# Data capture: issues, best practices, options, lessons learned
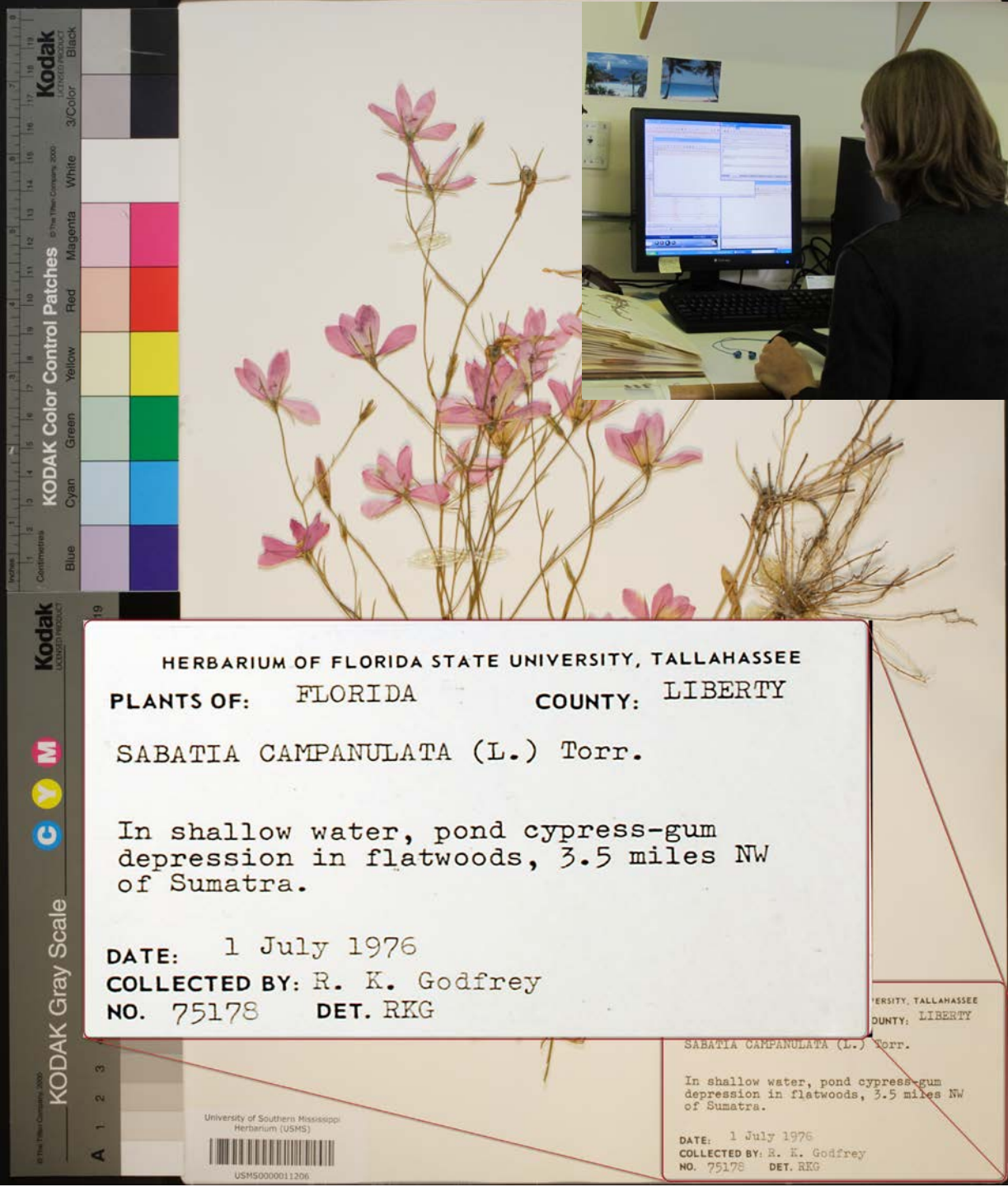
Deborah Paul
iDigBio, Florida State University
DigIn iDigBio Marine Invertebrate Digitisation Workshop
4-5 February 2019
@idbdeb @iDigBio

# Goals of data capture

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

- Topics
  - parsing
  - label variation, ancillary data, derived data
  - wet collections digitization tasks and resources
    - data capture is one part
  - data quality – make a plan
  - data organization and doer happiness
  - future data collection

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

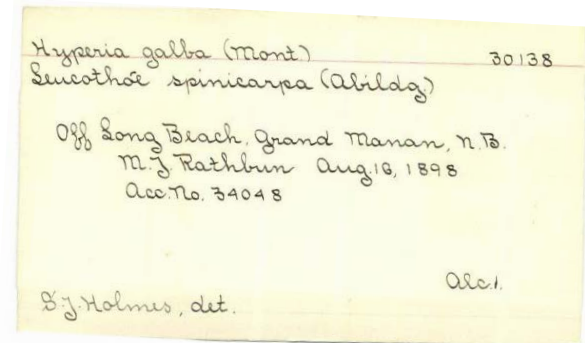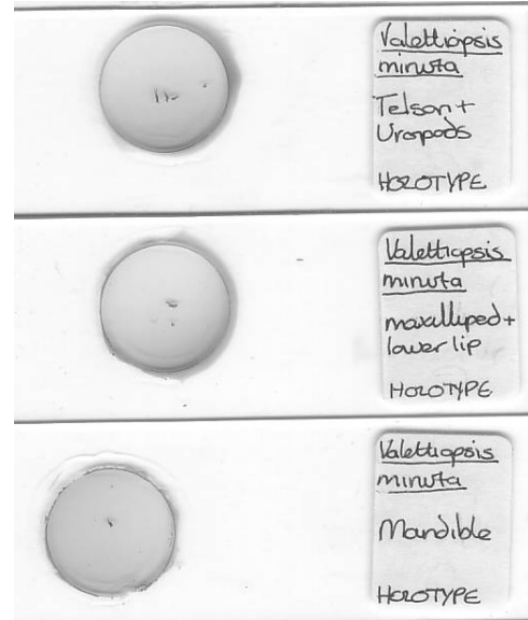University of Southern Mississippi Herbarium (USMS)
USMS0000011206

Occurrence Data

Long Form <<

Collector ?          Number ?   Date ?          Dupes?
R. K. Godfrey        75178      1976-07-01       ☐ 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
- **wet**
- typed
- pencil
- fragile
- printed
- curved
- stacked
- obscured
- handwritten
- uneven lines
- colored paper
- non-planar surfaces
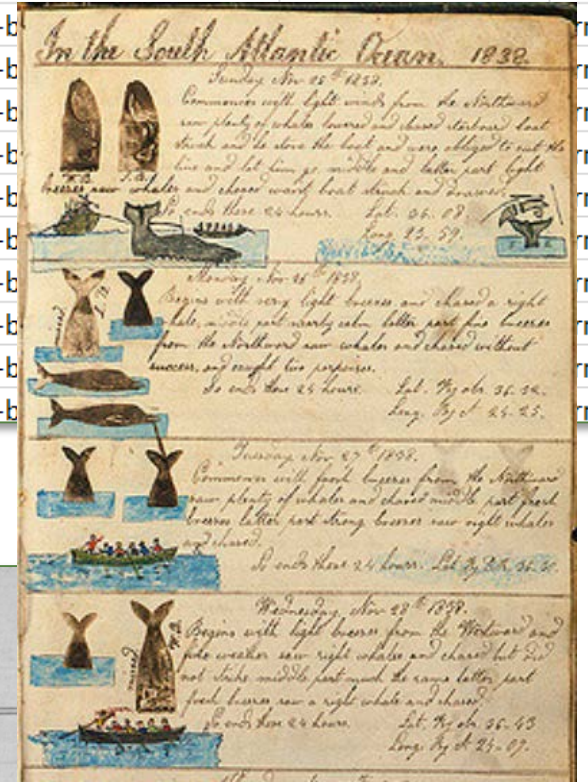- non-standard terms
- non-standard formats

# 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 / Depth
- DNA…

# Wet Collections Tasks: what to consider?

## Things in Spirits in Jars

- Module 0 Pre-digitization Curation Tasks

- Module 1A Imaging Ledgers, Cards, Field Notes

- Module 1B Imaging Specimen Labels

- Module 1C Specimen Imaging

- Module 1D Image Processing

- Module 1E Phototank Immersion Imaging Setup

- Module 1F Phototank Immersion Specimen Prep

- Module 1G Phototank Immersion Image Capture

- Module 1F Phototank Immersion Image Processing

- Module 2 Data Entry

- Module 3 Proactive Digitization

**Module 2: Data entry, Fluid-preserved**

**Module 2: Data Entry from Ledger, Card, Label, or Catalog Images**

| Task ID | Task Name | Explanations and Comments | Resources |
|---------|-----------|---------------------------|-----------|
| T1 | Navigate to the file folder in which image files are stored and load image of label to be transcribed. | | |
| T2 | Create verbatim data record from image file. | Focus of da... preservatio... when inte... In many i... governed... | ...an, ...ogy, ...ds-based ...se ...ement ...that ...es ...mic |

Guidelines for: tasks, resources, and decisions involved in digitizing wet collections

7

# 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
  - verbatim and / or interpreted
- do you have to move the specimens?
  - can you take the data capture to the shelves?

**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 | Explanations and Comments | Resources |
|---|---|---|---|
| T1 | Perform any preparatory steps | Determine application to be used for data capt... | Data entry application |

**Workflow Detail: Data Capture from Specimen**

**Module 4B: Data Capture from Specimen**

| Task ID | Task Name | | ources |
|---|---|---|---|
| T1 | Select and transport drawer(s) to database station. | | |
| T2 | Select specimen for data entry. | | |
| T3 | Space labels on pin to make the data visible or remove label(s) from pins and arrange in pin order, top to bottom to facilitate label re-insertion. | | |

*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 preparation steps. |

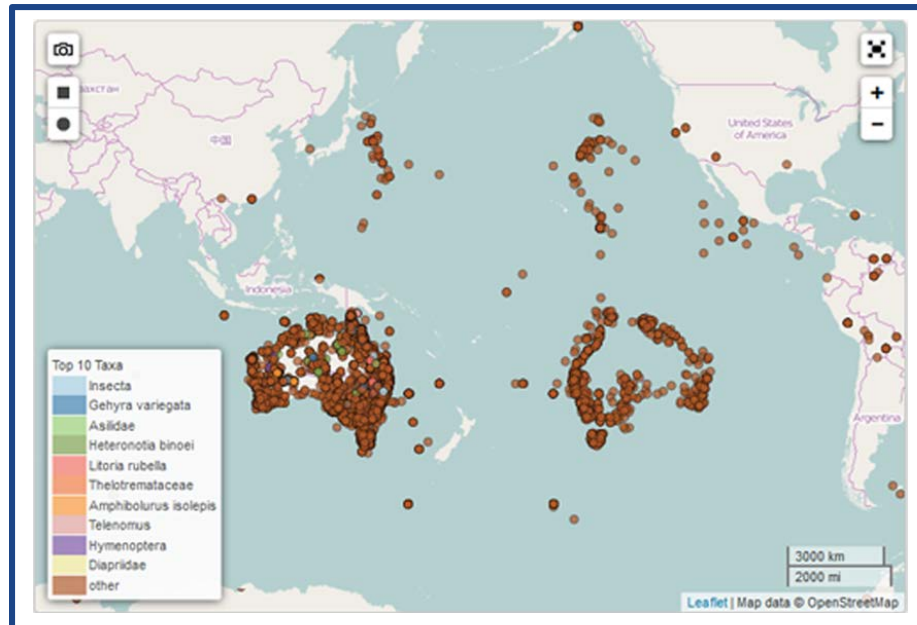*Guidelines for: tasks, resources, and decisions involved in data capture*

# Data quality: an issue at many levels

**Hannah Frost**
@feefifofannah

*Following*

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



Top 10 Taxa
- Insecta
- Gehyra variegata
- Asilidae
- Heteronotia binoei
- Litoria rubella
- Thelotremataceae
- Amphibolurus isolepis
- Telenomus
- Hymenoptera
- Diapriidae
- other

Country: united k

- united kindgom
- united king
- **united kingdom**
- united kingdom (england)
- united kingdom (scotland)
- united kingdom (wales)
- united kingdom [?]
- united kingdom of great b
- united kingdom?

196 Countries in the world, but 1100 distinct values in the country field
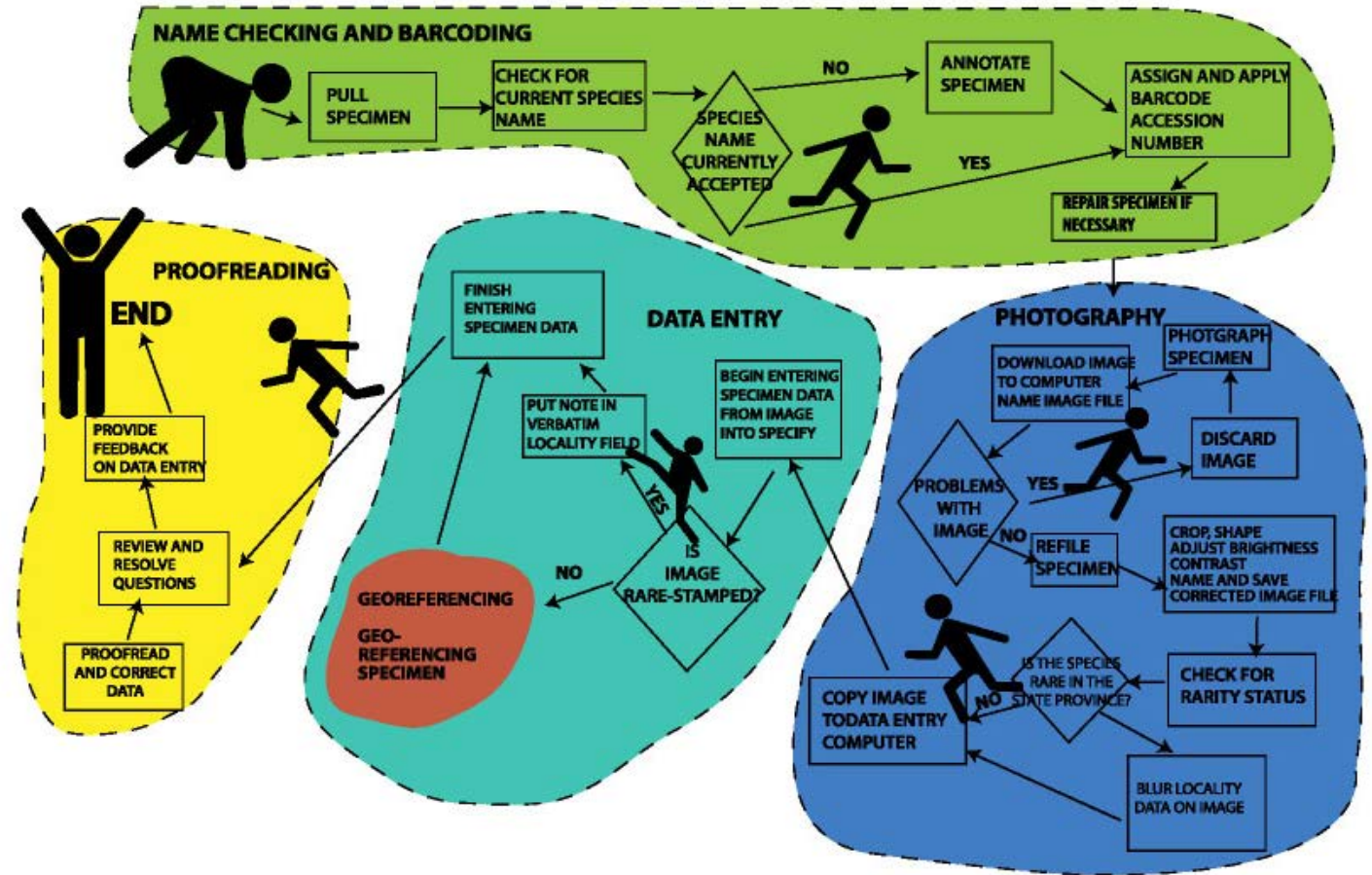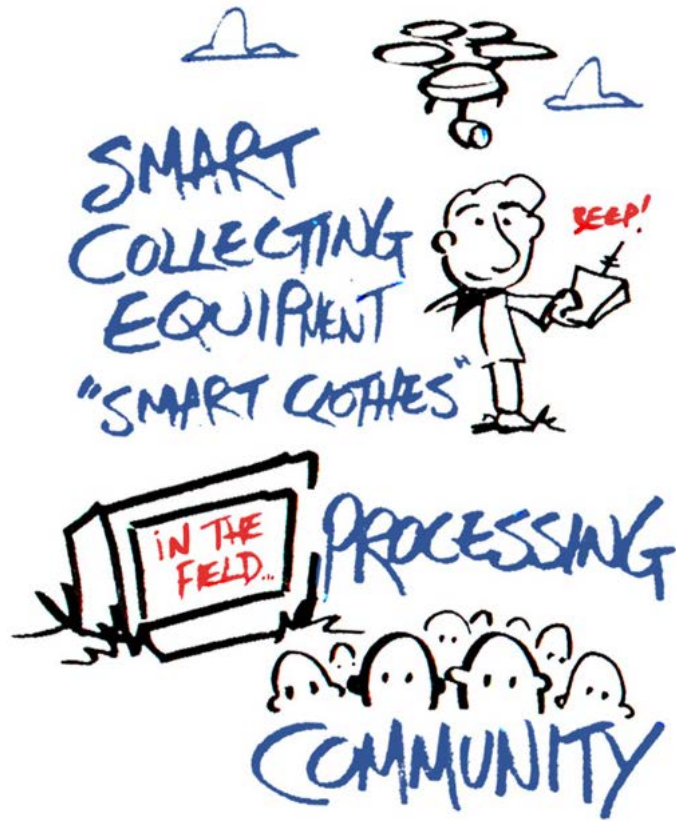
# Do-er happiness for
## *productivity and data quality*

# From the field ➡ *born digital*

# Thanks – insights?

**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