

Graph embedding and transfer learning can help predict potential species interaction networks despite data limitations

Tanya Strydom^{1,2,‡} Salomé Bouskila¹ Francis Banville^{1,2,3} Ceres Barros⁴ Dominique Caron^{2,5}
Maxwell J Farrell⁶ Marie-Josée Fortin⁶ Victoria Hemming⁷ Benjamin Mercier^{2,3} Laura
J. Pollock^{2,5} Rogini Runghen⁸ Giulio V. Dalla Riva⁹ Timothée Poisot^{1,2,‡}

¹ Département de Sciences Biologiques, Université de Montréal, Montréal, Canada ² Quebec Centre for Biodiversity Science, Montréal, Canada ³ Département de Biologie, Université de Sherbrooke, Sherbrooke, Canada ⁴ Department of Forest Resources Management, University of British Columbia, Vancouver, B.C., Canada ⁵ Department of Biology, McGill University, Montréal, Canada ⁶ Department of Ecology & Evolutionary Biology, University of Toronto, Toronto, Canada ⁷ Department of Forest and Conservation Sciences, University of British Columbia, Vancouver, Canada ⁸ Centre for Integrative Ecology, School of Biological Sciences, University of Canterbury, Canterbury, New Zealand ⁹ School of Mathematics and Statistics, University of Canterbury, Canterbury, New Zealand

‡ Equal contributions

Correspondance to:

Timothée Poisot — timothee.poisot@umontreal.ca

1. Metawebs (networks of potential interactions within a species pool) are a powerful abstraction to understand how large-scale species interaction networks are structured.
2. Because metawebs are typically expressed at large spatial and taxonomic scales, assembling them is a tedious and costly process; predictive methods can help circumvent the limitations in data deficiencies, by providing a first approximation of metawebs.
3. One way to improve our ability to predict metawebs is to maximize available information by using graph embeddings, as opposed to an exhaustive list of species interactions. Graph embedding is an emerging field in machine learning that holds great potential for ecological problems.

4. Here, we outline how the challenges associated with inferring metawebs line-up with the advantages of graph embeddings; followed by a discussion as to how the choice of the species pool has consequences on the reconstructed network, specifically as to the role of human-made (or arbitrarily assigned) boundaries and how these may influence ecological hypotheses.

1 The ability to infer *potential* interactions could serve as a significant breakthrough in our ability to
2 conceptualize species interaction networks over large spatial scales (Hortal et al., 2015). Reliable
3 inferences would not only boost our understanding of the structure of species interaction networks, but
4 also increase the amount of information that can be used for biodiversity management. In a recent
5 overview of the field of ecological network prediction, Strydom, Catchen, et al. (2021) identified two
6 challenges of interest to the prediction of interactions at large scales. First, there is a relative scarcity of
7 relevant data in most places globally – restricting the inference of interactions to locations where least
8 required (and leaves us unable to make inference in data scarce regions); second, accurate predictors are
9 important for accurate predictions, and the lack of methods that can leverage a small amount of *accurate*
10 data is a serious impediment to our predictive ability. In most places, our most reliable biodiversity
11 knowledge is that of a species pool where a set of potentially interacting species in a given area could
12 occur: through the analysis of databases like the Global Biodiversity Information Facility (GBIF) or the
13 International Union for the Conservation of Nature (IUCN), it is possible to construct a list of species for a
14 region of interest; however inferring the potential interactions between these species still remains a
15 challenge.

16 Following the definition of Dunne (2006), a metaweb is the ecological network analogue to the species
17 pool; specifically, it inventories all *potential* interactions between species for a spatially delimited area (and
18 so captures the γ diversity of interactions). The metaweb itself is not a prediction of local networks at
19 specific locations within the spatial area it covers: it will have a different structure, notably by having a
20 larger connectance (see e.g. Wood et al., 2015) and complexity (see e.g. Galiana et al., 2022), from any of
21 these local networks. These local networks (which capture the α diversity of interactions) are a subset of
22 the metaweb’s species and its realized interactions, and have been called “metaweb realizations” (Poisot et
23 al., 2015). Differences between local networks and their metawebs are due to chance, species abundance
24 and co-occurrence, local environmental conditions, and local distribution of functional traits, among
25 others. Specifically, although co-occurrence can be driven by interactions (Cazelles et al., 2016),
26 co-occurrence alone is not a predictor of interactions (Blanchet et al., 2020; Thurman et al., 2019), and
27 therefore the lack of co-occurrence cannot be used to infer the lack of a feasible interaction. Yet, recent
28 results by Saravia et al. (2021) strongly suggested that local (metaweb) realizations only respond weakly to
29 local conditions: instead, they reflect constraints inherited by the structure of their metaweb. This sets up
30 the core goal of predictive network ecology as the prediction of metaweb structure, as it is required to

accurately produce downscaled, local predictions.

Because the metaweb represents the joint effect of functional, phylogenetic, and macroecological processes (Morales-Castilla et al., 2015), it holds valuable ecological information. Specifically, it represents the “upper bounds” on what the composition of the local networks, given a local species pool, can be (see e.g. McLeod et al., 2021); this information can help evaluate the ability of ecological assemblages to withstand the effects of, for example, climate change (Fricke et al., 2022). These local networks may be reconstructed given an appropriate knowledge of local species composition and provide information on the structure of food webs at finer spatial scales. This has been done for example for tree-galler-parasitoid systems (Gravel et al., 2018), fish trophic interactions (Albouy et al., 2019), tetrapod trophic interactions (J. Braga et al., 2019; O’Connor et al., 2020), and crop-pest networks (Grünig et al., 2020). In this contribution, we highlight the power of viewing (and constructing) metawebs as *probabilistic* objects in the context of low-probability interactions, discuss how a family of machine learning tools (graph embeddings and transfer learning) can be used to overcome data limitations to metaweb inference, and highlight how the use of metawebs introduces important questions for the field of network ecology.

A metaweb is an inherently probabilistic object

Treating interactions as probabilistic (as opposed to binary) events is a more nuanced and realistic way to represent them. Dallas et al. (2017) suggested that most interactions (links) in ecological networks are cryptic, *i.e.* uncommon or hard to observe. This argument echoes Jordano (2016): sampling ecological interactions is difficult because it requires first the joint observation of two species, and then the observation of their interaction. In addition, it is generally expected that weak or rare interactions will be more prevalent in networks than common or strong interactions (Csermely, 2004), compared to strong, persistent interactions; this is notably the case in food chains, wherein many weaker interactions are key to the stability of a system (Neutel et al., 2002). In the light of these observations, we expect to see an over-representation of low-probability (hereafter rare) interactions under a model that accurately predicts interaction probabilities. Although defining an interaction as ‘rare’ is perhaps not as straight forward as one may assume (rare in the context of likelihood of occurrence or the context of biologically plausible?) but for the context of the discussion in this manuscript the core idea is that construction of probabilistic networks affords us the nuance of capturing this inherent variability of interaction occurrence.

69 Critically, the original metaweb definition, and indeed most past uses of metawebs, was based on the
60 presence/absence of interactions. Moving towards *probabilistic* metawebs, by representing interactions as
61 Bernoulli events (see *e.g.* Poisot et al., 2016), offers the opportunity to weigh these rare interactions
62 appropriately. The inherent plasticity of interactions is important to capture: there have been documented
63 instances of food webs undergoing rapid collapse/recovery cycles over short periods of time (*e.g.* Pedersen
64 et al., 2017). Furthermore, because the structure of the metaweb cannot be known in advance, it is
65 important to rely on predictive tools that do not assume a specific network topology for link prediction
66 (Gaucher et al., 2021), but are able to work on generalizations of the network. These considerations
67 emphasize why metaweb predictions should focus on quantitative (preferentially probabilistic)
68 predictions, and this should constrain the suite of models that are appropriate for prediction.

69 It is important to recall that a metaweb is intended as a catalogue of all potential (feasible) interactions,
70 which is then filtered for a given application (Morales-Castilla et al., 2015). It is therefore important to
71 separate the interactions that happen “almost surely” (repeated observational data), “almost never”
72 (repeated lack of evidence *or* evidence that the link is forbidden through *e.g.* trait mis-match), and
73 interactions with a probability that lays somewhere in between (Catchen et al., 2023). In a sense, that most
74 ecological interactions are elusive can call for a slightly different approach to sampling: once the common
75 interactions are documented, the effort required in documenting each rare interaction will increase
76 exponentially (Jordano, 2016). Recent proposals in other fields relying on machine learning approaches
77 emphasize the idea that algorithms meant to predict, through the assumption that they approximate the
78 process generating the data, can also act as data generators (Hoffmann et al., 2019). High quality
79 observational data can be used to infer core rules underpinning network structure, and be supplemented
80 with synthetic data coming from predictive models trained on them, thereby increasing the volume of
81 information available for analysis. Indeed, Strydom, Catchen, et al. (2021) suggested that knowing the
82 metaweb may render the prediction of local networks easier, because it fixes an “upper bound” on which
83 interactions can exist. In this context, a probabilistic metaweb represents an aggregation of informative
84 priors on the biological feasibility of interactions, which is usually hard to obtain yet has possibly the most
85 potential to boost our predictive ability of local networks (Bartomeus, 2013; Bartomeus et al., 2016). This
86 would represent a departure from simple rules expressed at the network scale (*e.g.* Williams & Martinez,
87 2000) to a view of network prediction based on learning the rules that underpin interactions *and* their
88 variability (Gupta et al., 2022).

Graph embedding offers promises for the inference of potential interactions

Graph (or network) embedding (fig. 1) is a family of machine learning techniques, whose main task is to learn a mapping function from a discrete graph to a continuous domain (Arsov & Mirceva, 2019; Chami et al., 2022). Their main goal is to learn a low dimensional vector representations of the graph (embeddings), such that its key properties (*e.g.* local or global structures) are retained in the embedding space (Yan et al., 2005). The embedding space may, but will not necessarily, have lower dimensionality than the graph. Ecological networks are promising candidates for the routine application of embeddings, as they tend to possess a shared structural backbone (see *e.g.* Bramon Mora et al., 2018), which hints at structural invariants in empirical data. Assuming that these structural invariants are common enough, they would dominate the structure of networks, and therefore be adequately captured by the first (lower) dimensions of an embedding, without the need to measure derived aspects of their structure (*e.g.* motifs, paths, modularity, ...).

Graph embedding produces latent variables (but not traits)

Before moving further, it is important to clarify the epistemic status of node values derived from embeddings: specifically, they are *not* functional traits, and therefore should not be interpreted in terms of effects or responses. As per the framework of Malaterre et al. (2019), these values neither derive from, nor result in, changes in organismal performance, and should therefore not be used to quantify *e.g.* functional diversity. This holds true even when there are correlations between latent values and functional traits: although these enable an ecological discussion of how traits condition the structure of the network, the existence of a statistical relationship does not elevate the latent values to the status of functional traits. Rather than directly predicting biological rules (see *e.g.* Pichler et al., 2020 for an overview), which may be confounded by the sparse nature of graph data, learning embeddings works in the low-dimensional space that maximizes information about the network structure. This approach is further justified by the observation, for example, that the macro-evolutionary history of a network is adequately represented by

115 some graph embeddings [Random dot product graphs (RDPG); see Dalla Riva & Stouffer (2016)]. In a
116 recent publication, Strydom et al. (2022) have used an embedding (based on RDPG) to project a metaweb
117 of trophic interactions between European mammals, and transferred this information to mammals of
118 Canada, using the phylogenetic distance between related clades to infer the values in the latent subspace
119 into which the European metaweb was projected. By performing the RDPG step on re-constructed values,
120 this approach yields a probabilistic trophic metaweb for mammals of Canada based on knowledge of
121 European species, despite a limited ($\approx 5\%$) taxonomic overlap, and illustrates how the values derived from
122 an embedding can be used for prediction without being “traits” of the species they represent.

123 **Ecological networks are good candidates for embedding**

124 Food webs are inherently low-dimensional objects, and can be adequately represented with less than ten
125 dimensions (J. Braga et al., 2019; M. P. Braga et al., 2021; Eklöf et al., 2013). Simulation results by Botella
126 et al. (2022) suggested that there is no dominant method to identify architectural similarities between
127 networks: multiple approaches need to be tested and compared to the network descriptor of interest on a
128 problem-specific basis. This matches previous results on graph embedding, wherein different embedding
129 algorithms yield different network embeddings (Goyal & Ferrara, 2018), calling for a careful selection of
130 the problem-specific approach to use. In tbl. 1, we present a selection of common graph and node
131 embedding methods, alongside examples of their use to predict interactions or statistical associations
132 between species. These methods rely largely on linear algebra or pseudo-random walks on graphs. All
133 forms of embeddings presented in tbl. 1 share the common property of summarizing their objects into
134 (sets of) dense feature vectors, that capture the overall network structure, pairwise information on nodes,
135 and emergent aspects of the network, in a compressed way (*i.e.* with some information loss, as we later
136 discuss in the illustration). Node embeddings tend to focus on maintaining pairwise relationships (*i.e.*
137 species interactions), while graph embeddings focus on maintaining the network structure (*i.e.* emergent
138 properties). Nevertheless, some graph embedding techniques (like RDPG, see *e.g.* Wu et al., 2021) will
139 provide high-quality node-level embeddings while also preserving network structure.

140 Graph embeddings *can* serve as a dimensionality reduction method. For example, RDPG (Strydom et al.,
141 2022) and t-SVD [truncated Singular Value Decomposition; Poisot et al. (2021)] typically embed networks
142 using fewer dimensions than the original network [the original network has as many dimensions as
143 species, and as many informative dimensions as trophically unique species; Strydom, Dalla Riva, et al.

(2021)]. However, this is not necessarily the case – indeed, one may perform a PCA (a special case of SVD) to project the raw data into a subspace that improves the efficacy of t-SNE [t-distributed stochastic neighbor embedding; Maaten (2009)]. There are many dimensionality reductions (Anowar et al., 2021) that can be applied to an embedded network should the need for dimensionality reduction (for example for data visualization) arise. In brief, many graph embeddings *can* serve as dimensionality reduction steps, but not all do, neither do all dimensionality reduction methods provide adequate graph embedding capacities. In the next section (and fig. 1), we show how the amount of dimensionality reduction can affect the quality of the embedding.

Graph embedding has been under-used in the prediction of species interactions

One prominent family of approaches we do not discuss in the present manuscript is Graph Neural Networks [GNN; Zhou et al. (2020)]. GNN are, in a sense, a method to embed a graph into a dense subspace, but belong to the family of deep learning methods, which has its own set of practices (see *e.g.* Goodfellow et al., 2016). An important issue with methods based on deep learning is that, because their parameter space is immense, the sample size of the data fed into them must be similarly large (typically thousands of instances). This is a requirement for the model to converge correctly during training, but this assumption is unlikely to be met given the size of datasets currently available for metawebs (or single time/location species interaction networks). This data volume requirement is mostly absent from the techniques we list below. Furthermore, GNN still have some challenges related to their shallow structure, and concerns related to scalability (see Gupta et al., 2021 for a review), which are mostly absent from the methods listed in tbl. 1. Assuming that the uptake of next-generation biomonitoring techniques does indeed deliver larger datasets on species interactions (Bohan et al., 2017), there is nevertheless the potential for GNN to become an applicable embedding/predictive technique in the coming years.

Table 1: Overview of some common graph embedding approaches, by type of embedded objects, alongside examples of their use in the prediction of species interactions. These methods have not yet been routinely used to predict species interactions; most examples that we identified were either statistical associations, or analogues to joint species distribution models. ^a: application is concerned with *statistical* interactions, which are not necessarily direct biotic interactions; ^b: application is concerned with joint-SDM-like approach, which is also very close to statistical associations as opposed to direct biotic interactions. Given the need to evaluate different methods on a problem-specific basis, the fact that a lot of methods have not been used on network problems is an opportunity for benchmarking and method development. Note that the row for PCA also applies to kernel/probabilistic PCA, which are variations on the more general method of SVD. Note further that tSNE has been included because it is frequently used to embed graphs, including of species associations/interactions, despite not being strictly speaking, a graph embedding technique (see e.g. Chami et al., 2022).

Method	Object	Technique	Reference	Application
tsNE	nodes	statistical divergence	Hinton & Roweis (2002)	(Cieslak et al., 2020, species-environment responses ^a) (Gibb et al., 2021, host-virus network representation)
LINE	nodes	stochastic gradient descent	Tang et al. (2015)	
SDNE	nodes	gradient descent	D. Wang et al. (2016)	
node2vec	nodes	stochastic gradient descent	Grover & Leskovec (2016)	
HARP	nodes	meta-strategy	H. Chen et al. (2017)	
DMSE	joint nodes	deep neural network	D. Chen et al. (2017)	(D. Chen et al., 2017, species-environment interactions ^b)
graph2vec	sub- graph	skipgram network	Narayanan et al. (2017)	
RDPG	graph	SVD	Young & Scheinerman (2007)	(Dalla Riva & Stouffer, 2016, trophic interactions) (Poisot et al., 2021, host-virus network prediction)

Method	Object	Technique	Reference	Application
GLEE	graph	Laplacian eigenmap	Torres et al. (2020)	
DeepWalk	graph	stochastic gradient descent	Perozzi et al. (2014)	(Wardeh et al., 2021, host-virus interactions)
GraphKGE	graph	stochastic differential equation	Melnyk et al. (2020)	(Melnyk et al., 2020, microbiome species associations a)
FastEmbed	graph	eigen decomposition	Ramasamy & Madhow (2015)	
PCA	graph	eigen decomposition	Surendran (2013)	(Strydom, Catchen, et al., 2021, host-parasite interactions)
Joint methods	multiple graphs	multiple strategies	S. Wang et al. (2021)	

The popularity of graph embedding techniques in machine learning is more than the search for structural invariants: graphs are discrete objects, and machine learning techniques tend to handle continuous data better. Bringing a sparse graph into a continuous, dense vector space (Xu, 2021) opens up a broader variety of predictive algorithms, notably of the sort that are able to predict events as probabilities (Murphy, 2022). Furthermore, the projection of the graph itself is a representation that can be learned; Runghen et al. (2021), for example, used a neural network to learn the embedding of a network in which not all interactions were known, based on the nodes' metadata. This example has many parallels in ecology (see fig. 1 C), in which node metadata can be represented by phylogeny, abundance, or functional traits. Using phylogeny as a source of information assumes (or strives to capture) the action of evolutionary processes on network structure, which at least for food webs have been well documented (M. P. Braga et al., 2021; Dalla Riva & Stouffer, 2016; Eklöf & Stouffer, 2016; Stouffer et al., 2012; Stouffer et al., 2007); similarly, the use of functional traits assumes that interactions can be inferred from the knowledge of trait-matching rules, which is similarly well supported in the empirical literature (Bartomeus, 2013; Bartomeus et al.,

179 2016; Goebel et al., 2023; Gravel et al., 2013). Relating this information to an embedding rather than a list
180 of network measures would allow to capture their effect on the more fundamental aspects of network
181 structure; conversely, the absence of a phylogenetic or functional signal may suggest that
182 evolutionary/trait processes are not strong drivers of network structure, therefore opening a new way to
183 perform hypothesis testing.

184 **An illustration of metaweb embedding**

185 In this section, we illustrate the embedding of a collection of bipartite networks collected by Hadfield et al.
186 (2014), using t-SVD and RDPG. Briefly, an RDPG decomposes a network into two subspaces (left and
187 right), which are matrices that when multiplied give an approximation of the original network. RDPG has
188 the particularly desirable properties of being a graph embedding technique that produces relevant
189 node-level feature vectors, and provides good approximations of graphs with varied structures (Athreya et
190 al., 2017). The code to reproduce this example is available as supplementary material (note, for the sake of
191 comparison, that Strydom, Catchen, et al., 2021 have an example using embedding through PCA followed
192 by prediction using a deep neural network on the same dataset). The resulting (binary) metaweb \mathcal{M} has
193 2131 interactions between 206 parasites and 121 hosts, and its adjacency matrix has full rank (*i.e.* it
194 represents a space with 121 dimensions). All analyses were done using Julia (Bezanson et al., 2017)
195 version 1.7.2, *Makie.jl* (Danisch & Krumbiegel, 2021), and *EcologicalNetworks.jl* (Poisot et al., 2019).

196 [Figure 2 about here.]

197 In fig. 2, we focus on some statistical checks of the embedding. In panel **A**, we show that the averaged L_2
198 loss (*i.e.* the sum of squared errors) between the empirical and reconstructed metaweb decreases when the
199 number of dimensions (rank) of the subspace increases, with an inflection at 39 dimensions (out of 120
200 initially) according to the finite differences method. As discussed by Runghen et al. (2021), there is often a
201 trade-off between the number of dimensions to use (more dimensions are more computationally
202 demanding) and the quality of the representation. In panel **B**, we show the increase in cumulative
203 variance explained at each rank, and visualize that using 39 ranks explains about 70% of the variance in
204 the empirical metaweb. This is a different information from the L_2 loss (which is averaged across
205 interactions), as it works on the eigenvalues of the embedding, and therefore captures higher-level

206 features of the network. In panel **C**, we show positions of hosts and parasites on the first two dimensions
 207 of the left and right subspaces. Note that these values largely skew negative, because the first dimensions
 208 capture the coarse structure of the network: most pairs of species do not interact, and therefore have
 209 negative values. Finally in panel **D**, we show the predicted weight (*i.e.* the result of the multiplication of
 210 the RDGP subspaces at a rank of 39) as a function of whether the interactions are observed, not-observed,
 211 or unknown due to lack of co-occurrence in the original dataset. This reveals that the observed
 212 interactions have higher predicted weights, although there is some overlap; the usual approach to identify
 213 potential interactions based on this information would be a thresholding analysis, which is outside the
 214 scope of this manuscript (and is done in the papers cited in this illustration). Because the values returned
 215 from RDGP are not bound to the unit interval, we performed a clamping of the weights to the unit space,
 216 showing a one-inflation in documented interactions, and a zero-inflation in other species pairs. This last
 217 figure crosses from the statistical into the ecological, by showing that species pairs with no documented
 218 co-occurrence have weights that are not distinguishable from species pairs with no documented
 219 interactions, suggesting that (as befits a host-parasite model) the ability to interact is a strong predictor of
 220 co-occurrence.

221 [Figure 3 about here.]

222 The results of fig. 2 show that we can extract an embedding of the metaweb that captures enough variance
 223 to be relevant; specifically, this is true for both L_2 loss (indicating that RDGP is able to capture pairwise
 224 processes) and the cumulative variance explained (indicating that RDGP is able to capture network-level
 225 structure). Therefore, in fig. 3, we relate the values of latent variables for hosts to different
 226 ecologically-relevant data. In panel **A**, we show that host with a higher value on the first dimension have
 227 fewer parasites. This relates to the body size of hosts in the *PanTHERIA* database (Jones et al., 2009), as
 228 shown in panel **B**: interestingly, the position on the first axis is only weakly correlated to body mass of the
 229 host; this matches well established results showing that body size/mass is not always a direct predictor of
 230 parasite richness in terrestrial mammals (Morand & Poulin, 1998), a result we observe in panel **C**. Finally,
 231 in panel **D**, we can see how different taxonomic families occupy different positions on the first axis, with
 232 *e.g.* Sciuridae being biased towards higher values. These results show how we can look for ecological
 233 informations in the output of network embeddings, which can further be refined into the selection of
 234 predictors for transfer learning.

235 **The metaweb merges ecological hypotheses and practices**

236 Metaweb inference seeks to provide information about the interactions between species at a large spatial
237 scale, typically a scale large enough to be considered of biogeographic relevance (indeed, many of the
238 examples covered in the introduction span areas larger than a country, some of them global). But as
239 Herbert (1965) rightfully pointed out, “[y]ou can’t draw neat lines around planet-wide problems”; any
240 inference of a metaweb must therefore contend with several novel, interwoven, families of problems. In
241 this section, we outline three that we think are particularly important, and can discuss how they may
242 addressed with subsequent data analysis or simulations, and how they emerge in the specific context of
243 using embeddings; some of these issues are related to the application of these methods at the
244 science-policy interface.

245 **Identifying the properties of the network to embed**

246 If the initial metaweb is too narrow in scope, notably from a taxonomic point of view, the chances of
247 finding another area with enough related species (through phylogenetic relatedness or similarity of
248 functional traits) to make a reliable inference decreases. This is because transfer requires similarity (fig. 1).
249 A diagnostic for the lack of similar species would likely be large confidence intervals during estimation of
250 the values in the low-rank space. In other words, the representation of the original graph is difficult to
251 transfer to the new problem. Alternatively, if the initial metaweb is too large (taxonomically), then the
252 resulting embeddings would need to represent interactions between taxonomic groups that are not present
253 in the new location. This would lead to a much higher variance in the starting dataset, and to
254 under-dispersion in the target dataset, resulting in the potential under or over estimation of the strength of
255 new predicted interactions. Llewelyn et al. (2022) provided compelling evidence for these situations by
256 showing that, even at small spatial scales, the transfer of information about interactions becomes more
257 challenging when areas rich with endemic species are considered. The lack of well documented metawebs
258 is currently preventing the development of more concrete guidelines. The question of phylogenetic
259 relatedness and distribution is notably relevant if the metaweb is assembled in an area with mostly
260 endemic species (*e.g.* a system that has undergone recent radiation or that has remained in isolation for a
261 long period of time might not have an analogous system with which to draw knowledge from), and as with
262 every predictive algorithm, there is room for the application of our best ecological judgement. Because

263 this problem relates to distribution of species in the geographic or phylogenetic space, it can certainly be
264 approached through assessing the performance of embedding transfer in simulated starting/target species
265 pools.

266 **Identifying the scope of the prediction to perform**

267 The area for which we seek to predict the metaweb should determine the species pool on which the
268 embedding is performed. Metawebs can be constructed by assigning interactions in a list of species within
269 specific regions. The upside of this approach is that information relevant for the construction of this
270 dataset is likely to exist, as countries usually set conservation goals at the national level (Buxton et al.,
271 2021), and as quantitative instruments are consequently designed to work at these scales (Turak et al.,
272 2017); specific strategies are often enacted at smaller scales, nested within a specific country (Ray et al.,
273 2021). However, there is no guarantee that these arbitrary boundaries are meaningful. In fact, we do not
274 have a satisfying answer to the question of “where does an ecological network stop?”, the answer to which
275 would dictate the spatial span to embed/predict. Recent results by Martins et al. (2022) suggested that
276 networks are shaped within eco-regions, with abrupt structural transitions from an eco-region to the next.
277 Should this trend hold generally, this would provide an ecologically-relevant scale at which metawebs can
278 be downscaled and predicted. Other solutions could leverage network-area relationships to identify areas
279 in which networks are structurally similar (see *e.g.* Fortin et al., 2021; Galiana et al., 2022, 2018). Both of
280 these solutions require ample pre-existing information about the network in space. Nevertheless, the
281 inclusion of species for which we have data but that are not in the right spatial extent *may* improve the
282 performance of approaches based on embedding and transfer, *if* they increase the similarity between the
283 target and destination network. This proposal can specifically be evaluated by adding nodes to the
284 network to embed, and assessing the performance of predictive models (see *e.g.* Llewelyn et al., 2022).

285 **Minding legacies shaping ecological datasets**

286 In large parts of the world, boundaries that delineate geographic regions are merely a reflection the legacy
287 of settler colonialism, which drives global disparity in capacity to collect and publish ecological data.
288 Applying any embedding to biased data does not debias them, but rather embeds these biases, propagating
289 them to the models using embeddings to make predictions. Furthermore, the use of ecological data itself is

not an apolitical act (Nost & Goldstein, 2021): data infrastructures tend to be designed to answer questions within national boundaries (therefore placing contingencies on what is available to be embedded), their use often drawing upon, and reinforcing, territorial statecraft (see *e.g.* Barrett, 2005). As per Machen & Nost (2021), these biases are particularly important to consider when knowledge generated algorithmically is used to supplement or replace human decision-making, especially for governance (*e.g.* enacting conservation decisions on the basis of model prediction). As information on networks is increasingly leveraged for conservation actions (see *e.g.* Eero et al., 2021; Naman et al., 2022; Stier et al., 2017), the need to appraise and correct biases that are unwittingly propagated to algorithms when embedded from the original data is immense. These considerations are even more urgent in the specific context of biodiversity data. Long-term colonial legacies still shape taxonomic composition to this day (Lenzner et al., 2022; Raja, 2022), and much shorter-term changes in taxonomic and genetic richness of wildlife emerged through environmental racism (Schmidt & Garroway, 2022). Thus, the set of species found at a specific location is not only as the result of a response to ecological processes separate from human influence, but also the result of human-environment interaction as well as the result legislative/political histories.

Conclusion: metawebs, predictions, and people

Predictive approaches in ecology, regardless of the scale at which they are deployed and the intent of their deployment, originate in the framework that contributed to the ongoing biodiversity crisis (Adam, 2014) and reinforced environmental injustice (Choudry, 2013; Domínguez & Luoma, 2020). The risk of embedding this legacy in our models is real, especially when the impact of this legacy on species pools is being increasingly documented. This problem can be addressed by re-framing the way we interact with models, especially when models are intended to support conservation actions. Particularly on territories that were traditionally stewarded by Indigenous people, we must interrogate how predictive approaches and the biases that underpin them can be put to task in accompanying Indigenous principles of land management (Eichhorn et al., 2019; No'kmaq et al., 2021). The discussion of “algorithm-in-the-loop” approaches that is now pervasive in the machine learning community provides examples of why this is important. Human-algorithm interactions are notoriously difficult and can yield adverse effects (Green & Chen, 2019; Stevenson & Doleac, 2021), suggesting the need to systematically study them for the specific purpose of, here, biodiversity governance. Improving the algorithmic literacy of decision makers is part of

the solution (e.g. Lamba et al., 2019; Mosebo Fernandes et al., 2020), as we can reasonably expect that model outputs will be increasingly used to drive policy decisions (Weiskopf et al., 2022). Our discussion of these approaches need to go beyond the technical and statistical, and into the governance consequences they can have. To embed data also embeds historical and contemporary biases that acted on these data, both because they shaped the ecological processes generating them, and the global processes leading to their measurement and publication. For a domain as vast as species interaction networks, these biases exist at multiple scales along the way, and a challenge for prediction is not only to develop (or adopt) new quantitative tools, but to assess the behavior of these tools in the proper context.

Acknowledgements: We acknowledge that this study was conducted on land within the traditional unceded territory of the Saint Lawrence Iroquoian, Anishinabewaki, Mohawk, Huron-Wendat, and Omàmiwininiwak nations. TP, TS, DC, and LP received funding from the Canadian Institute for Ecology & Evolution. FB is funded by the Institute for Data Valorization (IVADO). TS, SB, and TP are funded by a donation from the Courtois Foundation. CB was awarded a Mitacs Elevate Fellowship no. IT12391, in partnership with fRI Research, and also acknowledges funding from Alberta Innovates and the Forest Resources Improvement Association of Alberta. M-JF acknowledges funding from NSERC Discovery Grant and NSERC CRC. RR is funded by New Zealand's Biological Heritage Ngā Koiora Tuku Iho National Science Challenge, administered by New Zealand Ministry of Business, Innovation, and Employment. BM is funded by the NSERC Alexander Graham Bell Canada Graduate Scholarship and the FRQNT master's scholarship. LP acknowledges funding from NSERC Discovery Grant (NSERC RGPIN-2019-05771). TP acknowledges financial support from the Fondation Courtois, and NSERC through the Discovery Grants and Discovery Accelerator Supplement programs. MJF is supported by an NSERC PDF and an RBC Post-Doctoral Fellowship.

Conflict of interest: The authors have no conflict of interests to disclose

Authors' contributions: TS, and TP conceived the ideas discussed in the manuscript. All authors contributed to writing and editing the manuscript.

Data availability: There is no data associated with this manuscript.

References

- Adam, R. (2014). *Elephant treaties: The Colonial legacy of the biodiversity crisis*. UPNE.
- Albouy, C., Archambault, P., Appeltans, W., Araújo, M. B., Beauchesne, D., Cazelles, K., Cirtwill, A. R., Fortin, M.-J., Galiana, N., Leroux, S. J., Pellissier, L., Poisot, T., Stouffer, D. B., Wood, S. A., & Gravel, D. (2019). The marine fish food web is globally connected. *Nature Ecology & Evolution*, 3(8, 8), 1153–1161. <https://doi.org/10.1038/s41559-019-0950-y>
- Anowar, F., Sadaoui, S., & Selim, B. (2021). Conceptual and empirical comparison of dimensionality reduction algorithms (PCA, KPCA, LDA, MDS, SVD, LLE, ISOMAP, LE, ICA, t-SNE). *Computer Science Review*, 40, 100378. <https://doi.org/10.1016/j.cosrev.2021.100378>
- Arsov, N., & Mirceva, G. (2019). *Network Embedding: An Overview*. <http://arxiv.org/abs/1911.11726>
- Athreya, A., Fishkind, D. E., Levin, K., Lyzinski, V., Park, Y., Qin, Y., Sussman, D. L., Tang, M., Vogelstein, J. T., & Priebe, C. E. (2017). *Statistical inference on random dot product graphs: A survey* (No. arXiv:1709.05454). arXiv. <http://arxiv.org/abs/1709.05454>
- Barrett, S. (2005). *Environment and Statecraft: The Strategy of Environmental Treaty-Making* (1st ed.). Oxford University PressOxford. <https://doi.org/10.1093/0199286094.001.0001>
- Bartomeus, I. (2013). Understanding linkage rules in plant-pollinator networks by using hierarchical models that incorporate pollinator detectability and plant traits. *PloS One*, 8(7), e69200. <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0069200>
- Bartomeus, I., Gravel, D., Tylianakis, J. M., Aizen, M. A., Dickie, I. A., & Bernard-Verdier, M. (2016). A common framework for identifying linkage rules across different types of interactions. *Functional Ecology*, 30(12), 1894–1903. <http://onlinelibrary.wiley.com/doi/10.1111/1365-2435.12666/full>
- Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. B. (2017). Julia: A Fresh Approach to Numerical Computing. *SIAM Review*, 59(1), 65–98. <https://doi.org/10.1137/141000671>
- Blanchet, F. G., Cazelles, K., & Gravel, D. (2020). Co-occurrence is not evidence of ecological interactions. *Ecology Letters*.
- Bohan, D. A., Vacher, C., Tamaddoni-Nezhad, A., Raybould, A., Dumbrell, A. J., & Woodward, G. (2017).

371 Next-Generation Global Biomonitoring: Large-scale, Automated Reconstruction of Ecological
 372 Networks. *Trends in Ecology & Evolution*. <https://doi.org/10.1016/j.tree.2017.03.001>

373 Botella, C., Dray, S., Matias, C., Miele, V., & Thuiller, W. (2022). An appraisal of graph embeddings for
 374 comparing trophic network architectures. *Methods in Ecology and Evolution*, 13(1), 203–216.
 375 <https://doi.org/10.1111/2041-210X.13738>

376 Braga, J., Pollock, L. J., Barros, C., Galiana, N., Montoya, J. M., Gravel, D., Maiorano, L., Montemaggioli,
 377 A., Ficetola, G. F., Dray, S., & Thuiller, W. (2019). Spatial analyses of multi-trophic terrestrial vertebrate
 378 assemblages in Europe. *Global Ecology and Biogeography*, 28(11), 1636–1648.
 379 <https://doi.org/10.1111/geb.12981>

380 Braga, M. P., Janz, N., Nylin, S., Ronquist, F., & Landis, M. J. (2021). Phylogenetic reconstruction of
 381 ancestral ecological networks through time for pierid butterflies and their host plants. *Ecology Letters*,
 382 *n/a*(*n/a*). <https://doi.org/10.1111/ele.13842>

383 Bramon Mora, B., Gravel, D., Gilarranz, L. J., Poisot, T., & Stouffer, D. B. (2018). Identifying a common
 384 backbone of interactions underlying food webs from different ecosystems. *Nature Communications*,
 385 9(1), 2603. <https://doi.org/10.1038/s41467-018-05056-0>

386 Buxton, R. T., Bennett, J. R., Reid, A. J., Shulman, C., Cooke, S. J., Francis, C. M., Nyboer, E. A., Pritchard,
 387 G., Binley, A. D., Avery-Gomm, S., Ban, N. C., Beazley, K. F., Bennett, E., Blight, L. K., Bortolotti, L. E.,
 388 Camfield, A. F., Gadallah, F., Jacob, A. L., Naujokaitis-Lewis, I., ... Smith, P. A. (2021). Key
 389 information needs to move from knowledge to action for biodiversity conservation in Canada.
 390 *Biological Conservation*, 256, 108983. <https://doi.org/10.1016/j.biocon.2021.108983>

391 Catchen, M., Poisot, T., Pollock, L., & Gonzalez, A. (2023). *The missing link: Discerning true from false*
 392 *negatives when sampling species interaction networks* (Preprint No. 4929). EcoEvoRxiv.
 393 <https://doi.org/10.32942/X2DW22>

394 Cazelles, K., Araújo, M. B., Mouquet, N., & Gravel, D. (2016). A theory for species co-occurrence in
 395 interaction networks. *Theoretical Ecology*, 9(1), 39–48.
 396 <https://doi.org/10.1007/s12080-015-0281-9>

397 Chami, I., Abu-El-Haija, S., Perozzi, B., Ré, C., & Murphy, K. (2022). Machine Learning on Graphs: A
 398 Model and Comprehensive Taxonomy. *Journal of Machine Learning Research*, 23(89), 1–64.
 399 <http://jmlr.org/papers/v23/20-852.html>

400 Chen, D., Xue, Y., Fink, D., Chen, S., & Gomes, C. P. (2017). *Deep Multi-species Embedding*. 3639–3646.
 401 <https://www.ijcai.org/proceedings/2017/509>

402 Chen, H., Perozzi, B., Hu, Y., & Skiena, S. (2017). *HARP: Hierarchical Representation Learning for*
 403 *Networks*. <http://arxiv.org/abs/1706.07845>

404 Choudry, A. (2013). Saving biodiversity, for whom and for what? Conservation NGOs, complicity,
 405 colonialism and conquest in an era of capitalist globalization. In *NGOization: Complicity,*
 406 *contradictions and prospects* (pp. 24–44). Bloomsbury Publishing.

407 Cieslak, M. C., Castelfranco, A. M., Roncalli, V., Lenz, P. H., & Hartline, D. K. (2020). T-Distributed
 408 Stochastic Neighbor Embedding (t-SNE): A tool for eco-physiological transcriptomic analysis. *Marine*
 409 *Genomics*, 51, 100723. <https://doi.org/10.1016/j.margen.2019.100723>

410 Csermely, P. (2004). Strong links are important, but weak links stabilize them. *Trends in Biochemical*
 411 *Sciences*, 29(7), 331–334. <https://doi.org/10.1016/j.tibs.2004.05.004>

412 Dalla Riva, G. V., & Stouffer, D. B. (2016). Exploring the evolutionary signature of food webs' backbones
 413 using functional traits. *Oikos*, 125(4), 446–456. <https://doi.org/10.1111/oik.02305>

414 Dallas, T., Park, A. W., & Drake, J. M. (2017). Predicting cryptic links in host-parasite networks. *PLOS*
 415 *Computational Biology*, 13(5), e1005557. <https://doi.org/10.1371/journal.pcbi.1005557>

416 Danisch, S., & Krumbiegel, J. (2021). Makie.jl: Flexible high-performance data visualization for Julia.
 417 *Journal of Open Source Software*, 6(65), 3349. <https://doi.org/10.21105/joss.03349>

418 Domínguez, L., & Luoma, C. (2020). Decolonising Conservation Policy: How Colonial Land and
 419 Conservation Ideologies Persist and Perpetuate Indigenous Injustices at the Expense of the
 420 Environment. *Land*, 9(3, 3), 65. <https://doi.org/10.3390/land9030065>

421 Dunne, J. A. (2006). The Network Structure of Food Webs. In J. A. Dunne & M. Pascual (Eds.), *Ecological*
 422 *networks: Linking structure and dynamics* (pp. 27–86). Oxford University Press.

423 Eero, M., Dierking, J., Humborg, C., Undeman, E., MacKenzie, B. R., Ojaveer, H., Salo, T., & Köster, F. W.
 424 (2021). Use of food web knowledge in environmental conservation and management of living
 425 resources in the Baltic Sea. *ICES Journal of Marine Science*, 78(8), 2645–2663.
 426 <https://doi.org/10.1093/icesjms/fsab145>

427 Eichhorn, M. P., Baker, K., & Griffiths, M. (2019). Steps towards decolonising biogeography. *Frontiers of*
428 *Biogeography*, 12(1), 1–7. <https://doi.org/10.21425/F5FBG44795>

429 Eklöf, A., Jacob, U., Kopp, J., Bosch, J., Castro-Urgal, R., Chacoff, N. P., Dalsgaard, B., de Sassi, C., Galetti,
430 M., Guimarães, P. R., Lomáscolo, S. B., Martín González, A. M., Pizo, M. A., Rader, R., Rodrigo, A.,
431 Tylianakis, J. M., Vázquez, D. P., & Allesina, S. (2013). The dimensionality of ecological networks.
432 *Ecology Letters*, 16(5), 577–583. <https://doi.org/10.1111/ele.12081>

433 Eklöf, A., & Stouffer, D. B. (2016). The phylogenetic component of food web structure and intervality.
434 *Theoretical Ecology*, 9(1), 107–115. <https://doi.org/10.1007/s12080-015-0273-9>

435 Fortin, M.-J., Dale, M. R. T., & Brimacombe, C. (2021). Network ecology in dynamic landscapes.
436 *Proceedings of the Royal Society B: Biological Sciences*, 288(1949), rspb.2020.1889, 20201889.
437 <https://doi.org/10.1098/rspb.2020.1889>

438 Fricke, E. C., Ordonez, A., Rogers, H. S., & Svenning, J.-C. (2022). The effects of defaunation on plants’
439 capacity to track climate change. *Science*.
440 <https://www.science.org/doi/abs/10.1126/science.abk3510>

441 Galiana, N., Lurgi, M., Bastazini, V. A. G., Bosch, J., Cagnolo, L., Cazelles, K., Claramunt-López, B., Emer,
442 C., Fortin, M.-J., Grass, I., Hernández-Castellano, C., Jauker, F., Leroux, S. J., McCann, K., McLeod, A.
443 M., Montoya, D., Mulder, C., Osorio-Canadas, S., Reverté, S., ... Montoya, J. M. (2022). Ecological
444 network complexity scales with area. *Nature Ecology & Evolution*, 1–8.
445 <https://doi.org/10.1038/s41559-021-01644-4>

446 Galiana, N., Lurgi, M., Claramunt-López, B., Fortin, M.-J., Leroux, S., Cazelles, K., Gravel, D., & Montoya,
447 J. M. (2018). The spatial scaling of species interaction networks. *Nature Ecology & Evolution*, 2(5),
448 782–790. <https://doi.org/10.1038/s41559-018-0517-3>

449 Gaucher, S., Klopp, O., & Robin, G. (2021). Outlier detection in networks with missing links.
450 *Computational Statistics & Data Analysis*, 164, 107308.
451 <https://doi.org/10.1016/j.csda.2021.107308>

452 Gibb, R., Albery, G. F., Becker, D. J., Brierley, L., Connor, R., Dallas, T. A., Eskew, E. A., Farrell, M. J.,
453 Rasmussen, A. L., Ryan, S. J., Sweeny, A., Carlson, C. J., & Poisot, T. (2021). Data Proliferation,
454 Reconciliation, and Synthesis in Viral Ecology. *BioScience*, 71(11), 1148–1156.
455 <https://doi.org/10.1093/biosci/biab080>

Goebel, L. G. A., Vitorino, B. D., Frota, A. V. B., & Santos-Filho, M. dos. (2023). Body mass determines the role of mammal species in a frugivore-large fruit interaction network in a Neotropical savanna. *Journal of Tropical Ecology*, 39, e12. <https://doi.org/10.1017/S0266467422000505>

Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.

Goyal, P., & Ferrara, E. (2018). Graph embedding techniques, applications, and performance: A survey. *Knowledge-Based Systems*, 151, 78–94. <https://doi.org/10.1016/j.knosys.2018.03.022>

Gravel, D., Baiser, B., Dunne, J. A., Kopelke, J.-P., Martinez, N. D., Nyman, T., Poisot, T., Stouffer, D. B., Tylianakis, J. M., Wood, S. A., & Roslin, T. (2018). Bringing Elton and Grinnell together: A quantitative framework to represent the biogeography of ecological interaction networks. *Ecography*, 0(0). <https://doi.org/10.1111/ecog.04006>

Gravel, D., Poisot, T., Albouy, C., Velez, L., & Mouillot, D. (2013). Inferring food web structure from predator-prey body size relationships. *Methods in Ecology and Evolution*, 4(11), 1083–1090. <https://doi.org/10.1111/2041-210X.12103>

Green, B., & Chen, Y. (2019). Disparate Interactions: An Algorithm-in-the-Loop Analysis of Fairness in Risk Assessments. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 90–99. <https://doi.org/10.1145/3287560.3287563>

Grover, A., & Leskovec, J. (2016). Node2vec: Scalable Feature Learning for Networks. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 855–864. <https://doi.org/10.1145/2939672.2939754>

Grünig, M., Mazzi, D., Calanca, P., Karger, D. N., & Pellissier, L. (2020). Crop and forest pest metawebs shift towards increased linkage and suitability overlap under climate change. *Communications Biology*, 3(1, 1), 1–10. <https://doi.org/10.1038/s42003-020-0962-9>

Gupta, A., Furrer, R., & Petchey, O. L. (2022). Simultaneously estimating food web connectance and structure with uncertainty. *Ecology and Evolution*, 12(3), e8643. <https://doi.org/10.1002/ece3.8643>

Gupta, A., Matta, P., & Pant, B. (2021). Graph neural network: Current state of Art, challenges and applications. *Materials Today: Proceedings*, 46, 10927–10932. <https://doi.org/10.1016/j.matpr.2021.01.950>

- 484 Hadfield, J. D., Krasnov, B. R., Poulin, R., & Nakagawa, S. (2014). A Tale of Two Phylogenies: Comparative
485 Analyses of Ecological Interactions. *The American Naturalist*, 183(2), 174–187.
486 <https://doi.org/10.1086/674445>
- 487 Herbert, F. (1965). *Dune* (1st ed.). Chilton Book Company.
- 488 Hinton, G., & Roweis, S. T. (2002). Stochastic neighbor embedding. *NIPS*, 15, 833–840.
- 489 Hoffmann, J., Bar-Sinai, Y., Lee, L. M., Andrejevic, J., Mishra, S., Rubinstein, S. M., & Rycroft, C. H. (2019).
490 Machine learning in a data-limited regime: Augmenting experiments with synthetic data uncovers
491 order in crumpled sheets. *Science Advances*, 5(4), eaau6792.
492 <https://doi.org/10.1126/sciadv.aau6792>
- 493 Hortal, J., de Bello, F., Diniz-Filho, J. A. F., Lewinsohn, T. M., Lobo, J. M., & Ladle, R. J. (2015). Seven
494 Shortfalls that Beset Large-Scale Knowledge of Biodiversity. *Annual Review of Ecology, Evolution, and*
495 *Systematics*, 46(1), 523–549. <https://doi.org/10.1146/annurev-ecolsys-112414-054400>
- 496 Jones, K. E., Bielby, J., Cardillo, M., Fritz, S. A., O'Dell, J., Orme, C. D. L., Safi, K., Sechrest, W., Boakes, E.
497 H., Carbone, C., Connolly, C., Cutts, M. J., Foster, J. K., Grenyer, R., Habib, M., Plaster, C. A., Price, S.
498 A., Rigby, E. A., Rist, J., ... Purvis, A. (2009). PanTHERIA: A species-level database of life history,
499 ecology, and geography of extant and recently extinct mammals: Ecological Archives E090-184.
500 *Ecology*, 90(9), 2648–2648. <https://doi.org/10.1890/08-1494.1>
- 501 Jordano, P. (2016). Sampling networks of ecological interactions. *Functional Ecology*, 30(12), 1883–1893.
502 <https://doi.org/10.1111/1365-2435.12763>
- 503 Lamba, A., Cassey, P., Segaran, R. R., & Koh, L. P. (2019). Deep learning for environmental conservation.
504 *Current Biology*, 29(19), R977–R982. <https://doi.org/10.1016/j.cub.2019.08.016>
- 505 Lenzner, B., Latombe, G., Schertler, A., Seebens, H., Yang, Q., Winter, M., Weigelt, P., van Kleunen, M.,
506 Pyšek, P., Pergl, J., Kreft, H., Dawson, W., Dullinger, S., & Essl, F. (2022). Naturalized alien floras still
507 carry the legacy of European colonialism. *Nature Ecology & Evolution*, 1–10.
508 <https://doi.org/10.1038/s41559-022-01865-1>
- 509 Llewelyn, J., Strona, G., Dickman, C. R., Greenville, A. C., Wardle, G. M., Lee, M. S. Y., Doherty, S.,
510 Shabani, F., Saltr  , F., & Bradshaw, C. J. A. (2022). *Predicting predator-prey interactions in terrestrial*
511 *endotherms using random forest* [Preprint]. *Ecology*. <https://doi.org/10.1101/2022.09.02.506446>

512 Maaten, L. van der. (2009). Learning a Parametric Embedding by Preserving Local Structure. *Proceedings*
513 *of the Twelfth International Conference on Artificial Intelligence and Statistics*, 384–391.
514 <https://proceedings.mlr.press/v5/maaten09a.html>

515 Machen, R., & Nost, E. (2021). Thinking algorithmically: The making of hegemonic knowledge in climate
516 governance. *Transactions of the Institute of British Geographers*, 46(3), 555–569.
517 <https://doi.org/10.1111/tran.12441>

518 Malaterre, C., Dussault, A. C., Mermans, E., Barker, G., Beisner, B. E., Bouchard, F., Desjardins, E., Handa,
519 I. T., Kembel, S. W., Lajoie, G., Maris, V., Munson, A. D., Odenbaugh, J., Poisot, T., Shapiro, B. J., &
520 Suttle, C. A. (2019). Functional Diversity: An Epistemic Roadmap. *BioScience*, 69(10), 800–811.
521 <https://doi.org/10.1093/biosci/biz089>

522 Martins, L. P., Stouffer, D. B., Blendinger, P. G., Böhning-Gaese, K., Buitrón-Jurado, G., Correia, M., Costa,
523 J. M., Dehling, D. M., Donatti, C. I., Emer, C., Galetti, M., Heleno, R., Jordano, P., Menezes, Í.,
524 Morante-Filho, J. C., Muñoz, M. C., Neuschulz, E. L., Pizo, M. A., Quitián, M., ... Tylianakis, J. M.
525 (2022). Global and regional ecological boundaries explain abrupt spatial discontinuities in avian
526 frugivory interactions. *Nature Communications*, 13(1, 1), 6943.
527 <https://doi.org/10.1038/s41467-022-34355-w>

528 McLeod, A., Leroux, S. J., Gravel, D., Chu, C., Cirtwill, A. R., Fortin, M.-J., Galiana, N., Poisot, T., & Wood,
529 S. A. (2021). Sampling and asymptotic network properties of spatial multi-trophic networks. *Oikos*,
530 *n/a*(n/a). <https://doi.org/10.1111/oik.08650>

531 Melnyk, K., Klus, S., Montavon, G., & Conrad, T. O. F. (2020). GraphKKE: Graph Kernel Koopman
532 embedding for human microbiome analysis. *Applied Network Science*, 5(1), 96.
533 <https://doi.org/10.1007/s41109-020-00339-2>

534 Morales-Castilla, I., Matias, M. G., Gravel, D., & Araújo, M. B. (2015). Inferring biotic interactions from
535 proxies. *Trends in Ecology & Evolution*, 30(6), 347–356.
536 <https://doi.org/10.1016/j.tree.2015.03.014>

537 Morand, S., & Poulin, R. (1998). Density, body mass and parasite species richness of terrestrial mammals.
538 *Evolutionary Ecology*, 12(6), 717–727. <https://doi.org/10.1023/A:1006537600093>

539 Mosebo Fernandes, A. C., Quintero Gonzalez, R., Lenihan-Clarke, M. A., Leslie Trotter, E. F., & Jokar
540 Arsanjani, J. (2020). Machine Learning for Conservation Planning in a Changing Climate.

541 Sustainability, 12(18, 18), 7657. <https://doi.org/10.3390/su12187657>

542 Murphy, K. P. (2022). *Probabilistic machine learning: An introduction*. MIT Press. probml.ai

543 Naman, S. M., White, S. M., Bellmore, J. R., McHugh, P. A., Kaylor, M. J., Baxter, C. V., Danehy, R. J.,
 544 Naiman, R. J., & Puls, A. L. (2022). Food web perspectives and methods for riverine fish conservation.
 545 *WIREs Water*, n/a(n/a), e1590. <https://doi.org/10.1002/wat2.1590>

546 Narayanan, A., Chandramohan, M., Venkatesan, R., Chen, L., Liu, Y., & Jaiswal, S. (2017). *Graph2vec:*
 547 *Learning Distributed Representations of Graphs*. <http://arxiv.org/abs/1707.05005>

548 Neutel, A.-M., Heesterbeek, J. A. P., & de Ruiter, P. C. (2002). Stability in Real Food Webs: Weak Links in
 549 Long Loops. *Science*, 296(5570), 1120–1123. <https://doi.org/10.1126/science.1068326>

550 No'kmaq, M., Marshall, A., Beazley, K. F., Hum, J., Joudry, shalan, Papadopoulos, A., Pictou, S., Rabesca,
 551 J., Young, L., & Zurba, M. (2021). “Awakening the sleeping giant”: Re-Indigenization principles for
 552 transforming biodiversity conservation in Canada and beyond. *FACETS*, 6(1), 839–869.

553 Nost, E., & Goldstein, J. E. (2021). A political ecology of data. *Environment and Planning E: Nature and*
 554 *Space*, 25148486211043503. <https://doi.org/10.1177/25148486211043503>

555 O'Connor, L. M. J., Pollock, L. J., Braga, J., Ficetola, G. F., Maiorano, L., Martinez-Almoyna, C.,
 556 Montemaggiori, A., Ohlmann, M., & Thuiller, W. (2020). Unveiling the food webs of tetrapods across
 557 Europe through the prism of the Eltonian niche. *Journal of Biogeography*, 47(1), 181–192.
 558 <https://doi.org/10.1111/jbi.13773>

559 Pedersen, E. J., Thompson, P. L., Ball, R. A., Fortin, M.-J., Gouhier, T. C., Link, H., Moritz, C., Nenzen, H.,
 560 Stanley, R. R. E., Taranu, Z. E., Gonzalez, A., Guichard, F., & Pepin, P. (2017). Signatures of the
 561 collapse and incipient recovery of an overexploited marine ecosystem. *Royal Society Open Science*, 4(7),
 562 170215. <https://doi.org/10.1098/rsos.170215>

563 Perozzi, B., Al-Rfou, R., & Skiena, S. (2014). DeepWalk: Online learning of social representations.
 564 *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data*
 565 *Mining*, 701–710. <https://doi.org/10.1145/2623330.2623732>

566 Pichler, M., Boreux, V., Klein, A.-M., Schleuning, M., & Hartig, F. (2020). Machine learning algorithms to
 567 infer trait-matching and predict species interactions in ecological networks. *Methods in Ecology and*
 568 *Evolution*, 11(2), 281–293. <https://doi.org/10.1111/2041-210X.13329>

569 Poisot, T., Belisle, Z., Hoebeke, L., Stock, M., & Szefer, P. (2019). EcologicalNetworks.jl - analysing
570 ecological networks. *Ecography*. <https://doi.org/10.1111/ecog.04310>

571 Poisot, T., Cirtwill, A. R., Cazelles, K., Gravel, D., Fortin, M.-J., & Stouffer, D. B. (2016). The structure of
572 probabilistic networks. *Methods in Ecology and Evolution*, 7(3), 303–312.
573 <https://doi.org/10.1111/2041-210X.12468>

574 Poisot, T., Ouellet, M.-A., Mollentze, N., Farrell, M. J., Becker, D. J., Albery, G. F., Gibb, R. J., Seifert, S. N.,
575 & Carlson, C. J. (2021). *Imputing the mammalian virome with linear filtering and singular value*
576 *decomposition*. <http://arxiv.org/abs/2105.14973>

577 Poisot, T., Stouffer, D. B., & Gravel, D. (2015). Beyond species: Why ecological interaction networks vary
578 through space and time. *Oikos*, 124(3), 243–251. <https://doi.org/10.1111/oik.01719>

579 Raja, N. B. (2022). Colonialism shaped today's biodiversity. *Nature Ecology & Evolution*, 1–2.
580 <https://doi.org/10.1038/s41559-022-01903-y>

581 Ramasamy, D., & Madhow, U. (2015). Compressive spectral embedding: Sidestepping the SVD. In C.
582 Cortes, N. Lawrence, D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in neural information*
583 *processing systems* (Vol. 28). Curran Associates, Inc. [https://](https://proceedings.neurips.cc/paper/2015/file/4f6ffe13a5d75b2d6a3923922b3922e5-Paper.pdf)
584 proceedings.neurips.cc/paper/2015/file/4f6ffe13a5d75b2d6a3923922b3922e5-Paper.pdf

585 Ray, J. C., Grimm, J., & Olive, A. (2021). The biodiversity crisis in Canada: Failures and challenges of
586 federal and sub-national strategic and legal frameworks. *FACETS*, 6, 1044–1068.
587 <https://doi.org/10.1139/facets-2020-0075>

588 Runghen, R., Stouffer, D. B., & Dalla Riva, G. V. (2021). *Exploiting node metadata to predict interactions in*
589 *large networks using graph embedding and neural networks*.
590 <https://doi.org/10.1101/2021.06.10.447991>

591 Saravia, L. A., Marina, T. I., Kristensen, N. P., De Troch, M., & Momo, F. R. (2021). Ecological network
592 assembly: How the regional metaweb influences local food webs. *Journal of Animal Ecology*, n/a(n/a).
593 <https://doi.org/10.1111/1365-2656.13652>

594 Schmidt, C., & Garroway, C. J. (2022). Systemic racism alters wildlife genetic diversity. *Proceedings of the*
595 *National Academy of Sciences*, 119(43), e2102860119. <https://doi.org/10.1073/pnas.2102860119>

596 Stevenson, M. T., & Doleac, J. L. (2021). *Algorithmic Risk Assessment in the Hands of Humans* (SSRN

Scholarly Paper No. 3489440). <https://doi.org/10.2139/ssrn.3489440>

Stier, A. C., Samhouri, J. F., Gray, S., Martone, R. G., Mach, M. E., Halpern, B. S., Kappel, C. V., Scarborough, C., & Levin, P. S. (2017). Integrating Expert Perceptions into Food Web Conservation and Management. *Conservation Letters*, 10(1), 67–76. <https://doi.org/10.1111/conl.12245>

Stouffer, D. B., Camacho, J., Jiang, W., & Nunes Amaral, L. A. (2007). Evidence for the existence of a robust pattern of prey selection in food webs. *Proceedings of the Royal Society B: Biological Sciences*, 274(1621), 1931–1940. <https://doi.org/10.1098/rspb.2007.0571>

Stouffer, D. B., Sales-Pardo, M., Sirer, M. I., & Bascompte, J. (2012). Evolutionary Conservation of Species' Roles in Food Webs. *Science*, 335(6075), 1489–1492. <https://doi.org/10.1126/science.1216556>

Strydom, T., Bouskila, S., Banville, F., Barros, C., Caron, D., Farrell, M. J., Fortin, M.-J., Hemming, V., Mercier, B., Pollock, L. J., Runghen, R., Dalla Riva, G. V., & Poisot, T. (2022). Food web reconstruction through phylogenetic transfer of low-rank network representation. *Methods in Ecology and Evolution*, n/a(n/a). <https://doi.org/10.1111/2041-210X.13835>

Strydom, T., Catchen, M. D., Banville, F., Caron, D., Dansereau, G., Desjardins-Proulx, P., Forero-Muñoz, N. R., Higino, G., Mercier, B., Gonzalez, A., Gravel, D., Pollock, L., & Poisot, T. (2021). A roadmap towards predicting species interaction networks (across space and time). *Philosophical Transactions of the Royal Society B: Biological Sciences*, 376(1837), 20210063. <https://doi.org/10.1098/rstb.2021.0063>

Strydom, T., Dalla Riva, G. V., & Poisot, T. (2021). SVD Entropy Reveals the High Complexity of Ecological Networks. *Frontiers in Ecology and Evolution*, 9. <https://doi.org/10.3389/fevo.2021.623141>

Surendran, S. (2013). Graph Embedding and Dimensionality Reduction - A Survey. *International Journal of Computer Science & Engineering Technology*, 4(1). <https://www.semanticscholar.org/paper/Graph-Embedding-and-Dimensionality-Reduction-A-Surendran/3f413d591e4b2b876e033eeb9390e232ad4826ca>

Tang, J., Qu, M., Wang, M., Zhang, M., Yan, J., & Mei, Q. (2015). LINE: Large-scale Information Network Embedding. *Proceedings of the 24th International Conference on World Wide Web*, 1067–1077. <https://doi.org/10.1145/2736277.2741093>

Thurman, L. L., Barner, A. K., Garcia, T. S., & Chestnut, T. (2019). Testing the link between species

625 interactions and co-occurrence in a trophic network. *Ecography*, 0.
 626 <https://doi.org/10.1111/ecog.04360>

627 Torres, L., Chan, K. S., & Eliassi-Rad, T. (2020). GLEE: Geometric Laplacian Eigenmap Embedding.
 628 *Journal of Complex Networks*, 8(2), cnaa007. <https://doi.org/10.1093/comnet/cnaa007>

629 Turak, E., Brazill-Boast, J., Cooney, T., Drielsma, M., DelaCruz, J., Dunkerley, G., Fernandez, M., Ferrier,
 630 S., Gill, M., Jones, H., Koen, T., Leys, J., McGeoch, M., Mihoub, J.-B., Scanes, P., Schmeller, D., &
 631 Williams, K. (2017). Using the essential biodiversity variables framework to measure biodiversity
 632 change at national scale. *Biological Conservation*, 213, 264–271.
 633 <https://doi.org/10.1016/j.biocon.2016.08.019>

634 Wang, D., Cui, P., & Zhu, W. (2016). Structural Deep Network Embedding. *Proceedings of the 22nd ACM*
 635 *SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1225–1234.
 636 <https://doi.org/10.1145/2939672.2939753>

637 Wang, S., Arroyo, J., Vogelstein, J. T., & Priebe, C. E. (2021). Joint Embedding of Graphs. *IEEE*
 638 *Transactions on Pattern Analysis and Machine Intelligence*, 43(4), 1324–1336.
 639 <https://doi.org/10.1109/TPAMI.2019.2948619>

640 Wardeh, M., Baylis, M., & Blagrove, M. S. C. (2021). Predicting mammalian hosts in which novel
 641 coronaviruses can be generated. *Nature Communications*, 12(1, 1), 780.
 642 <https://doi.org/10.1038/s41467-021-21034-5>

643 Weiskopf, S. R., Harmáčková, Z. V., Johnson, C. G., Londoño-Murcia, M. C., Miller, B. W., Myers, B. J. E.,
 644 Pereira, L., Arce-Plata, M. I., Blanchard, J. L., Ferrier, S., Fulton, E. A., Harfoot, M., Isbell, F., Johnson,
 645 J. A., Mori, A. S., Weng, E., & Rosa, I. M. D. (2022). Increasing the uptake of ecological model results in
 646 policy decisions to improve biodiversity outcomes. *Environmental Modelling & Software*, 149, 105318.
 647 <https://doi.org/10.1016/j.envsoft.2022.105318>

648 Williams, R. J., & Martinez, N. D. (2000). Simple rules yield complex food webs. *Nature*, 404(6774),
 649 180–183. <https://doi.org/10.1038/35004572>

650 Wood, S. A., Russell, R., Hanson, D., Williams, R. J., & Dunne, J. A. (2015). Effects of spatial scale of
 651 sampling on food web structure. *Ecology and Evolution*, 5(17), 3769–3782.
 652 <https://doi.org/10.1002/ece3.1640>

- 653 Wu, D., Palmer, D. R., & Deford, D. R. (2021). *Maximum a Posteriori Inference of Random Dot Product*
654 *Graphs via Conic Programming* (No. arXiv:2101.02180). arXiv. <http://arxiv.org/abs/2101.02180>
- 655 Xu, M. (2021). Understanding Graph Embedding Methods and Their Applications. *SIAM Review*, 63(4),
656 825–853. <https://doi.org/10.1137/20M1386062>
- 657 Yan, S., Xu, D., Zhang, B., & Zhang, H.-J. (2005). Graph embedding: A general framework for
658 dimensionality reduction. *2005 IEEE Computer Society Conference on Computer Vision and Pattern*
659 *Recognition (CVPR'05)*, 2, 830–837 vol. 2. <https://doi.org/10.1109/CVPR.2005.170>
- 660 Young, S. J., & Scheinerman, E. R. (2007). Random Dot Product Graph Models for Social Networks. In A.
661 Bonato & F. R. K. Chung (Eds.), *Algorithms and Models for the Web-Graph* (pp. 138–149). Springer.
662 https://doi.org/10.1007/978-3-540-77004-6_11
- 663 Zhou, J., Cui, G., Hu, S., Zhang, Z., Yang, C., Liu, Z., Wang, L., Li, C., & Sun, M. (2020). Graph neural
664 networks: A review of methods and applications. *AI Open*, 1, 57–81.
665 <https://doi.org/10.1016/j.aiopen.2021.01.001>

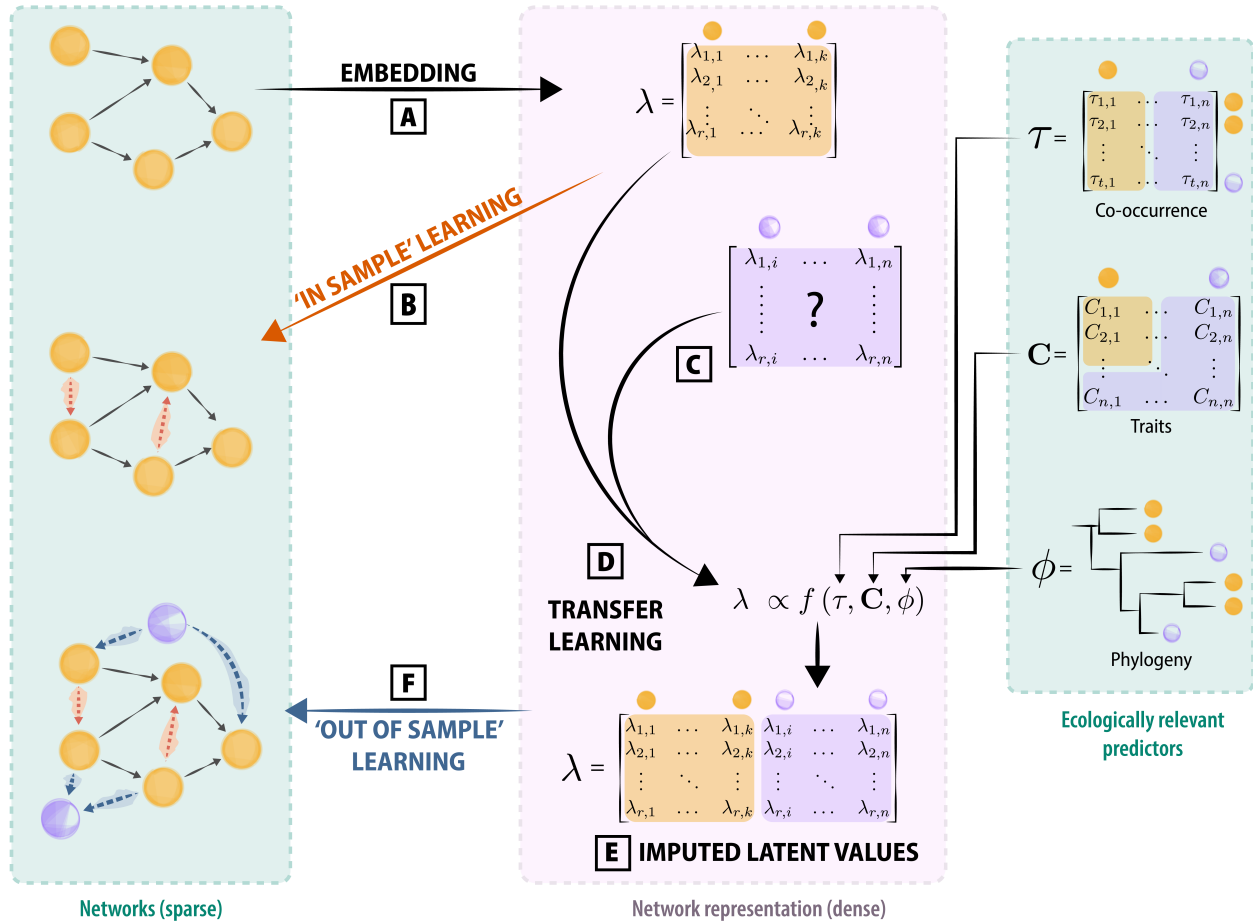


Figure 1: The embedding process (A) can help to identify links (interactions) that may have been missed within the original community (represented by the orange dashed arrows, B). Transfer learning (D) allows for the prediction links (interactions) even when novel species (C) are included alongside the original community. This is achieved by learning using other relevant predictors (e.g. traits) in conjunction with the known interactions to infer latent values (E). Ultimately this allows us to predict links (interactions) for species external from the original sample (blue dashed arrows) as well as missing within sample links (F). Within this context the predicted (and original) networks as well as the ecological predictors used (green boxes) are products that can be quantified through measurements in the field, whereas the embedded as well as imputed matrices (purple box) are representative of a decomposition of the interaction matrices onto the embedding space

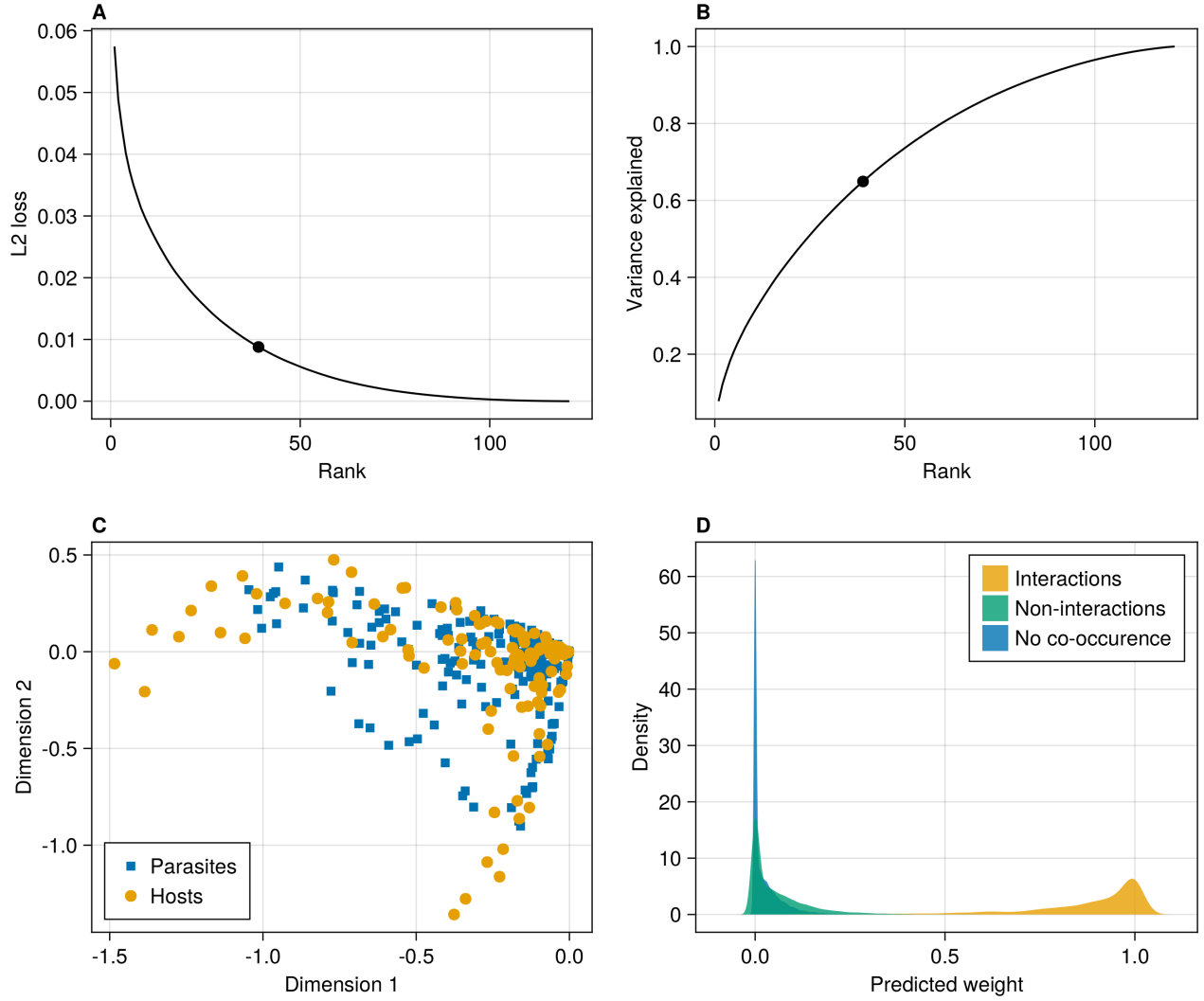


Figure 2: Validation of an embedding for a host-parasite metaweb, using Random Dot Product Graphs. **A**, decrease in approximation error as the number of dimensions in the subspaces increases. **B**, increase in cumulative variance explained as the number of ranks considered increases; in **A** and **B**, the dot represents the point of inflexion in the curve (at rank 39) estimated using the finite differences method. **C**, position of hosts and parasites in the space of latent variables on the first and second dimensions of their respective subspaces (the results have been clamped to the unit interval). **D**, predicted interaction weight from the RDPG based on the status of the species pair in the metaweb.

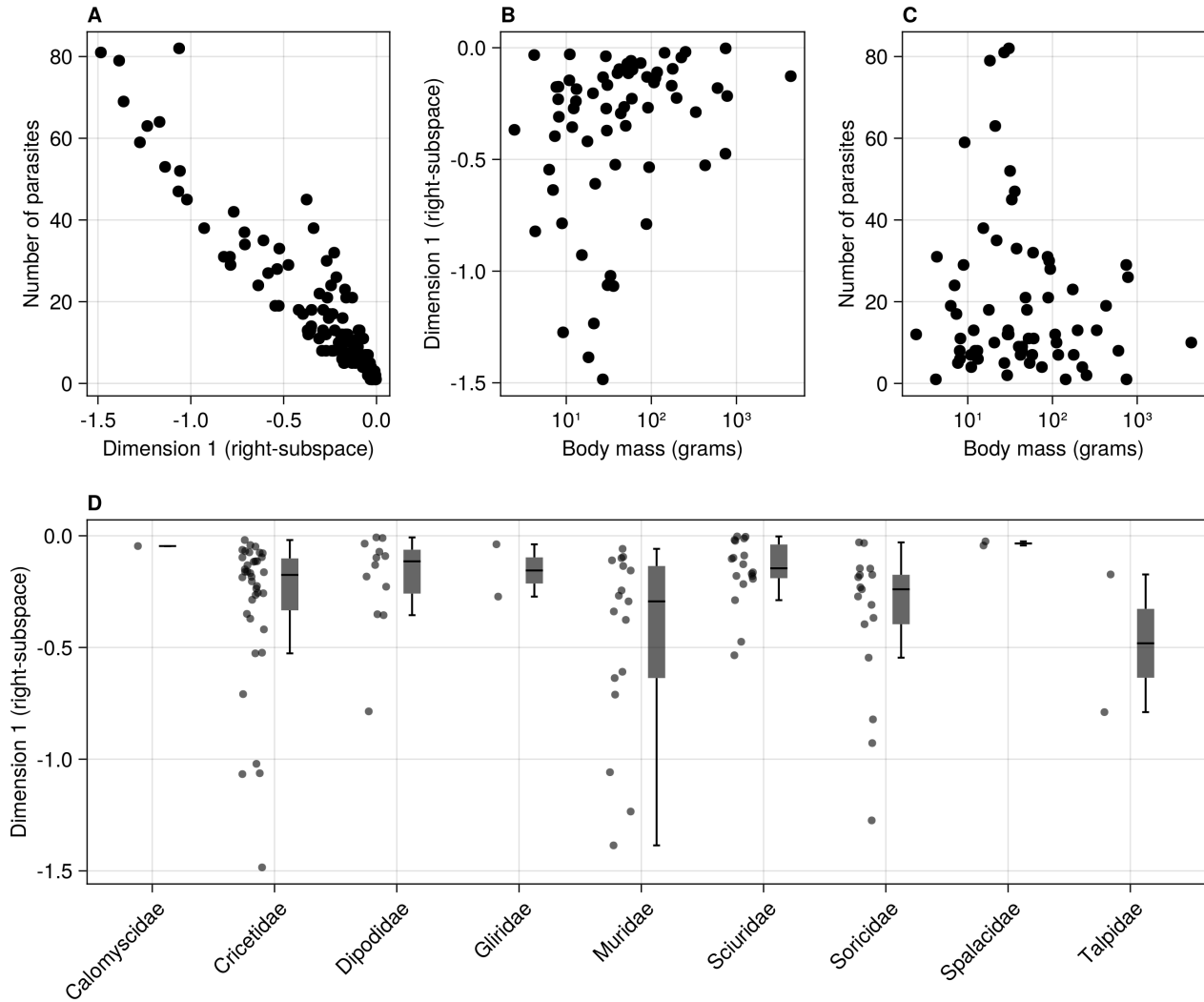


Figure 3: Ecological analysis of an embedding for a host-parasite metaweb, using Random Dot Product Graphs. **A**, relationship between the number of parasites and position along the first axis of the right-subspace for all hosts, showing that the embedding captures elements of network structure at the species scale. **B**, weak relationship between the body mass of hosts (in grams) and the position alongside the same dimension. **C**, weak relationship between body mass of hosts and parasite richness. **D**, distribution of positions alongside the same axis for hosts grouped by taxonomic family.