

# Travaux & problématique en métagénomique au CBiB

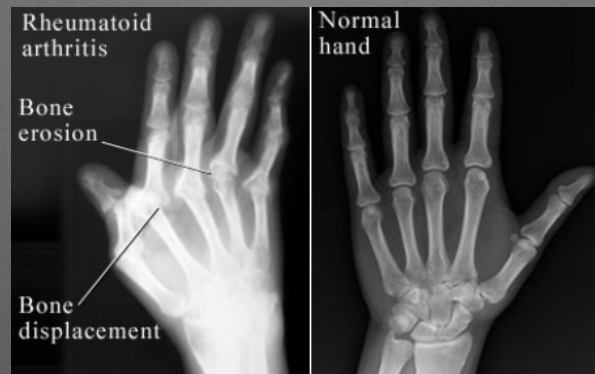
Groupe de travail « France Génomique Métagénomique »  
BARRE Aurélien - Paris - 30 juin 2015







Virus de plantes



Santé



Paléo métagénomique

### Questions :

- Connaitre la liste des espèces présentes
- Comparer des échantillons
- Analyser fonctionnellement le milieu



# **Classification génomique bactérienne**

2013-2015



# Polyarthrite rhumatoïde



- Etude de faisabilité avec un séquençage de reads courts (75 nt) avec la technologie Illumina
- Gestion des ambiguïtés résultant du mapping multiple des reads sur les séquences 16 S
- Prise en compte de la phylogénie des bactéries
- Recherche de lien entre la composition en bactéries & le développement de la maladie



## Etapes du projet

- Phase 1 : Etude de faisabilité (3 échantillons)
- Phase 2 : Analyse comparative entre patients sains et malades en aveugle (11 échantillons)
- Phase 3 : Analyse de variation de la composition au cours du traitement & de la prédisposition à la réponse au traitement (12 échantillons)



## Les taxonomies disponibles

- NCBI Taxonomy (<http://www.ncbi.nlm.nih.gov/Taxonomy/>)
- Greengenes (<http://greengenes.lbl.gov/>)
- Ribosomal Dataset Project (<http://rdp.cme.msu.edu/>)

## Formulation du problème

**In** A genomic reference  $\mathbf{S}$  (set of sequences)

A taxonomic reference  $\mathbf{T}$  (tree)

A leaf set  $\mathbf{L}$ , s.t.  $\forall Li \in \mathbf{L}$  has an associated sequence of  $\mathbf{S}$

A set  $\mathbf{R}$  of sequence reads

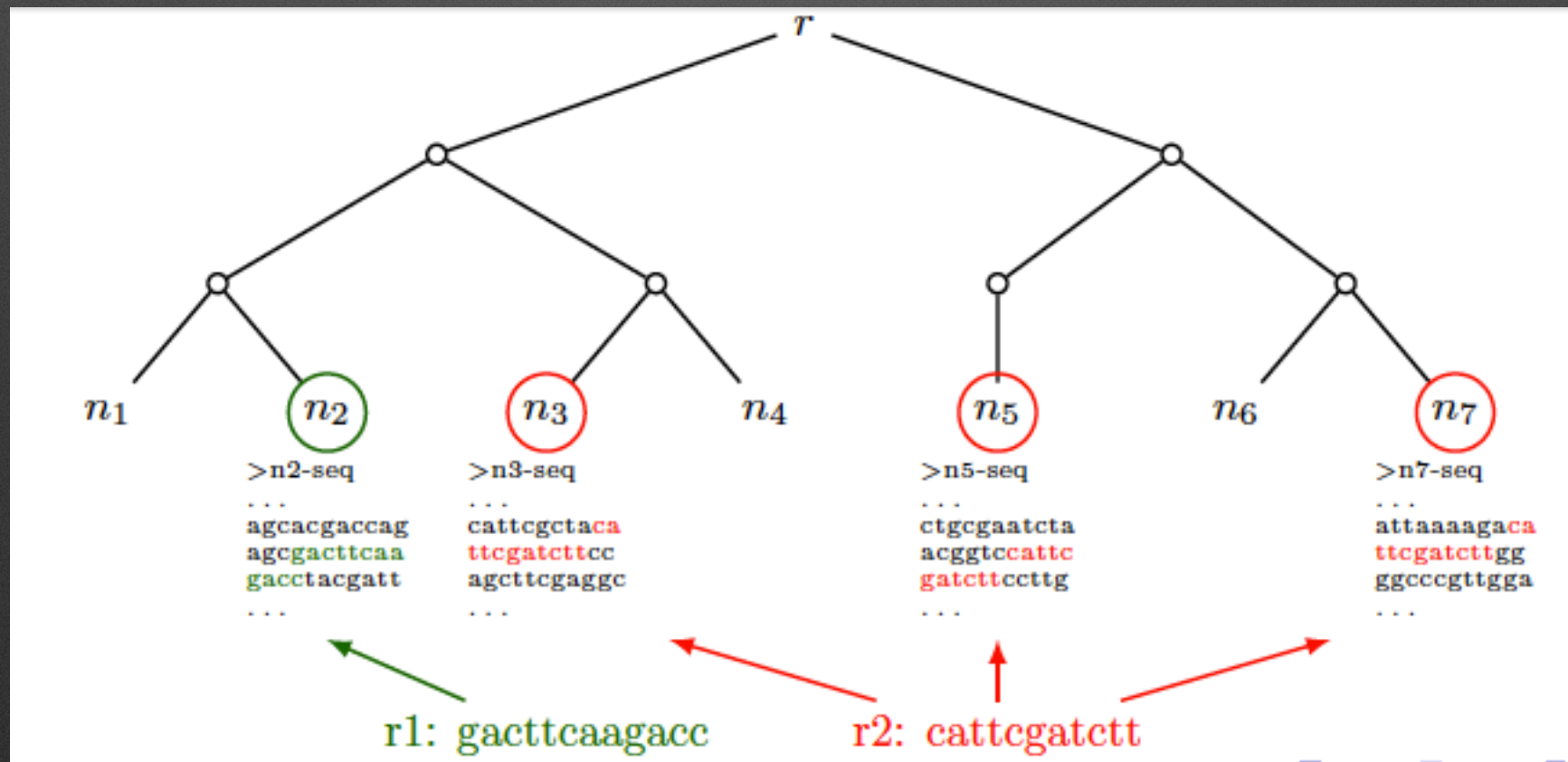
**Out** For each read  $Ri \in \mathbf{R}$ , a node in  $\mathbf{T}$  that represents a “good” subset of matches  $\mathbf{Mi} \in \mathbf{L}$



# 1 - Mapping

## Alignement des reads contre les séquences de référence

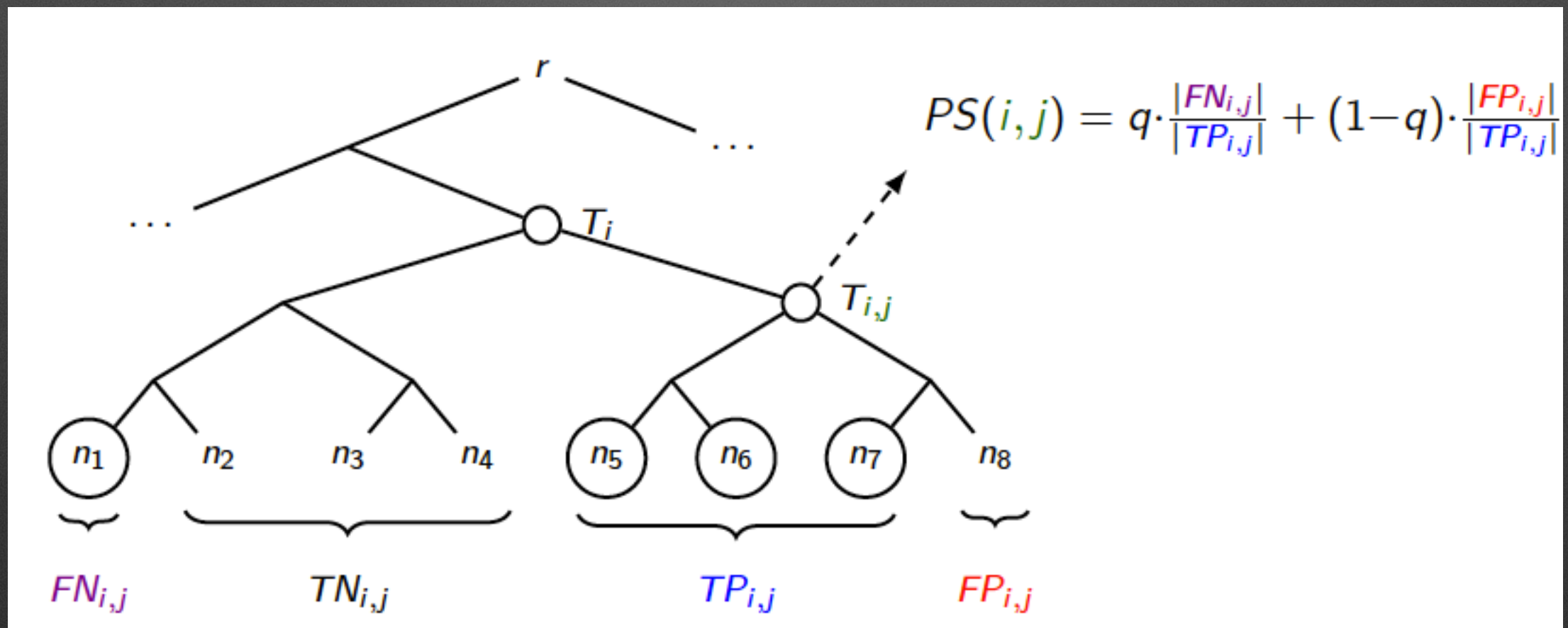
- single match -> **read non ambigu**
- multiple matches -> **read ambigu**





## 2 - Assignment taxonomique

Recherche de la meilleure assignation basée sur le calcul d'un *Penalty Score* (PS)





### 3 - Taxonomy equalizing

**Input:**

- Two taxonomic trees  $T_1, T_2$
- Leaf mapping  $\emptyset: L(T_1) \rightarrow L(T_2)$

**Output:**

- Node mapping  $\emptyset: T_1 \rightarrow T_2$

**Algorithmic Solution**

- *post-order* traversal of the tree  $T_1$
- for each node  $ni \in T_1$ :  
$$\emptyset(ni) = LCA(\{\emptyset(nk) \mid nk \in \text{desc}(ni)\})$$

**where:**

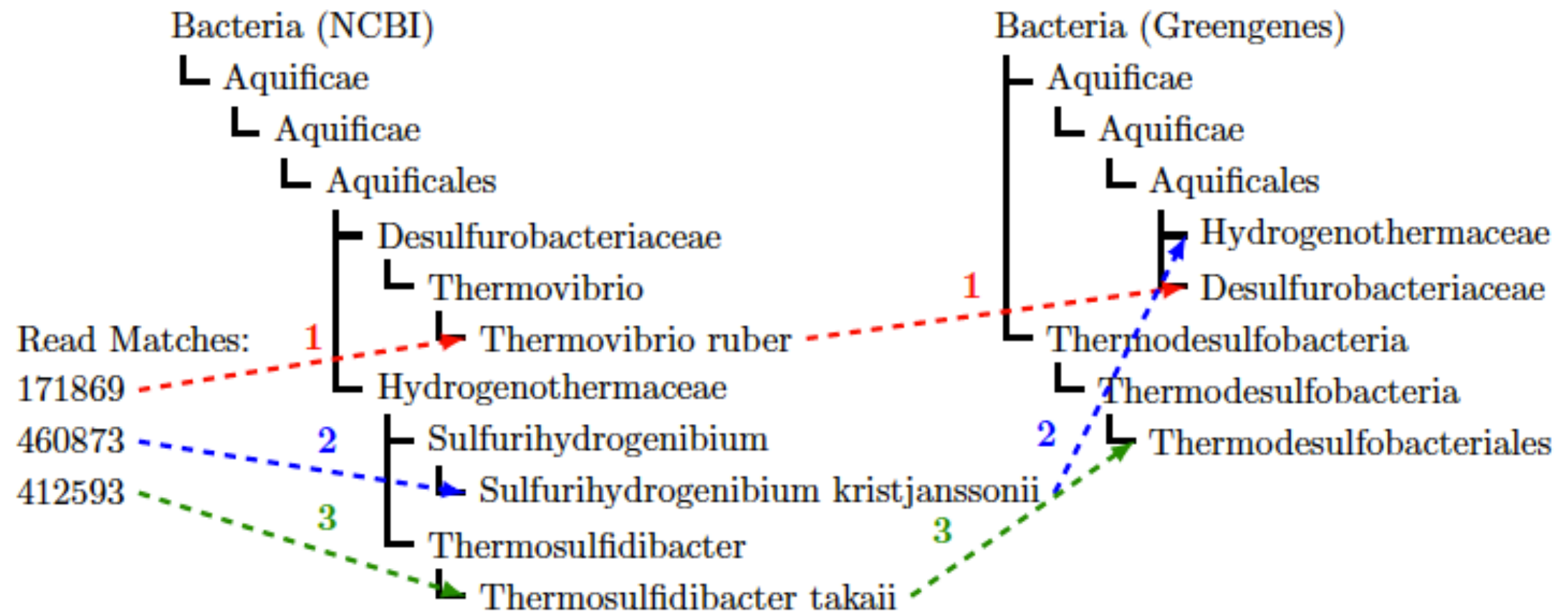
- $LCA(X)$  is the *Lowest Common Ancestor* of nodes  $xi \in X$
- $\text{desc}(x)$  is the set of descendants of the node  $x$





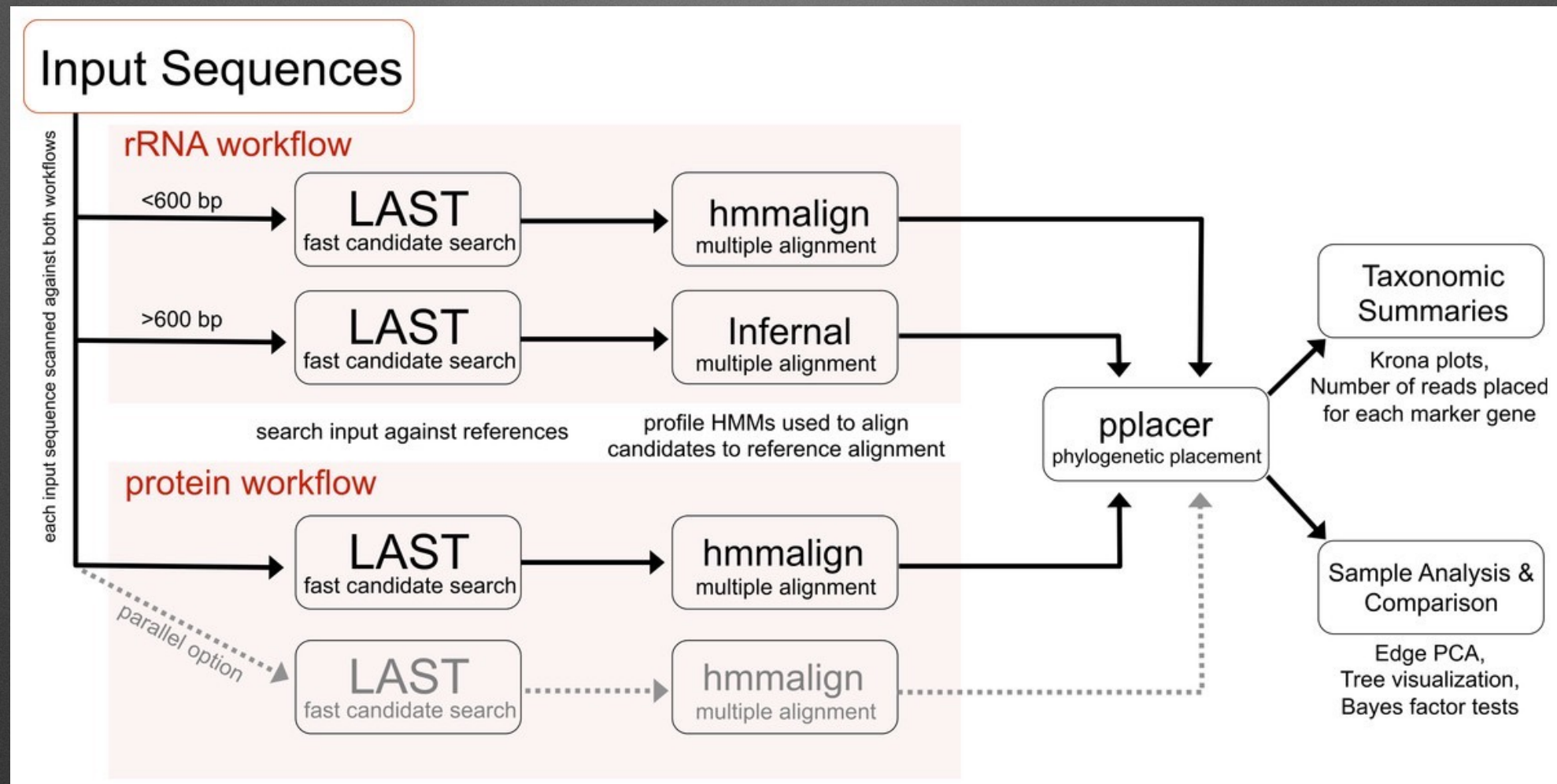


# Reads relabeling





# Alternative pour l'assignation taxonomique : Phylsift





# Pipeline

## Phases 1 et 2

Alignement avec GEM  
contre GreenGenes

puis

Assignation taxonomique  
avec Tango

## Phases 3

Alignement avec BWA  
contre GreenGenes

puis

Assignation taxonomique  
avec Tango et confirmation  
par Phylosift

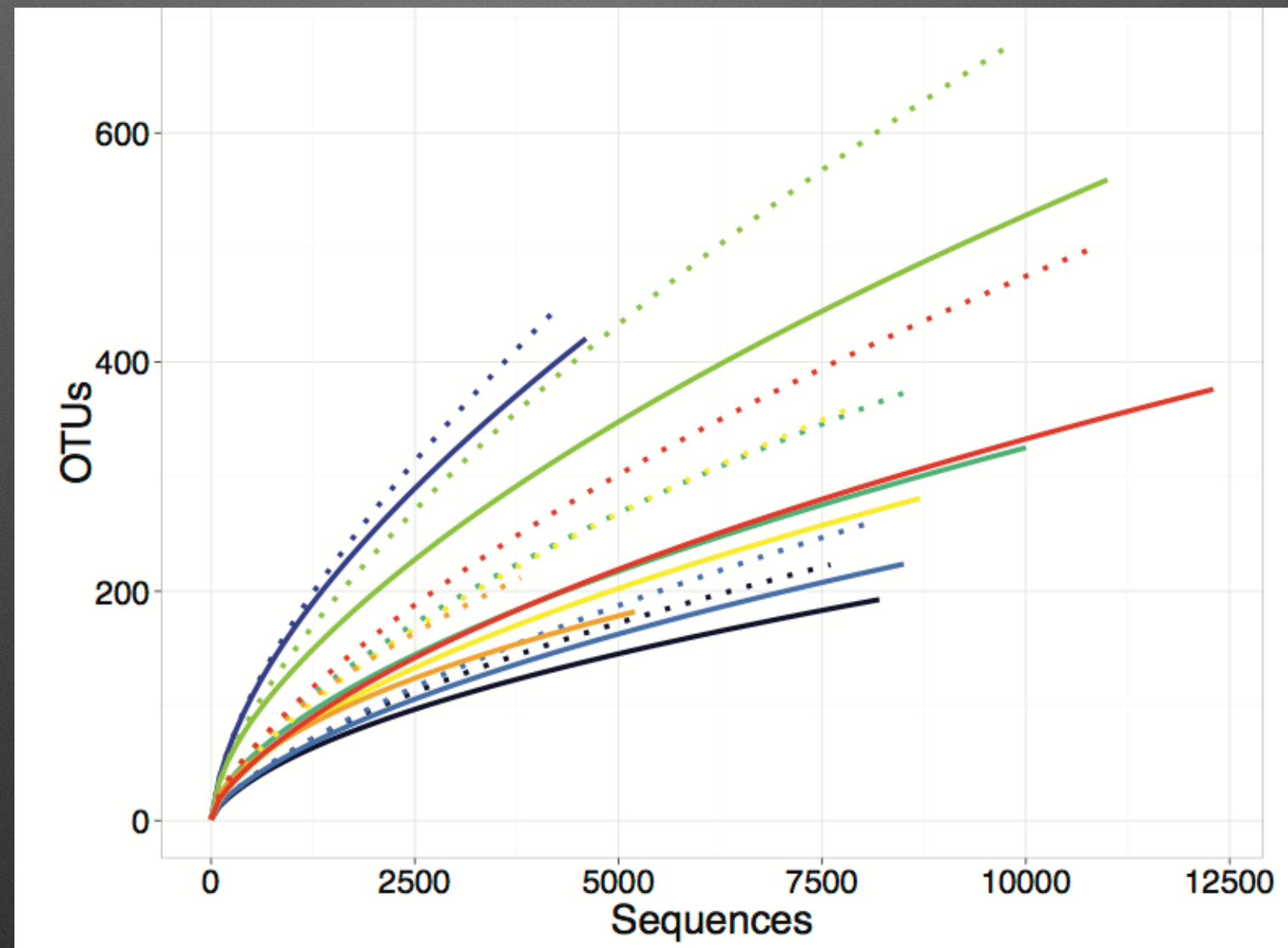


## Comparaison

- Normalisation basée sur les k-mer (logiciel khmer) et sur le nombre de reads produits dans chaque échantillon
- Estimation de la variabilité intra-condition et inter-condition
- Analyse statistique classique des noeuds basée sur une version améliorée du package bioconductor « metagenomeseq »
- Représentation graphique des résultats

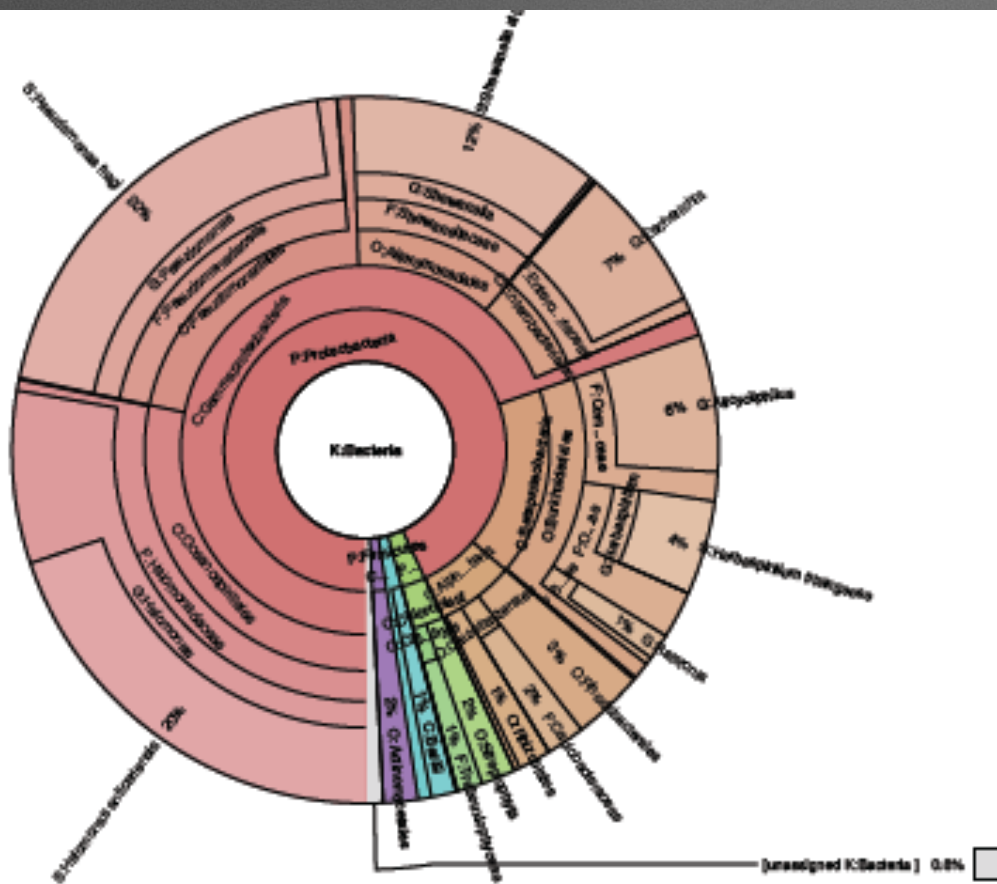


Les données de séquençage fournies pour l'analyse ne sont malheureusement pas toujours optimales et suffisantes

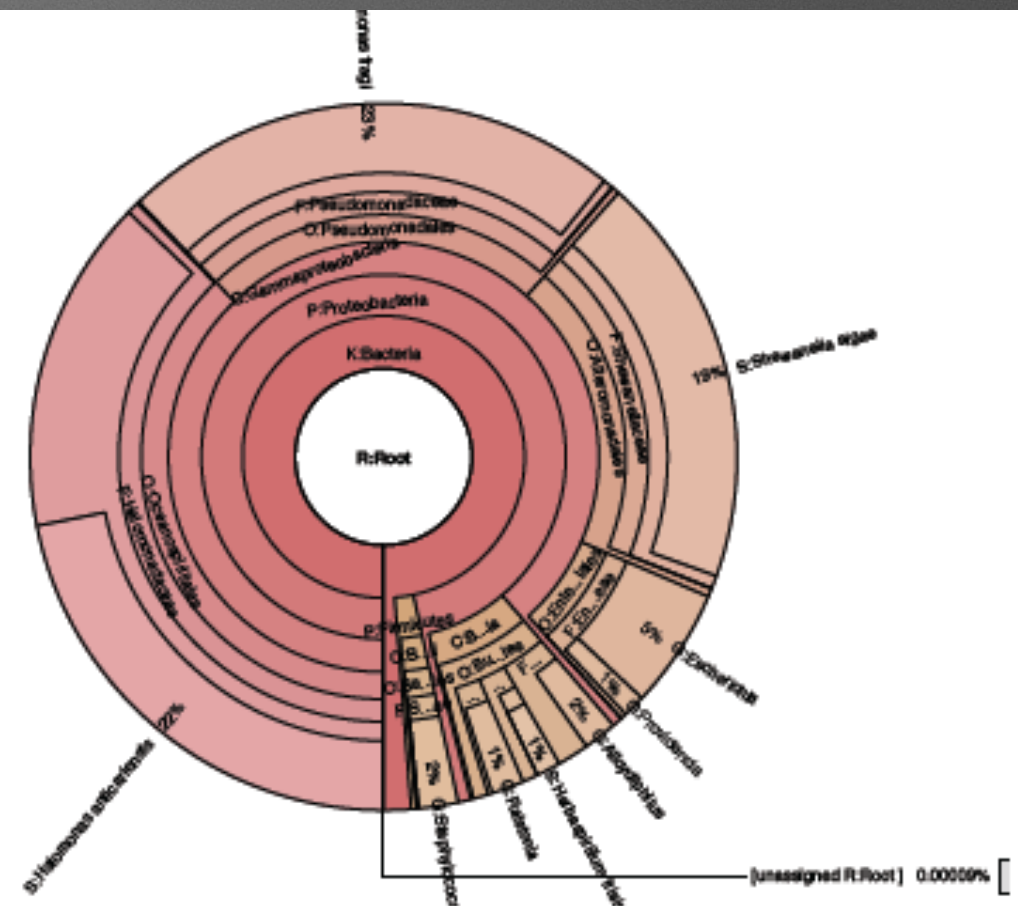




# Représentation type Krona produit par Phylosift et Tango



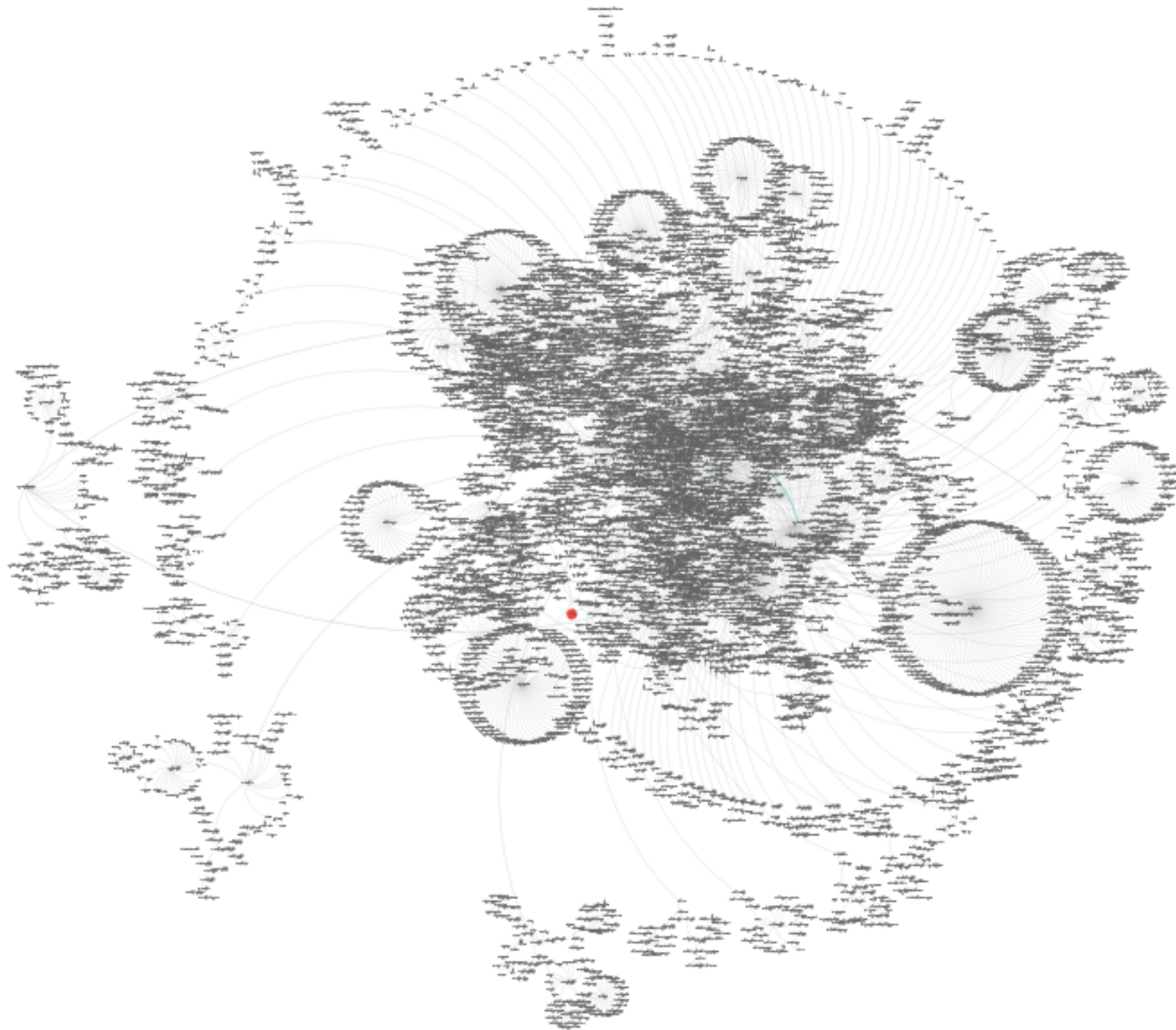
Healthy



## Rheumatoid arthritis



## GreenGenes DB (450.000 nodes)





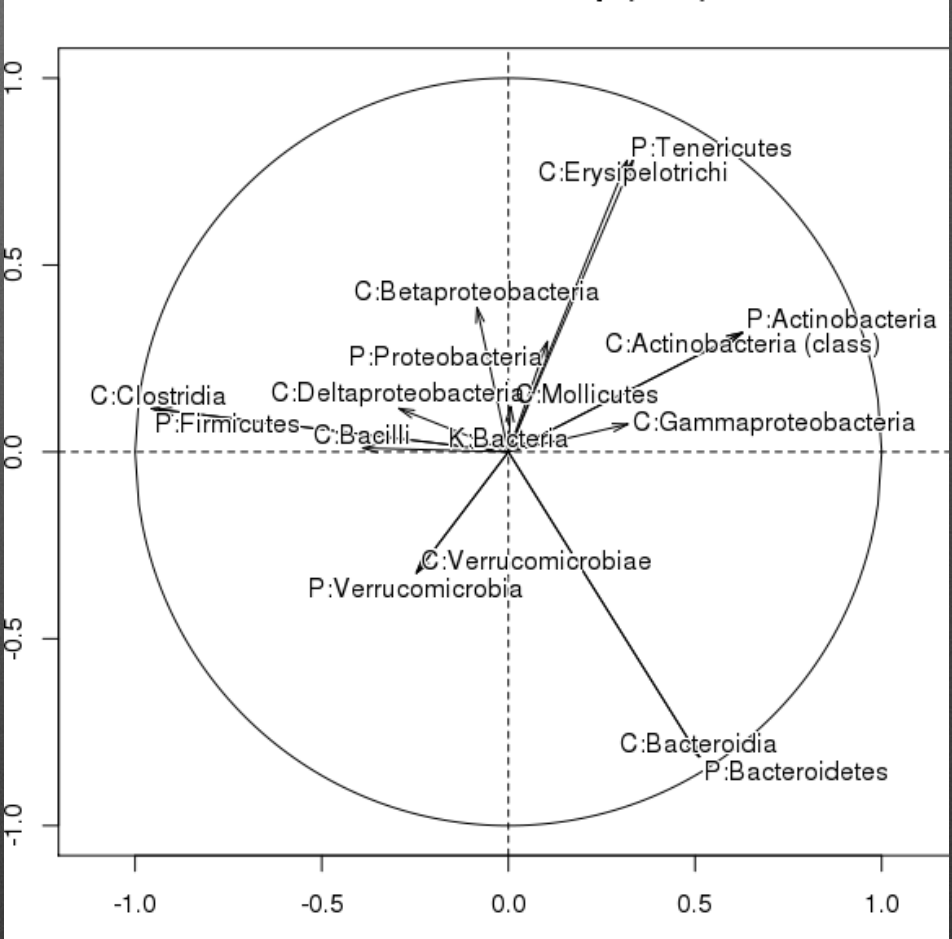
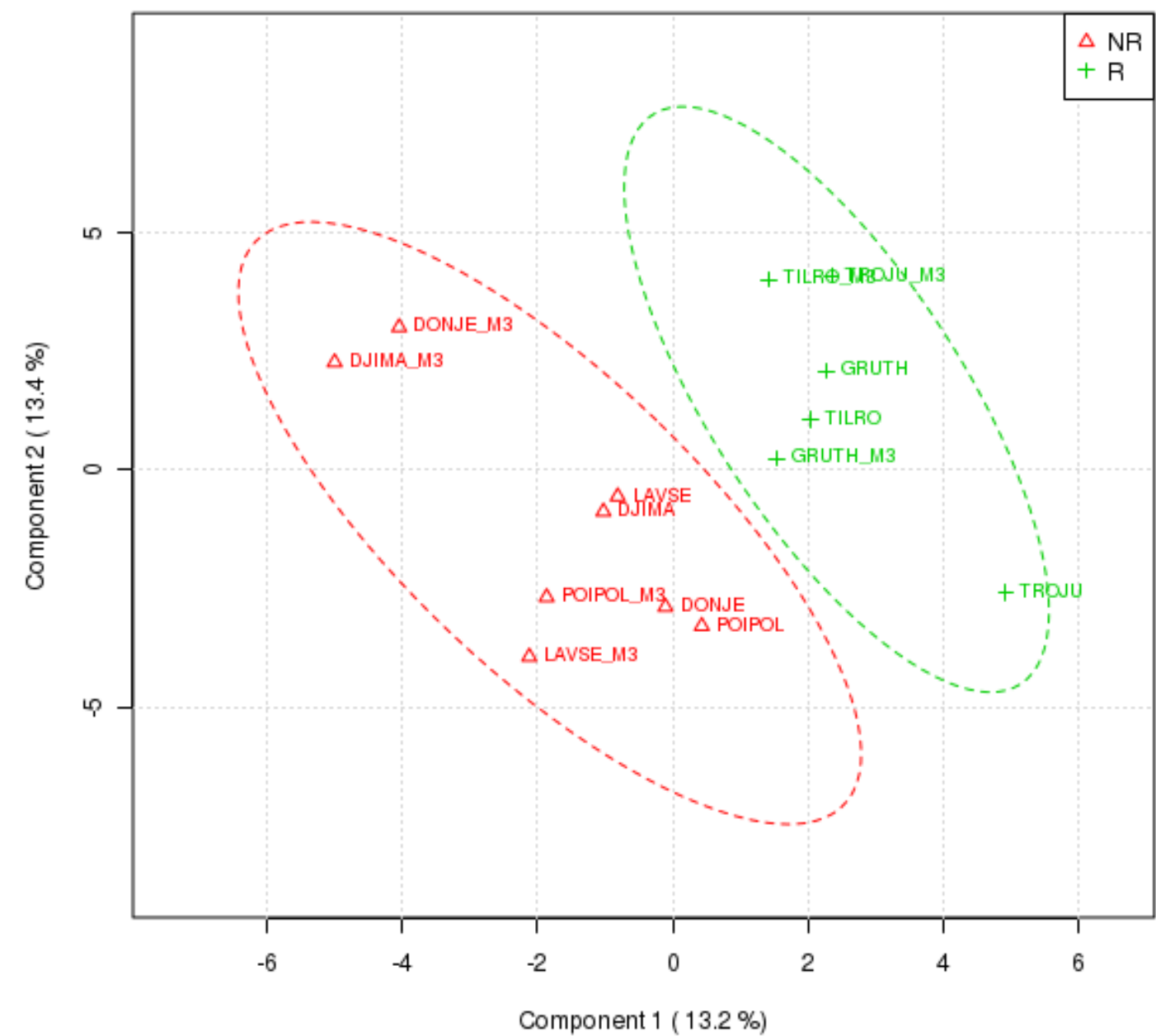
## Phase 2





Phase 3

Score Plot





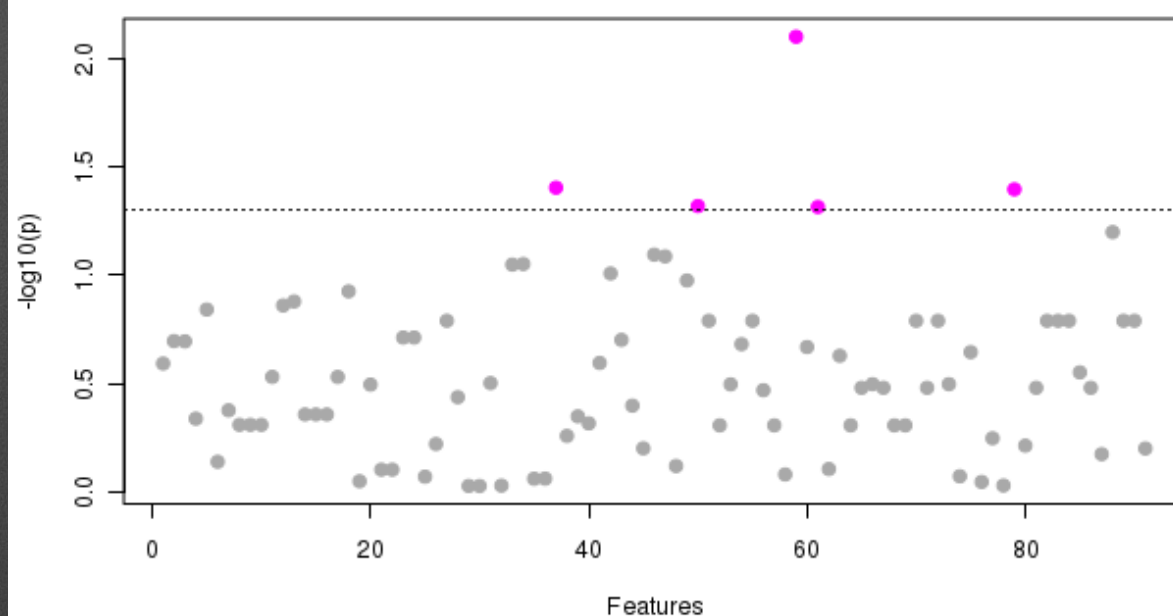


- Ensemble de méthodes pour l'importation, le normalisation, l'analyse et la visualisation des données métagénomiques
- Workflow très flexible
- Post traitement possible avec R (Rdata format)

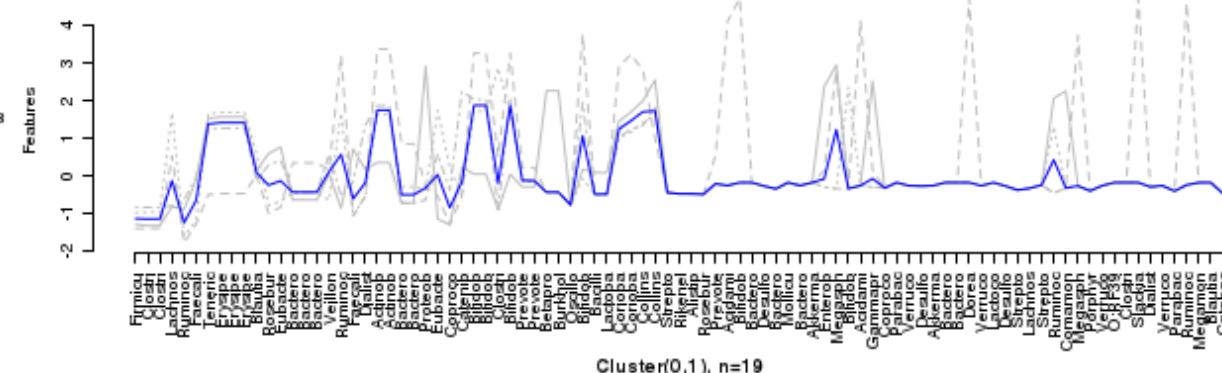
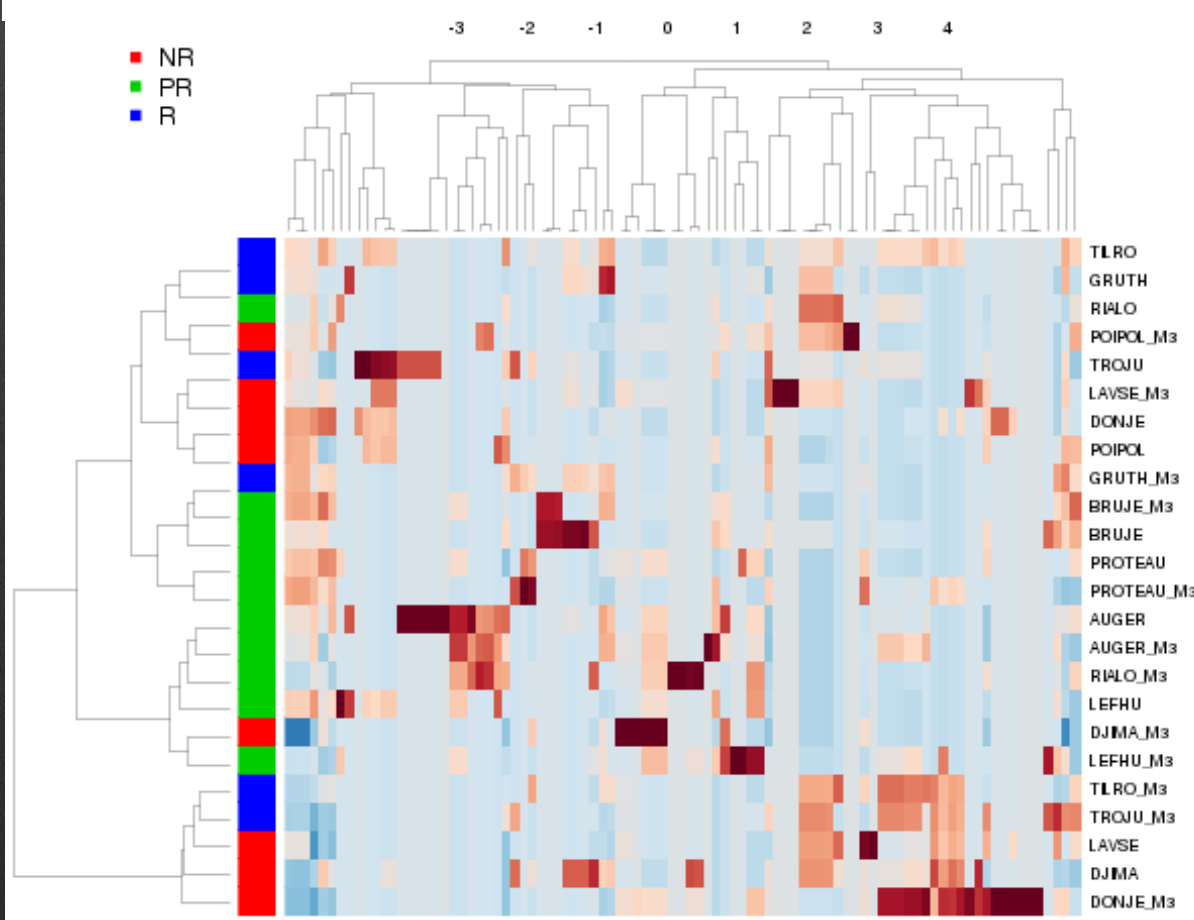
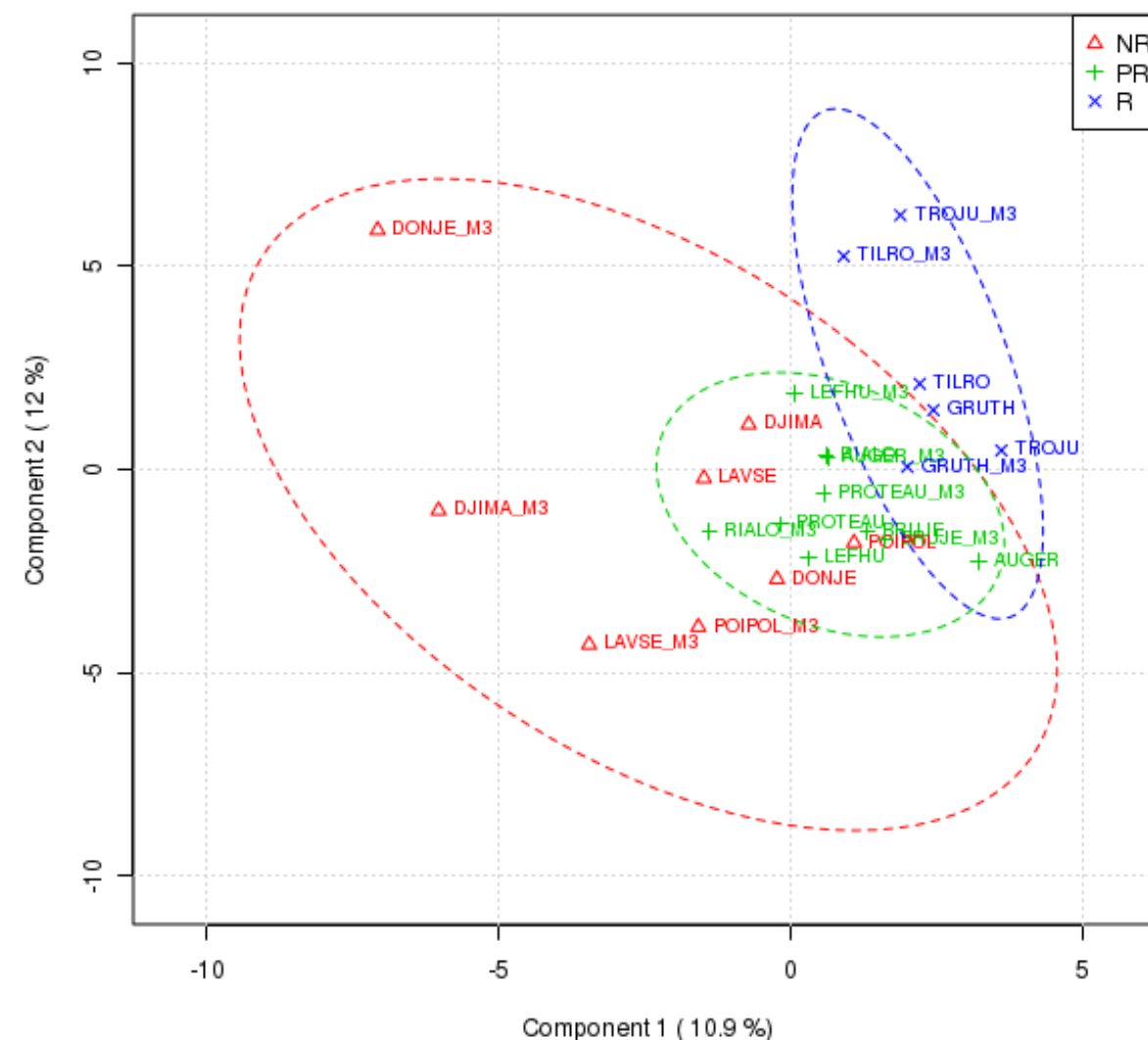




Student's t-tests



Score Plot







Développement en cours d'un workflow dans l'esprit de MetagenAssist mais qui intègre des étapes en amont (assignation taxonomique) & des étapes en aval (visualisation des résultats sur une phylogénie)

Intégration possible dans Galaxy

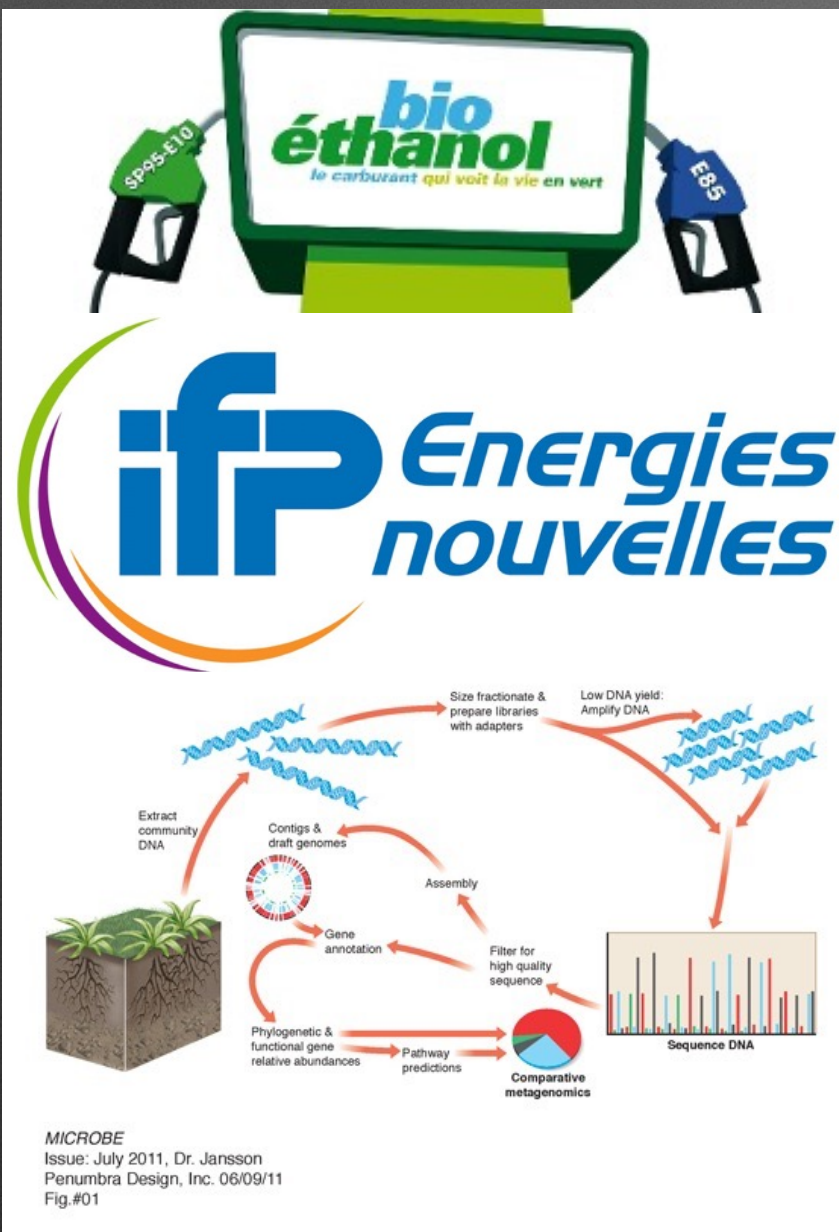


# Métagénomique fonctionnelle

2014-2015

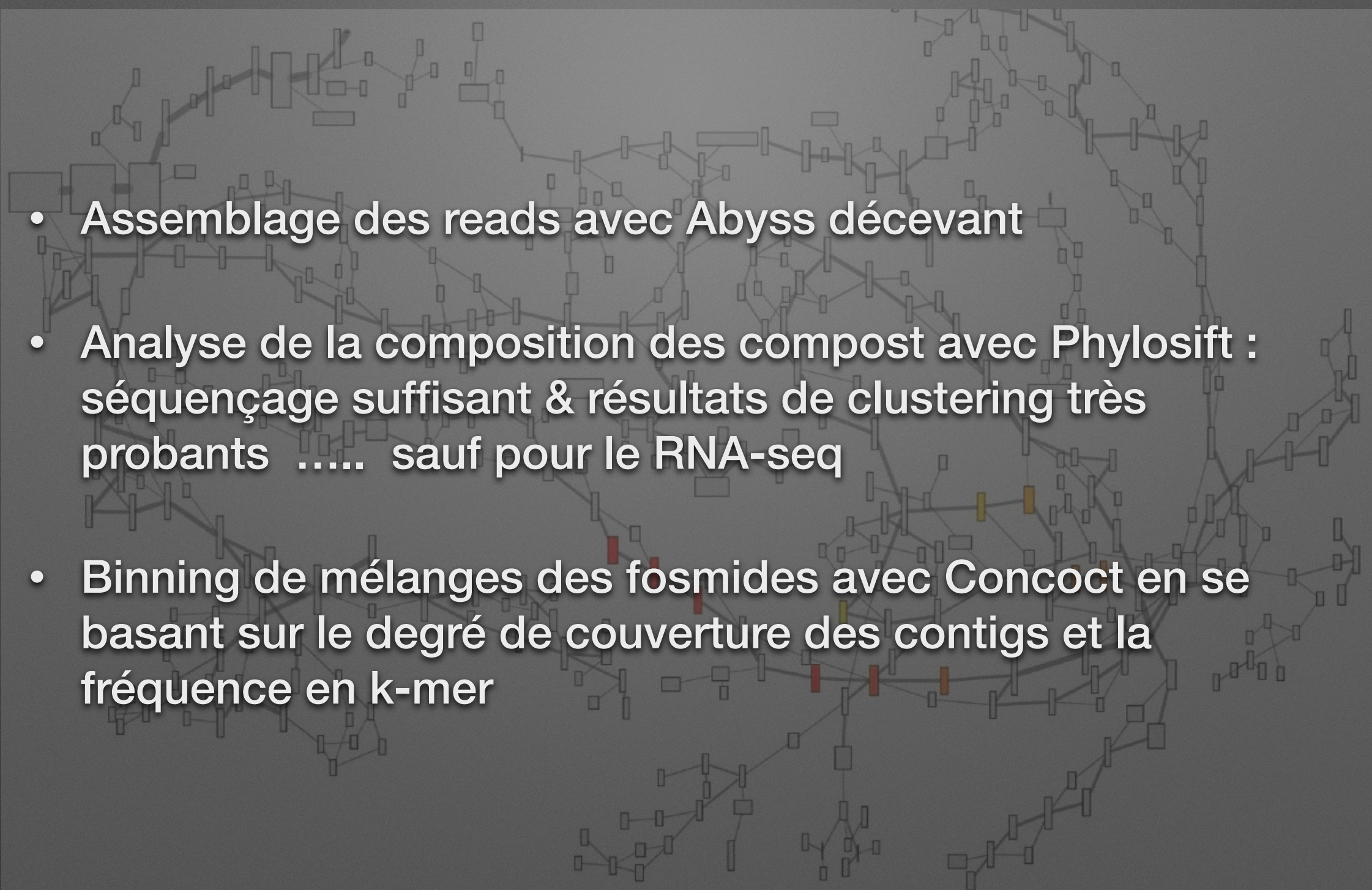


# Projet Biomines



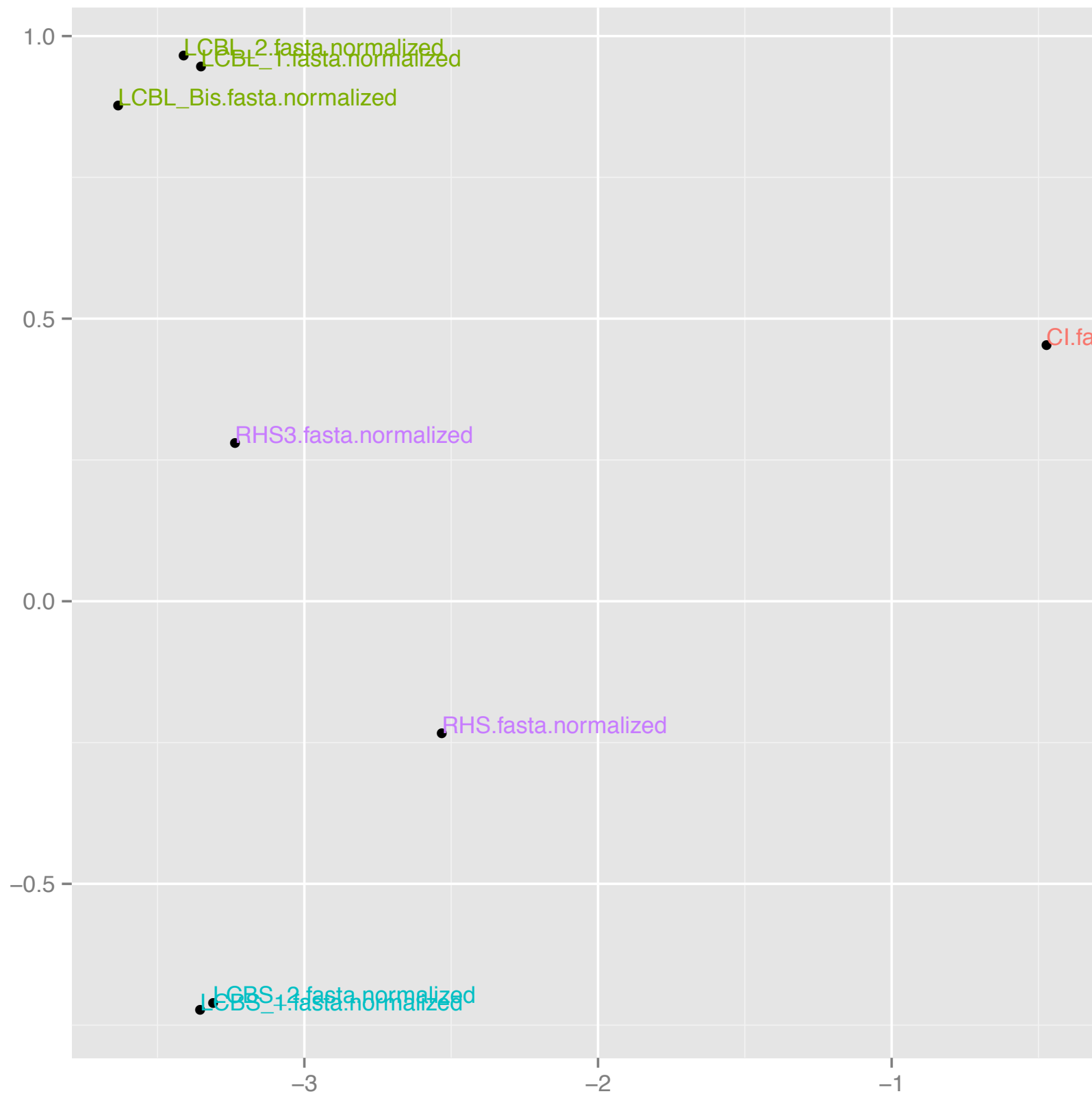
- Déterminer un nouveau bioprocess pour produire bio-ethanol
- Etude de compost pour en déterminer la composition en micro organismes et le potentiel enzymatique de 8 échantillons
- Séquençage principalement en 454
- Analyse DNA-seq et RNA-seq



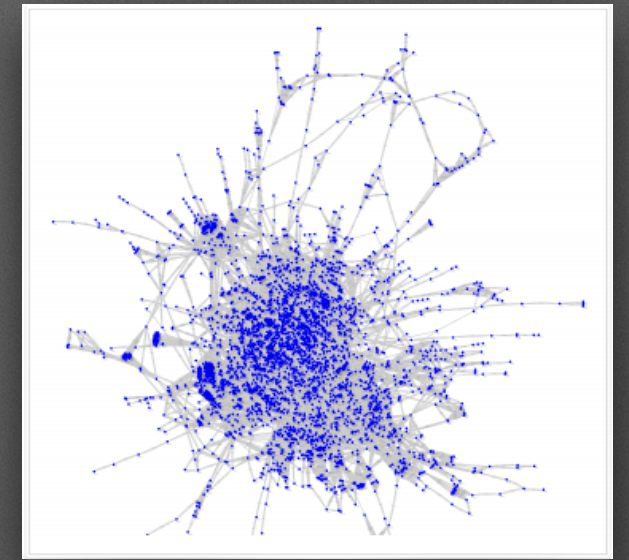
- 
- Assemblage des reads avec Abyss décevant
  - Analyse de la composition des compost avec Phyllosift : séquençage suffisant & résultats de clustering très probants ..... sauf pour le RNA-seq
  - Binning de mélanges des fosmides avec Concoct en se basant sur le degré de couverture des contigs et la fréquence en k-mer



# EPCA







- Test du package « mmnet » (microbiome metabolic network) qui intègre les données du KEGG et des prédictions enzymatiques issues de MGRast pour construire des réseaux qui peuvent être ensuite comparés.
- Annotation fonctionnelle des composés par recherche d'ORF sur les reads avec FragGeneScan puis aggrégation d'information (tigrfam et blast)
  - + Etape de réduction des données avec uclust



## Paleo - metagenomic



- Analyse d'échantillons d'ADN ancien de tombes
- Recherche de la composition en espèce de ces échantillons afin de déterminer les rites funéraires
- Etape de filtrage pour dissocier l'ADN ancien des contamination contemporaines grace aux dégradation spécifiques (substitution G-> A)





## Partenaires



M. Nikolski  
A. Groppi



G. Valiente  
D. Alonso-Alemany  
S. Beretta



Mesocentre de Calcul Intensif Aquitain  
(MCIA)