

1 SUPPLEMENTARY NOTE

2 FRAMEWORK OF *CHROMO.CRAWL*

3 *Chromo.Crawl* is one of the primary functions of *PhyloWGA* that seeks to test for the presence of
4 “supergenes” (i.e., contiguous loci with evidence of a shared genealogical history) across whole-
5 genome alignments (WGAs). In this section, we discuss several major benefits of the *Chromo.Crawl*
6 framework when compared to standard “statistical binning” procedures (i.e., Mirarab *et al.* 2014;
7 Bayzid *et al.* 2015; Liu and Edwards 2015; Adams and Castoe 2019a) that have been developed for a
8 similar purpose. *Chromo.Crawl* acts by “crawling” across a WGA from one genomic window to the
9 next adjacent one, and for each pair of windows, a likelihood ratio test is applied using
10 CONCADEPILLAR (Leigh *et al.*, 2008). These two features of *Chromo.Crawl* (crawling across
11 contiguous genomic windows and implementing a likelihood ratio test) represent important and distinct
12 advantages over standard statistical binning procedures, such as the pipelines described in (Mirarab *et*
13 *al.* 2014; Bayzid *et al.* 2015; Liu and Edwards 2015; Adams and Castoe 2019a,b).

14 First, the procedure of crawling from one genomic region to the next means that the specific genomic
15 location and context are considered by *Chromo.Crawl*. Recombination acts to effectively decouple the
16 genealogical histories of adjacent loci along a chromosome, such that the genealogical histories of loci
17 separated by a recombination event will be unlinked. Conversely, two loci will share the same
18 underlying gene tree in the absence of recombination. Because the probability of crossing over between
19 two loci is a function of their distance to one another, adjacent loci are more likely to share a common
20 genealogical history than more distant loci, and the “crawling” behavior of *Chromo.Crawl* is
21 effectively designed to incorporate this spatial information while testing for phylogenetic congruency
22 along a WGA. This procedure is in stark contrast to typical techniques for statistical binning that use

23 differences in estimated topologies among gene trees obtained across distant genomic regions and even
24 chromosomes (Bayzid *et al.*, 2015; Adams and Castoe, 2019). Indeed, the recent studies of Adams and
25 Castoe (2019a,b) found evidence that standard statistical binning is likely to bias species tree estimates
26 by incorrectly constructing supergenes from loci sampled from distant genomic regions and diverse
27 chromosomes.

28 Another primary advantage of the *Chromo.Crawl* framework when compared to traditional binning
29 techniques is the use of a formal likelihood ratio test for assessing evidence of phylogenetic congruency
30 among loci. CONCATEPILLAR is based on the principles of likelihood ratio tests of topological
31 congruency, such as Huelsenbeck and Bull (1996), and therefore benefits from the statistical foundation
32 provided by such tests, lending itself to predictable statistical properties. In comparison, standard
33 statistical binning pipelines typically use an *ad hoc* procedure based on differences in topologies among
34 bootstrap replicates. For standard statistical binning, if the fraction of bootstrap replicates leading to
35 phylogenetic incongruence is below a certain user-defined cutoff, then the loci are deemed to be
36 congruent. *Chromo.Crawl* instead uses CONCATEPILLAR to test whether a model of congruency
37 (i.e., single genealogical history) or incongruency (i.e., distinct genealogical histories) is a better fit for
38 a given alignment. If the data support a congruent model for adjacent genomic windows, then the set
39 of contiguous loci are concatenated together as a supergene, and conversely, the loci are deemed
40 independent if a model of discordance is supported based on the likelihood ratio.

41 BENCHMARKING

42 The computational complexity of the functions of *PhyloWGA* is highly dynamic and depends upon a
43 number of factors, including (but not limited to): (1) number of taxa in the WGA, (2) total length of
44 the WGA, (3) experimental design (i.e., number, length, and spatial distribution of genomic windows),

45 (4) evolutionary model complexity (e.g., JC69 versus GTR model versus model selection, and
46 recombination rate), and (5) particular analysis type (*Chromo.Phylome* versus *Chromo.Crawl*). Like
47 most standard phylogenetic analyses, the complexities of the algorithms scale with both the number of
48 taxa and the total length of the WGA, whereas the particular experimental design (i.e., distribution of
49 windows defined by the user) may have dynamic effects on computational speed. For example,
50 *Chromo.Phylome* using a single, long window is likely to be faster than if that window is partitioned
51 into multiple shorter windows because *Chromo.Phylome* will be applied to each window separately
52 (i.e., phylogenetic tree model will be fit to each window instead of a single window). The particular
53 evolutionary model used will also impact the efficiency of *PhyloWGA* by modulating the number of
54 parameters estimated for each window. *Chromo.Phylome* conducts phylogenetic inference across
55 genomic windows, whereas *Chromo.Crawl* applies tests of congruency among windows. Thus, the
56 *Chromo.Crawl* function will run substantially slower than *Chromo.Phylome* because it both infers trees
57 and implements a likelihood ratio test, and thus, the efficiency of *Chromo.Crawl* will also likely
58 fluctuate as it crawls along a chromosome. For example, we ran two versions of *Chromo.Phylome* on
59 the primate WGA using a single thread of a 2.8 GHz CPU that represent two different nucleotide model
60 settings, and we found the following run times: GTR model for all windows (run time of 10.35 hours)
61 and GTR+ Γ model for all windows (run time of 22.07 hours). That is, estimating the shape parameter
62 of the Γ model of among-site variation approximately doubled the running time for this example.
63 Additionally, we found that our *Chromo.Crawl* simulation demonstration (described in the next section
64 and results shown in Figures S1-S4) of 100 kb alignments ran for an average of approximately 2.4
65 hours each using a 2.8 GHz CPU with two threads. The *Chromo.Phylome* analyses for these 100 kb
66 simulated WGAs were far quicker (i.e., each ran under a minute using a single thread of a 1.7 GHz
67 Dual-Core Intel Core i7 CPU).

69 In any case, the functions *Chromo.Phylome* and *Chromo.Crawl* both provide indications of
70 computational time required for analyses. For example, *Chromo.Phylome* provides an estimate of the
71 total time needed for analysis (based on the first window), and the percentage of total windows that
72 have been analyzed as the algorithm proceeds along the WGA. Similarly, *Chromo.Crawl* tracks the
73 progress of the algorithm by printing the percentage of total windows that have been “crawled” over,
74 and *Chromo.Crawl* also allows the user to specify the number of cores with the
75 “numeric.NumberOfCores” argument. Finally, *PhyloWGA* now includes a function
76 *Organize.ParallelPhyloWGA* that streamlines and organizes WGAs and *PhyloWGA* scripts for parallel
77 analyses. This function partitions a WGA into a number of user-defined subsets placed within
78 directories for easy execution and parallel analysis.

79 EXPLORING THE ACCURACY OF *PHYLOWGA* ON SIMULATED WGAS

80 To explore the accuracy of *PhyloWGA*, we conducted an array of simulation analyses that are inspired
81 by the Primate dataset (shown in Figure 1), and that varied in recombination rate r . We simulated
82 genealogies along 100 kb alignments with the program *ms* (Hudson, 2002) using a 10-taxon species
83 tree inspired by the relationships of Human, Chimpanzee, Gorilla, Orangutan, Macaque, Marmoset,
84 Tarsier, Galago, Lemur, and Rat provided in a previous study (Song, *et al.*, 2012) and three different
85 recombination rate values ($r = 10^{-9}$, 10^{-8} , or 10^{-7} per site per generation). For each genealogy and
86 associated alignment block output by *ms* (i.e., subsets of the 100 kb alignment separated by
87 recombination events), we then simulated nucleotide sequence alignments using the program *Seq-Gen*
88 (Rambaut and Grass, 1997) and a HKY model (Hasegawa *et al.*, 1985) with the following parameters:
89 transition/transversion ratio of 4.6, and base equilibrium frequencies of 0.3 (A), 0.2 (C), 0.2 (G), and
90 0.3 (T). In our simulations, we used the population-scaled mutation rate $\theta = 4N\mu = 0.00104$ for a
91 diploid effective population size $N = 10^4$ (Takahata, 1993) and a mutation rate $\mu = 2.6 \times 10^{-8}$ per site

92 per generation (Narasimhan *et al.*, 2017). These HKY parameters were inspired by previous studies of
93 primate relationships and species tree analyses (Burgess and Yang, 2008; Koch and DeGiorgio, 2020).
94 These resulting simulated alignments were then concatenated together to form a single, 100 kb WGA
95 fasta file. We next conducted two *PhyloWGA* analyses: (1) *Chromo.Crawl* with one kb windows and
96 step sizes, and (2) using these *Chromo.Crawl* coordinates of concatenated windows to construct gene
97 trees using *Chromo.Phylome* with nucleotide substitution model selection. For each combination of
98 simulation parameters, we repeated the process nine times, and we plotted Robinson-Foulds (RF)
99 distances (Robinson and Foulds, 1981) for each nucleotide site between the true simulated genealogies
100 and their corresponding inferred trees from *Chromo.Phylome*. Additionally, we plotted the location of
101 recombination events that resulted in topology swaps (i.e., red lines and dots in Figures S2-S4) and the
102 location of breakpoints reconstructed with *Chromo.Crawl* (i.e., alternating light and dark gray blocks
103 that represent concatenated windows in Figures S2-S4).

104 As demonstrated with the simulations, the accuracy of *PhyloWGA* is dynamic in response to the
105 recombination rate r . For example, the mean RF distance is much lower under the low recombination
106 rate ($r = 10^{-9}$ per site per generation) simulations when compared with the high recombination rate (r
107 $= 10^{-7}$ per site per generation) scenarios (Figures S1a versus S1c). These results can be observed when
108 comparing the true, simulated recombination breakpoints that result in topology changes (i.e., red lines
109 in Figures S2-S4) with the inferred breakpoints recovered by *PhyloWGA* (alternating light and dark
110 gray blocks in Figures S2-S4). For example, the inferred and true recombination breakpoints appear to
111 be more accurately reconstructed with the low recombination simulations (Figure S2) compared to the
112 high recombination rate simulations (Figure S4). Under the large recombination rate ($r = 10^{-7}$ per site
113 per generation), the lengths of recombination-free nucleotide stretches are small because there are a
114 large number of observed recombination events across the alignment (Figure S4).

115 A fundamental goal of *Chromo.Crawl*, and *PhyloWGA* more generally, is to improve the accuracy of
116 chromosome-scale phylogenetic analyses by flexibly and adaptively incorporating information about
117 gene tree signal and variability of genomic regions. Thus, we sought to understand its performance at
118 achieving this goal by comparing its accuracy with three alternative strategies that do not consider
119 evidence for (or against) shared gene trees among adjacent windows: (1) trees inferred across an
120 alignment using a fixed window size (“FIXED”), (2) trees inferred across random windows without
121 regards to location (“RANDOM”), and (3) a single tree is inferred by concatenating the entire 100kb
122 alignment (“CONCAT”), such that gene tree variability is ignored. The RANDOM approach can be
123 described in three steps: (1) randomly sample a starting site (with uniform probability across the 100kb
124 alignment) to denote beginning of a window, (2) define the end of a window by the site position that is
125 located one, two, four, five, 10, or 20 kb downstream of the start site sampled in the previous step, and
126 (3) this process is repeated 25 times to obtain 25 windows that are positioned randomly throughout the
127 chromosome. *Chromo.Crawl* represents an adaptive, genome-informed approach designed to respond
128 to shared phylogenetic signal (or lack thereof) among adjacent loci, and thus, we predicted that it would
129 yield more accurate inferences when compared to the either the FIXED or RANDOM procedures, as
130 well as the CONCAT approach that ignores gene tree variation. We simulated nine 100kb replicate
131 datasets according to the same procedures as before (i.e., Figures S1-S4), and we applied each of the
132 four strategies (*Chromo.Crawl*, FIXED, RANDOM, and CONCAT) across a range of different window
133 sizes (one, two, four, five, 10, and 20 kb) for the low ($r = 10^{-9}$), medium ($r = 10^{-8}$), and high
134 recombination rates ($r = 10^{-7}$). We measured the mean RF distance (between the true and inferred tree
135 at each site along the alignment) across replicates to quantify differences between the three methods in
136 terms of phylogenetic accuracy. That is, we sought to understand whether the core function of
137 *Chromo.Crawl* (i.e., infer more accurate gene trees using longer, concatenated windows) was indeed
138 successful at improving phylogenetic accuracy across simulated chromosomes, and in particular, its

139 success when compared to the RANDOM, FIXED, and CONCAT approaches that are agnostic to
140 diverse spans of phylogenetic signal across an alignment.

141 Our results provide evidence that *Chromo.Crawl* does indeed provide more accurate inferences than
142 either the RANDOM, FIXED, or CONCAT approaches (Figure S5). For example, the minimum mean
143 RF distance of *Chromo.Crawl* was always smaller than or comparable to the corresponding window
144 size of both the FIXED and RANDOM in the majority of cases. This is perhaps best illustrated when
145 comparing the one kb results of *Chromo.Crawl* (mean RF = 0.09) to the one kb FIXED (mean RF =
146 0.37) and RANDOM (mean RF = 0.38) analyses, respectively, for the low recombination results
147 (Figure S5a). Accuracy tends to increase for both the FIXED and RANDOM strategies with larger
148 window sizes (i.e., approximately 10 to 20 kb), while in contrast, *Chromo.Crawl* appears to adaptively
149 adjust to the optimal window size across each recombination rate setting, finding a minimum mean RF
150 distance with windows of size two kb for the low recombination rate (Figure S5a), one kb for the
151 medium recombination rate (Figure S5b), and four to 20 kb for the high recombination rate (Figure
152 S5c). Moreover, even the smallest minimum window size for *Chromo.Crawl* of one kb performed
153 comparatively well, indicating that it would be sufficient for a user to just specify a minimum window
154 size of one kb, with *PhyloWGA* adaptively changing window sizes across a chromosome and providing
155 improved phylogenetic accuracy over fixed user-defined window approaches. The pervasive
156 abundance of recombination events (e.g., see Figure S4 as an illustration) reduced phylogenetic
157 accuracy across the board in our highest recombination rate simulations (Figure S5c). Yet, we found
158 evidence that *Chromo.Crawl* nonetheless outperformed the FIXED, RANDOM, and CONCAT
159 approaches in these challenging scenarios. These results suggest that the genome-informed approach
160 of *Chromo.Crawl* does indeed provide meaningful improvements in phylogenetic accuracy over
161 alternative approaches that do not attempt account for the tendency of adjacent genomic regions to
162 share a common genealogical history.

163 Collectively, these results (Figures S1-S5) promote the user-friendly framework *PhyloWGA* as a tool
164 for improving chromosome-scale phylogenetic accuracy. In general, we expect the accuracy of
165 *PhyloWGA* to be dynamic as a function of both experimental and evolutionary parameters (i.e., Figs.
166 S2-S5). As expected, these results suggest that genomic regions with low recombination rates may be
167 more accurately reconstructed with *PhyloWGA* (and most any other approach, including RANDOM
168 and FIXED strategies; Figure S5), when compared with regions with high recombination rate. In any
169 case, we encourage users to carefully consider the evolutionary context (e.g., recombination rates) of
170 their particular datasets when analyzing with *PhyloWGA* or any other phylogenetic analyses.

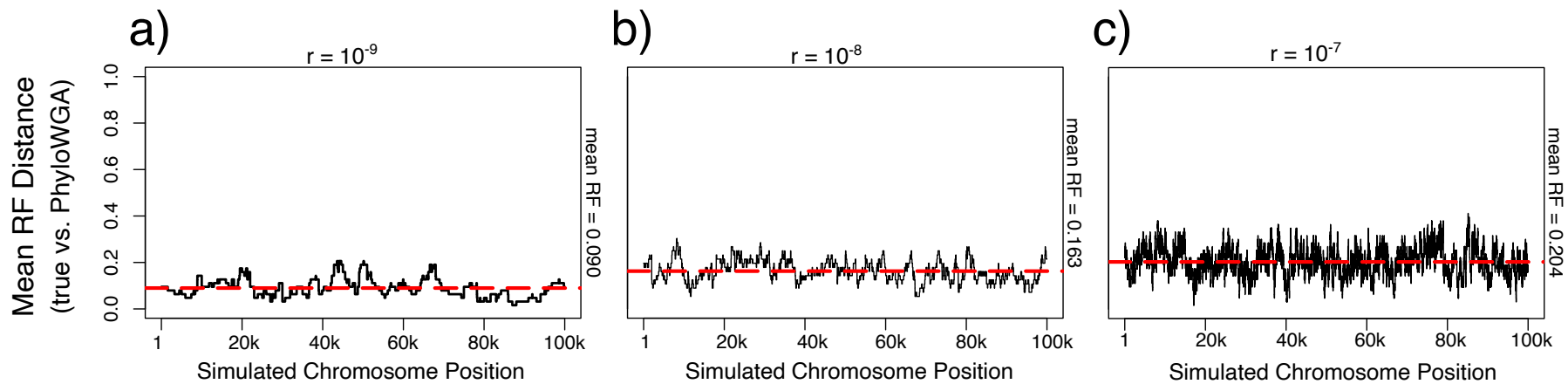
171 REFERENCES

- 172 Adams, Richard H. and Castoe, T.A. (2019) Statistical binning leads to profound model violation due to
173 gene tree error incurred by trying to avoid gene tree error. *Mol. Phylogenet. Evol.*, **134**, 164-171.
- 174 Adams, Richard H and Castoe, T.A. (2019) Supergene validation: A model-based protocol for assessing
175 the accuracy of non-model-based supergene methods. *MethodsX*, **6**, 2181–2188.
- 176 Bayzid, M.S. *et al.* (2015) Weighted statistical binning: enabling statistically consistent genome-scale
177 phylogenetic analyses. *PLoS One*, **10**, e0129183.
- 178 Burgess, R. and Yang, Z. (2008) Estimation of hominoid ancestral population sizes under Bayesian
179 coalescent models incorporating mutation rate variation and sequencing errors. *Mol. Biol. Evol.*,
180 **25**, 1979–1994.
- 181 Hasegawa, M. *et al.* (1985) Dating of the human-ape splitting by a molecular clock of mitochondrial
182 DNA. *J. Mol. Evol.*, **22**, 160–174.

- 183 Hudson,R.R. (2002) Generating samples under a Wright–Fisher neutral model of genetic variation.
184 *Bioinformatics*, **18**, 337–338.
- 185 Huelsenbeck,J.P. and Bull,J.J. (1996) A likelihood ratio test to detect conflicting phylogenetic signal.
186 *Syst. Biol.*, **45**, 92–98.
- 187 Koch,H. and DeGiorgio,M. (2020) Maximum likelihood estimation of species trees from gene trees in
188 the presence of ancestral population structure. *Genome Biol. Evol.*, **12**, 3977–3995.
- 189 Leigh,J.W. *et al.* (2008) Testing congruence in phylogenomic analysis. *Syst. Biol.*, **57**, 104–115.
- 190 Liu,L. and Edwards,S. V (2015) Comment on “Statistical binning enables an accurate coalescent-based
191 estimation of the avian tree”. *Science*, **350**, 171.
- 192 Mirarab,S. *et al.* (2014) Statistical binning enables an accurate coalescent-based estimation of the avian
193 tree. *Science*, **346**, 1250463.
- 194 Narasimhan, V. M., *et al.* (2017). Estimating the human mutation rate from autozygous segments
195 reveals population differences in human mutational processes. *Nat. com.*, **8**, 1–7.
- 196 Rambaut,A. and Grass,N.C. (1997) Seq-Gen: an application for the Monte Carlo simulation of DNA
197 sequence evolution along phylogenetic trees. *Bioinformatics*, **13**, 235–238.
- 198 Robinson,D.F. and Foulds,L.R. (1981) Comparison of phylogenetic trees. *Math. Biosci.*, **53**, 131–147.
- 199 Song, Sen, Liang Liu, Scott V. Edwards, and Shaoyuan Wu. (2012) "Resolving conflict in eutherian
200 mammal phylogeny using phylogenomics and the multispecies coalescent model." *Proc. Natl.*
201 *Acad. Sci.*, **37**, 14942–14947.

202 Takahata, Naoyuki. (1993) "Allelic genealogy and human evolution." *Mol. Biol. Evol.*, **10**: 2–22.

203



204

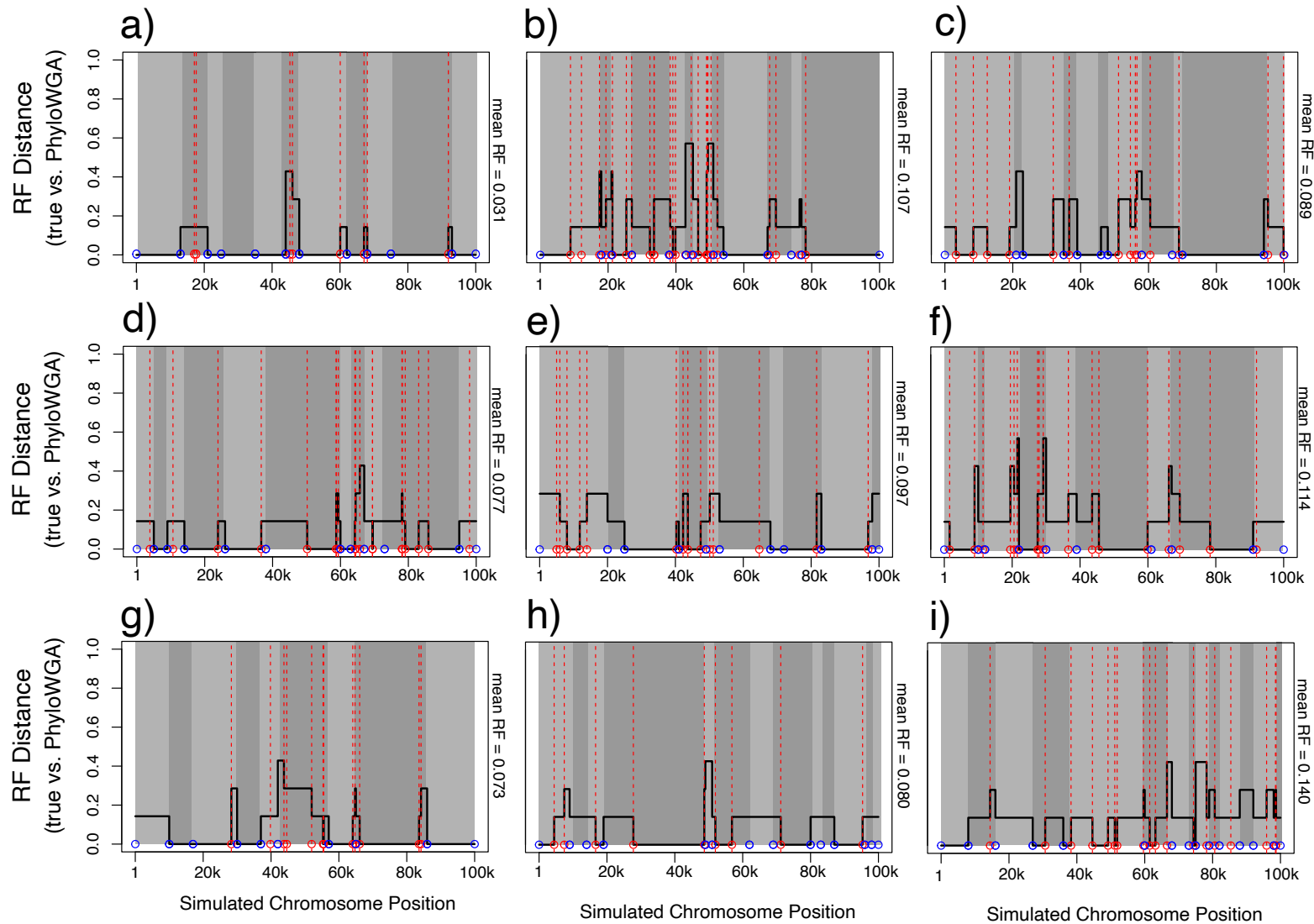
205

206

207

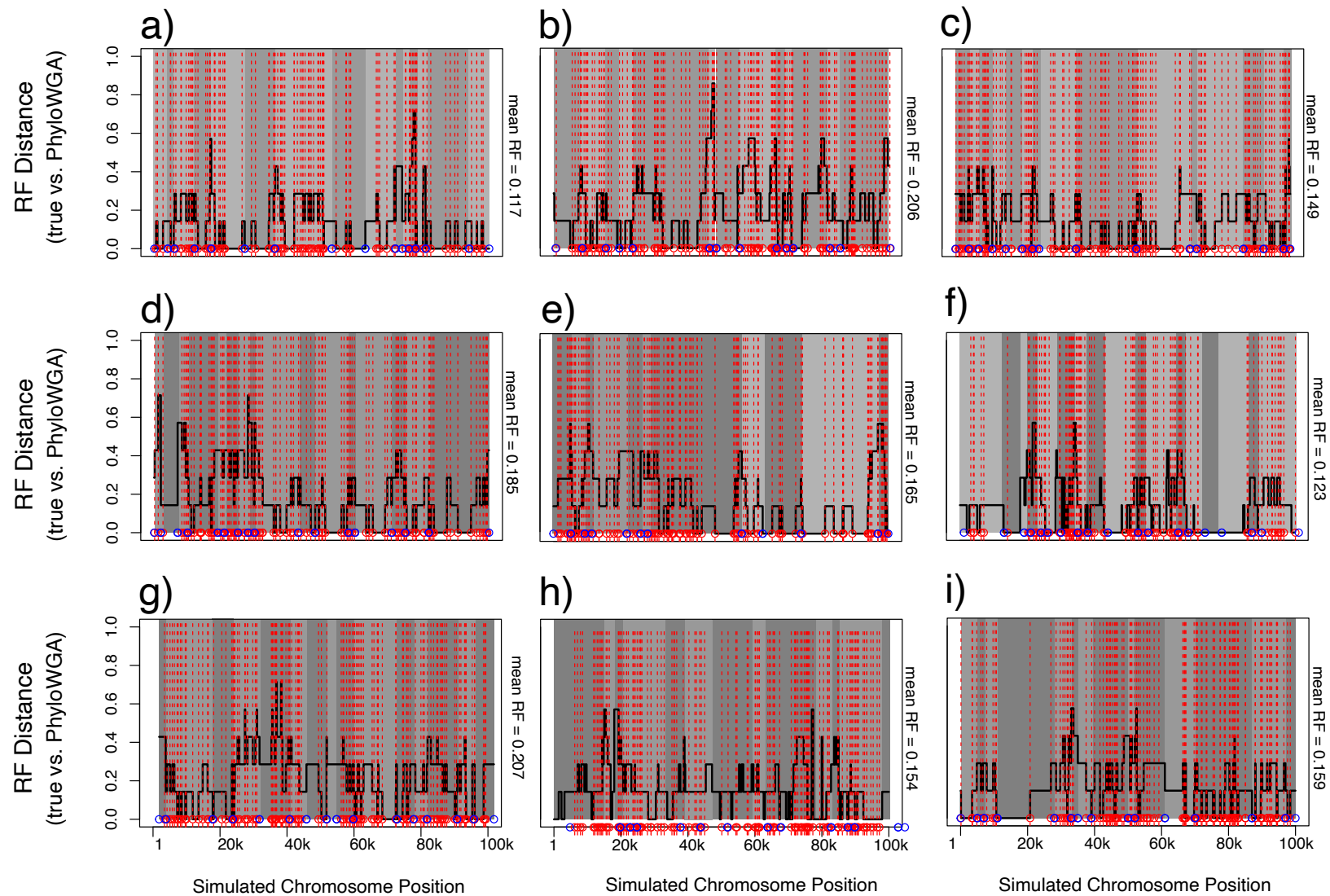
208

Figure S1. Results of simulation analyses for demonstrating the accuracy of *PhyloWGA* on chromosome alignments simulated with a recombination rate of $r = 10^{-9}$ (a), 10^{-8} (b), and 10^{-7} (c) per site per generation. Black lines indicate the mean Robinson-Foulds (RF) distance for each nucleotide site position in the 100 kb simulated alignments (i.e., site means measured across the respective nine replicates shown in Figures S2-S4). Red dashed line indicates the total mean across the entire alignment.



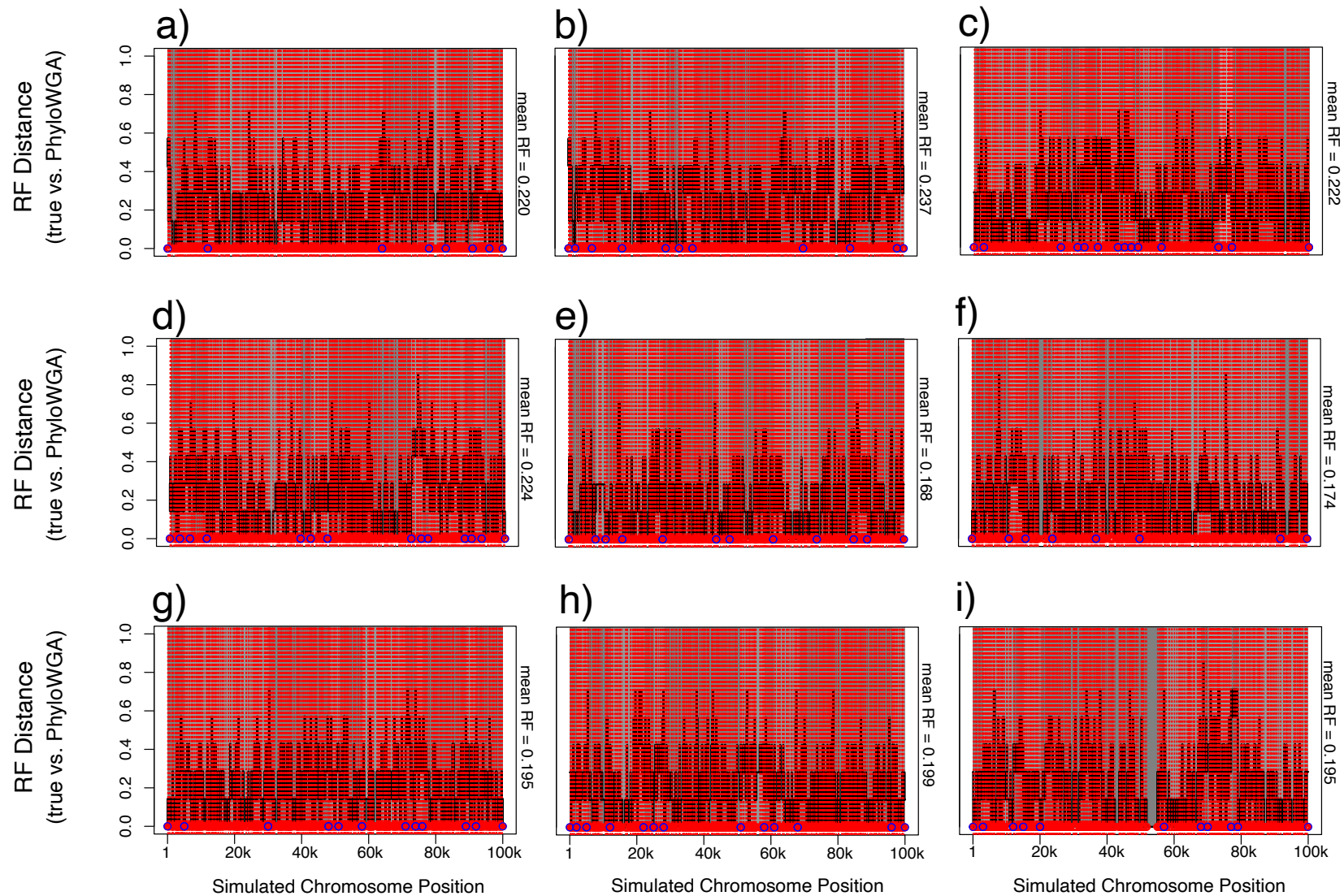
209
210
211
212
213
214

Figure S2. Results of the simulation analyses for demonstrating the accuracy of *PhyloWGA* on chromosome alignments simulated with a recombination rate of $r = 10^{-9}$ per site per generation. Each of nine replicate analyses are shown in each successive panels a-i. Red dashed lines indicate recombination events yielding differences in topology, while alternating light and dark gray blocks indicate concatenated windows from *Chromo.Crawl* with one kb windows and step sizes. Blue circles demarcate the boundaries of concatenated windows. Black lines measure the Robinson-Foulds (RF) distance between the true, simulated tree and the inferred tree from *Chromo.Phylome*.

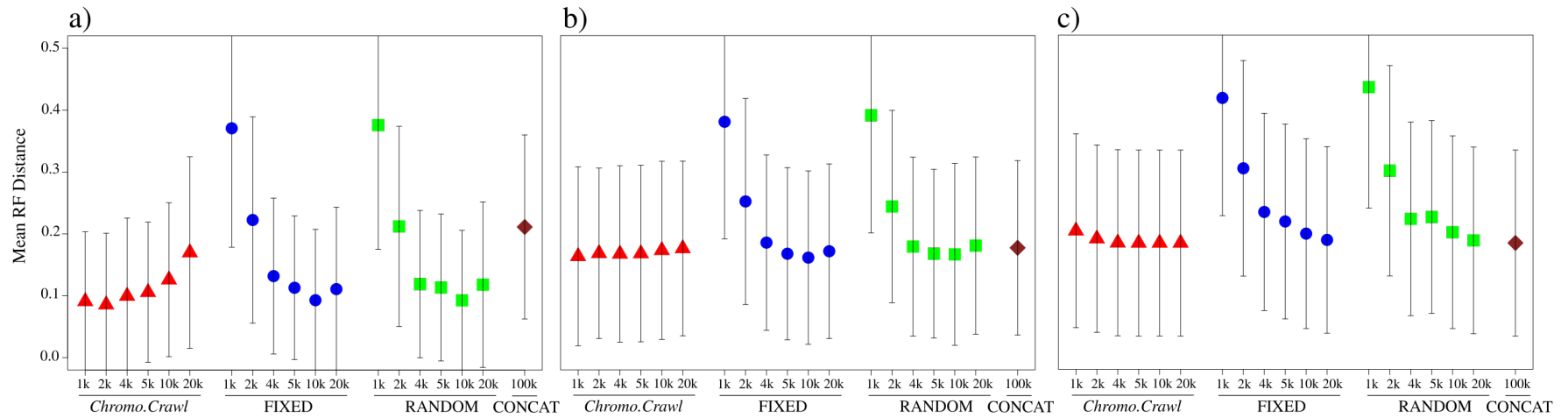


215
216
217
218
219
220

Figure S3. Results of the simulation analyses for demonstrating the accuracy of *PhyloWGA* on chromosome alignments simulated with a recombination rate of $r = 10^{-8}$ per site per generation. Each of nine replicate analyses are shown in each successive panels a-i. Red dashed lines indicate recombination events yielding differences in topology, while alternating light and dark gray blocks indicate concatenated windows from *Chromo.Crawl* with one kb windows and step sizes. Blue circles demarcate the boundaries of concatenated windows. Black lines measure the Robinson-Foulds (RF) distance between the true, simulated tree and the inferred tree from *Chromo.PhyloMe*.



221
 222 **Figure S4.** Results of the simulation analyses for demonstrating the accuracy of *PhyloWGA* on chromosome alignments simulated with a
 223 recombination rate of $r = 10^{-7}$ per site per generation. Each of nine replicate analyses are shown in each successive panels a-i. Red dashed
 224 lines indicate recombination events yielding differences in topology, while alternating light and dark gray blocks indicate concatenated
 225 windows from *Chromo.Crawl* with one kb windows and step sizes. Blue circles demarcate the boundaries of concatenated windows. Black
 226 lines measure the Robinson-Foulds (RF) distance between the true, simulated tree and the inferred tree from *Chromo.Phylome*.



227

228

229

230

231

232

233

234

235

Figure S5. The adaptive window size for estimating gene trees provides *Chromo.Crawl* with improved chromosome-scale phylogenetic accuracy compared to typical approaches that do not approximate optimal window sizes. Phylogenetic accuracy of *Chromo.Crawl* (red triangles) compared with three alternative strategies: the FIXED approach using a fixed sliding window size to infer gene trees (blue circles), the RANDOM approach with trees inferred across randomly sampled genomic regions also of fixed size (green squares), and the CONCAT approach that assumes all windows share the same tree by concatenating the entire 100kb window (brown diamonds). Results shown for mean Robinson-Foulds (RF) distance (bars indicate standard deviation) across nine replicates for the (a) low ($r = 10^{-9}$), (b) medium ($r = 10^{-8}$), and (c) high ($r = 10^{-7}$) recombination rate simulations using one, two, four, five, 10, and 20 kb window sizes, respectively (left to right in each panel).

236