

Sequence analysis

MNHN-Tree-Tools: a toolbox for tree inference using multi-scale clustering of a set of sequences

Thomas Haschka^{1,*}, Loic Ponger¹, Christophe Escudé¹ and Julien Mozziconacci ^{1,2,*}

¹Muséum National d'Histoire Naturelle, Structure et Instabilité des Génomes, UMR7196, Paris 75231, France and ²Institut Universitaire de 1, rue Descartes 75005 Paris, France

*To whom correspondence should be addressed.

Associate Editor: Teresa Przytycka

Received on December 23, 2020; revised on May 12, 2021; editorial decision on May 31, 2021; accepted on June 7, 2021

Abstract

Summary: Genomic sequences are widely used to infer the evolutionary history of a given group of individuals. Many methods have been developed for sequence clustering and tree building. In the early days of genome sequencing, these were often limited to hundreds of sequences but due to the surge of high throughput sequencing, it is now common to have millions of sampled sequences at hand. We introduce MNHN-Tree-Tools, a high performance set of algorithms that builds multi-scale, nested clusters of sequences found in a FASTA file. MNHN-Tree-Tools does not rely on multiple sequence alignment and can thus be used on large datasets to infer a sequence tree. Herein, we outline two applications: a human alpha-satellite repeats classification and a tree of life derivation from 16S/18S rDNA sequences.

Availability and implementation: Open source with a Zlib License via the Git protocol: <https://gitlab.in2p3.fr/mnhn-tools/mnhn-tree-tools>.

Manual: A detailed users guide and tutorial: <https://gitlab.in2p3.fr/mnhn-tools/mnhn-tree-tools-manual/-/raw/master/manual.pdf>.

Website and FAQ: <http://treetools.haschka.net>.

Contact: julien.mozziconacci@mnhn.fr or thomas.haschka@mnhn.fr

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Sequences are slowly diverging in the course of evolution. The similarity between genomic loci, as for instance specific gene sequences, can in principle be used to infer the evolutionary relationship between individuals. Clustering methods are often used to group sequences together into species, genus, families, orders, class, phylum's, kingdoms and domains. Different experimental methods, such as DNA barcoding (DeSalle and Goldstein, 2019; Hajibabaei *et al.*, 2007), are used to determine the set of sequences to be clustered. Sequences are then often curated and gathered into large databases (McDonald *et al.*, 2012; Munoz *et al.*, 2011). With the recent advances in DNA high-throughput sequencing (Goodwin *et al.*, 2016), specimen collections and storage capacities, it is now common to deal with datasets with millions of entries. Several computational approaches have been developed to keep up with the size of these datasets (Mahé *et al.*, 2015; Rognes *et al.*, 2016) but they all provide clusters rather than trees. We propose here a new and fast method that performs a multiple alignment free, multi-scale clustering of a set of sequences found in a FASTA (Lipman and Pearson, 1985) file, leveraging DBSCAN a density-based algorithm for

discovering clusters in large spatial databases with noise (Ester *et al.*, 1996). Nested clusters are then identified to build a tree.

Briefly sketched, the DBSCAN algorithm is a two parameter algorithm requiring a radius ϵ and a minimum number of objects *minpts*, in our case sequences, to be found within this radius. As such, this algorithm finds density $\rho = \frac{n_{\text{minpts}}}{V(\epsilon)}$ connected regions, i.e. clusters with a density $> \rho(\text{minpts}, \epsilon)$ (Ester *et al.*, 1996). The use of DBSCAN has been proposed by others as a guide to phylogenetic inference (Mahapatro *et al.*, 2012; Ruzgar and Erciyes, 2012). The novelty introduced by our multi-scale approach is that we perform the DBSCAN algorithm at various densities, and use these layered results to infer a clustering tree that can further be used as a guide for phylogenetics. Clustering for different ϵ values allows us to find dense sequence clusters embedded into diffuse clusters (Fig. 1a). We can then build a tree of density connected clusters by successive DBSCAN runs with increased ϵ parameters and cluster comparison as outlined in Figure 1b and in the supplementary document. The DBSCAN algorithm was chosen over newer density based methods as, contrary to OPTICS (Ankerst *et al.*, 1999), DBSCAN allows us to control the density of the clusters found and thus allows us to precisely build trees from layers of specific densities. Further DBSCAN

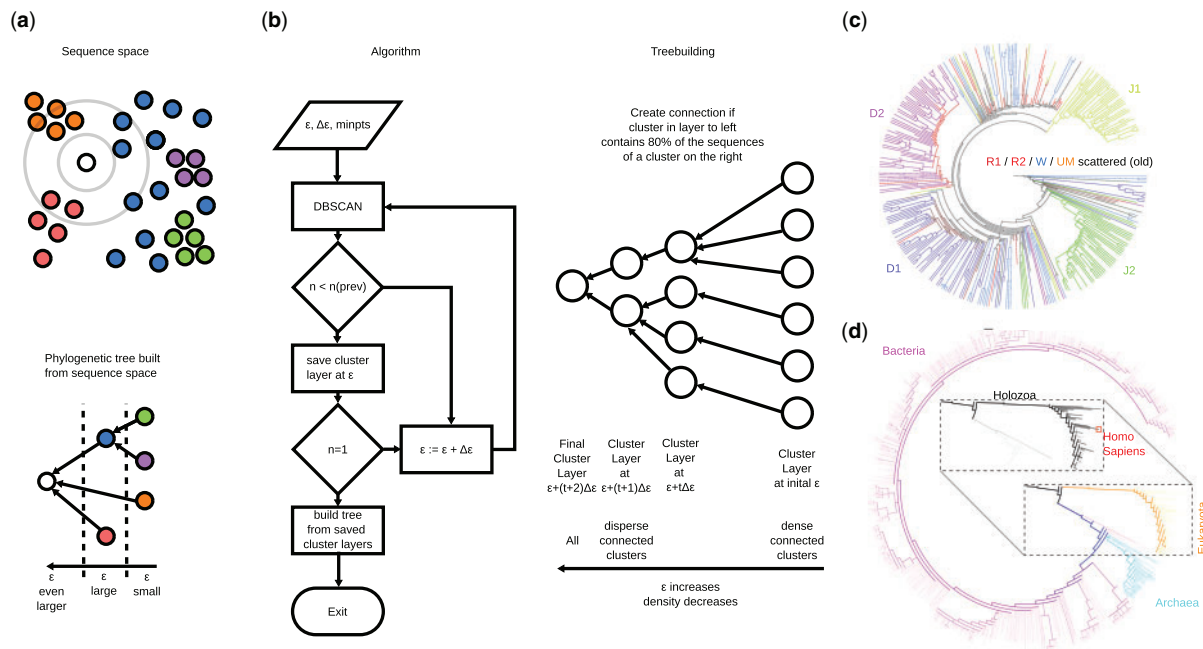


Fig. 1. Overview of MNHN-Tree-Tools. (a) Closely related sequences form dense clusters (in purple and green). These are embedded into a less dense cluster (in blue). The DBSCAN algorithm applied at various radius values (ϵ), can identify these nested clusters. A tree of the identified clusters can then be built. (b) Detailed computational workflow. (c) Tree build with human Alpha-satellites sequences. Colors correspond to the family annotations in the original dataset (Uralsky et al., 2019). (d) The tree of life built from 16S/18S RNA sequences. Bacteria, Archaea and Eukaryota are highlighted, with the color intensity corresponding to a logarithmic gradient of the number of sequences in the tree branches. A zoomed representation of Holozoa, clearly outlined as a subclass of Eukaryota, shows the Homo Sapiens branch

features a reduced algorithmic complexity and has a runtime advantage over algorithms, such as SUBCLU (Kailing et al., 2004). MNHN-Tree-Tools contains the utilities to cluster sequences using two different distance measures:

- the L_2 -norm operating on a principal component analysis (PCA) based subspace projection of the k-mer sequences representations (Chatterji et al., 2008).
- The Smith–Waterman distance (Smith and Waterman, 1981), which features parametric penalties for both substitutions and insertions.

In comparing the k-mer/PCA based approach to the use of the Smith–Waterman distance, we show that a k-mer/PCA based distance can yield better clusters and trees due to the inherent feature selection of the PCA (see supplementary document). The Smith–Waterman distance computation was implemented in OpenCL (Stone et al., 2010) allowing for execution on graphics processing units (GPUs). We also used the message passing interface (MPI) library (Forum, 1994) to distribute the workloads across different high performance computing cluster nodes.

2 Description of MNHN-Tree-Tools

MNHN-Tree-Tools is a modular suite of command line tools written in the C language. In this section, we outline the core utilities, which lead to a multilayered clustering with clusters organized into a tree.

Input data: MNHN-Tree-Tools uses as input a FASTA file format that gathers sequences that do not need to be aligned. Typical lengths can vary from 100 to 10000 bp, with length variations up to 10% within samples, but are ultimately only limited by k-mer length or PCA information retention.

fasta2kmer A utility to transform FASTA files into a k-mer representation.

kmer2pca Computes projections of a k-mer representation onto its first principal components.

adaptive_clustering_PCA Performs clustering at different densities, with the following variable input parameters: initial ϵ , $\Delta\epsilon$ and minpts . c.f. (Fig. 1b).

split_sets_to_newick Generates a Newick tree from the clusters obtained.

3 Performance and accuracy

We evaluated the accuracy of the algorithm presented herein in three different ways: at first, we compared our results to trees that were annotated by experts and as such provided us with valuable ground truth (see the case studies section below). We also compared the algorithm to partitions found by the SWARM2 tool (Mahé et al., 2015). Complementing these experiments we used MNHN-Tree-Tools on simulated datasets. The comparison to ground truth trees shows that our algorithm is able to find known partitions (see tables provided in the supplementary document). A comparison to SWARM2 (Mahé et al., 2015) clearly shows that our tree based approach yields, as we search for clusters at different densities, a sweet spot where the found partitions are in close correspondence to those annotated by experts. Classical partitioning tools, such as SWARM2 (Mahé et al., 2015), on the other hand, yield a single partition that may not correspond, for the application presented herein, to the expected results. We refer the reader to our Supplementary for further details on the accuracy and performance of MNHN-Tree-Tools where the outlined experiences are discussed in detail.

4 Case studies

Human alpha-satellites classification: Sequences were retrieved from (Uralsky et al., 2019). Our algorithm reconstructed a tree (Fig. 1c) from these 426 106 sequences which was colored according to their family annotation.

The tree of life—the SILVA dataset: Ribosomal RNA sequences from diverse species (2 225 272) were downloaded from (Munoz et al., 2011). Our algorithm was used to reconstruct a tree of life based on these sequences (Fig. 1d).

For these two applications, the run time for one clustering step ranges from 5 min (426 106 seq.) to 173 min (2 225 272 seq.) on a single Intel(R)i7-4771 3.50 GHz core. The clustering for different epsilon values can easily be run in parallel on several cores.

Acknowledgements

The authors thank the ‘Maison de la Simulation de Champagne Ardenne’ for allowing us to develop our algorithm on the ROMEO supercomputer and the members of the analysis hub of the UMS 2700 for stimulating discussions.

Funding

The project was funded by the Museum National d’Histoire Naturelle (MNHN) and the Institut Universitaire de France (IUF).

Conflict of Interest: none declared.

References

- Ankerst, M. *et al.* (1999) *Optics: Ordering Points to Identify the Clustering Structure*. ACM Press, New York, NY, USA. 49–60.
- Chatterji, S. *et al.* (2008) Compostbin: A DNA composition-based algorithm for binning environmental shotgun reads. In: *Annual International Conference on Research in Computational Molecular Biology*, p.17–28. Springer.
- DeSalle, R. and Goldstein, P. (2019) Review and interpretation of trends in DNA barcoding. *Front. Ecol. Evol.*, **7**, 302.
- Ester, M. *et al.* (1996) A density-based algorithm for discovering clusters a density-based algorithm for discovering clusters in large spatial databases with noise. In: *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, KDD’96, p.226–231. AAAI Press.
- Forum, M.P. (1994) MPI: a message-passing interface standard. *Technical report*, USA. <https://www.mpi-forum.org/>
- Goodwin, S. *et al.* (2016) Coming of age: ten years of next-generation sequencing technologies. *Nat. Rev. Genet.*, **17**, 333–351.
- Hajibabaei, M. *et al.* (2007) DNA barcoding: how it complements taxonomy, molecular phylogenetics and population genetics. *Trends Genet.*, **23**, 167–172.
- Kailing, K. *et al.* (2004) Density-connected subspace clustering for high-dimensional data. In: Karin Kailing, Hans-Peter Kriegel and Peer Kröger (eds.) *Proc. 4th SIAM Int. Conf. on Data Mining*, pp. 246–257, Lake Buena Vista, FL, 2004, p.246–257.
- Lipman, D. and Pearson, W. (1985) Rapid and sensitive protein similarity searches. *Science*, **227**, 1435–1441.
- Mahapatro, G. *et al.* (2012) Phylogenetic tree construction for DNA sequences using clustering methods. *Proc. Eng.*, **38**, 1362–1366.
- Mahé, F. *et al.* (2015) Swarm v2: highly-scalable and high-resolution amplicon clustering. *PeerJ*, **3**, e1420.
- McDonald, D. *et al.* (2012) An improved greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J.*, **6**, 610–618.
- Munoz, R. *et al.* (2011) Release ltps104 of the all-species living tree. *Syst. Appl. Microbiol.*, **34**, 169–170.
- Rognes, T. *et al.* (2016) Vsearch: a versatile open source tool for metagenomics. *PeerJ*, **4**, e2584.
- Ruzgar, E. and Erciyes, K. (2012) Clustering based distributed phylogenetic tree construction. *Expert Syst. Appl.*, **39**, 89–98.
- Smith, T. and Waterman, M. (1981) Identification of common molecular subsequences. *J. Mol. Biol.*, **147**, 195–197.
- Stone, J.E. *et al.* (2010) Opencl: a parallel programming standard for heterogeneous computing systems. *Comput. Sci. Eng.*, **12**, 66–73.
- Uralsky, L. *et al.* (2019) Classification and monomer-by-monomer annotation dataset of suprachromosomal family 1 alpha satellite higher-order repeats in hg38 human genome assembly. *Data Brief*, **24**, 103708. AAAI Press