1    TITLE
2
3    EvoWeaver: Large-scale prediction of gene functional associations from coevolutionary signals
4
5    AUTHORS
6
7    Aidan Lakshman[1] and Erik S. Wright[1,2,*]
8
9    [1]Department of Biomedical Informatics, University of Pittsburgh
10   [2]Center for Evolutionary Biology and Medicine, Pittsburgh, PA
11   *address correspondence to eswright@pitt.edu
12

13    ABSTRACT
14
15    The known universe of uncharacterized proteins is expanding far faster than our ability to annotate their functions
16    through laboratory study. Computational annotation approaches rely on similarity to previously studied proteins,
17    thereby ignoring unstudied proteins. Coevolutionary approaches hold promise for injecting new information into our
18    knowledge of the protein universe by linking proteins through 'guilt-by-association'. However, existing coevolutionary
19    algorithms have insufficient accuracy and scalability to connect the entire universe of proteins. We present
20    EvoWeaver, an algorithm that weaves together 12 signals of coevolution to quantify the degree of shared evolution
21    between genes. EvoWeaver accurately identifies proteins involved in protein complexes or separate steps of a
22    biochemical pathway. We show the merits of EvoWeaver by partly reconstructing known biochemical pathways
23    without any prior knowledge other than that available from genomic sequences. Applying EvoWeaver to 1,545 gene
24    groups from 8,564 genomes reveals missing connections in popular databases and potentially undiscovered links
25    between proteins.

INTRODUCTION

Our ability to capture the protein universe with genome sequencing far outpaces our ability to investigate individual proteins. A select few proteins have historically received a disproportionate amount of study[1-3]. This annotation inequality hinders biomedical progress by neglecting many proteins that could be important determinants of health[4]. Only a small fraction of uncharacterized proteins can be automatically annotated via similarity to experimentally investigated proteins of known function[5-7]. The sparsity of high-quality annotations exacerbates the problem of non-specific and low-confidence annotations that proliferate across genomes[8,9]. Thus, computational approaches to infer function without dependence on prior knowledge are acutely needed.

Computationally annotating the remainder of the protein universe requires establishing connections with characterized proteins to generate hypotheses about function through 'guilt by association'[10]. Shared function necessitates that protein-encoding genes coevolve in the same cell, thereby leaving behind a molecular signal of coevolution[11]. Four primary approaches are used to identify coevolution: phylogenetic profiling[12-14], phylogenetic structure[15-17], gene organization[18-20], and sequence level methods[21-23]. Each of these coevolutionary signals is an outcome of a shared selection pressure acting on groups of genes. To date, these four coevolutionary approaches have primarily been applied independently. Even large databases of functional associations, such as STRING, only consider evidence from a small subset of coevolutionary approaches[24].

Although coevolutionary analyses have shown great potential for predicting functional associations[25-32], scalability is a major impediment to comprehensive application on large datasets. The era of big data holds the promise of distinguishing coevolution from other drivers of molecular evolution[13]. Additionally, holistic evaluation of many coevolutionary signals offers a means of amplifying weaker signals to make higher accuracy predictions. For example, conserved genes may not display a phylogenetic profiling signal but can still show patterns of gene organization. Combining disparate coevolutionary signals and scaling to larger datasets requires inventing new approaches for discerning signal from noise.

Coevolutionary analyses have the potential to infer functional associations directly from sequencing data in a way that is agonistic to prior annotations, thereby overcoming the current reliance on extrapolating from existing knowledge that compounds annotation inequality. Here, we set out to develop a scalable approach to extract and combine coevolutionary signals for predicting functional associations between protein-coding genes. This required improving upon existing approaches to scale to larger input data and incorporate statistical testing. We unite these signals of coevolution using machine learning models to quantify the degree of functional association between genes. Our approach, named EvoWeaver, is available within the *SynExtend* package (v1.17.1) for R and serves as a high-quality hypothesis generator to help extend our knowledge of the protein universe.

RESULTS

Existing coevolutionary algorithms have widespread issues with scalability, interoperability, and interpretability[13]. We chose to implement all our coevolutionary analyses from scratch within a single software package to standardize user interaction and allow for easy application of ensemble methods. Our approach, named EvoWeaver, takes as input a set of phylogenetic gene trees and optional metadata (Fig. 1a). EvoWeaver then performs four types of coevolutionary analysis, comprised of 12 algorithms optimized for scalable performance. The output of EvoWeaver is 12 scores ranging from -1 to 1 that quantify the strength of coevolution between a pair of gene groups. These scores are combined using a machine learning classifier to generate novel inferences or hypotheses about gene function.

The first type of coevolutionary analysis, Phylogenetic Profiling, investigates patterns of presence/absence (P/A) or gain/loss (G/L) of genes, which manifest when multiple genes work in concert (Fig. 1b). While P/A analyses have been successfully used to predict gene function[12-14,33-35], existing approaches can be susceptible to biases from small sample sizes or low evolutionary divergence[36]. We addressed these biases with a novel algorithm (G/L Distance) that examines the distance between G/L events to measure compensatory changes rather than extant patterns. We also implement clade-wise phylogenetic profiling (P/A Jaccard), which corrects for bias from oversampled taxa[37]. Finally, we analyze the mutual information of ancestral state transitions (G/L MI), as well as the conservation of mutual presence in ancestral states (P/A Overlap). The end result is a category of algorithms for identifying coevolution between gene groups that are not highly conserved.

The second type of coevolutionary analysis, Phylogenetic Structure, uses the fact that functionally associated genes tend to evolve in tandem[38], giving rise to similar genealogies (Fig. 1c). Commonly used phylogenetic structure approaches include MirrorTree and ContextTree[39-41], although these approaches scale poorly due to high computational complexity. We addressed this issue by introducing novel algorithms (RP MirrorTree, RP ContextTree) that use random projection to decrease computational overhead and improve accuracy by reducing redundant information. Random projection provides the added advantage that computation can be distributed across computers, unlike in SVD-phy[42], allowing EvoWeaver to process very large datasets on compute clusters. Additionally, we introduce the use of tree distance metrics (Tree Distance) to analyze coevolution via topological differences in

85   genealogies[43]. Taken together, these algorithms facilitate inference of coevolution among more conserved gene
86   groups.
87      The third type of coevolutionary analysis, Gene Organization, leverages the fact that functionally linked genes
88   tend to colocate on the genome to facilitate gene regulation and horizontal gene transfer[44-46] (Fig. 1d). These
89   approaches most commonly employ profile hidden Markov models, such as antiSMASH[47-49]. While these methods
90   perform well on functional prediction, they rely on *a priori* knowledge about genes that colocalize. We circumvented
91   this limitation by introducing an algorithm that compares the number of coding regions separating genes (Gene
92   Distance). Our approach is similar to STRING's colocalization metric, which measures the number of nucleotides
93   separating genes[24], but STRING's approach fails to consider that low rates of evolutionary divergence can inflate
94   evidence of colocalization. We address this issue by using Moran's I to calculate the extent to which evolutionary
95   divergence affects the observed colocalization of genes. Additionally, EvoWeaver analyzes the conservation of
96   relative gene orientation (Orientation MI), since this also indicates functional association[50]. Collectively, these
97   algorithms provide evidence of coevolution among conserved gene groups on the same chromosome.
98      The last type of coevolutionary analysis, Sequence Level methods, looks at sequence patterns across gene
99   groups, which are sometimes indicative of physical interactions between gene products[51] (Fig. 1e). Direct coupling
100  analysis is a well-known approach in this category[52-54], but it suffers from high computational complexity. Instead, we
101  extended a prior approach based on mutual information to predict interacting sites between sequences[55]. EvoWeaver
102  analyzes the extent of these site-wise interactions to construct an overall score (Sequence Info). Additionally,
103  EvoWeaver compares gene sequence natural vectors (Gene Vector), which carry evidence of functional association
104  and can be quickly computed[56]. These algorithms provide additional evidence of coevolution for physically interacting
105  gene products.
106     The four categories of analysis span levels of coevolution from the organism (Phylogenetic Profiling) to the
107  genome (Gene Organization) to the gene (Phylogenetic Structure) to the sequence. Since our component analyses
108  individually capture different facets of coevolution, we sought to combine their strengths into a single comprehensive
109  estimate of evidence for functional association between gene pairs. To this end, we trained three machine learning
110  classifiers (logistic regression, random forest, and neural network) on sets of protein-coding gene pairs with known
111  functional associations (Fig. 1a). While these ensemble models require *a priori* knowledge to calibrate their
112  predictions, after training they permit the extension of this knowledge to gene pairs without previously known
113  associations. More details about each algorithm are provided in section **SI1** of the Supplemental Information.
114
115  **Ensemble methods accurately identify functionally associated genes**
116
117     Selection of high-quality ground truth datasets for coevolutionary analysis is a challenging task[13]. As with previous
118  studies[42,57], we relied upon the Kyoto Encyclopedia of Genes and Genomes database (KEGG) because it is well-
119  curated and experimentally validated[58,59]. KEGG provides a hierarchical ontology of biochemical pathways consisting
120  of orthologous gene groups (KO groups) participating in protein complexes (Fig. 1f) and/or enzymatic reactions within
121  modules (Fig. 1g). Modules are the building blocks of larger biochemical pathways. We first sought to validate the
122  performance of EvoWeaver at identifying KO groups that participate in the same complex. We anticipated a strong
123  coevolutionary signal for these pairs because of their mutual dependence. Each algorithm's performance was graded
124  on its ability to distinguish 867 pairs of KO groups that complex (positives) versus 867 randomly selected pairs of
125  unrelated KO groups (negatives). The negative set was constructed from a weighted random sample of 57,321
126  unrelated KO groups. Weighted sampling reduces risk of overfitting by matching the distribution of data features in
127  the negative set to the positive set.
128     Almost all coevolution algorithms performed well at identifying KO groups involved in the same complex (Fig. S1).
129  Sequence Level methods performed slightly worse than other categories of coevolutionary signal. This outcome was
130  expected because many non-interacting proteins appear to physically interface similarly to interacting proteins[60]. The
131  predictions of most algorithms were weakly correlated with each other, which suggests combining signals could
132  further improve performance (Fig. S1). To this end, we evaluated three ensemble methods (Logistic Regression,
133  Random Forest, and Neural Network) using five-fold cross-validation. All ensemble methods displayed predictive
134  power exceeding component coevolutionary signals, with Logistic Regression performing the best (Fig. S1).
135     Given EvoWeaver's excellent performance on the Complexes benchmark, we next sought to establish its ability to
136  identify functionally associated protein-coding genes that were not involved in the same protein complex. To this end,
137  we developed the Modules benchmark as a set of 899 pairs of gene groups acting in adjacent steps of a biochemical
138  pathway (positives) and 899 randomly selected pairs from disconnected pathways (negatives). This task is more
139  challenging because proteins involved in the same module need not physically interact (Fig. 1g). As shown in Figure
140  2, performance of component algorithms on the Modules benchmark was slightly worse than on the Complexes
141  benchmark. However, ensemble methods retained high performance (AUROC of 0.955 for Random Forest) and
142  outperformed individual coevolutionary signals. The gap between ensemble and component predictors highlights the
143  importance of using multiple coevolutionary signals to infer functional associations.

144  Next, we sought to determine whether EvoWeaver's ensemble predictions were transferrable to a different
145  database than KEGG. To this end, we used the CORUM database to test EvoWeaver's ability to identify human
146  proteins that participate in common complexes (Fig. S2). The 12 component algorithms were less accurate on
147  CORUM than KEGG, which was consistent with prior evaluations of coevolutionary algorithms on the CORUM
148  benchmark[12]. As expected, the ranking of component algorithms was different on CORUM than KEGG due to
149  CORUM's stronger focus on eukaryotes. EvoWeaver's ensemble method trained on the KEGG Modules benchmark
150  did not show an advantage over the best component algorithm. We attributed this discrepancy to differences between
151  the two databases and CORUM's sole emphasis on physical protein-protein interactions. However, retraining the
152  ensemble classifier on the CORUM database resulted in a substantial increase in accuracy above that of any
153  component algorithm. This result underscores the fact that transferability of accurate ensemble predictions relies
154  upon a shared prediction objective. More detail on the discrepancies between the CORUM and KEGG benchmarks is
155  discussed in sections **SI2** and **SI3** of the Supplemental Information.
156
157  **EvoWeaver infers hierarchical relationships among genes**
158
159  Functional relationships among genes exist at multiple levels, ranging from physically interacting to merely being
160  part of the same cellular environment. For this reason, it would be ideal to predict a strength of coevolution across a
161  hierarchy of multi-level relationships among gene groups. We evaluated the Random Forest model on pairs of KEGG
162  module blocks belonging to each of five classes: "Direct Connection", "Same Module", "Same Pathway", "Same
163  Global Pathway", and "Unrelated" module blocks. These classes are arranged in a hierarchy of decreasing functional
164  association. Accurate classification on this Multiclass benchmark would imply EvoWeaver can construct a hierarchical
165  classification scheme of genes and recapitulate the relationships in KEGG. We then used five-fold cross-validation to
166  predict class membership for 642,770 pairs of module blocks (Fig. 3). Notably, all 12 predictors contributed to the
167  ensemble classifier's accuracy (Fig. 3b). Most Random Forest predictions were assigned to the correct class or the
168  adjacent class (Fig. S3), even when requiring at least 50% confidence for prediction (Fig. 3a). Unsurprisingly, the
169  model frequently confused the "Same Global Pathway" and "Unrelated" classes, which are both expected to contain
170  weakly coevolving genes.
171  The Random Forest ensemble classifier was best at distinguishing the top two from bottom three hierarchical
172  classes. Hence, we tested whether these predictions could be used to recapitulate KEGG pathways by building a
173  network of module blocks with connections between pairs predicted as "Direct Connection" or "Same Module". We
174  applied Louvain clustering[61] to detect communities within this network. A randomly selected community is shown in
175  Figure 3c-d, which included all module blocks involved in the prodigiosin biosynthesis pathway. EvoWeaver correctly
176  identified most "Direct Connection" pairs within the pathway and properly distinguished the two modules within the
177  pathway. However, EvoWeaver incorrectly classified many "Same Module" pairs as "Direct Connection". This
178  analysis suggests EvoWeaver's predictions can be used to hypothesize biochemical pathways, although
179  EvoWeaver's predictions do not provide directionality to biochemical steps.
180
181  **EvoWeaver rivals STRING without reliance upon external data**
182
183  STRING is one of the most comprehensive databases of knowledge about functionally associated genes. One of
184  STRING's stated goals[57] is to predict genes belonging to the same pathway in KEGG, which corresponds to
185  EvoWeaver's "Direct Connection", "Same Module", and "Same Pathway" classifications. STRING's Total Score is a
186  composite of seven evidence streams[24]. We applied STRING's formula for Total Score to quantify the marginal
187  benefit of each evidence stream (Fig. S4). External data, including mining the literature for cooccurrence of terms
188  (Text Mining) and knowledge bases such as KEGG (Databases), provided the majority of STRING's predictive
189  performance (Fig. 4a). As expected, STRING's coevolutionary evidence streams (Cooccurrence, Gene
190  Neighborhood) were correlated with comparable signals derived by EvoWeaver (Fig. 4b). Excluding Text Mining,
191  EvoWeaver nearly matches STRING at its stated goal of predicting pairs of gene groups sharing a functional pathway
192  in KEGG (Fig. 4a). This is especially notable given that STRING's Database evidence incorporates KEGG itself as a
193  predictor, whereas EvoWeaver only relies on information extracted from genome sequences. This makes EvoWeaver
194  particularly powerful for identifying unknown functional associations without reliance on prior knowledge, which may
195  help to mitigate the problem of annotation inequality[1,2].
196
197  **EvoWeaver can inform novel hypotheses**
198
199  EvoWeaver's primary purpose is to serve as a generator for novel hypotheses about functional associations. As a
200  case study, we examined one of EvoWeaver's high confidence mispredictions, which was between human genes
201  *B3GNT5* and *ST6GAL1*. *B3GNT5* encodes an enzyme responsible for the synthesis of lactotriaosylceramide, the
202  primary precursor for lacto- and neolacto-series glycosphingolipids, and this enzyme is known to play a role in a

203  variety of human diseases[62,63]. *ST6GAL1* is responsible for the α2,6-sialylation of N-glycosylated proteins. Despite
204  *B3GNT5* and *ST6GAL1* having no common modules or pathways in the KEGG database (Fig. 5a), EvoWeaver
205  predicted this pair to be "Direct Connection" with probability 0.63 or "Same Module" with probability 0.36 (Fig. 5b).
206  This finding is consistent with experimental evidence showing mutations in glycosphingolipid biosynthetic enzymes
207  can cause changes in sialylation of N-glycosylated membrane-bound proteins[64] and, specifically, mutations in
208  *B3GNT5* modulate α2,6-sialylation of membrane-bound glycoproteins in ovarian cancer cells by directly silencing the
209  expression of *ST6GAL1* in several human cell lines[65]. EvoWeaver's prediction was supported by Phylogenetic
210  Profiling evidence because of the multiple inferred simultaneous gains of both genes (Fig. 5c) along with moderate
211  evidence for Gene Organization due to conservation in gene orientation and relative distance across the phylogeny
212  (Fig. 5d). *B3GNT5* and *ST6GAL1* also displayed strong similarity in their genealogies (Fig. 5e) and moderate
213  evidence for coevolutionary signal at the sequence level (Fig. 5f). While both *B3GNT5* and *ST6GAL1* have functional
214  associations with *B4GALT* family genes in KEGG (Fig. 5a), EvoWeaver's ensemble method did not identify a
215  connection between *ST6GAL1* and *B4GALT* family genes (Fig. S5), suggesting that the predicted linkage between
216  *B3GNT5* and *ST6GAL1* is unlikely to have resulted from transitivity.
217      To further substantiate EvoWeaver's power as a hypothesis generating tool, we investigated the top 100
218  mispredictions wherein a pair of genes were classified as "Same Pathway" in KEGG, but EvoWeaver predicted them
219  to be "Direct Connection". Many of these gene pairs were actually directly connected (19%) or separated by only a
220  few genes in a KEGG pathway (Fig. S6), but were categorized as "Same Pathway" because they lacked connections
221  in a common module. Therefore, the top mispredictions were partly artifacts of how the KEGG database defines
222  modules within pathways. We also investigated the top mispredictions wherein a pair was classified as "Same Global
223  Pathway" in KEGG, but EvoWeaver predicted "Direct Connection". Of the top five misclassifications, three involved
224  gene pairs between KEGG module M00892 (UDP-GlcNAc biosynthesis in eukaryotes) and KEGG module M00055
225  (N-glycan precursor biosynthesis). Coevolutionary (Fig. S5) and experimental evidence support the
226  interconnectedness of these modules: N-glycan branching is hypersensitive to UDP-GlcNAc concentrations in
227  mammals[66], UDP-GlcNAc transporters are involved in delivery of N-glycan substrates in plants[67], and components of
228  the UDP-GlcNAc biosynthetic pathway are required for complex N-glycan synthesis in *C. elegans*[68].
229      Next, we asked whether EvoWeaver can contribute in cases where a set of genes is implicated in a common
230  function but their interrelationships are unknown. We investigated EvoWeaver's predictions for six sets of genes
231  comprising discrete biochemical pathways (Fig. 6). The four categories of coevolutionary algorithms often disagreed
232  with each other and differed from the connections in KEGG. However, EvoWeaver's ensemble predictions generated
233  more accurate connections, which reinforces the notion that merging evidence streams improves predictions. Taken
234  together, these findings suggest that EvoWeaver can be used to augment existing biological knowledge by predicting
235  credible gene functional associations.
236
237  DISCUSSION
238
239      EvoWeaver represents a marked advancement in employing coevolutionary principles to the discovery of
240  functional associations. In this work, we showed that EvoWeaver can capitalize on multiple sources of coevolutionary
241  signal to generate a more complete understanding of the functional relationships between gene groups. Importantly,
242  EvoWeaver's ensemble predictions have the advantage that they do not require users to choose which
243  coevolutionary signals are appropriate for a particular context. EvoWeaver's accuracy permitted us to construct a
244  hierarchical model of functional associations that was able to partly recapitulate experimentally validated KEGG
245  pathways without any prior knowledge of the proteins other than their coding sequences and genomic locations.
246  Moreover, we demonstrated how EvoWeaver's predictions can be leveraged to infer novel functional associations
247  that are absent from large databases of biological knowledge.
248      EvoWeaver excels at three characteristics that are necessary for the practical application of coevolutionary
249  analyses on large-scale datasets. First, EvoWeaver is highly scalable owing to its optimized algorithms. We
250  demonstrated this by applying EvoWeaver to 1,545 gene groups from 8,564 genomes across the tree of life. To our
251  knowledge, this is the largest coevolutionary analysis to date in terms of number of genomes analyzed, exceeding the
252  2,167 genomes analyzed in previous work[12,13]. Unlike popular prior approaches, such as ContextTree or SVD-
253  phy[42,69], EvoWeaver's pairwise comparisons are independent and can be readily distributed across a cluster of
254  computers. Second, EvoWeaver's predictions are higher accuracy because they incorporate multiple sources of
255  coevolutionary signal, and each component algorithm incorporates statistical testing that mitigates spurious signals.
256  Third, EvoWeaver standardizes the application of multiple algorithms within a single software package with consistent
257  inputs and outputs. This addresses usability issues previously identified in reviews of coevolutionary analyses[13].
258      Coevolution differs from protein-protein interactions in that it does not require any physical interaction. Many prior
259  approaches exist for predicting protein-protein interactions, along with databases of known interactors[53,54,70,71].
260  Benchmarking functional association algorithms presents its own challenges, as proteins that do not physically
261  interact may nevertheless be functionally associated[17]. This renders common benchmarks for protein-protein

262  interactions insufficient for benchmarking coevolutionary algorithms[71-73]. We chose to rely on the KEGG database as
263  a source of experimentally validated functional associations within a multi-level hierarchy. Although KEGG is limited
264  in size (i.e., 26,418 orthology groups), it is one of the few comprehensive sources of genomes and genes linked
265  across pathways.
266      We anticipate EvoWeaver to be particularly useful for generating hypotheses that catalyze investigations into
267  understudied proteins. EvoWeaver allows users to search through millions of gene pairs to find a comparatively small
268  number of potential functional associations. EvoWeaver's predictions are particularly valuable when combined with
269  network analyses or expert insights. In the future, EvoWeaver will assist in curating and supplementing large
270  databases of biological knowledge to address errors and annotation inequality. We also expect EvoWeaver's
271  predictions to be useful for other sequence features, such as non-coding RNAs, although protein-coding genes were
272  the focus of this study. Most importantly, EvoWeaver empowers users to combat annotation inequality by predicting
273  functional associations for the rapidly expanding collection of sequences with unknown function.

274 ONLINE METHODS
275
276 **Code Availability and Experimental Details**
277
278 EvoWeaver is available as part of the *SynExtend* (v1.17.1) package[74] for R[75], which is distributed via the
279 Bioconductor[76] platform. A comprehensive description of input/output and examples of running each algorithm are
280 contained in the supplementary R Markdown file available on GitHub
281 (https://github.com/WrightLabScience/EvoWeaver-ExampleCode). Briefly, users first construct an EvoWeaver object
282 with the *EvoWeaver* function using input gene groups, and then run the *predict* method to generate predictions using
283 any of the 12 component algorithms. Depending on the algorithm, the input consists of a reference tree, gene trees,
284 positional data, or sequences. The output is a matrix of scores, representing the pairwise strength of coevolution
285 measured by each algorithm between each pair of gene groups. Scores range from -1 (strong negative association)
286 to +1 (strong positive association). Detailed information about individual algorithms is described below and in **SI1** of
287 the Supplemental Information.
288 All analyses were performed with R (v4.3.3). Algorithms were implemented in EvoWeaver using the R and C
289 programming languages, with user-exposed methods available in R via the *SynExtend* package (v1.17.1). *SynExtend*
290 is dependent on the *DECIPHER* package (v3.0.0) and is distributed via the Bioconductor software repository[76]. Area
291 under the receiver operator characteristic curves (AUROC) and precision-recall curves (AUPRC) were calculated with
292 the *AUC* function in the *DescTools* package (v0.99.49) for R. Scripts for reproducing all analyses are available on
293 GitHub (https://github.com/WrightLabScience/EvoWeaver-ExampleCode). Data and pretrained ensemble models
294 used in this work are also available on GitHub, and larger datafiles are available from Zenodo (DOI:
295 10.5281/zenodo.13256882).
296 Local analyses were performed on a MacBook Pro with M1 Pro CPU and 32GB of RAM. Runtimes were
297 measured on a Dell PowerEdge T650 with an Intel Xeon processor (E5-2690 v4 2.6 GHz) and 792 GB of memory
298 running Ubuntu 22.04.4 LTS. Distributed computing was performed on the Open Science Grid[77]. Phylogenetic tree
299 reconstruction used eight core nodes with 8 - 16 GB RAM and 8 GB disk space, and pairwise coevolutionary score
300 calculations with EvoWeaver used single core nodes with 2 - 4 GB RAM and 2 - 4 GB disk space. Computers
301 matching these node specifications varied based on availability and Open Science Grid scheduling.
302
303 **Coevolutionary Algorithms in EvoWeaver**
304
305 The goal of EvoWeaver is to capture a holistic view of coevolution for predicting functional associations between
306 groups of genes. To achieve this, we implemented 12 algorithms from scratch that quantify different sources of
307 coevolutionary signal. Each algorithm analyzes a pair of gene groups and returns a score between zero and one,
308 where zero represents an absence of signal and more positive scores imply greater coevolutionary signal. Some
309 algorithms can provide scores between -1 and 1, in which case rare negative scores represent an inverse
310 coevolutionary association. To correct for spurious signal resulting from insufficient information, we multiply all scores
311 by their significance (1 – *p-value*). The resulting final scores are combined into an overall prediction using an
312 ensemble machine learning method. When an algorithm cannot make a prediction for a particular pair, the final score
313 passed to the ensemble method for that algorithm is zero. For example, if a pair of genes do not cooccur in any
314 organisms, then their final score for all Gene Organization algorithms is zero. The 12 algorithms we implemented fall
315 into four categories: Phylogenetic Profiling, Phylogenetic Structure, Gene Organization, and Sequence Level
316 methods (Fig. 1a). Of these, four algorithms are completely novel (G/L Distance, P/A Overlap, RP ContextTree, and
317 RP MirrorTree), four are new applications of existing algorithms (TreeDistance, Moran's I, Orientation MI, Gene
318 Vector), and the remaining four are refinements on existing algorithms. Computational scaling for all algorithms in
319 terms of number of gene groups and size of each gene group is available in **SI3** of the Supplemental Information and
320 Figure S7.
321
322 *Phylogenetic Profiling*
323 Phylogenetic profiling is a common technique that uses presence/absence (P/A) profiles of genes to investigate
324 shared function. The approaches previously introduced in the literature use binary P/A profiles, where one represents
325 the presence of a gene and zero represents its absence[78]. The first P/A approach used Hamming distances on binary
326 profiles as a score[79]. Later, Jaccard index and mutual information (MI) were applied to score P/A profiles[12,80].
327 Subsequent work accounted for clade-wise conservation[25] or transformed P/A profiles into ancestral gain/loss (G/L)
328 events and scored the correlation between events[81]. These transformations reduce redundancy for sets of organisms
329 with low rates of gene gain and loss[36,81].
330 EvoWeaver includes four Phylogenetic Profiling algorithms (Figs. 1b & S8). The first algorithm, P/A Jaccard, uses
331 the centered Jaccard index[82] of P/A profiles with conserved clades collapsed to mitigate bias from closely related
332 closely organisms. The second algorithm, P/A Overlap, applies Fitch Parsimony[83] to infer ancestral states on the

333  reference tree from P/A profiles and calculates the proportion of the tree both genes appear together relative to their
334  overall prevalence. The third algorithm, G/L MI, calculates weighted MI of G/L events (G/L profiles). G/L profiles
335  include three states: -1 for gene loss, 0 for no change, and +1 for gene gain. G/L MI uses the weighted mutual
336  information of four cases: simultaneous concordant transitions (i.e., gain/gain or loss/loss), simultaneous gain in gene
337  one and loss in gene two, simultaneous gain in gene two and loss in gene one, and non-simultaneous transitions. MI
338  is calculated by weighting the first case with +1, the second and third cases with -1, and the fourth case with 0.
339      G/L MI fails to adequately measure compensatory changes that do not occur on the same branch of the reference
340  tree, which are common in sequence evolution[84]. The fourth algorithm, G/L Distance, complements the previous
341  algorithms by quantifying the evolutionary distance between G/L events assuming the time between gain or loss
342  events is exponentially distributed. Thus, the score between a pair of events for two gene groups is calculated
343  as $we^{-d(v_1, v_2)}$, where $w$ is +1 if the events are the same (i.e., both gain or both loss) and -1 if the events are different,
344  and $d(v_1, v_2)$ is the distance between events $v_1$ and $v_2$ on the reference tree. The distance between events on
345  separate branches is defined as the total distance between their branch midpoints. The distance between events on
346  the same branch is defined as zero. For each pair of genes, events are paired to their closest event from the other
347  group. The total score for the gene pair is the average score for all event pairs, and ranges from -1 to +1.
348      Statistical significance for P/A Jaccard is calculated using an empirical distribution of scores obtained from
349  bootstrapping P/A vectors. Significance for G/L MI is calculated using Fisher's Exact Test on the contingency table of
350  the four cases, and p-values for P/A Overlap and G/L Distance are calculated using empirical values from
351  permutation testing.
352
353  *Phylogenetic Structure*
354      Gene tree structural comparisons were pioneered by MirrorTree[40], which scores each pair of gene groups by the
355  correlation of their pairwise sequence distances. Subsequent improvements to MirrorTree attempted to correct for
356  background evolutionary signal prior to analysis[85]. These extensions, often referred to as ContextTree or
357  ContextMirror, use different approaches to remove shared signal represented by the reference tree[39,69,86]. More
358  recently, SVD-phy was introduced as an alternative approach using SIMAP[87] or BLAST to measure distance between
359  sequences[42,88]. SVD-phy uses singular value decomposition (SVD) to reduce redundant information contained in
360  pairwise distances, which removes signal shared across all genes and improves overall predictions. However, this
361  approach requires that all pairwise distances be simultaneously kept in memory.
362      EvoWeaver uses random projection in lieu of SVD for dimensionality reduction. Random projection (RP) is a
363  surjective mapping that approximately preserves distances between vectors[89]. While traditional RP uses a large
364  matrix of random values, this requirement can be circumvented by generating values of the matrix on demand with a
365  preset random seed. Hence, this dimensionality reduction can be done with negligible memory overhead, allowing for
366  efficient and replicable distribution across a compute cluster. The RP MirrorTree algorithm applies RP to patristic
367  distances and scores pairs of vectors using Spearman's correlation coefficient. The RP ContextTree algorithm also
368  subtracts the reference tree from each distance matrix prior to random projection and scoring. RP ContextTree's final
369  scores are multiplied by the Hamming distance of overlap in organism membership to correct for spurious
370  correlations caused by minimally overlapping sets. Statistical significance for both RP ContextTree and RP
371  MirrorTree are calculated using the closed form solution for significance of Spearman's correlation coefficient.
372      EvoWeaver also incorporates tree distance metrics to measure topological similarity. A variety of previously
373  benchmarked metrics[43] were implemented as measures of functional similarity, all of which were highly correlated in
374  their tree distances. By default, EvoWeaver's Tree Distance predictor uses normalized Robinson-Foulds Distance
375  due to its low memory requirement and closed form solution for significance[90]. The score for each pair of genes was
376  defined as one minus the tree distance of the gene trees pruned to their common leaves. If two gene groups do not
377  appear in any common genomes, their Tree Distance score is set to zero.
378
379  *Gene Organization*
380      Gene organization is commonly used as a signature of functional association. For example, *a priori* knowledge of
381  genes that colocalize can be used to find biosynthetic gene clusters. Existing programs, such as antiSMASH[47], use
382  profile hidden Markov models to search for clusters of genes with known functional associations. However, these
383  approaches cannot be used to find gene clusters *de novo*. STRING makes use of the distance in nucleotides
384  between genes as a *de novo* predictor of functional association[24]. To our knowledge, analysis of gene organization is
385  one of the most understudied approaches for *de novo* prediction of functional associations.
386      EvoWeaver incorporates three Gene Organization algorithms. Together, they provide a well-rounded view of gene
387  organization: the first algorithm looks at whether genes possibly share regulation, the second measures how closely
388  genes are located to each other, and the third quantifies the extent to which gene distances are preserved across
389  phylogenies. The first algorithm, Orientation MI, examines the relative orientation of paired genes. Conservation of
390  relative gene direction has been validated in prior work to be indicative of shared function[19]. The score for Orientation

391 MI is defined as the bidirectional mutual information[91] between the orientation of paired genes, with Fisher's Exact
392 Test used to determine statistical significance.
393     The second algorithm, Gene Distance, examines the separation between genes. For each pair of genes on the
394 same chromosome or contig, the distance $d$ is calculated as the absolute value of the difference in gene index. The
395 index of a gene is its gene order in the chromosome or contig, starting from one for the first gene. We used indices
396 rather than nucleotide locations to mitigate the effect of variability in gene lengths. The score for each pair of
397 sequences is defined as $e^{1-d}$, and the overall score for a pair of gene groups is the mean of their sequence pair
398 scores. In this way, Gene Distance is maximized (1) when two genes are always adjacent ($d = 1$). Statistical
399 significance is derived from the distribution of distances between two random points on a line segment[92]. If a pair of
400 gene groups never appear in the same organism on the same chromosome/contig, the score for the pairing is defined
401 as zero.
402     The third algorithm, Moran's I, measures spatial autocorrelation among gene distances. Moran's I requires
403 pairwise weights represented by the inverse exponential of the patristic distances[93] and values in the form of gene
404 distances ($d$). Moran's I measures the extent to which the relative distances between genes are correlated with the
405 evolutionary trajectories of their respective organisms on the reference tree. Statistical significance is calculated using
406 the closed form solution to the expected value and variance of Moran's I (ref. [94]).
407
408 *Sequence Level Methods*
409     Covariation of residues is a common signal of protein-protein interactions, and numerous methods have been
410 devised for this purpose. A popular approach is direct coupling analysis[54], which fits a Potts model to a multiple
411 sequence alignment in order to parse "direct effects" from "indirect effects." Other algorithms using deep learning
412 have been successfully applied to sequencing data for finding interaction sites between proteins[95,96]. While some
413 previously developed approaches improved scaling[97,98], many of these algorithms have prohibitively high
414 computational complexity for high-throughput analysis. Additionally, the focus of these algorithms is on finding
415 interaction sites between small numbers of proteins or proteins known *a priori* to have a high likelihood of interacting.
416     EvoWeaver implements two Sequence Level methods that support either amino acid or nucleotide sequences,
417 although amino acid sequences were utilized for all analyses in this work. The first of these, Gene Vector, uses the
418 gene sequence natural vector approach, developed to predict protein-protein interactions[56]. We extended this
419 algorithm to amino acids following the same theoretical model as the initial nucleotide-based method. We chose to
420 use the natural vector without 2-mers or 3-mers, since the full vector incurred higher computational overhead with a
421 negligible difference in scores. For each pair of gene groups, we subset the sequences to the intersection of the
422 organisms present in both groups. The natural vector for each group in the pair is the average of the natural vectors
423 for each of its constituent sequences. We centered each natural vector assuming a null model of equally distributed
424 nucleotides or amino acids. The final score and statistical significance for the pairing are calculated from Spearman's
425 correlation coefficient of the natural vectors.
426     The second approach, Sequence Info, extends a prior approach to measure the MI between sites within
427 sequence alignments of each gene group[55]. For every pair of gene groups, we subset the sequences to the genomes
428 that appear in both groups, and subset the sites to those with high information content (entropy ≥ 0.3 bits) using the
429 *MaskAlignment* function in DECIPHER[99]. Mutual information is calculated for each pair of sites (i.e., columns) across
430 both alignments after applying a background entropy correction along with an average product correction[100]. The final
431 score is calculated as the average of the highest scoring pairing for each site. Statistical significance is calculated by
432 applying Fisher's combined probability test to the distribution of p-values across sites.
433
434 *Ensemble Methods*
435     EvoWeaver combines the output of each of all 12 coevolutionary algorithms into a final prediction using an
436 ensemble machine learning method (Fig. 2). For ensemble methods, we tested logistic regression, random forest,
437 and neural network models in R[75]. Logistic regression was performed with the *glm* function with *family='binomial'*,
438 random forests using the *randomForest* package[101] (v4.7-1.1), and neural networks using the *neuralnet* package
439 (v1.44.2). The random forest model used *maxdepth=25* for binary classification and *maxdepth=100* for multiclass
440 classification to avoid overfitting trees of unlimited depth. The neural network architecture was a feed forward network
441 with 12 inputs, one hidden layer of matched size (i.e., 12), two output nodes (i.e., class=0 or class=1), and sigmoid
442 activation functions on each node. We intentionally chose relatively simple architectures with default parameters for
443 our ensemble models to maintain interpretability of the predictions and mitigate overfitting to the dataset. All models
444 were evaluated using 5-fold cross-validation without hyperparameter tuning.
445     Only random forest was used for hierarchical classification due to its better performance in the binary
446 classification benchmarks. Hierarchical classification was also evaluated using 5-fold cross-validation. Members of
447 each class were distributed equally among each train/test fold. To prevent overfitting from high class imbalance in the
448 complete dataset, we downsampled classes in each training set to match the size of the smallest class, Direct
449 Connection, with 899 members. This meant that each class in the train set for each fold had 719 members (i.e.,

450    80%). Testing was done on the complete set of data partitioned for testing, which comprised 128,552-128,557
451    members (i.e., ~20%) per fold. Each pair was in exactly one test set, and no pairs belonged to both the train and test
452    set for any fold. Feature importance for the random forest model was calculated using permutation importance, which
453    was chosen over mean decrease in Gini impurity since the latter has been shown to produce biased estimates[102].
454        To construct an example network, we first created a weighted adjacency matrix from the random forest
455    predictions. Each node represented a single gene group and was connected to its top two "Direct Connection"
456    predictions with edges of weight 1.0. All predicted "Same Module" pairs were connected with edges of weight 0.5.
457    Our basis for this approach is that most nodes in KEGG are directly connected to two neighbors, and other nodes in
458    the same module are less important than direct connections. We then used Louvain clustering implemented in the
459    *igraph* package[103] (v1.5.0.1) to perform community detection. The network in Fig. 3c was randomly chosen from the
460    resulting communities.
461        A possible concern with holding out pairs in cross-validation is that ensemble methods could use spurious signals
462    to simply distinguish highly connected gene groups from less connected groups. On binary benchmarks, we further
463    validated our results by reevaluating our ensemble classifier using 10-fold cross-validation with gene group holdouts
464    rather than pair holdouts. Within each fold, 10% of gene groups were randomly selected, and all pairs involving at
465    least one of these groups was taken as the test set. The resulting train/test sets each comprised roughly 80/20% of
466    the data (respectively), which forms a comparable scenario to 5-fold cross-validation with pair holdouts. We also
467    evaluated the impact of module/complex holdouts, which were performed similarly to gene group holdouts. The
468    results of these analyses were virtually identical to prior results (Figs. S3,9-10), implying that EvoWeaver is not
469    heavily relying on spurious signals when making predictions. More details on factors impacting EvoWeaver's
470    performance are available in section **SI3** of the Supplemental Information.
471
472    **Construction of Benchmark Datasets**
473
474        The goal of the Complexes benchmark is to judge each algorithm's ability to discern genes encoding proteins
475    involved a complex versus genes encoding unrelated proteins. To this end, we identified all orthology groups
476    belonging to a complex in KEGG[104], for a total of 372 gene groups. We computed pairwise coevolutionary scores
477    between orthology groups with at least three sequences that were involved in a complex, for a total of 358 orthology
478    groups. This resulted in 57,321 pairs of orthology groups that are not in the same pathway (unrelated pairs) and 867
479    pairs participating as required or optional components of the same complex. Importantly, there was negligible
480    similarity between distinct orthology groups (Fig. S11), which might have otherwise resulted in data leakage. Positive
481    pairs were defined as the 867 pairs from the same complex, and an equivalent number of negative pairs were drawn
482    to create a balanced dataset for benchmarking. Random sampling of negative pairs was weighted in order to match
483    the distribution in number of sequences per gene group to that of the positive pairs. This weighted sampling was used
484    to mitigate the ability of algorithms to use the number of sequences per group as a proxy for functional association.
485        Next, we constructed the Modules benchmark to test each algorithm's ability to discern proteins acting in
486    subsequent steps of a biochemical pathway versus unrelated proteins. We first identified all module blocks within the
487    KEGG MODULES database. Each module block is a set of one or more orthology groups that perform a discrete step
488    within a biochemical pathway (Fig. 1g). Each module was parsed from its definition on KEGG (Table S1), for a total of
489    369 modules. Positive test cases were defined as successive blocks in a module, and negative cases were defined
490    as module blocks in separate modules not sharing a pathway in KEGG. KEGG's "Global and Overview Pathways"
491    were not considered, since their broad definition encompasses most proteins in KEGG. Blocks containing complexes
492    were also excluded to prevent overlap with the Complexes benchmark. Since some orthology groups belong to
493    multiple blocks, only pairs of blocks without overlap in orthology groups were assessed. The final Modules
494    benchmark was comprised of 1,187 blocks with 899 positive pairs. An equivalent number of negative pairs were
495    sampled in the same manner as in the Complexes benchmark.
496        Having constructed two binary benchmarks, we constructed the Multiclass benchmark to explore EvoWeaver's
497    ability to distinguish interaction strengths among proteins. Accordingly, we used the relationships encoded in the
498    KEGG PATHWAYS database to define multiple hierarchical levels of functional association. We assigned all pairs of
499    module blocks into one of five categories: "Direct Connection", "Same Module", "Same Pathway", "Same Global
500    Pathway", or "Unrelated". The "Same Pathway" group comprises pairs of module blocks that share a pathway not in
501    the "Global and Overview Pathways" category in KEGG, and the "Unrelated" group comprises pairs with no modules
502    or pathways in common. We chose 50% confidence as the cutoff for classification (Fig. 3a) because these predictions
503    have higher probability assigned to their predicted class than their sum of probabilities across all other classes. The
504    confusion matrix at 0% confidence is shown in Figure S3. To look for novel connections (Fig. 5), we examined pairs
505    belonging to "Unrelated" and "Same Global Pathway" groups that EvoWeaver predicted as being "Direct Connection".
506    More detail on benchmark datasets and data preprocessing is available in section **SI2** of the Supplemental
507    Information. A list of all misclassifications is available in Supplemental Datafile 1.
508

509 **Preparing Gene Groups for Analysis**
510
511     EvoWeaver takes as input a set of two or more gene trees, which may include sequences, gene indexes, and/or a
512 reference tree. It then applies the set of component algorithms for which it has the necessary input data types. We
513 obtained amino acid sequences for each gene group from KEGG and used DECIPHER[99] to align sequences and
514 construct neighbor-joining gene trees. In total, there were 8,564 genomes with at least one gene present in the
515 benchmarks. Reference trees were estimated using the ASTRID algorithm[105]. Impact of error in the reference tree is
516 discussed in sections **SI1** and **SI3** of the Supplemental Information and shown in Figure S12. To find each gene's
517 index within its genome, we downloaded the set of all genes available for each organism from KEGG, along with their
518 chromosome/contig, orientation, and location. We mapped locations to indices by calculating the index of each gene
519 relative to all other genes on the same chromosome/contig available for that genome. Of the 8,564 genomes present
520 in the benchmarks, 8,136 had location data available in KEGG. A taxonomic breakdown of the genomes used and
521 their location data is available in Supplemental Datafile 2.
522
523 **Comparison with STRING**
524
525     Data for STRING's clusters of orthologous genes (COGs) and interactions were downloaded from STRING v12.0.
526 Since STRING's COG membership sometimes did not perfectly correspond to KEGG's KO groups, we tabulated the
527 KO group assignments for sequences belonging to each STRING COG. Overall, 6,849 COGs had at least one
528 sequence that could be mapped to a KO group in KEGG. Each STRING COG was mapped to KEGG Module blocks
529 using its majority (≥ 50%) KEGG KO group. A total of 6,311 COGs had a majority KO group, and 4,481 (71%) of
530 these COGs had perfect consensus. Only 538 STRING COGs lacked a consensus KO group, and these COGs were
531 excluded from analysis.
532     STRING's stated goal for its Total Score is to estimate how likely a reported functional linkage between two
533 proteins "is at least as specific as that between an average pair of proteins annotated on the same 'map' or 'pathway'
534 in KEGG"[57]. Therefore, EvoWeaver's analogous predictions were made by summing the probabilities predicted for
535 "Direct Connection", "Same Module", and "Same Pathway" in the hierarchical classification (Fig. 3). A total of 757
536 pairs of COGs in the matched dataset belonged to the "Same Pathway", "Same Module", or "Direct Connection"
537 categories in KEGG. An equivalent number of negatives were randomly sampled from the remaining pairs. STRING
538 provides its Total Score calculation within a Python script available on their website. We used this formula to calculate
539 the hypothetical Total Score using subsets of STRING's evidence streams. The sequence of AUROCs in Figure 4a
540 was obtained by sequentially adding evidence streams from lowest to highest marginal impact on AUROC to the
541 Total Score calculation (Fig. S4).
542
543 **KEGG Case Studies**
544
545     Case studies for Fig. 6 were manually constructed from KEGG data for biologically meaningful sets of KEGG
546 modules belonging to the same KEGG pathway. Gene groups were connected according to the directed connections
547 available in the corresponding KEGG pathway. Only gene groups used in the Multiclass benchmark were included.
548 For the ensemble network, we connected each node to its top connection, where each connection's ranking is
549 determined from its probability of "Direct Connection" according to the Multiclass classifier used in Fig. 3. The
550 component predictor networks were constructed similarly to the ensemble network, but used the average rank of the
551 component score ranks for calculating each node's top connections. For example, the Phylogenetic Profiling
552 connections are determined by the mean rank of P/A Jaccard, P/A Overlap, G/L MI, and G/L Distance.
553
554 **List of Abbreviations Used**
555
556 **AUPRC**: Area under the PRC
557 **AUROC**: Area under the ROC curve
558 **COG**: Cluster of orthologous genes
559 **G/L**: Gain/loss
560 **KEGG**: Kyoto Encyclopedia of Genes and Genomes
561 **KO group**: KEGG orthology group
562 **MI**: Mutual information
563 **P/A**: Presence/absence
564 **PRC**: Precision-recall curve
565 **ROC**: Receiver operating characteristic
566 **RP**: Random projection
567 **SVD**: Singular value decomposition

568
574

575     REFERENCES

576
577   1     Kustatscher, G. *et al.* Understudied proteins: opportunities and challenges for functional proteomics. *Nature*
578        *Methods* **19**, 774-779 (2022). https://doi.org/10.1038/s41592-022-01454-x
579   2     Kustatscher, G. *et al.* An open invitation to the Understudied Proteins Initiative. *Nature Biotechnology* **40**,
580        815-817 (2022). https://doi.org/10.1038/s41587-022-01316-z
581   3     Sinha, S., Eisenhaber, B., Jensen, L. J., Kalbuaji, B. & Eisenhaber, F. Darkness in the Human Gene and
582        Protein Function Space: Widely Modest or Absent Illumination by the Life Science Literature and the Trend
583        for Fewer Protein Function Discoveries Since 2000. *PROTEOMICS* **18**, 1800093 (2018).
584        https://doi.org/10.1002/pmic.201800093
585   4     Haynes, W. A., Tomczak, A. & Khatri, P. Gene annotation bias impedes biomedical research. *Scientific*
586        *Reports* **8**, 1362 (2018). https://doi.org/10.1038/s41598-018-19333-x
587   5     Salzberg, S. L. Next-generation genome annotation: we still struggle to get it right. *Genome Biology* **20**, 92
588        (2019). https://doi.org/10.1186/s13059-019-1715-2
589   6     Lobb, B., Tremblay, B. J.-M., Moreno-Hagelsieb, G. & Doxey, A. C. An assessment of genome annotation
590        coverage across the bacterial tree of life. *Microbial Genomics* **6** (2020).
591        https://doi.org/10.1099/mgen.0.000341
592   7     Stoeger, T., Gerlach, M., Morimoto, R. I. & Nunes Amaral, L. A. Large-scale investigation of the reasons why
593        potentially important genes are ignored. *PLOS Biology* **16**, e2006643 (2018).
594        https://doi.org/10.1371/journal.pbio.2006643
595   8     Schnoes, A. M., Ream, D. C., Thorman, A. W., Babbitt, P. C. & Friedberg, I. Biases in the Experimental
596        Annotations of Protein Function and Their Effect on Our Understanding of Protein Function Space. *PLOS*
597        *Computational Biology* **9**, e1003063 (2013). https://doi.org/10.1371/journal.pcbi.1003063
598   9     Gillis, J. & Pavlidis, P. The Impact of Multifunctional Genes on "Guilt by Association" Analysis. *PLOS ONE* **6**,
599        e17258 (2011). https://doi.org/10.1371/journal.pone.0017258
600   10    Aravind, L. Guilt by association: contextual information in genome analysis. *Genome Research* **10**, 1074-
601        1077 (2000).
602   11    Codoñer, F. M. & Fares, M. A. Why should we care about molecular coevolution? *Evol Bioinform Online* **4**,
603        29-38 (2008).
604   12    Moi, D., Kilchoer, L., Aguilar, P. S. & Dessimoz, C. Scalable phylogenetic profiling using MinHash uncovers
605        likely eukaryotic sexual reproduction genes. *PLOS Computational Biology* **16**, e1007553 (2020).
606        https://doi.org/10.1371/journal.pcbi.1007553
607   13    Moi, D. & Dessimoz, C. Phylogenetic profiling in eukaryotes comes of age. *Proceedings of the National*
608        *Academy of Sciences* **120** (2023). https://doi.org/10.1073/pnas.2305013120
609   14    Canavati, C. *et al.* Using multi-scale genomics to associate poorly annotated genes with rare diseases.
610        *Genome Medicine* **16** (2024). https://doi.org/10.1186/s13073-023-01276-2
611   15    Kann, M. G., Shoemaker, B. A., Panchenko, A. R. & Przytycka, T. M. Correlated Evolution of Interacting
612        Proteins: Looking Behind the Mirrortree. *Journal of Molecular Biology* **385**, 91-98 (2009).
613        https://doi.org/10.1016/j.jmb.2008.09.078
614   16    Chikina, M., Robinson, J. D. & Clark, N. L. Hundreds of Genes Experienced Convergent Shifts in Selective
615        Pressure in Marine Mammals. *Molecular Biology and Evolution* **33**, 2182-2192 (2016).
616        https://doi.org/10.1093/molbev/msw112
617   17    Little, J., Chikina, M. & Clark, N. L. Evolutionary rate covariation is a reliable predictor of co-functional
618        interactions but not necessarily physical interactions. *eLife* **12**, RP93333 (2024).
619        https://doi.org/10.7554/eLife.93333
620   18    Umemura, M., Koike, H. & Machida, M. Motif-independent de novo detection of secondary metabolite gene
621        clusters-toward identification from filamentous fungi. *Front Microbiol* **6**, 371-371 (2015).
622        https://doi.org/10.3389/fmicb.2015.00371
623   19    Korbel, J. O., Jensen, L. J., Von Mering, C. & Bork, P. Analysis of genomic context: prediction of functional
624        associations from conserved bidirectionally transcribed gene pairs. *Nature Biotechnology* **22**, 911-917
625        (2004). https://doi.org/10.1038/nbt988

| 626 | 20 | Cotroneo, C. E., Gormley, I. C., Shields, D. C. & Salter-Townshend, M. Computational modelling of |
| 627 | | chromosomally clustering protein domains in bacteria. *BMC Bioinformatics* **22**, 593 (2021). |
| 628 | | https://doi.org/10.1186/s12859-021-04512-x |
| 629 | 21 | Feinauer, C., Szurmant, H., Weigt, M. & Pagnani, A. Inter-Protein Sequence Co-Evolution Predicts Known |
| 630 | | Physical Interactions in Bacterial Ribosomes and the Trp Operon. *PLOS ONE* **11**, e0149166 (2016). |
| 631 | | https://doi.org/10.1371/journal.pone.0149166 |
| 632 | 22 | Clark, G. W., Ackerman, S. H., Tillier, E. R. & Gatti, D. L. Multidimensional mutual information methods for |
| 633 | | the analysis of covariation in multiple sequence alignments. *BMC bioinformatics* **15**, 1-12 (2014). |
| 634 | 23 | Bitbol, A.-F. Inferring interaction partners from protein sequences using mutual information. *PLOS* |
| 635 | | *Computational Biology* **14**, e1006401 (2018). https://doi.org/10.1371/journal.pcbi.1006401 |
| 636 | 24 | Szklarczyk, D. *et al.* The STRING database in 2021: customizable protein–protein networks, and functional |
| 637 | | characterization of user-uploaded gene/measurement sets. *Nucleic Acids Research* **49**, D605-D612 (2021). |
| 638 | | https://doi.org/10.1093/nar/gkaa1074 |
| 639 | 25 | Stupp, D. *et al.* Co-evolution based machine-learning for predicting functional interactions between human |
| 640 | | genes. *Nature Communications* **12** (2021). https://doi.org/10.1038/s41467-021-26792-w |
| 641 | 26 | Tabach, Y. *et al.* Human disease locus discovery and mapping to molecular pathways through phylogenetic |
| 642 | | profiling. *Molecular Systems Biology* **9**, 692 (2013). https://doi.org/10.1038/msb.2013.50 |
| 643 | 27 | Tabach, Y. *et al.* Identification of small RNA pathway genes using patterns of phylogenetic conservation and |
| 644 | | divergence. *Nature* **493**, 694-698 (2013). https://doi.org/10.1038/nature11779 |
| 645 | 28 | Sherill-Rofe, D. *et al.* Mapping global and local coevolution across 600 species to identify novel homologous |
| 646 | | recombination repair genes. *Genome Research* **29**, 439-448 (2019). https://doi.org/10.1101/gr.241414.118 |
| 647 | 29 | Andreo-Vidal, A., Binda, E., Fedorenko, V., Marinelli, F. & Yushchuk, O. Genomic Insights into the |
| 648 | | Distribution and Phylogeny of Glycopeptide Resistance Determinants within the Actinobacteria Phylum. |
| 649 | | *Antibiotics* **10** (2021). https://doi.org/10.3390/antibiotics10121533 |
| 650 | 30 | Ding, D. *et al.* Co-evolution of interacting proteins through non-contacting and non-specific mutations. |
| 651 | | *Nature Ecology & Evolution* **6**, 590-603 (2022). https://doi.org/10.1038/s41559-022-01688-0 |
| 652 | 31 | Fongang, B., Zhu, Y., Wagner, E. J., Kudlicki, A. & Rowicka, M. Co-evolutionary analysis accurately predicts |
| 653 | | details of interactions between the Integrator complex subunits. *bioRxiv*, 696583 (2019). |
| 654 | | https://doi.org/10.1101/696583 |
| 655 | 32 | Ramani, A. K. & Marcotte, E. M. Exploiting the Co-evolution of Interacting Proteins to Discover Interaction |
| 656 | | Specificity. *Journal of Molecular Biology* **327**, 273-284 (2003). https://doi.org/https://doi.org/10.1016/S0022- |
| 657 | | 2836(03)00114-1 |
| 658 | 33 | Fukunaga, T. & Iwasaki, W. Inverse Potts model improves accuracy of phylogenetic profiling. *Bioinformatics* |
| 659 | | **38**, 1794-1800 (2022). https://doi.org/10.1093/bioinformatics/btac034 |
| 660 | 34 | Cheng, Y. & Perocchi, F. ProtPhylo: identification of protein–phenotype and protein–protein functional |
| 661 | | associations via phylogenetic profiling. *Nucleic Acids Research* **43**, W160-W168 (2015). |
| 662 | | https://doi.org/10.1093/nar/gkv455 |
| 663 | 35 | Ji, F. *et al.* DEPCOD: a tool to detect and visualize co-evolution of protein domains. *Nucleic Acids Research* |
| 664 | | (2022). https://doi.org/10.1093/nar/gkac349 |
| 665 | 36 | Škunca, N. & Dessimoz, C. Phylogenetic profiling: how much input data is enough? *PLOS ONE* **10**, |
| 666 | | e0114701 (2015). https://doi.org/10.1371/journal.pone.0114701 |
| 667 | 37 | Shin, J. & Lee, I. Co-Inheritance Analysis within the Domains of Life Substantially Improves Network |
| 668 | | Inference by Phylogenetic Profiling. *PLOS ONE* **10**, e0139006 (2015). |
| 669 | | https://doi.org/10.1371/journal.pone.0139006 |
| 670 | 38 | Clark, N. L., Alani, E. & Aquadro, C. F. Evolutionary rate covariation reveals shared functionality and |
| 671 | | coexpression of genes. *Genome Research* **22**, 714-720 (2012). https://doi.org/10.1101/gr.132647.111 |
| 672 | 39 | Pazos, F., Ranea, J. A., Juan, D. & Sternberg, M. J. Assessing protein co-evolution in the context of the tree |
| 673 | | of life assists in the prediction of the interactome. *J Mol Biol* **352**, 1002-1015 (2005). |
| 674 | | https://doi.org/10.1016/j.jmb.2005.07.005 |
| 675 | 40 | Pazos, F. & Valencia, A. Similarity of phylogenetic trees as indicator of protein–protein interaction. *Protein* |
| 676 | | *Engineering, Design and Selection* **14**, 609-614 (2001). https://doi.org/10.1093/protein/14.9.609 |
| 677 | 41 | Clark, G. W. *et al.* in *Network Biology: Methods and Applications* Vol. 781 (eds Gerard Cagney & Andrew |
| 678 | | Emili) 237-256 (Humana Press, 2011). |
| 679 | 42 | Franceschini, A., Lin, J., von Mering, C. & Jensen, L. J. SVD-phy: improved prediction of protein functional |
| 680 | | associations through singular value decomposition of phylogenetic profiles. *Bioinformatics* **32**, 1085-1087 |
| 681 | | (2016). https://doi.org/10.1093/bioinformatics/btv696 |
| 682 | 43 | Smith, M. R. Information theoretic generalized Robinson–Foulds metrics for comparing phylogenetic trees. |
| 683 | | *Bioinformatics* **36**, 5007-5013 (2020). https://doi.org/10.1093/bioinformatics/btaa614 |

684    44    Rokas, A., Wisecaver, J. H. & Lind, A. L. The birth, evolution and death of metabolic gene clusters in fungi.
685          *Nature Reviews Microbiology* **16**, 731-744 (2018). https://doi.org/10.1038/s41579-018-0075-3
686    45    Periwal, V. & Scaria, V. Insights into structural variations and genome rearrangements in prokaryotic
687          genomes. *Bioinformatics* **31**, 1-9 (2014). https://doi.org/10.1093/bioinformatics/btu600
688    46    Rocha, E. P. The organization of the bacterial genome. *Annual Review of Genetics* **42**, 211-233 (2008).
689    47    Blin, K. *et al.* antiSMASH 6.0: improving cluster detection and comparison capabilities. *Nucleic Acids*
690          *Research* **49**, W29-W35 (2021). https://doi.org/10.1093/nar/gkab335
691    48    Kautsar, S. A., Suarez Duran, H. G., Blin, K., Osbourn, A. & Medema, M. H. plantiSMASH: automated
692          identification, annotation and expression analysis of plant biosynthetic gene clusters. *Nucleic Acids*
693          *Research* **45**, W55-W63 (2017). https://doi.org/10.1093/nar/gkx305
694    49    Kautsar, S. A., Blin, K., Shaw, S., Weber, T. & Medema, M. H. BiG-FAM: the biosynthetic gene cluster
695          families database. *Nucleic Acids Research* **49**, D490-d497 (2021). https://doi.org/10.1093/nar/gkaa812
696    50    Davila Lopez, M., Martinez Guerra, J. J. & Samuelsson, T. Analysis of gene order conservation in
697          eukaryotes identifies transcriptionally and functionally linked genes. *PLOS ONE* **5**, e10654 (2010).
698          https://doi.org/10.1371/journal.pone.0010654
699    51    Thomas, J., Ramakrishnan, N. & Bailey-Kellogg, C. Graphical models of protein-protein interaction
700          specificity from correlated mutations and interaction data. *Proteins: Structure, Function, and Bioinformatics*
701          **76**, 911-929 (2009). https://doi.org/10.1002/prot.22398
702    52    Morcos, F., Hwa, T., Onuchic, J. N. & Weigt, M. Direct coupling analysis for protein contact prediction.
703          *Methods Mol Biol* **1137**, 55-70 (2014). https://doi.org/10.1007/978-1-4939-0366-5_5
704    53    Morcos, F. *et al.* Direct-coupling analysis of residue coevolution captures native contacts across many
705          protein families. *Proceedings of the National Academy of Sciences* **108**, E1293-E1301 (2011).
706          https://doi.org/10.1073/pnas.1111471108
707    54    Weigt, M., White, R. A., Szurmant, H., Hoch, J. A. & Hwa, T. Identification of direct residue contacts in
708          protein-protein interaction by message passing. *Proceedings of the National Academy of Sciences* **106**, 67-
709          72 (2009). https://doi.org/10.1073/pnas.0805923106
710    55    Martin, L. C., Gloor, G. B., Dunn, S. D. & Wahl, L. M. Using information theory to search for co-evolving
711          residues in proteins. *Bioinformatics* **21**, 4116-4124 (2005). https://doi.org/10.1093/bioinformatics/bti671
712    56    Zhao, N., Zhuo, M., Tian, K. & Gong, X. Protein–protein interaction and non-interaction predictions using
713          gene sequence natural vector. *Communications Biology* **5** (2022). https://doi.org/10.1038/s42003-022-
714          03617-0
715    57    Szklarczyk, D. *et al.* The STRING database in 2011: functional interaction networks of proteins, globally
716          integrated and scored. *Nucleic Acids Research* **39**, D561-D568 (2011). https://doi.org/10.1093/nar/gkq973
717    58    Kanehisa, M. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research* **28**, 27-30
718          (2000). https://doi.org/10.1093/nar/28.1.27
719    59    Kanehisa, M., Sato, Y., Furumichi, M., Morishima, K. & Tanabe, M. New approach for understanding
720          genome variations in KEGG. *Nucleic Acids Research* **47**, D590-d595 (2019).
721          https://doi.org/10.1093/nar/gky962
722    60    Launay, G., Ceres, N. & Martin, J. Non-interacting proteins may resemble interacting proteins: prevalence
723          and implications. *Scientific Reports* **7**, 40419 (2017). https://doi.org/10.1038/srep40419
724    61    Raghavan, U. N., Albert, R. & Kumara, S. Near linear time algorithm to detect community structures in large-
725          scale networks. *Phys Rev E Stat Nonlin Soft Matter Phys* **76**, 036106 (2007).
726          https://doi.org/10.1103/PhysRevE.76.036106
727    62    Wang, Z. *et al.* High expression of lactotriaosylceramide, a differentiation-associated glycosphingolipid, in
728          the bone marrow of acute myeloid leukemia patients. *Glycobiology* **22**, 930-938 (2012).
729          https://doi.org/10.1093/glycob/cws061
730    63    Togayachi, A. *et al.* Molecular cloning and characterization of UDP-GlcNAc: lactosylceramide β1, 3-N-
731          acetylglucosaminyltransferase (β3Gn-T5), an essential enzyme for the expression of HNK-1 and Lewis X
732          epitopes on glycolipids. *Journal of Biological Chemistry* **276**, 22032-22040 (2001).
733          https://doi.org/10.1074/jbc.M011369200
734    64    Boccuto, L. *et al.* A mutation in a ganglioside biosynthetic enzyme, ST3GAL5, results in salt & pepper
735          syndrome, a neurocutaneous disorder with altered glycolipid and glycoprotein glycosylation. *Human*
736          *Molecular Genetics* **23**, 418-433 (2013). https://doi.org/10.1093/hmg/ddt434
737    65    Alam, S. *et al.* Altered (neo-) lacto series glycolipid biosynthesis impairs α2-6 sialylation on N-glycoproteins
738          in ovarian cancer cells. *Scientific Reports* **7**, 45367 (2017). https://doi.org/10.1038/srep45367
739    66    Lau, K. S. *et al.* Complex N-glycan number and degree of branching cooperate to regulate cell proliferation
740          and differentiation. *Cell* **129**, 123-134 (2007). https://doi.org/10.1016/j.cell.2007.01.049
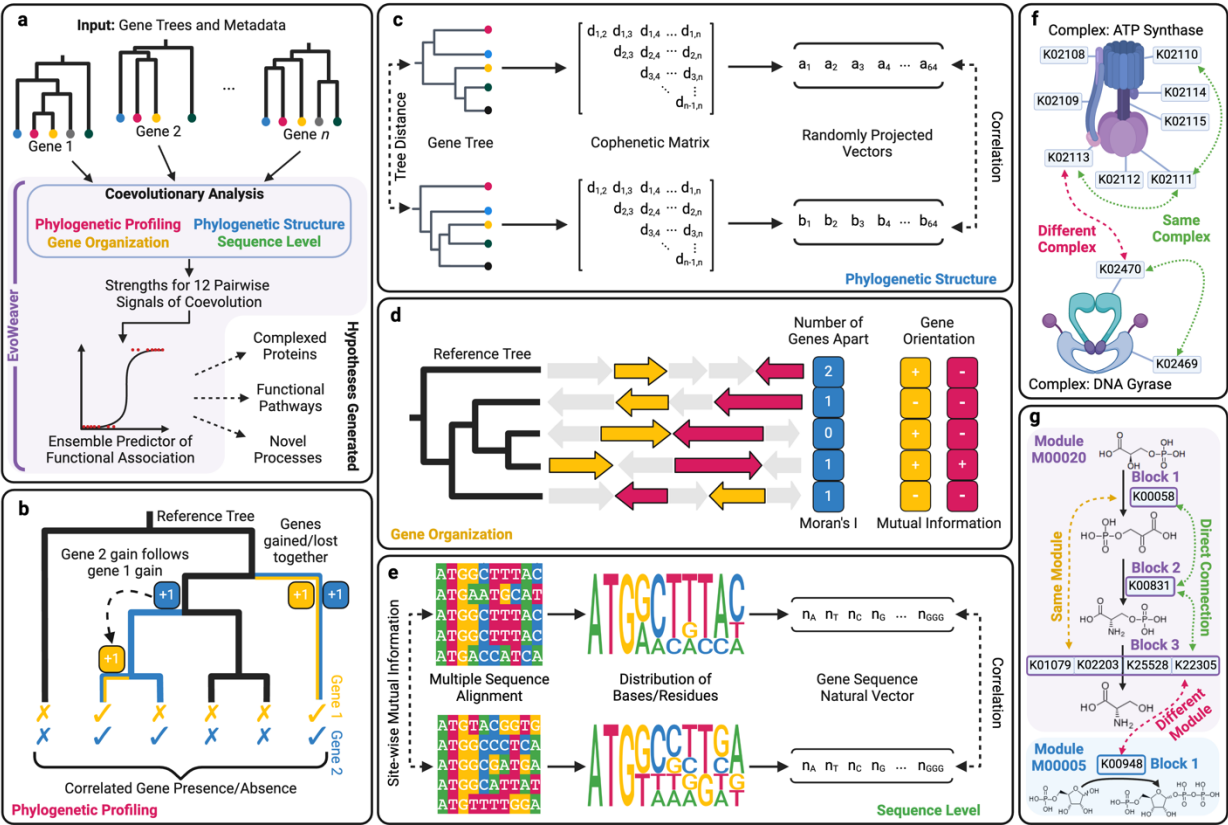741    67    Ebert, B. *et al.* A Golgi UDP-GlcNAc transporter delivers substrates for N-linked glycans and sphingolipids.
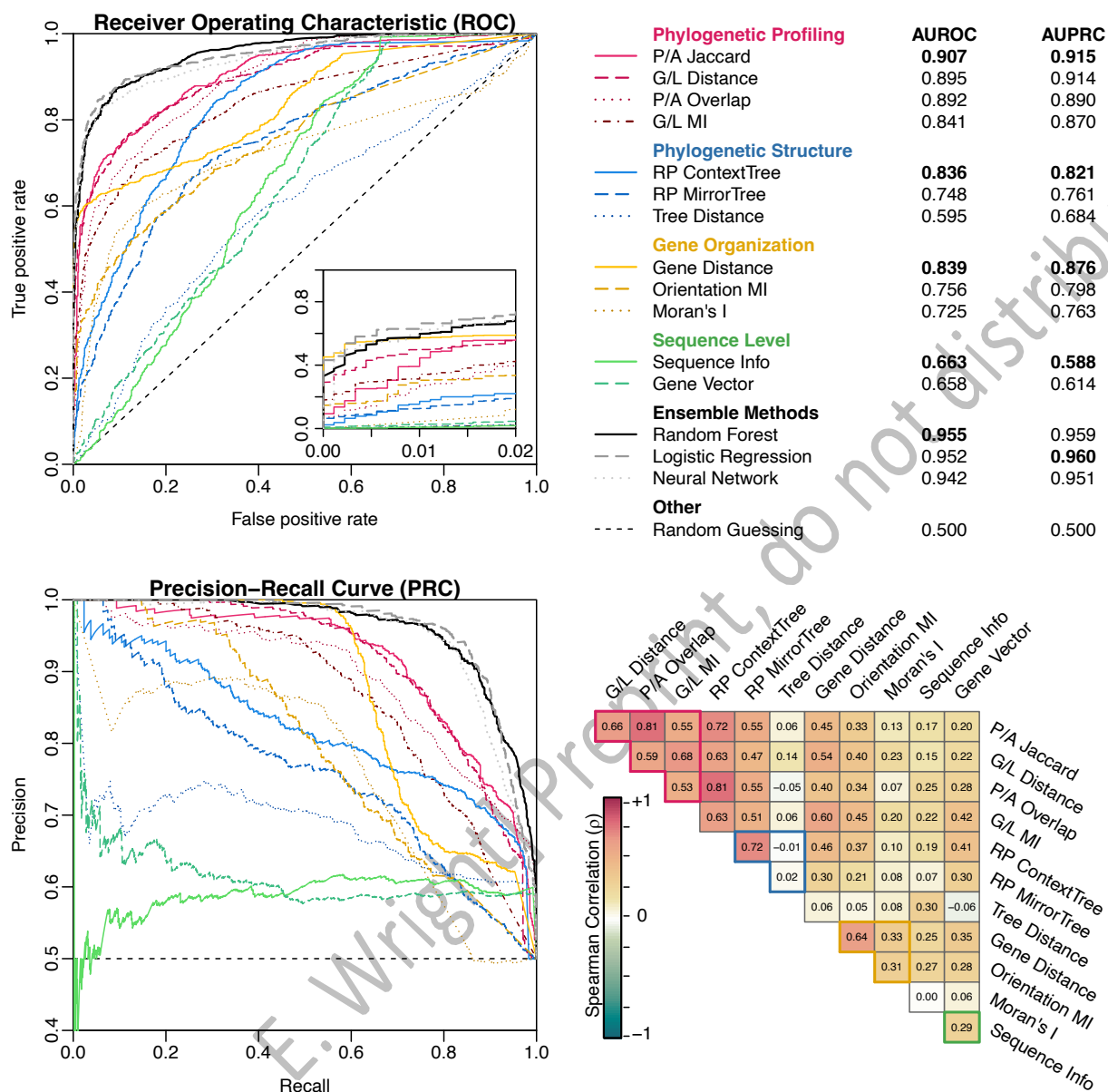742          *Nature Plants* **4**, 792-801 (2018). https://doi.org/10.1038/s41477-018-0235-5

743 68  Zhang, W. *et al.* Synthesis of paucimannose N-glycans by Caenorhabditis elegans requires prior actions of
744     UDP-N-acetyl-D-glucosamine: alpha-3-D-mannoside beta1, 2-N-acetylglucosaminyltransferase I, alpha3, 6-
745     mannosidase II and a specific membrane-bound beta-N-acetylglucosaminidase. *Biochemical Journal* **372**,
746     53-64 (2003). https://doi.org/10.1042/BJ20021931
747 69  Sato, T., Yamanishi, Y., Horimoto, K., Kanehisa, M. & Toh, H. Partial correlation coefficient between
748     distance matrices as a new indicator of protein-protein interactions. *Bioinformatics* **22**, 2488-2492 (2006).
749     https://doi.org/10.1093/bioinformatics/btl419
750 70  Luck, K. *et al.* A reference map of the human binary protein interactome. *Nature* **580**, 402-408 (2020).
751     https://doi.org/10.1038/s41586-020-2188-x
752 71  Blohm, P. *et al.* Negatome 2.0: a database of non-interacting proteins derived by literature mining, manual
753     annotation and protein structure analysis. *Nucleic Acids Research* **42**, D396-D400 (2014).
754     https://doi.org/10.1093/nar/gkt1079
755 72  Fabregat, A. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Research* **46**, D649-D655 (2018).
756 73  Oughtred, R. *et al.* TheBioGRIDdatabase: A comprehensive biomedical resource of curated protein, genetic,
757     and chemical interactions. *Protein Science* **30**, 187-200 (2021). https://doi.org/10.1002/pro.3978
758 74  Cooley, N., Lakshman, A. & Wright, E. S. SynExtend: Tools for Working with Synteny Objects. **v1.17.1**
759     (2024). https://doi.org/10.18129/B9.bioc.SynExtend
760 75  R Core Team. R: A language and environment for statistical computing. *R Foundation for Statistical*
761     *Computing*, Vienna, Austria (2021).
762 76  Gentleman, R. C. *et al.* Bioconductor: open software development for computational biology and
763     bioinformatics. *Genome Biology* **5**, 1-16 (2004). https://doi.org/10.1186/gb-2004-5-10-r80
764 77  OSG. Open Science Data Federation. *OSG* (2015). https://doi.org/https://doi.org/10.21231/0KVZ-VE57
765 78  Brilli, M. *et al.* Analysis of plasmid genes by phylogenetic profiling and visualization of homology
766     relationships using Blast2Network. *BMC Bioinformatics* **9**, 551 (2008). https://doi.org/10.1186/1471-2105-9-
767     551
768 79  Pellegrini, M., Marcotte, E. M., Thompson, M. J., Eisenberg, D. & Yeates, T. O. Assigning protein functions
769     by comparative genome analysis: Protein phylogenetic profiles. *Proceedings of the National Academy of*
770     *Sciences* **96**, 4285-4288 (1999). https://doi.org/10.1073/pnas.96.8.4285
771 80  Date, S. V. & Marcotte, E. M. Discovery of uncharacterized cellular systems by genome-wide analysis of
772     functional linkages. *Nature Biotechnology* **21**, 1055-1062 (2003). https://doi.org/10.1038/nbt861
773 81  Dembech, E. *et al.* Identification of hidden associations among eukaryotic genes through statistical analysis
774     of coevolutionary transitions. *Proceedings of the National Academy of Sciences* **120**, e2218329120 (2023).
775     https://doi.org/10.1073/pnas.2218329120
776 82  Chung, N. C., Miasojedow, B., Startek, M. & Gambin, A. Jaccard/Tanimoto similarity test and estimation
777     methods for biological presence-absence data. *BMC Bioinformatics* **20**, 644 (2019).
778     https://doi.org/10.1186/s12859-019-3118-5
779 83  Fitch, W. M. On the problem of discovering the most parsimonious tree. *The American Naturalist* **111**, 223-
780     257 (1977). https://doi.org/10.1086/283157
781 84  Kryazhimskiy, S., Dushoff, J., Bazykin, G. A. & Plotkin, J. B. Prevalence of Epistasis in the Evolution of
782     Influenza A Surface Proteins. *PLOS Genetics* **7**, e1001301 (2011).
783     https://doi.org/10.1371/journal.pgen.1001301
784 85  Juan, D., Pazos, F. & Valencia, A. High-confidence prediction of global interactomes based on genome-wide
785     coevolutionary networks. *Proceedings of the National Academy of Sciences* **105**, 934-939 (2008).
786     https://doi.org/10.1073/pnas.0709671105
787 86  Sato, T., Yamanishi, Y., Kanehisa, M. & Toh, H. The inference of protein-protein interactions by co-
788     evolutionary analysis is improved by excluding the information about the phylogenetic relationships.
789     *Bioinformatics* **21**, 3482-3489 (2005). https://doi.org/10.1093/bioinformatics/bti564
790 87  Arnold, R., Goldenberg, F., Mewes, H.-W. & Rattei, T. SIMAP—the database of all-against-all protein
791     sequence similarities and annotations with new interfaces and increased coverage. *Nucleic Acids Research*
792     **42**, D279-D284 (2013). https://doi.org/10.1093/nar/gkt970
793 88  Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J Mol*
794     *Biol* **215**, 403-410 (1990). https://doi.org/10.1016/s0022-2836(05)80360-2
795 89  Achlioptas, D. Database-friendly random projections. *Proceedings of the Twentieth ACM SIGMOD-SIGACT-*
796     *SIGART Symposium on Principles of Database Systems*, 274-281 (2001).
797 90  Steel, M. A. & Penny, D. Distributions of Tree Comparison Metrics—Some New Results. *Systematic Biology*
798     **42**, 126-141 (1993). https://doi.org/10.1093/sysbio/42.2.126
799 91  Beckley, A. M. & Wright, E. S. Identification of antibiotic pairs that evade concurrent resistance via a
800     retrospective analysis of antimicrobial susceptibility test results. *The Lancet Microbe* **2**, e545-e554 (2021).
801     https://doi.org/10.1016/s2666-5247(21)00118-x

802 92 Philip, J. The probability distribution of the distance between two random points in a box. *KTH Mathematics,*
803 *Royal Institute of Technology* (2007).
804 93 Gittleman, J. L. & Kot, M. Adaptation: Statistics and a Null Model for Estimating Phylogenetic Effects.
805 *Systematic Biology* **39**, 227-241 (1990). https://doi.org/10.2307/2992183
806 94 Cliff, A. D. & Ord, J. K. *Spatial Processes: Models and Applications*. (Pion Limited, 1981).
807 95 Pesaranghader, A. *et al.* deepSimDEF: deep neural embeddings of gene products and Gene Ontology
808 terms for functional analysis of genes. *Bioinformatics* **38**, 3051-3061 (2022).
809 https://doi.org/10.1093/bioinformatics/btac304
810 96 Soleymani, F., Paquet, E., Viktor, H. L., Michalowski, W. & Spinello, D. ProtInteract: a Deep Learning
811 Framework for Predicting Protein—Protein Interactions. *Computational and Structural Biotechnology Journal*
812 **21**, 1324-1348 (2023). https://doi.org/10.1016/j.csbj.2023.01.028
813 97 Ekeberg, M., Hartonen, T. & Aurell, E. Fast pseudolikelihood maximization for direct-coupling analysis of
814 protein structure from many homologous amino-acid sequences. *Journal of Computational Physics* **276**,
815 341-356 (2014). https://doi.org/10.1016/j.jcp.2014.07.024
816 98 Jones, D. T., Buchan, D. W. A., Cozzetto, D. & Pontil, M. PSICOV: precise structural contact prediction
817 using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics* **28**, 184-
818 190 (2011). https://doi.org/10.1093/bioinformatics/btr638
819 99 Wright, E. S. Using DECIPHER v2.0 to analyze big biological sequence data in R. *R Journal* **8**, 352-359
820 (2016). https://doi.org/10.32614/RJ-2016-025
821 100 Buslje, C. M., Santos, J., Delfino, J. M. & Nielsen, M. Correction for phylogeny, small number of
822 observations and data redundancy improves the identification of coevolving amino acid pairs using mutual
823 information. *Bioinformatics* **25**, 1125-1131 (2009). https://doi.org/10.1093/bioinformatics/btp135
824 101 Liaw, A. & Wiener, M. Classification and Regression by randomForest. *R News* **2**, 18-22 (2002).
825 102 Strobl, C., Boulesteix, A.-L., Zeileis, A. & Hothorn, T. Bias in random forest variable importance measures:
826 Illustrations, sources and a solution. *BMC Bioinformatics* **8**, 25 (2007). https://doi.org/10.1186/1471-2105-8-
827 25
828 103 Csárdi, G. & Nepusz, T. The igraph software package for complex network research. *InterJournal* **Complex**
829 **Systems**, 1695 (2006). https://doi.org/10.5281/zenodo.7682609
830 104 Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for
831 gene and protein annotation. *Nucleic Acids Research* **44**, D457-D462 (2015).
832 https://doi.org/10.1093/nar/gkv1070
833 105 Vachaspati, P. & Warnow, T. ASTRID: Accurate Species TRees from Internode Distances. *BMC Genomics*
834 **16**, S3 (2015). https://doi.org/10.1186/1471-2164-16-s10-s3
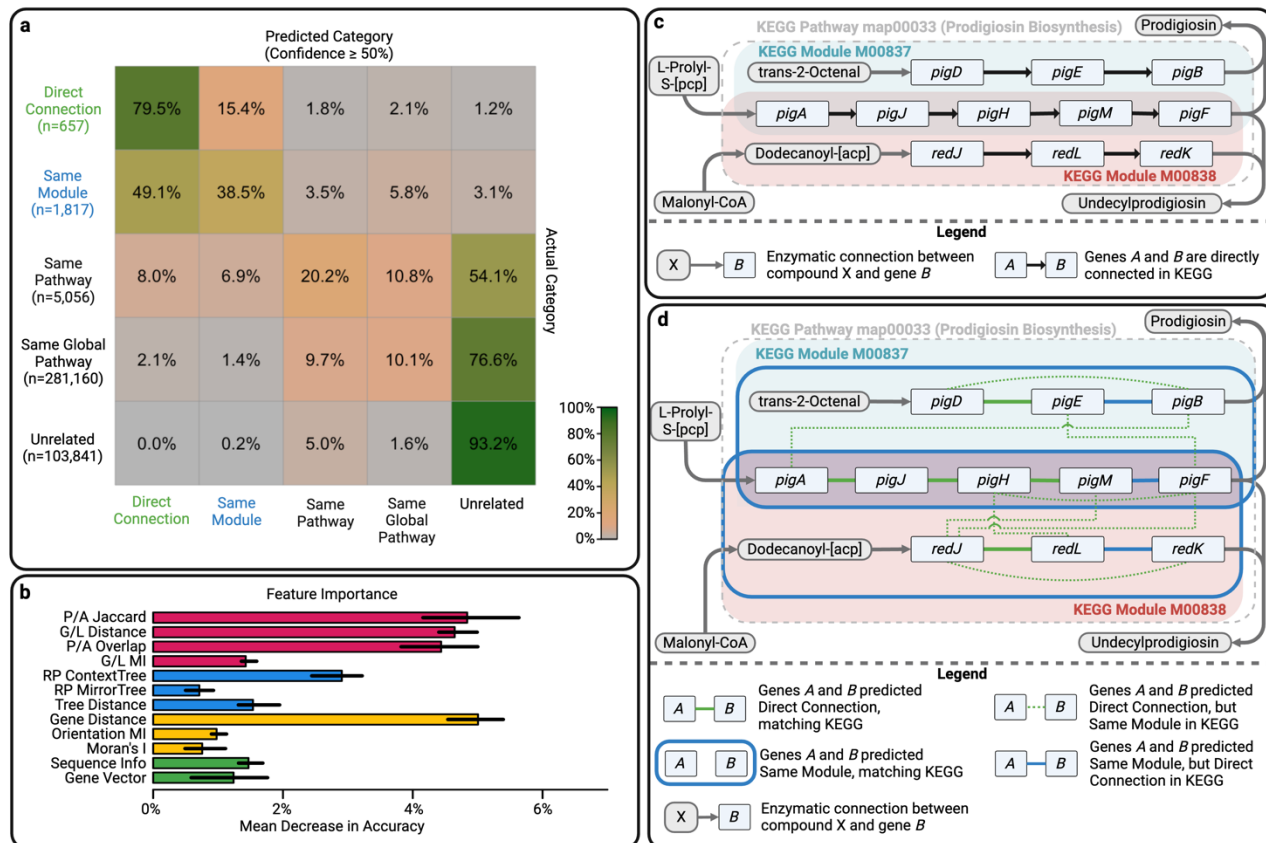835

838
**Figure 1. Overview of the EvoWeaver algorithm and benchmarking. (a)** Phylogenetic trees from groups of orthologous genes serve as the primary input to EvoWeaver. Four categories of coevolutionary signal are quantified for each pair of genes. These signals are combined in an ensemble classifier to predict functional relationships between gene pairs. EvoWeaver provides as output its 12 predictions for signals of coevolution, and can optionally provide an ensemble prediction using built-in pretrained models. **(b)** Functional associations often result in correlated gain/loss patterns on a reference phylogenetic tree (e.g., a species tree). EvoWeaver assesses the presence/absence patterns, correlation between gain/loss events, and distance between gain/loss events as signals of coevolution. **(c)** Similarity in phylogenetic structure is another indicator of coevolution between genes. EvoWeaver computes topological distance as well as correlation in patristic distances following dimensionality reduction using random projection. **(d)** Functionally associated genes sometimes cluster on the genome due to co-regulation or horizontal gene transfer. EvoWeaver derives signals from the conservation in gene orientation and the distance between gene pairs. **(e)** Functional associations sometimes cause concerted changes in sequences that are interrogated by EvoWeaver. EvoWeaver can analyze both nucleotide sequences or amino acid sequences, though nucleotide sequences are pictured here. **(f)** Proteins involved in the same complex are functionally associated and can be identified through signals of coevolution. The goal of the Complexes benchmark is to distinguish orthology groups in the same complex (i.e., positives) from those in different complexes (i.e., negatives). **(g)** Functional associations between proteins that are adjacent in the same module are stronger than those between different modules. The goal of the Modules benchmark is to distinguish adjacent proteins in the same module from independent modules.
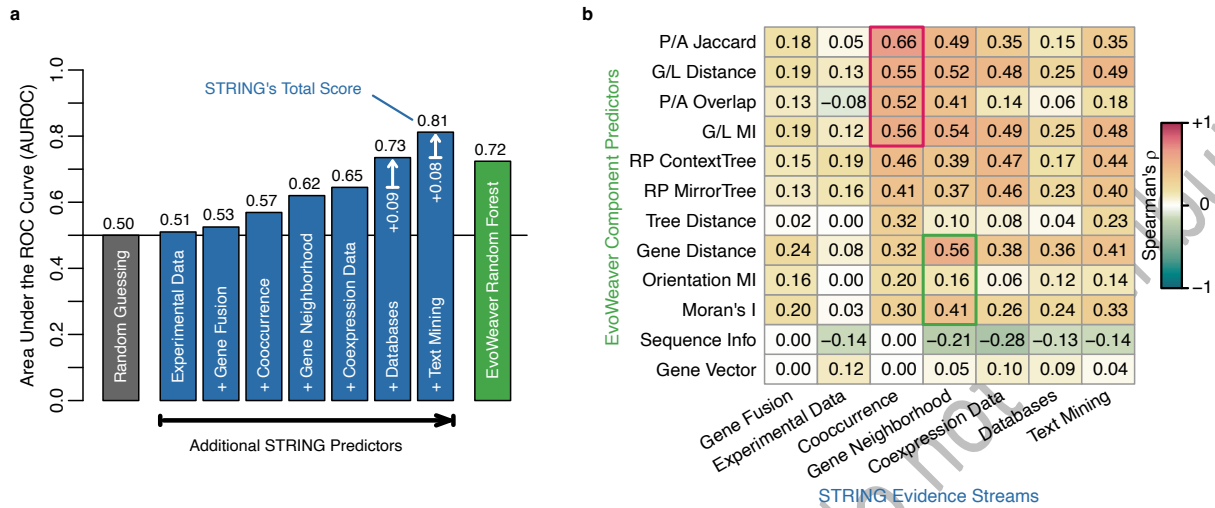
858
859 **Figure 2. EvoWeaver's ensemble predictions outperform individual algorithms on the Modules benchmark.**
860 Coevolutionary approaches were compared for their ability to discern adjacent proteins in KEGG modules (i.e., 899
861 positives) from proteins in distinct modules (i.e., 899 negatives). No single source of coevolutionary signal greatly
862 outcompeted all other sources. However, EvoWeaver's ensemble predictions that combine all component sources of
863 coevolutionary signal substantially improved predictive accuracy, as seen by larger areas under the curves. Inset of
864 the receiver operating characteristic highlights the region with low false positive rates. Scores from individual
865 algorithms tended to have low correlation except within similar categories of coevolutionary signal (i.e., boxed groups
866 in the heatmap), suggesting that the ensemble approach is superior because it combines semi-orthogonal
867 coevolutionary signals. Spearman's correlation from positive and negative sets is averaged to correct for artificial
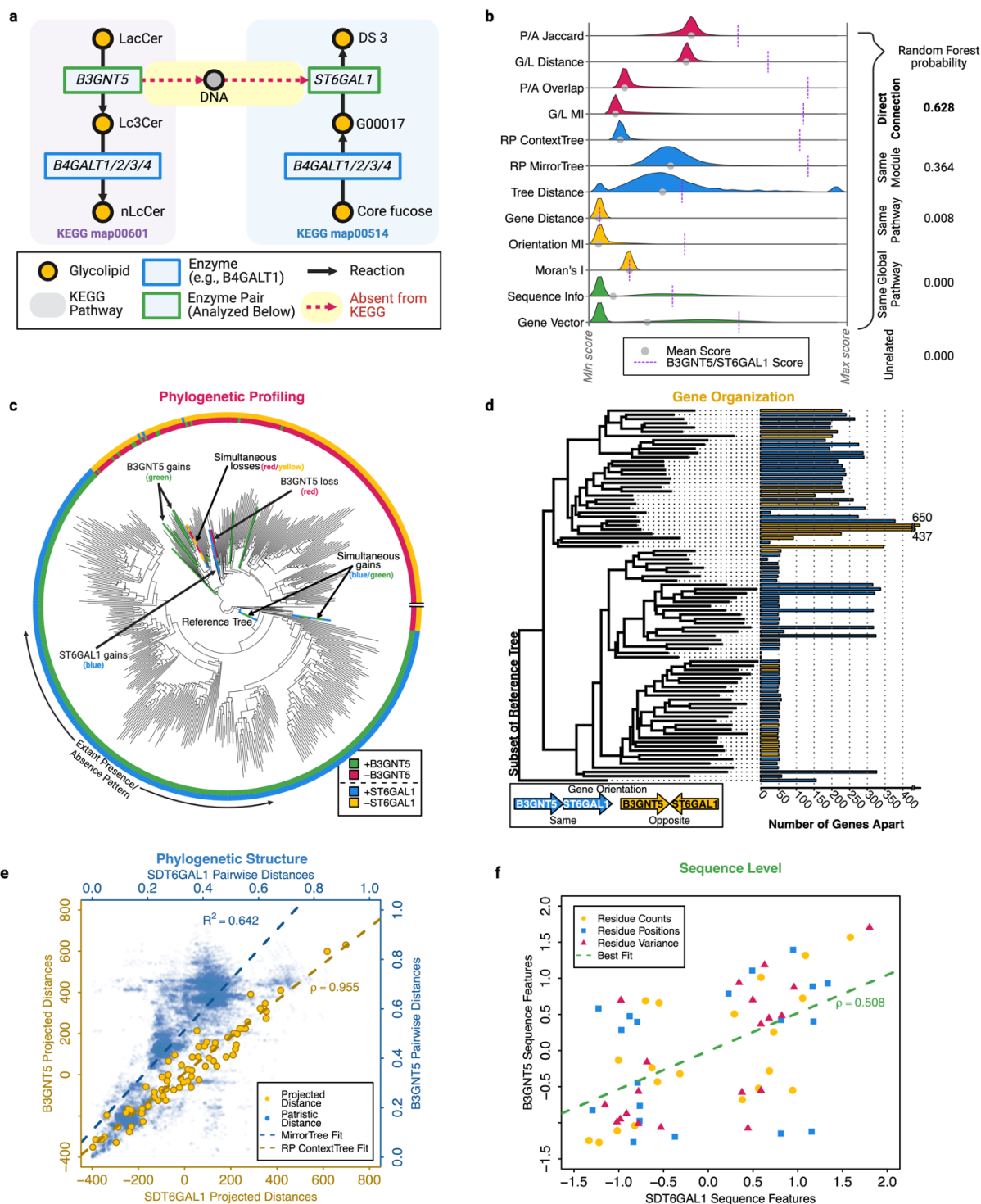868 correlation among high performing algorithms.

869

870 **Figure 3. EvoWeaver is sufficiently accurate to hierarchically classify functional associations. (a)** The
871 confusion matrix of five level classifications indicates that EvoWeaver's ensemble predictions (i.e., Random Forest)
872 rarely confuse proteins within the same module with those from different modules. Values represent the percent of
873 each actual class classified to each predicted class. **(b)** The best performing algorithm from each category on the
874 Modules benchmark also was assigned greater feature importance by the random forest model in hierarchical
875 classification. All features were important in the ensemble's predictions, further underscoring the benefit of using
876 multiple coevolutionary signals. Error bars denote the range of importances across all train/test folds. **(c)** A group of
877 proteins randomly selected from hierarchical clustering exactly matches an existing tightly linked set of modules from
878 KEGG. **(d)** EvoWeaver's ensemble predictions for genes involved in prodigiosin biosynthesis generally match
879 experimentally verified connections in KEGG. Note that *pigA*, *pigJ*, *pigH*, *pigM*, and *pigF* belong to both modules.
880

**a**

**b**

881
882
883 **Figure 4. EvoWeaver rivals STRING without reliance on external data. (a)** Predictive accuracy was compared on
884 1,514 pairs of gene groups that overlapped between STRING and the Multiclass benchmark. Area under the ROC
885 curve (AUROC) is shown for discerning between pairs sharing the same pathway in KEGG (i.e., positives) versus
886 pairs in different pathways (i.e., negatives). STRING's predictions are a composite of seven evidence streams,
887 including three coevolutionary evidence streams (i.e., Gene Fusion, Cooccurrence, Gene Neighborhood).
888 Sequentially incorporating evidence streams from least to most beneficial demonstrates their marginal impact on
889 STRING's reported Total Score. Text Mining and Databases were the most impactful STRING evidence streams.
890 Despite STRING's predictions incorporating KEGG into its Databases evidence stream, EvoWeaver's Random Forest
891 predictions roughly match STRING's predictions without Text Mining while only using sequence information. **(b)** As
892 expected, some of EvoWeaver's component predictors were modestly correlated with STRING's evidence streams.
893 For example, STRING's Cooccurrence score is correlated with EvoWeaver's Phylogenetic Profiling algorithms (red
894 box), and STRING's Gene Neighborhood score is correlated with EvoWeaver's Gene Organization algorithms (green
895 box). Spearman's correlation is calculated in the same manner as in Figure 2.
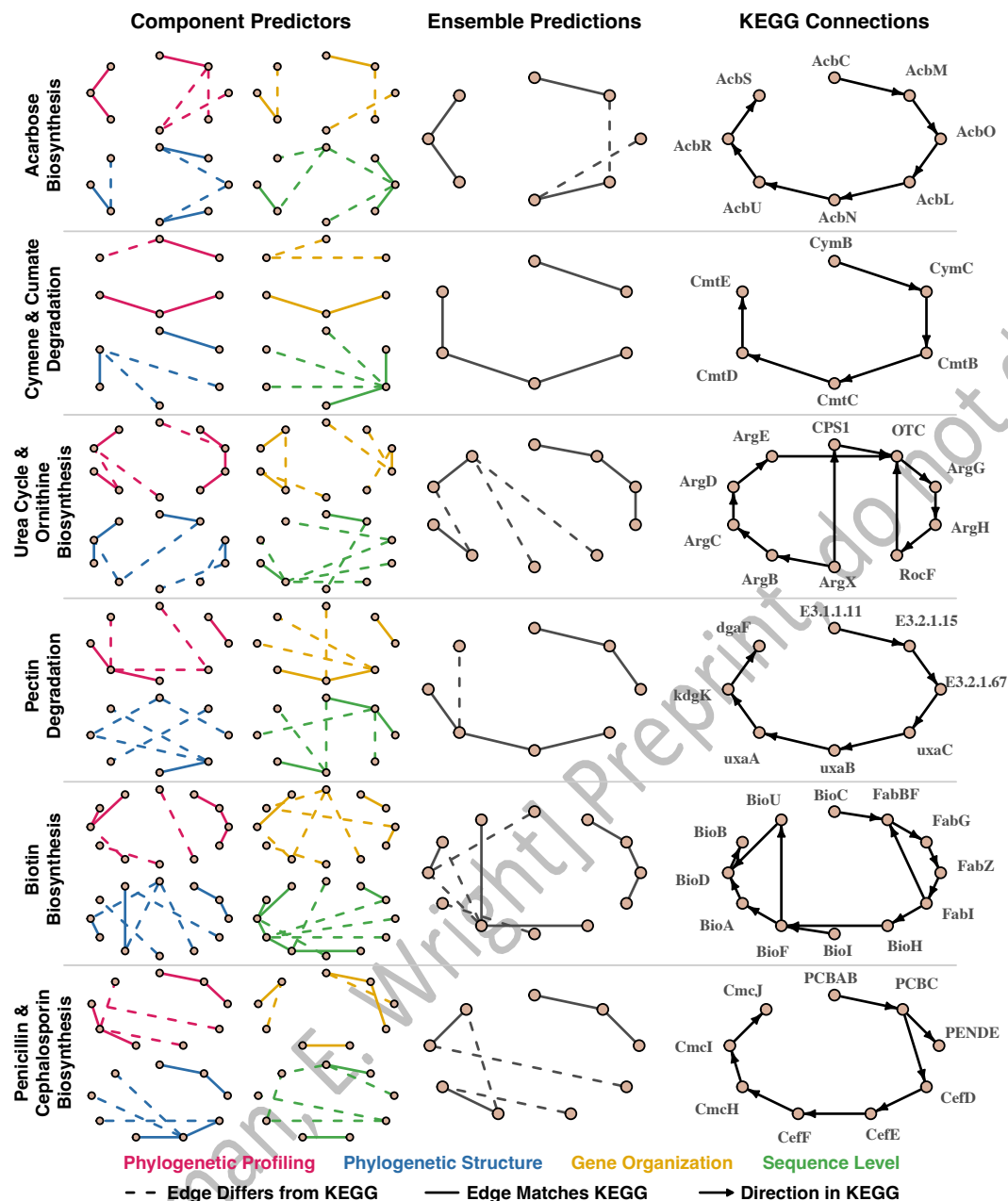
**Figure 5. EvoWeaver's ensemble predictions can generate high fidelity biological hypotheses. (a)** The protein product of *B3GNT5* promotes the expression of *ST6GAL1*[65], although this connection is missing in KEGG and STRING. **(b)** EvoWeaver's component and ensemble predictions indicate that *B3GNT5* and *ST6GAL1* are functionally associated, which is supported by experiments in human cell culture[65]. **(c)** Phylogenetic Profiling demonstrates a pattern of association between *B3GNT5* and *ST6GAL1*, although it is supported by relatively few gain/loss events on the reference tree. **(d)** Organisms with both *B3GNT5* and *ST6GAL1* on the same chromosome

903   display correlations in gene orientation and modest signal of colocalization. **(e)** Shared patristic distances from both
904   gene trees are correlated, especially after compression with random projection, suggesting a high degree of
905   coevolution between *B3GNT5* and *ST6GAL1*. **(f)** Gene sequence natural vectors for both *B3GNT5* and *ST6GAL1* are
906   moderately correlated, implying similar residue compositions and providing further signal of coevolution.
907

**Figure 6.** EvoWeaver partly identifies biochemical pathway connectivity. EvoWeaver's pairwise scores from component algorithms provide a ranking of functional association drawn from alternative categories of coevolutionary signal (colors). EvoWeaver combines 12 component scores into a single ensemble prediction for each pair of gene groups. The strongest predicted connection for every gene group shows high consistency (solid lines) with the actual connectivity of KEGG pathways (arrows). Discrepancies (dashed lines) between predicted and actual connections are often caused by EvoWeaver incorrectly linking consecutive gene groups, such as AcbO-AcbL-AcbN in Acarbose Biosynthesis or ArgB-ArgC-ArgD in Ornithine Biosynthesis. Component predictors are connected according to the gene group with the highest mean rank among all algorithms in a category. Ensemble predictions are determined by connecting each gene group to the gene group with the highest probability of "Direct Connection" in the Multiclass Random Forest model.