

1 **Sequencing platform shifts provide opportunities**  
2 **but pose challenges for combining genomic**  
3 **datasets**  
4

5 Running title

6 **Challenges for combining genomic datasets**

7 Authors

8 De-Kayne, Rishi<sup>1,2†</sup>, Frei, David<sup>1,2,†</sup>, Greenway, Ryan<sup>1</sup>, Mendes, Sofia L.<sup>3</sup>, Retel Cas<sup>1</sup>,  
9 Feulner, Philine G. D.<sup>1,2,\*</sup>

10 <sup>1</sup>Department of Fish Ecology and Evolution, Centre of Ecology, Evolution and Biogeochemistry,  
11 EAWAG Swiss Federal Institute of Aquatic Science and Technology, Seestrasse 79, 6047  
12 Kastanienbaum, Switzerland

13 <sup>2</sup>Division of Aquatic Ecology and Evolution, Institute of Ecology and Evolution, University of Bern,  
14 Baltzerstrasse 6, 3012 Bern, Switzerland

15 <sup>3</sup>Centre for Ecology, Evolution and Environmental Changes, Faculdade de Ciências, Universidade  
16 de Lisboa, Lisbon, Portugal

17 † These authors contributed equally

18 \* Correspondence: philine.feulner@eawag.ch

19 Abstract

20 Technological advances in DNA sequencing over the last decade now permit the  
21 production and curation of large genomic datasets in an increasing number of non-  
22 model species. Additionally, this new data provides the opportunity for combining  
23 datasets, resulting in larger studies with a broader taxonomic range. Whilst the  
24 development of new sequencing platforms has been beneficial, resulting in a  
25 higher throughput of data at a lower per-base cost, shifts in sequencing technology

26 can also pose challenges for those wishing to combine new sequencing data with  
27 data sequenced on older platforms. Here, we outline the types of studies where  
28 the use of curated data might be beneficial, and highlight potential biases that  
29 might be introduced by combining data from different sequencing platforms. As an  
30 example of the challenges associated with combining data across sequencing  
31 platforms, we focus on the impact of the shift in Illumina's base calling technology  
32 from a four-channel to a two-channel system. We caution that when data is  
33 combined from these two systems, erroneous guanine base calls that result from  
34 the two-channel chemistry can make their way through a bioinformatic pipeline,  
35 eventually leading to inaccurate and potentially misleading conclusions. We also  
36 suggest solutions for dealing with such potential artifacts, which make samples  
37 sequenced on different sequencing platforms appear more differentiated from one  
38 another than they really are. Finally, we stress the importance of archiving tissue  
39 samples and the associated sequences for the continued reproducibility and  
40 reusability of sequencing data in the face of ever-changing sequencing platform  
41 technology.

#### 42 Keywords

43 NGS, reproducibility, reusability, poly-G, NovaSeq, HiSeq

44 Opportunities: Combining and extending datasets across time and  
45 space

46 DNA sequencing data reflecting the diversity of life is accumulating, as  
47 technological developments continue to increase the basepair yield of sequencing  
48 runs, whilst lowering the per-basepair prices. This data continues to facilitate  
49 comparative studies of genome structure for more and more organisms, spanning  
50 the tree of life (Baker et al., 2020; Cheng et al., 2018; Leebens-Mack et al., 2019;  
51 Morris et al., 2018; Peter et al., 2018; Shen et al., 2018; Shi et al., 2018; Zhang et  
52 al., 2014). Further, the field of molecular ecology is flourishing, with more and  
53 more studies investigating the genetic variation within and among closely related  
54 groups of organisms (Brawand et al., 2014; Lamichhaney et al., 2015; Tollis et al.,  
55 2018). However, for molecular ecologists working on non-model species, budgets  
56 still limit the amount of sequence data that can be produced. As a result,  
57 exhaustive experimental designs, which include the sampling of many individuals  
58 from many different populations, are rare (but are emerging; (Feulner et al., 2015;  
59 Greenway et al., 2020; Martin et al., 2016; Soria-Carrasco et al., 2014; Stankowski  
60 et al., 2019; Vijay et al., 2016)). The effort to publicly archive sequence data that  
61 has already contributed to publications helps to maintain the reproducibility of  
62 sequencing studies, whilst prolonging the value of such sequence data in  
63 perpetuity. Additionally, this practice of sequence data storage provides the  
64 opportunity to expand datasets beyond those that one laboratory is capable of  
65 producing (in terms of time, labour, and finances) to increase the impact of studies  
66 despite a potentially limited budget. Repositories like the Short Read Archive  
67 (SRA) -- part of the International Nucleotide Sequence Database Collaboration

68 (INSDC) that includes the NCBI Sequence Read Archive (SRA), the European  
69 Bioinformatics Institute (EBI), and the DNA Database of Japan (DDBJ) -- are  
70 essential for both the reproducibility of genetic and genomic studies, and the  
71 reusability of sequencing data. Although combining datasets is challenging for  
72 many sequencing approaches, particularly those that sequenced anonymous  
73 reduced representations of the genome (i.e. microsatellites, amplified fragment  
74 length polymorphisms, and maybe even restriction site associated DNA  
75 sequencing and genotyping by sequencing; but see Leigh, Lischer, Grossen, &  
76 Keller (2018) for an example), the increasingly common approach of re-  
77 sequencing whole-genomes (even for a broader range of non-model organisms)  
78 makes the possibility of combining datasets more inviting.

79 Between the continued growth of sequencing data repositories and the continued  
80 ability to sequence more DNA quicker and cheaper the following types of studies  
81 are increasingly carried out:

82 (1) Broad macroevolutionary studies. Typically, such macroevolutionary studies  
83 benefit from a wide taxon sampling and few individuals suffice, making the  
84 combination of samples from different published datasets particularly useful. Often  
85 these analyses are restricted to more conserved regions of the genome. For  
86 example, Zhang et al. (2020) compiled a comprehensive dataset of 365 species of  
87 asterids representing all 17 orders containing published and newly sequenced  
88 whole genomes and transcriptomes to resolve the deep asterid phylogeny. In  
89 another example, Greenway et al. (2020) focus on the Poeciliidae family of fish, to  
90 demonstrate that adaptation to extreme, here sulfide-rich, environments has

91 evolved convergently in ten independent lineages, by combining already published  
92 and newly sequenced transcriptome sequences.

93 (2) Microevolutionary studies investigating spatial variation across populations or  
94 closely related taxa. Such studies typically focus on one study system but rely on a  
95 larger sampling to reflect the variation within species or populations. These studies  
96 may benefit from combining newly sequenced material with archived sequence  
97 data from previous projects to produce larger within-system datasets. By taking  
98 advantage of existing sequence data, these combined datasets facilitate analyses  
99 of genomic differentiation across a much broader geographic sampling or among  
100 more individuals than would be otherwise possible. Here, the curated data is used  
101 to evaluate patterns in comparable populations to widen the perspective, i.e. to  
102 show whether a pattern is general or specific to the population under investigation.  
103 For example, Ravinet, Kume, Ishikawa, & Kitano (2020) evaluated if patterns of  
104 divergence and introgression between Japan Sea and Pacific Ocean stickleback  
105 resemble patterns at other locations where these species co-occur. In a  
106 comprehensive study conducted by Samuk et al. (2017), the authors compiled  
107 multiple genotyping by sequencing and whole genome sequencing datasets to a  
108 global evaluation of 1300 stickleback individuals across 51 populations, to show  
109 that putative adaptive alleles tend to occur more often in regions of low  
110 recombination. Bergland, Behrman, O'Brien, Schmidt, & Petrov (2014) used  
111 curated data to check haplotypes under seasonal selection in *Drosophila*  
112 *melanogaster* for between-species divergence with a sister species (*D. simulans*).  
113 Most recently, Jones, Mills, Jensen, & Good (2020) combined new and published  
114 whole-genome and exome sequences with targeted genotyping of *Agouti*, a

115 pigmentation gene introgressed from black-tailed jackrabbits, to investigate the  
116 evolutionary history of local seasonal camouflage adaptation in Snowshoe hares  
117 from the Pacific Northwest.

118 (3) Studies investigating temporal variation within and between population and  
119 species. Such studies involve combining datasets across time scales and often  
120 contain sequencing data that originated from a variety of sample types including  
121 museum collections, long-term preserved fossils or hard tissues, and  
122 contemporary fresh samples. For example, the use of museum specimens  
123 facilitated the investigation of independent temporal genomic contrasts spanning a  
124 century of climate change for two co-distributed chipmunk species (Bi et al., 2019)  
125 and a paleogenomics approach investigated the temporal component of  
126 adaptation to freshwater in sticklebacks by sequencing the genomes of 11-13,000-  
127 year-old bones and comparing them with 30 modern stickleback genomes (Kirch,  
128 Romundset, Gilbert, Jones, & Foote, 2020). Experimental approaches combining  
129 previous sequencing efforts with new samples are also commonly used to  
130 increase our understanding of temporal variation. Tenailon et al. (2016) compiled  
131 sequence data from several other publications in addition to new sequences to  
132 strengthen their conclusions on the tempo and mode of *E. coli* genome evolution.  
133 Bottery, Wood, & Brockhurst (2019), after having shown that tetracycline  
134 resistance requires multiple mutations, used curated data to investigate if the  
135 mutation establishment order was repeatable. This by no means exhaustive  
136 selection of examples highlights that the growing amount of sequence data  
137 provides the opportunity for endless combinations of datasets to be analysed to  
138 address a multitude of questions.

139 Challenges: Biases change with technological developments

140 One technological advance which sped up the Illumina workflow and made it more  
141 cost-effective was a change from four-channel chemistry, where each of the four  
142 DNA bases is detected by a different fluorescent dye, to a two-channel chemistry,  
143 that uses only two different fluorescent dyes (Illumina). In these two-channel  
144 workflows, as implemented in the NextSeq and NovaSeq platforms, a guanine  
145 base (G) is called in the absence of fluorescence (Figure 1). Hence, it is difficult to  
146 differentiate between no signal and a G, resulting in an overrepresentation of poly-  
147 G strings in sequence data from both NextSeq and NovaSeq (Chen, Zhou, Chen,  
148 & Gu, 2018).

149 To most accurately capture biological variation in a given sample or population, it  
150 is important to differentiate between potentially erroneous and correct base calls,  
151 which is often done using base quality scores. However, erroneous poly-G base  
152 calls produced on the NextSeq and NovaSeq platforms can be difficult to detect,  
153 because, as a result of the two-colour chemistry, they are not always associated  
154 with reduced base qualities. Unfortunately, read trimming software packages that  
155 were written for the older four-colour systems do not flag or trim poly-G tails.

156 Although one might think that mapping should remove the effect of these  
157 overrepresented Gs without the need for read trimming, it has been shown that  
158 some may still trickle through a bioinformatics pipeline and influence variant calling  
159 steps. A comprehensive empirical study making use of cancer cell lines to  
160 benchmark systematic differences between technologies revealed that NovaSeq  
161 instruments produced more stretches of Gs than HiSeqX in both paired-end reads  
162 (Arora et al., 2019). Arora et al. (2019) further confirmed that the bias remained

163 detectable in the mapped reads and resulted in a relatively large number of T > G  
164 mutations among the variants unique to the NovaSeq instrument. To reduce the  
165 potential down-stream impact of these poly-G strings, newer trimming software  
166 packages such as fastp (Chen et al., 2018) check the source of the data and  
167 implement poly-G trimming by default for the two-colour systems. This not only  
168 improves the computational efficiency of sequence alignment, but should also  
169 reduce the impact of erroneous variant calling on these bases.

170 The impact of these changes in base calling and the subsequent erroneous G  
171 calls on the biological interpretation may vary with the chosen experimental design  
172 and other sources of variation such as for example DNA quality. Although the  
173 biases resulting from not trimming off or filtering out poly-G strings might be mild or  
174 irrelevant when analysing data produced from high quality input DNA from a single  
175 system, this may not be true when data from different technologies are combined  
176 across various biological units (e.g. across populations, species, treatments, or  
177 time points). On top of variation in the quality of input DNA, a range of variation in  
178 sequencing approaches exists, along with differences in library preparation,  
179 including variation in read length or whether reads are single-end or paired-end.  
180 Where different individuals within a single dataset have been sequenced with  
181 variation in these methodological factors biases may also be exacerbated,  
182 potentially producing misleading results. Variation in length of sequences reads  
183 across a dataset for example has been shown to lead to pronounced allele  
184 frequency differences between populations and subsequently suggested false  
185 biological trends (Leight et al. 2018). Metagenomic work suggested that both  
186 library preparation and sequencing platform had systematic effects on the

187 microbial community description (Poulsen, Pamp, Ekstrøm, & Aarestrup, 2019;  
188 Sato et al., 2019). In summary, attention should be paid to DNA quality, library  
189 preparation protocols, and the sequencing platform used when analysing and  
190 interpreting publicly available genomic data.

191 Although the prospect of combining datasets to improve our power to detect  
192 patterns is alluring, it is important to consider the ways in which these data may  
193 result in misleading conclusions. Combining datasets often means combining data  
194 from different sequencing platforms, as DNA sequencing technology continues to  
195 develop through time. Unfortunately, some of the developments (e.g. the change  
196 from four-channel to two-channel chemistry in Illumina sequencing machines)  
197 have changed the way in which uncertainties in base calling are presented in the  
198 sequencer's output files. If managed incorrectly, these changes hamper our ability  
199 to combine datasets obtained with different sequencing technologies, and the  
200 subsequent genotyping and analysis of these combined datasets may be biased  
201 (in the worst cases leading to erroneous conclusions). The most straightforward  
202 way to prevent this is a well-thought out experimental design, a step which can  
203 often be overlooked in a time where sequencing data is being produced so rapidly  
204 (see Mason (2017) for sound advice on experimental design). As has been shown  
205 for sequencing reduced-representation libraries, it is crucial for any type of  
206 sequencing experiment to carefully consider types of errors that may be  
207 introduced during laboratory work and data processing, and how to minimize,  
208 detect and remove these errors (O'Leary, Puritz, Willis, Hollenbeck, & Portnoy  
209 2018). However, it may be difficult to achieve the ideal or optimal study design  
210 when an investigation integrates new information with already existing data (e.g.

211 with individuals and treatments randomised across sequencing batches). Despite  
212 this limitation there are a number of approaches that can help to rectify some of  
213 these imbalances and allow the combination of multiple genomic datasets whilst  
214 minimising the impact of cross-platform biases.

215 Ways forward: Suggestions on how to minimise technological bias  
216 when integrating datasets

217 Despite the ease with which new datasets can be produced it is critical that  
218 researchers do not forgo project planning and experimental design steps and aim  
219 to understand and reduce the potential impact of intrinsic data biases. These  
220 planning steps should be similar to those carried out for the sequencing of new  
221 samples and could include an assessment of the dataset (1) and the pipeline for  
222 analysis (2):

223 (1) When compiling a combined dataset, it is important to consider the key  
224 question that is being addressed and to evaluate how many samples of each  
225 population, species, treatment, or time unit are needed to have the power to draw  
226 meaningful conclusions. It is also worth evaluating the trade-offs between  
227 sequencing new samples or using existing data (e.g. if only a handful of samples  
228 are missing could it be worthwhile to sequence more samples so that all  
229 individuals are sequenced the same way, reducing the likelihood that biases or  
230 batch effects will cause problems downstream in the analysis). If datasets will be  
231 combined to address a specific question then it is important to assess which  
232 specific sequenced samples are available and how many different datasets these  
233 samples come from. It is important to be conscious of, and carefully document, the

234 different technologies used for library preparation and sequencing across samples  
235 and datasets, and if possible, to glean an understanding of the origin and quality of  
236 the input DNA. Ideally, the dataset would be compiled in a way that minimizes the  
237 number of differences between samples from different sources. Further, it is  
238 critical to strive to randomise samples from different biological units across  
239 different sequencing batches (Meirmans 2015). It can be particularly beneficial to  
240 repeat sequencing of one or a few representatives from a curated dataset to  
241 evaluate and correct potential biases. If feasible, repeated sequencing of the same  
242 individual allows to identify problematic loci that are not genotyped identically or  
243 consistently across technologies despite originating from the same individual. We  
244 therefore urge researchers wherever possible to archive tissue and/or DNA. These  
245 collections can be of tremendous value, as they facilitate the repeated sequencing  
246 of past samples into newly compiled datasets to determine whether any variants or  
247 alleles may have been erroneously missed because of technological biases. Using  
248 archived tissue or DNA in this way is one of the only possibilities to verify new  
249 sequence variants found using future technologies.

250 (2) Once it is decided that integrating dataset from various sources provides the  
251 best power to answer a particular question, it is important to determine which  
252 checks should be implemented in the analysis pipeline to avoid misleading  
253 biological interpretation of the data. The ways in which biological and technological  
254 differences are distributed across the compiled dataset should be reported and  
255 critical steps that would identify potentially problematic sequence artifacts and  
256 biases should be implemented in the bioinformatic pipeline. It is also crucial to  
257 determine how potential artifacts and biases amongst datasets will be handled.

258 Figure 2 provides a suggestion for a pipeline evaluating known differences  
259 between sequencing data produced with four-channel chemistry (e.g. HiSeqX) and  
260 two-channel chemistry (e.g. NovaSeq). We suggest comparing the FastQC report  
261 (Andrews, 2010) between samples sequenced with the two technologies to each  
262 other. Any systematic difference across FastQC reports might be relevant,  
263 however, when samples sequenced with different sequence chemistry that affects  
264 the base calling are combined reports on per base sequence and k-mers content  
265 are particularly worth paying attention to (see Figure 1 for an example, illustrating  
266 differences in k-mer counts). To see whether mapping reduces sequencing  
267 artefacts, FastQC can be re-run on only the reads that mapped well and will be  
268 used for genotyping. If biases persist, read trimming should be considered. Here  
269 fastp (Chen et al., 2018) could be used to trim poly-G tails efficiently. Once reads  
270 have been mapped, variants have been called, and genotypes have been  
271 determined, genotypes should be evaluated for potential batch effects. Here, we  
272 recommend identifying individuals sampled using different datasets and/or  
273 technologies with specific symbols or colours allowing the possible differences  
274 between these artificial groups to be highlighted (see section above). For example,  
275 in a Principal Component Analysis (PCA) which represents the various  
276 technological and sample differences by different symbols and biological  
277 differences (i.e. populations or species) by colour, any PC axis separating symbols  
278 instead of colours suggests there might be some technological bias causing batch  
279 effects (Figure 1). However, biases might not always show up as batch effects and  
280 are especially problematic when one population or other biological unit is the only  
281 one sequenced with a different technology. In this scenario, artifacts and biological  
282 differences would be confounded and as a result artifacts and biases would be

283 hard to detect (not visible as a batch effect in a PCA) and correct for. For this  
284 reason, we suggest that researchers aim to sequence biological units (species,  
285 populations, treatments, or time points) across each batch to avoid confounding  
286 biological differences with library or other technical effects. Alternatively, a bias  
287 might (although not necessarily) show up as a mutational bias relative to the  
288 reference, which can be evaluated and compared to published biases resulting  
289 from sequencing platform shifts (see Arora et al. (2019)). To reduce biases and  
290 undesired batch effects, the filtering parameters for variant calls and genotypes  
291 will need to be adjusted. One way to find the optimal filtering settings could be to  
292 determine which filtering thresholds allow you to minimize the differences between  
293 the detected batches. Specifically, it may be useful to compare distributions of  
294 quality scores between reference and alternate allele, which should look very  
295 similar in the absence of batch effects. However, we do not recommend solely  
296 relying on this to remove biases in the reads (such as poly-Gs in NovaSeq data)  
297 but mention this as one option that might help to reduce other sources of  
298 undesired batch effects. If none of these approaches suffice to identify and remove  
299 biases, one potential solution could be to define variable sites in a subset of the  
300 data, which only represents one technology, and then call genotypes on the whole  
301 dataset for only those regions. This comes with a potential ascertainment bias  
302 depending on how broadly biological units are represented in such a subset, but  
303 should reduce spurious variation caused by technological differences. Such an  
304 approach is similar to defining a SNP panel and then using SNPchips or other  
305 technologies to genotype a larger sampling (Kim et al., 2018). As all datasets are  
306 different, different approaches might be needed to reduce any effects of  
307 technological differences in compiled datasets. Critically, in each of these

308 scenarios the identification and removal of biases associated with technological  
309 shifts serves to reduce the possibility of incorrectly or erroneously inferring  
310 biological patterns or processes.

311 Finally, we want to emphasise the huge value of community efforts to archive  
312 sequencing data that makes science reproducible and reusable. We hope that we  
313 have demonstrated not only how technological shifts may pose challenges for the  
314 meaningful reusability of data, but also that the removal of biases associated with  
315 such shifts allows us to address new and exciting biological questions. We  
316 highlight the importance and value of accurate documentation, archiving of tissue  
317 and DNA samples, and sequence data, and urge researchers to assess the  
318 experimental design of their research projects to ensure scientifically sound and  
319 robust results.

## 320 Acknowledgements

321 We thank David Marques and the Fish Ecology and Evolution group at Eawag for  
322 their helpful comments and fruitful discussions on the topic. We are grateful to  
323 three reviewers for their constructive comments on the manuscript.

324 RD is supported by an SNSF grant (31003A\_163446) awarded to PF. DF is  
325 supported by the grant “SeeWandel: Life in Lake Constance – the past, present  
326 and future” within the framework of the Interreg V programme “Alpenrhein-  
327 Bodensee-Hochrhein (Germany/Austria/Switzerland/Liechtenstein)” which funds  
328 are provided by the European Regional Development Fund as well as the Swiss  
329 Confederation and cantons. D.F. received additional financial support from the  
330 Swiss Federal Office for the Environment and Eawag (Eawag Discretionary Funds

331 2018-2022). SLM is supported by an FCT scholarship (SFRH/BD/145153/2019)

332 granted by the Portuguese National Science Foundation (Fundação para a

333 Ciência e a Tecnologia – FCT).

334 The funders had no role in study design, decision to publish, or preparation of the

335 manuscript.

336

337 References

- 338  
339 Andrews, S. (2010). FastQC: a quality control tool for high throughput sequence  
340 data. Available online at:  
341 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.  
342 Arora, K., Shah, M., Johnson, M., Sanghvi, R., Shelton, J., Nagulapalli, K., . . .  
343 Robine, N. (2019). Deep whole-genome sequencing of 3 cancer cell lines  
344 on 2 sequencing platforms. *Scientific Reports*, 9, 19123.  
345 doi:10.1038/s41598-019-55636-3  
346 Baker, B. J., De Anda, V., Seitz, K. W., Dombrowski, N., Santoro, A. E., & Lloyd,  
347 K. G. (2020). Diversity, ecology and evolution of Archaea. *Nature*  
348 *Microbiology*, 5(7), 887-900. doi:10.1038/s41564-020-0715-z  
349 Bergland, A. O., Behrman, E. L., O'Brien, K. R., Schmidt, P. S., & Petrov, D. A.  
350 (2014). Genomic evidence of rapid and stable adaptive oscillations over  
351 seasonal time scales in *Drosophila*. *Plos Genetics*, 10(11), e1004775.  
352 doi:10.1371/journal.pgen.1004775  
353 Bi, K., Linderoth, T., Singhal, S., Vanderpool, D., Patton, J. L., Nielsen, R., . . .  
354 Good, J. M. (2019). Temporal genomic contrasts reveal rapid evolutionary  
355 responses in an alpine mammal during recent climate change. *Plos*  
356 *Genetics*, 15(5), e1008119. doi:10.1371/journal.pgen.1008119  
357 Bottery, M. J., Wood, A. J., & Brockhurst, M. A. (2019). Temporal dynamics of  
358 bacteria-plasmid coevolution under antibiotic selection. *Isme Journal*, 13(2),  
359 559-562. doi:10.1038/s41396-018-0276-9  
360 Brawand, D., Wagner, C. E., Li, Y. I., Malinsky, M., Keller, I., Fan, S., . . . Di  
361 Palma, F. (2014). The genomic substrate for adaptive radiation in African  
362 cichlid fish. *Nature*, 513(7518), 375-381. doi:10.1038/nature13726  
363 Chen, S. F., Zhou, Y. Q., Chen, Y. R., & Gu, J. (2018). Fastp: An ultra-fast all-in-  
364 one FASTQ preprocessor. *Bioinformatics*, 34(17), 884-890.  
365 doi:10.1093/bioinformatics/bty560  
366 Cheng, S. F., Melkonian, M., Smith, S. A., Brockington, S., Archibald, J. M.,  
367 Delaux, P. M., . . . Wong, G. K. S. (2018). 10KP: A phylodiverse genome  
368 sequencing plan. *Gigascience*, 7(3). doi:10.1093/gigascience/giy013  
369 Feulner, P. G. D., Chain, F. J. J., Panchal, M., Huang, Y., Eizaguirre, C., Kalbe,  
370 M., . . . Milinski, M. (2015). Genomics of divergence along a continuum of  
371 parapatric population differentiation. *Plos Genetics*, 11(2), e1005414.  
372 doi:10.1371/journal.pgen.1004966  
373 Greenway, R., Barts, N., Henpita, C., Brown, A. P., Rodriguez, L. A., Pena, C. M.  
374 R., . . . Shaw, J. H. (2020). Convergent evolution of conserved  
375 mitochondrial pathways underlies repeated adaptation to extreme  
376 environments. *Proceedings of the National Academy of Sciences of the*  
377 *United States of America*, 117(28), 16424-16430.  
378 doi:10.1073/pnas.2004223117  
379 Jones, M. R., Mills, L. S., Jensen, J. D., & Good, J. M. (2020). The origin and  
380 spread of locally adaptive seasonal camouflage in snowshoe hares.  
381 *American Naturalist*, 196(3), 316-332. doi:10.1086/710022  
382 Kim, J. M., Santure, A. W., Barton, H. J., Quinn, J. L., Cole, E. F., Visser, M. E., . . .  
383 . Great Tit HapMap, C. (2018). A high-density SNP chip for genotyping  
384 great tit (*Parus major*) populations and its application to studying the

385 genetic architecture of exploration behaviour. *Molecular Ecology*  
386 *Resources*, 18(4), 877-891. doi:10.1111/1755-0998.12778

387 Kirch, M., Romundset, A., Gilbert, M. T. P., Jones, F. C., & Foote, A. D. (2020).  
388 Pleistocene stickleback genomes reveal the constraints on parallel  
389 evolution. *bioRxiv*, 2020.2008.2012.248427.  
390 doi:10.1101/2020.08.12.248427

391 Lamichhaney, S., Berglund, J., Almen, M. S., Maqbool, K., Grabherr, M., Martinez-  
392 Barrio, A., . . . Andersson, L. (2015). Evolution of Darwin's finches and their  
393 beaks revealed by genome sequencing. *Nature*, 518(7539), 371-375.  
394 doi:10.1038/nature14181

395 Leebens-Mack, J. H., Barker, M. S., Carpenter, E. J., Deyholos, M. K.,  
396 Gitzendanner, M. A., Graham, S. W., . . . One Thousand Plant, T. (2019).  
397 One thousand plant transcriptomes and the phylogenomics of green plants.  
398 *Nature*, 574(7780), 679-685. doi:10.1038/s41586-019-1693-2

399 Leigh, D. M., Lischer, H. E. L., Grossen, C., & Keller, L. F. (2018). Batch effects in  
400 a multiyear sequencing study: False biological trends due to changes in  
401 read lengths. *Molecular Ecology Resources* 18: 778-788. doi:  
402 10.1111/1755-0998.12779

403 Marques, D. A., Lucek, K., Sousa, V. C., Excoffier, L., & Seehausen, O. (2019).  
404 Admixture between old lineages facilitated contemporary ecological  
405 speciation in Lake Constance stickleback. *Nature Communications*, 10,  
406 4240. doi:10.1038/s41467-019-12182-w

407 Martin, S. H., Most, M., Palmer, W. J., Salazar, C., McMillan, W. O., Jiggins, F. M.,  
408 & Jiggins, C. D. (2016). Natural selection and genetic diversity in the  
409 butterfly *Heliconius melpomene*. *Genetics*, 203(1), 525-541.  
410 doi:10.1534/genetics.115.183285

411 Mason, C. C. (2017). Four study design principles for genetic investigations using  
412 next generation sequencing. *Bmj-British Medical Journal*, 359, j4069.  
413 doi:10.1136/bmj.j4069

414 Meirmans, P. G. (2015). Seven common mistakes in population genetics and how  
415 to avoid them. *Molecular Ecology* 24: 3223-3231. doi: 10.1111/mec.13243

416 Morris, J. L., Puttick, M. N., Clark, J. W., Edwards, D., Kenrick, P., Pressel, S., . . .  
417 Donoghue, P. C. J. (2018). The timescale of early land plant evolution.  
418 *Proceedings of the National Academy of Sciences of the United States of*  
419 *America*, 115(10), E2274-E2283. doi:10.1073/pnas.1719588115

420 O'Leary Shannon, J., Puritz Jonathan, B., Willis Stuart, C., Hollenbeck  
421 Christopher, M., & Portnoy David, S. (2018). These aren't the loci you're  
422 looking for: Principles of effective SNP filtering for molecular ecologists.  
423 *Molecular Ecology* 27:3193–3206. doi: 10.1111/mec.14792

424 Peter, J., De Chiara, M., Friedrich, A., Yue, J. X., Pflieger, D., Bergstrom, A., . . .  
425 Schacherer, J. (2018). Genome evolution across 1,011 *Saccharomyces*  
426 *cerevisiae* isolates. *Nature*, 556(7701), 339-344. doi:10.1038/s41586-018-  
427 0030-5

428 Poulsen, C. S., Pamp, S. J., Ekstrøm, C. T., & Aarestrup, F. M. (2019). Library  
429 preparation and sequencing platform introduce bias in metagenomics  
430 characterisation of microbial communities. *bioRxiv*, 592154.  
431 doi:10.1101/592154

432 Ravinet, M., Kume, M., Ishikawa, A., & Kitano, J. Patterns of genomic divergence  
433 and introgression between Japanese stickleback species with overlapping

434 breeding habitats. *Journal of Evolutionary Biology*, 00, 1-14.  
435 doi:10.1111/jeb.13664

436 Samuk, K., Owens, G. L., Delmore, K. E., Miller, S. E., Rennison, D. J., & Schluter,  
437 D. (2017). Gene flow and selection interact to promote adaptive divergence  
438 in regions of low recombination. *Molecular Ecology*, 26(17), 4378-4390.  
439 doi:10.1111/mec.14226

440 Sato, M. P., Ogura, Y., Nakamura, K., Nishida, R., Gotoh, Y., Hayashi, M., . . .  
441 Hayashi, T. (2019). Comparison of the sequencing bias of currently  
442 available library preparation kits for Illumina sequencing of bacterial  
443 genomes and metagenomes. *DNA Research*, 26(5), 391-398.  
444 doi:10.1093/dnares/dsz017

445 Shen, X. X., Opulente, D. A., Kominek, J., Zhou, X., Steenwyk, J. L., Buh, K. V., . . .  
446 . Rokas, A. (2018). Tempo and mode of genome evolution in the budding  
447 yeast subphylum. *Cell*, 175(6), 1533-1545. doi:10.1016/j.cell.2018.10.023

448 Shi, M., Lin, X. D., Chen, X., Tian, J. H., Chen, L. J., Li, K., . . . Zhang, Y. Z.  
449 (2018). The evolutionary history of vertebrate RNA viruses. *Nature*,  
450 561(7722), E6. doi:10.1038/s41586-018-0310-0

451 Soria-Carrasco, V., Gompert, Z., Comeault, A. A., Farkas, T. E., Parchman, T. L.,  
452 Johnston, J. S., . . . Nosil, P. (2014). Stick insect genomes reveal natural  
453 selection's role in parallel speciation. *Science*, 344(6185), 738-742.  
454 doi:10.1126/science.1252136

455 Stankowski, S., Chase, M. A., Fuiten, A. M., Rodrigues, M. F., Ralph, P. L., &  
456 Streisfeld, M. A. (2019). Widespread selection and gene flow shape the  
457 genomic landscape during a radiation of monkeyflowers. *Plos Biology*,  
458 17(7), e3000391. doi:10.1371/journal.pbio.3000391

459 Tenailon, O., Barrick, J. E., Ribbeck, N., Deatherage, D. E., Blanchard, J. L.,  
460 Dasgupta, A., . . . Lenski, R. E. (2016). Tempo and mode of genome  
461 evolution in a 50,000-generation experiment. *Nature*, 536(7615), 165-170.  
462 doi:10.1038/nature18959

463 Tollis, M., Hutchins, E. D., Stapley, J., Rupp, S. M., Eckalbar, W. L., Maayan, I., . . .  
464 . Kusumi, K. (2018). Comparative genomics reveals accelerated evolution  
465 in conserved pathways during the diversification of anole lizards. *Genome  
466 Biology and Evolution*, 10(2), 489-506. doi:10.1093/gbe/evy013

467 Vijay, N., Bossu, C. M., Poelstra, J. W., Weissensteiner, M. H., Suh, A., Kryukov,  
468 A. P., & Wolf, J. B. W. (2016). Evolution of heterogeneous genome  
469 differentiation across multiple contact zones in a crow species complex.  
470 *Nature Communications*, 7, 10. doi:10.1038/ncomms13195

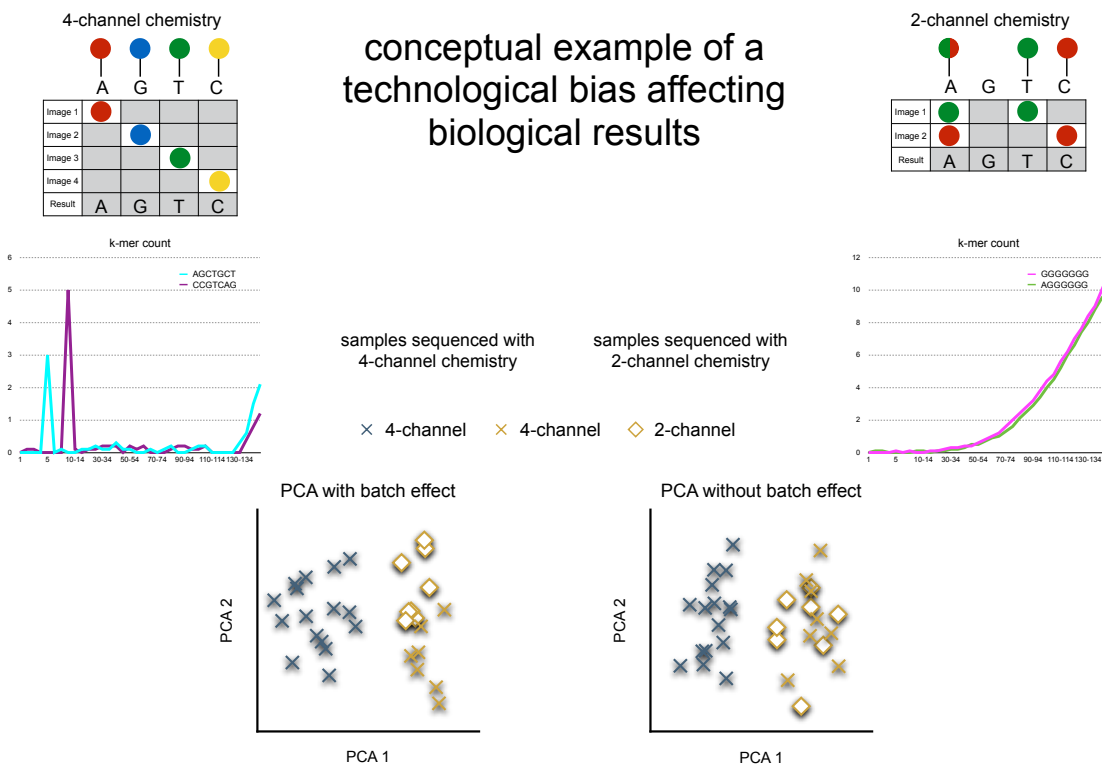
471 Zhang, G. J., Li, C., Li, Q. Y., Li, B., Larkin, D. M., Lee, C., . . . Avian Genome, C.  
472 (2014). Comparative genomics reveals insights into avian genome evolution  
473 and adaptation. *Science*, 346(6215), 1311-1320.  
474 doi:10.1126/science.1251385

475 Zhang, C., Zhang, T., Luebert, F., Xiang, Y., Huang, C.-H., Hu, Y., . . . Ma, H.  
476 (2020). Asterid phylogenomics/phylotranscriptomics uncover morphological  
477 evolutionary histories and support phylogenetic placement for numerous  
478 whole-genome duplications. *Molecular Biology and Evolution*.  
479 doi:10.1093/molbev/msaa160

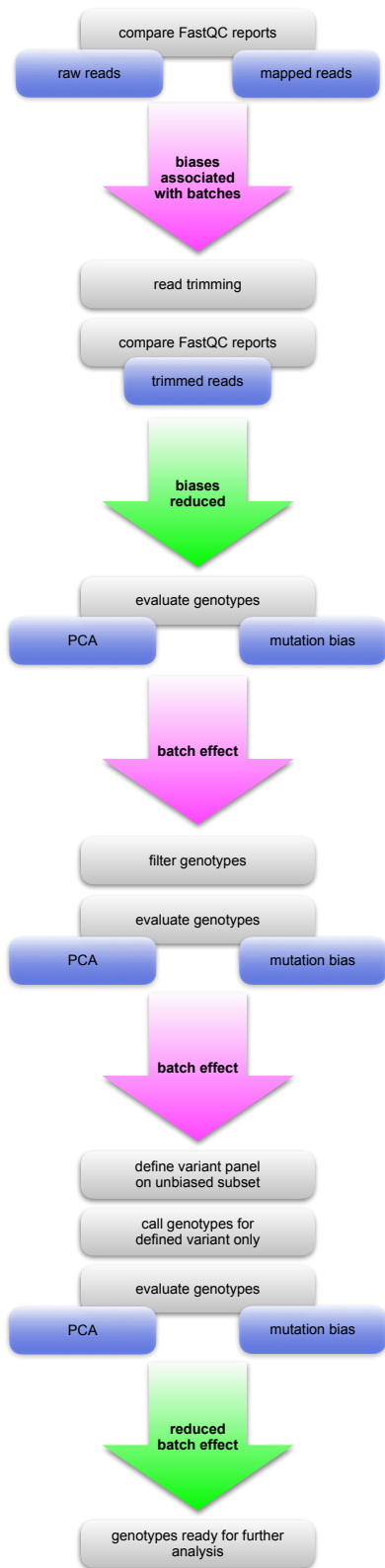
480

481 Author's contributions

482 RD, DF, and PF conceived of the presented ideas based on the experience and  
483 insights of DF. RD and PF drafted the manuscript. PF drafted the figures. All  
484 authors contributed to the discussion and critical revision of the final manuscript.  
485



486  
487 **Figure 1:** Example of a technological difference between sequencing chemistries, which  
488 introduces a bias (overrepresentation of G k-mers) in the sequenced reads and result in a batch  
489 effect visible when genotypes are evaluated in a principal component analysis (PCA).  
490 Top: Schematic redrawn from Illumina representing the differences between 4-channel chemistry  
491 evaluating each of the four bases by a distinct fluorescence label, and 2-channel chemistry  
492 representing the four bases with two dyes only.  
493 Middle: Redrawn examples of the one aspect of a typical FastQC (Andrews, 2010) report, which  
494 evaluates the count of each short nucleotide of length  $k$  (default = 7) starting at each position along  
495 the read. Any given  $k$ -mer should be evenly represented across the length of the read. The y axis  
496 reports the relative enrichment (log<sub>2</sub> observed over expected counts) of the 7-mers over the read  
497 length (x axis). The graph presents those  $k$ -mers which appear at specific positions with greater  
498 than expected frequency. In the left panel reads sequenced with 4-channel chemistry are  
499 represented which show a slight overrepresentation of two random 7-mers represented by different  
500 colours (typically the report would plot the first six hits). The overrepresentation is small and most  
501 pronounced at the beginning of the read (to the left of the x axis), a pattern often found in high  
502 quality sequencing libraries due to slight, sequence dependent efficiency of DNA shearing or a  
503 result of random priming. In the right panel, an overrepresentation of poly-G-mers toward the end  
504 of the reads is exemplified as typical for raw reads sequenced with 2-channel chemistry. Note the  
505 difference in the logarithmic scale between left and right panel.  
506 Bottom: Conceptual representation of a batch effect resulting from technological differences. Each  
507 sample's genotype, compiled of a large number of loci distributed across the whole genome, is  
508 represented as a coloured symbol in multivariate space, where PC axis one and two reflect two  
509 primary axes of variation in the dataset. The left panel would reflect a dataset with a batch effect.  
510 The fact that samples are separated by sequencing technology on PC axis 2 indicates the  
511 presence of a technological bias. In the right panel, batch effects have been reduced, e.g. by  
512 trimming off poly-G tails. Symbols in the PCA differentiate samples sequenced with either 2-  
513 channel (diamond) or 4-channel (cross) chemistry, colours differentiate different populations or  
514 species (biological differences). The left panel is imagined to be based on a data set of untrimmed  
515 reads, PC axis 2 separates samples due to technological differences. That effect is gone in the  
516 right panel, after read trimming was applied.



517  
518  
519

**Figure 2:** Flow diagram of an exemplified pipeline evaluating and accounting for biases caused by different sequencing technologies in a compiled data set. For more details see text.