

Detecting macroevolutionary genotype-phenotype associations using error-corrected rates of protein convergence

Kenji Fukushima (✉ kenji.fukushima@uni-wuerzburg.de)

University of Würzburg

David Pollock

University of Colorado

Article

Keywords: genotype-phenotype associations, convergent evolution, nonsynonymous per synonymous substitution rate ratio

Posted Date: April 6th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1509769/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Cover Page**

2

3 **Article Title:**

4 Detecting macroevolutionary genotype-phenotype associations using error-corrected rates of protein
5 convergence

6

7 **Short Title:**

8 Phylogenetic error correction of convergence

9

10 **Authors:**

11 Kenji Fukushima^{1,*} and David D. Pollock²

12

13 **Affiliations:**

14 ¹ Institute for Molecular Plant Physiology and Biophysics, University of Würzburg

15 ² Department of Biochemistry and Molecular Genetics, University of Colorado, School of Medicine, Aurora,
16 CO 80045, USA.

17

18 * Corresponding author: Kenji Fukushima

19

20 Kenji Fukushima

21 **Mailing address:** Julius-von-Sachs Platz 2, Würzburg, Germany, 97072

22 **Phone number:** +4915159128431

23 **E-mail address:** kenji.fukushima@uni-wuerzburg.de

24 **ORCID ID:** <https://orcid.org/0000-0002-2353-9274>

25

26 David D. Pollock

27 **Mailing address:** Department of Biochemistry & Molecular Genetics, Mail Stop 8101, Research Complex
28 1 South, 12801 17th Ave., Room #10111, PO Box 6508, Aurora, CO 80045 USA.

29 **Phone number:** +1-303-724-3234

30 **FAX number:** +1-303-724-3215

31 **E-mail address:** david.pollock@cuanschutz.edu

32 **ORCID ID:** <https://orcid.org/0000-0002-7627-4214>

33

34 **Keywords:** genotype-phenotype associations, convergent evolution, nonsynonymous per synonymous
35 substitution rate ratio

36

37 **Abstract**

38 On macroevolutionary timescales, extensive mutations and phylogenetic uncertainty mask the signals
39 of genotype-phenotype associations underlying convergent evolution. To overcome this problem, we
40 extended the widely used framework of nonsynonymous-to-synonymous substitution rate ratios and
41 developed the novel metric ω_C , which measures the error-corrected convergence rate of protein evolution.
42 While ω_C distinguishes natural selection from genetic noise and phylogenetic errors in simulation and real
43 examples, its accuracy allows an exploratory genome-wide search of adaptive molecular convergence
44 without phenotypic hypothesis or candidate genes. Using gene expression data, we explored over 20 million
45 branch combinations in vertebrate genes and identified the joint convergence of expression patterns and
46 protein sequences with amino acid substitutions in functionally important sites, providing hypotheses on
47 undiscovered phenotypes. We further extended our method with a heuristic algorithm to detect highly
48 repetitive convergence among computationally nontrivial higher-order phylogenetic combinations. Our
49 approach allows bidirectional searches for genotype-phenotype associations, even in lineages that diverged
50 for hundreds of millions of years.

51 **Introduction**

52 A central aim of modern biology is to differentiate the huge amount of nonfunctional genetic noise
53 from phenotypically important changes. Evolutionary processes at the molecular level are largely neutral
54 and stochastic, but natural selection can constrain evolutionary pathways available to the organism. If similar
55 environmental conditions recur in divergent lineages, the adaptive response may also be similar, leading to
56 convergence, the repeated emergence of similar features in distantly related organisms (Losos, 2017). The
57 prevalence of phenotypic convergence is demonstrated by various examples throughout the tree of life, such
58 as the camera eyes of vertebrates and cephalopods, powered flight of birds and bats, and trap leaves of
59 distantly related carnivorous plants. Because the repeated emergence of such complex traits by neutral
60 evolution alone is extremely unlikely, convergence at the phenotypic level is considered strong evidence for
61 natural selection.

62 Phenotypic convergence is necessarily caused by molecular events and often coincides with
63 detectably excess levels of convergent molecular changes in gene regulation, gene sequences, gene
64 repertoires, and other hierarchies of biological organization (Stern, 2013; Storz, 2016). A meta-analysis
65 reported that 111 out of 1,008 loci had been convergently modified to attain common phenotypic
66 innovations, sometimes even between different phyla (Martin and Orgogozo, 2013), demonstrating that
67 genotype-phenotype associations frequently occur on macroevolutionary scales. For example, several
68 lineages of mammals, reptiles, amphibians, and insects acquired resistance to toxic cardiac glycosides using
69 largely overlapping sets of amino acid substitutions in a sodium channel (Ujvari et al., 2015). Another
70 example illustrated how human cancer cells and plants employed common amino acid substitutions in
71 Topoisomerase I to cope with a common toxic cellular environment generated by plant-derived anticancer
72 drugs (Sirikantaramas et al., 2008).

73 Genome sequences are becoming more available for diverse lineages from the entire tree of life
74 (Lewin et al., 2022), making it possible to explore macroevolutionary genotype-phenotype associations on
75 large scales. However, because most molecular changes are nearly neutral or essentially nonfunctional in
76 nature (Ohta, 1973), false-positive convergence in the form of stochastic, nonadaptive, convergent events is
77 particularly problematic when conducting a genome-scale search. Furthermore, false positives can arise from
78 methodological biases. For molecular convergence, a major source of bias occurs because such inference is
79 sensitive to the topology of the phylogenetic tree on which substitution events are placed (Mendes et al.,
80 2016) (Fig. 1A), while alternative methods that do not place substitutions on phylogenetic trees suffer even
81 more severe rates of false positives (Foote et al., 2015; Thomas and Hahn, 2015; Zou and Zhang, 2015b)
82 (Supplementary Text 1). A correctly inferred tree avoids false positives due to phylogeny (Castoe et al.,
83 2009), but topological misinference due to technical errors, insufficient data, or biological factors such as
84 introgression, horizontal gene transfer (HGT), paralogy, incomplete lineage sorting, and within-locus
85 recombination, can all create substantial amounts of false convergence signals even when adaptive
86 convergence did not actually occur (Mendes et al., 2016, 2019; Stern, 2013; Thomas et al., 2017).
87 Importantly, false convergence events driven by topological errors tend to similarly affect both
88 nonsynonymous and synonymous substitutions (Fig. S1A). By contrast, truly adaptive convergence should
89 occur almost exclusively in nonsynonymous substitutions (amino acid-changing substitutions), as positive
90 selection on synonymous substitutions is negligible, or at least not prevalent (Yang, 2006) (Fig. S1B).
91 Therefore, synonymous convergence can potentially serve as a reliable reference for measuring the rate of
92 expected nonsynonymous convergence due to phylogenetic inference error.

93 A widely used framework to understand how functionally constrained proteins evolve compared to
94 neutral expectations is to contrast rates of nonsynonymous and synonymous substitutions. The ratio of these
95 rates within a protein-coding sequence accounts for mutation biases and is often denoted as ω , dN/dS , or
96 K_a/K_s (Zhang and Yang, 2015). Here, we extend this framework to derive the new metric ratio ω_c and
97 implement it to measure phylogenetic error-corrected rates of convergence. Simulation and empirical data
98 analysis show that this new metric has high sensitivity while suppressing false positives. We further show
99 its capability to detect factors that affect protein convergence rates and to identify likely adaptive protein
100 evolution in a genome-scale dataset by an exploratory analysis without a pre-existing hypothesis. We also

101 develop a heuristic algorithm to explore convergent signals with high signal-to-noise ratios in exponentially
102 increasing numbers of higher-order phylogenetic combinations.
103

104 **Results**

105

106 **Extending the framework of nonsynonymous per synonymous substitution rate ratio to molecular**
107 **convergence.** One of the most commonly accepted measures of the rate of protein evolution compared to
108 neutral expectations is the ratio between nonsynonymous and synonymous substitution rates, denoted as dN
109 and dS , or K_a and K_s , respectively (Yang, 2006). Using the ratio dN/dS to measure relative rates of protein
110 evolution is justified, as the selective pressure on synonymous sites is negligible compared to that on
111 nonsynonymous sites and thus remains fairly constant relative to the mutation rate (Yang, 2006). In a model-
112 based framework, this ratio is parameterized as ω .

113 Inspired by ω , we developed a similar metric, ω_C , that applies to substitutions that occurred
114 repeatedly on a combination of separate phylogenetic branches (combinatorial substitutions; Fig. S1C;
115 Supplementary Text 2). The metric ω_C estimates the relative rates of convergence obtained by contrasting
116 the rates of nonsynonymous and synonymous convergence (dN_C and dS_C , respectively). Using this ratio,
117 important biological fluctuations, such as among-site rate heterogeneity and codon equilibrium frequencies,
118 are taken into account (for details, see Supplementary Text 3 and Methods). Similar to previously proposed
119 convergence metrics (Castoe et al., 2009; Goldstein et al., 2015; Zou and Zhang, 2015a), ω_C is calculated
120 from substitutions at multiple codon sites across protein-coding sequences. As a result, one ω_C value is
121 obtained for each gene for each branch pair (or for a combination of more than two branches) in the
122 phylogenetic tree. A unique feature of ω_C setting it apart from other metrics is its error tolerance. For
123 example, if one of the branches in a branch combination is in error, ω_C is a measure of the ratio of false
124 convergence events of both kinds falsely attributed to a non-existent branch combination. In this way, the
125 ω_C values remain close to the neutral expectation of 1.0, even when topology errors are involved. Our method
126 is implemented in the python program CSUBST (<https://github.com/kfuku52/csubst>), which takes as input
127 a rooted phylogenetic tree and a codon sequence alignment (Fig. 1B and Fig. S2).

128

129 **The robustness of ω_C as a relative rate of molecular convergence.** Conventionally, observed levels of
130 convergent amino acid substitutions have been contrasted either to the amount of convergence expected
131 under a neutral model with no constraint (R (Zou and Zhang, 2015a)) or to other combinations of amino acid
132 substitution patterns that are similarly affected by site-specific constraint (i.e., double divergence; C/D
133 (Castoe et al., 2009; Goldstein et al., 2015)) (Table S1; Supplementary Text 4). The metric R , for example,
134 is intended to have an expectation of 1.0 under neutral evolution, but in practice is somewhat lower than 1.0,
135 even when the tree and substitution model are correct and exactly match simulation conditions (Zou and
136 Zhang, 2015a). Using $R > 1.0$ as a criterion to identify convergence is thus in principle conservative for
137 detecting convergence levels greater than fully neutral evolution. Furthermore, its accuracy depends on the
138 accuracy of the phylogenetic tree in various aspects, e.g., neutral substitution model, tree topology, branch
139 lengths, and reconstructed ancestral states. By contrast, the C/D comparison ratio, which compares
140 convergence levels to double divergence events between branch pairs, is not strongly dependent on neutral
141 substitution estimates (Castoe et al., 2009; Goldstein et al., 2015); however, it is dependent on the accuracy
142 of the reconstructed tree compared to the true tree that applies. The C/D ratio may vary among proteins due
143 to varying levels of constraint among proteins but is generally well below 1.0 (Goldstein et al., 2015).

144 Here, we focus on whether ω_C performs better as a measure of convergence between branches in
145 comparison to alternative metrics. Accordingly, we generated simulated sequences with 500 codons along a
146 balanced phylogenetic tree ending with 32 sequences at the tips (or leaves), in all cases comparing two deeply
147 separated tip lineages (shown as dots in Fig. 1C; Table S2). In this analysis, we compared C/D , dN_C , dS_C ,
148 and ω_C under four evolutionary scenarios of relationships between the two tips being compared: 1) full
149 neutral evolution along all branches (Neutral); 2) neutral evolution for nearly all branches but with
150 convergent selection along the two deeply separated tip lineages (Convergent); 3) neutral evolution with
151 phylogenetic tree topology error in the form of a copy-and-paste transfer from one of the two deeply
152 separated lineages to the other, overwriting its genetic information (Transfer); or 4) neutral evolution but

153 using a randomly reconstructed phylogenetic tree to detect convergence (Random). The metric dN_C is
154 obtained by dividing the observed value of nonsynonymous convergence (O_C^N) by the expected value (E_C^N)
155 and is essentially equivalent to the previously proposed metric called R (Zou and Zhang, 2015a), but we use
156 the dN_C notation here to clarify its relationship to dS_C , the ratio of observed to expected values of
157 synonymous convergence (O_C^S/E_C^S).

158 During neutral evolution, sequences evolved under a constant codon substitution model without any
159 adaptive convergence or constraint on amino acid substitutions other than those imposed by the structure of
160 the genetic code and relative codon frequencies. In the Neutral scenario (Fig. 1C), the trees used for
161 simulation and reconstruction were identical. C/D was much lower than 1.0, as expected, while the other
162 three metrics (dN_C , dS_C , and ω_C) were close to but lower than the neutral expectation of 1.0 (Fig. 1D). This
163 observation is likely due to the fact that the convergent events must be inferred and are not actually observed,
164 as investigated previously in R (Zou and Zhang, 2015a). In the Convergent scenario, adaptive convergence
165 on the focal pair of deeply separated branches (red branches in Fig. 1C) was mimicked by convergently
166 evolving 5% of codon sites (25 sites) in the two branches under substitution models biased toward codons
167 encoding the same randomly selected amino acid. This generated an average of four excess nonsynonymous
168 convergent substitutions on these two branch pairs (see O_C^N in Fig. 1C). In the Convergent scenario, the three
169 protein convergence metrics, C/D , dN_C , and ω_C , yielded values substantially higher than they did under the
170 Neutral scenario, while the synonymous change measure dS_C remained comfortably well below 1.0. Using
171 the distribution of metric values under the Neutral scenario as a reference, we see that 70–80% of the
172 detection metric values in the Convergent scenario are above the 95th percentile of the 1,000 simulations in
173 their respective neutral distributions, while only 3.5% of dS_C values are above this threshold, indicating that
174 this level of convergence is usually detected by all three of the protein convergence metrics (Fig. 1D). To be
175 thorough, we considered that ω_C metrics can in general be derived for nine types of combinatorial
176 substitutions (i.e., substitutions occurring at the same protein site in multiple independent branches;
177 Supplementary Text 2) based on whether the ancestral and descendant states are the same or different, or in
178 any state among multiple branches (Fig. S1C). In the Convergent scenario, only the ω_C metrics involved in
179 convergence (i.e., not divergence) showed a response, confirming its specificity (Fig. S3A).

180 We next considered Transfer and Random scenarios that include phylogenetic error. In the Transfer
181 scenario, we transferred one of the focal tip sequences to the other focal tip sequence in the simulation, but
182 the phylogenetic tree used in the analysis remained unchanged, as might happen with HGT events (Fig. 1C).
183 In the Random scenario, we fully randomized the entire reconstructed tree relative to the true tree (Fig. 1C).
184 Excess convergence detected in either of these scenarios is considered a false positive. We determined that
185 both C/D and dN_C are sensitive to the errors (Fig. 1D). By contrast, and as intended, ω_C values were close
186 to the neutral expectation because the rise in dN_C due to phylogenetic error is matched by a similar increase
187 of dS_C , and they cancel each other out in the ω_C metric (Fig. 1D). Further simulations supported the
188 robustness of ω_C against the rate of protein evolution, model misspecification, tree size, and protein size
189 (Fig. S3B–F; Supplementary Text 5). Furthermore, ω_C showed low false-positive rates in sister branches
190 that serve as a control for the focal branch pairs (Foote et al., 2015) (Fig. S3G). Taken together, our
191 simulation showed that ω_C effectively counteracts false positives caused by phylogenetic errors without loss
192 of power.

193
194 **ω_C distinguishes between adaptive and false convergence in empirical datasets.** To test whether ω_C
195 performs well with real data, we collected protein-coding sequence datasets from known molecular
196 convergence events in various pairs of lineages covering insects, tetrapods, and flowering plants (Fig. 1E,
197 Fig. S4, Fig. S5, and Table S3). Insects that feed on milkweed (Apocynaceae) harbor amino acid
198 substitutions in a sodium pump subunit (ATPalpha1) that confer cardiac glycoside resistance (Dobler et al.,
199 2012; Yang et al., 2019a; Zhen et al., 2012) (Fig. S4A). Echolocating bats and whales share amino acid
200 substitutions in the hearing-related motor protein Prestin to enable high-frequency hearing (Liu et al., 2010,
201 2014) (Fig. S4B). An extensive molecular convergence occurred in the mitochondrial genomes of agamid

202 lizards and snakes, presumably due to physiological adaptations for radical fluctuations in their aerobic
203 metabolic rates (Castoe et al., 2009). Specialized digestive physiology of herbivorous mammals and
204 carnivorous plants led to the molecular convergence of digestive enzymes (Fukushima et al., 2017; Stewart
205 et al., 1987; Zhang, 2006; Zhang and Kumar, 1997) (Fig. S4C–G). Phosphoenolpyruvate carboxylase
206 (PEPC), a key enzyme for carbon fixation in C_4 photosynthesis, shares multiple amino acid convergence
207 (Besnard et al., 2009; Christin et al., 2007) (Fig. S4H). In all these examples, ω_c successfully detected
208 convergent lineages, while it was always lower and in many cases close to the neutral expectation in the
209 branch pairs sister to the focal lineages, which serve as a negative control (Fig. 1E; Table S4). Moreover, the
210 ω_c values of the focal branch pairs tended to be high compared to background levels in the phylogenetic
211 trees (Fig. S4I). Analysis of different categories of combinatorial substitutions correctly recovered a trend
212 consistent with the action of intramolecular epistasis, which did not appear in the simulations
213 (Supplementary Text 6; Fig. S4J–K).

214 To test robustness against phylogenetic errors, we also employed reported cases of HGTs associated
215 with C_4 photosynthesis (Dunning et al., 2019) and plant parasitism (Yang et al., 2019b). We reconstructed
216 the phylogenetic trees of the HGT genes with a constraint that enforces species tree-like topologies (Fig. S6).
217 This operation separates the HGT donor and acceptor lineages and creates false convergence (Fig. S1A).
218 Consistent with the simulation results, ω_c values in HGTs were lower than the adaptive convergence events
219 (Fig. 1E). By contrast, C/D and dN_c showed values higher in HGTs than in the adaptive convergence events.
220 Together with the simulations, these results show that the consideration of synonymous substitutions is
221 essential for the accurate detection of molecular convergence in the presence of phylogenetic error and that
222 ω_c outperforms current alternative methods.

223
224 **Temporal variation of convergence rates.** The probability of protein convergence decreases over time,
225 with intramolecular epistasis among amino acid residues considered to be a primary biological source of
226 such an evolutionary pattern (Goldstein et al., 2015; Zou and Zhang, 2015a; Goldstein and Pollock, 2017).
227 Indeed, over a long timescale, the environment around any given focal site changes through substitutions at
228 other amino acid sites, thus altering which amino acid state at the focal site is suitable to maintain structure
229 and function (Goldstein and Pollock, 2017; Pollock et al., 2012) (Fig. S4L). However, gene tree discordance
230 due to biological and technical causes, including tree inference error, incomplete lineage sorting,
231 introgression, HGT, and intralocus recombination, can create a false convergence signal that similarly
232 decreases with the time since branches separated (Mendes et al., 2016, 2019) (Fig. 1A and Fig. S1A). While
233 the analysis of the mitochondrial genome (Goldstein et al., 2015) would not have been confounded by
234 recombination-mediated mechanisms, other factors would have as great an influence as for nucleus-encoded
235 genes. Nevertheless, all of the above problems would produce false convergence signals equally in
236 synonymous and nonsynonymous substitutions via errors in the phylogenetic tree topology; therefore, ω_c
237 should be a natural candidate to unbiasedly evaluate whether convergence rates in nucleus-encoded genes
238 also decrease with time.

239 We obtained 21 vertebrate genomes covering a range from fish to humans (Fig. 2A and Fig. S7A) and
240 calculated ω_c for all independent branch pairs in 16,724 orthogroups classified by OrthoFinder (Emms and
241 Kelly, 2015, 2019). CSUBST completed the analysis even for the largest orthogroup (OG0000000),
242 containing 682 genes encoding zinc finger proteins and 901,636 independent branch pairs (alignment length
243 including gaps: 31,665 bp). We obtained a total of 20,150,538 branch pairs from all orthogroups and further
244 analyzed 2,349,515 branch pairs with at least one synonymous and nonsynonymous convergence (i.e., $O_c^N \geq$
245 1.0 and $O_c^S \geq 1.0$). In all metrics (C/D , dN_c , and ω_c), protein convergence rates clearly decreased over time
246 (approximated by inter-branch genetic distance) (Fig. 2B). Notably, we observed no such pattern for the rate
247 of synonymous convergence (dS_c), making it more likely that the diminishing protein convergence is caused
248 by evolutionarily selected mechanisms (Goldstein et al., 2015; Zou and Zhang, 2015a). We also detected a
249 similarly decreasing pattern in the rates of divergent substitutions over time, which does not contradict the
250 effect of epistasis (Fig. S7B–C; Supplementary Text 7). Thus, the pattern of diminishing convergence

251 remains a clear trend in recombining nucleus-encoded genes, even after correcting for the rate of
252 synonymous convergence, and therefore is consistent with the action of intramolecular epistasis (Fig. S4L).

253

254 **Gene duplication decreases convergence rates.** Gene duplication generates new genetic building blocks
255 (Conant and Wolfe, 2008) and elevates the rate of protein evolution (Fukushima and Pollock, 2020).
256 However, it remains unknown whether substitution profile changes influence convergence rates following
257 gene duplication. Convergent substitutions in duplicates may indicate convergent functional changes in
258 independently duplicated genes, and our genome-scale dataset contains 90,028 duplication events, providing
259 an excellent opportunity to address this question. If independent duplications in a family of genes tend to
260 result in mutually similar derived pairs of proteins, the convergence rate should increase. Conversely, if the
261 new proteins tend to move into a divergent sequence space in which they do not overlap, gene duplication
262 would not increase convergence and may even decrease it. Accelerated non-adaptive change might not
263 change the convergence rate if gene duplication only causes an increase in the rate of protein evolution
264 without changing the substitution profiles. To distinguish these possibilities, we compared the convergence
265 rates of branch pairs after two separate speciation (SS) events and branch pairs after two independent gene
266 duplications (DD) (Fig. 2C). Strikingly, gene duplication significantly decreased convergence rates ($P \approx 0$,
267 $W = 23.0$, as determined by a two-sided Brunner–Munzel test; Fig. 2C). Again, the trend was evident in
268 nonsynonymous convergence (dN_C) but not in synonymous convergence (dS_C), implying a relaxation in
269 site-specific constraints or adaptive divergence in the duplicates. Notably, the effect of gene duplication was
270 stronger in closely related branch pairs (i.e., smaller bin numbers in Fig. 2C), and the ω_C distributions
271 became progressively indistinguishable between SS and DD pairs with increasing inter-branch distance. The
272 immediate drop of the convergence probability was consistent with the idea that gene duplication allows the
273 new gene copies to explore a new sequence space, potentially involving natural selection. We note that this
274 is an averaged trend across genes and does not exclude possible adaptive convergence in some genes.
275 However, it is likely that such convergence, if it does exist, is masked by the opposing, predominant signal
276 of relaxed or divergent constraints.

277 It is also noteworthy that the DD branch pairs show anomalously high synonymous convergence rates
278 (dS_C) in the smallest bin of genetic distance (bin 1 in Fig. 2C). This observation is probably due to the
279 difficulty of locating gene duplication events in the phylogenetic tree, especially when sequences are not
280 sufficiently diverged and lead to an extremely short branch length. Consistent with this idea, small genetic
281 distances were associated with low branch supports in the DD branch pairs (Fig. S7D). Additionally, we
282 detected similar anomalies in extremely distant branch pairs and attributed them to false orthogroup
283 inference (Supplementary Text 8; Fig. S7E). These examples illustrate how various aspects of phylogenetic
284 analysis can generate false patterns of convergence that are successfully captured by dS_C and corrected for
285 in ω_C .

286

287 **Extracting a high-confidence set of convergent lineages.** Discovering adaptive molecular convergence in
288 genome-scale datasets, which may be translated into genotype-phenotype associations, has been challenging
289 since it is a rare phenomenon and false positives are high (Foote et al., 2015; Thomas and Hahn, 2015; Zou
290 and Zhang, 2015b). To examine whether the application of ω_C can generate plausible hypotheses of adaptive
291 molecular convergence, we analyzed the 21 vertebrate genomes (Fig. S7A). We first extracted the branch
292 pairs with the top 1% of C/D , dN_C , or ω_C values with a cutoff for a minimum of three nonsynonymous and
293 synonymous convergence ($O_C^N \geq 3.0$ and $O_C^S \geq 3.0$) (Fig. 3A). The overlap between each set of branch pairs
294 was moderate, with 1,348 branch pairs satisfying all three criteria out of 5,659 pairs with the top 1% ω_C
295 values.

296 To examine which metrics better enrich for likely adaptive convergence, we compared the topological
297 confidence scores of the selected branches. If artifacts due to tree topology errors are included, low
298 confidence branches should be enriched. Analysis of the bootstrap-based confidence values (Hoang et al.,
299 2018; Minh et al., 2013) showed that ω_C selects branch pairs with higher confidence than the other two
300 metrics (Fig. 3A). Furthermore, we examined the synonymous convergence rate (dS_C), which is not expected

301 to be greater than the neutral expectation in the adaptive convergence, and established that only ω_C satisfies
302 such an assumption (Fig. 3A). These results indicate that ω_C has excellent properties for finding adaptive
303 protein convergence in genome-scale analyses.

304

305 **Identification of molecular convergence associated with a particular phenotype.** As convergence
306 metrics have been used to search for genes associated with phenotypes of interest, we next examined whether
307 ω_C might be used to discover candidate genes underlying phenotypic convergence. Here, we analyzed a pair
308 of herbivorous animal lineages as an example: ruminants (cattle [*Bos taurus*] and red sheep [*Ovis aries*]) and
309 rabbits (*Oryctolagus cuniculus*). Using minimum thresholds for the number of convergent amino acid
310 substitutions ($O_C^N \geq 3.0$) and protein convergence rate ($\omega_C \geq 3.0$), we obtained 352 candidate branch pairs
311 in a genome-scale analysis of the 21 vertebrates (Table S5). By mapping the positions of substitutions onto
312 known conformations of homologous proteins, we identified particularly compelling cases of likely adaptive
313 convergence (Fig. S8). Examples included olfactory receptors in which convergent substitutions are located
314 in the interior of the receptor barrel (ODORANT RECEPTOR 7A [OR7A], Olfactory Receptor Family 2
315 Subfamily M Member 2 [OR2M2], and OR1B1), where substitutions may change ligand preference
316 associated with herbivorous behavior.

317 Similarly, the barrel-like structure of some solute carriers harbored convergent substitutions in their
318 interior sides (Solute Carrier Family 5 Member 12 [SLC5A12], SLC51A, SLC22A, and SLC44A1),
319 suggesting their involvement in the uptake or transport of plant-derived compounds. Among these, SLC51A
320 (also known as Organic solute transporter α [OST α]) may be a particularly attractive candidate. This protein
321 plays a major role in bile acid absorption and, hence, in dietary lipid absorption (Ballatori et al., 2005). The
322 convergence in SLC51A may be coupled with another convergent event detected in CYP7A1, a cytochrome
323 P450 protein known to serve as a critical regulatory enzyme of bile acid biosynthesis (Chiang and Ferrell,
324 2020). CYP7A1 harbored two convergent substitutions in its substrate-binding sites (Fig. S8). While most
325 herbivores secrete bile acids mainly in a glycine-conjugated form, ruminant bile is mostly in the form of
326 taurine-conjugated bile acids, which remain soluble in highly acidic conditions (Noble, 1981). The
327 predominance of taurine-conjugated forms is also observed in rabbits, depending on species and
328 developmental stage (Hagey et al., 1998). Thus, convergence in these proteins may be related to such
329 nutritional physiology.

330 Other examples of detected convergence included two convergent substitutions in the DNA-binding
331 sites of a member of the zinc-finger protein family, which functions as a transcriptional regulator (Patel et
332 al., 2018) (Fig. S8). Convergence in the substrate-binding sites of pancreatic elastase (Mulchande et al.,
333 2007) and pancreatic DNase I (Weston et al., 1992) may be related to their specialized digestion (Fig. S8).
334 In DNase I, amino acid sites exposed on the surface of protein structures displayed additional convergent
335 substitutions that change the charge of their target amino acid residues (E124K, G172D, and H208N),
336 possibly resulting in convergent changes in the biochemical properties of the protein, such as optimal pH,
337 resistance to proteolysis, and posttranslational modifications. Consistent with this idea, bovine and rabbit
338 DNase I proteins are known to be more resistant to degradation by pepsin than their homologs in other
339 animals (Fujihara et al., 2012). Furthermore, E124K was shown to be important for the phosphorylation of
340 bovine DNase I (Nishikawa et al., 1997). Other convergent substitutions will be promising candidates for
341 future characterization. Taken together, these results show how our approach can detect genetic changes
342 associated with phenotypes on the macroevolutionary scale.

343

344 **Exploratory analysis of as-yet-uncharacterized molecular convergence.** We further exploited the 21
345 vertebrate genomes to examine whether ω_C might be used to discover adaptive molecular convergence that
346 may generate hypotheses of linked phenotypes. Since convergence at multiple levels of biological
347 organization can provide strong evidence for adaptive evolution, we searched for simultaneous convergence
348 in protein sequences and gene expression in an exploratory manner without a predefined hypothesis on
349 convergently evolved genes and lineages. Using the same thresholds applied to the analysis of herbivores
350 above ($O_C^N \geq 3.0$ and $\omega_C \geq 3.0$), we obtained 53,805 candidate branch pairs from all orthogroups (Fig. 3B).

351 Although this was an exploratory analysis in which all independent branch pairs were exhaustively
352 analyzed, many studies of convergent evolution involve only a few groups of focal species. If such a research
353 design is applied to this dataset (similar to the analysis of herbivores), the number of detected branch pairs
354 will be much smaller. For example, because there are 538 independent branch pairs in the species tree, on
355 average 100 cases of protein convergence will be obtained in our genome-scale dataset for any particular
356 analysis of two groups of species.

357 To detect convergent gene expression evolution, we employed the amalgamated transcriptomes for
358 six organs in the 21 vertebrate species (Fukushima and Pollock, 2020). Using this previously published
359 dataset, we subjected curated gene expression levels (SVA-log-TMM-FPKM) to multi-optima phylogenetic
360 Ornstein-Uhlenbeck (OU) modeling, in which expression evolution is inferred as regime shifts of estimated
361 optimal expression levels (Khabbazian et al., 2016). Phylogenetic positions and the numbers of expression
362 evolution were determined by a LASSO-based algorithm with Akaike Information Criterion, which was also
363 used for finding convergent shifts toward similar optimal values. In total, we detected 12,017 cases of
364 expression convergence in 4,308 orthogroups (Fig. 3B). Setting the thresholds for gene expression
365 specificity at $\tau \geq 0.67$ (Yanai et al., 2005) and expression levels at $\mu_{max} \geq 2.0$ (the maximum value of fitted
366 SVA-log-TMM-FPKM) (Fukushima and Pollock, 2020), we obtained a set of 2,917 high-confidence branch
367 pairs for potentially adaptive convergence of expression patterns.

368 By taking the intersection of protein convergence and expression convergence, we discovered 33 cases
369 of potentially adaptive joint convergence of expression patterns and protein sequences in 31 orthogroups
370 (Fig. 3B; Table S6). Gene duplication was frequently associated with joint convergence, with at least one
371 branch experiencing gene duplication in 23 out of the 33 branch pairs ($P = 3.11 \times 10^{-25}$, $\chi^2 = 107.7$, χ^2 test of
372 independence). While gene duplication generally reduced the convergence rate, as discussed earlier
373 (Fig. 2C), some of the independently generated duplicates may tend to evolve into the same sequence space
374 when similar expression evolution takes place. Convergence of testis-specific genes was most frequently
375 observed (19/33 orthogroups) and significantly enriched ($P = 1.36 \times 10^{-31}$, $\chi^2 = 136.8$, χ^2 test of
376 independence). The mechanism by which the testis serves as a major place for functional evolution of
377 duplicated genes has been explained by several factors, including the ease with which expression is acquired
378 in spermatogenic cells (Kaessmann, 2010; Kleene, 2005). This phenomenon is called the out-of-the-testis
379 hypothesis, and our results suggest that predictable protein evolution may be enriched in this evolutionary
380 pathway.

381 To infer the functional effect of convergent amino acid substitutions, we mapped the positions of
382 substitutions onto known conformations of homologous proteins. Strikingly, we observed convergently
383 evolved proteins where clusters of substitutions are localized to functionally important sites. They included
384 members of aldo-keto reductase family 1 (AKR1), which play essential roles in steroid metabolism (Rižner
385 and Penning, 2014). The OU analysis revealed that *AKR1* acquired preferential expression in the ovary after
386 repeated lineage-specific duplications in rabbits and mice (*Mus musculus*) (Fig. 3C). Among the paired
387 substitutions in the two lineages, F129I (convergence) and F306A/V (double divergence) located to the
388 positions that delineate the steroid-binding cavity (Fig. 3C). At residue 306, the size of the amino acid was
389 shown by targeted mutagenesis to be important for catalytic promiscuity in rabbits (Couture et al., 2004).
390 Similarly, D224C/E (double divergence) occurred in a loop that contributes to substrate specificity (Couture
391 et al., 2004). These results suggest that the phenotypic change related to substrate specificity might have
392 occurred not only in rabbits but also in mice and underscore how F129I, together with the other two
393 convergence cases (N11S and T/S289P, Fig. S9A), should be a major target for future characterization.

394 Similarly, *nudix hydrolase 16-like 1* (*NUDT16L1*, also known as *Tudor-interacting repair regulator*
395 [*TIRR*]), which is involved in cell migration (Gunaratne et al., 2011) and whose encoded protein binds to
396 RNA and P53-binding protein 1 (53BP1) (Botuyan et al., 2018), showed lineage-specific duplications in
397 chinchillas (*Chinchilla lanigera*) and another rodent lineage connected to mice and rats (*Rattus norvegicus*)
398 (Fig. 3D). The duplication events were followed by convergent regime shifts that resulted in testis-specific
399 expression. The expression evolution was coupled with convergent substitutions in the protein sites
400 corresponding to the substrate-binding pocket of the de-ADP-ribosylating homolog NUDT16

401 (Thirawatananond et al., 2019; Zhang et al., 2020). Protein convergence linked to testis-specific expression
402 was also observed in *myeloid-associated differentiation marker (MYADM)*, which encodes a transmembrane
403 protein that localizes to membrane rafts (Aranda et al., 2011), regulates eosinophil apoptosis through binding
404 to Surfactant protein A (SP-A) (Dy et al., 2021), and participates in cell proliferation and migration (Sun et
405 al., 2016). This orthogroup showed joint convergence in two pairs of branches, in both of which the
406 convergent amino acid substitutions were almost entirely confined to one side of the transmembrane domains
407 (Fig. 3E), suggesting altered interactions with other molecules through this portion of the protein.

408 Finally, an orthogroup of dihydrodiol dehydrogenase (DHDH) showed joint convergence of
409 expression and proteins (Fig. 3F). Possible physiological roles of this enzyme included the detoxification of
410 cytotoxic dicarbonyl compounds, such as 3-deoxyglucosone derived from glycation (Nakayama et al., 1991;
411 Sato et al., 1993). Although the domain structure of proteins was well conserved among species (Fig. S9A),
412 the gene expression patterns of the encoding genes tended to vary. *DHDH* is known to show distinct tissue-
413 specific expression patterns in mammals: kidney in monkeys (*Macaca mulatta*) (Nakagawa et al., 1989),
414 kidney and liver in dogs (*Canis lupus*) (Sato et al., 1994), liver and lens in rabbits (Arimitsu et al., 1999),
415 and various tissues in pigs (*Sus scrofa*) (Nakayama et al., 1991). Our amalgamated transcriptomes showed
416 largely consistent species-specific expression patterns (Fig. 3F). The OU analysis recovered four lineage-
417 specific regime shifts categorized into two pairs of convergent expression evolution. One of them, the
418 convergence of gene expression that occurred between frogs (*Xenopus*) and the blind cave fish (*Astyanax*),
419 which diverged approximately 435 million years ago (Hedges et al., 2015), is characterized by kidney-
420 specific expression. The *Xenopus* gene ENSXETG00000033613 appeared to have arisen from a more widely
421 expressed ancestral gene after a lineage-specific gene duplication. By contrast, the *Astyanax* gene
422 ENSAMXG00000005808 may have acquired kidney-specific expression without any detectable duplication.
423 In this branch pair, we detected a protein convergence rate that cannot be explained by neutral evolution,
424 with a convergence of five amino acid sites (Fig. S9A). These convergent substitutions localized around the
425 active site, while we did not observe such a trend for the double divergence (Fig. 3F). This result suggests
426 that the convergent substitutions may have occurred adaptively to change ancestral catalytic function.

427 *DHDH* has a broad substrate specificity for carbonyl compounds. This protein oxidizes *trans*-
428 cyclohexanediol, *trans*-dihydrodiols of aromatic hydrocarbons, and monosaccharides including D-xylose,
429 while it reduces dicarbonyl compounds, aldehydes, and ketones (Sato et al., 1994). Its active site is
430 predominantly formed by hydrophobic residues, suggesting their role in catabolizing aromatic hydrocarbons
431 (Carbone et al., 2008b, 2008a). Notably, the convergent substitutions in the substrate-binding sites tended to
432 increase amino acid hydrophobicity (Fig. S9B), suggesting that the remodeling of the active site may have
433 led to the acquisition of new substrates in *Xenopus* and *Astyanax*.

434 In summary, ω_C was not only robust against phylogenetic errors, outperforming other methods in
435 simulation and empirical data, but also allowed us to discover plausible adaptive convergence from a
436 genome-scale dataset without a pre-existing hypothesis. Molecular convergence revealed by our exploratory
437 analysis will provide a basis for understanding overlooked phenotypes that protein evolution led to in
438 corresponding lineages.

439
440 **Heuristic detection of highly repetitive adaptive convergence.** Convergent events observed on even more
441 than two independent lineages are exceptionally good signals of adaptive evolution, if they exist, because
442 three or more combined convergences should be extremely rare in random noise. Conventionally,
443 convergence in more than two branches has been analyzed as multiple pairwise comparisons for which there
444 is a prior hypothesis of convergence. The difficulty in analyzing higher-order combinatorial substitutions
445 without specific prior hypotheses lies in the need to explore a vast combinatorial space that exponentially
446 expands as the number of branches to be combined (K) increases. For example, an evenly branching tree
447 with 64 tips has 7,359 independent branch pairs (i.e., at $K = 2$), but the number of branch combinations
448 exponentially increases to 333,375 and 6,976,859 in triple ($K = 3$) and quadruple ($K = 4$) combinations,
449 respectively, making it impractical to exhaustively search highly repetitive convergence even in a single
450 phylogenetic tree when a hypothesis on focal lineages is unavailable.

451 To overcome this limitation, we developed an efficient branch-and-bound algorithm (Land and Doig,
452 1960) that progressively searches for higher-order branch combinations (Fig. 4A and Fig. S10A). For the
453 performance evaluation, we used the PEPC tree (Fig. 4B) because it has repeated adaptive convergence for
454 its use in C₄ photosynthesis (Fig. 1E). While the exhaustive search required 156 minutes with $K = 3$ to
455 analyze 307,432 branch combinations using two central processing units (CPUs), our branch-and-bound
456 algorithm required only 21 seconds. At $K = 4$, the exhaustive search completed within a practical time by
457 using 16 CPUs (46 hours for nearly 8 million combinations) but failed to complete at $K = 5$ (152 million
458 combinations). By sharp contrast, the heuristic search took about 5 minutes for the entire analysis, of which
459 the higher-order analysis with K ranging from 3 to 6 took only about 1 minute to analyze as few as 390
460 combinations with two CPUs (Table S7).

461 The analyzed tree covered nine independent origins of C₄-type PEPC, and the corresponding branch
462 pairs of C₄ lineages accounted for 1.1% of all possible pairs (94/8,308). Convergent branch pairs defined by
463 a threshold ($\omega_C \geq 5.0$ and $O_C^N \geq 2.0$) enriched for the C₄ lineages at $K = 2$ (29.9%, 26/87; Fig. 4C). The
464 convergence of non-C₄ lineages (61/87, including pairs of C₄ and non-C₄ branches) can be interpreted as
465 false positives or adaptive convergence associated with other currently unknown functions. The subsequent
466 higher-order analysis resulted in the discovery of highly repetitive convergence in combinations of as many
467 as six branches (i.e., $K = 3$ to $K = 6$). As the order increased, the lineages of C₄-type PEPCs rapidly
468 predominated and accounted for all the combinations detected at $K \geq 5$ (Fig. 4C), even though the heuristic
469 algorithm was not given any information about the C₄ lineages.

470 In the higher-order C₄ branch combinations, the detected convergence events were almost entirely
471 nonsynonymous (O_C^N), while synonymous convergence (O_C^S) was negligible (Fig. 4D). As a result, the rate
472 of synonymous convergence (dS_C) quickly approached zero (Fig. 4D). Notably, the higher-order convergent
473 substitutions were located at functionally important protein sites. In the convergent branch combinations
474 with $K = 6$, we identified three amino acid sites with a joint posterior probability of nonsynonymous
475 convergence greater than 0.5: V627I, H665N, and A780S (Fig. S10B–D). The H665N substitution generates
476 a putative N-glycosylation site that may be important for protein folding (Christin et al., 2007). The A780S
477 substitution, for which the signature of positive selection had been detected previously (Besnard et al., 2009;
478 Hermans and Westhoff, 1992; Poetsch et al., 1991), has been shown to change the enzyme kinetics related
479 to the first committed step of C₄ carbon fixation (Bläsing et al., 2000; DiMario and Cousins, 2019;
480 Engelmann et al., 2002) and is therefore considered a diagnostic substitution of C₄-type PEPC (Besnard et
481 al., 2009; Christin et al., 2007). The third substitution, C627I, might be a good focus for future
482 experimentation. These results demonstrate that higher-order analysis can substantially increase the signal-
483 to-noise ratio in convergence analysis when there is repeated selective pressure to evolve similar biochemical
484 functions.

485

486 Discussion

487 In this study, we introduced a measure of convergent protein evolution, ω_c , designed to account for
488 false signals due to phylogenetic error. We showed, through simulation and analysis of real biological data,
489 that ω_c mostly eliminates false positives without reduction in power to detect true signals. We also developed
490 an approach to estimate the rates of highly repetitive convergence (i.e., on more than two lineages) fully
491 accounting for phylogenetic combinatorics and demonstrated that the specificity of ω_c increases further in
492 the higher-order analysis. Because of its improved accuracy, ω_c should further drive macroevolutionary
493 analyses where uncorrected measures have been used to identify responsible genotypes for particular
494 phenotypes in a way similar to genome-wide association studies (GWASs). As in GWAS-identified alleles
495 (or genes in gene-level association tests (Wang et al., 2021)), genes with excess convergence serve as clues
496 to study macroevolutionary traits for which the molecular basis is unknown (Fig. 5). Furthermore, the
497 accuracy of ω_c even allows exploratory analysis (Fig. 5), as demonstrated here in vertebrate genomes
498 (Fig. 3). By conducting a genome-wide search of convergent branch combinations, we detected signatures
499 of likely adaptive convergence, which leads to hypothesis generation on responsible phenotypes. This
500 outcome was possible because ω_c , unlike P -values from GWASs, does not require phenotypic traits as input.
501 Convergently evolved genes identified by exploratory analysis will, in turn, lead to the discovery of
502 overlooked phenotypes through future experimentation.

503 Although ω_c is a powerful means to detect convergence while removing the effect of phylogenetic
504 error, there are other sources of stochastic error that can mask small signals. We successfully captured
505 multiple known convergence events here, even with only two or three amino acid substitutions involved in
506 small proteins (Fig. 1E and Table S3). However, a convergent amino acid substitution at a single site in only
507 two lineages may not reliably be identified as resulting from adaptation rather than random homoplasy, by
508 ω_c or any other measure. Therefore, the number of observed nonsynonymous convergence (O_c^N) should
509 always be considered in addition to the phylogenetic error-corrected convergence rate (ω_c), especially in a
510 genome-scale screening with only two or three focal lineages. If many amino acid sites and/or many separate
511 lineages are involved, true convergence is, in general, more easily detected.

512 Protein convergence has attracted a great deal of attention for its potential to associate long-term
513 genotypic variation with phenotypic change, from its first discovery (Stewart et al., 1987), subsequent
514 theoretical development (Castoe et al., 2009; Zhang and Kumar, 1997), the first claim of genome-wide
515 detection (Parker et al., 2013), to recent findings that highlighted epistatic effects (Goldstein and Pollock,
516 2017; Goldstein et al., 2015; Zou and Zhang, 2015a, 2017) and technical difficulties (Foote et al., 2015;
517 Mendes et al., 2016, 2019; Thomas and Hahn, 2015; Zou and Zhang, 2015b). Other types of convergence at
518 the molecular level beyond amino acid substitutions have also been considered, including convergent shifts
519 of site-wise substitution profiles (Rey et al., 2018), convergent shifts of evolutionary rates (i.e., number of
520 substitutions per time regardless of the amino acid state or substitution profile) (Kowalczyk et al., 2019),
521 convergent rate shifts of noncoding elements (Hu et al., 2019), convergent gene losses (Hiller et al., 2012;
522 Prudent et al., 2016), convergent losses of noncoding elements (Marcovitz et al., 2016), and functional
523 enrichments of convergently evolved loci (Marcovitz et al., 2019). Using transcriptome amalgamation,
524 which integrates multi-species gene expression data in a comparable manner (Fukushima and Pollock, 2020),
525 we developed a means to detect convergence in gene expression levels and to correlate the obtained results
526 with protein convergence rates. Further integration of these methods will allow us to examine how well
527 convergent patterns correlate across multiple hierarchies of biological organizations. Such analysis will
528 provide a quantitative perspective of the extent to which evolution at one hierarchical level causes predictable
529 changes in another.

530 Although it is well established that phenotypes are associated with genotypes, the genetic basis for
531 particular convergently evolved phenotypes may arise from distinct, non-convergent genetic changes
532 (Concha et al., 2019; Natarajan et al., 2016). These specific cases may sometimes occur because of
533 convergent mechanisms, such as the use of similar but not identical amino acids, and the use of similar
534 changes at adjacent residues in the protein structure (Castoe et al., 2008). The accumulation of knowledge

535 about which mutations are repeatedly selected and which are not during convergent evolution may provide
536 insight into the evolvability and constraints that govern the diversification of organisms.

537 While some evolutionary innovations may be unique, many traits arose convergently (Vermeij, 2006).
538 Fascinating examples not mentioned above include endothermy, hibernation, burrowing, diving, venom
539 injection, electrogenic organs, eusociality, anhydrobiosis, bioluminescence, biomineralization, plant
540 parasitism, mycoheterotrophy, and multicellularity. In the past, the observation of similar phenotypes in
541 multiple species led to the theory of evolution by natural selection (Darwin, 1859). The analysis of protein
542 sequences in multiple species gave rise to the formulation of the nearly neutral theory of molecular evolution
543 (Kimura, 1968; Ohta, 1973). Likewise, cross-species genotype-phenotype associations illuminated through
544 the analysis of molecular convergence, coupled with experimental evaluation of mutational effects
545 (Supplementary Text 9), may lead to new conceptual frameworks on the constraint and adaptive changes at
546 the molecular level that drive phenotypic change among species.

547

548 **Methods**

549

550 **Simulated codon sequence evolution.** With the input phylogenetic tree (Fig. 1C), codon sequences of
551 specified length (500 codons) were generated with the ‘simulate’ function of CSUBST
552 (<https://github.com/kfuku52/csubst>), which internally utilizes the python package pyvolve for simulated
553 sequence evolution (Spielman and Wilke, 2015). An empirical codon substitution model with multiple
554 nucleotide substitutions (Kosiol et al., 2007) was adjusted with observed codon frequencies (ECMK07+F)
555 in the vertebrate genes encoding phosphoglycerol kinases (PGKs, available from the ‘dataset’ function of
556 CSUBST). The conventional ω (dN/dS) was set to 0.2. In the Convergent scenario, 5% of codon sites were
557 evolved convergently in focal lineages (the pair of terminal branches in Fig. 1C). At convergent codon sites,
558 the frequency of nonsynonymous substitutions to codons encoding a single randomly selected amino acid
559 was increased so that nonsynonymous substitutions to the selected codons accounted for approximately 90%
560 of the total. This operation increases the probability of amino acid convergence without changing relative
561 frequencies among synonymous codons. The site-specific substitution rate at convergent codon sites was
562 also doubled (i.e., $r = 2$), and a higher nonsynonymous/synonymous substitution rate ratio was applied (i.e.,
563 $\omega = 5$) to mimic adaptive evolution. The simulation parameters for the other scenarios are summarized in
564 Table S2. For the Random scenario, randomized trees were generated in 1,000 simulations with the ‘shuffle’
565 function of NWKIT v0.10.0 and the --label option (<https://github.com/kfuku52/nwkit>).

566

567 **Animal gene sets.** A dataset of amalgamated cross-species transcriptomes (Fukushima and Pollock, 2020)
568 was generated for 21 vertebrate genomes in Ensembl 91 (Yates et al., 2016) (Table S8). To ensure
569 compatibility, the same versions of protein-coding sequences were also used for the convergence analysis.
570 Completeness of genome assembly was evaluated using BUSCO v4.0.5 (Simão et al., 2015) with the single-
571 copy gene set of ‘tetrapoda_odb10’ (Table S8). A species phylogenetic tree previously downloaded from
572 TimeTree (Hedges et al., 2006) was used (Fukushima and Pollock, 2020). Orthogroups were classified by
573 OrthoFinder v2.4.1 (Emms and Kelly, 2015, 2019). Orthogroups containing more than three genes were
574 analyzed further. During the analysis of this dataset, a protein size-dependent change in measured
575 convergence rates was observed (Fig. S11) but was determined to be an artifact; ω_c was shown to be more
576 robust to the bias than the other metrics (Supplementary Text 10).

577

578 **Sequence retrieval from public databases.** Gene sets for previously confirmed cases of molecular
579 convergence and horizontal gene transfer events (HGTs) were generated based on previous reports with
580 increased taxon sampling (Table S3; Supplementary Text 11). With GenBank accession numbers for
581 ATPalpha1, Prestin, PEPC, and PCK homologs (Supplementary Dataset), coding sequences (CDSs) were
582 retrieved using the ‘accession2fasta’ function of CDSKIT. Lysozyme sequences were downloaded as
583 GenBank files from NCBI and were converted to fasta files with the ‘parsegb’ function of CDSKIT. For the
584 retrieval of the mitochondrial genome, a custom python script was used to select balanced numbers and
585 lineages of foreground and background species (Supplementary Dataset). Orthogroup CDS files for og3737
586 (leucine-tRNA ligase), og9103 (pentatricopeptide repeat protein), and og9298 (pentatricopeptide repeat
587 protein) for the HGT events in *Cuscuta* were obtained from a previous report (Yang et al., 2019b), and genes
588 leading to unrealistically long branches were excluded. HGTs in the other parasitic lineage Orobanchaceae
589 were also analyzed in the same report, but HGTs in *Cuscuta* were used for performance evaluation because
590 the donor lineage was unequivocal in several genes.

591

592 **Sequence retrieval from plant gene sets.** Gene sets were downloaded from public databases for the retrieval
593 of CDSs encoding digestive enzyme homologs (Table S8). Transcriptome assemblies were used as a part of
594 gene sets. For *Drosera adela*, *Nepenthes* cf. *alata*, and *Sarracenia purpurea*, previously assembled
595 transcriptomes were used (Fukushima et al., 2017). The transcriptome assembly of *Rhododendron delavayi*
596 was generated from publicly available RNA-seq data (NCBI BioProject ID: PRJNA476831) with Trinity
597 v2.8.5 (Grabherr et al., 2011) after pre-processing with fastp v0.20.1 (Chen et al., 2018) (Supplementary

598 Dataset). Subsequently, open reading frames (ORFs) were obtained with TransDecoder v5.5.0
599 (<https://github.com/TransDecoder/TransDecoder>). The longest ORFs among isoforms were extracted with
600 the ‘aggregate’ function of CDSKIT v0.9.1 (<https://github.com/kfuku52/cdskit>). The completeness of
601 assembly was evaluated using BUSCO scores with the single-copy gene set of ‘embryophyta_odb10’
602 (Table S8). Finally, digestive enzyme homologs were retrieved by TBLASTX v2.9.0 searches against all
603 gene sets with an E-value cutoff of 0.01 and >50% query coverage (Camacho et al., 2009).

604

605 **Characterization of protein-coding sequences.** Coding sequences were used for RPS-BLAST v2.9.0
606 searches (Camacho et al., 2009) against Pfam-A families (El-Gebali et al., 2019) (released on April 30, 2020)
607 with an E-value cutoff of 0.01 to obtain protein domain architectures. The numbers of transmembrane
608 domains were predicted by TMHMM v2.0 (Krogh et al., 2001). The numbers of introns in protein-coding
609 sequences were extracted from GFF files downloaded from Ensembl. Further gene annotations were
610 obtained using Trinotate v3.2.1 (<https://github.com/Trinotate/Trinotate.github.io/wiki>).

611

612 **Plant species tree.** Orthogroup classification was performed with OrthoFinder v2.4.1 (Emms and Kelly,
613 2019). Stop codons and ambiguous codons were masked as gaps using CDSKIT. In-frame multiple sequence
614 alignments of single-copy orthologs were generated by MAFFT v7.455 with the --auto option (Katoh and
615 Standley, 2013) and tralign in EMBOSS v6.6.0 (Rice et al., 2000). Ambiguous codon sites were then
616 removed by ClipKIT v0.1.2 with the default parameters (Steenwyk et al., 2020). After the concatenation of
617 trimmed sequences, a maximum-likelihood phylogenetic tree was reconstructed by IQ-TREE v2.0.3 with
618 the GTR+G nucleotide substitution model (Minh et al., 2020; Nguyen et al., 2015). The tree was rooted using
619 *Amborella trichocarpa* as an outgroup. The divergence time of the species tree was estimated using
620 mcmctree in the PAML package v4.9 (Yang, 2007). The priors and parameters were chosen according to the
621 mcmctree tutorial (<http://abacus.gene.ucl.ac.uk/software/paml.html>). Fossil calibrations were adopted from
622 a previous study (Zhang et al., 2017).

623

624 **In-frame codon sequence alignment.** Retrieved coding sequences were formatted into in-frame sequences
625 using the ‘pad’ function of CDSKIT. Stop codons and ambiguous codons were replaced with gaps with the
626 ‘mask’ function of CDSKIT. Amino acid sequences from translated coding sequences were aligned using
627 MAFFT with the --auto option (Katoh and Standley, 2013), trimmed with ClipKIT with default parameters,
628 and reverse-translated with the ‘backtrim’ function of CDSKIT. Gappy codon sites were excluded with the
629 ‘hammer’ function of CDSKIT.

630

631 **Phylogenetic tree reconstruction.** The gene tree was first reconstructed using IQ-TREE with the general
632 time-reversible (GTR) nucleotide substitution model and four gamma categories of among-site rate
633 heterogeneity (ASRV). To suppress branch attraction in the trees containing HGTs, topological constraints
634 consistent with species classification were generated from the NCBI Taxonomy (Schoch et al., 2020) using
635 the ‘constrain’ function of NWKIT and used for tree search. Ultrafast bootstrapping with 1,000 replicates
636 was performed to evaluate the credibility of tree topology (Minh et al., 2013) with further optimization of
637 each bootstrapping tree (-bnni option) (Hoang et al., 2018). To improve tree topology, some datasets were
638 subjected to phylogeny reconciliation with the species tree using GeneRax v1.2.2 (Morel et al., 2020)
639 (Table S3). Branching events in gene trees were categorized into speciation or gene duplication by a species-
640 overlap method (Huerta-Cepas et al., 2007). *Arabidopsis thaliana* orthologs in each clade were inferred from
641 the tree topology. Minor differences in the methods applied to each dataset, from sequence retrieval to
642 phylogenetic analysis, are summarized in Table S3.

643

644 **Detecting convergent expression evolution.** Using the dated species tree and rooted gene trees as inputs,
645 the divergence time of individual gene trees was estimated by RADTE
646 (<https://github.com/kfuku52/RADTE>) as described previously (Fukushima and Pollock, 2020). Evolution of
647 gene expression levels (SVA-log-TMM-FPKM) (Fukushima and Pollock, 2020) in brain, heart, kidney,

648 liver, ovary, and testis samples was modeled on the dated gene tree with phylogenetic multi-optima Ornstein-
 649 Uhlenbeck models (i.e., Hansen models (Hansen, 1997)) with the ‘estimate_shift_configuration’ function in
 650 the R package *l1ou* (Khabbazian et al., 2016) as described previously (Fukushima and Pollock, 2020).
 651 Convergent regime shifts were then detected as multiple regime shifts that lead to similar expression levels,
 652 as judged by the ‘estimate_convergent_regimes’ function (Khabbazian et al., 2016).

653
 654 **Classification of combinatorial substitutions.** Combinatorial substitutions were collectively defined as
 655 substitutions at the same protein site that occur in multiple independent branches in a phylogenetic tree.
 656 When this occurs only in two branches, it is called a paired substitution. In unambiguous notation, we
 657 consider paired substitutions along two branches with the same specific state (spe), different states (dif), or
 658 any state (any) at the ancestral and derived nodes. The five combinatorial states that we discuss and that are
 659 frequently considered in the literature are paired substitutions (any→any), double divergence (any→dif),
 660 convergence (any→spe), discordant convergence (dif→spe), and congruent convergence (spe→spe)
 661 (Fig. S1C). Convergence is discussed throughout this report because it is of particular importance in testing
 662 evolutionary genotype-phenotype associations.

663
 664 **Ancestral state reconstruction and parameter estimation.** Our method estimates convergent substitution
 665 via ancestral reconstruction. Whereas ancestral amino acid reconstruction has been used in previous reports
 666 (Foote et al., 2015; Goldstein et al., 2015; Thomas and Hahn, 2015; Zou and Zhang, 2015a), here we used
 667 codon sequence reconstruction. Using the input phylogenetic tree and observed codon sequences, CSUBST
 668 internally uses IQ-TREE to estimate the posterior probabilities of ancestral sequences by the empirical
 669 Bayesian method (Minh et al., 2020). At the same time, the parameters used in CSUBST are estimated:
 670 codon equilibrium frequencies (π_i), ASRV (r_l), nonsynonymous per synonymous substitution ratio (ω), and
 671 transition per transversion substitution ratio (κ).

672
 673 **Multidimensional array structures for substitution history.** CSUBST stores the coding sequences and
 674 the reconstructed probable ancestral states in a three-dimensional array whose size is $M \times L \times 61$ for a
 675 phylogenetic tree with M nodes (excluding the root node) generated from an alignment of coding sequences
 676 with L codon sites, each of which can take a distribution of 61 different codon states (in the universal genetic
 677 code), excluding stop codons. We denote by $P_{mlj}(X|D, \theta)$ the posterior probability of codon X for codon
 678 state j at site l on node m . The three-dimensional array for codon states is then converted to a four-
 679 dimensional array that stores the probability of substitutions with the size of $B \times L \times 61 \times 61$, where B
 680 denotes the number of branches excluding the root branch. This array stores the posterior probability of
 681 substitution $P_{blj}(S|D, \theta)$ for single substitution S from ancestral codon state i to derived codon state j for a
 682 codon site l in branch b . For a site l in branch b connecting ancestral node n with codon state i and
 683 descendant node m with codon state j , the posterior probability substitution matrix $P_{ij}(S|D, \theta)$ is derived as

$$P_{ij}(S|D, \theta) = P_i(X|D, \theta) \times P_j(X|D, \theta)^T = \begin{pmatrix} i_1 \\ i_2 \\ \vdots \\ i_{61} \end{pmatrix} \times \begin{pmatrix} j_1 & j_2 & \dots & j_{61} \end{pmatrix} = \begin{pmatrix} i_1 j_1 & i_1 j_2 & \dots & i_1 j_{61} \\ i_2 j_1 & i_2 j_2 & \dots & i_2 j_{61} \\ \vdots & \vdots & \ddots & \vdots \\ i_{61} j_1 & i_{61} j_2 & \dots & i_{61} j_{61} \end{pmatrix} \quad (1)$$

684 As the transition between the same codon state is not considered a substitution, the diagonal elements ($ij_{i=j}$)
 685 are filled with 0. Although Equation 1 is an approximation that does not take into account the non-
 686 independence between nodes of a phylogenetic tree, we confirmed that the effect was negligible
 687 (Supplementary Text 12; Fig. S12). For efficient processing of nonsynonymous and synonymous
 688 substitution probabilities with the array operation of NumPy (Harris et al., 2020), the four-dimensional array
 689 is converted into a pair of five-dimensional arrays (A^N and A^S for nonsynonymous and synonymous
 690 substitutions, respectively) whose individual size is $B \times L \times G \times I \times J$, where codon states are grouped into

691 G categories (Fig. S2A). Stored values range between 0 and 1, denoted by $P_{blgij}(S|D, \theta)$, the probability of
 692 single substitution S from ancestral codon i to derived codon j ($i \neq j$) in codon group g at site l of branch
 693 b , given the observed sequence data D and model parameters θ that include the phylogenetic tree. The
 694 elements in the array A^N indicate $P_{blgij}(S^N|D, \theta)$, the probabilities of nonsynonymous substitutions (S^N),
 695 whereas those in the array A^S correspond to $P_{blgij}(S^S|D, \theta)$, the probabilities of synonymous substitutions
 696 (S^S). In A^N , a single 20×20 matrix records all the substitution probabilities, and therefore $G = 1$ and $I = J =$
 697 20 . Synonymous substitutions occur only between codons that code for the same amino acid. Since there are
 698 20 different amino acids, G equals 20 in A^S . In the case of the universal genetic code, the maximum number
 699 of codons encoding the same amino acid is six, for leucine, serine, and arginine, so $I = J = 6$. In the matrix
 700 corresponding to these three amino acids, all values are between 0 and 1, but for amino acids with a smaller
 701 number of codons, the out-of-range indices are filled with zero. Missing sites in the sequence alignment are
 702 also treated as zero. For simplicity, we explain the case where there is no missing site in the observed
 703 sequences and ancestral states in the following sections, but the implementation in CSUBST appropriately
 704 takes into account the missing sites by subtracting its numbers from L at every necessary step in individual
 705 branches or branch combinations.

706

707 **Tree rescaling.** During the ancestral state reconstruction, IQ-TREE estimates the branch length as the
 708 number of nucleotide substitutions per codon site. Since our model requires the number of codon
 709 substitutions rather than the number of nucleotide substitutions, and since branch lengths are required
 710 separately for both synonymous and nonsynonymous substitutions, we obtained rescaled branch length t_b
 711 of branch b as follows:

$$t_b = \frac{\sum_{l=1}^L \sum_{g=1}^G \sum_{i=1}^I \sum_{j=1}^J P_{blgij}(S|D, \theta)}{L} \quad (2)$$

712 t_b^N and t_b^S for nonsynonymous and synonymous substitutions were obtained with $P_{blgij}(S^N|D, \theta)$ and
 713 $P_{blgij}(S^S|D, \theta)$, respectively. For example, with the ECMK07+F+R4 model, the total branch lengths of the
 714 21 vertebrate PGK tree before and after rescaling are 7.57 nucleotide-substitutions/codon-site and 7.20
 715 codon-substitutions/codon-site (1.59 nonsynonymous and 5.62 synonymous codon substitutions per codon
 716 site).

717

718 **Observed number of combinatorial substitutions.** The only true observations are the gene sequences of
 719 the extant species, and the posterior probabilities of ancestral sequences and codon substitutions are
 720 estimates. However, we refer to the posterior probabilities as “observations” (Zou and Zhang, 2015a) to
 721 unambiguously distinguish them from the expected values described in the next section. Here, we denote by
 722 $P_l(S_C|D, \theta)$ the probability of combinatorial substitution S_C at codon site l given observed sequences D and
 723 model θ . The probabilities of nonsynonymous and synonymous combinatorial substitutions at site l are
 724 separately obtained as $P_l(S_C^N|D, \theta)$ and $P_l(S_C^S|D, \theta)$, respectively, with the following equations:

$$P_l^{any \rightarrow any}(S_C|D, \theta) = \sum_{g=1}^G \prod_{k=1}^K \sum_{i=1}^I \sum_{j=1}^J P_{klgij}(S|D, \theta) \quad (3)$$

for paired substitutions,

$$P_l^{any \rightarrow spe}(S_C|D, \theta) = \sum_{g=1}^G \sum_{j=1}^J \prod_{k=1}^K \sum_{i=1}^I P_{klgij}(S|D, \theta) \quad (4)$$

for convergence,

725 and

$$P_l^{spe \rightarrow spe}(S_C | D, \theta) = \sum_{g=1}^G \sum_{i=1}^I \sum_{j=1}^J \prod_{k=1}^K \underbrace{P_{klgij}(S | D, \theta)}_{k_1 i = k_2 i, k_1 j = k_2 j} \quad (5)$$

for concordant convergence,

726 where k represents a branch of interest. We denote by K the degree of combinatorial substitutions or the
 727 number of branches to be compared. Because two branches are often compared in conventional convergence
 728 analysis, we explain here the case of $K = 2$. Array operations in the underlined parts of Equation 3 to
 729 Equation 5 are illustrated in Fig. S2B. The total probabilities of observed substitution pairs across sites in
 730 the branch pair are calculated as

$$O_C = \sum_{l=1}^L P_l(S_C | D, \theta) \quad (6)$$

731 O_C is separately obtained for nonsynonymous and synonymous combinatorial substitutions (O_C^N and O_C^S ,
 732 respectively). By definition (Fig. S1C), the values of O_C for double divergence and discordant convergence
 733 are derived as follows at $K = 2$:

$$O_C^{any \rightarrow dif} = O_C^{any \rightarrow any} - O_C^{any \rightarrow spe} \quad \text{for double divergence} \quad (7)$$

734 and

$$O_C^{dif \rightarrow spe} = O_C^{any \rightarrow spe} - O_C^{spe \rightarrow spe} \quad \text{for discordant convergence.} \quad (8)$$

735 C/D (Goldstein et al., 2015) corresponds to $O_N^{any \rightarrow spe} / O_N^{any \rightarrow dif}$ in our notation.

736

737 **Applying codon substitution models for the expectation of combinatorial substitutions.** To estimate the
 738 rate of combinatorial substitutions, the observed number O_C is contrasted with the expected number E_C . E_C
 739 is derived from codon substitution models in a way similar to the previous application of amino acid
 740 substitution models (Zou and Zhang, 2015a). The tested codon substitution models include the empirical
 741 models ECMK07 and ECMrest (Kosiol et al., 2007) and the mechanistic models MG (Muse and Gaut, 1994)
 742 and GY (Goldman and Yang, 1994). The same model was consistently used in the ancestral state
 743 reconstruction and in deriving the model-based expectations of combinatorial substitutions. In the method
 744 described below, empirical equilibrium codon frequencies, the rescaled branch length, and ASRV are also
 745 taken into account. In the empirical models, the codon substitution rate matrix Q is derived according to
 746 previous literature (Kosiol et al., 2007; Whelan and Goldman, 2001) as follows:

$$Q = \{q_{ij}\} = \begin{pmatrix} - & s_{1,2} & \cdots & s_{1,61} \\ s_{1,2} & - & \cdots & s_{2,61} \\ \vdots & \vdots & \ddots & \vdots \\ s_{61,1} & s_{61,2} & \cdots & - \end{pmatrix} \times \text{diag}(\pi_1, \pi_2, \cdots, \pi_{61}) \quad (9)$$

747 where $s_{i,j}$ denotes the exchangeabilities of codon pairs i and j ($s_{ij} = s_{ji}$), and π_i represents the equilibrium
 748 frequencies of 61 codons estimated from the input alignment. In the mechanistic models, mechanistic
 749 substitution parameters are used instead of the exchangeabilities. In the MG model, q_{ij} is obtained with π_i
 750 and nonsynonymous per synonymous substitution ratio ω , whereas transition per transversion substitution
 751 ratio κ is also taken into account in the GY model. Q is then rescaled as

$$\sum_{i=1}^{61} \sum_{\substack{j=1 \\ j \neq i}}^{61} \pi_i q_{ij} = 1 \quad (10)$$

752 Finally, the diagonal elements of Q are completed as

$$q_{ii} = - \sum_{\substack{j=1 \\ j \neq i}}^{61} q_{ij} \quad (11)$$

753 With substitution rate r , the codon transition probability matrix $P_{ij}(t, r)$ after time t are obtained using
754 matrix exponentiation as

$$P_{ij}(t, r) = e^{Qtr}, \quad (12)$$

755 where CSUBST uses the site-wise substitution rate r_l pre-estimated by IQ-TREE and rescaled branch lengths
756 t_b^N and t_b^S in place of r and t , respectively. The distribution of expected substitutions at site l in branch k
757 connecting ancestral node n with codon state i and a descendant node is therefore given by

$$P_{ij}(S^{expected}|D, \theta) = P_i(X|D, \theta) \times P_{ij}(t, r). \quad (13)$$

758 Using $P_{klgij}(S^{expected}|D, \theta)$ in place of $P_{klgij}(S|D, \theta)$, the total probabilities of expected substitution pairs
759 across sites in the branch pair denoted by E_C are obtained by the same procedure used to obtain O_C
760 (Equation 3 to Equation 8). Similar to O_C , the expected numbers of combinatorial substitutions (E_C) are
761 separately calculated for nonsynonymous and synonymous substitution pairs (E_C^N and E_C^S , respectively). By
762 definition (Fig. S1C), the following relationships hold at $K = 2$:

$$E_C^{any \rightarrow dif} = E_C^{any \rightarrow any} - E_C^{any \rightarrow spe} \quad (14)$$

763 and

$$E_C^{dif \rightarrow spe} = E_C^{any \rightarrow spe} - E_C^{spe \rightarrow spe}. \quad (15)$$

764

765 **Nonsynonymous and synonymous combinatorial substitution rates.** With the observed and expected
766 numbers of combinatorial substitutions (O_C and E_C , respectively), the rates of nonsynonymous and
767 synonymous combinatorial substitutions are obtained, respectively, by

$$dN_C = O_C^N / E_C^N \quad (16)$$

768 and

$$dS_C = O_C^S / E_C^S. \quad (17)$$

769 dN_C can be regarded as equivalent to R with the per-gene equilibrium amino acid frequencies (their f_{gene}),
770 but note that some features are different from the corresponding parts for R (Zou and Zhang, 2015a). In
771 particular, we used the standard procedure to derive codon transition probabilities (Equation 13) (Equation
772 1.2 in (Yang, 2006)), whereas no matrix exponentiation is applied for R . In the 21-vertebrate genome dataset,
773 the total expected convergence ($E_C^{N, any \rightarrow spe} = 6,939,070$) corresponds to 87.2% of the total observed
774 convergence ($O_C^{N, any \rightarrow spe} = 6,051,985$). This expectation matches the observation with better accuracy than
775 the previously published results with the *Drosophila* genomes ($582.8/932 = 62.5\%$ with their JTT- f_{gene}
776 model) (Zou and Zhang, 2015a).

777

778 **Accounting for different range distributions of nonsynonymous and synonymous rates of**
779 **combinatorial substitutions.** Under purifying selection, which is the default evolutionary mode of many
780 proteins (Bustamante et al., 2005), the rate of synonymous substitutions is faster than that of nonsynonymous
781 substitutions. Therefore, saturation of synonymous substitutions becomes a potential problem, especially in
782 a counting method that cannot properly account for the effects of multiple substitutions. To account for this
783 issue, we applied a transformation using quantile values (U_p) as follows:

$$dS_C^{corrected} = \begin{cases} dS_C^{uncorrected}, & \text{if } dS_C^{uncorrected} \geq dN_C \\ U_p^{dN_C}, & \text{otherwise} \end{cases}, \quad (18)$$

784 where $U_p^{dN_C}$ denotes the quantile value of the empirical dN_C distribution at p^{dS_C} , the quantile rank of the
 785 dS_C value, among all branch combinations. This operation rescales dS_C to match its distribution range with
 786 that of dN_C , and the resulting ω_C becomes robust for outlier values (Fig. S13). Because of the need for
 787 quantile values, this transformation is only applicable when the branch combinations are exhaustively
 788 searched. In this work, $dS_C^{corrected}$ is used at $K = 2$ unless otherwise mentioned.

789
 790 **Nonsynonymous per synonymous combinatorial substitution rate ratio.** A nonsynonymous per
 791 synonymous combinatorial substitution rate ratio for K branches is given by

$$\omega_C = \frac{dN_C}{dS_C} = \frac{O_C^N/E_C^N}{O_C^S/E_C^S}. \quad (19)$$

792 ω_C can be separately calculated for different categories of combinatorial substitutions, e.g., $\omega_C^{any \rightarrow any}$ for
 793 paired substitutions, $\omega_C^{any \rightarrow spe}$ for double divergence, $\omega_C^{any \rightarrow dif}$ for convergence, $\omega_C^{dif \rightarrow spe}$ for discordant
 794 convergence, and $\omega_C^{spe \rightarrow spe}$ for concordant convergence. For simplicity, the derivation of ω_C was explained
 795 above for the combinatorial substitutions illustrated in Fig. S1C. However, our method can be applied to
 796 other categories of combinatorial substitutions as well. For example, phenotypic convergence may be
 797 associated with the same ancestral amino acid substituted to different amino acids (Konečná et al., 2021), in
 798 which case $\omega_C^{spe \rightarrow any}$ may be useful for analysis.

799
 800 **Branch combinations.** Combinatorial substitutions are a collection of independently occurring evolutionary
 801 events (Fig. S1C). Branch combinations containing an ancestor-descendant relationship did not satisfy the
 802 evolutionary independence and were therefore excluded from the analysis. Although convergent
 803 substitutions occurring in sister branch pairs satisfy the evolutionary independence, they are difficult to
 804 discriminate and are often treated as a single ancestral substitution. For this reason, sister branches were also
 805 excluded from the analysis (Fig. S10A).

806
 807 **A branch-and-bound algorithm for the higher-order signature of combinatorial substitutions.** O_C and
 808 E_C , and hence ω_C , can also be obtained for combinations of more than two branches ($K > 2$). The higher-
 809 order analysis is particularly useful when analyzing traits with extensively repetitive convergence, such as
 810 C_4 photosynthesis, which is thought to have evolved at least 62 times independently (Sage et al., 2011). To
 811 efficiently explore the higher-order dimensions of branch combinations, we devised a branch-and-bound
 812 algorithm that combines the convergence metric cutoff, and the generation of $K + 1$ branch combinations
 813 from the branch overlaps at $K - 1$ (Fig. 4A and Fig. S10A). The higher-order analysis starts with an
 814 exhaustive comparison of branch pairs (i.e., $K = 2$). Next, convergent branch pairs are extracted with an ω_C
 815 cutoff value (≥ 5.0 in Fig. 4). At this time, branch pairs with a small number of convergent substitutions are
 816 excluded by applying an O_C^N cutoff value (≥ 2.0 in Fig. 4). The convergent branch pairs are then subjected to
 817 the all-vs-all comparison. When a shared branch is found, their union is generated as a combination of three
 818 branches to be analyzed. Before proceeding to the analysis at $K = 3$, branch combinations containing a sister
 819 or ancestor-descendant relationship are discarded. In this way, K is sequentially increased by one at a time.
 820 As such, the algorithm searches only for higher-order branch combinations that are guaranteed to have
 821 sufficient convergence metrics in lower-order combinations. In each round, convergent branch combinations
 822 are first extracted by the cutoffs, and then the $K + 1$ combinations are generated by the $K - 1$ overlap, as in
 823 the analysis at $K = 2$. For example, two, three, and four branches should be shared at $K = 3$, $K = 4$, and

824 $K = 5$, respectively. The increase in K continues until the algorithm no longer finds a branch combination
825 that satisfies the criteria of ω_C and O_C^N .

826

827 **Implementation of CSUBST.** The proposed methods, including the calculation of ω_C and the branch-and-
828 bound algorithm for higher-order combinations, were implemented in the ‘analyze’ function of CSUBST,
829 which was written in Python 3 (<https://www.python.org/>). Phylogenetic tree processing was implemented
830 with the python package ETE 3 (Huerta-Cepas et al., 2016). Numpy (Harris et al., 2020), SciPy (Virtanen et
831 al., 2020), and pandas (<https://pandas.pydata.org/>) were used for array and table data processing. Parallel
832 computation was performed by multiprocessing with Joblib (<https://joblib.readthedocs.io/en/latest/>). The
833 intensive calculation was optimized with Cython (Behnel et al., 2011).

834

835 **Mapping combinatorial substitutions to protein structures.** For the analysis of protein structures, a
836 streamlined pipeline was implemented in the ‘site’ function of CSUBST. Using the ‘--pdb besthit’ option,
837 CSUBST requests an online MMseqs2 search (Steinegger and Söding, 2017) against the RSCB Protein Data
838 Bank (PDB) (Berman et al., 2000) to obtain three-dimensional conformation data of closely related proteins.
839 If no hit is obtained, a BLASTP search against the UniProt database is run on the QBLAST server to identify
840 the best-hit protein for which AlphaFold-predicted structure is available (Varadi et al., 2022; Jumper et al.,
841 2021). For some proteins, structural data were manually selected because more appropriate structures were
842 available (e.g., with substrate). Subsequently, CSUBST internally uses MAFFT to generate protein
843 alignments to determine the homologous positions of amino acids and write a PyMOL session file. The
844 protein structures were visualized using Open-Source PyMOL v2.4.0
845 (<https://github.com/schrodinger/pymol-open-source>).

846

847 **Data visualization.** Phylogenetic trees were visualized using the python package ETE 3 (Huerta-Cepas et
848 al., 2016) and the R package ggtree (Yu et al., 2017). General data visualization was performed with python
849 packages matplotlib (Hunter, 2007) and seaborn (Waskom, 2021) as well as the R package ggplot2
850 (Wickham, 2009, 2). Boxplot elements of all figures are defined as follows: center line, median; box limits,
851 upper and lower quartiles; whiskers, $1.5 \times$ interquartile range.

852

853 **Data availability**

854 Raw data and results are available in the Supplementary Dataset (<https://doi.org/10.5061/dryad.tx95x6b0v>
855 [\[Private URL for peer review:](#)
856 [https://datadryad.org/stash/share/BhwCL-1-YsgHFvYLPX76jMrwoD2AIQP8P84bFmMyFu8\]](https://datadryad.org/stash/share/BhwCL-1-YsgHFvYLPX76jMrwoD2AIQP8P84bFmMyFu8])).

857

858 **Code availability**

859 CSUBST is available from GitHub (<https://github.com/kfuku52/csubst>). The results reported in this study
860 can be reproduced with CSUBST v.0.20.17. Scripts used in this study are available in the Supplementary
861 Dataset (<https://doi.org/10.5061/dryad.tx95x6b0v> [\[Private URL for peer review:](#)
862 [https://datadryad.org/stash/share/BhwCL-1-YsgHFvYLPX76jMrwoD2AIQP8P84bFmMyFu8\]](https://datadryad.org/stash/share/BhwCL-1-YsgHFvYLPX76jMrwoD2AIQP8P84bFmMyFu8])).

863

864 **Acknowledgments**

865 We acknowledge the following sources for funding: MEXT/JSPS KAKENHI 18J00178 (K.F.), the Sofja
866 Kovalevskaja programme of the Alexander von Humboldt Foundation (K.F.), a Human Frontier Science
867 Program (HFSP) Young Investigators Grant RGY0082/2021 (K.F.), and NIH R01 GM083127 (D.D.P.).
868 Computations were partially performed on the National Institute of Genetics (NIG) supercomputer.

869

870 **Author Contributions**

871 K.F. designed the study. K.F. designed and wrote all programs and performed data analysis. D.P. contributed
872 to conceptualizing and helping guide the analysis. K.F. and D.D.P. wrote the paper.

873

874 **Competing Interests**

875 The authors declare no competing interests.

876

877 **References**

- 878 **Anzalone, A.V., Koblan, L.W., and Liu, D.R.** (2020). Genome editing with CRISPR–Cas nucleases,
879 base editors, transposases and prime editors. *Nat. Biotechnol.* **38**: 824–844.
- 880 **Aranda, J.F., Reglero-Real, N., Kremer, L., Marcos-Ramiro, B., Ruiz-Sáenz, A., Calvo, M., Enrich,**
881 **C., Correas, I., Millán, J., and Alonso, M.A.** (2011). MYADM regulates Rac1 targeting to
882 ordered membranes required for cell spreading and migration. *Mol. Biol. Cell* **22**: 1252–1262.
- 883 **Arendt, J. and Reznick, D.** (2008). Convergence and parallelism reconsidered: what have we learned
884 about the genetics of adaptation? *Trends Ecol. Evol.* **23**: 26–32.
- 885 **Arimitsu, E., Aoki, S., Ishikura, S., Nakanishi, K., Matsuura, K., and Hara, A.** (1999). Cloning and
886 sequencing of the cDNA species for mammalian dimeric dihydrodiol dehydrogenases. *Biochem. J.*
887 **342**: 721–728.
- 888 **Ballatori, N., Christian, W.V., Lee, J.Y., Dawson, P.A., Soroka, C.J., Boyer, J.L., Madejczyk, M.S.,**
889 **and Li, N.** (2005). OST α -OST β : A major basolateral bile acid and steroid transporter in human
890 intestinal, renal, and biliary epithelia. *Hepatology* **42**: 1270–1279.
- 891 **Behnel, S., Bradshaw, R., Citro, C., Dalcin, L., Seljebotn, D.S., and Smith, K.** (2011). Cython: The
892 best of both worlds. *Comput. Sci. Eng.* **13**: 31–39.
- 893 **Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N., and**
894 **Bourne, P.E.** (2000). The Protein Data Bank. *Nucleic Acids Res.* **28**: 235–242.
- 895 **Besnard, G., Muasya, A.M., Russier, F., Roalson, E.H., Salamin, N., and Christin, P.-A.** (2009).
896 Phylogenomics of C₄ photosynthesis in sedges (Cyperaceae): Multiple appearances and genetic
897 convergence. *Mol. Biol. Evol.* **26**: 1909–1919.
- 898 **Bläsing, O.E., Westhoff, P., and Svensson, P.** (2000). Evolution of C₄ phosphoenolpyruvate carboxylase
899 in *Flaveria*, a conserved serine residue in the carboxyl-terminal part of the enzyme is a major
900 determinant for C₄-specific characteristics. *J. Biol. Chem.* **275**: 27917–27923.
- 901 **Botuyan, M.V. et al.** (2018). Mechanism of 53BP1 activity regulation by RNA-binding TIRR and a
902 designer protein. *Nat. Struct. Mol. Biol.* **25**: 591–600.
- 903 **Brunner, E. and Munzel, U.** (2000). The nonparametric Behrens-Fisher problem: asymptotic theory and a
904 small-sample approximation. *Biom. J.* **42**: 17–25.
- 905 **Bustamante, C.D. et al.** (2005). Natural selection on protein-coding genes in the human genome. *Nature*
906 **437**: 1153–1157.
- 907 **Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T.L.**
908 (2009). BLAST+: architecture and applications. *BMC Bioinformatics* **10**: 421.
- 909 **Carbone, V., Endo, S., Sumii, R., Chung, R.P.-T., Matsunaga, T., Hara, A., and El-Kabbani, O.**
910 (2008a). Structures of dimeric dihydrodiol dehydrogenase apoenzyme and inhibitor complex:
911 Probing the subunit interface with site-directed mutagenesis. *Proteins Struct. Funct. Bioinforma.*
912 **70**: 176–187.
- 913 **Carbone, V., Hara, A., and El-Kabbani, O.** (2008b). Structural and functional features of dimeric
914 dihydrodiol dehydrogenase. *Cell. Mol. Life Sci.* **65**: 1464–1474.
- 915 **Castoe, T.A., Jiang, Z.J., Gu, W., Wang, Z.O., and Pollock, D.D.** (2008). Adaptive evolution and
916 functional redesign of core metabolic proteins in snakes. *PloS One* **3**: e2201.
- 917 **Castoe, T.A., de Koning, A.P., Kim, H.M., Gu, W., Noonan, B.P., Naylor, G., Jiang, Z.J., Parkinson,**
918 **C.L., and Pollock, D.D.** (2009). Evidence for an ancient adaptive episode of convergent
919 molecular evolution. *Proc. Natl. Acad. Sci. U. S. A.* **106**: 8986–8991.
- 920 **Chandler, C.H., Chari, S., and Dworkin, I.** (2013). Does your gene need a background check? How
921 genetic background impacts the analysis of mutations, genes, and evolution. *Trends Genet.* **29**:
922 358–366.
- 923 **Chen, S., Zhou, Y., Chen, Y., and Gu, J.** (2018). fastp: an ultra-fast all-in-one FASTQ preprocessor.
924 *Bioinformatics* **34**: i884–i890.
- 925 **Chiang, J.Y.L. and Ferrell, J.M.** (2020). Up to date on cholesterol 7 alpha-hydroxylase (CYP7A1) in bile
926 acid synthesis. *Liver Res.* **4**: 47–63.

- 927 **Christin, P.-A., Salamin, N., Savolainen, V., Duvall, M.R., and Besnard, G.** (2007). C₄ photosynthesis
928 evolved in grasses via parallel adaptive genetic changes. *Curr. Biol.* **17**: 1241–1247.
- 929 **Conant, G.C. and Wolfe, K.H.** (2008). Turning a hobby into a job: How duplicated genes find new
930 functions. *Nat. Rev. Genet.* **9**: 938–950.
- 931 **Concha, C. et al.** (2019). Interplay between developmental flexibility and determinism in the evolution of
932 mimetic *Heliconius* wing patterns. *Curr. Biol.* **29**: 3996–4009.e4.
- 933 **Couture, J.-F., Legrand, P., Cantin, L., Labrie, F., Luu-The, V., and Breton, R.** (2004). Loop
934 relaxation, a mechanism that explains the reduced specificity of rabbit 20 α -hydroxysteroid
935 dehydrogenase, a member of the aldo-keto reductase superfamily. *J. Mol. Biol.* **339**: 89–102.
- 936 **Darwin, C.R.** (1859). On the origin of species by means of natural selection, or the preservation of
937 favoured races in the struggle for life 1st ed. (John Murray: London).
- 938 **DiMario, R.J. and Cousins, A.B.** (2019). A single serine to alanine substitution decreases bicarbonate
939 affinity of phosphoenolpyruvate carboxylase in C₄ *Flaveria trinervia*. *J. Exp. Bot.* **70**: 995–1004.
- 940 **Dobler, S., Dalla, S., Wagschal, V., and Agrawal, A.A.** (2012). Community-wide convergent evolution
941 in insect adaptation to toxic cardenolides by substitutions in the Na,K-ATPase. *Proc. Natl. Acad. Sci. U. S. A.* **109**: 13040–13045.
- 942 **Dunning, L.T. et al.** (2019). Lateral transfers of large DNA fragments spread functional genes among
943 grasses. *Proc. Natl. Acad. Sci. U. S. A.* **116**: 4416–4425.
- 944 **Dy, A.B.C., Langlais, P.R., Barker, N.K., Addison, K.J., Tanyaratsrisakul, S., Boitano, S.,
945 Christenson, S.A., Kraft, M., Meyers, D., Bleecker, E.R., Li, X., and Ledford, J.G.** (2021).
946 Myeloid-associated differentiation marker is a novel SP-A-associated transmembrane protein
947 whose expression on airway epithelial cells correlates with asthma severity. *Sci. Rep.* **11**: 23392.
- 948 **El-Gebali, S. et al.** (2019). The Pfam protein families database in 2019. *Nucleic Acids Res.* **47**: D427–
949 D432.
- 950 **Emms, D.M. and Kelly, S.** (2019). OrthoFinder: phylogenetic orthology inference for comparative
951 genomics. *Genome Biol.* **20**: 238.
- 952 **Emms, D.M. and Kelly, S.** (2015). OrthoFinder: solving fundamental biases in whole genome
953 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol.* **16**: 157.
- 954 **Engelmann, S., Bläsing, O.E., Westhoff, P., and Svensson, P.** (2002). Serine 774 and amino acids 296 to
955 437 comprise the major C₄ determinants of the C₄ phosphoenolpyruvate carboxylase of *Flaveria*
956 *trinervia*. *FEBS Lett.* **524**: 11–14.
- 957 **Foote, A.D. et al.** (2015). Convergent evolution of the genomes of marine mammals. *Nat. Genet.* **47**: 272–
958 275.
- 959 **Fujihara, J., Yasuda, T., Ueki, M., Iida, R., and Takeshita, H.** (2012). Comparative biochemical
960 properties of vertebrate deoxyribonuclease I. *Comp. Biochem. Physiol. B Biochem. Mol. Biol.*
961 **163**: 263–273.
- 962 **Fukushima, K. et al.** (2017). Genome of the pitcher plant *Cephalotus* reveals genetic changes associated
963 with carnivory. *Nat. Ecol. Evol.* **1**: 0059.
- 964 **Fukushima, K. and Pollock, D.D.** (2020). Amalgamated cross-species transcriptomes reveal organ-
965 specific propensity in gene expression evolution. *Nat. Commun.* **11**: 4459.
- 966 **Goldman, N. and Yang, Z.** (1994). A codon-based model of nucleotide substitution for protein-coding
967 DNA sequences. *Mol. Biol. Evol.* **11**: 725–736.
- 968 **Goldstein, R.A., Pollard, S.T., Shah, S.D., and Pollock, D.D.** (2015). Nonadaptive amino acid
969 convergence rates decrease over time. *Mol. Biol. Evol.* **32**: 1373–1381.
- 970 **Goldstein, R.A. and Pollock, D.D.** (2017). Sequence entropy of folding and the absolute rate of amino
971 acid substitutions. *Nat. Ecol. Evol.* **1**: 1923–1930.
- 972 **Grabherr, M.G. et al.** (2011). Full-length transcriptome assembly from RNA-Seq data without a
973 reference genome. *Nat. Biotechnol.* **29**: 644–652.
- 974 **Gunaratne, J., Goh, M.X., Swa, H.L.F., Lee, F.Y., Sanford, E., Wong, L.M., Hogue, K.A.,
975 Blackstock, W.P., and Okumura, K.** (2011). Protein interactions of phosphatase and tensin
976

977 homologue (PTEN) and its cancer-associated G20E mutant compared by using stable isotope
978 labeling by amino acids in cell culture-based parallel affinity purification. *J. Biol. Chem.* **286**:
979 18093–18103.

980 **Hagey, L.R., Schteingart, C.D., Rossi, S.S., Ton-Nu, H.-T., and Hofmann, A.F.** (1998). An N-acyl
981 glycyltaurine conjugate of deoxycholic acid in the biliary bile acids of the rabbit. *J. Lipid Res.* **39**:
982 2119–2124.

983 **Hansen, T.F.** (1997). Stabilizing selection and the comparative analysis of adaptation. *Evolution* **51**:
984 1341–1351.

985 **Harris, C.R. et al.** (2020). Array programming with NumPy. *Nature* **585**: 357–362.

986 **Hedges, S.B., Dudley, J., and Kumar, S.** (2006). TimeTree: a public knowledge-base of divergence times
987 among organisms. *Bioinformatics* **22**: 2971–2972.

988 **Hedges, S.B., Marin, J., Suleski, M., Paymer, M., and Kumar, S.** (2015). Tree of life reveals clock-like
989 speciation and diversification. *Mol. Biol. Evol.* **32**: 835–845.

990 **Hermans, J. and Westhoff, P.** (1992). Homologous genes for the C₄ isoform of phosphoenolpyruvate
991 carboxylase in a C₃ and a C₄ *Flaveria* species. *Mol. Gen. Genet.* **234**: 275–284.

992 **Hiller, M., Schaar, B.T., Indjeian, V.B., Kingsley, D.M., Hagey, L.R., and Bejerano, G.** (2012). A
993 “Forward Genomics” approach links genotype to phenotype using independent phenotypic losses
994 among related species. *Cell Rep.* **2**: 817–823.

995 **Hoang, D.T., Chernomor, O., von Haeseler, A., Minh, B.Q., and Vinh, L.S.** (2018). UFBoot2:
996 Improving the ultrafast bootstrap approximation. *Mol. Biol. Evol.* **35**: 518–522.

997 **Hu, Z., Sackton, T.B., Edwards, S.V., and Liu, J.S.** (2019). Bayesian detection of convergent rate
998 changes of conserved noncoding elements on phylogenetic trees. *Mol. Biol. Evol.* **36**: 1086–1100.

999 **Huerta-Cepas, J. et al.** (2007). The human phylome. *Genome Biol.* **8**: 934–941.

1000 **Huerta-Cepas, J., Serra, F., and Bork, P.** (2016). ETE 3: Reconstruction, analysis, and visualization of
1001 phylogenomic data. *Mol. Biol. Evol.* **33**: 1635–1638.

1002 **Hunter, J.D.** (2007). Matplotlib: a 2D graphics environment. *Comput. Sci. Eng.* **9**: 90–95.

1003 **Jones, D.T., Taylor, W.R., and Thornton, J.M.** (1992). The rapid generation of mutation data matrices
1004 from protein sequences. *Comput. Appl. Biosci.* **8**: 275–282.

1005 **Jumper, J. et al.** (2021). Highly accurate protein structure prediction with AlphaFold. *Nature* **596**: 583–
1006 589.

1007 **Kaessmann, H.** (2010). Origins, evolution, and phenotypic impact of new genes. *Genome Res.* **20**: 1313–
1008 1326.

1009 **Karageorgi, M. et al.** (2019). Genome editing retraces the evolution of toxin resistance in the monarch
1010 butterfly. *Nature* **574**: 409–412.

1011 **Katoh, K. and Standley, D.M.** (2013). MAFFT multiple sequence alignment software version 7:
1012 improvements in performance and usability. *Mol. Biol. Evol.* **30**: 772–780.

1013 **Khabbazian, M., Kriebel, R., Rohe, K., and Ané, C.** (2016). Fast and accurate detection of evolutionary
1014 shifts in Ornstein-Uhlenbeck models. *Methods Ecol. Evol.* **7**: 811–824.

1015 **Kimura, M.** (1968). Evolutionary rate at the molecular level. *Nature* **217**: 624–626.

1016 **Kleene, K.C.** (2005). Sexual selection, genetic conflict, selfish genes, and the atypical patterns of gene
1017 expression in spermatogenic cells. *Dev. Biol.* **277**: 16–26.

1018 **Knott, G.J. and Doudna, J.A.** (2018). CRISPR-Cas guides the future of genetic engineering. *Science*
1019 **361**: 866–869.

1020 **Konečná, V., Bray, S., Vlček, J., Bohutínská, M., Požárová, D., Choudhury, R.R., Bollmann-Giolai,
1021 A., Flis, P., Salt, D.E., Parisod, C., Yant, L., and Kolář, F.** (2021). Parallel adaptation in
1022 autopolyploid *Arabidopsis arenosa* is dominated by repeated recruitment of shared alleles. *Nat.*
1023 *Commun.* **12**: 4979.

1024 **Kosiol, C., Holmes, I., and Goldman, N.** (2007). An empirical codon model for protein sequence
1025 evolution. *Mol. Biol. Evol.* **24**: 1464–1479.

1026 **Kowalczyk, A., Meyer, W.K., Partha, R., Mao, W., Clark, N.L., and Chikina, M.** (2019).

1027 RERconverge: an R package for associating evolutionary rates with convergent traits.
1028 Bioinformatics **35**: 4815–4817.

1029 **Krogh, A., Larsson, B., von Heijne, G., and Sonnhammer, E.L.L.** (2001). Predicting transmembrane
1030 protein topology with a hidden markov model: application to complete genomes. *J. Mol. Biol.*
1031 **305**: 567–580.

1032 **Land, A.H. and Doig, A.G.** (1960). An automatic method of solving discrete programming problems.
1033 *Econometrica* **28**: 497–520.

1034 **Lartillot, N. and Philippe, H.** (2004). A Bayesian mixture model for across-site heterogeneities in the
1035 amino-acid replacement process. *Mol. Biol. Evol.* **21**: 1095–1109.

1036 **Lewin, H.A. et al.** (2022). The Earth BioGenome Project 2020: Starting the clock. *Proc. Natl. Acad. Sci.*
1037 *U. S. A.* **119**: e2115635118.

1038 **Liu, Y., Cotton, J.A., Shen, B., Han, X., Rossiter, S.J., and Zhang, S.** (2010). Convergent sequence
1039 evolution between echolocating bats and dolphins. *Curr. Biol.* **20**: R53–R54.

1040 **Liu, Z., Qi, F.-Y., Zhou, X., Ren, H.-Q., and Shi, P.** (2014). Parallel sites implicate functional
1041 convergence of the hearing gene *prestin* among echolocating mammals. *Mol. Biol. Evol.* **31**:
1042 2415–2424.

1043 **Losos, J.B.** (2017). *Improbable Destinies: Fate, Chance, and the Future of Evolution* (Riverhead Books:
1044 New York).

1045 **Lyons, D.M., Zou, Z., Xu, H., and Zhang, J.** (2020). Idiosyncratic epistasis creates universals in
1046 mutational effects and evolutionary trajectories. *Nat. Ecol. Evol.* **4**: 1685–1693.

1047 **Marcovitz, A., Jia, R., and Bejerano, G.** (2016). “Reverse Genomics” predicts function of human
1048 conserved noncoding elements. *Mol. Biol. Evol.* **33**: 1358–1369.

1049 **Marcovitz, A., Turakhia, Y., Chen, H.I., Gloudemans, M., Braun, B.A., Wang, H., and Bejerano, G.**
1050 (2019). A functional enrichment test for molecular convergent evolution finds a clear protein-
1051 coding signal in echolocating bats and whales. *Proc. Natl. Acad. Sci. U. S. A.* **116**: 21094–21103.

1052 **Martin, A. and Orgogozo, V.** (2013). The loci of repeated evolution: a catalog of genetic hotspots of
1053 phenotypic variation. *Evol. Int. J. Org. Evol.* **67**: 1235–1250.

1054 **Mendes, F.K., Hahn, Y., and Hahn, M.W.** (2016). Gene tree discordance can generate patterns of
1055 diminishing convergence over time. *Mol. Biol. Evol.* **33**: 3299–3307.

1056 **Mendes, F.K., Livera, A.P., and Hahn, M.W.** (2019). The perils of intralocus recombination for
1057 inferences of molecular convergence. *Philos. Trans. R. Soc. B Biol. Sci.* **374**: 20180244.

1058 **Minh, B.Q., Nguyen, M.A.T., and von Haeseler, A.** (2013). Ultrafast approximation for phylogenetic
1059 bootstrap. *Mol. Biol. Evol.* **30**: 1188–1195.

1060 **Minh, B.Q., Schmidt, H.A., Chernomor, O., Schrempf, D., Woodhams, M.D., von Haeseler, A., and**
1061 **Lanfear, R.** (2020). IQ-TREE 2: New models and efficient methods for phylogenetic inference in
1062 the genomic era. *Mol. Biol. Evol.* **37**: 1530–1534.

1063 **Morel, B., Kozlov, A.M., Stamatakis, A., and Szöllösi, G.J.** (2020). GeneRax: A tool for species-tree-
1064 aware maximum likelihood-based gene family tree inference under gene duplication, transfer, and
1065 loss. *Mol. Biol. Evol.* **37**: 2763–2774.

1066 **Mulchande, J., Martins, L., Moreira, R., Archer, M., Oliveira, T.F., and Iley, J.** (2007). The efficiency
1067 of C-4 substituents in activating the β -lactam scaffold towards serine proteases and hydroxide ion.
1068 *Org. Biomol. Chem.* **5**: 2617–2626.

1069 **Muñoz-Clares, R.A., González-Segura, L., Juárez-Díaz, J.A., and Mújica-Jiménez, C.** (2020).
1070 Structural and biochemical evidence of the glucose 6-phosphate-allosteric site of maize C₄-
1071 phosphoenolpyruvate carboxylase: its importance in the overall enzyme kinetics. *Biochem. J.* **477**:
1072 2095–2114.

1073 **Muse, S.V. and Gaut, B.S.** (1994). A likelihood approach for comparing synonymous and
1074 nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. *Mol.*
1075 *Biol. Evol.* **11**: 715–724.

1076 **Nakagawa, M., Matsuura, K., Hara, A., Sawada, H., Bunai, Y., and Ohya, I.** (1989). Dimeric

1077 dihydrodiol dehydrogenase in monkey kidney. Substrate specificity, stereospecificity of hydrogen
1078 transfer, and distribution. *J. Biochem. (Tokyo)* **106**: 1104–1109.

1079 **Nakayama, T., Sawada, H., Deyashiki, Y., Kanazu, T., Hara, A., Shinoda, M., Matsuura, K., Bunai,**
1080 **Y., and Ohya, I.** (1991). Distribution of dimeric dihydrodiol dehydrogenase in pig tissues and its
1081 role in carbonyl metabolism. In *Enzymology and Molecular Biology of Carbonyl Metabolism 3*,
1082 H. Weiner, B. Wermuth, and D.W. Crabb, eds, *Advances in Experimental Medicine and Biology*.
1083 (Springer US: Boston, MA), pp. 187–196.

1084 **Natarajan, C., Hoffmann, F.G., Weber, R.E., Fago, A., Witt, C.C., and Storz, J.F.** (2016). Predictable
1085 convergence in hemoglobin function has unpredictable molecular underpinnings. *Science* **354**:
1086 336–339.

1087 **Nguyen, L.-T., Schmidt, H.A., von Haeseler, A., and Minh, B.Q.** (2015). IQ-TREE: A fast and effective
1088 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**: 268–
1089 274.

1090 **Nishikawa, A., Gregory, W., Frenz, J., Cacia, J., and Kornfeld, S.** (1997). The phosphorylation of
1091 bovine DNase I Asn-linked oligosaccharides is dependent on specific lysine and arginine residues.
1092 *J. Biol. Chem.* **272**: 19408–19412.

1093 **Noble, R.C.** (1981). Digestion, absorption and transport of lipids in ruminant animals. In *Lipid*
1094 *Metabolism in Ruminant Animals*, W.W. Christie, ed (Pergamon), pp. 57–93.

1095 **Ohta, T.** (1973). Slightly deleterious mutant substitutions in evolution. *Nature* **246**: 96–98.

1096 **Parker, J., Tsagkogeorga, G., Cotton, J.A., Liu, Y., Provero, P., Stupka, E., and Rossiter, S.J.** (2013).
1097 Genome-wide signatures of convergent evolution in echolocating mammals. *Nature* **502**: 228–231.

1098 **Patel, A., Yang, P., Tinkham, M., Pradhan, M., Sun, M.-A., Wang, Y., Hoang, D., Wolf, G., Horton,**
1099 **J.R., Zhang, X., Macfarlan, T., and Cheng, X.** (2018). DNA conformation induces adaptable
1100 binding by tandem zinc finger proteins. *Cell* **173**: 221–233.e12.

1101 **Poetsch, W., Hermans, J., and Westhoff, P.** (1991). Multiple cDNAs of phosphoenolpyruvate
1102 carboxylase in the C₄ dicot *Flaveria trinervia*. *FEBS Lett.* **292**: 133–136.

1103 **Pollock, D.D. and Pollard, S.T.** (2016). Parallel and convergent molecular evolution. In *Encyclopedia of*
1104 *Evolutionary Biology*, R.M. Kliman, ed (Academic Press: Oxford), pp. 206–211.

1105 **Pollock, D.D., Thiltgen, G., and Goldstein, R.A.** (2012). Amino acid coevolution induces an
1106 evolutionary Stokes shift. *Proc. Natl. Acad. Sci. U. S. A.* **109**: E1352–E1359.

1107 **Pond, S.L.K. and Frost, S.D.W.** (2005). Not so different after all: A comparison of methods for detecting
1108 amino acid sites under selection. *Mol. Biol. Evol.* **22**: 1208–1222.

1109 **Prudent, X., Parra, G., Schwede, P., Roscito, J.G., Hiller, M., DW, L., CR, B., JB, S., N, M.-P., and**
1110 **MA., M.** (2016). Controlling for phylogenetic relatedness and evolutionary rates improves the
1111 discovery of associations between species' phenotypic and genomic differences. *Mol. Biol. Evol.*
1112 **33**: 2135–2150.

1113 **Rey, C., Guéguen, L., Sémon, M., and Boussau, B.** (2018). Accurate detection of convergent amino-acid
1114 evolution with PCOC. *Mol. Biol. Evol.* **35**: 2296–2306.

1115 **Rice, P., Longden, I., and Bleasby, A.** (2000). EMBOSS: the European Molecular Biology Open
1116 Software Suite. *Trends Genet.* **16**: 276–277.

1117 **Rižner, T.L. and Penning, T.M.** (2014). Role of aldo–keto reductase family 1 (AKR1) enzymes in human
1118 steroid metabolism. *Steroids* **79**: 49–63.

1119 **Sage, R.F., Christin, P.-A., and Edwards, E.J.** (2011). The C₄ plant lineages of planet Earth. *J. Exp. Bot.*
1120 **62**: 3155–3169.

1121 **Sato, K., Inazu, A., Yamaguchi, S., Nakayama, T., Deyashiki, Y., Sawada, H., and Hara, A.** (1993).
1122 Monkey 3-deoxyglucosone reductase: Tissue distribution and purification of three multiple forms
1123 of the kidney enzyme that are identical with dihydrodiol dehydrogenase, aldehyde reductase, and
1124 aldose reductase. *Arch. Biochem. Biophys.* **307**: 286–294.

1125 **Sato, K., Nakanishi, M., Deyashiki, Y., Hara, A., Matsuura, K., and Ohya, I.** (1994). Purification and
1126 characterization of dimeric dihydrodiol dehydrogenase from dog liver. *J. Biochem. (Tokyo)* **116**:

1127 711–717.

1128 **Schoch, C.L. et al.** (2020). NCBI Taxonomy: a comprehensive update on curation, resources and tools.
 1129 Database **2020**: baaa062.

1130 **Shah, P., McCandlish, D.M., and Plotkin, J.B.** (2015). Contingency and entrenchment in protein
 1131 evolution under purifying selection. *Proc. Natl. Acad. Sci. U. S. A.* **112**: E3226–E3235.

1132 **Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., and Zdobnov, E.M.** (2015).
 1133 BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs.
 1134 *Bioinformatics* **31**: 3210–3212.

1135 **Sirikantaramas, S., Yamazaki, M., and Saito, K.** (2008). Mutations in topoisomerase I as a self-
 1136 resistance mechanism coevolved with the production of the anticancer alkaloid camptothecin in
 1137 plants. *Proc. Natl. Acad. Sci. U. S. A.* **105**: 6782–6786.

1138 **Spielman, S.J. and Wilke, C.O.** (2015). Pyvolve: A flexible python module for simulating sequences
 1139 along phylogenies. *PLOS ONE* **10**: e0139047.

1140 **Starr, T.N., Flynn, J.M., Mishra, P., Bolon, D.N.A., and Thornton, J.W.** (2018). Pervasive contingency
 1141 and entrenchment in a billion years of Hsp90 evolution. *Proc. Natl. Acad. Sci. U. S. A.* **115**: 4453–
 1142 4458.

1143 **Steenwyk, J.L., Iii, T.J.B., Li, Y., Shen, X.-X., and Rokas, A.** (2020). ClipKIT: A multiple sequence
 1144 alignment trimming software for accurate phylogenomic inference. *PLOS Biol.* **18**: e3001007.

1145 **Steinegger, M. and Söding, J.** (2017). MMseqs2 enables sensitive protein sequence searching for the
 1146 analysis of massive data sets. *Nat. Biotechnol.* **35**: 1026–1028.

1147 **Stern, D.L.** (2013). The genetic causes of convergent evolution. *Nat. Rev. Genet.* **14**: 751–764.

1148 **Stewart, C.B., Schilling, J.W., and Wilson, A.C.** (1987). Adaptive evolution in the stomach lysozymes
 1149 of foregut fermenters. *Nature* **330**: 401–404.

1150 **Storz, J.F.** (2016). Causes of molecular convergence and parallelism in protein evolution. *Nat. Rev. Genet.*
 1151 **17**: 239–250.

1152 **Sun, L., Bai, Y., Zhao, R., Sun, T., Cao, R., Wang, F., He, G., Zhang, W., Chen, Y., Ye, P., and Du,
 1153 G.** (2016). Oncological miR-182-3p, a novel smooth muscle cell phenotype modulator, evidences
 1154 from model rats and patients. *Arterioscler. Thromb. Vasc. Biol.* **36**: 1386–1397.

1155 **Taverner, A.M. et al.** (2019). Adaptive substitutions underlying cardiac glycoside insensitivity in insects
 1156 exhibit epistasis in vivo. *eLife* **8**: e48224.

1157 **Thirawatananond, P., McPherson, R.L., Malhi, J., Nathan, S., Lambrecht, M.J., Brichacek, M.,
 1158 Hergenrother, P.J., Leung, A.K.L., and Gabelli, S.B.** (2019). Structural analyses of NudT16–
 1159 ADP-ribose complexes direct rational design of mutants with improved processing of poly(ADP-
 1160 ribosyl)ated proteins. *Sci. Rep.* **9**: 5940.

1161 **Thomas, G.W.C. and Hahn, M.W.** (2015). Determining the null model for detecting adaptive
 1162 convergence from genomic data: A case study using echolocating mammals. *Mol. Biol. Evol.* **32**:
 1163 1232–1236.

1164 **Thomas, G.W.C., Hahn, M.W., and Hahn, Y.** (2017). The effects of increasing the number of taxa on
 1165 inferences of molecular convergence. *Genome Biol. Evol.* **9**: 213–221.

1166 **Tien, M.Z., Meyer, A.G., Sydykova, D.K., Spielman, S.J., and Wilke, C.O.** (2013). Maximum allowed
 1167 solvent accessibilities of residues in proteins. *PloS One* **8**: e80635.

1168 **Ujvari, B., Casewell, N.R., Sunagar, K., Arbuckle, K., Wüster, W., Lo, N., O’Meally, D., Beckmann,
 1169 C., King, G.F., Deplazes, E., and Madsen, T.** (2015). Widespread convergence in toxin
 1170 resistance by predictable molecular evolution. *Proc. Natl. Acad. Sci. U. S. A.* **112**: 11911–11916.

1171 **Varadi, M. et al.** (2022). AlphaFold Protein Structure Database: massively expanding the structural
 1172 coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res.* **50**: D439–
 1173 D444.

1174 **Vermeij, G.J.** (2006). Historical contingency and the purported uniqueness of evolutionary innovations.
 1175 *Proc. Natl. Acad. Sci. U. S. A.* **103**: 1804–1809.

1176 **Virtanen, P. et al.** (2020). SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat.*

1177 Methods **17**: 261–272.

1178 **Wang, L. et al.** (2019). A draft genome assembly of halophyte *Suaeda aralocaspica*, a plant that performs
1179 C₄ photosynthesis within individual cells. *GigaScience* **8**: giz116.

1180 **Wang, Q. et al.** (2021). Rare variant contribution to human disease in 281,104 UK Biobank exomes.
1181 *Nature* **597**: 527–532.

1182 **Waskom, M.L.** (2021). seaborn: statistical data visualization. *J. Open Source Softw.* **6**: 3021.

1183 **Weston, S.A., Lahm, A., and Suck, D.** (1992). X-ray structure of the DNase I-d(GGTATACC)₂ complex
1184 at 2.3Å resolution. *J. Mol. Biol.* **226**: 1237–1256.

1185 **Whelan, S. and Goldman, N.** (2001). A general empirical model of protein evolution derived from
1186 multiple protein families using a maximum-likelihood approach. *Mol. Biol. Evol.* **18**: 691–699.

1187 **Wickham, H.** (2009). ggplot2: Elegant Graphics for Data Analysis (Springer-Verlag: New York).

1188 **Yanai, I., Benjamin, H., Shmoish, M., Chalifa-Caspi, V., Shklar, M., Ophir, R., Bar-Even, A., Horn-
1189 Saban, S., Safran, M., Domany, E., Lancet, D., and Shmueli, O.** (2005). Genome-wide
1190 midrange transcription profiles reveal expression level relationships in human tissue specification.
1191 *Bioinformatics* **21**: 650–659.

1192 **Yang, L., Ravikanthachari, N., Mariño-Pérez, R., Deshmukh, R., Wu, M., Rosenstein, A., Kunte, K.,
1193 Song, H., and Andolfatto, P.** (2019a). Predictability in the evolution of Orthopteran cardenolide
1194 insensitivity. *Philos. Trans. R. Soc. B Biol. Sci.* **374**: 20180246.

1195 **Yang, Z.** (2006). Computational Molecular Evolution (Oxford University Press: Oxford, UK).

1196 **Yang, Z. et al.** (2019b). Convergent horizontal gene transfer and cross-talk of mobile nucleic acids in
1197 parasitic plants. *Nat. Plants* **5**: 991–1001.

1198 **Yang, Z.** (2007). PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**: 1586–
1199 1591.

1200 **Yates, A. et al.** (2016). Ensembl 2016. *Nucleic Acids Res.* **44**: D710–D716.

1201 **Yu, G., Smith, D.K., Zhu, H., Guan, Y., and Lam, T.T.-Y.** (2017). GGTREE: an R package for
1202 visualization and annotation of phylogenetic trees with their covariates and other associated data.
1203 *Methods Ecol. Evol.* **8**: 28–36.

1204 **Zhang, F., Lou, L., Peng, B., Song, X., Reizes, O., Almasan, A., and Gong, Z.** (2020). Nudix hydrolase
1205 NUDT16 regulates 53BP1 protein by reversing 53BP1 ADP-ribosylation. *Cancer Res.* **80**: 999–
1206 1010.

1207 **Zhang, G.-Q. et al.** (2017). The *Apostasia* genome and the evolution of orchids. *Nature* **549**: 379–383.

1208 **Zhang, J.** (2006). Parallel adaptive origins of digestive RNases in Asian and African leaf monkeys. *Nat
1209 Genet* **38**: 819–823.

1210 **Zhang, J. and Kumar, S.** (1997). Detection of convergent and parallel evolution at the amino acid
1211 sequence level. *Mol. Biol. Evol.* **14**: 527–536.

1212 **Zhang, J. and Yang, J.-R.** (2015). Determinants of the rate of protein sequence evolution. *Nat. Rev.
1213 Genet.* **16**: 409–420.

1214 **Zhen, Y., Aardema, M.L., Medina, E.M., Schumer, M., and Andolfatto, P.** (2012). Parallel molecular
1215 evolution in an herbivore community. *Science* **337**: 1634–1637.

1216 **Zou, Z. and Zhang, J.** (2015a). Are convergent and parallel amino acid substitutions in protein evolution
1217 more prevalent than neutral expectations? *Mol. Biol. Evol.* **32**: 2085–2096.

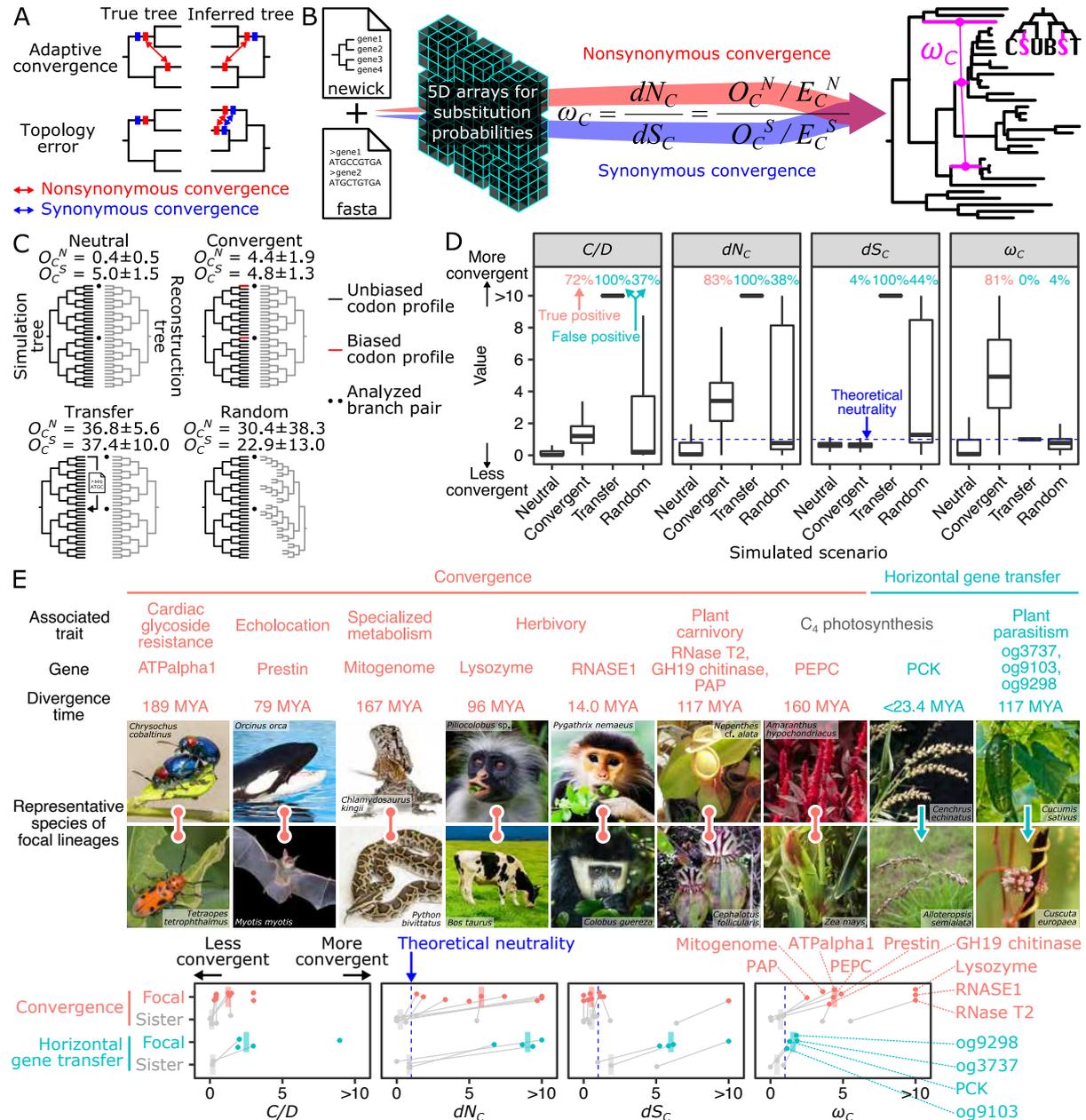
1218 **Zou, Z. and Zhang, J.** (2017). Gene tree discordance does not explain away the temporal decline of
1219 convergence in mammalian protein sequence evolution. *Mol. Biol. Evol.* **34**: 1682–1688.

1220 **Zou, Z. and Zhang, J.** (2015b). No genome-wide protein sequence convergence for echolocation. *Mol.
1221 Biol. Evol.* **32**: 1237–1241.

1222

1223

Figures



1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

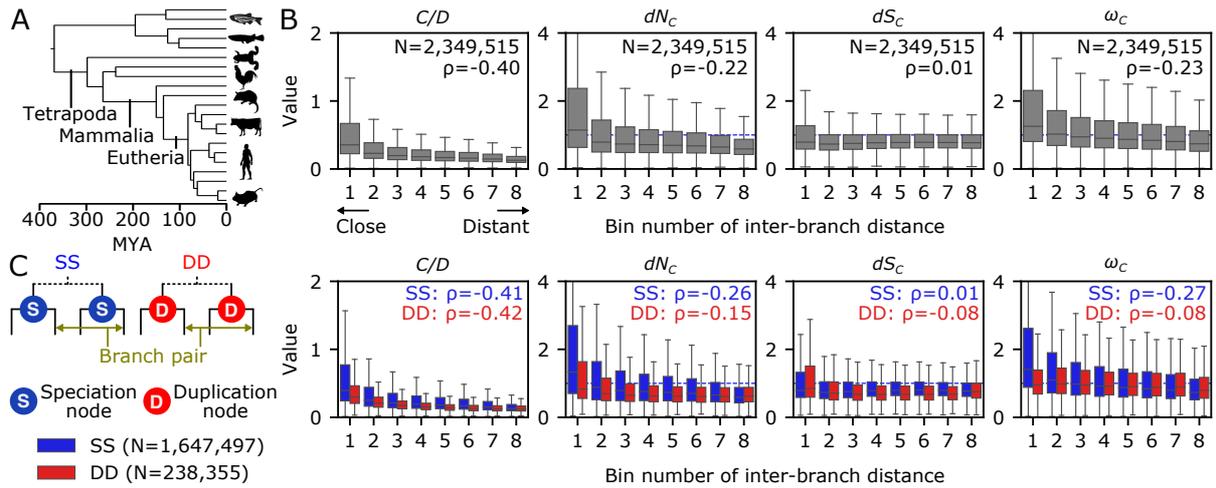
1237

1238

1239

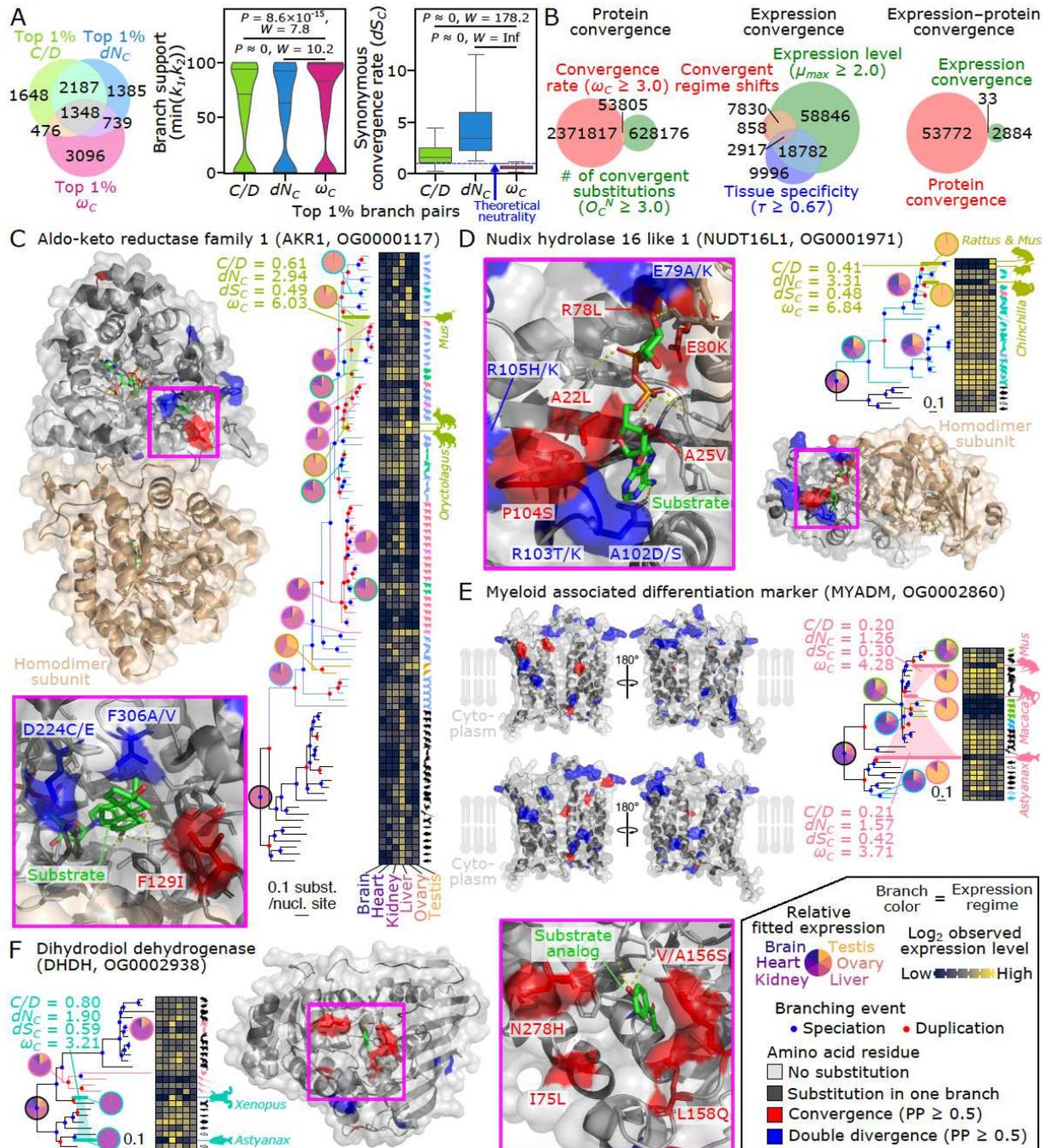
Figure 1. Challenges and solutions for the detection of molecular convergence. (A) False convergence is caused by tree topology errors. (B) The overview of CSUBST. This program processes substitution probabilities to derive observed (O_c^N and O_c^S) and expected (E_c^N and E_c^S) numbers of nonsynonymous and synonymous convergence and evaluate their rates (dN_c and dS_c) in branch combinations in a phylogenetic tree. A more detailed illustration is available in Fig. S2. (C) Generation of simulated datasets for performance evaluation in different evolutionary scenarios. The ECMK07+F codon substitution model was used to simulate the evolution of 500-codon sequences on a phylogenetic tree with 32 leaves 1,000 times. The numbers of observed nonsynonymous and synonymous convergence are indicated above trees (O_c^N and O_c^S , respectively; mean \pm standard deviation). (D) The estimated rates of protein convergence in different scenarios. Each box plot corresponds to the results of 1,000 simulations. Dashed lines indicate the neutral expectation (=1.0) except for C/D (Castoe et al., 2009; Goldstein et al., 2015), for which no theoretical expectation is available. dN_c is largely equivalent to the previously proposed metric called R (Zou and Zhang, 2015a). Values greater than the 95th percentile in the Neutral scenario are defined as true and false positives in Convergent and other scenarios, respectively, and are indicated at the top of the plot. The positive

1240 rate of dS_C is interpreted as a false positive rate even in the Convergent scenario because the probability of
1241 only nonsynonymous substitutions is manipulated. (E) Performance of convergence metrics in empirical
1242 datasets. Known examples of protein convergences and horizontal gene transfers (HGTs) are analyzed with
1243 C/D , dN_C , dS_C , and ω_C . Median values (bars) are overlaid on individual data points that correspond to gene
1244 trees. In trees where convergence occurred in more than two lineages, the median of all foreground branch
1245 pairs is reported. The branch pairs sister to the focal branches are shown as a control (Foote et al., 2015),
1246 except in cases where there is no substitution at all or the sister branches are phylogenetically not
1247 independent. Dataset and photographs of representative species are shown above the plot. The taxonomic
1248 range follows the NCBI Taxonomy database (Schoch et al., 2020), and the divergence time is according to
1249 timetree.org (Hedges et al., 2015). The lineages involving adaptive convergence or HGTs are referred to as
1250 focal lineages. The gene trees are illustrated in Fig. S5 and Fig. S6. The comparison with the background
1251 levels for each dataset is shown in Fig. S4. The characteristics of the datasets are summarized in Table S3.
1252 The photograph of *Alloteropsis semialata* is licensed under CC BY-SA 3.0
1253 (<https://creativecommons.org/licenses/by-sa/3.0/>) by Marjorie Lundgren.
1254



1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267

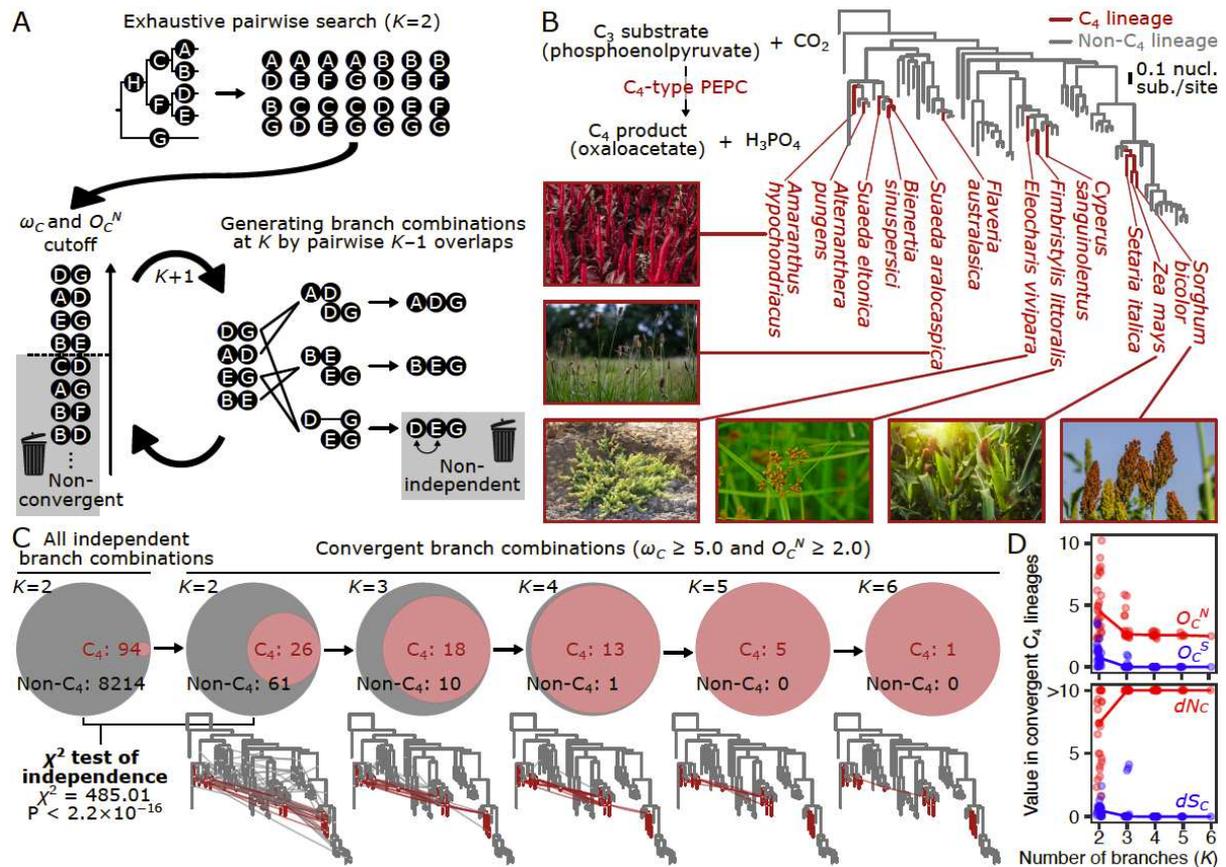
Figure 2. Biological variation of ω_c in a genome-scale dataset. (A) Phylogenetic relationships of the selected species. See Fig. S7A for the complete phylogeny. The tree and divergence time estimates were obtained from timetree.org (Hedges et al., 2015). Some animal silhouettes were obtained from PhyloPic (<http://phylopic.org>). (B) Temporal variation of convergence rates. The numbers of branch pairs (N) and Spearman's correlation coefficient (ρ) are shown. The bin range was determined to assign an equal number of branch pairs to each bin. To reduce the noise originating from branches where almost no substitutions occurred, branch pairs with both O_c^N and O_c^S greater than 1 were analyzed (i.e., at least one convergent substitution each). (C) Convergence rates depending on gene duplications. Branch pairs were categorized into speciation events (SS) and branch pairs after two independent gene duplications (DD) according to the presence of preceding gene duplications in no or both branches, respectively. Branch pairs with one preceding duplication were excluded from the analysis. Dashed lines indicate the neutral expectation (=1.0).



1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281

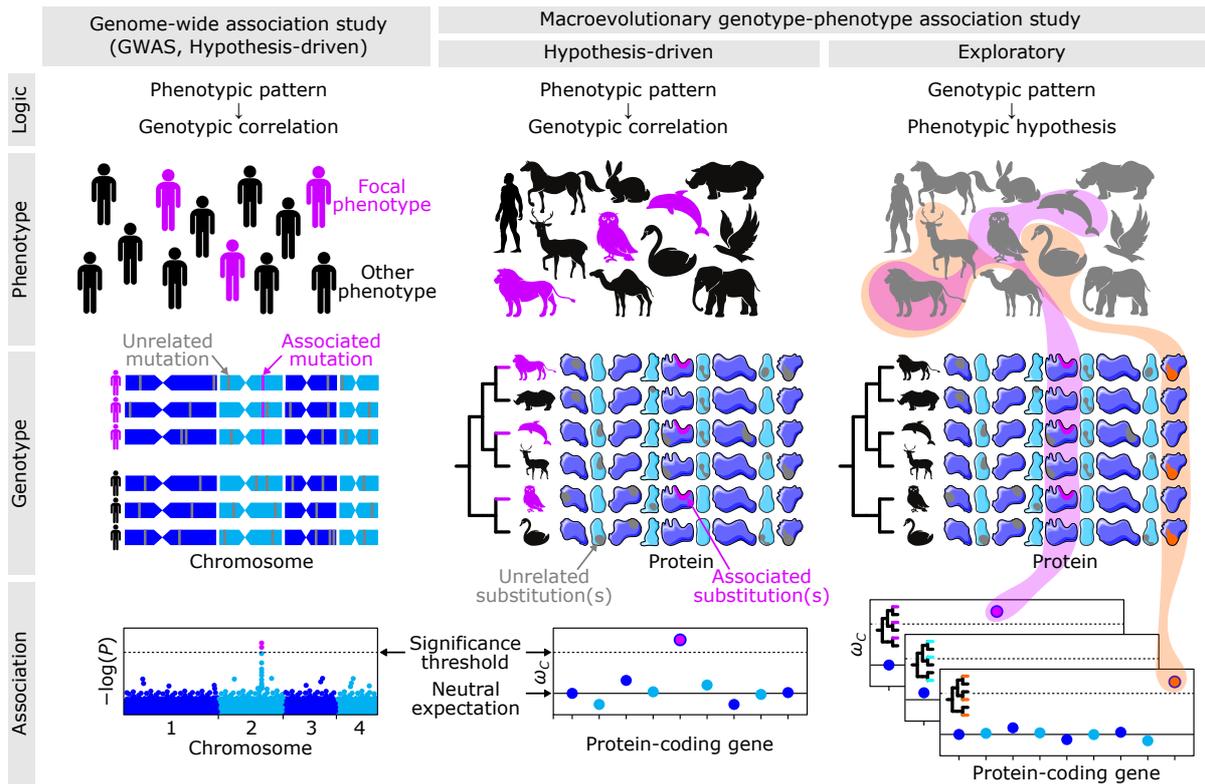
Figure 3. Joint convergence of gene expression patterns and protein sequences. (A) Comparison of convergent branch pairs obtained by different methods in the vertebrate dataset. Branch pairs with $O_C^N \geq 3.0$ and $O_C^S \geq 3.0$ were analyzed. The Venn diagram on the left shows the extent of overlap between the top 1% convergent branch pairs. The violin plot in the middle shows the lower bootstrap support of the parental branches of the convergent branch pairs. The boxplot on the right compares the rate of synonymous convergence (dS_C). The stochastic equality of data was tested by a two-sided Brunner–Munzel test (Brunner and Munzel, 2000). (B) Venn diagrams showing the extent of overlap between protein and expression convergence. Circles represent the sets of branch pairs. Shifts in tissue-specific expression regime were identified with the thresholds of expression levels (the maximum fitted SVA-log-TMM-FPKM among tissues (Fukushima and Pollock, 2020)) and tissue specificity (Yanai's τ (Yanai et al., 2005)). (C–F) Examples of the likely adaptive joint convergence. Aldo-keto reductase family 1 (AKR1, C), Nudix hydrolase 16 like 1 (NUDT16L1, D), Myeloid associated differentiation marker (MYADM, E), and Dihydrodiol dehydrogenase (DHDH, F) are shown (see Fig. S9A for complete trees). Node colors in the

1282 trees indicate inferred branching events of speciation (blue) and gene duplication (red). The heatmap shows
1283 expression levels observed in extant species. The silhouettes signify the species (see Fig. S7A) that carries
1284 the gene, and the clades involved in the joint convergence are indicated with an enlarged size. The colors of
1285 branches and animal silhouettes indicate expression regimes. Among-organ expression patterns are shown
1286 as a pie chart for each regime. Branches involved in joint convergence are highlighted with thick lines,
1287 connected by the color of the expression regime, and annotated with convergence metrics. Localization of
1288 convergent and divergent substitutions on the protein structure is shown along with a close-up view of
1289 functionally important sites. The surface representation of each protein is overlaid with a cartoon
1290 representation. Convergent and divergent amino acid loci shown in Fig. S9 are highlighted in red and blue,
1291 respectively. Substrates and their analogs are shown as green sticks. Side chains forming the substrate-
1292 binding site are also shown as sticks. Note that these are the side chains in the protein from databases, so
1293 amino acid substitutions in the convergent lineages may result in distinct structures and arrangements. Site
1294 numbers correspond to those in the PDB entry or the AlphaFold structure (from C to F: 1Q13, 5W6X, AF-
1295 Q6DFR5-F1-model_v2, and 2O48). The silhouettes of *Astyanax mexicanus* and *Oreochromis niloticus* are
1296 licensed under CC BY-NC-SA 3.0 (<https://creativecommons.org/licenses/by-nc-sa/3.0/>) by Milton Tan
1297 (reproduced with permission), and those of *Anolis carolinensis* (by Sarah Werning), *Ornithorhynchus*
1298 *anatinus* (by Sarah Werning), and *Rattus norvegicus* (by Rebecca Groom; with modification) are licensed
1299 under CC BY 3.0 (<https://creativecommons.org/licenses/by/3.0/>).
1300



1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315

Figure 4. Heuristic search of higher-order branch combinations for adaptive protein convergence. (A) Branch-and-bound algorithm for higher-order branch combinations. This method explores the higher-order combinatorial space until there are no more convergent branch combinations. (B) The maximum-likelihood phylogenetic tree of phosphoenolpyruvate carboxylases (PEPCs) in flowering plants. The catalytic function of PEPC, which is crucial in C_4 photosynthesis, is illustrated. Photographs of representative C_4 photosynthetic lineages are shown. The photograph of *Suaeda aralocaspica* is reproduced from the literature (Wang et al., 2019). The bar indicates 0.1 nucleotide substitutions per nucleotide site. The complete tree is shown in Fig. S5. (C) Higher-order convergence enriches C_4 -type PEPCs. The Venn diagrams show the proportion of convergent branch combinations of C_4 -type and non- C_4 -type lineages (red and gray, respectively). Branch combinations containing both were included in non- C_4 . In the phylogenetic trees, convergent branch combinations are shown as edges connecting branches. (D) Improvement of the signal-to-noise ratio in higher-order branch combinations. The line graph shows the median values of the total probabilities (O_c^N and O_c^S) and the rates (dN_c and dS_c) of nonsynonymous and synonymous convergence in the convergent branch combinations of C_4 lineages. Points correspond to branch combinations.



1316
 1317
 1318
 1319
 1320
 1321
 1322

Figure 5. Analysis of the genotype-phenotype association within and between species. The proposed method improves the accuracy of the hypothesis-driven approach in the macroevolutionary scale and enables exploratory approaches. Note that for visualization purposes, the number of individuals and species shown here is smaller than the actual number required for analysis. The icons of proteins are licensed under CC BY 3.0 (<https://creativecommons.org/licenses/by/3.0/>) by Smart Servier Medical Art.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supptables.xlsx](#)
- [supp.pdf](#)