

# Package 'BIGr'

September 18, 2025

**Title** Breeding Insight Genomics Functions for Polyploid and Diploid Species

**Version** 0.6.2

**Maintainer** Alexander M. Sandercock <sandercock.alex@gmail.com>

**Description** Functions developed within Breeding Insight to analyze diploid and polyploid breeding and genetic data. 'BIGr' provides the ability to filter variant call format (VCF) files, extract single nucleotide polymorphisms (SNPs) from diversity arrays technology missing allele discovery count (DArT MADC) files, and manipulate genotype data for both diploid and polyploid species. It also serves as the core dependency for the 'BIGapp' 'Shiny' app, which provides a user-friendly interface for performing routine genotype analysis tasks such as dosage calling, filtering, principal component analysis (PCA), genome-wide association studies (GWAS), and genomic prediction. For more details about the included 'breedTools' functions, see Funkhouser et al. (2017) <doi:10.2527/tas2016.0003>, and the 'updog' output format, see Gerard et al. (2018) <doi:10.1534/genetics.118.301468>.

**License** Apache License (>= 2)

**URL** <https://github.com/Breeding-Insight/BIGr>

**BugReports** <https://github.com/Breeding-Insight/BIGr/issues>

**Encoding** UTF-8

**RoxygenNote** 7.3.2

**Depends** R (>= 4.4.0)

**Imports** parallel, dplyr, Rdpack (>= 0.7), readr (>= 2.1.5), reshape2 (>= 1.4.4), rlang, tidyr (>= 1.3.1), vcfR (>= 1.15.0), Rsamtools, Biostrings, pwalgn, janitor, quadprog, tibble, stringr

**Suggests** covr, spelling, rmdformats, knitr (>= 1.10), rmarkdown, testthat (>= 3.0.0)

**RdMacros** Rdpack

**NeedsCompilation** no

**Author** Alexander M. Sandercock [cre, aut],  
 Cristiane Taniguti [aut],  
 Josue Chinchilla-Vargas [aut],  
 Shufen Chen [ctb],  
 Manoj Sapkota [ctb],  
 Meng Lin [ctb],  
 Dongyan Zhao [ctb],  
 Cornell University [cph] (Breeding Insight)

**Repository** CRAN

**Date/Publication** 2025-09-18 12:40:13 UTC

## Contents

allele_freq_poly . . . . .	2
calculate_Het . . . . .	3
calculate_MAF . . . . .	4
check_homozygous_trios . . . . .	5
check_ped . . . . .	7
check_replicates . . . . .	8
dosage2vcf . . . . .	9
dosage_ratios . . . . .	10
filterMADC . . . . .	11
filterVCF . . . . .	12
flip_dosage . . . . .	14
get_countsMADC . . . . .	15
imputation_concordance . . . . .	15
madc2gmat . . . . .	17
madc2vcf_all . . . . .	18
madc2vcf_targets . . . . .	20
merge_MADCs . . . . .	21
solve_composition_poly . . . . .	22
thinSNP . . . . .	24
updog2vcf . . . . .	25

**Index** **28**

---

allele\_freq\_poly      *Compute Allele Frequencies for Populations*

---

### Description

Computes allele frequencies for specified populations given SNP array data

### Usage

```
allele_freq_poly(geno, populations, ploidy = 2)
```

**Arguments**

geno	matrix of genotypes coded as the dosage of allele B {0, 1, 2, ..., ploidy} with individuals in rows (named) and SNPs in columns (named)
populations	list of named populations. Each population has a vector of IDs that belong to the population. Allele frequencies will be derived from all animals
ploidy	integer indicating the ploidy level (default is 2 for diploid)

**Value**

data.frame consisting of allele\_frequencies for populations (columns) for each SNP (rows)

**References**

Funkhouser SA, Bates RO, Ernst CW, Newcom D, Steibel JP. Estimation of genome-wide and locus-specific breed composition in pigs. *Transl Anim Sci.* 2017 Feb 1;1(1):36-44.

**Examples**

```
# Example inputs
geno_matrix <- matrix(
  c(4, 1, 4, 0, # S1
    2, 2, 1, 3, # S2
    0, 4, 0, 4, # S3
    3, 3, 2, 2, # S4
    1, 4, 2, 3), # S5
  nrow = 4, ncol = 5, byrow = FALSE, # individuals=rows, SNPs=cols
  dimnames = list(paste0("Ind", 1:4), paste0("S", 1:5))
)

pop_list <- list(
  PopA = c("Ind1", "Ind2"),
  PopB = c("Ind3", "Ind4")
)

allele_freqs <- allele_freq_poly(geno = geno_matrix, populations = pop_list, ploidy = 4)
print(allele_freqs)
```

---

 calculate\_Het

---

*Calculate Observed Heterozygosity from a Genotype Matrix*


---

**Description**

This function calculates the observed heterozygosity from a genotype matrix. It assumes that the samples are the columns, and the genomic markers are in rows. Missing data should be set as NA, which will then be ignored for the calculations. All samples must have the same ploidy.

**Usage**

```
calculate_Het(geno, ploidy)
```

**Arguments**

geno	Genotype matrix or data.frame
ploidy	The ploidy of the species being analyzed

**Value**

A dataframe of observed heterozygosity values for each sample

**Examples**

```
# example input for a diploid
geno <- data.frame(
  Sample1 = c(0, 1, 2, NA, 0),
  Sample2 = c(1, 1, 2, 0, NA),
  Sample3 = c(0, 1, 1, 0, 2),
  Sample4 = c(0, 0, 1, 1, NA)
)
row.names(geno) <- c("Marker1", "Marker2", "Marker3", "Marker4", "Marker5")

ploidy <- 2

# calculate observed heterozygosity
result <- calculate_Het(geno, ploidy)

print(result)
```

---

calculate\_MAF

*Calculate Minor Allele Frequency from a Genotype Matrix*

---

**Description**

This function calculates the allele frequency and minor allele frequency from a genotype matrix. It assumes that the Samples are the columns, and the genomic markers are in rows. Missing data should be set as NA, which will then be ignored for the calculations. All samples must have the same ploidy.

**Usage**

```
calculate_MAF(df, ploidy)
```

**Arguments**

df	Genotype matrix or data.frame
ploidy	The ploidy of the species being analyzed

**Value**

A dataframe of AF and MAF values for each marker

**Examples**

```
# example input for a diploid
geno <- data.frame(
  Sample1 = c(0, 1, 2, NA, 0),
  Sample2 = c(1, 1, 2, 0, NA),
  Sample3 = c(0, 1, 1, 0, 2),
  Sample4 = c(0, 0, 1, 1, NA)
)
row.names(geno) <- c("Marker1", "Marker2", "Marker3", "Marker4", "Marker5")

ploidy <- 2

# calculate allele frequency
result <- calculate_MAF(geno, ploidy)

print(result)
```

---

check\_homozygous\_trios

*Check Homozygous Loci in Trios*

---

**Description**

This function analyzes homozygous loci segregation in trios (parents and progeny) using genotype data from a VCF file. It calculates the percentage of homozygous loci in the progeny that match the expected segregation patterns based on the tested parents.

**Usage**

```
check_homozygous_trios(
  path.vcf,
  ploidy = 4,
  parents_candidates = NULL,
  progeny_candidates = NULL,
  verbose = TRUE
)
```

**Arguments**

`path.vcf` A string specifying the path to the VCF file containing genotype data.

`ploidy` An integer specifying the ploidy level of the samples. Default is 4.

`parents_candidates` A character vector of parent sample names to be tested. Must be provided.

progeny_candidates	A character vector of progeny sample names to be tested. Must be provided.
verbose	A logical value indicating whether to print the number of combinations tested. Default is TRUE.

### Details

This function is designed to validate the segregation of homozygous loci in trios, ensuring that the progeny genotypes align with the expected patterns based on the parental genotypes. It requires both parent and progeny candidates to be specified. The function validates the ploidy level and ensures that all specified samples are present in the VCF file. The results include detailed statistics for each combination of parents and progeny. Reciprocal comparisons (e.g., A vs. B and B vs. A) and self-comparisons (e.g., A vs. A) are removed to avoid redundancy. Missing genotype data is also accounted for and reported in the results.

### Value

A data frame with the following columns:

- parent1: The name of the first parent in the pair.
- parent2: The name of the second parent in the pair.
- progeny: The name of the progeny sample.
- homoRef\_x\_homoRef\_n: Number of loci where both parents are homozygous reference.
- homoRef\_x\_homoRef\_match: Percentage of matching loci in the progeny for homozygous reference parents.
- homoAlt\_x\_homoAlt\_n: Number of loci where both parents are homozygous alternate.
- homoAlt\_x\_homoAlt\_match: Percentage of matching loci in the progeny for homozygous alternate parents.
- homoRef\_x\_homoAlt\_n: Number of loci where one parent is homozygous reference and the other is homozygous alternate.
- homoRef\_x\_homoAlt\_match: Percentage of matching loci in the progeny for mixed homozygous parents.
- homoalt\_x\_homoRef\_n: Number of loci where one parent is homozygous alternate and the other is homozygous reference.
- homoalt\_x\_homoRef\_match: Percentage of matching loci in the progeny for mixed homozygous parents (alternate-reference).
- missing: The number of loci with missing genotype data in the comparison.

### Examples

```
# Example VCF file
example_vcf <- system.file("iris_DArT_VCF.vcf.gz", package = "BIGr")

parents_candidates <- paste0("Sample_",1:10)
progeny_candidates <- paste0("Sample_",11:20)

#Check homozygous loci in trios
```

```
check_tab <- check_homozygous_trios(path.vcf = example_vcf,  
                                   ploidy = 2,  
                                   parents_candidates = parents_candidates,  
                                   progeny_candidates = progeny_candidates)
```

---

check\_ped

*Evaluate Pedigree File for Accuracy*

---

### Description

Check a pedigree file for accuracy and output suspected errors

### Usage

```
check_ped(ped.file, seed = NULL, verbose = TRUE)
```

### Arguments

ped.file	path to pedigree text file. The pedigree file is a 3-column pedigree tab separated file with columns labeled as id sire dam in any order
seed	Optional seed for reproducibility
verbose	Logical. If TRUE, print the errors to the console.

### Details

check\_ped takes a 3-column pedigree tab separated file with columns labeled as id sire dam in any order and checks for:

- Ids that appear more than once in the id column
- Ids that appear in both sire and dam columns
- Direct (e.g. parent is a offspring of his own daughter) and indirect (e.g. a great grandparent is son of its grandchild) dependencies within the pedigree.
- Individuals included in the pedigree as sire or dam but not on the id column and reports them back with unknown parents (0).

When using check\_ped, do a first run to check for repeated ids and parents that appear as sire and dam. Once these errors are cleaned run the function again to check for dependencies as this will provide the most accurate report.

Note: This function does not change the input file but prints any errors found in the console.

### Value

A list of data.frames of error types, and the output printed to the console

**Examples**

```
##Get list with a dataframe for each error type
ped_file <- system.file("check_ped_test.txt", package="BIGr")
ped_errors <- check_ped(ped.file = ped_file,
                        seed = 101919)

##Access the "messy parents" dataframe result
ped_errors$messy_parents

##Get list of sample IDs with messy parents error
messy_parent_ids <- ped_errors$messy_parents$id
print(messy_parent_ids)
```

---

check\_replicates      *Compatibility Between Samples Genotypes*

---

**Description**

This function checks the compatibility between sample genotypes in a VCF file by comparing all pairs of samples.

**Usage**

```
check_replicates(path.vcf, select_samples = NULL, verbose = TRUE)
```

**Arguments**

path.vcf	A string specifying the path to the VCF file containing genotype data.
select_samples	An optional character vector of sample names to be selected for comparison. If NULL (default), all samples in the VCF file are used.
verbose	A logical value indicating whether to print the number of combinations tested. Default is TRUE.

**Details**

The function removes reciprocal comparisons (e.g., A vs. B and B vs. A) and self-comparisons (e.g., A vs. A) to avoid redundancy. Compatibility is calculated as the percentage of matching genotypes between two samples, excluding missing values. The percentage of missing genotypes is also reported for each pair.

**Value**

A data frame with four columns:

- sample1: The name of the first sample in the pair.
- sample2: The name of the second sample in the pair.
- %\_matching\_genotypes: The percentage of compatible genotypes between the two samples.
- %\_missing\_genotypes: The percentage of missing genotypes in the comparison.

## Examples

```
#Example VCF
example_vcf <- system.file("iris_DArT_VCF.vcf.gz", package = "BIGr")

# Checking for replicates
check_tab <- check_replicates(path.vcf = example_vcf, select_samples = NULL)
```

---

dosage2vcf

*Convert DArTag Dosage and Counts to VCF*

---

## Description

This function will convert the DArT Dosage Report and Counts files to VCF format

## Usage

```
dosage2vcf(dart.report, dart.counts, ploidy, output.file)
```

## Arguments

dart.report	Path to the DArT dosage report .csv file. Typically contains "Dosage Report" in the file name.
dart.counts	Path to the DArT counts .csv file. Typically contains "Counts" in the file name.
ploidy	The ploidy of the species being analyzed
output.file	output file name and path

## Details

This function will convert the Dosage Report and Counts files from DArT into a VCF file. These two files are received directly from DArT for a given sequencing project. The output file will be saved to the location and with the name that is specified. The VCF format is v4.3

## Value

A vcf file

## Examples

```
## Use file paths for each file on the local system

#The files are directly from DArT for a given sequencing project.
#The are labeled with Dosage_Report or Counts in the file names.

#Temp location (only for example)
output_file <- tempfile()
```

```

dosage2vcf(dart.report = system.file("iris_DaRT_Allele_Dose_Report_small.csv", package = "BIGr"),
           dart.counts = system.file("iris_DaRT_Counts_small.csv", package = "BIGr"),
           ploidy = 2,
           output.file = output_file)

# Removing the output for the example
rm(output_file)

##The function will output the converted VCF using information from the DaRT files

```

---

dosage\_ratios

*Calculate the Percentage of Each Dosage Value*


---

### Description

This function calculates the percentage of each dosage value within a genotype matrix. It assumes that the samples are the columns, and the genomic markers are in rows. Missing data should be set as NA, which will then be ignored for the calculations. All samples must have the same ploidy.

### Usage

```
dosage_ratios(data, ploidy)
```

### Arguments

data	Genotype matrix or data.frame
ploidy	The ploidy of the species being analyzed

### Value

A data.frame with percentages of dosage values in the genotype matrix

### Examples

```

# example numeric genotype matrix for a tetraploid
n_ind <- 5
n_snps <- 10

geno <- matrix(as.numeric(sample(0:4, n_ind * n_snps, replace = TRUE)), nrow = n_snps, ncol = n_ind)
colnames(geno) <- paste0("Ind", 1:n_ind)
rownames(geno) <- paste0("SNP", 1:n_snps)
ploidy <- 4

# ratio of dosage value (numeric genotypes) across samples in dataset
result <- dosage_ratios(geno, ploidy)

print(result)

```

---

filterMADC	<i>Filter MADC Files</i>
------------	--------------------------

---

### Description

Filter and process MADC files to remove low quality microhaplotypes

### Usage

```
filterMADC(
  madc_file,
  min.mean.reads = NULL,
  max.mean.reads = NULL,
  max.mhaps.per.loci = NULL,
  min.reads.per.site = 1,
  min.ind.with.reads = NULL,
  target.only = FALSE,
  n.summary.columns = NULL,
  output.file = NULL
)
```

### Arguments

<code>madc_file</code>	Path to the MADC file to be filtered
<code>min.mean.reads</code>	Minimum mean read depth for filtering
<code>max.mean.reads</code>	Maximum mean read depth for filtering
<code>max.mhaps.per.loci</code>	Maximum number of matching mhaps per target loci. Retains only the target Ref and Alt loci at the sites that exceeds the <code>max.mhaps.per.loci</code> threshold.
<code>min.reads.per.site</code>	Minimum number of reads per site for <code>min.ind.with.reads</code> . Otherwise, this parameter is ignored
<code>min.ind.with.reads</code>	Minimum number of individuals with <code>min.reads.per.site</code> reads for filtering
<code>target.only</code>	Logical indicating whether to filter for target loci only
<code>n.summary.columns</code>	(optional) Number of summary columns to remove from MADC file not including the first three. Otherwise, the columns will be automatically detected and removed.
<code>output.file</code>	Path to save the filtered data (if NULL, data will not be saved)

### Details

This function can filter raw MADC files or pre-processed MADC files with fixed allele IDs. Additionally, it can filter based on mean read depth, number of mhaps per target loci, and other criteria. Optionally, users can plot summary statistics and save the filtered data to a file.

**Value**

data.frame or saved csv file

**Examples**

```
#Example

#Example MADC
macd_file <- system.file("example_MADC_FixedAlleleID.csv", package="BIGr")

#Remove mhaps exceeding 3 per target region including the ref and alt target mhaps
filtered_df <- filterMADC(macd_file,
                          min.mean.reads = NULL,
                          max.mean.reads = NULL,
                          max.mhaps.per.loci = 3,
                          min.reads.per.site = 1,
                          min.ind.with.reads = NULL,
                          target.only = FALSE,
                          n.summary.columns = NULL,
                          output.file = NULL)
```

---

filterVCF

*Filter a VCF file*

---

**Description**

This function will filter a VCF file or vcfR object and export the updated version

**Usage**

```
filterVCF(
  vcf.file,
  filter.OD = NULL,
  filter.BIAS.min = NULL,
  filter.BIAS.max = NULL,
  filter.DP = NULL,
  filter.MPP = NULL,
  filter.PMC = NULL,
  filter.MAF = NULL,
  filter.SAMPLE.miss = NULL,
  filter.SNP.miss = NULL,
  ploidy,
  output.file = NULL
)
```

**Arguments**

<code>vcf.file</code>	vcfR object or path to VCF file. Can be unzipped (.vcf) or gzipped (.vcf.gz).
<code>filter.OD</code>	Updog filter
<code>filter.BIAS.min</code>	Updog filter (requires a value for both <code>BIAS.min</code> and <code>BIAS.max</code> )
<code>filter.BIAS.max</code>	Updog filter (requires a value for both <code>BIAS.min</code> and <code>BIAS.max</code> )
<code>filter.DP</code>	Total read depth at each SNP filter
<code>filter.MPP</code>	Updog filter
<code>filter.PMC</code>	Updog filter
<code>filter.MAF</code>	Minor allele frequency filter
<code>filter.SAMPLE.miss</code>	Sample missing data filter
<code>filter.SNP.miss</code>	SNP missing data filter
<code>ploidy</code>	The ploidy of the species being analyzed
<code>output.file</code>	output file name (optional). If no <code>output.file</code> name provided, then a vcfR object will be returned.

**Details**

This function will input a VCF file or vcfR object and filter based on the user defined options. The output file will be saved to the location and with the name that is specified. The VCF format is v4.3

**Value**

A gzipped vcf file

**Examples**

```
## Use file paths for each file on the local system

#Temp location (only for example)
output_file <- tempfile()

filterVCF(vcf.file = system.file("iris_DArT_VCF.vcf.gz", package = "BIGr"),
          filter.OD = 0.5,
          filter.MAF = 0.05,
          ploidy = 2,
          output.file = output_file)

# Removing the output for the example
rm(output_file)

##The function will output the filtered VCF to the current working directory
```

---

`flip_dosage`*Switch Dosage Values from a Genotype Matrix*

---

### Description

This function converts the dosage count values to the opposite value. This is primarily used when converting dosage values from reference based (0 = homozygous reference) to alternate count based (0 = homozygous alternate). It assumes that the Samples are the columns, and the genomic markers are in rows. Missing data should be set as NA, which will then be ignored for the calculations. All samples must have the same ploidy.

### Usage

```
flip_dosage(df, ploidy, is.reference = TRUE)
```

### Arguments

<code>df</code>	Genotype matrix or data.frame
<code>ploidy</code>	The ploidy of the species being analyzed
<code>is.reference</code>	The dosage calls value is based on the count of reference alleles (TRUE/FALSE)

### Value

A genotype matrix

### Examples

```
# example code

# example numeric genotype matrix for a tetraploid
n_ind <- 5
n_snps <- 10

geno <- matrix(as.numeric(sample(0:4, n_ind * n_snps, replace = TRUE)), nrow = n_snps, ncol = n_ind)
colnames(geno) <- paste0("Ind", 1:n_ind)
rownames(geno) <- paste0("SNP", 1:n_snps)
ploidy <- 4

# Output matrix with the allele count reversed
results <- flip_dosage(geno, ploidy, is.reference = TRUE)

print(results)
```

---

get\_countsMADC      *Obtain Read Counts from MADC File*

---

### Description

This function takes the MADC file as input and retrieves the ref and alt counts for each sample, and converts them to ref, alt, and size(total count) matrices for dosage calling tools. At the moment, only the read counts for the Ref and Alt target loci are obtained while the additional loci are ignored.

### Usage

```
get_countsMADC(madc_file)
```

### Arguments

madc\_file      Path to MADC file

### Value

A list of read count matrices for reference, alternate, and total read count values

### Examples

```
# Get the path to the MADC file
madc_path <- system.file("iris_DArT_MADC.csv", package = "BIGr")

# Extract the read count matrices
counts_matrices <- get_countsMADC(madc_path)

# Access the reference, alternate, and size matrices

# ref_matrix <- counts_matrices$ref_matrix
# alt_matrix <- counts_matrices$alt_matrix
# size_matrix <- counts_matrices$size_matrix

rm(counts_matrices)
```

---

imputation\_concordance      *Calculate Concordance between Imputed and Reference Genotypes*

---

### Description

This function calculates the concordance between imputed and reference genotypes. It assumes that samples are rows and markers are columns. It is recommended to use allele dosages (0, 1, 2) but will work with other formats. Missing data in reference or imputed genotypes will not be considered for concordance if the missing\_code argument is used. If a specific subset of markers should be excluded, it can be provided using the snps\_2\_exclude argument.

**Usage**

```
imputation_concordance(
  reference_genos,
  imputed_genos,
  missing_code = NULL,
  snps_2_exclude = NULL,
  verbose = FALSE
)
```

**Arguments**

reference_genos	A data frame containing reference genotype data, with rows as samples and columns as markers. Dosage format (0, 1, 2) is recommended.
imputed_genos	A data frame containing imputed genotype data, with rows as samples and columns as markers. Dosage format (0, 1, 2) is recommended.
missing_code	An optional value to specify missing data. If provided, loci with this value in either dataset will be excluded from the concordance calculation.
snps_2_exclude	An optional vector of marker IDs to exclude from the concordance calculation.
verbose	A logical value indicating whether to print a summary of the concordance results. Default is FALSE.

**Details**

The function identifies common samples and markers between the reference and imputed genotype datasets. It calculates the percentage of matching genotypes for each sample, excluding missing data and specified markers. The concordance is reported as a percentage for each sample, along with a summary of the overall concordance distribution.

**Value**

A list with two elements:

- `result_df`: A data frame with sample IDs and their concordance percentages.
- `summary_concordance`: A summary of concordance percentages, including minimum, maximum, mean, and quartiles.

**Examples**

```
# Example Input variables
ignore_file <- system.file("imputation_ignore.txt", package="BIGr")
ref_file <- system.file("imputation_reference.txt", package="BIGr")
test_file <- system.file("imputation_test.txt", package="BIGr")

# Import files
snps = read.table(ignore_file, header = TRUE)
ref = read.table(ref_file, header = TRUE)
test = read.table(test_file, header = TRUE)
```

```
#Calculations
result <- imputation_concordance(reference_genos = ref,
                                imputed_genos = test,
                                snps_2_exclude = snps,
                                missing_code = 5,
                                verbose = FALSE)
```

---

macd2gmat

*Convert MADC Files to an Additive Genomic Relationship Matrix*

---

### Description

Scale and normalize MADC read count data and convert it to an additive genomic relationship matrix.

### Usage

```
macd2gmat(
  macd_file,
  seed = NULL,
  method = "collapsed",
  ploidy,
  output.file = NULL
)
```

### Arguments

macd_file	Path to the MADC file to be filtered
seed	Optional seed for random number generation (default is NULL)
method	Method to use for processing the MADC data. Options are "unique" or "collapsed". Default is "collapsed".
ploidy	Numeric. Ploidy level of the samples (e.g., 2 for diploid, 4 for tetraploid)
output.file	Path to save the filtered data (if NULL, data will not be saved)

### Details

This function reads a MADC file, processes it to remove unnecessary columns, and then converts it into an additive genomic relationship matrix using the first method proposed by VanRaden (2008). The resulting matrix can be used for genomic selection or other genetic analyses.

### Value

data.frame or saved csv file

**Author(s)**

Alexander M. Sandercock

**References**

VanRaden, P. M. (2008). Efficient methods to compute genomic predictions. *Journal of Dairy Science*, 91(11), 4414-4423

**Examples**

```
#Input variables
mac_file <- system.file("example_MADC_FixedAlleleID.csv", package="BIGr")

#Calculations
temp <- tempfile()

# Converting to additive relationship matrix
gmat <- mac2gmat(mac_file,
                seed = 123,
                ploidy=2,
                output.file = NULL)
```

---

mac2vcf\_all

*Converts MADC file to VCF recovering target and off-target SNPs*

---

**Description**

This function processes a MADC file to generate a VCF file containing both target and off-target SNPs. It includes options for filtering multiallelic SNPs and parallel processing to improve performance.

**Usage**

```
mac2vcf_all(
  mac = NULL,
  botloci_file = NULL,
  hap_seq_file = NULL,
  n.cores = 1,
  rm_multiallelic_SNP = FALSE,
  multiallelic_SNP_dp_thr = 0,
  multiallelic_SNP_sample_thr = 0,
  alignment_score_thr = 40,
  out_vcf = NULL,
  verbose = TRUE
)
```

## Arguments

<code>madc</code>	A string specifying the path to the MADC file.
<code>botloci_file</code>	A string specifying the path to the file containing the target IDs designed in the bottom strand.
<code>hap_seq_file</code>	A string specifying the path to the haplotype database fasta file.
<code>n.cores</code>	An integer specifying the number of cores to use for parallel processing. Default is 1.
<code>rm_multiallelic_SNP</code>	A logical value. If TRUE, SNPs with more than one alternative base are removed. If FALSE, the thresholds specified by <code>multiallelic_SNP_dp_thr</code> and <code>multiallelic_SNP_sample_thr</code> are used to filter low-frequency SNP alleles. Default is FALSE.
<code>multiallelic_SNP_dp_thr</code>	A numeric value specifying the minimum depth by tag threshold for filtering low-frequency SNP alleles when <code>rm_multiallelic_SNP</code> is FALSE. Default is 0.
<code>multiallelic_SNP_sample_thr</code>	A numeric value specifying the minimum number of samples threshold for filtering low-frequency SNP alleles when <code>rm_multiallelic_SNP</code> is FALSE. Default is 0.
<code>alignment_score_thr</code>	A numeric value specifying the minimum alignment score threshold. Default is 40.
<code>out_vcf</code>	A string specifying the name of the output VCF file. If the file extension is not <code>.vcf</code> , it will be appended automatically.
<code>verbose</code>	A logical value indicating whether to print metrics and progress to the console. Default is TRUE.

## Details

The function processes a MADC file to generate a VCF file containing both target and off-target SNPs. It uses parallel processing to improve performance and provides options to filter multiallelic SNPs based on user-defined thresholds. The alignment score threshold can be adjusted using the `alignment_score_thr` parameter. The generated VCF file includes metadata about the processing parameters and the BIGr package version. If the `alignment_score_thr` is not met, the corresponding SNPs are discarded.

## Value

This function does not return an R object. It writes the processed VCF file v4.3 to the specified `out_vcf` path.

## Examples

```
# Example usage:
```

```

Sys.setenv("OMP_THREAD_LIMIT" = 2)

madc_file <- system.file("example_MADC_FixedAlleleID.csv", package="BIGr")
bot_file <- system.file("example_SNPs_DArTag-probe-design_f180bp.botloci", package="BIGr")
db_file <- system.file("example_allele_db.fa", package="BIGr")

#Temp location (only for example)
output_file <- tempfile()

madc2vcf_all(
  macd = macd_file,
  botloci_file = bot_file,
  hap_seq_file = db_file,
  n.cores = 2,
  rm_multiallelic_SNP = TRUE,
  multiallelic_SNP_dp_thr = 10,
  multiallelic_SNP_sample_thr = 5,
  alignment_score_thr = 40,
  out_vcf = output_file,
  verbose = TRUE
)

rm(output_file)

```

---

madc2vcf\_targets

*Format MADC Target Loci Read Counts Into VCF*


---

## Description

This function will extract the read count information from a MADC file target markers and convert to VCF file format.

## Usage

```
madc2vcf_targets(madc_file, output.file, botloci_file, get_REF_ALT = FALSE)
```

## Arguments

macd_file	Path to MADC file
output.file	output file name and path
botloci_file	A string specifying the path to the file containing the target IDs designed in the bottom strand.
get_REF_ALT	if TRUE recovers the reference and alternative bases by comparing the sequences. If more than one polymorphism are found for a tag, it is discarded.

**Details**

The DArTag MADC file format is not commonly supported through existing tools. This function will extract the read count information from a MADC file for the target markers and convert it to a VCF file format for the genotyping panel target markers only

**Value**

A VCF file v4.3 with the target marker read count information

A VCF file v4.3 with the target marker read count information

**Examples**

```
# Load example files
macd_file <- system.file("example_MADC_FixedAlleleID.csv", package="BIGr")
bot_file <- system.file("example_SNPs_DArTag-probe-design_f180bp.botloci", package="BIGr")

#Temp location (only for example)
output_file <- tempfile()

# Convert MADC to VCF
macd2vcf_targets(macd_file = macd_file,
                 output.file = output_file,
                 get_REF_ALT = TRUE,
                 botloci_file = bot_file)

rm(output_file)
```

---

merge\_MADCs

*Merge MADC files*


---

**Description**

If duplicated samples exist in different files, a suffix will be added at the end of the sample name. If `run_ids` is defined, they are used as suffix, if not, files will be identified from 1 to number of files, considering the order that was defined in the function.

**Usage**

```
merge_MADCs(..., macd_list = NULL, out_macd = NULL, run_ids = NULL)
```

**Arguments**

<code>...</code>	one or more MADC files path
<code>macd_list</code>	list containing path to MADC files to be merged
<code>out_macd</code>	output merged MADC file path
<code>run_ids</code>	vector of character defining the run ID for each file. This ID will be added as a suffix in repeated sample ID in case they exist in different files.

**Value**

A data frame containing the merged MADC data. The merged file is also written to the specified `out_madc` path in CSV format. Numeric columns are filled with zeros where data is missing.

**Examples**

```
# First generating example MADC files
temp_dir <- tempdir()
file1_path <- file.path(temp_dir, "macd1.csv")
file2_path <- file.path(temp_dir, "macd2.csv")
out_path <- file.path(temp_dir, "merged_madc.csv")

# Data for file 1: Has SampleA and SampleB
df1 <- data.frame(
  AlleleID = c("chr1.1_0001|Alt_0002", "chr1.1_0001|Ref_0001", "chr1.1_0001|AltMatch_0001"),
  CloneID = c("chr1.1_0001", "chr1.1_0001", "chr1.1_0001"),
  AlleleSequence = c("GGG", "AAA", "TTT"),
  SampleA = c(10, 8, 0),
  SampleB = c(5, 4, 9),
  stringsAsFactors = FALSE,
  check.names = FALSE
)
write.csv(df1, file1_path, row.names = FALSE, quote = FALSE)

# Data for file 2: Has SampleA (duplicate name) and SampleC, different rows
df2 <- data.frame(
  AlleleID = c("chr1.1_0001|Alt_0002", "chr1.1_0001|Ref_0001", "chr1.1_0001|AltMatch_0001"),
  CloneID = c("chr1.1_0001", "chr1.1_0001", "chr1.1_0001"),
  AlleleSequence = c("GGG", "AAA", "TTT"),
  SampleA = c(11, 7, 20),
  SampleC = c(1, 2, 6),
  stringsAsFactors = FALSE,
  check.names = FALSE
)
write.csv(df2, file2_path, row.names = FALSE, quote = FALSE)

# 2. Run the merge function
# Use default suffixes (.x, .y) for the duplicated "SampleA"
merge_MADCs(madc_list = list(file1_path, file2_path),
            out_madc = out_path)
```

---

solve\_composition\_poly

*Compute Genome-Wide Breed Composition*

---

**Description**

Computes genome-wide breed/ancestry composition using quadratic programming on a batch of animals.

**Usage**

```

solve_composition_poly(
  Y,
  X,
  ped = NULL,
  groups = NULL,
  mia = FALSE,
  sire = FALSE,
  dam = FALSE,
  ploidy = 2
)

```

**Arguments**

Y	numeric matrix of genotypes (columns) from all animals (rows) in population coded as dosage of allele B {0, 1, 2, ..., ploidy}
X	numeric matrix of allele frequencies (rows) from each reference panel (columns). Frequencies are relative to allele B.
ped	data.frame giving pedigree information. Must be formatted "ID", "Sire", "Dam"
groups	list of IDs categorized by breed/population. If specified, output will be a list of results categorized by breed/population.
mia	logical. Only applies if ped argument is supplied. If true, returns a data.frame containing the inferred maternally inherited allele for each locus for each animal instead of breed composition results.
sire	logical. Only applies if ped argument is supplied. If true, returns a data.frame containing sire genotypes for each locus for each animal instead of breed composition results.
dam	logical. Only applies if ped argument is supplied. If true, returns a data.frame containing dam genotypes for each locus for each animal instead of breed composition results.
ploidy	integer. The ploidy level of the species (e.g., 2 for diploid, 3 for triploid, etc.).

**Value**

A data.frame or list of data.frames (if groups is !NULL) with breed/ancestry composition results

**References**

Funkhouser SA, Bates RO, Ernst CW, Newcom D, Steibel JP. Estimation of genome-wide and locus-specific breed composition in pigs. *Transl Anim Sci.* 2017 Feb 1;1(1):36-44.

**Examples**

```

# Example inputs for solve_composition_poly (ploidy = 4)

# (This would typically be the output from allele_freq_poly)
allele_freqs_matrix <- matrix(

```

```

c(0.625, 0.500,
  0.500, 0.500,
  0.500, 0.500,
  0.750, 0.500,
  0.625, 0.625),
nrow = 5, ncol = 2, byrow = TRUE,
dimnames = list(paste0("SNP", 1:5), c("VarA", "VarB"))
)

# Validation Genotypes (individuals x SNPs)
val_genotype_matrix <- matrix(
  c(2, 1, 2, 3, 4, # Test1 dosages for SNP1-5
    3, 4, 2, 3, 0), # Test2 dosages for SNP1-5
  nrow = 2, ncol = 5, byrow = TRUE,
  dimnames = list(paste0("Test", 1:2), paste0("SNP", 1:5))
)

# Calculate Breed Composition
composition <- solve_composition_poly(Y = val_genotype_matrix,
                                       X = allele_freqs_matrix,
                                       ploidy = 4)

print(composition)

```

---

thinSNP

*Thin a dataframe of SNPs based on genomic position*


---

## Description

This function groups SNPs by chromosome, sorts them by physical position, and then iteratively selects SNPs such that no two selected SNPs within the same chromosome are closer than a specified minimum distance.

## Usage

```
thinSNP(df, chrom_col_name, pos_col_name, min_distance)
```

## Arguments

df	The input dataframe.
chrom_col_name	A string specifying the name of the chromosome column.
pos_col_name	A string specifying the name of the physical position column.
min_distance	A numeric value for the minimum distance between selected SNPs. The unit of this distance should match the unit of the pos_col_name column (e.g., base pairs).

## Value

A thinned dataframe with the same columns as the input.

**Examples**

```

# Create sample SNP data
set.seed(123)
n_snps <- 20
snp_data <- data.frame(
  MarkerID = paste0("SNP", 1:n_snps),
  Chrom = sample(c("chr1", "chr2"), n_snps, replace = TRUE),
  ChromPosPhysical = c(
    sort(sample(1:1000, 5)), # SNPs on chr1
    sort(sample(1:1000, 5)) + 500, # More SNPs on chr1
    sort(sample(1:2000, 10))      # SNPs on chr2
  ),
  Allele = sample(c("A/T", "G/C"), n_snps, replace = TRUE)
)
# Ensure it's sorted by Chrom and ChromPosPhysical for clarity in example
snp_data <- snp_data[order(snp_data$Chrom, snp_data$ChromPosPhysical), ]
rownames(snp_data) <- NULL

print("Original SNP data:")
print(snp_data)

# Thin the SNPs, keeping a minimum distance of 100 units (e.g., bp)
thinned_snps <- thinSNP(
  df = snp_data,
  chrom_col_name = "Chrom",
  pos_col_name = "ChromPosPhysical",
  min_distance = 100
)

print("Thinned SNP data (min_distance = 100):")
print(thinned_snps)

# Thin with a larger distance
thinned_snps_large_dist <- thinSNP(
  df = snp_data,
  chrom_col_name = "Chrom",
  pos_col_name = "ChromPosPhysical",
  min_distance = 500
)

print("Thinned SNP data (min_distance = 500):")
print(thinned_snps_large_dist)

```

---

updog2vcf

*Export Updog Results as VCF*


---

**Description**

This function will convert an Updog output to a VCF file

**Usage**

```
updog2vcf(
  multidog.object,
  output.file,
  updog_version = NULL,
  RefAlt = NULL,
  compress = TRUE
)
```

**Arguments**

<code>multidog.object</code>	updog output object with class "multidog" from dosage calling
<code>output.file</code>	output file name and path
<code>updog_version</code>	character defining updog package version used to generate the multidog object
<code>RefAlt</code>	optional data frame with four columns named "Chr", "Pos", "Ref", and "Alt" containing the reference and alternate alleles for each SNP in the same order as in the multidog object
<code>compress</code>	logical. If TRUE returns a vcf.gz file

**Details**

When performing dosage calling for multiple SNPs using Updog, the output file contains information for all loci and all samples. This function will convert the updog output file to a VCF file, while retaining the information for the values that are commonly used to filter low quality and low confident dosage calls.

**Value**

A vcf file

**References**

Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping polyploids from messy sequencing data. *Genetics*, 210(3), 789-807.

**Examples**

```
# Retrieving the updog output multidog object
load(system.file("extdata", "iris-multidog.rdata", package = "BIGr"))

temp_file <- tempfile()

# Convert updog to VCF, where the new VCF will be saved at the location specified in the output.file
updog2vcf(
  multidog.object = mout,
  output.file = temp_file,
  updog_version = "0.0.0",
  compress = TRUE
)
```

)

```
#Removing the example vcf  
rm(temp_file)
```

# Index

allele\_freq\_poly, [2](#)

calculate\_Het, [3](#)

calculate\_MAF, [4](#)

check\_homozygous\_trios, [5](#)

check\_ped, [7](#)

check\_replicates, [8](#)

dosage2vcf, [9](#)

dosage\_ratios, [10](#)

filterMADC, [11](#)

filterVCF, [12](#)

flip\_dosage, [14](#)

get\_countsMADC, [15](#)

imputation\_concordance, [15](#)

madc2gmat, [17](#)

madc2vcf\_all, [18](#)

madc2vcf\_targets, [20](#)

merge\_MADCs, [21](#)

solve\_composition\_poly, [22](#)

thinSNP, [24](#)

updog2vcf, [25](#)