

Alignment-free Genomic Analysis via a Big Data Spark Platform

Umberto Ferraro Petrillo*[†] Francesco Palini* Giuseppe Cattaneo[‡] §
Raffaele Giancarlo[¶] §

Abstract

Motivation: Alignment-free distance and similarity functions (AF functions, for short) are a well established alternative to two and multiple sequence alignments for many genomic, metagenomic and epigenomic tasks. Due to data-intensive applications, the computation of AF functions is a Big Data problem, with the recent Literature indicating that the development of fast and scalable algorithms computing AF functions is a high-priority task. Somewhat surprisingly, despite the increasing popularity of Big Data technologies in Computational Biology, the development of a Big Data platform for those tasks has not been pursued, possibly due to its complexity.

Results: We fill this important gap by introducing FADE, the first extensible, efficient and scalable Spark platform for Alignment-free genomic analysis. It supports natively eighteen of the best performing AF functions coming out of a recent hallmark benchmarking study. FADE development and potential impact comprises novel aspects of interest. Namely, (a) a considerable effort of distributed algorithms, the most tangible result being a much faster execution time of reference methods like MASH and FSWM; (b) a software design that makes FADE user-friendly and easily extendable by Spark non-specialists; (c) its ability to support data- and compute-intensive tasks. About this, we provide a novel and much needed analysis of how informative and robust AF functions are, in terms of the statistical significance of their output. Our findings naturally extend the ones of the highly regarded benchmarking study, since the functions that can really be used are reduced to a handful of the eighteen included in FADE.

Contact: umberto.ferraro@uniroma1.it

1 Introduction

Alignment-free distance and similarity functions (AF functions, for short) have been introduced as an alternative to traditional alignment-based methods, e.g. [1, 27], in order to assess how similar each pair of sequences in a collection are to each other. By now, their use has been widely investigated for sequence analysis in genomics [31], metagenomics [3], and epigenomics [9, 10]. The pros/cons of AF functions with respect to their alignment counterparts is well presented in [31]. One of the key pro features highlighted in that study is that their implementations offer data scalability, opportunities that alignment methods lack. Taking into account the throughput of HTS technologies, another compelling case on how important is to design and implement fast and scalable algorithms for the computation of AF functions is presented in [23]. It is of interest here to notice that the mentioned two studies clearly indicate that the computation of AF functions is now a Big Data problem.

As such, it needs algorithmic solutions that use Big Data Technologies to grant efficiency and scalability as a function of the available computational resources and of the amount of data to process. Somewhat surprisingly, although those technologies are finding more and more use in Computational Biology [6, 22], and Cloud Storage and Computing is the future of genomic data [14], only a few studies in that area are

*Dipartimento di Scienze Statistiche, Università di Roma - La Sapienza, Rome, 00185, Italy

[†]To whom correspondence should be addressed.

[‡]Dipartimento di Informatica, Università di Salerno, Fisciano (SA), 84084, Italy

[§]Those two authors contributed equally to the research.

[¶]Dipartimento di Matematica ed Informatica, Università di Palermo, Palermo, 90133, Italy

available, concentrating on AF methods for informational and linguistic analysis of genomic sequences [5]. Yet, an effective Big Data platform supporting both the computation of AF functions and based on a pervasive Big Data framework such as Spark, would place AF sequence comparison at a peer with other important domains in Data Science (see, e.g., [13]) that have such platforms, as well as contribute to increase their usage in the Life Sciences.

Our first contribution is to address this acute need. Indeed, we propose FADE, the first extensible and scalable Spark platform for the effective computation of AF functions. In its starting configuration, it offers eighteen implementations of highly performing AF functions according to results in [31]. The basic pipeline of our platform consists of five (possibly optional) stages, as outlined in Figure 1. Its key features are as follows.

- **FADE is User-Friendly and Easily Extensible** FADE comes as a ready-to-use Spark application that can be easily executed without writing any line of code. If needed, its standard processing can be customized by just editing a provided reference configuration file, so as to choose from an included library which statistics to extract and which AF functions to evaluate. Some examples of configuration files are reported in Section 3.1 of the Supplementary Material. The user can add support for a target statistic or a target AF function not originally included in the library, by writing the corresponding code using one of the provided class templates. This part is outlined in Section 2.3.
- **FADE is Scalable and Efficient.** In order to assess FADE ability to profitably support the implementation of AF functions, we have compared two of the most prominent ones in the shared memory category, i.e., Mash [23] and FSWM [18], vs their respective implementations supported by FADE (denoted with the prefix FADE). Being those latter based on a distributed framework, they are able to take advantage of the much higher number of processing cores available, achieving performances that are much better than those of shared memory tools and that scale as a function of the processing units available. It is to be noted that, given the wide range of application scenarios for AF sequence comparison, varying from the examination of a large collection of very small reads, up to the comparison of a few huge genomic sequences, a flexible workload assignment is fundamental to grant scalability and speed. To this end, FADE provides, and we have experimented with, three workload choices. Transparently from the user, each of the three is associated to a transformation of the logical basic pipeline into a suitable “run time” software pipeline, which is then executed. The Partial Aggregation strategy results to be the most appropriate, with varying workloads. All the results, and the corresponding experiments, regarding this part are presented in Section 3.1.

Our second contribution, introduced in Section 2.5 and detailed in Section 3.2, shows the ability of FADE to tackle data and compute intensive tasks and it is of interest in its own right. Indeed, we use novel ideas, that boil down to very costly Monte Carlo Simulations, to gain insights into the properties of AF functions, going a step further in the direction indicated by the mentioned benchmarking study. The end result is new guidelines on the use of AF functions in day-to-day genomic analysis tasks.

To this end, we consider reliability and robustness of AF functions in terms of the statistical significance of the distance/similarity matrices they produce on benchmark datasets. In details, we consider a p-value obtained via a Monte Carlo simulation with the Null Hypothesis being that the values obtained by a given AF function on two biological sequences is no better than the value obtained on two random sequences. A value is significant when the Null Hypothesis is rejected. We account for repeated tests by applying Bonferroni Correction. Our key findings are the following.

- **A Novel Class of AF functions: Consistently Significant.** Across benchmark datasets from [31], only a handful of the eighteen AF functions we have considered provide distance/similarity matrices for which the Null Hypothesis is rejected for the vast majority of their entries. We refer to those functions as *consistently significant*, since they behave consistently well, statistically, irrespective of the dataset, and it is well known that statistical relevance is a good indication of biological relevance [7, 11, 19]. They are members of the D2 family [29] (D2 and D2*), and FSWM [18]. Since the D2 family has been

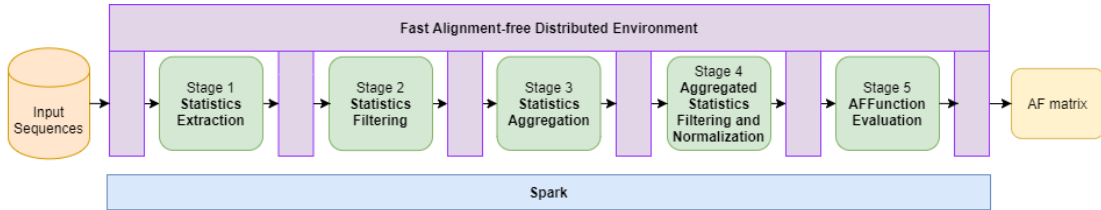


Figure 1: A layout of the logical architecture of the basic pipeline for the fast computation of AF functions.

studied extensively from the statistical point of view, our results confirm that it is an excellent choice. For FSWM, this analysis is completely new. This part is in Section 3.2.1.

- Sensitivity to Noise of Consistently Significant AF Functions.** In order for the consistently significant AF functions to be reliably useful for “everyday analysis”, it is also important to assess how sensitive they are to the presence of “noise”. That is, how the performance of an AF matrix varies as a function of the number of its entries that are not statistically significant. To the best of our knowledge, this sensitivity analysis is completely new and the findings are very informative. It is described in Section 3.2.2. Based on those results, a small amount of “noise” is enough to have a significant reduction in performance. Therefore, the operational range of those functions is limited to AF matrices that pass at least 95% of the statistical significance test. Interestingly, since we used phylogenetic tree construction for this sensitivity analysis, we find that UPGMA [28] is better than the Neighbor Joining method [26]) (NJ, for short), in dealing with small amounts of “noisy” entries in the AF matrices. This aspect of those two methods is new.

In the remainder of this paper, we assume that the reader is familiar with Apache Spark [2]. For the convenience of the reader, and due to space limitations, a short primer is reported in Section 1 of the Supplementary Material.

2 Methods

2.1 A Selection of Alignment-free Distances and Similarity Functions

The two most popular classes of AF functions are those based on *k-mer statistics* and on *word matches*, sometimes also referred to as *micro alignments*. For this research, we consider those two classes, providing also motivation for our choices.

An AF function in the first class is based on *k-mer statistics* (or *histogram statistics* [20]). and it can be used to analyze a set of sequences as follows. For each sequence in the set, the contiguous subwords of length k therein contained (i.e., *k-mers*) with their associated frequencies are counted. The result is a set of vectors. Then, sequences are compared pairwise by computing suitable distance/similarity functions between each pair of vectors. The interested reader can find in [31] a list of the ones that have been the object of a recent benchmarking study. One of the most surprising findings of that study is that those simple AF functions are among the best performing and most versatile in terms of application domain. We have chosen the best performing ones according to that benchmarking, representatives of all types of AF functions described in [20] and that can be broadly used in biological studies, e.g., metagenomics [3]. The complete list of the selected AF functions is in Section 2.1 of the Supplementary Material, together with their definitions.

An AF function in the second class is based on the notion of *match* between two sequences. This latter is usually encoded via a binary vector, where the one entries indicate the positions where two subsequences of the two sequences must be identical. Zero entries may not matter. Hence, the distance between two sequences is mainly estimated according to the length of their substring matches. The interested reader can find examples of those methods in [16, 17, 21]. For our study, we have chosen *spaced word methods* [18]. In

particular, the **FSMW** distance, since it has emerged as the most competitive in this class of AF functions [18, 31], details regarding its definition are in Section 2.2 of the Supplementary Material.

2.2 A Spark Platform for Fast Computation of AF Functions: the Basic Pipeline

In this section, we first provide a “user level” functional description of the proposed platform. Then we outline, again at a functional level, how the architectural issues influencing scalability have been dealt with.

2.3 A User-view of the Basic Pipeline as a General and Extensible Spark Programming Paradigm for Implementation of AF Functions

To a user, the basic pipeline appears as a succession of stages, described next. Assuming that the dataset to be processed is composed of n sequences, the output is an $n \times n$ matrix, in which entry (i, j) corresponds to the value of the chosen AF function on sequences i and j . It is also worth pointing out that since the input sequences are partitioned over different computing nodes, two steps are required to collect a global statistic. First, the desired statistic is partially evaluated on each node holding a part of a given sequence. Second, all partial statistics are aggregated to derive the global statistic.

- *Stage 1: Collection of Partial Statistics.* In this stage, the statistic that needs to be collected, e.g. k -mers, is extracted from each of the input sequences and provided as output. This is transparently done in a distributed way, so that each computing node extracts the partial statistic from the parts of the input sequences it stores. We anticipate that Stage 3 takes care of aggregating the different partial statistics extracted from a same sequence. User code can be provided to support more statistics, in addition to those already included in the platform.
- *Stage 2: Feature based Statistics Filtering.* The user implementing the AF algorithm may require the exclusion of a selected subset of features from the statistics it is computing, e.g., specific k -mers such as those containing the “N” character. To this end, this stage acts as a filter to exclude from the output of the previous stage the selected features, according to conditions specified by the user. The filtering occurs at this point, so as to (possibly) alleviate the workload of the following stages.
- *Stage 3: Statistics Aggregation.* All partial statistics extracted by different computing nodes during Stage 1 (and possibly Stage 2), but originating from the same input sequence are automatically and transparently gathered on a same node and aggregated. For instance, statistics about a particular k -mer and extracted from different parts of the same sequence are summed to obtain the overall k -mers statistics for that sequence. User code can be provided to support aggregation for statistics, in addition to those already included in the platform.
- *Stage 4: Value-based Aggregated Statistics Filtering and Normalization.* Stage 2 filters the features existing in a statistic, while this stage filters according to a user-defined condition targeting the aggregated value assumed by a feature in a statistic. For instance, one would want to exclude low frequency k -mers when collecting k -mers statistics. This stage also performs, if required, data normalization. Indeed, as well argued in [20], it is advisable to take the statistic of each sequence, e.g., k -mers counts, and transform it so that all the statistics refer to the same scale. Details are in [10, 20]).
- *Stage 5: AF matrix computation.* For each pair of different input sequences, their final aggregated statistics are sent by the platform to the same node. The AF function that has been chosen, from the ones available, is evaluated on each pair of sequences and the AF matrix is filled accordingly. In addition to the ones available, additional functions can be supported by providing user code.

The output AF matrix is encoded as a distributed data structure, whose content can be saved on file or used as input for further analysis.

Each of the aforementioned stages is modeled as one or more Spark distributed transformations. A general and extensible library of built-in basic functions implementing them is described in details in Section

3.1 of the Supplementary Material. The user interested in supporting a new statistic and/or implementing a variant of these functions can provide her code, as described in Section 3.2 of the Supplementary Material.

2.4 Architectural Engineering: Tuning the Pipeline as a Function of the Workload

We briefly highlight the different data partitioning strategies supported by FADE, designed with the aim of tuning the basic pipeline as a function of the input workload so as to allow for an efficient and scalable execution. Additional details regarding them are available in Section 4 of the Supplementary Material, while their comparative experimental evaluation is reported in Section 3.1.

- *Strategy 1: Total aggregation.* This strategy allows for very good execution times when extracting and processing statistics having an overall small size. This is possible because all statistics (either partial or aggregated) extracted during the pipeline are maintained and processed on a single node of the distributed system. The same occurs to the partial AF function evaluated on each statistic. On one side, this implies that no distributed computation occurs, apart from that of Stage 1. On the other side, this strategy allows to avoid the data transmission overhead required to transfer data to the nodes of the distributed system prior to their processing.
- *Strategy 2: No aggregation.* This strategy allows for a very good scalability when extracting and processing statistics from very large input data. This is possible because every single statistic (either partial or aggregated) extracted during the pipeline is managed as a stand-alone data object. The same occurs to the partial AF function evaluated on each statistic. The only aggregation occurs at the end of the pipeline when, for each pair of distinct sequences, partial AF function values are combined to return the overall value of the function. This ensures for a very fine scalability and load balancing as Spark tends to scatter these data objects uniformly at random on the different nodes of the distributed system. This holds because the amount of memory required to process single data objects is, typically, much smaller than the one required for processing collections of data objects.
- *Strategy 3: Partial aggregation.* This strategy allows for a good trade-off between efficiency and scalability when extracting and processing statistics from large input data. This is possible because all statistics (either partial or aggregated) extracted during the pipeline are partitioned into bins. The same occurs to the partial AF function evaluated on each statistic. Consequently, each node processes a smaller number of data records batches (i.e., the content of each bin) rather than a (potentially) much larger number of single data records. This has a positive effect both on the processing and the communication times.

2.5 Reliability and Robustness of AF Functions

2.5.1 Reliability of an AF Function via a Hypothesis Test Monte Carlo Simulation

Intuition suggests that the larger the number of entries in the AF matrix not due to “chance”, the more indicative of biological relevance the outcome, e.g. phylogenetic tree, of that function use is expected to be. However, experience suggests that AF functions may have a “behavior” that depends on the dataset being processed. Therefore, it is uncontroversially desirable to use functions that are consistently significant.

We formalize such an intuition by requiring that a consistent AF function must provide a high percentage of statistically significant entries in its corresponding AF matrix, with very little dependence on the input dataset. In regard to those tests, we resort to a Monte Carlo Simulation method, in which we assume as the Null Hypothesis H_0 that two biological sequences are as similar as two random ones. An entry (i, j) of an AF matrix is statistically significant when the Null Hypothesis involving the two sequences associated to that entry is rejected with a given significance level. The entire simulation consists of three steps, described next. The Spark algorithms implementing it are briefly described in Section 5 of the Supplementary Material. In what follows, we consider only the case of similarities, since the case of distances is analogous.

Step 1: Synthetic Datasets Generation via Bootstrapping. The first step is to devise a procedure that generates synthetic datasets, that are meant to represent “random data”. Such a task can be accomplished by choosing a *Null Model*, e.g., an Information Source emitting symbols uniformly and at random. However, in our case, it seems more appropriate to resort to *bootstrapping* [8], i.e., to generate the synthetic datasets from real ones, since it is desirable to preserve the biological origin of the input dataset also in the synthetic ones. To this end, we proceed as follows.

Let S be the input dataset and let q be a parameter. All q -mers of sequences in S are extracted and placed in a bin B . Then, in order to obtain a synthetic dataset \hat{S} , we extract uniformly and at random q -mers from B in order to form new sequences to be included in \hat{S} . This latter has the same number of sequences as in S and each sequence in \hat{S} corresponds to only one in S in terms of length.

It is to be noted that the parameter q allows us to generate synthetic datasets along a wide spectrum of subsequence statistics present in S , e.g., $q = 1$ corresponds to the case in which the synthetic dataset is generated according to the empirical probability distribution of symbols in S .

Step 2: Significance Test via a Monte Carlo Simulation for Two Sequences. The next step consists of the following simulation, adapted from [12]. It applies to sets of two sequences. C denotes the AF function to which the procedure is applied. The Null Hypothesis H_0 is that two input sequences are as similar as two random ones, when similarity is assessed via S . We want to reject the Null Hypothesis with confidence level α , with the performance of ℓ Monte Carlo Simulations.

Procedure MECCA(ℓ, C, S, α)

- (1) For $1 \leq i \leq \ell$, compute a new set of two sequences \hat{S}_i according to the procedure outlined in Step 1. Compute the similarity between the two sequences in \hat{S}_i via C . Let T_i be its value.
- (2) For $1 \leq i \leq \ell$, sort the T_i values in non-decreasing order and let SL be the corresponding list.
- (3) Let T denote the value of C computed on S . Let j be the maximal index such that $SL[j] < T$. Let $\delta = (j/\ell)$. The p -value is then $1 - \delta$ and, letting α be the desired significance level, the hypothesis that the two sequences in S are as similar as two randomly chosen ones is rejected with that significance level if $1 - \delta \leq \alpha$.

Step 3: AF Matrix Significance via Bonferroni Correction. Consider now a set S consisting of n sequences, labeled from 1 to n . Let F be the $n \times n$ AF matrix for S computed via C . In order to assess how statistically significant is F , with family-wise significance level α , we can resort to the pairwise application of the simulation procedure outlined in Step 2. Since we are performing $m = \frac{n(n-1)}{2}$ hypothesis tests, we have to correct for rejecting H_0 simply by chance. Since those tests may not be assumed to be independent, we use the well known Bonferroni correction. That is, for each test performed for each entry of F , H_0 is rejected with significance level α/m . Then, we reject the Null Hypothesis that the matrix is no better than one obtained on a set of random sequences, if all m entries pass the test. However, even if the full matrix does not pass the test, such a procedure outlines entries that are statistically significant in terms of similarity values of the corresponding sequences.

2.5.2 Robustness of an AF Function via a Matrix Perturbation Method

In addition to the reliability of an AF function, it is also important to assess how robust is the function with respect to “noise”, i.e., the number of entries in the AF matrix that are not statistically significant. Informally, we refer to this as the *operational range* of an AF function. In order to empirically estimate this latter, we informally proceed as follows: the larger is the amount of noise injected in a AF matrix returned by a given function, the worse should be the biological relevance of its outcome (e.g., the phylogenetic tree). Then, to measure the operational range of a function, we start from the AF matrix returned by that function on a given dataset and assume as a reference its performance score. Finally, we study its performance variation, assuming it will decrease while increasing the amount of noisy entries. More precisely, we assume

that the dataset under consideration has a gold standard solution. For this study, it is a phylogenetic tree associated to the species in the dataset. We also use a computational method G to build a phylogenetic tree from an AF matrix and a distance/similarity measure D to assess how different is a tree produced by G via an AF matrix with respect to the gold standard. Formal details of the basic perturbation step follow.

- Given the AF matrix computed by C on dataset S , we select uniformly and at random a given percentage of entries and substitute each value with a “noisy” one. For histogram-based functions, such a value is taken randomly among the ones appearing in the AF matrices generated via Monte Carlo simulation. As for FSWM, since the AF matrices it returns when evaluated on the synthetic datasets are likely to contain null values (see discussion about filtering in Section 3.2.1), the “noisy” value is obtained by increasing the original one by a value chosen uniformly and at random. For both types of AF functions, the new matrix so obtained is used to build a tree via G , and its distance from the gold standard is computed via D . The loss in performance at the given noise level is given by the difference in the D score obtained with the noisy vs the original matrix.

3 Results and Discussion

As a first preliminary step, we selected for our study datasets coming from [31], with the criterion that the AF functions chosen here would work on them (see again [31]). They are reported in Section 6 of the Supplementary Material. As a second preliminary step, we assessed that our AF functions implementations are in line with those used for the benchmarking of the AFproject [30] (details are omitted for brevity and available upon request). All our experiments have been executed on the hardware platform described in Section 6 of the Supplementary Material. The parameters of the algorithms have been set according to the procedures described in Section 7 of the Supplementary Material. Execution times have been measured by collecting the job execution elapsed time returned within the Spark framework.

3.1 Assessing the Scalability of Methods Supported by FADE and the Effectiveness of The Aggregation Strategies

3.1.1 Experiments

When considering a single computing platform where multiple computing units coexist on the same motherboard and communicate, at no cost, by using shared memory, nothing can beat the performance of a native application purposely developed to take advantage of that setting. On the other hand, the number of processing units and memory size on a single motherboard is a natural performance bottleneck.

Given the above scenario, the use of a distributed framework is to overcome such a bottleneck, by scaling performance as a function of the number of computing units far over the ones on a stand-alone computing platform. Therefore, it is important to assess the efficiency of methods supported by our software pipeline, when compared to analogous state-of-art shared memory alignment-free tools, while assessing the scalability of our framework. For such an analysis, we have chosen two reference software systems that provide shared memory parallel software for their evaluation: FSWM [16] and Mash [23].

We measure the execution time of FSWM and Mash, as well as that of their implementations supported by our framework, i.e., FADE-FSWM and FADE-Mash, on the Plants (assembled) dataset and while increasing the level of parallelism, as explained shortly. Mash and FSWM have been executed on a single machine equipped with 8 computing cores. The FADE versions have been executed on a distributed framework equipped with 24 worker nodes, for a total of 128 computing cores. Each of those machines is identical to the one on which the shared memory algorithms have been executed. The result of this experiment is visible in Figure 2, where we report the execution time of these applications as a function of the number of concurrent threads, for shared memory applications, or of Spark workers, for FADE. Each thread/worker executes on a single core.

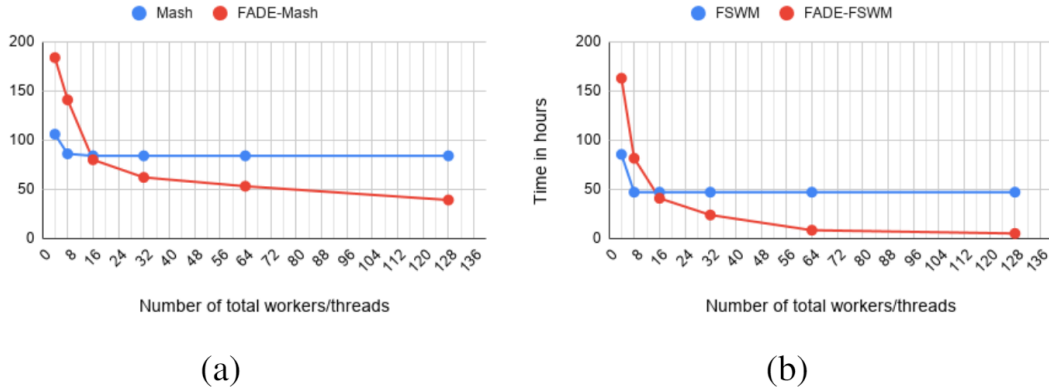


Figure 2: Execution time, on the ordinate axis, required by both versions of (a) Mash and (b) FSWM on the Plants (assembled), using an increasing number of workers/threads. The Partial Aggregation strategy has been used, since naturally suited to the algorithmic methods supporting Mash and FSWM, respectively.

In order to assess the effectiveness of the aggregation strategies supported by FADE, we executed a significance test for all of the histogram-based functions, with different data sizes. The results are reported in Table 1.

3.1.2 Results

The results of the first experiment show that the shared memory versions of Mash and FSWM are faster than the implementations supported by FADE, when using a very small number of threads compared to workers. This is expected. With the increase of the number of those units, the performance gap narrows. As soon as the use of 16 threads/workers is reached, FADE-Mash and FADE-FSWM require about the same execution time as that of their shared memory counterparts. However, those latter stop scaling, due to the bottleneck mentioned earlier, while the FADE methods continue to scale as more threads/workers are added.

The results of the second experiment clearly show that the Partial Aggregation strategy is the most flexible of the three, since it guarantees equal or better performance with respect to the other two, independently of the workload.

Dataset	No Aggregation	Partial Aggregation	Total Aggregation
Small	0.9	0.9	1.1
Large	2.73	2.82	NA
Very large	48.3	42.0	NA

Table 1: Execution time, in minutes, required by our framework to execute one instance of the AF significance test, for all of the histogram-based functions, on three reference datasets with different aggregation strategies. The NA value indicate that the test took too long to complete and was stopped. Small = Mitochondria (assembled). Large = Shigella (assembled). Very large = Plants (assembled).

3.2 Reliability and Robustness of AF Functions: An Application of FADE to Data/Compute Intensive Analysis

3.2.1 Reliability of AF Functions

The Experiment. For each of the benchmark datasets included in this study, we execute the AF statistical significance test described in Section 2.5.1. Specifically, for each matrix obtained via a dataset, we generate

Mitochondria (assembled)			
	q=1	q=7	q=10
Canberra	100	100	100
Chebyshev	3.3	3	3.7
ChiSquare	100	100	99.3
D2	100	100	99
D2S	100	99.7	81.7
D2Z	100	99.7	98.7
D2*	100	91.3	82
Euclidean	89.3	85.7	84.3
FSWM	100	100	100
Harmonic Mean	98.7	70.7	48
Intersection	65	40	35.3
Jaccard	27.3	28	43.3
Jeffreys	99.3	96	95
Jensen Shannon	99.7	96.3	95
Kulczynski2	100	99.3	99
Manhattan	100	100	99.3
Mash	78	86	83.7
Squared Chord	100	100	99.3

E.coli (unassembled, coverage=1)			
	q=1	q=7	q=10
Canberra	100	100	100
Chebyshev	0	0	0
ChiSquare	58.4	21.9	16.3
D2	100	100	100
D2S	100	100	100
D2Z	100	79.8	65
D2*	100	100	100
Euclidean	1.5	1.5	1.5
FSWM	100	100	100
Harmonic Mean	100	96.3	85
Intersection	100	15.5	2.5
Jaccard	0	0	0
Jeffreys	71.9	28.1	15.3
Jensen Shannon	82.5	34.7	17.5
Kulczynski2	13.3	3.7	3.7
Manhattan	66.5	28.1	18.2
Mash	0	0	0.9
Squared Chord	41.1	20.7	16.3

Plants (assembled)			
	q=1	q=7	q=10
Canberra	22	22	27.5
Chebyshev	0	0	0
ChiSquare	6.6	6.6	6.6
D2	98.9	97.8	97.8
D2S	14.3	7.7	7.7
D2Z	1.1	1.1	1.1
D2*	100	100	100
Euclidean	0	0	0
FSWM	100	100	100
Harmonic Mean	19.8	18.7	18.7
Intersection	42.9	29.7	27.5
Jaccard	1.1	3.3	3.3
Jeffreys	2.2	2.2	2.2
Jensen Shannon	3.3	2.2	2.2
Kulczynski2	2.2	2.2	2.2
Manhattan	13.2	11	11
Mash	2.2	2.2	2.2
Squared Chord	6.6	6.6	6.6

Figure 3: Summary of the Hypothesis Test results for the different AF functions considered in this research when executed on three different datasets with $q = 1, 7, 10$ and with significance level set to 1%. Datasets have been chosen as described in the Main Text.

additional ones via bootstrapping. Then, we reject the null hypothesis family-wise with p -value $\leq 1\%$, applying Bonferroni correction to all of its m entries. The number l of simulations for each test has been chosen according to the size of the dataset being processed, so as to guarantee the execution of the experiments in a reasonable amount of time. An outline of these parameters, including the choice of k and the values of q , is available in Section 7 of the Supplementary Material.

A summary of the results is shown in Figure 10 of the Supplementary Material. For each dataset, each AF function and each considered null model, the percentage of entries passing the test is reported. This value is drawn in green, if at least 75% of the entries passes the test, in red, if no entry passes the test, and in yellow, in the remaining case. Figure 3 reports a representative synopsis of those results: a dataset where most AF functions perform well, one in which many perform poorly and one in which most perform poorly.

Results and Insights

- *A novel class of AF functions:consistently significant.* With reference to Figure 3 and Figure 10 of the Supplementary Material, it is evident that there are AF functions returning matrices passing the family-wise significance test either fully or with a high percentage of entries, for all of the benchmark datasets. We denote them as *consistently significant*, since they behave consistently well, independently of the dataset and, as already stated in the Introduction, it is well known that statistical relevance is a good indication of biological relevance [7, 11, 19]. They are: D2, D2*, FSWM. The remaining ones either perform inconsistently or poorly (Chebyshev).
- *Filtering can be useful.* While the statistical guarantees of AF functions in the D2 family are well known and have been identified via deep investigations [29], we find a novel fact regarding FSWM: it is quite good in delivering matrices that pass the statistical significance test. This can be attributed to the filtering mechanism present in the algorithm and that was designed to “flush out” the parts of the statistics it is collecting and that are considered “weak”. Such a filtering is able to detect the “low relatedness” of the synthetic genomes ensuring that the original genomes “win” the test.

3.2.2 Robustness of AF functions: the case of consistently significant functions

The Experiment. We concentrate only on AF functions that have been classified as consistently significant in the previous section, since they are the most likely to be useful on a day-to-day basis. For each of them, this experiment is conducted by injecting an increasing percentage of noisy entries into an AF matrix, following the method outlined in Section 2.5.2. As a computational method G to build a phylogenetic tree from an AF matrix, we use both UPGMA and NJ. This choice is motivated by the fact that, with pros and cons, those two methods are reference points in the Literature. Moreover, as a distance/similarity measure D to assess how different is a tree produced by G via an AF matrix with respect to the gold standard, we use both the Matching Cluster [4] (MCM for short) and the Robinson Fould [25] (RF, for short) metrics. Such a choice is motivated by the fact that they were among the top performing ones in a classic benchmarking study [15], in particular the second. As for RF, it is quite sensitive to small changes in tree topology, i.e., it “saturates” rapidly, a fact that has been criticized in the past [24]. Yet, it is a standard in the Literature.

The results of this experiment are available in Figures 11-16 of the Supplementary Material, where we plot the “distance” from the gold standard as the number of noisy entries grows. Moreover, to gain more insights, we also report in Tables 2-4, the variation in distance only when 10% noisy entries are injected. Moreover, we do not consider changes in the distance from the gold standard that are in the interval $+/-3$, i.e., reasonably close to the case of “no noise”.

Results and Insights

- *AF matrices are sensitive to the injection of noisy entries.* Figures 11-16 of the Supplementary Material show a trend of deterioration in classification ability as the percentage of noisy entries increases. Such a trend is more pronounced when we use the RF metric to assess how close to the gold standard is a phylogenetic tree computed from a noisy matrix with respect to the one that uses an uncorrupted AF matrix. This is possibly due to the “saturation” effect, but MCM largely confirms those experiments.

The results in Tables 2-4 give more compelling evidence of the sensitivity to noise of the AF functions under analysis. In particular, FSWM is quite sensitive, with both NJ and UPGMA. As for the other two AF functions, they display such a sensitivity mostly when used with NJ. Or, better, UPGMA seems to be able to tolerate a small amount of noise in the input data. We believe that such an aspect of UPGMA is novel. It is to be noted that the entries with negative values, i.e., an improvement in classification with the addition of random entries, is justified by the fact that the original classification was so poor that even noise could do no harm.

In conclusion, the indication we receive from the above is to use D2, D2* and FSWM in conjunction with UPGMA and with a high number of statistically significant entries in the AF matrix.

		Mitochondria	Shigella	Yersinia			Mitochondria	Shigella	Yersinia			Mitochondria	Shigella	Yersinia
RF	NJ	6 (200%)	10 (90.9%)	4 (Inf%)	RF	NJ	5 (125%)	7 (63.6%)	-	RF	NJ	17 (425%)	20 (500%)	-
	UPGMA	-	-	-		UPGMA	-	-	-		UPGMA	10 (333.3%)	6 (120%)	-
MCM	NJ	28 (23.3%)	23 (26.4%)	14 (Inf%)	MCM	NJ	15 (14.6%)	7 (8.0%)	-	MCM	NJ	40 (40.8%)	69 (97.2%)	4 (33.3%)
	UPGMA	-31 (-31.0%)	-39 (-24.2%)	-		UPGMA	-	5 (6.3%)	-		UPGMA	49 (245%)	41 (146.4%)	-
D2					D2*					FSWM				

Table 2: Corruption tables relative to D2, D2* and FSWM functions and $q=1$, varying dataset, method for building the phylogenetic tree, metric to evaluate the distance between the reference and obtained tree. In each entry, it is reported the difference between the phylogenetic distance evaluated on the distance matrix with 10% of corrupted entries and the original distance matrix (between the brackets is reported the percentage difference, indicating with 'Inf' when the tree obtained from the original distance matrix is equal to the reference tree, i.e. the distance between the trees is 0). Dashed entries represent cases where the introduction of noise caused only a small change in distance, as specified in the Main Text.

		Mitochondria	Shigella	Yersinia
RF	NJ	4 (133.3%)	11 (100%)	4 (Inf%)
	UPGMA	-	-	-
MCM	NJ	22 (18.3%)	31 (35.6%)	14 (Inf%)
	UPGMA	-32 (-32.0%)	-28 (-17.4%)	-

D2

		Mitochondria	Shigella	Yersinia
RF	NJ	6 (150%)	7 (63.6%)	-
	UPGMA	-	-	-
MCM	NJ	21 (20.4%)	67 (77.0%)	-
	UPGMA	-	5 (6.3%)	-

D2*

		Mitochondria	Shigella	Yersinia
RF	NJ	17 (425%)	20 (500%)	-
	UPGMA	10 (333.3%)	4 (80%)	-
MCM	NJ	40 (40.8%)	100 (140.8%)	4 (33.3%)
	UPGMA	49 (245%)	32 (114.3%)	-

FSWM

Table 3: Corruption tables relative to $q=7$, The Legend is as in Table 2

		Mitochondria	Shigella	Yersinia
RF	NJ	4 (133.3%)	10 (90.9%)	4 (Inf%)
	UPGMA	-	-	-
MCM	NJ	22 (18.3%)	24 (27.6%)	14 (Inf%)
	UPGMA	-32 (-32%)	-24 (-14.9%)	-

D2

		Mitochondria	Shigella	Yersinia
RF	NJ	7 (175%)	7 (63.6%)	-
	UPGMA	-	-	-
MCM	NJ	22 (21.4%)	19 (21.8%)	-
	UPGMA	-	5 (6.3%)	-

D2*

		Mitochondria	Shigella	Yersinia
RF	NJ	17 (425%)	19 (475%)	-
	UPGMA	10 (333.3%)	4 (80%)	-
MCM	NJ	40 (40.8%)	87 (122.5%)	4 (33.3%)
	UPGMA	49 (245%)	33 (117.9%)	-

FSWM

Table 4: Corruption tables relative to $q=10$, The Legend is as in Table 2

4 Conclusion

We have provided the first Spark platform, supporting AF functions evaluation, that is extensible and that guarantees good scalability with respect to the workload and computational resources available. This is witnessed by the very good performance, both in terms of computing time and scalability, of FADE-Mash and FADE-FSWM when compared to their state-of-the-art shared memory counterparts. This is particularly relevant, since the shared memory version of Mash has a leadership position due to its speed.

In order to demonstrate the ability of FADE to support large data and compute intensive tasks, we have conducted a novel analysis of AF functions. Thanks to it, we gain insights into which functions are best suited for day-to-day AF genomic analysis: D2, D2* and FSWM. We also account for the “operational range” of those latter measures, i.e., circumstances in which a result due to those methods can be trusted. Such an “operational range” translates into the requirement to use the AF matrices corresponding to those methods when the vast majority of their entries is statistically significant. Moreover, the use of UPGMA is recommended, since it is less sensitive to “noise” in the input data, compared to NJ.

Acknowledgements

R.G. is grateful to Prof. Chiara Romualdi for helpful discussions. All authors would like to thank the Department of Statistical Sciences of University of Rome - La Sapienza for computing time on the TeraStat cluster and for other computing resources, and the GARR Consortium for having made available a cutting edge OpenStack Virtual Datacenter for this research.

Funding

G.C., R.G. and U.F.P. are partially supported by GNCS Project 2019 “Innovative methods for the solution of medical and biological big data”. R.G. is additionally supported by MIUR-PRIN project “Multicriteria Data Structures and Algorithms: from compressed to learned indexes, and beyond” grant n. 2017WR7SHH. U.F.P. and F.P. are partially supported by Università di Roma - La Sapienza Research Project 2018 “Analisi, sviluppo e sperimentazione di algoritmi praticamente efficienti”.

References

- [1] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410, 1990.

- [2] Apache Software Foundation. Apache Spark. (Available from: <http://spark.apache.org>), 2016.
- [3] G. Benoit, P. Peterlongo, M. Mariadassou, E. Drezen, S. Schbath, D. Lavenier, and C. Lemaitre. Multiple comparative metagenomics using multiset k-mer counting. *PeerJ Computer Science*, 2:1, 2016.
- [4] D. Bogdanowicz and K. Giaro. On a matching distance between rooted phylogenetic trees. *International Journal of Applied Mathematics and Computer Science*, 23(3):669–684, 2013.
- [5] G. Cattaneo, U. Ferraro Petrillo, R. Giancarlo, and G. Roscigno. An effective extension of the applicability of alignment-free biological sequence comparison algorithms with Hadoop. *The Journal of Supercomputing*, pages 1–17, 2016.
- [6] G. Cattaneo, R. Giancarlo, U. Ferraro Petrillo, and G. Roscigno. MapReduce in computational biology via Hadoop and Spark. In Ranganathan, N. S., S. C. K., and M. Gribskov, editors, *Encyclopedia of Bioinformatics and Computational Biology*, volume 1, pages 221–229. Elsevier, Oxford, 2019.
- [7] S. Dudoit and J. Fridlyand. A prediction-based resampling method for estimating the number of clusters in a dataset. *Genome Biology*, 3:36, 2002.
- [8] B. Efron. Bootstrap methods another look at the jackknife. *The Annals of Statistics*, 7:1–26, 1979.
- [9] R. Giancarlo, S. E. Rombo, and F. Utro. Epigenomic k-mer dictionaries: Shedding light on how sequence composition influences nucleosome positioning *in vivo*. *Bioinformatics*, 31:2939–2946, 2015.
- [10] R. Giancarlo, S. E. Rombo, and F. Utro. In vitro versus in vivo compositional landscapes of histone sequence preferences in eucaryotic genomes. *Bioinformatics*, 34:3454–3460, 2018.
- [11] R. Giancarlo, D. Scaturro, and F. Utro. Computational cluster validation for microarray data analysis: experimental assessment of cleft, consensus clustering, figure of merit, gap statistics and model explorer. *BMC Bioinformatics*, 9(1):462, 2008.
- [12] R. Giancarlo, D. Scaturro, and F. Utro. A tutorial on computational cluster analysis with applications to pattern discovery in microarray data. *Mathematics in Computer Science*, 1:655–672, 2008.
- [13] J. E. Gonzalez, Y. Low, H. Gu, D. Bickson, and C. Guestrin. Powergraph: Distributed graph-parallel computation on natural graphs. In *Presented as part of the 10th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 12)*, pages 17–30, 2012.
- [14] S. D. Kahn. On the future of genomic data. *Science*, 331(6018):728–729, 2011.
- [15] M. K. Kuhner and J. Yamato. Practical performance of tree comparison metrics. *Systematic Biology*, 64(2):205–214, 2015.
- [16] C.-A. Leimeister, M. Boden, S. Horwege, S. Lindner, and B. Morgenstern. Fast alignment-free sequence comparison using spaced-word frequencies. *Bioinformatics*, 30:1991–1999, 2014.
- [17] C.-A. Leimeister and B. Morgenstern. kmacs: the k-mismatch average common substring approach to alignment-free sequence comparison. *Bioinformatics*, 30:2000–2008, 2014.
- [18] C.-A. Leimeister, S. Sohrabi-Jahromi, and B. Morgenstern. Fast and accurate phylogeny reconstruction using filtered spaced-word matches. *Bioinformatics*, 33:971–979, 2017.
- [19] M. Y. Leung, G. M. Marsh, and T. P. Speed. Over- and underrepresentation of short DNA words in herpesvirus genomes. *Journal of Computational Biology*, 3(3):345–360, 1996.
- [20] B. B. Luczak, B. T. James, and H. Z. Girgis. A survey and evaluations of histogram-based statistics in alignment-free sequence comparison. *Briefings in Bioinformatics*, 20(4):1222–1237, 12 2017.

- [21] B. Morgenstern, B. Zhu, S. Horwege, and C. A. Leimeister. Estimating evolutionary distances between genomic sequences from spaced-word matches. *Algorithms for Molecular Biology*, 10(1):5, 2015.
- [22] H. Mushtaq and Z. Al-Ars. Cluster-based Apache Spark implementation of the GATK DNA analysis pipeline. In *Bioinformatics and Biomedicine (BIBM), 2015 IEEE International Conference on*, pages 1471–1477. IEEE, 2015.
- [23] 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. *Genome Biology*, 17:132, 2016.
- [24] D. Penny, L. R. Foulds, and M. D. Hendy. Testing the theory of evolution by comparing phylogenetic trees constructed from five different protein sequences. *Nature*, 297(5863):197–200, 1982.
- [25] D. F. Robinson and L. R. Foulds. Comparison of phylogenetic trees. *Mathematical biosciences*, 53(1-2):131–147, 1981.
- [26] N. Saitou and M. Nei. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular biology and evolution*, 4(4):406–425, 1987.
- [27] T. F. Smith and M. S. Waterman. Identification of common molecular subsequences. *Journal of molecular biology*, 147(1):195–197, 1981.
- [28] P. H. A. Sneath and R. R. Sokal. *Numerical taxonomy. The principles and practice of numerical classification*. WH Freeman, San Francisco, W.H. Freeman and Company., USA, 1973.
- [29] K. Song, J. Ren, G. Reinert, M. Deng, M. S. Waterman, and F. Sun. New developments of alignment-free sequence comparison: measures, statistics and next-generation sequencing. *Briefings in Bioinformatics*, 15(3):343–353, 09 2013.
- [30] A. Zieiezinski, H. Z. Girgis, G. Bernard, C.-A. Leimeister, K. Tang, T. Dencker, A. K. Lau, S. Röhling, J. J. Choi, M. S. Waterman, M. Comin, S.-H. Kim, S. Vinga, J. S. Almeida, C. X. Chan, B. T. James, F. Sun, B. Morgenstern, and W. M. Karlowski. Afproject. <http://afproject.org>, 2019.
- [31] A. Zieiezinski, H. Z. Girgis, G. Bernard, C.-A. Leimeister, K. Tang, T. Dencker, A. K. Lau, S. Röhling, J. J. Choi, M. S. Waterman, M. Comin, S.-H. Kim, S. Vinga, J. S. Almeida, C. X. Chan, B. T. James, F. Sun, B. Morgenstern, and W. M. Karlowski. Benchmarking of alignment-free sequence comparison methods. *Genome Biology*, 20(1):144, 2019.

Alignment-free Genomic Analysis via a Big Data Spark Platform

Supplementary Material

Umberto Ferraro Petrillo* Francesco Palini* Giuseppe Cattaneo†
Raffaele Giancarlo‡

Abstract

Additional details about the Main Manuscript are provided in this document.

1 Spark and Distributed Architectures

Apache Spark [2] is a framework used mainly to support programs with in-memory computing and acyclic data-flow model to be executed on a distributed computing architecture. In simple and intuitive terms, this latter can be described as a set of computing nodes that cooperate in order to solve a problem via local computing and by exchange of “messages” [11]. On that type of architecture, Spark can be used for applications that reuse a working set of data across multiple parallel operations (e.g., iterative algorithms) and it allows the combination of streaming and batch processing, as opposed to Hadoop [1], that can be only used for batch applications. In addition, Spark is not limited to support only the *MapReduce* [5] paradigm, although it preserves all the facilities of *MapReduce*-based frameworks, and it can have the Hadoop as the underlying middleware.

It provides high-level APIs which support multiple languages (Scala, its native language, Java and Python) and is made of two fundamental logic blocks: a programming model that creates a task dependency graph, and an optimized runtime system which exploits this graph to deploy code and data and schedule work units on the distributed system nodes. At the core of the Spark programming model is the *Resilient Distributed Dataset* (RDD) abstraction, a fault-tolerant, distributed data structure that can be created and manipulated using a rich set of operators: programmers start by defining one or more RDDs through *transformations* of data that originally resides on stable storage or other RDDs.

2 A Selection of Alignment-free Distances and Similarity Functions

In this section, we provide definitions of the AF functions included in this study.

2.1 Histogram Statistics Selected for this Study

In what follows, we adopt both the classification and the notation from [10]. It is to be remarked that some of the AF functions defined next appear in the Literature with different names or they are easy variants of the functions introduced here. For instance, FFP adopted in [16] is the well known Jensen-Shannon Divergence defined here, while Skmer [14] is a fast approximation method for the computation of the well known Jaccard Index defined here. For the interested reader, the original publication where the functions defined here have been introduced can be found in [10] for the less known cases.

*Dipartimento di Scienze Statistiche, Università di Roma - La Sapienza, Rome, 00185, Italy

†Dipartimento di Informatica, Università di Salerno, Fisciano (SA), 84084, Italy

‡Dipartimento di Matematica ed Informatica, Università di Palermo, Palermo, 90133, Italy

Given a set of sequences $S = \{s_1, \dots, s_n\}$, a k -mer histogram h_s for a sequence s in the set is defined as follows:

$$h_s = \langle c(w_1), c(w_2), \dots, c(w_{|K|}) \rangle \quad (1)$$

where $c(w_i)$ is the number of occurrences of the word w_i (i.e. the i -th k -mer) in the sequence s and K is the set of all possible words of length k over the alphabet $\{A, C, G, T\}$.

2.1.1 The Minkowski Family

Given two sequences s and t and their associated statistics h_s and h_t , the *Euclidean* distance is defined as:

$$Euclidean(h_s, h_t) = \sqrt{\sum_{w \in K} (h_s(w) - h_t(w))^2} \quad (2)$$

A widely adopted variant of the *Euclidean* distance is the *Manhattan* distance:

$$Manhattan(h_s, h_t) = \sum_{w \in K} |h_s(w) - h_t(w)| \quad (3)$$

Another member of this family, the *Chebyshev* distance, is based on the idea of applying the p -th root on the *Manhattan* distance, with $p \rightarrow \infty$, and it is defined as follows:

$$Chebyshev(h_s, h_t) = \max_{w \in K} |h_s(w) - h_t(w)| \quad (4)$$

2.1.2 The Match/mismatch Family

Let h_s^* and h_t^* be the sets of k -mers with non-zero entries in h_s and h_t , respectively. The *Jaccard* index measures the similarity between those two sets as follows:

$$Jaccard(h_s, h_t) = \frac{|h_s^* \cap h_t^*|}{|h_s^* \cup h_t^*|} \quad (5)$$

It is to be noted that here we use the “classic” definition of Jaccard index rather than the one given in [10]. Moreover, the computation of the Jaccard index has been the object of much attention in the Information Retrieval community and recently in Bioinformatics. Therefore, we also consider the MinHash approximation algorithm presented in [12]. The basic idea is to reduce the input sequences to sets of integers of small cardinality, each referred to as a *sketch*. To create a single sketch of cardinality s (i.e., the *sketch size*), each k -mer in a sequence is hashed. Such a process assigns k -mers to integers pseudorandomly. The s smallest values are retained. Then, in order to compute the Jaccard index of two sequences, the intersection of their sketches is computed. Indeed, the percentage of shared elements in the sketches can be shown to be an approximation of the Jaccard index (for details, see [12]).

2.1.3 The χ^2 Distance

It is defined as:

$$\chi^2(h_s, h_t) = \sum_{w \in K} \frac{(h_s(w) - h_t(w))^2}{(h_s(w) + h_t(w))} \quad (6)$$

2.1.4 The Canberra Distance

It is a mix between Manhattan and χ^2 distances:

$$Canberra(h_s, h_t) = \sum_{w \in K} \frac{|h_s(w) - h_t(w)|}{(h_s(w) + h_t(w))} \quad (7)$$

2.1.5 The D_2 Statistics

It expresses the similarity of two sequences in a very natural way as a sort of inner product between two histograms as shown in equation (8).

$$D_2(h_s, h_t) = \sum_{w \in K} h_s(w)h_t(w) \quad (8)$$

However, extensive studies of this function suggest that a standardized version of it is more useful in the AF sequence analysis setting. Such a variant is denoted D_2^Z . It is as D_2 , except that the histograms have been standardized via the well known z-score transformation. Details can be found in [10].

We also consider two other members of the D_2 family, denoted D_2^S and D_2^* and described in [17].

D_2^S is based on the finding [15] that if two independent random variables X and Y are normally distributed with mean zero, also $\frac{XY}{\sqrt{X^2+Y^2}}$ is normally distributed. We normalize $h_s(w)$ and $h_t(w)$ as follow:

$$\tilde{h}_s(w) = h_s(w) - np_s(w)$$

and

$$\tilde{h}_t(w) = h_t(w) - mp_t(w)$$

In which n and m are the number of k -mers, respectively, in S and T , while $p_s(w)$ and $p_t(w)$ are the probabilities of the k -mer w under the background model for, respectively, S and T .

The D_2^S statistics is defined as follow:

$$D_2^S = \sum_{w \in K} \frac{\tilde{h}_s(w)\tilde{h}_t(w)}{\sqrt{\tilde{h}_s(w)^2 + \tilde{h}_t(w)^2}} \quad (9)$$

The D_2^* statistic is based on the idea that, for relatively long k -mers, the number of occurrences can be approximated by a Poisson distribution, thus mean and variance are nearly the same.

The D_2^* statistics is defined as follow:

$$D_2^* = \sum_{w \in K} \frac{\tilde{h}_s(w)\tilde{h}_t(w)}{\sqrt{np_s(w)mp_t(w)}} \quad (10)$$

2.1.6 The Intersection Family

From this family, we selected two distances. The first is the *Intersection* distance, also known as *Czekanowski*. It is based on the intersection of the k -mers counts divided by their union. It is:

$$Intersection(h_s, h_t) = \sum_{w \in K} \frac{2 * \min(h_s(w), h_t(w))}{h_s(w) + h_t(w)} \quad (11)$$

The second is the *Kulczynski2* distance;

$$Kulczynski2(h_s, h_t) = A_\mu \sum_{w \in K} \min(h_s(w), h_t(w)) \quad (12)$$

where A_μ is equal to $\frac{4^k(\mu_s - \mu_t)}{2\mu_s\mu_t}$ (μ_s and μ_t denote the mean of the histograms h_s and h_t , respectively).

2.1.7 The Inner Product Family

The *Harmonic Mean* distance is:

$$HarmonicMean(h_s, h_t) = 2 \sum_{w \in K} \frac{h_s(w)h_t(w)}{h_s(w) + h_t(w)} \quad (13)$$

As opposed to the Euclidean distance that computes the square root over the summation value, the *Squared Chord* computes the square root over each value of the histograms independently:.

$$SquaredChord(h_s, h_t) = \sum_{w \in K} \left(\sqrt{h_s(w)} - \sqrt{h_t(w)} \right)^2 \quad (14)$$

This can be simplified as

$$= \sum_{w \in K} h_s(w) + h_t(w) - 2\sqrt{h_s(w)h_t(w)} \quad (15)$$

2.1.8 The Divergence Family

This family uses probabilities to compare two sequences measuring the distance in the log-probability space. From such a family, we selected the *Jeffrey* and the *Jensen-Shannon* distances. The former is defined as:

$$Jeffrey(h_s, h_t) = \sum_{w \in K} (p_s(w) - p_t(w)) \ln \frac{p_s(w)}{p_t(w)} \quad (16)$$

while the second (JSD for short) is defined in (17):

$$JSD(h_s, h_t) = \frac{1}{2} \sum_{w \in K} p_s(w) \log_2 \left(\frac{p_s(w)}{\frac{1}{2}(p_s(w) + p_t(w))} \right) + \frac{1}{2} \sum_{w \in K} p_t(w) \log_2 \left(\frac{p_t(w)}{\frac{1}{2}(p_s(w) + p_t(w))} \right) \quad (17)$$

where $p_s(w)$ is the empirical probability of k -mer w over all the strings of length k from the alphabet $\{A, C, G, T\}$ in the input sequence s and therefore $p_s(w) = \frac{h_s(w)}{4^k}$.

2.2 Spaced Words Methods Selected for this Study

2.2.1 The FSWM Distance

Filtered Spaced Word Matches (FSWM) is a method introduced in [9].

We need to define a spaced word match between two sequences s and t . Given a length l binary pattern P of *match* and *don't care* positions, there is a *spaced word* matching between s , t , in positions i_1 and i_2 , respectively, according to the pattern P , if for each match position m in P , it is true that $s[i_1 + m] = t[i_2 + m]$. This match is also weighted according to the Chiaromonte's substitution matrix [4], applied on the *don't care* positions.

All such matches are collected and the ones with a score lower than a given threshold are discarded. Finally, the distance between s and t is computed by applying *Jukes-Cantor* correction based on the computed spaced word matches:

$$FSWM(s, t) = -\frac{3}{4} \log \left(1 - \frac{4}{3} \frac{mm_{s,t}}{\delta_{s,t}} \right), \quad (18)$$

where $mm_{s,t}$ is the number of *don't care* characters that don't match in the collection of spaced word matches between s and t , while $\delta_{s,t}$ is the total number of *don't care* characters in that collection.

3 A General and Extensible Spark Programming Paradigm for Implementation of AF Functions: Details

3.1 User Programming: Software Library and Examples

3.1.1 The Primitives Library

This library consists of ready-to-use classes and modules for the computation of AF functions within our system. Each of the functions defined in Section 2 has its own class, identified by the function name. Each of them uses a set of modules, mentioned next, that refer to the statistics supporting a given function and that naturally correspond to the stages of the basic pipeline (see Figure 1 of

the Main Manuscript). The general library is outlined in Figure 1. The main three modules, i.e., `StatisticExtractor`, `StatisticAggregator`, `AFFunctionEvaluator`, correspond to stages 1, 3 and 5 of the basic pipeline and are described first. They are responsible for the layout of the software that is executed, depending on the aggregation strategy that the user specifies. Such a task is transparent to the user and it is described in detail in Section 4. In turn, they specialize with respect to the statistics that is relevant for the AF function to be executed, as we outline next.

- *k-mers Histograms*. For this type of functions, the library highlighted in Figure 1 is specialized as shown in Figure 2. The modules `FastKmerExtractorByBin` and `FastKmerAggregatorByBin` are responsible for *k*-mer collection and aggregation in accordance with the chosen aggregation strategy. The basic tool to perform this task is the `FastKmer` package [13]. In addition, the modules `MashExtractor` and `MashAggregator` are responsible for the execution of *k*-mers Histograms based AF functions using sketches, such as the `MinHash` function. This particular type of specialization is shown in Figure 4. It is worth pointing out that one can use a different package for this task, provided that it is properly interfaced with the mentioned classes. A list of modules implementing different distance/similarity evaluation functions by specializing the `AFFunctionEvaluator` module is reported in Figure 5.
- *spaced words*: For this type of functions, the library highlighted in Figure 1 is specialized as shown in Figure 3. The modules `SwExtractorByBin` and `SwAggregatorByBin` are responsible for *spaced word* collection and aggregation in accordance with the chosen aggregation strategy. The distance evaluation function is implemented by the module `FSWM` reported in Figure 5.

As for stages 2 and 4, they are implemented by two universal modules `StatisticFilter` and `AggregatedStatisticFilter`. They work with any of the statistics supported by our framework, as they operate by evaluating the content and/or the value of each statistic against a user-provided boolean condition encoded using a standard regular expression language.

3.1.2 Examples: Euclidean and Filtered SpacedWords Matches Distance

Assuming that one wants to evaluate the pairwise Euclidean distance between sequences originally contained as separate files in the directory `data`, Listing 1 provides the Java source code required in our framework to perform that task. It is also possible to achieve the same result by writing the following instructions in a properly formatted configuration file and use it to run FADE (see Listing 2). In such a case, no programming skill is required.

With reference to Listing 1, the user defines a new Java application where: (i) the configuration parameters required by the *k*-mers extraction module are provided, as well as the location of the input and of the output files (lines 5-15); (ii) the choice of which modules to use for the different stages of the pipeline is indicated (lines 19-22); (iii) the resulting pipeline is run on the underlying Spark computing cluster and the output is saved on the chosen location (line 26-28).

Assuming now that one wants to evaluate the pairwise Filtered Spaced Word Matches distance [9] between sequences originally contained as separate files in the directory `data`, Listing 3 provides the Java source code required in our framework to perform that task. As for the previous case, is also possible to achieve the same result by writing the following instructions in a properly formatted configuration file and use it to run FADE (see Listing 4). In such a case, no programming skill is required.

With reference to Listing 3, the user defines a new Java application where: (i) the configuration parameters required by the *spaced words* extraction module are provided, as well as the location of the input and of the output files (lines 5-16); (ii) the choice of which modules to use for the different stages of the pipeline is indicated (lines 20-23); (iii) the resulting pipeline is run on the underlying Spark computing cluster and the output is saved on the chosen location (line 27-29).

Listing 1: Code required to evaluate the k -mer Euclidean distance among all sequences stored in the data directory with $k=13$.

```
1 public class KmerTest {
2     public static void main(String[] args) {
3         # Step 1: Input, output and parameters definition
4
5         Configuration conf = new Configuration();
6
7         conf.setInt("k", 13);
8         conf.setInt("x", 3);
9         conf.setInt("m", 6);
10        conf.setInt("slices", 2048);
11
12        Fade f = new Fade(conf);
13
14        f.setInput("data/*.fasta");
15        f.setOutput("distances")
16
17        # Step 2: Modules definition
18
19        f.setStrategy(Strategy.PARTIALAGGREGATION);
20        f.setStatisticExtractor(FastKmerExtractorByBin.class);
21        f.setStatisticAggregator(FastKmerAggregatorByBin.class);
22        f.addAFFunctionEvaluator(Euclidean.class);
23
24        # Step 3: Pipeline execution
25
26        f.compute();
27        f.save();
28        f.close();
29    }
30 }
```

Listing 2: Configuration file equivalent to the code reported in Listing 1.

```
# Input, output and parameters definition
k=13
x=3
m=6
slices=2048
task=distance

input=data/*.fasta
output=distances

# Modules definition
strategy=partial_aggregation
extractor=fade.kmer.fast.FastKmerExtractorByBin
aggregator=fade.kmer.fast.FastKmerAggregatorByBin
evaluator=fade.affunction.Euclidean
```

Listing 3: Code required to evaluate the *spaced word* FSWM distance among all sequences stored in the data directory with pattern=100101000100011001 and threshold=0.

```

1 public class SpacedWordTest {
2     public static void main(String [] args) {
3         # Step 1: Input, output and parameters definition
4
5         Configuration conf = new Configuration ();
6
7         String pattern = "100101000100011001";
8         conf.setString(" pattern", pattern);
9         conf.setInt(" threshold", 0);
10        conf.setInt(" k", pattern.length ());
11        conf.setInt(" slices", 2048);
12
13        Fade f = new Fade(conf);
14
15        f.setInput(" data/*. fasta ");
16        f.setOutput(" distances ");
17
18        # Step 2: Modules definition
19
20        f.setStrategy (Strategy.PARTIALAGGREGATION);
21        f.setStatisticExtractor (SwExtractorByBin.class);
22        f.setStatisticAggregator (SwAggregatorByBin.class);
23        f.addAFFunctionEvaluator (FSWM.class);
24
25        # Step 3: Pipeline execution
26
27        f.compute ();
28        f.save ();
29        f.close ();
30    }
31 }

```

Listing 4: Configuration file equivalent to the code reported in Listing 3.

```

# Input, output and parameters definition
pattern=100101000100011001
k=18
threshold=0
slices=2048
input=data/*. fasta
output=distances
task=distance

# Modules definition
strategy=partial_aggregation
extractor=fade.sw.SwExtractorByBin
aggregator=fade.sw.SwAggregatorByBin
evaluator=fade.affunction.FSWM

```

Figure 1: FADE class diagram.

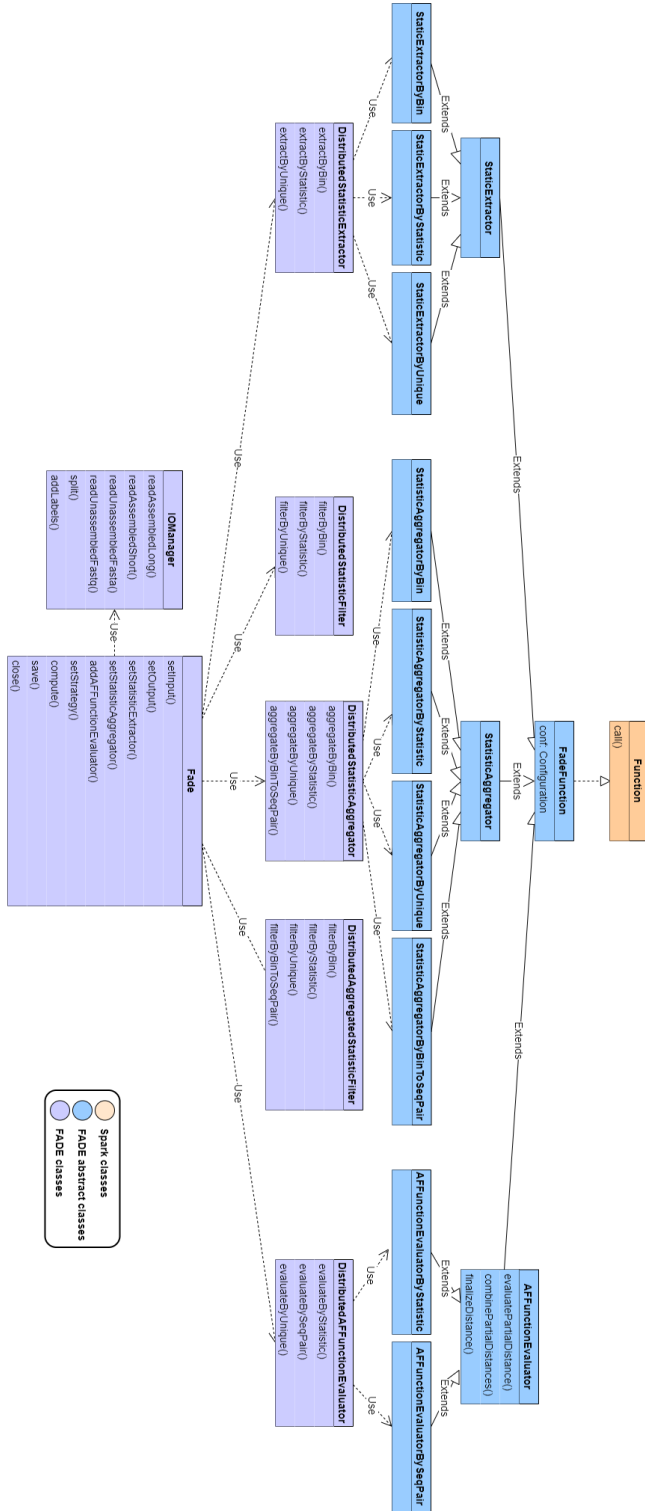


Figure 2: Kmer specialization class diagram.

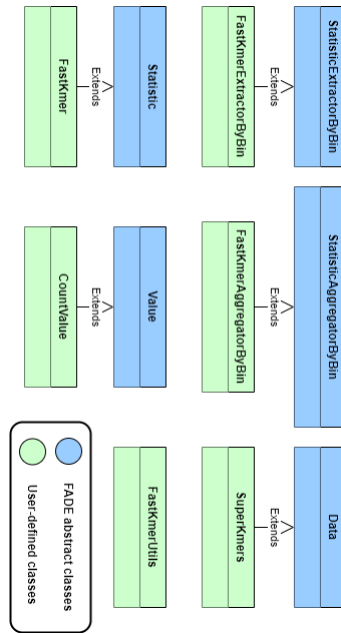


Figure 3: Spacedword specialization class diagram.

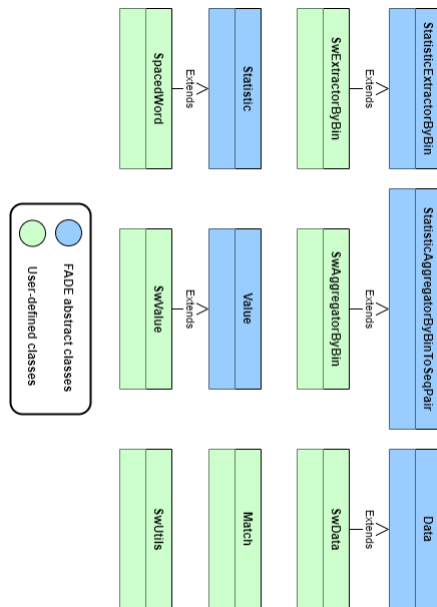


Figure 4: Mash specialization class diagram.

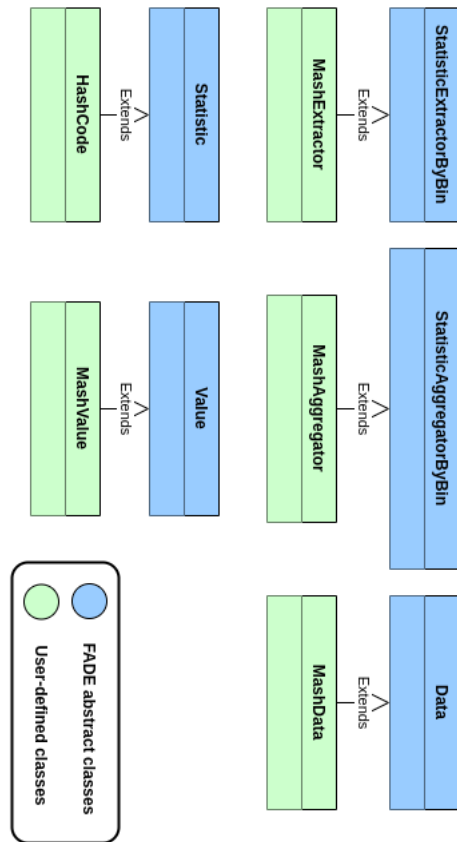
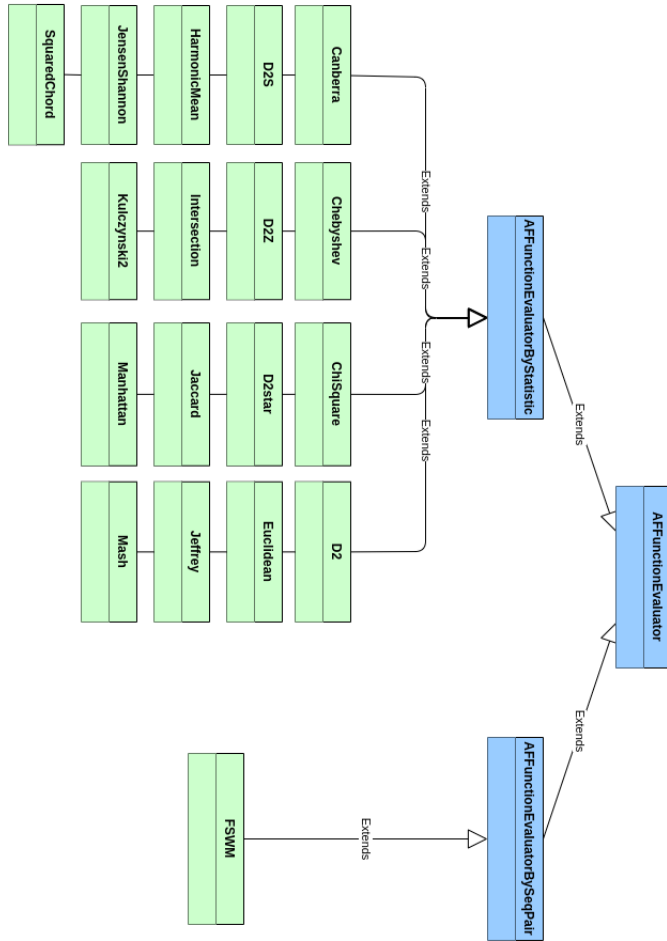


Figure 5: AF function specialization class diagram.



3.1.3 A General Paradigm

Our framework introduces a general paradigm for building AF analysis pipelines over a collection of input sequences, independently of the statistics being collected and of the AF function being used. Assuming that the required modules are already available in our standard primitives library or have been provided using the procedure described in Section 3.2, the paradigm is comprised of the following steps:

1. **Input, output and parameters definition:** the location of the input and of the output files is provided, as well as the definition of the parameters required by the pipeline modules.
2. **Modules definition:** the modules to use for the different stages of the pipeline are defined. The modules to be used in Stage 2 and in Stage 3 to filter partial and aggregated statistics are optional (i.e., if no filtering is required). More modules can be defined for Stage 5 (i.e., if two or more AF functions have to be evaluated).
3. **Pipeline execution:** the assembled pipeline is run over the sequences found in the input location. The AF matrices returned by the pipeline can be either saved on the output location or used, as input, for further processing.

3.2 Possible Extensions

Coherently with the architecture of our framework reported in Figure 1, it is possible to add support for more statistics and/or AF functions by properly deriving and specializing some standard classes, as described next.

3.2.1 Supporting More Statistics.

The user can add support for a target statistic not originally included in the library by extending and specializing a set of standard classes in the following way:

- Inherit and specialize the `Statistic` class to provide a Java representation for the target statistic.
- Inherit and specialize the `Value` class to provide a Java representation for the partial or aggregated value assumed by the target statistic.
- Inherit and specialize the `Data` class to provide a Java representation for a collection of occurrences of the target statistic.
- Inherit and specialize the `StatisticExtractor` class to provide a method to be used for extracting occurrences of the target statistic or collection of occurrences of the target statistic from an input sequence. From one to three versions of this method can be provided according to the partitioning strategy to support (see Section 4).
- Inherit and specialize the `StatisticAggregator` class to provide a method to be used for aggregating occurrences of the target statistic or collection of occurrences of the target statistic. From one to three versions of this method can be provided according to the partitioning strategy to support (see Section 4).

3.2.2 Supporting More AF Functions.

The user can add support for a target AF function not originally included in the library by extending and specializing the `AFFunctionEvaluator` class. This class provides a customizable implementation of a distributed algorithm for evaluating an AF function over a distributed collection \mathcal{C} of statistics. The algorithm is made of the following steps:

1. *Partial Distance/Similarity Evaluation.* Given a pair of statistics in \mathcal{C} , evaluates their partial distance/similarity by means of a user-provided associative and commutative two-arguments function.
2. *Partial Distances/Similarities Combination.* Given a pair of partial distances/similarities, evaluates their combination by means of a user-provided associative and commutative two-arguments function.
3. *Overall Distance/Similarity Finalization.* Given the overall distance/similarity resulting from the combination of all partial distances/similarities, finalizes its value by means of a user-provided function.

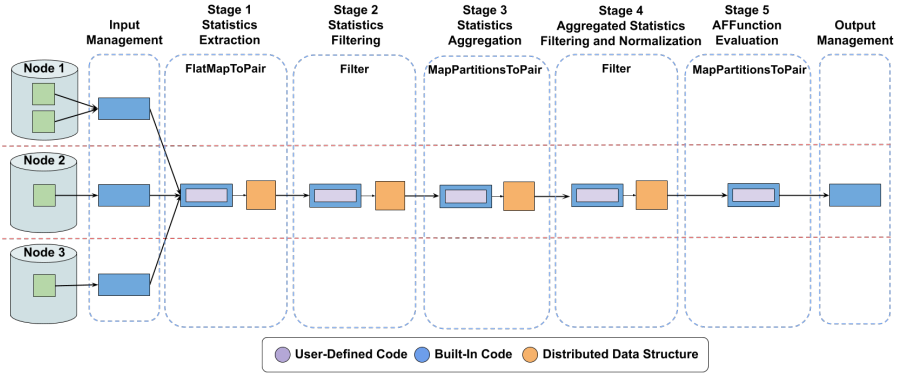


Figure 6: **Total Aggregation Strategy:** layout of the basic pipeline available in our framework when used to process a collection of input sequences on three different nodes of a distributed system by means of the total aggregation partitioning strategy.

Steps 1 and 2 of the algorithm are run in parallel on the different nodes of the distributed system holding a part of \mathcal{C} . Once each node has processed all of its local statistics, the resulting partial distance/similarity is combined with those of other nodes by means of step 2. The algorithm ends with the execution of step 3 on the value resulting from the combination of all partial distances/similarities.

An example of implementation for the *Euclidean* distance is shown in Listing 5. Here the `evaluatePartialAFValue` method is used to evaluate the partial euclidean distance, given a pair of statistics in input, as the square of the difference of their counts. The `combinePartialAFValues` method is used to combine a pair of partial distances/similarities as a unique distance/similarity value using the sum operator. Finally, the `finalizeAFValue` method is used to finalize the computed distance/similarity, by evaluating the square root over the sum of the square of the differences of all statistics.

Listing 5: Java code of the class implementing support for the Euclidean Distance in our framework

```

1 public class Euclidean extends AFFunctionEvaluatorByStatistic {
2     public AFValue evaluatePartialAFValue(Value v1, Value v2) {
3         long count1 = ((CountValue)v1).count;
4         long count2 = ((CountValue)v2).count;
5
6         return new AFValue(Math.pow(count1 - count2, 2));
7     }
8
9     public AFValue combinePartialAFValues(AFValue afv1,
10    AFValue afv2) {
11         return new AFValue(afv1.value + afv2.value);
12     }
13
14    public AFValue finalizeAFValue(AFValue afv) {
15        return new AFValue(Math.sqrt(afv.value));
16    }
17 }

```

4 Architectural Engineering: Details

The three different data partitioning strategies implemented by our basic pipelines are briefly introduced and discussed. Intuitively, they transform the basic logical pipeline introduced in Figure 1 of the Main

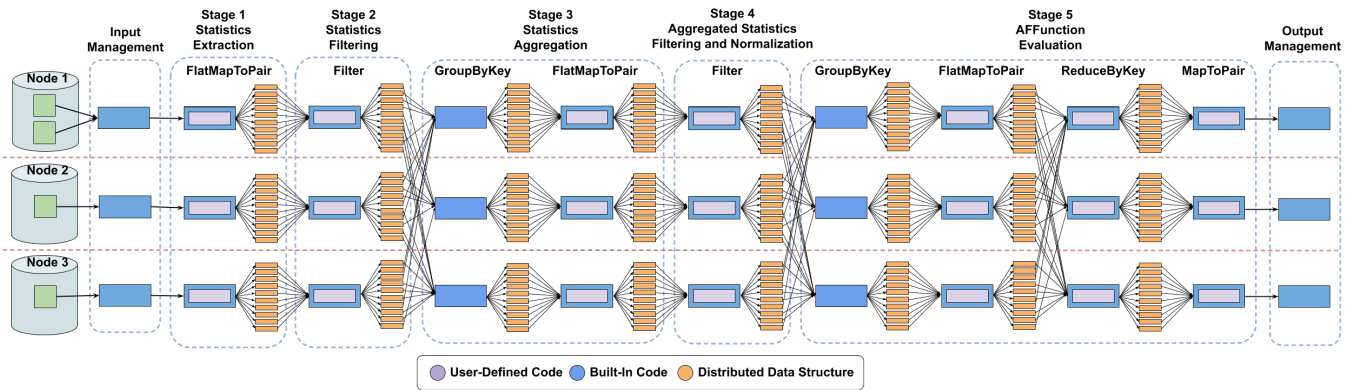


Figure 7: **No Aggregation Strategy:** layout of the basic pipeline available in our framework when used to process a collection of input sequences on three different nodes of a distributed system by means of the no aggregation partitioning strategy.

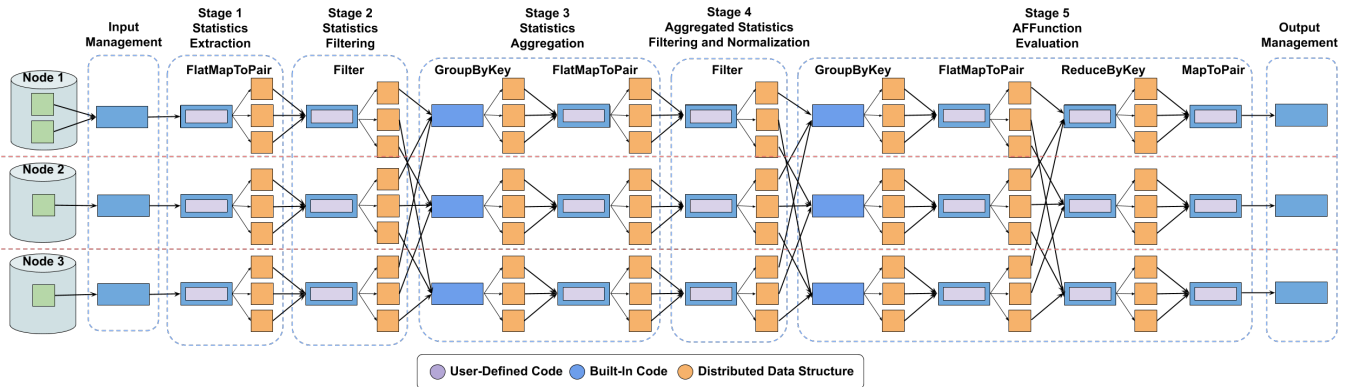


Figure 8: **Partial Aggregation Strategy:** layout of the basic pipeline available in our framework when used to process a collection of input sequences on three different nodes of a distributed system by means of the partial aggregation partitioning strategy.

Manuscript into one of the pipelines exemplified in Figures 6-8 for three nodes. The pipeline selected by the user via the choice of the aggregation strategy is then executed. The transformation from logical to “physical” is transparent to the user.

Strategy 1: Total Aggregation. This strategy implements the basic logical pipeline introduced in Figure 1 of the Main Manuscript as the physical pipeline reported in Figure 6. Among the proposed strategies, this is the one more closely resembling the logical pipeline as all stages are executed in a sequential way, by aggregating all data records on a single node of the distributed system. In the following list, details are provided about the way each stage of the basic pipeline is implemented by this strategy.

- **StatisticExtractor.** It takes a genomic sequence, or a chunk of it, as input, and returns a single list with an associated value. The resulting distributed data structure has the form $\langle unique, (idSeq, V) \rangle$, where *unique* is the id used to group all the records into a single node, *idSeq* is the id of the sequence where the statistic comes from and *V* is its partial record. This data structure contains, at least, a single record associated to the statistics extracted from all the tasks and from all the sequences. This stage is implemented as a Spark `flatMapToPair` transformation.
- **StatisticFilter.** Filters from the list of partial statistics the ones satisfying a provided boolean condition. It is implemented as a Spark `filter` transformation.
- **StatisticAggregator.** Combines all partial statistics according to a user-provided function. The resulting distributed data structure has the form $\langle unique, (idSeq, S, V) \rangle$, where *unique* is the id used to group all the records into a single node, *S* is the extracted statistic, *idSeq* is the id of the sequence where *S* comes from, and *V* is its aggregated value. It is implemented as a Spark `mapPartitionsToPair` transformation.
- **AggregatedStatisticFilter.** Filters from the list of aggregated statistics the ones satisfying a provided boolean condition. It is implemented as a Spark `filter` transformation.
- **AFFunctionEvaluator.** Given the statistics aggregations, it computes the partial distances/similarities between all the pairs of sequences. The partial distances/similarities are then combined and, if necessary, finalized. It is implemented as a Spark `mapPartitionsToPair` transformation.

Strategy 2: No Aggregation. This strategy implements the logical pipeline introduced in Figure 1 of the Main Manuscript as the physical pipeline reported in Figure 7. Differently from the previous strategy, this one aims at maximizing the degree of parallelism as each data record is processed independently of the others in a separate task. In the following list, details are provided about the way each stage of the basic pipeline is implemented by this strategy.

- **StatisticExtractor.** It takes a genomic sequence, or a chunk of it, as input, and returns a list of partial statistics with an associated value. Each partial statistic is encoded as a stand-alone data record, distinct from the others. The resulting distributed data structure has the form $\langle S, (idSeq, V) \rangle$, where *S* is the extracted statistic, *idSeq* is the id of the sequence where *S* comes from, and *V* is its partial value. At the end of the stage, this data structure contains, at least, as many records as the number of distinct statistics extracted from all the tasks and from all the sequences. This stage uses a Spark `flatMapToPair` transformation.
- **StatisticFilter.** As in Strategy 1.
- **StatisticAggregator.** Groups on a node all the occurrences of a same statistic for the same sequence, then combines them according to a customizable function. The resulting distributed data structure has the form $\langle S, (idSeq, V) \rangle$, where *S* is the extracted statistic, *idSeq* is the id of the sequence where *S* comes from, and *V* is its aggregated value. This stage uses a Spark `groupByKey` and a Spark `flatMapToPair` transformations.
- **AggregatedStatisticFilter.** As in Strategy 1.
- **AFFunctionEvaluator.** For each of the aggregated statistics resulting from the previous stage, it computes the partial distances/similarities between all the pairs of sequences. Then, for each pair of sequences, the resulting overall distance/similarity is obtained by combining their partial distances/similarities. It is implemented using, respectively, a Spark `groupByKey` transformation (used

to gather on a same node all the occurrences of a same aggregated statistic), a Spark `flatMapToPair` transformation (used to compute the partial distances/similarities), a Spark `reduceByKey` transformation (used to combine the partial distances/similarities) and a Spark `mapToPair` transformation (used to finalize the distances/similarities).

Strategy 3: Partial Aggregation. This strategy implements the logical pipeline introduced in Figure 1 of the Main Manuscript as the physical pipeline reported in Figure 8. It provides a sort of balance between the previous two strategies as it allows to improve performance by processing large data aggregations in a distributed way.

In the following list, details are provided about the way each stage of the basic pipeline is implemented by this strategy.

- **StatisticExtractor.** It takes a genomic sequence, or a chunk of it, as input, and returns a list of statistics with an associated value and partitioned into bins. The resulting distributed data structure has the form $\langle idBin, (idSeq, V) \rangle$, where $idBin$ is the id of the bin, $idSeq$ is the id of the sequence where the statistics comes from and V is a set of statistics with partial values. It is possible to customize the size of this data structure by modifying the number of the bins used for the partitioning. This stage uses a Spark `flatMapToPair` transformation.
- **StatisticFilter.** As in Strategy 1.
- **StatisticAggregator.** Groups on a node all the partial statistics stored in a same bin, then combines them according to the sequence they come from and using a user-provided function. The resulting distributed data structure has the form $\langle idBin, (idSeq, S, V) \rangle$, where $idBin$ is the id of the bin, $idSeq$ is the id of the sequence where the statistics comes from, S is the aggregated statistic and V is its value. This stage uses a Spark `groupByKey` and `flatMapToPair` transformations.
- **AggregatedStatisticFilter.** As in Strategy 1.
- **AFFunctionEvaluator.** Given the statistics aggregations, it computes the partial distances/similarities between all the pairs of sequences. The partial distances/similarities are first gathered according to the bin they belong to, then are combined and, if necessary, finalized. It is implemented as Spark `groupByKey` transformation, `flatMapToPair` transformation (used to compute the partial distances/similarities), `reduceByKey` transformation (used to combine the partial distances/similarities), `mapToPair` transformation (used to finalize the distances/similarities).

5 A Test of Consistent Significance of an AF Function: Spark Implementation Details

With reference to the Hypothesis Test method outlined in Section 2.5.1 of the Main Manuscript, its Spark modules are summarized in Figure 9. The upper pipeline processes the set S , while repeated executions of the lower pipeline execute the simulation. The module **Randomizer** generates in each run the random datasets and it is executed in a distributed way over the architecture. Finally, for each entry of the AF matrix being tested, a ranking is produced accounting for the corresponding entry in the simulated AF matrices. The module **Ranker** is responsible for that process. Once again, it is executed in a distributed way over the architecture. The final result is a matrix, where each entry contains the rank of the original entry with respect to the simulated ones. The desired confidence level and Bonferroni correction determine which entries pass the test, i.e., for which the null hypothesis can be rejected, and whether the AF matrix passes the test as a whole.

6 Datasets and Hardware.

The datasets used in our study are:

- **E.coli/Shigella.** It contains 27 genomic sequences of the *Bacteria* taxonomic group, having an average length of 4,905,896 bp. For further details, see [3, 19].
- **Mitochondria.** It contains a collection of 25 different genomes of fish species of the suborder *Labroidae*, having an average size of 16,618 bp. For further details, see [6, 19].

- **Plants.** It contains 14 assembled very-large genomes of the *Plants* taxonomic group, having an average size of 337,515,688 bp each. For further details, see [19].
- **Unassembled E.coli strains.** The dataset is the unassembled version of the **E.coli strains** taxonomic group, with an average read length of 150 bp. The sequencing coverages considered are: 0.03125 (29,557 reads), 0.125 (118,266 reads), 1 (946,169 reads). For further details, see [18, 19].
- **Unassembled Plants.** The dataset is the unassembled version of the **Plants** taxonomic group, with an average read length of 150 bp. The sequencing coverages considered are: 1 (30,903,727 reads). For further details, see [19].
- **Yersinia.** It contains 8 genomic sequences of the *Bacteria* taxonomic group. The average sequence length is 4,605,552 bp. For further details, see [3, 19].

Hardware Platform. All the experiments have been performed on a 25 nodes Linux-based cluster, with one node acting as *resource manager* and the remaining nodes being used as workers. The cluster is installed with Hadoop 2.8.1 and the Spark 2.2 software distributions. Each node of this cluster is equipped with one 8-core Intel Xeon E3-12@2.40 GHz processor and 32GB of RAM. Moreover, each node has a 200 GB virtual disk reserved to HDFS, for an overall capacity of about 6 TB.

	<i>l</i>	<i>k</i>	<i>w</i>	<i>L</i>	<i>s</i>
E.coli/Shigella	100	11	12	112	1000
Mitochondria	100	7	12	112	1000
Plants	80	14	12	112	1000
Unass. E.coli (cov. 0.03125)	100	8	12	72	1000
Unass. E.coli (cov. 0.125)	100	9	12	72	1000
Unass. E.coli (cov. 1)	100	11	12	72	1000
Unass. Plants (cov. 1)	80	14	12	72	1000
Yersinia	100	11	12	112	1000

Table 1: Outline of the parameters used in our experiment for assessing the statistical significance of several AF functions on the considered datasets. *l* denotes the number of runs used for the Monte Carlo simulation. *k* denotes the length of the *k*-mers used for running AF functions based on histogram statistics. *w* and *L* are respectively the length and the weight of the pattern used for running spaced words AF functions. *s* is the size of the sketches used when running the MinHash approximation algorithm.

7 AF Functions Parameters and Monte Carlo Simulation Iterations

For each experiment, the number of Monte Carlo Simulations has been chosen so as to obtain the results of the test in a reasonable amount of time, given the size of the dataset. The criteria for the choice of relevant parameters for AF functions is provided next, while a summary of all the parameters used in our experiments is available in Table 1.

7.1 Histogram Statistics: choice of *k*

As well argued in [10], for histogram-based AF functions, the choice of the *k*-mers *k* is crucial for the success of those methods. It is a heuristic choice. The one that seems to work best is the one shown in equation (19).

$$k = \lceil \log_4 \left(\frac{1}{|S|} \sum_{i \in S} \text{len}(i) \right) \rceil - 1 \quad (19)$$

in which *n* is the total number of sequences. It is also to be remarked that the choice of *k* should be also in a range that preserves the “information theoretic content” of the sequences to be analyzed. To the

best of our knowledge, two closely related approaches have been proposed in the literature [7, 16]. Here we follow the approach in [7]. Then, the value of k for the AF histogram-based function to be used for the analysis is chosen as the minimum between the k' returned by this approach and the k coming out of equation (19).

As for the choice of the sketch size to be used when running the MinHash approximation algorithm Mash, we considered the default value of 1,000 [12].

7.2 Word matches: spaced word methods

When considering *spaced words* methods, there are two parameters that need to be fixed: the length of the pattern \mathcal{P} and its weight w . In our experiments with assembled genomes, we have used the same default values considered in [9], i.e., a pattern with $w = 12$ and 100 don't care positions. Instead, when analyzing collections of reads we resorted to the default values used by the authors of [8] to analyze short sequencing reads: a pattern with $w = 12$ and 60 don't care positions.

8 Additional Figures

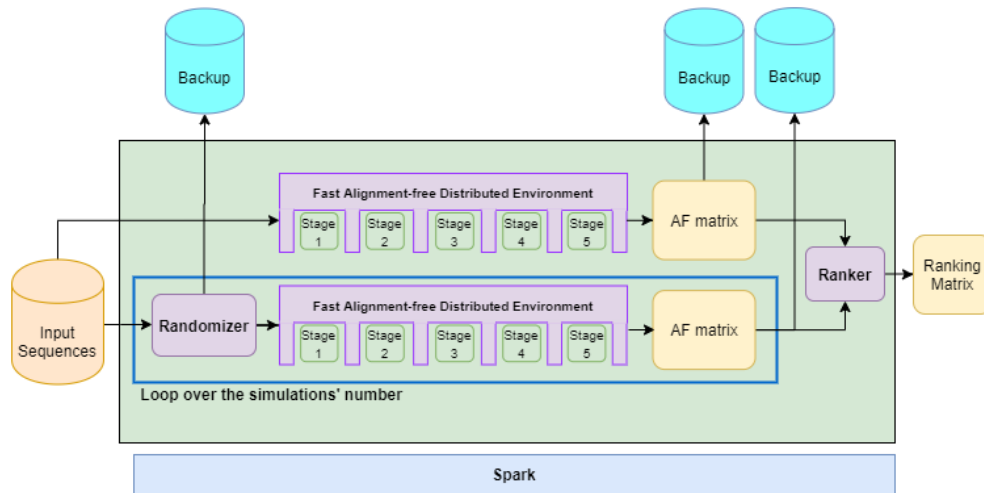


Figure 9: A layout of the architecture of the pipeline for the fast hypothesis testing of alignment-free algorithms. It repeatedly uses, as subroutine, the basic pipeline for the fast computation of AF functions. Backup primitives are used to save the progresses of a test so as to recover an interrupted computation or to allow a very long execution to be broken in shorter parts.

Yersinia (assembled)			
	q=1	q=7	q=10
Canberra	100	100	100
Chebyshev	0	0	0
ChiSquare	100	100	100
D2	100	100	100
D2S	100	100	100
D2Z	100	100	100
D2*	100	100	100
Euclidean	89.3	89.3	89.3
FSWM	100	100	100
Harmonic Mean	0	0	0
Intersection	0	0	0
Jaccard	100	100	100
Jeffreys	96.4	96.4	96.4
Jensen Shannon	96.4	96.4	96.4
Kulczynski2	78.6	57.1	53.6
Manhattan	100	100	100
Mash	78	86	83.6
Squared Chord	100	100	100

Shigella/E.coli (assembled)			
	q=1	q=7	q=10
Canberra	100	100	100
Chebyshev	0.6	0.6	0.9
ChiSquare	100	100	100
D2	100	100	100
D2S	100	100	100
D2Z	100	89.5	77.2
D2*	100	100	100
Euclidean	56.1	56.1	56.1
FSWM	100	100	100
Harmonic Mean	100	97.2	94.9
Intersection	100	100	100
Jaccard	0	0	3.1
Jeffreys	100	95.7	94
Jensen Shannon	100	100	100
Kulczynski2	100	100	100
Manhattan	100	100	100
Mash	19.9	99.7	100
Squared Chord	100	100	100

Mitochondria (assembled)			
	q=1	q=7	q=10
Canberra	100	100	100
Chebyshev	3.3	3	3.7
ChiSquare	100	100	99.3
D2	100	100	99
D2S	100	99.7	81.7
D2Z	100	99.7	98.7
D2*	100	91.3	82
Euclidean	89.3	85.7	84.3
FSWM	100	100	100
Harmonic Mean	98.7	70.7	48
Intersection	65	40	35.3
Jaccard	27.3	28	43.3
Jeffreys	99.3	96	95
Jensen Shannon	99.7	96.3	95
Kulczynski2	100	99.3	99
Manhattan	100	100	99.3
Mash	78	86	83.7
Squared Chord	100	100	99.3

E.coli (unassembled, coverage=1)			
	q=1	q=7	q=10
Canberra	100	100	100
Chebyshev	0	0	0
ChiSquare	58.4	21.9	16.3
D2	100	100	100
D2S	100	100	100
D2Z	100	79.8	65
D2*	100	100	100
Euclidean	1.5	1.5	1.5
FSWM	100	100	100
Harmonic Mean	100	96.3	85
Intersection	100	15.5	2.5
Jaccard	0	0	0
Jeffreys	71.9	28.1	15.3
Jensen Shannon	82.5	34.7	17.5
Kulczynski2	13.3	3.7	3.7
Manhattan	66.5	28.1	18.2
Mash	0	0	0.9
Squared Chord	41.1	20.7	16.3

Plants (assembled)			
	q=1	q=7	q=10
Canberra	22	22	27.5
Chebyshev	0	0	0
ChiSquare	6.6	6.6	6.6
D2	98.9	97.8	97.8
D2S	14.3	7.7	7.7
D2Z	1.1	1.1	1.1
D2*	100	100	100
Euclidean	0	0	0
FSWM	100	100	100
Harmonic Mean	19.8	18.7	18.7
Intersection	42.9	29.7	27.5
Jaccard	1.1	3.3	3.3
Jeffreys	2.2	2.2	2.2
Jensen Shannon	3.3	2.2	2.2
Kulczynski2	2.2	2.2	2.2
Manhattan	13.2	11	11
Mash	2.2	2.2	2.2
Squared Chord	6.6	6.6	6.6

Plants (unassembled, coverage=1)			
	q=1	q=7	q=10
Canberra	22	53.8	54.9
Chebyshev	0	0	0
ChiSquare	6.6	4.4	4.4
D2	98.9	97.8	95.6
D2S	14.3	7.7	7.7
D2Z	1.1	1.1	0
D2*	100	100	100
Euclidean	0	0	0
FSWM	100	100	100
Harmonic Mean	19.8	17.6	17.6
Intersection	42.9	5.5	3.3
Jaccard	1.1	1.1	1.1
Jeffreys	0	0	0
Jensen Shannon	0	0	0
Kulczynski2	1.1	0	0
Manhattan	13.2	11	11
Mash	1.1	1.1	1.1
Squared Chord	6.6	2.2	2.2

E.coli (unassembled, coverage=0.125)			
	q=1	q=7	q=10
Canberra	0	0	0.2
Chebyshev	0	0	0
ChiSquare	19.2	17.7	7.9
D2	100	100	100
D2S	100	100	99.7
D2Z	100	100	93.3
D2*	100	100	100
Euclidean	3.9	3.4	3
FSWM	100	100	100
Harmonic Mean	79.8	64.8	52.7
Intersection	100	98.8	92.6
Jaccard	0	0	0
Jeffreys	93.1	63.6	59.1
Jensen Shannon	93.4	66.3	59.1
Kulczynski2	68	55.7	42.9
Manhattan	71.7	42.1	25.9
Mash	0	0	0
Squared Chord	0	0	0.2

E.coli (unassembled, coverage=0.03125)			
	q=1	q=7	q=10
Canberra	0	0	0.2
Chebyshev	0	0	0
ChiSquare	19	16.7	6.4
D2	100	100	96.8
D2S	100	100	100
D2Z	100	100	100
D2*	100	100	100
Euclidean	5.9	4.9	4.4
FSWM	100	100	100
Harmonic Mean	74.9	61.3	44.3
Intersection	100	98.5	91.9
Jaccard	0	0	0
Jeffreys	100	87.2	74.4
Jensen Shannon	100	87.9	75.4
Kulczynski2	88.2	58.9	49
Manhattan	77.8	60.1	27.8
Mash	0	0	0
Squared Chord	0	0	0

Figure 10: Hypothesis Test complete results for the different AF functions considered in this research when executed on three different datasets with $q = 1, 7, 10$ and with significance level set to 1%.

Figure 11: Robinson Fould corruption graph of Mitochondria dataset, varying AF function and q . Each graph reports the distance of the phylogenetic trees, obtained using the Neighbor-joining and the UPGMA methods, from the gold standard, measured using the Robinson Fould metric, for the considered AF function and q , and as the number of noisy entries grows.

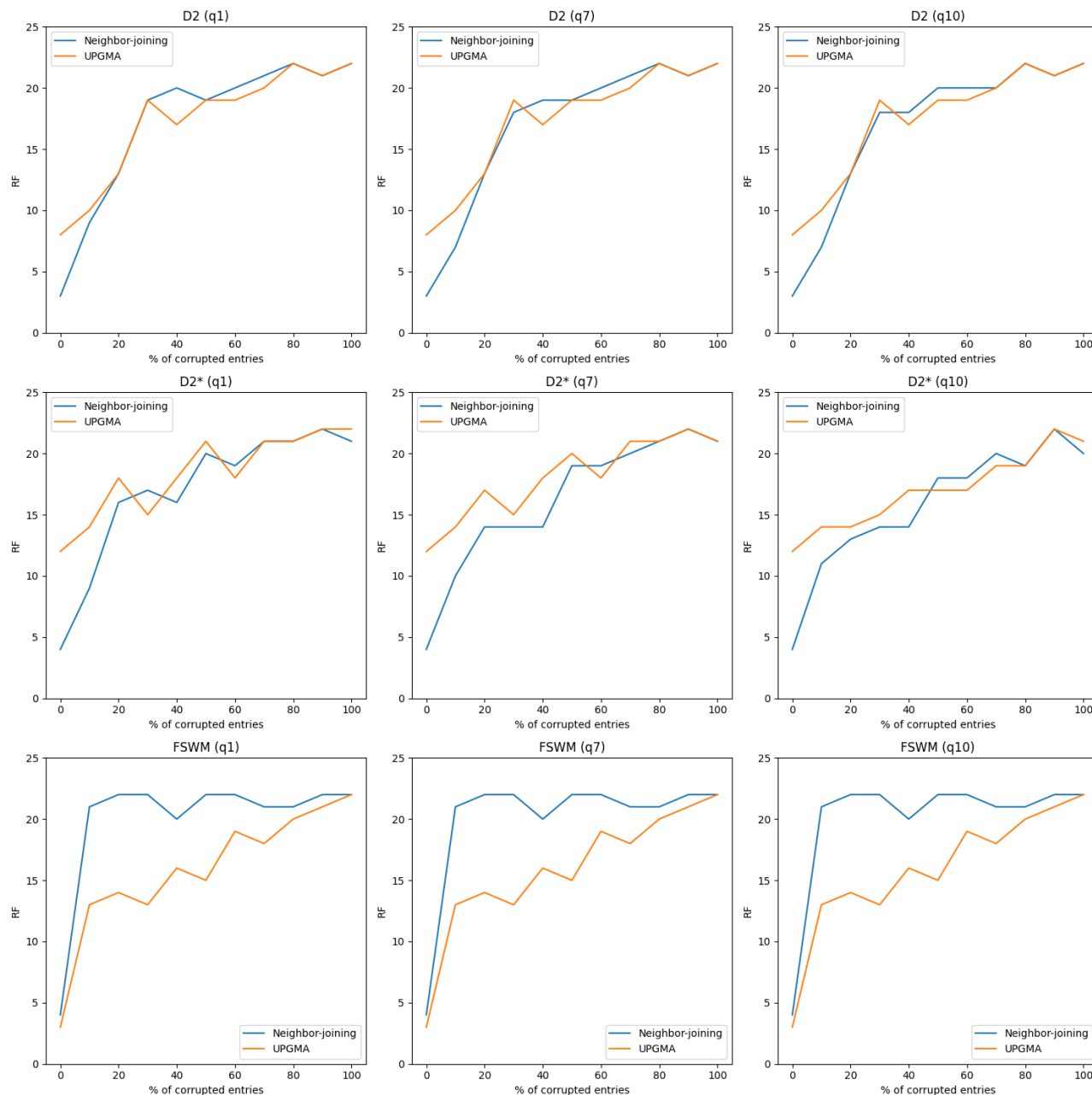


Figure 12: Robinson Fould corruption graph of Shigella dataset, varying AF function and q . Each graph reports the distance of the phylogenetic trees, obtained using the Neighbor-joining and the UPGMA methods, from the gold standard, measured using the Robinson Fould metric, for the considered AF function and q , and as the number of noisy entries grows.

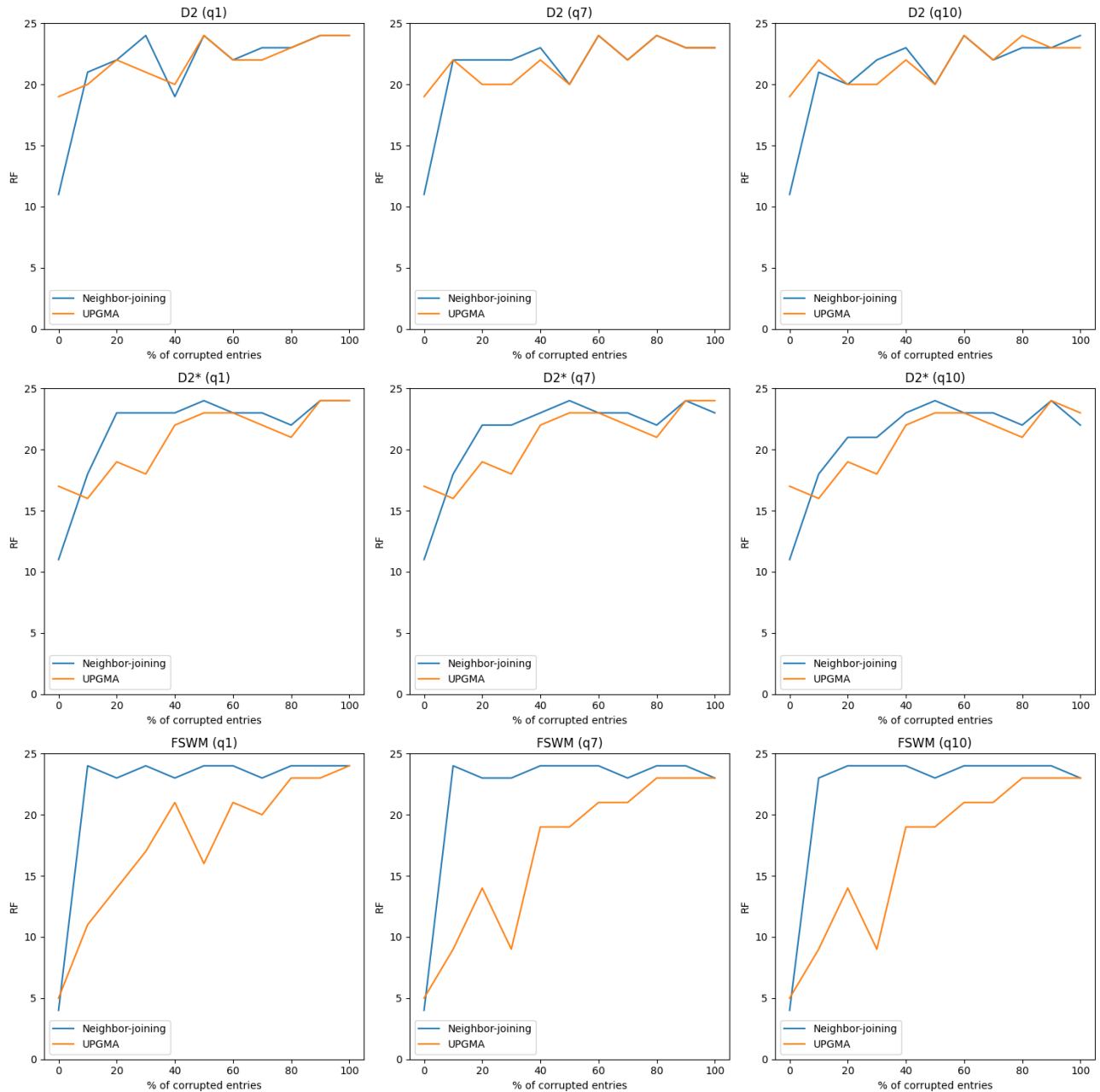


Figure 13: Robinson Fould corruption graph of Yersinia dataset, varying AF function and q . Each graph reports the distance of the phylogenetic trees, obtained using the Neighbor-joining and the UPGMA methods, from the gold standard, measured using the Robinson Fould metric, for the considered AF function and q , and as the number of noisy entries grows.

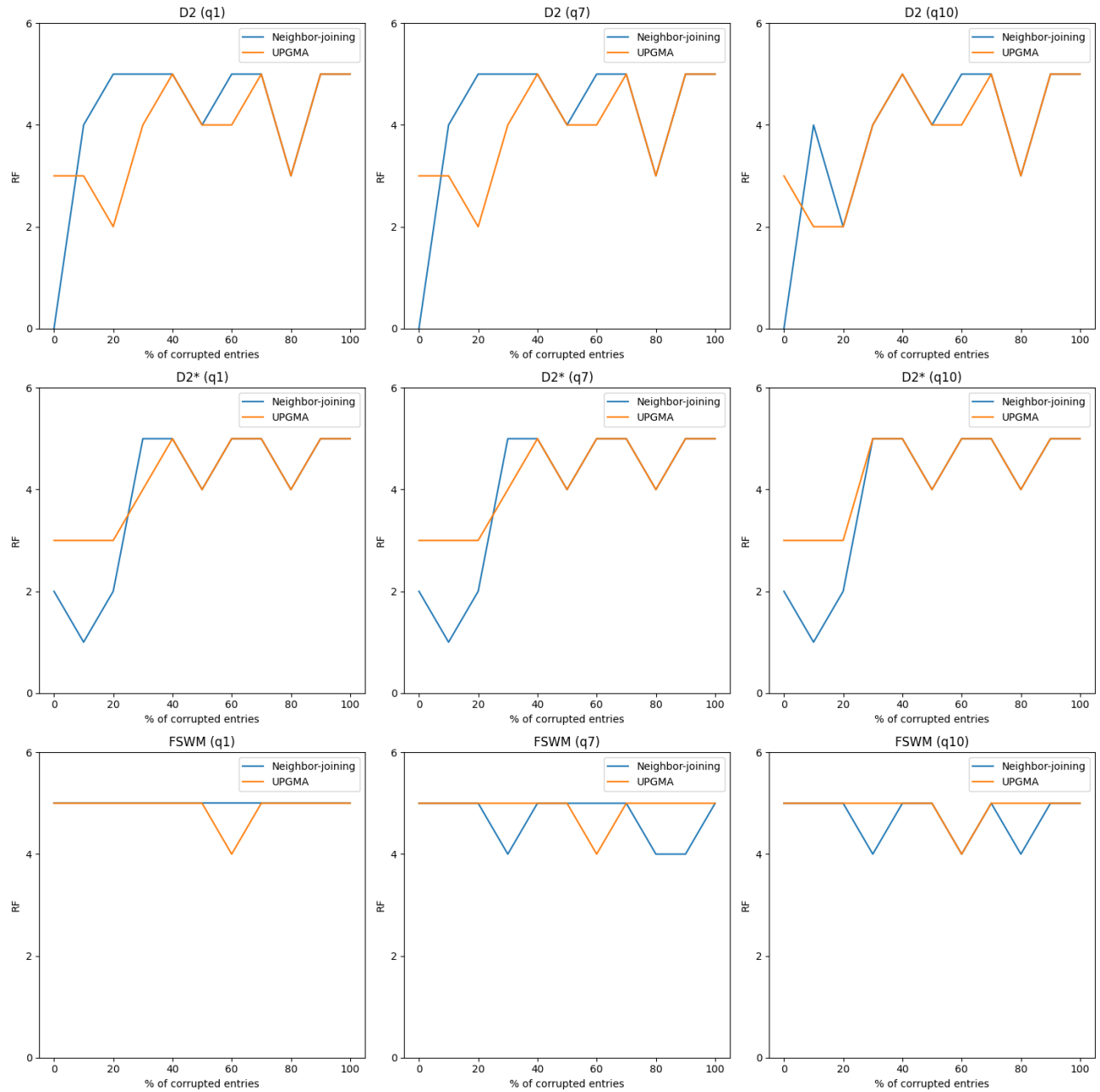


Figure 14: Matching Cluster metric corruption graph of Mitochondria dataset, varying AF function and q . Each graph reports the distance of the phylogenetic trees, obtained using the Neighbor-joining and the UPGMA methods, from the gold standard, measured using the MCM metric, for the considered AF function and q , and as the number of noisy entries grows.

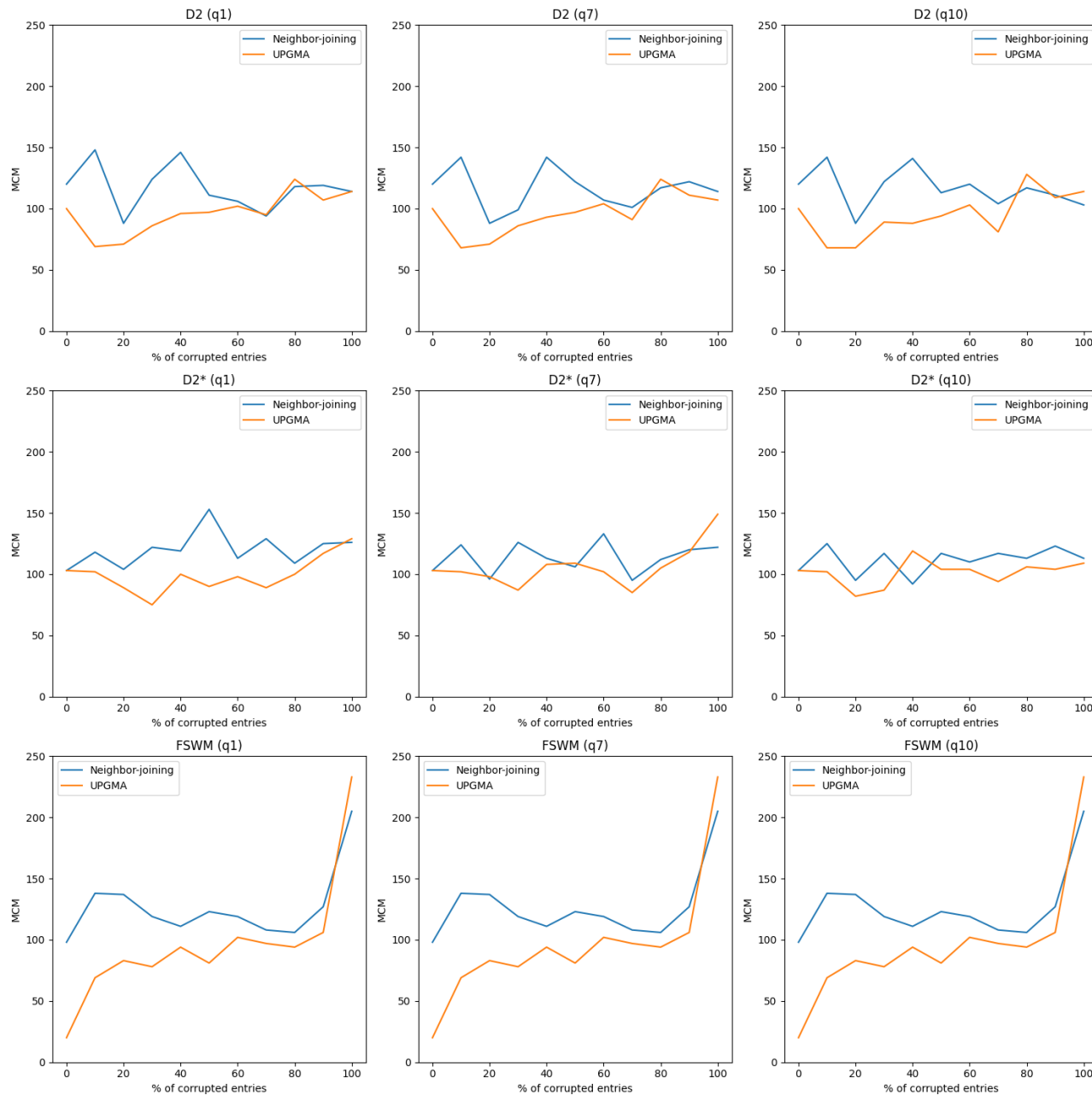


Figure 15: Matching Cluster metric corruption graph of Shigella dataset, varying AF function and q . Each graph reports the distance of the phylogenetic trees, obtained using the Neighbor-joining and the UPGMA methods, from the gold standard, measured using the MCM metric, for the considered AF function and q , and as the number of noisy entries grows.

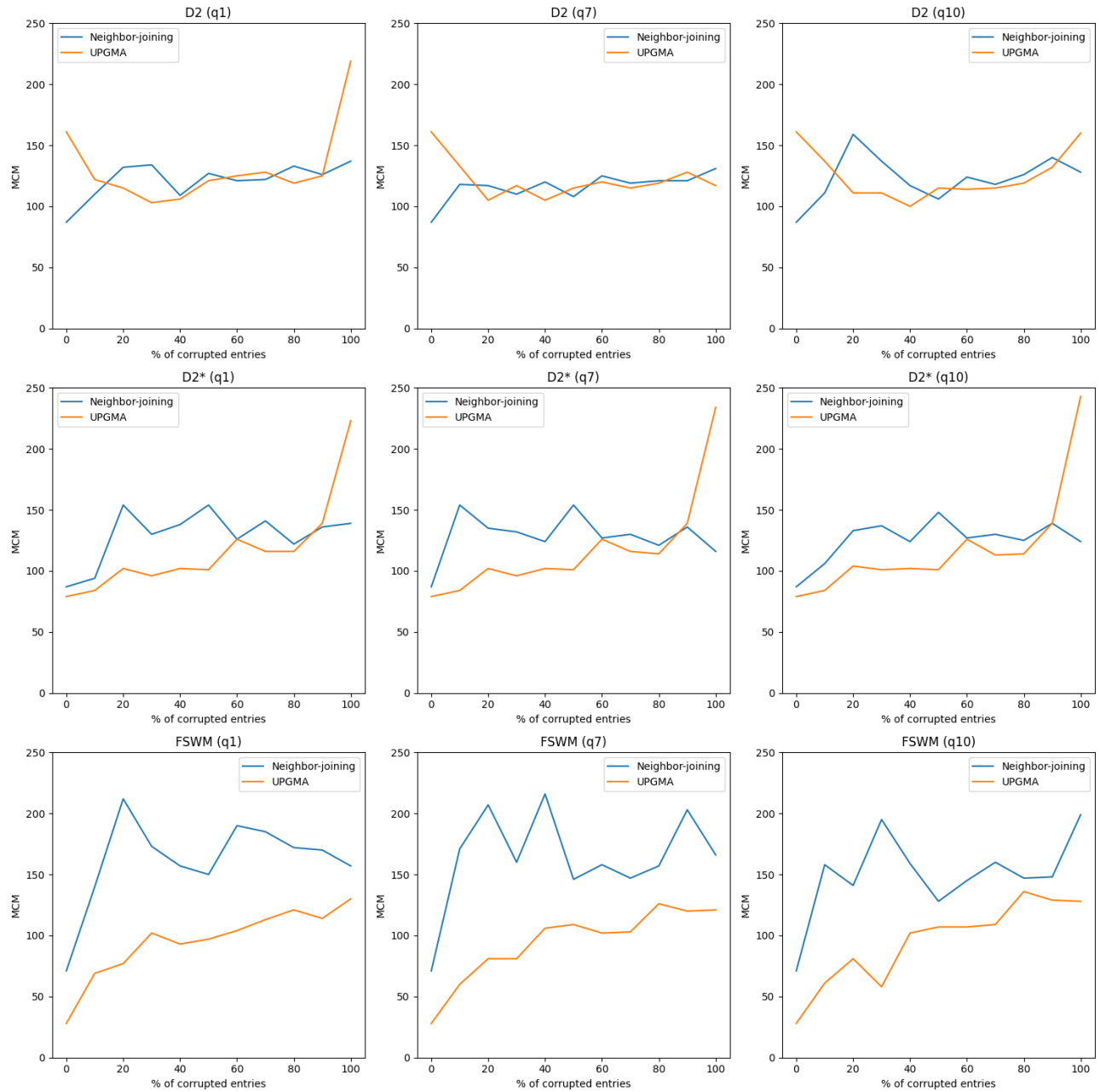
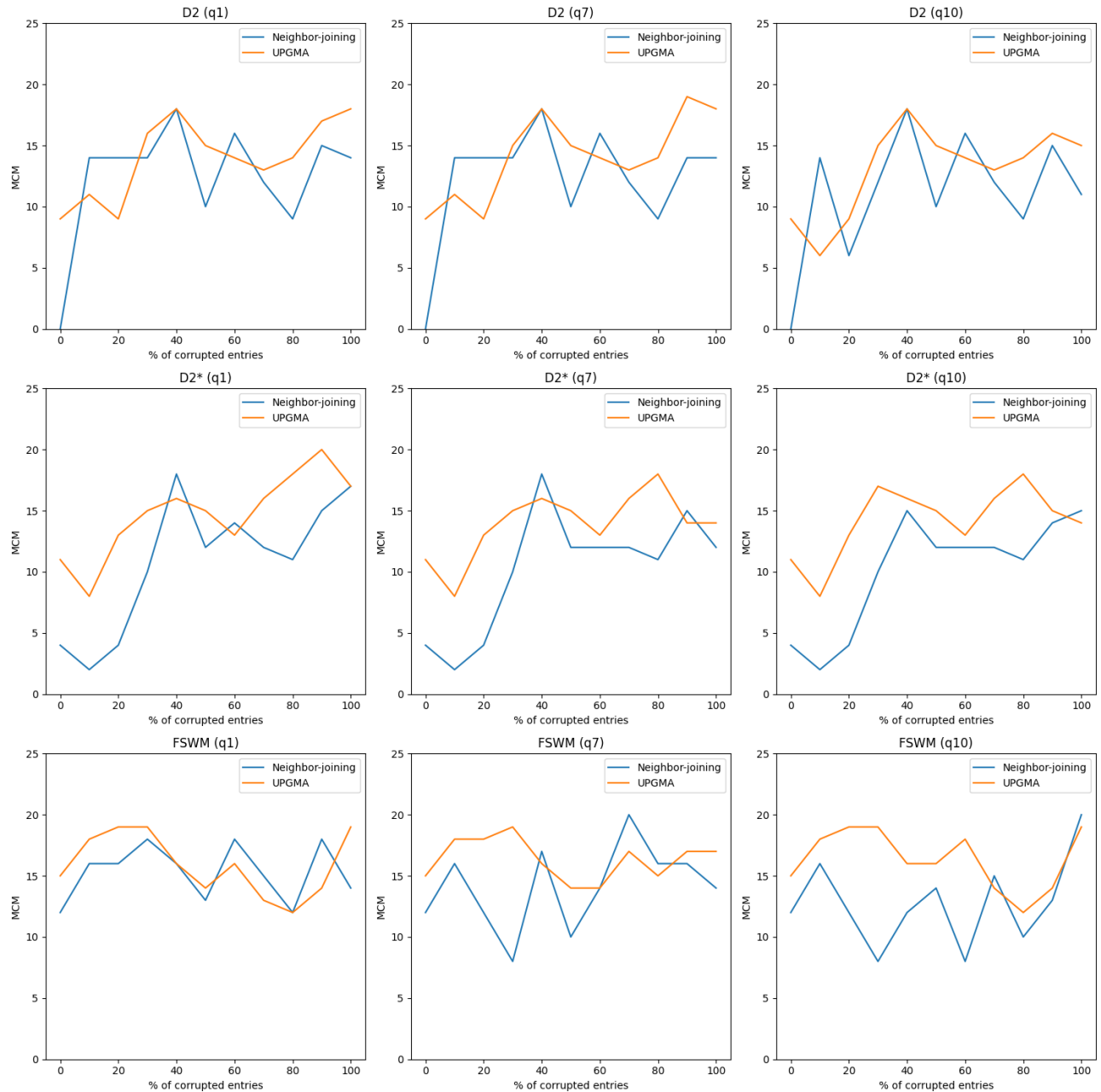


Figure 16: Matching Cluster metric corruption graph of Yersinia dataset, varying AF function and q. Each graph reports the distance of the phylogenetic trees, obtained using the Neighbor-joining and the UPGMA methods, from the gold standard, measured using the MCM metric, for the considered AF function and q, and as the number of noisy entries grows.



References

- [1] Apache Software Foundation. Apache Hadoop. (Available from: <http://hadoop.apache.org>), 2006.
- [2] Apache Software Foundation. Apache Spark. (Available from: <http://spark.apache.org>), 2016.
- [3] G. Bernard, C. X. Chan, and M. A. Ragan. Alignment-free microbial phylogenomics under scenarios of sequence divergence, genome rearrangement and lateral genetic transfer. *Scientific reports*, 6:28970, 2016.
- [4] F. Chiaromonte, V. Yap, and W. Miller. Scoring pairwise genomic sequence alignments. In *Biocomputing 2002*, pages 115–126. World Scientific, 2001.
- [5] J. Dean and S. Ghemawat. MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 51:107–113, 2008.
- [6] C. Fischer, S. Koblmüller, C. Güllly, C. Schlötterer, C. Sturmbauer, and G. G. Thallinger. Complete mitochondrial dna sequences of the threadfin cichlid (*petrochromis trewavasae*) and the blunthead cichlid (*tropheus moorii*) and patterns of mitochondrial genome evolution in cichlid fishes. *PLoS One*, 8(6):e67048, 2013.
- [7] R. Giancarlo, S. E. Rombo, and F. Utro. Epigenomic k-mer dictionaries: Shedding light on how sequence composition influences nucleosome positioning *in vivo*. *Bioinformatics*, 31:2939–2946, 2015.
- [8] A. K. Lau, C.-A. Leimeister, S. Dörrer, C. Bleidorn, and B. Morgenstern. Read-spam: assembly-free and alignment-free comparison of bacterial genomes with low sequencing coverage. *BioRxiv*, page 550632, 2019.
- [9] C.-A. Leimeister, S. Sohrabi-Jahromi, and B. Morgenstern. Fast and accurate phylogeny reconstruction using filtered spaced-word matches. *Bioinformatics*, 33:971–979, 2017.
- [10] B. B. Luczak, B. T. James, and H. Z. Girgis. A survey and evaluations of histogram-based statistics in alignment-free sequence comparison. *Briefings in Bioinformatics*, 20(4):1222–1237, 12 2017.
- [11] N. A. Lynch. *Distributed algorithms*. Morgan Kaufmann, 1996.
- [12] 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. *Genome Biology*, 17:132, 2016.
- [13] U. F. Petrillo, M. Sorella, G. Cattaneo, R. Giancarlo, and S. E. Rombo. Analyzing big datasets of genomic sequences: fast and scalable collection of k-mer statistics. *BMC bioinformatics*, 20(4):138, 2019.
- [14] S. Sarmashghi, K. Bohmann, M. T. P. Gilbert, V. Bafna, and S. Mirarab. Skmer: assembly-free and alignment-free sample identification using genome skims. *Genome biology*, 20(1):1–20, 2019.
- [15] L. Shepp. Normal functions of normal random variables. *Siam Review*, 4(3):255, 1962.
- [16] G. E. Sims and S.-H. Kim. Whole-genome phylogeny of *Escherichia coli*/Shigella group by feature frequency profiles (FFPs). *Proceedings of the National Academy of Sciences*, 108(20):8329–8334, 2011.
- [17] K. Song, J. Ren, G. Reinert, M. Deng, M. S. Waterman, and F. Sun. New developments of alignment-free sequence comparison: measures, statistics and next-generation sequencing. *Briefings in Bioinformatics*, 15:343–353, 2013.
- [18] H. Yi and L. Jin. Co-phylog: an assembly-free phylogenomic approach for closely related organisms. *Nucleic Acids Research*, 41:e75, 2013.
- [19] A. Zielezinski, H. Z. Girgis, G. Bernard, C.-A. Leimeister, K. Tang, T. Dencker, A. K. Lau, S. Röhling, J. J. Choi, M. S. Waterman, M. Comin, S.-H. Kim, S. Vinga, J. S. Almeida, C. X. Chan, B. T. James, F. Sun, B. Morgenstern, and W. M. Karlowski. Benchmarking of alignment-free sequence comparison methods. *Genome Biology*, 20(1):144, 2019.