

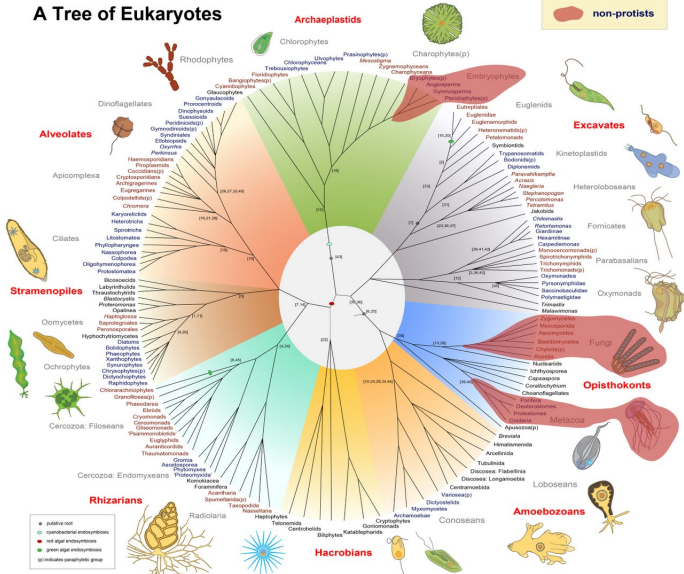
L'exploration de la diversité des protistes: l'apport du calcul intensif

J.-M. Frigerio, P. Chaumeil, F. Rué, S. Thérond, V. Louvet, O.
Coulaud & A. Franc

INRA BioGeCo; INRIA Pleiade, Hiepacs, SED; CNRS IDRIS & GRICAD

JCAD, Lyon
26 octobre 2018

L'immense diversité des Eucaryotes !



- putative root
- cyano-bacterial endosymbiosis
- not-algal endosymbiosis
- green-algal endosymbiosis
- [(p)] indicates paraphyletic group

References:
1. Surana & Surana 2009 The evolutionary phylogeny of ctenophores: insights gained from studies of ribosomal genetics of algae and invertebrates. In
24. Parveen & Burt 2000 Integrating the phylogeny of animal protists. J Mol Evol 50: 16-25
25. Pelletier & Cavalier-Smith 2001 Mitosis, flagellar motility and the systematics of eukaryotes

- 1 Les données
- 2 Calcul des distances
- 3 Format et transfert des fichiers
- 4 Image euclidienne et réduction de la dimension
- 5 Projection aléatoire et librairie *Diodon*
- 6 Construction d'OTUs = classification non supervisée

Pour chaque échantillon

- un fichier de format dit fasta
- ascii
- de texte (`strings`)
- avec pour chaque séquence
 - un identifiant `>taratata`
 - une séquence `acgtgcatgcatgc...`

Tailles

- un fichier: entre 20 *k* et 150 *k* séquences
- taille ≤ 54 Mo
- typiquement 100 échantillons donc fichiers par projets

Distance entre séquences

Distance d'édition (Levenstein, 1965)

Etant données deux séquences s et s' , le plus petit nombre d'opérations élémentaires parmi

- substitution
- insertion
- délétion

nécessaires pour transformer s en s' .

Algorithme de Needleman-Wunsch (NW, 1970)

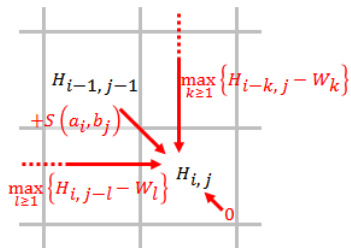
- algorithme de programmation dynamique

Algorithme de Smith-Waterman (SW, 1983)

- alignement local: une chaîne peut être incluse dans une autre avec $d = 0$
- de complexité $\mathcal{O}(ll')$

Algorithme de SW

-	-	A	T	C	G	A	A
-	0	0	0	0	0	0	0
C	0	0	0	5	1	0	0
A	0	5	1	1	2	5	5
T	0	1	10	6	2	1	2
A	0	5	6	7	3	7	6
C	0	1	2	11	7	3	4



Algorithm 1 pseudocode pour construire la matrice de distances

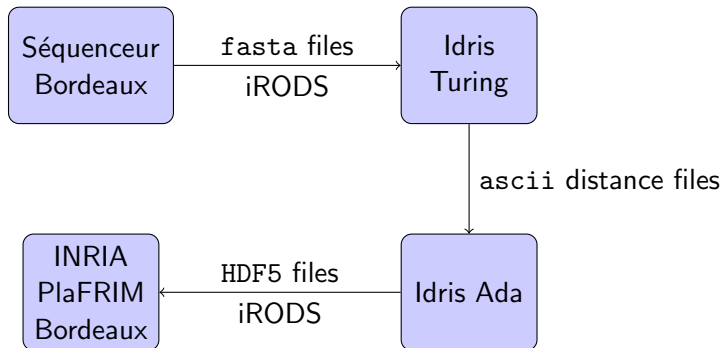
```
1: for  $s \in \llbracket 1, n-1 \rrbracket$  do  
2:   for  $s' \in \llbracket s+1, n \rrbracket$  do  
3:      $D[s, s'] = \text{SW}(s, s')$   
4:   end for  
5: end for
```

Construit par blocs

$$\begin{pmatrix} B_{11} & \dots & B_{1n} \\ \vdots & & \vdots \\ B_{m1} & \dots & B_{mn} \end{pmatrix} \quad \& \quad \forall i, j, \quad B_{ij} \longrightarrow \text{SW}(B_{ij})$$

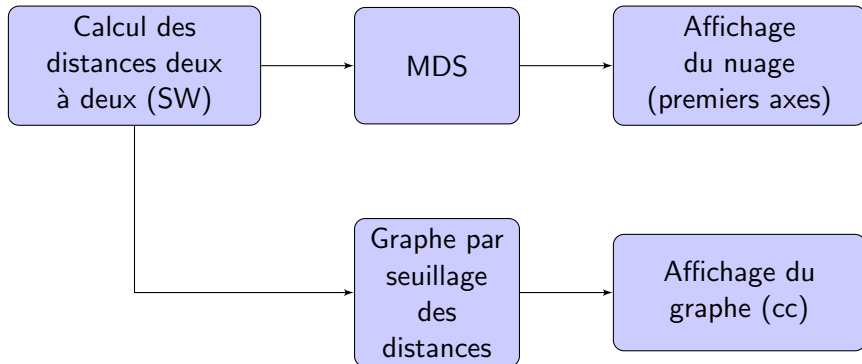
- programme en C, parallélisation en MPI
- passe parfaitement à l'échelle sur une machine hyperparallèle (Turing, Blue Gene Q, $16\,384 = 2^{14}$ cœurs)

Transfert pour post-traitement



Multidimensional Scaling

Tout commence par une distance ...



Concrètement ...

Multidimensional Scaling

Etant donné $D_{ij} = d(i, j)$
 $r < n$
trouver $\mathcal{X} = (x_1, \dots, x_n)$
tel que $x_i \in \mathbb{R}^r$
 $\|x_i - x_j\| \simeq d(i, j)$

Graphe

Etant donné $D_{ij} = d(i, j)$
 $\alpha > 0$
Construire $G = (V, E)$
tel que $V = \{1, n\}$
 $i \sim j \Leftrightarrow d(i, j) \leq \alpha$

Algorithme de la MDS: cubique

Algorithm 2 pseudocode pour la MDS

- 1: **input:** $D, r, D[i, j] = d_{ij}$
- 2: build Gram matrix G from D : ($g_{ij} = \langle x_i, x_j \rangle$)
- 3: **SVD of G :** $G = U\Sigma V^T$
- 4: coordinates: $X = U\Sigma^{1/2}$
- 5: **return** X

Projection aléatoire

Lemme de Johnson-Lindenstrauss (1984)

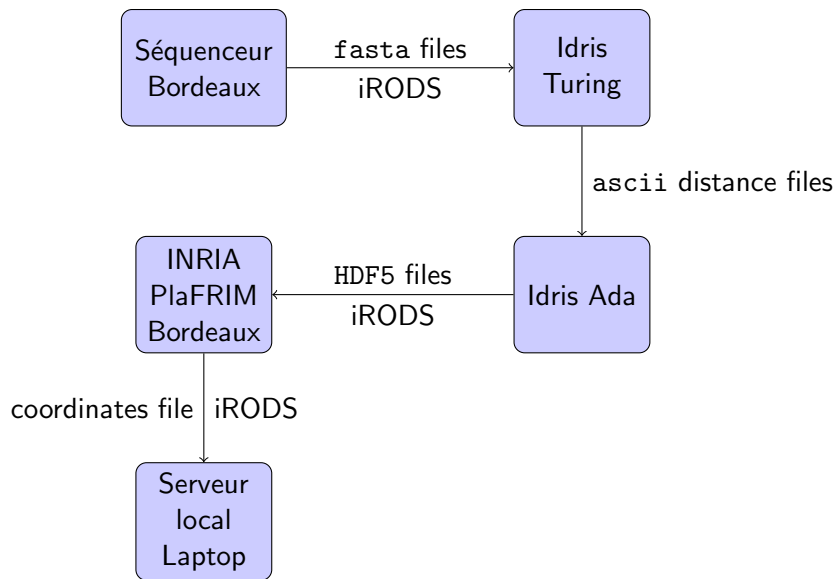
- Donnés: $n \in \mathbb{N}$, $\epsilon \in [0, 1]$, $k \geq \frac{8 \text{Log } n}{\epsilon^2}$
- Alors, $\forall X = [x_1, \dots, x_n]$, $x_i \in \mathbb{R}^n$
- $\exists \mathbb{R}^n \xrightarrow{f} \mathbb{R}^k$
- tel que $\forall x, y \in X$

$$(1 - \epsilon)\|x - y\|^2 \leq \|f(x) - f(y)\|^2 \leq (1 + \epsilon)\|x - y\|^2$$



Méthode de projection aléatoire

- voir Halko & al., SIAM Review, 2011
- si ϵ fixé, se ramène à une SVD dans la dimension k
- avec indicateurs de contrôle de la qualité



Pont entre l'apprentissage et l'inférence statistique (Tibshirani)

Machine learning

weights

learning

supervised learning

unsupervised learning

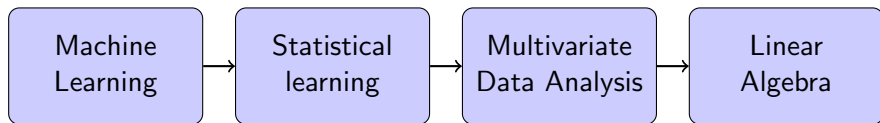
Statistics

parameters

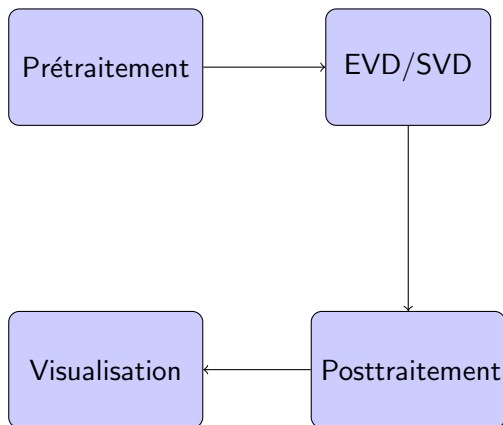
fitting

regression/classification

density estimation, clustering



L'essentiel des méthodes de réduction de la dimension



Quelques exemples

ACP

$$(U, \Sigma, V) = \text{SVD}(A)$$

$$\Lambda = \Sigma^2$$

$$Y = U\Lambda$$

MDS

$$G = \text{GRAM}(D)$$

$$(U, \Sigma, U) = \text{SVD}(G)$$

$$Y = U\Sigma^{1/2}$$

Canonical Analysis

$$T = A'B$$

$$D_A = (A'A)^{-1}, \quad D_B = (B'B)^{-1}$$

$$M_{A,B} = D_A^{1/2} T D_B^{1/2}$$

$$(W_B, \Psi) = \text{EVD}(M_{A,B})$$

$$W_A = M W_B$$

$$U_A = D_A^{1/2} W_A, \quad U_B = D_B^{1/2} W_B$$

$$Y_A = A U_A, \quad Y_B = B U_B$$

AFC

$$r = A \mathbf{1}_m$$

$$c = A^T \mathbf{1}_m$$

$$M = \text{diag}(1/\sqrt{r_i})$$

$$Q = \text{diag}(1/\sqrt{c_j})$$

$$A_{M,Q} = M(A - r \otimes c)$$

$$(Z, X, \Lambda) = \text{PCA}(A_{M,Q})$$

...

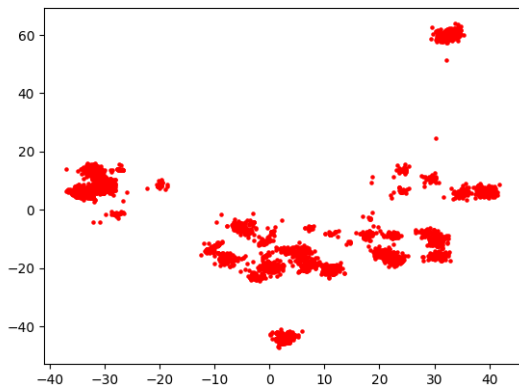
...

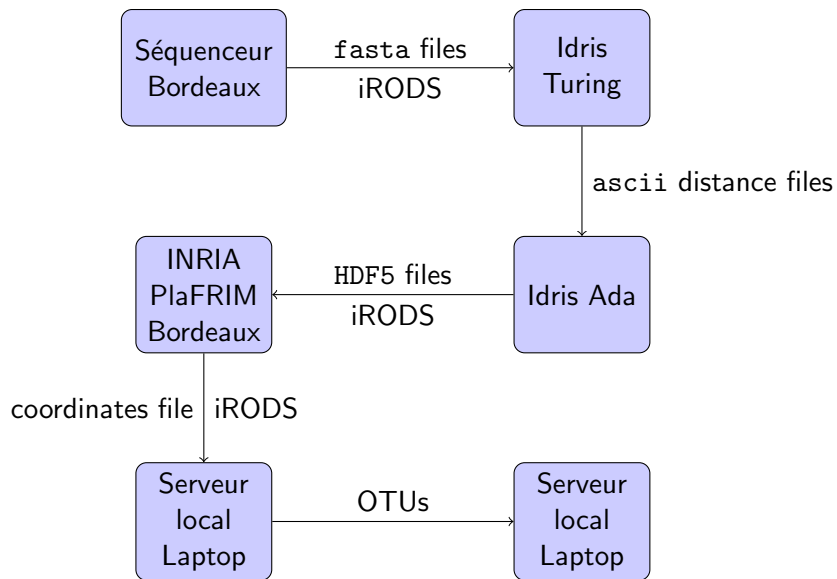


Un projet de librairie ...

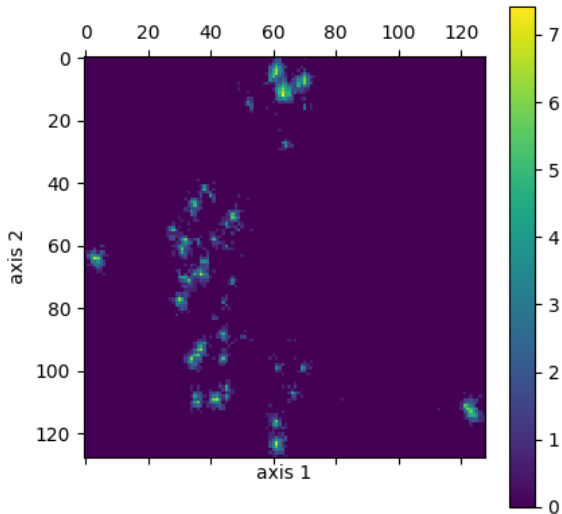
- englobant l'essentiel des méthodes dites "multivariate Analysis"
- basée sur un noyau de SVD par projection aléatoire
- avec vectorisation et optimisation des phases de prétraitement
- écrite en python/numpy, julia, C++
- actuellement en mémoire partagée sur un nœud
- en cours de développement vers une mémoire distribuée (MPI)
- avec la pile logicielle Inria Bordeaux (deux projets démarrent en 2019)

Nuage de points

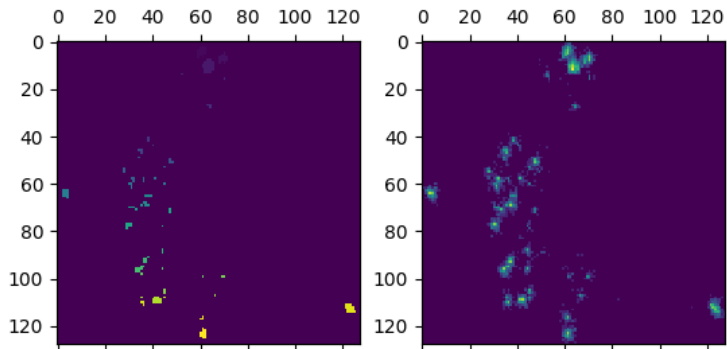




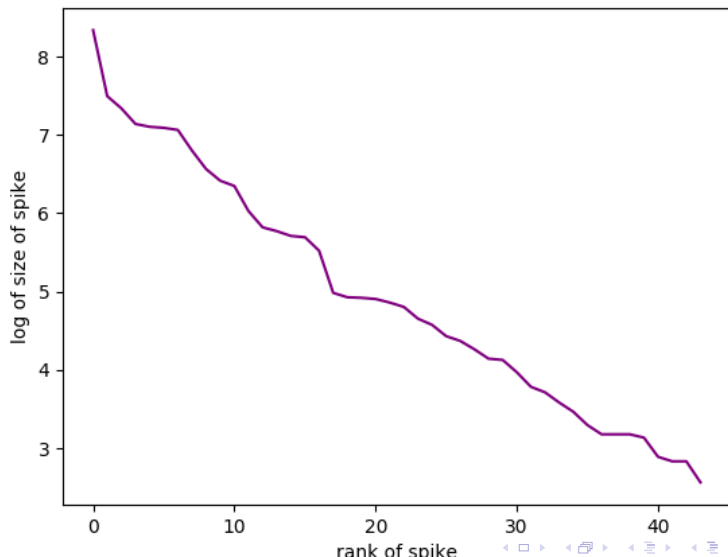
Nuage de points, en densité



Nuage de points, spikes



Distribution rang/taille

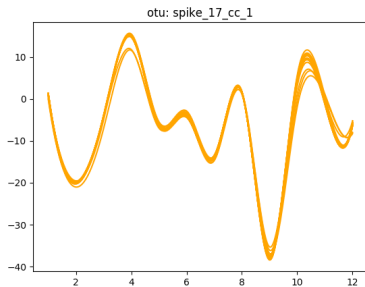
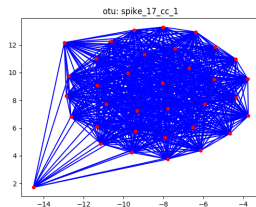
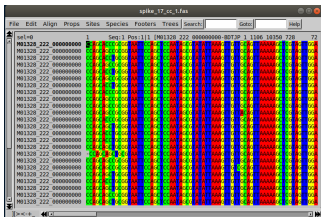


Construction des OTUs

Algorithm 3 pseudocode to build OTUs from MDS coordinates

- 1: **input** Coordinate file, loaded as X , distance file, loaded as D , threshold α , minimum size of an OTU m
 - 2: build spikes as sets of pixels
 - 3: **for** each spike **do**
 - 4: list all sequences projected on the spike
 - 5: extract the pairwise distances for this spike
 - 6: build the threshold graph
 - 7: build the connected components for this graph
 - 8: **if** the number of sequences $\geq m$ **then**
 - 9: add the cc in the list of OTUs
 - 10: **end if**
 - 11: **end for**
 - 12: **return** list of OTUs (sequences, distances, coordinates)
-

Un exemple d'OTU



Prochaines étapes ?

- Optimiser chaque maillon
- GPU ?
- passer à l'échelle en développant Diodon en mémoire/calcul distribué (pile logicielle Inria Bordeaux)
- Construire des bases de référence d'OTUs (intercalibrations des OTUs entre projets et échantillons)

Remerciements

- Les différentes infrastructures de calcul
 - IDRIS (projet DARI Biodiversiton 2015-2016-2018)
 - PlaFRIM (INRIA Bordeaux Sud-Ouest)et pour les deux le support aux utilisateurs
- Le projet européen EOSCpilot (2017-2018)
 - le WP6 "Interopérabilité" piloté par l'IN2P3
 - le "tesbed" pico2 sur l'interopérabilité
- Le COST DNAqua.Net (groupe de travail sur le traitement des données)
- Le réseau R-Syst (taxonomie moléculaire, très nombreux échanges)
- Les collègues biologistes ayant partagé leurs données pour la mise au point de ces outils (arbres de Guyane, diatomées du lac Léman et de Scandinavie, champignons phytopathogènes, ...) notamment via le projet Malabar (Bassin d'Arcachon)
- Last but not least: la plateforme PGTB (Franck Salin & Emilie Chancerel) de Pierroton