

RESEARCH

Open Access



Resolving biology's dark matter: species richness, spatiotemporal distribution, and community composition of a dark taxon

Emily Hartop^{1,2*}, Leshon Lee^{3,4}, Amrita Srivathsan^{3,5}, Mirkka Jones^{6,7}, Pablo Peña-Aguilera⁸, Otso Ovaskainen^{6,9}, Tomas Roslin^{6,8} and Rudolf Meier^{5,10*}

Abstract

Background Zoology's dark matter comprises hyperdiverse, poorly known taxa that are numerically dominant but largely unstudied, even in temperate regions where charismatic taxa are well understood. Dark taxa are everywhere, but high diversity, abundance, and small size have historically stymied their study. We demonstrate how entomological dark matter can be elucidated using high-throughput DNA barcoding ("megabarcoding"). We reveal the high abundance and diversity of scuttle flies (Diptera: Phoridae) in Sweden using 31,800 specimens from 37 sites across four seasonal periods. We investigate the number of scuttle fly species in Sweden and the environmental factors driving community changes across time and space.

Results Swedish scuttle fly diversity is much higher than previously known, with 549 putative species detected, compared to 374 previously recorded species. Hierarchical Modelling of Species Communities reveals that scuttle fly communities are highly structured by latitude and strongly driven by climatic factors. Large dissimilarities between sites and seasons are driven by turnover rather than nestedness. Climate change is predicted to significantly affect the 47% of species that show significant responses to mean annual temperature. Results were robust regardless of whether haplotype diversity or species-proxies were used as response variables. Additionally, species-level models of common taxa adequately predict overall species richness.

Conclusions Understanding the bulk of the diversity around us is imperative during an era of biodiversity change. We show that dark insect taxa can be efficiently characterised and surveyed with megabarcoding. Undersampling of rare taxa and choice of operational taxonomic units do not alter the main ecological inferences, making it an opportune time to tackle zoology's dark matter.

Keywords Dark taxa, Megabarcoding, DNA barcoding, Biodiversity discovery, Hierarchical Modelling of Species Communities, Diptera, Phoridae

*Correspondence:

Emily Hartop
emily.hartop@ntnu.no

Rudolf Meier
Rudolf.Meier@hu-berlin.de

Full list of author information is available at the end of the article



© The Author(s) 2024, corrected publication 2024. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

Background

If we go on the way we have, the fault is our greed and if we are not willing to change, we will disappear from the face of the globe, to be replaced by the insect. Jacques Yves Cousteau.

Biodiversity loss in the Anthropocene is driven by changes in climate, land use, but also species introductions [14, 19]. Such loss can result in the concomitant decline in ecosystem services vital to society [47, 49] and may ultimately disrupt global supply chains and food security [81]. Accurate monitoring of biodiversity is therefore a global priority [54].

A crucial first step to monitoring biodiversity is obtaining robust quantitative baseline data. Most biologists would consider this to be data on species diversity, abundance, and biomass for those taxa that contribute substantially to these quantitative metrics. However, most biodiversity studies cover only a few well-studied groups that are relatively easily identified and quantified (e.g. birds, mammals, amphibians, bees, and butterflies). Such charismatic taxa are then used as proxies for all taxa in a region [7, 8, 33, 56, 57, 61, 70, 79] instead of basing our understanding of global biodiversity on a broad and unbiased representation of biodiversity covering a wide range of traits and responses to environmental change [2, 21, 37].

A key component of global biodiversity are “dark taxa”, i.e. taxonomic groups for which less than 10% of the diversity is described and the species diversity is estimated to be upwards of 1000 [40]. Such poorly known groups do not just inhabit inaccessible realms like the deep sea [60] but also the terrestrial habitats in which we live. A recent study [66] revealed that 20 insect families (of which 10 belong to Diptera) account for > 50% of local species diversity. Alarming, the very same families suffer from extreme taxonomic neglect and are therefore poorly represented in biodiversity surveys. Identifying and tackling the diversity of these dark taxa with scalable techniques thus emerges as an urgent priority for biodiversity science.

Large-scale studies of dark taxa have only become feasible in recent years due to advancements in sequencing technologies coupled with efficient single-specimen DNA barcoding workflows [16, 40, 51, 68, 69, 82]: “megabarcoding”. Such workflows allow large numbers of specimens to be processed and sorted to putative species (mOTUs), while providing exact specimen counts, and vouchers for subsequent morphological, taxonomic, and biological work.

Scuttle flies in the family Phoridae (Diptera) have been considered the seventh most speciose and abundant insect family globally [66]. However, to date only

ca. 4000 species have been described, although their actual diversity may be two orders of magnitude larger [68]. In addition to their extreme species richness, scuttle flies occupy a wide range of ecological niches, containing species that are herbivores and predators to scavengers, parasitoids, and parasites (reviewed by [20]). Nonetheless, previous ecological studies focusing on scuttle flies [23–30] have been of limited scope due to time-consuming morphological identification methods. It is essential that we assess taxa like these to ensure that ecological analyses reflect the bulk of biodiversity and represent a broad range of ecological niches.

In this study, we use sorting with megabarcoding to generate the data for answering fundamental questions about this dark taxon. We ask how many species of scuttle flies occur in Sweden, how are they distributed across time and space, and what environmental variables drive their distribution. To test whether the choice of species and species delimitation method will affect the results, we carry out the same analyses using both mOTUs as species-proxies and haplotype diversity, and test whether rare species influence the overall inferences. We show how “dark taxon zoology” can quickly yield the answers urgently needed in the Anthropocene.

Results

Diversity

We obtained a total of 31,739 *COI* barcodes belonging to 2697 haplotypes from scuttle fly samples from 37 sites (Fig. 1, Additional file 1: Fig. S1) and four seasonal time periods (Additional file 1: Fig. S2) across Sweden (Additional file 1: Table S1). At the species threshold (1.7%) [40], we detected 549 mOTUs. Species accumulation curves revealed that scuttle fly species diversity was incompletely sampled overall and across sites, horticultural zones, and time periods (Additional file 1: Figs. S3, S4). Between 38 and 145 species were observed per site, with Chao1 richness estimates per site ranging from 83 to 244 species (Fig. 2b, c). Regional richness estimates varied between nearly 400 species estimated in the southernmost coastal zone (zone 1), to fewer than 200 estimated for the alpine zone (Additional file 1: Fig. S3). Midsummer and late-summer time periods were both characterised by richness estimates of around 450 species, while the late spring and offseason time periods showed lower species richness estimates of around 300 species (Additional file 1: Fig. S4). Total species richness in Sweden was estimated at between 652 species (for Chao1, Fig. 2a) and 713 species (combined non-parametric estimator (CNE) in [63]), suggesting that 100–160 species of scuttle flies remain to be sampled.

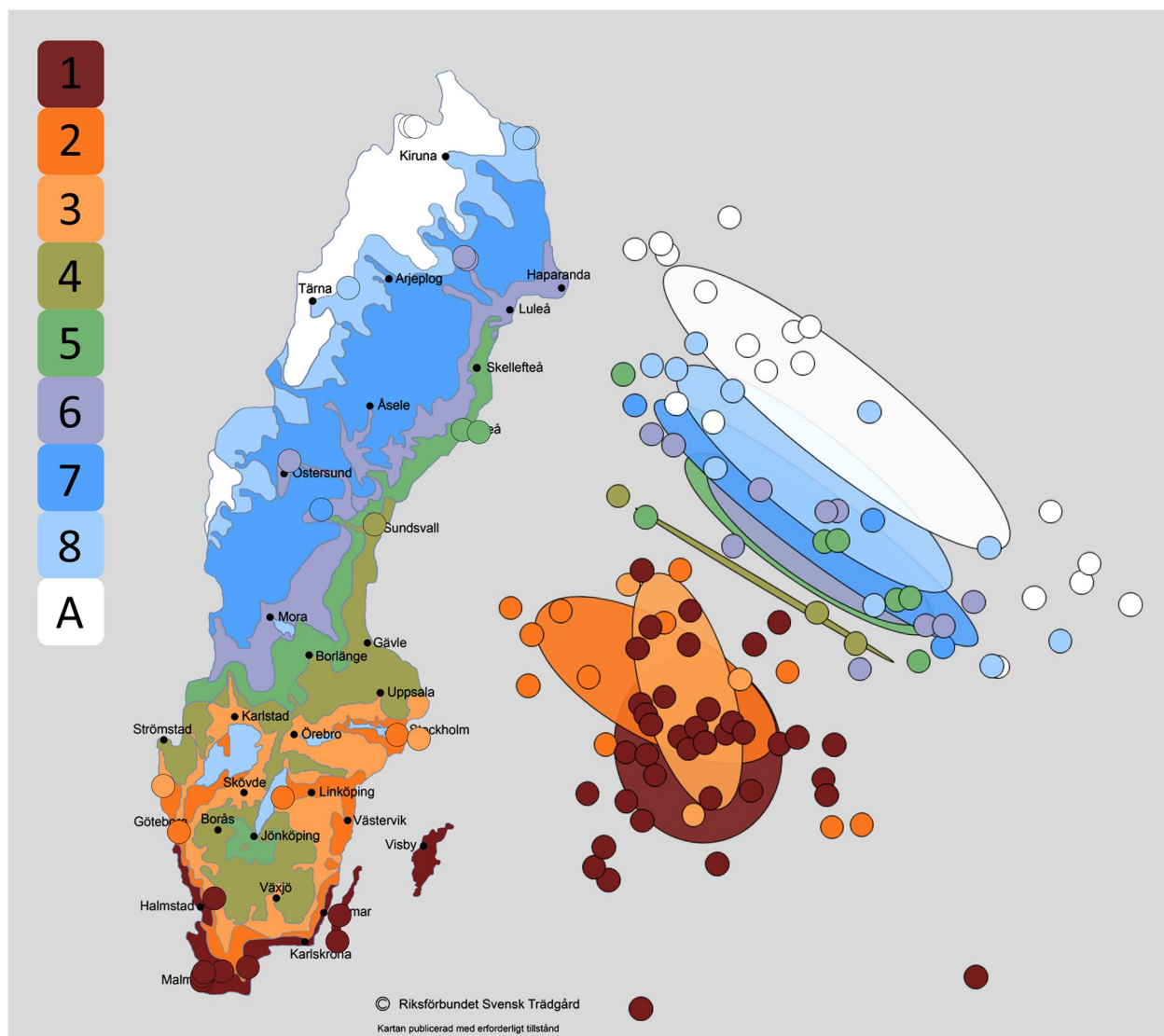


Fig. 1 The location of the 37 study sites of the Swedish Insect Inventory Project colour-coded according to the plant hardiness zones (odlingszoner, 1–8 and alpine) of the Swedish Horticultural Society (Riksförbundet Svensk Trädgård) (map used with permission) next to an NMDS plot of study samples colour-coded with the same zones

Ordinations of community composition

NMDS plots revealed a clear distinction between scuttle fly communities in the southern (zones 1–3) and northern (zones 4–alpine) plant hardiness zones (Fig. 1, Additional file 1: Fig. S5). ANOSIM supported significant differences between southern and northern zones at this threshold ($R=0.58$, $p=0.001$), and SIMPER revealed that north–south similarity was just 13.2%, as compared to similarities of 24.6% and 23.2%, respectively, among samples within the northern and southern zones (Additional file 1: Tables S2–S3).

The separation between zones was consistent across clustering thresholds ranging from haplotypes-as-such to

a threshold of 1.7% sequence similarity, with stress values around 0.21 (Additional file 1: Fig. S5 top row). Above a threshold of 3%, however, the patterns were increasingly blurred and the stress values were higher (0.25–0.27) (Additional file 1: Fig. S5 bottom row). These patterns were also evident in ANOSIM analyses, where the sample statistic decreased from 0.54 for haplotypes to 0.27 for 5% mOTUs, indicating a decreasing relationship between scuttle fly community composition and plant hardiness zones at higher clustering thresholds (Additional file 1: Table S2). Similarly, in SIMPER analyses, average similarity between zones increased from 20.2% for haplotypes to 41.4% for 5% mOTUs (Additional file 1: Table S3). For

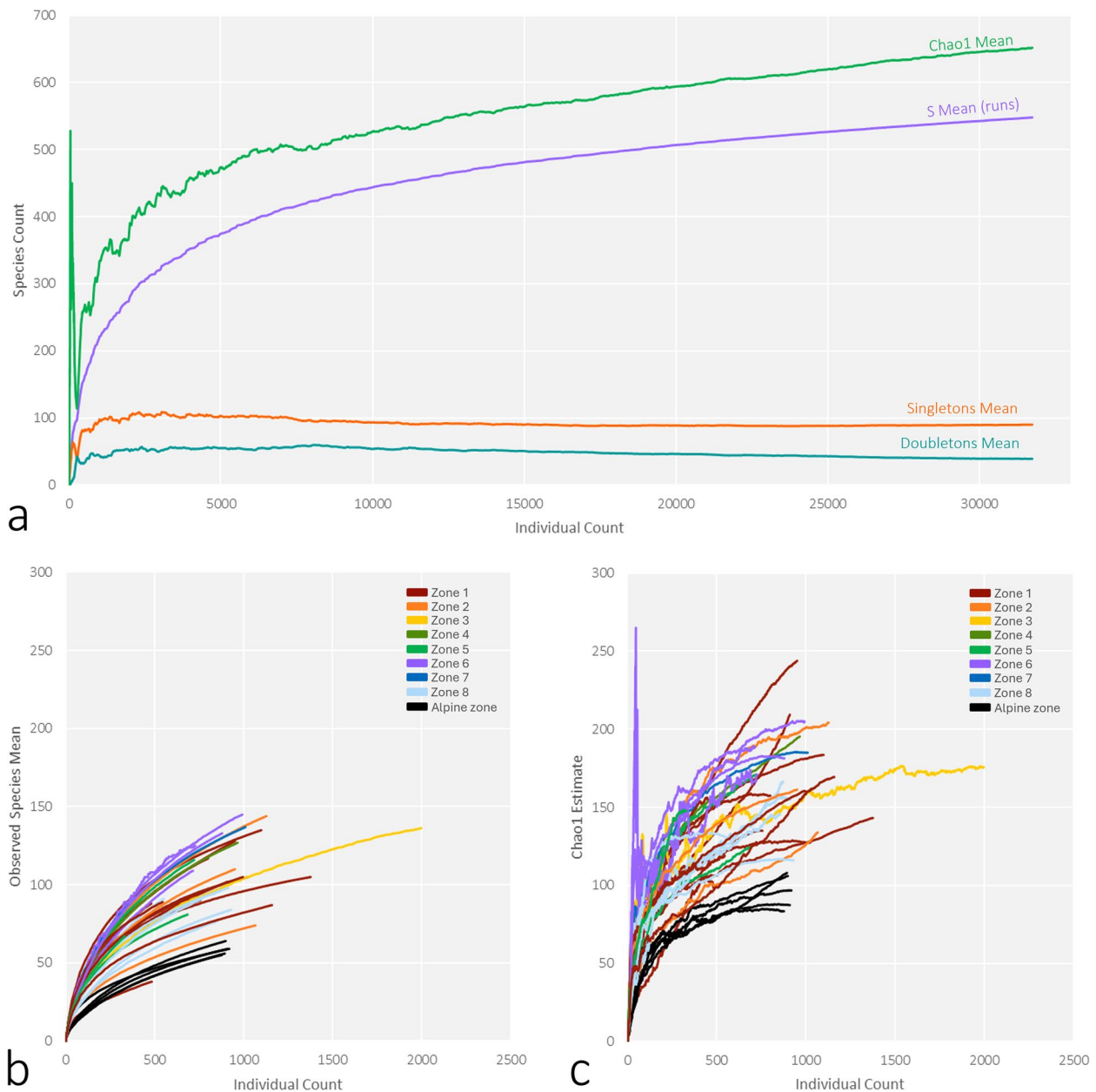


Fig. 2 **a** Species accumulation and Chao1 estimate curves for the scuttle fly dataset across Sweden. Notably, current sampling is far from exhaustive, and numbers of singletons and doubletons in our dataset are high regardless of sample size, **b** species accumulation curves by sampling sites, colour-coded by zone, and **c** Chao1 estimate curves by sampling sites, colour-coded by zone. For a map of the zones using the same colour codes, see Fig. 1

species, scuttle fly communities were found to be distinct across most zones ($R=0.51$, $p=0.001$), with higher (average=28.5%) similarity within zones than between (average=19.2%) zones (Additional file 1: Tables S2, 3).

Samples from the late spring, midsummer, and late-summer time periods showed a clear progression along the ranked plant hardiness zones, while off-season samples appeared more randomly distributed

(Additional file 1: Fig. S6). Northern sites (IV–alpine) showed higher distinctness across seasons ($R=0.59$) than when considering all sites ($R=0.27$) (Additional file 1: Table S4). The between-season similarity of all sites averaged 16.9% and within time period similarity averaged 23.7% (Additional file 1: Table S5, top). However, for northern sites only (IV–alpine), between time period similarity was 21.3% and within time

period similarity was 37.2% (Additional file 1: Table S5, bottom).

Hierarchical Modelling of Species Communities (HMSC)

Models of the four response matrices (species-occurrence, haplotypes-occurrence, species-abundance, haplotypes-abundance) showed relatively good MCMC convergence with potential scale reduction factors of the models' beta and omega parameters close to one (Additional file 1: Fig. S7).

Explained variance averaged c. 30% for the occurrences and c. 60% for the abundances of both species and haplotypes, but there were large differences in model fit among taxa (Additional file 1: Fig. S8), haplotype occurrence mean \pm sd $T_{\text{jur}} R^2 = 0.30 \pm 0.15$ (range 0.05–0.82); haplotype abundance $R^2 = 0.61 \pm 0.26$ (range 0.00–1.00); species presence-absence $T_{\text{jur}} R^2 = 0.30 \pm 0.14$ (range 0.05–0.72); and species abundance $R^2 = 0.59 \pm 0.26$ (range 0.02–1.00).

Sampling time period explained the largest fraction of variance, on average, in all four HMSC models (mean 10–11% in the occurrence models; mean 15–16% in the abundance models; Fig. 3, Additional file 1: Fig. S8). Mean annual temperature explained almost as much variance as sampling time period in the occurrence models (mean 7% and 8% in the species vs haplotype models), but clearly less than sampling season in the abundance models (mean 6% in the species models and 5% in the haplotype models) (Fig. 3, Additional file 1: Fig. S8). All response variables showed strong climatic (including temporal) community structure and there was also strong spatial structure in species and haplotype occurrences linked with the latitudinal temperature gradient across Sweden. Tree cover explained more variance in the abundance models (mean 6% for both species and haplotypes) than in the occurrence models (mean 1% for both species and haplotypes). Trapping effort also explained more variance on average in the abundance models (mean 7% and 8%) than in the occurrence models (mean 1.5% for both), as did the effect of having (vs not having) sequenced the full trap sample (2% mean for occurrence, 6% mean for abundance).

While the occurrences of taxa and haplotypes showed statistically supported responses to all model covariates, their abundances showed responses less frequently, and responses with strong support were primarily with seasonal covariates (Fig. 4).

Most species and haplotypes with a statistically supported seasonal abundance trend peaked in the late spring. Compared to late spring, the midsummer fauna showed a reduction of 28% in species counts and 16% in haplotype counts, respectively, with a further reduction of 24% in both species counts and haplotype counts

towards the late summer. Taxa and haplotypes were also usually more prevalent in late spring than in mid- or especially late summer. However, a minority of species and haplotypes (9% and 7%, respectively) showed the reverse pattern, being more prevalent in late summer than in the spring. Offseason captures in the late fall to early winter were consistently low, and the number of off-season samples included in the models was the smallest ($n=20$ samples). Nonetheless, 33% of the species modelled and 29% of the haplotypes modelled were occasionally detected in off-season samples.

Taxon occurrences showed a mixture of positive (29% vs 30% for species and haplotypes, respectively) and negative (18% vs 23%) responses to the annual temperature gradient ("bio1") across Sweden (Fig. 4), reflecting the broad-scale compositional changes from southern towards northern Sweden seen in the NMDS ordinations. Where detected, the occurrence responses of taxa to forest or woodland cover were more often positive than negative (10% positive vs 4% negative for species; 8% vs 2% for haplotypes; Fig. 4). Longer trapping periods did not consistently result in higher detection probabilities of taxa, presumably due to seasonal differences in trapping duration (4% positive vs 5% negative responses for species; 8% positive vs 4% negative for haplotypes; Fig. 4). As expected, most taxa (81–82% in both the species and haplotype models) showed a negative occurrence response to the binary variable indicating whether all specimens in a sample were sequenced or not ("Full-Sample") (Fig. 4). The mean annual temperature gradient across Sweden was predicted to affect the prevalence of 38% of species, but does not appear to be a main driver of species abundance (Fig. 4). The predicted effect of the temperature gradient on species prevalence during late summer was positive in 22% of taxa and negative in 16% of taxa (Additional file 1: Fig. S9). Beyond spatial patterns explained by these climatic predictors, there was evidence of localised spatial autocorrelation in species and haplotype site occupancies at scales of less than 40 km. Neither species nor haplotype abundances were spatially autocorrelated, nor did we detect statistical support for temporal autocorrelation in any model.

Residual correlations in the distributions of taxa were detected among sites and samples in both the species and especially the haplotype occurrence models (Additional file 1: Fig. S10). Residual associations of taxa over time were also evident, but less frequently (Additional file 1: Fig. S10). Residual associations between taxon/haplotype occurrences likely indicate that our models either lack or imperfectly represent some of the variables that structure the occurrences of scuttle fly taxa and haplotypes in space and time. No residual associations among taxa were evident among sites or over time in the species or haplotype

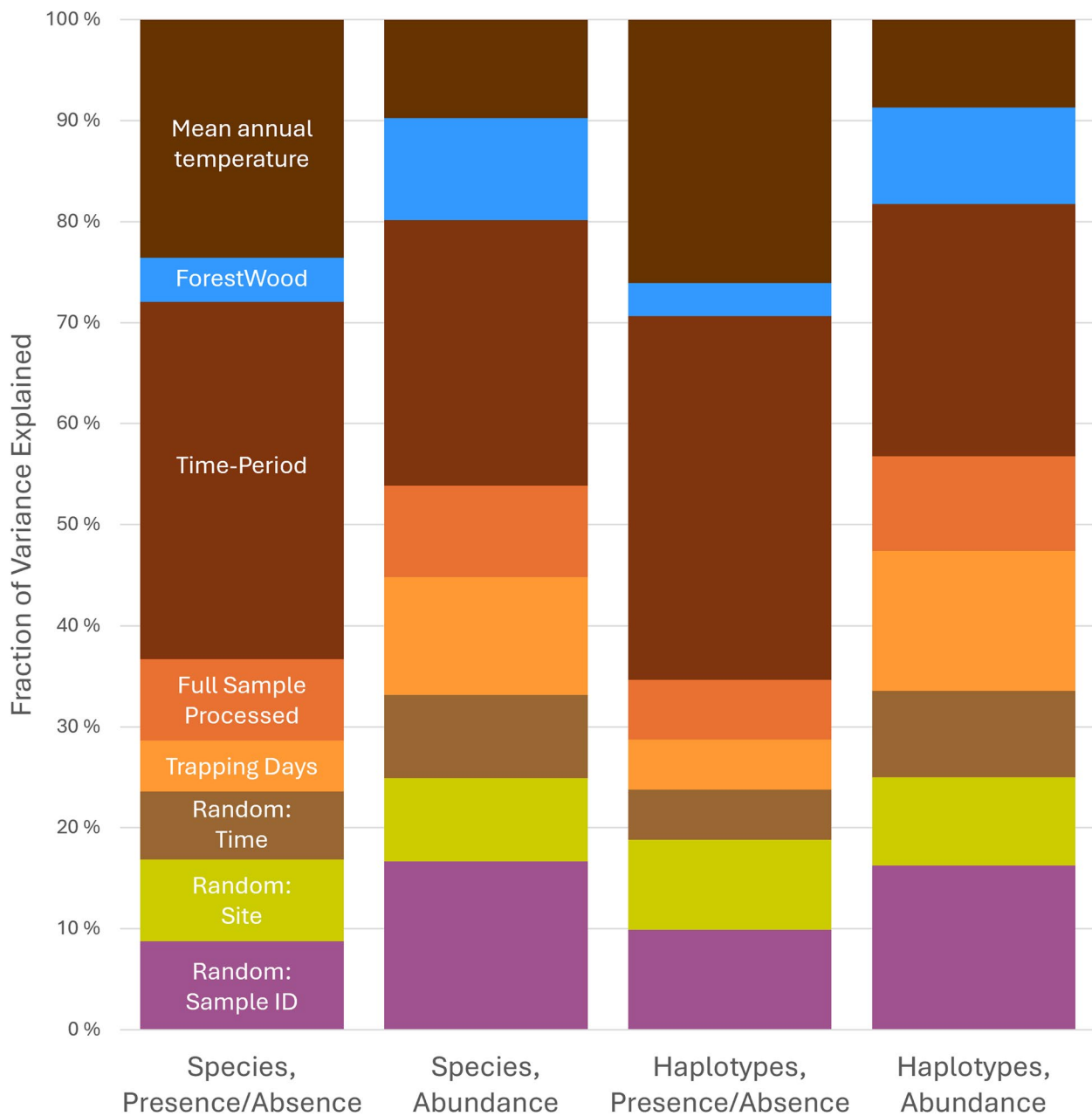


Fig. 3 Summary of explained variance in the occurrences and abundances of scuttle fly haplotypes and species across samples (fractions within bars represent the average percentage of variance explained by each fixed or random effect in the models). All four models show strong climatic structure (in shades of brown) on scuttle fly communities, as captured by fixed variables describing sampling season and mean annual temperature and a random effect based on median sampling date. The spatial fraction of explained variance (lime green) reflects community structural differences among sites that were not captured by the fixed effect covariates. Differences in sampling effort, i.e. whether or not trapped flies were all sequenced or not and the number of field trapping days per sample, also affect the predictability of community structure (shades of orange). The abundance of species and haplotypes, and to a lesser extent their occurrence, was also strongly structured by habitat type as described by forest and woodland cover (blue). Finally, we included a categorical random effect representing sample identity (purple)

abundance models, and very few were detected among samples (Additional file 1: Fig. S10). Hence, covariance in the abundances of taxa and haplotypes across occupied samples was well modelled by the environmental and other covariates in these models.

The observed richness of the excluded rare vs modelled common species and haplotypes was strongly positively correlated ($R=0.59$ for species and $R=0.80$ for haplotypes). Compositional differences between samples in terms of the rare vs common species and haplotypes were

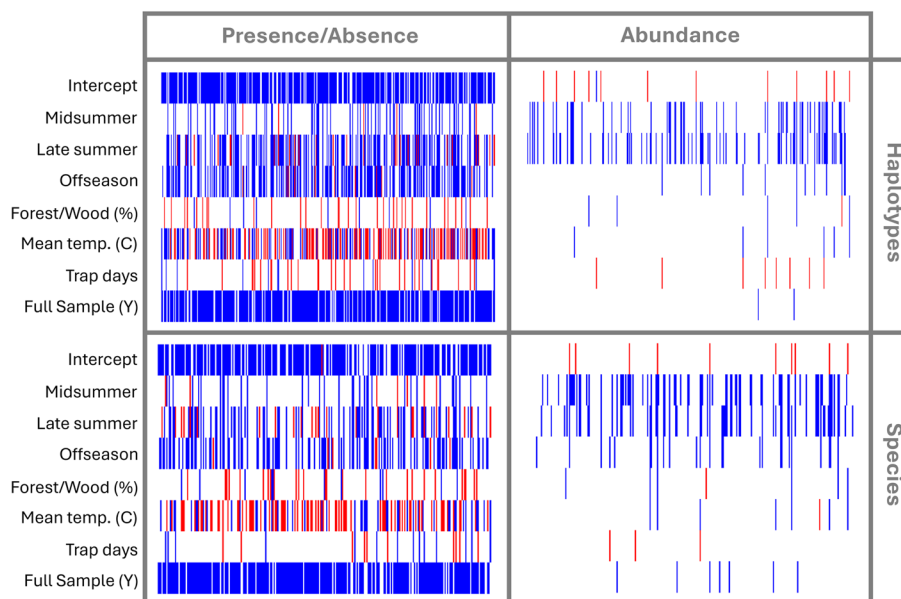


Fig. 4 Predicted occurrence and abundance responses of 193 scuttle fly haplotypes and 162 species to HMSC model covariates. Positive (red) and negative (blue) estimated responses with a posterior probability of 0.95 are illustrated. The three rows below the intercept illustrate the estimated effects of three levels of a categorical variable representing sampling time period (midsummer, late summer, offseason) relative to the baseline (late spring). The subsequent three rows represent responses to % forest or woodland cover and mean annual temperature (“bio1”) at sampling sites. The final two rows represent the effects of two sampling-related differences among samples: the number of trapping days and whether all specimens in the trap sample were sequenced, and hence available for HMSC analysis, or not

also positively correlated (for taxon presence-absence Mantel $R=0.31$ for species and Mantel $R=0.64$ for haplotypes; for taxon abundance Mantel $R=0.30$ for species and Mantel $R=0.60$ for haplotypes).

Community dissimilarity

Turnover accounted for the bulk of spatial, temporal, and spatiotemporal variation in community dissimilarity (Fig. 5). In the spatial analyses, the mean turnover of species and haplotypes between sample pairs was 0.75 and 0.87, respectively, while the corresponding mean values of nestedness were 0.07 and 0.03. In temporal analyses, similarly, the mean values of turnover were 0.45 and 0.74 for species and haplotypes, respectively, and the corresponding means of nestedness were 0.13 vs 0.07. Finally, mean turnover values for the spatiotemporal analyses were 0.75 and 0.84 vs a mean nestedness of 0.07 and 0.04 for species and haplotypes, respectively.

Regardless of the time period and operational taxonomic units used, we found a significant positive correlation between turnover and geographical distances, meaning that communities in closer proximity to each other are more similar in composition (Additional file 1: Fig. S11). The mean spatial turnover of species communities was highest in the summer time periods

(midsummer: 0.82, late summer: 0.84) and lower in late spring and offseason (0.78 and 0.72, respectively). However, we did not find any significant correlation between temporal distances (difference in mean week) and turnover (see Table S6). Patterns of nestedness showed no detectable correlation with any of the distances explored (except for the grouped temporal distance) (Additional file 1: Table S6). Within each time period, scuttle fly communities become more distinct from neighbouring communities from late spring to late summer (Additional file 1: Table S6 and Figs. S11–13).

Discussion

Our study marks the first country-wide examination of a dark taxon’s diversity and distribution. It revealed more than 500 species of scuttle flies based on processing ca. 31,800 specimens. The ecological analyses suggest that climate change will have profound (and quickly apparent) effects on communities of scuttle flies that could serve as early indicators of future shifts in the environment. We demonstrate that armed with recent advancements in sequencing technologies, bioinformatics pipelines, and molecular barcoding workflows [40, 52, 68, 69, 75, 82], we are now able to resolve patterns of alpha diversity, spatial and temporal

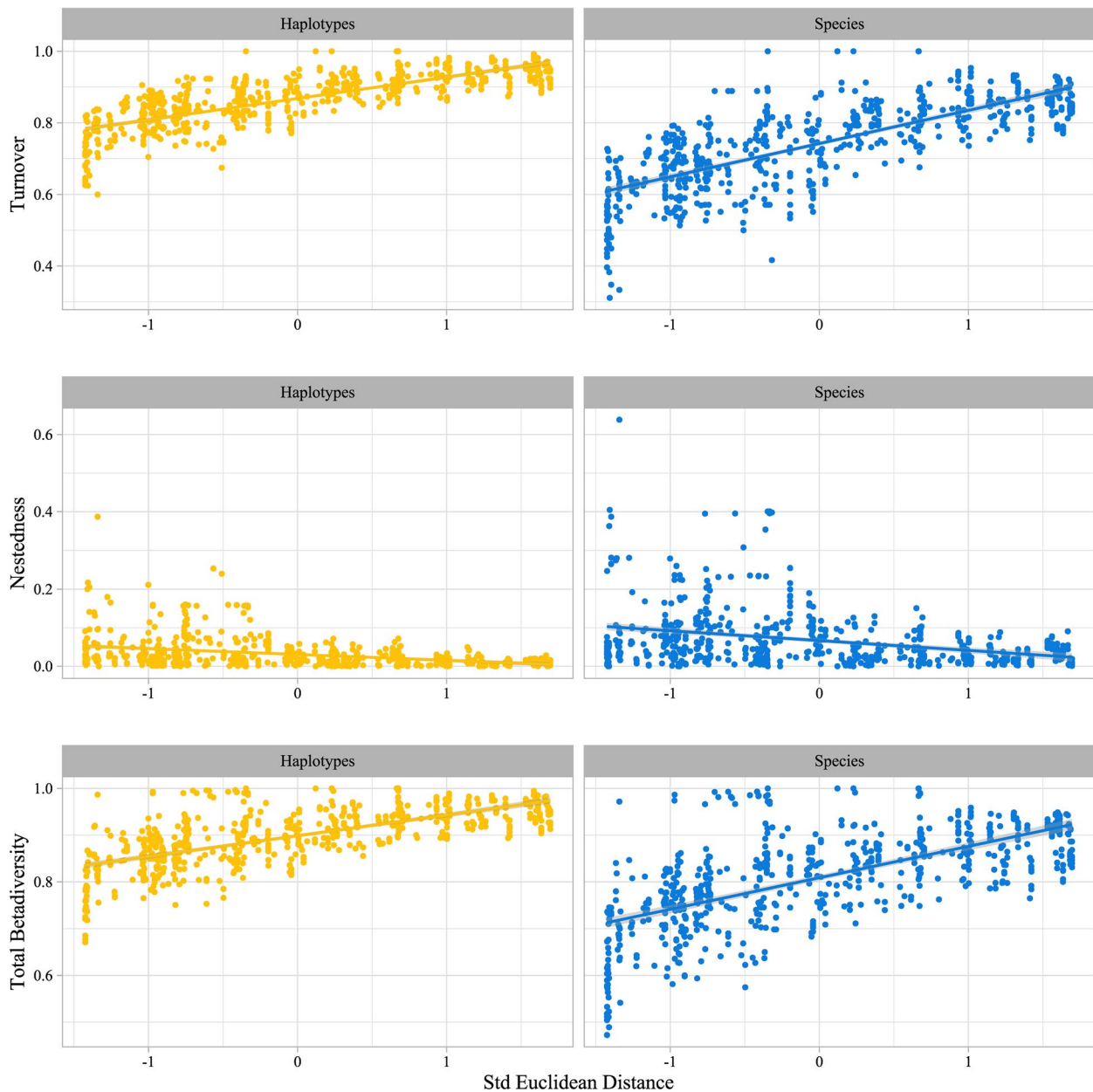


Fig. 5 Pairwise community dissimilarity among scuttle fly species and haplotype samples. Following Baselga and Orme [5], we partitioned overall dissimilarity into its turnover and nestedness components and illustrate these as well as overall dissimilarity as a function of distance (standardised Euclidean distance calculated from sample site coordinates)

turnover, and species communities of challenging dark taxa. We show that for the scuttle flies in Sweden, many species remain undiscovered, that local communities show major turnover in space and time, and that ecological patterns are largely robust to the finer details of species delimitation. Below, we will discuss each of these findings in further detail—and the importance of these patterns in dark taxon biology to biodiversity science.

Diversity

Our sample only scratched the surface of the scuttle fly fauna of Sweden (Fig. 2, Additional file 1: Figs. S3, S4). While the true Swedish fauna of scuttle flies is still hard to estimate, the 549 putative species found (and 652–713 predicted based on current sampling) greatly exceed the 374 previously documented from the country. Previous estimates for the scuttle fly fauna have ranged from 1100 to nearly 2000 species, suggesting

that the true numbers may be even higher [63]. This would not be surprising, given that our study included samples from just 37 traps in a country of over 450,000 km²; i.e. many habitats remained unsampled. However, the available data confirm that for dark taxa the basic pattern in ecology holds that most species are rare and that finding all species would require massive sampling [13].

Where the remaining species will be lurking is unknown. While Malaise traps are remarkably effective at capturing flies [42, 64], other trapping methods would certainly reveal additional species. An “All Diptera Biodiversity Inventory” conducted in Costa Rica utilised a wide variety of methods [9], revealing that 59% of fly species were unique to a specific collection method [10]. Intriguingly, species-rich taxa with particularly interesting and often impactful biologies (e.g. parasitoids) may be underrepresented in Malaise trap samples. For example, Brown [12] reported that nearly half of the species of ant-parasitising scuttle flies were caught exclusively over army ant raids. Similar results were obtained by Disney et al. [22], who found some species of scuttle flies were uniquely or primarily collected in water traps rather than Malaise traps. Another possible frontier is the canopy. Ongoing research in upper levels of the canopy in the Amazon suggests that it contains many scuttle fly species not found in lower sampling elevations [36], contrasting with previous data showing scuttle flies as a taxon primarily collected at ground level [44]. Based on a limited sample (159 specimens) from canopy trap prototypes run concurrently with the Malaise trap sampling for this study, we found a single unique mOTU not seen in the 31,794 scuttle flies from ground level, indicating that even in a temperate environment we may have more to find in the canopy.

Thus, no single sampling method will suffice to reveal all taxa and additional trapping sites that target regions or habitats un- or under-represented in the current study will also be needed for an exhaustive inventory of scuttle flies. Some of this complementary sampling may have already been carried out by the Swedish Malaise Trap Project years ago [42]. An excellent follow-up to this study would be sequencing scuttle fly samples from that sampling campaign, or from another more recent insect campaign in Sweden—the Insect Biome Atlas project (www.insectbiomeatlas.org).

In addition to focusing future efforts on revealing more species, the analysis of other life stages may offer a more nuanced understanding of scuttle fly communities. Our study is based on adult flies. Patterns for the longer-living larval stages are at this point unknown, as they are not easily collected.

Lessons from ecological analyses

Our ordinations of community composition suggested strong structuring of scuttle fly communities by Swedish horticultural zones (Fig. 1) and season (Additional file 1: Fig. S6). Clear north–south structuring of the Swedish insect fauna has also been observed in damselflies, mosquitos, and caddisflies [38, 48, 77], and strong phenological patterning is a hallmark of any high-latitude fauna [80]. Our ordinations also indicated that spatial structuring was largely independent of mOTU clustering threshold at or below the species proxy level (Additional file 1: Fig. S5).

To address these findings more rigorously and without the constraints of pre-determined zones, we implemented HMSC. This confirmed that scuttle fly community variation is both highly seasonal and strongly tied to the latitudinal temperature gradient over Sweden, with all models strongly predicted by climatic covariates (Fig. 3, Additional file 1: Fig. S8). The presence of significant spatial autocorrelation in the HMSC occurrence models reflects the gradation of scuttle fly distributions across space. Scuttle fly communities also showed clear compositional changes over time during the warm season, from late spring through mid- and late summer, while offseason sampling was too inconsistent to reveal any patterns (Additional file 1: Figs. S2, S6). Consistent with our findings that scuttle fly communities are driven by climatic covariates, we found more rapid spatial turnover in taxa at the species than at the haplotype level (Fig. 5). Recent work has predicted that when dispersal limitation is the dominant driver of species distributions, the rate of spatial turnover of biological communities will be similar at both the haplotype and species levels [4]. Conversely, when environmental conditions strongly constrain species ranges, community similarity is predicted to decay at different rates across genealogical scales.

Our results have several implications in the face of climate change. Individual taxa show a mixture of positive and negative responses to mean annual temperature, and to a lesser extent to seasonality. This implies that we may see a substantial number of both climate change “winners and losers” in the future, as species ranges and phenology expand, contract, or shift (Fig. 4, Additional file 1: Fig. S9). Adult scuttle flies are ephemeral—with short lifespans and high turnover and mobility—they may respond rapidly and serve as indicators for future shifts of other taxa and in the environment more broadly. Our observation of steeply declining abundance from late spring through the summer into the offseason may partially reflect the extreme temperatures and drought in Europe in summer 2018 [1, 6]. While this suggests that our seasonal results may be atypical, they are perhaps also a sign

of summers to come, as climate change increases the frequency and severity of these events [39].

With our initial results confirmed by HMSC, the predictive value of the simple plant hardiness map for scuttle fly distributions offers excellent news. It suggests that even such a relatively coarse-grained tool can be used as a reliable indicator of the regional compositions of scuttle flies. That this is the case is only intuitive: Many scuttle flies exist close to ground level and therefore, like plants, their distributions may be closely tied to soil temperature, acidity, moisture, or composition (all of which would be interesting covariates to explore in future modelling). Additionally, some species are known to have direct interactions with plants, fungi, and ground-dwelling insects—factors which may again tie them to the microclimate at ground level. Future studies might focus on these microenvironmental variables, to address the unexplained variance in scuttle fly communities from this study.

Dark taxon zoology

Dark taxa have historically been largely ignored due to the complexities involved in studying groups of highly diverse and abundant organisms of small size. How, then, do we make the study of dark taxa efficient and start “dark taxon zoology” to respond to the need for quantitative data on biodiversity?

To address this, we assessed whether the taxa analysed need to exactly match taxonomically validated species for ecological hypotheses to be valid. If not, we may avoid endless discussions about species delimitation and proceed with ecological analyses. Promisingly, our current results suggest that ecological patterns are largely robust to the molecular clustering thresholds used. Overall biodiversity patterns, patterns of community turnover, and drivers of distribution were virtually identical using thresholds from 0% (haplotype data) to 1.7% (Figs. 3, 4, 5, Additional file 1: Figs. S5, 8, 11–13), which was the best threshold for obtaining congruence between barcode clusters and morphospecies in a previous study [40]. It is important for researchers to calculate the appropriate species proxy threshold specific to their own taxa, as this value may vary depending on the group being analysed. Exceeding an appropriate species proxy threshold can obscure patterns by lumping species together (Additional file 1: Fig. S5, bottom row). The promising results using haplotype data are convenient because they do not require any taxonomic decisions, which may be appealing to molecular ecologists.

A second potential stumbling block in the study of dark taxa is the large numbers of rare species. If understanding basic ecological patterns is dependent on these rare species, dark taxon zoology will be difficult. Our results

indicate that both richness and compositional differences between samples of rare versus common species and haplotypes were both positively correlated. This suggests that more common and rarer species respond similarly to the same drivers. Again, this is excellent news, since it suggests that ecological inferences regarding the drivers of species distribution and community composition can be based on the more common species—which are much easier to detect.

Conclusions

We here argue for a dark taxon zoology and illustrate that it is not a hopeless undertaking by targeting scuttle flies as a typical dark taxon in Sweden. We sample across the entire country; we estimate the species richness, resolve patterns of distribution across time and space, and pinpoint environmental features that drive these distributions. Our results suggest that such assessments will be insensitive to specific taxonomic cut-offs and robust to undersampling of rare taxa. Overall, the study hopes to contribute to a more quantitative approach to biodiversity. In the future, advances in molecular workflows, bioinformatics, robotics, and automation will make these groups increasingly efficient to study [51]. We hope that our case study will serve to propel forward dark taxon zoology, bringing the main part of diversity into the realm of biodiversity science.

Methods

Target taxon: the scuttle flies of Sweden

Sweden has one of the best-known animal faunas in the world due to efforts dating from Carl Linnaeus to the Swedish Taxonomy Initiative and Malaise Trap Project [42, 46, 53], but see [63]. While 374 species of scuttle flies have been documented in the country [65], this is an underestimate. For comparison, a single suburban garden in Cambridge, UK, yielded nearly 100 species of scuttle flies [11], while backyards in Los Angeles, CA, can support up to 82 species [11]. Previous estimates of scuttle fly diversity in Sweden have proposed that the true fauna may approach 2000 species (CNE estimate in [63]), but this may be an overestimate based on an error-prone process of morphological identification [63].

Sampling

To start resolving the species richness, spatiotemporal distribution, and community composition of Swedish scuttle flies, we sampled communities of flying insects at 37 locations across Sweden (Fig. 1, Additional file 1: Fig. S1, Table S1). These samples were collected by the Swedish Insect Inventory Project (<https://www.stationlinne.se/sv/forskning-research/the-swedish-insect-inven>

tory-project-siip/) [42] in ~80% ethanol using Townes-style Malaise traps [73]. Scuttle flies were sorted from the trap samples and preserved in ~80% ethanol at -20°C .

Sampling started in May 2018 and continued into the following year. Although trapping was continuous, the sample periods varied across sites due to their extensive distribution across Sweden. Three distinct time periods that aligned with the warm season phenology, along with one offseason period, were established sequentially for each site. These periods corresponded as closely as possible to late spring (sampled in May), midsummer (sampled in late June/over the summer solstice), late summer (sampled in late July and early August), and the offseason (sampled from September onward) (for specifics and exceptions see Additional file 1: Fig. S2, Table S1). Some offseason samples could not be retrieved in late 2018 and were collected in 2019. To ensure uniformity in the number of weeks per sampling period and avoid introducing biases, we attributed the latest sampling date for the off-season traps in 2018 to those samples collected in 2019. Additionally, late spring samples were only obtained from 25 sites, as 12 sites were not installed until later in the season (see Additional file 1: Table S1 for sampling details). To compare richness estimates across seasons, we plotted accumulation curves for each time period excluding the 12 traps that were not yet installed in late spring.

Sequencing and bioinformatics

A total of 136 samples were selected for analysis. Most samples (94) contained thousands of scuttle fly specimens and subsampling was needed. Individuals were randomly selected with a minimum of two 96-well microplates of scuttle flies processed per sample. This resulted in a total of at least 190 specimens extracted per site and time period for most samples. Thirty-one samples containing fewer than 190 specimens and were thus processed completely, and two samples were found to contain no scuttle flies. Additional plates of specimens from nine samples that had been processed in an earlier study [40] were also included (Additional file 1: Table S1). All sample information, including numbers of barcodes obtained per sample and whether a sample was fully processed, is found in Additional file 1: Table S1.

DNA were obtained from the specimens using 10 μl of HotSHOT lysis buffer [74]. The extraction was carried out in a thermocycler at 65°C for 18 min, then for 2 min at 98°C . Amplification was carried out on a 313-bp fragment of cytochrome oxidase 1 (*COI*) using primers m1COLintF: 5'-GGWACWGGWTGAACWGTWTAYCCYCC-3' [45] and modified jgHCO2198: 5'-TANACYTCNGGRTGNCCRAARAAYCA-3' [34]. For a small subset of samples, *COI* was amplified using the

primer pair jgHCO2198 and LCO1490 [32, 34]. Amplifications were conducted with tagged primers following Wang et al. [75] for Illumina and Srivathsan et al. [68, 69] for MinION. PCR reactions contained 4 μl Mastermix from CWBio, 1 μl of 1 mg/ml BSA, 1 μl of 10 μM of each primer, and 1 μl of DNA. PCR conditions were a 5 min initial denaturation at 94°C followed by 35 cycles of denaturation/annealing/extension (94°C (1 min)/ 47°C (2 min)/ 72°C (1 min)), and a final extension at 72°C (5 min). A subset of wells ($N=8-12$ with negative control) from each PCR plate was run on an agarose gel to check for plate-wide failure, before products were pooled and then purified using AMPure XP beads (Beckman Coulter Life Sciences, IN, USA). DNA concentration was quantified with a QubitTM dsDNA HS Assay Kit (Invitrogen, CA, USA).

Illumina and MinION sequencing were used to sequence the amplicons for this study. Illumina libraries were prepared using a TruSeq DNA PCR-free kit to obtain 250-bp PE sequences using Illumina HiSeq 2500; the sequencing was outsourced. Nanopore sequencing using MinION was conducted in house following Srivathsan et al. [68]. Library preparation was conducted using either SQK-LSK109 or SQK-LSK110 ligation sequencing kits (Oxford Nanopore Technologies, Oxford, UK) with 200-ng pooled and purified PCR products. The manufacturer's instructions were followed except for use of 1 \times AMPure beads instead of 0.4 \times as suggested in the instructions because the minibarcode amplicons in our experiments were short (~391 bp with primers and tags), and a modified protocol for end-repair as described in Srivathsan et al. [69]. The sequencing was carried out using a MinION sequencer with either R9.4.1 or R10.3 flow cells for a maximum of 72 h. Basecalling was conducted using Guppy versions 2.3.5+53a111f and 4.2.3+f90bd04. FastQ files were demultiplexed to obtain individual *COI* sequences identifiable back to individual specimens with unique specimen codes. This was done following Srivathsan et al. [67–69] for MinION and Wang et al. [75] for Illumina.

To exclude contaminants and potential mis-sorts, *COI* sequences were matched using BLAST to NCBI GenBank's nucleotide database and sequences with >97% similarity to non-scuttle fly taxa were removed. Sequences were then aligned using MAFFT v7 [43] and clustered into molecular operational taxonomic units (mOTUs) using objective clustering (part of "TaxonDNA", see [52]), a distance-based method that groups sequences based on a user-defined threshold for minimum interspecific uncorrected p-distance.

A species-proxy clustering threshold of 1.7% was used for the primary analyses as this was the distance that maximised cluster congruence to species-level

morphology for a subset of 18,000 scuttle flies from this dataset [40]. References to “species” are to mOTUs at this threshold. To assess the impact of the selected threshold on distribution patterns, we examined other thresholds ranging from 0 to 5% and ran models using both species and haplotype data.

Diversity

To characterise local and national species richness, we used accumulation curves and Chao1 richness, as implemented in EstimateS [18] and R (R Development Core Team) package iNEXT [15] for total diversity, diversity per zone, and diversity per site. Following the approach of Ronquist et al. (2020), a combined non-parametric estimator (CNE) was used to obtain an alternative estimate of total species richness.

Characterisation of community variation

To visualise scuttle fly community composition in the context of geographic and climatic variation, we used non-metric multi-dimensional scaling (NMDS) plotted with R packages ggplot2 [78] and Vegan v2.5–4 [55] using both abundance-based (Bray–Curtis) and incidence-based (Jaccard) dissimilarity indices. Since both metrics yielded nearly identical results, the results in the paper are based on Bray–Curtis indices. Sites were classified according to the nine plant hardiness zones of the Swedish Horticultural Society that synthesise climatic variables of particular importance to horticultural plants into horticultural zones: strong and rapid temperature changes, very low temperatures sustained for long periods, evaporation occurring when the sun is shining but the ground is still frozen, and temporal considerations of climatic and environmental conditions (Fig. 1, [62], used with permission). As these zones have proven useful in identifying the survival and growth of horticultural plants across Sweden, we hypothesised that they may offer a relevant description of the environment for any organism sensitive to similar climatic variables—including scuttle flies. They also offer a more regional approach than the broad biogeographic classification previously used for spatial analysis of scuttle flies in Sweden [63]. We therefore tested whether the zones could predict scuttle fly richness and distributions [20, 44, 50].

To quantify differences in scuttle fly community composition between plant hardiness zones and to assess the significance of those differences, we used analysis of similarities (ANOSIM) and similarity percentage analysis (SIMPER). These tests were run with PRIMER v7 [17].

We excluded samples that contained fewer than 100 specimens from the ordination analyses, as small sample sizes can artificially generate large distances in the NMDS

plots, thus obscuring structured variation in community composition arising from responses to the environment.

Hierarchical Modelling of Species Communities

To relate variation in community composition to continuous environmental drivers without the assumptions of zones used in our NMDS visualisations, we used Hierarchical Modelling of Species Communities (HMSC [58, 59], a type of joint species distribution model [76]). Due to the zero-inflated nature of the data, we fitted hurdle models, i.e. one model for the occurrence of taxa across samples (probit regression), and a second model for their abundance conditional on presence (linear regression for log-transformed count data, with zeros masked as missing data), henceforth referred to simply as an “abundance model”.

To evaluate the impact of the criterion used in species delimitation, we used four different response variables: species presence-absence, species abundance (where present), haplotype occurrence, and haplotype abundance (where present), thus resulting in four HMSC models being fit. As these analyses are uninformative for taxa with very sparse data, we included only haplotypes or species with at least five occurrences ($n = 391$ haplotypes and 273 species, out of the 2697 haplotypes and 549 species observed; see “Results”). To be able to relate each sample to specific climatic conditions, we also excluded 15 samples for which the trapping duration exceeded approximately a month (> 34 days).

We included five predictor variables, coded as fixed effects, in our HMSC models. To test for seasonal changes in scuttle fly communities among the sampling time periods, we included the four-category variable “Time.Period”. To assess the effect of the spatial gradient in mean climatic conditions across Sweden on scuttle fly communities, we included the continuous variable mean annual temperature (“bio1” from the Worldclim database, [31]). The importance of tree cover for scuttle fly distributions was assessed by including a continuous variable describing the percentage of forest or woodland cover within a 50-m buffer (“ForestWood”) derived from the Swedish National Land Cover Database (<https://www.naturvardsverket.se/en/services-and-permits/maps-and-map-services/national-land-cover-database/>). Site values for “bio1” and “ForestWood” were extracted using the R package raster [41]. To account for the effect of sampling effort on taxon detection, we further included a continuous variable quantifying the total number of trapping days per sample (“TrapDays”) and a binary variable indicating whether all specimens in a sample were sequenced or not (“FullSample”). The latter defines whether the sample was fully sequenced or not, and was included to account for the higher likelihood of encountering taxa in samples that

had very high individual abundances (i.e. samples that could not be fully sequenced). We moreover included a spatially explicit random effect based on sample site coordinates (“Random: site”) and a temporally explicit random effect based on median sampling date (“Random: time”) to model any spatial or temporal autocorrelation in the scuttle fly dataset, and a categorical random effect representing sample identity (“Random: sample”).

We fitted the model with the R package HMSC [71] assuming the default prior distributions. We sampled the posterior distribution of four MCMC chains, each of which was run for 37,500 iterations, of which the first 12,500 were removed as burn-in. The iterations were thinned by 100 to yield 250 posterior samples per chain, and thus 1000 posterior samples in total. To explore the rate of Markov chain Monte Carlo (MCMC) convergence, we also fitted otherwise identical models but with 375 iterations (burn-in 125, thin 1) and 3750 iterations (burn-in 1250, thin 10). Convergence was assessed by examining the potential scale reduction factor (PSRF) distribution over the fixed effect (β) and random effect (Ω) parameters, equivalent to the Gelman–Rubin statistic [35].

We examined the explanatory power of the models through species-specific coefficients of discrimination ($T_{\text{jur}} R^2$) for the occurrence models, which measure how well the model discriminates those samples in which a taxon occurs from those in which it does not occur, and R^2 for the abundance models. $T_{\text{jur}} R^2$ is defined as the difference in average model-predicted probability of occurrence for samples in which the species is present vs absent [72].

To assess whether the distributions of rare taxa are likely to be driven by similar factors to the common ones, we calculated Pearson and Mantel correlations, respectively, between the observed species richness and community composition (Bray–Curtis dissimilarity) of the modelled more common taxa and the excluded rare taxa (those with < 5 occurrences).

Community dissimilarity

To further partition variation in scuttle fly community composition in time (between seasons) and space (across Sweden), we dissected overall community dissimilarity (β -diversity) into its turnover (i.e. species replacement) and nestedness (differences of species richness between sites) components [3].

Our analyses encompassed three aspects: (1) spatial differences: differences in community composition between each sampling site considering their geographical distance, (2) temporal differences: differences in community composition between each sampling period, accounting for their temporal distance in weeks, and (3) spatiotemporal differences: community differences between each

sampling site and sampling period including their joint effect. We computed pairwise Jaccard community dissimilarity values using both the species and haplotypes datasets for comparison. Given the observed geographical distance between sampling sites, these distances were rescaled to a mean of zero and a standard deviation of one before analyses. To characterise the temporal difference between samples in each sample pair, we used the difference in the mean week of sampling (for sample-specific details, see Additional file 1: Table S1). Values of total β -diversity, turnover, and nestedness were calculated for each pairwise comparison of sampling sites, sampling period, and sample pairs. We excluded self-pairs and included data for each pair only once. To test the correlation between each beta diversity component and the distances in space and time, we used Mantel tests based on Pearson moment correlations. All calculations were implemented in R package “betapart” [5].

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12915-024-02010-z>.

Supplementary Material 1: Additional file 1: Figures S1–S13, Tables S1–S6. Fig. S1 Map of study sites. Fig. S2 Timeline of analysed samples. Fig. S3 Chao1 estimates of scuttle fly species richness per zone. Fig. S4 Species accumulation curve of scuttle flies by season. Fig. S5 Consistency in community patterning, as based on different clustering thresholds for species delimitation. Fig. S6 NMDS plot of all samples (species, threshold of 100 specimens) with samples colour-coded according to time period. Fig. S7 Violin plots showing potential scale reduction factors (PSRF) for the beta and omega parameters of four HMSC models: species occurrence (1), species abundance (2), haplotype occurrence (3), and haplotype abundance (4). Fig. S8 Variance explained by fixed and random effects in HMSC models of the presence/absence or abundance of either haplotypes or species proxies. Fig. S9 Predicted mean prevalence and log(abundance) of species proxies as a function of forest or woodland cover (%) and mean annual temperature (°C) during late summer in HMSC models. Fig. S10 Pairwise residual associations among haplotypes and species proxies in space and time as detected in HMSC models. Fig. S11 Pairwise compositional turnover as a function of intervening geographic distance for each time period in the species and haplotype datasets. Fig. S12 Nestedness values scored along geographic distance gradients for each time period among the species and haplotype datasets. Fig. S13 Total beta-diversity values scored along geographic distance gradients for each time period among the species and haplotype datasets. Table S1 Sample information. Table S2 Distinctness of scuttle fly communities across zones, as based on different clustering thresholds for species delimitation. Table S3 Similarity of scuttle fly communities across zones, as based on different clustering thresholds for species delimitation. Table S4 Distinctness of scuttle fly communities across time periods. Table S5 Similarity of scuttle fly communities across time periods. Table S6 Mantel tests relating pairwise community dissimilarity to pairwise differences in space, time, or both.

Acknowledgements

We thank the Scuttle Fly Sorting Party crew at Station Linné – Carina Romero Ugarph, Harald Havnäs, Johan Ennerfelt, Marianne “Mia” Blomqvist, Nino Pettersson, Robert Ennerfelt, and especially Dave Karlsson. We thank the members of the Evolutionary Biology Lab at the National University of Singapore for help with the many hours of wetlab work. We thank Darren Yeo for assisting us with visualisations. We thank Inger Ekrem at Riksförbundet Svensk Trädgård for helping with the plant hardiness zone map, Tomas Lagerström for further information on the map, and Eva Ronquist for first bringing our attention to

the map. We thank the Swedish Taxonomy Initiative for the support to investigate the scuttle flies of Sweden, this study brings us one step closer to the ultimate goal of describing all of these species.

Authors' contributions

EH and RM designed the study. EH and LL generated the data. EH, LL, AM, MJ, and PPA analysed the data. EH, LL, and RM wrote the initial draft of the manuscript. All authors contributed to manuscript revision, read and approved the submitted version.

Funding

Open access funding provided by NTNU Norwegian University of Science and Technology (incl St. Olavs Hospital - Trondheim University Hospital) EH was funded by Swedish Taxonomy Initiative grant 2016–203 4.3. TR and OO were funded by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (ERC-synergy grant 856506—LIFEPLAN). OO was funded by Academy of Finland (grant no. 336212 and 345110). MJ was supported by the Academy of Finland's "Thriving Nature" research profiling action.

Availability of data and materials

Data and scripts are available on the project GitHub page at <https://github.com/leshonlee/DocumentingPhorids>. General HMSC pipeline scripts are available at <https://www2.helsinki.fi/en/researchgroups/statistical-ecology/hmsc>.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare no competing interests.

Author details

¹Department of Natural History, NTNU University Museum, Norwegian University of Science and Technology, Trondheim NO-7491, Norway. ²Zoology Department, Stockholm University, Stockholm 106 91, Sweden. ³Department of Biological Sciences, National University of Singapore, Science Drive 4, Singapore 117558, Singapore. ⁴National University of Singapore, Lee Kong Chian Natural History Museum, 2 Conservatory Dr, Singapore 117377, Singapore. ⁵Center for Integrative Biodiversity Discovery, Leibniz Institute for Evolution and Biodiversity Science, Museum Für Naturkunde, Invalidenstraße 43, Berlin 10115, Germany. ⁶Faculty of Biological and Environmental Sciences, University of Helsinki, P.O. Box 65, Helsinki 00014, Finland. ⁷Institute of Biotechnology, HILIFE Helsinki Institute of Life Science, University of Helsinki, P.O. Box 65, Helsinki 00014, Finland. ⁸Department of Ecology, Swedish University of Agricultural Sciences (SLU), Ulls Väg 18B, Uppsala 75651, Sweden. ⁹Department of Biological and Environmental Science, University of Jyväskylä, P.O. Box 35, Jyväskylä 40014, Finland. ¹⁰Institute for Biology, Humboldt University of Berlin, Unter Den Linden 6, Berlin 10117, Germany.

Received: 28 May 2024 Accepted: 5 September 2024

Published: 27 September 2024

References

- Bakke SJ, Ionita M, Tallaksen LM. The 2018 northern European hydrological drought and its drivers in a historical perspective. *Hydrol Earth Syst Sci*. 2020;24(11):5621–53. <https://doi.org/10.5194/hess-24-5621-2020>.
- Bar-On YM, Phillips R, Milo R. The biomass distribution on Earth. *Proc Natl Acad Sci*. 2018;115(25):6506–11. <https://doi.org/10.1073/pnas.1711842115>.
- Baselga A. Partitioning the turnover and nestedness components of beta diversity. *Glob Ecol Biogeogr*. 2010;19(1):134–43.
- Baselga A, Gómez-Rodríguez C, Araújo MB, Castro-Insua A, Arenas M, Posada D, et al. Joint analysis of species and genetic variation to quantify the role of dispersal and environmental constraints in community turnover. *Ecography*. 2022;2022(5):e05808. <https://doi.org/10.1111/ecog.05808>.
- Baselga A, Orme CDL. betapart: an R package for the study of beta diversity. *Methods Ecol Evol*. 2012;3(5):808–12.
- Bastos A, Ciais P, Friedlingstein P, Sitch S, Pongratz J, Fan L, et al. Direct and seasonal legacy effects of the 2018 heat wave and drought on European ecosystem productivity. *Sci Adv*. 2020;6(24):eaba2724. <https://doi.org/10.1126/sciadv.aba2724>.
- Bell JR, Botham MS, Henrys PA, Leech DI, Pearce-Higgins JW, Shortall CR, et al. Spatial and habitat variation in aphid, butterfly, moth and bird phenologies over the last half century. *Glob Change Biol*. 2019;25:1982–94.
- Bogoni JA, Peres CA, Ferraz KMPMB. Extent, intensity and drivers of mammal defaunation: a continental-scale analysis across the Neotropics. *Sci Rep*. 2020;10:14750.
- Borkent A, Brown BV. How to inventory tropical flies (Diptera)-one of the megadiverse orders of insects. *Zootaxa*. 2015;3949:301–22.
- Borkent ART, Brown BV, Adler PH, Amorim DDS, Barber K, Bickel D, et al. Remarkable fly (Diptera) diversity in a patch of Costa Rican cloud forest: why inventory is a vital science. *Zootaxa*. 2018;4402:53.
- Brown BV, Hartop EA. Big data from tiny flies: patterns revealed from over 42,000 phorid flies (Insecta: Diptera: Phoridae) collected over one year in Los Angeles, California, USA. *Urban Ecosyst*. 2017;20:521–34.
- Brown BV. Phorid newsletter 5 [Internet]. 1996. Available from: <https://phorid.net/newsletters/pnews5.pdf>.
- Callaghan CT, Borda-de-Água L, van Klink R, Rozzi R, Pereira HM. Unveiling global species abundance distributions. *Nat Ecol Evol*. 2023;7:1600–9. <https://doi.org/10.1038/s41559-023-02173-y>.
- Cardinale BJ, Duffy JE, Gonzalez A, Hooper DU, Perrings C, Venail P, et al. Biodiversity loss and its impact on humanity. *Nature*. 2012;486:59–67.
- Chao A, Ellison AM, Colwell RK, Gotelli NJ, Sander EL, Hsieh TC, et al. Rarefaction and extrapolation with Hill numbers: a framework for sampling and estimation in species diversity studies. *Ecol Monogr*. 2014;84:45–67.
- Chua PYS, Bourlat SJ, Ferguson C, Korlevic P, Zhao L, Ekrem T, et al. Future of DNA-based insect monitoring. *Trends Genet*. 2023;39(7):531–44. <https://doi.org/10.1016/j.tig.2023.02.012>. Epub 2023 Mar 10 PMID: 36907721.
- Clarke K, Gorley RN. PRIMER v6: user manual/tutorial. PRIMER-E, Plymouth. 2006;29:1060–5.
- Colwell RK. Estimates: Statistical Estimation of Species Richness and Shared Species from Samples. Version 9. User's Guide and Application [Internet]. 2013. Available from: <http://purl.oclc.org/estimates>.
- Díaz S, Fargione J, Chapin FS, Tilman D. Biodiversity loss threatens human well-being. *PLoS Biol*. 2006;4:1300–5.
- Disney RHL. Scuttle flies: the Phoridae. London: Chapman and Hall; 1994.
- Disney RHL, Durska E. Conservation evaluation and the choice of faunal taxa to sample. *Biodivers Conserv*. 2008;17:449–51.
- Disney RHL, Erzincinoglu YZ, de C Henshaw DJ, Unwin DM, Withers P, Woods A. Collecting methods and the adequacy of attempted fauna surveys, with reference to the Diptera. *Field Stud*. 1982;5:607–621.
- Durska E. The species composition and structure of scuttle fly communities (Diptera: Phoridae) in mature tree stands in pine forests at different stages of habitat degradation. *Fragm Faun*. 1996;39:267–85.
- Durska E. Secondary succession of scuttle fly communities (Diptera: Phoridae) in moist pine forest in Białowieża Forest. *Fragm Faun*. 2001;44:79–128.
- Durska E. The phenology of dominant scuttle-fly (Diptera: Phoridae) species in the Białowieża Forest. *Entomol Fenn*. 2002;13:123–7.
- Durska E. The phenology of *Triphleba Rondani* species (Diptera: Phoridae) in moist pine forests in the Białowieża Forest. *Entomol Fenn*. 2003;14:177–82.
- Durska E. Diversity of scuttle fly (Diptera: Phoridae) communities in the plantations of moist pine forests of the Białowieża Primeval Forest and the Tuchola Forest (Poland). *Biodivers Conserv*. 2006;15:385–93.
- Durska E. Effects of disturbances on scuttle flies (Diptera: Phoridae) in pine forests. *Biodivers Conserv*. 2013;22:1991–2021.
- Durska E. Preliminary data of the scuttle flies (Diptera: Phoridae) in the linden-oak-hornbeam forest of the Wigry National Park. *North East Poland Fragm Faun*. 2020;63:89–98.

30. Durska E, Kaczorowska E, Disney RHL. Scuttle flies (Diptera: Phoridae) of saline habitats of the Gulf of Gdansk. *Poland Entomol Fenn.* 2005;16:159–64.
31. Fick SE, Hijmans RJ. Worldclim 2: new 1-km spatial resolution climate surfaces for global land areas. *Int J Climatol.* 2017;37:4302–15.
32. Folmer O, Black M, Hoeh W, Lutz R, Vrijenhoek R. DNA primers for amplification of mitochondrial cytochrome c oxidase subunit I from diverse metazoan invertebrates. *Mol Mar Biol Biotechnol.* 1994;3:294–9.
33. Garda AA, Stein MG, Machado RB, Lion MB, Juncá FA, Napoli MF. Ecology, biogeography, and conservation of amphibians of the Caatinga. In: Silva JMC, Leal IR, Tabarelli M, editors. *Caatinga*. Cham: Springer; 2017. p. 133–149. Available from: https://doi.org/10.1007/978-3-319-68339-3_5.
34. Geller J, Meyer C, Parker M, Hawk H. Redesign of PCR primers for mitochondrial cytochrome c oxidase subunit I for marine invertebrates and application in all-taxa biotic surveys. *Mol Ecol Resour.* 2013;13:851–61.
35. Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. *Bayesian data analysis*. 3rd ed. London: CRC Press; 2014.
36. Gilliland HC. Hundreds of new and unusual insects discovered in the Amazon's canopy. Available from: <https://www.nationalgeographic.com/magazine/article/hundreds-of-new-and-unusual-insects-discovered-in-the-amazon-canopy-feature>. 2021.
37. Goodsell R, Tack A, Ronquist F, van Dijk L, Iwazkiewicz-Eggebrecht E, Miraldo A, et al. The rarity of Invertebrates prevents reliable application of IUCN Red-List criteria. *EcoEvoRxiv.* 2024. Available from: <https://doi.org/10.32942/X23G71>.
38. Gullefors B. Limes norrlandicus - a natural biogeographical border for caddisflies (Trichoptera) in Sweden. *Ferrantia.* 2008;55:61–5.
39. Hari V, Rakovec O, Markonis Y, Hanel M, Kumar R. Increased future occurrences of the exceptional 2018–2019 Central European drought under global warming. *Sci Rep.* 2020;10:12207. <https://doi.org/10.1038/s41598-020-68872-9>.
40. Hartop E, Srivathsan A, Ronquist F, Meier R. Towards large-scale integrative taxonomy (LIT): resolving the data conundrum for dark taxa. *Syst Biol.* 2022. <https://doi.org/10.1093/sysbio/syac033>.
41. Hijmans R. raster: geographic data analysis and modeling. R package version 3.6–27, Available from: <https://rspatial.org/raster>. 2024.
42. Karlsson D, Hartop E, Forshage M, Jaschhof M, Ronquist F. The Swedish Malaise Trap Project: A 15 Year Retrospective on a Countrywide Insect Inventory. *Biodivers Data J.* 2020;8:e47255. <https://doi.org/10.3897/BDJ.8.e47255>.
43. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol.* 2013;30:772–80.
44. Kitching RL, Bickel DJ, Boulter S. The evolutionary biology of flies. In: Guild analyses of dipteran assemblages: a rationale and investigation of seasonality and stratification in selected rainforest faunas. New York: Columbia University Press; 2005. p. 388–415.
45. Leray M, Yang JY, Meyer CP, Mills SC, Agudelo N, Ranwez V, et al. A new versatile primer set targeting a short fragment of the mitochondrial COI region for metabarcoding metazoan diversity: application for characterizing coral reef fish gut contents. *Front Zool.* 2013;10:34.
46. Linnaeus C. *Systema naturæ, sive regna tria naturæ systematice proposita per classes, ordines, genera, & species*. 10th ed. Leiden: Lugduni Batavorum; 1758.
47. Lu Y, Yang Y, Sun B, Yuan J, Yu M, Stenseth NC, Bullock JM, Obersteiner M. Spatial variation in biodiversity loss across China under multiple environmental stressors. *Sci Adv.* 2020;6(47). <https://doi.org/10.1126/sciadv.abd0952>.
48. Lundström J, Schäfer M, Hesson J, Blomgren E, Lindström A, Wahlqvist P, Halling A, Hagelin A, Ahlm C, Evander M, Broman T, Forsman M, Persson Vinnersten T. The geographic distribution of mosquito species in Sweden. *J Eur Mosq Control Assoc.* 2013;31:35.
49. Mace GM, Barrett M, Burgess ND, Cornell SE, Freeman R, Grooten M, et al. Aiming higher to bend the curve of biodiversity loss. *Nat Sustain.* 2018;1:448–51.
50. McGlynn TP, Meineke EK, Bahlai CA, Li E, Hartop EA, Adams BJ, et al. Temperature accounts for the biodiversity of a hyperdiverse group of insects in urban Los Angeles. *Proc R Soc B.* 2019;286:20191020.
51. Meier R, Hartop E, Pylatiuk C, Srivathsan A. Towards holistic insect monitoring: species discovery, description, identification, and traits for all insects. *Phil Trans R Soc B.* 2024;379:20230120. <https://doi.org/10.1098/rstb.2023-0120>.
52. Meier R, Shiyang K, Vaidya G, Ng PK. DNA barcoding and taxonomy in Diptera: a tale of high intraspecific variability and low identification success. *Syst Biol.* 2006;55:715–28.
53. Miller G. Linnaeus's legacy carries on. *Science.* 2005;307:1038–9.
54. Naeem S, Chazdon R, Duffy JE, Prager C, Worm B. Biodiversity and human well-being: an essential link for sustainable development. *Proc R Soc B Biol Sci.* 2016;283:20162091.
55. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, et al. *vegan*. Community ecology package. 2019;2:5–6.
56. Ollerton J. Pollinator diversity: distribution, ecological function, and conservation. *Annu Rev Ecol Evol Syst.* 2017;48:353–76.
57. Orr MC, Hughes AC, Chesters D, Pickering J, Zhu CD, Ascher JS. Global patterns and drivers of bee distribution. *Curr Biol.* 2020;30:843–8.
58. Ovaskainen O, Abrego N. *Joint species distribution modelling*. Cambridge: Cambridge University Press; 2020. p. 372. (Ecology, Biodiversity and Conservation). ISBN: 978-1-108-71678-9. eISBN: 9781108591720.
59. Ovaskainen O, Tikhonov G, Norberg A, Guillaume Blanchet F, Duan L, Dunson D, et al. How to make more out of community data? A conceptual framework and its implementation as models and software. *Ecol Lett.* 2017;20:561–76.
60. Rabone M, Wiethase JH, Simon-Lledó E, Emery AM, Jones DOB, Dahlgren TG, Bribeasa-Contreras G, Wiklund H, Horton T, Glover AG. How many metazoan species live in the world's largest mineral exploration region? *Curr Biol.* 2023;33(12):2383–2396.e5. <https://doi.org/10.1016/j.cub.2023.04.052>.
61. Ricklefs RE. Evolutionary diversification and the origin of the diversity–environment relationship. *Ecology.* 2006;87(S3–S13). [https://doi.org/10.1890/0012-9658\(2006\)87\[3:EDATOO\]2.0.CO;2](https://doi.org/10.1890/0012-9658(2006)87[3:EDATOO]2.0.CO;2).
62. Riksförbundet Svensk Trädgård. Zonkarta. 2018. Available from: http://www.tradgard.org/svensk_tradgard/zonkarta.html. Cited 2024 May 28.
63. Ronquist F, Forshage M, Häggqvist S, Karlsson D, Hövmler R, Bergsten J, et al. Completing Linnaeus's inventory of the Swedish insect fauna: only 5,000 species left? *PLOS One.* 2020;15:e0228561.
64. Skvarla M, Larson J, Fisher R, Dowling A. A review of terrestrial and canopy Malaise traps. *Ann Entomol Soc Am.* 2020;114. <https://doi.org/10.1093/aesa/saaa044>.
65. SLU Artdatabanken. Famij: Phoridae - puckelflugor. 2021. Available from: <https://www.dyntaxa.se/Taxon/Info/2001326>. Cited 2024 May 28.
66. Srivathsan A, Ang Y, Heraty JM, Hwang WS, Jusoh WFA, Kutty SN, et al. Convergence of dominance and neglect in flying insect diversity. *Nat Ecol Evol.* 2023;7(7):1012–21.
67. Srivathsan A, Balogh B, Wang W, et al. A MinION™-based pipeline for fast and cost-effective DNA barcoding. *Mol Ecol Resour.* 2018;18:1035–49. <https://doi.org/10.1111/1755-0998.12890>.
68. Srivathsan A, Hartop E, Puniamoorthy J, Lee WT, Kutty SN, Kurina O, et al. Rapid, large-scale species discovery in hyperdiverse taxa using 1D MinION sequencing. *BMC Biol.* 2019;17:96.
69. Srivathsan A, Hartop EA, Puniamoorthy J, Lee WT, Kutty SN, Kurina O, et al. MinION barcodes: biodiversity discovery and identification by everyone, for everyone. *BMC Biol.* 2021;19:217.
70. Svenning JC, Borchsenius F, Bjorholm S, Balslev H. High tropical net diversification drives the New World latitudinal gradient in palm (Arecaceae) species richness. *J Biogeogr.* 2008;35:394–406.
71. Tikhonov G, Opedal ØH, Abrego N, Lehtikoinen A, Jonge MMJ, Oksanen J, et al. Joint species distribution modelling with the R-package HMSC. *Methods Ecol Evol.* 2020;11:442–7.
72. Tjur T. Coefficients of determination in logistic regression models - a new proposal: the coefficient of discrimination. *Am Stat.* 2009;63:366–72.
73. Townes H. A light-weight Malaise trap. *Entomol News.* 1972;83:239–47.
74. Truett GE, Heeger P, Mynatt RL, Truett AA, Walker JA, Warman ML. Preparation of PCR-quality mouse genomic DNA with hot sodium hydroxide and tris (HotSHOT). *Biotechniques.* 2000;29:52–4.
75. Wang WY, Srivathsan A, Foo M, Yamane S, Meier R. Sorting specimen-rich invertebrate samples with cost-effective NGS barcodes: validating a reverse workflow for specimen processing. *Mol Ecol Resour.* 2018;18:490.
76. Warton DI, Blanchet FG, O'Hara RB, Ovaskainen O, Taskinen S, Walker SC, et al. So many variables: joint modeling in community ecology. *Trends Ecol Evol.* 2015;30:766–79.
77. Wellenreuther M, Larson KW, Svensson EI. Climatic niche divergence or conservatism? Environmental niches and range limits in ecologically similar damselflies. *Ecology.* 2012;93:1353–66.

78. Wickham H. ggplot2: Elegant Graphics for Data Analysis [Internet]. New York: Springer-Verlag; 2016. Available from: <https://ggplot2.tidyverse.org>.
79. Wiens JJ. Global patterns of diversification and species richness in amphibians. *Am Nat*. 2007;170:S86–106.
80. Wolda H. Insect seasonality: why? *Annu Rev Ecol Syst*. 1988;19(1):1–18.
81. World Economic Forum. The Global Risks Report 2020 [Internet]. 15th ed. Geneva: World Economic Forum; 2020. [cited 2024 Sep 17]. Available from: <https://www.weforum.org/reports/the-global-risks-report-2020>.
82. Yeo D, Srivathsan A, Meier R. Longer is not always better: optimizing barcode length for large-scale species discovery and identification. *Syst Biol*. 2020;0:1–16.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.