# Analysing Genomic Data with **dartRverse**: Accessible Tools for Conservation

ICCB

JUNE 15-20 2025



UNIVERSITY OF **CANBERRA**

**DArT** | diversity arrays technology

MONASH University

Centre for Biodiversity Analysis

# Session 6: SNP Panel Selection

Elise Furlan

Andrzej Kilian

Bernd Gruber

Elise.Furlan@canberra.edu.au

# Aim

Provide the knowledge and practical tools to reduce large SNP datasets into smaller, targeted SNP panels for conservation monitoring.

# Background

- Targeted sets of SNP markers – 10s to100s

- Reproducible

- Cost-effective for high sample volume

- Suitable for low-quality or low-quantity DNA samples

# Purpose

SNP panels can be used to address either specific questions or to span multiple conservation genetic applications
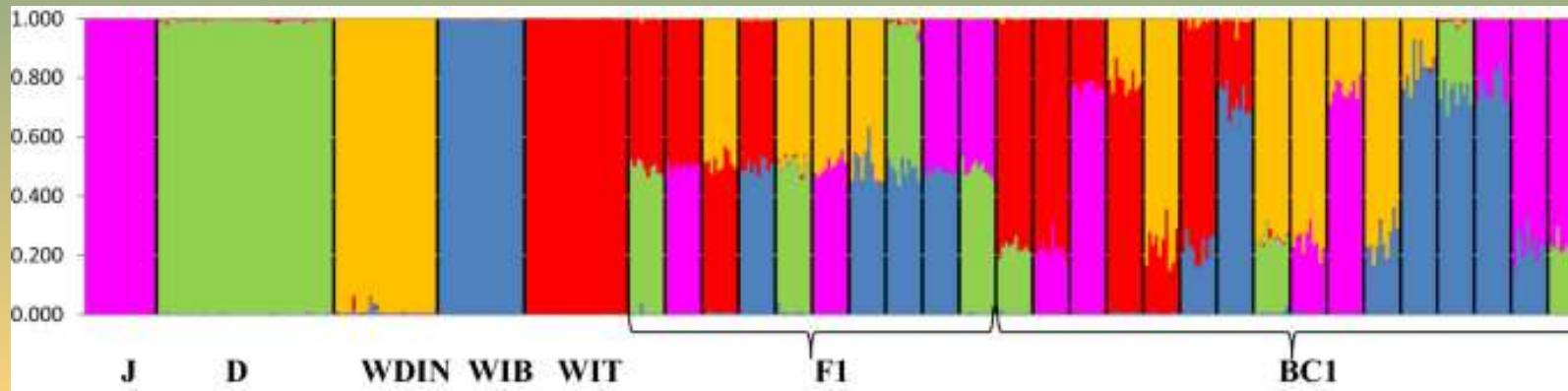
- **Population assignment**

- **Parentage or relatedness**

- **Individual ID**

- **Hybridisation**

- **Additional metrics**
  - Sex-linked SNPs
  - Candidate adaptive markers
  - Diagnostic SNPs for population ID
  - Phenotypic markers



Cheetah Conservation Fund, Magliolo et al. 2021

# Example - Hybridisation

- 192 SNP genotypes

- Differentiate 5 canid species.
  - Jackals (J), dogs (D), Dinaric wolves (WDIN), Iberian wolves (WIB) and Italian wolves (WIT)

- Identify hybrids
  - First-generation (F1) hybrid and first-generation backcross (BC1) genotypes

- Included 3x phenotypic markers relating to coat colour, nail colour and dewclaw presence (absent in wild canids)
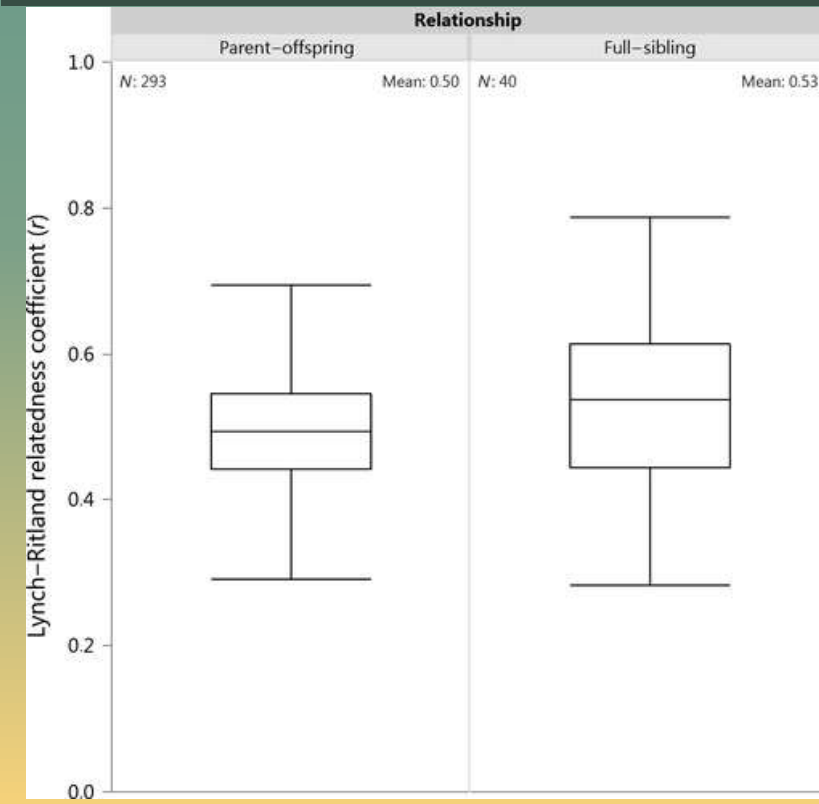
Stronen et al., 2022

# Example - Non-invasive samples

- 96-SNPs used on faecal samples in brown bears

- Estimated population size
  - Fell within the 95% CI of Capture-Mark-Recapture estimates

- Estimated relatedness

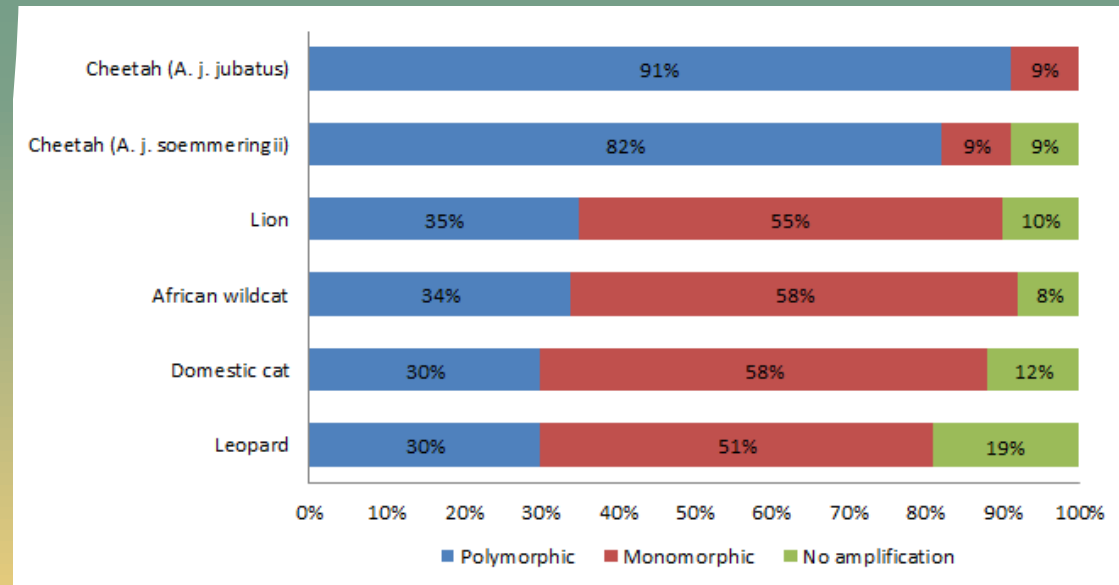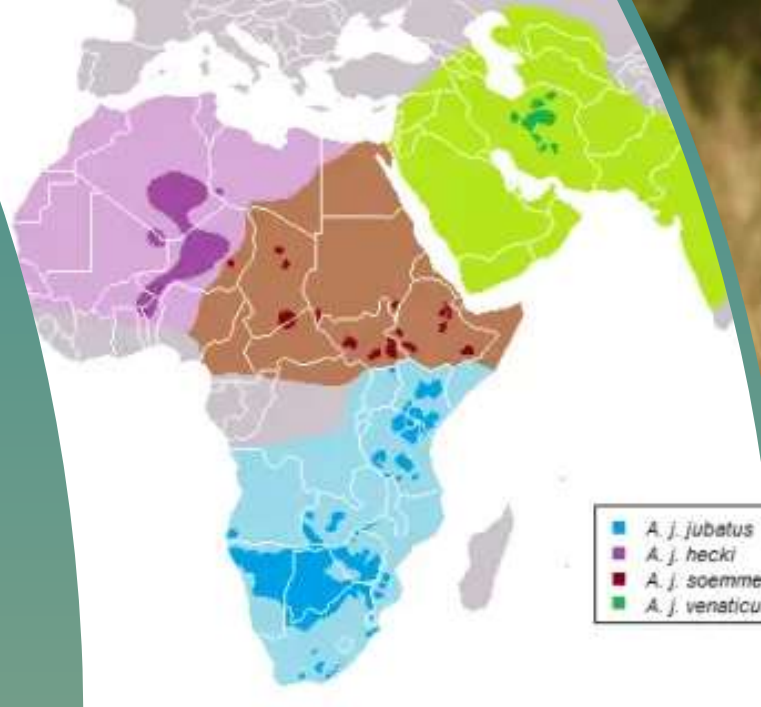- Determined sex
  - using sex-linked markers



Spitzer et al., 2016, Source: Nyhetsbyrån

# Marker Selection



- SNP panels must be carefully selected to maximise informativeness for the specific application
  - Less room for redundancy

- Requires genome-wide SNP data from individuals spanning the full distribution
  - Avoid ascertainment bias and loss of power in other species/populations
  - SNP panels can be expanded later if needed



Map legend:
- A. j. jubatus
- A. j. hecki
- A. j. soemmeri
- A. j. venaticus

Chart:
- Cheetah (A. j. jubatus): Polymorphic 91%, Monomorphic 9%
- Cheetah (A. j. soemmeringii): Polymorphic 82%, Monomorphic 9%, No amplification 9%
- Lion: Polymorphic 35%, Monomorphic 55%, No amplification 10%
- African wildcat: Polymorphic 34%, Monomorphic 58%, No amplification 8%
- Domestic cat: Polymorphic 30%, Monomorphic 58%, No amplification 12%
- Leopard: Polymorphic 30%, Monomorphic 51%, No amplification 19%

Legend: Polymorphic · Monomorphic · No amplification

Magliolo et al. 2021

# Considerations

- Targeted SNP panels address specific applications
  - may not support broader analyses
- Reduced representation
  - some genetic signals may be lost (e.g., selection, subtle structure)
- A new measure of genetic diversity
  - Not comparable to genome-wide diversity

# Requirements:

- Existing genome-wide SNP data
- Good geographic coverage

# *When to use SNP Panels*

✅ Use SNP panels when:

- monitoring large numbers of individuals,
- long-term surveillance,
- DNA samples are low quality or degraded (e.g., scats, feathers, eDNA)

❌ Avoid SNP panels when:

- sample sizes are small,
- genome-wide resolution is required (e.g., adaptation studies)
- genetically distant populations will be targeted

# Using dartR for SNP panel selection

- Purpose: To select a subset of informative SNPs

- Versatility
  - Modify the *number* of SNPs
  - Find the best panel to address one *specific* conservation question or *multiple* conservation questions

# Key metrics for SNP selection

- Population structure, population assignment
    - **dapc**: Select loci contributing most to discrimination between populations using DAPC (Discriminant Analysis of Principal Components).
    - **pahigh**: Select loci with private alleles having high frequency i.e., diagnostic
    - **monopop**: Select monomorphic loci within populations i.e., fixed
- Individual-level resolution e.g., individual IDs, parentage, relatedness
    - **PIC**. Select loci with high Polymorphic Information Content i.e., high minor allele frequency.
    - **PICdart**. Similar to PIC but based on allele presence/absence rather than frequencies.
- Heterozygosity estimates e.g., diversity, inbreeding
    - **hafall**: Select loci with the highest minor allele frequencies across all populations. These are likely to be more polymorphic and informative across all populations.
    - **hafpop**: Select loci with the highest minor allele frequencies within each population. Increases within-population informativeness e.g., within-population diversity
- Genome-wide diversity
    - **random**: Randomly select loci. Provides an unbiased snapshot of diversity across the genome
    - **stratified**: Stratified sampling of loci based on allele frequencies. Similar to random but ensures broad coverage of genetic variation

# Evaluate panel performance

- The final panel can be checked for concordance with:
  - $F_{ST}$ - genetic differentiation
  - $F_{IS}$ - inbreeding
  - $N_{ALL}$ - number of alleles
  - $H_E$ – expected heterozygosity
  - $H_O$ – observed heterozygosity
  - $N_E$ – effective population size



SNP panel estimate

$R^2 = 0.87$

Genome-wide SNP estimate

# Next steps

- Select sequence provider and prepare data according to their requirements

- Primer design
  - Bioinformatics to avoid primer interactions

- Lab testing
  - Identify over-amplified or under-amplified loci. Filter as necessary and retest

# Genomic Services @ DArT

Started with DNA array-based methods but moved to using Next Generation **Sequencing (NGS) supported by DArTdb/LIMS application**

## 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
- "De novo" and "SNP recall" analytical pipelines – reference free

## 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

## DArTreseq and WG sequencing – gene cloning and pangenome construction
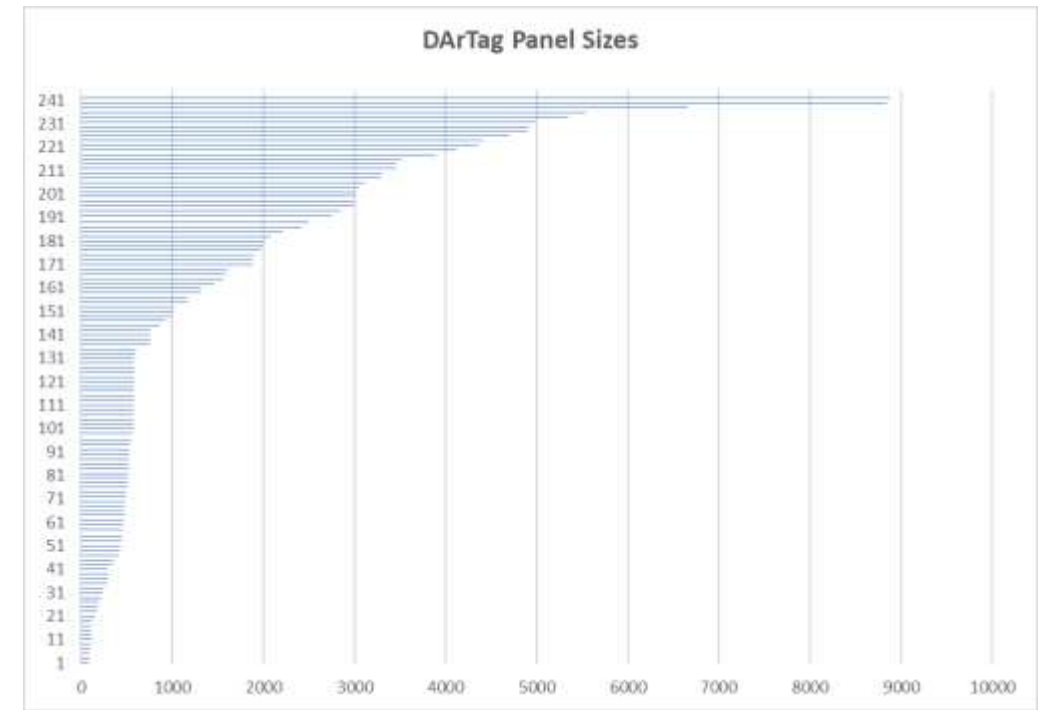
# DArTag platform

- **Adopted padlock molecule concept to capture target region**

- **Three step process:**

  1. capture of region with SNP/INDEL via padlock oligos
  2. addition of sample barcode and flowcell attchement via PCR
  3. Sequencing of DArT libraries and marker data extraction

- **Dramatically simplified previous attempts at utilising padlocks in genotyping (MIPs)**

  - Eliminated some molecular "features", several steps and some expensive enzymes
  - Flexible sequence capture window (mostly 70-110 bp range)
  - Moved assay to 384 plate format reducing assay volume and therefore the cost while increasing throughput
  - Tested scalability beyond 10,000 markers in a single assay

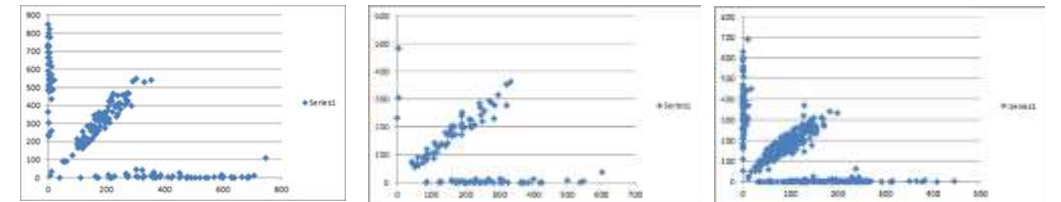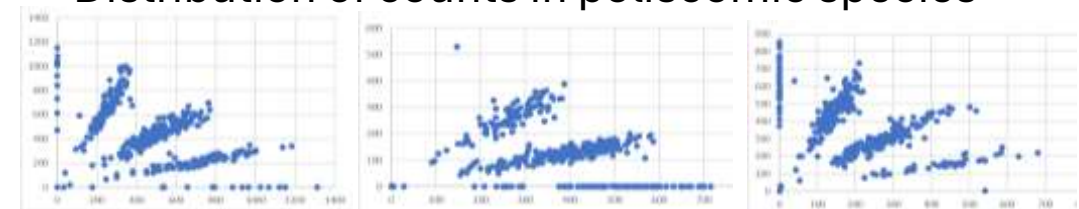**DArTag molecule final structure**

# Technical performance

- **Conversion rate ~97% even** in species with limited sequencing resources

- **Technology works well from a few markers to over 10,000**

- **Outperformed other technologies on the market, including in species with polysomic inheritence (potato, blueberry, alfalfa…)**

- **Typical call rate: >98%**

- **Calling reproducibility: >99.95%**

- Average marker read depth ~100 X for diploids and 200- 500 X for polysomic species

- Cost dependent on the number of markers and required sequencing depth: -----> application!

- Main use at the moment in Genomic Selection of crops and animals

- In Ecology mostly large volume monitoring/Close Kin Mark Recapture applications



DArTag Panel Sizes

Distribution of counts in discomic species



Distribution of counts in poliscomic species

# Panel design and $ considerations

- **Fully automated panel design process** when the submission file is formatted and filled properly

- **Detailed description of data and format requirements** downloadable from https://www.diversityarrays.com/services/targeted-genotyping/

- **Over 250 panels established** since technology launched in 2015

- **Median panel size: 577**

- **Average panel size: 1640**

- **Reference genome** and marker data quality **very important for design success**

- **DArTag outperforms other technologies in $ in medium-to-large scale applications**

- **Panel development cost depended on service volume (synthesis scale of oligos)**
  - **Cost per marker between $5-$15 for 20 K - 2 M assays**

- **Pricing strictly "per plate" as cost the same for full (94 samples) and partial plates**

- **Genotyping costs between $750/plate (small panels, very large service volume) to $3,000 (large panels, small volume)**
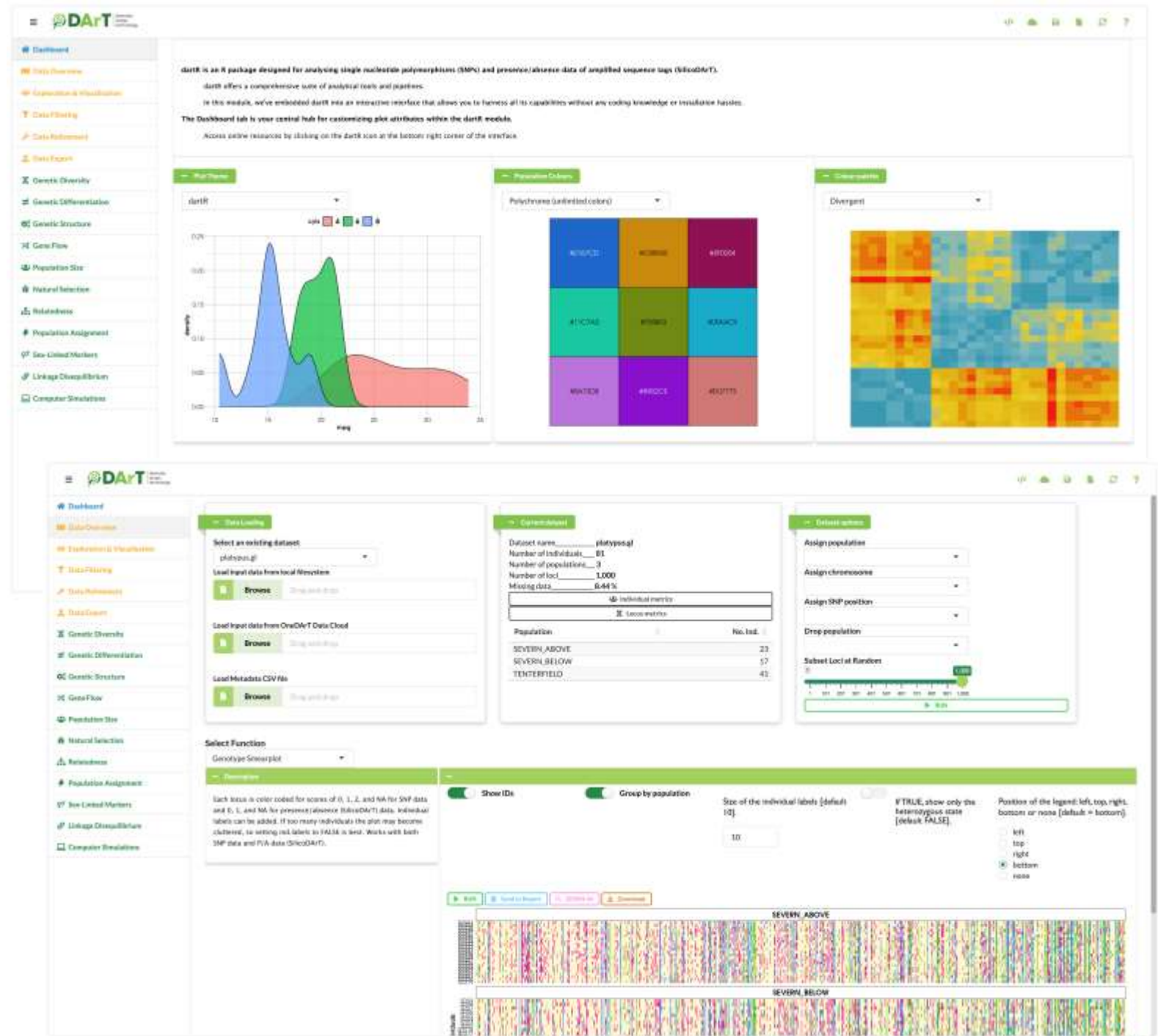


**Automated pipeline for panel design consisting of 6 plugins, with dedicated application (TagGen) for oligonucleotides (padlocks) developed @DArT**

# OneDArT release in July!
## dartR integration in Analytics+

- ✓ Partnership between DArT and Australian academics
- ✓ Extending user base of dartR
- ✓ Expanding from ecology-focused application to more general utility
- ✓ A broad range of analytical functions including complex modelling
- ✓ Accessible in OneDArT for people with no skill (or interest) in using R
- ✓ Providing expandable compute resources at low cost
- ✓ Bringing genomic and environmental data together with mobile app collected sample metadata

# Example - Redfin blue eye

SNP panel development required for species monitoring
from non-destructive, trace DNA samples

# Example - Redfin blue eye

- Decide on aim and number of SNPs in the panel

- Filter for quality of SNPs
- Filter for sequence quality and suitability
- Select and check panel

# Example - Redfin blue eye