
BIOINTERACT: A Large-Scale Multimodal Dataset for Evaluating Fine-Grained Semantic Understanding of Biotic Interactions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 While recent advances in vision–language models (VLMs) have spurred the devel-
2 opment of domain-specific datasets and benchmarks, these often fail to assess fine-
3 grained semantic understanding, allowing models to achieve high scores without
4 robust visual grounding. We address this evaluation gap through the lens of biotic
5 interactions: directional, asymmetric relationships between organisms (e.g., wasp
6 parasitizes caterpillar vs. caterpillar parasitizes wasp). This relational complexity
7 yields naturally adversarial instances that expose superficial reasoning in current
8 VLMs. To this end, we introduce BIOINTERACT, the largest multimodal dataset for
9 evaluating VLM robustness on real-world biodiversity challenges. Curated from
10 iNaturalist and validated against scientific literature, the dataset contains 15.4K
11 unique interactions spanning 6.5K taxa across 256K images. Each interaction is
12 structured as a source-relation-target triplet, enabling controlled semantic pertur-
13 bations. We further introduce BIOINTERACT100, an adversarial image retrieval
14 benchmark revealing that state-of-the-art VLMs suffer from severe consistency
15 gaps and are highly brittle to relation-direction reversals. BIOINTERACT provides
16 a faithful evaluation of multimodal AI, while encouraging the development of
17 robust systems for ecological research. Data and code are available at the project
18 website.

19 1 Introduction

20 Recent advances in artificial intelligence (AI) technologies have led to the widespread adoption of
21 foundation models (FMs) across diverse domains, enabling scalable and increasingly automated
22 analysis of complex data [34]. In particular, multimodal FMs trained on heterogeneous data (e.g., text,
23 images) exhibit emergent capabilities that extend beyond their original training objectives, opening
24 new opportunities for data-driven discovery. This progress has led to a growing class of domain-
25 specific datasets and benchmarks aimed at evaluating such capabilities in realistic settings, including
26 health, medicine, and chemistry [6, 35, 61, 65]. However, most domain-specific benchmarks are
27 treated as monolithic performance targets, their design often emphasizing answer correctness, while
28 fine-grained semantic understanding is typically assessed on datasets centered around common objects
29 and everyday scenes. This disconnect allows models to achieve high performance without robust
30 visual grounding, thereby providing an incomplete picture of multimodal competence. Addressing
31 this evaluation gap is essential for ensuring reliable deployments of such systems in specialized fields
32 [22].

33 In ecology, advances in vision and machine learning have accelerated biodiversity knowledge through
34 targeted visual tasks ranging from basic visual interpretation (e.g., taxonomic classification, attribute
35 prediction, trait inference) to higher-level ecological context understanding (e.g., disease detection)

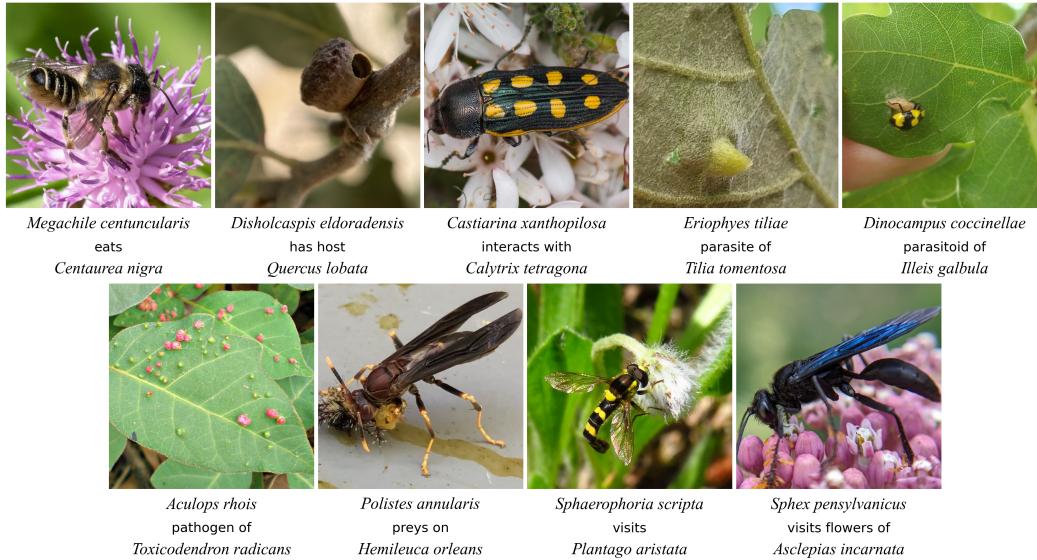


Figure 1: Example of images depicting biotic interactions in the BIOINTERACT dataset. Biotic interactions are annotated triplets—source taxon, interaction type, target taxon (e.g., *Megachile centuncularis* eats *Centaurea nigra*). Interactions are standardized to nine types defined by the Open Biological and Biomedical Ontologies (OBO) Relations Ontology (RO) [39] (see Table 2). Images originate from *research-graded* iNaturalist [28] observation records, where the source taxon occurrence is indexed (i.e., taxonomic identity, location, and time), while interaction type and target taxon are optionally annotated.

36 [31, 36, 50, 58]. However, understanding biodiversity and coping with the biodiversity crisis in human-
 37 modified ecosystems, requires knowledge about biotic interactions between multiple organisms.
 38 Biotic interactions are highly opportunistic in nature and, therefore, difficult to study at scale; direct
 39 observations are labor-intensive, and often restricted to specific habitats or timeframes, limiting
 40 their generalizability. Recent studies highlight the potential of community science data as a scalable
 41 resource for analyzing interactions. Platforms such as iNaturalist [28] and Observation.org [40] now
 42 host hundreds of millions of geo-referenced observations paired with images. Although these datasets
 43 are primarily used to document species occurrences, the images are captured *in situ* and frequently
 44 contain substantially richer ecological information as a by-product [44], including morphological
 45 traits, behavior, habitat, and interactions. For instance, pollinator images may reveal host flower
 46 preferences [45], while images of predatory species can provide evidence of trophic interactions
 47 [25]. Despite the rapid growth of such repositories, extracting richer ecological descriptors at scale
 48 remains a significant challenge. Fewer than 25% of over 300 million observations on iNaturalist
 49 include detailed annotations [55]. Although indispensable, manual labeling is both time-consuming
 50 and dependent on expert knowledge, rendering it impractical at this scale and limiting the ability of
 51 researchers to fully exploit these data. As ecosystem degradation accelerates [1], there is an urgent
 52 need for methods capable of leveraging these vast resources to automatically extract latent ecological
 53 information from existing observations.

54 Vision-language models (VLMs) offer a promising approach for zero-shot multimodal understanding
 55 of ecologically relevant images through tasks such as visual question answering (VQA) and image
 56 retrieval [15, 23, 37, 54]. By design, these tasks rely on unconstrained, free-form natural language
 57 inputs to enable flexible, user-friendly image queries. However, biotic interactions are particularly
 58 challenging, as their meaning is inherently relational and directional, encoding asymmetric effects
 59 between organisms where small semantic changes can completely alter the meaning, yielding naturally
 60 adversarial interaction instances (e.g., *wasp parasitizes caterpillar* and *caterpillar parasitized by wasp*
 61 are semantically equivalent, whereas *wasp parasitizes caterpillar* and *caterpillar parasitizes wasp* are
 62 syntactically equivalent, yet represent distinct interactions). Current benchmarks remain limited in
 63 scale and scope, typically evaluating answer accuracy using a single query per image in VQA such as
 64 CoralVQA [23] and AgMMU [15], or in text-to-image retrieval such as INQUIRE [54]. While these
 65 benchmarks assess whether a model produces the correct answer on fixed image-query pairs, they
 66 fail to capture the complexities of open-world settings, leaving a substantial gap in understanding

Table 1: Current benchmarks provide limited insight into a model’s ability to reason about biotic interactions under input perturbation—conditions that commonly arise in open-world settings. BIOINTERACT (ours) scales both in size and annotation richness (see Table 4, enabling adversarial and stress-based evaluation. Studies showcasing the potential of using images to study biotic interactions for advancing ecological research, outside of the vision and machine learning community, are listed in Appendix B.

| Tasks | Benchmark | # Images | Query | Description |
|--|---------------|----------|----------|---|
| <i>visual question answering (VQA)</i> | CoralVQA [23] | 12K | single | ecological and health-related conditions questions for coral reef analysis, in open-ended and closed-ended formats, including symbiotic relationship between corals and algae in underwater images |
| | AgMMU[15] | 50.2K | single | agriculture knowledge multiple-choice and open-ended questions, including insect/pest identification |
| <i>image retrieval</i> | INQUIRE [54] | 33K | single | 200 unique queries extracted from expert interviews and academic literature, including 12 parasitism and symbiosis relationships |
| <i>multi-task</i> | BIOINTERACT | 256K | multiple | 15.4K unique biotic interactions across five kingdoms (<i>Animalia</i> , <i>Plantae</i> , <i>Fungi</i> , <i>Chromista</i> , and <i>incertae sedis</i>) denoted with both scientific and vernacular names, standardized to nine interaction types, grounded in spatiotemporal metadata |

67 how effectively VLMs operate in realistic scenarios. A system capable of fine-grained semantic
68 understanding in an open-world setting would unlock meaningful progress in AI-driven ecological
69 research.

70 In this work, we challenge the biodiversity AI research by curating and releasing BIOINTERACT,
71 a large-scale multimodal dataset containing 15.4K unique biotic interactions spanning 6.5K taxa
72 across 256K images. Specifically, BIOINTERACT captures who interacts with whom, how, and
73 under what context, by grounding interactions in structured triplets—source taxa, interaction type,
74 target taxa. Examples from BIOINTERACT are shown in Figure 1, while Table 1 summarizes how
75 BIOINTERACT contrasts with existing datasets in scale and annotation richness. By focusing on
76 real-world biodiversity challenges, BIOINTERACT provides a faithful evaluation of multimodal AI
77 progress, encouraging the development of robust systems that can assist with accelerating ecological
78 research. Our main contributions include: (1) the largest multimodal dataset of biotic interactions to
79 date; (2) rich annotations enabling controlled natural language descriptions for fine-grained semantic
80 understanding; (3) BIOINTERACT100, a novel image retrieval task on a fixed set, comprising 281
81 unique interactions with 100 images; (4) controlled adversarial benchmark revealing that state-of-the-
82 art VLMs suffer from severe compositional shortcomings, failing drastically under relation-direction
83 reversals despite high baseline accuracy.

84 2 Related work

85 Vision-language models (VLMs) have emerged as a powerful paradigm for learning joint representa-
86 tions across visual and textual modalities. Early contrastive approaches such as CLIP [47] learn a
87 shared embedding space where data from two modalities can be encoded jointly. Building on this
88 foundation, VLMs, including GPT-4o [26], BLIP [26], LLaVA[33], and Qwen2-VL [56], extend
89 these capabilities by connecting the outputs of visual encoders directly into language models. Recent
90 expert-level benchmarks further challenge multimodal models in scenarios where expert-level knowl-
91 edge is required [6, 61, 65]. For instance, GMAI-MMBench [6] and BenchX [65] feature a collection
92 of healthcare and medically relevant questions, while MMMU-Pro [61], comprises visual-questions
93 related to problems from diverse fields such as business, arts, and science. However, the evaluation
94 landscape remains fragmented: domain-specific benchmarks are treated as monolithic performance
95 targets, while fine-grained semantic understanding is typically assessed only in general-domain
96 datasets depicting common objects and everyday scenes (e.g., MS-COCO) [2, 52, 59, 62]. Nearly all
97 domain-specific benchmarks are *static*, with performance gains increasingly reflecting memorization
98 rather than capability in pursuit of state-of-the-art performance [7]. It is crucial to develop evaluation
99 protocols that reflect true capabilities, otherwise we run the risk of deploying models that perform

Table 2: Interaction types in the BIOINTERACT dataset, defined by Open Biological and Biomedical Ontologies (OBO) Relations Ontology (RO) [39]. Interactions follow the hierarchy of biological process types. Each process type is mapped to a corresponding interaction relation, and finer-grained relations inherit from broader ones through a subrelation structure (see Appendix A).

| Type | Definition |
|-------------------|--|
| eats | A biotic interaction where one organism consumes a material entity through a type of mouth or other oral opening. |
| has host | X 'has host' y if and only if: x is an organism, y is an organism, and x can live on the surface of or within the body of y. |
| interacts with | An interaction relationship in which at least one of the partners is an organism and the other is either an organism or an abiotic entity with which the organism interacts. |
| parasite of | A parasite-host relationship where an organism benefits at the expense of another. |
| parasitoid of | A parasite that kills or sterilizes its host. |
| pathogen of | Inverse of has pathogen. |
| preys on | An interaction relationship involving a predation process, where the subject kills the target in order to eat it or to feed to siblings, offspring or group members. |
| visits | - |
| visits flowers of | - |

Table 3: Key statistics of the BIOINTERACT dataset. Biotic interactions are recorded between organisms at different taxonomic levels. The table includes the number of unique taxa (source and target taxa) recorded at the *species* and *genus* levels.

| Statistics | eats | has host | interacts with | parasite of | parasitoid of | pathogen of | preys on | visits | visits flowers of | All |
|------------|-------|----------|----------------|-------------|---------------|-------------|----------|--------|-------------------|--------|
| # Images | 27778 | 6926 | 116483 | 2558 | 13 | 30 | 2878 | 1636 | 98456 | 256758 |
| # Species | 2058 | 961 | 3924 | 219 | 4 | 4 | 370 | 198 | 1628 | 9366 |
| # Genus | 742 | 230 | 1452 | 36 | - | 2 | 109 | 42 | 523 | 3136 |

100 sub-optimally in real-life applications, potentially leading to incorrect decisions, user dissatisfaction,
 101 or even serious consequences in high-stakes applications such as healthcare [22].

102 BIOINTERACT is different from existing expert-level benchmarks for the following reasons. (1) *Focus*
 103 *on biotic interactions*: While prior work has addressed a range of visual tasks [31, 36, 37, 50, 58],
 104 these efforts largely focus on individual organisms or ecological context understanding . In contrast,
 105 BIOINTERACT targets biotic interactions, capturing relationships between organisms (e.g., predation,
 106 symbiosis, competition), which are fundamental to understanding biodiversity. The closest related
 107 efforts are found in VQA tasks, which evaluate multimodal understanding through both open-
 108 ended and closed-ended questions [15, 23] (e.g., *Is there a symbiotic relationship between the coral*
 109 *depicted in the center of the image and zooxanthellae?*), and in image retrieval, where models learn
 110 joint embeddings to align images with queries describing relationships between entities, leveraging
 111 multimodal large language models as out-of-the-box rerankers [54] (e.g., *Does this image show {some*
 112 *query}? Answer with "Yes" or "No" and nothing else.*). (2) *Systematic generation of semantically*
 113 *controlled natural language description*: Unlike domains such as mathematics or programming,
 114 where ground truth is unambiguous and logic is formal, biotic interactions are particularly challenging,
 115 as subtle semantic differences can fundamentally alter their meaning, yielding naturally adversarial
 116 interaction instances. To address this, we enable fine-grained semantic evaluation through the
 117 systematic generation of semantically equivariant statements from interaction triplets [10]. This
 118 approach provides precise control over semantically similar and dissimilar queries, without relying
 119 on human or prompt-based biases [9, 14, 32, 48].

120 3 BIOINTERACT

121 3.1 Data curation

122 We retrieve biotic interaction data from GloBI [10], using the archive hosted on Zenodo (version
 123 0.8) [11]. GloBI periodically collects biotic interactions from a wide range of independent datasets,
 124 such as iNaturalist [28] and Encyclopedia of Life [13], indexed in tabular form. Biotic interactions
 125 are semantic relationships between entities represented as knowledge graphs, stored in a triplet format
 126 consisting of a source taxon, an interaction type, and a target taxon. To construct our dataset, we focus
 127 on records that are linked to iNaturalist observations, which include images of organisms observed *in*
 128 *situ*. In these records, the source taxon corresponds to the organism documented in the iNaturalist

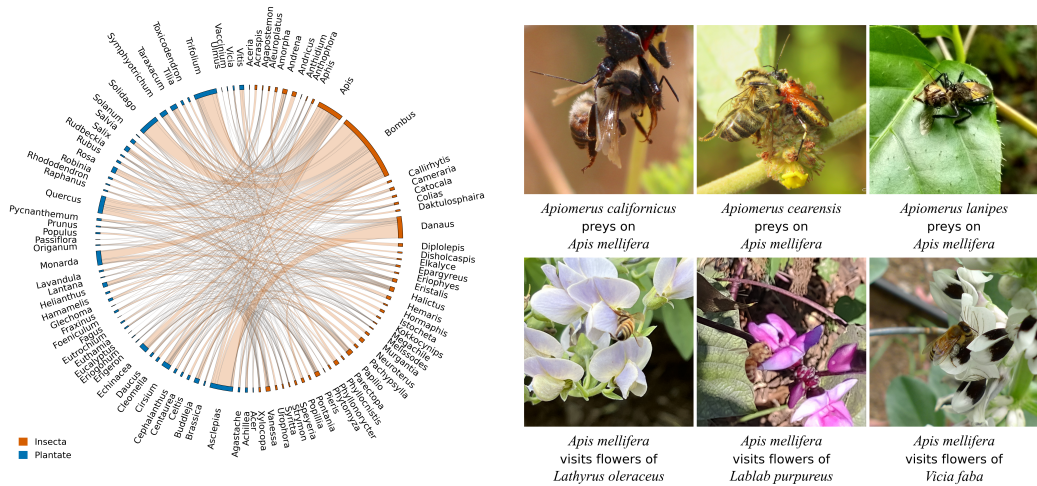


Figure 2: (left) BIOINTERACT top 100 most frequent taxa and their interactions. The dataset is centered around *Insecta-Plantae* interactions. (right) Examples illustrating ecological networks in BIOINTERACT. Top row showcases predatory interactions of the honeybee (*Apis mellifera*) with rove beetles (genus *Apiomerus*) [25]; while bottom row showcases its plant preferences, including associations with families such as *Fabaceae* [4].

129 observation, while additional annotations describe its ecological context: interaction type with a
 130 target taxon. We filter the GloBI data to retain only interaction records that reference a *research-grade*
 131 [27] iNaturalist observation, yielding 1.27M records. To these records we aggregate corresponding
 132 iNaturalist observation data from GBIF [17], including taxonomic identity, spatiotemporal information
 133 (location, and time) and corresponding images.

134 Interacting taxa are recorded at heterogeneous taxonomic levels. We align all records to a standardized
 135 seven-level Linnaean taxonomy (*kingdom, phylum, class, order, family, genus, species*) using the
 136 GBIF Backbone Taxonomy [16] (the most widely used, open-access, community-developed standard
 137 for sharing biodiversity data, providing a stable, standardized language for describing biological
 138 specimens and observations), and discard records whose taxonomic resolution falls outside this
 139 hierarchy. From the 1.27M biotic interaction records, we retains 636K records, corresponding to
 140 215K unique biotic interactions. We further enrich the dataset with corresponding English vernacular
 141 names from GBIF Backbone Taxonomy [16].

142 Biotic interaction with image data are documented *a priori* via iNaturalist; however, these annotations
 143 are not routinely verified as the *research-grade* data [27]. Consequently, we do not treat iNaturalist as
 144 a primary source of valid biotic interactions. Instead, we retain only interactions that are independently
 145 supported by at least one additional resource (a separate curated dataset, other than iNaturalist) indexed
 146 in the GloBI database [10]. This cross-validation procedure ensures high-confidence annotations and
 147 yields a final set of 15.4K unique biotic interactions represented by 256K image records.

148 3.2 Dataset overview

149 BIOINTERACT standardizes interactions to nine types defined by OBO RO [39] (see interaction
 150 type definitions in Table 2), while interacting taxa are reported at heterogeneous taxonomic levels
 151 across five kingdoms (*Animalia, Plantae, Fungi, Chromista, and incertae sedis*), with over 99% taxa
 152 recorded at the fine-grained *species* levels (see dataset key statistics in Table 3). Based on the most
 153 recent IUCN Red List assessment, BIOINTERACT covers 87 threatened species [30] and 23 of the
 154 world’s worst invasive species [18], enabling evaluation on ecologically critical and conservation-
 155 relevant taxa (see Appendix A) [20]. The spatiotemporal metadata associated with each image is
 156 particularly valuable for studying biotic interactions, as it supports region-specific and seasonal
 157 analyses that are fundamental to ecological research, especially in the study of invasive species [38].

158 Interactions are primarily recorded between insects and plants, indicating a strong skew toward
 159 plant-associated ecological relationships (see Figure 2). This focus is well justified, as recent studies

Table 4: BIOINTERACT annotations enumerate 27 fields provided in Parquet format; the dataset is openly available (for download and browsing) on [HuggingFace Datasets](#).

| Type | Description |
|----------------------------------|--|
| Source/target scientific name | Scientific taxon name represented as a string in fields <code>sourceTaxonName</code> and <code>targetTaxonName</code> . |
| Interaction type | Interactions standardized according to OBO RO [39]. |
| Source/target rank | The taxonomic rank at which the interacting taxa is identified (<i>kingdom, phylum, class, order, family, genus</i>). |
| Source/target vernacular name(s) | English common or vernacular name(s) represented as list of strings in fields <code>sourceVernacularName</code> and <code>targetVernacularName</code> . |
| Taxonomic hierarchy | Taxonomic hierarchy deterministically derived from taxa scientific name, represented as strings in separate fields for each source and target (<i>kingdom, phylum, class, order, family, genus</i>). |
| Spatiotemporal context | Date when the interaction was observed (separated DD, MM, YYYY), and latitude, longitude coordinates where the interaction was observed (decimals). |
| Licence | Image usage licence associated with each observation. |
| URL | Downloadable image link from Naturalist represented as a string in field <code>imageURL</code> . |

160 report strong declines worldwide in both insects and plants, threatening plant diversity, food security
 161 and ecosystem stability [42]. Interaction types are hierarchically organized, in which the general
 162 relation *interacts with* is recursively refined into more specific relations, such as *has host*, and further
 163 into fine-grained sub-relations like *parasite of*. Notably, this hierarchy is not strictly tree-structured:
 164 relations at the same level may overlap or exhibit cross-links, reflecting the compositional nature of
 165 biological interactions. For example, a bee can both visit a plant and visit flowers of it; parasitoid
 166 wasps can be both parasitoid of and functionally parasite of a host, illustrating cross-links between
 167 branches.

168 **License.** The dataset contains a heterogeneous mixture of Creative Commons licenses, predominantly
 169 non-commercial licenses (CC BY-NC 4.0), all of which are provided and should be considered
 170 accordingly in downstream usage.

171 **Geoprivacy.** We include geolocation metadata for all records, relying on the source platforms’
 172 automated and user-specified obscuration for sensitive species [29]. This means in practice that
 173 endangered or protected taxa have deliberately imprecise coordinates, in line with geoprivacy best
 174 practices.

175 **Responsible use.** Models trained on this data should not be used to infer, exploit, or manipulate
 176 sensitive biotic interactions in ways that could harm species, disrupt ecosystems, or facilitate unlawful
 177 activities such as wildlife persecution, invasive species spread, or targeted ecological interference;
 178 the dataset is provided to support ecological research, conservation, and responsible environmental
 179 monitoring.

180 3.3 Semantic perturbations

181 Biotic interactions are inherently directional, and can be expressed through multiple equivalent formu-
 182 lations. Let $\mathcal{D} = \{(x_i, s_i, r_i, t_i)\}_{i=1}^N$ denote a dataset of N samples, where x_i is an image, annotated
 183 with interaction triplets (s_i, r_i, t_i) , where s_i and t_i are the source and target taxa, respectively, and r_i
 184 is the interaction type. For each image x , BIOINTERACT enables controlled generation of meaning
 185 preserving queries, and contradictory queries that alter salient attributes such as taxa or interaction
 186 type leveraging rich image annotations (see Table 4). Controlled perturbation is preferred to ensure
 187 semantic equivalence and avoid hallucinations or unintended biases that may arise from generative
 188 approaches [57].

189 Given an image x and an interaction (s, r, t) , we construct queries through five perturbation types:
 190 (1) *entity removal*: the source or the target is removed; (2) *entity replacement*: the source or the
 191 target is randomly replaced by with another entity; (3) *relation removal*: the relation is removed,
 192 while preserving the source and the target; (4) *relation replacement*: the relation is replaced, while
 193 preserving the source or the target; (5) *relation inversion*: reverse the meaning of the relation, while

Table 5: Performance on fine-grained semantic understanding of biotic interaction queries for different models. We report the mean Average Precision at k (mAP@k), mean Reciprocal Rank (mMRR), and mean Recall at k (mRecall@10) computed over semantically equivalent query variants. **Bold and underlined** entries indicate the **best** and the **second best** results, respectively.

| Model | Method | Params (M) | mAP@50 | mMRR | mRecall@10 |
|------------------|----------------------|------------|--------------------|--------------------|--------------------|
| CLIP | ViT-L-14 | 427 | 0.14 ± 0.02 | 0.34 ± 0.06 | 0.56 ± 0.05 |
| MetaCLIP | ViT-L-14-quickgelu | 427 | 0.25 ± 0.06 | 0.49 ± 0.10 | 0.72 ± 0.11 |
| SigLIP | SO400M-14-SigLIP | 877 | 0.35 ± 0.10 | 0.60 ± 0.11 | 0.78 ± 0.08 |
| SigLIP2 | ViT-L-16-SigLIP2-256 | 881 | 0.37 ± 0.10 | 0.63 ± 0.10 | 0.83 ± 0.11 |
| BioTrove-BioCLIP | ViT-B-16 | 149 | 0.30 ± 0.08 | 0.57 ± 0.11 | 0.78 ± 0.09 |
| BioCAP | ViT-B-16 | 149 | 0.50 ± 0.06 | 0.78 ± 0.04 | 0.93 ± 0.03 |
| BioCLIP | ViT-B/16 | 149 | 0.30 ± 0.06 | 0.63 ± 0.07 | 0.82 ± 0.07 |
| BioCLIP2 | ViT-L/14 | 427 | 0.61 ± 0.08 | 0.83 ± 0.05 | 0.95 ± 0.02 |

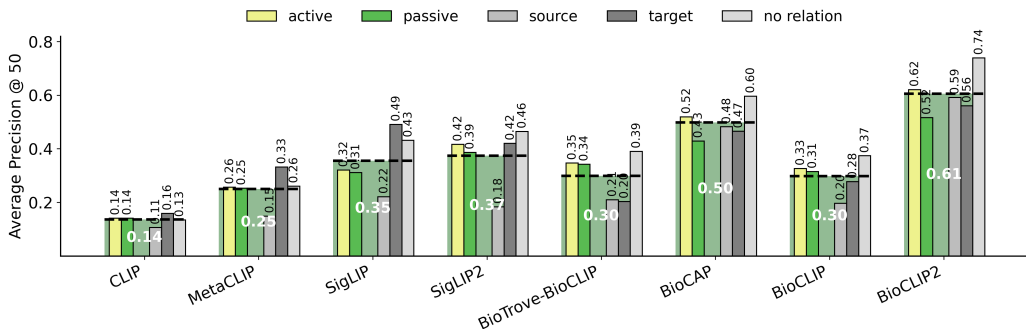


Figure 3: Performance per query reported under AP. Queries are based on interaction triples (s, r, t) : **active** queries specify the full triple (s, r, t) , **passive** queries invert the interaction (t, r', s) , **source** queries condition on (s, r) , **target** queries on (t, r') , and **no relation** queries specify only entity pairs without the relation (s, t) .

194 preserving the source and the target. To preserve semantic validity, we ensure that replacements
 195 of entities and relations remain ecologically plausible by leveraging hierarchical structures in both
 196 taxonomy and interaction types. Specifically, we generate alternative variants of s_i and t_i by sampling
 197 across different taxonomic levels, and substitute r_i with relations drawn from the same interaction-
 198 type hierarchy (see Appendix A). Conversely, to construct contradictory statements, we deliberately
 199 replace entities and relations with elements that fall outside these hierarchies, thereby violating
 200 ecological and semantic constraints.

201 Each perturbation can be augmented with multiple levels of semantic granularity by incorporating
 202 both scientific and vernacular names, as well as hierarchical variants of taxonomy and interactions
 203 types.

204 4 Evaluations

205 Our evaluation considers how predictions change under semantic perturbations, while the visual input
 206 remains constant. We introduce a ranking task similar to INQUIRE [54], which fixes images for each
 207 query and uses models like GPT-4o [26] to improve over initial text-to-image CLIP-style retrieval.
 208 We, thus, measure fine-grained semantic understanding in embedding similarity and binary (yes/no)
 209 questions.

210 **BioInteract100.** We construct a subset dataset via stratified sampling, selecting 100 image samples
 211 per interaction triplet group. Samples are prioritized based on annotation quality, favoring those
 212 with available vernacular names for both taxa and *species*-level annotations. Within each triplet
 213 group, we promote visual diversity by selecting at most one sample per image in the first pass,
 214 and subsequently filling remaining slots if necessary. This procedure yields BIOINTERACT100, a

Table 6: Performance on different general-purpose proprietary models grouped per query perturbation type. First we group semantically equivalent queries. Next we group contradictory queries which deliberately reverse entities and relations. We report the accuracy and flip rate for each group. Lower flip rates indicate greater consistency. **Bold** entries indicate the **best** results.

| Model | Equivalent | | Contradictory | | All | |
|----------------------------|--------------------|-------------|--------------------|-------------|--------------------|-------------|
| | Acc | Flip | Acc | Flip | Acc | Flip |
| <i>Text-only baselines</i> | | | | | | |
| GPT 5.2 mini | 0.97 ± 0.02 | 0.10 | 1.0 ± 0.0 | 0.0 | 0.98 ± 0.02 | 0.10 |
| Gemini 3 Flash Preview | 0.05 ± 0.03 | 0.23 | 0.60 ± 0.14 | 0.35 | 0.21 ± 0.26 | 0.81 |
| Claude Sonnet 4.6 | 0.90 ± 0.08 | 0.41 | 0.99 ± 0.0 | 0.71 | 0.93 ± 0.07 | 0.41 |
| <i>Image-text</i> | | | | | | |
| GPT 5.2 mini | 0.74 ± 0.12 | 0.54 | 0.52 ± 0.14 | 0.35 | 0.68 ± 0.17 | 0.95 |
| Gemini 3 Flash Preview | 0.77 ± 0.15 | 0.67 | 0.59 ± 0.20 | 0.44 | 0.71 ± 0.20 | 0.97 |
| Claude Sonnet 4.6 | 0.81 ± 0.03 | 0.44 | 0.65 ± 0.26 | 0.60 | 0.77 ± 0.09 | 0.71 |

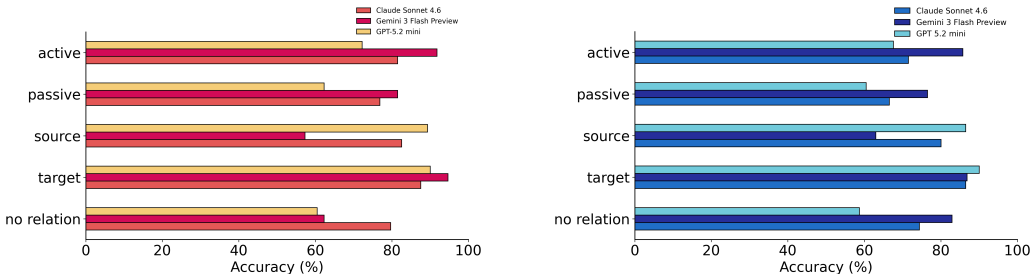


Figure 4: Performance on natural language understanding of biotic interaction queries grouped by language specificity: (left) scientific terminology; (right) vernacular names. For each query type we report accuracy.

subset of 28K images, depicting 281 unique interactions of taxa from *Animalia*, *Plantae*, and *Fungi*, balancing data quality, taxonomic specificity, and visual diversity across seven interaction types. We then generate five semantically equivalent statements using the aforementioned perturbations. To stress-test robustness, we also generate two contradictory statements that deliberately invert entity roles and relations. Appendix C lists all query variants for each of the 281 unique interactions.

Models. We evaluate CLIP-style models for text-to-image retrieval, such as OpenAI [47], MetaCLIP [8], SigLIP [63], SigLIP2 [53], as well as specialized models, such as BioTrove-BioCLIP [60], BioCAP [64], BioCLIP [51] and BioCLIP2 [21]. We adopt proprietary multimodal language models (MLLMs) such as GPT-5.4 mini [41], Claude Sonnet 4.6 [3], and Gemini 3 Flash Preview [19] for ranking prompting for each image and corresponding queries per image: *Does this image show {some query}? Answer with "Yes" or "No" and nothing else.* All proprietary models were evaluated via their APIs with deterministic decoding (temperature = 0.0).

Metrics. We report the Average Precision at k (AP@50), Mean Reciprocal Rank (MRR), and Recall at k (Recall@10). We also report accuracy and flip rate over MLLMs performance. Beyond accuracy, we measure whether models answer consistently across query paraphrases: if the model produces inconsistent predictions across this set—i.e., at least one prediction differs—we count this as a flip. The flip rate is the proportion of such instances over the dataset.

Results. We report retrieval evaluation of meaning preserving queries on BIOINTERACT100 in Table 5 and Figure 3 (see also extensive results in Appendix D). High-quality training data and scale is crucial for expert-level queries. While models like BioCLIP and BioTrove-BioCLIP were pre-trained on images and taxonomic label queries using TreeOfLife-10M [51], they underperform comparatively to larger general models like SigLIP and SigLIP2. BioCLIP2, the largest specialized model (trained on TreeOfLife-200M [21]) achieves the highest overall retrieval score. Although trained on the same dataset as BioCLIP and BioTrove-BioCLIP, BioCAP’s complementary synthetic descriptive caption alignment beyond taxonomic labels achieves competitive performance with BioCLIP2.

Table 7: Performance on adversarial relation-direction reversal queries. We group meaning-preserving queries with their syntactically similar but semantically contradictory counterparts: **active** queries specifying the full triple (s, r, t) , contrasted with (t, r, s) ; **passive** queries inverting the relation (t, r', s) , contrasted with (s, r', t) . **Bold** entries indicate the **best** results.

| Model | By active | | | By passive | | |
|----------|-------------|-------------|-------------|-------------|-------------|-------------|
| | AP@50 | MRR | Recall@10 | AP@50 | MRR | Recall@10 |
| SigLIP2 | 0.30 | 0.41 | 0.47 | 0.27 | 0.39 | 0.47 |
| BioCAP | 0.30 | 0.43 | 0.41 | 0.27 | 0.41 | 0.48 |
| BioCLIP2 | 0.32 | 0.44 | 0.49 | 0.28 | 0.41 | 0.48 |

240 Different queries present challenges of varying difficulty. Figure 3 illustrates difference in perfor-
 241 mance across the five semantically equivalent statements. We observe that, for specialized models,
 242 queries containing only the interacting entities—without explicitly specifying the relation—yield the
 243 best performance. This is likely because images in BIOINTERACT100 often depict small organisms,
 244 with the visual signal dominated by environmental context providing cues about other relevant
 245 organisms.

246 High accuracy on paraphrasing does not imply consistency. We report performance of general-
 247 purpose proprietary models in Table 6 and Figure 4. First, we evaluate a text-only control condition in
 248 which prompts refer to an image despite no image being supplied. Strong performance in this setting
 249 reflects correct detection of absent visual input rather than visual recognition capability, with Gemini
 250 variant suggesting weaker modality-awareness and lower overall calibration compared to GPT and
 251 Claude variants. In image-text settings, MLLMs struggle to maintain stable predictions as reflected
 252 by their high flip rates. Claude Sonnet 4.6 seems to achieve the best overall balance between accuracy
 253 and consistency, exhibiting both strong performance and lower flip rates compared to other models. In
 254 this setting, apparent gains in accuracy can reflect overfitting to benchmark adversarial perturbations
 255 rather than genuine semantic understanding. Furthermore, as shown in Figure 4, domain-specific
 256 language has no significant impact on the performance of general-purpose proprietary models.

257 Biotic interactions are direction-sensitive, and even the best models are affected by reversal. We
 258 report image retrieval evaluation of relation-direction reversals adversarial queries in Table 7. We find
 259 that these limitations are particularly pronounced in real-world, direction-sensitive settings, where
 260 models frequently fail to distinguish between natural adversarial instances.

261 **Limitations.** While our labels rely on robust community consensus and literature validation, they
 262 remain susceptible to occasional human error. We report significant zero-shot results, showing
 263 that pretrained embeddings can effectively navigate this complex data. However, applying few-
 264 shot learning—where models are calibrated using a handful of gold-standard examples—presents a
 265 clear path to further enhance performance through task-specific adaptation. Secondly, in-the-wild
 266 community science imagery lacks explicit visual bounding boxes, and interacting organisms may
 267 be heavily occluded or out of focus. Future work is required to disentangle whether the observed
 268 performance gaps stem from a failure in fine-grained multi-entity localization or a genuine deficit in
 269 relational reasoning. Consequently, BIOINTERACT100 naturally serves as a challenging, real-world
 270 testbed for both few-shot semantic adaptation and future models with localized visual grounding.

271 5 Conclusion

272 We introduce BIOINTERACT, the largest publicly available multimodal dataset designed to better
 273 evaluate biotic interactions understanding and ultimately, accelerate trustworthy AI solutions for
 274 biodiversity. This dataset focuses on fine-grained semantic understanding, surpassing existing datasets
 275 in scale and annotation richness. Our controlled adversarial benchmark over BIOINTERACT100, an
 276 image retrieval task on a fixed set, underscore the importance of stress-testing multimodal under-
 277 standing, particularly for tasks that require inferring meaning from visual evidence and expressing it
 278 through diverse, semantically equivalent natural language forms. With BIOINTERACT, we aim to
 279 accelerate the development of multimodal AI models for biodiversity that are robust to the nuanced
 280 challenges of species interactions in real-world environments.

References

- 281
- 282 [1] Luis Abdala-Roberts, Adriana Puentes, Deborah L Finke, Robert J Marquis, Marta Montserrat,
283 Erik H Poelman, Sergio Rasmann, Arnaud Sentis, Celia C Symons, Nicole M van Dam, et al.
284 Connecting the dots: Managing species interaction networks to mitigate the impacts of global
285 change. *eLife*, 14:e98899, 2025.
- 286 [2] Kumail Alhamoud, Shaden Alshammari, Yonglong Tian, Guohao Li, Philip HS Torr, Yoon Kim,
287 and Marzyeh Ghassemi. Vision-language models do not understand negation. In *Proceedings*
288 *of the Computer Vision and Pattern Recognition Conference*, pages 29612–29622, 2025.
- 289 [3] Anthropic. Claude sonnet 4.6. <https://www.anthropic.com/claude/sonnet>, 2026.
290 Accessed 2026-05-06.
- 291 [4] Fernanda Baena-Díaz and Wesley Dáttilo. Plant traits driving visits by the honeybee *apis*
292 *mellifera* across native and introduced ranges: A global analysis. *Ecological Entomology*, 51
293 (1):191–201, 2026.
- 294 [5] Florencia Baudino, Ramiro R Ripa, Jorgelina Franzese, and Victoria Werenkraut. Using citizen
295 science as a research prioritization tool to detect co-occurrences of the invasive species *Harmonia*
296 *axyridis*. *Insect Conservation and Diversity*, 2025.
- 297 [6] Pengcheng Chen, Jin Ye, Guoan Wang, Yanjun Li, Zhongying Deng, Wei Li, Tianbin Li,
298 Haodong Duan, Ziyang Huang, Yanzhou Su, et al. Gmai-mmbench: A comprehensive multimodal
299 evaluation benchmark towards general medical ai. *Advances in Neural Information Processing*
300 *Systems*, 37:94327–94427, 2024.
- 301 [7] Zerui Cheng, Stella Wohnig, Ruchika Gupta, Samiul Alam, Tassallah Abdullahi, João Alves
302 Ribeiro, Christian Nielsen-Garcia, Saif Mir, Siran Li, Jason Orender, et al. Benchmarking is
303 broken—don’t let ai be its own judge. *arXiv preprint arXiv:2510.07575*, 2025.
- 304 [8] Yung-Sung Chuang, Yang Li, Dong Wang, Ching-Feng Yeh, Kehan Lyu, Ramya Raghavendra,
305 James Glass, Lifei Huang, Jason Weston, Luke Zettlemoyer, et al. Meta clip 2: A worldwide
306 scaling recipe. *arXiv preprint arXiv:2507.22062*, 2025.
- 307 [9] Alessio Cocchieri, Luca Ragazzi, Giuseppe Tagliavini, and Gianluca Moro. Remedqa: Are we
308 done with medical multiple-choice benchmarks? In *Proceedings of the 19th Conference of the*
309 *European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*,
310 pages 2706–2738, 2026.
- 311 [10] GloBI Community. Global Biotic Interactions: Interpreted Data, January 2025. URL <https://doi.org/10.5281/zenodo.14640564>. Accessed via inaturalist.org on 2025-12-13.
312
- 313 [11] GloBI Community. Global Biotic Interactions: Interpreted Data
314 Products hash://md5/e76bf914309ad27dce6ab911d8854590 hash://
315 /sha256/ba79836caab5b7ba2d7d659123d27c89f4ad990bd50f97ded935edee9fbe9f87 ,
316 January 2025. URL <https://doi.org/10.5281/zenodo.14640564>.
- 317 [12] Jean-François Doherty, Antoine Filion, Jerusha Bennett, Upendra Raj Bhattarai, Xuhong
318 Chai, Daniela de Angeli Dutra, Erica Donlon, Fátima Jorge, Marin Milotic, Eunji Park, et al.
319 The people vs science: can passively crowdsourced internet data shed light on host–parasite
320 interactions? *Parasitology*, 148(11):1313–1319, 2021.
- 321 [13] Encyclopedia of Life. Encyclopedia of life (eol), 2025. URL <https://eol.org>. Accessed:
322 2026-03-25.
- 323 [14] Federico Felizzi, Olivia Riccomi, Michele Ferramola, Francesco Andrea Causio, Manuel
324 Del Medico, Vittorio De Vita, Lorenzo De Mori, Alessandra Piscitelli, Pietro Eric Risuleo,
325 Bianca Destro Castaniti, et al. Are large vision language models truly grounded in medical im-
326 ages? evidence from italian clinical visual question answering. *arXiv preprint arXiv:2511.19220*,
327 2025.

- 328 [15] Aruna Gauba, Irene Pi, Yunze Man, Ziqi Pang, Vikram S Adve, and Yu-Xiong Wang. Ag-
329 mmu: A comprehensive agricultural multimodal understanding benchmark. *arXiv preprint*
330 *arXiv:2504.10568*, 2025.
- 331 [16] GBIF Secretariat. GBIF Backbone Taxonomy, 2023. URL [https://doi.org/10.15468/3](https://doi.org/10.15468/39omei)
332 9omei. Accessed via GBIF.org on 2025-12-30.
- 333 [17] GBIF.org. Global biodiversity information facility. <https://www.gbif.org>, 2025. Accessed
334 on on 2025-12-25.
- 335 [18] Global Invasive Species Database. 100 of the world’s worst invasive alien species. [http:](http://www.iucngisd.org/gisd/100_worst.php)
336 [://www.iucngisd.org/gisd/100_worst.php](http://www.iucngisd.org/gisd/100_worst.php), 2026. Accessed: 26-03-2026.
- 337 [19] Google. Gemini 3 flash preview. [https://ai.google.dev/gemini-api/docs/models/g](https://ai.google.dev/gemini-api/docs/models/gemini-3-flash-preview)
338 [emini-3-flash-preview](https://ai.google.dev/gemini-api/docs/models/gemini-3-flash-preview), 2025. Gemini API model documentation, accessed 2026-05-06.
- 339 [20] Quentin Groom, Nadja Pernat, Tim Adriaens, Maarten de Groot, Sven D. Jelaska, Diana
340 Marčiulynienė, Angeliki F. Martinou, Jiří Skuhrovec, Elena Tricarico, Ernst C. Wit, and
341 Helen E. Roy. Species interactions: next-level citizen science. *Ecography*, 44(12):1781–1789,
342 2021. doi: <https://doi.org/10.1111/ecog.05790>. URL [https://nsojournals.onlinelibr](https://nsojournals.onlinelibrary.wiley.com/doi/abs/10.1111/ecog.05790)
343 [ary.wiley.com/doi/abs/10.1111/ecog.05790](https://nsojournals.onlinelibrary.wiley.com/doi/abs/10.1111/ecog.05790).
- 344 [21] Jianyang Gu, Samuel Stevens, Elizabeth G Campolongo, Matthew J Thompson, Net Zhang,
345 Jiaman Wu, Andrei Kopanev, Zheda Mai, Alexander E White, James Balhoff, et al. Bio-
346 clip 2: Emergent properties from scaling hierarchical contrastive learning. *arXiv preprint*
347 *arXiv:2505.23883*, 2025.
- 348 [22] Yu Gu, Jingjing Fu, Xiaodong Liu, Jeya Maria Jose Valanarasu, Noel CF Codella, Reuben Tan,
349 Qianchu Liu, Ying Jin, Sheng Zhang, Jinyu Wang, et al. The illusion of readiness in health ai.
350 *arXiv preprint arXiv:2509.18234*, 2025.
- 351 [23] Hongyong Han, Wei Wang, Gaowei Zhang, Mingjie Li, and Yi Wang. Coralvqa: A large-
352 scale visual question answering dataset for coral reef image understanding. *arXiv preprint*
353 *arXiv:2507.10449*, 2025.
- 354 [24] Madison Hernandez, Paul Masonick, and Christiane Weirauch. Crowdsourced online images
355 provide insights into predator-prey interactions of putative natural enemies. *Food Webs*, 21:
356 e00126, 2019.
- 357 [25] Fang-Shuo Hu, Yun Hsiao, and Alexey Solodovnikov. A global citizen science effort via
358 inaturalist reveals food webs of large predatory rove beetles. *Food Webs*, 43:e00399, 2025.
- 359 [26] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark,
360 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv*
361 *preprint arXiv:2410.21276*, 2024.
- 362 [27] iNaturalist. We estimate the accuracy of Research Grade observations to be 95% correct!,
363 January 2024. URL [https://www.inaturalist.org/blog/89255-we-estimate-t](https://www.inaturalist.org/blog/89255-we-estimate-the-accuracy-of-research-grade-observations-to-be-95-correct)
364 [he-accuracy-of-research-grade-observations-to-be-95-correct](https://www.inaturalist.org/blog/89255-we-estimate-the-accuracy-of-research-grade-observations-to-be-95-correct). Accessed:
365 2026-01-02.
- 366 [28] iNaturalist. <https://www.inaturalist.org/>, 2025. Accessed via inaturalist.org on
367 2025-12-13.
- 368 [29] iNaturalist. What is geoprivacy? what does it mean for an observation to be obscured? [https:](https://help.inaturalist.org/en/support/solutions/articles/151000169938-what-is-geoprivacy-what-does-it-mean-for-an-observation-to-be-obscured-)
369 [://help.inaturalist.org/en/support/solutions/articles/151000169938-wha](https://help.inaturalist.org/en/support/solutions/articles/151000169938-what-is-geoprivacy-what-does-it-mean-for-an-observation-to-be-obscured-)
370 [t-is-geoprivacy-what-does-it-mean-for-an-observation-to-be-obscured-](https://help.inaturalist.org/en/support/solutions/articles/151000169938-what-is-geoprivacy-what-does-it-mean-for-an-observation-to-be-obscured-),
371 2025. iNaturalist Help. Modified on March 12, 2025. Accessed April 6, 2026.
- 372 [30] IUCN. The iucn red list of threatened species, 2025. URL <https://www.iucnredlist.org>.
373 Downloaded on 2025-04-07; accessed via GBIF.org on 2026-03-25.

- 374 [31] Faizan Farooq Khan, Xiang Li, Andrew J Temple, and Mohamed Elhoseiny. Fishnet: A large-
375 scale dataset and benchmark for fish recognition, detection, and functional trait prediction. In
376 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 20496–20506,
377 2023.
- 378 [32] Mario Koddenbrock, Rudolf Hoffmann, David Brodmann, and Erik Rodner. On the domain ro-
379 bustness of contrastive vision-language models. In *German Conference on Artificial Intelligence*
380 (*Künstliche Intelligenz*), pages 62–76. Springer, 2025.
- 381 [33] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan
382 Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint*
383 *arXiv:2408.03326*, 2024.
- 384 [34] Fan Liu, Jindong Han, Tengfei Lyu, Weijia Zhang, Zhe-Rui Yang, Lu Dai, Cancheng Liu, and
385 Hao Liu. Foundation models for scientific discovery: From paradigm enhancement to paradigm
386 transition. *arXiv preprint arXiv:2510.15280*, 2025.
- 387 [35] Xuan Liu, Siru Ouyang, Xianrui Zhong, Jiawei Han, and Huimin Zhao. Fgbench: A dataset and
388 benchmark for molecular property reasoning at functional group-level in large language models.
389 *arXiv preprint arXiv:2508.01055*, 2025.
- 390 [36] Georgiana Manolache, Gerard Schouten, and Joaquin Vanschoren. Crypticbio: A large multi-
391 modal dataset for visually confusing species. *arXiv preprint arXiv:2505.14707*, 2025.
- 392 [37] M Maruf, Arka Daw, Kazi S Mehrab, Harish B Manogaran, Abhilash Neog, Medha Sawhney,
393 Mridul Khurana, James P Balhoff, Yasin Bakış, Bahadır Altintas, et al. Vlm4bio: A benchmark
394 dataset to evaluate pretrained vision-language models for trait discovery from biological images.
395 *Advances in Neural Information Processing Systems*, 37:131035–131071, 2024.
- 396 [38] Peter Mikula, Pavel Pipek, Martin Bulla, María L Castillo, Shawan Chowdhury, Łukasz
397 Dylewski, Franz Essl, Josh A Firth, Jérôme MW Gippet, Theresa Henke, et al. Harness-
398 ing ecology data to uncover invasive species behaviour. *Methods in Ecology and Evolution*,
399 2026.
- 400 [39] OBO Relations Ontology Consortium. Relations ontology (ro). <https://obofoundry.org/ontology/ro.html>, 2024. Accessed on on 2025-12-15.
- 402 [40] Observation International and local partners. Observation.org dataset. <https://observation.org>, 2026. Accessed via observation.org on 2025-12-13.
- 404 [41] OpenAI. Gpt-5.4 mini. [https://developers.openai.com/api/docs/models/gpt-5.4](https://developers.openai.com/api/docs/models/gpt-5.4-mini)
405 -mini, 2026. OpenAI API model documentation, accessed 2026-05-06.
- 406 [42] Kaixuan Pan, Leon Marshall, Geert R de Snoo, and Jacobus C Biesmeijer. Dutch landscapes
407 have lost insect-pollinated plants over the past 87 years. *Journal of Applied Ecology*, 61(6):
408 1323–1333, 2024.
- 409 [43] Connor T Panter and Arjun Amar. Sex and age differences in the diet of the eurasian spar-
410 rowhawk (*accipiter nisus*) using web-sourced photographs: exploring the feasibility of a new
411 citizen science approach. *Ibis*, 163(3):928–947, 2021.
- 412 [44] Nadja Pernat, Susan Canavan, Marina Golivets, Jasmijn Hillaert, Yuval Itescu, Ivan Jarić,
413 Hjalte MR Mann, Pavel Pipek, Cristina Preda, David M Richardson, et al. Overcoming
414 biodiversity blindness: Secondary data in primary citizen science observations. *Ecological*
415 *Solutions and Evidence*, 5(1):e12295, 2024.
- 416 [45] Nadja Pernat, Daniyar Memedemin, Tom August, Cristina Preda, Lien Reyserhove, Jens
417 Schirmel, and Quentin Groom. Extracting secondary data from citizen science images reveals
418 host flower preferences of the mexican grass-carrying wasp *isodontia mexicana* in its native and
419 introduced ranges. *Ecology and Evolution*, 14(6):e11537, 2024.
- 420 [46] Jorrit H Poelen, James D Simons, and Chris J Mungall. Global biotic interactions: An open
421 infrastructure to share and analyze species-interaction datasets. *Ecological informatics*, 24:
422 148–159, 2014.

- 423 [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
424 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
425 models from natural language supervision. In *International conference on machine learning*,
426 pages 8748–8763. PmLR, 2021.
- 427 [48] Binesh Sadanandan and Vahid Behzadan. Psf-med: Measuring and explaining paraphrase
428 sensitivity in medical vision language models. *arXiv preprint arXiv:2602.21428*, 2026.
- 429 [49] Manu E Saunders, Emma K Goodwin, Karen CBS Santos, Carolyn A Sonter, and Romina Rader.
430 Cavity occupancy by wild honey bees: need for evidence of ecological impacts. *Frontiers in*
431 *Ecology and the Environment*, 19(6):349–354, 2021.
- 432 [50] Davinder Singh, Naman Jain, Pranjali Jain, Pratik Kayal, Sudhakar Kumawat, and Nipun Batra.
433 Plantdoc: A dataset for visual plant disease detection. In *Proceedings of the 7th ACM IKDD*
434 *CoDS and 25th COMAD*, pages 249–253. 2020.
- 435 [51] Samuel Stevens, Jiaman Wu, Matthew J Thompson, Elizabeth G Campolongo, Chan Hee Song,
436 David Edward Carlyn, Li Dong, Wasila M Dahdul, Charles Stewart, Tanya Berger-Wolf, et al.
437 Bioclip: A vision foundation model for the tree of life. In *Proceedings of the IEEE/CVF*
438 *conference on computer vision and pattern recognition*, pages 19412–19424, 2024.
- 439 [52] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela,
440 and Candace Ross. Winoground: Probing vision and language models for visio-linguistic
441 compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
442 *Recognition*, pages 5238–5248, 2022.
- 443 [53] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alab-
444 dulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip
445 2: Multilingual vision-language encoders with improved semantic understanding, localization,
446 and dense features. *arXiv preprint arXiv:2502.14786*, 2025.
- 447 [54] Edward Vendrow, Omiros Pantazis, Alexander Shepard, Gabriel Brostow, Kate E Jones, Oisín
448 Mac Aodha, Sara Beery, and Grant Van Horn. Inquire: A natural world text-to-image retrieval
449 benchmark. *Advances in Neural Information Processing Systems*, 37:126500–126514, 2024.
- 450 [55] Edward Vendrow, Julia Chae, Rupa Kurinchi-Vendhan, Isaac Eckert, Jazlynn Hall, Marta
451 Jarzyna, Reymond Miyajima, Ruth Oliver, Laura Pollock, Lauren Shrack, et al. Inquire-search:
452 A framework for interactive discovery in large-scale biodiversity databases. *arXiv preprint*
453 *arXiv:2511.15656*, 2025.
- 454 [56] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing
455 Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception
456 of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- 457 [57] Qiyao Wei, Edward Morrell, Lea Goetz, and Mihaela van der Schaar. Semantic-kg: Using
458 knowledge graphs to construct benchmarks for measuring semantic similarity. *arXiv preprint*
459 *arXiv:2511.19925*, 2025.
- 460 [58] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a
461 comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern*
462 *analysis and machine intelligence*, 41(9):2251–2265, 2018.
- 463 [59] Zhenlin Xu, Yi Zhu, Siqi Deng, Abhay Mittal, Yanbei Chen, Manchen Wang, Paolo Favaro,
464 Joseph Tighe, and Davide Modolo. Benchmarking zero-shot recognition with vision-language
465 models: Challenges on granularity and specificity. In *Proceedings of the IEEE/CVF Conference*
466 *on Computer Vision and Pattern Recognition*, pages 1827–1836, 2024.
- 467 [60] Chih-Hsuan Yang, Benjamin Feuer, Talukder Jubery, Zi Deng, Andre Nakkab, Md Zahid Hasan,
468 Shivani Chiranjeevi, Kelly Marshall, Nirmal Baishnab, Asheesh Singh, et al. Biotrove: A large
469 curated image dataset enabling ai for biodiversity. *Advances in Neural Information Processing*
470 *Systems*, 37:102101–102120, 2024.

- 471 [61] Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan
472 Sun, Botao Yu, Ge Zhang, Huan Sun, et al. Mmmu-pro: A more robust multi-discipline
473 multimodal understanding benchmark. In *Proceedings of the 63rd Annual Meeting of the*
474 *Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15134–15186, 2025.
- 475 [62] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When
476 and why vision-language models behave like bags-of-words, and what to do about it? *arXiv*
477 *preprint arXiv:2210.01936*, 2022.
- 478 [63] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for
479 language image pre-training. In *Proceedings of the IEEE/CVF international conference on*
480 *computer vision*, pages 11975–11986, 2023.
- 481 [64] Ziheng Zhang, Xinyue Ma, Arpita Chowdhury, Elizabeth G Campolongo, Matthew J Thomp-
482 son, Net Zhang, Samuel Stevens, Hilmar Lapp, Tanya Berger-Wolf, Yu Su, et al. Biocap:
483 Exploiting synthetic captions beyond labels in biological foundation models. *arXiv preprint*
484 *arXiv:2510.20095*, 2025.
- 485 [65] Yang Zhou, Tan L Faith, Yanyu Xu, Sicong Leng, Xinxing Xu, Yong Liu, and Rick S Goh.
486 Benchx: A unified benchmark framework for medical vision-language pretraining on chest
487 x-rays. *Advances in Neural Information Processing Systems*, 37:6625–6647, 2024.

488 **NeurIPS Paper Checklist**

489 **1. Claims**

490 Question: Do the main claims made in the abstract and introduction accurately reflect the
491 paper’s contributions and scope?

492 Answer: [Yes]

493 Justification: We present BioInteract, the largest publicly available multimodal dataset
494 comprising richly annotated images depicting interactions between organisms, or biotic
495 interactions, providing a natural testbed for tasks involving images and unconstrained, free-
496 form natural language, as interacting organisms are discerned from images alone and their
497 relationship can be expressed through multiple linguistic forms.

498 **2. Limitations**

499 Question: Does the paper discuss the limitations of the work performed by the authors?

500 Answer: [Yes]

501 Justification: We discuss Limitations in Section 4.

502 **3. Theory assumptions and proofs**

503 Question: For each theoretical result, does the paper provide the full set of assumptions and
504 a complete (and correct) proof?

505 Answer: [N/A]

506 Justification: We submit a dataset & evaluation paper.

507 **4. Experimental result reproducibility**

508 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
509 perimental results of the paper to the extent that it affects the main claims and/or conclusions
510 of the paper (regardless of whether the code and data are provided or not)?

511 Answer: [Yes]

512 Justification: All assets are reproducible via our curation pipeline available in our git page
513 <https://github.com/georgianagmanolache/biointeract>.

514 **5. Open access to data and code**

515 Question: Does the paper provide open access to the data and code, with sufficient instruc-
516 tions to faithfully reproduce the main experimental results, as described in supplemental
517 material?

518 Answer: [Yes]

519 Justification: All experimental results are reproducible, the data is avail-
520 able in <https://huggingface.co/datasets/gmanolache/BioInteract>; and script in
521 <https://github.com/georgianagmanolache/biointeract>.

522 **6. Experimental setting/details**

523 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
524 rameters, how they were chosen, type of optimizer) necessary to understand the results?

525 Answer: [Yes]

526 Justification: See section 4 and supplementary material D.

527 **7. Experiment statistical significance**

528 Question: Does the paper report error bars suitably and correctly defined or other appropriate
529 information about the statistical significance of the experiments?

530 Answer: [Yes]

531 Justification: See section 4 and supplementary material D.

532 **8. Experiments compute resources**

533 Question: For each experiment, does the paper provide sufficient information on the com-
534 puter resources (type of compute workers, memory, time of execution) needed to reproduce
535 the experiments?

536 Answer: [Yes]
537 Justification: See section 4 and supplementary material D.

538 **9. Code of ethics**

539 Question: Does the research conducted in the paper conform, in every respect, with the
540 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

541 Answer: [Yes]
542 Justification: See Section 3.2 License, Geoprivacy and Responsible use and supplementary
543 material A.

544 **10. Broader impacts**

545 Question: Does the paper discuss both potential positive societal impacts and negative
546 societal impacts of the work performed?

547 Answer: [Yes]
548 Justification: See Section 3 and supplementary material A.

549 **11. Safeguards**

550 Question: Does the paper describe safeguards that have been put in place for responsible
551 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
552 image generators, or scraped datasets)?

553 Answer: [Yes]
554 Justification: See Section 3.2 & supplementary material A.

555 **12. Licenses for existing assets**

556 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
557 the paper, properly credited and are the license and terms of use explicitly mentioned and
558 properly respected?

559 Answer: [Yes]
560 Justification: See Section 3.

561 **13. New assets**

562 Question: Are new assets introduced in the paper well documented and is the documentation
563 provided alongside the assets?

564 Answer: [Yes]
565 Justification: Code and dataset are accessible through the project website
566 <https://georgianagmanolache.github.io/biointeract>.

567 **14. Crowdsourcing and research with human subjects**

568 Question: For crowdsourcing experiments and research with human subjects, does the paper
569 include the full text of instructions given to participants and screenshots, if applicable, as
570 well as details about compensation (if any)?

571 Answer: [N/A]

572 **15. Institutional review board (IRB) approvals or equivalent for research with human
573 subjects**

574 Question: Does the paper describe potential risks incurred by study participants, whether
575 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
576 approvals (or an equivalent approval/review based on the requirements of your country or
577 institution) were obtained?

578 Answer: [N/A]

579 **16. Declaration of LLM usage**

580 Question: Does the paper describe the usage of LLMs if it is an important, original, or
581 non-standard component of the core methods in this research? Note that if the LLM is used
582 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
583 scientific rigor, or originality of the research, declaration is not required.

584 Answer: [N/A]

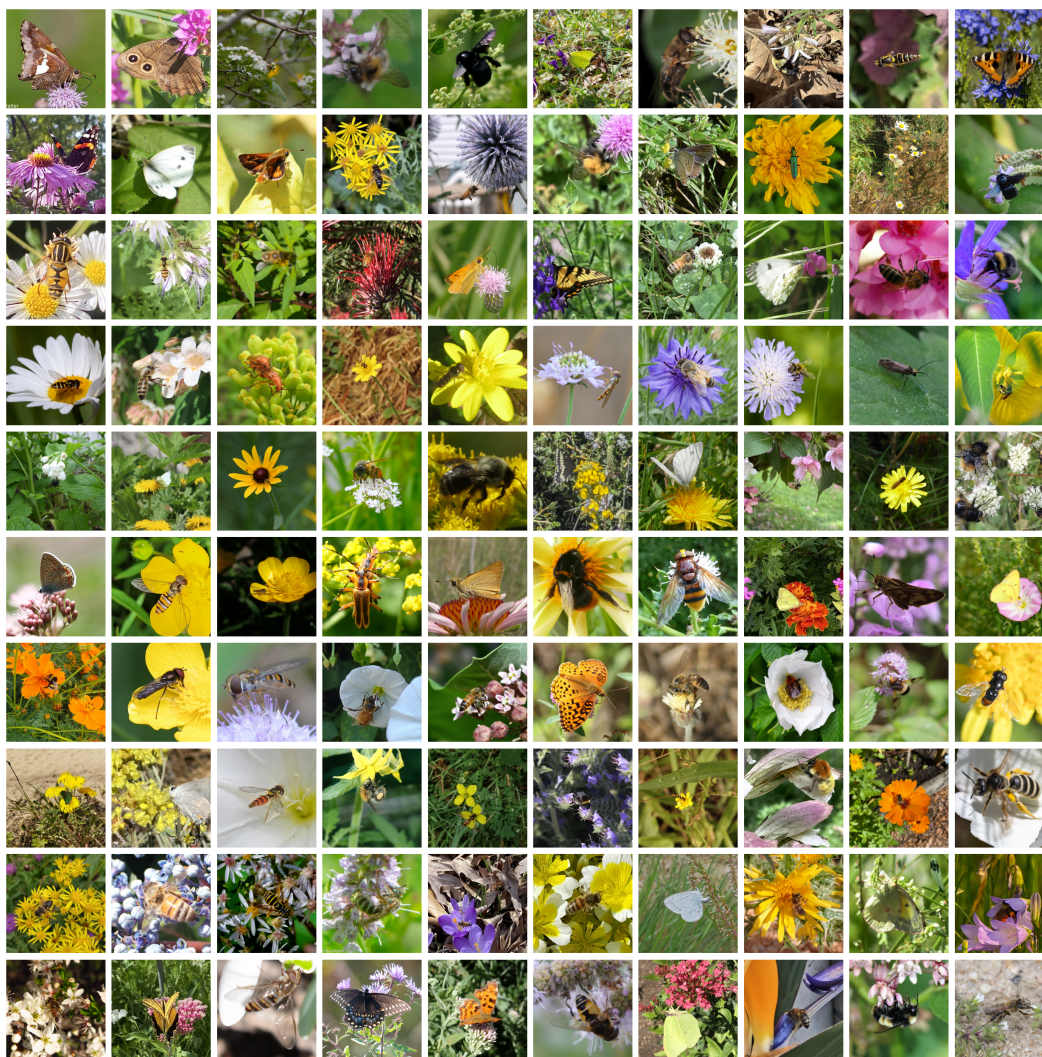


Figure 5: Sample interactions from BIOINTERACT depicting insects visiting plants.

586 BIOINTERACT contains 256K ecological images depicting 15.4K unique biotic interactions between
 587 organisms across five kingdoms, primarily between *Animalia* and *Plantae* (see Figure 6). Biotic
 588 interactions are represented as **triplets** consisting of a source taxon, an interaction type, and a **target**
 589 **taxon**, reported at heterogeneous taxonomic levels [46]. Interactions are standardized to nine types
 590 defined by OBO RO [39] (see Figure 7). Biotic interaction triplets are documented *a priori* via the
 591 iNaturalist platform and cross-validated with at least one additional GloBI-indexed data source [10],
 592 establishing them as reliable ground-truth annotations.

593 Spatio-temporal data plays a crucial role in biodiversity datasets because species distributions and
 594 appearances are strongly influenced by both geographic location and time [36]. Figure 9 illustrates
 595 the spatiotemporal distribution. This information is particularly important for monitoring threatened
 596 species, whose populations are often restricted to small geographic areas or sensitive ecosystems.
 597 Subsequently, it is critically important not to disclose the precise locations of threatened species
 598 because doing so can inadvertently put them at even greater risk. Many vulnerable species face
 599 threats from poaching, illegal wildlife trade, habitat disturbance, and over-collection. Sharing exact
 600 geographic coordinates, especially online or in open databases, can make it possible to locate and
 601 exploit these species. To mitigate the risks associated with the disclosure of sensitive biodiversity
 602 data, iNaturalist implements automatic geoprivacy measures for taxa listed on the global IUCN.

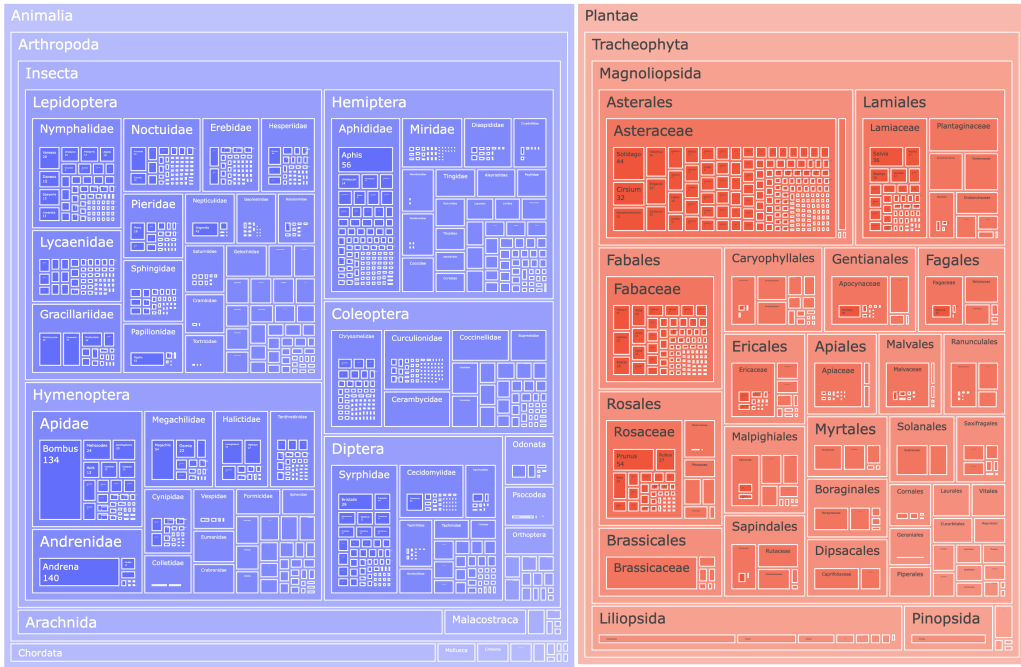


Figure 6: Treemap visualization of the species covered in BIOINTERACT.

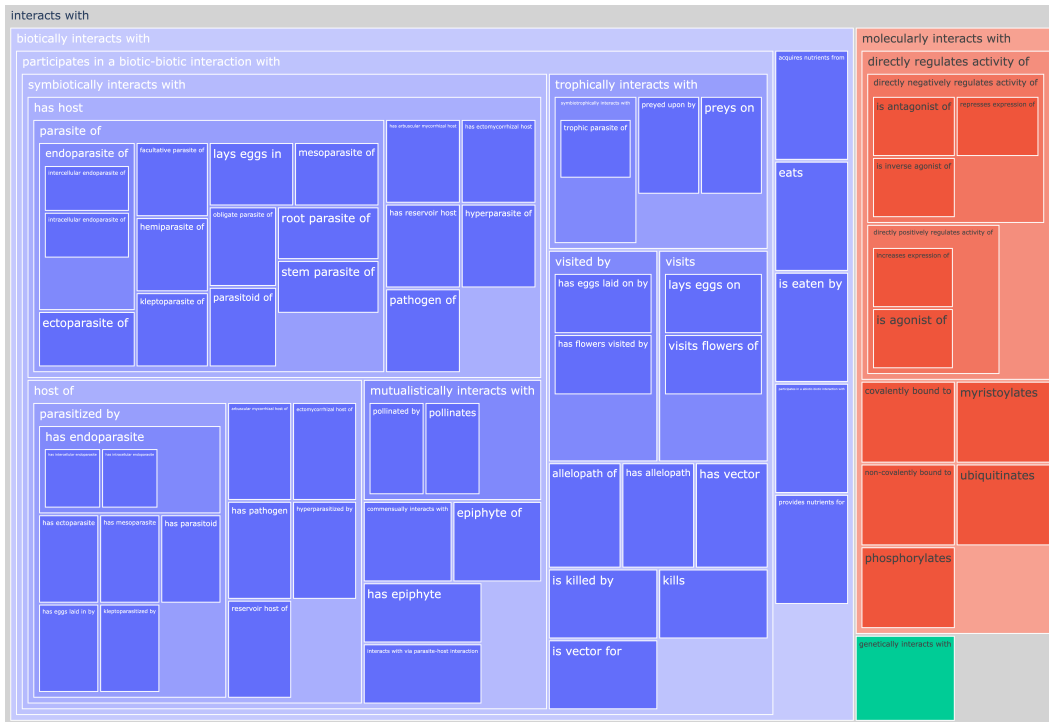


Figure 7: Treemap visualization of the interaction hierarchy used in OBO RO [39].

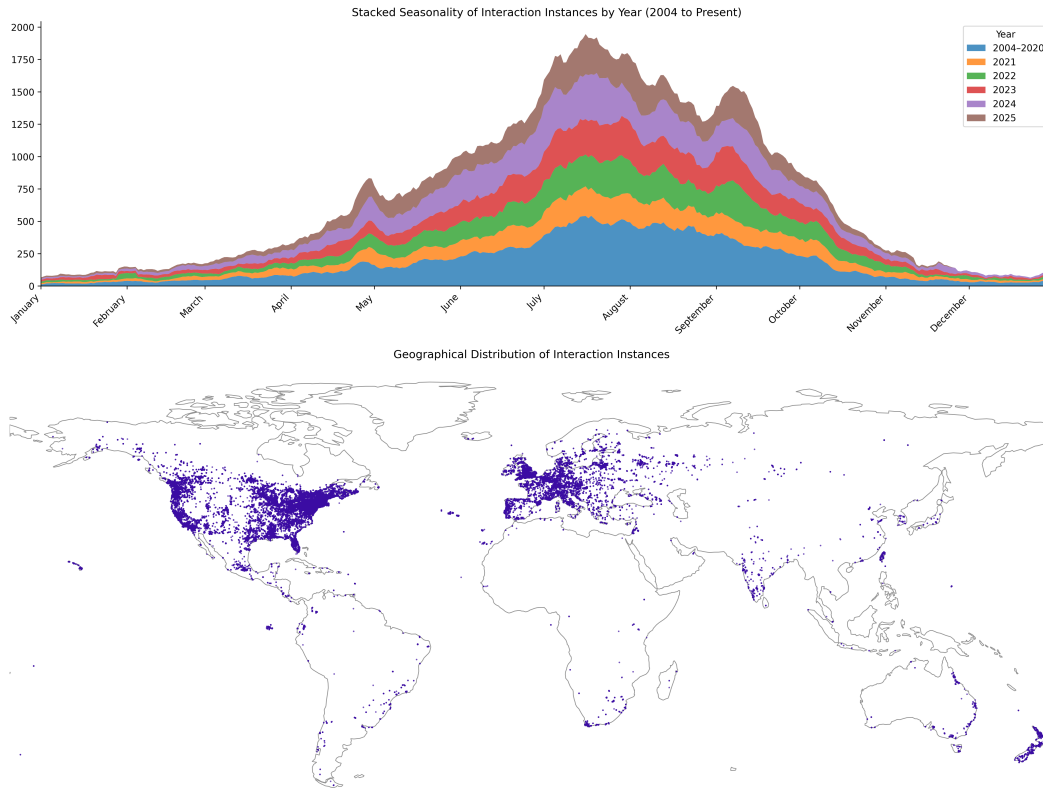


Figure 8: Spatiotemporal distribution of interaction instances in BIOINTERACT: (top) stacked seasonality; (bottom) geographical distribution. Majority of records are concentrated in Europe and North America, with a seasonal peak in observations during July.

603 B Related Work

604 A small body of work leverages online image repositories to manually extract biotic interaction data;
 605 however, these datasets are typically interaction- or taxon-specific and are not released as ML-ready
 606 datasets at scale (see Table 8).

Table 8: BIOINTERACT comparable datasets. Data is manually extracted from sources and is interaction- and taxon-specific.

| Interaction Type | Description | # Images | Sources | Reference |
|------------------------|---|----------|---|-----------|
| Predation | Prey in relation to age and sex of sparrowhawk (<i>Accipiter nisus</i>). | 843 | iNaturalist, Google Images, Macaulay Library, BirdGuides, Facebook, Twitter | [43] |
| Predation | Prey spectrum of assassin bugs (<i>Hemiptera: Reduviidae</i>) | 832 | Flickr, iSpot Nature, BugGuide, Nature-Watch, Google Images | [24] |
| Predation / parasitism | Detect co-occurrences and potential novel interactions between the invasive ladybird (<i>Harmonia axyridis</i>) and local biota (plants, arthropods and fungi). | 783 | iNaturalist, ArgentinNat, WhatsApp | [5] |
| Competition | Competition for nest cavities between bee species and wild honey bee colonies in native and introduced ranges. | 326 | iNaturalist | [49] |
| Flower visitation | Host flowers of the introduced wasp <i>Isodontia mexicana</i> . | 250 | iNaturalist | [45] |
| Host preferences | Host preferences and phenological peak of hairworms inferred from infected hosts. | 20 | iNaturalist | [12] |

607 **C Controlled diagnostic benchmarks**

608 Queries are based on interaction triples (s, r, t) : **active** queries specify the full triple (s, r, t) , **passive**
 609 queries invert the interaction (t, r', s) , **source** queries condition on (s, r) , **target** queries on (t, r') ,
 610 and **no relation** queries specify only entity pairs without the relation (s, t) , **active reversal** queries
 611 swap the source and target from active queries (t, r, s) , and **passive reversal** queries swap the source
 612 and target from active passive (t, r, s) .

Table 9: BIOINTERACT100 query examples for each interaction type.

| Interaction type | Perturbations | Query example |
|-------------------|---|---|
| eats | active passive source target no relation active reverse passive reverse | <i>Tetranychus lintearius</i> eats <i>Ulex europaeus</i> <i>Ulex europaeus</i> is eaten by <i>Tetranychus lintearius</i> <i>Tetranychus lintearius</i> eats another organism <i>Ulex europaeus</i> is eaten by another organism <i>Tetranychus lintearius</i> and <i>Ulex europaeus</i> <i>Ulex europaeus</i> eats <i>Tetranychus lintearius</i> <i>Tetranychus lintearius</i> is eaten by <i>Ulex europaeus</i> |
| has host | active passive source target no relation active reverse passive reverse | <i>Trichilogaster acaciaelongifoliae</i> has host <i>Acacia longifolia</i> <i>Acacia longifolia</i> is hosted by <i>Trichilogaster acaciaelongifoliae</i> <i>Trichilogaster acaciaelongifoliae</i> has host another organism <i>Acacia longifolia</i> is hosted by another organism <i>Trichilogaster acaciaelongifoliae</i> and <i>Acacia longifolia</i> <i>Acacia longifolia</i> has host <i>Trichilogaster acaciaelongifoliae</i> <i>Trichilogaster acaciaelongifoliae</i> is hosted by <i>Acacia longifolia</i> |
| interacts with | active passive source target no relation active reverse passive reverse | <i>Aceria erinea</i> interacts with <i>Juglans regia</i> <i>Juglans regia</i> interacts with <i>Aceria erinea</i> <i>Aceria erinea</i> interacts with another organism <i>Juglans regia</i> interacts with another organism <i>Aceria erinea</i> and <i>Juglans regia</i> <i>Juglans regia</i> interacts with <i>Aceria erinea</i> <i>Aceria erinea</i> interacts with <i>Juglans regia</i> |
| parasite of | active passive source target no relation active reverse passive reverse | <i>Cotesia glomerata</i> parasite of <i>Pieris brassicae</i> <i>Pieris brassicae</i> is parasitized by <i>Cotesia glomerata</i> <i>Cotesia glomerata</i> parasite of another organism <i>Pieris brassicae</i> is parasitized by another organism <i>Cotesia glomerata</i> and <i>Pieris brassicae</i> <i>Pieris brassicae</i> parasite of <i>Cotesia glomerata</i> <i>Cotesia glomerata</i> is parasitized by <i>Pieris brassicae</i> |
| prey on | active passive source target no relation active reverse passive reverse | <i>Philanthus triangulum</i> preys on <i>Apis mellifera</i> <i>Apis mellifera</i> is preyed on by <i>Philanthus triangulum</i> <i>Philanthus triangulum</i> preys on another organism <i>Apis mellifera</i> is preyed on by another organism <i>Philanthus triangulum</i> and <i>Apis mellifera</i> <i>Apis mellifera</i> preys on <i>Philanthus triangulum</i> <i>Philanthus triangulum</i> is preyed on by <i>Apis mellifera</i> |
| visits | active passive source target no relation active reverse passive reverse | <i>Bombus pensylvanicus</i> visits <i>Helianthus annuus</i> <i>Helianthus annuus</i> is visited by <i>Bombus pensylvanicus</i> <i>Bombus pensylvanicus</i> visits another organism <i>Helianthus annuus</i> is visited by another organism <i>Bombus pensylvanicus</i> and <i>Helianthus annuus</i> <i>Helianthus annuus</i> visits <i>Bombus pensylvanicus</i> <i>Bombus pensylvanicus</i> is visited by <i>Helianthus annuus</i> |
| visits flowers of | active passive source target no relation active reverse passive reverse | <i>Agapostemon virescens</i> visits flowers of <i>Heliopsis helianthoides</i> <i>Heliopsis helianthoides</i> is visited by <i>Agapostemon virescens</i> <i>Agapostemon virescens</i> visits flowers of another organism <i>Heliopsis helianthoides</i> is visited by another organism <i>Agapostemon virescens</i> and <i>Heliopsis helianthoides</i> <i>Heliopsis helianthoides</i> visits flowers of <i>Agapostemon virescens</i> <i>Agapostemon virescens</i> is visited by <i>Heliopsis helianthoides</i> |

613 **D Evaluation details**

Table 10: Performance on natural language understanding of biotic interaction queries across BIOINTERACT100 on scientific terminology queries.

| Model | Method | Params (M) | mAP@50 | mMRR | mRecall@10 |
|------------------|----------------------|------------|--------------------|--------------------|--------------------|
| CLIP | ViT-L-14 | 427 | 0.19 ± 0.03 | 0.40 ± 0.04 | 0.59 ± 0.05 |
| MetaCLIP | ViT-L-14-quickgelu | 427 | 0.29 ± 0.04 | 0.51 ± 0.05 | 0.68 ± 0.06 |
| SigLIP | SO400M-14-SigLIP | 877 | 0.34 ± 0.06 | 0.57 ± 0.06 | 0.76 ± 0.05 |
| SigLIP2 | ViT-L-16-SigLIP2-256 | 881 | 0.35 ± 0.07 | 0.59 ± 0.08 | 0.78 ± 0.07 |
| BioTrove-BioCLIP | ViT-B-16 | 149 | 0.19 ± 0.07 | 0.38 ± 0.12 | 0.58 ± 0.14 |
| BioCAP | ViT-B-16 | 149 | 0.45 ± 0.07 | 0.72 ± 0.07 | 0.88 ± 0.03 |
| BioCLIP | ViT-B/16 | 149 | 0.16 ± 0.07 | 0.41 ± 0.12 | 0.60 ± 0.14 |
| BioCLIP2 | ViT-L/14 | 427 | 0.47 ± 0.06 | 0.70 ± 0.06 | 0.85 ± 0.07 |

Table 11: Performance on natural language understanding of biotic interaction queries across BIOINTERACT100 on vernacular terminology queries.

| Model | Method | Params (M) | mAP@50 | mMRR | mRecall@10 |
|------------------|----------------------|------------|--------------------|--------------------|--------------------|
| CLIP | ViT-L-14 | 427 | 0.07 ± 0.02 | 0.24 ± 0.06 | 0.42 ± 0.09 |
| MetaCLIP | ViT-L-14-quickgelu | 427 | 0.14 ± 0.05 | 0.36 ± 0.11 | 0.60 ± 0.13 |
| SigLIP | SO400M-14-SigLIP | 877 | 0.21 ± 0.10 | 0.46 ± 0.14 | 0.67 ± 0.15 |
| SigLIP2 | ViT-L-16-SigLIP2-256 | 881 | 0.23 ± 0.09 | 0.48 ± 0.13 | 0.71 ± 0.13 |
| BioTrove-BioCLIP | ViT-B-16 | 149 | 0.20 ± 0.06 | 0.46 ± 0.09 | 0.69 ± 0.12 |
| BioCAP | ViT-B-16 | 149 | 0.35 ± 0.10 | 0.67 ± 0.14 | 0.87 ± 0.13 |
| BioCLIP | ViT-B/16 | 149 | 0.21 ± 0.06 | 0.53 ± 0.09 | 0.75 ± 0.12 |
| BioCLIP2 | ViT-L/14 | 427 | 0.43 ± 0.12 | 0.70 ± 0.14 | 0.88 ± 0.14 |

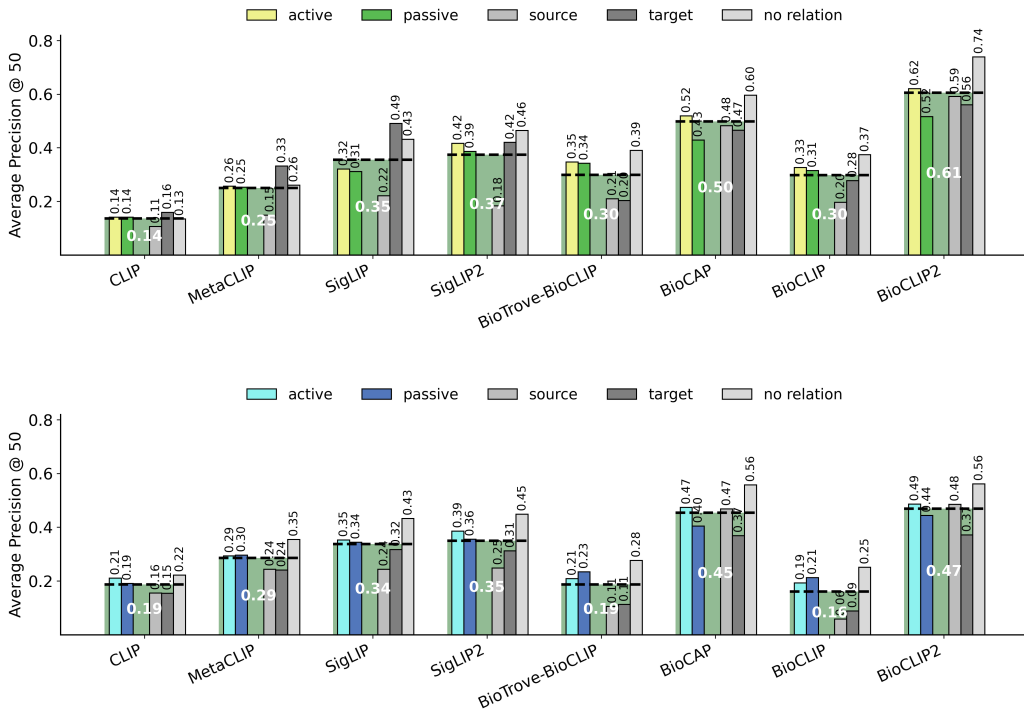


Figure 9: Retrieval performance for each semantic variation query for (top) scientific terminology (bottom) vernacular terminology. Queries are based on interaction triples (s, r, t) : **active** queries specify the full triple (s, r, t) , **passive** queries invert the interaction (s, r', t) , **source** queries condition on (s, r) , **target** queries on (t, r') , and **no relation** queries specify only entity pairs without the relation (s, t) .