



Article

Phylogenetic weighting does little to improve the accuracy of evolutionary coupling analyses

Adam J. Hockenberry¹* and Claus O. Wilke¹

¹ Department of Integrative Biology, The University of Texas at Austin

* Correspondence: adam.hockenberry@utexas.edu

Version August 14, 2019 submitted to Journal Not Specified

Abstract: Homologous sequence alignments contain important information about the constraints that shape protein family evolution. Correlated changes between different residues, for instance, can be highly predictive of physical contacts within three-dimensional structures. Detecting such co-evolutionary signals via direct coupling analysis is particularly challenging given the shared phylogenetic history and uneven sampling of different lineages from which protein sequences are derived. Current best practices for mitigating such effects include sequence-identity-based weighting of input sequences and *post-hoc* re-scaling of evolutionary coupling scores. However, numerous weighting schemes have been previously developed for other applications, and it is unknown whether any of these schemes may better account for phylogenetic artifacts in evolutionary coupling analyses. Here, we show across a dataset of 150 diverse protein families that the current best practices out-perform several alternative sequence- and tree-based weighting methods. Nevertheless, we find that sequence weighting in general provides only a minor benefit relative to *post-hoc* transformations that re-scale the derived evolutionary couplings. While our findings do not rule out the possibility that an as-yet-untested weighting method may show improved results, the similar predictive accuracies that we observe across distinct weighting methods suggests that there may be little room for further improvement on top of existing strategies.

Keywords: direct coupling analysis; evolutionary coupling analysis; contact prediction; phylogenetic bias

1. Introduction

Correlated evolution of amino acid positions within a sequence alignment can be leveraged to inform structural models of proteins, predict mutational effects, and identify protein binding partners [1–5]. The ability to detect correlated evolution has been revolutionized by direct coupling analyses and other related methods that seek to re-construct one- and two-site marginal amino acid probabilities based on the observed distribution of sequence data [6–11]. Inference of two-site coupling parameters from a multiple sequence alignment is technically challenging, however, and numerous related approaches have been developed in recent years [9,10,12–17]. This intense focus on related methodologies stems from the fact that the highest scoring evolutionary coupling values are highly enriched in residue-residue pairs whose side-chains physically interact within three dimensional structures [18]. Evolutionary couplings can thus provide valuable information about structural constraints within and between protein families, while only requiring sequence information as inputs [15,19–22].

All methods to detect correlated evolution between different positions in a protein family require large numbers of representative sequences and therefore start by finding—and subsequently aligning—homologous sequences from large sequence databases [5]. An oft-remarked upon fact is that sequence databases are composed of a highly biased sample of life on earth; some species are

36 much more densely sampled than others (as are some genera, families, orders, *etc.*) [23–27]. Even if all
37 extant life were equally well sampled and represented in sequence databases, species are related by
38 complicated historical patterns and cannot be considered as independent observations [28].

39 Statistical issues arising from this shared phylogenetic history and biased sampling have long been
40 noted by biologists [28]. The problem can be most clearly summarized by a toy example. In Figure 1A,
41 we show a hypothetical sequence alignment and ask the question: What amino acid is preferred at the
42 indicated site? At first glance, a phylogenetically agnostic method would simply count the frequency of
43 different amino acids and conclude that valine (V, four occurrences) is preferred. However, accounting
44 for phylogenetic relationships, a different perspective could reasonably conclude that threonine (T,
45 three occurrences) is more highly preferred given that it occupies a substantially larger fraction of
46 the phylogenetic tree and therefore dominates the evolutionary history of the protein family; the
47 abundance of valines in the alignment is an apparent result of over-sampling one closely related
48 lineage (which may represent numerous representatives of the same species, for example). Naively, the
49 problem can be solved by simply selecting a single member from each species to prevent over-sampling.
50 However, the issue remains equally problematic at other taxonomic levels (i.e. sampling numerous
51 species from the same genus, numerous genera from the same family, *etc.*) and it is clear that a more
52 general solution is required.

53 Prior research has shown that the best way to account for phylogenetic effects is to explicitly
54 incorporate an evolutionary model into the statistical methods whenever possible [29–36]. However,
55 this strategy can be challenging for certain problems [37] and simpler methods that differentially weight
56 taxa according to their overall similarity to other taxa in a given dataset have been developed and
57 applied for decades [38–46]. In the context of the toy example from Figure 1A, the choice of valine as the
58 preferred amino acid comes from a model that weights each sequence uniformly. By down-weighting
59 highly similar sequences, however, weighted frequencies could be used to come to the conclusion that
60 threonine is instead the preferred amino acid. Instead of looking at preferred amino acid residues
61 (one-site probabilities), evolutionary coupling analyses use sequence alignments to infer co-evolving
62 positions via their two-site marginal probabilities. The current best practice for evolutionary coupling
63 analyses is to down-weight sequences that are highly similar to one-another when inferring parameters
64 from the multiple sequence alignment data. While this strategy appears in numerous methods, a
65 systematic analysis of the benefit that sequence weighting provides in comparison to uniform weights,
66 and an evaluation of different conceptually distinct strategies for assigning weights to sequences has
67 not been performed to our knowledge.

68 Here, we evaluate existing weighting strategies alongside alternative tree- and sequence-based
69 methods that have been proposed and used in various biological applications. We define the
70 accuracy of a given method according to how well the resulting evolutionary couplings are able
71 to predict residue–residue contacts within known representative structures of protein families [18].
72 Despite potential theoretical disadvantages, we find that the current best practice method of 80%
73 sequence-identity-based weighting outperforms alternative methods that explicitly incorporate
74 knowledge of phylogenetic relatedness. We show that a modification of this method provides a
75 slight but insignificant improvement, and more broadly show that several methodologically distinct
76 methods produce accuracies that are nearly indistinguishable both from one-another and from uniform
77 weights.

78 2. Results

79 2.1. An explanation of weighting methods

There are many variants of evolutionary coupling analysis methods that have been developed, and most methods implement a sequence-identity-based correction to mitigate the effect of phylogenetic relatedness [10,11,13]. Specifically, given n sequences in an alignment, the pairwise similarity of all sequences is calculated and the weight $W(i)$ of a given sequence i within an alignment equals

the inverse of the total number of sequences j whose distance $d(i, j)$ to sequence i is less than some parameter λ :

$$W(i) = 1 / \sum_{j=1}^n I(i, j), \quad (1)$$

where n is the number of sequences in the alignment and $I(i, j)$ is an indicator variable defined as

$$I(i, j) = \begin{cases} 0 & \text{if } d_{i,j} < \lambda, \\ 1 & \text{if } d_{i,j} \geq \lambda. \end{cases} \quad (2)$$

80 The distance $d(i, j)$ and the cutoff λ are usually measured as percent sequence identity: the number of
81 identical residues between two aligned sequences divided by their total length.

82 Under this weighting scheme, highly unique sequences are given a weight value of 1, whereas
83 sequences that are similar to others are assigned weights between 0 and 1 according to how many
84 such similar sequences are in the alignment. Given this strategy, the effective number of sequences is
85 simply the sum the weights assigned to all sequences, which takes a value between 0 and n .

86 Several possible issues arise from this weighting scheme. First, it is not immediately apparent
87 what value of λ is most appropriate to use as a sequence identity threshold. While this parameter
88 can be optimized for practical utility (the field has coalesced largely around a value of 80%), it is
89 unclear what this value tells us about the co-evolutionary process or *why* it works so well. Second,
90 this weighting scheme can produce some counter-intuitive results. Given an 80% sequence identity
91 threshold, two otherwise independent sequences in an alignment sharing 99% sequence identity
92 will each be assigned a weight of 0.5 reflecting their relative similarity to one another. In the same
93 alignment, two sequences sharing 81% sequence identity will similarly each be assigned a weight of
94 0.5 despite being much more distinct from one another compared to the former pair. Yet two sequences
95 sharing 79% sequence identity will be assigned a weight of 1.0. Finally, the underlying phylogenetic
96 history of the sequence evolution is ignored by this sequence-based comparison method which may
97 inhibit its overall effectiveness.

98 Our goal here is not to exhaustively evaluate all possible strategies for assigning weights to
99 sequences or tips on a phylogeny but rather to test several popular methods that represent logical
100 starting points for possible improvements for use in evolutionary coupling analyses. Specifically, we
101 decided to implement and test three algorithms: one sequence-based method and two conceptually
102 distinct tree-based methods. The sequence-based method was proposed in Henikoff and Henikoff
103 [44] and proceeds across each position by first awarding each observed residue at given position in
104 an alignment an equal share of the weight for that position (where each position in the alignment
105 has a starting weight of 1). The weights at that position for each sequence in the alignment are then
106 assigned by dividing the weight assigned to each residue equally among all sequences sharing the
107 same residue. Finally, the weight of a given sequence is simply the sum of the weights assigned to
108 each position/residue. The method gives intuitively correct results for toy examples and has been
109 used in numerous popular applications including HMMER and PSI-BLAST, with several different
110 modifications for dealing with gap sequences [47,48].

111 We additionally implemented two tree-based methods that were initially proposed in Altschul
112 *et al.* [38] (hereafter referred to as “ACL” weights) and Gerstein *et al.* [43] (hereafter referred to as
113 “GSC” weights). The ACL method is equivalent to a model of electricity where a power source is
114 plugged into the root of the tree, each branch provides resistance proportional to its length, and the
115 current flowing out of each tip is used to determine the weights [38]. By contrast, the GSC method
116 is a way of partitioning the branch lengths of a tree where the final weight of each tip is a weighted
117 sum of all the branch lengths leading up to it [38,43]. Conceptually, ACL and GSC weights are quite
118 distinct with GSC weights assigning a higher weight to tips that have particularly long branch lengths
119 (and thus occupy a larger proportion of the tree) and ACL weights assigning the highest weights to
120 sequences with particularly short branch lengths that reside closest to the root. We note that both

121 metrics explicitly account for the underlying tree topology and thus require a previously constructed
122 rooted evolutionary tree.

123 A notable caveat to the HH, ACL, and GSC weighting methods is that they do not provide
124 intuitive *absolute* scales. The sum of all HH weights in their original formulation is equivalent to the
125 length of the alignment, ACL weights are relative and sum to 1, and GSC weights are in units of branch
126 length (substitutions per unit time) [38,43,44]. Thus, for each of these three methods we employ two
127 re-scaling strategies: First, we divide each weight value by the mean for that alignment, such that the
128 weights for a given alignment will sum to n , where n is the number of sequences. Second, we divide
129 each weight by the maximum observed weight in an alignment, such that the largest relative weight
130 will be assigned a value of 1 and all other weights are some fraction of this.

131 For an example protein (PDB:1AOE), assigning weights to a sequence alignment / tree
132 demonstrates that the methods vary substantially in how uniformly they distribute weights (Figure 1B).
133 The GINI coefficient is a measurement of uniformity where values of zero correspond to uniform
134 weights and values approaching 1 illustrate the case where a small number of sequences have very
135 large weights while the remainder have very small weights. This relationship can be visualized by a
136 Lorenz curve, which in this case plots the cumulative fraction of weights (y-axis) against the cumulative
137 fraction of sequences (x-axis, sorted from lowest to highest weights). The Lorenz curves in Figure 1B
138 show that ACL weights in particular result in a highly uneven distribution of weights. This finding
139 holds more broadly across a dataset of 150 diverse protein families; the tree-based methods produce a
140 more un-even distribution of weights, with ACL weights being particularly highly skewed (Figure 1C).

141 In general, the different weighting schemes (when applied to the same multiple sequence
142 alignment) are only modestly correlated with one-another. Figure 1D shows the median correlation
143 (across the 150 protein families) observed between HH, GSC, and ACL as well as the commonly used
144 80% sequence-identity-based re-weighting method. In general, the weights produced by different
145 methods on the same protein family are significantly positively correlated with one-another, but the
146 correlations are fairly low, demonstrating that the weighting methods themselves are distinct.

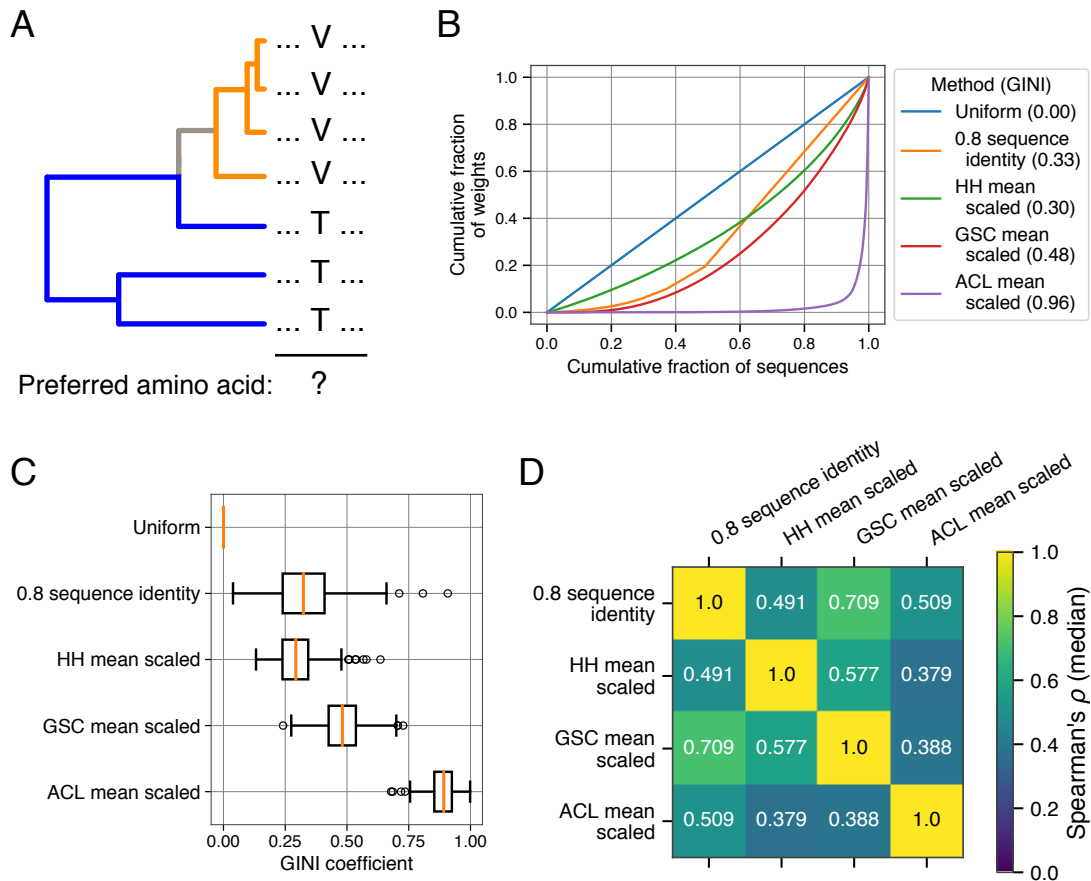


Figure 1. Weighting methods and their relationships in empirical datasets. (a) A toy example illustrating the problem of biased sampling and phylogenetic relatedness. Judging by their frequency (i.e. uniform weighting), valine (V) is the preferred amino acid at the indicated position. However, threonine (T) occupies a substantially larger proportion of the inferred evolutionary history. (b) For an example protein sequence alignment (PDB:1AOE), different weighting strategies produce a more- and less-uniform distribution of weights as visualized by the Lorenz curve. (c) The distribution of GINI coefficients for 150 protein families (higher coefficients correspond to a less uniform distribution of weights) using different weighting strategies (boxes span the 25th through 75th percentiles, red line indicates the median). (d) The median correlation coefficient (Spearman's ρ) of different weighting methods observed across the same 150 protein families.

147 2.2. Sequence weighting does little to improve contact predictions

148 To test the effectiveness of different weighting methods, we calculated evolutionary couplings
 149 using the program CCMPredPy—a Python-based implementation of one of the most popular
 150 pseudo-likelihood based methods (CCMPred), which we modified to accept weights from externally
 151 supplied files—for 150 unique protein families with known structural representatives [13,16]. We next
 152 tested what fraction of the top L couplings for a given protein family (where L is the length of the
 153 reference sequence with a known three-dimensional structure) are true intramolecular residue–residue
 154 contacts—a metric known as the Positive Predictive Value (PPV) (see Materials and Methods for details)
 155 [18]. We separately quantified accuracies from the raw evolutionary couplings, entropy-corrected
 156 couplings, and Average Product Corrected (APC) couplings. The latter two *post-hoc* corrections have
 157 been shown to improve the accuracy of evolutionary couplings by accounting for uneven sequence
 158 entropies across positions in the alignment and perhaps the underlying phylogenetic structure [16,49].

159 As expected, we found that across all weighting schemes, the APC (and to a slightly lesser extent,
 160 the entropy-corrected) evolutionary couplings produce substantially more accurate results compared
 161 to raw coupling scores (Figure 2). In nearly all cases, sequence-identity-based weighting resulted

162 in the highest accuracy. For the best performing APC coupling scores (Figure 2A), the commonly
163 used λ parameter representing an 80% sequence identity threshold resulted in significantly higher
164 accuracies compared to the uniform weight controls (Wilcoxon signed-rank test, $p < 0.001$). One
165 phylogeny-based weighting method (GSC) and the HH sequence-based method were slightly more
166 accurate than uniform weights provided that they were mean-scaled but the improvement was not
167 significant in either case ($p = 0.09$ and $p = 0.1$, respectively); both methods were significantly less
168 accurate than the 80% sequence-identity-based method ($p < 0.001$ for both cases). ACL weights by
169 contrast generally performed poorly in all cases.

170 We note that even in the best case scenario the increase in PPV due to sequence weighting is
171 comparatively small when compared to the large improvements in accuracy that result from the
172 *post-hoc* APC and entropy corrections: median PPV for uniform weights are more than twice as high
173 for APC couplings relative to raw couplings. Interestingly, the best performing weighting schemes
174 substantially improve the accuracy of raw evolutionary couplings relative to the uniform weight
175 control (Figure 2C, 44% median increase in PPV for max-scaled GSC weights, $p < 0.001$), but do
176 comparatively little in the case of the more accurate APC couplings (Figure 2A, 2% median increase in
177 PPV for 80% sequence-identity-based weights, $p < 0.001$).

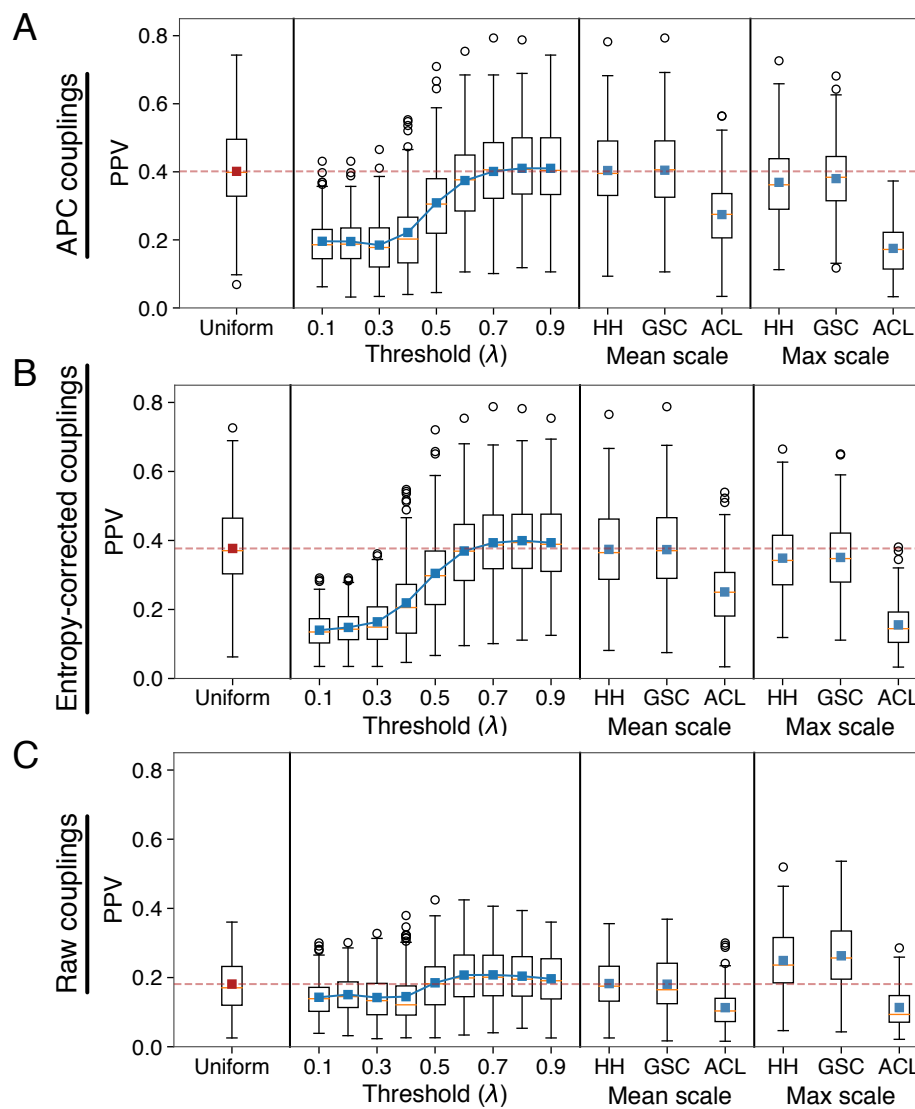


Figure 2. Testing the ability of evolutionary couplings to predict residue–residue contacts in representative structures. “Uniform” refers to the use of uniform weights for all sequences when fitting evolutionary coupling parameters (red dashed line indicates the mean of this distribution and represents a baseline performance that methods should improve upon). “Threshold (λ)” refers to sequence-identity based weighting with different parameters, and “Mean scale”, “Max scale” refer to two different scalings of the indicated weighting methods (HH, GSC, and ACL). (a) Using APC couplings, the mean positive predictive values (PPVs) of the top L couplings vary across different weighting schemes used to infer evolutionary couplings. However, the only methods that significantly improve performance is sequence-identity-based re-weighting with $\lambda=0.8$ or 0.9 (Wilcoxon signed-rank test, $p < 0.001$), but the magnitude of the improvement is modest (1.9% and 1.1% median improvement over uniform). (b) Using entropy-corrected evolutionary coupling values leads to similar conclusions that no weighting scheme substantially outperforms uniform weights. (c) Using raw evolutionary coupling values results in substantially higher accuracies for certain weighting methods relative to uniform, but the overall accuracies remain low compared to (a) and (b).

178 2.3. Weighting on time-scaled trees

179 In Figure 1, we noted that tree-based weighting methods produced a more un-even distribution
 180 of weights compared to the sequence-based weighting methods that we tested. A potential issue with
 181 both of the tree-based weighting methods that we consider here is that the rates of evolution vary
 182 across phylogenetic trees and thus species are not equidistant from the root sequence. Phylogenetic

183 trees reflect both the relationship between species and the rate of evolution along each branch. For
184 trees consisting solely of extant species, numerous methods can re-scale trees to produce tips that
185 are contemporaneous and equidistant from the root (Figure 3A) [50]. Since GSC and ACL weighting
186 methods are significantly influenced by the overall distance from the root for individual tips, we
187 reasoned that computing these weights on scaled-trees may produce less variable weights and perhaps
188 more accurate results. We thus used the RelTime algorithm to transform each raw tree into a time-scaled
189 tree and re-computed the weights for the two tree-based weighting methods on these RelTime trees
190 [50].

191 For a given protein alignment, weights constructed in this manner display significantly less
192 heterogeneity than weights calculated from the raw trees (Wilcoxon signed-rank test, $p < 0.001$). The
193 PPVs of mean- and max-scaled weighting methods were significantly improved in all cases relative to
194 weights computed on the raw trees (Figure 3B, results shown for APC couplings). The improvements
195 were again comparatively small and no method out-performed 80% sequence-identity-based weights.
196 However, PPVs with mean-scaled GSC weights calculated from RelTime trees were significantly
197 higher than PPVs from uniform weighting (Wilcoxon signed-rank test, $p = 0.003$) and the difference in
198 PPV between these weights and the best performing 80% sequence-identity-based weights was not
199 significant ($p = 0.14$).

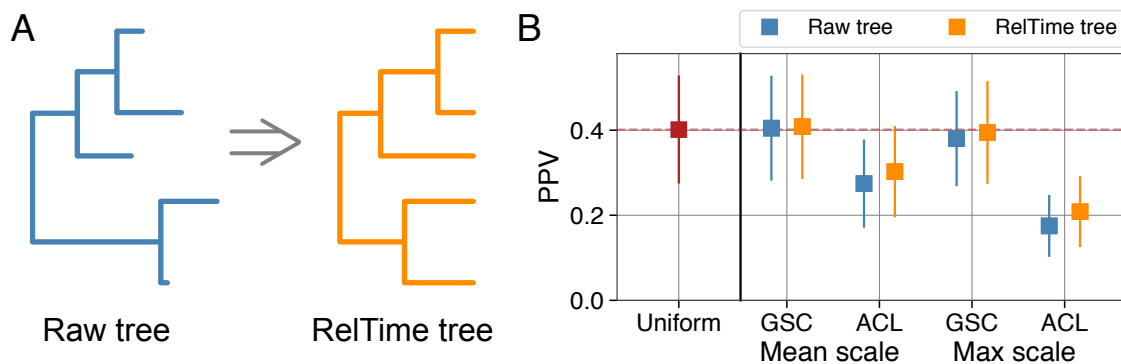


Figure 3. Tree re-scaling prior to calculation of weights slightly improves accuracies. (a) Raw, rooted phylogenetic trees can be converted to time-scaled trees with contemporaneous tips using the RelTime algorithm. (b) Sequence weights calculated from RelTime trees result in slightly better residue-residue contact prediction for the two tree-based weighting methods that we consider (and the two separate scalings of those weights). Shown is the mean PPV for 150 protein families using APC couplings, with error bars showing the standard deviation.

200 2.4. An altered sequence-identity-based method that accounts for sequence similarity.

Thus far we have shown that the current best practice of using sequence-identity-based weighting within a 80% sequence similarity neighborhood results in evolutionary couplings that have the highest power to predict intra-molecular residue-residue contacts. However, we also discussed some potentially counter-intuitive properties of this sequence-identity-based method. We thus developed and tested a variant of the sequence-identity-based method that down-weights sequences according to pairwise similarity and an identity threshold, but does so by accounting for the actual similarity between the sequences. Whereas the original method assigns each sequence a value of 1 and divides by the raw number of similar sequences (defined according to the λ parameter), our modification instead divides by the sum of a similarity-adjusted value for each sequence. Specifically,

$$W(i) = 1 / \sum_{j=1}^n I_{\text{adj}}(i, j). \quad (3)$$

In contrast to Equation (2), $I_{\text{adj}}(i, j)$ produces a continuous range of values between 0 and 1:

$$I_{\text{adj}}(i, j) = \begin{cases} 0 & \text{if } d_{i,j} < \lambda, \\ (d_{i,j} - \lambda)/(1 - \lambda) & \text{if } d_{i,j} \geq \lambda. \end{cases} \quad (4)$$

201 As in Equations (1,2), the distance $d_{i,j}$ and the cutoff λ are measured as percent sequence identity.

202 Using this method with a λ value of 0.8, two otherwise independent sequences in an alignment
203 with 99% sequence identity will each be assigned a weight of 0.513 [or $1/(1 + 0.95)$, where $0.95 =$
204 $(0.99 - 0.8)/(1 - 0.8)$], reflecting their high similarity to one another. In the same alignment, two
205 sequences sharing only 81% sequence identity will by contrast each be assigned only a slightly
206 decreased weight of 0.95 [or $1/(1 + 0.05)$, where $0.05 = (0.81 - 0.8)/(1 - 0.8)$]. All else being equal,
207 the more similar sequences are, the more they will be down-weighted up to the given sequence identity
208 threshold, at which point no further down-weighting occurs.

209 Comparing this similarity-adjusted sequence-identity-based method to the original method
210 shows that the similarity-based adjustment produces more robust results across the range of possible
211 values for λ (Figure 4). Across all of the different variants that we tested, similarity-adjusted
212 sequence-identity-based weights with an identity parameter of 0.8 (and the APC, Figure 4A) produced
213 evolutionary couplings with the highest median and mean PPV for the 150 protein families. PPVs
214 resulting from this method were significantly higher than results from uniform weights (1.9% median
215 and 3.7% mean increase in PPV, Wilcoxon signed-rank test $p < 0.001$) but the increase compared to
216 80% sequence-identity weights calculated in the original manner was slight and not significant (0%
217 median and 0.3% mean increase in PPV, $p = 0.11$).

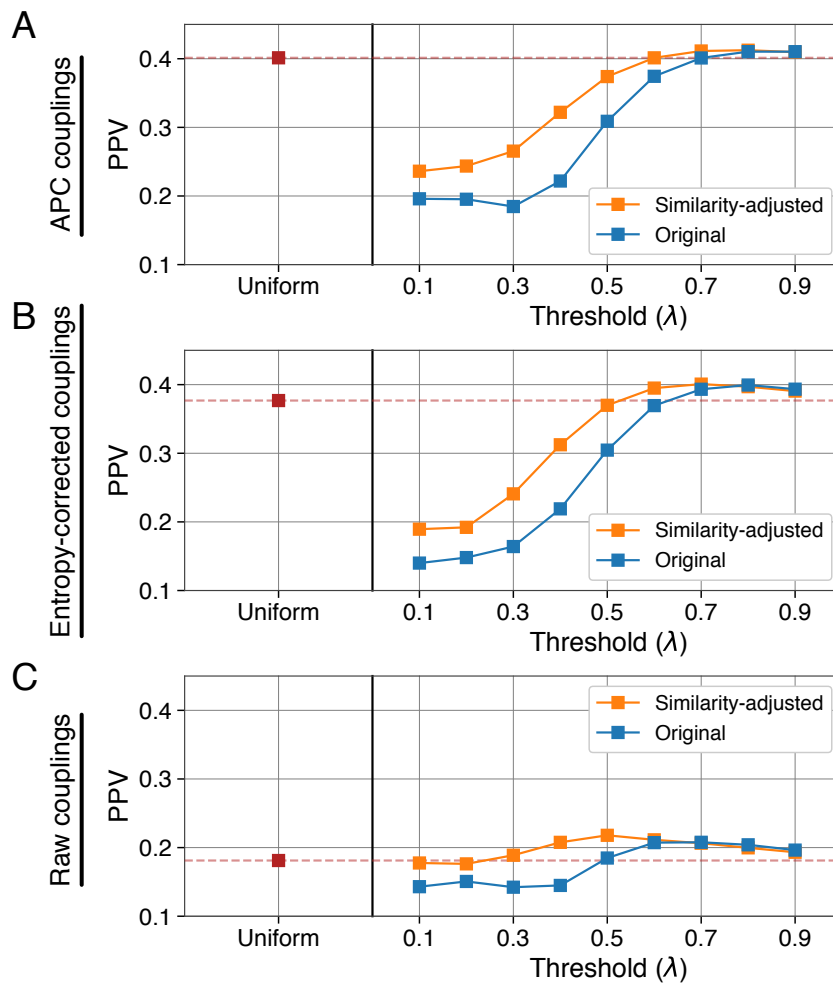


Figure 4. An altered sequence-identity-based method is more robust to parameter choice. (a) Using APC couplings, mean PPVs for similarity-adjusted sequence-identity-based weights are equal-to or higher-than PPVs calculated with the commonly used sequence-identity-based weights. (b) Same as in (a), using entropy-corrected evolutionary coupling values. (c) Same as in (a) and (b), using raw evolutionary coupling values.

218 3. Discussion

219 Natural sequence alignments are not composed of independently evolved lineages and instead
220 have an unknown pattern of relationships that can be inferred and visualized as a phylogenetic tree.
221 Statistical methods that fail to account for these relationships are expected to be biased, but in the case
222 of direct coupling analyses a phylogenetically agnostic model has nevertheless proven valuable at
223 predicting residue-residue contacts within protein structures [5,10,11]. Differential sequence weighting
224 is commonly employed in such analyses as a way to partially mitigate phylogenetic effects, but the
225 overall benefit that such weights provide has yet to be systematically interrogated. We have shown here
226 that numerous (and conceptually distinct) weighting methods produce evolutionary couplings with
227 a roughly equivalent ability to predict residue-residue contacts—given that the coupling values are
228 transformed *post-hoc* via the average product correction (APC). We found that uniform, HH, GSC, and
229 two variants of 80% sequence-identity-based weights all produce nearly indistinguishable accuracies
230 from one another. While we have only evaluated a few different weighting methods and variants, the
231 similar predictive power of top-performing weighting strategies (despite being substantially different
232 from one-another, Figure 1D) suggests that there may be little room for improvement on top of current
233 best practices.

234 Intuitively, uneven sampling and phylogenetic biases are *expected* to introduce spurious effects into
235 statistical models. Indeed, this is known to be the case in numerous contexts, such as when assessing
236 the strength of correlations between discrete and continuous traits [28,34,36]. Nevertheless, we have
237 shown here that using variable sequence weights to correct for these problems provides little (if any)
238 practical benefit when attempting to predict residue–residue contacts. Why might this be the case?
239 We caution that weights alone are an imperfect method of accounting for shared phylogenetic history,
240 and in other contexts achieving accurate true and false positive rates from statistical tests requires
241 more than simple re-weighting of data points [29,31,36,51,52]. In the context of evolutionary couplings,
242 it is unclear whether uneven sampling and phylogenetic biases do not affect the fitting of coupling
243 parameters as much as one might initially think, whether the APC (a *post-hoc* re-scaling procedure)
244 largely corrects for any such factors, or whether weighting in general is simply an inadequate solution
245 to the problem of phylogeny.

246 While we found that numerous weighting methods produce roughly equivalent end results,
247 our findings raise several potential issues that may be worthy of further study moving forward. We
248 noted that many weighting methods do not clearly provide an intuitive absolute scale and instead
249 assign weights to sequences (or tips in a phylogenetic tree) that are either relative or in irrelevant units.
250 This can be problematic from a practical standpoint because most methods for inferring evolutionary
251 coupling parameters between residue–residue pairs rely on some form of prior and the weight given to
252 observed data relative to this prior may affect results. For the HH, GSC, and ACL methods we found
253 that two different scaling procedures (which maintain relative weights within a dataset but change
254 their absolute values) produced varying accuracies (Figure 2). With the exception of star phylogenies,
255 the effective sample size from phylogenetically structured data is strictly less than the number of
256 sequences/data points analyzed. More accurately estimating the effective sample size and scaling
257 weights accordingly may improve the performance of different weighting schemes beyond what we
258 observed here.

259 Additionally, the HH, GSC, and ACL methods do not include a free parameter that can be tuned
260 to improve results. We validated that an 80% sequence identity neighborhood is optimal using the
261 currently accepted method and a similarity-adjusted variant, but this 80% value is a free parameter
262 that has been optimized to produce the highest accuracy for sequence-identity-based weighting.
263 What we believe the optimality of this parameter represents in practice is that once two sequences
264 diverge past approximately 80% similarity, their evolution is effectively independent. If this is the
265 case, down-weighting sequences that for instance share 50% sequence identity would make little sense
266 (and indeed, doing so produces less accurate results). By contrast, the HH, GSC, and ACL methods
267 all inherently compare each sequence to every other sequence in a global manner. It seems possible
268 that some phylogenetic tree transformation may be able to introduce the same intuition of ignoring
269 evolutionary relatedness past some threshold level into tree-based weighting methods [30,32]. The best
270 way to perform such re-scaling, or how to perform something conceptually similar for HH weights, is
271 a promising area for future research.

272 Despite being weakly correlated with one another, uniform, 80% sequence identity, HH, and GSC
273 weights perform roughly equivalently at predicting residue–residue contacts. We recommend that any
274 method with improved performance should become the standard (provided it does not substantially
275 increase computational run-time), and found that a slightly modified sequence-identity-based
276 re-weighting method that accounts for sequence similarity actually performs the best of any method
277 that we tested. However, using either the original or similarity-adjusted sequence-identity-based
278 weighting can be expected to offer less than a few percent improvement in accuracy compared
279 to uniform weights which completely ignore phylogeny. We therefore speculate that substantial
280 improvements to evolutionary coupling analyses will require the explicit incorporation of phylogenies
281 and time-dependent sequence evolution, but how to do so remains elusive.

282 4. Materials and Methods

283 4.1. Description of the dataset.

284 For all of our analyses, we used the so-called “psicov” dataset—an existing set of 150 distinct
285 protein structures with corresponding multiple sequence alignments that have been used in numerous
286 benchmark studies for predicting residue–residue contacts from evolutionary couplings [14,53,54].
287 All sequence and structure data were taken directly from Jones and Kandathil [54], but given the
288 large number of different analyses that we ran, we first randomly down-sampled each alignment to a
289 maximum of 1001 sequences (1000 sequences plus the mandated inclusion of the reference protein
290 sequence).

291 4.2. Phylogenetic tree construction.

292 For each sequence alignment in our dataset, we constructed a rough phylogenetic tree using the
293 double precision version of FastTree2 (v.2.1.10; LG model, gamma distributed rate variation, pseudo
294 flag) [55]. We next adjusted the branch lengths on each guide tree by running the alignment and the
295 template tree through the more accurate IQtree software (v1.6.9; LG model, Gamma-distributed rate
296 variation with 20 categories) [56]. Finally, we rooted the resulting trees using the mid-point method
297 [57].

298 For RelTime trees, we implemented our own version of the RelTime algorithm as described in the
299 original manuscript while ensuring that our method produced similar results [50]. We note here only
300 that our implementation does not perform a statistical test (and subsequent alteration of rates) at the
301 end of the algorithm to ensure that rate changes are significant.

302 4.3. Weighting methods.

303 We developed all of our weighting methods from scratch using custom python programs that
304 heavily leveraged tools from the Biopython package [57]. For sequence identity weighting and the
305 novel similarity-adjusted version we propose here, details are presented in the main text, Equations
306 (1–4). We ensured that our own version of sequence-identity-based weighting was equivalent to the
307 method implemented within CCMpredPy by comparing the resulting effective number of sequences
308 metrics and accuracies and finding them to be identical.

309 For HH based weights, we followed the procedure outlined in the initial paper and ensured that
310 our implementation gave the desired results on the toy examples presented therein [44]. Researchers
311 have pointed out subsequent modifications to this method [47,48] concerning how to effectively
312 treat gap sequences. Rather than treating these as a 21st character as some implementations have
313 done, our implementation assigns gap sequences a weight value of zero. Further, each column in
314 the alignment is weighted from 0 to 1 according to the fraction of non-gapped positions. In this
315 manner, alignment positions with more gaps are assigned lower weights and the positions with
316 gaps themselves contribute a weight of zero. Summation and calculation of final weights follows
317 the published procedure [44]. However, since the units and absolute value of these weights are
318 not intuitive, we finally re-scaled the weights via separate mean- and max-scaling procedures. In
319 mean-scaling, we calculate the mean of all weights determined via the HH algorithm for a particular
320 sequence alignment and then divide the weight of each sequence in the alignment by this value. This
321 ensures that the sum of all final weights will be equal to the number of sequences in the alignment
322 (n). In the separate max-scaling procedure, we find the maximum weight observed for a particular
323 sequence alignment, and subsequently divide all weights in the alignment by this value. The sum
324 of all weights following this procedure is guaranteed to be some value less than the total number of
325 sequences (n).

326 For ACL and GSC weights, we again followed the procedures outlined in the respective
327 manuscripts [38,43] and ensured that our implementations produced identical results to the examples

328 presented therein. As with HH, calculation of final weights occurred by (separately) scaling the weight
329 values via their mean and maximum values as noted above.

330 4.4. Evolutionary coupling analysis.

331 We chose to use CCMpredPy (v1.0.0, contained as part of the CCMgen package) [13,16] for
332 all evolutionary coupling analyses since we were able to modify the source code for this popular
333 method to accept externally supplied weights in the form of a simple text file where the weight value
334 for each sequence corresponded to its line in the input sequence file. We used the default values
335 with the `ofn-pll` flag corresponding to the pseudo-likelihood optimization of coupling parameters.
336 For each different weighting method that we tested, we outputted files corresponding to the raw,
337 entropy-corrected, and average product corrected coupling matrices.

338 4.5. Structural analysis and accuracy determination.

339 We used the .PDB files provided as part of the psicov dataset and for each structure computed a
340 matrix of residue–residue distances. Each distance value is measured according to the geometric center
341 for all side-chain heavy atoms for a particular residue (including the $C\beta$ atom, excluding the $C\alpha$ atom)
342 [18]. In the case of glutamine, the side-chain center coordinates were assigned to the $C\alpha$ atom. We
343 determined residue–residue contacts according to a uniform 7.5 angstrom threshold for all proteins.

344 We determined the accuracy of evolutionary couplings by determining how well they were able
345 to predict residue–residue contacts within a reference structure. We first selected the top L -ranked
346 couplings for each dataset, where L corresponds to the length of the reference protein sequence (i.e.
347 the sequence for which we have a known structure). The PPV for a particular dataset corresponds to
348 the fraction of those top L -ranked couplings that are classified as residue–residue contacts according to
349 the above definition.

350 **Author Contributions:** Conceptualization, A.J.H. and C.O.W.; methodology, A.J.H. and C.O.W.; software, A.J.H.;
351 validation, A.J.H. and C.O.W.; formal analysis, A.J.H.; investigation, A.J.H.; resources, A.J.H. and C.O.W.; data
352 curation, A.J.H.; writing–original draft preparation, A.J.H.; writing–review and editing, A.J.H. and C.O.W.;
353 visualization, A.J.H. and C.O.W.; supervision, C.O.W.; project administration, C.O.W.; funding acquisition, A.J.H.
354 and C.O.W.

355 **Funding:** This work was funded by National Institutes of Health grant R01 GM088344 and by F32 GM130113.

356 **Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the
357 study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to
358 publish the results.

359 Abbreviations

360 The following abbreviations are used in this manuscript:

361 HH weights derived via the method of Henikoff and Henikoff [44]
GSC weights derived via the method of Gerstein *et al.* [43]
362 ACL weights derived via the method of Altschul *et al.* [38]
APC Average Product Correction/ed
PPV Positive Predictive Value

363 References

- 364 1. Gobel, U.; Sander, C.; Schneider, R.; Valencia, A. Correlated Mutations and Residue Contacts in Proteins.
365 *Proteins* **1994**, *18*, 309–317.
- 366 2. Hopf, T.A.; Scharfe, C.P.I.; Rodrigues, J.P.G.L.M.; Green, A.G.; Kohlbacher, O.; Sander, C.; Bonvin, A.M.J.J.;
367 Marks, D.S. Sequence co-evolution gives 3D contacts and structures of protein complexes. *eLife* **2014**,
368 *3*, 1–45. doi:10.7554/eLife.03430.
- 369 3. Hopf, T.A.; Ingraham, J.B.; Poelwijk, F.J.; Schärfe, C.P.; Springer, M.; Sander, C.; Marks, D.S. Mutation
370 effects predicted from sequence co-variation. *Nature Biotechnology* **2017**, *35*, 128–135. doi:10.1038/nbt.3769.

- 371 4. Kamisetty, H.; Ovchinnikov, S.; Baker, D. Assessing the utility of coevolution-based residue-residue contact
372 predictions in a sequence- and structure-rich era. *Proceedings of the National Academy of Sciences* **2013**,
373 *110*, 15674–15679. doi:10.1073/pnas.1314045110.
- 374 5. Ovchinnikov, S.; Park, H.; Varghese, N.; Huang, P.S.; Pavlopoulos, G.A.; Kim, D.E.; Kamisetty, H.;
375 Kyrpides, N.C.; Baker, D. Protein structure determination using metagenome sequence data. *Science* **2017**,
376 *355*, 294–298. doi:10.1126/science.aah4043.
- 377 6. Lapedes, A.S.; Giraud, B.G.; Liu, L.; Stormo, G.D. Correlated mutations in models of protein sequences:
378 phylogenetic and structural effects. *Statistics in molecular biology and genetics* **1999**, *33*, 236–256.
379 doi:10.1214/lnms/1215455556.
- 380 7. Burger, L.; Van Nimwegen, E. Accurate prediction of protein-protein interactions from sequence alignments
381 using a Bayesian method. *Molecular Systems Biology* **2008**, *4*. doi:10.1038/msb4100203.
- 382 8. Weigt, M.; White, R.A.; Szurmant, H.; Hoch, J.A.; Hwa, T. Identification of direct residue contacts in
383 protein-protein interaction by message passing. *Proceedings of the National Academy of Sciences* **2009**,
384 *106*, 67–72. doi:10.1073/pnas.0805923106.
- 385 9. Burger, L.; Van Nimwegen, E. Disentangling direct from indirect co-evolution of residues in protein
386 alignments. *PLoS Computational Biology* **2010**, *6*. doi:10.1371/journal.pcbi.1000633.
- 387 10. Marks, D.S.; Colwell, L.J.; Sheridan, R.; Hopf, T.A.; Pagnani, A.; Zecchina, R.; Sander, C. Protein 3D structure
388 computed from evolutionary sequence variation. *PLoS ONE* **2011**, *6*. doi:10.1371/journal.pone.0028766.
- 389 11. Morcos, F.; Pagnani, A.; Lunt, B.; Bertolino, A.; Marks, D.S.; Sander, C.; Zecchina, R.; Onuchic,
390 J.N.; Hwa, T.; Weigt, M. Direct-coupling analysis of residue coevolution captures native contacts
391 across many protein families. *Proceedings of the National Academy of Sciences* **2011**, *108*, E1293–E1301.
392 doi:10.1073/pnas.1111471108.
- 393 12. Ekeberg, M.; Lövkvist, C.; Lan, Y.; Weigt, M.; Aurell, E. Improved contact prediction in proteins: Using
394 pseudolikelihoods to infer Potts models. *Physical Review E* **2013**, *87*, 1–16. doi:10.1103/PhysRevE.87.012707.
- 395 13. Seemayer, S.; Gruber, M.; Söding, J. CCMpred - Fast and precise prediction of protein residue-residue
396 contacts from correlated mutations. *Bioinformatics* **2014**, *30*, 3128–3130. doi:10.1093/bioinformatics/btu500.
- 397 14. Jones, D.T.; Singh, T.; Kosciölek, T.; Tetchner, S. MetaPSICOV: Combining coevolution methods for accurate
398 prediction of contacts and long range hydrogen bonding in proteins. *Bioinformatics* **2015**, *31*, 999–1006.
399 doi:10.1093/bioinformatics/btu791.
- 400 15. Figliuzzi, M.; Barrat-Charlaix, P.; Weigt, M. How pairwise coevolutionary models capture the
401 collective residue variability in proteins? *Molecular Biology and Evolution* **2018**, *35*, 1018–1027.
402 doi:10.1093/molbev/msy007.
- 403 16. Vorberg, S.; Seemayer, S.; Söding, J. Synthetic protein alignments by CCMgen quantify noise in
404 residue-residue contact prediction. *PLoS Computational Biology* **2018**, *14*, e1006526.
- 405 17. Hopf, T.A.; Green, A.G.; Schubert, B.; Mersmann, S.; Schärfe, C.P.; Ingraham, J.B.; Toth-Petroczy, A.; Brock,
406 K.; Riesselman, A.J.; Palmedo, P.; Kang, C.; Sheridan, R.; Draizen, E.J.; Dallago, C.; Sander, C.; Marks, D.S.
407 The EVcouplings Python framework for coevolutionary sequence analysis Thomas. *Bioinformatics* **2018**, p.
408 bty862.
- 409 18. Hockenberry, A.J.; Wilke, C.O. Evolutionary couplings detect side-chain interactions. *PeerJ* **2019**, e7280, 1–22.
410 doi:10.7717/peerj.7280.
- 411 19. Morcos, F.; Jana, B.; Hwa, T.; Onuchic, J.N. Coevolutionary signals across protein lineages help capture
412 multiple protein conformations. *Proceedings of the National Academy of Sciences* **2013**, *110*, 20533–20538.
413 doi:10.1073/pnas.1315625110.
- 414 20. Bitbol, A.F.; Dwyer, R.S.; Colwell, L.J.; Wingreen, N.S. Inferring interaction partners from protein sequences.
415 *Proceedings of the National Academy of Sciences* **2016**, *113*, 12180–12185. doi:10.1073/pnas.1606762113.
- 416 21. Uguzzoni, G.; John Lovis, S.; Oteri, F.; Schug, A.; Szurmant, H.; Weigt, M. Large-scale identification of
417 coevolution signals across homo-oligomeric protein interfaces by direct coupling analysis. *Proceedings of*
418 *the National Academy of Sciences* **2017**, *114*, E2662–E2671. doi:10.1073/pnas.1615068114.
- 419 22. Cong, Q.; Anishchenko, I.; Ovchinnikov, S.; Baker, D. Protein interaction networks revealed by proteome
420 coevolution. *Science* **2019**, *365*, 185–189. doi:10.1126/science.aaw6718.
- 421 23. Bonnet, X.; Shine, R.; Lourdais, O. Taxonomic chauvinism. *Trends in Ecology & Evolution* **2002**, *17*, 1–3.

- 422 24. Chen, C.; Natale, D.A.; Finn, R.D.; Huang, H.; Zhang, J.; Wu, C.H.; Mazumder, R. Representative Proteomes:
423 A Stable, Scalable and Unbiased proteome set for sequence analysis and functional annotation. *PLoS ONE*
424 **2011**, *6*, e18910. doi:10.1371/journal.pone.0018910.
- 425 25. Rinke, C.; Schwientek, P.; Sczyrba, A.; Ivanova, N.N.; Anderson, I.J.; Cheng, J.F.; Darling, A.; Malfatti,
426 S.; Swan, B.K.; Gies, E.A.; Dodsworth, J.A.; Hedlund, B.P.; Tsiamis, G.; Sievert, S.M.; Liu, W.T.;
427 Eisen, J.A.; Hallam, S.J.; Kyrpides, N.C.; Stepanauskas, R.; Rubin, E.M.; Hugenholtz, P.; Woyke, T.
428 Insights into the phylogeny and coding potential of microbial dark matter. *Nature* **2013**, *499*, 431–437.
429 doi:10.1038/nature12352.
- 430 26. Troudet, J.; Grandcolas, P.; Blin, A.; Vignes-Lebbe, R.; Legendre, F. Taxonomic bias in biodiversity data and
431 societal preferences. *Scientific Reports* **2017**, *7*, 1–14. doi:10.1038/s41598-017-09084-6.
- 432 27. Titley, M.A.; Snaddon, J.L.; Turner, E.C. Scientific research on animal biodiversity is systematically biased
433 towards vertebrates and temperate regions. *PLoS ONE* **2017**, *12*, 1–14. doi:10.1371/journal.pone.0189577.
- 434 28. Felsenstein, J. Phylogenies and the comparative method. *The American Naturalist* **1985**, *125*, 1–15.
- 435 29. Grafen, A. The phylogenetic regression. *Philosophical Transactions of the Royal Society B* **1989**, *326*, 119–157.
- 436 30. Pagel, M. Inferring historical patterns of biological evolution. *Nature* **1999**, *401*, 877–884.
- 437 31. Rohlf, F.J. Comparative methods for the analysis of continuous variables: geometric interpretations.
438 *Evolution* **2001**, *55*, 2143–2160. doi:10.1111/j.0014-3820.2001.tb00731.x.
- 439 32. Blomberg, S.P.; Garland Jr., T.; Ives, A.R. Testing for phylogenetic signal in comparative data: behavioral
440 traits are more labile. *Evolution* **2003**, *57*, 717–745. doi:10.1111/j.0014-3820.2003.tb00285.x.
- 441 33. Ives, A.R.; Midford, P.E.; Garland Jr., T. Within-species variation and measurement error in phylogenetic
442 comparative methods. *Systematic Biology* **2007**, *56*, 252–270. doi:10.1080/10635150701313830.
- 443 34. Ives, A.R.; Garland Jr., T. Phylogenetic Regression for Binary Dependent Variables. *Systematic Biology* **2010**,
444 *59*, 9–26. doi:10.1007/978-3-662-43550-2_9.
- 445 35. Revell, L.J. Size-correction and principal components for interspecific comparative studies. *Evolution* **2009**,
446 *63*, 3258–3268. doi:10.1111/j.1558-5646.2009.00804.x.
- 447 36. Revell, L.J. Phylogenetic signal and linear regression on species data. *Methods in Ecology and Evolution*
448 **2010**, *1*, 319–329. doi:10.1111/j.2041-210x.2010.00044.x.
- 449 37. Uyeda, J.C.; Zenil-Ferguson, R.; Pennell, M.W. Rethinking phylogenetic comparative methods. *Systematic*
450 *Biology* **2018**, *67*, 1091–1109. doi:10.1093/sysbio/syy031.
- 451 38. Altschul, S.F.; Carroll, R.J.; Lipman, D.J. Weights for data related by a tree. *Journal of Molecular Biology* **1989**,
452 *207*, 647–653. doi:10.1016/0022-2836(89)90234-9.
- 453 39. Vingron, M.; Argos, P. A fast and multiple sequence alignment algorithm. *Bioinformatics* **1989**, *5*, 115–121.
- 454 40. Sibbald, P.R.; Argos, P. Weighting aligned protein or nucleic acid sequences to correct for unequal
455 representation. *Journal of Molecular Biology* **1990**, *216*, 813–818. doi:10.1016/S0022-2836(99)80003-5.
- 456 41. Vingron, M.; Sibbald, P.R. Weighting in sequence space: A comparison of methods in terms of generalized
457 sequences. *Proceedings of the National Academy of Sciences* **1993**, *90*, 8777–8781. doi:10.1073/pnas.90.19.8777.
- 458 42. Thompson, J.D.; Higgins, D.G.; Gibson, T.J. Improved sensitivity of profile searches through the use of
459 sequence weights and gap excision. *Bioinformatics* **1994**, *10*, 19–29.
- 460 43. Gerstein, M.; Sonnhammer, E.L.; Chothia, C. Volume changes in protein evolution. *Journal of Molecular*
461 *Biology* **1994**, *236*, 1067–1078. doi:10.1016/0022-2836(94)90012-4.
- 462 44. Henikoff, S.; Henikoff, J.G. Position-based sequence weights. *Journal of Molecular Biology* **1994**, *243*, 574–578.
463 doi:10.1016/0022-2836(94)90032-9.
- 464 45. Krogh, A.; Mitchison, G. Maximum entropy weighting of aligned sequences of proteins or DNA. *Proceedings*
465 *of the International Conference on Intelligent Systems for Molecular Biology* **1995**, *3*, 215–21.
- 466 46. Stone, E.A.; Sidow, A. Constructing a meaningful evolutionary average at the phylogenetic center of mass.
467 *BMC Bioinformatics* **2007**, *8*, 1–13. doi:10.1186/1471-2105-8-222.
- 468 47. Altschul, S.F.; Madden, T.L.; Schäffer, A.A.; Zhang, J.; Zhang, Z.; Miller, W.; Lipman, D.J. Gapped BLAST
469 and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research* **1997**,
470 *25*, 3389–3402. doi:10.1093/nar/25.17.3389.
- 471 48. Eddy, S.R. Profile hidden Markov models. *Bioinformatics* **1998**, *14*, 755–763.
472 doi:10.1093/bioinformatics/14.9.755.

- 473 49. Dunn, S.D.; Wahl, L.M.; Gloor, G.B. Mutual information without the influence of phylogeny
474 or entropy dramatically improves residue contact prediction. *Bioinformatics* **2008**, *24*, 333–340.
475 doi:10.1093/bioinformatics/btm604.
- 476 50. Tamura, K.; Battistuzzi, F.U.; Billing-Ross, P.; Murillo, O.; Filipinski, A.; Kumar, S. Estimating divergence
477 times in large molecular phylogenies. *Proceedings of the National Academy of Sciences* **2012**, *109*, 19333–19338.
478 doi:10.1073/pnas.1213199109.
- 479 51. Bruno, W.J. Modeling residue usage in aligned protein sequences via maximum likelihood. *Molecular*
480 *Biology and Evolution* **1996**, *13*, 1368–1374. doi:10.1093/oxfordjournals.molbev.a025583.
- 481 52. Newberg, L.A.; McCue, L.A.; Lawrence, C.E. The Relative Inefficiency of Sequence Weights Approaches in
482 Determining a Nucleotide Position Weight Matrix. *Statistical Applications in Genetics and Molecular Biology*
483 **2005**, *4*. doi:10.2202/1544-6115.1135.
- 484 53. Jones, D.T.; Buchan, D.W.; Cozzetto, D.; Pontil, M. PSICOV: Precise structural contact prediction
485 using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics* **2012**,
486 *28*, 184–190. doi:10.1093/bioinformatics/btr638.
- 487 54. Jones, D.T.; Kandathil, S.M. High precision in protein contact prediction using fully
488 convolutional neural networks and minimal sequence features. *Bioinformatics* **2018**, *34*, 3308–3315.
489 doi:10.1093/bioinformatics/bty341.
- 490 55. Price, M.N.; Dehal, P.S.; Arkin, A.P. FastTree 2 - Approximately maximum-likelihood trees for large
491 alignments. *PLoS ONE* **2010**, *5*. doi:10.1371/journal.pone.0009490.
- 492 56. Nguyen, L.T.; Schmidt, H.A.; Von Haeseler, A.; Minh, B.Q. IQ-TREE: A fast and effective stochastic
493 algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution* **2015**,
494 *32*, 268–274. doi:10.1093/molbev/msu300.
- 495 57. Cock, P.J.; Antao, T.; Chang, J.T.; Chapman, B.A.; Cox, C.J.; Dalke, A.; Friedberg, I.; Hamelryck, T.; Kauff,
496 F.; Wilczynski, B.; De Hoon, M.J. Biopython: Freely available Python tools for computational molecular
497 biology and bioinformatics. *Bioinformatics* **2009**, *25*, 1422–1423. doi:10.1093/bioinformatics/btp163.

498 **Sample Availability:** Raw and processed data used in this manuscript have been deposited at: [10.5281/zenodo.](https://zenodo.org/record/3368652)
499 [3368652](https://zenodo.org/record/3368652). All necessary code to replicate the analyses presented here have been deposited at: [https://github.com/](https://github.com/adamhockenberry/dca-weighting)
500 [adamhockenberry/dca-weighting](https://github.com/adamhockenberry/dca-weighting).

501 © 2019 by the authors. Submitted to *Journal Not Specified* for possible open access
502 publication under the terms and conditions of the Creative Commons Attribution (CC BY) license
503 (<http://creativecommons.org/licenses/by/4.0/>).