

1 **FGTpartitioner: A rapid method for parsimonious delimitation of ancestry**
2 **breakpoints in large genome-wide SNP datasets**

3 Tyler K. Chafin

4 Department of Biological Sciences, University of Arkansas, Fayetteville, AR 72701

5 **Keywords:** Recombination, Ancestry blocks, SNP analysis, phylogenomics

6 **Type:** Applications

7 **Word count:** 1464

8 **Number of figures and tables:** 0

9

10 *Corresponding author and person to whom reprint requests should be addressed:*

11 Tyler K. Chafin

12 Department of Biological Sciences

13 University of Arkansas

14 Fayetteville, AR 72701

15 Voice: +1-970-631-2100

16 e-mail: tkchafin@uark.edu

17

18 **Disclosure statement:** Authors have nothing to disclose.

19

20

21 **Abstract**

22 1. Partitioning large (e.g. chromosomal) alignments into ancestry blocks is a common step in
23 phylogenomic analyses. However, current solutions require complicated analytical assumptions,
24 or are difficult to implement due to excessive runtimes and unintuitive documentation.

25 Additionally, most methods require haplotype phasing, which is often intractable for non-model
26 studies.

27 2. Here, I present an efficient and rapid solution for partitioning large genome alignments into
28 ancestry blocks, which better accommodates non-model diploid organisms in that phasing
29 information is not required *a priori*.

30 3. FGTpartitioner processes a full-chromosome alignment orders of magnitude faster than
31 alternative solutions, while recovering identical results, both via algorithmic improvements and
32 the use of native parallelization.

33 4. FGTpartitioner provides a means for simple and rapid block delimitation in genome-wide
34 datasets as a pretext for phylogenomic analysis. It thus widens the potential uses for researchers
35 studying phylogenetic processes across large, non-model genomes. Complete code and
36 documentation for FGTpartitioner are available as an open-source repository on GitHub:
37 <https://github.com/tkchafin/FGTpartitioner>

38

39 **Acknowledgements**

40 I would like to thank vonHoldt *et al.* (2016) and Kukekova *et al.* (2018) for making their
41 raw datasets available via the NCBI SRA, and the Arkansas High Performance Computing
42 Center and my Ph. D. advisors Drs. Michael and Marlis Douglas for access to computational
43 resources which I used for developing and benchmarking this program.

44

45

46 **Introduction**

47 Inferring genome-wide variation in localized ancestry is a critical step in reconstructing
48 species trees (Springer & Gatesy, 2018), and as such has become a major goal in phylogenomic
49 studies utilizing whole-genome datasets. Analyzing this variation is necessary to understand
50 processes such as horizontal transfer and hybridization, and has led to insights into adaptive
51 introgression (Fontaine et al., 2015), genome-wide selection (Sabeti et al., 2002), disease
52 susceptibility (Seldin, Pasaniuc, & Price, 2011), and ancestral human demographics
53 (Sankararaman et al., 2014). However, partitioning such alignments is a non-trivial task,
54 particularly for researchers utilizing non-model organisms for which limited genomic reference
55 data exists.

56 Multiple approaches have been proposed for delimiting ancestry blocks in genomes (i.e.
57 establishing recombination breakpoints), which generally fall into one of two categories: those
58 which require dense or phased genotypic data (Liu et al., 2013); and those with complex
59 analytical assumptions which require the definition of informative prior probability distributions
60 and are computationally intensive (Dutheil et al., 2009). Both conditions are problematic for
61 genome-scale studies of non-model diploid organisms, where large-scale resequencing and
62 phased reference data are unavailable, and genomes are often sequenced at low coverage.

63 I here describe a solution, FGTpartitioner, which is specifically designed for use with
64 non-model genomic data without the need for high-quality phased reference data or dense
65 population-scale sampling. FGTpartitioner delimits chromosome scale alignments using a fast
66 interval-tree approach which detects pairwise variants which violate the four-gametes

67 assumption (Hudson & Kaplan, 1985), and rapidly resolves a most parsimonious set of
68 recombination events to yield non-overlapping intervals which are both unambiguously defined
69 and consistent regardless of processing order. These sub-alignments are then suitable for separate
70 phylogenetic analysis, or as a ‘first pass’ which may facilitate parallel application of finer-
71 resolution (yet more computationally intensive) methods.

72

73

74 **Program Description**

75 For ease of application, inputs are required to follow the widely used VCF format
76 (Danecek et al., 2011). Users may provide parameter settings as arguments in the command-line
77 interface which can restrict block delimitation to a certain chromosome (<-c> flag), with the
78 option to additionally target a region via start (<-s>) and end (<-e>) coordinates. Parallel
79 computation is also possible (<-t>) for particularly large alignments. After parsing user-inputs,
80 the workflow of FGTPartitioner is as follows:

- 81 (1) For each SNP, perform four-gamete tests sequentially for rightward neighboring records,
82 up to a maximal physical distance (if defined; <-d>) and stopping when a conflict
83 (='interval') is found. Intervals are stored in a self-balancing tree. When using
84 multiprocessing (<-t>), daughter processes are each provided an offset which guarantees
85 a unique pairwise SNP comparison for each iteration
- 86 (2) Merge interval trees of daughter processes (if <-t 2 or greater>)
- 87 (3) Assign rank k per-interval, defined as the number of SNP records (indexed by position)
88 spanned by each interval

89 (4) Order intervals by k ; starting at $\min(k)$, resolve conflicts as follows: For each candidate
90 recombination site (defined as the mid-point between SNPs), compute the depth d of
91 spanning intervals. The most parsimonious breakpoint is that which maximizes d
92

93 These algorithm choices have several implications: indexing SNPs by physical position
94 guarantees that the same recombination sites will be chosen given any arbitrary ordering of
95 SNPs; and defining breakpoints as physical centerpoints between nodes means that
96 monomorphic sites will be evenly divided on either side of a recombination event. Because
97 monomorphic sites by definition lack phylogenetic information, they cannot be unambiguously
98 assigned to any particular ancestry block, thus my solution is to evenly divide them.

99 Heterozygous sites in diploid genomes are dealt with in multiple ways. By default,
100 FGTpartitioner will randomly resolve haplotypes. The user can select an alternate resolution
101 strategy (<-r>) which will either treat a SNP pair as failing if any resolution meets the four-
102 gamete condition, or as passing if any possible resolution passes (i.e. the ‘pessimistic’ and
103 ‘optimistic’ strategies of Wang *et al.*, 2010).

104

105 **Benchmarking and comparison to other software**

106 Performance benchmarking was conducted FGTpartitioner and similar methods using
107 publicly available data from the NCBI SRA. I ran tests using a 4-taxon alignment of *Canis lupus*
108 chromosome 1 (~120 Mb), generated from raw reads (*C. lupus*: SRR7107787; *C. rufus*:
109 SRR7107783; *C. latrans*: SRR1518489; *Vulpes vulpes*: SRR5328101-115). Reads were first
110 aligned with bowtie2 (Langmead & Salzberg, 2012) and subsequently genotyped using the
111 HaplotypeCaller and best practices from GATK (McKenna *et al.*, 2010). Tests were all

112 conducted within a 64-bit Linux environment, on machines equipped with dual 8-core Intel Xeon
113 E5-4627 3.30GHz processors and 265GB RAM.

114 Memory usage for processing a full alignment with FGTPartitioner on 16 cores peaked at
115 18GB (with less memory required when using fewer cores). Total runtime scales linearly with
116 the maximum allowable physical distance between SNPs to check for four-gamete conditions,
117 with <-d 250000> taking ~50 minutes and <-d 100000> taking 36 minutes to process an
118 alignment comprising >2 million variants. Increasing number of cores offer diminishing returns,
119 with 1, 2, 4, 8, and 16 cores taking 39, 20, 12, 10, and 9 minutes, respectively (to process a 1
120 million base subset of the 4-taxon chromosome 1 alignment). Runtimes also scale linearly with
121 dataset size (in total variants). Haplotype resolution strategy (random, pessimistic, or optimistic)
122 was not found to impact runtimes.

123 BA3-SNPs-autotune was used to find optimal mixing parameters for each run, with
124 exploratory analyses employing 10,000 MCMC generations in length. A maximum of 10
125 exploratory analyses were conducted for each data file. The number of repetitions required to
126 find optimal mixing parameters was recorded for each, and mixing parameters verified to
127 produce adequate MCMC acceptance rates (i.e., $0.2 < \text{acceptance rate} < 0.6$). All tests were
128 performed on a computer equipped with dual Intel Xeon E5-4627 3.30GHz processors, 265GB
129 RAM, and with a 64-bit Linux environment. Since neither program is multithreaded, a single
130 processor core was used per analysis.

131 For comparison, I attempted to delimit recombinatorial genes using an alternate method
132 based on the four-gamete test, RminCutter.pl (https://github.com/RILAB/rmin_cut). Because
133 RminCutter does not handle diploid data, I created pseudohaplotypes by randomly resolving
134 heterozygous sites (pseudoHaploidize.py; <https://github.com/tkchafin/scripts>). I also ran the

135 MDL approach (Ané, 2011), which uses a parsimony criterion to assign breakpoints separating
136 phylogenetically homogenous loci, while penalizing for the number of breakpoints. Neither
137 MDL nor RminCutter could complete an analysis on the full chromosome 1 dataset in a
138 reasonable time span (capped at 10 days), so I first divided the alignment into blocks of ten-
139 thousand parsimony-informative sites each (as these determine runtime, not total length). This
140 produced 21 sub-alignments. MDL took 24 hours across 48 cores to complete a single 5 Mb
141 segment, thus was not further explored. RminCutter.pl took an average of 22 hours per segment
142 (single-threaded). Of note, RminCutter and FGTpartitioner yielded identical results when using
143 comparable settings restricted to a pseudo-haploid dataset. FGTpartitioner thus offers an efficient
144 and methodologically simple solution for ancestry delimitation in non-model diploid organisms
145 with large genomes.

146

147 Conclusion

148 FGTpartitioner has several advantages over similar methods: 1) algorithmic and
149 performance enhancements allow it to perform orders of magnitude faster, thus extending
150 application to larger genomes; and 2) the flexibility of diploid resolution strategies precludes the
151 need for haplotype phasing *a priori*. Validation using empirical data indicated the suitability of
152 FGTpartitioner for highly distributed work on high-performance computing clusters, with
153 parallelization easily facilitated by built-in options in the command-line interface. Additionally,
154 runtime and memory profiling indicate its applicability on modern desktop workstations as well,
155 when applied to moderately sized datasets. Thus, it provides an efficient and user-friendly
156 solution to alignment pre-processing for phylogenomic studies, or as a method of breaking up

157 large alignments in order to efficiently distribute computation for more rigorous recombination

158 tests.

159

160 **References**

- 161
- 162 Ané, C. (2011). Detecting phylogenetic breakpoints and discordance from genome-wide
163 alignments for species tree reconstruction. *Genome Biology and Evolution*, 3(1), 246–258.
164 doi:10.1093/gbe/evr013
- 165 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., ... Durbin, R.
166 (2011). The variant call format and VCFtools. *Bioinformatics*.
167 doi:10.1093/bioinformatics/btr330
- 168 Dutheil, J. Y., Ganapathy, G., Hobolth, A., Mailund, T., Uyenoyama, M. K., & Schierup, M. H.
169 (2009). Ancestral population genomics: The coalescent hidden Markov model approach.
170 *Genetics*, 183(1), 259–274. doi:10.1534/genetics.109.103010
- 171 Fontaine, M. C., Pease, J. B., Steele, A., Waterhouse, R. M., Neafsey, D. E., Sharakhov, I. V., ...
172 Besansky, N. J. (2015). Extensive introgression in a malaria vector species complex
173 revealed by phylogenomics. *Science*, 347(6217), 1258524. doi:10.1126/science.1258522
- 174 Hudson, R. R., & Kaplan, N. L. (1985). Statistical properties of the number of recombination
175 events in the history of a sample of DNA sequences. *Genetics*, 111(1), 147–164.
- 176 Kukekova, A. V., Johnson, J. L., Xiang, X., Feng, S., Liu, S., Rando, H. M., ... Graphodatsky, A.
177 S. (2018). Red fox genome assembly identifies genomic regions associated with tame and
178 aggressive behaviours. *Nature Ecology & Evolution*, 2(September). doi:10.1038/s41559-
179 018-0611-6
- 180 Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature*
181 *Methods*. doi:10.1038/nmeth.1923
- 182 Liu, Y., Nyunoya, T., Leng, S., Belinsky, S. A., Tesfaigzi, Y., & Bruse, S. (2013). Softwares and
183 methods for estimating genetic ancestry in human populations. *Human Genomics*, 7(1), 1–7.

- 184 doi:10.1186/1479-7364-7-1
- 185 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., ... DePristo,
186 M. A. (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
187 generation DNA sequencing data. *Genome Research*. doi:10.1101/gr.107524.110
- 188 Sabeti, P. C., , David E. Reich*, J. M. H., Haninah Z. P. Levine*, Daniel J. Richter*, S. F. S.,
189 Stacey B. Gabriel*, Jill V. Platko*, Nick J. Patterson*, G. J. M., Hans C. Ackerman‡, Sarah
190 J. Campbell‡, D. A., Richard Cooperk, Dominic Kwiatkowski‡, R. W. & E. S. L., & Eisner,
191 T. (2002). Detecting recent positive selection in the human genome from haplotype
192 structure. *Nature*, 419(October). doi:10.1038/nature01027
- 193 Sankararaman, S., Mallick, S., Dannemann, M., Prufer, K., Kelso, J., Paabo, S., ... Reich, D.
194 (2014). The Genomic Landscape of Denisovan and Neanderthal Ancestry in Present-Day
195 Humans. *Nature*, 507, 354–357. doi:10.1016/j.cub.2016.03.037
- 196 Seldin, M. F., Pasaniuc, B., & Price, A. L. (2011). New approaches to disease mapping in
197 admixed populations. *Nature Reviews Genetics*, 12(8), 523–528. doi:10.1038/nrg3002
- 198 Springer, M. S., & Gatesy, J. (2018). Delimiting Coalescence Genes (C-Genes) in
199 Phylogenomic Data Sets, 1–19. doi:10.3390/genes9030123
- 200 vonHoldt, B. M., Cahill, J. A., Fan, Z., Gronau, I., Robinson, J., Pollinger, J. P., ... Wayne, R. K.
201 (2016). Whole-genome sequence analysis shows that two endemic species of North
202 American wolf are admixtures of the coyote and gray wolf. *Science Advances*, 2(7),
203 e1501714–e1501714. doi:10.1126/sciadv.1501714
- 204 Wang, J., Moore, K. J., Zhang, Q., de Villena, F. P.-M., Wang, W., & McMillan, L. (2010).
205 Genome-wide compatible SNP intervals and their properties. *Proceedings of the First ACM
206 International Conference on Bioinformatics and Computational Biology*, 43–52.

207 doi:10.1145/1854776.1854788

208