

# **The future of next generation sequencing datasets: technological shifts provide opportunities but pose challenges for reproducibility and reusability**

## Running title

## **Challenges for reusability of NGS data**

## Authors

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>,  
Feulner, Philine G. D.<sup>1,2,\*</sup>

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

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

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

† These authors contributed equally

\* Correspondence: philine.feulner@eawag.ch

## Abstract

Technological advances in DNA sequencing over the last decade now permit the  
production and curation of large genomic datasets in an increasing number of non-  
model species. Additionally, this new data provides the opportunity for combining  
datasets, resulting in larger studies with a broader taxonomic range. Whilst the

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

#### Keywords

NGS, reproducibility, reusability, polyG, NovaSeq, HiSeq

## Opportunities: Combining and extending data sets across time and space

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

66 (INSDC) that includes the NCBI Sequence Read Archive (SRA), the European  
67 Bioinformatics Institute (EBI), and the DNA Database of Japan (DDBJ) -- are  
68 essential for both the reproducibility of genetic and genomic studies, and the  
69 reusability of sequencing data. Although reusability is challenging for many  
70 sequencing approaches, particularly those that sequenced a reduced  
71 representation of the genome (i.e. restriction site associated DNA sequencing,  
72 genotyping by sequencing, amplified fragment length polymorphism,  
73 microsatellites; but see Marques, Lucek, Sousa, Excoffier, & Seehausen (2019)),  
74 the increasingly common approach of re-sequencing whole-genomes (even for a  
75 broader range of non-model organisms) makes the possibility of combining  
76 datasets more inviting.

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

80 (1) Broad macroevolutionary studies. Typically, such macroevolutionary studies  
81 benefit from a wide taxon sampling and few individuals suffice, making the  
82 combination of samples from different published datasets particularly useful. Often  
83 these analyses are restricted to more conserved regions of the genome. For  
84 example, Hug et al. (2016) compile a phylogenomic data set containing published  
85 and newly sequenced whole genomes to build a phylogeny including Bacteria,  
86 Archaea and Eukarya using conserved sequences. In another example, Greenway  
87 et al. (2020) focus on the Poeciliidae family of fish, to demonstrate that adaptation  
88 to extreme, here sulfide-rich, environments has evolved convergently in ten

independent lineages, by combining already published and newly sequenced transcriptome sequences.

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

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

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

## Challenges: Biases change with technological developments

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

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

Although one might think that mapping should remove the effect of these overrepresented Gs without the need for read trimming, it has been shown that some may still trickle through a bioinformatics pipeline and influence variant calling steps. For example, cancer genomics demonstrated using cell lines the existence of systematic differences between the reads produced by HiSeqX and by NovaSeq as they noted a mild enrichment of T > G mutations in the variants called uniquely in NovaSeq and not in HiSeqX data (Arora et al., 2019). To reduce the

potential down-stream impact of these poly-G strings, newer trimming software packages such as fastp (Chen et al., 2018) check the source of the data and implement poly G trimming by default for the two-color systems. This not only improves the computational efficiency of sequence alignment, but should also reduce the impact of erroneous variant calling on these bases.

The impact of these changes in base calling and the subsequent erroneous G calls may also be affected by the experimental design or DNA quality. Although the biases resulting from not trimming off or filtering out poly-G strings might be mild or irrelevant when analysing data produced from high quality input DNA from a single system, this may not be true when data from different technologies are combined. On top of this, variation in the quality of input DNA may also amplify biases, potentially producing misleading results. Metagenomic work revealed that both library preparation and sequencing platform had systematic effects on the microbial community description (Poulsen, Pamp, Ekstrøm, & Aarestrup, 2019; Sato et al., 2019). In summary, attention should be paid to DNA quality, library preparation protocols, and sequencing platform used when analysing and interpreting publicly available genomic data.

Although the prospect of combining datasets to improve our power to detect patterns is alluring, it is important to consider the ways in which these data may result in misleading conclusions. Combining datasets often means combining data from different sequencing platforms, as DNA sequencing technology continues to develop through time. Unfortunately, some of the developments (e.g. the change from four-channel to two-channel chemistry in Illumina sequencing machines) have changed the way in which uncertainties in base calling are presented in the



sequencer's output files. If managed incorrectly, these changes hamper our ability to combine datasets obtained with different sequencing technologies, and the subsequent genotyping and analysis of these combined datasets may be biased (in the worst cases leading to erroneous conclusions). The most straightforward way to prevent this is a well-thought out experimental design, a step which can often be overlooked in a time where sequencing data is being produced so rapidly (see Mason (2017) for sound advice on experimental design). However, it may be difficult to achieve the ideal or optimal study design when an investigation integrates new information with already existing data (e.g. with individuals and treatments randomised across sequencing batches). Despite this limitation there are a number of approaches that can help to rectify some of these imbalances and allow the combination of multiple genomic datasets whilst minimising the impact of cross-platform biases.

#### How to minimise technological bias when combining datasets

Despite the ease with which new datasets can be produced it is critical that researchers do not forgo project planning and experimental design steps and aim to understand and reduce the potential impact of intrinsic data biases. These planning steps should be similar to those carried out for the sequencing of new samples and could include an assessment of:

(1) What is the key question that is being addressed and how many samples of each treatment or population are needed to have the power to draw meaningful conclusions? What might the tradeoffs be between sequencing new or using existing data (e.g. if only a handful of samples are missing could it be worthwhile

208 to sequence more samples so everything is sequenced similarly and sequence  
209 artifacts will not be problematic)? If we are to combine datasets then which  
210 individuals/populations are available to allow us to address our question?

211 (2) How many different datasets are combined? What technologies were used for  
212 library preparation and sequencing across the data sets? What is known about the  
213 origin and quality of the input DNA? Can we minimize the number of differences  
214 between data sets being compiled? Can we randomise biological  
215 samples/treatments across different sequencing batches? Do we have the option  
216 to repeat sequencing of one or a few representatives from a curated data set to  
217 evaluate potential biases? We also urge researchers wherever possible to archive  
218 tissue and/or DNA samples. These collections can be of tremendous value for  
219 future research, as they allow one to include repeated sequences of past samples  
220 into newly compiled data sets to determine whether any variants or alleles may  
221 have been erroneously missed because of technological biases. Using archived  
222 tissue or DNA is one of the only ways it is possible to verify new sequence variants  
223 found using future technologies.

224 (3) How are genetic differences, including those potentially causing biases,  
225 distributed across the compiled data set? What are the critical steps in an  
226 envisioned bioinformatic pipeline that would identify problematic sequence  
227 artifacts? How will we address known artifacts if they are present in our data  
228 and/or could confound our results? Figure 2 provides a suggestion for a pipeline  
229 evaluating known differences between sequencing data produced with four-  
230 channel chemistry (e.g. HiSeqX) and two-channel chemistry (e.g. NovaSeq). We  
231 suggest comparing the fastqc report (Andrews, 2010) between samples

232 sequenced with the two technologies to each other (see Figure 1 for an example,  
233 revealed by differences in kmer counts). To see if mapping reduces sequencing  
234 artefacts, fastqc can be re-run on only the reads that mapped well and will be used  
235 for genotyping. If biases persist, read trimming should be considered. Here fastp  
236 (Chen et al., 2018) could be used to trim polyG tails efficiently. Once reads have  
237 been mapped, variants have been called, and genotypes have been determined,  
238 genotypes should be evaluated for potential batch effects. Here, we recommend  
239 identifying individuals sampled using different data sets and/or technologies with  
240 specific symbols or colors allowing the possible differences between these artificial  
241 groups to be highlighted (see section above). For example, in a PCA which  
242 represents the various technological and sample differences by different symbols  
243 and biological differences (i.e. populations or species) by color, any PC axis  
244 separating symbols instead of colors suggests there might be some technological  
245 bias causing batch effects (Figure 1). Batch effects might be especially  
246 problematic when one population, timepoint, or treatment is the only one  
247 sequenced with a different technology. In this scenario artifacts and biological  
248 differences would be confounded and as a result would be hard to detect and  
249 correct for. For this reason, we suggest that researchers aim to sequence  
250 biological units, species, populations, or treatments across each batch to avoid  
251 confounding treatments/timesteps/populations with library or other technical  
252 effects. Alternatively, any mutational bias relative to the reference can be  
253 evaluated and be compared to the results established due to difference in  
254 sequencing technology only (see Arora et al. (2019)). To reduce batch effects  
255 once detected, filtering variant calls and genotypes will need to be adjusted. One  
256 way to find the critical filtering settings could be to see which filtering thresholds

257 allow you to minimize the differences between the detected batches. One  
258 promising approach might be to compare distributions of quality scores between  
259 reference and alternate allele, which should look very similar in the absence of  
260 batch effects. However, we do not recommend solely relying on this to remove  
261 detectable biases in the reads (such as poly Gs in NovaSeq data) but mention this  
262 option as it might help to reduce other sources of undesired batch effects. If none  
263 of these approaches suffice to identify and remove biases, one potential solution  
264 could be to define variable sites in a subset of the data, which only represents one  
265 technology, and then call genotypes on the whole data set for only those regions.  
266 This comes with a potential ascertainment bias depending on which biological  
267 units are represented in such a subset, but should allow to limit variation due to  
268 technological differences. Such an approach is similar to defining a SNP panel and  
269 then using SNPchips or other technologies to genotype a larger sampling (Kim et  
270 al., 2018). As all data sets are different, different approaches might be needed to  
271 reduce any effects of technological differences in compiled data sets. Critically, in  
272 each of these scenarios the identification and removal of biases associated with  
273 technological shifts serves to reduce the possibility of incorrectly or erroneously  
274 inferring biological patterns or processes.

275 Finally, we want to emphasise the huge value of our community efforts to archive  
276 sequencing data to make our science reproducible and reusable. We hope that we  
277 have demonstrated how technological shifts may pose challenges for the  
278 meaningful reusability of data, but also that the removal of biases associated with  
279 such shifts allows us to address new and exciting biological questions. We  
280 highlight the importance and value of accurate documentation, archiving of tissue

281 and DNA samples, and sequence data, and urge researchers to assess the  
282 experimental design of their research projects to ensure scientifically sound and  
283 robust results.

## 284 Acknowledgements

285 We thank David Marques and the Fish Ecology and Evolution group at Eawag for  
286 their helpful comments and fruitful discussions on the topic.

287 RD is supported by an SNSF grant (31003A\_163446) awarded to PF. DF is  
288 supported by the grant “SeeWandel: Life in Lake Constance - the past, present  
289 and future” within the framework of the Interreg V programme “Alpenrhein-  
290 Bodensee-Hochrhein (Germany/Austria/Switzerland/Liechtenstein)” which funds  
291 are provided by the European Regional Development Fund as well as the Swiss  
292 Confederation and cantons. SLM is supported by an FCT scholarship  
293 (SFRH/BD/145153/2019) granted by the Portuguese National Science Foundation  
294 (Fundação para a Ciência e a Tecnologia - FCT).

295 The funders had no role in study design, decision to publish, or preparation of the  
296 manuscript.

297

## 298 References

- 299  
300 Andrews, S. (2010). FastQC: a quality control tool for high throughput sequence  
301 data. Available online at:  
302 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.  
303 Arora, K., Shah, M., Johnson, M., Sanghvi, R., Shelton, J., Nagulapalli, K., . . .  
304 Robine, N. (2019). Deep whole-genome sequencing of 3 cancer cell lines  
305 on 2 sequencing platforms. *Scientific Reports*, 9, 19123.  
306 doi:10.1038/s41598-019-55636-3  
307 Baker, B. J., De Anda, V., Seitz, K. W., Dombrowski, N., Santoro, A. E., & Lloyd,  
308 K. G. (2020). Diversity, ecology and evolution of Archaea. *Nature*  
309 *Microbiology*, 5(7), 887-900. doi:10.1038/s41564-020-0715-z  
310 Bergland, A. O., Behrman, E. L., O'Brien, K. R., Schmidt, P. S., & Petrov, D. A.  
311 (2014). Genomic evidence of rapid and stable adaptive oscillations over  
312 seasonal time scales in *Drosophila*. *Plos Genetics*, 10(11), e1004775.  
313 doi:10.1371/journal.pgen.1004775  
314 Bi, K., Linderoth, T., Singhal, S., Vanderpool, D., Patton, J. L., Nielsen, R., . . .  
315 Good, J. M. (2019). Temporal genomic contrasts reveal rapid evolutionary  
316 responses in an alpine mammal during recent climate change. *Plos*  
317 *Genetics*, 15(5), e1008119. doi:10.1371/journal.pgen.1008119  
318 Bottery, M. J., Wood, A. J., & Brockhurst, M. A. (2019). Temporal dynamics of  
319 bacteria-plasmid coevolution under antibiotic selection. *Isme Journal*, 13(2),  
320 559-562. doi:10.1038/s41396-018-0276-9  
321 Brawand, D., Wagner, C. E., Li, Y. I., Malinsky, M., Keller, I., Fan, S., . . . Di  
322 Palma, F. (2014). The genomic substrate for adaptive radiation in African  
323 cichlid fish. *Nature*, 513(7518), 375-381. doi:10.1038/nature13726  
324 Chen, S. F., Zhou, Y. Q., Chen, Y. R., & Gu, J. (2018). Fastp: An ultra-fast all-in-  
325 one FASTQ preprocessor. *Bioinformatics*, 34(17), 884-890.  
326 doi:10.1093/bioinformatics/bty560  
327 Cheng, S. F., Melkonian, M., Smith, S. A., Brockington, S., Archibald, J. M.,  
328 Delaux, P. M., . . . Wong, G. K. S. (2018). 10KP: A phylodiverse genome  
329 sequencing plan. *Gigascience*, 7(3). doi:10.1093/gigascience/giy013  
330 Feulner, P. G. D., Chain, F. J. J., Panchal, M., Huang, Y., Eizaguirre, C., Kalbe,  
331 M., . . . Milinski, M. (2015). Genomics of divergence along a continuum of  
332 parapatric population differentiation. *Plos Genetics*, 11(2), e1005414.  
333 doi:10.1371/journal.pgen.1004966  
334 Greenway, R., Barts, N., Henpita, C., Brown, A. P., Rodriguez, L. A., Pena, C. M.  
335 R., . . . Shaw, J. H. (2020). Convergent evolution of conserved  
336 mitochondrial pathways underlies repeated adaptation to extreme  
337 environments. *Proceedings of the National Academy of Sciences of the*  
338 *United States of America*, 117(28), 16424-16430.  
339 doi:10.1073/pnas.2004223117  
340 Hug, L. A., Baker, B. J., Anantharaman, K., Brown, C. T., Probst, A. J., Castelle,  
341 C. J., . . . Banfield, J. F. (2016). A new view of the tree of life. *Nature*  
342 *Microbiology*, 1(5), 16048. doi:10.1038/nmicrobiol.2016.48  
343 Jones, M. R., Mills, L. S., Jensen, J. D., & Good, J. M. (2020). The origin and  
344 spread of locally adaptive seasonal camouflage in snowshoe hares.  
345 *American Naturalist*, 196(3), 316-332. doi:10.1086/710022

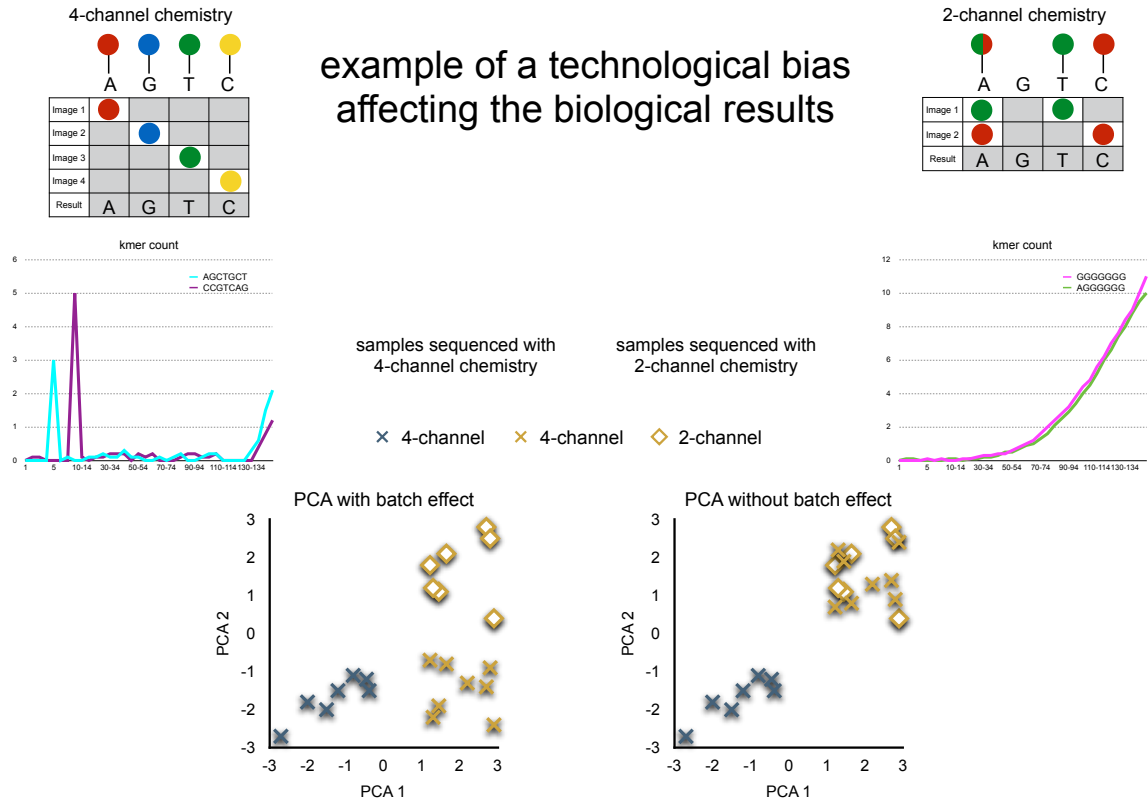
- Kim, J. M., Santure, A. W., Barton, H. J., Quinn, J. L., Cole, E. F., Visser, M. E., . . .  
 . Great Tit HapMap, C. (2018). A high-density SNP chip for genotyping  
 great tit (*Parus major*) populations and its application to studying the  
 genetic architecture of exploration behaviour. *Molecular Ecology  
 Resources*, 18(4), 877-891. doi:10.1111/1755-0998.12778
- Kirch, M., Romundset, A., Gilbert, M. T. P., Jones, F. C., & Foote, A. D. (2020).  
 Pleistocene stickleback genomes reveal the constraints on parallel  
 evolution. *bioRxiv*, 2020.2008.2012.248427.  
 doi:10.1101/2020.08.12.248427
- Lamichhaney, S., Berglund, J., Almen, M. S., Maqbool, K., Grabherr, M., Martinez-  
 Barrio, A., . . . Andersson, L. (2015). Evolution of Darwin's finches and their  
 beaks revealed by genome sequencing. *Nature*, 518(7539), 371-375.  
 doi:10.1038/nature14181
- Leebens-Mack, J. H., Barker, M. S., Carpenter, E. J., Deyholos, M. K.,  
 Gitzendanner, M. A., Graham, S. W., . . . One Thousand Plant, T. (2019).  
 One thousand plant transcriptomes and the phylogenomics of green plants.  
*Nature*, 574(7780), 679-685. doi:10.1038/s41586-019-1693-2
- Marques, D. A., Lucek, K., Sousa, V. C., Excoffier, L., & Seehausen, O. (2019).  
 Admixture between old lineages facilitated contemporary ecological  
 speciation in Lake Constance stickleback. *Nature Communications*, 10,  
 4240. doi:10.1038/s41467-019-12182-w
- Martin, S. H., Most, M., Palmer, W. J., Salazar, C., McMillan, W. O., Jiggins, F. M.,  
 & Jiggins, C. D. (2016). Natural selection and genetic diversity in the  
 butterfly *Heliconius melpomene*. *Genetics*, 203(1), 525-541.  
 doi:10.1534/genetics.115.183285
- Mason, C. C. (2017). *Four study design principles for genetic investigations using  
 next generation sequencing*. *Bmj-British Medical Journal*, 359, j4069.  
 doi:10.1136/bmj.j4069
- Morris, J. L., Puttick, M. N., Clark, J. W., Edwards, D., Kenrick, P., Pressel, S., . . .  
 Donoghue, P. C. J. (2018). The timescale of early land plant evolution.  
*Proceedings of the National Academy of Sciences of the United States of  
 America*, 115(10), E2274-E2283. doi:10.1073/pnas.1719588115
- Peter, J., De Chiara, M., Friedrich, A., Yue, J. X., Pflieger, D., Bergstrom, A., . . .  
 Schacherer, J. (2018). Genome evolution across 1,011 *Saccharomyces  
 cerevisiae* isolates. *Nature*, 556(7701), 339-344. doi:10.1038/s41586-018-  
 0030-5
- Poulsen, C. S., Pamp, S. J., Ekstrøm, C. T., & Aarestrup, F. M. (2019). Library  
 preparation and sequencing platform introduce bias in metagenomics  
 characterisation of microbial communities. *bioRxiv*, 592154.  
 doi:10.1101/592154
- Ravinet, M., Kume, M., Ishikawa, A., & Kitano, J. Patterns of genomic divergence  
 and introgression between Japanese stickleback species with overlapping  
 breeding habitats. *Journal of Evolutionary Biology*, 00, 1-14.  
 doi:10.1111/jeb.13664
- Samuk, K., Owens, G. L., Delmore, K. E., Miller, S. E., Rennison, D. J., & Schluter,  
 D. (2017). Gene flow and selection interact to promote adaptive divergence  
 in regions of low recombination. *Molecular Ecology*, 26(17), 4378-4390.  
 doi:10.1111/mec.14226

- Sato, M. P., Ogura, Y., Nakamura, K., Nishida, R., Gotoh, Y., Hayashi, M., . . . Hayashi, T. (2019). Comparison of the sequencing bias of currently available library preparation kits for Illumina sequencing of bacterial genomes and metagenomes. *DNA Research*, 26(5), 391-398. doi:10.1093/dnares/dsz017
- Shen, X. X., Opulente, D. A., Kominek, J., Zhou, X., Steenwyk, J. L., Buh, K. V., . . . Rokas, A. (2018). Tempo and mode of genome evolution in the budding yeast subphylum. *Cell*, 175(6), 1533-1545. doi:10.1016/j.cell.2018.10.023
- Shi, M., Lin, X. D., Chen, X., Tian, J. H., Chen, L. J., Li, K., . . . Zhang, Y. Z. (2018). The evolutionary history of vertebrate RNA viruses. *Nature*, 561(7722), E6. doi:10.1038/s41586-018-0310-0
- Soria-Carrasco, V., Gompert, Z., Comeault, A. A., Farkas, T. E., Parchman, T. L., Johnston, J. S., . . . Nosil, P. (2014). Stick insect genomes reveal natural selection's role in parallel speciation. *Science*, 344(6185), 738-742. doi:10.1126/science.1252136
- Stankowski, S., Chase, M. A., Fuiten, A. M., Rodrigues, M. F., Ralph, P. L., & Streisfeld, M. A. (2019). Widespread selection and gene flow shape the genomic landscape during a radiation of monkeyflowers. *Plos Biology*, 17(7), e3000391. doi:10.1371/journal.pbio.3000391
- Tenaillon, O., Barrick, J. E., Ribick, N., Deatherage, D. E., Blanchard, J. L., Dasgupta, A., . . . Lenski, R. E. (2016). Tempo and mode of genome evolution in a 50,000-generation experiment. *Nature*, 536(7615), 165-170. doi:10.1038/nature18959
- Tollis, M., Hutchins, E. D., Stapley, J., Rupp, S. M., Eckalbar, W. L., Maayan, I., . . . Kusumi, K. (2018). Comparative genomics reveals accelerated evolution in conserved pathways during the diversification of anole lizards. *Genome Biology and Evolution*, 10(2), 489-506. doi:10.1093/gbe/evy013
- Vijay, N., Bossu, C. M., Poelstra, J. W., Weissensteiner, M. H., Suh, A., Kryukov, A. P., & Wolf, J. B. W. (2016). Evolution of heterogeneous genome differentiation across multiple contact zones in a crow species complex. *Nature Communications*, 7, 10. doi:10.1038/ncomms13195
- Zhang, G. J., Li, C., Li, Q. Y., Li, B., Larkin, D. M., Lee, C., . . . Avian Genome, C. (2014). Comparative genomics reveals insights into avian genome evolution and adaptation. *Science*, 346(6215), 1311-1320. doi:10.1126/science.1251385

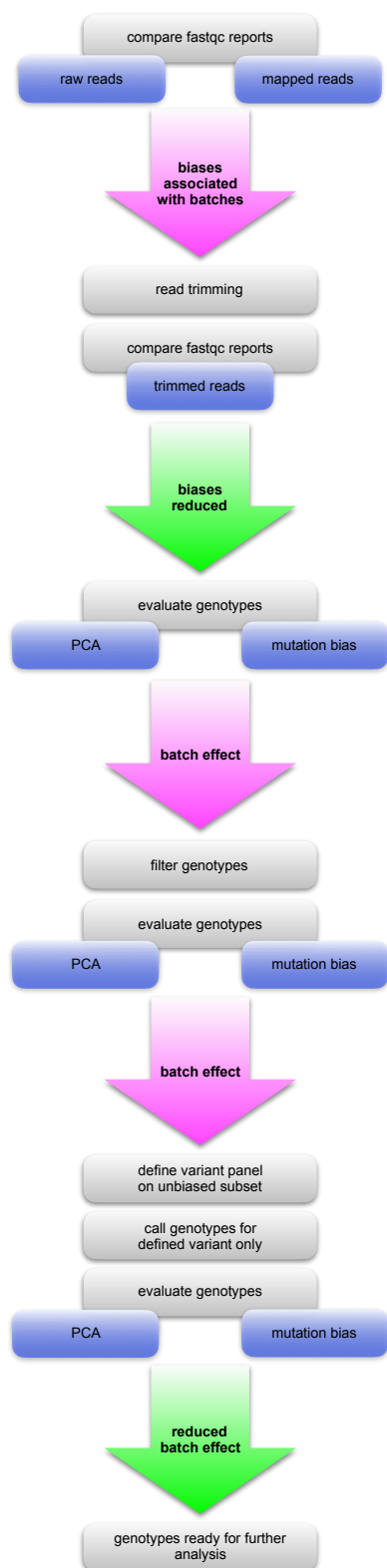
## Author's contributions

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





**Figure 1:** Example of a technological difference between sequencing chemistries, which introduces a bias (overrepresentation of G kmers) in the sequenced reads and result in a batch effect visible when genotypes are evaluated in a principal component analysis (PCA). Top: Schema redrawn from Illumina representing the differences between 4-channel chemistry evaluating each of the four bases by a distinct fluorescence label, and 2-channel chemistry representing the four bases with two dyes only. Middle: Redrawn examples of the one aspect of a typical fastqc (Andrews, 2010) report, which evaluates the count of each short nucleotide of length  $k$  (default = 7) starting at each position along the read. Any given Kmer should be evenly represented across the length of the read. The y axis reports the relative enrichment (log2 observed over expected counts) of the 7mers over the read length (x axis). The graph presents those kmers which appear at specific positions with greater than expected frequency. In the left panel reads sequenced with 4-channel chemistry are represented which show a slight overrepresentation of two random 7mers represented by different colors (typically the report would plot the first six hits). The overrepresentation is small and most pronounced at the beginning of the read (to the left of the x axis), a pattern often found in high quality sequencing libraries due to slight, sequence dependent efficiency of DNA shearing or a result of random priming. In the right panel, an overrepresentation of poly Gmers toward the end of the reads is exemplified as typical for raw reads sequenced with 2-channel chemistry. Note the difference in the logarithmic scale between left and right panel. Bottom: Each sample's genotype, compiled of a large number of loci distributed across the whole genome, is represented as a colored symbol in multivariate space, where PC axis one and two are presented here which explain some majority of variation across genotypes. Symbols in the PCA differentiate samples sequenced with either 2-channel (diamond) or 4-channel (cross) chemistry, colors differentiate different populations or species (biological differences). The left panel is imagined to be based on a data set of untrimmed reads, PC axis 2 separates samples due to technological differences. That effect is gone in the right panel, after read trimming was applied.



**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.