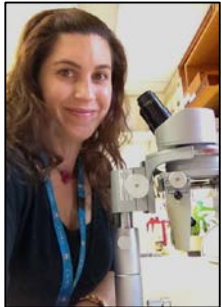# Implementing Collections Data Quality (DQ) Feedback
## *a survey and your community experience stories to shed some light into the data integration abyss*
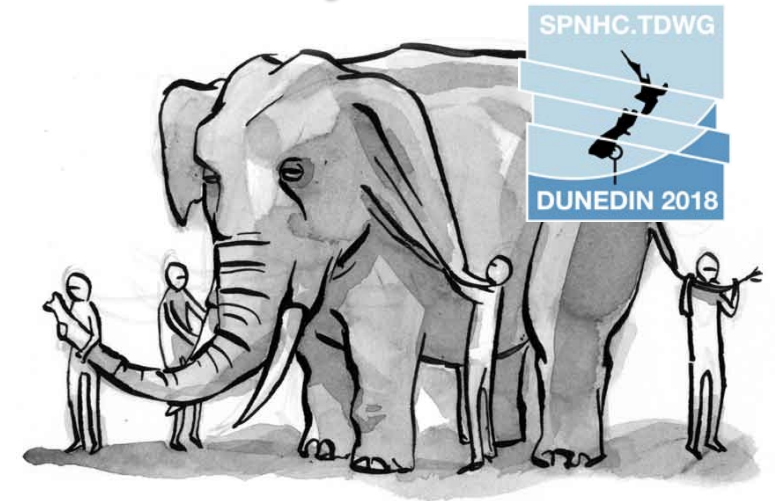
Deborah Paul, iDigBio, Florida State University
Nicole Fisher, CSIRO
@SPNHC-TDWGNZ Thursday 30 August 2018
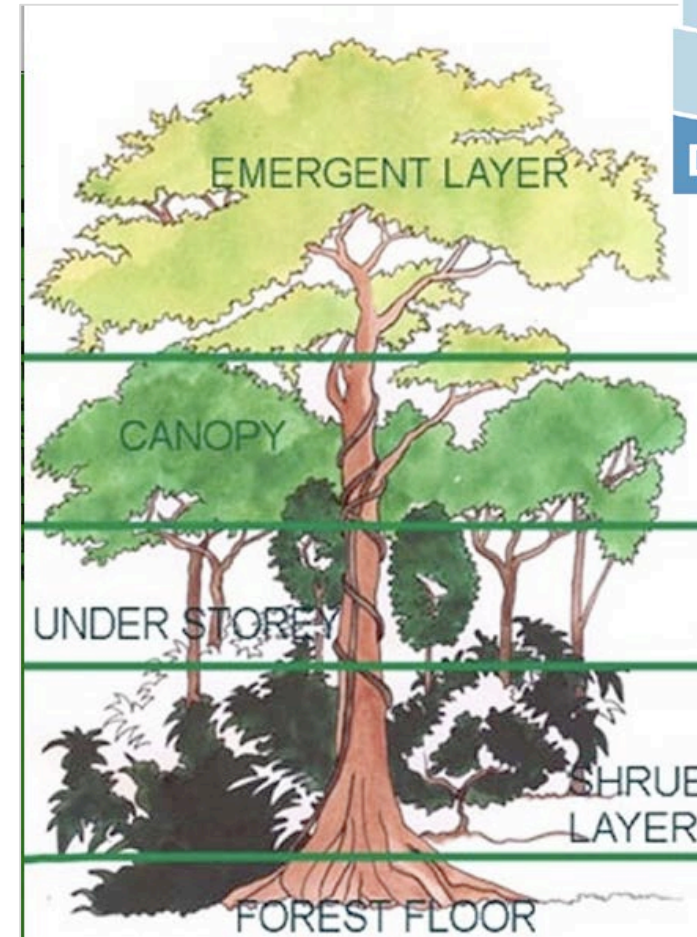@idbdeb @fisher_anic http://bit.ly/spnhcdq2018
Collections and Data in an Uncertain World

# Topics

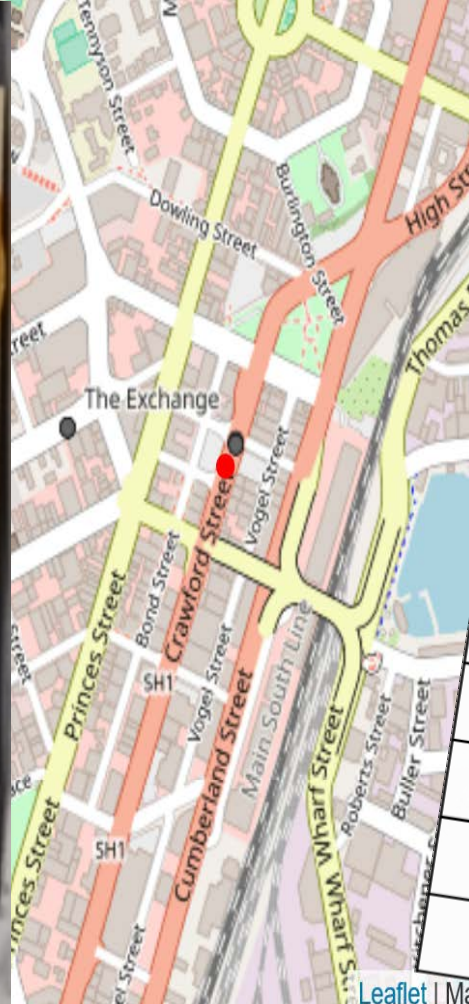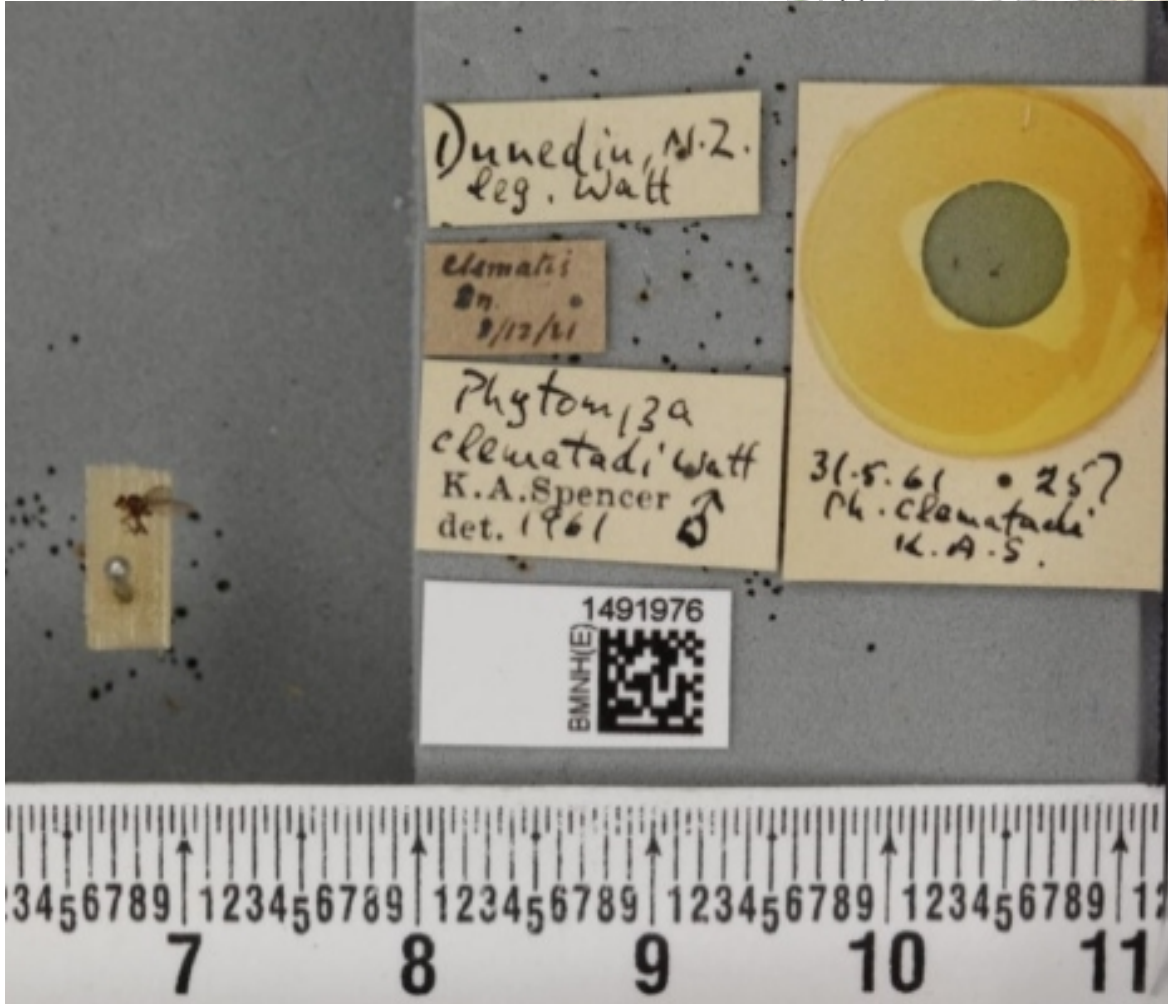- ## What we did and why
  - Survey, SIG, Report
- ## Highlights of Some Results
- ## What next?
  - Our speakers
  - SIG tomorrow 11-1230
    - your abyss?
  - Survey report *in progress*
  - Darwin Core Hour *follow-up*

From the Collection of: NHMUK

| From the Collection of: NHMUK | |
| --- | --- |
| Collector: Watt | |
| *Phytomyza clematadi* Watt | |
| Found on: *Clematis* | |
| Date Collected: 1921-12-09 | |
| Date Due | Borrower's Name |
| 1961 | det. K. A. Spencer ♂ |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |

Australasia; New Zealand; Otago; Dunedin, -45.8788, 170.5028

Biodiversity Data Pipeline: where are (some of) the bottlenecks (now)?
What enables use of data quality feedback and supports data integration?

**Collection Catalog**

**Applications**

**Researcher**

from the field and / or from the specimen

**Manage Data**

encoding

standardization

outlier detection

taxonomy

duplicate detection

annotation    georeferencing

**Possibilities**

species ranges, outlier discovery, new species, gaps in collecting, traits, relationships, predictive niche models, collector maps,...

# Survey to evaluate community experience
*integrating DQ Feedback*

| Left | Term | Right |
|---|---|---|
| believe that using the ALA's | data | quality assertions is an invaluable |
| for validating and improving our | data | . We would like to prioritise |
| actually deciding which of the | data | quality assertions the ALA adds |
| and gave us feedback on | data | flow and quality of coordinates |
| do a thorough job of | data | evaluation and cleaning Some the |
| different ways to standardize the | data | , of which there are many |
| emails in order to clean | data | . These messages are minimal. Once |
| be more of a priority. | data | quality improvement posterior to data |
| Data quality improvement posterior to | data | entry has not been a |
| taxonomies). I have found some | data | entry errors such as GPS |
| the benefit of cleaning their | data | or do not know to |



Data Corrected    Data Use    Raw

This table shows any data corrections that were per
represents the correction performed. The last two co
the data quality flags and their descriptions can be fo
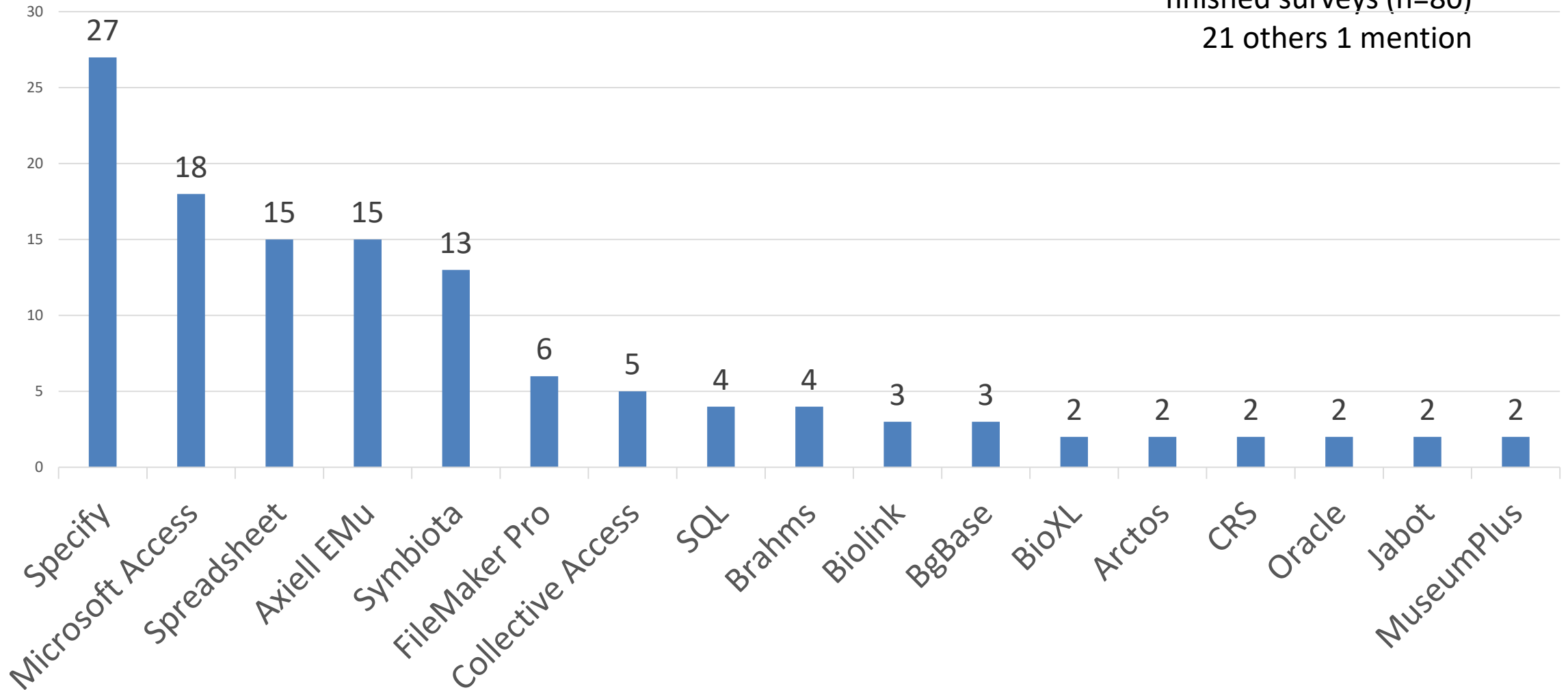this recordset.

**Flag**

dwc_datasetid_added ⓘ

dwc_parentnameusageid_added ⓘ

dwc_taxonid_added ⓘ

dwc_taxonomicstatus_added ⓘ

gbif_canonicalname_added ⓘ

gbif_genericname_added ⓘ

gbif_taxon_corrected ⓘ

dwc_taxonrank_added ⓘ

dwc_kingdom_added ⓘ

idigbio_isocountrycode_added ⓘ

gbif_reference_added ⓘ

gbif_vernacularname_added ⓘ

dwc_phylum_added ⓘ

dwc_multimedia_added ⓘ

# Collection Management Software

n=104 responders
software mentioned > 1 time
finished surveys (n=80)
21 others 1 mention

Bar chart values:
- Specify: 27
- Microsoft Access: 18
- Spreadsheet: 15
- Axiell EMu: 15
- Symbiota: 13
- FileMaker Pro: 6
- Collective Access: 5
- SQL: 4
- Brahms: 4
- Biolink: 3
- BgBase: 3
- BioXL: 2
- Arctos: 2
- CRS: 2
- Oracle: 2
- Jabot: 2
- MuseumPlus: 2

Not (yet) using dq feedback.
*Why? Which are tractable? (or not)*

- top selections were
  - lack of resources (time, staff, funds),
  - *not aware of feedback,*
  - *software challenges,*
  - *job is too massive,*
  - *not knowing where to find this feedback, and*
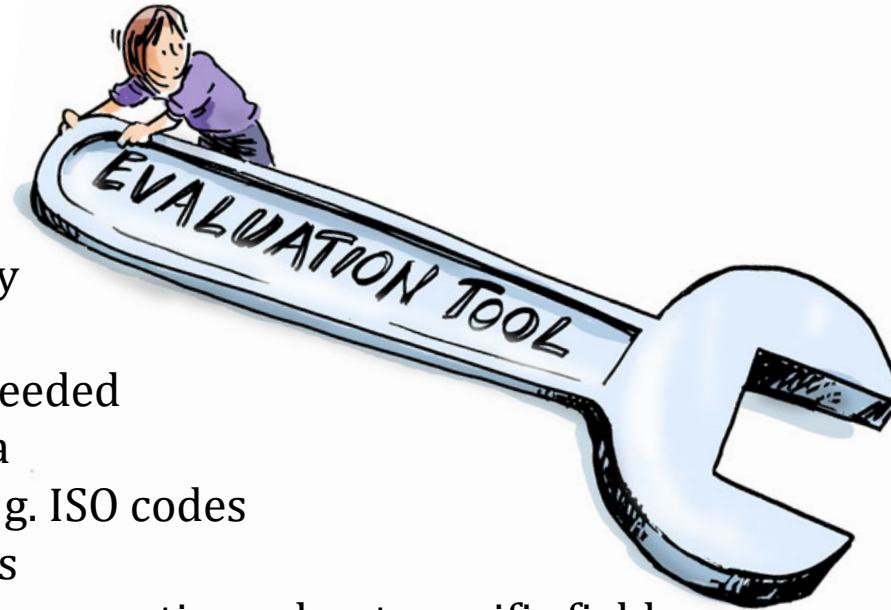  - *need more biodiversity informatics skills*

# Using dq feedback.
## *What challenges noted? Which are tractable? (or not)*

- lack of resources (time, staff)
- data quality feedback
  - organization priority differences
  - massive scope
    - have to prioritize or not a priority
  - workflow issues
    - assigning / tracking / approval needed
    - no one to manage geospatial data
    - can't make changes requested, e. g. ISO codes
  - erroneous feedback, e. g. taxon names
  - data standards knowledge missing, e. g. questions about specific fields
  - difficult to interpret
- Skills missing, or software impeded
  - changing by hand
- Prefer to work on curation rather than digitization DQ
  - An opportunity here

# DQ workflows to reveal skills, literacy, software

**Data and Collections Literacy**

- collection knowledge
- data entry
- formatting
- import functions
- knowledge of related biodiversity data
- knowing which records need attention
- parsing
- querying
- scripting
- taxonomic, geography, geology skills
- track collection / collector
- track down correct date
- understanding relational databases
- understanding feedback

**Software and skills**

- Emu
- how "the Atlas" works
- advanced spreadsheet skills
- database software
- postgreSQL
- able to use SQL / MySQL
- FileMaker
- Specify
- Symbiota
- Open Refine

# DQ feedback – benefits and changes

- dates
- distribution data
  – understanding scope of DQ issues
- fixing inaccuracies
- georeferences
- misspellings
- taxonomic name "insights"
- using dq feedback for prioritization
  – example: biosecurity and trade

*"…crowdsourcing data error detection is a better way to go rather than hiding data until you are sure it is all correct."*

## Comments and requests

- "county" boundary checks
- embed DQ tests in collections software?
- pest - host data quality checks
- community-proofing
- re-format of dq feedback (data downloads)
- metrics tracking
- more workshops, webinars, (data mgmt., open refine, …)

*"Making changes that are suggested is a big task and it is scary to make large batch changes."*
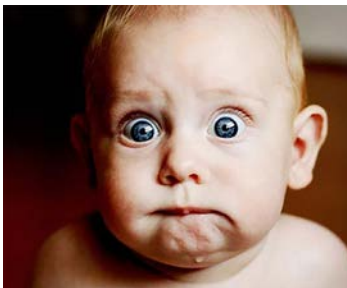
Key points or
*why it's important to shed light on the abysses*

- awareness of DQ issues

[now we are] "aware of how many different terms are used for the same things ... I hope as a community we can fix this issue soon."

https://bit.ly/niceuglydata

# Key points

- asking for and getting help
  - do you know what to do? is it working?
- need for transparency in DQ processes
  - manage expectations

## Key points

- software bottlenecks
- skills bottlenecks
  - … need for biodiversity data mobilization skills
    - source and development of skills?
    - changing roles in collections?
    - expectations for the future?

*"… huge need for the community to proof on-line data of all kinds. Extra eyes regularly find things that have [been] missed …"*

# What's next for data integration at the DQ feedback step?

- Algorithms can't do it all
- AI can't do it all
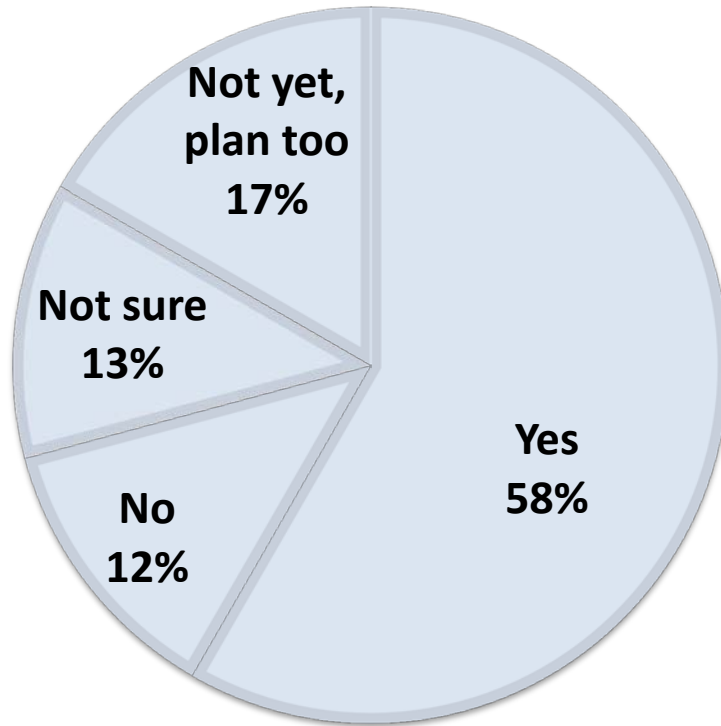- Data standards can't do it all

- What's your role?
- What's your abyss?

# Integrating Crowdsourced Data - community feedback

- 64 responses
- 54.69% use crowdsourced data



Integrate crowdsourced data

Pie chart:
- Yes 58%
- Not yet, plan too 17%
- Not sure 13%
- No 12%



Word cloud: increased-awareness, biodiversity, complex, volunteers, of, tool, collections, increased, data, time-consuming, more, public-awareness, challenging, volume, evaluating, validity, data-mapping, benefits, "cost"
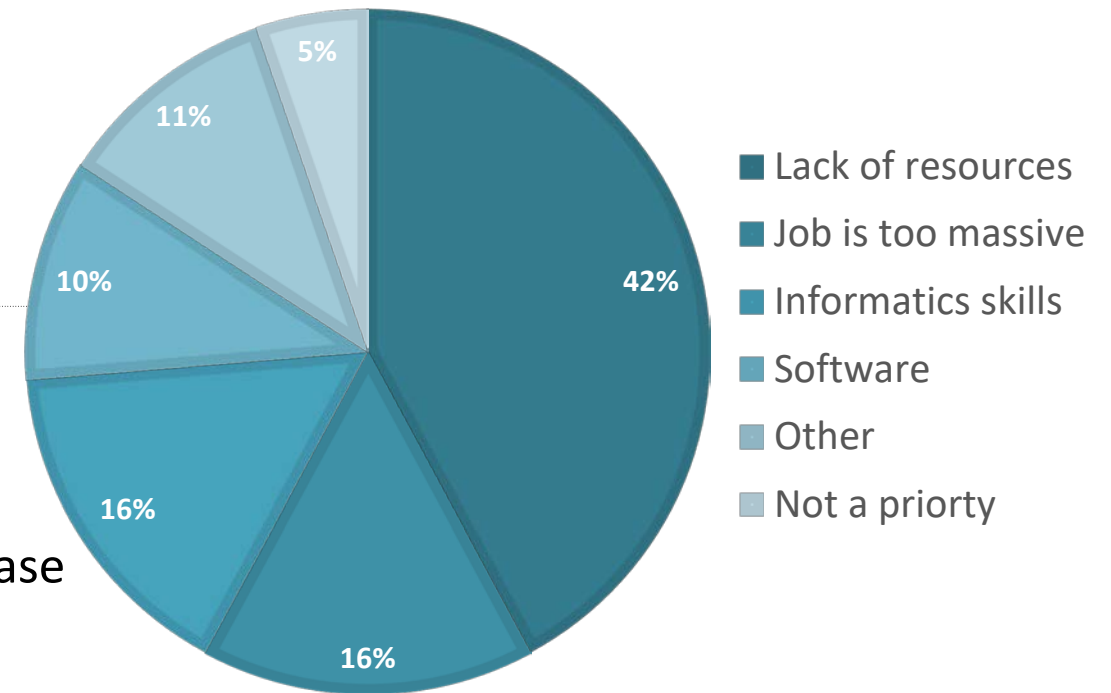
# Challenges - integrating crowdsourced transcription data?

- 42% responses - "Lack of resources"
- 16% responses – "Job is too massive"

## Other Challenges listed ……

- Complex data-mapping (not standardised DC)
- Data needs formatting
- Data validation
- More staff for validation and integration in database
- Procedure challenges



Legend:
- Lack of resources
- Job is too massive
- Informatics skills
- Software
- Other
- Not a priorty

Pie chart values: 42%, 16%, 16%, 10%, 11%, 5%

# Integrating crowdsourced transcriptions – what would help?

more focus on the issue from the "top”

broader ability to rate "trustworthiness" of transcribers

crowdsourcing portals integrated with local database
*note – Symbiota and @NfromN

standardization (process & data)

more expert validators

18

# Summary & future thoughts ...

- Does crowdsourcing actually save time?

- Many other benefits … increase in awareness of collection

- ICEDIG / DiSSCo : evaluating costs/benefits of different transcription choices.
  https://icedig.eu/content/deliverables

  - D4.2  Data quality in transcription                                                                      January 2019
  - D4.3  Data standards in transcription                                                                    February 2019
  - D4.4  Interoperability with institutional collection management systems    April 2019
  - D4.5  Cost analysis of transcription methods                                                   December 2019
  - D5.1  Recommendations for volunteer transcription systems and a source repository    April 2019

# Publications and Activities

- Belbin L, Daly J, Hirsch T, Hobern D, Salle JL. A specialist's audit of aggregated occurrence records: An "aggregator"s' perspective. *ZooKeys*. 2013;(305):67-76. doi:10.3897/zookeys.305.5438.

- Mesibov R. An audit of some processing effects in aggregated occurrence records. ZooKeys. 2018;(751):129-146. doi:10.3897/zookeys.751.24791.

- SPNHC_TDWGNZ W08 - Standardizing data to Darwin Core using R: A hands-on workshop with lessons learned from the TrIAS project. (2 – 3.30pm, Thursday)

- Project Paleo: Citizen Curation and Community Science at the Natural History Museum of Los Angeles County. - Elizabeth R Ellwood (4 – 4.20pm, Tuesday)

"… now on to the stories!"

# Challenges For Implementing Collections Data Quality Feedback:
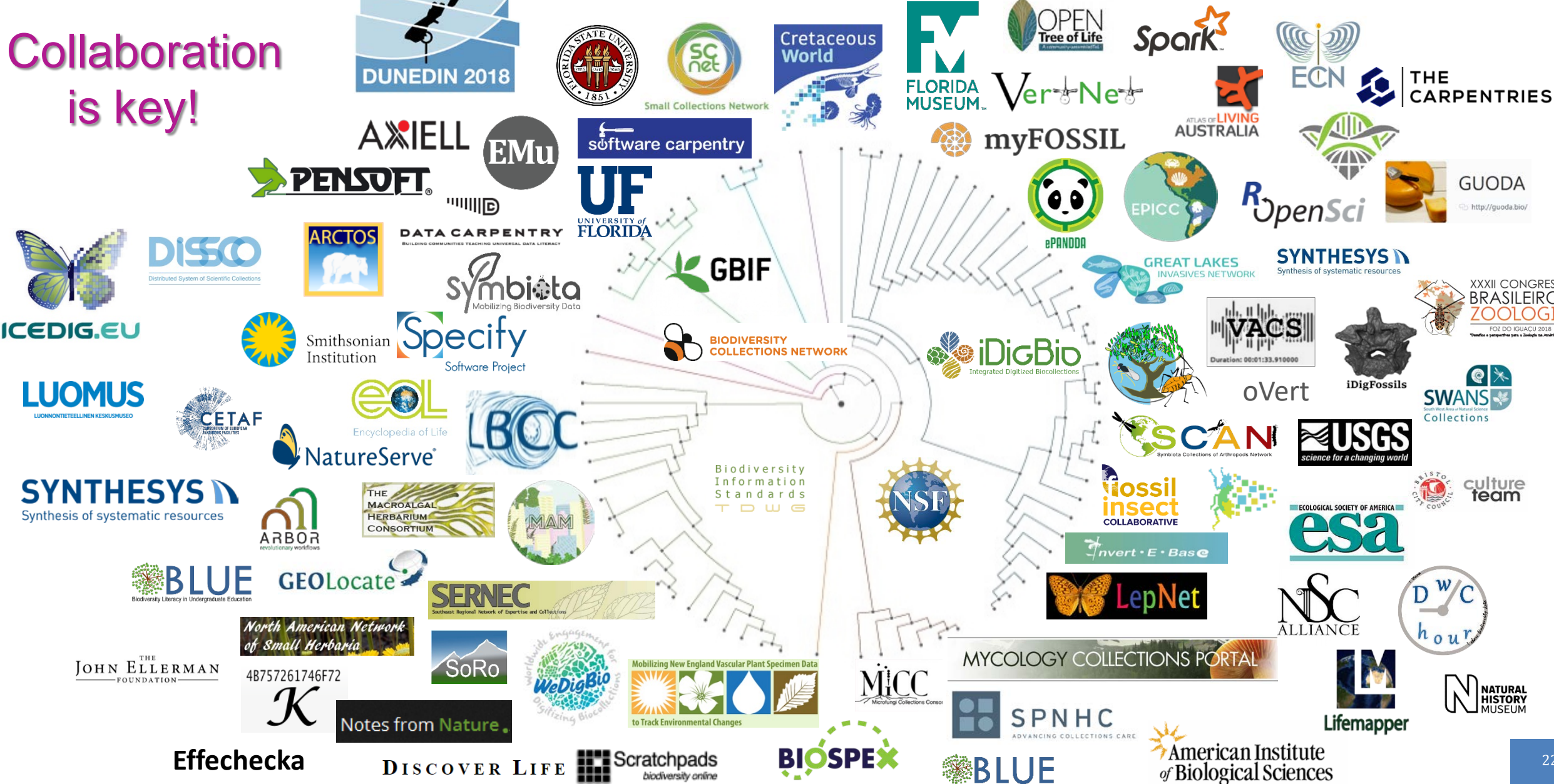## Synthesizing the community experience.

**2.20 – 2.40pm** : **Arthur Chapman**   Data Quality – Whose responsibility is it?

**2.40 – 3.00pm** : **Mare Nazaire**   Integrating Data Quality Feedback: a Data Provider's Perspective.

**3.00 – 3.20pm** : **Robert Cubey**   Label Transcript is Done – Now what do we do with that Data?

**3.30 – 4.00pm**   **Coffee Break**

**4.00 – 4.20pm** : **Andrew Bentley**   Practical use of aggregator data quality metrics in a collection scenario.

**4.20 – 4.40pm** : **Teresa Mayfield**   Who Has Time for Biological Collections Data Quality Feedback? Maybe a Community Can Help.

**4.40 – 5.00pm** : **Sharon Grant**   Repatriation of Augmented Information to an Institutional Database.

**From our speakers' data integration stories to yours**

- SPNHC #SIG on DQ Feedback is tomorrow for your part of the *#biodiversity #dataIntegration story*  21

Collaboration is key!

# Kia ora

## from Nicole Fisher and Deborah Paul

see you tomorrow too at the SPNHC #SIG Share your data integration stories – successes and snafus too!

Special thanks to Shari Ellis, iDigBio Project Evaluator, for her guidance when developing our ideas for this work. Very kind thanks to all our speakers for being ready, willing, and yes, even eager to share their data stories – including the juicy bits.

facebook.com/iDigBio

twitter.com/iDigBio

vimeo.com/idigbio

idigbio.org/rss-feed.xml

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