

Supplementary information

Next-generation data filtering in the genomics era

In the format provided by
the authors and unedited

1 **Supporting Information for: Next-generation data filtering in the genomics era**

2 William Hemstrom*, Jared A. Grummer*, Gordon Luikart, Mark R. Christie*

3

4 **Corresponding author information:**

5 * William Hemstrom; email: whemstro@purdue.edu

6 * Jared Grummer; email: jared.grummer@flbs.umt.edu

7 * Mark Christie; email: christ99@purdue.edu

8

9 **This PDF includes:**

10 Supplementary Text

11 Supplementary Tables 1-4

12 Supplementary Figures 1-7

13 **Supplementary Methods:**

14 **Empirical data preparation**

15 We acquired empirical data for ten taxa from published studies with taxonomic coverage
16 including mammals (killer whales [*Orcinus orca*]¹, deer mice [*Peromyscus maniculatus*]²,
17 humans [*Homo sapiens*]³, mountain goats [*Oreamnos americanus*]⁴, and white-tailed deer
18 [*Odocoileus virginianus*]⁵), arthropods (water fleas [*Daphnia pulex*]⁶, stoneflies [*Sweltsa*
19 *coloradensis*]⁷, and monarch butterflies [*Danaus plexippus*]⁸), fish (yellow perch [*Perca*
20 *flavescens*]⁹), and plants (*Arabidopsis thaliana*)¹⁰ (Table S2). The data types included whole-
21 genome sequencing, low-coverage whole-genome sequencing, exome capture, and restriction
22 site-associated DNA (RAD) sequencing. When more than two populations were sampled, we
23 randomly selected two with 30 individuals per population for filtering. Populations with fewer
24 than 30 individuals were not sub-sampled. For datasets where it was possible, we also applied
25 GATK's suggested hard-filters¹¹ (QD > 2, FS < 60, SOR > 3, MQ > 40, MQRankSum > -12.5,
26 ReadPosRankSum < -8) and a genotype quality (GQ) cutoff of 13 prior to filtering. The *D. pulex*
27 data was also filtered to remove loci with very high heterozygosities (>60%) to remove probably
28 paralogous loci.

29 **Simulated data preparation**

30 We used the *scrm* coalescent simulator¹² via the *coala* R package¹³ to simulate three genomic
31 datasets under three different demographic histories: a neutral (static) scenario, a recent
32 population bottleneck, and a recent population expansion (Table S3). For each model, we
33 simulated three populations, all of which descended from a common ancestral population which
34 split 1,000 generations before present to form populations A and (B + C). Populations B and C
35 then split from each other 500 generations later. Population C then remained static for 450

36 generations, after which it either continued without change (neutral/static model), declined
37 exponentially over five generations to 1/20th its original size (bottleneck model), or
38 exponentially expanded ten-fold over the same time-frame (expansion model). Prior to
39 demographic changes, all populations were held at a constant effective population size of 10,000.
40 Gene flow between populations B and C was allowed following the population split at a rate of
41 0.1 migrants per generation. For each model, we sampled 30 individuals from populations B and
42 C at the end of the simulation for 10 chromosomes, each with a length of 10mb and a
43 recombination rate averaging at one per chromosome per cross per generation. Population C was
44 used for all further analyses except F_{ST} , for which both B and C were used. An R markdown
45 document with the code used to perform these simulations is available in the Supplemental
46 Materials.

47 To simulate selection for a range of recombination rates (Table S3), we used the *msms*
48 simulator¹⁴, also via *coala*¹³. We used the same parameters as the neutral/static model, but each
49 chromosome had a different recombination rate with selection on a single new mutation
50 beginning 50 generations in the past. We varied recombination rates between 0.1 and 10 (results
51 are reported for recombination rates of 0.1 and 1 in Box 1 of the main text). In all cases, we used
52 a selection coefficient of 0.2 against the ancestral allele during selection. An R script with the
53 code used to perform these simulations is available in the Supplemental Materials
54 (Supplementary Notebook 3).

55 **Filtering**

56 We filtered each empirical and simulated dataset with a range of different filters and thresholds
57 using the *filter_snps* function in the *snpR* R package¹⁵. Specifically, we used the following filters
58 and thresholds (function arguments listed in italics in parentheses):

- 59 ● MAF (*maf*): 0.02–0.01 in 0.01 increments using within-group filtering (*maf_facets* =
60 “*pop*”) such that any locus with a MAF less than the threshold in *all* groups was
61 removed. We did not filter at 0.01 because only 30 diploid individuals were included
62 from each population, resulting in a minimum observable MAF of ~0.017 for
63 polymorphic loci.
- 64 ● HWP (*hwe*): 1×10^{-6} – 1×10^{-2} and 0.05 in increments of factors of ten (1×10^{-6} , 1×10^{-5} , and
65 so on) using within-group filtering (*hwe_facets* = “*pop*”) such that loci were removed
66 only if they were significantly out of HWP in any individual sample group. HWP was
67 assessed using an exact test¹⁶. No corrections for multiple testing were conducted to
68 ensure that identical filtering thresholds were used for all loci and to ensure conservative
69 removal of loci out of HWP.
- 70 ● Required % individuals genotyped (*min_ind*): 10–90% in increments of 10% such that
71 loci were removed if they were not genotyped in at least the given percentage of
72 individuals.
- 73 ● Required % loci genotyped (*min_loci*): 10–90% in increments of 10% such that
74 individuals were removed if they were not genotyped in at least the given percentage of
75 individuals.

76 When testing different thresholds for a given parameter value, we generally held all other
77 parameter values constant at these values:

78 • MAF = 0, MGC = 1. Note: MGC = 1 ($mgc = 1$) removes any loci sequenced in only one
79 individual, regardless of the genotypic state of that individual such that loci observed in a
80 single homozygous individual were still removed.

81 • HWP = 1×10^{-6}

82 • Required % individuals genotyped: 70%

83 • Required % loci genotyped: 70%

84 In addition to the solitary filter variation iterations, we also varied required % individuals and
85 loci genotyped together for their ranges (both values 10-90%) for all datasets and performed a
86 full factorial comparison of our filter thresholds for HWP and MAF for the mountain goat
87 (RAD) and stonefly datasets specifically.

88 Following filtering, we computed expected heterozygosity (H_E), observed heterozygosity
89 (H_O), nucleotide diversity (π), F_{IS} and pairwise F_{ST} according to Weir and Cockerham¹⁷,
90 Tajima's D ¹⁸, Watterson's θ ¹⁹, Tajima's θ ¹⁸, a rarefaction-corrected measure of the number of
91 segregating sites (Hemstrom and Christie, *in prep*), and a rarefaction-corrected estimate of the
92 number of private alleles²⁰ for each population or pair of populations (where applicable). Lastly,
93 we also estimated effective population sizes (N_e) for each population in each dataset using the
94 LD method in the NeEstimator software²¹ using only loci pairs on different chromosomes or
95 scaffolds. All analyses were performed using the *snpR* R package¹⁵. Filtering R scripts, bash
96 (shell) handling scripts, and parameter files are available at
97 https://github.com/ChristieLab/filtering_simulation_paper.

98 **Supplementary References:**

- 99 1. Kardos, M. *et al.* Inbreeding depression explains killer whale population dynamics. *Nature*
100 *Ecology & Evolution* **7**, 675–686 (2023).
- 101 2. Schweizer, R. M. *et al.* Broad Concordance in the Spatial Distribution of Adaptive and
102 Neutral Genetic Variation across an Elevational Gradient in Deer Mice. *Molecular Biology*
103 *and Evolution* **38**, 4286–4300 (2021).
- 104 3. Lowy-Gallego, E. *et al.* Variant calling on the GRCh38 assembly with the data from phase
105 three of the 1000 Genomes Project [version 2; peer review: 2 approved]. *Wellcome Open*
106 *Research* **4**, (2019).
- 107 4. Martchenko, D. & Shafer, A. B. A. Contrasting whole-genome and reduced representation
108 sequencing for population demographic and adaptive inference: an alpine mammal case
109 study. *Heredity* **131**, 273–281 (2023).
- 110 5. Kessler, C., Wootton, E. & Shafer, A. B. A. Speciation without gene-flow in hybridizing
111 deer. *Molecular Ecology* **32**, 1117–1132 (2023).
- 112 6. Maruki, T., Ye, Z. & Lynch, M. Evolutionary Genomics of a Subdivided Species. *Molecular*
113 *Biology and Evolution* **39**, msac152 (2022).
- 114 7. Malison, R. L. *et al.* Landscape Connectivity and Genetic Structure in a Mainstem and a
115 Tributary Stonefly (Plecoptera) Species Using a Novel Reference Genome. *Journal of*
116 *Heredity* **113**, 453–471 (2022).

- 117 8. Hemstrom, W. B., Freedman, M. G., Zalucki, M. P., Ramírez, S. R. & Miller, M. R.
118 Population genetics of a recent range expansion and subsequent loss of migration in monarch
119 butterflies. *Molecular Ecology* **31**, 4544–4557 (2022).
- 120 9. Schraidt, C. E. *et al.* Dispersive currents explain patterns of population connectivity in an
121 ecologically and economically important fish. *Evolutionary Applications* **n/a**, (2023).
- 122 10. Alonso-Blanco, C. *et al.* 1,135 Genomes Reveal the Global Pattern of Polymorphism in
123 *Arabidopsis thaliana*. *Cell* **166**, 481–491 (2016).
- 124 11. Van der Auwera, G. A. & O’Connor, B. D. *Genomics in the Cloud: Using Docker, GATK,*
125 *and WDL in Terra*. (O’Reilly Media, 2020).
- 126 12. Staab, P. R., Zhu, S., Metzler, D. & Lunter, G. scrm: efficiently simulating long sequences
127 using the approximated coalescent with recombination. *Bioinformatics* **31**, 1680–1682
128 (2015).
- 129 13. Staab, P. R. & Metzler, D. Coala: an R framework for coalescent simulation. *Bioinformatics*
130 **32**, 1903–1904 (2016).
- 131 14. Ewing, G. & Hermisson, J. MSMS: a coalescent simulation program including
132 recombination, demographic structure and selection at a single locus. *Bioinformatics* **26**,
133 2064–2065 (2010).
- 134 15. Hemstrom, W. & Jones, M. snpR: User friendly population genomics for SNP data sets with
135 categorical metadata. *Molecular Ecology Resources* **23**, 962–973 (2023).

- 136 16. Wigginton, J. E., Cutler, D. J. & Abecasis, G. R. A Note on Exact Tests of Hardy-Weinberg
137 Equilibrium. *The American Journal of Human Genetics* **76**, 887–893 (2005).
- 138 17. Weir, B. S. & Cockerham, C. C. Estimating F-Statistics for the Analysis of Population
139 Structure. *Evolution* **38**, 1358–1370 (1984).
- 140 18. Tajima, F. Statistical method for testing the neutral mutation hypothesis by DNA
141 polymorphism. *Genetics* **123**, 585 LP – 595 (1989).
- 142 19. Watterson, G. A. On the number of segregating sites in genetical models without
143 recombination. *Theoretical Population Biology* **7**, 256–276 (1975).
- 144 20. Kalinowski, S. T. Counting Alleles with Rarefaction: Private Alleles and Hierarchical
145 Sampling Designs. *Conservation Genetics* **5**, 539–543 (2004).
- 146 21. Do, C. *et al.* NeEstimator v2: re-implementation of software for the estimation of
147 contemporary effective population size (N_e) from genetic data. *Molecular Ecology*
148 *Resources* **14**, 209–214 (2014).
- 149 22. Li, M.-X., Yeung, J. M. Y., Cherny, S. S. & Sham, P. C. Evaluating the effective numbers of
150 independent tests and significant p-value thresholds in commercial genotyping arrays and
151 public imputation reference datasets. *Human Genetics* **131**, 747–756 (2012).

152 **Supplementary Tables:**153 **Supplementary Table 1: Different types of filters available for genomic sequencing data.**

Filter	Stage	Description
Base quality scores	<i>i</i>	Removal of reads with many poor-quality (likely mis-read) bases.
Poly-G tails	<i>i</i>	Removal of guanines ("G"s) erroneously called at the ends of reads on certain sequencing platforms.
Adapter/Barcode/Cut-site trimming	<i>i</i>	Removal of adapter, barcode, or cut-site sequences from the reads.
Adapter/Barcode/Cut-site mismatches	<i>i</i>	Removal of reads with sequences that do not match known adapter, barcode, or cut-site sequences.
Read K-mer distribution	<i>i, ii</i>	Removal of reads with too many very common or rare runs of base-pairs (K-mers).
Technical/PCR duplicates	<i>i, ii</i>	Thinning of technical or PCR duplicates down to a single representative read.
Alignment/Mapping scores	<i>ii</i>	Removal of reads that have mapping scores below a user-defined threshold.
Improperly paired reads (orientation and distance)	<i>ii</i>	Removal of paired-reads that are improperly paired (unexpectedly far apart or incorrectly oriented)
Stack depth of coverage	<i>ii</i>	Removal of loci "stacks" that have too low of a sequencing depth across samples; usually for reduced-representation sequencing.
Stack mismatches	<i>ii</i>	Removal of loci "stacks" that have too many mismatched base-pairs across samples; usually for reduced-representation sequencing.
Number of Alleles	<i>ii, iii</i>	Removal of genotypes, haplotypes, or "stacks" with too many possible alleles (usually > 2 for SNPs). Usually for computational efficiency, but also to remove potential errors.
Low coverage/Quality-by-depth	<i>iii</i>	Removal of individual called genotypes with coverage below a user-defined threshold. Joint "Quality-by-depth" often alternatively used.
Genotype Quality/Confidence	<i>iii</i>	Removal of individually called genotypes with genotyping confidence below a user-defined threshold. Joint "Quality-by-depth" often alternatively or additionally used.
High coverage	<i>iii</i>	Removal of individual called genotypes with coverage above a user-defined threshold (usually indicating errors in the reference, paralogs, or copy-number variants, all of which require additional investigation).

Supplementary Table 1: Continued from previous page.

Filter	Stage	Description
Insertion-deletions (Indels)	<i>iii</i>	Removal of insertions or deletions (indels), often required by many down-stream applications
Non-biallelic loci	<i>iii</i>	Removal of non-biallelic loci (for example, monomorphic or tri-allelic SNPs); required by many down-stream applications.
Allow/deny-listed variants	<i>iii</i>	Removal or inclusion of a set of user-defined loci. Common for methods that target specific loci or where specific variants are known to be problematic.
Variant Read Position	<i>iii</i>	Removal of variants that tend to occur in biased positions on shotgun-sequenced reads.
Missing data - per individual	<i>iii, iv</i>	Removal of individuals with called genotypes at fewer than a user-defined number of loci.
Missing data - per locus	<i>iii, iv</i>	Removal of loci with called genotypes at fewer than a user-defined number of individuals.
Minor allele frequency	<i>iii, iv</i>	Removal of loci with minor allele frequencies below a user-defined threshold.
Minor allele count	<i>iii, iv</i>	Removal of loci with a count of the minor allele below a user-defined threshold across samples.
Hardy-Weinberg proportions	<i>iii, iv</i>	Removal of loci out of Hardy-Weinberg proportions, typically below a user-defined p-value.
Strand Bias	<i>iii, iv</i>	Removal of loci where specific alleles are detected primarily on only the forward or reverse DNA strand.
Copy number variation	<i>iii, iv</i>	Removal of copy number variants. Often remain undiscovered.
Structural variants	<i>iii, iv</i>	Removal of structural variants, such as inversions. Often remain undiscovered.
Sex-linked loci	<i>iii, iv</i>	Removal of sex-linked loci, which may behave in unexpected ways or have biased statistical outcomes due to sex-specific sampling.
Paralogs - allelic imbalance/depth/heterozygosity	<i>ii, iii, iv</i>	Removal of reads aligned to paralogous genomic regions, where for recently diverged paralogs it can be unclear from which of the gene copies the read was sequenced. Additional analyses are required.

Filter	Stage	Description
Mislabeled/Contamination	<i>iv</i>	Removal of individuals or loci that are likely mislabeled, contaminated, or have similar issues. Can often be identified via PCA and other comparative analyses.
Transition-transversion bias	<i>iv</i>	Removal of loci from genomic regions with unexpected transition:transversion ratios.
F_{ST} /Selection Outliers	<i>iv</i>	Removal of outlier loci likely to be under selection. Useful for cases where putatively neutral processes specifically are of interest (for example, gene flow).

i = sequence QC (Quality control), *ii* = alignment to a reference, *iii* = variant discovery, and *iv* = data analysis. Note that stages *i* and *ii* constitute pre-variant filtering and stages *iii* and *iv* constitute post-variant filtering.

157 **Supplementary Table 2:** Empirical datasets used for filtering simulations.

158

Organism	Dataset Type	Number of SNPs	Number of Individuals	Reference
<i>Arabidopsis</i>	WGS	3,135,226	60	10 161 162
Monarch butterflies	RAD	238,368	54	8 163 164
<i>Daphnia</i>	WGS			6 165 166
Stoneflies	RAD	279,496	60	7 167
Yellow perch	WGS	6,586,547	57	9 168 169
Humans	WGS	17,458,468	60	3 170
Deer mice	Exome capture	5,373,633	55	2 171 172 173
Mountain goats	RADseq/ WGS	48,192/8,113,114	60/20	4 174 175
White-tailed deer	WGS	48,441,150	20	5 176 177
Killer whales	WGS	3,015,993	54	1 178 179

180

Supplementary Table 3: Simulated datasets used for filtering simulations.

Selection	Demographic History	Relevant Parameters	Number of SNPs	Number of Individuals
No Selection	Neutral/Static	$N_{curr} = N_{anc}$	237,698	60
No Selection	Large expansion in one population	$N_{curr} = 10N_{anc}$	247,922	60
No Selection	Large bottleneck in one population	$N_{curr} = 0.05N_{anc}$	198,783	60
Hard sweep at one locus in one population	Static	$s = 0.2$ $f_s = 1/2N$ $r = 0.1$ or 1 $t_s = 50$	21228/21744	60/60

N_{curr} : sample size at present

N_{anc} : ancestral sample size

s : selection coefficient against ancestral allele

f_s : frequency of the new mutation

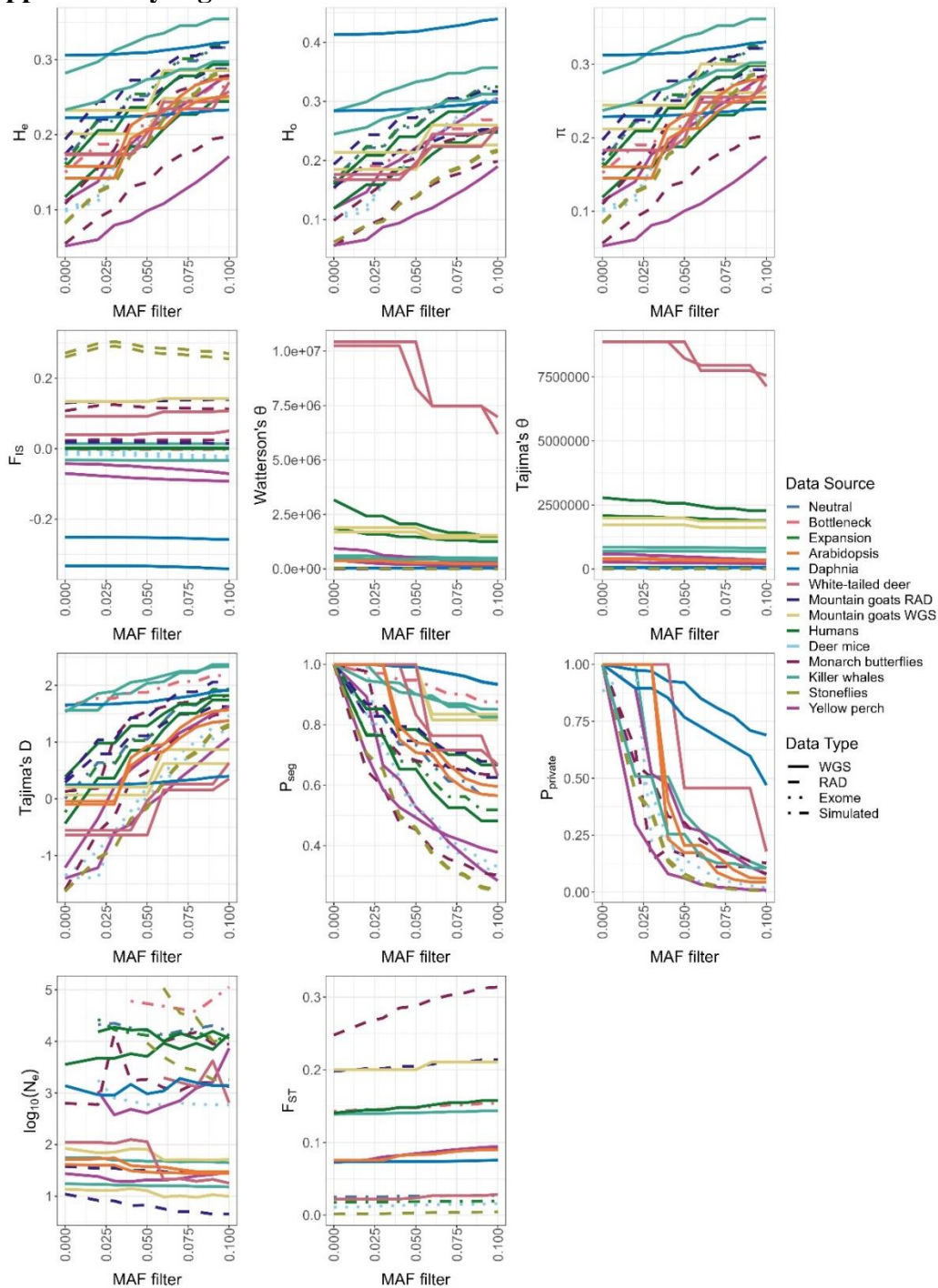
r : recombination rate, average number per chromosome per generation

t_s : time (in generations) before present at which selection began

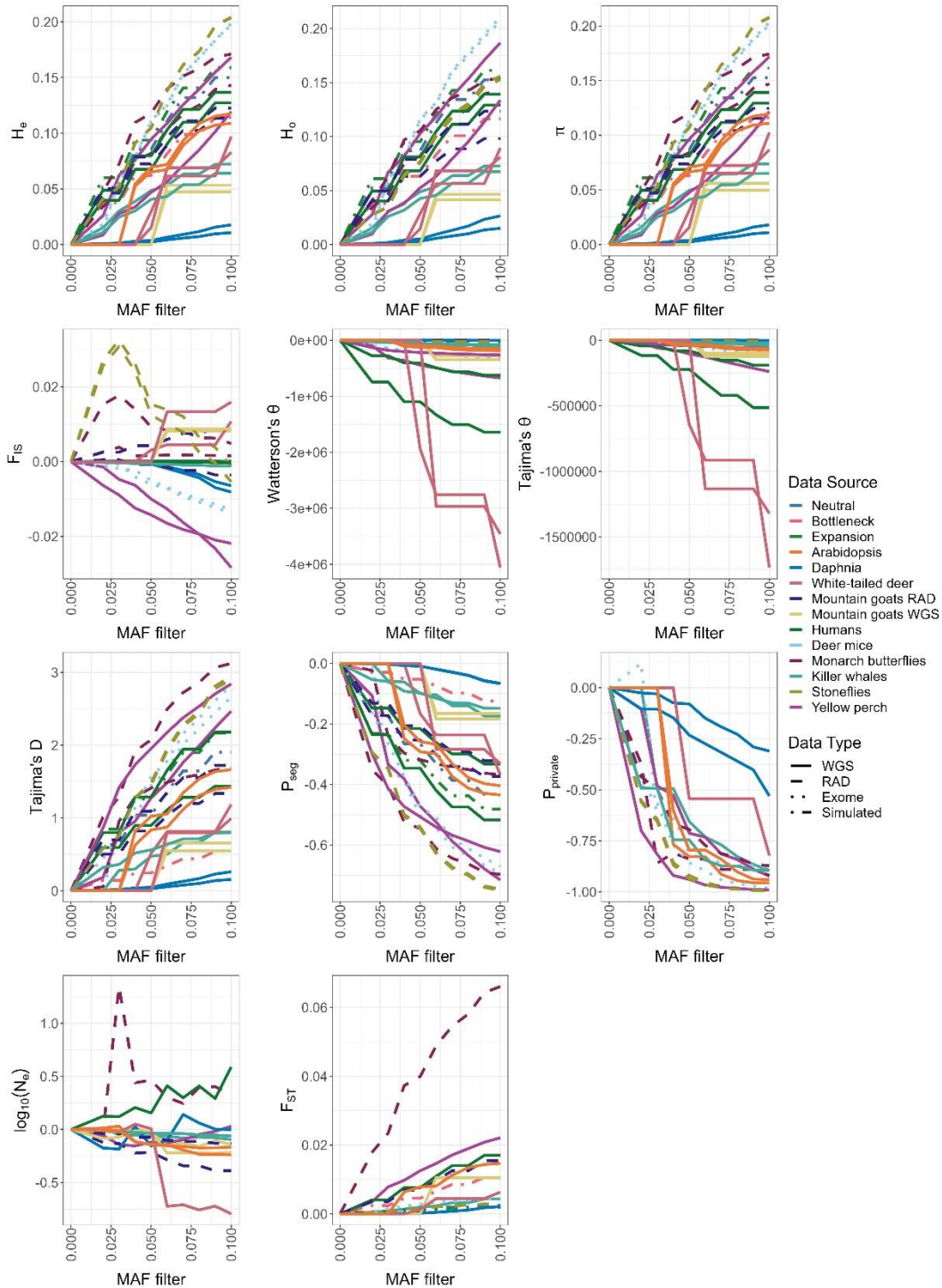
182 **Supplementary Table 4:** Justifications for filtering threshold recommendations in Table 1. Note: suggested MAC values may actually
 183 be stricter than the suggested MAF values for datasets with low sample sizes and should be adjusted accordingly. Some methods for
 184 specific questions may require filter values that differ from these suggestions.
 185

Question/Approach	Individual missing data; <X% missing loci	Loci missing data; <X% missing individuals	MAC/MAF	LD	HWP
Population Structure	Poor quality individuals or loci may mask structure	Poor quality individuals or loci may mask structure	Higher MAF values can reveal additional population structure	Probably has no effect unless inversions or other factors drive clustering.	Need to keep a low pass filter for cryptic populations
Demography	Many demographic estimation approaches function well with small sample sizes but can be misled by poor quality individuals	Projection can reduce the impact of missing data.	Any measures which depend on SFS estimates are extremely biased by the removal of low frequency variants.	Non-independent loci can create misleading site frequency spectra. Not as essential for Tajima's D as it is for ABC, etc.	Paralog removal
Selection	Don't want to remove signal	Don't want to remove signal	Don't want to remove signal, but some methods require or recommend higher MAF (0.05)	Removing loci in LD can remove signals of hitchhiking/selective sweeps	Don't want to remove signal; less strict filtering to keep more loci
Genetic Diversity	Some metrics (# seg sites) can drop fast as you exclude too many individuals/loci (see Box 1 fig.)	Some metrics (# seg sites) can drop fast as you exclude too many individuals/loci (see Box 1 fig.)	Note that some metrics (like # seg sites) are impacted differently from others (like H_E)	Should not <i>usually</i> cause major impacts, but can if regions in high LD vary (chromosomal inversion). Can skew confidence intervals.	This range will usually capture most of the changes due to filtering (see Box 1 fig).

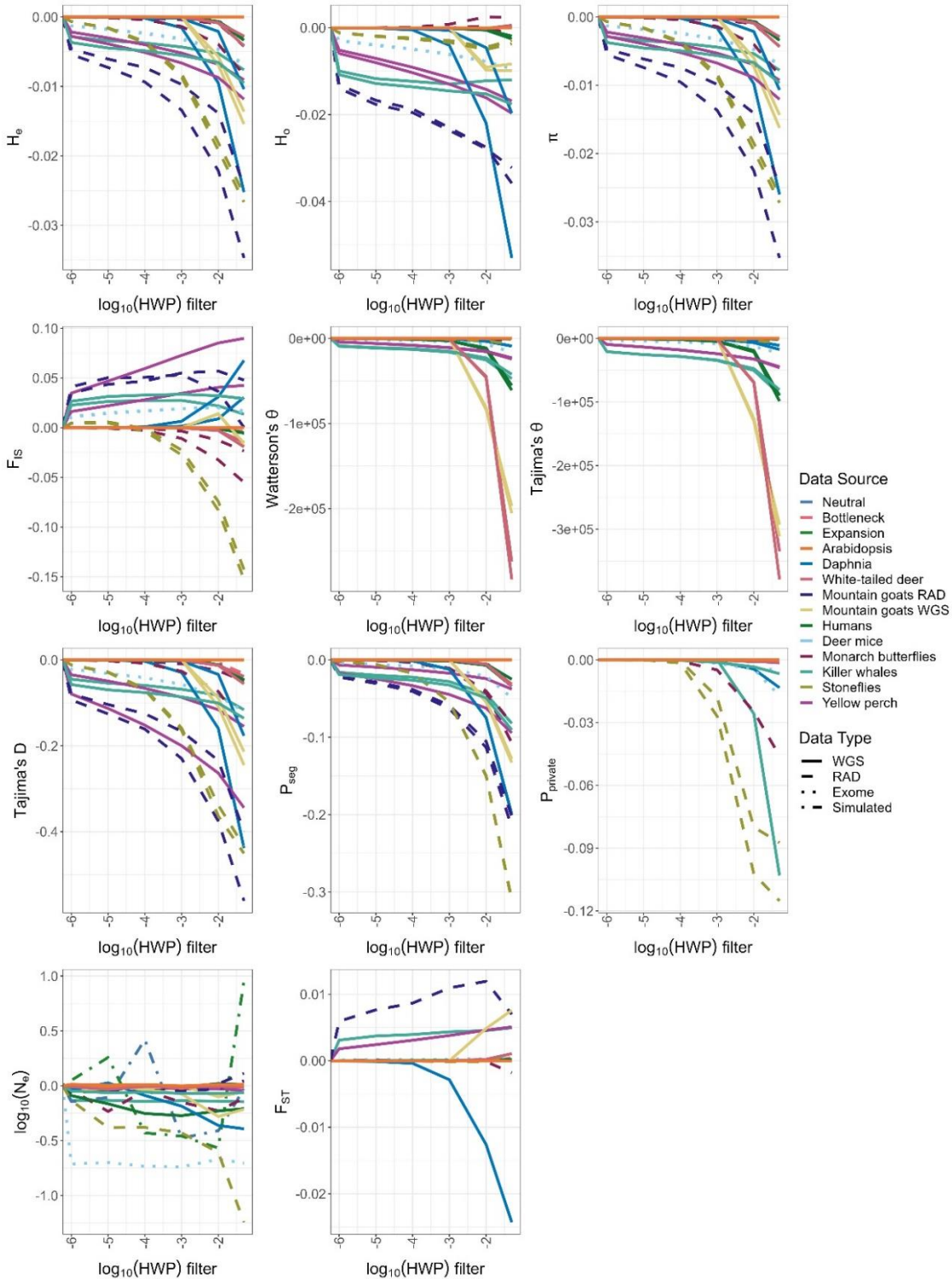
Question/Approach	Individual missing data; <X% missing loci	Loci missing data; <X% missing individuals	MAF/MAC	LD	HWP
Phylogenetics	High amounts of missing data within a locus can cause long branch attraction and incorrect topological inferences	For species-tree inference, must assure adequate representation of individuals in each species across gene trees	Autapomorphies/Singletons are not informative	Need independent evolutionary histories (e.g., unlinked)	Don't want anything with odd behavior throwing off signal
GWAS	If your method requires no missing data, use a high filter to avoid extra imputation	If your method requires no missing data, use a high filter to avoid extra imputation	Low frequency variants are typically uninformative unless sample sizes are very large.	Removal can cloud signals from linkage around causal genes. <i>p</i> -value correction is needed for multiple testing—see ²²	Selection on causal may rarely cause deviations. No filter to check, then do a permissive filter to keep most loci.
Mutation Detection	Cannot have missing data in parents, but skip over missing data in offspring (can't detect mutations but won't cause problems)	Cannot have missing data in parents, but skip over missing data in offspring (can't detect mutations but won't cause problems)	Many (new) mutations have very low frequencies	Irrelevant; don't want to remove potential mutations.	One test to get rid of paralogs, otherwise don't filter
Metagenomics/eDNA	NA	NA	Need multiple reads from a region for confidence (5+ or so)	Depends on context: no filter for a mix of species but consider otherwise	NA
Relatedness/ Pedigree Construction	Including all individuals is important	Only a few loci needed to infer a relationships but they need to be high quality.	Singletons don't help	Should not bias mean outcomes, but could change confidence intervals.	Only a few loci needed to infer a relationships but they need to be high quality.



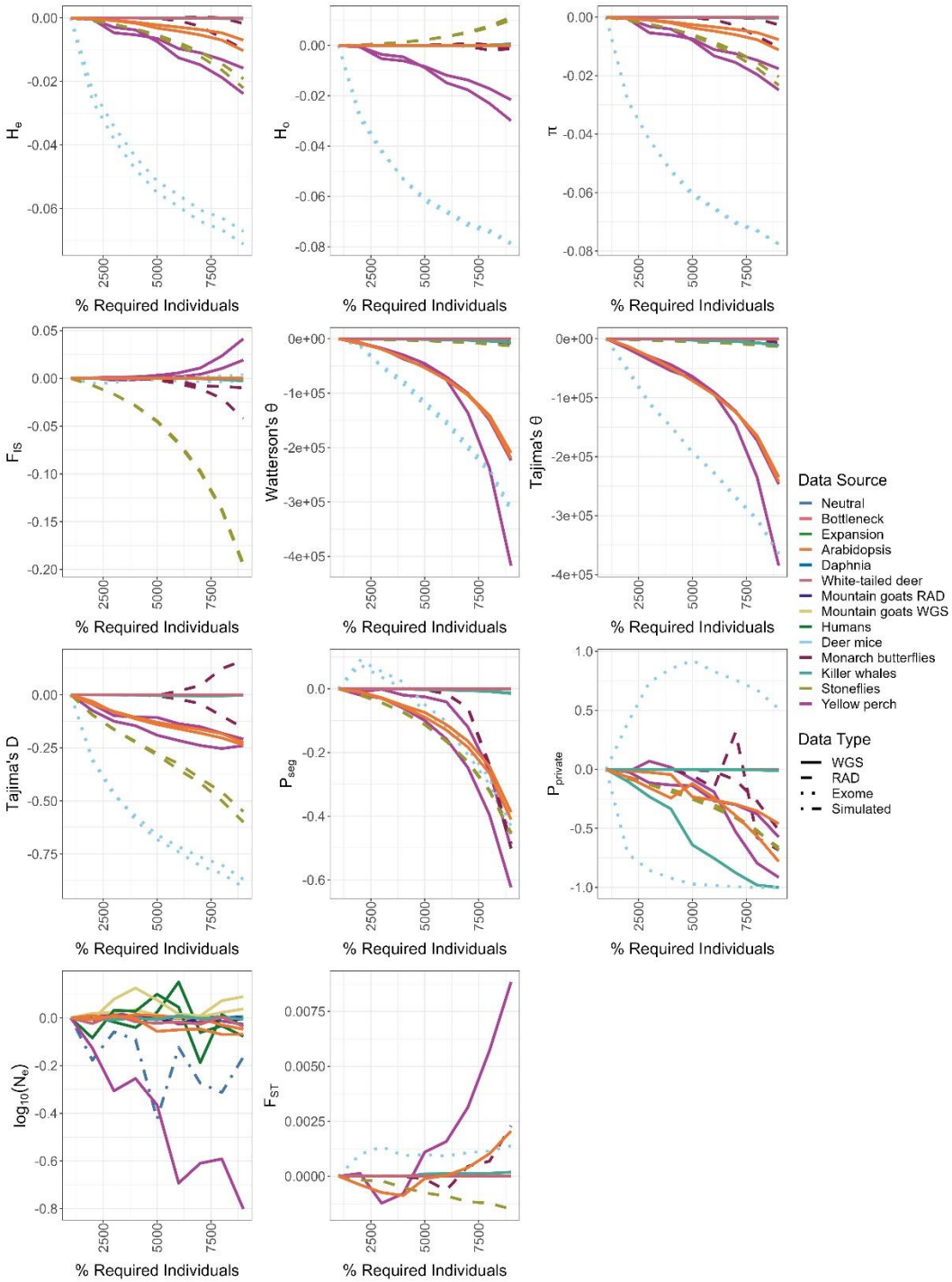
189
 190 **Supplementary figure S1:** Effects of **MAF filter** threshold on **non-standardized** change in F_{ST} ,
 191 H_o , F_{IS} , Tajima's D, H_e , π , Watterson's θ , Tajima's θ , the proportion of the total (highest)
 192 number of segregating sites, the proportion of the total (highest) number of private alleles, and
 193 N_e , observed in each population and dataset. Note the drastic changes in F_{ST} in the monarch
 194 dataset, driven by strong differences in the site frequency spectra between populations. The
 195 removal of rare alleles substantially changes the distribution of allele frequencies and increases
 196 F_{ST} in one population which had undergone a substantial demographic expansion.



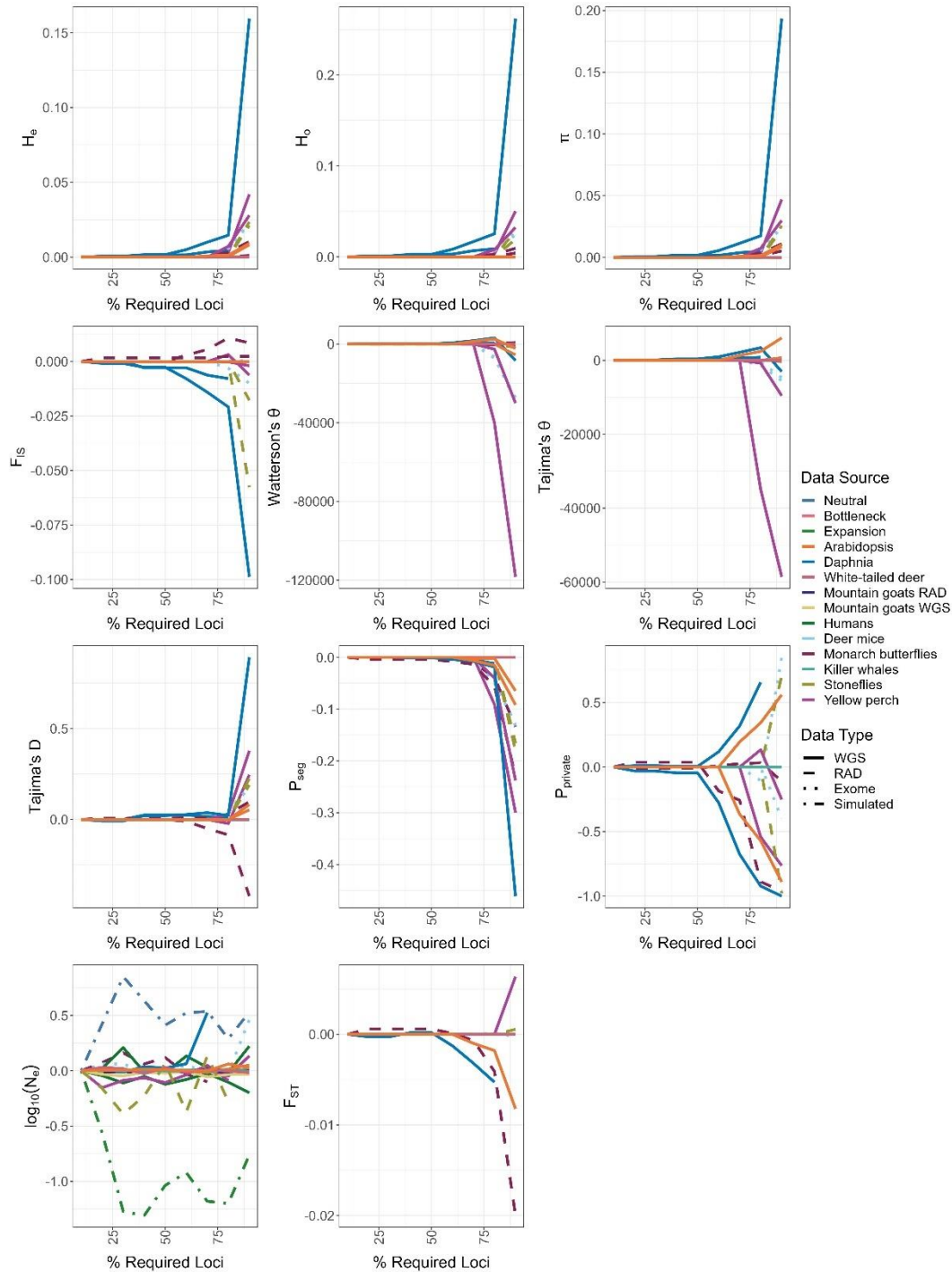
197
 198 **Figure S2:** Effects of **MAF filter** threshold on **standardized** F_{ST} , H_o , F_{IS} , Tajima's D, H_e , π ,
 199 Watterson's θ , Tajima's θ , the proportion of the total (highest) number of segregating sites, the
 200 proportion of the total (highest) number of private alleles, and N_e , observed in each population
 201 and dataset depending on MAF filtering stringency. Parameter values have been standardized to
 202 a value of 0 for the first MAF filter value (also 0) to highlight differences among studies. Some
 203 panels (H_o , F_{IS} , F_{ST} , Tajima's D, and P_{seg}) are also shown in Box 2 (main text), but are retained
 204 here for ease of comparison.



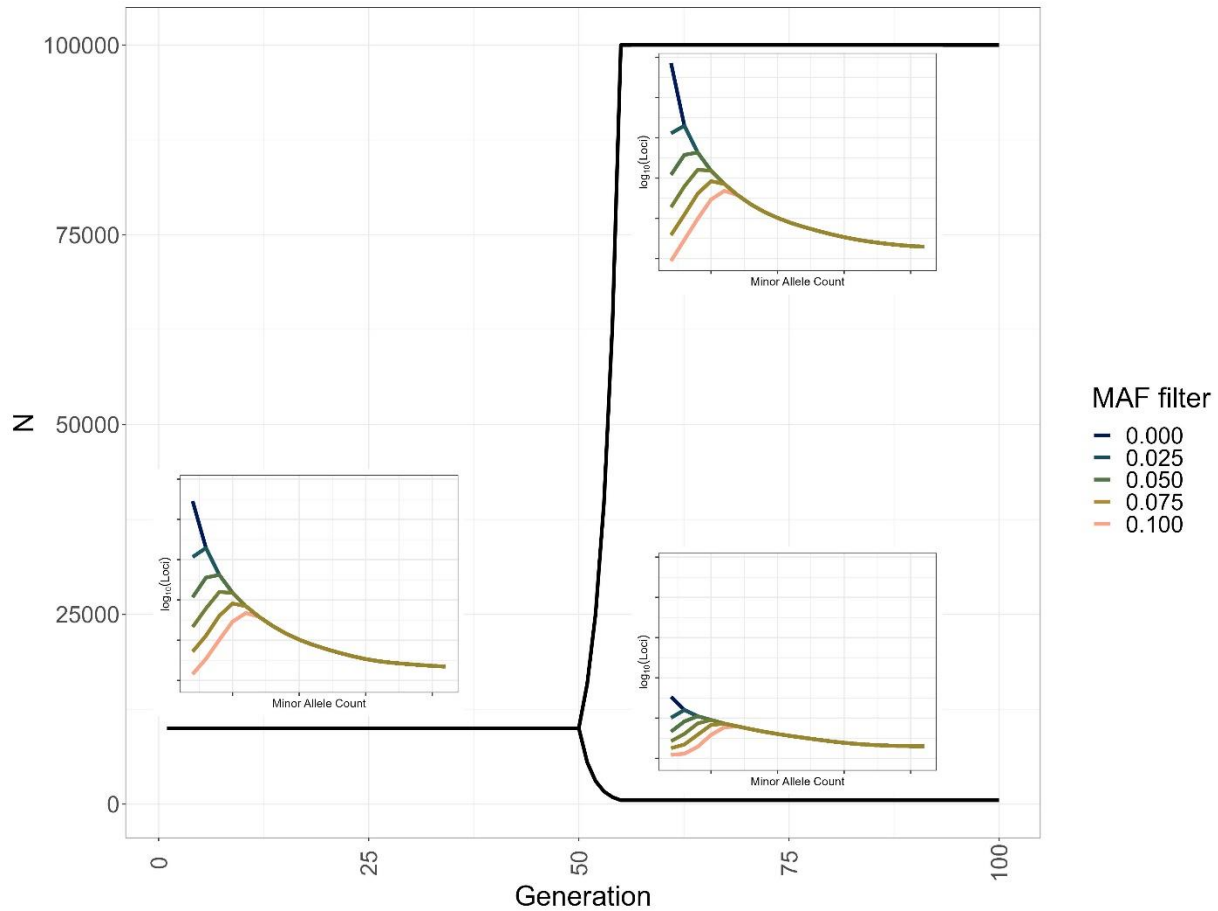
205
 206 **Figure S3:** Effects of **HWP filtering** threshold on standardized F_{ST} , H_o , F_{IS} , Tajima's D , H_e , π ,
 207 Watterson's θ , Tajima's θ , the proportion of the total (highest) number of segregating sites, the
 208 proportion of the total (highest) number of private alleles, and N_e . Parameter values have been
 209 normalised to show change by subtracting off the value observed with no filter. Note that RAD
 210 datasets tend to see more changes with increasing HWP filter stringency for many statistics,
 211 including substantial Tajima's D and F_{IS} changes in some cases.



212
 213 **Figure S4:** Effects of different thresholds for **filtering out loci based on missing data** on F_{ST} ,
 214 H_o , F_{IS} , Tajima's D, H_e , π , Watterson's θ , Tajima's θ , the proportion of the total (highest)
 215 number of segregating sites, the proportion of the total (highest) number of private alleles, and
 216 N_e , observed in each population and dataset depending on missing data filtering stringency. Loci
 217 were removed if genotyped in too few individuals. Parameter values have been normalised to
 218 show change by subtracting off the value observed with no filter. Note that the deer mouse
 219 exome sequencing data in particular shows very large changes across most statistics with stricter
 220 missing data filters, including a very large decrease in Tajima's D.

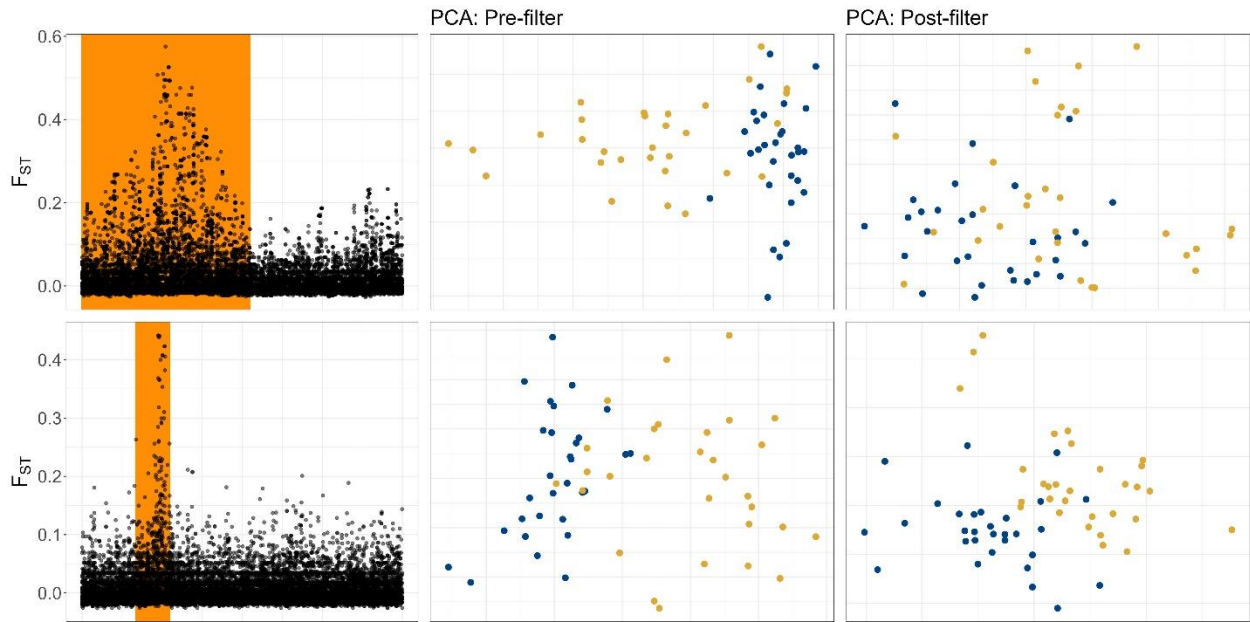


221
 222 **Figure S5: Effects of different thresholds for filtering out individuals based on missing data**
 223 on F_{ST} , H_o , F_{IS} , Tajima's D , H_e , π , Watterson's θ , Tajima's θ , the proportion of the total
 224 (highest) number of segregating sites, the proportion of the total (highest) number of private
 225 alleles, and N_e , observed in each population and dataset depending on missing data filtering
 226 stringency. Individuals were removed if genotyped at too few loci. Parameter values have been
 227 normalised to show change by subtracting off the value observed with no filter. Note that the
 228 *Daphnia* and yellow perch data show the largest change for most statistics, including substantial
 229 changes in Tajima's D , F_{IS} , and H_o/H_e .



230
 231
 232
 233
 234
 235
 236
 237
 238

Figure S6: Effect of **MAF filtering** on the site frequency spectrum depending on demographic history. Primary plot displays demographic history, insets show the effect of different MAF filters on the site frequency spectra of neutral, historically expanded, and historically bottlenecked populations. Filtering for MAF impacts “flat” site frequency spectra, like those that have undergone bottlenecks, much less strongly than those with many rare alleles, like in those that have undergone expansions. Note that the y-axes in each spectra plot are scaled to the same minimum and maximum values. Full demographic model parameters are available in Supplementary Table 3.



239
 240 **Figure S7:** Effect of **filtering out high F_{ST} regions** on subsequent principal component analyses
 241 depending on local linkage patterns. Selection or other factors which generate localized sections
 242 of elevated F_{ST} can have a strong effect on population structure estimates, particularly in areas
 243 with lower recombination rates and/or higher rates of linkage. Blocks of elevated F_{ST} (left,
 244 marked in orange) were generated via simulated selection according to the parameters noted in
 245 Supplementary Table 3 using either a recombination rate of 0.1 (top) or 1 (bottom).