**DATA and KNOWLEDGE INFRASTRUCTURE** ⟷ **ANALYTICS, KNOWLEDGE DISSEMINATION**
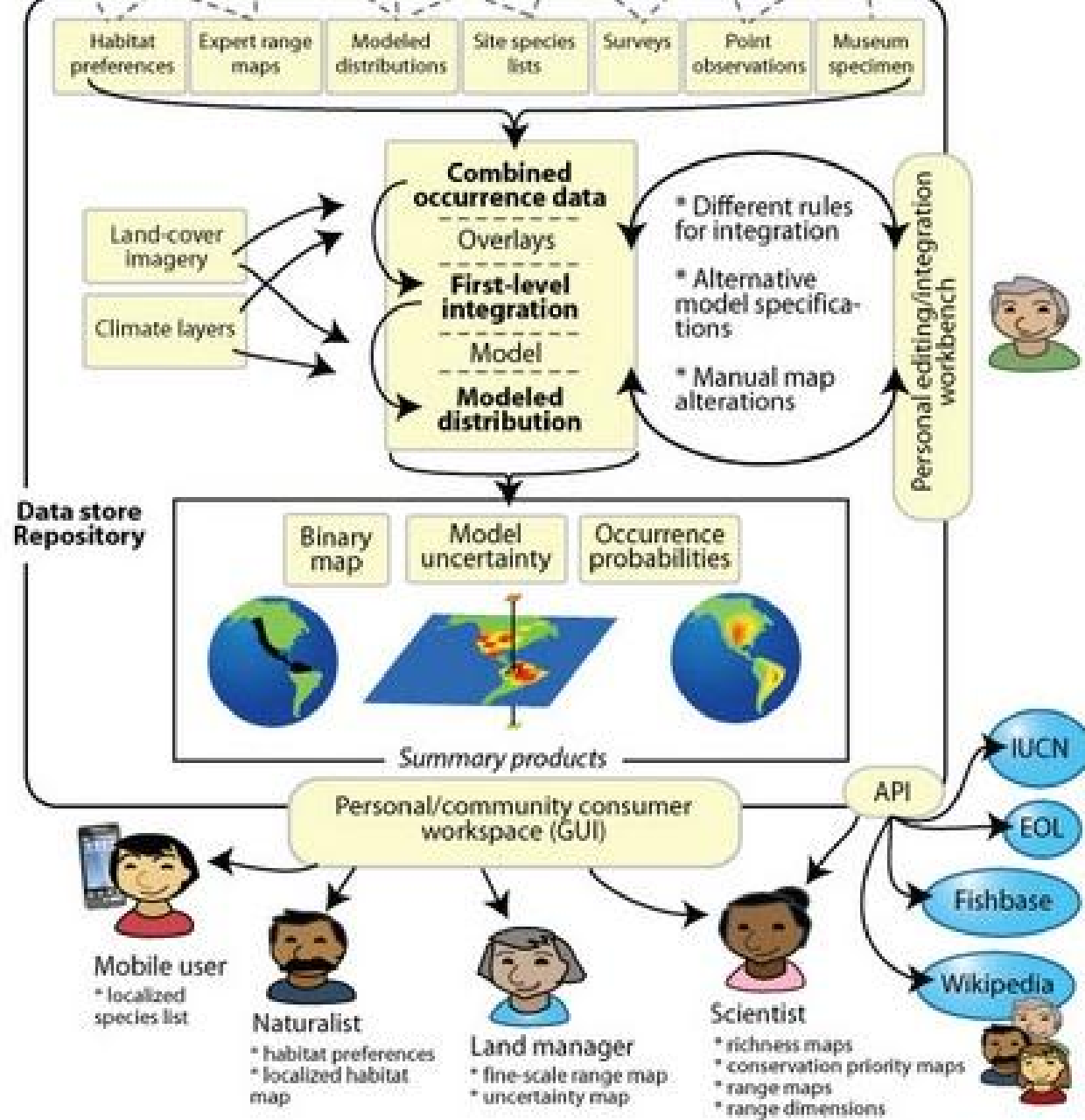
**Contributors from multiple sectors add to Map of Life**

**Consumers use Map of Life knowledge for societal needs**

# BUT TO PROPERLY USE THESE DATA WE NEED TO UNDERSTAND DIFFERENT SOURCES OF DATA AND HOW TO DESCRIBE THEM

| Description | Example | Contribut., Quality | Proto-cols | Effort Report | Source | Raw data | Temp-oral Scope | Geogr-aphic scope | Reporting basis | Complete &suited absence inference | Suited input occpy. Model |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Summary inventory (limited protocol & effort report) | Protected area species list | many, heterogen. | multi, unclear | Very limited | Literature | no | long (years) | Clear (often >1km) | Multi (observation, Photo, Lit,) | Yes | No |
| Summary inventory (some protocol and effort reporting) | Standardized area survey (e.g. atlas grid cell) | many, heterogen. | multiple single, clear | Possible | Literature, Project reporting | no | long (months, years) | Clear (often >1km) | Single (e.g. Observation) | yes | no |
| Single person/group inventory: observation | Standardized area survey (e.g. transect count) | single, high & vetted | single, clear | Yes | Project data | yes | short (hours, days) | Clear (m to 1km) | Single (e.g. Observation) | no | Potent-ially |
| Inventory following protocol: stationary trapping | Camera traps & more typical trappings | single, high & vetted | singe, some-what clear | Yes | Project data | yes | short (hours, days) | Very small (meters) | Single (e.g. Observation) | somewhat (over very small extent) | Yes |
| Inventory following protocol: active sampling campaigns | fish, zooplancton netting, algal sampling | single, high & vetted | single, clear | Yes | Project data | yes | usually short | Small (e.g. meters) | Single (e.g. Observation) | no | Yes |
| Full inventory following very defined protocol | CTFS forest plot, Revelle plots | single, high & vetted | single, clear | Perfect: full coverage | Project data | yes | short | V. small (e.g. meters) | Single (e.g. Observation) | yes (over very small extent) | NA |
| Inventory following loose protocol: citizen science observation | transect | single, heterogen., unvetted | single, clear | Yes | Project data | yes | short (hours, days) | Clear (m to 1km) | Single (e.g. Observation) | no | Potent-ially |

# A metadata schema for collating data from inventories

| Humboldt Core Version 1 | Area species checklists | | Geographically restricted surveys | | |
|---|---|---|---|---|---|
| | Gridded Atlas survey | Protected area species list | Transect count | Trapping and netting | CTFS forest, Revelle plots |
| **General dataset & identification terms** | inventory perfomed by; dataset name, identifier, publisher, licensing, rights holders; metadata recorded by; citation reference and id; taxa identifier 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/lon 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; Compiled data Y/N & type; abundances and/or absences reported? Absence list | | Inventory type, protocol name, detail, citation , reference, abundances reported Y/N &cap; absences reported? | | |
| **Completeness & effort terms** | Completeness reported and how; Inferred taxonomic completeness Upper/lower bound and how. | | Effort reporting & lower/upper bound and granular breakdown; effort method; Vouchers or samples taken and how? | | |

Team of Boulder and Yale developers and students assembled metadata for (so far) 143 area checklists and collated information about area checklists characteristics

| Humbolt Core Term | Possible Values | Percents |
|---|---|---|
| Compilation effort | Low, medium, high, na | Low – 38.7%, Medium – 6.2%, High – 8.5%, n.a – 45.7% |
| Abundances reported | Yes/No | Yes-57.5%  No-42.6 |
| Absences reported | Yes/No | Yes-41.8 , No-58.2% |
| Completeness assessment | Scale in 25 percent increments from 0-100 | 50-100% complete – 30.3%, 75-100% - 27.9%, other- 41.8% |

# A tool for the long-tail data on YOUR computer

# Point (and list) Uploader Key Points

- The point uploader is the first of many upload tools

- Metadata about datasets provides critical ownership and rights data

- As well as essential content for better use of the datasets e.g. probabilistic assessment of absences

- Data provided may be kept private for use

- Or made available for broader use, curation, improvement

- Metadata also provides value for downstream modeling (more expert versus novice data)

## TAXONOMY "TRIVIAL BUT TERRIFYING" ISSUES

**Whether citizen science data or data from museum records, data cleaning before import is important and provides value for providers and consumers**

Based on a gold standard, hand vetted set of 500 museum digitized label data in VertNet:

- 7.8% of scientific names could not be resolved at all
- 32% of names are unaccepted but could be resolved to accepted names
- 2.6% are misspelled and unaccepted names
- 10% are misspellings of accepted names
- 47.6% are current accepted names

Take home:   Huge issues with ingested data, requiring novel solutions

# The non-trival problem-
## tracking how naming **meanings** change

Species X    Species Y — Lump → **Species X**

Initial circumscription      New circumscription

Species X / Sp. X ssp1 — Split → **Species X**   Species Y (ssp1)

Initial circumscription      New circumscription

Geographic range outcome

# This is not a problem of "old records" or just "some groups"



**Number of species definitions created in each decade between 1885 and 2014**

(Y axis: Number of species definitions; X axis: Decade — 1885, 1905, 1925, 1945, 1965, 1985, 2005)

Taxonomic effort
In **North American birds** 1885-2015

Primary descriptions

Redescriptions that change taxon concepts

\* Based on AOU checklist (a conservative assessment)

# Big Challenges working with Names –
## Reconciling Names on Ingest
**(for those resources where name validation has been inconsistent/problematic)**

Scrape list to text

Import list to resolver

Check name is accepted?
If no, get accepted name

Check names in MOL authority files

Use new algorithmic approaches to find misspelling, revalidate, store location geom. & validated names

**Birds of El Malpais**
National Monument & National Conservation Area

Binding Guide to El Malpais

By Ken Jones
Illustrations by Zachary Zdinak

eol
Encyclopedia of Life

GBIF

ITIS

| Feature | Total |
|---|---|
| **Total records** | **200183168 (100.00%)** |
| (A) Geospatial issues | |
| Coordinates equal to zero | 1456654 (0.73%) |
| Impossible coordinates | 10117 (0.01%) |
| Low precision | 16252100 (8.12%) |
| Out of the specified country | 14040820 (7.01%) |
| Transposed coordinates | 322734 (0.16%) |
| Negated Latitude | 272829 (0.14%) |
| Negated Longitude | 383919 (0.19%) |
| | |
| (B) Spatio-taxonomic issues | |
| Inside range map | 146877631 (73.37%) |
| Less than 55Km | 18408468 (9.20%) |
| 55-111Km | 2228994 (1.11%) |
| 111- 555Km | 3783218 (1.89%) |
| More than 555Km | 5417136 (2.71%) |
| Without range map assessment | 23467721 (11.72%) |
| Without RM assessment - taxon issues | 12912456 (6.45%) |

A global assessment of terrestrial vertebrates using GBIF records. Data from Otegui and Guralnick

**CLEANING IS NOT A ONE STEP PROCESS**
**... It is a constant process of further refining ...**

VertNet beta    ×
portal.vertnet.org/o/royal-ontario-museum/mammalogy-collection-royal-ontario-museum?id=0f7e0451-ec0a-4ace-aa96-1770064bccd8

Apps  M Gmail  Feedly  Projects  VertNet  VNComm  IPT  iDigBio  Cal  Georef  Portals  LatLongX  B VertNet  Henderson Field Not  Citizen Archivist Da  CartoDB + Time  »

VertNet beta    Search    Publishers    About    Feedback

dbloom@github    logout

# ROM Mammals 6

Summary    Data Rights    Darwin Co

Submit data issue

## Submit data issue ×

Help data publishers track data issues using GitHub!

Sauromys not in Aves

This record lists Sauromys as Class = Aves.  Bats are mammals - even in Canada.

Submit data issue

| Term | Value |
| --- | --- |
| InstitutionCode | Royal Ontario Museum: ROM |
| CollectionCode | Mammals |
| CatalogNumber | 65044 |
| BasisOfRecord | PreservedSpecimen |
| Year | 1965 |
| Country | BOTSWANA |

🔒 GitHub, Inc. [US] https://github.com/mvz-vertnet/mvz-herp/issues

Apps | M Gmail | Feedly | Projects | VertNet | VNComm | IPT | iDigBio | Cal | Georef | Portals | LatLongX | B VertNet | Henderson Field No | Citizen Archivist Das | CartoDB + Time | »

This repository ▾ | Search or type a command | ? | Explore | Gist | Blog | Help | dbloom

PUBLIC | mvz-vertnet / **mvz-herp**

Unwatch ▾ | 8 | ★ Star | 0 | Fork | 1

Browse Issues | Milestones

New Issue

Everyone's Issues | 5

5 Open | 0 Closed | Sort: Newest ▾

Assigned to you | 0

Created by you | 0

Mentioning you | 0

Close | Label ▾ | Assignee ▾ | Milestone ▾

[MVZ Herp 195816] Batrachoseps wrighti - Test: Look Batrachoseps in Oregon!
Opened by atrox10 19 days ago
#5

[MVZ Herp 195816] Batrachoseps wrighti - Test: Look Batrachoseps in Oregon!
Opened by atrox10 19 days ago
#4

[MVZ Herp 195816] Batrachoseps wrighti - Test: Look Batrachoseps in Oregon!
Opened by atrox10 19 days ago
#3

[MVZ Herp 65979] Hyla eximia - see if CArol gets this
Opened by atrox10 a month ago
#2

[MVZ Herp 15073] Elgaria coerulea - wow it's in oregon
Opened by mkoo a month ago
#1

No milestone selected ⚙▾

Labels

■ bug | 0
▌duplicate | 0
■ enhancement | 0
▌invalid | 0
■ question | 0
wontfix | 0

Keyboard shortcuts available ⌨

Manage Labels

New label

New label name

- Reintegration of disparate data critical

  (but so is improving those data)

- The data and communities assembling data are

  *highly* heterogeneous and disconnected

- The data sciences components are *not trivial*.

- Map of Life provides tools for **ALL** to provision data, metadata and provide innovative tools to help curate & improve it

- To better serve needs for monitoring and assessment