1    **Title:** Genomic evolutionary analysis in R with geaR.

2    Christopher M. Ward[1*], Alastair J. Ludington[1], James Breen[2,3,4], Simon W. Baxter[5]

3    [1] School of Biological Sciences, University of Adelaide, Australia

4    [2] Bioinformatics Platform, South Australian Health & Medical Research Institute (SAHMRI), Australia

5    [3] Robinson Research Institute, University of Adelaide, Australia

6    [4] Faculty of Health & Medical Sciences, University of Adelaide, Australia

7    [5] Bio21 Institute, School of BioSciences, University of Melbourne, Australia

8    * Corresponding author: christopher.ward@adelaide.edu.au

9

10   **Abstract:**

11   The analysis and interpretation of datasets generated through sequencing large numbers of

12   individual genomes is becoming commonplace in population and evolutionary genetic studies. Here

13   we introduce geaR, a modular R package for evolutionary analysis of genome-wide genotype data.

14   The package leverages the Genomic Data Structure (GDS) format, which enables memory and time

15   efficient querying of genotype datasets compared to standard VCF genotype files. geaR utilizes

16   GRange object classes to partition an analysis based on features from GFF annotation files, select

17   codons based on position or degeneracy, and construct both positional and coordinate genomic

18   windows. Tests of genetic diversity (eg. $d_{XY}$, $\pi$, $F_{ST}$) and admixture ($f_4$, $\widehat{f_d}$) along with tree building and

19   sequence output, can be carried out on partitions using a single function regardless of sample ploidy

20   or number of observed alleles. The package and associated documentation are available on GitHub

21   at https://github.com/CMWbio/geaR.

22   Keywords: Evolution, Population Genomics, R Package, Admixture

23   **Introduction:**

24   Improvements in genome sequencing technologies has led to increased production of data at lower

25   relative cost per base (Schwarze et al. 2020). Genome-wide sequencing datasets with hundreds of

26   samples can be produced for population genomic analysis, allowing researchers to investigate

1

27 population and evolutionary history at an unprecedented scale. However, due to file size and data

28 complexity downstream problems during data storage and analysis can arise. The most common

29 format for handling genome-wide SNP data is the Variant Call Format (VCF), which has historically

30 had a large memory overhead when being read into an R environment. To resolve this, the Genomic

31 Data Structure (GDS) format has allowed all genotype and metadata to be compressed into a

32 queriable, on-disk file that substantially reduce memory requirements and decrease analysis time

33 (Zheng et al. 2017). The GDS format provides an efficient format for filtering SNP data in order to

34 perform Principal Component Analysis, estimate genetic relatedness and tests for genetic

35 association (Zheng et al. 2012).

36 GDS files use GRange objects from the GenomicRanges package (Lawrence et al. 2013) to define loci

37 to query from file and import into R. In their most basic form, GRange class objects define genomic

38 loci based on reference position. Although widely used throughout Bioconductor, GRange objects, to

39 our knowledge, have not been utilized in the same manner to define loci for evolutionary analyses.

40 Few R packages attempt to carry out genome-wide investigation of genotype data. Most packages

41 focus on the analysis of single or multi-locus data, with the notable exception of PopGenome (Pfeifer

42 et al. 2014). However, one limitation of PopGenome is customizability of how the target genome is

43 partitioned, and which sites are selected for analysis. Most tools, including PopGenome, allow

44 datasets to be partitioned into sliding or tiled windows based on reference or SNP position.

45 PopGenome also provides methods to split data into GFF attributes, however selection of bespoke

46 partitions not possible. This makes calculating population metrics on specific codon positions (eg.

47 four-fold or zero-fold degenerate sites) or analysing many non-contiguous loci difficult and time

48 consuming.

49 To overcome these issues, here we present the R package geaR, which leverages the GDS format, to

50 efficiently construct GRanges containing genome-wide or local loci of interest and to carry out

51 common tasks for evolutionary analysis on genome-wide genotype data using a single function.

2

52    Furthermore, we provide methods to partition the genome based on annotation, codon position or

53    degeneracy through utilizing data in GFF files, or using reference genome coordinates or genotype

54    position.

55    **Features:**

56    *Input data*

57    Genotype input files are required to be in GDS format, enabling high compressibility compared to

58    gzipped VCFs (>5X smaller on disk), efficient querying and the capability to work on large datasets

59    with a reduced memory footprint (Zheng et al. 2017). Conversion of sample genotypes in the VCF

60    format to a Genomic Data Structure (GDS) format can be performed using the SeqArray package

61    (Zheng et al. 2017) before analysis with geaR. Genotypes called at any level of ploidy can be utilized

62    in geaR, which includes whole genome sequence data generated from pools of two or more

63    individuals.

64    *Partitioning the genome using GRanges*

65    The geaR package utilizes GRange objects to define partitions for the analysis, for example,

66    segmenting a genome into 10-kb windows. This allows users to define their own GRanges for the

67    analysis or build them with provided functions. Currently, users are able to generate both coordinate

68    (based on reference coordinate) and positional (based on genotype number) windows using

69    *makeWindows()* or *makeSnpWindows()* functions. Sequence features, such as protein coding

70    regions, can be extracted from a GFF with *getFeatures().*

71    Many evolutionary analyses seek to calculate population metrics over different codon positions. To

72    make this as simple as possible, geaR provides methods to index a reference genome according to

73    codon position with *buildCodonDB(),* which can either be stored in memory as a GRangesList object

74    or an SQLite database (DB) on disk to limit static memory usage. Users may then filter codons based

75    on degeneracy (0-fold or 2-fold) and position using the function *filterCodonDB().* A codon DB can also

3

76    be passed to the function *validate4fold()* to select 4-fold degenerate sites across the genome that

77    are empirically supported in the GDS file. This is done by querying the GDS to i) remove codons with

78    missing data, ii) select 4-fold degenerate codons, iii) remove all those where codon positions one or

79    two have variation and iv) select third positions.

80    GRange objects generated using geaR can then be combined using *mergeLoci()* to further customize

81    partitions. For example, genome-wide tiled windows can be combined with four-fold degenerate

82    sites to output either genomic windows that contain only 4-fold sites or all sites excluding 4-fold
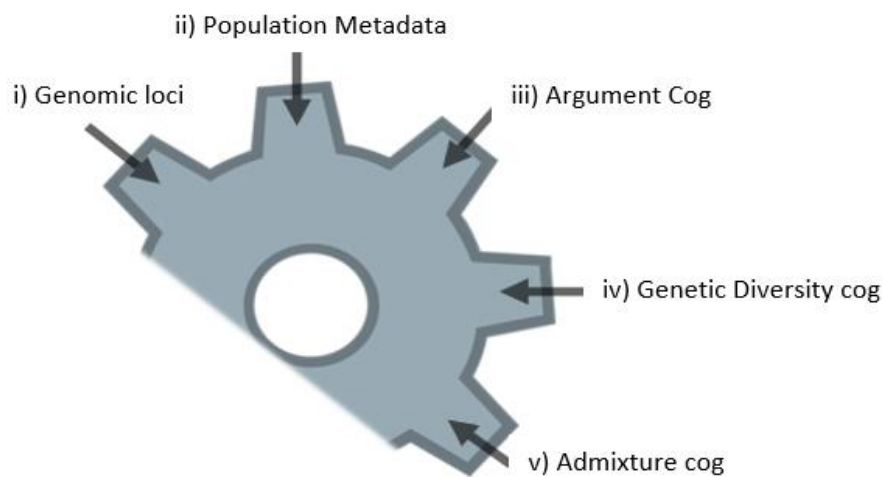
83    degenerate sites.



84    Figure 1: Structure of the of the gear S4 object: i) Genomic loci (GRanges) to carry out analyses
85    across, ii) Population metadata encoding sample names to the population/species they belong to, iii)
86    A cog containing general arguments for all analysis, iv) a cog specifying that the Genetic Diversity
87    module should be carried out and v) a cog specifying that the Admixture cog should be carried out
88    on the dataset.

89

90    ***Setting up an analysis: cogs and gears***

91    geaR operates through two S4 classes, the 'cog' and 'gear' (Figure 1). Cogs, built using *makeCog(),*

92    specify multiple analyses to carry out (see Table 1) setting parameters specific to each analysis. A

93    single gear object can then be constructed, using *makeGear(),* which contains all of the specified

94    cogs for analysis, along with the genomic loci and population metadata (Figure 2A). The

95   *analyzeGear()* function then performs all analyses on the same set of genomic loci and samples,

96   greatly reducing run time compared to sequential execution.
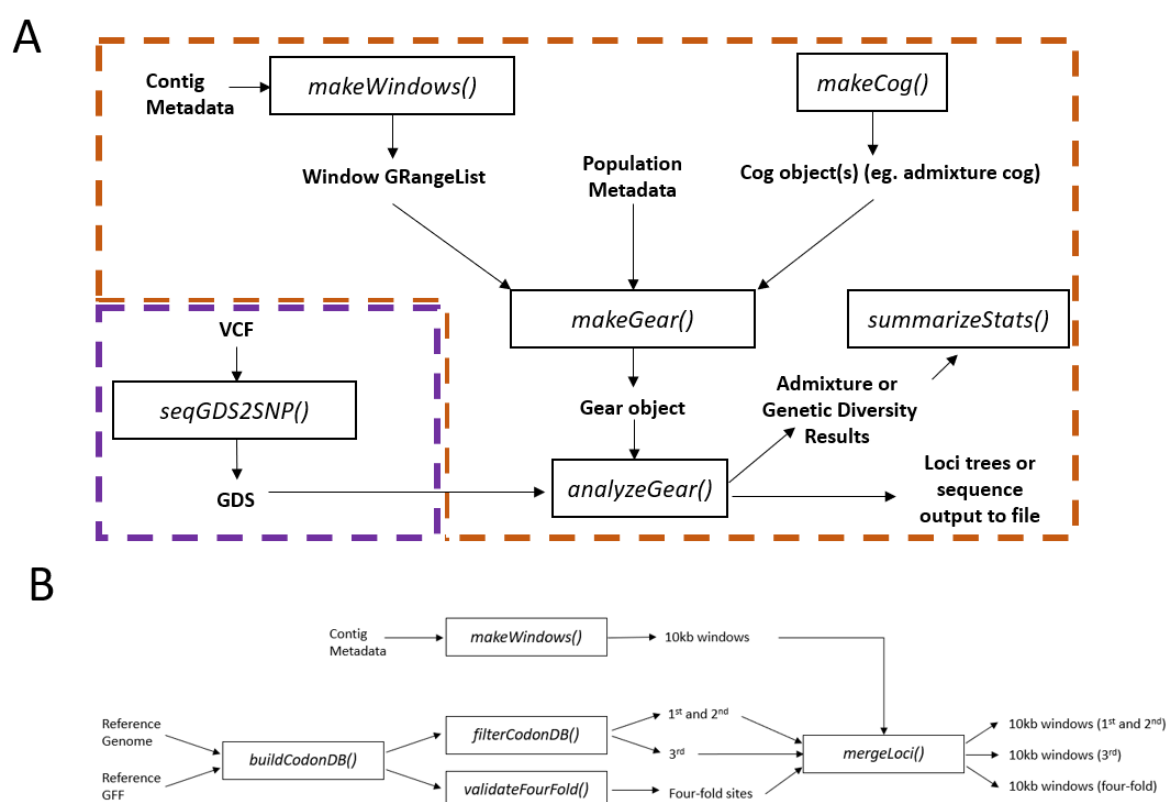


98   Figure 2: A) Basic analysis workflow in geaR to carry out analysis on windowed genomic loci.
99   Functions specific to geaR are within the orange boundary and external functions in purple. After
100  converting the VCF to GDS format using SeqArray, contig metadata (contig length) is used to
101  construct windows across the genome. A dataframe containing population metadata defining
102  population to sample grouping is then constructed. This is used, along with windows for the analysis
103  and cogs, to construct the gear class object. The analysis is then carried out on the gear object and
104  outputs depend on which cogs were specified. B) Workflow used to generate partition schemes for
105  examples in Figure 3. First 10kb windows were generated from contig metadata. This was followed
106  by generation of a codon database that indexes codon position in a reference genome. The codon
107  database was then passed to the function *filterCodonDB* to output separate loci-sets for four codon
108  partition schemes: $1^{st}+2^{nd}$; $3^{rd}$; 0-fold and 2-fold. *validateFourFold()* was also used to select 4-fold
109  degenerate sites that are supported by genotypes in the GDS file. Each of these codon loci-sets can
110  then be passed to the function *mergeLoci(),* along with the 10kb windows, to combine loci into 10kb
111  windows that contain only the selected codon types.

112

113

114    *Analysis types*

115    Four different cogs can be generated to carry out an analysis: i) genetic diversity, ii) admixture, iii)

116    outputLoci and iv) outputTrees. Genetic diversity allows the calculation of a range of population

117    metrics (Table 1), most of which rely on genetic distance which is calculated based on the hamming

118    distance between haplotypes at all sites within the locus. The admixture cog utilizes outgroup

119    polarized allele frequency at all biallelic sites within the locus to calculate $f_4$ (Patterson et al. 2012)

120    and $\widehat{f_d}$ (Martin et al. 2015) statistics. The package also enables users to output data in fasta format

121    for each individual (or sample pool) using outputLoci or as distance trees using outputTrees.

122    Haplotypes are used in diversity calculations and are output to file according to the phase within the

123    supplied GDS file, not calculated by geaR.

124    Outputs of both genetic diversity and admixture cogs can be summarized using *summarizeStats()*

125    which calculates a mean and median values for each statistic across all loci using a block jack-knife

126    approach.

127    Table 1: Analysis types and functionality available to apply at each locus.

| Cog type | Functionality |
|---|---|
| **Genetic diversity** | Nucleotide diversity ($\pi$), genetic distance ($d_{XY}$), maximum distance ($d_{max}$), minimum distance ($d_{min}$), ancestral distance ($d_a$), $\gamma_{ST}$, relative node distance (RND), minimum relative node distance ($RND_{min}$), and $G_{min}$ |
| **Admixture** | $f_4$ and $\widehat{f_d}$ |
| **Output loci** | *fasta* format output to file |
| **Output trees** | *newick* format distance trees and metadata output to file |

128

129

130

131    *Parallelization*

132    All functions allow operations to be run in parallel by leveraging methods in the furrr

133    (https://github.com/DavisVaughan/furrr) and parallel R packages.

134    **Carrying out an analysis using gear:**

135    geaR has successfully been used to calculate genome wide diversity metrics between populations

136    containing 532 moth genomes (You et al. 2020) and to identify introgressed regions between two

137    *Bactrocera* fly species (Ward et al. 2020). Below we outline two example analyses using geaR. Code

138    for each of examples and other common workflows can be found on the wiki

139    (https://github.com/CMWbio/geaR/wiki).

140    In our first example we use a subset of the data from Ward and Baxter (2018) containing three

141    populations of diamondback moth collected from Australian Capital Territory, Australia; South

142    Australia, Australia and Hawaii, USA. Second. Following the general workflow shown in Figure 2A, we

143    converted the called genotypes to the GDS format using SeqArray, constructed partitions, built cogs,

144    combined those cogs into a gear and then carried out the analysis. We constructed our partitions

145    for the analysis by generating a GRange object containing only scaffold_4. The analysis will use this

146    GRange object to construct six different partition schemes based on 10kb tiled windows (workflow

147    shown in Figure 2B): i) all sites, ii) only $1^{st}+2^{nd}$ codon positions, iii) windows only $3^{rd}$ codon positions,

148    iv) only 0-fold degenerate sites, v) only 2-fold degenerate sites and vi) only 4-fold degenerate sites.

149    Partition i) was then used to calculate pairwise genetic distance ($d_{XY}$) between each population

150    across the scaffold (Figure 3A) and partitions ii-vi) were used to calculate within population

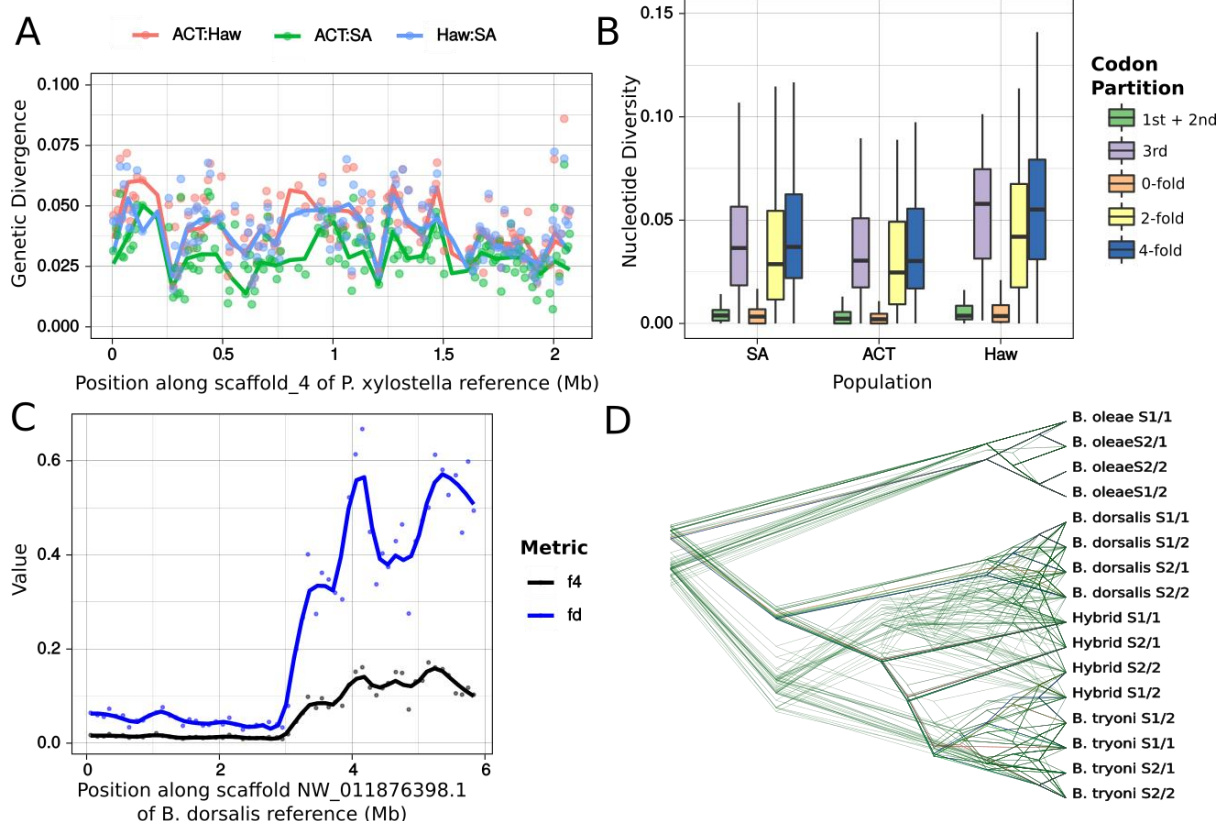151    nucleotide diversity across the whole genome (Figure 3B).

152



153

Figure 3: Example workflows to carry out with geaR: Panels A and B use *P. xylostella* data from Ward and Baxter (2018), panels C and D use *Bactrocera* from Ward *et al* (2020). A) Absolute genetic distance ($d_{XY}$) was calculated between pairwise comparisons of three populations for 10kb tiled windows across scaffold_4 of the diamondback moth reference genome. C) The five loci-sets constructed using the workflow in Figure 2B were used to calculate nucleotide diversity (π) at $1^{st}+2^{nd}$, $3^{rd}$, 0-fold, 2-fold and 4-fold codon sites across scaffold_4 of the diamondback moth reference genome C) Admixture metrics $f_4$ and $\widehat{f_d}$ calculated on 100kb windows across scaffold NW_011876398.1 of the *B. dorsalis* reference genome. D) Distance trees for each 100kb window across NW_011876398.1 output using the *outputTrees* cog showing a mixture of discordant and concordant topologies. Plots A), B) and C) were generated using ggplot2 (Wickham 2009) and D) using densitree (Bouckaert 2010).

For a second example we will identify one of the introgressed regions from Ward et al. (2020). This will use data from a single scaffold (NW_011876398.1) of the *B. dorsalis* reference genome (GCF_000789215.1) for two samples of *B. tryoni, B. dorsalis, B. oleae* and a *B. dorsalis/B.tryoni* hybrid line. Using the same methodology as the first example, we constructed a 100kb tiled window partition scheme. However, for this analysis we used the admixture cog to calculate $f_4$ and $\widehat{f_d}$ admixture metrics showing clear evidence for introgression at the 3' end of the scaffold (Figure 3C).

171    We also used the outputTrees cog to output distance trees for each of these windows to illustrate

172    both the congruent and incongruent topologies resulting from partial admixture on

173    NW_011876398.1 (Figure 3D).

174

175    **Conclusion**

176    Genome-wide datasets with many individuals are becoming the norm in population genetic studies,

177    increasing the need for tools to efficiently carry out analyses on genotype data. The functional

178    programming capabilities of the R programming language provide an intuitive environment for users

179    to carry out calculation and visualization of population and evolutionary genomics metrics. The

180    methods provided in geaR allow users easily and effectively partition the genome for generic and

181    bespoke analysis of genome-wide genotype data regardless of sample ploidy and number of

182    observed alleles.

183

184    **Acknowledgements**

188    **Data Availability**

189    All data is available in the referenced publications.

190

191

192    **References**

193    Bouckaert, R. R. 2010. DensiTree: making sense of sets of phylogenetic trees. Bioinformatics
194         **26**:1372-1373.

195 Lawrence, M., W. Huber, H. Pagès, P. Aboyoun, M. Carlson, R. Gentleman, M. T. Morgan, and V. J.
196          Carey. 2013. Software for computing and annotating genomic ranges. PLoS computational
197          biology **9**:e1003118-e1003118.
198 Martin, S. H., J. W. Davey, and C. D. Jiggins. 2015. Evaluating the use of ABBA-BABA statistics to
199          locate introgressed loci. Molecular biology and evolution **32**:244-257.
200 Patterson, N., P. Moorjani, Y. Luo, S. Mallick, N. Rohland, Y. Zhan, T. Genschoreck, T. Webster, and D.
201          Reich. 2012. Ancient Admixture in Human History. Genetics **192**:1065.
202 Pfeifer, B., U. Wittelsbürger, S. E. Ramos-Onsins, and M. J. Lercher. 2014. PopGenome: an efficient
203          Swiss army knife for population genomic analyses in R. Molecular biology and evolution
204          **31**:1929-1936.
205 Schwarze, K., J. Buchanan, J. M. Fermont, H. Dreau, M. W. Tilley, J. M. Taylor, P. Antoniou, S. J. L.
206          Knight, C. Camps, M. M. Pentony, E. M. Kvikstad, S. Harris, N. Popitsch, A. T. Pagnamenta, A.
207          Schuh, J. C. Taylor, and S. Wordsworth. 2020. The complete costs of genome sequencing: a
208          microcosting study in cancer and rare diseases from a single center in the United Kingdom.
209          Genetics in Medicine **22**:85-94.
210 Ward, C., R. Aumann, M. Whitehead, K. Nikolouli, G. Leveque, G. Gouvi, E. Fung, S. Reiling, H.
211          Djambazian, M. Hughes, S. Whiteford, C. Caceres-Barrios, T. Nguyen, A. Choo, P. Crisp, S.
212          Sim, S. Geib, F. Marec, I. Häcker, J. Ragoussis, A. Darby, K. Bourtzis, S. Baxter, and M.
213          Schetelig. 2020. White pupae genes in the Tephritids Ceratitis capitata, Bactrocera dorsalis
214          and Zeugodacus cucurbitae: a story of parallel mutations. BioRxiv:2020.2005.2008.076158.
215 Ward, C. M., and S. W. Baxter. 2018. Assessing Genomic Admixture between Cryptic Plutella Moth
216          Species following Secondary Contact. Genome biology and evolution **10**:2973-2985.
217 Wickham, H. 2009. ggplot2: Elegant Graphics for Data Analysis. Springer Publishing Company,
218          Incorporated.
219 You, M., F. Ke, S. You, Z. Wu, Q. Liu, W. He, S. W. Baxter, Z. Yuchi, L. Vasseur, G. M. Gurr, C. M. Ward,
220          H. Cerda, G. Yang, L. Peng, Y. Jin, M. Xie, L. Cai, C. J. Douglas, M. B. Isman, M. S. Goettel, Q.
221          Song, Q. Fan, G. Wang-Pruski, D. C. Lees, Z. Yue, J. Bai, T. Liu, L. Lin, Y. Zheng, Z. Zeng, S. Lin,
222          Y. Wang, Q. Zhao, X. Xia, W. Chen, L. Chen, M. Zou, J. Liao, Q. Gao, X. Fang, Y. Yin, H. Yang, J.
223          Wang, L. Han, Y. Lin, Y. Lu, and M. Zhuang. 2020. Variation among 532 genomes unveils the
224          origin and evolutionary history of a global insect herbivore. Nature Communications
225          **11**:2321.
226 Zheng, X., S. M. Gogarten, M. Lawrence, A. Stilp, M. P. Conomos, B. S. Weir, C. Laurie, and D. Levine.
227          2017. SeqArray—a storage-efficient high-performance data format for WGS variant calls.
228          Bioinformatics **33**:2251-2257.
229 Zheng, X., D. Levine, J. Shen, S. M. Gogarten, C. Laurie, and B. S. Weir. 2012. A high-performance
230          computing toolset for relatedness and principal component analysis of SNP data.
231          Bioinformatics **28**:3326-3328.

232