

Research Proposal: *Decentralized indices for genomic data*

Luiz Irber¹

¹University of California, Davis

11 April 2019

Introduction

Since the announcement of the sequencing of the human genome in 2001, technological breakthroughs in sequencing lowered the cost per megabase from thousands of dollars to fractions of cents, opening the way to new classes of experiments and deeper exploration of biological fields, from studying diseases to crop improvements. The data generation rate increased with lower costs and public data repositories for genomic data are reaching multiple petabytes of data, leading to discussions on how to search them and allow exploratory analysis and hypothesis generation from these resources.

Bioinformatics methods that once were viable for the available data don't scale to this new situation, creating new challenges both for method developers as well as scientists doing data analysis. Analysis that was once possible in a workstation now demand clusters and increasingly access to data in many different locations. Storing and transferring this data is an additional burden, especially since most resources are maintained in centralized locations and made available in ways that don't explore local caches in the network, leading to congestion and low transfer speeds. The centralization aspect also lead to resources disappearing when funding for maintaining it ends, since most of this resources were not developed with resilience and federation in mind, making them harder or impossible to deploy without access to the original location.

This proposal describe methods for working with large scale sequencing datasets and databases, considering data acquisition, analysis and distribution solutions for problems encountered by biologists during their experiments. Special focus is given for searching for similar datasets in large sequence databases, and taxonomic profiling of metagenomes.

Background

Problem Description

Searching biological databases

The quintessential bioinformatics tool is BLAST [2], a method that performs local alignment between sequences and supports building a database for searching multiple datasets. NCBI offers a web portal to search some public databases, but it is a subset of all the data available. For very large databases like the Sequence Read Archive [16] it is only possible to search in specific experiments, especially considering that even a subset like the microbial datasets contains more than 800 thousand experiments, making the full resource opaque for discovery and exploratory analysis.

For tools that support building local databases, the main issue is downloading the required data. Going back to the microbial subset of the SRA, the 800 thousand experiments need more than 400 TB of data transferring, and storing it is also a non-trivial problem.

Taxonomic profiling

Metagenomics is the study of the community from an environmental sample. This has applications both in many contexts, with special interest in health and clinical areas, where it is more commonly referred as the microbiome (the microorganisms composing the microbial community inside a person).

Taxonomic profiling is the characterization of the organisms present in a metagenomics sample, including their relative abundance. Methods for taxonomic profiling can be divided in three groups: marker genes, alignment and k -mer composition.

Marker genes methods use specific genes (like the 16S ribosomal RNA gene for the microbial case) or a combination of them [19] to classify reads and aggregate it to create a summary of the relative abundance in the sample. Drawbacks include throwing out information from reads that don't align to the marker gene, and by focusing on one or a small set of genes it discards relevant information from all the other genes. This approach also doesn't work across long evolutionary distances, since marker genes diverged too much to be comparable.

Alignment methods are not limited to specific genes [14], but are computationally more expensive due to use of larger sequence reference databases and more traditional local alignment methods (like BLAST). Finally, k -mer composition methods work by preparing a database with assignments from a k -mer to a specific taxon in a phylogenetic tree [25], and then using all k -mer assignments for a read to decide from what taxon the read came. The k -mer composition of a genomic sequence is a set of all k -length substrings, computed as a sliding window over the genomic sequence. For a genomic dataset the k -mer composition is the set of all k -length substrings for all sequences in it. Further aggregation is also possible by combining information from each read to do the taxonomic profiling of the full sample. Profiling is very fast compared to other methods, but building the database involves more preprocessing. Another drawback is k -mer composition methods do exact matching, while other methods (being based on alignment) can account for mutation and other relevant biological processes.

Data sketches

A data sketch is a representative proxy for the original data focused on queries for specific properties. It can also be viewed as a probabilistic data structure (in contrast with deterministic data structures), since it uses hashing techniques to provide statistical guarantees on the precision of the answer for a query. This allows a memory/accuracy trade-off: using more memory leads to more accurate results, but in memory-constrained situations it still bound results to an expected error rate.

MinHash sketch: similarity and containment

The MinHash sketch [6] was developed at Altavista in the context of document clustering and deduplication. It provides an estimate of the Jaccard similarity $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ (called **resemblance** in the original article) and the **containment** of two documents $C(A, B) = \frac{|A \cap B|}{|A|}$, estimating how much of document A is contained in document B . These estimates depend on two processes: Converting documents to sets ("Shingling"), and transforming large sets into short signatures, while preserving similarity ("Min-Hashing"). In the original use case the w -shingling Ω of a document D is defined as the set of all continuous subsequence of w words contained in D . *Min-hashing* is the process of creating $W = \{h(x) | \forall x \in \Omega\}$, where $h(x)$ is a uniform hash function, and then either

- keeping the n smallest hash values as a representative sketch of the original document (**MIN** $_n(W)$)
- keeping elements that are 0 mod m (**MOD** $_m(W)$).

MIN $_n(W)$ if fixed-sized (length n) and supports similarity estimation, but doesn't support containment. **MOD** $_m(W)$ supports both similarity and containment, with the drawback of not being fixed-sized anymore, growing with the complexity of the document.

Mash [17] was the first implementation of MinHash in the genomic context, relating the w -shingling of a document to the k -mer composition of a genomic sequence, and using the **MIN** $_n(W)$ fixed-sized formula-

tion for the signature. Mash needs extra information (the genome size for the organism being queried) to account for genomic complexity in datasets. This extra information is required because using a fixed-size MinHash leads to different degrees of accuracy when comparing across highly-diverged organisms (bacteria to animals, for example), and it is even more extreme when taking more complex datasets into account (like metagenomes).

Bloom Filter: Set membership

The Bloom Filter [5] allows set membership queries. Inserting an element in a Bloom Filter is a two-step process: first use multiple hash functions on the element, and then for each bit set in the hashed value update the same position in the bit array. Querying an element involves calculating the multiple hash values for the element and then checking if the bits are set in the bit array. This guarantees that a false negative is impossible (if the element was inserted, a bit would be set), but can report false positives if there are collisions in the hash values.

khmer [8] implements a variation where longer (and multiple) bit arrays are used, and hash functions are derived from a composed hashing strategy $h_i(x) = h(x) \bmod p_i$, where each bit array has a distinct length p_i (and p is a prime number), with $h(x)$ being a more CPU-intensive hash function.

Bloom Filters are used extensively in bioinformatics, including lossy representation of assembly graphs [18] and as a filtering step in processing pipelines [17].

HyperLogLog: Cardinality estimation

The HyperLogLog sketch [11] estimates the number of unique elements in a dataset. It is designed to lower the variability of a more basic estimator: given a run of ρ zeros in a binary sequence, it estimates the cardinality of the dataset to be 2^ρ . It achieves this by splitting the binary sequence in two: the lower bits define an index for m registers, and each register contain the longest run of zeros seen for that index. The HyperLogLog estimator $E(D)$ is an harmonic mean of the registers M , with a correction factor α_m for the number of registers:

$$E(D) = \alpha_m m^2 \left(\sum_{j=0}^{m-1} 2^{-M[j]} \right)^{-1}$$

More recent methods [13] improve the cardinality estimator by further refining the estimate based on empirical data.

In genomic contexts, the khmer [8] implementation of HyperLogLog [15] uses the k -mer composition of a genomic dataset and the *murmurhash3* hash function, together with the improved estimator from [13]. Dashing [3] is a recent method that supports both similarity and cardinality estimation using HyperLogLog, based on a better estimator for union and intersection of HyperLogLog sketches by [10].

Hierarchical index

Searching for matches in large collection of datasets is not viable when hundreds of thousands of them are available, especially if they are partitioned and not all present at the same place.

Bloofi [7] is a hierarchical index structure that extends the Bloom Filter basic query to collections of Bloom Filters. Instead of calculating the union of all Bloom Filters in the collection (which would allow answering if an element is present in any of them) it defines a tree structure where the original Bloom Filters are leaves, and internal nodes are the union of all the Bloom Filters in their subtrees. Searching is based on a breadth-first search, with pruning when no matches are found at an internal level. Bloofi can also be partitioned in a network, with network nodes containing a subtree of the original tree and only being accessed if the search requires it.

The Sequence Bloom Tree [20] adapts Bloofi for genomic contexts, rephrasing the problem as experiment discovery: given a query sequence Q and a threshold θ , which experiments contain at least θ of the original

query Q ? Experiments are encoded in Bloom Filters containing the k -mer composition of transcriptomes, and queries are transcripts.

Further developments focused on clustering similar datasets to prune search early [22] and developing more efficient representations for the internal nodes [21] [12] to use less storage space and memory.

Decentralized data access

Public genomic databases are traditionally maintained by governmental agencies, including international consortiums with public funding. The Sequence Read Archive [16] is mirrored in Europe (the European Nucleotide Archive) and Japan (the DNA Data Bank of Japan), but they are not connected: if one resource is down, there is no fallback to the other mirrors. They also support processed data formats which might not be available in the other mirrors, despite being derived from the same original data. Similar architectural decisions are common in other public genomic databases.

Despite some discussions [1] of alternatives in case funding or support for the SRA or similar large archives is cancelled, they usually don't take into account decentralized solutions for data access. IPFS [4] is one such system, based on the idea of locating resources in the network based on their content and not a specific location (like URLs). This allows querying peers for a specific object based on its content hash instead of depending on a specific URL, making it a good solution for mirrored data and local caches.

Aims

1. **Adapting genomic MinHash for containment and cardinality estimation.** The original MinHash article [6] defines two estimates for similarity resemblance and containment, with two variations of the MinHash sketch, one with fixed length, but only supporting resemblance, and another with variable length, supporting both resemblance and containment. Mash [17] uses the fixed length sketch, and defines a new distance metric called Mash distance to account for the size of the genomes being compared. The resemblance estimate works well for genomes of similar size, but when dealing with datasets of highly-diverged organisms or even more complex datasets (like metagenomes) the containment estimate is more appropriate and closer to the sort of problems that biologists need to solve.

In sourmash [23] I implemented a variable length MinHash sketch (with fallback to the fixed length case for compatibility with Mash) called the scaled MinHash. The scaled MinHash supports containment estimation, allowing new applications not available on previous methods.

Sketches are data structures planned for specific classes of queries, but it is possible to support additional use cases. For example, dashing [3] is a method that supports both similarity and cardinality estimation using HyperLogLog sketches [11], a data structure initially derived for supporting cardinality estimation and later extended to similarity cases [10]. Even so, it doesn't support containment estimation.

Scaled MinHash can be adapted to also support cardinality estimation. Due to how the scaled MinHash is constructed the same estimator from HyperLogLog can be used for it.

2. **Fast queries on many MinHash sketches using MinHash Bloom Trees** [20] introduces the Sequence Bloom Tree, an hierarchical index data structure for finding a query sequence in large dataset collections. It represents the k -mer composition of a dataset using a Bloom Filter [5], and the hierarchical aspect comes from organizing multiple dataset Bloom Filters into a tree structure, where the leaves are the dataset Bloom Filters and the internal nodes are Bloom Filters containing all k -mers below it. A query is evaluated by doing a breadth-first search of the tree, and truncating the search when the query is not present over a pre-determined threshold.

I adapted the Sequence Bloom Tree to use MinHash sketches as dataset representations, referred as MinHash Bloom Tree from now on to highlight how they are distinct. Since a MinHash is a subset of

the k -mer composition of a dataset, internal nodes are still Bloom Filters, but this time containing all the k -mers present in the MinHash.

On top of supporting the search method defined by the Sequence Bloom Tree (a breadth-first search with early truncating on a similarity or containment threshold), the MinHash Bloom Tree index also supports another search method called **gather**, a variation of Best-First search using containment estimation. **gather** can be used for doing taxonomic profiling of genomic datasets, finding all organisms present in a sample and how abundant they are. This is especially important when working with metagenomes, a dataset containing sequenced genomic data from an environmental sample (be it soil, ocean or human gut, for example) representing a community of organisms. **gather** does iterative Best-First searches, at each step removing matches from the original query, until there are no more hashes to search or a detection threshold is reached.

3. **Decentralized indices for genomic data.** The MinHash Bloom Tree can be viewed as a persistent data structure [9], since leaves never change once added and internal nodes only change if a new leaf is added to its subtree. This makes it a very good fit for content-addressable storage methods, and I'm exploring two different decentralized data storage systems (IPFS and dat) as ways of storing and interacting with MinHash Bloom Tree indices.

On top the data storage aspects, another important point is how researchers can interact with these indices (both querying and updating it) in a way that centralized storage is not essential (but operations are optimized if it is [24]).

This is important in the context of long term sustainability of projects, something often overlooked in bioinformatics systems. The initial implementation is centralized for simplicity and prototyping the user interaction, supporting data submission and querying, with a data pipeline based on the experience downloading 800k+ microbial datasets from NCBI SRA and also the preparation of indices for the IMG database from JGI (60k+ datasets).

Once this prototype is functional and we find what features are desirable, I plan to start decentralizing this process by defining updates in terms of CRDT operations, publishing a feed of these changes and moving to PubSub messaging between remote sites, allowing them to update their indices and keep them consistent with ours.

References

- [1] Closure of the NCBI SRA and implications for the long-term future of genomics data storage. 12:402. ISSN 1474-760X. doi: 10.1186/gb-2011-12-3-402. URL <http://dx.doi.org/10.1186/gb-2011-12-3-402>.
- [2] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. 215(3):403–410.
- [3] D. N. Baker and B. Langmead. Dashing: Fast and accurate genomic distances with HyperLogLog. page 501726. doi: 10.1101/501726. URL <https://www.biorxiv.org/content/early/2018/12/20/501726>.
- [4] J. Benet. IPFS - content addressed, versioned, p2p file system. URL <http://arxiv.org/abs/1407.3561>.
- [5] B. H. Bloom. Space/time trade-offs in hash coding with allowable errors. 13(7):422–426. ISSN 0001-0782. doi: 10.1145/362686.362692. URL <http://doi.acm.org/10.1145/362686.362692>.
- [6] A. Z. Broder. On the resemblance and containment of documents. In *Compression and Complexity of Sequences 1997. Proceedings*, pages 21–29. IEEE. URL <http://ieeexplore.ieee.org/abstract/document/666900/>.
- [7] A. Crainiceanu and D. Lemire. Bloofi: Multidimensional bloom filters. 54:311–324. ISSN 03064379. doi: 10.1016/j.is.2015.01.002. URL <http://arxiv.org/abs/1501.01941>.
- [8] M. R. Crusoe, H. F. Alameldin, S. Awad, E. Boucher, A. Caldwell, R. Cartwright, A. Charbonneau, B. Constantinides, G. Edverson, S. Fay, and others. The khmer software package: enabling efficient nucleotide sequence analysis. 4. URL <https://f1000research.com/articles/4-900/v1>.
- [9] J. R. Driscoll, N. Sarnak, D. D. Sleator, and R. E. Tarjan. Making data structures persistent. 38(1):86–124. ISSN 0022-0000. doi: 10.1016/0022-0000(89)90034-2. URL <http://www.sciencedirect.com/science/article/pii/0022000089900342>.
- [10] O. Ertl. New cardinality estimation algorithms for HyperLogLog sketches. URL <http://arxiv.org/abs/1702.01284>.
- [11] P. Flajolet, . Fusy, O. Gandouet, and F. Meunier. HyperLogLog: the analysis of a near-optimal cardinality estimation algorithm. (1). URL <http://www.dmtcs.org/dmtcs-ojs/index.php/proceedings/article/viewArticle/914>.
- [12] R. S. Harris and P. Medvedev. Improved representation of sequence bloom trees. page 501452. doi: 10.1101/501452. URL <https://www.biorxiv.org/content/early/2018/12/19/501452>.
- [13] S. Heule, M. Nunkesser, and A. Hall. HyperLogLog in practice: Algorithmic engineering of a state of the art cardinality estimation algorithm. In *Proceedings of the 16th International Conference on Extending Database Technology*, pages 683–692. ACM.
- [14] D. H. Huson, A. F. Auch, J. Qi, and S. C. Schuster. MEGAN analysis of metagenomic data. 17(3): 377–386.
- [15] L. C. Irber Junior and C. T. Brown. Efficient cardinality estimation for k-mers in large DNA sequencing data sets. page 056846. URL <http://www.biorxiv.org/content/early/2016/06/07/056846.abstract>.
- [16] Y. Kodama, M. Shumway, and R. Leinonen. The sequence read archive: explosive growth of sequencing data. 40:D54–D56.
- [17] B. D. Ondov, T. J. Treangen, P. Melsted, A. B. Mallonee, N. H. Bergman, S. Koren, and A. M. Phillippy. Mash: fast genome and metagenome distance estimation using MinHash. 17:132. ISSN 1474-760X. doi: 10.1186/s13059-016-0997-x. URL <http://dx.doi.org/10.1186/s13059-016-0997-x>.
- [18] J. Pell, A. Hintze, R. Canino-Koning, A. Howe, J. M. Tiedje, and C. T. Brown. Scaling metagenome sequence assembly with probabilistic de bruijn graphs. 109(33):13272–13277. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1121464109. URL <http://www.pnas.org/content/109/33/13272>.

- [19] N. Segata, L. Waldron, A. Ballarini, V. Narasimhan, O. Jousson, and C. Huttenhower. Metagenomic microbial community profiling using unique clade-specific marker genes. 9(8):811–814. ISSN 1548-7105. doi: 10.1038/nmeth.2066. URL <https://www.nature.com/articles/nmeth.2066>.
- [20] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. 34(3): 300–302, . ISSN 1087-0156. doi: 10.1038/nbt.3442. URL <https://www.nature.com/nbt/journal/v34/n3/full/nbt.3442.html>.
- [21] B. Solomon and C. Kingsford. Improved search of large transcriptomic sequencing databases using split sequence bloom trees. In *International Conference on Research in Computational Molecular Biology*, pages 257–271. Springer, .
- [22] C. Sun, R. S. Harris, R. Chikhi, and P. Medvedev. AllSome sequence bloom trees. In *Research in Computational Molecular Biology*, Lecture Notes in Computer Science, pages 272–286. Springer, Cham. ISBN 978-3-319-56969-7 978-3-319-56970-3. doi: 10.1007/978-3-319-56970-3_17. URL https://link.springer.com/chapter/10.1007/978-3-319-56970-3_17.
- [23] C. Titus Brown and L. Irber. sourmash: a library for MinHash sketching of DNA. 1(5). doi: 10.21105/joss.00027. URL <http://joss.theoj.org/papers/10.21105/joss.00027>.
- [24] J. N. Tsitsiklis and K. Xu. On the power of (even a little) centralization in distributed processing. 39 (1):121–132.
- [25] D. E. Wood and S. L. Salzberg. Kraken: ultrafast metagenomic sequence classification using exact alignments. 15(3):R46. ISSN 1474-760X. doi: 10.1186/gb-2014-15-3-r46. URL <https://doi.org/10.1186/gb-2014-15-3-r46>.