# A story of data won, data lost and data re-found: the realities of ecological data preservation

Alison Specht

School of Earth and Environmental Sciences
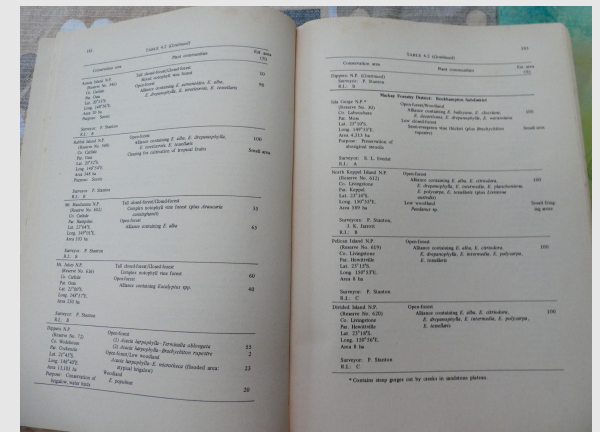
THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA

with collaborators

Matt Bolton, Corymbia Ecospatial Consultants,
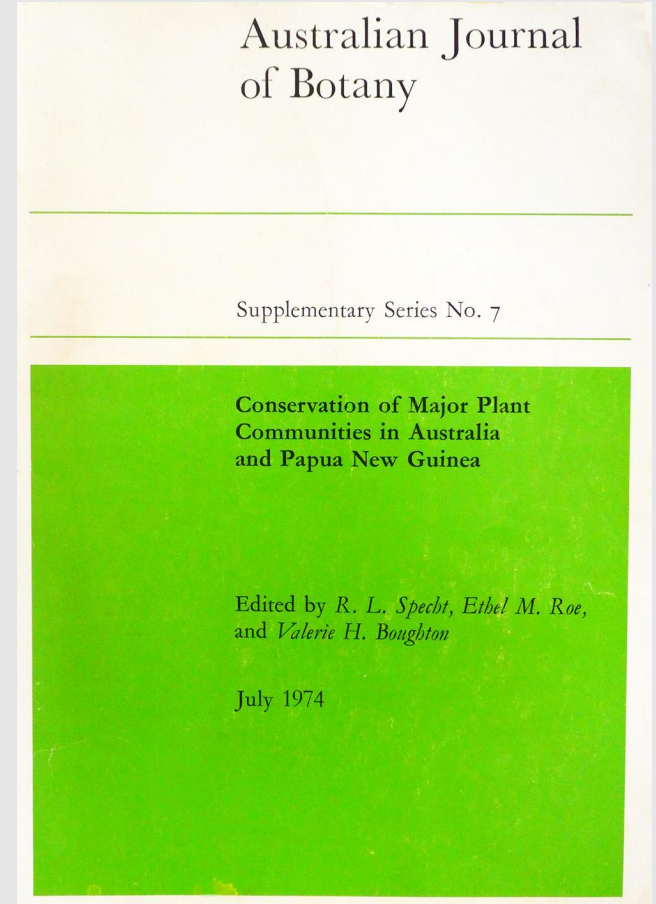
Lee Belbin, Atlas of Living Australia.

# The story…



- In the beginning…

  - the advent of the big computer age: heralding new possibilities and vision for data manipulation

- Potential disaster! Risk of imminent data loss

- But rescue was in sight! someone cared…

  - A program of retrieval and recovery commenced

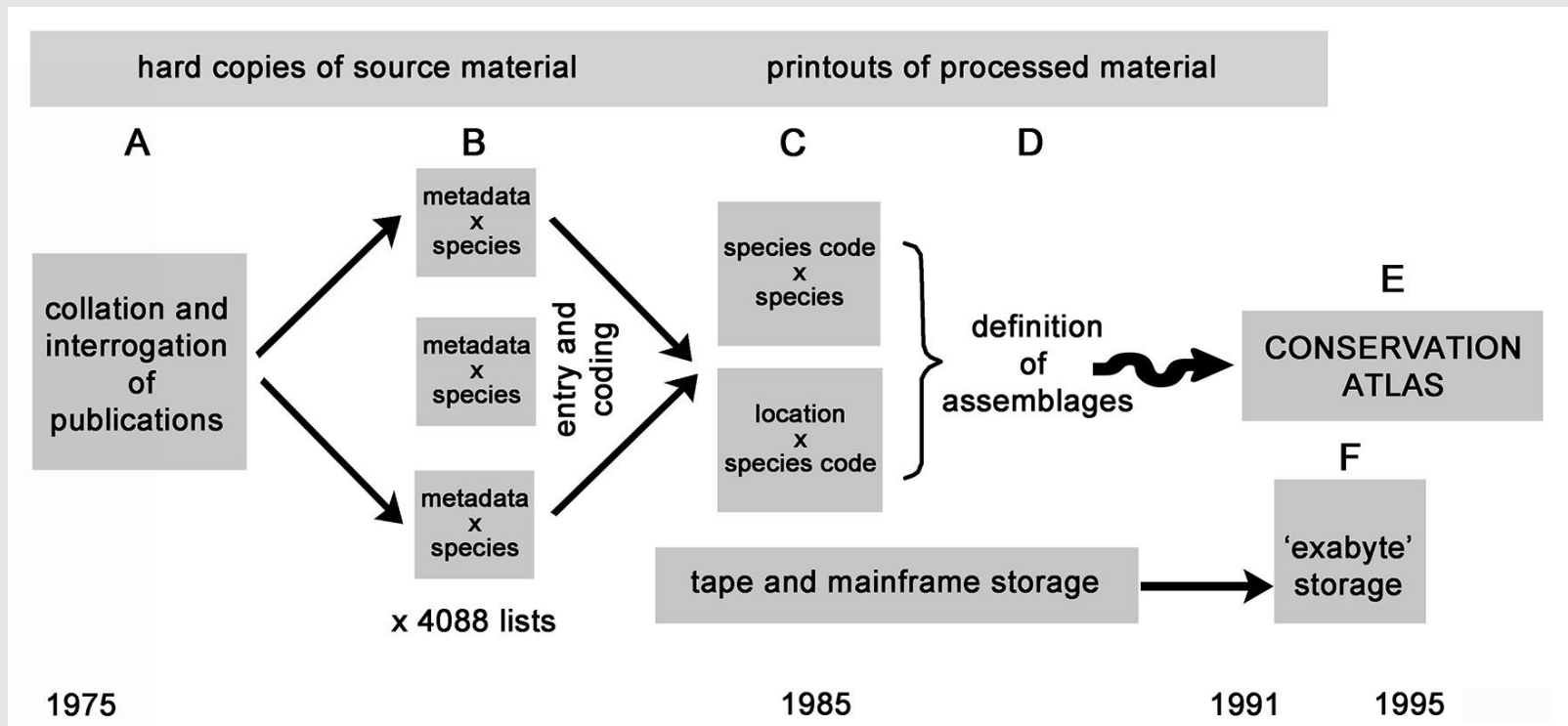- This talk: what we learnt…how might this help others?

# In the beginning – Phase 1

- Conservation of ecosystems and their biota demands knowledge: what are they composed of, how unique are they, where are they?

- In 1974 a conservation survey of Australian plant communities was published. This was based on expert opinion and although comprehensive and innovative, it was imbalanced as expertise varied across systems.



Australian Journal of Botany

Supplementary Series No. 7

Conservation of Major Plant Communities in Australia and Papua New Guinea

Edited by R. L. Specht, Ethel M. Roe, and Valerie H. Boughton

July 1974

# In the beginning – Phase 2

- The advent of big computing meant an objective assessment could be made.

- Research project started, led by R.L. Specht

# Collation and extraction of information from publications (A, B) for entry into computers

# Data organization and entry (B & C)

Due to computing capacity, the data were organized into state x formation datasets to be read by FORTRAN programmes.

**State:** N = New South Wales, V = Victoria, T = Tasmania etc.
**Formation:** Closed forests, chenopod shrubland, desert acacia etc.

| LINE ID | Information |
|---------|-------------|
| 800000 | N |
| 503200 | LOCATION N032 = CENTRAL COAS |
| 903200 | 33 51 1 |
| 503201 | COMMUNITY 01 = FRESHWATER |
| 003201 | UTRIAUST UTRIEXOL UTRIBILO VA |
| # | |
| 003201 | NAJAMARI MYRIPROP PHRAAUST |
| TRIGSTRI # | |
| 003201 | JUNCPAUC JUNCPALL JUNCPLAN |
| MELASTYP # | |
| 003201 | CALLSALI EUCAROBU EUCAAMPL |
| 003201 | GRATPUBE GOODPAN HYDRREDI |
| SCHOAPUG # | |
| 003201 | ORLIIMBE BLECINDI AD |
| 503202 | COMMUNITY 02 = FRES |
| 003202 | BAUMTERE BAUMART |
| 003202 | ISOLINUN GRATPEDU DROSSPAT |
| # | |
| 00 | A BOROPARV |
| SPHAGNUM* # | |
| 003202 | VIOLHEDE # |
| 500000 | ---------------------------- |

*Annotations:* Unique identifier · Latitude longitude · program (e.g. spe... EUCARO... · Two lists for this location

| Formation | Locations | Communities | Species* |
|-----------|-----------|-------------|----------|
| Closed forests | n/a | 644 | 1,418 |
| Dry scrubs – SE Queensland | 232 | 232 | 475 |
| Dry scrubs – Northern Territory | n/a | 1,219 | 559 |
| Eucalypt open-forests and woodlands (tree species) | 201 | 1,275 | 276 |
| Sclerophyll vegetation SW Western Australia | 64 | 172 | 1,761 |
| Sclerophyll vegetation Central and Eastern Australia | 188 | 549 | 2,581** |
| Sclerophyll vegetation – heathland and tall shrubland | 136 | 312 | 2,071** |
| Alpine vegetation | 73 | 61 | 556 |
| Savanna understorey | 56 | 198 | 1,313 |
| Mallee open-scrub | 28 | 41 | 395 |
| Desert Acacia | 54 | 148 | 1,229 |
| Chenopod shrubland | 30 | 68 | 410 |
| Forested wetlands (including brigalow) | 31 | 36 | 193 |
| Arid wetlands | 20 | 42 | 642 |
| Freshwater swamp vegetation | 80 | 80 | 139 |
| Coastal dune vegetation | 45 | 56 | 315 |
| Coastal wetland vegetation (mangroves and saltmarshes) | n/a | 15 | 74 |

\* Not including introduced species or singletons within the formation.
\*\* Not including tree species > 10m tall

# Data

- **Entry and storage**
  - Punch cards then desktop computers were used for data entry to UQ's PDP-10. 9-track magnetic tapes used as regular backup.

- **For analysis**
  - Analysis on CSIRONET mainframe computer (TAXON & TWINSPAN).
  - Hard copies (as in print-outs for proofing and run outputs) obtained throughout.

- **Data processing**
  - Described in a procedures manual (CAVE: Bolton)

# Product (D & E): 911 objectively-defined plant communities, mapped, keys for their identification, their conservation status, and biogeographic regionalization…

Specht R.L.. and Specht A. (2013) Australia: Biodiversity of Ecosystems. In, *The Encyclopedia of Biodiversity* Vol. 1 (ed. B. Levin, et al.) pp 291-306. Waltham, MA: Academic Press.

Specht R.L. and Specht A. (2002) Objective classification of plant communities in tropical and subtropical Australia. *Proceedings of the Royal Society of Queensland* 110: 65-82.

# But what about the data?

- The primary objective of phase 2 was the research, secondarily to find a home for public access.

- In 1995 there was no 'home', so the data were 'saved' on the magnetic tapes and subsequently exabyte tapes when the main frame reader was de-commissioned. The print-outs were conserved.

- High-level data for biogeographical analysis (by PATN) was saved on excel

- *So there they sat…until someone cared…*

# Why should we care?

**Value proposition**

These are heritage data. They were collected on field trips from 1879-1989, and provide unique records for comparison.

Repeating initial project work would be painful, if not impossible

## Opportunities in the 2010s

new data repositories were emerging; the Terrestrial Ecosystem Research Network (TERN) and the Atlas of Living Australia (ALA) linked globally to DataONE, GBIF, KNB etc

AND, key members of the team were still alive, personally invested and new team members identified.

☐ *Can we save the data and finally make it available?*

# Retrieval – with support from TERN and the ALA

- Recover available data

- Design an appropriate structure

- Update the species codes/names to current nomenclature

- Update georeferencing and check errors

- Map the fields used in the Conservation Atlas project to the Darwin Core standard

- Deliver the data in an open repository



TERN

Terrestrial Ecosystem Research Network

ATLAS OF LIVING AUSTRALIA
sharing biodiversity knowledge

# Recover available data

- As mentioned, the data were moved from mag. tapes to exabyte tapes in 1991

- **Challenge: (a)** find exabyte tapes, and **(b)** an exabyte tape reader.



- Started on the print-outs:

  - Master sites file (location, formation & community data)

  - [Reference file](#) (source data)



- Finally, the last exabyte tape reader in captivity was found (and about to be de-commissioned)!

- Two major challenges remained: updating [georeferences](#) and [species names](#)

# Master sites file

**From original printouts (slightly updated)**

1. The formation, location and community number (1,2 etc)
2. Locality: general description (soil type, landscape etc)
3. The source reference (link to reference file)
4. Latitude and longitude (degrees minutes)
5. Broad community description
6. Additional information such as dominant species or association
7. Notes

**From retrieval team**

8. Decimal latitude and longitude
9. Coordinate uncertainty in metres
10. Comments (using a consistent vocabulary)

# Reference file

| ID | Author(s) | Date | Title | Journal etc. | Volume No. | Page numbers |
|----|-----------|------|-------|--------------|------------|--------------|
| 1 | Abbott, J. | 1977 | Species richness, turnover and equilibrium in insular floras near Perth, Western Australia. | Aust. J. Bot. | 25 | 193-208 |
| 8 | Adams, L. D. & Craven, L. A. | 1976 | Checklist of vascular plants in a study area of the South Coast of N.S.W. | C.S.I.R.O. Land Use Res. Tech. Mem. | 76/16 | |
| 387 | McMahon, A.R.G., Carr, G.W., Todd, J.A. & Race, G.J. | 1990 | The Conservation Status of Major Plant Communities in Australia: Victoria. | Ecological Horticulture Pty Ltd, Clifton Hill, Vic. | | |
| 474 | Pye, K. | 1982 | Morphology and sediments of the Ramsay Bay sand dunes, Hinchinbrook Island, North Queensland. | Proc. R. Soc. Qld | 93 | 31-47 |
| 560 | Tate, R. | 1880 | On the geological and botanical features of southern Yorke Peninsula, South Australia. | Trans. R. Soc. S. Aust. | 13 | 112-120 |
| 705 | Willis, J.H. | 1967 | Systematic arrangement of vascular plants noted on the slopes and summit of the peak: The Rocks Nature Reserve, New South Wales. | Nat. Pks & Wildl. Serv., N.S.W. | 705 | |

# Georeferences

Original locations were accurate to half a degree which was unacceptable in the present day so the team did four things:

- Reviewed original documents and where possible contacted authors to update locations

- Checked locations on google maps

- Checked locations on the ALA's Spatial Portal so vegetation and soil type could be displayed for checking

- Mapped data repeatedly on the ALA sandbox site.

Co-ordinate precision was then estimated to reflect confidence in the range of the community.

On-line resources such as the Biodiversity Heritage Library, the National Library of Australia



Original articles





Maps in Appendices often not scanned in digital copies of old journals

# Species names



## 1. CODES to NAMES

- apply master species conversion file
- blend across formations (with caution as some species names are location- and formation-specific)

| Sequential row number | Validity and Growth habit flag | species code | Original scientific name | Scientific names updated during Conservation Atlas project |
|---|---|---|---|---|
| 2 | L G | ABELMOSC | Abelmoschus moschatus | |
| 19 | LMG | ACACARGY | Acacia argyrodendron | |
| 20 | SZG | ACACARMA -> ACACPARA | Acacia armata | Acacia paradoxa |
| 21 | MLG | ACACASHA -> ACACOSHA | Acacia ashanesii | Acacia oshanesii |
| 174 | S G | ACAKEMP | Acacia sp. aff. A. sibirica | Acacia sp. aff. A. kempeana |
| 466 | S G | BORRCARP/ -> SPERSTEN/ | Borreria sp. aff. carpentariae | Spermacoce sp. aff. stenophylla |
| 704 | S G | CARPAEQU -> CARPMODE | Carpobrotus aequilaterus | Carpobrotus modestus |
| 705 | L G | CARPMODE | Carpobrotus modestus | |

# Update to current nomenclature

**Stage 1.** Current name check

Due to the size of the data set, the Atlas of Living Australia web service lookup (BIE) was employed, with codes allocated for follow-up (or not).

**Stage 2.** Validation

**Stage 3.** Reference to an expert

Resources used included:
1. On-line national species records
2. State species records
3. Books and papers
4. Experts

| CODE | Meaning | action |
|---|---|---|
| **MATCH** | Near-exact match or better | accept |
| **PARTIAL-L and PARTIAL-R** | A significant substring match | manual check |
| **FUZZY** | Fuzzy matching algorithm built on the score from the web service using a 'letter-pair similarity' score | manual check |
| **WEAK** | A weak match falling below thresholds; the best match is retained | manual check |
| **TAXM** | No match or major problem with original or subsequent species name | refer to expert |

# Map the fields used to the Darwin Core standard

| row # | Target DwC Field | ALA field | Source of Field Contents | Remarks |
|---|---|---|---|---|
| 1 | datasetID | DataResource | ALA-generated | |
| 3 | catalogNumber | Catalog number | Concatenation of CAVE data: formation dataset-location number-community number-line number-position in the line | Allowable values for position in the line are 1-8, inclusive. |
| 4 | occurrenceID | Occurrence ID | Concatenation of CAVE data: species alphacode-formation dataset-line number-position in the line (allowable values 1-8) | Allowable values for position in the line are 1-8, inclusive. |
| 23 | scientificName | Scientific name | Scientific name as CAVE data matched to current name by ALA BIE facility. (Unless the name match was overridden manualy.) | Overrides, where present, were made by authors MB and/or RLS. See also identificationFlag. |
| 24 | taxonRank | Taxon rank | Generated from scientificName by ALA, unless overriden by taxon master file in cases of genus-level taxa. | |
| 39 | habitat | Habitat | Derived from Vegetation_Type in master sites file and CAVE data prefixed with 1, 2 or 3 and expanded via lookup tables. | |
| 43 | locationRemarks | Location remarks | Field Veg2Association from master sites file plus text from CAVE comment lines for relevant location and vegetation community. | |
| 44 | coordinatePrecision | Coordinate precision | "0.000278" (nearest second), "0.01667" (nearest minute) | |
| 45 | coordinateUncertaintyInMeters | Coordinate uncertainty in meters | from master sites file | Estimated manually, mostly by AS. |
| 46 | georeferenceVerificationStatus | Georeference verification status | from field: "comments - all locations verified using google maps." in master sites file | |

# Data delivery

- Ingested into the Atlas of Living Australia as a collection, discoverable through species records with associated metadata:

    - https://collections.ala.org.au/public/show/dr8212

- Delivered as excel with associated code for replication of the process in the Knowledge Network for Biocomplexity:

    - http://doi.org/10.5063/F1QC01QK

- In the future, discoverable as plot information on other sites (e.g. TERN).

# How did we do?

✓ Data saved, updated and deposited for future use in two stable repositories.

✓ 9450 taxa found in 1390 communities at 461 locations across the continent of Australia, between 1879 and 1989. This is a lot!

But this represents only around half of the original resource. Why?

The primary cause was loss of data on transfer from magnetic tape to exabyte tape back in 1991. And it appears in some instances those data cannot be found elsewhere.

# So what?

## Challenges

*Lots of talk but too little action – I propose*

- We neglect our valuable and hard-won data because of the dominant research imperative and lack of funding and rewards for data management

- Technological change

- Metadata (what are those rows and columns, the units the dates etc.)

- Curated, stable and accessible repositories

## Lessons learnt

- we need to deposit data and metadata for future re-use as soon as possible after creation,

- We need to have repositories that are open but secure, and are properly managed for technological change in the long term

- For data archiving, don't work individually or at the small scale, team with others

*Without this more data will be lost than were ever gathered and analysed.*

# Thankyou!



**Contact:** a.specht@uq.edu.au
School of Earth and Environmental Sciences
The University of Queensland, Australia

Biodiversity Data Journal  https://bdj.pensoft.net/article/28073/
Knowledge Network for Biocomplexity https://knb.ecoinformatics.org/#view/doi:10.5063/F1QC01QK
Atlas of Living Australia https://collections.ala.org.au/public/show/dr8212

TERN

ATLAS OF LIVING
AUSTRALIA
sharing biodiversity knowledge