

METHOD OPEN ACCESS

Challenges in Detecting Ecological Interactions Using Sedimentary Ancient DNA Data

Fiona Margaret Callahan¹  | Jacky Kaiyuan Li² | Rasmus Nielsen^{1,3,4,5}

¹Center for Computational Biology, University of California Berkeley, Berkeley, California, USA | ²Biostatistics Division, School of Public Health, University of California Berkeley, Berkeley, California, USA | ³Department of Integrative Biology, University of California Berkeley, Berkeley, California, USA | ⁴Department of Statistics, University of California Berkeley, Berkeley, California, USA | ⁵Center for GeoGenetics, University of Copenhagen, Copenhagen, Denmark

Correspondence: Fiona Margaret Callahan (fiona_callahan@berkeley.edu) | Jacky Kaiyuan Li (li_jacky@berkeley.edu) | Rasmus Nielsen (rasmus_nielsen@berkeley.edu)

Received: 23 August 2024 | **Revised:** 21 January 2025 | **Accepted:** 27 January 2025

Funding: The authors received no specific funding for this work.

Keywords: animal distribution | biodiversity | DNA, ancient | DNA, environmental | ecosystem | metagenomics | microbial interactions | spatio-temporal analysis

ABSTRACT

With increasing availability of ancient and modern environmental DNA technology, whole-community species occurrence and abundance data over time and space is becoming more available. Sedimentary ancient DNA data can be used to infer associations between species, which can generate hypotheses about biotic interactions, a key part of ecosystem function and biodiversity science. Here, we have developed a realistic simulation to evaluate five common methods from different fields for this type of inference. We find that across all methods tested, false discovery rates of interspecies associations are high under simulation conditions where the assumptions of the methods are violated in a variety of ecologically realistic ways. Additionally, we find that for more realistic simulation scenarios, with sample sizes that are currently realistic for this type of data, models are typically unable to detect interactions better than random assignment of associations. Different methods perform differentially well depending on the number of taxa in the dataset. Some methods (SPIEC-EASI, SparCC) assume that there are large numbers of taxa in the dataset, and we find that SPIEC-EASI is highly sensitive to this assumption while SparCC is not. Additionally, we find that for many methods, default calibration can result in high false discovery rates. We find that for small numbers of species, no method consistently outperforms logistic and linear regression, indicating a need for further testing and methods development.

1 | Introduction

Having a better understanding of how and why ecological communities change over time facilitates informed decisions regarding management and environmental protection (Beng and Corlett 2020; Alsos et al. 2024; Williams et al. 2023). One of the key components of this is to understand the influence of species on each other as communities assemble and during periods of environmental change (Akesson et al. 2021; Dussex et al. 2021). Inclusion of species interactions in ecological modeling changes

predictions about the ecological effects of climate change (Akesson et al. 2021; Bascompte et al. 2019), extinction events (Dussex et al. 2021), and the dynamics of whole ecosystems (Alsos et al. 2024). Species distribution data up to this point have been largely limited to spatial data or data over short time periods relative to many of the ecological processes at play (Beng and Corlett 2020). However, as sedimentary ancient DNA (sedaDNA) data become more available, it is becoming possible to survey populations across large spatial and temporal extents and to simultaneously capture data across a wide range of taxa (Beng

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2025 The Author(s). *Environmental DNA* published by John Wiley & Sons Ltd.

and Corlett 2020; Williams et al. 2023). These data present the opportunity to study how different taxa have co-occurred over large spatiotemporal scales and to make inferences about association networks among species (Beng and Corlett 2020; Alsos et al. 2024).

Several categories of methods have been used to infer associations between species using data from various proxies for species presence or abundance, including sedaDNA (Chen and Ficetola 2020; Kurtz et al. 2015; Popovic et al. 2019). These associations may arise from interactions among the species or from other sources such as shared responses to the environment (Popovic et al. 2019; Dormann et al. 2012). One of the most popular correlative models for spatiotemporal inference using presence–absence data are referred to as species distribution models (SDM) or joint species distribution models (JSDM) (Elith and Graham 2009). These are generalized linear mixed models with either a logit or probit link function and can include a random effect that accounts for spatiotemporal autocorrelation (Wang et al. 2021; Schliep et al. 2018; Pollock et al. 2014). This class of models is highly varied and can include components that account for a variety of factors, but in most cases, they are designed for very few species in a single study and therefore may not scale to the number of species often seen in sedaDNA data sets (except see Pichler and Hartig 2021). An important feature of (some of) these methods is that they can account for autocorrelation between samples in space and time (Wang et al. 2021). Failing to account for this can result in high rates of false inference (Dormann et al. 2007). However, there are still some potentially important dynamics of this data type not accounted for by these models. Many JSDMs, including those examined here, do not account for uncertainty in covariates and time points, non-linear effects, and false detections of species (but see Hui et al. 2023; Clark and Wells 2023).

Though SDMs have primarily been used to detect correlations between species and their abiotic environment, some have also been used to detect both biotic and abiotic interactions (Wang et al. 2021). Wang et al. used a SDM method in a study that investigated ecological interactions over the last 50,000 years in the Arctic using sedaDNA data (Wang et al. 2021). They make inferences about the ecological dynamics between humans and megafauna, plants and megafauna, and the co-occurrence of different species over a large spatiotemporal scale (Wang et al. 2021). Here, we will refer to the methods they used as SDM-INLA. Another study implemented a JSDM for spatiotemporal ordinal abundance data using a Markov chain Monte Carlo (MCMC) algorithm (Pollock et al. 2014). They modeled the co-occurrence patterns of several frog species and several Eucalypts, concluding that frog species had positive residual correlation not accounted for by measured environmental variables, whereas the Eucalypts had negative residual co-occurrence patterns after accounting for the measured environmental variables (Pollock et al. 2014). We will refer to this method here as JSDM-MCMC.

Network analysis methods are another set of correlative methods for detecting associations between taxa using sedaDNA data (Banerjee et al. 2018). They are most often applied to microbial data, though not exclusively (Zimmermann et al. 2023). These methods include many different approaches, but in general, they use aspects of mathematical network theory in the inference and interpretation of associations. These methods are often used for

time series data (i.e., a single sediment core subsampled vertically, representing sampling through time), although the methods considered here do not explicitly model time or space (Kurtz et al. 2015; Popovic et al. 2019; Friedman and Alm 2012). SPIEC-EASI and SparCC (Kurtz et al. 2015; Friedman and Alm 2012) use sedaDNA read abundances per taxon as input. EcoCopula (Popovic et al. 2019), which is not explicitly designed for sedaDNA, can accommodate count, biomass, or presence–absence data, and incorporates covariates and species interactions. These *network analysis* methods have been used to investigate microbial associations in the human gut microbiome (Kurtz et al. 2015), link microbial network complexity to ecosystem functionality (Wagg et al. 2019), and investigate associations among a broad range of taxa in a marine ecosystem over a period of vast environmental change (Zimmermann et al. 2023). In the microbial context, DNA sequences are generally not assigned to species or taxa. Instead, reads are grouped into operational taxonomic units (OTUs) or amplicon sequence variants (ASVs) using sequence similarity (Bharti and Grimm 2021; Chiarello et al. 2022). On the other hand, in studies involving plants and megafauna, sequences are often assigned at a species or genus level (Wang et al. 2021). While we recognize that there are differences between these data types in real data, since all data are simulated in this study, we will use species, taxon, OTU, and ASV interchangeably here.

The associations inferred by these methods may arise from (1) direct interactions between species such as trophic, mutualistic, or competitive interactions, (2) indirect interactions such as similar responses to the same environment, or (3) falsely inferred associations. If species are associated through an unmeasured aspect of the environment or through an unmeasured species, associations inferred through correlative models may not be causal (Dormann et al. 2012). Therefore, we do not expect all inferred associations to be causal, but rather consider the aim of these models to be to generate hypotheses about causal interactions, which can then later be tested experimentally (Dormann et al. 2012). Additionally, false inferences will always be expected to occur at some low rate due to stochastic effects, but violations of modeling assumptions such as non-equilibrium dynamics, differing statistical distributions of the data and the models, or uncertainty in covariates can cause much higher rates of false inference (Dormann et al. 2012; Landi et al. 2018; Yuan et al. 2021).

One of the difficulties of inferring associations between species in this setting is that the number of possible pairs of taxa scales with the square of the number of taxa in the dataset (Weiss et al. 2016). In studies focusing on megafauna or plants the number of taxa is most often in the tens (Wang et al. 2021; Pollock et al. 2014), but in microbiome studies, the number of taxa can be in the hundreds (Kurtz et al. 2015; Wagg et al. 2019). This means that the number of parameters that are being inferred (interactions between pairs of species) can be much larger than the number of data points (Kurtz et al. 2015). In these cases, additional assumptions in the methods are necessary to make inference possible (Kurtz et al. 2015). Additionally, some methods consider sedaDNA as binary occurrence (presence-absence) data, while others consider relative read abundance data (Wang and Marshall 2016). These methods assume that relative read abundance is a proxy to organismal abundance or relative biomass, but there are many potential confounders that may affect this metric (Giguet-Covex et al. 2019).

Understanding the relationships between these models and the circumstances under which they succeed or fail to infer species associations will be informative as the availability of sedaDNA increases. In particular, although we find that many of these methods do not perform well with realistic data and currently available sample sizes (see Section 3), we are hopeful that this will be informative as data availability increases. Additionally, one of the areas of promise presented by sedaDNA is the fact that we can potentially study the effects of taxa across kingdoms using the same samples (Beng and Corlett 2020), but in order to maximize this potential, we will need to understand the performance of different methods across a broad range of ecological contexts.

Models to infer species associations have been tested against each other within the category of (J)SDMs (Elith and Graham 2009; Zurell et al. 2018) and within the microbial network modeling literature (Kurtz et al. 2015; Weiss et al. 2016), but little is known about the relative performance of methods in the two categories. Often, these methods have been tested using data simulated under simple models that do not directly incorporate ecological processes (Kurtz et al. 2015; Elith and Graham 2009; Zurell et al. 2018). Due to the complexity of real ecosystems, this may severely underestimate false inference rates of these models

because the simulated data likely meet their assumptions much better than real data does. Therefore, it is important to also test methods under more challenging conditions that are based on ecological theory. To this end, we have developed a novel simulation model that uses ecological theory to simulate species abundances and simulates the sedaDNA sampling process. We compare this to a simpler simulation model that simulates sedaDNA read counts for multiple species without simulating ecological processes. We specify inter-species and species–environment interactions to create datasets for which the true interactions are known and apply different inference methods to the simulated data to test accuracy across a range of models, parameters, and datasets (Figure 1).

2 | Materials and Methods

2.1 | Simulation Models

Simulated data are generated under several different models. The first set of simulations uses a multivariate probit regression model, which is modified to accommodate read counts. This is referred to here as *covariance matrix simulations*. In this

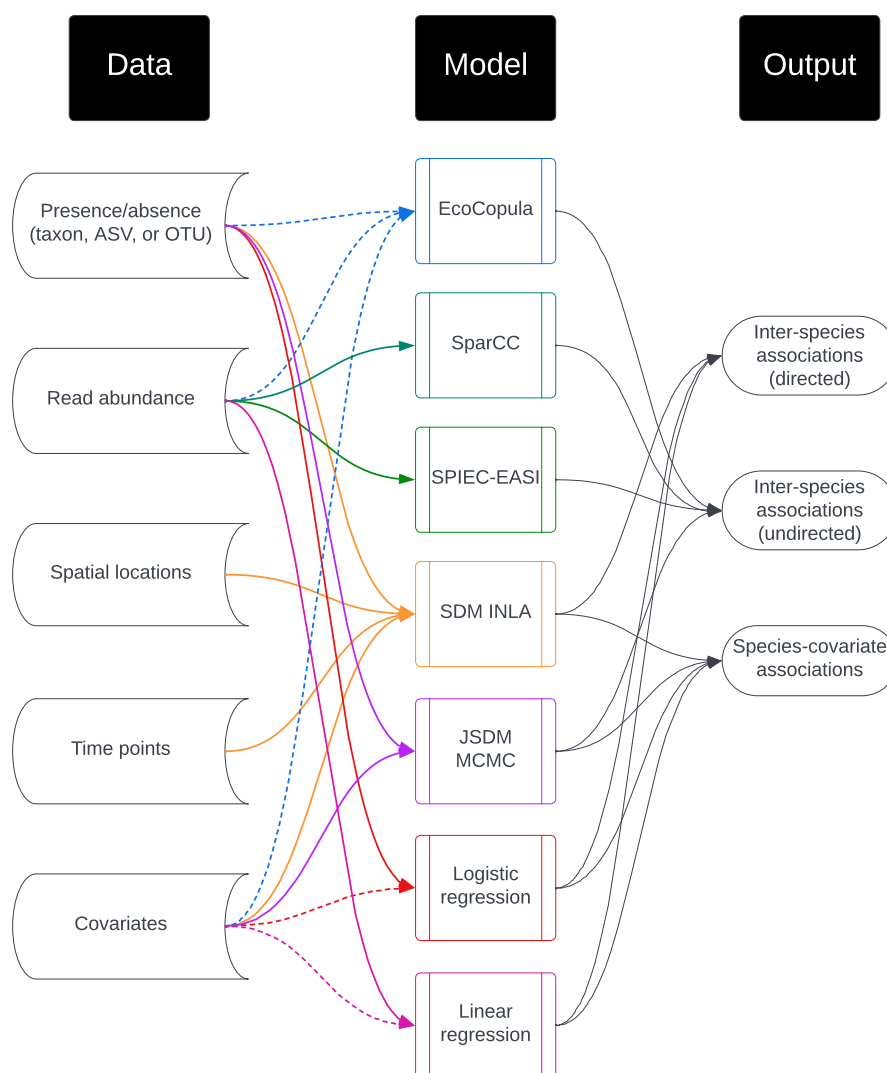


FIGURE 1 | Inputs and outputs of inference methods tested. Dotted lines indicate optional inputs.

model, species are simply associated via covariances in a latent Gaussian variable that is used to generate read abundances and presence–absence data. However, because we believe that ecological data can be much more complex than this model, we also simulated data under a more complex model that incorporates ecological theory, which we refer to as the *ecological simulation* model. We allowed parameters in this model to vary widely, but since we recognize that this may create unrealistic scenarios, replicates with a single set of parameters that generated more realistic data were also performed.

The *covariance matrix simulations* are designed to be as favorable as possible to the inference methods by violating their assumptions as little as possible. On the other hand, the *ecological simulations* introduce dynamics that violate many assumptions of the inference methods, so we expect that the methods will not perform as well under these conditions. The purpose of these models is to explore the robustness of these methods to a variety of realistic violations of their modeling assumptions.

2.1.1 | Notation

Throughout, j indicates species (used interchangeably with taxon, OTU, or ASV), \mathbf{s} indicates a location in two-dimensional space, and t indicates time. An omitted species index indicates the vector of values for all species. Bold symbols indicate vectors or matrices. Notation specifics can be found in Tables 1–3.

2.1.1.1 | Covariance Matrix Simulation Model. In one set of simulations, species interactions are encoded through the covariance matrix of a latent multivariate normal variable. The covariance matrix Σ for this set of simulations is defined as a $J \times J$ matrix with clusters of interacting and non-interacting species. In the clusters of interacting species, the interspecies covariances are defined using a latent factor model (Wilkinson et al. 2019), such that all covariances are nonzero in a cluster but the correlation values vary as described below (empirical distribution of covariances in Appendix S14: Figures S27, S28).

Let J_b be the number of species in the interacting cluster, and let F be the number of latent variables. Then we define Λ_b as a $J_b \times F$ matrix with standard normally distributed components. For each cluster b , the entries in Λ_b are

$$\Lambda_{ji} \stackrel{\text{iid}}{\sim} \text{Normal}(0, 1)$$

for all $j \in \{1, \dots, J_b\}, i \in \{1, \dots, F\}$.

Now let Γ_b be the intermediate covariance matrix before rescaling. Define Γ_b as

$$\Gamma_b = \Lambda_b \Lambda_b^T + \mathbf{I} \quad (1)$$

where \mathbf{I} is the identity matrix with dimension $J_b \times J_b$.

Then we re-scale Γ_b to a correlation matrix as follows:

Define

TABLE 1 | Variables and parameters in the *covariance matrix simulation* model.

b	Interacting block identifier
J	Total number of species
J_b	Number of species in interacting cluster b
F	Number of latent variables for each block
Λ_b	Latent factor matrix for block b
Γ_b	Intermediate covariance matrix before rescaling for block b
Σ_b	Rescaled covariance matrix for block b
Σ	Full covariance matrix for all species
\mathbf{X}	Matrix of covariate values for all sampling locations/times
β	Matrix of covariate effects for all species and covariates
$\mu^{(z)}$	Matrix of means of \mathbf{z}
$z_j(\mathbf{s}, t)$	Latent Gaussian variable used to define data (species j , location \mathbf{s} , time t)
$a_j(\mathbf{s}, t)$	Per-species read abundance data (species j , location \mathbf{s} , time t)
$y_j(\mathbf{s}, t)$	Binary species presence-absence data (species j , location \mathbf{s} , time t)

$$\mathbf{D}_b = \sqrt{\text{diag}(\Gamma_b)} = \begin{bmatrix} \sqrt{\Gamma_{11}} & 0 & \dots & 0 \\ 0 & \sqrt{\Gamma_{22}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\Gamma_{J_b J_b}} \end{bmatrix}$$

Finally, let the covariance matrix (also a correlation matrix in this case) for block b be defined as

$$\Sigma_b = \mathbf{D}_b^{-1} \Gamma_b \mathbf{D}_b^{-1}. \quad (2)$$

Now, for clusters $1, \dots, B$ we define the final covariance matrix Σ as a block matrix where $\Sigma_1, \dots, \Sigma_B$ fill independent blocks on the diagonal and all other off-diagonal entries are 0. All entries on the diagonal are 1. Note that by defining Σ in this way, we have some species in clusters that are correlated to one another, and others which are statistically independent.

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & \Sigma_2 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & \ddots & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \Sigma_B & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \mathbf{I} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & \dots & \mathbf{I} \end{bmatrix}$$

TABLE 2 | Variables in *ecological simulation* model.

\mathbf{s}	Point in two-dimensional space (unspecified units)
t	Point in time (unspecified units)
j	Taxon/species/ASV/OTU identifier
$n_j(\mathbf{s}, t)$	Number of individuals that deposited DNA
$a_j(\mathbf{s}, t)$	sedaDNA read abundance data (species j , location \mathbf{s} , time t)
$y_j(\mathbf{s}, t)$	Binary species presence-absence data (species j , location \mathbf{s} , time t)
$N_j(\mathbf{s}, t)$	Abundance (species j , location \mathbf{s} , time t)
$K_j(\mathbf{s}, t)$	Carrying capacity of the environment (species j , location \mathbf{s} , time t)
$M_j(\mathbf{s}, t)$	Number of migrants (species j , (to) location \mathbf{s} , time t)
$\mathbf{x}^*(\mathbf{s}, t)$	Measured covariates

Note: Variables change within a single simulation, in contrast to parameters, which are chosen only once per simulation.

We define a set of clusters such that there are approximately the same number of interactions as species to maintain sparsity of interactions (as is assumed by some of the methods tested) (Kurtz et al. 2015): $\sum_b (J_b^2 - J_b) \approx J$.

We also include covariates, representing the background environment. We simulate covariates and their effects following Wilkinson et al. (2019). We first simulate a matrix of four covariates plus an intercept, which vary over time and space. Call this matrix \mathbf{X} , which has dimension $5 \times N$, where N is the number of samples we are simulating (number of points in space-time, which for this model are statistically independent of each other). For the four non-intercept columns, we simulate matrix entries as $x_{ik} \stackrel{\text{iid}}{\sim} \text{Normal}(0, 1)$ for $i \in \{1, \dots, N\}$ and $k \in \{1, \dots, 4\}$. The intercept column is filled with 1's. Now we simulate the effect of each covariate on each species, which we store in the matrix β . This will be a $J \times 5$ matrix. For some simulations we will fill β with 0's, which denotes no effect of the environmental covariates on the species. For others, all entries of β will be independent and identically distributed standard normal variables. Now define $\mu^{(2)} = \beta\mathbf{X}$ as the $J \times N$ matrix of means for the latent multivariate normal variable. Each $\mu^{(2)}(\mathbf{s}, t)$ will be a single column of this matrix, and will therefore be a vector of length J . Note that when the entries in β are nonzero, this introduces correlation structure between the species that is different from the correlations defined in Σ . For example, if by chance two species have similar responses to the environment, then their presence and read abundance will be positively correlated even if the corresponding entry in Σ is 0.

Now we simulate the latent variable \mathbf{z} , which will be used to produce both presence-absence and read abundance data. Let $\mathbf{z}(\mathbf{s}, t)$ be a length- J vector of latent multivariate normal variables which are correlated across species according to Σ .

$$\mathbf{z}(\mathbf{s}, t) \sim \text{MVN}(\mu^{(2)}(\mathbf{s}, t), \Sigma) \quad (3)$$

TABLE 3 | Parameters in *ecological simulation* model.

J	Number of species in the model
P	Number of covariates (environmental/abiotic factors)
β	$J \times P$ matrix of covariate effects on the different species
α	$J \times J$ matrix of species effects on each other. Note: diagonal is 0
c_i	Tuning parameters
λ_j	Migration rate for species j
r	Intrinsic population growth rate
μ_M	Mean of migration rates across species
σ	Constant controlling amount of noise in population growth
d	Migration neighborhood radius for all species
p_n	Organism sampling rate
λ_a	Read sampling rate per organism
$\mathbf{x}(\mathbf{s}, t)$	Vector of environmental explanatory variables (covariates)
σ_x	Standard deviation for the covariate measurement noise
R	Detection threshold for presence-absence
V	Scaling parameter for the spatial covariance function for $\mathbf{x}(\mathbf{s}, t)$
T_x	Temporal period of the mean of for $\mathbf{x}(\mathbf{s}, t)$

Note: Parameters change between simulations but are constant within a simulation.

Now we will use these latent variables to generate presence-absence data and read abundance data.

Let $\mathbf{a}(\mathbf{s}, t)$ be the vector of simulated read abundances for all species at location \mathbf{s} and time t . We then define the abundance as a Poisson random variable with mean proportional to the latent variable $z_j(\mathbf{s}, t)$ if $z_j(\mathbf{s}, t)$ is positive, or 0 otherwise. For all j, \mathbf{s}, t :

$$a_j(\mathbf{s}, t) \sim \text{Poisson}(\mathbb{1}(z_j(\mathbf{s}, t) > 0) \cdot 100z_j(\mathbf{s}, t)) \quad (4)$$

Define the presence-absence data \mathbf{y} as 1 if $z_j(\mathbf{s}, t)$ is positive or 0 otherwise. For all j, \mathbf{s}, t :

$$y_j(\mathbf{s}, t) = \mathbb{1}(z_j(\mathbf{s}, t) > 0) \quad (5)$$

For the presence-absence data, this is a multivariate probit regression model. Due to the mean-zero truncated Gaussian, approximately half of the read counts are 0, which may be more or less realistic depending on the dataset. In this case, this is chosen to create optimal conditions for presence-absence data, as these methods lose power when percent presence is very low or high (Appendix S13: Figures S23–S26). Additionally, for read abundance data the latent variables are multiplied by 100 to put the mean of the Poisson distribution on the same order of

magnitude as several example datasets, although we also recognize that this may vary between studies and depend on the level of classification of reads (e.g., OTU vs. ASV vs. species vs. family levels) (Kurtz et al. 2015; Zimmermann et al. 2023).

2.1.1.2 | Ecological Simulation Model. We constructed a simulation model that takes as input (1) interspecies and species–environment interactions, (2) environment covariates, and (3) simulation hyperparameters (e.g., life history traits, detection rates). The output of the simulation is sedaDNA read abundance and presence–absence data for each species at each point in time and space. Presence–absence data are created by setting a threshold and mapping the read abundance data to binary data.

This is a population-level model in which all individuals of the same species share the same dynamics and traits. It includes space explicitly, over which abiotic environmental covariates can vary in both space and time and an arbitrary number of species whose populations vary over space and time.

At each time step, we model migration and logistic population growth that depend on a time-varying carrying capacity. The carrying capacity at each time point and location is a function of abiotic covariates and the abundances of other species. We also model the detection process, including modeling DNA deposition, covariate measurement uncertainty, and varying numbers of sampled points (Figure 2).

Interactions are assumed to have an effect on carrying capacity. For example, competition between two species is represented through a lower carrying capacity for one species when the abundance of the other species is higher. Trophic interactions between species can also be represented in this way since a higher abundance of prey may increase the carrying capacity of the environment for the predator. Conversely, a higher abundance of a predator may decrease the carrying capacity for its prey. Similarly, mutualistic relationships may be represented as a positive relationship between the abundance of one species and the carrying capacity of the other. Associations between species arise as an emergent property of these simulated interactions.

Note that variables (Table 2) change within a simulation, whereas parameters (Table 3) can be changed between simulations but are constant in a given simulation.

2.1.1.2.1 | Mechanistic Model of Abundance. We model species abundance as logistic growth in discrete time with noise and migration from a neighborhood of locations. Simulations were initialized with species abundances at time zero of $N_j(\mathbf{s}, 0) = 10$ for all species and locations, where $N_j(\mathbf{s}, t)$ is the species abundance of species j in location \mathbf{s} at time t . The first 100 time points are discarded before analysis.

The algorithm then proceeds by simulating the abundances as a two-step process as follows. At every time point $t > 0$:

Migration: At location \mathbf{s} , the number of new immigrants at time t for species j is $M_j(\mathbf{s}, t)$, which is assumed to be Poisson distributed at a rate that depends on the number of individuals of species j that were in a neighborhood of location \mathbf{s} at time t , and a species-specific migration rate λ_j . Note that migrating individuals are not subtracted from the population they come from, which may be unrealistic for some scenarios but realistic for some others where the dispersal mechanism uses propagules rather than individual movement. The migration rate varies across species, with the actual value for each species drawn from a Gamma distribution parameterized with the mean across species set to μ_M . The radius of the neighborhood, d , and mean migration rate are simulation parameters (allowed to vary between simulations but constant for a given simulation). The abundance in each location is then adjusted as follows:

$$M_j(\mathbf{s}, t) \sim \text{Poisson} \left(\lambda_j \cdot \underbrace{\sum_{\mathbf{s}^* \in S} N_j(\mathbf{s}^*, t-1)}_{\text{individuals in neighborhood}} \right) \quad (6)$$

where S is the set of spatial locations within radius d of location \mathbf{s} and $\lambda_j \sim \text{Gamma}(\text{shape} = \mu_M / 0.01, \text{scale} = 0.01)$. Then the

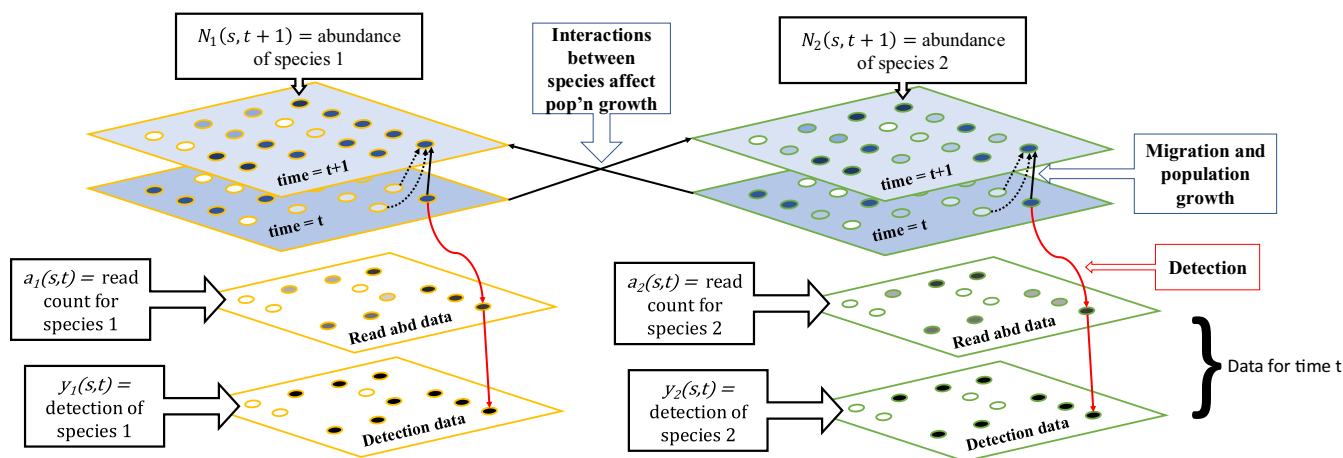


FIGURE 2 | Ecological simulation diagram. At each time step, we model local migration and logistic population growth, which depends on the abundance of other species with specified interactions and a set of abiotic factors. We also model a sedaDNA detection process, resulting in two types of data: read abundance and presence–absence (detection) data. Arrows indicate the flow of information through the simulation. Dotted black arrows indicate migration, solid black arrows indicate population growth and species interactions, solid red arrows indicate data simulation.

migrating individuals are added to the population to form an intermediate population level, $N_j^*(\mathbf{s}, t)$.

$$N_j^*(\mathbf{s}, t) = N_j(\mathbf{s}, t) + M_j(\mathbf{s}, t) \quad (7)$$

Population growth: The change in population size between time t and time $t + 1$, is modeled using discrete logistic growth with Gaussian noise and with a carrying capacity that depends on other biotic and abiotic factors. The growth rate of all species, r , and the noise of the process, scaled by σ , can be adjusted for each simulation but are constant within a simulation. We then assume that the local change in population size is

$$\Delta N_j(\mathbf{s}, t) = \underbrace{r N_j^*(\mathbf{s}, t) \left[1 - \frac{N_j^*(\mathbf{s}, t)}{K_j(\mathbf{s}, t)} \right]}_{\text{Logistic population growth}} + \underbrace{\sigma N_j^*(\mathbf{s}, t) W_j(\mathbf{s}, t)}_{\text{Gaussian noise}} \quad (8)$$

where $W_j(\mathbf{s}, t) \sim \text{Normal}(0, 1)$.

The variance of the Gaussian noise term in the equation above scales with the population size as in (Dennis 1989).

We model carrying capacity as a linear function of covariates and a non-linear function of other species abundance. The vector of carrying capacities for all species, $\mathbf{K}(\mathbf{s}, t)$, in the logistic growth equation depends on simulated abiotic environmental conditions, $\mathbf{x}(\mathbf{s}, t)$, and the abundance of the other species at time $t - 1$. The direction of the effect of each covariate and each species on other species is defined by matrices β and α . β is a J by P matrix where each entry in the matrix indicates the sign and direction of influence of a particular covariate on a particular species. α is a J by J matrix where each entry in the matrix indicates the sign and direction of influence of a species on another species. In order to avoid interspecies effects going to infinity, the arctangent function is applied to species effects on each other. This function has a horizontal asymptote, creating an upper bound on the effect of one species on another, representing saturation of the interspecies effect. c_1 and c_2 are constants that control the relative strength of the effects of other species versus the abiotic environment on the growth rate of each species. Carrying capacities are truncated at 0 since they should not be negative.

$$\mathbf{K}^*(\mathbf{s}, t) = \underbrace{c_1 \cdot \beta \cdot \mathbf{x}(\mathbf{s}, t)}_{\text{abiotic effects}} + \underbrace{c_2 \cdot \alpha \cdot \arctan(\mathbf{N}(\mathbf{s}, t - 1))}_{\text{inter - species effects}} \quad (9)$$

$$K_j(\mathbf{s}, t) = \max(K_j^*(\mathbf{s}, t), 0) \quad (10)$$

2.1.1.2.2 | Covariates (Ecological Simulation). Covariates are simulated in two different ways for different *ecological simulation* sets (See Section 2.1.2 of Methods).

Method 1: For the set-parameter simulations, covariates are simulated using Gaussian random walks through time. Covariates are therefore autocorrelated in time but not in space in this case. The vector of covariates for each time and location, $\mathbf{x}(\mathbf{s}, t)$, is an input to the population simulation model.

$$x_i(\mathbf{s}, t + 1) = x_i(\mathbf{s}, t) + \epsilon_t \quad (11)$$

where $\epsilon_t \sim \text{Normal}(0, 0.01^2)$.

Method 2: In the random-parameter simulations, covariates were simulated using a spatial Gaussian random field with an exponential covariance function, where the global mean varied in time according to a deterministic sinusoidal function. The period of this function, T_x , and the spatial covariance parameter for the covariance function, V , are set uniformly at random within a set range for each simulation (see Table 4).

For each time point, all covariate values across space are drawn from a spatial random field with an exponential covariance function as follows:

$$\mathbf{x}_i(\cdot, t) \sim \text{MVN}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_x) \quad (12)$$

with the mean vector filled with

$$\mu_t = \sin\left(\frac{2\pi t}{T_x}\right).$$

The mean value of the spatial random field (global mean of the covariate), vary in time according to a sine function. T_x is the period of the sine function in time, and it is drawn uniformly between 100 and 10,000 once per simulation:

$$T_x \sim \text{Unif}(\text{min} = 100, \text{max} = 10,000) \quad (13)$$

TABLE 4 | Parameters in random simulation model.

Parameter	Settings in random simulations
J	10 or 100 based on simulation set
P	Uniform($\{3, 4, 5, \dots, J\}$)
β ($J \times P$ matrix)	Uniform randomly selected J entries set as 1 or -1 , all others set as 0
α ($J \times J$ matrix)	Uniform randomly selected J entries set as 1 or -1 , others set as 0
c_1	200
c_2	Uniform(min = 100, max = 500)
λ_j	Gamma(shape = $\mu_M/0.01$, scale = 0.01) (For each species)
μ_M	Uniform(min = 0, max = 0.1)
r	Uniform(min = 0.01, max = 1)
σ	Uniform(min = 0, max = 0.2)
d	Uniform(min = 12, max = 50)
p_n	Uniform(min = 0.001, max = 0.02)
λ_a	Uniform(min = 0.5, max = 100)
$\mathbf{x}(\mathbf{s}, t)$	See covariate section
σ_x	Uniform(min = 0, max = 0.2)
R	Uniform($\{1, 2, 3, \dots, 50\}$)

Note: Parameters without a species index were set to be the same for all species, times, and locations in a given simulation.

Σ_x is a spatial covariance matrix which is calculated as follows: for spatial points \mathbf{s}_1 and \mathbf{s}_2 :

$$\Sigma_{\mathbf{s}_1, \mathbf{s}_2} = \exp\left(\frac{-\|\mathbf{s}_1 - \mathbf{s}_2\|_2}{V}\right) \quad (14)$$

where $\|\mathbf{s}_1 - \mathbf{s}_2\|_2$ is the Euclidean distance between the points and V is drawn once per simulation as:

$$V \sim \text{Unif}(1, 100) \quad (15)$$

The covariance between two points decays exponentially as the distance increases, but it decays slower with increasing V .

2.1.1.2.3 | sedaDNA Read Abundance and Species Detection Model. We model sedaDNA by assuming that first, $n_j(\mathbf{s}, t)$ individuals from the population are selected to deposit DNA. $n_j(\mathbf{s}, t)$ is binomially distributed with parameters $N_j(\mathbf{s}, t)$ and p_n , the individual sampling probability. The mean of this distribution is therefore proportional to the true species abundance, $N_j(\mathbf{s}, t)$. This is equivalent to flipping a weighted coin with weight p_n for each individual present in the location to decide whether it deposits DNA. Then, the number of sedaDNA reads, $a_j(\mathbf{s}, t)$, is Poisson distributed with a rate proportional to the number of DNA-depositing individuals, $n_j(\mathbf{s}, t)$. The parameter λ_a dictates the rate at which each individual deposits DNA. In other words:

$$n_j(\mathbf{s}, t) \sim \text{Binom}(N_j(\mathbf{s}, t), p_n) \quad (16)$$

$$a_j(\mathbf{s}, t) \sim \text{Poisson}(\lambda_a \cdot n_j(\mathbf{s}, t)) \quad (17)$$

Presence-absence data, $y_j(\mathbf{s}, t)$, are created by truncating the read abundances at some threshold R .

$$y_j(\mathbf{s}, t) = \mathbb{1}\{a_j(\mathbf{s}, t) > R\} \quad (18)$$

In the ecological model simulations, under many parameter settings, some species went extinct quickly, so species were filtered for at least 10% presence in the dataset. Therefore, the number of species actually analyzed is less than the number of species simulated in many cases. However, if the number of species after filtering was less than half the number originally simulated, the dataset was discarded.

2.1.1.3 | Covariate Measurement Uncertainty. Covariates are measured with Gaussian noise with variance σ_x^2 :

$$x_i^*(\mathbf{s}, t) = x_i(\mathbf{s}, t) + \epsilon \quad (19)$$

and $\epsilon \sim \text{Normal}(0, \sigma_x^2)$.

2.1.2 | Simulation Sets

Simulation sets were designed with three main axes of stratification: simulation design (*covariance matrix simulation* and *ecological simulation*), number of species (~ 10 and ~ 100), and sample size (100, 250, and 10,000). The different simulation designs and three sample sizes represent different levels of realism. Ecological data is inherently complex, so all methods must

make assumptions that are potentially violated by real data. However, methods may be differentially robust to these violations. If a method performs poorly on realistic data, this may be attributable to several causes. The assumptions of the method and the actual data generating process may be too different for the method to perform well, or the amount of data may be too small. We can differentiate between these different types of problems by examining performance at a realistic sample size (100 samples), a large but potentially feasible sample size (250 samples), and an unrealistically large sample size (10,000 samples). If a method is statistically consistent, when its assumptions are met it should converge to the true solution as the sample size gets large. Therefore we can attribute large errors at a large sample size to violations of modeling assumptions.

2.1.2.1 | Covariance Matrix Simulations Without Covariate Effects. We simulated 100 datasets with 10 species and 100 datasets with 100 species. For the datasets with 10 species, there was one interacting block with 4 species for a total of 12 one-way interactions. In the datasets with 100 species, there were 5 interacting blocks, each with 5 species, for a total of 100 one-way interactions. For this simulation, the covariate effects β were all set to 0. Although the covariates did not affect the species data, environmental variables themselves were still simulated. Although time and space are not included in this simulation algorithm, spatial and temporal labels were included in the final dataset to allow all methods to be tested. Each simulation was subsampled to 100 samples, 250 samples, and 10,000 samples for analysis.

2.1.2.2 | Covariance Matrix Simulations With Covariate Effects. We simulated 100 datasets with 10 species and 100 datasets with 100 species. Everything was the same as the set without covariate effects except that the covariate effects β were each set to be independent standard normal variables. In other words, for all $j \in \{1, \dots, J\}$, $i \in \{1, \dots, P\}$, $\beta_{ji} \stackrel{\text{iid}}{\sim} \text{Normal}(0, 1)$.

2.1.2.3 | Random-Parameter Ecological Simulations. We ran 100 simulations under the ecological model described above, with all parameters set uniformly at random in what we assessed to be a reasonable range (ranges and distributions specified in Table 4). In each simulation, 100 spatial locations arranged in a 10 by 10 grid were simulated for 10,000 time points, but the full simulated dataset was not analyzed. We simulated 100 replicates, each with 10 species and 100 species. We randomly re-sampled each dataset to 100 samples, 250 samples, and 10,000 samples for analysis.

2.1.2.4 | Set-Parameter Ecological Simulations. As these random-parameter simulations often lead to scenarios with poor performance of most methods (see Section 3), we also used a model with fixed parameters for which we expect somewhat better performance as many sources of noise are set to a low level. For example, measurement noise in the covariates was set at 0, the detection rate was relatively high, and noise in population growth was set low. Specific parameter values are described in Table 5. Using these parameter settings, we simulated 100 replicates, each with 10 species and 100 species. We randomly subsampled each dataset to 100, 250, and 10,000 samples.

2.1.3 | Testing Inference Models

All methods were used with default settings in the papers cited, with the exception of threshold adjustment to make receiver operating characteristic (ROC) curves. All methods that accept but do not require covariates (logistic and linear regression, SDM-INLA, EcoCopula) were tested both with and without covariates included in the analysis, regardless of whether the covariates had an effect on the simulated species data (Figure 1).

2.1.3.1 | Logistic Regression. Logistic regression was performed using the `glm` function from the `stats` package version 4.4.0 (R Core Team 2024). Separate regressions were run for each species as a function of all other species (and covariates for runs with covariates). Interactions were then considered significant based on a p -value threshold. For reported false discovery rates at a single p -value threshold, the threshold was chosen according to the Benjamini-Hochberg procedure with an expected false discovery rate of 0.05 (Benjamini and Hochberg 1995). For ROC curves, the threshold was varied from 0 to 1. Presence-absence data (and covariates, where applicable) were used as input to the model.

2.1.3.2 | Linear Regression. Linear regression was performed using the `lm` function from the `stats` package version 4.4.0 (R Core Team 2024). Analyses thereafter proceeded as described for logistic regression. Read abundance data (and covariates, where applicable) were used as input to the model.

2.1.3.3 | JSMD-MCMC. JSMD-MCMC was performed as in Pollock et al. (2014) using JAGS version 4.3.2 (Plummer 2003) and package R2jags version 0.7-1 (Su and Yajima 2021). It decomposed the species co-occurrence patterns into components describing shared environmental responses and residual patterns of co-occurrence, including species interactions. For each species pair, it returned a set of MCMC samples of the association parameter, and significance was determined based on whether a q -percent credible interval contained 0. To calculate the false discovery rate at a single threshold, the chosen q -percent was 95%. For ROC curves, the q -percent was varied from 0% to 100%. The model was run with covariates and with presence-absence data as the input. Convergence of the MCMC algorithm was evaluated by examining trace plots of multiple chains and using Gelman-Rubin statistics in the R package Coda (Gelman and Rubin 1992; Plummer et al. 2006) (Appendix S3: Figures S2-S5).

2.1.3.4 | SDM-INLA. For SDM-INLA, analysis was performed as in Wang et al. (2021) using INLA version 24.2.9 (Lindgren and Rue 2015). For model selection, WAIC was used to choose between four models: (1) using all other species and covariates as predictors, (2) using all other species as predictors without environmental covariates, (3) only environmental covariates as predictors without other species' data, and (4) no predictors. 95% posterior credible intervals were used to determine associations. Following Wang et al. (2021), all models also included a spatiotemporal effect (Wang et al. 2021). For ROC curves, no model selection was used (Models (1) and (2) form two separate curves) and the cutoff for the quantile

TABLE 5 | Parameters in set-parameters simulation.

Parameter	Setting in set-parameter simulations
J	10 or 100 based on simulation set
P	J
$\beta (J \times P = J$ matrix)	I_J (Identity matrix)
$\alpha (J \times J$ matrix)	Randomly selected J entries set as -1 or 1 , others set as 0
c_1	200
c_2	300
λ_j	Gamma (shape = $\mu_M / 0.01$, scale = 0.01)
μ_M	0.01
r	0.05
σ	0.005
d	16 (only direct neighbors)
p_n	0.01
λ_a	1
$\mathbf{x}(s, t)$	See covariate section
σ_x	0
R	5

Note: Parameters without a species index were set only once per simulation for all species, times, and locations.

of the posterior credible interval was varied. Presence-absence data, covariates (where applicable), time points, and locations were used as input to the model.

JSMD-MCMC and SDM-INLA methods were not used in the simulations with a larger numbers of species because it took a prohibitively long time for the methods to run on the large number of datasets simulated here (Appendix S1). Additionally, for SDM-INLA, the method often failed, but failure was not consistently repeatable, even using the same dataset. This could be avoided for individual datasets by re-running multiple times or optimizing settings specifically for individual datasets, but doing this for every dataset was not practical for this study.

2.1.3.5 | SPIEC-EASI. SPIEC-EASI was run using the R package SpieEasi version 1.1.2 (Kurtz et al. 2023). All parameters were set to default, and both the `mb` version and the `glasso` version of the model were run for comparison. The input to the model was simulated read abundances. For ROC curves, the regularization parameter, `lambda`, was adjusted to 100 different values automatically by the SpieEasi package (`nLambda = 100`) and then results were averaged across simulations for the ordered `lambda` values. The `lambda` values were ordered but not necessarily the same values across runs of the simulation, since by default the model automatically selects the actual values. For the model-selected results (reported false discovery rates), a specific `lambda` value was selected through the default calibration procedure.

2.1.3.6 | SparCC. The SpiecEasi package also implements the model SparCC, which was originally published by Friedman and Alm (2012) (Friedman and Alm 2012). This model, as implemented in SpiecEasi, was also tested on simulated read counts. For ROC curves, the threshold for estimated covariance, which is used to call interactions, was adjusted. For the false discovery rate after model selection, a pseudo p -value was generated using the bootstrapping procedure described by Friedman and Alm (2012) and implemented in SpiecEasi (Friedman and Alm 2012; Kurtz et al. 2023). The cutoff for this pseudo p -value was set at 0.05.

2.1.3.7 | EcoCopula. EcoCopula was run using the R package EcoCopula version 1.0.2 (Popovic et al. 2019). The model was run with and without covariates. This method, unlike the others, can take flexible types of input data. Therefore, we tested it using the mode where it takes binary (presence-absence) data and where it takes count (read abundance) data. By default, the method selects 100 values to test for the regularization parameter, λ . Then it chooses one using BIC. The default selection procedure for λ was used for FDR results. ROC curves were produced by varying the λ values (selected from those chosen by the default model) and averaging the resulting true and false positive rates across simulations. The λ values were ordered but not necessarily the same values across runs of the simulation. The λ parameter was also set at 0 to complete the ROC curve, although this was not among the values selected by the model by default.

2.1.4 | Counting Mistakes

The performance of these models was evaluated in several ways in order to get a more complete picture of their performance. First, we evaluated the success of the methods at detecting direct, causal interactions. Since this type of data is generally observational, and therefore can (and likely does) have significant confounding variables, we do not claim that any of these methods would detect direct, causal relationships consistently in real data. However, we believe that it is still useful to examine whether they are able to detect causal relationships when all variables are observed. For this metric, the sign of interactions was ignored in the calculation of false-positive and -negative rates. If an interaction exists and one was inferred, this was counted as correct regardless of the interaction being positive or negative. Interactions were considered directed, so A influencing B does not imply B influencing A. Therefore, for example, the false-positive count was calculated as the number of times where an interaction was inferred from A to B, but no interaction exists from A to B, regardless of the sign being positive or negative (Figure 3).

Second, we evaluated the success of the different methods at detecting whether there are any interactions between two species, whether they are direct or indirect through other species (indirect interactions). For this metric, we considered an interaction to exist if the two species are connected by interactions, regardless of the direction of the interactions. We likewise consider an interaction to be inferred between two species if any inferred associations connect them (Figure 3).

Third, we evaluate the success of the models at detecting direct associations between species without direction (direct, symmetric interactions). Some methods have the theoretical potential to infer an interaction in one direction but not the other (SDM-INLA, JSMD-MCMC, logistic regression, linear regression), although this should not happen often since they are all correlative methods, while other methods always infer symmetric interactions (SPIEC-EASI, SparCC, EcoCopula) (Figure 1). Using this metric, if an interaction exists from A to B, then we automatically assume one exists from B to A. Likewise, with the inferred associations, we assume that an inferred association in one direction implies one in the other direction.

Fourth, we consider indirect interactions to include interactions through covariates (indirect interactions, covariate). For example, if a covariate affects two species, these species are considered to interact. However, as interactions between covariates and species are not inferred, the inferred interactions were defined based on whether two species are connected by interactions between species, without regard for covariates. This metric did not cause different results than the indirect interactions metric without considering covariates for the *ecological simulations* and is not an informative metric for the *covariance matrix simulations* since all species are connected by covariates (See Section 3), but we include it here for completeness.

False discovery rate was calculated as false-positive-count/total-inferred-associations, or the probability that an inferred interaction was incorrect. This is distinct from false-positive rate, which was calculated as false-positive-count/total-actual-interactions. These metrics can differ enormously, especially when there is severe class imbalance, which is the case here with many more pairs with no actual interaction compared to the number of actual interactions (sparsity of the interaction matrix). Here, we present the false discovery rate after model selection and the false-positive rate at many thresholds as part of ROC curves.

2.1.5 | Analysis of Effect of Individual Simulation Parameters on Predictive Performance

A random forest model was used to predict false discovery rate for linear and logistic regression using the simulation parameters that were set at random in the simulations. Random forest models were run in R using Ranger version 0.16.0 (Wright and Ziegler 2017). The model was trained on 1000 simulations, with 10 species each and with 100 or 10,000 samples per simulation. Linear and logistic regression were both tested, using no covariates for either model, and with corrected p -values (Benjamini-Hochberg correction at a 0.05 false discovery control level). FDR was evaluated for direct/symmetric interactions only.

The following formula was used for random forest analysis (See Table 3 for definitions and Table 4 for ranges of parameter settings):

$$\text{FDR} \sim d + \sigma_x + r + \sigma + c_2 + \mu_M + p_n + \lambda_a + P + V + R \quad (20)$$

Variable importance in the random forest prediction was evaluated using permutation importance (`importance = "permutation"` in Ranger model). The random forest predictive

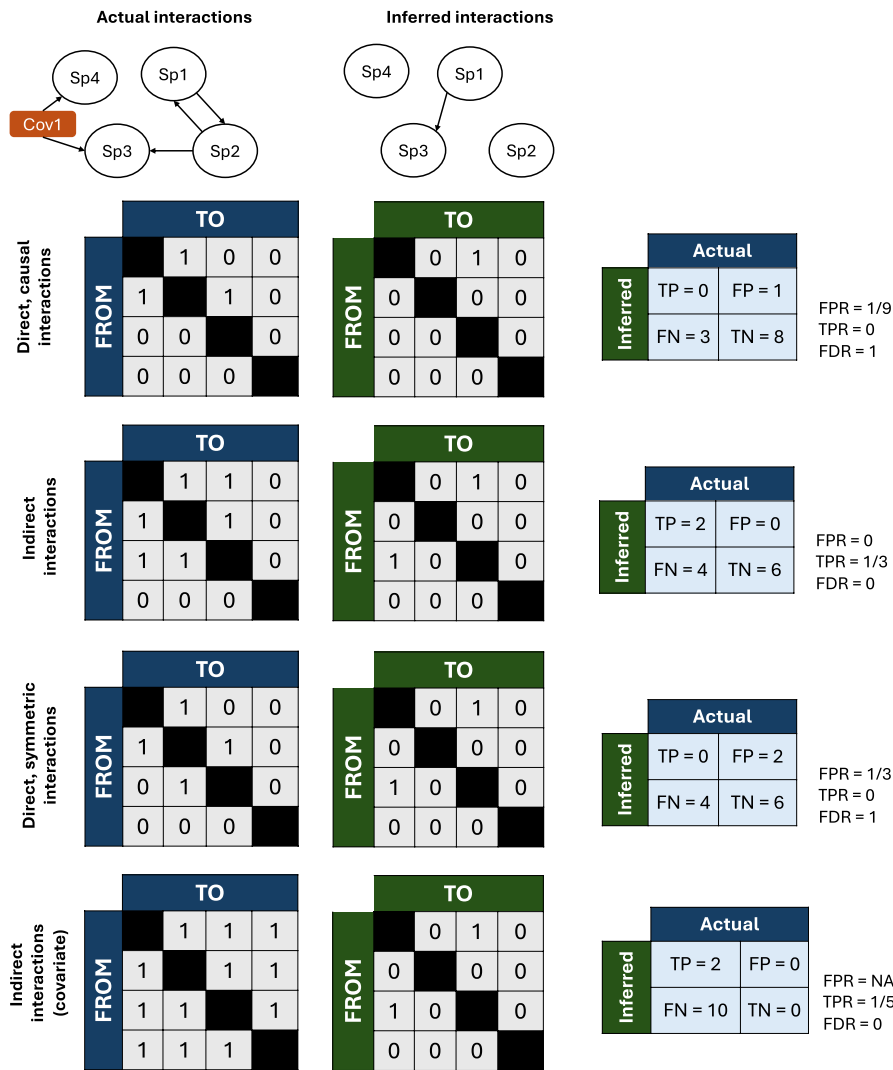


FIGURE 3 | Schematic explaining the calculations of false-positive rate (FPR), true-positive rate (TPR), and false discovery rate (FDR), when interactions are considered direct (symmetric or asymmetric) versus indirect (with or without considering shared covariate interactions).

performance was evaluated using root mean square error (RMSE) for an independent test set of 100 simulations. The RMSE of the random forest was compared to a naive predictor to evaluate how much predictive power the simulation parameters have for the resulting FDR. The naive predictor was the median FDR value for the dataset.

2.1.6 | Code

Simulations and SDM-INLA analyses were run in R version 4.3. All other analyses were run in R version 4.4 (R Core Team 2024).

Code is available at <https://github.com/Fiona-MC/eDNA-sims-pub>.

3 | Results

Several methods that test for associations between different taxa (used interchangeably with species, OTU, or ASV) were tested for their effectiveness at detecting these associations. Data were

simulated using two models: a mechanistic *ecological simulation* model and a simpler *covariance matrix simulation* model. The *ecological simulation* model was designed to be as realistic as possible, and as the complexity of ecological data necessitates making assumptions that real data may not follow, this means that these data violate many of the assumptions of the inference methods. On the other hand, the *covariance matrix simulations* are less realistic but also violate fewer of the assumptions of the inference methods, so we expect performance on this data to be better. These different simulation types examine the robustness of the methods to differing levels of violations of modeling assumptions.

3.1 | False Discovery Rates Are High for Most Inference Models and Simulations When Assumptions Are Violated

For the *ecological simulations*, the false discovery rate (FDR) for direct, causal interactions (Figure 3) in all cases was over 50% for 100 samples (Figure 4). For 250 samples, the same general trends apply, though the FDRs are lower in general, as expected, with the lowest being 38%. We also include simulations with 10,000 samples to see

which mistakes are caused by an insufficient amount of data. For 10,000 samples in the *ecological simulations*, the lowest false discovery rate was still 50% (Figure 4). The only simulation for which the false discovery rates are generally low is the *covariance matrix simulation* with no covariate effects, which is the simulation that minimally violates the assumptions of the methods.

Some of the mistakes are caused by including the direction of the inferred associations in counting mistakes. For the *ecological simulation* sets, the false discovery rate decreases when interactions are considered symmetric, and for the *covariance matrix simulation* sets, they are approximately the same. Interactions in ecology are often stronger in one direction than the other (i.e., species A affects species B but species B has little effect on species A), so we wanted to examine the effectiveness of these methods at detecting directed versus undirected interactions. In the *ecological simulation* model, some species have a unidirectional influence on other species, whereas in the *covariance matrix simulation* model, all interspecies effects are bidirectional and the species are organized in clusters, so there is little difference between the three metrics shown. SDM-INLA, JSMD-MCMC, and the regression methods can in theory predict asymmetric associations, but all other methods predict associations symmetrically. These associations are nonetheless often interpreted as potentially having causal meaning. Therefore, we have included the direct, causal interactions to illustrate one of many reasons that these associations should not be interpreted as directional or causal without additional information.

False discovery rates using the metric of indirect interactions (Figure 3) vary dramatically between methods. For some methods and simulations, all inferred interactions are correct (FDR=0), while for others, the false discovery rate is as high as 0.99 (nearly all inferred interactions incorrect) (Figure 4). EcoCopula and SPIEC-EASI claim to be able to avoid inferring indirect interactions by estimating conditional dependence of species presences (read abundances) given all other species presences (read abundances) (Kurtz et al. 2015; Popovic et al. 2019). However, here we see that the false discovery rate for SPIEC-EASI is either comparable or goes down when interactions are considered indirect (Figure 4). Notably, this remains true when the model is simplified to the *alternative covariance matrix simulation*, which has a different simulation mechanism specified in Appendix S7, including having interactions that are not organized in clusters (Appendix S7: Figures S13, S14). For EcoCopula, the story is less clear between considering interactions as direct but symmetric versus indirect. There are many cases where indirect interactions are inferred much better. In fact, using similar data, the number of interactions inferred after calibration varies considerably. It is possible that for a single dataset, this problem could be mitigated by fine-tuning the parameters based on the specifics of the data.

There are a few outliers in the FDR results for 100 samples that may be caused by very low overall rates of inferred interactions. For example, the 0% FDR for EcoCopula in the direct symmetric interactions case was caused by only four interactions that were all inferred correctly in one dataset (out of 100) (Figure 4). In all other datasets, no interactions were inferred. Similarly, for EcoCopula with no covariates for the set of simulations with

100 species, only two interactions were inferred in total, but both were incorrect, resulting in an FDR of 100% (Figure 4 and Appendix S2: Figure S1). When such low numbers of total interactions are inferred, the estimates of FDR may rely on just a few unusual cases and, therefore, may have high variance.

For the *covariance matrix simulations*, the three interaction metrics shown are quite similar (Figures 4–8 and Figures S8, S9). Any differences arise from transforming the inferred interactions to be symmetric or indirect, as the ground truth interactions are the same for all three. Whether interactions through the environment are considered correct or incorrect has almost no effect on FDR for *ecological simulations* (Appendix S5: Figures S7, S10, S11). For the *covariance matrix simulations*, all species are connected through covariates, and therefore this metric is uninformative.

3.1.1 | Model Calibration (Model Selection or p -Value Cutoff) Dramatically Affects Performance

False discovery rates are calculated using the default calibration for each method, but these calibration methods (model selection or choosing a p -value cutoff) vary between methods, which can cause large differences in the number of inferred interactions and the FDR. ROC curves do not depend on calibration, though the calibrated value is shown on the curves as the larger dot (Figures 5–8). The number of total inferred interactions is as low as 0 for some (seen as NA in Figure 4) and up to over 9000 inferred interactions per simulation for others (Appendix S2: Figure S1). In many cases, ROC curves look relatively good, but the FDR is very high, indicating that the underlying model is able to discriminate between interacting and non-interacting pairs of species, but the calibration method is selecting a point on the ROC curve that results in a high FDR.

As expected for lower amounts of data, when only 100 points were sampled, many methods often inferred very few interactions. On the other hand, with 10,000 samples, hundreds or thousands of total interactions were inferred across the 100 simulations per dataset (Appendix S2: Figure S1). One exception was SPIEC-EASI, which did not consistently infer fewer interactions when fewer points were sampled. In additional tests with as few as 10 samples, SPIEC-EASI continues to infer large numbers of interactions (Appendix S4: Figure S6), but with so few samples, it is unlikely that there is enough information in the data to infer this, indicating potential problems with calibration. Also, EcoCopula sometimes inferred drastically different numbers of interactions based on whether presence–absence or read abundance data were used as input. Interestingly, for some of these cases, the ROC curves do not differ significantly, indicating a problem with calibration rather than the underlying model. For example, for the *covariance matrix simulation* with no covariates and 100 species, over 6000 direct associations per simulation were inferred for read abundance data (with and without covariates) and around 40 direct associations per simulation with presence-absence data (with and without covariates). In general, we would expect that calibrated values would be somewhere on the ROC curve line. However, we see that this is not always the case for EcoCopula. The ROC curves are produced by varying the regularization parameter, but we find that manually setting

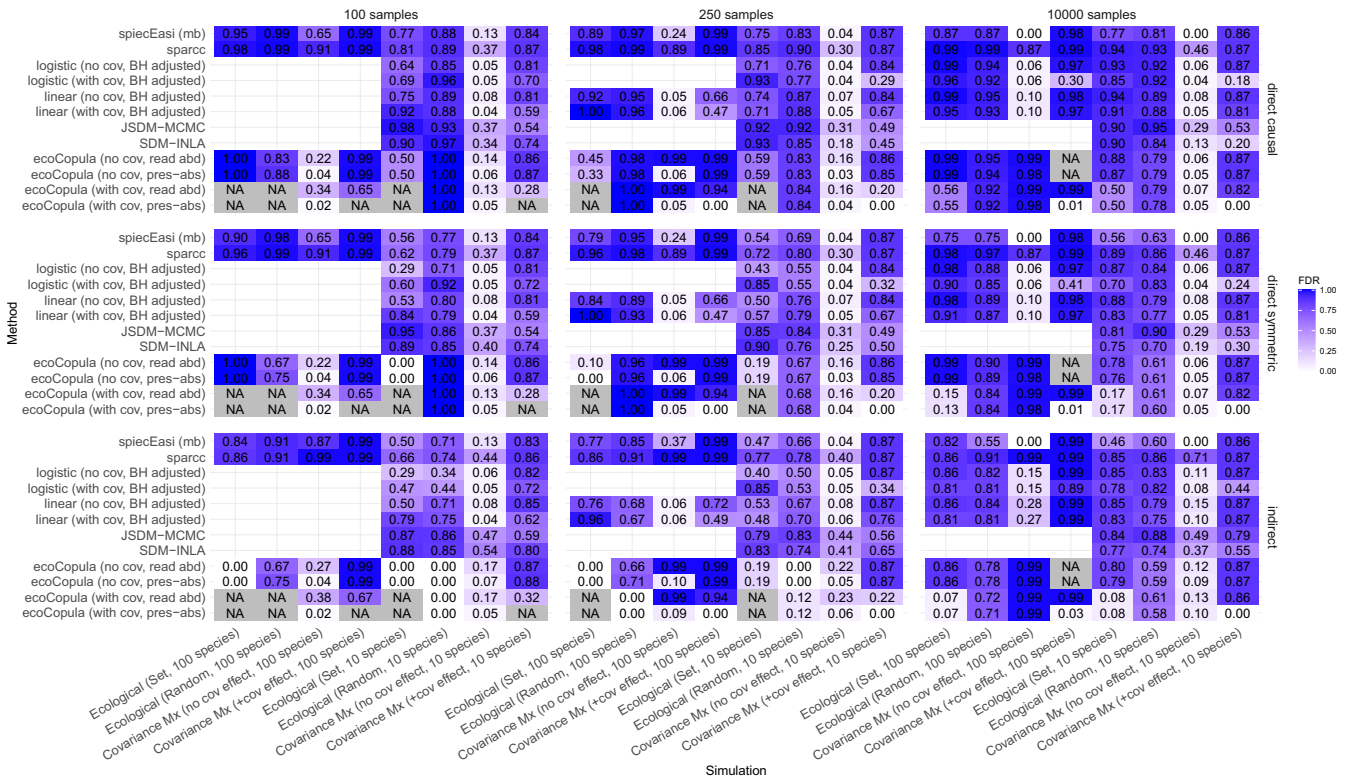


FIGURE 4 | False discovery rates of direct/causal interactions, direct/symmetric interactions, and indirect interactions with 100, 250, and 10,000 samples. The false discovery rate is defined as the number of false positive detections of species interactions divided by the total number of inferred interactions, with interactions defined in three ways. NA indicates that no interactions were inferred, and therefore the false discovery rate is undefined. Where there is an option, *cov* refers to whether covariates were used as predictors in the inference methods or as environmental effects in the simulation. Similarly, where adjustable, *pres-abs* and *read abd* refer to presence-absence data versus read abundance data inputs, respectively.

the value of this regularization parameter sometimes has a different result than allowing the model to select the value using BIC, even when we attempt to manually select the same value as the optimal value that the model selects. Additionally, curves are sometimes not concave due to the effect of averaging across runs while changing the regularization parameter. We therefore recommend careful examination of sensitivity to calibration metrics and parameters in these methods.

We also examined different calibration methods for the regression methods, including no false discovery rate correction (least conservative), Benjamini–Hochberg adjustment, and Bonferroni adjustment (most conservative). As expected, more interactions were inferred, and false discovery rates were generally higher with the less conservative methods (Appendix S15: Figure S29). All of the results shown for regression use the Benjamini–Hochberg adjustment, with the results for other correction methods shown in Appendix S15. For methods like SDM-INLA, the model is run separately for each species, but no FDR correction is applied (instead, the model is selected using WAIC), which may contribute to high false discovery rates.

3.1.2 | High Numbers of Samples Are Needed to Infer Associations Between Species in More Realistic Scenarios

Increasing sample size is expected to improve the performance of statistical methods in most cases, although model misspecification

can result in incorrect inference even as sample size gets very large (Blanchet et al. 2020). Here, we have tested each model at a relatively low sample size, which we believe is realistic for sedaDNA studies at this time (100 samples) and a larger but still potentially feasible sample size (250 samples). We have also tested all models with 10,000 samples, which is a higher number of samples than would currently be reasonable in sedaDNA studies. However, it is useful to examine performance at a very high sample size as large errors can then be attributed to insufficient robustness rather than insufficient data. As expected, we find that model performance is better for most methods and scenarios using a higher sample size, although FDR remains high in many cases (Figure 4).

Using only 100 samples, on average for random parameter settings in the *ecological simulation* model and for the *covariance matrix simulations* with covariates, no method performs better than a random classifier (defined as choosing to infer interactions completely at random, which is expected to follow the diagonal on a ROC curve plot) (Figures 5–8). The only set of simulations where the methods all consistently perform better than random using this small sample size was the simplest set of simulations where species have no response to environmental covariates, which is the least realistic scenario. For the *ecological simulation* with set parameters, some methods perform better than random, and some still perform poorly with low sample sizes.

With model selection, it is often the case that very few interactions were inferred with only 100 samples, which indicates low power

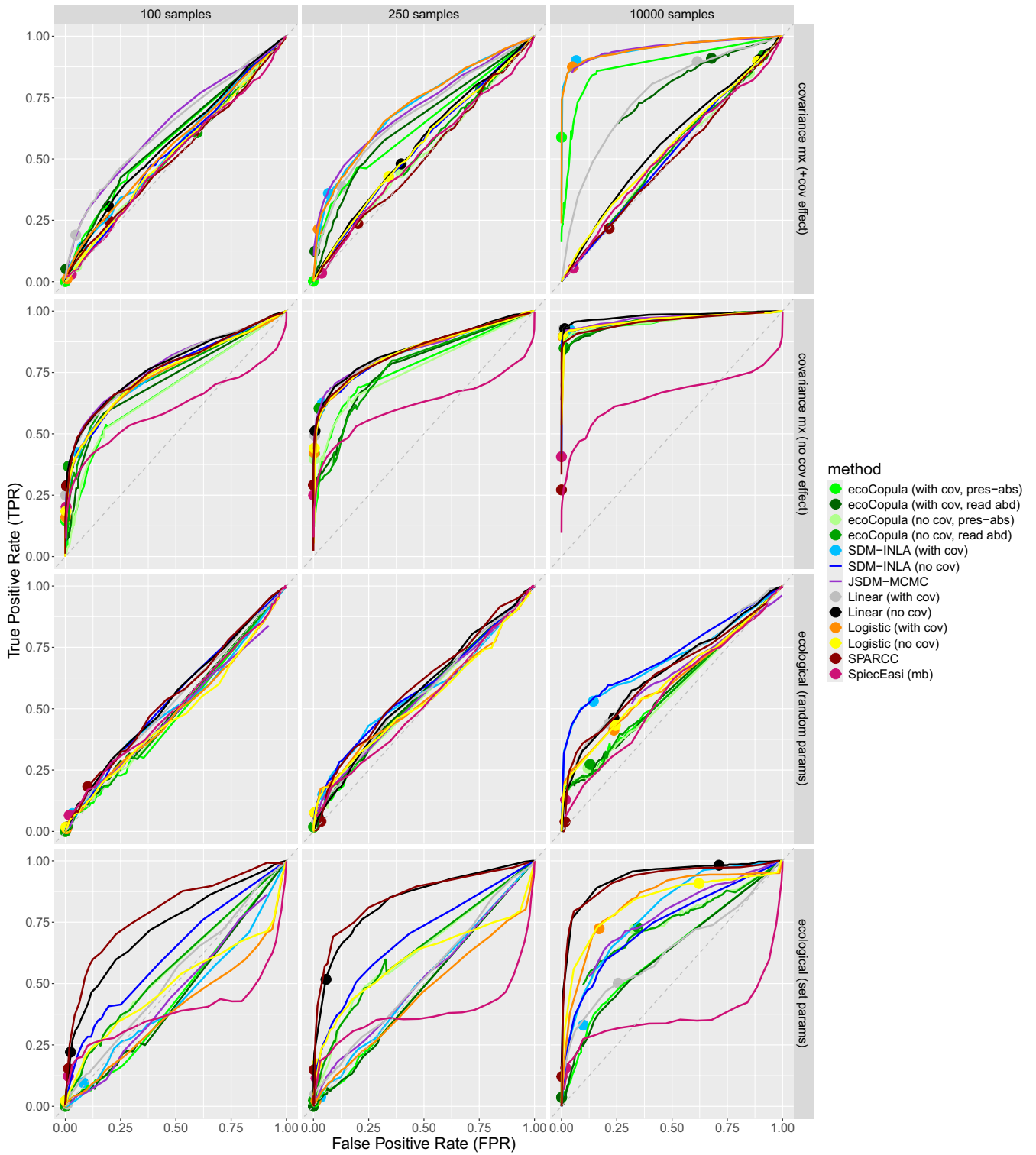


FIGURE 5 | ROC curves for inference of direct, symmetric associations between 10 simulated species (5–10 species actually present). *cov* refers to whether covariates were included in the model, but in all cases, species interactions were inferred. *read abd* and *pres-abs* are specified for EcoCopula because this method allows for different inputs (presence-absence data or read abundance data). Solid points are the points chosen by the default model selection of each method.

to detect these interactions if the method is calibrated correctly (Appendix S2). When interactions are inferred, they are most often incorrect for *ecological simulations* and *covariance matrix simulations* with covariate responses (Figure 4). Therefore, we find it to be a good sign when methods consistently infer very few interactions

since no method is able to detect the correct set of interactions. As expected, interactions are inferred more often with a higher number of samples and are more often correct (Figures 4–6). However, there is still a high chance that inferred direct interactions are incorrect under many scenarios (Figure 4). Using higher sample

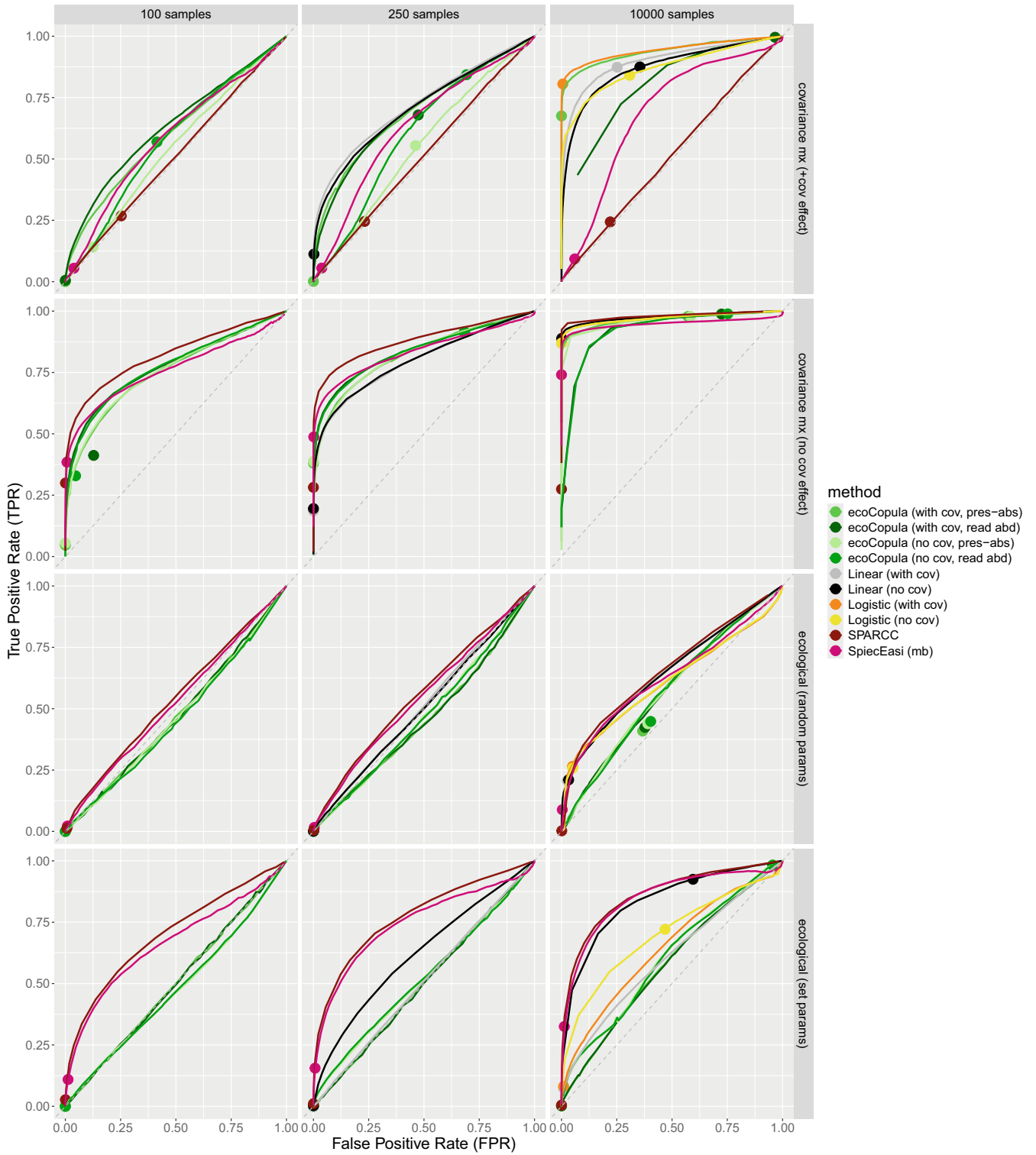


FIGURE 6 | ROC curves for inference of direct, symmetric associations between 100 simulated species (50–100 species actually present). *cov* refers to whether covariates were included in the model, but in all cases, species interactions were inferred. *read abd* and *pres-abs* are specified for EcoCopula because this method allows for different inputs (presence-absence data or read abundance data). Solid points are the points chosen by the default model selection of each method.

sizes, the false discovery rate is only consistently low in very favorable scenarios (*covariance matrix simulation* without covariate effects; note exception of set-parameters with covariates under the EcoCopula model, though this relies on very few inferred interactions, and a few exceptions for the *covariance matrix simulations* with covariates) (Figure 4).

The level of correlation induced by species interactions in real data, which is expected to influence the necessary number of samples, is unknown and likely highly variable between systems. The correlation level set here in the *covariance matrix simulation* was between -0.88 and 0.89 (Appendix S14: Figures S27, S28). The mean empirical correlation of the samples

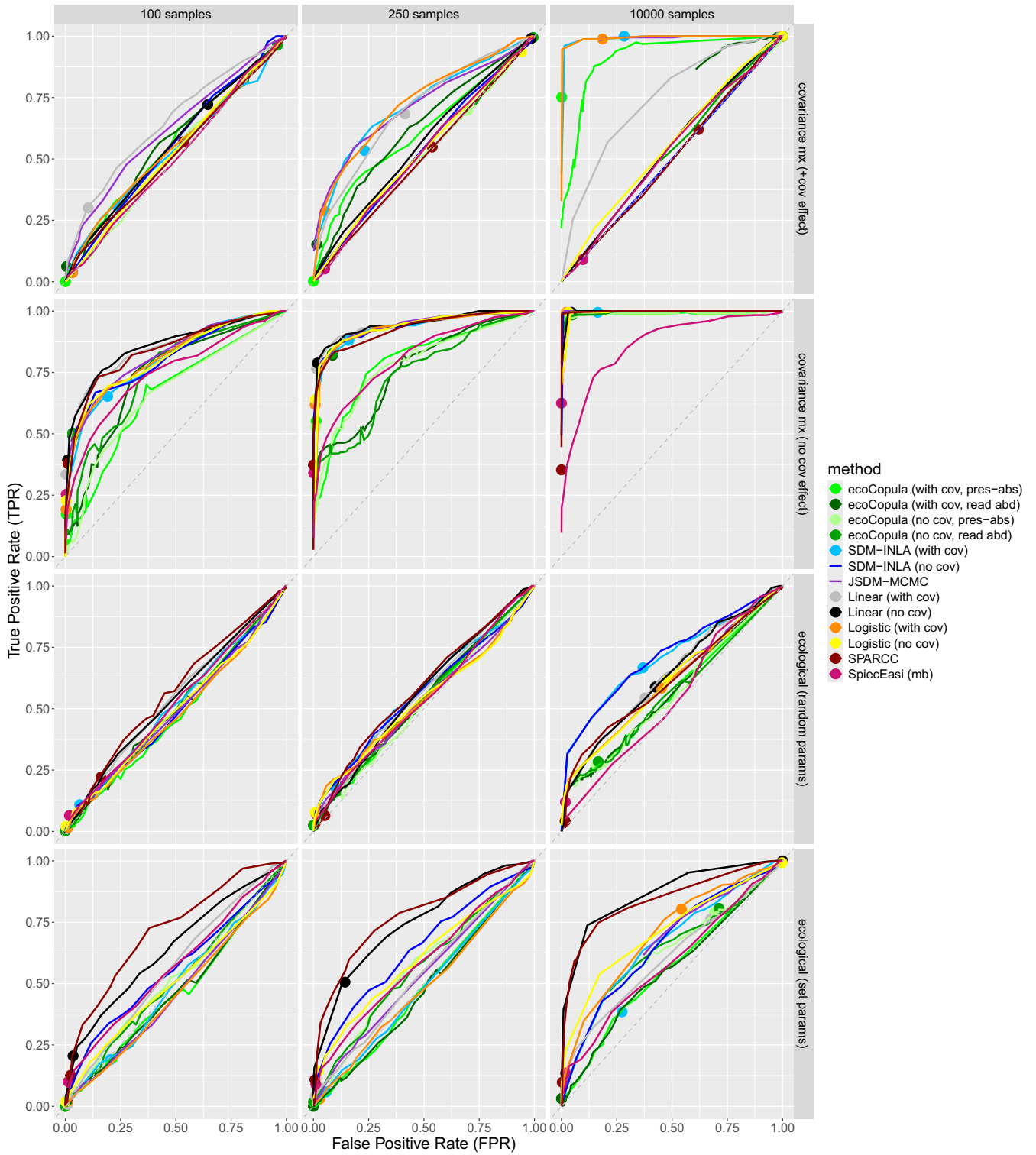


FIGURE 7 | ROC curves for inference of indirect associations between 10 simulated species (5–10 species actually present). *cov* refers to whether covariates were included in the model, but in all cases, species interactions were inferred. *read abd* and *pres-abs* are specified for EcoCopula because this method allows for different inputs (presence-absence data or read abundance data). Solid points are the points chosen by the default model selection of each method.

in the *ecological simulations* for set-parameters simulations when there is a positive interaction is 0.168, and the mean for negative interactions is -0.170 for simulations with 10 species. For 100 species, it is 0.118 for positive and -0.105 for negative interactions.

We also explored several additional *alternative covariance matrix simulations*, which have a different simulation mechanism specified in Appendix S7; these include some simulation sets with a lower and intermediate level of correlation. In one case, we ran an additional test of linear and logistic regression on *alternative*

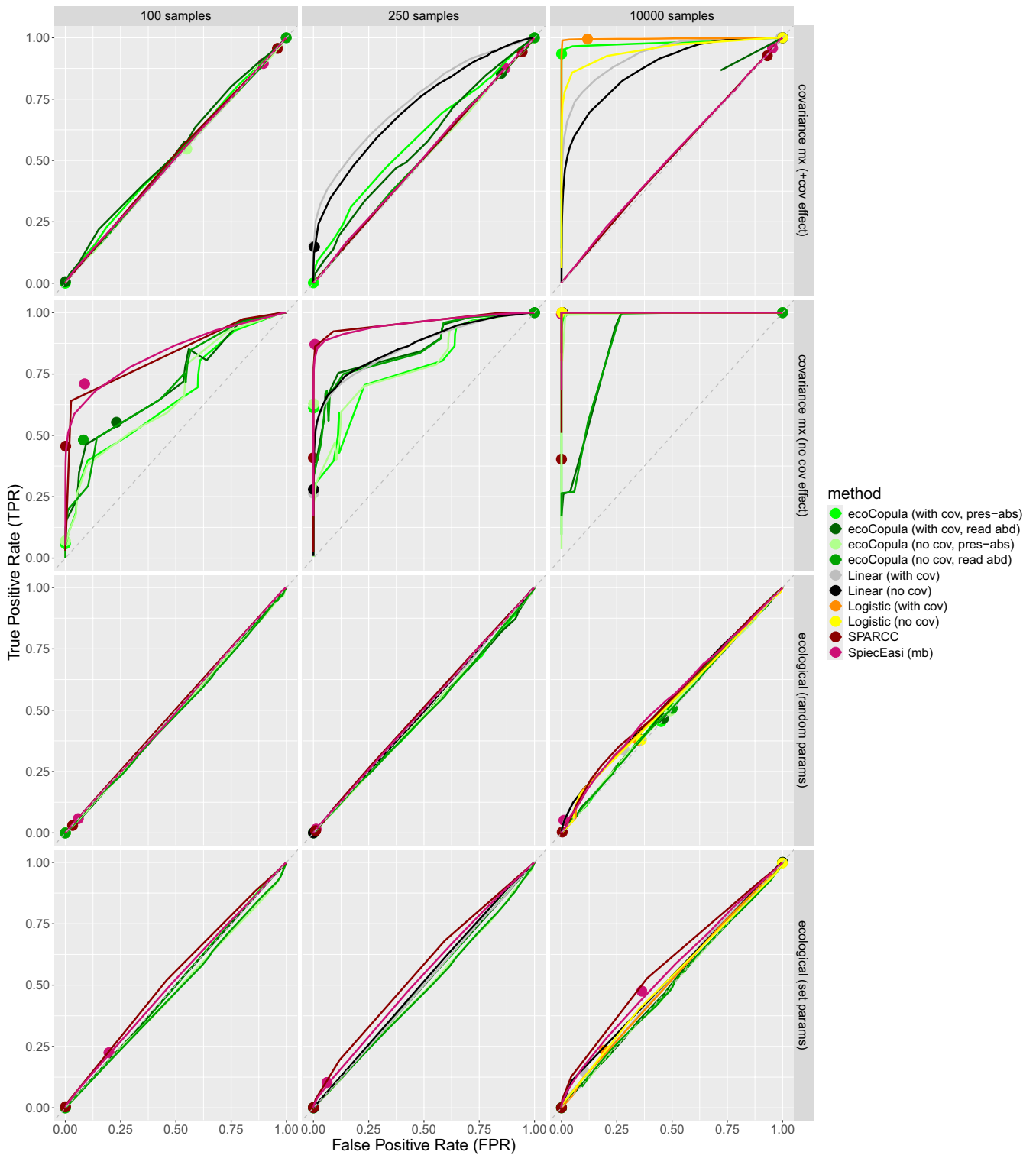


FIGURE 8 | ROC curves for inference of indirect associations between 100 simulated species (50–100 species actually present). *cov* refers to whether covariates were included in the model, but in all cases, species interactions were inferred. *read abd* and *pres-abs* are specified for EcoCopula because this method allows for different inputs (presence-absence data or read abundance data). Solid points are the points chosen by the default model selection of each method.

covariance matrix simulations where 100 species were associated in 50 independent pairs, and the correlation level was set at 0.5. In these simulations, 250 samples were sufficient for linear regression without covariates to classify the samples nearly perfectly according to ROC, with logistic regression as a close second (Appendix S7: Figures S15, S16). Other methods were not tested

on these data, although based on performance with large sample sizes in the *covariance matrix simulation* model, we believe that they would perform similarly well under this idealized scenario. Many factors, including covariate responses, were not modeled here, so we believe this is a less realistic scenario. We also ran *alternative covariance matrix simulations* with correlation values

set at 0.1 to correspond more closely to the correlation of the *ecological simulations*. In these cases, 100 samples were not sufficient for any method to perform well, but they performed well with 10,000 samples (Appendix S7: Figures S13, S14).

3.1.3 | Information About Species Abundance has a Mixed Effect on Inference Results

Some methods tested here take sedaDNA read abundance as input, which is considered a proxy for relative species abundance or relative biomass. Others take presence-absence data, which is derived from read abundances by setting a threshold for how many reads indicate species presence. Using presence-absence data is justified if read abundances are biased and therefore are not a good proxy for organism abundance. Additionally, this can be considered a way to give low-abundance taxa a higher weight in the model (Claussen et al. 2017). We consider these methods to fundamentally take the same data but process it differently.

We find that information about species abundance improves inference of species associations in the *ecological simulations* with set parameters, but not on average across all sets of parameters when they are set at random. In the simplest case, we compare logistic regression (presence-absence data) to linear regression (read abundance data). For the simulations where any methods perform better than a random classifier, linear regression performs better than logistic regression according to ROC curves in the set-parameter *ecological simulations* (Figures 5–8). EcoCopula is the only other method that flexibly takes presence-absence or read abundance data. The performance of EcoCopula varies, but does not seem to change based on data type for the *ecological simulations* (Figure 4).

On the other hand, for the *covariance matrix simulations*, in several cases we see that the performance of logistic regression (presence-absence data) exceeds that of linear regression (read abundance data) (Figures 4, 5, 7). Additionally, EcoCopula in several cases has a lower FDR with presence-absence data than with read abundance data (Figure 4). In the *covariance matrix simulations*, read abundance data are not a direct proxy for species abundance since no species abundances were simulated. More assumptions are violated by the count data than the binary data in this case, which may explain the differences in performance.

Here we have not added species-level biases to simulated read abundance, which means that read abundance is correlated to species abundance in the *ecological simulations*, although often with significant levels of noise. More biases may exist in read abundance in real data, so further investigation is needed to determine the circumstances under which sedaDNA read abundances are a good proxy for species abundances.

3.1.4 | Measuring Environmental Covariates Can Improve Inference Performance, but It Can Also Have a Negative Effect When There Is Multicollinearity

In the *covariance matrix simulations* where there is a covariate effect, we see that the only cases where the FDR is low is

for the methods that correct for the effect of covariates. That is because without correcting for the covariates, the methods are detecting shared responses to similar environments. This could still be defined as an association in some studies, so this is mainly a comment about interpretation (Figure 3). For example, SpiecEasi and SparCC cannot separate the environmental covariate effects from residual correlation between species because they do not accept covariates as input. This demonstrates that, as expected, these methods are often detecting shared responses to the environment rather than species interactions. We see the same effect for other methods when the covariates are unobserved (Figure 4). However, even in the cases where all covariates are observed perfectly, the false discovery rate of interactions is still much higher than the cases with no covariate effects in the simulation (Figure 4). Even if covariates are measured and included as input in the inference method, they can still cause errors. For example, in the *covariance matrix simulations* with covariate effects, the methods that take the Poisson distributed read count data still perform poorly even when they attempt to correct for the environment, even though all variables are measured. We believe this is because of violations of assumptions of the methods. In support of this conclusion, linear regression performs relatively well on the intermediate latent Gaussian variable $z_j(\mathbf{s}, t)$, but poorly on the Poisson reads derived from that variable (Appendix S16: Figures S30, S31). Another factor is that drawing the Poisson variable adds higher variance and thus the signal to noise ratio is lower, but as this issue persists at large sample sizes (if anything, it is exacerbated), we believe it is more likely attributable to mismatching assumptions.

On the other hand, contrary to expectations, in the *ecological simulations*, inference of species associations is often worse with environmental covariates included in the model than without them (Figures 5–8). We would generally expect the inclusion of many covariates to cause loss of statistical power (this would not affect ROC curves), but improve accuracy by correcting for confounding variables (Blanchet et al. 2020). However, if many covariates are included in a model and they do not affect multiple species, correlations caused by random fluctuations in the data can cause false discovery rates to be higher (Appendix S8: Figure S17).

Additionally, in the set of simulations with set parameters, there is a high level of correlation between some of the environmental covariates and the species that they affect. In the methods that take both species read abundances or occurrences and covariates as inputs, this can cause issues with identifiability between the effect of the covariate and the effect of the species (Appendix S9: Figure S18). Therefore, it is important to check for collinearity between species and covariates if both covariates and species are going to be included as explanatory variables in these models.

3.1.5 | SPIEC-EASI Is Highly Sensitive to Assumptions About Number of Species

SPIEC-EASI assumes that the number of species is large, although the actual number of species required is not specified. There is an approximation in their model that is more accurate

as the number of species gets large (Kurtz et al. 2015). We find that at the 5–10 species level, SPIEC-EASI is highly sensitive to this assumption according to ROC curves when looking at direct associations (Figure 5). With 50–100 species, this sensitivity disappears, and it diminishes when looking at indirect associations (Figures 6 and 7). Specifically, with fewer than 10 species, we find that SPIEC-EASI seems to have strong evidence for actual negative associations between species and strong evidence against actual positive associations between species (Appendix S10: Figures S19, S20). This effect is not observed after model selection because the model never selects the regularization parameter in this region of the ROC curve (Figures 5 and 6). Therefore, the false discovery rates after model selection are not noticeably affected (Figure 4).

3.1.6 | Regression Performs Similarly to Other Methods With the Same Data

Contrary to expectations, we find that in most scenarios, linear and logistic regression perform, on average, similarly to other methods with the same data input (Figures 5–8). We do not interpret this to mean that regression is the correct model for this data, but rather that the more complex methods do not improve modeling of the data structure in these scenarios. Each of the more complex methods attempts to account for different aspects of the data better than linear or logistic regression, but we find that for a variety of simulation scenarios, they do not succeed in better modeling the data generating process. It may be the case that under certain circumstances, each model does perform better, but we find that on average across many scenarios, no method recovers associations consistently better than regression. Additionally, regression methods are much faster than any of the other methods. In the case of the *covariance matrix simulations*, very few of the assumptions of logistic regression are violated. Therefore, it is expected that it would perform well given enough data. In the simulations shown here, we use the Benjamini-Hochberg method for false discovery rate correction. However, we also explored the use of no correction (less conservative) or the Bonferroni correction (more conservative). We found that the overall trends were similar, although the exact values of the false discovery rates vary considerably (Appendix S15: Figure S29). The ROC curves would not change based on the correction method. Rather, each correction method chooses one point on the corresponding curve. Linear and logistic regression will not work with the number of predictors (number of species plus covariates in this case) being greater than or equal to the number of samples. Therefore, for simulations with 100 species, the 100 samples case was omitted for the regression methods. Additionally, for logistic regression, when the number of predictors is close to the number of samples, we encounter numerical problems because the samples provided can be separated perfectly by the regressors. Therefore, results for logistic regression with 100 species and 250 samples are also not shown.

For methods that take presence–absence data, SDM-INLA performs best when averaged across random simulation parameters at a high number of samples in the *ecological simulation* and performs similarly to other methods in the *covariance matrix simulation* (Figures 5 and 7). However, this result does not generalize to our set-parameters *ecological simulation*, where

logistic regression performs better (Figures 5 and 7). The main difference between these models is that SDM-INLA models spatiotemporal autocorrelation. The effect of this may depend on whether a dataset has spatiotemporal structure and whether this structure is correctly modeled by the SDM-INLA method.

In simulations with random parameters, we see less separation of the different methods in the ROC curves than when the parameters were set at a single level (Figures 5 and 6). This is the effect of averaging performance across many different ecological scenarios. In some ecological contexts, some methods may perform better than others, but this may not be consistent across all contexts. As expected, the variance of the results within a method when parameters were varied was much higher than when parameters in the simulation were held constant.

3.1.7 | Many Dimensions of Simulation Parameter Space Affect the Success of Regression Models in Predicting Direct, Symmetric Interactions

In the *ecological simulations* produced with random parameters, there was a great variety of resulting false discovery rates between different parameter sets. In order to better understand how the simulation parameters are affecting the success of the inference models, we created an additional 1000 data sets with random parameters (drawn from distributions specified in Table 4), and used these simulation parameters as predictors for a random forest model that predicts false discovery rates of direct, symmetric interactions. Due to the large number of simulations in this set, this was only performed for logistic and linear regression, which are much faster than other methods. We believe this is justified as no other methods consistently outperformed these methods in tests on smaller simulation sets; though we do not know if the same simulation parameters would affect the performance of all methods. Random forest models were verified by testing for predictive accuracy on an independent test set of 100 simulations. RMSE for the test sets was evaluated against a naive predictor of FDR (median of observed FDRs). The random forest predicts FDR between 8% and 15% better than the naive predictor (Appendix S11: Tables S1, S2), indicating that the simulation parameters have some predictive power, though a lot of the variation in FDR is not well predicted.

For 10,000 samples, three simulation parameters seem to affect FDR significantly for logistic and linear regression. For larger r (population growth rate for all species) models tend to have lower FDR (perform better) (Appendix S12: Figures S22, S21). Larger population growth rates will cause the populations of species to change more rapidly and to reach their carrying capacity faster, causing less lag between a change in environment and the change in the species population. Non-equilibrium dynamics caused by lag between changing conditions and population growth can cause temporal and spatial autocorrelation not accounted for by covariates. This violates the assumption of conditional independence of the samples in linear and logistic regression and can lead to a loss of information.

With a larger migration radius (d), the models tend to have a higher FDR (perform worse) (Appendix S12: Figures S22, S21). In this simulation model, migration rates are not affected by

species interactions and may obscure the effect of one species on another.

The third simulation parameter that seems to affect false discovery rate in this set of simulations is the individual sampling probability (p_n) (Figure 9). When this parameter is higher, FDR tends to be higher (Appendix S12: Figures S22, S21). This effect is stronger for logistic regression than linear regression. This is a counter-intuitive result because we would expect a higher probability of sampling individuals to result in better performance of the models. However, we also observed that high observed rate of species presence can result in difficulty with inference because there is less power to infer interactions if species are present almost everywhere, which we believe to be the cause of this effect (Appendix S13: Figures S23–S26).

4 | Discussion

We have tested a range of methods that have been used to detect associations between taxa from sedaDNA data. We simulated data under a variety of models, including a simple model where all data points are independent over time and space, and a custom ecological model that leverages principles of ecological theory to create more realistic data. Explicitly simulating population dynamics through time in the presence of a changing abiotic environment will create patterns in the data that mimic those in nature. We find that in all but the most idealized scenarios, false discovery rates of species associations are high for all methods tested.

Learning how past ecosystems have changed in response to environmental change will be an essential part of understanding and mitigating present-day global climate change and making informed resource-management decisions (Landi et al. 2018). Biotic interactions are a key piece of this puzzle (Aksesson et al. 2021) and sedaDNA is a new frontier in understanding ecological interactions. It has the potential to uncover interactions within and between trophic levels, model species distributions over large spatiotemporal scales, and show how biodiversity has changed over periods of massive environmental change (Beng and Corlett 2020; Alsos et al. 2024; Williams et al. 2023; Wang et al. 2021). Robust computational methods are needed to come to accurate and reproducible conclusions, and it is important to understand the relative performance of methods under different conditions.

Many species distribution models have been developed in ecology (Elith and Graham 2009), but since these methods were developed with traditional ecological survey data in mind, we find that they are often not scalable to the number of taxa that is common in sedaDNA datasets. Many of these methods account for temporal and/or spatial information in the data and allow for information about abiotic covariates (Elith and Graham 2009; Wang et al. 2021). However, they do not account for the compositional structure of sedaDNA read abundance data. On the other hand, several methods have been developed specifically for sedaDNA data and therefore account for non-independence due to compositional sedaDNA read data (Kurtz et al. 2015; Friedman and Alm 2012). However, many of the most commonly used methods were developed specifically for microbiome sedaDNA

data and therefore assume large numbers of taxa (or OTUs/ASVs) in a single data set. We find that SPIEC-EASI is very sensitive to this assumption, exhibiting very poor performance on data with small numbers of taxa (10 or fewer). This is important, even in the sedaDNA field, because depending on the taxa under study and the taxonomic level of assignment, many studies may want to do similar analyses with small numbers of taxa (Wang et al. 2021; Pollock et al. 2014).

It has been previously established that inferring species interactions from spatial presence–absence data are a challenging statistical problem, and that there are many potential causes for false inferences from these data (Blanchet et al. 2020). It has also been documented, theoretically for co-occurrence (presence-absence) data (Blanchet et al. 2020) and empirically for read abundance data (Kurtz et al. 2015), that high sample sizes are needed to accurately estimate associations between taxa. The challenge is in part because a large number of taxa create an even larger set of potential interactions (Kurtz et al. 2015), but also due to the complexity of the ecological networks involved and the abiotic factors that influence them, among other reasons (Blanchet et al. 2020). From a mathematical perspective, for linear regression with all assumptions met, when the correlation is 0.1 (close to the average in our *ecological simulations*) using a p -value cutoff of 0.05 with Bonferroni correction, the necessary sample size to detect 80% of true interactions is nearly 2000 samples for 10 species and nearly 3000 samples for 100 species (Appendix S6: Figure S12). If the correlation induced by species interactions is higher, the expected sample size needed would go down (Appendix S6), which we also observe for simulated data (Appendix S7: Figures S15 and S16), so lower sample sizes may be sufficient under some circumstances. However, even in these cases, false discovery rates may remain high due to model misspecification. Since very few studies using sedaDNA currently have more than a few hundred sampled points, high sample sizes are currently unrealistic, but we are hopeful that future data will rise to meet this challenge. Additionally, shared responses to the environment may cause associations that are separate from direct interactions between species (Popovic et al. 2019), and therefore in real data we would caution users to interpret inferred associations carefully. We also observe in this study that even when the environmental variables are fully observed and statistically independent across time and space (thus minimally violating the assumptions of the methods), they still invariably cause higher false discovery rates than were seen in the *covariance matrix simulations* without a species response to the environment.

False discovery rates reported here depend heavily on the calibration of various methods (such as model selection or choosing a p -value cutoff). FDR reported here reflects one point on the corresponding ROC curve, which measures the performance of the model independent of calibration. As calibration is an intrinsic part of these methods, we consider these calibrated values to have substantial significance. The dramatic differences in calibrated FDR values between different versions of the same method and between methods are often more a consequence of the calibration than the underlying model. For example, the ROC curves for all methods may look similar across all models in many cases, reflecting similar success of the underlying model, but the average FDR varies considerably due to poor

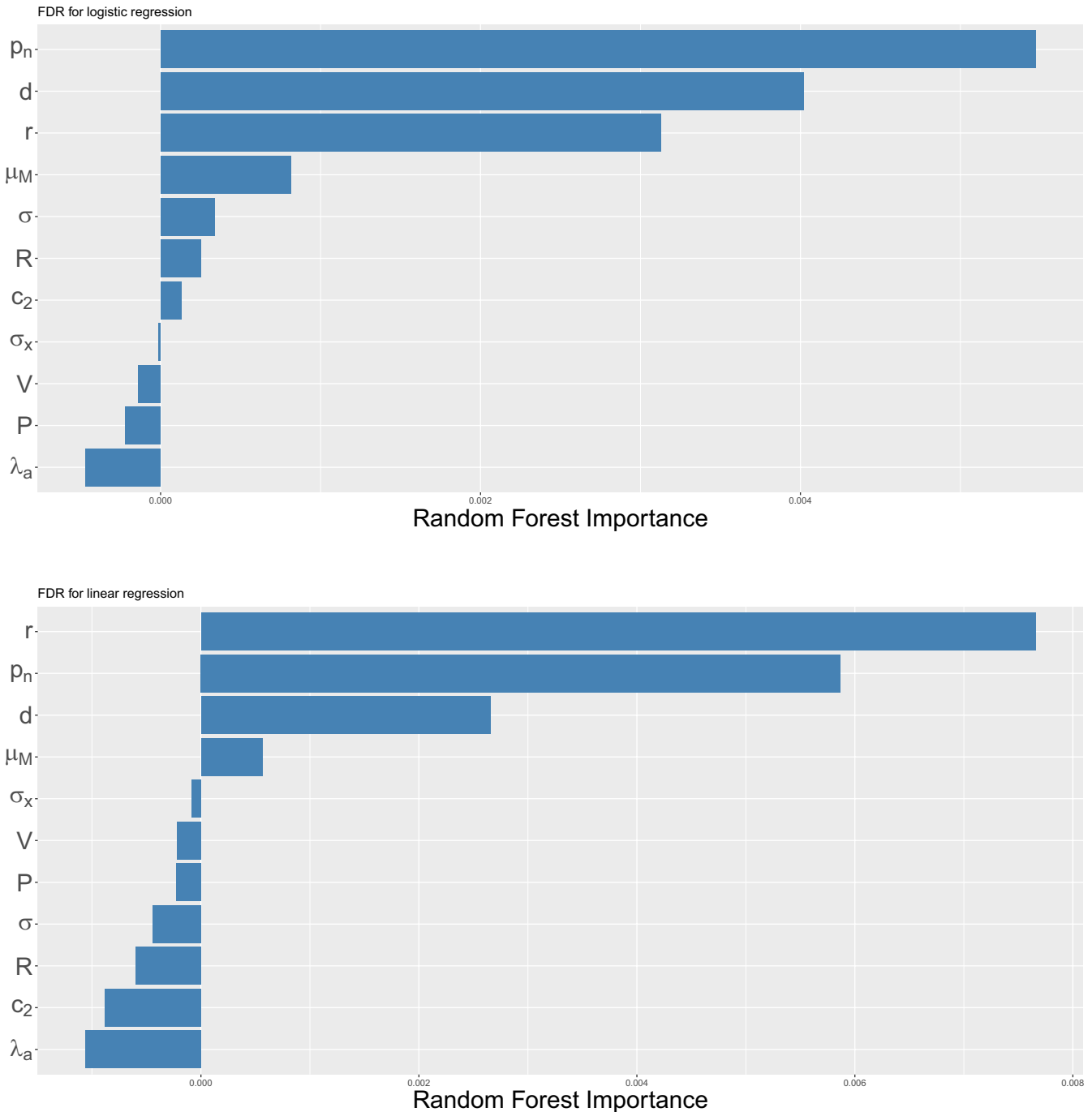


FIGURE 9 | Importance of predictor variables in random forest predicting FDR as a function of simulation parameters. Analysis was performed with linear and logistic regression of each species' data as a function of all other species data. Simulations: *Ecological simulations* with 10 species, random parameter settings, and 10,000 samples per simulation. Benjamini–Hochberg correction was applied with a false discovery control level of 0.05, and FDR was evaluated for direct, symmetric interactions. Random forest was trained on results from 1000 simulations.

calibration of some methods compared to others. Even within the same method and simulation set, the calibration can cause very different results depending on whether covariates are included. Additionally, the number of total inferred interactions varies dramatically between methods and simulations, so some average FDR values may have high variance between simulation runs (Appendix S2: Figure S1).

Sample size is not the only potential concern, since we observe high false discovery rates in our *ecological simulations* and in

the *covariance matrix simulations* with covariate responses even with high sample sizes. With very high sample sizes, statistical methods should perform well when the data meet their assumptions perfectly, but when the true data generating process is different from the assumptions of the method, the error rate may not go down with more data. The assumptions about the data generating process vary across methods. For example, none of the methods tested here account for uncertainty in covariates. Some are parametric methods that rely, for example, on residuals being independent and normal (Dormann et al. 2007). Many

of them do not account for residual autocorrelation in space and/or time, which may be caused either by non-equilibrium distributions or by autocorrelation of an unmeasured abiotic factor (Dormann et al. 2007; Thibaud et al. 2014). Some methods cannot detect non-linear interactions between species and covariates or between species. These are all examples of assumptions that we have violated to varying degrees in the simulations. We have shown here that given enough samples, some methods perform moderately well with some assumptions violated, but the only case in which any method is able to reliably recover interactions with a low false discovery rate is under an extremely idealized scenario.

Many simulation studies use data that assumes the same data generating process as the methods they are testing. As real data are not expected to follow the assumptions of the methods exactly, this procedure should be expected to overestimate performance. As in all simulation models, we make many assumptions about the data generating process, but in our case, some of the simulations significantly deviate from the assumptions of the inference methods. In our *ecological simulations*, we assume that populations in this simulation will grow according to a logistic growth curve with a changing carrying capacity through time. This may not be the optimal way to calculate growth rates for some environments and species (Hatton et al. 2024). Additionally, as actual ecological network structure is likely variable in different systems (Bascompte et al. 2019), we assign interactions between species randomly in our *ecological simulations*. This is likely a conservative assumption because graphs with evenly distributed neighborhoods are easier to recover than those with large hub nodes (Kurtz et al. 2015). In the *covariance matrix simulations*, interactions are instead organized in clusters of interacting species. Differences in performance due to the structure of the interaction networks between species have been explored (Landi et al. 2018), and the inference methods discussed here are often used to estimate overall characteristics of the network (Kurtz et al. 2015), so these assumptions may be significant. Additional assumptions of our *ecological simulation* model include assuming that some species traits and data characteristics are constant across time, space, and species, which is certainly not always the case. For example, we expect that covariate uncertainty, average read abundances, and species detection rates may vary in space and time and across species (Capo et al. 2021; Dolenz et al. 2024). Our simulation model also assumes that interactions between species are constant over time and space, which may not be true over evolutionary timescales. All of these assumptions are conservative, but we recognize that we also made assumptions about the data that are less conservative. For example, strong species interactions in real data may induce higher levels of correlation than we see in our models. Additionally, other characteristics of the simulation may create unrealistic dynamics that affect the performance of the inference models. However, within our chosen simulation framework, we attempt to simulate many ecological dynamics by varying the parameters widely in our simulations. Additionally, we test the methods under a simpler simulation framework (*covariance matrix simulations*) that is not based on ecological theory but is expected to minimally violate assumptions of the inference methods.

The input to some models was sedaDNA read abundances, whereas other methods take only presence–absence data. EcoCopula was the only model that takes flexible input. We find that information about species abundances improves the performance of models at detecting species interactions in our *ecological simulations*. On the other hand, in the *covariance matrix simulations*, we see more mixed results. In a few cases, methods with presence–absence data perform better than their read abundance counterpart. We believe this could be due to differences in the distribution assumed by the models and the actual distribution of read abundances. For example, linear regression assumes the residuals are Gaussian, and the actual read abundances are Poisson distributed. Additionally, the Poisson abundances have a second source of random variability from drawing the latent Gaussian variable first and then the Poisson variable for the reads. We believe that this may be realistic since the read counts may also introduce extra variability in the data as compared to treating the data as presence–absence data. Additionally, in real data, read abundance data may be biased by taxon, which may introduce confounding factors that result in spurious correlations. For example, many factors may create biases such as species-specific DNA deposition rates (Giguët-Covex et al. 2019), PCR bias (Krehenwinkel et al. 2017), or bias in species assignments (Dolenz et al. 2024). In our simulations, we did not include per-species bias in rates of DNA deposition, so it is not surprising that simulated read abundances have more information about species interactions than presence–absence data, though we do not know whether this is the case in real data. These concerns would still be applicable to some extent in presence–absence data, as detection rates may also be biased across taxa, space, and time, but the influence of these biases may be lower. We also note that in microbiome studies, reads are not assigned to taxa but rather grouped into OTUs or ASVs, which may have different properties with respect to read abundance biases (Chiarello et al. 2022). Further investigation will be needed to fully understand these effects and the circumstances under which it is effective to use read abundances as a proxy for species abundance.

The increasing availability of whole-community ecology data has enormous potential to uncover ecological dynamics, parameterize models, and make predictions about the future of ecosystems. However, data of this complexity must also be approached with caution, and benchmarking computational methods under a variety of scenarios is an essential step in understanding and interpreting results. Further investigations into the relative success of these methods and others will be needed for a complete understanding of their performance. For example, it is possible that these methods would be more successful at recovering overall network properties as opposed to specific interactions (Kurtz et al. 2015). Additionally, many other methods exist that were not tested here due to differences in data requirements and outputs, including some mechanistic models, which may perform better under certain conditions (Kearney and Porter 2009; Overcast et al. 2021). An area where sedaDNA has a lot of potential is in its ability to recover species distribution data across a wide range of taxa, including microbes, plants, and megafauna, from the same samples. However, tools to analyze data on differing taxonomic scales have been developed with differing assumptions (perhaps for good reason). In order to fully take advantage

of the potential of these data, we will need analysis frameworks that are able to handle many different ecological scenarios.

Author Contributions

F.M.C.: conceptualization, acquisition, analysis, and interpretation of the data, writing of the initial manuscript; J.K.L.: analysis of the data (JSDM-MCMC analysis), editing of the initial manuscript; R.N.: conceptualization, interpretation of the data, resources, editing of the initial manuscript.

Conflicts of Interest

The authors declare no conflicts of interest.

Data Availability Statement

The paper contains no new unpublished data. Code used to simulate data is available at <https://github.com/Fiona-MC/eDNA-sims-pub>.

References

- Akesson, A., A. Curtsdotter, A. Eklöf, B. Ebenman, J. Norberg, and G. Barabás. 2021. "The Importance of Species Interactions in Eco-Evolutionary Community Dynamics Under Climate Change." *Nature Communications* 12: 4759.
- Alsos, I. G., V. Boussange, D. P. Rijal, et al. 2024. "Using Ancient Sedimentary DNA to Forecast Ecosystem Trajectories Under Climate Change." *Philosophical Transactions of the Royal Society, B: Biological Sciences* 379: 20230017.
- Banerjee, S., K. Schlaeppli, and M. G. A. Van Der Heijden. 2018. "Keystone Taxa as Drivers of Microbiome Structure and Functioning." *Nature Reviews Microbiology* 16: 567–576.
- Bascompte, J., M. B. García, R. Ortega, E. L. Rezende, and S. Pironon. 2019. "Mutualistic Interactions Reshuffle the Effects of Climate Change on Plants Across the Tree of Life." *Science Advances* 5: eaav2539.
- Beng, K. C., and R. T. Corlett. 2020. "Applications of Environmental DNA (eDNA) in Ecology and Conservation: Opportunities, Challenges and Prospects." *Biodiversity and Conservation* 29: 2089–2121.
- Benjamini, Y., and Y. Hochberg. 1995. "Controlling the False Discovery Rate – A Practical and Powerful Approach to Multiple Testing." *Journal of the Royal Statistical Society, Series B* 57: 289–300.
- Bharti, R., and D. G. Grimm. 2021. "Current Challenges and Best-Practice Protocols for Microbiome Analysis." *Briefings in Bioinformatics* 22: 178–193.
- Blanchet, F. G., K. Cazelles, and D. Gravel. 2020. "Co-Occurrence Is Not Evidence of Ecological Interactions." *Ecology Letters* 23, no. 7: 1050–1063. <https://doi.org/10.1111/ele.13525>.
- Capo, E., C. Giguët-Covex, A. Rouillard, et al. 2021. "Lake Sedimentary DNA Research on Past Terrestrial and Aquatic Biodiversity: Overview and Recommendations." *Quaternary* 4: 6.
- Chen, W., and G. F. Ficetola. 2020. "Numerical Methods for Sedimentary-Ancient-DNA-Based Study on Past Biodiversity and Ecosystem Functioning." *Environmental DNA* 2: 115–129.
- Chiarello, M., M. McCauley, S. Villéger, and C. R. Jackson. 2022. "Ranking the Biases: The Choice of OTUs vs. ASVs in 16S rRNA Amplicon Data Analysis has Stronger Effects on Diversity Measures Than Rarefaction and OTU Identity Threshold." *PLoS One* 17: e0264443.
- Clark, N. J., and K. Wells. 2023. "Dynamic Generalised Additive Models (DGAMs) for Forecasting Discrete Ecological Time Series." *Methods in Ecology and Evolution* 14: 771–784.

- Claussen, J. C., J. Skievecičienė, J. Wang, et al. 2017. "Boolean Analysis Reveals Systematic Interactions Among Low-Abundance Species in the Human Gut Microbiome." *PLoS Computational Biology* 13: e1005361.
- Dennis, B. 1989. *Estimation and Analysis of Insect Populations*, 219–238. Springer.
- Dolenz, S., T. van der Valk, C. Jin, et al. 2024. "Unravelling Reference Bias in Ancient DNA Datasets." *Bioinformatics* 40: btae436.
- Dormann, C. F., S. J. Schymanski, J. Cabral, et al. 2012. "Correlation and Process in Species Distribution Models: Bridging a Dichotomy." *Journal of Biogeography* 39: 2119–2131.
- Dormann, F., J. M. McPherson, M. B. Araújo, et al. 2007. "Methods to Account for Spatial Autocorrelation in the Analysis of Species Distributional Data: A Review." *Ecography* 30: 609–628.
- Dussex, N., N. Bergfeldt, V. de Anca Prado, et al. 2021. "Integrating Multi-Taxon Palaeogenomes and Sedimentary Ancient DNA to Study Past Ecosystem Dynamics." *Proceedings of the Royal Society B: Biological Sciences* 288: 20211252.
- Elith, J., and C. H. Graham. 2009. "Do They? How Do They? WHY Do They Differ? On Finding Reasons for Differing Performances of Species Distribution Models." *Ecography* 32: 66–77.
- Friedman, J., and E. J. Alm. 2012. "Inferring Correlation Networks From Genomic Survey Data." *PLoS Computational Biology* 8: e1002687.
- Gelman, A., and D. B. Rubin. 1992. "Inference From Iterative Simulation Using Multiple Sequences." *Statistical Science* 7: 457–472. <https://projecteuclid.org/journals/statistical-science/volume-7/issue-4/Inference-from-Iterative-Simulation-Using-Multiple-Sequences/10.1214/ss/1177011136.full>.
- Giguët-Covex, C., G. F. Ficetola, K. Walsh, et al. 2019. "New Insights on Lake Sediment DNA From the Catchment: Importance of Taphonomic and Analytical Issues on the Record Quality." *Scientific Reports* 9: 14676.
- Hatton, I. A., O. Mazzarisi, A. Altieri, and M. Smerlak. 2024. "Diversity Begets Stability: Sublinear Growth and Competitive Coexistence Across Ecosystems." *Science* 383: eadg8488.
- Hui, F. K. C., D. I. Warton, S. D. Foster, and C. R. Haak. 2023. "Spatiotemporal Joint Species Distribution Modelling: A Basis Function Approach." *Methods in Ecology and Evolution* 14: 2150–2164.
- Kearney, M., and W. Porter. 2009. "Mechanistic Niche Modelling: Combining Physiological and Spatial Data to Predict Species' Ranges." *Ecology Letters* 12: 334–350.
- Krehenwinkel, H., M. Wolf, J. Y. Lim, A. J. Rominger, W. B. Simison, and R. G. Gillespie. 2017. "Estimating and Mitigating Amplification Bias in Qualitative and Quantitative Arthropod Metabarcoding." *Scientific Reports* 7: 17668.
- Kurtz, Z., C. Mueller, E. Miraldi, and R. Bonneau. 2023. "SpiecEasi: Sparse Inverse Covariance for Ecological Statistical Inference R Package Version 1.1.2, Commit d6bc1273211fef632b86f65b4b98290805b9ab5b." <https://github.com/zdk123/SpiecEasi>.
- Kurtz, Z. D., C. L. Müller, E. R. Miraldi, D. R. Littman, M. J. Blaser, and R. A. Bonneau. 2015. "Sparse and Compositionally Robust Inference of Microbial Ecological Networks." *PLoS Computational Biology* 11: e1004226.
- Landi, P., H. O. Minoarivelo, Å. Brännström, C. Hui, and U. Dieckmann. 2018. "Complexity and Stability of Ecological Networks: A Review of the Theory." *Population Ecology* 60: 319–345.
- Lindgren, F., and H. Rue. 2015. "Bayesian Spatial Modelling With R – INLA." *Journal of Statistical Software* 63: 1–25. <https://doi.org/10.18637/jss.v063.i19>.
- Overcast, I., M. Ruffley, J. Rosindell, et al. 2021. "A Unified Model of Species Abundance, Genetic Diversity, and Functional Diversity Reveals the Mechanisms Structuring Ecological Communities." *Molecular Ecology Resources* 21: 2782–2800.

- Pichler, M., and F. Hartig. 2021. "A New Joint Species Distribution Model for Faster and More Accurate Inference of Species Associations From Big Community Data." *Methods in Ecology and Evolution* 12: 2159–2173.
- Plummer, M. 2003. "JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs Sampling." In *Proceedings of the 3rd international workshop on distributed statistical computing* 124 (Vienna, Austria), 1–10.
- Plummer, M., N. Best, K. Cowles, and K. Vines. 2006. "CODA: Convergence Diagnosis and Output Analysis for MCMC." *R News* 6: 7–11.
- Pollock, L. J., R. Tingley, W. K. Morris, et al. 2014. "Understanding Co-Occurrence by Modelling Species Simultaneously With a Joint Species Distribution Model (JSDM)." *Methods in Ecology and Evolution* 5: 397–406.
- Popovic, G. C., D. I. Warton, F. J. Thomson, F. K. C. Hui, and A. T. Moles. 2019. "Untangling Direct Species Associations From Indirect Mediator Species Effects With Graphical Models." *Methods in Ecology and Evolution* 10: 1571–1583.
- Schliep, E. M., N. K. Lany, P. L. Zarnetske, et al. 2018. "Joint Species Distribution Modelling for Spatio-Temporal Occurrence and Ordinal Abundance Data." *Global Ecology and Biogeography* 27: 142–155.
- Su, Y.-S., and M. Yajima. 2021. "R2jags: Using R to Run 'JAGS' R Package Version 0.7-1." <https://cran.r-project.org/package=R2jags>.
- R Core Team. 2024. "R: A Language and Environment for Statistical Computing." Vienna, Austria. <https://www.R-project.org/>.
- Thibaud, E., B. Petitpierre, O. Broennimann, A. C. Davison, and A. Guisan. 2014. "Measuring the Relative Effect of Factors Affecting Species Distribution Model Predictions." *Methods in Ecology and Evolution* 5: 947–955.
- Wagg, C., K. Schlaeppli, S. Banerjee, E. E. Kuramae, and M. G. A. Van Der Heijden. 2019. "Fungal-Bacterial Diversity and Microbiome Complexity Predict Ecosystem Functioning." *Nature Communications* 10: 4841.
- Wang, S. C., and C. R. Marshall. 2016. "Estimating Times of Extinction in the Fossil Record." *Biology Letters* 12: 20150989.
- Wang, Y., M. W. Pedersen, I. G. Alsos, et al. 2021. "Late Quaternary Dynamics of Arctic Biota From Ancient Environmental Genomics." *Nature* 600: 86–92.
- Weiss, S., W. van Treuren, C. Lozupone, et al. 2016. "Correlation Detection Strategies in Microbial Data Sets Vary Widely in Sensitivity and Precision." *ISME Journal* 10: 1669–1681.
- Wilkinson, D. P., N. Golding, G. Guillera-Arroita, R. Tingley, and M. A. McCarthy. 2019. "A Comparison of Joint Species Distribution Models for Presence–Absence Data." *Methods in Ecology and Evolution* 10: 198–211.
- Williams, J. W., T. L. Spanbauer, P. D. Heintzman, et al. 2023. "Strengthening Global-Change Science by Integrating aeDNA With Paleocoinformatics." *Trends in Ecology & Evolution* 38: 946–960.
- Wright, M. N., and A. Ziegler. 2017. "Ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R." *Journal of Statistical Software* 77: 1–17. <https://doi.org/10.18637/jss.v077.i01>. <http://www.jstatsoft.org/v77/i01/>.
- Yuan, M. M., X. Guo, L. Wu, et al. 2021. "Climate Warming Enhances Microbial Network Complexity and Stability." *Nature Climate Change* 11: 343–348.
- Zimmermann, H. H., K. R. Stooß-Leichsenring, V. Dinkel, et al. 2023. "Marine Ecosystem Shifts With Deglacial Sea-Ice Loss Inferred From Ancient DNA Shotgun Sequencing." *Nature Communications* 14: 1650.
- Zurell, D., L. J. Pollock, and W. Thuiller. 2018. "Do Joint Species Distribution Models Reliably Detect Interspecific Interactions From Co-Occurrence Data in Homogenous Environments?" *Ecography* 41: 1812–1819.

Supporting Information

Additional supporting information can be found online in the Supporting Information section.