

1 Skip-mers: increasing entropy and 2 sensitivity to detect conserved genic 3 regions with simple cyclic q-grams

4 **Bernardo J. Clavijo¹, Gonzalo Garcia Accinelli¹, Luis Yanes¹, Katie Barr¹,
5 and Jonathan Wright¹**

6 ¹Earlham Institute, Norwich Research Park, Norwich, NR4 7UZ, UK

7 Corresponding author:

8 Bernardo J. Clavijo

9 Email address: bernardo.clavijo@earlham.ac.uk

10 ABSTRACT

11 Bioinformatic analyses and tools make extensive use of k-mers (fixed contiguous strings of k nucleotides)
12 as an informational unit. K-mer analyses are both useful and fast, but are strongly affected by single-
13 nucleotide polymorphisms or sequencing errors, effectively hindering direct-analyses of whole regions
14 and decreasing their usability between evolutionary distant samples.

15 We introduce a concept of skip-mers, a cyclic pattern of used-and-skipped positions of k nucleotides
16 spanning a region of size $S \geq k$, and show how analyses are improved compared to using k-mers. The
17 entropy of skip-mers increases with the larger span, capturing information from more distant positions
18 and increasing the specificity, and uniqueness, of larger span skip-mers within a genome. In addition,
19 skip-mers constructed in cycles of 1 or 2 nucleotides in every 3 (or a multiple of 3) lead to increased
20 sensitivity in the coding regions of genes, by grouping together the more conserved nucleotides of the
21 protein-coding regions.

22 We implemented a set of tools to count and intersect skip-mers between different datasets. We used
23 these tools to show how skip-mers have advantages over k-mers in terms of entropy and increased
24 sensitivity to detect conserved coding sequence, allowing better identification of genic matches between
25 evolutionarily distant species. We also highlight potential applications to problems such as whole-genome
26 alignment and multi-genome evolutionary analyses.

27 **Software availability:** the skm-tools implementing the methods described in this manuscript are available
28 under MIT license at <http://github.com/bioinfologics/skm-tools/>

29 1 INTRODUCTION

30 Genomes are not random strings, but are the product of millions of years of evolution and selection
31 pressure which imparts unique characteristics to the sequence of nucleotides. These characteristics need
32 to be considered in order to better analyse genomic datasets. Here we exploit the increase in entropy
33 (mean amount of information) from positions that are further away in the genome (Chaisson et al., 2009),
34 and the uneven conservation of coding sequence due to synonymous mutations and the neutral model
35 (Kimura, 1977). The new concept of skip-mers extends the familiar concept of k-mers by taking these
36 two properties into account.

37 First, we harness the increased entropy of nucleotides that are further apart by introducing gaps.
38 This has previously been explored to predict regulatory sequences (Ghandi et al., 2014) and to classify
39 sequences taxonomically (Hahn et al., 2016) with q-grams. Instead of general q-gram patterns, we define
40 simple cycles of nucleotide skips which preserve more of the useful properties of k-mers. Second, we
41 take advantage of the increased conservation present in the first two nucleotides of every trinucleotide
42 codon by analysing the skip-mer content of genomes in cycles of three. The consideration of these two
43 concepts allowed us to design skip-mers that improve genic region matches for syntenic analyses.

1.1 From k-mers to skip-mers

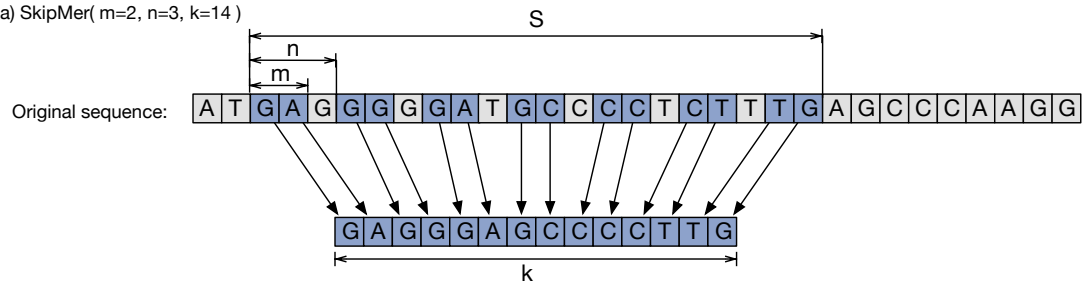
Bioinformatic analyses make extensive use of k-mers (contiguous strings of k nucleotides of sequence) as an informational unit, a concept popularised by short read assemblers (Zerbino and Birney, 2008). Analyses within the k-mer space benefit from a simple formulation of the sampling problem and direct in-hash comparisons (Mapleson et al., 2017). For some analyses, the contiguous nature of k-mers imposes limitations. A single base difference, due to real biological variation or a sequencing error, affects all k-mers crossing that position thus impeding direct analyses by identity. Also, given the strong interdependence of local sequence, contiguous sections capture less information about genome structure and are thus more affected by sequence repetition (Chaisson et al., 2009).

Q-grams are strings of nucleotides constructed from a pattern of used-and-skipped positions and have been applied to the sequence matching problem (Burkhardt and Kärkkäinen, 2003). The increased entropy due to a larger span and the higher tolerance to single base differences make q-grams a better tool than k-mers for many bioinformatics tasks. However, q-gram analyses are complicated by the inherent flexibility of the concept and the loss of some useful properties of k-mers like reverse complementability.

We define skip-mers as a cyclic pattern of used-and-skipped positions which achieves increased entropy and tolerance to nucleotide substitution differences by following some simple rules (see Figure 1 and the next section). Skip-mers preserve many of the elegant properties of k-mers such as reverse complementability and existence of a canonical representation. Also, using cycles of three greatly increases the power of direct intersection between the genomes of different organisms by grouping together the more conserved nucleotides of the protein-coding regions.

1.2 Skip-mer definition

a) SkipMer($m=2, n=3, k=14$)



b) SkipMer($m=2, n=4, k=14$)

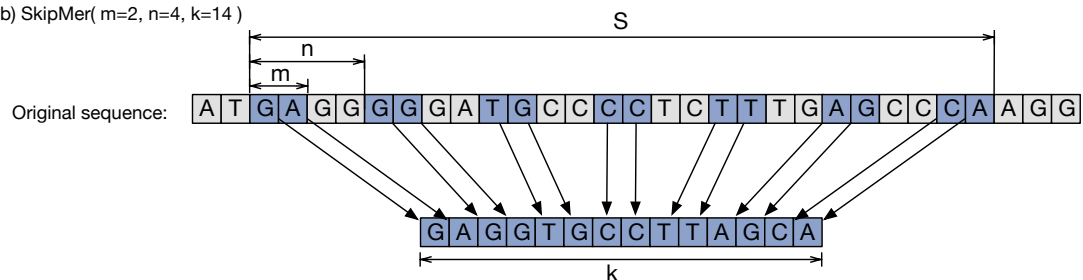


Figure 1. Different $SkipMer(m, n, k)$ cycles defined over the same sequence region, resulting in different combinations of bases. The shape of the underlying cyclic q-gram is defined by the variables m (used bases per cycle), n (cycle length), and k (total number of bases).

A skip-mer is a simple cyclic q-gram that includes m out of every n bases until a total of k bases is reached. Its shape is defined by a function $SkipMer(m, n, k)$, as shown in Figure 1. To maintain cyclic properties and the existence of the reverse-complement as a skip-mer defined by the same function, k must be a multiple of m . This also enables the existence of a canonical representation for each skip-mer, defined as the lexicographically smaller of the forward and reverse-complement representations.

Defining m , n and k fixes a value for S , the total span of the skip-mer, given by:

$$S = n \times \left(\frac{k}{m} - 1 \right) + m$$

71 It is important to note that k-mers are a sub-class of skip-mers. A skip-mer with $m = n$ will use all
72 contiguous k nucleotides, which makes it a k-mer. Throughout this manuscript we often use $m = 1 \wedge n = 1$,
73 or the shorter form notations $1-1$ or $1-1-k$ to refer to k-mers.

74 2 MATERIALS & METHODS

75 2.1 Genome sequences and annotations

76 To evaluate the properties of skip-mers in a genomic context we used publically available genome
77 assemblies.

78 Hexaploid bread wheat, *Triticum aestivum* (Clavijo et al., 2017), is a highly repetitive and complex
79 genome, and we used it to investigate the effect of the increased entropy when using larger skip-mer
80 cycles. *Oryza sativa* (Kawahara et al., 2013) and *Brachypodium distachyon* (Vogel et al., 2010) were used
81 as a typical example of synteny in plants, with *Arabidopsis thaliana* (Lamesch et al., 2012) providing a
82 well annotated and distant genome for the 3-way comparisons. *Homo sapiens* (Schneider et al., 2017),
83 *Mus musculus* (Waterston and Pachter, 2002) and *Canis familiaris* (Lindblad-Toh et al., 2005) were used
84 for 2-way and 3-way comparisons between mammal genomes. *Drosophila melanogaster* and eleven other
85 fly genomes (Clark et al., 2007) were used for the multi-way comparison and the presence score analysis.

86 In all analyses where gene regions were used, we downloaded the current GFF3 annotations from
87 Ensembl (Yates et al., 2015) and used a a minimum gene size of 100bp.

88 2.2 skm-tools: skip-mer intersection and coverage analyses

89 All skip-mer intersection analyses and skip-mer spectra were computed with our skm-tools, available
90 at <http://github.com/bioinfologics/skm-tools/>. The implementation is based on sorted lists of canonical
91 skip-mers with added attributes such as position on the reference genomes or number of occurrences in a
92 dataset.

93 In particular, the following tools have been used in the preparation of this manuscript:

94 **skm-count** counts the number of occurrences of each distinct canonical skip-mer in a fasta input and
95 outputs a spectra histogram.

96 **skm-multiway-coverage** receives a reference fasta, optionally alongside a GFF3 file and a feature name,
97 and any number of extra datasets. The intersection of skip-mers from all the extra datasets is
98 computed versus the reference dataset, shared-by-all skip-mers statistics are reported as each
99 genome is processed (See Figure 4 for an example progression). If a GFF3 and a feature name is
100 provided, the output will classify the skip-mers according to their presence in regions annotated
101 with the feature, and a file with details of coverage for each feature by each of the extra datasets
102 will be produced.

103 The current implementation of the skm-multiway-coverage tool includes a coverage cut-off that
104 defaults to 1 as this is appropriate for the current study. All skip-mers that are at a higher frequency than
105 the cut-off are eliminated before any analysis. To consider candidate matches for alignment of conserved
106 sequence it is appropriate to discard skip-mers with a higher copy number than your expected number of
107 matches as this will filter repetitive matches including background noise. While our current choice of
108 cut-off at 1 makes sense in a general analysis as the one presented in this manuscript, care needs to be
109 taken to make reasonable choices for future applications.

110 2.3 Coverage score

111 The coverage score is used as a proxy for sequence conservation. To approximate a measure of conserved
112 nucleotides, the coverage is projected over individual nucleotides rather than directly counting shared
113 skip-mers which would introduce redundancy from phased matches. The score for each feature (i.e. gene)
114 versus each genome in the multi-way analyses is calculated as the total number of bases that are included
115 in matching skip-mers from that genome divided by the total number of bases that are covered by valid
116 (i.e. copy number below threshold in the reference) skip-mers from the reference:

$$\text{Coverage score} = \frac{\text{Bases covered by matching skip-mers}}{\text{Bases covered by valid reference skip-mers}}$$

117 The coverage cut-off is applied before any analyses are performed. When using the default cut-off
118 of 1, skip-mers that have a higher copy number in the reference will not be evaluated for scoring and
119 skip-mers that have a single copy in the reference but more than one copy in the scoring genome will not
120 be counted as covered.

121 3 RESULTS

122 3.1 Increasing a skip-mer cycle length and span increases specificity

123 We analysed a genome assembly of *Triticum aestivum* to investigate the effect of the cycle size n in the
124 multiplicity of the skip-mers in a genome. Figure 2 shows how increasing n , and thus the total span of a
125 skip-mer (S), increases the entropy for each skip-mer. The increased entropy decreases the number of
126 copies of each distinct skip-mer in the genome. This ultimately results in more unique skip-mers. In the
127 wheat genome, there are more than twice as many unique skip-mers using $SkipMer(1, 16, 31)$ as there are
128 using $SkipMer(1, 1, 31)$ which corresponds to a 31-mers.

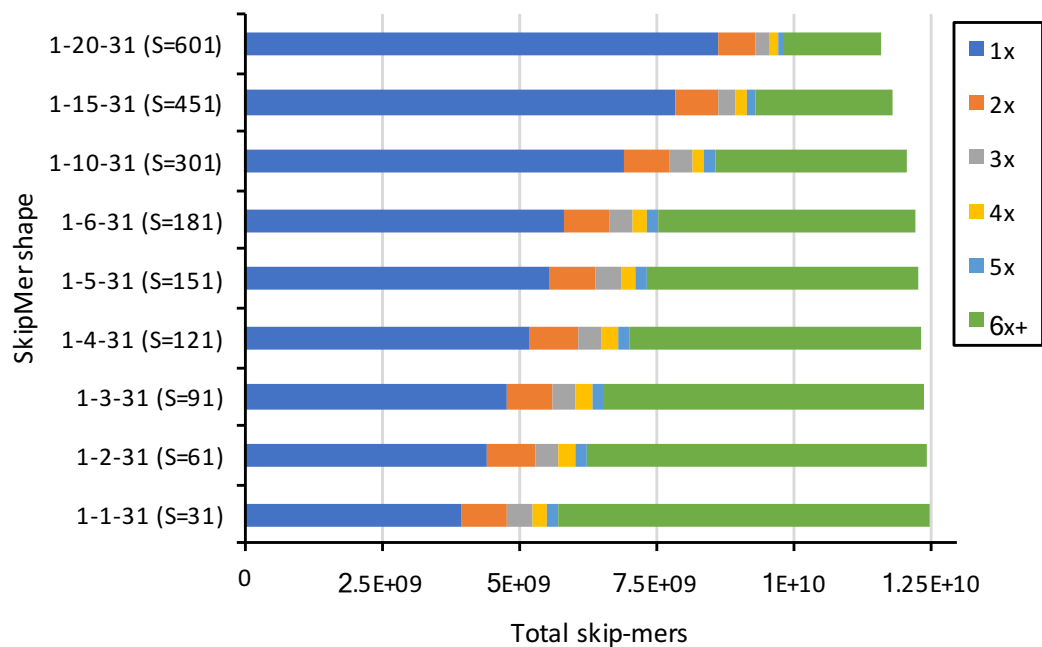


Figure 2. Multiplicity as a function of n for skip-mers in *Triticum aestivum*. Each bar represents the total number of skip-mers in the assembly, so a skip-mer appearing x times contributes x units in the x component. All skip-mers use 31bp ($k = 31$) and 1bp per cycle ($m = 1$).

129 3.2 Using triplet-based cycles increases perfect skip-mer matches in conserved genic 130 sequence between species

131 Synonymous mutations which are not removed by purifying selection, because they do not affect the
132 amino-acid encoded by a trinucleotide codon, produce a cycle-3 modulation in conserved coding regions
133 (Kimura, 1977). Skip-mers with cycle lengths that are a multiple of 3 ($n = 3c$) group first and/or second
134 nucleotides in subsequent in-frame codons, to increase sensitivity on $SkipMer(m, n = 3c, k)$ to detect
135 conserved coding regions.

136 Figure 3 shows how, for the 2-way intersections in (a) and (c), the shared skip-mers in non-genic
137 regions decrease as the span increases, in agreement with the increase of entropy and thus uniqueness. In
138 genic regions, this higher entropy is combined with increased sensitivity for coding sequence resulting in
139 increased matches when $n = 3c$. This effect is further accentuated in the 3-way intersections in (b) and
140 (d), due to independent synonymous mutations in the 3 genomes.

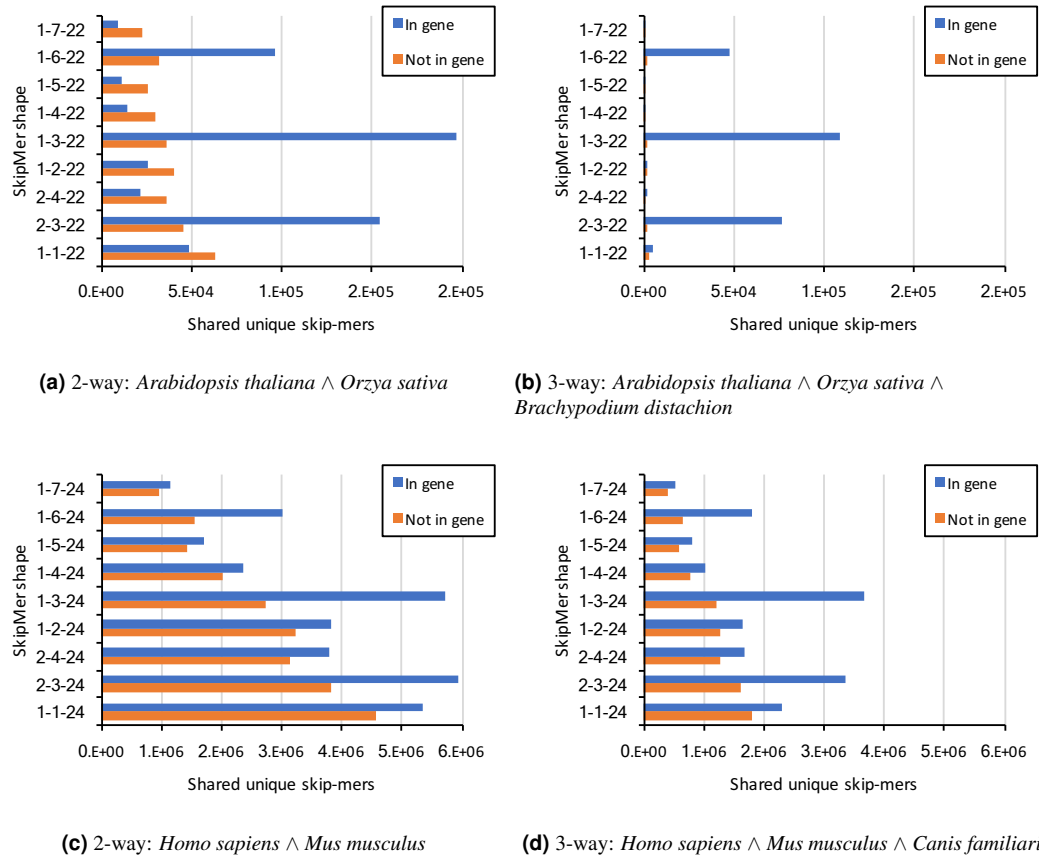


Figure 3. Effect of different combinations of m and n , while keeping k constant, for 2-way and 3-way skip-mer intersections. Only unique skip-mers are considered and skip-mers originating from sequence annotated with gene features on the first genome are classified as "In gene". The skip-mer shapes are sorted along the vertical axis according to total skip-mer span (S), with the largest span on top.

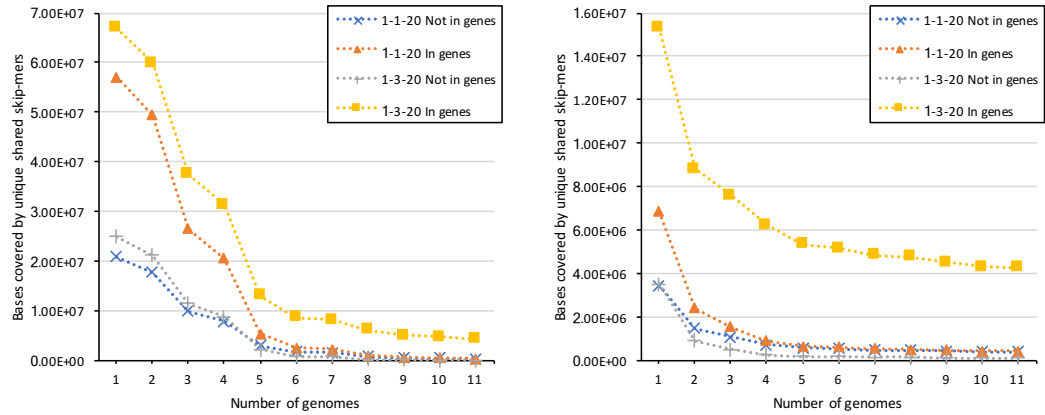
141 The best result in terms of sensitivity for three of the four examples is produced by $SkipMer(1, 3, k)$
 142 which groups the nucleotides according to their position in the codon. While $SkipMer(2, 3, 24)$ presents a
 143 slightly larger number of "In gene" matches in (c), $SkipMer(1, 3, 24)$ with its larger span decreases the
 144 number of "Not in gene" matches, which makes it a better choice.

145 3.3 Conserved sequence from $n=3$ skip-mers enables direct intersection analyses across 146 many samples at diverse evolutionary distances

147 One of the limitations of k -mers for direct intersection analyses among many samples is the decrease in
 148 probability of finding k -mers that are shared across all of the samples. The results from our 2-way and
 149 3-way analyses show that skip-mers are more sensitive to conserved coding sequence. We intersected the
 150 twelve *Drosophila* genomes from Clark et al. (2007) to explore how coverage over the reference from
 151 unique shared skip-mers decreases both in genic and non-genic regions as we progressively include more
 152 samples.

153 In Figure 4 (a) the eleven other genomes are incorporated into the *Drosophila melanogaster* based
 154 analysis starting from the closest to *D. melanogaster* in the phylogeny proposed by Clark et al. (2007):
 155 *D. simulans*, *D. sechellia*, *D. yakuba*, *D. ananassae*, *D. erecta*, *D. pseudoobscura*, *D. willistoni*, *D. virilis*,
 156 *D. mojavensis*, and *D. grimshawi*. In Figure 4 (b) this order is reversed, starting from *D. grimshawi*
 157 and ending on *D. simulans*.

158 The genic intersection computed by $SkipMer(1, 3, 20)$ is less affected by the introduction of extra
 159 genomes due to the increased sensitivity in the conserved coding regions. The k -mer intersection is in
 160 both cases negligible after the fifth genome is introduced into the analyses. The rates of decrease in



(a) Bases covered by unique shared skip-mers when adding more genomes, closer genomes first.

(b) Bases covered by unique shared skip-mers when adding more genomes, more distant genomes first.

Figure 4. Bases covered by unique shared skip-mers shared across all genomes, for sets of different numbers of genomes from the 12 *Drosophila* datasets.

161 (a) and (b) illustrate the effect of genome proximity on the intersection analysis. In (a), starting from
 162 evolutionary proximal genomes, there is more shared content initially, followed by a pronounced drop
 163 as the first 5 genomes are incorporated, and then a smaller decrease. In (b), starting from evolutionary
 164 distant genomes, there is less shared content initially, followed by an exponential decrease. While in (a)
 165 the effect of the higher conservation of skip-mers is reflected in higher shared content throughout, in (b)
 166 a stable point in the decay has been reached, with skip-mers having a significantly larger conservation.

167 3.4 Coverage of matching sequence across many samples using skip-mers with $n=3$ 168 shows higher correlation than using k-mers

169 To explore the advantages of the *SkipMer*(1, 3, k) analyses across divergent genomes we compared the
 170 properties of sequence coverage for *Drosophila melanogaster* to the 11 other *Drosophila* genomes using
 171 *SkipMer*(1, 1, 20), which is equivalent to a 20-mer, and *SkipMer*(1, 3, 20). We implement a base coverage
 172 score as described in section 2.3 and assigned each gene with a length of 100bp or more in *Drosophila*
 173 *melanogaster* a coverage score between 0 and 1 for each of the genomes.

174 A distribution analysis for the scores per genome (Supplementary Figure S1) shows the more divergent
 175 genomes increase their scores for sequence coverage in genes when using *SkipMer*(1, 3, 20). This reflects
 176 the increased sensitivity of cycle-3 skip-mers within coding regions.

177 We computed correlations between the gene scores for *D. melanogaster* from every pair of the other
 178 11 genomes for both *SkipMer*(1, 1, 20) and *SkipMer*(1, 3, 20) (See Supplementary Material Tables S1
 179 and S2, and Figures S2 and S3). Figure 5 shows the comparison between each genome-pair correlation on
 180 *SkipMer*(1, 1, 20) and *SkipMer*(1, 3, 20). There is increased correlation when using cycle-3 skip-mers,
 181 with larger relative improvements on the less correlated genome pairs. This suggests cycle-3 skip-mers
 182 provide a more robust coverage score which can be better used as a proxy for evolutionary pressure and
 183 selection.

184 4 DISCUSSION

185 Increasing the span of skip-mers increases their entropy when sampled from a genome. Using this
 186 increased-entropy analysis unit rather than k-mers will enable more informative analyses with small
 187 adaptations to existing techniques. We expect this key feature of the data points having the same amount
 188 of data (bp) but increased entropy to enable more exhaustive or significant analyses in roughly similar
 189 computational space and time.

190 The analysis of sequence in cyclic groups of $n = 3$ increases sensitivity to detect conserved coding
 191 sequence by grouping the nucleotides in synchronisation with the codon positions. In the typical case
 192 of *SkipMer*($m = 1, n = 3, k$) there will be, for the same group of k contiguous codons, a skip-mer

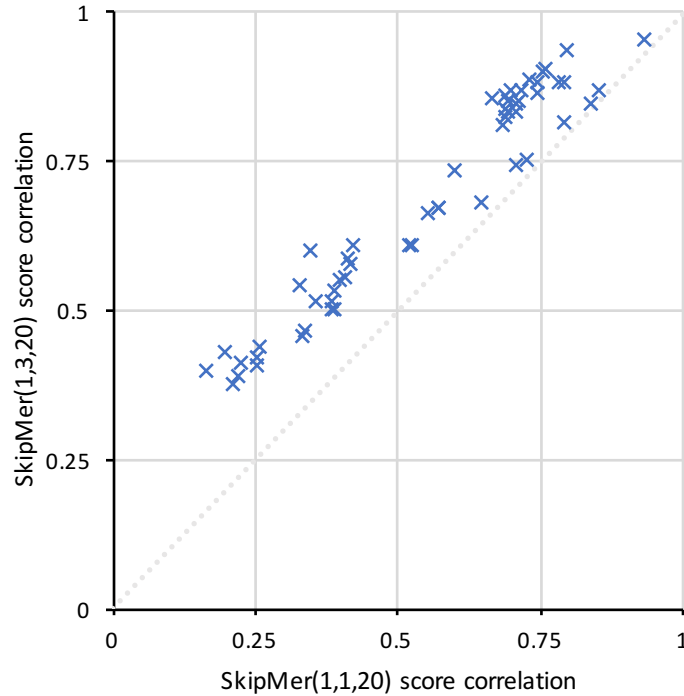


Figure 5. Comparison of correlations for the coverage scores of the *Drosophila melanogaster*. Points above the 1:1 line represent genome pairs with better score correlation in *SkipMer*(1,3,20) than in *SkipMer*(1,1,20) (equivalent to 20-mers).

193 containing all first nucleotides, a skip-mer containing all second nucleotides and a skip-mer containing all
194 third nucleotides. This grouping increases perfect matches in genes from first-nucleotide position and
195 second-nucleotide position skip-mers, providing an alternative to the use of protein translation analyses.

196 When computing multi-genome intersections, there is a stronger signal of conservation across multiple
197 divergent genomes from $n = 3$ skip-mers than from contiguous sequences such as k-mers. Because
198 of the better correspondence between sequence that is actively conserved and the set of matches, the
199 reduction in matches due to the addition of extra genomes in intermediate positions in the phylogeny is
200 less pronounced. Also, the coding sequence coverage by direct matches is a more robust metric, which
201 enables the direct comparison of results produced from different sets of genomes.

202 A complementary effect to the concentration of more conserved sequence, from first and second
203 nucleotides, in the cycle-3 skip-mers is the concentration of more variable sequence, from third nucleotides,
204 in a small number of skip-mers. In the preceding analyses, these more variable nucleotides have been
205 discarded with the noise and repetitions. For applications where a weak signal for variation needs to be
206 analysed, skip-mers can be leveraged to provide a very high entropy set of sequences to give increased
207 discrimination power.

208 Our results suggest skip-mers will have a wide range of applications in bioinformatic analyses. For
209 whole-genome and multi-genome alignment, skip-mers will provide accurate conserved seeds, and more
210 specific matches in complex regions. For evolutionary analyses, skip-mers will allow improved detection
211 of functionally equivalent regions. For RNA-seq and exome analyses, skip-mers will provide a meaningful
212 set of starting seeds or a projection base, thus enabling more distant samples to be analysed together either
213 against a reference or in a reference-free manner.

214 Skip-mers will also be useful in raw read analyses. For classification of sequences, or species detection,
215 skip-mers will provide better clustering of coding regions from a common origin, and could even be
216 used to estimate conservation scores for single reads. Aligning skip-mers from raw reads to one or many
217 references will guide the reconstruction of conserved regions while considering novel variants. These
218 conserved region intersected representations can then be used to quickly characterise the genic space of a

219 genome.

220 CONCLUSIONS

221 We have shown how skip-mer based analyses benefit from extra entropy and sensitivity to outperform
222 k-mer based analyses given the non-random nature of genomic sequence. These principles stand across
223 a wide range of eukaryotic genomes and in different multi-genome scenarios, improving the analysis
224 of conserved coding regions. Common k-mer based techniques can easily adopt skip-mers, due to their
225 many shared properties. Both constructions are reversible strings of nucleotides that can be made strand-
226 agnostic with canonical representations. In general, with a genomic landscape that is shifting to in-field
227 sampling and exploring more diversity than ever before, we expect skip-mers and other evolution-friendly
228 information units to provide the basis for a new generation of biological analyses.

229 AUTHORS CONTRIBUTIONS

230 BJC and GG developed the initial concepts and discussed implications and refinements over time. BJC
231 implemented the first version of the *skm-tools*, ran the analyses and produced the first draft of the
232 manuscript. LY contributed optimisations and improvements for the *skm-tools*. All authors tested the
233 *skm-tools*; discussed analyses, results and improvements; and contributed to the final version of the
234 manuscript.

235 ACKNOWLEDGMENTS

236 The authors would like to thank Wilfried Haerty for his support on the analysis of the *Drosophila* genomes,
237 and Erik Garrison, Mark McMullan, Neil Hall, Norma Paniego, Manfred Grabherr and Federica di Palma
238 for their useful comments and feedback.

239 FUNDING

240 This work was strategically funded by the BBSRC, Core Strategic Programme Grant BB/CSP17270/1 at
241 the Earlham Institute.

242 REFERENCES

- 243 Burkhardt, S. and Kärkkäinen, J. (2003). Better filtering with gapped q-grams. *Fundamenta informaticae*,
244 56(1-2):51–70.
- 245 Chaisson, M. J., Brinza, D., and Pevzner, P. A. (2009). De novo fragment assembly with short mate-paired
246 reads: Does the read length matter? *Genome research*, 19(2):336–346.
- 247 Clark, A. G., Eisen, M. B., Smith, D. R., Bergman, C. M., Oliver, B., Markow, T. A., Kaufman, T. C.,
248 Kellis, M., Gelbart, W., Iyer, V. N., et al. (2007). Evolution of genes and genomes on the drosophila
249 phylogeny. *Nature*, 450(7167):203.
- 250 Clavijo, B. J., Venturini, L., Schudoma, C., Accinelli, G. G., Kaithakottil, G., Wright, J., Borrill, P.,
251 Kettleborough, G., Heavens, D., Chapman, H., Lipscombe, J., Barker, T., Lu, F.-H., McKenzie, N.,
252 Raats, D., Ramirez-Gonzalez, R. H., Coince, A., Peel, N., Percival-Alwyn, L., Duncan, O., TrÄ¶sch,
253 J., Yu, G., Bolser, D. M., Namaati, G., Kerhornou, A., Spannagl, M., Gundlach, H., Haberer, G., Davey,
254 R. P., Fosker, C., Palma, F. D., Phillips, A., Millar, A. H., Kersey, P. J., Uauy, C., Krasileva, K. V.,
255 Swarbreck, D., Bevan, M. W., and Clark, M. D. (2017). An improved assembly and annotation of the
256 allohexaploid wheat genome identifies complete families of agronomic genes and provides genomic
257 evidence for chromosomal translocations. *Genome Research*.
- 258 Ghandi, M., Lee, D., Mohammad-Noori, M., and Beer, M. A. (2014). Enhanced regulatory sequence
259 prediction using gapped k-mer features. *PLoS computational biology*, 10(7):e1003711.
- 260 Hahn, L., Leimeister, C.-A., Ounit, R., Lonardi, S., and Morgenstern, B. (2016). Rasbhari: optimizing
261 spaced seeds for database searching, read mapping and alignment-free sequence comparison. *PLoS
262 computational biology*, 12(10):e1005107.
- 263 Kawahara, Y., de la Bastide, M., Hamilton, J. P., Kanamori, H., McCombie, W. R., Ouyang, S., Schwartz,
264 D. C., Tanaka, T., Wu, J., Zhou, S., et al. (2013). Improvement of the oryza sativa nipponbare reference
265 genome using next generation sequence and optical map data. *Rice*, 6(1):4.

- 266 Kimura, M. (1977). Preponderance of synonymous changes as evidence for the neutral theory of molecular
267 evolution. *Nature*, 267(5608):275–276.
- 268 Lamesch, P., Berardini, T. Z., Li, D., Swarbreck, D., Wilks, C., Sasidharan, R., Muller, R., Dreher, K.,
269 Alexander, D. L., Garcia-Hernandez, M., Karthikeyan, A. S., Lee, C. H., Nelson, W. D., Ploetz, L.,
270 Singh, S., Wensel, A., and Huala, E. (2012). The arabidopsis information resource (tair): improved
271 gene annotation and new tools. *Nucleic Acids Research*, 40(D1):D1202–D1210.
- 272 Lindblad-Toh, K., Wade, C. M., Mikkelsen, T. S., Karlsson, E. K., et al. (2005). Genome sequence,
273 comparative analysis and haplotype structure of the domestic dog. *Nature*, 438(7069):803.
- 274 Mapleson, D., Garcia Accinelli, G., Kettleborough, G., Wright, J., and Clavijo, B. J. (2017). Kat: a k-mer
275 analysis toolkit to quality control ngs datasets and genome assemblies. *Bioinformatics*, 33(4):574–576.
- 276 Schneider, V. A., Graves-Lindsay, T., Howe, K., Bouk, N., Chen, H.-C., Kitts, P. A., Murphy, T. D.,
277 Pruitt, K. D., Thibaud-Nissen, F., Albracht, D., et al. (2017). Evaluation of grch38 and de novo haploid
278 genome assemblies demonstrates the enduring quality of the reference assembly. *Genome research*,
279 27(5):849–864.
- 280 Vogel, J. P., Garvin, D. F., Mockler, T. C., Schmutz, J., Rokhsar, D., Bevan, M. W., Barry, K., Lucas,
281 S., Harmon-Smith, M., Lail, K., et al. (2010). Genome sequencing and analysis of the model grass
282 *brachypodium distachyon*. *Nature*, 463(7282):763–768.
- 283 Waterston, R. H. and Pachter, L. (2002). Initial sequencing and comparative analysis of the mouse genome.
284 *Nature*, 420(6915):520–562.
- 285 Yates, A., Akanni, W., Amode, M. R., Barrell, D., Billis, K., Carvalho-Silva, D., Cummins, C., Clapham,
286 P., Fitzgerald, S., Gil, L., et al. (2015). Ensembl 2016. *Nucleic acids research*, 44(D1):D710–D716.
- 287 Zerbino, D. R. and Birney, E. (2008). Velvet: algorithms for de novo short read assembly using de bruijn
288 graphs. *Genome research*, 18(5):821–829.