

1 *Article Paper*

2 **The utility of data transformation for alignment, *de*** 3 ***novo* assembly and classification of short read virus** 4 **sequences**

5 **Avraam Tapinos^{1*}, Bede Constantinides^{1,3}, My VT Phan², Samaneh Kouchaki^{1,4}, Matthew Cotten²,**
6 **and David L Robertson^{1,5}**

7 ¹ School of Biological Sciences, The University of Manchester, Manchester, UK

8 ² Department of Viroscience, Erasmus Medical Centre, Rotterdam, the Netherlands

9 ³ Nuffield Department of Medicine, University of Oxford, Oxford, UK

10 ⁴ Department of Engineering Science, University of Oxford, Oxford, UK

11 ⁵ MRC-University of Glasgow Centre for Virus Research, Glasgow, UK

12

13 * Correspondence avraam.tapinos@manchester.ac.uk; Tel.: +44 (0) 161 701 7563

14

Received: date; Accepted: date; Published: date

15 **Abstract:** Advances in DNA sequencing technology are facilitating genomic analyses of
16 unprecedented scope and scale, widening the gap between our abilities to generate and fully exploit
17 biological sequence data. Comparable analytical challenges are encountered in other data-intensive
18 fields involving sequential data, such as signal processing, in which dimensionality reduction (i.e.,
19 compression) methods are routinely used to lessen the computational burden of analyses. In this
20 work we explore the application of dimensionality reduction methods to numerically represent
21 high-throughput sequence data for three important biological applications of virus sequence data:
22 reference-based mapping, short sequence classification and *de novo* assembly. Despite using highly
23 compressed sequence transformations to accelerate the processes, our sequence processing approach
24 yielded comparable accuracy to existing approaches, and are ideally suited for sequences originating
25 from highly diverse virus populations. We demonstrate the application of our methodology to both
26 synthetic and real viral pathogen sequence data. Our results show that the use of highly compressed
27 sequence approximations can provide accurate results and that useful analytical performance can be
28 retained and even enhanced through appropriate dimensionality reduction of sequence data.

29 **Keywords:** Alignment, assembly, taxonomic classification, time series, data transformation, DWT,
30 DFT, PAA, data compression, compressive genomics

31

32 **1. Introduction**

33 Next generation sequencing (NGS) enables massively parallel determination of nucleotide
34 order within genetic material, making it possible to rapidly sequence the genomes of individuals,
35 populations and metagenomic samples [1-5]. However, the sequences generated by these
36 instruments are almost always considerably shorter in length than the genomic regions studied.
37 Genomic analyses often begin with the process of sequence assembly, where sequence fragments
38 (reads) are reconstructed into the larger sequences from which they originated. Computational
39 methods play a vital role in the assembly of short reads, and a variety of assemblers and related tools
40 have been developed in tandem with emerging sequencing platforms [6]. All subsequent analyses
41 and investigations depend upon the quality, accuracy and speed of this crucial sequence assembly
42 process.

43 There are many computational methods to generate consensus sequences representing the
44 genomes of coexisting species in a sample. Such approaches includes seed-and-extend alignment
45 methods using suffix array derivatives such as the Burrows-Wheeler Transform (BWT) for aligning
46 short reads informed by a known reference sequence [7,8], and graph-based methods employing
47 Overlap Layout Consensus (OLC) [9,10] and *de Bruijn* graphs of *k*-mers [11-13] for reference-free *de*
48 *nov*o sequence assembly. However, for sequencing projects to characterise genetic variation within
49 populations (deep sequencing), metagenomics and pathogen discovery, the effectiveness of the
50 aforementioned approaches varies considerably [14].

51 Samples with mixed viral infections, especially those comprising divergent variants, present a
52 number of analytical and computational problems. The use of a reference sequence, even the use of a
53 data specific generated sequence, can lead to valuable read information being discarded during the
54 alignment process [15]. On the other hand, while *de novo* approaches require little *a priori* knowledge
55 of target sequence composition, the methods are computationally intensive and their performance
56 scales poorly with datasets of increasing size [9]. Aggressive heuristics must be employed, to
57 traverse graphs and deal with mismatches, to reduce the running time of *de novo* assemblers, which
58 in turn can compromise assembly quality. Indexing structures such as the BWT and FM-index are
59 widely used to reduce the burden of pairwise sequence comparison, for both reference base
60 mapping and *de novo* assembly. However, they cannot process mismatches within reads,
61 necessitating the use of aggressive and computationally expensive heuristics to establish
62 relationships on divergent reads. The increasing NGS read length will affect the performance of
63 these approaches [16].

64 A major challenge working with high-throughput sequencing data for metagenomics and
65 within-host variation analysis is the substantially great diversity of biological data. Also the amount
66 of sequences generated challenge many computational system for a feasible working solution in
67 terms of time and the computational resources typically available in biological laboratories. For
68 biologists working in outbreak responses or pathogen discovery, both the accuracy of the assembly
69 results and the speed of sequence analyses (e.g. assembly, alignment and pathogen classification) are
70 crucial for crisis response and management. The ability to run analyses in the field on portable
71 computer systems without internet connectivity is also important. Here, we explore the utility of
72 data transform methods (used in data intensive fields such as signal and image processing) to extract
73 major features from viral NGS sequence data and use the features to analyse data in a lower
74 dimensional space.

75 Similar analytical challenges involving high dimensional sequential data are encountered in
76 other data-intensive fields such as signal and image processing, and time series analysis, where data
77 transforms and approximation techniques are used for data dimensionality reduction. Data
78 transform/approximation techniques include the discrete Fourier transform (DFT) [17], the discrete
79 wavelet transform (DWT) [18,19], and piece-wise aggregate approximation (PAA) [20,21]. The DFT
80 or DWT are used to transform data to their frequency domains, allowing feature extraction [22], and
81 PAA is used as a data approximation approach. In data intensive fields, data
82 transformations/approximations are commonly used as dimensionality reduction approaches, for
83 obtaining fast approximate solutions for a given problem. Due to the ordered nature of genetic data,
84 many of these transformation approaches can be applied to sequences of nucleotides [23] or amino
85 acids [24]. An example of a successful implementation of a Fourier transform in computational
86 biology is the alignment algorithm MAFFT [25] where the physiochemical properties of amino acids
87 are used to represent sequences for fast matching of homologous sequence regions for alignment.
88 Since most transformation approaches are suitable only for numerical sequences, the strings of
89 letters representing genetic sequences must be mapped into numerical space using a numerical
90 sequence representation method [26].

91 In addition to the DFT, the DWT and PAA, suitable methods for measuring the pairwise
92 similarity of sequential data or transformations include the L_p -norms [27], dynamic time warping
93 (DTW) [28], longest common subsequence (LCS) [29] and alignment approaches such as the
94 Needleman-Wunsch and Smith-Waterman algorithms. Euclidean distance is arguably the most

95 widely used Lp-norm method for sequential data comparison but can only be used on sequences
 96 with same lengths. Furthermore, Lp-norm methods do not accommodate shifts in the x-axis (time or
 97 position) and are thus limited in their ability to identify similar features within offset data. Elastic
 98 similarity/dissimilarity methods such as LCS, unbounded DTW and various alignment algorithms
 99 permit comparison of data with different dimensions and tolerate shifts in the x-axis. These
 100 properties of elastic similarity methods can be very useful in the analysis of speech signals, for
 101 example, but can be computationally expensive [30]. Several approaches have been proposed to
 102 permit fast search with DTW, including the introduction of global constraints (wrapping path) or the
 103 use of lower bounding techniques such as LB_keogh [28].

104 While pairwise comparison methods may be used for clustering, classification and similarity
 105 searches, they are very time consuming for large datasets ($O(n^2)$ time complexity). Indexing
 106 structures such as the R*-tree, KD-tree, VP-tree and MVP-tree have significantly lower time
 107 complexity ($O(n \log(n))$) for similarity search [31] and are more appropriate for efficient analysis of
 108 large datasets. The R*-tree [32,33] and KD-tree [34] indexing structures are very accurate for low
 109 dimensional datasets. However, their performance deteriorates significantly in high dimensional
 110 space [31], known as the 'curse of dimensionality' [35,36]. Metric trees such as the VP-tree [37] and
 111 MVP-tree [38] are less prone to this limitation. Metric space indexing structures make use of
 112 geometric properties for partitioning data and work efficiently on both low and high dimensional
 113 data [39]. The curse of dimensionality can be further mitigated using data approximations such as
 114 the DFT, the DWT, and the PAA to partition a dataset in an approximated space without loss of
 115 generality [21].

116 Here we investigate the performance of three established dimensionality reduction techniques,
 117 on three common analysis tasks involving viral short read sequence data: classification, reference
 118 based mapping/alignment, and *de novo* assembly. We benchmarked the accuracy of our proposed
 119 methodology against existing tools, and demonstrated the applicability of time series and signal
 120 processing data mining techniques for the analysis of viral NGS data.

121 2. Materials and Methods

122 2.1. Symbolic to numeric sequence representations

123 Various numeric sequence representation methods can be used for symbolising a nucleotide
 124 sequence to a numerical space (see Table 1). Depending on the chosen numerical representation,
 125 each nucleotide is associated with a specific numerical value or vector. The specific values are
 126 assigned to the position of each nucleotide indicating the presence of a nucleotide at each sequence
 127 position (Equation 1). R_i is the indicator for a specific nucleotide in the i^{th} position of the sequence S
 128 with a length of n nucleotides. Values $v_1 \dots v_5$ correspond to the numerical value or numerical vector
 129 associated with each nucleotide.

$$R_i = \begin{cases} v_1 & \text{if } S_i = A \\ v_2 & \text{if } S_i = T \\ v_3 & \text{if } S_i = C \\ v_4 & \text{if } S_i = G \\ v_5 & \text{otherwise} \end{cases}, \forall i \in S_n \quad (1)$$

130

131

132 Methods such as the electron-ion interaction pseudo potentials (EIIP) [40] and the atomic
 133 representation approach [41] aim to mimic the biochemical properties of nucleic acids, but introduce
 134 some mathematical bias that does not exist in reality [26]. Other methods, like the Voss inductor [42]
 135 and the Tetrahedron approach, do not introduce internucleotide mathematical bias, meaning the
 136 pairwise distances between each non-identical transformed nucleotide is the same (for example, the
 137 distance between A and T is equal to the distances between A and C as well as A and G).

138 Furthermore, the cumulative sum of a numerical representation R can be used to indicate the
139 trajectory of a sequence in nucleotide space. Table 1 indicates the associate values used for different
140 representation methods [26].

141 2.2. Sequence Transformation

142 Effective methods for transforming/approximating sequential data should: *i)* accurately
143 transform/approximate data without loss of useful information, *ii)* have low computational
144 overheads, *iii)* facilitate rapid comparison of data, and *iv)* provide lower bounding—where the
145 distance between data representations is always smaller than or equal to that of the original
146 data—guaranteeing against false negative results [43]. We employ the DFT and the DWT
147 transformation methods and PAA approximation method as they satisfy the above requirements,
148 are widely used for analysing discrete signals [44], and can be used to transform/approximate
149 nucleotide sequence numerical representations to different levels of resolution, permitting reduced
150 dimensionality sequence analysis.

151 Figure 1A illustrates an example of the DFT and DWT transformations and PAA approximation
152 of a short nucleotide sequence. DFT and the fast Fourier transform (FFT) transform data from their
153 original domain to a frequency domain. In principle, the DFT decomposes a numerically represented
154 nucleotide sequence with n positions (dimensions) into a series of n frequency components ordered
155 by their frequency. A subset of the resulting Fourier frequencies are used to approximate the original
156 sequence in a lower dimensional space [17], and the tradeoff between analytical speed and accuracy
157 can be varied according to the number of frequencies considered [45]

158 DWT transform data from their original domain to their time-frequency, accommodating for
159 changes in signal frequency over time [18,46,47]. DWT is a set of averaging and differencing
160 functions that may be used recursively to represent sequential data at different resolutions and each
161 resolution can be used as an approximation of the original data. Figure 1B depicts an example of the
162 DWT transformations of a short nucleotide sequence.

163 In PAA a numerical sequence is divided into n equally sized windows, the mean values of
164 which together form a compressed sequence representation [20,21]. The selection of n determines the
165 resolution of the compressed or approximate representation. While PAA is faster and easier to
166 implement than the DFT and the DWT, unlike these two methods it is irreversible, meaning that the
167 original sequence cannot be recovered from the approximation. Figure 1C depicts an example of the
168 PAA transformations of a short nucleotide sequence.

169 2.3. Similarity search approaches for sequential data

170 Here we adopt the Euclidian distance and VP-indexing tree to partition our data based on their
171 approximate space distances and perform a fast k -nearest neighbor (k -NN) similarity search for
172 aligning the reads to the reference genome.

173 In a VP-tree indexing structure, data partitioning is implemented in metric space. A data point
174 which is used as a vantage point is selected (either randomly or by applying some heuristic to find
175 and use the furthest point in the dataset [37]), and the rest of the data points are partitioned into two
176 nodes based on their distance to that point. Data found to be closer to the vantage point than a given
177 threshold (the median distance between all the data points and the vantage point) are assigned to the
178 same node, and the rest of the data points to a different node. This function is repeated recursively in
179 order to complete the partitioning process. The resulting indexing structure can then be used for fast
180 identification of a k -nearest neighbour (k -NN) search. A k -NN-search returns the data points that are
181 closest to a query q . Initially, the distance between the query q and the vantage point in the top node
182 is calculated. If the distance between q and the vantage point satisfies a set of given conditions (the
183 distance is smaller or larger than a given threshold – this threshold being the median distance
184 between the vantage point and other data points within the node), a decision is made to visit either
185 one or both of the nodes. This process is repeated until the entire tree has been traversed. The k data
186 points—in this case reads— found closest to our query are the k -nearest neighbours to the query q .

187

188 2.4. Proposed short reads processing methodology

189 Our methodology for taxonomic classification, reference based mapping and *de novo* assembly
190 of short reads, used time series and digital signal processing data transformation techniques. Figure
191 2 illustrates the fundamental concept of our approach. The short reads and reference genomes are
192 mapped to a numerical space using an appropriate method from table 1. Subsequently, lower
193 dimensional approximations are generated for all data using the appropriate data transformation
194 method, such as DFT, DWT, and PAA. A VP-tree is constructed to allow fast data comparison.
195 Depending on the application, the VP-tree is constructed either by using *k*-mer transformations
196 obtained from the reference genomes or by using the short reads' transformations. Consequently,
197 the best matches for our short reads' transformations are identified using a *k*-NN search approach on
198 the VP-tree. As a final step, the results obtained from the *k*-NN search are re-evaluated in the original
199 space to remove potential false positive results.

200 2.5. Data

201 The implementations of our proposed methodologies were assessed with both simulated and
202 real virus datasets. The simulated datasets were generated using CuReSim [48] and WGSIM
203 (<https://github.com/lh3/wgsim>). Simulated data included information such as the reference genome
204 used, the alignment position, and alignment direction for each read, enabling rigorous evaluation of
205 the proposed techniques. We used two simulators to generate different datasets to test our approach
206 under every eventuality. CuReSim can generate long Ion Torrent reads, allowing the user to control
207 the type of variation (insertion, deletion and substitution) to simulate. The user can also define the
208 extent of variation for each individual type of variation. WGSIM can simulate diploid genomes with
209 uniform insertion, deletion and substitution sequencing variation.

210 CuReSim was used to generate 16 HIV-HXB2 simulated datasets with different levels and
211 types of variation. WGSIM was used to generate 4 mixed virus datasets with different levels of
212 variation. Each simulation contained 200,000 reads generated using 5 Norovirus genomes, 5 Ebola
213 virus genomes, and 5 Respiratory syncytial virus (RSV) genomes, with various types and extents of
214 simulated variation. Table 2 contains detailed information about the simulated datasets.

215 Furthermore, 15 publicly available real virus datasets were used for the evaluation of our
216 methodology. The real datasets comprise 5 Norovirus, 5 Ebola virus, and 5 human respiratory
217 syncytial virus (RSV) short read datasets. Norovirus NGS datasets (ERR225628, ERR225629,
218 ERR225631, ERR225632, ERR225633) were generated from diarrhoeal patients in Vietnam [49].
219 Group A rotavirus datasets were obtained from human and pig samples from Vietnam [50]. Human
220 coronavirus NL63 datasets were obtained from Kenya [51]. The Ebola virus datasets (SRR3107337,
221 SRR3107338, SRR3107340, SRR3107342, SRR3107343) were retrieved from the bioproject
222 PRJNA309162, generated during the outbreaks in West Africa in 2013-2016 [52]. The human
223 respiratory syncytial virus (RSV) datasets, ERR303259, ERR303260, ERR303261, ERR303262,
224 ERR303263 [53], were generated from humans in Kenya. Further information for the real data can be
225 found in Table 3.

226 The HIV-HXB2 genome (K03455) was used as a reference index to align and or run the
227 taxonomic classification analysis for the HIV-HXB2 simulated dataset. The Norovirus genome
228 KM198509, the Ebola virus genome KM034562, and the RSV genome KP317934, were used as a
229 reference index to align and or run the taxonomic classification analysis for the mixed virus datasets.
230 The Norovirus genome KM198509 was used to run the taxonomic classification analysis on the real
231 Norovirus datasets, the Ebola virus genome KM034562 was used to run the taxonomic classification
232 analysis on the real Ebola datasets, and the RSV genome KP317934, was used to perform the
233 taxonomic classification analysis on the real RSV datasets. Further information for the reference
234 genomes can be found in Table 4.

235 2.6. Classification and alignment evaluation

236 The accuracy of a classification and an alignment tool can be quantified in terms of F-measure
237 [48], a balanced measure of precision and recall. With $\text{precision} = \frac{\text{true positive}}{\text{true positive} + \text{false positive}}$
238 and $\text{recall} = \frac{\text{true positive}}{\text{true positive} + \text{false negative}}$, then $\text{F-measure} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$ [48]. In the case of simulated data, information concerning the position of
239 the read on the reference and alignment direction can be used to establish the correctness of an
240 alignment, and thereby provide a more informative F-measure score. Unclassified reads are
241 considered a false negative result. Any hit on the correct section of the genome and in the correct
242 direction report is considered a true positive result. However, if the alignment position and direction
243 information are not available, the F-measure can be calculated from the number of hits reported for a
244 read, or the absence of a hit. Again, unclassified reads are considered false negative results, and
245 classified reads are considered true positive results. In the case of mixed genome data, the F-measure
246 score can be calculated by taking in consideration the number of hits that were reported for a read,
247 as well as if a read was assigned to a reference genome from the same family. If a read is assigned to
248 a genome from a different virus family, it is considered a false positive result, and unclassified reads
249 are considered a false negative result.
250

251 3. Results

252 3.1. Classification by numbers (CBN)

253 For the taxonomic classification analysis, a classification tool was implemented in C++
254 (<https://github.com/Avramis/ClassificationByNumber>). The implementation was developed to
255 evaluate our methodology but is not optimised for speed. Users may specify parameters such as the
256 representation method, transformation method, search stringency and the k -mer length. A
257 vantage-point tree (VP-Tree) indexing structure, containing information from a set of given genomic
258 references, is used to classify reads. The initial step to build the VP-Tree is to extract all of the unique
259 k -mers, of a user-specified size k , from a set of given reference genomes. Each unique k -mer is
260 represented in numerical sequence, and then transformed to a lower dimensional space. The
261 transformed data are then used to generate the VP-Tree indexing structure. Subsequently, each short
262 read from a query set is converted into numerical space and transformed to a lower dimensional
263 space and evaluated against the VP-tree. The approximate solution is then evaluated using the
264 original data to identify false positive results. The CBN algorithm generates two output files. The
265 first output is a text file providing detailed information on all of the classification hits generated for
266 each read. The information generated for each classification includes the reference name, the
267 direction in which the query read was aligned to the reference, the start and end position of the
268 query on the reference, the alignment score, the CIGAR string, and the actual alignment of the query
269 read on the reference genome. The second output file provides a brief review of the alignment. In
270 every line, it contains the short read's name, the number of classifications generated for that
271 particular read, the highest classification score obtained, the name of the reference which provided
272 the highest classification score, the alignment direction and starting position on the reference.

273 The CBN tool was evaluated against NCBI-BLAST 2.8.1 BLASTn [54] and Kaiju 1.6.3 [55]
274 classifier tools. BLASTn performs the analysis in nucleotide space, whereas Kaiju translates
275 nucleotide sequences into every possible protein frame and performs the analysis in protein space.
276 Figures 3, 4, and 5 illustrate the results of the classification evaluation process. Figure 3 shows the
277 results obtained from the classification process on the HIV-HXB2 data. Figure 4 illustrates the results
278 of the mixed virus datasets, and Figure 5 illustrates the results obtained from the real data.

279 In the taxonomic classification of HIV-HXB2 simulated data, where the short reads were
280 classified against the genome that was used to generate them, Kaiju reported the highest accuracy
281 scores. CBN outperformed BLASTn in most cases and only following behind in on the high variation
282 rates datasets. For the mixed viruses simulated datasets, where reads were classified against
283 species strains related to the ones used to generate the reads, BLASTn manage to correctly classify
284 data to the correct species, followed closely by CBN, and Kaiju falling last. In the evaluation of the
285 tools on the real data, where reads were classified against genome strains related to their respecting

286 species, CBN generated more accurate results compared to the other tools, followed by Kaiju, and
287 BLASTn in third.

288 3.2. Alignment by numbers (ALBN)

289 To test the applicability of sequential data transformations and feature selection for read
290 alignment, we implemented a prototype k -NN read aligner (Table 5) in C++ (available at
291 https://github.com/Avramis/Alignment_by_numbers). As with the CBN classification analysis, the
292 ALBN code was not optimised for speed, and users may specify parameters such as the
293 representation method, transformation method, search stringency and the k -mer length used for
294 seeding alignments. The algorithm's output was used to construct gapped alignments in the widely
295 used Sequence Alignment/Map (SAM) file format.

296 The ALBN tool was evaluated against a set of well established, widely used, stated of the art
297 tools such as, bowtie2 (version 2.3.1) [56], BWA-MEM (version 0.7.16) [7], Graphmap (version 0.5.2)
298 [57], and Segmehl (version 0.3.4) [58]. Each aligner's accuracy was quantified in terms of F-measure
299 [48]. CuReSim provides information such as the simulated read's origin on the reference genome
300 and its alignment direction, enabling evaluation of each aligner's output an calculation of alignment
301 accuracy in terms of F-measure. For mixed virus datasets, tool performance was evaluated in terms
302 of ability to match and align reads to the appropriate virus reference genome. For the real data,
303 F-measures were calculated according to the number of reads aligned to the given genome or
304 otherwise.

305 Figures 6, 7 and 8 illustrate the F-measures obtained by evaluating alignments from each
306 aligner. Figure 6 illustrates alignment performance for each of the 16 datasets simulated using the
307 K03455 HIV-HXB2 reference genome. Figure 7 illustrates alignment performance for virus reads
308 simulated with Norovirus genome KM198509.1, Ebola genome KM034562.1, and the RSV genome
309 KP317934.1. Figure 8-i to 8-iii illustrates alignment performance (F-measure) for alignments of real
310 Norovirus, Ebola virus and RSV sequences against the same corresponding reference genomes used
311 in the simulations.

312 ALBN provided highly accurate results in all data cases. Regarding the HIV-HXB2 data, where
313 the short reads were aligned to genome that was used to generate them, ALBN provided the most
314 accurate results in all 16 cases, followed by Bowtie2 in terms of accuracy. Also, in the case of the
315 mixed virus datasets, where reads were aligned against genome strains related to those used to
316 generate the dataset, ALBN provided the most accurate results, with Graphmap and BWA-MEM
317 third and fourth respectively. ALBN also generated the most accurate alignment results using real
318 data, where reads where aligned against a species-specific reference genome.

319 3.3. De novo assembly by numbers

320 Lastly, to test the applicability of this approach to the *de novo* assembly of short reads, we
321 implemented a naïve algorithm for all-against-all k -mer comparison using data
322 transformations/approximation. Figure 9 illustrates the main concept of our *de novo* assembly
323 approach. For the ASBN tool, reads are represented as numerical sequences using an appropriate
324 numerical representation method (Table 1). Here we use the tetrahedron numerical representation
325 approach. Every k -mer of each numerically represented read is identified and transformed to lower
326 dimensional space using the chosen transformation method. All k -mers' transformations are used to
327 build a VP-tree, to allow for fast data comparison. Afterwards, all k -mers are compared to the rest of
328 the data using the VP-tree index. Information from the data comparison is used to construct a
329 weighted graph similar to that shown in Figure 9A. The shortest path on the weighted graph is
330 identified with a breadth-first search (BFS) (Figure 9B). Reads overlaps are used to generate an OLC
331 alignment of short reads (Figure 9C).

332 The ASBN assembler was compared with Megahit (version 1.1.3) [59], and SPAdes (version
333 3.13.0) [60] *de novo* assemblers on the HIV-HXB2 and mixed virus simulated datasets accordingly.
334 The derived contigs from each assembler were evaluated against the reference genomes used to
335 generate the data simulations with BLASTn [54]. Contig positions across the genome were collected

336 and plotted in a histogram. Secondly, a measure of assembly progress was plotted on an X-Y matrix
337 with X being the total coverage of the genomes generated and Y being the total number of gaps in
338 the coverage. A perfect assembly would have $X = \text{full genome length}$ and $Y = 0$, implicating, that the
339 contig matches the genome in terms of size and nucleotide sequences, without any mismatches or
340 gaps. For the HIV-HXB2 datasets, the contigs were evaluated against the K03455 genome, and the
341 contigs obtained from the mixed virus datasets were evaluated against the 15 different genomes,
342 KM198529, KM198528, KM198511, KM198500, KM198486, KU296608, KU296553, KU296549,
343 KU296528, KU296416, KP317952, KP317946, KP317934, KP317923, and KP317922.

344 Figure 10 illustrates the assembly results of SPAdes, Megahit and all three variants of ASBN on
345 the 16 simulated HIV-HXB2 datasets, and Figure 11 illustrates the assembly results on the mixed
346 virus simulated databases. Although ASBN processes data and assembles short reads in a lower
347 dimensional space, it nevertheless generated contigs that collectively cover the expected genome
348 length and provided comparable results to both state of the art *de novo* assemblers in this
349 experiment (Figure 10, Figure 11). In all cases, ASBN generated contigs spanning the whole genomes
350 of their respective viral species.

351 4. Discussion

352 Although well-established data compression methods for reversible compression of
353 one-dimensional and multivariate signals, images, text and binary exist [61-63]; there are very few
354 examples of their application to biological sequence data. We developed algorithms incorporating
355 signal compression methods for three common biological sequence analysis problems: classification,
356 alignment and *de novo* assembly of NGS short read virus data. Our results show that our approach
357 permits accurate classification and reference alignment in spite of high rates of sequence variation or
358 the use of a divergent reference genome. Data approximation/summarisation techniques such as the
359 DFT, the DWT and the PAA can be used to extract major features of sequence data and to suppress
360 noise or low-level variation. This allows data comparison in terms of major characteristics, thus
361 enabling the identification of similarities among data that might otherwise be concealed by minor
362 variation or noise.

363 Collectively these results demonstrate that complete nucleotide-level sequence resolution is not
364 a prerequisite of accurate sequence analysis, and that analytical performance can be preserved or
365 even enhanced through appropriate dimensionality reduction (compression) of sequences. While
366 our implementations use *k*-mers, the nature of the transformation/compression methods used
367 showed that optimal *k*-mer selection is far less important than with conventional exact *k*-mer
368 matching methods. The inherent error tolerance of the approach also permits use of larger *k* values
369 than are normally used with conventional sequence comparison algorithms, reducing the
370 computational burden of pairwise comparison, and thus, in *de novo* assembly specifically, the
371 complexity of building and searching an assembly graph.

372 Efficient mining of terabase-scale biological sequence datasets requires looking beyond
373 substring-indexing algorithms towards more versatile methods of compression for both data storage
374 and analysis. The use of probabilistic data structures can considerably reduce the computer memory
375 required for in-memory sequence lookups at the expense of a few false positives, and Bloom filters
376 and related data structures have seen broad application in *k*-mer centric tasks such as error
377 correction [64], *in silico* read normalisation [65] and *de novo* assembly [66,67]. However, while these
378 hash-based approaches perform well on datasets with high sequence redundancy, for large datasets
379 with many distinct *k*-mers, large amounts of memory are still necessary [65]. Lower bounding
380 transformations and approximation methods (such as the DFT, the DWT and PAA) can exhibit the
381 same attractive one-sided error offered by these probabilistic data structures, but instead of hash
382 tables use concrete and thus reusable sequence representations.

383 Furthermore, transformations allow compression of standalone sequence composition,
384 enabling flexible reduction of sequence resolution according to analytical requirements, so that
385 redundant sequence precision need not hinder analysis. While the problem of read alignment to a
386 known reference sequence is largely considered solved, assembly of large genomes remains a

387 formidable problem in computing. Moreover, consideration of the metagenomic composition of
 388 mixed biological samples, as demonstrated, further extends the scope and scale of the assembly
 389 problem beyond what is tractable using conventional sequence comparison approaches. By
 390 implementing a reference-based aligner and *de novo* assembler, we have demonstrated that using
 391 compressed numerical representations represents a tractable and versatile approach for
 392 reconstructing genomes and metagenomes sequenced with short reads.

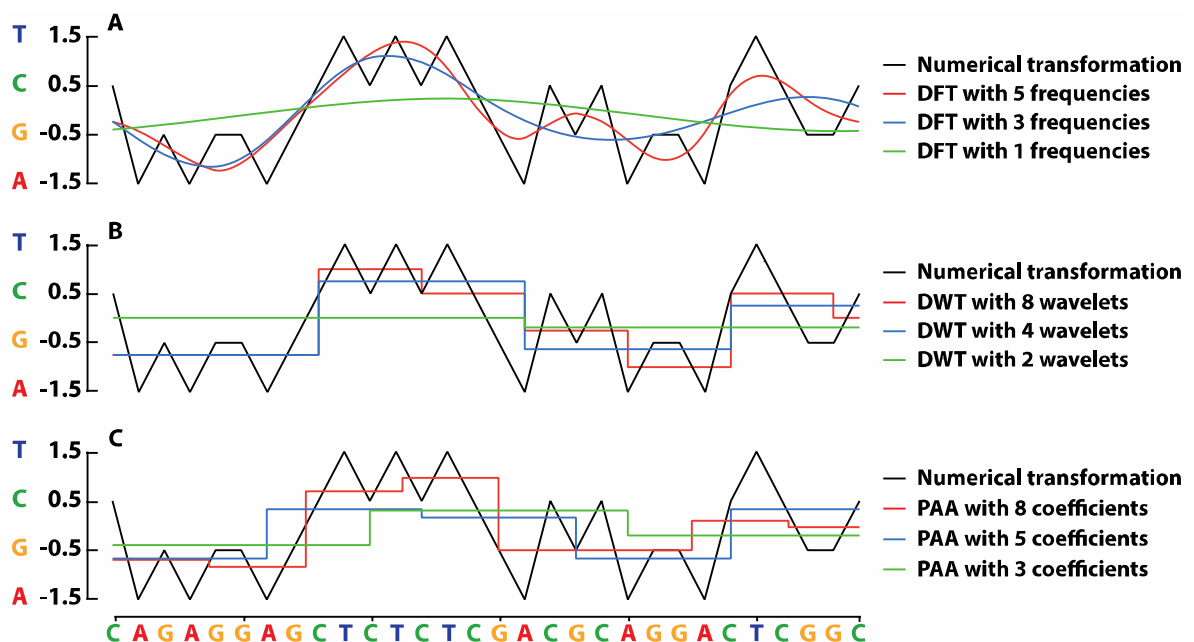
393 In conclusion, short nucleotide sequences may be effectively represented as numerical series,
 394 enabling the application of existing analytical methods from a variety of mathematical and
 395 engineering fields for the purposes of sequence alignment and assembly. By applying established
 396 signal decomposition methods, compressed representations of nucleotide sequences can be created,
 397 permitting reductions in the spatiotemporal complexity of their analysis, without necessarily
 398 compromising analytical accuracy.
 399

400 Authors' Contribution

401 AT designed and wrote the methods and software, and performed the data analysis with help
 402 from BC, MP, SK, and MC. BC and MC, generated the simulated data and BC, MP and MC help on
 403 the data evaluation. AT and DLR conceived the study. AT and BC wrote the manuscript with
 404 comments from MP, MC and DLR. All authors read and approved the final manuscript.
 405

406 Funding

407 This work has been supported by the Wellcome Trust [097820/Z/11/B]; the BBSRC
 408 [BB/H012419/1 and BB/M001121/1 and BC by a BBSRC DTP studentship to DLR]; and the
 409 VIROGENESIS project which receives funding from the European Union's Horizon 2020 research
 410 and innovation programme under grant agreement No 634650.
 411

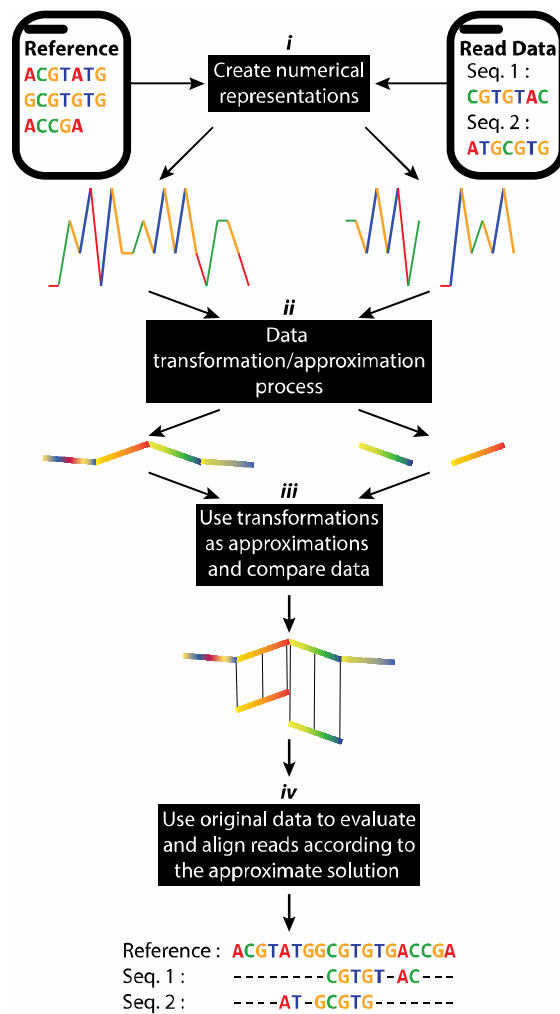


412

413 **Figure 1.** A numerically represented DNA sequence transformed at various levels of spatial
 414 resolution using the discrete Fourier transform (DFT) of the whole sequence (A), the Haar discrete
 415 wavelet transform (DWT) (B), and piecewise aggregate approximation (PAA) (C). A 30 nucleotide
 416 sequence (x-axis) is represented as a numerical sequence (black lines) using the real number
 417 representation method (y-axis where T=1.5, C=0.5, G=-0.5 and A=-1.5) for DFT approximations of the
 418 sequence with 5 (red), 3 (blue) and 1 (green) Fourier frequencies (A); DWT approximations of the

419
420

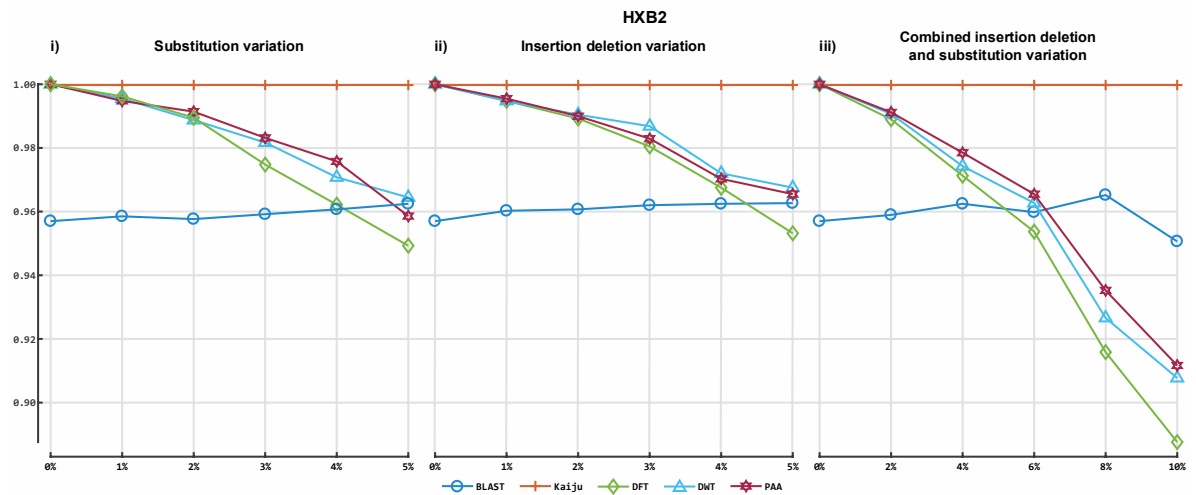
same sequence with 8 level wavelets (red), 4 level wavelets (blue), and 2 level wavelets (green) (B); and PAA approximations of the same sequence with 8 (red), 5 (blue) and 3 (green) coefficients (C).



421

422
423
424
425
426
427
428
429

Figure 2. Overview of our proposed methodology using time series transformation/approximation methods: (i) Creation of numerical representations of input sequences. (ii) Application of an appropriate signal decomposition method to transform sequences into their feature space. (iii) Use of approximated transformations to perform rapid data analysis in lower dimensional space. (iv) Validation of inferences against original, full-resolution input sequences. In the case of reference-based alignment, and taxonomic classification, approximated read transformations were compared with a reference sequence. In our *de novo* implementation, pairwise comparisons were performed between all of the approximated read transformations.



430

431

432

433

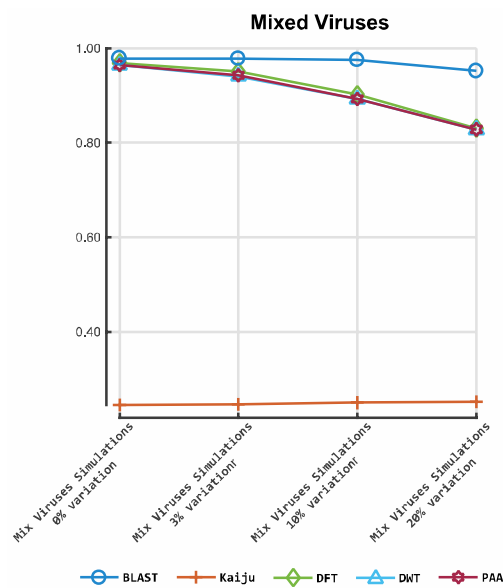
434

435

436

437

Figure 3. Accuracy of our prototype classification implementation and two established tools on HIV-HXB2 simulated datasets. All plots illustrate the F-measures obtained on the 16 different HIV datasets. The Y axis indicates the F-measure score, and the X axis depicts the reads data files. Plot 3-i depicts the F-measures obtained for each classifier on the simulations with 0% to 5% of substitution variation rate. Plot 6-ii illustrates the F-measures obtained for each classifier on the simulations with 0% to 5% uniform insertion/deletion variation and plot 3-iii illustrates the F-measures obtained for each tool on simulations of uniform 0% to 10% insertion/deletion and substitution variation.



438

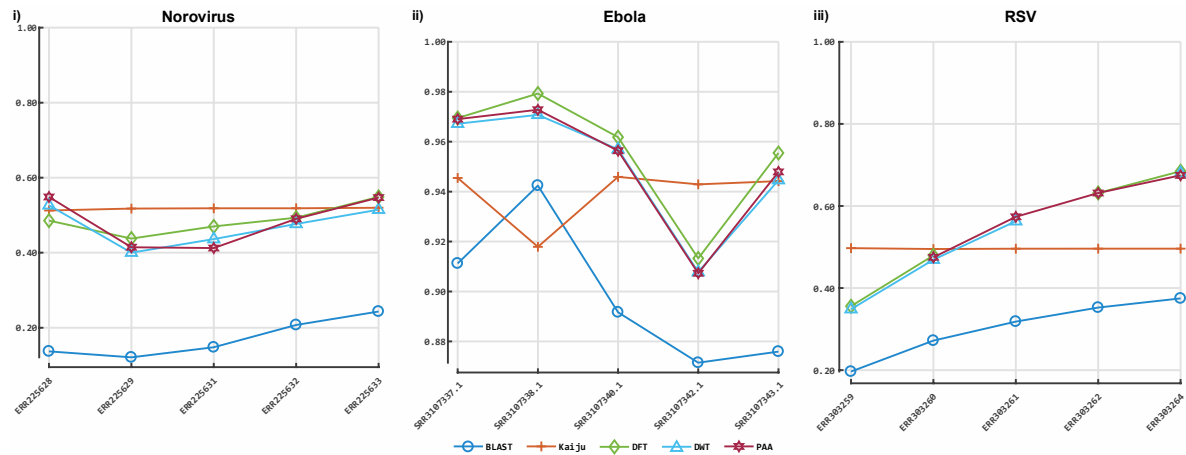
439

440

441

442

Figure 4. Accuracy of our prototype classification implementation and two established tools on mixed viruses simulated datasets. The Y axis indicates the F-measure score, and the X axis depicts the reads data files. The plot depicts the F-measures obtained for each classifier on the mixed virus simulations.



443

444

445

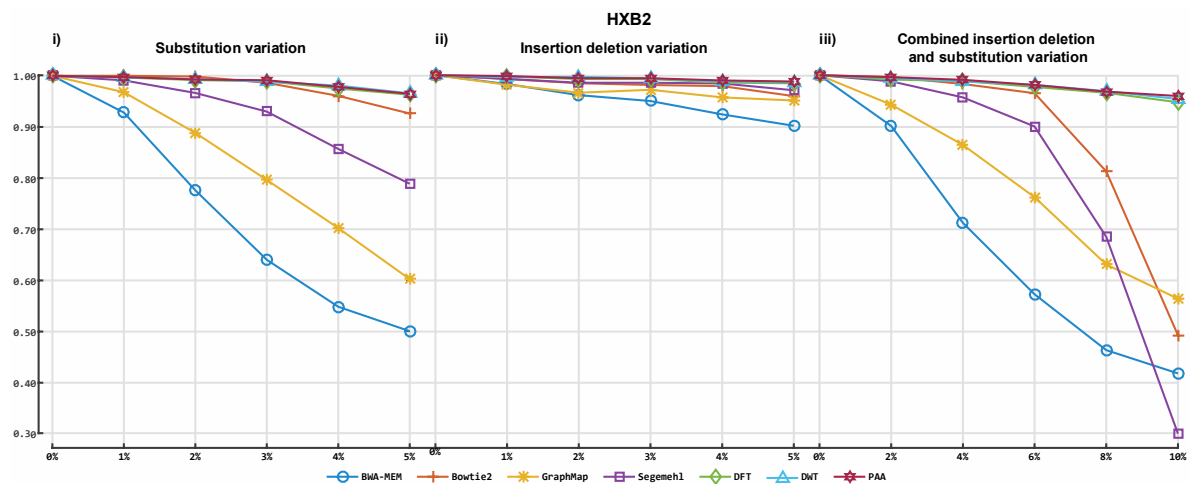
446

447

448

449

Figure 5. Accuracy of our prototype classification implementation and two established tools on real sequences. The Y axis indicates the F-measure score, and the X axis depicts the reads data files. The Y axis indicates the F-measures obtained for each classifier on the Norovirus sequences data. Plot 5-ii illustrates the F-measures obtained for each classifier on the Ebola sequence data. Plot 5-iii illustrates the F-measures obtained for each tool on Respiratory syncytial virus (RSV) sequence data.



450

451

452

453

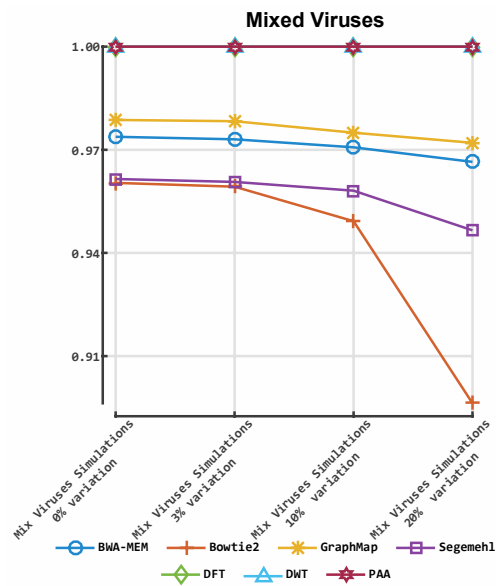
454

455

456

457

Figure 6. Accuracy of our prototype reference alignment implementation and four established tools on HIV-HXB2 simulated datasets. Figure 6 illustrates the F-measures obtained on the 16 different HIV datasets. Plot 6-i depicts the F-measures obtained for each aligner on the simulations with 0% to 5% of substitution variation rate. Plot 6-ii illustrates the F-measures obtained for each aligner on the simulations with 0% to 5% uniform insertion/deletion variation and plot 6-iii illustrates the F-measures obtained for each tool on simulations of uniform 0% to 10% insertion/deletion and substitution variation.



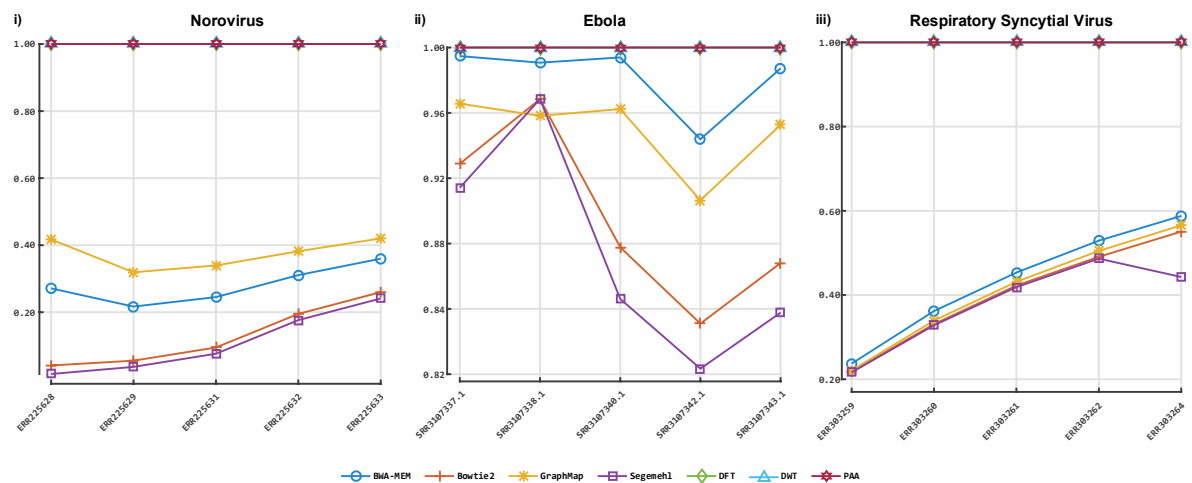
458

459

460

461

Figure 7. Accuracy of our prototype aligner implementation and four established tools on mixed viruses simulated datasets. The Y axis indicates the F-measure score, and the X axis depicts the reads data files. The plot depicts the F-measures obtained for each aligner on the mixed virus simulations



462

463

464

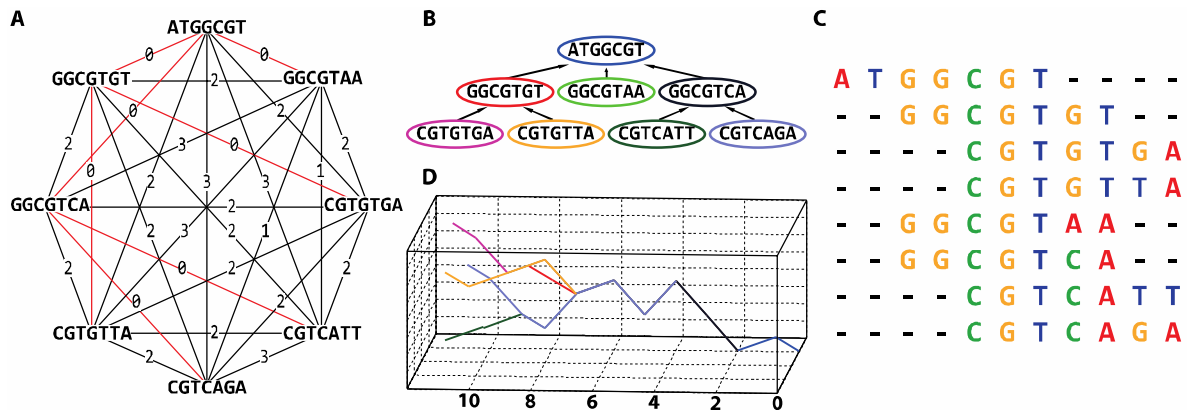
465

466

467

468

Figure 8. Accuracy of our prototype aligner implementation and four established tools on real sequences datasets. The Y axis indicates the F-measure score, and the X axis depicts the reads data files. The Y axis indicates the F-measure score, and the X axis depicts the reads data files. Plot 8-i depicts the F-measures obtained for each aligner on the Norovirus sequences data. Plot 8-ii illustrates the F-measures obtained for each aligner on the Ebola sequences data. Plot 8-iii illustrates the F-measures obtained for each tool on the Respiratory syncytial virus (RSV) sequences data.



469

470

471

472

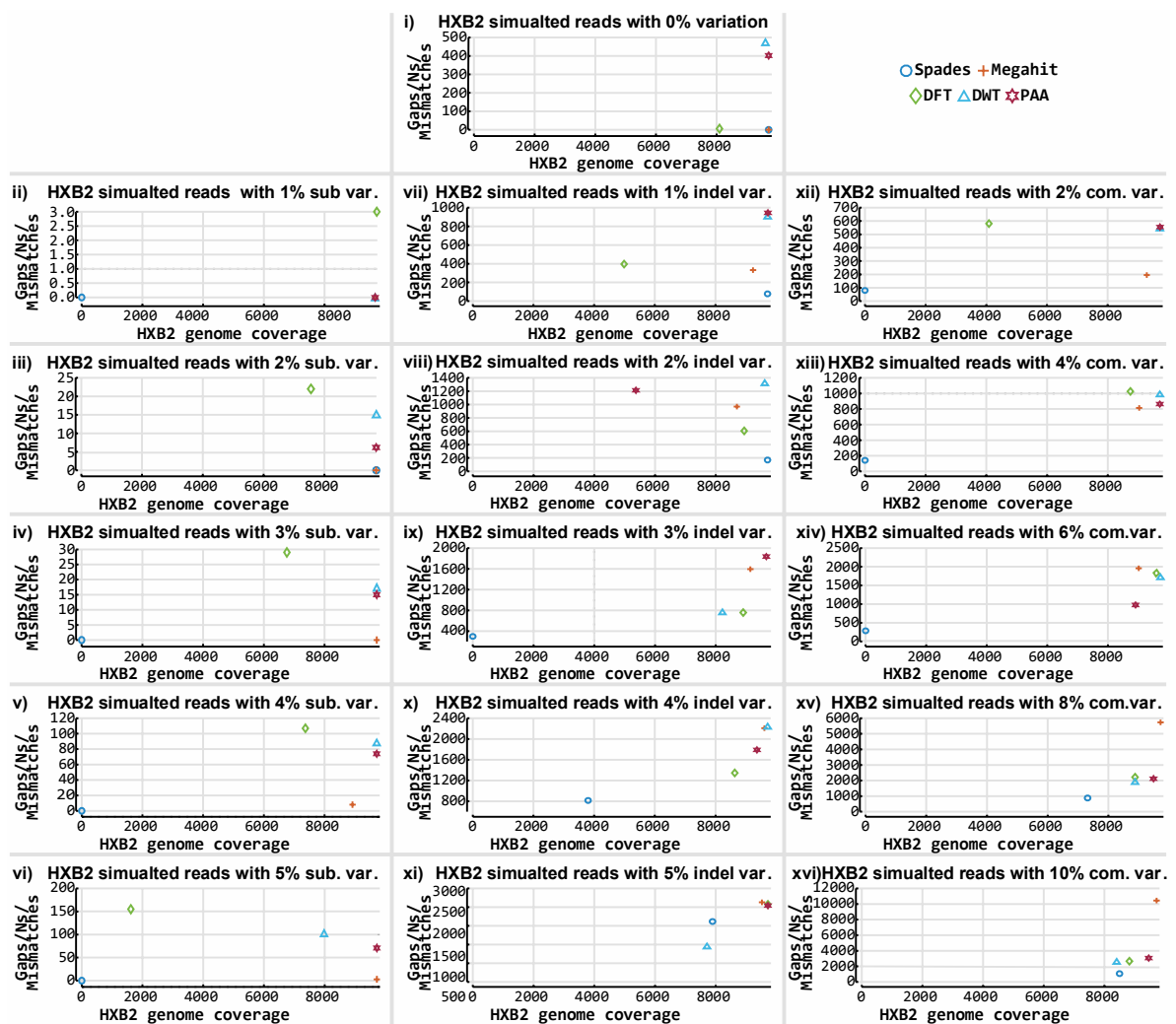
473

474

475

476

Figure 9. A *de novo* assembly methodology for numerically represented nucleotide reads. All-against-all sequence comparison (A) enables construction of a read graph with weighted edges. The weight assigned to each edge is the smallest pairwise distance between every possible k -mer representation of the two reads. (B) The shortest path in the graph is identified with a breadth-first search algorithm (red coloured edges) thereby (C) enabling read alignment. A DNA walk representation of the overlapped reads (D) may subsequently be used as a three-dimensional graphical portrayal of the reads, illustrating alignment characteristics.



477

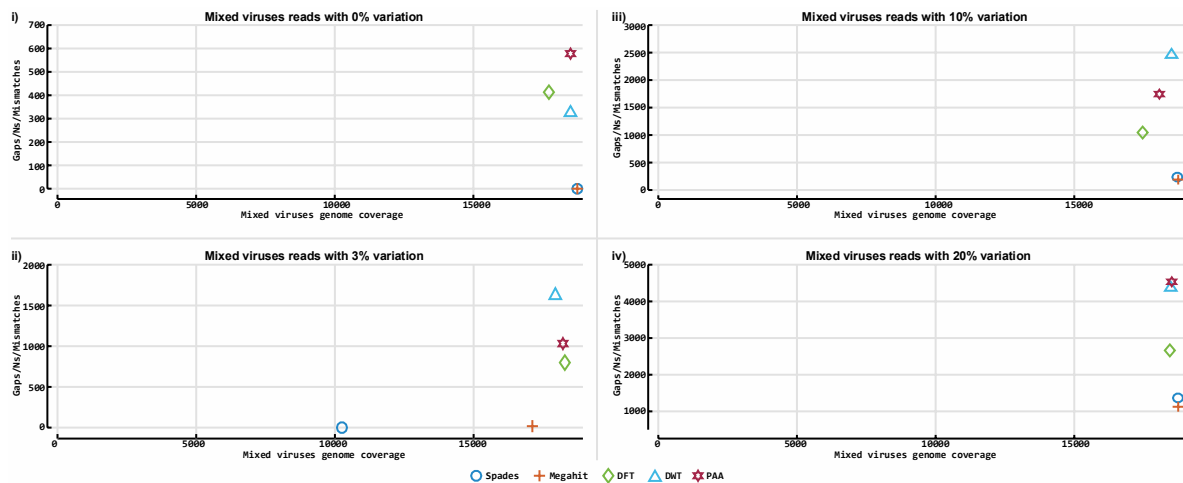
478

479

480

Figure 10. Accuracy of our prototype *de novo* assembly implementation and two established tools on HIV-HXB2 simulated datasets. The contigs obtained for each assembler were evaluated against the reference genome used to generate the simulated data. BLASTn was used to align all contigs to a

481 HIV-HXB2 reference genome and determine genome coverage. The Y axis indicates the number of
 482 gaps and mismatches that exist in the contigs obtained for each tool, and the X axis depicts the length
 483 of the genome the reported contigs cover. The contigs obtained from the assembly of the HIV-HXB2
 484 simulated short read data were evaluated against the K03455 reference genome. Plot 10-i illustrates
 485 results obtained from all assemblers on variation free data. Plots 10-ii to 10-vi illustrate results
 486 obtained from all assemblers on data with different levels of substitution variation. Plots 10-vii to
 487 10-xi illustrate results obtained from all assemblers on data with different levels of insertion/deletion
 488 variation. Plots 10-xiii to 10-xvi illustrate results obtained from all assemblers on data with different
 489 levels of combined insertion/deletion and substitution variation.



490

491 **Figure 11.** Accuracy of our prototype *de novo* assembly implementation and two established tools on
 492 mixed viruses simulated datasets. The contigs obtained for each assembler were evaluated against
 493 the reference genome that was used to generate the simulated data. BLASTn was used to align all
 494 contigs to a HIV-HXB2 reference genome, and determine how much of the particular genome they do
 495 cover. The Y axis indicates the number of gaps and mismatches that exist in the contigs obtained for
 496 each tool, and the X axis depicts the length of the genome the reported contigs cover. The contigs
 497 obtained from the mixed virus simulated dataset were evaluated against the, KM198529, KM198528,
 498 KM198511, KM198500, KM198486, KU296608, KU296553, KU296549, KU296528, KU296416,
 499 KP317952, KP317946, KP317934, KP317923, and KP317922 references genomes. Plots 11-i to 11-iv
 500 illustrate results obtained from all assemblers on data with 0%, 3%, 10% and 20% variation levels
 501 accordingly.

502

503 **Table 1** Numerical nucleotide sequences representation methods.

Integer number	$A = 1, C = -1, G = 2, T = -2, N = 0$
Real number	$A = -1.5, C = 0.5, G = -0.5, T = 1.5, N = 0.0$
EIIP	$A = 0.1260, C = 0.1340, G = 0.0806, T = 0.1335, N = \emptyset$
Atomic	$A = 70, C = 58, G = 78, T = 66, N = \emptyset$
Pair	$A \text{ or } T = 1, C \text{ or } G = -1, N = 0$
Complex number	$A = 1 + li, C = -1 + li, G = -i - li, T = 1 - li, N = 0 + 0i$
DNA Walk	$A = [1, 0], C = [0, 1], G = [0, -1], T = [-1, 0], N = [0, 0]$
Tetrahedron	$A = [0, 0, 1], C = \left[-\frac{\sqrt{2}}{3}, -\frac{\sqrt{6}}{3}, \frac{1}{3} \right], G = \left[-\frac{\sqrt{2}}{3}, -\frac{\sqrt{6}}{3}, -\frac{1}{3} \right],$ $T = \left[2 \times \frac{\sqrt{2}}{3}, 0, -\frac{1}{3} \right], N = [0, 0, 0]$
Voss indicator	$A = [0, 0, 1, 0], C = [1, 0, 0, 0], G = [0, 1, 0, 0],$ $T = [0, 0, 0, 1], N = [0, 0, 0, 0]$

504

505

506
507
508**Table 2.** Simulated read data. Each row contains details for each simulated dataset (i.e. virus family, virus, GenBank id, variation type, variation level, number of reads, and simulator used to generate data).

Family	Virus	GenBank genome ID	Variation Type (%)			Reads	Simulator
			Ins	Del	Sub		
HIV	HXB2	K03455	0.0	0.0	0.0	2133	Curesim
HIV	HXB2	K03455	0.0	0.0	1.0	2133	Curesim
HIV	HXB2	K03455	0.0	0.0	2.0	2133	Curesim
HIV	HXB2	K03455	0.0	0.0	3.0	2133	Curesim
HIV	HXB2	K03455	0.0	0.0	4.0	2133	Curesim
HIV	HXB2	K03455	0.0	0.0	5.0	2133	Curesim
HIV	HXB2	K03455	0.5	0.5	0.0	2133	Curesim
HIV	HXB2	K03455	1.0	1.0	0.0	2133	Curesim
HIV	HXB2	K03455	1.5	1.5	0.0	2133	Curesim
HIV	HXB2	K03455	2.0	2.0	0.0	2133	Curesim
HIV	HXB2	K03455	2.5	2.5	0.0	2133	Curesim
HIV	HXB2	K03455	0.5	0.5	1.0	2133	Curesim
HIV	HXB2	K03455	1.0	1.0	2.0	2133	Curesim
HIV	HXB2	K03455	1.5	1.5	3.0	2133	Curesim
HIV	HXB2	K03455	2.0	2.0	4.0	2133	Curesim
HIV	HXB2	K03455	2.5	2.5	5.0	2133	Curesim
Mixed Viruses: Caliciviridae, Filoviridae, Pneumoviridae	Norovirus, Ebola virus, RSV	KM198529, KM198528, KM198511, KM198500, KM198486, KU296608, KU296553, KU296549, KU296528, KU296416, KP317952, KP317946, KP317934, KP317923, KP317922	0.0	0.0	0.0	200000	WGSIM
Mixed Viruses: Caliciviridae, Filoviridae, Pneumoviridae	Norovirus, Ebola virus, RSV	KM198529, KM198528, KM198511, KM198500, KM198486, KU296608, KU296553, KU296549, KU296528, KU296416, KP317952, KP317946, KP317934, KP317923, KP317922	1.0	1.0	1.0	200000	WGSIM
Mixed Viruses, Caliciviridae, Filoviridae, Pneumoviridae	Norovirus, Ebola virus, RSV	KM198529, KM198528, KM198511, KM198500, KM198486, KU296608, KU296553, KU296549, KU296528, KU296416, KP317952, KP317946, KP317934, KP317923, KP317922	3.33	3.33	3.33	100000	WGSIM
Mixed Viruses, Caliciviridae, Filoviridae, Pneumoviridae	Norovirus, Ebola virus, RSV	KM198529, KM198528, KM198511, KM198500, KM198486, KU296608, KU296553, KU296549, KU296528, KU296416, KP317952, KP317946, KP317934, KP317923, KP317922	6.66	6.66	6.66	200000	WGSIM

509

510

511

512

Table 3. Real short reads data. Rows contain information for each real reads' dataset (i.e. virus family, virus, genome strain GenBank id, SRA project ID, number of reads, and technology used to sequence data).

Family	Virus	Amplicon/random primer	GenBank genome ID	ENA/SRA_ID	Reads	Sequencing Technology
Caliciviridae	Norovirus	Amplicon	KM198486	ERR225628	2126502	Illumina MiSeq
Caliciviridae	Norovirus	Amplicon	KM198500	ERR225629	3037674	Illumina MiSeq
Caliciviridae	Norovirus	Amplicon	KM198511	ERR225631	3285078	Illumina MiSeq
Caliciviridae	Norovirus	Amplicon	KM198528	ERR225632	4361884	Illumina MiSeq
Caliciviridae	Norovirus	Amplicon	KM198529	ERR225633	5187234	Illumina MiSeq
Filoviridae	Ebola virus	Amplicon	KU296608	SRR3107337	522968	Ion Torrent PGM
Filoviridae	Ebola virus	Amplicon	KU296549	SRR3107338	771031	Ion Torrent PGM
Filoviridae	Ebola virus	Amplicon	KU296416	SRR3107340	186657	Ion Torrent PGM
Filoviridae	Ebola virus	Amplicon	KU296553	SRR3107342	478346	Ion Torrent PGM
Filoviridae	Ebola virus	Amplicon	KU296528	SRR3107343	42410	Ion Torrent PGM
Pneumoviridae	RSV	Amplicon	KP317934	ERR303259	7275032	Illumina MiSeq
Pneumoviridae	RSV	Amplicon	KP317922	ERR303260	9278070	Illumina MiSeq
Pneumoviridae	RSV	Amplicon	KP317946	ERR303261	11111114	Illumina MiSeq
Pneumoviridae	RSV	Amplicon	KP317923	ERR303262	13293226	Illumina MiSeq
Pneumoviridae	RSV	Amplicon	KP317952	ERR303263	15237848	Illumina MiSeq

513

514

515 **Table 4.** Reference genomes used during classification and reference based alignment

Family	Virus	GenBank id:	Length (nt)
Retroviridae	Human immunodeficiency virus 1 (HXB2)	K03455	9179
Caliciviridae	Norovirus	KM198509.1	7425
Filoviridae	Zaire ebolavirus	KM034562.1	18957
Pneumoviridae	Human orthopneumovirus (Respiratory Syncytial Virus)	KP317934.1	15233

516
517
518**Table 5.** Pseudocode for the alignment procedure

- 1) Represent short reads and reference genome as numerical sequences.
 - 2) Select k -mer length.
 - 3) Create transformations of each reference sequence k -mer, build VP-tree, and create transformations of the initial k -mer of each short reads.
 - 4) Identify candidate alignments using data transformations.
- for each read i
- candidate_alignments[i] = $VPtree.k$ -NNSearch(query i)
- end
- 5) Align approximate results with original data using the Smith-Waterman (SW) algorithm:
- for each read i
- best_score = null
- best_aln = []
- for each k neighbour in candidate_alignments[i]
- if SW_score(k neighbour, read i)
- best_score = SW_score(k neighbour, read i)
- best_aln = SW_aln(k neighbour, read i)
- end
- end
- end
- 6) Output alignment in Sequence Alignment/Map (SAM) format.

519

520

521 **References**

- 522 1. Margulies, M.; Egholm, M.; Altman, W.E.; Attiya, S.; Bader, J.S.; Bemben, L.A.; Berka, J.; Braverman,
523 M.S.; Chen, Y.-J.; Chen, Z. Genome sequencing in microfabricated high-density picolitre reactors.
524 *Nature* **2005**, *437*, 376-380.
- 525 2. Bentley, D.R.; Balasubramanian, S.; Swerdlow, H.P.; Smith, G.P.; Milton, J.; Brown, C.G.; Hall, K.P.;
526 Evers, D.J.; Barnes, C.L.; Bignell, H.R. Accurate whole human genome sequencing using reversible
527 terminator chemistry. *Nature* **2008**, *456*, 53-59.

- 528 3. Rothberg, J.M.; Hinz, W.; Rearick, T.M.; Schultz, J.; Mileski, W.; Davey, M.; Leamon, J.H.; Johnson, K.;
529 Milgrew, M.J.; Edwards, M. An integrated semiconductor device enabling non-optical genome
530 sequencing. *Nature* **2011**, *475*, 348-352.
- 531 4. Eid, J.; Fehr, A.; Gray, J.; Luong, K.; Lyle, J.; Otto, G.; Peluso, P.; Rank, D.; Baybayan, P.; Bettman, B.
532 Real-time DNA sequencing from single polymerase molecules. *Science* **2009**, *323*, 133-138.
- 533 5. Salipante, S.J.; Roach, D.J.; Kitzman, J.O.; Snyder, M.W.; Stackhouse, B.; Butler-Wu, S.M.; Lee, C.;
534 Cookson, B.T.; Shendure, J. Large-scale genomic sequencing of extraintestinal pathogenic *Escherichia*
535 *coli* strains. *Genome research* **2014**, gr. 180190.180114.
- 536 6. Rose, R.; Constantinides, B.; Tapinos, A.; Robertson, D.L.; Prosperi, M. Challenges in the analysis of
537 viral metagenomes. *Virus Evolution* **2016**, *2*, vew022.
- 538 7. Li, H.; Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform.
539 *Bioinformatics* **2009**, *25*, 1754-1760.
- 540 8. Shrestha, A.M.S.; Frith, M.C.; Horton, P. A bioinformatician’s guide to the forefront of suffix array
541 construction algorithms. *Briefings in bioinformatics* **2014**, *15*, 138-154.
- 542 9. Myers, E.W. Toward simplifying and accurately formulating fragment assembly. *Journal of*
543 *Computational Biology* **1995**, *2*, 275-290.
- 544 10. Kececioglu, J.D.; Myers, E.W. Combinatorial algorithms for DNA sequence assembly. *Algorithmica*
545 **1995**, *13*, 7-51.
- 546 11. Earl, D.; Bradnam, K.; John, J.S.; Darling, A.; Lin, D.; Fass, J.; Yu, H.O.K.; Buffalo, V.; Zerbino, D.R.;
547 Diekhans, M. Assemblathon 1: A competitive assessment of de novo short read assembly methods.
548 *Genome research* **2011**, *21*, 2224-2241.
- 549 12. Iqbal, Z.; Caccamo, M.; Turner, I.; Flicek, P.; McVean, G. De novo assembly and genotyping of variants
550 using colored de Bruijn graphs. *Nature genetics* **2012**, *44*, 226-232.
- 551 13. Pevzner, P.A.; Tang, H.; Waterman, M.S. An Eulerian path approach to DNA fragment assembly.
552 *Proceedings of the National Academy of Sciences* **2001**, *98*, 9748-9753.
- 553 14. Bradnam, K.R.; Fass, J.N.; Alexandrov, A.; Baranay, P.; Bechner, M.; Birol, I.; Boisvert, S.; Chapman,
554 J.A.; Chapuis, G.; Chikhi, R. Assemblathon 2: evaluating de novo methods of genome assembly in
555 three vertebrate species. *GigaScience* **2013**, *2*, 1-31.
- 556 15. Archer, J.; Rambaut, A.; Taillon, B.E.; Harrigan, P.R.; Lewis, M.; Robertson, D.L. The evolutionary
557 analysis of emerging low frequency HIV-1 CXCR4 using variants through time—an ultra-deep
558 approach. *PLoS computational biology* **2010**, *6*, e1001022.
- 559 16. Clement, N.L.; Thompson, L.P.; Miranker, D.P. ADaM: augmenting existing approximate fast
560 matching algorithms with efficient and exact range queries. *BMC bioinformatics* **2014**, *15*, S1.
- 561 17. Agrawal, R.; Faloutsos, C.; Swami, A. *Efficient similarity search in sequence databases*; Springer: 1993.
- 562 18. Chan, K.-P.; Fu, A.-C. Efficient time series matching by wavelets. In *Proceedings of Data Engineering*,
563 1999. Proceedings., 15th International Conference on; pp. 126-133.
- 564 19. Woodward, A.M.; Rowland, J.J.; Kell, D.B. Fast automatic registration of images using the phase of a
565 complex wavelet transform: application to proteome gels. *Analyst* **2004**, *129*, 542-552.
- 566 20. Geurts, P. Pattern extraction for time series classification. In *Principles of Data Mining and Knowledge*
567 *Discovery*, Springer: 2001; pp. 115-127.
- 568 21. Keogh, E.; Chakrabarti, K.; Pazzani, M.; Mehrotra, S. Locally adaptive dimensionality reduction for
569 indexing large time series databases. *ACM SIGMOD Record* **2001**, *30*, 151-162.

- 570 22. Shumway, R.H.; Stoffer, D.S.; Stoffer, D.S. *Time series analysis and its applications*; Springer New York:
571 2000; Vol. 3.
- 572 23. Silverman, B.; Linsker, R. A measure of DNA periodicity. *Journal of theoretical biology* **1986**, *118*, 295-300.
- 573 24. Cheever, E.; Searls, D.; Karunaratne, W.; Overton, G. Using signal processing techniques for DNA
574 sequence comparison. In Proceedings of Bioengineering Conference, 1989., Proceedings of the 1989
575 Fifteenth Annual Northeast; pp. 173-174.
- 576 25. Katoh, K.; Misawa, K.; Kuma, K.i.; Miyata, T. MAFFT: a novel method for rapid multiple sequence
577 alignment based on fast Fourier transform. *Nucleic acids research* **2002**, *30*, 3059-3066.
- 578 26. Kwan, H.K.; Arniker, S.B. Numerical representation of DNA sequences. In Proceedings of
579 Electro/Information Technology, 2009. eit'09. IEEE International Conference on; pp. 307-310.
- 580 27. Yi, B.-K.; Faloutsos, C. Fast time sequence indexing for arbitrary Lp norms.
- 581 28. Keogh, E.; Ratanamahatana, C.A. Exact indexing of dynamic time warping. *Knowledge and information
582 systems* **2005**, *7*, 358-386.
- 583 29. Vlachos, M.; Kollios, G.; Gunopulos, D. Discovering similar multidimensional trajectories. In
584 Proceedings of Data Engineering, 2002. Proceedings. 18th International Conference on; pp. 673-684.
- 585 30. Kotsakos, D.; Trajcevski, G.; Gunopulos, D.; Aggarwal, C.C. Time-Series Data Clustering. 2013.
- 586 31. Chávez, E.; Navarro, G.; Baeza-Yates, R.; Marroquín, J.L. Searching in metric spaces. *ACM computing
587 surveys (CSUR)* **2001**, *33*, 273-321.
- 588 32. Beckmann, N.; Kriegel, H.-P.; Schneider, R.; Seeger, B. *The R*-tree: an efficient and robust access method for
589 points and rectangles*; ACM: 1990; Vol. 19.
- 590 33. Lin, R.; King-Ip, A.; Shim, H.S.S.K. Fast similarity search in the presence of noise, scaling, and
591 translation in time-series databases. In Proceedings of Proceeding of the 21th International Conference
592 on Very Large Data Bases; pp. 490-501.
- 593 34. Bingham, S.; Kot, M. Multidimensional trees, range searching, and a correlation dimension algorithm
594 of reduced complexity. *Physics Letters A* **1989**, *140*, 327-330.
- 595 35. Bellman, R.; Bellman, R.E.; Bellman, R.E.; Bellman, R.E. *Adaptive control processes: a guided tour*;
596 Princeton university press Princeton: 1961; Vol. 4.
- 597 36. Verleysen, M.; François, D. The curse of dimensionality in data mining and time series prediction. In
598 *Computational Intelligence and Bioinspired Systems*, Springer: 2005; pp. 758-770.
- 599 37. Yianilos, P.N. Data structures and algorithms for nearest neighbor search in general metric spaces. In
600 Proceedings of SODA; pp. 311-321.
- 601 38. Bozkaya, T.; Ozsoyoglu, M. Indexing large metric spaces for similarity search queries. *ACM
602 Transactions on Database Systems (TODS)* **1999**, *24*, 361-404.
- 603 39. Uhlmann, J.K. Satisfying general proximity/similarity queries with metric trees. *Information processing
604 letters* **1991**, *40*, 175-179.
- 605 40. Nair, A.S.; Sreenadhan, S.P. A coding measure scheme employing electron-ion interaction
606 pseudopotential (EIIP). *Bioinformatics* **2006**, *1*, 197.
- 607 41. Holden, T.; Subramaniam, R.; Sullivan, R.; Cheung, E.; Schneider, C.; Tremberger, G.; Flamholz, A.;
608 Lieberman, D.; Cheung, T. ATCG nucleotide fluctuation of *Deinococcus radiodurans* radiation genes.
609 In Proceedings of Instruments, Methods, and Missions for Astrobiology X; p. 669417.
- 610 42. Voss, R.F. Evolution of long-range fractal correlations and 1/f noise in DNA base sequences. *Physical
611 review letters* **1992**, *68*, 3805.

- 612 43. Faloutsos, C.; Ranganathan, M.; Manolopoulos, Y. *Fast subsequence matching in time-series databases*;
613 ACM: 1994; Vol. 23.
- 614 44. Mitsa, T. *Temporal data mining*; CRC Press: 2010.
- 615 45. Mörchen, F. Time series feature extraction for data mining using DWT and DFT. Univ.: 2003.
- 616 46. Jensen, A.; la Cour-Harbo, A. *Ripples in mathematics: the discrete wavelet transform*; Springer: 2001.
- 617 47. Wu, Y.-L.; Agrawal, D.; El Abbadi, A. A comparison of DFT and DWT based similarity search in
618 time-series databases. In Proceedings of Proceedings of the ninth international conference on
619 Information and knowledge management; pp. 488-495.
- 620 48. Caboche, S.; Audebert, C.; Lemoine, Y.; Hot, D. Comparison of mapping algorithms used in
621 high-throughput sequencing: application to Ion Torrent data. *BMC genomics* **2014**, *15*, 264.
- 622 49. Cotten, M.; Petrova, V.; Phan, M.V.; Rabaa, M.A.; Watson, S.J.; Ong, S.H.; Kellam, P.; Baker, S. Deep
623 sequencing of norovirus genomes defines evolutionary patterns in an urban tropical setting. *Journal of*
624 *virology* **2014**, *88*, 11056-11069.
- 625 50. Phan, M.V.; Anh, P.H.; Cuong, N.V.; Munnink, B.B.O.; van der Hoek, L.; My, P.T.; Tri, T.N.; Bryant,
626 J.E.; Baker, S.; Thwaites, G. Unbiased whole-genome deep sequencing of human and porcine stool
627 samples reveals circulation of multiple groups of rotaviruses and a putative zoonotic infection. *Virus*
628 *evolution* **2016**, *2*.
- 629 51. Kiyuka, P.K.; Agoti, C.N.; Munywoki, P.K.; Njeru, R.; Bett, A.; Otieno, J.R.; Otieno, G.P.; Kamau, E.;
630 Clark, T.G.; van der Hoek, L. Human coronavirus NL63 Molecular epidemiology and evolutionary
631 patterns in rural coastal Kenya. *The Journal of infectious diseases* **2018**, *217*, 1728-1739.
- 632 52. Arias, A.; Watson, S.J.; Asogun, D.; Tobin, E.A.; Lu, J.; Phan, M.V.; Jah, U.; Wadoun, R.E.G.; Meredith,
633 L.; Thorne, L. Rapid outbreak sequencing of Ebola virus in Sierra Leone identifies transmission chains
634 linked to sporadic cases. *Virus Evolution* **2016**, *2*.
- 635 53. Agoti, C.N.; Otieno, J.R.; Munywoki, P.K.; Mwihuri, A.G.; Cane, P.A.; Nokes, D.J.; Kellam, P.; Cotten,
636 M. Local evolutionary patterns of human respiratory syncytial virus derived from whole-genome
637 sequencing. *Journal of virology* **2015**, *89*, 3444-3454.
- 638 54. Altschul, S.F.; Gish, W.; Miller, W.; Myers, E.W.; Lipman, D.J. Basic local alignment search tool. *Journal*
639 *of molecular biology* **1990**, *215*, 403-410.
- 640 55. Menzel, P.; Ng, K.L.; Krogh, A. Fast and sensitive taxonomic classification for metagenomics with
641 Kaiju. *Nature communications* **2016**, *7*, 11257.
- 642 56. Langmead, B.; Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nature Methods* **2012**, *9*, 357,
643 doi:10.1038/nmeth.1923 <https://www.nature.com/articles/nmeth.1923#supplementary-information>.
- 644 57. Sović, I.; Šikić, M.; Wilm, A.; Fenlon, S.N.; Chen, S.; Nagarajan, N. Fast and sensitive mapping of
645 nanopore sequencing reads with GraphMap. *Nature communications* **2016**, *7*, 11307.
- 646 58. Otto, C.; Stadler, P.F.; Hoffmann, S. Lacking alignments? The next-generation sequencing mapper
647 segemehl revisited. *Bioinformatics* **2014**, *30*, 1837-1843.
- 648 59. Li, D.; Liu, C.-M.; Luo, R.; Sadakane, K.; Lam, T.-W. MEGAHIT: an ultra-fast single-node solution for
649 large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* **2015**, *31*,
650 1674-1676.
- 651 60. Anton, B.; Sergey, N.; Dmitry, A.; Alexey, A.; Mikhail, D.; Alexander, S.; Valery, M.; Sergey, I.; Son, P.;
652 Andrey, D. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing.
653 *Journal of Computational Biology* **2012**, *19*, 455.

- 654 61. Tapinos, A.; Mendes, P. A method for comparing multivariate time series with different dimensions.
655 *PloS one* **2013**, *8*, e54201.
- 656 62. Sheybani, E.O. An Algorithm for Real-Time Blind Image Quality Comparison and Assessment.
657 *International Journal of Electrical and Computer Engineering (IJECE)* **2011**, *2*, 120-129.
- 658 63. Hendriks, R.C.; Gerkmann, T.; Jensen, J. DFT-domain based single-microphone noise reduction for
659 speech enhancement: a survey of the state of the art. *Synthesis Lectures on Speech and Audio Processing*
660 **2013**, *9*, 1-80.
- 661 64. Shi, H.; Schmidt, B.; Liu, W.; Müller-Wittig, W. A Parallel Algorithm for Error Correction in
662 High-Throughput Short-Read Data on CUDA-Enabled Graphics Hardware. *Journal of Computational*
663 *Biology* **2010**, *17*, 603-615, doi:10.1089/cmb.2009.0062.
- 664 65. Zhang, Q.; Pell, J.; Canino-Koning, R.; Howe, A.C.; Brown, C.T. These Are Not the K-mers You Are
665 Looking For: Efficient Online K-mer Counting Using a Probabilistic Data Structure. *PLoS ONE* **2014**, *9*,
666 e101271, doi:10.1371/journal.pone.0101271.
- 667 66. Salikhov, K.; Sacomoto, G.; Kucherov, G. Using cascading Bloom filters to improve the memory usage
668 for de Bruijn graphs. In Proceedings of WABI; pp. 364-376.
- 669 67. Berlin, K.; Koren, S.; Chin, C.-S.; Drake, J.P.; Landolin, J.M.; Phillippy, A.M. Assembling large genomes
670 with single-molecule sequencing and locality-sensitive hashing. *Nature biotechnology* **2015**, *33*, 623-630.
- 671 68. Camera, A.; Palpanas, T.; Shieh, J.; Keogh, E. iSAX 2.0: Indexing and mining one billion time series. In
672 Proceedings of Data Mining (ICDM), 2010 IEEE 10th International Conference on; pp. 58-67.
673

674