

Benchmarking of alignment-free sequence comparison methods

Andrzej Zielezinski¹, Hani Z. Girgis², Guillaume Bernard³, Chris-Andre Leimeister⁴, Kujin Tang⁵, Thomas Dencker⁴, Anna K. Lau⁴, Sophie Röhling⁴, JaeJin Choi^{6,7,8,9}, Michael S. Waterman^{5,10}, Matteo Comin¹¹, Sung-Hou Kim^{6,7,8,9}, Susana Vinga^{12,13}, Jonas S. Almeida¹⁴, Cheong Xin Chan¹⁵, Benjamin T. James², Fengzhu Sun^{5,10}, Burkhard Morgenstern⁴, Wojciech M. Karlowski^{1*}

Authors affiliations

¹ Department of Computational Biology, Faculty of Biology, Adam Mickiewicz University in Poznan, Umultowska 89, 61-614 Poznan, Poland

² Tandy School of Computer Science, The University of Tulsa, 800 South Tucker Drive, Tulsa, OK 74104, USA

³ Sorbonne Université, UMR 7205 ISYEB, Paris 75005 France

⁴ University of Göttingen, Institute of Microbiology and Genetics, Department of Bioinformatics, Goldschmidtstr. 1, 37077 Göttingen, Germany

⁵ Quantitative and Computational Biology Program, Department of Biological Sciences, University of Southern California, CA 90089, USA

⁶ Department of Chemistry, University of California, Berkeley, CA 94720

⁷ Molecular Biophysics & Integrated Bioimaging Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720

⁸ Department of Integrated Omics for Biomedical Sciences, Yonsei University, Seoul 03722, Republic of Korea

⁹ Korea Research Institute of Bioscience and Biotechnology, Daejeon 34141, Republic of Korea

¹⁰ Centre for Computational Systems Biology, School of Mathematical Sciences, Fudan University, Shanghai, 200433, China

¹¹ Department of Information Engineering, University of Padova, Padova, Italy

¹² INESC-ID, Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais 1, 1049-001 Lisbon, Portugal

¹³ IDMEC, Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais 1, 1049-001 Lisbon, Portugal

¹⁴ National Cancer Institute (NIH/NCI), Division of Epidemiology and Genetics (DCEG)

¹⁵ Institute for Molecular Bioscience, and School of Chemistry and Molecular Biosciences, The University of Queensland, Brisbane, QLD 4072, Australia

* corresponding author: wmk@amu.edu.pl

ABSTRACT

Alignment-free (AF) sequence comparison is attracting persistent interest driven by data-intensive applications. Hence, many AF procedures have been proposed in recent years, but a lack of a clearly defined benchmarking consensus hampers their performance assessment. Here, we present a community resource (<http://afproject.org>) to establish standards for comparing AF methods across different areas of sequence-based research. We characterize 74 AF methods available in 24 software tools for five research applications, namely, protein sequence classification, gene tree inference, regulatory element detection, genome-based phylogenetic inference and reconstruction of species trees under horizontal gene transfer and recombination events. The interactive web service allows researchers to explore the performance of AF tools relevant to their data types and analytical goals. It also allows method developers to assess their own algorithms and compare them with the current state-of-the-art tools, accelerating the development of new, more accurate AF solutions.

Keywords: alignment-free, sequence comparison, benchmark, whole-genome phylogeny, horizontal gene transfer, web service

INTRODUCTION

Comparative analysis of DNA and amino acid sequences is of fundamental importance in biological research, particularly in molecular biology and genomics. It is the first and key step in the study of molecular evolutionary analysis, gene function and regulatory region

prediction, sequence assembly, homology searching, molecular structure prediction, gene discovery and protein structure-function relationships analysis. Traditionally, sequence comparison was based on pairwise or multiple sequence alignment (MSA). Software tools for sequence alignment, such as BLAST¹ and CLUSTAL², are the most widely used bioinformatics methods. Although alignment-based approaches generally remain the references for sequence comparison, MSA-based methods do not scale with the very large data sets that are available today. Additionally, alignment-based techniques have been shown to be inaccurate in scenarios of low sequence identity³ (e.g., gene regulatory sequences^{4,5} and distantly related protein homologs^{3,6}). Moreover, alignment algorithms assume that the linear order of homologies is preserved within the compared sequences, so these algorithms cannot be directly applied in the presence of sequence rearrangements (e.g., recombination and protein domain swapping⁷) or horizontal transfer⁸ in cases where large-scale sequence data sets are processed, e.g., for whole-genome phylogenetics⁹. Therefore, as an alternative to sequence alignment, many so-called alignment-free (AF) approaches to sequence analysis have been developed³, with the earliest works dating back to the mid 1970s¹⁰, although the concept of the alignment-independent sequence comparison gained increased attention only in the beginning of the 2000s¹¹. Most of these methods are based on word statistics or word comparison, and their scalability allows them to be applied to much larger data sets than conventional MSA-based methods.

A wide array of AF approaches to sequence comparison have been developed. These approaches include methods based on word or k -mer counts¹²⁻¹⁶, the length of common substrings¹⁷⁻²⁰, micro-alignments²¹⁻²⁵, sequence representations based on chaos theory^{26,27},

moments of the positions of the nucleotides²⁸, Fourier transformations²⁹, information theory³⁰ and iterated-function systems^{30,31}. Currently, the most widely used AF approaches are based on *k*-mer counts³². These methods are very diverse, providing a variety of statistical measures that are implemented across different software tools^{3,33–35} (**Table 1**). Many *k*-mer methods work by projecting each of the input sequences into a feature space of *k*-mer counts, where sequence information is transformed into numerical values (e.g., *k*-mer frequencies) that can be used to calculate distances between all possible sequence pairs in a given data set.

Despite the extensive progress achieved in the field of AF sequence comparison³, developers and users of AF methods face several difficulties. New AF methods are usually evaluated by their authors, and the results are published together with these new methods. Therefore, it is difficult to compare the performance of these tools since they are based on inconsistent evaluation strategies, varying benchmarking data sets and variable testing criteria. Moreover, new methods are usually evaluated with relatively small data sets selected by their authors, and they are compared with a very limited set of alternative AF approaches. As a consequence, the assessment of new algorithms by individual researchers presently consumes a substantial amount of time and computational resources, compounded by the unintended biases of partial comparison. To date, no comprehensive benchmarking platform has been established for AF sequence comparison to select algorithms for different sequence types (e.g., genes, proteins, regulatory elements, or genomes) under evolutionary scenarios (e.g., high mutability or horizontal gene transfer (HGT)). As a result, users of these methods cannot easily identify appropriate tools for the problems at hand and are instead often confused by a plethora of existing programs of unclear applicability to their study. Finally, as for other

software tools in bioinformatics, the results of most AF tools strongly depend on the specified parameter values. For many AF methods, the word length k is a crucial parameter. Note, however, that words are used in different ways by different AF methods, so there can be no universal optimal word length k for all AF programs. Instead, different optimal word lengths have to be identified for the different methods. In addition, optimal parameter values may depend on the data-analysis task at hand, for instance, whether a set of protein sequences is to be grouped into protein families or superfamilies.

To address these problems, we developed AFproject (<http://afproject.org>), a publicly available web-based service for comprehensive and unbiased evaluation of AF tools. The service is based on eight well-established and widely used reference sequence data sets as well as four new data sets and can be used to comprehensively evaluate AF methods under five different sequence analysis scenarios: protein sequence classification, gene tree inference, regulatory sequence identification, genome-based phylogenetics and HGT (**Table 2**). To evaluate the existing AF methods with these data sets, we asked the developers of 24 AF tools to run their software on our data sets or to recommend suitable input parameter values appropriate for each data set. In total, our study involved 10,181 program runs, resulting in 1,020,463,773 pairwise sequence comparisons (**Table 1**; [Supplementary Table 1](#)). All benchmarking results are stored and can be downloaded, reproduced and inspected with the AFproject website. Thus, any future evaluation results can be seamlessly compared to the existing ones obtained using the same reference data sets with precisely defined software parameters. By providing a way to automatically include new methods and to disseminate their results publicly, we aim to

maintain an up-to-date and comprehensive assessment of state-of-the-art AF tools, allowing contributions and continuous updates by all developers of AF-based methods.

RESULTS

Benchmarking service

To automate AF method benchmarking with a wide range of reference data sets, we developed a publicly available web-based evaluation framework (**Figure 1**). Using this workflow, an AF method developer who wants to evaluate their own algorithm first downloads sequence data sets from one or more of the five categories (e.g., data set of protein sequences with low identity from the protein sequence classification category) from the server. The developer then uses the downloaded data set to calculate pairwise AF distances or dissimilarity scores between the sequences of the selected data sets. The benchmarking service accepts the resulting pairwise distances in tab-separated value (TSV) format or as a matrix of pairwise distances in standard PHYLIP format. In addition, benchmarking procedures in two categories (genome-based phylogeny and horizontal gene transfer) also support trees in Newick format to allow for further comparative analysis of tree topologies.

Once the output file is uploaded to the AFproject web server, the service starts the benchmarking procedure, which is typically completed in a few seconds. Finally, the raw data and the time-stamped benchmark report are stored and provided to the submitter. The report shows the performance of the evaluated method and compares it with the performance of other methods that have been previously evaluated through the AFproject web server. In the report, the performance of the compared methods is ordered by a statistical measure specific

to the respective benchmark category (e.g., the Robinson-Foulds distance measure³⁶ in the categories of gene trees, genome-based phylogeny and horizontal gene transfer). By default, the report is private (visible only to the submitter), and the developer can choose if and when to make the report publicly available. Similar to other benchmarking platforms³⁷, we have released the source code of the web service to facilitate transparency and encourage feedback and improvements from the community (<https://github.com/afproject-org/afproject>).

Alignment-free method catalog

To evaluate the performance of currently available AF tools and create a reference data for future comparisons, we benchmarked 24 standalone tools (**Table 1**), covering a large proportion of the currently available AF methods. Some tools offer multiple related methods to measure pairwise distances (or dissimilarity) between sequences; for instance, jD2Stat³⁵ supports three different distance measures based on the D_2 statistic: jD2Stat--d2n, jD2Stat--d2s and jD2Stat--d2st. In this study, we included these different distance measures, resulting in a total of 74 tested tool variants (**Figure 2**). Each of these tool variants was run with various combinations of parameter values ([Supplementary Table 1](#)). The values yielding the best performance for a given method were selected and saved in the AFproject database; if multiple parameters produced the same best-performing results for a tool, we selected only the values that were least computationally demanding (e.g., the shortest word length for word-counting methods or the smallest sketch size). Full information about the benchmarking results, including all combinations of parameter values of the evaluated tools, can be downloaded from <http://afproject.org/download/>.

Only three tools (Alignment-Free-Kmer-Statistics (AFKS)³², FFP³⁸ and mash⁹) are sufficiently generic to be applied to all 12 benchmarking data sets; the remaining tools can handle only subsets of our reference data sets, either because they have been designed only for a specific purpose (e.g., to handle only certain sequence types, such as nucleotides, proteins, and unassembled or assembled genomic sequences) or — less frequently — because of some unexpected software behavior (e.g., a program stops functioning, does not terminate in a reasonable amount of time or produces invalid results; [Supplementary Table 1](#)). Hence, one of the results of our benchmarking study is an extensive and annotated catalog of tools (<http://afproject.org/tools/>), which constitutes a resource not only for users of AF methods, but also for the developers of these methods, as it should help identify which aspects of existing software code may be in need of further development.

Protein sequence classification

Recognition of structural and evolutionary relationships among amino acid sequences is central to the understanding of the function and evolution of proteins. Historically, the first comprehensive evaluation of AF methods⁶ investigated the accuracy of the tools for protein structure classification at four hierarchical levels used in the Structural Classification of Proteins (SCOP) database, namely, family, superfamily, class, and fold³⁹. The original protocol tested six *k*-mer-based distance measures against a subset of the SCOP database, containing protein family members sharing less than 40% sequence identity⁶. In the present study, we extend the original analysis⁶ to test the accuracy of 56 tool-measure variants in recognition of structural relationships of protein sequences sharing both low (<40%) and high (≥40%) sequence identity (**Figure 2**).

The area under the receiver operating characteristics (ROC) curve (AUC), which indicates whether a method is able to discriminate between homologous and nonhomologous protein sequences (**Methods**), showed the favorable performance of AFKS³² software. AFKS with parameters set to the *simratio*³² distance and a word length of $k = 2$ is the best performing tool for both low- and high-sequence-identity data sets (**Figure 2**). For the latter type of the data set, the method produces the highest AUC values across all four structural levels, with an average AUC of 0.798 ± 0.139 ([Supplementary Table 2](#)). When considering the low-sequence-identity data set ([Supplementary Table 3](#)), AFKS--*simratio* also has the highest average AUC of 0.742 ± 0.079 but lower performance at the superfamily and family levels than *alfpy*³ (set to the Google distance and $k = 1$). *alfpy*--*google* is ranked 2nd (0.738 ± 0.091) and fourth (0.778 ± 0.142) for the low- and high-sequence-identity data sets, respectively. Notably, the top-seven-ranking positions in both the low- and high-sequence-identity data sets are occupied, though in a different order, by the same measures from AFKS and *alfpy* software (**Figure 2**).

In general, the tested tools achieve greater discriminatory power in recognizing structural relationships (higher average AUCs) in our high-sequence-identity data set than in the low-sequence-identity data set (**Figure 2**; Wilcoxon signed rank test: p -value = $2.602e-11$). Almost all tool variants, except AFKS--*afd* (AUC: 0.492 ± 0.016) for the low-sequence-identity data set, achieved higher overall performance than the random classifier (AUC > 0.5). As expected and previously reported^{3,6}, the tools lose discriminatory power from the family to the class level for both data sets (the AUC decreases;

[Supplementary Table 2-3](#)), as the sequence similarity is lower within higher hierarchical groups. As a result, all methods tested (except AFKS--*harmonic_mean*) achieve their best accuracy at the family level. The AUC values at the family, superfamily and fold levels are higher (Wilcoxon signed rank tests: p -value $< 1e-05$) for data sets with high sequence similarity than for data sets with low sequence similarity. The greatest difference in performance was observed at the family level, where the maximum AUC obtained by the tools with the high- and low-sequence-identity data sets was 1.0 and 0.84, respectively. The methods result in more similar AUCs at the class level for the low-sequence-identity data set than for the high-sequence-identity data set (Wilcoxon signed rank tests: p -value = 0.0285). Protein sequences at the class level lack conserved segments, and the median AUC values obtained by the methods in high- and low-sequence-identity data sets are similar to those obtained with the random classifier (median AUC: 0.57 in both data sets).

Gene tree inference

Only a few studies^{40,41} have evaluated AF methods in the construction of gene trees. Because of the limited amount of sequence information available, gene trees are typically more difficult to reconstruct than species trees⁴². We assessed the accuracy of 11 AF tools (55 tool variants) in inferring phylogenetic relationships of homologous sequences based on a collection of high-confidence SwissTree phylogenies representing different types of challenges for homology prediction, e.g., numerous gene duplications and HGT^{37,43}. Similar to SwissTree, we assessed the gene families at the protein-sequence level to minimize the impact of codon degeneracy. We thus interpret an inferred phylogenetic tree based on a homologous family of protein sequences as the tree for the gene family (i.e., the gene tree). As a measure

of accuracy, we computed the normalized Robinson-Foulds (nRF) distance³⁶ between the trees reconstructed by the AF methods under study and the reference trees. The nRF distance has values between 0 and 1, with 0 indicating identical tree topologies and 1 indicating the most dissimilar topologies.

None of the AF methods that we tested were able to perfectly infer the respective reference tree topology for any of the 11 gene families. jD2Stat³⁵ (D_2^n with parameter values $n=1$ and $k=5$) was the most accurate tool in our test (**Figure 2**). This method achieved the lowest nRF values (highest accuracy) among all the tested methods averaged across all 11 reference gene families (nRF = 0.3296 ± 0.1511 ; [Supplementary Table 4](#)), which can be interpreted as 33% ($\pm 15\%$) of incongruent bipartitions between the inferred and the reference tree. To put this number into perspective, the corresponding gene trees based on MSA (i.e., neighbor-joining trees inferred using ClustalW alignments generated with default parameters) yielded a similar average accuracy (nRF = 0.2995 ± 0.1511). In general, the nRF distances obtained by the tested methods vary greatly across the gene families (Friedman rank sum test: p -value < $2.2e-16$, $df=10$, Friedman chi-square = 463.88) due to different complexities of the encoded protein families (e.g., evolutionary distance between proteins, domain architecture, and structural and functional affiliations). Consequently, the tools obtain their best accuracy in phylogenetic inference of the eukaryotic protein family of sulfatase modifying factor (SUMF) proteins, which are characterized by a single protein domain and the smallest number of gene duplications; four distance measures in AFKS software generated trees (nRF = 0.077) with minor topological differences in the speciation order of three proteins ([Supplementary Figure 1](#)). The AF methods achieved the second-best accuracy (median nRF = 0.178) for the

eukaryotic NOX family NADPH oxidases, a gene family coding for transmembrane enzymes with 10 gene duplications and 3-4 protein domains. However, the examined tools produced highly inaccurate phylogenetic trees of two other transmembrane protein families, namely, Bambi and Asterix (median nRFs: 0.615 and 0.611, respectively), where more than 60% of tree topologies differed from the reference tree.

Regulatory elements

Analysis of gene regulatory sequences is another domain where AF methods are popular, as the similarity between these elements is usually low and alignments typically fail to detect it properly⁴. We adopted a benchmarking procedure and a reference data set of cis-regulatory modules (CRMs) introduced by Kantarovitz et al.⁴, which was further used in other studies⁴⁴, showing that alignment algorithms lag behind AF methods in recognizing functionally related CRMs. A CRM can be broadly defined as a contiguous noncoding sequence that contains multiple transcription factor binding sites and regulates the expression of a gene. The Kantarovitz protocol assesses to what extent AF tools are capable of capturing the similarities between functionally-related CRMs, expressed in the tissues of fly and human (see **Methods**).

However, none of the AF methods produced perfect results for any of the seven tissues/species data set combinations (i.e., all functionally related CRM pairs ranked ahead of all random DNA pairs), and alfpv software³ set to three distance measures — Canberra, Chebyshev and Jensen-Shannon divergence — captured the largest number (averaged across 7 tissue samples) of functionally related regulatory elements (**Figure 2**). The selection of Canberra distance (word length of $k = 2$) correctly recognized $73.6\% \pm 10.54\%$ of CRMs,

capturing the highest functional relatedness in three out of seven data sets (tracheal system: 97%, eye: 78% and blastoderm-stage embryo: 76% in fly; [Supplementary Table 5](#)). The Chebyshev distance ($k = 7$) obtained the second-highest average performance of 67.59% and the highest performance variation across seven data sets (standard deviation = 20.14%) among all methods in the ranking; this measure had the highest performance for two tissues (peripheral nervous system in fly and HBB complex in human) and relatively low performance in human liver tissue. The third measure, Jensen-Shannon divergence ($k = 2$), achieved more stable performance across the data sets than the Canberra and Chebyshev distances ($63.16\% \pm 8.22\%$). Overall, 51 out of 63 methods showed average performance better than that of the random classifier ($> 50\%$).

Genome-based phylogeny

AF methods are particularly popular in genome-based phylogenetic studies^{9,12,13,38}, because of the considerable size of the input data and complex correspondence of the sequence parts, often resulting from genome rearrangements⁴⁵. Additionally, no statistical substitution models are currently available for assessing the evolution of complete genomes. We assessed the ability of AF methods to infer species trees using benchmarking data from different taxonomic groups, including bacteria, animals and plants. Here, we used completely assembled genomes as well as simulated unassembled next-generation sequencing reads at different levels of coverage.

Assembled genomes

As many studies have applied AF methods to whole mitochondrial genomes^{46,47}, we tested the performance of 23 AF software tools (70 tool variants in total) in phylogenetic inference using complete mtDNA from 25 fish species of the suborder Labroidei⁴⁸. The best accuracy was achieved by nine AF tools (19 tool variants), which generated tree topologies that were almost identical to the reference Labroidei tree (nRF = 0.05; **Figure 2**; [Supplementary Table 6](#)). The results differ only in the speciation order of three closely related fish species belonging to the Tropheini tribe of the Pseudocrenilabrinae family ([Supplementary Figure 2](#)). The same species were misplaced in the topologies generated by another 39 tool variants that all occupied the second place in the benchmark ranking (nRF = 0.09). These methods additionally misplace species within the Pomacentridae and Embiotocidae families. These results indicate that most AF methods infer trees in general agreement with the reference tree of mitochondrial genomes^{18,46,49,50}.

We further tested the performance of AF methods in phylogenetic inference with larger, bacterial genomes of *Escherichia coli/Shigella* and with nuclear genomes of plant species (**Figure 2**). Seven tools (nine tool variants) could not be tested on all three sets of complete genomes since the programs did not complete analyses ([Supplementary Table 1](#)). The remaining 16 tools (61 tool variants) lead to greater nRF distances, i.e., lower performance for the phylogeny of the *E. coli/Shigella* and plant nuclear genomes than for the phylogeny of mitochondrial genomes (**Figure 2**; one-way analysis of variance (ANOVA) with repeated measures: p -value < $2e-16$; *post hoc* pairwise paired t test: p -value < $2e-16$). Although the tools that we tested show similar nRF distances for bacterial and plant genomes in general (pairwise paired t test: p -value = 0.073), the top-performing tools are different between the

two data sets. For example, phylonium⁵¹ and andi²², which were developed for phylogenetic comparison of closely related organisms, are the best performing tools for the *E. coli/Shigella* data sets, whereas on the plant data sets, both tools perform poorly (**Figure 2**). Phylonium almost perfectly reproduced the reference tree for the *E. coli/Shigella* group with an nRF = 0.04 ([Supplementary Table 7](#); there was only a single error in the placement of two closely related *E. coli* K-12 substrains: BW2952 and DH10B; [Supplementary Figure 3](#)), while the plant trees obtained by these tools showed very low topological similarity to the reference tree (nRF = 0.64; [Supplementary Table 8](#)).

The best-performing tools for the plant data set are co-phylog²¹, mash⁹ and Multi-SpaM²³, all of which almost perfectly recovered the reference tree topology of the plant species (with an nRF = 0.09 for all three programs). In each of the trees produced by these programs, there was exactly one species placed at an incorrect position compared to in the reference tree, namely, in the speciation order in the Brassicaceae family for co-phylog ([Supplementary Figure 4](#)), mash ([Supplementary Figure 5](#)) and for Multi-SpaM, the last of which placed *Carica papaya* outside the Brassicales order ([Supplementary Figure 6](#)). Additionally, co-phylog is the third-best-performing tool in reconstructing the *E. coli/Shigella* tree topology (nRF = 0.12), while mash and Multi-SpaM are at the fourth and sixth positions, respectively, in this ranking (nRF = 0.15 and nRF = 0.27, respectively). As a result, co-phylog, mash, FFP³³, Skmer⁵² and FSWM²⁴ are among the top 5 best-performing tools for both data sets (**Figure 2**), .

Raw sequencing reads

We also tested the accuracy of AF tools in phylogenetic inference based on unassembled sequencing reads, represented by seven different levels of sequencing coverage, from *E. coli/Shigella* and from a set of plant species (**Table 1**; see **Methods**). No differences in nRF values were observed between the results based on the unassembled and assembled *E. coli/Shigella* genomes (Wilcoxon signed rank test: p -value = 0.06067), indicating that the AF tools exhibited equal performance for the unassembled and assembled genomes. In contrast, the tested tools showed lower performance (i.e., higher nRF values) in assembly-free phylogenetic reconstruction of the plant species (Wilcoxon signed rank test: p -value = $2.344e-05$). co-phylog²¹ achieved the minimum nRF for five out of seven coverage levels in the *E. coli/Shigella* data set (i.e., coverage from 0.0625 to 0.25 and from 1 to 5). This tool is thus the most accurate one (**Figure 2**), with an average nRF distance of 0.21 ± 0.14 ([Supplementary Table 9](#)). However, the accuracy of co-phylog for the unassembled plant data sets is drastically reduced (nRF = 0.4 ± 0.3 ; [Supplementary Table 10](#)), which places the tool at the 6th position in the ranking for the plant data set (**Figure 2**).

For the unassembled plant data sets, mash is the most accurate tool, i.e., the tool with the shortest nRF distance between the inferred trees and the reference tree. For the lowest coverage level (0.015625), mash still allows us to infer trees with average nRF distances of 0.27 from the reference tree ([Supplementary Table 10](#)). In general, mash and Skmer show a constantly high performance at all coverage levels. For the unassembled *E. coli/Shigella* data set, mash is ranked at the 4th position, with an average nRF distance of 0.34 ± 0.17 , and Skmer is placed at the 5th position (0.34 ± 0.19). Notably, while mash and Skmer have the highest performance at the lowest coverage levels for the plant data set, CAFE set to d_2^S — the

third-best-performing tool — achieved the highest accuracy (nRF distance = 0.18) at a medium level of sequencing coverage (coverage: 0.125) and maintained that accuracy for higher-coverage data sets.

When considering the most universal tools applied to all the tested reference data sets, mash ranks first and the fourth for the assembly-free phylogeny of plants and *E. coli/Shigella*, respectively (**Figure 2**). In addition to mash, two other methods designed specifically for phylogenetic reconstruction from next-generation sequencing data - Skmer and AAF⁵³ - are the only tools ranked among the top 5 methods tested on both unassembled data sets (**Figure 2**).

Horizontal gene transfer

To assess the accuracy of the AF methods in phylogenetic reconstruction of sequences that underwent frequent HGT events and genome rearrangements, we used sets of simulated genomes with different levels of HGT⁵⁴ as well as two real-world data sets of microbial species, namely, 27 genomes of *E. coli* and *Shigella*⁵⁴⁻⁵⁶ and eight *Yersinia* genomes^{54,57} (**Table 1**). Similarly to previous tests, we applied the nRF distance between the obtained and trusted reference trees as a measure of accuracy.

We simulated five sets of 33 genomes each with different extents of HGT as determined by the mean number of HGT events per iteration ($l = 0, 250, 500, 750, \text{ and } 1,000$; l is the number of HGT events attempted in the set at each iteration of the simulation process of genome evolution; for details, see **Methods**). The tool AFKS (Markov measure, with a word length of

$k = 12$) achieved the highest general accuracy (**Figure 2**) by obtaining the lowest average nRF (0.05 ± 0.05) and perfect topological agreement with the reference trees at the two lowest frequencies of simulated HGT ($l = 0$ and 250; [Supplementary Table 11](#)). As expected, for most AF methods, the accuracy of phylogenetic inference declines with an increase in the extent of HGT. Nevertheless, the six best-performing software applications - AFKS, CAFE, mash, jD2Stat, alphy, and FFP - were capable of reconstructing the reference tree with little incongruence at almost all HGT frequency levels ($\text{nRF} \leq 0.1$ at $l \leq 750$), except for the highest frequencies of HGT simulated, where the nRF distance was in the range of 0.13-0.17 ([Supplementary Table 11](#)). Interestingly, the basic AF distance measures (Euclidean, Manhattan, Canberra and LCC distances) implemented in alphy achieve a lower average nRF (0.07 ± 0.06) and minimum nRF at a higher HGT frequency level ($\text{nRF} = 0.13$) than AF tools designed for phylogenetic reconstruction of whole genomes (co-phylog, FSWM, Multi-SpaM and kr), which surprisingly were relatively inaccurate ($\text{nRF} > 0.2$ for different values of l). As has been reported before⁵⁴, the accuracy of kr generally increased (nRF : from 0.73 to 0.33) with increasing l .

To assess the performance of AF methods with real-world sequence data, we first used a reference supertree of 27 genomes of *E. coli* and *Shigella* that was generated based on thousands of single-copy protein trees⁵⁴⁻⁵⁶. For this data set, the tools designed for whole-genome phylogenetics achieved lower nRF values than did basic AF distance measures; eleven tools for whole-genome phylogenetics occupied the first six positions in the ranking list (**Figure 2**). Three such methods — andi, co-phylog and phylonium — achieved the highest accuracy (**Figure 2**), with a minimum nRF of 0.08 ([Supplementary Table 12](#)). andi

and co-phylog yielded topologically equivalent trees that were very similar to the reference tree, misplacing only two closely related *E. coli* strains in the D and B1 reference groups ([Supplementary Figure 7](#)), while phylonium showed two minor topological differences in *E. coli* reference group D ([Supplementary Figure 8](#)). Most AF measures implemented in AFKS, alfpy and CAFE were ranked at the 10th position (**Figure 2**) and led to reconstruction of inaccurate species trees where half of the bipartitions were not present in the reference tree (nRF = 0.5). Interestingly, the opposite result was obtained for phylogenetic inference of 8 *Yersinia* genomes, where almost all basic measures (42 tool variants) recovered the reference tree topology (nRF = 0) while whole-genome phylogenetic tools obtained relatively incongruent trees (nRF > 0.2) compared to the reference (**Figure 2**, [Supplementary Table 13](#)).

DISCUSSION

We have addressed key challenges in assessing methods for AF sequence comparison by automating the application of multiple AF methods to a range of reference data sets. This automated approach critically benefits from extensive work described in the previous section to identify optimal parameter values for all combinations of methods and data sets. Finally, the resulting open platform for a standardized evaluation of new methods is provided with an interactive web-based interface and a reporting functionality designed to ensure reproducibility. We believe that the uniform framework for testing AF algorithms with common data sets and procedures will be beneficial to both developers and users of these methods. The benchmarking results will guide users in choosing the most effective tool tailored to their project needs and for finding optimal parameter settings. For developers, the

interactive platform speeds up benchmarking and provides a reference data set with which to compare new AF methods to existing approaches.

Our results showed that no single method performed best across all the data sets tested. Nevertheless, some tools were among the top five performers more often than others. For example, when considering genomic-scale benchmarks, encompassing 8 datasets from the whole-genome phylogeny and horizontal gene transfer categories, the tools developed for genomic comparisons were among the top-5-performing tools: mash (8 times), Skmer (7 times), co-phylog and FFP (6 times), and FSWM/Read-SpaM (5 times; **Figure 2**). Since mash is the only method that is placed among the top five-best performing tools on all genome-scale benchmarking data sets, it is particularly well suited for genome sequence comparisons, regardless of the phylogenetic range and technology that were used to obtain the data (e.g., short reads or assembled contigs). Most AF approaches (14 out of 21 software applications or, more specifically, 56 out of 68 tool variants) performed particularly well — although not perfectly — in phylogenetic inference of mitochondrial genomes from different fish species, yielding trees generally consistent ($nRF < 0.1$) with the reference phylogeny (**Figure 2**, [Supplementary Table 6](#)). These results, however, are in contrast to our results on whole-genome sequence comparison for prokaryotes and eukaryotes. Thus, novel AF methods should not be benchmarked with mitochondrial sequences alone. Considering the evolutionary and structural relationships among the protein sequences and inferred gene trees, we were surprised by the highest performance of very simple AF distance measures implemented in AFKS and alphy (i.e., intersection, simratio, Kulczynski, Bray-Curtis, Google, Canberra, Squared_chord, chi_squared, and Manhattan). Overall, methods based on

conventional statistics performed better than approaches using more complex statistics such as state-of-the-art D_2 -related metrics implemented in jD2Stat (D_2^S , D_2^* , and D_2^n) and AFKS (D_2^z , D_2^* , and D_2^S), the Markov metric in AFSK (sim_mm, rr_k_r, and markov), and the N_2 metric in AFKS (n_{2r}) ([Supplementary Table 14](#)). Interestingly, the basic Canberra distance implemented in alphy is the most effective distance measure in recognizing functionally related regulatory sequences ([Supplementary Table 5](#)), greatly exceeding the D_2^S and D_2^* statistics from CAFE and jD2Stat.

Another surprising observation in our study is that different implementations of the same AF algorithm, run with the same input parameter values, can deliver different results. For example, two implementations of the Canberra distance from AFKS and alphy achieve different performances in almost all data sets (**Figure 2**). The discrepancy in the Canberra distance with a word length of $k = 2$ between the two tools is apparent for the CRM data set, where AFKS--*Canberra* obtained a performance score of 54, while alphy--*Canberra* had a performance score of 74, which was the highest performance score among the tools that we evaluated ([Supplementary Table 5](#); see the **Methods** section for the definition of “performance score”). The differences observed were due to the different methods of sequence data preprocessing applied by both the tools — alphy projects sequences into a vector of k -mer frequencies, while AFKS represents sequences as k -mer count vectors with the inclusion of pseudocounts. This sequence data preprocessing in alphy and AFKS has the highest impact on the performance of methods based on the Canberra distance in the case of nucleotide data sets of regulatory elements, whole genomes of plants and simulated genomes that underwent HGT ([Supplementary Figure 9](#)). For other data sets, the same distance

measures in alfpy and AFKS, run on common word lengths, produce results with very similar performances, and the observed differences between the tools in this study are the results of different ranges of k . Similarly, the D_2^* and D_2^S metrics implemented in AFKS, CAFE and jD2Stat produce slightly different results.

When assessing the accuracy of AF methods in inferring phylogenetic relationships, we compared the inferred phylogenetic tree topologies to trusted reference tree topologies. However, the assumption that evolutionary relationships are generally tree-like is known to be unrealistic because genome evolution is shaped by both vertical and lateral processes^{56,58,59}. Although the signal of vertical descent (e.g., for ribosomal rRNAs) can be described adequately using a phylogenetic tree, horizontal transfer of genetic material between different taxa and genome rearrangements can obscure this signal. A classic example involves the *Yersinia* genomes, which are well-known to have undergone extensive structural rearrangements⁵⁷. We have shown in this study that reconstructing phylogenetic trees of these taxa from whole-genome sequences is difficult with AF methods. The same is true for more conventional approaches that are based on MSA⁵⁷, and finding a trusted reference tree for these taxa has been problematic. In such cases, a non-tree-like network representation of genome evolution is more appropriate. Recent studies^{60,61} have demonstrated the scalability and applicability of AF methods to quickly infer networks of relatedness among microbial genomes. Although we did not consider networks in this study, the curated benchmarking data sets can be easily extended to AF phylogenetic analysis beyond a tree-like structure in the future.

We acknowledge that the presented data sets do not cover all possible applications of AF tools. The data sets include only the most typical sequence comparison tasks, where all-versus-all sequence comparisons need to be computed. Although the AF project is extendable and new data sets can be seamlessly added in the future, for more specific applications such as orthology prediction, genome assembly, RNA-seq aligners or metagenomics analyses, we recommend using other web-based benchmarking services developed for these purposes^{37,62–65}. Nevertheless, AFproject can be used to evaluate any sequence comparison tool — not necessarily AF — that produces dissimilarity scores between sequence pairs. Since similarity scores can be easily converted into dissimilarity scores, our benchmarking system can also be used to evaluate methods that generate similarity scores, e.g., alignment scores. We thus invite developers and users of sequence comparison methods to submit and evaluate their results with the AFproject benchmarking platform. The ability to rapidly, objectively and collaboratively compare computational methods for sequence comparison should be beneficial for all fields of DNA and RNA sequence analysis, regardless of whether the analysis is alignment-based or alignment-free.

METHODS

Data sets. Twelve sequence data sets were used to evaluate AF methods across five research areas (**Table 1**).

Protein homology. The reference data sets of protein family members sharing a high ($\geq 40\%$) and low ($< 40\%$) sequence identity were constructed based on two sections of the SCOPE database v. 2.07³⁹, namely, ASTRAL95 and ASTRAL40 v. 2.07⁶⁶, respectively. The SCOPE

database provides a structural classification of proteins at four levels: classes, folds, superfamilies and families. According to previous studies^{3,6}, the ASTRAL data sets were subsequently trimmed to exclude sequences with unknown amino acids and families with fewer than 5 proteins and included only the four major classes (i.e., α , β , α/β , and $\alpha+\beta$). To minimize the requirements for AF method submission related to performing all-versus-all sequence comparisons and uploading the output to the AFproject server, we further reduced the data sets by randomly selecting only two protein members in each family. As ASTRAL95 also contains protein family members sharing a sequence identity lower than 40%, the Needleman-Wunsch alignment was performed (using needle software in the EMBOSS package⁶⁷) to select proteins with a sequence identity $\geq 40\%$ to acquire a reference data set of proteins with high sequence identity.

Gene trees. Reference trees and corresponding protein sequences of eleven gene families were downloaded from SwissTree release 2017.0 (<https://swisstree.vital-it.ch/>): Popeye domain-containing protein family (49 genes), NOX 'ancestral-type' subfamily NADPH oxidases (54 genes), V-type ATPase beta subunit (49 genes), Serine incorporator family (115 genes), SUMF family (29 genes), Ribosomal protein S10/S20 (60 genes), Bambi family (42 genes), Asterix family (39 genes), Cited family (34 genes), Glycosyl hydrolase 14 family (159 genes), and Ant transformer protein (21 genes).

Gene regulatory elements. The data set of CRMs known to regulate expression in the same tissue and/or developmental stage in fly or human was obtained from Kantorovitz et al.⁴. The data set was specifically selected to test the capacity of AF measures to identify functional

relationships between regulatory sequences (e.g., enhancers or promoters). The data set contains 185 CRM sequences taken from *D. melanogaster* - blastoderm-stage embryo ($n = 82$), eye ($n = 17$), peripheral nervous system ($n = 23$), and tracheal system ($n = 9$) - and *Homo sapiens* - HBB complex ($n = 17$), liver ($n = 9$) and muscle ($n = 28$).

Genome-based phylogeny. The sequences of 25 whole mitochondrial genomes of fish species from the suborder Labroidei and the species tree were taken from Fischer et al.⁴⁸. The set of 29 *E. coli* genome sequences was originally compiled by Yin and Jin²¹ and has been used in the past by other groups to evaluate AF programs^{22,23,68}. Finally, the set of 14 plant genomes is from⁶⁹. This set was also used in the past to evaluate AF methods. To simulate unassembled reads from these data sets, we used the program ART⁷⁰.

Horizontal gene transfer. The 27 *E. coli* and *Shigella* genomes, and the 8 *Yersinia* genomes, were taken from Bernard et al.⁵⁴. We used EvolSimulator⁷¹ to simulate HGT in microbial genomes, adopting an approach similar to that described in Bernard et al.⁵⁴. Each set of genomes was simulated under a birth-and-death model at speciation rate = extinction rate = 0.5. The number of genomes in each set was allowed to vary from 25 to 35, with each containing 2,000–3,000 genes 240–1,500 nucleotides in length. HGT receptivity was set at a minimum of 0.2, mean of 0.5 and maximum of 0.8, with a mutation rate $m = 0.4$ –0.6 and a number of generations $i = 5,000$. The varying extent of HGT was simulated using the mean number of HGT events attempted per iteration $l = 0, 250, 500, 750$ and 1000, and divergence factor $d = 2,000$ (transferred genes that are of high sequence divergence, i.e., >2,000 iterations apart, will not be successful). All other parameters in this simulation followed Beiko et al.⁷¹.

Alignment-free tools

AAF⁵³ reconstructs a phylogeny directly from unassembled next-generation sequencing reads. Specifically, AAF calculates the Jaccard distance between sets of k -mers of two samples of short sequence reads. This distance between samples or species is based on the estimate of the rate parameter from a Poisson process for a mutation occurring at a single nucleotide. The phylogeny is constructed using weighted least squares with weights proportional to the expected variance of the estimated distances. AAF provides features for correcting tip branches and bootstrapping of the obtained phylogenetic trees, directly addressing the problems of sequencing error and incomplete coverage.

AFKS³² is a package for calculating 33 k -mer-based dissimilarity/distance measures between nucleotide or protein sequences. AFKS categorizes the measures into nine families:

Minkowski (e.g., Euclidean), Mismatch (e.g., Jaccard), Intersection (e.g., Kulczynski), D2 (e.g., D2s), Squared Chord (e.g., Hellinger), Inner Product (e.g., normalized vectors), Markov (e.g., SimMM), Divergence (e.g., KL Conditional), and Others (e.g., length difference). The tool determines the optimal k -mer size for given input sequences and calculates dissimilarity/distance measures between k -mer counts that include pseudocounts (adding 1 to each k -mer count). The obtained distance is standardized to between 0 and 1.

alfpy³ provides 38 AF dissimilarity measures with which to calculate distances among given nucleotide or protein sequences. The tool includes 25 k -mer-based measures (e.g., Euclidean, Minkowski, Jaccard, and Hamming), eight information-theoretic measures (e.g., Lempel–Ziv

complexity and normalized compression distance), three graph-based measures, and two hybrid measures (e.g., Kullback-Leibler divergence and W-metric). `alfpy` is also available as a web application and Python package. In this study, the results based on 14 dissimilarity measures are evaluated.

ALFRED-G⁷² uses an efficient algorithm to calculate the length of maximal k -mismatch common substrings between two sequences. Specifically, to measure the degree of dissimilarity between two nucleic acid or protein sequences, the program calculates the length of maximal word pairs — one word from each of the sequences — with up to k mismatches.

`andi`²² estimates phylogenetic distances between genomes of closely related species by identifying pairs of maximal unique word matches a certain distance from each other and on the same diagonal in the comparison matrix of two sequences. Such word matches can be efficiently found using enhanced suffix arrays. The tool then uses these gap-free alignments to estimate the number of substitutions per position.

CAFE³⁴ is a package for efficient calculation of 28 AF dissimilarity measures, including 10 conventional measures based on k -mer counts, such as Chebyshev, Euclidean, Manhattan, uncentered correlation distance, and Jensen-Shannon divergence. It also offers 15 measures based on the presence/absence of k -mers, such as Jaccard and Hamming distances. Most importantly, it provides a fast calculation of background-adjusted dissimilarity measures including CVTree, `d2star` and `d2shepp`. CAFE allows for both assembled genome sequences and unassembled next-generation sequencing shotgun reads as inputs. However, it does not

deal with amino acid sequences. In this study, the results based on CVTree, d2star and d2shepp are evaluated.

co-phylog²¹ estimates evolutionary distances among assembled or unassembled genomic sequences of closely related microbial organisms. The tool finds short, gap-free alignments of a fixed length and consisting of matching nucleotide pairs only, except for the middle position in each alignment, where mismatches are allowed. Phylogenetic distances are estimated from the fraction of such alignments for which the middle position is a mismatch.

EP-sim⁷³ computes an AF distance between nucleotide or amino acid sequences based on entropic profiles^{74,75}. The entropic profile is a function of the genomic location that captures the importance of that region with respect to the whole genome. For each position, it computes a score based on the Shannon entropies of the word distribution and variable-length word counts. EP-sim estimates a phylogenetic distance, similar to D_2 by summing the entropic profile scores over all positions, or similar to D_2^* , with the sum of normalized entropic profile scores.

FFP^{33,38} estimates the distances among nucleotide or amino acid sequences. The tool calculates the count of each k -mer and then divides the count by the total count of all k -mers to normalize the counts into frequencies of a given sequence. This process leads to the conversion of each sequence into its feature frequency profile (FFP). The pairwise distance between two sequences is then calculated by the Jensen-Shannon divergence between their respective FFPs.

FSWM²⁴ estimates the phylogenetic distance between two DNA sequences. The program first defines a fixed binary pattern P of length l representing “match positions” and “don’t care positions”. Then, it identifies all “Spaced-word Matches” (*SpaM*) w.r.t. P , i.e., gap-free local alignments of the input sequences of length l , with matching nucleotides at the “match positions” of P and possible mismatches at the “don’t care” positions. To estimate the distance between two DNA sequences, *SpaMs* with low overall similarity are discarded, and the remaining *SpaMs* are used to estimate the distance between the sequences, based on the mismatch ratio at the “don’t care” positions. There is a version of FSWM that can compare sets of unassembled sequencing reads to each other called *Read-SpaM*⁷⁶.

jD2Stat³⁵ utilizes a series of D_2 statistics^{15,16} to extract k -mers from a set of biological sequences and generate pairwise distances for each possible pair as a matrix. For each sequence set, we generated distance matrices (at the defined k ; [Supplementary Table 1](#)), each using D_2^s (D2S; exact k -mer counts normalized based on the probability of occurrences of specific k -mers), D_2^* (d2St; similar to D_2^s but normalized based on means and variance), and D_2^n (d2n; extension of D_2 that expands each word w recovered in the sequences to its neighborhood n , i.e., all possible k -mers with n number of wildcard residues, relative to w).

kmacs¹⁸ compares two DNA or protein sequences by searching for the longest common substrings with up to k mismatches. More precisely, for each position i in one sequence, the program identifies the longest pair of substrings with up to k mismatches, starting at i in the

first sequence and somewhere in the second sequence. The average length of these substring pairs is then used to define the distance between the sequences.

kr^{49} estimates the evolutionary distance between genomes by calculating the number of substitutions per site. The estimator for the rate of substitutions between two unaligned sequences depends on a mathematical model of DNA sequence evolution and average shortest unique substring (shustring) length.

kSNP3⁷⁷ identifies single nucleotide polymorphisms (SNPs) in a set of genome sequences without the need for genome alignment or a reference genome. The tool defines a SNP locus as the k -mers surrounding a central SNP allele. kSNP3 can analyze complete genomes, draft genomes at the assembly stage, genomes at the raw reads stage, or any combination of these stages. Based on the identified SNPs, kSNP3.0 estimates phylogenetic trees by parsimony, neighbor-joining and maximum-likelihood methods and reports a consensus tree with the number of SNPs unique to each node.

kWIP⁷⁸ estimates genetic dissimilarity between samples directly from next-generation sequencing data without the need for a reference genome. The tool uses the weighted inner product (WIP) metric, which aims to reduce the effect of technical and biological noise and elevate the relevant genetic signal by weighting k -mer counts by their informational entropy across the analysis set. This procedure downweights k -mers that are typically uninformative (highly abundant or present in very few samples).

LZW-Kernel⁷⁹ classifies protein sequences and identifies remote protein homology via a convolutional kernel function. LZW-Kernel exploits code blocks detected by the universal Lempel-Ziv-Welch (LZW) text compressors and then builds a kernel function out of them. LZW-Kernel provides a similarity score between sequences from 0 to 1, which can be directly used with support vector machines (SVMs) in classification problems. LZW-Kernel can also estimate the distance between protein sequences using normalized compression distances (LZW-NCD).

mash⁹ estimates the evolutionary distance between nucleotide or amino acid sequences. The tool uses the MinHash algorithm to reduce the input sequences to small 'sketches', which allow fast distance estimations with low storage and memory requirements. To create a 'sketch', each k -mer in a sequence is hashed, which creates a pseudorandom identifier (hash). By sorting these hashes, a small subset from the top of the sorted list can represent the entire sequence (min-hashes). Two sketches are compared to provide an estimate of the Jaccard index (i.e., the fraction of shared hashes) and the Mash distance, which estimates the rate of sequence mutation under an evolutionary model.

Multi-SpaM²³, similarly to FSWM, starts with a binary pattern P of length l representing “match positions” and “don’t care positions”. It then searches for four-way Spaced-word Matches (*SpaMs*) w.r.t. P , i.e., local gap-free alignments of length l involving four sequences each and with identical nucleotides at the “match positions” and possible mismatches at the “don’t care positions”. Up to 1,000,000 such multiple SpaMs with a score above some threshold are randomly sampled, and a quartet tree is calculated for each of them with

RAxML⁸⁰. The program *Quartet Max-Cut*⁸¹ is used to calculate a final tree of all input sequences from the obtained quartet trees.

phylonium⁵¹ estimates phylogenetic distances among closely related genomes. The tool selects one reference from a given set of sequences and finds matching sequence segments of all other sequences against this reference. These long and unique matching segments (anchors) are calculated using an enhanced suffix array. Two equidistant anchors constitutes homologous region, in which SNPs are counted. With the analysis of SNPs, phylonium estimates the evolutionary distances between the sequences.

RTD-Phylogeny⁸² computes phylogenetic distances among nucleotide or protein sequences based on the time required for the reappearance of k -mers. The time refers to the number of residues in successive appearance of particular k -mers. Thus, the occurrence of each k -mer in a sequence is calculated in the form of a return time distribution (RTD), which is then summarized using the mean (μ) and standard deviation (σ). As a result, each sequence is represented in the form of a numeric vector of size $2 \cdot 4^k$ containing the μ and σ of 4^k RTDs. The pairwise distance between sequences is calculated using Euclidean distance.

Skmer⁵² estimates phylogenetic distances between samples of raw sequencing reads. Skmer runs mash⁹ internally to compute the k -mer profile of genome skims and their intersection, and estimates the genomic distances by correcting for the effect of low coverage and sequencing error. The tool can estimate distances between samples with high accuracy from

low-coverage and mixed-coverage genome skims with no prior knowledge of the coverage or the sequencing error.

Slope-SpaM⁸³ estimates the phylogenetic distance between two DNA sequences by calculating the number N_k of k -mer matches for a range of values of k . The distance between the sequences can then be accurately estimated from the *slope* of a certain function that depends on N_k . Instead of exact word matches, the program can also use *SpaMs* w.r.t. a predefined binary pattern of “match positions” and “don’t care positions”.

spaced^{84–86} is similar to previous methods that compare the k -mer composition of DNA or protein sequences. However, the program uses so-called “spaced words” instead of k -mers. For a given binary pattern P of length l representing “match positions” and “don’t care positions”, a spaced word w.r.t. P is a word of length l with nucleotide or amino acid symbols at the “match positions” and “wildcard characters” at the “don’t care positions”. The advantage of using spaced words instead of exact k -mers is that the obtained results are statistically more stable. This idea has been previously proposed for database searching^{87,88}. The original version of Spaced⁸⁴ used the Euclidean or Jensen-Shannon⁸⁹ distance to compare the spaced-word composition of genomic sequences. By default, the program now uses a distance measure introduced by Morgenstern et al. 2015⁸⁶ that estimates the number of substitutions per sequence position.

Underlying Approach⁹⁰ estimates phylogenetic distances between whole genomes using matching statistics of common words between two sequences. The matching statistics are

derived from a small set of independent subwords with variable lengths (termed *irredundant common subwords*). The dissimilarity between sequences is calculated based on the length of the longest common subwords, such that each region of genomes contributes only once, thus avoiding counting shared subwords multiple times (i.e., subwords occurring in genomic regions covered by other more significant subwords are discarded).

Benchmarks.

Evaluation of structural and evolutionary relationships among proteins. To test the capacity of AF distance measures to recognize SCOPe relationships (i.e., family, superfamily, fold, and class), we used a benchmarking protocol from previous studies^{6,3}. Accordingly, the benchmarking procedure takes the distances between all sequence pairs present in the data set file. The distances between all protein pairs are subsequently sorted from minimum to maximum (i.e., from the maximum to minimum similarity). The comparative test procedure is based on a binary classification of each protein pair, where 1 corresponds to the two proteins sharing the same group in the SCOPe database and 0 corresponds to other outcomes. The group can be defined at one of the four different levels of the database (family, superfamily, fold, and class), exploring the hierarchical organization of the proteins in that structure. Therefore, each protein pair is associated with four binary classifications, one for each level. At each SCOP level, ROC curves and AUC values computed in scikit-learn⁹¹ are obtained to give a unique number of the relative accuracy of each metric and level according to the SCOP classification scheme. The overall assessment of method accuracy is an average of AUC values across all four SCOP levels.

Evaluation of functionally related regulatory sequences. To test how well AF methods can capture the similarity between sequences with similar functional roles, we used the original benchmarking protocol introduced by Kantorovitz et al.⁴. Briefly, a set of CRMs known to regulate expression in the same tissue and/or developmental stage is taken as the ‘positive’ set. An equally sized set of randomly chosen noncoding sequences with lengths matching the CRMs is taken as the ‘negative’ set. Each pair of sequences in the positive set is compared, as is each pair in the negative set. The test evaluates if functionally-related CRM sequence pairs (from the positive half) are better scored by a given AF tool (i.e., have lower distance/dissimilarity values) than unrelated pairs of sequences (from the negative half). This procedure is done by sorting all pairs, whether they are from the positive set or the negative set, in one combined list and then counting how many of the pairs in the top half of this list are from the positive set. The overall assessment of method accuracy is the weighted average of the positive pairs across all seven subsets.

Evaluation of phylogenetic inference. The accuracy of AF methods for data sets from three categories - protein homology, genome-based phylogeny and horizontal gene transfer - was evaluated by a comparison of topology between the method’s tree and the reference tree. The pairwise sequence distances obtained by the AF method were used as input for the neighbor-joining algorithm (fneighbor in the EMBOSS package⁶⁷, version: EMBOSS:6.6.0.0 PHYLIPNEW:3.69.650) to generate the corresponding method tree. To assess the degree of topological (dis)agreement between the inferred and reference trees, we calculated the nRF distance³⁶ using the Tree.compare function in the ETE3⁹² toolkit for phylogenetic trees with the option unrooted=True. When nRF = 0, the test and reference topologies are identical,

implying the highest accuracy for the method. Conversely, at $nRF = 1$, no bipartition in the reference is recovered.

Performance summary criteria. Figure 2 shows the color-coded performance of the evaluated AF methods across 12 reference data sets.

Performance score. For our benchmarking data sets, we use different measures to assess the performance of each method for a given data set, for example, nRF or AUC. To make our benchmarking results from different data sets comparable, we converted these measures to a performance score with values between 0 and 100. For the protein sequence classification data sets, this score is defined as $AUC \times 100$; for data sets from gene trees, genome-based phylogeny and horizontal gene transfer categories, we define the performance score as $(1 - nRF) \times 100$. For the regulatory elements data set, the performance score is already a number between 0 and 100, namely, the weighted average performance across seven data subsets.

Moreover, we define an *overall performance score* that assesses each method across the data sets and that also takes values between 0 and 100. For a given method, we calculate revised scores for each data set, on which the method was tested as $(S - min_score) / (max_score - min_score) \times 100$, where S is the performance score obtained by the method and min_score and max_score are the minimum and maximum scores obtained with all methods for a given data set, respectively. In this way, the best-performing method in a given data set receives a score of 1, and the worst performer receives a score of 0. The overall performance is an average of the revised scores across the data sets on which the given method was tested.

Data and code availability. All data sets and results discussed in the paper are freely available from our website (<http://afproject.org>). The source code of the AFproject service is available under an open source license (Mozilla Public License Version 2.0) at <https://github.com/afproject-org/afproject>.

ACKNOWLEDGMENTS

We thank Svenja Dörrer for providing benchmarking data sets of unassembled sequencing reads. This work was funded by National Science Centre Poland [2017/25/B/NZ2/00187] to A.Z. and W.M.K.; The Oklahoma Center for the Advancement of Science and Technology [PS17-015] to H.Z.G. and B.T.J.; US National Science Foundation (NSF) [DMS-1518001] and National Institutes of Health (NIH) [R01GM120624] to K.T., M.S.W. and F.S.; VW Foundation [VWZN3157] to T.D.; FCT [UID/CEC/50021/2019], [UID/EMS/50022/2019], [PTDC/EMSSIS/0642/2014], [PTDC/CCI-CIF/29877/2017] to S.V.; Australian Research Council [DP150101875] and [DP190102474] to CX.C; The Oklahoma Center for the Advancement of Science and Technology [PS17-015] to B.T.J;

AUTHOR CONTRIBUTIONS

A.Z. and W.M.K. conceived the project. B.M., A.K.L., CX.C, G.B., A.Z. and W.M.K contributed the reference data sets. A.Z. and W.M.K. designed and implemented the benchmarking service. A.Z., H.Z.G., B.T.J., G.B., C.L., K.T., T.D., J.C., M.C., S.K, S.R. and M.S.W. contributed the benchmarking results. A.Z., B.M., F.S., S.V. and W.M.K. analyzed

the results. A.Z., B.M., S.V., J.S.A., C.X.C., H.Z.G., J.C. and W.M.K. prepared the manuscript, with feedback from all other coauthors. W.M.K. coordinated the project.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

FIGURE LEGENDS

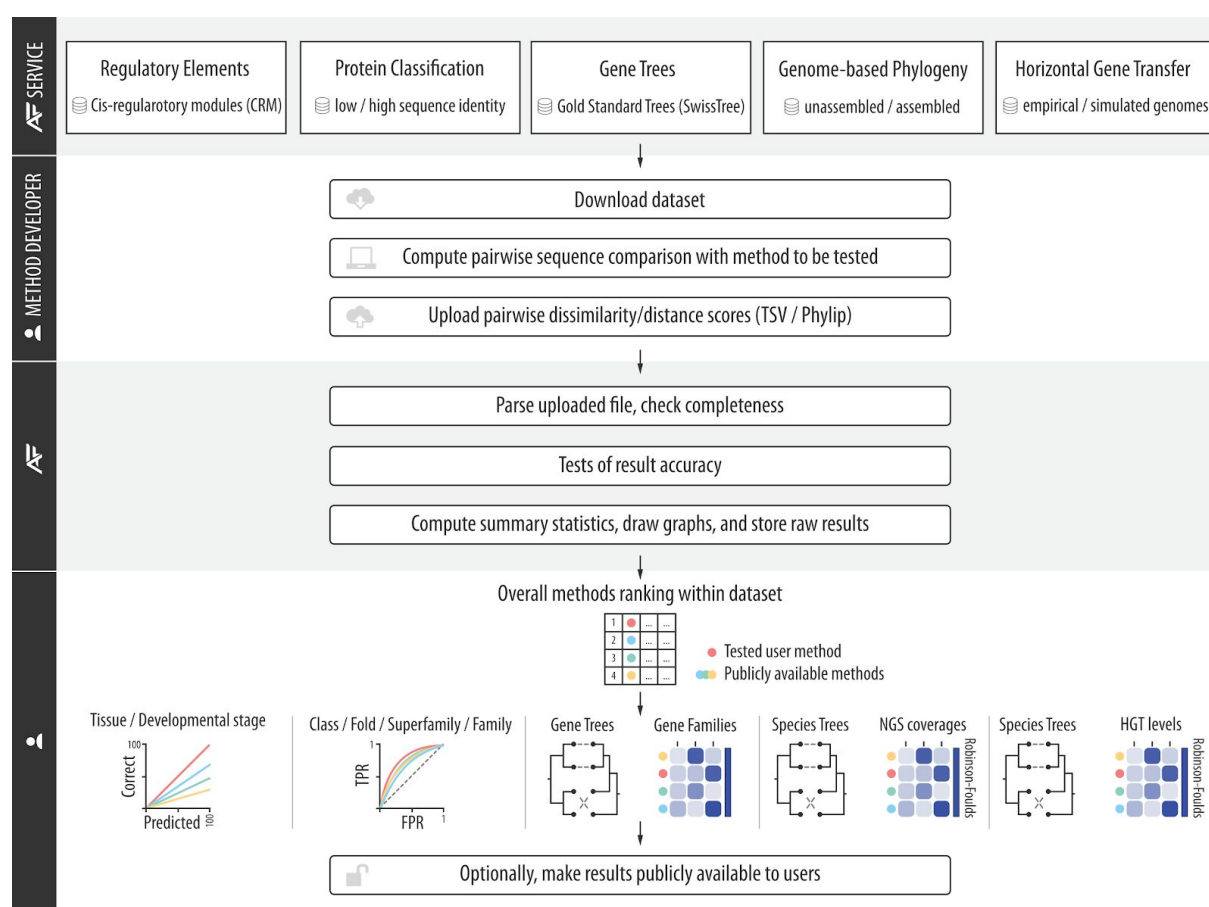


Figure 1. Overview of the AFproject benchmarking service facilitating assessment and comparison of AF methods. AF method developers run their methods on a reference sequence set and submit the computed pairwise sequence distances to the service. The submitted

distances are subjected to a test specific to given data sets, and the results are returned to the method developer, who can choose to make the results publicly available.

Figure 2. Summary of AF tool performance across all reference data sets. The numbers in the fields indicate the performance ranks of a given AF method for a given data set. Fields are color-coded by performance values from 0 to 100 (see **Methods**). An empty field indicates the corresponding tool’s inability to be run on a data set. Horizontal bars show the overall performance of a given tool across data sets on which the tool was tested (dark blue bar) and the number of such data sets (gray bar). An extended version of this figure including values of the performance score is provided in [Supplementary Table 14](#).

TABLES

Table 1. Alignment-free sequence comparison tools included in this study.

Software	Approach class	Software version	Availability
AAF ⁵³	exact k-mer count	10/01/2017	https://github.com/fanhuan/AAF
AFKS ³²		1.0	https://github.com/TulsaBioinformaticsToolsmith/Alignment-Free-Kmer-Statistics
alfpy ³		1.0.6	https://github.com/aziele/alfpy
CAFÉ ³⁴		1.0.0	https://github.com/younglululu/CAFE
FFP ^{33,38}		2v.2.1	https://github.com/jaejinchoi/FFP
jD2Stat ³⁵		1.0	http://bioinformatics.org.au/tools/jD2Stat/
LZW-Kernel ⁷⁹	information theory	NA	https://github.com/kfattila/LZW-Kernel
spaced ^{84–86}	inexact <i>k</i> -mer count	1.0	http://spaced.gobics.de
kWIP ⁷⁸	k-mer count	0.2.0-13-g3cf8a9e	https://github.com/kdmurray91/kWIP
ALFRED-G ⁷²	maximal length of exact common substrings	NA	https://alurulab.cc.gatech.edu/phylo
kmacs ^{18,85}		1.0	http://kmacs.gobics.de
kr ⁴⁹		2.0.2	http://guanine.evolbio.mpg.de/cgi-bin/kr2/kr.cgi.pl
Underlying Approach ⁹⁰		NA	http://www.dei.unipd.it/~ciompin/main/underlying.html
andi ²²	micro-alignments	0.02	https://github.com/EvolBioInf/andi
co-phylog ²¹		NA	https://github.com/yhg926/co-phylog
FSWM ²⁴ /Read-SpaM ⁷⁶		1.0	http://fswm.gobics.de

Multi-SpaM ²³		1.0	https://github.com/tdencker/multi-SpaM
phylum		0.3	https://github.com/kloetzi/phylum
mash ⁹	number of word matches	2.1	https://github.com/marbl/Mash
Slope-SpaM		0.1	https://github.com/burkhard-morgenstern/Slope-SpaM
Skmer ⁵²		2.0.2	https://github.com/shahab-sarmashghi/Skmer
RTD-Phylogeny ⁸²	return time distribution	1.0.1	https://github.com/pandurang-kolekar/rtd-phylogeny
kSNP3 ⁷⁷	SNP count	3.1	https://sourceforge.net/projects/ksnp/files/
EP-sim ⁷³	variable-length word counts	1.0	http://www.dei.unipd.it/~ciompin/main/EP-sim.html

Detailed information about the tools parameter values used in this study for different reference data sets is provided in [Supplementary Table 1](#). A concise description of the listed tools is provided in the **Methods**.

Table 2. Overview of the reference data sets.

Category	Name	# Sequences	Average sequence length	# Files	# Sequence comparisons
Regulatory elements detection	Cis-regulatory modules (CRMs) ⁴	370	764 nt	370	68,256
Protein sequence classification	Low sequence identity (<40%) ⁹³	1,066	180 aa	1,066	567,645
	High sequence identity (≥40%) ⁹³	2,128	184 aa	2,128	2,263,128
Gene trees Inference	SwissTree ⁴³	651	398 aa	651	211,575
Genome-based phylogeny	<i>Assembled genomes:</i>				
	29 <i>E. coli/Shigella</i> strains	29	4,895,247 nt	29	406
	14 plant species	14	337,515,688 nt	14	91
	25 fish mitochondrial genomes ⁴⁸	25	16,623 nt	25	300
	<i>Unassembled genomes</i>				
	29 <i>E. coli/Shigella</i> strains				
	coverage 0.03125	29,557	150 nt	29	406
	coverage 0.0625	59,116	150 nt	29	406
	coverage 0.125	118,266	150 nt	29	406
	coverage 0.25	236,541	150 nt	29	406
	coverage 0.5	473,081	150 nt	29	406
	coverage 1	946,169	150 nt	29	406
coverage 5	4,730,778	150 nt	29	406	

	14 plant species				
	coverage 0.015625	48,274	150 nt	14	91
	coverage 0.03125	96,489	150 nt	14	91
	coverage 0.0625	1,931,268	150 nt	14	91
	coverage 0.125	3,862,905	150 nt	14	91
	coverage 0.25	7,725,928	150 nt	14	91
	coverage 0.5	15,461,718	150 nt	14	91
	coverage 1	30,903,727	150 nt	14	91
Horizontal gene transfer	27 <i>E. coli/Shigella</i> genomes ⁵⁵	27	4,905,896 nt	27	351
	8 <i>Yersinia</i> species ⁵⁷	8	4,605,553 nt	8	28
	33 simulated genomes ⁵⁴				
	HGT level 0	33	2,205,524 nt	33	528
	HGT level 250	33	2,149,620 nt	33	528
	HGT level 500	33	2,230,317 nt	33	528
	HGT level 750	33	2,263,926 nt	33	528
	HGT level 1,000	33	2,238,661 nt	33	528

An interactive visualization of all results for all data sets can be found online (<http://afproject.org>).

REFERENCES

1. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
2. Thompson, J. D., Higgins, D. G. & Gibson, T. J. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**, 4673–4680 (1994).
3. Zielezinski, A., Vinga, S., Almeida, J. & Karlowski, W. M. Alignment-free sequence comparison: benefits, applications, and tools. *Genome Biol.* **18**, 186 (2017).
4. Kantorovitz, M. R., Robinson, G. E. & Sinha, S. A statistical method for alignment-free comparison of regulatory sequences. *Bioinformatics* **23**, i249–55 (2007).
5. Ivan, A., Halfon, M. S. & Sinha, S. Computational discovery of cis-regulatory modules in *Drosophila* without prior knowledge of motifs. *Genome Biol.* **9**, R22 (2008).
6. Vinga, S., Gouveia-Oliveira, R. & Almeida, J. S. Comparative evaluation of word composition distances for the recognition of SCOP relationships. *Bioinformatics* **20**, 206–215 (2004).

7. Terrapon, N., Weiner, J., Grath, S., Moore, A. D. & Bornberg-Bauer, E. Rapid similarity search of proteins using alignments of domain arrangements. *Bioinformatics* **30**, 274–281 (2014).
8. Cong, Y., Chan, Y.-B. & Ragan, M. A. A novel alignment-free method for detection of lateral genetic transfer based on TF-IDF. *Sci. Rep.* **6**, 30308 (2016).
9. Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using MinHash. *Genome Biol.* **17**, 132 (2016).
10. Fox, G. E., Magrum, L. J., Balch, W. E., Wolfe, R. S. & Woese, C. R. Classification of methanogenic bacteria by 16S ribosomal RNA characterization. *Proc. Natl. Acad. Sci. U. S. A.* **74**, 4537–4541 (1977).
11. Vinga, S. & Almeida, J. Alignment-free sequence comparison--a review. *Bioinformatics* **19**, 513–523 (2003).
12. Jun, S.-R., Sims, G. E., Wu, G. A. & Kim, S.-H. Whole-proteome phylogeny of prokaryotes by feature frequency profiles: An alignment-free method with optimal feature resolution. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 133–138 (2010).
13. Sims, G. E. & Kim, S.-H. Whole-genome phylogeny of Escherichia coli/Shigella group by feature frequency profiles (FFPs). *Proc. Natl. Acad. Sci. U. S. A.* **108**, 8329–8334 (2011).
14. Blaisdell, B. E. A measure of the similarity of sets of sequences not requiring sequence alignment. *Proc. Natl. Acad. Sci. U. S. A.* **83**, 5155–5159 (1986).
15. Reinert, G., Chew, D., Sun, F. & Waterman, M. S. Alignment-free sequence comparison (I): statistics and power. *J. Comput. Biol.* **16**, 1615–1634 (2009).
16. Wan, L., Reinert, G., Sun, F. & Waterman, M. S. Alignment-free sequence comparison (II): theoretical power of comparison statistics. *J. Comput. Biol.* **17**, 1467–1490 (2010).
17. Ulitsky, I., Burstein, D., Tuller, T. & Chor, B. The average common substring approach to phylogenomic reconstruction. *J. Comput. Biol.* **13**, 336–350 (2006).
18. Leimeister, C.-A. & Morgenstern, B. Kmacs: the k-mismatch average common substring

- approach to alignment-free sequence comparison. *Bioinformatics* **30**, 2000–2008 (2014).
19. Yang, L., Zhang, X., Fu, H. & Yang, C. An estimator for local analysis of genome based on the minimal absent word. *J. Theor. Biol.* **395**, 23–30 (2016).
 20. Yang, L., Zhang, X. & Zhu, H. Alignment free comparison: similarity distribution between the DNA primary sequences based on the shortest absent word. *J. Theor. Biol.* **295**, 125–131 (2012).
 21. Yi, H. & Jin, L. Co-phylog: an assembly-free phylogenomic approach for closely related organisms. *Nucleic Acids Res.* **41**, e75 (2013).
 22. Haubold, B., Klötzl, F. & Pfaffelhuber, P. andi: fast and accurate estimation of evolutionary distances between closely related genomes. *Bioinformatics* **31**, 1169–1175 (2015).
 23. Dencker, T. *et al.* Multi-SpaM: A Maximum-Likelihood Approach to Phylogeny Reconstruction Using Multiple Spaced-Word Matches and Quartet Trees. in *Lecture Notes in Computer Science* 227–241 (2018).
 24. Leimeister, C.-A., Sohrabi-Jahromi, S. & Morgenstern, B. Fast and accurate phylogeny reconstruction using filtered spaced-word matches. *Bioinformatics* **33**, 971–979 (2017).
 25. Leimeister, C.-A. *et al.* Prot-SpaM: fast alignment-free phylogeny reconstruction based on whole-proteome sequences. *Gigascience* **8**, (2019).
 26. Almeida, J. S., Carrico, J. A., Marezek, A., Noble, P. A. & Fletcher, M. Analysis of genomic sequences by Chaos Game Representation. *Bioinformatics* **17**, 429–437 (2001).
 27. Jeffrey, H. J. Chaos game representation of gene structure. *Nucleic Acids Res.* **18**, 2163–2170 (1990).
 28. Yau, S. S.-T., -T. Yau, S. S., Yu, C. & He, R. A Protein Map and Its Application. *DNA Cell Biol.* **27**, 241–250 (2008).
 29. Yin, C. & Yau, S. S.-T. An improved model for whole genome phylogenetic analysis by Fourier transform. *J. Theor. Biol.* **382**, 99–110 (2015).
 30. Vinga, S. Information theory applications for biological sequence analysis. *Brief. Bioinform.* **15**,

- 376–389 (2014).
31. Almeida, J. S. Sequence analysis by iterated maps, a review. *Brief. Bioinform.* **15**, 369–375 (2014).
 32. Luczak, B. B., James, B. T. & Girgis, H. Z. A survey and evaluations of histogram-based statistics in alignment-free sequence comparison. *Brief. Bioinform.* (2017).
doi:10.1093/bib/bbx161
 33. Sims, G. E., Jun, S.-R., Wu, G. A. & Kim, S.-H. Alignment-free genome comparison with feature frequency profiles (FFP) and optimal resolutions. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 2677–2682 (2009).
 34. Lu, Y. Y. *et al.* CAFE: aCcelerated Alignment-FrEe sequence analysis. *Nucleic Acids Res.* **45**, W554–W559 (2017).
 35. Chan, C. X., Bernard, G., Poirion, O., Hogan, J. M. & Ragan, M. A. Inferring phylogenies of evolving sequences without multiple sequence alignment. *Sci. Rep.* **4**, 6504 (2014).
 36. Robinson, D. F. & Foulds, L. R. Comparison of phylogenetic trees. *Math. Biosci.* **53**, 131–147 (1981).
 37. Altenhoff, A. M. *et al.* Standardized benchmarking in the quest for orthologs. *Nat. Methods* **13**, 425–430 (2016).
 38. Choi, J. & Kim, S.-H. A genome Tree of Life for the Fungi kingdom. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 9391–9396 (2017).
 39. Fox, N. K., Brenner, S. E. & Chandonia, J.-M. SCOPE: Structural Classification of Proteins--extended, integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Res.* **42**, D304–9 (2014).
 40. Wu, T. J., Burke, J. P. & Davison, D. B. A measure of DNA sequence dissimilarity based on Mahalanobis distance between frequencies of words. *Biometrics* **53**, 1431–1439 (1997).
 41. Hide, W., Burke, J. & Davison, D. B. Biological evaluation of d2, an algorithm for

- high-performance sequence comparison. *J. Comput. Biol.* **1**, 199–215 (1994).
42. Rokas, A., Williams, B. L., King, N. & Carroll, S. B. Genome-scale approaches to resolving incongruence in molecular phylogenies. *Nature* **425**, 798–804 (2003).
 43. Boeckmann, B. *et al.* Taxon sampling unequally affects individual nodes in a phylogenetic tree: consequences for model gene tree construction in SwissTree. (2017). doi:10.1101/181966
 44. Dai, Q., Yang, Y. & Wang, T. Markov model plus k-word distributions: a synergy that produces novel statistical measures for sequence comparison. *Bioinformatics* **24**, 2296–2302 (2008).
 45. Chan, C. X. & Ragan, M. A. Next-generation phylogenomics. *Biol. Direct* **8**, 3 (2013).
 46. Haubold, B. Alignment-free phylogenetics and population genetics. *Brief. Bioinform.* **15**, 407–418 (2014).
 47. Li, M. *et al.* An information-based sequence distance and its application to whole mitochondrial genome phylogeny. *Bioinformatics* **17**, 149–154 (2001).
 48. Fischer, C. *et al.* 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**, e67048 (2013).
 49. Haubold, B., Pfaffelhuber, P., Domazet-Lošćo, M. & Wiehe, T. Estimating Mutation Distances from Unaligned Genomes. *J. Comput. Biol.* **16**, 1487–1500 (2009).
 50. Lin, J., Adjero, D. A., Jiang, B.-H. & Jiang, Y. K2 and K2*: efficient alignment-free sequence similarity measurement based on Kendall statistics. *Bioinformatics* **34**, 1682–1689 (2018).
 51. Fabian, K. & Bernard, H. Phylonium - fast and accurate estimation of evolutionary distances. *GitHub* Available at: <https://github.com/kloetzl/phylonium>. (Accessed: 10th February 2019)
 52. Sarmashghi, S., Bohmann, K., P Gilbert, M. T., Bafna, V. & Mirarab, S. Skmer: assembly-free and alignment-free sample identification using genome skims. *Genome Biol.* **20**, 34 (2019).
 53. Fan, H., Ives, A. R., Surget-Groba, Y. & Cannon, C. H. An assembly and alignment-free method of phylogeny reconstruction from next-generation sequencing data. *BMC Genomics* **16**, 522

- (2015).
54. Bernard, G., Chan, C. X. & Ragan, M. A. Alignment-free microbial phylogenomics under scenarios of sequence divergence, genome rearrangement and lateral genetic transfer. *Sci. Rep.* **6**, 28970 (2016).
 55. Skippington, E. & Ragan, M. A. Within-species lateral genetic transfer and the evolution of transcriptional regulation in *Escherichia coli* and *Shigella*. *BMC Genomics* **12**, 532 (2011).
 56. Beiko, R. G., Harlow, T. J. & Ragan, M. A. Highways of gene sharing in prokaryotes. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 14332–14337 (2005).
 57. Darling, A. E., Miklós, I. & Ragan, M. A. Dynamics of genome rearrangement in bacterial populations. *PLoS Genet.* **4**, e1000128 (2008).
 58. Doolittle, W. F. & Bapteste, E. Pattern pluralism and the Tree of Life hypothesis. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 2043–2049 (2007).
 59. Dagan, T. & Martin, W. Getting a better picture of microbial evolution en route to a network of genomes. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **364**, 2187–2196 (2009).
 60. Bernard, G., Greenfield, P., Ragan, M. A. & Chan, C. X. -mer Similarity, Networks of Microbial Genomes, and Taxonomic Rank. *mSystems* **3**, (2018).
 61. Bernard, G., Ragan, M. A. & Chan, C. X. Recapitulating phylogenies using -mers: from trees to networks. *F1000Res.* **5**, 2789 (2016).
 62. Earl, D. *et al.* Assemblathon 1: a competitive assessment of de novo short read assembly methods. *Genome Res.* **21**, 2224–2241 (2011).
 63. Bradnam, K. R. *et al.* Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species. *Gigascience* **2**, 10 (2013).
 64. Baruzzo, G. *et al.* Simulation-based comprehensive benchmarking of RNA-seq aligners. *Nat. Methods* **14**, 135–139 (2017).
 65. Sczyrba, A. *et al.* Critical Assessment of Metagenome Interpretation—a benchmark of

- metagenomics software. *Nat. Methods* **14**, 1063–1071 (2017).
66. Chandonia, J.-M. *et al.* The ASTRAL Compendium in 2004. *Nucleic Acids Res.* **32**, D189–92 (2004).
 67. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open Software Suite. *Trends Genet.* **16**, 276–277 (2000).
 68. Tran, N. H. & Chen, X. Comparison of next-generation sequencing samples using compression-based distances and its application to phylogenetic reconstruction. *BMC Res. Notes* **7**, 320 (2014).
 69. Hatje, K. & Kollmar, M. A phylogenetic analysis of the brassicales clade based on an alignment-free sequence comparison method. *Front. Plant Sci.* **3**, 192 (2012).
 70. Huang, W., Li, L., Myers, J. R. & Marth, G. T. ART: a next-generation sequencing read simulator. *Bioinformatics* **28**, 593–594 (2012).
 71. Beiko, R. G. & Charlebois, R. L. A simulation test bed for hypotheses of genome evolution. *Bioinformatics* **23**, 825–831 (2007).
 72. Thankachan, S. V., Chockalingam, S. P., Liu, Y., Krishnan, A. & Aluru, S. A greedy alignment-free distance estimator for phylogenetic inference. *BMC Bioinformatics* **18**, 238 (2017).
 73. Comin, M. & Antonello, M. On the comparison of regulatory sequences with multiple resolution Entropic Profiles. *BMC Bioinformatics* **17**, 130 (2016).
 74. Fernandes, F., Freitas, A. T., Almeida, J. S. & Vinga, S. Entropic Profiler - detection of conservation in genomes using information theory. *BMC Res. Notes* **2**, 72 (2009).
 75. Comin, M. & Antonello, M. Fast Entropic Profiler: An Information Theoretic Approach for the Discovery of Patterns in Genomes. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **11**, 500–509 (2014).
 76. Lau, A. K., Leimeister, C.-A. & Morgenstern, B. Read-SpaM: assembly-free and alignment-free

- comparison of bacterial genomes with low sequencing coverage. *bioRxiv* (2019).
doi:10.1101/550632
77. Gardner, S. N., Slezak, T. & Hall, B. G. kSNP3.0: SNP detection and phylogenetic analysis of genomes without genome alignment or reference genome. *Bioinformatics* **31**, 2877–2878 (2015).
 78. Murray, K. D., Webers, C., Ong, C. S., Borevitz, J. & Warthmann, N. kWIP: The k-mer weighted inner product, a de novo estimator of genetic similarity. *PLoS Comput. Biol.* **13**, e1005727 (2017).
 79. Filatov, G., Bauwens, B. & Kertész-Farkas, A. LZW-Kernel: fast kernel utilizing variable length code blocks from LZW compressors for protein sequence classification. *Bioinformatics* **34**, 3281–3288 (2018).
 80. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
 81. Snir, S. & Rao, S. Quartet MaxCut: a fast algorithm for amalgamating quartet trees. *Mol. Phylogenet. Evol.* **62**, 1–8 (2012).
 82. Kolekar, P., Kale, M. & Kulkarni-Kale, U. Alignment-free distance measure based on return time distribution for sequence analysis: applications to clustering, molecular phylogeny and subtyping. *Mol. Phylogenet. Evol.* **65**, 510–522 (2012).
 83. Röhling, S. & Morgenstern, B. The number of spaced-word matches between two DNA sequences as a function of the underlying pattern weight. *bioRxiv* 527515 (2019).
doi:10.1101/527515
 84. Leimeister, C.-A., Boden, M., Horwege, S., Lindner, S. & Morgenstern, B. Fast alignment-free sequence comparison using spaced-word frequencies. *Bioinformatics* **30**, 1991–1999 (2014).
 85. Horwege, S. *et al.* Spaced words and kmacs: fast alignment-free sequence comparison based on inexact word matches. *Nucleic Acids Res.* **42**, W7–11 (2014).
 86. Morgenstern, B., Zhu, B., Horwege, S. & Leimeister, C. A. Estimating evolutionary distances

- between genomic sequences from spaced-word matches. *Algorithms Mol. Biol.* **10**, 5 (2015).
87. Ma, B., Tromp, J. & Li, M. PatternHunter: faster and more sensitive homology search. *Bioinformatics* **18**, 440–445 (2002).
 88. Li, M., Ma, B., Kisman, D. & Tromp, J. PATTERNHUNTER II: HIGHLY SENSITIVE AND FAST HOMOLOGY SEARCH. *J. Bioinform. Comput. Biol.* **02**, 417–439 (2004).
 89. Lin, J. Divergence measures based on the Shannon entropy. *IEEE Transactions on Information Theory* **37**, 145–151 (1991).
 90. Comin, M. & Verzotto, D. Alignment-free phylogeny of whole genomes using underlying subwords. *Algorithms Mol. Biol.* **7**, 34 (2012).
 91. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
 92. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).
 93. Brenner, S. E., Koehl, P. & Levitt, M. The ASTRAL compendium for protein structure and sequence analysis. *Nucleic Acids Res.* **28**, 254–256 (2000).