中华人民共和国国家知识产权局

State Intellectual Property Office
of the People's Republic of China

# SIPO's Efforts on Improving Quality of Chinese-English Patent Machine Translation Service

Dan WANG

SIPO

August 30, 2009

www.sipo.gov.cn

# Outline

www.sipo.gov.cn

# Chinese-English Patent MT Roadmap

2005.6     Chinese-English MT project started in SIPO

2007.4     Single-user version prototype **CPMT** (China Patent Machine Translation)

2008.4     CPMT launched for patent gisting

2008.4 -  Test and improvement for higher quality

# Current Status of CPMT

◆ Integrated use with patent search services of SIPO and CPIC

➲ available on http://www.sipo.gov.cn/sipo_English

➲ available on http://www.cnpat.com.cn

➲ English search of bibliographic data and abstracts, and on-the-fly MT for full text

# MT Entrance: Abstracts by Manual Translation



The present invention relates to a ADSL router, formed from wide area network interface, local area network interface and signal processing portion. The described signal processing portion is formed from ADSL user terminal equipment router and exchange portion, and the described ADSL user terminal equipment router and exchange portion are made up by adopting integrated chip, so the integration level of the system is high, its power consumption is low and it can support several protocols.

Entrance to Machine Translation

# Quality of MT: a Critical Topic

- Problems of uncertainty and ambiguities are common in patent document translation? (e.g. Long NPs and multiple verbs.)

- Linguistic issue in MT:

  - it is very hard to cover all the rules and adjust them to all possible variations

- But in patent document, the language is very specific, thus MT is useful

# Strategies for Tackling the Problems

- Customization of existed general domain Chinese-English MT engines

- Integration of different machine translation paradigms (rule-based, statistical, and example-based)

- Facilitating and speeding up acquiring and handling of language resources in patent documents
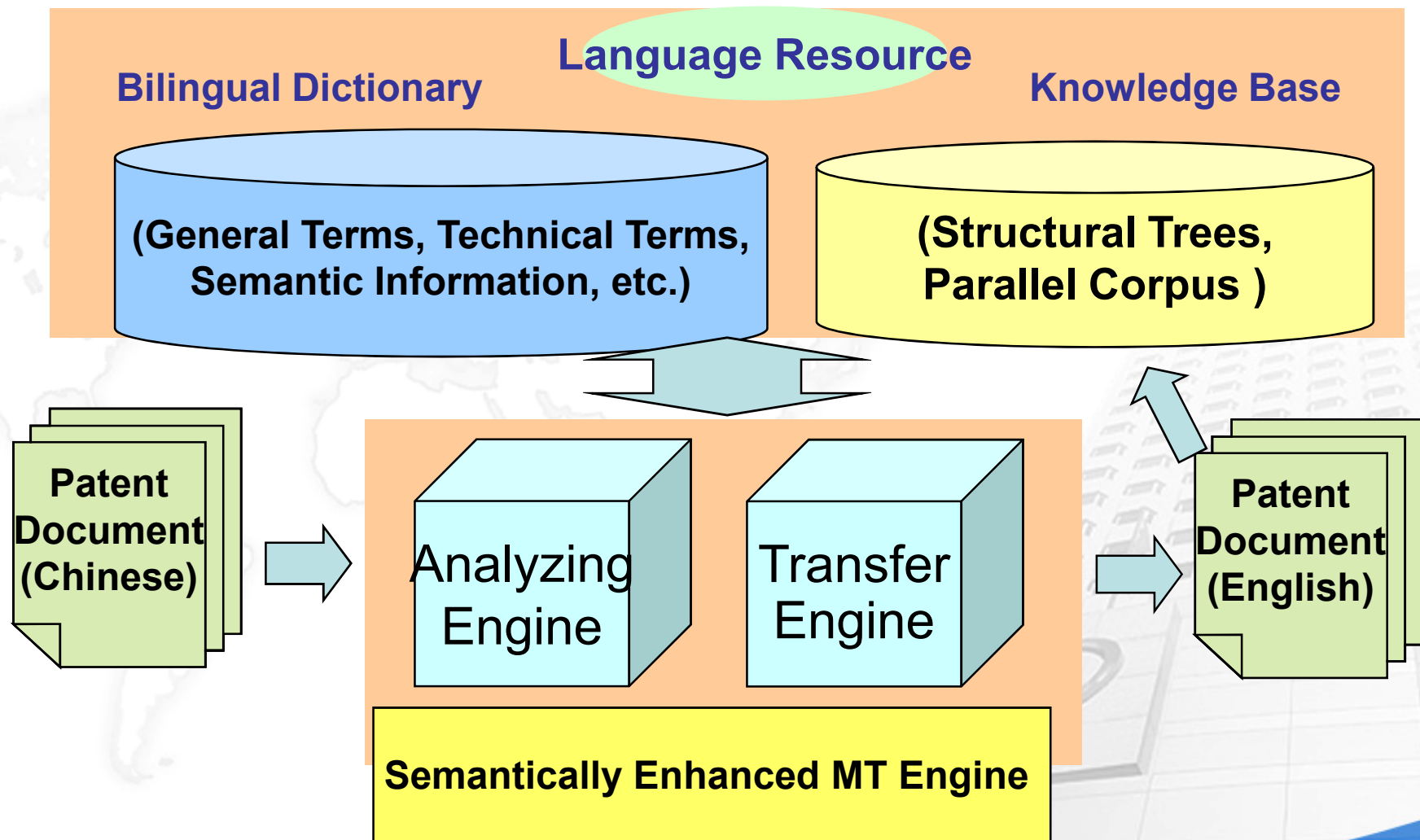
# MT System Overview

- Semantic features ad hoc: high-quality semantic-based MT engine

- Semantic analyzing: effective for disambiguation

- Representing the meaning of the source sentence and then generate the translation from the meaning
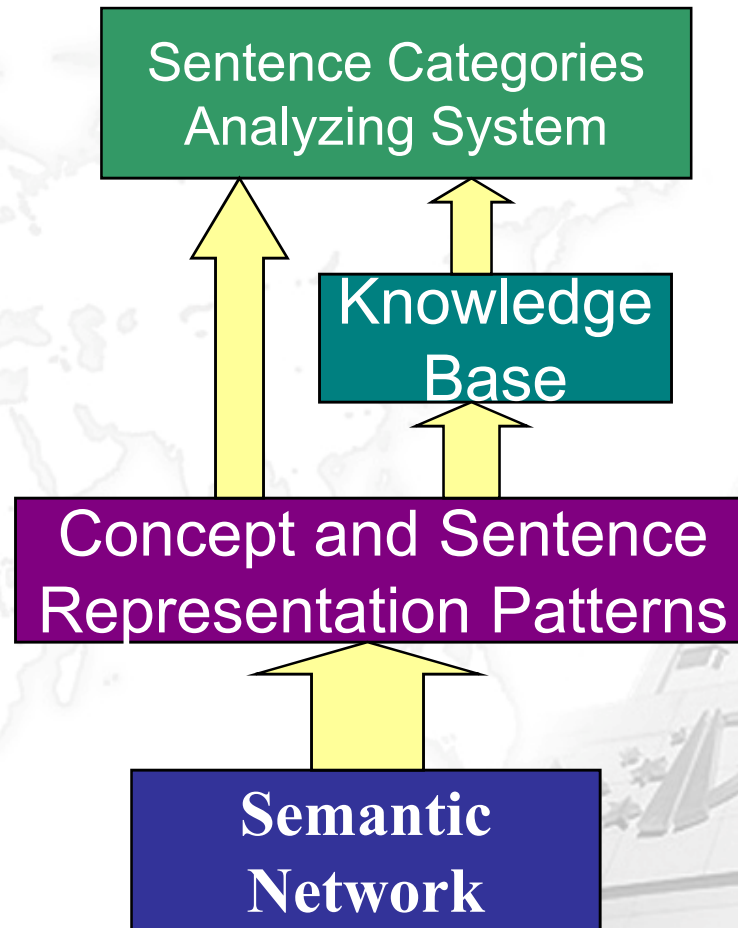  - Fix cases of syntactic mismatch

# Patent MT System Architecture

## Language Resource

**Bilingual Dictionary**

**Knowledge Base**

(General Terms, Technical Terms, Semantic Information, etc.)

(Structural Trees, Parallel Corpus )

**Patent Document (Chinese)**

Analyzing Engine

Transfer Engine

**Patent Document (English)**

**Semantically Enhanced MT Engine**

# Hierarchical Network of Concepts (HNC)



Sentence Categories Analyzing System

Knowledge Base

Concept and Sentence Representation Patterns

Semantic Network

# Sentence Categories in HNC

- Sentences are classified into 57 categories, this significantly raises the likelihood that the MT system correctly translates different sentence structures and meanings of words

- The 57 categories cover over 90 percent of the real-text sentences to be translated from Chinese into English

- Grammatical outputs are produced in most cases

# HNC Symbols
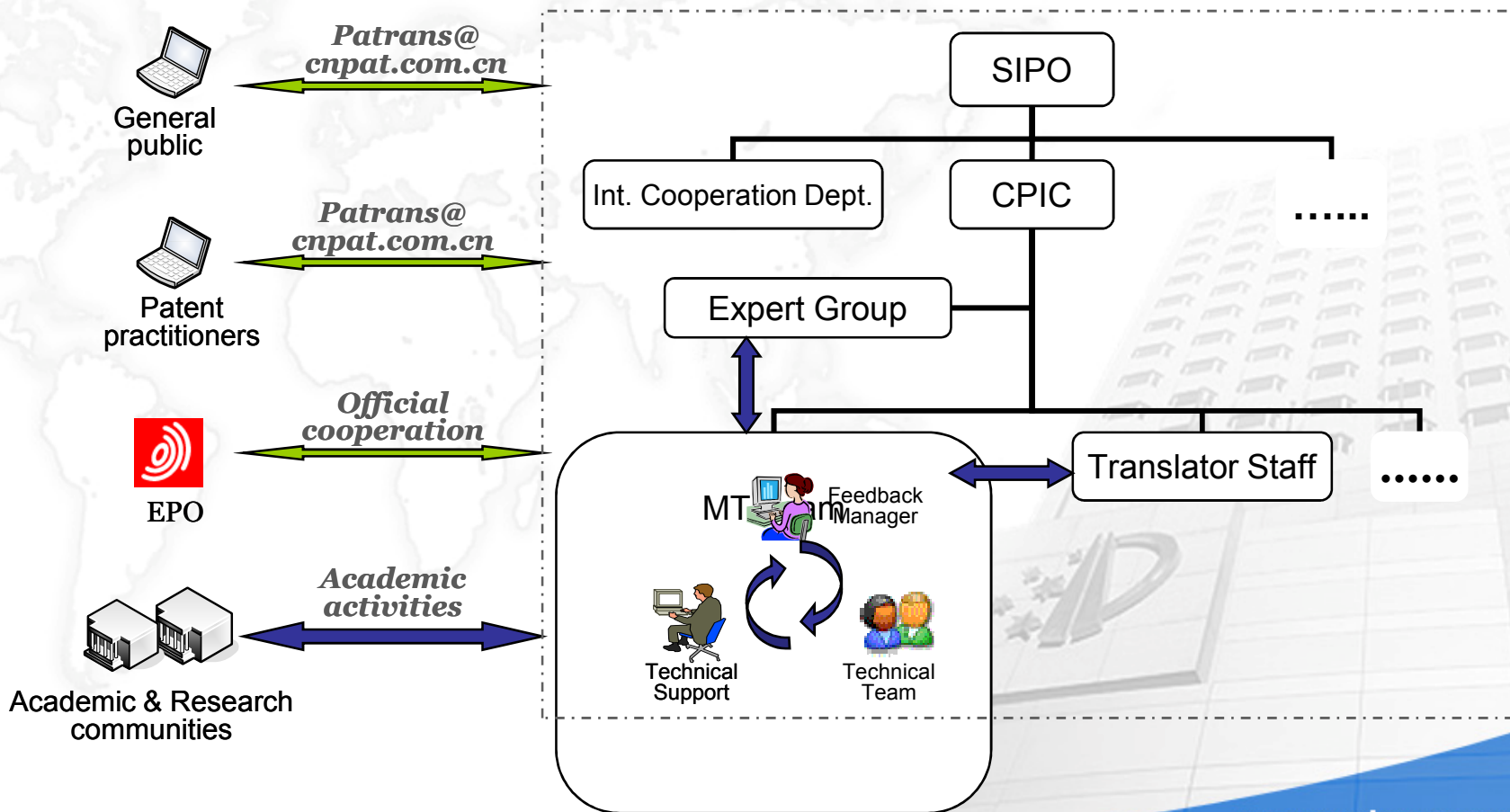
思考 v80　　　产生 v311　　情感 g713　　思维 g80　　　消除 v31

想法 r80　　　推动 v361　　承担 v901

力　g008　　　抑制 v362　　保护 v3219　责任 rc010

力量 gz00　　　调节 v360　　照顾 v653219　　　圆满 u30a

力度 z00　　　年　wj10-　　维护 v93219　完成 v30a8

弱　u00c21　月　wj10-0　保卫 vc3219　精力 gz655098

强　u00c22　日　wj10-00写作 va31　　旺盛 zu5098e71

萌芽 gv10ac41　　　体　j20-　　幸福 gu50a9ae81

成长 v10ac42　　　面　j20-0　治疗 va82　　生活 gv50a9

成熟 vu10ac43　　　线　j20-00　处方 gwa82

衰亡 v10ac44　　　点　j20-000　　药物 wa82

达成 vc249a$(v308|(jlv001/v810))

# Improvement via. Analyzing User Feedbacks and Evaluations

# Response to User Feedback

| Issues | | Responses |
|---|---|---|
| **Error Reports** | - Terminologies<br>- Inappropriate EN expression<br>- Linguistic errors | Integrated analysis for quality improvement |
| **Thoughts & Recommendations** | - MT technique<br>- Evaluation methodology | Strategy research & investigation |
| **Comments & Questions** | - Accessibility<br>- Usability | System performance optimization |

# Standardization of Test & Evaluation Workflow

◆ Test suites created based on check-point definition

◆ Full trace of each test case ensured by test->evaluation->retest and evaluation database

**Parallel corpora** → **Check-point extraction** → **Test suites** → **Test** → **MT system**

**Check-point Definition**

**Retest** **Evaluation**

**Evaluation database**

**Totally 3,200 records**

# Check-point Extraction

◆ **Check-points are now manually extracted**

➲ Source corpora :

bilingual versions of PCT documents and priority documents;

parallel abstracts from daily translation work;

existing sentence-aligned parallel corpus

➲ Principle of extraction :

adequate number of test cases per check-point to ensure system stability;
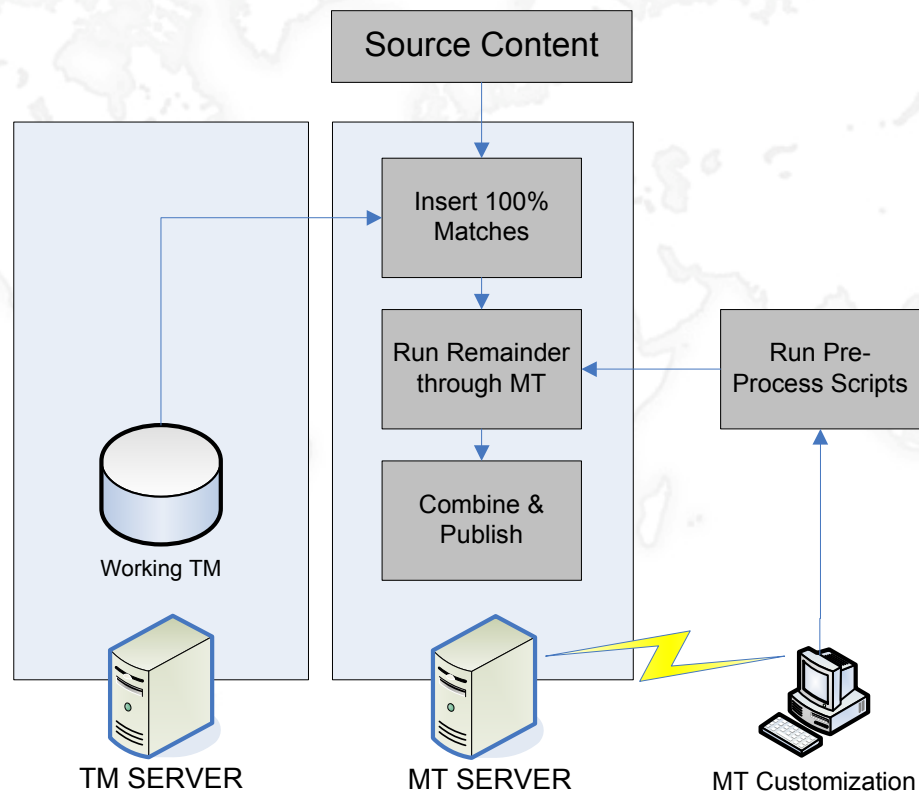
diversity of test cases for one single check-point

# Feedbacks and Check-points Definition

| Terminological | Idiomatic | Syntactic | Semantic |
|---|---|---|---|
| Technical term | Patent specific idiom | Indefinite article | Definite article |
| Patent-specific term | Technical idiom | Ordinal numbers | Word/sentence segmentation |
| UNK | Fixed collocation | Plurality | POS recognition and conversion |
| …… | Discrete structure | Tense and voice | Physical interrelationship |
| | …… | Adverbial phrase of time | Logical interrelationship |
| | | Adverbial phrase of location | Word/phrase order |
| | | Verb-subject/object collocation | Special-style CN sentence |
| | | Inappropriate EN expression | Long and complex sentence |
| | | …… | …… |

# Language Resources in Patent Documents



- Improvement for translation quality is achieved through introduction of large-scale parallel corpus
- MT system output benefited from human translations
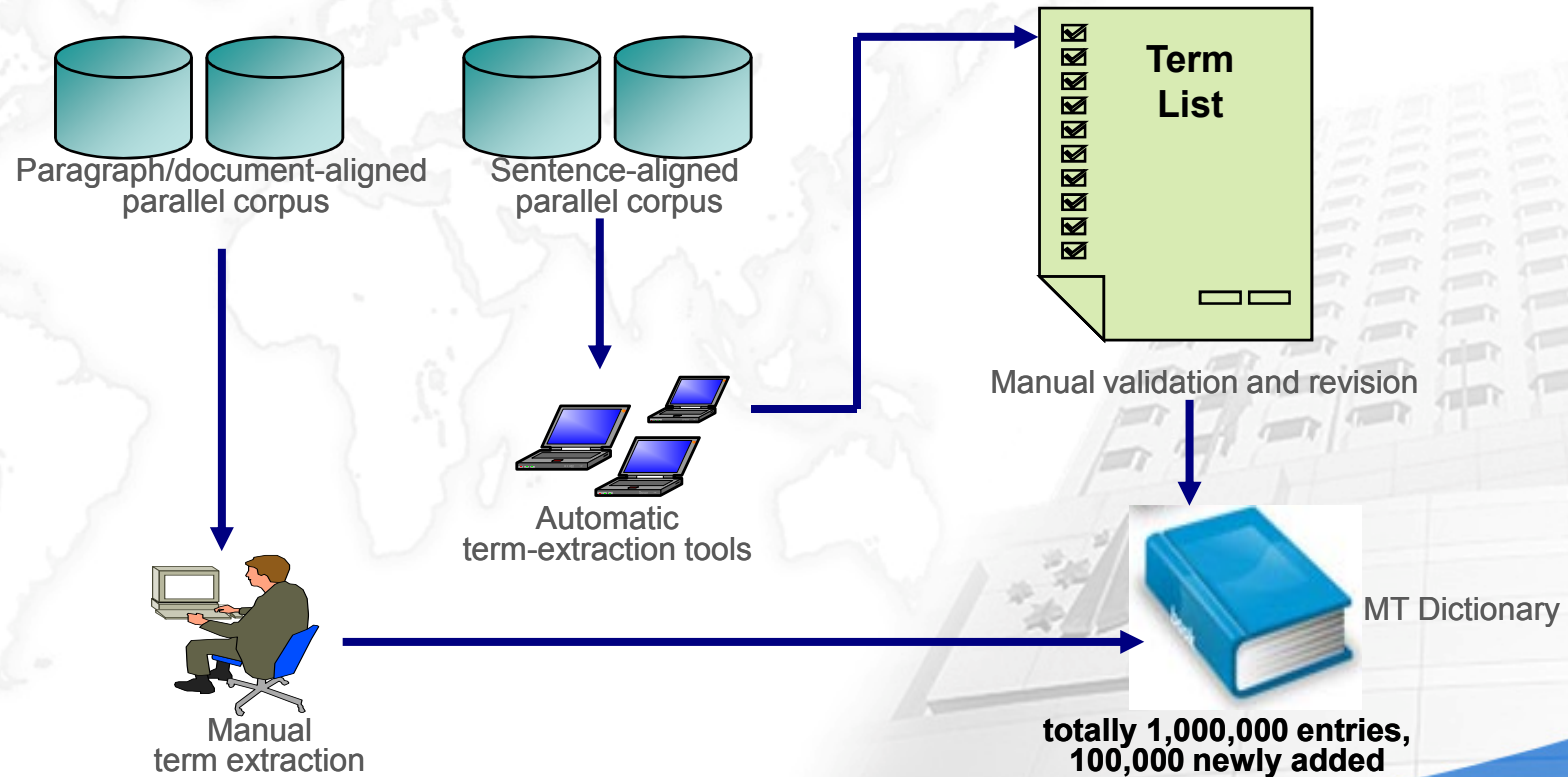- Improvement over time, also through terminology additions

Source Content

Insert 100% Matches

Run Remainder through MT

Run Pre-Process Scripts

Combine & Publish

Working TM

TM SERVER

MT SERVER

MT Customization

# Language Resources in Patent Documents

- Establishing and maintaining bilingual dictionaries for various sub-domains

| Chinese | English | | |
|---------|---------|---|---|
| 隔板 | board | ⇒ | baffle |
| 先验信息 | transcendent information | ⇒ | prior information |
| 立体 | stereo | ⇒ | three-dimensional |
| 中间结果 | middle results | ⇒ | intermediary results |
| : | : | | |

# Knowledge Acquisition Techniques

◆ Statistical-manual approach for term extraction



Paragraph/document-aligned
parallel corpus

Sentence-aligned
parallel corpus

Term
List

Manual validation and revision

Automatic
term-extraction tools

Manual
term extraction

MT Dictionary

**totally 1,000,000 entries,
100,000 newly added**

# New Words Detection Based on Large-scale Corpus

# Attempts towards a Multi-engine MT System

◆ Selection and comparison of two candidate engines
◆ Customization of the semantic engine adopting Chinese NLP techniques

Current results of MT engines comparison and customization

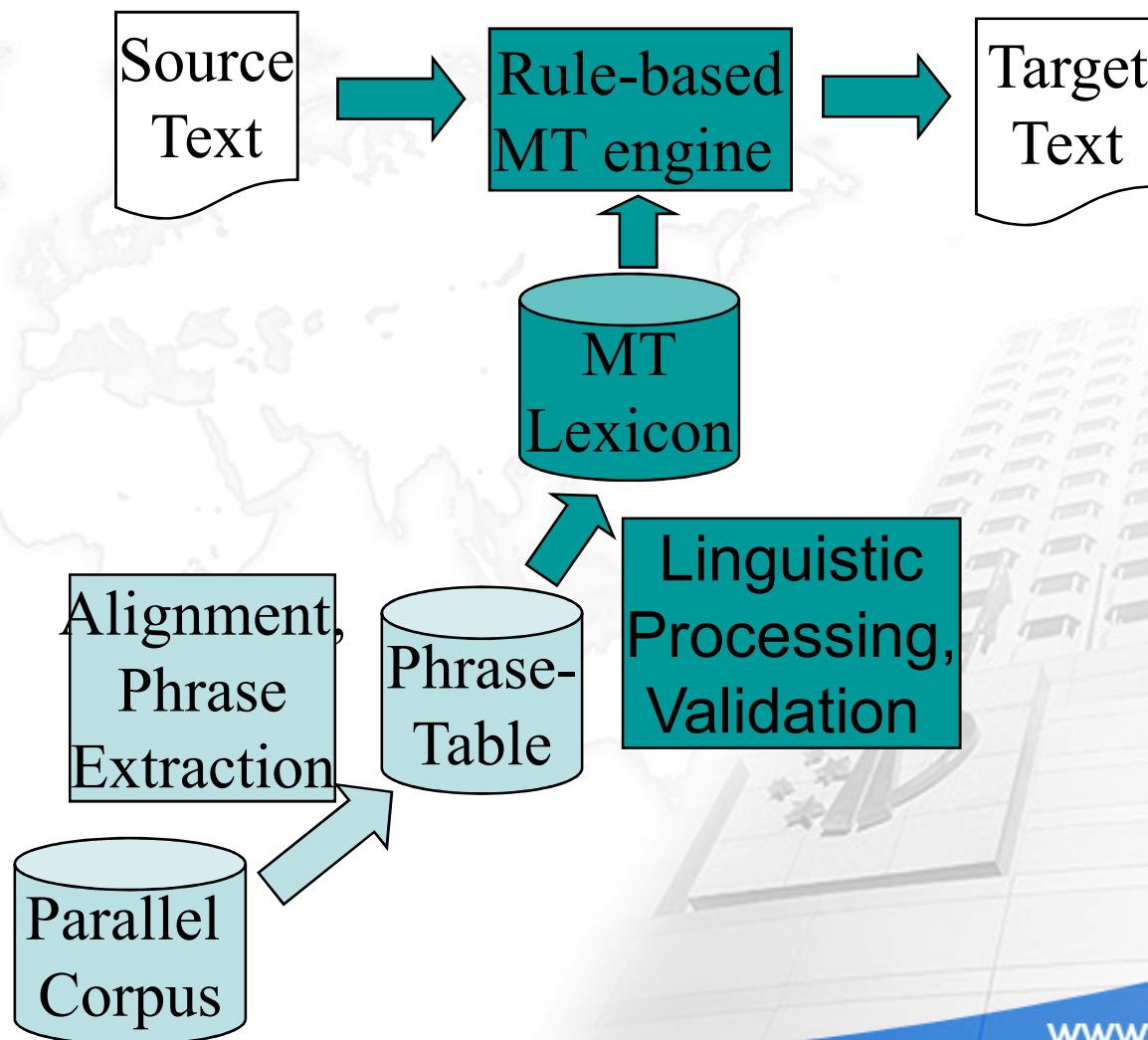| RBMT engine A | Semantic engine B |
|---|---|
| 👍 grammatical constituent combination<br>👍 English expression generation | 👍 segmentation<br>👍 logical chunk recognition and  scheduling<br>👍 CN-EN sentence style transformation<br>👍 parsing of long and complex sentences |

◆ Future work: to implement multi-engine MT system through exploring system combination strategies.
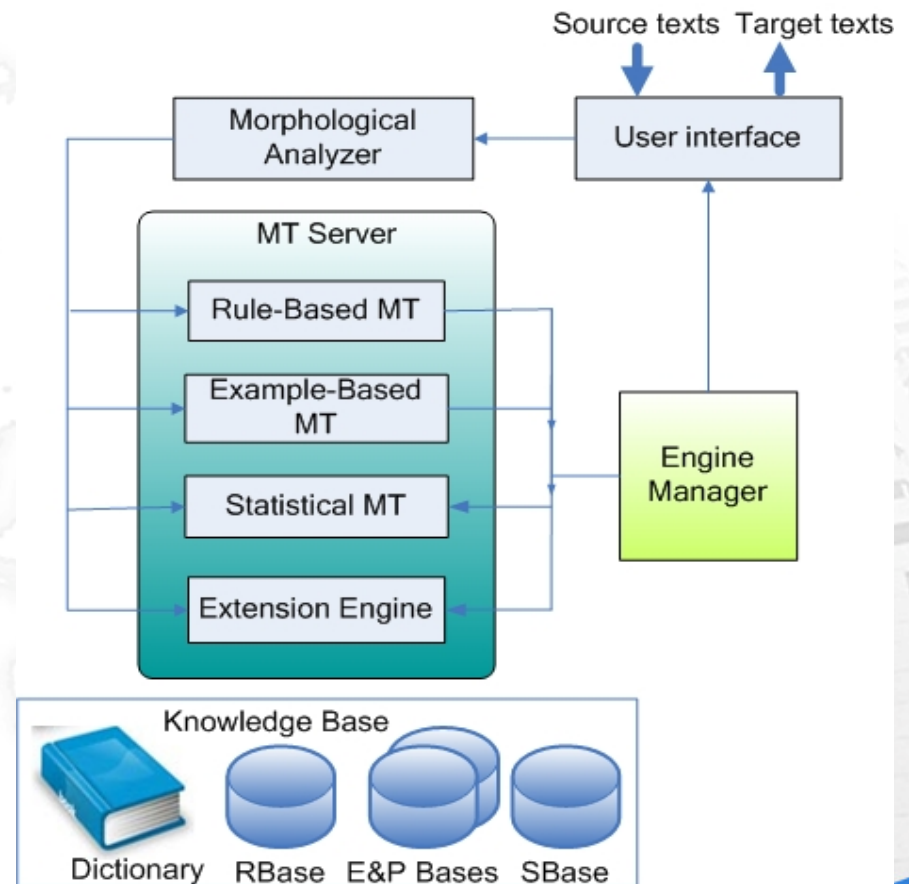
# Statistical Approach Feeds Rule-based MT

# MT Engine Combination, How to?

◆ Rule-based MT and SMT co-exist and show complementary strengths possible?

◆ System combination, how to combine the advantages?



**Future Vision of CPMT**