

Index distribué de k -mers permettant leur comptage et leur localisation

libGkArrays-MPI/gkampi

<https://gite.lirmm.fr/doccy/libGkArrays-MPI>



05 08 JUIL
Université de Rennes 1

Clément AGRET

<clement.agret@lirmm.fr>

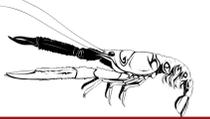
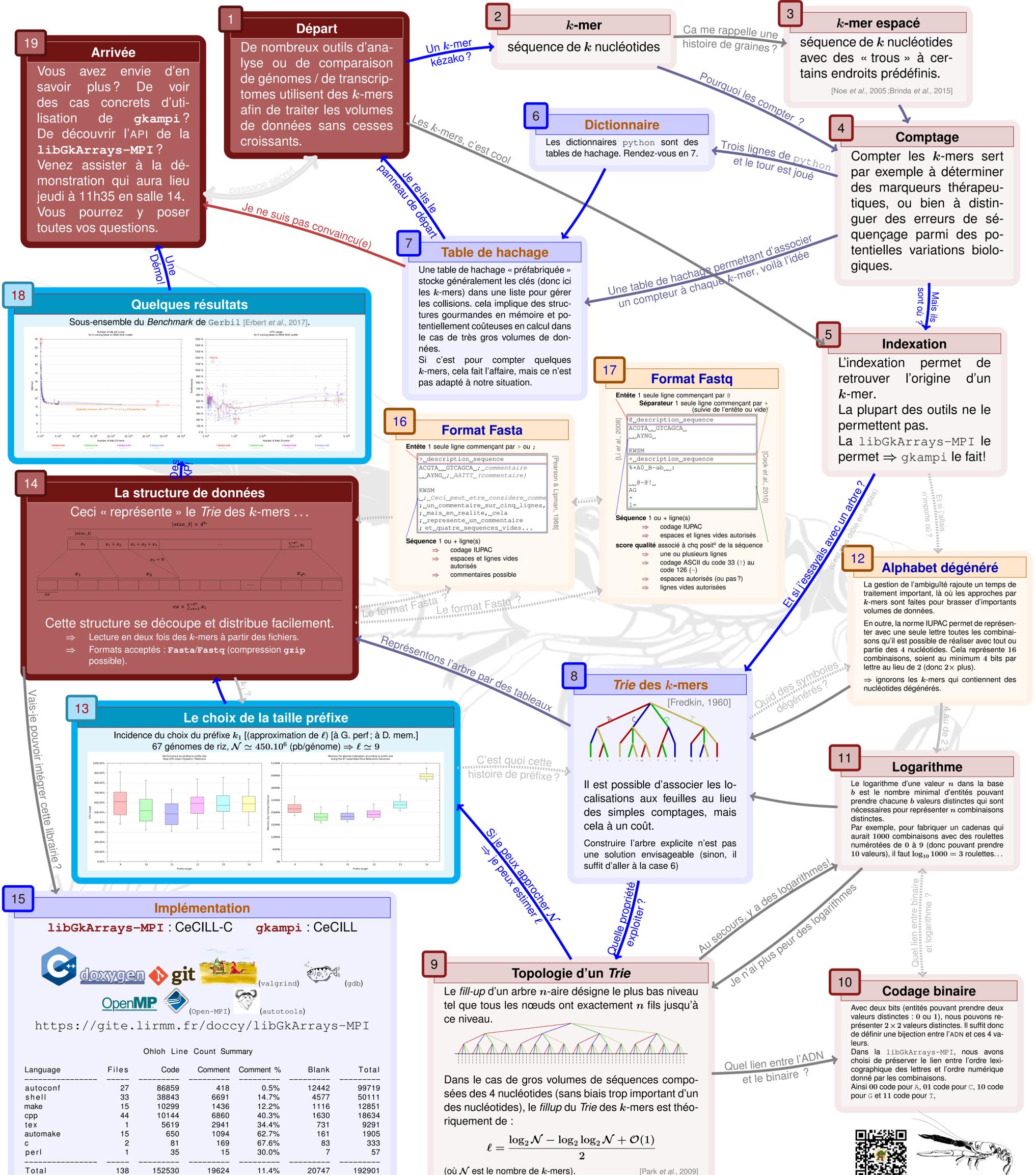
Annie CHATEAU

<annie.chateau@lirmm.fr>

Alban MANCHERON

<alban.mancheron@lirmm.fr>

Ceci est un poster dont vous êtes le héros. Parcourez les chemins au gré de vos envies, notez quelque part votre itinéraire, et découvrez jusqu'où il vous aura mené...



Index distribué de k -mers permettant leur comptage et leur localisation

libGkArrays-MPI/gkampi

Clément AGRET

<clement.agret@lirmm.fr>

Annie CHATEAU

<annie.chateau@lirmm.fr>

Alban MANCHERON

<alban.mancheron@lirmm.fr>

LIRMM, CNRS Université Montpellier 2 - CC 477
161, rue ADA 34095 Montpellier CEDEX 5

<https://gite.lirmm.fr/doccy/libGkArrays-MPI>

Le comptage de k -mers

Les progrès des 15 dernières années en matière de séquençage d'ADN et d'ARN, associés à la baisse de leurs coûts ont eu comme effet direct une production massive de données à analyser. Ce changement d'échelle du volume de données a induit l'apparition de nouvelles méthodologies et notamment celles basées sur le comptage des k -mers—fragments de longueur k —présents dans les séquences. Ces comptages peuvent être utilisés de différentes manières. Par exemple, rechercher des marqueurs spécifiques de certaines populations ou bien pour discriminer les erreurs de séquençage des variations biologiques. . .

Bien que le comptage de k -mers consiste à associer à une séquence de k nucléotides une valeur entière, il existe de multiples manières de structurer cette information et de l'interroger [1] et les choix algorithmiques et méthodologiques ont une incidence forte sur les performances et la fiabilité des méthodes. De nombreux outils ont été développés pour effectuer ces comptages, tels que Jellyfish [2], DSK [3], KMC3 [4], . . . Cependant, l'inconvénient commun à la plupart des méthodes existantes actuellement est qu'elles sont limitées par les capacités matérielles de la machine sur laquelle elles sont exécutées. Aussi pour compter les k -mers sur de très gros volumes de données, est-il nécessaire de disposer de machines surpuissantes ou d'adapter les méthodes existantes afin de distribuer les calculs (stratégies *MapReduce* [5]). Nous avons développé une méthode originale permettant de distribuer le calcul sur plusieurs machines, repoussant *de facto* les limitations de ces autres outils.

Comptage et indexation massivement parallélisés de k -mers

Nous avons développé une librairie en C++ (intitulée `libGkArrays-MPI` et distribuée sous la licence libre `CeCILL-C`), exploitant le parallélisme léger (*multithreading*) mais également le calcul distribué, permettant de compter les k -mers des séquences décrites dans un ou plusieurs fichiers (*fasta*, *fastq*, compressés ou non). Outre le simple comptage, cette librairie permet également de les indexer (donc de pouvoir retrouver leurs séquences d'origine). Sur la base de cette librairie, nous avons également développé un outil (intitulé `gkampi` et distribué sous licence libre `CeCILL`) pouvant s'exécuter sur une simple machine comme sur un *cluster* de calcul.

L'outil `gkampi` et la librairie `libGkArrays-MPI` permettent également de compter/indexer des k -mers espacés [6], proposent les mêmes fonctionnalités que les outils standards (*Jellyfish*, *KMC*, . . .) et sont documentés. Leur installation est conforme aux standards des `GNU autotools` et le code respecte strictement la norme ISO 2011 du C++.

Remerciements

Ce travail a été en partie financé par l'Institut de Biologie Computationnelle de Montpellier, le projet *GenomeHarvest (Agropolis foundation)* et le projet *Coalab* (région Languedoc-Roussillon). Nous tenons également à remercier les rapporteurs pour leurs commentaires.

Bibliographie

- [1] Hilde VINJE, Kristian Hovde LILAND, Trygve ALMØY et Lars SNIPEN : Comparing k -mer based methods for improved classification of 16S sequences. *BMC Bioinformatics*, 16:205, 2015.
- [2] Guillaume MARÇAIS et Carl KINGSFORD : A fast, lock-free approach for efficient parallel counting of occurrences of k -mers. *Bioinformatics*, 27(6):764–770, mars 2011.
- [3] Guillaume RIZK, Dominique LAVENIER et Rayan CHIKHI : DSK : k -mer counting with very low memory usage. *Bioinformatics*, 29(5):652–653, 2013.
- [4] Marek KOKOT, Maciej DŁUGOSZ et Sebastian DEOROWICZ : KMC 3 : counting and manipulating k -mer statistics. *Bioinformatics*, 33(17):2759–2761, 2017.
- [5] Tao GAO, Yanfei GUO, Yanjie WEI, Bingqiang WANG, Yutong LU, Pietro CICOTTI, Pavan BALAJI et Michela TAUFER : Bloomfish : A Highly Scalable Distributed k -mer Counting Framework. *In 2017 IEEE 23rd International Conference on Parallel and Distributed Systems (ICPADS)*, pages 170–179, décembre 2017.
- [6] Karel BŘINDA, Maciej SYKULSKI et Gregory KUCHEROV : Spaced seeds improve k -mer-based metagenomic classification. *Bioinformatics*, 31(22):3584–3592, novembre 2015.

