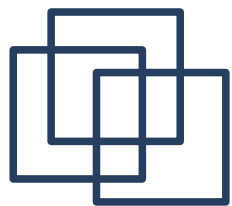


Search Engine ภาษาไทย

วิรัช ศรีเลิศล้ำวาณิช

Thai Computational Linguistics
Laboratory (TCL), NICT
virach@tcllab.org



Search Engine สำหรับภาษาอังกฤษ

➤ Lemmatization

➤ work:- work, working, works, worked

➤ go:- go, going, goes, went, gone

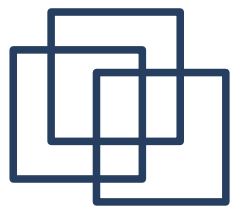
➤ Windows:- WINDOWS, Windows

➤ Scoring สำหรับการจัดลำดับบทความ

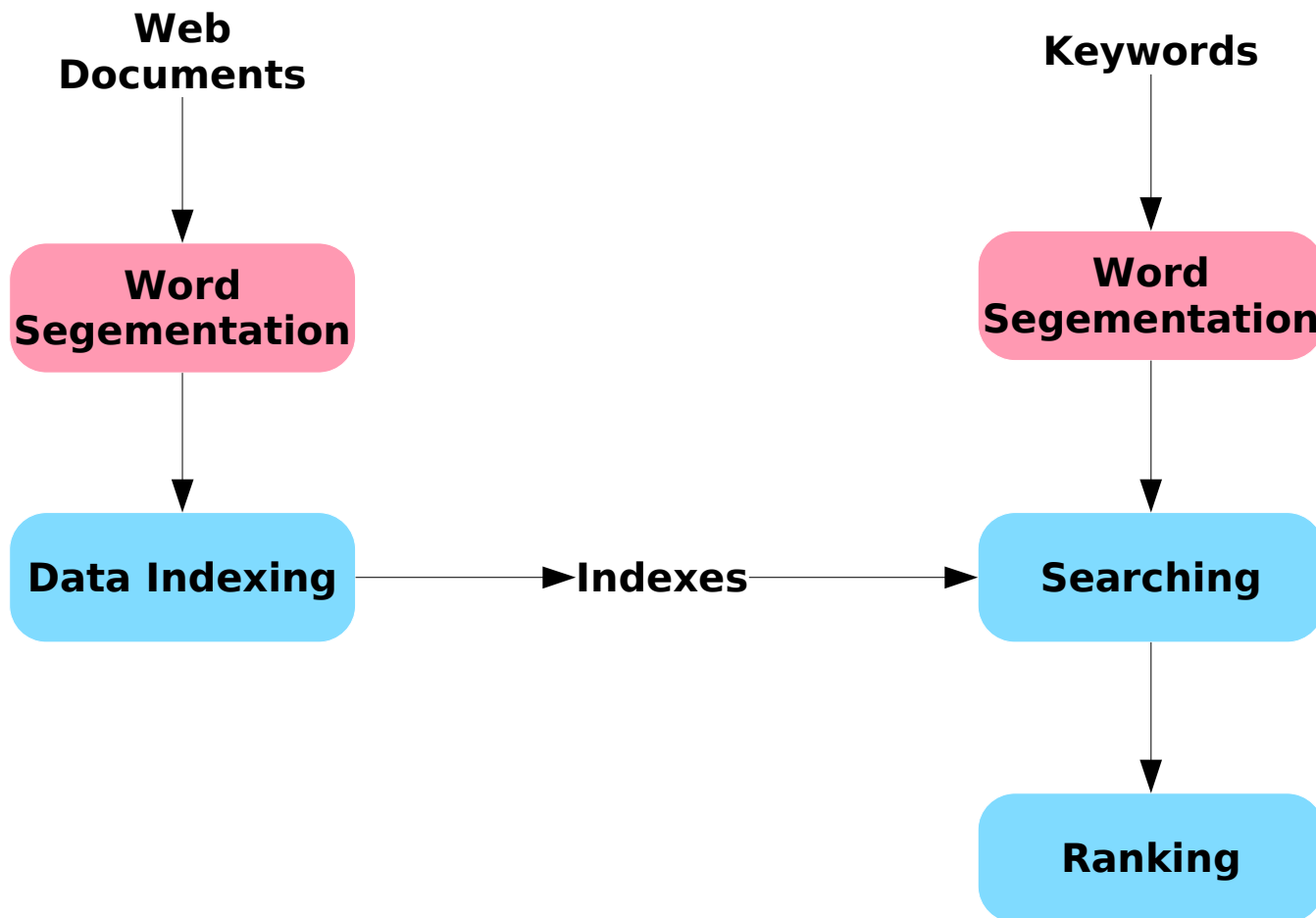
➤ ความสำคัญของคำในบทความ เช่น คำที่ปรากฏใน Title, คำสำคัญ (key word), ความถี่ของคำในไฟล์ เป็นต้น

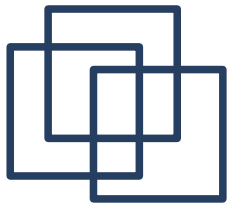
➤ ความนิยมของบทความ เช่น จำนวน click ที่เข้าแเวะชม, จำนวน link จากคำที่ใช้สืบค้น, page rank, เป็นต้น

➤ สร้าง Inverted Index File โดยอาศัยโปรแกรม Database



Dictionary-based Search Engine --Architecture--





ความยากในการตัดคำ

● ความกำกวม

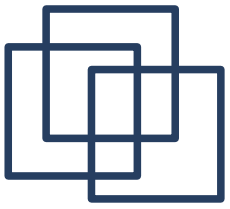
| | | |
|------------------|---------------------------|------------------------------|
| แบบนอก | -> แบบ นอก | <> แบ บน ออก |
| มีที่นา | -> มี ที่นา | <> มี ที่ นา |
| ขอบอกขอบใจ | -> ขอบอก ขอบใจ | <> ขอ บอก ขอบใจ |
| ขนมอบกรอบ | -> ขนม อบ กรอบ | <> ขน มอบ กรอบ |
| ร้านข้าวซอยลำดวน | -> ร้าน ข้าวซอย ลำดวน | <> ร้าน ข้าว ซอย ลำดวน |

● คำที่ไม่ปรากฏในพจนานุกรม

| | |
|---------------------|---------------------------------|
| นาตาลี่ | -> นา ตา ลี่ |
| อยุธยาอะลิอันซ์ซีพี | -> อยุธยา อะลิอันซ์ ซี พี |
| กาลิเลโอ | -> กา ลีเล โอ |

● การจัดเก็บข้อมูลในพจนานุกรม

อาทิ บัญชีคำ, หน้าที่ของคำ, เป็นต้น



Search Engine

อาศัยความน่าจะเป็นของคำ

กรอบ
ใน
ที่
ห้อง
ทำ
:

ครัว → "ครัว"

ความหลากหลายของ
คำที่อยู่ข้างเคียง

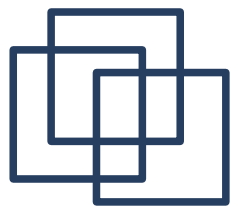
ค + รอบ → "กรอบ"

กรอบ + ครัว → "กรอบครัว"

ความบ่อยของการอยู่เคียง
ข้างกันของอักขระ

Ranking

- เปรียบเทียบความน่าจะเป็นของอันดับคำ
- Weight ตามค่าความสำคัญของคำ (key word, title, ...) และความถี่ (term frequency)



ความยากในการสืบค้นกรณีภาษาไทย

ครัว

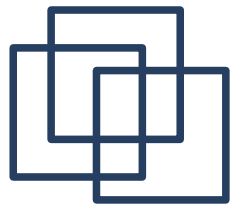
...การสมรสหมู่โดยสหพันธ์ครอบครัวเพื่อความสามัคคี...
...เปลี่ยนเป็นห้องน้ำ...ห้องครัว...ห้องรับแขก...

ประชา

...กองประชาสัมพันธ์การสื่อสารแห่งประเทศไทย...
...นายแพทย์ประชา เป็นประธานคณะกรรมการ...

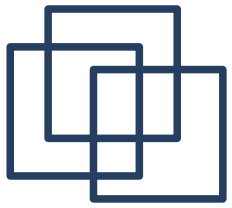
ธาตุ

...ประวัติวัดมหาธาตุวรวิหาร...
...โปรแกรมช่วยสอนเคมีเบื้องต้น และตารางธาตุ...

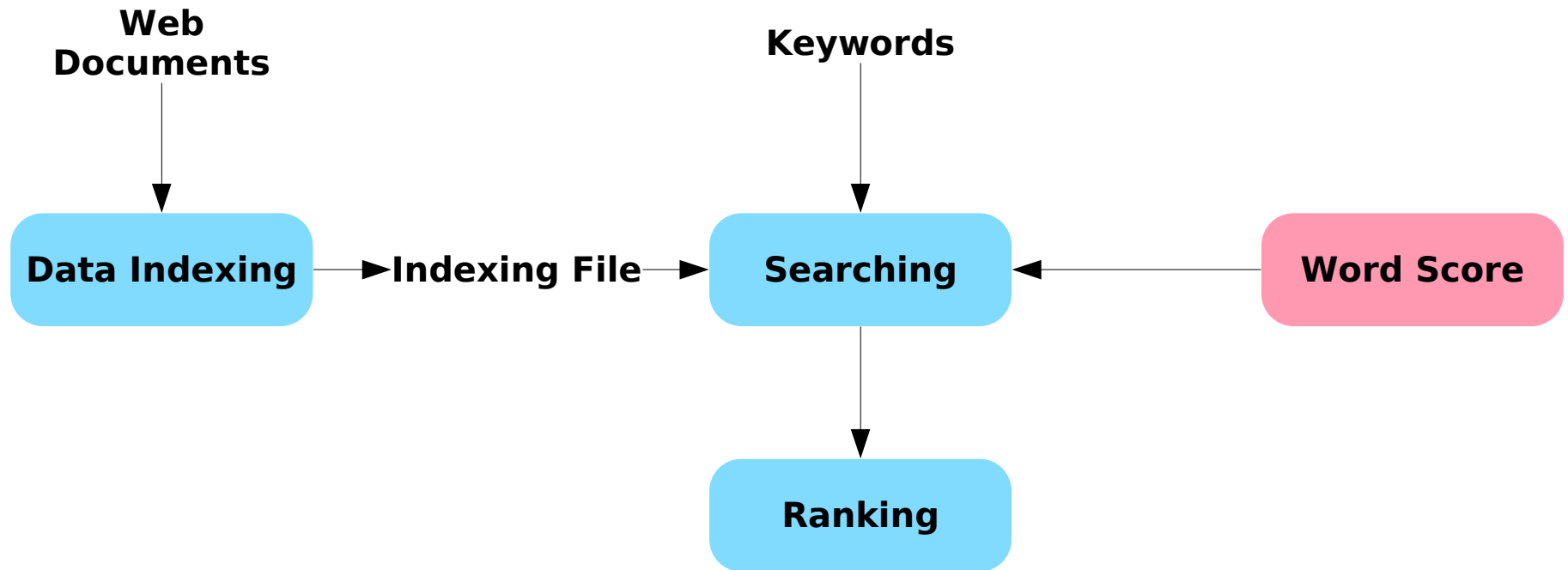


Search Engine ที่ไม่อาศัยพจนานุกรม

- To overcome the limitation of vocabulary for making index
- To deal with the out-of-vocabulary problem
- To avoid the incomplete word segmentation result
- To avoid multiple search in case of phrase search
- To make it extensible for multi-lingual search



Dictionary-less Search Engine --Architecture--



http://www.sansarn.com



all [input type="text" value="ครัว"] ค้นหา

ในกรณีค้นหามากกว่า 1 คำ และ หรือ

ผลลัพธ์จากการค้นหาคำว่า ครัว 1-10 จาก 1539

★★★★★

[E LIB : HOME BUILD](#)

...องอะไร ติดอากาศภายนอกดี ระหว่างห้องนอน ห้องรับแขก ห้องน้ำ และห้อง**ครัว**? HIGH ZONE และ LOW ZONE คืออะไร คอนกรีตผสมไม่มีความแข็งแรงซักเท่าไร กัน? อย่างสู่มสี่ สู่มห้า เปลี่ยนพื้นสำเร็จ กับ พื้นหล่อกับที่... พังแน่นอน เสาะ ซ่อมมีก็อย่าง แล้วจะเลือกใช้อย่างไร? ระวางหล่อระดับเสาคณิตทีเดียว พื้น flat Slab ของท่าน... จะ พังทลาย วางพื้นสำเร็จ หากไม่ย...

http://ite.nectec.or.th/%7Eelib/homes/home_build500.html - 14k

★★★★★

[ห้องสมุด E-LIB : Health Library for Thai : การดูแลสุขภาพ สิวหรือสิ่ว](#)

...มีครั้งหนึ่งที่เรากล่องให้ลูกวางไว้ต่อหน้าต่อตาเรา โดยไม่หยิบไปไว้ที่**ครัว**และเราก็ไม่ได้ดักเตือน เขา ต่อไปเราจะพบว่ามันครั้งที่ 2 และ 3 ตามมาเพราะเด็กจะตามใจตนเองและเลือกสิ่งที่ง่ายที่สุดไว้ก่อน แม้ทำผิดก็จะตามใจตนเอง ไม่เพียงแต่เด็ก แม้ในคนที่เป็นผู้ใหญ่แล้ว หากไม่มีการลงโทษ ดักเตือน ในไม่ช้าก็จะเริ่มทำตามใจตัวเอง และทำสิ่งผิดอยู่เสมอ เพราะความเคยชิน...

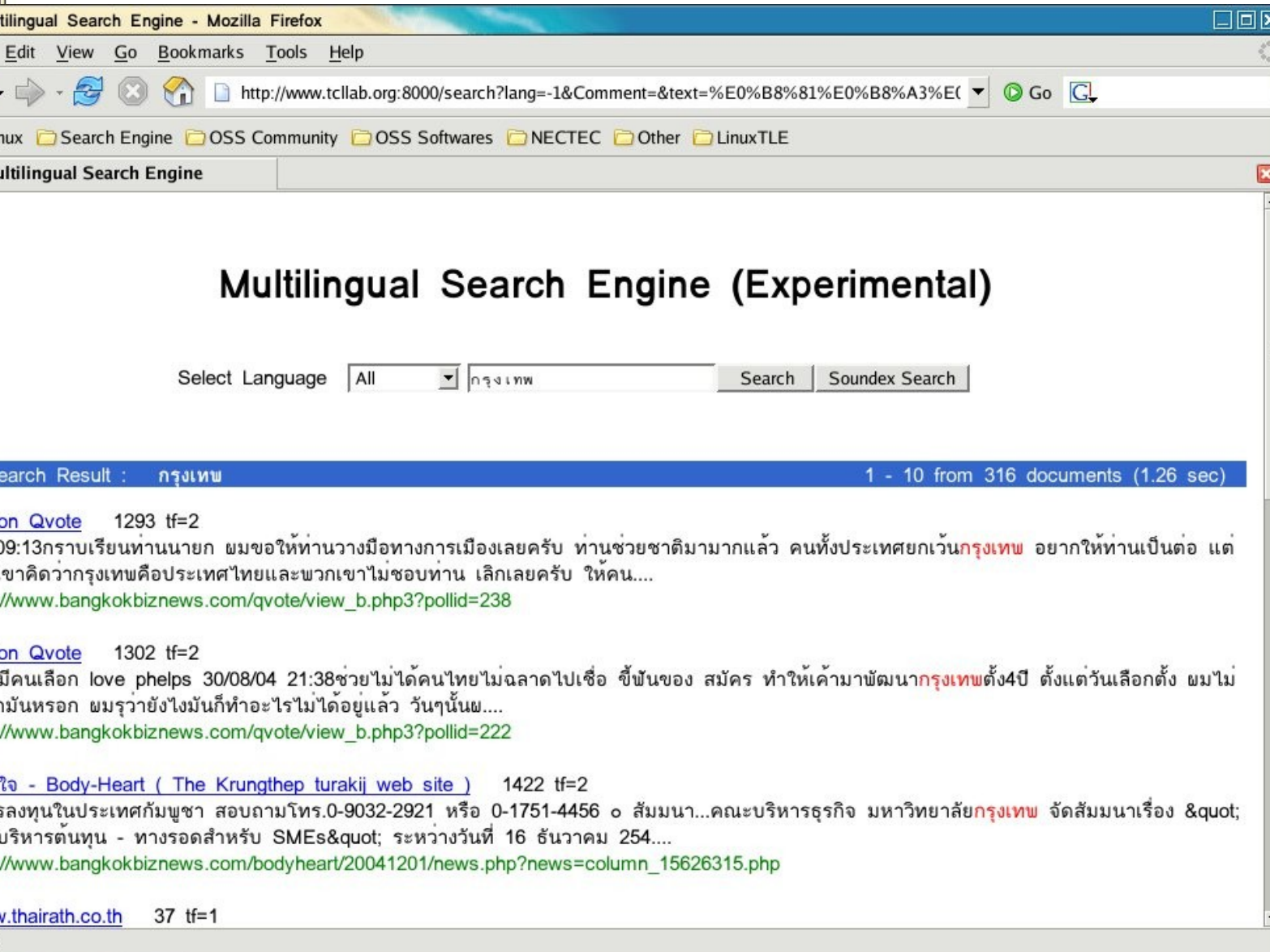
http://ite.nectec.or.th/%7Eelib/doctors/child_hit01.html - 7k

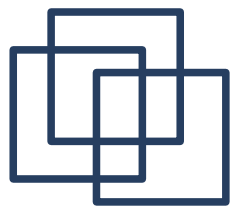
★★★★★

[ห้องสมุด E-LIB : Health Library for Thai : ลีลาแอนด์เทคนิค](#)

...เพื่อความแปลกใหม่... และสร้างความตื่น ตื่น... เปลี่ยนเป็นห้องน้ำ...ห้อง**ครัว**...ห้องรับแขก... ระเบียงบ้านไม่ต้องนะครับ...เดี่ยวข้างบ้านตกใจ... ทำทางที่เหมาะสมสำหรับห้องรับแขกคือ... "ท่าเก้าอี้นายพราน" อีกท่าหนึ่ง ตามสัญญาณครับ...ท่านี้เรียกว่า... "ท่าลิงอุ้มแดง" ท่านี้เหมาะสำหรับบริเวณที่มีพื้นที่จำกัด...เช่นใน Office ขณะพักเที่ยง...หรือในห้องน้ำ... แ...

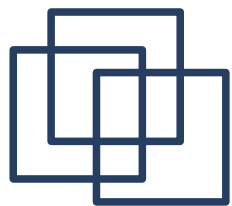
http://ite.nectec.or.th/%7Eelib/doctors/sexed_technic02.html - 37k



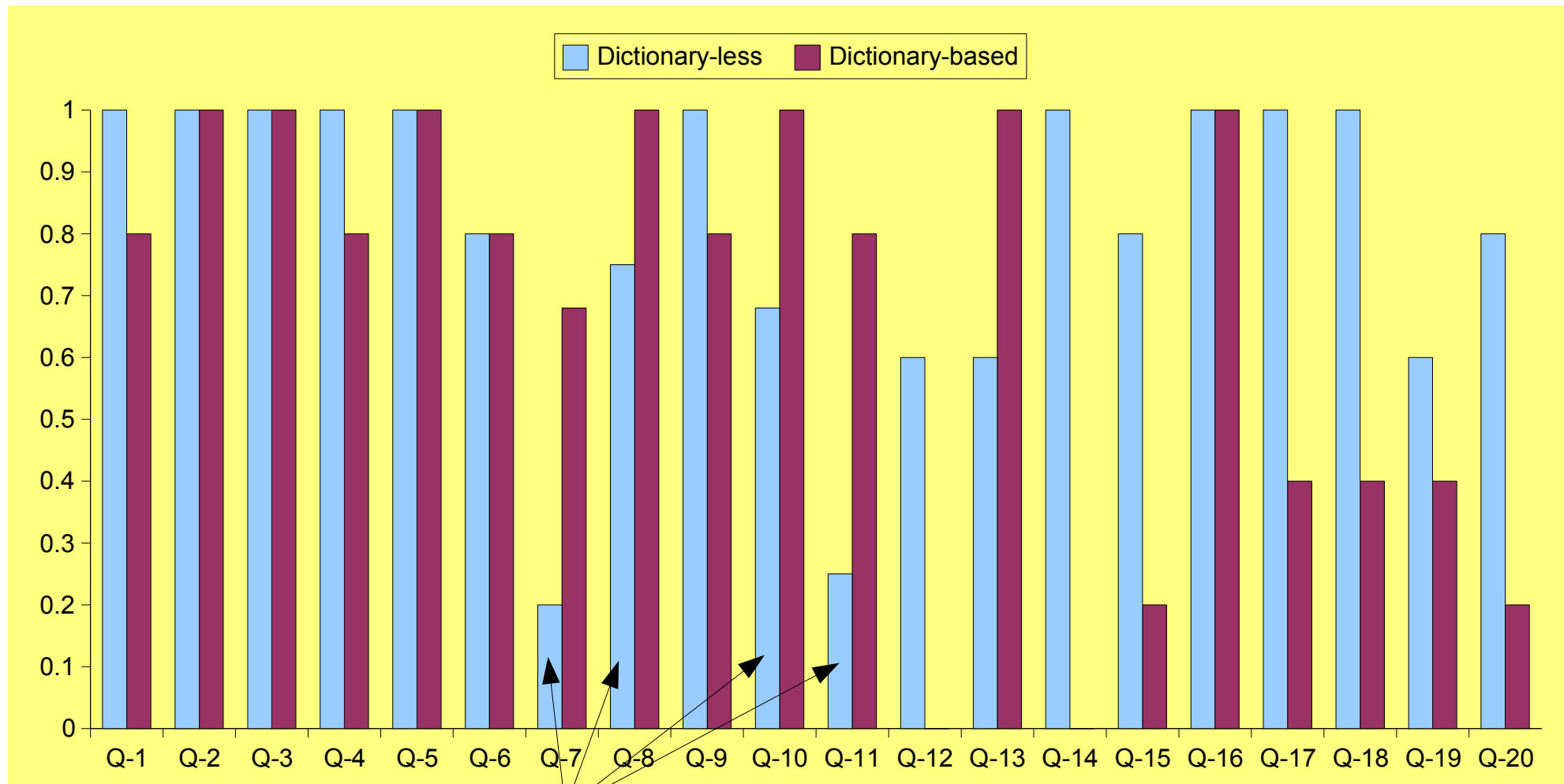


การประเมินผล

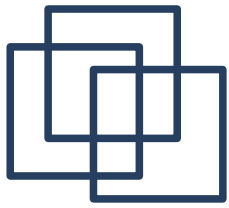
- Evaluate top 10 results of 20 queries by 5 evaluators
- Relevant if 3 out of 5 evaluators agree on each result
- Satisfaction on the result for each query is the average on the relevant



Dictionary-based VS Dictionary-less --Evaluation--



Inferior



List of Queries

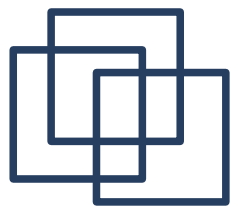
Uncommon and large

| | <i>Unsegmented queries</i> | <i>Segmented queries</i> |
|---|-------------------------------|-----------------------------------|
| 1 | บริจาค, สีนามิ | บริจาค, สีนามิ |
| 2 | เส้นตาย, ชัดคม | เส้นตาย, ชัดคม |
| 3 | มันส์, โชว์ตัว, จีน | มันส์, โชว์ ตัว, จีน |
| 4 | ผลกระทบ, ราคาน้ำมันแพง | ผล กระทบ, ราคา น้ำมัน แพง |
| 5 | ใช้หัวदनก | ใช้หัว दनก |
| 6 | ทุจริต, การเลือกตั้ง | ทุจริต, การ เลือกตั้ง |
| 7 | จับกุม, ผู้ก่อการร้าย, ภาคใต้ | จับ กุม, ผู้ ก่อการ ร้าย, ภาค ใต้ |
| 8 | นโยบาย, แก้ไข, ปัญหา ยาเสพติด | นโยบาย, แก้ ไข, ปัญหา ยา เสพติด |
| 9 | ทดลอง, ลดค่า ทางด่วน | ทดลอง, ลด ค่า ทาง ด่วน |

| | <i>Unsegmented queries</i> | <i>Segmented queries</i> |
|----|------------------------------------|-------------------------------------|
| 11 | ลงทุน, ในพม่า | ลงทุน, ใน พม่า |
| 12 | ส่งเสริม, การท่องเที่ยว, ไทย | ส่งเสริม, การ ท่องเที่ยว, ไทย |
| 13 | เลือกตั้ง, ประธานาธิบดี, ปาเลสไตน์ | เลือกตั้ง, ประธานาธิบดี, ปา เลสไตน์ |
| 14 | เลือกตั้ง, ผู้ว่า, กทม. | เลือกตั้ง, ผู้ ว่า, กทม . |
| 15 | สินค้าไทย, การส่งออก | สินค้า ไทย, การ ส่ง ออก |
| 16 | แปรรูป รัฐวิสาหกิจ | แปรรูป รัฐ วิสาหกิจ |
| 17 | อุ้ม, นายสมชาย | อุ้ม, นาย สม ชาย |
| 18 | ฉลองปีใหม่ | ฉลอง ปี ใหม่ |
| 19 | พรทิพย์, ลาออก | พร ทิพย์, ลา ออก |
| 20 | สวนสนุก | สวน สนุก |

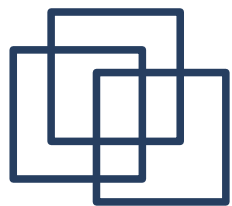
Common and small

=> Dictionary-based approach is good at searching uncommon words but not the common words.



การประเมินผล

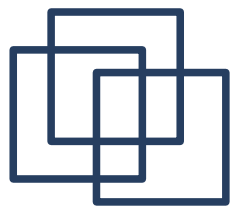
- Dictionary-based gets inferior results in case of excessive segmentation to be common and small words i.e. การส่งออก (การ|ส่ง|ออก), but not in the case of being uncommon and large words i.e. ผลกระทบ (ผลก|ระ|ทบ)



Web language engineering for Open collaborative archiving (WLE-OCA)

Motivation:

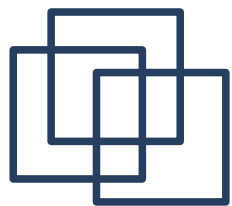
- Establishing Thailand-Japan collaborative archiving on broadband Internet (NECTEC-MIC 45M bps high speed link)
- Collaboration between TCL-Thailand and NUT-Japan to support Language Observatory Project and Asian Language Resource Network
- Real-time monitoring population of web document



Web language engineering for Open collaborative archiving (WLE-OCA)

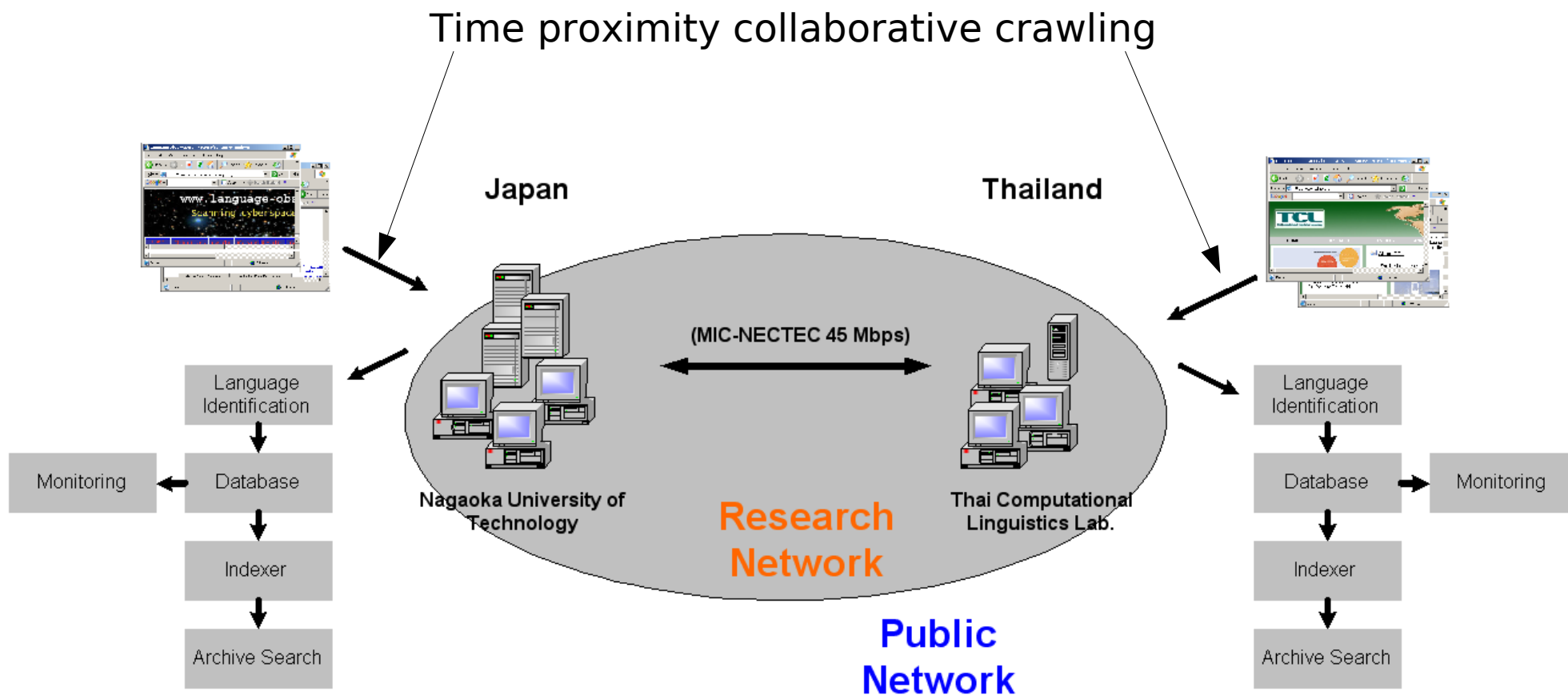
Goal:

- Archiving web document at 1m pages/day/site
- Evaluating Multi-lingual and Cross-lingual search engine



Web language engineering for Open collaborative archiving (WLE-OCA)

Architecture:





Collaborative Crawler

Log out

Server List



Language
Encoding
Charset



Language
Encoding
Charset

Host : TH-ORCHID @ TCL
Start crawl date : 2005-11-18



Language
Encoding
Charset

Host : JP-KIKU @ NUT
Start crawl date : 2005-11-21

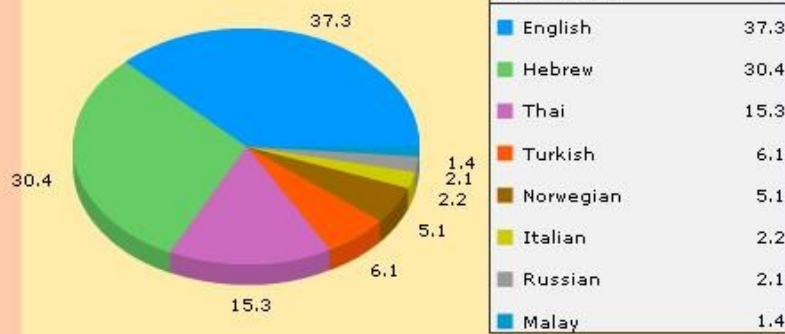


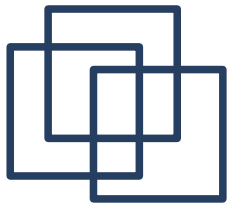
Language
Encoding
Charset

Crawler Status

Language Distribution

of 347,875 pages from : all sites





เสมือนกระจกเงาสำหรับสำรวจความรู้ในสังคม
ให้รู้ถึง แหล่งที่มา ภาษา และ เนื้อหา

WLE-OCA

www.tcilab.org/weblang