# Metadata Development and Documentation for a Research Data Repository

Huang-Sin Syu, Cheng-Jen Lee, Yao-Hsien Yeh, Tyng-Ruey Chuang

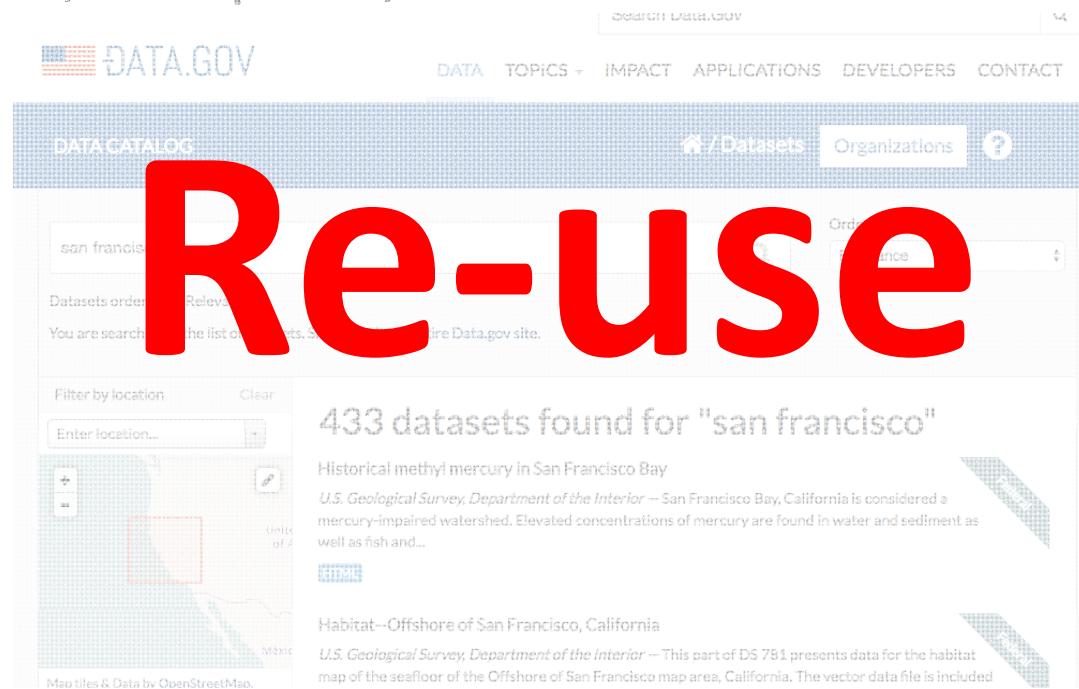Institute of Information Science, Academia Sinica, Taiwan

# Outline

- Research Data Repository

- Why we need metadata?

- Metadata schema and the quality of metadata

- Creating metadata for research datasets

- Documentation for Research Data Repository

- Documentation workflow for metadata schema

- Current results and future works

# Research Data Repository (RDR)

- Active management, adding value and maintaining access to research data

- Accurate, complete, retrievable

# Why we need metadata?

- Reduce dataset duplication
- Use datasets more effectively
- Make research data searchable

# Metadata schema and the quality of metadata

- **General standard**
  - Dublin Core Metadata Initiative(DCMI)

- **Geospatial**
  - ISO 19115:2003 Geographic information metadata

# Metadata schema and the quality of metadata

Institute of Information Science, Academia Sinica, Taiwan

# Metadata schema and the quality of metadata

**Buttom-up** ➡ ⬅ **Top-Down**

Expert

User

UGC

Expert

standards

standards

Research Data Repository

# Metadata schema and the quality of metadata

- **Continuous revisions of metadata schema**

- **Customized metadata schema**

# Creating metadata for research datasets

- Use CKAN to build Research Data Repositories
  - Taijiang Project
  - Ka-lam Project
  - More

- Metadata authoring via CKAN webpage

- Dataset/Metadata Bulk Upload (Customized process for Taijiang project)

- Dataset/Metadata Harvest API

# Creating metadata for research datasets

- Metadata authoring via CKAN webpage

# Creating metadata for research datasets

- Dataset/Metadata Bulk Upload (Customized process for Taijiang project)  https://taijiang.tw/en/help

**Institute of Information Science, Academia Sinica, Taiwan**
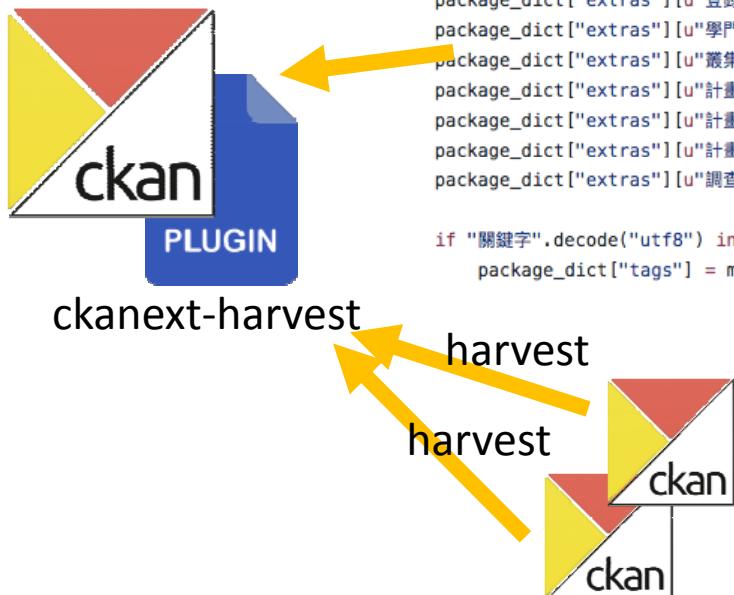
# Creating metadata for research datasets

- Dataset/Metadata Harvest API



```
           meta[key] = value
package_dict["title"] = meta[u"計畫名稱"]
package_dict["author"] = meta[u"計畫主持人"]
package_dict["notes"] = meta[u"摘要"]
#package_dict["metadata_modified"] = datetime.today().strftime(
package_dict["extras"][u"資料集網址"] = self.PREFIX_URL + harvest
package_dict["extras"][u"登錄號"] = meta[u"登錄號"]
package_dict["extras"][u"學門類型"] = meta[u"學門類型"]
package_dict["extras"][u"叢集名稱"] = meta[u"叢集名稱"]
package_dict["extras"][u"計畫執行單位"] = meta[u"計畫執行單位"]
package_dict["extras"][u"計畫委託單位"] = meta[u"計畫委託單位"]
package_dict["extras"][u"計畫執行期間"] = meta[u"計畫執行期間"]
package_dict["extras"][u"調查執行期間"] = meta[u"調查執行期間"]

if "關鍵字".decode("utf8") in meta.keys():
    package_dict["tags"] = meta[u"關鍵字"].split(u"、")
```

ckanext-harvest

harvest

harvest

◎ 分類查詢結果（共 16 筆）

| 學門類型 | 藝術學 |
| --- | --- |
| 登錄號 | E87001 |
| 計畫名稱 | 生活型態與飲食文化對廚具設計開發影響之研究 |
| 計畫主持人 | 胡祖武 |
| 叢集名稱 | |
| 計畫執行單位 | 大葉大學工業設計系 |
| 計畫執行期間 | 1997-08-01 ~ 1998-07-31 |
| 調查方式 | 郵寄問卷 |
| 樣本數 | 154 |
| 中英文關鍵字 | 生活型態、飲食文化、廚具設計<br>Dietary Culture、Kitchen Unit Design、Life Style |

| 學門類型 | 藝術學 |
| --- | --- |
| 登錄號 | E87006 |
| 計畫名稱 | 台灣地區藝術科系與非藝術科系學生對「錯視圖形」的 |

https://srda.sinica.edu.tw/search/field/2
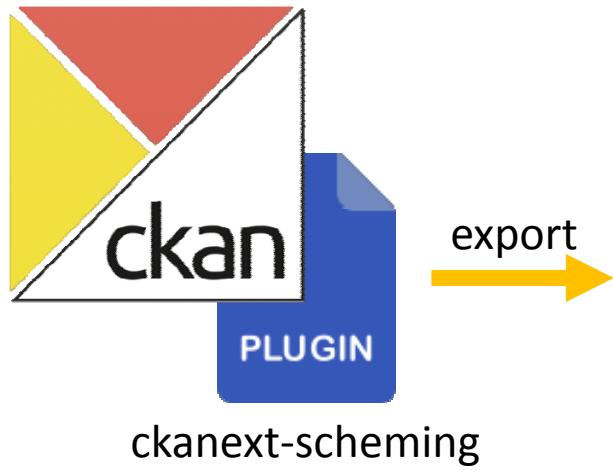
# Documentation for RDR

- Explain metadata schema
- Aid to interpretation

- Workshop
- User guide
- Documentation for different usage scenarios

# Documentation for RDR

- **One source document for all metadata schema**

- **Rapid revision to metadata schema**

# Documentation workflow for metadata schema

https://taijiang.tw/dataset/mongolia-spectral-data



ckanext-scheming

export

JSON

# Documentation workflow for metadata schema



export

ckanext-scheming

```
5    "dataset_fields": [
6      {
7        "field_name": "title",
8        "label": {
9          "en": "Title",
10   "zh_TW": "標題"
11       },
12       "preset": "title",
13       "form_placeholder": {
14         "en": "eg. A descriptive title",
15         "zh_TW": "例如：一個描述性的標題"
16       }
17     },
18     {
19       "field_name": "name",
20       "label": {
21         "en": "URL"
```

```
1    {
2      "scheming_version": 1,
3      "dataset_type": "dataset"
4      "about_url": "https://github.com/u10313335/ckanext-taijiang",
5      "dataset_fields": [...],
1829   "resource_fields": [...]
1896 }
1897
```

# Documentation workflow for metadata schema

# Documentation workflow for metadata schema



export

JSON

ckanext-scheming

design

MD

UML

publish

GitBook

# Documentation workflow for metadata schema

**Schema** **+** **Description**

# Documentation workflow for metadata schema



Type to search

Introduction

Metadata (UML diagram)

zh-TW

en

Metadata (Table view)

Data_type.md

Encoding_type.md

Hist_material_type.md

Language_type.md

License_code.md

Loc_keyword_type.md

Metadata.md

Ref_book.md

Scan_pics.md

**GitBook**

C Ref_book

○ 國際標準書[isbn][0..1]:CharacterString
○ 國際標準叢刊號[issn][0..1]:CharacterString
○ 期刊名稱[0..1]:CharacterString
○ 卷期[0..1]:Interger
○ 論文集名稱[0..1]:CharacterString
○ 出版地[0..1]:CharacterString
○ 出版社[0..1]:CharacterString
○ 出版年份[0..1]:Interger
○ 書目查詢[0..1]:CharacterString
○ 網址[0..1]:CharacterString
○ 使用史料[0..*]:Hist_material_type
○ 備註[0..1]:CharacterString

○ 最南緯度值[0..1]:Angle
○ 最北緯度值[0..1]:Angle
○ 主題分類[0..1]:Theme_type
○ 授權[1..1]:License_code
○ 參考來源[0..3]:CharacterSt
○ 主題分類[0..*]:CharacterSt
○ 主題分類[0..1]:CharacterSt

C SpatialMetad

○ 坐標系統[0..1]:Integer
○ 空間解析度[0..1]:Real
○ 比例尺[0..1]:Integer
○ 處理歷程/*Lineage[0..1]:Ch

○ 原件來
○ 原件尺
○ 掃描解

add some text here

# Documentation workflow for metadata schema

Type to search

Introduction

Metadata (UML diagram)

   zh-TW

   en

Metadata (Table view)

   Data_type.md

   Encoding_type.md

   Hist_material_type.md

   Language_type.md

   License_code.md

   Loc_keyword_type.md
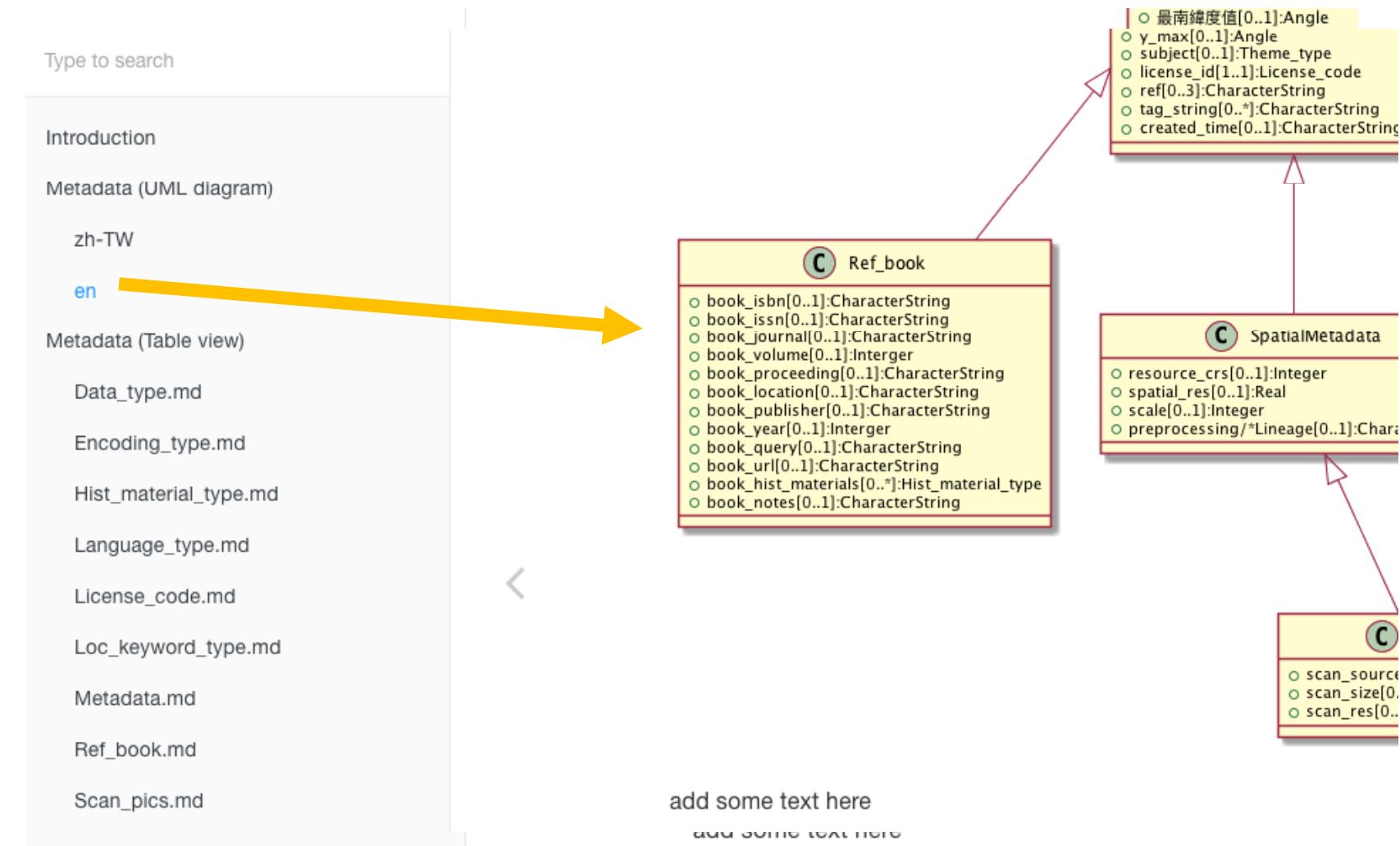
   Metadata.md

   Ref_book.md

   Scan_pics.md

```
                                                  ○ 最南緯度值[0..1]:Angle
                                                  ○ y_max[0..1]:Angle
                                                  ○ subject[0..1]:Theme_type
                                                  ○ license_id[1..1]:License_code
                                                  ○ ref[0..3]:CharacterString
                                                  ○ tag_string[0..*]:CharacterString
                                                  ○ created_time[0..1]:CharacterString

                 C  Ref_book

○ book_isbn[0..1]:CharacterString
○ book_issn[0..1]:CharacterString
○ book_journal[0..1]:CharacterString          C  SpatialMetadata
○ book_volume[0..1]:Interger
○ book_proceeding[0..1]:CharacterString   ○ resource_crs[0..1]:Integer
○ book_location[0..1]:CharacterString     ○ spatial_res[0..1]:Real
○ book_publisher[0..1]:CharacterString    ○ scale[0..1]:Integer
○ book_year[0..1]:Interger                ○ preprocessing/*Lineage[0..1]:Chara
○ book_query[0..1]:CharacterString
○ book_url[0..1]:CharacterString
○ book_hist_materials[0..*]:Hist_material_type                 C
○ book_notes[0..1]:CharacterString
                                             ○ scan_source
                                             ○ scan_size[0.
                                             ○ scan_res[0..
```

add some text here

add some text here

**Institute of Information Science, Academia Sinica, Taiwan**

21

# Documentation workflow for metadata schema

**Institute of Information Science, Academia Sinica, Taiwan**

# Current results and future works

- **Rapid documentation for** CKAN-based repository and publish to GitBook.

- Metadata schema designers **only need to maintain UML files.**

- Publishing to GitBook **makes schema and documents searchable.**

# Current results and future works

**Metadata as Linked Data for Research Data Repositories** (Mr. Cheng-Jen LEE)