



# Humboldt Core – toward a standardized capture of biological inventories for biodiversity monitoring, modeling and assessment

Robert Guralnick, Ramona Walls and Walter Jetz

R. Guralnick (<http://orcid.org/0000-0001-6682-1504>) ([robgur@gmail.com](mailto:robgur@gmail.com)), Univ. of Florida Museum of Natural History, Univ. of Florida at Gainesville, Gainesville, FL, USA. – R. Walls, Cyverse, Univ. of Arizona, Tucson, AZ, USA. – W. Jetz, Dept of Ecology and Evolutionary Biology, Yale Univ., New Haven, CT, USA, and Dept of Life Sciences, Imperial College London, Ascot, Berkshire, UK.

Species inventories, i.e. the recording of multiple species for a specific place and time, are routinely performed and offer particular value for characterizing biodiversity and its change. However, reporting standards allowing these inventories to be re-used, compared to one another, and further integrated with other sources of biodiversity data are lacking, impeding their broadest utility. Here we provide a conceptual informatics framework for capturing, in a standardized and general way, core information about processes underpinning inventory work. We dub this set of terms Humboldt Core, and demonstrate its utility. This proposed framework is based on a community input process, followed by extensive refinement and testing using published inventories. We first develop a typology of inventories and inventory processes, distinguishing between single, elementary inventories, extended and summary inventories, representing increasing levels of sampling event aggregation. We then further describe typical reporting content for inventory processes, along with their value for inferring absence and use in occupancy modeling. Next we provide an overview of the Humboldt Core terms for capture of geospatial, temporal, taxonomic, and environmental conditions, along with methodological descriptors related to the assessment of sampling effort and inventory completeness. Finally, we introduce a pilot implementation of Humboldt Core for the ingestion and provision of inventory process metadata into Map of Life, demonstrating standardized mobilization of metadata from several hundred previously published summary inventories. Humboldt Core helps facilitate integration of different types of inventories and their use in model-supported assessment of spatial biodiversity and its change, critical for meeting global monitoring goals. Humboldt Core will benefit from further enhancements based on community testing and input, but represents a step toward significantly expanding biodiversity dataset discovery, interoperability, and modeling utility for a type of data essential to the assessment of biodiversity variation in space and time.

Both basic inference in ecology and assessments of biodiversity are reliant on a sound, general, and representative empirical knowledge base. Many regional evaluations for biodiversity conservation and especially global science-policy efforts such as IPBES (Intergovernmental Platform on Biodiversity and Ecosystem), GEO BON (Group on Earth Observation Biodiversity Observation Network), and Future Earth are strongly limited by available data (Pereira et al. 2013, Meyer et al. 2015, Proença, et al. 2016). A multitude of data types can inform about the distributions of species and their changes (Jetz et al. 2012), but not all are equally available and used. One type, incidental point records, such as those generated from museum specimens or citizen science contributions, has seen strong recent growth, thanks to technological advancements, especially mobile devices (Newman et al. 2012), extensive digitization efforts by the biocollections community (Beaman and Cellinese 2012), and increased incentives for use of this information in species distribution modeling, threat assessment, and change analysis (Elith and Leathwick 2009, Anderson 2013, Maes et al. 2015). A key facilitator for this growth and associated

use has been the development of the Darwin Core data standard (Wieczorek et al. 2012) for the integration of point record data (e.g. Symbiota, Gries et al. 2014 and the Integrated Publishing Toolkit, Robertson et al. 2014). Due to their greater complexity and more multi-faceted communities of providers and users, other biodiversity data types from typically more formal monitoring efforts are, to date, lacking widely applicable standards. This is likely a key factor behind their lesser mobilization, integration, and re-use, although we recognize that lack of incentives for such efforts are also a significant hindrance (Enke et al. 2012). With increased urgency and recognition for bringing more harmonized and well-described biodiversity data and metadata into global assessments, now is an important moment to enable contributions covering a wider gamut of biodiversity evidence.

One core biodiversity data type that to date has lacked a more general integration is the taxonomic inventory. Inventories are related to incidental point records, but differ in at least one key aspect: within a sampling event or period, they address multiple rather than single biological entities,

for example, multiple species. Inventories (sometimes called local species lists, surveys, or samples) set out to catalog a location (a defined area of land or a volume of water) for a particular group of organisms over a defined time, using a specified approach. Sometimes, but not always, they include some form of abundance or biomass estimates. Inventories have a defined temporal and spatial scope, for example characterizing survey plots or transects, volumetric samples, grid cells, or even counties or islands over timeframes from minutes (and sometimes seconds) to years (Leon-Cortes et al. 1998, Nakamura and Soberón 2009, Isaac and Pocock 2015).

Thanks to increasing recognition of national monitoring efforts (Pereira et al. 2013, Schmeller et al. 2015, Proença et al. 2016) and especially citizen science schemes with at least minimal protocols (Dickinson et al. 2012, Pescott et al. 2015), inventory data have seen dramatic and still increasing growth, but without concomitant attention to their most effective capture, sharing, and use. Inventories have the potential to directly inform about which species co-exist or may be absent in a given area and time period, but this utility depends on the suitability of a given sampling protocol to address a given taxonomic and spatiotemporal scope, which in turn affects ‘inventory completeness’ – the proportion of present species successfully detected. Even if not reliably complete, inventories can vitally inform about potential absences through advanced modeling and can help overcome the limitations of presence-only models (Lobo et al. 2010).

The lack of means to describe inventories processes in a standard way has, to date, strongly impeded their integration and re-use for biodiversity distribution and change assessments, despite the obvious and documented value of such integration (Kremen 1994, Corona et al. 2011). Instead, inventory data outputs and descriptions of how inventories were performed are most often stored in local databases, making it extremely challenging to reconcile them with other outputs. More recently, inventory data are also sometimes provisioned in flat file formats, such as Darwin Core Archives, which were not developed with inventory processes in mind and cannot easily capture their complexity, although, as discussed more below, new publishing approaches begin to close that gap. At worst, inventory data are left in printed text in the pages of journals, the gray literature, or file cabinets.

Although inventory reporting and outputs are often tied to the particulars of the taxa of interest and science or conservation outcomes, we aim to show that inventories and the sampling processes that generate them can be placed in a logical and general framework built around a core set of information. This core includes reporting on the scope (spatial, temporal, taxonomic, and environmental), methods, and measures of effort that are most critical to capture for re-use and re-integration in new downstream analyses. These analyses include the more informed use of repeated inventories to assess species-specific detection probabilities in occupancy models (Dorazio et al. 2006, Kéry et al. 2010, Wintle et al. 2012, van Strien et al. 2013, Iknayan et al. 2014) and, for cases of representative and well-documented sampling, assessing probabilistic species absence (Tyre et al. 2003, Lobo et al. 2010, Szabo et al. 2010, Sadoti et al. 2013, Lahoz-Monfort et al. 2014).

Building on existing standards, we provide an initial step towards the capture of standardized information about inventory processes, thus providing a critical starting point for harmonizing reporting and re-use of this type of data. We dub this effort the Humboldt Core, for one of the pioneering compilers of spatial biodiversity knowledge, Alexander Von Humboldt, and show its utility for capturing content from a wide variety of exemplar inventory efforts. We begin by developing a typology of different inventories, focusing on differences in spatial and temporal scope, their implications for reporting requirements, and the ultimate utility of such inventories for inferring absence or use in occupancy models. While the Humboldt Core is meant to serve as a singular resource that can encompass describing inventories and survey process, we show how its application to different types of inventories may require subsets of terms and provide example ‘profiles’. Finally, we demonstrate how the proposed terms and capture method may be operationalized by way of example in a web-based submission process implemented in Map of Life (<<http://mol.org>>). A unified set of terms addressing taxonomic inventory processes, combined with mobilization and sharing mechanisms for inventory data, has the potential to enable a significantly enhanced knowledge base for macroecology, biogeography, and global change research.

## **Toward community metadata standards for biological inventories – setting the stage**

The importance of community standards in biodiversity research has grown in the past decade with the rise of data sharing platforms such as the Global Biodiversity Information Facility (GBIF; [gbif.org](http://gbif.org)), the Ocean Biogeographic Information System (OBIS; [iobis.org](http://iobis.org)), Map of Life (MOL, [mol.org](http://mol.org)), and others. These platforms, which republish digitized content with support from a network of providing scientists, data managers, and organizations, rely on standards that assure data and their descriptions are in common formats, thus allowing users to discover content from multiple providers using the same search terms. Darwin Core has been the key standard for describing species point records and has served as the basis for interoperability of taxonomic and occurrence-based datasets. Darwin Core describes species point records by formally defining a set of terms with clearly defined meanings (Wieczorek et al. 2012). However, it has its basis in the natural history collections community and was not initially intended to capture metadata about multi-species sampling processes.

Although Darwin Core has limited expressiveness for capturing a full accounting of inventory scope and processes, recent efforts have begun to develop an ‘event core’, as well as new terms that to capture certain inventory aspects. The terms support capture of sizes of sample from a sample event (e.g. `dwc:sampleSizeValue` and `dwc:sampleSizeUnit`) and the ability to publish related data from multiple events (via the new `dwc:parentEventID` term) to represent certain forms of plot data and nested plot designs. The use of ‘dwc:’ in front of terms specifies that they are from the Darwin Core vocabulary. More details about the event core in relation to sampling designs is discussed in Wieczorek et al. (2014). Similar approaches in marine systems, but focused

more on linking in-situ measurements with sampling events, have also been modeled as part of the Ocean Biogeographic Information System (OBIS, Costello pers. comm.). Darwin Core and its extensions are much needed steps for aggregating data but still **provide limited ability to express detailed reporting of scope** (e.g. prospective taxonomic scope and exclusions), as well as a whole suite of commonly measured aspects of inventory sampling processes shared among inventories (e.g. direct or inferred measures of multiple types of sampling effort). Description of sampling processes is needed for later data reuse, even in cases where they result in no reported outcome taxa. Limitations of current approaches can be overcome using existing mechanisms, and an extension to the Darwin Core event model, that further captures **inventory process metadata related to those events, may provide a logical path for such integration.**

The Ecological Metadata Language (EML), on the other hand, was explicitly designed to capture data and metadata about ecological studies. It is an XML-based method for formalizing and standardizing the set of terms and their related concepts that are essential for describing ecological data (<[http://sbc.lternet.edu/external/EML/EML\\_documents/eml\\_metadata\\_guide.html](http://sbc.lternet.edu/external/EML/EML_documents/eml_metadata_guide.html)>) including inventories or observational data sets. However, EML is broader in intent, as are allied efforts that have a stronger basis in semantics, such as OBOE (Extensible Observation Ontology; (Madin et al. 2007)) and O&M (Observations and Measurements; Cox 2013). As such, none of these models capture the needed granularity to describe inventory processes, but rather rely on extensions. Some more specific vocabularies have been developed for special use cases such as the Bird Monitoring Data Exchange (BMDE; Kelling et al. 2009), but with limited generality of the metadata terms addressed. Much more common than an explicit standardized vocabulary are project specific standardized datasheets used by broad-scale projects, such as the Christmas Bird Count (<<http://netapp.audubon.org/cbcoobservation/>>), or VegBank (<<http://vegbank.org/vegbank/index.jsp>>). We note that there are efforts underway to more fully standardize vegetation plot data (Wiser 2016).

At the same time, efforts to develop standard vocabularies in the arena of formal ontology development are widely expanding with more directed impact on the biodiversity science community. For taxonomic inventories, the most relevant ontology is the Biological Collections Ontology (BCO; Walls et al. 2014a). Building off the Ontology for Biomedical Investigations (OBI, Bandrowski et al. 2016), the BCO refines the ‘event model’ used in Darwin Core to explicitly specify any type of planned processes used in biodiversity studies, thereby encompassing more types of taxonomic inventory processes. Given its remit, BCO’s multi-level structure provides a vital area of overlap between descriptions of taxonomic inventory processes and other types of biological research and their data outputs such as sequence or phenotype data (Walls et al. 2014b).

If combined by data holders and mobilizers in an appropriate and formalized way, efforts such as BCO, the Darwin Core event model, and related observation and measurement models have the potential to provide standardized vocabulary and formats for recording and describing a broad array of biodiversity and ecological metadata in a flexible framework.

However, each on its own covers a limited set of use cases and, even if used jointly, gaps remain to make data most (re)-usable for many specific applications, including downstream inference and modeling based on careful descriptions of biodiversity inventories. There remains a pressing need for a broadly applicable vocabulary for the description of biodiversity inventory processes.

## Methods and results

### Developing the taxonomic inventory metadata vocabulary

Development of the list of terms that make up Humboldt Core was seeded through an expert workshop (see Supplementary material Appendix 1 for a list of participants) dedicated to capturing domain-specific inventorying practices. Experts clarified the semantics of taxonomic inventory processes, utilizing existing terminology often discussed in inventory literature but not rigorously defined. The output of this effort was a set of definitions of ‘taxonomic inventory processes’ in the BCO.

An expert-provided list of example inventory cases allowed the development of a more formal inventory typology, recognizing variation in spatial and temporal coverage as a key distinguishing factor (Fig. 1, Table 1). Finally, inventory experts provided a list of terms needed to define the intended and realized scope of a survey, along taxonomic, geographic, environmental, and temporal dimensions. These were captured by break-out groups that worked with existing inventory reports covering a wide gamut of inventory processes and dimensions, with the goal of identifying the essential and consistent elements. The same groups were asked to develop inventory metadata required to capture methodologies used during surveys, including how those methodologies, along with scope, provide a means to assess sampling effort when not directly reported.

### Evaluation and implementation

Over the course of 24 months, the utility of the initial Humboldt Core term list was iteratively tested and refined as part of a mobilization of published inventory studies into **Map of Life** (<<http://mol.org/datasets>>). Student workers checked each study for basic quality metrics (e.g. is there a well-defined species list, a defined areal extent, etc.), assessed metadata, and staged species lists and metadata for full ingest into Map of Life. This included digitizing the species lists and capturing the polygon(s) describing the shape of the surveyed area. These data were then made available along with other kinds of species distribution data, so that they are discoverable through the Map of Life website and services and can be integrated with other data, information, and knowledge.

Utilizing an initial set of inventories, a small team made needed improvements to the draft set of Humboldt Core terms. A full recording of best practices was produced as part of that process ([goo.gl/0J9PAM](http://goo.gl/0J9PAM)). Once refinements were complete, a finalized list of terms was generated, and

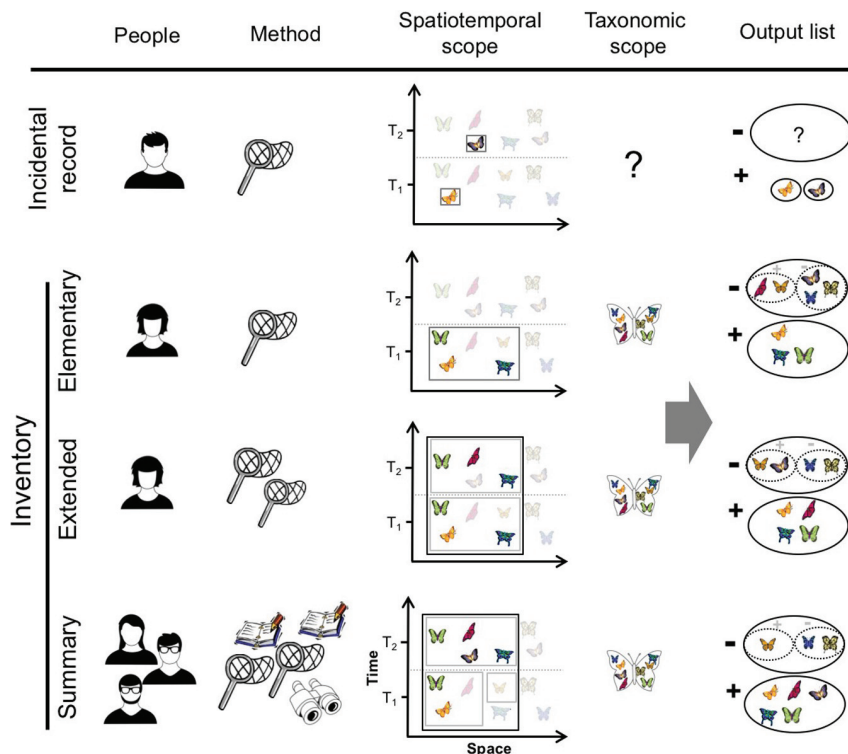


Figure 1. Overview of main inventory types. Incidental sampling (row 1) lacks a taxonomic scope and only provides presence information (+). In contrast, inventories (row 2–4) have a defined taxonomic scope (e.g. a global or regional list of butterflies, here eight species), and a spatiotemporal scope larger than a single point. The taxonomic scope enables identification of non-detections (–), composed of true absences (grey–) and false absences (gray+). Increased duration/sampling effort (row 3) may increase presences and decrease true and false absences. While ‘elementary inventories’ address a single sampling event (and are therefore often limited in spatiotemporal scope), ‘extended inventories’ combine multiple single, component events (grey boxes). ‘Summary inventories’ (row 4) represent a greater aggregation yet and involve multiple sampling parties and/or methods, often extending over larger spatial and temporal scales. As the sampling effort for a given spatiotemporal scope increases with the number and type of sampling events aggregated for it, false absence rates are expected to decrease from elementary toward extended and summary inventories – depending on protocol and event-level effort (Table 1).

we further tested Humboldt Core by selecting a diversity of exemplar inventories covering a wide range of habitats and taxa, representing both national-scale and long-term inventories. Examples include the Christmas Bird Count and the North American Butterfly Association Counts, as well as replicated vegetation plot surveys across North America that are part of the VegBank project. Seven case-examples showing metadata capture using Humboldt Core are provided in Supplementary material Appendix 2 and available online at <https://mol.org/humboldtdcore/>.

## Definition, hierarchy and typology of inventories

Essential for the harmonization of data and metadata about ‘biological inventories’ is an appreciation of its generality and role as an operational umbrella for a vast variety of biodiversity data. Using the BCO term [http://purl.obolibrary.org/obo/BCO\\_0000048](http://purl.obolibrary.org/obo/BCO_0000048), we define a **taxonomic inventory** as: ‘A list of names ascribed to biological entities of specified organismal scope recorded over a defined spatial and temporal scope following a described sampling protocol and sampling effort, potentially including values indicating abundance or biomass of the biological entities’. The biological entities are usually species and the organismal scope

is typically a taxon, but potentially also an ecologically or functionally defined grouping (e.g. ‘water birds’ or ‘trees’) or other operational unit (e.g. genes). The output list would typically contain multiple entities, but may be as short as a single entity or, if a sampling effort focused on a certain taxon yielded no results, even contain **none**.

Inventories are differentiated by key characteristics that have implications for their suitability for spatial biodiversity inference (Fig. 1, Table 1) such as the number of people conducting the survey, methods, spatiotemporal and taxonomic scope, and type of output. At the highest level, all inventories may be contrasted with the incidental report of a taxon presence, which, like an inventory, may be well defined in spatiotemporal scope, but lacks a taxonomic scope and thus, by necessity, also lacks information about non-detections and potential absences (first row, Fig. 1). Inventories, by contrast, have a defined, prospective taxonomic scope, but vary in level of aggregation, and accordingly, spatial and temporal scope (rows 2–4, Fig. 1). How inventory biologists aggregate and report outputs from inventory processes is fully dependent on the intent of those performing the inventory, as discussed more below. An ‘elementary inventory’ represents a single sampling event that covers an often small spatiotemporal extent, such as a single sweep of a net, a single trap, or a single survey plot. Aggregated reporting of



Table 1. Summary typology of different types of inventories, specific examples, typical reporting content based on that typology, and the value of those different inventory types for inferring absence and use in occupancy modeling. Abbreviations: BBS – Breeding Bird Survey; CBC – Christmas Bird Count; CTFS – Center for Tropical Forest Science; GBIF – Global Biodiversity Information Facility; MOL mobile – Map of Life Mobile app; NABA – North American Butterfly Association; NAWQA – National Water-Quality Assessment Program; TEAM – Tropical Ecology Assessment and Monitoring; WWF – World Wildlife Fund. Scoring for ‘raw data’ represents whether data about observations of individuals or populations are available, as opposed to summary lists. ‘Somewhat’ refers to variable and often limited reporting of direct observations.

Description	Protocol	General examples	Specific examples	Temporal scope	Geographic scope	Taxonomic scope	Raw data	Contributors, competency	Inventory search protocols	Inference about effort	Suited for absence inference	Suited for occupancy models
Incidental observation	No protocol given or recorded	Typical museum specimen, incidental citizen science data	GBIF, MOL mobile, iSpot, iNaturalist, Observado	Very short	Very small	Unclear	Yes	Single, heterogeneous competency	Usually unclear	No	No	No
Elementary inventory	Limited protocol	Citizen science data, natural history collecting campaigns	eBird ‘all’ reporting, BirdTrack, Firefly Watch	Short to medium	Small to medium	Usually clear	Yes	Single, heterogeneous competency	Single, clear, open or opportunistic searches	Usually OK	Potentially	Potentially
		Observations within delineated area: standard area survey report (e.g. single atlas grid cell, transect count)	BBS, CBC, NABA Counts, PollardBase, National Plant Monitoring Scheme	Usually short	Small	Usually clear	Yes	Single, often high competency	Single, clear, open or restricted searches	Very good	Potentially	Potentially
	Explicit protocol	Stationary, somewhat mobile organisms: camera traps, insect traps, small mammal traps; marine, freshwater traps	TEAM camera trap network, pitfall trapping	Short to medium	Small to medium	Clear	Somewhat	Single, often high competency	Single, clear, open or restricted searches	Usually OK	More limited	Potentially
Extended inventory		Multi-point or transect sampling campaigns: fish, zooplankton netting	International Bottom Trawl Survey, Reef Life Survey, Census of Marine Zooplankton	Short to medium	Very small	Usually clear	Somewhat	Single, often high competency	Single, somewhat clear, typically restricted	Usually OK	Potentially	Potentially
		Stationary, little mobile organisms: vegetation plots, Revell plots	CTFS forest plots, Vegbank	Usually short	Small	Usually clear	Somewhat	Single, often high competency	Single, clear, restricted searches	Usually OK	More limited	Potentially
Summary inventory	Limited protocol	Multiple parties and sources, limited effort reporting: species lists for delineated areas (countries, reserves, etc.).	WWF ecoregions, Biological Inventories of World Protected Areas	Often long	Often large	Clear	No	Often many, heterogeneous competency	Heterogeneous, often unclear	Limited	Yes	No
	Explicit protocol	Multiple parties, some effort reporting: area survey report (e.g. atlas grid cell summary)	Atlas of South African Birds, Atlas of amphibians and reptiles in Europe	Medium to long	Often large	Clear	No	Often many, often high competency	Somewhat heterogeneous, clear	Possible	Yes	Possibly

Humboldt Core Version 1	Single Sourced Inventories (e.g., transect count, trapping and netting, gridded atlas survey)	Summary Inventories (e.g., CTFs forest, relevé plots, protected area species lists)
General Dataset & Identification Terms	inventory performed by; dataset name, identifier, publisher, license, rights holders; metadata recorded by; citation reference and id; taxa identified by; identification quality; cited taxonomic authority	
Geospatial & Habitat Scope Terms	geospatial scope; areal extent; total area inventoried; number of sites; site names and details; lat/long by site; elevation range and units; habitats included and excluded	
Temporal Scope Terms	survey time blocks; start and end year, month day; time units spent in blocks; daily start, end time; study diurnality, study season	
Taxonomic Scope Terms	prospective taxonomic scope inclusion and exclusion; distribution status included and excluded; developmental stage included and excluded, size classes included and excluded	
Methodology Description Terms	inventory type; protocol name, detail, citation; reference; abundances reported Y/N; absences reported Y/N?	inventory type; compiled data Y/N, type; abundances and/or absences reported?; absence list
Completeness & Effort Terms	effort reporting, lower/upper bound, granular breakdown; effort method; vouchers/samples taken, if so, how? completeness reporting if necessary	completeness reported, how assessed; inferred taxonomic completeness upper/lower bound, how assessed

Figure 2. Humboldt Core Version 1 summary organized into main categories e.g. scope, methodology and effort with a summary overview of the kinds of information captured, such as granular aspects of temporal scope e.g. study season and diurnality. Different application profiles usable for different inventory types (e.g. single-sourced elementary and extended versus summary inventories) may use subsets of Humboldt Core terms. This list is informative only, and the full term list is available in Supplementary material Appendix 3.

multiple such events, i.e. repeated for the same location in time or extended to another locale, extends the spatiotemporal scope to an inventory, but in this ‘extended inventory’ a consistency is retained through sampling that is identical in approach, and/or done by the same, single party. We refer to these as ‘single source’ inventories in Fig. 2. Table 1 provides further details of key protocols, data resources, and exemplar projects.

In both the elementary and extended inventory case, usefulness for spatial biodiversity inference is affected by the strictness to which a standardized protocol is adhered. Thus we distinguish between single, opportunistic searches that lack rigorous sampling methodologies or clear spatial delineations, such as those performed in some citizen science efforts, and inventories that specify a more explicit protocol, such as those used in vegetation surveys, plankton tows and algal sampling used in water monitoring programs, or transect-based observational methods such as Pollard walks.

At a higher level of aggregation, summary inventories (row 4, Fig. 1) may combine studies using multiple protocols, processes and observers, with often variable reporting of the methods employed and sources that include direct project data, other compiled data sources, and literature. These inventories are typically aggregates of multiple broad surveys

performed within well-defined, often broad spatial sampling units (usually grids or circle counts) and clear, replicable methods. In some cases, such inventories may only address select species of economic or other interest, such as expert- or multi-source-based presence–absence of pests or disease vectors for whole countries. These types of inventories typically do not provide needed information for occupancy-based approaches and need careful vetting of sampling effort in order to assess completeness or provide reliable absence data, as portions of geographic space and time periods within their scope may not be sampled.

### Humboldt Core: toward a new metadata standard for inventories

We present a summary list of inventory terms in Fig. 2, separated into profiles for terms more relevant for elementary and extended inventories and those with specific application in summary inventories. Supplementary material Appendix 3 provides more information on the full list of terms and their relations to existing standards. We provide ready to use template spreadsheets in Excel and CSV format at <http://mol.org/humboldtcare>.

We fully support the re-use of existing standards where appropriate, and terms in the Humboldt Core are related to terms in Dublin Core (<http://dublincore.org/documents/dcmi-terms/>) Darwin Core, and EML. The columns for related terms in Supplementary material Appendix 3 clearly show that while no single existing standard is adequate for describing taxonomic inventories, many existing terms are available to reuse as part of Humboldt Core. Some terms allow multiple entries, e.g. taxonomic authority allows more than one authority to be used to identify species or higher-level units. Others are only relevant depending on conditional prior information. For example, 'compilation types' is only relevant if the 'Were compiled data included?' term had 'yes' as a value. 'Compilation types' is also a good example of a term that has a controlled vocabulary (that includes museum specimens, expert knowledge, literature, or other sources). Not all fields do, but those with such constraints are noted in the 'Description' field in Supplementary material Appendix 3. As Humboldt Core is still in an early form, we note that terms labels in Supplementary material Appendix 3 correspond directly to labels that were be used in data collection questionnaires (e.g. 'Did authors provide information about effort?'). In future work, as Humboldt Core becomes formalized into a metadata vocabulary, each term will be assigned a brief label (e.g. 'is effort reported') and an URL identifier for use in RDF and ontology-based applications. Below, we provide further detail on Humboldt Core, focusing on terms describing: inventory geospatial, temporal, taxonomic scope along with environmental conditions; methodological descriptors include types of inventory processes; and terms for assessment of effort and completeness.

### **Reporting on scope**

Scope covers geospatial, temporal, taxonomic scope, and environmental conditions. The main geospatial characteristics are the textual description of the area surveyed and a quantification of the total areal extent surveyed, along with the granular information on site locations and areal extents (e.g. latitude and longitude ranges). Temporal scope terms include information about total length in reported time units of a checklist or survey process, along with more granular information such as the daily start and end time of inventorying, whether the inventory included sampling in day, night, or both, and the study seasons. Summary surveys over repeated years and seasons, and that span activities of diurnal and nocturnal species, are likely to be more complete.

Taxonomic scope includes information about both the higher taxonomic units that were explicitly of interest to the surveyors and those taxa that might have been intentionally excluded. This yields two key terms, 'prospective taxonomic scope' and 'taxonomic group(s) excluded from study'. Although values for these terms often are simply taxon units such as 'Aves', many surveys explicitly mention other characteristics that are critical for defining scope such as 'non-volant' or 'large', and these are fully captured in the exhaustive list of terms (e.g. the terms 'size class' or 'size class excluded'). In some cases, the survey scope is simply a list of target species, typically (but not always) covering a specific higher-level group.

Environmental conditions refer to specific environmental descriptions within the spatial scope considered. 'Habitat

inclusion' and 'habitat exclusion' offers fields for biological conditions linked to offers a two- or even three-dimensional detail about the search area (e.g. 'forest canopy', 'tree holes', 'scrub', 'littoral', etc.) that may be dynamic in time and in practice could not easily be captured with a pure geographic delineation. Humboldt Core also includes terms for reporting on the ground conditions such as weather that help provide full understanding of presences and absences (e.g. some species will not be visible if it is too hot or raining).

### **Inventory process descriptions**

Simple, broad descriptions of inventory search processes provide an essential mechanism to communicate key characteristics of how inventories were performed. These are reported as part of the methods utilized during an inventory, under the term 'inventory type' (Supplementary material Appendix 3), representing a high-level summary of inventory search process type. For this term, we provide a controlled vocabulary for core classes of taxonomic inventory search processes as elements shared by both the Biological Collections Ontology (BCO) and Humboldt Core. The controlled vocabulary terms are.

1) Restricted search: a taxonomic inventory process that is restricted to plots, transects, or points, in which a person or group of people is comprehensively covering the entire area, usually with a well-described survey time or pace. The search is restricted to a defined and human-scale geospatial area (usually traversable within a time course of less than a day) within which there is an expectation of a comprehensive accounting of the taxonomic items of interest.

2) Open search: a taxonomic inventory process in which the search is restricted within a larger defined geographic area, but where effort isn't even or complete across the region, and thus not a comprehensive accounting of taxa of interest. Temporal duration is typically longer than restricted searches, lasting hours to several days.

3) Opportunistic search: a taxonomic inventory process that is a more casual reporting of occurrences of taxa of interest, still intended to be a comprehensive accounting of the taxa of interest, but with no pre-specified investment of effort nor planned trajectory for discovery, thus of often idiosyncratic length or spatial scope.

4) Trap or sample inventory: a taxonomic inventory process that is typically restricted in geospatial extent that involves either the physical extraction of some evidence of the presence of the taxa of interest, such as a whole organisms, scat, fur, other material samples or information artifacts such as photographs or sound recordings

5) Incidental/adventitious: a taxonomic inventory process in which taxon occurrences are recorded as co-variables of another study, or by happenstance, and later compiled as a taxonomic inventory.

These terms are already available in the BCO through their shared superclass 'taxonomic inventory process' ([http://purl.obolibrary.org/obo/BCO\\_0000047](http://purl.obolibrary.org/obo/BCO_0000047)) which is defined as 'A planned process by which a taxonomic inventory is created'. We use the attribute 'inventory search process' in Humboldt Core and specify that researchers should supply a subclass of BCO's 'taxonomic inventory process' as the value. This use of BCO for Humboldt Core provides a path for

further integration of Humboldt Core-annotated data with other types of biodiversity data in the future.

### ***Reporting on methodology and effort***

Methodology includes not only the type of inventory process but also protocol details, which may be named protocols (e.g. Pollard walks, Carolina Vegetation Survey) or cited references to the protocol. The other needed aspect of methodology is whether absences and abundances were explicitly reported. If absences are reported, a related field captures further information on taxa listed as absent. Abundances are sometimes capped at a maximum value, and if so, a separate attribute captures that information (e.g. 'reported abundance cap = 100'). Information about whether vouchers were taken, whether and where those specimens were deposited, and information about what else might have been measured from the specimens is also recorded. Finally, data quality assurance steps reported by authors are also captured (e.g. checking localities for accuracy, reporting uncertainties).

Effort reporting is heterogeneous across all inventories. In only a few cases is this directly reported or inferred via species accumulation curves or other statistical assessment for summary inventories. In some cases, effort can be inferred by metadata curators from the time periods (e.g. detailed temporal scope) when surveys were conducted. For spatially restricted, elementary surveys, quantification of effort is often much more detailed. In such cases, it is much more common to have direct reporting of the upper and lower estimate effort bounds and how those estimates were calculated. Humboldt Core provides a means to report whether effort was directly reported by the surveyor or author of a survey report or derived post-hoc by an analyst.

### ***Determining overall compilation and survey effort***

Many summary inventories make use of existing compilations of data when producing species lists. These compiled data sources may include museum data, literature sources, and expert knowledge. Compilation effort is particularly important to at least qualify, because it is used either directly or indirectly to assess completeness of the inventory. A combination of effort and quality of the assembled compiled information determines the reliability of inventory completeness assessments. We note some published inventories were only compilations – that is, no new data were collected. But instead effort was put towards compiling existing, often scattered resources, reconciling them, and producing a new knowledge product.

We are not aware of prior attempts to measure 'compilation effort' in a methodologically sound way. Given the varying ways previously compiled data are used in inventories, we chose to qualify effort as low, medium or high. High compilation effort is when multiple types of sources are consulted, each source type has more than one and preferably many resources consulted within that type, and the source quality is considered high (e.g. peer reviewed papers versus grey literature, well curated specimen collections in known museums versus amateur collections). Low effort is where only one type of resource is moderately or poorly consulted, or that resource quality is relatively low. We also considered effort low in cases where reporting is poor about the sources used. Medium effort reflects either excellent use

of only one type of resource, or shallow use of multiple types. We recognize that the current term definition mixes quality and quantity of effort in an attempt to capture a summary output, and further refinement may be required.

### ***Reporting survey completeness***

Completeness reporting ranges from highly incomplete (i.e. most focal taxa were missed) to presumed complete coverage of all target species. Although such completeness estimates must always be inferred, quantifiable ways to assess it are well developed when performing single, elementary inventories. Occupancy modeling approaches utilize multiple samples from different locations, conditions, and species to predict the variation in detection (Dorazio and Royle 2005, MacKenzie 2006). This approach can be applied to whole assemblages of given taxonomic scope in multi-species occupancy modeling and provides both an estimate of survey completeness as well as presence probabilities for species not recorded (Iknayan et al. 2014, Jarzyna and Jetz 2016). Species accumulation or rarefaction curves provide another potential means to estimate inventory completeness, without requiring multiple samples in space or addressing single species explicitly (Gotelli and Colwell 2001, Chao et al. 2009), but these require a large number of raw samples from the same location, and a sampling event limited to specific habitats, places, or time periods may inflate completeness estimates. Thus, especially for larger spatial or temporal scales, the use of species accumulation approaches is limited. The strongest constraint overall remains the availability of sufficiently detailed raw sampling data. Therefore, ancillary information that allows even approximate estimates on the completeness of a single (in case of raw inventories) or aggregated (in case of extended or summary inventories) sampling events can offer vital information. In many cases, authors provide text descriptions of taxa that were absent during the interval and more general aspects of completeness, given other evidence (e.g. previous work within the geographic area) and their own assessment of their methods. More commonly, however, authors provide no concrete assessment of completeness at all. However, some qualification can still be inferred by an analyst given all the other metadata collected, i.e. the measures of the scope of inventory and the survey and compilation effort. This qualification may be bounded estimates of percentage completeness, with range of values depending on reporting quality as discussed further below in the implementation section.

## **Humboldt Core implementation for summary inventories in Map of Life**

Humboldt Core enables an effective integration and model-based use of vastly different sources and types of inventories for the evaluation of biodiversity distribution and change by standardizing inventory process metadata. An additional benefit is that, especially for summary inventories, it provides a framework for a quantitative mobilization and use of data sources that hitherto remained underutilized. One example is literature-based summary inventories following limited protocols, such as local or regional area checklists. Tens of thousands of such



checklists have been published in one form or another, representing millions of species presence and putative absence records, but they remain unused or at best applied simply as incidental observations.

The mobilization and model-based use of summary inventory data and metadata is a core goal of Map of Life, and a dedicated user-interface and data store for the Humboldt Core-based capture of inventory metadata has now been implemented. This proof of concept implementation also captures the inventory data from such published accounts, reported as a species list and geographic extent. To illustrate the potential of the mobilization and eventual use of summary inventories in practice, we have used the Core to capture metadata from several hundred individual summary inventories from the literature. Figure 3 shows the workflow for assembling inventory data and metadata into Map of Life. Data curators check inventories for suitability for ingestion based on initial quality assessments, with criteria for suitability including enough details on taxonomic scope, species list reporting, and areal extents for re-assembly as digital objects. Qualifying species lists are digitized, with taxonomic names validated against an extensive set of curated taxonomic authority files, in order to assure that source names of organismal entities (e.g. species) match known and valid

scientific names from different curated source lists. Different lists are set as default 'master' lists for different taxa e.g. AmphibiaWeb (<[www.amphibiaweb.org/](http://www.amphibiaweb.org/)>) for amphibians. More documentation on taxonomy tools is available at <<https://api.mol.org/0.x/docs>>. Data curators also digitize the geographic extents described as polygons or use existing gazetteers if those areas are already known, e.g. protected areas or administrative units.

At the same time, inventory reports are passed to metadata curators and analysts. The metadata curators carefully read and assess the reported contents and extract all values related to Humboldt Core fields utilizing Google Spreadsheets. In implementing capture of metadata using Humboldt Core terms, we chose labels that make it the easiest for those curating content to provide the correct information (see Supplementary material Appendix 3 for labels). If these terms are eventually formalized into a metadata vocabulary using RDF, we will use noninformative strings (i.e. numbers) as term identifiers, following contemporary best practice in identifiers. In cases where inference is required when completing reporting, metadata curators report their inference and, in the case of completeness assessments, their rationale. All metadata are then double-checked by a second metadata curator for accuracy.

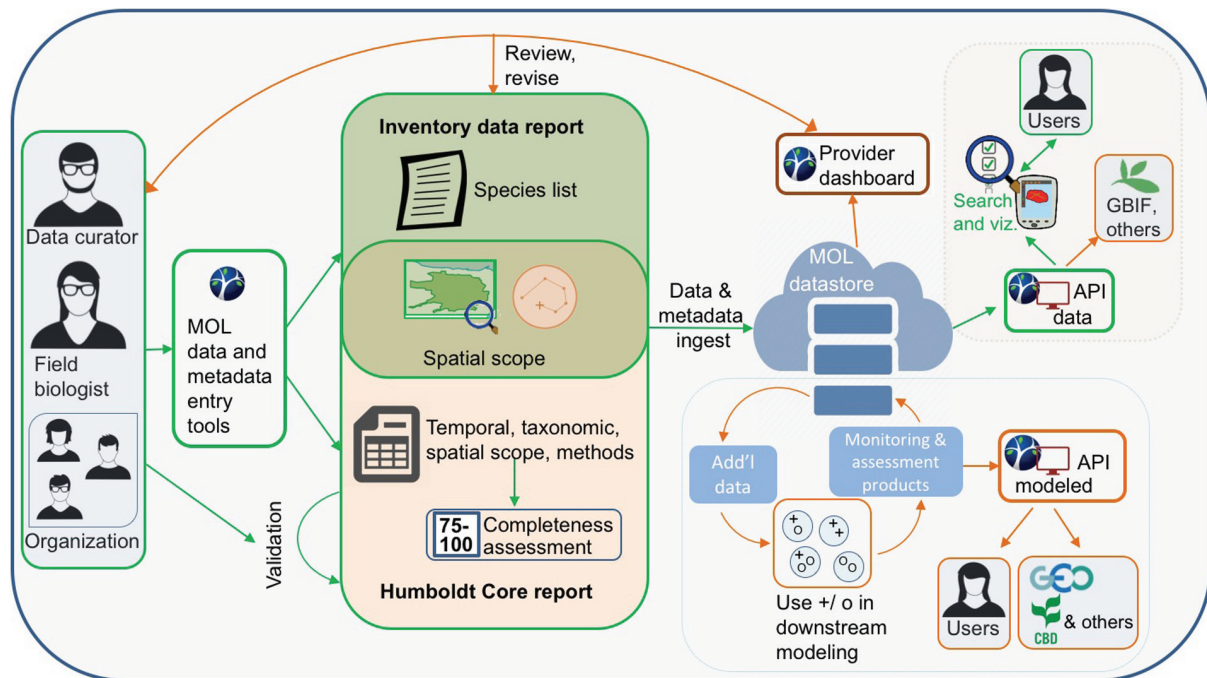


Figure 3. A workflow for the assembly, provision and use of species inventory data and metadata into Map of Life, exemplified here for summary inventories, but applicable more generally. Workflow steps that are operational are marked with green outlined boxes and arrows, and those under implementation with orange colors. Different actors, including data curators working on retrospective capture of inventory data, field biologists with mobile apps, and organizations managing inventory data, all use Map of Life inventory data reporting tools to provide species lists and determine, e.g. via drawing or selection of pre-existing boundaries, the specific spatial scope. Additional metadata define e.g. taxonomic, temporal scopes. Data and metadata are then provisioned into the Map of Life back-end cloud datastore for use by those actors and to expose those data broadly via Map of Life search and discovery mechanisms, and if they wish, to the broader community, such as GBIF, via APIs (Application Programming Interfaces). APIs provide direct access to users via the Web (denoted with screen icons in the API boxes) and for integration programmatically for uses in other infrastructures. While not operational, prototype user dashboards are being developed, so that actors who upload content can manage their data resources, including revising, deleting, or adding further inventories. Presence and probabilistic absence data along with additional, ancillary data are intended to feed into downstream modeling steps, as shown in a modeling loop, producing new knowledge products critical for monitoring and assessment that are themselves stored in Map of Life and made broadly available.

Literature-based summary inventories usually lack detailed completeness assessments but often provide indications or quantify effort in formal (e.g. species accumulation curves) or informal ways (e.g. a self-reported, colloquial, overview assessment of completeness). In our example implementation, we inferentially quantify completeness as a measure with upper and lower bounds that ranged between 0 and 100 in 25 percent increments. We recognize this estimate is imprecise, but finer measures would necessitate a more formal model specification that is ultimately likely to be impractical given reporting heterogeneity. Rather, the goal was to have a repeatable measure that was accurate when scored by any trained analyst. To facilitate training, we developed a best practices manual for this effort ([goo.gl/0J9PAM](http://goo.gl/0J9PAM)). Because all inventories were scored and then carefully checked by another worker, we were able to determine cases of discordance in scoring of inventory completeness. Once workers were trained and heuristics fully developed for scoring completeness, workers rarely had inconsistent scores, with over 90% overlap across inventories and fast resolution of cases of scoring discrepancies.

Those inventories scored as mostly complete were cases where inventory effort was well quantified and reported, where species accumulation curves were produced and shown to reach saturation, or where the habitats in the geographic scope were well sampled over multiple time periods, and where other compiled resources were included where warranted. Scores between 0–25% were often given to rapid or poorly documented assessments over large geographic areas or diverse habitats. Poor reporting also resulted in higher uncertainty about assessed completeness i.e. 25–75% completeness.

Once mobilized in this form, inventory data in Map of Life (MOL) is stored on cloud-based spreadsheets for any further editing, with complete and validated content automatically pushed into their own tables in the MOL PostgreSQL-based datastore. The metadata content is then associated with a globally unique data object identifier that links to the polygon describing inventory location and extent, and species lists that make up the inventory data assets and, where applicable, the literature citation. When a user clicks on the map where an inventory has been performed or checks available datasets (<http://mol.org/datasets>), they can access all the metadata directly or easily download these data for further use. Equally importantly, all metadata can be discovered via application programming interfaces (APIs), which allow them to be shared for broadest use. Future plans for sharing these resources as broadly as possible, including GBIF and GEO-BON among other actors, are highlighted in Fig. 3.

As of July 2017, 402 inventories have been assessed for metadata completeness and have metadata available on Map of Life; a subset of 240 have passed the full data and metadata assembly pipeline, and these are available at <http://mol.org/datasets>. Based on our current sample, nearly 30% of these summary inventories were reported to be between 75–100% complete and many (> 50%) have compiled content and reported abundances, along with nearly 35% directly reporting absence information (see Supplementary material Appendix 4 for more detailed summaries). Exemplar URLs for accessing data, and metadata formatted web and API outputs are provided in Supplementary material Appendix 5.

## Data deposition

Data and metadata for inventories efforts discussed here are available from Map of Life <https://mol.org/datasets/?dt=localinv>.

## Discussion

Taxonomic inventories are one of the most common and critical sources of biodiversity distribution data (Jetz et al. 2012), yet they have rarely been assembled in a way that provides value beyond their original intended scope. Unlocking those data for broadest re-use has been hampered by a lack of reporting standards and of consistent methods to provide an initial assessment of completeness. Our work provides the necessary first steps in this process. Through a rigorous process of community input and refinement, the set of terms provided here has been shown to be detailed enough to capture the elements that describe scope and effort, while general enough to be inclusive of most inventory types. That content in turn yields a novel, albeit initially heuristic, approach to produce completeness assessments.

Although we measured completeness coarsely, such assessment opens the door for developing, for example, prior knowledge (e.g. in a Bayesian framework) at the inventory-level scale and informing about false absence rate at that level. Weightings or priors derived from completeness scores could then be used in modeling the absence of species that are regionally expected but locally non-detected. Such evidence can be combined with information on species-level detection probabilities (MacKenzie et al. 2002) to inform the probabilistic assessment of absence rate. On a more basic level, inventory completeness estimates may assist in establishing a threshold of single species distribution models (Jiménez-Valverde and Lobo 2007) or deriving more realistic estimates of species richness from stacked models (Calabrese et al. 2014). The most complete surveys provide strongest ability to inform about absences, but only ~30% of surveys reached an expert-assessed 75–100% completeness level given the sample of inventories currently available in Map of Life. Thus, our initial assessment of inventory completeness suggests that considerable work is still required to assess absences in most locations.

There are a number of needed next steps for further development of the full term list we have produced. First and foremost, the current version and example implementation of Humboldt Core will benefit from continued community input and involvement. This community includes both experts in standards and ontologies who can specify the formal logic of the taxonomic inventory domain and practicing inventory biologists who can test the standard and determine how to adapt and promote its growth. We argue that a critical next step is to further coordinate inventory data and metadata standards. The Taxonomic Database Working Group (TDWG) – charged with biodiversity standards development, maintenance, and governance – is a logical convener to help coordinate such efforts. However accomplished, further efforts to broadly test and improve all existing approaches to mobilize inventory data and metadata for use in downstream modeling, along with formal

mechanisms for further vetting, community curation, and further technical implementations, are needed.

Our work clearly demonstrates the value of Humboldt Core for retrospective assembly of published inventories, but, perhaps more importantly, the standard can strongly support proper capture of key information about inventories as they are performed in the field. Folding Humboldt Core terms into field information management systems and having both the data and metadata pushed to repositories where they can be most effectively used is desperately needed. We argue that while the full term list should be as complete as needed to be fit for purpose and work broadly, application profiles can allow field information management systems to capture subsets as required for the inventory work in question.

In closing, the power of standardized reporting of inventory metadata, such as might be based on Humboldt Core terms, is providing the community a framework and clear semantics for reporting critical aspects of biodiversity work in a way that is findable, accessible, interoperable, and reusable (Wilkinson et al. 2016). Such activities have been clearly called out as needed in order to meet the goals of the Convention on Biological Diversity (Schmeller et al. 2015) and promise to lower, not increase, reporting burden while assuring clear communication and enhanced re-use value. While Humboldt Core will only benefit from input from the community and deeper connections to other standards efforts, it already represents a useful step toward expanding biodiversity dataset interoperability and producing modeling-relevant metadata for a data type critical to the assessment of biodiversity change.

*Acknowledgements* – We thank the other participants of the Yale workshop in 2013 where Humboldt Core was initially developed: G. Amatulli, J. Belmaker, S. Blum, J. Deck, P. Goldstein, H. Kreft, J. Malczyk, S. Meiri, L. Ries, T. Robertson, N. Robinson, M. Schildhauer, K. Triantis, D. Vieglais, P. Weigelt, A. Wilson, and J. Wiczorek (see Supplementary material Appendix 1 for affiliations). Special thanks to Shai Meiri, Carsten Meyer, and Patrick Weigelt for comments on a previous draft. Ben Carlson, Michelle Duong, Ajay Ranipeta, Jeremy Malczyk, Raphael LaFrance and Luis Villanueva have been instrumental in implementing the Humboldt Core in Map of Life. Rebecca Pederson, Helena De Souza Brasil Barreto, Ethan Linck – all undergraduates from Univ. of Colorado at the time – were essential in developing a working Humboldt Core implementation, and were the first curators of analysts of inventory datasets.

*Funding* – We acknowledge support from Encyclopedia of Life and BioSync for funding to hold that meeting. Support for the work here was provided by a BioSync working group grant to WJ and RPG. We also acknowledge support from NSF DBI-1062148 and DBI-1535793 to RPG, NSF DBI-0735191 and DBI-1265383 to RLW, and NSF grants DBI-1262600, DEB-1441737, and NASA Grant NNX11AP72G to WJ.

## References

Anderson, R. P. 2013. A framework for using niche models to estimate impacts of climate change on species distributions. – *Ann. N. Y. Acad. Sci.* 1297: 8–28.  
 Bandrowski, A. et al. 2016. The ontology for biomedical investigations. – *PLoS One* 11: e0154556.

Beaman, R. S. and Cellinese, N. 2012. Mass digitization of scientific collections: new opportunities to transform the use of biological specimens and underwrite biodiversity science. – *ZooKeys* 209: 7–17.  
 Calabrese, J. M. et al. 2014. Stacking species distribution models and adjusting bias by linking them to macroecological models. – *Global Ecol. Biogeogr.* 23: 99–112.  
 Chao, A. et al. 2009. Sufficient sampling for asymptotic minimum species richness estimators. – *Ecology* 90: 1125–1133.  
 Corona, P. et al. 2011. Contribution of large-scale forest inventories to biodiversity assessment and monitoring. – *For. Ecol. Manage.* 262: 2061–2069.  
 Cox, S. J. D. 2013. An explicit OWL representation of ISO/OGC observations and measurements. – <<http://dl.acm.org/citation.cfm?id=2874543.2874544>>.  
 Dickinson, J. L. et al. 2012. The current state of citizen science as a tool for ecological research and public engagement. – *Front. Ecol. Environ.* 10: 291–297.  
 Dorazio, R. M. and Royle, J. A. 2005. Estimating size and composition of biological communities by modeling the occurrence of species. – *J. Am. Stat. Assoc.* 100: 389–398.  
 Dorazio, R. M. et al. 2006. Estimating species richness and accumulation by modeling species occurrence and detectability. – *Ecology* 87: 842–854.  
 Elith, J. and Leathwick, J. R. 2009. Species distribution models: ecological explanation and prediction across space and time. – *Annu. Rev. Ecol. Evol. Syst.* 40: 677–697.  
 Enke, N. et al. 2012. The user's view on biodiversity data sharing – investigating facts of acceptance and requirements to realize a sustainable use of research data. – *Ecol. Inform.* 11: 25–33.  
 Gotelli, N. J. and Colwell, R. K. 2001. Quantifying biodiversity: procedures and pitfalls in the measurement and comparison of species richness. – *Ecol. Lett.* 4: 379–391.  
 Gries, C. et al. 2014. Symbiota – a virtual platform for creating voucher-based biodiversity information communities. – *Biodivers. Data J.* 2: e1114.  
 Iknayan, K. J. et al. 2014. Detecting diversity: emerging methods to estimate species diversity. – *Trends Ecol. Evol.* 29: 97–106.  
 Isaac, N. J. B. and Pocock, M. J. O. 2015. Bias and information in biological records. – *Biol. J. Linn. Soc.* 115: 522–531.  
 Jarzyna, M. A. and Jetz, W. 2016. Detecting the multiple facets of biodiversity. – *Trends Ecol. Evol.* 31: 527–538.  
 Jetz, W. et al. 2012. Integrating biodiversity distribution knowledge: toward a global map of life. – *Trends Ecol. Evol.* 27: 151–159.  
 Jiménez-Valverde, A. and Lobo, J. M. 2007. Threshold criteria for conversion of probability of species presence to either-or presence-absence. – *Acta Oecol.* 31: 361–369.  
 Kelling, S. et al. 2009. Data-intensive science: a new paradigm for biodiversity studies. – *BioScience* 59: 613–620.  
 Kéry M. et al. 2010. Site-occupancy distribution modeling to correct population-trend estimates derived from opportunistic observations. – *Conserv. Biol.* 24: 1388–1397.  
 Kremen, C. 1994. Biological inventory using target taxa: a case study of the butterflies of Madagascar. – *Ecol. Appl.* 4: 407–422.  
 Lahoz-Monfort, J. J. et al. 2014. Imperfect detection impacts the performance of species distribution models. – *Global Ecol. Biogeogr.* 23: 504–515.  
 Leon-Cortes, J. L. et al. 1998. Assessing completeness of Mexican sphinx moth inventories through species accumulation functions. – *Divers. Distrib.* 4: 37–44.  
 Lobo, J. M. et al. 2010. The uncertain nature of absences and their importance in species distribution modelling. – *Ecography* 33: 103–114.  
 MacKenzie, D. I. 2006. Occupancy estimation and modeling: inferring patterns and dynamics of species occurrence. – Elsevier/Academic Press.

- MacKenzie, D. I. et al. 2002. Estimating site occupancy rates when the detection probabilities are less than one. – *Ecology* 83: 2248–2255.
- Madin, J. et al. 2007. An ontology for describing and synthesizing ecological observation data. – *Ecol. Inform.* 2: 279–296.
- Maes, D. et al. 2015. The use of opportunistic data for IUCN Red List assessments. – *Biol. J. Linn. Soc.* 115: 690–706.
- Meyer, C. et al. 2015. Global priorities for an effective information basis of biodiversity distributions. – *Nat. Commun.* 6: 8221.
- Nakamura, M. and Soberón, J. 2009. Use of approximate inference in an index of completeness of biological inventories. – *Conserv. Biol.* 23: 469–474.
- Newman, G. et al. 2012. The future of citizen science: emerging technologies and shifting paradigms. – *front. Ecol. Environ.* 10: 298–304.
- Pereira, H. M. et al. 2013. Ecology. Essential biodiversity variables. – *Science* 339: 277–278.
- Pescott, O. L. et al. 2015. Ecological monitoring with citizen science: the design and implementation of schemes for recording plants in Britain and Ireland. – *Biol. J. Linn. Soc.* 115: 505–521.
- Proença, V. L. J. et al. 2016. Global biodiversity monitoring: from data sources to essential biodiversity variables. – *Biol. Conserv.* in press.
- Robertson, T. et al. 2014. The GBIF integrated publishing toolkit: facilitating the efficient publishing of biodiversity data on the internet. – *PLoS One* 9: e102623.
- Sadoti, G. et al. 2013. Applying occupancy estimation and modelling to the analysis of atlas data. – *Divers. Distrib.* 19: 804–814.
- Schmeller, D. S. et al. 2015. Towards a global terrestrial species monitoring program. – *J. Nat. Conserv.* 25: 51–57.
- Szabo, J. K. et al. 2010. Regional avian species declines estimated from volunteer-collected long-term data using list length analysis. – *Ecol. Appl.* 20: 2157–2169.
- Tyre, A. J. et al. 2003. Improving precision and reducing bias in biological surveys: estimating false-negative error rates. – *Ecol. Appl.* 13: 1790–1801.
- van Strien, A. J. et al. 2013. Opportunistic citizen science data of animal species produce reliable estimates of distribution trends if analysed with occupancy models. – *J. Appl. Ecol.* 50: 1450–1458.
- Walls, R. L. et al. 2014a. Semantics in support of biodiversity knowledge discovery: an introduction to the biological collections ontology and related ontologies. – *PLoS One* 9: e89606.
- Walls, R. L. et al. 2014b. Meeting report: advancing practical applications of biodiversity ontologies. – *Stand. Genomic Sci.* 9: 17.
- Wieczorek, J. et al. 2012. Darwin Core: an evolving community-developed biodiversity data standard. – *PLoS One* 7: e29715.
- Wieczorek, J. et al. 2014. Meeting report: GBIF Hackathon-workshop on Darwin Core and sample data. – *Stand. Genomic Sci.* 9: 585–598.
- Wilkinson, M. D. et al. 2016. The FAIR guiding principles for scientific data management and stewardship. – *Sci. Data* 3: 160018.
- Wintle, B. A. et al. 2012. Designing occupancy surveys and interpreting non-detection when observations are imperfect. – *Divers. Distrib.* 18: 417–424.
- Wiser, S. K. 2016. Achievements and challenges in the integration, reuse and synthesis of vegetation plot data. – *J. Veg. Sci.* 27: 868–879.

Supplementary material (Appendix ECOG-02942 at <[www.ecography.org/appendix/ecog-02942](http://www.ecography.org/appendix/ecog-02942)>). Appendix 1–5.