

R-evolutionary insights

Deep dive into
SNP-based
population
genomics

ANU Guest Conference WiFi
Username: PopGenR
Password: Kioloa_PopGen_24



SESSION 1

Introduction

- Olly Berry – Welcome and Global Perspective [10 mins]
- Andrzej Kilian – Introduction to Diversity Arrays Technology [15 mins]
- Bernd Gruber – The dartRverse [15 mins]
- Arthur Georges – dartR fundamentals -- HANDS ON [55 mins]
- Renee Catullo – More on filtering and the need for a flexible experimental approach to analysis [15 mins]
- Arthur Georges -- Discussion [10 mins]

Welcome!

Workshop on population
genomics with the
dartR platform

Kioloa, March 11-15, 2024

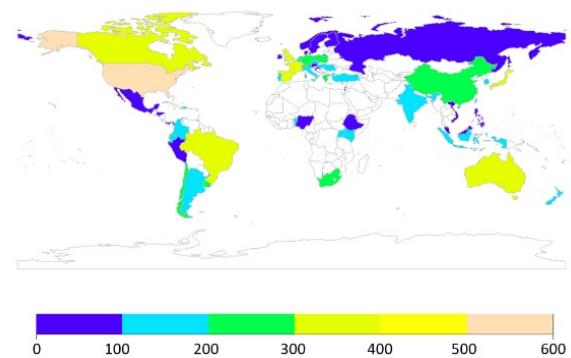


Oliver Berry
CSIRO Environomics Future Science Platform

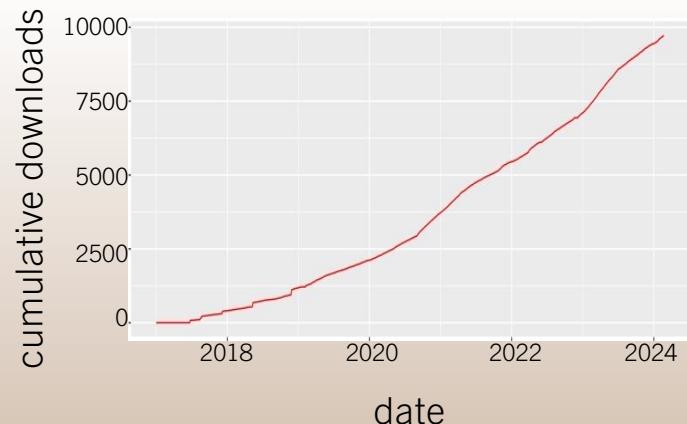




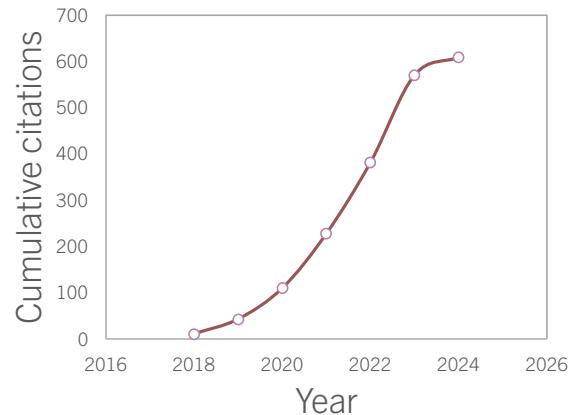
dartR downloads (geographic)



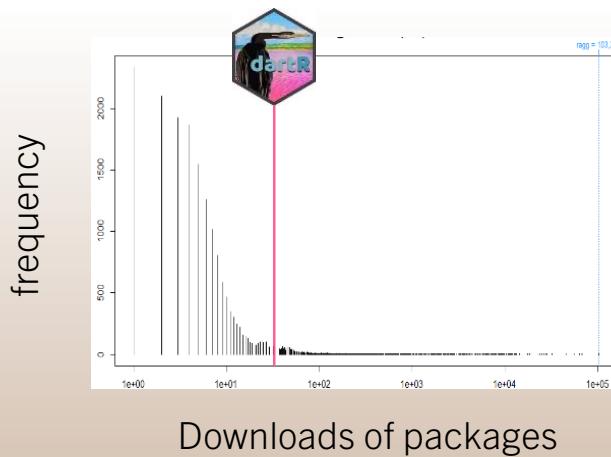
dartR downloads numbers



dartR Citations



dartR in context





Popgen in R Workshop | 2019 | Hobart





Main developers



Luis



Diana



Arthur



Bernd



Carlo



Workshop contributors



Welcome!

Workshop on population
genomics with the
dartR platform

Kioloa, March 11-15, 2024



Introduction to DArT

Andrzej Kilian
Kioloa PopGen
11 March 2024



for sustainable future



 **DArT** diversity arrays technology

Since 2001

Our Mission

Actively support more
effective and
sustainable use of
natural resources in
equitable manner.

A social enterprise
operating as a
successful business



About Us

Genome profiling services and IT support
for agriculture and ecology.

→ Cost-effective solutions

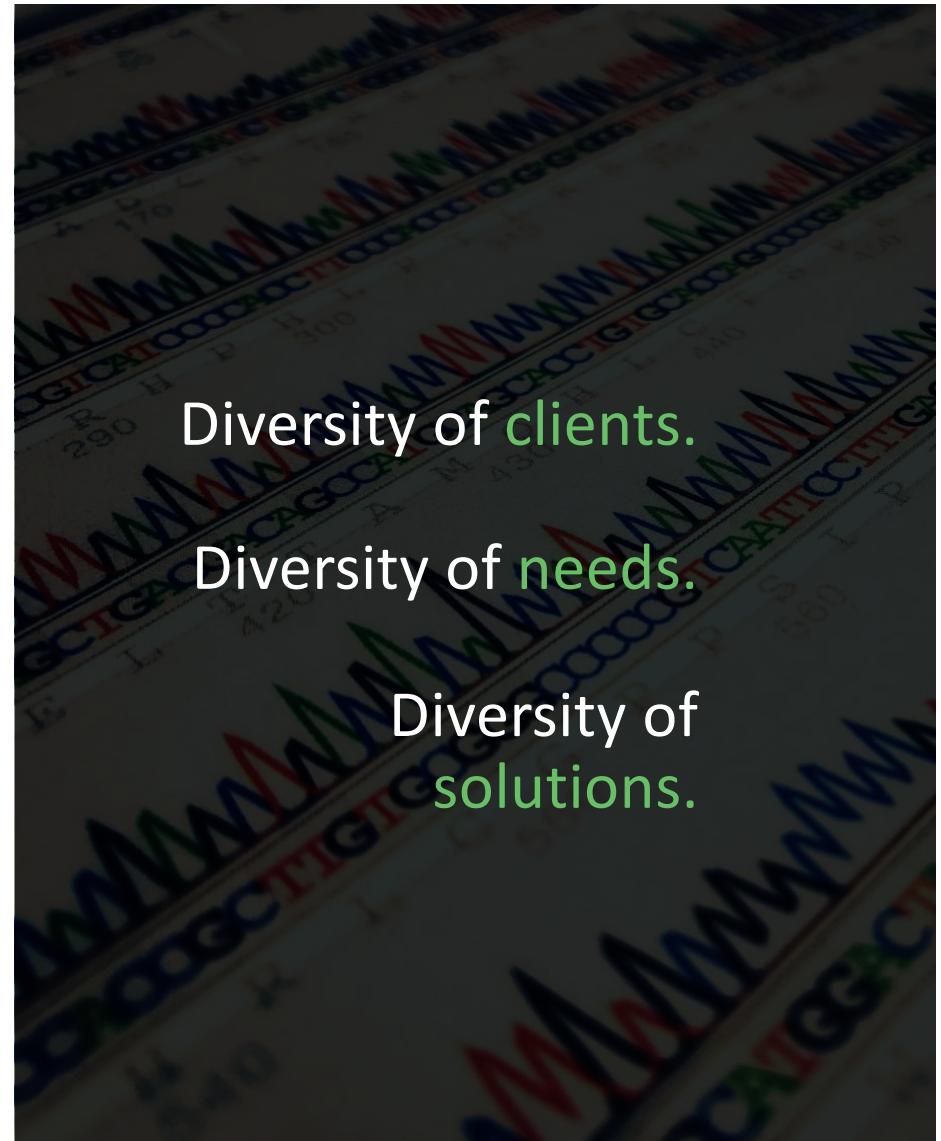
We are the **first affordable** whole genome profiling service based on DArT, providing end-to-end IT and service management throughout the **entire breeding lifecycle**.

→ Scalable & independent

We offer **high throughput solution** for any organism due to a high level of automation, sequence information & platform independence.

→ Integrated and Big Data ready

OneDArT **integrates all data types** relevant for crop improvement, ecology and agricultural practice



Diversity of clients.

Diversity of needs.

Diversity of
solutions.



Our Team

A diverse, international team with a mixture of molecular biology, chemistry, electronics, mathematics, stats and IT skills.

65
staff & growing



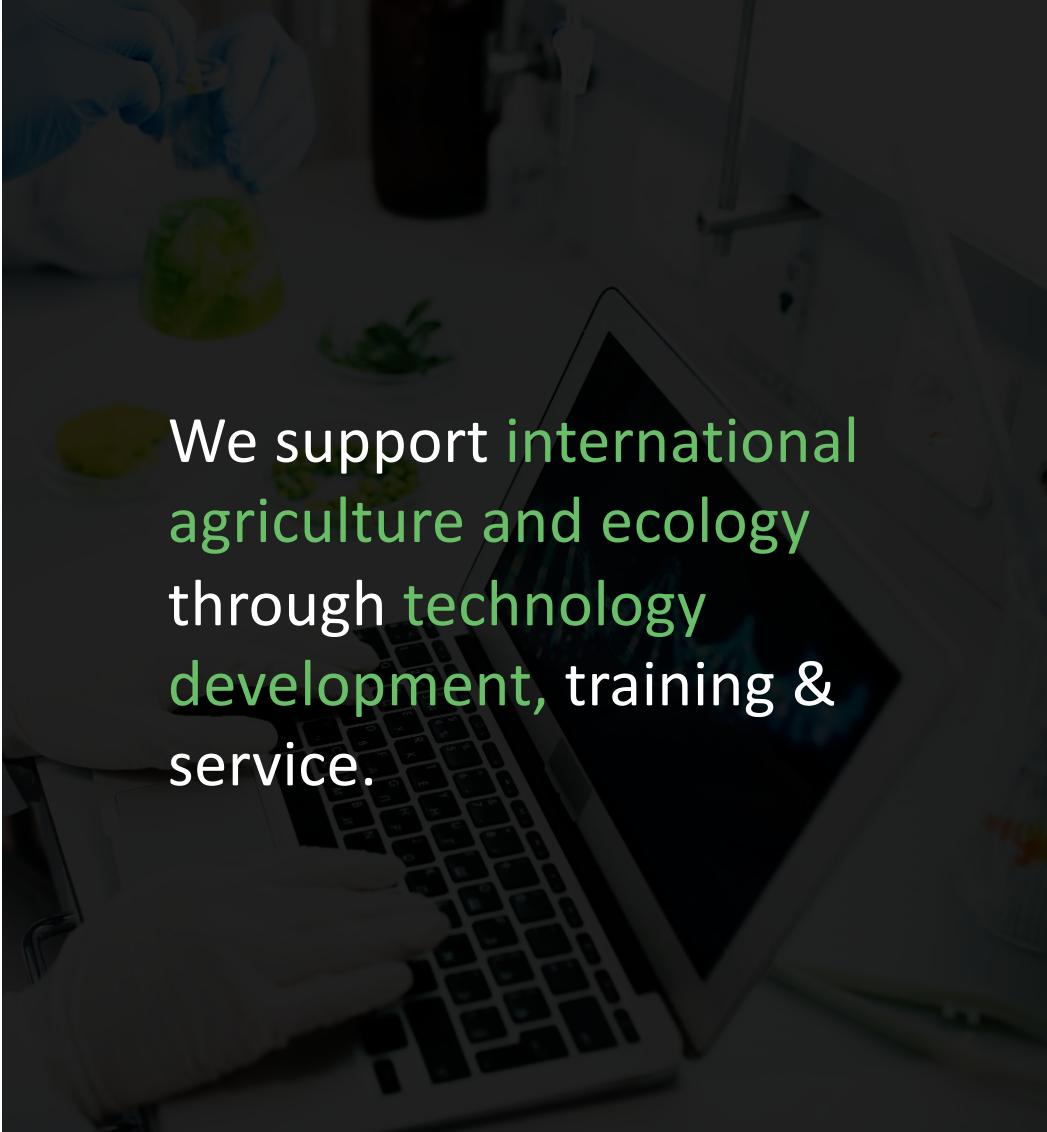
Corporate Structure

Organised into three main domains of operations with significant overlap of staff across domains



Our culture

Operating on “chaordic” principles and family-like atmosphere



We support **international agriculture** and ecology through **technology development**, training & service.

3.7 million + genome profiles

1000+ organisms
(thousands of species)

2500+ clients

Global Reach

SAGA, Mexico

Our international footprint spans every continent and over 80 countries.

South America

North America

United Kingdom

European Union

Middle East

C&E Africa
IGSS/SeqArt

Russia

Asia

South East Asia

ACT, Australia
KDDart HQ



Genomic analysis Landscape

Started with DNA array-based methods but moved to using Next Generation Sequencing (NGS)

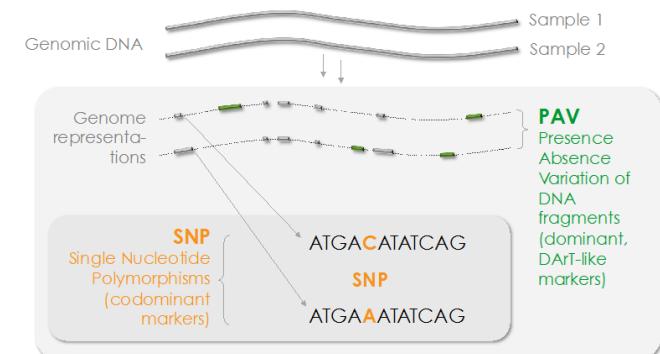
DArTseq - leading genotyping by sequencing technology

- Sequencing random genome fragments creating “genome representations”
- Highly scalable – adjust marker number
- High data quality
- Ability to resolve closely related material

Targeted Genotyping - amplicon sequencing

- **DArTag** – 300- 10,000 selected SNPs
 - Predominantly breeding tool, increasingly adopted in ecology
- **DArTcap** – 100- 10,000 selected SNPs used heavily in agriculture and in ecology
- **DArTmp** – up to 300 amplicons sequenced, used for genetic identification and paternity testing

DArTseqMet - DNA methylation analysis both at specific loci and genome-wide

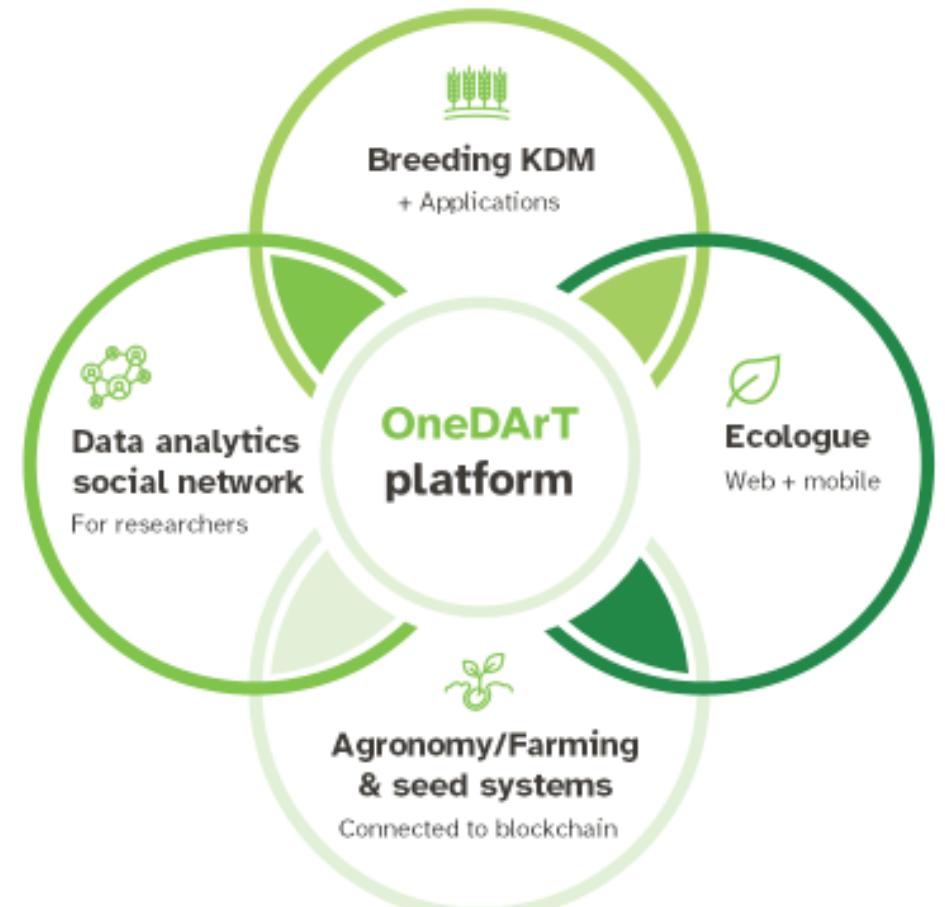


OneDArT platform hub

Meet the client-facing,
super-applications of the
OneDArT platform.

- ✓ Genotyping & sequencing services
- ✓ Ecommerce platform for **all services**
- ✓ **Shared** data sources, schema elements & APIs (DNA, environment, performance)

*Some domain specific



Architecture

Data

UI

Analytics

Configuration

Ecologue Overview

Target Customer Fields

- ✓ Ecology and popgen
- ✓ Biosecurity
- ✓ Biodiversity conservation
- ✓ Ecological Compliance
- ✓ Environmental Impact Assessment
- ✓ Captive breeding/genetic rescue

Correlate ecological data with genomic data

Ecologue Mobile App

- Collect Ecological Data, including Geographic data, time and place of sample collection
- Run App offline in remote locations, upload data when online later

Ecologue Web App

- Manage, curate & export data
- View data collection location information on maps
- Data analysis using integrated **dartR package**



Stakeholder Engagement

Started as a project for Chevron – platform for biodiversity monitoring on Barrow Island in Western Australia

Last three years supported by government with \$3M Priority Investment Program grant

Envisioned as a component of OneDAR^T platform: from sample and its metadata collection to analytics!

Initial testing and engagement in our target customer fields in preparation for release in second half 2024



External persona
Steven
Ecologist



Steven is an ecologist, at an ecological & biodiversity consultancy company in Canberra. He manages a team of ecologists in conducting field surveys, assessments and impact analysis of the biodiversity and ecological impact of construction & infrastructure projects in extractive industries to help clients comply with Australian legislation and mitigate negative impact to the environment.

Goals

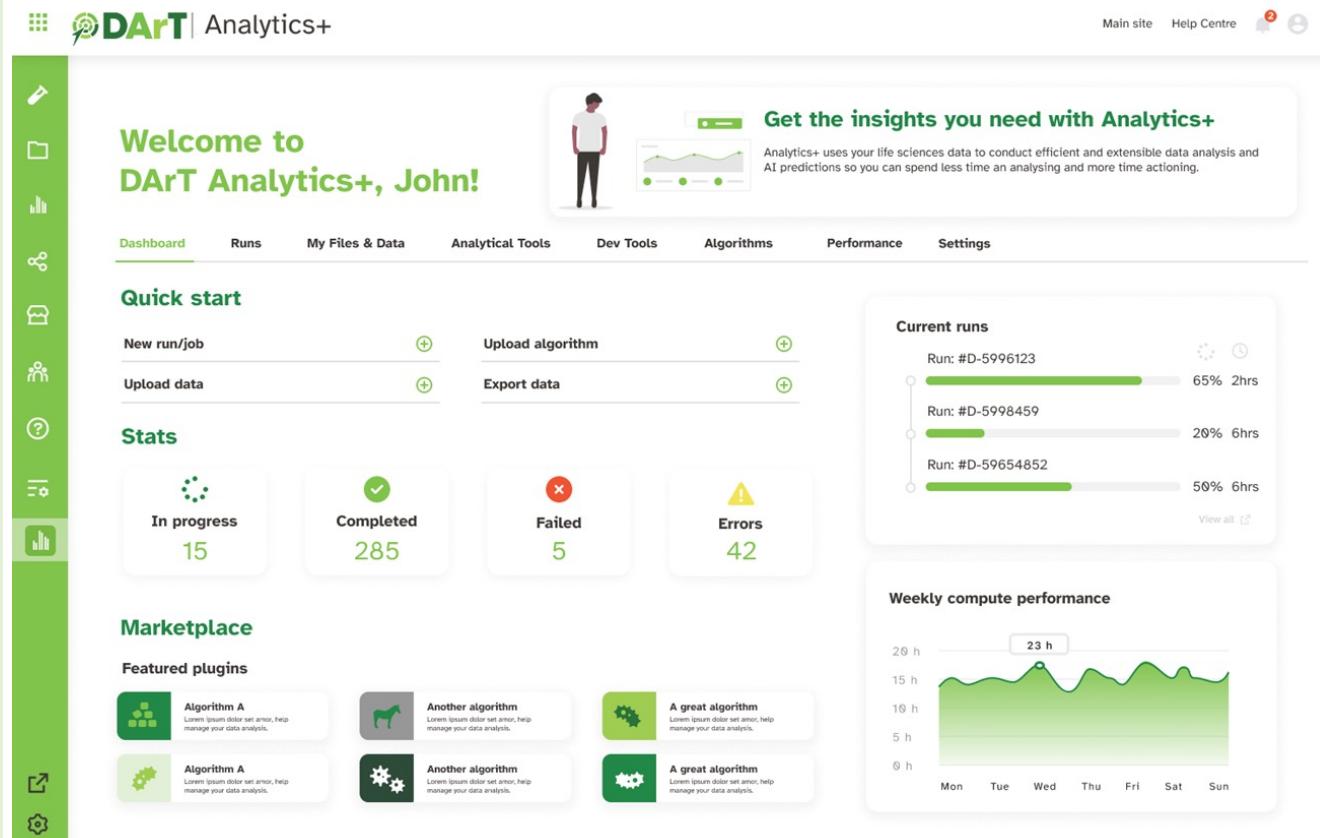
- To have a streamlined, central, cloud based data management tool where all field survey data is collated and can be exported and analysed
- To provide accurate assessments of development sites to provide advice to clients through a unified analysis framework
- Reduce data loss and increase data re-usability and reproducibility
- Being able to link genotypic information and DNA barcodes to ecological survey data

Pain Points

- Not being able to integrate multiple data sources including geographic data e.g. GEOME
- Using multiple, disparate platforms
- Challenges using and managing 'big data' and having the computational power to manage large data sets from different sources e.g. optical sensors, satellite sensors, apps, etc.
- Being able to extract usable information from complex and diverse data sources
- Connecting patterns in the data to ecological processes
- Reliability of traditional genomic tools e.g. phylogenetic trees
- The re-usability of different datasets and survey processes

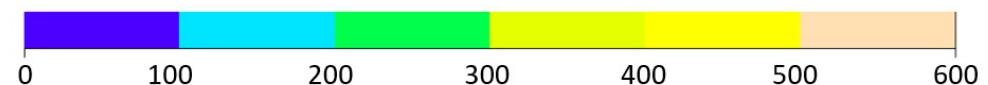
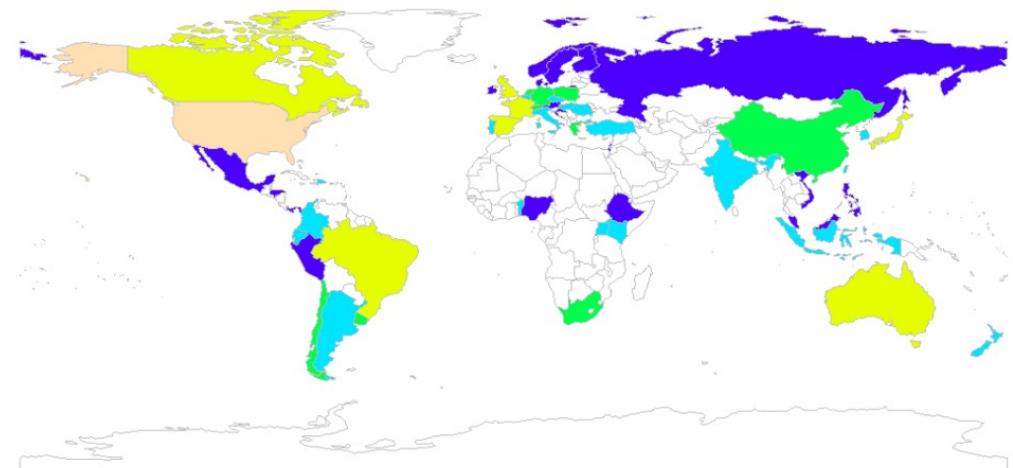
- **Analytics+ features**

- Users can take data directly from files/database storage
- **dartR** – one of the key applications, especially for ecologists
- Replication of all dartRverse functionalities – complementing and expanding dartR reach
- Reporting Hub – access to reports and format conversion
- DS14 application for marker data extraction from DArTseq libraries + reference genome-based tools
- Universal Test Bench (UTB): ML benchmarking platform for GWAS and genomic prediction tools



The dartRverse

Downloads [dartR](#) [July 2023]



The dartRverse



From 1 to 7 (8)

Why?

- Easier maintenance.
- Faster development.
- Less dependencies on other packages.
- Easier to become a contributor.
- Have your own package developed that can be branded as part of the **dartRverse**.
- [How to install is in the manual]



Overview

- + Input/output
- + filter and report
- + Basic analysis
- + Simple visualisation



Overview

- + Sample data for examples



Overview

- + Simple time forward simulations
- + Mutations, offspring
- + Simple dispersal scenarios
- + Complex processes (overdominance)
- + integrates with slimR



Overview

- + Comprehensive id and filter of sexlinked markers



Overview

- + Functions to study kinship
- + Support captive breeding



Overview

- + Studying diversity and population structure
- + Run Structure, Faststructure, Newhybrids, blast
- + Create an SFS
- + Study LD
- + Estimate Ne (Neestimator2)



Overview

- + IBD
- + spatial autocorrelation
- + Resistance surfaces
- + EEMS (estimating effective migration surfaces)



- + To come:
- + Estimate historic population sizes (Epos, Stairways2, Gone)
- + More “genomic” features



How become a contributor?

- + Simply send a script to us
[via email]

How become a contributor?

- + Simply send a script to us
- + Develop a fully-fledged gl function
 - some requirements
 - good help page
 - tutorial
 - <http://georges.biomatix.org/dartR>
 - coding evening
- + Get cited and become a package contributor

```
1 gl.sfs <- function(x,
2                           minbinsize = 0,
3                           folded = TRUE,
4                           singlepop = FALSE,
5                           plot.out = TRUE,
6                           plot.file = NULL,
7                           plot.dir = NULL,
8                           verbose = NULL) {
9
10  # SET VERBOSITY
11  verbose <- gl.check.verbosity(verbose)
12
13  # SET WORKING DIRECTORY
14  plot.dir <- gl.check.wd(plot.dir, verbose = 0)
15
16  # CHECK DATATYPE
17  datatype <- utils.check.datatype(x, verbose = verbose)
18
19  # FUNCTION SPECIFIC ERROR CHECKING
20
21  # DO THE JOB
22
23  # optionally save the plot -----
24
25  # FLAG SCRIPT END
26
27  return(sfs)
28
29  # + a good help page
```

How become a contributor?

- + Simply send a script to us
- + Develop a fully-fledged gl function
 - some requirements,
 - use github

THIS IS GIT. IT TRACKS COLLABORATIVE WORK ON PROJECTS THROUGH A BEAUTIFUL DISTRIBUTED GRAPH THEORY MODEL.

COOL. HOW DO WE USE IT?

NO IDEA. JUST MEMORIZIZE THESE SHELL COMMANDS AND TYPE THEM TO SYNC UP. IF YOU GET ERRORS, SAVE YOUR WORK ELSEWHERE, DELETE THE PROJECT, AND DOWNLOAD A FRESH COPY.



How become a contributor?

- + Simply send a script to us
- + Develop a fully-fledged gl function
- + Create your own package in the dartRverse and become a permanent member of the developer team

```
> citation("dartR.sexlinked")
```

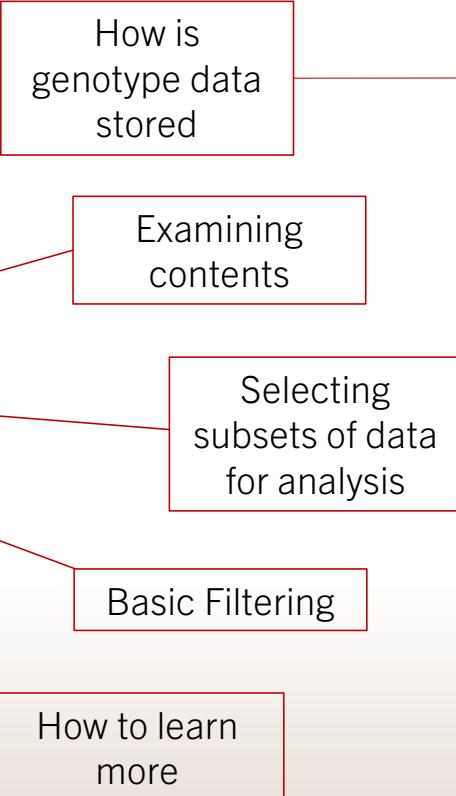
> To cite dartR in publications use:

Robledo-Ruiz, D. A., Austin, L., Amos, J. N., Castrejón-Figueroa, J., Harley, D. K. P., Magrath, M. J. L., Sunnucks, P., & Pavlova, A. (2023). Easy-to-use R functions to separate reduced-representation genomic datasets into sex-linked and autosomal loci, and conduct sex assignment. *Molecular Ecology Resources*, 00, 1–21. <https://doi.org/10.1111/1755-0998.13844>



dartR Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials



dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

INDIVIDUALS	LOCI																														
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	
AA010915	2	0	0	2	1	2	2	2	0	0	2	0	0	2	1	1	2	2	2	2	0	0	2	2	0	0	1	0	0	0	
UC_00126	2	-	2	0	0	2	1	2	2	2	0	0	2	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	0	
AA032760	0	0	-	0	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	
AA013214	0	2	0	0	0	2	2	0	0	0	1	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0	0	
AA011723	0	2	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	
AA012411	2	0	2	2	0	2	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0
AA019237	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	0
AA019238	0	0	0	2	2	2	0	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2
AA019239	0	2	0	0	0	-	0	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0
AA019235	0	2	0	0	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	
AA019240	1	0	-	0	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	0	
AA019241	2	0	2	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	
AA019242	0	0	0	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	
AA019243	0	1	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	
AA019251	0	0	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	
AA019252	2	0	0	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	
AA012405	2	-	0	1	0	0	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0
AA012406	0	0	0	2	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	
AA012409	0	0	2	0	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	
AA012499	0	2	2	2	2	0	0	0	2	2	0	0	1	-	-	2	0	0	0	2	2	0	0	0	1	-	2	-	2	0	
AA012422	1	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	0	
AA012434	2	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	
AA012469	0	0	0	2	2	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	
AA012500	2	0	1	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	0	1	
AA032799	2	0	0	2	2	2	1	2	0	0	2	0	0	0	2	0	2	2	2	1	1	0	0	0	2	2	0	0	0	1	

0 Homozygous reference allele
1 Heterozygous
2 Homozygous alternate allele
- Missing

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

Dataframe		
INDIVIDUAL METADATA		
	Latitude	Longitude
	Maturity	Sex

Dataframe

LOCUS METADATA

AlleleID, CloneID, AlleleSequence, SNP, SnpPosition, CallRate, OneRatioRef, OneRatioSnp, FreqHomRef, FreqHomSnp, FreqHets, PICRef, PICsnp, AvgPIC, AvgCountRef, AvgCountSnp, RepAvg

INDIVIDUALS	LOCI																												
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29
AA010915	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	2	2	0	1	0	0
UC_00126	2	-	2	0	0	2	1	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0
AA032760	0	0	-	0	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0
AA013214	0	2	0	0	0	2	2	0	0	0	1	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0
AA011723	0	2	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2
AA012411	2	0	2	2	0	2	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0
AA019237	2	0	2	0	0	2	1	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0
AA019238	0	0	0	2	2	2	0	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	2	1	1	2	2	2
AA019239	0	2	0	0	0	-	0	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0
AA019235	0	2	0	0	0	0	2	0	0	2	1	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	
AA019240	1	0	-	0	0	2	2	0	0	2	1	2	2	2	0	0	2	1	1	2	2	2	0	0	0	2	0	0	
AA019241	2	0	2	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0
AA019242	0	0	0	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	
AA019243	0	1	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0	0	2	2	0	
AA019251	0	0	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	
AA019252	2	0	0	0	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0	0	2	2	
AA012405	2	-	0	1	0	0	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	2	1	1	2	2	2	0	
AA012406	0	0	0	2	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	
AA012409	0	0	2	0	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	
AA012499	0	2	2	2	2	0	0	2	2	0	0	0	1	-	-	2	0	0	0	2	2	0	0	0	1	-	2	-	
AA012422	1	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	0	0	
AA012434	2	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	
AA012469	0	0	0	2	2	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	2	1	1	2	2	2	0	
AA012500	2	0	1	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	
AA032799	2	0	0	2	2	2	1	2	0	0	2	0	0	2	0	2	2	2	1	1	0	0	0	2	2	0	0	0	

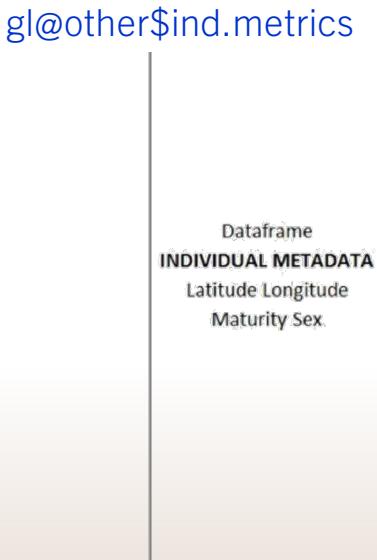
0 Homozygous reference allele
1 Heterozygous
2 Homozygous alternate allele
- Missing

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

Flags
History
Verbosity Setting
etc



dartR Object

INDIVIDUALS	LOCI																													
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
AA010915	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	2	2	0	1	0	0	0
UC_00126	2	-	2	0	0	2	1	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	0	2	2	0	0
AA032760	0	0	-	0	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2
AA013214	0	2	0	0	0	2	2	0	0	0	1	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0	0
AA011723	0	2	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0
AA012411	2	0	2	2	0	2	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0	0
AA019237	2	0	2	0	0	2	1	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	0
AA019238	0	0	0	2	2	2	0	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2
AA019239	0	2	0	0	0	-	0	-	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0
AA019235	0	2	0	0	0	0	2	0	0	2	1	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	0
AA019240	1	0	-	0	0	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0	0	0
AA019241	2	0	2	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0
AA019242	0	0	0	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2
AA019243	0	1	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0	0	2	2	0	0	0
AA019251	0	0	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2
AA019252	2	0	0	0	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0	0	2	2	0	0
AA012405	2	-	0	1	0	0	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	2	1	1	2	2	2	0	0	0
AA012406	0	0	0	2	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2
AA012409	0	0	2	0	2	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2
AA012499	0	2	2	2	0	0	0	2	2	0	0	1	-	-	2	0	0	0	2	2	0	0	0	1	-	2	-	2	-	2
AA012422	1	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	0	0	0	0
AA012434	2	2	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0
AA012469	0	0	0	2	2	2	0	0	2	0	0	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	0	0
AA012500	2	0	1	2	1	2	2	2	0	0	2	0	0	0	2	1	1	2	2	2	2	0	0	0	2	2	0	0	0	1
AA032799	2	0	0	2	2	2	1	2	0	0	2	0	0	2	0	2	1	1	0	0	0	2	2	0	0	0	1	0	0	0

0 Homozygous reference allele
1 Heterozygous
2 Homozygous alternate allele
- Missing

NA

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials



Follow along if you like.
Maybe use your own
genlight object?

Quick and Nasty



testset.gl

/// GENLIGHT OBJECT //////////

// 250 genotypes, 255 binary SNPs, size: 752 Kb
7868 (12.34 %) missing data

// Basic content

@gen: list of 250 SNPbin
@ploidy: ploidy of each individual (range: 2-2)

// Optional content

@ind.names: 250 individual labels
@loc.names: 255 locus labels
@loc.all: 255 alleles
@position: integer storing positions of the SNPs
@pop: population of each individual (group size range: 1-11)
@other: a list containing: loc.metrics latlong ind.metrics history
loc.metrics.flags

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

{adegenet} Accessors



```
gl <- testset.gl
nInd(gl)
[1] 250
nLoc(gl)
[1] 250
nPop(gl)
[1] 30
popNames(gl) or indNames(gl) or locNames(gl)
```

table(pop(gl))

```
EmmacBrisWive EmmacBurdMist EmmacBurnBara EmmacClarJack EmmacClarYate EmmacCoopAvin EmmacCoopCully
10          10          11          5          5          10          10
EmmacCoopEulb EmmacFitzAllig EmmacJohnWari EmmacMacIGeom EmmacMaryBoru EmmacMaryPetr EmmacMDBBown
10          10          10          11          6          4          10
EmmacMDBCond EmmacMDBCudg EmmacMDBForb EmmacMDBGwyd EmmacMDBMaci EmmacMDBMurrMung EmmacMDBSanf
10          10          11          9          10          10          10
EmmacNormJack EmmacNormLeic EmmacNormSalt EmmacRichCasi EmmacRoss EmmacRussEube EmmacTweeUki
6           1           1           10          10          10          10
EmsubRopeMata EmvicVictJasp
5           5
```

as.matrix(gl)[1:7,1:5]

```
100049687-12-C/T 100049698-16-G/A 100049728-23-A/G 100049805-56-T/A 100049816-51-A/G
AA010915          2          NA          0          0          0
UC_00126          2          NA          0          0          0
AA032760          NA         NA          0          0          0
AA013214          2          NA          0          0          0
AA011723          2          NA          0          0          0
AA012411          2          NA          0          0          0
AA019237          2          NA          0          0          0
```

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

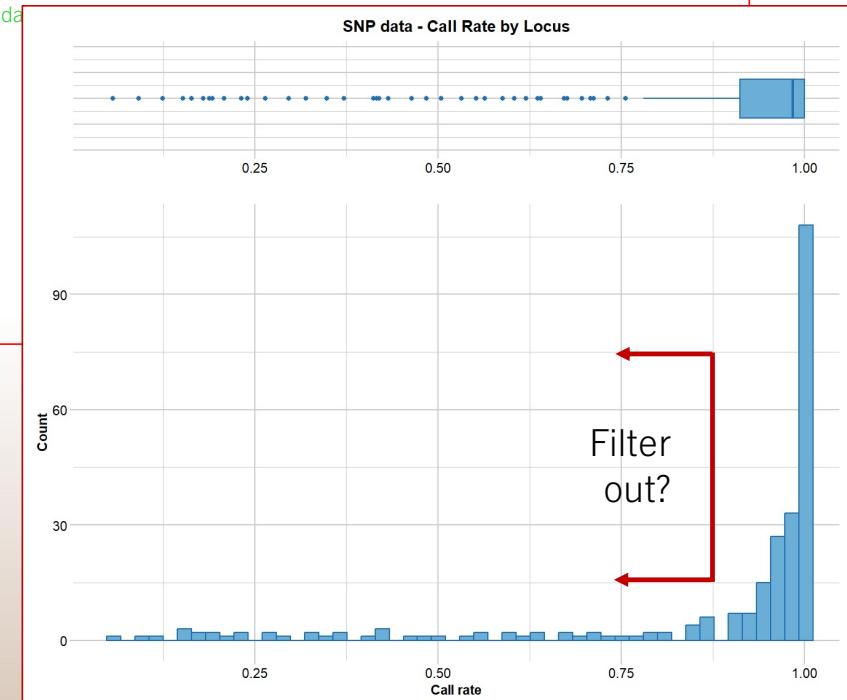
gl.report functions



```
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)
```

```
[Starting gl.report.callrate  
Processing genlight object with SNP data  
Reporting Call Rate by Locus  
No. of loci = 250  
No. of individuals = 250  
Minimum : 0.056  
1st quartile : 0.912  
Median : 0.984  
Mean : 0.8765804  
3rd quartile : 1  
Maximum : 1  
Missing Rate Overall: 0.1234
```

```
Completed: gl.report.callrate
```



dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

Who
will be
lost?



```
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")
```

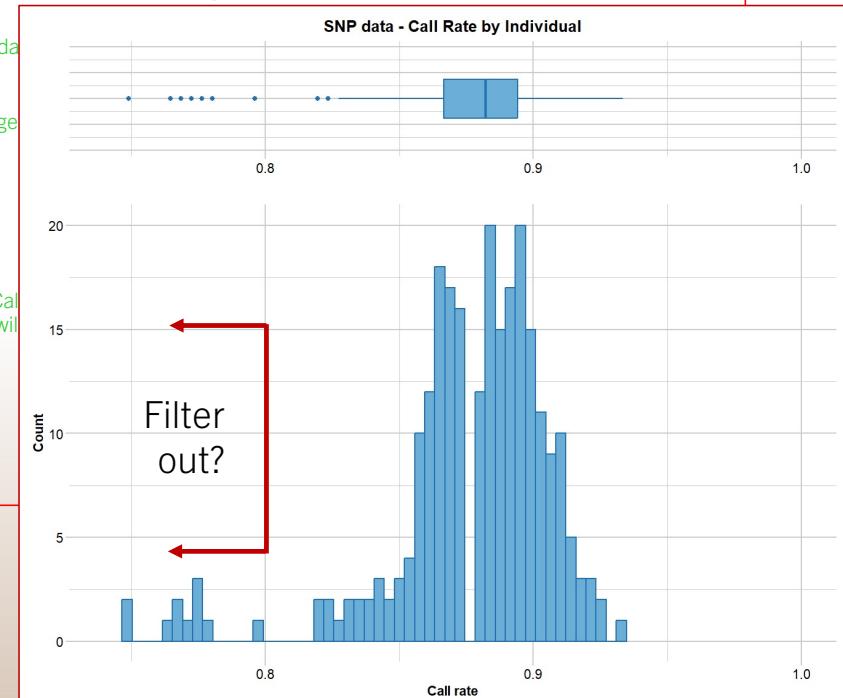
Starting gl.report.callrate
Processing genlight object with SNP da
Reporting Call Rate by Individual

Listing 30 populations and their average
Monitor again after filtering
Population CallRate N
1 EmmacBrisWive 0.8839 10
2 EmmacBurdMist 0.8808 10
3 EmmacBurnBara 0.8859 11
.....

Listing 20 individuals with the lowest CallRate
Use this list to see which individuals will be lost
Individual CallRate
1 AA063722 0.7490196
2 AA063726 0.7490196
3 AA063732 0.7647059
.....

Completed: gl.report.callrate

SNP data - Call Rate by Individual



dartR

Fundamentals

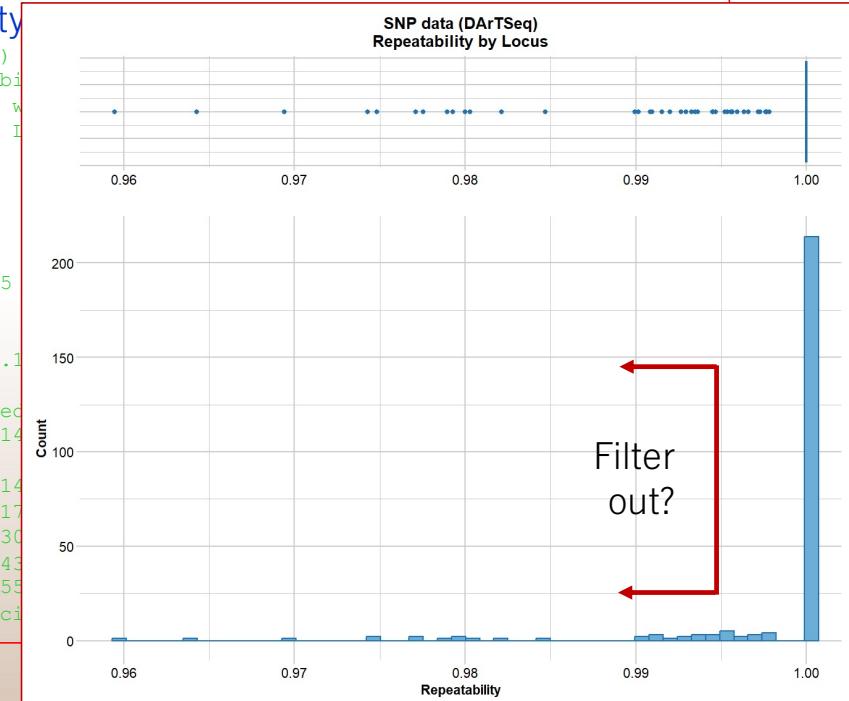
- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

gl.report functions



```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility
gl.report.reproducibility(gl)
Starting gl.report.reproducibility
Processing genlight object w...
Reporting Repeatability by L...
No. of loci = 255
No. of individuals = 250
  Minimum      : 0.959459
  1st quartile : 1
  Median       : 1
  Mean         : 0.9981525
  3rd quartile : 1
  Maximum      : 1
  Missing Rate Overall: 0.1

  Quantile Threshold Retained
1    100% 1.000000    214
.....
17   20% 1.000000    214
18   15% 0.997674    217
19   10% 0.994536    230
20    5% 0.984694    243
21    0% 0.959459    255
Completed: gl.report.reproducibility
```



dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

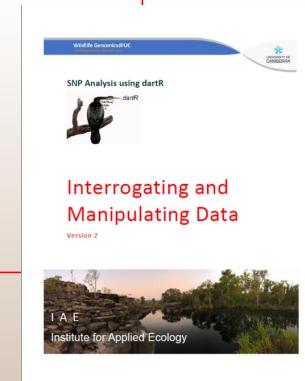


Try some of these in your own time



You get the gist

- `gl.report.callrate()`
- `gl.report.reproducibility()`
- `gl.report.secondaries()`
- `gl.report.rdepth()`
- `gl.report.monomorphs()`
- `gl.report.overhang()`
- `gl.report.hamming()`
- `gl.report.overshoot()`
- `gl.report.locmetric()`
- etc



Keep up ...

`testset.gl`

```
gl <- testset.gl
nlnd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

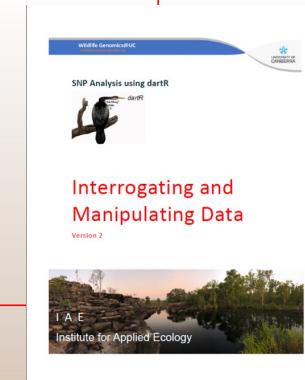


Try some of these in your own time

Subsetting your data

- `gl.keep.ind()`
- `gl.drop.ind()`
- `gl.keep.loc()`
- `gl.drop.loc()`
- `gl.keep.pop()`
- `gl.drop.pop()`
- `gl.merge.pop()`
- `gl.subsample.ind()`
- `gl.subsample.loc()`

..... etc



Keep up ...

`testset.gl`

```
gl <- testset.gl
nInd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials



Recall

- **gl.report.callrate()**
- **gl.report.reproducibility()**
- **gl.report.secondaries()**
- **gl.report.rdepth()**
- **gl.report.monomorphs()**
- **gl.report.overhang()**
- **gl.report.hamming()**
- **gl.report.overshoot()**
- **gl.report.locmetric()**
- etc



Keep up ...

testset.gl

```
gl <- testset.gl
nInd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials



Basic Filtering

- `gl <- gl.filter.callrate()`
- `gl <- gl.filter.reproducibility()`
- `gl <- gl.filter.secondaries()`
- `gl <- gl.filter.rdepth()`
- `gl <- gl.filter.monomorphs()`
- `gl <- gl.filter.overhang()`
- `gl <- gl.filter.hamming()`
- `gl <- gl.filter.overshoot()`
- `gl <- gl.report.locmetric()`
..... etc



Keep up ...

`testset.gl`

```
gl <- testset.gl
nlnd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

Who will be lost?



```
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")
```

```
Starting gl.report.callrate  
Processing genlight object with SNP da  
Reporting Call Rate by Individual
```

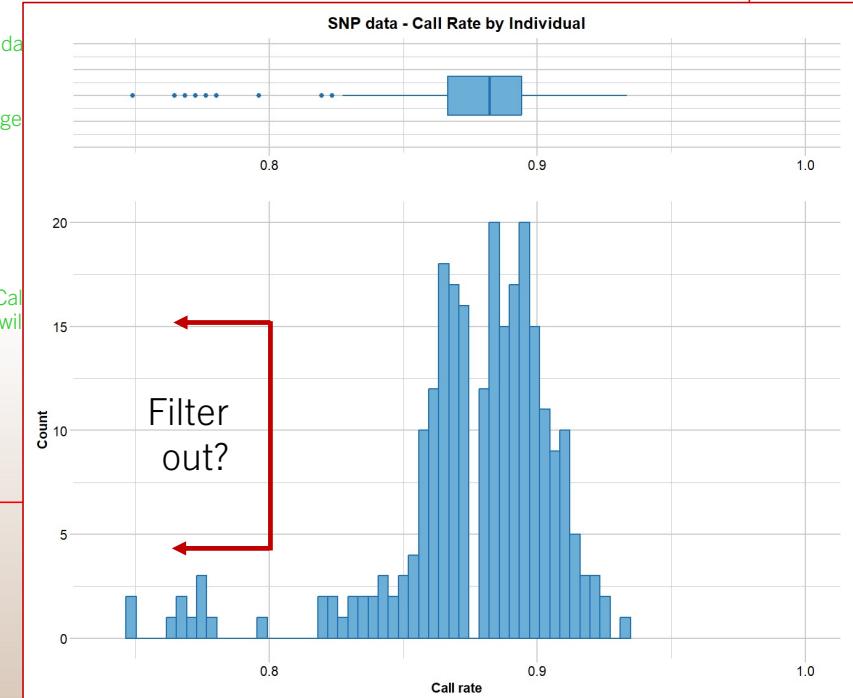
```
Listing 30 populations and their average  
Monitor again after filtering
```

```
Population CallRate N  
1 EmmacBrisWive 0.8839 10  
2 EmmacBurdMist 0.8808 10  
3 EmmacBurnBara 0.8859 11  
.....
```

```
Listing 20 individuals with the lowest CallRate  
Use this list to see which individuals will be lost  
Individual CallRate  
1 AA063722 0.7490196  
2 AA063726 0.7490196  
3 AA063732 0.7647059  
.....
```

```
Completed: gl.report.callrate
```

Filter out?



dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

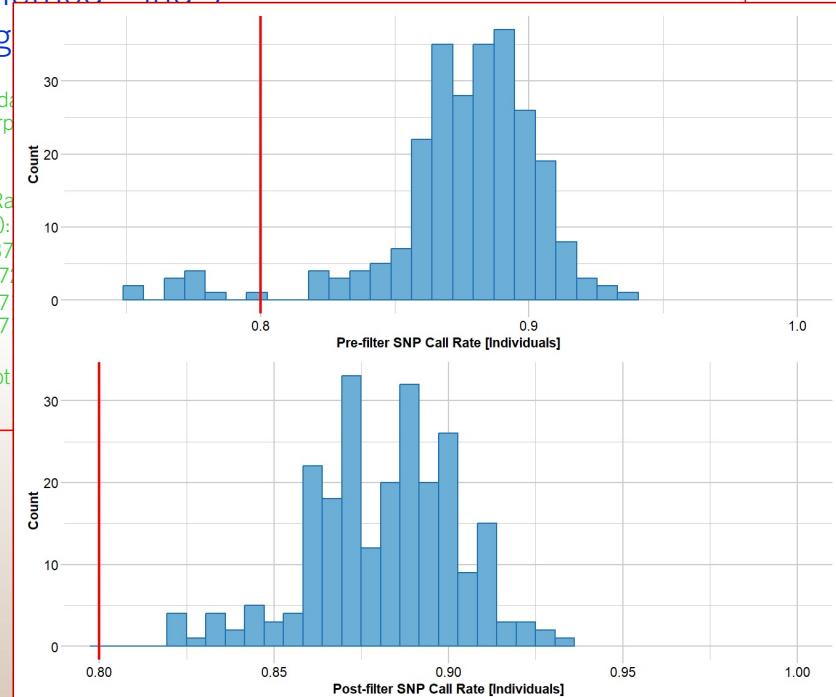
Who has
been
discarded?

gl.filter functions



```
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")  
gl <- gl.filter.callrate(g
```

```
Starting gl.filter.callrate  
Processing genlight object with SNP data  
Warning: Data may include monomorphic loci - these will not be included in calculations for filtering  
Recalculating Call Rate  
Removing individuals based on Call Rate  
Individuals deleted (CallRate <= 0.8):  
AA032760[EmmacMDBMacil], AA063714[EmmacCoopAvin], AA063715[EmmacCoopAvin], AA063716[EmmacCoopAvin], AA063717[EmmacCoopAvin], AA063718[EmmacCoopAvin], AA063719[EmmacCoopAvin], AA063720[EmmacCoopAvin], AA063721[EmmacCoopAvin], AA063722[EmmacCoopAvin], AA063723[EmmacCoopAvin], AA063724[EmmacCoopAvin], AA063725[EmmacCoopAvin], AA063726[EmmacCoopAvin], AA063727[EmmacCoopAvin], AA063728[EmmacCoopAvin], AA063729[EmmacCoopAvin], AA063730[EmmacCoopAvin], AA063731[EmmacCoopAvin], AA063732[EmmacCoopAvin], AA063733[EmmacCoopAvin], AA063734[EmmacCoopAvin], AA063735[EmmacCoopAvin], AA063736[EmmacCoopAvin], AA063737[EmmacCoopAvin], AA063738[EmmacCoopAvin], AA063739[EmmacCoopAvin], AA063740[EmmacCoopAvin], AA063741[EmmacCoopAvin], AA063742[EmmacCoopAvin], AA063743[EmmacCoopAvin], AA063744[EmmacCoopAvin], AA063745[EmmacCoopAvin], AA063746[EmmacCoopAvin], AA063747[EmmacCoopAvin], AA063748[EmmacCoopAvin], AA063749[EmmacCoopAvin], AA063750[EmmacCoopAvin], AA063751[EmmacCoopAvin], AA063752[EmmacCoopAvin], AA063753[EmmacCoopAvin], AA063754[EmmacCoopAvin], AA063755[EmmacCoopAvin], AA063756[EmmacCoopAvin], AA063757[EmmacCoopAvin], AA063758[EmmacCoopAvin], AA063759[EmmacCoopAvin], AA063760[EmmacCoopAvin], AA063761[EmmacCoopAvin], AA063762[EmmacCoopAvin], AA063763[EmmacCoopAvin], AA063764[EmmacCoopAvin], AA063765[EmmacCoopAvin], AA063766[EmmacCoopAvin], AA063767[EmmacCoopAvin], AA063768[EmmacCoopAvin], AA063769[EmmacCoopAvin], AA063770[EmmacCoopAvin], AA063771[EmmacCoopAvin], AA063772[EmmacCoopAvin], AA063773[EmmacCoopAvin], AA063774[EmmacCoopAvin], AA063775[EmmacCoopAvin], AA063776[EmmacCoopAvin], AA063777[EmmacCoopAvin], AA063778[EmmacCoopAvin], AA063779[EmmacCoopAvin], AA063780[EmmacCoopAvin], AA063781[EmmacCoopAvin], AA063782[EmmacCoopAvin], AA063783[EmmacCoopAvin], AA063784[EmmacCoopAvin], AA063785[EmmacCoopAvin], AA063786[EmmacCoopAvin], AA063787[EmmacCoopAvin], AA063788[EmmacCoopAvin], AA063789[EmmacCoopAvin], AA063790[EmmacCoopAvin], AA063791[EmmacCoopAvin], AA063792[EmmacCoopAvin], AA063793[EmmacCoopAvin], AA063794[EmmacCoopAvin], AA063795[EmmacCoopAvin], AA063796[EmmacCoopAvin], AA063797[EmmacCoopAvin], AA063798[EmmacCoopAvin], AA063799[EmmacCoopAvin], AA063710[EmmacCoopAvin], AA063711[EmmacCoopAvin], AA063712[EmmacCoopAvin], AA063713[EmmacCoopAvin], AA063714[EmmacCoopAvin], AA063715[EmmacCoopAvin], AA063716[EmmacCoopAvin], AA063717[EmmacCoopAvin], AA063718[EmmacCoopAvin], AA063719[EmmacCoopAvin], AA063720[EmmacCoopAvin], AA063721[EmmacCoopAvin], AA063722[EmmacCoopAvin], AA063723[EmmacCoopAvin], AA063724[EmmacCoopAvin], AA063725[EmmacCoopAvin], AA063726[EmmacCoopAvin], AA063727[EmmacCoopAvin], AA063728[EmmacCoopAvin], AA063729[EmmacCoopAvin], AA063730[EmmacCoopAvin], AA063731[EmmacCoopAvin], AA063732[EmmacCoopAvin], AA063733[EmmacCoopAvin], AA063734[EmmacCoopAvin], AA063735[EmmacCoopAvin], AA063736[EmmacCoopAvin], AA063737[EmmacCoopAvin], AA063738[EmmacCoopAvin], AA063739[EmmacCoopAvin], AA063740[EmmacCoopAvin], AA063741[EmmacCoopAvin], AA063742[EmmacCoopAvin], AA063743[EmmacCoopAvin], AA063744[EmmacCoopAvin], AA063745[EmmacCoopAvin], AA063746[EmmacCoopAvin], AA063747[EmmacCoopAvin], AA063748[EmmacCoopAvin], AA063749[EmmacCoopAvin], AA063750[EmmacCoopAvin], AA063751[EmmacCoopAvin], AA063752[EmmacCoopAvin], AA063753[EmmacCoopAvin], AA063754[EmmacCoopAvin], AA063755[EmmacCoopAvin], AA063756[EmmacCoopAvin], AA063757[EmmacCoopAvin], AA063758[EmmacCoopAvin], AA063759[EmmacCoopAvin], AA063760[EmmacCoopAvin], AA063761[EmmacCoopAvin], AA063762[EmmacCoopAvin], AA063763[EmmacCoopAvin], AA063764[EmmacCoopAvin], AA063765[EmmacCoopAvin], AA063766[EmmacCoopAvin], AA063767[EmmacCoopAvin], AA063768[EmmacCoopAvin], AA063769[EmmacCoopAvin], AA063770[EmmacCoopAvin], AA063771[EmmacCoopAvin], AA063772[EmmacCoopAvin], AA063773[EmmacCoopAvin], AA063774[EmmacCoopAvin], AA063775[EmmacCoopAvin], AA063776[EmmacCoopAvin], AA063777[EmmacCoopAvin], AA063778[EmmacCoopAvin], AA063779[EmmacCoopAvin], AA063780[EmmacCoopAvin], AA063781[EmmacCoopAvin], AA063782[EmmacCoopAvin], AA063783[EmmacCoopAvin], AA063784[EmmacCoopAvin], AA063785[EmmacCoopAvin], AA063786[EmmacCoopAvin], AA063787[EmmacCoopAvin], AA063788[EmmacCoopAvin], AA063789[EmmacCoopAvin], AA063790[EmmacCoopAvin], AA063791[EmmacCoopAvin], AA063792[EmmacCoopAvin], AA063793[EmmacCoopAvin], AA063794[EmmacCoopAvin], AA063795[EmmacCoopAvin], AA063796[EmmacCoopAvin], AA063797[EmmacCoopAvin], AA063798[EmmacCoopAvin], AA063799[EmmacCoopAvin]
```



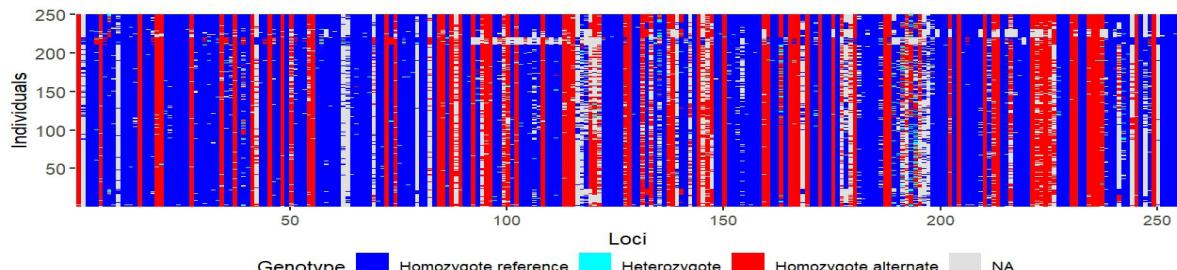
dartR

Fundamentals

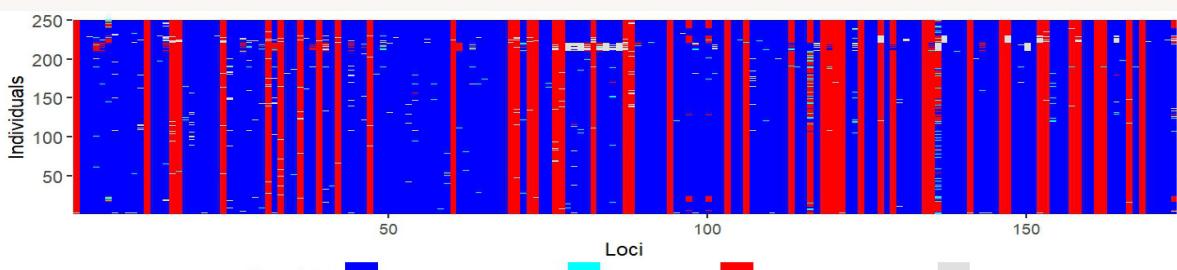
- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials

Quick graphical evaluation

 `gl <- testset.gl`
`gl.smearplot(gl)`



 `gl <- gl.filter.callrate(gl,verbose=0)`
`gl <- gl.filter.callrate(gl,method="ind",threshold=0.80,verbose=0)`
`gl.smearplot(gl)`



dartR Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials



Try some of these in your own time



Basic Filtering

- `gl <- gl.filter.callrate()`
- `gl <- gl.filter.reproducibility()`
- `gl <- gl.filter.secondaries()`
- `gl <- gl.filter.rdepth()`
- `gl <- gl.filter.monomorphs()`
- `gl <- gl.filter.overhang()`
- `gl <- gl.filter.hamming()`
- `gl <- gl.filter.overshoot()`
- `gl <- gl.report.locmetric()`
..... etc



Keep up ...

`testset.gl`

```
gl <- testset.gl
nlnd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

`gl <- testset.gl`

```
gl.set.verbosity(3)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

```
gl.report.callrate(gl,method="ind")
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80)
```

`gl <- testset.gl`

```
gl.smearplot(gl)
gl <- gl.filter.callrate(gl,verbose=0)
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80,
verbose=0)
gl.smearplot(gl)
```

dartR

Fundamentals

- Data structure
- Interrogation
- Subsetting
- Basic Filtering
(refer Renee's sessions)
- Tutorials
 - inline help
 - ?function
 - CRAN documentation

Tutorials

The collage consists of several dartR tutorial slides arranged in a grid-like pattern. The top row contains two slides: 'RStudio Refresher' (Version 3) on the left and 'Interrogating and Manipulating Data' (Version 2) on the right. The middle row contains two slides: 'How dartR stores your data and data input' on the left and 'Guide to Basic Filtering' (Version 3) on the right. The bottom row contains two slides: 'Installing dartR' on the left and 'SNP Analysis using dartR' on the right. Each slide features a blue header bar with the 'Wildlife Genomics' logo and the 'dartR' icon. The body of each slide contains text and images related to the specific tutorial topic.



testset.gl

```
gl <- testset.gl
nlnd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

```
gl.report.callrate(gl,method="ind")
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80)
```

```
gl <- testset.gl
gl.smearplot(gl)
gl <- gl.filter.callrate(gl,verbose=0)
gl <- gl.filter.callrate(gl,verbose=0,
method="ind", threshold=0.80,
verbose=0)
gl.smearplot(gl)
```

Filtering strategies (or the Art of Filtering)

- One set of SNPs and one set of filters cannot answer all your questions about a single species
- It will give you an answer but there is a good chance it is wrong

This week I will cover:

1. What are the basic filters and considerations? (this session)
2. What are some key considerations for filtering for specific conservation questions? (session 2)
3. What do you need to think about when calling SNPs? (session 3)

Population genomics and sexual signals support reproductive character displacement in *Uperoleia* (Anura: Myobatrachidae) in a contact zone

Frederick R. Jaya, Jessie C. Tanner, Michael R. Whitehead, Paul Doughty, J. Scott Keogh, Craig C. Moritz, Renee A. Catullo 

Genotyping using commercial *DarTseq™* 1.0 and SNP calling algorithm
Diversity Arrays Ltd., Canberra

SNP calling 1: Target species together including hybrids
Purpose: Identification of interspecific hybrid individuals and assessing level of interspecific admixture
Total unfiltered SNPs: 233,691
Missing data: 43.95%
Avg loci call rate: 0.560 [0.212-1]
Mean reproducibility: 0.997 [0.875-1]
Mean read depth: 7.171 [2.5-154.3]

Analysis: Structure
Filtering:

1. Filter loci by call rate (0.95)
2. Remove singletons
3. Filter individuals by call rate (0.9)
4. Filter by reproducibility (0.99)
5. Filter by average fragment read depth (lower = 5, upper = 19)
6. Remove monomorphic loci
7. Filter to one SNP per fragment (method = best)

Total filtered SNPs: 3,956

Analysis: NewHybrids
Filtering:

1. Further filter SNPs from Structure using gl.nhybrids

Total filtered SNPs: 186

SNP calling 2: Each target species called separately excluding interspecific hybrids
Purpose: Intraspecific population structure and diversity analyses
Total unfiltered SNPs: 112,381 (*U. borealis*), 227,145 (*U. crassa*)
Missing data: 37.47% (*U. borealis*), 25.36% (*U. crassa*)
Avg loci call rate: 0.625 [0.21-1] (*U. borealis*), 0.746 [0.22-1] (*U. crassa*)
Mean reproducibility: 1 (*U. borealis*), 0.998 (*U. crassa*)
Mean read depth: 8.161 [2.5-149.7] (*U. borealis*), 6.718 [2.5-121.7] (*U. crassa*)

Analysis: Structure
Filtering:

1. Filter loci by call rate (0.95)
2. Remove singletons
3. Filter individuals by call rate (0.9)
4. Filter by reproducibility (0.99)
5. Filter by average fragment read depth (lower = 5, upper = 19)
6. Remove monomorphic loci
7. Filter to one SNP per fragment (method = best)

Total filtered SNPs: 1,938

SNP calling 3: Target species individuals plus additional *Uperoleia* species in clade
Purpose: Phylogenetic analyses
Total unfiltered SNPs: 269,303
Missing data: 60.5%
Avg loci call rate: 0.395 [0.146-0.476]
Mean reproducibility: 0.996 [0-1]
Mean read depth: 6.348 [1.5-182.7]

Analysis: IQtree (ML phylogeny)
Filtering:

1. Remove hybrid individuals
2. Remove monomorphic loci
3. Filter by reproducibility (0.99)
4. Filter by average fragment read depth (lower = 11, upper = 17)
5. Filter loci by call rate (0.9)
6. Remove singletons
7. Select one SNP per fragment (method = best)

Total filtered SNPs: 1,938

&

Analysis: SNAPP (Species tree)
Filtering:

1. Subset target species to four individuals per clade
2. Steps as per IQtree filtering

Total filtered SNPs: 1,162

Analyses: Heterozygosity & Tajima's D
Filtering:

1. Filter loci by call rate (0.95)
2. Filter individuals by call rate (0.9)
3. Filter by reproducibility (0.99)
4. Filter by average fragment read depth (lower = 8, upper = 19)
5. Remove monomorphic loci

Total filtered SNPs: 21,603 (*U. crassa*)

Filtering for call rate (method = “loc”)



What is it?

- Removing SNPs because they didn't work across individuals to a set level

Why is it important?

- Missing loci can add “noise” and computation time to an analysis

Key Considerations:

- Trade-off between strength of filter and number of loci
- Call rate filters matter for different metrics
 - Strong filters disproportionately remove heterozygous sites
- Many population structure methods work better without high levels of missing data

		A/A	C/G	T/T	-	A/T
		A/C	C/G	T/T	G/G	A/A
		A/A	C/C	T/T	-	A/A
		A/C	C/G	T/T	G/A	A/A
		A/A	C/C	T/T	-	A/T
		A/A	C/C	T/T	G/G	A/A
		A/A	C/G	A/T	-	A/A
		A/C	C/C	A/T	G/A	A/T

```
testset.gl  
  
gl <- testset.gl  
nlnd(gl)  
nLoc(gl)  
nPop(gl)  
popNames(gl)  
indNames(gl)  
locNames(gl)  
table(pop(gl))  
as.matrix(gl)[1:7,1:5]  
  
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")  
gl.report.reproducibility(gl)  
  
gl.report.callrate(gl,method="ind")  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80)  
  
gl <- testset.gl  
gl.smearplot(gl)  
gl <- gl.filter.callrate(gl,verbose=0)  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80,  
verbose=0)  
gl.smearplot(gl)
```

Filtering for call rate (method = “ind”)



What is it?

- Removes *individuals* that did not sequence to the specified level of completeness

Why is it important?

- Deletes the key thing you need to do a study
- Removes individuals that may be misleading as they are sequencing outliers

Key Considerations:

- The filtering step that costs you the very most amount of money
- Generally not a first step – recommend filtering lightly here in the first go and seeing if other filters improve the individuals completeness

		A/A	C/G	T/T	G/A	A/T
		A/C	C/G	T/T	G/G	A/A
		A/A	C/C	T/T	G/A	A/A
		A/C	C/G	T/T	G/A	A/A
		-	-	T/T	G/G	-

		A/A	C/C	T/T	G/G	A/A
		A/A	C/G	A/T	G/G	A/A
		A/C	C/C	A/T	G/A	A/T

testset.gl

```
gl <- testset.gl
nlnd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

```
gl.report.callrate(gl,method="ind")
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80)
```

```
gl <- testset.gl
gl.smearplot(gl)
gl <- gl.filter.callrate(gl,verbose=0)
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80,
verbose=0)
gl.smearplot(gl)
```

Filtering for reproducibility (gl.filter.reproducibility)



Keep up ...

```
testset.gl  
gl <- testset.gl  
nInd(gl)  
nLoc(gl)  
nPop(gl)  
popNames(gl)  
indNames(gl)  
locNames(gl)  
table(pop(gl))  
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")  
gl.report.reproducibility(gl)
```

```
gl.report.callrate(gl,method="ind")  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80)
```

```
gl <- testset.gl  
gl.smearplot(gl)  
gl <- gl.filter.callrate(gl,verbose=0)  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80,  
verbose=0)  
gl.smearplot(gl)
```

What is it?

- Diversity Arrays duplicates individuals during library prep, to assess whether the same answer is found for every locus
- A control

Why is it important?

- Provides confidence in your base calls which is pretty fundamental

Key Considerations:

- If reproducibility is low, you need your SNPs re-assessed
- If doing sequencing through AGRF or other facilities, make sure to include duplicated individuals (this is an additional cost but do it anyway)

Filtering for read depth (gl.filter.rdepth)



What is it?

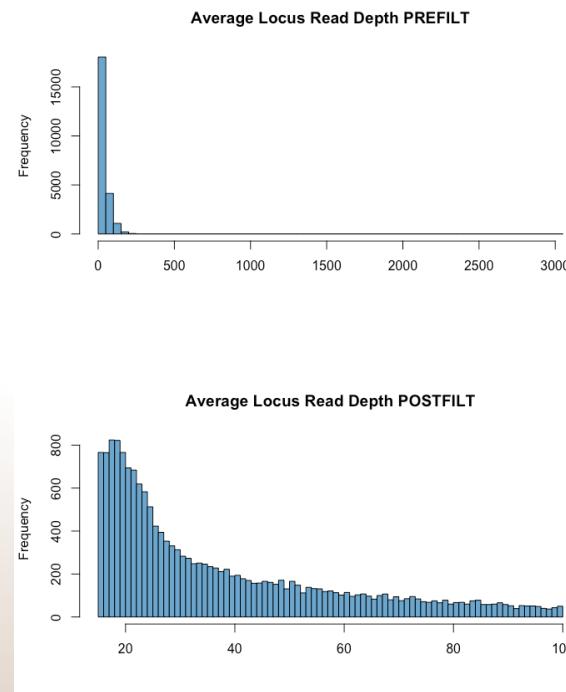
- The mean number of sequencing reads for a particular locus, across all individuals

Why is it important?

- Number of reads in a stack tells you how confident you can be in your base calls
- A low read depth means that your heterozygous sites have pretty low coverage

Key Considerations:

- If your read depth is not high enough, there are questions you can't confidently answer (e.g., heterozygosity)
- Very high read depth suggests paralogs (genes with multiple copies) being assembled in to one fragment



testset.gl

```
gl <- testset.gl
nlnd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl)
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

```
gl.report.callrate(gl,method="ind")
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80)
```

```
gl <- testset.gl
gl.smearplot(gl)
gl <- gl.filter.callrate(gl,verbose=0)
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80,
verbose=0)
gl.smearplot(gl)
```

Filtering out minor alleles (gl.filter.maf)

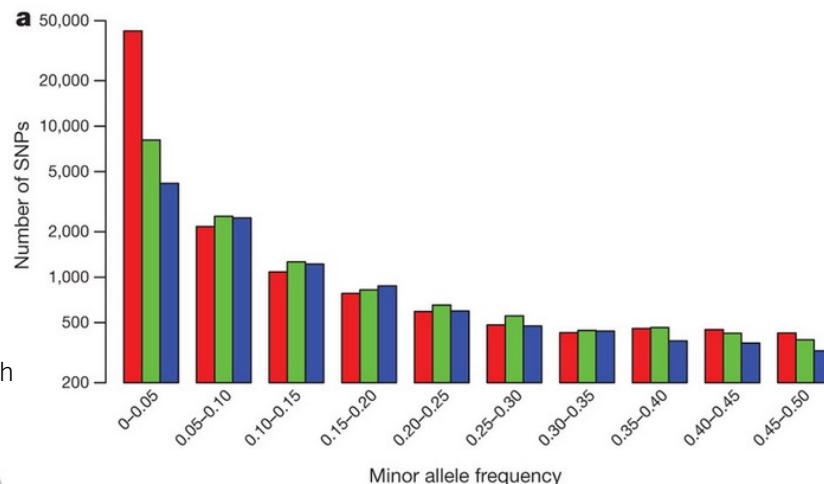


What is it?

- Minor allele frequency is the frequency of the second most common allele in a given population
- Removing SNPs based on their relative proportion
- In the literature, often 0.05 BUT...

Key Considerations:

- 0.05 can be high if you have lot of individuals
 - If you have 300 diploid individuals, then you delete SNPs with less than 30 copies of the minor allele
- Or singleton if you have a small number of individuals
- Should depend on your question
 - Rare alleles can be very important to certain questions
 - Most SNPs are rare alleles, as are most heterozygous sites
 - Rare alleles are important in expansion processes
 - NOT important in structure or phylogenetic analyses



```
testset.gl
```

```
gl <- testset.gl  
nlnd(gl)  
nLoc(gl)  
nPop(gl)  
popNames(gl)  
indNames(gl)  
locNames(gl)  
table(pop(gl))  
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl,method="ind")  
gl.report.callrate(gl,method="ind")  
gl.report.reproducibility(gl)  
gl.report.callrate(gl,method="ind")  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80)
```

```
gl <- testset.gl  
gl.smearplot(gl)  
gl <- gl.filter.callrate(gl,verbose=0)  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80,  
verbose=0)  
gl.smearplot(gl)
```

Filtering secondaries (gl.filter.secondaries)



What is it?

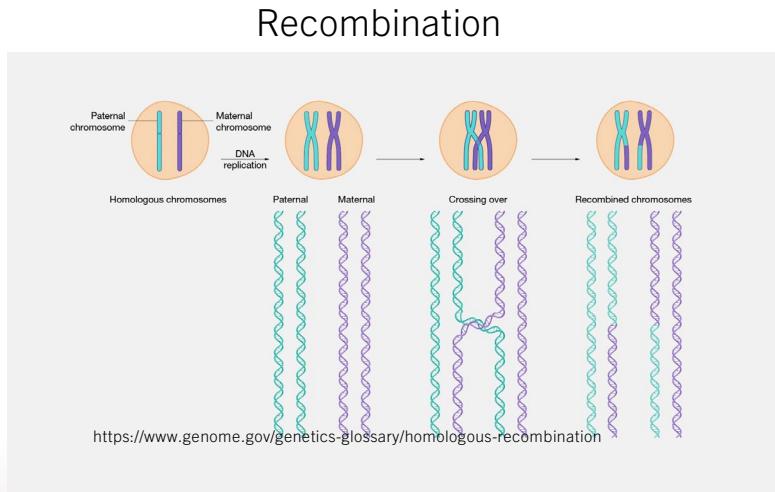
- When there are two SNPs on a single fragment, choosing to keep only one.

Why is it important?

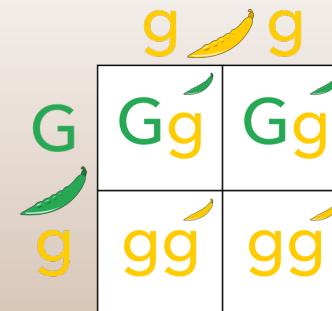
- Loci that are very close together in the genome are not independently inherited as they are too close together for recombination to split.

Key Considerations:

- Can mislead structure and phylogenetic analyses
- Can improve representations of heterozygosity



Genetic diversity



```
testset.gl  
gl <- testset.gl  
nlnd(gl)  
nLoc(gl)  
nPop(gl)  
popNames(gl)  
indNames(gl)  
locNames(gl)  
table(pop(gl))  
as.matrix(gl)[1:7,1:5]  
  
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")  
gl.report.reproducibility(gl)  
  
gl.report.callrate(gl,method="ind")  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80)  
  
gl <- testset.gl  
gl.smearplot(gl)  
gl <- gl.filter.callrate(gl,verbose=0)  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80,  
verbose=0)  
gl.smearplot(gl)
```

What's the right filtering order?

- There isn't
- Be iterative – test different options
- Might do the same filter twice (call rate)

As a starting point I would get rid of loci I don't believe in and then individuals that didn't work properly at all:

- gl.filter.callrate(gl, method = "loc", threshold = 0.7)

Get rid of really poorly sequenced loci

But don't cut hard

- gl.filter.callrate(gl, method = "ind", threshold = 0.25)

Very low filter – this is only to get rid of your really bad individuals

- gl.filter.monomorphs(gl)

Always run this after removing individuals – removes loci that are no longer variable

- gl.filter.reproducibility(gl)

Get rid of unreliable loci

- gl.filter.rdepth(gl, lower = X, higher = X)

Get rid of low and super high read depth loci



testset.gl

```
gl <- testset.gl  
nlnd(gl)  
nLoc(gl)  
nPop(gl)  
popNames(gl)  
indNames(gl)  
locNames(gl)  
table(pop(gl))  
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")  
gl.report.reproducibility(gl)
```

```
gl.report.callrate(gl,method="ind")  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80)
```

```
gl <- testset.gl  
gl.smearplot(gl)  
gl <- gl.filter.callrate(gl,verbose=0)  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80,  
verbose=0)  
gl.smearplot(gl)
```

What's the right filtering order?

- There isn't
- Be iterative – test different options
- Might do the same filter twice (call rate)

Then I would filter more strongly as appropriate for my question. For a population structure analysis I would,

- gl.filter.callrate(gl, method = "loc", threshold = 0.95)

Structure dislikes missing data

- gl.filter.maf(gl, threshold = 1/(2*nInd(gl))

I usually set up the threshold so it is just removing singletons to improve computation time

- gl.filter.secondaries(gl)

Always do this as the last loci filter so that you've cut for quality before you cut because there are two SNPs

- gl.filter.callrate(gl, method = "ind", threshold = .9)

Filter on individuals. You can usually be a bit flexible at this point.

- gl.filter.monomorphs(gl)

Always run this after removing individuals



testset.gl

```
gl <- testset.gl
nInd(gl)
nLoc(gl)
nPop(gl)
popNames(gl)
indNames(gl)
locNames(gl)
table(pop(gl))
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl
gl.set.verbosity(3)
gl.report.callrate(gl,method="ind")
gl.report.callrate(gl,method="ind")
gl.report.reproducibility(gl)
```

```
gl.report.callrate(gl,method="ind")
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80)
```

```
gl <- testset.gl
gl.smearplot(gl)
gl <- gl.filter.callrate(gl,verbose=0)
gl <- gl.filter.callrate(gl,
method="ind", threshold=0.80,
verbose=0)
gl.smearplot(gl)
```

What is right for one situation is probably wrong for a different situation

- Be flexible
- Know what each filter is doing to your data
- Think carefully about whether the filter is appropriate to the test you want to run
- Analyse your data many different ways
- Don't over-interpret your PCoA



```
testset.gl  
gl <- testset.gl  
nInd(gl)  
nLoc(gl)  
nPop(gl)  
popNames(gl)  
indNames(gl)  
locNames(gl)  
table(pop(gl))  
as.matrix(gl)[1:7,1:5]  
  
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")  
gl.report.reproducibility(gl)  
  
gl.report.callrate(gl,method="ind")  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80)  
  
gl <- testset.gl  
gl.smearplot(gl)  
gl <- gl.filter.callrate(gl,verbose=0)  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80,  
verbose=0)  
gl.smearplot(gl)
```



Conclusion

Where have we come?

Discussion



Keep up ...

```
testset.gl
```

```
gl <- testset.gl  
nInd(gl)  
nLoc(gl)  
nPop(gl)  
popNames(gl)  
indNames(gl)  
locNames(gl)  
table(pop(gl))  
as.matrix(gl)[1:7,1:5]
```

```
gl <- testset.gl  
gl.set.verbosity(3)  
gl.report.callrate(gl)  
gl.report.callrate(gl,method="ind")  
gl.report.reproducibility(gl)
```

```
gl.report.callrate(gl,method="ind")  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80)
```

```
gl <- testset.gl  
gl.smearplot(gl)  
gl <- gl.filter.callrate(gl,verbose=0)  
gl <- gl.filter.callrate(gl,  
method="ind", threshold=0.80,  
verbose=0)  
gl.smearplot(gl)
```