

Social Movement Understanding by Keyword Tracking

Virach Sornlertlamvanich
virach@gmail.com

Kobkrit Viriyayudhakorn
kobkrit@gmail.com

OCT
2012

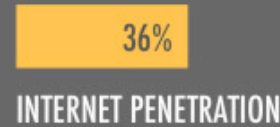
THAILAND



67,091,000
TOTAL POPULATION



24,000,000
INTERNET USERS



16,834,140
USERS ON TOP SOCIAL NETWORK



78,667,910
MOBILE SUBSCRIBERS

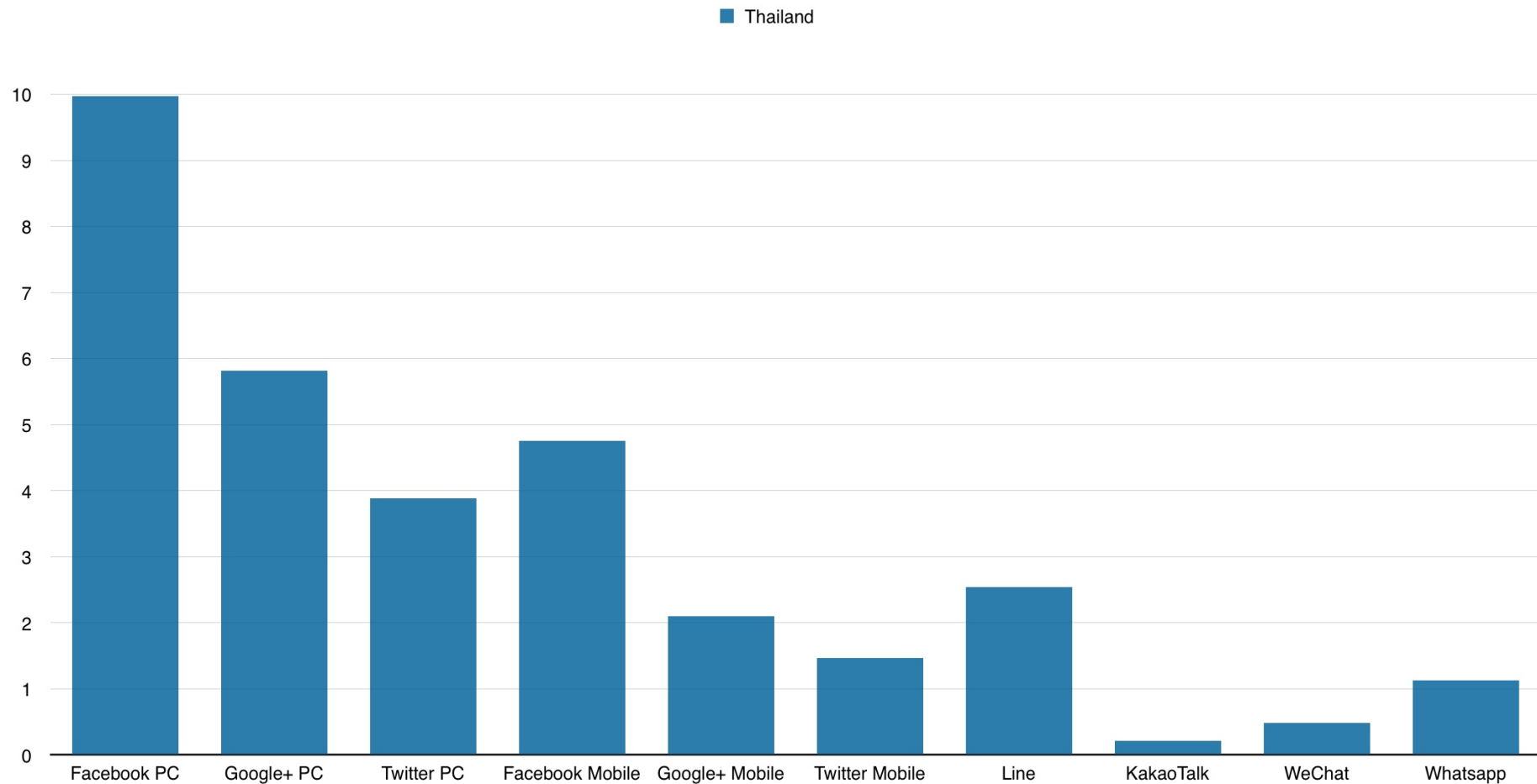


COMPILED BY @WEARESOCIALSG. FOR MORE SOCIAL MEDIA ANALYSIS AND INSIGHTS, VISIT OUR WEBSITE AT WEARESOCIALSG. SOURCES: POPULATION: BASED ON US CENSUS BUREAU (ACCESSED SEP 2012); URBANISATION: UN (2011); INTERNET: INTERNETWORLDSTATS (ACCESSED SEP 2012); SOCIAL NETWORKS: FACEBOOK (SEP 2012); MOBILE: ITU SUBSCRIBER FIGURES (LATEST AVAILABLE DATA, SEP 2012)

101

<http://tulaneict4d.wordpress.com/2013/04/05/social-media-in-thailand/>

Thailand Monthly Active Users on Social Network and Messenger App in 2013



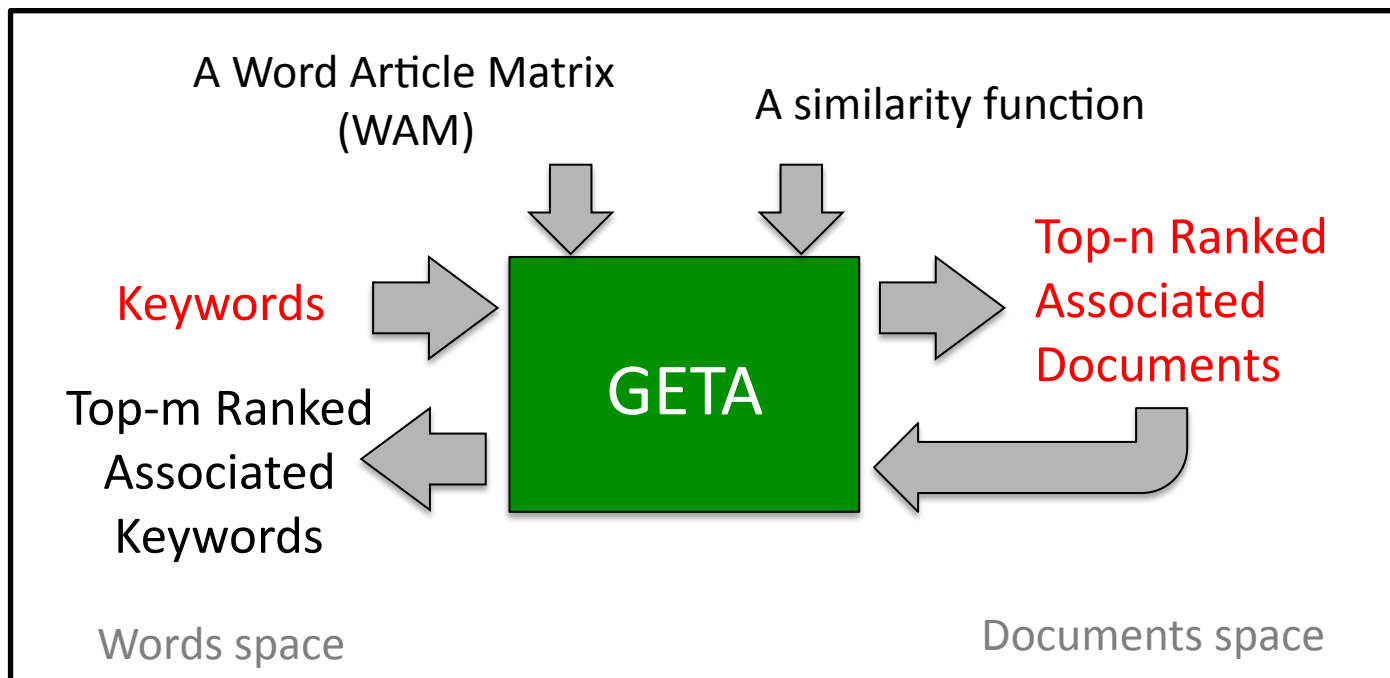
<http://www.chandlernguyen.com/2013/12/state-of-social-networks-in-southeast-asia.html>

Data Preparation

- Word segmentation
- Keyword
extracted from topic related documents
(training set)
- Tweeter inquiry
using the prepared topic related list of
keywords
- Text similarity
using GETA algorithm

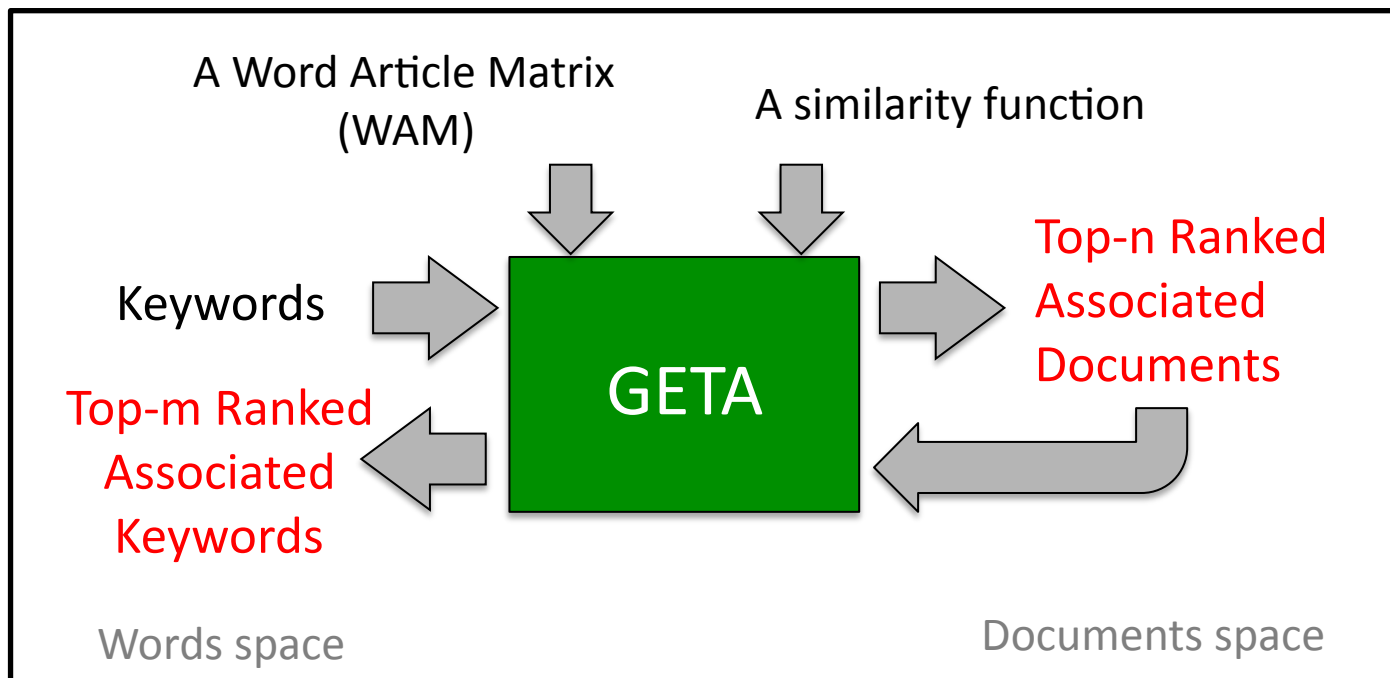
GETA: Association Search Engine

- Proposed by Akihiko Takano (Takano, 2003)
- Back-and-forth searches between words and documents spaces
 - Set of **keywords** → The top-n ranks of **associated documents**
 - Set of **documents** → The top-m ranks of **associated keywords**
- Require a **Word Article Matrix (WAM)**, and a **similarity function**



GETA: Association Search Engine

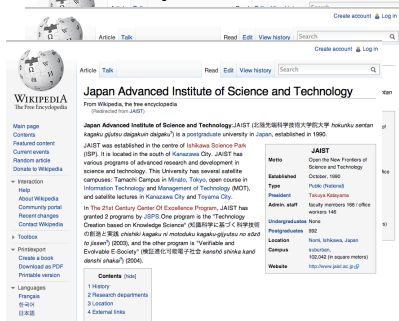
- Proposed by Akihiko Takano (Takano, 2003)
- Back-and-forth searches between words and documents spaces
 - Set of **keywords** → The top-n ranks of **associated documents**
 - Set of **documents** → The top-m ranks of **associated keywords**
- Require a **Word Article Matrix (WAM)**, and a **similarity function**



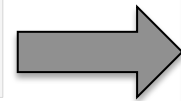
Word Article Matrix (WAM)

Building WAM is indexing

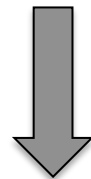
All Wikipedia Articles



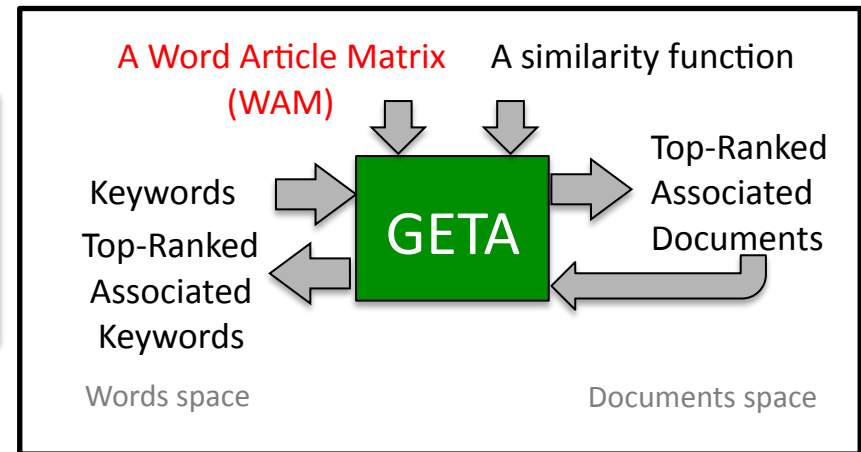
All Pages' Contents



Word seg. & Lemmatization
Thinking -> Think



All Words List



All Pages' Titles

Pages\Words	"Twitter"	"Tennis"	"Dollar"	"Google"	...
IT	2	0	1	4	
Sport	0	2	1	0	
Economics	0	0	2	0	
...					

A Wikipedia WAM

Similarity Functions

- **Ranking weight** in both words and documents spaces
- **Equation form**

$$\text{SIM}(d, q) = \sum_{t \in q} \frac{wq(t, q) \cdot wd(t, d)}{\text{norm}(d, q)}$$

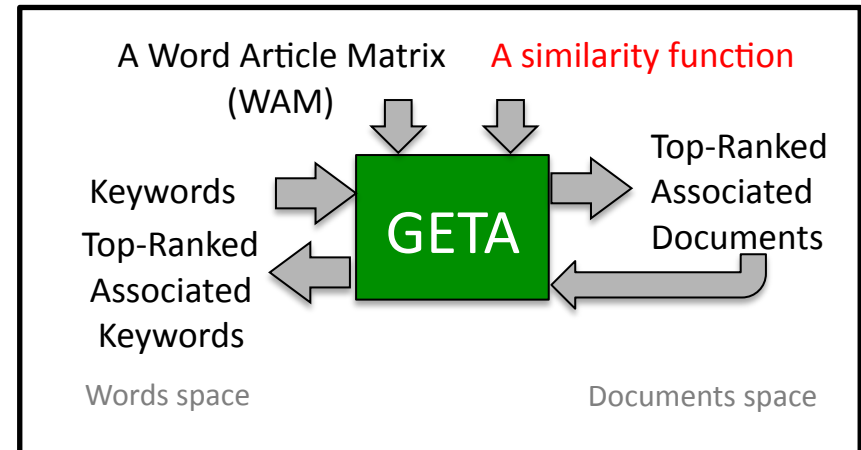
- $wd(t, d)$ = Weight of term t in document d , $wd(t, d) = 0$ if $t \notin d$
- $wq(t, q)$ = Weight of term t in query sentence q , $wq(t, q) = 0$ if $t \notin q$.
- $\text{norm}(d, q)$ = Normalization function due to the different length of d and q .

- **Smart measure (Singhal et.al., 1996)**

$$\frac{1}{\text{avg}(f_d) + \theta(f_d - \text{avg}(f_d))} \sum_{t \in q \wedge t \in d} \log\left(\frac{N}{f_t}\right) \cdot \frac{1 + \log(f_{d,t})}{1 + \log(\text{avg}_{w \in d}(f_{d,w}))} \cdot \frac{1 + \log(f_{q,t})}{1 + \log(\text{avg}_{w \in q}(f_{q,w}))}$$

- **Dot Product (Wilkinson et.al., 1996)**

$$\sum_{t \in q \wedge t \in d} (w_{q,t} \cdot w_{d,t})$$



Example GETA Calculation

SIM Function: Dot Product $\sum_{t \in q \wedge t \in d} (w_{q,t} \cdot w_{d,t})$

1. Twitter has 800M dollars.

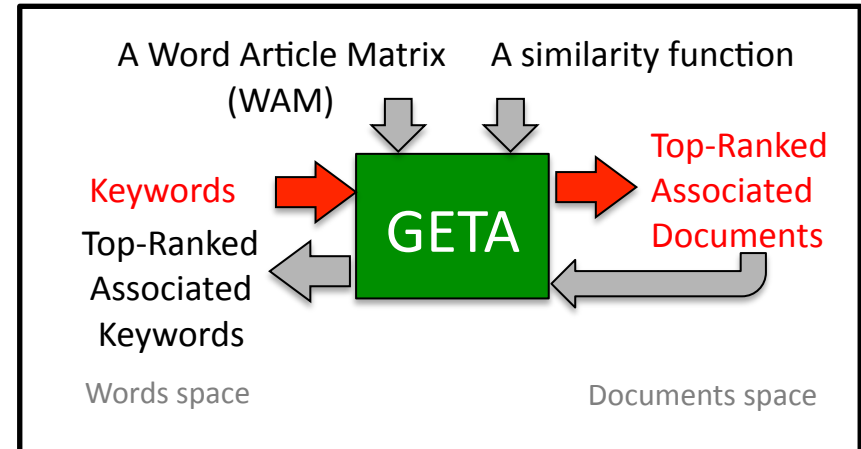
A idea label
"Twitter has 800M dollars"

Word seg. & Lemmatization
dollars -> dollar

"Twitter"	"Tennis"	"Dollar"	"Google"
1	0	1	0

Pages\Words	"Twitter"	"Tennis"	"Dollar"	"Google"	...
IT	2	0	1	4	
Sport	0	2	1	0	
Economics	0	0	2	0	
...					

Wikipedia WAM



(n=2) Select Top-n
Most Associated Documents Ranking

Pages	Dot Product Score
IT	3
Sport	1
Economics	2
...	

Example GETA Calculation

SIM Function: Dot Product $\sum_{t \in q \wedge t \in d} (w_{q,t} \cdot w_{d,t})$

1. Twitter has 800M dollars.

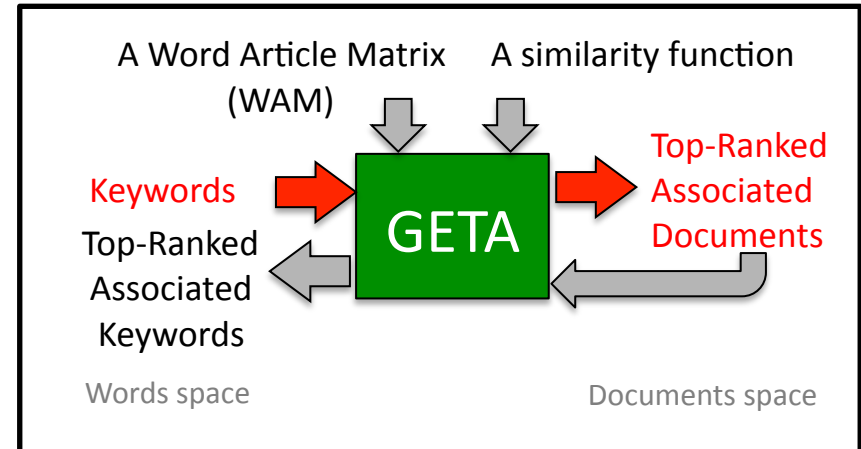
A idea label
"Twitter has
800M dollars"

Word seg. &
Lemmatization
dollars -> dollar

"Twitter"	"Tennis"	"Dollar"	"Google"
1	0	1	0

Pages\Words	"Twitter"	"Tennis"	"Dollar"	"Google"	...
IT	2	0	1	4	
Sport	0	2	1	0	
Economics	0	0	2	0	
...					

Wikipedia WAM



(n=2) Select Top-n
Most Associated Documents Ranking

Pages	Dot Product Score
IT	3
Sport	0
Economics	2
...	

Example GETA Calculation (Cont.)

SIM Function: Dot Product $\sum_{t \in q \wedge t \in d} (w_{q,t} \cdot w_{d,t})$

Top-2 Most Associated Documents

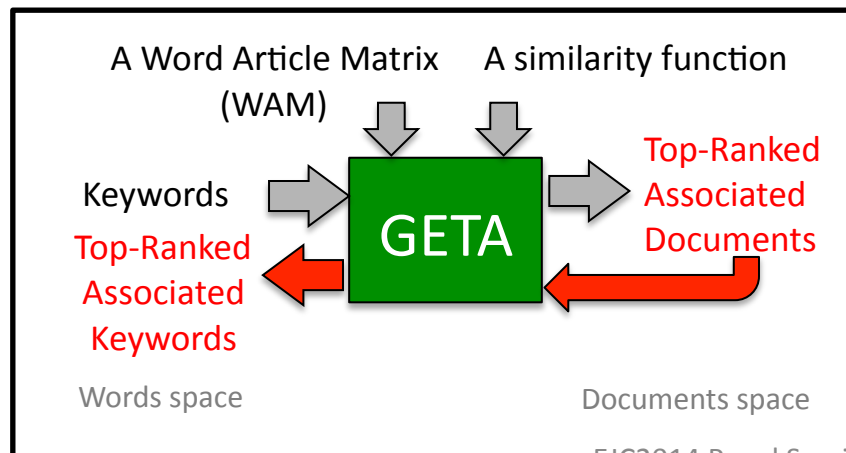
Pages	Dot Product Score	Pages\Words	"Twitter"	"Tennis"	"Dollar"	"Google"	...
IT	3	IT	2	0	1	4	
Sport	0	Sport	0	2	1	0	
Economics	2	Economics	0	0	2	0	
		...					

Output of GD

Wikipedia WAM

"Twitter"	"Tennis"	"Dollar"	"Google"
6	0	7	12

(m=3) Top-m Most Associated Keywords



Example GETA Calculation (Cont.)

SIM Function: Dot Product $\sum_{t \in q \wedge t \in d} (w_{q,t} \cdot w_{d,t})$

A idea label
"Twitter has
800M dollars"

1. Twitter has 800M dollars.

Google,
Dollar,
Twitter

Top-2 Most Associated Documents

Pages	Dot Product Score
IT	3
Sport	0
Economics	2



*Sport Pages = 0, since only Top-2 is selected from previous slide.

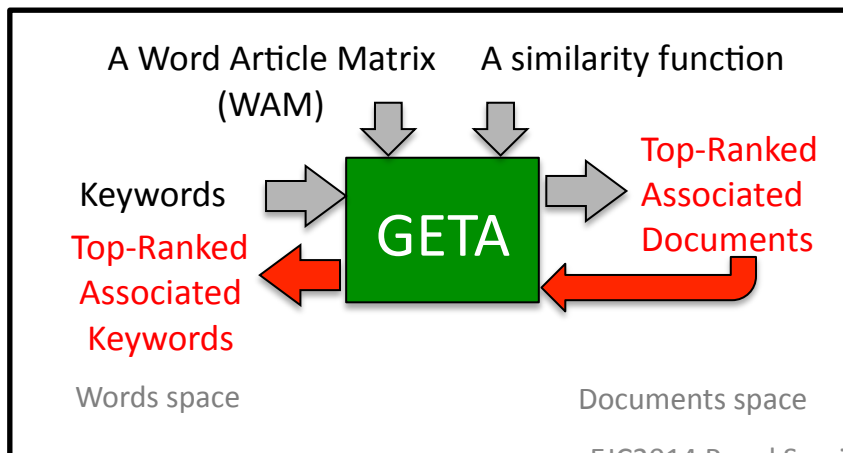
Pages\Words	"Twitter"	"Tennis"	"Dollar"	"Google"	...
IT	2	0	1	4	
Sport	0	2	1	0	
Economics	0	0	2	0	
...					

Wikipedia WAM

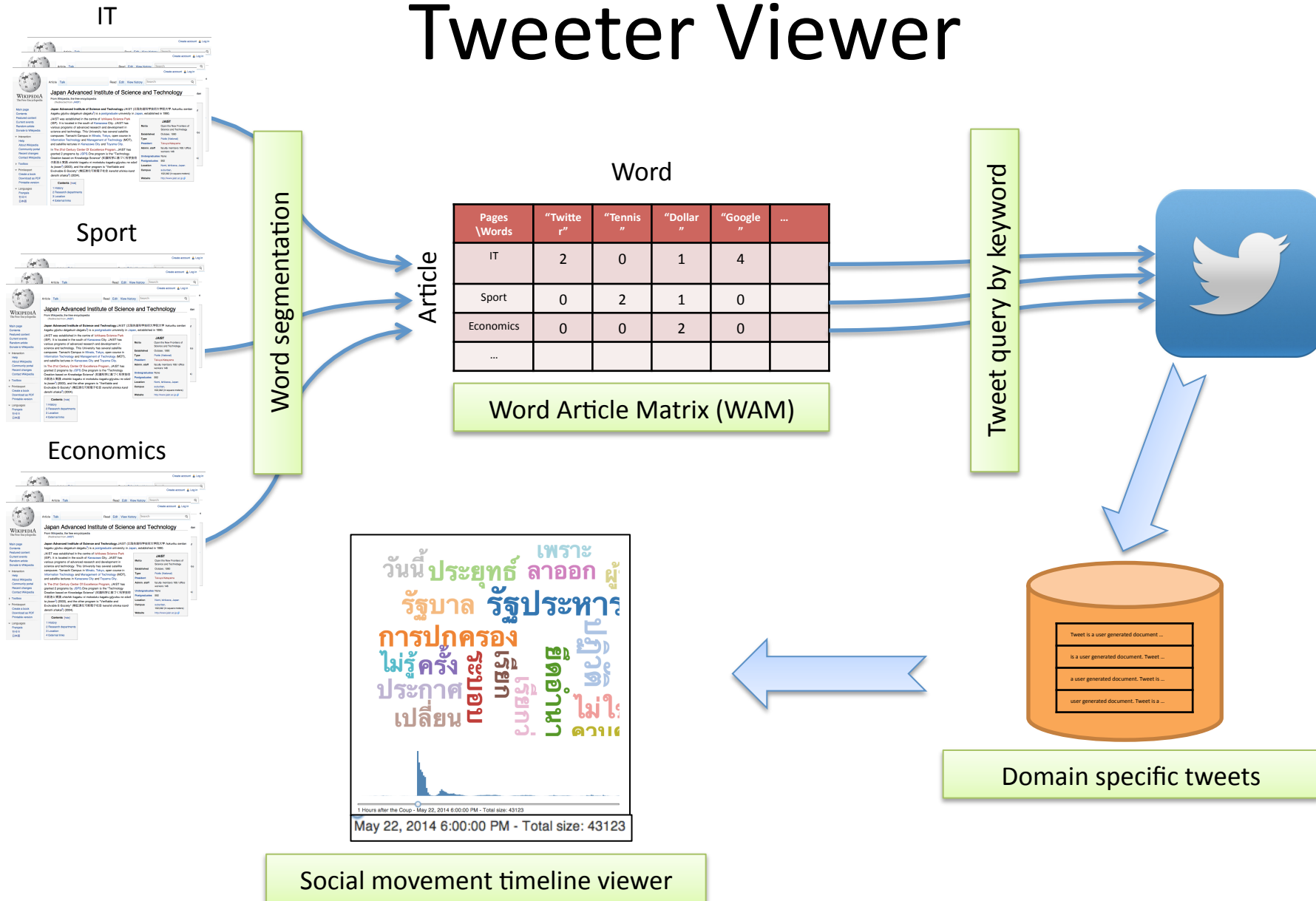


"Twitter"	"Tennis"	"Dollar"	"Google"
6	0	7	12

(m=3) Top-m Most Associated Keywords



Tweeter Viewer



Coup on May 22, 2014

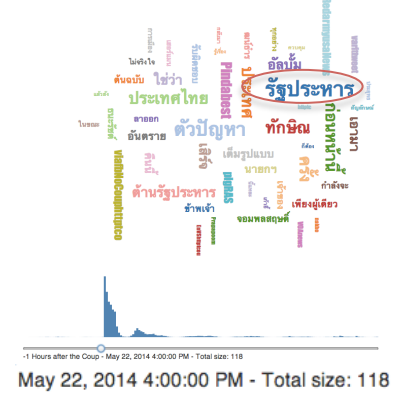
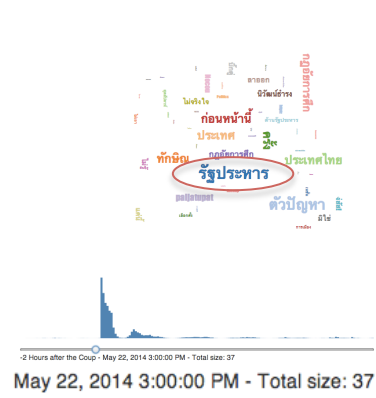
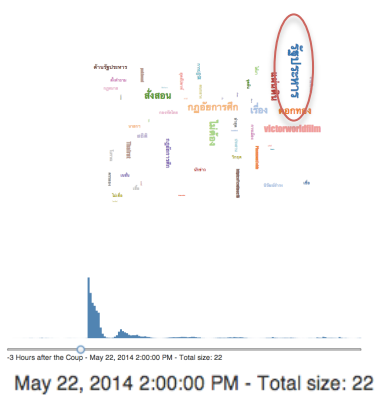
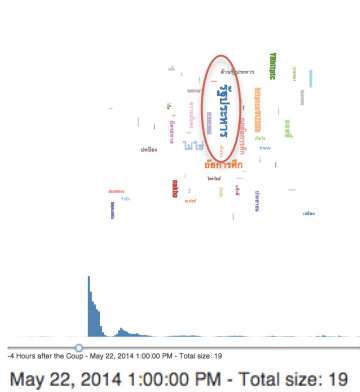
- ทหาร, คสช., ประเทศ, ประกาศ, สงบ, อำนาจ, รัฐบาล, รัฐประหาร, ชุมนุม, ตำรวจ, สถานการณ์, นายก, ควบคุม, ยึด, ประชุม, เศรษฐกิจ, กฎหมาย, คีท, แกนนำ, รัฐมนตรี, เลือกตั้ง, ประชาธิปไตย, ปฏิวัติ. ยึดอำนาจ. เคอร์ฟิว. กฎอัยการศึก
- military, NCPO, country, announce, peace, power, government, coup d'etat, gathering, police, situation, PM, control, seize, meeting, economy, law, war, leader, minister, election, democracy, revolution, seize the power, curfew, martial law

Tweet Query

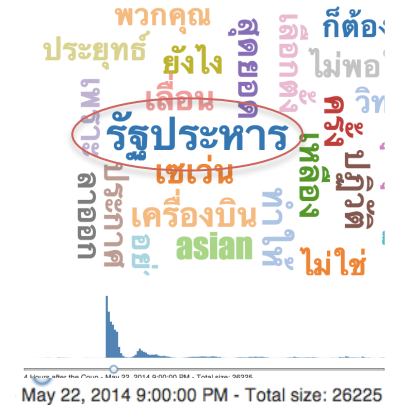
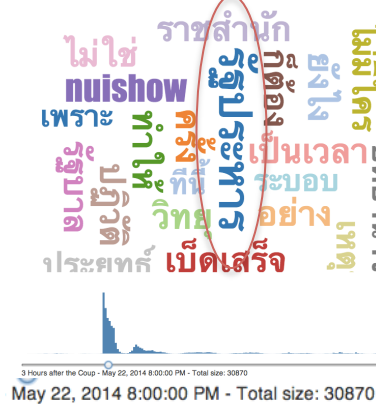
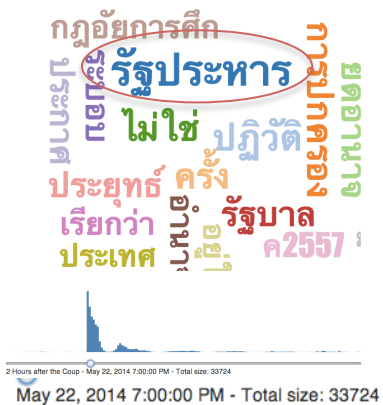
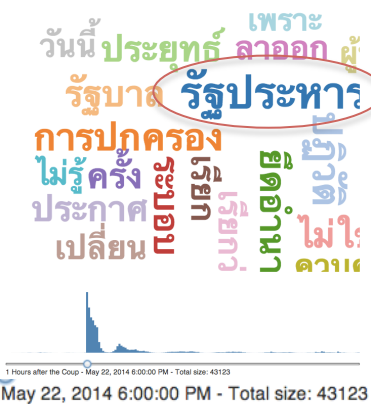
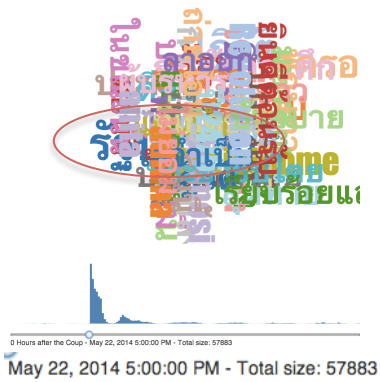
- Search Tweets by using Restful API
 - GET search/tweets. Set q = the keyword set
 - 100 tweets/search limited
 - Repeatedly fetch data until all tweets in the coup periods are discovered
- Be able to search back to 7 days

May 22, 2014 Coup-related tweet : 339,148 tweets

Timeline Word Cloud



Coup D'etat



Conclusion

- Key word expansion is effective to understand the short message i.e. tweet
- Key word expansion can be done using the known training corpus and GETA algorithm
- Timeline word cloud shows the development of the social movement e.g. some events are predictable