

「中文分詞語料庫」說明

(一) 簡介

中文分詞語料庫為一包含兩百萬詞、不含詞類標記的語料庫，每個文句皆根據分詞原則來斷詞。而此分詞原則，乃是中華民國計算語言學學會在經濟部中央標準局委辦的「資訊處理用中文分詞規範調查研究及草案研擬」計畫中所訂定的。本語料庫來源包括書面語和口語兩部分，其中資訊類佔21%。

(二) 內容說明

1. 檔案說明

共有 94 筆 *.seg 檔案，每個檔案為 UTF-8 編碼。

2. 語料庫樣本內容

64. 。我真正認識錢院長還是在一九八三年被選為中央研究院院士之後。
65. 。那時為了在中研院籌設新的原子與分子科學研究所，
66. ，我們曾有過頻繁的接觸與深入的討論，