

Project 1

Group 12

施昱竹 104078502 莊翊鈞 104078514

Outline

- Goal
- Tool
- Framework
- Code
- Index of word cloud
- Word cloud of each department

Goal

- Analytics goal:

Generating word cloud for each department.
(47 departments in total)

- Business goal:

For senior high students or researchers in different fields could have a big picture that what is the main topic in each department.

Tool

1. Python for text-mining
2. Word cloud using wordcloud2 website
(<https://timdream.org/wordcloud2.js/#les-miz>)

Framework

Degree of completion:
100%

步驟說明

1. 先將所有的論文，以“系所名稱”欄位做區分，個別分成小群集。目的是為個別系所建立文字雲；原因是以系所建立，將比較有參考價值，同類所研究較為相近。

*使用方法或套件:
using index (“系所名稱”)

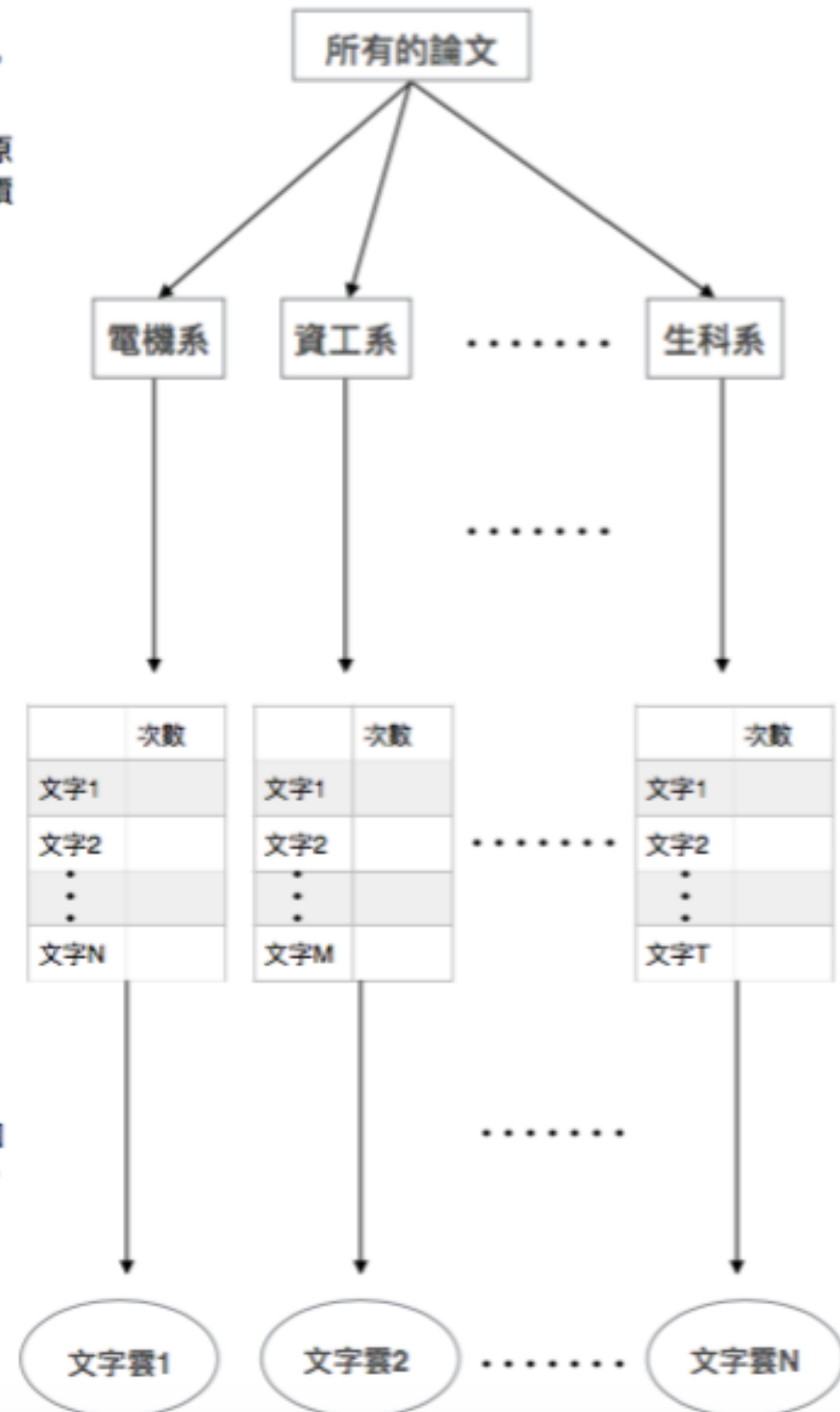
2. 因為論文多為中文撰寫，我們將以“中文關鍵詞”為目標，並計算該關鍵字出現在多少文件的次數(也就是Document Frequency)，來當作權重。

*使用方法或套件:
using nltk count the DF

3. 最後每個系所都會個別統計出該系所論文中，有哪些關鍵字，和各個關鍵字的次數統計，用這些次數統計，每個系所都將產生一個文字雲。

*方法或套件:
using wordcloud or
tool Timdream's wordcloud2

步驟示意



Code

```
import json
theses = json.load(open('nthu_thesis20170330.json'))
theses[0].index('系所名稱')
theses[0].index('中文關鍵詞')

#extract department and Chinese keyword
data_list = [(thesis[2], thesis[17].split()) for thesis in theses[1:]]

#store all the department into a list
index = []
i=0
while i < 5444:
    if data_list[i][0] not in index:
        index.append(data_list[i][0])
    i+=1
```

#Take EE as an example

```
ee = []
```

#store all the EE's keywords into a list

```
j=0
```

```
while j < len(data_list):
```

```
    if data_list[j][0] == index[0]:
```

```
        ee.append(data_list[j][1])
```

```
    j+=1
```

#remove list of list into one list

```
flat_ee = []
```

```
for sublist in ee:
```

```
    for val in sublist:
```

```
        flat_ee.append(val.lower())
```

#count the word and rank it

```
from collections import Counter
```

```
counter_ee = Counter(flat_ee)
```

```
sort_ee = sorted(counter_ee.items(), key=lambda x: -x[1])
```

#reverse the word's and number's position

```
m=0
```

```
while m < len(sort_ee):
```

```
    print (sort_ee[m][1], sort_ee[m][0])
```

```
    m=m+1
```

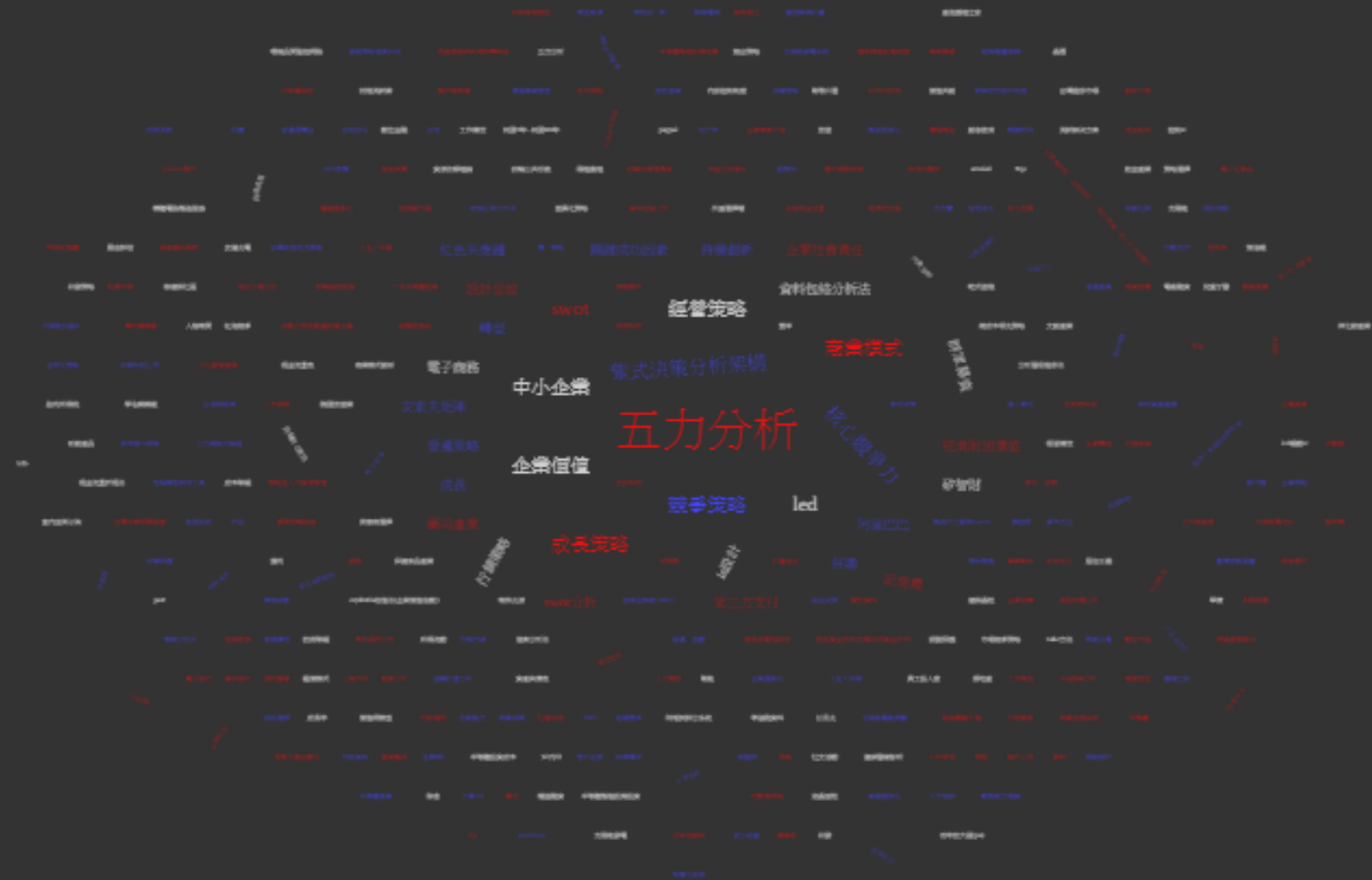
#copy the result and paste to wordcloud2 website

Index of word cloud

0. 電機工程學系
1. 高階經營管理碩士在職專班
2. 生醫工程與環境科學系
3. 材料科學工程學系
4. 分子與細胞生物研究所
5. 通訊工程研究所
6. 資訊工程學系
7. 經濟學系
8. 生物資訊與結構生物研究所
9. 工業工程與工程管理學系碩士在職專班
10. 歷史研究所
11. 台灣文學研究所
12. 電子工程研究所
13. 資訊系統與應用研究所
14. 奈米工程與微系統研究所
15. 物理系
16. 計量財務金融學系
17. 光電工程研究所
18. 數學系
19. 天文研究所
20. 動力機械工程學系
21. 外國語文學系
22. 工程與系統科學系
23. 化學系
24. 社會學研究所
25. 人類學研究所
26. 學習科學研究所
27. 服務科學研究所
28. 核子工程與科學研究所
29. 化學工程學系
30. 統計學研究所
31. 工業工程與工程管理學系
32. 分子醫學研究所
33. 中國文學系
34. 哲學研究所
35. 科技法律研究所
36. 生物科技研究所
37. 經營管理碩士在職專班
38. 科技管理研究所
39. 語言學研究所
40. 先進光源科技學位學程
41. 國際專業管理碩士班
42. 系統神經科學研究所
43. 臺灣研究教師在職進修碩士學位班
44. 生物醫學工程研究所
45. 半導體元件及製程產業研發碩士專班
46. 亞際文化研究國際碩士學位學程

0.電機工程學系

1.高階經營管理碩士在職專班



2.生醫工程與環境科學系



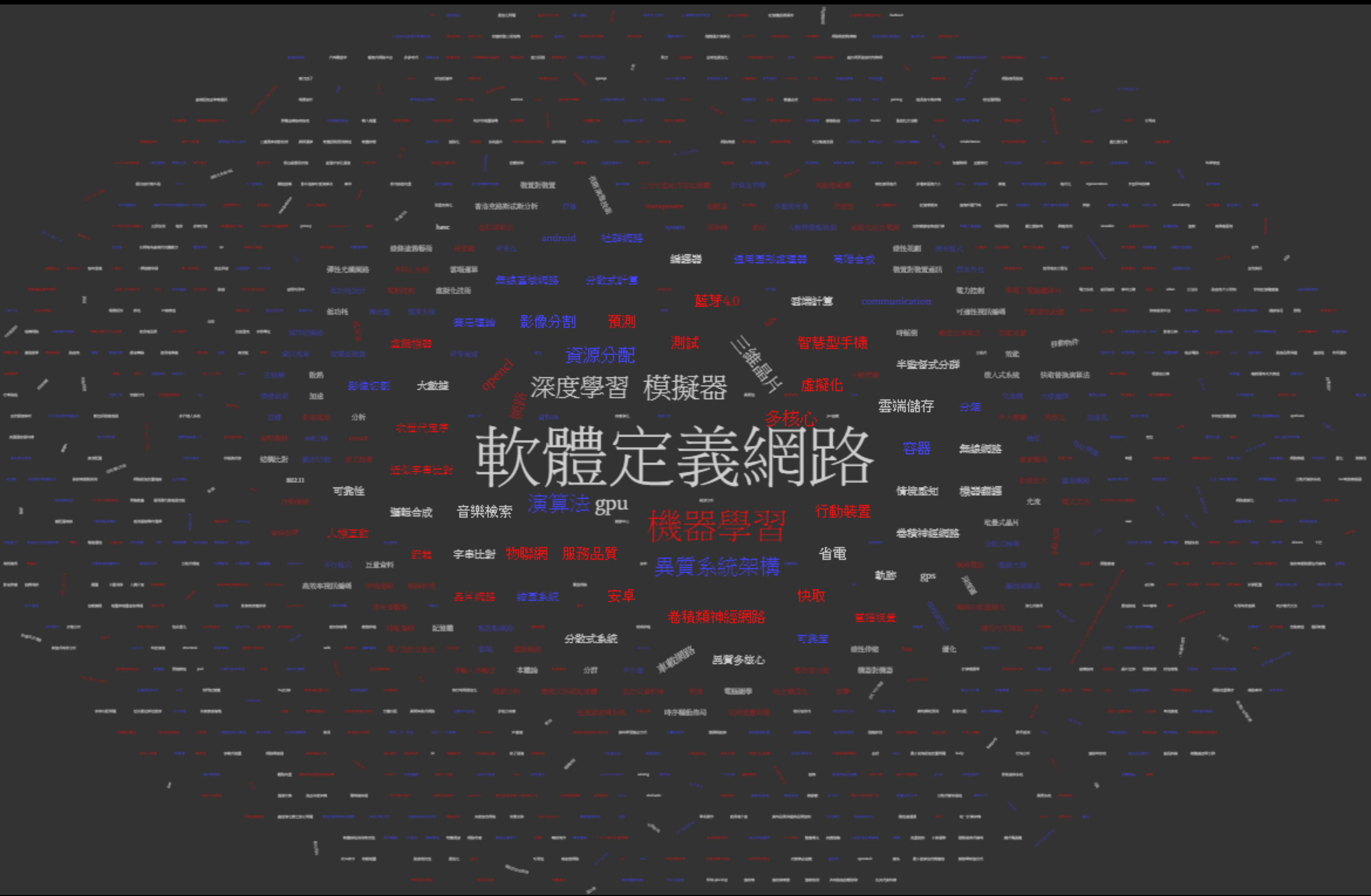
3.材料科學工程學系



4. 分子與細胞生物研究所

5. 通訊工程研究所

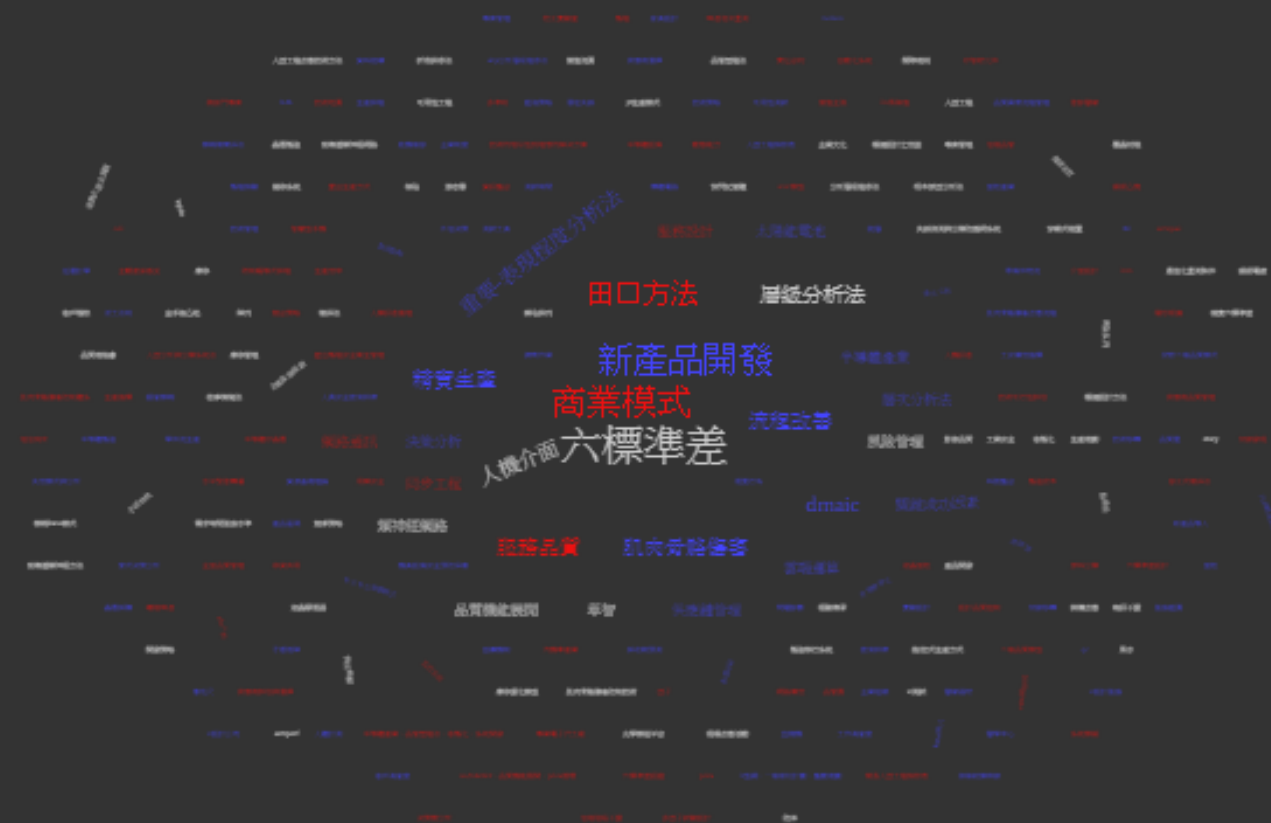
6.資訊工程學系



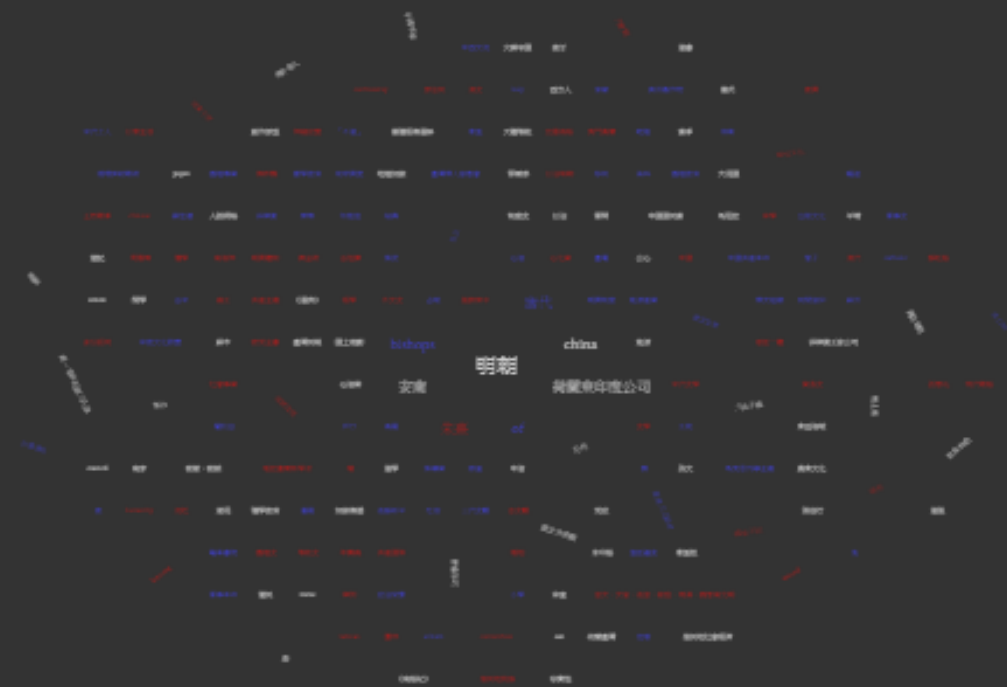
7.經濟學系

8. 生物資訊與結構生物研究所

9.工業工程與工程管理學系碩士在職專班



10. 歷史研究所



11.台灣文學研究所

12. 電子工程研究所

13.資訊系統與應用研究所

14. 奈米工程與微系統研究所

15.物理系

16.計量財務金融學系

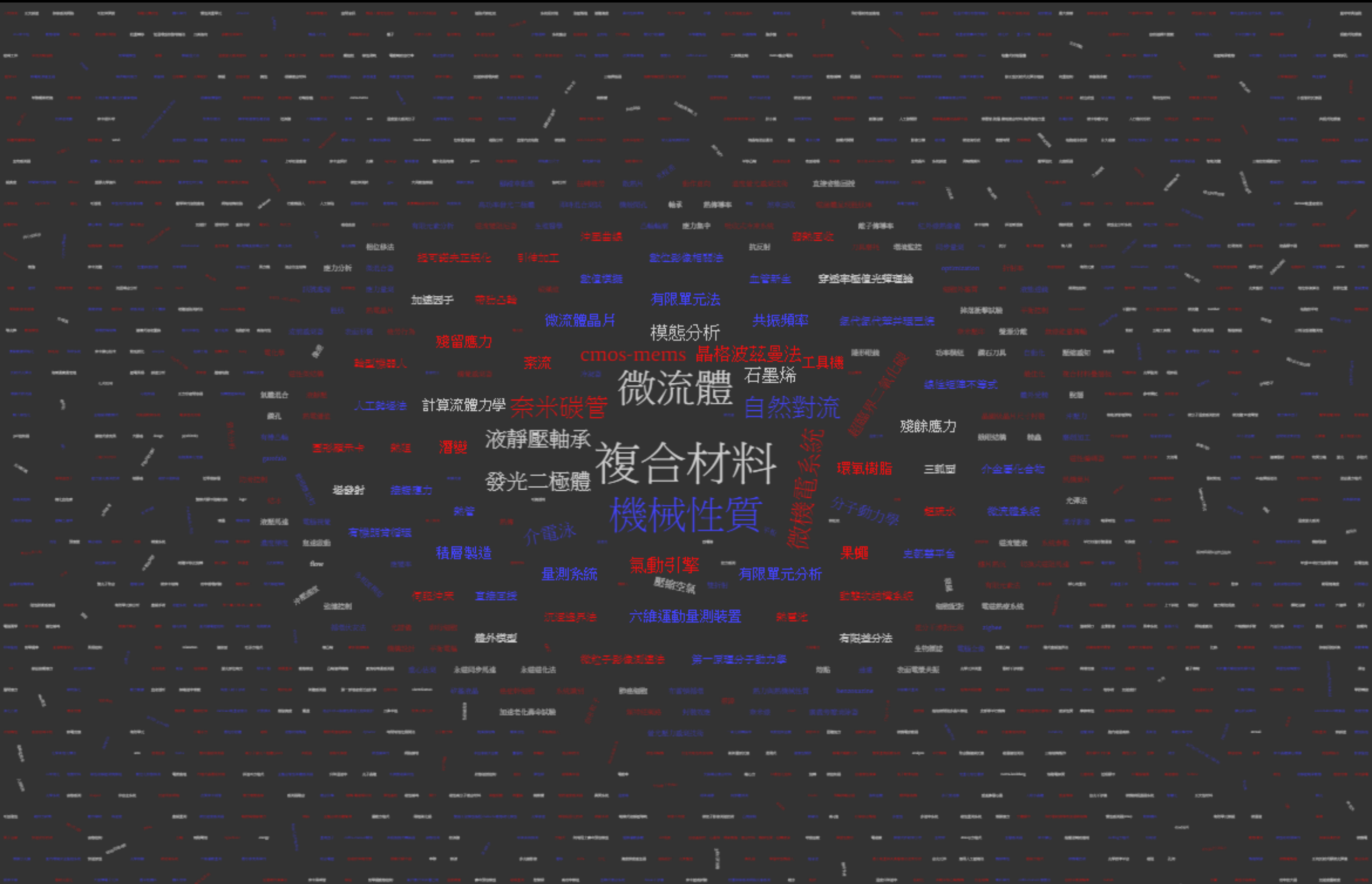
17.光電工程研究所

18.數學系

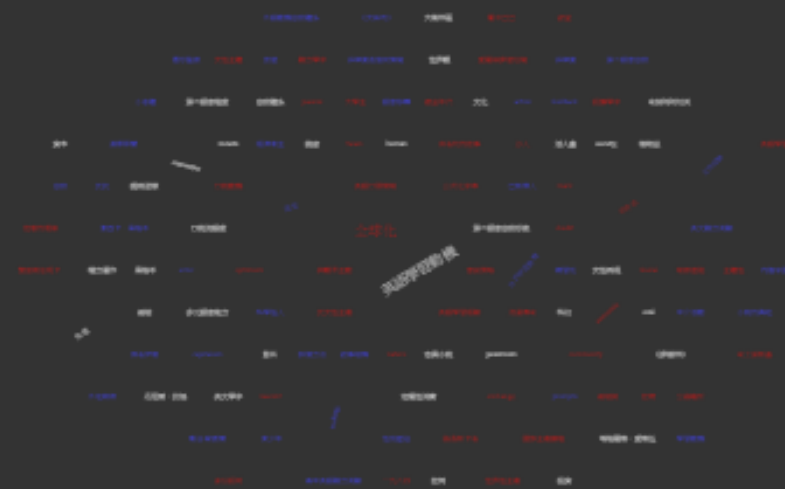


19.天文研究所

20.動力機械工程學系



21.外國語文學系



22. 工程與系統科學系

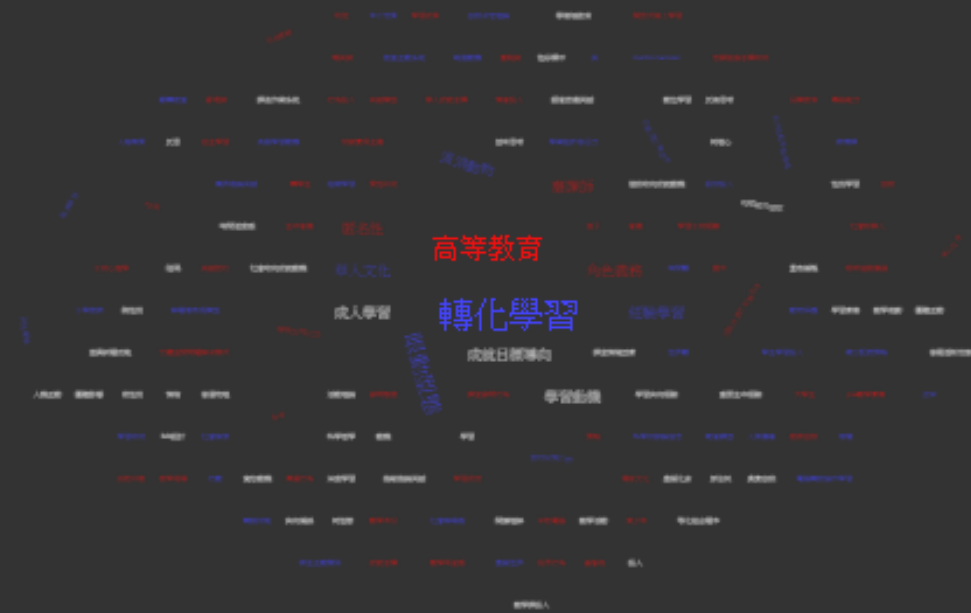
23.化學系

24.社會學研究所

25.人類學研究所

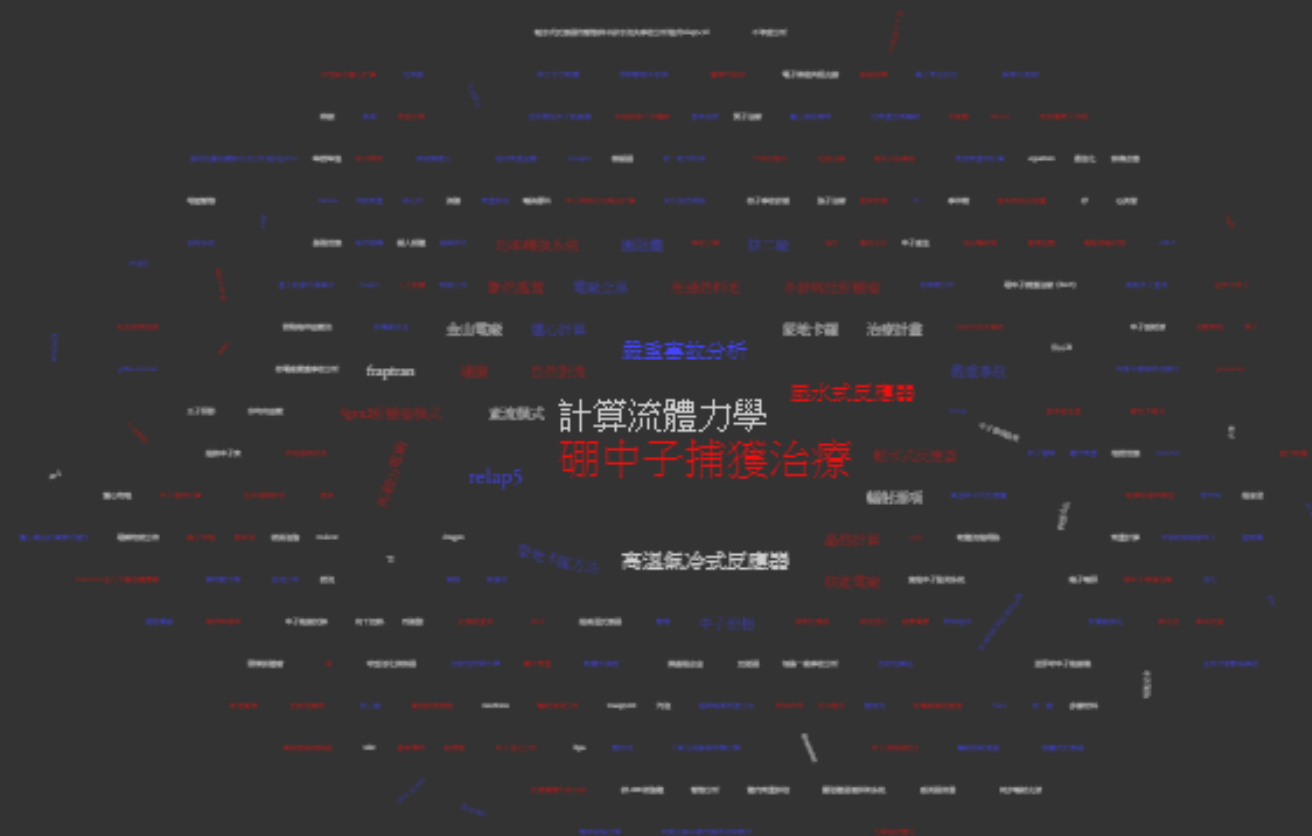


26. 學習科學研究所



27.服務科學研究所

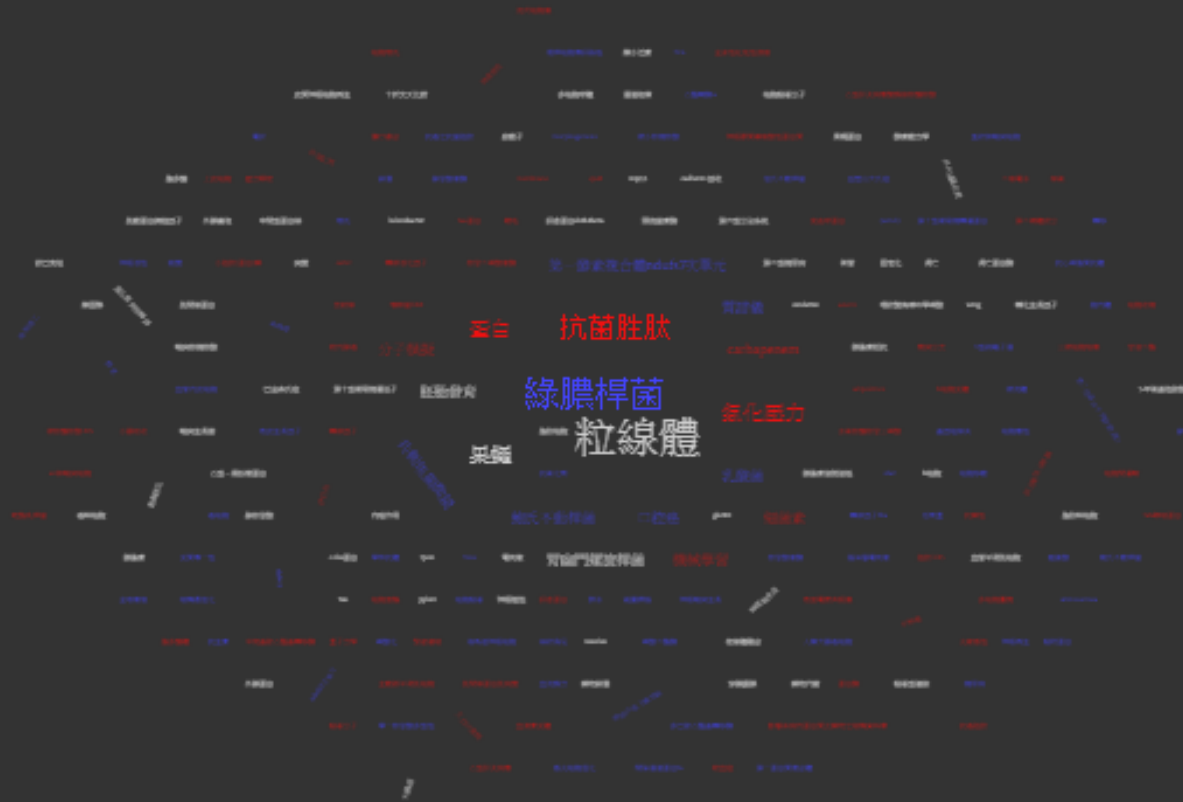
28.核子工程與科學研究所



29.化學工程學系

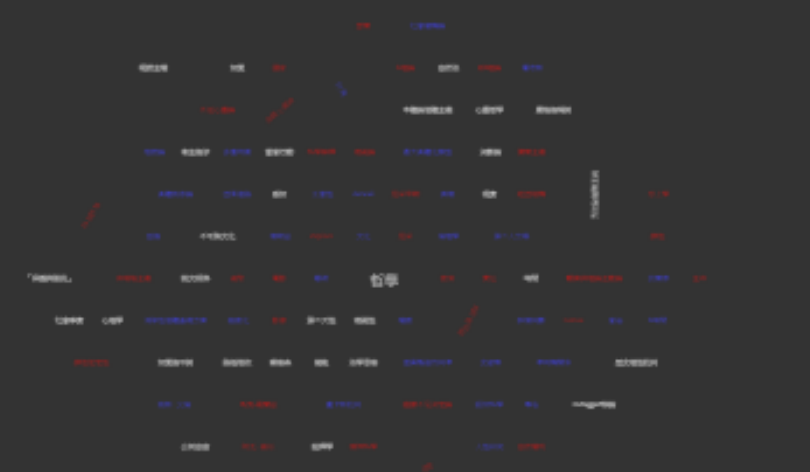
30.統計學研究所

32.分子醫學研究所



33.中國文學系

34. 哲學研究所

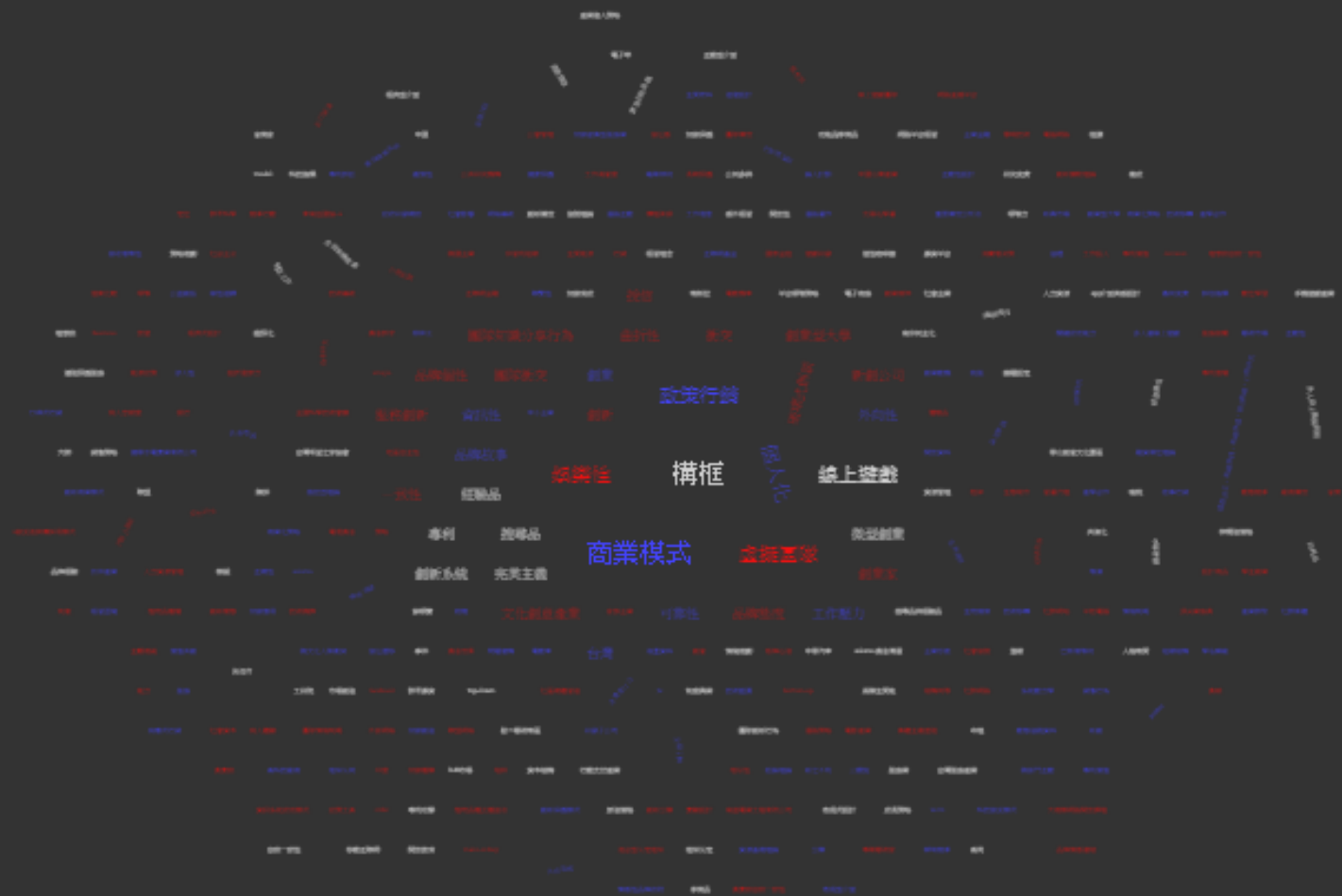


35.科技法律研究所

36. 生物科技研究所

37. 經營管理碩士在職專班

38.科技管理研究所



39. 語言學研究所

40.先進光源科技學位學程

41.國際專業管理碩士班



42.系統神經科學研究所

43.臺灣研究教師在職進修碩士學位班

44.生物醫學工程研究所



45.半導體元件及製程產業研發碩士專班



46.亞際文化研究國際碩士學位學程

