



Garbage In Garbage Out

開放資料的
前置準備與清理

楊富鈞

2024生物多樣性資料發布與應用工作坊

關於我

自然史博物館典藏管理
生物多樣性資料開放 (與應用)

楊富鈞 Fu-Chun Yang

典藏管理組 助理研究員

yuukumo0312@gmail.com

文化部 國立臺灣博物館



動機/誘因

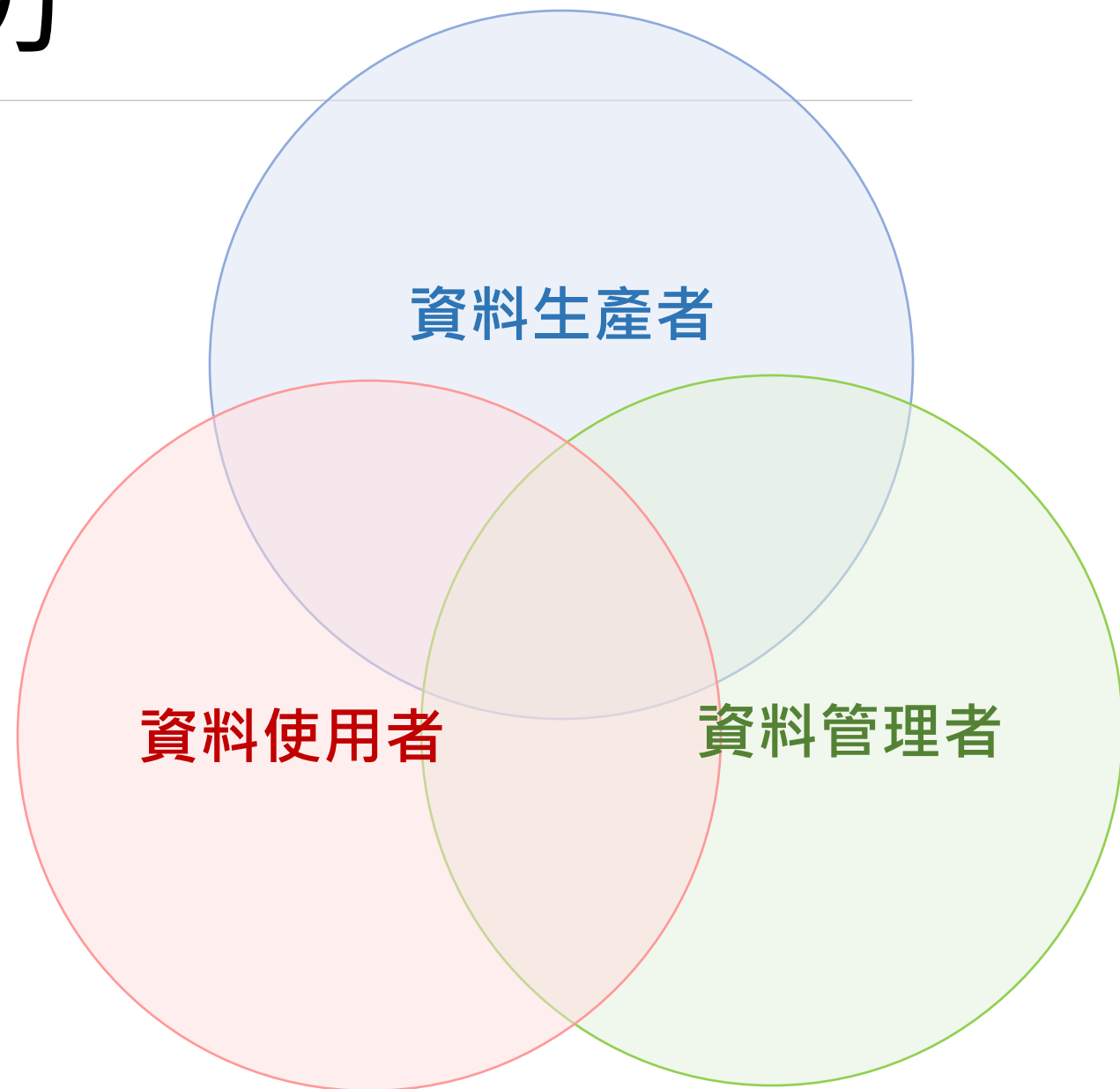
為何要花時間 整理資料？

讓資料變得更好用

動機/誘因 + 角色/人力

讓誰好用？為誰整理？

- 1 週後的自己
- 1 月後的自己
- **1 年後的自己**
- 長官或老闆
- 同行好友們 (專業同溫層)
- 非同行好友們 (人生同溫層)
- 同行冤家們
- **社會大眾 (終極目標)**



我想使用
或
由我負責管理
的資料

最**完美**的資料狀態

一千零一夜後...



@ Luis Prado from
Noun Project

已足夠**合理**的資料狀態

開放資料

已足夠合理成為開放資料的資料狀態

- 內容維護具可行性與持續性 (重要)
- 資料集內部邏輯一致具整體性 (強烈建議)
- 合法且合乎倫理 (重要)
- ...

一筆快樂的原始資料

！！請先備份！！

雲端1份、本機1份、外接硬碟1份 (這三個大家都有吧)

確認對上的共通欄位

Darwin Core Quick Reference Guide

(定下所有的資料表格式)

資料清理

Excel、OpenRefine、R

確認：

- (1) Core & extension分配
- (2) 各表單格式、單位、表示法
- (3) 有沒有還可以填的欄位

建議：留下文件給內部共通用
可藉此寫下MetaData初稿
(防止自己金魚腦，做事好利利)

目標：符合定下的DwC格式

學名對不對

Nomenmatch

ID生成&配對

自設ID 或 GUID

引入分類API

GBIF Backbone

座標轉換

TaiBIF上有小工具、R 或
Canadensys coordinate conversion

敏感資料處理

可跟著同步處理MetaData、
與coordinateUncertaintyInMeters

抓錯字、多餘空格、ID重複

GBIF Data Validator

記得要做版本控制囉

檢核完成 & 蓄勢待發

別忘了還有MetaData囉

開始處理資料前的 重要概念

開始處理資料前的重要概念

- 資料備份與版本維護永遠重要
- 原始資料須妥善保存並儘可能使其可被系統性檢索與存取
- 以原始資料為整理資料時的根本對象與依據
- 創建適用的表單

開始處理資料前的重要概念

- 資料**備份**與**版本維護**永遠重要
- 原始資料須妥善保存並儘可能使其可被系統性檢索與存取
- 以原始資料為整理資料時之**對象與依據**

你本人
(磨拳擦掌準備整理資料)

資料清理

生成複本

執行

提供

上游來源



資料處理前重要的基本概念

- 資料備份與版本維護永遠重要
- **原始資料**須妥善保存並儘可能使其可被系統性檢索與存取
- 以**原始資料**為整理資料時的根本對象與依據
- 實際動手清理 —— 收東西前應先確保你預計用於收納的櫃架
本身井然有序

養成思考哪一套內容版本才是
「原始資料」的習慣！

壹博館...

檔案 常用 插入 版面 公式 資料 校閱 檢視 其他資訊 揚音

AI64429 : X ✓ f 41°30'N, 128°10'E

	B	AI	AJ	AK	AL
1	catalogNumber	verbatim	coordinates	country	country
64429	TAIMB001764	41°30'N, 128°10'E			
64430	TAIMB001765	41°40'N, 127°50'E			
64431	TAIMB001766	41°40'N, 128°00'E			
64432	TAIMB001767	41°30'N, 128°10'E			
64433	TAIMB001768	41°40'N, 128°00'E			
64434	TAIMB001769	41°50'N, 127°50'E			

NTM_BioCo ...

就緒 100%

2919. **TMBAE001764**
 BOTANICAL MUSEUM, UNIVERSITY OF HELSINKI
 Helsinki, Finland Mniaceae
Mnium ambiguum H. Mill. 提灯藓
 det. T. Koponen -82

CHINA. Jilin (Kirin) Prov. An-tu Co.:
 Mt. Chang Bai, Valley of R. Er-do
 Bai-hu. Upper oroboreal Betula ermanii
 forest at small lake, alt. 1700 m,
41°30'N, 128°10'E, on boulder

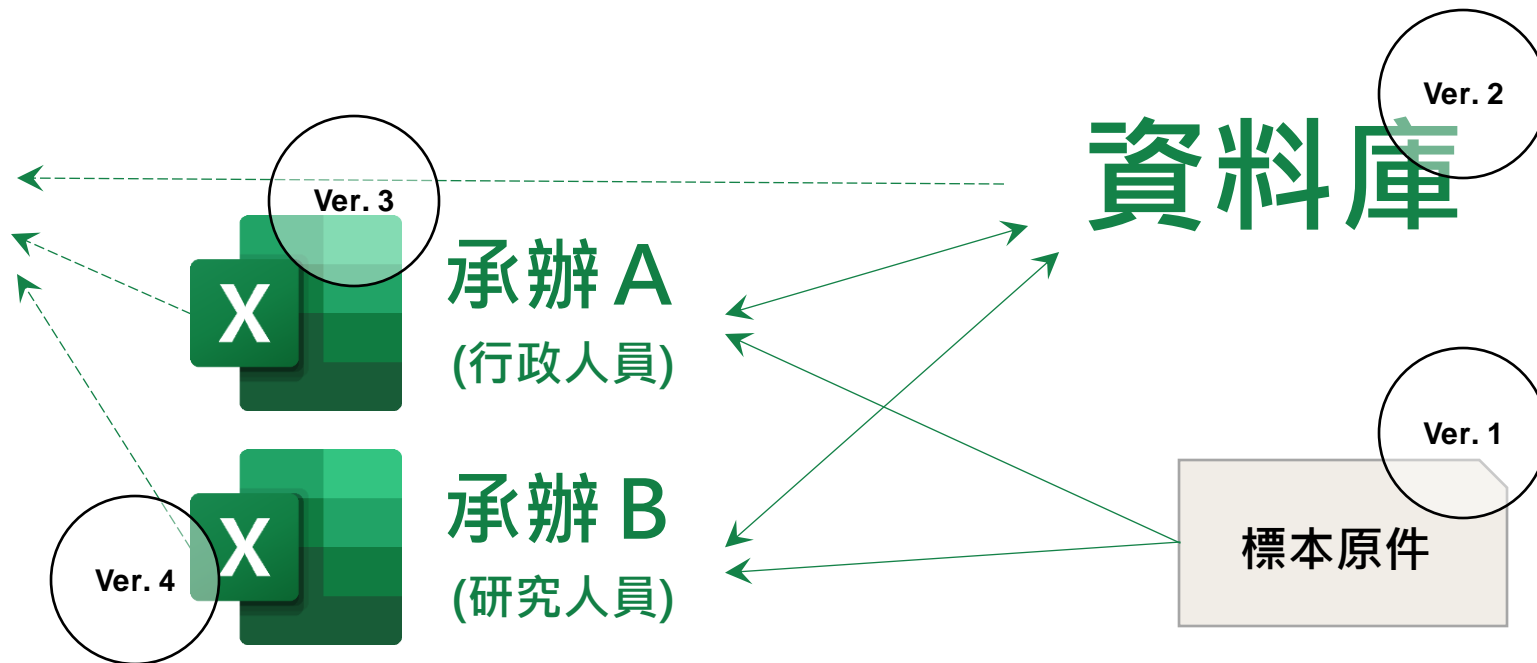
22 September 1981 Timo Koponen 36870
 (collection site no. 6)

標本原件

數位化很重要！

你本人

(磨拳擦掌準備整理資料)



個案說明與分析

CASE: P0853



典藏資源檢索系統
國立臺灣博物館

貝類/P0853江西巴蝸牛

江西巴蝸牛 詮釋資料

藏品記錄

編目號：P0853

中文名稱：江西巴蝸牛

數量：1

單位：件

來源說明：碧山巖寺99.11.23捐贈

採集地理資訊

採集地：廣州市廣西壯族自治區

物種分類資訊

綱名：腹足綱 Gastropoda

目名：柄眼目 Stylommatophora

科名：扁蝸牛科 Bradybaenidae

學名：*Bradybaena kiangsinensis*

命名者：(Martens)

命名年代：1875

藏品描述

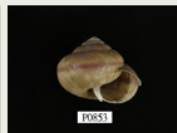
屬較大型的蝸牛。螺殼約有6個螺層、右旋，殼體呈圓球形，殼質厚且堅固，螺塔呈圓錐形，殼頂尖，各螺層略膨脹，體層膨大，縫合線深。殼面具有許多細密的生長紋。殼表呈黃褐色具光澤，體層的周緣有1條紅褐色的寬色帶。殼口呈橢圓形，外唇緣薄而銳利，臍孔深、部分被反捲的的軸唇所遮蓋。本種為中國大陸的特有物種，目前記錄於黑龍江、北京、河北、河南、湖南、湖北、四川、江西、廣西。主要棲息在樹林中，喜好陰暗潮濕的環境，也可見於農田田埂、壟間的雜草叢、灌木叢或亂石堆里、甚至住家公園也可發現到他們的蹤跡。本科物種的卵屬於小型且柔軟，產於鬆軟的泥土中或是地表上。

殼頂附近的殼面輕微磨損、殼皮脫落，殼體有幾處裂損，殼表和殼口內有髒污，殼唇緣有缺損。



P0853

廣州市廣西壯族自治區???



個案說明與分析

CASE: P0854



典藏資源檢索系統
國立臺灣博物館

貝類/P0854江西巴蝸牛



江西巴蝸牛 詮釋資料



藏品記錄

編目號：P0854

中文名稱：江西巴蝸牛

數量：1

單位：件

來源說明：碧山巖寺99.11.23捐贈



採集地理資訊

採集地：廣州市廣西壯族自治區



物種分類資訊

綱名：腹足綱 Gastropoda

目名：柄眼目 Stylommatophora

科名：扁蝸牛科 Bradybaenidae

學名：*Bradybaena kiangsinensis*

命名者：(Martens)

命名年代：1875



藏品描述

屬較大型的蝸牛。螺殼約有6個螺層、右旋，殼體呈圓球形，殼質厚且堅固，螺塔呈圓錐形，殼頂尖，各螺層略膨脹，體層膨大，縫合線深。殼面具有許多細密的生長紋。殼表呈黃褐色具光澤，體層的周緣有1條紅褐色的寬色帶。殼口呈橢圓形，外唇緣薄而銳利，臍孔深、部分被反捲的的軸唇所遮蓋。本種為中國大陸的特有物種，目前記錄於黑龍江、北京、河北、河南、湖南、湖北、四川、江西、廣西。主要棲息在樹林中，喜好陰暗潮濕的環境，也可見於農田田埂、壟間的雜草叢、灌木叢或亂石堆里、甚至住家公園也可發現到他們的蹤跡。本科物種的卵屬於小型且柔軟，產於鬆軟的泥土中或是地表上。

殼頂附近的殼皮磨損脫落，殼體有幾處裂損，殼表和殼口內有髒污，殼唇緣有缺損。



廣州市廣西壯族自治區???



- 已儲存
- 最近
- 廣州市
- 廣西壯族自治區

中國廣東省廣州市



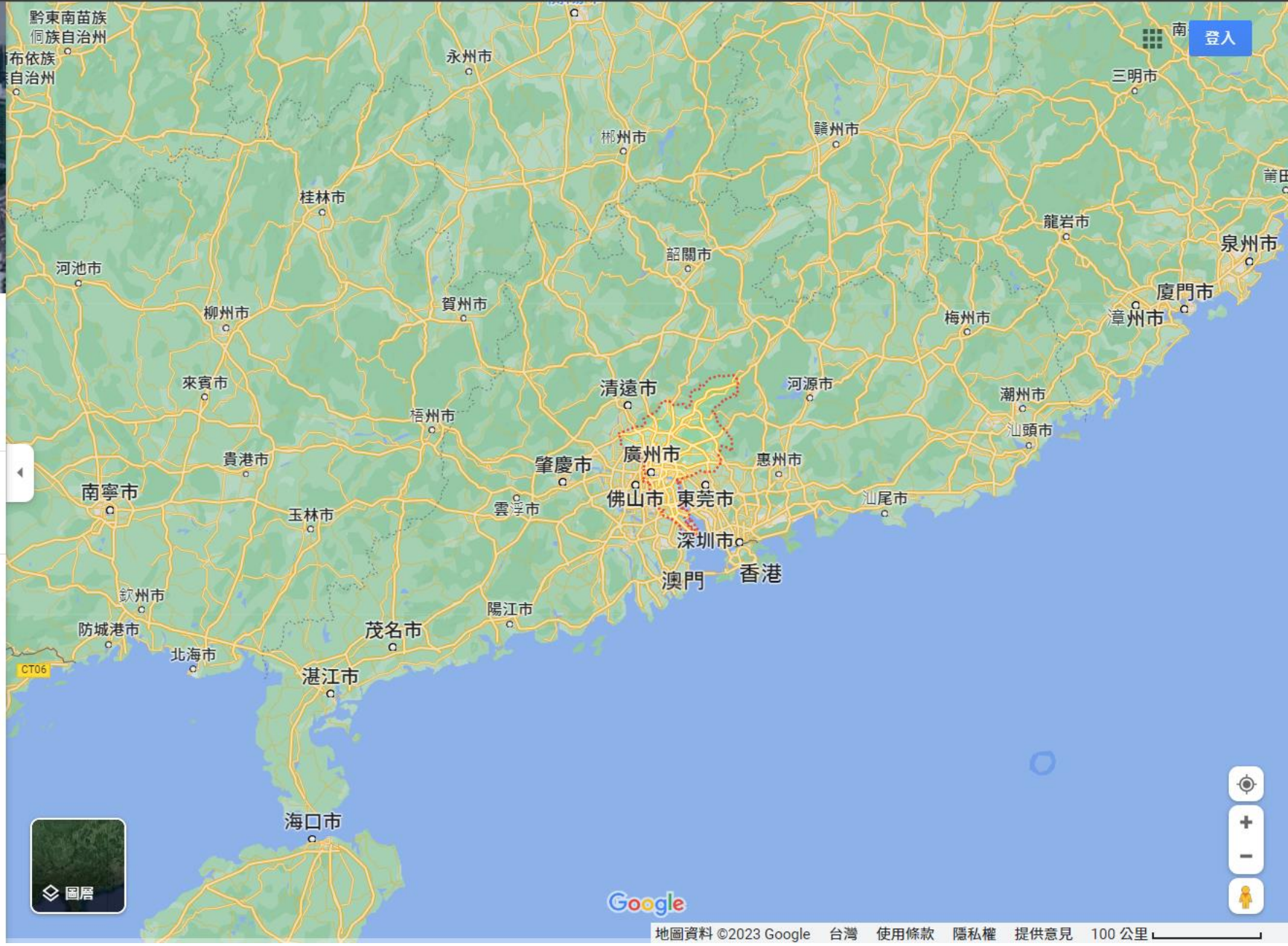
廣州市

广州市
中國
廣東省

多雲時陰 · 29°C
凌晨1:41

- 規劃路線
- 儲存
- 附近
- 傳送到手機
- 分享

廣州市，通稱廣州，簡稱廣、穗，別稱羊城，是中華人民共和國廣東省省會、副省級市、首批沿海開放城市。廣州市為中國大陸和廣東對外的商貿中心兼綜合交通樞紐，是中國大陸的一線城市之一，也是粵港澳大灣區的中心城市之一，中國人民解放軍南部戰區聯合指揮部亦駐紮該地。
[維基百科](#)



- ☰
- 📌 已儲存
- 🕒 最近
- 🏠 廣西壯族自治區
- 🏠 廣州市

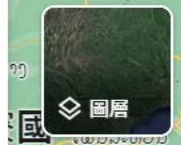
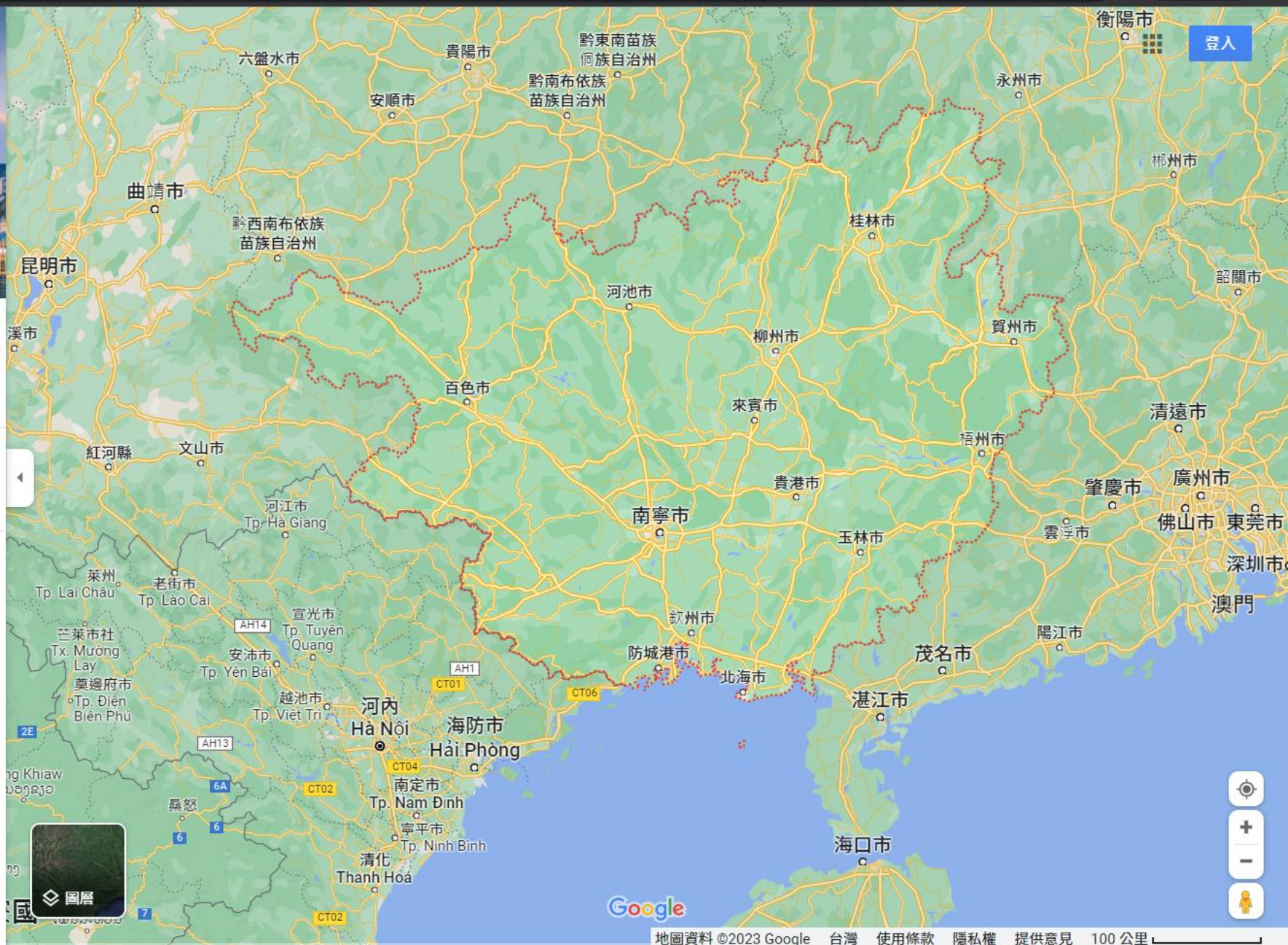
中國廣西壯族自治區



廣西壯族自治區
 广西壮族自治区
 中國

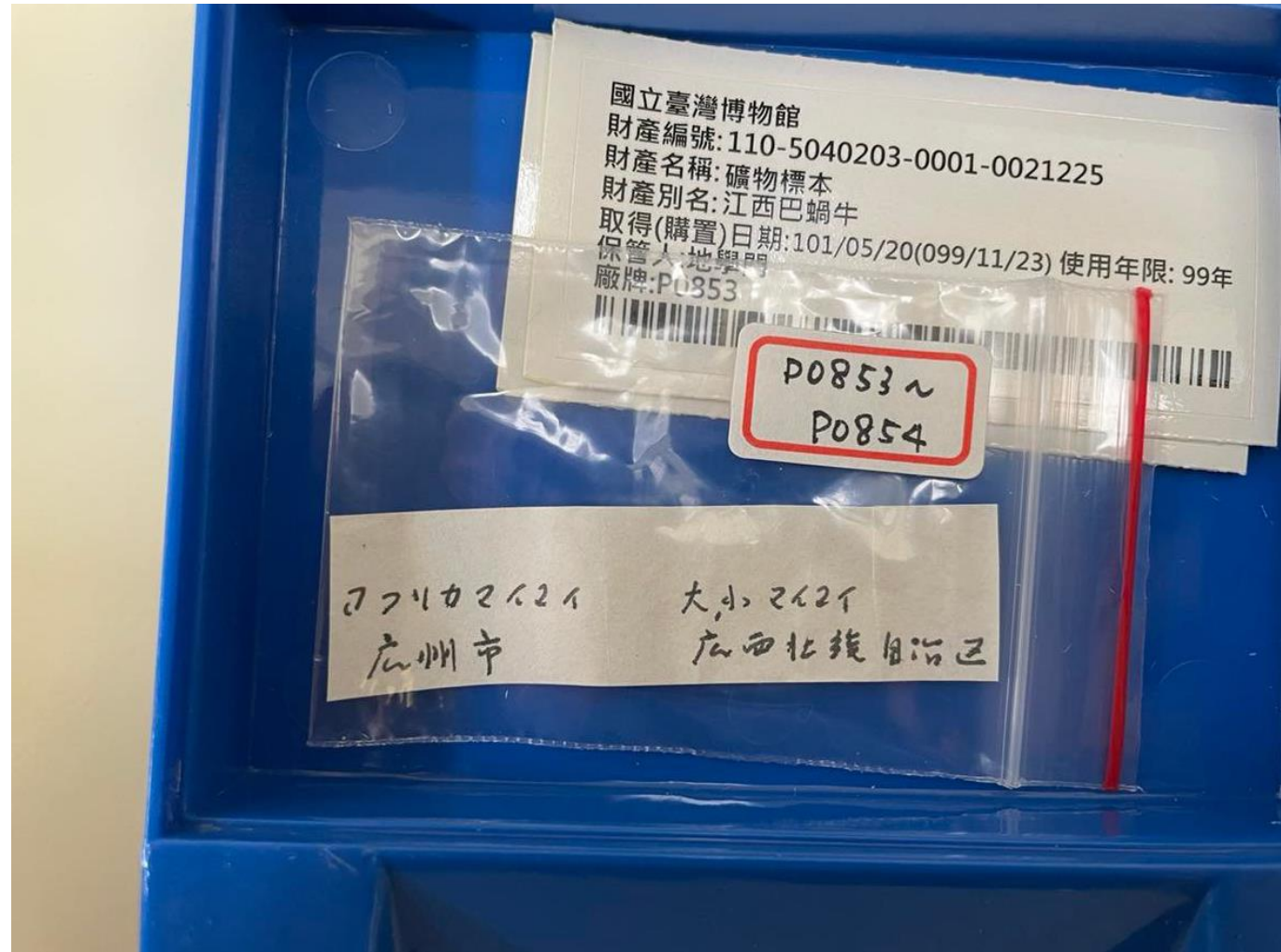
- 📍 規劃路線
- 📌 儲存
- 📍 附近
- 📱 傳送到手機
- 🔗 分享

相片



個案說明與分析

CASE: P0853 & P0854



資料處理前重要的基本概念

- 資料備份與版本維護永遠重要
- 原始資料須妥善保存並儘可能使其可被系統性檢索與存取
- 以原始資料為整理資料時的根本對象與依據
- 創建適用的表單

資料的櫃架是表格 (Table)

以由直行 (column)、橫列 (row) 組成，且遵循一套標準規則的表格來存儲及整理你的資料。

欄位 & 欄位名稱

欄位值

列

行

身分證字號			
DADA000271			
SEVEN340			
QURASO303			

欄位定義與欄位分工

欄位定義放哪裡？

欄位分工

- 所有欄位都有明確功能。
- 欄位之間的功能**不重複**。

欄位定義

- 所有欄位都有明確定義（用來載錄什麼資料、資料的值域、是否使用控制詞彙等），**且該定義最好是可被使用者透過網際網路系統性查詢的**。



欄位定義的**維護**頗麻煩，所以較合理的做法是

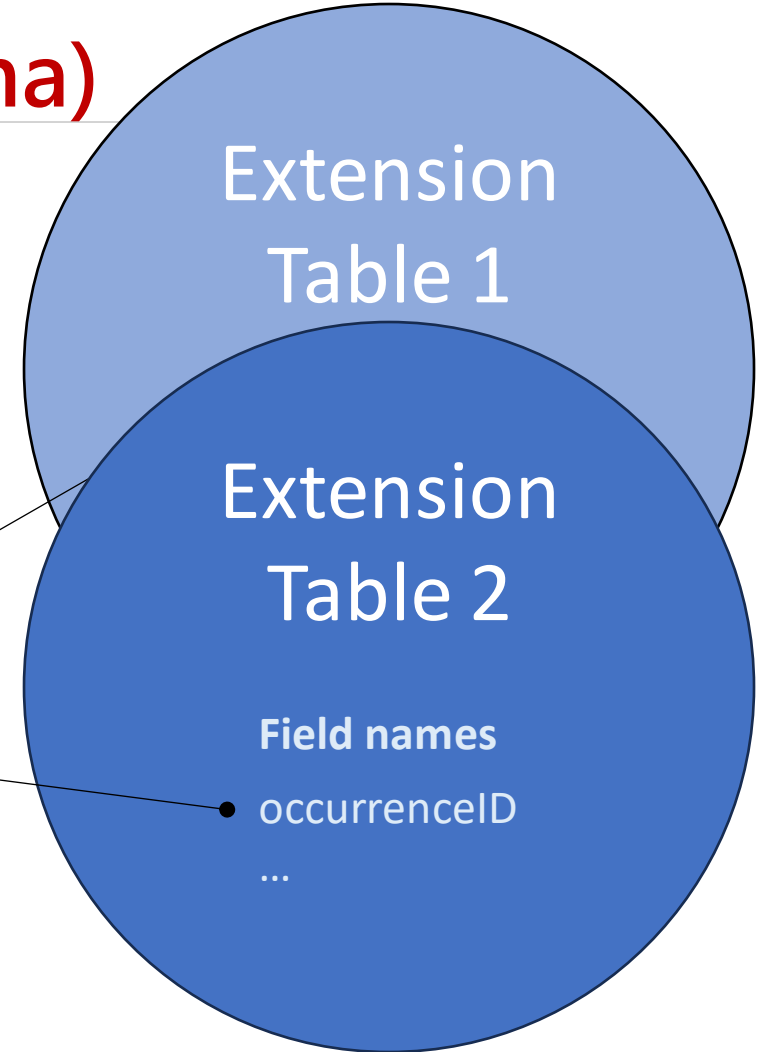
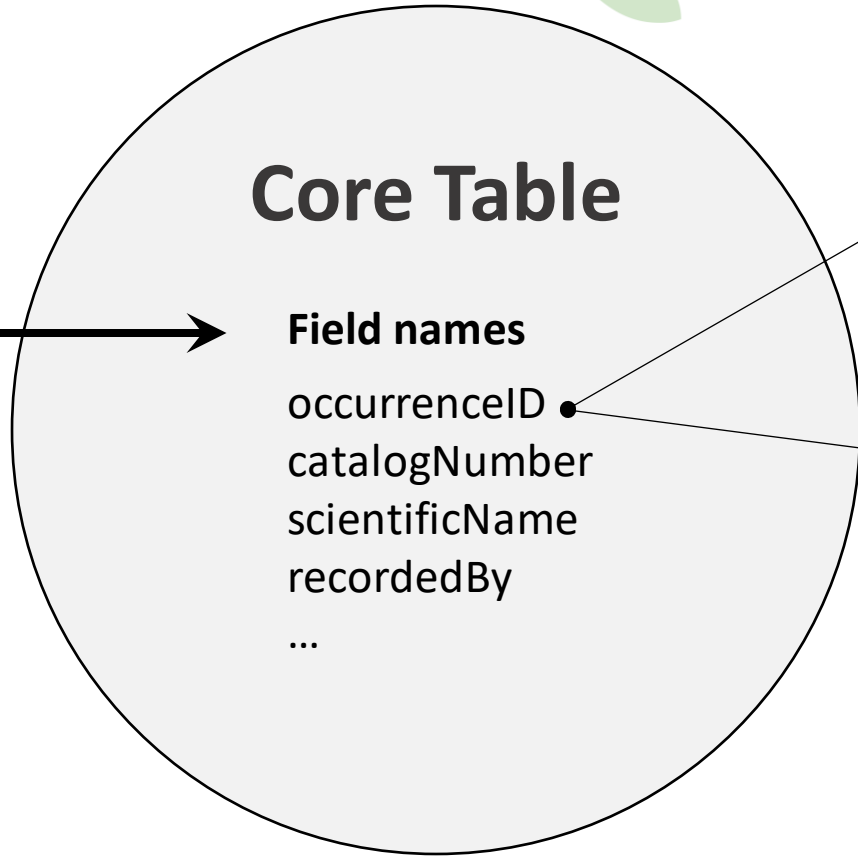
引用人家已發展成熟的方案

——至少以之為發展起點

GBIF 及 TBIA 倡議的資料結構 (Data Schema)



cited →



<https://rs.gbif.org/extensions.html>

概念/專業知識

以臺博館為例

occurrenceID	catalogNumber	parentCatalog	otherCatalog	otherCollection	basisOfRecord	occurrenceStatus	typeStatus	preparation	collection
4510	TAIMBOC_TAIMB005825				Preserved	present			
4511	TAIMBOC_TAIMB005826				Preserved	present			
4512	TAIMBOC_TAIMB005827				Preserved	present			
4513	TAIMBOC_TAIMB005828				Preserved	present			
4514	TAIMBOC_TAIMB005829				Preserved	present			乾燥標本 Distrit
4515	TAIMBOC_TAIMB005830				Preserved	present			乾燥標本 藏品
4516	TAIMBOC_TAIMB005831				Preserved	present			乾燥標本 藏品
4517	TAIMBOC_TAIMB005832				Preserved	present			乾燥標本 藏品
4518	TAIMBOC_TAIMB005833				Preserved	present			乾燥標本 贈葉

occurrenceID

Darwin Core Occurrence

主表處理館號、採集者、採集日期、採集地、館方選定的鑑定結果等資料

ONLINE
(in OneDrive)

selectIndicator	catalogNumber	identificationID	verification	scientificName
TRUE	NTMP0001-01	NTMP0001-01_meta_00101	Champeo	Champsodon guentheri Regan, 1908
TRUE	NTMP0001-02	NTMP0001-02_meta_00101	Champeo	Champsodon guentheri Regan, 1908
TRUE	NTMP0001-03	NTMP0001-03_meta_00101	Champeo	Champsodon guentheri Regan, 1908
TRUE	NTMP0001-04	NTMP0001-04_meta_00101	Champeo	Champsodon guentheri Regan, 1908
TRUE	NTMP0001-05	NTMP0001-05_meta_00101	Champeo	Champsodon guentheri Regan, 1908
TRUE	NTMP0001-06	NTMP0001-06_meta_00101	Champeo	Champsodon guentheri Regan, 1908
TRUE	NTMP0001-07	NTMP0001-07_meta_00101	Champeo	Champsodon guentheri Regan, 1908
TRUE	NTMP0001-08	NTMP0001-08_meta_00101	Champeo	Champsodon guentheri Regan, 1908
TRUE	NTMP0002	NTMP0002_meta_00101	Crossosoma	Crossosoma lacustre Steindachner, 1906

Darwin Core Identification History

處理各標本的鑑定脈絡

occurrenceID	measurementType	measurement
TMMA0001	snout-vent length [吻肛長]	1660
TMMA0002	head length [頭長]	47
TMMA0003	snout-vent length [吻肛長]	375
TMMA0003	tail length [尾長]	50
TMMA0004	snout-vent length [吻肛長]	630
TMMA0004	tail length [尾長]	50
TMMA0006	head length [頭長]	40.5
TMMA0007	snout-vent length [吻肛長]	360
TMMA0007	tail length [尾長]	140

Extended Measurement Or Facts

處理各標本的形質測量資料

catalogNumber	DNA	otherNucleotide
TMA0234	CAAAATTCCTATTTTGGGGCCGTGAGCAGGATGGTAGGAAGCTGCCCTTAGCCCTCTTATCC	
TMA0250	CAAGGAAMWAAACCTTCGGCCTGTMTCGGCATAATCGGAACARSCCTAARCCCTTCT	
TMA0266	CTGGACTAATCTTTGGYGCCTGGCCGGATAATCGGGACAGCCTTAAGCCTGCTAATTCG	
TMA0268	AGGTCCWAMTCTTTGGCGCCTGGCCGGATATCGGGACAGCCTTAAGCCTACTAATTCG	
TMA0273	TAGGSCWYTRCWWTATTTTGGAGCTGAGCAGGATAGTAGGAAGCTGCCCTTAGCCTC	
TMA0282	CCRKGTWCTACTSTTGGYGCCTGRGSGGATAGTCGGGACTGCCCTTAGCCTTCTMATY	
TMA0285	CTACAGCTCCTSAOCACTAGTCCTCTGCGCCTAGASCTCTAACCKSCCCAYCTTATCC	
TMA0331	CTGGCTATCTTTGGCGCWGGCCGGATATCGGGACAGCCTTAAGCCTGCTAATTCGAGC	
TMA0332	AMGGRAMTTAATCTTTRCGCAWGGCCRGATAATCGGACAGCCTTAAGCCTGCTAAT	

DNA derived data

處理各標本的定序資料

catalogNumber	filename	type	format	identifier	description

Simple Multimedia

處理各標本的影像資料

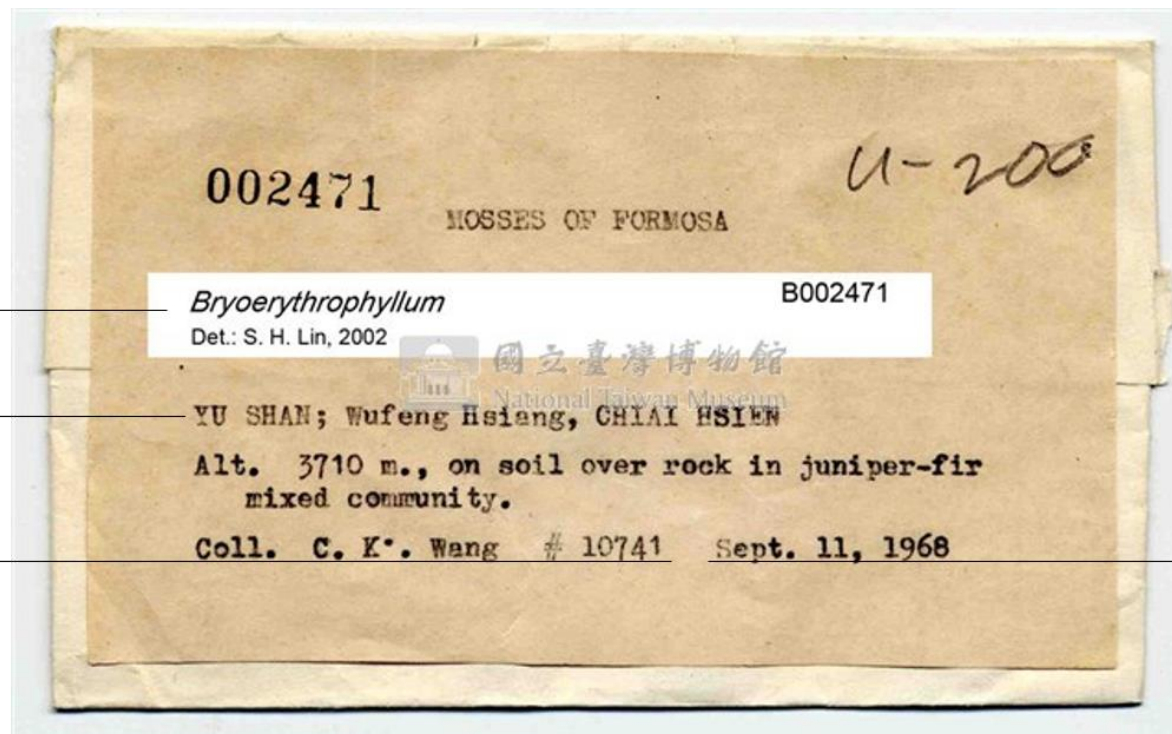
我對欄位的分類邏輯

基於內容類型

分類鑑定資訊

採集地資訊

採集者資訊



採集日期資訊

我對欄位的分類邏輯

基於長期維護的方式

- 辨識碼 (ID)
- 逐字稿欄位 ← 原封不動地將內容轉錄進去！
- 標準化內容欄位 ← 基於特定規則之內容；字數一般較少。
 - 根源於逐字稿欄位
 - 獨立於逐字稿欄位
- 註記欄位 ← 用於收納任意型式的註記；字數可多可少。
- 其他

表格型制穩定後

再來就是正式面對欄位值清理

生物多樣性資料有哪些常見的待清理樣態？

以「欄位值」錯誤樣態舉例：

- 違反資料結構（如欄位值未落於該欄位值域範圍、資料型態錯誤）
- 不合理的資料重複（如 unique field 出現重複值）
- 拼字錯誤
- 缺失值
- 異常值
- 冗餘的「空格」
- 欄位值缺乏一致性或整體性
- 相關欄位彼此矛盾或關聯性出現斷裂

...

條列一些個人對資料整理的"原則性建議或經驗"

- 善用各種線上服務與開源軟體
- 根基於常識/專業知識的**聯想力**或**想像力**
- 先建立**權威檔**
- 善用**自己**已建立的資料
- 把握與每個錯誤樣態的**不期而遇**
- 與 **Excel 函式 & 內建功能**當好朋友
- 透過與其他資料庫相互比對來校正內容
- 學習**正規表達式 (regular expression)**

資料取得 實務作法 工具 社群 關於

GBIF | Global Biodiversity Information Facility

自由、開放、可存取

出現紀錄 搜尋

什麼是 GBIF? 關於 GBIF 臺灣

發布者	使用者	GBIF 實驗室
IPT 整合式發布工具	Hosted portals	物種學名對應
資料驗證工具	資料處理	學名解析
科學典藏	衍生的資料集	序列識別碼
建議資料集	rgbif	相對觀測趨勢
New data model ★	pygbif	GBIF 資料部落格
	MAXENT	
	工具目錄	

Sawfly parasitic wasp (*Dahlbominus fuscipennis* (Zetterstedt, 1838)) collected in Nationaal Park Veluwezoom, Netherlands. Photo via Naturalis Biodiversity Center.



2,575,559,246

出現紀錄



89,923

資料集



2,129

發布機構



9,316

使用資料的同儕審查論文

從案例中學習

概念/專業知識 + 技巧/工具

- 善用各種線上服務與開源軟體
- 根基於常識/專業知識的**聯想力**或**想像力**

• 先建立權威檔

catalogNumber	eventDate
NTMP0479-13	1977-06-40

• 善用自己

• 把握與每個錯誤樣態的不期而遇

• 與 Excel 函式 & 內建功能當好朋友

• 透過與其他資料庫相互比對來校正內容

• 學習正規表達式 (regular expression)

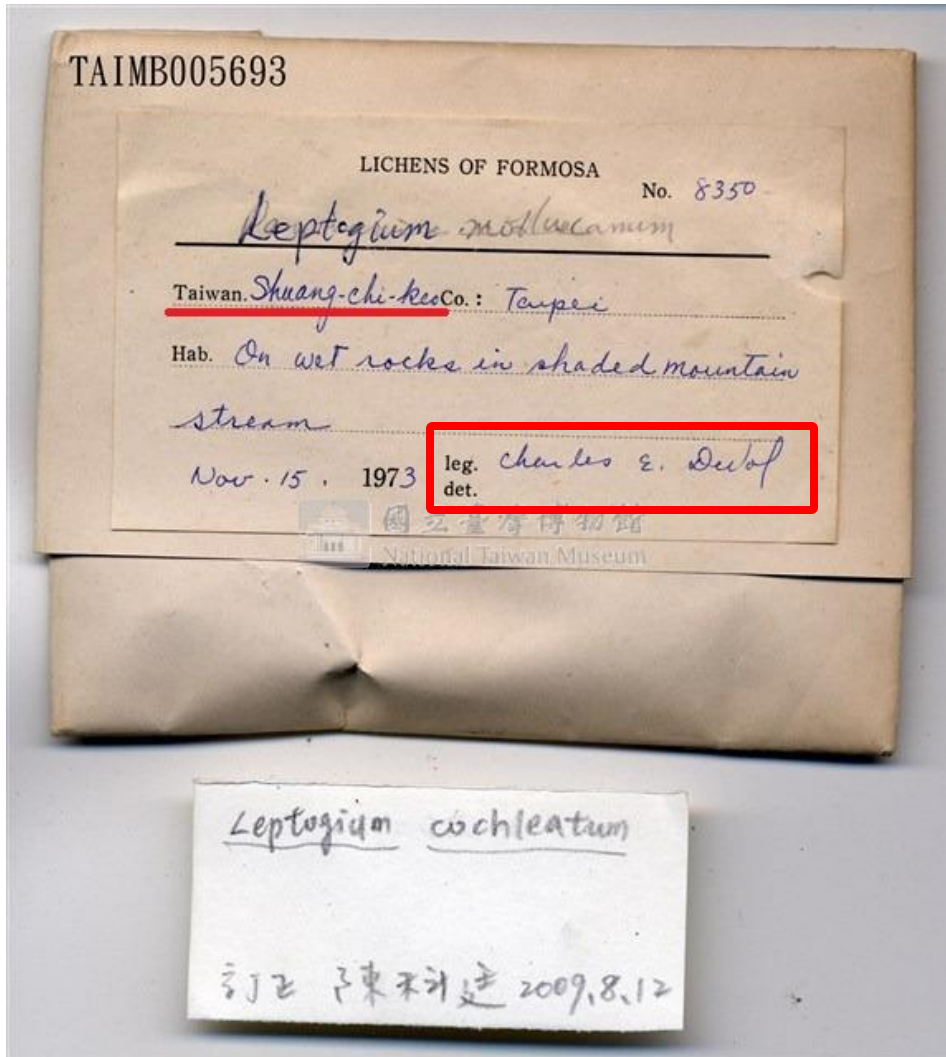
概念/專業知識 + 技巧/工具

- 善用各種線上服務與開源軟體
- 根基於常識/專業知識的**聯想力**或**想像力**
- 先建立**權威檔**
- 善用**自己**已建立的資料
- 把握與每個錯誤樣態的**不期而遇**
- 與 **Excel 函式** & **內建功能**當好朋友
- 透過與其他資料庫相互比對來校正內容
- 學習**正規表達式 (regular expression)**

概念/專業知識 + 技巧/工具

- 善用各種線上服務與開源軟體
- 根基於常識/專業知識的**聯想力**或**想像力**
- 先建立**權威檔**
- 善用**自己**已建立的資料
- 把握與每個錯誤樣態的**不期而遇**
- 與 **Excel 函式 & 內建功能**當好朋友
- **透過與其他資料庫相互比對來校正內容**
- 學習**正規表達式 (regular expression)**

CASE: TAIMB005693



verbatimLocality

Taiwan. Shuang-Chi-{Keo[?]} Co.: Taipei. Hab. On wet rocks in shaded mountain stream

Shuang-Chi-{Keo[?]} 在哪？

Shuang-Chi = 雙溪？哪個雙溪？

Charles E. DeVol

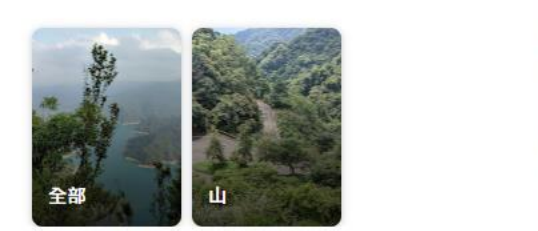
棣慕華 (1903-1989)





WHCM+CM 新店區 新北市

相片



新增相片

找不到所需的答案嗎？ 社群通常會在 20 分鐘內回答問題。

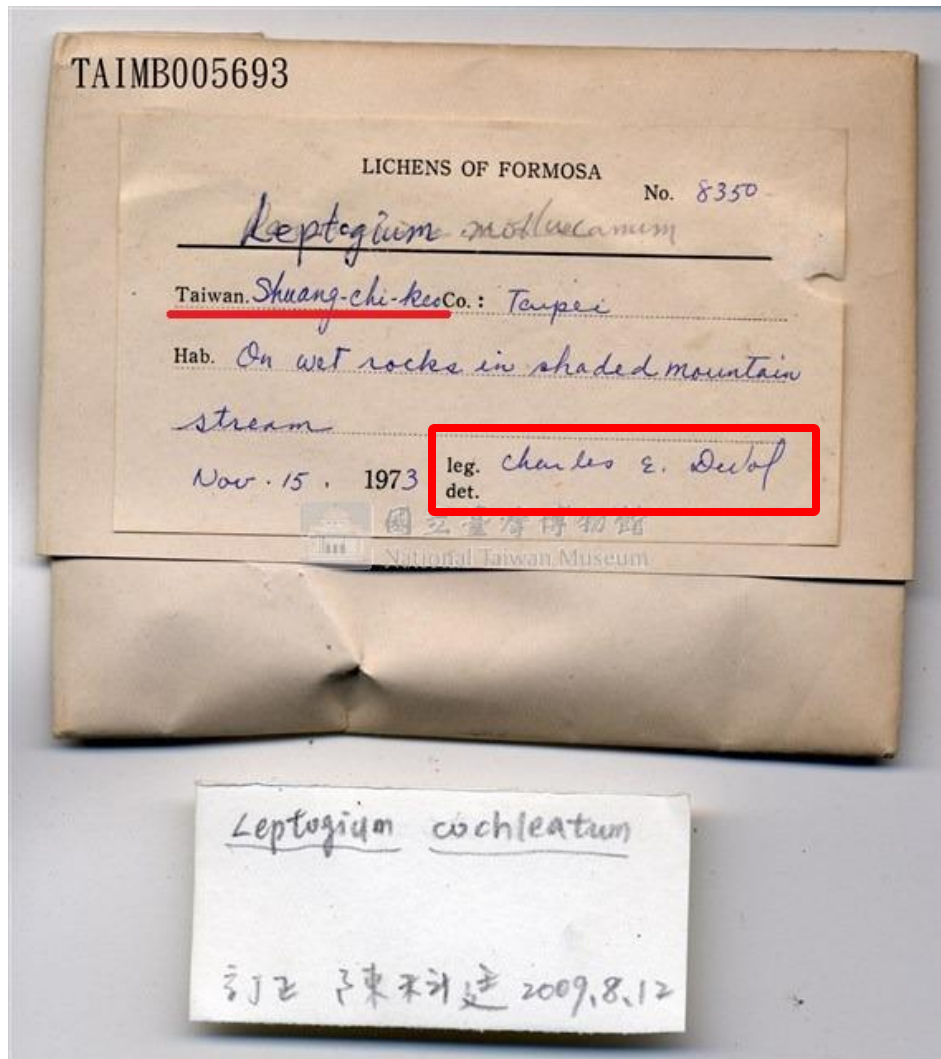


詢問社群成員



第 1 - 20 筆結果

CASE: TAIMB005693



verbatimLocality

Taiwan. Shuang-Chi-~~{Keo[?]}~~ Co.: Taipei. Hab. On wet rocks in shaded mountain stream

locality

New Taipei City | Xindian District | Shuangxikou, near Niaozuijian Mountain [新北市|新店區|雙溪口, 鄰近鳥嘴尖山]

結語

經你整理的資料不再是檔案而是品牌

身為資料管理者的自尊

就是時刻以中立之姿幫助使用者取得好資料

加油！敬祝工作顺利！