# Genome-wide association and prediction at the population level using Bayesian hierarchical models

Mathieu Gautier

UMR INRAE/CIRAD/IRD/SupAgro CBGP

23$^{rd}$ January 2025

# Introduction

## Forces driving the evolution of genetic diversity in populations

- Mutation : generates variability
- Drift : introduces stochasticity (Finite Population Size)
- Migration (gene flow)
- <u>Selection</u>

## Different Influences of the evolutionary forces

- Demographic Factors (genetic drift, gene flow) expected to be common to all loci
  $\implies$ Global (genomic) effect $\rightarrow$ correlation structure of pop. allele frequencies
- Selection (mutation and recombination) expected to vary across loci
  $\implies$ Local (genomic) effect

# Introduction

## General assumption

- Diversity (pop. allele freq.) at loci underlying (genetic) adaptation of populations co-vary with fitness-related traits (but see Lotterhos, 2022)
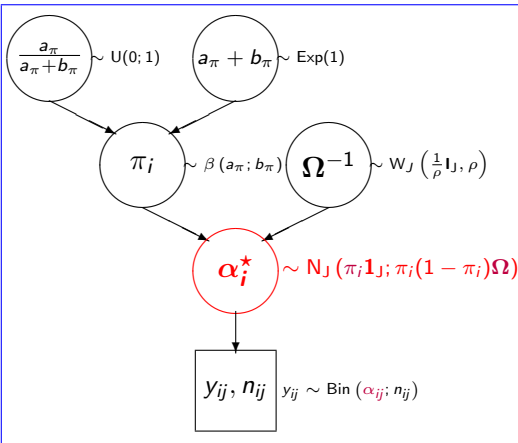
## Genome-wide association with population-specific covariables

- Modelling the relationship between genetic diversity and population covariables of interest across several (differentiated) populations may allow
  - uncovering the nature of adaptive traits and their genetic architecture
  - predicting covariate value from genomic information
- Different covariables of interest
  - Environmental (e.g., bioclimatic covariates, host plant, etc.) $\Rightarrow$ GEA
  - Phenotypic (e.g., mean height, mean weight, coat color) $\Rightarrow$ "pGWAS"

## Demographic history : a critical confounding factor

- Shared population history $\Rightarrow$ covariance structure of allele freq.

# The BAYPASS core model (Gautier, 2015)

$\frac{a_\pi}{a_\pi + b_\pi} \sim U(0; 1)$

$(a_\pi + b_\pi) \sim \text{Exp}(1)$

$\pi_i \sim \beta(a_\pi; b_\pi)$

$\Omega^{-1} \sim W_J\left(\frac{1}{\rho} I_J, \rho\right)$

$\alpha_i^\star \sim N_J\left(\pi_i 1_J; \pi_i(1-\pi_i)\Omega\right)$

$y_{ij}, n_{ij}$ $\quad y_{ij} \sim \text{Bin}(\alpha_{ij}; n_{ij})$

- **Multivariate Gaussian** prior on pop. (reference) allele frequencies (see Coop et al., 2010) of the I SNPs on J pops

- "instrumental" allele freq. $\alpha_{ij}^\star$ defined over the **real line support** :

$$\alpha_{ij} = \begin{cases} \alpha_{ij}^\star, & \text{if } \alpha_{ij}^\star \in (0,1), \\ 0, & \text{if } \alpha_{ij}^\star < 0 \,(\text{allele "lost"}), \\ 1, & \text{if } \alpha_{ij}^\star > 1 \,(\text{allele "fixed"}). \end{cases}$$

- $\pi_i$ might be interpreted as the "ancestral" ref. allele freq. of SNP $i$

- $\Omega = J \times J$ <u>scaled</u> covariance matrix of allele freq.

- $\Omega \Leftrightarrow$ "population relationship matrix" (captures the global effect of the demography)

- **Scaled allele frequencies** (i.e., corrected for pop. demographic history) :
  $\mathbf{X}_i = \{\widetilde{\alpha}_{ij}\}_{1..J} = \Gamma^{-1} \frac{\alpha_i^\star - \pi_i}{\sqrt{\pi_i(1-\pi_i)}}$ with $\Omega = \Gamma'\Gamma$ (Guenther & Coop, 2013 ; Olazcuaga et al., 2020)

Introduction
○○○

pGWAS/GEA
●○○

Genetic Offset
○○○

Genomic Prediction
○○○

Conclusions
○○

# BAYPASS models for association studies (GEA/pGWAS)

## General Principles

- Equivalent to a multivariate linear regression of the <u>scaled</u> allele frequencies $\widetilde{\alpha}_{ij}$ (SNP $i$; pop. $j$) on $K$ pop. covariate vectors $\boldsymbol{Z}_{\boldsymbol{k}}^{(k)} = \left\{ z_{jk} \right\}_{1..J}$ ($\Leftrightarrow$ "fixed" effect) :

$$\widetilde{\alpha}_{ij} = \sum_{k=1}^{K} \beta_{ik} z_{jk} + \epsilon_{ij} \text{ with } \epsilon_{ij} \sim \mathsf{N}(0, 1)$$

- Accounts for the confounding ($\Leftrightarrow$ "random") effect of shared population history by the modeling of $\widetilde{\alpha}_{ij}$ (instead of $\alpha_{ij}$)

- If $\widehat{\beta}_{ik} \neq 0$, SNP $i$ is deemed associated with the $k^{\text{th}}$ covariate

## In BAYPASS : 3 procedures to estimate the $\beta_i$'s and/or BF's

- From $\widetilde{\alpha}_{ij}$'s sampled under the core model with MCMC :
  - Importance Sampling <u>approximation</u> of the $\beta_i$'s and BF
  - "quick and dirty" and $\Leftrightarrow$ <u>univariate</u> regression on each covariable in turn

- MCMC sampling of the $\beta_i$'s $\Rightarrow$ accurately estimated but decision harder

- Penalized regression $\Rightarrow$ BF estimation (but some $\beta_i$'s shrinked towards 0)
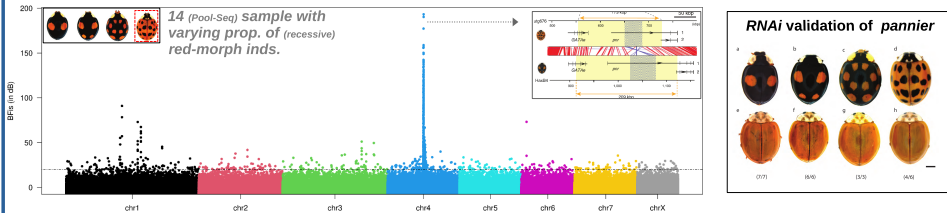
Introduction
○○○

pGWAS/GEA
○●○

Genetic Offset
○○○

Genomic Prediction
○○○

Conclusions
○○

# The "AUX" covariate model (i.e., with 'auxiliary variable')



- The binary variable $\delta_i$ specifies whether the SNP is associated ($\delta_i = 1$) or not ($\delta_i = 0$)

- Integrating over $P$ (prop. of associated SNPs) allows dealing with multiple testing issues

- From $\mathrm{P}[\delta_i = 1 | data]$ (a.k.a. PIP), $\mathbb{BF}_{mc} = \frac{\text{Post. odds}}{\text{Prior odds}} = \frac{\text{PIP}}{1 - \text{PIP}} \times \frac{1 - \mathbb{E}[P]}{\mathbb{E}[P]}$ (with $\mathbb{E}[P] = \frac{a_P}{a_P + b_P}$)

Introduction
○○○

pGWAS/GEA
○○●

Genetic Offset
○○○

Genomic Prediction
○○○

Conclusions
○○

# Example of applications



**A) pGWAS and color morphs in the ladybird beetle *H. axyridis*** (Gautier et al., 2018)

*14 (Pool-Seq) sample with varying prop. of (recessive) red-morph inds.*

*RNAi validation of pannier*

**B) GEA and climate adaptation in *A. thaliana*** (Frachon et al., 2018)

GEA with 6 non correlated env. Covariates *(e.g. Mean Annual Temperature)*

**168 (Pool-Seq) samples (micro-geographic scale)**

# GEA models : beyond the hunt for genes...

Simple (but efficient) modeling of the relationship (across populations) between adaptive genomic composition and the environment

- In GEA linear models (e.g., BAYPASS) : the $\beta$'s quantify the effect of (env.) covariates on the genetic diversity of adaptive variants

$$\widetilde{\alpha}_{ij} = \beta_i^{(1)} z_j^{(1)} + \ldots + \beta_i^{(K)} z_j^{(K)} + \epsilon_{ij}$$

- The ($n_{\text{snps}} \times n_{\text{cov}}$) matrix $\boldsymbol{B} = \{\beta_{ik}\}$ summarizes (linearly) the relationship between adaptive genetic diversity and environment (on a genome-wide basis)

Some assumptions to gain insights from $\boldsymbol{B}$ (Gain et al., 2023)

- Genotyped SNPs capture the whole-genome adaptive genetic diversity
- Sampled populations are representative of species diversity (for the geographical scale of interest) and locally adapted
- (some) covariables are (co)related to the (main) selective pressure
  $\boldsymbol{B}$ may then give insights into those driving adaptation (e.g., via s.v.d.)

# Evaluating population maladaptation to a new environment

## The (geometric) Genetic Offset (Gain et al., 2023)

- If $e_o$ (resp. $e^\star$) is the vector of the K covariable values (e.g., bioclim variables) for the original (resp. new) environment :

$$\text{GO} = \frac{1}{I} \left( e_o - e^\star \right)' B' B \left( e_o - e^\star \right) = \frac{1}{I} \sum_{i=1}^{I} \left( \tilde{e}_i - \tilde{e}_i^\star \right)^2$$

- $\tilde{e} = Be = \left\{ \sum_{k=1}^{K} \beta_{ik} e_k \right\}_i$ is the $n_{\text{snp}}$-length vector of global effect of environment on genetic div. at each SNP (NB : $\tilde{e}_i = 0$ if SNP $i$ is "neutral")

- GO $\Leftrightarrow$ (squared) euclidean environmental distance ("genetically") weighted by the env. effect on adaptive genetic diversity)

## Properties of (geometric) GO

- $GO \propto -\log \left( w \left( x, x^\star \right) \right)$ where $w \left( x, x^\star \right) < 1$ is the relative fitness value of traits at equilibrium in $e$ when placed in $e^\star$

- Supported by simulated and empirical data (e.g., Laruson et al., 2022, Gain et al., 2023)

# GO to predict population invasiveness (Camus et al., 2024)



**Simulation Study**

3 "native" environment grid (5 x 5 pop.) with **two environmental variables**, polygenic local adaptation during 3000 generations.

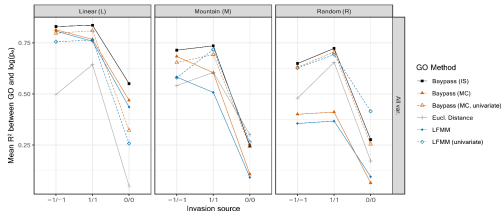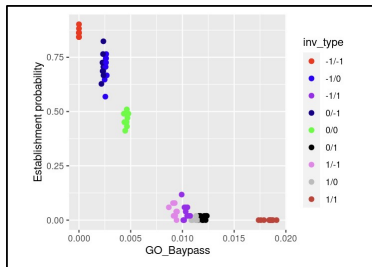**10 individuals** from a source population are randomly chosen to **invade a new environment**.

Each invasion is replicated 50 times under a non-WF model

$$PE = \frac{n_{establishment}}{n_{replicates}}$$

3 possible source populations :

9 invaded environments

**Predicting species invasiveness with genomic data: Is genomic offset related to establishment probability?** Evolutionary Applications

# Genomic prediction of population covariate

## Rationale

- Rely on GEA modeling of the relationship between genetic and covariate variation across populations to estimate population covariate values
  ⇒ pop-specific covariate is treated as a random variable

- Interpretation : pop. mean phenotype or tolerance range (e.g., for env. covariable)
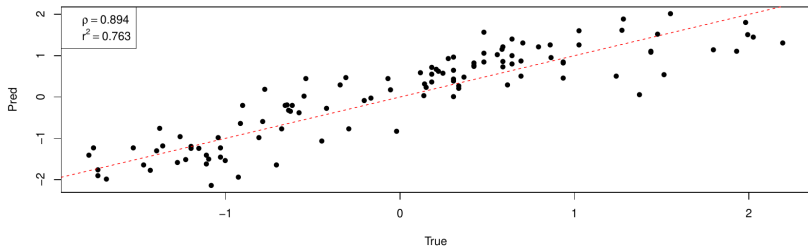
## Extending the BAYPASS model for genomic prediction

- Modeling uncertainty of the population covariate values

- full uncertainty ⇒ prediction

# The 'AUX' genomic prediction model (univariate case)

# Empirical evaluation : dog breeds weight (Gautier, in prep)

- Data (Hayward et al., 2016)
  - Genotypes : 155,609 SNPs genotyped on 111 dog breeds (n=6–636)
  - Phenotypes : mean male weight of each breed (American Kennel Club)

- 'Leave-one out' analysis (1 predicted pheno. vs. 110 known $\pm$0.01)

# Conclusions

## Linear models : not as trendy as AI but still useful !

- Flexible, robust (to non-linearity)
- Competitive esp. with limited number of pop. samples (bias-variance trade-off)

## Why bother with (old-school) Bayesian modeling as in BAYPASS ?

- Versatility makes it easy (but more computationally expensive) to account for
  - neutral structuring of genetic diversity (demographic history)
  - unbalanced designs, missing data, additional source of variation (e.g., Pool-Seq, pop. covariables)
  - combined data sets (Pool-Seq + Ind-Seq GL + count data in BAYPASS 3.0)
- Yet, urgent need to accelerate MCMC (subsampling, HMC)

## Predictive approaches are promising but still need

- Further evaluation on real (e.g., D. melanogaster) and simulated data (SLiM)
  - GO : robustness to genetic architecture, demographic history (e.g., admixture), genetic load, etc.
  - Genomic Prediction : sensitivity to the nb. of SNPs (LD), genetic architecture, etc.

- New developments esp. for (pop-level) genomic prediction :
  - BAYPASS : extend to categorical variable (e.g., fruit) ; multivariate GP
  - Comp. with other (machine/deep learning) approaches (e.g., Random Forest or CNN)

Introduction
ooo
pGWAS/GEA
ooo
Genetic Offset
ooo
Genomic Prediction
ooo
Conclusions
o•

# Acknowledgements



CBGP

Louise Camus

Simon Boitard

Maria Bogaerts-Marquez

Laure Olazcuaga

Nicolas Rode

Arnaud Estoup

Anne Loiseau

IBDML *(CNRS, Marseille)*

Benjamin Prud'homme

Junichi Yamaguchi

IBE *(Barcelona)*; DrosEU

Josefa Gonzalez

LIPME *(INRAE, Toulouse)*

Fabrice Roux

Léa Frachon

*... and you for your attention*