

韓國語文章 處理시스템에 관한 考察

柳 京 熙
〈KORSTIC 電算室長〉

案을 생각해 보겠다.

1. 머리말

컴퓨터는 英數學을 쓰는 나라에서 만들어진 것이기 때문에 英數字의 처리에 있어 아주 效果的이다. 当初에는 컴퓨터가 數字처리를 위하여 만든 것이기는 하나 이것을 英文字處理에 応用하는 技術이 向上되어서 오늘날에는 文字처리가 一般化되었다. 그러나 文法이란 것은 數의 世界의 各種 法則보다 훨씬 어려우므로 文章處理技術이 確立되기까지는 상당히 오래 걸린 셈이다.

데이터베이스(DB)에도 數值DB와 非數值DB가 있는데 非數值인 경우에는 그 實用化가 늦었기 때문에 數值DB만이 데이터 베이스를 代表하고 있는 것처럼 알려지고 있다. 그러나 데이터베이스의 效用性은 데이터의 共同利用과 情報檢索에 있다고 볼 수 있기 때문에 數值DB만으로는 滿足할 만한 情報檢索을 期할 수가 없다. 特定한 事物에 관한 情報를 찾으려 할 때 그 事物의 分類記号나 符號를 일일이 외워두기가 힘들기 때문에 一般化하기 어려우며 오히려 말로서 찾는 것만 못한 경우도 많다. 따라서 情報檢索에서 制限된 能力을 가지고 있는 數值DB보다 文章處理가 될 수 있는 DB가 바람직하며 경우에 따라서는 둘다 할 수 있는 것이 좋다.

筆者는 우리말 文章을 處理할 수 있는 “우리말 DB시스템”을 構想함에 있어서 尙存하고 있는 몇가지 問題點을 들어서 이에 관한 解決方

2. 우리글의 制限

한글專用과 漢字混用의 問題는 오래묵은 論争의 対象이며 앞으로도 언제 해결될지 모르는 커다란 課題이다. 한글의 字素와 모아쓴 字種에 관한 問題는 이미 解決된 것으로 보아도 좋을 것이다. 그러나 漢字의 使用制限問題는 여전히 남아 있다.

요컨대 常用漢字만으로 하느냐 擴張한다면 어떤 基準으로 하느냐 하는 것이다. 특히 人名, 地名과 같은 固有名詞로 된 漢字語때문에 字種의 數가 거의 無制限이 되어버린다. 筆者는 使用漢字를 確定지은 다음에 機械化한다는 것은 너무 長期間이 所要되기 때문에 使用頻度에 따라 第一水準 3,200字 第二水準 3,200字…… 등으로 確定해 나가면서 機械化를 서둘러야 할 것으로 생각된다.

3. 情報의 最小單位

컴퓨터가 文章處理를 함에 있어서 가장 基本이 되는 것은 字素와 字種別로 處理하거나 간에 同一한 結果를 얻는 것이 바람직하지만 이것이 必須要件이라고 하기는 어렵다. 漢字인 경우에는 字素(部首)別과 字種別 處理를 区分할 理由를 찾기가 어렵다. 왜냐하면 漢字의 字素(部首)의 數와 順序가 많고 複雜하기 때문에 한글, 알

要한 일이다.

5.2 字種의 識別

前述한대로 字種은 情報의 最小單位로 삼는 것이 좋을 듯하다. 다만 한글인 경우에는 字素가 最小單位이기는 해도 “모아쓰기 프로그램”에 의해 字種으로 바꾸는 것이 바람직하다. 다만 한글字素데이터의 處理와 한글字種데이터의 處理에서 同一하게 結果를 내어야 한다는 固定觀念은 없어도 좋을 것이다. 字種의 處理에서는 目的에 따라서 다르게 設計할 수 있으나 여기에서는 한글과 漢字의 区分이 自動적으로 되도록 하는 것으로 充分하다고 본다.

5.3 낱말의 識別

英文의 Text Processing을 위한 單語識別은 빈칸으로 하고 있다. 英文은 每單語마다 띄어쓰도록 되어 있기 때문에 單語識別을 自動化하기 쉽다. 그러나 우리말인 경우는 토씨를 붙여서 쓰도록 되어 있어서 英文과 같이 完全 띄어쓰기가 아니라 半띄어쓰기이다. 이러한 경우의 單語識別을 위하여서는 데이터 自体의 加工이 必要하다. 즉 데이터를 完全 띄어쓰기로 바꾸거나 또는 每토씨앞에 識別子를 追加人力토록 한다면 自動化가 可能하다. 筆者는 識別子의 機能으로서 슬래시(斜線)를 쓰는 것이 좋을 것이라고 본다. 컴퓨터가 빈칸, 슬래시 또는 기타 特殊文字로서 單語를 識別하고 出力시킬 때에는 슬래시를 無視해 버리면 된다.

5.4 不用語 또는 索引語의 識別

KWIC(Key Word in Context)이란 文章處理를 함에 있어서 重要語(Key Word)인지 아닌지를 識別하기 위하여 重要하지 않으면서 出現頻도가 높은 單語를 모아서 不用語(Stop Word) 目錄을 컴퓨터에 담아두고, 識別된 모든 單語 가운데서 不用語가 아닌 單語를 重要語로 取扱하도록 만든 시스템을 말한다.

이와 反對로 KWOC(Key Word out of Context)이란 重要한 單語를 미리 選定하여 컴퓨터에 담아두고 識別된 모든 單語 가운데서 重要語와

一致하는 單語만 가려내는 시스템을 말한다.

어느 것이나 重要도가 問題가 되며 重要도를 가리는 基準이 曖昧하다. KWIC인 경우에는 別로 重要하지 않은 單語가 重要語로 取扱될 우려가 있으며 KWOC인 경우의 重要度判定은 너무 主觀的이다. 그러나 重要度判定을 위한 資料로서 KWIC부터 着手하는 것이 올바른 順序일 것이다.

KWIC시스템인 경우에는 不用語인지 아닌지를 識別하게 되고 KWOC시스템인 경우에는 重要度인지 아닌지를 識別하게 된다.

5.5 文章의 識別

英文이나 우리글의 文章은 句讀點으로 区分하고 있다. 마찬가지로 컴퓨터도 句讀點을 文章의 識別子로 삼으면 된다. 그러나 每文章마다 所在를 가리키는 別途의 識別子를 두는 것도 생각해 볼만한 일이다. 특히 데이터베이스를 만드는 경우에는 文章마다 꼬리표의 役割을 하는 識別子가 情報檢索을 위하여 반드시 必要하다.

5.5 5.5 段落의 識別

새로운 段落를 쓸 때에는 一般的으로 줄을 바꾸고 몇칸의 빈칸을 둔 다음 시작한다. 시스템도 물론 이것을 識別할 수 있어야 하며 追加해서 매 段落마다 그 所在에 관한 情報, 그 段落가 包含하고 있는 表現方法을 考慮해 볼만하다.

6. 데이터·코우드와 文章處理

情報處理能力을 極大化하기 위하여 意味를 가진 데이터를 略号나 數字로서 縮小表現하려는 試圖는 오래전부터 있었다. 이에 따라 각 시스템間的 데이터互換性を 向上시키기 위한 “데이터·코우드의 統一制定” 事業이 先進各國에서 相當히 推進되고 있으며 ISO의 여러가지 勸告도 있다.

우리나라에서도 內務部가 主管하여 制定 實施하고 있는 住民登錄番号는 데이터 코우드의 役

割을 充分히 感當하고 있으며 이로써 個人의 識別을 하고 있다.

그밖에도 모든 意味를 가지는 事物의 데이터·코우드化는 必要하다.

그러나 데이터·코우드 制定을 위한 分類体系 確立의 어려움 그리고 分類体系變動에 따른 更新의 어려움 등으로 早速한 實現을 期하기는 어려운 實情이다.

情報檢索이란 側面에서 볼 때 데이터·코우드에 의한 檢索은, 시스템은 索引語 - 데이터·코우드 参照表를 만든다든가 하는 追加的인 機能이 必要하고 利用者는 여기에 익숙하여야 한다든가 하는 要件이 수반된다. 대부분의 DBMS(데이터베이스 管理시스템)에는 이른바 文章處理 機能이 없다. 간혹 있는 것도 있으나 그나마 英語文章에서 單語程度에 不過하다.

理想的인 情報檢索시스템은 檢索하고자 하는 實體의 全体를 모르더라도 皮相的인 프로필(편모)만 가지고도 實體를 알아낼 수 있어야 한다. 왜냐하면 情報要求者의 대부분은 必要한 情報를 分明하게 表現하지 못하기 때문이다. 코우드檢索보다 重要語 또는 自然語에 의한 檢索이 훨씬 利用하기 便利하다.

7. 맺음

文章處理技法은 모든 出版社와 新聞社와 같이 編輯校正을 必須的으로 遂行하여야 하는 곳의 電算化에는 반드시 考慮되어야 한다. 뿐만 아니라 索引作業에 莫大한 人力과 經費를 投入하고 있는 곳도 마찬가지이다. 또한 電算化된 資料를 순식간에 찾아 낼 수 있는 데이터·뱅크에는 더욱더 必要하다. 즉 編輯, 校正, 索引, 檢索을 自動化할 수 있는 道具가 된다.

우리말 데이터베이스를 製作함에 있어서 컴퓨터 可讀型으로 製作된 資料를 活用하면 製作費가 아주 싸진다. 데이터更新을 해야 할 경우에는 더욱더 값싸게 된다.

오늘날 우리가 必要로 하는 데이터의 管理가 너무 소홀한 傾向이 있다. 모든 데이터를 早速히 컴퓨터 可讀型으로 바꾸는 作業에 着手하였으면 한다. 컴퓨터 시스템, 데이터通信 등에 관한 問題는 데이터가 있고 나서부터 考慮되어도 늦지않다. 우리의 情報는 우리말로 入力處理, 出力될 수 있는 시스템과 더불어 "National Data Base"가 構築되도록 關係當局에 建議하는 바이다.