

# Sélection génomique, WGDs, Réseaux

Charles-Elie Rabier

ISEM, Institut des Sciences de l'Evolution de Montpellier

LIRMM, Laboratoire d'informatique, de Robotique et de Microélectronique



## Première partie : Sélection Génomique

# Prédiction en sélection génomique

Brigitte Mangin

LIPM, Laboratoire des Interactions Plantes et Microorganismes, Toulouse

Simona Grusea

Institut National des Sciences Appliquées de Toulouse



# Plan

- 1 Introduction
- 2 Formules pour la précision de la prédiction
- 3 Optimisation du panel
- 4 Amélioration de la prédiction

# Plan

- 1 Introduction
- 2 Formules pour la précision de la prédiction
- 3 Optimisation du panel
- 4 Amélioration de la prédiction

# Introduction

## Sélection Génomique

- Objectif = prédire les valeurs génétiques des candidats à la sélection
- Plus besoin de détecter les QTLs!!!

## Nouvelles technologies de séquençage

- Milliers de marqueurs disponibles  $\Rightarrow$  tous les QTLs fortement corrélés (fort Déséquilibre de liaison) avec au moins un marqueur

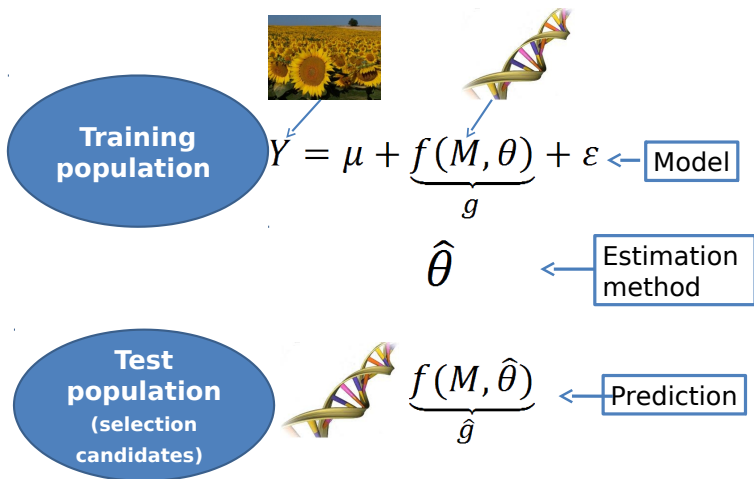
## Outils statistiques

- Tous les marqueurs analysés simultanément  $\Rightarrow$  Régression sur tout le génome

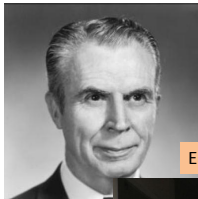
Problème de haute dimension

$$K \gg n$$

# Genomic Selection



# Cadres statistiques



**C.R. Henderson**  
Modèle mixte

Endelman



**T. Bayes**  
Modèle bayésien

Kärkkäinen and Sillanpää



**R. Tibshirani**  
Régression pénalisée

Li and Sillanpää



# Classement

- En général, le classement des méthodes est Bayes  $\geq$  Régressions pénalisées  $>$  Modèle mixte
- les méthodes ont moins d'influence que la densité de marquage, la taille de la population Training, l'héritabilité
- ou que la "distance" entre la population Training et la population TEST

Cet exposé : focus sur GBLUP, RRBLUP, Ridge (Pénalité L2)

# Modèle causal vs modèle de prédiction

Objectif : Prédiction en grande dimension

Modèle causal (Q QTLs)

$\theta^*$  effets QTLs,  $M^*$  matrice de génotypes aux QTLs,  
 $n$  individus d'entraînement

$$Y = M^* \theta^* + e$$

où  $Y = (Y_1, \dots, Y_n)'$ ,  $\theta^* = (\theta_1^*, \dots, \theta_Q^*)'$ ,  $e \sim N(0, \sigma_e^2 I_n)$

Modèle Bayésien de prédiction (K marqueurs, où  $K \gg n$ )

$\theta$  effets marqueurs,  $M$  matrice de génotypes aux marqueurs des "Training"

$$Y = M\theta + \varepsilon$$

où  $Y = (Y_1, \dots, Y_n)'$ ,  $\theta = (\theta_1, \dots, \theta_K)'$   $\sim N(0, \sigma_\theta^2 I_K)$ ,  $\varepsilon \sim N(0, \sigma_\varepsilon^2 I_n)$ ,  $\varepsilon_j \perp \theta_k$

Apprentissage  $\hat{\theta} = \mathbb{E}(\theta | Y) = M' (MM' + \lambda I_n)^{-1} Y = (M' M + \lambda I_K)^{-1} M' Y$   
 Régression Ridge (Pénalité L2) avec  $\lambda = \sigma_\varepsilon^2 / \sigma_\theta^2$

# Introduction d'un individu TEST

- Soit un individu TEST numéroté  $n + 1$

$$Y_{n+1} = M_{n+1}^* \theta^* + e_{n+1} \quad \text{où} \quad e_{n+1} \sim N(0, \sigma_e^2)$$

et  $M_{n+1}^*$  génotypes aux QTLs de l'individu  $n + 1$

- Prédiction de la valeur phénotypique

$$\begin{aligned} \hat{Y}_{n+1} = M_{n+1} \hat{\theta} &= M_{n+1} M' (M M' + \lambda I_n)^{-1} Y \\ &= M_{n+1} (M' M + \lambda I_K)^{-1} M' Y \end{aligned}$$

⇒ Critère d'accuracy

$$\rho = \frac{\text{Cov}(\hat{Y}_{n+1}, Y_{n+1})}{\sqrt{\text{Var}(\hat{Y}_{n+1}) \text{Var}(Y_{n+1})}} \quad \text{avec } M_{n+1} \text{ aléatoire et } M \text{ fixe}$$

# Plan

- 1 Introduction
- 2 Formules pour la précision de la prédiction
- 3 Optimisation du panel
- 4 Amélioration de la prédiction
- 5 Conclusion

# Formule générale pour l'accuracy

Théorème (R., Barre, ..., Mangin, Plos One 2016)

Conditionnellement à  $M$  et  $M^*$ ,

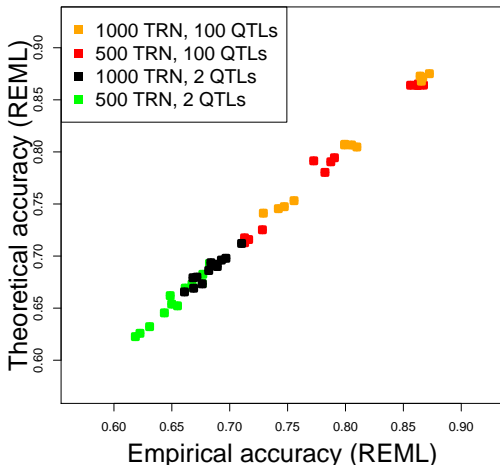
$$\rho = \frac{\theta^{*\prime} \mathbb{E} (M'_{n+1} M_{n+1}) M' V^{-1} M^* \theta^*}{\left\{ \sigma_e^2 \mathbb{E} \left( \|M_{n+1} M' V^{-1}\|^2 \right) + \theta^{*\prime} M^* V^{-1} M \text{Var} (M'_{n+1}) M' V^{-1} M^* \theta^* \right\}^{1/2} \Omega^{1/2}}$$

où  $V = MM' + \lambda I_n$  et  $\Omega = \text{Var} (M^*_{n+1} \theta^*) + \sigma_e^2$

**Facteurs** agissant sur l'accuracy :

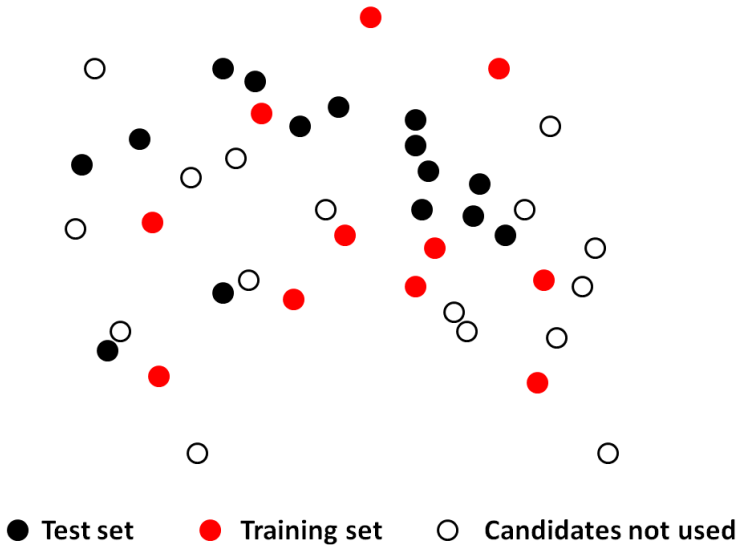
- Colonne  $q$  de  $M' V^{-1} M^*$  : **DL (utilisant la métrique V)** entre chaque marqueur et le QTL  $q$
- $\mathbb{E} \left( \|M_{n+1} M' V^{-1}\|^2 \right)$  : **similarité** entre TRN et TEST
- $\text{Var} (M'_{n+1})$  : matrice de covariance de taille  $K \times K$ , contenant l'ensemble des **DL classiques** du TEST

# Les QTLs sont situés sur quelques marqueurs (DL parfait)



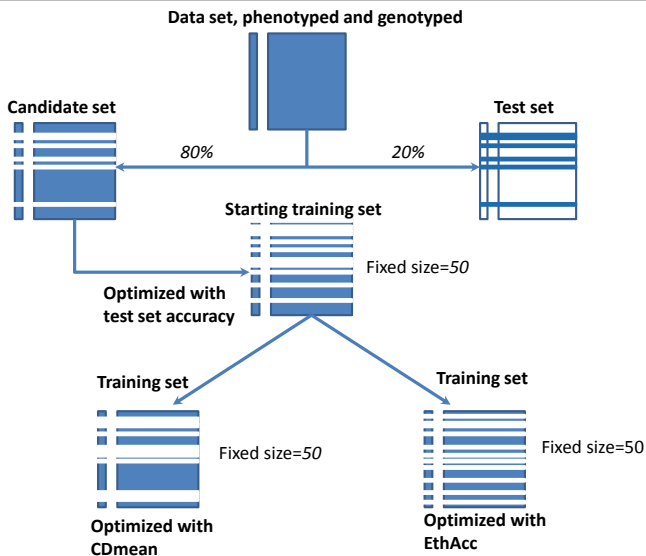
# Plan

- 1 Introduction
- 2 Formules pour la précision de la prédiction
- 3 **Optimisation du panel**
- 4 Amélioration de la prédiction
- 5 Conclusion

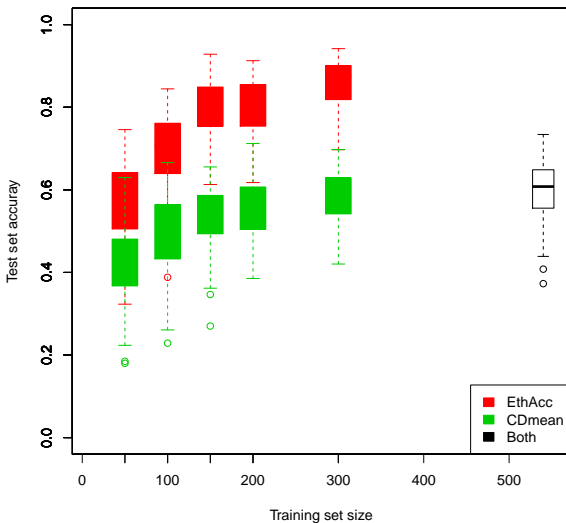


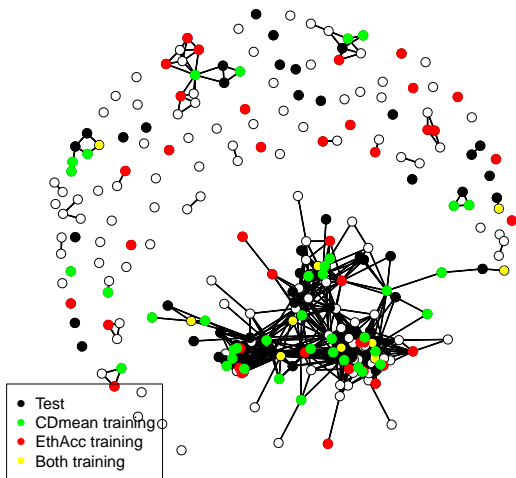


# Training set optimization : EthAcc versus CDmean (Rincant et al, Genetics 2012)



# Training set optimization (Sugar beet) : EthAcc versus CDmean





- CDmean optimization gave an accuracy of 0.07
- EthAcc optimization gave an accuracy of 0.76

# Plan

- 1 Introduction
- 2 Formules pour la précision de la prédiction
- 3 Optimisation du panel
- 4 Amélioration de la prédiction
- 5 Conclusion

# Pour aller plus loin ...

Théorème (R., Barre, ..., Mangin, Plos One 2016)

Conditionnellement à  $M$  et  $M^*$ ,

$$\rho = \frac{\theta^{*\prime} \mathbb{E} (M_{n+1}^{*\prime} M_{n+1}) M' V^{-1} M^* \theta^*}{\left\{ \sigma_e^2 \mathbb{E} \left( \|M_{n+1} M' V^{-1}\|^2 \right) + \theta^{*\prime} M^{*\prime} V^{-1} M \text{Var} (M'_{n+1}) M' V^{-1} M^* \theta^* \right\}^{1/2} \Omega^{1/2}}$$

où  $V = MM' + \lambda I_n$  et  $\Omega = \text{Var} (M_{n+1}^* \theta^*) + \sigma_e^2$ .

Décomposition SVD de  $M$

$$M = U D W'$$

où

- $D$  matrice diagonale de taille  $r \times r$ , de plein rang
- $U$  matrice de taille  $n \times r$ , telle que  $U' U = I_r$
- $W$  matrice de taille  $p \times r$ , telle que  $W' W = I_r$

# En injectant la décomposition SVD dans notre formule

Par abus de notation :

- $\theta^*$  vecteur sparse de dimension  $K$  contenant les effets QTLs

Théorème (R., Mangin, Grusea (en révision pour Scand J. Statistics))

$$\hat{\rho} \geq \frac{\|WW'\theta^*\|^2 \min \frac{d_s^4}{d_s^2 + \lambda}}{\sqrt{\sigma_e^2 r + \|WW'\theta^*\|^2 \max d_s^2} \sqrt{\|WW'\theta^*\|^2 \max d_s^2 + \sigma_e^2}}$$

$$\hat{\rho} = \frac{\sum_{s=1}^r \frac{d_s^4}{d_s^2 + \lambda} \|W^{(s)} W^{(s)'} \theta^*\|^2}{\left( \sigma_e^2 \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} + \sum_{s=1}^r \frac{d_s^6}{(d_s^2 + \lambda)^2} \|W^{(s)} W^{(s)'} \theta^*\|^2 \right)^{1/2} \left( \sum_{s=1}^r d_s^2 \|W^{(s)} W^{(s)'} \theta^*\|^2 + \sigma_e^2 \right)^{1/2}}$$

On aimerait avoir  $\sigma_e^2 \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2}$  petit, et  $\sum_{s=1}^r \frac{d_s^4}{d_s^2 + \lambda} \|W^{(s)} W^{(s)'} \theta^*\|^2$  grand

# Vers une amélioration du RRBLUP, GBLUP, Ridge

Idée : considérer un espace de dimension plus faible

Rappel :  $U = (U^{(1)}, \dots, U^{(r)})$  base orthonormale de l'espace engendré par les colonnes de  $M$ .

On choisit  $\tilde{r}$  colonnes de  $U$ .

Soit l'estimateur

$$\tilde{\theta} = M' V^{-1} \tilde{U} \tilde{U}' Y \quad \text{où} \quad \tilde{U} = (U^{\sigma(1)}, \dots, U^{\sigma(\tilde{r})})$$

où  $\tilde{U} \tilde{U}' Y$  est la projection de  $Y$  sur  $\text{Vect} \{U^{\sigma(1)}, \dots, U^{\sigma(\tilde{r})}\}$ .

⇒ Prédiction à l'aide de  $\tilde{\theta}$

# Accuracy basée sur ce nouvel estimateur

Théorème (R., Mangin, Grusea (en révision pour Scand J. Statistics))

$$\tilde{\rho} \geq \frac{\left\| \tilde{W} \tilde{W}' \theta^* \right\|^2 \min_{1 \leq s \leq \tilde{r}} \frac{d_{\sigma(s)}^4}{d_{\sigma(s)}^2 + \lambda}}{\sqrt{\sigma_e^2 \tilde{r} + \left\| \tilde{W} \tilde{W}' \theta^* \right\|^2 \max_{1 \leq s \leq \tilde{r}} d_{\sigma(s)}^2} \sqrt{\left\| \tilde{W} \tilde{W}' \theta^* \right\|^2 \max_{1 \leq s \leq \tilde{r}} d_{\sigma(s)}^2 + \sigma_e^2}}$$

$$\tilde{\rho} = \frac{\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{d_{\sigma(s)}^2 + \lambda} \left\| \tilde{W}^{(\sigma(s))} \tilde{W}^{(\sigma(s))'} \theta^* \right\|^2}{\left( \sigma_e^2 \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} + \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^6}{(d_{\sigma(s)}^2 + \lambda)^2} \left\| \tilde{W}^{(\sigma(s))} \tilde{W}^{(\sigma(s))'} \theta^* \right\|^2 \right)^{1/2} (\Omega)^{1/2}}$$

avec  $\Omega = \sum_{s=1}^r d_{\sigma(s)}^2 \left\| W^{(\sigma(s))} W^{(\sigma(s))'} \theta^* \right\|^2 + \sigma_e^2$



# Dans quelles conditions améliore-t-on l'accuracy ?

- **Estimateur Ridge**  $\hat{\theta}$  basé sur toutes les colonnes de  $U$ 
  - accuracy  $\hat{\rho}$ , prédiction  $\hat{Y}_{n+1}$
- **Nouvel estimateur**  $\tilde{\theta}$  basé sur  $\tilde{r}$  colonnes de  $U \Rightarrow \tilde{\beta}$ 
  - accuracy  $\tilde{\rho}$ , prédiction  $\tilde{Y}_{n+1}$
- **Complémentaire**  $\vec{\theta}$  de notre nouvel estimateur basé sur les  $r - \tilde{r}$  colonnes restantes de  $U$ 
  - accuracy  $\vec{\rho}$ , prédiction  $\vec{Y}_{n+1}$

$\tilde{\rho} \geq \hat{\rho}$  si et seulement si

$$\frac{\widehat{\text{Cov}}(\tilde{Y}_{n+1}, Y_{n+1})}{\widehat{\text{Cov}}(\vec{Y}_{n+1}, Y_{n+1})} \geq \frac{\widehat{\text{Var}}(\tilde{Y}_{n+1})}{\widehat{\text{Var}}(\vec{Y}_{n+1})} \left( 1 + \sqrt{1 + \frac{\widehat{\text{Var}}(\vec{Y}_{n+1})}{\widehat{\text{Var}}(\tilde{Y}_{n+1})}} \right).$$

# Un exemple où $\theta$ appartient à l'espace engendré par les lignes de $M$

- Chromosome de longueur 1M
- 1000 marqueurs
- $\theta = 0.3W^{(1)} + 0.3W^{(2)} + 0.3W^{(3)}$
- $\tilde{r}$  et les colonnes de  $U$  choisies par validation croisée
- 100 TESTS

$\sigma_e^2$	$n$	Méthode	200 générations	400 générations
1	500	$\hat{\rho}$	0.7550	0.6721
		$\hat{\tilde{\rho}}$	0.7810	0.7041
	800	$\hat{\rho}$	0.7487	0.7037
		$\hat{\tilde{\rho}}$	0.7728	0.7312
9	500	$\hat{\rho}$	0.3370	0.2623
		$\hat{\tilde{\rho}}$	0.3809	0.3079
	800	$\hat{\rho}$	0.3317	0.2904
		$\hat{\tilde{\rho}}$	0.3734	0.3330

# Remerciements

Simona Grusea (INSA),

Brigitte Mangin (LIPM)



Collaborateurs sur CROPDL

- Muriel Tavaud / Jacques David (SupAgro Montpellier)
- Gilles Charmet / Delphine Ly / François Balfourier (INRA Clermont Ferrand)
- Philippe Barre (INRA Lusignan)/ Torben Asp (Aarhus university, Danemark)



## Deuxième partie : Duplications Entières du Génome

# Probabilistic approaches for detecting and locating whole genome duplications

Charles-Elie Rabier, Tram Ta, Cécile Ané

UW - Madison

Departments of Statistics and of Botany



THE UNIVERSITY  
*of*  
**WISCONSIN**  
MADISON

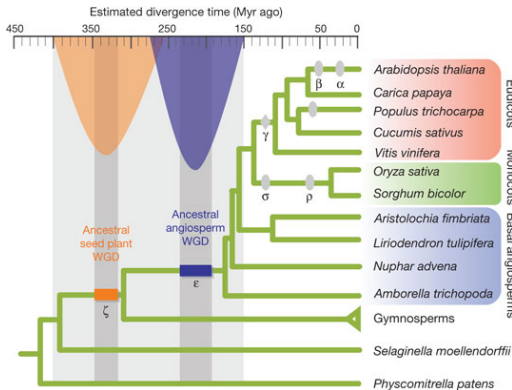
# Whole Genome Duplication (WGD)

“Ancestral polyploidy in seed plants and angiosperms”, Jiao et al. (Nature 2009)

“Whole-genome duplication followed by gene loss and diploidization has long been recognized as an important evolutionary force in animals, fungi and other organisms, especially plants”

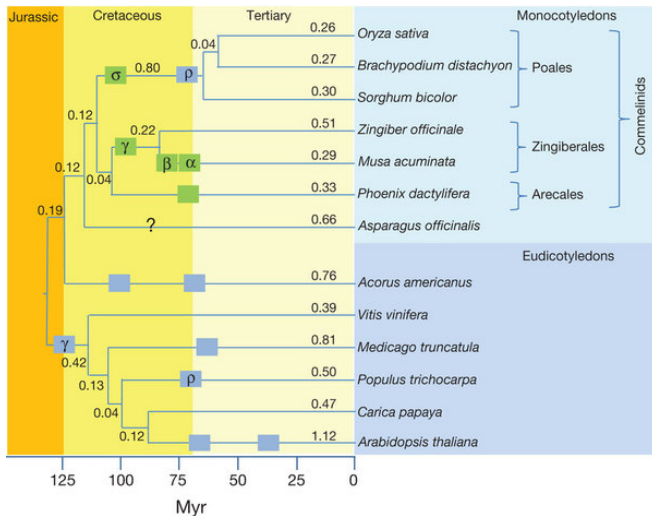
# WGD in seed plants and angiosperms

Jiao et al. (Nature 2009)



# WGD in bananas

D'Hont et al. (Nature 2012)





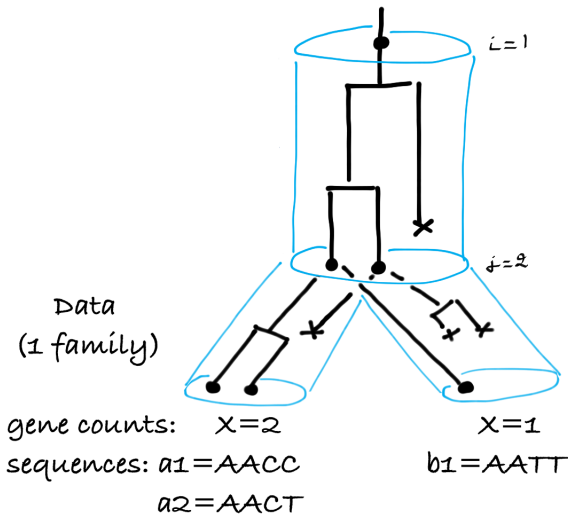
# Probabilistic model for gene family evolution

- phylogenetic framework : multiple species
- probabilistic model to avoid ad-hoc filtering of families or nodes
- requires : genes clustered into families. No synteny.

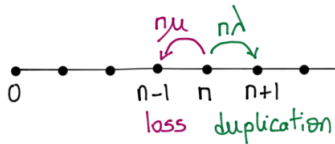
Birth-death model for small-scale events, and  
WGD model for large-scale events.

$$\text{likelihood} = \prod_{\text{families } f} \text{likelihood}(f)$$

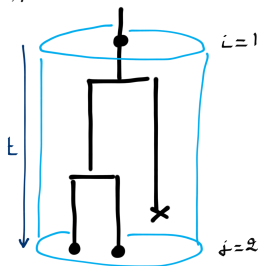
# Birth - death process for small scale events



Birth rate  $\lambda$ , death rate  $\mu$



## Likelihood of gene counts, birth - death process

 $\lambda, \mu$  : birth & death rates

$$p_t(i, j) = \mathbb{P}(X_t = j | X_0 = i)$$

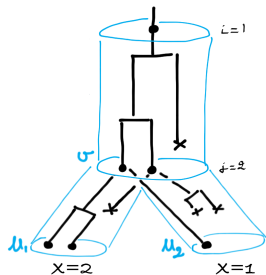
$$p_t(1, 0) = \gamma_t = \frac{\mu(e^{(\lambda-\mu)t} - 1)}{\lambda e^{(\lambda-\mu)t} - \mu},$$

$$p_t(1, 1) = (1 - \gamma_t)(1 - \psi_t) \text{ with } \psi_t = \frac{\lambda}{\mu} \gamma_t$$

$$p_t(i, j) = \sum_{k=0}^{i \wedge j} \binom{i}{k} \binom{i+j-k-1}{i-1} \gamma_t^{i-k} \psi_t^{j-k} (1 - \gamma_t - \psi_t)^k$$

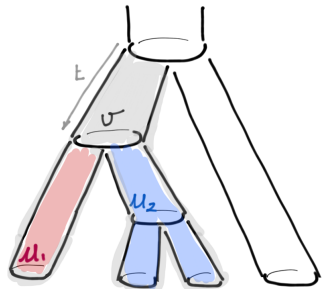
Bailey (1964)

## Likelihood of gene counts, birth - death process



Conditional likelihood  $L_v(i)$  at node  $v$  : probability of gene count data below  $v$  given  $X = i$  at parent of  $v$ , calculated recursively :

$$L_v(i) = \sum_j p_t(i, j) L_{u_1}(j) L_{u_2}(j)$$

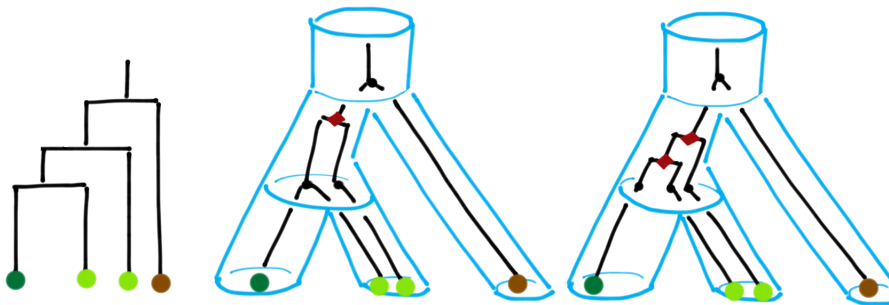


Geometric prior  $\pi$  for # at the root :

$$\text{likelihood} = \sum_j \pi(j) L_{u_1}(j) L_{u_2}(j)$$

or Csűrös & Miklós (2009)

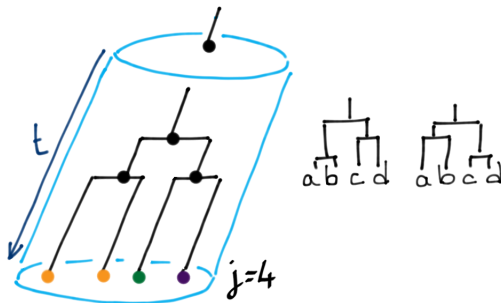
# Likelihood of gene tree reconciliations, BD process



Problem 1 : each gene tree has many "reconciliations" : to map gene tree inside species tree.

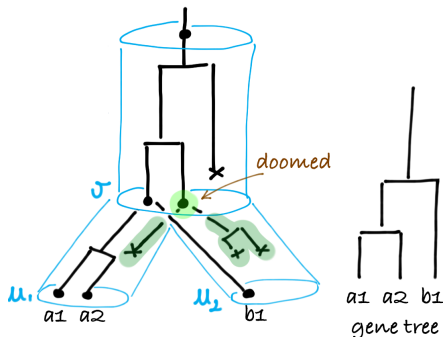
## Likelihood of gene tree reconciliations, BD process

Problem 2 : labels

For a reconciled subtree within a 'slice',  $j$  tips, 3 colors

Arvestad et al. (2009), Rasmussen &amp; Kellis (2011)

## Likelihood of gene trees reconciliations, BD process

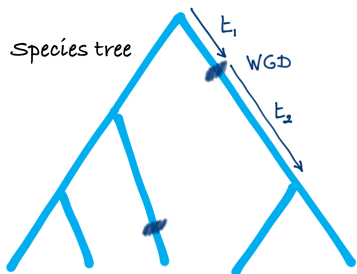


Problem 3 : gene trees lack **doomed** lineages

$d_v$  : probability that a lineage starting at node  $v$  leaves no descendent (or : is doomed). Recursively :

$$d_v = \left( \sum_j p_{t_1}(1, j) d_{u_1}^j \right) \left( \sum_j p_{t_2}(1, j) d_{u_2}^j \right)$$

# WGD model for large-scale events



At the WGD :

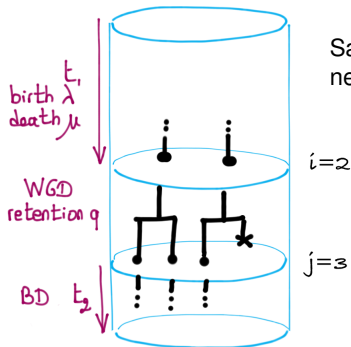
- each gene is duplicated
- second copy lost immediately with probability  $1 - q$ .

Each WGD has its own retention rate  $q$ , to explain :

- Large-scale events
- fragmentation : tendency to lose the extra copy, increased background loss rate shortly after WGD
- extension to whole genome triplications



## Likelihood : birth-death + WGD model



Same recursive algorithm through the tree, but new transition probabilities along WGD edges :

$$p_{\text{WGD}}(i, j) = \binom{i}{j-i} q^{j-i} (1-q)^{2i-j} \quad (i \leq j \leq 2i)$$

# Two methods to detect WGDs

Using **gene counts** only :

- **fast** (< 1s)
- exact likelihood
- optimize  $\lambda, \mu$  and separate  $q$ 's at each WGD
- but : **limited** information

R package `WGDgc`

# Two methods to detect WGDs

Using full [sequences](#) :

- **rich** information and model, but
- **slow** (e.g. 1h/family) : integrate over tree, reconciliation, branch lengths (gene-specific and species-specific rates).
- approximate likelihood
  - search over gene trees, but most parsimonious reconciliation.
  - new algorithm to find MP reconciliation with WGDs
- fixed  $\hat{\lambda}, \hat{\mu}$

C++ program `spimapWGD`, based on `SPIMAP` (Rasmussen & Kellis 2011)

# If you are interested in the gene tree ...

## Some notations

- $S$  : species tree
- $D$  : data (ie. alignment data)
- $T$  : gene tree topology
- $\ell$  : branch length
- $R$  : reconciliation

## Bayesian framework

- $\mathbb{P}(T, R|S)$  : topology prior
- $\mathbb{P}(\ell|T, R, S)$  : branch length prior
- $\mathbb{P}(T, R, \ell|D, S)$  : posterior

⇒ Markov Chain Monte Carlo (Hasting Metropolis) to estimate posterior distribution  $\mathbb{P}(T, R, \ell|D, S)$

# Approximate versus exact likelihood

## Exact Likelihood

$$\begin{aligned}\mathbb{P}(D|S) &= \sum_{T,R} \int_I \mathbb{P}(D, I, T, R|S) \\ &= \sum_{T,R} \int_I \mathbb{P}(D|I, T, S) \mathbb{P}(I|T, R, S) \mathbb{P}(T, R|S)\end{aligned}$$

## Approximate Likelihood

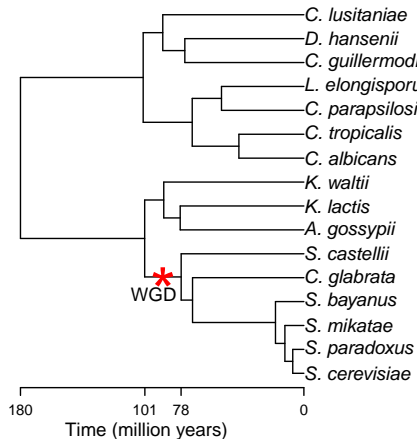
$$\mathbb{P}(D|S) \approx \mathbb{P}(D, \hat{\ell}, \hat{T}, \hat{R}|S)$$

with  $\hat{\ell}$ ,  $\hat{T}$ ,  $\hat{R}$  maximum a posteriori estimators of  $\ell$ ,  $T$ ,  $R$  given the data

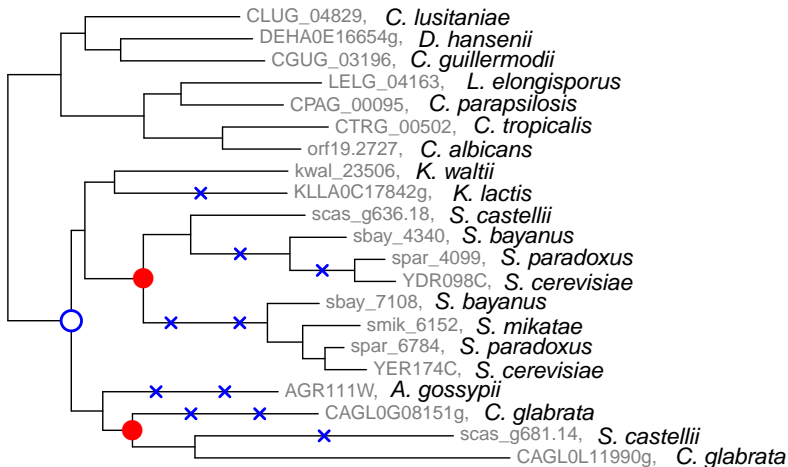
# Yeast genome evolution

Kellis et al. (2004), from synteny on *Kluyveromyces waltii* and *S. cerevisiae* :  
 "12% of the paralogous gene pairs were retained in each doubly conserved synteny block"

- 9209 gene families (Butler et al 2009)
- filter : 3932 families with  $\geq 1$  gene in both *Candida* and *Saccharomyces* subclades



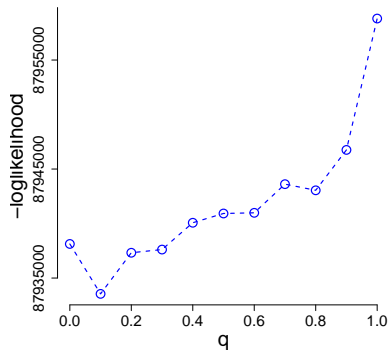
# A phylogenetic tree of gene family 1306



2 duplications at the WGD (red circles), 0 loss at the WGD  
 1 duplication, 10 losses (blue crosses)

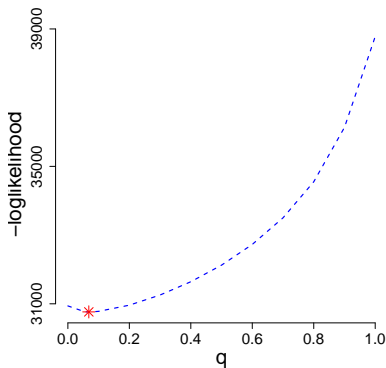
# Testing the Yeast WGD

from sequences



LRT : 9159.5

from gene counts



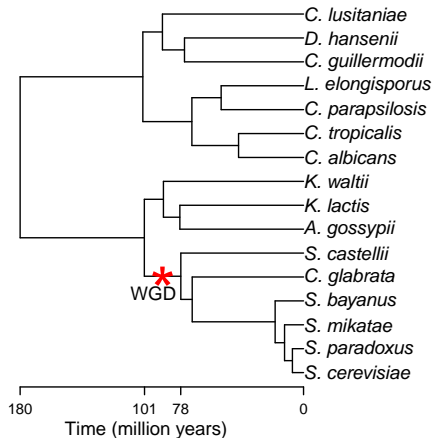
LRT : 348.1

retention rate :  $\hat{q} = 6.81\%$ , in  $[0.058, 0.079]$  with 95% confidence



# Yeast WGD timing

$\hat{t} = 0$  : immediately before speciation,  
 $\hat{t} \leq 5.04$  My with 95% confidence.



# Thanks to

Cécile Ané  
Tram Ta

Matt Rasmussen  
Bill Taylor



DEB-0949121



## Troisième partie : Réseaux phylogénétiques

# Apport des approches phylogénétiques pour expliquer l'origine des génomes mosaïques, exemple chez le Riz

Charles-Elie Rabier, Vincent Berry, Fabio Pardi et Céline Scornavacca

ISEM, Institut des Sciences de l'Evolution de Montpellier

LIRMM, Laboratoire d'informatique, de Robotique et de Microélectronique  
Genome Harvest

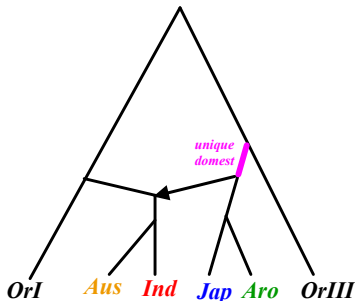
Jean-Christophe Glaszmann, Joao Santos

AGAP, Amélioration Génétique et Adaptation des Plantes, CIRAD



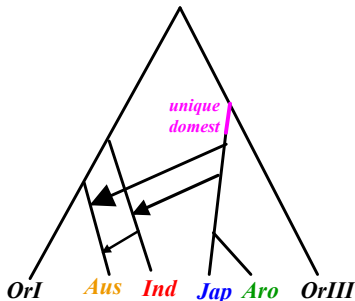
# Quelques thèses sur la domestication

- Huang et al. (Nature, 2012) : japonica domestiqué à partir d'un riz sauvage dans le sud de la Chine, puis croisé à un sauvage dans le sud est de l'Asie, générant indica



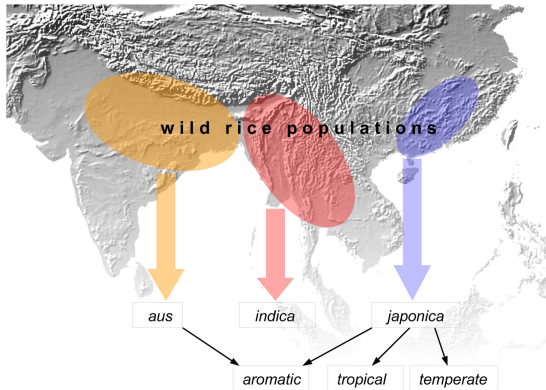
# Quelques thèses sur la domestication

- Choi et al. (MBE, 2017) soutiennent aussi **un seul évènement de domestication (japonica)**. Introgression par hybridation de japonica et proto-indica et proto-aus, générant indica et aus



# Quelques thèses sur la domestication

- Civan et al. (Nature Plants, 2015) : *indica*, *japonica* et *aus* domestiqués **séparément** dans différentes parties d'Asie



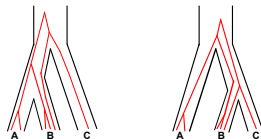
# Notre approche méthodologique

On s'intéresse à un modèle qui, outre le **tri de lignées**, considère explicitement les **mutations et hybridation**. Modélisation Bayésienne plus fine.

## Nos pistes :

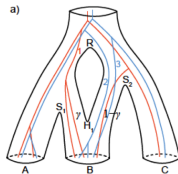
- 1 Inférence d'arbres d'espèces + arbres résumés en réseaux phylogénétiques

**SNAPP** (Bryant et al. 2012, MBE) + **SplitsTree**



- 2 Inférence directe de réseaux

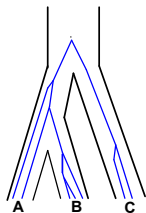
Extension de **SNAPP**  
aux réseaux





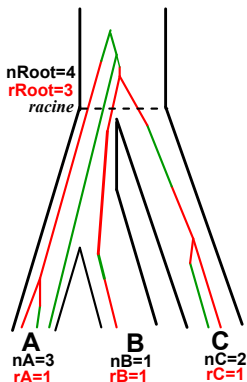
# Logiciel SNAPP pour l'inférence Bayésienne d'arbres (Bryant et al. 2012, MBE)

- Marqueurs bialléliques (SNPs) **indépendants** sachant l'arbre d'espèces
- Modélisation de l'arbre de locus (backward)
  - Processus de **coalescence** évoluant à l'intérieur d'un arbre d'espèces (**MultiSpecies Coalescent**)
  - Processus autorisant la **discordance** entre arbres de locus et arbres d'espèces (**tri de lignées incomplet**)



# Les mutations interviennent au cours du temps

- Modélisation des séquences (forward)
  - mutation (rouge  $\leftrightarrow$  vert) : modèle markovien évoluant le long des branches de l'arbre de locus
  - $u$  : taux de mutation rouge  $\rightarrow$  vert
  - $v$  : taux de mutation vert  $\rightarrow$  rouge



- V.a. :  $r_{Root}$ ,  $n_{Root}$ ,  $r_A$ ,  $r_B$ ,  $r_C$
- pas d'aléa dans  $n_A$ ,  $n_B$ ,  $n_C$
- $Data = (r_A, r_B, r_C)$
- Vraisemblance :  $\mathbb{P}(Data | S)$

# Calcul de vraisemblance dans un arbre (1)

$$\begin{aligned} & \mathbb{P}(\text{Data}) \\ &= \sum_i \sum_j \mathbb{P}(\text{Data} \mid \text{Count}, n_{\text{root}} = i, r_{\text{root}} = j) \mathbb{P}(n_{\text{root}} = i, r_{\text{root}} = j \mid \text{Count}) \\ &= \sum_i \sum_j \mathbb{P}(\text{Data} \mid \text{Count}, n_{\text{root}} = i, r_{\text{root}} = j) \mathbb{P}(r_{\text{root}} = j \mid n_{\text{root}} = i) \\ & \quad \times \mathbb{P}(n_{\text{root}} = i \mid \text{Count}) \end{aligned}$$

## Calcul de vraisemblance dans un arbre (2)

- $\mathbb{P}(n_{root} = i \mid Count)$  calculé récursivement en remontant dans le temps (postorder)

Tavaré (Theor Pop Biol, 1984), Watterson (Theor Pop Biol, 1984), Takahata and Nei (Genetics, 1985) ...

- $\mathbb{P}(Data \mid Count, n_{root} = i, r_{root} = j)$  calculé récursivement en remontant dans le temps (postorder)

Slatkin (Genetics, 1996) vs. Griffiths and Tavaré (Springer, 1997)

- $\mathbb{P}(r_{root} = j \mid n_{root} = i)$  calculé par
  - la loi Binomiale :  $\mathbb{P}(r_{root} = j \mid n_{root} = i) = C_i^j p^j (1-p)^{i-j}$
  - la loi  $\beta(\theta, \theta)$  sur le paramètre  $p$  de la Binomiale :  

$$\mathbb{P}(r_{root} = j \mid n_{root} = i) = C_i^j B(j + \theta, i - j + \theta) / B(\theta, \theta)$$
- Astuces afin de raccourcir les calculs : **Vraisemblances partielles...**

# La statistique Bayésienne dans SNAPP

- $S$  : arbre d'espèces (topologie, longueurs de branches, tailles de populations)
- $X_i$  : alignements pour le locus  $i$
- $G_i$  : arbre de locus pour le locus  $i$
- $m$  loci

$$\begin{aligned}\mathbb{P}(S|X_1, \dots, X_m) &\propto \left( \prod_{i=1}^m \int_{\psi} \mathbb{P}(X_i|G_i)\mathbb{P}(G_i|S)dG_i \right) P(S) \\ &\propto \mathbb{P}(\text{Data} | S) P(S)\end{aligned}$$

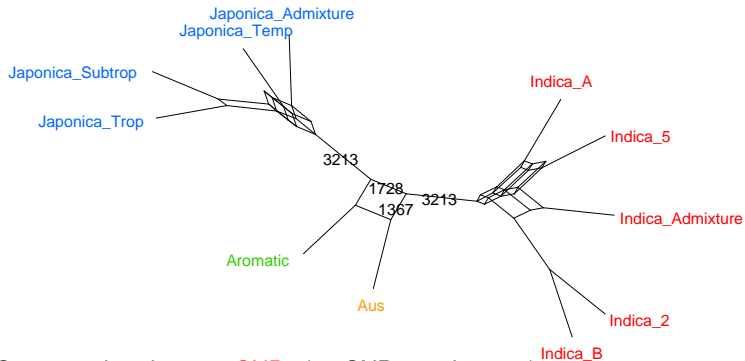
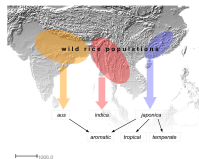
SNAPP intègre sur tous les arbres de locus

Calcul de la *prior*  $P(S)$  par le processus de **naissances**

⇒ **Markov Chain Monte Carlo** (MCMC) afin d'estimer la distribution à posteriori de  $\mathbb{P}(S|X_1, \dots, X_m)$

Implémenté dans **BEAST**

## Chromosome 6 (données J. Santos, J-C. Glaszmann)



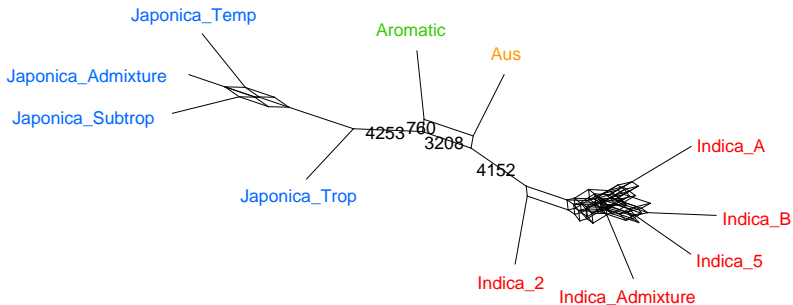
Conservation de **1550 SNPs** (un SNP tous les 500)

# Chromosome 10 (données J. Santos, J-C. Glaszmann)

Conservation de **1089 SNPs** (un SNP tous les 500)

- **JDD2** (1er SNP= 50ème SNP du chromosome 10)

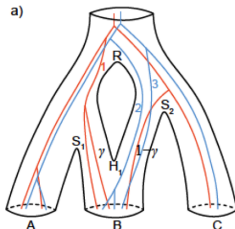
10000



# Simulateur basé sur un réseau (Genome Harvest)

SNAPPSimNet construit sur la base du simulateur SNAPPSim de Bryant et al. (2012)

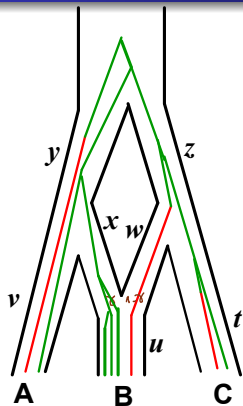
- Génération d'arbres de locus évoluant à l'intérieur d'un réseau selon un processus de coalescence



- Snapp est fortement attiré par un scénario sous-jacent au réseau



# Piste 2 : une méthode Bayésienne directe d'inférence de réseaux



$Data_z$  : proportion de rouge/vert dans les espèces sous la branche  $z$

$Data_y$  : proportion de rouge/vert dans les espèces sous la branche  $y$

$Data_{zT}$  et  $Data_{yT}$  ne sont plus indépendantes...

$Data_{zT}$  et  $Data_{yT}$  comprennent les allèles rouges et verts de l'espèce hybride B!!!

On ne peut plus effectuer le produit des probabilités

$\mathbb{P}(Data_{zT}) \times \mathbb{P}(Data_{yT})$  !!!

# Calcul de la vraisemblance dans un réseau

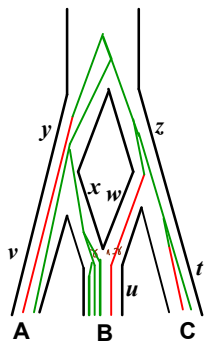
Pour raccourcir, *Count* est implicite dans les probabilités ...

$$\begin{aligned}
 & \mathbb{P}(\text{Data}) \\
 = & \sum_i \sum_j \mathbb{P}(\text{Data} \mid n_{\text{root}} = i, r_{\text{root}} = j) \mathbb{P}(r_{\text{root}} = j \mid n_{\text{root}} = i) \\
 & \mathbb{P}(n_{\text{root}} = i) \\
 = & \sum_i \sum_j \sum_{i'} \sum_{j'} \mathbb{P}(\text{Data}_{zT} \text{Data}_{yT} \mid n_{yT} = i', n_{zT} = i - i', r_{yT} = j', \\
 & r_{zT} = j - j') \mathbb{P}(r_{yT} = j', r_{zT} = j - j' \mid n_{yT} = i', n_{zT} = i - i', r_{\text{root}} = j) \\
 & \mathbb{P}(n_{yT} = i', n_{zT} = i - i' \mid n_{\text{root}} = i) \mathbb{P}(r_{\text{root}} = j \mid n_{\text{root}} = i) \mathbb{P}(n_{\text{root}} = i)
 \end{aligned}$$

On ne peut plus effectuer le produit des probabilités

$\mathbb{P}(\text{Data}_{zT}) \times \mathbb{P}(\text{Data}_{yT}) !!!$

# Notre algorithme : calcul des lois jointes



Quantités calculées successivement

- (1)  $\mathbb{P}(\text{Data}_{uT} \mid n_{uT}, r_{uT})$
- (2)  $\mathbb{P}(\text{Data}_{xB} \text{Data}_{wB} \mid n_{xB}, r_{xB}, n_{wB}, r_{wB})$
- (3)  $\mathbb{P}(\text{Data}_{xT} \text{Data}_{wB} \mid n_{xT}, r_{xT}, n_{wB}, r_{wB})$
- (4)  $\mathbb{P}(\text{Data}_{xT} \text{Data}_{wT} \mid n_{xT}, r_{xT}, n_{wT}, r_{wT})$
- (5)  $\mathbb{P}(\text{Data}_{vT} \mid n_{vT}, r_{vT})$
- (6)  $\mathbb{P}(\text{Data}_{yB} \text{Data}_{wT} \mid n_{yB}, r_{yB}, n_{wT}, r_{wT})$
- (7)  $\mathbb{P}(\text{Data}_{tT} \mid n_{tT}, r_{tT})$
- (8)  $\mathbb{P}(\text{Data}_{yB} \text{Data}_{zB} \mid n_{yB}, r_{yB}, n_{zB}, r_{zB})$
- (9)  $\mathbb{P}(\text{Data}_{yT} \text{Data}_{zB} \mid n_{yT}, r_{yT}, n_{zB}, r_{zB})$
- (10)  $\mathbb{P}(\text{Data}_{yT} \text{Data}_{zT} \mid n_{yT}, r_{yT}, n_{zT}, r_{zT})$
- (11)  $\mathbb{P}(\text{Data} \mid n_{root}, r_{root})$

- Optimisation combinatoire des calculs en cours (dimension des matrices)

# Conclusion

- Implémentation de la méthode Bayésienne pour les réseaux (en cours)
- L'inférence de réseau est un **sujet compétitif** : **Tanja Stadler** (ETH Zurich), **Luay Nakhleh** (Rice University, USA). Notre approche devrait être plus performante sur des réseaux aux nombreuses hybridations
- Jusqu'alors travail sur le **riz** → **autre plante d'intérêt**? Genome Harvest : **Banane, Citrus, Caféier, Riz, Tomate, Canne à sucre ...**
- Afin de comprendre l'histoire des **riz cultivés**, nécessité de disposer de **riz sauvages**, à l'instar de Choi et al. (MBE, 2017), Wang et al (Genome Research, 2017) ...
- SNAPP disponible sur <http://snapp.otago.ac.nz>







# Descriptif des simulations

- Population panmictique évoluant pendant 30, 50 or 70 generations
- Taille de la population : 100 TESTS +
  - $n = 500$  Trainings
  - $n = 1000$  Trainings
- Espèce haploïde
- Recombination modélisée selon Haldane (processus de Poisson d'intensité 1)
- Longueur du génome  $L = 1M$
- Marqueurs équidistants
- Densité de marqueurs : 100 SNPs, 1000 SNPs, 5000 SNPs ou 10000 SNPs
- Nombre de QTLs
  - 2 QTLs à 3cM et 80cM avec effet +1 et -2
  - 100 QTLs équidistants avec effet +0.15



# Un software intégrant différentes formules pour l'accuracy en selection génomique

**ShinyGPAS** par Morota (Université de Nebraska Lincoln)  
disponible sur <https://chikudaisei.shinyapps.io/shinygpas/>

Formules implémentées :

- Daetwyler et al. (Plos One 2008, Genetics 2010)
- Goddard et al. (Genetica 2009, Journal Of Animal Breeding And Genetics 2011)
- Nous :) (Plos One, 2016)
- de los Campos et al. (Plos Genetics, 2013)

$$\rho \leq h\sqrt{\{1 - (1 - b)^2\} h^2}$$

$b$  : coeff de régression entre relations génomiques aux marqueurs et relations génomiques aux génomes causaux

- Karaman et al. (Plos One, 2016),  $\tilde{h}^2$  proportion de variance expliquée par les marqueurs

$$\rho = \tilde{h}\sqrt{\frac{n\tilde{h}^2}{n\tilde{h}^2 + M_e}}$$

# De nouvelles proxies utilisant les effets QTLs

En **injectant** l'estimation des **causaux** dans notre formule ( **LASSO**, **Adaptative LASSO**, **Group LASSO**) ...

$$\hat{\rho} = \frac{\hat{\theta}^{*'} \mathbb{E} (M_{n+1}' M_{n+1}) M' V^{-1} \hat{M}^* \hat{\theta}^*}{\left\{ \sigma_e^2 \mathbb{E} \left( \|M_{n+1} M' V^{-1}\|^2 \right) + \hat{\theta}^{*'} \hat{M}^{*'} V^{-1} M \text{Var} (M'_{n+1}) M' V^{-1} \hat{M}^* \hat{\theta}^* \right\}^{1/2} \Omega^{1/2}}$$

où  $V = MM' + \lambda I_n$  et  $\Omega = \text{Var} (M_{n+1}' \theta^*) + \sigma_e^2$

$\mathbb{E} (M_{n+1}' M_{n+1})$ ,  $\mathbb{E} \left( \|M_{n+1} M' V^{-1}\|^2 \right)$  et  $\text{Var} (M'_{n+1})$  peuvent être estimés à l'aide :

- **uniquement** des Trainings → estimation de l'accuracy **avant** le génotypage des TESTS ...  $\hat{\rho}_{before}$
- **à la fois** des Trainings et TESTS → estimation de l'accuracy **après** le génotypage des TESTS ...  $\hat{\rho}_{after}$

# Régressions pénalisées

LASSO (Tibshirani, JRSSB 1996)

$$\hat{\theta}_{LASSO} = \operatorname{argmin} \|Y - X\theta\|^2 + \lambda \sum_{k=1}^K |\theta_k|$$

Adaptative LASSO (Zou, JASA 2006)

$$\hat{\theta}_{ADLASSO} = \operatorname{argmin} \|Y - X\theta\|^2 + \lambda \sum_{k=1}^K w_k |\theta_k|$$

Group LASSO (Yuan and Lin, JRSSB 2006)

sparsité par groupes,  $L$  nombre de groupes  
 $n_\ell$  nombre de marqueurs dans le groupe  $\ell$

$$\hat{\theta}_{GPLASSO} = \operatorname{argmin} \left\| Y - \sum_{\ell=1}^L X_\ell \vec{\theta}_\ell \right\|^2 + \lambda \sum_{\ell=1}^L \sqrt{n_\ell} \left\| \vec{\theta}_\ell \right\|^2$$

# Accuracy moyenne (100 simulations)

- LD parfait
- 100 QTLs avec le même effet +0.15
- 2 QTLs à 3cM et 80cM avec effets +1 and -2
- Training et TESTS basés sur 50 générations
- 1000 marqueurs, 500 Trainings, 100 TESTS

Méthode	100 QTLs	2 QTLs
Emp. Acc.	0.8143 (0.0036)	0.6683 (0.0053)
$\hat{\rho}_{before}(\theta^*)$	0.8086 (0.0001)	0.6597 (0.0001)
$\hat{\rho}_{before}(\hat{\theta}_{LASSO}^*)$	0.7635 (0.0017)	0.5354 (0.0031)
$\hat{\rho}_{before}(\hat{\theta}_{ADLASSO}^*)$	0.7627 (0.0012)	0.6488 (0.0017)
$\hat{\rho}_{before}(\hat{\theta}_{GPLASSO}^*)$	0.7581 (0.0014)	0.5471 (0.0029)
$\hat{\rho}_{after}(\theta^*)$	0.8045 (0.0019)	0.6576 (0.0021)
$\hat{\rho}_{after}(\hat{\theta}_{LASSO}^*)$	0.7502 (0.0026)	0.5347 (0.0037)
$\hat{\rho}_{after}(\hat{\theta}_{ADLASSO}^*)$	0.7489 (0.0023)	0.6454 (0.0027)
$\hat{\rho}_{after}(\hat{\theta}_{GPLASSO}^*)$	0.7479 (0.0024)	0.5495 (0.0034)

# Predictive ability : $\rho$

There are several predictive ability estimators

- assuming the mixed model is correct
  - based on coefficient of determination (CD) :

$$\hat{\rho}^{\text{CD}} = \frac{h}{n_{\text{test}}} \sum_i \sqrt{\text{CD}(g_{\text{test},i})}$$

- based on prediction error variance (PEV) :

$$\hat{\rho}^{\text{PEV}} = \frac{h}{n_{\text{test}}} \sum_i \sqrt{1 - \frac{\text{PEV}(g_{\text{test},i})}{\sigma_g^2}}$$

- assuming a fixed linear QTL model is known and correct and using the mixed model as "instrumental" model : the theoretical accuracy Rabier et al., PlosOne, 2016

# Thank you for your attention.

A special thank to Fanny Bonnafous, Prune Pegot-Espagnet, Charles-Elie Rabier, Ellen Goudemand, Renaud Rincent



# Equipe de Tanja Stadler (ETH Zurich)

## Statistique Bayésienne dans le cadre d'un réseau

- $N$  : réseau phylogénétique
- $X_i$  : alignements pour le SNP  $i$
- $G_i$  : arbre de gènes pour le SNP  $i$
- $m$  SNPs

$$\mathbb{P}(N, G_1, \dots, G_m | X_1, \dots, X_m) \propto \left( \prod_{i=1}^m \mathbb{P}(X_i | G_i) \mathbb{P}(G_i | N) \right) \mathbb{P}(N)$$

Calcul de l'a priori  $\mathbb{P}(N)$  par un processus de **naissance/hybridation**

⇒ **Markov Chain Monte Carlo** afin d'estimer la distribution à posteriori de  $\mathbb{P}(N, G_1, \dots, G_m | X_1, \dots, X_m)$ .

**Ils n'intègrent pas sur tous les arbres de gènes**

Zhang et al (MBE, Décembre 2017)

# Equipe de Luay Nakhleh (Université de Rice, USA)

- Wen et al. (Plos Genetics, 2016)
  - Arbres de gènes inférés lors d'un étape préliminaire
  - **Données = arbres de gènes !!!**

$$\mathbb{P}(N|G_1, \dots, G_m) \propto \left( \prod_{i=1}^m \mathbb{P}(G_i|N) \right) \mathbb{P}(N)$$

- Zhu et al. (Plos Computational Biology, Janvier 2018)
  - Données = alignements
  - **Intégration sur tous les arbres de gènes ...**
  - Algorithme inspiré de SNAPP mais très couteux en temps de calcul

$$\mathbb{P}(N|X_1, \dots, X_m) \propto \left( \prod_{i=1}^m \int_{\psi} \mathbb{P}(X_i|G_i) \mathbb{P}(G_i|S) dG_i \right) \mathbb{P}(N)$$