

The 4th International Conference on
Literature and Information Technology

第四屆文學與資訊科技國際研討會

淺談文章內容的標誌 與意義的處理

謝清俊

銘傳大學講座教授

2008年11月13日

大 綱

- ❁ 從數位化看文章生態的變遷
- ❁ 後設資料 (metadata) 的省思
- ❁ 內容的標誌
- ❁ 內容標誌舉例
- ❁ 結語

從數位化看文章生態的變遷

從數位化看文章生態的變遷

- ❁ 數位化改變了溝通的生態，擔任溝通中介的文字紀錄或文章，其生態也風行草偃。如：
 - * 網際網路上「輕薄短小」的文章。
 - * 迎合青少年口味「圖多字少，膚淺花俏」的圖文夾雜。
 - * 有關閱讀習慣和認知行為變遷的研究報告。
 - * 電腦中數位化的文章。
 - ❖ 引起文章定義（界定）和範疇的問題。

文章定義與範疇的問題

❁ 文章經數位化存在電腦中時，只存文章的信息不夠，必需把一些有關背景的情境資料也存起來，並與文章作適當的连接。

❁ 互為文本(Inter-textuality) 理論

● *Julia Kristeva*

❁ 文章內容與外界的聯繫

❖ *hyperlinks*

內容之外化

- ❖ 標點、句讀、標題、章節段落…
- ❖ 位置、字體、色彩、加網加邊、美工加工…
- ❖ 標誌 (markup) 與內容標誌 (content markup)
 - ❖ 通用結構 (DTD)
 - 版面的、結構的、語文的、內容上的……
 - ❖ 標籤集 (tag set)
- ❖ 後設語言 (meta-language) 和後設資料 (metadata)

情境 (context)

許多人以為情境只是「上下文」，其實「上下文」是context 在語言修辭「情境」下的意義。

❁ 情境對文章而言，可泛指文章作成時所有相關的背景，包括：

❁ 作者相關的背景

❖ 如：作者生平、成書時間、著作時的身心狀態……

❁ 時代相關的背景

❖ 如：政治背景、經濟背景、軍事背景……

➢ 是承平還是戰亂、天災人禍、社會重大事件……

❁ 文化相關的背景……

❁ 與其他文章相關的背景

情境 (context)

- ❁ 一旦作品完成，情境信息即已固定，且恆久不變。
 - ✱ 此所以文物為文化之記錄。
- ❁ 意義是依情境而定的。情境既已固定，則作者創作的原意亦隨之固定。
 - ✱ 但閱聽者可依閱聽的情境作詮釋。
 - ❖ 『作者已死』 ● 羅蘭·巴特 (釋義學)
 - ❖ 詮釋是一種創新。
- ❁ 若無法描述情境，則無法真正處理文章的意義。

情境的處理

- ❖ 因為情境信息必需不分國家、種族，甚至於不分電腦機種都要能夠處理，所以需要一種電腦會處理的通用人工語言 (artificial language)，也就是後設語言，來描述。
- ❖ 後設語言不僅僅可以描述情境信息，文章內容的注疏、註釋，以及文章之間彼此的參照，甚至於文章內容與實物之間的聯繫關係等，也都可以用後設語言描述。

數位化文章的表達

❁ 文章的結構在電腦中產生了根本的改變：變成以自然語言和後設語言相輔表達的雙重結構：

✱ 以自然語言寫文章本身

✱ 以後設語言描述數位化的文章與外界的各種關係。

數位化文章信息之表達

數化之文章		表現系統
文章本身		自然語言
文章與外界的關係	情境描述(<i>metadata</i>)	後設語言
	參照聯繫(<i>hyperlinks</i>)	
	內容詮注(<i>content markup</i>)	

後設資料的省思

後設資料的認知

❁ 有人引據國外的文章，說後設資料就是「資料的資料」。

❁ 有了這樣的說法，許多人便認定：

「除了文物數位化的本身之外，所有其他的資料都屬後設資料」。

❁ 其實，這樣的認知是有問題的。

後設資料的認知

- ❖ 說後設資料是「資料的資料」，只是為了闡明後設資料這個概念的性質，並不是將後設資料定義為「資料的資料」。
- ❖ 因為，後設資料固然是「資料的資料」，可是並不是所有的「資料的資料」都是後設資料。
- ❖ 將後設資料界定為「資料的資料」這種認知，與「不吃豬肉的都是回教徒」犯了同樣的邏輯錯誤。

後設資料的範疇

- ❁ 現行的任何後設資料，其表達的方式、訂定的規格，以及標籤 (tag) 或欄位 (field) 的選擇和數目... 等都限制了後設資料的範疇。
- ❁ 這很明顯表示：
 - * 不是所有的「資料的資料」都是後設資料。要明白數位化的後設資料，不能把資料二分為資料和「資料的資料」這樣籠統的概念去理解。

時下後設資料的性質

❁ 為了某類文物而訂定

✱ 比方說，書目資料是一般書籍的後設資料，玉器、青銅器、畫作、雕刻……等都有各自的後設資料。

❁ 既然是描述「某類」文物的資料，那麼就有它的特徵和它的侷限。

✱ 它適合敘述某類文物的共同現象（共相）

✱ 無法顧及個別現象（別相）

❖ 後設資料充其量只能摘錄文本的一部份，而無法深入觸及文本的內容。

時下後設資料的性質

- ✿ 一般而言後設資料敘述的多屬事實、屬性這類較客觀可考的資料(共相)，不涉及文本內容的理解、感受、比較、批評，以及詮釋等(別相)。所以：
 - ✿ 後設資料是可以由具技術專業人士查訪、考證
 - ✿ 但是，它不可以作詮釋。
 - ❖ 比方說，我們可以考證《紅樓夢》的作者是誰，卻不能詮釋《紅樓夢》的作者是誰。

時下後設資料的性質

- ❁ 後設資料的內容是依應用的目的而異，一件數位文物的後設資料可以有許多種。例如新聞稿：
 - * 對記者而言，有一種後設資料；
 - * 對報社來說，同一則新聞有編輯用的、管理用的，甚至於是與其他通訊社交流用的各種後設資料。
 - * 這些後設資料之間，會有些重複，但也有獨特之處
- ❁ 所以，應用時會要求後設資料的獨特者能彼此互通，重複者需彼此一致。

時下後設資料的性質

- ❁ 人世間的事情常有變化，所以後設資料不會是固定不變的，它會與時遷移，需要花很多力氣更新、維護和保養。
- ❁ 後設資料既然如此複雜，就不是電腦常用的欄位結構可以處理的。所以，描述後設資料現在都用後設語言(meta-language)，如：HTML、XML。
- ❁ 只有語言才有能力描述後設資料的種種規格和後設資料之間的相容關係，以符合應用的需求。

内容標誌

內容標誌 (content markup)

- ❁ 內容標誌要照顧的正是後設資料無法觸及的一關於文物個別內容描述的這一部份。
- ❁ 以文章而言，對文本內容的理解、解釋、感受、比較、批評、詮釋等，正是內容標誌的主要工作。
 - ✧ 這些工作觸及人文、歷史、社會、美學、哲學…等學門的核心問題，需要真正了解內容的專業人士為之。

內容標誌

- ❁ 內容標誌，無論作理解、解釋、感受、比較、批評或詮釋，均觸及一個人文方面最根本的問題—意義(meaning)和了解(understanding)。
- ❁ 這是認知科學、語言學、記號學等近幾年來致力研究的重點，也是電腦迄今未能有效處理的痛處。
- ❁ 內容標誌正是為了解決這個困局而設

電腦標誌與傳統漢語文獻標誌

❁ 目前的電腦標誌多對文章的外形或形式作標誌：

❖ 通用結構 (DTD)

➤ 版面的、結構的、語文辭彙的……

➤ 標籤集 (tag set)

❁ 傳統漢語文獻的標誌則側重於文章內容的、意義上的標誌。

例：標點符號

- ❁ 標點符號改變了文章內容的表現方式：
 - ❁ 無標點符號的文章內容較為隱晦(implicit)—需經分析、理解的過程才能窺見原意。有了標點符號，則內容較外顯(explicit)，諸如：
 - ❖ 私名號的使用已明顯的標出姓名或機構名稱，減少了斷詞的工作
 - ❖ 句點、逗點、分號等則已將斷句標明。
- ❁ 標點符號將部份文章內容由隱晦轉為外顯，這就是一個「意義表達」的例子。

例

『民可使由之不可使知之』

『民可，使由之；不可，使知之』

『民，可使由之，不可使知之』



「閒雜人等不得在此小便」

「閒雜人等，不得在此小便」

「閒雜人，等不得，在此小便」



例：句讀

❁ 古文雖然不用現代的標點符號，然而有另一套常用的標誌系統：句讀。

❁ 句讀主要用途是作文章內容的標誌：

❖ 標明文中之美辭、佳句、警句，或文中之不佳之處等。

❖ 對詩詞韻文，也有用於標示韻腳和朗誦時的間歇者。

❁ 標點符號或句讀這類的標誌，都是設計來幫助讀者理解文章內容的；它也幫做標誌人，把他們對文章的理解用標誌記錄下來。

意義處理的問題

- ❁ 近年來，計算語言學和人工智能均致力於處理意義的研究，也取得一些成果。例如，詞網(word net)、主題圖(topic map)、知識本體結構(ontology)等。
 - * 它們將詞彙間的關係在電腦中作了適當的表達，並構成資料庫和研發為數位工具。
 - * 詞彙間的關係是語意中的一種，將它數位化，對意義的處理是有助益。
 - ❖ 可是助益有限，原因是囿於形式和內容(意義)是一對一前提，因此，並無能力處理意義的癥結—多義問題。

意義處理的問題

- ❁ 多義問題，簡單說，就是當一種形式可能對應到好幾種意義時，如何作正確選擇的問題。此即「義隨境轉」語意隨情境而轉移的現象。
- ❁ 例如，作數目字時，「十、拾」通用，可是情境變為「路不拾遺」時，就不可以作「路不十遺」。
- ❁ 人面對多義或義隨境轉問題並無太大難色，所有的自然語言都有濃厚的義隨境轉色彩，因為人多半了解情境，對「意義」會作適當的「了解」。所以，電腦處理意義問題的先決條件，是要會表達情境。可是目前學界在這方面的努力，還沒有顯著的成績。

內容標誌之例 《心經》

雛型軟體與圖一和圖二的說明

- ✿ 本文介紹的軟體是中央研究院資訊科學研究所文獻處理實驗室所發展的一個雛型（proto-type）（註2）。它能顯現兩種視窗：一個是文件的視窗，稱為「文件夾」；另一個是文件內容結構的視窗，稱為「知識結構夾」。文件夾有兩個（名為文件夾一與文件夾二，如圖一右側上下的兩個視窗），每個最多可呈現四篇文件。所以文件夾總共可同時呈現八篇文件備用。
- ✿ 知識結構夾也有兩個（如圖二左側上下的兩個視窗），每個亦可最多呈現四種對文件夾中文件內容的結構描述。所以文件夾總共可同時呈現八篇文件的內容結構備用。
- ✿ 這樣的工具是可以協助使用者處理文件意義的。比方說，在文件夾中選擇的是同一文章的不同版本，在知識結構夾中選的是不同作者對此文章的解釋或詮釋。像這樣的安排，就可以處理不同作者對此文章各種解釋或詮釋的比對，同時也可以參照各版本做內容異同的對應。
- ✿ 在我們舉的例子裡，文章選的是《心經》；另分別選周止庵先生和印順導師對《心經》的註解，作為心經的知識結構（請參考圖二）。

圖二

知識結構夾一

1. 心經科文(改自印順心經講記)

- 心經科文(印順)
 - 懸論
 - 釋經題
 - 釋譯題
 - 正釋
 - 序分
 - 正宗分
 - 標宗
 - 顯義
 - 正為利根示常道
 - 法說般若體
 - 喻讚般若德
 - 曲為鈍根說方便
 - 流通分

文件夾一

1. 心經(玄奘梵文漢譯) 2. 心經(藏文英譯)

般若波羅蜜多心經

唐三藏法師玄奘譯

觀自在菩薩。行深般若波羅蜜多時。照見五蘊皆空。度一切苦厄。舍利子。色不異空。空不異色。色即是空。空即是色。受想行識亦復如是。舍利子。是諸法空相。不生不滅。不增不減。不垢不淨。不異不離。無受。無想。無行。無識。無眼。耳鼻舌身。無意。無意識界。無無明。亦無集。滅。道。乃至無老死。亦無得。以無所礙。無罣礙。無恐怖。依羅蜜多故。心無罣礙。究竟涅槃。三世諸佛。遠離顛倒夢想。究竟涅槃。三世諸佛。依般若

知識結構夾二

1. 心經科文(改自周止庵心經註)

- 心經科文(周止庵)
 - 總釋名題
 - 正釋經文
 - 序分
 - 正宗分
 - 顯說般若
 - 因人顯法
 - 正示法空
 - 顯彰妙果
 - 結讚功能
 - 密說般若
 - 流通分

文件夾二

1. 心經(鳩摩羅什梵文漢譯)

觀世音菩薩。行深般若波羅蜜時。照見五陰空。度一切苦厄。舍利弗。色空故。無惱壞相。受空故。無受相。想空故。無知相。行空故。無作相。識空故。無覺相。何以故。舍利弗。是諸法非色異空。非空異色。色即是空。空即是色。受。想。行。識。亦復如是。舍利弗。是諸法空相。不生不滅。不垢不淨。不增不減。不異不離。無眼耳鼻舌身。無意。無意識界。無無明。亦無集。滅。道。乃至無老死。亦無得。以無所礙。無罣礙。無恐怖。依羅蜜多故。心無罣礙。究竟涅槃。三世諸佛。遠離顛倒夢想。究竟涅槃。三世諸佛。依般若

圖四說明

什無覺五空質，時觀是印
羅，無「異性義。同的然，印
摩故，文不聯意。會順此顯較，此
鳩空故前色關的也。擊庵，這相
覺受空對「相文，下點止」與此
察，識為字的經，兩，點若。與
以相，視文義會。目，如般的。
可壞相以後意體。項，譬以說行
們惱作可其是來。文。譬以說行
我無無，釋都，中。科。釋。顯度若
中故，一為，係的註釋。屬的修
文空空的作係關夾中的的段主「
譯色行多可關係應中。的經屬的修
的「相，這也些對構相相為」它
同了相。這，這些結他現認？說
不多知串釋。這。知識其出庵麼度
什大致無一解樑的知和會止什角
羅大，這的橋示。擊文，周是的
摩，故「義的顯點經」。容行的
鳩多空故意「色」從若的行四內修
與較想以「色」以，關若圖文重
奘文，何空異可的，相般如經偏
玄譯相，皆不者樣現修點。「經則
在的受相蘊空讀同呈「修點。從順

圖五說明

❁ 再舉一例，讓我們看看周止庵說的「結讚功能」，其中有「無等等咒」一語。如果我們弄不清楚它是什麼意思，那麼，我們可以把藏文英譯的《心經》叫出來（如圖五右上角視窗），看看對等的英譯文字「the mantra equalizes whatever is unequal」，「無等等咒」一語的意義也就豁然開朗了。

圖五

1. 心經科文(改自印順心經講記)

- 正宗分
 - 標宗
 - 顯義
 - 正為利根示常道
 - 法說般若體
 - 修般若行
 - 廣觀蘊空
 - 融相即性觀(加行)
 - 泯相證性觀(正證)
 - 略觀處界等空
 - 結顯空義
 - 得般若果
 - 喻讚般若德
- 曲為鈍根說方便
- 流通分

1. 心經(玄奘梵文漢譯) 2. 心經(藏文英譯)

and no non-attainment. Therefore, Shariputra, since bodhisattvas have nothing to attain they abide in reliance upon Prajnaparamita. Without obscurations of mind, they have no fear. Completely transcending false views They go to the ultimate of nirvana. All the buddhas of the three times by relying on the Prajnaparamita awaken completely to the perfect, unsurpassable enlightenment. Therefore, the mantra of Prajnaparamita is the mantra of great awareness; it is the unsurpassed mantra, the mantra that equalizes whatever is unequal, and the mantra that totally pacifies all suffering. Since it does not deceive, it should be known as truth. The Prajnaparamita mantra is uttered

知識結構夾二

1. 心經科文(改自周止庵心經註注)

- 心經科文(周止庵)
 - 總釋名題
 - 正釋經文
 - 序分
 - 正宗分
 - 顯說般若
 - 因人顯法
 - 正示法空
 - 明蘊空
 - 顯空德
 - 顯彰妙果
 - 結讚功能
 - 密說般若
 - 流通分

文件夾二

1. 心經(鳩摩羅什梵文漢譯)

受空故。無受相。想空故。無知相。行空故。無作相。識空故。無覺相。何以故。舍利弗。非色異空。非空異色。色即是空。空即是色。受·想·行·識。亦復如是。舍利弗。是諸法空相。不生。不滅。不垢。不淨。不增。不減。是空法非過去。非未來。非現在。是故空中無色。無受·想·行·識。無眼·耳·鼻·舌·身·意。無色·聲·香·味·觸·法。無眼界。乃至無意識界。無無明。亦無無明盡。乃至無老死。亦無老死盡。無苦·集·滅·道。無智。亦無得。以無所得故。菩薩依般若波羅蜜故。心無罣礙。無罣礙故。無有恐怖。遠離一切顛倒夢想苦惱。究竟涅槃。三世諸佛。依般若波羅蜜故。得阿耨多羅三藐三菩提。故知般若波羅蜜。是大明咒。是無上明咒。是無等等明咒。能除一切苦。真實不虛。故說般若波

圖六說明

❁ 最後的例子，讓我們了解一下《心經》的咒語。我們點擊《心經》最後的咒語那一段，就會看到說它是為了鈍根說方便，印順《心經》的方便法門，（「曲為鈍根說方便」，請見圖六。

❁ 再者，由於譯文的文字古老，用現在的發音唸不準，那怎麼辦呢？我們可以在右邊把窗所是「揭是確」的《心經》「揭」字發還至少知道「揭」的發音。我們發音是「揭」？我們發音是「揭」？我們發音是「揭」？

圖六

1. 心經科文(改自印順心經講記)

- 正宗分
 - 標宗
 - 顯義
 - 正為利根示常道
 - 法說般若體
 - 修般若行
 - 廣觀蘊空
 - 融相即性觀(加行)
 - 泯相證性觀(正證)
 - 略觀處界等空
 - 結顯空義
 - 得般若果
 - 喻讚般若德
 - 曲為鈍根說方便
 - 流通分

1. 心經(玄奘梵文漢譯) 2. 心經(藏文英譯)

無受·想·行·識。無眼·耳·鼻·舌·身
 ·意·無色·聲·香·味·觸·法。無眼界。乃至無意識界。無無明。亦無無盡。乃至無智
 老死。亦無老死盡。無苦·集·滅·道。無智波
 。亦無得。以無所得故。菩提薩埵。依般若波
 羅蜜多故。心無罣礙。無罣礙故。無有恐怖。若
 遠離顛倒夢想。究竟涅槃。三世諸佛。依般若
 波羅蜜多故。得阿耨多羅三藐三菩提。是無上
 。是無等等咒。能除一切苦。真實不虛。故說
 般若波羅蜜多咒。即說咒曰。揭諦揭諦。波羅
 揭諦。波羅僧揭諦。菩提薩婆訶。

1. 心經科文(改自周止庵心經註)

- 心經科文(周止庵)
 - 總釋名題
 - 正釋經文
 - 序分
 - 正宗分
 - 顯說般若
 - 因人顯法
 - 正示法空
 - 明蘊空德
 - 顯空德
 - 顯彰妙果
 - 結讚功能
 - 密說般若
 - 流通分

1. 心經(鳩摩羅什梵文漢譯) 2. 心經(梵音漢字)

吠怛多三迦南散娑三迦喃。尾迦南伊賀。舍利
 補怛囉。薩囉縛達磨成涅多洛叉拏。阿耨多婆
 左阿哩也嚧駄。阿麼●。阿尾也嚧●。阿怒那
 阿婆哩布●怛。多娑磨舍哩布怛囉。成涅哆耶
 成涅多薩嚧吽。曩吠怛多。散迦。南三散娑迦
 。喃尾●南。囊斫●助嚧囉迦囉拏即賀縛迦野
 頗那悉莎嚧吽。攝那彥駄囉娑娑播囉瑟吒尾演
 達摩。曩斫●駄都。也縛那麼怒尾迦南駄都。
 那尾南。那尾南。吒喻那尾南。又欲夜縛。那
 樣囉麼囉南。那左囉麼囉拏叉喻。那耨法娑敏
 那野室嚧駄縛室哩●尾演●羅●婆囉明多。縛
 室利●尾賀囉野枳陀婆囉●枳陀婆囉●。那悉
 帝達摩那●囉尾地那備。娑曼也駄三曼。曩●
 囉尾耶奴答覽。三藐三牟備牟備。娑曼駄縛。
 室里●尾演。●囉●娑囉彌陀。摩訶曼但羅。
 摩訶尾欲曼怛羅。阿耨多羅曼怛羅。阿娑婆三
 備曼怛羅。娑婆奴迦●囉。舍婆那娑帝地
 演達摩。●囉●婆囉明多。目訖姤滿但羅。怛囉
 嚧。盧諦盧諦。播羅盧諦。播羅僧盧諦。菩提
 娑婆訶。

《心經》內容與意義的標誌例舉

❁ 心經各版本與其科文間內容的對應

- ❁ 給人用的標誌介面。

- ❁ 標誌的彈性：

 - ❖ 無固定的標籤集。

- ❁ 內容與結構(ontology)間的對應：

 - ❖ 對本文的詮釋。

- ❁ 版本間內容的對應：

 - ❖ 對本文呈現的不同形式。

- ❁ 結構間的對應：

 - ❖ 對本文不同的詮釋。

- ❁

《心經》內容與意義的標誌例舉

- ✿ 以上展現的內容處理，都沒有用到檢索常用的詞語聯繫(morphological linking)結構。
- ✿ 以上的展現可呈現有個別差異的「別相」
- ✿ 請注意意義聯繫間的彈性：
 - ✿ 跨語文
 - ✿ 跨作者
 - ✿ 跨版本
 - ✿ 知識結構的延伸
 - ✿ 各種意義之間的聯繫關係……

❁ 綜觀內容標誌的特徵，至少包含下列項目：

- ✧ 它是給人用的標誌介面。
- ✧ 標誌的彈性：它無固定的標籤集。
- ✧ 它可往復的顯示文字與內容結構（ontology）間的對應，以及對本文的詮釋。
- ✧ 它可呈現本文不同形式、版本間內容的對應。
- ✧ 它可呈現結構間的對應
- ✧ 因為「語意」的表達是跨越國界、跨越文化的。它可跨越不同的語言文字、不同的作者、不同的版本等。

- ❁ 在此舉例的雛型軟體，可以證明電腦是可以處理語意的。
 - ✱ 展現的內容處理，都是以內容標誌為工具所做的意義檢索，完全沒有用到檢索常用的詞語聯繫（morphological linking）結構。而意義檢索可呈現「共相」，亦可呈現有個別差異的「別相」（如印順與周止庵對《心經》個別的註釋）。
- ❁ 做文章內容的標誌，不是電腦工程師可以做的，需要了解文章內容的專家來做。這情形正好提供人文學者一個絕佳的機會加入文獻數位化的行列。
 - ✱ 如果人文學者能用標籤把他們的知識，也就是對文章的理解、真知灼見，表達給電腦知道，那麼，久而久之電腦將匯集大量的人文知識。果真如此，那麼，一種嶄新形式的人工智能（artificial intelligence）即將誕生，且讓我們拭目以待。

結語

後設資料和內容標誌

- ❖ 後設資料和內容標誌並不相互排擠，它們是兩種類型完全不一樣的工作。換言之，後設資料和內容標誌兩者都是不可缺的，且彼此相輔相成、相得益彰。
- ❖ 若認為：除了文物數位化的本身之外，所有其他的資料都屬後設資料，那麼就犯了大錯—扼殺了內容標誌生存的空間。

意義處理的問題

- ❁ 意義(meaning)是可以用電腦處理的。
 - ✱ 內容標誌即文章意義的標誌。
 - ✱ 從《心經》的例子可知：情境在電腦中可以表達。據此可以設法處理多義或歧義的問題。
- ❁ 未來，電腦可能以兩種方式來處理意義問題：
 - ✱ 其一是將所有的多義關係轉化為單義的語法關係。例如，建立「常識庫」讓電腦能辨識「情境」。
 - ✱ 其次是與人合作，以人機共建的系統來做「了解」和處理「意義」問題。
 - ❖ 這就是內容標誌要做的事。

謝謝聆聽

歡迎批評指教！