


Towards the Cross- Language Technology



Virach Sornlertlamvanich
Information Research and Development Division

NECTEC, Thailand

virach@nectec.or.th

APAN 2003 Conference in Fukuoka, 21-24 Jan 2003

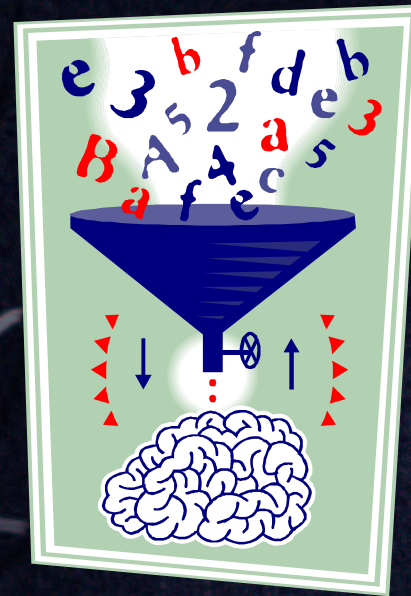
- ParSit */parse-it/*
 - ภาษิต (proverb)
 - English-Thai cross-language web navigator.
 - www.suparsit.com



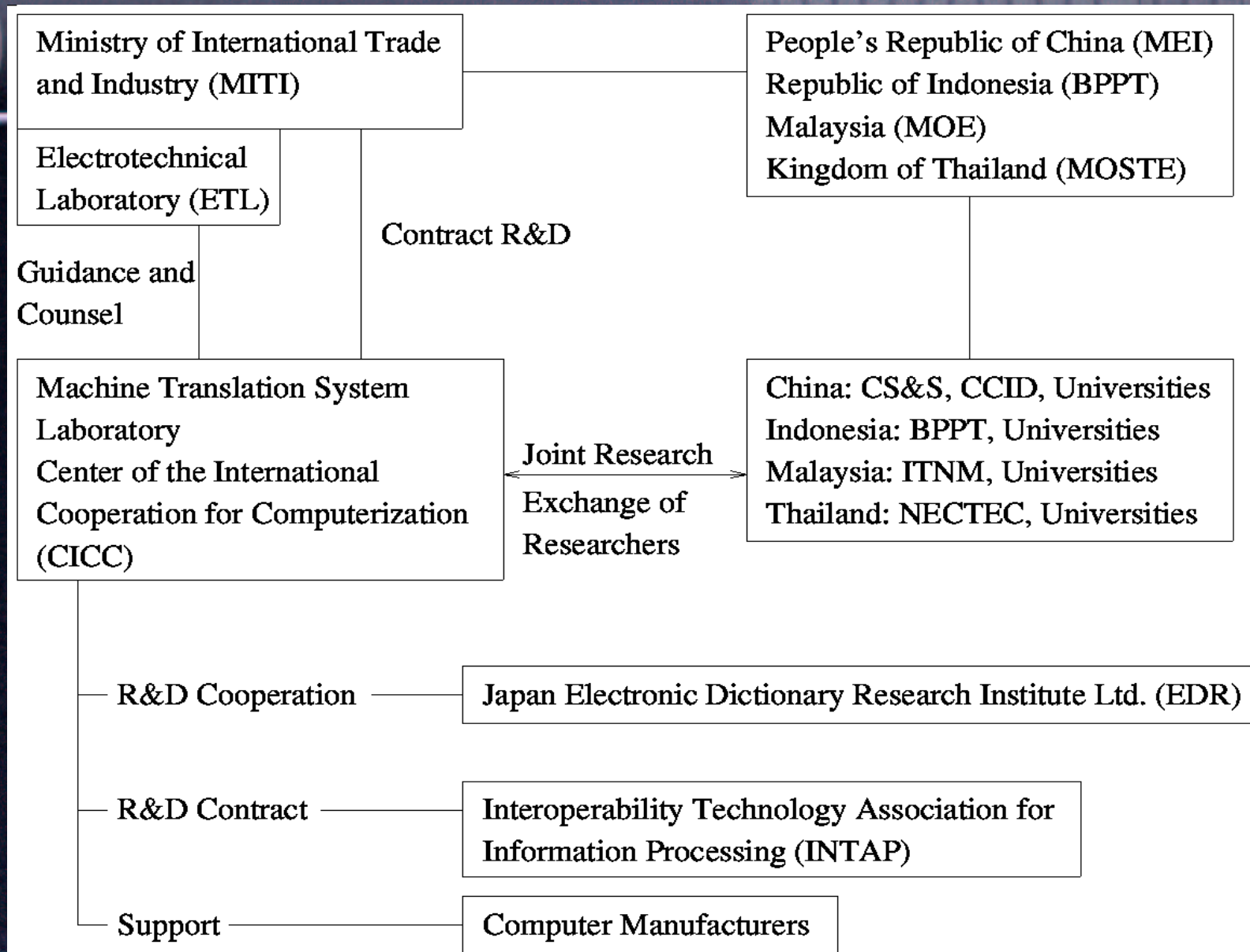
2 Big Machine Translation Projects

- ARIANE Project (1981-)
- Multilingual Machine Translation Project (1987-94) C, I, M, J, T

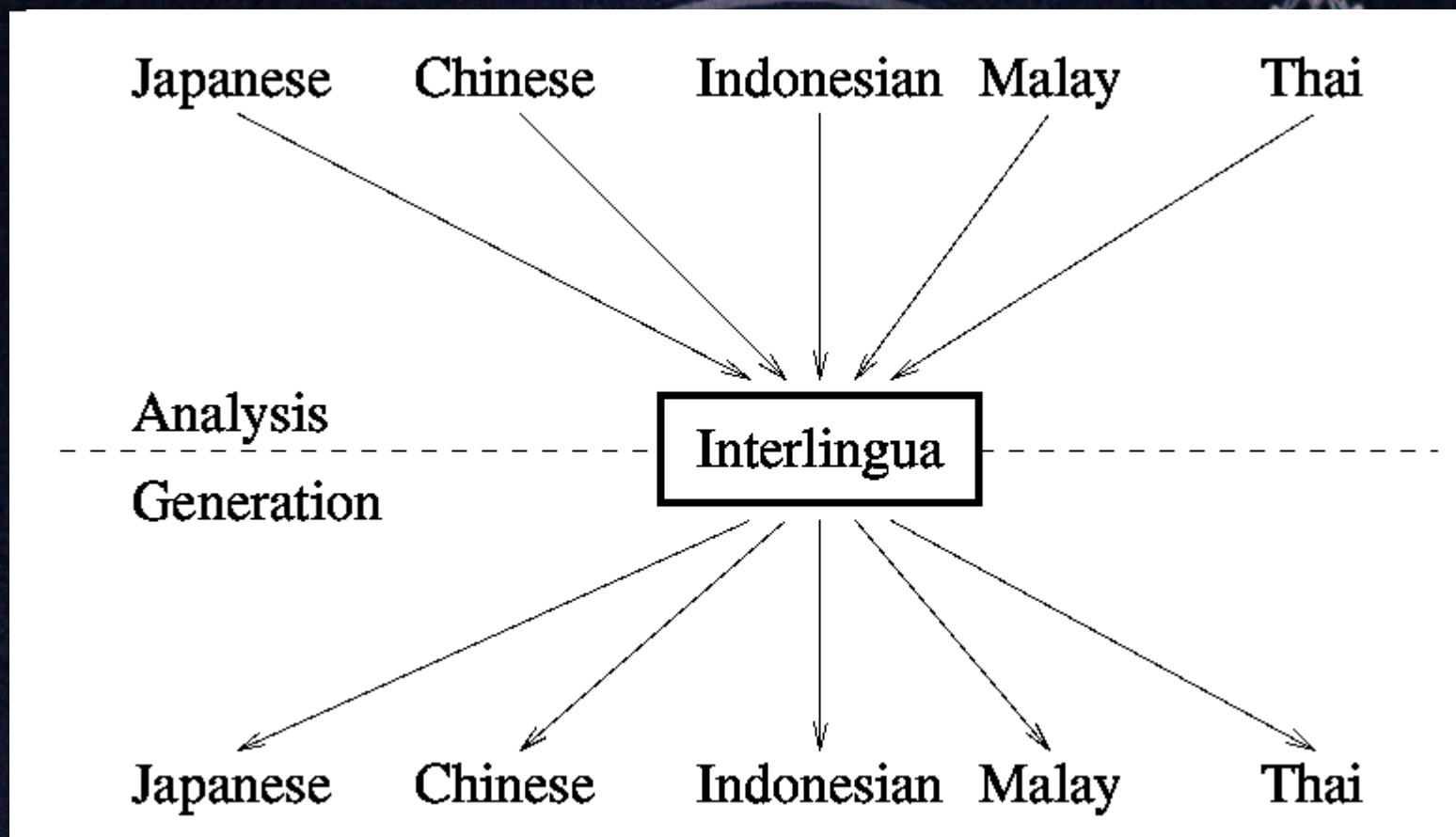
Learn from
experienced
countries



Multilingual Machine Translation Project



Interlingual Approach



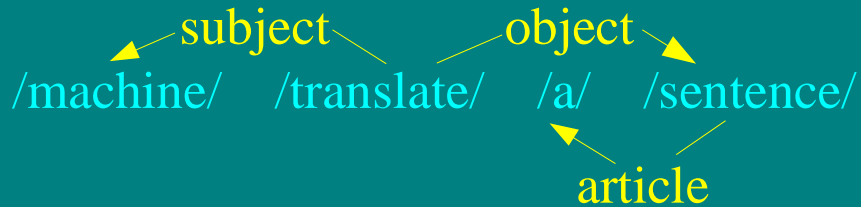
Interlingual Approach

Machine translates a sentence.

Morphological analysis

/machine/ /translate/ /a/ /sentence/

syntactic analysis



เครื่องแปลประโยค

morphological generation

/เครื่อง/ /แปล/ /ประโยค/

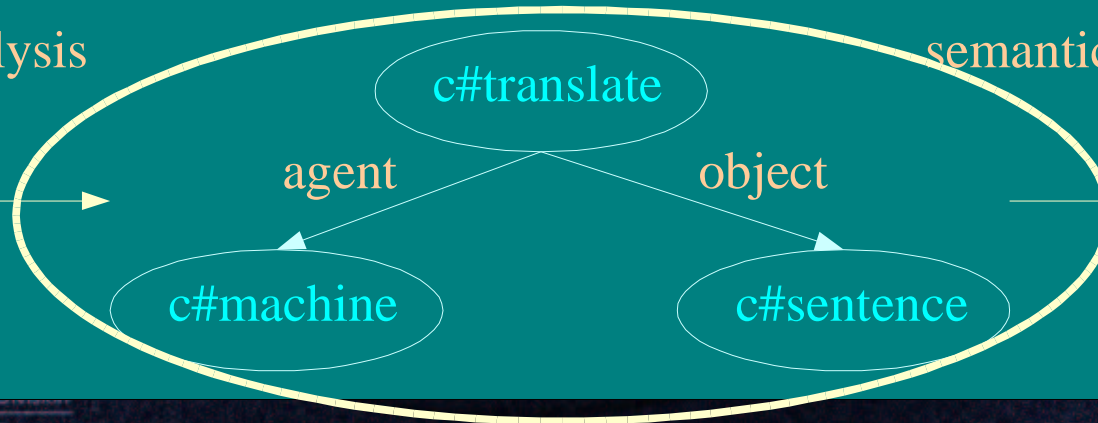
syntactic generation



Interlingua

semantic analysis

semantic generation



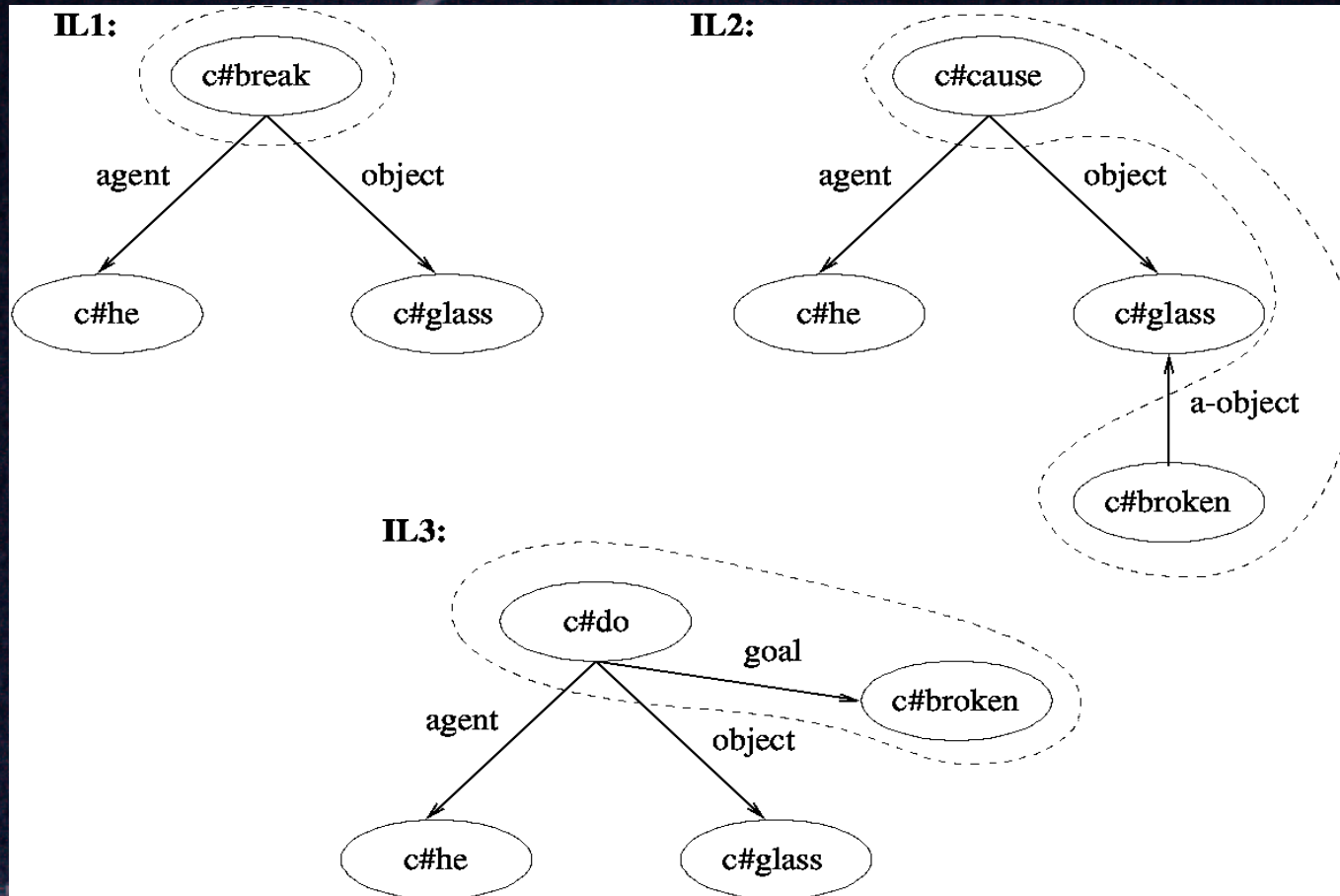
Difficulties in MT

- Grammatical issues
 - ➔ Inflection
 - ➔ Article
 - ➔ Grammatical case marker
 - ➔ Word and sentence boundary, Punctuation
- Semantic issues
 - ➔ Definitions of concept, set of semantic case, ontology
 - ➔ Concept decomposition, concept divergency, concept granularity



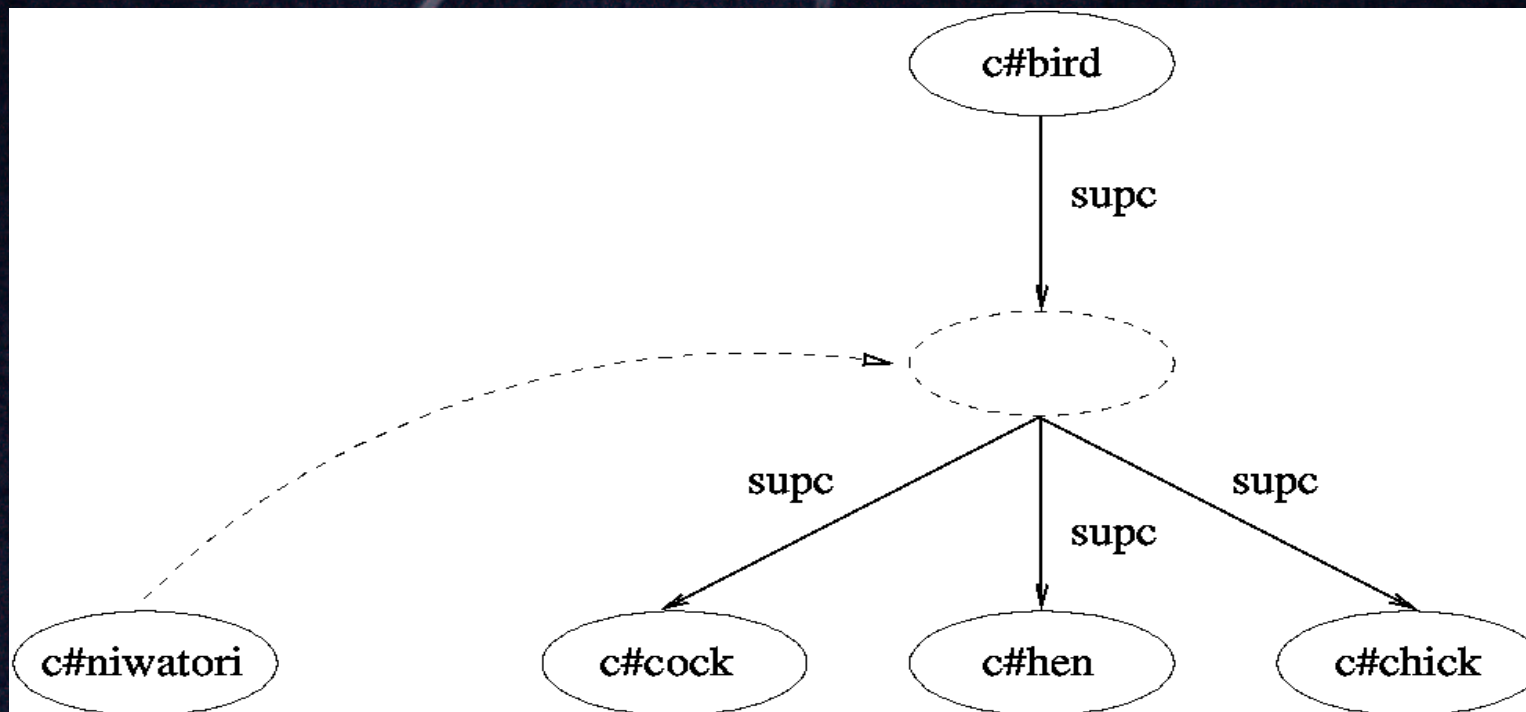
Concept decomposition

He breaks the glass. [He causes the glass broken.]



Concept divergency

A concept of Language A is corresponding to a meta-concept in Language B.



Concept granularity

A concept of Language A is corresponding to a multiple concepts in Language B.

Japanese		English	
sakura	(tree)	cherry tree	
	(flower)	cherry blossom	
sakuranbo	(fruit)	cherry	(fruit)
			(tree)



Concept hierarchy



Ontologies

- **EDR**
 - ➔ Approach: Word description as employed in dictionaries
 - ➔ Problem: Ambiguities and incomputability
- **Wordnet**
 - ➔ Approach: Synonym set and simple semantic relations to other words
 - ➔ Problem: Ambiguities
- **UW**
 - ➔ Approach: Headwords and semantic restrictions
 - ➔ Advantage: Computability and no ambiguity



Ontologies

Representation of concept 'tired' in different schemes

EDR	Wordnet 1.5	UW
<ul style="list-style-type: none">- having or displaying a need for rest- having lost of interest- lack of imagination	<ul style="list-style-type: none">- A1: tired (vs. rested)- A2: bromidic, commonplace, hackneyed, ...- V1: tire, pall, grow weary, fatigue- V2: tire, wear upon, fag out- V3: run down, exhaust, sap, ...- V4: bore, tire, ...	<ul style="list-style-type: none">- tired- tired(icl>physical)- tired(icl>mental)



Design of Ontology

- Computational concept
- Concept insertion/deletion;
composition/decomposition
- Expressive concept



www.suparsit.com

English-Thai Translation - Parsit - Mozilla [Build ID: 2002121223]

File Edit View Search Go Bookmarks Tasks Help

Back Forward Reload Stop <http://suparsit.com/> Search Print

Home Bookmarks Red Hat Network Support Shop Products Training

Parsit

ภาษิต

SWATH LEXITRON SANGKARN ARNTHAI LinuxTLE OfficeTLE NECTEC SERVICES

ท่องเว็บทั่วโลกด้วยภาษาไทยผ่าน "ภาษิต"
บริการแปลภาษาอังกฤษเป็นไทยด้วยคอมพิวเตอร์ทางอินเทอร์เน็ต

ข้อมูลข่าวสารภาษาอังกฤษบนอินเทอร์เน็ตจะถูกส่งมายังระบบที่อาศัยฐานความรู้ทางไวยากรณ์และความหมาย "ภาษิต" จะแปลเว็บเพจเป็นภาษาไทยและส่งผลไปยังผู้ใช้ด้วยโครงสร้างเดิมของต้นฉบับ

พัฒนาโดยผ่านกลุ่มวิจัยและพัฒนาศาสตร์สารสนเทศ ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ

ได้รับความร่วมมือจากบริษัท เอ็นอีซี ประเทศญี่ปุ่น

Surf the webs worldwide with Thai via "Parsit: the English-to-Thai computer-based translation service on the Internet".

English information on the Internet is transmitted to the system operated on syntactic and semantic knowledge base. "Parsit" translates English web pages into Thai web pages and then sends the result in the original layout back to users.

Developed by Information Research and Development Division, National Electronics and Computer Technology Center.

In Cooperation with NEC Corporation (Japan).

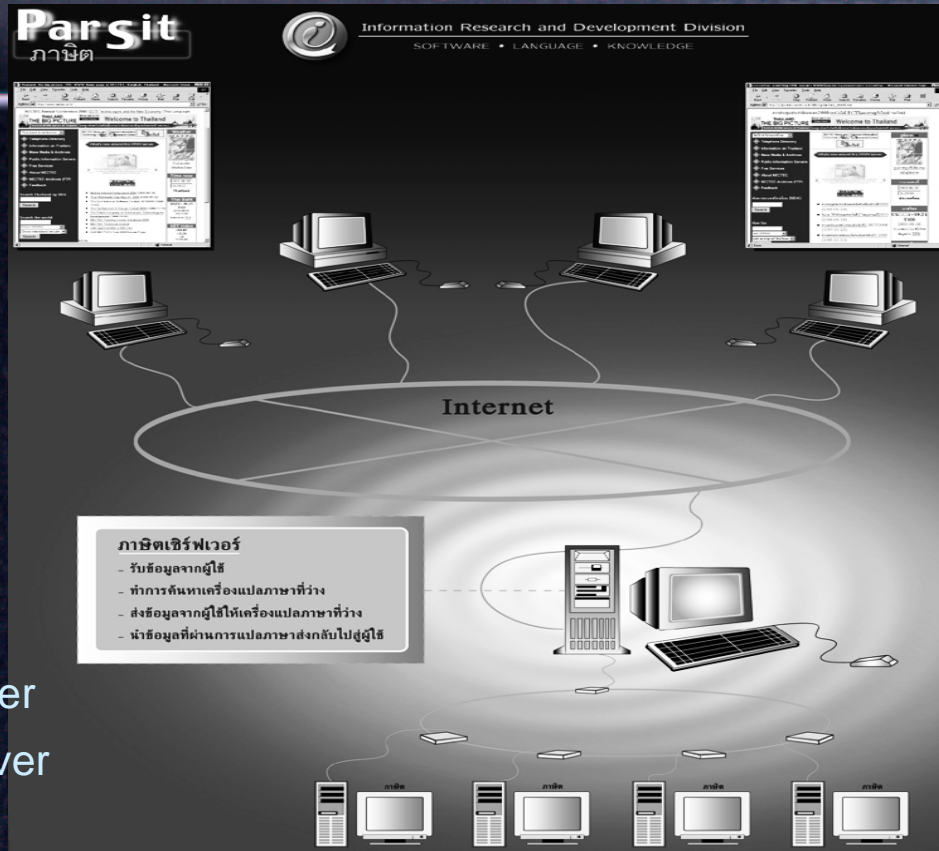
Since June 24, 2000

Document: Done (10.131 secs)

ParSit Service Process

1
User send URL/TEXT
to ParSit server

2
ParSit Server
- receive data from user
- send data to MT server

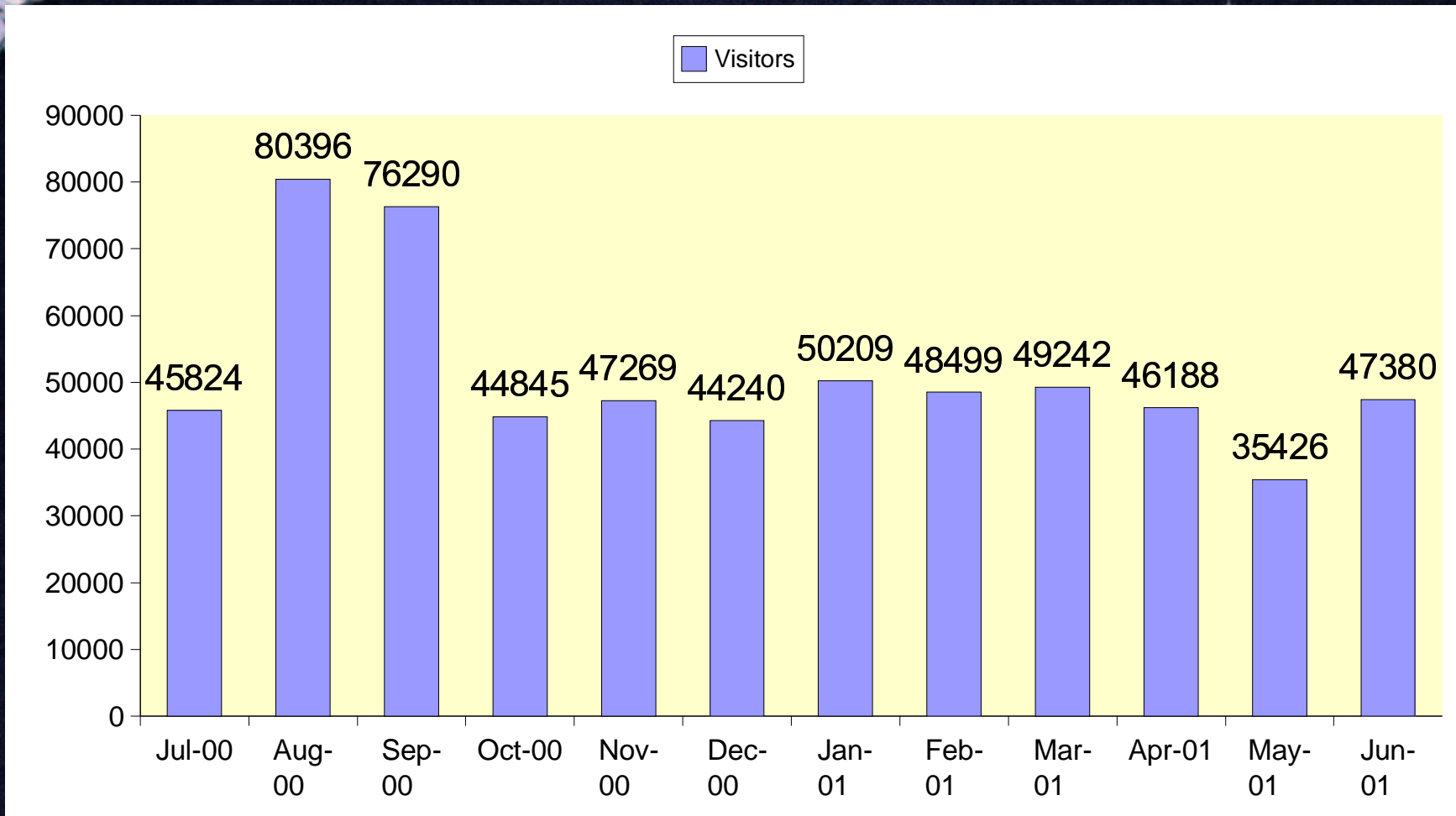


3
MT Server (E-T)

5
User receive Thai
WP from server

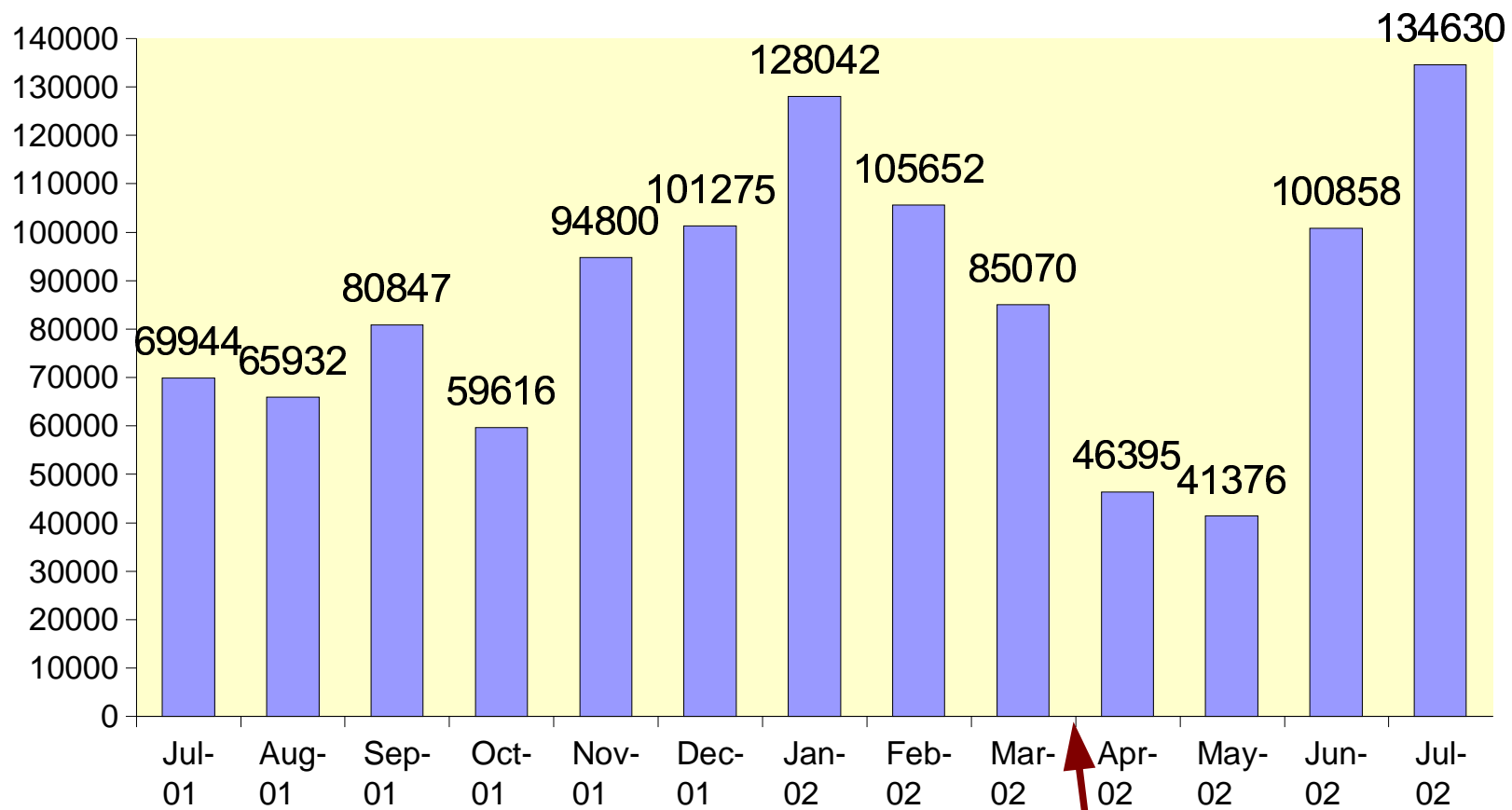
4
ParSit Server
- receive result from MT
- send result to User

Number of visitor/month



Number of unique session/month

Translation

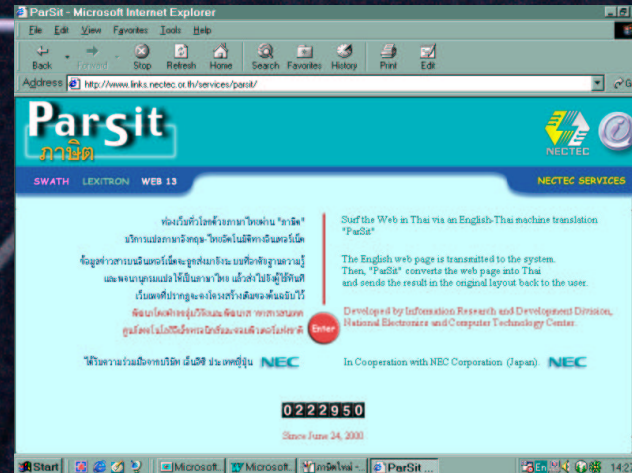


Move to Science Park

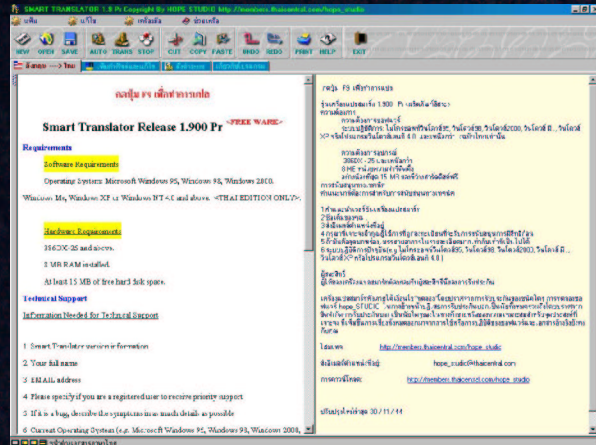


Available MT in Thailand

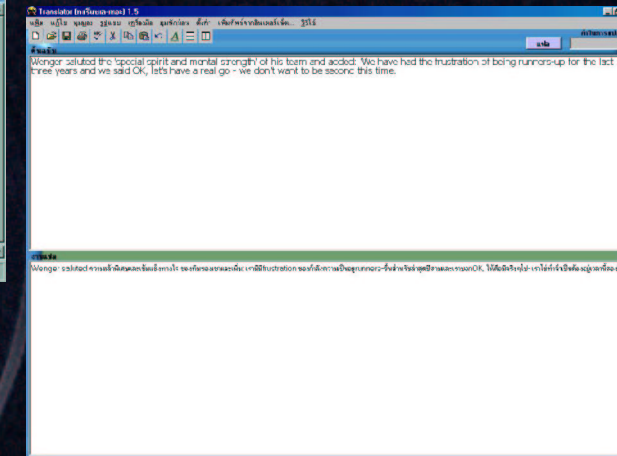
ParSit



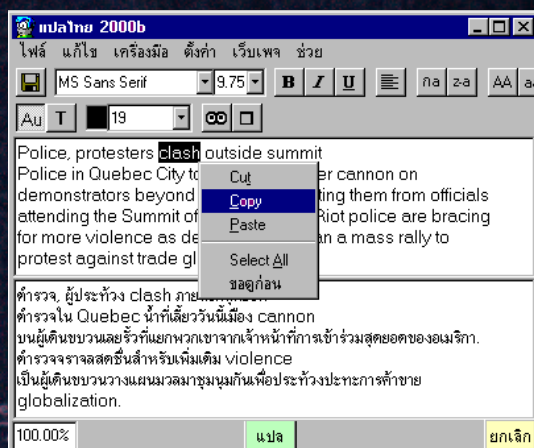
SMART TRANSLATOR 1.9



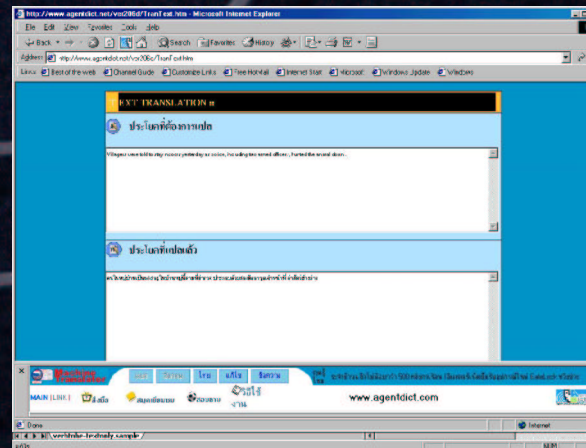
TRANSLATOR 1.5



PLAETHAI



AGENT DICT



Available MT in Thailand

- **ParSit**
www.suparsit.com
- **Smart Translator 1.9**
- **Translator 1.5**
www.the-kstudio.com
- **PlaeThai 2000b**
www.palthai.com
- **AgentDict Translation v.2.06**
www.dinosoft.co.th/software/AgentDict/AgentDict.htm



MT Systems in Thailand

	Parsit	Smart Translator	PlaeThai	Translator 1.5	AgentDict
Language	Eng-Thai	Eng-Thai	Eng-Thai	Eng-Thai	Eng-Thai
Lexicon	80,000	33,400	11,000	8,000	> 200,000
Type of document	-Web page -Text	-File (Text) -Text	-Web page -File (Text) -Text	-File (Text) -Text	-Web page -Text



Evaluation Corpus

- Sentence Corpus for Evaluation
 - 770 Sentences
 - Designed by Japan Electronic Industry Development Association (JEIDA)
 - Have the characteristics for testing the following linguistic problems.
 - Word Level
 - Concept mismatching
 - Word absence
 - etc.
 - Sentence Level
 - Grammar
 - Modifier misplacement
 - etc.



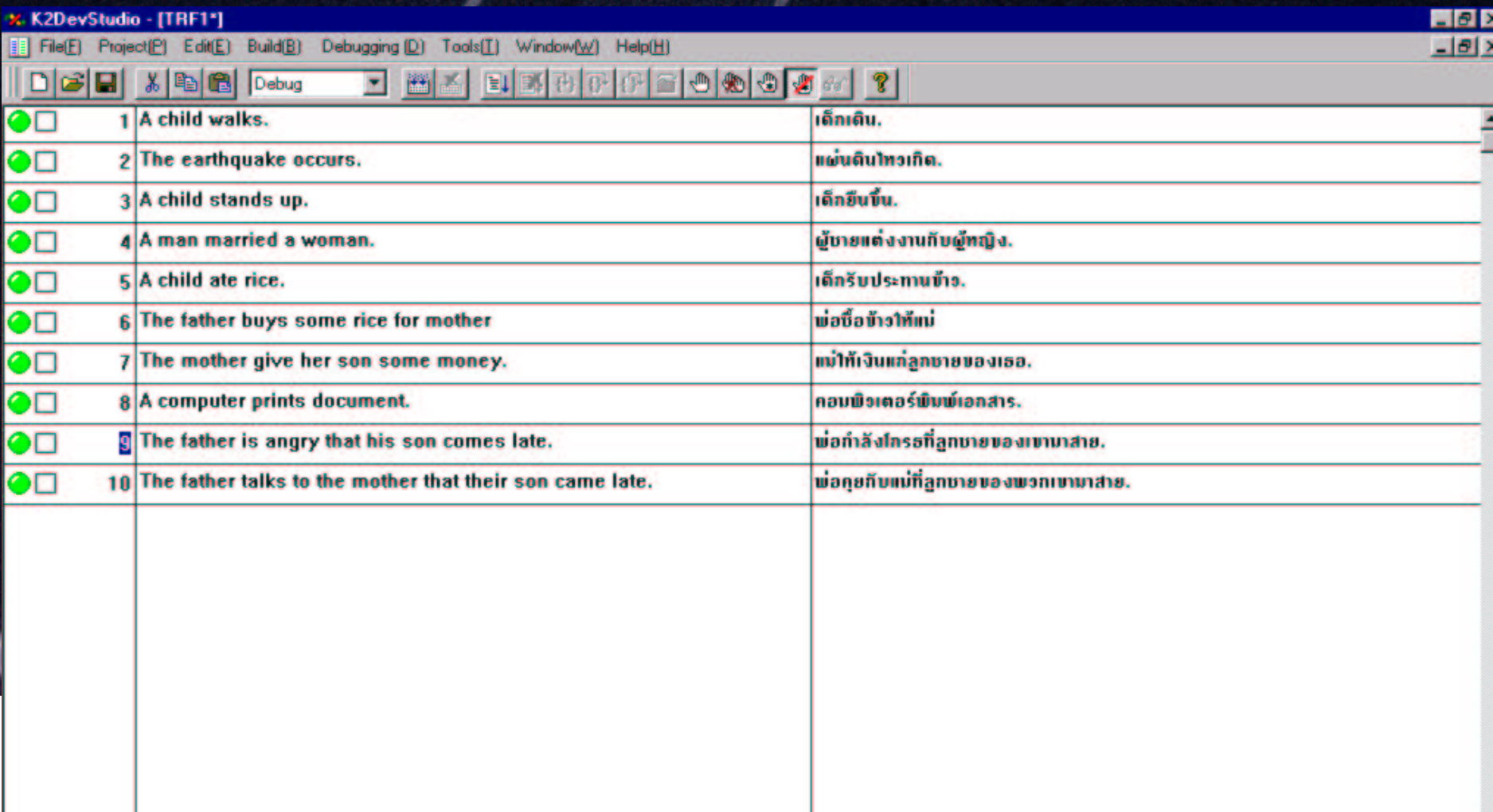
Evaluation Result (number of sentence)

	Parsit	Smart Translator	AgentDict	PlaeThai	Translator 1.5
Perfect Translation Sentence	359	200	180	151	128
Partial Translation Sentence	351	504	486	465	567
False Translation Sentence	60	66	104	154	85
Total number of sentence	770	770	770	770	770
Accuracy	47-92%	26-91%	23-86%	20-80%	17-90%



ParSit 2 (Thai-English)

- Translation from Thai to English
- 13 Verb Patterns; 4,000 Rules; more than 30,000 words



The screenshot shows the K2DevStudio interface with a table of Thai-English translations. The table has two columns: English sentences and Thai translations. Each row has a green circle and a checkbox on the left. The table is as follows:

English Sentence	Thai Translation
1 A child walks.	เด็กเดิน.
2 The earthquake occurs.	แผ่นดินไหวเกิด.
3 A child stands up.	เด็กยืนขึ้น.
4 A man married a woman.	ผู้ชายแต่งงานกับผู้หญิง.
5 A child ate rice.	เด็กรับประทานข้าว.
6 The father buys some rice for mother	พ่อซื้อข้าวให้แม่
7 The mother give her son some money.	แม่ให้เงินแก่ลูกชายของเธอ.
8 A computer prints document.	คอมพิวเตอร์พิมพ์เอกสาร.
9 The father is angry that his son comes late.	พ่อกำลังโกรธที่ลูกชายของเขามาสาย.
10 The father talks to the mother that their son came late.	พ่อคุยกับแม่ที่ลูกชายของพวกเขามาสาย.

- Sansarn */san-sarn/*
 - สรรสาร (information search)
 - Language-independent probabilistic full-text search.
 - www.sansarn.com



Meaningful Bits

ADLTSUG**KNOWLEDGE**BWGWZKTILA

ปรัจตเสีศฐาดี**ความรู้**ษะฎุกฮเศฉ



Mutual Information

$$Lm(xyz) = \frac{P(xyz)}{P(x)P(yz)}$$

$$Rm(xyz) = \frac{P(xyz)}{P(xy)P(z)}$$

x y z

x 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.

High mutual information implies that xyz co-occurs more than expected by chance. If xyz is a word then its Lm and Rm must be high.

...Efunction... vs ...Function...



Entropy

$$LEnt(y) = - \sum P(xy/y) \cdot \log P(xy/y)$$

x

y

$$REnt(y) = - \sum P(yz/y) \cdot \log P(yz/y)$$

y

z

where A is the set of characters
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 then its left and right entropy must be high.

...?function... vs ...?unction...





ภาพยนต์ ค้นหา ค้นหาทั่วโลกด้วย

Language-independent full-text search

ผลลัพธ์จากการค้นหาคำว่า ภาพยนต์ 1 - 10 จาก 374 เอกสาร

[การไฟฟ้าฝ่ายผลิตแห่งประเทศไทย](#) ★★★★★

....ข้อมูลเกี่ยวกับการไฟฟ้าฝ่ายผลิตแห่งประเทศไทย....

<http://www.egat.or.th>

[สพร.กฟผ.](#) ★★★★★

....สพร.กฟผ. This page uses frames, but your browser doesn't support t....

<http://www.seea.egat.or.th/main.html>

[สพร.กฟผ.](#) ★★★★★

....สพร.กฟผ. This page uses frames, but your browser doesn't support t....

<http://www.seea.egat.or.th/main.html>

[คลิกเพื่อดูข้อมูลเพิ่มเติมจากเว็บไซต์นี้.....]



ค้นในสถาบัน

ทุกสถาบัน

ผลงานตาม

ความสัมพันธ์กับคำ

ผู้เชี่ยวชาญ

ซอฟต์แวร์

ค้นหา มากกว่า 1 คำ เลือก และ หรือ

เว้นวรรคหว่านคำ ในกรณีที่ต้องการค้นหา

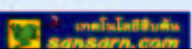
มากกว่า

ฐานข้อมูลนักวิจัย หน้าแรก

ค้นหา ทั่วไป | หัวข้อ

ค้นหา

ค้นหาใกล้เคียง



National Research Database

1 - 10 จาก 35 โครงการ

- [เนคเทค] [โปรแกรมแบบจำลองการไหลของน้ำในลำน้ำสายน้ำสาขาของเขื่อนศรีนครินทร์](#) ***** จำนวนผู้เข้าชม : 1
2546 นาย สแตนนิสลาส มาคานอฟ บทคัดย่อ [text]
- [เนคเทค] [ซอฟต์แวร์จำลองตลาดและฮาร์ดแวร์สำหรับคนพิการแทนเพื่อใช้งานคอมพิวเตอร์](#) ***** จำนวนผู้เข้าชม : 0
2546 นาย มานะ ศรียุทธศักดิ์ บทคัดย่อ [text]
- [เนคเทค] [โครงการพัฒนาระบบพูดแทนพิมพ์ภาษาไทย](#) ***** จำนวนผู้เข้าชม : 0
2545 นาย มนตรี กาญจนะเดชะ บทคัดย่อ [text]
- [เอ็มเทค] [การพัฒนาโปรแกรมเพื่อวิเคราะห์ลักษณะสมบัติของวัสดุที่มีพฤติกรรมเป็นแฟรคทัล](#) ***** จำนวนผู้เข้าชม : 0
**** ชวิชัย ชรินพานิชกุล, วิวัฒน์ ตัณฑะพานิชกุล, เกษม สัตยาวุฒิพงษ์ บทคัดย่อ [text]
- [สกว] [ระบบตรวจสอบคุณภาพของผลไม้อัตโนมัติโดยการมองเห็นของคอมพิวเตอร์](#) ***** จำนวนผู้เข้าชม : 2
2541 ดร.ธนาชาติ น่มนนท์ บทคัดย่อ [text]

ฐานข้อมูลนักวิจัย

- หน้าแรก
- รายชื่อทั้งหมด
- รายชื่อตามสาขา

ล็อกอินเข้าระบบ

ชื่อ Login

รหัสผ่าน

[ลืมรหัสผ่าน](#) | [ลงทะเบียนใหม่](#)

สาขางานวิจัย

- Agricultural Science
 - Agricultural Product
 - Agroindustrial Com
 - Agronomy
 - Animal Breeding I
 - Animal Feed
 - Animal Husbandry
 - Animal Nutrition
 - Apiculture
 - Aquaculture

ค้นหา:

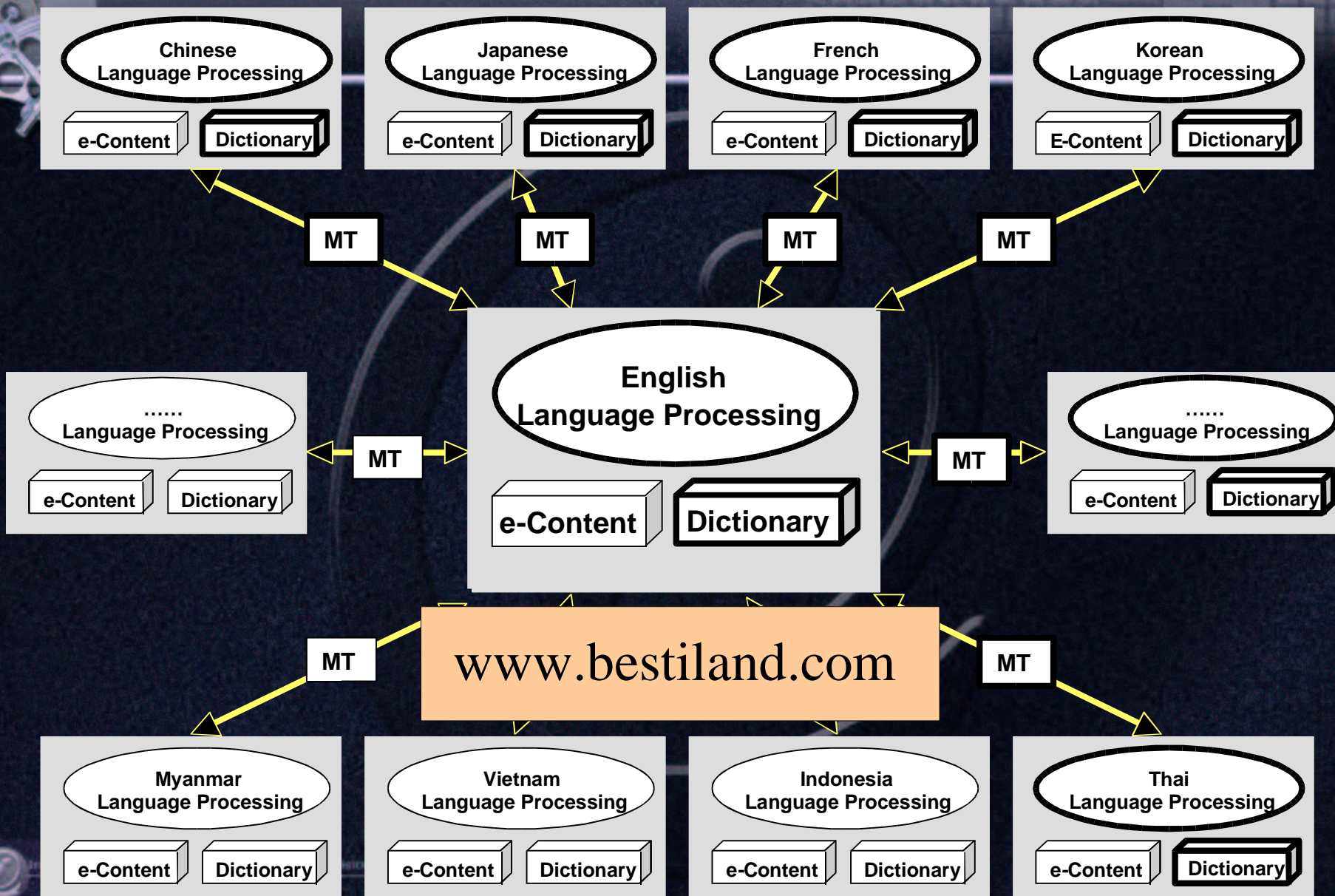
[ค้นหา](#) [ทั่วไป](#) | [หัวข้อ](#)

Researcher Database

รายชื่อนักวิจัยสาขา Agricultural Science click ที่ชื่อเพื่อดูรายละเอียด

ชื่อไทย	ชื่ออังกฤษ	สาขาวิจัยที่เกี่ยวข้อง
กัญญารัตน์ สุโพธิ์วัฒน		PLANT BIOTECHNOLOGY/ HORTICULTURE/ AGRICULTURE
กิจการ สุขมาตย์		P IMMUNITY INFECTIOUS MALS
จิระชัย กาญจนพฤทธิพงศ์	Kanjansaprompong	NUTRITION / ANIMAL NUTRITION / RUMEN ECOSYSTEM / RUMEN EFFICIENCY
ชัยฤทธิ์ มณีพงษ์		TISSUE CULTURE / AGRONOMY / PLANT BREEDING / PLANT GENETICS / GENETIC ENGINEERING / PROTOPLASTS
ประสาทพร สมิตะมาน	Prasartporn Smitamana	BOTANY / PLANT GENETICS TISSUE CULTURE / GENETIC ENGINEERING PROTOPLAST SHORTICULTURE TEMPERATE CUT FLOWER PLANT PROPAGATION PLANT IMPROVEMENT
ปรานอม พุดมพงษ์		HORTICULTURE / TISSUE CULTURE
เปี่ยมศักดิ์ เมนะเศวต	Piamsak Menasveta	MARINE / AQUATIC / AQUACULTURE / AQUACULTURE BIOTECHMISTOGY / MARINE / OLLUTION / WATER RECIRCULATING SYSTEM / SHRIMP
พันทิพา พงษ์เพียรจันทร์	Puntipa Pongpiachan	ANIMAL NUTRITION
ไพเราะ กิพย์ทัศน์		NITROGEN FIXATION BIOCHEMISTRY / VITAMIN METABOLISM STEROID HORMONES CANCER REPRODUCTIVE SYSTEM PICENITROGEN FIXATION
ยุวดี มานะเกษม	Yuvadee Manakasem	HORTICULTURE / TISSUE CULTURE

Cross Language Technology

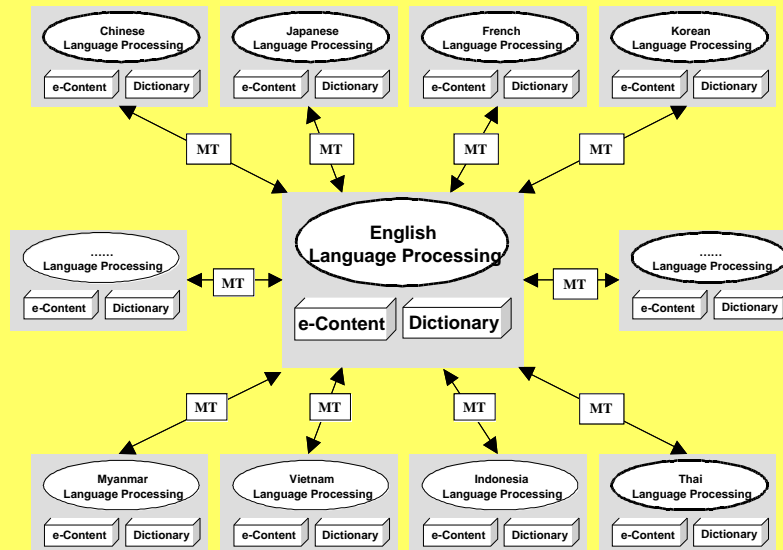


Application over the Cross Language Tech.

E-services

Presentation | Extraction | Retrieval | Summarization | MT | Mining | Visualization

Cross Language Technology



Cross Language Information Retrieval

