

# 利用關鍵字分析建立自閉症論文知識庫

張峻彬<sup>1</sup>、謝愛家<sup>2</sup>、盛夢徽<sup>3</sup>

嘉南藥理大學{<sup>1,3</sup>資訊管理系<sup>2</sup>資訊多媒體應用系}

{cpchang, qajhsieh, dream}@mail.cnu.edu.tw

## 摘要

本研究中利用 PDF 論文檔做為知識的主要來源，建立一個以自閉症主題的線上知識庫管理系統。系統主要由 PDF 文件上傳、文件驗證、文件內容解析、關鍵字關聯分析和主題知識查詢等五部分所組成。在系統特性上，使用者與管理者均不需具備特定領域的專有知識，論文內容的收集、分析、組織均是自動進行，這個過程主要是利用論文作者自定關鍵字來進行即時性的關聯分析與整理，而在知識的搜尋與應用上，系統中採用主題式的導引，使用者可在不受本身專業知識的限制下，自主性地定位，找出其所需的文獻與相關內容。

**關鍵詞：**關鍵字、知識庫、主題地圖

## Abstract

In this Study, an online Knowledge Base Management System is established by using PDF documents as the source of knowledge. System consists of five parts, include: PDF file upload, validation, content analysis, keyword correlation analysis and knowledge query. In this system, users and manager do not need to have any domain knowledge in specific areas, thesis collection, because the analysis and organizations are automatically. In the search of knowledge, the system uses the Topic Map guide. Users can locate their required documents and relevant content by autonomous positioning function.

**Keywords:** Keyword, Knowledge Base, Topic Map

## 1. 前言

由於網路應用的發達，各類資訊的取得相當的容易，也正因为資訊的取得太容易，可能會陷入“過多的資訊等於沒有資訊”的困境之中，由於網際網路並未具有組織高重複性資料的特性，使得資料擷取產生困難，亦有資訊負載的問題[1]。因此在資訊超載的環境中，如何針對特定主題快速地收集、分析、組織、應用資訊變成一個很重要的課題。為了應付前述的狀況，建置一個系統化地保存、管理與分享知識的知識庫是必要的，此知識庫可以協助一般使用者取得知識、了解專門領域的概念。在此知識庫(Knowledge Base)指的是針對某一特定領域的經驗與知識，採用某種知識表示方式在資料庫中儲

存、組織、管理和使用的知識集合。

可攜式文件格式(Portable Document Format, PDF)是由 Adobe 在 1993 年所公開出來的格式，目前是最被廣泛應用的電子文件格式。使用 PDF 有兩項主要的優點，一是它跨越各種作業系統平台(WINDOWS、MAC、UNIX...)均有閱讀器可讀取，另一個是它可以保存原始文件樣貌在各種作業系統平台上閱讀或列印，也正因為這兩項優點，許多的學術論文都採用 PDF 做為傳送與保存的格式。

以 Airiti Library 華藝線上圖書館而言，他所提供下載的論文均是使用 PDF 的的檔案格式。若以自閉症為主題(關鍵字)進行搜尋，可以找到 305 篇相關的學術論文(2014/8/28, 圖 1)，這雖然不是相當巨大的文件量，但若是使用者想在自閉症這一主題下找到自己有興趣的資訊是相當不容易的，因為這類電子資源搜尋引擎只能針對論文題目、關鍵字及摘要內容進行搜尋，但是現實上論文作者不可能將論文的所有重要內容都放在論文題目、關鍵字或摘要中，所以不在搜尋結果中的論文並不代表論文中沒有使用者需要的訊息。而且就算是使用者將 305 篇相關的學術論文都下載，他會需要花大量的時間才有辦法(或沒有辦法)去進行整理，由於這些論文主要是以個別文件檔案作為處理單位，面對大量的檔案資料，如何運用有效的原則來組織這些論文，以便有效的管理以從中取的資訊與知識，將是一個很重要的工作。



圖 1 以自閉症為關鍵字進行搜尋





圖 4 匯入與刪除等管理功能

### 3. 文件驗證：

在使用者匯入論文後，檢驗此論文是否符合本知識主題(主題關鍵字至少出現 10 次以上)，自動抓取論文標題與關鍵字等訊息，並於畫面上顯示論文的相關內容以供使用者確認。在使用者確定要匯入此論文後，檢驗知識庫中是否已存在此論文。



圖 5 確認與檢驗論文

### 4. 文件內容解析：

將本文中若有原本不存在知識庫中的關鍵字加入知識庫中，對論文進行基本資訊分析，統計原生關鍵字與延伸關鍵字於本文中出現的次數。分析統計知識庫中原生關鍵字與延伸關鍵字於知識庫中出現的總次數。統計原生關鍵字與延伸關鍵字於知識庫其他論文中出現的次數。

### 5. 主題知識查詢：

主要提供關鍵字於知識庫中的相關訊息；依關鍵字查詢相關論文；單一論文資訊展示，包括本文、關鍵字類型與出現資訊、與此論文具關聯性的其他論文顯示；單一論文相關查詢、關鍵字查詢、關聯性查詢等相關功能。

## 4. 主題知識查詢

知識庫管理系統在經過 PDF 文件上傳、文件驗證、文件內容解析與關鍵字關聯分析就可提供使用

者進行主題知識查詢。進入後有區分為兩種搜尋方式：關鍵字關聯搜尋(Keyword Relational Search)，論文排序搜尋(Article Rank Search)。有別於一般查詢系統會讓使用者輸入想查詢的字串，再根據此字串去找出相關的訊息，在這個自閉症主題關聯搜尋系統中，完全不需要使用者輸入任何資訊，而是由使用者去選擇他有興趣的主題，由這個主題透過關聯分析的導引去找出使用者想看的論文。

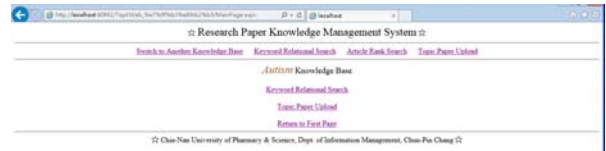


圖 6 主功能畫面

以關鍵字關聯搜尋來說，由頁面(圖 7)上可以了解目前在自閉症知識庫中有幾篇論文(目前有 45 篇跟自閉症有關的中文論文)，這些論文總共有幾個原生與延伸關鍵字(目前有 111 個)。在學術論文中，作者都會對他自己的文章定義五個左右，最能代表這篇論文的關鍵字，像這種原作者定義的關鍵字在系統中稱做"原生"關鍵字(Original Keyword)，而在 A 論文中作者定義到的關鍵字，若是出現在 B 論文中，但是並沒有被 B 論文作者定義為關鍵字，像這樣的關鍵字在系統中稱為 B 論文的"延伸"關鍵字(Extended Keyword)。



圖 7 關鍵字關聯搜尋畫面

在關鍵字關聯搜尋中會依關鍵字出現於論文篇數的多寡排序，列出關鍵字供使用者選擇。在畫面的右側有表列出各關鍵字於知識庫中出現的總次數，以及出現的論文篇數。關聯搜尋功能中，一開始"列出包含選定關鍵字的論文"是沒辦法按的，一定要選擇 1 至 3 關鍵字，這個按鈕功能才可以使



圖 8 選擇一個關鍵字(問題行為)啟用搜尋功能

以選擇一個關鍵字(問題行為)為例，可以知道知識庫有幾篇論文中出現了這個關鍵字，關鍵字出現了幾次，以及在這篇論文中這個關鍵字屬於何種類型。



圖 9 選擇一個關鍵字的搜尋結果畫面

以兩個關鍵字為例，可由知識庫中找出包含這兩個特定關鍵字的相關論文。



圖 9 選擇二個關鍵字的搜尋結果畫面

以三個關鍵字為例，可由知識庫中找出包含這三個特定關鍵字的相關論文。

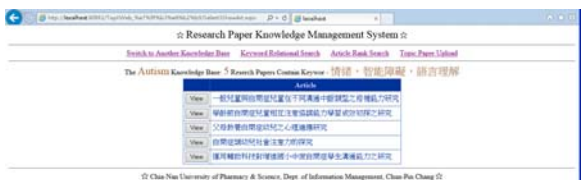


圖 10 選擇三個關鍵字的搜尋結果畫面

不管是用一個、兩個或三個關鍵字去搜尋，都可以文章列表前的 View 按鈕去看特定文章的內文。



圖 11 特定文章內文與相關資訊

若想看看與本論文關聯性比較高的論文有那一些，可利用觀看本文關聯列表(圖 11 中左邊的列表)，此列表主要是針對本篇論文中的原生與延伸關鍵字依出現的次數(Sort by Count)、知識庫中出現的總次數(Sort by TotalCount)或出現的論文篇數來排序(Sort by PaperCount)，以前三高的關鍵字來找出關聯度較高的文章。Sort by Count, Sort by TotalCount, Sort by PaperCount 三個功能可以直接點擊切換呈現的方式。

若是對圖 11 中"Top 3 Keywords in This Research Paper"標題底下所指具有"相互注意協調能力 & 自閉症 & 自閉症兒童"三個關鍵字的相關文章有興趣，可點擊前面+號將它打開，出現文章項目，點擊文章名稱，就可觀看其內容與關聯文章(Sort by Count)。

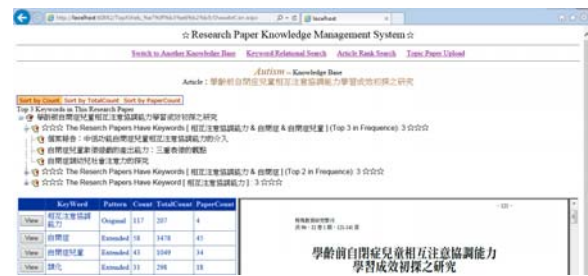


圖 12 具有相似關鍵字的高關聯論文(Sort by Count)

依不同的排序方式，可能會找出不同的關聯文章。



圖 13 具有相似關鍵字的高關聯論文(Sort by TotalCount)

