

1 **asymptoticMK: A web-based tool for the asymptotic McDonald–Kreitman test**

2

3 Benjamin C. Haller and Philipp W. Messer

4

5 Department of Biological Statistics and Computational Biology, Cornell University, Ithaca, NY  
6 14853, USA

7

8 **Corresponding author:** Philipp W. Messer, 102J Weill Hall, Cornell University, Ithaca, NY  
9 14853, phone: 607-255-3984, email: messer@cornell.edu

10

11 **Keywords:** molecular evolution, positive selection, web service

12

13 **ABSTRACT**

14

15 The McDonald–Kreitman (MK) test is a widely used method for quantifying the role of positive  
16 selection in molecular evolution. One key shortcoming of this test lies in its sensitivity to the  
17 presence of slightly deleterious mutations, which can severely bias its estimates. An asymptotic  
18 version of the MK test was recently introduced that addresses this problem by evaluating  
19 polymorphism levels for different mutation frequencies separately, and then extrapolating a  
20 function fitted to that data. Here we present asymptoticMK, a web-based implementation of this  
21 asymptotic McDonald–Kreitman test. Our web service provides a simple R-based interface into  
22 which the user can upload the required data (polymorphism and divergence data for the genomic  
23 test region and a neutrally evolving reference region). The web service then analyzes the data  
24 and provides plots of the test results. This service is free to use, open-source, and available at  
25 <http://benhaller.com/messerlab/asymptoticMK.html>. We provide results from simulations to  
26 illustrate the performance and robustness of the asymptoticMK test under a wide range of model  
27 parameters.

## 28 INTRODUCTION

29

30 The extent to which molecular evolution is driven by positive selection, rather than neutral  
31 evolutionary processes such as random genetic drift, is one of the central questions of modern  
32 evolutionary biology. This question can be studied quantitatively by estimating the parameter  $\alpha$ ,  
33 which specifies the fraction of nucleotide substitutions in a given genomic region that were  
34 driven to fixation by positive selection (Eyre-Walker 2006). Values of  $\alpha$  close to one indicate  
35 that most substitutions in the region were indeed the result of positive selection, whereas values  
36 close to zero indicate neutral evolution.

37 One of the most widely used approaches for inferring  $\alpha$  from polymorphism and  
38 divergence data is the McDonald–Kreitman (MK) test (McDonald and Kreitman 1991; Eyre-  
39 Walker 2006), which compares levels of divergence between a genomic test region and a  
40 neutrally evolving reference region with the levels of polymorphism in the two regions. Early  
41 applications of the MK test typically focused on nonsynonymous sites in protein-coding regions  
42 as the test region, while synonymous sites were used as the neutral reference. However, the  
43 approach can also be applied to arbitrary genomic compartments or classes of mutations  
44 (Andolfatto 2005).

45 The original MK test makes several critical assumptions about the nature of the  
46 evolutionary process. First, it assumes that the positively selected mutations that ultimately  
47 contribute to divergence in the test region go to fixation quickly, such that they do not contribute  
48 noticeably to polymorphism levels. Second, it assumes that deleterious mutations in the test  
49 region are sufficiently deleterious to be lost quickly, such that they contribute to neither  
50 polymorphism nor divergence. Finally, neutral mutations in the test region are assumed to be  
51 subject to drift similar to the mutations in the neutral reference region and can therefore  
52 contribute to both polymorphism and divergence. Under these assumptions, it holds that

53

$$54 \quad (1) \quad \alpha = 1 - \frac{d_0 p}{d p_0},$$

55

56 where  $d$  and  $d_0$  are substitution rates per site in the test region and neutral reference region,  
57 respectively, while  $p$  and  $p_0$  specify the respective levels of polymorphism per site in the two  
58 regions (Eyre-Walker 2006). Note that if polymorphism and divergence levels are estimated over the

59 same region, the total number of sites cancels out in the ratios  $p/d$  and  $d_0/p_0$  and one may then simply use  
60 the actual counts of observed substitutions ( $D$  and  $D_0$ ) and polymorphic sites ( $P$  and  $P_0$ ) instead of rates  
61 per site (Eyre-Walker 2006).

62 With the growing availability of genome-level polymorphism and divergence datasets,  
63 the MK test has become a popular method for inferring positive selection in various organisms  
64 (Fay 2011). Several software tools and web services with implementations of the test have also  
65 been developed (Egea *et al.* 2008; Librado and Rozas 2009; Eyre-Walker and Keightley 2009;  
66 Stoletzki and Eyre-Walker 2011; Vos *et al.* 2013). The estimates of  $\alpha$  obtained in these studies  
67 range from as high as  $\sim 0.5$  for nonsynonymous substitutions in *Drosophila* (Sella *et al.* 2009), to  
68 close to zero in organisms such as yeast (Elyashiv *et al.* 2010) or many plants (Gossmann *et al.*  
69 2010). Indeed, estimates of  $\alpha$  obtained from Equation (1) are often negative, indicating that at  
70 least some of the assumptions of the test were likely not met (since negative values of  $\alpha$  have no  
71 biological meaning – estimates of  $\alpha$  may be negative, but the true value cannot be).

72 One major problem with the original MK test lies in its assumption that deleterious  
73 mutations do not contribute to polymorphism in the test region. This stands in contrast to the  
74 frequent observation of weakly deleterious mutations in many organisms, and the fact that such  
75 mutations can substantially affect the site frequency spectrum (SFS) of polymorphisms in  
76 functional genomic regions (Bustamante *et al.* 2005; Eyre-Walker *et al.* 2006). In the presence of  
77 weakly deleterious mutations, the observed level of polymorphism in the test region ( $p$ ) in  
78 Equation (1) will overestimate the rate at which neutral polymorphisms are expected to go to  
79 fixation in this region, which will bias estimates of  $\alpha$  downwards (providing one possible  
80 explanation for the frequent observation of negative  $\alpha$  estimates).

81 As one strategy to address this problem, it has been proposed to only consider  
82 polymorphisms for which the derived allele is above a certain threshold frequency when  
83 estimating  $p$  and  $p_0$  (Charlesworth and Eyre-Walker 2008). This is because the fraction of  
84 weakly deleterious mutations among all polymorphisms should be lower for higher derived-  
85 allele frequencies. Ideally, one would wish to set this cutoff high, to minimize the bias due to  
86 weakly deleterious mutations; however, the higher this cutoff, the fewer polymorphisms will  
87 actually remain in the dataset, thus increasing statistical noise. To circumvent this problematic  
88 tradeoff, more sophisticated extensions of the original MK test first attempt to infer the actual  
89 distribution of fitness effects among new mutations in the test region from the SFS, and then

90 correct fixation probabilities accordingly (Boyko *et al.* 2008; Eyre-Walker and Keightley 2009).  
91 Yet these approaches can still suffer from unknown effects of demography or linked selection  
92 that are also expected to affect the shape of the SFS. The most sophisticated extensions of the  
93 test therefore additionally incorporate basic demographic models to improve their estimates  
94 (Keightley and Eyre-Walker 2007; Boyko *et al.* 2008; Eyre-Walker and Keightley 2009), which  
95 requires additional (and often uncertain) assumptions about the demographic history of the  
96 population of interest.

97 In contrast to such model-based approaches, a considerably simpler, heuristic approach  
98 was recently proposed by Messer and Petrov (2013). This approach generalizes the frequency-  
99 cutoff approach described above, without the need to discard polymorphism data. Instead of  
100 setting a specific frequency cutoff, it separately estimates  $\alpha$  for each of a set of discrete  
101 mutational frequency classes:

102

$$103 \quad (2) \quad \alpha(x) = 1 - \frac{d_0 p(x)}{d p_0(x)}.$$

104

105 Here  $p(x)$  and  $p_0(x)$  specify the levels of polymorphism per site in the test and reference regions,  
106 respectively, considering only those polymorphisms for which the derived allele is present at  
107 frequency  $x$  in the population (estimated from a population sample, for example). In the presence  
108 of deleterious mutations,  $\alpha(x)$  will underestimate the true value of  $\alpha$  for small  $x$ , yet should  
109 converge to the correct value as  $x$  approaches one (Messer and Petrov 2013). The asymptotic  
110 estimate of  $\alpha$  is then obtained by fitting a function  $\alpha_{\text{fit}}(x)$  to the empirical  $\alpha(x)$  values and  
111 extrapolating this function to  $x = 1$ :

112

$$113 \quad (3) \quad \alpha_{\text{asymptotic}} = \alpha_{\text{fit}}(x = 1).$$

114

115 One key advantage of this approach is that because  $\alpha(x)$  does not depend on the  
116 individual functions  $p(x)$  and  $p_0(x)$  but only on their ratio, any biases due to demography or  
117 linked selection that affect the SFS in the test and reference regions in the same way will  
118 effectively cancel out (Messer and Petrov 2013). Another advantage over model-based  
119 approaches is that the asymptotic McDonald–Kreitman approach is much more computationally  
120 efficient, as it requires only fitting a simple curve to the data.

121 In this paper, we present asymptoticMK, a web-based tool for executing the asymptotic  
122 McDonald–Kreitman test quickly and easily in any web browser. After the necessary values are  
123 entered, asymptoticMK generates analyses and plots that are directly usable in publications. It is  
124 based internally on R, but no knowledge of R is needed to use it, nor does the user of  
125 asymptoticMK need to have R installed on their computer. For those who do wish to run the test  
126 themselves in R, the necessary code is freely available online. The asymptoticMK service can  
127 also be run in an automated fashion at the command line, for bulk analysis in script-based  
128 workflows. Finally, we present results from forward genetic simulations to illustrate the  
129 performance and robustness of the asymptotic MK test in various scenarios.

130

## 131 **MATERIALS AND METHODS**

132

### 133 **Implementation**

134 The asymptoticMK web service is implemented in R (R Development Core Team 2016). It uses  
135 the package FastRWeb (Urbanek 2008) to parse HTTP requests and generate responses, and uses  
136 the package Rserve (Urbanek 2003) as the lower-level interface that communicates with the web  
137 server through the standard CGI mechanism. A version of asymptoticMK that runs in R on the  
138 user’s local machine is also provided. Source code and additional resources related to  
139 asymptoticMK are posted on GitHub at <https://github.com/MesserLab/asymptoticMK>.

140

### 141 **Usage**

142 The web service is free to use, without license restrictions of any kind, and is available at  
143 <http://benhaller.com/messerlab/asymptoticMK.html>. That URL displays an entry page (Figure 1)  
144 with an input form in which the user may enter the necessary data for the test:  $d$  (the substitution  
145 rate in the test region),  $d_0$  (the substitution rate in the neutral reference region), and an uploaded  
146 file containing tab-delimited rows of data with values for  $x$  (the derived allele frequency),  $p(x)$   
147 (the polymorphism level in the test region at that frequency), and  $p_0(x)$  (the polymorphism level  
148 in the neutral reference region at that frequency). A sample polymorphism file is provided on the  
149 website. In practice, it is often advisable to combine polymorphism levels into a smaller number  
150 of frequency bins, where  $x$  then specifies the central frequency of the bin. This is particularly  
151 relevant when the data includes frequencies at which no polymorphisms are actually present in

152 the neutral region, in which case  $\alpha(x)$  would be undefined for those particular frequencies  
153 according to Equation (2). The frequency bins supplied to asymptoticMK do not need to be  
154 equally spaced, but to obtain the best possible  $\alpha$  estimate it is preferable to have bins providing  
155 good coverage across the full frequency spectrum. The input form also allows entry of minimum  
156 and maximum values defining a cutoff interval for  $x$ , such that the test is run using only the  
157 polymorphisms whose frequencies fall within that cutoff interval; this is usually desirable as a  
158 means of excluding the lowest- and highest-frequency polymorphisms, where SNP quality issues  
159 and polarization errors are generally most pronounced. This frequency cutoff is set to [0.1, 0.9]  
160 by default, but should be adjusted as needed.

161       Upon submission of the web form, asymptoticMK conducts its analysis and then opens a  
162 results page in a new browser tab, presenting a summary of the input data and the results from  
163 the analysis. The first plot on this results page shows binned polymorphism counts,  $p_0(x)$  and  
164  $p(x)$ , for the submitted data; the second plot shows that same data normalized (i.e., the  
165 normalized SFS in the test and reference regions). A third plot shows the calculated empirical  
166  $\alpha(x)$  as a function of  $x$ , estimated from the input data according to Equation (2). The fourth plot  
167 shows the same  $\alpha(x)$  data, with the best-fitting model and the asymptotic estimate of  $\alpha$  from  
168 Equation (3) superimposed upon the data.

169       Below these plots, the results of the analysis are presented in two tables. The first table  
170 provides the coefficients  $a$ ,  $b$ , and (for exponential fits)  $c$  of the model yielding the best fit to the  
171 data. The second table provides the estimated  $\alpha_{\text{asymptotic}}$  according to Equation (3), and the upper  
172 and lower limits of the 95% confidence interval around that estimate, as well as the estimated  $\alpha$   
173 from the original non-asymptotic McDonald–Kreitman test ( $\alpha_{\text{original}}$ ) for comparison (also  
174 estimated from all polymorphisms falling within the frequency cutoff interval specified on the  
175 input page).

176       For purposes of automation, the asymptoticMK web service can also be run at the  
177 command line using the Linux/Unix `curl` command. For example, the command

178

```
179 curl -F"d=593" -F"d0=930" -F"xlow=0.1" -F"xhigh=0.9" -  
180 F"datafile=@polymorphisms.txt" -F"reply=table" -o "MK_table.txt"  
181 http://benhaller.com/cgi-bin/R/asymptoticMK_run.html
```

182

183 would run asymptoticMK with the given values of  $d$  and  $d_0$ , the given  $x$  cutoff interval, and  
184 polymorphism data uploaded from the local file `polymorphisms.txt`, and would output a simple  
185 table of results to the file `MK_table.txt`. Further documentation on the use of this feature is  
186 provided on the asymptoticMK web page.

187 Finally, it is also possible to run asymptoticMK in R on the user's local machine. The R  
188 code for doing so can be found on asymptoticMK's GitHub repository. In addition to allowing the  
189 user to modify asymptoticMK's analysis as desired, this option also allows PDF plots to be  
190 created, rather than the PNG plots provided by the web-based service.

191

### 192 **Fitting and analysis procedure**

193 The asymptotic McDonald–Kreitman test first involves calculating values of  $\alpha(x)$  by applying  
194 Equation (2) to each frequency bin provided, as described by Messer and Petrov (2013). The  
195 next step involves fitting a function  $\alpha_{\text{fit}}(x)$  to these empirical  $\alpha(x)$  values. For greater robustness,  
196 asymptoticMK fits two functions to the data. The first function is exponential, of the form  $\alpha_{\text{fit}}(x)$   
197  $= a + b \exp(-cx)$ , and is fitted using the `nls2()` function, from the R package `nls2` (Grothendieck  
198 2013). This fit is done in two steps. First, a brute-force scan for the closest fit is conducted across  
199 the likely portion of the three-dimensional parameter space defined by  $a$ ,  $b$ , and  $c$ , by exhaustive  
200 search. This supplies reasonably good starting values for the second step, which refines those  
201 starting values using standard nonlinear least-squares regression. While this two-step procedure  
202 generally works well, it can fail to converge if the data is not exponential in form.

203 To address this possibility of nonconvergence of the exponential fit, asymptoticMK also  
204 fits a linear function of the form  $\alpha_{\text{fit}}(x) = a + bx$ , with the `lm()` function that is part of the `stats`  
205 package included in R. This fit always converges, and thus provides a backstop that allows the  
206 test to complete even when given irregular or extremely noisy data; however, it is always  
207 recommended that the results of the analysis be inspected visually to confirm that they are in fact  
208 meaningful.

209 Once these two models have been fitted, asymptoticMK chooses which model will be  
210 used for the remainder of the analysis. If the exponential fit failed to converge, then the linear  
211 model is chosen; if both fits succeeded, then the better model is chosen using the Akaike  
212 Information Criterion (AIC). Occasionally, in pathological cases, the exponential fit will have

213 the better AIC but will have extremely large coefficient standard error(s); in this case, the linear  
214 fit is chosen since predictions from the exponential model would be effectively worthless.

215 The chosen model is then used to provide an estimate of the value of  $\alpha_{\text{asymptotic}}$  according  
216 to Equation (3), by evaluating the fitted function  $\alpha_{\text{fit}}(x)$  at  $x = 1$ ; this is the primary result of the  
217 test, and provides the test's estimate of the true value of  $\alpha$  within the test region. A 95%  
218 confidence interval around this estimate is also calculated. For the exponential model, this is  
219 done using Monte Carlo simulation based upon the fitted model, using the `predictNLS()` function  
220 published online by Spiess (2013); for the linear model, it is done using the standard R function  
221 `predict()`.

222

### 223 **Test datasets**

224 To provide a test of asymptoticMK using empirical data, we used the same *Drosophila*  
225 *melanogaster* dataset that Messer and Petrov (2013) used in their Figure 3C. This data set  
226 consists of SNPs obtained from the genome sequences of 162 inbred fly lines generated by the  
227 *Drosophila* genetic reference panel (Mackay *et al.* 2012). Divergence data was obtained from  
228 genome alignments between *D. melanogaster* and *D. simulans*, extracted from the 12 *Drosophila*  
229 genomes data (Clark *et al.* 2007). The test data in the asymptoticMK analysis ( $d$  and  $p$ ) are  
230 genome-wide nonsynonymous mutations, while synonymous sites were used as the neutral  
231 reference ( $d_0$  and  $p_0$ ). The polymorphism data is available online at on asymptoticMK's GitHub  
232 repository, with associated values  $d = 59570$  and  $d_0 = 159058$ . The default frequency cutoff  
233 interval of  $[0.1, 0.9]$  was used in the analysis of this dataset with asymptoticMK.

234 We also tested asymptoticMK on simulated data, using the forward genetic simulation  
235 framework SLiM 2 (Haller and Messer 2017). A population of 1000 diploid individuals was  
236 simulated to evolve in a total of 13 different scenarios, with 20 replicates for each scenario.  
237 Simulation runs depended upon six free parameters ( $T, L, \mu, r_b, s_d, s_b$ ) as described hereafter.  
238 After an initial burn-in period of 10,000 generations to equilibrate the model, runs executed for  $T$   
239 additional generations. The simulated chromosome was  $L$  base pairs long. Nucleotide mutations  
240 occurred uniformly at a rate of  $\mu$  per base per generation, and recombination occurred uniformly  
241 at a rate of  $10^{-7}$  per base per generation. Each new mutation was either of neutral type "m1"  
242 (relative proportion of 0.5 of all new mutations), of functional non-beneficial type "m2" (relative  
243 proportion of 0.5 of all new mutations), or of functional beneficial type "m3" (relative proportion



244 of  $r_b$  of all new mutations); these relative proportions were automatically rescaled by SLiM to be  
245 absolute proportions. The neutral m1 mutations always had a selection coefficient of  $s = 0.0$ ; the  
246 selection coefficients of m2 mutations were drawn from a gamma distribution with a mean of  $s_d$   
247 and a shape parameter of 0.2; and m3 mutations always had a selection coefficient of  $s_b$ . Fitness  
248 effects were assumed to be codominant. Every 500 generations after the burn-in period, all  
249 polymorphisms were recorded in the population by dividing them according to their frequency  
250 into 50 equal-width frequency bins, and then adding them to an ongoing binned tabulation. The  
251 SLiM configuration script used for these simulations is provided on asymptoticMK's GitHub  
252 repository.

253 The “baseline” parameterization of this model utilized parameter values of  $L = 10^7$ ,  
254  $\mu = 10^{-9}$ ,  $r_b = 0.0005$ ,  $s_b = 0.1$ ,  $s_d = -0.02$ , and  $T = 2 \times 10^5$ . The other twelve parameterizations  
255 involved either a “high” or a “low” value of one of the six parameters, replacing the “central”  
256 value used in the baseline scenario:  $L = 10^8$  or  $10^6$ ,  $\mu = 10^{-8}$  or  $10^{-10}$ ,  $r_b = 0.001$  or  $0.0001$ ,  
257  $s_b = 0.2$  or  $0.02$ ,  $s_d = -0.2$  or  $-0.002$ , and  $T = 2 \times 10^6$  or  $2 \times 10^4$ . At the end of each model run, we  
258 obtained binned values for  $p(x)$  and  $p_0(x)$ , where  $p_0$  was estimated from all polymorphisms  
259 involving mutations of type m1, while  $p$  was estimated from the combined mutations of types m2  
260 and m3. Values for  $d$  and  $d_0$  were obtained from the set of mutations fixed during the simulation;  
261 as with  $p_0$  and  $p$ ,  $d_0$  was estimated from all mutations of type m1, while  $d$  was estimated from the  
262 combined mutations of types m2 and m3. These values, output by the model, were used in  
263 asymptoticMK with the default  $x$  cutoff interval of  $[0.1, 0.9]$  to calculate an  $\alpha$  estimate. The  $\alpha$   
264 estimate from the original MK test was also calculated using the data within the same interval.  
265 Finally, the true value of  $\alpha$  was estimated from the simulation run as the fraction  $d_3 / (d_2 + d_3)$ ,  
266 where  $d_2$  is the number of m2 mutations fixed and  $d_3$  is the number of m3 mutations fixed. This  
267 value provides a metric for the accuracy of the  $\alpha$  estimates – a benefit of using simulated data,  
268 where the true  $\alpha$  can be calculated.

269 From this raw data provided by each set of 20 replicates for a given parameterization,  
270 summary statistics for that parameterization were computed. In particular, we calculated (i) the  
271 mean and standard deviation of the true  $\alpha$  values, (ii) the mean and standard deviation of the  
272 asymptoticMK  $\alpha$  estimates, (iii) the mean of the absolute differences between the true  $\alpha$  and the  
273 asymptoticMK estimate (i.e., the mean estimation error for asymptoticMK), (iv) the mean and  
274 standard deviation of the estimates of  $\alpha$  using the original non-asymptotic MK test, (v) the mean

275 of the absolute differences between the true  $\alpha$  and the original MK test estimate (i.e., the mean  
276 estimation error for the original MK test), and (vi) the fraction of the 20 replicates for which  
277 asymptoticMK chose an exponential (as opposed to linear) fit (Table 1).

278

## 279 RESULTS AND DISCUSSION

280

281 Results from our test of asymptoticMK with the empirical *D. melanogaster* dataset are shown in  
282 Figures 2A and 2B. The fitted exponential function is:  $\alpha_{\text{fit}}(x) = 0.585 - 0.622 \exp(-3.80x)$ . The  
283 asymptotic McDonald–Kreitman estimate provided by this model is  $\alpha_{\text{asymptotic}} = 0.571$ . These  
284 results match those obtained by Messer and Petrov (2013) using the same dataset (their Figure  
285 3C), as expected. The estimate provided by the original McDonald–Kreitman test is  $\alpha_{\text{original}} =$   
286 0.407, by comparison (shown in Figure 2B).

287 The results from the analysis of the SLiM simulations are shown in Table 1. In 12 of the  
288 13 parameterizations, the mean estimation error of asymptoticMK was markedly lower than that  
289 of the original MK test; in the other parameterization ( $T = 2 \times 10^4$ ) the tests performed similarly  
290 (mean estimation errors of 0.126 and 0.120). For 3 of the 13 parameterizations ( $L = 10^6$ ,  
291  $\mu = 10^{-10}$ , and  $T = 2 \times 10^4$ ), however, the mean estimation error of asymptoticMK was greater than  
292 0.1, indicating that  $\alpha$  estimates were relatively inaccurate for those simulations. These three  
293 parameterizations involved a shorter chromosome, a lower mutation rate, or a shorter duration,  
294 and thus all provided approximately ten times less polymorphism data upon which to base  
295 estimations than our baseline scenario. Accordingly, parameterizations that provided more  
296 polymorphism data ( $L = 10^8$ ,  $\mu = 10^{-8}$ , and  $T = 2 \times 10^6$ ) provided more accurate  $\alpha$  estimates (mean  
297 estimation errors below 0.03). This pattern was weak or absent for the original MK test; even for  
298 the high-data parameterizations the original MK test always showed a mean estimation error  
299 greater than 0.1, and its mean estimation error for the high-data  $L = 10^8$  case was actually  
300 substantially higher than for the low-data  $L = 10^6$  case (0.151 versus 0.120). This is consistent  
301 with the fact that the original MK test systematically underestimates  $\alpha$  in the presence of  
302 deleterious mutations (as discussed in the Introduction). The asymptotic MK test may still have a  
303 tendency toward underestimation as well, but errors are much smaller.

304 Another noteworthy observation is that in the high-data parameterizations the exponential  
305 fit was chosen by asymptoticMK in 100% of cases, whereas in the low-data parameterizations

306 the linear fit was chosen a majority of the time (Table 1). It is not the case that the linear model  
307 always produces a poor  $\alpha$  estimate; we observed many runs where the linear fit performed well.  
308 However, it may indicate that a poor cutoff interval was chosen, that the binning of the  
309 polymorphism data ought to be done differently, or that the data is simply too noisy. We suggest  
310 that the result of the asymptotic MK test should always be inspected visually to verify that the fit  
311 is reasonable and that appropriate cutoff intervals and bin sizes were used.

312 To illustrate how such manual inspection can help improve estimates, we examine two of  
313 the simulation runs from Table 1 in more detail. Figure 2C shows the automated fit for one of the  
314 simulations in the low-data  $L = 10^6$  scenario. A linear fit function produced an asymptotic  $\alpha$   
315 estimate of 0.0313, which is quite distant from the true value of 0.3462. The binned  
316 polymorphism data within the cutoff interval of [0.1, 0.9] is rather flat, but the data appears  
317 reasonable across the whole frequency spectrum in this case, and the upward trend of the data is  
318 much more visible outside of the cutoff interval used. Changing the cutoff interval to [0.0, 1.0]  
319 produces the fit shown in Figure 2D, with an asymptotic  $\alpha$  estimate of 0.2829 – much closer to  
320 the true value. Figure 2E shows the automated fit for another  $L = 10^6$  scenario run. This fit also  
321 used a linear fit function, producing an asymptotic  $\alpha$  estimate of 0.0103 compared to the true  
322 value of 0.2462. Here, the data are very noisy, which could be an indication that more bins have  
323 been used than can be robustly supported by the data. Re-binning the polymorphism data into  
324 half as many bins provides a less noisy dataset that results in a much better fit (Figure 2F), with  
325 an  $\alpha$  estimate of 0.1813 – again, a substantial improvement. These examples illustrate that  
326 automated fits can be particularly problematic in low-data situations such as the  $L = 10^6$  scenario,  
327 but that hand inspection and tailoring of the fitting process can sometimes improve the result  
328 noticeably.

329

## 330 **CONCLUSIONS**

331

332 In this paper, we presented asymptoticMK, a new web-based tool for executing the  
333 asymptotic McDonald–Kreitman test. To demonstrate its functionality, we analyzed both  
334 empirical and simulation-generated datasets. Our results illustrate the greater power of the  
335 asymptotic McDonald–Kreitman test to estimate the true value of  $\alpha$ , compared to the original  
336 non-asymptotic test. However, our results also underline the need for a large dataset to obtain

337 reasonably accurate results from the asymptotic test; estimates of  $\alpha$  from a single gene, or from a  
338 system with a very short divergence time, are unlikely to be meaningful. In addition, visual  
339 inspection of the quality of the fit used to estimate  $\alpha$  is necessary for accuracy. With attention to  
340 these caveats, the asymptoticMK service presented here allows the user to obtain  $\alpha$  estimates  
341 quickly and easily through any web browser, or using R on any machine.

342

## 343 **ACKNOWLEDGMENTS**

344

345 The authors thank G. Grothendieck, J. Horner, and S. Urbanek for their contributions to the R  
346 packages used here. We also thank S. Urbanek for his unstinting help with the installation and  
347 use of `FastRWeb`, A.-N. Spiess for his `predictNLS()` function, and most of all the R Core Team for  
348 R itself. `asymptoticMK` is only possible because of the free software developed by these and  
349 other contributors. This work was supported by funds from the College of Agriculture and Life  
350 Sciences at Cornell University to PWM. We also thank Dmitri Petrov for his important role in  
351 the original development of the asymptotic MK test, and the editor and two anonymous  
352 reviewers for their helpful comments.

353

## 354 **LITERATURE CITED**

355

356 Andolfatto, P., 2005 Adaptive evolution of non-coding DNA in *Drosophila*. *Nature* 437: 1149–  
357 1152.

358 Boyko, A. R., S. H. Williamson, A. R. Indap, J. D. Degenhardt, R. D. Hernandez *et al.*, 2008  
359 Assessing the evolutionary impact of amino acid mutations in the human genome. *PLoS*  
360 *Genet.* 4: e1000083.

361 Bustamante, C. D., A. Fledel-Alon, S. Williamson, R. Nielsen, M. T. Hubisz *et al.*, 2005 Natural  
362 selection on protein-coding genes in the human genome. *Nature* 437: 1153–1157.

363 Charlesworth, J., and A. Eyre-Walker, 2008 The McDonald-Kreitman test and slightly  
364 deleterious mutations. *Mol. Biol. Evol.* 25: 1007–15.

365 Clark, A. G., M. B. Eisen, D. R. Smith, C. M. Bergman, B. Oliver *et al.*, 2007 Evolution of genes  
366 and genomes on the *Drosophila* phylogeny. *Nature* 450: 203–218.

367 Egea, R., S. Casillas, and A. Barbadilla, 2008 Standard and generalized McDonald-Kreitman

368 test: a website to detect selection by comparing different classes of DNA sites. *Nucleic*  
369 *Acids Res.* 36: W157-62.

370 Elyashiv, E., K. Bullaughey, S. Sattath, Y. Rinott, M. Przeworski *et al.*, 2010 Shifts in the  
371 intensity of purifying selection: an analysis of genome-wide polymorphism data from two  
372 closely related yeast species. *Genome Res.* 20: 1558–73.

373 Eyre-Walker, A., 2006 The genomic rate of adaptive evolution. *Trends Ecol. Evol.* 21: 569–575.

374 Eyre-Walker, A., and P. D. Keightley, 2009 Estimating the rate of adaptive molecular evolution  
375 in the presence of slightly deleterious mutations and population size change. *Mol. Biol.*  
376 *Evol.* 26: 2097–108.

377 Eyre-Walker, A., M. Woolfit, and T. Phelps, 2006 The distribution of fitness effects of new  
378 deleterious amino acid mutations in humans. *Genetics* 173: 891–900.

379 Fay, J. C., 2011 Weighing the evidence for adaptation at the molecular level. *Trends Genet.* 27:  
380 343–9.

381 Grothendieck, G., 2013 nls2: Non-linear regression with brute force. Available at:  
382 <https://CRAN.R-project.org/package=nls2>. Accessed: December 14, 2016.

383 Gossman, T. I., B.-H. Song, A. J. Windsor, T. Mitchell-Olds, C. J. Dixon *et al.*, 2010 Genome  
384 wide analyses reveal little evidence for adaptive evolution in many plant species. *Mol. Biol.*  
385 *Evol.* 27: 1822–1832.

386 Haller, B. C., and P. W. Messer, 2017 SLiM 2: Flexible, interactive forward genetic simulations.  
387 *Mol. Biol. Evol.* 34: 230–240.

388 Keightley, P. D., and A. Eyre-Walker, 2007 Joint inference of the distribution of fitness effects  
389 of deleterious mutations and population demography based on nucleotide polymorphism  
390 frequencies. *Genetics* 177: 2251–2261.

391 Librado, P., and J. Rozas, 2009 DnaSP v5: a software for comprehensive analysis of DNA  
392 polymorphism data. *Bioinformatics* 25: 1451–1452.

393 Mackay, T. F., S. Richards, E. A. Stone, A. Barbadilla, J. F. Ayroles *et al.*, 2012 The *Drosophila*  
394 *melanogaster* genetic reference panel. *Nature* 482: 173–178.

395 McDonald, J. H., and M. Kreitman, 1991 Adaptive protein evolution at the *Adh* locus in  
396 *Drosophila*. *Nature* 351: 652–4.

397 Messer, P. W., and D. A. Petrov, 2013 Frequent adaptation and the McDonald-Kreitman test.  
398 *Proc. Natl. Acad. Sci. U. S. A.* 110: 8615–20.

399 R Development Core Team, 2016 *R: A language and for statistical computing*. Vienna, Austria.

400 Sella, G., D. A. Petrov, M. Przeworski, and P. Andolfatto, 2009 Pervasive natural selection in the  
401 *Drosophila* genome? PLoS Genet. 5: e1000495.

402 Spiess, A.-N. predictNLS (Part 1, Monte Carlo simulation): confidence intervals for “nls”  
403 models. R-bloggers. Available at: [https://www.r-bloggers.com/predictnls-part-1-monte-](https://www.r-bloggers.com/predictnls-part-1-monte-carlo-simulation-confidence-intervals-for-nls-models/)  
404 [carlo-simulation-confidence-intervals-for-nls-models/](https://www.r-bloggers.com/predictnls-part-1-monte-carlo-simulation-confidence-intervals-for-nls-models/). Accessed: December 14, 2016.

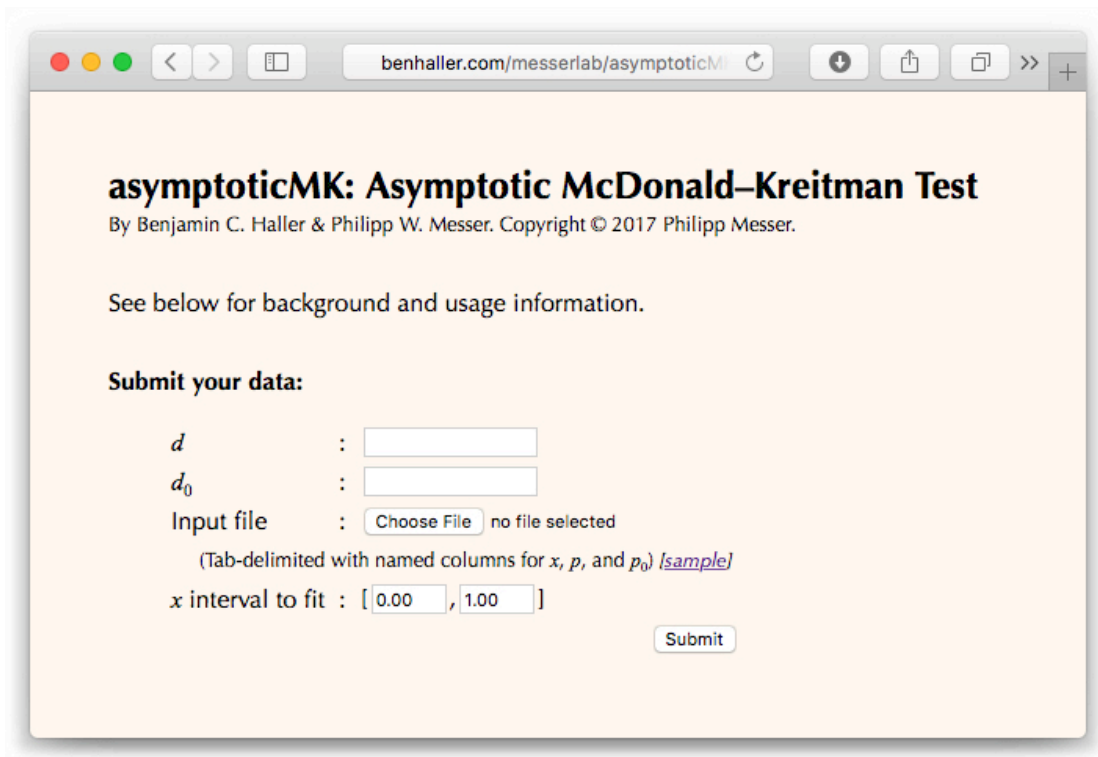
405 Stoletzki, N., and A. Eyre-Walker, 2011 Estimation of the neutrality index. Mol. Biol. Evol. 28:  
406 63–70.

407 Urbanek, S., 2008 FastRWeb: Fast interactive web framework for data mining using R, in *IASC*  
408 *2008 World Congress*. Available at: <https://rforge.net/FastRWeb/>. Accessed: December 14,  
409 2016.

410 Urbanek, S., 2003 Rserve - A fast way to provide R functionality to applications. In *Proceedings*  
411 *of the 3rd International Workshop on Distributed Statistical Computing*. Available at:  
412 <https://rforge.net/Rserve/>. Accessed: December 14, 2016.

413 Vos, M., T. A. H. te Beek, M. A. van Driel, M. A. Huynen, A. Eyre-Walker *et al.*, 2013 ODoSE:  
414 a webserver for genome-wide calculation of adaptive divergence in prokaryotes. PLoS One  
415 8: e62447.

416



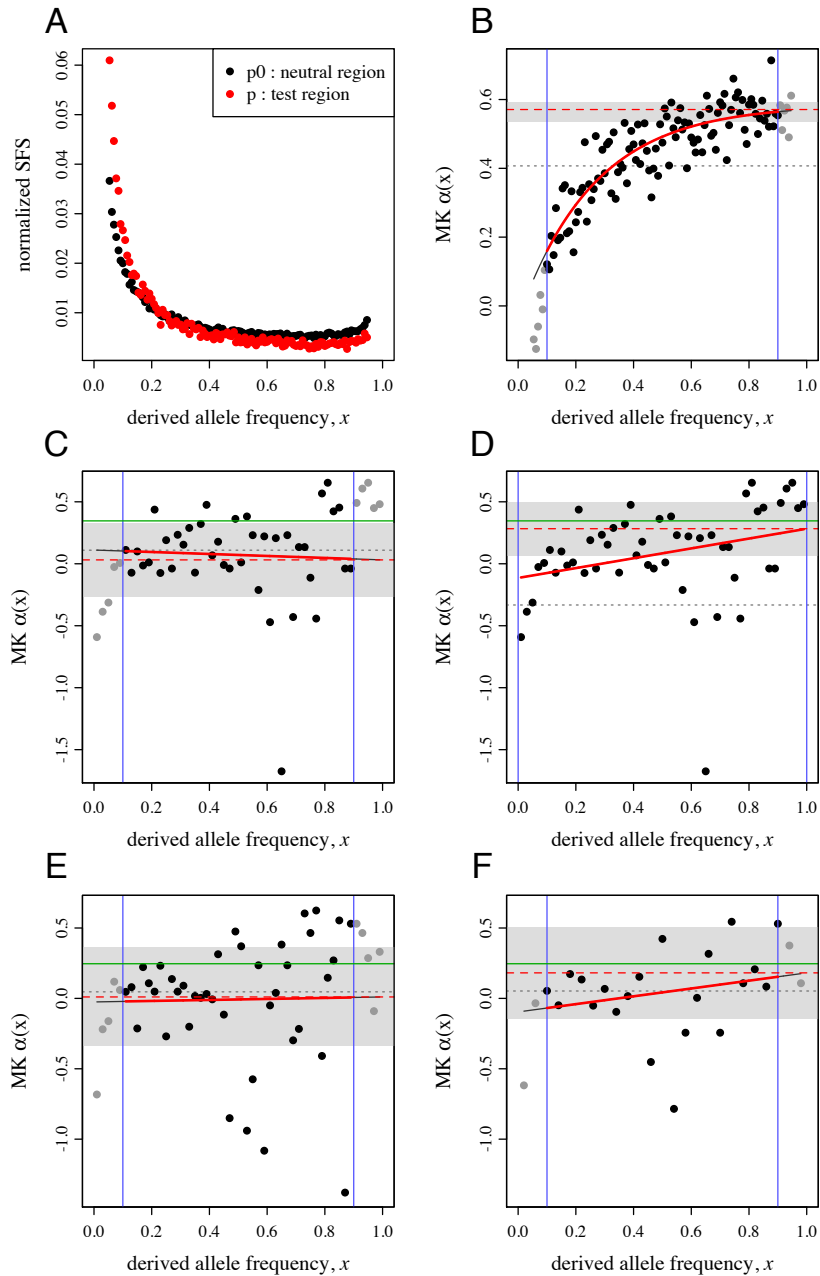
The screenshot shows a web browser window with the URL `benhaller.com/messerlab/asymptoticMK`. The page title is "asymptoticMK: Asymptotic McDonald–Kreitman Test" by Benjamin C. Haller & Philipp W. Messer. The page contains a form for submitting data. The form fields are:

- $d$  :
- $d_0$  :
- Input file :  no file selected
- (Tab-delimited with named columns for  $x$ ,  $p$ , and  $p_0$ ) [\[sample\]](#)
- $x$  interval to fit : [  ,  ]

A "Submit" button is located at the bottom right of the form.

418

419 Figure 1. A screenshot of the web page for asymptoticMK. After entering values for  $d$  and  $d_0$ ,  
420 choosing an input file with binned values for  $x$ ,  $p$ , and  $p_0$ , and choosing the  $x$  interval to fit, the  
421 user can click the Submit button and asymptoticMK will provide its results in a new browser  
422 window or tab.



423

424 Figure 2. Results from asymptoticMK for three test datasets. (A) Normalized site frequency  
 425 spectrum (SFS) for the *Drosophila* dataset used in Messer and Petrov (2013). Points show  
 426 normalized binned polymorphism frequencies for the neutral region (black) and the test region  
 427 (red). (B) Result of asymptoticMK's analysis of that dataset. The two vertical blue lines show the  
 428 limits of the frequency cutoff interval used for fitting. Points indicate binned values of  $\alpha(x)$ ,  
 429 estimated according to Equation 2; points are gray if they are outside the cutoff interval (and thus  
 430 not used in fitting). The solid red curve shows the fitted  $\alpha_{\text{fit}}(x)$  (here, exponential). The dashed



431 red line shows the estimate of  $\alpha_{\text{asymptotic}}$ , obtained from the fitted function according to Equation  
432 3. The gray band indicates the 95% confidence intervals around this  $\alpha_{\text{asymptotic}}$  estimate. The  
433 dotted gray line shows the estimate of  $\alpha_{\text{original}}$ , obtained from the original (non-asymptotic) MK  
434 test, for comparison (also calculated using only the data within the cutoff interval). (C) and (D)  
435 show corresponding results from one SLiM simulation run, and (E) and (F) show results from  
436 another SLiM simulation run; in each case, the first panel shows the result of an automated fit  
437 using asymptoticMK, whereas the second shows the improvement after hand tailoring of the fit  
438 (see Results and Discussion). Note that in all four cases the linear fit was deemed more  
439 appropriate by asymptoticMK. The solid green horizontal lines, finally, show the true value of  $\alpha$   
440 in the simulation runs for comparison.

441 **TABLES**

442

443 Table 1. Results from asymptoticMK for simulation runs conducted with SLiM 2. The first row  
 444 shows the averaged results from 20 replicate runs of the baseline SLiM model supplied on  
 445 GitHub (see text). These runs used parameter values of mutation rate  $\mu = 10^{-9}$  per base position  
 446 per generation, chromosome length  $L = 10^7$  base positions, beneficial mutation rate  $r_b = 0.0005$ ,  
 447 beneficial mutation selection coefficient  $s_b = 0.1$ , deleterious mutation selection coefficient  
 448  $s_d = -0.02$ , and time after burn-in  $T = 2 \times 10^5$  generations. Each subsequent row shows the results  
 449 from 20 replicate runs using the non-baseline parameter value shown.  $\alpha_{\text{true}}$  specifies the true  
 450 value of  $\alpha$  averaged across the 20 replicates in each row;  $\alpha_{\text{asymptotic}}$  and  $\alpha_{\text{original}}$  specify the  
 451 asymptoticMK estimate and the estimate from the original test, respectively. Standard deviations  
 452 across the 20 replicates of each row are shown as  $\pm$  values.  $\Delta_{\text{asymptotic}} = |\alpha_{\text{asymptotic}} - \alpha_{\text{true}}|$  and  
 453  $\Delta_{\text{original}} = |\alpha_{\text{original}} - \alpha_{\text{true}}|$  specify the mean absolute errors between true  $\alpha$  values and the estimates  
 454 from asymptoticMK and the original test, respectively, in each run, averaged over the 20  
 455 replicates.  $\rho_{\text{exp}}$  specifies the fraction of runs in which the exponential fit was chosen.

456

<b>Model</b>	$\alpha_{\text{true}}$	$\alpha_{\text{asymptotic}}$	$\alpha_{\text{original}}$	$\Delta_{\text{asymptotic}}$	$\Delta_{\text{original}}$	$\rho_{\text{exp}}$
baseline	0.329 $\pm$ 0.015	0.307 $\pm$ 0.058	0.164 $\pm$ 0.035	0.045	0.165	0.75
$L = 10^8$	0.327 $\pm$ 0.008	0.301 $\pm$ 0.013	0.174 $\pm$ 0.012	0.025	0.152	1.00
$L = 10^6$	0.321 $\pm$ 0.067	0.246 $\pm$ 0.134	0.142 $\pm$ 0.141	0.120	0.191	0.15
$\mu = 10^{-8}$	0.306 $\pm$ 0.005	0.287 $\pm$ 0.016	0.173 $\pm$ 0.009	0.019	0.132	1.00
$\mu = 10^{-10}$	0.317 $\pm$ 0.057	0.288 $\pm$ 0.169	0.145 $\pm$ 0.074	0.134	0.173	0.05
$r_b = 0.0010$	0.493 $\pm$ 0.018	0.481 $\pm$ 0.045	0.378 $\pm$ 0.025	0.041	0.114	0.70
$r_b = 0.0001$	0.091 $\pm$ 0.014	0.115 $\pm$ 0.080	-0.103 $\pm$ 0.053	0.071	0.194	0.55
$s_b = 0.20$	0.477 $\pm$ 0.016	0.451 $\pm$ 0.032	0.366 $\pm$ 0.025	0.029	0.111	0.70
$s_b = 0.02$	0.096 $\pm$ 0.011	0.090 $\pm$ 0.068	-0.119 $\pm$ 0.047	0.057	0.215	0.50
$s_d = -0.200$	0.424 $\pm$ 0.024	0.422 $\pm$ 0.042	0.289 $\pm$ 0.036	0.032	0.135	0.60
$s_d = -0.002$	0.233 $\pm$ 0.011	0.234 $\pm$ 0.057	0.104 $\pm$ 0.039	0.045	0.129	0.50
$T = 2 \times 10^6$	0.324 $\pm$ 0.006	0.302 $\pm$ 0.014	0.173 $\pm$ 0.012	0.022	0.151	1.00
$T = 2 \times 10^4$	0.345 $\pm$ 0.063	0.369 $\pm$ 0.183	0.225 $\pm$ 0.113	0.126	0.120	0.05

457