

# 文獻標注與文史研究

祝平次

清華大學中國文學系副教授

「古籍全文資料庫的回顧與展望」工作坊

中央研究院數位文化中心  
清華大學中國文學、科技部數位人文籌備小組

2014-06-05

(部分投影片為萊登大學何浩洋所提供，特此誌謝)

# 報告大綱

- 目的：示範文獻標注以後對於文史研究的幫助(如何擴充全文資料的用途)
- 文獻標注的用處
- 標注的方法
- 配合CBDB的自動標注系統Markus
- 個案示範：黃榦《勉齋集》和陳淳《北溪大全集》：計量、GIS與社會網絡比較
- 結論



# 文獻標注的用處

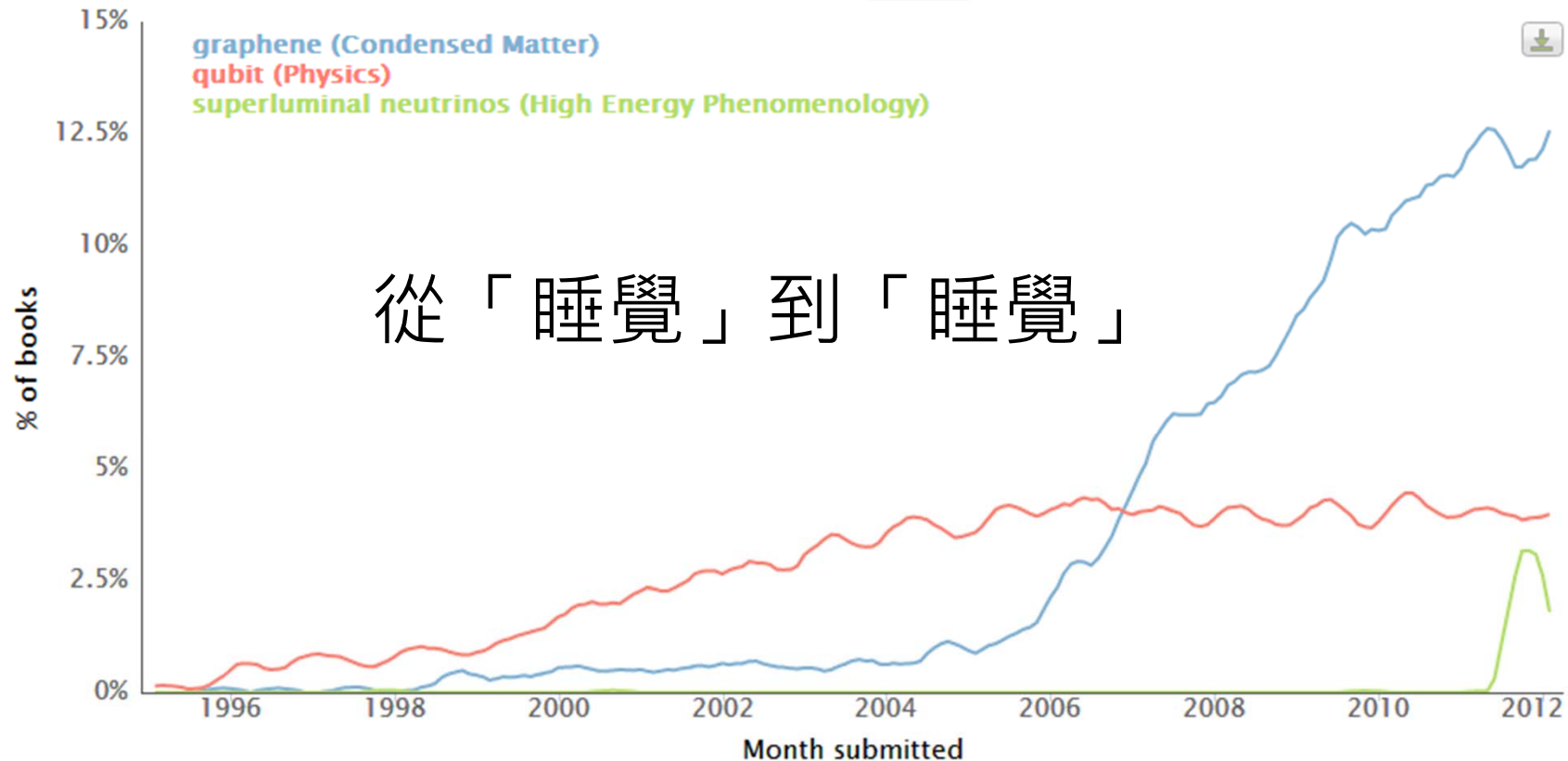
# 利用大量文本分析特定詞彙不同時代的情況

## bookworm ArXiv

Search for trends in hundreds of thousands of articles on [arxiv.org](http://arxiv.org)



graphene in Subject Class: Condensed Matter x +  
qubit in Category: Physics x +  
superluminal neutrinos in Subject Class: High Energy Phenomenology x + Search

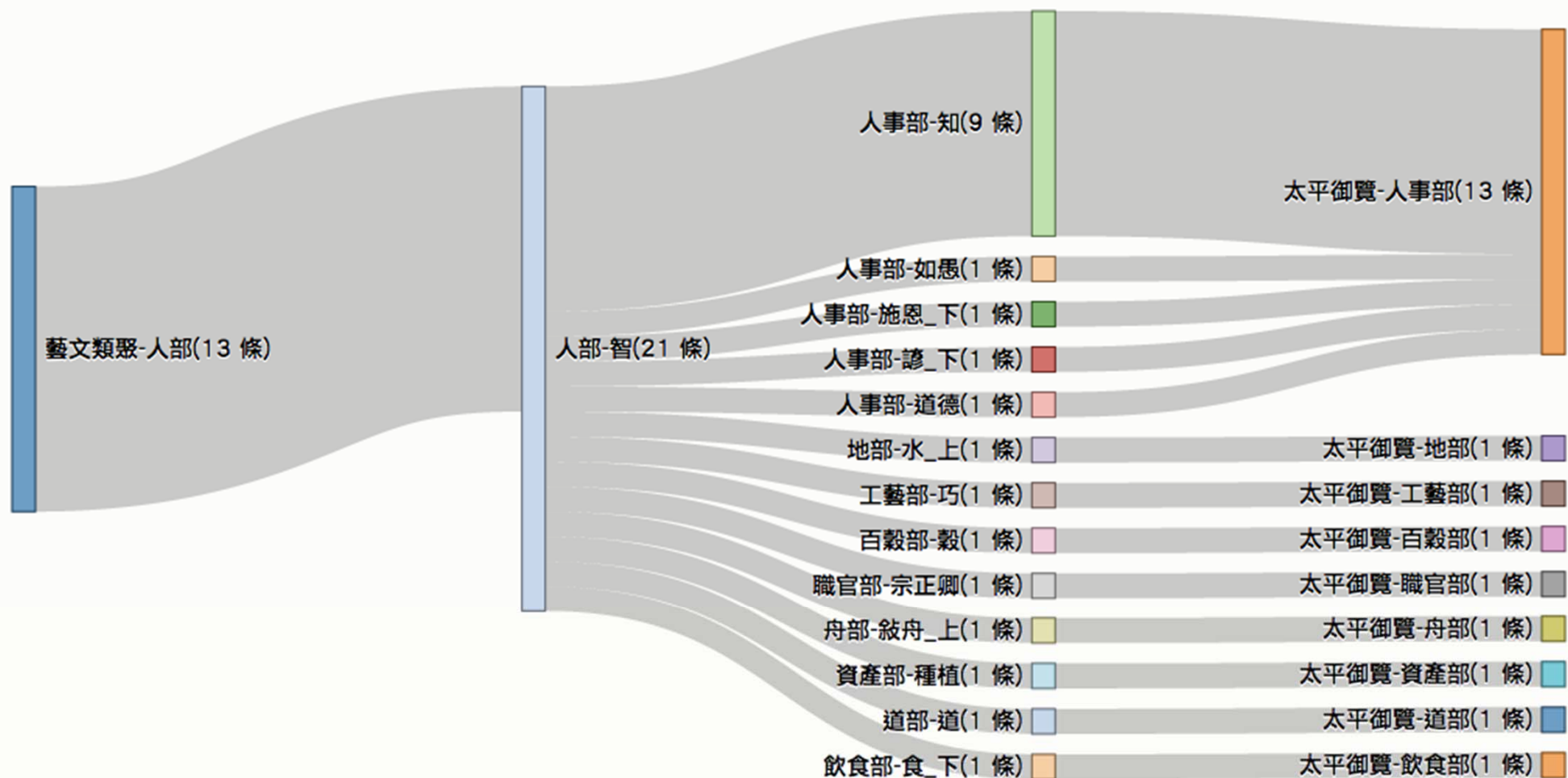


Culturomics (google n-gram viewer)



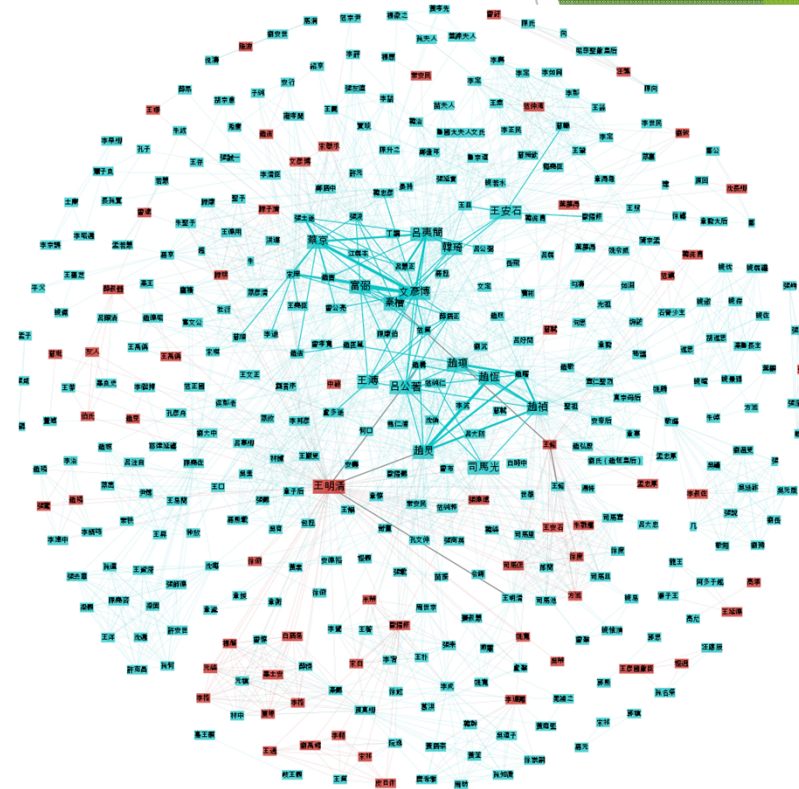
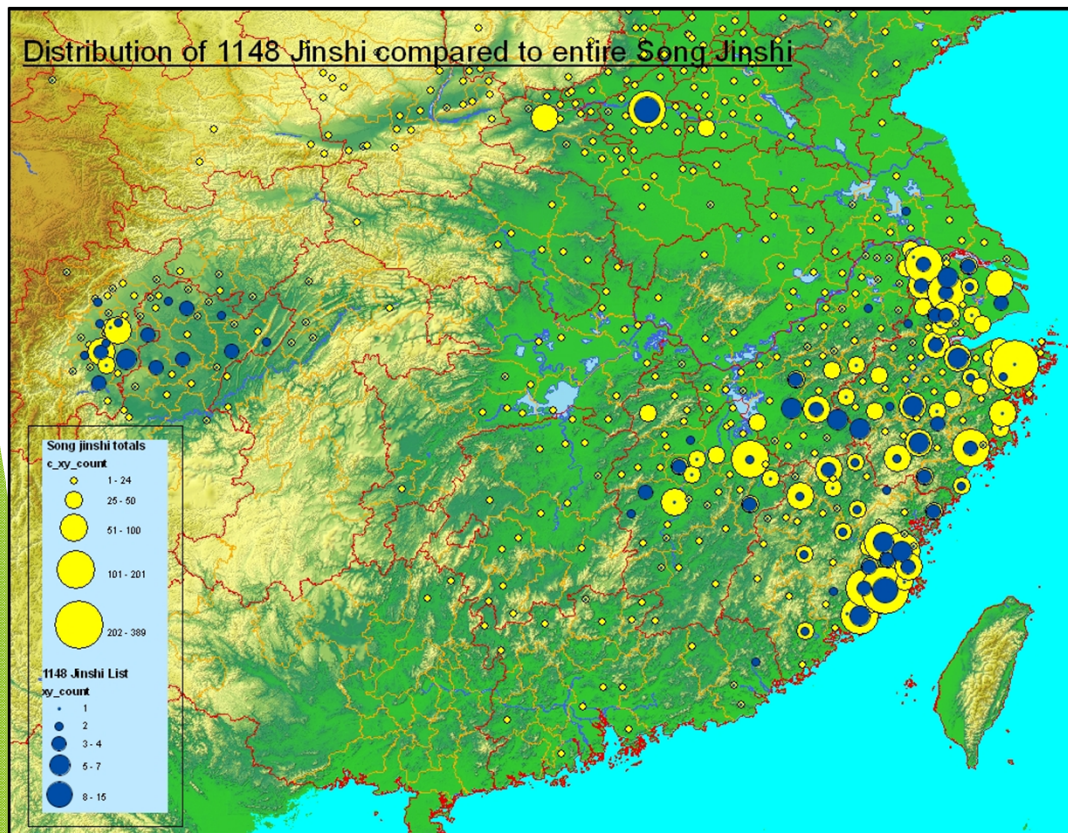
# 分析不同類書之間的相關性

## ◦ 統計





# 利用地理、人際網絡 觀察中國歷代人物

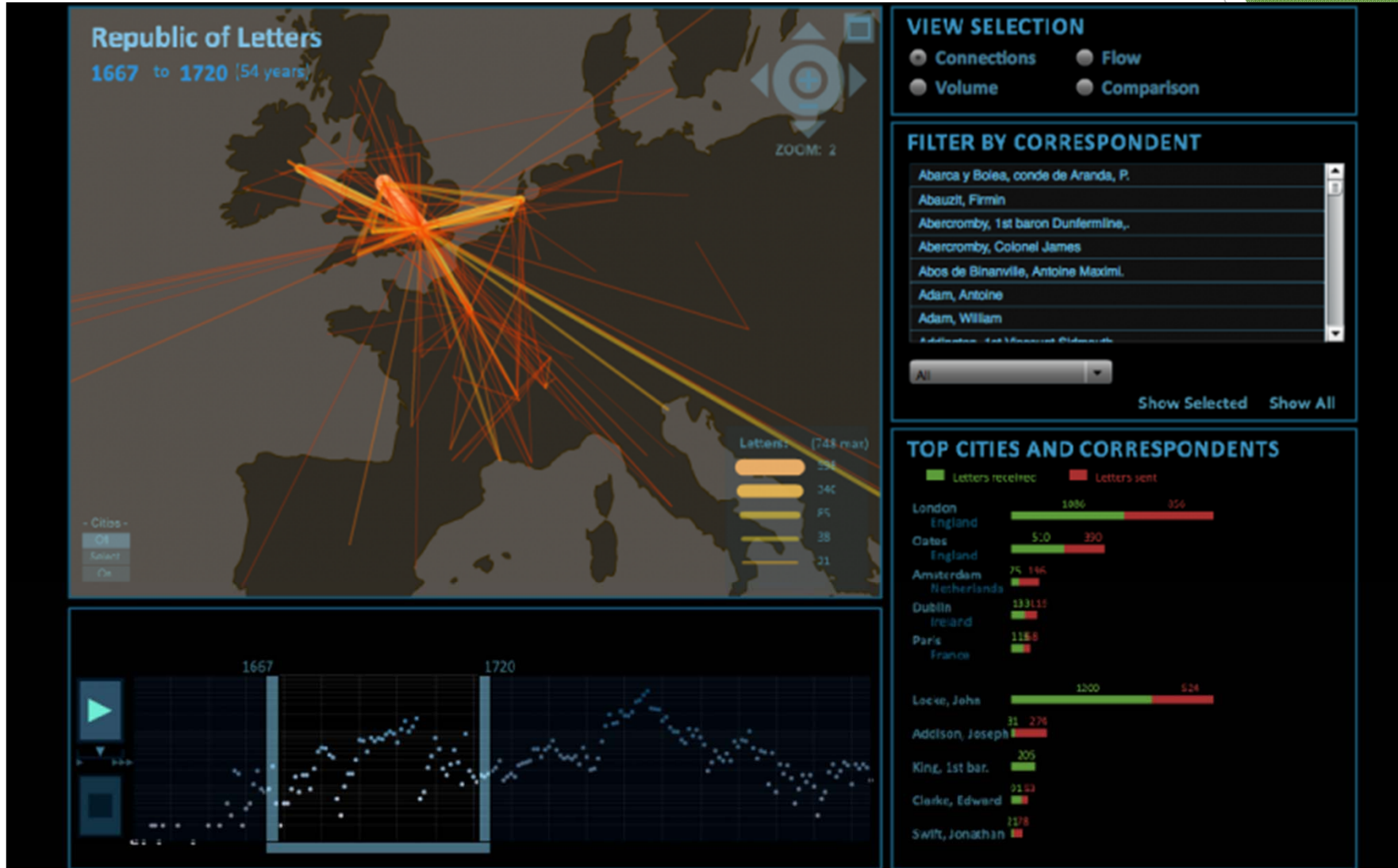


Communication and empire

China biographical database (CBDB)



# 利用地理、時間、空間分析 信件的流動

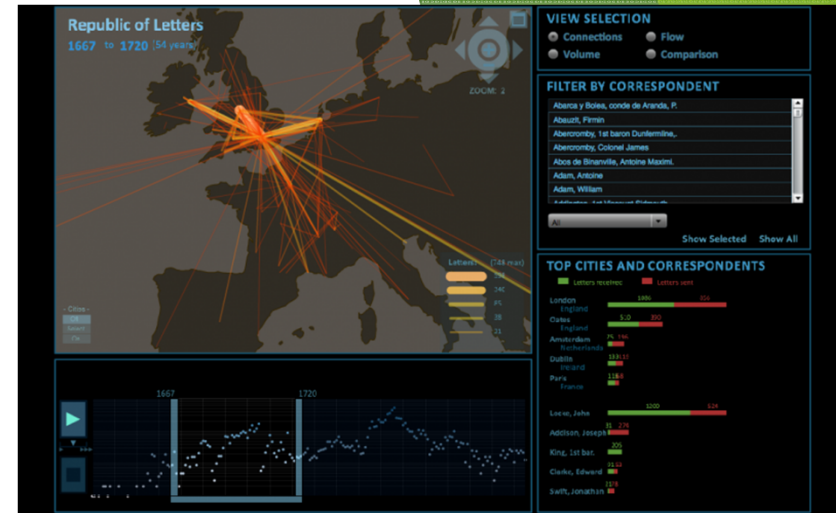
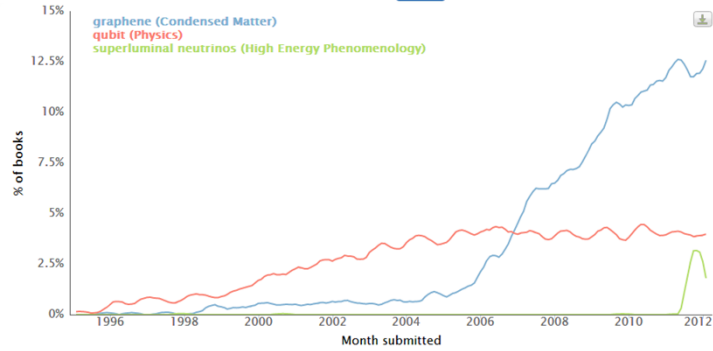


Republic of letters

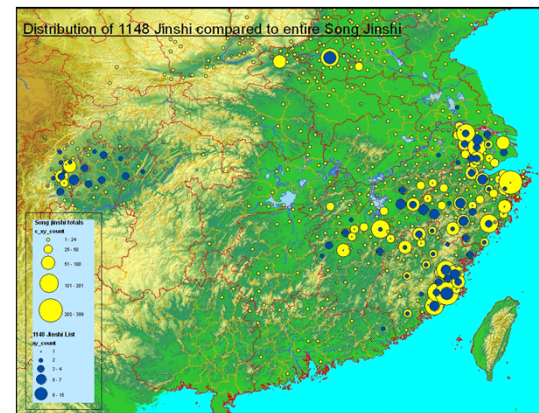
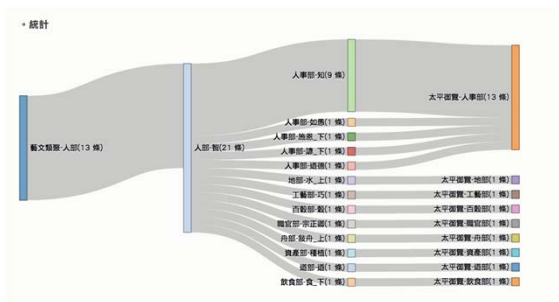
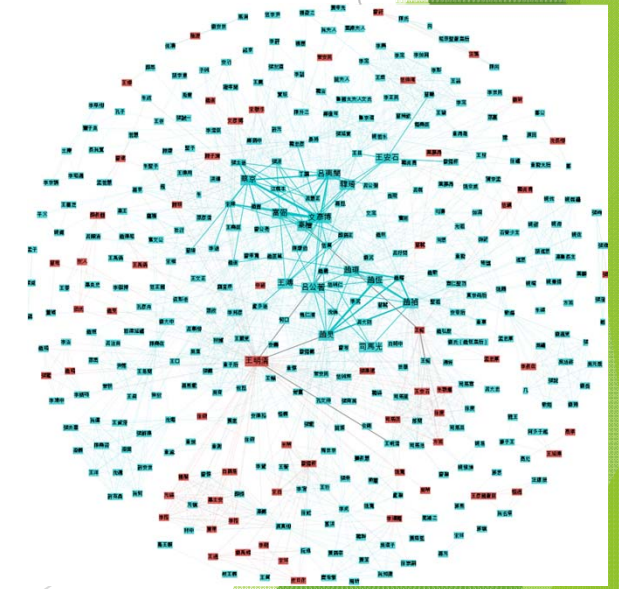
**bookworm ArXiv**

Search for trends in hundreds of thousands of articles on [arxiv.org](http://arxiv.org)

graphene in Subject Class: Condensed Matter x +  
 qubit in Category: Physics x +  
 superluminal neutrinos in Subject Class: High Energy Phenomenology x + **Search**



# 數據





# 宋會要輯稿

文獻集: 宋會要輯稿

文件檢索: 邵雍

測試版

最近幾次檢索	次數
邵雍	3
婺州	208
{TM:****}.all	18878
.all	80396

呈現模式: 檢索結果 | 逐篇檢視 | 詞頻與全文 | 列印全文 | 排列方式: 檔名順序 | 年代

找到筆數: 3 頁次: 1  
相關人名: 邵雍<sup>生平</sup> · 蘇洵<sup>生平</sup> · 魏漢津<sup>生平</sup> · 韓琦<sup>生平</sup> · 劉光祖<sup>生平</sup> · 常安民<sup>生平</sup>  
相關地名: 霸州 · 燕 · 潁州 · 衛州 · 河南 · 共城 · 文安

依 df 降序排列 | 依 t→q 降序排列

人名	Term t	(tf)	df	t→q
邵雍	生平	(3)	2	1.000
徐積	生平	(2)	1	0.250
蘇洵	生平	(3)	1	0.167
李壁	生平	(1)	1	0.091
劉光祖	生平	(1)	1	0.083
魏漢津	生平	(1)	1	0.083
常安民	生平	(1)	1	0.077
方時	生平	(1)	1	0.071
范祖禹	生平	(1)	1	0.033
韓琦	生平	(1)	1	0.033

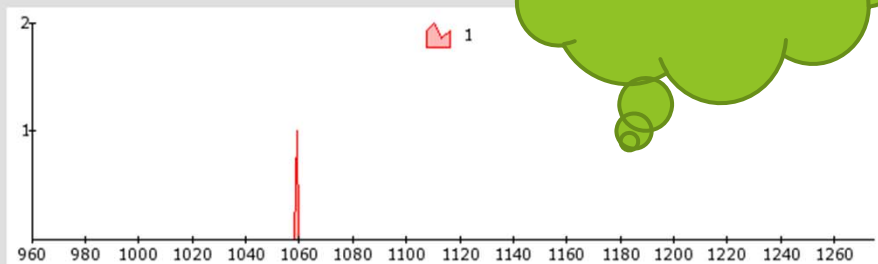
依 df 降序排列 | 依 t→q 降序排列

地名	Term t	(tf)	df	t
共城		(1)	1	0.083
文安		(2)	1	0.019
潁州		(1)	1	0.018
霸州		(1)	1	0.015
衛州		(1)	1	0.011
河南		(1)	1	0.004
燕		(1)	1	0.002

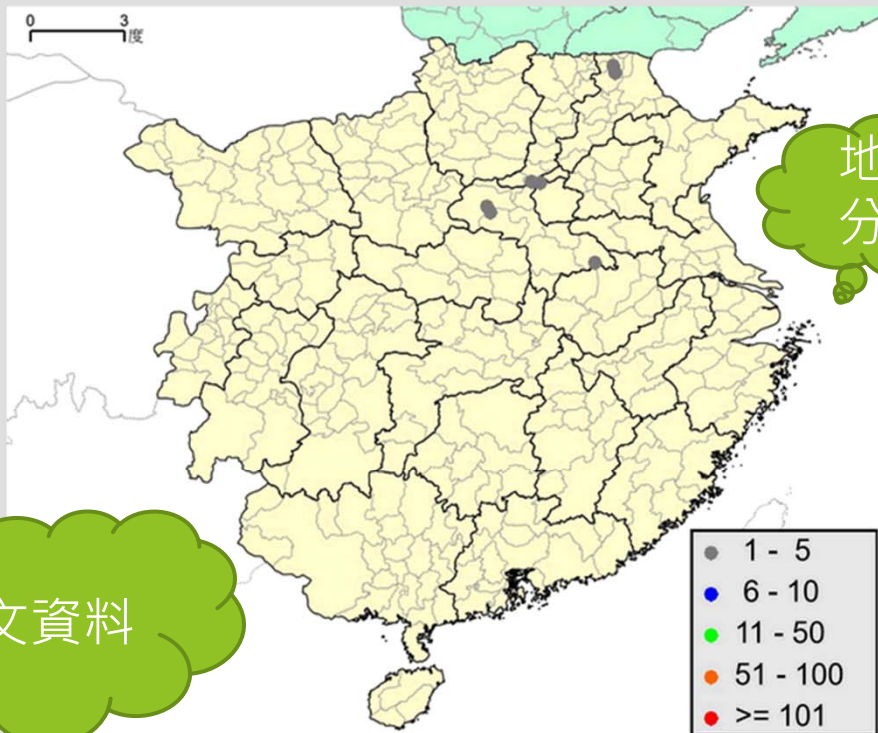
依 df 降序排列 | 依 t→q 降序排列

官名	Term t	(tf)	t
太常		(3)	0.005
團練推官		(1)	0.005
將作監主簿		(1)	0.005
宣德郎		(1)	0.011
光祿寺丞		(1)	0.009
著作郎		(1)	0.007
縣主簿		(2)	0.006
推官		(1)	0.005

年代分布



檢索結果分佈圖 (此圖僅計入含有年代資訊的文件; 橫軸: 西元年, 縱軸: 文件數量)



地理分布

全文資料

地名文件數量分佈圖 (此圖僅計入含有地名資訊的文件)





# 文獻標注的例子

Chinese Text Project Dict x

ctext.org/dictionary.pl?if=en&id=1102

應用程式 Gmail eBooks 3000 他的資料 共享 四庫 字典 文件 webmaster 論壇 論理學與宋代詩學中... 線上文獻 南方快報 文本分析

中文版 簡體

百諸家子 Chinese Text Project

Confucianism -> The Analects -> Xue Er -> 1

子曰：「學而<sub>4,2</sub>時習之，不<sub>3,3</sub>亦說乎？有朋自遠方來，不<sub>3,3</sub>亦樂乎？人不知而<sub>4,4</sub>不<sub>3,1</sub>慍，不<sub>3,3</sub>亦君子乎？」

Character Composition	Variants	Reading	Meaning	Hanyu	Kangxi	Cihai
子	子+0	學 巛 孛 孛 孛 只 zǐ 尸	(2.1) 古代對男子的美稱。也用以尊稱對方。 Respectful term for a man.	v2,p1006#06 p.277#01	p.395r1c01	
曰	曰+0	yuē 凵 卩	(1) 說。 Say, speak.	v2,p1482#02 p.502#01	p.648r3c03	
學	子+13	幸 学 xué 丁 卩 卩	(1) 學習。 Study, learn.	v2,p1019#11 p.280#28	p.403r3c05	
而	而+0	ér 儿	(4.2) 連詞：表示順承，相當於「就」、「才」。 Connective: and, then.	v4,p2810#01 p.961#18	p.1080r2c03	
時	日+6	豈 时 shí 尸	(8) 時常，時時。 Always, constantly.	v2,p1505#05 p.494#22	p.639r3c04	
習	羽+5	习 xí 丁 一	(1.2) 反復練習，復習。 Practice, try repeatedly.	v5,p3345#01 p.956#23	p.1074r3c09	
之	丿+3	出 zhī 出	(3) 第三人稱代詞：他、她、它。 Third person pronoun: him, her, it.	v1,p0043#01 p.82#04	p.47r2c02	
不	一+3	不 丕 堯 bù ㄅㄨˋ	(3.3) 否定副詞：表示反問，常與「乎」相呼應。 Negational adverb expressing rhetorical question: is not...?	v1,p0011#06 p.76#15	p.30r1c01	
亦	宀+4	亦 yì 一	(2) 也。 Also.	v1,p0281#03 p.88#09	p.78r2c02	
說	言+7	說 说 yuè 凵 卩	(8.1) Same as 「悅(1)」：高興，喜歡。 Happy, pleased, like.	v6,p3979#03 p.1164#08	p.1246r4c02	
乎	丿+4	虞 hū 厂 乂	(1.1) 語氣詞：表示疑問或反問，相當於「嗎」或「呢」。 Modal particle expressing questioning or a rhetorical	v1,p0036#08 p.82#07	p.47r3c03	



Gmail 企業版

利用來自 Google Apps 的自訂電子郵件，讓電子郵件看起來更專業

開始免費試用

# 文獻標注的用處

- ▶ 可以用不同的方式審視(有研究的意味)與展示(分享知識的意味)文獻
  - ▶ 包括轉化成其它媒體；
  - ▶ 包括顯題化文獻的某些特性；
  - ▶ 成為「數」據。
- ▶ 與其它數據的串連
  - ▶ 其為文獻也不同，其為數據也則一；
  - ▶ 同類文獻的整合(四庫全書、類書)；
  - ▶ 跨文類的整合(Google圖片搜尋、將音樂轉化成圖示)。
- ▶ 與電腦進行互動(人機互動)
  - ▶ 數位工具的產生：時間線圖、族譜
  - ▶ 對於人工研究過程的省視：年譜



# 可能的研究課題舉例

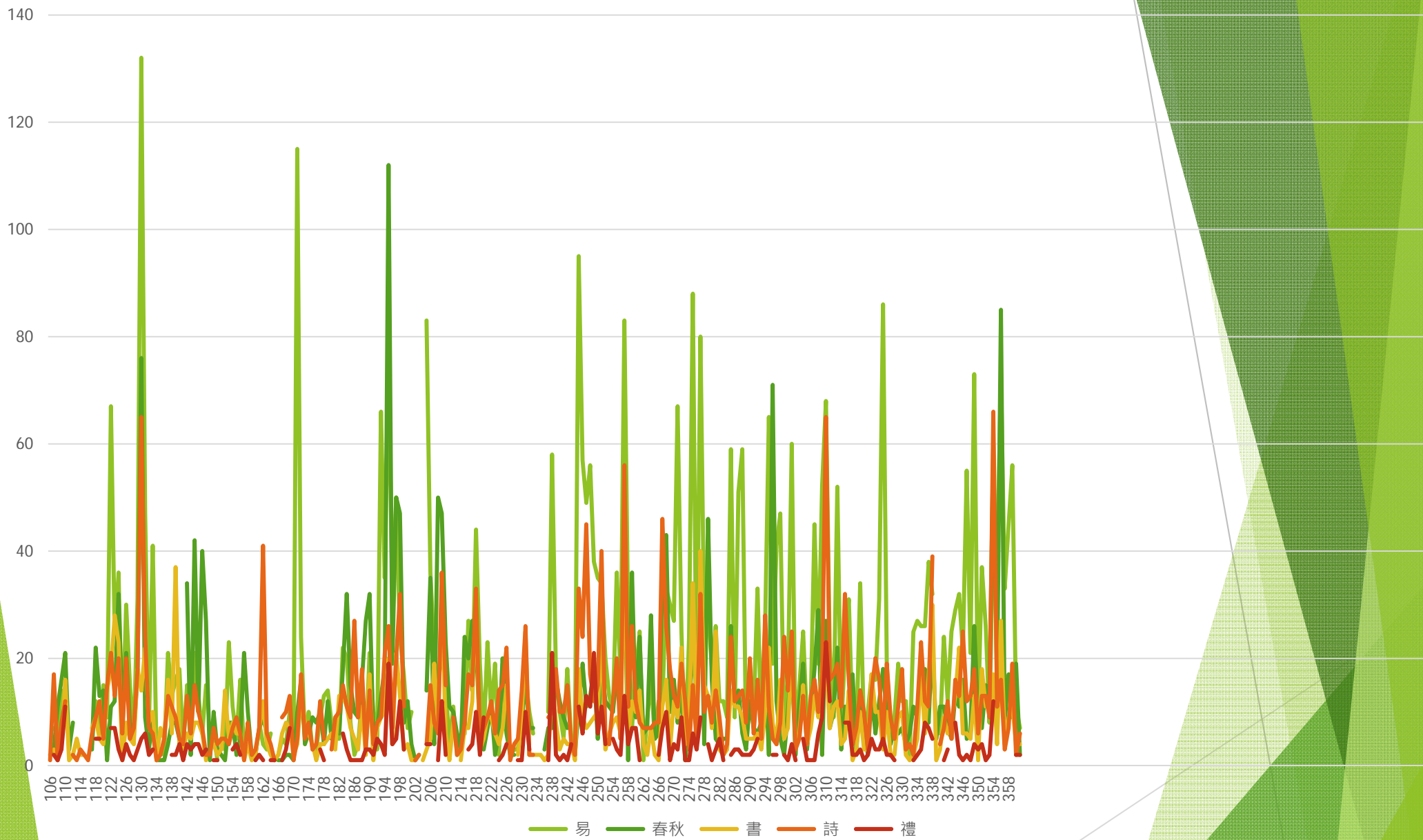
- 全宋文提到的文獻研究
  - 有多少還存在？
  - 有沒有時段前後的差別？
  - 不同文類的引用情形一樣不一樣？意味著什麼？
  - 不同作者(例如理學家、非理學)引用的情形一樣不一樣？意味著什麼？
  - 文獻彼此之間的關係(使用在同一個單位出現來建構文獻之間的關係)呈現什麼樣的狀況？意味著什麼？

書名	冊	ID1	ID1	篇ID	作者	卷數
祭法	1			qswi_001006	宋太祖	qswj_0010001
國風	1			qswi_001015	宋太祖	qswj_0010001
國風	1			qswi_001018	宋太祖	qswj_0010001
歸妹	1			qswi_001018	宋太祖	qswj_0010001
易	1			qswi_001018	宋太祖	qswj_0010001
三禮圖	1			qswi_001030	宋太祖	qswj_0010001
祠令	1			qswi_001066	宋太祖	qswj_0010002
重詳定刑統節文	1			qswi_001099	宋太祖	qswj_0010003
詩	1			qswi_001105	宋太祖	qswj_0010003
國風	1			qswi_001105	宋太祖	qswj_0010003
內則	1			qswi_001105	宋太祖	qswj_0010003
薤露	1			qswi_001106	宋太祖	qswj_0010003
禮	1			qswi_001106	宋太祖	qswj_0010003
書	1			qswi_001106	宋太祖	qswj_0010003



易6775	周官808	孝經391
春秋4538	坤777	史記360
詩4261	周禮729	尚書345
書3158	孟子655	坎341
禮1289	洪範540	震335
中庸1244	雅497	艮323
論語1105	記492	左氏322
乾1012	象401	孟321
大學1003	離397	禮記318
傳975	語394	頌295

圖表標題





易	6775	詩	4261
朱熹	405	朱熹	201
李綱	152	魏了翁	99
魏了翁	142	歐陽修	88
晁說之	132	周必大	72
李石	100	陳傅良	71
楊簡	97	劉克莊	68
楊萬里	95	晁說之	64
曾丰	92	司馬光	63
胡銓	92	薛季宣	60
黃裳	87	蘇軾	55
薛季宣	79	陳舜俞	50
陽枋	77	楊萬里	47
陳淳	74	真德秀	46
葉適	69	王安石	44
真德秀	68	蘇轍	43
劉克莊	66	周紫芝	40
薛軾	61	薛適	30

# 全文資料是資料庫嗎？

- ▶ 全文資料的資料庫性質是隱性的。
- ▶ 以前的人就有隱性資料庫的概念；也就是分類的概念—篇章節段、分類、類書。



# 標注的方法

如何讓不會中文的人  
分辨出人名、地名

紹興年間紹興在紹興喝紹興酒



# 如何讓不會中文的人 分辨出人名、地名

紹興年間紹興在紹興喝紹興酒

—

# 標注(Markup)

## 讓電腦了解文本

<nianhao>紹興</nianhao>年間

<name type="person"><family\_name>紹

</family\_name><given\_name>興

</given\_name></name>

在<name type="place">紹興</name>喝

<object type="alcohol">紹興酒</object>



# 過溫寄鞏縣宰吳秘丞 原注： 皇祐元年。

- ▶ `<head><date when="1049"/>`
- ▶ `<c>過</c><placeName><c>溫  
</c></placeName><c>寄  
</c><persName><placeName><c>鞏</c><c>縣  
</c></placeName><addName  
type="office"><c>宰  
</c></addName><surname><c>吳  
</c></surname><addName type="office"><c>  
秘</c><c>丞  
</c></addName></persName><hi>原注：  
<date notBefore="1049-02-05">皇祐元年  
</date>。 </hi></head>`

# 過溫寄鞏縣宰吳秘丞 原注： 皇祐元年。

- ▶ `<head><date when="1049"/>`
- ▶ `<hi>原注：<date notBefore="1049-02-05">皇祐元年</date>。</hi></head>`
- ▶ 兩個時間標記：  
一個標詩的著成年代，  
一個標詩的較準確的年代。



過溫寄鞏縣宰吳秘丞 原注：  
皇祐元年。

- ▶ `<placeName><c>溫</c></placeName>`
- ▶ 一個地名標記

# 過溫寄鞏縣宰吳秘丞 原注： 皇祐元年。

- ▶ `<persName><addName type="職"><placeName><c>鞏</c><c>縣</c></placeName><addName type="office"><c>宰</c></addName></addName><surname><c>吳</c></surname><addName type="官"><c>秘</c><c>丞</c></addName></persName>`
- ▶ 一個人名，其中還包括一個包括地名、職名的額外名、一個姓、一個官名的額外名。



# XML的標註原則

▶ 標註不得交叉，必須要有個根標註。

▶ `<root>`

▶ `<lloveyou>`

▶ `<lhateyou>`

▶ `</lhateyou>`

▶ `</lloveyou>`

▶ `<home>`

▶ `<room>Once upon a time</room>`

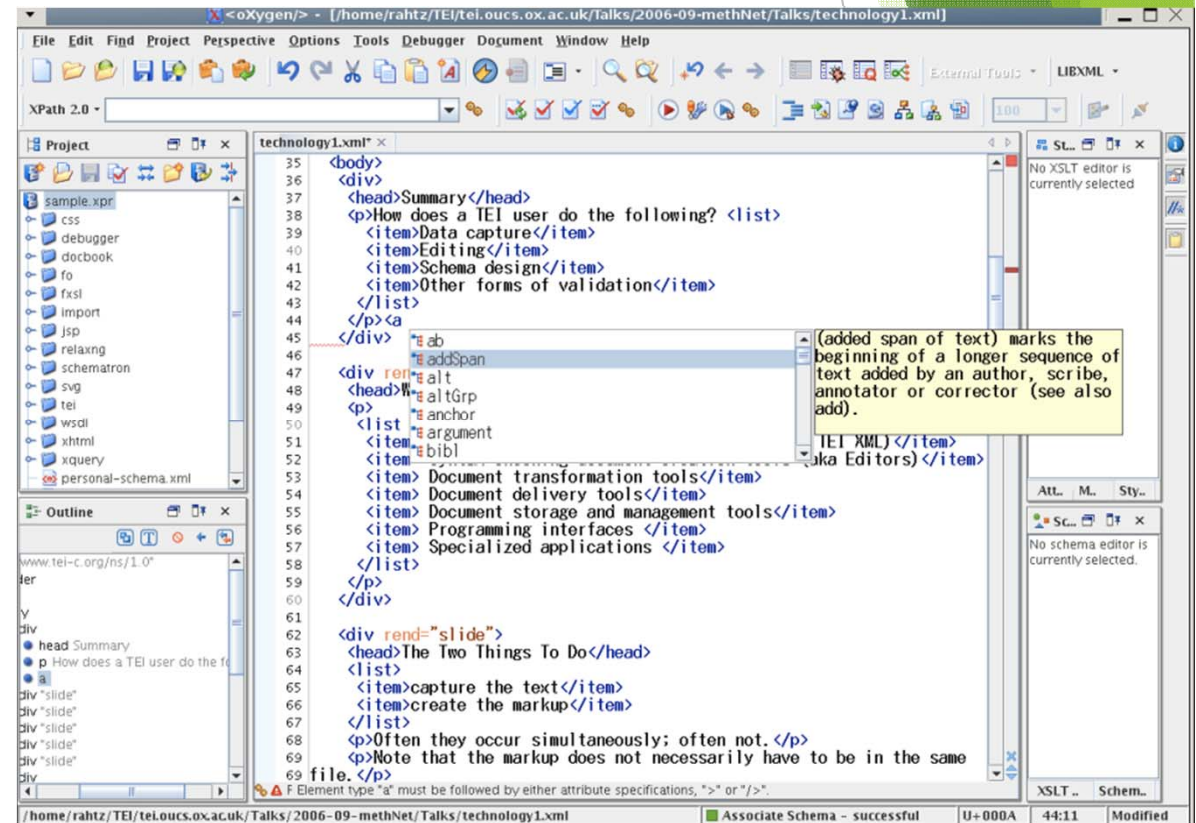
▶ `<room>There was a prince</room>`

▶ `</home>`

▶ `</root>`

# TEI - Text Encoding Initiative

- ▶ 第五版標準有1500 頁的說明
- ▶ 《佛教傳記文學視覺化平台》標記工作手冊(法鼓山)
- ▶ 標注工具：  
Oxygen XML editor





# 自動標注的必要性 (程式設計)

- ▶ 標注很有用，  
但人工標注很繁瑣，  
費時費力，  
很難進行。

# 配合CBDB的自動標注系統

## Markus



# Markus

The screenshot shows a web browser window with the URL `dh.chinese-empires.eu/beta/`. The browser's address bar and tabs are visible at the top. The page content includes a navigation bar with the text "MARKUS Communication and Empire" and an "About" link. The main heading is "MARKUS" in large, bold letters. Below the heading is a paragraph describing the application's features: "With MARKUS you can upload a file in classical Chinese (and perhaps in the future other languages) and tag personal names, place names, temporal references, and bureaucratic offices automatically. You can also upload your own list of key terms for automated tagging. You can then read a document while checking a range of reference works at the same time, or compare passages in which the same names or keywords appear. Or, you can extract the information you have tagged and use it for further analysis in our visualization platform and other tools." Below this text is a "Step 1" section with two blue buttons: "Upload a txt (UTF-8) or saved MARKUS file" and "Paste your txt here", separated by the word "or". At the bottom of the page, there is a footer with the text: "MARKUS was developed as part of the project 'Communication and Empire: Chinese Empires in Comparative Perspective,' funded by the European Research Council. All Rights Reserved." The browser's taskbar at the very bottom shows the file `antcon3.2.4w.exe` and a download icon with the text "顯示所有下載...".

Regi x 正蒙 x 課網 x 2014 x 祝平 x Goo x 15 國立 x ptc.c x Z 控管 x Z 台灣 x f 胡淑: x MAF x cP How x

dh.chinese-empires.eu/beta/

應用程式 Gmail eBooks 3000 他的資料 共享 四庫 字典 文件 webmaster 論壇 論理學與宋代詩學中... 線上文獻 南方快報 文本分析

MARKUS Communication and Empire About

# MARKUS

With MARKUS you can upload a file in classical Chinese (and perhaps in the future other languages) and tag **personal names**, **place names**, **temporal references**, and **bureaucratic offices** automatically. You can also upload your own list of key terms for automated tagging. You can then read a document while checking a range of reference works at the same time, or compare passages in which the same names or keywords appear. Or, you can extract the information you have tagged and use it for further analysis in our visualization platform and other tools.

**Step 1 :**  or

MARKUS was developed as part of the project "Communication and Empire: Chinese Empires in Comparative Perspective," funded by the European Research Council. All Rights Reserved.

antcon3.2.4w.exe 顯示所有下載...

# Markus (Mark us)

- ▶ 解決人工標注的問題：讓電腦進行自動標注。
- ▶ 利用CBDB(China Biography Database 哈佛大學、北京大學、史語所合作的中國歷代人物傳記資料庫)來進行人名、地名、官名、時間名的標註



# Markus

http://dh.chinese-empires.eu/beta/



Regi x 正蒙 x 課網 x 2014 x 祝平 x Goo x 15 國立 x ptc.c x Z 控管 x Z 台灣 x f 胡淑 x MAF x cP How x

dh.chinese-empires.eu/beta/automarkup.html?file=bxqdj.xml

應用程式 Gmail eBooks 3000 他的資料 共享 四庫 字典 文件 webmaster 論壇 论理学与宋代诗学中... 線上文獻 南方快報 文本分析

MARKUS Communication and Empire Save Back to last save Export HTML About

陳淳 2007-08-30  
unknown

Converted from a Word document

DOCX to TEI 2014-03-06T12:33:53Z **Ying-Izu Chu**

提要  
欽定四庫全書集部四

**北溪**大全集 別集類三宋

提要

臣等謹案：《**北溪**大全集》五十卷，《外集》一卷，宋·**陳淳**撰。淳有《**北溪**字義》已著錄，其生平不以文章名，故其詩其文皆如語錄。然淳於朱門弟子之中最為篤實，故發為文章，亦多質樸真摯無所修飾。元·王環翁序以為「讀其文者當如布帛菽粟，可以濟平人之飢寒，苟律以古文律度，聯篇累牘，風形霧狀，能切日用乎否」云云。是雖矯枉過直之詞，要之，儒家實有此一派，不能廢也。又淳以朱子終身與**陸九淵**如水火，故生平大旨在於力申儒釋之辨，以鍼砭**金谿**一派之失。集中如《道學體統》等四篇，《似道》、《似學》二辨，皆在嚴陵時所作，反覆詰辨，務闡明**鷺湖**會講之緒論，亦可謂堅守師傳、不失尺寸者矣。集為其子渠所編，末有《外集》一卷，載其祭文、誌銘、敘述五篇，亦渠所輯，附淳祐戊辰郡倅**薛季良**為錢板**龍溪書院**，歲久散佚。元·**至元**乙亥、明·宏治庚戌又兩經翻刻，今所傳者蓋猶宏治本云。**乾隆**四十四年三月恭校上。

總纂官臣紀昀臣陸錫熊臣孫士毅

總校官臣陸費墀

**北溪**大全集原序

道之顯者謂之文，惟辭類亦通語除怪字乎哉？《文經》、《禮記》也；四書、日月也；各口成文，下筆成書，惟《詩》、

標注的東西

查閱工具

姓名 別名 年號 地名 官職

markup selection

CBDB CHGIS ZDict Wikipedia

Search Search here

# 在標注之後

個案示範：

黃榦《勉齋集》

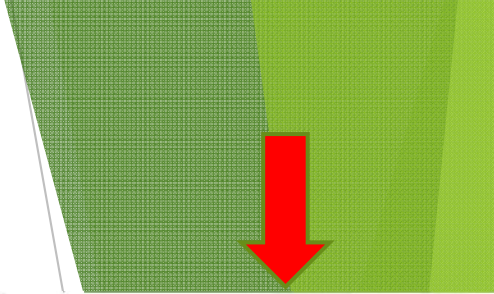
陳淳《北溪大全集》



# 儲存結果

陳淳			黃榦		
Type	Total	Unique	Unique	Total	Type
<u>fullName</u>	1018	329	624	1638	<u>fullName</u>
<u>partialName</u>	41	11	13	23	<u>partialName</u>
<u>placeName</u>	604	230	496	2072	<u>placeName</u>
<u>officialTitle</u>	651	162	436	2010	<u>officialTitle</u>
<u>nianhao</u>	460	75	98	527	<u>nianhao</u>

# 顯示結果



bxdqj.xml.data.html - Excel

檔案 常用 插入 版面配置 公式 資料 校閱 檢視 增益集

A1 type

	A	B	C	D
1	type	tag	id	
2	fullName	陳淳	10882136449	
3	fullName	北溪	10882152937	
4	fullName	北溪	10882152937	
5	fullName	陳淳	10882136449	
6	fullName	北溪	10882152937	
7	officialTitle	弟子		
8	fullName	陸九淵		3632
9	placeName	金谿		
10	fullName	鷺湖		68133
11	nianhao	淳祐		
12	fullName	薛季良		11345

bxdqj.xml.data



# 與CBDB連結

計量、GIS與社會網絡比較

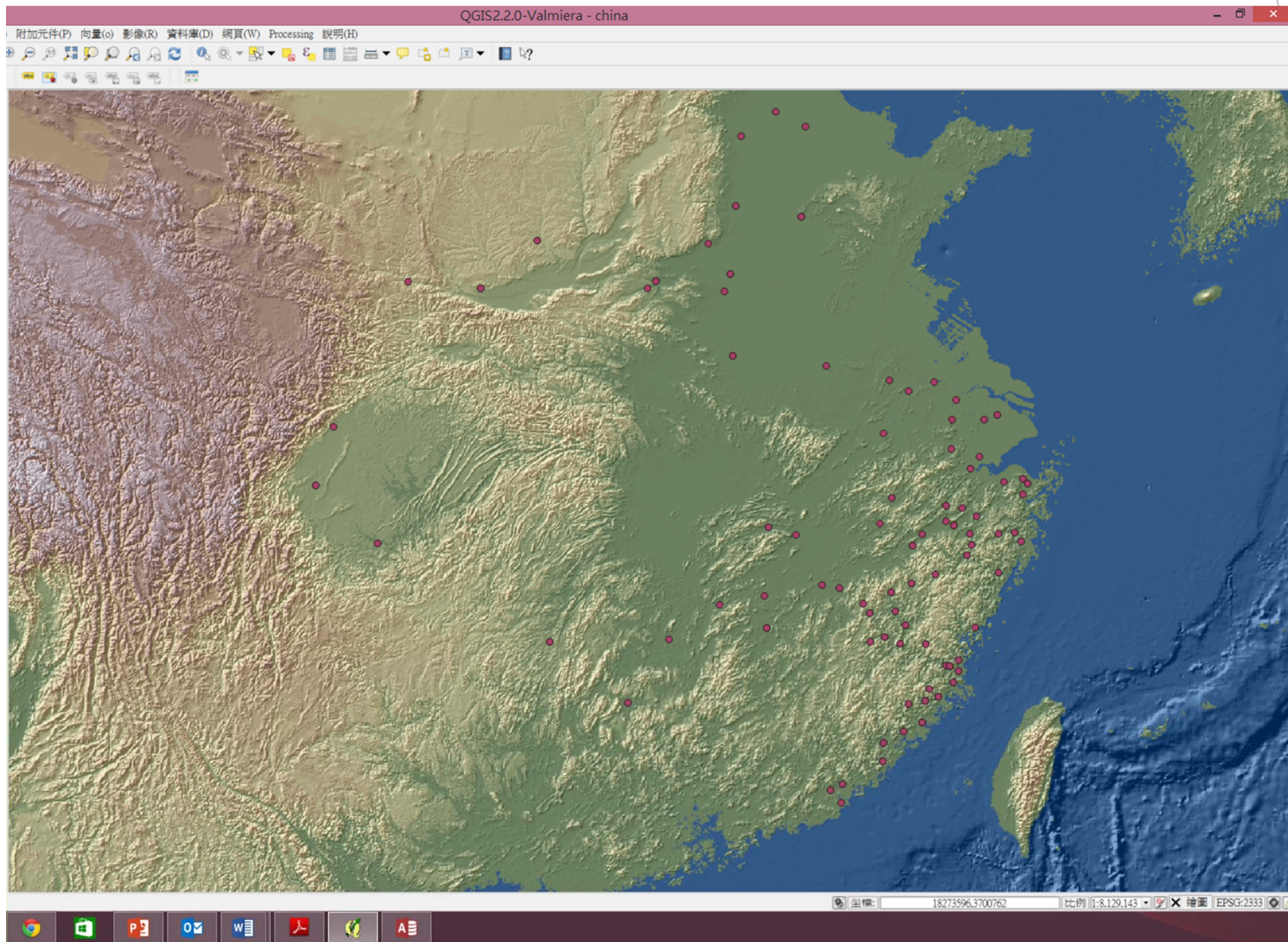
# 比較顯示結果



3	列標籤	計數 - tag	列標籤	計數 - tag	1	列標籤	計數 - tag	列標籤	計數 - tag
4	10882	181	10882	181	2	11134	199	11134	199
5	3257	73	北溪	124	3	26991	36	直卿	2
6	43538	44	安卿	3	4	7123	35	勉齋	85
7	7164	16	陳淳	54	5	3257	26	黃榦	112
8	10892	9	3257	73	6	20027	11	26991	36
9	504	8	文公	20	7	26590	11	陳安國	36
10	25123	7	朱子	1	8	11127	11	7123	35
11	19162	6	朱文公	6	9	11827	10	陳安節	35
12	46581	5	考亭	6	10	42629	10	3257	26
13	42981	5	晦翁	18	11	12408	9	朱文公	10
14	3718	5	晦庵	7	12	38978	8	朱熹	1
15	51981	5	紫陽	15	13	10892	7	考亭	9
16	49517	4	43538	44	14	14600	7	晦翁	4
17	42758	4	伯澡	1	15	25535	7	紫陽	1
18	28168	4	陳伯澡	40	16	40880	6	邀翁	1
19	52892	3	陳沂	3	17	40723	6	20027	11
20	11134	3	7164	16	18	7164	6	李貫之	10
21	48607	3	南軒	15	19	11904	6	李道傳	1
22	1216	3	南軒先生	1	20	51938	6	26590	11
23	52959	3	10892	9	21	42031	5	伯量	1
24	20041	3	陳宓	4	22	30454	5	胡伯量	10
25	53695	3	復齋	5	23	2129	5	11127	11
26	22959	3	504	8	24	22389	5	潘謙之	11
27	37813	3	濂溪	7	25	37056	5	11827	10
28	42755	2	濂溪先生	1	26	22487	5	屏山	10
29	40992	2	25123	7	27	49814	5	42629	10
30	47602	2	廖子晦	3	28	22524	5	存齋	6
31	11871	2	槎溪	4	29	3597	5	存齋先生	1
32	42702	2	19162	6	30	44600	4	林公度	3
33	12165	2	李公晦	6	31	16209	4	12408	9
34	44128	2	46581	5	32	51319	4	張潛	9
35	14421	2	郭子從	5	33	7055	4	38978	8
36	3105	2	42981	5	34	44128	4	楊志仁	8
37	15206	2	一之	2	35	10304	4	10892	7
					36	47842	4	陳宓	1
					37	12053	4	陳師復	6
					38	16039	4	14600	7

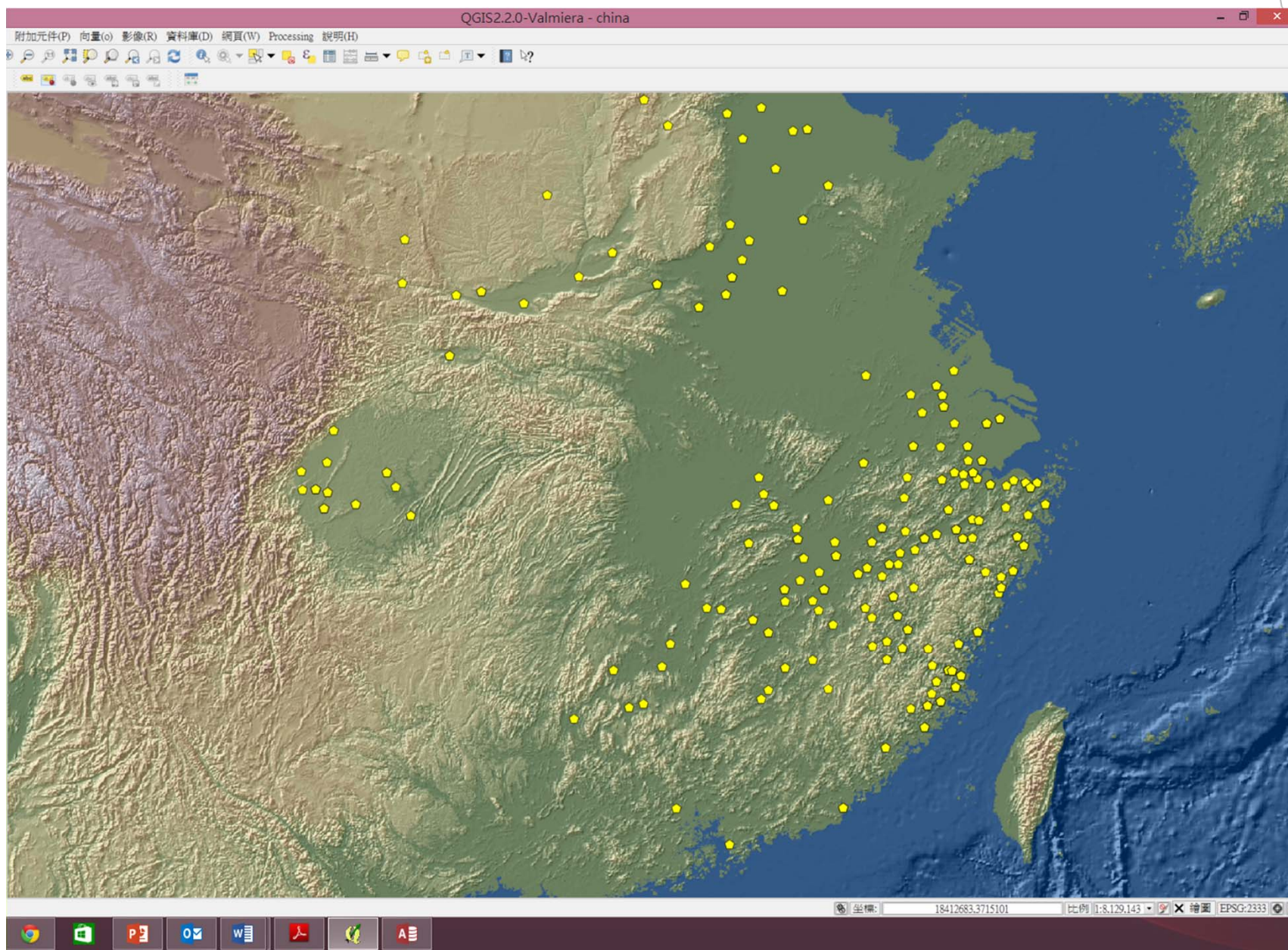


# 陳淳北溪大全集中人名的地理分布



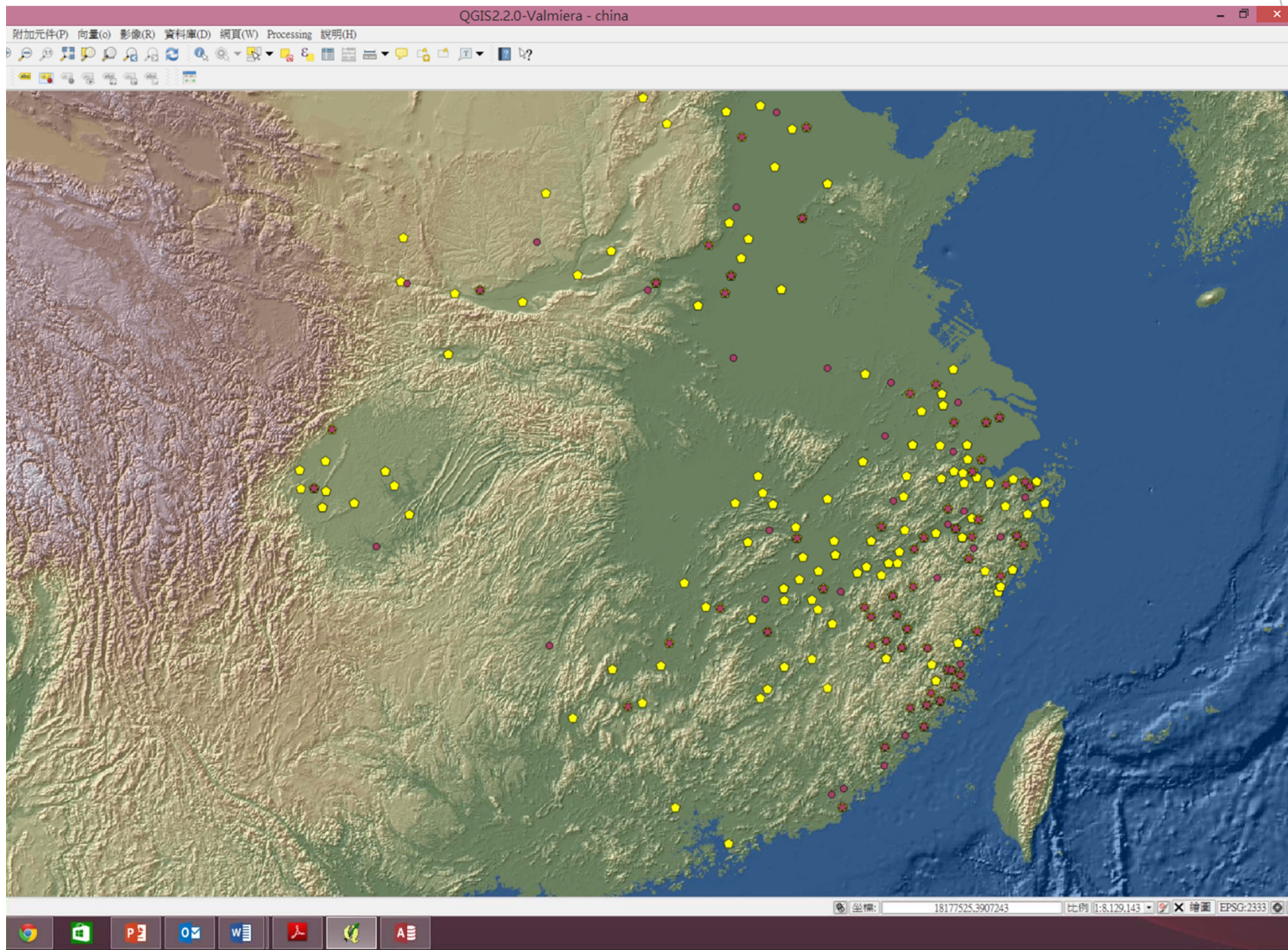


# 黃榦勉齋集中人名的地理分布



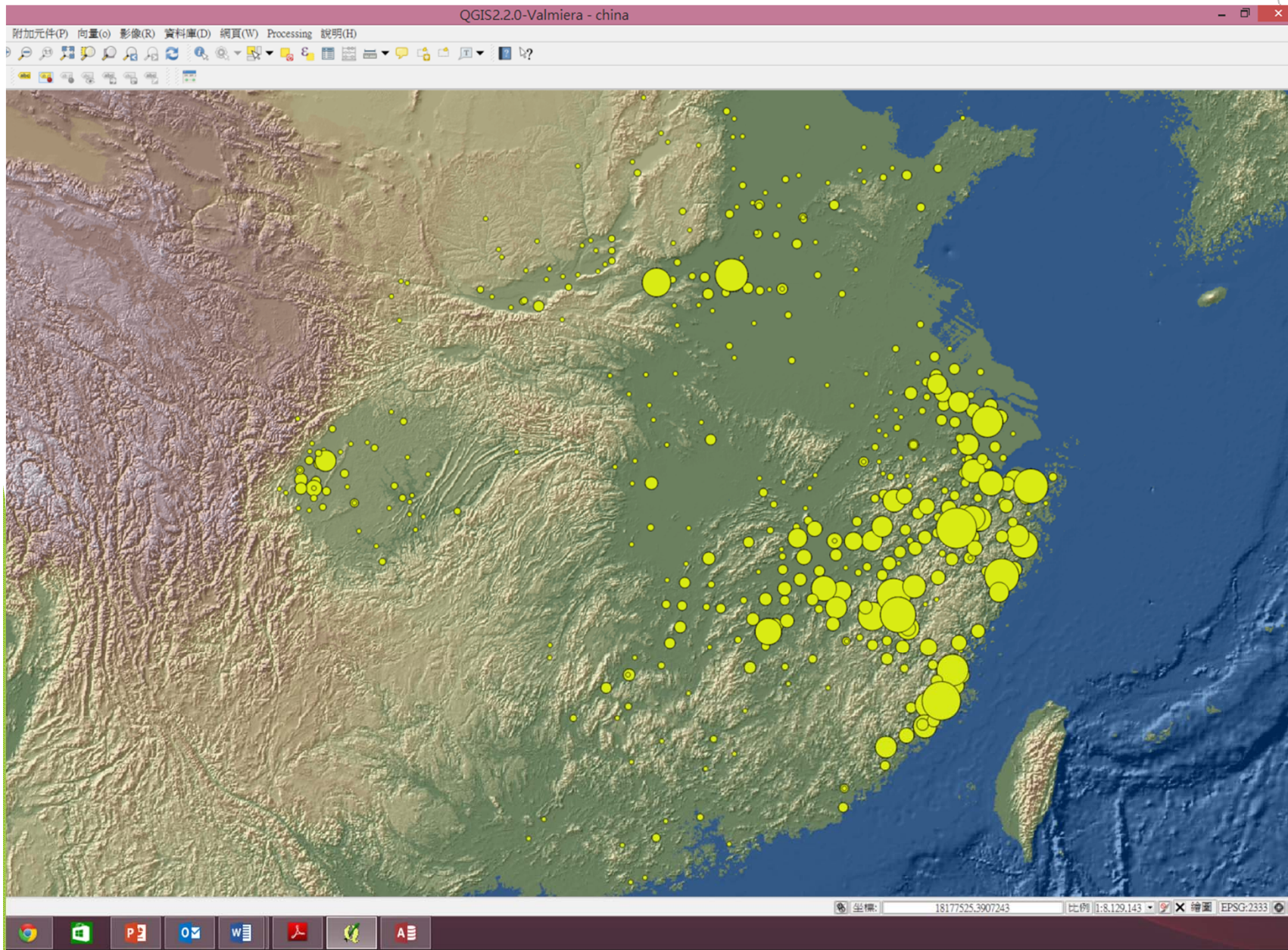


# 集中人名地理分布比較



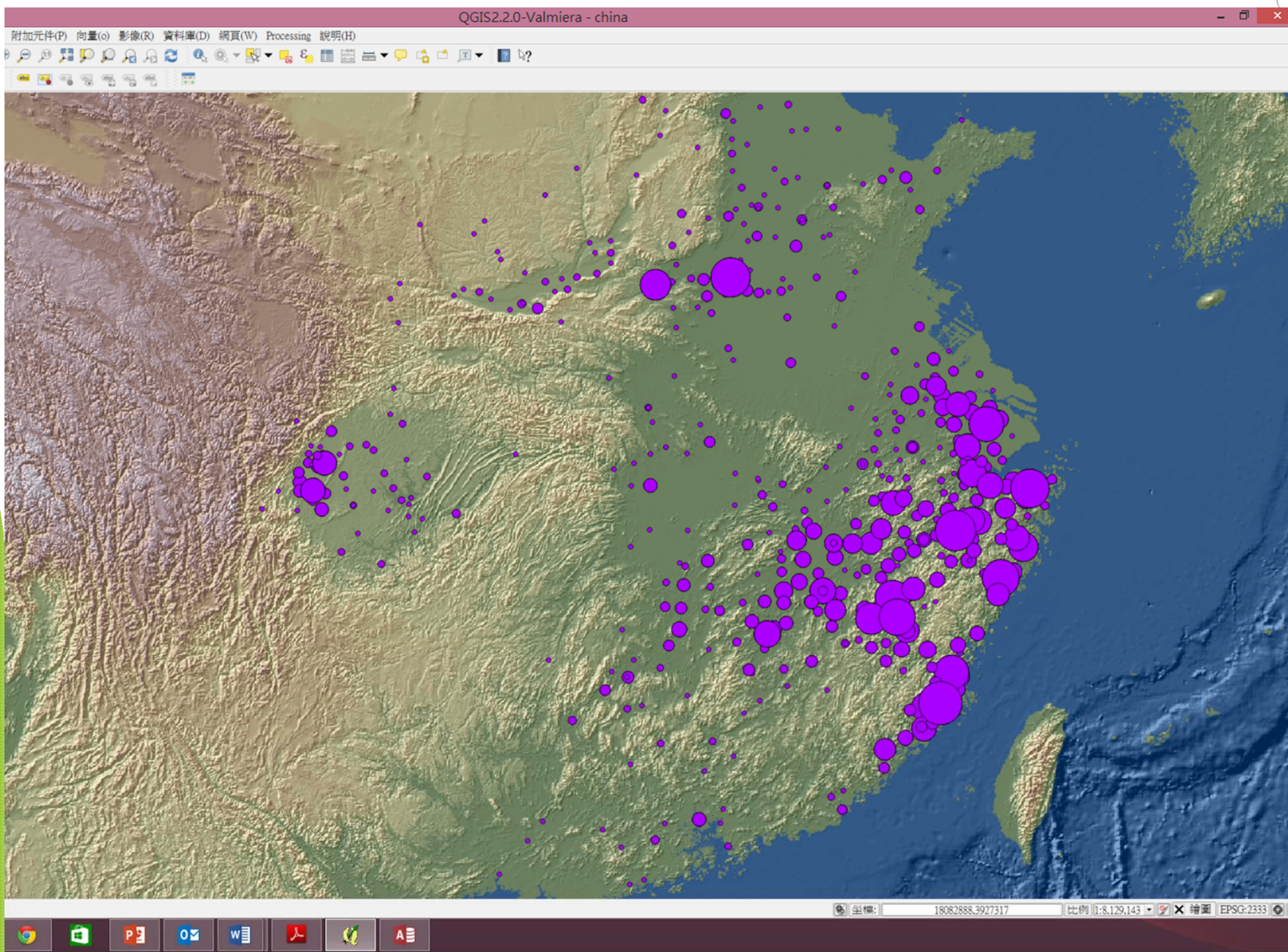


# 陳淳集中人名社會網絡地理分布



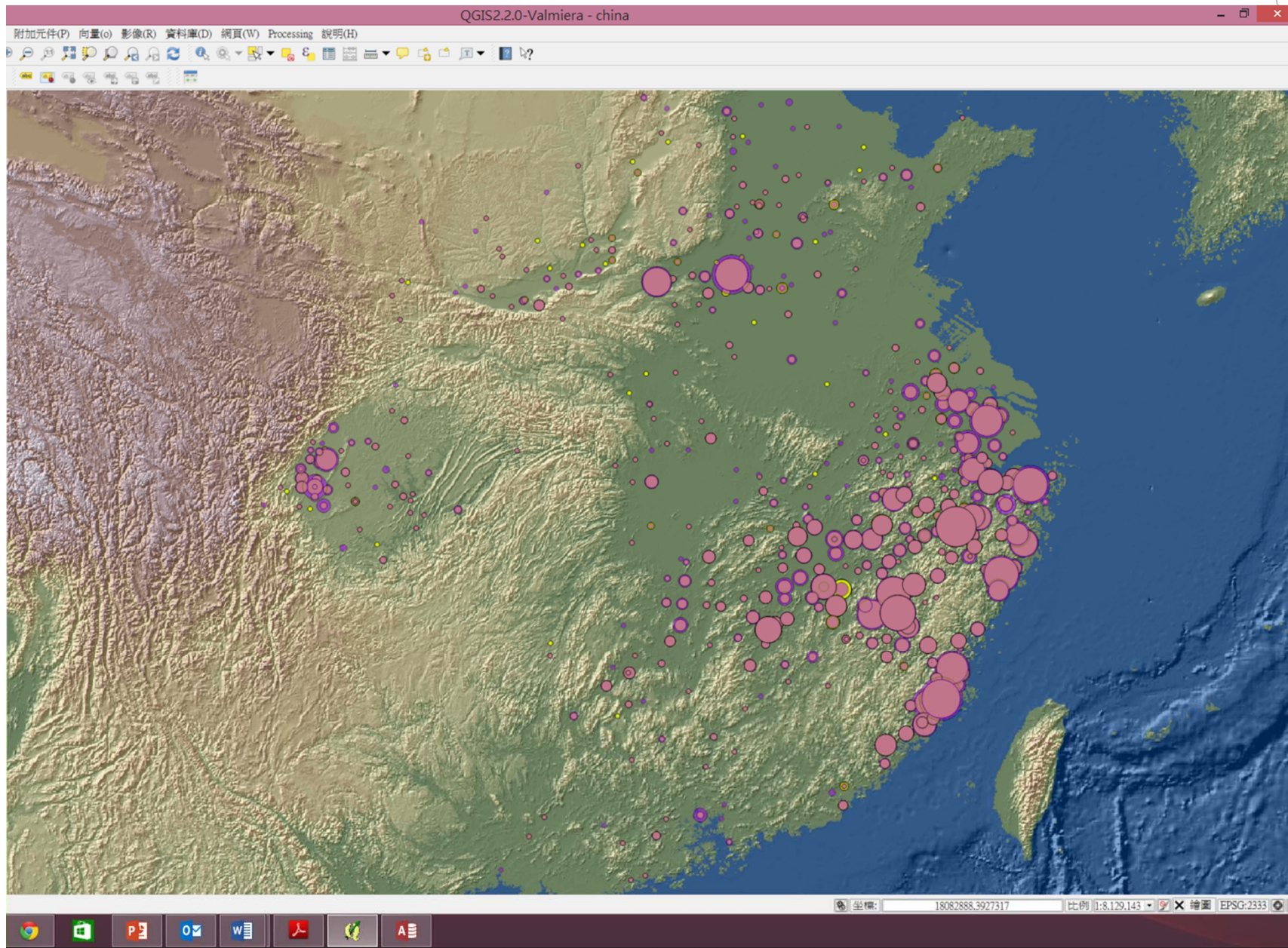


# 黃軫集中人名社會網絡地理分布

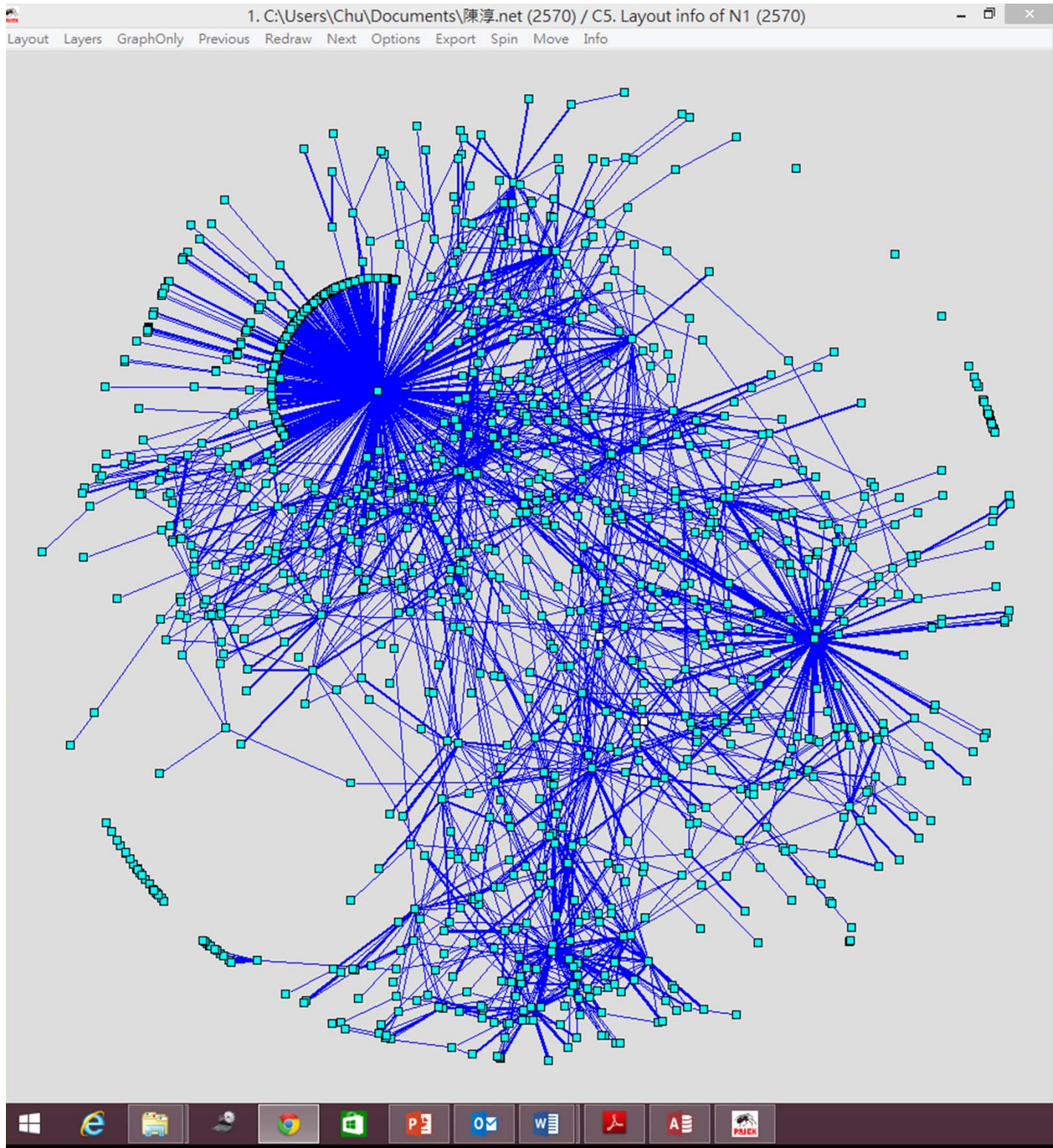




# 兩人集中人名社會網絡地理分布



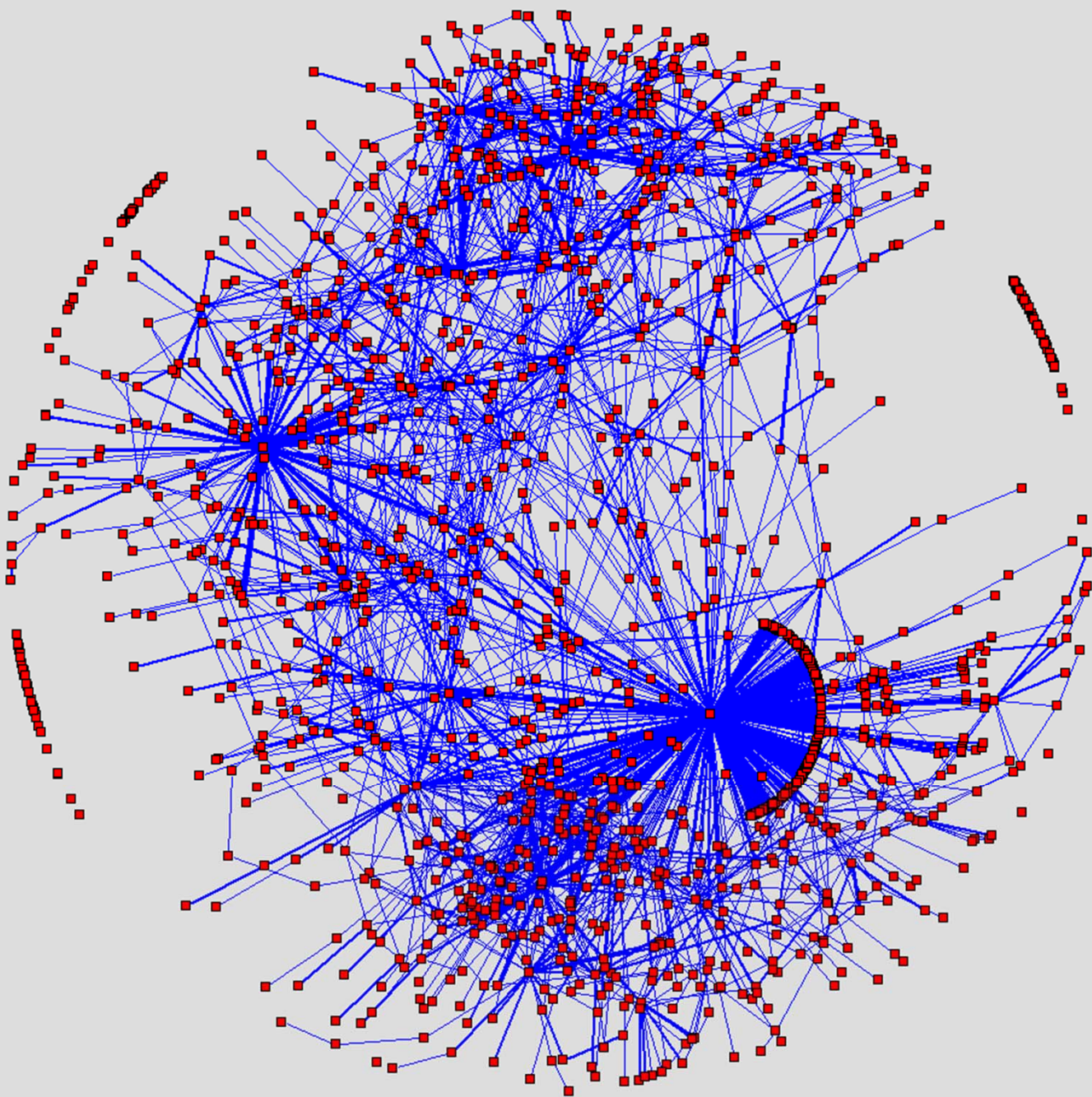




# 陳淳集中人名學術政治網絡





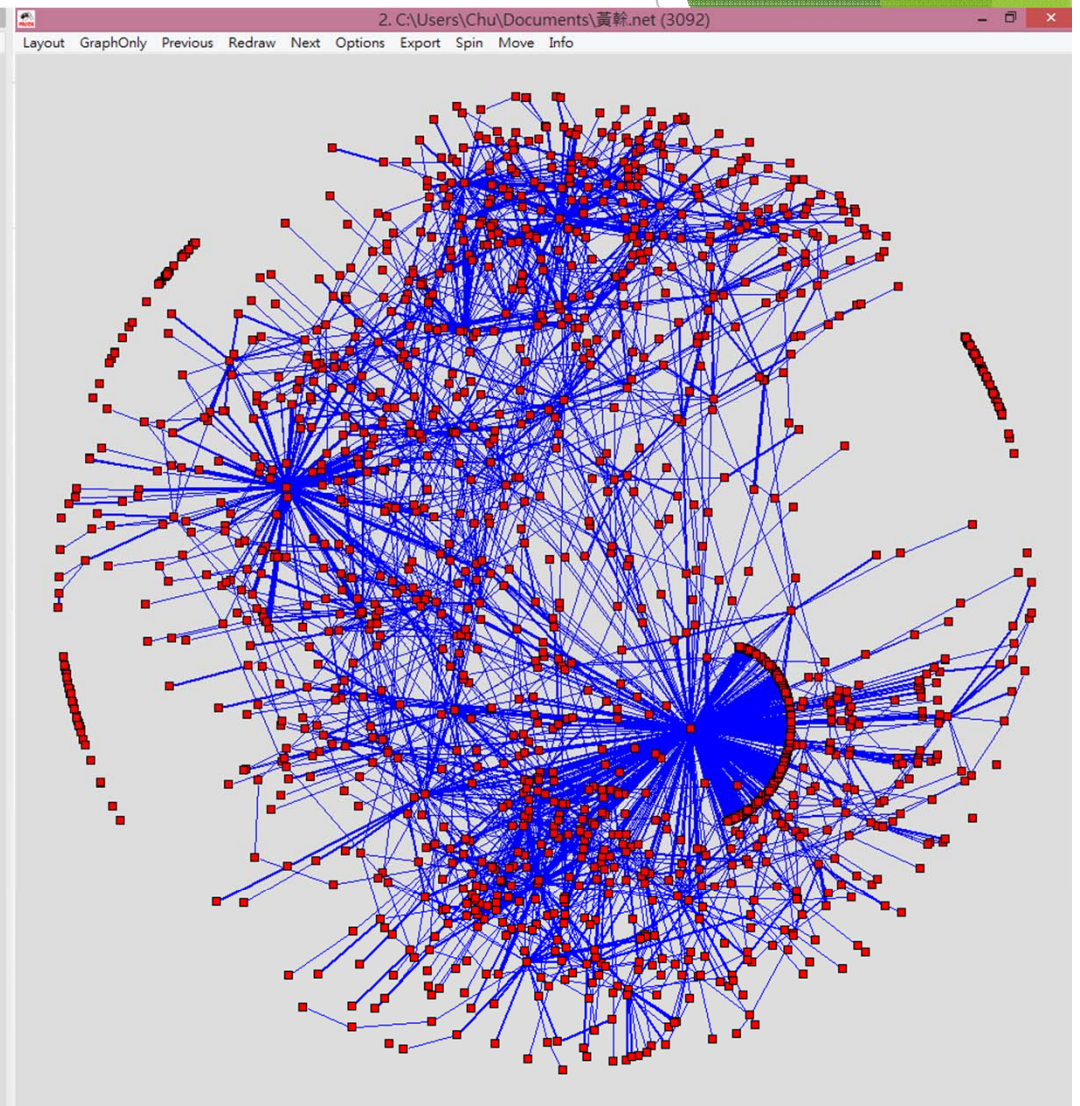
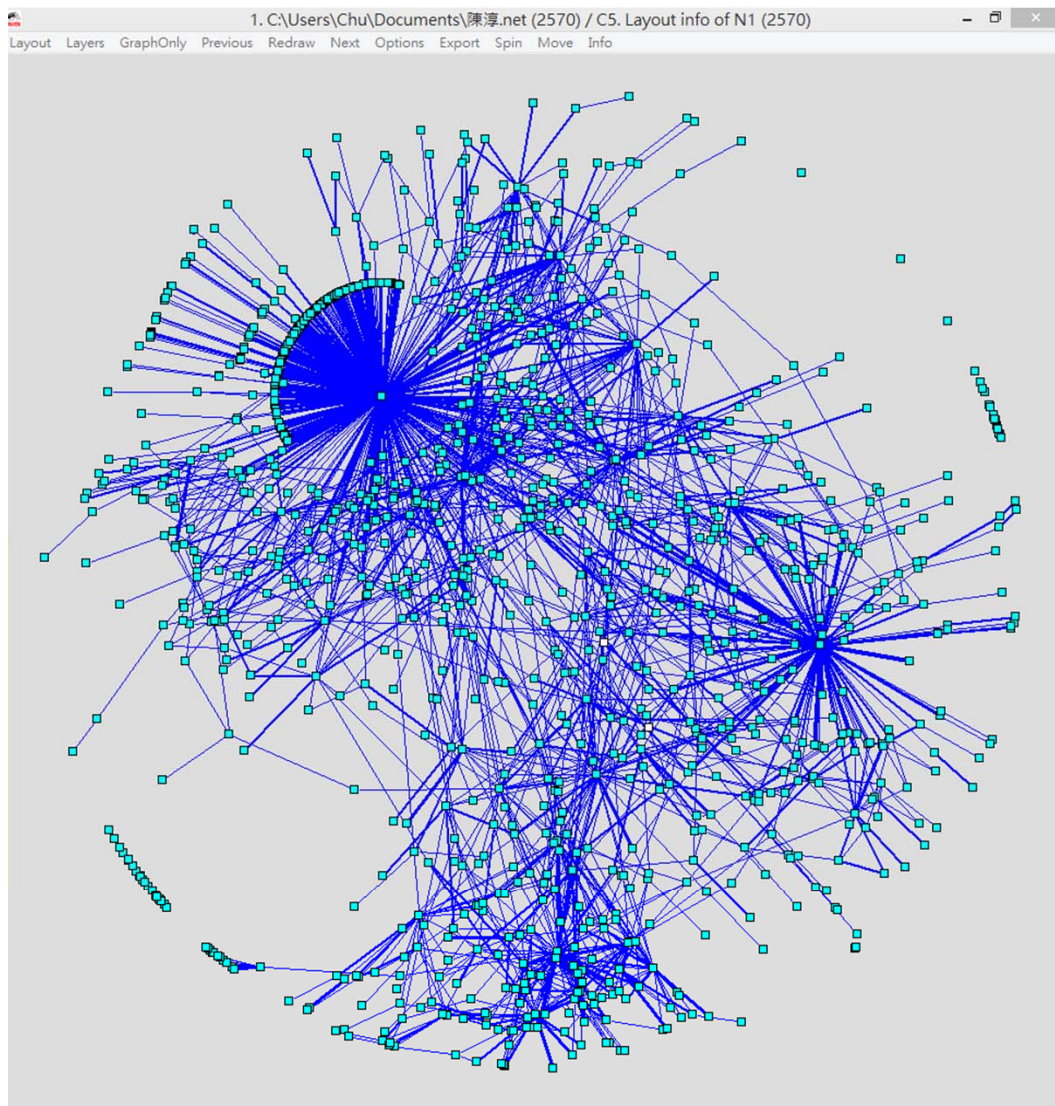


黃軫集中人名學術政治網絡





# 比較





# 結論



- 一般的全文資料經過XML/TEI標注後可以把原來潛在文本中的資料庫結構以及數據性質明示化。
- 明示化資料庫結構後，可以有各式各樣利用全文資料的方式，亦即它的利用價值會一下子提升很多。
- 但全文資料各式各樣利用的前提是：研究者要能取得全文資料。  
因為不能取得全文資料的話，其利用總是有限的。這樣子也就無法促成對於相同資料進行不同數位解讀的經驗、也難以促成數位人文研究群體的出現。也會使得有全文資料和沒有全文資料變成一種學術不平等。

# 置入性行銷

- ▶ <https://sites.google.com/a/ptc.cl.nthu.edu.tw/dhintaiwan/>
- ▶ <http://tinyurl.com/dhintaiwan>



終