively the total number of annotated  
814 genes (15,845) without replacement and calculated the intersections for all possible comparisons. We  
815 compared the resulting intersections from 1,000 replicates with the observed values (Figure S5).

### 816 [Inferring introgression](#)

817 To identify admixture among three populations we calculated the  $f_3$  statistic with ADMIXTOOLS v 7.0  
818 (Patterson, *et al.* 2012) implemented in the admixr package in R (Petr, *et al.* 2019). We used two parental  
819 source populations (A and B) and the target population (C) in the form (C; A, B). Significantly negative  $f_3$   
820 statistic indicates that population C is a mixture of populations A and B or closely related populations.

### 821 [Local ancestry inference](#)

822 To prepare the SNP set, Beagle v4.1 was used to phase and impute genotypes with 10,000 bp step size and  
823 1000 bp overlapping sliding windows (Browning and Browning 2009). Local ancestry inference was  
824 conducted with Loter (Dias-Alves, *et al.* 2018) which infers the origin of each SNP in an admixed individual  
825 from two ancestral source populations and doesn't require additional biological parameters. The respective  
826 two parental species populations were used to reconstruct the ancestry tracts of the three putative *galeata* x  
827 *cucullata* hybrid individuals and eleven putative *galeata* x *longispina* hybrid individuals using Loter with  
828 default settings.

### 829 [Authors contributions](#)

830 MC, JN, AT, SD and MHM performed the sampling and wet lab work. TS, JN and AT conducted the  
831 genome assembly and annotation. JN conducted the population genomic analysis and mitochondrial genome

832 reconstruction. TH and MHM conducted the phylogenetic analyses. JN, TS, MC, MP, MHM and TH wrote  
833 the manuscript, and all authors edited and contributed to the final version. All authors gave final approval  
834 for preprint deposition and publication.

## 835 Data Availability

836 Genome assembly, annotation and read data (Illumina and PacBio) for the genotype M5 are stored under  
837 accession number PRJEB42807. Short read data from resequencing are available in the European  
838 Nucleotide Archive under accession numbers ERS5080327- ERS5080375, ERS4993274 and ERS4993282.  
839 The annotation used in the present analysis is deposited in Zenodo (DOI: 10.5281/zenodo.4479324),  
840 together with supplementary information on the mitochondrial tree (alignment file).

## 841 Acknowledgments

842 We thank LOEWE-TBG for providing sequencing funds. TH and MM were supported by the grant  
843 “SeeWandel: Life in Lake Constance - the past, present and future” within the framework of the Interreg V  
844 programme “Alpenrhein-Bodensee-Hochrhein (Germany/Austria/Switzerland/Liechtenstein)”, which  
845 funds are provided by the European Regional Development Fund as well as the Swiss Confederation and  
846 cantons. MHM was supported by the Austrian Science Fund (FWF): P29667-B25 and J 3774. The funders  
847 had no role in study design, data collection and analysis, decision to publish, or preparation of the  
848 manuscript. The computational results presented here have been achieved in (part) using the LEO HPC  
849 infrastructure of the University of Innsbruck. Some of the data produced and analyzed in this paper were  
850 generated in collaboration with the Genetic Diversity Centre (GDC), ETH Zurich.

851 We thank Jae-Seong Lee and Zhiqiang Ye for giving us access to genome annotations for *Daphnia magna*  
852 (SK strain) and *Daphnia pulex* (PA42 strain), respectively. Miklós Bálint provided sediment samples and  
853 isotope data for the Arendsee lake. We thank Michael Matschiner for his help with ancestry painting and  
854 two anonymous reviewers for their comments on a previous version of the manuscript.



## 855 References

- 856 Abbott R, *et al.* 2013. Hybridization and speciation. *J. Evol. Biol.* 26: 229-246. doi: 10.1111/j.1420-  
857 9101.2012.02599.x
- 858 Adamowicz SJ, Petrusek A, Colbourne JK, Hebert PDN, Witt JDS 2009. The scale of divergence: A  
859 phylogenetic appraisal of intercontinental allopatric speciation in a passively dispersed freshwater  
860 zooplankton genus. *Mol. Phylogenet. Evol.* 50: 423-436. doi: 10.1016/j.ympev.2008.11.026
- 861 Alexander DH, Lange K. 2011. Enhancements to the ADMIXTURE algorithm for individual ancestry  
862 estimation. *BMC Bioinform.* 12: 246
- 863 Alric B, *et al.* 2016. Local human pressures influence gene flow in a hybridizing *Daphnia* species complex.  
864 *J. Evol. Biol.* 29: 720-735. doi: 10.1111/jeb.12820
- 865 Andrews S. 2010. FastQC: A quality control tool for high throughput sequence data.
- 866 Arnold ML, Martin NH 2009. Adaptation by introgression. *J. Biol.* 8.
- 867 Barrett RD, Schluter D 2008. Adaptation from standing genetic variation. *Trends Ecol. Evol.* 23: 38-44.
- 868 Barth JM, *et al.* 2020. Stable species boundaries despite ten million years of hybridization in tropical eels.  
869 *Nat. Commun.* 11: 1-13.
- 870 Barton NH, Hewitt GM 1985. Analysis of hybrid zones. *Annu. Rev. Ecol. Syst.* 16: 113-148.
- 871 Beaton MJ, Hebert PD 1994. Variation in chromosome numbers of *Daphnia* (Crustacea, Cladocera).  
872 *Hereditas* 120: 275-279.
- 873 Bernt M, *et al.* 2013. MITOS: improved de novo metazoan mitochondrial genome annotation. *Mol.*  
874 *Phylogenet. Evol.* 69: 313-319.
- 875 Billiones R, Brehm M, Klee J, Schwenk K 2004. Genetic identification of *Hyalodaphnia* species and  
876 interspecific hybrids. *Hydrobiologia* 526: 43-53.
- 877 Boetzer M, Pirovano W 2014. SSPACE-LongRead: scaffolding bacterial draft genomes using long read  
878 sequence information. *BMC Bioinform.* 15: 211.
- 879 Bolger AM, Lohse M, Usadel B 2014. Trimmomatic: a flexible trimmer for Illumina sequence data.  
880 *Bioinformatics* 30: 2114-2120. doi: 10.1093/bioinformatics/btu170
- 881 Brede N, *et al.* 2009. The impact of human-made ecological changes on the genetic architecture of *Daphnia*  
882 species. *P Natl Acad Sci USA* 106: 4758-4763. doi: 10.1073/pnas.0807187106
- 883 Brede N, *et al.* 2006. Microsatellite markers for European *Daphnia*. *Mol. Ecol. Notes* 6: 536-539. doi:  
884 10.1111/j.1471-8286.2005.01218.x
- 885 Browning BL, Browning SR 2009. A unified approach to genotype imputation and haplotype-phase  
886 inference for large data sets of trios and unrelated individuals. *Am. J. Hum. Genet.* 84: 210-223.
- 887 Burri R, *et al.* 2015. Linked selection and recombination rate variation drive the evolution of the genomic  
888 landscape of differentiation across the speciation continuum of *Ficedula* flycatchers. *Genome Res.* 25:  
889 1656-1665.

- 890 Butlin RK, *et al.* 2014. Parallel evolution of local adaptation and reproductive isolation in the face of gene  
891 flow. *Evolution* 68: 935-949.
- 892 Camacho C, *et al.* 2009. BLAST+: architecture and applications. *BMC Bioinform.* 10: 421.
- 893 Canestrelli D, *et al.* 2017. Climate change promotes hybridisation between deeply divergent species. *PeerJ*  
894 5: e3072
- 895 Chang CC, *et al.* 2015. Second-generation PLINK: rising to the challenge of larger and richer datasets.  
896 *Gigascience* 4: s13742-015-0047-8
- 897 Chen S, Zhou Y, Chen Y, Gu J 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*  
898 34: i884-i890.
- 899 Chernomor O, Von Haeseler A, Minh BQ 2016. Terrace aware data structure for phylogenomic inference  
900 from supermatrices. *Syst. Biol.* 65: 997-1008.
- 901 Colbourne JK, *et al.* 2011. The Ecoresponsive Genome of *Daphnia pulex*. *Science* 331: 555-561. doi:  
902 10.1126/science.1197761
- 903 Cordellier M, Wojewodzic MW, Wessels M, Kuster CJ, Von Elert E. 2021. Next-generation sequencing of  
904 DNA from resting eggs: signatures of eutrophication in a lake's sediment. *Zoology* 145: 125895 doi:  
905 10.1016/j.zool.2021.125895
- 906 Cornetti L, Fields PD, Van Damme K, Ebert D 2019. A fossil-calibrated phylogenomic analysis of *Daphnia*  
907 and the Daphniidae. *Mol. Phylogenet. Evol.* 137: 250-262.
- 908 Cousyn C, *et al.* 2001. Rapid, local adaptation of zooplankton behavior to changes in predation pressure in  
909 the absence of neutral genetic changes. *P Natl Acad Sci USA* 98: 6256-6260. doi: 10.1073/pnas.111606798
- 910 Danecek P, *et al.* 2011. The variant call format and VCFtools. *Bioinformatics* 27: 2156-2158.
- 911 Dias-Alves T, Mairal J, Blum MG 2018. Loter: A software package to infer local ancestry for a wide range  
912 of species. *Mol. Biol. Evol.* 35: 2318-2326.
- 913 Dlouha S, *et al.* 2010. Identifying hybridizing taxa within the *Daphnia longispina* species complex: a  
914 comparison of genetic methods and phenotypic approaches. *Identifying hybridizing taxa within the Daphnia*  
915 *longispina* species complex: a comparison of genetic methods and phenotypic approaches 643: 107-122.  
916 doi: DOI 10.1007/s10750-010-0128-8
- 917 Doellman MM, *et al.* 2018. Genomic differentiation during speciation-with-gene-flow: Comparing  
918 geographic and host-related variation in divergent life history adaptation in *Rhagoletis pomonella*. *Genes* 9:  
919 262.
- 920 Dohmen E, Kremer LP, Bornberg-Bauer E, Kemena C 2016. DOGMA: domain-based transcriptome and  
921 proteome quality assessment. *Bioinformatics* 32: 2577-2581.
- 922 Doyle JJ, Doyle JL 1987. A rapid DNA isolation procedure for small quantities of fresh leaf tissue.  
923 *Phytochem. Bull.* 19: 11-15.
- 924 Dziuba MK, *et al.* 2020. Temperature increase altered *Daphnia* community structure in artificially heated  
925 lakes: a potential scenario for a warmer future. *Sci. Rep.* 10: 13956
- 926 Edgar RC 2004. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic*  
927 *Acids Res* 32: 1792-1797. doi: 10.1093/nar/gkh340

- 928 Ewels P, Magnusson M, Lundin S, Källér M 2016. MultiQC: summarize analysis results for multiple tools  
929 and samples in a single report. *Bioinformatics* 32: 3047-3048.
- 930 Figuerola J, Green AJ, Michot TC 2005. Invertebrate eggs can fly: evidence of waterfowl-mediated gene  
931 flow in aquatic invertebrates. *Am. Nat.* 165: 274-280.
- 932 Flaxman SM, Wacholder AC, Feder JL, Nosil P 2014. Theoretical models of the influence of genomic  
933 architecture on the dynamics of speciation. *Mol. Ecol.* 23: 4074-4088.
- 934 Folk RA, Mandel JR, Freudenstein JV 2017. Ancestral gene flow and parallel organellar genome capture  
935 result in extreme phylogenomic discord in a lineage of angiosperms. *Syst. Biol.* 66: 320-337.
- 936 Fraïsse C, Roux C, Welch JJ, Bierne N 2014. Gene-flow in a mosaic hybrid zone: is local introgression  
937 adaptive? *Genetics* 197: 393-951.
- 938 Frisch D, Green AJ, Figuerola J 2007. High dispersal capacity of a broad spectrum of aquatic invertebrates  
939 via waterbirds. *Aquat. Sci.* 69: 568-574.
- 940 Frisch D, *et al.* 2014. A millennial-scale chronicle of evolutionary responses to cultural eutrophication in  
941 *Daphnia*. *Ecol. Lett.* 17: 360-368. doi: 10.1111/ele.12237
- 942 Gannon JE, Stemberger RS 1978. Zooplankton (Especially Crustaceans and Rotifers) as Indicators of Water  
943 Quality. *Trans. Am. Microsc. Soc.* 97: 16-35. doi: 10.2307/3225681
- 944 Gliwicz ZM. 1990. Why do cladocerans fail to control algal blooms? In. *Bio-manipulation Tool for Water*  
945 *Management*: Springer. p. 83-97.
- 946 Griebel J, Gießler S, Yin M, Wolinska J 2016. Parental and hybrid *Daphnia* from the *D. longispina* complex:  
947 long-term dynamics in genetic structure and significance of overwintering modes. *J. Evol. Biol.* 29: 810-  
948 823.
- 949 Guindon S, *et al.* 2010. New algorithms and methods to estimate maximum-likelihood phylogenies:  
950 assessing the performance of PhyML 3.0. *Syst. Biol.* 59: 307-321.
- 951 Gurevich A, Saveliev V, Vyahhi N, Tesler G 2013. QUAST: quality assessment tool for genome assemblies.  
952 *Bioinformatics* 29: 1072-1075.
- 953 Hahn C, Bachmann L, Chevreux B. 2013. Reconstructing mitochondrial genomes directly from genomic  
954 next-generation sequencing reads—a baiting and iterative mapping approach. *Nucleic Acids Res* 41: e129
- 955 Harrison RG 1986. Pattern and process in a narrow hybrid zone. *Heredity* 56: 337-349.
- 956 Herrmann M, Henning-Lucass N, Cordellier M, Schwenk K 2017. A genotype–phenotype association  
957 approach to reveal thermal adaptation in *Daphnia galeata*. A genotype–phenotype association approach to  
958 reveal thermal adaptation in *Daphnia galeata* 327: 53-65. doi: 10.1002/jez.2070
- 959 Hoang DT, Chernomor O, Von Haeseler A, Minh BQ, Vinh LS 2018. UFBoot2: improving the ultrafast  
960 bootstrap approximation. *Mol. Biol. Evol.* 35: 518-522.
- 961 Holt C, Yandell M. 2011. MAKER2: an annotation pipeline and genome-database management tool for  
962 second-generation genome projects. *BMC Bioinform.* 12: 491
- 963 Huang DI, Hefer CA, Kolosova N, Douglas CJ, Cronk QC 2014. Whole plastome sequencing reveals deep  
964 plastid divergence and cytonuclear discordance between closely related balsam poplars, *Populus*  
965 *balsamifera* and *P. trichocarpa* (Salicaceae). *New Phytol.* 204: 693-703.

- 966 Huylmans AK, López Ezquerro A, Parsch J, Cordellier M 2016. De Novo Transcriptome Assembly and  
967 Sex-Biased Gene Expression in the Cyclical Parthenogenetic *Daphnia galeata*. *Genome Biol. Evol* 8: 3120-  
968 3139. doi: 10.1093/gbe/evw221
- 969 Jin J-J, *et al.* 2020. GetOrganelle: a fast and versatile toolkit for accurate de novo assembly of organelle  
970 genomes. *Genome Biol.* 21: 1-31.
- 971 Jones P, *et al.* 2014. InterProScan 5: genome-scale protein function classification. *Bioinformatics* 30: 1236-  
972 1240.
- 973 Kaiser TS, von Haeseler A, Tessmar-Raible K, Heckel DG 2021. Timing strains of the marine insect *Clunio*  
974 *marinus* diverged and persist with gene flow. *Mol. Ecol.* n/a. doi: 10.1111/mec.15791
- 975 Kalyaanamoorthy S, Minh BQ, Wong TK, von Haeseler A, Jermin LS 2017. ModelFinder: fast model  
976 selection for accurate phylogenetic estimates. *Nat. Methods* 14: 587-589.
- 977 Keller B, Wolinska J, Manca M, Spaak P 2008. Spatial, environmental and anthropogenic effects on the  
978 taxon composition of hybridizing *Daphnia*. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 363: 2943-2952. doi:  
979 10.1098/rstb.2008.0044
- 980 Kim D, Paggi JM, Park C, Bennett C, Salzberg SL 2019. Graph-based genome alignment and genotyping  
981 with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* 37: 907-915.
- 982 Klüttgen B, Dülmer U, Engels M, Ratte HT 1994. ADaM, an artificial freshwater for the culture of  
983 zooplankton. *Water Res.* 28: 743-746.
- 984 Kong S, Kubatko LS 2021. Comparative Performance of Popular Methods for Hybrid Detection using  
985 Genomic Data. *Syst. Biol.* 70: 891-907. doi: 10.1093/sysbio/syaa092
- 986 Korf I. 2004. Gene finding in novel genomes. *BMC Bioinform.* 5: 59
- 987 Korneliusson TS, Albrechtsen A, Nielsen R. 2014. ANGSD: analysis of next generation sequencing data.  
988 ANGSD: analysis of next generation sequencing data 15: 356
- 989 Lack JB, Weider LJ, Jeyasingh PD 2018. Whole genome amplification and sequencing of a *Daphnia* resting  
990 egg. *Mol. Ecol. Resour.*: 1-10. doi: 10.1111/1755-0998.12720
- 991 Laetsch DR, Blaxter ML. 2017. BlobTools: Interrogation of genome assemblies. *F1000Research*: 6:1287
- 992 Lampert W, Sommer U. 2007. Limnoecology: The Ecology of Lakes and Streams. Oxford: Oxford  
993 University Press.
- 994 Lawson DJ, Van Dorp L, Falush D 2018. A tutorial on how not to over-interpret STRUCTURE and  
995 ADMIXTURE bar plots. A tutorial on how not to over-interpret STRUCTURE and ADMIXTURE bar plots  
996 9: 1-11.
- 997 Lee-Yaw JA, Grassa CJ, Joly S, Andrew RL, Rieseberg LH 2019. An evaluation of alternative explanations  
998 for widespread cytonuclear discordance in annual sunflowers (*Helianthus*). *New Phytol.* 221: 515-526.
- 999 Lee B-Y, *et al.* 2019. The genome of the freshwater water flea *Daphnia magna*: A potential use for  
1000 freshwater molecular ecotoxicology. *Aquat. Toxicol.* 210: 69-84.
- 1001 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv:  
1002 1303.3997

- 1003 Li H 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34: 3094-3100.
- 1004 Li H 2011. A statistical framework for SNP calling, mutation discovery, association mapping and population  
1005 genetical parameter estimation from sequencing data. *Bioinformatics* 27: 2987-2993.
- 1006 Li H, Durbin R 2009. Fast and accurate short read alignment with Burrows–Wheeler transform.  
1007 *Bioinformatics* 25: 1754-1760.
- 1008 Limburg PA, Weider LJ 2002. ‘Ancient’DNA in the resting egg bank of a microcrustacean can serve as a  
1009 palaeolimnological database. *P Roy Soc B-Biol Sci* 269: 281-287.
- 1010 Llopart A, Herrig D, Brud E, Stecklein Z 2014. Sequential adaptive introgression of the mitochondrial  
1011 genome in *Drosophila yakuba* and *Drosophila santomea*. *Mol. Ecol.* 23: 1124-1136.
- 1012 Lomsadze A, Ter-Hovhannisyanyan V, Chernoff YO, Borodovsky M 2005. Gene identification in novel  
1013 eukaryotic genomes by self-training algorithm. *Nucleic Acids Res* 33: 6494-6506.
- 1014 Ma X, Hu W, Smilauer P, Yin M, Wolinska J 2019. *Daphnia galeata* and *D. dentifera* are geographically  
1015 and ecologically separated whereas their hybrids occur in intermediate habitats: A survey of 44 Chinese  
1016 lakes. *Mol. Ecol.* 28: 785-802.
- 1017 Marçais G, Kingsford C 2011. A fast, lock-free approach for efficient parallel counting of occurrences of k-  
1018 mers. *A fast, lock-free approach for efficient parallel counting of occurrences of k-mers* 27: 764-770.
- 1019 Marková S, Dufresne F, Manca M, Kotlík P. 2013. Mitochondrial capture misleads about ecological  
1020 speciation in the *Daphnia pulex* complex. *PLoS ONE* 8: e69497
- 1021 Martin SH, *et al.* 2013. Genome-wide evidence for speciation with gene flow in *Heliconius* butterflies.  
1022 *Genome Res.* 23: 1817-1828.
- 1023 Martin SH, Davey JW, Salazar C, Jiggins CD. 2019. Recombination rate variation shapes barriers to  
1024 introgression across butterfly genomes. *PLoS Biol.* 17: e2006288 doi: 10.1371/journal.pbio.2006288
- 1025 Martin SH, *et al.* 2020. Whole-chromosome hitchhiking driven by a male-killing endosymbiont. *Whole-  
1026 chromosome hitchhiking driven by a male-killing endosymbiont* 18: e3000610
- 1027 Martins Ribeiro M, *et al.* 2019. Mitogenome of *Daphnia laevis* (Cladocera, Daphniidae) from Brazil.  
1028 *Mitochondrial DNA Part B* 4: 194-196.
- 1029 McKenna A, *et al.* 2010. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-  
1030 generation DNA sequencing data. *Genome Res.* 20: 1297-1303. doi: 10.1101/gr.107524.110
- 1031 Meier JI, *et al.* 2017. Ancient hybridization fuels rapid cichlid fish adaptive radiations. *Nat. Commun.* 8: 1-  
1032 11.
- 1033 Melo-Ferreira J, *et al.* 2014. The elusive nature of adaptive mitochondrial DNA evolution of an arctic  
1034 lineage prone to frequent introgression. *Genome Biol. Evol* 6: 886-896.
- 1035 Meng G, Li Y, Yang C, Liu S. 2019. MitoZ: a toolkit for animal mitochondrial genome assembly, annotation  
1036 and visualization. *Nucleic Acids Res* 47: e63
- 1037 Miner BE, De Meester L, Pfrender ME, Lampert W, Hairston NG 2012. Linking genes to communities and  
1038 ecosystems: *Daphnia* as an ecogenomic model. *P Roy Soc B-Biol Sci* 279: 1873-1882. doi:  
1039 10.1098/rspb.2011.2404

- 1040 Möst MH 2013. Environmental change and its impact on hybridising *Daphnia* species complexes. [Doctoral  
1041 Thesis]. [Zürich]: ETH.
- 1042 Nadkarni MA, Martin FE, Jacques NA, Hunter N 2002. Determination of bacterial load by real-time PCR  
1043 using a broad-range (universal) probe and primers set. *Microbiology* 148: 257-266. doi: 10.1099/00221287-  
1044 148-1-257
- 1045 Nguyen L-T, Schmidt HA, Von Haeseler A, Minh BQ 2015. IQ-TREE: a fast and effective stochastic  
1046 algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32: 268-274.
- 1047 Orsini L, *et al.* 2013. The evolutionary time machine: using dormant propagules to forecast how populations  
1048 can adapt to changing environments. *Trends Ecol. Evol.* 28: 274-282. doi: 10.1016/j.tree.2013.01.009
- 1049 Patterson N, *et al.* 2012. Ancient admixture in human history. *Genetics* 192: 1065-1093.
- 1050 Petr M, Vernot B, Kelso J 2019. admixr—R package for reproducible analyses using ADMIXTOOLS.  
1051 *Bioinformatics* 35: 3194-3195.
- 1052 Petrussek A, *et al.* 2008a. A taxonomic reappraisal of the European *Daphnia longispina* complex (Crustacea,  
1053 Cladocera, Anomopoda). *Zool. Scr.* 37: 507-519. doi: 10.1111/j.1463-6409.2008.00336.x
- 1054 Petrussek A, Seda J, Machacek J, Ruthova S, Smilauer P 2008b. *Daphnia* hybridization along ecological  
1055 gradients in pelagic environments: the potential for the presence of hybrid zones in plankton. *Philos. Trans.*  
1056 *R. Soc. Lond. B Biol. Sci.* 363: 2931-2941. doi: 10.1098/rstb.2008.0026
- 1057 Petrussek A, Thielsch A, Schwenk K 2012. Mitochondrial sequence variation suggests extensive cryptic  
1058 diversity within the Western Palearctic *Daphnia longispina* complex. *Limnol Oceanogr* 57: 1838-1845. doi:  
1059 10.4319/lo.2012.57.6.1838
- 1060 Pietrzak B, Slusarczyk M 2006. The fate of the ephippia - *Daphnia* dispersal in time and space. *Pol. J. Ecol.*  
1061 54: 709-714.
- 1062 Pinsky ML, *et al.* 2021. Genomic stability through time despite decades of exploitation in cod on both sides  
1063 of the Atlantic. Genomic stability through time despite decades of exploitation in cod on both sides of the  
1064 Atlantic 118.
- 1065 Poplin R, *et al.* 2018. Scaling accurate genetic variant discovery to tens of thousands of samples. Scaling  
1066 accurate genetic variant discovery to tens of thousands of samples: 201178 doi: 10.1101/201178
- 1067 Rafajlović M, Emanuelsson A, Johannesson K, Butlin RK, Mehlig B 2016. A universal mechanism  
1068 generating clusters of differentiated loci during divergence-with-migration. *Evolution* 70: 1609-1621.
- 1069 Ranwez V, Douzery EJ, Cambon C, Chantret N, Delsuc F 2018. MACSE v2: toolkit for the alignment of  
1070 coding sequences accounting for frameshifts and stop codons. *Mol. Biol. Evol.* 35: 2582-2584.
- 1071 Reich D, Thangaraj K, Patterson N, Price AL, Singh L 2009. Reconstructing Indian population history.  
1072 *Nature* 461: 489-494.
- 1073 Rellstab C, Keller B, Girardclos S, Anselmetti FS, Spaak P 2011. Anthropogenic eutrophication shapes the  
1074 past and present taxonomic composition of hybridizing *Daphnia* in unproductive lakes. *Science* 56: 292-  
1075 302. doi: 10.4319/lo.2011.56.1.0292
- 1076 Revell LJ 2012. phytools: an R package for phylogenetic comparative biology (and other things). *Methods*  
1077 *Ecol. Evol.* 3: 217-223.

- 1078 Riesch R, *et al.* 2017. Transitions between phases of genomic differentiation during stick-insect speciation.  
1079 Nat. Ecol. Evol. 1: 82
- 1080 Runemark A, Eroukhmanoff F, Nava-Bolanos A, Hermansen JS, Meier JI 2018a. Hybridization, sex-  
1081 specific genomic architecture and local adaptation. Hybridization, sex-specific genomic architecture and  
1082 local adaptation 373. doi: 10.1098/Rstb.2017.0419
- 1083 Runemark A, *et al.* 2018b. Variation and constraints in hybrid genome formation. Nat. Ecol. Evol. 2: 549-  
1084 556. doi: 10.1038/s41559-017-0437-7
- 1085 Rusek J, *et al.* 2015. New possibilities arise for studies of hybridization: SNP-based markers for the multi-  
1086 species *Daphnia longispina* complex derived from transcriptome data. J. Plankton Res. 37: 626-635. doi:  
1087 10.1093/plankt/fbv028
- 1088 Sarver BA, *et al.* 2021. Diversification, Introgression, and Rampant Cytonuclear Discordance in Rocky  
1089 Mountains Chipmunks (Sciuridae: *Tamias*). Syst. Biol.
- 1090 Schreiber D, Pfenninger M 2020. Genomic divergence landscape in recurrently hybridizing *Chironomus*  
1091 sister taxa suggests stable steady state between mutual gene flow and isolation. Evol. Lett. 5: 86-100.
- 1092 Schwenk K 1993. Interspecific hybridization in *Daphnia*: distinction and origin of hybrid matriline. Mol.  
1093 Biol. Evol. 10: 1289-1302.
- 1094 Schwenk K, Bijl M, Menken SBJ 2001. Experimental interspecific hybridization in *Daphnia*. Hydrobiologia  
1095 442: 67-73. doi: 10.1023/A:1017594325506
- 1096 Schwenk K, Posada D, Hebert PDN 2000. Molecular systematics of European Hyalodaphnia: the role of  
1097 contemporary hybridization in ancient species. P Roy Soc B-Biol Sci 267: 1833-1842. doi:  
1098 10.1098/rspb.2000.1218
- 1099 Schwenk K, *et al.* 1998. Genetic markers, genealogies and biogeographic patterns in the cladocera. Aquat.  
1100 Ecol 32: 37-51.
- 1101 Sedlazeck FJ, Rescheneder P, Von Haeseler A 2013. NextGenMap: fast and accurate read mapping in highly  
1102 polymorphic genomes. Bioinformatics 29: 2790-2791.
- 1103 Seidendorf B, Boersma M, Schwenk K 2007. Evolutionary stoichiometry: The role of food quality for clonal  
1104 differentiation and hybrid maintenance in a *Daphnia* species complex. Limnol Oceanogr 52: 385-394.
- 1105 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM 2015. BUSCO: assessing genome  
1106 assembly and annotation completeness with single-copy orthologs. Bioinformatics 31: 3210-3212. doi:  
1107 10.1093/bioinformatics/btv351
- 1108 Skage M, *et al.* 2007. Intra-specific rDNA-ITS restriction site variation and an improved protocol to  
1109 distinguish species and hybrids in the *Daphnia longispina* complex. Hydrobiologia 594: 19-32. doi:  
1110 10.1007/s10750-007-9090-5
- 1111 Slager DL, *et al.* 2020. Cryptic and extensive hybridization between ancient lineages of American crows.  
1112 Mol. Ecol. 29: 956-969.
- 1113 Sloan DB, Havird JC, Sharbrough J 2017. The on-again, off-again relationship between mitochondrial  
1114 genomes and species boundaries. Mol. Ecol. 26: 2212-2236.
- 1115 Smit AFA, Hubley R. 2015. RepeatModeler Open. Version 1.0.

- 1116 Smit AFA, Hubley R, Green P. 2013-2015. RepeatMasker Open-4.0.
- 1117 Spaak P 2004. Spatial and temporal patterns of sexual reproduction in a hybrid *Daphnia* species complex.  
1118 Spatial and temporal patterns of sexual reproduction in a hybrid *Daphnia* species complex 26: 625-635. doi:  
1119 10.1093/plankt/fbh064
- 1120 Spaak P, Fox J, Hairston NG 2012. Modes and mechanisms of a *Daphnia* invasion. P Roy Soc B-Biol Sci  
1121 279: 2936-2944. doi: 10.1098/rspb.2012.0280
- 1122 Spaak P, Hoekstra JR 1997. Fish predation on a *Daphnia* hybrid species complex: A factor explaining  
1123 species coexistence? Limnol Oceanogr 42: 753-762. doi: 10.4319/lo.1997.42.4.0753
- 1124 Stanke M, Diekhans M, Baertsch R, Haussler D 2008. Using native and syntenically mapped cDNA  
1125 alignments to improve de novo gene finding. Bioinformatics 24: 637-644.
- 1126 Stephens JD, Rogers WL, Mason CM, Donovan LA, Malmberg RL 2015. Species tree estimation of diploid  
1127 *Helianthus* (Asteraceae) using target enrichment. Am. J. Bot. 102: 910-920.
- 1128 Taylor DJ, Hebert PDN, Colbourne JK 1996. Phylogenetics and evolution of the *Daphnia longispina* group  
1129 (Crustacea) based on 12S rDNA sequence and allozyme variation. Mol. Phylogenet. Evol. 5: 495-510. doi:  
1130 10.1006/mpev.1996.0045
- 1131 Team RC. 2017. R: A Language and Environment for Statistical Computing: R Foundation for Statistical  
1132 Computing, Vienna, Austria.
- 1133 Thielsch A, Brede N, Petrussek A, De Meester L, Schwenk K 2009. Contribution of cyclic parthenogenesis  
1134 and colonization history to population structure in *Daphnia*. Mol. Ecol 18: 1616-1628. doi: 10.1111/j.1365-  
1135 294X.2009.04130.x
- 1136 Thielsch A, Knell A, Mohammadyari A, Petrussek A, Schwenk K. 2017. Divergent clades or cryptic species?  
1137 Mito-nuclear discordance in a *Daphnia* species complex. BMC Evol. Biol. 17: 227
- 1138 Thielsch A, Volker E, Kraus RHS, Schwenk K 2012. Discrimination of hybrid classes using cross-species  
1139 amplification of microsatellite loci: methodological challenges and solutions in *Daphnia*. Mol. Ecol. Resour.  
1140 12: 697-705. doi: 10.1111/j.1755-0998.2012.03142.x
- 1141 Tollrian R, Harvell CD. 1999. The ecology and evolution of inducible defenses: Princeton University Press.
- 1142 Trentini M 1980. Chromosome numbers of nine species of Daphniidae (Crustacea, Cladocera). Genetica  
1143 54: 221-223.
- 1144 Vaidya G, Lohman DJ, Meier R 2011. SequenceMatrix: concatenation software for the fast assembly of  
1145 multi-gene datasets with character set and codon information. Cladistics 27: 171-180.
- 1146 Van der Auwera GA, *et al.* 2013. From FastQ data to high-confidence variant calls: the genome analysis  
1147 toolkit best practices pipeline. Curr. Protoc. Bioinformatics 43: 11-10.
- 1148 Vergilino R, Markova S, Ventura M, Manca M, Dufresne F 2011. Reticulate evolution of the *Daphnia pulex*  
1149 complex as revealed by nuclear markers. Mol. Ecol. 20: 1191-1207.
- 1150 Waldvogel AM, *et al.* 2018. The genomic footprint of climate adaptation in *Chironomus riparius*. Mol.  
1151 Ecol. 27: 1439-1456.
- 1152 Walker BJ, *et al.* 2014. Pilon: an integrated tool for comprehensive microbial variant detection and genome  
1153 assembly improvement. PLoS ONE 9: e112963



- 1154 Wingett SW, Andrews S 2018. FastQ Screen: A tool for multi-genome mapping and quality control.  
1155 F1000Research 7.
- 1156 Wolf HG. 1987. Interspecific hybridization between *Daphnia hyalina*, *D. galeata*, and *D. cucullata* and  
1157 seasonal abundances of these species and their hybrids. In: Forró L., D.G. F, editors. Cladocera. Dordrecht:  
1158 Springer. p. 213-217.
- 1159 Wolf HG, Mort MA 1986. Interspecific hybridization underlies phenotypic variability in *Daphnia*  
1160 populations. *Oecologia* 68: 507-511.
- 1161 Ye ZQ, *et al.* 2017. A New Reference Genome Assembly for the Microcrustacean *Daphnia pulex*. *G3* 7:  
1162 1405-1416. doi: 10.1534/g3.116.038638
- 1163 Yeaman S, Whitlock MC 2011. The genetic architecture of adaptation under migration–selection balance.  
1164 *Evolution* 65: 1897-1911.
- 1165 Yin MB, Giessler S, Griebel J, Wolinska J 2014. Hybridizing *Daphnia* communities from ten neighbouring  
1166 lakes: spatio-temporal dynamics, local processes, gene flow and invasiveness. *Hybridizing Daphnia*  
1167 communities from ten neighbouring lakes: spatio-temporal dynamics, local processes, gene flow and  
1168 invasiveness 14. doi: 10.1186/1471-2148-14-80
- 1169 Zehnder A, Gorham PR 1960. Factors influencing the growth of *Microcystis aeruginosa* Kütz. Emend.  
1170 Elenkin. *Can. J. Microbiol* 6: 645-660.
- 1171 Zheng X, *et al.* 2012. A high-performance computing toolset for relatedness and principal component  
1172 analysis of SNP data. *Bioinformatics* 28: 3326-3328.
- 1173