# Data capture: aka data digitization
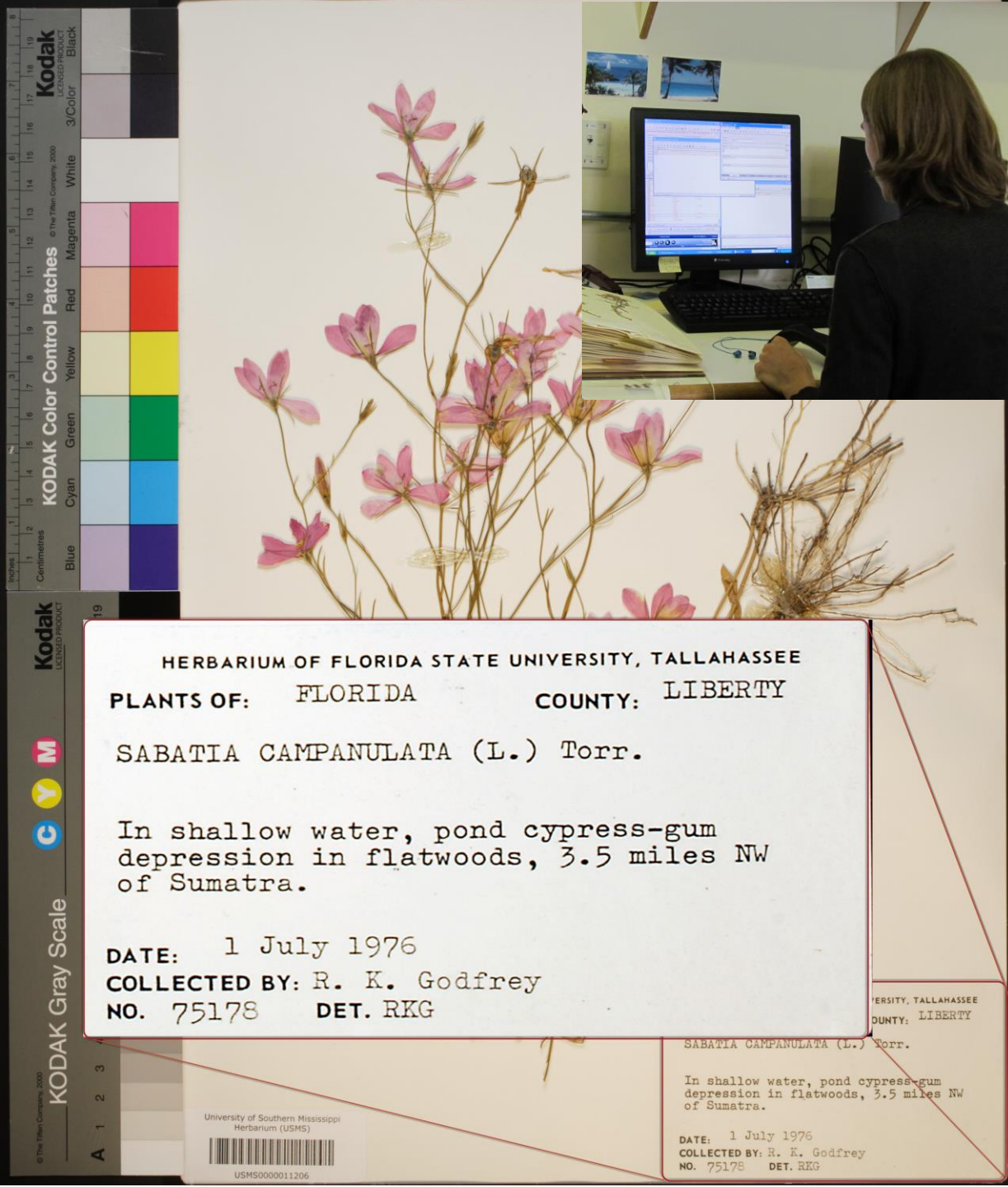
## some best practices, options, lessons learned

Deborah Paul
iDigBio, Florida State University
Bristol UK Swans Practical Digitisation Workshop
9 March 2018
@idbdeb @iDigBio

**Goals of data capture**

- Read and transcribe written materials
- Move accurate data into database

HERBARIUM OF FLORIDA STATE UNIVERSITY, TALLAHASSEE

PLANTS OF: FLORIDA    COUNTY: LIBERTY

SABATIA CAMPANULATA (L.) Torr.

In shallow water, pond cypress-gum depression in flatwoods, 3.5 miles NW of Sumatra.

DATE: 1 July 1976
COLLECTED BY: R. K. Godfrey
NO. 75178    DET. RKG

Occurrence Data

Long Form <<

Collector ?    Number ?   Date ?
R. K. Godfrey    75178    1976-07-01    Dupes?
☐ Auto search

Associated Collectors ?    Verbatim Date ?
1 July 1976

Exsiccati Title    Number

Scientific Name ?
Sabatia campanulata (L.) Torr.

Country    State/Province    County
United States    Florida    Liberty

Locality
3.5 miles NW of Sumatra.

Latitude    Longitude    Uncertainty ?    Verbatim Coordinates    Tools
<<

Elevation in Meters    Verbatim Elevation
-    <<

Habitat
In shallow water, pond cypress-gum depression in flatwoods

Substrate

Notes

Save Edits
Status Auto-Set: Pending Review

# Data Capture Challenges

- ink
- typed
- pencil
- printed
- stacked
- handwritten
- uneven lines
- colored paper
- non-planar surfaces
- non-standard terms
- non-standard formats

# Mobilization

**other data formats needing capture and standardization in order to share**

- spreadsheets
- log books
- field notes
- other derivative objects
- storage formats

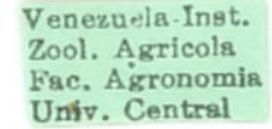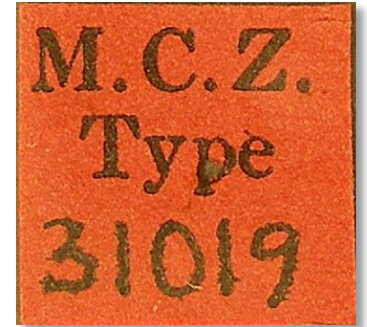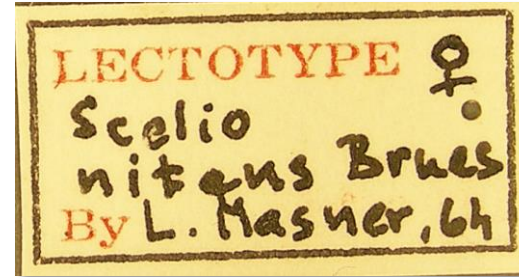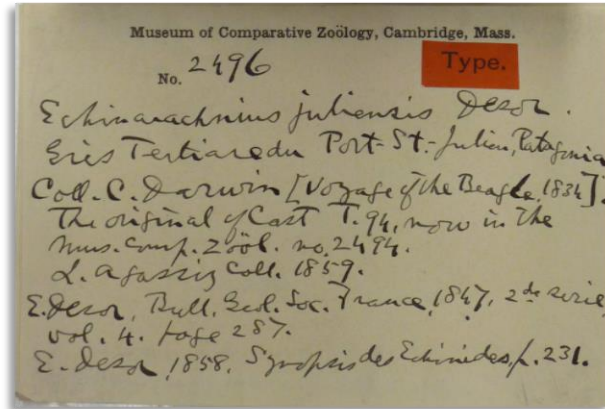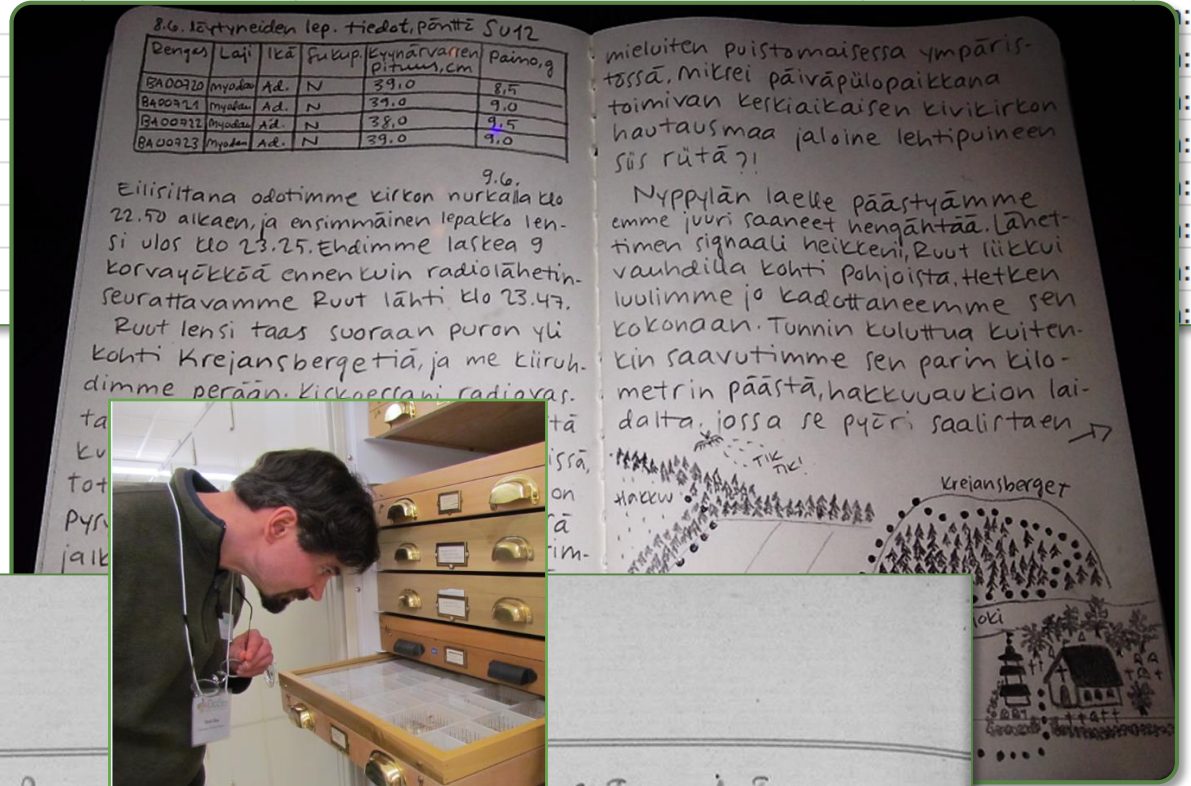# Extract and Derive

- Geolocation
- Phenology
- Habitat
- Ecology
- Morphology
- Stratigraphy
- DNA…

HERBARIUM OF FLORIDA STATE UNIVERSITY, TALLAHASSEE

PLANTS OF: FLORIDA    COUNTY: LIBERTY

SABATIA CAMPANULATA (L.) Torr.

In shallow water, pond cypress-gum depression in flatwoods, 3.5 miles NW of Sumatra.

DATE: 1 July 1976
COLLECTED BY: R. K. Godfrey    DET. RKG
NO. 75178

1 of 3
Epistenia coeruleata, Westwood, 1832, 1991-05-14
1.00 mm

2 of 3
Epistenia coeruleata, Westwood, 1832, 1991-05-14
2.0 mm

3 of 3
Epistenia coeruleata, Westwood, 1832, 1991-05-14
1.00 mm

GEOLocate Web Application
GEOLocate
Workbench    Results
Select File    or load an existing file using a retrieval code:    Load

# Data Capture: what to consider?

- data from image or data from label

- identifier for the object
  - local to global
  - never reuse

- how much data to capture?
  - all or some?

- is there useful existing digitized data?
  - taxonomy, geography, collector names

- do you have the database fields you need?
  - where to put the data

### Module 11: Data Capture

...nderlying focus of the steps throughout these digitization modules is to encourage institutions to follow an object to image to data workflow through which all specimens are first imaged and data recorded from these images. Nevertheless, some institutions choose, for various justifiable reasons, to pursue a specimen to data workflow and we try to accommodate both approaches below.

| Task ID | Task Description | Explanations and Comments | Resources |
|---------|------------------|---------------------------|-----------|
| T1 | Perform any preparatory steps. | Determine application to be used for data capt... considera... (especial... and proje... informati... requiren... | Data entry application. ...s. ...al data. ...software or ...integrated ...entry ...cation. |

Guidelines for: tasks, resources, and decisions involved in data capture

# Data Capture: what to consider? Part 2

- ## transcription issues
  - parsing (what goes in which field), implicit values
  - text missing from authority file

- ## data quality checks
  - transcription errors, erroneous information on labels
  - human and automated checks

- ## written protocols
  - iterative improvements, updates when equipment or software changes

iDigBio DROID Working Group product

**Module 11: Data Capture**

The underlying focus of the steps throughout these digitization modules is to encourage institutions to follow an object to image to data workflow through which all specimens are first imaged and data recorded from these images. Nevertheless, some institutions choose, for various justifiable reasons, to pursue a specimen to data workflow and we try to accommodate both approaches below.

| Task ID | Task Description |
| --- | --- |
| T1 | Perform any prepara... steps. |

*Guidelines for: tasks, resources, and decisions involved in data capture*

# Data quality: an issue at many levels



From a @HydraInABox interview: "People will put anything and their dog in the date field. It's absolutely astonishing."
— Hannah Frost @feefifofannah





196 Countries in the world, but 1100 distinct values in the country field

# Data cleaning

country=
"united kingdom"

Or does it?



Country: united kingdom
☐ Present   ☐ Missing

Top 10 Taxa
- Rallus recessus
- Grus latipes
- Anolis conspersus
- Porzana piercei
- Rallus ibycus
- Undet. productida
- Rhynchotreta cuneata
- Braconidae
- Gastropoda
- Meristina tumida
- other

Leaflet | Map data © OpenStreetMap

# Data Quality:
# Grooming and tics (**adapted from Joanna McCaffrey at iDigBio**)
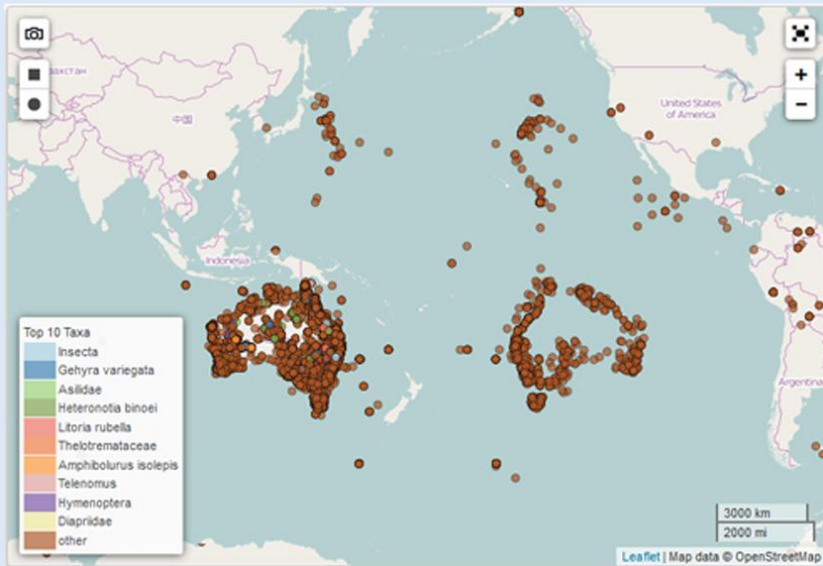
Your dataset **is no longer just for making labels**, there are other considerations for being digital, and out in the wild:

1. Put dates in ISO 8601 format, i.e., YYYY-MM-DD, e.g., 2015-09-17
2. Parse out scientific name
3. Conversely, put the piece parts into a scientific name
4. Provide as much higher taxonomy as your feel comfortable with, fill in tribe, sub+super family, kingdom, division, class, order) get out of 'family' land.
5. Make sure lat and lon coordinates are in degrees decimal, and no N, S, E, W
6. Do not export '0' in fields to represent no value, e.g., lat or lon
7. Put elevation in METERS units in the elevation field without the units (e.g., the fields dwc:minimumElevationInMeters and dwc:maximumElevationInMeters already assume the numeric values are in meters, so there no need to include the units with the data)
8. And not to get too esoteric, do not use un-escaped newline characters or embedded tabs
9. Watch out for diacritics (à á â ã ä å, save in UTF-8)

# iDigBio Data Quality (DQ) Flags enhance Digitization Workflows

Example: spot and fix georeferencing issues.



| Flag |
|------|
| idigbio_isocountrycode_added ⓘ |
| dwc_continent_added ⓘ |
| dwc_country_replaced ⓘ |
| geopoint_datum_missing ⓘ |
| dwc_class_replaced ⓘ |
| dwc_phylum_replaced ⓘ |
| dwc_order_replaced ⓘ |
| geopoint_low_precision ⓘ |
| rev_geocode_eez ⓘ |
| dwc_stateprovince_replaced ⓘ |
| rev_geocode_mismatch ⓘ |
| dwc_order_added ⓘ |
| datecollected_bounds ⓘ |
| dwc_class_added ⓘ |
| dwc_kingdom_added ⓘ |
| dwc_phylum_added ⓘ |
| dwc_country_added ⓘ |
| rev_geocode_corrected ⓘ |
| rev_geocode_lon_sign ⓘ |

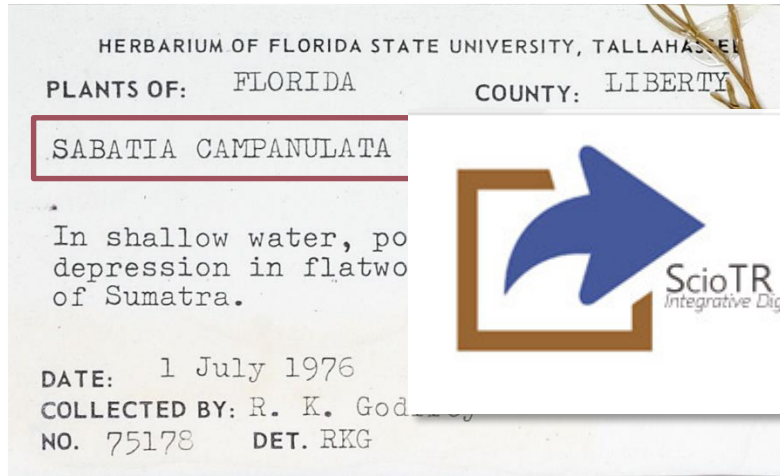- Value of data aggregation

- Planning for feedback integration

14

# Some Newer Data Capture Options

- Voice technology

- Touch screen technology

- OCR + ML + NLP

HERBARIUM OF FLORIDA STATE UNIVERSITY, TALLAHASSEE

PLANTS OF: FLORIDA    COUNTY: LIBERTY

SABATIA CAMPANULATA

In shallow water, po
depression in flatwo
of Sumatra.

DATE:    1 July 1976
COLLECTED BY: R. K. Godfrey
NO. 75178    DET. RKG

ScioTR
Integrative Digitization

Occurrence (Verbatim)

Verbatim Institution
HERBARIUM OF FLORIDA STATE UNIVERSITY,

Catalog Number
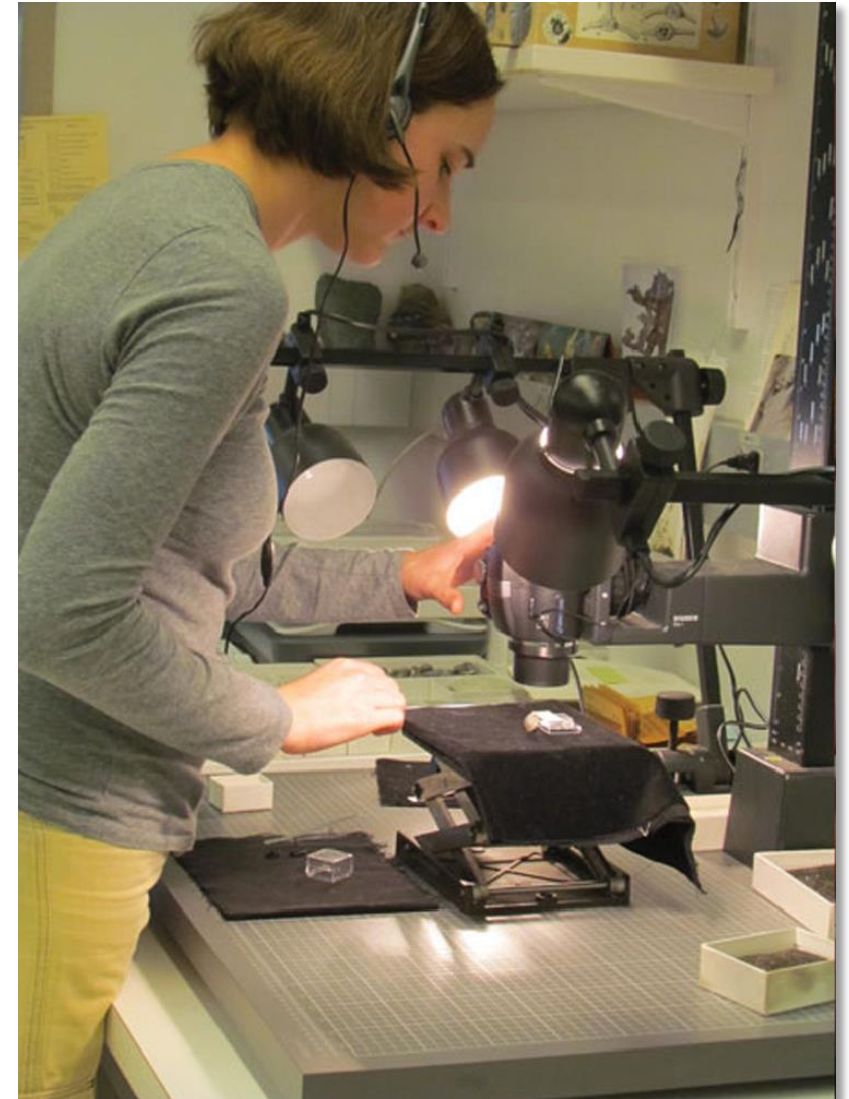75178

Verbatim Scientific Name
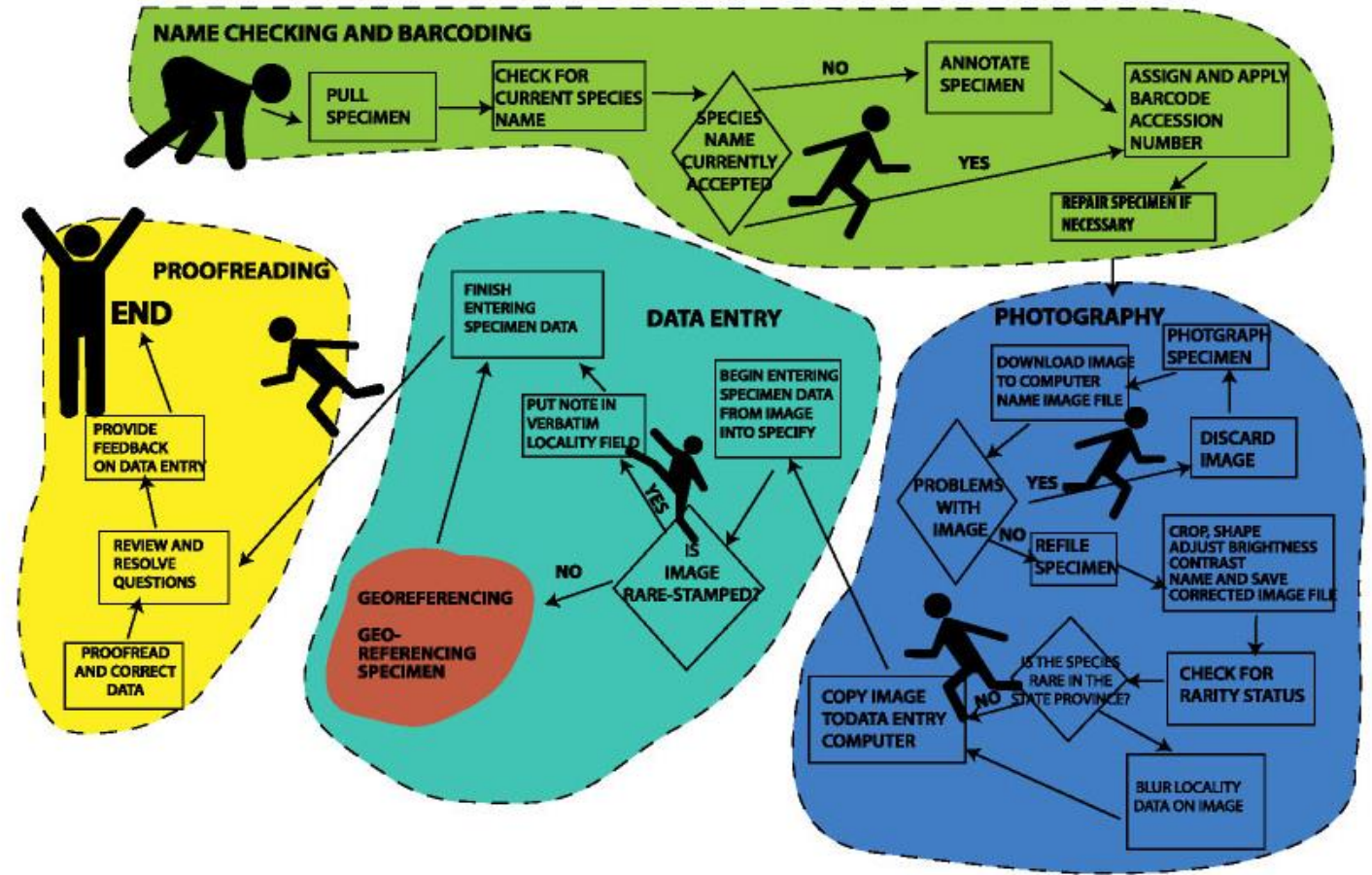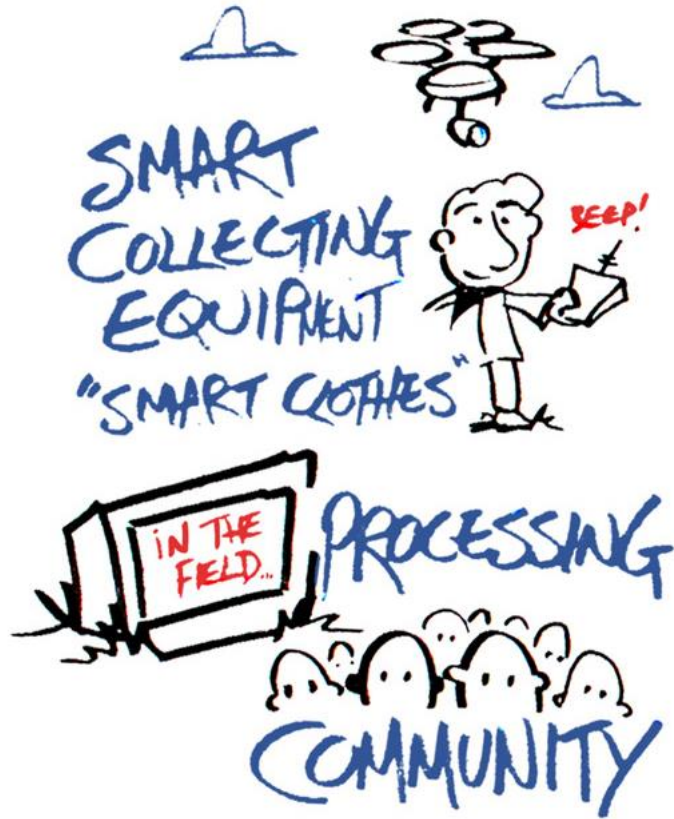SABATIA CAMPANULATA (L.) Torr.

Recorded By...

# Do-er happiness for
*productivity and data quality*

# From the field ➡️ *born digital*

By Dorothy Allard

17

# Thanks – questions?

**www.idigbio.org**

facebook.com/iDigBio

twitter.com/iDigBio

vimeo.com/idigbio

idigbio.org/rss-feed.xml

webcal://www.idigbio.org/events-calendar/export.ics