

## (甲) 切詞規則 Segmentation Guidelines

規 則	例
<b>1</b> 名詞部分	
.1 古今中外各族人名的姓名、譯名概不切分。  附加「二世」之類的也不切分。	張勝利   歐陽海   江澤民 卡爾·馬克思   埃弗斯--威廉斯  約翰保祿二世
.2 對帶有普遍性或半確定性的稱謂詞，帶稱謂或職稱者，不切分。  一般性的不合。	李氏   王太   李伯   小張   華仔   陳x偉   陳總   張導  陳妻   李嫌   陳童   呂女   林父   陳犯   李翁
.3 國家全名不切分。	中華人民共和國
.4 民族名、地名中的「族、省、市、州、鄉、區、江、河、山」等不切分。	哈薩克族 北京市   浙江省   福田區
.5 區域名是雙音節或以上的切分。	香港 地區   廣西 自治區
.6 地理名稱不切分。	九龍半島   松嫩平原   青海高原
.7 村名、街道名全稱不切分。	黃大仙上村   皇后大道東

<p>.8 樓宇名稱不切分， 但由複合詞組成者則切分。</p> <p>.9 地名含單音詞的不切分，由複合詞組成者則切分。</p> <p>.10 商標名、公司名、機構名由單音詞殿後的不切分； 由複合詞組成的切分。</p> <p>.11 報刊、書籍、影劇、通訊社等專用名稱為單音詞的合，為複合詞的切分。</p> <p>.12 “單字名詞+單字方位詞”用來表示地點、機構或時間的，即使是表達本義的，只要有一定的頻率，一概從合。  所表達超出本義或已虛化者合。  “兩字名詞+單字方位詞”表實義的，一律從分。</p> <p>.13 名詞性的並列結構概不切分。</p> <p>.14 經濟名詞或經濟術語不切分。</p>	<p>東翠樓 華興大廈   東方海外商業大廈</p> <p>中港城   黃埔花園   皇后碼頭</p> <p>樂聲牌   國泰 AK-47 型   彩虹戰士號   天弓二型</p> <p>明報   聯合早報   金玉集   鄧小平文選   法新社   新華通訊社</p> <p>家裏   屋前   床上   村裏   水下   河邊   黨內   選前   賽後   生前   婚後   會後</p> <p>手上   身上   手下   根底</p> <p>學校裏   世界上</p> <p>中英   中港台</p> <p>息率   減息   護盤   跌幅</p>
<p><b>2 數詞</b></p> <p>.1 由基數詞及數位詞組成的各種數詞概不切分，數量結構則予切分。</p> <p>.2 序數詞切分。</p> <p>.3 百分數以及含有「數、幾」二字的約數不切分。</p> <p>.4 其餘約數切分。</p>	<p>一千八百   一九九七   34   1億3502萬3千   十七八歲</p> <p>第一   第三十四</p> <p>百分之三十   幾千分之幾   1又12分之1.67   數萬   十幾二十</p> <p>十來個   三千餘人   近三千多人</p>

<p><b>3 時間詞</b></p> <p>.1 年、日、時、分、秒切分。</p> <p>.2 星期、月、日的名稱合。</p> <p>.3 兩個字的合。</p> <p>.4 三個字: A.屬套合詞者不切分。</p> <p style="padding-left: 40px;">B.非套合詞者應切分</p> <p style="padding-left: 40px;">C.含「個」的時間詞應切分。</p> <p style="padding-left: 40px;">D.「月份」「年度」不切分。</p> <p style="padding-left: 40px;">E.「來」字應切分。</p> <p style="padding-left: 40px;">F.其他</p> <p>.5 四個字的按以上規則切分。</p>	<p>1997 年   15 日   三 時 廿五 分</p> <p>周一   星期日   一月   八月   十二月   初一   初八   年初二 大年初一</p> <p>上月   下周   去年   前天   後年   明年   周末   月初</p> <p>上月底   下周一   八月初   前年中</p> <p>本 月初</p> <p>上 個 月   下 個 月   半 個 月</p> <p>八月份   下月份   上年度</p> <p>多 年 來   幾 天 來   近 年 來</p> <p>較 早 時   大 前 天   大 清 早   霎 那 間   一 時 間</p> <p>上 星 期 三   上 星 期 末   下 世 紀 初</p>
<p><b>4 略語</b></p> <p>.1 簡略語概不切分。</p> <p>.2 合併詞不切分。</p> <p>.3 套合詞合。</p>	<p>四化   亞運會   男單   臨立會 港人   通脹   江八條   三違反</p> <p>國內外   中小學   大中型   中低收入</p> <p>直升機場   心臟病發   羽毛球拍 精神病患   國民黨籍   大連市長 生態學家   國防部長</p>

## 5 二字結構

- 1 在兩個語素中其中有一個語素是粘著的(B)，其結合具有詞匯意義的，有一定使用頻率的，應視為一個詞而不予切分。
- 木欄 | 鋼板 | 虎骨 | 雞爪 |  
母羊 | 益蛇 | 豆油 | 菜油 |  
校門 | 併軌 | 供氣 | 排毒 |  
排石 | 設攤 | 貪睡 | 修機 |  
節電 | 拔拳 | 搞混 | 刮擦
- 某些粘著的(B)的語素結合面極廣，幾乎最高 法院  
都可與另一語素結合成詞，為免造成此類最後(做連詞用)  
詞語太多，故只收錄具詞匯意義的詞條或你 最好 別管 這 件 事  
文言詞，餘者不予收錄(在處理上即予切全 公 司 中 張 三 身 材 最 高 李 四 最  
分)。  
矮  
他的成績最好
- 兩個語素均是自由的(F)，而其組合又不吵 醒 | 走 往 |  
具詞匯意義的，或使用頻率不高的，原則飛 往 | 花 香  
上予以切分。
- a 兩個語素均是自由的，但其組合後引生豬 | 撞破  
申出新的詞匯意義，即詞義產生了轉義，(有詞匯意義或轉義)  
一加一已不等於二，則不予切分而視為  
詞。
- b 兩個語素的結合具顯著性，往往表達熊膽 | 蛇膽 | 紅花(中藥名)  
特別意義，且已相當普遍在使用，此類詞上崗 | 下崗 | 到位 | 待崗  
應視為新詞而予收錄。上證 (大陸新詞)
- c 有不少兩個自由語素組合的雙音詞，打壓 | 情治 | 研採 (台灣用語) |  
結合緊密，使用頻率十分高，已由詞組慢停車 | 回國  
慢演變成詞，則應考慮予以收錄。
- d 某些與英文單詞有相對應的雙音詞組牛肉(beef) | 豬肉(pork)  
合與否，仍應視乎使用頻率，對應因素僅鹿 肉(venison) (使用頻率不高)  
作參考。
- 3.1 凡屬文言詞、古語、縮略語，一般叩門 | 易主 | 擒兇 | 拔拳 |  
不予切分。 十佳
- 某些動賓關係的雙音詞，在具體語境中本 趟 列 車 不 掛 車 運 行 。  
會演變成偏正結構的詞語。對這種偏正結本 趟 列 車 共 有 八 卡 掛 車 。  
構的雙音詞，不予切分。 沉 艦 | 用 電 | 用 字

.2	並列、主謂、聯動、偏正、動賓、動補等結構:	
	A.二個字均為自由語素、或從屬關係者切分。	寫信   綠葉
	B.其中一個黏著語素是相當自由的，切分。	本黨   李某
	C.其中一個或兩個都是黏著語素的合。	手部   鞋類   逃離   打蠟   頭痛   飛抵   打壞   受訪   做好   看成   女童   力證
	D.結構緊密、使用穩定的合。	回國   出院
	E.切分後意思基本不變的切分，切分後意思有變的不切分。	本國   全黨 綠表   白菜
.3	動(形)詞 + 趨向動詞(介):	
	A.下列五種情況不予切分: 1.趨向動詞(介詞)已虛化的，不表示具體方向、方位的，或是表示結果、程度的，或表示引申意的; 2.動詞是黏著語素的; 3.動詞是文言詞的; 4.動詞是方言的; 5.使用穩定、結構緊密的高頻詞。	敢於   歸於   過於   急於   樂於   達到   感到   遭到   受到   超出   發出   付出   退出   駁回   挽回   召回   進去   看去   前去   失去   揭開   解開   半開   避開   岔開   落入   流入   介入   併入   步入   給予   賜予   賦予   寄予   交予
	B.下列兩種情況切分: 1.趨向動詞(介詞)有實義的，表示具體方向、方位、地點、時間的; 2.動詞是自由語素的。	定於   生於   用於   走到 帶到   回到   調到   跪下 坐下   彈出   浮出   爬出 流出   調回   退回   收回 帶來   借來   喚來   吹去 走進   衝進   住進   放進

<p>6 三字複合詞</p> <p>.1 動(形)詞 + 趨(介)切分。</p> <p>.2 從屬關係切分。</p> <p>.3 接頭詞下列各字合，餘者酌分： 准 總 半 後 禁 無 核 性 零 軟 硬</p> <p>.4 接尾詞盡量不分，包括： 班 表 兵 波 部 廠 場 車 處 袋 燈 地 點 店 動 度 段 額 法 發 犯 房 費 工 觀 館 光 棍 漢 戶 花 會 火 機 劑 件 界 酒 局 科 庫 礦 力 量 令 樓 率 論 迷 面 民 年 派 片 品 器 氣 切 人 生 師 史 士 式 室 屬 所 稅 獎 司 台 堂 體 天 廳 團 網 尾 物 系 線 箱 形 型 學 炎 業 邑 儀 意 園 院 症 質 種 眾 組</p> <p>.5 接尾詞「上、下」的用法已超出本義或已虛化者，不予切分；餘者分。(這類結構前面如能加“在”而不影響其接受性者，就需切分。)</p>	<p>走 過來   退 回去   優 越 於   接 觸 到</p> <p>美 國 人</p> <p>准 新 郎   總 參 謀   核 威 嚇   性 開 放   零 事 故   軟 環 境   硬 道 理</p> <p>超 時 代</p> <p>人 權 法   介 紹 費   午 飯 錢   手 續 費   支 持 票   主 席 台   召 集 人   外 圍 馬   伙 食 費   伙 食 錢   休 息 室   光 榮 感   在 野 黨   收 錢 者   住 宿 費   助 選 團   戒 嚴 令   抗 議 信   投 注 人   投 注 額   投 訴 信   私 家 車   和 平 獎   承 建 商   社 會 課   保 險 界   軍 政 團   重 案 組   特 別 法   託 運 人   退 休 金   贊 成 票</p> <p>實 際 上   事 實 上   大 致 上 天 空 上   手 續 上</p>
<p>7 四字詞</p> <p>.1 凡屬下列類型的四字結構，不予拆分： a. 成語； b. 使用穩定的古語； c. 四個字不能任意替換的； d. 有一定格式的短語。</p>	<p>欣 欣 向 榮   愚 不 可 及   不 爭 之 論 坐 視 不 理   下 落 不 明   相 持 不 下 由 此 可 見   易 燃 易 爆   一 心 一 意 半 生 不 熟   大 吃 一 驚   猶 豫 不 決</p>

<p>.2 常見的政治術語、慣用語合。</p>	<p>一中一台   一國兩制   第一時間 一個中國</p>
<p><b>8 短語</b></p> <p>超過四字以上的短語概予切分。</p> <p>可兩分可三分的一概兩分。</p> <p>不可兩分的、並列結構的、專名等則三分。</p>	<p>財政 預算案   公園 管理處</p> <p>新聞 自由 獎   服裝 首飾 展   高等 教育 部   民主 進步 黨</p>
<p><b>9 疊詞</b></p> <p>.1 AA 或 AABB 式疊詞合。</p> <p>.2 ABB 式疊詞合。</p> <p>.3 AAB、ABAB、A 一 A、A 了 A、A 了一 A 式疊詞切分。</p> <p>.4 AAB=BAA,合</p>	<p>看看   個個   人人   明明白白 清清楚楚</p> <p>綠油油   眼睜睜</p> <p>說說 看   研究 研究   想 一 想</p> <p>哈哈笑   笑哈哈 (但:睜睜眼 ≠ 眼睜睜)</p>
<p><b>10 非漢字部份</b></p> <p>.1 所有非漢字符號，包括其它語言的字母、詞串等，保留原有形式，中間不再切分。(註2)</p> <p>.2 略語、縮寫語均予收錄。</p> <p>.3 人名不再切分。(註2)</p> <p>.4 地名不再切分。(註2)</p> <p>.5 機構名稱不再切分。(註2)</p> <p>.6 由英文字母、數目字或「/」、「-」等組成的用以表示型號的詞語不再切分。但複合詞組仍須切分。</p>	<p>DOS   H<sub>2</sub>O   cc   4<sup>2</sup>=16   World News Asia</p> <p>RAM   SOS   BNO   APEC   亞太 經合 會議 (APEC)</p> <p>Peter Hongda Wu   Franjo Tudjman</p> <p>St. Clare Ave. W</p> <p>Osaka University</p> <p>AK-47   AE-100   波音 A2-747 型 飛機</p>

<p>.7 由英文字母、數目字或「/」、「-」等組成的用以表示序號的詞語不再切分。</p> <p>.8 詞串中含有「x」(叉)者，不再切分。</p> <p>.9 一般英語單詞或短句，在文中無特別含意者，視為一個完整的詞語，中間不再切分。（註2）</p> <p>.10 不完整詞語，仍予切分。</p>	<p>A4 版   三樓 A 座   第 96M 1253 號   車牌 BT4523</p> <p>車牌 HK19x7</p> <p>Supreme Governor of the Church of England   gun   abc</p> <p>澳 (大利亞)</p>
<p><b>11 其他</b></p> <p>.1 古語不切分。</p> <p>.2 方言不切分。</p> <p>.3 形式相當確定的熟語不切分。</p>	<p>是次   甚為   事發   不果   為免   問及</p> <p>唔使指意   唔通   玩</p> <p>踏破鐵鞋無覓處</p>

## (乙) 命名實體標記 **Named Entities Annotation**

類別	說明	例子
<b>人名</b> <b>PERSON</b>	1. 古今中外各族人名的全名或名字、譯名 2. 筆名 3. 虛構的人名 4. 歷史名人外號、帝號 5. 暱稱、綽號	胡錦濤、劉德華、克林頓、雅子、小鳳  金庸、瓊瑤、小草、阿寬 令狐沖、步驚雲、則卷千平、大雄 秦始皇、武則天、康熙、宋太祖  華仔、阿姐
<b>地名</b> <b>LOCATION</b>	1. 國家或地區的全名、簡稱、或特指 2. 省、市、州、鄉、區、江、河、山的名稱 3. 大廈、樓宇、屋苑名稱 4. 街道名 5. 洲名、海洋	中華人民共和國、港、英、上海、滬、大陸（指中國）、兩岸 北京市、福田區、峨嵋山  東翠樓、英明苑、中港城  彌敦道、皇后大道東、鴨寮街、長安大街 美洲、太平洋
<b>機構/其他專名</b> <b>ORGANIZATION</b>	1. 公司及機構名 2. 團體及隊伍名 3. 大型會議及賽事 4. 特別日子及事件	長實、證監會、路透社、教統局 皇馬、自由黨、湖人 奧運、十五大 九一一、十一、九七