

EVOL TME 02

WANG Jinxin 3404759
Spécialité Bioinformatique et Modélisation
MASTER D'INFORMATIQUE Niveau 2
Département Master d'Informatique
Université Pierre et Marie CURIE

A : Récupérer les génomes à étudier

Une python script *gb2aa.py* est implémenté pour extraire CDS et traduire les gènes aux séquences protéiques.

La script prend trois paramètres obligatoires :

paramètres	Explications
<i>-gb</i> ou <i>--geneBank-file-name</i>	Gene bank file name to read
<i>-f</i> ou <i>--fasta-file-name</i>	Fasta file name to read
<i>-o</i> ou <i>--out-fasta-file-name</i>	Fasta file name to write

B : Construction des familles de gènes

Effectuer les comparaisons nécessaires pour :

Dresser la liste des Reciprocal Best Hits entre les 3 paires de génomes.

Pour réaliser le sujet, la base de donnée sont construit d'abord en utilisant *makeblastdb*, les données sont alignés par *blastp*.

Pour sélectionner les meilleurs hits, *select_best_blast_hits.py* a plusieurs paramètres :

paramètres	Explications
<i>-f</i> ou <i>--hits-file-name</i>	input blastp hits file name
<i>-o</i> ou <i>--output-file-name</i>	output blastp hits file name
<i>-i</i> ou <i>--iden-percent</i>	% identity
<i>-a</i> ou <i>--align-len</i>	alignment length
<i>-m</i> ou <i>--mismat</i>	mismatches
<i>-g</i> ou <i>--gaps</i>	gap opens
<i>-e</i> ou <i>--evaluate</i>	evaluate
<i>-b</i> ou <i>--bit-score</i>	bit score

J'ai choisi *evaluate* comme la seuil à 1, par conséquent, tous les hits qui ont une *evaluate* supérieur que 1 sont rejetés.

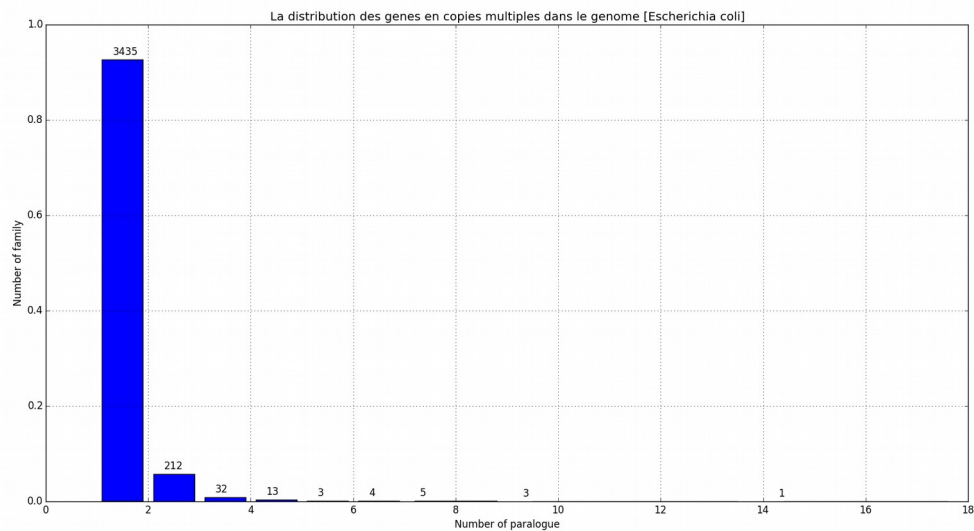
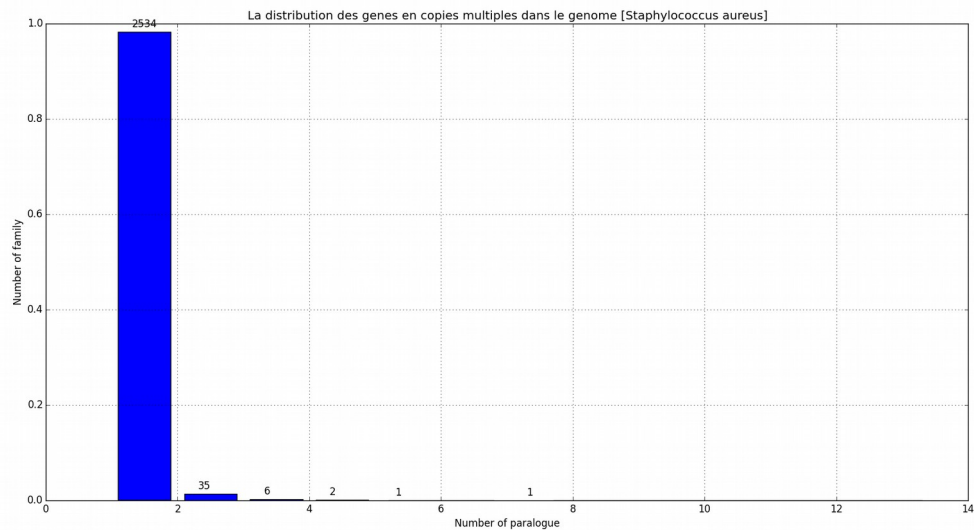
Dresser la liste des familles de gènes homologues partagés par ces 3 génomes.

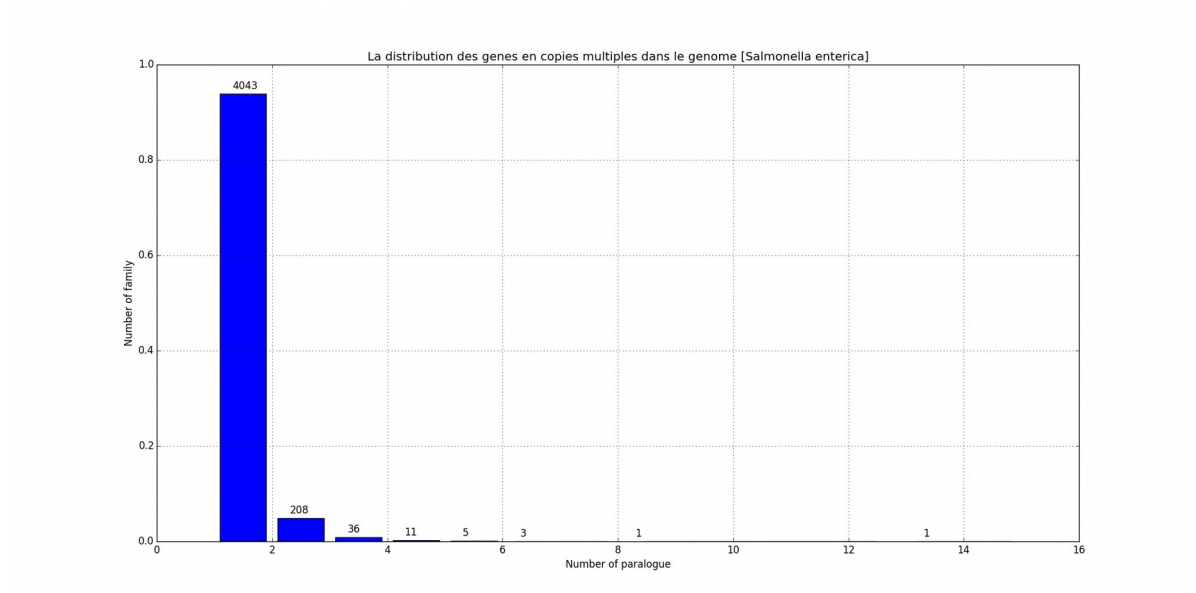
Ensuite, les paramètres par défaut de *silix* sont choisis.

C : Analyse

- Pour chaque génome, faire un histogramme de la distribution des gènes en copies multiples dans le génome.

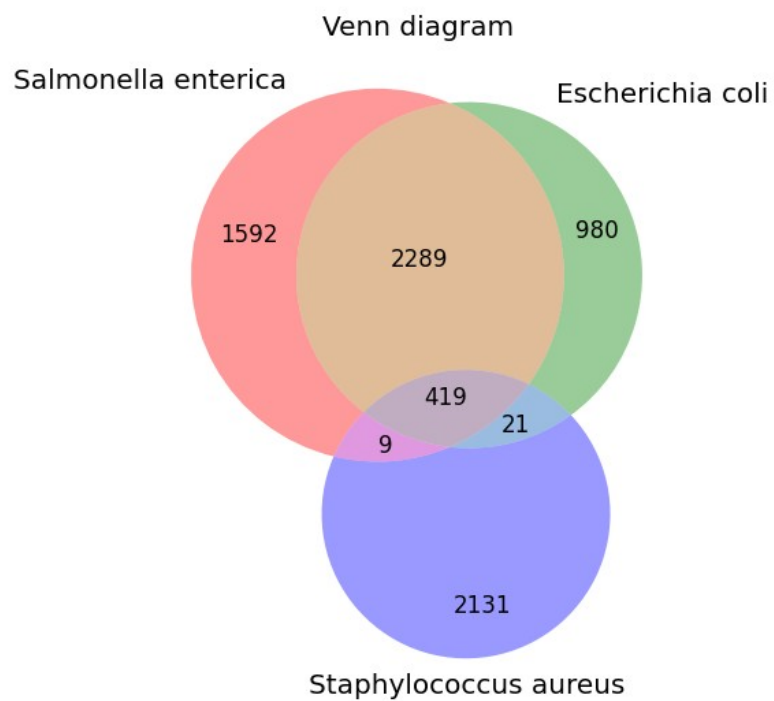
La script `dist_hist.py` est implémentée pour ploter les histogrammes.





- Réaliser un diagramme de Venn afin de visualiser :

Le plot est généré par `venn_list`.



- Proposer une première analyse biologique de ces résultats

Dans les histogrammes de la distribution des gènes en copies multiples dans le génome, on observe que les orthologues sont très et diminuer exponentiellement. On peut aussi conclure que dans les gènes, il y a des gènes très exprimés.

Le diagramme de Venn montre que Salmonella Enterica et Escherichia Coli ont beaucoup d'orthologues en commun, mais un peu éloigné de Staphylococcus Aureus. Pour cela, nous pouvons dire que les activités des Salmonella Enterica et Escherichia Coli sont plus proches que Staphylococcus Aureus.