1	Filling Monitoring Gaps for Data-deficient Species Using Annual
2	Occupancy Predictions from Co-occurrence Models
3	
4	HY Chung ^a , DK Lee ^b and JE Losey ^{a1}
5	^a Department of Entomology at Cornell University, Ithaca, NY 14853
6	^b Multi Campus, Seoul, South Korea 06220
7	
8	Article impact statement
9 10	Co-occurrence-based annual distribution modeling enables IUCN assessment for rare taxa, filling monitoring gaps without structured surveys.
11	Keywords
12	Data Deficient, Extinction Risk, Citizen Science, Multi-source, Data Biases, Rare Species,
13	Insect Extinction, IUCN Red List
14	Word counts
15	5,239
16	Conflict of interest statement
17	The authors declare no conflict of interest.
18	Author contribution
19	Conceptualization Hyun Yong Chung, Dae Kyung Lee, John Losey; Data curation Hyun
20	Yong Chung, Dae Kyung Lee; Formal analysis Hyun Yong Chung; Funding acquisition John
21	Losey; Investigation Hyun Yong Chung, Dae Kyung Lee; Methodology Hyun Yong Chung,
22	Dae Kyung Lee; Project administration Hyun Yong Chung; Resources Hyun Yong Chung,
23	John Losey; Software Dae Kyung Lee, Hyun Yong Chung; Supervision Hyun Yong Chung;
24	Validation Hyun Yong Chung; Visualization Hyun Yong Chung; Writing – original draft
25	Hyun Yong Chung; Writing – review and editing Hyun Yong Chung, John Losey.
26	Acknowledgements

¹Corresponding author's email: jel27@cornell.edu

- We extend our gratitude to LM Guzman, BW Suh, and JY You for their thoughtful reviews and
- 28 comments.
- 29 ORCID

30 Hyun Yong Chung https://orcid.org/0000-0001-7698-8105

Filling Monitoring Gaps for Data-deficient Species Using Annual Occupancy Predictions from Co-occurrence Models

Abstract

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

Fragmented surveys and limited monitoring have excluded most invertebrates from conservation policy. We present a fill-in framework that uses species distribution models (SDMs) to reconstruct missing annual trends—not to extrapolate trends, but to fill them in. Instead of filtering data-sparse regions or years or relying on static environmental variables, we used co-occurrence patterns (COP) as variables, to capture year-to-year assemblage shifts. COP variables enabled annual prediction at all recorded sites from multisource, presence-only, sparse data. When applied to four rare native ladybugs across North America (2007–2021), COP models exceeded reliability benchmarks (Accuracy ≥ 0.70 , AUC ≥ 0.70 , Kappa ≥ 0.40 , Brier ≤ 0.25) across standard 7:3 splits, cross-source and cross-period validations. Annual predictions were robust to temporal biases from variation in data volume and source composition. Multiple regression indicated negligible effects of those biases on reconstructed trends. Predicted decadal declines (9-31%) closely aligned with independent regional longterm monitoring, operationalizing IUCN Red List classification (from least concern to vulnerable) in the absence of standardized monitoring. By converting fragmented observations—primarily from citizen science—into reliable annual trend estimates, the fill-in approach extends extinction-risk assessment to data-deficient taxa long excluded from conservation frameworks.

Introduction

Most invertebrate species—despite ongoing declines—remain invisible to conservation action, not because they are safe but because they are silent in the data. The International Union for Conservation of Nature (IUCN) Red List requires trend estimates over the past decade or three generations for declining species (IUCN, 2024). However, standardized monitoring programs are scarce or limited in scope for these species (Estes et al., 2018; Bayraktarov et al., 2019) and available records are typically sparse, presence-only, and opportunistic. Such data even fail to meet the requirements of current trend models (Harvey et al., 2020), which assume ordinal abundance (e.g., Newson et al., 2015; Inamine et al., 2016; Schultz et al., 2017; Martín et al., 2021), checklist-based surveys (Walker & Taylor, 2017; LeCroy et al., 2020), metrics of survey effort (e.g., Szabo et al., 2010; Isaac et al., 2014; Kamp et al., 2016; Horns et al., 2018;

Fink et al., 2020), and repeat visit protocols (e.g., MacKenzie et al., 2002; 2006; Kéry et al., 2010; van Strien et al., 2013; Altwegg & Nichols, 2019). For example, Bayesian occupancy models (OM) are effective tools for trend estimation, yet in North America they often entail ≥ 10,000 km² grid cells and 10–20-year intervals to achieve adequate data density for bees and dragonflies (Soroye et al., 2020; Jackson et al., 2022). Without methods capable of estimating temporal trends from the data available, underrepresentation of invertebrates in global conservation frameworks will persist (Montgomery et al., 2020; Jönsson et al., 2021).

This study reconstructs annual occupancy trajectories by filling monitoring gaps with annual predictions from species distribution models (SDMs) using variables derived from cooccurrence patterns (COP). Previous studies have used SDMs to predict occupancy at unsurveyed locations by modeling correlative relationships between known occurrence records and environmental or biotic variables (Olden et al., 2008; Zimmermann et al., 2010; Franklin, 2013). These relationships have been modeled via recursive binary partitioning (decision-tree algorithms), entropy maximization under environmental constraints (maximum entropy modeling), or multivariate probit regression (joint species distribution models; JSDMs), among other approaches. Such predictions can help fill monitoring gaps, and SDMs perform well with small sample sizes (Hernandez et al., 2006; Wisz et al., 2008; Luan et al. 2020) and presenceonly data (Hernandez et al., 2006; Wisz et al., 2008; Mi et al., 2017; Robinson et al., 2018). However, most applications target single-time distributions or long-term range shifts (Tingley & Beissinger, 2009; Svancara et al., 2019). They primarily rely on static variables, such as environmental and geographical proxies, making them ill-suited to trend estimation. Rare attempts at fine-temporal prediction require structured data and extensive data thinning (Svancara et al., 2019; Fink et al., 2020), which is impractical for under-monitored invertebrates.

To predict annual occupancy at historical sites, we novelly operationalized co-occurrence patterns (COP) as numerical variables for machine learning classifiers that use decision trees (COP-ML). In community ecology, COP denotes the observed frequency with which species co-occur across locations. Prior SDM studies suggest that such patterns can encode biotic interactions (Pollock et al., 2014) and community-level environmental responses (Kissling et al., 2012). Extending these ideas, we propose that annually updated COP vectors could also track habitat and community shifts at a site more rapidly than static or slowly varying variables. Community composition may change before environmental changes are detectable and can

immediately respond to non-environmental drivers. Here, we represented COP as a site–year vector of the annual frequencies of non-target species within a fixed radius. Decision-tree classifiers then learned rule-based partitions linking target species occupancy to both individual variables and their joint patterns (e.g., "if variable $A \ge a$ and variable $B \le b$, then ...").

This study evaluates whether COP-ML can produce accurate predictions of annual occupancy. It further tests whether the model remains robust to temporal and structural biases in multisource, sparse and opportunistic data typical of under-monitored species. With the expansion of citizen science, opportunistic observations increasingly dominate in those taxa (Kissling et al., 2018; Knape et al., 2022) but also introduce temporal and structural biases (Isaac et al., 2014; Guzman et al., 2021; Larsen & Shirey, 2021). Mitigation strategies, such as data thinning or quality-based filtering (Wisz et al., 2008; Isaac et al., 2014; Kamp et al., 2016; Zizka et al., 2021; Van Eupen et al., 2021), can exacerbate sparsity and also preclude estimation of the absolute trend across the entire range as required under the IUCN Red List framework. Reliable annual trend estimation therefore requires methods that tolerate biased data while preserving temporal signals.

To test robustness to temporal bias, we examine whether models trained on one period generalize to others (Martínez-Minaya et al., 2018), a necessity given that citizen-science data are disproportionately concentrated in recent years (Geldmann et al., 2016). To test robustness to structural bias, we assess whether models trained on one survey type can predict others. Multisource integration is often indispensable for data-poor species (Fletcher et al., 2019; Miller et al., 2019; Isaac et al., 2020), yet inconsistent methods, even among citizen-science platforms (Gardiner et al., 2012), can introduce systematic errors (Cheney et al., 2013). For example, citizen science tends to target urban areas, whereas ecologists' surveys more often cover natural or semi-natural areas (Geldmann et al., 2016). Without such robustness, temporal trends risk reflecting shifts in prevailing survey types rather than biological change (Pagel et al., 2014; Knape et al., 2022).

We test three hypotheses: first, can COP-ML distinguish target species' presence from absence using annual COP variables? Second, can annual COP variables generalize across time periods? Third, can these variables generalize across survey methods? Accordingly, we ran three tests: standard 7:3 split evaluation, temporal generalization, and structural generalization. We then applied COP-ML to four native North American ladybugs and generated annual

occupancy estimates for 2007–2021 at all recorded sites. From these, we reconstructed trajectories and quantified 10-year reductions to operationalize IUCN Red List categories.

Methods

We asked whether COP (co-occurrence pattern) variables can yield accurate annual occupancy despite temporal and structural biases. In the step for data assembly and labeling, we compiled presence-only records from multiple sources (2007–2021) for all North American ladybugs with minimal filtering, labeling four native species as presences and two invasive competitors as pseudo-absences. In the step for variable construction, we built COP variables for each site—year labeled as presence or pseudo-absence by counting non-target species within a fixed radius and scaling counts. In the step for model training, we trained XGBoost decisiontree classifiers on these variables, which served as predictors in the model. In the step for generalization and evaluation, we assessed model performance with six metrics under three settings: cross-source (structural) generalization, cross-period (temporal) generalization, and standard 7:3 split evaluation. For each test, we summarized the distribution of performance from 2,500 resampled datasets. Finally, in the annual prediction and trend analysis step, COP-ML trained on the full dataset predicted annual occupancy of the four target species; an ensemble of resampled models (majority vote) produced the final yearly presence-absence at all historical sites for 2007–2021. From these predictions, we estimated decadal reduction rates under IUCN Criterion A and tested their significance.

Target species

Four native ladybug species—Coccinella novemnotata Herbst, 1793; Coccinella transversoguttata Faldermann, 1835; Adalia bipunctata (Linnaeus, 1758); and Hippodamia parenthesis (Say, 1824)—once dominated North American ladybug communities, thriving across diverse habitats and prey types (Losey 2007; 2012; taxonomy hereafter follows Gordon (1985) and Gordon & Vandenberg (1991)). Since the mid-1980s, their relative abundance in collections has dropped to 1–5% of former levels (Harmon et al., 2007) due to the newly-established competitors Coccinella septempunctata Linnaeus, 1758 and Harmonia axyridis (Pallas, 1773) (Wheeler & Hoebeke 1995; Harmon et al., 2007). These introduced species now dominate and alter community structure, reducing diversity and abundance continent-wide (Petersen & Losey, 2024). Estimating reduction rates and extinction risks for the natives is challenging given their current low densities and broad distributions (Wheeler & Hoebeke,

1995; Hesler et al., 2004; Harmon et al., 2007). Addressing this requires integrating multisource data across periods, regions, and survey methods while addressing inherent biases.

Occurrence data

We compiled Ladybug records from multiple sources: three citizen-science sources (The lost ladybug project, 2021; iNaturalist, 2021; BugGuide.Net, 2021), one museum website (NCSU, 2021), and three metadata platforms (GBIF, 2021; BISON, 2021; IdigBio, 2021; Appendix S1). Of the citizen-science, iNaturalist and BugGuide.Net relied on user identifications, whereas The Lost Ladybug Project relied on experts. We further verified identifications of the target species from iNaturalist and BugGuide.Net.

To assess how COP variables address biases, we applied minimal preprocessing: Records were restricted to the U.S. (excluding Alaska and Hawaii) and parts of Canada (Manitoba, Ontario, Saskatchewan, British Columbia, Alberta, Quebec) for 2007–2021. Only adult forms identified to species level were retained. Where available (89% of records), we restricted GPS accuracy to 1 km. We removed duplicate records identical in species, year, and GPS. We then computed descriptive statistics to reveal temporal and structural inconsistencies in the compiled dataset.

The dataset included 188,644 records of 353 ladybug species from 85 sources, with 324 records for *C. novemnotata*, 510 for *C. transversoguttata*, 732 for *H. parenthesis*, and 1,426 for *A. bipunctata*, which were labeled as presence.

Pseudo-absence

When explicit absence records are unavailable, pseudo-absences are often drawn by randomly sampling coordinates from all other species' records in a dataset (Robinson et al., 2018). Here we instead used records of the introduced competitors, *C. septempunctata* and *H. axyridis*. First, they competitively exclude the target species. Their occurrence within an 18-km radius without the targets was assumed to be a logical proxy for absence, reflecting a reshaped COP after local displacement. Second, their dominance (61% of our dataset) meant that traditional random sampling would largely select them anyway, minimizing methodological deviation. Therefore, we used competitors as pseudo-absences to encode a local-displacement hypothesis in COP variables.

We pooled 10,000 pseudo-absence points by subsampling from states or provinces in proportion to the regional frequency of the four targets' presence. Without this proportionality

filter, apparent accuracy increased but models relied more on geographical variables (e.g., *Coleomegilla maculata* (De Geer, 1775) concentrated in eastern regions), which are insensitive to temporal changes. In the variable construction step, we combined the entire pool with presence records. In the training and testing step, absences were resampled from this pool to balance classes.

Variables

Direct and indirect competitions shape ladybug assemblages, where the dominance of newly-established species drives niche differentiation (Petersen & Losey, 2024) and avoidance behaviors in native species (Elliott et al., 1996; Hesler & Kieckhefer 2008; Mukwevho et al., 2017). These changes can be immediate since ladybugs, as highly mobile predators, engage in long-distance interactions (e.g., H. axyridis, 442 km/year; McCorquodale, 1998) and actively forage across habitats (Woltz & Landis, 2013). 18 km is the commonly reported dispersal distance of this group (Jeffries et al., 2013; COSEWIC, 2016a; 2016b). Here, we represented COP as counts $c_{s,t,j}$ of cooccurring non-target ladybug species j within an 18 km radius r of site s in year t, where each site-year (s,t) was labeled as presence (= 1) or absence (= 0) of the target species. This is expressed as Equation (1):

$$c_{s,t,j} = \sum_{u} 1$$
, if
$$\begin{cases} \text{species}(u) = j \\ \text{year}(u) = t \\ \text{distance}(s, u) \le r \end{cases}$$
 (1)

Where u denotes a single georeferenced record (coordinate, year, species). Next, for each species-year (j,t), we trimmed outliers in the distribution of $c_{s,t,j}$ outside [Q1 - 1.5IQR, Q3 + 1.5 IQR] to reduce distributional bias during scaling. To maintain consistency of variable vectors, outlier detection in a single variable resulted in exclusion of the entire site—year record. This mainly affected the pseudo-absence pool, removing only 0–12 presence records per target species ($\leq 0.008\%$). We then min-max scaled $c_{s,t,j}$ per (j,t) to adjust for species-specific overreporting and temporal variations in observation efforts, as expressed as Equation (2):

$$x_{s,t,j} = \begin{cases} \frac{c_{s,t,j} - \min_{s \in S_{t,j}} c_{s,t,j}}{\max_{s \in S_{t,j}} c_{s,t,j} - \min_{s \in S_{t,j}} c_{s,t,j}}, & \text{if } \max_{s \in S_{t,j}} c_{s,t,j} > \min_{s \in S_{t,j}} c_{s,t,j}, \\ 0, & \text{otherwise.} \end{cases}$$

We excluded environmental variables to avoid multicollinearity with COP variables, as their effects are expected to be partially embedded in co-occurrence patterns (Kissling et al., 2012). This also aligns with our focus on short-term temporal interpolation. We retained 85

species with at least 30 co-occurrences with a target species, excluding unidentified 'sp.' We applied additional screening using forward regressions (p < 0.05) and variance inflation factors (< 10). Although prior feature selection rarely improves performance of decision trees, we applied prescreening to limit ecologically implausible variables and to improve interpretability. To prevent leakage, we excluded target species from their own variables. Finally, we ranked the top 15 key variables using SHapley Additive exPlanations (SHAP) values, which assess feature importance in predictions from preliminary loops. The selected predictors $x_{s,t}$ were stacked into the matrix X and fed into models as:

$$X_{(s,t)} = (x_{s,t,j_1}, \dots, x_{s,t,j_{15}}), y_{\{(s,t)\}} \in \{0,1\}$$
 (3)

To analyze associations between each variable and target species occupancy, we calculated averages of point-biserial correlations by resampling pseudo-absence points 50 times to match presence record counts.

Development and characterization of models

We modeled associations between COP variables and targets' occupancy with XGBoost as an ensemble of decision trees:

$$F(x_i) = \sum_{m=1}^{M} f_m(x_i), \qquad \hat{p}i = \frac{1}{1 + \exp(-F(x_i))}$$
 (4)

Here, $F(x_i)$ is the raw logit score and $\hat{p}i$ is predicted probability of presence at site-year i. Each tree f_m is a set of if—then split rules to the predictor vector x_i (the COP variables); for example, x_a < Threshold split_a and x_b > Threshold split_b. Thus, tree depth ≥ 2 naturally encodes interactions among predictors. Training minimizes logistic loss with a tree-complexity penalty (Chen & Guestrin 2016).

We implemented the XGBoost package in Python. From the default hyperparameter settings, we only adjusted objective='binary:logistic' and n_estimators=1000, as our aim was to evaluate COP variables rather than optimize the model.

Across all analyses, we balanced presence and pseudo-absence records 1:1 by undersampling pseudo-absences. We then generated 50 independent datasets by resampling pseudo-absences and, within each, drew 50 random train–test splits, yielding 2,500 runs (50×10^{-2})

50). For each test scenario, we summarized and evaluated performance based on the mean and distribution across runs.

We assessed model performance with six metrics: Accuracy (correct response rate), Kappa (agreement adjusted for random chance; Cohen, 1960), Recall (true positive rate), and Precision (positive predictive rate) to measure ability to predict binary presence-absence, plus Brier score (mean squared discrepancy; Brier, 1950) and AUC (ability to rank presence over absence; Fielding & Bell, 1997) for probability quality.

Generalization

Generalization tests evaluate a model's ability to predict data distinct from training data in temporal, geographical, or source aspects (Vaughan & Ormerod, 2005), minimizing traintest autocorrelation, and demonstrate robustness when ground truth comparisons are limited (Justice et al., 1999). Our tests assessed whether COP-ML could generalize across structurally or temporally distinct data pools.

Structural Generalization: To evaluate generalizability across survey types, we trained models on presence and pseudo-absence data from citizen science (LLP, iNaturalist, BugGuide.Net) to predict institutional data from 28 institutes. We assessed differences between their COP structures using ANOSIM with Manhattan distance (Appendix S2). Presence records comprised 280 citizen-driven versus 44 institutional for *C. novemnotata*, 485 versus 25 for *C. transversoguttata*, 626 versus 116 for *H. parenthesis*, and 1,338 versus 88 for *A. bipunctata*, with institutional pseudo-absence points ranging from 416 to 510. In a separate test, we also trained models on a group dominated by open-ended citizen science and evaluated them on the program emphasizing rare species (LLP).

Temporal Generalization: For forward testing, we trained models on presence data from 2007 until the year when approximately 70% of presence was accumulated, testing on the remaining about 30%. For backward testing, we reversed this, training from 2021 backward (Appendix S2). Pseudo-absence points were selected using the same cutoff year.

Evaluation

To evaluate COP-ML's annual prediction performance, we followed standard 7:3 split test by training models on 70% of presence data and testing on the remaining 30%. Unlike the generalization tests, which restricted the scope of records, this split used the full range of

presences, providing a baseline measure of COP-ML's predictive reliability for subsequent reduction rate estimation.

Prediction on annual distributions and reduction rates

To enable consistent temporal comparisons, COP-ML predicted annual presence of target species at all historical sites in our dataset since 2007, addressing yearly observation gaps.

Prediction: We developed models as described in Development and characterization of models, but we trained them on all available presence data to improve prediction accuracy given the sparsity of records (Fielding & Bell, 1997; Rencher, 1995). A site-year was classified as occupied if a majority of the 2,500 model runs (50 pseudo-absence resamples × 50 random seeds for repeated fitting) predicted presence.

Analysis: We evaluated distributional trends under IUCN Red List Criterion A, based on changes in Area of Occupancy (AOO) and Extent of Occurrence (EOO). The AOO, calculated as the number of 4 km2 grid cells occupied by a species, reflects occupancy extent and, indirectly, population size (IUCN, 2024). The EOO, defined as the polygon enclosing all known occurrences, indicates risk dispersion across a species' range (IUCN, 2024). For Criterion A, we fitted linear regressions to predicted AOO (2007–2021) and estimated the 10-year decline (2012–2021), interpreting it as a proxy for population trends (IUCN, 2024). We tested for heteroskedasticity with Breusch–Pagan and White tests, and identified influential outliers with Cook's distance. We applied robust standard errors (HC3) to assess trend significance, followed by robust regression to estimate final AOO decline. Because robust regression may downweight abrupt changes that could represent real ecological signals, we also ran ordinary least squares (OLS) regression for comparison and to improve the reliability of trend interpretation.

Validation: To verify that AOO changes predicted by COP-ML reflect consistent temporal trends despite varying data availability, we regressed predicted AOO against time (year) while including annual volumes of each citizen-science source as covariates.

Results

Biases in multisource data

The compiled dataset from multiple sources showed both structural and temporal bias. Structural bias, arising from heterogeneous effort and methods across sources (Figure 1), appeared as unequal Efficiency (defined as the percentage of all records represented by target species): institutional data (3.5% of the compiled dataset) recorded target species 2.79 times more often than citizen-science data (96.5%). Even among citizen-science sources, LLP (5%) had Efficiency of 6.6%, 6 times that of iNaturalist where it was 1.1% (89%). Excluding LLP, the remaining dataset averaged 1.3% Efficiency (corresponding to the "lower efficiency group" in the second structural generalization). Temporal bias reflected an exponential increase in annual observations (Figure 1), with post-2014 volume being 9.61 times higher than pre-2014.

Structural and temporal generalization

We tested COP-ML for annual predictive performance and generalizability against biases through structural (cross sources) and temporal (forward/backward) generalizations. All test results met or exceeded established reliability benchmarks from prior studies: Accuracy \geq 0.70 (rule of thumb), AUC \geq 0.70 (Hosmer et al., 2013), Kappa \geq 0.40 (Landis & Koch, 1977), and Brier \leq 0.25 (Brier, 1950; Figure 2). Models trained on the citizen-science group accurately predicted presence-absence in the institution group. Likewise, models trained on the low-efficiency group generalized to the highest efficiency source, and both forward and backward temporal generalizations satisfied these benchmarks.

In the generalization from citizen-science to institutional groups, C. transversoguttata model achieved the highest performance (Accuracy, AUC, Kappa, Brier = 0.87, 0.94, 0.75, 0.11), followed by C. novemnotata (0.81, 0.85, 0.61, 0.16), H. parenthesis (0.78, 0.84, 0.55, 0.17), and A. bipunctata (0.73, 0.84, 0.46, 0.19). Analysis of similarities (ANOSIM) indicated small dissimilarities in COP predictor structures between these groups (R < 0.25, p < 0.005; Appendix S2), with species-specific R-values of 0.22 (R. transversoguttata), 0.12 (R. parenthesis), 0.06 (R. novemnotata), 0.05 (R. bipunctata), and 0.06 (pseudo-absences). Comparable results also held between the low-efficiency group, dominated by 96.3% of openended citizen science (plus 3.7% institutional), and the highest-efficiency source, a single citizen-science program emphasizing rare species.

In temporal generalization, *C. transversoguttata*, *C. novemnotata*, and *H. parenthesis* maintained consistent performance scores in forward and backward tests. In contrast, *A. bipunctata* showed a 7% decrease in backward Accuracy, with Recall (true positives among actual positives) increasing 2% and Precision (true positives among predicted positives) decreasing 13%. This indicates that the model trained on recent data classified a wider set of habitat conditions as occupied than were historically. In a supplementary analysis, we shifted

the split to an earlier point (pre-2012) to further exclude recent records from the test set, expanding the training set to 87% of records. This intensified the tendency: Recall rose 11% and Precision fell 18%, while accuracy remained unchanged. This suggests that its recent co-occurrence patterns encompassed conditions not occupied in the past, unlike the other three species.

Evaluation of the developed models

We assessed COP-ML classifiers, trained on 70% of the dataset and tested on 30%, against the benchmarks. All models exceeded them; even the lowest-performing species achieved Accuracy, Precision, and Recall ≥ 0.75 , AUC ≥ 0.87 , Kappa ≥ 0.57 , and Brier ≤ 0.15 (Figure 2). *C. transversoguttata* (with 510 presence datapoints) performed best, followed by *C. novemnotata* (324), *A. bipunctata* (1,438), and *H. parenthesis* (742).

Predicted reduction rates and conservation status

Given COP-ML's demonstrated performance in tests, we generated annual occupancy predictions for 2007–2021 at all historical sites of the target species in our dataset to fill monitoring gaps and enable consistent temporal comparisons (Figure 3; Appendix S6).

Area of occupancy (AOO) declined significantly over time for all species (p < 0.05 for OLS and robust SE; B: -56.2 to -13.0; R^2 : 0.50 to 0.83; all 95% CIs excluded zero). Although we detected heteroskedasticity in A. bipunctata and influential points in A. bipunctata and H. parenthesis, slope (B), R^2 and CIs were similar across OLS, OLS with robust SE, and robust regression, indicating that regression-estimated reduction rates are reliable (Appendix S3).

Predictions indicated three species were threatened by declines in North America (Appendix S4). Occupancy measured as AOO, an indicator of occupied area and an acceptable proxy of population size (IUCN, 2024), decreased from 2007 to 2021: 1,962 km² for *H. parenthesis*, 584 km² for *A. bipunctata*, and 480 km² each for *C. novemnotata* and *C. transversoguttata*. Under IUCN Red List Criterion A (10-year reduction), estimated rates were 31% for *H. parenthesis* (vulnerable), 15% for *A. bipunctata* (near threatened), 15% for *C. novemnotata* (near threatened), and 9% for *C. transversoguttata* (least concern; Figure 4).

Extent of occurrence (EOO), a proxy for spatial buffering against extinction risk (IUCN, 2024), declined most in *C. transversoguttata*. Despite its least concern classification here, this contraction suggests reduced spatial resilience with ongoing decline.

Multiple linear regression confirmed that time (year) was a significant predictor of AOO declines across all species, whereas annual volumes of each citizen-science source showed no evidence of statistically or practically meaningful effect (Appendix S5).

Variable importance and correlation

The importance metrics of variables (SHAP values) and point-biserial correlations showed positive co-occurrence patterns among *C. novemnotata*, *C. transversoguttata*, and *H. parenthesis*. Their occurrences were positively correlated and mutually informative in each model (Figure 5). In contrast, the two competitors *H. axyridis* and *C. septempunctata* were negatively correlated with these natives, and they were ranked among the top variables by SHAP. *A. bipunctata* was the exception: it showed positive correlations with both newlyestablished species. The common native *Hippodamia convergens* Guérin-Méneville, 1842 (third most abundant in the dataset) also correlated positively with three natives and ranked highly in all models, except *H. parenthesis*.

Discussion

Rationale for estimated reduction rates

This study provides the first continent-wide estimates of decadal declines for *C. novemnotata*, *C. transversoguttata*, *A. bipunctata*, and *H. parenthesis* based on annual occupancy predictions. Earlier studies from the 1980s–1990s reported steep relative abundance declines of 95–99% (rescaled from Harmon et al., 2007). By contrast, our more moderate reduction rates from 2007–2021 likely reflect the post-establishment phase in which the new competitors had already become dominant.

Several lines of evidence support the plausibility of these more moderate rates. Historical records indicate that the most acute declines occurred shortly after the establishment of *C. septempunctata* and *H. axyridis* in North America (Colunga-Garcia & Gage, 1998; Bahlai et al., 2015). Meanwhile, subsequent regional studies suggest that declines have plateaued or transitioned into a chronic, low-intensity phase (Turnock et al., 2003; Elton, 2000; Strayer et al., 2006; Harmon et al., 2007; Hesler & Kieckhefer, 2008), with no further sharp reductions observed (Alyokhin & Sewell, 2004; Bahlai et al., 2015). Such stabilization may reflect community-level reequilibration, resistance in remnant populations, or the persistence of spatial refuges (Evans, 2000; 2004; Evans et al., 2011).

Standardized long-term monitoring in Michigan (2007–2019) corroborates this interpretation, indicating 10-year declines of 37% for *H. parenthesis* and 20% for *A. bipunctata* (KBS LTER; https://lter.kbs.msu.edu/datatables/67). These trends are obtained from linear regressions of sticky-trap captures normalized by the number of survey spots to control for effort. Although limited to a single site, these local declines align closely with our continent-wide estimates (31% and 15%, respectively), suggesting that our predictions with COP-ML provide an ecologically realistic baseline for conservation assessments.

Interpretation of COP variables

Annual prediction accuracy may imply that yearly COP variables capture time-responsive ecological signals (interaction structure, habitat turnover) that static climate or land-cover variables may miss at an annual resolution. Despite attenuation bias in our noisy, heterogeneous dataset that likely damp effect sizes, correlations between key variable species and targets generally aligned with known ecological associations and were also prioritized by the models (Figure 5).

For instance, *C. novemnotata* and *C. transversoguttata* showed among the strongest positive correlations, consistent with overlapping habitat use and resource preferences (Hesler et al., 2009). In contrast, the competitors *H. axyridis* and *C. septempunctata* were negatively correlated with the three target natives and ranked among the top SHAP variables, consistent with well-documented competitive displacement (Wheeler & Hoebke, 1995; Harmon et al., 2007; Petersen & Losey, 2024).

By contrast, *A. bipunctata* showed positive correlations with both newly-established species, likely reflecting macro-scale overlap in arboreal habitat use with *H. axyridis* (Coderre et al., 1995; Koch, 2003; Omkar & Pervez, 2005; Hentley et al., 2016) and competitive coexistence with *C. septempunctata* in Europe, where both are native (Honěk 1985; Nedvěd 1999). However, this does not rule out competitive exclusion at finer spatial scales that may fall below our 18 km COP radius (Kajita et al., 2000; Kajita et al., 2006; Soares & Serpa, 2007). Notably, coexistence signals persisted despite competitor-based pseudo-absences. This alleviates concerns that pseudo-absence sampling introduced shortcut artifacts—for example, if neighborhoods of pseudo-absence points enriched with competitor records might automatically imply target absence.

Lower-ranked COP variables involved species pairs with little documented interaction and may act as proxies for geography or environment.

Strength and limitation of COP

The robustness of COP-ML across periods and sources suggests that COP variables could encode latent distributional constraints within noisy, opportunistic datasets—particularly when interactions strongly shape them (trophic dynamics, habitat filtering, or competition; Pollock et al., 2014). In our system, the prolonged competition between native and newly-established ladybugs has reshaped communities (Harmon et al., 2007; Petersen & Losey, 2024), and COP variables appear to capture these patterns.

For three species, performance metrics were similar in forward and backward predictions, indicating limited sensitivity to temporal fluctuations in data quantity and quality. Analysis of similarities (ANOSIM) reached consistent conclusions (Appendix S2). Consistent co-occurrence signals within noisy datasets provide a stable basis for annual occupancy prediction. One possible explanation is that COP variables capture relational signals, which are less sensitive to sampling noise than single-species occurrence rates (Tikhonov et al., 2017; Johnston et al., 2017).

However, this strength depends on the temporal stability of co-occurrence patterns (Tikhonov et al., 2017). COP-ML declined in performance when backcasting *A. bipunctata*, whose habitat selection was reported to have shifted under post-invasion habitat compression (Bahlai et al., 2015). The model trained on recent COP tended to overpredict past suitability by classifying historically unoccupied conditions as suitable—evidenced by higher recall than precision. This underscores a limitation: when biotic interactions change, the assumption of time-invariant COP may fail.

Our COP-ML generalized across heterogeneous datasets—from open-ended to targeted rare-species citizen science—while showing small dissimilarities in COP values (Appendix S2). Although opportunistic data are often viewed skeptically and has limited its utility (Isaac & Pocock, 2015; Steen et al., 2019), multisource integration is increasingly essential for datapoor species (Miller et al., 2019; Isaac et al., 2020). Our results indicate that COP variables, largely driven by commonly recorded species, can indirectly inform distributions of rarer taxa. These findings highlight the conservation value of citizen science and suggest that its rapidly growing data volumes can be productively leveraged.

Fill-in approach with annual predictions

We propose a fill-in approach that generates predictions of annual occupancy to bridge monitoring gaps. By tracking year-to-year occupancy across North America since 2007, we evaluated 10-year reductions under the IUCN Red List Criterion A (Figure 3). Traditional time-series workflows often filter datasets to well-monitored regions, narrowing the spatiotemporal scope of inference and precluding absolute-extent assessments (e.g., IUCN Red List). Current trend models require extensive filtering or structured surveys—resources structurally inaccessible to the taxa most in need of conservation insight. Our approach produces fine-temporal estimations from sparse, presence-only, multisourced datasets, directly benefitting them.

For under-monitored species, annual COP-ML predictions could complement application of Bayesian occupancy models (OM) that track temporal change but typically require at least two revisits per period or high spatiotemporal density in data (Royle, 2006; Kamp et al., 2016; Outhwaite et al., 2018; Perkins-Taylor & Frey, 2020; Jha et al., 2022). In North America, limited data density has often forced coarse spatial (~10,000 km²) and temporal (10–20 years) resolutions for insects such as bees and dragonflies (Soroye et al., 2020; Jackson et al., 2022; Shirey et al., 2023). One way to recover resolution is to commission additional, targeted surveys (Xue et al., 2016; Tulloch et al., 2013), but this is costly; our approach offers an alternative by producing annual predictions without new field effort. However, incorporating ML-based predictions into occupancy modeling frameworks—as pseudoobservations—remains largely untested. Occupancy models explicitly model detection and survey processes, so predicted probabilities must be reconciled with those components. Although recent OM advances have explored non-ideal data (e.g., assuming random-walk observation processes, using pseudo-absence instead of checklist absence, or treating opportunistic records as revisits; Outhwaite et al., 2018), the statistical compatibility of ML predictions within OM frameworks has yet to be demonstrated.

The fill-in and filtering strategies are complementary. However, we deliberately did not apply spatial thinning or filtering to address spatial autocorrelation. First, our targets are data-deficient, and our goal is to cover the entire known range, so filtering is impractical. Second, our aim was to test whether COP-ML, devised for such taxa, remains reliable without deep filtering. Third, our case is temporal interpolation at the same sites, rather than prediction to new locations, spatial leakage is less relevant in our setting. Finally, state-matched pseudo-absences further mitigate overfitting by emphasizing within-region discrimination. Future

study will be needed to determine filtering levels that optimally balance performance, bias reduction, and data retention in the fill-in framework.

This study presents a scalable method to bridge monitoring gaps for data-deficient species, using sparse, largely opportunistic, presence-only records to generate annual occupancy estimates. Predicted trends aligned with long-term trends from independent regional monitoring and operationalized IUCN Red List criteria for species previously excluded due to lack of data. COP-ML demonstrated robust performance across heterogeneous sources and time periods, showing that reliable signals of extinction risk can emerge even from unstructured datasets. By converting fragmented observations into interpretable trends, the fill-in approach with annual predictions provides a practical pathway to extend assessment of extinction risk and strengthen conservation decisions where standardized monitoring is absent.

References

487

488

489

490

491

492

493

494

495

- 497 Altwegg, R., & Nichols, J. D. (2019). Occupancy models for citizen-science data. *Methods in*
- 498 *Ecology and Evolution*, **10(1)**, 8-21.
- 499 Alyokhin, A., & Sewell, G. (2004). Changes in a lady beetle community following the
- establishment of three alien species. *Biological Invasions*, **6(4)**, 463–471.
- Bahlai, C. A., Colunga-Garcia, M., Gage, S. H., & Landis, D. A. (2015). The role of exotic
- ladybeetles in the decline of native ladybeetle populations: evidence from long-term
- monitoring. *Biological Invasions*, **17**, 1005-1024.
- Bayraktarov, E., Ehmke, G., O'Connor, J., Burns, E. L., Nguyen, H. A., McRae, L., Possingham,
- 505 H. P., & Lindenmayer, D. B. (2019). Do big unstructured biodiversity data mean more
- knowledge? Frontiers in Ecology and Evolution, **6**, 239.
- 507 Biodiversity Information Serving Our Nation (BISON). (2021, September 16). Coccinellidae
- 508 records, 2007–2021, United States and Canada [Data set]. Retrieved from
- 509 http://bison.usgs.ornl.gov/
- Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly*
- 511 *weather review*, **78(1)**, 1-3.
- 512 BugGuide.Net. (2021, September 26). Coccinellidae search results [Web-scraped data].
- Retrieved from https://bugguide.net/node/view/15740

- 514 Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In B.
- Krishnapuram, M. Shah, A. J. Smola, C. Aggarwal, D. Shen, & R. Rastogi (Eds.), *Proceedings*
- of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data
- 517 *Mining* (pp. 785–794). ACM.
- 518 Cheney, B., Thompson, P. M., Ingram, S. N., Hammond, P. S., Stevick, P. T., Durban, J. W.,
- 519 Culloch, R. M., Elwen, S. H., Mandleberg, L., Janik, V. M., Quick, N. J., Islas-Villanueva, V.,
- Robinson, K. P., Costa, M., Eisfeld, S. M., Walters, A., Philips, C., Weir, C. R., Evans, P. G.
- 521 H., ... Wilson, B. (2013). Integrating multiple data sources to assess the distribution and
- abundance of bottlenose dolphins Tursiops truncatus in Scottish waters. *Mammal Review*, **43(1)**,
- 523 71–88.
- 524 Coderre, D., Lucas, É., & Gagné, I. (1995). The occurrence of Harmonia axyridis
- 525 (Pallas)(Coleoptera: Coccinellidae) in Canada. *The Canadian Entomologist*, **127(4)**, 609-611.
- 526 Cohen, J. (1960). A coefficient of agreement for nominal scales. Educational and psychological
- 527 *measurement*, **20(1)**, 37-46.
- 528 Colunga-Garcia, M., & Gage, S. H. (1998). Arrival, establishment, and habitat use of the
- 529 multicolored Asian lady beetle (Coleoptera: Coccinellidae) in a Michigan landscape.
- 530 Environmental Entomology, **27(6)**, 1574-1580.
- 531 COSEWIC. (2016a). COSEWIC assessment and status report on the Transverse Lady Beetle
- 532 (Coccinella transversoguttata) in Canada. Committee on the Status of Endangered Wildlife in
- 533 Canada.
- 534 COSEWIC. (2016b). COSEWIC assessment and status report on the Nine-spotted Lady Beetle
- 535 (Coccinella novemnotata) in Canada. Committee on the Status of Endangered Wildlife in
- 536 Canada.
- Elliott, N., Kieckhefer, R., & Kauffman, W. (1996). Effects of an invading coccinellid on native
- coccinellids in an agricultural landscape. *Oecologia*, **105(4)**, 537–544.
- Elton, C. S. (2000). The ecology of invasions by animals and plants. *Methuen*. (Original work
- 540 published 1958)
- Estes, L., Elsen, P. R., Treuer, T., Ahmed, L., Caylor, K., Chang, J., Choi, J. J., & Ellis, E. C.
- 542 (2018). The spatial and temporal domains of modern ecology. *Nature Ecology & Evolution*,
- **2(5)**, 819–826.

- 544 Evans, E. W. (2000). Morphology of invasion: body size patterns associated with establishment
- of Coccinella septempunctata (Coleoptera: Coccinellidae) in western North America.
- *European Journal of Entomology*, **97(4)**, 469-474.
- 547 Evans, E. W. (2004). Habitat displacement of North American ladybirds by an introduced
- 548 species. *Ecology*, **85(3)**, 637-647.
- Evans, E. W., Soares, A. O., & Yasuda, H. (2011). Invasions by ladybugs, ladybirds, and other
- predatory beetles. *BioControl*, **56**, 597-611.
- Fielding, A. H., & Bell, J. F. (1997). A review of methods for the assessment of prediction
- errors in conservation presence/absence models. *Environmental conservation*, **24(1)**, 38-49.
- 553 Fink, D., Auer, T., Johnston, A., Ruiz-Gutierrez, V., Hochachka, W. M., & Kelling, S. (2020).
- Modeling avian full annual cycle distribution and population trends with citizen science data.
- 555 Ecological Applications, **30(3)**, e02056.
- Fletcher Jr, R. J., Hefley, T. J., Robertson, E. P., Zuckerberg, B., McCleery, R. A., & Dorazio,
- R. M. (2019). A practical guide for combining data to model species distributions. *Ecology*,
- 558 **100(6)**, e02710.
- Franklin, J. (2013). Species distribution models in conservation biogeography: developments
- and challenges. *Diversity and distributions*, **19(10)**, 1217-1223.
- 561 Gardiner, M. M., Allee, L. L., Brown, P. M., Losey, J. E., Roy, H. E., & Smyth, R. R. (2012).
- Lessons from lady beetles: accuracy of monitoring data from US and UK citizen-science
- programs. Frontiers in Ecology and the Environment, **10(9)**, 471-476.
- 564 GBIF. (2022, October 31). Coccinellidae occurrence dataset, 2007–2021, United States and
- Canada [Data set]. GBIF Occurrence Download. https://doi.org/10.15468/dl.vgen57
- Geldmann, J., Heilmann-Clausen, J., Holm, T. E., Levinsky, I., Markussen, B. O., Olsen, K.,
- Rahbek, C., & Tøttrup, A. P. (2016). What determines spatial bias in citizen science? Exploring
- four recording schemes with different proficiency requirements. *Diversity and Distributions*,
- **22(11)**, 1139–1149.
- 570 Gordon, R. D. (1985). The Coccinellidae (Coleoptera) of America North of Mexico. *Journal*
- of the New York Entomological Society, **93(1)**, 1-912.

- 572 Gordon, R. D., & Vandenberg, N. J. (1991). Field guide to recently introduced species of
- 573 Coccinellidae (Coleoptera) in North America, with a revised key to North American genera of
- 574 Coccinellini. Proceedings of the Entomological Society of Washington, 93(4), 845-867.
- Guzman, L. M., Johnson, S. A., Mooers, A. O., & M'Gonigle, L. K. (2021). Using historical
- data to estimate bumble bee occurrence: Variable trends across species provide little support
- for community-level declines. *Biological Conservation*, **257**, 109141.
- Harmon, J. P., Stephens, E., & Losey, J. (2007). The decline of native coccinellids (Coleoptera:
- 579 Coccinellidae) in the United States and Canada. *Journal of Insect Conservation*, **11(2)**, 85–94.
- Harvey, J. A., Heinen, R., Armbrecht, I., Basset, Y., Baxter-Gilbert, J. H., Bezemer, T. M.,
- Böhm, M., Christie, A. P., Cornelisse, T., Crone, E. E., Dicke, M., Dicks, L. V., Elder, M.,
- Fartmann, T., Forister, M. L., Gaston, K. J., Jepsen, S. J., Jones, T. H., Kaydan, M. B., ... de
- 583 Kroon, H. (2020). International scientists formulate a roadmap for insect conservation and
- recovery. *Nature Ecology & Evolution*, **4(2)**, 174–176.
- Hentley, W. T., Vanbergen, A. J., Beckerman, A. P., Brien, M. N., Hails, R. S., Jones, T. H., &
- Johnson, S. N. (2016). Antagonistic interactions between an invasive alien and a native
- coccinellid species may promote coexistence. *Journal of Animal Ecology*, **85(4)**, 1087-1097.
- Hernandez, P. A., Graham, C. H., Master, L. L., & Albert, D. L. (2006). The effect of sample
- size and species characteristics on performance of different species distribution modeling
- 590 methods. *Ecography*, **29(5)**, 773-785.
- Hesler, L. S., & Kieckhefer, R. W. (2008). Status of exotic and previously common native
- 592 coccinellids (Coleoptera) in South Dakota landscapes. Journal of the Kansas Entomological
- 593 *Society*, **81(1)**, 29-49.
- Hesler, L. S., Catangui, M. A., Losey, J. E., Helbig, J. B., & Mesman, A. (2009). Recent records
- of Adalia bipunctata (L.), Coccinella transversoguttata richardsoni Brown, and Coccinella
- 596 novemnotata Herbst (Coleoptera: Coccinellidae) from South Dakota and Nebraska. The
- 597 *Coleopterists Bulletin*, **63(4)**, 475-484.
- Hesler, L. S., Kieckhefer, R. W., & Catangui, M. A. (2004). Surveys and field observations of
- 599 Harmonia axyridis and other Coccinellidae (Coleoptera) in eastern and central South Dakota.
- 600 Transactions of the American Entomological Society, **130(1)**, 113–133.

- 601 Honěk, A. (1985). Habitat preferences of aphidophagous coccinellids [Coleoptera].
- 602 Entomophaga, **30(3)**, 253-264.
- Horns, J. J., Adler, F. R., & Şekercioğlu, Ç. H. (2018). Using opportunistic citizen science data
- to estimate avian population trends. *Biological conservation*, **221**, 151-159.
- Hosmer, D. W., Jr., Lemeshow, S., & Sturdivant, R. X. (2013). Applied logistic regression (3rd
- 606 ed.). Wiley.
- iDigBio. (2021, September 16). Coccinellidae records, 2007–2021, United States and Canada
- 608 [Data set]. Retrieved from https://www.idigbio.org/
- Inamine, H., Ellner, S. P., Springer, J. P., & Agrawal, A. A. (2016). Linking the continental
- 610 migratory cycle of the monarch butterfly to understand its population decline. Oikos, 125(8),
- 611 1081-1091.
- 612 iNaturalist. (2021, September 28). Coccinellidae observations, 2007–2021, United States and
- 613 Canada [Data set]. Retrieved from https://www.inaturalist.org/
- Isaac, N. J. B., Jarzyna, M. A., Keil, P., Dambly, L. I., Boersch-Supan, P. H., Browning, E.,
- 615 Freeman, S. N., Golding, N., Guillera-Arroita, G., Henrys, P. A., Jarvis, S., Lahoz-Monfort, J.,
- Pagel, J., Pescott, O. L., Schmucki, R., Simmonds, E. G., & O'Hara, R. B. (2020). Data
- 617 integration for large-scale models of species distributions. Trends in Ecology & Evolution,
- 618 **35(1)**, 56–67.
- 619 Isaac, N. J., & Pocock, M. J. (2015). Bias and information in biological records. Biological
- 620 *Journal of the Linnean Society*, **115(3)**, 522-531.
- Isaac, N. J., van Strien, A. J., August, T. A., de Zeeuw, M. P., & Roy, D. B. (2014). Statistics
- 622 for citizen science: extracting signals of change from noisy ecological data. *Methods in Ecology*
- 623 and Evolution, **5(10)**, 1052-1060.
- 624 IUCN. (2024). The IUCN Red List of Threatened Species. Version 16. Available at:
- 625 https://www.iucnredlist.org/resources/redlistguidelines.
- Johnston, A., Fink, D., Hochachka, W. M., & Kelling, S. (2018). Estimates of observer
- 627 expertise improve species distributions from citizen science data. Methods in Ecology and
- 628 Evolution, **9(1)**, 88-97.

- Jackson, H. M., Johnson, S. A., Morandin, L. A., Richardson, L. L., Guzman, L. M., &
- 630 M'Gonigle, L. K. (2022). Climate change winners and losers among North American
- 631 bumblebees. *Biology letters*, **18(6)**, 20210551.
- 632 Jeffries, D. L., Chapman, J., Roy, H. E., Humphries, S., Harrington, R., Brown, P. M., &
- Handley, L. J. L. (2013). Characteristics and drivers of high-altitude ladybird flight: insights
- from vertical-looking entomological radar. *PloS one*, **8(12)**, e82278.
- Jha, A., Praveen, J., & Nameer, P. O. (2022). Contrasting occupancy models with presence-
- only models: does accounting for detection lead to better predictions?. Ecological Modelling,
- **472**, 110105.
- Jönsson, G. M., Broad, G. R., Sumner, S., & Isaac, N. J. (2021). A century of social wasp
- occupancy trends from natural history collections: spatiotemporal resolutions have little effect
- on model performance. *Insect Conservation and Diversity*, **14(5)**, 543-555.
- Justice, A. C., Covinsky, K. E., & Berlin, J. A. (1999). Assessing the generalizability of
- prognostic information. Annals of internal medicine, 130(6), 515-524.
- Kajita, Y., Takano, F., Yasuda, H., & Agarwala, B. K. (2000). Effects of indigenous ladybird
- species (Coleoptera: Coccinellidae) on the survival of an exotic species in relation to prey
- abundance. *Applied Entomology and Zoology*, **35(4)**, 473-479.
- Kajita, Y., Yasuda, H., & Evans, E. W. (2006). Effects of native ladybirds on oviposition of the
- exotic species, Adalia bipunctata (Coleoptera: Coccinellidae), in Japan. Applied Entomology
- 648 and Zoology, **41(1)**, 57-61.
- Kamp, J., Oppel, S., Heldbjerg, H., Nyegaard, T., & Donald, P. F. (2016). Unstructured citizen
- science data fail to detect long-term population declines of common birds in Denmark.
- 651 *Diversity and Distributions*, **22(10)**, 1024-1035.
- Kéry, M., Gardner, B., & Monnerat, C. (2010). Predicting species distributions from checklist
- data using site-occupancy models. *Journal of Biogeography*, **37(10)**, 1851-1862.
- Kissling, W. D., Ahumada, J. A., Bowser, A., Fernandez, M., Fernández, N., García, E. A.,
- 655 Guralnick, R. P., Isaac, N. J. B., Kelling, S., Los, W., McRae, L., Mihoub, J. B., Obst, M.,
- 656 Santamaria, M., Skidmore, A. K., Williams, K. J., Agosti, D., Amariles, D., Arvanitidis, C.,
- Bastin, L., ... Hardisty, A. R. (2018). Building essential biodiversity variables (EBVs) of species
- distribution and abundance at a global scale. *Biological Reviews*, **93**, 600–625.

- Kissling, W. D., Dormann, C. F., Groeneveld, J., Hickler, T., Kühn, I., McInerny, G. J., Montoya,
- J. M., Römermann, C., Schiffers, K., Schurr, F. M., Singer, A., Svenning, J.-C., Zimmermann,
- N. E., & O'Hara, R. B. (2012). Towards novel approaches to modelling biotic interactions in
- multispecies assemblages at large spatial extents. *Journal of Biogeography*, **39(12)**, 2163–2178.
- Knape, J., Coulson, S. J., van der Wal, R., & Arlt, D. (2022). Temporal trends in opportunistic
- citizen science reports across multiple taxa. *Ambio*, **51(1)**, 183–198.
- Koch, R. L. (2003). The multicolored Asian lady beetle, Harmonia axyridis: a review of its
- biology, uses in biological control, and non-target impacts. *Journal of insect Science*, **3(1)**, 32.
- Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical
- data. *Biometrics*, **33(1)**, 159–174.
- 669 Larsen, E. A., & Shirey, V. (2021). Method matters: Pitfalls in analysing phenology from
- 670 occurrence records. *Ecology Letters*, **24(6)**, 1287-1289.
- 671 LeCroy, K. A., Savoy-Burke, G., Carr, D. E., Delaney, D. A., & Roulston, T. A. H. (2020).
- Decline of six native mason bee species following the arrival of an exotic congener. Scientific
- 673 reports, **10(1)**, 18745.
- Losey, J. E., Perlman, J. E., & Hoebeke, E. R. (2007). Citizen scientist rediscovers rare nine-
- 675 spotted lady beetle, Coccinella novemnotata, in eastern North America. Journal of Insect
- 676 *Conservation*, **11(4)**, 415–417.
- Losey, J., Allee, L., & Smyth, R. (2012). The Lost Ladybug Project: Citizen spotting surpasses
- 678 scientist's surveys. *American Entomologist*, **58(1)**, 22-24.
- Luan, J., Zhang, C., Xu, B., Xue, Y., & Ren, Y. (2020). The predictive performances of random
- 680 forest models with limited sample size and different species traits. Fisheries Research, 227,
- 681 105534.
- MacKenzie, D. I., Nichols, J. D., Lachman, G. B., Droege, S., Andrew Royle, J., & Langtimm,
- 683 C. A. (2002). Estimating site occupancy rates when detection probabilities are less than one.
- 684 *Ecology*, **83(8)**, 2248-2255.
- MacKenzie, D. I., Nichols, J. D., Royle, J. A., Pollock, K. H., Bailey, L., & Hines, J. E. (2017).
- Occupancy estimation and modeling: Inferring patterns and dynamics of species occurrence
- 687 (2nd ed.). Academic Press.

- Martín, B., González-Arias, J., & Vicente-Vírseda, J. A. (2021). Machine learning as a
- 689 successful approach for predicting complex spatio-temporal patterns in animal species
- abundance. Machine learning, 44, 289-301.
- Martínez-Minaya, J., Cameletti, M., Conesa, D., & Pennino, M. G. (2018). Species distribution
- modeling: a statistical review with focus in spatio-temporal issues. Stochastic environmental
- 693 research and risk assessment, **32**, 3227-3244.
- McCorquodale, D. B. (1998). Adventive lady beetles (Coleoptera: Coccinellidae) in eastern
- Nova Scotia, Canada. *Entomological News*, **109**, 15–20.
- 696 Mi, C., Huettmann, F., Guo, Y., Han, X., & Wen, L. (2017). Why choose Random Forest to
- 697 predict rare species distribution with few samples in large undersampled areas? Three Asian
- 698 crane species models provide supporting evidence. *PeerJ*, **5**, e2849.
- 699 Miller, D. A., Pacifici, K., Sanderlin, J. S., & Reich, B. J. (2019). The recent past and promising
- future for data integration methods to estimate species' distributions. Methods in Ecology and
- 701 Evolution, **10(1)**, 22-37.
- Montgomery, G. A., Dunn, R. R., Fox, R., Jongejans, E., Leather, S. R., Saunders, M. E.,
- Shortall, C. R., Tingley, M. W., & Wagner, D. L. (2020). Is the insect apocalypse upon us? How
- to find out. *Biological Conservation*, **241**, 108327.
- Mukwevho, V. O., Pryke, J. S., & Roets, F. (2017). Habitat preferences of the invasive
- 706 harlequin ladybeetle *Harmonia axyridis* (Coleoptera: Coccinellidae) in the Western Cape
- 707 Province, South Africa. *African Entomology*, **25(1)**, 86-97.
- Nedvěd O. 1999. Host complexes of predaceous ladybeetles (Coleoptera: Coccinellidae).
- 709 *Journal of Applied Entomology*, **123**, 73–76.
- Newson, S. E., Evans, H. E., & Gillings, S. (2015). A novel citizen science approach for large-
- scale standardised monitoring of bat activity and distribution, evaluated in eastern England.
- 712 *Biological Conservation*, **191**, 38-49.
- 713 North Carolina State University Insect collection (NCSU). (2021, September 26).
- 714 Coccinellidae records [Web-scraped data]. Retrieved from
- 715 http://specimens.insectmuseum.org/public/specimen
- 716 Olden, J. D., Lawler, J. J., & Poff, N. L. (2008). Machine learning methods without tears: a
- primer for ecologists. The Quarterly review of biology, 83(2), 171-193.

- Omkar, & Pervez, A. (2005). Ecology of two-spotted ladybird, *Adalia bipunctata*: a review.
- 719 *Journal of Applied Entomology*, **129(9-10)**, 465-474.
- Outhwaite, C. L., Chandler, R. E., Powney, G. D., Collen, B., Gregory, R. D., & Isaac, N. J.
- 721 (2018). Prior specification in Bayesian occupancy modelling improves analysis of species
- occurrence data. *Ecological Indicators*, **93**, 333-343.
- Pagel, J., Anderson, B. J., O'Hara, R. B., Cramer, W., Fox, R., Jeltsch, F., Roy, D. B., Thomas,
- 724 C. D., & Schurr, F. M. (2014). Quantifying range-wide variation in population trends from local
- abundance surveys and widespread opportunistic occurrence records. *Methods in Ecology and*
- 726 Evolution, **5(8)**, 751–760.
- Perkins-Taylor, I. E., & Frey, J. K. (2020). Predicting the distribution of a rare chipmunk
- 728 (Neotamias quadrivittatus oscuraensis): comparing MaxEnt and occupancy models. *Journal of*
- 729 *Mammalogy*, **101(4)**, 1035-1048.
- Petersen, M. J., & Losey, J. E. (2024). Niche overlap with an exotic competitor mediates the
- abundant niche-centre relationship for a native lady beetle. *Diversity and Distributions*, **30(5)**,
- 732 e13825.
- Pollock, L. J., Tingley, R., Morris, W. K., Golding, N., O'Hara, R. B., Parris, K. M., Vesk, P. A.,
- 8 McCarthy, M. A. (2014). Understanding co-occurrence by modelling species simultaneously
- with a Joint Species Distribution Model (JSDM). Methods in Ecology and Evolution, 5(5),
- 736 397–406.
- 737 Rencher, A. C. (1995). Methods of multivariate analysis. Wiley.
- Robinson, O. J., Ruiz-Gutierrez, V., & Fink, D. (2018). Correcting for bias in distribution
- modelling for rare species using citizen science data. Diversity and Distributions, 24(4), 460-
- 740 472.
- Royle, J. A. (2006). Site occupancy models with heterogeneous detection probabilities.
- 742 *Biometrics*, **62(1)**, 97-102.
- Schultz, C. B., Brown, L. M., Pelton, E., & Crone, E. E. (2017). Citizen science monitoring
- demonstrates dramatic declines of monarch butterflies in western North America. *Biological*
- 745 *Conservation*, **214**, 343-346.

- Shirey, V., Khelifa, R., M'Gonigle, L. K., & Guzman, L. M. (2023). Occupancy-detection
- models with museum specimen data: Promise and pitfalls. *Methods in Ecology and Evolution*,
- 748 **14(2)**, 402-414.
- Soares, A. O., & Serpa, A. (2007). Interference competition between ladybird beetle adults
- 750 (Coleoptera: Coccinellidae): Effects on growth and reproductive capacity. *Population Ecology*,
- **49(1)**, 37–43.
- Soroye, P., Newbold, T., & Kerr, J. (2020). Climate change contributes to widespread declines
- among bumble bees across continents. *Science*, **367(6478)**, 685-688.
- Steen, V. A., Elphick, C. S., & Tingley, M. W. (2019). An evaluation of stringent filtering to
- 755 improve species distribution models from citizen science data. Diversity and Distributions,
- **25(12)**, 1857-1869.
- 757 Strayer, D. L., Eviner, V. T., Jeschke, J. M., & Pace, M. L. (2006). Understanding the long-term
- 758 effects of species invasions. *Trends in ecology & evolution*, **21(11)**, 645-651.
- 759 Svancara, L. K., Abatzoglou, J. T., & Waterbury, B. (2019). Modeling current and future
- potential distributions of milkweeds and the monarch butterfly in Idaho. Frontiers in Ecology
- 761 and Evolution, 7, 168.
- Szabo, J. K., Vesk, P. A., Baxter, P. W., & Possingham, H. P. (2010). Regional avian species
- 763 declines estimated from volunteer-collected long-term data using List Length Analysis.
- 764 *Ecological Applications*, **20(8)**, 2157-2169.
- 765 The Lost Ladybug Project. (2021, September 28). Lost Ladybug Project dataset [Data set].
- Retrieved from http://www.lostladybug.org/
- 767 Tikhonov, G., Abrego, N., Dunson, D., & Ovaskainen, O. (2017). Using joint species
- 768 distribution models for evaluating how species-to-species associations depend on the
- 769 environmental context. *Methods in Ecology and Evolution*, **8(4)**, 443-452.
- 770 Tingley, M. W., & Beissinger, S. R. (2009). Detecting range shifts from historical species
- occurrences: new perspectives on old data. Trends in ecology & evolution, 24(11), 625-633.
- Tulloch, A. I., Possingham, H. P., Joseph, L. N., Szabo, J., & Martin, T. G. (2013). Realising
- the full potential of citizen science monitoring programs. *Biological Conservation*, **165**, 128-
- 774 138.

- 775 Turnock, W. J., Wise, I. L., & Matheson, F. O. (2003). Abundance of some native coccinellines
- 776 (Coleoptera: Coccinellidae) before and after the appearance of *Coccinella septempunctata*. The
- 777 *Canadian Entomologist*, **135(3)**, 391-404.
- Van Eupen, C., Maes, D., Herremans, M., Swinnen, K. R., Somers, B., & Luca, S. (2021). The
- impact of data quality filtering of opportunistic citizen science data on species distribution
- 780 model performance. Ecological Modelling, 444, 109453.
- Van Strien, A. J., Van Swaay, C. A., & Termaat, T. (2013). Opportunistic citizen science data
- of animal species produce reliable estimates of distribution trends if analysed with occupancy
- 783 models. *Journal of Applied Ecology*, **50(6)**, 1450-1458.
- Vaughan, I. P., & Ormerod, S. J. (2005). The continuing challenges of testing species
- distribution models. *Journal of applied ecology*, **42(4)**, 720-730.
- Walker, J., & Taylor, P. D. (2017). Using eBird data to model population change of migratory
- 787 bird species. Avian Conservation and Ecology, 12(1), 4.
- Wheeler, A. G., Jr., & Hoebeke, E. R. (1995). Coccinella novemnotata in northeastern North
- America: Historical occurrence and current status (Coleoptera: Coccinellidae). *Proceedings of*
- 790 the Entomological Society of Washington, **97**, 701–716.
- Willig, M. R., Woolbright, L., Presley, S. J., Schowalter, T. D., Waide, R. B., Heartsill Scalley,
- 792 T., Zimmerman, J. K., González, G., & Lugo, A. E. (2019). Populations are not declining and
- food webs are not collapsing at the Luquillo Experimental Forest. *Proceedings of the National*
- 794 Academy of Sciences, **116(25)**, 12143–12144.
- 795 Wisz, M. S., Hijmans, R. J., Li, J., Peterson, A. T., Graham, C. H., Guisan, A., & NCEAS
- 796 Predicting Species Distributions Working Group. (2008). Effects of sample size on the
- 797 performance of species distribution models. *Diversity and distributions*, **14(5)**, 763-773.
- Woltz, J. M., & Landis, D. A. (2013). Coccinellid immigration to infested host patches
- 799 influences suppression of Aphis glycines in soybean. *Biological Control*, **64(3)**, 330-337.
- Xue, Y., Davies, I., Fink, D., Wood, C., & Gomes, C. P. 2016. Avicaching: A two stage game
- 801 for bias reduction in citizen science. Proceedings of the 15th International Conference on
- *Autonomous Agents and Multiagent Systems*, 776–785.
- Zimmermann, N. E., Edwards Jr, T. C., Graham, C. H., Pearman, P. B., & Svenning, J. C.
- 804 (2010). New trends in species distribution modelling. *Ecography*, **33(6)**, 985-989.

Zizka, A., Antonelli, A., & Silvestro, D. (2021). Sampbias, a method for quantifying geographic
sampling biases in species distribution data. *Ecography*, 44(1), 25-32.

808 Figures

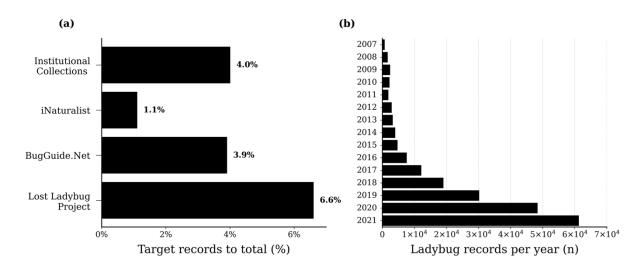


Figure 1. Structural and temporal inconsistencies in our multisourced dataset; (a) difference in detection efficiency (target records per 100 ladybug records) in institutional collections and three citizen science sources (The Lost Ladybug Project, iNaturalist, BugGuide.Net); (b) exponential increase in annual ladybug observations across the USA and Canada (2007–2021; $y = 758.89e^{0.29x}$, $R^2 = 0.88$).

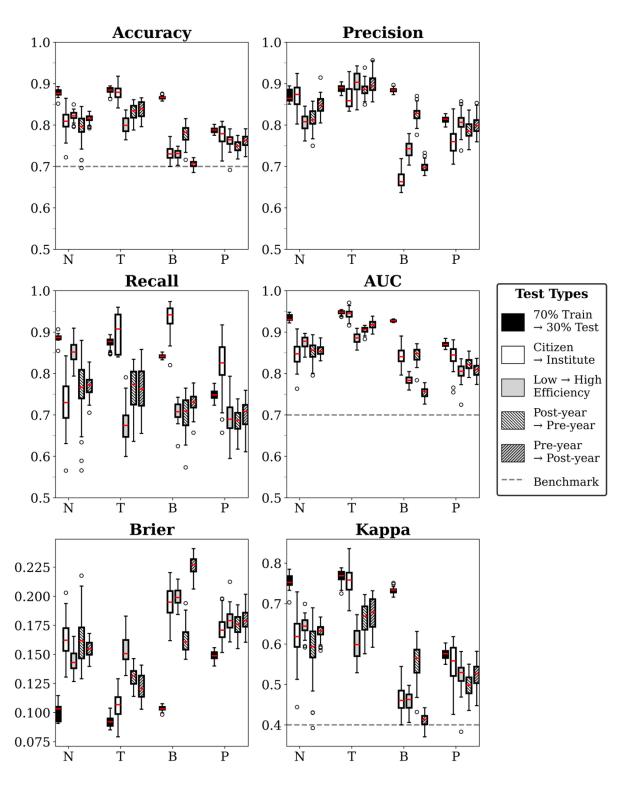


Figure 2. Performance of species distribution models using co-occurrence pattern variables for C. novemnotata (N), C. transversoguttata (T), A. bipunctata (B), and H. parenthesis (P) across standard evaluation, structural generalizations, and temporal generalizations; black plots, standard 70% training and 30% testing; white plots, citizen science training and institutional data testing; gray plots, low target-density group (1.3%) training and high density source (6.6%) testing; left hatch, post-2007 training until \sim 70% coverage and later-year testing; right hatch,

Pre-2021 traing until ~70% coverage and earlier-year testing; red lines, mean performance across 2,500 runs; dash, reliability benchmarks—Accuracy \geq 0.70, AUC \geq 0.70 (Hosmer et al., 2013), Kappa \geq 0.40 (Landis & Koch, 1977), and Brier \leq 0.25 (Brier, 1950).

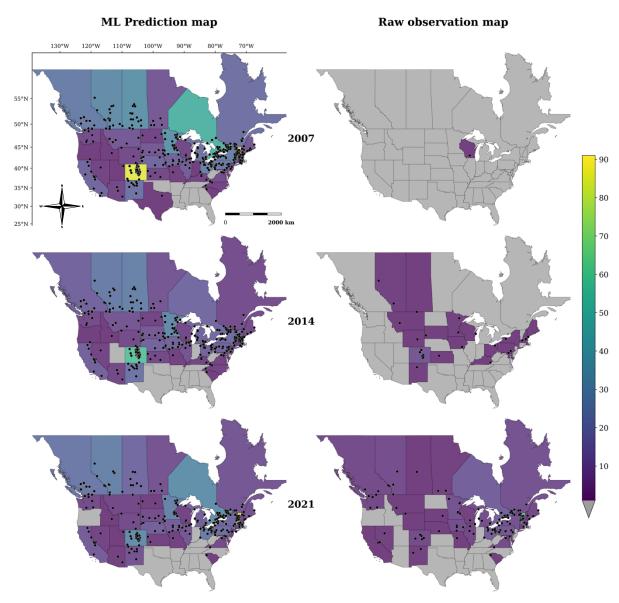


Figure 3. Annual occupancy of *Hippodamia parenthesis* (2007–2021) in the USA and Canada; Left, co-occurrence pattern model predictions at sites with prior records of the species; Right, raw observations from the compiled dataset; Dots, occupancy at previously recorded coordinates in the compiled dataset; color gradients, density of occupancy within each state or province.

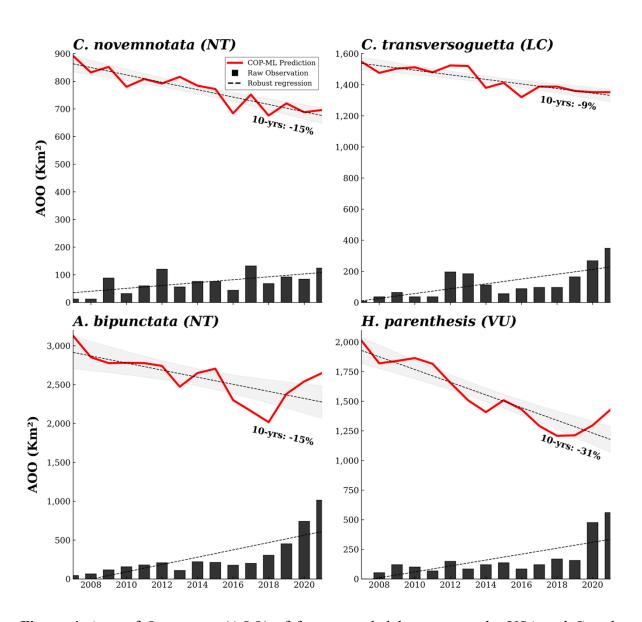


Figure 4. Area of Occupancy (AOO) of four target ladybugs across the USA and Canada (2007–2021); red lines, annual model predictions; dashed lines, robust regression fits showing declines with 95% confidence intervals; IUCN Red List categories, derived from 10-year reduction rates; bars, raw observations showing increases.

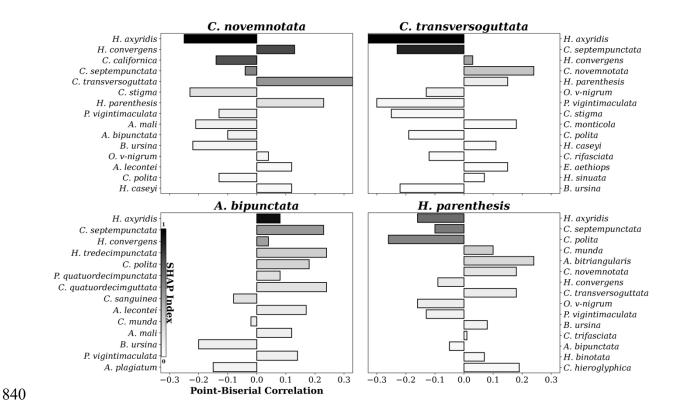


Figure 5. Predictor importance and correlation of co-occurrence pattern variables with target species' occupancy (presence = 0, absence = 1); y-axis, SHapley Additive exPlanations (SHAP) indices quantifying variable contributions in each model; x-axis, point-biserial correlation between target occupancy and each variable species' occurrence within 18 km radius.

Appendix

Table S1. Ladybug records used were compiled from seven digital platforms (three citizen science, one museum collection, and three metadata sources); Regional abbreviations: AK = Alaska, HI = Hawaii, MB = Manitoba, ON = Ontario, SK = Saskatchewan, BC = British Columbia, AB = Alberta, QC = Quebec.

Source	Size	Type	Data Download & Refinement Criteria
Lost Ladybug Project	32,905	Citizen science	• Years 2007-2021
iNaturalist	197,990	Citizen science	• U.S. (excluding AK & HI) and Canadian provinces (MB, ON, SK, BC, AB, QC)
bugGuide.Net	27,018	Citizen science	• Positional accuracy < 1 km (if applicable)
GBIF	143,000	Metadata	• Species level
GDII	143,000	source	 Adult records or images
BISON	109,834	Metadata source	• Drop duplicates at year-GPS-species
IdigBio	99,723	Metadata source	
NCSU Insect Museum	5,425	Institute	Final Dataset: 188,644

Table S2. ANOSIM statistics assess differences in co-occurrence pattern variables across data groups in the generalization test sets.

Species	Citizen science vs Institutions		Post-year (train) vs Pre-year (test)			Pre-year (train) vs Post-year (test)				
Species	R-value	<i>p</i> -value	Test period	Train (% of presence)	R-value	<i>p</i> -value	Test period	Train (% of presence)	<i>R</i> -value	<i>p</i> -value
C. transversoguttata	0.22	0.001	after 2019	65%	0.03	0.001	before 2014	70%	0.07	0.001
C. novemnotata	0.06	0.001	after 2018	70%	0.06	0.001	before 2013	72%	0.02	0.001
H. parenthesis	0.12	0.001	after 2020	74%	0.07	0.001	before 2014	74%	0.05	0.001
A. bipunctata	0.05	0.001	after 2019+	55%	0.02	0.001	before 2017	67%	0.02	0.001
Absence datapoints	0.06	0.001			0.03	0.001			0.01	0.001

⁺To assess COP model generalization, the training period for A. bipunctata—the species with most presence records—was reduced to extend the testing period.

Table S3. Results of regression estimates, diagnostic tests, and 2012-2021 reduction rates (OLS: ordinary least squares regression, Huber: robust regression).

	C. novemnotata	C. transversoguetta	A. bipunctata	H. parenthesis
B (OLS)	-13.36	-14.73	-45.59	-53.63
R ² (OLS)	0.818	0.733	0.500	0.834
p (OLS)	0.0000****	0.0000****	0.0032***	0.0000****
95% CI (OLS)	-17.13, -9.59	-20.06, -9.4	-72.91, -18.26	-67.97, -39.29
Reduction (10-yr, OLS)	-15%	-9%	-15%	-29%
Breusch- Pagan p	0.9001	0.7992	0.0226*	0.0804
White p	0.7238	0.1968	0.0402*	0.0736
p (Robust SE)	0.0000****	0.0000****	0.0078*	0.0000****
95% CI (Robust SE)	-16.77, -9.95	-18.14, -11.31	-79.18, -11.99	-72.55, -34.71
Max Cook's Distance	0.2056	0.2028	0.6505	1.0651
B (Huber)	-13.00	-14.38	-45.98	-56.24
R ² (Huber)	0.811	0.731	0.500	0.831
95% CI (Huber)	-16.42, -9.58	-17.39, -11.38	-72.32, -19.65	-68.23, -44.26
Reduction (10-yr, Huber)	-15%	-9%	-15%	-31%

(p* < 0.05, p** < 0.05, p*** < 0.005, p**** < 0.005)

Table S4. Predicted distribution trends (2007–2021) and corresponding IUCN Red List categories for four target species, based on reductions in area of occupancy (AOO) and extent of occurrence (EOO).

C	Reduction	IUCN category	AOO (km²)		EOO (km²)		
Species	in 10-yrs		2007	2021	2007	2021	
H. parenthesis	31%	VU	1,548	1,352	8,450,469	7,749,070	
A. bipunctata	15%	NT	3,128	2,648	11,538,691	10,817,443	
C. novemnotata	15%	NT	2,012	1,428	5,480,067	5,399,901	
C. transversoguttata	9%	LC	892	696	9,820,525	9,146,848	

		C. novemnotata	C. transversoguttata	H. parenthesis	A. bipunctata	
F-statistic (DF Model, DF Residual)		12.96 ** (4, 10)	7.80 ** (4, 10)	34.72*** (4, 10)	12.86 ** (4, 10)	
F	2 ²	0.83	0.76	0.93	0.84	
			B coefficier	nt (± SE)		
Inte	rcept	32224.9** ± 7617.4	36645.9* ± 10517.8	149163.4*** ±18885.1	160044.6*** ± 32274.6	
95% CI	Upper Lower		60081.1 13210.7	191242.0 107084.8	231956.8 88132.4	
Ye	ear	-15.6 ** ± 3.8	-17.5* ± 5.2	-73.2*** ± 9.4	-78.2*** ± 16.1	
95% CI	Upper Lower		-5.8 -29.2	-52.3 -94.3	-42.2 -114.0	
	adybug ject	$0.0003 \\ \pm 0.0236$	0.0221 ± 0.0326	$0.0870 \\ \pm 0.0586$	0.2587* ± 0.1001	
iNatu	ıralist	0.0008 ± 0.0013	0.0012 ± 0.0018	0.0061 ± 0.0032	0.0105 ± 0.0055	
bugGuide.Ne		0.0680 ± 0.2646	-0.0031 ± 0.3653	-0.7286 ± 0.6560	-1.9637 ± 1.1212	

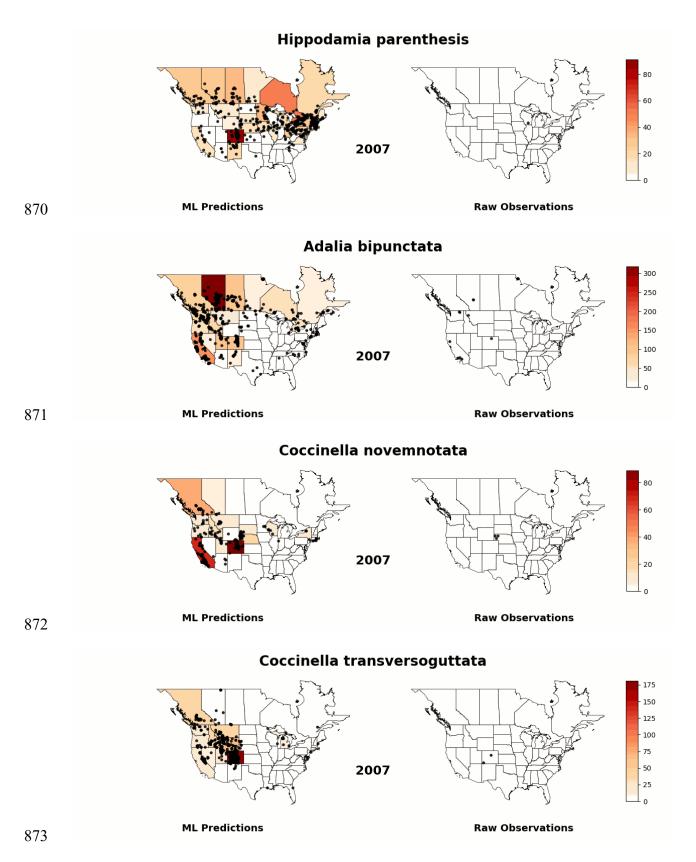


Figure S6. Annual occupancy maps (2007–2021) for each species; Left maps, occupancy predicted by co-occurrence–based models. Right maps, occupancy based on reported observations; Heatmap colors, the number of occupied coordinates per state, with temporal changes in color intensity reflecting shifts in occupancy (Active figures are available at: https://figshare.com/s/17cef8ef530f0a4f7b99).