

1 **Model Adequacy Tests for Likelihood Models of**
2 **Chromosome-Number Evolution**

3 Anna Rice and Itay Mayrose[†]

4 School of Plant Sciences and Food Security, Tel Aviv University, Tel Aviv, Israel

5 [†] Corresponding author: Itay Mayrose.

6 Address: School of Plant Sciences and Food Security, George S. Wise Faculty of Life
7 Sciences, Tel Aviv University, Tel Aviv 69978, Israel.

8 Phone: +972-3-640-7212.

9 Email: itaymay@post.tau.ac.il.

10 Total word count of main text: 6,353

11 Word count for: Introduction 1,220; Description 1,926; Results 1,894, Discussion
12 1,313

13 3 Figures (2 in color); 3 Tables.

14 Supporting information: includes supplementary text and five tables.

15 Summary

- 16 • Chromosome number is a central feature of eukaryote genomes. Deciphering
17 patterns of chromosome-number change along a phylogeny is central to the
18 inference of whole genome duplications and ancestral chromosome numbers.
19 ChromEvol is a probabilistic inference tool that allows the evaluation of several
20 models of chromosome-number evolution and their fit to the data. However,
21 fitting a model does not necessarily mean that the model describes the empirical
22 data adequately. This vulnerability may lead to incorrect conclusions when model
23 assumptions are not met by real data.
- 24 • Here, we present a model adequacy test for likelihood models of chromosome-
25 number evolution. The procedure allows to determine whether the model can
26 generate data with similar characteristics as those found in the observed ones.
- 27 • We demonstrate that using inadequate models can lead to inflated errors in several
28 inference tasks. Applying the developed method to 200 angiosperm genera, we
29 find that in many of these, the best-fitted model provides poor fit to the data. The
30 inadequacy rate increases in large clades or in those in which hybridizations are
31 present.
- 32 • The developed model adequacy test can help researchers to identify phylogenies
33 whose underlying evolutionary patterns deviate substantially from current
34 modelling assumptions and should guide future methods developments.

35 **Key words:** chromEvol, chromosome number, dysploidy, model adequacy, model
36 selection, phylogenetics, polyploidy.

37 Introduction

38 Chromosome number is widely recognized as a key feature of eukaryote genomes. Its
39 popularity in cyto-taxonomical and evolutionary studies has been attributed to its
40 ability to provide a concise description of the karyotype, the ease by which it can be
41 recorded, and its stable phenotype across repeated measurements. Processes that lead
42 to changes in chromosome numbers have direct consequences on central evolutionary
43 processes related to reproductive isolation and speciation, thus providing important
44 information for species determination and phylogenetic relationships (Guerra, 2008;
45 Weiss-Schneeweiss & Schneeweiss, 2013). While chromosome numbers generally
46 exhibit strong phylogenetic signal (e.g. Vershinina & Lukhtanov, 2017; Carta *et al.*,
47 2018), they are also highly dynamic. This variability has been particularly well
48 acknowledged in plants, with counts ranging from $n = 2$ to $n = 720$ (Khandelwal,
49 1990; Ruffini Castiglione & Cremonini, 2012), and records showing intraspecific
50 variation in 23% of angiosperm species (Rice *et al.*, 2015). Understanding the
51 underlying processes that gave rise to these changes allows inference of major
52 genomic events that have occurred in the history of a clade of interest and the
53 processes that have shaped its diversification.

54 Of the various mechanisms underlying chromosome-number change, polyploidy, or
55 whole genome duplication (WGD) has received significant attention because of the
56 profound impacts such an event has on the organism. Polyploids often differ markedly
57 from their progenitors in morphological, physiological, or life history characteristics,
58 which may contribute to their establishment in novel ecological settings (Stebbins,
59 1971; Levin, 1983; Ramsey & Schemske, 2002; Soltis *et al.*, 2007; Leitch & Leitch,
60 2008; Ramsey & Ramsey, 2014; Spoelhof *et al.*, 2017; Rice *et al.*, 2019). Polyploidy
61 is thus recognized as one of the major processes that has driven and shaped the
62 evolution of higher organisms. A more subtle change in chromosome number is
63 dysploidy, leading to step-wise changes in the number of chromosomes, but typically
64 does not immediately alter the genomic content. Dysploidy occurs via several types of
65 genome rearrangements, leading to ascending or descending dysploidy through
66 chromosome fission or fusion (Weiss-Schneeweiss & Schneeweiss, 2013).
67 Deciphering the pattern of chromosome-number change within a clade allows
68 inferring the number and type of transitions that have occurred along branches of a

69 phylogeny, to estimate ancestral chromosome numbers, and to categorize extant
70 species as diploids or polyploids.

71 In the last decade, several tools that infer changes in chromosome numbers along a
72 phylogeny were developed (Mayrose *et al.*, 2010; Hallinan & Lindberg, 2011; Glick
73 & Mayrose, 2014; Freyman & Höhna, 2017; Zenil-Ferguson *et al.*, 2017, 2018;
74 Blackmon *et al.*, 2019). Among these, the chromEvol probabilistic framework
75 (Mayrose *et al.*, 2010) was the first to incorporate a continuous time Markov process
76 that describes the instantaneous rate of change from a genome with i haploid
77 chromosomes to a genome with j haploid chromosomes via specific types of
78 dysploidy and polyploidy transitions. Further development of this framework allowed
79 for more intricate types of chromosome-number transitions (Glick & Mayrose, 2014),
80 to differentiate between transitions that coincide with speciation events and those that
81 occur continuously in time along branches of the phylogeny (Freyman & Höhna,
82 2017), and to associate patterns of chromosome-number change with the evolution of
83 a discrete character trait (Zenil-Ferguson *et al.*, 2017; Blackmon *et al.*, 2019).

84 In the chromEvol model, each type of transition is represented by a parameter
85 describing its rate of change. The inclusion (or exclusion) of different parameters
86 entails different hypotheses regarding the pathways by which the evolution of
87 chromosome number proceeded in the clade under study. In a regular application of
88 the chromEvol framework, different models are fitted to the data and the best one is
89 chosen by comparing the relative fit of each model to the data at hand using
90 established model selection criteria, such as the likelihood ratio test or Akaike
91 Information Criterion (AIC; Akaike, 1974). In reality, however, no empirical dataset
92 will meet all the assumptions of any model and thus relying on the best model (or set
93 of models) may be vulnerable to incorrect conclusions in datasets whose underlying
94 evolutionary process deviate substantially from current modelling assumptions. To
95 prevent such errors, here we develop a model adequacy test that allows determining
96 whether a given model of chromosome-number evolution provides a realistic
97 description of the evolutionary process for reliable inferences.

98 Several assumptions made by existing models of chromosome-number evolution may
99 be violated when empirical data are analyzed. For example, all models rely on a
100 memory-less Markovian process, in which the transition rates are only dictated by the

101 current number of chromosomes of the lineage. Thus, for example, the transition rate
102 from $n = 10$ to $n = 9$ is not affected by the duration of time the lineage possessed 10
103 chromosomes, nor by the sequence of events that had led to it. However, because
104 rates of descending dysploidy may increase following WGD (Wood *et al.*, 2009;
105 Wendel, 2015; Soltis *et al.*, 2016), the transition from $n = 10$ to $n = 9$ is more probable
106 if $n = 5$ was the ancestral state compared to $n = 11$. Additionally, most models assume
107 that the transition rates are similar across the phylogeny, although in practice the
108 transition patterns may be rather different in some sub-clades compared to others, as
109 has been demonstrated, for example, in Cyperaceae (Márquez-Corro *et al.*, 2019).
110 Finally, all current models are based on a phylogenetic structure and thus ignore the
111 possibility of hybridizations. Notably, allopolyploidy, one of the main types of
112 polyploidy, is defined by such reticulate evolutionary events and the biases caused by
113 their presence is rather unexplored.

114 One aspect of understanding the reliability of a model and interpreting its results is to
115 quantify its adequacy for the data and the question at hand. The aim of model
116 adequacy tests is to determine the absolute fit of a model to the data, rather than to
117 compare its relative fit among a set of models. With some variations, the general
118 procedure of such tests is composed of several steps: first, given an empirical dataset,
119 obtain the best-fitting model and its parameter values. Next, use that model to
120 generate multiple simulated datasets. Then, compute several test statistics that
121 describe various characteristics of the data on each simulated dataset and on the
122 empirical dataset. If the empirical values of the test statistics fall outside the range of
123 variation encompassed by the simulated data, then it may be concluded that the model
124 cannot provide an adequate description of the data at hand. To date, model adequacy
125 approaches are established for several types of data and inference tasks, including
126 those related to sequence evolution (Bollback, 2002; Brown, 2014; Duchêne *et al.*,
127 2015; Chen *et al.*, 2019) and continuous valued organismal traits (Slater & Pennell,
128 2013; Pennell *et al.*, 2015). However, both are inappropriate for data and analyses
129 concerning the evolution of chromosome numbers as the former rely on statistics
130 derived from many sites, while the latter rely on Brownian motion statistics.

131 In the following, we first provide the details of the developed model adequacy
132 framework for likelihood models of chromosome-number evolution. We then use

133 simulations to assess the type I error rate and to explore the consequences of using
134 inadequate models in several common inference tasks, such as ancestral
135 reconstructions of chromosome numbers and ploidy-level inference. Finally, we apply
136 the developed procedure to a large cohort of angiosperm genera, as well as to clades
137 that are expected to violate model assumptions.

138 **Methodological Description**

139 **Model adequacy framework for chromosome-number evolution**

140 Given chromosome count data and a compatible phylogeny (together denoted as D),
141 chromEvol can be used to assess the fit of various models (M_1, M_2, \dots, M_N ; N denotes
142 for the number of models) to D . Each model differs with respect to the included rate
143 parameters or the constraints placed on them [$\theta(M_1), \theta(M_2), \dots, \theta(M_N)$]. The most
144 general model considered here includes six free parameters (Glick & Mayrose, 2014)
145 and assumes that five types of events are possible: a single chromosome-number
146 increase (ascending dysploidy with rate λ) or decrease (descending dysploidy with
147 rate δ), WGD (i.e. exact duplication of the number of chromosomes with rate ρ),
148 demi-polyploidy (multiplications of the number of chromosomes by 1.5 with rate μ),
149 and base-number transitions (the addition to the genome by any multiplication of an
150 inferred base number, where β , is the inferred base number and ν is its respective
151 transition rate). A combination of these parameters allows a range of models to be
152 evaluated (Table 1 shows the various models considered here). We note that the
153 chromEvol software also allows the ascending and descending dysploidy rates to
154 depend on the current number of chromosomes, but this option was not evaluated
155 here.

156 In a common application of chromEvol, several models are fitted to D , the optimal
157 model is selected based on its relative fit using established model selection criteria
158 (e.g. AIC), and subsequent inference tasks are performed based on this model. The
159 model adequacy test can be carried out to any model of interest, whether or not it is
160 the most fitted one. The general aim of this test is to examine whether a specified
161 model, M_x , is able to generate data that are similar to D . Our model adequacy
162 procedure is based on parametric bootstrapping (Goldman, 1993; Efron & Tibshirani,
163 1994), where the observed data are compared to a background distribution generated

164 from simulations. These simulations are generated under the specified model, whose
165 parameters, $\hat{\theta}(M_x)$, were optimized with respect to D and the respective probabilities
166 of chromosome-numbers inferred at the root of the phylogeny (exact details of the
167 simulation procedure are given in the Supporting Information). Comparing between
168 true and simulated data is performed using a set of test statistics, which reflects
169 various characteristics of the data. First, the test statistics (T_1, T_2, \dots, T_m ; m denotes for
170 the number of statistics) are computed for the true data D . Second, multiple datasets
171 are simulated under the specified model and its inferred parameters. For each
172 simulated dataset, the same set of test statistics is computed, resulting in a distribution
173 for each test statistic ($T_{s1}, T_{s2}, \dots, T_{sm}$). If the empirical value of the test statistic falls
174 within the range of variation encompassed by the simulated data (herein defined as
175 the 2.5th and 97.5th percentiles), the model is considered as capable of generating data
176 similar to the original ones and is thus inferred as adequate. Otherwise, it is inferred
177 as inadequate. A schematic illustration of the developed model adequacy test is
178 presented in Fig. 1.

179 In our implementation, four test statistics were calculated given the chromosome-
180 number data of extant taxa and the corresponding phylogeny: (1) Variance; higher
181 values in the simulated data relative to the observed ones may point to some
182 constraints that were not accounted for by the model (e.g. hard bounds on the number
183 of chromosomes in the genome), or to errors in the parameter estimation process. (2)
184 Shannon's entropy (Shannon, 1948); Lower entropy of the observed data than
185 predicted by the model is indicative of higher-than-expected concentration of
186 genomes with certain haploid numbers. This could be due to selective constraints, or
187 to a very low variability exhibited in certain subclades of the phylogeny, such that
188 specific states are clumped into large blocks of the tree more than expected. (3)
189 Parsimony score; the most parsimonious number of character transitions across the
190 phylogeny is calculated based on Fitch (1971). If the parsimony score of the observed
191 data are lower than expected it means that the model assumes more transitions than
192 actually occurred. This could occur due to rate heterogeneity across the tree. For
193 example, if chromosome-number transitions occur more frequently in one subclade
194 relative to the rest of the phylogeny, this could be accommodated by inferring higher
195 values of the transition rates. (4) Parsimony versus time ($\text{Pars}^{\text{Time}}$); the parsimonious
196 number of transitions are computed per branch using the accelerated transformation

197 criterion (ACCTAN; Farris, 1970). The regression line between the divergence
198 times (computed from the root to the end of the branch) and their parsimony scores is
199 calculated, and the slope of this line is taken as the test statistic. This statistic is
200 similar in spirit to that employed by Pennell *et al.* (2015) for testing the adequacy of
201 models for continuous trait evolution. Under a time-homogenous model, as
202 implemented in chromEvol, we expect no relationship between the divergence times
203 and the number of transitions. Violations of this assumption suggests that transitions
204 are either concentrated around the root or occur more frequently towards the tips. We
205 note that aside from these four statistics, two additional ones were computed (the
206 range and the number of unique counts). These two statistics were found to be highly
207 correlated with the other test statistics ($r^2 = 0.85$ between range and variance and $r^2 =$
208 0.74 between unique counts and entropy, when computed over the 200 empirical
209 datasets; detailed below), and thus we chose to discard them from further analyses.
210 The coefficient of determinations between all pairs of the four remaining test statistics
211 was below 0.40 (Supporting Information Table S1). Because the four test statistics are
212 not independent and researchers might be interested in revealing the specific aspects
213 of the data that differ from expectations, we followed Pennell *et al.* (2015) and did not
214 apply a multiple testing correction. Thus, in all analyses presented here a model is
215 considered as adequate only if all four statistics fall within the boundaries of the
216 simulated distribution.

217 **Performance assessment using simulations**

218 Simulations were conducted to examine the performance of the model adequacy
219 procedure. Given an input phylogeny and a set of model parameters, simulated
220 chromosome numbers were generated as previously described in Mayrose *et al.*
221 (2010). As the number of simulation conditions is infinite, we concentrated on eight
222 scenarios that vary in terms of data size (the number of tips in the phylogeny and the
223 observed chromosome-number distribution) and the inferred pattern of chromosome-
224 number change (Table 2). The phylogenies, chromosome counts, and model
225 parameters were taken from empirical datasets previously analyzed using chromEvol
226 (Glick *et al.*, 2016; Rice *et al.*, 2019), thus representing realistic data characteristics.
227 For each simulation scenario, a total of 100 replicates were generated. Each simulated
228 dataset was then fitted to a set of four models: D_{ys} , $D_{ys}D_{up}$, $D_{ys}B_{num}$, $D_{ys}D_{up}B_{num}$. For

229 each simulation scenario, one of these models was the generating model (i.e. the
230 model that was used to simulate the data) and three were non-generating models. We
231 note that these models share both common and distinct aspects of the parameter space,
232 such that some, but not all, models are nested within each other (Table 1). Finally, the
233 adequacy of each model to the simulated data was assessed.

234 **Inference errors of adequate and inadequate models**

235 The consequences of using an adequate versus inadequate model were evaluated by
236 comparing the errors of four common inference tasks: (1) the chromosome number at
237 the root of the phylogeny; calculated as the deviation from 1.0 of the posterior
238 probability assigned to the true (i.e. simulated) chromosome number at the root. (2)
239 The total number of dysploidy events across the phylogeny; calculated as the relative
240 error between the inferred and simulated number of events: $2 * \frac{|x_1 - x_2|}{x_1 + x_2}$, where x_1 and
241 x_2 are the simulated and inferred number of dysploidy events, respectively. In case
242 both x_1 and x_2 equal zero the error was assigned as zero. (3) The total number of
243 polyploidization events across the phylogeny; the relative error was calculated similar
244 to the total number of dysploidy events. Duplication events, demi-duplications, and
245 base-number transitions were regarded as polyploidization events. (4) Ploidy level
246 assignments. The ploidy-level inference of tip taxa, as either diploids or polyploids,
247 was based on the procedure described in Glick & Mayrose (2014). The assignments of
248 all tips were compared between the inferred and true values. The number of falsely
249 inferred taxa, divided by the total number of taxa, was used as the error measure.

250 In this analysis, six of the eight simulation scenarios in Table 2 were examined. The
251 two scenarios excluded were those generated under the simple D_{ys} model, for which
252 not all inference tasks are relevant. To eliminate possible confounding effects between
253 the specific model used for inference and the magnitude of the error, in this evaluation
254 a single non-generating model ($D_{ys}D_{up}$ or $D_{ys}B_{num}$) was fitted to the data per
255 simulation scenario (Supporting Information Table S2). For each simulation scenario,
256 300 replicates were generated. For each replicate, the phylogeny and the simulated
257 chromosome counts were given as input to the model adequacy test and the dataset
258 was determined as either adequate or inadequate. A one-sided t -test was conducted to

259 determine whether the error of a certain inference task is significantly larger in the
260 inadequate set compared to the adequate set.

261 **Application to empirical datasets**

262 To demonstrate the usability of the model adequacy framework, we applied it to a
263 dataset of 200 angiosperm genera, which were randomly selected from a large
264 database consisting of thousands of plant genera, excluding genera with no variations
265 in chromosome numbers as well as those with less than 5 species with both
266 phylogenetic and chromosome-numbers information. The initial database was used, in
267 part or as a whole, in several previous analyses (e.g. Glick *et al.*, 2016; Zhan *et al.*,
268 2016; Salman-Minkov *et al.*, 2016; Rice *et al.*, 2019). From this database we also
269 selected 40 angiosperm genera that each contains at least one allopolyploid species,
270 based on data from Barker *et al.* (2016). Due to overlaps between these two sets, a
271 total of 233 unique datasets were analyzed. Full details of the reconstruction of the
272 original database are described in Rice *et al.* (2019). Briefly, for each genus, the
273 OneTwoTree pipeline (Drori *et al.*, 2018) was used to automatically reconstruct the
274 phylogeny using publicly available sequence data as appear in GenBank (Benson *et*
275 *al.*, 2013). Chromosome numbers for all species were retrieved from the Chromosome
276 Counts Database (CCDB; Rice *et al.*, 2015). These data were given as input to
277 chromEvol, which was executed on the six models detailed in Table 1. Additionally,
278 we applied similar procedures to seven clades of higher taxonomical ranks, including
279 five families, one subfamily, and one tribe. The evolution of chromosome numbers in
280 these clades using chromEvol was previously examined in several studies (Supporting
281 Information Table S3).

282 **Implementation and availability**

283 The model adequacy procedure was implemented in Python and R (R Core Team,
284 2013). The source codes and running instructions are available at
285 https://github.com/MayroseLab/chromEvol_model_adequacy. The obligatory inputs
286 are three files obtained through a chromEvol run of the examined model: the
287 summary results file, the tree with the inferred ancestral reconstruction in a NEWICK
288 format, and the original counts file in FASTA format. The program outputs, for each
289 test statistic examined, its value computed from the empirical data, the percentile in

290 which it falls within the simulated distribution, and the 2.5th and 97.5th percentiles of
291 the simulated distribution as indicative for the upper and lower bounds expected under
292 the modelling assumptions. The model adequacy test is also available for on-line use
293 through the chromEvol web-server (<http://chromevol.tau.ac.il/>), which is currently in
294 a Beta version.

295 **Results**

296 In this work we developed a statistical framework for testing the adequacy of
297 likelihood models of chromosome-number evolution. In essence, the method tests
298 whether a specified model is capable of generating data that are similar to the data at
299 hand. If not, the model is considered as providing inadequate description of the data,
300 suggesting that other processes than those modeled have driven the evolution of
301 chromosome numbers along the examined phylogeny. We first evaluated the
302 performance of the model adequacy framework using simulations. We then applied it
303 to a large number of real datasets derived from dozens of angiosperm genera, as well
304 as to seven clades of higher taxonomic ranks, that together vary greatly in their extent
305 of divergence time and patterns of chromosome number variation.

306 **Framework validation**

307 Simulations were used to validate the developed model adequacy approach. Several
308 simulation scenarios were examined, whose phylogenies and simulated parameters
309 were derived from real data analyses and cover various data characteristics (Table 2).
310 In each scenario, a single model was used to generate the data. Given the simulated
311 data, the generating model and three additional models were fitted to the data, and
312 their adequacies were examined. The four examined models are indicated by the type
313 of transitions they allow for: D_{ys} , $D_{ys}D_{up}$, $D_{ys}B_{num}$, and $D_{ys}D_{up}B_{num}$ (Table 1). In total,
314 eight different simulation scenarios were examined; two for each type of generating
315 model.

316 We first examined the type I error rate, i.e. inferring the generating model as
317 inadequate. Our results indicated that when considering a single test statistic
318 independently, the error rate is around the expected value of 0.05 (average = 0.02,
319 across the eight simulation scenarios and four test statistics; Supporting Information

320 Table S4). Combining multiple test statistics together, we consider a model as
321 inadequate if one or more of the statistics fell outside the margins of the simulated
322 distributions (see Methodological Description). Under this definition, the percentage
323 of generating models that were inferred as inadequate varied between 0.04 and 0.13
324 across the eight simulation scenarios (Table 3). When Bonferroni correction for
325 multiple testing was applied, the type I error rate dropped to an average of 0.008. We
326 note however, that the four test statistics are not independent, violating the assumption
327 of this correction.

328 We next examined the capability of the adequacy test to detect models that deviate
329 from that of the generating models. Three types of model misspecification were
330 examined: over-parameterization, under-parameterization, and miss-parameterization.
331 In the case of over-parameterization, the tested model allows for additional types of
332 chromosome-number change (as represented by extra free parameters) than those used
333 to generate the data. This corresponds to cases where the generating model is nested
334 within the tested model (e.g. the generating model is $D_{ys}D_{up}$ while the tested model is
335 $D_{ys}D_{up}B_{num}$). Our results indicated that the performances of over-parameterized
336 models are very similar to that of the generating models (Table 3). The few
337 discrepancies were the result of either (1) inaccurate parameter estimates of the more
338 general model due to the extra degrees of freedom; (2) the optimization procedure
339 reaching suboptimal regions of the parameter space (we note that while chromEvol
340 allow for more thorough likelihood optimization search, which should reduce such
341 instances, this was not attempted here due to the large number of simulations
342 employed); (3) very similar parameter estimates obtained using the two models, but
343 slight deviations of the test statistics led one model to be inferred as inadequate while
344 the other one as adequate.

345 In the case of under-parameterized models, the tested model allows for fewer types of
346 transitions than the generating model (e.g. the generating model is $D_{ys}D_{up}$ while the
347 tested model is D_{ys}). As may be expected, in all simulation scenarios the under-
348 parameterized models were more frequently inferred as inadequate compared to the
349 generating models. The adequacy rate was very low when the tested model allowed
350 only for dysploid transitions while in reality polyploid transitions (either WGD and/or
351 base-number transitions) have occurred (Table 3; all cases where the tested model is
352 D_{ys}). The adequacy rates were higher when the generating model allowed for multiple

353 types of polyploid transitions (i.e. $D_{ys}D_{up}B_{num}$ allowing for both exact duplications
354 and base-number transitions), while the tested model allowed for a subset of these
355 ($D_{ys}D_{up}$ and $D_{ys}B_{num}$ that allow only for duplications or base-number transitions,
356 respectively). Comparing the adequacy of the two under-parametrized models ($D_{ys}D_{up}$
357 and $D_{ys}B_{num}$), the $D_{ys}B_{num}$ model that incorporated base-number transitions had higher
358 adequacy rates compared to the $D_{ys}D_{up}$ model that allowed for exact duplications, as
359 the former allows for several transitions that frequently include also exact
360 duplications (e.g. in case the base number is 8, both $8 \rightarrow 16$ and $8 \rightarrow 24$ transitions are
361 allowed).

362 In the case of miss-parametrization, the tested and generated models are not nested
363 within each other and thus their parameters only partially overlap. For the set of
364 models examined here, this fits the case where the generating model is $D_{ys}D_{up}$ while
365 the tested model is $D_{ys}B_{num}$, or vice versa. When the tested model was $D_{ys}B_{num}$, it
366 obtained similar adequacy rates to those of the generating $D_{ys}D_{up}$ model. In contrast,
367 and similar to the results detailed in the case of under-parameterized models, the
368 $D_{ys}D_{up}$ model was inferred as inadequate a large number of times when the generating
369 model was $D_{ys}B_{num}$.

370 **Inference errors of adequate and inadequate models**

371 A central usage of probabilistic models of chromosome number evolution is their
372 inference capabilities, such as ancestral reconstructions of chromosome numbers, or
373 predicting the branches in which dysploidy and polyploidy events have most likely
374 occurred. Still, it is unclear whether the use of inadequate models would deteriorate
375 the performance of such inference tasks. To this end, simulations were used to
376 compare the errors of the following four common inference tasks when adequate and
377 inadequate models are employed: (1) the chromosome number at the root of the
378 phylogeny; (2) the total number of inferred dysploidy events; (3) the total number of
379 inferred polyploidization events, and (4) inferring the ploidy level of tip taxa as either
380 diploid or polyploidy (see Methodological Description for details regarding the error
381 computed for each inference task).

382 Our results demonstrated that the use of inadequate models frequently leads to larger
383 inference errors, although under some simulation scenarios the inference errors of

384 inadequate models were similar to that obtained using adequate models. For example,
385 the error in the inference of the root chromosome number was significantly larger in
386 the case of inadequate models under two simulation scenarios, but was non-
387 significantly different in the other four (Fig. 2). Similarly, in two out of the six
388 simulation scenarios, the error of inferring the ploidy level of extant taxa was
389 significantly larger when computed using inadequate versus adequate models. In this
390 case, the magnitude of the error was relatively low whether adequate or inadequate
391 models were applied: when inadequate models were applied, the mean error was 4.6%
392 across all simulation scenarios, reaching up to 12% under the *Brassica* simulation
393 scenario. In comparison, the mean error was 2% when adequate models were applied,
394 reaching up to 6% of erroneous inferences under the *Hordeum* simulation scenario.
395 Larger differences in the errors between adequate and inadequate models were
396 observed in inferring the total number of polyploidizations, and even more so in
397 inferring the total number of dysploidy events. For both these inference tasks,
398 significant differences between adequate and inadequate models were obtained for
399 three out of the six simulation scenarios. Generally, the relative error in inferring the
400 total number of dysploidy events was larger compared to that of inferring the total
401 number of polyploidizations (the mean relative error was roughly twice for dysploidy
402 compared to polyploidy transitions, both in the adequate set and the inadequate set;
403 Fig. 2).

404 **Application to empirical datasets**

405 We applied the model adequacy framework to 200 datasets, each corresponding to a
406 single randomly-selected angiosperm genus. First, we performed a standard model
407 selection procedure based on the AIC (Akaike, 1974) to evaluate the relative fit of
408 each of the six chromEvol models to the data. In 24% of the datasets, the simple D_{ys}
409 model, which allows for dysploid transitions only, was selected. The model that was
410 most frequently selected was $D_{ys}D_{up}$ (28%), while models that allow for demi-
411 polyploidy transitions and those that allow for base-number transitions were selected
412 in 27% and 21% of the datasets, respectively (Fig. 3a). Next, we applied the model
413 adequacy test to the best model identified for each dataset. We found that in 74% of
414 the genera, the model that was chosen as best by the AIC was inferred to provide an
415 adequate description of the data. Applying the model adequacy test to all six models

416 per dataset (whether or not selected as best), we found that models that allow for
417 fewer types of transitions were more frequently predicted as inadequate (Fig. **3b**). For
418 example, the D_{ys} model that allows only for dysploidy transitions was adequate in
419 only 28% of the 200 datasets, models that additionally allow for one type of
420 polyploidy, either duplication or base-number transition, were adequate 60% and 64%
421 of the cases, respectively, while the three models that incorporate two types of
422 polyploidy transitions ($D_{ys}D_{up}D_{em}$, $D_{ys}D_{up}D_{em}^*$, and $D_{ys}D_{up}B_{num}$) were inferred as
423 adequate most frequently. The adequacy rates of all models were generally related to
424 the complexity of the model that was selected as optimal. Thus, when the most
425 complex models were selected ($D_{ys}D_{up}D_{em}$ and $D_{ys}D_{up}B_{num}$), the adequacy rates of all
426 models – including that of the chosen model – were low (33% and 47%,
427 respectively), while when the least complex model was selected, the adequacy rates of
428 all models was high (70%; Supporting Information Table S5).

429 Next, we examined the model adequacy procedure in groups that have evolved via
430 reticulate evolution at some point in their histories. In these clades, the underlying
431 assumption of the chromEvol framework, in which evolution proceeds along a
432 phylogenetic structure, is violated, at least to some extent. This analysis was
433 performed on 40 genera that were identified in the literature to include allopolyploid
434 species, and thus hybridizations were reported to occur (data taken from Barker *et al.*,
435 2016). In the majority of these genera (24 out of 40), the model that was selected as
436 optimal according to the AIC was found by our model adequacy procedure as
437 inadequate. This adequacy rate is significantly lower ($p \ll 0.05$; χ^2 test) compared to
438 a random set of 193 genera in which allopolyploidy was not reported (the 200 genera
439 analyzed above, omitting seven that include a reported allopolyploid species).

440 Finally, we evaluated the model adequacy procedure on a set of seven groups whose
441 taxonomic rank is higher than the genus level, thus representing clades whose
442 divergence time is generally older than those inspected above. The evolution of
443 chromosome numbers in these clades likely violates the time homogeneity assumption
444 of chromEvol, in which the transition pattern is similar across the phylogeny. For four
445 of these seven clades, the model that was chosen as optimal according to AIC did not
446 provide adequate description of the data according to the model adequacy test
447 (Supporting Information Table S3) and in one additional case the empirical values of

448 two test statistics were placed close to the lower boundaries of the simulated
449 distributions (falling in the 0.027 and 0.043 percentiles). Taken together, the last two
450 analyses indicate that the model adequacy procedure can identify cases in which the
451 evolution of chromosome numbers is driven by processes that deviate from the basic
452 modelling assumptions of the chromEvol framework.

453 Discussion

454 For over a century, the determination of chromosome numbers has played a vital role
455 in studying evolutionary and genomic processes in plants. Probabilistic models of
456 chromosome-number change are a relatively recent addition to the research toolbox
457 available to study the evolution of major genomic processes. As the usage of such
458 models increases, so does the need to assess their validity when applied to real data.
459 Here, we developed a model adequacy test for likelihood models of chromosome-
460 number evolution. We focused our analysis on those models implemented in the
461 chromEvol software (Glick & Mayrose, 2014), but the procedures are general and can
462 be implemented in other platforms that use variations to the chromEvol model
463 (Freyman & Höhna, 2017; Zenil-Ferguson *et al.*, 2017; Blackmon *et al.*, 2019). The
464 developed test is based on the parametric bootstrapping approach (Goldman, 1993;
465 Efron & Tibshirani, 1994) in which observed data are compared to a simulated
466 distribution generated by the examined model. Using multiple test statistics that
467 describe various characteristics of the data, the test allows to determine whether the
468 model can generate data that are similar to those found in the observed ones.

469 Our simulation results indicate that the model adequacy framework has an acceptable
470 type I error rate (i.e. inferring as inadequate a model that was used to generate the
471 data). However, higher type I errors were found in models that allow for base-number
472 transitions ($D_{ys}B_{num}$ and $D_{ys}B_{num}D_{up}$). This suggests that these models might not be
473 appropriate in all cases. The current implementation of such models assumes the same
474 rate for all possible base-number transitions (e.g. given a base number of $\beta = 7$, the
475 additions of 7, 14, or 21 chromosomes are equally likely). Alternatively, it may be
476 more appropriate to place a probability distribution over the possible base-number
477 transitions. This will allow, for the example of $\beta = 7$, higher rates for additions by 7
478 chromosomes compared to those by 21.

479 Our simulation results also demonstrate that the adequacy rate of over-parameterized
480 models, which allows for more types of transitions than those that truly occurred, is
481 similar to that of the generating models. While it is expected that the accuracy of
482 inferring the model parameters will decrease as overly-complexed models are
483 evaluated, in many cases the auxiliary parameters were optimized to very low values,
484 resulting in a process that is nearly identical to the generating model. Thus, it seems
485 that the flexibility offered by complex models does not necessarily lead to their
486 disadvantage, at least for some inference tasks, as has been recently demonstrated for
487 models of nucleotide sequence evolution (Abadi *et al.*, 2019). In other cases of model
488 violations, either for under-parameterized or miss-parameterized models, when the
489 rate parameters deviated substantially from the original ones (e.g. dysploidy rates an
490 order of magnitude larger than the simulated rates), the model adequacy framework
491 detected such cases as inadequate. This suggests that the adequacy test is capable of
492 detecting models that are completely wrong. In other cases, the nature of model
493 misspecification affected the outcome. In the simulations examined here, $D_{ys}B_{num}$ was
494 more frequently adequate than $D_{ys}D_{up}$, both in the case of under-parameterization (i.e.
495 when the generating model was $D_{ys}D_{up}B_{num}$ such that both models miss one type of
496 transition) and miss-parameterization. Nevertheless, we note that the $D_{ys}B_{num}$ model
497 may not fit well in large phylogenies with high dysploidy rates. In its current
498 implementation, the model assumes that a single base number typifies a clade.
499 However, if there is a high dysploidy rate, each subclade of the phylogeny may be
500 characterized by its own base number or by multiple base numbers, which will
501 necessitate more complex modelling options.

502 We further tested the consequences of using an inadequate model by examining the
503 errors of several inference tasks. First, we found that the difference in inference error
504 between adequate and inadequate models depended on the simulation scenarios: in
505 some simulation scenarios the use of inadequate models resulted in significantly
506 inflated inference errors compared to the use of adequate models, in some scenarios it
507 affected only certain inference tasks and not others, while in others the difference was
508 negligible for all tasks. Second, we found that some inference tasks are much more
509 sensitive to model misspecification than others. The assignment of extant taxa as
510 diploids or polyploids was the inference task that was least affected from using an
511 inadequate model, and in general, the error of this inference task was very low (in all

512 scenarios, the ploidy level of 88% or more of the taxa were correctly identified). This
513 indicates that determining the ploidy levels of extant taxa is generally robust to model
514 misspecification. On the other hand, the error of determining the number of events
515 that had occurred – either dysploid or polyploid transitions – can be substantial when
516 inadequate models are employed.

517 Applying the model adequacy test to hundreds of angiosperm genera, we found that in
518 the majority of the cases the best-fitted model provided sufficient approximation to
519 the evolutionary processes underlying the data and was determined as adequate.
520 However, in roughly one fourth of the examined genera, this selection turned out to be
521 inadequate, suggesting that there is ample room for future modelling improvements.
522 Indeed, we found high rates of model inadequacy when applying the developed
523 procedures to two types of clades that are expected to violate basic modelling
524 assumptions: first, clades in which allopolyploidy events are known to occur, thus
525 violating the assumption that evolution proceeds via a phylogenetic structure; second,
526 in the case of large and diverse clades in which a single transition process is fitted to
527 the entire phylogeny, following the time homogeneity assumption, is insufficient.
528 These results thus indicate that promising future developments would be to focus on
529 analytical procedures based on phylogenetic networks (Nakhleh, 2010), rather than on
530 bifurcating phylogenies, and to further incorporate time-heterogeneous processes.

531 Phylogenetic model adequacy tests have been previously developed for other data
532 types and inference tasks, although their use has not been widely adopted. This could
533 be due to the apparent limited benefit offered to a researcher when all examined
534 models are deemed inadequate when applied to a clade of interest. We argue,
535 however, that model adequacy tests are of practical use to methods developers and
536 end users alike, and should thus be regularly practiced as part of a broader model
537 assessment routine. For researchers interested in data analysis, inadequate models can
538 hint on errors in the input data, which should thus be more carefully inspected. In the
539 case studied here, possible sources of errors include those in the assumed
540 phylogenetic hypothesis, in the collection of chromosome counts, or in taxa sampling.
541 Inadequacy could also point to additional attributes that should be considered in the
542 analysis. For example, if all models that assume a time-homogenous transition
543 process fail, it could suggest that patterns of chromosome-number change are
544 dependent on an organismal trait (e.g. the plant growth form), that if accounted for,

545 using more complex models (e.g. Zenil-Ferguson *et al.*, 2017; Blackmon *et al.*, 2019)
546 would enhance the analysis. For researchers interested in large scale analyses that
547 include multiple datasets, where the in-depth examination of each inadequate dataset
548 is not feasible, the filtration of such clades is one obvious possible direction. For some
549 inference tasks, such as the identification of ploidy levels of extant taxa, the effect of
550 using an inadequate model is rather negligible, indicating that the treatment of the
551 flagged clades should be tuned to the analysis in question. For the developers, the
552 frequent application of model adequacy tests should provide interesting test cases on
553 which new models are trained. Moreover, when a model is deemed inadequate, the
554 test statistics that fail to align may point to processes absent from existing models,
555 which could be included in the future. Model adequacy should thus take a vital part in
556 this recurrent chain of scientific progress in which new methods are developed,
557 regularly used, and then replaced by more advanced alternatives.

558 **Acknowledgements**

559 A.R. is supported by a fellowship from the Edmond J. Safra Center for Bioinformatics
560 at Tel Aviv University and by the NA'AMAT Professional Scholarship. This study
561 was supported by grant # 961/17 from the Israel Science Foundation to I.M.

562 **Author Contribution**

563 IM and AR conceived the study; AR built the tool and analyzed the data; AR and IM
564 wrote the manuscript; IM supervised the study.

565 **References**

- 566 **Abadi S, Azouri D, Pupko T, Mayrose I. 2019.** Model selection may not be a
567 mandatory step for phylogeny reconstruction. *Nature communications* **10**: 1–11.
- 568 **Akaike H. 1974.** A new look at the statistical model identification. *IEEE transactions*
569 *on automatic control* **19**: 716–723.
- 570 **Barker MS, Arrigo N, Baniaga AE, Li Z, Levin DA. 2016.** On the relative
571 abundance of autopolyploids and allopolyploids. *New Phytologist* **210**: 391–398.
- 572 **Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J,**
573 **Sayers EW. 2013.** GenBank. *Nucleic acids research* **41**: 36–42.
- 574 **Blackmon H, Justison J, Mayrose I, Goldberg EE. 2019.** Meiotic drive shapes rates
575 of karyotype evolution in mammals. *Evolution* **73**: 511–523.
- 576 **Bollback JP. 2002.** Bayesian model adequacy and choice in phylogenetics.
577 *Molecular Biology and Evolution* **19**: 1171–1180.
- 578 **Brown JM. 2014.** Detection of implausible phylogenetic inferences using posterior
579 predictive assessment of model fit. *Systematic Biology* **63**: 334–348.
- 580 **Carta A, Bedini G, Peruzzi L. 2018.** Unscrambling phylogenetic effects and
581 ecological determinants of chromosome number in major angiosperm clades.
582 *Scientific Reports* **8**: 1–14.
- 583 **Chen W, Kenney T, Bielawski J, Gu H. 2019.** Testing adequacy for DNA
584 substitution models. *BMC Bioinformatics* **20**: 349.
- 585 **Drori M, Rice A, Einhorn M, Chay O, Glick L, Mayrose I. 2018.** OneTwoTree:
586 An online tool for phylogeny reconstruction. *Molecular Ecology Resources* **18**: 1492–
587 1499.
- 588 **Duchêne DA, Duchêne S, Holmes EC, Ho SYW. 2015.** Evaluating the Adequacy of
589 Molecular Clock Models Using Posterior Predictive Simulations. *Molecular Biology*
590 *and Evolution* **32**: 2986–2995.
- 591 **Efron B, Tibshirani RJ. 1994.** *An introduction to the bootstrap*. CRC press.

- 592 **Farris JS. 1970.** Methods for computing Wagner trees. *Systematic Biology* **19**: 83–
593 92.
- 594 **Fitch WM. 1971.** Toward Defining the Course of Evolution: Minimum Change for a
595 Specific Tree Topology. *Systematic Zoology* **20**: 406.
- 596 **Freyman WA, Höhna S. 2017.** Cladogenetic and anagenetic models of chromosome
597 number evolution: a Bayesian model averaging approach. *Systematic biology* **67**:
598 195–215.
- 599 **Glick L, Mayrose I. 2014.** ChromEvol: Assessing the Pattern of Chromosome
600 Number Evolution and the Inference of Polyploidy along a Phylogeny. *Molecular*
601 *biology and evolution* **31**: 1914–1922.
- 602 **Glick L, Sabath N, Ashman T-L, Goldberg E, Mayrose I. 2016.** Polyploidy and
603 sexual system in angiosperms: Is there an association? *American Journal of Botany*
604 **103**: 1223–1235.
- 605 **Goldman N. 1993.** Statistical tests of models of DNA substitution. *Journal of*
606 *Molecular Evolution* **36**: 182–198.
- 607 **Guerra M. 2008.** Chromosome numbers in plant cytotaxonomy: concepts and
608 implications. *Cytogenetic and Genome Research* **120**: 339–350.
- 609 **Hallinan NM, Lindberg DR. 2011.** Comparative Analysis of Chromosome Counts
610 Infers Three Paleopolyploidies in the Mollusca. *Genome Biology and Evolution* **3**:
611 1150–1163.
- 612 **Khandelwal S. 1990.** Chromosome evolution in the genus *Ophioglossum* L.
613 *Botanical Journal of the Linnean Society* **102**: 205–217.
- 614 **Leitch AR, Leitch IJ. 2008.** Genomic plasticity and the diversity of polyploid plants.
615 *Science* **320**: 481–483.
- 616 **Levin D. 1983.** Polyploidy and novelty in flowering plants. *American Naturalist* **122**:
617 1–25.
- 618 **Márquez-Corro JI, Martín-Bravo S, Spalink D, Luceño M, Escudero M. 2019.**
619 Inferring hypothesis-based transitions in clade-specific models of chromosome

- 620 number evolution in sedges (Cyperaceae). *Molecular phylogenetics and evolution*
621 **135**: 203–209.
- 622 **Mayrose I, Barker MS, Otto SP. 2010.** Probabilistic models of chromosome number
623 evolution and the inference of polyploidy. *Systematic Biology* **59**: 132–144.
- 624 **Nakhleh L. 2010.** Evolutionary phylogenetic networks: models and issues. In:
625 Problem solving handbook in computational biology and bioinformatics. Springer,
626 125–158.
- 627 **Pennell MW, FitzJohn RG, Cornwell WK, Harmon LJ. 2015.** Model adequacy
628 and the macroevolution of angiosperm functional traits. *The American Naturalist* **186**:
629 E33–E50.
- 630 **R Core Team. 2013.** R: A language and environment for statistical computing.
- 631 **Ramsey J, Ramsey TS. 2014.** Ecological studies of polyploidy in the 100 years
632 following its discovery. *Philosophical Transactions of the Royal Society B: Biological*
633 *Sciences* **369**: 1–20.
- 634 **Ramsey J, Schemske DW. 2002.** Neopolyploidy in flowering plants. *Annual Review*
635 *of Ecology and Systematics* **33**: 589–639.
- 636 **Rice A, Glick L, Abadi S, Einhorn M, Kopelman NM, Salman-Minkov A,**
637 **Mayzel J, Chay O, Mayrose I. 2015.** The Chromosome Counts Database (CCDB)—a
638 community resource of plant chromosome numbers. *New Phytologist* **206**: 19–26.
- 639 **Rice A, Šmarda P, Novosolov M, Drori M, Glick L, Sabath N, Meiri S, Belmaker**
640 **J, Mayrose I. 2019.** The global biogeography of polyploid plants. *Nature ecology &*
641 *evolution* **3**: 265–273.
- 642 **Ruffini Castiglione M, Cremonini R. 2012.** A fascinating island: 2n=4. *Plant*
643 *Biosystems—An International Journal Dealing with all Aspects of Plant Biology* **146**:
644 711–726.
- 645 **Salman-Minkov A, Sabath N, Mayrose I. 2016.** Whole-genome duplication as a key
646 factor in crop domestication. *Nature Plants* **2**: 1–4.
- 647 **Shannon CE. 1948.** A mathematical theory of communication. *Bell system technical*

- 648 *journal* **27**: 379–423.
- 649 **Slater GJ, Pennell MW. 2013.** Robust regression and posterior predictive simulation
650 increase power to detect early bursts of trait evolution. *Systematic Biology* **63**: 293–
651 308.
- 652 **Soltis D, Soltis P, Schemske D. 2007.** Autopolyploidy in angiosperms: have we
653 grossly underestimated the number of species? *Taxon* **56**: 13–30.
- 654 **Soltis DE, Visger CJ, Marchant DB, Soltis PS. 2016.** Polyploidy: Pitfalls and paths
655 to a paradigm. *American Journal of Botany* **103**: 1146–1166.
- 656 **Spoelhof JP, Soltis PS, Soltis DE. 2017.** Pure polyploidy: Closing the gaps in
657 autopolyploid research. *Journal of Systematics and Evolution* **55**: 340–352.
- 658 **Stebbins GL. 1971.** *Chromosomal evolution in higher plants*. London, UK: Edward
659 Arnold Ltd.
- 660 **Vershinina AO, Lukhtanov VA. 2017.** Evolutionary mechanisms of runaway
661 chromosome number change in *Agrodiaetus* butterflies. *Scientific Reports* **7**: 1–9.
- 662 **Weiss-Schneeweiss H, Schneeweiss GM. 2013.** Karyotype diversity and
663 evolutionary trends in angiosperms. In: Greilhuber J, Dolezel J, Wendel JF, eds. *Plant*
664 *Genome Diversity Volume 2*. Vienna, Austria: Springer, 209–230.
- 665 **Wendel JF. 2015.** The wondrous cycles of polyploidy in plants. *American Journal of*
666 *Botany* **102**: 1753–1756.
- 667 **Wood TE, Takebayashi N, Barker MS, Mayrose I, Greenspoon PB, Rieseberg**
668 **LH. 2009.** The frequency of polyploid speciation in vascular plants. *Proceedings of*
669 *the National Academy of Sciences of the United States of America* **106**: 13875–13879.
- 670 **Zenil-Ferguson R, Burleigh JG, Ponciano JM. 2018.** chromploid: An R package
671 for chromosome number evolution across the plant tree of life. *Applications in Plant*
672 *Sciences* **6**: e1037.
- 673 **Zenil-Ferguson R, Ponciano JM, Burleigh JG. 2017.** Testing the association of
674 phenotypes with polyploidy: An example using herbaceous and woody eudicots.
675 *Evolution* **71**: 1138–1148.

676 **Zhan SH, Drori M, Goldberg EE, Otto SP, Mayrose I. 2016.** Phylogenetic
677 evidence for cladogenetic polyploidization in land plants. *American journal of botany*
678 **103:** 1252–8.

679 **Tables**

680 **Table 1.** The set of chromEvol models examined in this study, together with their rate
681 parameters.

Model	Model parameters ¹	Description	Nested in ²
D_{ys}	λ, δ	Dyploidy (descending or ascending)	$D_{ys}D_{up}, D_{ys}D_{up}D_{em}^*, D_{ys}D_{up}D_{em}, D_{ys}B_{num}, D_{ys}D_{up}B_{num}$
$D_{ys}D_{up}$	λ, δ, ρ	Dyploidy and duplication	$D_{ys}D_{up}D_{em}^*, D_{ys}D_{up}D_{em}, D_{ys}D_{up}B_{num}$
$D_{ys}D_{up}D_{em}^*$	$\lambda, \delta, \rho = \mu$	Dyploidy, constraining equal rates of duplication and demi-polyploidy	
$D_{ys}D_{up}D_{em}$	$\lambda, \delta, \rho, \mu$	Dyploidy, duplication, and demi-polyploidy	$D_{ys}D_{up}D_{em}^*$
$D_{ys}B_{num}$	$\lambda, \delta, \beta, v$	Dyploidy and base number transition	$D_{ys}D_{up}B_{num}$
$D_{ys}D_{up}B_{num}$	$\lambda, \delta, \rho, \beta, v$	Dyploidy, base number transition, and duplication	

682 ¹ The model parameters are the base number (β), and rates of ascending dyploidy (λ),
683 descending dyploidy (δ), duplication (ρ), demi-duplication (μ), and base number
684 transition (v).

685 ² In case all parameters of the model are a subset of other models, the more complex
686 models are indicated.

687 **Table 2.** The eight simulation scenarios examined in this study.

Genus	Number of taxa	Generating model	Model parameters ¹				
			λ	δ	ρ	β	ν
<i>Aloe</i>	120	D _{ys} D _{up}	0 (0)	0.34 (1)	2.61 (8)		
<i>Phacelia</i>	53	D _{ys} D _{up}	0.20 (2)	2.33 (21)	0.67 (6)		
<i>Lupinus</i>	77	D _{ys}	0.85 (7)	9.53 (76)			
<i>Hypochoeris</i>	38	D _{ys}	1.14 (5)	0.43 (2)			
<i>Brassica</i>	36	D _{ys} B _{num}	1.24 (11)	0.70 (6)		8	0.55 (5)
<i>Pectis</i>	49	D _{ys} B _{num}	0 (0)	0.40 (2)		12	0.55 (3)
<i>Crepis</i>	81	D _{ys} D _{up} B _{num}	2.41 (19)	0.99 (8)	0.26 (2)	8	0.18 (1)
<i>Hordeum</i>	36	D _{ys} D _{up} B _{num}	0 (0)	0 (0)	1.80 (5)	7	1.36 (4)

688 ¹ In parentheses: average number of simulated events across the tree.

689 **Table 3.** The inadequacy rates of the four tested models in the various simulation
 690 scenarios examined (100 simulations per tested model per scenario).

Simulation scenario	Generating Model	Tested Models ¹			
		$D_{ys}D_{up}$	D_{ys}	$D_{ys}B_{num}$	$D_{ys}D_{up}B_{num}$
Aloe	$D_{ys}D_{up}$	0.07	1.00	0.06	0.08
Phacelia	$D_{ys}D_{up}$	0.04	0.99	0.04	0.03
Lupinus	D_{ys}	0.06	0.08	0.06	0.07
Hypochaeris	D_{ys}	0.03	0.06	0.04	0.05
Brassica	$D_{ys}B_{num}$	0.14	0.98	0.05	0.06
Pectis	$D_{ys}B_{num}$	0.86	1.00	0.12	0.07
Crepis	$D_{ys}D_{up}B_{num}$	0.27	0.95	0.11	0.08
Hordeum	$D_{ys}D_{up}B_{num}$	0.77	1.00	0.30	0.13

691 ¹ The diagonal (white cells) are cases where the generating model is also the tested
 692 model. Dark grey represents over-parametrized models, light grey under-parametrized
 693 models, and patterned cells miss-parametrized models.

694 **Figure Legends**

695 **Fig 1.** A schematic illustration of the model adequacy framework for likelihood
696 models of chromosome-number evolution. In the case illustrated here, the model is
697 adequate because none of the test statistics lies in the tails of the simulated
698 distribution.

699 **Fig. 2.** The mean inference errors obtained under adequate and inadequate models for
700 each simulated scenario. Each row presents the error of a different inference task.
701 From top to bottom: inferring the total number of polyploid events across the tree,
702 inferring the total number of dysploid events across the tree, ploidy level assignments
703 of extant taxa, the probability of the chromosome number at the root of the
704 phylogeny. Each column denotes a different simulation scenario. For each scenario,
705 300 simulations were conducted and runs were partitioned to adequate and inadequate
706 models. The violin plots represent the distribution of the errors obtained for the
707 adequate (light grey, right) and inadequate (dark grey, left) sets. The black dot within
708 each distribution denotes its mean. Asterisk indicates significant difference between
709 the two groups (*, $p < 0.05$ and ***, $p < 0.01$).

710 **Fig. 3.** Application of the model adequacy test to 200 angiosperm genera. **(a)** A bar
711 plot representing the frequency of selection according to the AIC of each of the six
712 tested models in the 200 examined angiosperm genera. The height of each bar is
713 partitioned according to the percentage of genera that were determined as adequate
714 (light blue) or inadequate (red). **(b)** The adequacy rate of each model when applied to
715 all genera, regardless of whether the model was selected ($n = 200$).

716 **Supplementary information**

717 **Supplementary Information Methods:**

718 Item 1: Description of the simulation procedures.

719 **Supplementary Information Tables:**

720 Table S1: Pearson's r coefficient between each pair of statistics.

721 Table S2: The generating and fitted model for each simulation scenario used in the
722 comparison of inference error between adequate and inadequate models.

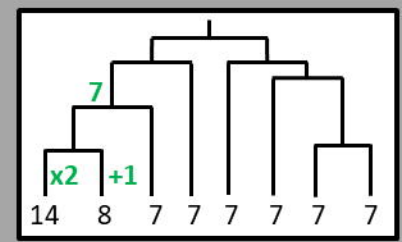
723 Table S3: Details of the seven plant clades, whose taxonomic rank is above the genus
724 level, examined in this study.

725 Table S4: Type I error rates for each test statistic per simulation scenario.

726 Table S5: Adequacy rates of all models, including those of the chosen models.

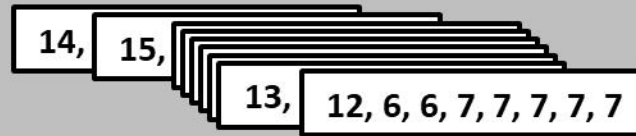
MODEL FITTING

1. Given a tree and chromosome counts, fit a model to the data.



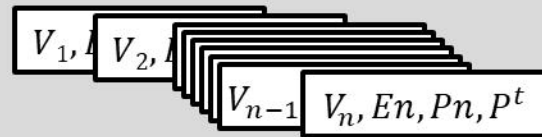
SIMULATE DATASETS

2. Simulate n datasets of chromosome counts using the optimized parameters.



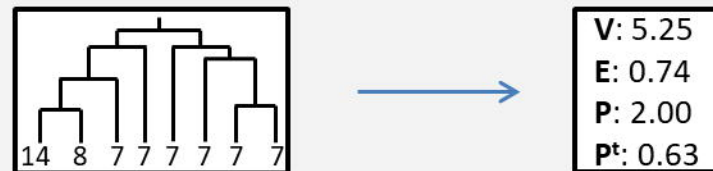
SIMULATIONS STATISTICS

3. Calculate the test statistics for the simulated datasets: Variance, Entropy, Parsimony score, and Parsimony vs. time.



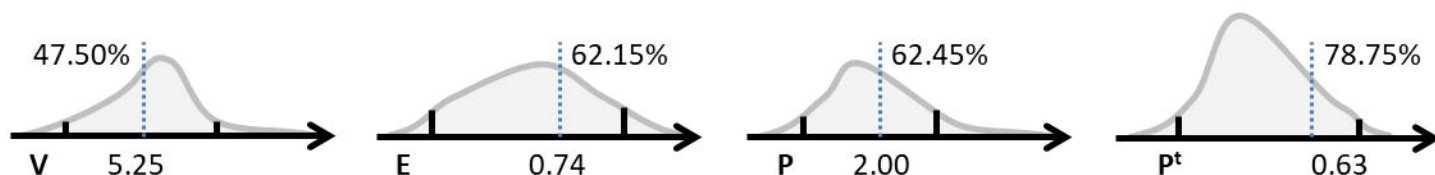
ORIGINAL STATISTICS

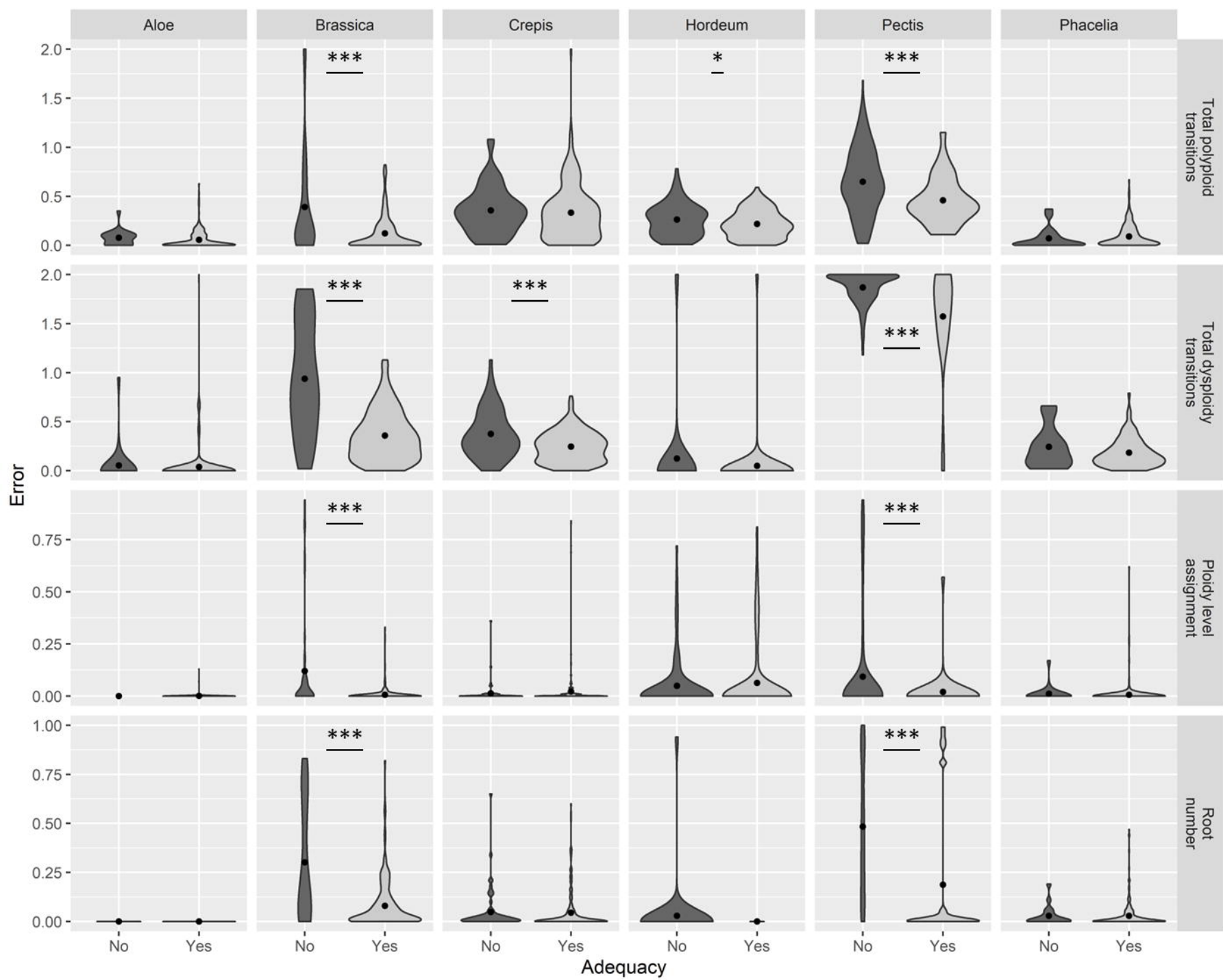
4. Calculate the same statistics for the empirical data:

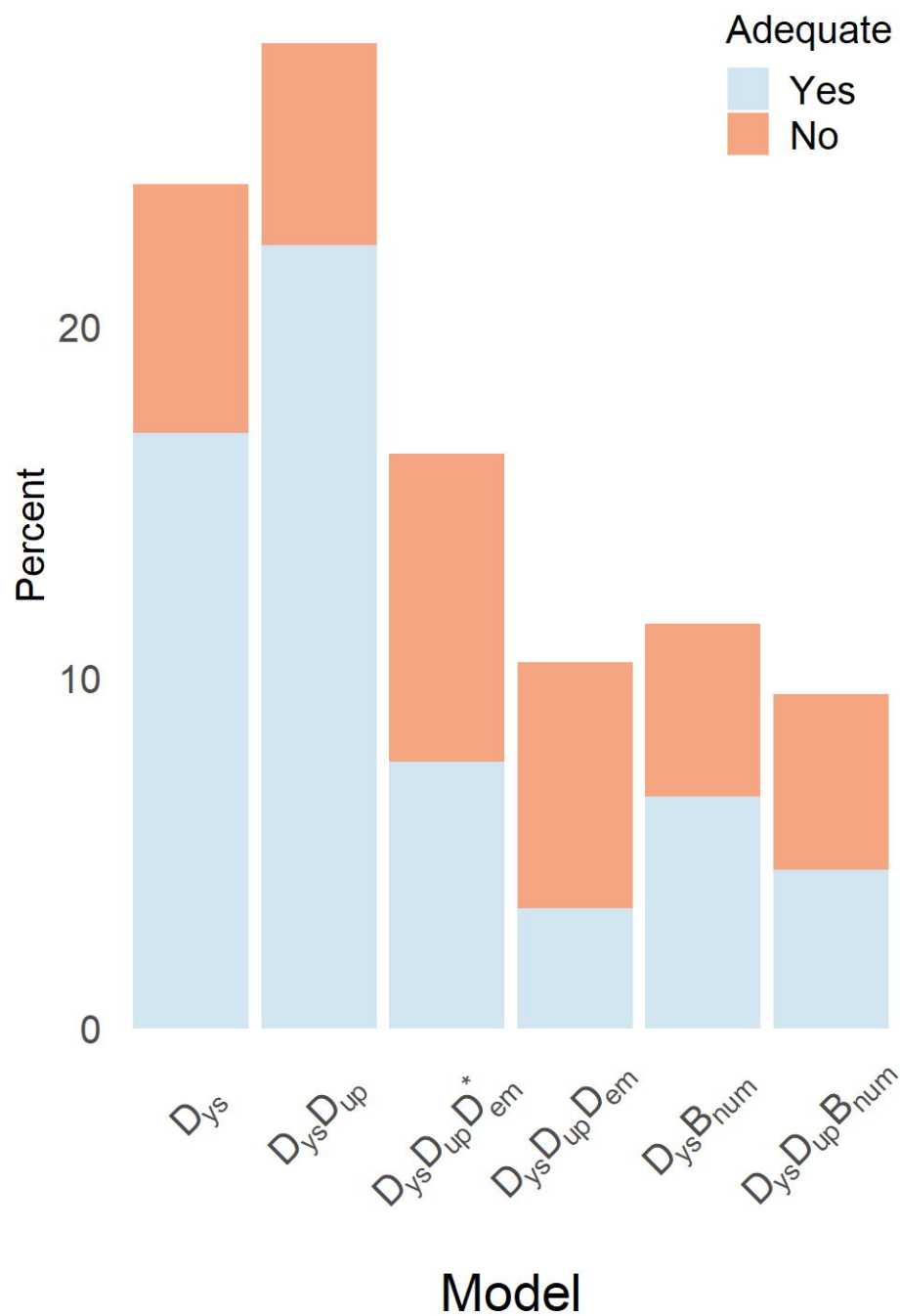


INFER ADEQUACY

5. Compare test statistics of the empirical and simulated data.





(a)**(b)**