

Using Information Structure to Improve Transfer-based MT

Sanghoun Song (UW)
Emily M. Bender (UW)

HPSG2011
Workshop on Information Structure and Formal Grammar
08-23-2011

Introduction

Puzzle

Information Structure in HPSG/MRS

Information Structure in English and Japanese/Korean

Translation

Conclusion

Outline

Introduction

Puzzle

Information Structure in HPSG/MRS

Information Structure in English and Japanese/Korean

Translation

Conclusion

Hypothesis

- ▶ The essential part of translation is reshaping the means of conveying information, instead of simply changing the words or reordering phrases.

Working Hypothesis

- ▶ Transfer-based MT systems can be improved
 - ▶ by encoding information structure in both the source and target grammars
 - ▶ by preserving information structure in the transfer stage

Goal

- ▶ How information structure can be represented in HPSG/MRS
- ▶ How information structure can help refine multilingual MT
- ▶ A sample translation between English & Japanese/Korean

Outline

Introduction

Puzzle

Information Structure in HPSG/MRS

Information Structure in English and Japanese/Korean

Translation

Conclusion

Allosentences

- ▶ close paraphrases which share truth conditions (Lambrecht, 1996)
- ▶ They are not always in the same felicity conditions.

(1) a. I am Kim.

b. watashi-ga/wa Kim desu.

I-NOM/TOP Kim COP [jpn]

(2) Q: Who are you?

A: watashi-#ga/wa Kim desu.

(3) Q: Who is Kim?

A: watashi-ga/#wa Kim desu.

Active/Passive Pairs

- (4) a. Kim tore the book.
 b. The book was torn by Kim.
- (5) a. Kim-ga sono hon-o yabut-ta.
 Kim-NOM DET book-ACC tear-PST
 ‘Kim tore the book.’
- b. ?sono hon-ga Kim-ni yabu-rare-ta.
 DET book-NOM Kim-DAT tear-PASS-PST
 ‘The book was torn by Kim.’ [jpn]

Other Issues

- ▶ (4b) has at least eight allosentences in Japanese ($2 \times 2 \times 2$)
 - ▶ actives vs. passives
 - ▶ case markers (e.g. *ga*) vs. topic markers (e.g. *wa*)
 - ▶ scrambling: SOV vs. OSV (Choi, 1999; Ishihara, 2001)

- (6) a. Kim-ga/wa sono hon-o/wa yabut-ta.
 Kim-NOM/TOP DET book-ACC/TOP tear-PST
- b. sono hon-o/wa Kim-ga/wa yabut-ta.
 DET book-ACC/TOP Kim-NOM/TOP tear-PST

- ▶ two more options: 32 allosentences, in total ($8 \times 2 \times 2$)
 - ▶ pro-drop
 - ▶ null NP markings
- ▶ not felicitous in the same contexts.

Information Structure

- ▶ In order to solve the mismatch and refine the translations, we need to use Information Structure.
 - ▶ The difference in felicity conditions between allosentences is the subject of study of information structure.
 - ▶ Information structure is hypothesized to be universal.
 - ▶ All languages have some way to mark topics and foci.
 - ▶ pitch accent (e.g. A/B-accent)
 - ▶ specific word order (e.g. scrambling)
 - ▶ morphological marking (e.g. *wa*)
 - ▶ some combination of them
- ▶ This study looks at the particular case of translating English passives into Japanese/Korean.

Outline

Introduction

Puzzle

Information Structure in HPSG/MRS

Information Structure in English and Japanese/Korean

Translation

Conclusion

Basic Assumptions

1. Sentences always have at least one focus, but they do not always have a topic; further, constituents may be 'background' (i.e. neither topic nor focus).
2. 'contrast' as a cross-cutting information structure category (Molnár, 2002)

(7) Kim thì đi Hanoi
Kim CT go Hanoi

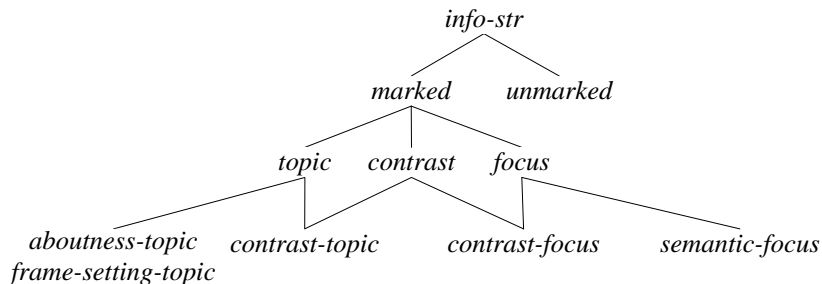
'Kim goes to Hanoi(, but nobody else).' [vie]

Basic Assumptions (cont'd)

3. Semantically empty categories (e.g. complementizers, expletives) are informatively empty as well (i.e. assigned no category).
 - (8) a. It is Kim ~~who~~ tore the book.
 - b. The book ~~was~~ torn ~~by~~ Kim.
4. It is assumed that the canonical position of topics is sentence-initial at least in our sample of languages (English, Japanese, and Korean).

Information Structure in HPSG/MRS

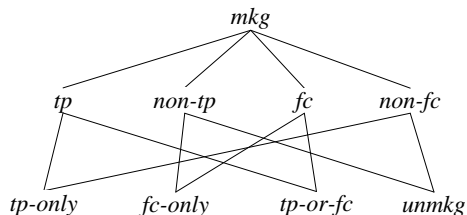
- ▶ INFO-STR: a semantic feature in the MRS
- ▶ MKG: a categorical feature encoding the lexical marking
- ▶ SFORM: a set of constraints on phrase structure rules.

MRS: *info-str*Figure: Type Hierarchy of *info-str*

- ▶ representing information structure with a feature on indices directly in the MRS

Marking: *mkg*

$$\left[\text{MKG} \begin{bmatrix} \text{TP} & \text{bool} \\ \text{FC} & \text{bool} \end{bmatrix} \right]$$

Figure: Type Hierarchy of *mkg*

Sentential Forms: *sform*

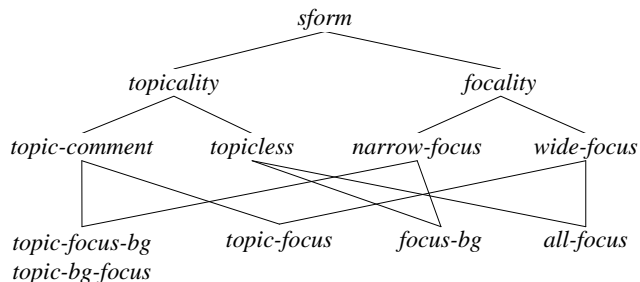


Figure: Type Hierarchy of *sform*: adapted from Paggio (2009), etc.

Sentential Forms: *sform* (cont'd)

- ▶ *topicality*

- ▶ *topic-comment*

(9) As for **the book**, KIM tore it.

- ▶ *topicless*

(10) It is KIM ~~who~~ tore the book.

Sentential Forms: *sform* (cont'd)

▶ *focality*

- ▶ *narrow-focus*
- ▶ *wide-focus*

- (11) a. ínam-é yimi
boy-F came
'THE BOY came.'
- b. ínam á-yimi
boy F-came
'The boy CAME.' (Rendille [red] (Lecarme, 1999))

Sentential Forms: *sform* (cont'd)

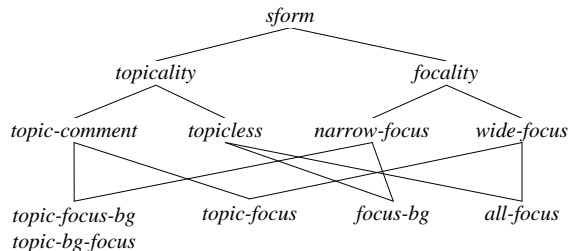


Figure: Type Hierarchy of *sform*

- (12) a. **The book** was torn by [_f KIM] (*topic-bg-focus*).
 b. **The book** [_f was torn] by Kim (*topic-focus-bg*).
 c. **The book** [_f was torn by Kim] (*topic-focus*).
 d. [_f THE BOOK] was torn by Kim (*focus-bg*).
 e. [_f The book was torn by Kim] (*all-focus*).

Outline

Introduction

Puzzle

Information Structure in HPSG/MRS

Information Structure in English and Japanese/Korean

Translation

Conclusion

Information Structure in English

- ▶ Information structure is normally constrained by pitch accents.
 - ▶ A-accented phrases (H*) take *focus*

$$\begin{array}{l} fp\text{-lex-rule} \rightarrow \\ \left[\begin{array}{l} \text{PROSODY } A\text{-accent} \\ \text{INFO-STR } focus \end{array} \right] \end{array}$$

- ▶ B-accented phrases (L+H*) take *topic*

$$\begin{array}{l} tp\text{-lex-rule} \rightarrow \\ \left[\begin{array}{l} \text{PROSODY } B\text{-accent} \\ \text{INFO-STR } topic \end{array} \right] \end{array}$$

Information Structure in Japanese/Korean

- ▶ topic markers: aboutness vs. contrastiveness

(13) Q: Who is a student?

A: Kim-ga/wa gakusei desu.

Kim-NOM/TOP student COP

'Kim is a student.' [jpn]

Information Structure in Japanese/Korean (cont'd)

▶ scrambling

- (14) a. Kim-wa sono hon-o yabut-ta.
Kim-TOP DET book-ACC tear-PST
(topic)
- b. sono hon-o Kim-wa yabut-ta.
DET book-ACC Kim-TOP tear-PST
(contrastive-focus)
- c. Kim-ga sono hon-wa yabut-ta.
Kim-NOM DET book-TOP tear-PST
(contrastive-focus)
- d. sono hon-wa Kim-ga yabut-ta.
DET book-TOP Kim-NOM tear-PST
(contrastive-topic) [jpn]

Information Structure in Japanese/Korean (cont'd)

Table: IS of topic-marked NP (adapted from Choi (1999))

	in-situ	scrambling
subject	<i>topic</i>	<i>contrast-focus</i>
non-subject	<i>contrast-focus</i>	<i>contrast-topic</i>

<i>nom-marker</i> →	<i>topic-marker</i> →
$\left[\begin{array}{l} \text{ORTH } \langle ga \rangle \\ \text{MKG } unmkg \\ \text{CASE } nom \end{array} \right]$	$\left[\begin{array}{l} \text{ORTH } \langle wa \rangle \\ \text{MKG } tp \\ \text{CASE } case \end{array} \right]$

Information Structure in Japanese/Korean (cont'd)

$$\left[\begin{array}{ll} \textit{topic-comment} & \\ \text{MKG} & \textit{tp} \\ \text{HD} \mid \text{MKG} & \textit{fc} \\ \text{NON-HD} \mid \text{MKG} & \textit{tp} \end{array} \right]$$

- ▶ INFO-STR in Japanese and Korean is specified at the phrasal level (i.e. each grammatical rule, such as *subj-head* and *comp-head*) unlike English

(15) sono hon-wa Kim-ga yabut-ta.
 DET book-TOP Kim-NOM tear-PST
 'The book was torn by [_f KIM].' [jpn]

$$\left[\begin{array}{l} \textit{top-scr-subj-head} \\ \text{HD} \mid \text{VAL} \mid \text{COMPS} \langle \square \rangle \\ \text{NON-HD} \mid \text{INFO-STR} \textit{contrast-focus} \end{array} \right]$$

$$\left[\begin{array}{l} \textit{top-scr-comp-head} \\ \text{HD} \mid \text{VAL} \mid \text{COMPS} \langle \rangle \\ \text{NON-HD} \mid \text{INFO-STR} \textit{contrast-topic} \end{array} \right]$$

Outline

Introduction

Puzzle

Information Structure in HPSG/MRS

Information Structure in English and Japanese/Korean

Translation

Conclusion

A Sample Translation

The book was torn by Kim.

1. Kim-ga sono hon-o yabut-ta.
2. sono hon-o Kim-ga yabut-ta.
3. Kim-wa sono hon-o yabut-ta.
4. sono hon-o Kim-wa yabut-ta.
5. Kim-ga sono hon-wa yabut-ta.
6. sono hon-wa Kim-ga yabut-ta.
7. Kim-wa sono hon-wa yabut-ta.
8. sono hon-wa Kim-wa yabut-ta.

A Sample Translation (cont'd)

The book was torn by [_f KIM].

1. Kim-ga sono hon-o yabut-ta.
2. sono hon-o Kim-ga yabut-ta.
3. Kim-wa sono hon-o yabut-ta.
4. sono hon-o Kim-wa yabut-ta.
5. Kim-ga sono hon-wa yabut-ta.
6. sono hon-wa Kim-ga yabut-ta.
7. Kim-wa sono hon-wa yabut-ta.
8. sono hon-wa Kim-wa yabut-ta.

A Sample Translation (cont'd)

The book was torn by [_f KIM].

~~1. Kim-ga sono hon-o yabut-ta.~~

~~2. sono hon-o Kim-ga yabut-ta.~~

~~3. Kim-wa sono hon-o yabut-ta.~~

~~4. sono hon-o Kim-wa yabut-ta.~~

~~5. Kim-ga sono hon-wa yabut-ta.~~

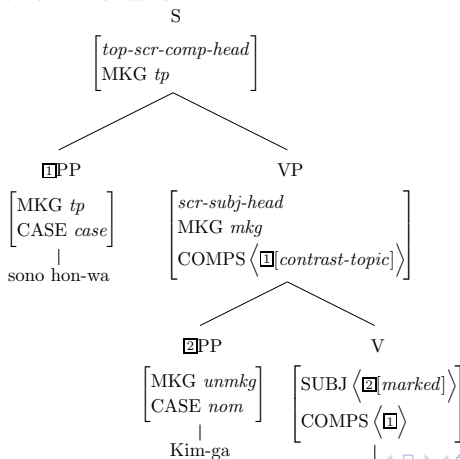
~~6. sono hon-wa Kim-ga yabut-ta.~~

~~7. Kim-wa sono hon-wa yabut-ta.~~

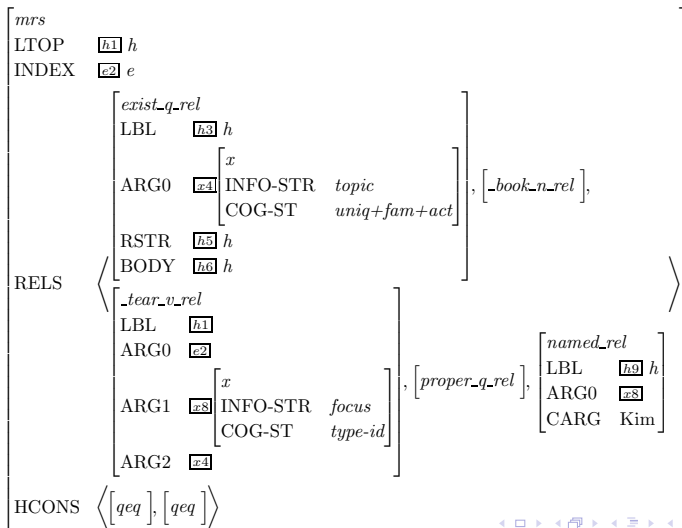
~~8. sono hon-wa Kim-wa yabut-ta.~~

A Sample Translation (cont'd)

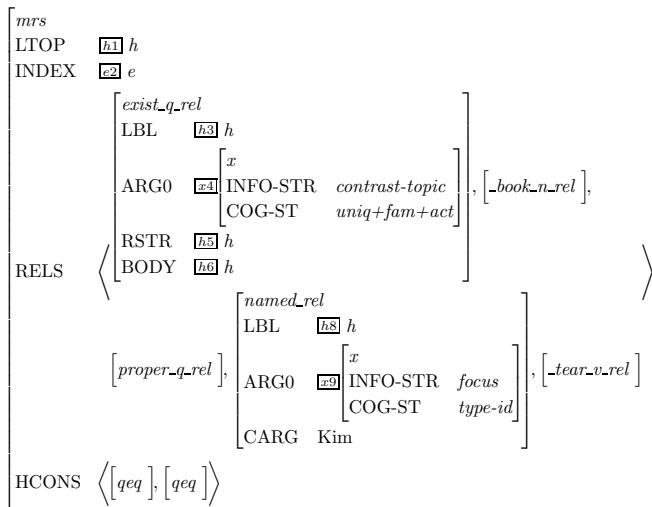
- (16) a. **The book** was torn by [_f KIM].
 b. sono hon-wa Kim-ga yabut-ta.
 DET book-TOP Kim-NOM tear-PST



An Input MRS (English)



An Output MRS (Japanese)



Experiment: Translating Passives

▶ process

1. constructing toy grammars for English, Japanese, and Korean
2. adding other rules to produce allosentences
3. implementing information structure into each grammar
4. creating the mapping between internal and external features of indices (semi.vpm)

Experiment: Translating Passives (cont'd)

- ▶ eight types of allosentences in English, for each of the three verbal types (i.e. 24 input sentences)
 - ▶ 'tear'
 - ▶ 'chase'
 - ▶ 'hit'

Experiment: Translating Passives (cont'd)

- ▶ hypothetical suffixes
 - ▶ '-FP': A-accent for foci
 - ▶ '-TP': B-accent for topic

- (17) a. **The book** was torn by [_f KIM].
b. The book-TP was torn by Kim-FP.

Results

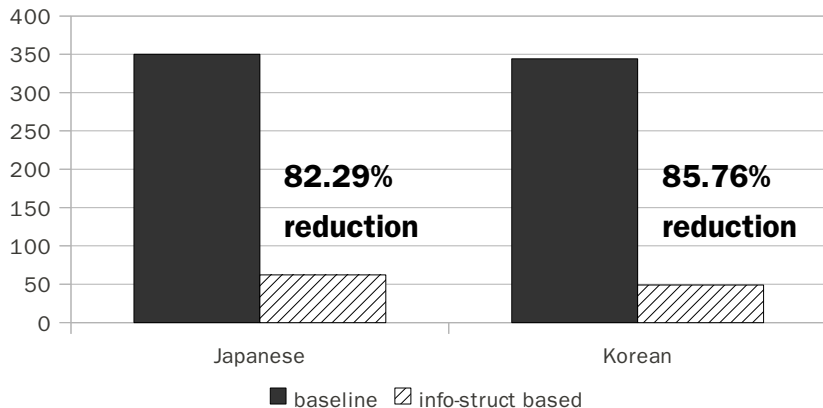


Figure: Evaluation

Contribution

- ▶ latency
 - ▶ The processing burden of MT component which ranks the translations and select only suitable results can be greatly lightened
- ▶ accuracy
 - ▶ Though it is still necessary to harness a re-ranking model for choosing translations, we can start from once-refined sets of translations.

Outline

Introduction

Puzzle

Information Structure in HPSG/MRS

Information Structure in English and Japanese/Korean

Translation

Conclusion

Conclusion

- ▶ So far, we have covered
 - ▶ how Information Structure can be represented in HPSG/MRS
 - ▶ how the representation can be used to refine translations
- ▶ Implications
 - ▶ Type hierarchies in this proposal are constructed almost language-independently.
 - ▶ We effectively move further up the MT pyramid (Vauquois, 1968).

Future Work

- ▶ evaluating with various sentences
 - ▶ dropped elements
 - ▶ clefting
- ▶ corpus studies: using multilingual texts
- ▶ other languages
 - ▶ Japanese/Korean to English
 - ▶ Chinese, Spanish, Russian
- ▶ a Grammar Matrix library for information structure