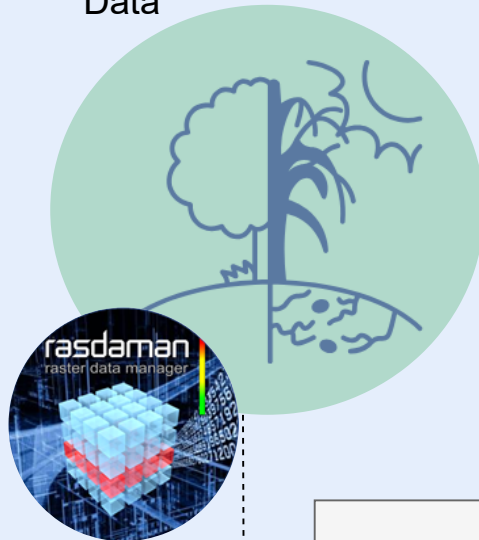
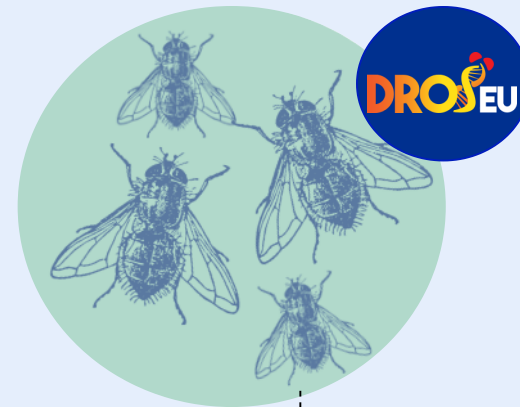


FairiCube UseCase3: Drosophila Genetics

Geospatial Data / Environmental Data



Genetic Data for *Drosophila melanogaster*



NHM
"Landscape Genomics Pipeline"

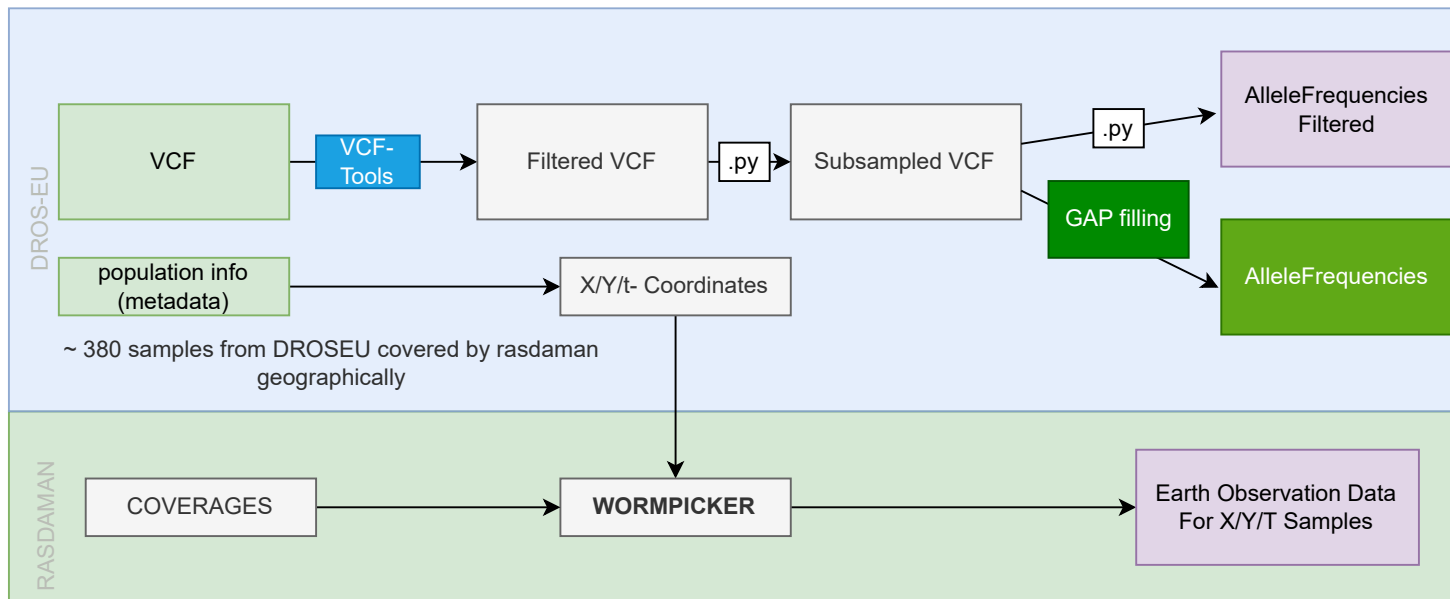
[uc3-drosophola-genetics/projects/LandscapeGenomicsPipeline](https://github.com/uc3-drosophola-genetics/projects/LandscapeGenomicsPipeline) at main · FAIRiCUBE/uc3-drosophola-genetics (github.com)

So far pipeline is adjusted to work with only a small subset due to computational efforts (10/737 potential samples, 10.000 Positions).
Shellscript / Python / R

Step 1: Obtain and Filter Data

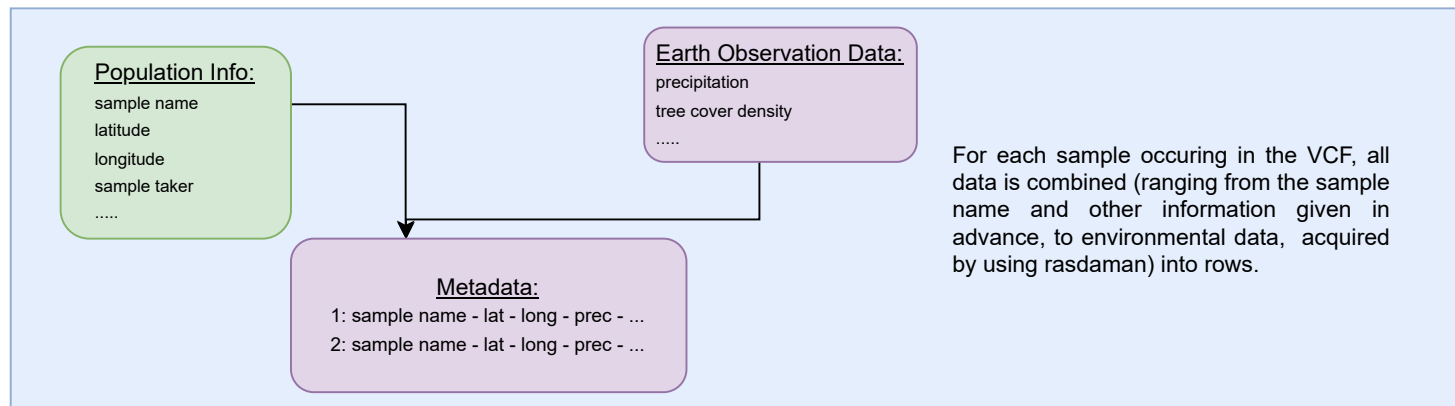
DrosEU: Download data as VCF, Subsample, convert to allele frequency table

Rasdaman: Install Rasdaman Software, get environmental data for spatial points covered in sampling



Step 2: Intersect / Combine Data

Match environmental data with metadata table (samplename, ID, population size,...)



Step 3: Apply Statistical Methods / Hypothesis Testing

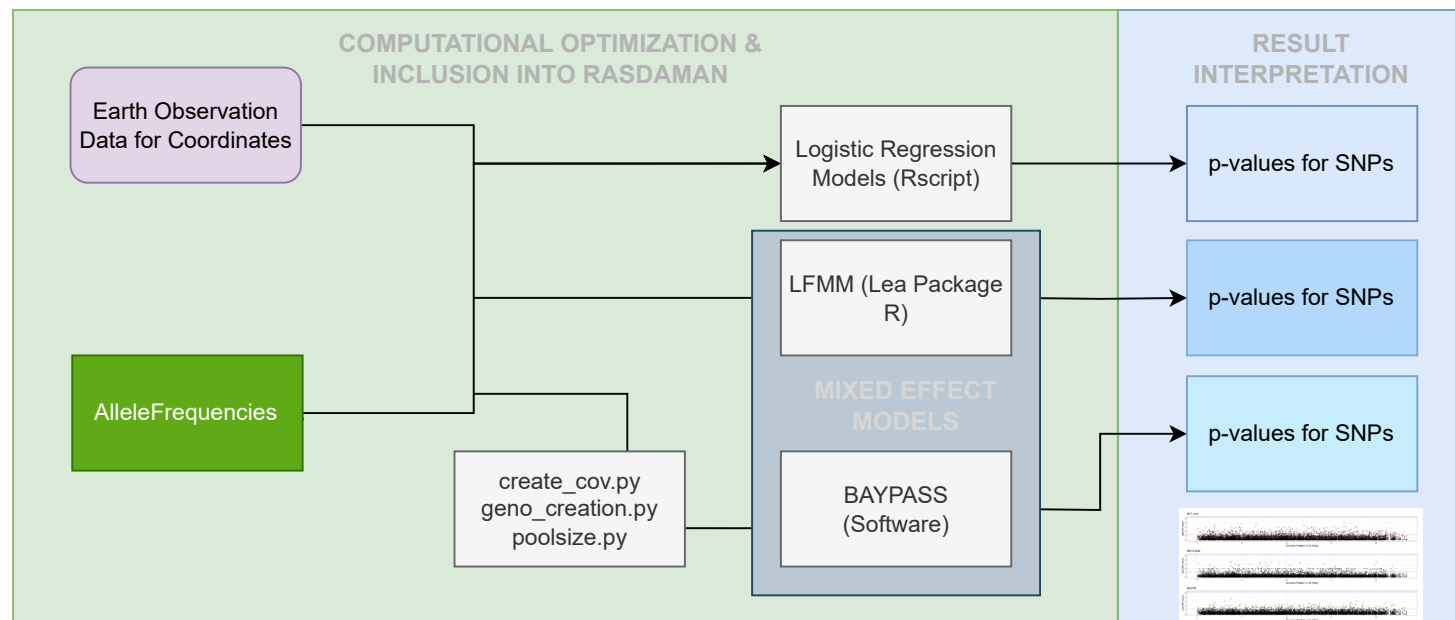
Test for correlation of polymorphism-frequency with variable:

- Linear Regression, Logistic Regression, Latent Factor mixed Models

Usage of specific programs written for (population)_genetic analyses:

- LEA, BAYPASS,...

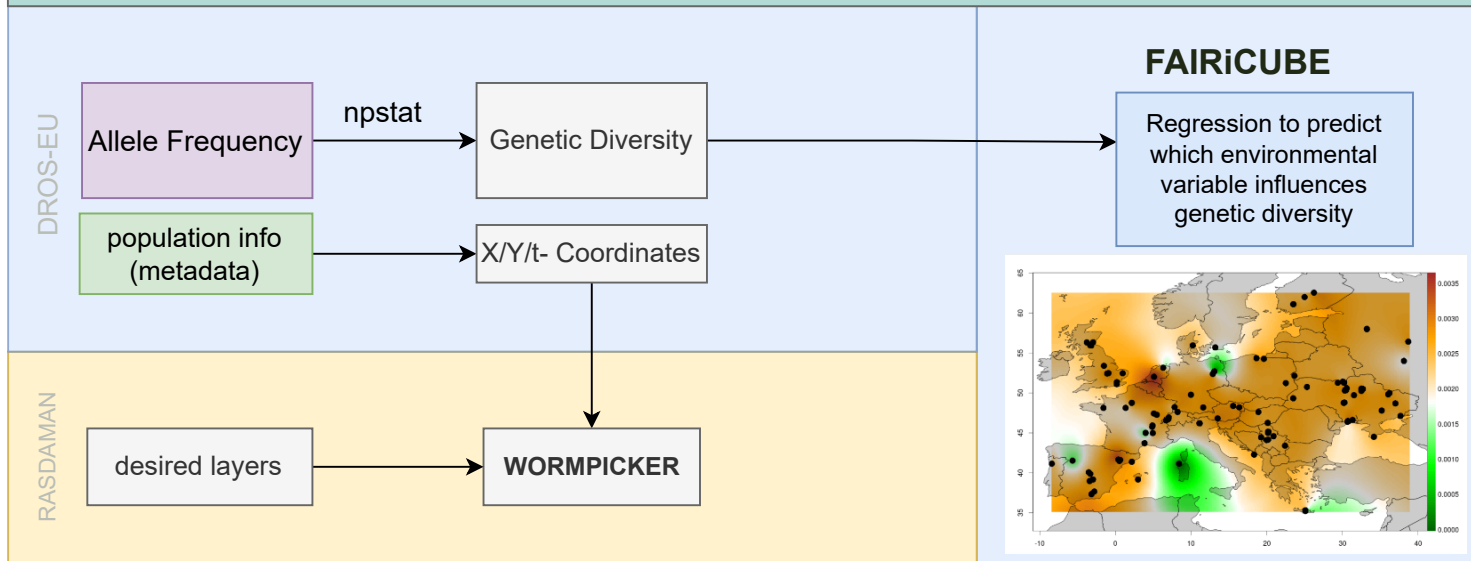
--> Computationally heavy for multiple 100 samples, pipeline so far not fully adjusted to computationally intense number of samples



LFMM: include a strategy to correctly estimate number of latent factors (strucutre)

BAYPASS: Calculation of p-values not clear yet, filter??? Include more models?
Way of comparing (not visually) the results?? Overlapping ?

SPATIAL DISTRIBUTION OF GENETIC VARIATION



THE INFLUENCE OF ENVIRONMENTAL VARIATION ON SPECIES DISTRIBUTION

