

# 結構資料的再次使用：語意、連結與實作

黃韋菁

中央研究院資訊科學研究所專案經理

E-mail: andreaHG@iis.sinica.edu.tw

李承鑫

中央研究院資訊科學研究所研究助理

E-mail: cjlee@iis.sinica.edu.tw

莊庭瑞

中央研究院資訊科學研究所副研究員

E-mail: trc@iis.sinica.edu.tw

關鍵詞[1]：CKAN；資料溯源；資料品質；知識庫；開放資料連結（LOD）；知識本體；語意再現

---

## 【摘要】

持續創造資料的語意與連結，藉由全球資訊網散布同時可由常人和機器處理並理解的結構性資料，進而增進資料集的「再次使用價值」（reuse value）是目前廣受重視的課題，也是本研究由理論探討邁向系統實作的動力與目的。本文簡述與「開放資料連結」（Linked Open Data, LOD）相關國際計畫與技術發展，介紹以「開放資料連結」方式建置的五項跨領域知識庫和七項專業知識庫，並解析資料品質、後設資料（Metadata）及資料溯源（Provenance）的關聯脈絡。

本研究同時進行實作網站 data.odw.tw，收納典藏品目錄資料，並設計知識本體（voc4odw）轉換半結構式資料為富語意結構的連結式資料。一方面擴充 CKAN（The Comprehensive Knowledge Archive Network）資料集管理系統，作為連結式資料的儲存與展示平台，進而強調從原始目錄資料到語意連結資料的分段轉換步驟，最後將各步驟轉換程式以及 CKAN 軟體程式碼以「開放原始碼」（Open Source）方式釋出。另一方面，由於研究資料來源採「創用 CC」（Creative Commons）公眾授權，因此研究成果亦以相同方式釋出，在開放基礎上促使資料與程式碼的保存與發展，可被自由再次使用與擴散。

# 前言

「資料連結」(Linked Data) 關聯了資料 (data)、常人 (human) 以及機器 (machine) 三者。在知識呈現與語意處理的共有課題。例如，若常人提出問題：臺灣國寶級植物-臺灣一葉蘭的地理分布為何？可由典藏臺灣聯合目錄中[2]紀錄該植物標本後設資料，得知採集地點與相關描述資訊。該項資料進一步經「資料連結」方法進行知識呈現及語意處理後[3]，不僅可連接外部地理資源與相關藏品知識以提供臺灣一葉蘭之地理分布資訊 (圖 1) [4]，同時也可經由外部地名資料庫顯示採集地如宜蘭大同等資訊。透過機器可處理的典藏目錄標本的地理分佈資料與採集脈絡，並可得知該植物採集時間前後跨越 1983-2010 近三個世代，標本資料建置在不同時期由不同機構人員處理與再次使用等資訊。

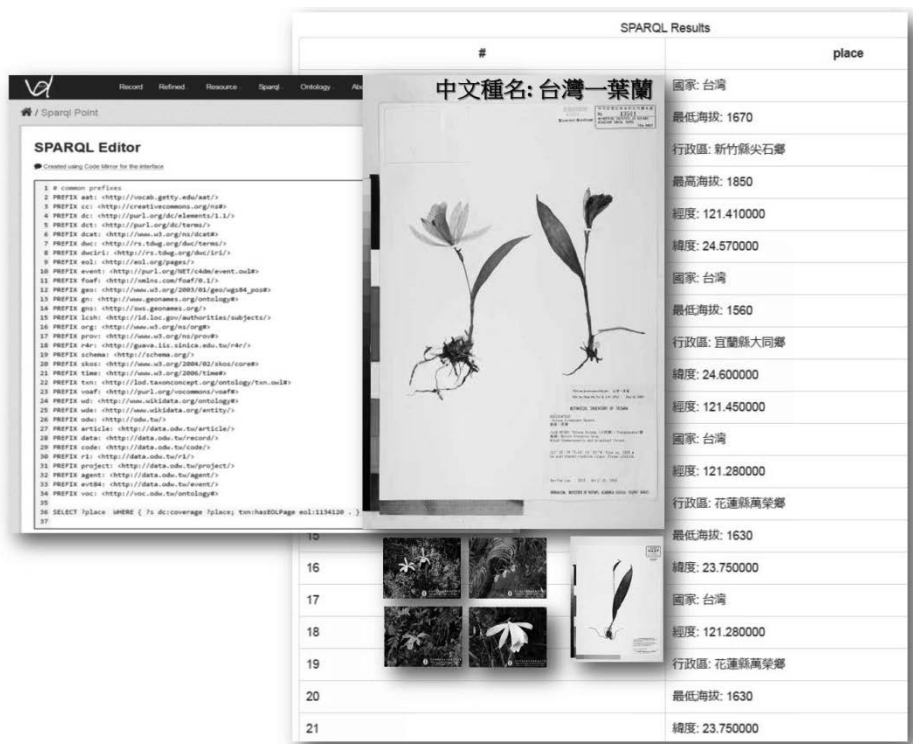


圖 1 以 SPARQL 查詢[3]臺灣一葉蘭標本藏品之地理分布示意圖

本研究使用 CKAN (The Comprehensive Knowledge Archive Network) 這項「資料藏庫」(Data Repository) 系統軟體工具[5]，在上面架構了「開放資料連結」(Linked Open Data, LOD) 新功能，整合提供可被常人以及機器可使用的連結資料服務 (Lee, Huang, & Chuang, 2016)。如今，典藏目錄裡臺灣一葉蘭的資訊，在此資料連結架構下，提供了同時適合機器與常人皆能理解的資料連結介面，進而能自資料連結的語意關係進行相對應的知識查詢。換言之，「開

放資料連結」的目的是支持語意網的運作，而本文所介紹的方法與實作將有助於常人對資料進行理解（interpretation of data）、資料被機器操作（machine actionable）、資料間建立語意連結（semantics linked）、且資料能被語意檢索（semantic query）等工作，最終走向人機語意連結與互通的語意網願景。

在此語意網與資料連結的願景中，圖書館、檔案館、博物館（Libraries, Archives and Museums, LAM）領域的後設資料語意標示的傳統優勢；例如，後設資料使用「都柏林核心集」（Dublin Core, DC）15 項欄位整合異質資料來源、積極建構特定學科主題控制語彙（Control Vocabularies）與知識索引典（Thesaurus）、重視服務使用者應用的實務操作經驗等特點。依此願景促使本研究使用典藏臺灣聯合目錄後設資料集為案例研究，一方面得利於上述三項優點；另一方面，該目錄的後設資料原始檔案為半結構化 XML 形式，亦提供機器自動化大量處理的便捷管道。

儘管如此，開放式資料連結方法仍須面臨許多挑戰。誠如 Hallo、Luján-Mora、Maté 與 Trujillo（2016）研究指出，圖書檔案博物館界在「開放資料連結」發展所面臨的困難包括六項：（1）技術工具的支持、（2）資料品質控制機制、（3）資料模型與語彙的實作、（4）人性化的瀏覽與查詢界面、（5）定義資料開放授權的困難，以及（6）缺少新技術知識的技術人員等。因此本文首先回顧近期「開放資料連結」相關計畫與技術發展、探究全球開放式連結資料知識庫的差異與品質、分析資料品質與資料溯源的關係，主要目的是一方面整理「開放資料連結」國際發展趨勢，便利進一步觀察與思考，另一方面也釐清目前實作案例 data.odw.tw 中許多設計原則的選擇因素與脈絡。

透過實作五大步驟：（1）探究共享脈絡下的資料再次使用的關聯性、（2）設計不同情境下的模式系統架構、（3）資料剖析、清理與比對、（4）以使用者為核心的資料庫知識平台技術架構、（5）透過知識本體以協助理解資料語意的再現與再次使用[6]。期望藉由實作案例分享，對目前正投入或預備投入「開放資料連結」研究者有所助益。而對於尚未接觸「開放資料連結」議題者，也能因此思考將已有的資料集附予新價值，讓單一或片段知識透過語意連結，成為全球知識網的連結點。

## 「開放資料連結」以及「開放資料連結知識庫」的發展

始創資料連結的 Tim Berners-Lee 於 2006 年指出：「驚人數量的資料呈現未連結的狀態」[7]。十年後的今天或可去掉「未」字，改為：「驚人數量的資料呈現連結的狀態」[8]。

觀察圖書館界的進展，Marden、Li-Madeo、Whysel 與 Edelstein（2013）分析當時 15 個文化資產的「開放資料連結」計畫後指出，文化資產資料無法廣泛被使用的最大障礙是：大多數機構尚未以開放資料連結方式發佈與使用資料。三年後此令人擔憂的狀況獲得改善，主要誘因包括為增加資料曝光度吸引更多使用者、示範資料集能完成資料連結程度、普遍聽聞

「資料連結」趨勢而嘗試、測試資料連結是否能優化搜索引擎效能等 (Mitchell, 2016; Godby, 2016)。例如「線上電腦圖書館中心」(Online Computer Library Center, OCLC) 針對 20 個國家 90 個圖書博物館機構的報告指出[9]，相對於 2014 年的調查資料連結計畫已快速成長兩倍。另一方面，歐盟計畫 Europeana[10]自 2008 年 11 月至 2016 年 4 月整合歐盟約 3500 個機構，五千二百萬筆藏品物件的後設資料，於 2010 年發展 Europeana Data Model (EDM) 邁向資料連結 (Haslhofer & Isaac, 2011) 工作，目前則積極建構 Europeana Semantic Enrichment Framework 架構[11]針對語意加強、資料品質以及資料評估三大方向進行 (Charles, 2016)。

2014 年起美國 LD4L Labs (Linked Data for Libraries Labs) 計畫[12]，由哈佛大學圖書館創新研究室、史丹佛大學圖書館、康乃爾大學圖書館三所機構共同合作發展超越傳統後設資料的全新蒐尋方法於 2016 年起擴大研究[13]並規劃將技術成果提供另外三所機構[14]共同進行 LD4P (Linked Data for Production) [15]。該計畫主要目的為發展超越傳統後設資料的全新蒐尋方法，針對學術資源如傳統專題著作、期刊發表、檔案資料、研究資料集、圖檔影音媒體、文化器物、新聞雜誌、甚至網路典藏等，進行資料脈絡與關係之語意網路平台系統整合與建置，學術資源語意資訊倉儲 (Scholarly Resource Semantic Information Store, SRSIS) 以及知識本體的建置與維護。預期將不同單位間的資料連結的工作流程標準化，並藉由 LD4P 六所機構所產生的連結資料，同步發展技術服務。

「開放資料連結」在其他學術應用層面上，若以歷史研究領域為例，Meroño-Peñuela 等學者指出三個方向：(1) 資料連結與語意網技術所提供的控制語彙與知識本體，可解決史料欠缺正規化及其隱含知識推理的問題；(2) 資料整合則提供散落各處的獨立史料的連結機會、資料互通則提供史學家新的資料搜尋與資料擷取的機會；(3) RDF (Resource Description Framework) 資料模型提供了史料不論採取「來源導向的再現模式」(source-oriented representation) 或「模式導向或是目標導向的再現方式」(model-oriented or goal-oriented representation) 一個更彈性且易於因應不同情境脈絡變化的設計選擇 (Meroño-Peñuela et al., 2014)。

例如，在資料轉換與更新過程中，史學家主要面臨的挑戰是保持原始資料完整性，以及能追溯資料轉換的過程，而 SPARQL 這種 RDF 查詢語言，其所提供的 CONSTRUCT (根據查詢結果自動建構 RDF 圖) 與 SELECT (選擇顯示查詢結果欄位值) 等資料網絡建構與選取方式，可提供呈現不同問題觀點的需求與結果。SPARQL 查詢方式不需要改變知識庫原先的狀態即可根據不同觀點需求進行，也因此提供了傳統知識庫中使用較無彈性所建構的資料模型的替代方案。而對於資料轉換的追溯，連結資料所重視的「資料溯源」(Provenance) 則提供了解決方案 (Meroño-Peñuela et al., 2014)。

近幾年資料連結與語意網技術與巨量資料的結合不僅在定義上密切相關，同時「開放資料連結」也提供了整合異質性資料成為可理解的巨量資料 (Understandable Big Data) 的角度，

用以偵測資料不一致性，並可透過推理引擎或連結外部資料產生新知識，皆賦予巨量資料更多資料處理與利用價值 (Emani, Cullot, & Nicolle, 2015)。以南加州大學處理文化機構巨量資料為例，當面臨資料差異與資料異質整合的問題時，透過其所發展的開放原始碼工具 Karma，可針對不同資料來源與格式的資料進行整合，並利用知識本體語彙進行語意對應 (Knoblock & Szekely, 2015)。另一方面，研究亦發現自 2014 年起，行動裝置結合「開放資料連結」與語意網技術的 APP 大量快速發展，技術方面也從早期行動裝置僅扮演用戶端，語意資料處理有賴遠端伺服器，提昇至近日發展為行動裝置用戶端也具備語意推理功能 (Yus & Pappachan, 2015)，「開放資料連結」貼近常人的每日生活為期不遠。

在「知識庫」(Knowledge Base) 或稱「知識圖」(Knowledge Graph) 方面，近期也逐步採「開放資料連結」方式建置，提供如前述 Europeana 等計畫所著重的可豐富資料之連結對象。因此本研究針對以「開放資料連結」方式提供，對全球知名的五項跨領域知識庫專案，以及七項專業領域知識庫進行考察。這項考察主要為地理資訊為主，但亦簡略介紹因文化典藏需求而使用的文化藝術類索引典，以及生物主題知識庫。主要目的是探索這些知識庫因開放連結而展現一體兩面的效果：一方面積極對外連結以豐富自身知識庫語意，另一方面因豐富的知識資源亦成為其它資料庫及知識庫的連結對象。表 1 列出這些「開放資料連結」知識庫的基本資訊。

表 1 五項全球「開放資料連結」跨領域知識庫與七項專業「開放資料連結」知識庫基本資訊 (2016/11/06)

LOD 知識庫	起始	組織	資料性質	主要來源	個體量	三元組量	更新頻率	
專家建構	OpenCyc	2008	商業	跨領域	自建	41,029	2,412,520	超過一年未更新
	Getty AAT	2014	商業	文化藝術	自建	45,327	13,259,890	LOD 後一年 3-5 次
	Getty TGN	2014	商業	地名	自建	2,495,100	204,614,290	視需求
	Ordnance Survey	2010	政府	地理資訊	自建 (二者視為同一專業知識庫)	2,938,707	58,377,209	一年兩次
	Open Names	2015	政府	地名		925,157	21,360,688	
混合	EOL (TraitBank)	2014	學會	生物	整合專業資料庫/ 資料協作為輔	10,753,384	359,292,712	約一週
協同合作	Freebase	2008	商業	跨領域	Wikipedia	49,947,799	3,124,791,156	2015 關閉
	YAGO	2007	大學	跨領域	Wikipedia	5,130,031	1,001,461,786	超過一年
	DBpedia	2007	大學	跨領域	Wikipedia	5,109,890	402,086,316	約一年/ 部分即時
	DBpediaPlace	2007	大學	地名	Wikipedia	816,252	53,895,946	
	Wikidata	2012	NGO	跨領域	Wikipedia	19,367,201	1,371,170,022	即時
	LinkedGeoData	2009	大學	地理資訊	OpenStreetMap	>3 billions	1,384,887,500	約一年
	GeoNames	2010	NGO	地名	資料協作為主/ 整合地名資料庫	>6.2 millions	93,896,732	即時

### 以開放資料連結的跨領域知識庫：五項知名專案的簡介

早期即由學術資源資料的開放而成為「開放資料連結」知識庫之先驅代表者為 OpenCyc[16]：自 2002 年起以開發人工智慧的企業公司 Cyc，使用開放原始碼軟體建置知識

本體與常識知識庫 OpenCyc，而後於 2008 年邁向連結資料[17]另以「開放資料連結」版本釋出[18]，其資料來源即為該公司提供給學術社群免費使用所創建的 ResearchCyc[19]。另外 2007 年美國軟體公司 Metaweb，開發 Freebase[20]以 HTTP/JSON 為基礎的 API 和以 RDF 為端點[21]，提供機器可抽取的資料庫，開放給一般使用者自由編修資料，並提供結構化的使用者參與介面 (Bollacker, Evans, Paritosh, Sturge, & Taylor, 2008)。2010 年 Google 購買 Freebase 以此基礎建立 Google 知識圖庫，2014 年底宣布將 Freebase 資料匯入 Wikidata[22]，在 Google 支持合作下藉由 WikiProject Freebase 陸續將資料與 Wikidata 整合[23]。

相對於商業公司在「開放資料連結」知識庫的進展，學術界方面則以德國學術圈發展最為活躍。「開放資料連結」知識庫的知名代表為德國馬克斯普朗克研究所的 YAGO (Yet Another Great Ontology) [24]，以及德國萊比錫大學、Mannheim 大學與開放連結軟體公司合作的 DBpedia[25]。YAGO 源於網路搜尋引擎尚未對知識單位進行搜尋的年代，即以自動抽取維基百科及 WordNet 知識單位方式，建立大規模實體與關係連結的知識本體 (Suchanek, Kasneci & Weikum, 2007, 2008)。其中最引起關注的是 YAGO 以時空為其語意資料模型的首要元素 (first-class citizen)，以事件和情境脈絡進行物件的語意再現 (Hoffart et al., 2011; Hoffart, Suchanek, Berberich, & Weikum, 2013)，是本研究在實作時設計物件語意強化版 (Refined Versions, R 版) 的啟蒙雛型 (於實作步驟五詳述)。近年來 YAGO 亦和其他知識庫同步發展連結資料的多語系統 (Mahdisoltani, Biega, & Suchanek, 2015)，主要採取的方法是自然語言處理，因不在現階段研究範圍內，故不詳述。

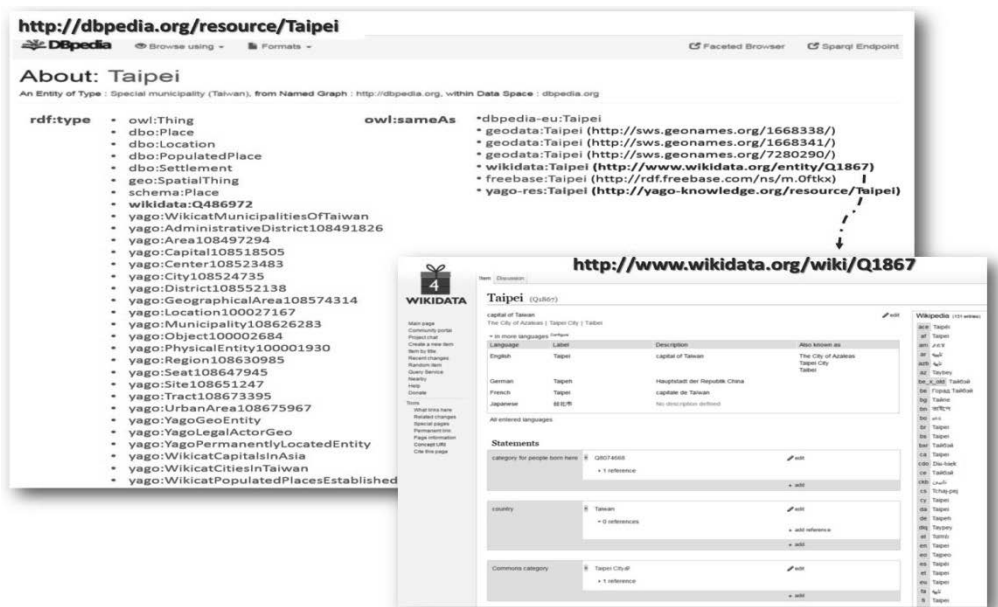


圖 2 以台北為例，DBpedia 中 Taipei 對外連結以及 Wikidata 連結外部知識庫的情形

事實上最能反映「開放資料連結」與語意網的知識庫，非 DBpedia 莫屬。DBpedia 以「維基百科語意網的反射鏡」[26]、「開放資料網核心」，「資料網結晶點」著稱，不僅與上述 OpenCyc、Freebase 與 YAGO 連結，自 2007 年 1 月初版之後約每年釋放一個版本，而後 2011 年也經 DBpedia Live[27]即時反應部分維基百科的資訊更新，相關研究更明確指出 DBpedia 在資料互通性、外部資料連結、相互連結[28]，以及多語機制設計等方面，提供了「開放資料連結」在解決問題與技術框架的領先示範(Auer et al., 2007; Bizer et al., 2009; Lehmann et al., 2015)。表 1 中亦列出 DBpedia 子資料集 DBpedia Place 作為實作連結地名時的參考(參看圖 2)，並將在探討地名「開放資料連結」知識庫時進行更詳細的分析。

以「開放資料連結」方式建置跨領域知識庫的許多專案中，目前最引起關注的是 Wikidata。在維基百科(Wikipedia)建立十年後，由維基媒體基金會於 2012 年推動以資料知識庫成為維基百科的知識架構基礎。Wikidata 清理維基百科裡的事實性資訊，整合集中使其成為可重新利用的知識庫，提供多種資料格式如 JSON、XML、RDF 等。一方面 Wikidata 類同 DBpedia 或 Freebase 一樣，抽取維基百科的結構化資訊(如 Infobox 區塊內的資訊)，另一方面也抽取資訊來源以及資料情境例如時間有效性，使其資料溯源的機制更加完備。並且 Wikidata 設計每個實體(entity)具其概念網址(concept URI)、以及其屬性名稱與屬性值所構成的陳述(statement)。這些陳述的設計相對其他知識庫更為彈性，例如可描述其屬性值是未知，或是「無/沒有」，例如澳洲「沒有鄰國」等(Erxleben, Günther, Krötzsch, Mendez, & Vrandečić, 2014; Vrandečić & Krötzsch, 2014)。其未來潛力可自上述 Freebase 的加入，以及 Europeana 的語意策略運用(Charles, Manguinhas, Alexiev, Charles, & Dammers, 2015)、DBpedia 增加比對連結(Ismayilov, Kontokostas, Auer, Lehmann, & Hellmann, 2016)、或是與專業知識庫—如 VIAF、GeoNames—的高度連結(Voß, 2016)等方面，已得多方使用與期待。以上均是本研究在實作目標連結知識庫時考量的因素。

不僅如此，若自資料品質角度觀察比較這五項知識庫，Färber、Bartscherer、Menne 與 Rettinger (2016) 以正確性、可信度、一致性、相關性、完整性、適時性、易了解性、互通性、可取得性、授權、相互連結等十一項指標研究後發現[29]：在正確性方面，在 RDF 文件驗證、文字語法驗證及「三元組」(Triple) 語意，五大知識庫大致表現優良。YAGO 在易了解性，如資源描述、多語標籤、提供多樣可理解 RDF 格式等表現最佳。Freebase 則是一致性及相互連結兩項指標的冠軍；而 DBpedia 在可取得性以及互通性上，不僅避免使用「空白節點」(blank node) 所造成無法根據 URI 參照取得資源(dereference)的問題、同時提供多種資料格式、大量使用外部語彙、且幾乎所有第二層類別(Class)資料均連結外部資源的類別，因此在此兩項指標中領先，並與 Wikidata 同列授權指標表現優異者。

然而，最引起關注的 Wikidata 在可信度、相關性、完整性、適時性、授權等五項指標上

比其他知識庫表現更佳。換言之，綜觀十一項指標，其中除正確性為五大知識庫持平外，Wikidata 在十項指標中具有五項指標最佳的優勢，若待後續 Freebase 資料陸續匯入後，亦可能延續 Freebase 在二項指標優勢而大幅超越其他知識庫。此研究結果亦呼應本研究第一階段評估綜合性知識庫作為連結目標時，選擇 Wikidata 而非 DBpedia 的方向。以下簡要討論為何在可信度方面，由專家建置的 OpenCyc 並未勝出，且 DBpedia 在資料的一致性上也未能超越 Freebase 與 Wikidata 的可能原因。

持平而論，何種指標為合適的檢驗項目，是所有品質研究可討論的議題，但進一步探討也可從 Färber 等人的指標設計上發現，以協同參與的機制建置資料、以及資料溯源資訊的完備，將會深遠影響資料品質。例如，該研究評估在一致性指標中，由於 Freebase 和 Wikidata 可由參與成員編輯，因此在用戶端界面中增加新陳述 (statement) 時，即可針對一致性進行簡易檢驗。另外，在可信度面向中，針對知識庫層級資料的匯入與策展，OpenCyc 與 Wikidata 可信度得到最佳評價，其原因為 OpenCyc 得利於專家建置而獲高可信度品質檢驗，而 Wikidata 有雙重的大眾參與機制為品質把關(資料匯入前通過維基百科社群檢驗，匯入後通過 Wikidata 社群檢視)。在陳述層級的可信度上，具「資料溯源」陳述機制為基準的專案，如 Freebase、Wikidata 與 YAGO 都能有較突出的表現，其中又以 YAGO 能儲存每一陳述的資料來源與資訊擷取技術[30]，是五項跨領域知識庫中最為獨特的代表。

### 以開放資料連結的專業領域知識庫：地名資訊專案簡介

以上探究的是跨領域知識整合的「開放資料連結」知識庫，然而根據資料特性，專業領域知識庫常是語意資料連結的主要目標。鑒於本研究產出 data.odw.tw 尚處於建置初期，基本的後設資料如時空資訊必須先組織整理，才能初步連結各項典藏品的語意。也因此地理資訊相關的知識庫如 GeoNames[31]、LinkedGeoData[32] (OpenStreetMap 資料的 RDF 呈現[33])、DBpediaPlace (DBpedia 子資料集) 與 Getty Thesaurus of Geographic Names (TGN) [34]，以及英國 Ordnance Survey 的 Open Names Linked Data[35]成為本文首要的觀察對象。目前研究方法主要以文獻探討，並實際查詢各知識庫所 (經由其 SPARQL 查詢端點) 查證比對資料。Getty TGN 為文化資產專業領域知識庫的連帶產品，在此僅引用文獻比較其與其他地名知識庫的差異，具體介紹則於介紹 Getty 語彙時詳述。

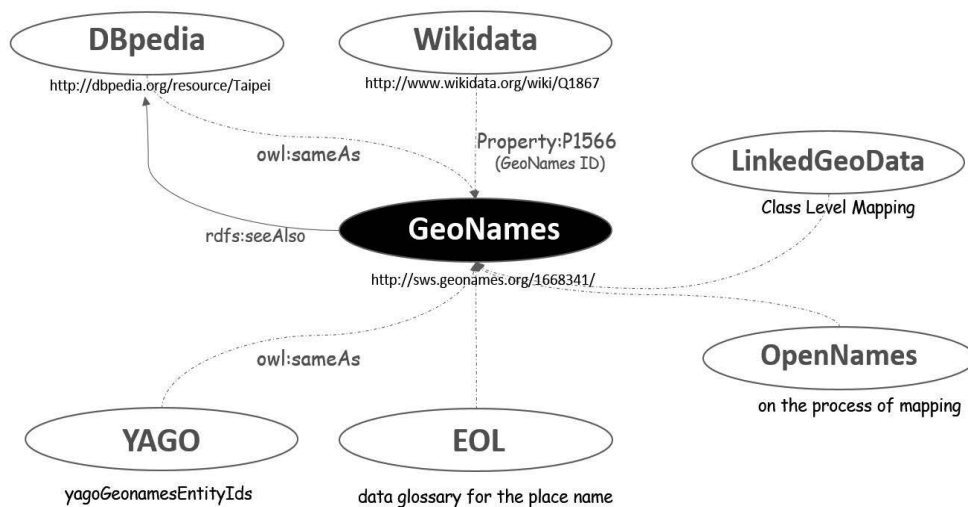


圖 3 GeoNames 與知識庫互連情形

從一項「開放資料連結」統計資料研究[36]的不同領域間最佳「開放資料連結」實作分析結果中得知，GeoNames 是第二大連結對象 (Schmachtenberg, Bizer, & Paulheim, 2014)。換言之，本研究可以合理假設大多數「開放資料連結」發佈時，均需地名資料庫作為空間參照對象，而多數的連結均選擇 GeoNames。反之，地名資料庫也需「開放資料連結」與語意網技術正規化地理相關資訊的語意關係、以及連結外部知識庫資料（地理相關或非地理相關），進而互補單一知識庫的不足（圖 3）。例如 2014 年統計顯示，約 95% 的 GeoNames 資料與 DBpedia Place 的資料並未重複、33% 的 DBpedia Place 資料與 GeoNames 資料無雷同之處 (Moura & Davis, 2014)；分析資料特性差異的研究指出 GeoNames 以地理特徵為主軸、DBpediaPlace 主要描述都市區域資訊、Getty TGN 則偏重歷史文化相關的地名資訊 (Zhu, Hu, Janowicz, & McKenzie, 2016)。

地名知識庫的連結選擇，也必須考量資料集的区域特性。若資料集所描述的資源多屬當地資源，該地區所屬的官方權威知識庫佔優勢的可能性較大。以英國為例，Ordnance Survey 是英國政府官方製圖單位，2010 年 4 月開始開放地理資料[37]同年 10 月以「開放資料連結」釋出資料集[38]，藉自建地理行政區域知識本體 (Goodwin, Dolbear, & Hart, 2008)，陸續發展郵政編碼知識本體、幾何關係知識本體等[39]。2015 年釋出基於空間關係本體所建置的「開放地名連結資料」(Open Names Linked Data) 於 2016 年 8 月約釋出約 92 萬筆開放地名所產生的二千多萬筆地名三元組[3]。地理學家的近期研究分析指出(De Sabbata & Acheson, 2016)：比較 GeoNames、Getty TGN 以及 Ordnance Survey 的 Open Names 的結果顯示，由於英國地名在 GeoNames 與 TGN 中相對於其他國家地名資料量而言，皆屬於高密度資料，因此以英國

地區性地理資料的對比代表性，應可適用其他地區。然而比較三個知識庫後，不論是地名數量、空間分佈、或是建置地名資訊創造者的不同釋義，Open Names 均呈現出較佳的表現。換言之，有地緣關係的地理知識庫理論上是有地區特性資料集考慮進行連結的較佳選項。然而以 data.odw.tw 實作為例，初期雖考量選擇臺灣地名資料庫，但本地地名知識庫品質、已否發佈為「開放資料連結」等考量，均是作為連結首選的限制因素。

OpenStreetMap[40]的特點是使用者數量與資料精細度等方面均較 GeoNames 優勢。且 LinkedGeoData 對空間資訊特徵如道路、結構關係、地貌等已進行連結資料的建構，對外連結 GeoNames, DBpediaPlace 以及聯合國農糧署[41]等知識庫，並已發佈「開放資料連結」的空間服務 (Stadler, Lehmann, Höffner, & Auer, 2012)。反觀 data.odw.tw 實作中使用的資料集，其空間部分處理的資訊較為單純 (地名萃取自藏品項目後設資料中都柏林核心集的 coverage 欄位)，因此雖在地圖視覺界面選擇使用 OpenStreetMap，但現階段僅實作連結 GeoNames。完整嚴謹的學術比較各地名「開放資料連結」知識庫的研究，目前仍是開放的研究議題。

### 以開放資料連結的專業領域知識庫：文化資產與生物資訊專案簡介

基於 data.odw.tw 實作所用資料集多為文化典藏品資料，且臺灣推動中文藝術與建築索引典 (AAT-Taiwan) [42]多年有成，本節亦嘗試解讀美國文化藝術專業機構 Getty 所發佈的索引典「開放資料連結」現況。藝術與建築索引典 (Art & Architecture Thesaurus, AAT) 主要針對文化資產物件提供詞彙、常用概念以及相關資訊；地名索引典 (Getty Thesaurus of Geographic Names, TGN) 則提供居住地、地理特徵、考古區域的地名與相關資訊；藝術家名稱聯合列表 (Union List of Artist Names, ULAN) 則提供藝術家和其他文化相關代理者的結構化人名與傳記類型資訊。

Getty 語彙自 1980 年代開始至 2014-2015 年陸續釋出開放連結資料，一方面該項工作保持其索引典在語彙結構上具有每一筆紀錄均有唯一概念的特性，以及概念間完整的相似、上下位、附屬等關係的語意架構，同時每個詞彙來源皆依照文獻保證原則以確保品質，而目前該三語彙之間的整合，也透過「開放資料連結」的方法進行比對連結，如 TGN 中的地方類型與 ULAN 中藝術家的國籍資訊的整合，近期也與 Europeana 在「開放資料連結」上合作 (Baca & Gill, 2015)。另一方面 Getty 索引典的「開放資料連結」化，也化解控制語彙被批評為一過時的知識與經驗產品、或被質疑無法順應網路時代資訊可取得性等問題。更進一步分析，語彙索引典的語意結構以開放連結方式之後，或可提供網際網路時代處理巨量資料中許多統計方法無法解決的問題，也因此「開放資料連結」方法成為索引典資源永續再生的契機 (Dextre Clarke, 2016)。

另外，基於本次研究實作採用的資料有超過三分之一為生物主題，本研究選擇了「生物大百科」(Encyclopedia of Life, EOL) [43]作為生物主題的連結目標。這裡也討論 EOL 近期所發展的開放連結物種特徵知識庫 TraitBank[44]，以及其強化 EOL 語意連結的能力。

不一致性 (inconsistency) 一直是資料品質與整合的巨大難題，然而資料語意的不一致性卻也可能是追求語意豐富的新契機。以 EOL 為例，全球生物物種的分類學系統因年代、地域、學派不同等因素，各有其獨立架構。基於不同觀點的解釋，EOL 允許多重分類，單一物種可被不同命名與分類系統所定義 (Parr et al., 2014)。例如，臺灣一葉蘭在 EOL 網頁收錄了 17 種不同的分類架構[45]，而從其獨特的分佈地域角度觀察[46]，即使是臺灣本地分類也因資料策展原始單位的不同而有差異[47]。而此臺灣特有物種的當地分類架構目前尚未包括在 EOL，對於國際生物物種的分類解釋層面而言，透過「開放資料連結」方法，是否也能成為填補全球生物知識的缺口，亦是值得繼續觀察的重點。

實際探討 EOL 在發展「開放資料連結」層面上可知，允許多重分類架構同時存在，成為發展開放式連結資料的重要環節。2014 年 EOL 開始建立物種特徵資料庫 (TraitBank)，在物種語意分類架構上，不是設計完整的語意架構整合資料，而是在物種語意分類架構延續多重分類方法。EOL 一方面採用許多不同生物知識本體以解釋單一複雜物種特徵，一方面也因現有國際語彙的不足而適時新增自訂語彙。EOL 也因此預期朝此趨勢發展，不僅增加生物領域的知識本體更廣泛與深入的研究應用，同時也能互補特定物種特徵資料庫分類與屬性資料的不足，並進一步增加新資料類型[48]與促進跨領域知識的整合 (Parr et al., 2016)。換言之，與其長期陷於眾多語彙無法取得共識的困境，對於資料的語意與連結，設計允許百家爭鳴的機制反而是最符合現實與應用的需求。而此機制也將反映到 data.odw.tw 允許多重語意精煉版本同時存在的設計理念。

### 後設資料與資料溯源的資料品質議題

品質雖是所有人對資料的基本要求，事實卻是本研究很難駁斥 Van Hooland 與 Verborgh (2014)「沒有完全乾淨的後設資料」[49]這項論點。data.odw.tw 實作案例的來源資料集，理論上雖已採用都柏林核心集十五項欄位，以統一整合異質資料來源[50]，並以 XML 結構化格式為內部資料儲存，但亦是無法完全避免標題亂碼、欄位空值、屬性值矛盾等錯誤[51]。後設資料品質議題不僅牽涉到資料建置時期的時空背景，如早期資訊技術處理資料與語意的限制、不同時期對資料要求與需求不同等因素，同時也會根據不同使用者情境，對資料品質有不同解讀。Yasser (2011) 指出，最常見的後設資料品質問題，包括不正確的資料屬性、屬性值以及系統功能性所造成的資訊遺漏，或因資料比對所造成語意資訊的喪失，以及資料再現格式不一致等。

若更進一步分析，不同領域研究者對資料品質研究的定義，以及其指標檢驗項目亦是各自不同。例如表 2 所整理的資訊、資料、後設資料、以及連結資料品質的不同觀察面向[52]。資訊管理學者的資訊物件包含資料內容與後設資料，並視後設資料為提供資訊物件的程序工具 (Stvilia, Gasser, Twidale, & Smith, 2007)；在廣義資料品質方面則注重方法論、資料種類、

以及系統面等因素，也因此其所需考量的角度（28 種面向）也更為廣泛（Batini, Cappiello, Francalanci, & Maurino, 2009）。

表 2 資料品質檢驗的不同面向

Information Quality	Data Quality	Metadata Quality	Linked Data Quality
Stvilia et al.(2007): 22 dimensions	Batini et al. (2009): 28 dimensions	Tani et al. (2013): 10 parameters	Zaveri et al. (2016): 18 dimensions
Naturalness (I)			Interoperability (RP)
Accessibility (R)	Accessibility	Accessibility	Availability (A)
Accuracy (R)	Accuracy	Accuracy (S)	Semantic Accuracy (I)
Accuracy/Validity (I)	Applicability	Pertinence	Syntactic Validity (I)
	Appropriate amount of data		
Complexity (R)	Clarity		
Precision/Completeness(R)	Completeness	Completeness(S)	Completeness (I)
Informativeness/Redundancy(R)	Comprehensiveness		Understandability (C)
Informativeness/Redundancy(I)	Conciseness		Conciseness (I)
Structural Consistency (I)	Consistency	Similarity	Consistency (I)
	Convenience		
Structural Consistency(R)	Correctness		
Verifiability (R)	Credibility		Trustworthiness (C)
Currency (I)	Currency		
Semantic Consistency(I)	Derivation Integrity		
	Ease of operation		
Naturalness (R)	Interactivity	Conformance(S)	Interlinking (A)
Semantic Consistency(R)	Interpretability		Interpretability (RP)
Precision/Completeness(I)	Maintainability	Preservability	
Complexity(I)	Objectivity		
Relevance/Aboutness(R)	Relevancy	Relevance	Relevancy (C)
Authority (Reputational)	Reputation		
Security(R)	Security		Security (A)
	Speed		Performance (A)
	Timeliness	Timeliness	Timeliness (C)
	Traceability		RP Conciseness (RP)
Cohesiveness (I)	Uniqueness	Significance	
	Usability		Licensing (A)
Volatility(R)	Volatility		
			Versatility (RP)

(I): Intrinsic; (R): Relational; (S): Metadata Spec.; (RP): Representational; (A):Accessibility; (C): Contextual

註：類似指標趨向併排

另外，圖書館界對於後設資料品質的討論包括書目資料簡易儲存、目錄資料的元素設計、以至遠端整合異質資料庫的資料脈絡重要性等，關於後設資料品質的判定，亦是眾說紛紜，有學者嘗試根據數位圖書館品質架構（Digital Library Quality Frameworks）歸納出適合後設資料語意以及數位物件的十個資料品質參數（Tani, Candela, & Castelli, 2013）。相對之下，「開

放資料連結」的資料品質則是新興議題，Zaveri 等人（2016）提出語意再現與連結是連結資料品質的基本要素，其中包括互通、語意正確、互連、可被解釋、再現的簡明、以及資料的多功能性等，均是連結資料品質所著重的面向。同時為達「開放資料連結」再次使用目的，授權資訊明確與否亦成為資料品質指標判別要素，前人這項研究所歸納 18 個連結資料品質參考面向，亦成為目前 W3C 資料品質語彙（Data Quality Vocabulary）比對 ISO/IEC 25012 資料品質模式的主要對象[53]。

即使是國際通用後設資料語彙如都柏林核心集，在語意層面也面臨概念模糊的問題，如 source 與 coverage（包含時間資訊）定義易混淆、欄位語意重疊（semantic overlaps）如 creator, contributor, publisher 間定義可相互套用、或不同單位對 relation 欄位解釋不同，因而使用方式呈現高度差異等（Park & Childress, 2009）。在此前提下，若以柏林核心集 15 項欄位為資料模型所產生的後設資料，在 Chuttur (2014) 的實證研究下指出，資料品質零錯誤的可能性為微乎其微。換言之，以上所探討品質定義多樣化、資料生成時空脈絡迥異、國際通用語彙具先天缺陷等問題，都密切攸關實作上資料語意再現的一大目標：結構性資料的再次使用以及永續價值。而對於資料品質與價值之間的衝突是否能藉由「開放資料連結」方法，轉化為和諧並存，也因實作將典藏品目錄以「開放資料連結」方式再現的過程，促使探究後設資料的資料溯源議題。

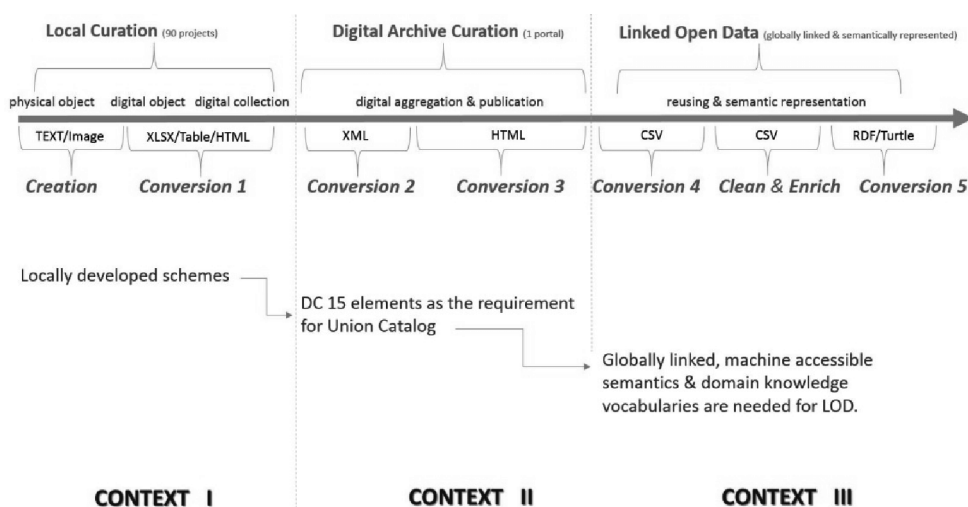


圖 4 結構資料在不同情境下的再次使用與資料轉換過程[54]

在資料價值、常人信賴、以及機器自動處理的多方期待，以及重視資料品質[55]的時代中，「資料溯源」顯然是資料品質和資料再次使用的通關護照。歐洲最古老圖書館之一的牛津博德利圖書館的積極採用可為例證（Burgess, 2016）。資料溯源也已成爲數位策展、資料引用的必要環節（Poole, 2016）。這同時，地理學家也擔憂若缺乏資料溯源，機器所提供的知識將會降低常

人解讀地方資訊的能力 (Ford & Graham, 2016)。若再就 data.odw.tw 個案為例 (圖 4)，若想達成資料再次使用以及開放連結的目標 (Context III)，首先必須確認使用的資料來源 (Context II)，因此需了解此資料的原始資料 (Context I)。再次使用而產生新資料時 (Context III)，為確保此新資料能再次被他人使用，data.odw.tw 也必須提供他人確認新資料來源的資訊、以及轉換資料的過程 (Context I + II)。追溯資料的歷史與脈絡的資訊即是所謂的資料溯源 (Carata et al., 2014)。

實際進入語意、連結、資料溯源的作法，需要考量包括資料與工作流程二種形式的資料溯源 (data and workflow provenance)、資料溯源的資料模型設計、以及資料溯源的儲存與再現 (Storage and Representation) (Omitola, Gibbins, & Shadbolt, 2010)。在資料溯源的資料模型設計方面，本研究體認案例的資料來源脈絡，將與其數位資源再次使用的效果相關。本文使用之前建立的「再次使用關聯性知識本體」(Relations for Reusing Ontology, R4R) 為基礎 (Huang & Chuang, 2014)：資料溯源資訊和要被再次使用的物件以同時打包的方式，一起提供常人與機器，以面對該數位物件被再次使用的不同情境。例如臺灣一葉蘭 (代碼 data:d2148340) [56] 藉由 r4r:hasProvenance 將資料溯源以 W3C 資料溯源知識本體 (Prov-O) [57] 描述，指出該數位資源分別在 1993、2011、2016 三個時間點，被不同地方不同單位進行了再次使用、格式轉換、以及語意連結的工作，因此也提供使用者資訊來源佐證 (圖 5、圖 6) [58]。

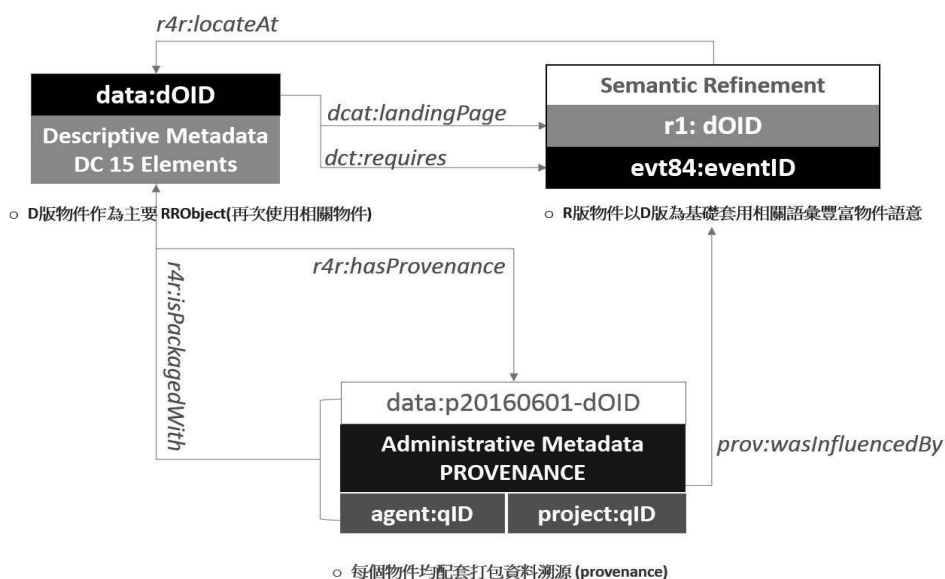


圖 5 資料模型包含 D 版 R 版資料與相對應的資料溯源

註：1. D 版：主要以 DC 15 欄位描述並富含資料溯源的再次使用相關物件的基礎版本。

2. R 版：以 D 版為基礎套用相關語彙、豐富物件語意的 Refined 版本，根據不同語彙可能在不同時期不同詮釋脈絡下，會有不同 R 版 如 R1, R2, R3...

3. 資料模型與 D 版 R 版參看實作案例步驟五更詳盡的說明。



圖 6 臺灣一葉蘭 (data:d2148340) 之資料溯源示意圖

## 實作案例：data.odw.tw

步驟一：探究共享脈絡下的資料再次使用的關聯性

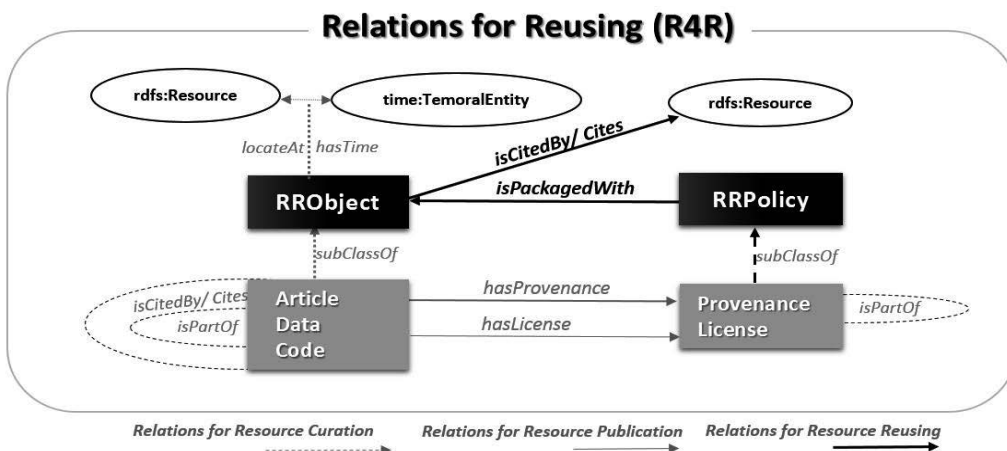


圖 7 再次使用之關聯性本體 (R4R Ontology)

以臺灣一葉蘭的標本典藏品為例，該品項歷經了不同機構的資料策展、資料發佈、資料再次使用 (curation, publication, reusing) 三個動態脈絡，以及相對應的三個描述資料語意的

階段：資料再現、資料保存、資料釋義（Representation, Preservation, Interpretation），在此不同的共享脈絡之下，「再次使用的關聯性」，亦即再次使用之關聯性知識本體 R4R( Relations for Reusing) [59]是本研究採用的理論基礎（圖 7）。

R4R 是一個簡易知識本體，以描述資源發佈（Publication）和再次使用（Reusing）的一般性關係。R4R 由兩個分立概念 RRObjcet 和 RRPolicy 組成。其中 RRObjcet 包含可分別獨立或可關聯的三組件：文件（Article）、資料（Data）、軟體碼（Code）；RRPolicy 包含可分別獨立或可關聯的二組件：資料溯源資訊（Provenance）和授權資訊（License）。關係陳述主要由 *r4r:isPackagedWith* 和 *r4r:isCitedBy* 兩個基本關係以定義。前者藉由資料套裝資料溯源與授權資訊，進行宣告資源是處於可再次使用的狀態。後者則對資源間的引用關係做描述。

## 步驟二：設計不同情境下的模式架構

我們以資料策展者角度出發，學習探索以關聯式資料庫、及開放檔案格式與開放程式碼二種不同情境的思考架構，探究發佈資料連結不同模式的運作，茲分別敘述如下：

### 一、模式一：使用關聯式資料庫進行「開放資料連結」的發佈

研究初期嘗試使用關聯式資料庫工具 D2RQ 發佈連結資料[60]，並試作數位典藏索引典與 AAT 語意描述典藏品的連結資料，結合 W3C 資料溯源知識本體，描述資料重整活動並設計 dat ontology[61]，提供機器可操作資料格式，並測試 SPARQL 語意查詢能力，如回答以下這類問題：銅琺瑯方瓶有哪些語意概念？概念侈口（器口向外張）描述了哪些器物？器物 A 和器物 B 有哪些相似的特質？

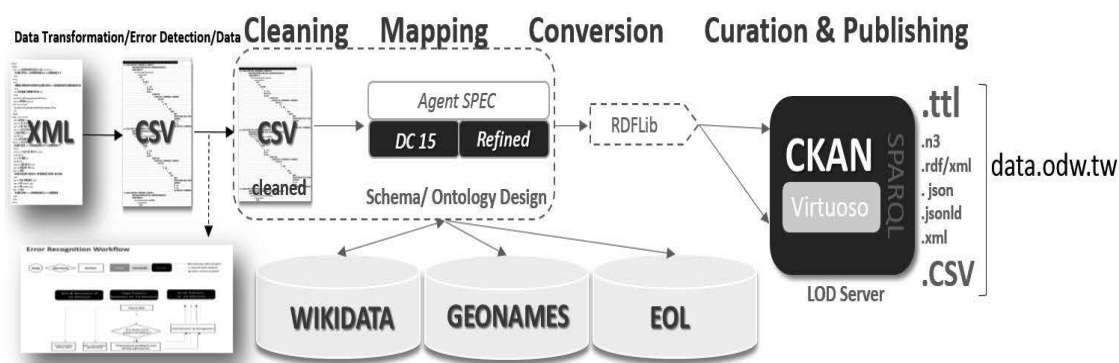


圖 8 模式二開放式連結資料的模式系統架構

## 二、模式二：開放式連結資料模式的系統架構

「開放資料連結」最具魅力的核心觀念是資料藉由開放與連結的方法，即使是其最小單位，以三元組方式呈現的單一敘述，只要能自由被常人和機器理解與操作，都是資料永續使用。也因此我們採用以開放檔案格式以及開放程式碼為基礎的資料釋出策略，將可開放的資料，同時整理為批次大量下載的結構資料檔案，如 CSV 格式檔案。並以此為基礎，使用開放原始碼程式工具，進行資料整理、清理、發佈等各階段工作。資料連結的發佈，亦採用開放檔案格式如 JSON、Turtle、XML 等供常人與機器下載使用，並同時提供常人使用的網站瀏覽介面，以及機器介接資料使用的 SPARQL 端點（圖 8）。

### 步驟三：資料剖析、清理與比對

#### 一、資料格式轉檔、清整與除錯

依模式二架構對 84 萬筆以創用 CC 授權的藏品目錄資料，進行由資料庫匯出為 XML 文件，再轉換為 CSV 格式資料表單。採用 CSV 為中介表單格式的優點包括：可參照其他資源表單、表單可人工勘誤、表單增修歷程可管理、軟體工具多、資料連結的產出方式有彈性等。過程中遇到由 XML 至 CSV 資料轉換可能遇到的問題，如 XML 樹狀結不易轉換 CSV 扁平結構、資訊遺失、或過多欄位等，因此測試多種版本後，目前採用的 CSV 版本為類似 XML 結構的格式如圖 9 所示。

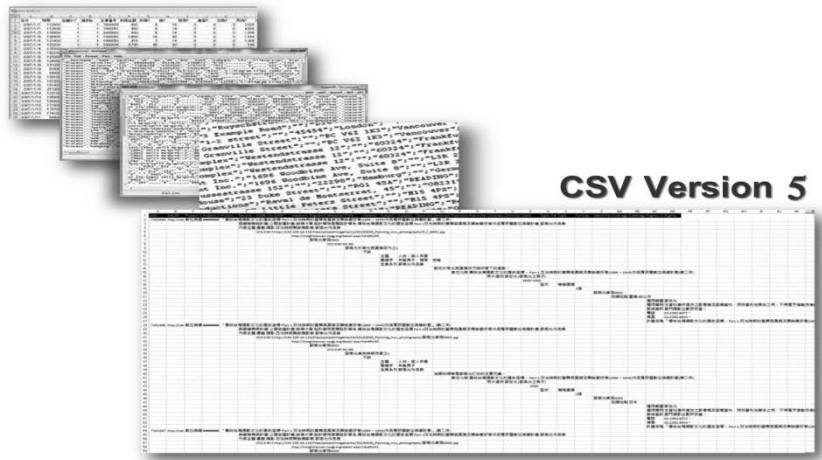


圖 9 CSV 轉置第五版

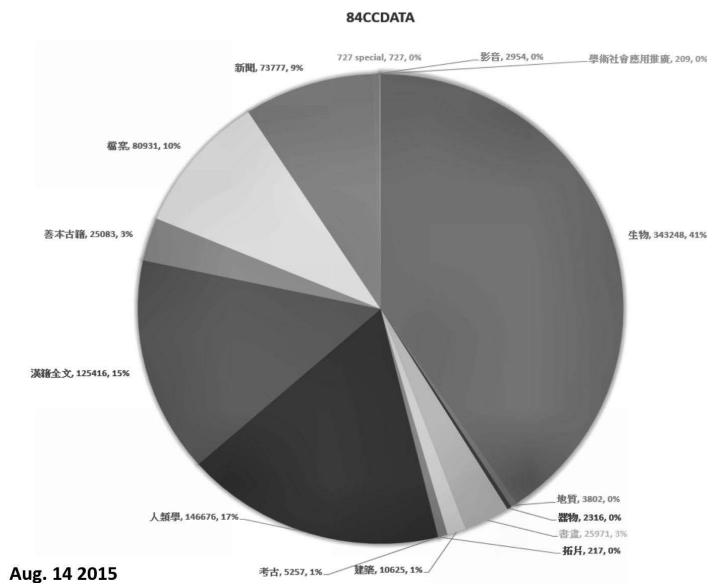


圖 10 84 萬筆資料剖析

另外分析這批資料的特性如圖 10 所示。其包括 14 個主題，內容中生物佔全部資料 41%，其次為人類學及漢籍全文。而這也是最初選擇外部專業知識庫 EOL 及 AAT 的背景因素。誠如前面所討論的資料品質問題，當資料已轉換為 CSV 後再進行設計程式檢驗歸納資料錯誤模式[62]，並產出錯誤藏品清單則相對容易。在此過程中我們發現資料具有如同都柏林核心集所遇到的問題：定義混淆、名稱模糊、編碼不一致（如時間欄位中，時間表示法的不一致）、語意超載（如 Subject 中包括 creator、contributor 欄位值[63]）、資料重複、或來自資料輸入程式錯誤等[64]。然而，資料品質牽涉廣泛，況且本研究成員並非原始資料創造者，許多資料脈絡無法取得與判別，或若因專業知識不足亦可能導至除錯反而出錯的結果（如生物命名規則中，斜體表示、加底線、問號等是允許的），若不慎亦可能在資料清理過程中，將正確資料視為錯誤或亂碼而過度清理。當然，無可避免的問題是，現有資源如時間人力經費等因素是否能支持資料清理的考量？亦可能成為促進永續資料再次使用的障礙。然而，若要達成解決前述資料品質與價值的衝突，除運用後設資料溯源外，如何才能同時保持資料品質、且減少資料清理的替代方法成為新的挑戰。

首先，一般對無效連結（broken/dead links）的看法是錯誤連結或連結資源消失，因此若非更正連結資訊，就是清除連結資料。在「開放資料連結」脈絡下，無效連結似乎更是必要的清理對象。然而美國國會圖書館的 Susan Manus 卻認為保存無效連結有其正當性[65]：一方面無效連結可協助搜尋推定原始典藏品不同版本的位置，另一方面無效連結的網址（URL）本身即是一種描述資源的後設資料。網址傳達的訊息包括網站結構，特定資源發佈日期、文

檔標題、作者、描述性關鍵字等資訊。即使主機僅為 IP 地址 URL，亦可能表示託管該域的地理區域設置。無效連結的網址包含如此豐富的語意資訊，因此成為在「資料溯源」設計原則中保留所有無效連結的主要理由。但如何以豐富語意的陳述來描述無效連結，使機器互通資料時不會回傳錯誤（404 訊息），或是在「開放資料連結」群體中在品質檢視時被視為錯誤，仍是需要研究的課題。

其次，開放資料具有協助改善資料品質如資料永續保存、增進外部驗證資料機會等益處（Janssen, Charalabidis, & Zuiderwijk, 2012），因此我們採取保留原資料 CSV，以此基礎在 data.odw.tw 以連結方式發佈原始（亦稱 D 版）資料集，以不更動原始資料為原則，僅增加資料溯源資訊，讓任何使用者欲再次使用該資源時，可根據使用者認定的資料品質定義與應用的需求，自行進行資料清理。或如同本文前面所提，經資料清理之後的 R 版資料集，允許多重分類架構原則，在此資料修正的脈絡下，R 版的另一功能為提供多重修正版本語意陳述，設計使資料清理版本也可因不同清理時間、方法、或對品質解釋不同而提供不同資料清理版本。

## 二、資料語意比對

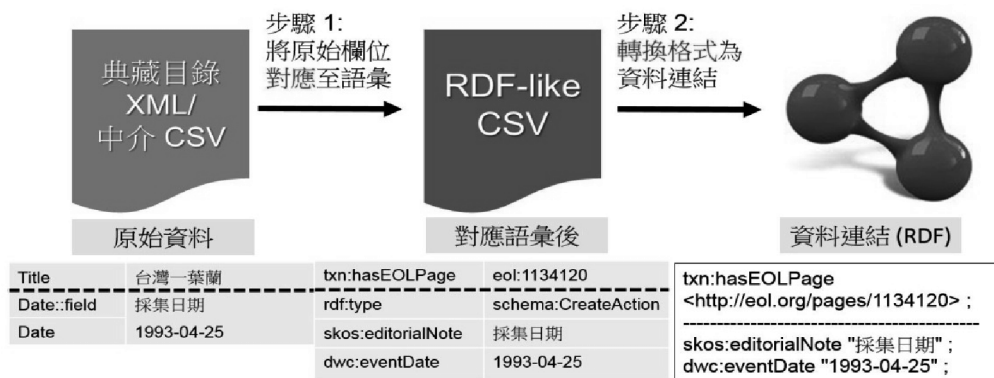


圖 11 語彙對應與格式轉換

資料品質的改善亦可借由「開放資料連結」方法達成，例如資料的語意定義使用語意網資料作為可信賴的參照對象、連結外部知識庫與協同合作再現語意、根據要求可自動驗證資料衝突、或以知識本體整合資料內容等（Fürber & Hepp, 2013）。關鍵在於語彙的運用以及資料語意的比對。實作中主要進行的工作包括時間資料正規化（如 date 欄位值若為西元時間以 ISO8601 為標準、民國年或朝代則對應到 Wikidata）、生物主題資料連結到 EOL 等。語彙對應與格式轉換包含兩步驟流程如圖 11 所示，由於以設定檔（profile）定義對應關係，因此更換語彙時僅需調整對應的設定檔。從空間資料（coverage）欄位值抽取到的地名資訊則對應到 GeoNames，使其原始資料的空間資訊理解度與查詢度大幅提升（如可回答：採集於大同

鄉的臺灣—葉蘭標本物件有那些?) [4]，同時也完成典藏機構與計劃名稱在 Wikidata 的系統代號建置與比對。完成 CSV 比對後，利用 Python 程式語言搭配函式庫 RDFLib[66]進行 Turtle 資料格式的檔案轉換。

#### 步驟四：以使用者為核心的資料庫知識平台技術架構

##### 一、CKAN 資料平台軟體

本研究使用開放原始碼套件 CKAN 建立 data.odw.tw 網站，以資料連結形式儲存與呈現典藏品資訊。CKAN 是目前開發最活躍、使用組織最多的開放原始碼資料平台軟體，包括英美澳及我國多個地方政府均以其作為開放資料平台之基礎。據官方網站於 2016 年 9 月統計 [67]，全球已有超過 140 個政府機構、社群或學術單位使用 CKAN 建置資料平台。CKAN 由開放知識國際(Open Knowledge International, OKI)於 2005 年發展，目前由 OKI 成立之 CKAN Association 維護，透過 GNU AGPL 3.0 授權條款釋出程式原始碼，本研究使用版本為 2.5.5。

因其開放原始碼之特性，機構可自行建置 (self-hosted) 系統提供服務，同時可避免被特定專有軟體 (proprietary software) 所套牢 (lock-in)。在功能方面，除基本資料發佈與存取外，CKAN 亦支援資料應用程式介面 (Application Programming Interface, API)、搜尋與條件篩選器、標籤、版本控制、分享與權限控制等功能，而可直接將資料以圖表形式呈現之「資料視覺化」功能更是其一大特色。以 Pylons 網頁開發框架寫成的 CKAN 具有現代網頁應用程式架構與極佳的自訂彈性，而 CKAN 更有為數眾多的擴充套件 (extension)，提供包含自訂後設資料、自動生成數位物件識別碼 (Digital Object Identifier, DOI)、通用資料採集介面 (harvesting，即大量匯入資料之機制) 及資料連結輸出等研究資料管理所需之各項功能。

##### 二、CKAN 資料連結支援

CKAN 在 2010 年發行的初期版本 [68] 即具有將資料集 (及所包含之資料) 之後設資料發佈為資料連結之功能 (支援 RDF/XML 與 Notation 3 格式)。而為進一步完善資料連結功能，OKI 於 2013 年啟動 ckanext-dcat 擴充套件 [69] 的開發，不僅提供更多資料連結格式 (RDF/XML、Notation 3、Turtle 與 JSON-LD) 輸出，更對應 CKAN 的資料採集介面以支援大量資料輸入 (輸入格式與輸出格式相同)。使用者除可自瀏覽器瀏覽以網頁形式呈現之藏品資料連結外，亦可於網址後方加上對應格式之副檔名 [70]，即可取得該筆藏品之資料連結，相當方便。

另值得一提的是，雖由 ckanext-dcat 名稱可知該套件對應語彙以 W3C 制定之 DCAT (Data Catalog Vocabulary) 為主，不過因其輸出入功能採用 CKAN 標準資料採集介面，而保留了擴展的彈性。本研究在此基礎之上，設計自訂資料採集介面，加上些許調整原 ckanext-dcat 套件之採集邏輯後，使其具備匯入以多種語彙描述之資料連結的能力。

### 三、操作流程與系統架構

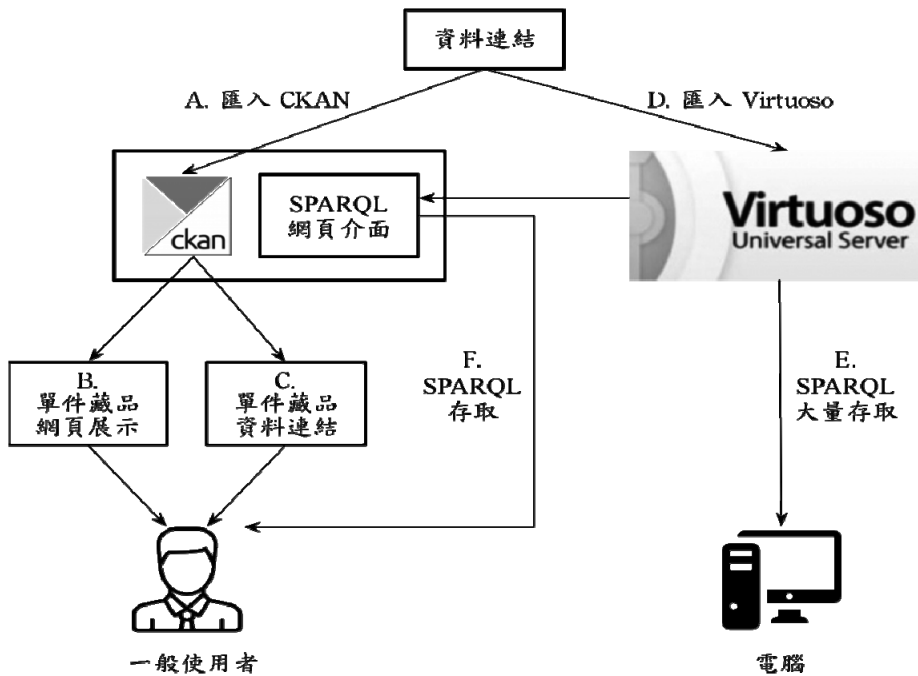


圖 12 CKAN「開放資料連結」系統架構[71]

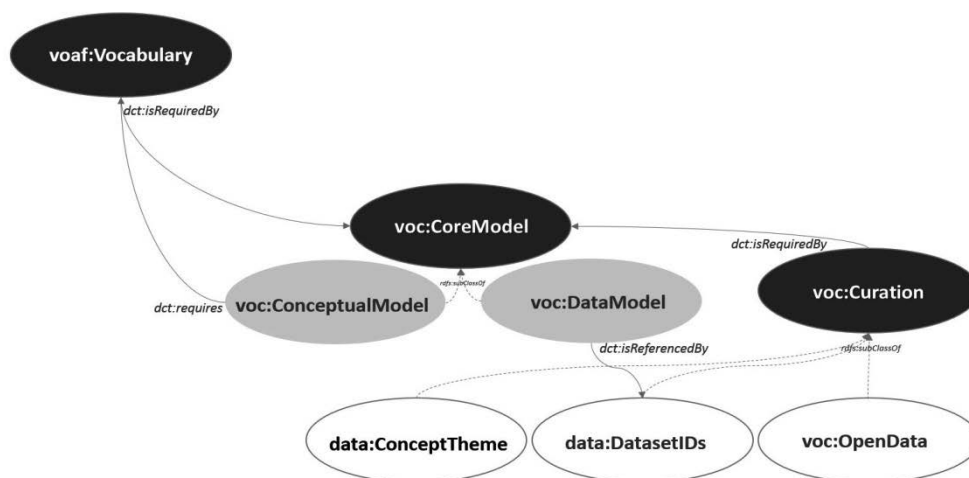
本研究建置之資料連結平台架構如圖 12 所示。轉換為資料連結之藏品後設資料，會先經由自訂資料採集介面匯入至 CKAN 平台（如圖中 A），使用者便可自網頁瀏覽單筆藏品之資料連結（如圖中 B），瀏覽介面係由 ckanext-scheming[72]與 ckanext-repeating[73]兩套件加以定義。本系統亦結合 ckanext-spatial[74]與自行開發之 ckanext-tempsearch[75]套件以分別支援空間、時間搜尋，並可與 CKAN 既有的過濾條件與關鍵字搜尋功能進行整合查詢。使用者亦可透過網頁介面，獲取單筆藏品之資料連結檔案（如圖中 C，包含前述之 RDF/XML、Notation 3 等格式）。另一方面，本系統同時將資料連結檔案匯入 OpenLink Virtuoso Open-Source Edition[76]（版本 07.20.3217）（如圖中 D），以提供 CKAN 目前缺少的 SPARQL 語意查詢功能（主要用於機器大量查詢，如圖中 E），並整合網頁查詢介面[77]於 CKAN 平台供使用者執行 SPARQL 查詢測試（如圖中 F）。相關程式碼均以 MIT 或 GNU AGPL 3.0 授權條款開放，可於 <https://gitlab.com/iislod> 取得。

### 四、系統限制與發展方向

如此搭建之平台功能雖尚屬完整，修改既有程式範圍亦不致過大，但使用較多擴充套件 [78]，且新增語彙等操作均須直接修改程式，提升維護難度；未來規劃將部分較常變動之設

定獨立為描述檔案，以降低程式複雜度。而一藏品對應產生一 CKAN 資料集的設計，所產生的大量資料集亦對匯入工作造成負擔（於 Intel E5-2620 2.1GHz、16GB 主記憶體伺服器實測，匯入 84 萬件藏品約需 2 個月時間），未來將朝改以多筆藏品彙整於一 CKAN 資料集方式匯入。

### 步驟五：透過知識本體以協助理解資料語意的再現與再次使用



Main Components of the Ontology for Open Data Web (voc4odw)

圖 13 voc4odw 知識本體主要架構

本研究實作一項知識本體 voc4odw[79]由核心主模型（Core）、策展（Curation）與國際語彙（voaf:Vocabulary）三大模型組成（圖 13）。主模型為該知識本體主要架構，並作為策展與國際語彙間的橋樑。策展模型是目前該知識本體中，連結資料、常人和機器三者溝通的管道。國際語彙則是關聯主模型參照外部常見國際語彙，並提供概念模型引用外部語彙的知識參照。

策展模型的課題主要包括資料識別、分類及資料發佈。如臺灣一葉蘭資料識別為 data:d2148340、事件 ID 為 evt84:phyCre-d2148340、策展分類（*dcat:themeTaxonomy*）為 data:Biology。而為回應全球開放資料與高規格要求資料與過程的可複製性，因此 R4R 設計再次使用機制包括 articles、data 與 code 的打包組合、以及臺灣一葉蘭多樣資料發佈格式，如 XML、JSON-LD、Turtle 或 SPARQL 端點等，均由策展模型結合 R4R 描述。

其次，關於主模型中的二大元素：概念模型與資料模型，可再細看臺灣一葉蘭的情境：1993 年 4 月 25 日（*dwc:eventDate*），臺灣一葉蘭（data:d2148340）此實體物件（*dct:PhysicalResource*）在一次採集活動（*evt84:phyCre-d2148340*）中，於地點（*gn:parentCountry*）臺灣（*gns:1668284*）

的 (*gn:parentFeature*) 宜蘭 (*gns:1674197*) 大同 (*gns:1667637*)，被製成 (*event:product*) 標本 (*dwc:PreservedSpecimen*)。此標本採集活動，若用常用語彙描述，是一個物件被創造的事件 (*schema:CreateAction*)；採集過程若用生物領域語彙 Darwin Core 來描述此脈絡 (*voc:Context*) 就是常人觀察的活動 (*dwc:HumanObservation*)。

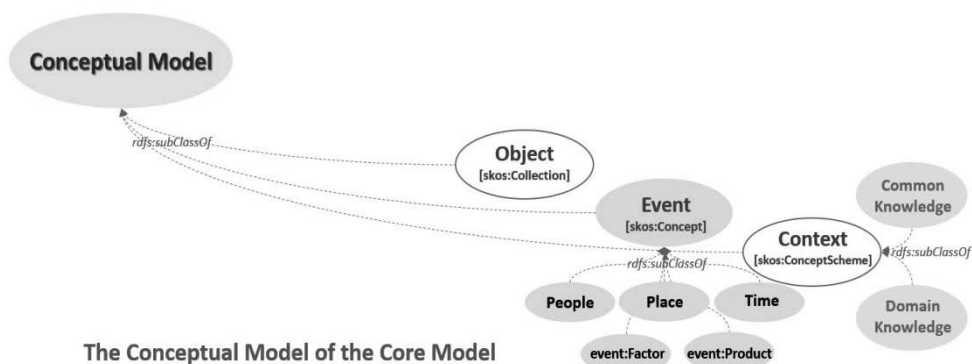


圖 14 核心模型之概念模型

依附在 *voc4odw* 知識本體下，臺灣一葉蘭透過概念模型中事件、物件、脈絡三大元素描述資源重現其語意。事件，由人時地三概念組成，而脈絡則由常用語彙如 *schema.org* 以描述，Darwin Core 則用來描述專業知識。如圖 14 所示，概念模型包括 SKOS 簡易知識組織系統的概念化模型，並闡釋由專業知識或一般性知識所關連的事件，藉此提供概念成形的架構。

觀察臺灣一葉蘭 (*data:d2148340*) 的數位化進展，會幫助讀者了解資料模型 (圖 15)。一葉蘭從實體標本物件於 (*prov:generatedAtTime*) 2011 年 5 月 13 日由 (*prov:wasGeneratedBy*) 中研院臺灣本土植物數位化典藏計畫 (*project:q21095859*) 被數位化，並在 (*prov:hadPrimarySource*) 該原始計畫網站上呈現其後設資料的資訊。這些描述是資料模型主要敘述一葉蘭數位演化的過程，藉由後設資料溯源資訊的追溯，如前圖 4 所示，臺灣一葉蘭該項資料在不同階段的脈絡所代表的不同角色：

- Context I：原始資料 (*prov:PrimarySource*)；
- Context II：目錄呈現 (*dcat:Catalog, prov:Revision*)；
- Context III：開放連結 (D 版為 *data:Reused* 與 *r4r:RRObjct*；R 版為 *data:Refined* 與 *r4r:Data*)

在資料模型中 (圖 15)，除資料溯源外另一重要模型的機制設計是衍生資料的兩個子類別設計：都柏林核心集描述版本 D 版的 *data:Reused*、以及強化語意 R 版的 *data:Refined*。

首先，D 版的 `data:Reused` 運用 R4R 語意描述模組化機制，提供基礎的都柏林後設資料 15 項欄位描述。在此所描述的資料為自原始資料中（`voc:PrimaryData`）所抽取的衍生資料（`voc:DerivationData`）。抽取資料的目的是再次使用該資源，因此宣告為 `data:Reused`，同時為此物件在 `data.odw.tw` 中賦予唯一 URI 而定義為 `r4r:RRObjct`（再次使用相關物件）。

例如臺灣一葉蘭（`data:d2148340`）在 D 版中 `rdf:type` 為 `data:Reused` 與 `r4r:RRObjct`，使用都柏林後設資料的 11 個欄位，在 RDF 描述架構中 Subject 為此資源的 URI，Property 為都柏林後設資料所對應的 URI，三元組的前二者均為連結，最後 Object 在 D 版則預設為文字。雖然國際「開放資料連結」與 Semantic Web 社群並不鼓勵發佈「開放資料連結」為文字值，但考量保存與策展原始資料最初原型的必要，我們將典藏品按原始資料不添加 Object 語意前提下，只增加該藏品資料溯源資訊，達成提供第三方後續可不被本研究語意框架限制而自由再次使用該資源的目的。

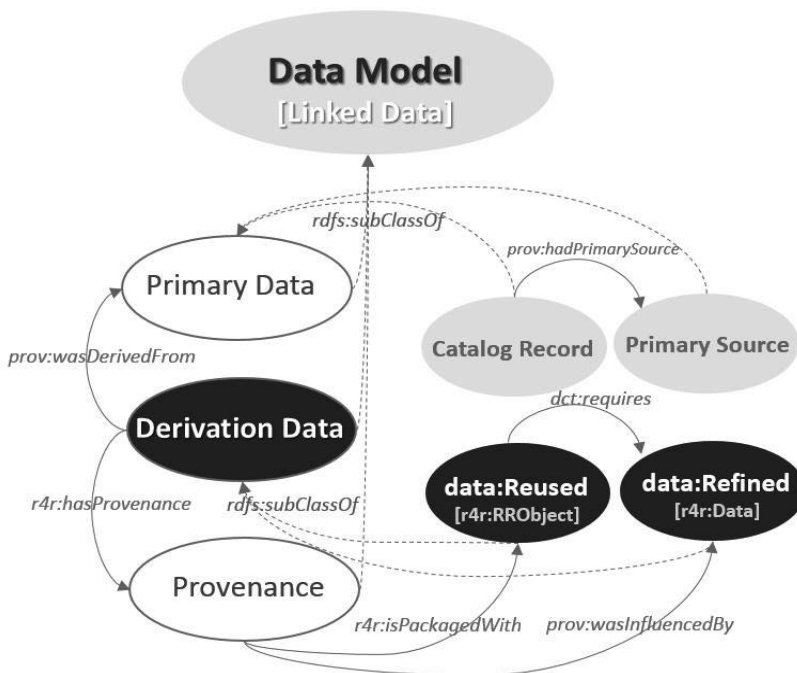


圖 15 資料模型

其次，R 版語意資料首先為自 D 版中所抽取的再次使用資料（`r4r:Data`），因此和 D 版共享同一資源 URI。D 版（`data:Reused`）需要（`dct:requires`）R 版（`data:Refined`）的語意強化與連結，因此 R 版三元組特色為三者主要是 URI 連結、或正規化後的資源如時間。進一步說明 R 版和 D 版在 R4R 的關係可知，二者為 `r4r:RRObjct` 資源唯一 URI 下之 `r4r:Data`，因此 R 版和 D

版二者並無上下位關係，而這也是現階段設計每一物件藏品的資源唯一 URI 在 CKAN 中宣告為 `dc:Dataset` 描述此物件藏品的資料集是集合不同版本、該資源不同檔案格式的各式資料。

例如臺灣一葉蘭 (`data:d2148340`) 在 D 版中 *coverage* 的欄位值為文字“行政區：宜蘭縣大同鄉”在目前 R1 版中則借由概念模型的事件描述，以 GeoNames ontology 與地名資源 URI 描述[80]，這二者所包含的三元組各均屬於 `data:d2148340` 資料集，而目前提供的資料集格式則包含 XML、JSON-LD 與 Turtle。換言之，臺灣一葉蘭雖在目前 R 版中連結 EOL 豐富語意，但目前因時間人力資源等限制只試作生物類三計畫，且回顧比對工程之時間點 EOL TraitBank 尚未釋出「開放資料連結」資料，也因此僅以連結網頁而非資源 URI 方式處理[81]，後續則可在不同 R 版完成所有 EOL 比對連結，或連結 TraitBank URI，或連結其他知識庫。更詳盡的臺灣一葉蘭 (`data:d2148340`) 在 D 版與 R 版的應用案例，可參看該資源物件在 LOD 脈絡前後不同版本的說明 (Lee, Huang, & Chuang, 2017)。簡言之，不同比對連結時機、不同語意解釋架構、或不同推導過程，則可以不同語意化資料版本策展 (例如，R1, R2, R3...)。以下歸納 R 版的多重機制主要有三：

### 一、多重清理版本機制

資料清理本身即是一種語意解釋過程，例如前述資料品質問題中提到的標題亂碼案例，在 D 版有一案例為 `data:d4653940`。我們並未對該標題進行清理或刪除原因有三：(1) 標題應由資料創造者定義，(2) 該資源原始資料[82]顯示此為一植物學遺傳研究頁面，若技術資源足夠時，因 `data.odw.tw` 已提供完整的後設資料溯源，大量自動化訂正錯誤是可能的，(3) 若未來採用開放大眾參與編修，此資源亦可能被修正，因此未來亦可能在不同時間由不同修正者提供不同修正版本。

### 二、多重知識連結機制

連結知識庫的目的不在於為特定資源提出完全正確的語意，各專業知識內自有理論解釋，而是要提供檢視資源的多種面向，了解其中有哪些概念能有助於處理當下使用者問題。再以台北的描述為例，TGN 知識樹分類以台北上層為國家／島嶼／特別都市、Wikidata 連結 14 種不同知識庫的 Taipei 編號、而 DBpedia 則描述台北有 34 種 `rdf:type` 知識分類 (圖 2)。這種多重知識的觀點，為連結至不同知識庫的典藏資源賦予了強化語意知識的基礎。另外也因知識本身會隨著時間演化，今日的事實可能是明日的謬論，對典藏品而言，連結的是固定 URI 而不是知識內容，一當知識庫知識內容更新，典藏品不需更新知識本身，透過資料連結，知識重新發現重新表達。

### 三、多重語意架構機制

對語意架構、語彙選擇、漸進式增加資料的語意深度等不同的需求劃分開來，並且各自獨立，即是 R 版多重語意架構機制主要目的。舉例而言，目前 R 版只針對生物類資源使用 dwc 語彙，然而描述物種語彙至少有 30 種[83]，目前 R 版只處理都柏林核心 15 項欄位中時空資訊以及部分生物類標題，因此類似生物類中 subject 所描述「界門綱目科屬種」等知識架構語意，如前文所述，多重 R 版機制的設計使其更具彈性。另外，目前 84 萬筆資料由 14 個內容主題構成，後續若時機許可亦可套用不同專業領域的語彙發佈不同 R 版，同時若有其他非典藏目錄新資料集加入，亦可據資料特性調整，套用現有 D 版架構、根據使用者需求調整 R 版，同時發佈該資料集之 R 版。此多版精鍊機制允許互相獨立，甚至是互不相容的語意需求，顯示了這種資料連結新式架構，達成了最適合策展者與使用者其雙方需求最終可調和的本質。

#### 國際語彙的運用 (voaf:Vocabulary)

鑑於 Schaible、Gottron 與 Scherp(2014)實證研究指出，目前語彙再次使用策略有六：(1) 再次使用常用語彙，(2) 自訂語彙並使其與外部常用語彙連結，(3) 語彙再次使用最大化，(4) 語彙再次使用最小化，(5) 語彙個別概念再次使用最小化，(6) 僅再次使用專業領域語彙。同時並建議與其限制語彙數目，不如採取再次使用共用語彙的策略。這與 Srinivasan、Becvar、Boast 與 Enot (2010) 提出的「多重知識本體」(multiple ontologies) 觀點相互輝映。該觀點認為對於一物件的不同了解和詮釋間的張力是應該被接受的[84]，因此使用 voaf:Vocabulary (在資料連結雲中所使用的語彙) 作為主要類別，借此關聯主模型 (Core Model) 至外部常見國際語彙 (圖 13)。

可再次使用的語彙資源可透過 Linked Open Vocabulary (LOV) [85] 了解與選擇目前國際語彙的使用。此資源網站目前並未囊括已「開放資料連結」發佈所有資料集中所有使用語彙以及索引典資源。然而對使用者而言，LOV 可依據語彙作者或單位、語彙名稱、語彙單詞如 class 或 property 名稱、專業知識分類查詢、SPARQL 查詢、以及蒐錄語彙時必須通過機器與 LOV 專業人員的審查等優點 (Vandenbussche, Atemezing, Poveda-Villalón, & Vatan, 2015)。本研究設計 voc4odw 時，得利於使用此國際開放語彙服務，進而採用 25 個國際語彙如表 3 所示：包含 W3C 標準語彙如：csvw、dcat、org、prov、skos、time；一般常用語彙如：cc、dc、dct、event、foaf、r4r、schema.org；以及專業知識語彙如：aat、dwc、geo、gn、txn。

表 3 **voc4odw** 知識本體命名空間

Common Knowledge		
Prefix	Namespace	Description
cc	<a href="http://creativecommons.org/ns#">http://creativecommons.org/ns#</a>	1. Creative Commons Rights Expression Language
csvw	<a href="http://www.w3.org/ns/csvw#">http://www.w3.org/ns/csvw#</a>	2. W3C CSVW Namespace Vocabulary Terms
dc	<a href="http://purl.org/dc/elements/1.1/">http://purl.org/dc/elements/1.1/</a>	3. DC 15 (Dublin Core Metadata Element Set)
dcat	<a href="http://www.w3.org/ns/dcat#">http://www.w3.org/ns/dcat#</a>	4. W3C Data Catalog Vocabulary
dct	<a href="http://purl.org/dc/terms/">http://purl.org/dc/terms/</a>	5. DCMI Metadata Terms
dctype	<a href="http://purl.org/dc/dcmitype/">http://purl.org/dc/dcmitype/</a>	6. DCMI Type Vocabulary
event	<a href="http://purl.org/NET/c4dm/event.owl#">http://purl.org/NET/c4dm/event.owl#</a>	7. Event Ontology
foaf	<a href="http://xmlns.com/foaf/0.1/">http://xmlns.com/foaf/0.1/</a>	8. FOAF Vocabulary Specification
geo	<a href="http://www.w3.org/2003/01/geo/wgs84_pos#">http://www.w3.org/2003/01/geo/wgs84_pos#</a>	9. W3C WGS84 Geo Positioning: an RDF vocabulary
gn	<a href="http://www.geonames.org/ontology#">http://www.geonames.org/ontology#</a>	10. GeoNames Ontology
gns	<a href="http://sws.geonames.org/">http://sws.geonames.org/</a>	11. GeoNames Entity
lcsb	<a href="http://id.loc.gov/authorities/subjects">http://id.loc.gov/authorities/subjects</a>	12. Library of Congress Subject Headings
org	<a href="http://www.w3.org/ns/org#">http://www.w3.org/ns/org#</a>	13. W3C Organization Ontology
prov	<a href="http://www.w3.org/ns/prov#">http://www.w3.org/ns/prov#</a>	14. W3C Provenance Ontology (PROV)
r4r	<a href="http://guava.iis.sinica.edu.tw/r4r/">http://guava.iis.sinica.edu.tw/r4r/</a>	15. Relations for Reusing Ontology (r4r)
schema	<a href="http://schema.org/">http://schema.org/</a>	16. Schema.org
skos	<a href="http://www.w3.org/2004/02/skos/core#">http://www.w3.org/2004/02/skos/core#</a>	17. W3C Simple Knowledge Organization System (SKOS)
time	<a href="http://www.w3.org/2006/time#">http://www.w3.org/2006/time#</a>	18. W3C Time Ontology
voaf	<a href="http://purl.org/vocommons/voaf#">http://purl.org/vocommons/voaf#</a>	19. Vocabulary of a Friend (VOAF)
wde	<a href="http://www.wikidata.org/entity/">http://www.wikidata.org/entity/</a>	20. Wikidata Entity
Domain Knowledge		
aat	<a href="http://vocab.getty.edu/aat/">http://vocab.getty.edu/aat/</a>	1. Art and Architecture Thesaurus
dwc	<a href="http://rs.tdwg.org/dwc/terms/">http://rs.tdwg.org/dwc/terms/</a>	2. Darwin Core Terms
dwciri	<a href="http://rs.tdwg.org/dwc/iri/">http://rs.tdwg.org/dwc/iri/</a>	3. Darwin Core terms
eol	<a href="http://eol.org/pages/">http://eol.org/pages/</a>	4. The Encyclopaedia of Life (EOL)
txn	<a href="http://lod.taxonconcept.org/ontology/txn.owl#">http://lod.taxonconcept.org/ontology/txn.owl#</a>	5. Taxon Concept OWL Ontology
Local Namespace		
voc	<a href="http://voc.odw.tw/ontology#">http://voc.odw.tw/ontology#</a>	1. Ontology for data.odw.tw (voc4odw)
agent	<a href="http://data.odw.tw/agent/">http://data.odw.tw/agent/</a>	2. Organization/Agent Entity
article	<a href="http://data.odw.tw/article/">http://data.odw.tw/article/</a>	3. Textual Description with <i>rdf:type</i> r4r:Article
code	<a href="http://data.odw.tw/code/">http://data.odw.tw/code/</a>	4. Code Description with <i>rdf:type</i> r4r:Code
data	<a href="http://data.odw.tw/record/">http://data.odw.tw/record/</a>	5. Linked Data for data.odw.tw
evt84	<a href="http://data.odw.tw/event/">http://data.odw.tw/event/</a>	6. Event Entity in data.odw.tw
project	<a href="http://data.odw.tw/project/">http://data.odw.tw/project/</a>	7. Project Entity in data.odw.tw
r1 (n)	<a href="http://data.odw.tw/r1/">http://data.odw.tw/r1/</a> (r2, r3...)	8. Refined Version(s) of data.odw.tw Entity
refined	<a href="http://data.odw.tw/refined/">http://data.odw.tw/refined/</a>	9. Directory of the Refined Versions
catdat	<a href="http://catalog.digitalarchives.tw/">http://catalog.digitalarchives.tw/</a>	10. Union Catalog of Digital Archives Taiwan

## 結論與建議

總結現階段實作案例的成果面向有三：

### 一、資料面向

提供使用者可選擇的多種連結資料格式、以及相關的轉換程式碼，協助使用者進行符合自身需求的資料語意處理與再次使用，本研究實驗產生的結構資料三元組約四千五百萬筆。

### 二、技術面向

使用並擴充 CKAN 軟體套件，發展為「開放資料連結」的儲存與展示平台，收納整合前述之連結式資料，建置常人與機器皆可瀏覽與操作的「開放資料連結」系統 data.odw.tw。

### 三、語意面向

已建置精鍊語意版的資料集（R 版），並連結知名的知識庫（如 Wikidata、GeoNames、EOL 等），此資料集約有二千五百萬筆。並透過知識本體 voc4odw 的設計，運用語意描述的模組化機制，提供基礎都柏林核心集的描述版本（D 版）與資料溯源，更進一步針對時空語意描述加強、設計 R 版彈性多重機制（資料多重清理、知識多重連結、語意多重架構），讓資料策展者或使用者能重新建構資料語意連結，走向資料近用（Access）與再次使用（Reuse），最終達成達「物盡其用」之資料盡用、語意整合、知識連結的語意網世界。

最後回應本文前言所討論之「開放資料連結」發展種種挑戰：本研究在 CKAN 技術的運用與開發不僅解決（1）技術工具支持與（4）人性化瀏覽與查詢界面等問題；另一方面知識本體 voc4odw 的建構與實作也提供面對（2）資料品質控制機制問題時，主要資料的溯源（Provenance）與資料多重清理機制；而在（3）資料模型與語彙的實作挑戰下，設計允許國際語彙彈性使用、以及不同語彙同時並存的語意多重架構。另外在（5）定義資料開放授權的困難，以及（6）缺少新技術知識的技術人員等方面，則分別提供建議如下：

就資料開放授權而言，Ermilov 與 Pellegrini（2015）提出「開放資料連結」資料需注意多層（multi-layers）授權、綜合（compositive）授權、以及授權必須符合人類、機器皆可讀、且能適用真實法律案件的複雜性等因素。在此撰寫本文的時間點，本研究雖尚未將資料集正式在 LOD Cloud 發佈，但目前 data.odw.tw 已針對 D 版藏品物件按照原資料 9 種 CC 授權方式釋出，R 版日後亦將增加機器可讀的 CC0 授權聲明，提供未來資料使用者可根據使用者情境，賦予不同資料再現的自由與彈性。

就新技術知識與人員問題，臺灣的實際現況實處於缺乏「開放資料連結」與 CKAN 人才的窘境，然身處一個目前個人手機裝置已能處理資料連結與語意技術的同時，預期不遠的未來將是數十億人指尖下語意連結即時、富語意、多樣態的資源與快速互動。也因此我們不免

在此提出對臺灣此現況的建議：盡速培養可同時關照資料、常人和機器三者語意連結技術的人才，才能反應臺灣在全球「開放資料連結」的行動與位置。

## 致謝

特別感謝中央研究院資訊科學研究所陳克健研究員所帶領的典藏臺灣聯合目錄團隊過去幾年的支持與協助。本研究尤其感謝洪崇熙先生在資料收集與處理的幫助、以及曹晉豪先生與陳心萍小姐在資料清理與比對的協力。

## 附註

[1]名詞翻譯說明

**開放資料連結 (Linked Open Data, LOD)**：本文視前後文語意，將 Linked Data 翻譯為「資料連結」或「連結資料」。「資料連結」意重資料的性質，亦即資料間是連結的。將動詞置於名詞之後，以修飾之前名詞的性質，常見於華語。例：「四輪傳動」、「峰峰相連」。Linked Open Data 則可翻譯為「開放資料連結」，即「開放資料」後加以「連結」動詞為修飾，描述其彼此間有連結。若文意著重於動作本身（而非資料本身），視情形將 Linked Data 翻譯為「連結資料」，意指將資料彼此間加以連結的動作。

**語意再現 (Semantic Representation)**：依照國家教育研究院雙語詞彙、學術名詞暨辭書資訊網 (<http://terms.naer.edu.tw/detail/538572/>) 中 “representation” 在不同學術領域有不同的翻譯名詞，然在語意網技術中，以 RDF 為基礎的 semantic representation 強調 “machine-processable representations”，因此本文作者採用具有以 RDF model “再現” 原始資源／資料含意的名詞。

[2]<http://catalog.digitalarchives.tw/>

[3]本文所有 SPARQL 查詢請參考 <http://data.odw.tw/examples/JLIS-2016-query.html>

[4]包括宜蘭三星大同棲蘭山林道；花蓮秀林萬榮；苗栗泰安南庄，新竹尖石鴛鴦湖、及嘉義縣等地。

[5]官方網站位於 <http://ckan.org>。

[6]若以實作機器可讀取可操作的知識本體建構而言，在本研究中至少包括二個階段

(1) 步驟一階段：設計 R4R ontology 階段，此階段為純理論探討，定義資料再次使用的基本架構，尚無資料集實作考量；(2) 步驟一至五階段：設計 data.odw.tw 所需的知識本體 voc4odw 階段，若就時間點而言實則涵蓋步驟一、二、三、四、五所有階段。如知識本體中主模型分別在步驟一、二、三前期成型（僅限 D 版，並實作 D 版 triple），步驟三後期、步驟四則針對語意加強的 R 版進行更進一步國際語彙模型的加強與根據 CKAN 提供的策展能力，進行資料策展模型的設計。因此將完整的知識本體列於最後，才能完整介紹。又若以非機器

可讀取可操作的人類理解資料模型而言，本研究則是在進行格式轉換試時，先設計資料模型草稿（schema draft）進行資料與 CKAN 交互測試，測試過程中不斷修改，待 CKAN 測試完成時，一方面人工隨機抽樣檢視、另一方面輔以 SPARQL 查詢驗證資料結果，確認資料品質後再根據測試最終版的 data schema 進行建構機器可讀取可操作的知識本體建構。簡言之，不論是一直在過程中變動的 schema，或變動的知識本體均是「開放資料連結」實作時不可避免的挑戰，而此經驗亦促成最後 CKAN 能支援對同一物件運用不同 schema 描述的功能，亦是 voc4odw 中策展模型下，實作多重精鍊版本機制誕生的可能。

[7]“a surprising amount of data isn't linked in 2006” at <https://www.w3.org/DesignIssues/LinkedData.html>, 2007 年約 12 個資料集。

[8]超過 149,423,660,620 的「開放資料連結」Triples 來自 2973 資料集，而此數字不包括總「開放資料連結」資料集超過 9960 與各機構組織未開放的 Linked Data 的統計內。數字來自：<http://stats.lod2.eu/>，<http://lodstats.aksw.org/>，2017 最新公布 The Linking Open Data Cloud Diagram (<http://lod-cloud.net/>)，本文最後更新時間：2017/03/28。

[9]OCLC 視 LOD 與知識庫計畫為該機構 Data Science 要項 (<http://www.oclc.org/research/themes/data-science.html>)，此研究報告為 Karen Smith-Yoshimura 於在 2016 年 4 月 CNI Spring Membership Meeting 簡報：Linked Data Implementations—Who, What and Why? <http://www.oclc.org/content/dam/research/presentations/smith-yoshimura/oclcresearch-linked-data-implementations-cni-2016.pptx>

[10]<http://www.europeana.eu/>

[11]<https://docs.google.com/document/d/1JvjrWMTpMIH7WnuieNqcT0zpJAXUPo6x4uMBj1pEx0Y/>

[12]<http://www.ld4l.org/>

[13]<https://www.ld4l.org/ld4l-labs/>

[14]Columbia University, Library of Congress, Princeton University

[15]<http://www.ld4l.org/ld4p/>

[16][www.opencyc.org](http://www.opencyc.org) ; <http://sw.opencyc.org/> (最後使用日期：2016/11/18)

[17][www.dataversity.net/opencyc-hooks-into-linked-data-web/](http://www.dataversity.net/opencyc-hooks-into-linked-data-web/)

[18]透過 Open Source Texai Project 發佈 RDF 相容的格式。 <https://sourceforge.net/projects/opencyc/files/> (最後使用日期：2016/11/18)

[19][www.cyc.com/platform/researchcyc/](http://www.cyc.com/platform/researchcyc/)

[20][www.freebase.com](http://www.freebase.com)

[21][rdf.freebase.com/](http://rdf.freebase.com/)

[22][www.wikidata.org](http://www.wikidata.org)

[23][www.wikidata.org/wiki/Wikidata:WikiProject\\_Freebase](http://www.wikidata.org/wiki/Wikidata:WikiProject_Freebase)

[24][www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/](http://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/)

[25][dbpedia.org](http://dbpedia.org)

[26][wiki.dbpedia.org/online-access/DBpediaLive](http://wiki.dbpedia.org/online-access/DBpediaLive)

[27][live.dbpedia.org](http://live.dbpedia.org)

[28]超過 30 個以上外部連結資料庫如 Amsterdam Museum, BBC, Eurostat Linked Statistics, CIA World Factbook, GeoNames, GeoSpecies, LinkedGeoData, New York Times, OpenCyc, WordNet, YAGO...等。 <http://wiki.dbpedia.org/Downloads2015-10>

[29]正確性(Accuracy),可信度(Trustworthiness),一致性(Consistency),相關性(Relevancy),完整性(Completeness),適時性(Timeliness),易了解性(Ease of understanding),互通性(Interoperability),可取得性(Accessibility),授權(Licensing),相互連結(Interlinking)

[30]見 [yago:extractionTechnique](http://wiki.cfcl.com/Projects/YAGO/Predicates) 與 [yago:extractionSource](http://wiki.cfcl.com/Projects/YAGO/Predicates) at [wiki.cfcl.com/Projects/YAGO/Predicates](http://wiki.cfcl.com/Projects/YAGO/Predicates)

[31][sws.geonames.org](http://sws.geonames.org)

[32][linkedgeodata.org](http://linkedgeodata.org)

[33][www.openstreetmap.org](http://www.openstreetmap.org)

[34][vocab.getty.edu/tgn/](http://vocab.getty.edu/tgn/)

[35][data.ordnancesurvey.co.uk/datasets/opennames](http://data.ordnancesurvey.co.uk/datasets/opennames)

[36]<http://lod-cloud.net/>

[37]<https://www.ordnancesurvey.co.uk/blog/2010/04/os-opendata-goes-live/>

[38]<http://data.ordnancesurvey.co.uk/datasets/os-linked-data/about>;但連結資料則自 2009 年 10 月對外公布 <http://lists.w3.org/Archives/Public/public-lod/2009Oct/0136.html>

[39]<http://data.ordnancesurvey.co.uk/ontology>

[40][openstreetmap.org](http://openstreetmap.org)

[41][aims.fao.org/standards/agrovoc/linked-open-data](http://aims.fao.org/standards/agrovoc/linked-open-data)

[42]<http://aat.teldap.tw/>

[43]<http://eol.org/>

[44]<http://eol.org/traitbank>

[45]<http://eol.org/pages/1134120/>

[46]<http://eol.org/pages/1134120/maps>

[47]參見臺灣—葉蘭 [data:d2148340](http://data.d2148340) 與 [data:d4542169](http://data.d4542169) 中 [dc:subject](http://dc.subject) 描述差異。

[48]EOL 新增資料類型方法尚包括使用文字探勘、公民科學參與、標本資料數位化等項目。

[49]“All metadata is dirty.”

[50]全部典藏品整合超過 90 計畫單位，就 84 萬筆 CC 授權資料而言，則來自 74 個跨領域單位。

[51]<http://catalog.digitalarchives.tw/item/00/47/03/74.html>

[52]此表源自本研究於 2015 年 12 月調查報告 (<http://goo.gl/pPUXcd>)，當時並比較 W3C 的 Data Quality Vocabulary 所提十面向 (<https://www.w3.org/TR/2015/WD-vocab-dqv-20150625/>)，目前最新 dqv 尚未進入正式推薦標準語彙， (<https://www.w3.org/TR/2016/NOTE-vocab-dqv-20160830/>)，且觀察不同版本間的變動差異甚多，因此先不列入此表。

[53]<https://www.w3.org/TR/vocab-dqv/#mapping-ISOZaveri>

[54]圖 4 為本文作者 2015 年所製 (<http://goo.gl/pPUXcd>) 簡報中第十頁的再製。

[55]<https://www.w3.org/DesignIssues/UI.html>

[56]<http://data.odw.tw/record/d2148340> 本文所使用的命名空間請參照表三。

[57]<http://www.w3.org/TR/prov-o/>

[58]目前臺灣一葉蘭 (data:d2148340) 後設資料溯源為 data:p20160530-d2148340 與 data:p20160912-d2148340，若有新的版本，後設資料溯源會根據資源產生日期生成。

[59]R4R 知識本體中英文以及 RDF/Turtle 檔案下載可參看：<http://guava.iis.sinica.edu.tw/r4r/>

[60]<http://dat.digitalarchives.tw/>

[61]<http://dat.digitalarchives.tw/ontology.html>

[62] <https://gitlab.com/iislod/pattern-statistics-and-error-summary>

[63]<http://catalog.digitalarchives.tw/item/00/00/46/14.html>

[64]<http://catalog.digitalarchives.tw/item/00/3a/3d/14.html>

[65]<http://blogs.loc.gov/thesignal/2012/03/the-value-of-a-broken-link/>

[66]<https://rdflib3.readthedocs.io>。

[67]「CKAN instances around the world」頁面：<http://ckan.org/instances/>。

[68]版本 0.11 於 2010 年 1 月釋出。見 <http://docs.ckan.org/en/latest/changelog.html#v0-11-2010-01-25>

[69]<https://github.com/ckan/ckanext-dcat>。最早的提交 (commit) 時間為 2013 年 7 月 3 日。

[70]以「鋼鐵沈思少女」 (<http://data.odw.tw/record/d4502674>) 藏品為例，若欲取得 Turtle 格式之資料連結，只需在其後加上.ttl (<http://data.odw.tw/dataset/d4502674.ttl>) 即可。

[71]Icon made by SimpleIcon (<http://www.flaticon.com/authors/simpleicon>) and Freepi (<http://www.flaticon.com/authors/freepik>)

[72]<https://github.com/ckan/ckanext-scheming>。

[73]<https://github.com/open-data/ckanext-repeating>。

[74]<https://github.com/ckan/ckanext-spatial>。

[75]<https://gitlab.com/iislod/ckanext-tempsearch>。

[76]<http://virtuoso.openlinksw.com/>。

[77]該查詢介面位於 <http://data.odw.tw/sparql>。

[78]ckanext-dcat、ckanext-harvest、ckanext-scheming、ckanext-repeating、ckanext-spatial、ckanext-tempsearch，計有六個擴充套件。

[79]voc4odw Ontology: <http://voc.odw.tw/ontology/>；本文所使用的 namespace 請參照表 3。

[80]{`evt84:phyCre-d2148340 gn:parentFeature gns:1667637`}

[81]{`data:d2148340 txn:hasEOLPage <http://eol.org/pages/1134120>`}

[82]<http://literature.tfri.gov.tw/atlas/content1.jsp?item=45600>

[83][http://lov.okfn.org/dataset/lov/terms?q=species&vocab\\_limit=0](http://lov.okfn.org/dataset/lov/terms?q=species&vocab_limit=0)

[84]“which accepts the tensions that lie between different interpretations and understandings of an object.”

[85]<http://lov.okfn.org/>

## 參考文獻

- Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., & Ives, Z. (2007). DBpedia: A nucleus for a web of open data. In K. Aberer et al. (eds.), *The Semantic Web, Lecture Notes in Computer Science* (vol. 4825, pp. 722-735). Springer, Berlin, Heidelberg.
- Baca, M., & Gill, M. (2015). Encoding multilingual knowledge systems in the digital age: The getty vocabularies. *Knowledge Organization*, 42(4), 232-243.
- Batini, C., Cappiello, C., Francalanci, C., & Maurino, A. (2009). Methodologies for data quality assessment and improvement. *ACM Computing Surveys (CSUR)*, 41(3), 16.
- Bizer, C., Lehmann, J., Kobilarov, G., Auer, S., Becker, C., Cyganiak, R., & Hellmann, S. (2009). DBpedia-A crystallization point for the web of data. *Web Semantics: Science, Services and Agents on the World Wide Web*, 7(3), 154-165.
- Bollacker, K., Evans, C., Paritosh, P., Sturge, T., & Taylor, J. (2008). *Freebase: A collaboratively created graph database for structuring human knowledge*. In Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data (pp. 1247-1250). ACM.
- Burgess, L. C. (2016). *Provenance in digital libraries: Source, context, value and trust*. In Building Trust in Information (pp. 81-91). Springer International Publishing.
- Carata, L., Akoush, S., Balakrishnan, N., Bytheway, T., Sohan, R., Seltzer, M., & Hopper, A. (2014). A primer on provenance. *Communications of the ACM*, 57(5), 52-60.
- Charles, V. (2016) *Linked data for Europeana cultural heritage: The Europeana approach*. Presentation given on April 28th in Paris at International Conference organised by ISSN IC: “Bibliographic metadata getting linked...”. Retrieved from <http://www.slideshare.net/ValentineCharles/linked-data-for-europeancultural-heritage-the-europeana-approach>.
- Charles, V., Manguinhas, H., Alexiev, V., Charles, V., & Dammers, M. (2015). *Wikidata, a Target for*

- Europeana's Semantic Strategy*. Glam-Wiki 2015. Retrieved from [https://nl.wikimedia.org/wiki/GLAM-WIKI\\_2015/Proposals/Wikidata,\\_a\\_target\\_for\\_Europeana%E2%80%99s\\_semantic\\_strategy%3F](https://nl.wikimedia.org/wiki/GLAM-WIKI_2015/Proposals/Wikidata,_a_target_for_Europeana%E2%80%99s_semantic_strategy%3F)
- Chuttur, M. Y. (2014). Investigating the effect of definitions and best practice guidelines on errors in Dublin Core metadata records. *Journal of Information Science*, 40(1), 28-37.
- De Sabbata, S., & Acheson, E. (2016, April). *Geographies of gazetteers in Great Britain*. In 24th GIS Research UK (GISRUK 2016) conference, University of Greenwich. Retrieved from <http://hdl.handle.net/2381/38182>
- Dextre Clarke, S. G. (2016). Origins and trajectory of the long thesaurus debate. *Knowledge Organization*, 43(3), 138-144.
- Emani, C. K., Cullot, N., & Nicolle, C. (2015). Understandable big data: A survey. *Computer Science Review*, 17, 70-81.
- Ermilov, I., & Pellegrini, T. (2015). *Data licensing on the cloud: Empirical insights and implications for linked data*. Paper presented at the Proceedings of the 11th International Conference on Semantic Systems (pp. 153-156). ACM.
- Erxleben, F., Günther, M., Krötzsch, M., Mendez, J., & Vrandečić, D. (2014). *Introducing Wikidata to the linked data web*. Paper presented at the International Semantic Web Conference (ISWC) (pp. 50-65). Springer International Publishing.
- Färber, M., Bartscherer, F., Menne, C., & Rettinger, A. (2016). Linked data quality of DBpedia, freebase, OpenCyc, Wikidata, and YAGO. *Semantic Web*. Retrieved from <http://www.semantic-web-journal.net/system/files/swj1465.pdf>
- Ford, H., & Graham, M. (in press). Provenance, power and place: Linked data and opaque digital geographies. *Environment and Planning D: Society and Space*.
- Fürber, C., & Hepp, M. (2013). *Using semantic web technologies for data quality management*. In Handbook of data quality (pp. 141-161). Springer, Berlin, Heidelberg.
- Godby, C. J. (2016). *Seeding the linked data cloud: The present and future of library identifiers*. Days of Knowledge organization, Oslo and Akershus University. Retrieved from [http://edu.hioa.no/korg2016/korg2016\\_godby.pdf](http://edu.hioa.no/korg2016/korg2016_godby.pdf)
- Goodwin, J., Dolbear, C., & Hart, G. (2008). Geographical linked data: The administrative geography of Great Britain on the semantic web. *Transactions in GIS*, 12(s1), 19-30.
- Hallo, M., Luján-Mora, S., Maté, A., & Trujillo, J. (2016). Current state of linked data in digital libraries. *Journal of Information Science*, 42, 117-127.
- Haslhofer, B., & Isaac, A. (2011). *data. europeana. eu: The europeana Linked Open Data pilot*. In International Conference on Dublin Core and Metadata Applications (pp. 94-104). Retrieved from <http://dcpapers.dublincore.org/pubs/article/view/3625/1851>
- Hoffart, J., Suchanek, F. M., Berberich, K., & Weikum, G. (2013). YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia. *Artificial Intelligence*, 194, 28-61.
- Hoffart, J., Suchanek, F. M., Berberich, K., Lewis-Kelham, E., De Melo, G., & Weikum, G. (2011, March). *YAGO2: exploring and querying world knowledge in time, space, context, and many languages*. Paper presented at the Proceedings of the 20th international conference companion on World Wide Web (pp. 229-232). ACM.

- Huang, A. W. C., & Chuang, T. R. (2014). *Relations for Reusing (R4R) in a Shared Context: An Exploration on Research Publications and Cultural Objects*. In Proceedings of the 4th International Workshop on Semantic Digital Archives (SDA)@ JCDL/TPDL (pp. 49-60). London, UK.
- Ismayilov, A., Kontokostas, D., Auer, S., Lehmann, J., & Hellmann, S. (2016). Wikidata through the eyes of DBpedia. *Semantic Web*. Retrieved from <http://www.semantic-web-journal.net/system/files/swj1462.pdf>
- Janssen, M., Charalabidis, Y., & Zuiderwijk, A. (2012). Benefits, adoption barriers and myths of open data and open government. *Information Systems Management*, 29(4), 258-268.
- Knoblock, C. A., & Szekely, P. A. (2015). Exploiting semantics for big data integration. *AI Magazine*, 36(1), 25-38.
- Lee, C.J., Huang, A.W.C., & Chuang, T.R. (2016, October). *A linked open data repository built with CKAN*. Paper presented at the CKANCon 2016, Madrid, Spain.
- Lee, C.J., Huang, A.W.C., & Chuang, T.R. (2017, March). *Metadata as Linked Data for Research Data Repositories*. Paper presented at the International Symposium on Grids and Clouds (ISGC) 2017. Retrieved from <http://m.odw.tw/u/odw/m/metadata-as-linked-data-for-research-data-repositories/>
- Lehmann, J., Isele, R., Jakob, M., Jentzsch, A., Kontokostas, D., Mendes, P. N., ... & Bizer, C. (2015). DBpedia—a large-scale, multilingual knowledge base extracted from Wikipedia. *Semantic Web*, 6(2), 167-195.
- Mahdisoltani, F., Biega, J., & Suchanek, F. (2015). *Yago3: A knowledge base from multilingual wikipedias*. In 7th Biennial Conference on Innovative Data Systems Research, CIDR Conference. Asilomar, CA.
- Marden, J., Li-Madeo, C., Whysel, N., & Edelstein, J. (2013). *Linked open data for cultural heritage: evolution of an information technology*. In Proceedings of the 31st ACM international conference on Design of communication (pp. 107-112). ACM.
- Meroño-Peñuela, A., Ashkpour, A., Van Erp, M., Mandemakers, K., Breure, L., Scharnhorst, A., ... & Van Harmelen, F. (2014). Semantic technologies for historical research: A survey. *Semantic Web*, 6(6), 539-564.
- Mitchell, E. T. (2016). The current state of linked data in libraries, archives, and museums. *Library Technology Reports*, 52(1), 5-13.
- Moura, T. H., & Davis Jr, C. A. (2014). *Integration of linked data sources for gazetteer expansion*. In Proceedings of the 8th Workshop on Geographic Information Retrieval. ACM.
- Omitola, T., Gibbins, N., & Shadbolt, N. (2010, December). *Provenance in linked data integration*. In Proceedings of the Workshop on Linked Data in the Future Internet at the Future Internet Assembly (LDFI-2010), Ghent, Belgium.
- Park, J. R., & Childress, E. (2009). Dublin Core metadata semantics: An analysis of the perspectives of information professionals. *Journal of Information Science*, 35(6), 727-739.
- Parr, C. S., Schulz, K. S., Hammock, J., Wilson, N., Leary, P., Rice, J., & Corrigan Jr, R. J. (2016). TraitBank: Practical semantics for organism attribute data. *Semantic Web*, 7(6), 577-588.
- Parr, C. S., Wilson, N., Leary, P., Schulz, K., Lans, K., Walley, L., ... & Holmes, J. (2014). The encyclopedia of Life v2: Providing global access to knowledge about life on earth. *Biodiversity Data Journal*, 2, e1079.
- Poole, A. H. (2016). The conceptual landscape of digital curation. *Journal of Documentation*, 72(5), 961-986.
- Schaible, J., Gottron, T., & Scherp, A. (2014). *Survey on common strategies of vocabulary reuse in linked open data modeling*. In European Semantic Web Conference (pp. 457-472). Springer International Publishing.
- Schmachtenberg, M., Bizer, C., & Paulheim, H. (2014). *Adoption of the linked data best practices in different*

- topical domains*. In International Semantic Web Conference (pp. 245-260). Springer International Publishing.
- Srinivasan, R., Becvar, K., Boast, R., & Enote, J. (2010). Diverse knowledges and contact zones within the digital museum. *Science, Technology, & Human Values*, 35(5), 735-768.
- Stadler, C., Lehmann, J., Höffner, K., & Auer, S. (2012). *Linkedgeodata: A core for a web of spatial open data*. *Semantic Web*, 3(4), 333-354.
- Stvilia, B., Gasser, L., Twidale, M. B., & Smith, L. C. (2007). A framework for information quality assessment. *Journal of the American Society for Information Science and Technology*, 58(12), 1720-1733.
- Suchanek, F. M., Kasneci, G., & Weikum, G. (2007). *Yago: A core of semantic knowledge*. In Proceedings of the 16th International Conference on World Wide Web (pp. 697-706). New York, NY: ACM.
- Suchanek, F. M., Kasneci, G., & Weikum, G. (2008). Yago: A large ontology from Wikipedia and Wordnet. *Web Semantics: Science, Services and Agents on the World Wide Web*, 6(3), 203-217.
- Tani, A., Candela, L., & Castelli, D. (2013). Dealing with metadata quality: The legacy of digital library efforts. *Information Processing & Management*, 49(6), 1194-1205.
- Van Hooland, S., & Verborgh, R. (2014). *Linked Data for Libraries, Archives and Museums: How to clean, link and publish your metadata*. London: Facet Publishing.
- Vandenbusche, P. Y., Atemezing, G. A., Poveda-Villalón, M., & Vatan, B. (2015). Linked Open Vocabularies (LOV): A gateway to reusable semantic vocabularies on the Web. *Semantic Web*, 8(3), 437-452.
- Voß, J. (2016, September). *Classification of Knowledge Organization Systems with Wikidata*. In Proceedings of the 15th European Networked Knowledge Organization Systems Workshop (NKOS 2016), Hannover, Germany.
- Vrandečić, D., & Krötzsch, M. (2014). Wikidata: A free collaborative knowledgebase. *Communications of the ACM*, 57(10), 78-85.
- Yasser, C. M. (2011). An analysis of problems in metadata records. *Journal of Library Metadata*, 11(2), 51-62.
- Yus, R., & Pappachan, P. (2015). *Are Apps Going Semantic? A Systematic Review of Semantic Mobile Applications*. Paper presented at the 1st International Workshop on Mobile Deployment of Semantic Technologies (MoDeST 2015), co-located with the 14th International Semantic Web Conference (ISWC 2015), Bethlehem, PA.
- Zaveri, A., Rula, A., Maurino, A., Pietrobon, R., Lehmann, J., & Auer, S. (2016). Quality assessment for linked data: A survey. *Semantic Web*, 7(1), 63-93.
- Zhu, R., Hu, Y., Janowicz, K., & McKenzie, G. (2016). Spatial signatures for geographic feature types: Examining gazetteer ontologies using spatial statistics. *Transactions in GIS*, 20(3), 333-355.

# ***Reuse of Structured Data: Semantics, Linkage, and Realization***

## **Andrea Wei-Ching Huang**

Project Manager (Research)

Institute of Information Science, Academia Sinica, Taiwan

E-mail: andrea~~h~~g@iis.sinica.edu.tw

## **Cheng-Jen Lee**

Research Assistant

Institute of Information Science, Academia Sinica, Taiwan

E-mail: cjlee@iis.sinica.edu.tw

## **Tyng-Ruey Chuang**

Associate Research Fellow

Institute of Information Science, Academia Sinica, Taiwan

E-mail: trc@iis.sinica.edu.tw

Keywords: CKAN; Data Provenance; Data Quality; Knowledge Base; Linked Open Data (LOD); Ontology; Semantic Representation

---

### **【Abstract】**

In order to increase the reuse value of existing datasets, it is now becoming a general practice to add semantic links among the records in a dataset, and to link these records to external resources. The enriched datasets are published on the web for both human and machine to consume and re-purpose. In this paper, we make use of publicly available structured records from a digital archive catalogue, and we demonstrate a principled approach to converting the records into semantically rich and interlinked resources for all to reuse. While exploring the various issues involved in the process of reusing and re-purposing existing datasets, we review the recent progress in the field of Linked Open Data (LOD), and examine twelve well-known knowledge bases built with a Linked Data approach. We also discuss the general issues of data quality, metadata vocabularies, and data provenance. The concrete outcome of this research work is the following: (1) a website data.odw.tw that hosts more than 840,000 semantically enriched catalogue records across multiple subject areas, (2) a lightweight ontology voc4odw for describing data reuse and provenance, among others, and (3) a set of open source

software tools available to all to perform the kind of data conversion and enrichment we did in this research. We have used and extended CKAN (The Comprehensive Knowledge Archive Network) as a platform to host and publish Linked Data. Our extensions to CKAN is open sourced as well. As the records we drawn from the originally catalogue are released under the Creative Commons licenses, the semantically enriched resources we now re-publish on the Web are free for all to reuse as well.

## **【Long Abstract】**

### **Introduction**

In order to enhance the reuse value of existing datasets, it is now becoming a general practice to add semantic links among the records in a dataset, and to link these records to external resources. The enriched datasets are published on the Web for both the human and the machine to consume and re-purpose. In the paper, we make use of publicly available structured records from a digital archive catalogue, and we demonstrate a principled approach to converting the records into semantically rich and interlinked resources for all to reuse. While exploring the various issues involved in the process of reusing and re-purposing existing datasets, we review the recent progress in the field of Linked Open Data (LOD), and examine twelve well-known knowledge bases built with a Linked Data approach. We also discuss the general issues of data quality, metadata vocabularies, and data provenance.

The concrete outcome of this research work is the following: (1) a website that hosts more than 840,000 semantically enriched catalogue records across multiple subject areas, (2) a lightweight ontology voc4odw for describing data reuse and provenance, among others, and (3) a set of open source software tools available to all to perform the kind of data conversion and enrichment we did in this research. We have used and extended CKAN (The Comprehensive Knowledge Archive Network) as a platform to host and publish Linked Data. Our extensions to CKAN is open sourced as well. As the records we have drawn from the originally catalogue are released under the Creative Commons licenses, the semantically enriched resources we now re-publish on the Web are free for all to reuse as well.

### **Review of Twelve Knowledge Bases**

We begin by first examine twelve knowledge bases built with a Linked Data approach. Five of them are built by domain knowledge experts (OpenCyc, Getty Art & Architecture Thesaurus, Getty Thesaurus of Geographic Names, and Ordnance Survey), six of them are collaborative databases (Freebase, YAGO, DBpedia, Wikidata, LinkedGeoData, GeoNames), and the last one is about ecological observations based on expert and community collaborations (Encyclopedia of Life). We further compare datasets

about geospatial entities with controlled vocabularies: Getty TGN, Open Names (Ordnance Survey), DBpediaPlace, LinkedGeoData, and GeoNames.

To make good reuse of structured data, ones need to first deal with the problem of data quality. Currently there exist different evaluation criteria, with various techniques for measuring the quality of information, data, metadata, and Linked Data. We review four papers on data quality and systematically compare their evaluation criteria. Moreover, data provenance --- contextual metadata about the source and use of data --- has proven to be fundamental for assessing authenticity, enabling trust, and allowing reproducibility. Thus, we examine key mechanisms of data provenance before we move forward to discussing LOD applications.

## Practices

We then make use of structured records from a digital archive catalogue, and convert the records into semantically rich and interlinked resources on the Web. This is realized as a unified Linked Data catalogue to several digital archive collections. Our work results in a LOD catalogue available to the public at the website <<http://data.odw.tw>>. The following five parts are involved in realizing this website. A catalogue record, about a species of *Pleione Formosana*, is used throughout in the paper as an example to demonstrate the way we model, convert, and represent the semantics of a structured record.

Part 1: Exploring data reuse relations in a shared context -- We review our previous research about the Relation for Reuse Ontology (R4R). In particular, we provide mechanisms for reusing article, data, and code with some flexibility of encoding provenance and license information.

Part 2: Comparing two different data conversion approaches to providing LOD for an archive catalogue -- We show two different scenarios: (1) The LOD catalogue is converted directly from a relational database, and (2) the LOD catalogue is generated from a series of format conversions --- from XML to CSV, and then to RDF.

Part 3: Data profiling, cleaning and mapping -- We demonstrate format conversion processes, and we discuss the pros and cons of various ways in handling broken links in source datasets. In addition, we mapped and linked catalogue records to three external knowledge bases: GeoNames, Wikidata, and Encyclopedia of Life.

Part 4: Using CKAN (The Comprehensive Knowledge Archive Network) as a Linked Data platform -- We briefly introduce CKAN, an open source web-based data portal software package for curating and publishing datasets. CKAN provides data preview, search, and discovery, especially with regard to geospatial datasets. We built several extensions to CKAN in order to deposit, publish, browse, and

search Linked Data. Various Linked Data representations of a catalogue record --- Turtle, RDF/XML, and JSON-LD --- can all be downloaded and reused.

Part 5: Designing ontologies for data representation and reuse -- We design an ontology voc4odw which includes the following 3 modules:

- (1) The Core Model. It is comprise of a data model and a conceptual model. The data model represents key data structure and relation. It is a framework to illustrate data source, derivation, and provenance. The conceptual model incorporates SKOS Simple Knowledge Organization System; it also connects to key event concepts. The conceptual model allows for data contextualization using common and domain knowledge vocabularies.
- (2) The Curation Model. It is responsible for disclosing the identification, classification, and publication of structured records at a curation platform, such as the classification of themes, the assignment of data identifiers, and the publication of datasets.
- (3) A vocabulary voaf:Vocabulary. It is defined as "A vocabulary used in the Linked Data cloud", from the Vocabulary of a Friend <<http://purl.org/vocommons/voaf>>. This module is to relate the Core Model to external common vocabularies. Some hierarchy relations between different external vocabularies can be traced with this vocabulary.

**【Romanization of Chinese references is offered in the paper.】**