

Data Mining

Indicators for advances in knowledge building – Application of content analysis tools to two sets of CACL discourse data from two comparable classes

Professor Nancy Law
& the Learning Community Project Team

Workshop sponsored by the
Strategic Research Theme on Information Technology
University of Hong Kong
27.10.2006

CSCCL lead to effective learning?

We (teachers, researchers & other interested parties need evidence on students' learning outcomes:

- Cognitive/conceptual developments – what have they learnt? Have students overcome common misconceptions?
- Metacognitive developments – are they better learners? Can they reflect on and monitor their own learning?
- Socio-metacognitive developments – do they know how to work productively with others?
- Can we identify learning progress at individual, group and community levels?

Existing methods for analyzing CACL discourse

Common machine supported methods include

- social network analysis to look at students' participation structure/pattern in the discourse
- determining nature of the discourse transactions (e.g. the level of argumentation/critical thinking exhibited , etc.) using syntactic analysis/ build-in scaffolds such as sentence openers.

But, these methods per se

- **Cannot reveal changes at the cognitive level without performing analysis at the semantic level**
- **Vocabulary growth** has been used as one form of semantic analysis, but \neq **“growth of knowledge”** (a lot of information can be posted without thoughtful consideration or understanding)

Our research to-date

TOOL development:

VINCA - Visual INtelligent Content Analyzer - content analysis tool jointly developed by CITE, HKU and CKSER, BNU

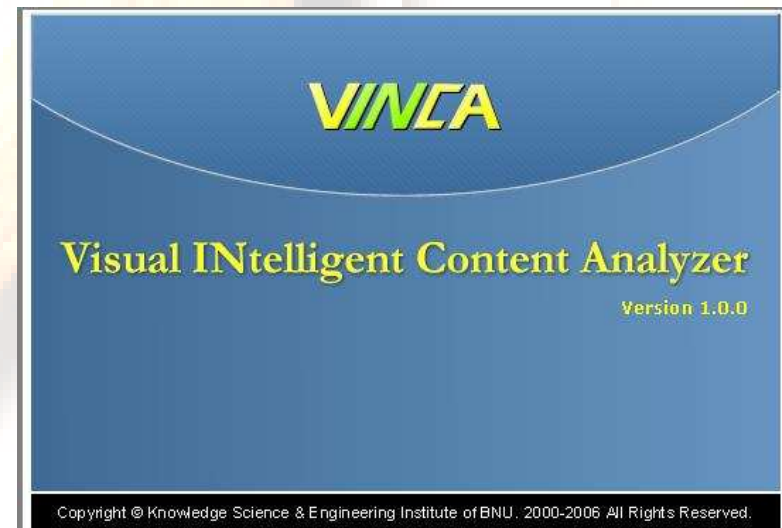
Goals:

- To develop a tool that can support semantic analysis, interaction analysis, social network analysis and a combination of the above to assess knowledge building outcomes at individual and group levels
- To conduct further mining of the multidimensional coding to develop models of learning in CACL contexts
- To develop online tools (learning facilitation agents) to support teachers and learners in CACL learning situations

VINCA - Visual INtelligent Content Analyzer - content analysis tool jointly developed by CITE, HKU and CKSER, BNU

Currently, it includes the following functions:

- Data preparation to convert Knowledge Forum® discourse in html to database format
- Keywords retrieval
- Manual coding support
- User-improvable semi-automatic semantic coding
- Social network analysis
- Novelty and similarity analysis



Examining knowledge building outcomes using conventional & data mining methods

Background

- Ho Lap College, Form 3 Design & Technology Curriculum
- Teacher wanted to develop students' critical thinking through discussion slimming
- Total 5 classes. Each class was split into two groups that took turn to study this subject in 2 different school terms (Oct – Dec, 04 ; Jan – May, 05)
- The classes met roughly twice a month



Research Tools Developed

- **Weight-loss & nutrition concept test** (aimed at assessing students' relevant (mis)conceptions & understanding)
- **Daily food intake assessment sheet** (to understand students' dietary habits)
- **Weight-loss, exercise & body image survey** (to understand students' perceptions and beliefs in such issues)

Data Collected

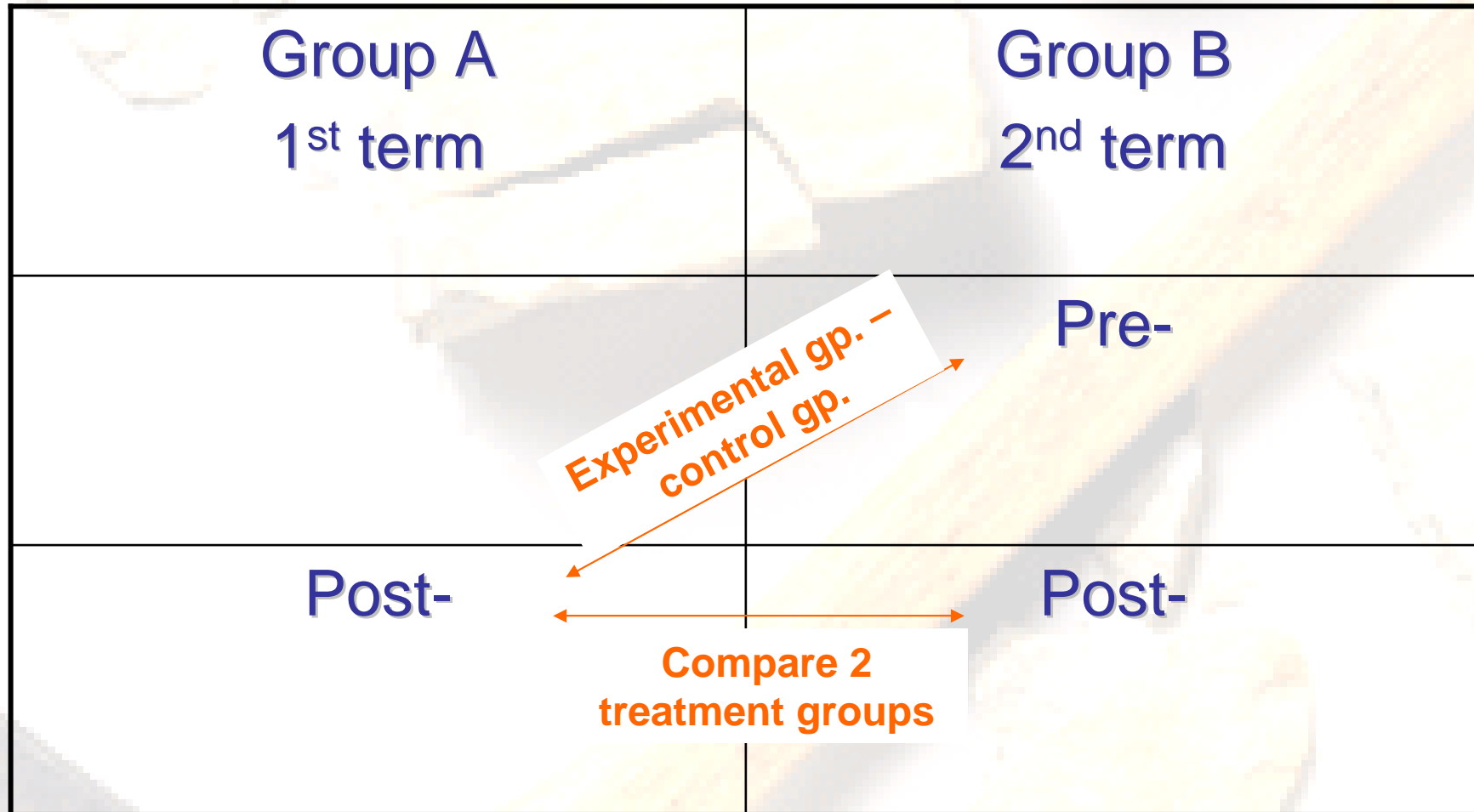
Via conventional instruments

1. Misconception test
2. Food intake assessment
3. Slim-up survey

Qualitative data from discussion process

1. Knowledge Forum® discussion contents
2. Class observation field notes
3. Student focus group interviews
4. Teacher reflections
5. Video recordings of selected classes

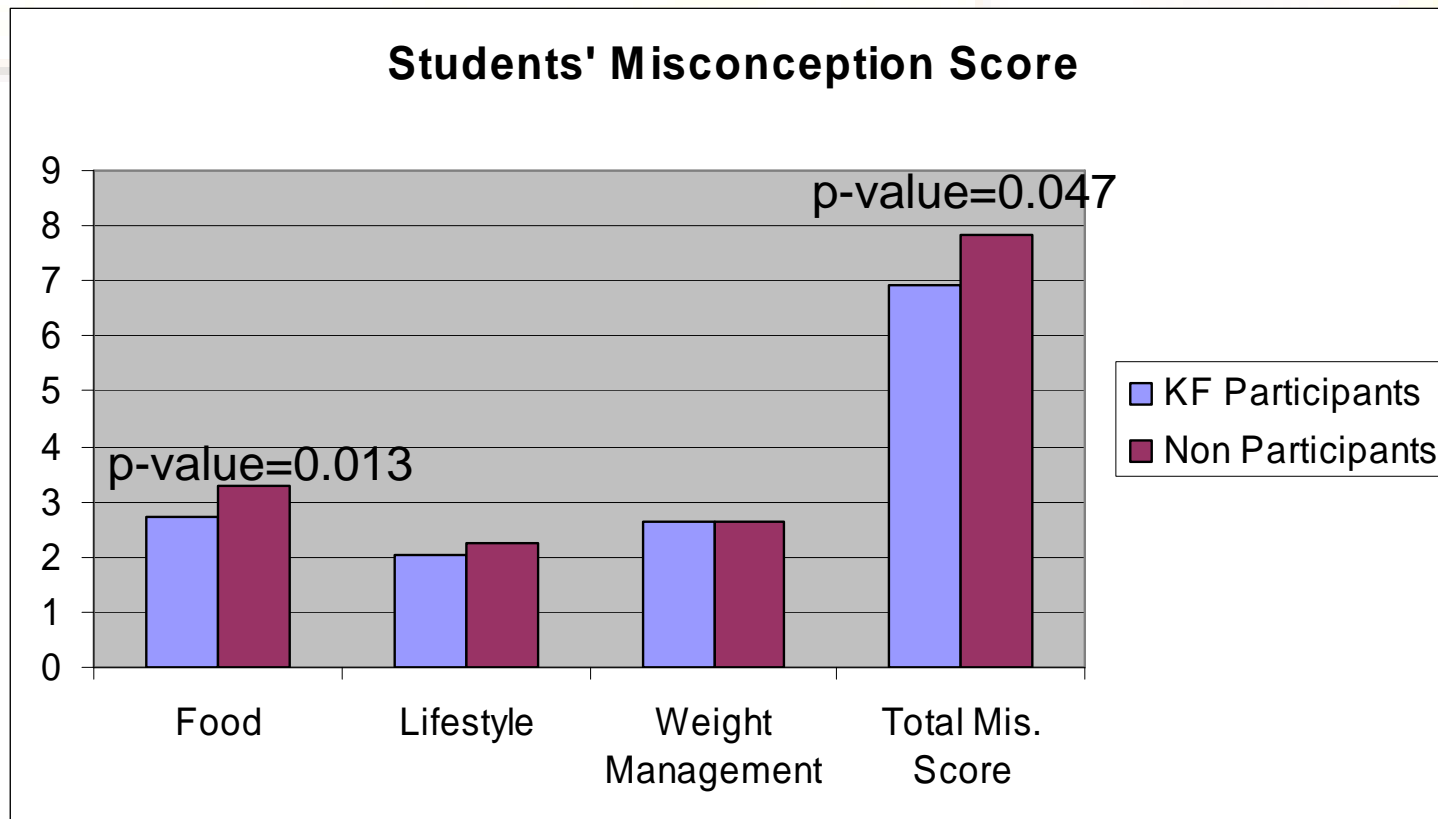
Data collected



Quantitative data findings

Misconceptions Test

1st Round Study (Control gp - expt gp)

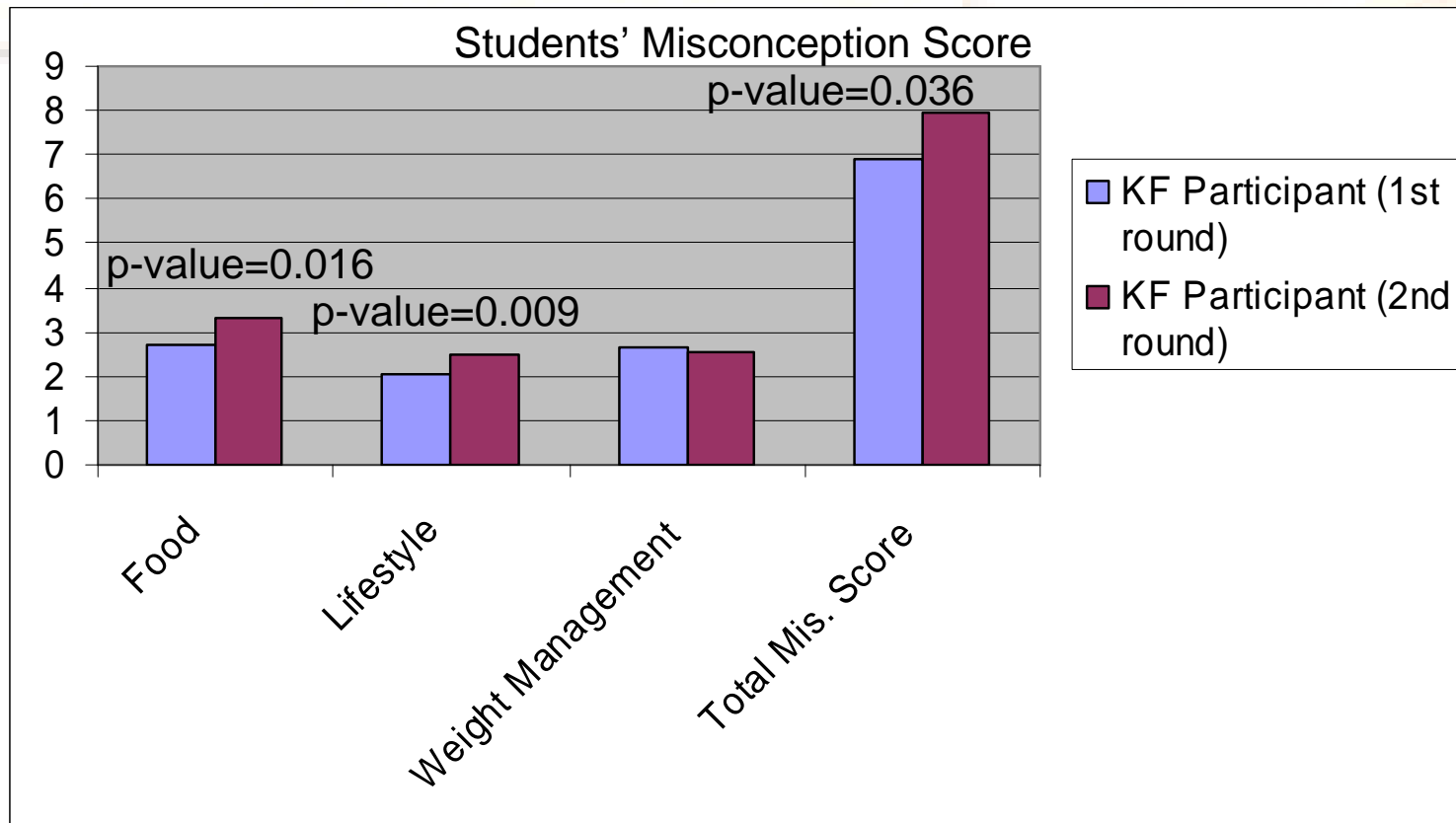


- **Term 1 treatment group has fewer misconceptions than control group**

Quantitative data findings

Misconception Test

1st Round & 2nd Round (post-) Comparison

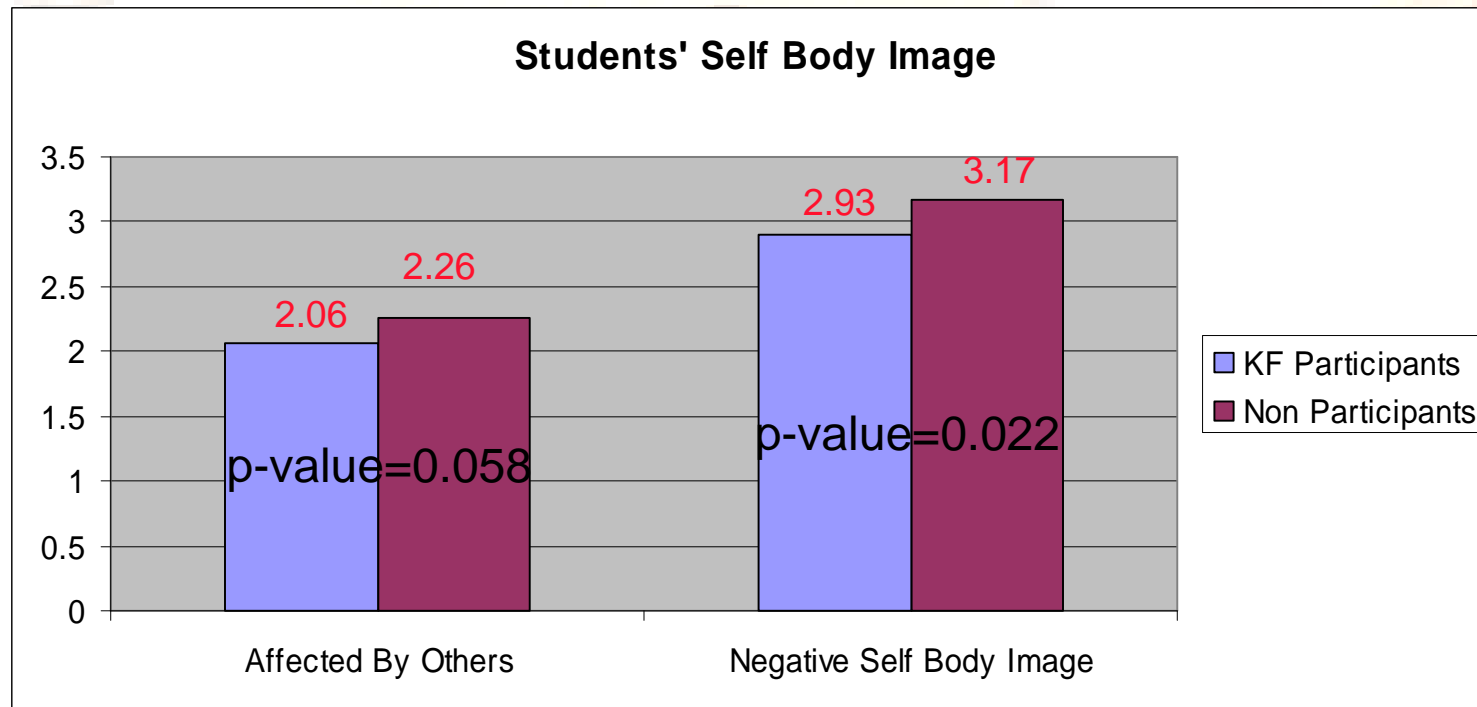


- **Term 1 treatment group has fewer misconceptions than term 2 treatment group**

Quantitative data findings

Slim-Up Survey

1st Round Study (Control gp - expt gp)



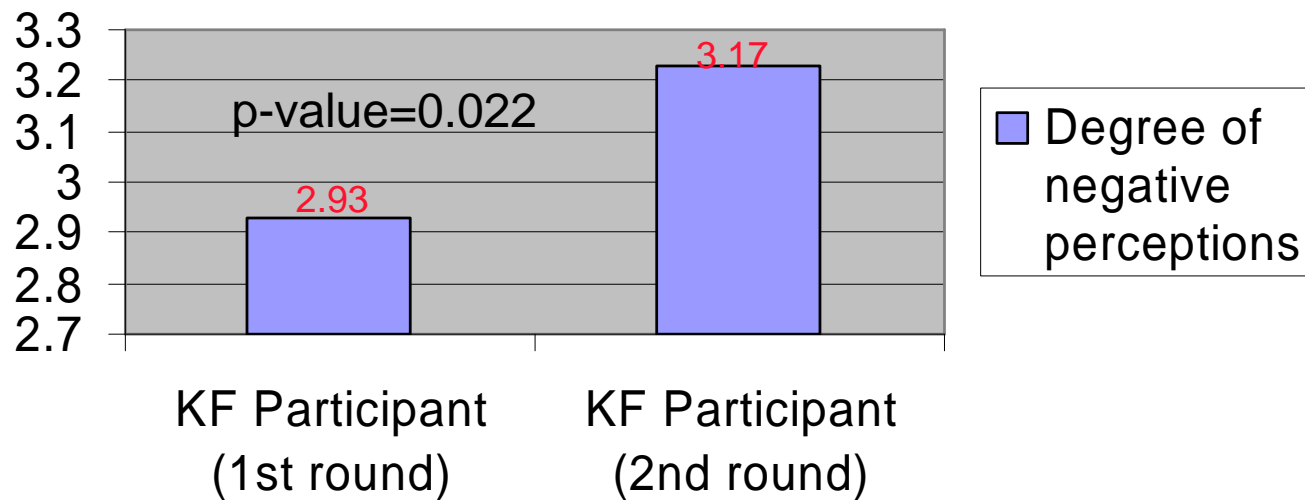
- **Term 1 treatment group has higher self-image than control group**

Quantitative data findings

Slim-Up Survey

1st Round & 2nd Round (post-) Comparison

Students' Self Body Image Perception

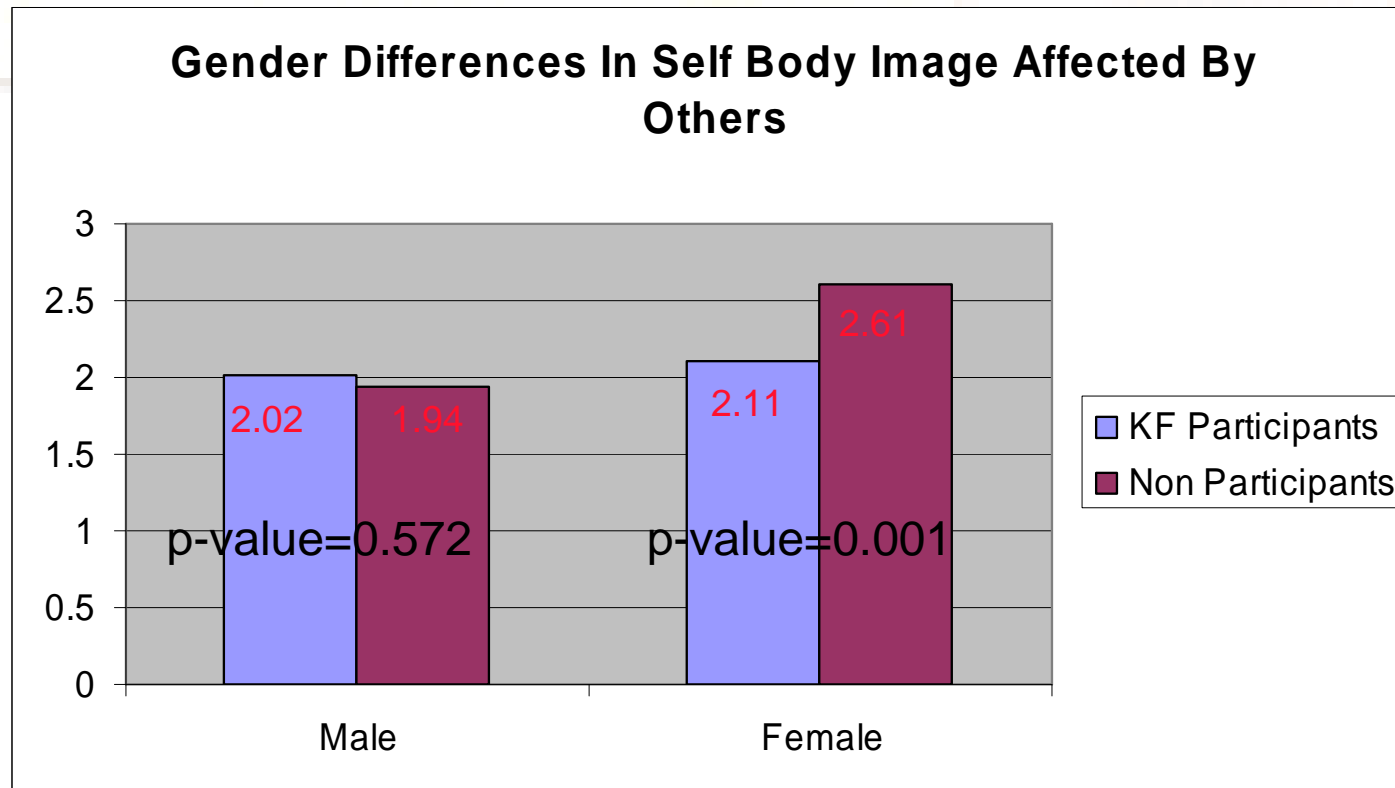


- **Term 1 treatment group has higher self-image than term 2 treatment group**

Quantitative data findings

Slim-Up Survey

1st Round Study (Control gp - expt gp)



- Improved self-image of 1st term treatment group only found in girls
- In control group, self-image of girls sign. Lower than boys
- In 1st term treatment group, no statistical gender difference in self-image

Learning outcomes are very different though both involve same kind of discussion task

- Why are there such big differences between the two treatment groups?
- What contributes to better learning through collaborative learning discussions?
- Can we identify features of more productive discussions?

A case study of discourse analysis: Slim up discussions on Knowledge Forum®

- Duration span: 1 term
- Number of students: 2 groups of Grade 9 students, ~ 20 for each group, randomly assigned
- Which group is better at knowledge building?

	Total no. of notes	No. of threads	Notes/thread	Threads with > 6 notes	No. of keywords
Group A	123	25	4.92	5	1552
Group B	298	86	3.46	2	5396



A 3-step semantic analysis

Step 1: **Keyword extraction** to identify focal ideas

- VINCA was used to generate the frequencies of all keywords found in the KF discussion.
- From the output, researchers were able to identify a number of key terms with high frequencies from the slim up discussion, such as “*lose weight*”, “*slimming*”, “*beauty*”, “*thin*” and “*I*”

Step 2: Extraction of discourse text around selected keywords using **concordance** technique

Step 3: **Further keyword analysis**

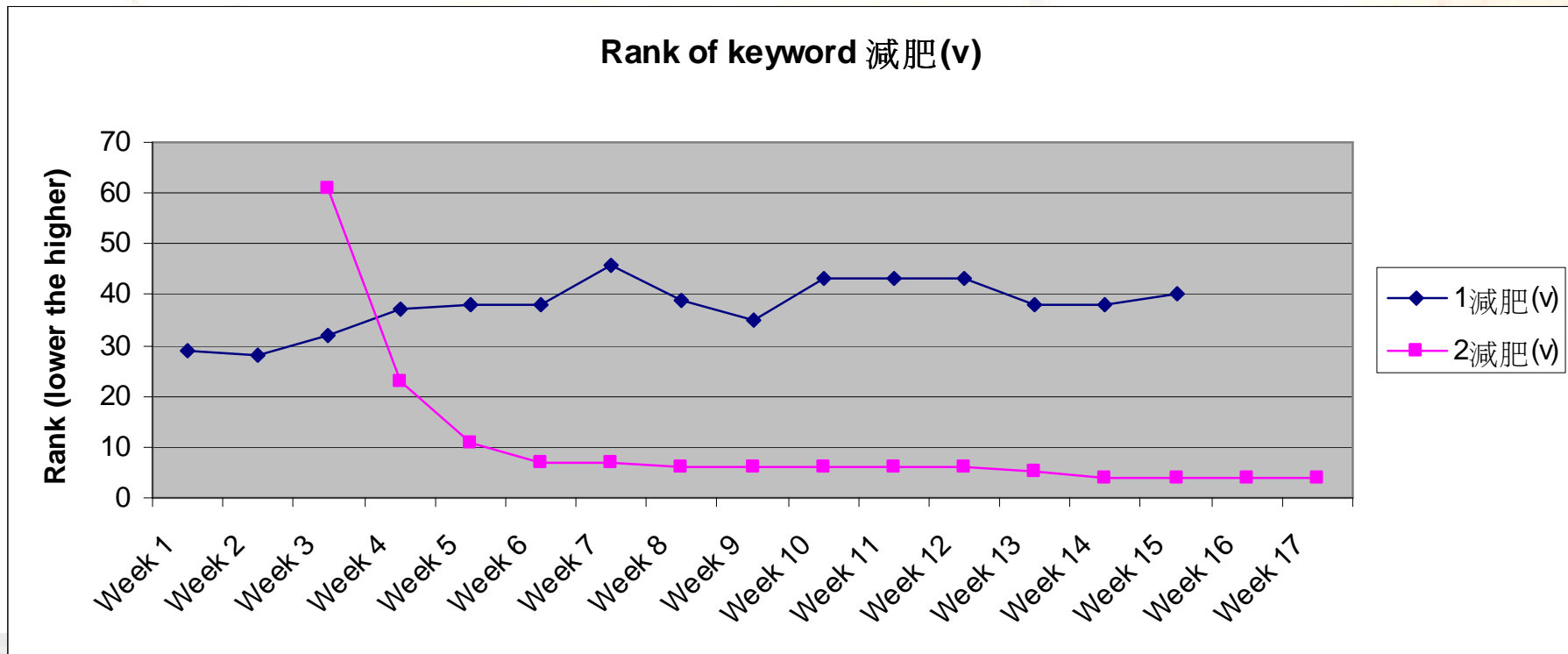
- The text extracted by VINCA from stage 2 was analyzed using VINCA again to generate a list of frequencies of keywords in close proximity to selected key terms.

Nouns

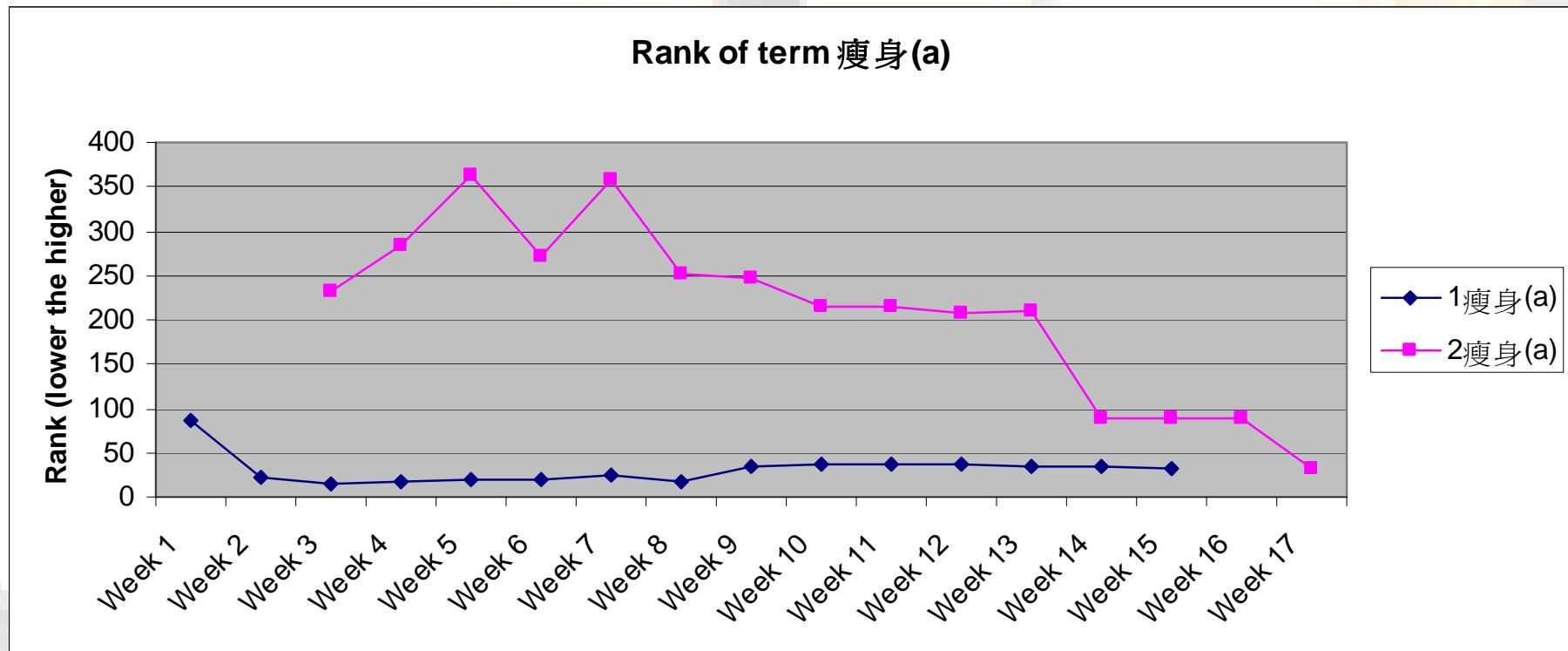
Mining

- On the other hand, nouns are more frequently used in Group B's Discourse
 - Group A: 49 diff. nouns, total freq. 98
 - Group B: 1948 diff. nouns, total freq. 6717!!!

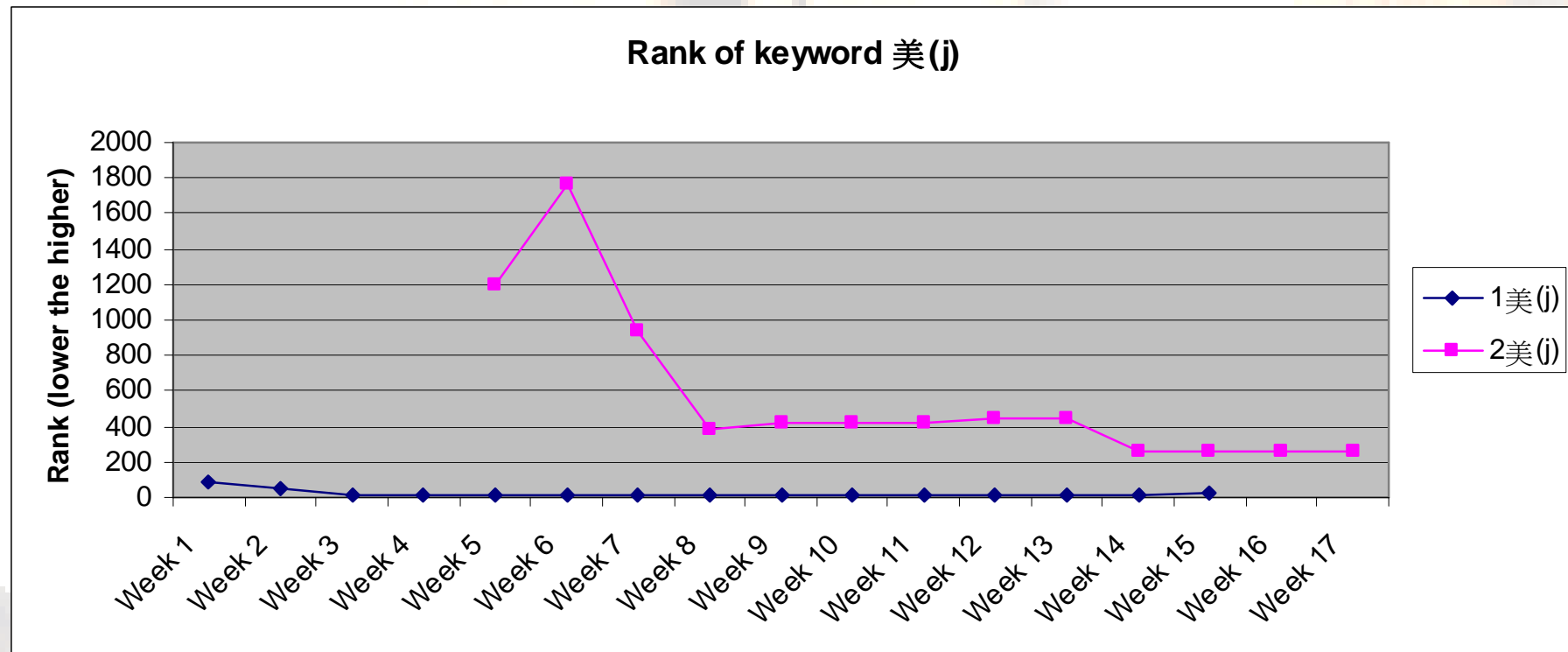
Keyword ranking



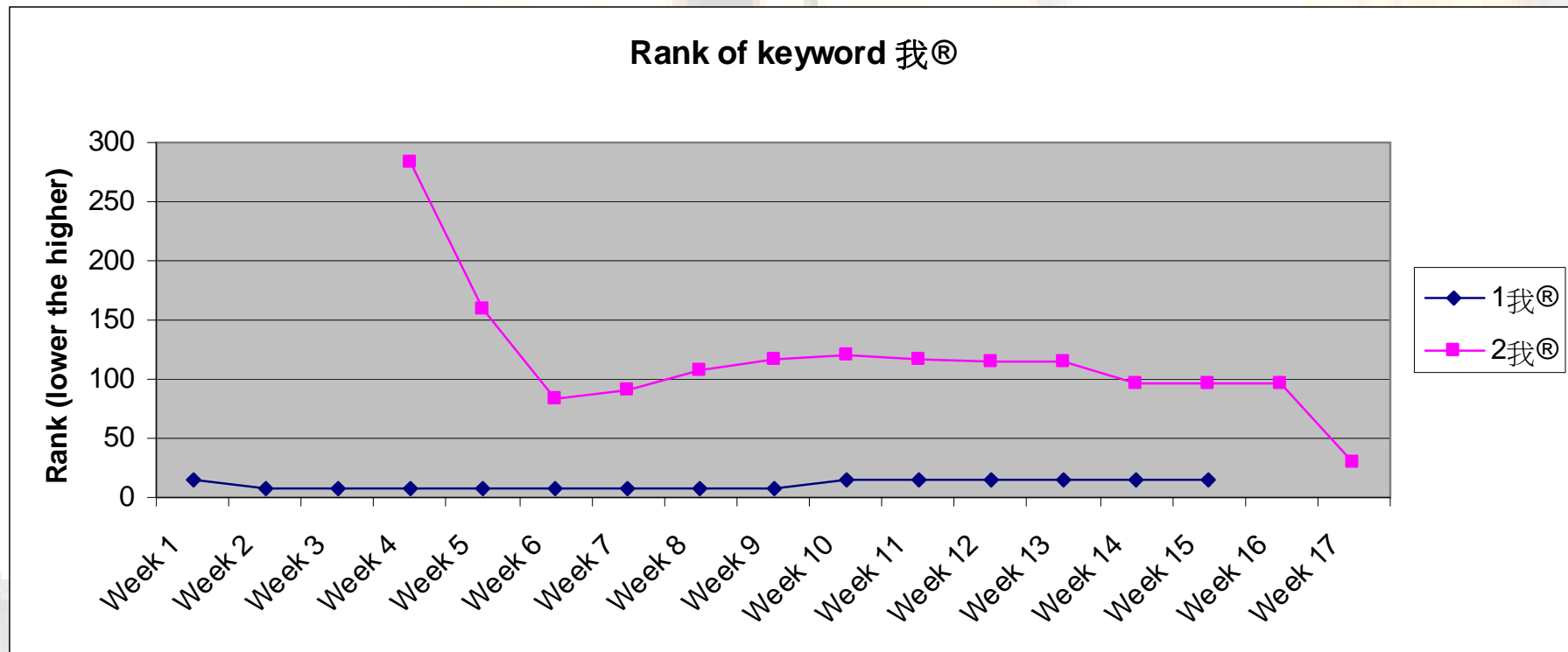
Keyword ranking



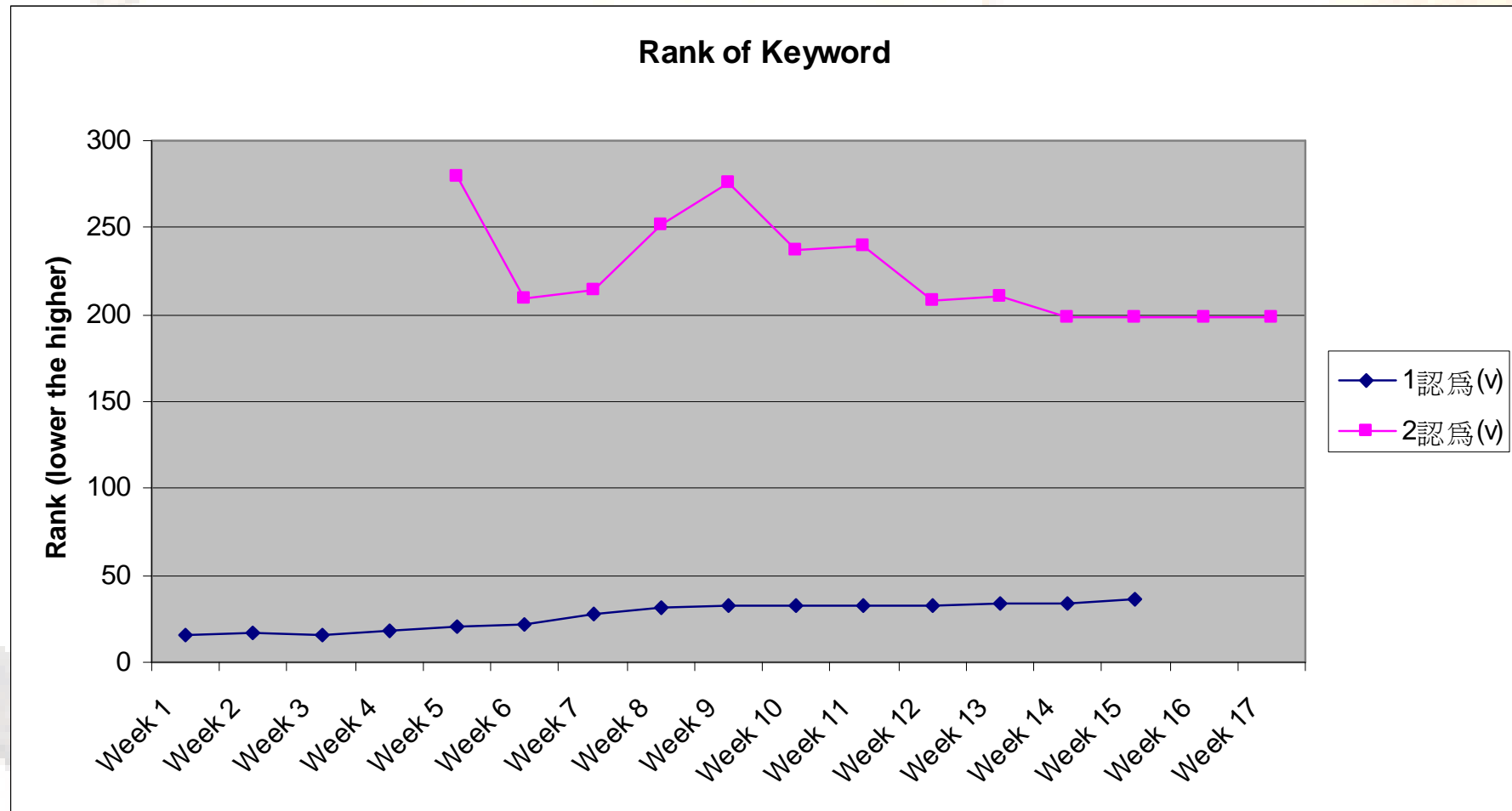
Keyword ranking



Keyword ranking



Keyword ranking



VINCA Text Analysis screen

Winca

File (F) View (V) AnnotationAid (A) TextAnalysis (C) DataExport (D) Window (W) Help (H)

Step1 Select Data... Step2 Get KeyWords... Step3 Word Category... Step4 Export...

Word List & Distribution of Frequency

ID	Word	Tag	Count	firstUser	hlc 3a cheng yui lung	hlc 3a cheng yui lung	hlc 3a cheng yui lung
1013	beauty	n	45	hlc 3a fok suk fun	0	2	2
69	我	r	45	hlc 3a lau lai pui	0	6	0
18	i	n	41	hlc 3a cheng yui lung	1	1	1
1009	keep	n	37	hlc 3a fok suk fun	0	2	1
1017	forever	n	36	hlc 3a fok suk fun	0	2	1
61	了	y	35	hlc 3a lau lai pui	0	1	0
15	my	n	34	hlc 3a cheng yui lung	0	1	1
20	a	n	33	hlc 3a cheng yui lung	0	5	3
175	是	v	31	hlc 3a fok suk fun	0	2	0
43	...	n	20	hlc 3a cheng yui lung	0	2	0

List of keywords & frequency counts

Concordance context of keywords

Concordance Result

User	Time	Context
A slim up	hlc 3a lau lai pui	oct 10 2004 (09: 4:52)
A slim up	hlc 3a lau lai pui	oct 10 2004 (09: 2:12)
A slim up	hlc 3a chan chi shing	oct 29 2004 (22: 4:45)
A slim up	hlc 3a fok suk fun	sep 30 2004 (15: 5:02)
A slim up	hlc 3a lau lai pui	oct 10 2004 (19: 5:11)
A slim up	hlc 3a choi yan kit	oct 12 2004 (21: 7:00)
A slim up	hlc 3a choi chi hin	oct 13 2004 (14: 0:58)
A slim up	hlc 3a fok suk fun	oct 28 2004 (11: 6:07)

Outcome of the second stage semantic analysis on slim up

		Group A		Group B	
減肥-Lose weight	Type				
	Personal	我	15	我	12
	Opinion	認為	7	認為	8
	Opinion	覺得	1	覺得	0
	Opinion	意見	0	意見	5
	Ideas	了解	4	了解	0
	Ideas	想法	3	想法	0
	Explain	因為	2	因為	10
	Social	認同	0	認同	1
	Question	想知道	2	想知道	0
	Question	嗎	6	嗎	0
	Question	麼	1	麼	1
	Personal	個人	1	個人	0

		Group A		Group B	
美-Beauty	Type				
	Personal	我	13	我	1
	Opinion	認為	10	認為	3
	Opinion	覺得	7	覺得	0
	Opinion	意見	2	意見	0
	Ideas	了解	0	了解	0
	Ideas	想法	4	想法	0
	Explain	因為	2	因為	10
	Social	認同	1	認同	1
	Question	想知道	0	想知道	0
	Question	嗎	0	嗎	2
	Question	麼	1	麼	0
	Personal	個人	1	個人	0

		Group A		Group B	
瘦身-Body slimming	Type				
	Personal	我	9	我	2
	Opinion	認為	2	認為	1
	Opinion	覺得	4	覺得	0
	Opinion	意見	0	意見	1
	Ideas	了解	1	了解	0
	Ideas	想法	0	想法	0
	Explain	因為	0	因為	0
	Social	認同	0	認同	0
	Question	想知道	0	想知道	0
	Question	嗎	1	嗎	0
	Question	麼	3	麼	0
	Personal	個人	0	個人	0

		Group A		Group B	
我-I	Type				
	Personal	我	51	我	55
	Personal	個人	7	個人	3
	Opinion	認為	16	認為	5
	Opinion	覺得	9	覺得	1
	Opinion	意見	7	意見	1
	Ideas	了解	5	了解	0
	Ideas	想法	2	想法	0
	Explain	因為	5	因為	1
	Social	認同	4	認同	1
	Question	想知道	2	想知道	0
	Question	嗎	1	嗎	1
	Question	麼	3	麼	0

	Group A	per 1000 Kws Count	Group B	per 1000 Kws Count
老師	0	0	1	0.03767
醫生	0	0	31	1.167784
營養師	0	0	22	0.82875
營養	1	0.207297	25	0.941761
護士	0	0	1	0.03767
脂肪	13	2.694859	108	4.06841
Reflective	認為	19	34	1.280796
	想	5	27	1.017102
	覺得	8	11	0.414375
	相信	2	7	0.263693
	知道	2	7	0.263693
	感到	2	7	0.263693
	認同	4	1	0.03767
	想到	0	1	0.03767
Claims	其實	15	17	0.640398
	所以	6	33	1.243125
	而	18	96	3.616364
	而且	5	16	0.602727
	就是	10	22	0.82875
	解釋	1	1	0.03767
	根據	4	10	0.376705
Queries	嗎	9	9	0.339034
	可否	1	1	0.03767
	有沒有	1	1	0.03767
	怎樣	3	2	0.075341
	如何	1	7	0.263693
	甚麼	5	0	0
	呢	10	6	0.226023
為甚麼	2	0	0	

A preliminary Interpretation

- Group A, seems to be more engaged in reflecting, making claims, and putting questions forward.
- While Group B students seems to do less reflections, claims and queries, while having many many nouns.
- Can we seek deeper understanding of the difference between the 2 groups' discourses?

Personal cognitive engagement

Examples from group A

- 我認同的的說法，的確肥胖的人進行纖體的確是健康，但纖瘦的人也依樣葫蘆，照著幹那就有問題了。
- 我應[認]為減肥是指把原先肥胖的身軀減至正常體重而瘦身則是把一個正常體重的身軀減至更瘦
- 根據我的理解，減肥就是透過一些方法來減輕體重從而做到控制體重。
- All these contents contain the word “我” to indicate some forms of cognitive engagement

Personal cognitive engagement

Non-examples from group B

- 有人說：“**我**吃了減肥藥不是瘦掉了？”快速減肥的確能使人瘦
- 但是，如果光靠睡覺減肥，這一點兒**我**卻不敢認同。從健美的理論來看，運動、休息、營養就拿**我**個人來說，**我**每天都會保持很大的運動量：早晨5點多，**我**就起來練功，每天練一個小時左右。睡眠時間也一定要保證，... **我**晚上睡眠時間雖然不多，但是**我**白天往往會補上一覺的，這樣就能保證充足的睡眠了。
- ...中風和心臟衰竭等，「已是一個病，纖體中心無法處理，**我**只有一個建議，就是請他們去睇醫生」。他又特別指出 ...
- All these contents which contain the word “**我**” are actually **quoted speech**.

Concordancing Mining

- Examining words in close proximity of selected keywords will reveal the semantic context when those keywords are used, thus revealing whether there is deep cognitive engagement or only casual sharing of information.
- Concordancing of “我” in the two discourse thus reveal the depth of engagement of the students when they discussed slimming in Knowledge Forum®.
- This indicates that some text mining of selected keywords in close proximity would be better at identifying significant features of CACL discourse.

我

你

	Group A Counts	per 1000 Kws count	Group B Counts	per 1000 Kws count	Group A Counts	per 1000 Kws count	Group B Counts	per 1000 Kws count
老師	0	0.0000	0	0.0000	0	0.0000	0	0.0000
醫生	0	0.0000	11	1.1015	0	0.0000	8	1.1532
營養師	0	0.0000	4	0.4006	0	0.0000	3	0.4325
營養	0	0.0000	14	1.4020	0	0.0000	13	1.8740
護士	0	0.0000	0	0.0000	0	0.0000	0	0.0000
脂肪	4	2.1810	37	3.7052	4	5.2219	57	8.2168
認為	15	8.1788	10	1.0014	3	3.9164	8	1.1532
想	4	2.1810	19	1.9027	1	1.3055	6	0.8649
覺得	6	3.2715	4	0.4006	1	1.3055	3	0.4325
相信	1	0.5453	3	0.3004	1	1.3055	2	0.2883
知道	2	1.0905	2	0.2003	0	0.0000	3	0.4325
感到	2	1.0905	2	0.2003	1	1.3055	2	0.2883
認同	3	1.6358	1	0.1001	2	2.6110	0	0.0000
想到	0	0.0000	0	0.0000	0	0.0000	0	0.0000
其實	9	4.9073	9	0.9013	3	3.9164	4	0.5766
所以	3	1.6358	17	1.7024	0	0.0000	8	1.1532
而	13	7.0883	43	4.3060	5	6.5274	33	4.7571
而且	3	1.6358	10	1.0014	3	3.9164	5	0.7208
就是	6	3.2715	12	1.2017	1	1.3055	13	1.8740
解釋	0	0.0000	1	0.1001	0	0.0000	1	0.1442
根據	1	0.5453	2	0.2003	1	1.3055	3	0.4325
嗎	3	1.6358	3	0.3004	4	5.2219	4	0.5766
可否	0	0.0000	0	0.0000	0	0.0000	0	0.0000
有沒有	0	0.0000	0	0.0000	0	0.0000	0	0.0000
怎樣	1	0.5453	0	0.0000	0	0.0000	0	0.0000
如何	0	0.0000	4	0.4006	0	0.0000	2	0.2883
甚麼	3	1.6358	0	0.0000	0	0.0000	0	0.0000
呢	7	3.8168	5	0.5007	2	2.6110	4	0.5766
為甚麼	2	1.0905	0	0.0000	0	0.0000	0	0.0000

Reflective

Claims

Queries

In the 我 concordance

- Comparing the “我” concordance between Group A and Group B, reflections, claims, and queries are still more frequently used in Group A’s discourse
- Data supports Group A do more reflections, claims and queries than Group B in the “我” concordance

Examples of cognitive engagement that fail to conform to the same pattern

- 怎樣才是肥,怎樣才是纖瘦?如何得到1個標準?BMI已不能如現在人們所想的標準了!
- 其實瘦就是美是一個非常錯誤的想法。因為美，不只是注重外表，有內在美都是美，無論是男性或是女性，兩者都是一樣。
- 而且吃藥減肥的話,又可能會引起副作用,就會帶來本來不必要的麻煩和煩惱.
- In the above quotes, part of speech information are hidden, but they still reflect active cognitive agency to push ideas to evolve within the text

Knowledge augmented text mining - Cognitive Linguistic Markers (1)

Part of speech indicators

1st person	我	本人			
1st person plural	我們	大家			
2nd person	你	您			
2nd person plural	你們	您們			
3rd person	他	她	它	別人	
3rd person plural	他們	她們	它們	有些人	

Time indicators

Past	古時				
Current	現時	現在	現代	現今	
Future	將來	未來			

Knowledge augmented text mining - Cognitive Linguistic Markers (2)

Claims

Explanations	因為 是指 就是 由於 其實 解釋 所謂
Contrasting	但 卻 而 則是 但是 雖然 另外 另一方面 可是 然而 雖 一方面 反而
Affirmative	對 的確 讚成 讚同 認同 就是
Negative	不 不能
Quoting	根據 說 稱 表示 好像
Concluding	所以 因此 故 故此
Relating	還 及 和 再者 或 並 而且 既 或者 此外 以及 並且 同時 不僅 況且
Conditioning	如果 若 倘若 只要 即使 若果 無論 否則 不然 除非 假如 雖說 何況
Sequencing	然後 其次
Possessive	的
Quantity	一些
Targeting	為了
Exaggeration	太 非常

Knowledge augmented text mining - Cognitive Linguistic Markers (3)

Queries

General 嗎 呢 怎樣 甚麼 有沒有 是不是 是否 難道

Seeking Information 想知道

Seeking Help 可否

Seeking Instructions 如何

Non-linguistic indicators ?

Knowledge augmented text mining - Meta-Cognitive Linguistic Markers

Affective indicators

Non-linguistic
indicators

=) :) :(^_^ ^_!

Emotion indicators

難過 幸福 自卑 開心 不開心 快樂 高興 討厭

Reflectives

Knowing

知道 記得 想 想到

Personal beliefs

相信 認為 覺得 理解 看法

Building up further text patterns using intelligent text encoding dictionaries

- 根據我的理解，減肥就是透過一些方法來減輕體重從而做到控制體重。
 - 根據 我 的 理解，
 - Quoting 1stperson possessive personal beliefs
 - 減肥就是透過一些方法來減輕體重從而做到控制體重。
 - explanations>quantity>relating

Next Steps

- Develop better indicators of knowledge building outcomes through text pattern identification from the following perspectives: domain ontology, social interaction patterns, discourse types, emotional affects
- Examine pattern changes over time & membership to identify developmental trajectories & emergence of group/community characteristics

Our next developments will be guided by the following general principles:

- Building up of ontological knowledge bases through user defined text patterns and machine learning
- Customizable knowledge bases
- Visualization tools
- Deployment of multidimensional cluster analysis and other mining methods

Data Mining

Thank you!

lcp@cite.hku.hk

<http://lcp.cite.hku.hk>