

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

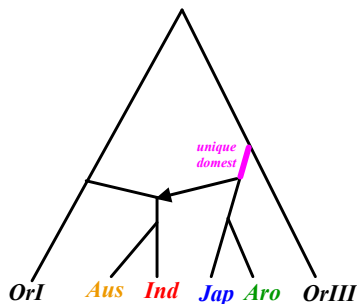
LIRMM, Laboratoire d'informatique, de Robotique et de Microélectronique
Genome Harvest

Jean-Christophe Glaszmann, Joao Santos

AGAP, Amélioration Génétique et Adaptation des Plantes, CIRAD

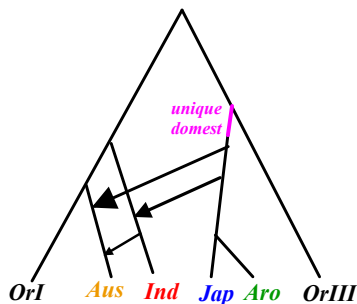
Quelques thèses sur la domestication

- 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



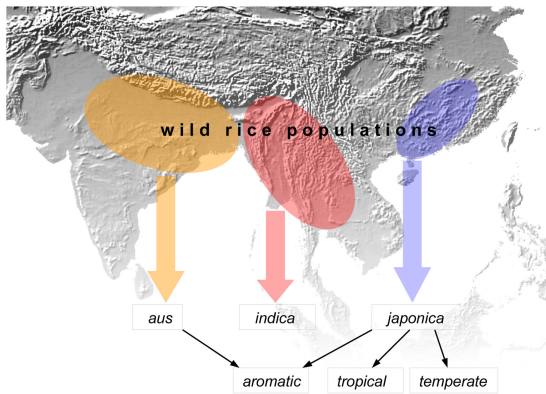
Quelques thèses sur la domestication

- 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



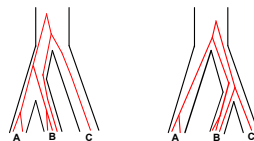
Notre approche méthodologique

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

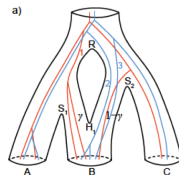
SNAPP (Bryant et al. 2012, MBE) + **SplitsTree**



a)

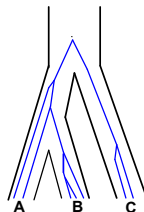
- 2 Inférence directe de réseaux

Extension de **SNAPP**
aux réseaux



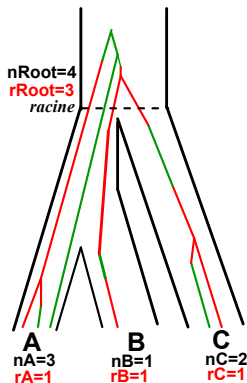
Logiciel SNAPP pour l'inférence Bayésienne d'arbres (Bryant et al. 2012, MBE)

- 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



- 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

- $\mathbb{P}(n_{root} = i \mid Count)$ 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

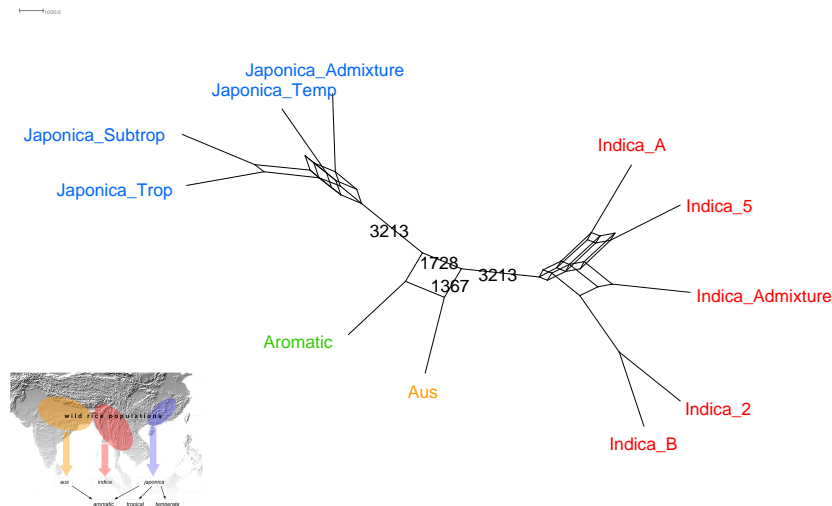
Calcul de la *prior* $P(S)$ par le processus de naissances

⇒ 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)

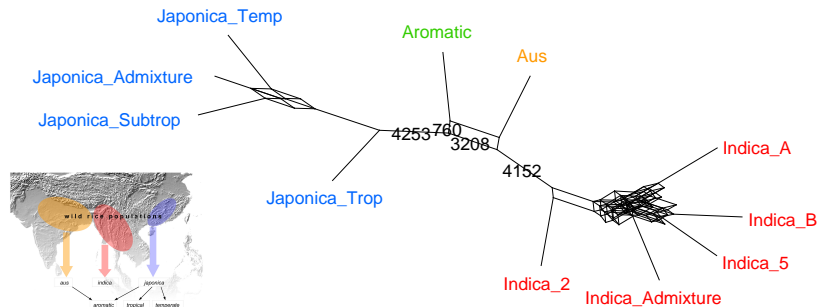


Chromosome 10 (données J. Santos, J-C. Glaszmann)

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

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

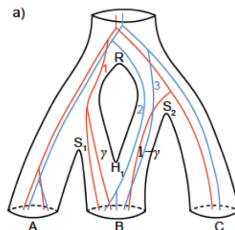
1000.0



Simulateur basé sur un réseau (Genome Harvest)

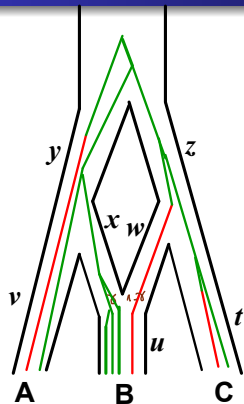
SNAPPSimNet construit sur la base du simulateur SNAPPSim de Bryant et al. (2012)

- 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

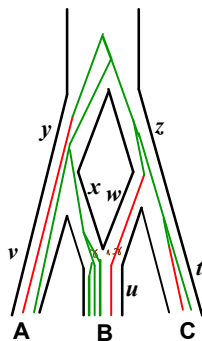
$Data_{zT}$ et $Data_{yT}$ ne sont plus indépendantes...

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

On ne peut plus effectuer le produit des probabilités

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

Notre algorithme : calcul des lois jointes



Quantités calculées successivement

- (1) $\mathbb{P}(\text{Data}_{u^T} \mid n_{u^T}, r_{u^T})$
- (2) $\mathbb{P}(\text{Data}_{x^B} \text{Data}_{w^B} \mid n_{x^B}, r_{x^B}, n_{w^B}, r_{w^B})$
- (3) $\mathbb{P}(\text{Data}_{x^T} \text{Data}_{w^B} \mid n_{x^T}, r_{x^T}, n_{w^B}, r_{w^B})$
- (4) $\mathbb{P}(\text{Data}_{x^T} \text{Data}_{w^T} \mid n_{x^T}, r_{x^T}, n_{w^T}, r_{w^T})$
- (5) $\mathbb{P}(\text{Data}_{v^T} \mid n_{v^T}, r_{v^T})$
- (6) $\mathbb{P}(\text{Data}_{y^B} \text{Data}_{w^T} \mid n_{y^B}, r_{y^B}, n_{w^T}, r_{w^T})$
- (7) $\mathbb{P}(\text{Data}_{t^T} \mid n_{t^T}, r_{t^T})$
- (8) $\mathbb{P}(\text{Data}_{y^B} \text{Data}_{z^B} \mid n_{y^B}, r_{y^B}, n_{z^B}, r_{z^B})$
- (9) $\mathbb{P}(\text{Data}_{y^T} \text{Data}_{z^B} \mid n_{y^T}, r_{y^T}, n_{z^B}, r_{z^B})$
- (10) $\mathbb{P}(\text{Data}_{y^T} \text{Data}_{z^T} \mid n_{y^T}, r_{y^T}, n_{z^T}, r_{z^T})$
- (11) $\mathbb{P}(\text{Data} \mid n_{\text{root}}, r_{\text{root}})$

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

Conclusion

- Implémentation de la méthode Bayésienne pour les réseaux
- 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
- 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) ...

Les données de riz

- ① Données disponibles à l'état brut sur le site de l'IRRI
- ② Prétraitement de Joao (données manquantes ...)
- ③ 3023 variétés avec 895 977 marqueurs disponibles sur le chromosome 6 (Merci Joao !)
- ④ 2 jeux de données proposés par JC Glazmann (core collections)
 - 20 variétés
 - 50 variétés (7 aromatic, 7 aus, 13 indica, 17 japonica, 4 indéterminés)

Constitutions de 4 jeux de données de manière empirique

- Conservation de **1550 SNPs** (un SNP tous les 500)
 - **JDD1** (1er SNP= 1er SNP du chromosome 6)
 - **JDD2** (1er SNP= 50e SNP du chromosome 6)
- Conservation de **7749 SNPs** (un SNP tous les 100)
 - **JDD3** (1er SNP= 1er SNP du chromosome 6)
 - **JDD4** (1er SNP= 50e SNP du chromosome 6)

Analyse des chromosomes 2 et 10

Chromosome 2 :

- 1 129 426 marqueurs
- Conservation de **2026 SNPs** (un SNP tous les 500)
 - **JDD1** (1er SNP= 1er SNP du chromosome 2)
 - **JDD2** (1er SNP= 50e SNP du chromosome 2)

Chromosome 10 :

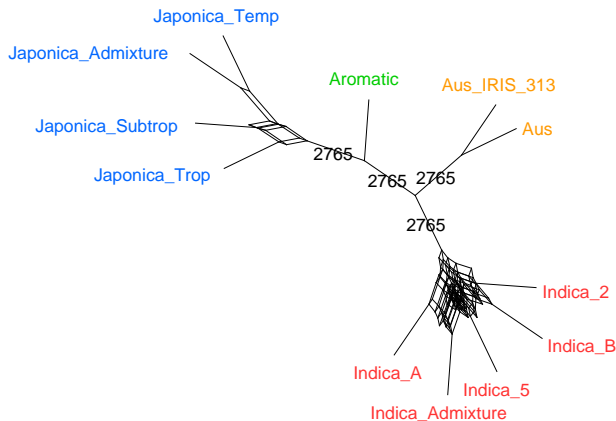
- 635 037 marqueurs
- Conservation de **1089 SNPs** (un SNP tous les 500)
 - **JDD1** (1er SNP= 1er SNP du chromosome 10)
 - **JDD2** (1er SNP= 50e SNP du chromosome 10)

Chromosome 2 (données J. Santos, J-C. Glaszmann)

Conservation de **2026 SNPs** (un SNP tous les 500)

- **JDD1** (1er SNP= 1er SNP du chromosome 2)

1000.0

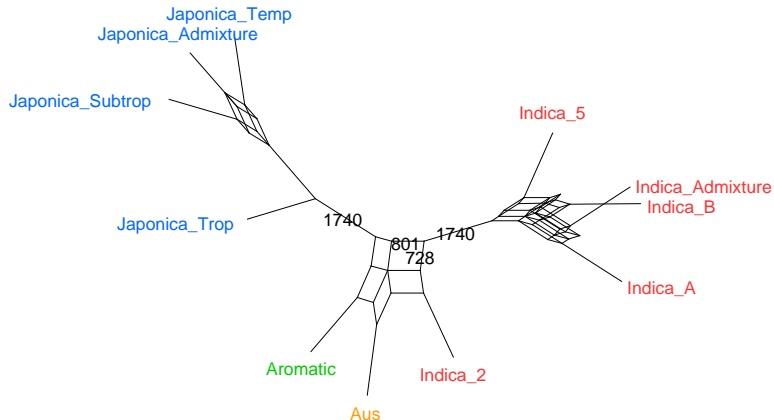


Chromosome 10 (données J. Santos, J-C. Glaszmann)

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

- **JDD1** (1er SNP= 1er SNP du chromosome 10)

1000.0



Calcul de vraisemblance dans un arbre

Où se situe l'aléatoire dans le modèle ?

- Modèle de coalescence \rightarrow pas d'aléatoire dans le nombre de lignées (Count) dans chaque espèce !
- L'aléatoire réside dans la répartition d'allèles **rouges** et **verts** dans chaque espèce

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