

## 摘要

中研院臺灣當代華語語料庫收錄多達7千多萬則臺灣正體華語詞彙(73007511)。取得授權的語料來源為2015至2020年間出版的報章雜誌(總共86267296個字，50486010個詞及238361篇文章)及兒童讀物(總共37162935個字，22297533個詞及112204篇文章)。內容涵蓋了生活、社會、金融、科學、文學、文化、教育以及哲學等多種類型的文章。本語料庫所有文本的內容都依照詞彙斷開並標示詞類標記及實體辨識標記。進一步彙整每個詞彙的詞類標記，以及計算其詞彙頻率(出現次數、百萬詞頻、百萬詞頻取對數)，語境變異(出現次數、出現次數取對數、百萬詞頻、百萬詞頻取對數、千詞頻)，以及語意變異等三項指標。

## 技術優勢

1. 2015至2020年間當代臺灣華語文本大數據分析
2. 提供詞彙頻率、語境變異及語義變異等三項指標
3. 依據不同年齡層閱讀文本的類型，可分別建立適合小孩、成人及老年人等閱讀經驗的語料庫

A	B	C	D	E	F	G	H	I	J	K
word	WF_ppm	log(WF_ppm)	CD	log(CD)	CD_ppm	log(CD_ppm)	CD_ppk	SD	POS	
詞彙	百萬詞頻	百萬詞頻取對數	詞類變異	詞類變異(取對數)	(百萬詞頻)	詞類變異(取對數)	(百萬詞頻取對數)	(千詞頻)	詞類標記	
1 向	55781	764.04	2.88	43711	4.64	124376.85	5.09	124.38	1.53	{P, 55256, (D, 79), (V, 441), (N, 5)}
2 前	55694	762.85	2.88	38621	4.59	109895.58	5.04	109.89	1.28	{C, 5334, (P, 48825), (V, 1835)}
3 誰	55330	760.86	2.89	23988	4.53	96566.40	4.98	96.57	1.47	{N, 5330}
4 國家	54991	751.85	2.88	41544	4.62	118210.79	5.07	118.21	1.27	{(C, 23824), (P, 28739), (N, 328)}
5 發現	54821	750.99	2.88	40648	4.61	115961.28	5.06	115.66	1.43	{(V, 53889), (N, 528), (N, 4)}
7 成為	54330	744.17	2.87	43946	4.64	125045.53	5.10	125.05	1.52	{(V, 54330)}
8 未	53356	730.83	2.86	39256	4.59	111700.43	5.05	111.70	1.56	{(P, 53553), (N, 3)}
9 指出	53020	726.23	2.86	41279	4.62	117456.75	5.07	117.46	1.58	{(V, 53020)}
10 升	52575	720.13	2.86	39938	4.58	108007.06	5.03	108.01	1.55	{(N, 52200), (N, 61), (N, 308)}
11 曾	52545	719.72	2.86	41501	4.62	118088.44	5.07	118.09	1.47	{(D, 52304), (N, 241)}
12 發展	52187	714.82	2.85	33716	4.53	95936.72	4.98	95.94	1.51	{(V, 48184), (N, 2400), (N, 1033)}
13 許多	52119	713.89	2.85	42569	4.63	121127.36	5.08	121.13	1.45	{(N, 50411), (D, 1708)}
14 事	52016	712.47	2.85	39482	4.60	112343.50	5.05	112.34	1.31	{(N, 51970), (V, 46)}
15 世界	51375	703.69	2.85	34874	4.54	92317.72	5.00	92.32	1.43	{(N, 51375)}
16 因此	50909	697.31	2.84	42598	4.63	122324.24	5.09	122.23	1.49	{(C, 50909)}
17 北	50544	692.31	2.84	35022	4.54	99652.86	5.00	99.65	1.51	{(P, 38716), (C, 2984), (V, 8392), (N, 439), (N, 9), (N, 2), (D, 2)}
18 學習	50110	686.37	2.84	26298	4.42	74829.27	4.87	74.83	1.34	{(V, 48070), (N, 2040)}
19 回	50078	685.93	2.84	35323	4.55	105309.33	5.00	105.31	1.58	{(N, 49116), (C, 36), (N, 296), (N, 52), (N, 35), (N, 12)}
20 未來	49330	675.68	2.83	37118	4.57	105166.89	5.02	105.62	1.55	{(N, 49330)}
21 總共	49260	674.73	2.83	22428	4.35	63817.44	4.80	63.82	1.32	{(N, 49260)}
22 公司	49025	671.51	2.83	29044	4.40	71261.10	4.85	71.26	1.42	{(N, 49025)}
23 該	48801	668.44	2.83	36654	4.56	104286.61	5.02	104.30	1.54	{(N, 29759), (D, 19008), (V, 1), (V, 34)}
24 還	48795	668.36	2.83	34565	4.54	98555.34	4.99	98.36	1.54	{(D, 47712), (V, 360), (P, 387), (V, 97), (V, 39)}
25 大家	48779	668.14	2.82	26220	4.56	103390.14	5.01	103.09	1.37	{(N, 48712), (N, 67)}
26 位	48647	666.33	2.82	34387	4.54	97846.01	4.99	97.85	1.41	{(N, 48111), (V, 1), (N, 295)}
27 教育	48418	663.19	2.82	24605	4.39	70011.95	4.85	70.01	1.41	{(N, 47378), (V, 1), (V, 1040)}
28 現在	47679	653.07	2.81	38774	4.59	110328.93	5.04	110.33	1.38	{(N, 47679)}
29 本	47378	648.95	2.81	33750	4.53	96033.46	4.98	96.03	1.61	{(N, 12306), (N, 31810), (D, 2038), (N, 389), (V, 65)}
30 提供	47256	646.65	2.81	36186	4.56	102737.31	5.01	102.74	1.57	{(V, 47256), (N, 24), (V, 1)}
31 知識	46941	642.96	2.81	30553	4.56	102596.50	5.01	102.59	1.22	{(V, 46940), (N, 1)}
32 也	46759	640.47	2.81	37139	4.57	109733.55	5.02	109.73	1.62	{(N, 3676), (N, 30388), (D, 8664), (V, 2851), (V, 247), (V, 867), (V, 59), (V, 7)}
33 去年	46737	640.17	2.81	32850	4.52	93472.57	4.97	93.47	1.57	{(N, 46737)}
34 社會	46611	638.44	2.81	29492	4.47	83917.60	4.92	83.92	1.42	{(N, 46611)}
35 種	46430	635.96	2.80	20231	4.31	57566.01	4.76	57.57	1.14	{(N, 46430)}

圖1. 中研院臺灣當代華語語料庫:詞彙頻率、語境變異、語義變異、詞類標示 \

## 本院覽號

24T-1120207

## 公告日期

2023-01-03

## 智財權狀態

know-how

## 應用範圍

- 1. 提供中文相關學術研究之實驗材料篩選語料
- 2. 提供發展學童、醫院臨床或長照機關病人及長者中文語言能力測驗題庫的參考語料
- 3. 提供學校教師編制給學童或是華語文第二外語學習者的中文教材參考語料
- 4. 其他中文字詞選取之可能應用

## 創作人

李佳穎