

# Building a Bilingual Lexicon Using Phrase-based Statistical Machine Translation via a Pivot Language

Takashi Tsunakawa\*

Naoaki Okazaki\*

Jun'ichi Tsujii\*/\*\*

\*Department of Computer Science, Graduate School of Information Science and Technology, University of Tokyo, Japan

\*\*School of Computer Science, University of Manchester / National Centre for Text Mining, UK

## Motivation

- Bilingual lexicons from/to English are usually richer than those within non-English languages  
→ **Connecting non-English terms via English (pivotal approach)** could be a reasonable approach
- We apply an **SMT framework** for the pivotal approach  
→ Characteristics between source and target languages can be directly modeled as the feature functions

## Problems of existing methods

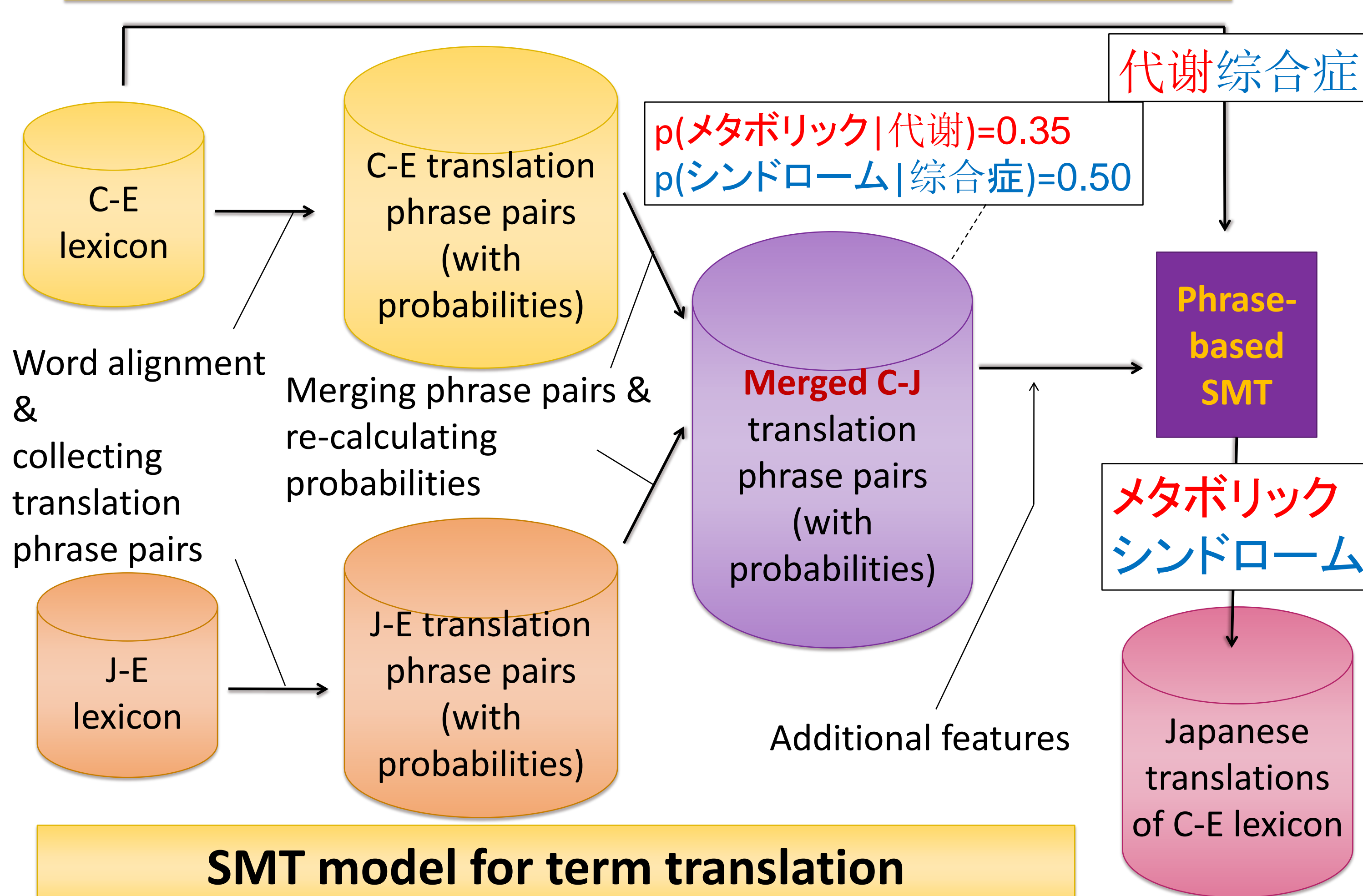
- Merging two bilingual lexicon via identical English terms can associate a few bilingual lexicons (**low utilization ratio**)
  - In the right example, only *Down's syndrome* can be associated
- Characteristics between languages are unused
  - Chinese and Japanese terms share identical/similar characters.

## Merging two bilingual lexicons

Chinese	English
代谢综合症	metabolic syndrome
道恩综合症	Down's syndrome
抗体缺乏综合症	antibody-deficiency syndrome
Japanese	English
代謝異常	metabolic disorder
ダウン症候群	Down's syndrome
メタボリックシンドローム	metabolic syndrome
抗体欠乏	antibody deficiency

Chinese	Japanese
道恩综合症	ダウン症候群
代谢综合症	メタボリックシンドローム
代谢综合症	代謝症候群
抗体缺乏综合症	抗体欠乏症候群

## Framework



## SMT model for term translation

## Features for the log-linear SMT modeling

- Merged translation probabilities** (Utiyama & Isahara, 2007)
  - Apply morphological analyzers, and obtain word alignments by GIZA++ (Och and Ney, 2003) for J-E and C-E lexicons
  - Collect phrase pairs by *grow-diag-final* method (using Moses, Koehn et al., 2007), and calculate the translation probabilities by relative frequencies

$$h_p(\bar{w}_j, \bar{w}_c) = p(\bar{w}_j | \bar{w}_c) = \frac{1}{Z_C} \sum_{\bar{w}_E} p(\bar{w}_j | \bar{w}_E) p(\bar{w}_E | \bar{w}_c),$$

$$Z_C = \sum_{\bar{w}'_j} \sum_{\bar{w}_E} p(\bar{w}'_j | \bar{w}_E) p(\bar{w}_E | \bar{w}_c)$$

- Normalized character-level edit distance** of two phrases

$$h_{ed}(w_j, w_c) = 1 - \frac{\text{Edit distance of Unicode characters}}{\text{Max. of the number of characters}}$$

$$h_{ed}(\text{後天性免疫不全症候群}, \text{後天免疫缺乏症候群}) = 1 - 3/10 = 0.7$$

- Additional lexicon**

— The number of translation word pairs included in the additional bilingual lexicon

$$h_{lex}(\text{後天性免疫不全症候群}, \text{获得性免疫缺陷综合症}) = 3$$

if (免疫, 免疫), (不全, 缺陷), and (症候群, 综合症) are in the additional lexicon

## Experiment

**Utilization ratio:** The ratio of Chinese terms translated into the other language

C-to-J	Utilization ratio
Exact matching	26.2%
Our method	72.8%

**Translation performance on the test set\***

Features	BLEU	NIST	Acc.
Base features	0.4519	7.4060	0.676
with edit distance	0.4670	7.4963	0.682
with additional lexicon	0.4800	7.5907	0.674
All	0.4952	7.7046	0.685

\*Test set consists of 500 Chinese and Japanese term pairs

**Used lexicons:** All lexicons consist of technical terms

C-E: Wanfang Data Dictionary (C: 375,990 terms / E: 429,807 terms)

J-E: JST MT Dictionary (J: 465,563 terms / 418,044 terms)

Additional lexicon: EDR J-E-C lexicon (C: 90,605 terms / J: 94,928 terms)

## Conclusion

- The experiment demonstrated that our method improved the utilization ratio drastically and performed reasonable translations of given lexicons
- We also showed that features between the source and target languages are effective for improving the performance

- Future work

- Improve Chinese word segmentation
- Introduce a sophisticated Chinese character similarity model
- Extract a source-target phrase table from corpora

## Translation examples

C	耻骨联合	→	C-to-J	恥骨 結合	○
E	pubis symphysis	→	E-to-J	結合 恥	×
J	恥骨結合			(symphysis shame)	

C	理想流体动力学	→	C-to-J	理想 流体 力学	○
E	ideal fluid dynamics	→	E-to-J	理想 液 力学	×
J	理想流体力学			(ideal fluid dynamics)	

C	中间腰淋巴结	→	C-to-J	中間 節 腰淋	×
E	intermediate lumbar lymph nodes	→	E-to-J	(int. nodes [mistokenized])	
J	中間腰リンパ節			中間 腰 リンパ 節	○