

# Challenges

- Many challenges – include:
  - Finding and accessing data
  - Making data useful
  - Dealing with legacy data
  - Improving Data Quality
  - Dealing with sensitive data



# The data equation

Oceans of  
Data



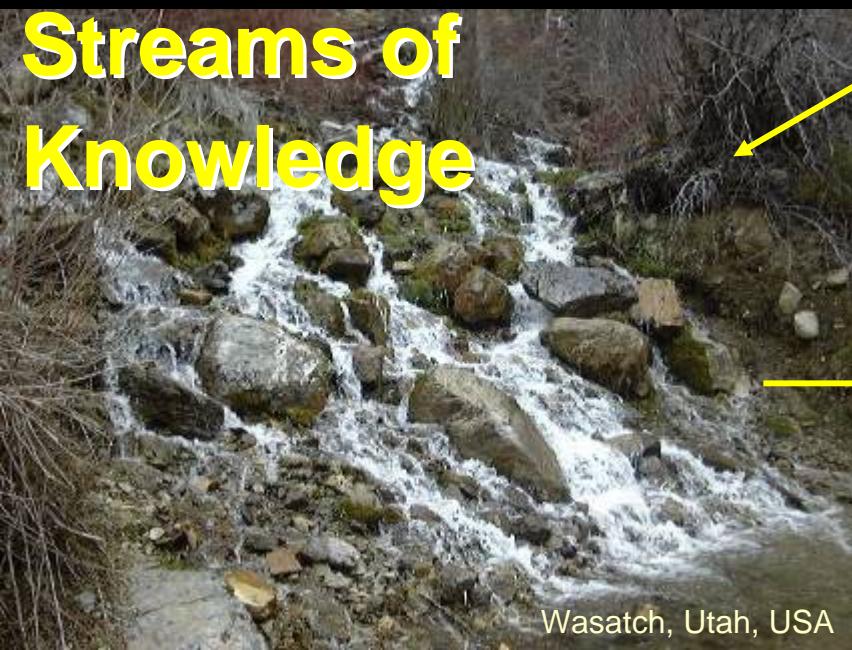
Praia de Forte, Brazil

Rivers of  
Information



Doubtful Sound, New Zealand

Streams of  
Knowledge



Wasatch, Utah, USA

Drops of  
Understanding



17 Nov

(Nix 1984)

# Challenges

- Finding and accessing data



# How many species are there?



Australian Government  
Department of the Environment,  
Water, Heritage and the Arts



australia's nature  
*there is more  
still to be discovered...*



## Numbers of Living Species in Australia and the World

2<sup>nd</sup> edition

Arthur D. Chapman  
Australian Biodiversity Information Services  
Toowoomba, Australia

Report for the Australian Biological Resources Study  
Canberra, Australia  
September 2009

Taxon	World Descr./ Accepted	Australia Descr./ Accepted	Austral. Percent.	Estimate World	Estimate Australia	World Threat. <sup>14</sup>	World Threat. Percent.	Aust Threat. <sup>15</sup>	Austral Threat. Percent	% of World's Threat.	Percent. Endemic
Chordates	64,788	~8,128	12.5%	~80,500	~9,088	5,966	9.2%	246	3.0%	4.1%	41.3%
Invertebrates	1,359,365	98,703	7.3%	~6,755,830	~320,465	2,524	0.2%	32	0.04%	1.3%	unknown
Plants	310,129	24,716 <sup>16</sup>	7.9%	~390,800	26,845	8,457	2.7%	1,263	5.1%	14.9%	86%
Fungi	98,998	11,846	11.9%	1,500,000	50,000	3	>0%	0	0%	0%	unknown
Others	~66,307	>4,186	6.2%	2,600,500	~160,000	6	0.01%	0	0%	0%	unknown
<b>TOTAL 2009</b>	<b>1,899,587</b>	<b>147,579</b>	<b>7.8%</b>	<b>~11,327,630</b>	<b>~566,398</b>	<b>16,956</b>	<b>0.9%</b>	<b>1,541</b>	<b>1.1%</b>	<b>9.1%</b>	<b>unknown</b>

# Challenges

- Making data useful



# Taking data to information

Species Data



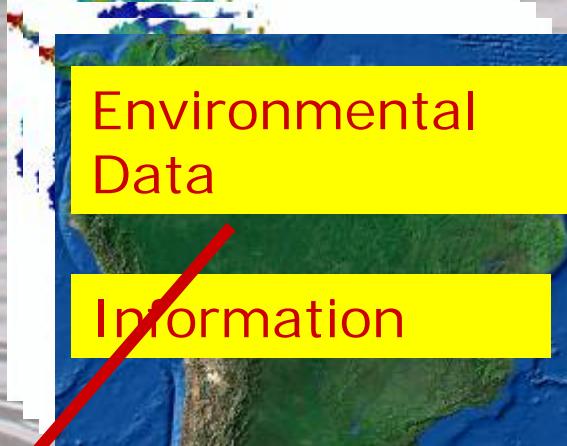
Stick Insect  
Campinas, Brazil

Species Data



Eucalyptus nicholii  
California

Environmental Data



Information



Decisions

Policy

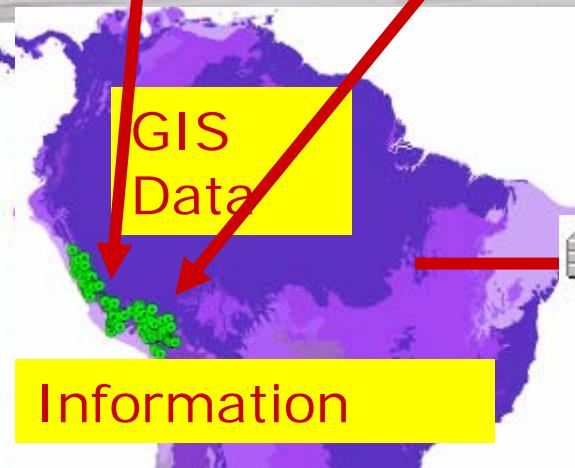
Conservation

Management



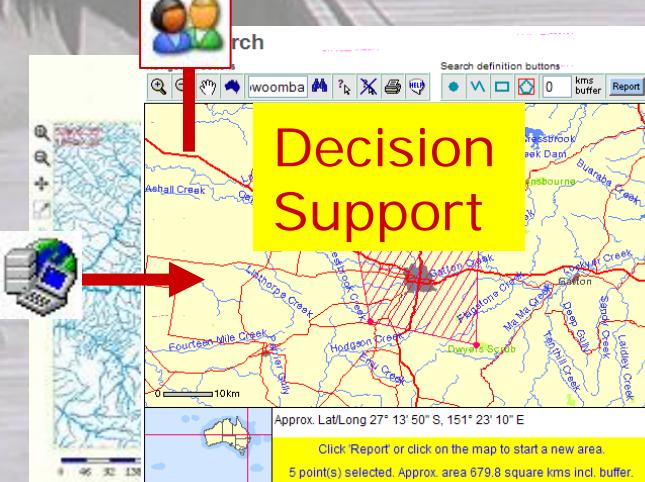
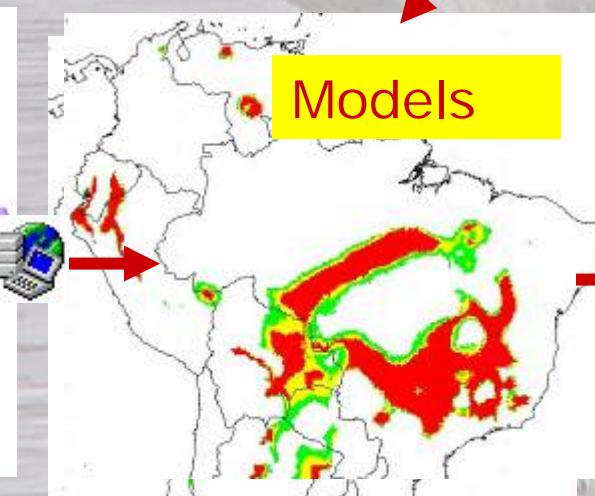
Decision Support

GIS Data

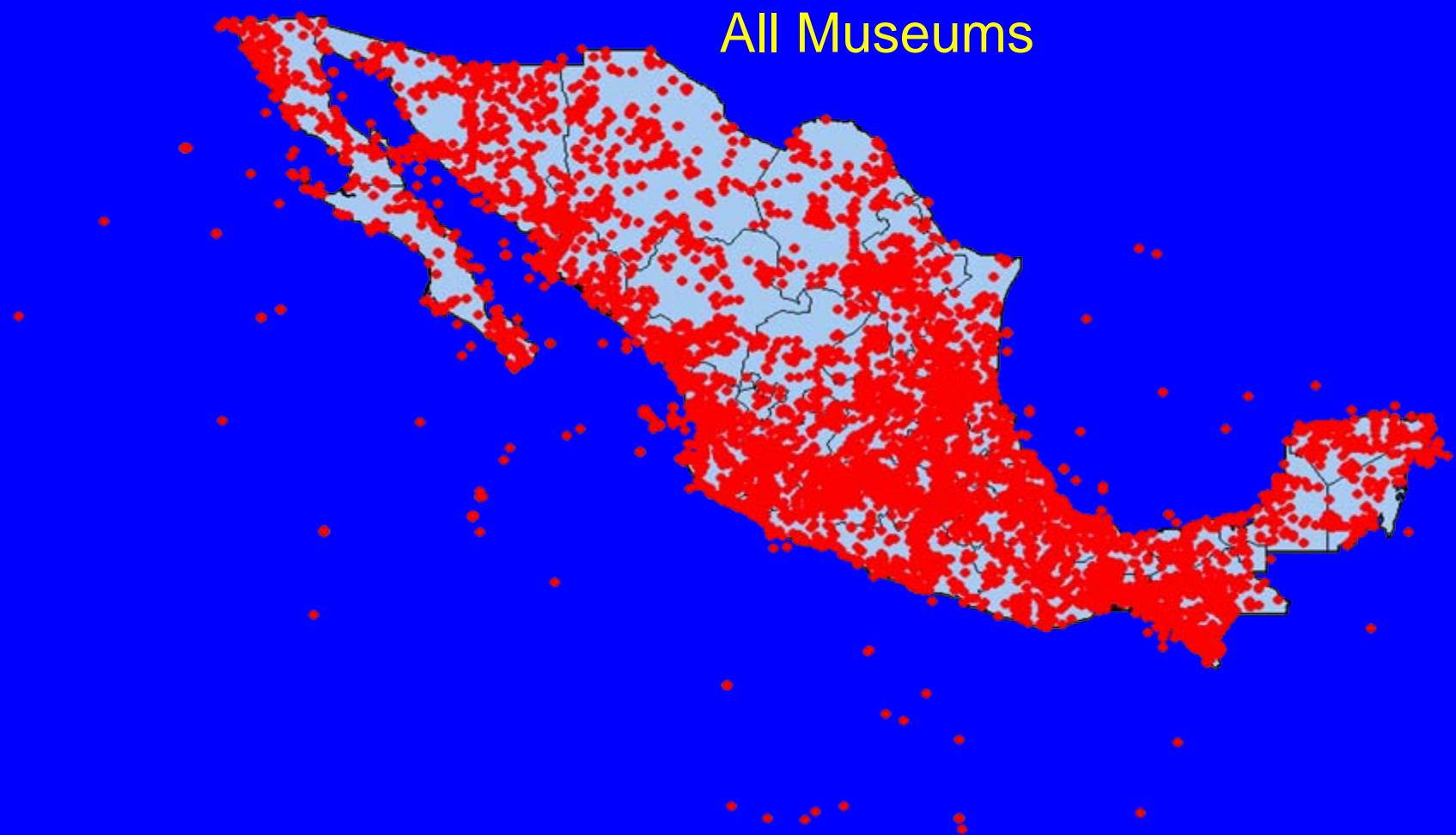


Information

Models



# Distributed studies using mexican birds



# Using species data

Taxonomic Studies, Ecological Biogeography, Phylogenies

Biogeographic Studies, Species Modelling

Species Diversity and Population studies

Life Histories and Phenologies

Studies of Threatened and Migratory species

Climate Change Impacts

Ecology, Ecosystems, Evolution and Genetics

Environmental Regionalisations

Conservation Planning

Natural Resource Management

Agriculture, Forestry, Fisheries and Mining

Health and Public Safety

Bioprospecting

Forensics

Border Control and Wildlife Trade

Education and Public Outreach

Ecotourism

Art and History, Science and Politics

Recreation

Human Infrastructure Planning



Arthur D. Chapman<sup>1</sup>

## Abstract:

This paper examines uses for primary species-occurrence data in research, education and in other areas of human endeavour, and provides examples from the literature of many of these uses. The paper examines not only data from labels, or from observational notes, but the data inherent in museum and herbarium collections

themselves, which are long-term storage receptacles of information and data that are still largely untouched. Projects include the study of the species and their distributions through both time and space; their use for education, both formal and public; for conservation and scientific research; use in medicine and forensic studies; in natural resource management and climate change; in art, history and recreation; and for social and political use. Uses are many and varied and may well form the basis of much of what we do as people every day.



# Challenges

- Dealing with legacy data



# Georeferencing - The Past

- Legacy data
  - Lots of data
  - Very little data with georeferencing
  - Georeferencing added from paper maps
  - Little accuracy recording or data quality checking

*Plant and animal specimen data held in museums provide a vast information resource, providing not only present day information on the locations of these entities, but also historic information going back several hundred years*

(Chapman and Busby 1994).

# Georeferencing – the present

- Databasing the backlog:
  - Many herbaria and museums well advanced
  - Many more collections still to do
- Use of GPS
- Distributed data (e.g. GBIF, AVH, OZCAM)
- Use of automation and new technologies
  - Filtered Push
- Recording accuracy and error information

# Recording Accuracy and Uncertainty

## Additional Uncertainty Fields

- Preferably in meters (Point-Radius)
- Remarks



## Documenting Validation tests

- Who
- What
- How

# MaNiS Georeferencing Calculator

Version 020411

## Georeferencing Calculator

Calculation Type: Coordinates and error - enter the Lat/Long for the named place or starting point

Locality Type: Distance at a heading (e.g., 10 mi E (by air) Bakersfield)

**Step 3) Enter all of the parameters for the locality.**

Coordinate Source:	USGS map: 1:24,000	Offset Distance:	10
Coordinate System:	degrees minutes seconds	Extent of Named Place:	3
Latitude:	35° 22' 24" N	Distance Units:	mi
Longitude:	119° 1' 4" W	Distance Precision:	1 mi
Datum:	(NAD27) North American 1927	Direction:	E
Coordinate Precision:	nearest second		

Decimal Latitude: 35.37333    Decimal Longitude: -118.84068    Maximum Error Distance: 9.930 mi    **Calculate**

degrees minutes seconds    nearest second    1 mi    35.37333    -118.84068    (NAD27) North American

**http://www.manisnet.org/gc.html**

Georef Calculator

# Locality Types- 1

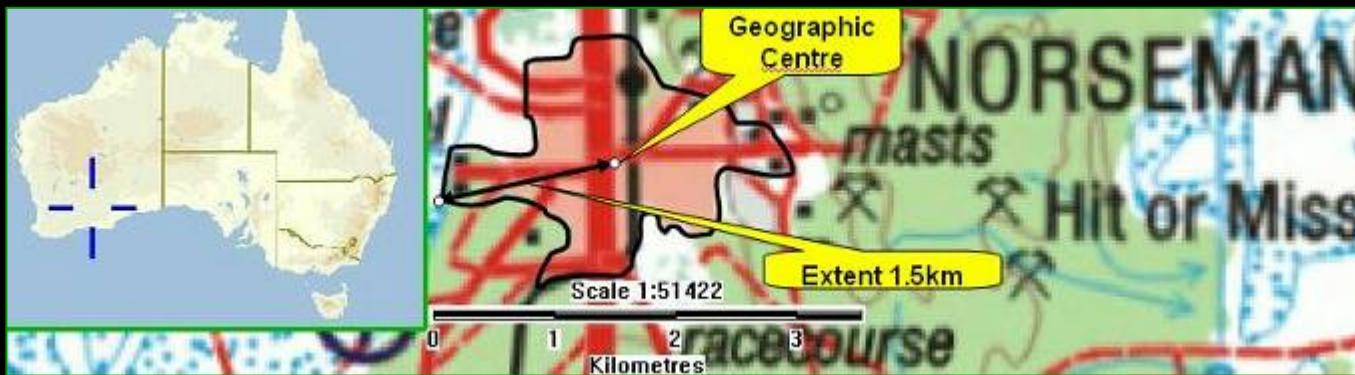
## Named Place

- » Example 1: “Point Lookout”
- » Example 2: “Bennetts Waterhole”
- » Example 3: “Isla Tiburon”
- » Example 4: “Lorne Reef”
- » Example 5: “Mt Hypipamee”

**Georeference:** Use centre of named place

**Extent:** Use the distance from the coordinates of the named place to the furthest point within the named place

**Uncertainty:** Use the *MaNIS Georeference Calculator*



# Example

## **Locality: “Bakersfield”**

Suppose the coordinates for Bakersfield came from the GNIS database (a gazetteer) and the distance from the center of Bakersfield to the furthest city limit is 3 km.

**Coordinate System:** degrees minutes seconds

**Latitude:** 35° 22' 24" N

**Longitude:** 119 ° 1' 4" W

**Datum:** not recorded; 79 m uncertainty

**Coordinate Precision:** nearest second; 40 m uncertainty

**Coordinate Source:** gazetteer

**Extent of Named Place:** 3 km

**Distance Units:** km

**Decimal Latitude:** 35.37333

**Decimal Longitude:** -119.01778

**Maximum Uncertainty Distance:** 3.119 km

# Locality Types - 2

## Between two Named Places

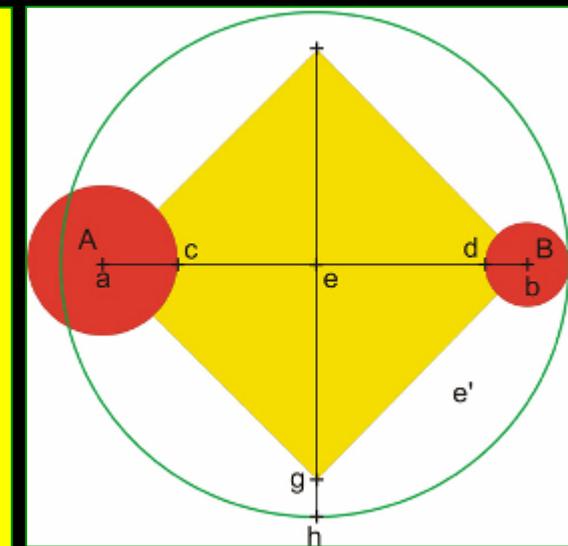
» Example 1: "between Point Reyes and Inverness"

**Georeference:** Find the coordinates of the midpoint between the centres of the two named places (**e**)

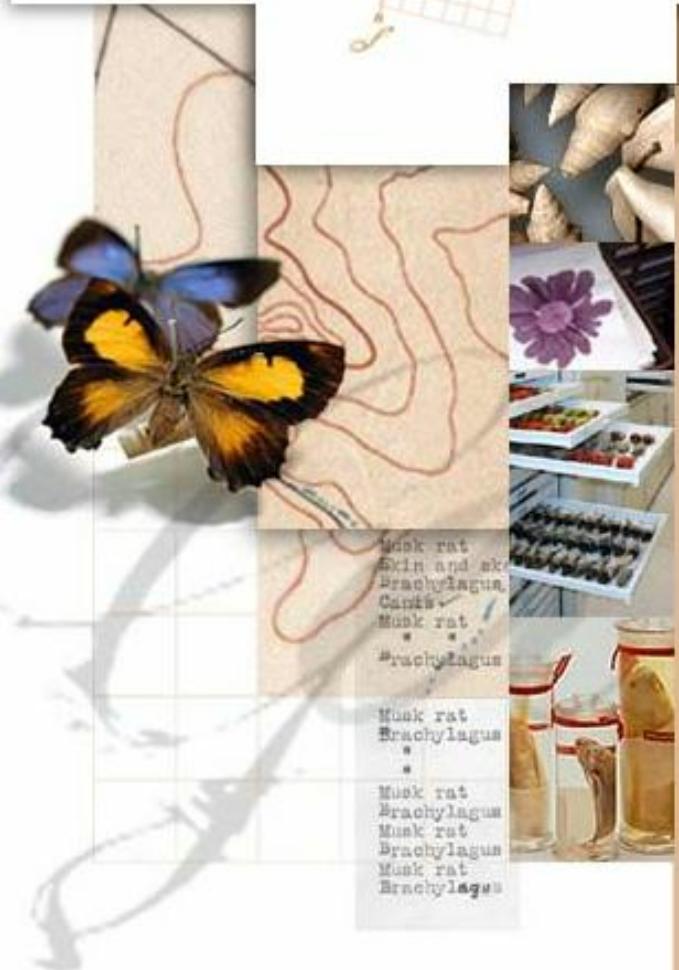
**Extent:** Use the extent of A or B, whichever is greater, plus one-half the distance between the centres of A and B.

**Uncertainty:** Use the **MaNIS Georeference Calculator**

Calculate the same as for '**Named Place**'.



# Bio Geomancer



## APPLICATIONS

[Home](#)

## STANDARDS

[About](#)

## LIBRARY

[Feedback](#)

## NEWS

### WHAT IS THE BIOGEOMANCER PROJECT?

The BioGeomancer (BG) Project is a worldwide collaboration of natural history and geospatial data experts. The primary goal of the project is to maximize the quality and quantity of biodiversity data that can be mapped in support of scientific research, planning, conservation, and management. The project promotes discussion, manages geospatial data and data standards, and develops software tools in support of this mission. [Learn more about us!](#)

### WHAT IS GEOREFERENCING?

Georeferencing is the process of converting text descriptions of locations to computer-readable geographic locations, such as a GIS system uses. [More about georeferencing...](#)

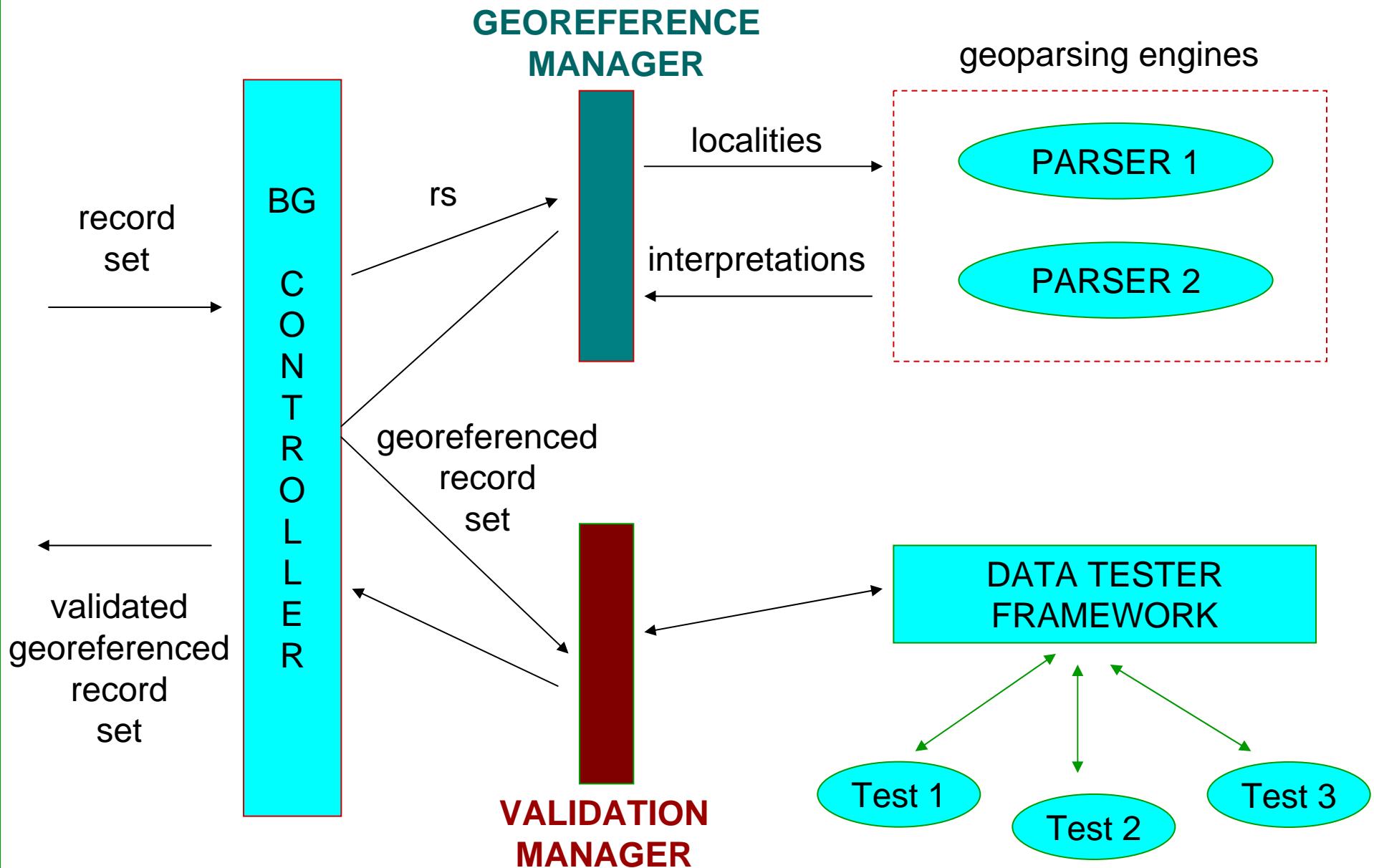
### LATEST BIOGEOMANCER NEWS

BioGeomancer features in August edition of GEOWORLD!

[More News...](#)

[www.biogeomancer.org](http://www.biogeomancer.org)

# BioGeomancer – Georeferencing



# BioGeomancer Workbench

8 km NW of Leige, Belg

Georeference [Submit issue](#)

Map Satellite Hybrid

Powered by Google

2000 mi  
2000 km

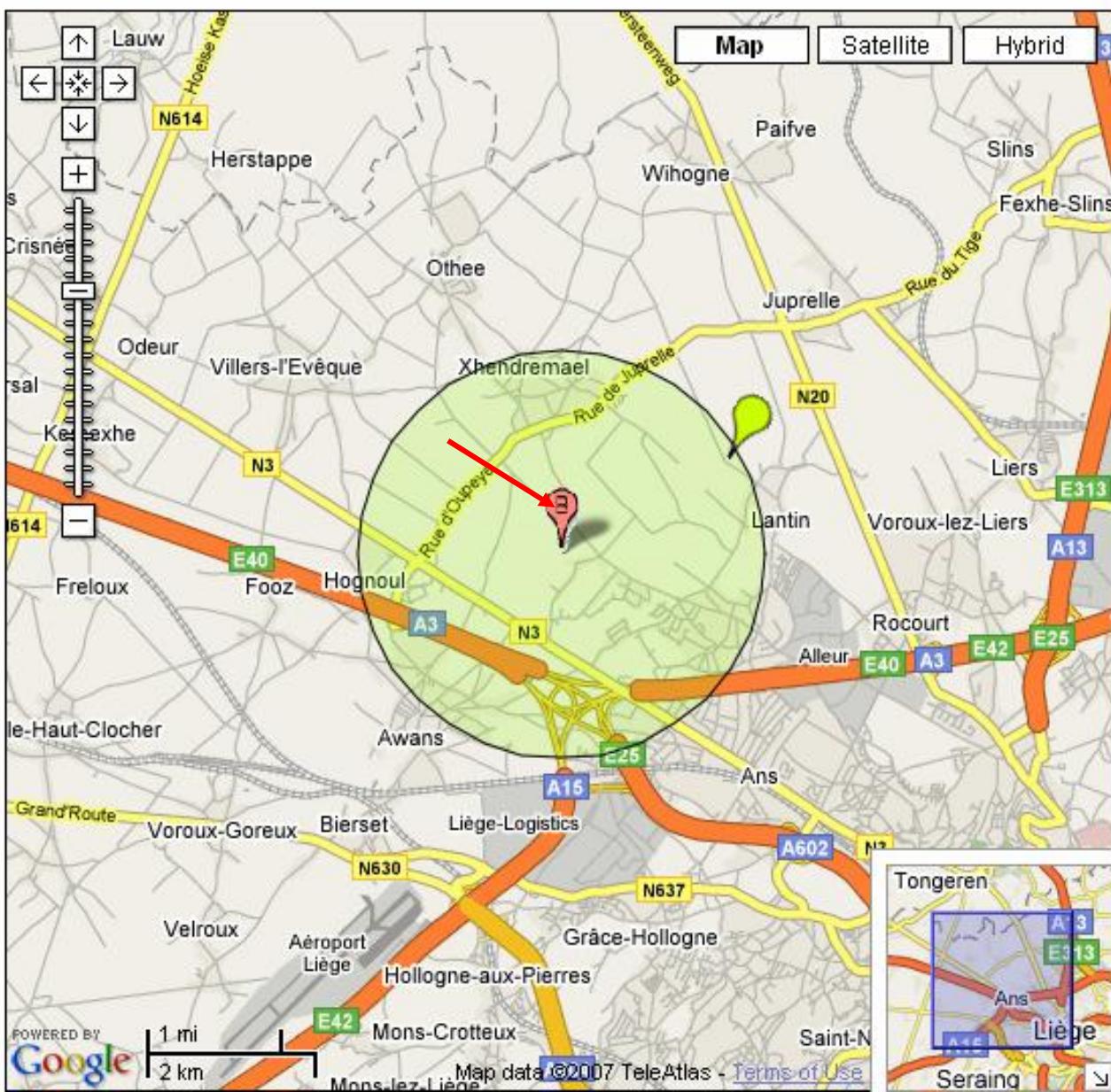
Southern Ocean

Terms of Use [View](#)

Australian Biodiversity Information Services



Belgium; 8 km NW of Liège  
50.6841853, 5.4867103  
2447 meter uncertainty



# BioGeomancer Workbench – Batch Georeferencing

- Need to register and Log in
- Submit a project
  - Use the “Help” to see how to go about submitting files
    - (See next slide for formatting)
    - NB – files must be in UTF-8
- Results downloaded as XML file

Records	Project ID	Project name	Download	Create Features Project	Delete
<a href="#">32 records</a>	5629	Training-NL			

# Format for batch georeferencing

## How do I format a file for batch georeferencing?

Save your data in a tab-delimited text file. The first row of the must contain column names. BioGeomancer can understand concept names from the [DarwinCore 1.4](#) and its [Geospatial Extension](#).

Specifically, BioGeomancer currently interprets the following fields (the order and case are not relevant):

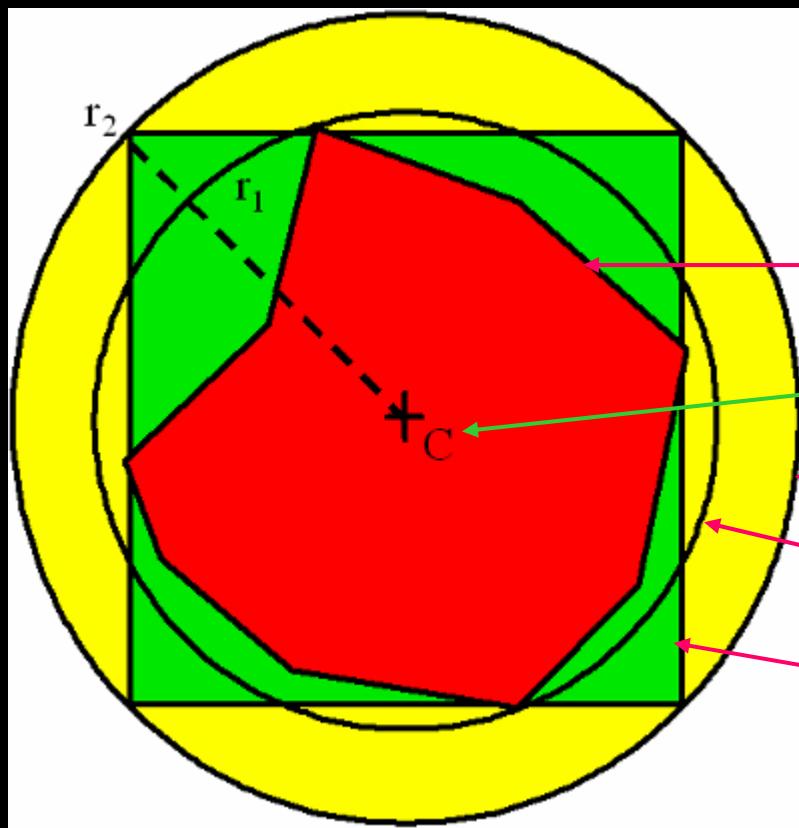
- Locality
- HigherGeography
- Continent
- WaterBody
- IslandGroup
- Island
- Country
- StateProvince
- County
- VerbatimLatitude
- VerbatimLongitude
- VerbatimCoordinates
- VerbatimElevation

Any other fields submitted in the upload are retained, but uninterpreted. Of those in the list above, the Verbatim fields are currently unused in the spatial description, and the WaterBody, IslandGroup, and Island are likely to be heavily under-represented in the gazetteer.

Pay attention to the encoding system used. You should use [UTF-8](#), a character encoding system for the [Unicode](#) standard to represent any of the world's scripts. Text encoded in the standard windows/English "[Latin-1](#)" will not be corrupted if it contains pure [ASCII](#) characters.

# Spatial Fit

*A measure of how well the geometric representation matches the original spatial representation.*



For an area where the original spatial representation of a locality is the **red polygon** with area ' $A$ '. The spatial fit is:

1.0

0

$(\pi * r_2^2)/A$

$(\pi * r_1^2)/A$

$(2 * r_2^2)/A$

# Challenges

- Improving Data Quality



# Users need quality information

So what do we mean by 'Data Quality'?

*An essential or distinguishing characteristic necessary for [spatial] data to be fit for use.*

SDTS 02/92

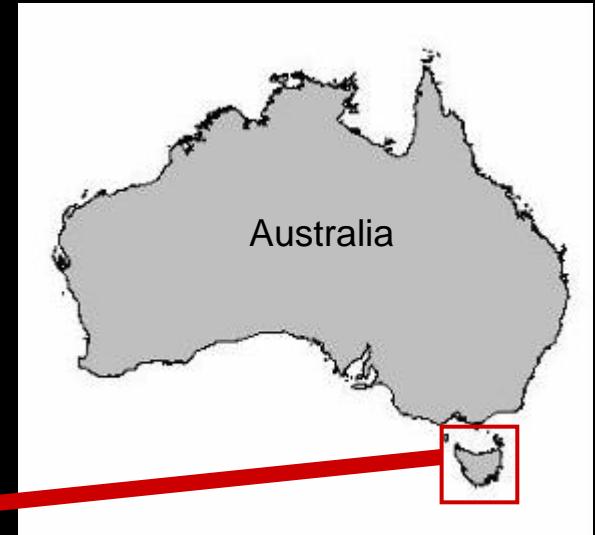
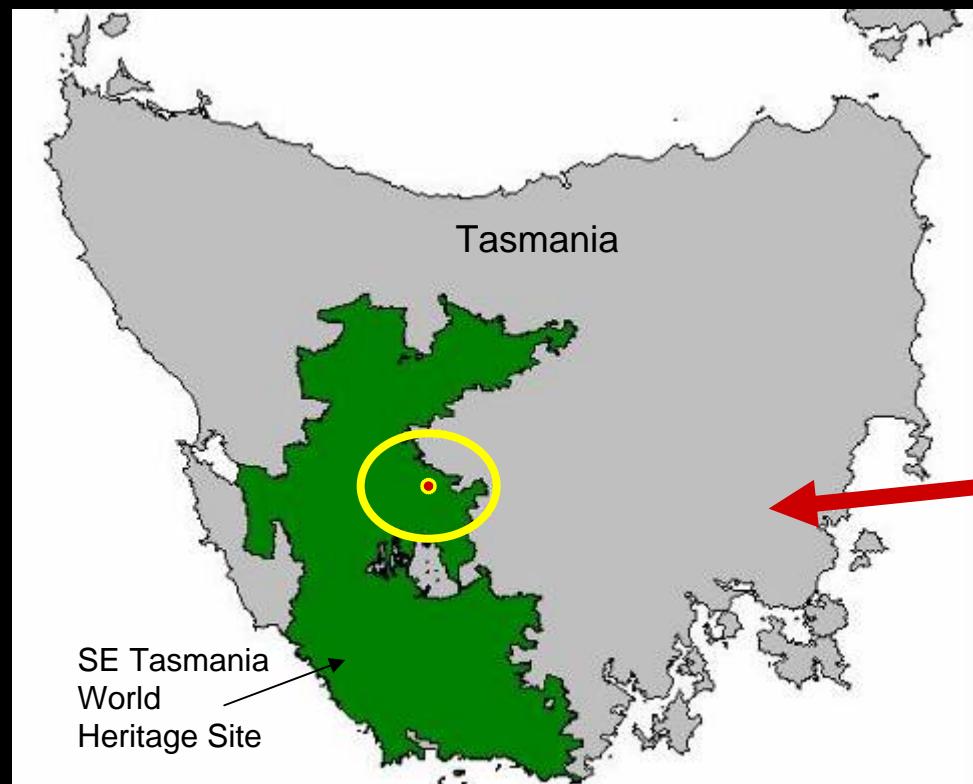
*The general intent of describing the quality of a particular dataset or record is to describe the fitness of that dataset or record for a particular use that one may have in mind for the data.* (Chrisman 1991)

# Data quality - fitness for use?

Fitness for use

Does species 'A' occur in Tasmania?

Does species 'A' occur in National Park 'y'



# Loss of data quality

Loss of data quality can occur at many stages:

- At the time of collection
- During digitization
- During documentation
- During storage and archiving
- During analysis and manipulation
- At time of presentation
- And through the use to which they are put

*Don't underestimate the simple elegance of quality improvement. Other than teamwork, training, and discipline, it requires no special skills. Anyone who wants to can be an effective contributor.*

(Redman 2001).

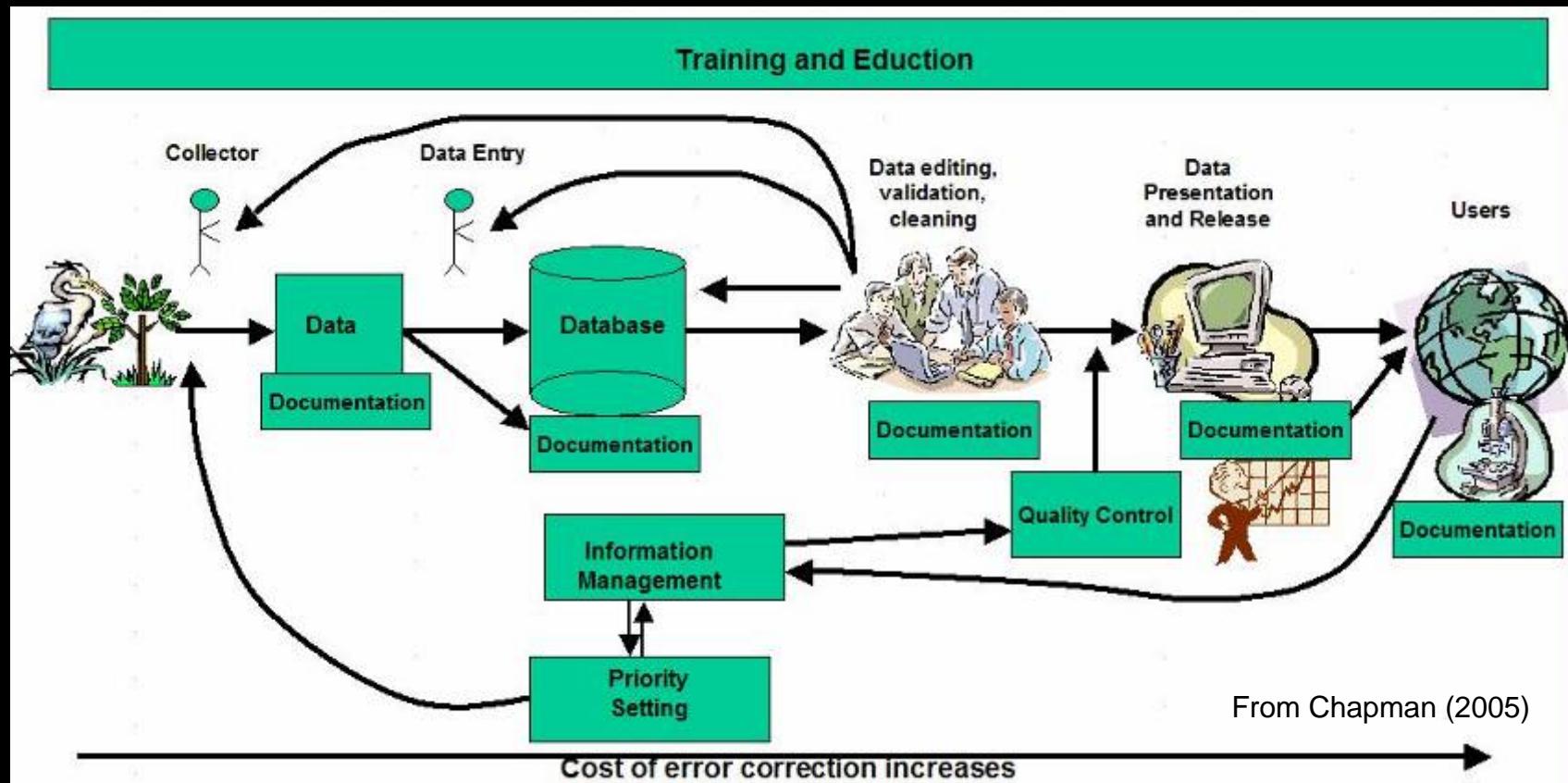
# Principles of data quality

It is important for organizations to have  
a vision with respect to having good quality data.  
a policy to implement that vision, and  
a strategy for implementation

*Experience has shown that treating data as a long-term asset and managing it within a coordinated framework produces considerable savings and ongoing value.*

(NLWRA 2003).

# Data Quality Information Chain



*Assign responsibility for the quality of data to those who create them. If this is not possible, assign responsibility as close to data creation as possible*

(Redman 2001)

# Data validation

Two key sources of error  
are:

- Taxonomic names
- Georeferences (lat's and long's)

Methods for identifying error

Documented here ----->

available via GBIF web site

<http://www.gbif.org>



PRIMARY SPECIES AND  
SPECIES-OCCURRENCE  
DATA

Arthur D. Chapman<sup>1</sup>

*Error qui non resistit, approbatur.*

An error not resisted is approved.

(Ref. *Doct. & Stud. c. 770*).





A. odorante. — A. suaveolens.

à rameaux triangulaires et très glabre. — Feuilles longues, un peu obtuses au sommet, tranchées facilement au milieu, avec une ligne de pétiole sur les bord, dont

**Australian Plant Name Index**

**A-C**

ngue, relevées diagonales indiquant la direction



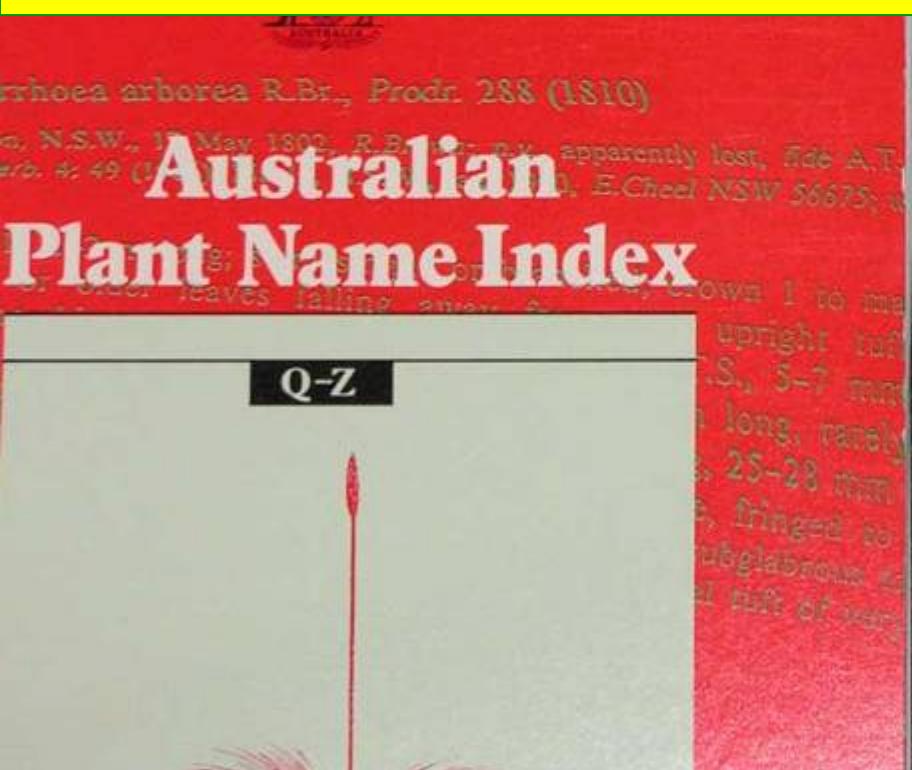
Gesneriaceae fragilis Labill.  
Tridaceae d. Bot.  
D. Gosselin Nov.  
**Australian Plant Name Index**

**K-P**

droide

04

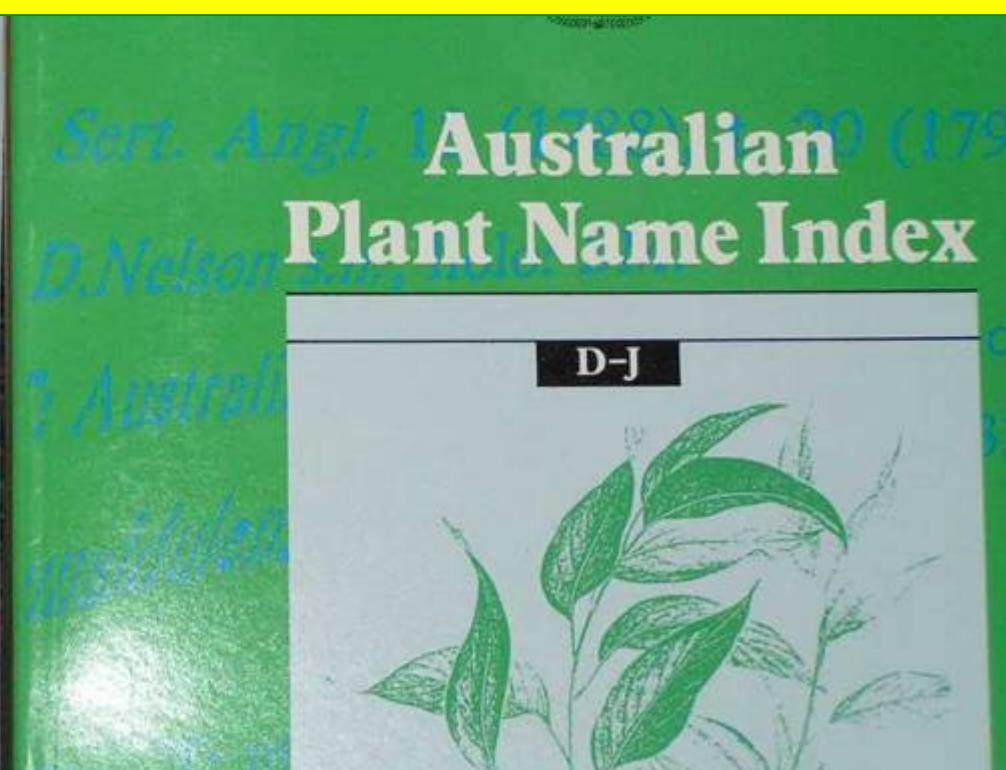
# 1. Taxonomic and Nomenclature Data



**Q-Z**



Chapman (1991)



**D-J**

<http://www.cpbr.gov.au/cpbr/databases/apni.html>

# Taxonomic Data

Consists of: (not all are always present):

- Name (scientific, common, hierarchy, rank)
- Nomenclatural status (synonym, accepted, typification)
- Reference (author, place and date of publication)
- Determination (by whom and when the record was identified)
- Type specimen citation
- Quality fields (accuracy of determination, qualifiers)

# Documenting Taxonomic Data Quality

- Several methods exist for documenting taxonomic verification - none are completely satisfactory
  - Herbarium Information Standards and Protocols for the Interchange of Data (HISPID)
  - Australian National Fish Collection (1993)
  - Several others restricted to one or two institutions
- Proposal – four level:
  - Who determined the specimen and when
  - What was used (type specimen, local flora, monograph, etc.)
  - Level of expertise of the determiner
  - What confidence did the determiner have in the determination.

# Taxon Verification Status - proposed

Name of determinor:

Date of determination:

Source of determination: (e.g. compared with holotype, used national flora)

- identified by **World expert** in the taxon with **high certainty**
- identified by **World expert** in the taxon with **reasonable certainty**
- identified by **World expert** in the taxon with **some doubt**
- identified by **regional expert** in the taxon with **high certainty**
- identified by **regional expert** in the taxon with **reasonable certainty**
- identified by **regional expert** in the taxon with **some doubt**
- identified by **non-expert** in the taxa taxon **high certainty**
- identified by **non-expert** in the taxa taxon **reasonable certainty**
- identified by **non-expert** in the taxa taxon **some doubt**
- identified by **the collector** with **high certainty**
- identified by **the collector** with **reasonable certainty**
- identified by **the collector** with **some doubt**.

From: Chapman (2005) Principles of Data Quality. GBIF

# Error checking

- **Missing Data Values**
  - empty fields where values should occur  
(e.g. if a species epithet is present, then a generic name **MUST** be present)

# Error checking

- **Incorrect Data Values**

- typographic errors,
- transposition of key strokes,
- data entered in the wrong place

(e.g. a species epithet present in a generic name field)

Can often be identified using Soundex/Phonex techniques

# Error checking

- **Nonatomic Data Values**
  - More than one fact entered into a single field  
(e.g. a species bionomial or trinomial present in a single field)

# Error checking

- **Domain Schizophrenia**
  - Fields used for purposes for which they weren't intended

Family	Genus	Species
Myrtaceae	Eucalyptus	globulus?
Myrtaceae	Eucalyptus	? globulus
Myrtaceae	Eucalyptus	aff. globulus
Myrtaceae	Eucalyptus	sp. nov.
Myrtaceae	Eucalyptus	?
Myrtaceae	Eucalyptus	sp. 1
Myrtaceae	Eucalyptus	To be determined

e.g.

Good reference:

Dalcin, E.C. 2004. *Data Quality Concepts and Techniques Applied to Taxonomic Databases*.

Thesis for the degree of Doctor of Philosophy, School of Biological Sciences, Faculty of Medicine, Health and Life Sciences, University of Southampton. November 2004. 266 pp.

# CRIA Data Cleaning

*species* link

<http://splink.cria.org.br/dc>

data & tools

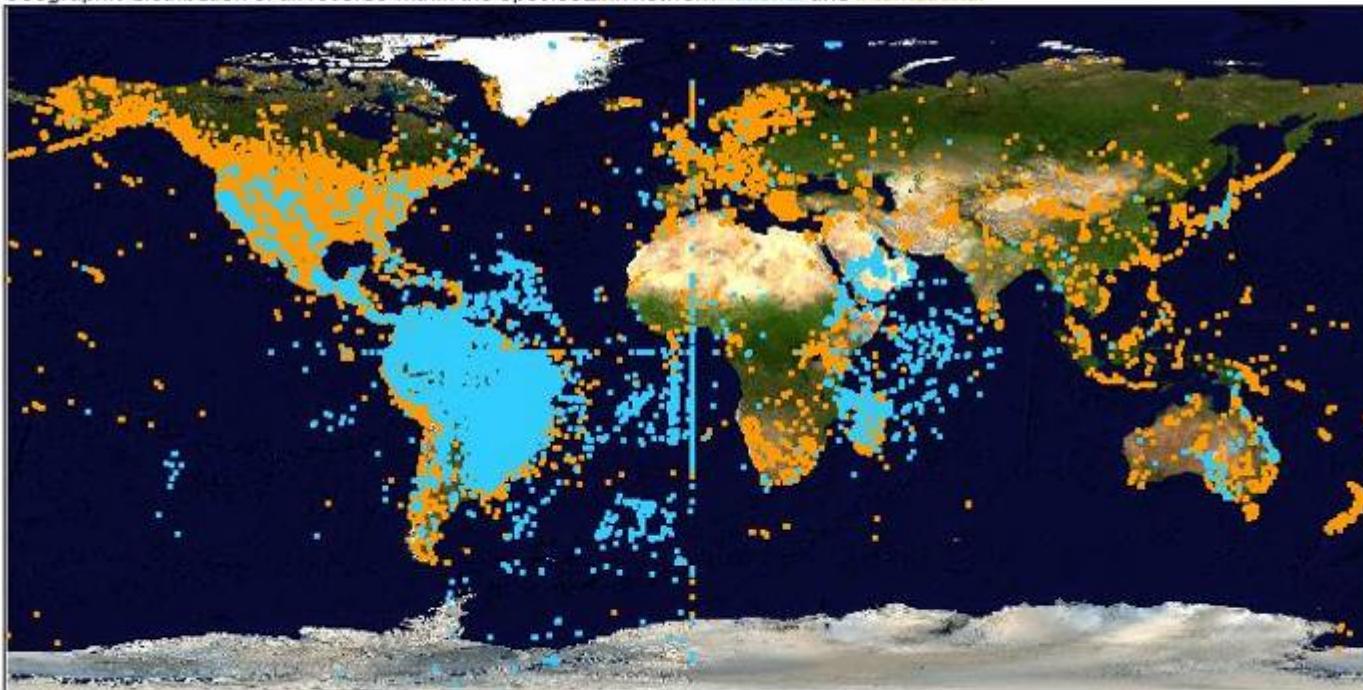
data cleaning

This tool aims at helping curators in identifying possible errors and to standardize data. Records are not modified. The system just presents "suspect" records, recommending that they be checked by each author or curator. The tool is under constant development, so any suggestion is more than welcome.

Select a collection

HSJRP

Geographic distribution of all records within the speciesLink network [national](#) and [international](#)



graphic representation of families

graphic representation of Brazilian states

origin of the records

main collectors

collection events by year

general graphic of all national collections

general graphic of all international collections

# CRIA Data Cleaning

taxonomic data	
inventory	scientific name - collected
family	397 suspect rec
genus	355 suspect rec
species	273 suspect rec
subspecies	not found
author	161 suspect rec
duplicate	618 suspect rec

family	genus
[Amaranthaceae]	<i>sp</i> [Alternanthera]
[Amaranthaceae]	<i>sp</i> [Alternanthera]

Diagram illustrating the CRIA Data Cleaning process flow:

```
graph TD; A[inventory, family, genus, species, subspecies, author, duplicate] --> B[Species 2000]; A --> C[RECURSOS]; C --> D[Species 2000]; C --> E[Externos]; E --> F[Species 2000]; E --> G[Externos]; G --> H[Species 2000];
```

The process starts with taxonomic data (inventory, family, genus, species, subspecies, author, duplicate) which feeds into both the Species 2000 search and the internal CRIA resources.

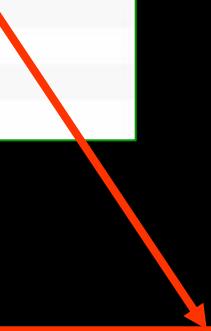
The CRIA resources section contains links to various databases and services, categorized into NO CRIA and EXTERNOS:

- NO CRIA:**
  - Atlas Biota (checked)
  - Banco de Imagens (unchecked)
  - Bioline International (unchecked)
  - Biota Neotropica (unchecked)
  - Neofrug (unchecked)
  - SICol (unchecked)
  - SinBiota (checked)
  - speciesLink (checked)
- EXTERNOS:**
  - Biblioteca Digital da Unicamp (checked)
  - Botanical Type (unchecked)
  - Specimens at US (unchecked)
  - Google Images (checked)
  - IPIII (checked)
  - ITIS (checked)
  - NCBI GenBank (unchecked)
  - NCBI PubMed (checked)
  - HYBG (checked)
  - SciElo Brasil (checked)
  - Taxonomic Name Server (checked)
  - W3Tropicos (unchecked)

The Species 2000 search results for *Alternanthera brasiliiana* (L.) Kuntze are displayed, showing its classification under the family Amaranthaceae. The result is marked as the accepted name.

# CRIA Data Cleaning

taxonomic data	
inventory	scientific name - collector - types
family	397 suspect records
genus	355 suspect records
species	273 suspect records
subspecies	not found
author	161 suspect records
duplicate	618 suspect records



SP	[Inga]	[cylindrica]	[]	2	<a href="#">see</a>	51	accepted name	32009 64192
SP	[Inga]	[cilindrica]	[]	1	<a href="#">see</a>	2		50335
SP	[Inga]	[cilyndrica]	[]	0		1		
SP	[Ipomoea]	[regnellii]	[]	2	<a href="#">see</a>	7		29867 31014
SP	[Ipomoea]	[regmnelli]	[]	1	<a href="#">see</a>	1		115815

# CRIA Data Cleaning

taxonomic data	
inventory	scientific name - collector - types
family	907 suspect records
genus	1074 suspect records
specie	935 suspect records
subspecie	not found
author	3060 suspect records
duplicate	197 suspect records



genus	species	subspecies	author	ocor_col	ocor_total
sp [Acacia]	[martusiana]	[]	[Burk.]	3	3
sp [Acacia]	[martusiana]	[]	[(Steud.) Burk.]	1	1
sp [Acacia]	[martusiana]	[]	[(Steud.) Burk.]	0	2
sp [Actinostemon]	[concolor]	[]	[Müll.Arg.]	0	1
sp [Actinostemon]	[concolor]	[]	[(Spreng.) Muell.Arg.]	0	22
sp [Actinostemon]	[concolor]	[]	[(Spreng) Muell Arg]	45	45
sp [Actinostemon]	[concolor]	[]	[(Spreng.) Mull. Arg.]	1	2
sp [Actinostemon]	[concolor]	[]	[(spreng.) Müll. Arg.]	0	2
sp [Actinostemon]	[concolor]	[]	[(Spreng.) Müll. Arg.]	0	23
sp [Actinostemon]	[concolor]	[]	[(Spreng.) Müller.Arg.]	0	65
sp [Actinostemon]	[concolor]	[]	[(Spreng) Müller. Arg.]	0	6
sp [Actinostemon]	[concolor]	[]	[(spr.) Muell. Arg.]	0	1
sp [Actinostemon]	[concolor]	[]	[(Spr.) Muell.Arg.]	0	2

# CRIA Data Cleaning

taxonomic data	
inventory	scientific name - collector - types
family	907 suspect records
genus	1074 suspect records
specie	935 suspect records
subspecie	not found
author	3060 suspect records
duplicate	197 suspect records



collector	collector number	collectioncode	genus	species	subspecies	identified by	ocor_col
Aguiar, O.T.	193	ESA	sp Sapium	glandulatum		Cordeiro, I.	1
Aguiar, O. T.	193	UEC	sp Sapium	longifolium		J. A Pastore	1
Aguiar, O.T.	193	SP	sp Sapium	glandulatum		Cordeiro, I.	1
Amaral Jr, A	118	UEC	sp Miconia	cubatanensis			1
Amaral Jr, A	118	UEC	sp Miconia	cubatanensis		Goldenberg, R	1
Amaral Jr, A.	118	SPF	sp Picramnia	sellowii	subsp. sellowii		1
Amaral Jr., A.	118	SPF	sp Picramnia	sellowii	subsp. sellowii		1

# CRIA data cleaning

*species* link

data & tools

<http://splink.cria.org.br/dc>

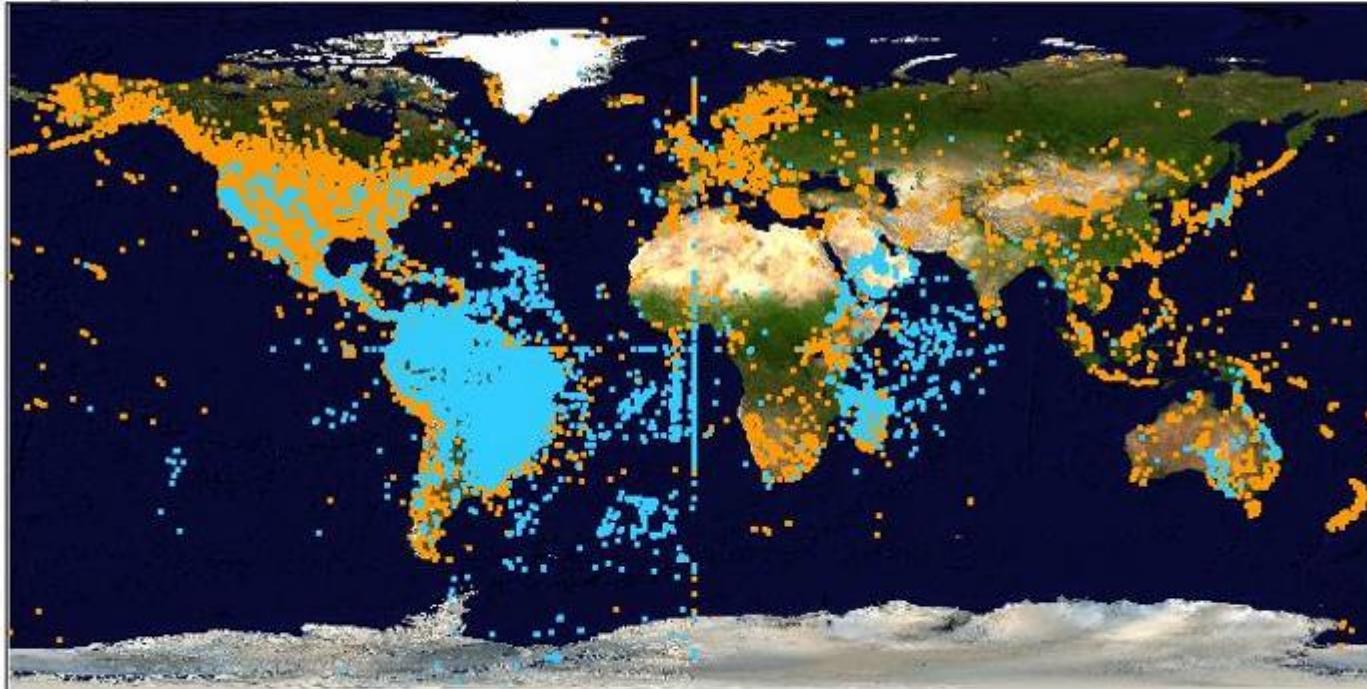
data cleaning

This tool aims at helping curators in identifying possible errors and to standardize data. Records are not modified. The system just presents "suspect" records, recommending that they be checked by each author or curator. The tool is under constant development, so any suggestion is more than welcome.

Select a collection

UEC

Geographic distribution of all records within the speciesLink network [national](#) and [international](#)



graphic representation of families

graphic representation of Brazilian states

origin of the records

main collectors

collection events by year

general graphic of all national collections

general graphic of all international collections

# The Brazil example

<http://splink.cria.org.br/dc>

taxonomic data	
inventory	scientific name - collector
family	not found
genus	337 suspect records
species	957 suspect records
subspecies	not found
author	3450 suspect records
duplicate	not found

family	genus
[Acanthaceae]	sp [Sa]
[Amaranthaceae]	sp [Al]
[Amaranthaceae]	sp [Al]
[Amaranthaceae]	sp [Al]

**Busca Nome Científico**

Digite o nome da espécie procurada  
gênero  espécie  infra-espécie   
[sugestões](#) [help](#) [buscar](#)

**RECURSOS**

**NO CRIA**

- [Atlas Biota](#)
- [Banco de Imagens](#)
- [Bioline International](#)
- [Biota Neotropica](#)
- [Neofrug](#)
- [SiCol](#)
- [SinBiota](#)
- [speciesLink](#)

**EXTERNIOS**

- [Biblioteca Digital da Unicamp](#)
- [Botanical Type Specimens at US](#)
- [Google Images](#)
- [IPIII](#)
- [ITIS](#)
- [NCBI GenBank](#)
- [NCBI PubMed](#)
- [NYBG](#)
- [SciElo Brasil](#)
- [Taxonomic Name Server](#)
- [W3Tropicos](#)

**Species 2000** Informações taxonômicas extraídas do **Catálogo da Vida**, versão 2005.

**Alternanthera brasiliiana** (L.) Kuntze  
reino: *Plantae*, filo: *Magnoliophyta*, classe: *Magnoliopsida*, ordem: *Caryophyllales*, família: *Amaranthaceae*

**Lista de Espécies**

**sp** *Alternanthera brasiliiana* (L.) Kuntze  
nome aceito

# The Brazil example

<http://splink.cria.org.br/dc>

taxonomic data	
inventory	scientific name - collector - types
family	not found
genus	337 suspect records
species	957 suspect records
subspecies	not found
author	3450 suspect records
duplicate	not found



genus	species	subspecies	ocor_col	ocor_total	status_sp2000
<i>sp</i> [Aechmea]	[ <i>distichanta</i> ]	[]	15	16	
<i>sp</i> [Aechmea]	[ <i>distichantha</i> ]	[]	1	11	
<i>sp</i> [Aeschynomene]	[ <i>brasiliiana</i> ]	[var. <i>brasiliiana</i> ]	13	13	
<i>sp</i> [Aeschynomene]	[ <i>brasiliiana</i> ]	[var. <i>brasiliiana</i> ]	1	1	
<i>sp</i> [Acalypha]	[ <i>gracilis</i> ]	[]	8	48	
<i>sp</i> [Acalynhal]	[ <i>gracillisl</i> ]	[]	8	8	
<i>sp</i> [Acacia]	[ <i>polyphylla</i> ]	[]	82	292	accepted name
<i>sp</i> [Acacia]	[ <i>poliphylia</i> ]	[]	1	1	
<i>sp</i> [Acacia]	[ <i>polyphilla</i> ]	[]	0	1	
<i>sp</i> [Acacia]	[ <i>polyphylla</i> ]	[]	1	1	

# The Brazil example

<http://splink.cria.org.br/dc>

taxonomic data	
inventory	scientific name - collector - types
family	not found
genus	337 suspect records
species	957 suspect records
subspecies	not found
author	3450 suspect records
duplicate	not found

genus	species	subspecies	author	ocor_col	ocor_total
sp [Acacia]	[martusiana]	[]	[Burk.]	3	3
sp [Acacia]	[martusiana]	[]	[(Steud.) Burk]	1	1
sp [Acacia]	[martusiana]	[]	[(Steud.) Burk.]	0	2
sp [Actinostemon]	[concolor]	[]	[Müll.Arg.]	0	1
sp [Actinostemon]	[concolor]	[]	[(Spreng.) Muell.Arg.]	0	22
sp [Actinostemon]	[concolor]	[]	[(Spreng) Muell Arg]	45	45
sp [Actinostemon]	[concolor]	[]	[(Spreng.) Mull. Arg.]	1	2
sp [Actinostemon]	[concolor]	[]	[(spreng.) Müll. Arg.]	0	2
sp [Actinostemon]	[concolor]	[]	[(Spreng.) Müll. Arg.]	0	23
sp [Actinostemon]	[concolor]	[]	[(Spreng.) Müll.Arg.]	0	65
sp [Actinostemon]	[concolor]	[]	[(Spreng) Müll. Arg.]	0	6
sp [Actinostemon]	[concolor]	[]	[(spr.) Muell. Arg.]	0	1
sp [Actinostemon]	[concolor]	[]	[(Spr.) Muell.Arg.]	0	2

# The Brazil example

<http://splink.cria.org.br/dc>

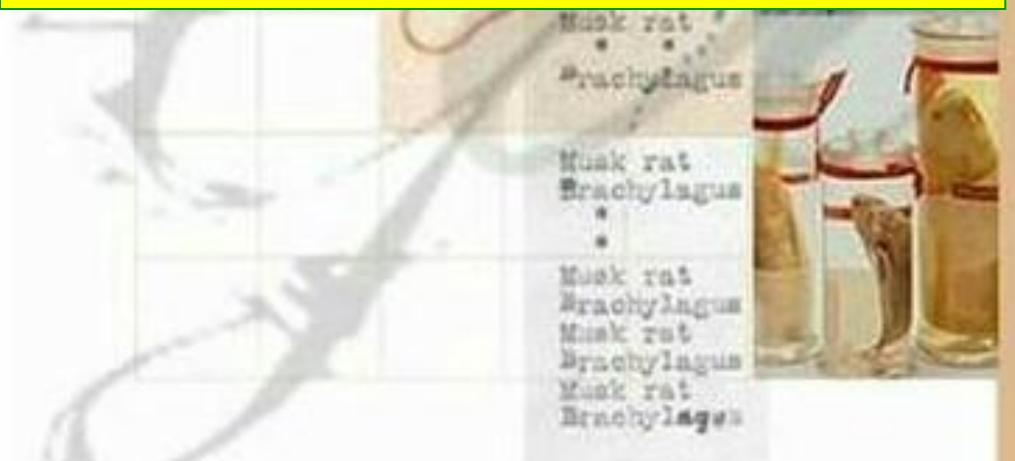
taxonomic data	
inventory	scientific name - collector - types
family	907 suspect records
genus	1074 suspect records
specie	935 suspect records
subspecie	not found
author	3060 suspect records
duplicate	197 suspect records



collector	collector number	collectioncode	genus	species	subspecies	identified by	ocor_col
Aguiar, O.T.	193	ESA	sp Sapium	glandulatum		Cordeiro, I.	1
Aguiar, O. T.	193	UEC	sp Sapium	longifolium		J. A Pastore	1
Aguiar, O.T.	193	SP	sp Sapium	glandulatum		Cordeiro, I.	1
Amaral Jr, A	118	UEC	sp Miconia	cubatanensis			1
Amaral Jr, A	118	UEC	sp Miconia	cubatanensis		Goldenberg, R	1
Amaral Jr, A.	118	SPF	sp Picramnia	sellowii	subsp. sellowii		1
Amaral Jr., A.	118	SPF	sp Picramnia	sellowii	subsp. sellowii		1



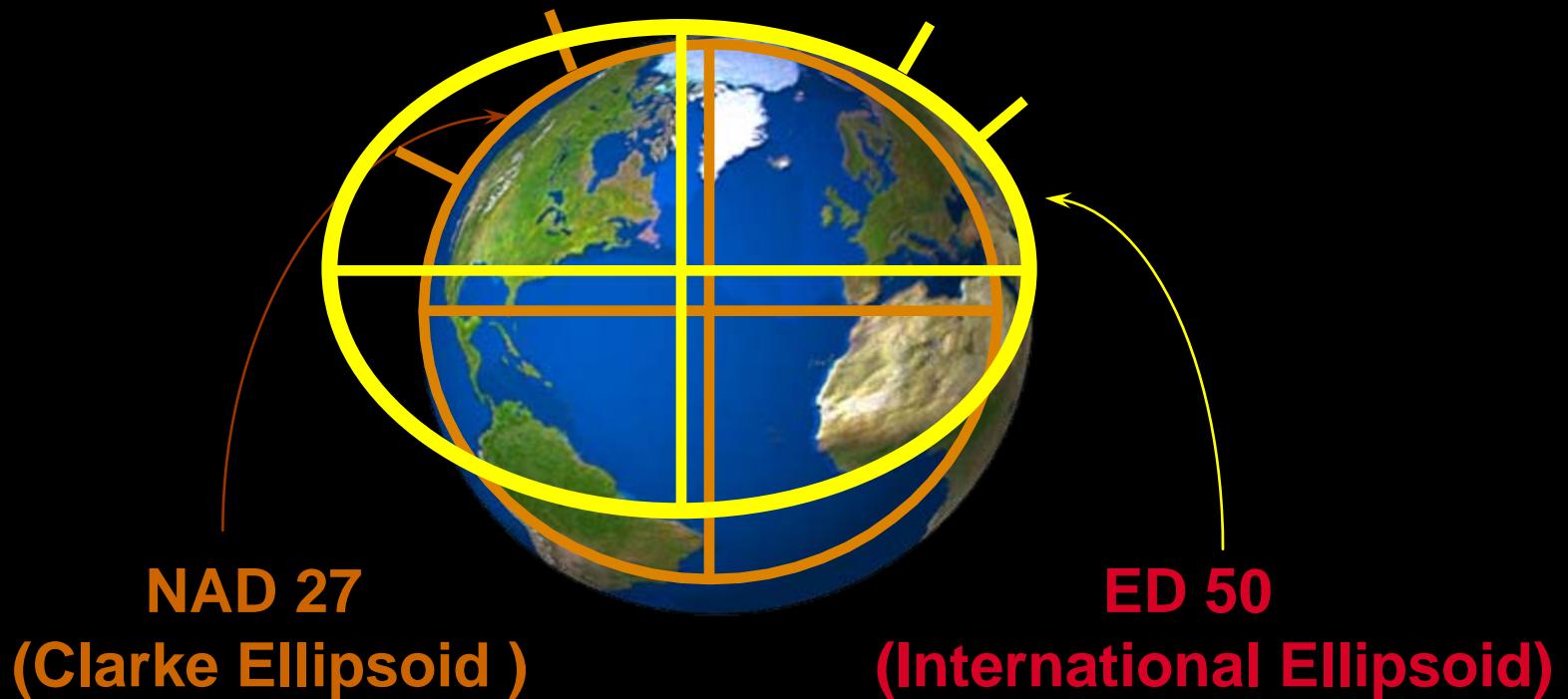
## 2. Geographic Data



# Guide to Best Practices for Georeferencing

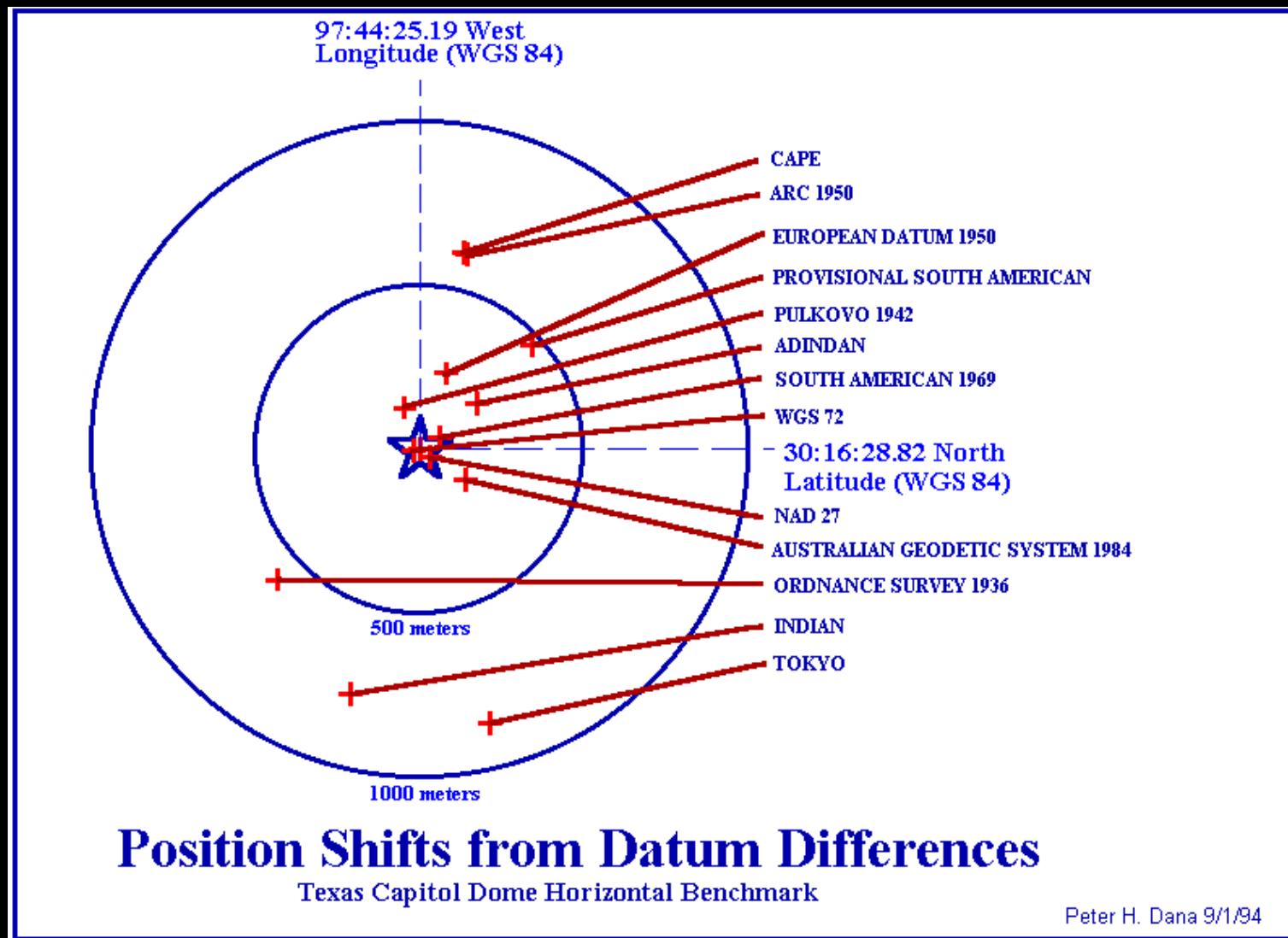
# Traditional Datums

## Traditional Horizontal Datums



From US Navy (n.dat.)

# Datum Shifts



## Position Shifts from Datum Differences

Texas Capitol Dome Horizontal Benchmark

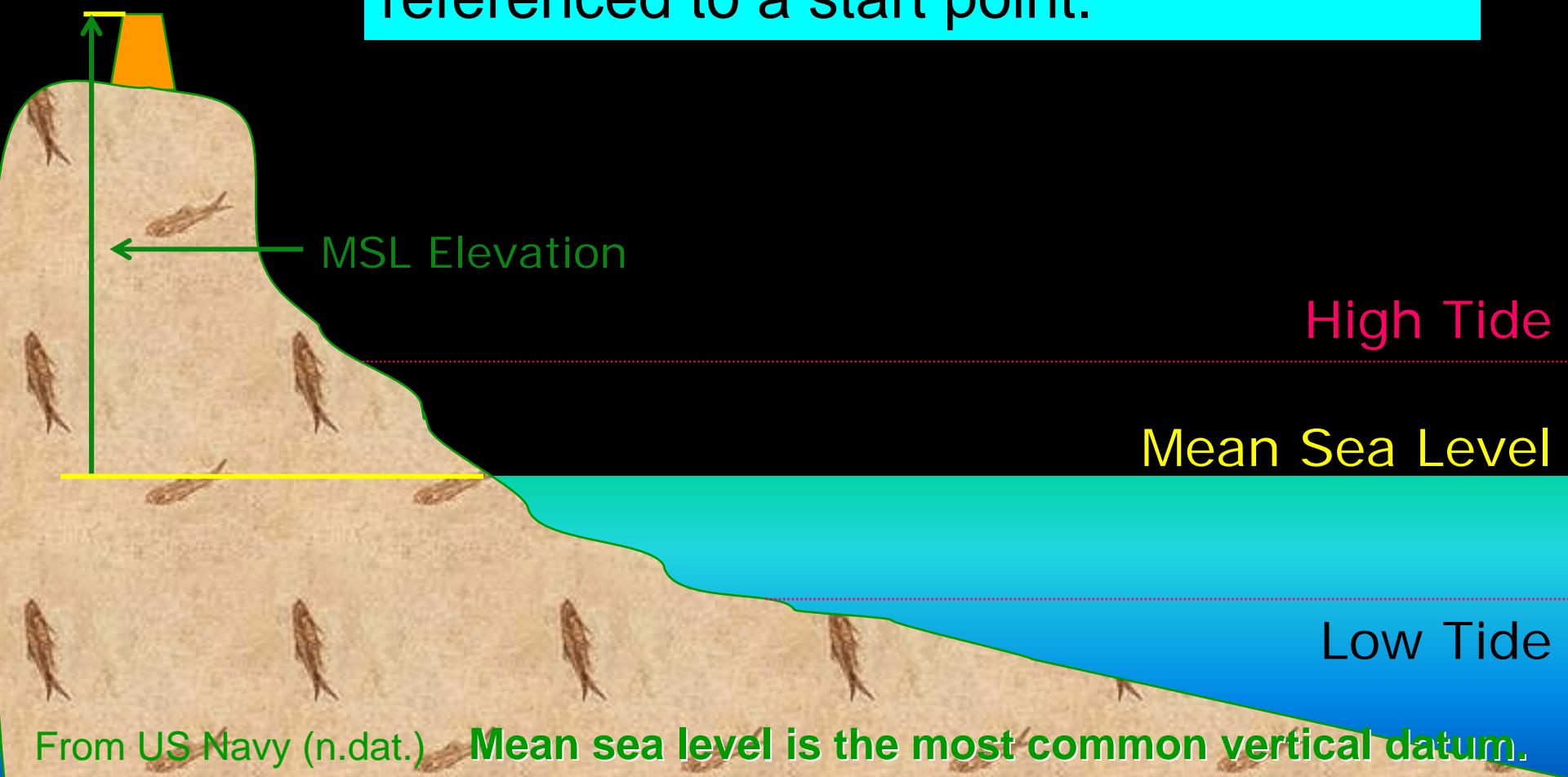
Peter H. Dana 9/1/94

# Differences between Datums

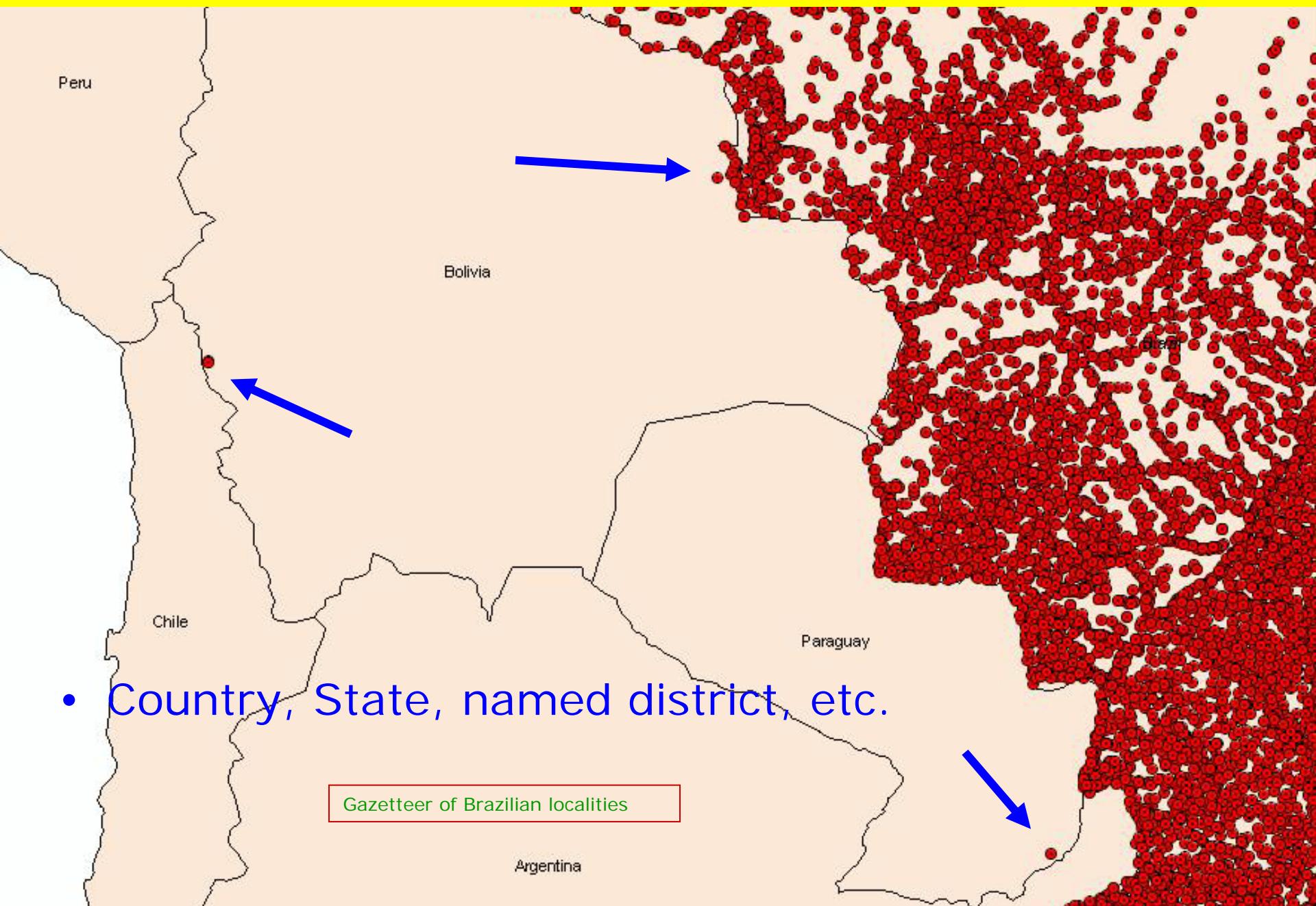
Datum from	Region or Location	Datum to	Difference
AGD66	Australia	AGD84	Max $\pm$ 0-5 m
AGD66/84	Australia	GDA94	Max $\pm$ 200 m
AGD66/84	Australia	WGS84	Max $\pm$ 200 m
GDA94	Australia	WGS84	Max $\pm$ <1 m
NAD 1983	North America	WGS84	Max $\pm$ <1 m
NAD27	North America	WGS84	Max $\pm$ 200 m
NAD 27	Contiguous USA	WGS84	Max $\pm$ 105 m
NAD 27	Aleutian Islands, Alaska	WGS84	Max $\pm$ 235 m
NAD 27	Hawaii	WGS 84	~ 500 m
TOKYO	Japan	WGS84	Max $\pm$ 750 m
ED-50	Europe	WGS84	Max $\pm$ 175 m
ARC-50	Africa	WGS84	Max $\pm$ 265 m
INDIAN 1975	Bangkok, Thailand	WGS84	~ 405 m
INDIAN 1956	Delhi, India	WGS84	~ 135 m
INDIAN 1956	Mumbai, India	WGS84	~ 120 m
HONG KONG 1973	Hong Kong	WGS84	~ 320 m
LUZON	Manila, The Philippines	WGS84	~ 225 m
TOKYO-KOREA	Seoul, South Korea	WGS84	~380 m
KERTAU 1948	Singapore	WGS84	~190 m

# Vertical Datums

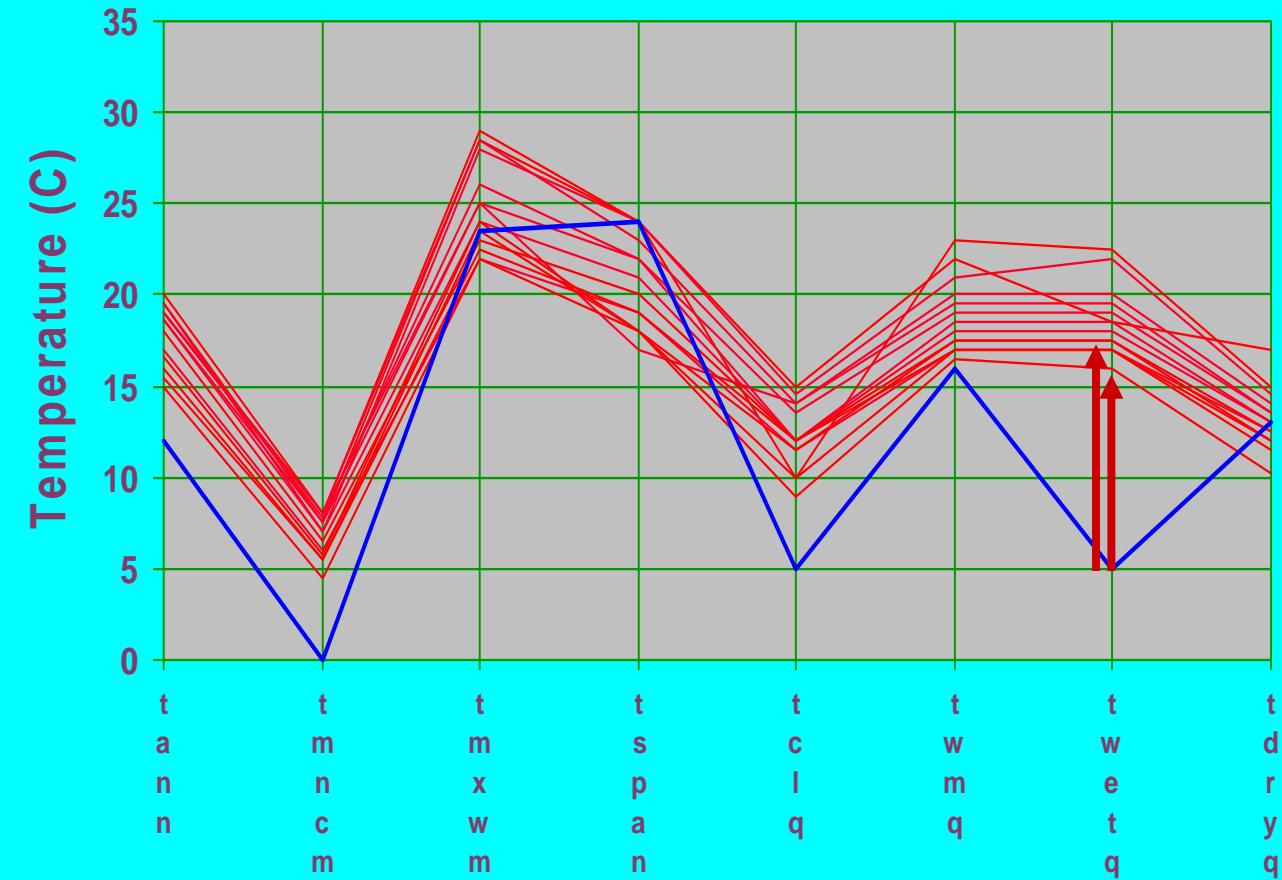
Like horizontal measurements, elevation only has meaning when referenced to a start point.



# Geographic outliers - GIS



# Using Climate to Identify Outliers



Acacia orites - 19 records -  
9 Temperature parameters

$$x < \bar{x}$$

if

$$y_{(i)} = (x_{(i+1)} - x_{(i)}) (\bar{x} - x_{(i)})$$

else

$$y_{(i)} = (x_{(i+1)} - x_{(i)}) (x_{(i+1)} - \bar{x})$$

then

$$C = \sqrt{\frac{\sum_{i=1}^n (y_{(i)} - \bar{y})^2}{n-1}}$$

Reverse Jack-knife

NB. Because the value of 'C' relates to its nearest point, successive values may be very small, so we ensure that if 'x[i]' is an outlier, then all points beyond are outliers too (even if they are clustered)

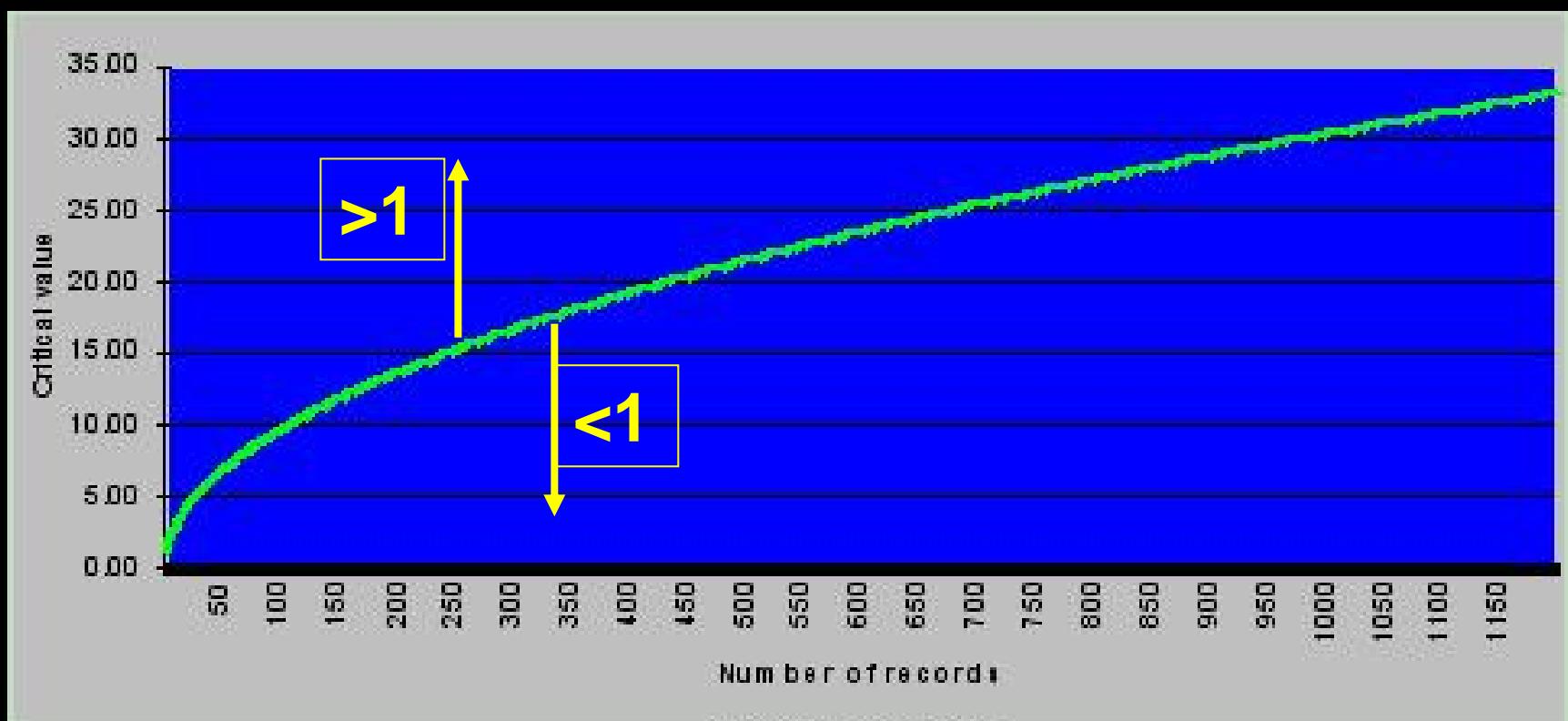
# Concept of "Outlierness"

$$T = ((0.95(\sqrt{n}) + 0.2) \times (\text{Range}/50))$$

where 'n' is the number of records

"Outlierness" is the degree to which a record is an outlier

$$\text{Outlierness} = c[i] / T$$



# *Araucaria bidwillii* (Bunya Pine)



## Outlierness (11/19)

11.0 (Isothermality)  
10.0 (Precipitation Seasonality)  
8.6 (Temp Seasonality)  
8.6 (Precip Wettest Q)  
8.3 (Precip Wettest Month)  
7.7 (Precip Warmest Q)  
5.7 (Ann Precip)  
5.0 (Ann Temp Range)  
4.5 (Mean Monthly Temp Range)  
3.3 (Max Temp Warmest Month)  
2.1 (Min Temp Coldest Month)  
1.1 (Mean Temp Driest Q)

## Outlierness (15/19)

12.0 (Isothermality)  
10.9 (Precipitation Seasonality)  
8.2 (Temp Seasonality)  
6.3 (Mean Temp Driest Q)  
5.8 (Min Temp Coldest Month)  
5.4 (Mean Ann Temp)  
5.2 (Mean Temp Coldest Q)  
5.2 (Precipitation Wettest Q)  
5.1 (Precip Wettest Month)  
4.3 (Precipitation Warmest Q)  
3.8 (Ann Temp Range)  
2.9 (Mean Monthly Temp Range)  
2.1 (Ann Precip)  
1.8 (Precip Driest Month)  
1.6 (Mean Temp Warmest Q)

After validation

# *Araucaria bidwillii* (Bunya Pine)



**Outlierness** (3/19)  
3.6 (Max Temp Warmest Month)  
3.4 (Mean Temp Warmest Q)  
3.1 (Mean Temp Wettest Q)

- Outlierness (6/19)
  - 6.0 (Precip Coldest Q)
  - 4.0 (Precipitation Driest Q)
  - 2.7 (Mean Ann Temp Range)
  - 3.2 (Ann Precip)
  - 2.7 (Mean Temp Driest Q)
  - 1.2 (Isothermality)

# Diva-GIS

- Free
- Simple GIS
- Modelling (BIOCLIM/Domain)
- Data Cleaning Tools



Brown Algae, Argentina

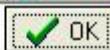
## About DIVA-GIS

DIVA-GIS  
VERSION 3



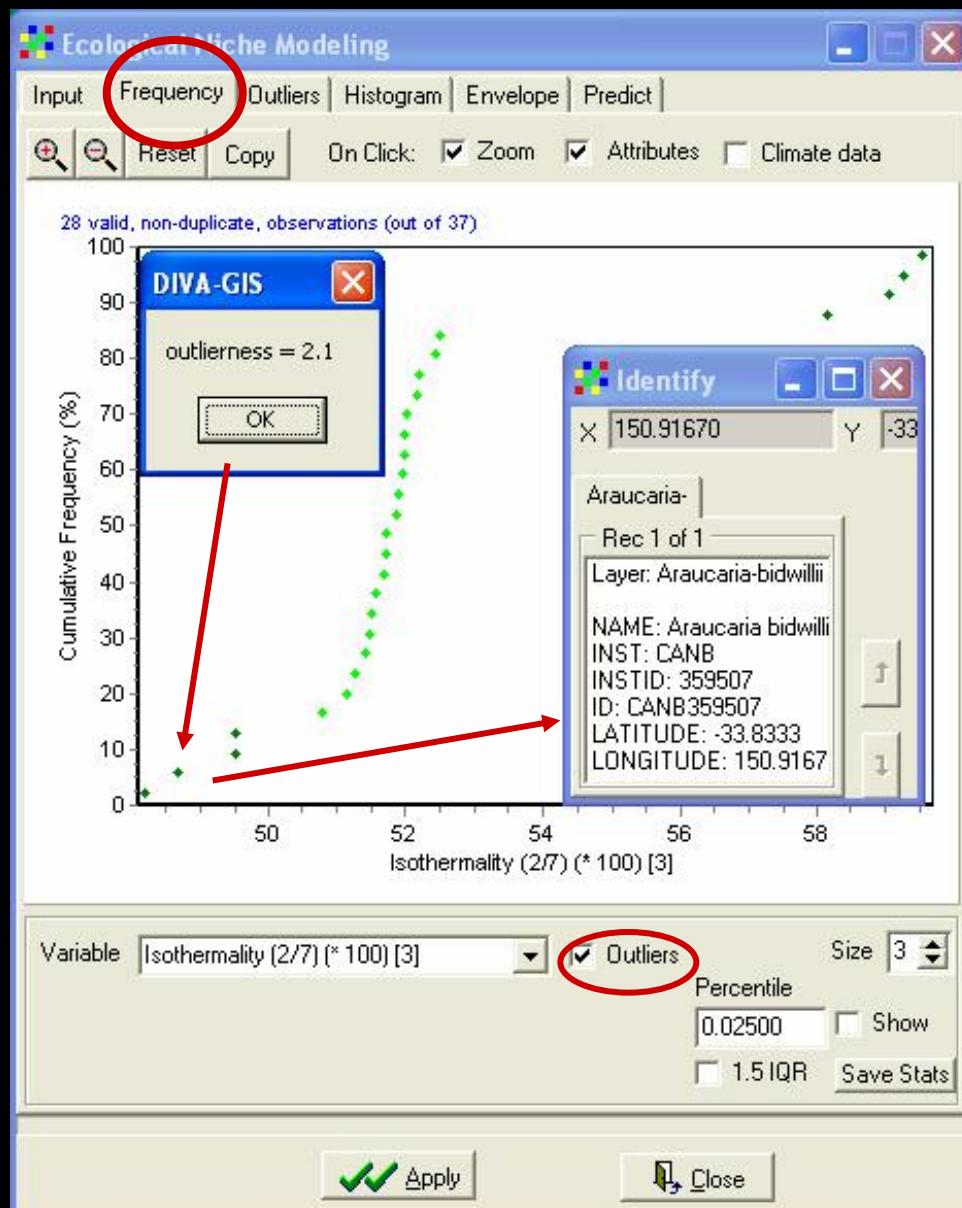
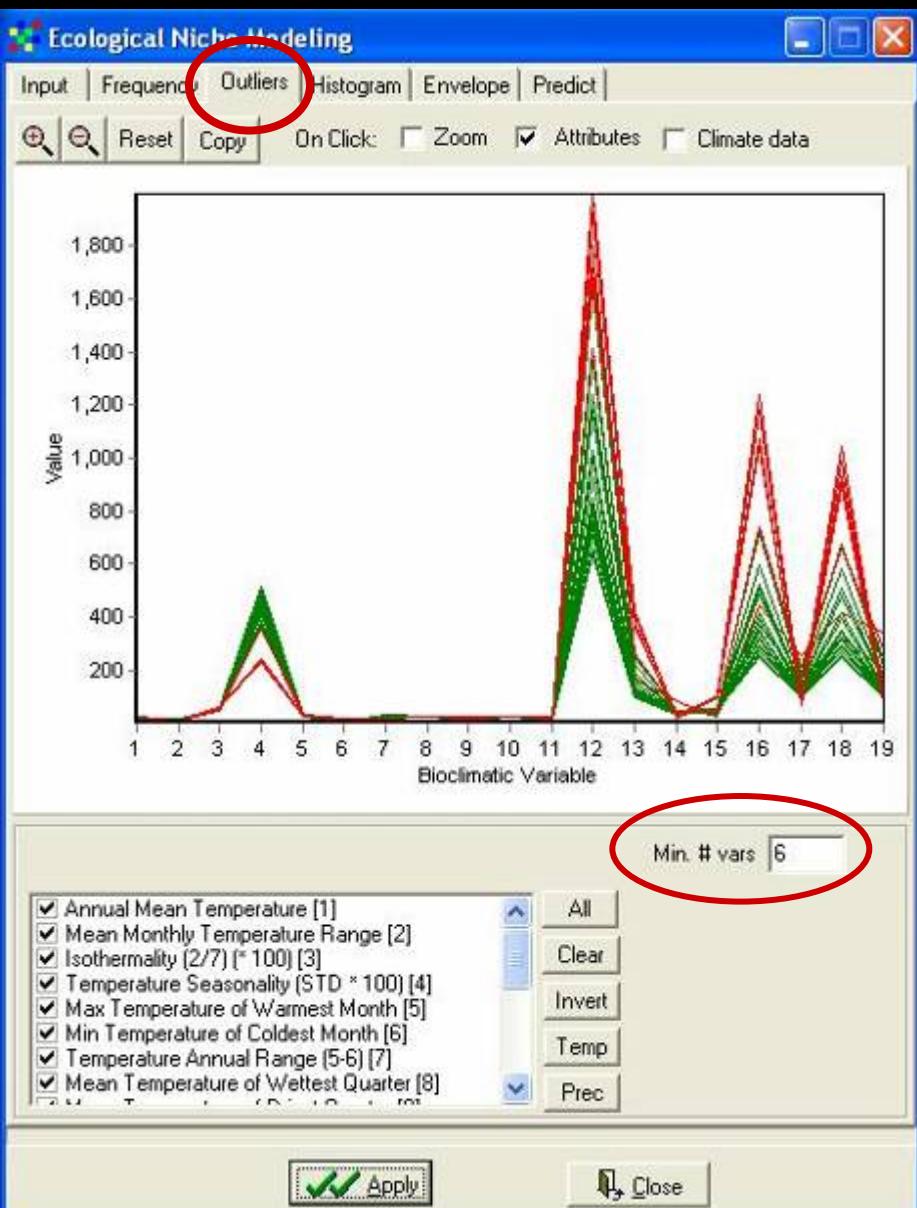
<http://www.diva-gis.org/> Version 3.1

DIVA-GIS was developed by Robert J. Hijmans, Edwin Rojas, Mariana Cruz, Luigi Guarino and Israel Barrantes. The development was partly supported by the International Plant Genetic Resources Institute, the International Potato Center, SINGER/SGRP, the UC Berkeley Museum of Vertebrate Zoology, and FAO, with additional support from USDA, SENASA, BMZ, and ESRI.

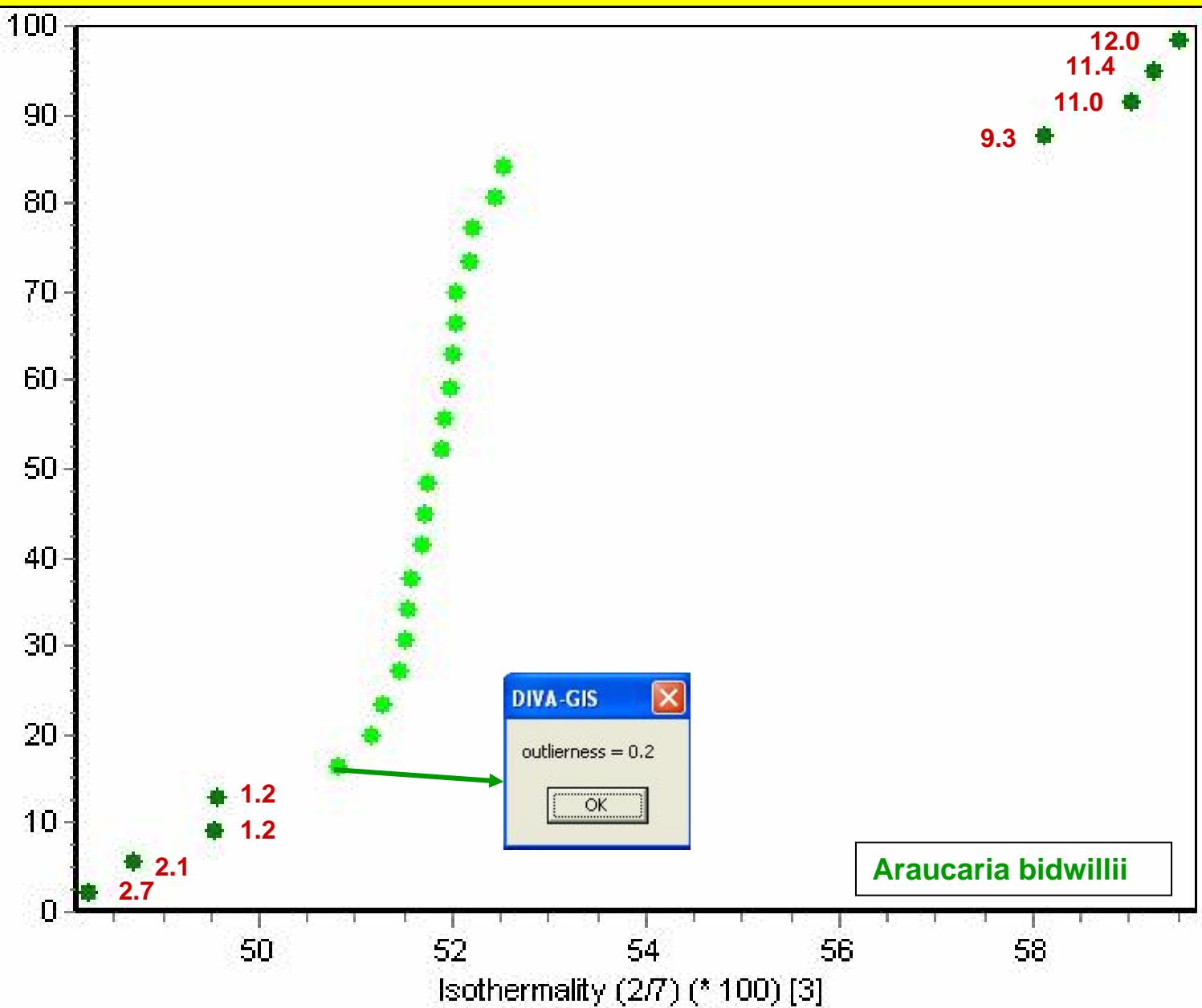


The following persons made additional software available that was used for the development of DIVA-GIS: M. Sawada (Rook's case), Gerald Everden and Frank Warmerdam (PROJ4); Andrew Williamson (Shapechk); the contributors to

# Diva-GIS



# Outlierness values for Isothermality



Cumulative frequency curves can be used to demonstrate the “outlierness” concept

# Errors in data

*In general, error must not be treated as a potentially embarrassing inconvenience, because error provides a critical component in judging fitness for use.*

Chrisman, 1991

*Although most data gathering disciplines treat error as an embarrassing issue to be expunged, the error inherent in (spatial) data deserves closer attention and public understanding.*

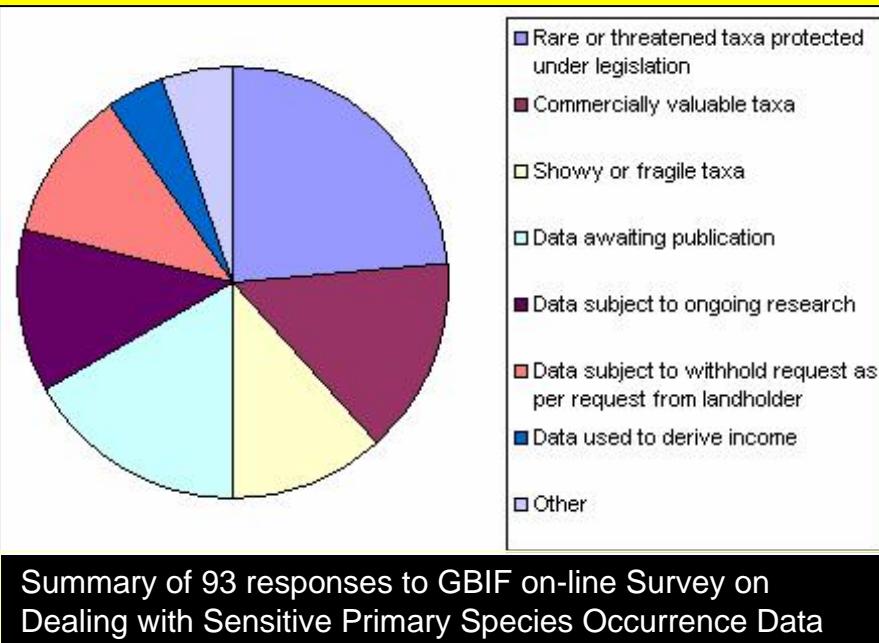
Chrisman, 1991

# Challenges

- Dealing with sensitive data



# What constitutes sensitive taxa?



## The issues:

- Institutions unknowingly making details of sensitive data available
- No universal global/regional list of sensitive taxa
- Duplicates of sensitive taxa from other countries

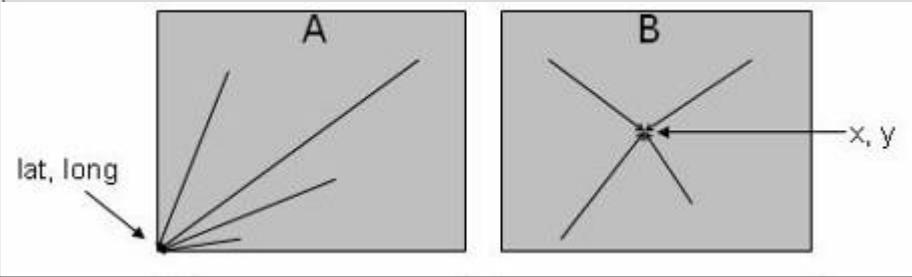
## The Solutions?

- Create an agreed list of sensitive taxa (by region)

# Problems with generalizing

Institutions are generalizing in many different ways and not documenting what they are doing.

	<b>Yes</b>	<b>Generalize (explain below)</b>	<b>Randomize (explain below)</b>	<b>Respondent Total</b>
a. Remove altogether	<b>94% (16)</b>	18% (3)	0% (0)	<b>17</b>
b. Round to 1 minute	<b>86% (6)</b>	57% (4)	29% (2)	<b>7</b>
c. Round to 10 minutes	<b>83% (5)</b>	67% (4)	0% (0)	<b>6</b>
d. Round to 30 minutes	<b>60% (3)</b>	<b>60% (3)</b>	0% (0)	<b>5</b>
e. Round to 1 degree	25% (1)	<b>75% (3)</b>	0% (0)	<b>4</b>
f. Move to nearest named place	<b>67% (2)</b>	<b>67% (2)</b>	0% (0)	<b>3</b>
g. Report by geographic region	<b>65% (11)</b>	53% (9)	6% (1)	<b>17</b>
h. Report by bioregion	<b>86% (6)</b>	29% (2)	0% (0)	<b>7</b>
i. Report by standard grid	50% (7)	<b>71% (10)</b>	14% (2)	<b>14</b>
j. Report by map sheet (explain scale, etc. below)	<b>88% (7)</b>	62% (5)	12% (1)	<b>8</b>
k. Combination of >1 above (note which, below)	<b>100% (4)</b>	50% (2)	25% (1)	<b>4</b>
l. Some other method (explain below)	<b>89% (8)</b>	67% (6)	44% (4)	<b>9</b>
<b>Total Respondents</b>				<b>46</b>
(skipped this question)				<b>56</b>



- Two generalization methods.
- A. a geographic grid where all records are referenced to the bottom right-hand corner.
  - B. a metric grid where all records are referenced to the centroid.

# Sociological Issues

- One case doesn't fit all
- Political issues
  - Endangered species (eg. Woller)
  - National legislation
  - Piracy
  - Trade and Quarantine
- Privacy
  - Names of observers, determiners
- Legal requirements
  - Permits
  - Observations in protected areas
  - Collections vis à vis permits
- Duplicate Collections
  - Filtered Push



*Solutions Still to be worked out*

A wide-angle photograph of a sunset or sunrise. The sky is filled with dramatic, wispy clouds colored in shades of orange, red, yellow, and purple. The horizon line is visible at the bottom, showing a dark silhouette of trees and land. The overall atmosphere is peaceful and scenic.

Questions  
please?

Thank You!