

# 大數與求真

## 如何由十億字語料中 學習語言知識

黃居仁

中央研究院語言學研究所

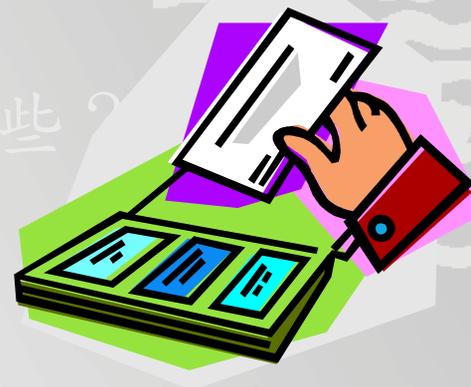
<http://cwn.ling.sinica.edu.tw/huang/huang.htm>

# 楔子：語言知識的規律化描述

## 「發票」怎麼用？

• 線上同義詞搜尋

「美白」相關詞彙有哪些？





Home **Concordance** Word Sketch Thesaurus Sketch-Diff

### Word Sketch Entry Form

Corpus: chinese\_all\_trd

Word Form: 發票 ←

Sort grammatical relations:

Minimum frequency: 5

Minimum salience: 0.0

Maximum number of items in a grammatical relation: 25

Show Word Sketch



Home Concordance **Word Sketch** Thesaurus Sketch-Diff

# 發票

chinese\_all\_trd:taiwan-only freq = 8408 [change options](#)

object_of 928 5.2	subject_of 692 0.3	a_modifier 276 0.5	n_modifier 562 -12.3	modifies 607 -12.8
開立 <a href="#">253</a> 73.78	收執 <a href="#">15</a> 43.41	不實 <a href="#">85</a> 57.41	增值稅 <a href="#">103</a> 47.76	逃漏稅 <a href="#">16</a> 30.08
虛開 <a href="#">50</a> 51.4	給 <a href="#">69</a> 32.61	假 <a href="#">48</a> 46.83	銷貨 <a href="#">17</a> 34.96	存根聯 <a href="#">7</a> 29.53
偽造 <a href="#">70</a> 40.88	對獎 <a href="#">7</a> 29.59	空白 <a href="#">24</a> 41.9	收銀機 <a href="#">16</a> 34.19	金額 <a href="#">55</a> 27.38
虛購 <a href="#">12</a> 40.47	中獎 <a href="#">12</a> 25.43	中獎 <a href="#">22</a> 39.32	式 <a href="#">32</a> 30.07	面額 <a href="#">13</a> 25.7
漏開 <a href="#">16</a> 40.07	逃漏 <a href="#">9</a> 24.03	普通 <a href="#">19</a> 28.32	小額 <a href="#">23</a> 29.64	日期 <a href="#">18</a> 21.74
開具 <a href="#">30</a> 39.48	充當 <a href="#">9</a> 23.46	領用 <a href="#">5</a> 27.62	進項 <a href="#">9</a> 27.8	獎金 <a href="#">18</a> 20.77
變造 <a href="#">14</a> 27.67	兌換 <a href="#">16</a> 23.27	可疑 <a href="#">5</a> 15.34	聯 <a href="#">29</a> 27.75	偷稅 <a href="#">7</a> 20.53
虛設 <a href="#">7</a> 22.83	換版工作 <a href="#">3</a> 22.92	增值 <a href="#">3</a> 12.82	虛立 <a href="#">4</a> 26.96	助創世 <a href="#">2</a> 19.68
使用 <a href="#">49</a> 22.27	捐贈 <a href="#">13</a> 22.15	原始 <a href="#">3</a> 11.24	加副聯 <a href="#">3</a> 22.14	號碼 <a href="#">12</a> 19.14
開出 <a href="#">15</a> 21.51	換好 <a href="#">4</a> 21.78	填開式 <a href="#">1</a> 11.17	愛心 <a href="#">18</a> 21.81	影本 <a href="#">7</a> 18.24
購買 <a href="#">24</a> 20.21	抬頭 <a href="#">7</a> 20.97	免用 <a href="#">1</a> 10.42	盜竊 <a href="#">11</a> 21.63	案件 <a href="#">22</a> 17.59
取得 <a href="#">34</a> 19.33	犯罪 <a href="#">21</a> 20.69	全額 <a href="#">2</a> 10.38	票載 <a href="#">3</a> 20.76	憑證 <a href="#">8</a> 17.17
募集 <a href="#">9</a> 17.95	膨脹 <a href="#">7</a> 19.8	小小 <a href="#">2</a> 10.3	六獎 <a href="#">3</a> 20.36	人因 <a href="#">8</a> 17.02
印製 <a href="#">7</a> 16.2	傳情 <a href="#">4</a> 19.13	正規 <a href="#">2</a> 10.2	開假 <a href="#">3</a> 20.36	管理員 <a href="#">7</a> 16.76
持 <a href="#">13</a> 15.4	冒領 <a href="#">5</a> 18.31	原 <a href="#">5</a> 9.97	電腦版 <a href="#">3</a> 19.02	收據 <a href="#">5</a> 15.47
假造 <a href="#">3</a> 14.93	盜領 <a href="#">5</a> 18.16	欣榮 <a href="#">1</a> 8.77	票券 <a href="#">11</a> 17.95	魔方 <a href="#">2</a> 15.45
發出去 <a href="#">2</a> 14.13	捐給 <a href="#">5</a> 17.81	作廢 <a href="#">1</a> 7.66	開具假 <a href="#">2</a> 17.84	普獎 <a href="#">2</a> 14.52
買 <a href="#">9</a> 13.66	開立 <a href="#">6</a> 15.83	足額 <a href="#">1</a> 7.0	預前 <a href="#">2</a> 17.84	婚紗秀 <a href="#">2</a> 14.07

Corpus: chinese\_all\_trd:taiwan-only  
Hits: 253  
[conc description](#)

- [CNA19910122.0254](#)
- [CNA19910122.0254](#)
- [CNA19910125.0259](#)
- [CNA19910128.0108](#)
- [CNA19910204.0191](#)
- [CNA19910403.0246](#)
- [CNA19910403.0246](#)
- [CNA19910617.0116](#)
- [CNA19910717.0256](#)
- [CNA19910719.0096](#)
- [CNA19910719.0096](#)
- [CNA19910724.0128](#)
- [CNA19910724.0128](#)
- [CNA19910724.0128](#)
- [CNA19910724.0220](#)
- [CNA19910724.0220](#)
- [CNA19910731.0239](#)

將是一大損失，商家最好誠實開立**發票**，以免因小失大。<p><p>財政部自十五日  
 <p><p>凡遭查到三次以上未誠實開立**發票**的商家，將遭停業處置，王建(火宣  
 進入各類商店購物，如當場未開立**發票**，即以「現行犯」成為處罰列管對象  
 趁機提高貨品價格，卻又沒有據實開立**發票**，可隨時向稅捐處檢舉。依據營業稅法  
 近，公司行號銷貨時務必主動開立**發票**，若漏開統一發票達三次者，無論  
 引起消費大眾認同，進而促使商店開立**發票**，建立良好納稅風氣。<p><p>在蘇法昭  
 捕頭」專線連絡，以對三次不開立**發票**的商店，處以停業處分。<p><p>至於  
 逃漏稅資料，在銷售車輛時所開立的**發票**，每輛車短開新台幣十萬元以上，  
 消費者未索取的統一發票，或未主動開立**發票**而於事後虛開大量小額發票，以冒領  
 : <p><p>- 經營買賣業平均每月開立**發票**申報銷售額未達二十萬元，應進行  
 清查列管，並輔導營業人依法誠實開立**發票**。<p><p>- 選擇重點行業進行深入查核  
 今天重申，所屬加油站員工應誠實開立**發票**，否則將移送法辦。<p><p>中油在高雄  
 稍早，中油公司所屬加油站因沒有開立**發票**，無法與民營加油站競爭，加油站營運  
 ，經台灣石油工會爭取後才奉准開立**發票**，不料又發生員工違法損及公司形象  
 中獎機會，要求營業人按件逐筆開立**發票**，使業者增加開立發票的手續，購用  
 按件逐筆開立發票，使業者增加開立**發票**的手續，購用發票成本也因而提高  
 的執行重點，初步選定平均每月開立**發票**申報銷售額未達新台幣二十萬元的營業人

[expand left](#)  
 )今天表示，凡遭查到漏開發票三次的營業單位，必將遭停業處分，在此舊曆春節之際，對於商家將是一大損失，商家  
 最好誠實開立**發票**，以免因小失大。財政部自十五日展開全省查緝漏開發票行動，王建(火宣)今天並特別邀集全省稅捐稽徵  
 基層單位人員及主管，進行一整天的查緝行動  
[expand right](#)

Home Concordance **Word Sketch** Thesaurus Sketch-Diff

# 發票

chinese\_all\_trd:taiwan-only freq = 8408

[change options](#)

object_of 928 5.2	subject_of 692 0.3	a_modifier 276 12.5	n_modifier 562 -12.3	modifies 607 -12.8
開立 <a href="#">253</a> 73.78	收執 <a href="#">15</a> 43.41	不實 <a href="#">85</a> 57.41	增值稅 <a href="#">103</a> 47.76	逃漏稅 <a href="#">16</a> 30.08
虛開 <a href="#">50</a> 51.4	給 <a href="#">69</a> 32.61	假 <a href="#">48</a> 46.83	銷貨 <a href="#">17</a> 34.96	存根聯 <a href="#">7</a> 29.53
偽造 <a href="#">70</a> 40.88	對獎 <a href="#">7</a> 29.59	空白 <a href="#">24</a> 41.9	收銀機 <a href="#">16</a> 34.19	金額 <a href="#">55</a> 27.38
虛購 <a href="#">12</a> 40.47	中獎 <a href="#">12</a> 25.43	中獎 <a href="#">22</a> 39.32	式 <a href="#">32</a> 30.07	面額 <a href="#">13</a> 25.7
漏開 <a href="#">16</a> 40.07	逃漏 <a href="#">9</a> 24.03	普通 <a href="#">19</a> 28.32	小額 <a href="#">23</a> 29.64	日期 <a href="#">18</a> 21.74
開具 <a href="#">30</a> 39.48	充當 <a href="#">9</a> 23.46	領用 <a href="#">5</a> 27.62	進項 <a href="#">9</a> 27.8	獎金 <a href="#">18</a> 20.77
變造 <a href="#">14</a> 27.67	兌換 <a href="#">16</a> 23.27	可疑 <a href="#">5</a> 15.34	聯 <a href="#">29</a> 27.75	偷稅 <a href="#">7</a> 20.53
虛設 <a href="#">7</a> 22.83	換版工作 <a href="#">3</a> 22.92	增值 <a href="#">3</a> 12.82	虛立 <a href="#">4</a> 26.96	助創世 <a href="#">2</a> 19.68
使用 <a href="#">49</a> 22.27	捐贈 <a href="#">13</a> 22.15	原始 <a href="#">3</a> 11.24	加副聯 <a href="#">3</a> 22.14	號碼 <a href="#">12</a> 19.14
開出 <a href="#">15</a> 21.51	換好 <a href="#">4</a> 21.78	填開式 <a href="#">1</a> 11.17	愛心 <a href="#">18</a> 21.81	影本 <a href="#">7</a> 18.24
購買 <a href="#">24</a> 20.21	抬頭 <a href="#">7</a> 20.97	免用 <a href="#">1</a> 10.42	盜竊 <a href="#">11</a> 21.63	案件 <a href="#">22</a> 17.59
取得 <a href="#">34</a> 19.33	犯罪 <a href="#">21</a> 20.69	全額 <a href="#">2</a> 10.38	票載 <a href="#">3</a> 20.76	憑證 <a href="#">8</a> 17.17
募集 <a href="#">9</a> 17.95	膨脹 <a href="#">7</a> 19.8	小小 <a href="#">2</a> 10.3	六獎 <a href="#">3</a> 20.36	人因 <a href="#">8</a> 17.02
印製 <a href="#">7</a> 16.2	傳情 <a href="#">4</a> 19.13	正規 <a href="#">2</a> 10.2	開假 <a href="#">3</a> 20.36	管理員 <a href="#">7</a> 16.76
持 <a href="#">13</a> 15.4	冒領 <a href="#">5</a> 18.31	原 <a href="#">5</a> 9.97	電腦版 <a href="#">3</a> 19.02	收據 <a href="#">5</a> 15.47
假造 <a href="#">3</a> 14.93	盜領 <a href="#">5</a> 18.16	欣榮 <a href="#">1</a> 8.77	票券 <a href="#">11</a> 17.95	魔方 <a href="#">2</a> 15.45
發出去 <a href="#">2</a> 14.13	捐給 <a href="#">5</a> 17.81	作廢 <a href="#">1</a> 7.66	開具假 <a href="#">2</a> 17.84	普獎 <a href="#">2</a> 14.52
買 <a href="#">9</a> 13.66	開立 <a href="#">6</a> 15.83	足額 <a href="#">1</a> 7.0	預前 <a href="#">2</a> 17.84	婚紗秀 <a href="#">2</a> 14.07



Home Concordance Word Sketch Thesaurus Sketch-Diff **Frequency** Collocation

KWIC/Sentence View options **Sample** Filter Sort Save

Corpus: chinese\_all\_trd:taiwan-only  
Hits: 85  
[conc description](#)

Page 1 of 5 Go Next | Last

- [CNA19911103.0009](#)
- [CNA19911119.0281](#)
- [CNA19930224.0091](#)
- [CNA19930817.0082](#)
- [CNA19931129.0190](#)
- [CNA19931130.0329](#)
- [CNA19940302.0282](#)
- [CNA19940302.0282](#)
- [CNA19950324.0326](#)
- [CNA19950324.0326](#)
- [CNA19950324.0326](#)
- [CNA19950327.0266](#)
- [CNA19951228.0287](#)
- [CNA19960312.0210](#)
- [CNA19960510.0176](#)
- [CNA19960626.0318](#)

申報案件抽查作業、營業人開立不實**發票**的查緝作業、進口大宗或高價值貨物顯示，有多起虛設行號開立的不實**發票**，由於案情重大且牽連甚廣，於是已達新台幣二十多億元。而開立不實**發票**的面額也達到四億五千四百三十六萬元和廠商協議開發**發票**，取得不實面額**發票**，向工業局請領公款，而以此種台灣再生公司涉嫌偽造文書，以不實**發票**向汽水公會領取寶特瓶清除處理費和，外傳該公司回收廢寶特瓶有以不實**發票**領取費用情形，是完全不了解回收身分證，申請設立公司行號並開出不實**發票**，金額高達新台幣五十餘億元。

台各地串連，形成供應和販賣不實**發票**的銷售網，開出發票的金額高達五十多億四百多萬元及十六萬多元面額的不實**發票**給民間業者，也被檢方併案偵辦。蔡盛財等四名中油員工和一名使用不實**發票**的民間業者范植森等五人收押禁見表示，這六名中油員工是否交付不實**發票**而取得民間業者的期約或收受賄賂120件為涉嫌繳驗偽造、變造或不實**發票**虛報完稅價格逃漏進口稅捐案件，其另被告等製造假買賣，以不實**發票**、支票向銀行融資，這八十四年五月公司財物經理侯秋香，相互開立不實**發票**向金融行庫辦理貼現，經檢方偵辦買賣**發票**逃漏稅捐，以往廠商開立不實**發票**都是採分散方式，人工查核不易，調查單位查獲十八家虛設行號開立不實**發票**、營業額十五億多元；徵收積欠稅收

expand left  
expand right

，財政部指出，十八個項目中，最重要的防止營業人漏開統一發票稽查作業、繁華地段房屋租金偏低案件查核作業、簽證不實的營所稅申報案件抽查作業、營業人開立不實**發票**的查緝作業、進口大宗或高價值貨物銷售流程追查作業、扣繳異常單位檢查作業、營利事業取得小店戶收據查核作業等七項，都由財政部列管，稽徵機關

# 前言

## 數大就是美

— 徐志摩

- 語料與大數
  - 語言學家的任務就是要在無窮盡的語言資料中找出規律與解釋的真理來。
  - 語料庫語言學研究的挑戰之一，是當語料庫不夠大時，有些語言現象無法呈現，或呈現不出全貌。
  - 但是語料量太大時，人又無法有效消化
- 從大數中求真
  - 語料量愈大，愈較能呈現語言使用的完整面貌
  - 更重要的是能呈現詞語與詞語間的關係

# 語料量多大才算夠大？

## 關鍵詞詞彙的觀察

- 第10,000名常用詞的詞頻大約在0.001%左右
  - 在五百萬詞語料庫中出現大約50次
- 
- 足夠觀察到一個詞彙使用的大致狀況
  - 也正好是人們可以從容處理語料的規模
  - 五百萬到一千萬詞，是關鍵詞檢索語料（KWIC）觀察的最佳規模



# 古典語料庫： 數量級在百萬的中文語料庫

Corpus Name	Online Year	Data	Duration/ Content
<b>Sinica 4.0 (Taiwan)</b>	<b>1996</b>	<b>5.2 M words 7.9 M characters</b>	<b>1990-1996 Fully Tagged</b>
<b>Sinica 5.0 (Taiwan)</b>	<b>2006</b>	<b>10 M words</b>	<b>1990-2004 Fully Tagged</b>
<b>Sinorama (Taiwan)</b>	<b>2003</b>	<b>3.2 M English words 5.3 M Chinese characters</b>	<b>1976 – 2000 (1999-2000) Aligned</b>
<b>CCL (Peking)</b>	<b>2003</b>	<b>85 M simplified characters</b>	<b>1919 -2003 Partially tagged (1 million)</b>

**M= million**

# 古典語料庫中如何將知識規律化

## The coloured pens method

1 arity, which will be used to take a party of under-privileged children to  
2 from outside. You are invited to a party and after a couple of drinks you  
3 tion, we believe politicians of all parties will listen to our views. &eq  
4 ould be reaching agreement with all parties concerned, as to which event  
5 lack people. I have certainly been party to one or two discussions amongs  
6 . These should be discussed by both parties before entering into the relat  
7 presents They had hosted a cocktail party at Kensington palace, for examp  
8 akes. By midnight the end-of-course party is in full swing, but most cad  
9 e should be a right for the injured party to terminate the contract. A ma  
10 by the Safran Peoples ' Liberation Party. This presents the powerful nei  
11 s. Ahead I could see the rest of my party plodding towards the final slope  
12 cial ethic. The two main political parties - the Tories and the Liberals  
13 ritish successes in Perth The small party of British players competing in  
14 to help control. One member of the party went to summon the rescue team a  
15 rket society fashion magazine. The party was held at his flat which was a  
16 security and secrecy than any Tory Party Conference : it seems that bootl

1 political association

2 social event

3 group of people

4 person in an agreement/dispute

5 to be party to something...

from Kilgarriff et al. 2005

# 自動觀察預測語言規律 所需要的語料量預估

- 以單一詞彙為中心的詞彙規律
- 呈現詞與詞之間的共現關係，如
  - V+N：「開立」+「發票」
  - A+N：「不實」+「發票」
- 以第一萬名中頻詞為切分點，希望能看到以這些詞為標的語法規律。
- 在關鍵詞前後**5**個詞內出現（比緊鄰多**10**倍機率）
- 估計需要十億詞的語料，在機率上才能與其看到資料

語言

# 當代鉅量語料庫的挑戰

鉅量：一億詞（100,000,000）以上，  
到十億詞（1,000,000,000 giga）的規模

- 第一萬名中頻詞約有1,000-10,000個例句
  - 足夠做自動預測
  - 並提供有效樣本作分佈分析，**但是**  
人無法有效閱讀
    - 50個例句: 輕輕鬆鬆
    - 500個例句: 得花點時間與力氣
    - 5000個例句: 無法集中精神，頭昏眼花
- 更難想像由人來標記/檢查十億個詞的詞類

# 處理鉅量語料庫的解決方案

- 採取機器自動分詞與詞類標記
  - 前提：計算語言學研究成績進步，有95%以上的準確率

*Ma and Huang 2006 LREC*

- 鉅量語料庫的處理模式
  - 必須以自動機器標記為主
  - 語言知識必須以自動方式取得並由機器處理
  - 人工方式評估取樣結果
  - 回到自動標記系統中改進

# 由鉅量語料跨越到語言規律

- 利用語法結構知識自動抽取詞與詞間的關係訊息，並進行統計
  - 跨越人類記憶力與注意力集中的限制
  - 由詞彙例句為本，轉為語法語意關係的分佈趨勢為本
  - 幫電腦加上初步歸納規律的能力
  - 觀察規律與趨勢，有大量資料時；少量的例外與誤差（如**5%**以下）不影響大局

# 以觀察人爲比喻

- 可以觀察其身高，體重，體重，衣著。◦◦◦
  - 還是不太瞭解這個人是什麼樣的人
- 觀察他的朋友，親人，他如何吃飯，走路，交談，如何處理人際關係。◦◦◦
  - 原來他是這樣的人！
- 詞彙也是如此

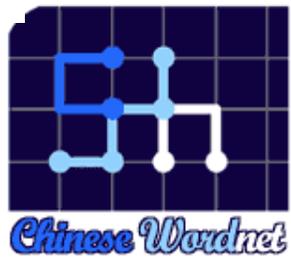
You shall know a word by the company it keeps. -M. K. Halliday

# 新一代的語料庫處理工具

- 直接呈現共現句型 (collocation)
- 以統計工具計算並共現句型的顯著性
  - 觀察處理約20-40個共現句型
  - 每個句型的共現詞(collocates)，都已按顯著性排列
    - 分析由幾千具中歸納出的規律性模式，而不是由幾千句素語料開始分析

# 由例句到詞彙關係：以BNC為例

- British National Corpus (BNC)
  - 一億 (100,000,000) 詞 POS-標注完成
- lemmatized
  - *assisting* => *assist* (v)
- Parsed (部分剖析) 得到
- 七千萬 (70,000,000) 組關係
  - <object, sip, coffee>    <subject, arrive, coffee>
  - <and-or, tea, coffee>    <modifier, coffee, instant>



# 語言結構知識的應用：BNC爲例

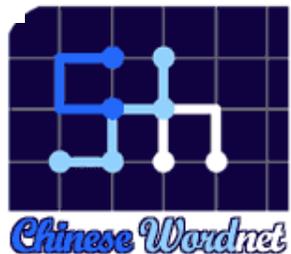
如何定義詞與詞間的關係？

- Uses CQL (Corpus query language)
  - Christ and Schulze, U. Stuttgart, 1994
- defining an object:

$v (adj | n | det | num | adv)^* n$

rewriting in CQL with BNC/CLAWS-5 tags

`[tag="VV.*"] [tag="(A[JTV]|D|O).*"]* [tag="NN.*"]`



# 方案平台：WordSketch Engine

- WordSketch Engine <http://sketchengine.co.uk>  
Adam Kilgarriff 的研究團隊開發
- 中文詞彙特性素描系統 Chinese WordSketch (CWS) 的開發
  - **Academia Sinica, Taiwan** (Huang, Smith, Ma, Simon 黃居仁，史尙明，馬偉雲，石穆)
  - **Masaryk University, Czech** (Rychly)
  - **Technical University, Budapest** (Tugwell)
  - **Lexical Computing Ltd** (Kilgarriff)

# CWS 所採用的語料庫

十四億字的Chinese Giga-word Corpus)  
含兩岸三地語料，當今最大中文語料庫，無詞類標記  
➤ CGW Corpus (LDC 2005) 的內容分布

CGW Corpus	K-characters	Documents
CNA(Taiwan)	792,195	1,769,953
Xinhua(PRC)	471,110	992,261
Zaobao(Sing)	28,066	41,418
TOTAL	1,291,371	2,803,632

# CGW 語料庫基本資料

	來源	字數	詞數	文件數
<b>First Edition</b>	<b>CNA</b>	735	462	1,649
	<b>Xinhua</b>	382	252	817
	<b>TOTAL</b>	1,118	714	2,466
<b>Second Edition</b>	<b>CNA</b>	792	497	1,769
	<b>Xinhua</b>	471	310	992
	<b>Zaobao</b>	28	18	41
	<b>TOTAL</b>	1,291	825	2,803

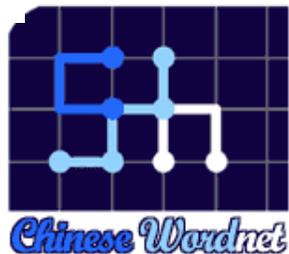
Unit: **Million**

# CWS 所用的語法知識

- 詞義語與法共現的自動抽取需要大量語言知識與理論的累積
  - 1: "V[BCJ]" "Di"? "N[abc]"? "DE"? "N[abc]"?
  - 2: "Na" [tag!= "Na"]
- ("XXX" represents XXX is a regular expression, "XXX"? represents XXX appears zero or one time, "XXX"{a,b} represents XXX appears a~b times.)

# 整合語料標記，統計，與語法知識

- Collocations are identified with Context free rules in Word Sketch Engine
  - Collocating Pattern for Object from CSE I
  - 1: "V[BCJ]" "Di"? "N[abc]"? "DE"? "N[abc]"? 2: "Na" [tag!= "Na"]
- Challenge: Long-distance relations
- 全穀麵包，吃了很健康。  
*quan.gu mian.bao, chi le hen jian.kang*
- 有人嘗試要將這荷花分類，卻越分越累。  
*you ren chang.shi yao jiang zhe he.hua fen.lei,  
que yue fen yue lei*



# 在處理平台中引進 語法知識 – 來源

- Information-based Case Grammar (ICG, Chen and Huang 1992)
- Encoded on over 40,000 verbs in Sinica Lexicon
  - ICG Basic Patterns for Stative Pseudo-transitive Verb (VI)

EXPERIENCER<GOAL[PP[對]]<VI

EXPERIENCER<VI<<GOAL[PP[於]]

THEME<GOAL[PP{對、以}]<VI

THEME<VI<<GOAL[PP[於]]

THEME<VI<<SOURCE[PP{自、於}]

THEME< SOURCE[PP{歸、爲}]<VI

# 在處理平台中引進 語法知識 – 應用

- 村莊(object) 明天將 被 夷為平地(VB11)  
*cunzhuang mingtian jiang bei yiweipingdi*  
– begin time1 location time1 adv? passive\_prep  
adv\_string 1:"V[BCJ].\*" [tag!="DE"]
- 大量的 遊客 破壞(VC2) 公園 景觀(object)  
*daliang de youke pohuai gongyuan jingguan*  
– 1:"VC.\*" (particle|prep)? NP not\_noun  
– (NP is defined as "...noun\_modifier{0,2}  
2:noun...".

# 在處理平台中引進語法知識 – 結果

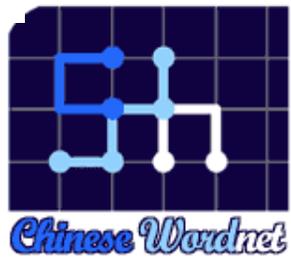
- Object Recall Comparison

	CSE I	CSE II
<i>hong2 (red)</i>	0	0
<i>pao3 (run)</i>	0	8,704
<i>kan4 (look)</i>	32,350	64,096
<i>da3 (hit)</i>	26,016	47,182
<i>song4 (give)</i>	0	76,378
<i>shuo1 (say)</i>	0	20,350
<i>xiang1xin4 (believe)</i>	0	52,373
<i>quan4 (persuade)</i>	0	3,852

語音

# 在處理平台中引進語法知識 – 結果 II

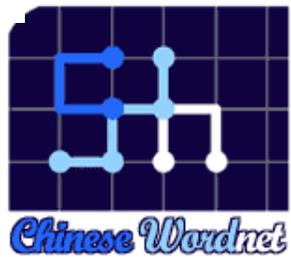
- **CSEII** 中「吃」的最顯著賓語  
**CSEII** 前**20**名，但不在**CSEI**前**20**名中的
- 飯 *fan4 rice* 802 70.96 (4),
- 虧 *kui disadvantage* 329 59.24 (12)
- 苦頭 *ku3tou2 suffering* 194 58.71 (14)



# 中文詞彙素描/CGW語料庫重要 資料統計(中央社語料)

- 規模：約5億詞（約497,000,000）
- 斷詞後詞數：455,526,796
- 使用語法規律數：約200 條
- 詞彙關係種類：20 個
- 找出的詞彙關係數59,183,238

語音



# 中文詞彙特性速描系統介紹

- Concordance
- WordSketch
- Sketch Difference
- Thesaurus
- 如何申請使用

語音

# Concordance

- 比傳統KWIC更強的功能

語彙



網址(D) http://corpora.fi.muni.cz/chinese\_all/

Google G Bookmarks PageRank 49 blocked Check AutoLink AutoFill Send to Settings

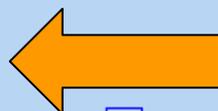
Home **Concordance** Word Sketch Thesaurus Sketch-Diff

Corpus: chinese\_all\_trd

Make Concordance

### Keyword(s)

Phrase: 發票



Word Form: PoS: unspecified Match case:

CQL:

Default attribute: word [Tagset summary](#)

### Context

Query Type: All of these items.

Left context

Right context

Window Size: 5 tokens. 5 tokens.

Word Form:

PoS: noun verb adjective adverb

noun verb adjective adverb

Make Concordance

Home **Concordance** Word Sketch Thesaurus Sketch-Diff Frequency Collocation

KWIC/Sentence View options Sample Filter Sort [List Icon] [List Icon]

121 Go Next | Last

Corpus: chinese\_all\_trd  
Hits: 8408  
[conc description](#)

- [CNA19910102.0102](#)
- [CNA19910102.0102](#)
- [CNA19910102.0102](#)
- [CNA19910102.0102](#)
- [CNA19910102.0102](#)
- [CNA19910102.0102](#)
- [CNA19910104.0131](#)
- [CNA19910104.0131](#)
- [CNA19910108.0104](#)
- [CNA19910122.0254](#)
- [CNA19910122.0254](#)
- [CNA19910122.0254](#)
- [CNA19910125.0259](#)
- [CNA19910125.0259](#)
- [CNA19910127.0102](#)
- [CNA19910127.0102](#)
- [CNA19910127.0102](#)

民眾可蒐集金額新台幣二千元以上的**發票**，向稅捐單位領取紀念品。

活動，除可繼續提醒民眾保持購物索取**發票**的習慣外，更配合年度所得稅的申報宣導將舉辦的活動包括：

- 集**發票**兌換紀念品，凡集滿本月份發票金額
- 集**發票**兌換紀念品，凡集滿本月份**發票**金額在二千元以上者，可憑發票於
- 月份發票金額在二千元以上者，可憑**發票**於本月三十日及三十一日向各稅捐處或摸彩活動。
- 辦理餐飲業開立**發票**模範商店選拔。

八名減刑獲釋秦人勒令停業一週。

稅捐處已將漏開**發票**與大額欠稅戶以電腦列管，一年內欠稅戶以電腦列管，一年內有三次漏開**發票**紀錄及欠稅超過十萬元，都將依規定八十年元月份二聯式統一發票或收銀機**發票**收執聯金額達二千元以上者，均可將是一大損失，商家最好誠實開立**發票**，以免因小失大。

財政部自十五日

凡遭查到三次以上未誠實開立**發票**的商家，將遭停業處置，王建(火宣一百萬者，在誠實開發票前提下，使**發票**數額提高至兩百萬元以上，財政部的進入各類商店購物，如當場未開立**發票**，即以「現行犯」成為處罰列管對象包含了自己所繳的稅額，如不索取**發票**，就等於被不肖商家吃掉，為維護打算再增加四人。

上海查獲非法**發票**交易(中央社台北二十七日電)

香港報導，上海工商執法人員查獲一批非法**發票**交易市場，當場擒獲多名票販。

用三元「人民幣」買來一張空白非法**發票**，任意填寫金額便可回去報假賬；而

expand left

台省各縣市舉辦統一發票宣導活動(中央社中興新村二日電)本月份購物時別忘了索取統一發票。臺灣省各縣市正舉辦統一發票宣導，民眾可蒐集金額新台幣二千元以上的**發票**，向稅捐單位領取紀念品。省稅務局聯合各縣市稅捐處，將在農曆年過年前，共同舉辦多項統一發票宣導活動，除可繼續提醒民眾保持購物索取**發票**的習慣外

# WordSketch

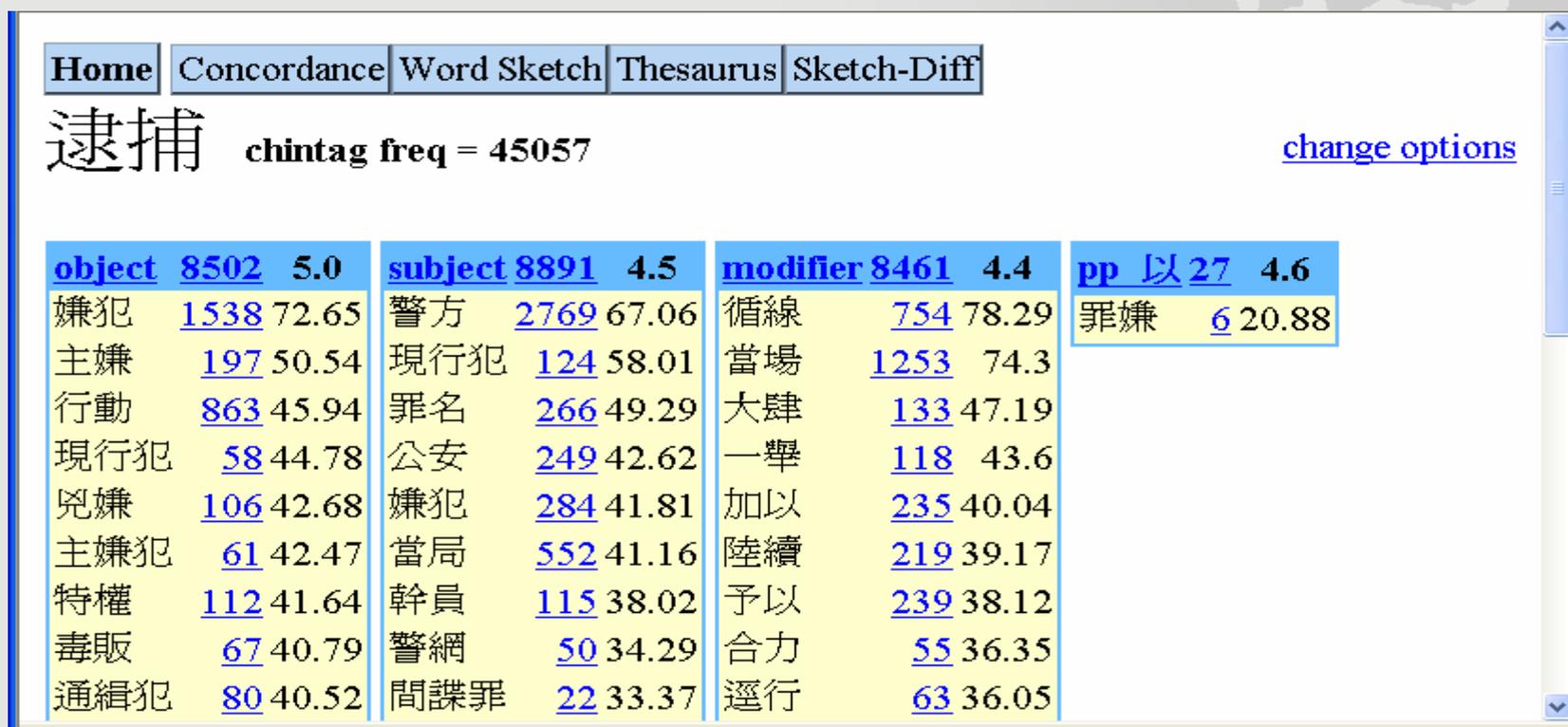
自動產生具體而微的簡易語法語意典

- 關鍵詞實際使用的完整描述
- 建立在不同的共現詞彙關係上
- 按照關係的顯著性排序

語音

# WordSketch:以「逮捕」爲例

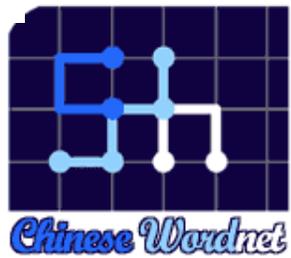
- It does produce the expected results with an easy to use interface.



Home Concordance Word Sketch Thesaurus Sketch-Diff

逮捕 chintag freq = 45057 [change options](#)

object 8502 5.0	subject 8891 4.5	modifier 8461 4.4	pp 以 27 4.6
嫌犯 <a href="#">1538</a> 72.65	警方 <a href="#">2769</a> 67.06	循線 <a href="#">754</a> 78.29	罪嫌 <a href="#">6</a> 20.88
主嫌 <a href="#">197</a> 50.54	現行犯 <a href="#">124</a> 58.01	當場 <a href="#">1253</a> 74.3	
行動 <a href="#">863</a> 45.94	罪名 <a href="#">266</a> 49.29	大肆 <a href="#">133</a> 47.19	
現行犯 <a href="#">58</a> 44.78	公安 <a href="#">249</a> 42.62	一舉 <a href="#">118</a> 43.6	
兇嫌 <a href="#">106</a> 42.68	嫌犯 <a href="#">284</a> 41.81	加以 <a href="#">235</a> 40.04	
主嫌犯 <a href="#">61</a> 42.47	當局 <a href="#">552</a> 41.16	陸續 <a href="#">219</a> 39.17	
特權 <a href="#">112</a> 41.64	幹員 <a href="#">115</a> 38.02	予以 <a href="#">239</a> 38.12	
毒販 <a href="#">67</a> 40.79	警網 <a href="#">50</a> 34.29	合力 <a href="#">55</a> 36.35	
通緝犯 <a href="#">80</a> 40.52	間諜罪 <a href="#">22</a> 33.37	逕行 <a href="#">63</a> 36.05	



# 同義詞典：Thesaurus

- 比對素描的相似度，找出使用上最接近的詞彙來

語音



網址(D) http://corpora.fi.muni.cz/chinese\_all/

Google G Go Bookmarks PageRank 141 blocked Check AutoLink AutoFill Send to

Home Concordance Word Sketch **Thesaurus** Sketch-Diff

### Thesaurus Entry Form

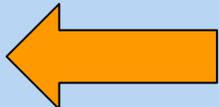
**Corpus:** chinese\_all\_ttd

**Word Form:** 快樂

Maximum number of items: 60

Minimum similarity between cluster items: 0.15

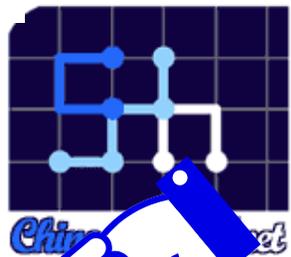
Show Similar Words



# 快樂 chinese\_all\_trd freq = 15438

愉快	0.308	輕鬆 0.242 活潑 0.205 溫馨 0.204 溫暖 0.199 樂觀 0.199 平靜 0.193 熱鬧 0.187 忙碌 0.185
幸福	0.295	美好 0.227 痛苦 0.207 美麗 0.2
歡樂	0.293	浪漫 0.186
開心	0.279	興奮 0.249 高興 0.243 勇敢 0.242 懂得 0.235 安心 0.225 關心 0.217 熱愛 0.216 辛苦 0.215 難過 0.215 期待 0.214 失望 0.214 在一起 0.213 害怕 0.212 喜歡 0.212 覺得 0.206 激動 0.205 放心 0.204 感受到 0.201 驕傲 0.199 用心 0.194 明白 0.189 喜愛 0.188 冷靜 0.187 瞭解 0.186 讚 0.185 滿意 0.185 相處 0.184 大聲 0.183 聽 0.183
有錢	0.221	聰明 0.211 富裕 0.209 舒適 0.194 乾淨 0.187 親近 0.186 賺錢 0.183
幸運	0.21	可愛 0.194 漂亮 0.188
讀書	0.197	體會 0.188 好好 0.183
歡喜	0.187	
健康	0.184	

Word Sketch Engine (ver:WSE-1.12-1.44)



# 快樂

chinese\_all\_trd freq = 15438

愉快	0.308	輕鬆	0.242	活潑	0.205	溫馨	0.204	溫暖	0.199	樂觀	0.199	平靜	0.193	熱鬧	0.187	忙碌	0.185																				
幸福	0.295	美好	0.227	痛苦	0.207	美麗	0.2																														
歡樂	0.293	浪漫	0.186																																		
開心	0.279	興奮	0.249	高興	0.243	勇敢	0.242	懂得	0.235	安心	0.225	關心	0.217	熱愛	0.216	辛苦	0.215	難過	0.215	激動	0.205	放心	0.204	感受到	0.201	驕傲	0.199	用心	0.194	明白	0.189	喜愛	0.188	冷靜	0.187	瞭解	0.187
有錢	0.221	聰明	0.211	富裕	0.209	舒適	0.194	乾淨	0.187	親近	0.186	賺錢	0.183																								
幸運	0.21	可愛	0.194	漂亮	0.188																																
讀書	0.197	體會	0.188	好好	0.183																																
歡喜	0.187																																				
健康	0.184																																				



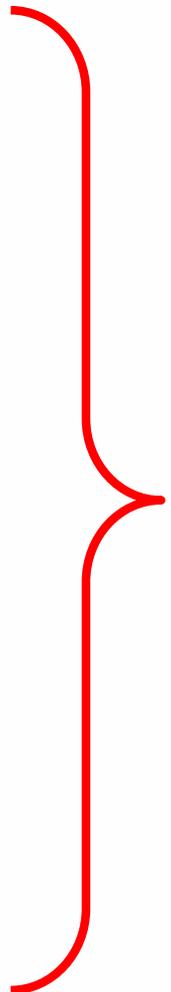
# 快樂/愉快 chinese\_all\_trd freq = 15438/11572

[change options](#)

## Common patterns

**快樂** 21 14 7 0 -7 -14 -21 **愉快**

modifier	2640	4052	5.4	7.1	subject	1200	1656	2.9	3.4
不	<a href="#">411</a>	<a href="#">1690</a>	30.9	48.7	佳節	<a href="#">242</a>	<a href="#">241</a>	78.2	75.5
非常	<a href="#">95</a>	<a href="#">420</a>	28.2	46.6	心情	<a href="#">8</a>	<a href="#">182</a>	15.7	58.4
很	<a href="#">354</a>	<a href="#">466</a>	39.2	39.6	秋節	<a href="#">53</a>	<a href="#">22</a>	53.7	36.7
十分	<a href="#">21</a>	<a href="#">225</a>	13.3	38.5	節日	<a href="#">58</a>	<a href="#">44</a>	40.3	34.4
相當	<a href="#">26</a>	<a href="#">274</a>	13.2	38.0	生活	<a href="#">38</a>	<a href="#">85</a>	19.0	26.2
更	<a href="#">189</a>	<a href="#">51</a>	30.6	13.1	人	<a href="#">52</a>	<a href="#">77</a>	15.5	17.4
最	<a href="#">275</a>	<a href="#">54</a>	30.4	9.8	們	<a href="#">31</a>	<a href="#">12</a>	15.6	6.8
又	<a href="#">117</a>	<a href="#">50</a>	27.7	15.6	民眾	<a href="#">22</a>	<a href="#">5</a>	12.5	2.5
過	<a href="#">69</a>	<a href="#">10</a>	22.7	4.6	<b>modifies</b>	<b>1541</b>	<b>631</b>	<b>1.0</b>	<b>0.4</b>
能	<a href="#">106</a>	<a href="#">37</a>	22.0	9.4	時光	<a href="#">90</a>	<a href="#">9</a>	58.8	25.3
還	<a href="#">8</a>	<a href="#">132</a>	2.2	21.6	心情	<a href="#">21</a>	<a href="#">34</a>	25.3	36.2
著	<a href="#">71</a>	<a href="#">33</a>	18.7	9.4	氣氛	<a href="#">35</a>	<a href="#">37</a>	27.2	32.6
太	<a href="#">8</a>	<a href="#">39</a>	7.4	18.5	地	<a href="#">6</a>	<a href="#">38</a>	6.2	26.2
並	<a href="#">15</a>	<a href="#">98</a>	5.3	18.3	生活	<a href="#">67</a>	<a href="#">10</a>	23.8	10.1
比較	<a href="#">26</a>	<a href="#">12</a>	17.6	9.2	假期	<a href="#">17</a>	<a href="#">12</a>	23.3	22.7
能夠	<a href="#">27</a>	<a href="#">8</a>	17.1	6.1	教育	<a href="#">8</a>	<a href="#">43</a>	3.9	20.7
就	<a href="#">47</a>	<a href="#">9</a>	14.6	1.8					
可以	<a href="#">36</a>	<a href="#">9</a>	14.0	2.9					
都	<a href="#">36</a>	<a href="#">29</a>	11.5	7.7					
也	<a href="#">40</a>	<a href="#">42</a>	8.6	6.6					
會	<a href="#">28</a>	<a href="#">10</a>	8.2	0.8					
更加	<a href="#">7</a>	<a href="#">6</a>	7.7	5.5					
卻	<a href="#">9</a>	<a href="#">7</a>	7.3	4.5					
了	<a href="#">69</a>	<a href="#">83</a>	5.7	4.4					
是否	<a href="#">8</a>	<a href="#">6</a>	5.6	3.0					



## "快樂" only patterns

modifier 2640 5.4		subject 1200 2.9		modifies 1541 1.0	
天天	<a href="#">25</a> 32.3	軍人節	<a href="#">20</a> 42.5	公主輪	<a href="#">37</a> 64.1
好不	<a href="#">9</a> 25.1	老人節	<a href="#">9</a> 32.9	捐血人	<a href="#">25</a> 45.8
真	<a href="#">27</a> 24.6	孩子	<a href="#">35</a> 28.9	童年	<a href="#">34</a> 43.1
愈	<a href="#">29</a> 24.0	小朋友	<a href="#">26</a> 27.9	炊事員	<a href="#">13</a> 36.6
一點	<a href="#">13</a> 20.6	榮民節	<a href="#">5</a> 23.4	聯播網	<a href="#">14</a> 34.9
極了	<a href="#">6</a> 19.5	兒童	<a href="#">41</a> 23.1	人生	<a href="#">40</a> 34.2
永遠	<a href="#">11</a> 17.4	青春	<a href="#">11</a> 21.5	夏令營	<a href="#">25</a> 32.8
來得	<a href="#">8</a> 16.7	媽咪	<a href="#">5</a> 21.4	拚財金	<a href="#">5</a> 31.3
一起	<a href="#">21</a> 16.2	孩子們	<a href="#">10</a> 20.5	天使	<a href="#">20</a> 30.9
不見得	<a href="#">5</a> 15.8	客輪	<a href="#">7</a> 20.5	電腦營	<a href="#">6</a> 29.0
起來	<a href="#">15</a> 15.4	學生	<a href="#">45</a> 20.4	暑假	<a href="#">5</a> 28.9
真正	<a href="#">15</a> 15.1	小孩	<a href="#">11</a> 19.1	捐血站	<a href="#">8</a> 28.7

## "愉快" only patterns

modifier 4052 7.1		subject 1656 3.4		modifies 631 0.4	
頗為	<a href="#">16</a> 20.4	神情	<a href="#">205</a> 73.2	經驗	<a href="#">95</a> 39.3
甚為	<a href="#">8</a> 17.8	旅途	<a href="#">77</a> 56.9	事件	<a href="#">124</a> 39.0
格外	<a href="#">13</a> 17.7	氣氛	<a href="#">60</a> 34.0	事情	<a href="#">19</a> 26.6
甚	<a href="#">18</a> 17.1	感	<a href="#">21</a> 32.9	笑容	<a href="#">9</a> 25.4
極為	<a href="#">16</a> 15.7	精神	<a href="#">120</a> 32.6	回憶	<a href="#">11</a> 24.6
首先	<a href="#">10</a> 10.8	身心	<a href="#">23</a> 22.5	地步	<a href="#">7</a> 22.0
極	<a href="#">6</a> 4.8	賓主	<a href="#">6</a> 19.4	情事	<a href="#">10</a> 20.0
已	<a href="#">15</a> 0.5	生日	<a href="#">8</a> 18.1	神情	<a href="#">6</a> 19.5
		表情	<a href="#">6</a> 17.6	感覺	<a href="#">5</a> 14.0
		國王	<a href="#">14</a> 17.5	婚姻	<a href="#">5</a> 13.9
		天皇	<a href="#">6</a> 17.1	場面	<a href="#">5</a> 13.6
		總理	<a href="#">33</a> 16.5	風波	<a href="#">5</a> 13.4

# 提示

Home Concordance Word Sketch **Thesaurus** Sketch-Diff

目標

chinese\_all\_trd freq = 189229

任務	0.343	方面 0.284	情況 0.265	行動 0.261	過程 0.259	條件 0.258	狀況 0.238	工作 0.235	項目 0.234	需求 0.232	結果 0.229	問題 0.227	情形 0.219													
目的	0.327	方案 0.318	方向 0.314	原則 0.311	計劃 0.31	措施 0.305	政策 0.303	計畫 0.29	策略 0.283	模式 0.26	立場 0.257	作法 0.252	標準 0.25	類 0.249	機制 0.249	體系 0.249	基礎 0.245	方法 0.241	制度 0.24	做法 0.239	體制 0.234	協議 0.23	辦法 0.228	經費 0.219	內容 0.217	式 0.215
指標	0.27	能力 0.243	水平 0.226	效果 0.223	效益 0.22	功能 0.22	力量 0.218	規模 0.216																		
理念	0.258	戰略 0.25	成果 0.249	方針 0.245	構想 0.23	進程 0.229	領域 0.223	經驗 0.218	觀念 0.218																	
重點	0.256	對象 0.249	範圍 0.249																							
要求	0.218																									



# SketchDifference

詞義與使用分析對比的最佳利器

- 比對兩個詞的素描，
- 從共有的關係詞中看出分布使用趨勢
- 從特有共現詞中看出特性

# 辭典的定義

「明星」和「演員」  
有甚麼不同？



# 重編國語辭典修訂本

辭典

教育部國語推行委員會編纂

- 明星
- 比喻傑出的人物。





# 重編國語辭典修訂本

辭典

教育部國語推行委員會編纂

- 演員
- 從事演藝工作的人員。



# 語言的使用

「明星」和「演員」

- 什麼時候可以互換？
- 什麼時候互換會改變意義？
- 什麼時候不能互換？

Home Concordance Word Sketch Thesaurus **Sketch-Diff**

### Word Sketch Differences Entry Form

Corpus:

First lemma:

Second lemma:

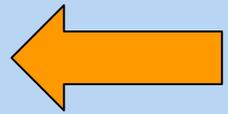
Sort grammatical relations:

Separate blocks:  all in one block  common/exclusive blocks

Minimum frequency:

Maximum number of items in a grammatical relation of the common block:

Maximum number of items in a grammatical relation of the exclusive block:



# 明星/演員 chinese\_all\_trd freq = 23923/23213 [change options](#)

## Common patterns

明星	21	14	7	0	-7	-14	-21	<b>演員</b>						
<b>a modifier</b>	1729	4116	1.7	4.2	<b>measure</b>	688	1532	1.1 2.5	<b>n_modifier</b>	13964	9897	1.6	1.2	
著名	84	738	39.0	70.4	位	310	576	50.3	54.0	大牌	162	13	57.7	21.6
武打	65	14	62.6	31.1	名	79	617	26.4	50.4	偶像	297	8	57.1	11.7
老牌	42	110	44.5	56.4	批	13	39	17.3	25.7	喜劇	50	144	33.8	52.8
資深	6	237	11.2	53.5	個	79	159	21.3	23.9	影視	248	64	50.8	32.3
知名	29	232	26.7	52.4	屆	13	7	12.6	5.7	好萊塢	190	70	50.5	37.0
當紅	43	30	49.4	38.7	場	11	5	12.5	4.9	新生代	25	133	24.0	50.3
大	468	19	46.6	5.1	次	9	17	8.1	9.5	男	41	355	16.1	46.6
年輕	8	136	13.3	43.8	<b>possessor</b>	459	921	1.0	2.1	演技派	31	37	39.8	44.5
最佳	27	135	23.2	39.5	知名度	6	10	17.6	21.1	歌仔戲	11	83	15.4	42.6
小	38	232	20.5	37.9	國家	9	26	4.8	8.6	芭蕾舞	7	63	13.6	42.5
老	21	112	19.6	35.5	<b>and/or</b>	336	1517	0.4	1.7	青年	26	564	5.4	41.7
眾多	44	16	33.8	17.5	導演	11	381	21.3	68.4	舞蹈	12	214	7.9	41.4
一流	9	22	18.5	25.4	歌星	16	50	32.2	42.6	電影	343	301	38.5	39.4
名	11	198	4.6	25.3	歌手	7	50	17.0	34.8	歌劇	11	64	13.4	35.2
新一代	12	13	22.9	20.6	藝術家	6	36	15.5	30.3	男女	28	112	14.1	31.7
已故	11	7	22.5	14.9						女	234	18	30.9	7.2

# 明星/演員 chinese\_all\_trd freq = 23923/23213 [change options](#)

## Common patterns

明星	21	14	7	0	-7	-14	-21	演員	
a modifier	1729	4116	1.7	4.2	measure	688	1532	1.1	2.5
著名	84	738	39.0	70.4	位	310	576	50.3	54.0
武打	65	14	62.6	31.1	名	79	617	26.4	50.4
老牌	42	110	44.5	56.4	批	13	39	17.3	25.7
資深	6	237	11.2	53.5	個	79	159	21.3	23.9
知名	29	232	26.7	52.4	屆	13	7	12.6	5.7
當紅	43	30	49.4	38.7	場	11	5	12.5	4.9
大	468	19	46.6	5.1	次	9	17	8.1	9.5
年輕	8	136	13.3	43.8	possessor	459	921	1.0	2.1
最佳	27	135	23.2	39.5	知名度	6	10	17.6	21.1
小	38	232	20.5	37.9	國家	9	26	4.8	8.6
老	21	112	19.6	35.5	and/or	336	1517	0.4	1.7
眾多	44	16	33.8	17.5	導演	11	381	21.3	68.4
一流	9	22	18.5	25.4	歌星	16	50	32.2	42.6
名	11	198	4.6	25.3	歌手	7	50	17.0	34.8
新一代	12	13	22.9	20.6	藝術家	6	36	15.5	30.3
已故	11	7	22.5	14.9	n modifier	13964	9897	1.6	1.2
					大牌	162	13	57.7	21.6
					偶像	297	8	57.1	11.7
					喜劇	50	144	33.8	52.8
					影視	248	64	50.8	32.3
					好萊塢	190	70	50.5	37.0
					新生代	25	133	24.0	50.3
					男	41	355	16.1	46.6
					演技派	31	37	39.8	44.5
					歌仔戲	11	83	15.4	42.6
					芭蕾舞	7	63	13.6	42.5
					青年	26	564	5.4	41.7
					舞蹈	12	214	7.9	41.4
					電影	343	301	38.5	39.4
					歌劇	11	64	13.4	35.2
					男女	28	112	14.1	31.7
					女	234	18	30.9	7.2

“演員” only patterns

possession 693 1.6	
演技	<a href="#">15</a> 33.7
演出	<a href="#">18</a> 21.9
手	<a href="#">21</a> 21.8
功力	<a href="#">6</a> 20.8
服裝	<a href="#">13</a> 19.6
積極性	<a href="#">9</a> 17.6
特質	<a href="#">5</a> 15.9
表現	<a href="#">13</a> 15.4
肢體	<a href="#">5</a> 14.9
角色	<a href="#">8</a> 14.4
技巧	<a href="#">5</a> 13.8
動作	<a href="#">6</a> 12.3

subject_of 3924 1.6	
謝幕	<a href="#">20</a> 44.9
擔綱	<a href="#">33</a> 41.7
親切	<a href="#">23</a> 30.1
跳起	<a href="#">11</a> 29.5
演戲	<a href="#">9</a> 29.4
演	<a href="#">22</a> 29.1
握手	<a href="#">19</a> 28.2
齊聚一堂	<a href="#">14</a> 27.3
載歌載舞	<a href="#">10</a> 26.3
清唱	<a href="#">8</a> 24.4
拍戲	<a href="#">6</a> 24.3
精湛	<a href="#">13</a> 24.1

n_modifier 9897 1.2	
京劇	<a href="#">246</a> 53.7
一級	<a href="#">424</a> 50.8
相聲	<a href="#">120</a> 50.4
實力派	<a href="#">71</a> 49.1
雜技	<a href="#">129</a> 48.5
替身	<a href="#">33</a> 43.7
特技	<a href="#">92</a> 43.1
特型	<a href="#">21</a> 42.4
越劇	<a href="#">45</a> 41.0
話劇	<a href="#">59</a> 38.1
舞台劇	<a href="#">45</a> 37.4
歌舞伎	<a href="#">21</a> 36.9

modifies 8784 1.0	
郎雄	<a href="#">51</a> 50.1
李康生	<a href="#">29</a> 42.4
柯俊雄	<a href="#">29</a> 42.3
金士傑	<a href="#">22</a> 39.9
柯受良	<a href="#">21</a> 37.7
胡軍	<a href="#">18</a> 36.6
姜昆	<a href="#">19</a> 36.0
鄭亞雲	<a href="#">15</a> 35.9
楊貴媚	<a href="#">20</a> 35.3
戴立忍	<a href="#">19</a> 34.2
楊呈偉	<a href="#">14</a> 33.9
濮存昕	<a href="#">13</a> 33.8

object_of 1535 0.7	
獨唱	<a href="#">50</a> 54.1
客串	<a href="#">22</a> 41.9
演	<a href="#">23</a> 34.0

"明星" only patterns

a_modifier	1729	1.7
超級	<a href="#">219</a>	63.7
演藝	<a href="#">78</a>	51.2
美式	<a href="#">46</a>	48.6
耀眼	<a href="#">26</a>	46.1
閃亮	<a href="#">18</a>	40.3
過氣	<a href="#">10</a>	36.9
超冷	<a href="#">5</a>	30.6
三棲	<a href="#">6</a>	29.8
佳	<a href="#">21</a>	29.3
頭號	<a href="#">19</a>	27.0
溫	<a href="#">12</a>	24.9
簡	<a href="#">10</a>	24.6

n_modifier	13964	1.6
籃球	<a href="#">433</a>	46.3
職籃	<a href="#">179</a>	45.6
足球	<a href="#">675</a>	44.9
聯隊	<a href="#">232</a>	43.0
卡	<a href="#">99</a>	41.9
夢幻	<a href="#">68</a>	40.8
網球	<a href="#">293</a>	39.2
抗癌	<a href="#">50</a>	34.9
偶像級	<a href="#">16</a>	34.7
恆康	<a href="#">15</a>	34.7
演藝圈	<a href="#">36</a>	33.8
演藝界	<a href="#">34</a>	33.8



modifies	13046	1.4
球員	<a href="#">1095</a>	59.7
對抗賽	<a href="#">346</a>	59.3
聯隊	<a href="#">301</a>	54.2
排名賽	<a href="#">108</a>	49.0
辛浦森	<a href="#">56</a>	48.4
籃球隊	<a href="#">167</a>	46.9
馬拉杜納	<a href="#">30</a>	43.0
白隊	<a href="#">47</a>	41.2
後衛	<a href="#">141</a>	40.0
前鋒	<a href="#">149</a>	38.8
馬拉多納	<a href="#">39</a>	38.4
王貞治	<a href="#">30</a>	37.8

possession	506	1.1
風采	<a href="#">26</a>	37.4
架子	<a href="#">12</a>	32.2
架勢	<a href="#">10</a>	32.1
搖籃	<a href="#">12</a>	30.2
名字	<a href="#">12</a>	24.4
架式	<a href="#">5</a>	23.8
丰采	<a href="#">5</a>	21.8
光環	<a href="#">5</a>	20.2
魅力	<a href="#">7</a>	18.7
照片	<a href="#">9</a>	18.5
故事	<a href="#">7</a>	16.3
青少年	<a href="#">6</a>	10.7

measure	688	1.1
類	<a href="#">79</a>	51.3
隻	<a href="#">7</a>	15.8
路	<a href="#">6</a>	10.8
所	<a href="#">15</a>	10.2
家	<a href="#">7</a>	8.2

possessor	459	1.0
世界級	<a href="#">11</a>	25.3
名氣	<a href="#">5</a>	19.9
比賽	<a href="#">21</a>	17.1
聯隊	<a href="#">6</a>	16.9
氣	<a href="#">5</a>	14.3

subject_of	2280	0.9
薈萃	<a href="#">53</a>	49.3
雲集	<a href="#">36</a>	39.5
三缺一	<a href="#">8</a>	37.4
開店	<a href="#">12</a>	34.6
評選	<a href="#">26</a>	32.8

object_of	1536	0.7
啓	<a href="#">44</a>	48.8
崇拜	<a href="#">19</a>	34.0
棲	<a href="#">11</a>	33.2
考上	<a href="#">16</a>	31.8
服務	<a href="#">93</a>	29.0

Home Concordance Word Sketch Thesaurus Sketch-Diff **Frequency** Collocation

KWIC/Sentence View options **Sample** Filter Sort Save

Corpus: chinese\_all\_trd  
Hits: 433  
[conc description](#)

First | Previous Page 2 of 22 Go Next | Last

- [CNA19920506.0058](#)
- [CNA19920506.0058](#)
- [CNA19920511.0029](#)
- [CNA19920605.0110](#)
- [CNA19920712.0164](#)
- [CNA19920727.0158](#)
- [CNA19920728.0285](#)
- [CNA19920801.0284](#)
- [CNA19920804.0252](#)
- [CNA19920805.0039](#)
- [CNA19920926.0159](#)
- [CNA19920930.0144](#)
- [CNA19920930.0161](#)
- [CNA19930111.0159](#)
- [CNA19930621.0111](#)
- [CNA19930718.0199](#)
- [CNA19930724.0200](#)
- [CNA19930814.0062](#)
- [CNA19930814.0062](#)
- [CNA19930815.0130](#)

鍾錦隆 台北 六日 電) <p></p>韓國 前 女子 籃球 **明星** 朴贊淑 認為，企業界對球隊的支持不  
 球員，獲觀眾票選為最受歡迎的籃球**明星**。<p></p>會多次代表南韓女籃出賽的  
 去年賺得一千五百萬美元。<p></p>籃球**明星** 喬丹的收入為一千六百萬美元，使他  
 愛滋病毒而退出洛杉磯湖人隊的前美國籃球**明星**「魔術」強生，昨晚喜獲麟兒。<p></p>  
 第十五屆威廉瓊斯盃國際籃球邀請賽男子**明星**球員，除二名地主球員外，其餘十  
 奧運第一天的比賽結果，美國男子籃球**明星**代表隊第一場對安哥拉隊的比賽，以  
 巴塞隆納二十八日法新電) <p></p>美國籃球**明星**魔術強森，右膝的韌帶受傷之後，  
 讓美國那些百萬美元身價的超級籃球**明星**，參加奧運競賽。但是施密特不以為然  
 了。<p></p>昨天這個由美國職業籃球**明星**組成的球隊，在奧運會場舉行了與賽  
 奧林匹克籃球隊的喬丹與另外五位NBA籃球**明星**卻與耐吉公司有約，使用耐吉公司  
 合眾國際電) <p></p>感染愛滋病毒的美國籃球**明星**「魔術」強生，昨天退出美國國家  
 二十九日合眾國際電) <p></p>美國職業籃球**明星**「魔術」強生今天宣布，他將在本  
 二十九日合眾國際電) <p></p>美國職業籃球**明星**「魔術」強生今天宣布，他將在本  
 王同禹 台北 十一日 電) <p></p>今年的甲組籃球**明星**對抗賽將於本月十七日下午二時，  
 葛瑞菲絲·卓納與前國會議員兼籃球**明星**湯姆·麥米倫主持白宮體育委員會的業務  
 周揚力 林口 十八日 電) <p></p>瓊斯盃籃球邀請賽**明星**球員出爐，獲大會最有價值球員的  
 <p></p>尤恩是第三位來華訪問的美國籃球**明星**、也是第一位訪華的NBA現役球員  
 聯邦調查局官員今天說，美國超級籃球**明星**邁克喬登的父親，可能在綁票事件中  
 歲的詹姆斯喬登，芝加哥公牛隊籃球**明星**邁克喬登的父親。<p></p>詹姆斯胸部中彈  
 <p></p>職業棒球選手謝長亨、王光熙與籃球**明星**鄭志龍、錢薇娟，今天聯手在台中主持

**吃喝** chinese\_all\_trd freq = 53654/19561

Common patterns

吃 21 14 7 0 -7 -14 -21 喝

pp_假 2 6 3.4 25.4	object 23421 13015 6.9 9.5	modifier 21712 6092 4.5 3.1
酒 1 3 9.0 18.0	酒 8 4756 5.4 106.0	少 400 87 49.5 35.3
	花酒 4 639 9.5 91.1	多 1173 281 47.5 36.7
	茶 1 748 0.6 74.8	一口 69 85 35.0 46.2
	藥 1430 8 72.0 7.2	同 379 4 45.3 4.7
	春酒 5 107 14.6 62.7	過 1181 200 45.3 30.1
	牛奶 18 291 17.4 61.2	不 3429 1011 42.0 36.2
	早餐 375 1 59.3 1.6	常 227 97 38.4 35.2
	東西 596 16 56.1 12.5	沒 324 81 38.0 28.2
	消夜 75 1 53.5 4.8	天天 77 23 37.0 27.3
	喜酒 6 34 18.8 46.3	連 177 21 36.4 17.9
	湯圓 84 1 45.8 3.3	了 4425 1125 36.2 29.3
	奶 143 92 45.7 42.7	倒 116 1 35.9 1.7
	飲料 4 183 3.2 43.6	一起 318 158 35.3 35.2
	食物 293 4 42.6 3.5	不要 226 109 31.1 30.6
	湯 1 59 2.3 40.9	怎麼 78 6 29.2 9.4
	口水 2 51 5.1 40.1	不准 38 38 22.2 28.9
	稀飯 36 14 39.3 28.1	給他 56 10 28.6 15.4
	奶茶 1 36 3.3 38.8	不能 284 71 28.2 20.6
	農藥 4 134 2.9 37.8	著 563 247 27.4 27.4
	習慣 155 149 33.2 36.9	邊 146 26 27.2 15.8
	巧克力 64 1 36.1 2.4	去 190 94 26.8 27.0
	麵條 40 1 36.1 3.6	誤 15 18 19.4 26.8
	零嘴 19 1 33.6 5.4	要 752 186 26.7 20.2
	豆漿 3 24 9.0 33.3	很少 44 15 26.4 20.0
	奶水 18 19 29.9 33.2	不用 48 4 26.2 8.0



# Objects Shared by Chi1 and He1

吃/喝 chinese\_all\_trd\_test freq = 53654/19561

## Common patterns

吃 21 14 7 0 -7 -14 -21 喝

SentObject_of	3859	1065	7.2	4.9
喜歡	<u>557</u>	<u>173</u>	69.7	57.7
試	<u>371</u>	<u>22</u>	68.2	29.2
愛	<u>571</u>	<u>185</u>	66.0	55.5
拒	<u>167</u>	<u>6</u>	55.5	15.4
嗜	<u>70</u>	<u>9</u>	55.5	27.4
顧不上	<u>63</u>	<u>12</u>	53.0	31.3
敢	<u>168</u>	<u>36</u>	47.1	31.9
捨不得	<u>39</u>	<u>8</u>	42.8	24.6
請	<u>216</u>	<u>44</u>	39.8	26.5
喜愛	<u>45</u>	<u>24</u>	32.7	30.6
怕	<u>49</u>	<u>10</u>	32.0	18.6
放心	<u>29</u>	<u>6</u>	30.5	16.6
喜	<u>36</u>	<u>7</u>	30.0	16.2
涉嫌	<u>9</u>	<u>44</u>	8.2	29.4
知道	<u>53</u>	<u>24</u>	25.2	22.6
喜好	<u>10</u>	<u>10</u>	20.7	25.1

Object	33038	16684	3.7	4.6
酒	<u>13</u>	<u>5198</u>	7.5	106.9
茶	<u>7</u>	<u>825</u>	7.0	75.5
藥	<u>1558</u>	<u>8</u>	73.0	7.3
牛奶	<u>24</u>	<u>386</u>	19.4	65.6
春酒	<u>5</u>	<u>119</u>	13.6	63.0
東西	<u>639</u>	<u>16</u>	53.3	11.1
食物	<u>610</u>	<u>6</u>	52.9	5.0
喜酒	<u>6</u>	<u>43</u>	17.7	48.8
奶	<u>160</u>	<u>106</u>	46.5	44.6
頓	<u>167</u>	<u>5</u>	43.9	8.0
稀飯	<u>47</u>	<u>23</u>	41.5	33.9
水	<u>16</u>	<u>320</u>	2.4	35.5
習慣	<u>181</u>	<u>165</u>	31.5	35.3
碗	<u>75</u>	<u>19</u>	35.2	21.4
奶水	<u>18</u>	<u>20</u>	29.8	34.3
母乳	<u>24</u>	<u>20</u>	32.3	32.7

Modifier	13757	4501	4.5	3.7
少	<u>440</u>	<u>95</u>	65.6	46.5
多	<u>1289</u>	<u>304</u>	61.7	46.6
同	<u>384</u>	<u>5</u>	52.2	7.4
不	<u>1885</u>	<u>909</u>	45.7	45.1
一起	<u>317</u>	<u>159</u>	44.1	42.0
大口	<u>5</u>	<u>24</u>	17.6	44.1
常	<u>198</u>	<u>94</u>	43.3	39.7
天天	<u>76</u>	<u>23</u>	42.8	30.9
沒	<u>307</u>	<u>79</u>	42.8	31.2
邊	<u>145</u>	<u>27</u>	41.7	25.1
連	<u>179</u>	<u>20</u>	41.1	19.4
只	<u>363</u>	<u>128</u>	37.0	31.1
不要	<u>214</u>	<u>110</u>	36.8	35.3
給他	<u>54</u>	<u>10</u>	34.9	18.9
不能	<u>258</u>	<u>71</u>	34.3	25.3
著	<u>177</u>	<u>36</u>	34.0	21.2

# “Chi1” only patterns

## “吃” only patterns

### SentObject\_of 3859 7.2

愁	<a href="#">103</a>	57.9
講究	<a href="#">27</a>	32.5
嚐試	<a href="#">13</a>	26.5
忌	<a href="#">8</a>	25.9
寧可	<a href="#">16</a>	25.3
捨得	<a href="#">11</a>	24.8
擔心	<a href="#">46</a>	24.1
拒絕	<a href="#">45</a>	22.9
寧願	<a href="#">11</a>	21.2
討厭	<a href="#">7</a>	21.0
忘	<a href="#">13</a>	20.8
記得	<a href="#">12</a>	20.3

### Modifier 13757 4.5

倒	<a href="#">128</a>	41.6
津津有味	<a href="#">19</a>	36.8
怎麼	<a href="#">78</a>	35.0
硬	<a href="#">41</a>	32.5
有得	<a href="#">20</a>	31.8
不用	<a href="#">47</a>	31.0
按時	<a href="#">35</a>	29.4
常年	<a href="#">28</a>	27.4
該	<a href="#">47</a>	25.2
年年	<a href="#">24</a>	24.5
一律	<a href="#">31</a>	23.3
著實	<a href="#">13</a>	23.3

### Subject 11519 4.3

飯	<a href="#">718</a>	78.7
啞巴	<a href="#">30</a>	42.0
最愛	<a href="#">72</a>	39.4
柿子	<a href="#">27</a>	33.5
糖	<a href="#">45</a>	31.5
金飯碗	<a href="#">14</a>	31.5
全家	<a href="#">44</a>	31.1
魚	<a href="#">73</a>	30.0
東西	<a href="#">78</a>	29.4
啞吧	<a href="#">11</a>	29.4
們同	<a href="#">10</a>	29.2
金碗	<a href="#">9</a>	27.7

### Object 33038 3.7

敗仗	<a href="#">326</a>	72.3
晚飯	<a href="#">310</a>	71.6
飯	<a href="#">802</a>	71.0
定心丸	<a href="#">211</a>	68.0
午飯	<a href="#">241</a>	67.5
大鍋飯	<a href="#">245</a>	66.6
閉門羹	<a href="#">173</a>	66.5
年夜飯	<a href="#">270</a>	65.8
狗肉	<a href="#">190</a>	61.5
肉	<a href="#">488</a>	60.1
虧	<a href="#">329</a>	59.2
頓飯	<a href="#">84</a>	59.2

### PP\_在 165 1.5

工地	<a href="#">17</a>	32.3
----	--------------------	------

### Modifies 1527 0.1

東西	<a href="#">109</a>	47.6
----	---------------------	------

# “He1” only patterns

## “喝” only patterns

### SentObject\_of 1065 4.9

涉	<u>8</u>	17.7
盛行	<u>5</u>	17.6
疑	<u>5</u>	15.0

### Object 16684 4.6

花酒	<u>666</u>	90.8
咖啡	<u>969</u>	74.3
啤酒	<u>421</u>	58.5
下午茶	<u>101</u>	57.1
白開水	<u>67</u>	54.3
口水	<u>114</u>	52.5
開水	<u>108</u>	51.2
杯	<u>237</u>	50.3
飲料	<u>283</u>	49.6
紅酒	<u>83</u>	49.6
花酒案	<u>33</u>	49.3
綠茶	<u>88</u>	48.7

### Modifier 4501 3.7

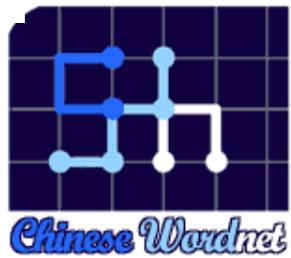
厲聲	<u>13</u>	35.5
猛	<u>17</u>	21.5
成天	<u>5</u>	19.7
獨自	<u>7</u>	16.8
有沒有	<u>7</u>	16.7
一口	<u>5</u>	16.2
盡情	<u>6</u>	15.9
其實	<u>5</u>	13.0
太	<u>7</u>	12.8
難免	<u>5</u>	12.7
總共	<u>5</u>	9.0
並	<u>12</u>	3.6

### Subject 3235 3.0

開水	<u>19</u>	32.3
口水	<u>16</u>	30.0
胡吃海	<u>5</u>	28.5
馬永成	<u>7</u>	19.9
蕭敦仁	<u>5</u>	19.3
友人	<u>17</u>	18.9
校長	<u>30</u>	18.1
自來水	<u>17</u>	17.7
粥	<u>5</u>	17.6
督學	<u>6</u>	15.8
牛奶	<u>7</u>	15.8
檢察官	<u>20</u>	14.8

### Modifies 306 0.1

水	<u>44</u>	34.5
飲料	<u>15</u>	28.6



# 素描對比也可用於 兩岸詞彙對比研究

- 利用CGW包含來自大陸與台灣兩個次語料庫的特性，

語音

# 大陸和台灣用語比較

• 出租車

台灣

109 筆

大陸



• 公安

台灣

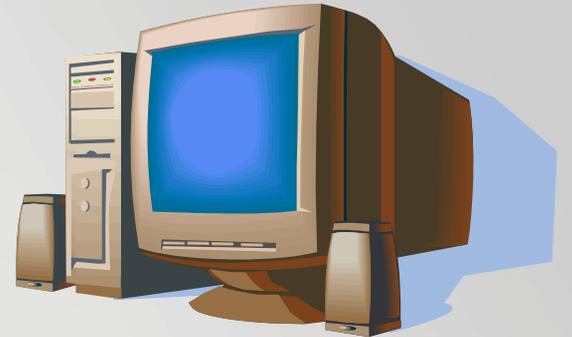
18,641 筆

大陸



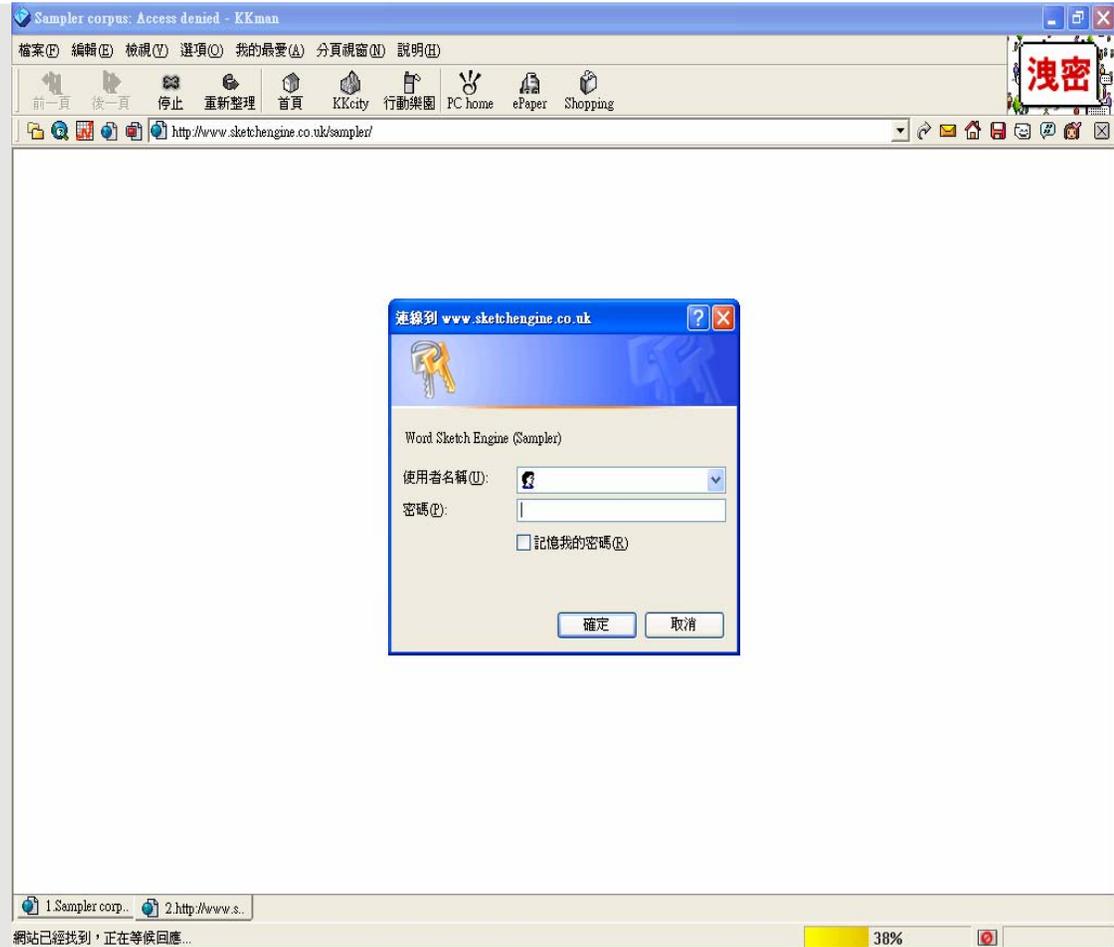
# 系統運用

- 國文學習
- 提供查詢功能
- 減少詞彙的誤用
- 輔助外國人學國文
- 大陸和台灣用語比較



# 帳號申請

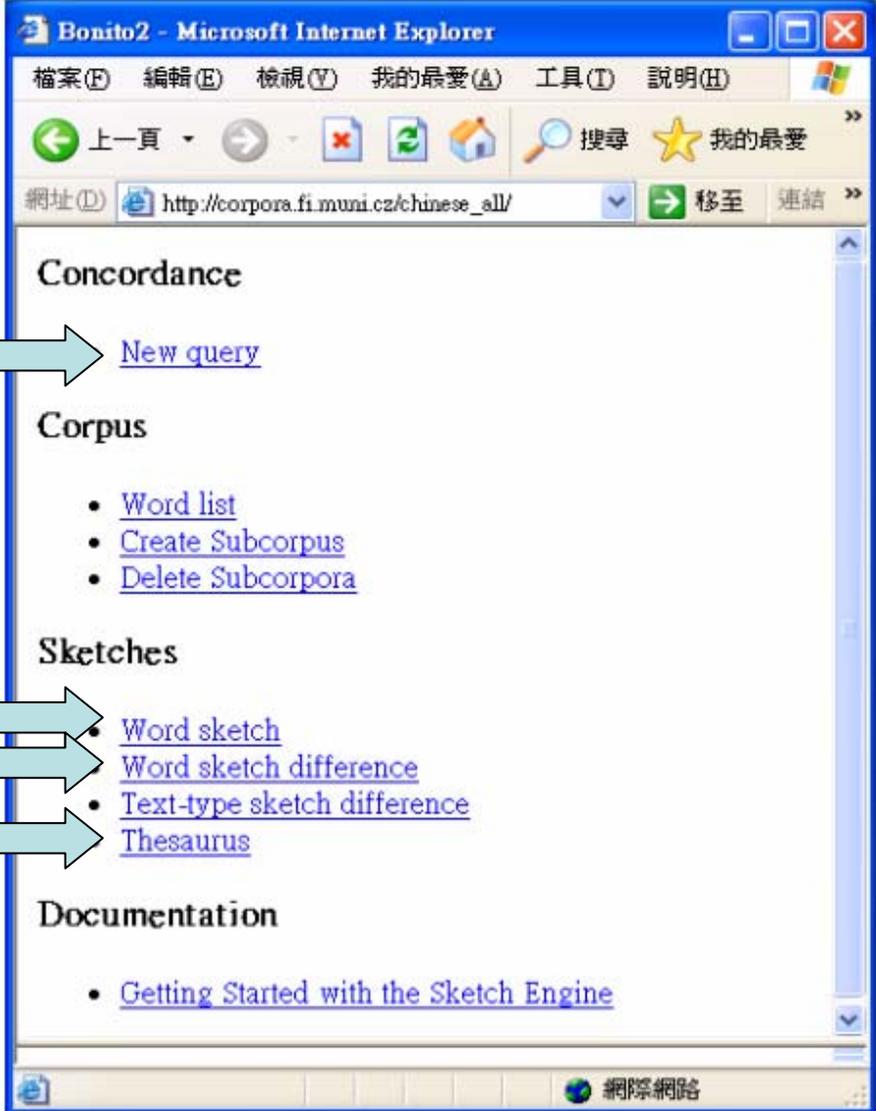
- 免費
- 使用者帳號和密碼
- 請跟我們聯絡



- 中研院版本
- <http://140.109.150.73/wse/>
- 帳號: guest
- 密碼: try2guest
- 捷克版本
- <http://corpora.fi.muni.cz/chinese/>
- 帳號: chinese
- 密碼: chinese

語音

# 操作流程- 功能列表



**Concordance**

- [New query](#)

**Corpus**

- [Word list](#)
- [Create Subcorpus](#)
- [Delete Subcorpora](#)

**Sketches**

- [Word sketch](#)
- [Word sketch difference](#)
- [Text-type sketch difference](#)
- [Thesaurus](#)

**Documentation**

- [Getting Started with the Sketch Engine](#)

檢索關鍵詞或詞組的語料

對語料庫的進階處理

檢索詞彙的進階相關訊息

其他相關訊息與使用說明

中央研究院語言學研究所

「中文詞彙特性速描系統」帳號使用切結書

【使用內容及操作方式說明】

軟體名稱：中文詞彙特性速描系統

【使用對象】

學齡以上民眾得申請使用

【使用方式】

1. 免費提供學術研究或學習使用，但不包括違法及營利之使用
2. 中央研究院語言學研究所保留授權帳號之所有權
3. 論文發表或引用時，請註明資料取得來源：

Huang, Chu-Ren, Adam Kilgarriff, Yicing Wu, Chih-Min Chiu, Simon Smith, Pavel Rychlý, Ming-Hong Bai, and Keh-Jiann Chen. 2005. Chinese Sketch Engine and the Extraction of Collocations. Proceedings of the Fourth SIGHAN Workshop on Chinese Language Processing, 48-55. October 14-15. Jeju, Korea.

Kilgarriff, Adam, Pavel Rychlý, Pavel Smrz and David Tugwell. 2004. The Sketch Engine. Proceedings of EURALEX, Lorient, France.

Chinese Sketch Engine

<http://wordsketch.ling.sinica.edu.tw>

禁止部份或全部轉移給其它第三者使用

申請單位(所屬單位)：

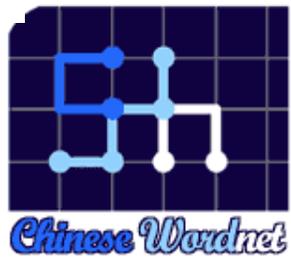
申請人：

計畫名稱(申請目的)：

聯絡電話：

日 期： 年 月 日

語音



# 中央研究院語言學研究所

## 「中文詞彙特性速描系統」帳號申請表

帳號 (本欄由承辦人輸入) 中文姓名 英文姓名 單位  個人申請 職別 **E-mail Address** 聯絡電話 異動備註 (本欄由承辦人輸入) (本欄由系統管理員輸入) ※申請核准後將由系統自動產生一組帳號密碼，並以電子郵件寄送至您的信箱，帳號密碼皆無法修改且僅供單人使用，請妥善保存

申請人簽名： \_\_\_\_\_ 承辦人： \_\_\_\_\_

申請日期 \_\_\_\_\_ 年 \_\_\_\_\_ 月 \_\_\_\_\_ 日

### 異動紀錄

承辦人員： \_\_\_\_\_ 處理日期： \_\_\_\_\_  
\_\_\_\_\_ 處理結果： \_\_\_\_\_

簡報完畢

謝謝



- SINICA BOW
- <http://bow.sinica.edu.tw>
- Chinese Wordnet
- <http://ling.sinica.edu.tw/cwn>

語音

# 美白/解酒

chinese\_all\_trd freq = 542/144

[change options](#)

## Common patterns

美白	21	14	7	0	-7	-14	-21	解酒
----	----	----	---	---	----	-----	-----	----

modifies	228	76	2.1	2.5
效果	<u>30</u>	<u>5</u>	33.4	15.6
功效	<u>8</u>	<u>5</u>	24.9	21.9
產品	<u>20</u>	<u>18</u>	18.8	22.4

## "美白" only patterns

subject	104	3.4
牙齒	<u>23</u>	46.0
義	<u>13</u>	32.5
雷射	<u>10</u>	28.9
皮膚	<u>7</u>	22.9

modifies	228	2.1
乳霜	<u>12</u>	48.1
保養品	<u>17</u>	45.2
面膜	<u>6</u>	31.6
化妝品	<u>12</u>	29.8
成份	<u>11</u>	26.4
霜	<u>5</u>	23.8
用品	<u>5</u>	14.7
商品	<u>5</u>	10.1