

# Thai Text Processing and Its Applications

Virach Sornlertlamvanich

Sirindhorn International Institute of Technology (SIIT),

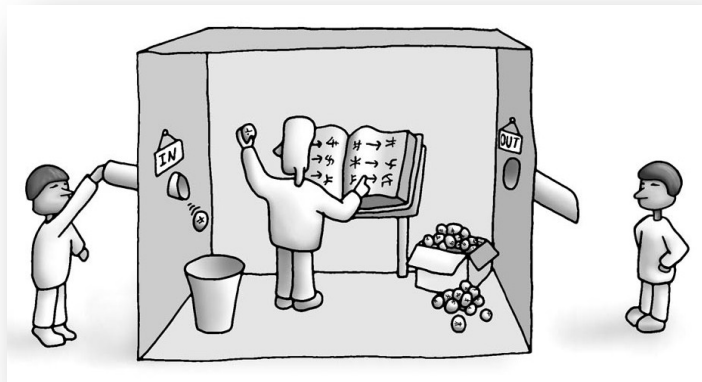
Thammasat University, Thailand

[virach@siit.tu.ac.th](mailto:virach@siit.tu.ac.th)

# NLP Challenges

- **Natural Language Processing** is no more for the **language study** only.

**Turing Test** as a tool to understand whether "**machines can think**".



NLP@Sapienza



AlphaGo

AlphaGo's algorithm uses a combination of **neural networks, machine learning and Monte Carlo tree search** techniques, combined with extensive training, both from human and computer play.

**Searle's Chinese Room argument** 1980 in "Minds, Brains, and Programs"

**Accurate NLP:** machine translation, summarization, machine reading, question answering, information retrieval, social text understanding, opinion mining, etc.

# NLP Challenges

- **Internet, Big Data, Machine Learning, Deep Learning** have brought along the possibilities.



## Facebook:-

Adds 0.5 petabyte ( $10^{15}$ ) of data every 24 hours

## Twitter:-

Adds 340 million tweets per day

## Youtube:-

Adds 100 hours of new videos every minute

Germin8, Social Intelligence

# NLP Challenges

## **NLP NEEDS BIG DATA**

Using Hadoop with NLTK

- Computational Linguistics methodologies are stochastic
- Examples are easier to create than rules
- Rules and Logic miss frequency and language dynamics
- Humans use lots of data for the same task- It's AI!
- More data is better - relevance is in the long tail
- If you don't have enough data - hire a knowledge engineer

## **BIG DATA WILL NEED NLP**

Using NLTK with Hadoop

- Hadoop is great at massive amounts of text data
- However, current methods aren't really NLP
- Indexing, Co-Occurrence, even N-Gram Modeling is search
- We haven't exhausted frequency analysis yet
- But when we do, we're going to want semantic analyses

Bird Steven, Edward Loper and Ewan Klein (2009), Natural Language Processing with Python. O'Reilly Media Inc.  
Data Community DC (DC2)



## Three NLP Fundamental Problems

- Word Segmentation/Tokenization
- Named Entity Recognition/Keyword Extraction
- Semantic Relation Extraction

## Three NLP Constructive Applications

- Linked Data  Knowledge Map
- Keyword Tracking  Social Movement
- Hyper Local News  Urban-Rural Info Gap

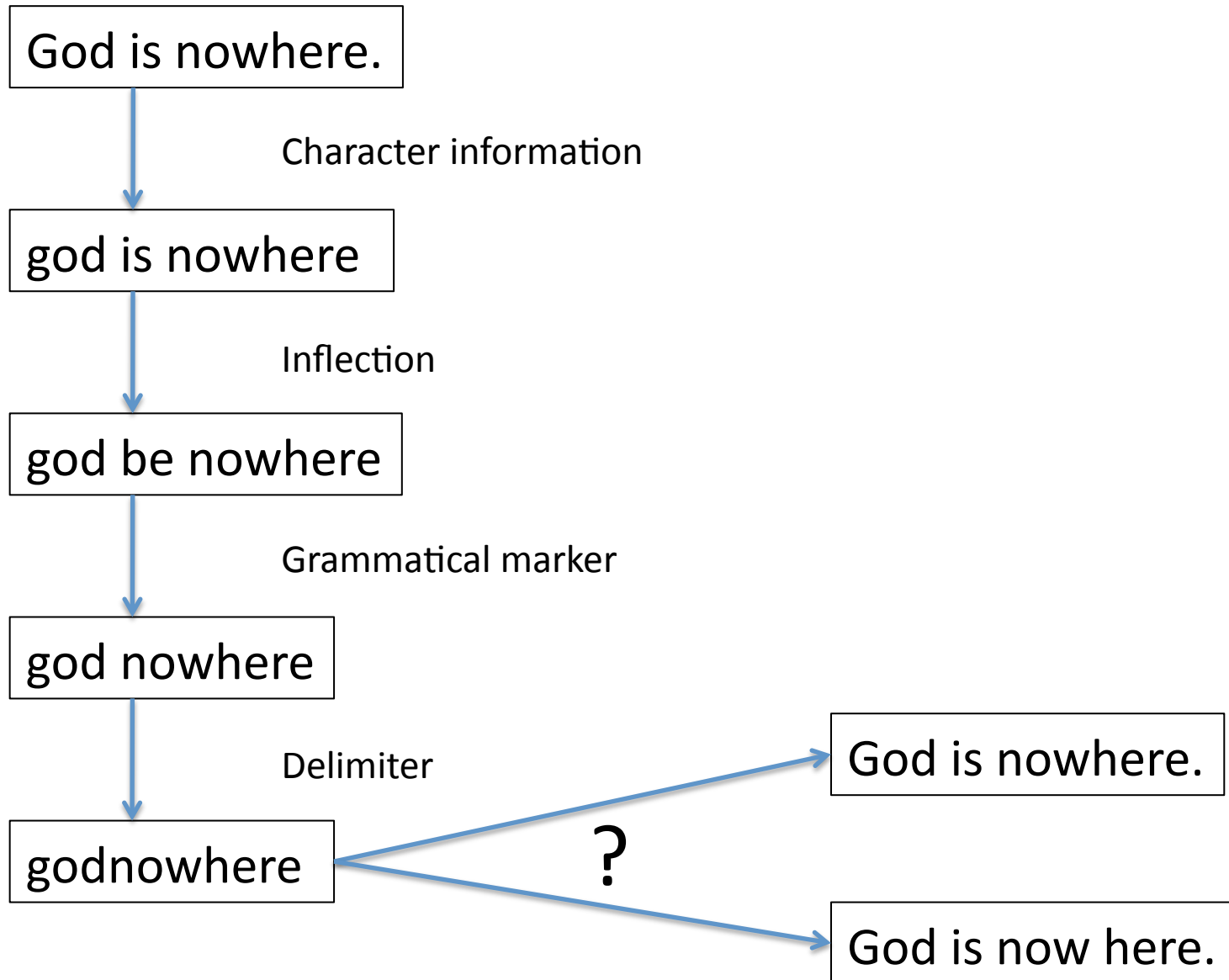
# **WORD SEGMENTATION**

# Thai Language as a Non-Segmenting Language

- No explicit word boundary marker e.g. capital letter, space character, punctuation mark, etc.
- No inflection
- No grammatical marker

How to determine word and sentence boundary?!?!

# Simulating a Non-Segmenting Language



# Word-Based Approach

## **Word segmentation (accuracy for Thai)**

- Longest matching: 92%
- Maximal matching: 93%
- POS tri-gram: 96%
- Machine learning: 97%

## **Sentence segmentation (accuracy for Thai)**

- POS tri-gram: 84.57%
- Feature-based approach (Winnow): 89.13%

# Character-Based Approach Term Extraction

- Automatic Corpus-Based Thai Word Extraction with the **C4.5** Learning Algorithm
- C4.5-trained decision tree for determining potential word boundary from **MI**, **Entropy** and some **linguistic information**
- Capable of discovering new words in document without assistance from static dictionary

# Attributes(1) : Left and Right Mutual Information

$$MI_L(xyz) = \frac{p(xyz)}{p(x)p(yz)}$$
$$MI_R(xyz) = \frac{p(xyz)}{p(xy)p(z)}$$

$x$   $yz$                        $xy$   $z$

where

$x$  is the leftmost character of string  $xyz$

$y$  is the middle substring of  $xyz$

$z$  is the rightmost character of string  $xyz$

$p( )$  is the probability function.

High mutual information implies that  $xyz$  co-occurs more than expected by chance. If  $xyz$  is a word, its  $MI_L$  and  $MI_R$  must be high.

...efunction... and ...function...

## Attributes(2) : Left and Right Entropy

$$H_L(\mathbf{y}) = - \sum_{\text{all } x \in \mathbf{A}} p(xy | \mathbf{y}) \cdot \log_2 p(xy | \mathbf{y})$$

$x$   $y$

$$H_R(\mathbf{y}) = - \sum_{\text{all } z \in \mathbf{A}} p(yz | \mathbf{y}) \cdot \log_2 p(yz | \mathbf{y})$$

$y$   $z$

where

$x$  is the leftmost character of string  $xyz$

$y$  is the middle substring of  $xyz$

$z$  is the rightmost character of string  $xyz$

$p( )$  is the probability function.

Entropy shows the variety of characters before and after a word.

If  $y$  is a word, its left and right entropy must be high.

...**?function**... and ...**?unction**...



# Attributes(3) : Frequency, Length Function Word

- **Frequency**  
Words tend to be used more often than non-word string sequences.
- **Length**  
Short strings are likely to happen by chance. The long and short strings should be treated differently.
- **Function Word**  
Function words are used mostly in phrases. They are useful to disambiguate words and phrases.

$$Func(s) = \begin{cases} 1 & \text{if } s \text{ contains a function word} \\ 0 & \text{otherwise} \end{cases}$$

## Attributes(4): First Two and Last Two Characters

- **Frequency of the first-two characters** of the considered string which appears in the first-two characters of words in the dictionary

high frequency -> the beginning of the considered string conforms to the Thai spelling system.

Ex.

*function*: how likely *fu* can be the beginning of word.

This idea can be also applied to **the last-two characters**.

# Experimental Results (1)

## The Precision of Word Extraction

	No. of strings extracted by the decision tree	No. of words extracted	No. of non-word strings extracted
Training Set	1882 (100%)	1643 (87.3%)	239 (12.7%)
Test Set	1815 (100%)	1526 (84.1%)	289 (15.9%)

## The Recall of Word Extraction

	No. of words that has more than 2 occurrences in corpus	No. of words extracted by the decision tree	No. of words in corpus that are found RID
Training Set	2933 (100%)	1643 (56.0%)	1833 (62.5%)
Test Set	2720 (100%)	1526 (56.1%)	1580 (58.1%)

Remark: These precision and recall are measured against 30,000 strings that occur more than 2 times in the corpus and conform to some simple Thai spelling rules.

# Experimental Results (2)

## Word Extraction VS. a Dictionary

	No. of words extracted by the decision tree	No. of words extracted by the decision tree which is in RID	No. of words extracted by the decision tree which is not in RID
Training Set	1643 (100.0%)	1082 (65.9%)	561 (34.1%)
Test Set	1526 (100.0%)	1046 (68.5%)	480 (31.5%)

Remark: RID is referred to the Thai-Thai dictionary published by The Royal Institute in 1982.

# **KEYWORD AND SEMANTIC RELATION EXTRACTION**

ICICTES 2016, Pullman Hotel, March 20-22, 2016

# Semantic Link Generation

- Semantic Representation of the description
  - **Keyword Extraction**
    - Extract keywords in text documents and link them to appropriate articles
  - **Semantic Relation Extraction**
    - Extract commons syntactic patterns between two keywords and generalize them to a triple  $(e_i, r_{ij}, e_j)$
- Linked Data
  - Set of triple  $(e_i, r_{ij}, e_j)$

# Keyword Extraction

- Some keywords are readily available in the set **tags**, but many of them are still missing.
- Our task is to extract those missing keywords from the **description** and **title**.

## Focused NE in Cultural Domain Database

- Cultural attraction (Place)
- Cultural person (Person)
- Cultural artifact (Object)

# Training Data Preparation

- Generate a keyword list from **tags** and **titles** that are not shorter than 5 characters and not longer than 30 characters
- POS tagged **descriptions**
- Label the POS tagged descriptions with the keyword list



# Labeling

- Apply BIO tagging
  - B: beginning position of a keyword
  - I: intermediate (or end) position of a keyword
  - O: other words
- If several matches are possible, select the longest one

# Training Data

- Description

ผ้าชั้นลายมัดหมี่บ้านปทุมแก้ว เป็นงานฝีมือพื้นบ้าน .....

# Training Data

- Description

ผ้าชั้นลายมัดหมี่บ้านปทุมแก้ว เป็นงานฝีมือพื้นบ้าน .....

- Segmented/Tagged/Labeled Description

Word	POS tag	Label
ผ้าชั้น	N	B-K
ลายมัดหมี่บ้านปทุมแก้ว	N	I-K
<space>	P	O
เป็น	V	O
งานฝีมือพื้นบ้าน	N	O
.....		.....

- Keyword List

(extracted from tag and title)

.....

ผ้า

.....

ผ้าชั้น

ผ้าชั้นลายมัดหมี่บ้านปทุมแก้ว

.....

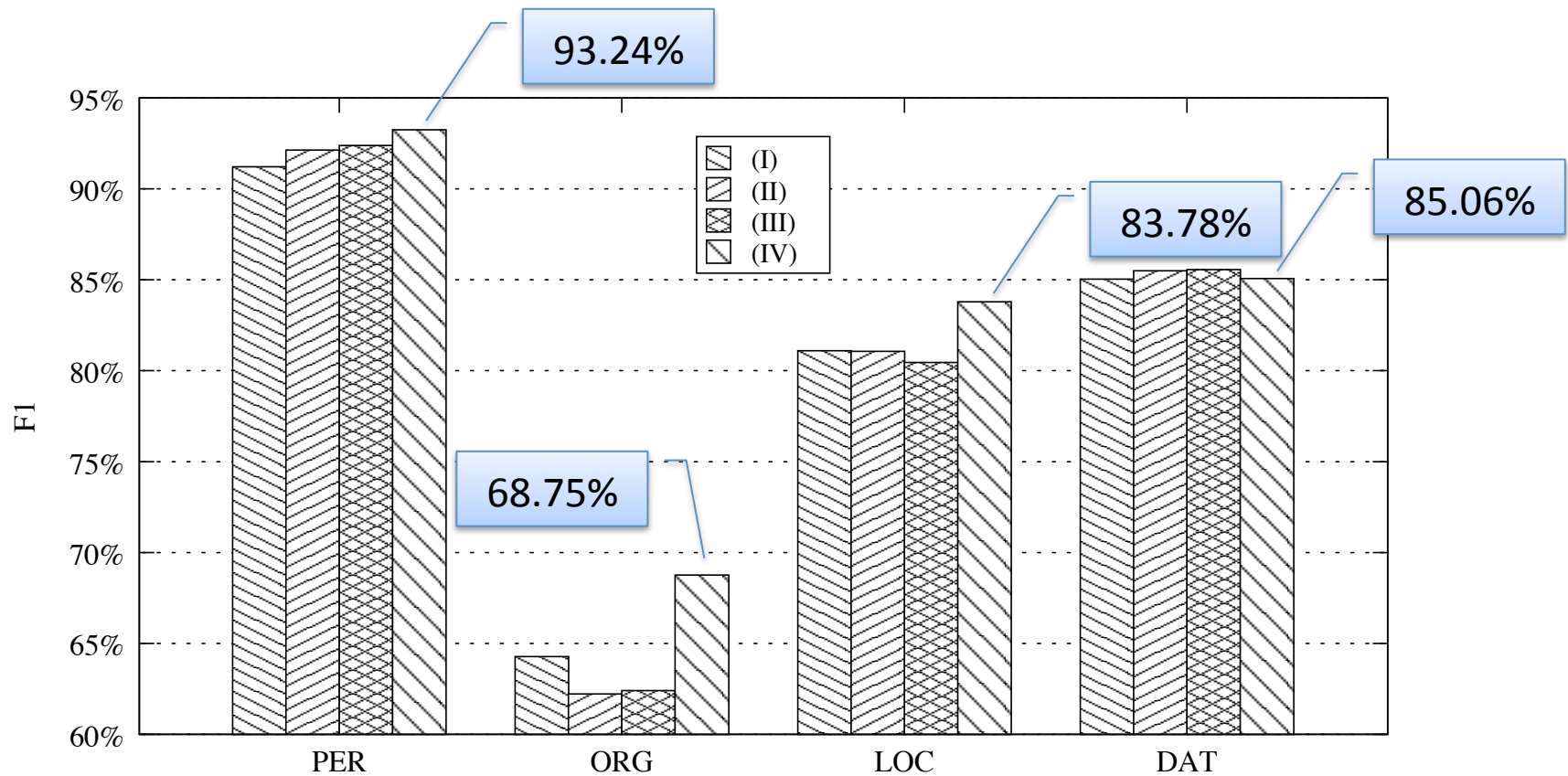
.....

# Preliminary Experiment Result

- Based on Margin Infused Relaxed Algorithm (MIRA), Crammer et al., 2005

<b>(I): word 1, 2 grams + label bigrams</b> $\langle w_j \rangle, j \in [-2, 2] \times y_0$ $\langle w_j, w_{j+1} \rangle, j \in [-2, 1] \times y_0$ $\langle y_{-1}, y_0 \rangle$	<b>(III): (II) + POS 3 grams</b> $\langle p_j, p_{j+1}, p_{j+2} \rangle, j \in [-2, 0] \times y_0$
<b>(II): (I) + POS 1,2 grams</b> $\langle p_j \rangle, j \in [-2, 2] \times y_0$ $\langle p_j, p_{j+1} \rangle, j \in [-2, 1] \times y_0$	<b>(IV): (III) + k-char prefixes/suffixes</b> $\langle P_k(w_0) \rangle, k \in [2, 3] \times y_0$ $\langle S_k(w_0) \rangle, k \in [2, 3] \times y_0$ $\langle P_k(w_0), S_k(w_0) \rangle, k \in [2, 3] \times y_0$

# F-Measure



Average evaluation result on NE annotated online news  
Recall=0.8256, Precision=0.9061, F1=0.8640

# Semantic Relation Acquisition

- Extract common syntactic patterns between two nouns
- Our task is to acquire triples  $(e_i, r_{ij}, e_j)$ , where
  - $e_i$  and  $e_j$  are entities (keywords)
  - $r_{ij}$  is a relationship between them

# Semantic Relation Template

Domain	Relation	Surface	Argument
Cultural attraction	ISLOCATEDAT	ตั้งอยู่ที่	LOC
	ISBUILTIN	สร้าง(ขึ้น)*ใน สร้าง(ขึ้น)*เมื่อ ตั้ง(ขึ้น)*เมื่อ	DATE
	ISBUILTBY	สร้าง(ขึ้น)*โดย ตั้ง(ขึ้น)*โดย	PER, ORG
	HASOLDNAME	เดิมชื่อ ชื่อเดิม	LOC, ORG
Cultural person	MARRIEDWITH	สมรสกับ	PER
	HASFATHERNAME	บิดาชื่อ	PER
	HASMOTHERNAME	มารดาชื่อ	PER
	HASOLDNAME	เดิมชื่อ ชื่อเดิม	PER
	HASBIRTHDATE	เกิด(เมื่อ)*	DATE
	BECOMEMONKIN	อุปสมบทเมื่อ	DATE
Cultural artifact	ISMADEBY	ผลิต(ขึ้น)*โดย ทำ(ขึ้น)*โดย ผลงานโดย	PER, ORG
	ISSOLDAT	จำหน่ายที่	LOC, ORG

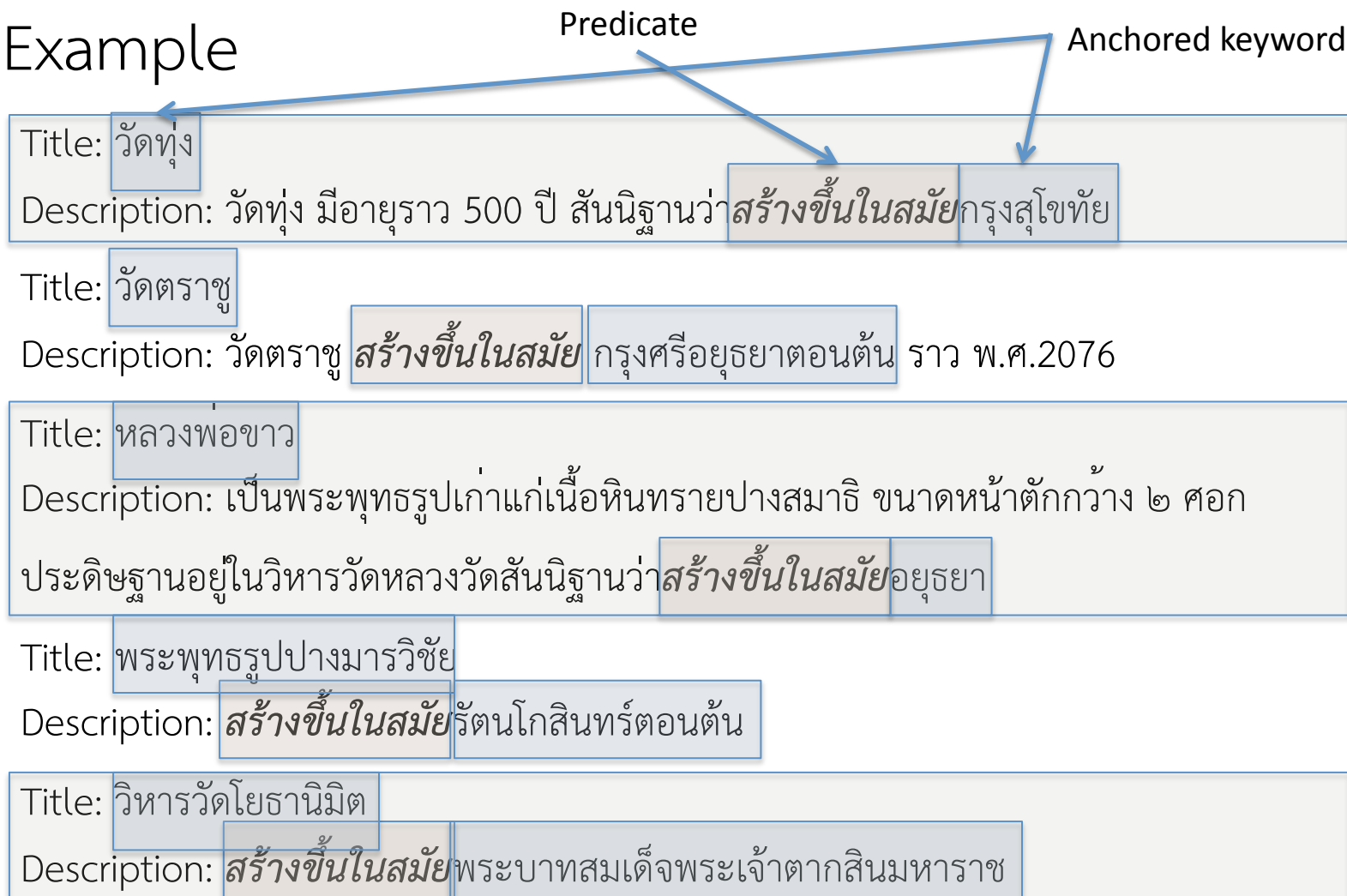
# Relation Instances Found in Word Distance

Relation	Argument	Distance					
		0	1	2	3	4	5
Cultural attraction							
ISLOCATEDAT	LOC	356	574	591	624	678	757
ISBUILDIN	DATE	3825	11487	11538	11573	11633	11667
ISBUILDBY	PER, ORG	131	202	218	234	249	257
HASOLDNAME	LOC, ORG	0	9	21	26	27	29
Cultural person							
MARRIEDWITH	PER	132	177	177	177	177	177
HASFATHERNAME	PER	120	372	372	373	373	373
HASMOTHERNAME	PER	97	383	383	383	383	383
HASOLDNAME	PER	51	259	273	277	277	283
HASBIRTHDATE	DATE	4122	4745	4801	4947	4966	5075
BECOMEMONKIN	DATE	346	435	435	436	436	436
Cultural artifact							
ISMADEBY	PER, ORG	62	107	109	125	129	130
ISSOLDAT	LOC, ORG	31	31	56	59	62	64



# Extract Common Syntactic Pattern of a Predicate between Two Keywords

- Example



# Extract Common Syntactic Pattern of a Predicate between Two Keywords

- Example

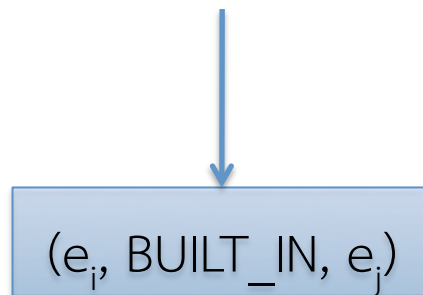
(วัดทุ่ง, *สร้างขึ้นในสมัย*, กรุงเทพมหานคร)

(วัดตราขู, *สร้างขึ้นในสมัย*, กรุงเทพมหานคร)

(หลวงพ่อกว, *สร้างขึ้นในสมัย*, ออยุธยา)

(พระพุทธรูปปางมารวิชัย, *สร้างขึ้นในสมัย*, รัตนโกสินทร์ตอนต้น)

(วิหารวัดโยธานิมิต, *สร้างขึ้นในสมัย*, พระบาทสมเด็จพระเจ้าตากสินมหาราช)



# Extract Common Syntactic Pattern of a Predicate between Two Keywords

- Example

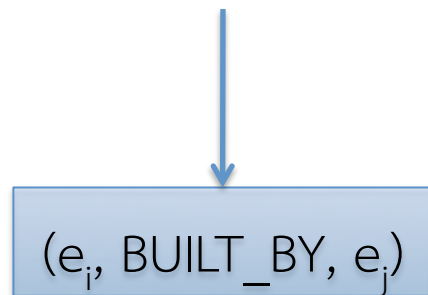
(กระโจมไฟบ้านโรงถ่าน, *สร้างโดย*, อพท.)

(ศาลเจ้าตากสิน วัดบ้านค่าย, *สร้างโดย*, พระครูพิพัฒน์ชยาภรณ์)

(วัดทุ่งไธ้, *สร้างโดย*, กลุ่มชาวลาวพวน)

(ศาลพระพรหม, *สร้างโดย*, เทศบาลนครเชียงราย)

(วงเวียนนิมิตร, *สร้างโดย*, เทศบาลนครภูเก็ต)



# Accuracy of Relation Extraction

Relation	Argument	#Sample	#Correct	#Incorrect	Accuracy
Cultural attraction					
ISLOCATEDAT	LOC	50	49	1	98%
ISBUILDIN	DATE	50	48	2	96%
ISBUILDBY	PER, ORG	50	48	2	96%
HASOLDNAME	LOC, ORG	27	23	4	85%
Cultural person					
MARRIEDWITH	PER	50	49	1	98%
HASFATHERNAME	PER	50	48	2	96%
HASMOTHERNAME	PER	50	49	1	98%
HASOLDNAME	PER	50	47	3	94%
HASBIRTHDATE	DATE	50	48	2	96%
BECOMEMONKIN	DATE	50	50	0	100%
Cultural artifact					
ISMADEBY	PER, ORG	50	44	6	88%
ISSOLDAT	LOC, ORG	50	49	1	98%

NLP Applications

# **KNOWLEDGE MAP GENERATION**

I hope we will  
use the Net to  
cross barriers  
and connect  
cultures.

*Tim Berners Lee*

meetville.com

# Types of Semantic Relation

description

title



**พระเจดีย์กลางน้ำ** <sup>2</sup>  
**รายละเอียด** <sup>3</sup>  
 เจดีย์กลางน้ำตั้งอยู่ที่ตำบลปากน้ำ อำเภอเมืองระยอง จังหวัดระยอง มีลักษณะเป็นเจดีย์ทรงระฆังฐานกลม กว้าง 4 เมตร สูง 10 เมตร มีกำแพงรอบฐานเจดีย์สองชั้น ตั้งอยู่บนเกาะกลางแม่น้ำระยอง ท่ามกลางป่าชายเลนที่ยาวเหยียด มีน้ำล้อมรอบ เนื้อที่ประมาณ 52 ไร่ เทศบาลนครระยองได้สร้างสะพานเชื่อมพระเจดีย์กับฝั่ง เจดีย์กลางน้ำเป็นสถานที่ประกอบประเพณีท้องถิ่นของชาวระยองมาแต่โบราณคือ ประเพณีทอดกฐินและห่มผ้าองค์เจดีย์ ในกลางเดือน 12 ของทุกปี ผู้สร้างเจดีย์ คือ เจ้าเมืองระยอง ในสมัยรัชกาลที่ 4 สันนิษฐานว่าสร้างในช่วง พ.ศ.2403 - 2404 ...

**หมวดหมู่** <sup>4</sup>  
 โบราณสถาน, แหล่งท่องเที่ยว

tag

Domain	Relation	Surface	Argument
Cultural attraction	ISLOCATEDAT	ตั้งอยู่ที่	LOC
	ISBUILTIN	สร้าง(ขึ้น)*ใน สร้าง(ขึ้น)*เมื่อ ตั้ง(ขึ้น)*เมื่อ	DATE
	ISBUILTBY	สร้าง(ขึ้น)*โดย ตั้ง(ขึ้น)*โดย	PER, ORG
	HASOLDNAME	เดิมชื่อ ชื่อเดิม	LOC, ORG
Cultural person	MARRIEDWITH	สมรสกับ	PER
	HASFATHERNAME	บิดาชื่อ	PER
	HASMOTHERNAME	มารดาชื่อ	PER
	HASOLDNAME	เดิมชื่อ ชื่อเดิม	PER
	HASBIRTHDATE	เกิด(เมื่อ)*	DATE
	BECOMEMONKIN	อุปสมบทเมื่อ	DATE
Cultural artifact	ISMADEBY	ผลิต(ขึ้น)*โดย ทำ(ขึ้น)*โดย ผลงานโดย	PER, ORG
	ISSOLDAT	จำหน่ายที่	LOC, ORG

# Knowledge Map

ISBUILTIN(พระเจดีย์กลางน้ำ, พ.ศ.2403)

ISLOCATEDAT(พระเจดีย์กลางน้ำ, ตำบลปากน้ำ)



**พระเจดีย์กลางน้ำ** 2

**รายละเอียด** 3


เจดีย์กลางน้ำตั้งอยู่ที่ตำบลปากน้ำ อำเภอเมืองระยอง จังหวัดระยอง มีลักษณะเป็นเจดีย์ทรงระฆังฐานกลม กว้าง 4 เมตร สูง 10 เมตร มีกำแพงรอบฐานเจดีย์สองชั้น ตั้งอยู่บนเกาะกลางแม่น้ำระยอง ท่ามกลางป่าชายเลนที่ยาวเหยียด มีน้ำล้อมรอบ เนื้อที่ประมาณ 52 ไร่ เทศบาลนครระยองได้สร้างสะพานเชื่อมพระเจดีย์กับฝั่ง เจดีย์กลางน้ำเป็นสถานที่ประกอบประเพณีท้องถิ่นของชาวระยองมาแต่โบราณคือ ประเพณีทอดกฐินและห่มผ้าองค์เจดีย์ ในกลางเดือน 12 ของทุกปี ผู้สร้างเจดีย์ คือ เจ้าเมืองระยอง ในสมัยรัชกาลที่ 4 สันนิษฐานว่าสร้างในช่วง พ.ศ.2403 - 2404 ...

**หมวดหมู่** 4

โบราณสถาน, แหล่งท่องเที่ยว

Infobox

ชื่อ
พระเจดีย์กลางน้ำ



พระเจดีย์กลางน้ำ

---

สร้างใน

พ.ศ.2403

---

ที่ตั้ง

ตำบลปากน้ำ

Knowledge map



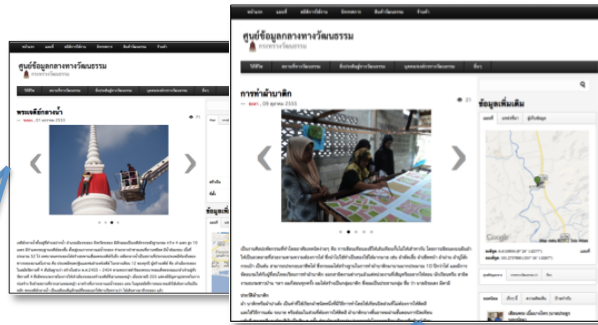


# Knowledging Semantically Enhanced Cultural Database [Place, Person, Artifact]

## Infobox



## Knowledge map



## Creator



## Making



## Product



## Shop



# **SOCIAL MOVEMENT UNDERSTANDING**

ICICTES 2016, Pullman Hotel, March 20-22, 2016

# Data Data Data!!!

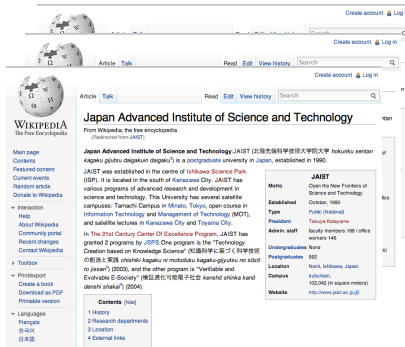
- Drastically increase number of users on social network
- Keywords in the contents express the concepts of the talk
- Social media texts are input in a time sequence
- But, social media texts are normally short, incomplete and diverse



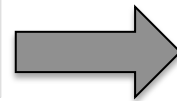
# Word Article Matrix (WAM)

Creating WAM

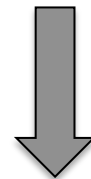
Wikipedia Articles



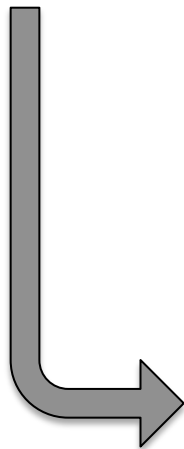
Page Contents



Word seg. & Lemmatization  
Thinking -> Think



Word list



Page Title

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

Wikipedia WAM

ICICTES 2016, Pullman Hotel, March 20-22, 2016

# Text Similarity

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

1. Twitter has 800M dollars.

“Twitter has 800M dollars”

Word seg. & Lemmatization  
dollars -> dollar

“Twitter”	“Tennis”	“Dollar”	“Google”
1	0	1	0

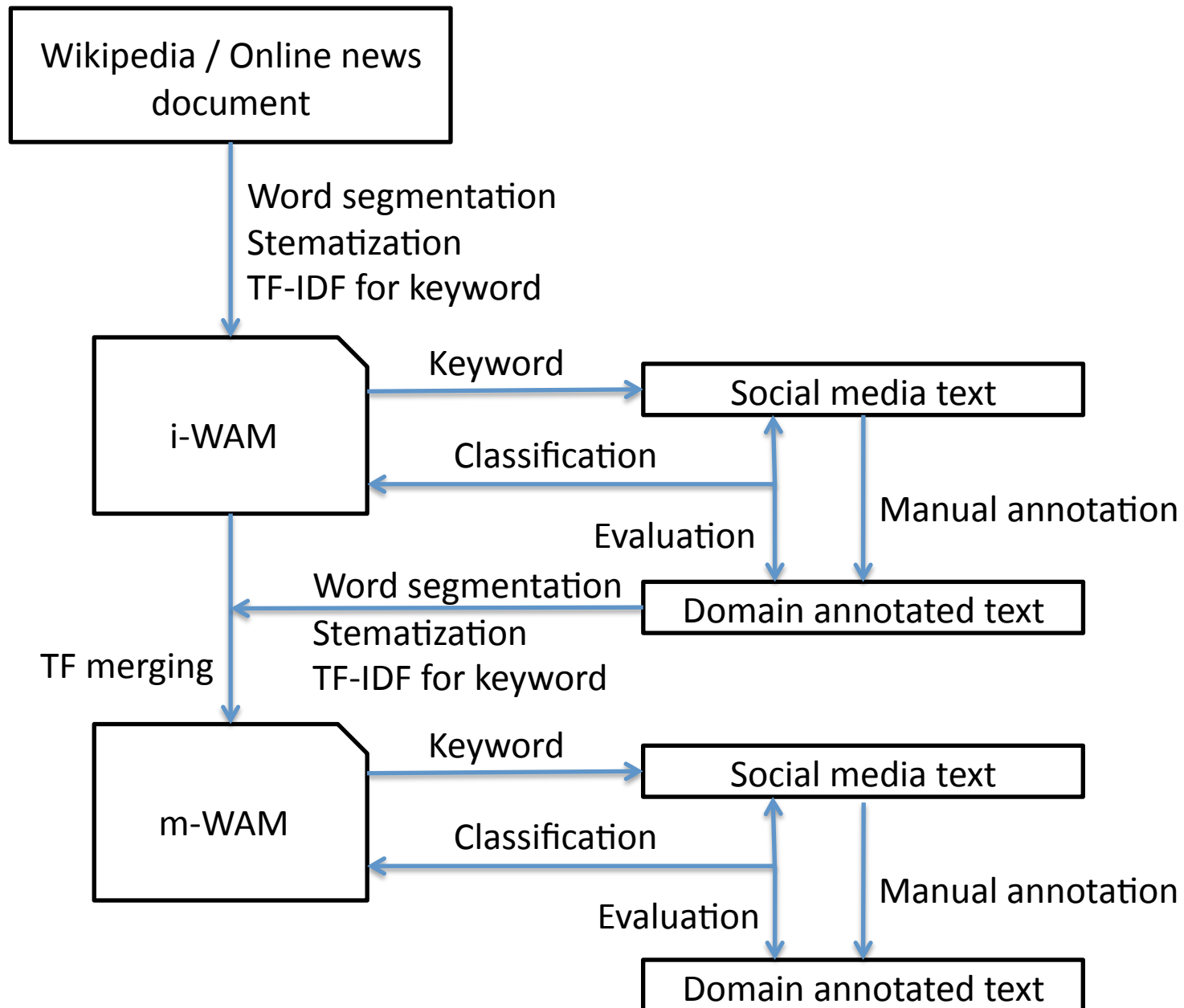
Pages\Words	“Twitter”	“Tennis”	“Dollar”	“Google”	...
IT	2	0	1	4	
Sport	0	2	1	0	
Economics	0	0	2	0	
...					

Wikipedia WAM

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

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

# Modified WAM for Social Media Text Classification



# A Part of Modified WAM

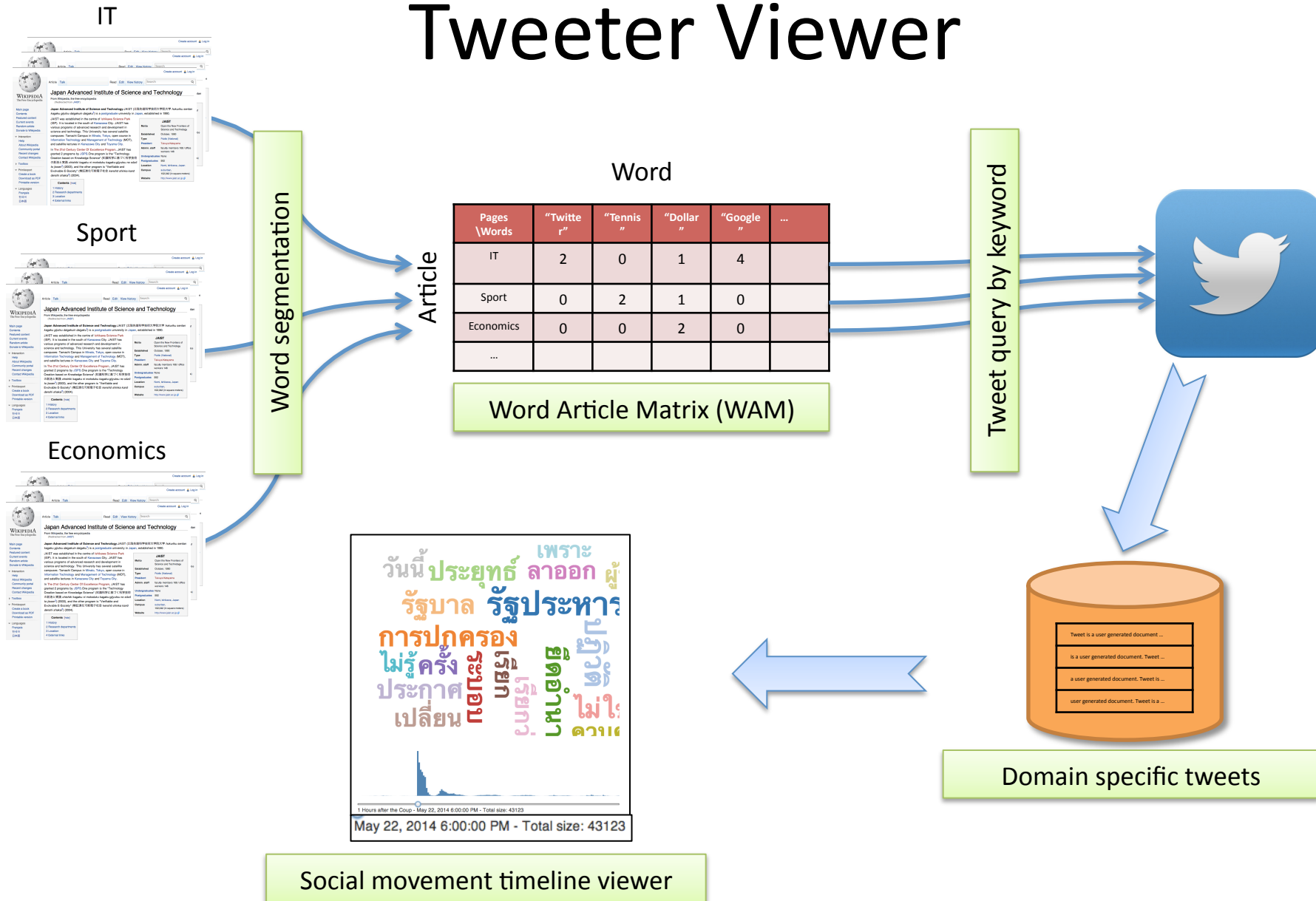
i-WAM	อิเทลด 'intel'	ศัลยกรรม 'surgery'	เมตตา 'mercy'	คลาวด์ 'cloud'	คลิป 'clip'	เน็ต 'network'
Life	0.052	0.205	0.160	0.168	0	0
Education	0	0.025	0.036	0	0	0
Technology	0.103	0	0	0.230	0	0
m-WAM	อิเทลด 'intel'	ศัลยกรรม 'surgery'	เมตตา 'mercy'	คลาวด์ 'cloud'	คลิป 'clip'	เน็ต 'network'
Life	0.026	0.275	0.268	0.177	0.036	0.040
Education	0	0.013	0.018	0	0	0
Technology	0.164	0	0	0.281	0.009	0.187
m-WAM2	อิเทลด 'intel'	ศัลยกรรม 'surgery'	เมตตา 'mercy'	คลาวด์ 'cloud'	คลิป 'clip'	เน็ต 'network'
Life	0.013	0.255	0.241	0.174	0.119	0.034
Education	0	0.009	0.011	0	0.009	0.019
Technology	0.107	0	0	0.238	0.019	0.143

# F-Measure

	Life	Education	Technology
i-WAM on text1	86.79%	36.09%	55.79%
m-WAM on text1	86.64%	34.57%	56.01%
m-WAM on text2	91.33%	29.91%	43.17%
m-WAM2 on text2	94.93%	29.47%	43.72%



# 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, เลือกลง, 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

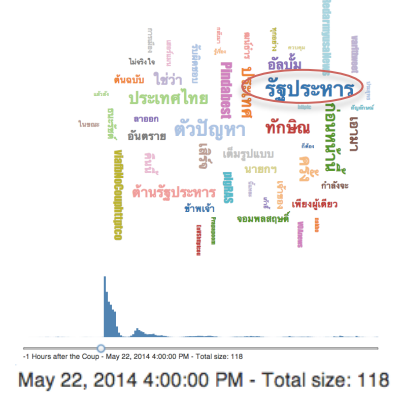
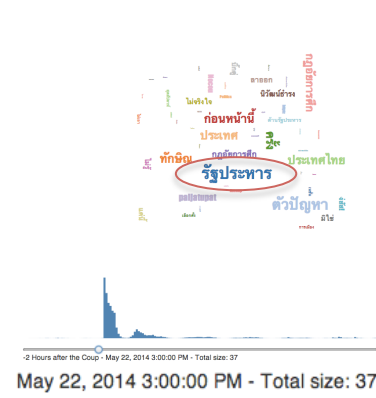
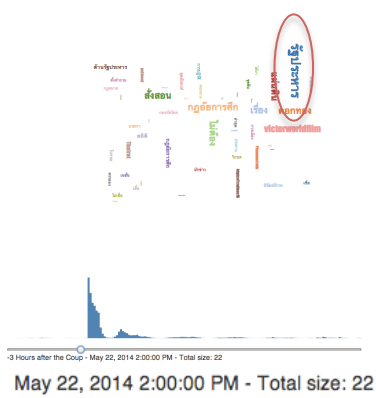
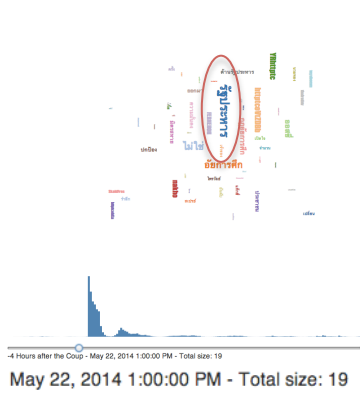
# Tweet Query

- 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

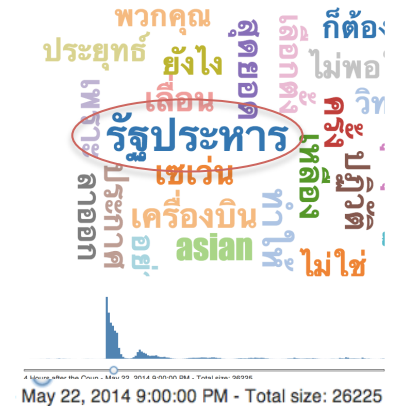
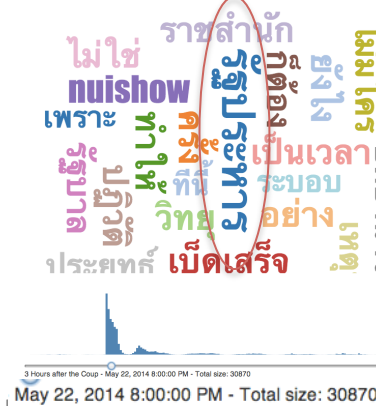
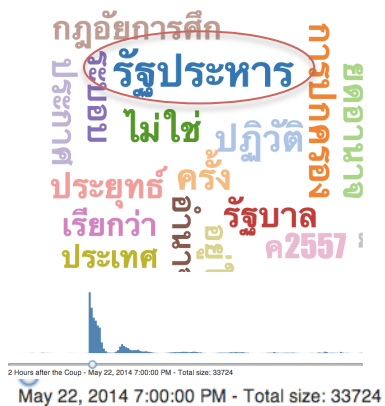
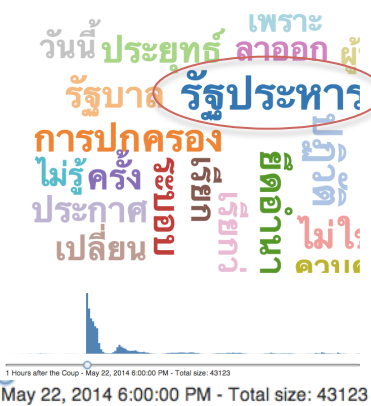
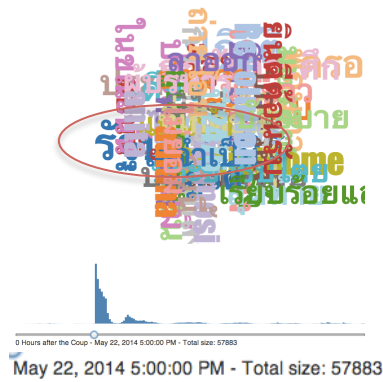
May 22, 2014 Coup-related tweet : 339,148 tweets

<http://sn.iisilab.org/>

# Timeline Word Cloud



Coup D'etat



# **URBAN-RURAL INFORMATION GAP FILLER**

ICICTES 2016, Pullman Hotel, March 20-22, 2016

# Purpose

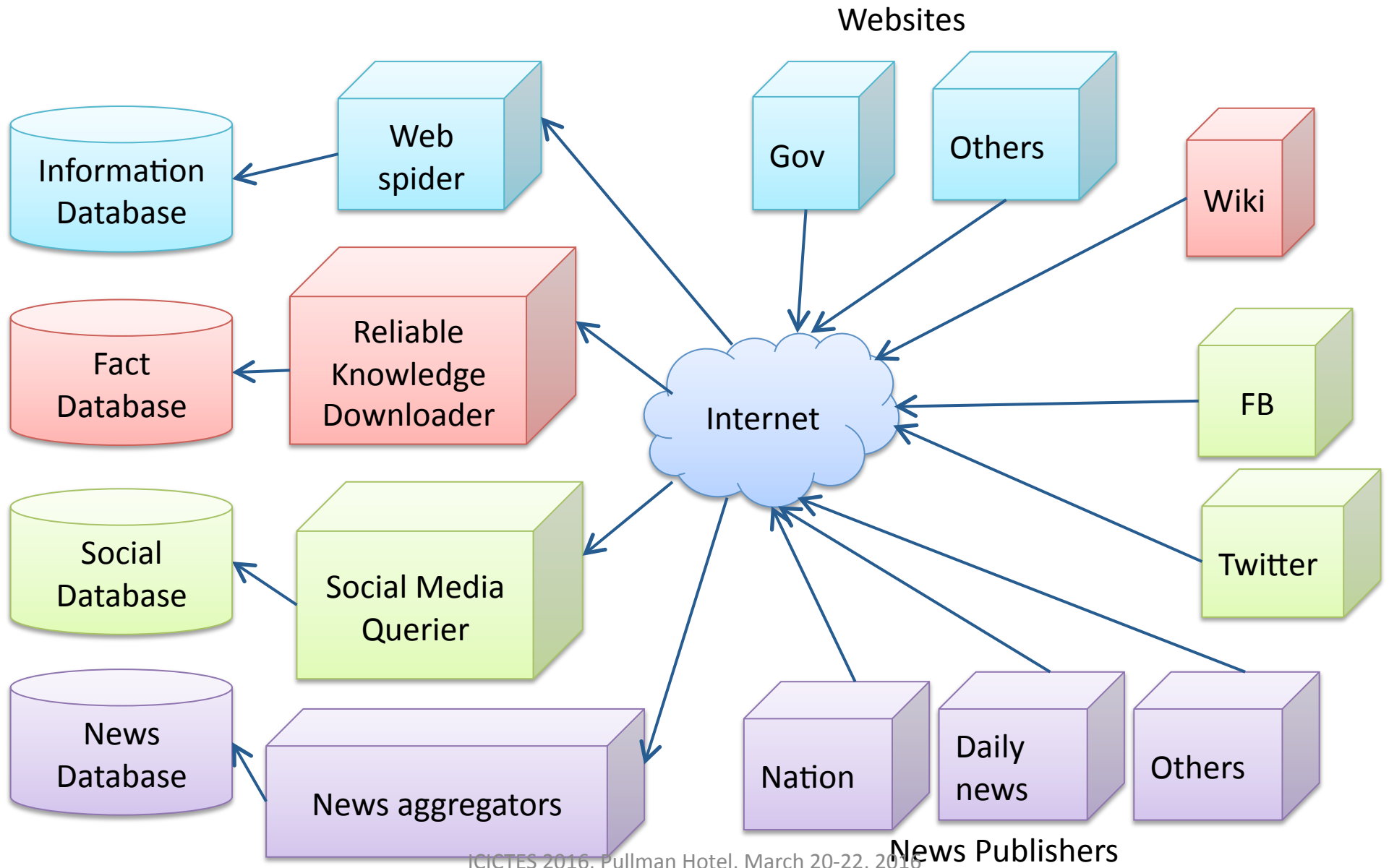
- To deliver local news to complement the deteriorating local newspaper
- To promote data usage in rural area



To improving rural livelihood through its reputation extraction!

<http://nation.iisilab.org>

# Collecting information

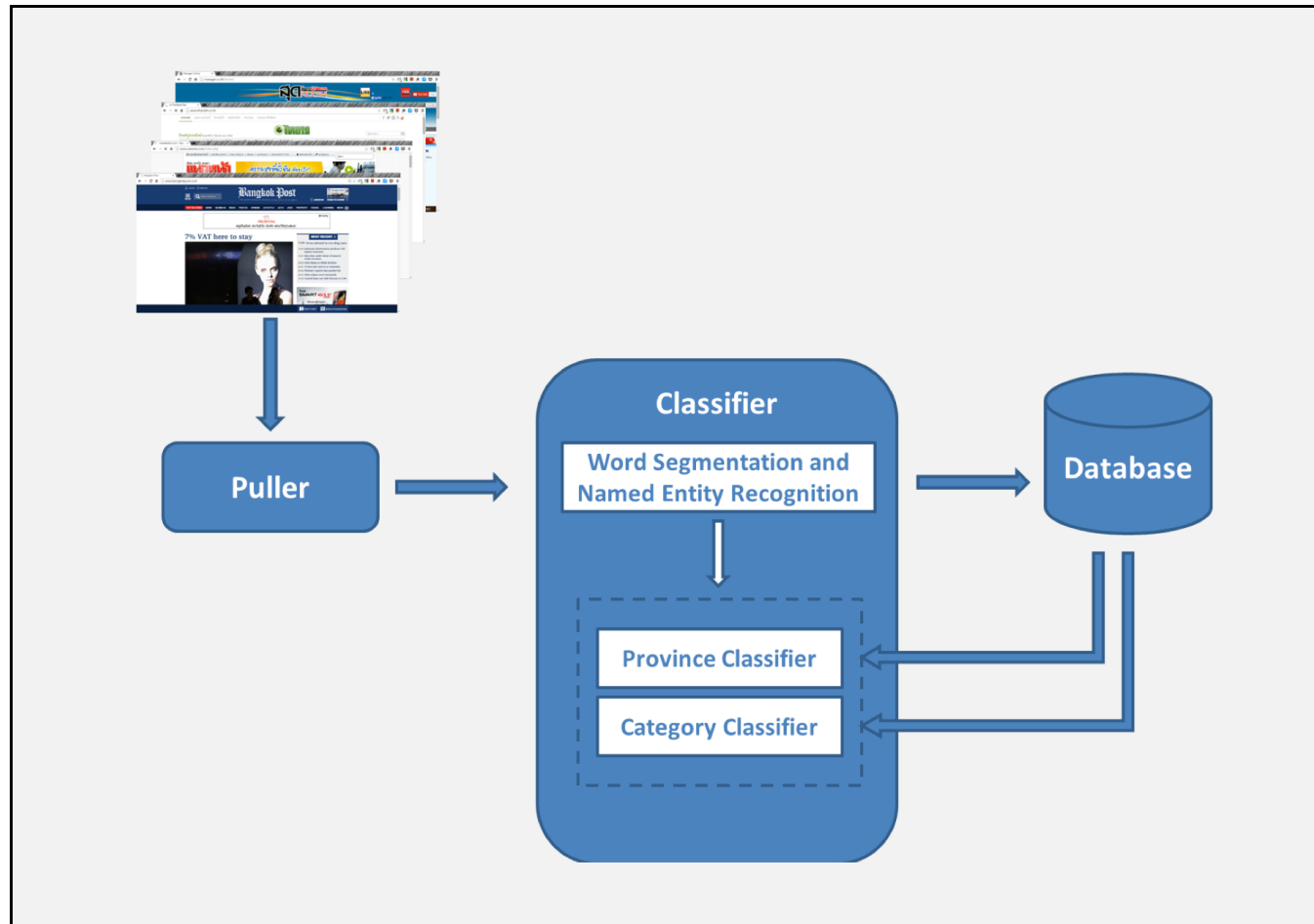


# Analyze

- **NEWS Classification**
  - Word segmentation
  - Named Entity Recognition
  - NEWS domain and province classification and ranking based on TF/IDF technique
- **Information Extraction**
  - Template based extraction
- **Infobox**
  - Celebrity i.e. politician, idol
  - Province
- **Social Media**
  - Trending topic
  - Leader and follower
  - Opinion polarity



# News Classification





31 พฤษภาคม 2555

**อุบลราชธานี**

เมนู (๓) อุบลราชธานี

ชาวทั้งหมด

ในหลวงทรงรับในพระฉายา

อ่านข่าวอื่น

สถานที่ที่เกี่ยวข้อง

MY HOME

น้ำตกสวยสุดสวย

เป็นน้ำตกขนาดใหญ่เกิดจากลำธาร 2 สาย คือ ห้วยสร้อย และห้วยไม้ที่ไหลจากหน้าผาตกลงมาบรรจบกัน...

ร้านอาหารทั้งหมด

PROMOTION

**Tesco Lotus**

โปรโมชั่น Te...  
ประจำสัปดาห์

**JUMBO**

โปรโมชั่น Big C

**Tops Supermarket**

โปรโมชั่น Tops Market Plus

ร่วมสนุกกิจกรรมเปิดตัว Zamborobic โดตัสศรีนครินทร์ สัมแบรนต์ ลูกโซคบอล 2014 ถิ่นรับ iPod Touch และกระเป๋าเงินใส่จาก No.1 Denim World ยิ่งซื้อ ยิ่งลุ้นไปกันเอง

บุคคลเด่น

เพชรประดับ อ.สิงห์ ประไพ (ชื่อเล่น: 30) ทัศนีย์ 14 กุมภาพันธ์ พ.ศ. 2537 เป็นนักจากอำเภอพิบูลมังสาหาร จังหวัดอุบลราชธานี (รับ)

ประกาศขอ-ขายที่ดิน บ้าน กวอนเฮาส์ กตม

<p>ขาย - กวอนโฮม คุณภาพราคา โครงการ ตระเวนสา สิ้น 11 แสมสุข, วารินชำราบ, อุบลราชธานี</p> <p><b>1,850,000 บาท</b></p>	<p>ขาย - กวอนโฮม ชยางกูร 99/79 ในเมือง, เมือง อุบลราชธานี, อุบลราชธานี</p> <p><b>2,300,000 บาท</b></p>	<p>ขาย - house for sale ubon ทาบใหญ่, เมือง อุบลราชธานี, อุบลราชธานี</p> <p><b>1,800,000 บาท</b></p>	<p>ขาย - บ้านสร้างใหม่ 2 ชั้น ราคาไม่เกินล้าน สมเด็จพระ (อุบล-ศาลสูง) 8 ปทุม, เมือง อุบลราชธานี, อุบลราชธานี</p> <p><b>990,000 บาท</b></p>
<p>ขาย - ปล่อยเช่า ทายกานต์ งามนัฒน์</p> <p>สุภาพพัฒนา ซอยไชยสิทธิ์ 1448/3 หมู่ที่ 18 ตำบล ทาบใหญ่, เมืองอุบลราชธานี, อุบลราชธานี</p> <p><b>1,500,000 บาท</b></p>	<p>ขาย - ขายกวอนโฮม ใกล้เคียงกับโรงเรียน อุดมศึกษา เพียง 1.6 กิโลเมตร</p> <p>พัลลภชัย ในเมือง, เมืองอุบลราชธานี, อุบลราชธานี</p> <p><b>1,600,000 บาท</b></p>	<p>ขาย - กวอนโฮม งามนัฒน์ งามนัฒน์</p> <p>แสมสุข, วารินชำราบ, อุบลราชธานี</p> <p><b>900,000 บาท</b></p>	<p>ขาย - บ้านสาวิตรี 10 ต.การสาร 31 ในเมือง, เมืองอุบลราชธานี, อุบลราชธานี</p> <p><b>3,500,000 บาท</b></p>

ดูประกาศซื้อ-ขายอื่นๆ ในหมวดทั้งหมด

วิดีโอ

<p>โรงเรียนพลาญทอง อ.สิรินธร จ.อุบลราชธานี</p> <p>30/5/2015 18:24:57</p>	<p>HL D2 อุบลราชธานี 0-1 อุบลเยี่ยมที่ 30-05-58</p> <p>30/5/2015 17:07:16</p>	<p>กิจกรรมวันงดสูบบุหรี่โลก ม.อุบลราชธานี</p> <p>30/5/2015 11:50:08</p>	<p>ปั่น ปลุค ป่าเมือง สวองงาม ครั้งที่ 15</p> <p>30/5/2015 03:12:08</p>
--	---	---	---

- ... (text partially obscured)
- ... (text partially obscured)
- ... (text partially obscured)
- ... (text partially obscured)

ดูงานอื่นๆ ในหมวดเอกชน

dtac Like

1,814,044 people like dtac.

Facebook social plugin

**จัดส่งฟรี**

เต็มสปีดทั่วไทย

ถึงทันใจภายใน 3 วันทำการ

# Summary

- **Word segmentation, Keyword extraction, NER, WSD, Semantic relation extraction etc.** are the basic issues of NLP in the current language environment.
- **Efficient computational scenario** is crucial under the drastic growth of data and coverage of the Internet.
- **User-generated contents** are the unlimited language resources.
- NLP technology helps in filling the **information gap** between the haves and the have-nots, filling the **technology gap** between the resource-rich and under-resourced languages.

# References

- Virach Sornlertlamvanich, Tanapong Potipiti and Thatsanee Charoenporn. **Automatic Corpus-based Thai Word Extraction with the C4.5 Learning Algorithm**. Proceedings of the 18th International Conference on Computational Linguistics (COLING2000), Saarbrucken, Germany, pp 802-807, July-August 2000.
- Virach Sornlertlamvanich, Thatsanee Charoenporn, Canasai Kruengkrai, and Hitoshi Isahara. **Statistical-based Approach to Non-segmented Language Processing**, IEICE Transactions on Information and Systems, Vol. E90-D, No.10, pp.1565-1573, October 2007.
- Canasai Kruengkrai, Virach Sornlertlamvanich, Watchira Buranasing, and Thatsanee Charoenporn. **Semantic Relation Extraction from a Cultural Database**, Proceedings of Workshop on South and Southeast Asian NLP, COLING2012, Mumbai, India, December 8-15, 2012.
- Virach Sornlertlamvanich and Kobkrit Viriyayudhakorn. **Social Movement Understanding by Keyword Tracking**, In Book: Information Modelling and Knowledge Bases XXVI, Volume 272, Eds. Thalheim, B. and et al., IOS Press, 2014.
- Virach Sornlertlamvanich and Canasai Kruengkrai. **Effectiveness of Keyword and Semantic Relation Extraction for Knowledge Map Generation**, Proceedings of The Second International Workshop on Worldwide Language Service Infrastructure (WLSI), Kyoto University, Kyoto, Japan, January 22-23, 2015.
- Virach Sornlertlamvanich. **Hyper Local NEWS Publishing. ---Collect, Analyze, Visualize---**, International Conference on Information Modelling and Knowledge Bases (EJC2015), Maribor, Slovenia, June 8-12, 2015.

# Acknowledgement

- Eakasit Pacharawongsakda
- Kobkrit Viriyayudhakorn
- Canasai Kruengkrai
- Thatsanee Charoenporn
- Tanapong Potipiti
- Sitdhibhong Laokok
- Watchira Buranasing
- Thanasan Tanhermhong
- Wirat Chinnan
- Ministry of Culture
- Claudia Soria (hint of images)