

語料庫應用於教材編纂

語文教育及編譯研究中心 吳鑑城

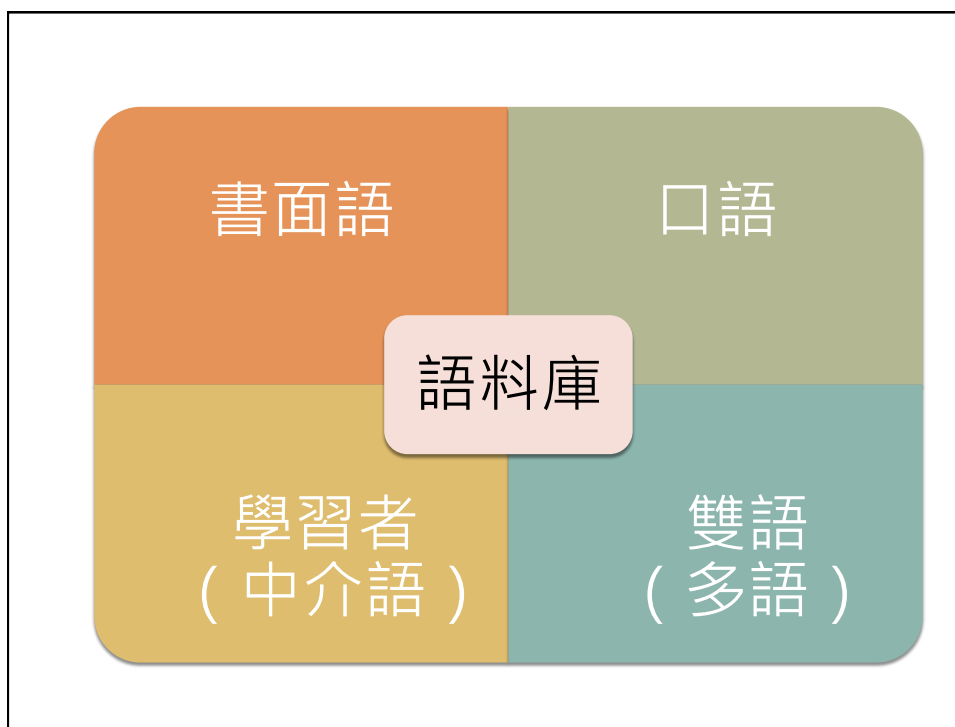
語料庫應用於教材編纂

- 語料庫概念介紹
- 語料庫技術介紹
- 語料庫工具應用於教材編纂介紹
 - ▣ 中文分詞工具
 - ▣ 華語文錯字自動偵測系統(雛型)
 - ▣ 索引典系統
 - ▣ 語義場工具

語料庫？有人用嗎？

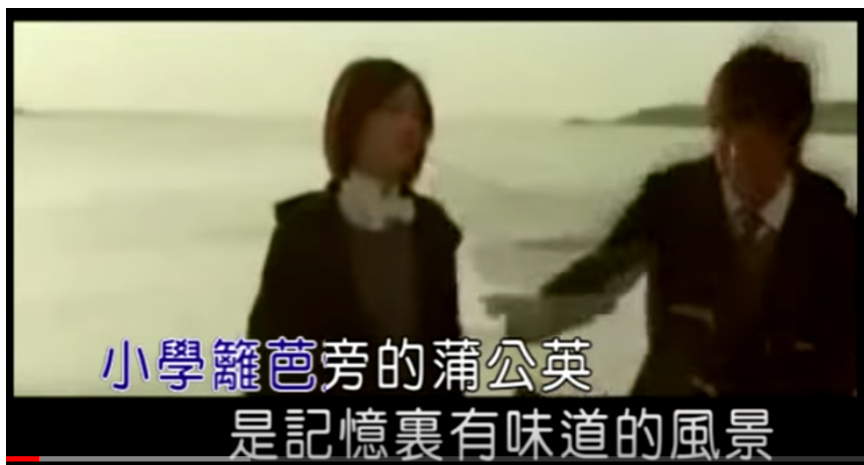
- 藉學術語料庫提出中文學術常用詞表: 以人文社會科學為例
 - ▣ 劉貞好·陳浩然·楊惠媚·2016
- 基於語料庫的英語專業教材詞彙研究
 - ▣ 宋曉舟·2016
- 使用專業語料庫增加籃球英語教材詞彙及搭配詞的豐富度
 - ▣ 高睿禹·2017

“語料庫是為了特定目的，根據特定原則搜集或取樣，並分門別類集合起來的一批語言材料”





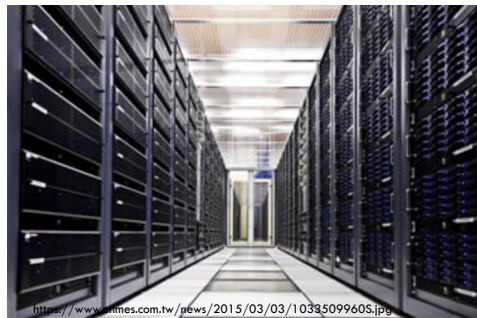
若我們所知，亞洲在世界上就有豐富多彩文化與歷史而且有明，也有很多人欣賞。說到文化和歷史方面，亞洲的確無人不知，尤其在手工藝的部分，自從一前手工藝的技藝-非常讓人不可界議，什麼都可做得到



小學籬笆旁的蒲公英
是記憶裏有味道的風景

圖片擷取自：<https://www.youtube.com/watch?v=DYyE68llx0U>

語料庫技術！？



Google

Baidu 百度

百度一下，你就知道！

bing

機會

想要成功，要培養實力並把握機會。

我坐飛機頭暈。



人參

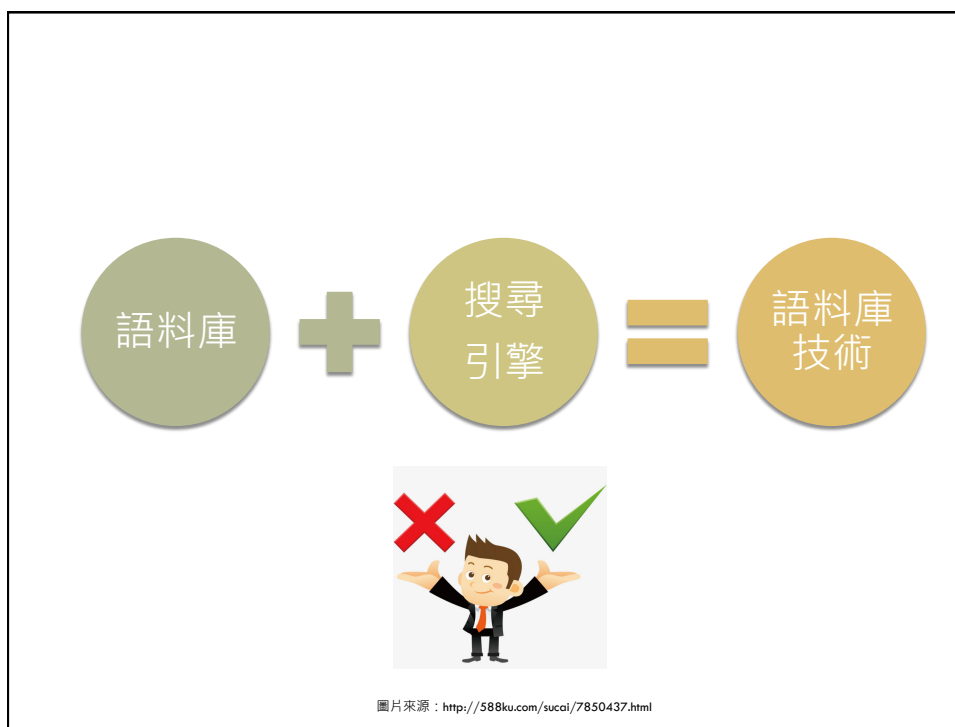
人參是許多家庭選購的拜年伴手禮。

老師說明天每個人參加大隊接力時要加油。



My dog also likes eating meat.

my dog also like eat meat .



國教院分詞系統

請輸入需要分詞的文本

市面上很少有「教科書設計」的專書，因為我們總覺得那是出版社的事！然而，真的是這樣嗎？教科書設計其實與課程綱要、教師的教學、學生的學習息息相關，是課程、教學、學習三位一體中一個重要的環節。除了有教育學與學科專業等內容涵納其中，也與編輯、版式等視覺設計元素的概念有關。有鑒於此議題的重要，本院教科書發展中心邀請淡江大學課程與教學研究所陳麗華所長，於8月27日上午進行「教科書設計研究」專題演講，除了院內同仁，也邀請出版社編輯企劃相關人員參與。

市面(Nc)上(Ncd)很少(D)有(V_2)「(PUNC)教科書(Na)設計(Na)」(PUNC)的(DE)專書(Na)，(PUNC)因為(Cbb)我們(Nh)總(D)覺得(VK)那(Nep)是(SHI)出版社(Nc)的(DE)事(Na)！(PUNC)然而(Cbb)，(PUNC)真的(D)是(SHI)這樣(VH)嗎(T)？(PUNC)教科書(Na)設計(Na)其實(D)與(P)課程(Na)綱要(Na)、(PUNC)教師(Na)的(DE)教學(Na)、(PUNC)學生(Na)的(DE)學習(Na)息息相關(VH)，(PUNC)是(SHI)課程(Na)、(PUNC)教學(VA)、(PUNC)學習(VC)三位一體(Na)間(Ng)一(Neu)個(Nf)重要(VH)的(DE)環節(Na)，(PUNC)除了(P)有(V_2)教育學(Na)與(Caa)學科(Na)專業(VH)等(Cab)內容(Na)涵納(VJ)其中(Nep)，(PUNC)也(D)與(P)編輯(Na)、(PUNC)版式(Na)等(Cab)視覺(Na)設計(Na)元素(Na)的(DE)概念(Na)有關(VJ)。(PUNC)有(V_2)鑒於此(VH)議題(Na)的(DE)重要(Na)，(PUNC)本(Nes)院(Nc)教科書(Na)發展(Na)中心(Nc)邀請(VC)淡江(Nb)大學(Nc)課程(Na)與(Caa)教學(Na)研究所(Nc)陳麗華(Nb)所長(Na)，(PUNC)於(P)8(FW)月(Na)27(Neu)日(Nf)上午(Nd)進行(VC)「(PUNC)教科書(Na)設計(Na)研究(Na)」(PUNC)專題(Na)演講(Na)，(PUNC)除了(P)院(Nc)內(Ncd)同仁(Na)，(PUNC)也(D)邀請(VC)出版社(Nc)編輯(VC)企劃(Na)相關(VH)人員(Na)參與(VC)。(PUNC)

http://coct.naer.edu.tw/Segmentor/

自拍是年輕人回應生活的方式。



中文分詞

自 拍 是 年 輕 人 回 應 生 活 的 方 式 。

新詞發現：自拍、都更、網軍....

騎機車要戴安全帽，以確保生命安全。

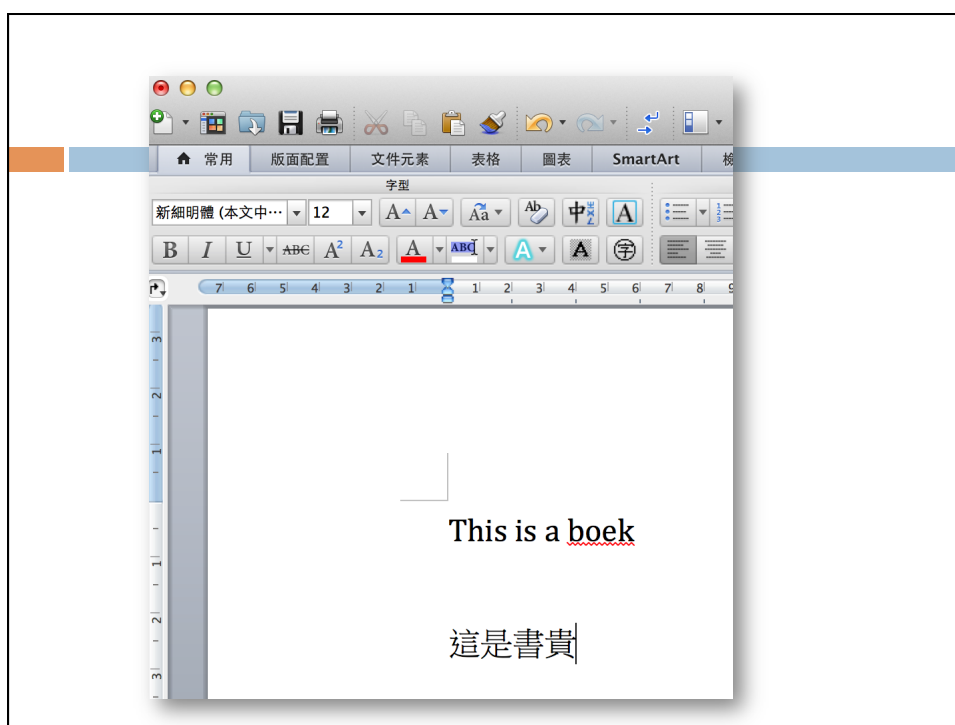
騎(VC) 機車(Na) 要(D) 戴(VC) 安全帽(Na) ，(PUNC) 以
(Cbb) 確保(VE) 生命(Na) 安全(Na) 。



真的很機車耶，幹嘛偷看我的LINE？

真的(D) 很(Dfa) 機車(VH) 耶(T) ，(PUNC) 幹嘛(D) 偷看(VC)
我(Nh) 的(DE) LINE(FW) ？(PUNC)

新義發現！？



五、改錯(二題共 5 分)

1. 大仁具有高瞻遠矚的見識，做事不會猶豫不決更不會推謝責任。
2. 美美發奮練舞，和同伴切磋舞藝，不敢半途而費，盼望自己能有精統的演出，出人頭地。

新北市板橋區江翠國中 101 學年度上學期七年級國文科第三次段考試題題目卷

教育部常用國字辨似

【芭】ㄅㄚˊ

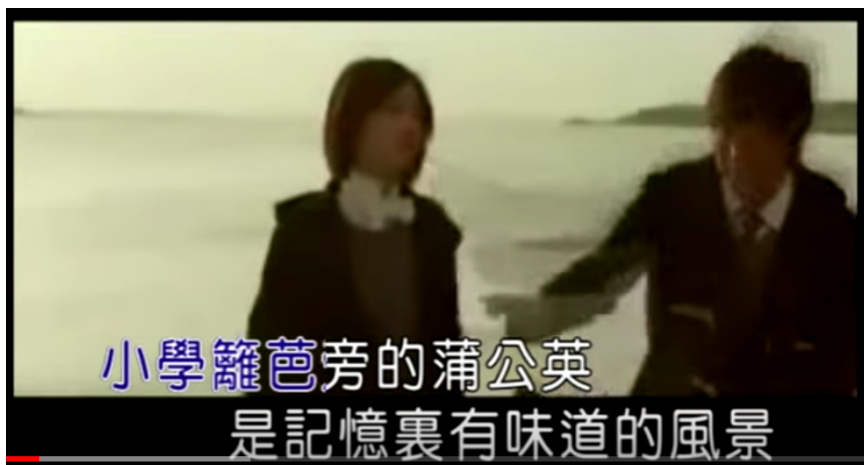
〔芭蕉〕ㄅㄚˊㄓㄠ

- 芭蕉
 - × 芭蕉
- 「芭蕉」是多年生草本亞熱帶植物；「芭」是指用竹子或柳條等編成的東西。所以「芭蕉」當用「芭」而非「芭」。

【篡】ㄘㄨㄢˋ

〔篡〕ㄘㄨㄢˋ 〔竄〕ㄘㄨㄢˋ

- 篡位
 - × 篡位 竄位
- 「篡」是奪取的意思，形近於「纂」；「纂」音ㄘㄨㄢˋ，是編輯的意思；「竄」是逃跑的意思。「篡位」是指古時臣下非分奪取王位，所以當用「篡」而非「纂」、「竄」。



圖片擷取自：<https://www.youtube.com/watch?v=DYyE68llx0U>



國教院索引典系統

Your query "開心" returned 3,402 matches in 2,528 different texts (in 68,809,790 words [40,960 texts]; frequency: 49.44 instances per million words) [0.163 seconds - retrieved from cache]

No	Filename	Solution 1 to 50	Page 1 / 69
1	56f4fa66326d1b19b79ecebb	，孩子們將會有許多困難，可能經常不開心，也可能因被留級而離開學校。相反的	
2	56f4fa66326d1b19b79ecebe	孩童高高舉起來或倒提舉高，心想孩童會開心，是一種寵愛的做法。一些對感覺	
3	56f4fa66326d1b19b79ecebe	地心引力獲得大量的感覺輸入，這些輸入會產生「開心」，對運動和情緒的發展非常重要。例如	
4	56f4fa66326d1b19b79ecbf	墊子，孩童在上面滾來滾去，常發出開心的笑聲，在小孩要匍匐、攀登、爬上、	
5	56f4fa66326d1b19b79ecf6b	原則71：在教室中如果有人大聲拿別人開心，以致干擾別人，教師可以讓這樣的孩子到	
6	56f4fa66326d1b19b79ecfa0	大人也能參與他們的活動，相信孩子會更開心。記得曾經有一位小孩扮演醫生，跑來向	
7	56f4fa66326d1b19b79ecfd2	再問他有何感覺，如何才能使小牛開心。要盡快採取行動，找出鬧彗扭的真正原因，	
8	56f4fa66326d1b19b79ecff2	跟別人鬧得面紅耳赤。結果，不但使自己不開心，而且會帶來挫敗的厄運。許多怨偶，彼此	
9	56f4fa66326d1b19b79ecff6	貧窮的人未必痛苦；社會地位高的人未必開心，市井小民未必不快樂。真正造成不快樂的原因	
10	56f4fa66326d1b19b79ed0c6	的自我。這一長串的說明，足以為開心生活畫出一幅色彩光鮮的圖譜。把佛學想像成	
11	56f4fa66326d1b19b79ed0c8	「心存感恩」，也都是佛家過著開心生活的修持法門，退一步，往往海闊天空。	
12	56f4fa66326d1b19b79ed0d2	朋友看待，把最喜歡的東西與人分享，開心的對待宇宙的萬事萬物，知道這一點，	
13	56f4fa66326d1b19b79ed0e4	不應侷限在醫療的定位，而是使人「開心」，這卻是宗教的本質，始於一種	
14	56f4fa66326d1b19b79ed0e5	用內省的方式，把她的關心找出來：「開心的活著」、「健康的生存」、「	

<http://coect.naer.edu.tw/cqpweb/>

Menu

- Corpus queries
- Standard query
- Restricted query
- Word lookup
- Frequency lists
- Keywords
- Analyse corpus
- Export corpus
- Saved query data
- Query history
- Saved queries
- Categorised queries
- Upload a query
- Create/edit subcorpora

中研院平衡語料庫

Standard Query

開心

Query mode: Simple query language syntax

Number of hits per page:

Restriction:

Your query "開心" returned 338 matches in 271 different texts (in 11,202,927 words [19,247 texts]; frequency: 30.17 instances per million words) [0.188 seconds - retrieved from cache]

< << >> > | Show Page: 1 Line View Show in random order New query Go!

No	Filename	Solution 1 to 50	Page 1 / 7
1	100548	工作可能是現在唸起來索然無味的電子、回控…如何	開心的起來？高一時就在疑惑自己出生是為
2	100622	的、白色的、黃色的，你一定會非常	開心…再比這兒好得多了…」李文秀緩緩了
3	100653	幾日後，在樓下的中庭相遇，她	開心極了，說：「還是新鮮的桂花才對呀
4	100654	工作，領文書的薪水。一開始，她很	開心，認為泡茶的工作簡單，又可以領文書的
5	100654	老是做著低三下四的泡茶工作，心裡很不	開心，不但端茶時表情鬱鬱，連泡出的茶
6	100662	讀者讀了歡喜，不喝茶的人讀了	開心，喝茶的人讀了也開心。不管學
7	100662	人讀了開心，喝茶的人讀了也	開心。不管學裡不學裡的人，都能
8	100663	適當的時機加以印證練習，這麼做難道不會覺得	開心嗎？」孔子在談到「學」，亦即做人
9	100664	時拿一毛錢買一枝冰棒，就很	開心了。上了大學之後，偶爾回去時都會
10	100732	問，在你小時候，吃飽、穿暖、玩得	開心之餘，想得最多的是什麼？答案之一
11	100764	所求，有人去探訪他，他都很	開心，不會抱著你們來看我，一定要帶
12	100789	所吸引，於是即興地舞動起來，跳得好	開心！跳完之後，觀賞的路人丟下了賞錢，黑人
13	100789	一切安頓得很好，根本不用我操心。我	開心極了，於是積極地往今後要走的路進行

詞彙基本搜尋

符號	意義	範例	符合之搜尋結果
?	任意的單一字	車?	車輛、車床、...
*	零~多字	*書	書、教書、保證書、...

如何查到：一板一眼、一心一意、...的例句？

如何查到：教學者、研究者、學者、...的例句？

詞彙基本搜尋 (詞性)

詞_詞性：搜尋符合特定詞性的詞彙

請比較下列差尋的差異性

工作

工作_VA

工作_Na

如何找到所有『臺○』的專有地名？

搜尋特定詞彙列表

Menu	中研院平衡語料庫	
Corpus queries	Word lookup	
Standard query	You can use this search to find out how many words matching the form you look up occur in the corpus, and the different tags that they have.	
Restricted query	Enter the word-form you want to look up	者 (NB. you can use the normal wild-cards of Simple Query language)
Word lookup	Show only words ...	<input type="radio"/> starting with <input checked="" type="radio"/> ending with ... the pattern you specified <input type="radio"/> containing <input type="radio"/> matching exactly
Frequency lists		List results by word-form, or by word-form AND tag? <input type="button" value="List by word-form and tag ↓"/>
Keywords		Number of items shown per page: <input type="text" value="50"/>
Analyse corpus		
Export corpus		
Saved query data		
Query history		
Saved queries		
Categorised queries		
Upload a query		
Create/edit subcorpora		

搜尋結果

Showing frequency breakdown of words in this query, at the query node; there are 265 different types and 30,810 tokens at this concordance position.

|< << >> >| Breakdown position: Node Frequency breakdown of words only Go!

No.	Search result	No. of occurrences	Percent
1	者	13695	44.45%
2	業者	3185	10.34%
3	或者	1821	5.91%
4	學者	1711	5.55%
5	記者	1539	5%
6	消費者	1423	4.62%
7	讀者	1120	3.64%
8	作者	845	2.74%
9	患者	778	2.53%

詞頻統計分析

Menu	中研院平衡語料庫
<ul style="list-style-type: none"> Standard query Restricted query Word lookup Frequency lists Keywords Analyse corpus Export corpus 	<p>Frequency lists</p> <p>You can view the frequency lists of the whole corpus and frequency lists for subcorpora you have created. Click here to create/view subcorpus frequency lists.</p> <p>View frequency list for ... <input type="text" value="Whole of 中研院平衡語料庫"/></p> <p>View a list based on ... <input type="text" value="Word forms"/></p> <p>Frequency list option settings</p> <p>Filter the list by <i>pattern</i> - show only words/tags ... <input type="text" value="starting with"/></p> <p>Filter the list by <i>frequency</i> - show only words/tags ... <input type="text" value="with frequency between"/> and <input type="text" value=""/></p>
<p>Saved query data</p> <ul style="list-style-type: none"> Query history 	

統計結果

Frequency list: **Word frequencies in entire “中研院平衡語料庫”, ending with “者”**

|< << >> New Frequency List ↓ Go!

No.	Word	Frequency
1	者	13,695
2	業者	3,185
3	或者	1,821
4	學者	1,711
5	記者	1,539
6	消費者	1,423
7	讀者	1,120
8	作者	845

搭配詞

“「搭配詞是某語言中字彙合併的方式，以產出自然的口語和文字」

(Collocation is the way words combine in a language to produce natural-sounding speech and writing) ”

Oxford Collocations Dictionary for Students of English (2002:vii)




1 Line View Show in random order Collocations... Go!

Solution 1 to 50 Page 1 / 12

不會。嫌她不懂得在美國生活，連開支票都不會。嫌她不會開車，什麼事都要
 著書包，她得了獎學金，拿著一張支票到銀行的櫃檯去領錢，櫃檯小姐問
 領錢，櫃檯小姐問她：「這張支票是你本人的嗎？」「是的，是
 的書包，戴著很深的眼鏡，拿了支票就跑到銀行去領錢。走到櫃台，向櫃台
 小姐說：「我要領錢。」然後把支票給小姐。小姐問：「這是你本人嗎
 冊第三課？」這裡的背書指的是在支票後面寫上你的名字，不是背課文的背書
 相信她甚於她父親。為了方便，大部分的支票都是以她的名字開出。婚後三年生
 在工作上不敬業且不負責任，像亂開空頭支票。7 缺乏道德、良知、冷漠無情，即使目睹
 選票，不加稅，但要加福利，選舉支票亂開，結果債留子孫。其實道德的問題又
 他賣出一批為數可觀的建材，對方以定期支票付款，於是陳姓商人要求對方另覓一位可靠
 館子裡吃了一餐。這件事一直到支票遭到退票時，才發現那張支票並沒有背書
 事一直到支票遭到退票時，才發現那張支票並沒有背書，他們只顧著高高興興的吃吃喝喝
 都不是，不過是一張未填日子的支票，看著它固教人心跳，但能不能憑
 恩恩愛愛，一齊兌了現，但又何必各拿支票一紙，等得那麼久耶？有些人

Choose settings for proximity-based collocations:

Include annotation: Part-of-speech tag Include Exclude

Maximum window span: + / - 

Collocation controls

Collocation based on: Statistic:

Collocation window from: Collocation window to:

Freq(node, collocate) at least: Freq(collocate) at least:

Filter results by: and/or tag:

Extra information: Log-likelihood scores collocations by significance: the higher the score, the more evidence you have that the association is not due to chance. More frequent words tend to get higher log-likelihood scores, because there is more evidence for such words.

There are 1,566 different words in your collocation database for "[word="支票"%c]". (Your query "支票" returned 566 matches in 346 different texts) (0.321 seconds - retrieved from cache)

No.	Word	Total no. in whole corpus	Expected collocate frequency	Observed collocate frequency	In no. of texts	Log-likelihood
1	空頭	236	0.012	51	36	766.336
2	張	20,788	1.026	103	88	749.114
3	開	13,611	0.672	44	39	282.052
4	美元	5,384	0.266	32	24	243.648
5	兌現	169	0.008	17	17	227.111
6	給	54,016	2.666	38	33	131.66
7	銀行	4,561	0.225	18	17	122.35
8	簽	589	0.029	12	9	120.872

Collocation based on: Statistic:

Collocation window from: Collocation window to:

Freq(node, collocate) at least: Freq(collocate) at least:

Filter results by: and/or tag:

No.	Word	Total no. in whole corpus	Expected collocate frequency	Observed collocate frequency	In no. of texts	Log-likelihood
1	開	13,611	0.56	49	43	342.382
2	兌現	169	0.007	15	15	201.562
3	簽	589	0.024	14	11	150.545
4	偽造	400	0.017	7	7	70.877
5	收到	2,614	0.108	9	9	61.968
6	簽署	303	0.013	5	3	50.031
7	收下	353	0.015	5	4	48.539
8	開出	470	0.019	5	5	45.671
9	拿到	2,626	0.108	7	7	44.652
10	拿	15,450	0.635	11	11	42.046
11	寫	25,311	1.041	12	10	36.803
12	使用	17,548	0.722	5	5	10.807

國教院華英雙語索引系統(試用版)

語料庫 **光華雜誌** 關鍵詞：

Sample Panel

- Even Belgium had installed a sinologist in a university before this . At Cambridge , Sir Thomas Wade , who had served as an interpreter and diplomat in China and is **famous for** the Wade-Giles system of romanization , became professor of Chinese in 1888 after donating a collection of over 650 books to the university .
就是他 和 繼任者 翟理斯 共同 創造了 Wade-Giles 羅馬 拼音法 。 總的來說 ， 英國 在 國富民強 的 帝國 巔峰 ， 也只 **提供** 過 五個 漢學 教授 席 而已 。
- Unlike the multiplicity of objectivity seen in the "Taiwan · Taiwan : Facing : Faces " exhibit in the Taiwan Pavilion , the exhibition " Segmentation / Multiplication " **provided** Taiwanese artists with opportunities to exhibit overseas and , through the careful planning of the curators , also **provided** the art lovers from around the world who attended the Biennale de Venezia a chance to understand another side of contemporary art in Taiwan .
不同於 台灣 國家 館 「 面 · 目 · 全 · 非 」 客觀 多元 的 呈現 ， 另一個 由 策展人 以 主觀 美學 所 呈現 的 「 裂合 與 聚生 」 展出 ， 不僅 **提供** 藝術家 有 更多 國際 展出 的 機會 ， 也 藉由 策展人 的 經營 ， **提供** 前往 威尼斯 朝聖 的 各國 藝術 愛好者 ， 瞭解 台灣 當代 藝術 的 另一 種 面目 。
- The Commission annually **provides** any where from US \$ 200-2000 to the majority of these schools in financial assistance . The Commission also **provides** textbooks free of charge to those who need them .

Translation Panel

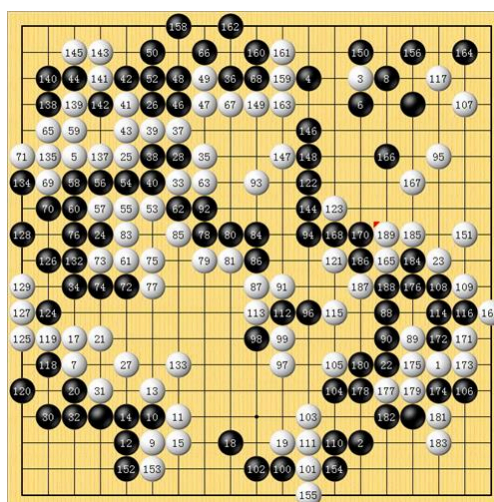
provide (531)
providing (185)
provided (184)
provides (164)
offer (128)
courtesy (107)
offering (60)
for (59)
offers (51)

Collocation Panel

提供 [V]
[C-]
並 ~,
[N-]
政府 ~,
[P-]
為 ~,
[Vi-]
免費 ~.

<http://coct.naer.edu.tw/bc/>

電腦進步到能
下圍棋，
結果它還是只
會算頻率跟當
搜尋引擎嗎？



圖片來源：<http://file.xdao.com/2016/03/20/f3d7c7ff-d585-41b7-befa-88f1e005cf70.jpg>

『依此類推』

蘋果

圖源: nlp1.com/


<http://a3.att.hudong.com/47/87/20300000636999135190876699345.jpg>

國教院近義詞觀察系統

Word2Vec 近義詞觀察系統

選擇語料庫	比較詞(組)相似度	查詢詞(組)近義詞
<input checked="" type="radio"/> Chinese Giga Words <input type="radio"/> 中國時報 <input type="radio"/> 中國時報(詞性) <input type="radio"/> 遠流語料 <input type="radio"/> 遠流語料(詞性) <input type="radio"/> 國語日報 <input type="radio"/> 平衡語料庫 <input type="radio"/> 平衡語料庫(詞性)	輸入 正相關詞(組): 蘋果 負相關詞(組): 請輸入負相關詞, 例如: 櫻桃 輸出詞數: 50 <input type="button" value="送出"/>	輸出 近義詞: 蘋果 0.5129565000534058 櫻桃 0.5031850934028625 柑橘 0.48714378476142883 梨子 0.48301517963409424 橘子 0.4727131128311157 香蕉 0.46589866280555725 荔枝 0.4639115631580353 葡萄 0.46065250039100647 奇異果 0.441455602645874 馬鈴薯 0.44132381677627563 雞肉 0.4407174289226532 水果 0.43828028440475464 戴爾 0.4361230432987213 牛肉 0.43600136041641235 昇陽 0.43259522318840027

<http://coct.naer.edu.tw/word2vec/>

國教院近義詞觀察系統

比較詞(組)相似度 查詢詞(組)近義詞

輸入

正相關詞(組):

蘋果

負相關詞(組):

牛肉

輸出詞數: 50

送出

輸出

近義詞:

宏碁 0.37290841341018677
昇陽 0.3706173598766327
精技 0.346422016620636
研宇 0.3399192988872528
郭士納 0.3358774483203888
康柏 0.33453407883644104
老貓 0.3339369595050812
Thiz 0.33262884616851807
IBM 0.32719194889068604
劉門 0.3197169005870819
interactive 0.3182205259799957
筆記簿型 0.3147069215774536
星島日 0.3138140141963959

國教院近義詞觀察系統

比較詞(組)相似度 查詢詞(組)近義詞

輸入

正相關詞(組):

蘋果 草莓

負相關詞(組):

牛肉

輸出詞數: 50

送出

輸出

近義詞:

梨子 0.5411585569381714
櫻桃 0.5038571357727051
蓮霧 0.5021500587463379
西瓜 0.5018879175186157
椪柑 0.49865397810935974
鳳梨 0.49444371461868286
枇杷 0.4891597032546997
茂谷柑 0.484711856499481
水蜜桃 0.48459604382514954
橘子 0.46837592124938965
葡萄柚 0.46056604385375977
高接梨 0.4559977948665619
酪梨 0.45508047938346863
火龍果 0.4546027183532715
芭樂 0.4524976313114166



索引典系統進階語法





Case Study

同中求異，異中求同



	解決	處理
問題	27491	5186
困境	971	69
爭議	926	349
困擾	710	119
難題	675	54
危機	624	560
水患	606	20
紛爭	428	64
困難	408	86

	解決	處理
問題	27491	5186
垃圾	111	2020
廢棄物	22	1756
案	167	1361
事	405	1073
事件	105	1072
事宜	69	1054
案件	60	1018
事務	40	1015

	燃眉之急	314	
	淹水	283	
	瓶頸	177	
	亂象	147	
	打者	102	
	不便	61	
	懸案	61	
	疑難	53	
	血荒	51	
	對手	47	
	後事	446	
公務	193		
屍體	151		
遺體	105		
郵件	74		
油污	67		
傷口	65		