

1 **DISCOMARK: Nuclear marker discovery from orthologous**  
2 **sequences using low coverage genome data**

3 Harald Detering\*<sup>1,2,3</sup>, Sereina Rutschmann\*<sup>1,2,3</sup>, Sabrina Simon<sup>4,5</sup>, Jakob Fredslund<sup>6</sup>, Michael  
4 T. Monaghan<sup>1,2</sup>

5 **Addresses:**

6 <sup>1</sup>*Leibniz-Institute of Freshwater Ecology and Inland Fisheries (IGB), Müggelseedamm 301,*  
7 *12587 Berlin, Germany*

8 <sup>2</sup>*Berlin Center for Genomics in Biodiversity Research, Königin-Luise-Straße 6-8, 14195*  
9 *Berlin, Germany*

10 <sup>3</sup>*Department of Biochemistry, Genetics and Immunology, University of Vigo, 36310 Vigo,*  
11 *Spain*

12 <sup>4</sup>*Sackler Institute for Comparative Genomics, American Museum of Natural History, Central*  
13 *Park West and 79<sup>th</sup> St., New York, NY 10024, USA*

14 <sup>5</sup>*Biosystematics Group, Wageningen University, Droevendaalsesteeg 1, 6708 PB*  
15 *Wageningen, The Netherlands*

16 <sup>6</sup>*Alexandra Institute, Åbogade 34, 8200 Aarhus, Denmark*

17 **Keywords:**

18 marker discovery, mayfly, non-model organism, primer design, phylogenetics

19 **Correspondence:**

20 Sereina Rutschmann, Phylogenomics Lab, Department of Biochemistry, Genetics and  
21 Immunology, University of Vigo, 36310 Vigo, Spain. E-Mail:  
22 [sereina.rutschmann@gmail.com](mailto:sereina.rutschmann@gmail.com)

23 \* these authors contributed equally

24 **Running title:**

25 DISCOMARK - Phylogenetic marker development

26 **Abstract**

27 High-throughput sequencing has laid the foundation for fast and cost-effective development  
28 of phylogenetic markers. Here we present the program DISCOMARK, which streamlines the  
29 development of nuclear DNA (nDNA) markers from whole-genome (or whole-transcriptome)  
30 sequencing data, combining local alignment, alignment trimming, reference mapping and  
31 primer design based on multiple sequence alignments in order to design primer pairs from  
32 input orthologous sequences. In order to demonstrate the suitability of DISCOMARK we  
33 designed markers for two groups of species, one consisting of closely related species and one  
34 group of distantly related species. For the closely related members of the species complex of  
35 *Cloeon dipterum* s.l. (Insecta, Ephemeroptera), the program discovered a total of 77 markers.  
36 Among these, we randomly selected eight markers for amplification and Sanger sequencing.  
37 The exon sequence alignments (2,526 base pairs (bp)) were used to reconstruct a well  
38 supported phylogeny and to infer clearly structured haplotype networks. For the distantly  
39 related species we designed primers for several families in the insect order Ephemeroptera,  
40 using available genomic data from four sequenced species. We developed primer pairs for 23  
41 markers that are designed to amplify across several families. The DISCOMARK program will  
42 enhance the development of new nDNA markers by providing a streamlined, automated  
43 approach to perform genome-scale scans for phylogenetic markers. The program is written in  
44 Python, released under a public license (GNU GPL v2), and together with a manual and  
45 example data set available at: <https://github.com/hdetering/discomark>.

## 46 **Introduction**

47 The inference of phylogenetic relationships has benefited profoundly from the availability of  
48 nuclear DNA (nDNA) sequences for an increasing number of organism groups. The  
49 development of new phylogenetic markers has provided unprecedented insight into the  
50 evolutionary relationships of non-model organisms in particular (Ellegren 2014). Large sets of  
51 nDNA markers (single copy genes) have recently been designed for taxonomic groups for  
52 which genomic resources were available, e.g. cichlid fish (Meyer *et al.* 2015), ray-finned fish  
53 (Near *et al.* 2012), reptiles (Ruane *et al.* 2014), birds (Kerr *et al.* 2014) and flowering plants  
54 (Zeng *et al.* 2014). However, for many taxonomic groups there are only a handful of nDNA  
55 markers available that are suitable for phylogenetic reconstruction. Other approaches, such as  
56 ultra-conserved element (UCE) sequencing (Faircloth *et al.* 2012), anchored hybrid  
57 enrichment (Lemmon and Lemmon 2012), restriction site-associated DNA (RAD) sequencing  
58 (Baird *et al.* 2008) or genotyping by sequencing (GBS, Elshire *et al.* 2011) have become  
59 popular for addressing specific questions in systematics or population genetics; however,  
60 these methods are still cost-intensive, require a comparatively high amount of starting DNA  
61 material and can depend on the availability of reference genomes (e.g. anchored hybrid  
62 enrichment). Consequently, standard Sanger sequencing approaches are still in high demand  
63 for various research questions.

64 Identification of novel phylogenetic markers has been a predominantly manual process,  
65 which impedes their large-scale development, and comprehensive primer design based on  
66 large sets of multiple sequence alignments remains challenging. Recently, tools have been  
67 developed for (1) specific primer design such as for automated primer design from  
68 transcriptome data (SCRIMER, Morkovsky *et al.* 2015), for individual degenerate primers  
69 (GEMI, Sobhy *et al.* 2012; PRIMER3, Untergasser *et al.* 2012; CEMASUITE, Lane *et al.* 2015),

70 for highly variable DNA targets (PRIMERDESIGN, Brodin *et al.* 2013; PRIMERDESIGN-M,  
71 Yoon and Leitner 2015), viral genomes (PRISM, Yu *et al.* 2015), multiple primer design  
72 (BATCHPRIMER3, You *et al.* 2008; PRIMERVIEW, O’Halloran 2015) and (2) the discovery of  
73 specific markers, including single nucleotide polymorphism (SNP) markers (POLYMARKER,  
74 Ramirez-Gonzalez *et al.* 2015), and putative single copy nuclear loci (MARKERMINER,  
75 Chamala *et al.* 2015). In addition, the challenge of developing new markers lies both in the  
76 discovery of conserved regions, the design of primer pairs and an estimation of their  
77 suitability as phylogenetic markers.

78 Our aim was to develop a flexible, user-friendly program that works with FASTA-  
79 formatted files of putative orthologous sequences from whole-genome or whole-transcriptome  
80 data, identified conserved regions and designs primers based on these multiple sequence  
81 alignments. Here we present DISCOMARK (=Discovery of Markers), a program for the  
82 discovery of phylogenetically suitable nDNA markers and design of primer pairs. The  
83 program can be used to easily screen for nDNA markers and design primers that can be used  
84 for Sanger sequencing as well as high-throughput sequencing. The program is structured into  
85 several steps that can be individually optimized by the user and run independently. In terms of  
86 input the program can be applied on large and small sets of taxa, including both closely and  
87 distantly related species. Ideally, orthologous sequences in combination with a whole-genome  
88 reference sequence are used. Thus, exon/intron boundaries can be inferred using the reference  
89 for each marker. Under the default settings, the program will design several primer pairs that  
90 anneal in conserved regions. The visualization of the alignments with potential primers allows  
91 the user to choose between primers targeting exons or introns (e.g. exon-primed intron-  
92 crossing (EPIC) markers). Additionally, information about the suitability as phylogenetic  
93 markers is provided by an estimate of the number of SNPs per marker and the applicability

94 across species. Finally, we demonstrate the utility of DISCOMARK for (1) closely related  
95 species (i.e. *Cloeon dipterum* s.l. species complex) using whole-genome data, and (2)  
96 distantly related species (i.e. insect order Ephemeroptera) using whole-genome data derived  
97 from genome sequencing projects.

## 98 **Materials and Methods**

### 99 *DISCOMARK implementation*

100 The program DISCOMARK is written in Python and is developed to design primer pairs in  
101 conserved regions of predicted orthologous genes. Orthologs are most suited for phylogenetic  
102 studies. The ortholog identification step is not part of the DISCOMARK workflow but  
103 DISCOMARK is designed to directly work with the output of several ortholog prediction  
104 programs, e.g. HAMSTR (Ebersberger *et al.* 2009), or Orthograph  
105 (<https://github.com/mptksen/Orthograph>, last accessed March 25, 2016). Orthologous groups  
106 may be derived from genomic or transcriptomic sequencing data. In addition to the  
107 orthologous genes, genomic data such as whole-genome sequencing data can be provided to  
108 DISCOMARK as a guide to detect exon/intron boundaries. DISCOMARK performs seven steps,  
109 combining Python scripts with widely used bioinformatics programs (Fig. 1). The steps: (1)  
110 combine orthologous groups of sequences, (2) align sequences of each orthologous group  
111 using MAFFT v.7.205 (Katoh and Standley 2013), (3) trim sequence alignments with TRIMAL  
112 v.1.4 (Capella-Gutierrez *et al.* 2009), (4) align sequences against a reference (e.g. whole-  
113 genome dataset from the same or closely related taxa) with BLASTN v.2.2.29 (Altschul *et al.*  
114 1997; Camacho *et al.* 2009) and re-alignment using MAFFT, (5) design primer pairs on  
115 single-gene alignments using a modified version of PRIFi (Fredslund *et al.* 2005), adapted by  
116 us into a Python package that uses BioPython v.1.65 and Python v.3.4.3, (6) check primer

117 specificity with BLASTN, and (7) generate output in several formats (visual HTML report,  
118 tabular data and FASTA files of the primers). The results of each step can be inspected in the  
119 respective output folders.

120 *1 Combine sequences.* In the first step, the putative orthologous sequences of different taxa  
121 are combined according to the orthologous groups. The input files are expected to be  
122 nucleotide sequences in FASTA format. We recommend using putative orthologous exon  
123 sequences (e.g. CDS) in combination with whole-genome data (e.g. a draft genome  
124 assembly). Each input file is expected to contain the sequences of one orthologous group;  
125 orthologs of each input taxon are to be organized into a taxon folder. Importantly, file names  
126 represent the ortholog identifiers used to combine orthologous sequences of the various input  
127 taxa; by default, ortholog prediction tools follow that convention.

128 *2 Align sequences.* Orthologous sequences combined according to the orthologous groups are  
129 separately aligned with the multiple sequence alignment (MSA) program MAFFT. Alignment  
130 parameters can be specified by the user via a configuration file (discomark.conf, located in the  
131 program folder). Default parameters are the following: ‘--localpair --maxiterate 16 --  
132 inputorder --preserve-case --quiet’ (L-INS-i alignment method). We chose MAFFT as multiple  
133 alignment tool because it combines accuracy and efficiency and has been adopted widely in  
134 the scientific community (Pais *et al.* 2014; Szitenberg *et al.* 2015).

135 *3 Trim alignments.* In order to remove poorly aligned regions, sequence alignments are  
136 trimmed using TRIMAL. The program TRIMAL analyzes the distribution of gaps and  
137 mismatches in the alignment and discard alignment positions and sequences of low quality.  
138 By default, DISCOMARK calls TRIMAL with the ‘-strictplus’ method. The preset is used by  
139 TRIMAL to derive the specific thresholds for alignment trimming (minimum gap score,

140 minimum residue similarity score, conserved block size). Since alignment trimming largely  
141 depends on the input data and influences the downstream results, TRIMAL can also be run with  
142 different settings (e.g. ‘-gappyout’, ‘-strict’, ‘-automated1’; but see Capella-Gutierrez *et al.*  
143 (2009). Alternatively, there is also the option to deactivate the alignment trimming with the  
144 DISCOMARK option ‘--no-trim’ or use alternative trimming programs such as GBLOCKS  
145 (Castresana 2000; Talavera and Castresana 2007) or GUIDANCE2 (Landan and Graur 2008;  
146 Sela *et al.* 2015).

147 *4 Blast and alignment to reference.* In this step a genomic reference sequence for each input  
148 ortholog is identified and added to the trimmed alignment. This step is particularly important  
149 when working with coding sequences which do not contain intron sequences; thus, a genomic  
150 sequence is needed to infer intron/exon boundaries. Working with coding sequences is  
151 advisable for more distantly related taxa which may include intron length polymorphisms, or  
152 to target EPIC markers. Any whole-genome data set (from one of the included taxa or a  
153 closely related taxa) can be used as reference for mapping the ortholog sequences. Here,  
154 mapping means that the input sequences are compared to the reference sequences, which are  
155 defined by the user using the local alignment program BLASTN. The best locally aligning  
156 reference sequence (the one that yields the longest alignment among all input sequences) for  
157 each orthologous group is added to the corresponding sequence alignment. Reference  
158 sequences are cut to 100 base pairs (bp) upstream and downstream of the first, respectively  
159 last, BLAST hit to avoid alignment length inflation. Then, the extended alignments are re-  
160 aligned with MAFFT. The reference alignment step is optional; however, the inclusion of  
161 whole-genome data is essential for estimating intron/exon boundaries. Given that information,  
162 the focus of target sequences to be amplified can be on entire exon markers, EPIC markers, or  
163 a combination.

164 *5 Design primers.* The single-gene alignments, after trimming, mapping and re-aligning to a  
165 reference, are used as input to design primer pairs. We integrated the webtool PRiFi  
166 (<http://cgi-www.daimi.au.dk/cgi-chili/PriFi/main>, last accessed December 20, 2015) as a  
167 Python package that provides a comprehensive set of parameters. As default settings for  
168 DISCOMARK we chose the following: estimated product length between 200-1,000 bp  
169 ('OptimalProductLength = [400, 600, 800, 1000], MinProductLength = 200,  
170 MaxProductLength = 1000'), maximum number of ambiguity positions within the primer  
171 sequences ('MaxMismatches = 2'), primer length between 20-30 bp ('MinPrimerLength = 20,  
172 MaxPrimerLength = 30, OptimalPrimerLength = [20, 25]'), melting temperature of the primer  
173 pairs between 50-60°C ('MinTm = 50.0, MinTmWithMismatchesAllowed = 58.0,  
174 SuggestedMaxTm = 60.0'), and we set the maximum number of primer pairs per alignment to  
175 six (note: only settings different from the PRiFi default are mentioned above). The program  
176 PRiFi was originally developed to design intron-spanning markers (but see Fredslund et al.  
177 2005). Here we use it because it enables primer design based on MSA input. Parameters for  
178 PRiFi can be specified in the DISCOMARK configuration file ('discomark.conf').

179 *6 Check marker specificity.* To ensure the specificity of the designed primer pairs, we  
180 compare their sequences against the NCBI database ('refseq\_mrna'). Primer sequences are  
181 searched in the NCBI database ('refseq\_mrna') using the online BLASTN interface. The  
182 default search settings are restricted to human and bacterial targets using the Entrez query  
183 'txid2[ORGN] OR txid9606[ORGN]' because these are most likely to be present as  
184 contaminants in sequencing libraries. The result hits of the BLAST search are indicated to the  
185 user in the HTML output.

186 *7 Visualize results.* As final step, the program produces a HTML report containing the list of

187 designed primers, an alignment viewer and plots visualizing the discovered set of markers.  
188 Besides the primer sequences the report lists several features such as the melting  
189 temperatures, predicted sequence length, and the number of taxa amplified by each primer set.  
190 Selected primer pairs and primer lists can be downloaded as FASTA or CSV files,  
191 respectively. In order to provide a measure of the suitability of the markers for phylogenetic  
192 reconstruction the program calculates the number of SNPs between a primer pair by  
193 comparing the aligned input sequences against each other. The number of SNPs between each  
194 primer pair is visualized in relation to the estimated product length (see Fig. 2 for an example)  
195 and reported in the tabular output. Furthermore, the report highlights the species coverage  
196 achieved by the discovered markers, i.e. how many species' sequences each primer set is  
197 expected to amplify, as an estimate of how universal each primer set can be applied.  
198 Additionally, functional annotations are reported, if available, to guide the user in the  
199 selection of markers of interest. Annotations can be supplied in form of a tab-delimited file  
200 with the '-a' option. In principle, any kind of annotations can be used depending on the  
201 desired research objective. In our usage scenarios, we used gene ontology (GO) terms which  
202 were retrieved by mapping the gene IDs contained in the HAMSTR core ortholog set via the  
203 UniProt website (<http://www.uniprot.org/>, last accessed December 20, 2015).

#### 204 *Usage cases*

205 *Closely related species - Cloeon dipterum s.l. species complex.* To test the suitability of  
206 DISCOMARK for closely related species, we designed primer pairs for the species complex of  
207 *Cloeon dipterum* s.l. (Ephemeroptera: Baetidae). The species complex consists of several  
208 closely related species, including *Cloeon peregrinator* GATTOLLIAT & SARTORI, 2008 from  
209 Madeira (Gattolliat *et al.* 2008; Rutschmann *et al.* 2014; Table 1). As input to design the

210 primer pairs data we used whole-genome sequencing data of *Cloeon dipterum* L. 1761  
211 (Baetidae; Sequence Read Archive SRP050093) and expressed sequence tags (EST) of *Baetis*  
212 sp. (Baetidae; FN198828-FN203024). The sequence reads of *C. dipterum* were trimmed and  
213 *de novo* assembled using NEWBLER v.2.5.3 (454 Life Science Corporation) under the default  
214 settings for large datasets. Ortholog sequences prediction of both data sets was performed  
215 with HAMSTR v.9 using the insecta\_hmmer3-2 core reference taxa set ([http://www.deep-](http://www.deep-phylogeny.org/hamstr/download/datasets/hmmer3/insecta_hmmer3-2.tar.gz)  
216 [phylogeny.org/hamstr/download/datasets/hmmer3/insecta\\_hmmer3-2.tar.gz](http://www.deep-phylogeny.org/hamstr/download/datasets/hmmer3/insecta_hmmer3-2.tar.gz), last accessed  
217 December 20, 2015), including 1,579 orthologous genes. We ran the program DISCOMARK  
218 with default settings ('python run\_project.py -i input/Cloeon -i input/Baetis -r  
219 input/reference/Cloeon.fa -a input/co2go.ixosc.csv -d output/cloeon\_baetis'), using the  
220 predicted orthologs from HAMSTR and the whole-genome *Cloeon*-data as reference (step 4).  
221 The Pearson correlation between the number of SNPs between primer pairs and  
222 corresponding estimated product length was calculated using the function cor within the stats  
223 package for R (R Development Core Team, 2016). A t-test for significance was performed  
224 using the function cor.test.

225 From the total of designed primer pairs (77 markers, 338 primer pairs, see results) we  
226 selected eight and amplified them for four species of the *C. dipterum* species complex (Table  
227 1) in the laboratory. We used standardized polymerase chain reactions (PCR; 35-40 PCR  
228 cycles with annealing temperature of 55°C), followed by Sanger sequencing. Forward and  
229 reverse sequences were assembled and edited with GENEIOUS R7 v.7.1.3 (Biomatters Ltd.),  
230 indicating ambiguous positions following the IUPAC nucleotide codes. Heterozygous  
231 sequences were decoded with CODONCODEALIGNER v.3.5.6 (CodonCode Corporation) using  
232 the find and split heterozygous function. Multiple sequence alignments were created for all

233 sequences per marker. The predicted orthologous sequences of *Baetis* sp. were used as  
234 reference to infer the exon-intron splicing boundaries (canonical and non-canonical splice site  
235 pairs). The final sequence alignments were checked for the occurrence of stop codons and  
236 indels, and split into exon and intron parts using a custom Python script  
237 ([https://github.com/srutschmann/python\\_scripts](https://github.com/srutschmann/python_scripts), last accessed March 28, 2016). Sequence  
238 alignments were phased using the program PHASE v.2.1.1 (Stephens *et al.* 2001; Stephens  
239 and Donnelly 2003) with a cutoff value of 0.6 (Harrigan *et al.* 2008; Garrick *et al.* 2010),  
240 whereby input and output files were formatted using the Perl scripts included in SEQPHASE  
241 (Flot 2010). Heterozygous sites that could not be resolved were coded as ambiguity codes for  
242 subsequent analyses. After phasing, all alignments were re-aligned with MAFFT. The number  
243 of variable and informative sites, and the nucleotide diversity per exon alignment was  
244 calculated with a custom script.

245 To investigate the heterogeneity of each marker's DNA sequences, we reconstructed  
246 haplotype networks, using FITCHI (Matschiner 2015). As input for each marker we inferred a  
247 gene tree using the program RAXML v.8 (Stamatakis 2014) with the GTRCAT model and  
248 1,000 bootstrap replicates under the rapid bootstrap algorithm. The phylogenetic relationships  
249 were calculated with Bayesian inference, using MRBAYES v.3.2.3 (Ronquist *et al.* 2012)  
250 based on a concatenated nDNA matrix that consisted of the exon sequences from all 15  
251 nDNA markers. The best-fitting model of molecular evolution for each sequence alignment  
252 was selected via a BIC criterion in JMODELTEST v.2.1 (Guindon and Gascuel 2003; Darriba *et*  
253 *al.* 2012). We calculated  $10^6$  generations with random seed, a burn-in of 25% and four  
254 MCMC chains. As an outgroup we used the predicted orthologous sequences of *Baetis* sp..

255 *Distantly related species - insect order Ephemeroptera*. In this test case, we used contigs  
256 derived from whole-genome sequencing projects of the species *Baetis* sp. (Baetidae;  
257 BioProject PRJNA219528), *Ephemera danica* MÜLLER 1764 (Ephemeridae; BioProject  
258 PRJNA219552), *Eurylophella* sp. (Ephemerellidae; BioProject PRJNA219556), and  
259 *Isonychia bicolor* WALKER 1853 (Isonychiidae; BioProject PRJNA219568). The contigs from  
260 each species were used for ortholog predicting with HAMSTR v.13.2.4  
261 (<http://sourceforge.net/projects/hamstr/files/hamstr.v13.2.4.tar.gz>, last accessed December 20,  
262 2015). We ran DISCOMARK with the default settings, using the *Baetis* sp. data as reference  
263 ('python run\_project.py -i input/Baetis -i input/Ephemera -i input/Eurylophella -i  
264 input/Isonychia -r input/references/Baetis.fa -a input/co2go.ixosc.csv -d output/mayflies').

## 265 **Results**

266 *Closely related species - species complex of Cloeon dipterum s.l.*

267 DISCOMARK identified a total of 804 nDNA markers and 77 alignments with 338 primer pairs  
268 for orthologous sequences of both species (*Baetis* sp. and *C. dipterum* s.l.). Ortholog  
269 prediction yielded 403 orthologous sequences for the *Baetis* sp. EST-data and 1,211 for *C.*  
270 *dipterum*. For the individual species, DISCOMARK identified 790 markers for *C. dipterum* and  
271 123 for *Baetis* sp. The lengths of the markers including both species were between 201 and  
272 925 bp with median length of 451.5 bp. The number of SNPs per marker ranged from zero to  
273 37 (median: 5) with an average of one SNP per 68 bp. Marker length and number of SNPs  
274 were correlated with a Pearson's correlation coefficient of 0.35 (Pearson's product-moment  
275 correlation  $P < 0.001$ ). The total run time for this data set on a local Linux machine (quad-  
276 core Intel i5, 8 GB RAM) was 24 min.

277 The haplotype networks based on the eight selected markers showed a clear structure for  
278 all markers, including two markers with shared haplotypes for the two species from the U.S.  
279 and Madeira (Fig. 3 and Fig. S1, Supporting information). The length of the concatenated  
280 sequence alignment of the eight markers was 3,530 bp (2,526 bp exon sequence, Table S1,  
281 Supporting information). The exon sequence matrix contained 78 variable sites, 27  
282 informative sites, and was 92.6% complete. The nucleotide diversity ranged between 0.009  
283 and 0.028 (median: 0.013). Phylogenetic tree reconstruction based on these eight markers  
284 resulted in a phylogeny with fully resolved nodes (Bayesian posterior probability (PP)  $\geq$  95%;  
285 Fig. 3). The species *C. dipterum* sp1 was found as outgroup to a clade containing the species  
286 *C. dipterum* sp2 from Switzerland and the two species from the U.S and Madeira. The latter  
287 two species formed a monophyletic clade.

#### 288 *Distantly related species - insect order Ephemeroptera*

289 In total, we found 22 orthologs with a total of 48 primer pairs for all four species (Table S2,  
290 Supporting information). The input files per species (i.e. putative orthologous sequences)  
291 ranged from 1,445 to 1,523. We detected 41 markers that covered three of the species (99  
292 primer pairs), 81 markers covering two species (210 primer pairs), and 117 markers that  
293 covered any single species (478 primer pairs). For the individual species, *Baetis* sp. had the  
294 most markers available (214) of the single- and multi-species markers. There were 138  
295 markers for *Eurylophella* sp., 107 markers for *I. bicolor*, and 88 markers for *E. danica*. The  
296 lengths for all markers covering all four species varied between 216 and 997 bp with median  
297 of 398.5 bp, containing between 39 and 298 SNPs per marker (Fig. 2,) with a SNP every 4.1  
298 bp on average. Marker length and number of SNPs were correlated with a Pearson's  
299 correlation coefficient of 0.97 (Pearson's product-moment correlation  $P < 0.001$ ). Run time

300 for this data set on a Linux client (quad-core Intel i5, 8 GB RAM) was 46 min.

## 301 **Discussion**

302 To our knowledge, the program DISCOMARK is the first stand-alone program with the aim of  
303 designing primer pairs based on multiple sequence alignments on a genome-wide scale. The  
304 visual output gives guidance on the suitability of each marker (i.e. variability within and  
305 between species measured as number of SNPs, and information about the included species of  
306 each marker. Using this approach, primers can be specifically chosen to match the  
307 ‘phylogenetic scale’ (i.e. for closely related species many markers with intermediate number  
308 of SNPs and for distantly related species fewer markers with generally higher number of  
309 SNPs can be selected. The automatic processing, including combining, aligning, trimming  
310 and blasting sequences of any nucleotide FASTA sequences together with the produced  
311 graphical output significantly facilitate the design of primer pairs for a large number of nDNA  
312 markers. Nevertheless, users retain a high degree of flexibility by the stepwise nature of the  
313 workflow. DISCOMARK is free, open-source software to assist the development of markers for  
314 non-model species on the genome scale. We demonstrated the efficacy of our approach for  
315 closely related species as well as for members of divergent families within an order of insects.  
316 Using a reference genome enabled resolution of intron-exon boundaries but is not a strict  
317 requirement for marker design.

### 318 *Markers development within the order Ephemeroptera*

319 The usage of DISCOMARK adds an extensive set of new potential nDNA markers to the ones  
320 that have been used to date for mayfly phylogenies based on individual genes (histone 3,  
321 elongation factor 1 alpha, phosphoenolpyruvate carboxykinase (Vuataz *et al.* 2011; Pereira-da-

322 Conceicoa *et al.* 2012; Vuataz *et al.* 2013). Most recent phylogenetic reconstructions are still  
323 mostly based on the information of mitochondrial DNA markers (e.g. Rutschmann *et al.*  
324 2014; Macher *et al.* 2016). The availability of more genome data will be very valuable in  
325 order to increase the number of markers suitable for phylogenetic studies. The use of the  
326 larger marker set for *C. dipterum* developed here resulted in a fully resolved phylogenetic tree  
327 in contrast to Rutschmann *et al.* (2014). The availability of more markers promote fine-scaled  
328 phylogenetic studies, which are needed to resolve the phylogenetic relationships of so-called  
329 morphologically cryptic species that can not be resolved with standard markers (Dijkstra *et al.*  
330 2014).

### 331 **Acknowledgements**

332 We are thankful to our research groups, in particular the Phylogenomics Lab at the University  
333 of Vigo for constructive discussion that improved this project. Research was partially  
334 supported by the Leibniz Association (PAKT für Forschung und Innovation “FREDIE”  
335 project) and the Swiss National Science Foundation (Early PostDoc.Mobility grant  
336 P2SKP3\_15869 to S.R.). This is publication number ### of the Berlin Center for Genomics in  
337 Biodiversity Research.

### 338 **References**

- 339 Altschul SF, Madden TL, Schaffer AA *et al.* (1997) Gapped BLAST and PSI-BLAST: a new  
340 generation of protein database search programs. *Nucleic Acids Research*, **25**, 3389-3402.  
341 Baird NA, Etter PD, Atwood TS *et al.* (2008) Rapid SNP Discovery and Genetic Mapping  
342 Using Sequenced RAD Markers. *PLoS ONE*, **3**, e3376.  
343 Brodin J, Krishnamoorthy M, Athreya G *et al.* (2013) A multiple-alignment based primer  
344 design algorithm for genetically highly variable DNA targets. *BMC Bioinformatics*, **14**,  
345 255.  
346 Camacho C, Coulouris G, Avagyan V *et al.* (2009) BLAST+: architecture and applications.  
347 *BMC Bioinformatics*, **10**, 421.  
348 Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T (2009) trimAl: a tool for automated  
349 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*, **25**, 1972-1973.

- 350 Castresana J (2000) Selection of conserved blocks from multiple alignments for their use in  
351 phylogenetic analysis. *Molecular Biology and Evolution*, **17**, 540-552.
- 352 Chamala S, García N, Godden GT *et al.* (2015) MarkerMiner 1.0: new application for  
353 phylogenetic marker development using angiosperm transcriptomes. *Applications in Plant*  
354 *Sciences*, **4**, 1400115.
- 355 Dijkstra KD, Monaghan MT, Pauls SU (2014) Freshwater biodiversity and aquatic insect  
356 diversification. *Annual Review of Entomology*, **59**, 143-163.
- 357 Darriba D, Taboada GL, Doallo R, Posada D (2012) jModelTest 2: more models, new  
358 heuristics and parallel computing. *Nature Methods*, **9**, 772.
- 359 Ebersberger I, Strauss S, von Haeseler A (2009) HaMStR: profile hidden markov model  
360 based search for orthologs in ESTs. *BMC Evolutionary Biology*, **9**, 157.
- 361 Ellegren H (2014) Genome sequencing and population genomics in non-model organisms.  
362 *Trends in Ecology & Evolution*, **29**, 51-63.
- 363 Elshire RJ, Glaubitz JC, Sun Q *et al.* (2011) A robust, simple genotyping-by-sequencing  
364 (GBS) approach for high diversity species. *PLoS ONE*, **6**, e19379.
- 365 Faircloth BC, McCormack JE, Crawford NG *et al.* (2012) Ultraconserved elements anchor  
366 thousands of genetic markers spanning multiple evolutionary timescales. *Systematic*  
367 *Biology*, **61**, 717-726.
- 368 Flot J-F (2010) seqphase: a web tool for interconverting phase input/output files and fasta  
369 sequence alignments. *Molecular Ecology Resources*, **10**, 162-166.
- 370 Fredslund J, Schauser L, Madsen LH, Sandal N, Stougaard J (2005) PriFi: using a multiple  
371 alignment of related sequences to find primers for amplification of homologs. *Nucleic*  
372 *Acids Research*, **33**, W516-520.
- 373 Garrick RC, Sunnucks P, Dyer RJ (2010) Nuclear gene phylogeography using PHASE:  
374 dealing with unresolved genotypes, lost alleles, and systematic bias in parameter  
375 estimation. *BMC Evolutionary Biology*, **10**, 118.
- 376 Gattolliat J-L, Hugher SJ, Monaghan MT, Sartori M (2008) Revision of Mdeiran mayflies  
377 (Insecta, Ephemeroptera). *Zootaxa*, **1957**, 69-80.
- 378 Guindon S, Gascuel O (2003) A simple, fast, and accurate algorithm to estimate large  
379 phylogenies by maximum likelihood. *Systematic Biology*, **52**, 696-704.
- 380 Harrigan RJ, Mazza ME, Sorenson MD (2008) Computation vs. cloning: evaluation of two  
381 methods for haplotype determination. *Molecular Ecology Resources*, **8**, 1239-1248.
- 382 Katoh K, Standley DM (2013) MAFFT multiple sequence alignment software version 7:  
383 improvements in performance and usability. *Molecular Biology and Evolution*, **30**, 772-  
384 780.
- 385 Kerr KCR, Cloutier A, Baker AJ (2014) One hundred new universal exonic markers for birds  
386 developed from a genomic pipeline. *Journal of Ornithology*, **155**, 561-569.
- 387 Landan G, Graur D (2008) Local reliability measures from sets of co-optimal multiple  
388 sequence alignments. *Pacific Symposium on Biocomputing*, 15-24.
- 389 Lane CE, Hulgán D, O'Quinn K, Benton MG (2015) CEMAsuite: open source degenerate  
390 PCR primer design. *Bioinformatics*, **31**, 3688-3690.
- 391 Lemmon AR, Lemmon EM (2012) High-throughput identification of informative nuclear loci  
392 for shallow-scale phylogenetics and phylogeography. *Systematic Biology*, **61**, 745-761.
- 393 Macher JN, Salis RK, Blakemore KS, Tollrian R, Matthaehi CD *et al.* (2016) Multiple- stressor  
394 effects on stream invertebrates: DNA barcoding reveals contrasting responses of cryptic  
395 mayfly species. *Ecological Indicators*, **61**, 159-169.

- 396 Matschiner M (2015) Fitchi: haplotype genealogy graphs based on the Fitch algorithm.  
397 *Bioinformatics*, doi:10.1093/bioinformatics/btv717.
- 398 Meyer BS, Matschiner M, Salzburger W (2015) A tribal level phylogeny of Lake Tanganyika  
399 cichlid fishes based on a genomic multi-marker approach. *Molecular Phylogenetics and*  
400 *Evolution*, **83**, 56-71.
- 401 Morkovsky L, Paces J, Ridl J, Reifova R (2015) Scrimmer: designing primers from  
402 transcriptome data. *Molecular Ecology Resources*, **15**, 1415-1420.
- 403 Near TJ, Eytan RI, Dornburg A, Kuhn KL, Moore JA *et al.* (2012) Resolution of ray-finned  
404 fish phylogeny and timing of diversification. *Proceedings of the National Academy of*  
405 *Sciences*, **109**, 13698-13703.
- 406 O'Halloran DM (2015) PrimerView: high-throughput primer design and visualization. *Source*  
407 *Code for Biology and Medicine*, **10**, 8.
- 408 Pais FS, Ruy Pde C, Oliveira G, Coimbra RS (2014) Assessing the efficiency of multiple  
409 sequence alignment programs. *Algorithms for Molecular Biology*, **9**, 4.
- 410 Pereira-da-Conceicao LL, Price BW, Barber-James HM, Barker NP, de Moor FC *et al.* (2012)  
411 Cryptic variation in an ecological indicator organism: mitochondrial and nuclear DNA  
412 sequence data confirm distinct lineages of *Baetis harrisoni* Barnard (Ephemeroptera:  
413 Baetidae) in southern Africa. *BMC Evolutionary Biology*, **12**, 26.
- 414 R Core Team (2016) *R: A Language and Environment for Statistical Computing*. R  
415 Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, available at:  
416 <https://www.R-project.org> (last accessed 26 March 2016).
- 417 Ramirez-Gonzalez RH, Uauy C, Caccamo M (2015) PolyMarker: A fast polyploid primer  
418 design pipeline. *Bioinformatics*, **31**, 2038-2039.
- 419 Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A *et al.* (2012) MrBayes 3.2:  
420 efficient Bayesian phylogenetic inference and model choice across a large model space.  
421 *Systematic Biology*, **61**, 539-542.
- 422 Ruane S, Bryson RW, Jr., Pyron RA, Burbrink FT (2014) Coalescent species delimitation in  
423 milksnakes (genus *Lampropeltis*) and impacts on phylogenetic comparative analyses.  
424 *Systematic Biology*, **63**, 231-250.
- 425 Rutschmann S, Gattolliat JL, Hughes SJ, Báez M, Sartori M *et al.* (2014) Evolution and  
426 island endemism of morphologically cryptic *Baetis* and *Cloeon* species (Ephemeroptera,  
427 Baetidae) on the Canary Islands and Madeira. *Freshwater Biology*, **59**, 2516-2527.
- 428 Sela I, Ashkenazy H, Katoh K, Pupko T (2015) GUIDANCE2: accurate detection of  
429 unreliable alignment regions accounting for the uncertainty of multiple parameters. *Nucleic*  
430 *Acids Research*, **43**, W7-14.
- 431 Sobhy H, Haitham S, Philippe C (2012) Gemi: PCR primers prediction from multiple  
432 alignments. *Comparative and Functional Genomics*, **2012**, 1-5.
- 433 Stamatakis A (2014) RAxML version 8: a tool for phylogenetic analysis and post-analysis of  
434 large phylogenies. *Bioinformatics*, **30**, 1312-1313.
- 435 Stephens M, Donnelly P (2003) A comparison of bayesian methods for haplotype  
436 reconstruction from population genotype data. *The American Journal of Human Genetics*,  
437 **73**, 1162-1169.
- 438 Stephens M, Smith NJ, Donnelly P (2001) A new statistical method for haplotype  
439 reconstruction from population data. *The American Journal of Human Genetics*, **68**, 978-  
440 989.
- 441 Szitenberg A, John M, Blaxter ML, Lunt DH (2015) ReproPhylo: An Environment for  
442 Reproducible Phylogenomics. *PLoS Computational Biology*, **11**, e1004447.

- 443 Talavera G, Castresana J (2007) Improvement of phylogenies after removing divergent and  
444 ambiguously aligned blocks from protein sequence alignments. *Systematic Biology*, **56**,  
445 564-577.
- 446 Untergasser A, Cutcutache I, Koressaar T *et al.* (2012) Primer3—new capabilities and  
447 interfaces. *Nucleic Acids Research*, **40**, e115.
- 448 Vuataz L, Sartori M, Gattolliat JL, Monaghan MT (2013) Endemism and diversification in  
449 freshwater insects of Madagascar revealed by coalescent and phylogenetic analysis of  
450 museum and field collections. *Molecular Phylogenetics and Evolution*, **66**, 979-991.
- 451 Vuataz L, Sartori M, Wagner A, Monaghan MT (2011) Toward a DNA taxonomy of Alpine  
452 *Rhithrogena* (Ephemeroptera: Heptageniidae) using a mixed Yule-coalescent analysis of  
453 mitochondrial and nuclear DNA. *PLoS ONE*, **6**, e19728.
- 454 Yoon H, Leitner T (2015) PrimerDesign-M: a multiple-alignment based multiple-primer  
455 design tool for walking across variable genomes. *Bioinformatics*, **31**, 1472-1474.
- 456 You FM, Huo N, Gu YQ, Luo MC, Ma Y *et al.* (2008) BatchPrimer3: a high throughput web  
457 application for PCR and sequencing primer design. *BMC Bioinformatics*, **9**, 253.
- 458 Yu L, Barakat E, Di Francesco J, Herzig HP (2015) Two-dimensional polymer grating and  
459 prism on Bloch surface waves platform. *Optics Express*, **23**, 31640-31647.
- 460 Zeng L, Zhang Q, Sun R, Kong H, Zhang N *et al.* (2014) Resolution of deep angiosperm  
461 phylogeny using conserved nuclear genes and estimates of early divergence times. *Nature*  
462 *Communications*, **5**, 4956.

463 **Data Accessibility**

464 The program, user manual and example data sets are freely available at:  
465 <https://github.com/hdetering/discomark> (last accessed March 28, 2016). Scripts used for the  
466 analyses are available at: [https://github.com/srutschmann/python\\_scripts](https://github.com/srutschmann/python_scripts) (last accessed March  
467 28, 2016). All DNA sequences from this study are available under GenBank accessions:  
468 KU987258-KU987260, KU987265- KU987268, KU987273- KU987276, KU987285-  
469 KU987288. GenBank accession numbers for sequences included in previous studies are the  
470 following: KU971838-KU971840, KU971851, KU972090-KU972092, KU972104,  
471 KU972490-KU972492, KU972503, KU973060-KU973061, KU973074.

472 **Author Contributions**

473 S.R., H.D., S.S., and M.T.M. conceived the study. S.R. coordinated the project and performed  
474 the empirical analyses. H.D. implemented the program in Python. S.R. and H.D. drafted the  
475 manuscript. S.S. gave guidance for the ortholog prediction. J.F. provided the code of the PriFi  
476 web tool. All authors gave helpful comments to the manuscript and approved the final  
477 version.

478 **Tables**

479 **Table 1** List of species used for the usage examples of the closely related species; *Cloeon*  
480 *diptherum* s.l. species complex.

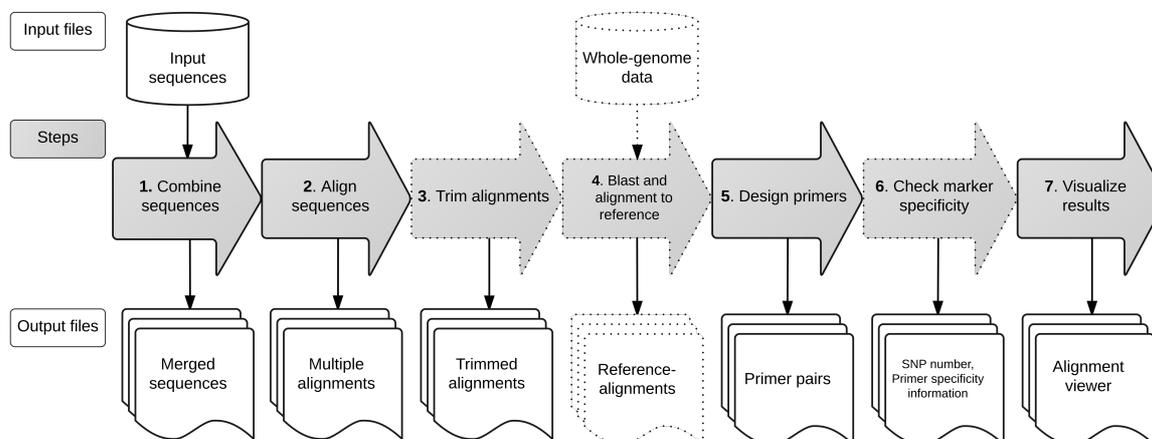
481

<b>Species</b>	<b>Voucher</b>	<b>Acc. no. <i>cox1</i></b>	<b>Geographical origin</b>
<i>Cloeon diptherum</i> sp1	SR21B07	KJ631626	Switzerland
<i>Cloeon diptherum</i> sp2	SR21B06	KJ631625	Switzerland
<i>Cloeon diptherum</i> sp3	US	KU757184	U.S.
<i>Cloeon peregrinator</i>	SR23A10	KU757122	Madeira

482

483 **Figures**

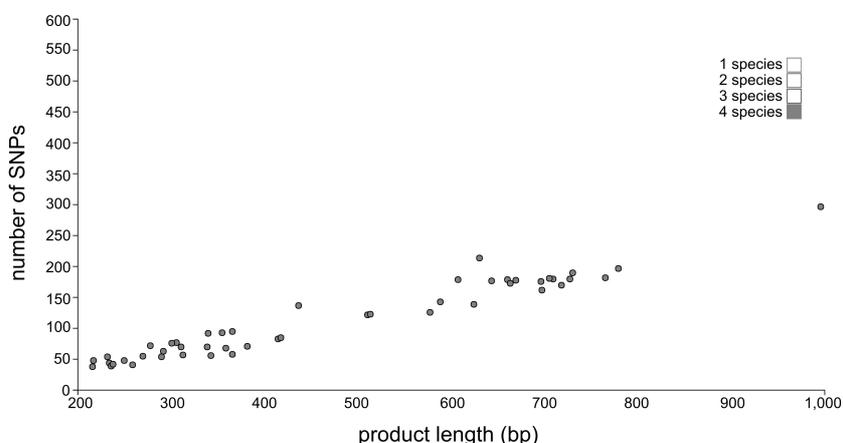
484



485

486 **Fig. 1** Overview of the DISCOMARK workflow and processing steps. Arrows with a broken  
487 outline indicate optional steps (for details see Materials and Methods section).

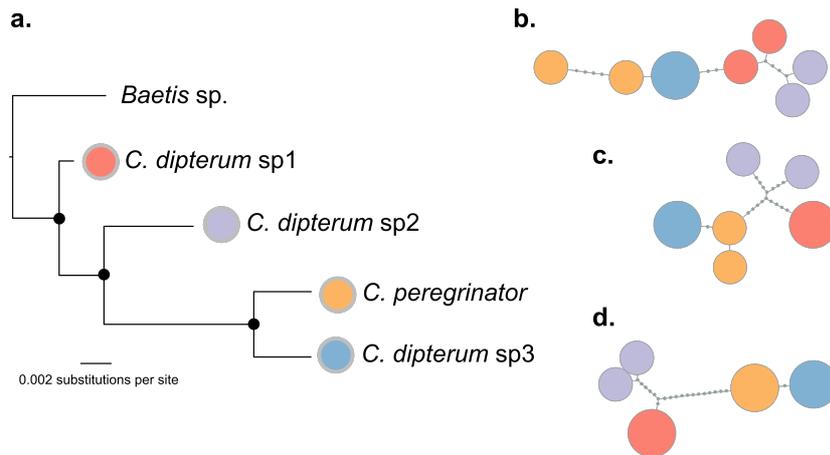
488



489

490 **Fig. 2** Visualization of DISCOMARK results: Scatter plot displaying the number of single  
491 nucleotide polymorphisms (SNPs) versus product length for each marker of the four mayfly  
492 species: *Baetis* sp., *Eurylophella* sp., *Ephemera danica*, and *Isonychia bicolor*. Shown are the  
493 markers for all four species (for details see Materials and Methods section).

494



495

496 **Fig. 3** Phylogenetic reconstruction and haplotype networks for the empirical data. **a**,  
497 Phylogenetic reconstruction of four representatives of the species complex *Cloeon dipterum*  
498 s.l., including *C. peregrinator*, based on the exon sequences of the eight newly developed  
499 nuclear DNA markers (2,526 base pairs). Bayesian inference was used to reconstruct the tree  
500 based on the concatenated supermatrix alignment. Bayesian posterior probabilities  $\geq 95\%$  are  
501 indicated by filled circles. *Baetis* was used as an outgroup. Scale bar represents substitutions  
502 per site. **b-d**, Haplotype networks of three amplified markers, **b**, marker 412045, **c**, marker  
503 412741, **d**, marker 412048 (full set of haplotype networks is available in Fig. S1, Supporting  
504 information). Circles are proportional to haplotype frequencies. Small circles along the branch  
505 indicate missing or unsampled haplotypes. Colors correspond to the four putative species.