

1 **CAGEE: computational analysis of gene expression evolution**  
2 Jason Bertram<sup>1,2,\*</sup>, Ben Fulton<sup>1,3</sup>, Jason P. Tourigny<sup>1,4</sup>, Yadira Peña-Garcia<sup>1</sup>, Leonie C.  
3 Moyle<sup>1</sup>, and Matthew W. Hahn<sup>1,4,\*</sup>

4

5 <sup>1</sup>Department of Biology, Indiana University, Bloomington, IN 47405, USA

6 <sup>2</sup>Department of Mathematics, Western University, London, ON N6A 5B7, Canada

7 <sup>3</sup>University Information Technology Services, Indiana University, Bloomington, IN

8 47405, USA

9 <sup>4</sup>Department of Computer Science, Indiana University, Bloomington, IN 47405, USA

10

11 \*Corresponding authors: E-mail: jason.bertram@uwo.ca, mwh@indiana.edu

12

13

14 **Abstract**

15 Despite the increasing abundance of whole transcriptome data, few methods are  
16 available to analyze global gene expression across phylogenies. Here, we present a  
17 new software package (CAGEE) for inferring patterns of increases and decreases in  
18 gene expression across a phylogenetic tree, as well as the rate at which these changes  
19 occur. In contrast to previous methods that treat each gene independently, CAGEE can  
20 calculate genome-wide rates of gene expression, along with ancestral states for each  
21 gene. The statistical approach developed here makes it possible to infer lineage-specific  
22 shifts in rates of evolution across the genome, in addition to possible differences in rates  
23 among multiple tissues sampled from the same species. We demonstrate the accuracy  
24 and robustness of our method on simulated data, and apply it to a dataset of ovule gene  
25 expression collected from multiple self-compatible and self-incompatible species in the  
26 genus *Solanum* to test hypotheses about the evolutionary forces acting during mating  
27 system shifts. These comparisons allow us to highlight the power of CAGEE,  
28 demonstrating its utility for use in any empirical system and for the analysis of most  
29 morphological traits. Our software is available at <https://github.com/hahnlab/CAGEE/>.

30

31

32

33 Key words: RNA-seq; phylogenetic comparative methods; Brownian motion; *Solanum*  
34

35 **Introduction**

36 Early studies of gene expression in single genes revealed widespread and  
37 frequent changes in the levels, timing, and breadth of expression across species  
38 (reviewed in Wray et al. 2003; Fay and Wittkopp 2008; Hill et al. 2021). Such changes in  
39 gene expression have been shown to be responsible for many differences between  
40 species, and may be a major driver of evolution (King and Wilson 1975). Advances in  
41 sequencing technologies (i.e. RNA-seq) have transformed research into gene  
42 expression, allowing researchers to cheaply and accurately measure transcript levels  
43 for every gene in a genome, in multiple tissues, and across several timepoints or  
44 conditions (Wang et al. 2009). There is now a flood of interest in applying RNA-seq to  
45 whole clades of organisms in order to identify the genetic changes and evolutionary  
46 forces driving species differences (e.g. Brawand et al. 2011; Meisel et al. 2012; Coolon  
47 et al. 2014; Harrison et al. 2015; Berthelot et al. 2018; Catalan et al. 2019; Blake et al.  
48 2020; El Taher et al. 2021).

49 To better understand the importance of changes in gene expression, researchers  
50 must be able to characterize the mechanisms and modes by which gene expression  
51 evolves. Such work entails understanding the role of natural selection in driving species  
52 differences, the stages of development or the tissues that evolve most rapidly, as well  
53 as the environments most likely to generate changes in gene expression (Dunn et al.  
54 2013; Hill et al. 2021; Price et al. 2022). Phylogenetic comparative methods enable the  
55 rigorous study of traits like gene expression across a species tree (Revell and Harmon  
56 2022). These methods can be used for testing hypotheses about natural selection, the  
57 inference of ancestral states (allowing us to polarize the direction of changes), and the  
58 estimation of evolutionary rates. Multiple software packages are available that  
59 implement a wide variety of comparative methods (e.g. Pennell et al. 2014), including  
60 models specifically intended for studying gene expression across a tree (Bedford and  
61 Hartl 2009; Rohlf et al. 2014; Rohlf and Nielsen 2015; Catalán et al. 2019; Chen et al.  
62 2019; Yang et al. 2019).

63 However, as far as we are aware, all existing comparative methods for analyzing  
64 gene expression implement fundamentally single-gene analyses. Each gene is

65 considered a separate trait, such the evolutionary parameters for each gene are  
66 estimated separately. Single-gene analyses can be used to identify tissue-specific or  
67 lineage-specific shifts in evolutionary rates, but their power is quite low (Beaulieu et al.  
68 2012). As a result, identifying trends in evolution must be carried out *post hoc* by  
69 summing the number of genes found to be individually significant (e.g. Harrison et al.  
70 2015; El Taher et al. 2021). This approach is less than ideal, especially when carrying  
71 out comparisons between branches of different lengths or between tissues with different  
72 average expression levels (both of which can result in differential statistical power).

73 Therefore, to better characterize the forces affecting gene expression evolution,  
74 we must be able to model effects shared along a lineage, experienced by many genes  
75 in the same tissue, or experienced by all genes found in the same environment. In this  
76 article, we present a genome-scale platform for the analysis of gene expression data  
77 that allows for such shared factors. Our software, CAGEE (Computational Analysis of  
78 Gene Expression Evolution), provides a robust set of methods for analyzing expression  
79 data across a species tree. CAGEE estimates ancestral states and rates, with rates  
80 shared by all or subsets of genes (single-gene analyses can also be carried out). We  
81 show that lineage-specific and tissue-specific (or condition-specific) rates can be  
82 accurately inferred, and provide principled statistical approaches for model selection.  
83 Our current implementation uses a bounded Brownian motion model and assumes  
84 expression data are accurate, but the architecture and codebase will easily allow for  
85 future extensions that relax these and other assumptions.

86

## 87 **New Approaches**

88 We model gene expression evolution as a bounded Brownian motion (BBM)  
89 process on a known species tree (cf. Boucher and Démery 2016). Our model has a  
90 single bound: trait values must be greater than or equal to zero; there is no upper bound  
91 (Figure 1). Previous researchers have often modeled gene expression using an  
92 Ornstein-Uhlenbeck (OU) process (e.g. Bedford and Hartl 2009; Rohlf et al. 2014;  
93 Rohlf and Nielsen 2015; Chen et al. 2019), a model that includes a force constraining  
94 traits about the mean. However, to our knowledge, the OU model has only been

95 compared against an unbounded Brownian motion model (i.e. one that allows negative  
96 expression values), making fair comparisons difficult. In addition, OU models may be  
97 frequently and incorrectly favored over simpler models due to several biases (e.g.  
98 measurement error), especially when the number of tips in a tree is small (Pennell et al.  
99 2015; Silvestro et al. 2015; Boucher and Démery 2016; Cooper et al. 2016; Catalán et  
100 al. 2019). Therefore, the initial version of our software models gene expression with the  
101 BBM process, which naturally bounds possible values without invoking an additional  
102 constraining force.

103 Let  $E_{ij} \geq 0$  be the expression level of gene  $i$  in species  $j$ . We assume that log-  
104 transformed expression  $X_{ij} = \ln(E_{ij} + e_{\min})$  evolves as a Brownian motion process with  
105 variance  $\sigma^2$  per unit time, where  $e_{\min}$  is a small offset (constant across genes and  
106 species) that prevents  $X_{ij}$  from taking infinite values if measured values of  $E_{ij}$  are zero.  
107 We log-transform before assuming Brownian motion because we expect the variance in  
108 the evolutionary process to scale with expression level. Assuming that  $E_{ij}$  is itself  
109 Brownian would unrealistically assume that the rate of evolution is constant across  
110 expression levels, even though expression levels vary by many orders of magnitude.  
111 We impose a reflecting lower boundary at  $x_{\min} = \ln(e_{\min})$ , meaning that the Brownian  
112 walk immediately bounces back if it reaches  $x_{\min}$ . Expression can therefore effectively  
113 never reach zero, our theoretical lower bound (Figure 1).

114 The second major feature of our model (as implemented in CAGEE) is that many  
115 genes can share the evolutionary rate parameter,  $\sigma^2$ . This rate may be shared among  
116 genes expressed in the same tissue or sample, among genes located on the same  
117 chromosome, or among genes evolving along the same lineage of the phylogenetic  
118 tree. The simplest model allows  $\sigma^2$  to be shared among all genes, providing an average  
119 rate of evolution across the genome and over time; this average may include genes that  
120 vary in their individual rates of evolution. We explain this model briefly here, with more  
121 detail provided in the Materials and Methods.

122 CAGEE infers the most likely value(s) of  $\sigma^2$  consistent with an ultrametric tree,  $T$ ,  
123 and a set  $E_{\{ij\}}$  of measured expression values at the tips of the tree; i.e. it maximizes

124 the likelihood  $L(\sigma^2 | E_{\{ij\}}, T)$ . Each gene is assumed to evolve independently, and so the  
 125 likelihood for each gene  $L_i(\sigma^2 | E_{\{ij\}}, T)$  is computed independently. The overall likelihood  
 126 is obtained as the product  $L(\sigma^2 | E_{\{ij\}}, T) = \prod_i L_i(\sigma^2 | E_{\{ij\}}, T)$  across genes. The likelihood  
 127 for each gene  $L_i(\sigma^2 | E_{\{ij\}}, T)$  is computed using the pruning algorithm (Felsenstein  
 128 1973). The key ingredient needed to apply the pruning algorithm is the transition  
 129 probability density  $p(x_t | x_{t_0}) = \Pr [X(t) = x_t | X(t_0) = x_{t_0}]$  for log-expression at time  $t$   
 130 conditional on having log-expression  $x_{t_0}$  at time  $t_0$  along a lineage. CAGEE computes  
 131 the transition density by solving the standard Brownian diffusion equation with reflecting  
 132 boundary conditions (Materials and Methods). The transition density is used to  
 133 propagate expression probabilities along the tree: if the probability density of log-  
 134 expression at time  $t_0$  is  $f(x_{t_0})$ , then the probability density at time  $t$  on the same lineage  
 135 is  $f(x_t) = \int p(x_t | x_{t_0})f(x_{t_0})dx_{t_0}$ . At each tip the probability density  $f(x_{t_0})$  is a delta  
 136 function centered at the corresponding measured value of  $X_{ij}$ .

137 Starting with the known tip distributions, the pruning algorithm propagates back  
 138 to the tips' parent nodes. The distribution at the parent node is then the product of the  
 139 two backward-propagated child node distributions. Proceeding iteratively across the  
 140 tree, we ultimately obtain the gene-specific probability density for expression value at  
 141 the root  $f_i(x_R)$ . Viewed as a likelihood for  $\sigma^2$ ,  $f_i(x_R)$  is the gene-specific likelihood  
 142 conditional on the unknown ancestral root value; i.e.  $f_i(x_R) = L_i(\sigma^2 | E_{\{ij\}}, T, x_R)$ .  
 143 Therefore, we integrate over all possible  $x_R$  to obtain,

$$144 \quad L_i(\sigma^2 | E_{\{ij\}}, T) = \int L_i(\sigma^2 | E_{\{ij\}}, T, x_R) \rho(x_R) dx_R, \quad (1)$$

145 where  $\rho(x_R)$  is the prior distribution for the root value of a randomly selected gene.  
 146 The default prior  $\rho(x_R)$  is assumed to be a gamma distribution with  $k = 0.375$  and  $\theta =$   
 147 1600, though this distribution can also be set by the user in CAGEE. This choice is  
 148 based on estimated expression distributions across genes in individual species, which  
 149 we take as our baseline for the ancestral distribution. CAGEE uses the Nelder-Mead  
 150 simplex method to find the optimal value(s) of  $\sigma^2$ .

151

152 **Results**

153 *Using CAGEE*

154 The required inputs for CAGEE are a Newick-formatted, rooted, ultrametric tree  
155 (with branch lengths) and a tab-delimited data file containing the expression levels of all  
156 species or taxa being studied. The data file can consist of data on one gene/transcript  
157 or thousands of different genes. The first line of the data file should contain the species'  
158 names (matching those used in the Newick tree). In addition, headers for gene names,  
159 gene descriptions, and sample IDs (see next section for an explanation of "samples" in  
160 CAGEE) can be used. Subsequent lines each correspond to a single gene and contain  
161 expression levels for each species. Missing data can be denoted using multiple  
162 characters (-/?/N). Examples of Newick trees and corresponding data files can be found  
163 in the online user manual  
164 ([https://github.com/hahnlab/CAGEE/docs/manual/cagee\\_manual.md](https://github.com/hahnlab/CAGEE/docs/manual/cagee_manual.md)).

165 We expect that CAGEE will most often be used to calculate the following outputs:  
166 one or more  $\sigma^2$  values, ancestral states at each internal node (including 95% credible  
167 intervals around these states), and the final likelihood associated with a model.  
168 However, users do not have to search for  $\sigma^2$ : if a value for this parameter is specified,  
169 then the output of CAGEE will just be the ancestral states and a likelihood. In addition to  
170 the raw outputs provided in multiple formats (both tab-delimited files and NEXUS-  
171 formatted files), CAGEE computes basic statistics about changes in expression levels  
172 by comparing values at parent and child nodes. Summaries of these inferred changes  
173 for every gene and for every branch of the tree are output, so that the evolutionary  
174 history of gene expression changes in every gene are accessible to users. To avoid  
175 over-interpretation of small changes in inferred expression levels—especially when  
176 there is uncertainty in ancestral states—CAGEE will also compare the credible intervals  
177 at parent and child nodes to note if a change is "credible" (i.e. the intervals do not  
178 overlap). Credible intervals are calculated by summing the probabilities across possible  
179 ancestral states at each node, so that 95% of the probability density is included.  
180 Credible changes on each branch are annotated as such in the output.

181        We most often expect that an ultrametric species tree will be used as the input  
182 topology, but this is not required by CAGEE. If users wish to specify a gene tree, or  
183 some other bifurcating tree, as input, those can be used in CAGEE as well. However,  
184 the major advantage of CAGEE—incorporating information from multiple genes to  
185 accurately estimate genome-wide rates—will rapidly diminish for trees that represent  
186 the history of only a minority of the genome. Trees that include duplication events  
187 should provide suitable estimates for any genes that follow this topology, but CAGEE  
188 does not have a way to further combine disparate gene trees.

189        There are multiple options available for running CAGEE. Users who can take  
190 advantage of multiple threads can specify the number to use on the command line.  
191 Complex models can also take a long time to converge; by default, CAGEE runs a  
192 maximum of 300 iterations of the Nelder-Mead search, but users can increase this  
193 number in subsequent runs if the likelihood is still improving when the limit is hit. As  
194 mentioned above, the default prior distribution for the root state is a gamma distribution  
195 with  $k = 0.375$  and  $\theta = 1600$ . This distribution can also be specified by the user if  
196 desired. Information on how to run more complex evolutionary models, beyond a single  
197  $\sigma^2$ , is given in the next section.

198 *Estimating evolutionary rates in CAGEE*

199        We tested CAGEE’s ability to accurately estimate  $\sigma^2$  by varying this rate  
200 parameter and the number of genes used for inference, as well as the amount of  
201 missing data in each dataset. We simulated different single values of  $\sigma^2$  across a tree  
202 with constant branch lengths (Supplementary Figure 1) using the simulation tool  
203 available within CAGEE. (Note that the total amount of evolution in a tree is determined  
204 by the product  $\sigma^2 \cdot t$ , such that changes in branch lengths will have an effect  
205 commensurate with changes in  $\sigma^2$ .) Figure 2 shows the average error associated with  
206 estimates of different  $\sigma^2$  values and using different numbers of genes within each  
207 dataset. As can be seen, the error across all parameter values and dataset sizes is  
208 quite small (generally less than 2.5%), and is less variable for larger dataset sizes.  
209 Fortunately, we expect that most empirical datasets will contain closer to 10,000 genes

210 than 1,000 genes. The results in Figure 2 are for an ancestral state vector of length  
211  $N=200$  (the default setting in CAGEE; Materials and Methods); we also estimated  $\sigma^2$   
212 when allowing the ancestral state vector to have length  $N=500$  (Supplementary Figure  
213 2A). There appears to be minimal gain from increasing the resolution in this vector,  
214 though the computational time is greatly increased (similar to results in Boucher and  
215 Démery 2016). We evaluated the accuracy of CAGEE when different amounts of data  
216 were randomly missing: from 0% to 75% for a dataset of 1,000 genes. As shown in  
217 Supplementary Figure 2B, CAGEE has high accuracy even when large amounts of data  
218 are missing (at random) from a dataset.

219 One major advantage of using CAGEE is that it combines information from  
220 multiple genes to infer a rate of evolution: this is why it can return estimates with high  
221 accuracy even when a large fraction of the data are missing. To further demonstrate this  
222 advantage, we simulated evolution in 1,000 genes using the same parameter value  
223 ( $\sigma^2=1$ ) and then estimated  $\sigma^2$  for each of the 1,000 genes individually. Supplementary  
224 Figure 2C shows that these individual estimates of  $\sigma^2$  are quite error-prone: although  
225 the mean of all genes is close to the true value, individual estimates can be 7-8X higher  
226 or lower and there is a large amount of variance. Although we have not shown it here,  
227 we do expect that the accuracy of  $\sigma^2$  will be greater for trees with larger numbers of  
228 tips, even for estimates derived from single genes (cf. O'Meara et al. 2006). On the  
229 other hand, CAGEE is combining information from multiple genes to infer an *average*  
230 rate of evolution, even when the underlying rate may be quite variable. To explore any  
231 effect of underlying rate variation, we carried out further simulations that combined three  
232 simulations of 1,000 genes each with  $\sigma^2$  equal to 0.5, 3, and 9, respectively (we  
233 repeated these simulations 10 times). When analyzed as single datasets with 3,000  
234 genes total, the average  $\sigma^2$  inferred was 3.76, approximately 9% lower than the  
235 arithmetic mean rate (Supplementary Figure 2D). It is well-known that single-rate  
236 phylogenetic likelihood models tend to underestimate rates of evolution when there is  
237 underlying variation (Golding 1983; Gillespie 1986; Yang 1996; Mendes et al. 2020),  
238 and we see this effect here. Fortunately, the bias is small, and can be corrected in the  
239 future by including gamma-distributed rate variation into CAGEE. Overall, inferences of

240  $\sigma^2$  should be quite accurate when a single rate parameter is shared across the tree and  
241 across all genes and lineages.

242 Variation in the rate of expression can currently be accommodated by CAGEE in  
243 a number of ways, using multi-rate  $\sigma^2$  models. One type of model allows users to  
244 specify that their data come from different “samples”: these samples can represent  
245 tissues, conditions, timepoints, and even subsets of the genome (e.g. the X  
246 chromosome, or a specific functional class of genes). In the input data file, the  
247 “SAMPLETYPE” column is used to indicate which sample each gene is a member of; a  
248 separate  $\sigma^2$  value will be calculated for each sample or set of samples (these values  
249 are assumed to be shared among all lineages in the tree). Specifying more than one  
250 sample means that an individual gene or transcript name can be used more than once  
251 (i.e. once for each sample), but there is no requirement that genes are measured in  
252 each sample. For instance, assigning all autosomal genes to sample 1 and all X-linked  
253 genes to sample 2 would not permit for any overlap in gene assignment, but is perfectly  
254 allowable in CAGEE.

255 Each additional sample requires another  $\sigma^2$  parameter to be estimated, and often  
256 researchers would like to know if fitting this extra parameter is justified by the data.  
257 Under standard information-theoretic criteria (Burnham and Anderson 2002), twice the  
258 difference in log-likelihoods between nested models should be  $\chi^2$ -distributed with  
259 degrees of freedom equal to the difference in the number of parameters between  
260 models. To test this expectation, we simulated 1000 datasets with a single  $\sigma^2$  value, but  
261 fit models with two  $\sigma^2$  values (assigning 1000 genes to two equal-sized samples at  
262 random; the relative size of the samples should not affect the false positive rate). As  
263 anticipated, the results fit a  $\chi^2$  distribution with one degree of freedom, with 4.4% of  
264 datasets having a difference in  $2 \times \text{log-likelihood}$  greater than 3.84 (5% are expected by  
265 chance). This indicates that standard statistical procedures should adequately control  
266 the false positive rate when fitting multi-sample  $\sigma^2$  models.

267 CAGEE also allows models in which  $\sigma^2$  varies across branches of the species  
268 tree. It does so by fitting separate  $\sigma^2$  parameters for different parts of the tree. On the

269 command line, CAGEE enables users to specify how multiple  $\sigma^2$  parameters should be  
270 assigned to branches. For  $n$  taxa, from 1 to  $2n-2$  parameters can be specified, and  
271 branches can be grouped together in any way. For instance, a two-parameter model  
272 can have all branches that share a rate adjacent to one another in the tree  
273 (Supplementary Figure 3A) or spread out across the tree (Supplementary Figure 3B).  
274 Similar to the analyses carried out above for the false positive rate associated with  
275 multiple samples, we simulated data with a single  $\sigma^2$  value and then fit models with  
276 multiple  $\sigma^2$  parameters. Regardless of how we distributed the two rate classes across  
277 the tree we observed good control of the false positive rate: 4.5% and 5.4% of 1000  
278 simulated datasets were significant at the  $P=0.05$  level (for the trees shown in  
279 Supplementary Figures 3A and 3B, respectively). More limited simulations also showed  
280 that we could accurately estimate multiple  $\sigma^2$  parameters when the data were simulated  
281 with multiple rates (Supplementary Table 1). Together, our results suggest that we can  
282 estimate multiple types of multi-rate models, and can accurately control the false  
283 positive rate when doing so.

284 *Analysis of wild tomato transcriptome data*

285 To demonstrate the utility of CAGEE in an empirical system, we analyzed data  
286 from a clade that includes domesticated tomato, *Solanum lycopersicum*. This dataset  
287 contains gene expression levels in unfertilized ovules from the flowers of six species,  
288 one of which (*S. pennellii*) has two different populations represented (Figure 3). There  
289 are 14,556 genes with expression levels measured in all seven accessions. RNA-seq  
290 data for five of the seven accessions have been published previously (Moyle et al. 2021;  
291 Hibbins and Hahn 2021), while two others are presented here for the first time  
292 (Materials and Methods). Note, however, that all data were collected from all samples at  
293 the same time (Materials and Methods).

294 Most species within the tomato clade are self-incompatible (SI), the ancestral  
295 state in the family Solanaceae (Igić et al. 2006). Self-incompatibility means that plants  
296 must outcross in order to successfully fertilize ovules. However, self-compatibility (SC)  
297 has evolved multiple times both within the Solanaceae and within the genus *Solanum*

298 (Goldberg et al. 2010; Bedinger et al. 2011). Self-compatible individuals are able to  
299 successfully fertilize ovules using their own pollen, though many also still outcross  
300 (Whitehead et al 2018; including in *Solanum*: Vosters et al. 2014 and references  
301 therein). Importantly, we have *a priori* expectations about the rate at which reproductive  
302 traits—including ovule gene expression—might evolve between groups with different  
303 mating systems. Due to conflict within and between the sexes, it is generally expected  
304 that reproductive traits in species that outcross more (i.e. SI taxa) should evolve more  
305 rapidly than in species that inbreed more (i.e. SC taxa; Clark et al. 2006). Such patterns  
306 are found in some analyses of the rate of protein evolution (e.g. Gossman et al. 2016;  
307 Harrison et al. 2019), but are equivocal in other comparisons (e.g. Gossman et al.  
308 2014, Moyle et al. 2021). These complex patterns might reflect additional effects that  
309 also accompany mating system shifts; for instance, such shifts often lead to reductions  
310 in effective population size in more selfing lineages (Charlesworth and Wright 2001).  
311 Mating system shifts could also alter global patterns of molecular evolution (including  
312 gene expression) by changing the strength and pattern of purifying selection, as  
313 morphological changes often accompany mating system changes. The exact effect of  
314 shifts in mating system on molecular evolution remains an open question.

315 The *Solanum* species sampled here represent two independent transitions from  
316 SI to SC, with one of the transitions (in accession *S. pennellii* LA0716) occurring  
317 recently enough that different populations within this species have different  
318 incompatibility systems (Figure 3). We therefore fit a series of nested models within  
319 CAGEE to test two related hypotheses about ovule gene expression evolution. First, we  
320 would like to know whether the rate of evolution of ovule gene expression is different in  
321 SI species than in SC species. Second, given the recent transition to SC within  
322 accession *S. pennellii* LA0716, we wanted to know if it shows a pattern of evolution  
323 more similar to SI or to SC species. In total, we fit four separate evolutionary models  
324 (Table 1; Figure 3). Model A has a single rate parameter for the entire tree. Model B has  
325 two rate parameters, one for SI species and one for SC species. This model assigns the  
326 branch leading to *S. pennellii* LA0716 as SC. Model C also has two rate parameters,  
327 one for SI and one for SC, but assigns *S. pennellii* LA0716 as SI. Model D has three

328 rate parameters: one for SI species, one for longer-term SC species, and one for *S.*  
329 *pennellii* LA0716.

330 Estimated results from the different models are shown in Table 1. Model A has a  
331 worse fit than any other model, with a single  $\sigma^2$  value of 0.102. For context, this value  
332 means that the bounded Brownian motion process the data are fit to has a variance of  
333 0.102 per million years (of log-transformed expression values). This is the average rate  
334 across all 14,556 genes and across all branches of the tree. In contrast to a single-rate  
335 model, both models B and C are significantly better fits to the data. Contrary to some  
336 hypotheses, both models find that SI lineages ( $\sigma_1^2$ ) have a lower rate of evolution than  
337 SC lineages ( $\sigma_2^2$ ; Table 1). There is also a difference between the models, with model C  
338 (the one in which *S. pennellii* LA0716 shares a rate with SI species) fitting significantly  
339 better. To further examine the evolution of *S. pennellii* LA0716, model D fits a three-  
340 parameter model, with this lineage assigned its own rate of evolution. This model is a  
341 significantly better fit than model C ( $P<0.00001$ ;  $\chi^2$  test with 1 degree of freedom), and  
342 demonstrates that *S. pennellii* LA0716 has a rate of evolution ( $\sigma_3^2$  in Table 1) that is  
343 slightly *lower* than SI species. This highly similar rate to SI species implies that it has  
344 only recently transitioned to self-compatibility, which is consistent with previous  
345 inferences about the timing of transition to SC in this particular accession (e.g. Rick and  
346 Tanksley 1981).

347 CAGEE also allows users to infer the number and direction of changes in gene  
348 expression levels along each branch of the tree. Figure 3 reports the number of genes  
349 that had “credible” increases and decreases in expression level under model D.  
350 Credible changes require that the credible intervals around states at parent and  
351 daughter nodes do not overlap, in order to account for uncertainty in our inferences.  
352 However, because of this, fewer credible changes will be inferred deeper in the tree,  
353 where credible intervals get wider. Therefore, while inferences about the identity of the  
354 genes changing along each branch is greatly strengthened by using credible changes  
355 (these genes are noted in the raw output from CAGEE), the absolute numbers of  
356 credible changes cannot be compared across branches, except for sister branches of

357 equal length. For completeness, we show the total numbers of increases and decreases  
358 of gene expression in Supplementary Figure 4; as expected, these total numbers are  
359 more uniformly distributed across older and younger branches.

360 We assessed whether the genes identified as having credible increases or  
361 decreases in expression specifically on any SC branch (solid red branches in Figure 3)  
362 were significantly enriched for any biological process or molecular function gene  
363 ontology (GO) categories compared to genes with credible changes on any SI branch  
364 (black branches in Figure 3). This comparison specifically assesses gene expression  
365 evolution associated with a transition to SC, over and above “background” rates of  
366 expression evolution across the rest of the clade. Although fold enrichment was modest  
367 1.20-1.36X; Supplementary Table 2), there were 11 terms significantly enriched  
368 (FDR<0.05) specifically on SC branches; these terms primarily focused on regulation of  
369 transcription, metabolic processes, and biosynthesis (Supplementary Table 2). Among  
370 the genes in these over-represented categories, a large fraction are transcription factors  
371 associated with development (e.g. WRKY and MADS Box), hormonal responses  
372 (including ethylene- and auxin-responsive transcription factors), and regulation of cell  
373 cycle (e.g. cyclins), in addition to protein kinases (Supplementary Table 2). This  
374 enrichment is consistent with increased expression changes in genes involved in cell  
375 division, differentiation, and development, that could follow transitions to SC.

376

377

378 **Discussion**

379 Here, we have developed a new software package that enables the estimation of  
380 rates of gene expression evolution across a tree, CAGEE. Gene expression levels are  
381 much like many other continuous traits, and multiple papers have introduced  
382 phylogenetic comparative methods for studying gene expression (Bedford and Hartl  
383 2009; Rohlf et al. 2014; Rohlf and Nielsen 2015; Catalán et al. 2019; Chen et al.  
384 2019). However, as far as we are aware none of these methods allows genes to share  
385 evolutionary parameters, which precludes the analysis of genome-wide trends, either  
386 along the branches of a tree or between tissues/samples/conditions. To overcome this  
387 limitation, CAGEE calculates the likelihood of the data using the pruning algorithm  
388 (Felsenstein 1973) to facilitate the sharing of evolutionary parameters along branches of  
389 the species tree, providing more statistical power to test evolutionary hypotheses.  
390 Fortunately, we were able to take advantage of much of the codebase of our existing  
391 software, CAFE (Hahn et al. 2005; De Bie et al. 2006; Hahn et al. 2007; Han et al. 2013;  
392 Mendes et al. 2020), which implements the pruning algorithm for the analysis of gene  
393 family sizes across a tree. While gene expression levels and gene family sizes differ in  
394 the type of data they represent (continuous vs. discrete) and their underlying  
395 evolutionary models (bounded Brownian motion vs. birth-death), many of the required  
396 likelihood calculations and software components are the same.

397 An important thing to consider for the input to CAGEE is the normalization used  
398 to make gene expression levels comparable across species. The data from wild  
399 tomatoes used here was normalized using TPM (transcripts per million; Wagner et al.  
400 2012); other published datasets also use this normalization (Berthelot et al. 2018; Chen  
401 et al. 2019; El Taher et al. 2021). However, multiple other normalizations have also  
402 been used in comparative analyses, including RPKM (Brawand et al. 2011), FPKM  
403 (Catalán et al. 2019), and both TMM and CPM (Blake et al. 2020). Each normalization  
404 approach has its advantages and disadvantages, and we cannot yet strongly  
405 recommend one specific approach as input to CAGEE. The normalization method used  
406 will likely depend on the conditions under which samples are collected: if all species can  
407 be raised simultaneously in a greenhouse, vivarium, or growth chamber, we expect

408 many fewer batch effects than in samples collected from the field, which will therefore  
409 necessitate different normalizations. However, even animals raised in a common  
410 environment—but fed different diets—can show many differences in gene expression  
411 not due to heritable change (e.g. Somel et al. 2008). Conversely, many between-sample  
412 normalization approaches (e.g. TMM, trimmed mean of M values; Robinson and  
413 Oshlack 2010) make the assumption that differences in gene expression between  
414 samples are rare. While such normalization is sensible in the context of testing for  
415 differential expression between samples from the same species, for a set of species  
416 that have been evolving independently for millions of years this is likely not an  
417 appropriate assumption.

418 CAGEE currently has multiple limitations, both in the available models that can  
419 be applied and in the types of data that can be analyzed. As mentioned earlier, many  
420 researchers have modeled gene expression using an OU process (Bedford and Hartl  
421 2009; Rohlf et al. 2014; Chen et al. 2019; Yang et al. 2019). Although OU models may  
422 be artifactually preferred over unbounded Brownian motion models due to a number of  
423 non-biological factors (see discussion in “New Approaches” above), it would still be  
424 helpful to be able to compare such a model to the bounded Brownian motion model  
425 used here. However, fitting such a model to genome-wide data is non-trivial: each gene  
426 must have its own mean expression value ( $\mu$ ), but possibly shared constraint  
427 parameters ( $\alpha$ ) across genes. We have the goal of implementing such a model in the  
428 near future, as well as other models commonly used in comparative methods research  
429 (e.g. Landis and Schraiber 2017; Boucher et al. 2018). Implementation of multiple  
430 models will not only allow for the analysis of different types of traits—each of which may  
431 be evolving under different regimes—but will also allow users to test the sensitivity of  
432 their analyses to model choice. For instance, it is not currently clear how different the  
433 inferred ancestral states or rates of evolution will be under different models (e.g. BBM  
434 vs. OU), and therefore how different the conclusions drawn from any such analyses  
435 might be. Ideally, qualitative results will be similar, even when there are slight  
436 quantitative differences.

437        Beyond the evolutionary model applied to any dataset, there are multiple  
438    additional sources of variation that could be modeled. For instance, we have previously  
439    accounted for measurement error in a likelihood framework, using an empirically  
440    parameterized error model (Han et al. 2013). We can imagine both applying a similar  
441    model here to RNA-seq data, as well as extending CAGEE to more error-prone data  
442    such as single-cell sequencing. Such an extension would treat the level of expression in  
443    each cell within a cell type as an error-prone draw from an underlying distribution; one  
444    would then be able to infer the rate of evolution within and across cell-types across  
445    multiple species. The biggest obstacle to this approach may be in identifying  
446    homologous cell types across species (e.g. Tarashansky et al. 2021). In addition, not all  
447    genes necessarily share the same average rate of evolution; gamma-distributed rate  
448    categories can be used to model this variation among genes (cf. Ames et al. 2012;  
449    Mendes et al. 2020). As shown above, not accounting for this rate variation leads to a  
450    slight underestimate of  $\sigma^2$ , but also obscures interesting patterns of evolution among  
451    genes. Finally, the gene tree discordance found in many phylogenomic datasets implies  
452    that complex traits (such as expression levels) will also be controlled by discordant gene  
453    trees (Hahn and Nakhleh 2016; Hibbins and Hahn 2021). This underlying discordance  
454    can cause evolutionary rates to be overestimated (Mendes et al. 2018), and should be  
455    taken into account when seeking accurate parameter estimates (see discussion of wild  
456    tomato data below). Our goal is to include methods for dealing with all these sources of  
457    variation in future versions of CAGEE.

458        In terms of the types of data that can be analyzed, at present CAGEE is limited to  
459    positive, continuously varying traits (i.e. the BBM model). However, we also envision  
460    different ways to represent and model gene expression data, including as a ratio (e.g.  
461    male/female expression). Such a ratio, after log2-transformation, would be most  
462    appropriately modeled by an unbounded Brownian motion model since both negative  
463    and positive values are possible. This and other data types will be supported in future  
464    releases. Moreover, CAGEE does not have to analyze whole-genome or even  
465    molecular data: it can be applied to any single trait for which the BBM model is  
466    appropriate, even morphological traits. One intriguing application of CAGEE could be to  
467    suites of morphological traits that are hypothesized to share a common evolutionary

468 rate parameter. If, for instance, there is a shift in body plan along some lineages, then  
469 multiple traits may all increase or decrease their rate of evolution at once, and CAGEE  
470 can be used to estimate these shared parameters. Even in the context of single-trait  
471 analyses, the pruning algorithm has been hailed as a solution for large-scale  
472 comparative analyses (Freckleton 2012). Importantly, the number of branches in a  
473 rooted, bifurcating tree with  $n$  tips is  $2n-2$ , so that the number of calculations scales  
474 linearly with the number of species. This makes the pruning algorithm ideal for  
475 comparative datasets with large numbers of taxa (e.g. Hahn et al. 2005; FitzJohn 2012;  
476 Hiscott et al. 2016; Caetano and Harmon 2018; Mitov et al. 2020).

477 The analysis of data from a clade of wild tomatoes revealed a possibly  
478 unexpected result: the rate of ovule gene expression evolution among self-compatible  
479 (SC) species is twice as high as the rate among self-incompatible (SI) species (Table  
480 1). This finding is contrary to some prior expectations— informed by research focused  
481 on male-female interactions, especially between interacting proteins in the reproductive  
482 tract (e.g. Swanson and Vacquier 2002; Clark et al. 2006)— that suggest that lineages  
483 might experience slower evolution after transitioning to self-compatibility. However, it is  
484 possible that global gene expression levels do not evolve in the same sort of tit-for-tat  
485 manner as interacting protein sequences, such that increases/decreases in male-  
486 expressed genes are not matched by increases/decreases in interacting female-  
487 expressed genes (or vice versa). Alternatively, only a very small subset of genes may  
488 evolve in this manner. Indeed, even prior studies comparing protein evolution have  
489 failed to find clear evidence of slower global evolutionary rates in more inbreeding  
490 species (e.g. Wong 2011). One caveat to the observed rate differences in our data is  
491 that underlying gene tree discordance, whether due to incomplete lineage sorting or  
492 introgression, can lead to artifactual higher rate estimates (Mendes et al. 2018;  
493 Hibbins and Hahn 2021). However, there is in fact less discordance among the SC  
494 lineages sampled here (Pease et al. 2016), which is the reverse of the pattern that  
495 would be required to explain our results.

496 If not due to underlying bias in our estimates, these findings still raise the  
497 question: why is ovule gene expression evolving more rapidly in SC than SI species?  
498 One possibility is that this increased rate is due to a relaxation of selection in SC

499 species, possibly because genes involved in male-female interactions are no longer  
500 needed. If this were the case, we might expect to see a general decrease in expression  
501 levels in SC species; however, there appears to be no consistent directionality to the  
502 changes along SC branches (Figure 3, Supplementary Figure 4). Instead, an alternative  
503 hypothesis is that transitions to SC involve adaptation to new optima of ovule gene  
504 expression, compared to SI species that tend to maintain ancestral optima. For  
505 example, transitions to SC might favor greater investment in fewer ovules, because self-  
506 compatibility decreases the probability that each ovule within a flower will go  
507 unfertilized—an otherwise wasted investment under conditions (like SI) where receiving  
508 sufficient compatible pollen to fertilize each ovule is less predictable (Burd et al. 2009).  
509 The nature of these new optima might be even more complex, as traits like ovule size  
510 and number can vary with multiple reproductive and ecological conditions, and often  
511 trade-off with each other (Greenway and Harder 2007). Of the species examined here,  
512 for example, two SC lineages (*S. pimpinellifolium*, and *S. lycopersicon*—domesticated  
513 tomato) have significantly larger seeds than most of the SI lineages and SC *S. pennellii*  
514 (unpubl. data). Indeed, individual genes identified in our GO analysis are known to  
515 directly influence ovule and/or seed size in *Solanum* (e.g. *NOR-like1*  
516 [SOLYC07G063420.3.1; Han et al, 2014], *GRAS2* [SOLYC07G063940.2.1; Li et al.  
517 2018], and *CRY2* [SOLYC09G090100.3.1; Fantini et al. 2019]). Some of our  
518 hypotheses could be evaluated with matching gene expression data from other (non-  
519 ovule) reproductive tissues. Analyses including pollen in the same SI and SC lineages,  
520 and/or data addressing alternative constraints and conditions shaping ovule evolution  
521 including ovule size and number (e.g. Mione and Anderson 1992), would be useful in  
522 teasing apart these hypotheses.

523

524

525 **Material and Methods**526 *Bounded Brownian motion model of expression evolution*

527 The probability density of expression,  $p(x, t)$ , at time  $t$  for evolutionary  
 528 trajectories following a Brownian motion process starting at value  $x_{t_0}$  at time  $t_0$  is  
 529 governed by the diffusion equation

$$530 \quad \frac{\partial p(x,t)}{\partial t} = \frac{\sigma^2}{2} \frac{\partial^2 p(x,t)}{\partial x^2}, \quad (2)$$

531 with initial condition  $p(x, t_0) = \delta(x - x_{t_0})$  where  $\delta$  is the Dirac delta function. The  
 532 reflective boundary condition at  $x = x_{\min}$  implies that the probability fluxes into and out  
 533 of the boundary are balanced, imposing the boundary condition

$$534 \quad \frac{\partial p(x=x_{\min}, t)}{\partial x} = 0. \quad (3)$$

535 Note that  $p(x, t)$  is identical to the transition density  $p(x_t | x_{t_0})$ .

536 Without the reflecting boundary,  $p(x, t) \propto e^{-(x-x_{t_0})^2/2\sigma^2(t-t_0)}$  is a normal  
 537 distribution with variance  $\sigma^2(t - t_0)$ . The variance therefore scales linearly with elapsed  
 538 time,  $t - t_0$ . With the reflecting boundary,  $p(x, t)$  is the sum of this spreading normal and  
 539 its mirror image centered at  $2x_{\min} - x_{t_0}$ . The analytical solution to this bounded process  
 540 is helpful for understanding the behavior of  $p(x, t)$ , but is not used in CAGEE. In  
 541 anticipation of implementing additional (and possibly more complicated) processes into  
 542 CAGEE, we instead solve Eq. (2) numerically using the approach described in Boucher  
 543 and Démery (2016). Briefly, the continuous diffusion equation is converted into a matrix  
 544 equation by discretizing expression values into  $N$  equal bins of width  $\delta = \frac{x_{\max} - x_{\min}}{N-1}$ .  
 545 Following Boucher and Démery (2016), we have used a default  $N=200$ , but this number  
 546 can be set by the user (see Results). This approach gives

$$547 \quad \frac{\partial P(t)}{\partial t} = \frac{\sigma^2}{2\delta^2} M \cdot P(t), \quad (4)$$

548 where  $P(t)$  is the vector obtained by discretizing  $p(x, t)$  and  $x_{\max}$  is the largest  
 549 expression value accounted for. The matrix  $M$  is tridiagonal with  $-2$  on the diagonal

550 except at the first and last diagonal entries which are  $-1$ . The sub- and supra-diagonal  
551 entries are  $1$ . This equation has the matrix exponential solution

552 
$$P(t) = \exp\left(\frac{\sigma^2(t-t_0)}{2\delta^2} M\right), \quad (5)$$

553 which is evaluated by diagonalizing  $M$ .

554 *Implementation of CAGEE*

555 CAGEE is written in C++ and is compatible with the C++11 standard. A  
556 comprehensive manual and extensive unit tests facilitate further code development and  
557 maintenance. CAGEE is organized into modular components. A *clade* class, with  
558 references to a parent clade and any number of descendant clades, represents a tree  
559 structure, and a *gene\_transcript* class represents the expression levels observed in the  
560 various species. These two classes comprise the fundamental data structures upon  
561 which CAGEE performs its analysis (Supplementary Figure 5).

562 Calculations are carried out by additional classes. The *optimizer* class has the  
563 responsibility of determining the  $\sigma^2$  value with the highest likelihood, by comparing the  
564 likelihood of candidate values and searching the likelihood surface using the Nelder-  
565 Mead optimization algorithm. The work of computing the likelihood of a given  $\sigma^2$  value is  
566 performed by a subclass of the *model* class, which for now is limited to a single *Base*  
567 model (allowing for further development in the future). After appropriate estimated  
568 values are found, the *transcript\_reconstructor* class builds a possible set of transcript  
569 values for the entire tree (Supplementary Figure 5).

570 Performing the likelihood calculations requires extensive matrix operations; it is  
571 recommended (though not required) that these be passed off to a specialized library  
572 such as Intel's MKL or Nvidia's CUBLAS. If no external library is available, CAGEE will  
573 carry out these calculations (slowly) by itself. Creating the diffusion matrix ( $M$ ) requires  
574 calculation of eigenvalues and eigenvectors, and is computationally expensive. This  
575 work is performed by the Eigen linear algebra library (<https://eigen.tuxfamily.org>);  
576 various internal data structures also take advantage of Eigen classes. To enable faster  
577 searching, the matrix for an ancestral state vector of length 200 (the default in CAGEE)

578 has been pre-computed and is included with CAGEE. Users who wish to use vectors of  
579 different lengths can specify this as an option.

580       Unit-testing is performed using the Doctest testing framework  
581 (<https://github.com/doctest/doctest> ). At the time of writing more than 200 unit tests had  
582 been created, comprising more than 1200 individual assertions. For complex logging  
583 and debugging cases, CAGEE uses the EasyLogging framework  
584 (<https://github.com/amrayn/easyloggingpp>). C++ development is always made easier by  
585 using the Boost C++ libraries (<https://www.boost.org/>), so we include them as well in  
586 CAGEE.

587 *RNA-seq data from wild tomatoes*

588       We briefly describe here the data collected from seven accessions of wild  
589 tomatoes (*S. lycopersicum* LA3475, *S. chmielewskii* LA1316, *S. pimpinellifolium*  
590 LA1589, *S. habrochaites* LA1777, *S. chilense* LA4117A, *S. pennellii* LA3778, and *S.*  
591 *pennellii* LA0716; all accession ID numbers from tgrc.ucdavis.edu). Further details are  
592 given in Moyle et al. (2021). Ovule RNA-seq was performed on between one to four  
593 (usually three) biological replicates (individual plants) from each accession. Plants were  
594 germinated from seed, and cultivated until flowering. For each replicate individual,  
595 ovules were dissected from mature, unpollinated flowers, flash frozen, and maintained  
596 at -80C until extraction. For each individual, all ovule collections were pooled into a  
597 single sample prior to library construction and sequencing on an Illumina HiSeq 2000.  
598 Reads were mapped against the tomato reference genome (iTAG 2.4) and the number  
599 of reads mapped onto genic regions were estimated with featureCounts (Liao et al.,  
600 2014). We normalized the read counts from each library by calculating TPM (transcripts  
601 per million; Wagner et al. 2012) and then calculated the mean normalized read counts  
602 across all samples (individuals) within each accession. These means per accession  
603 were used as input to CAGEE.

604       To construct a species tree for use with CAGEE, we started with the topology  
605 given in Pease et al. (2016). Specifically, we used the tree found in the supplementary  
606 file Pease\_etal\_TomatoPhylo\_RAxMLConcatTree\_no1360\_Fig2A.nwk, and pruned it to  
607 include only the accessions in our study using the software ETE (Huerta-Cepas et al.

608 2016). Using the “extend” method found in ETE, we converted this tree to ultrametric  
609 (same root-to-tip distance for all taxa). Setting the root age to 2.48 million years ago  
610 (following Pease et al. 2016), we were able to express all branches in millions of years.  
611 Analyses of GO enrichment were carried out using ShinyGO (Ge et al. 2020) with a  
612 false discovery rate of 0.05.

613

614

615 **Supplementary Material**

616 Supplementary data are available at .

617

618 **Acknowledgements**

619 We thank Mark Hibbins for assistance with the tomato phylogeny, Matthew Gibson for  
620 putting together the tomato gene expression data, and especially Dan Vanderpool for  
621 invaluable help in the initial development of CAGEE. Two reviewers provided helpful  
622 comments, and Scott Edwards pointed out relevant work that we had previously missed.  
623 This work was supported by National Science Foundation grants DEB-1856469 to  
624 L.C.M. and DBI-2146866 to M.W.H.

625

626 **Data Availability**

627 Raw reads for each sample library are available at NCBI BioProject PRJNA714065. The  
628 CAGEE software is available at <https://github.com/hahnlab/CAGEE>.

629 **References**

630 Ames RM, Money D, Ghatge VP, Whelan S, Lovell SC. 2012. Determining the  
631 evolutionary history of gene families. *Bioinformatics* 28:48-55.

632 Beaulieu JM, Jhwueng DC, Boettiger C, O'Meara BC. 2012. Modeling stabilizing  
633 selection: expanding the Ornstein–Uhlenbeck model of adaptive evolution. *Evolution*  
634 66:2369-2383.

635 Bedford T, Hartl DL. 2009. Optimization of gene expression by natural selection.  
636 *Proceedings of the National Academy of Sciences* 106:1133-1138.

637 Bedinger PA, Chetelat RT, McClure B, Moyle LC, Rose JK, Stack SM, van der Knaap E,  
638 Baek YS, Lopez-Casado G, Covey PA. 2011. Interspecific reproductive barriers in  
639 the tomato clade: opportunities to decipher mechanisms of reproductive isolation.  
640 *Sexual Plant Reproduction* 24:171-187.

641 Berthelot C, Villar D, Horvath JE, Odom DT, Flicek P. 2018. Complexity and  
642 conservation of regulatory landscapes underlie evolutionary resilience of mammalian  
643 gene expression. *Nature Ecology & Evolution* 2:152-163.

644 Blake LE, Roux J, Hernando-Herraez I, Banovich NE, Perez RG, Hsiao CJ, Eres I,  
645 Cuevas C, Marques-Bonet T, Gilad Y. 2020. A comparison of gene expression and  
646 DNA methylation patterns across tissues and species. *Genome Research* 30:250-  
647 262.

648 Boucher FC, Démery V. 2016. Inferring bounded evolution in phenotypic characters  
649 from phylogenetic comparative data. *Systematic Biology* 65:651-661.

650 Boucher FC, Démery V, Conti E, Harmon LJ, Uyeda J. 2018. A general model for  
651 estimating macroevolutionary landscapes. *Systematic Biology* 67:304-319.

652 Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, Weier M, Liechti  
653 A, Aximu-Petri A, Kircher M, et al. 2011. The evolution of gene expression levels in  
654 mammalian organs. *Nature* 478:343-348.

655 Burd M, Ashman T-L, Campbell DR, Dudash MR, Johnston MO, Knight TM, Mazer SJ,  
656 Mitchell RJ, Steets JA, Vamosi JC. 2009. Ovule number per flower in a world of  
657 unpredictable pollination. *American Journal of Botany* 96:1159-1167.

658 Burnham KP, Anderson DR. 2002. Model selection and multimodel inference: A  
659 practical information-theoretic approach. New York: Springer.

660 Caetano DS, Harmon LJ. 2018. Estimating correlated rates of trait evolution with uncertainty.  
661 *Systematic Biology* 68:412-429.

662 Catalán A, Briscoe AD, Höhna S. 2019. Drift and directional selection are the  
663 evolutionary forces driving gene expression divergence in eye and brain tissue of  
664 *Heliconius* butterflies. *Genetics* 213:581-594.

665 Charlesworth D, Wright SI. 2001. Breeding systems and genome evolution. *Current  
666 Opinion in Genetics & Development* 11:685-690.

667 Chen J, Swofford R, Johnson J, Cummings BB, Rogel N, Lindblad-Toh K, Haerty W, Di  
668 Palma F, Regev A. 2019. A quantitative framework for characterizing the  
669 evolutionary history of mammalian gene expression. *Genome Research* 29:53-63.

670 Clark NL, Aagaard JE, Swanson WJ. 2006. Evolution of reproductive proteins from  
671 animals and plants. *Reproduction* 131:11-22.

672 Coolon JD, McManus CJ, Stevenson KR, Graveley BR, Wittkopp PJ. 2014. Tempo and  
673 mode of regulatory evolution in *Drosophila*. *Genome Research* 24:797-808.

674 Cooper N, Thomas GH, Venditti C, Meade A, Freckleton RP. 2016. A cautionary note  
675 on the use of Ornstein Uhlenbeck models in macroevolutionary studies. *Biological  
676 Journal of the Linnean Society* 118:64-77.

677 De Bie T, Demuth JP, Cristianini N, Hahn MW. 2006. CAFE: a computational tool for the  
678 study of gene family evolution. *Bioinformatics* 22:1269-1271.

679 Dunn CW, Luo X, Wu Z. 2013. Phylogenetic analysis of gene expression. *Integrative  
680 and Comparative Biology* 53:847-856.

681 El Taher A, Böhne A, Boileau N, Ronco F, Indermaur A, Widmer L, Salzburger W. 2021.  
682 Gene expression dynamics during rapid organismal diversification in African cichlid  
683 fishes. *Nature Ecology & Evolution* 5:243-250.

684 Fantini E, Sulli M, Zhang L, Aprea G, Jiménez-Gómez JM, Bendahmane A, Perrotta G,  
685 Giuliano G, Facella P. 2018. Pivotal roles of cryptochromes 1a and 2 in tomato  
686 development and physiology. *Plant Physiology* 179:732-748.

687 Fay JC, Wittkopp PJ. 2008. Evaluating the role of natural selection in the evolution of  
688 gene regulation. *Heredity* 100:191-199.

689 Felsenstein J. 1973. Maximum likelihood and minimum-steps methods for estimating  
690 evolutionary trees from data on discrete characters. *Systematic Biology* 22:240-249.

691 FitzJohn RG. 2012. Diversitree: comparative phylogenetic analyses of diversification in  
692 R. *Methods in Ecology and Evolution* 3:1084-1092.

693 Freckleton RP. 2012. Fast likelihood calculations for comparative analyses. *Methods in  
694 Ecology and Evolution* 3:940-947.

695 Ge SX, Jung D, Yao R. 2020. ShinyGO: a graphical gene-set enrichment tool for  
696 animals and plants. *Bioinformatics* 36:2628-2629.

697 Gillespie JH. 1986. Variability of evolutionary rates of DNA. *Genetics* 113:1077-1091.

698 Goldberg EE, Kohn JR, Lande R, Robertson KA, Smith SA, Igić B. 2010. Species  
699 selection maintains self-incompatibility. *Science* 330:493-495.

700 Golding G. 1983. Estimates of DNA and protein sequence divergence: an examination  
701 of some assumptions. *Molecular Biology and Evolution* 1:125-142.

702 Gossmann TI, Saleh D, Schmid MW, Spence MA, Schmid KJ. 2016. Transcriptomes of  
703 plant gametophytes have a higher proportion of rapidly evolving and young genes  
704 than sporophytes. *Molecular Biology and Evolution* 33:1669-1678.

705 Gossmann TI, Schmid MW, Grossniklaus U, Schmid KJ. 2014. Selection-driven  
706 evolution of sex-biased genes is consistent with sexual selection in *Arabidopsis*  
707 *thaliana*. *Molecular Biology and Evolution* 31:574-583.

708 Greenway CA, Harder LD. 2007. Variation in ovule and seed size and associated size–  
709 number trade-offs in angiosperms. *American Journal of Botany* 94:840-846.

710 Hahn MW, De Bie T, Stajich JE, Nguyen C, Cristianini N. 2005. Estimating the tempo  
711 and mode of gene family evolution from comparative genomic data. *Genome*  
712 *Research* 15:1153-1160.

713 Hahn MW, Demuth JP, Han S-G. 2007. Accelerated rate of gene gain and loss in  
714 primates. *Genetics* 177:1941-1949.

715 Hahn MW, Nakhleh L. 2016. Irrational exuberance for resolved species trees. *Evolution*  
716 70:7-17.

717 Han MV, Thomas GWC, Lugo-Martinez J, Hahn MW. 2013. Estimating gene gain and  
718 loss rates in the presence of error in genome assembly and annotation using CAFE  
719 3. *Molecular Biology and Evolution* 30:1987-1997.

720 Han QQ, Song YZ, Zhang JY, Liu LF. 2014. Studies on the role of the *S/NAC3* gene in  
721 regulating seed development in tomato (*Solanum lycopersicum*). *The Journal of*  
722 *Horticultural Science and Biotechnology* 89:423-429.

723 Harrison MC, Mallon EB, Twell D, Hammond RL. 2019. Deleterious mutation  
724 accumulation in *Arabidopsis thaliana* pollen genes: a role for a recent relaxation of  
725 selection. *Genome Biology and Evolution* 11:1939-1951.

726 Harrison PW, Wright AE, Zimmer F, Dean R, Montgomery SH, Pointer MA, Mank JE.  
727 2015. Sexual selection drives evolution and rapid turnover of male gene expression.  
728 *Proceedings of the National Academy of Sciences* 112:4393-4398.

729 Hibbins MS, Hahn MW. 2021. The effects of introgression across thousands of  
730 quantitative traits revealed by gene expression in wild tomatoes. *PLoS Genetics*  
731 17:e1009892.

732 Hill MS, Zande PV, Wittkopp PJ. 2021. Molecular and evolutionary processes  
733 generating variation in gene expression. *Nature Reviews Genetics* 22:203-215.

734 Hiscott G, Fox C, Parry M, Bryant D. 2016. Efficient recycled algorithms for quantitative  
735 trait models on phylogenies. *Genome Biology and Evolution* 8:1338-1350.

736 Huerta-Cepas J, Serra F, Bork P. 2016. ETE 3: reconstruction, analysis, and  
737 visualization of phylogenomic data. *Molecular Biology and Evolution* 33:1635-1638.

738 Igić B, Bohs L, Kohn JR. 2006. Ancient polymorphism reveals unidirectional breeding  
739 system shifts. *Proceedings of the National Academy of Sciences* 103:1359-1363.

740 King M-C, Wilson AC. 1975. Evolution at two levels in humans and chimpanzees.  
741 *Science* 188:107-116.

742 Landis MJ, Schraiber JG. 2017. Pulsed evolution shaped modern vertebrate body sizes.  
743 Proceedings of the National Academy of Sciences 114:13224-13229.

744 Li M, Wang X, Li C, Li H, Zhang J, Ye Z. 2018. Silencing *GRAS2* reduces fruit weight in  
745 tomato. Journal of Integrative Plant Biology 60:498-513.

746 Liao Y, Smyth GK, Shi W. 2013. featureCounts: an efficient general purpose program  
747 for assigning sequence reads to genomic features. Bioinformatics 30:923-930.

748 Meisel RP, Malone JH, Clark AG. 2012. Disentangling the relationship between sex-  
749 biased gene expression and X-linkage. Genome Research 22:1255-1265.

750 Mendes FK, Fuentes-González JA, Schraiber JG, Hahn MW. 2018. A multispecies  
751 coalescent model for quantitative traits. eLife 7:e36482.

752 Mendes FK, Vanderpool D, Fulton B, Hahn MW. 2020. CAFE 5 models variation in  
753 evolutionary rates among gene families. Bioinformatics 36:5516-5518.

754 Mione T, Anderson GJ. 1992. Pollen-ovule ratios and breeding system evolution in  
755 *Solanum* section *Basarthrum* (Solanaceae). American Journal of Botany 79:279-  
756 287.

757 Mitov V, Bartoszek K, Asimomitis G, Stadler T. 2020. Fast likelihood calculation for  
758 multivariate Gaussian phylogenetic models with shifts. Theoretical Population  
759 Biology 131:66-78.

760 Moyle LC, Wu M, Gibson MJ. 2021. Reproductive proteins evolve faster than non-  
761 reproductive proteins among *Solanum* species. Frontiers in Plant Science  
762 12:635990.

763 O'Meara BC, Ané C, Sanderson MJ, Wainwright PC. 2006. Testing for different rates of  
764 continuous trait evolution using likelihood. Evolution 60:922-933.

765 Pease JB, Haak DC, Hahn MW, Moyle LC. 2016. Phylogenomics reveals three sources  
766 of adaptive variation during a rapid radiation. PLoS Biology 14:e1002379.

767 Pennell MW, Eastman JM, Slater GJ, Brown JW, Uyeda JC, FitzJohn RG, Alfaro ME,  
768 Harmon LJ. 2014. geiger v2. 0: an expanded suite of methods for fitting  
769 macroevolutionary models to phylogenetic trees. Bioinformatics 30:2216-2218.

770 Pennell MW, FitzJohn RG, Cornwell WK, Harmon LJ. 2015. Model adequacy and the  
771 macroevolution of angiosperm functional traits. The American Naturalist 186:E33-  
772 E50.

773 Price PD, Palmer Droguett DH, Taylor JA, Kim DW, Place ES, Rogers TF, Mank JE,  
774 Cooney CR, Wright AE. 2022. Detecting signatures of selection on gene expression.  
775 Nature Ecology & Evolution 6:1035-1045.

776 Revell LJ, Harmon LJ. 2022. Phylogenetic Comparative Methods in R: Princeton  
777 University Press.

778 Rick CM, Tanksley SD. 1981. Genetic variation in *Solanum pennellii*: Comparisons with  
779 two other sympatric tomato species. Plant Systematics and Evolution 139:11-45.

780 Robinson MD, Oshlack A. 2010. A scaling normalization method for differential  
781 expression analysis of RNA-seq data. *Genome Biology* 11:R25.

782 Rohlf RV, Harrigan P, Nielsen R. 2014. Modeling gene expression evolution with an  
783 extended Ornstein–Uhlenbeck process accounting for within-species variation.  
784 *Molecular Biology and Evolution* 31:201-211.

785 Rohlf RV, Nielsen R. 2015. Phylogenetic ANOVA: the expression variance and  
786 evolution model for quantitative trait evolution. *Systematic Biology* 64:695-708.

787 Silvestro D, Kostikova A, Litsios G, Pearman PB, Salamin N. 2015. Measurement errors  
788 should always be incorporated in phylogenetic comparative analysis. *Methods in  
789 Ecology and Evolution* 6:340-346.

790 Somel M, Creely H, Franz H, Mueller U, Lachmann M, Khaitovich P, Pääbo S. 2008.  
791 Human and chimpanzee gene expression differences replicated in mice fed different  
792 diets. *PLoS ONE* 3:e1504.

793 Swanson WJ, Vacquier VD. 2002. The rapid evolution of reproductive proteins. *Nature  
794 Reviews Genetics* 3:137-144.

795 Szövényi P, Ricca M, Hock Z, Shaw JA, Shimizu KK, Wagner A. 2013. Selection is no  
796 more efficient in haploid than in diploid life stages of an angiosperm and a moss.  
797 *Molecular Biology and Evolution* 30:1929-1939.

798 Tarashansky AJ, Musser JM, Khariton M, Li P, Arendt D, Quake SR, Wang B. 2021.  
799 Mapping single-cell atlases throughout Metazoa unravels cell type evolution. *eLife*  
800 10:e66747.

801 Vosters SL, Jewell CP, Sherman NA, Einterz F, Blackman BK, Moyle LC. 2014. The  
802 timing of molecular and morphological changes underlying reproductive transitions in  
803 wild tomatoes (*Solanum* sect. *Lycopersicon*). *Molecular Ecology* 23:1965-1978.

804 Wagner GP, Kin K, Lynch VJ. 2012. Measurement of mRNA abundance using RNA-seq  
805 data: RPKM measure is inconsistent among samples. *Theory in Biosciences*  
806 131:281-285.

807 Wang Z, Gerstein M, Snyder M. 2009. RNA-Seq: a revolutionary tool for  
808 transcriptomics. *Nature Reviews Genetics* 10:57-63.

809 Whitehead MR, Lanfear R, Mitchell RJ, Karron JD. 2018. Plant mating systems often  
810 vary widely among populations. *Frontiers in Ecology and Evolution* 6:38.

811 Wong A. 2011. The molecular evolution of animal reproductive tract proteins: What  
812 have we learned from mating-system comparisons? *International Journal of  
813 Evolutionary Biology* 2011:908735.

814 Wray GA, Hahn MW, Abouheif E, Balhoff JP, Pizer M, Rockman MV, Romano LA.  
815 2003. The evolution of transcriptional regulation in eukaryotes. *Molecular Biology  
816 and Evolution* 20:1377-1419.

817 Yang J, Ruan H, Xu W, Gu X. 2019. TreeExp2: an integrated framework for  
818 phylogenetic transcriptome analysis. *Genome Biology and Evolution* 11:3276-3282.

819 Yang Z. 1996. Among-site rate variation and its impact on phylogenetic analyses.  
820 Trends in Ecology & Evolution 11:367-372.

821

822

823

824 **Figures and Tables**

825 **Figure 1.** Bounded Brownian motion model. An example trait is shown in the bottom  
826 graph, evolving along the tree shown above. Although the data input to CAGEE are  
827 linear expression levels, internally it logs expression to ensure higher variance among  
828 more highly expressed genes. There is also a minimum value,  $x_{min}$ , added to all tips.

829

830 **Figure 2.** Accuracy of CAGEE. For five different values of  $\sigma^2$  we simulated 1000  
831 datasets, with each dataset comprised of either 1000 genes or 10000 genes. All genes  
832 in a dataset shared the same  $\sigma^2$ , but their values at the root were drawn independently  
833 from the prior. We then provided each simulated dataset to CAGEE in order to infer  $\sigma^2$ .  
834 Each box-and-whisker plot shows the mean (horizontal line), 50% interquartile range  
835 (box), 1.5X the interquartile range (vertical lines), and outliers (dots).

836

837 **Figure 3.** Changes in gene expression along the tomato phylogeny. Given the set of  
838 relationships among the seven *Solanum* accessions used here, we tested multiple  
839 models that had branches assigned as different  $\sigma^2$  parameters (Table 1). In model A, all  
840 branches share  $\sigma_1^2$ . In model B, all black branches share  $\sigma_1^2$ , while all red branches  
841 share  $\sigma_2^2$ . In model C, all black branches and the dashed red branch share  $\sigma_1^2$ , while all  
842 solid red branches share  $\sigma_2^2$ . In model D, all black branches share  $\sigma_1^2$ , all solid red  
843 branches share  $\sigma_2^2$ , and the dashed red branch is assigned  $\sigma_3^2$ . Using the results from  
844 model D, we inferred the number of genes that had credible increases or decreases in  
845 expression level along each branch (results for all changes are shown in Supplementary  
846 Figure 4). Numbers are reported as +increases/-decreases for each branch.

847 **Table 1.** Model parameters estimated from the tomato data.  
848

Model	Number of rates	-lnL	$\sigma^2_1$	$\sigma^2_2$	$\sigma^2_3$
A	1	67252.4	0.102		
B	2	65883.9	0.074	0.134	
C	2	65124.5	0.075	0.152	
D	3	65108.6	0.077	0.152	0.067

849

850

851 **Supplementary Figures**

852 **Supplementary Figure 1.** The tree used for simulations. The Newick-formatted tree  
853 string with branch lengths is:

854 `((sp1:1,sp2:1):1,sp3:2):1,sp4:3):1,((sp5:2,sp6:2):2):1,sp7:5):1,sp8:6)`

855

856 **Supplementary Figure 2.** Accuracy of CAGEE. A) This figure is the same as Figure 2  
857 in the main text, but the ancestral state vector has length  $N=500$  (Figure 2 uses  $N=200$ ).  
858 B) For each of three different simulated values of  $\sigma^2$ , we randomly removed different  
859 amounts of data from an input dataset with 1,000 genes (the tree is the same as in all  
860 other simulations). C) For 1,000 genes simulated with  $\sigma^2=1$  (dashed vertical line), we  
861 ran CAGEE independently on each one to estimate  $\sigma^2$ . D) We combined three datasets  
862 of 1,000 genes each simulated with three different values of  $\sigma^2$  (we repeated these  
863 simulations 10 times). The 10 estimates of  $\sigma^2$  on the combined datasets were slightly  
864 downwardly biased compared to the expected value (dashed horizontal line). Each dot  
865 represents each of the 10 estimates, with jitter added for clarity,

866

867 **Supplementary Figure 3.** Trees used for simulations with lineage-specific values of  $\sigma^2$ .

868 A) All black branches share a rate parameter ( $\sigma_1^2$ ), and all red branches share a rate  
869 parameter ( $\sigma_2^2$ ). This “sigma\_tree” is specified in CAGEE with the Newick string:  
870 `((sp1:2,sp2:2):2,sp3:2):2,sp4:2):2,((sp5:1,sp6:1):1):1,sp7:1):1,sp8:1)`

871 B) All black branches share a rate parameter ( $\sigma_1^2$ ), and all red branches share a rate  
872 parameter ( $\sigma_2^2$ ). This “sigma\_tree” is specified in CAGEE with the Newick string:  
873 `((sp1:2,sp2:1):1,sp3:2):1,sp4:1):1,((sp5:1,sp6:2):1):1,sp7:1):1,sp8:1)`

874

875 **Supplementary Figure 4.** Changes in gene expression along the tomato phylogeny.

876 This figure is the same as Figure 3 in the main text, but all increases and decreases are  
877 reported, regardless of whether they are “credible”.

878

879 **Supplementary Figure 5.** Component diagram for the CAGEE software.