# echinopscis:
## an extensible notebook for open science on specimens

Nicky Nicolson
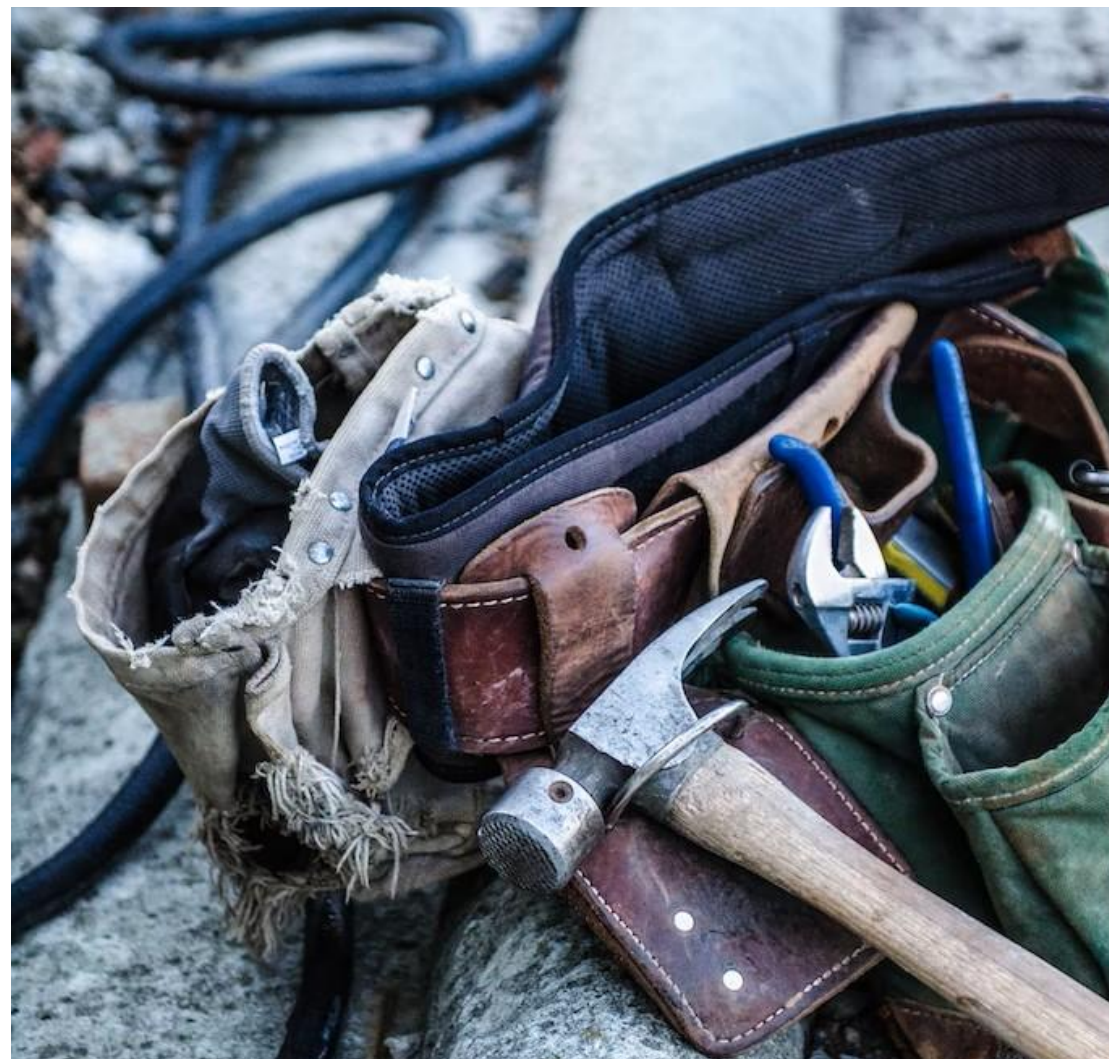Kew Science

**BHL**
Biodiversity Heritage Library

#BHLDay2023

Muséum national d'Histoire naturelle (MNHN)
Paris, France, 19 April 2023

# Context: personal & institutional

- Transitioned from software development into research

- Open science, take-up & how we design & build for participation

- How we can use software development practices in research:

  - Reuse

  - Automation

  - Version control

  - Dependency management

  - Continuous integration

- Also processes about communication, design & inclusion
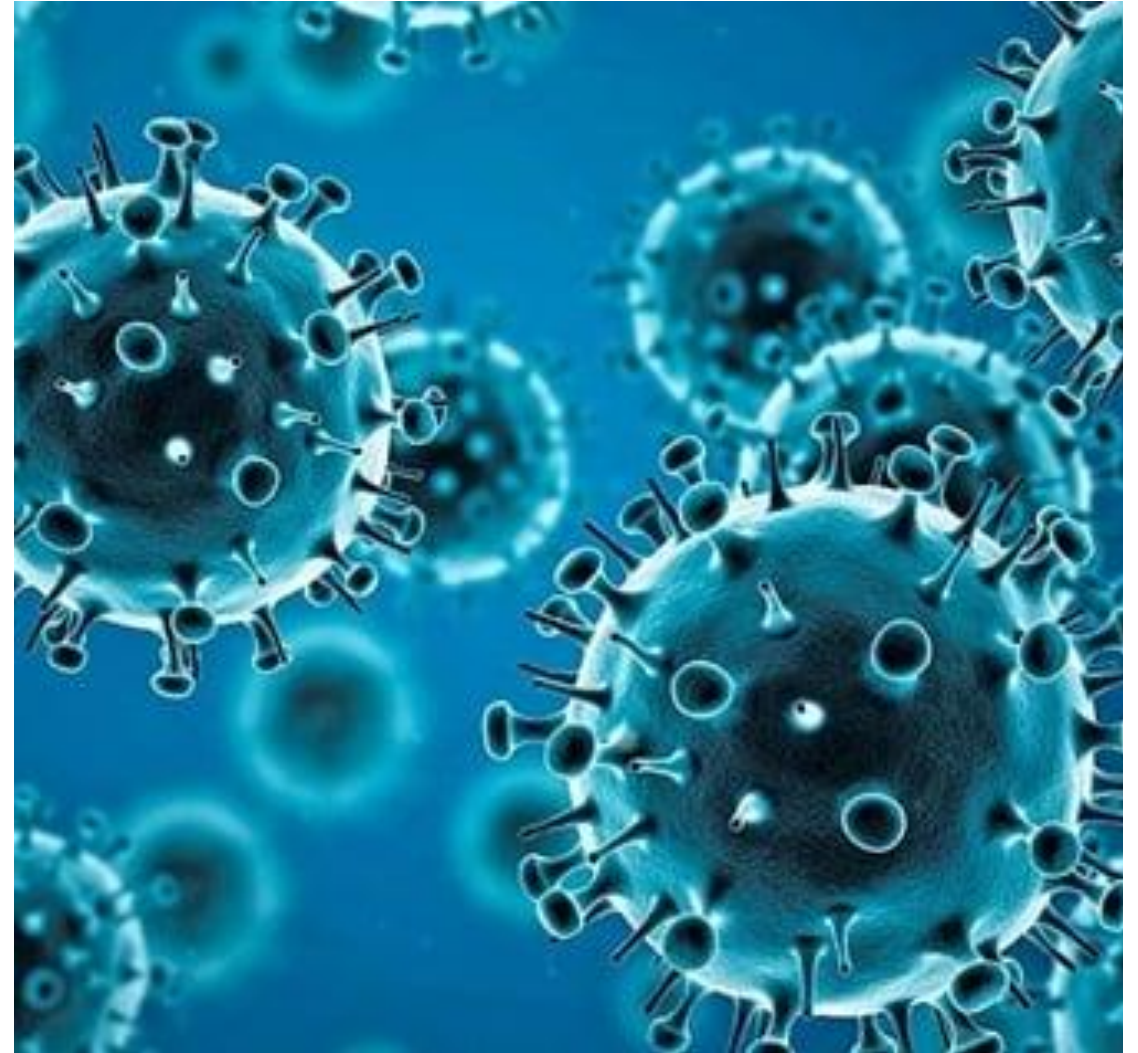
**Royal Botanic Gardens Kew**

Explosion of data availability online:

- Specimens online (Tracheophyta)
    - 88 million metadata records
    - 38 million images
- Comprehensive taxonomies with distribution
- **Born digital and digitised literature provides context**
- Machine learning / AI: text / images
- Recognised different roles in research
- Bring your idea to the data – compute available online

# **Where we work:** we're still learning how to transition online

- COVID made us move *everything* online

- When we work online, have we got a usable space, or is it "just another tab"

# Who does the work: evolving & inclusive research culture

- Skills development (Carpentries)
- Awareness of different roles in research
  - recognition of the research software engineer role
- Recognition of different activities required for successful research
- Open science: **open data**
- Online & remote collaboration

# How we work: building working environments online

- Inclusive design with researchers
- Prototype working environment:
  - Streamline access
  - Manual processes build training data
  - Mobilise rich data to publications & portals
  - A place to plugin AI approaches as these come onstream

- A personal knowledge manager: for creating & linking research notes
- Emphasises linking
- Data stored locally, using open formats
  - Markdown and optional structured data frontmatter
- Works offline
- Extensible architecture – plugins for data access and citation processing
- Active user and developer community

# OBSIDIAN

- A personal knowledge manager: for creating & linking research notes
- Emphasises linking
- Data stored locally, using open formats
  - Markdown and optional structured data frontmatter
- Works offline
- Extensible architecture – plugins for data access and citation processing
- Active user and developer community

# OpenRefine

# Extend Obsidian for specimen research

- Access of relevant data
    - Specimens (GBIF)
    - Names (International Plant Names Index)
    - Collections (Global Registry of Scientific Collections)
    - People (Bionomia)
    - Literature (crossref)
- Creation of links, spatial and network exploration
- Citation in new work
- Open science working practices

# echinopscis.github.io



**echinopscis**

Home  Team  Blog  Project ▾

**echinopscis**

An extensible notebook for open science

"echinopscis" is an experiment in creating an "extensible notebook for open science" - a working environment that allows researchers to write, access data and create links between literature, specimens, names, institutions, people, traits etc.

## Key principles:

- **Control of your data**: as a researcher, you remain in control of your data. The data is stored in text format, on your local machine. Text files are an open format, they will always be accessible without any need for specialised software.
- **Open to choose your working practices**: we've provided small pieces of functionality that can be combined in many different ways, enabling researchers to be "open to choose" how to organise their work.
- **Re-usable skills**: any skills necessary to work with this toolkit should be transferable to other open science tools and practices. If you invest in time exploring this prototype software, the things you learn (markdown formatting, bibliography / citation management, document production etc) could also be applied elsewhere in your work, or in other working environments.
- **Open science**: All code and documentation (and this project site) are managed on github - contributions are welcome.

# Demo: Unstructured text data to specimen links

# Demo: "Workbench" specimen examination

# Provisional roadmap (open to influence)

1. **Personal research environment** based on Markdown authoring and linking

2. **Web publication** using static site generators (conceptual similarities with the [GBIF hosted portal](#) work)

3. **Document production**: with structured bibliographic/specimen references

4. **Dataset production**: mobilisation of content and links into DarwinCore archives for aggregator harvesting

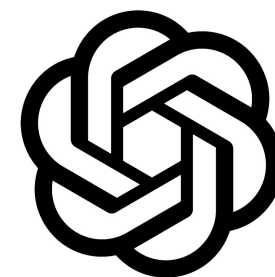# BHL relevance (1) data: integration of "meta" works like TL-2

- **TL-2 is entity rich:**
  - People
  - Places
  - Institutions
  - Timelines
  - Expeditions
  - Bibliographic works
  - Eponyms
- We have corrected text data in addition to BHL OCR (Smithsonian data package)
- Reference resource for researchers working on Bionomia

Royal Botanic Gardens Kew

mobilise

Bionomia

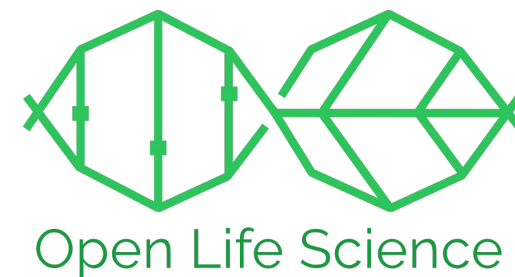**Link natural history specimens to the world's collectors**

- echinopscis relies on "strings to things" style APIs - pass a piece of text like might be found in unstructured paragraph, get a reference to an entity back
- BHL API requires atomised input
- Role for BHL to provide reference parsing (whatever the approach, BHL can learn from usage traffic):
  - anystyle.io
  - ChatGPT

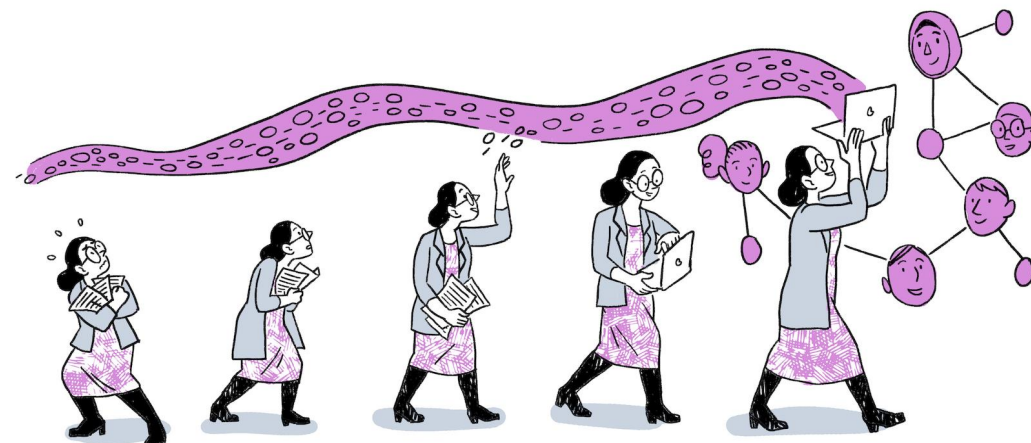# BHL relevance (3) approach: community development

Promoting use of generic tools through skills development, investing in development of research culture

- **Library Carpentry** (short training courses)
  - Use of reference managers
  - Automated bibliographic production
  - Article & author identification (DOI / ORCID)
  - Deposition
- **Open Life Sciences** (16 week mentoring):
  - Project management: "open by design"
  - Community building, inclusivity

# Conclusions: "e-taxonomy" lives(!) as "open science"

- We've built an open project based on participation & reusable skills

- Include people where they are now, show a relevant path to destination

- Extensible with AI approaches (entity identification, clustering, summarisation, link prediction)

- Usable today with generic skills - try it out: https://echinopscis.github.io



EVOLVING TOWARDS AN ERA OF OPEN RESEARCH

Scriberia

*n.nicolson@kew.org* / *@nickynicolson* / *@nickynicolson@mastodon.social*

# Image credits

- Slide 2: Leather toolbelt by jesse orrico (jessedo81) on Unsplash
- Slides 3 & 6: RBG Kew
- Slide 4: WHO
- Slides 5 & 16: The Turing Way project illustration by Scriberia. Used under a CC-BY 4.0 licence. DOI: 10.5281/zenodo.3332807

# Useful links

- echinopscis: https://echinopscis.github.io
- Obsidian: https://obsidian.md/  & Obsidian roundup (weekly newsletter): https://www.eleanorkonik.com/tag/roundup/
- Open Refine: https://www.openrefine.org
- Training resources:
  - Carpentries (including Library Carpentry): https://carpentries.org/
  - The Turing Way: https://the-turing-way.netlify.app/
  - Open Life Sciences: https://openlifesci.org/

# Contact

- Email: n.nicolson@kew.org
- Github: @nickynicolson
- Twitter: @nickynicolson
- Mastodon: @nickynicolson@mastodon.social
- Web: https://www.kew.org/science/our-science/people/nicky-nicolson