



HAL
open science

Combiner arbres phylogénétiques et visualisation d'ensembles

Jean-Baptiste Lamy, Flora Jay

► **To cite this version:**

Jean-Baptiste Lamy, Flora Jay. Combiner arbres phylogénétiques et visualisation d'ensembles. Atelier Visualisation d'informations, interaction et fouille de données, Jan 2020, Bruxelles, Belgique. hal-02968296

HAL Id: hal-02968296

<https://hal.science/hal-02968296>

Submitted on 15 Oct 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Combiner arbres phylogénétiques et visualisation d'ensembles

Jean-Baptiste Lamy^{*,**}, Flora Jay^{*,***}

^{*}Laboratoire de Recherche en Informatique,
CNRS/Université Paris-Sud/Université Paris-Saclay, Orsay, France
flora.jay@lri.fr

^{**}LIMICS, Université Paris 13, Sorbonne Paris Cité, 93017 Bobigny,
INSERM UMRS 1142, Sorbonne Universités
jean-baptiste.lamy@univ-paris13.fr

^{***}Laboratoire EcoAnthropologie et Ethnobiologie,
CNRS/MNHN/Université Paris Diderot, Paris, France

Les arbres sont très largement utilisés en phylogénie. Cependant, un arbre avec n feuilles présente les similarités entre $n - 1$ sous-ensembles de feuilles, alors qu'il existe $2^n - n - 1$ sous-ensembles possibles d'au moins deux feuilles. Par exemple, pour 3 feuilles A , B et C , 4 similarités peuvent être mesurées, entre les sous-ensembles $\{A, B\}$, $\{B, C\}$, $\{A, C\}$ et $\{A, B, C\}$, mais un arbre n'en montrera que 2, *e.g.* $\{A, B\}$ et $\{A, B, C\}$ si une branche rassemble A et B . Plus n augmente, plus la vision donnée par l'arbre deviendra réductrice.

Ce problème est particulièrement important en génétique des populations (Pickrell et Pritchard, 2012), lorsque l'on étudie la diversité génétique des populations d'êtres vivants et leurs relations. Dans ce contexte, il est nécessaire de tenir compte des processus biologiques telle l'apparition de mutations dans le génome mais aussi des processus démographiques, tels que la séparation des populations (aussi appelée divergence), la variation de leur taille effective, et les migrations (mouvement d'un groupe d'individus qui quittent une population pour en rejoindre ou en créer une autre, apportant au passage leur matériel génétique). Un arbre généalogique peut représenter la composante évolutive et une partie de la composante démographique (divergence simple des populations, dérive génétique plus forte dans les populations de petite taille, etc) mais pas la composante migratoire post divergence qui induit des "flux de gènes" entre branches de l'arbre.

Afin de résoudre ce problème, nous proposons l'utilisation de visualisation d'ensembles, et notamment des boîtes arc-en-ciel (Lamy et al., 2017) et de leur variante proportionnelle (Lamy et Tsopra, 2019), et l'illustrons par une application à un sous-ensemble de données extraites du *1000 Genomes Project* (Auton et al., 2015). Une première approche consiste à visualiser les similarités comme des ensembles (sous-ensemble des populations avec les mêmes mutations).

Une seconde approche consiste à superposer arbre phylogénétique et boîtes arc-en-ciel (voir Figure 1). L'arbre (en noir) représente l'histoire démographique "prépondérante" des populations, lesquelles sont présentées en colonne (les couleurs dans les en-têtes de colonne identifient les continents). La longueur des branches de l'arbre indique le nombre de mutations depuis l'ancêtre commun. Les boîtes rectangulaires permettent de visualiser les similarités entre branches éloignées de l'arbre : par exemple la boîte bleu ciel à gauche montre une similarité

