



**VERB SUBCATEGORISATION
IN ENGLISH-JAPANESE
MACHINE TRANSLATION**

Alexis Kauffmann

CONTENTS

- Introduction and related research about subcategorisation lexicons
- Context: Improvement of an English-Japanese LBMT system for better complex sentence translation
- Purpose and method for subcategorisation translation
- Application: Modifying LATL databases:
 - decription
 - difficult points
 - experiment
- Results
- Conclusion
- Next Steps of my research



INTRODUCTION AND RELATED RESEARCH ABOUT SUBCATEGORISATION LEXICONS

- As you already know (Sasano, 夏勉強会 2010)
- Famous databases: FrameNet, Propbank, SUCAT, 格フレーム
- LATL databases (Nerima, Werhli, Scherrer, 2009): -handmade monolingual and bilingual lexical databases. For example, 11544 English verb subcats, 7079 French ones, 6843 German ones.
 - transitivity generated bilingual databases
 - Japanese lexical databases



Bilingual

Source Language

Target Language

- go : V [111048081] - [NP _ S]
- go : V [111048581] - [NP _]
- go : V [111048668] - [NP _ FP]
- go (into) : V [111048692] - [NP _ PP]
- go (from) (to) : V [111049200] - [NP _ PP PP]**
- go (with) : V [111054031] - [NP _ PP]
- go (against) : V [111057199] - [NP _ PP]

- 行く (から) (まで) : V [611219792] - [NP _ PP P**
- なる : V [611161400] - [NP _]
- 行く : V [611207272] - [NP _]
- 行く (に) : V [611219793] - [NP _ PP]
- 足を運ぶ : V [611167307] - [NP _]
- 出かける : V [611174359] - [NP _]
- 作動する : V [611192433] - [NP _]

Argument Correspondence

Target Language
Correspondence

	1	2	3	No Arg
Source Language Correspondence 1	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Source Language Correspondence 2	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>
Source Language Correspondence 3	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>

Features

Preference

Translation
Context

Description

INTRODUCTION AND RELATED RESEARCH ABOUT SUBCATEGORISATION LEXICONS

- LATL databases contain data for every category of word. For verbs: syntactic subcat and semantic case frame and other syntactic and semantic information about the verb, and possibly about each argument.
- Quality and drawbacks of LATL verb databases
 - syntactically very correct (hand-made)
 - semantically, could do better,
 - could be bigger.



CONTEXT: IMPROVEMENT OF AN ENGLISH-JAPANESE LBMT SYSTEM FOR BETTER COMPLEX SENTENCE TRANSLATION

- A need to translate better complex structures
- Work on modality translation
- Work on infinitives, gerundives and object sentence translation
 - => need for subcategorisation data about verbal/sentential objects
- Also a need for less mistakes in nominal argument postpositional particles handling
 - => need for global subcategorisation frame data



PURPOSE AND METHOD FOR SUBCATEGORISATION TRANSLATION

- Purpose: Handling subcategorisation translation:
 - to know Japanese verb subcategorisation frames
 - to choose the right case frame and to transfer the arguments correctly
- Method: Using Kawahara sensei's Japanese subcategorisation frame data



APPLICATION: MODIFYING LATL DATABASES

- LATL Japanese monolingual Data
 - Taken from a Japanese dictionary file:
 - 34700 Verb entries. One entry per subcategorisation frame.
 - Transitive/intransitive only
 - => problems in bilingual correspondences
 - About 50 entries entered or corrected by hand



APPLICATION: MODIFYING LATL DATABASES: DESCRIPTION

- LATL grammatical functions:
 - Subject
 - Direct Object
 - Indirect/Prepositional Object
(+ preposition type)
 - Sentencial/verbal
(+ verb tense information)
 - Adjectival predicate



APPLICATION: MODIFYING LATL DATABASES : DESCRIPTION

- Kawahara sensei's data Japanese particles:
 - が
 - を
 - が2
 - に
 - と
 - から、まで、より、へ



APPLICATION: MODIFYING LATL DATABASES : DIFFICULT POINTS

- In order to make the right bilingual correspondences, matching of corresponding subcat frames

ex: がを ---> S O

が から まで ---> S PO(from) PO(to)

- Problem1

がに ---> S PO or O S



APPLICATION: MODIFYING LATL DATABASES : DIFFICULT POINTS

- In order to make the right bilingual correspondences, matching of corresponding subcat frames

ex: がを ---> S O

が から まで ---> S PO(from) PO(to)

- Problem1

がに ---> S PO or O S

Answer: with additional case structure data, showing which one is the agent. The agent is more likely to be the subject.



APPLICATION: MODIFYING LATL DATABASES : DIFFICULT POINTS

- Problem2 : selection of the right correspondence when several possible ones
- Problem3 : how to predict asymmetrical correspondences?
ex: 見る(がを) ---> to look (S PO(at))



APPLICATION: MODIFYING LATL DATABASES : DIFFICULT POINTS

- Problem2 : selection of the right correspondence when several possible ones

Answer: additional data indicating number of occurrences can help.

- Problem3 : how to predict asymmetrical correspondences?

ex: 見る(がを) ---> to look (S PO(at))



APPLICATION: MODIFYING LATL DATABASES : DIFFICULT POINTS

- Problem2 : selection of the right correspondence when several possible ones

Answer: additional data indicating number of occurrences can help.

- Problem3 : how to predict asymmetrical correspondences?

ex: 見る(がを) ---> to look (S PO(at))

Partial Answer: generation of an asymmetrical correspondence only when no possible symmetrical correspondence exist.



APPLICATION: MODIFYING LATL DATABASES : EXPERIMENT

- Creation of the Japanese verb entries for LATL monolingual database, formatting Kawahara sensei's data.
- Insertion of the part of the data that was not already there.
- Knowing LATL bilingual database verb correspondences (without considering subcategorisation), creation of the possible bilingual correspondences, with symmetric subcat most often, and asymmetric subcat sometimes
- Insertion of absent bilingual correspondences
- Lowering correspondence scores of unprobable existing bilingual entries



RESULTS

- About 5000 monolingual verb subcat entries inserted.
- About 2500 were already there.
- About 8000 bilingual correspondences entered and many old correspondences bilingual almost deleted.



CONCLUSION

- Japanese subcategorisation knowledge much better in the database
- Need to test the results with the MT system
- Need to improve by hand subcat for verbs with 3 arguments or more, verbal objects and adjectival objects.



9) NEXT STEPS OF MY RESEARCH

- Finishing improving subcat data
- Finishing programming translation of modality and verbal/sentential objects
- Improving collocation/multi-word expressions translation

