

Syntax and Semantics of Korean Numeral Classifier
Constructions

Jong-Bok Kim and Jaehyung Yang

Kyung Hee University and Kangnam University

Proceedings of the HPSG07 Conference

Stanford Department of Linguistics and CSLI's LinGO Lab

Stefan Müller (Editor)

2007

CSLI Publications

<http://csli-publications.stanford.edu/>

Abstract

The so-called floating quantifier constructions in languages like Korean display intriguing properties whose successful processing can prove the robustness of a parsing system. This paper shows that a constraint-based analysis, in particular couched upon the framework of HPSG, can offer us an efficient way of analyzing these constructions together with proper semantic representations. It also shows how the analysis has been successfully implemented in the LKB (Linguistic Knowledge Building) system.

1 Issues

One of the most salient features in languages like Korean is the complex behavior of numeral classifiers (Num-CL) linked to an NP they classify. Among several types of Num-CL constructions, the most complicated type includes the one where the Num-CL floats away from its antecedent:

- (1) *pemin-i cengmal sey myeng-i/*-ul te iss-ta*
criminal-NOM really three CL-NOM/ACC more exist-DECL
'There are three more criminals.'

There also exist constraints on which arguments can 'launch' floating quantifiers (FQ). Literature has proposed that the antecedent of the FQ needs to have the identical case marking as in (1). However, issues become more complicated with raising and causative constructions where the two do not agree in the case value:

- (2) a. *haksayng-tul-ul sey myeng-i/ul chencay-i-lako mit-ess-ta.*
student-PL-ACC three-CL-NOM/*ACC genius-COP-COMP believed
'(We) believed three students to be genius.'
- b. *haksayng-tul-ul sey-myeng-i/ul/*eykey ttena-key hayessta*
student-PL-ACC three-CL-NOM/ACC/*DAT leave-COMP did
'(We) made three students to leave.'

As given in the raising (2a) and causative (2b), the Num-CL *sey myeng* 'three CL' can have a different case marking from its antecedent, functioning as the matrix object. In a sense, it is linked to the original grammatical function of the raised object and the causee, respectively.

Central issues in deep-parsing numeral classifier constructions thus concern how to generate such FQ constructions and link the FQ with its remote antecedent together with appropriate semantics (cf. Kang 2002). This paper provides a typed feature structure grammar, HPSG, together with Minimal Recursion Semantics (MRS), is well-suited in providing the syntax and semantics of these constructions for computational implementations.

[†]We thank three anonymous reviewers for the constructive comments and suggestions. This work was supported by the Korea Research Foundation Grant (KRF-2005-042-A00056).

2 An Analysis

2.1 Forming a Numeral-Classifier Sequence and its Semantics

The starting point of our analysis is forming well-formed Num-CL expressions.¹ Syntactically, numeral classifiers are a subclass of nouns (for Japanese see Bond and Paik (2000), Bender and Siegel (2004)). However, unlike common nouns, they cannot stand alone and must combine with a numeral or a limited set of determiners as in **(twu) kay* ‘two CL’ (Numeral) and **(myech) kay* ‘how many’ (Interrogative).² Semantically, there are tight sortal constraints between the classifiers and the nouns (or NPs) they modify. For example, *pen* can classify only events, *tay* machinery, and *kwuen* just books. Such sortal constraints block classifiers like *tay* from modifying thin entities like books as in **chayk twu tay* ‘book two-CL’. Reflecting these syntactic and semantic properties, we can assign the following lexical information to numerals (*num-det*) and classifiers (*cl-n*) within the feature structure system of HPSG and MRS (cf. Copestake et al. 2006).

$$(3) \text{ a. } \left[\begin{array}{l} \text{num-det} \\ \text{ORTH } \langle \text{sey '세'} \rangle \\ \text{SYN} \mid \text{HEAD} \left[\begin{array}{l} \text{POS } \textit{det} \\ \text{NUM } + \end{array} \right] \\ \text{SEM} \left[\begin{array}{l} \text{HOOK} \left[\begin{array}{l} \text{INDEX } i \\ \text{LTOP } h2 \end{array} \right] \\ \text{RELS} \left\langle \left[\begin{array}{l} \text{PRED } \textit{card_rel} \\ \text{LBL } h2 \\ \text{ARG0 } i \\ \text{CARG } 3 \end{array} \right] \right\rangle \end{array} \right] \end{array} \right]$$

¹We have inspected the Sejong Treebank Corpus to figure out the distributional frequency of Korean numeral classifiers in real texts. From the corpus of total 378,689 words (33,953 sentences), we identified 694 occurrences of numeral classifier expressions. Of these 694 examples, we identified 36 FQ examples.

²A limited set of common nouns such as *salam* ‘person’, *kulus* ‘vessel’, *can* ‘cup’, *khep* ‘cup’, and *thong* ‘bucket’ can also function as classifiers.

$$\begin{array}{l}
\text{b.} \left[\begin{array}{l}
cl-n \\
\text{ORTH } \langle \text{myeng '명'} \rangle \\
\text{SYN} \left[\begin{array}{l}
\text{HEAD} \left[\begin{array}{l} \text{POS } noun \\ \text{CLTYPE } + \end{array} \right] \\
\text{VAL} | \text{SPR } \langle \left[\begin{array}{l} \text{NUM } + \\ \text{INDEX } i \end{array} \right] \rangle \\
\text{HOOK} \left[\begin{array}{l} \text{INDEX } i \\ \text{LTOP } h1 \end{array} \right] \\
\text{SEM} \left[\begin{array}{l} \text{RELS } \langle \left[\begin{array}{l} \text{PRED } person_rel \\ \text{LBL } h1 \\ \text{ARG0 } i \end{array} \right] \rangle \end{array} \right]
\end{array} \right]
\end{array} \right]
\end{array}$$

The feature structure in (3a) represents that there exists an individual x whose CARG (constant argument) value is “3”. The feature NUM is assigned to the numerals as well as to determiners like *yele* ‘several’ and *myech* ‘some’ which combine with classifiers. Meanwhile, (3b) indicates that syntactically a classifier selects a NUM element through the SPR, whereas semantically it belongs to the ontological category *person_rel*. The feature CLTYPE differentiates classifiers from common nouns. An independent grammar rule then ensures that only [NUM +] elements can combine with the [CLTYPE +] expression, ruling out unwanted forms such as **ku myeng* ‘the CL’.

2.2 Dealing with FQ Constructions

As noted earlier, the Num-CL can float away from the NP it classifies. There exist several supporting phenomena indicating that the FQ modifies the following verbal expression. One phenomenon is the substitution by the proverb *kule-* ‘do so’. As noted in (4), unlike the NI type, only in the NC type, an FQ and the following main verb can be together substituted by the proverb *kulay-ss-ta*:

- (4) a. *namca-ka [sey myeng o-ass-ko], yeca-to kulay-ss-ta*
 man-NOM three CL come-PST-CONJ woman-also do-PST-DECL.
 ‘As for man, three came, and as for woman, the same number came.’
 b. **[namca sey myeng-i] o-ass-ko, yeca-to [kulay-ss-ta]*

This means that the FQ in the NC type is a VP modifier, though it is linked to a preceding NP.

Coordination data also support a VP modifier analysis:

- (5) *[namhaksayng-kwa] kuliko [yehaksayng-i] [sey myeng-i] oassta*
 boy student-and and girl student-NOM three CL-NOM came
 ‘The total 3 of boys and girls came.’

The FQ ‘three-CL’ cannot refer to only the second conjunct ‘girl students’: its antecedent must be the total number of boys and girls together. This means the FQ refers to the whole NP constituent as its reference. This implies that an analysis in which the FQ forms a constituent with the preceding NP then cannot ensure the reading such that the number of boys and girls is in total three.

Given this VP-modifier treatment, the following question is how to link an FQ with its appropriate antecedent. There exist several constraints in identifying the antecedents. When the floating quantifier is case-marked, it seems to be linked to an argument with the same case marking. However, further complication arises from examples in which either the antecedent NP or the FQ are not marked with a case marker, but a delimiter or topic marker:

- (6) a. haksayng-tul-i/un sakwa-lul sey kay-lul mekessta
 student-PL-NOM/TOP apple-ACC three CL-ACC eat
 ‘As for the students, they ate three apples.’
 b. sakwa-lul haksayng-tul-i/un sey kay-lul mekessta

The data suggest that a surface case marking cannot be a sole indicator for the linking relation, and that we need to refer to grammatical functions. What we can observe is that, regardless of the location, the NOM-marked FQ is linked to the subject whereas the ACC-marked FQ is linked to the object. This observation is reflected in the following lexical information given to the type *num-cl-mw* (*numeral-classifier-multiword*):³

- (7) a.
$$\left[\begin{array}{l} \text{num-cl-mw} \\ \text{ORTH } \langle \text{sey myeng-i} \rangle \\ \text{HEAD} \left[\begin{array}{l} \text{POS } \textit{noun} \\ \text{CASE} \mid \text{GCASE } \textit{nom} \\ \text{MOD} \left\langle \left[\begin{array}{l} \text{POS } \textit{verb} \\ \text{SUBJ } \langle \text{NP}_i \rangle \end{array} \right] \right\rangle \end{array} \right] \\ \text{SEM} \mid \text{HOOK} \mid \text{INDEX } i \end{array} \right]$$
- b.
$$\left[\begin{array}{l} \text{num-cl-mw} \\ \text{ORTH } \langle \text{sey myeng-ul} \rangle \\ \text{HEAD} \left[\begin{array}{l} \text{POS } \textit{noun} \\ \text{CASE} \mid \text{GCASE } \textit{acc} \\ \text{MOD} \left\langle \left[\begin{array}{l} \text{POS } \textit{verb} \\ \text{COMPS } \langle \text{NP}_{i,\dots} \rangle \end{array} \right] \right\rangle \end{array} \right] \\ \text{SEM} \mid \text{HOOK} \mid \text{INDEX } i \end{array} \right]$$

³When the FQ has a delimiter marker (rather than a case marker) or no marker at all, it will refer to one of the elements in the ARG-ST (argument structure). Its antecedent will be determined in context.

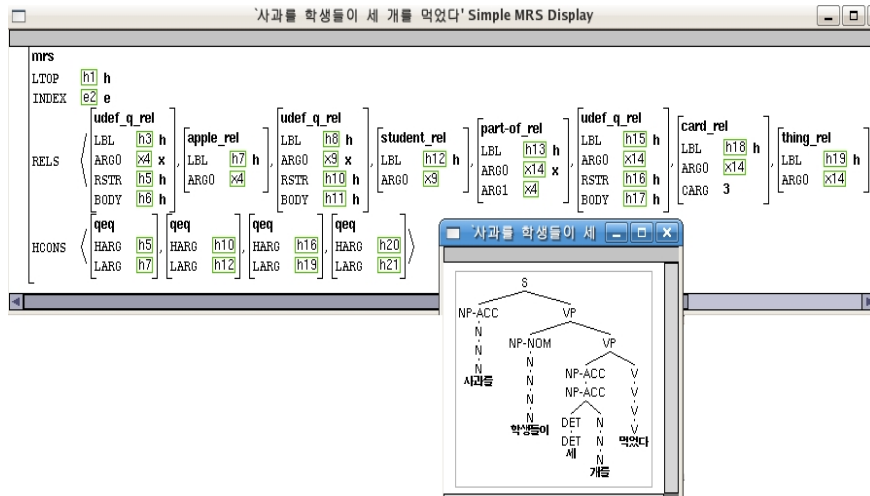


Figure 1: Parsed Tree and MRS for ‘As for the students, they ate three apples.’

As given in (7), the NOM-marked *num-cl-mw* modifies a verbal element whose SUBJ has the same index value, whereas the ACC-marked *num-cl-mw* modifies a verbal element which has at least one unsaturated COMPS element whose INDEX value is identical with its own INDEX value. What this means is that the NOM or ACC marked *num-cl-mw* is semantically linked to the SUBJ or COMPS element through the INDEX value.

Figure 1 is the parsing results for (6b) that our system yields. As seen from the parsed syntactic structure in Figure 1, the FQ *sey kay-lul* ‘three CL-ACC’ (NP-ACC) modifies the verbal expression *mek-ess-ta* ‘eat-PST-DECL’. However, as noted from the output MRS, this modifying FQ is linked with its antecedent *sakwa-lul* ‘apple-ACC’ through the relation *part-of_rel*. Leaving aside the irrelevant semantic relations, let’s see *card_rel* and *apple_rel*. As noted, the ARG0 value (x14) of *part-of_rel* is identified with that of *card_rel* whereas its ARG1 value (x4) is identified with the ARG0 value of the *apple_rel*. We thus can have the interpretation that there are three individuals x14s which belongs to the set x4.

3 Case Mismatches

Further complication in parsing FQ constructions comes from raising, causatives, and topicalization where the FQ and its antecedent have different case values. In such examples, the two need not have an identical case value. For example, as given in (8b), the ACC-marked raised object can function as the antecedent of either the NOM-marked or ACC-marked FQ:

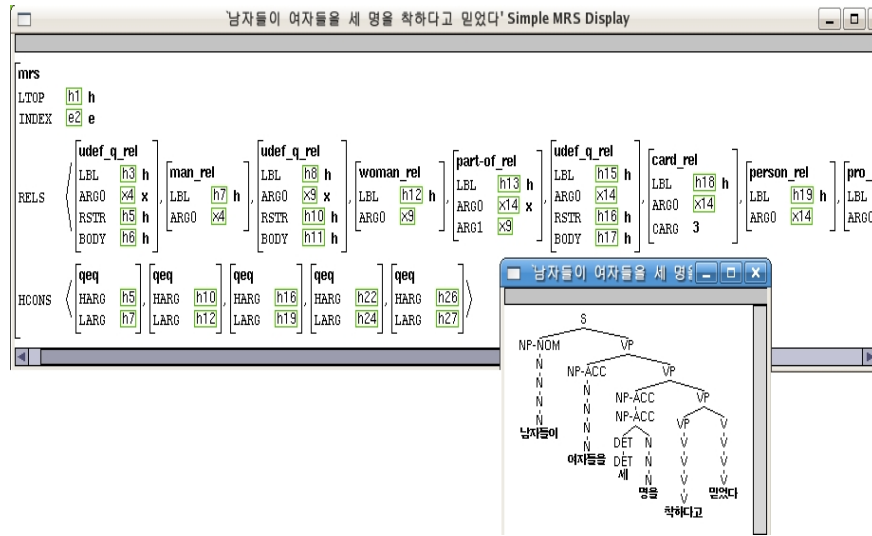


Figure 2: Parsed Tree and MRS for ‘As for the students, they ate three apples.’

- (8) a. namcatul-i [yecatul-i sey myeng-i/*ul chakhata-ko] mitessa.
 men-NOM women-NOM three-CL-NOM/*ACC honest-COMP thought
 ‘Men thought that three women are honest.’
 b. namcatul-i yecatul-ul sey myeng-ul chakhata-ko mitessa.
 c. namcatul-i yecatul-ul sey myeng-i chakhata-ko mitessa.

In the present analysis in which the case-marked FQ is linked to either the SUBJ or a COMPS element, we can expect these variations. Let us consider the lexical entry for the raising verb *mitessa* ‘believed’:

- (9) a.
$$\left[\begin{array}{l} \text{HEAD} \mid \text{POS } verb \\ \text{SUBJ } \langle 1 \text{NP} \rangle \\ \text{COMPS } \langle 2 \text{S} \rangle \\ \text{ARG-ST } \langle 1, 2 \rangle \end{array} \right]$$
 b.
$$\left[\begin{array}{l} \text{HEAD} \mid \text{POS } verb \\ \text{SUBJ } \langle 1 \text{NP} \rangle \\ \text{COMPS } \langle 2 \text{NP}_i, 3 \text{VP}[\text{SUBJ } \langle \text{NP}_i \rangle] \rangle \\ \text{ARG-ST } \langle 1, 2, 3 \rangle \end{array} \right]$$

(9a) represents the lexical entry for *mitessa* ‘believed’ in (8a) selecting a sentential complement. Meanwhile, (9b) represents the raising verb ‘believed’ in (8b, c) in which the subject of the embedded clause is raised as the object. This lexical element allows *yecatul-ul* ‘women-ACC’ to function as the syntactic object of the verb even though it is the semantic subject of the lower predicate.

Equipped with these, our grammar generates Figure 2 as the parsing results for (8b). Syntactically, as noted from the parsed structure, the ACC-marked FQ *sey myeng-ul* ‘three CL-ACC’ (NP-ACC) modifies the VP *chakhata-ko mitessa* ‘honest-COMP believed’.⁴ Meanwhile, semantically, the ACC-marked FQ is linked

⁴Our grammar allows only binary structures for the language. One strong advantage of assuming binary structures comes from scrambling facts. See Kim and Yang (2004).

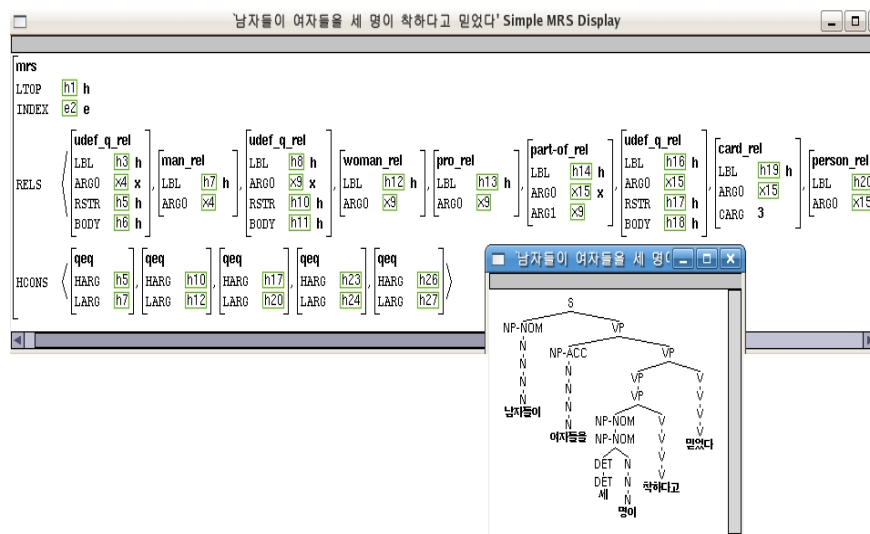


Figure 3: Parsed Tree and MRS for ‘Men (NOM) thought three (NOM) women (ACC) are honest.’

to the ACC-marked object *yecatul-ul* ‘woman-ACC’. This is because in our grammar the antecedent of the ACC-marked FQ must be an unsaturated complement of the VP it modifies. As noted from the semantic relations *part-of_rel*, *card_rel* and *woman_rel* in the parsed MRS, this linking relation is attested. That is, the ARG0 value (x9) of *woman_rel* is identified with the ARG1 value of *part-of_rel* whereas the ARG0 value of *card_rel* is identical with the ARG0 value of *part-of_rel*. Thus, the semantic output correctly indicates that the individuals denoted by the FQ is a subset of the individuals denoted by the antecedent.

For the mismatch example (8c), our grammar correctly produces two structures. Let's see Figure 3 first. As seen from the parsed syntactic structure here, the FQ *sey myeng-i* 'three CL-NOM' (NP-NOM) modifies the complex VP *chakhatako mitessta* 'honest-COMP believed'. However, in terms of semantics, the FQ is linked to the subject of the VP that it modifies.⁵ This linking relation is once again attested by the MRS structure here. As noted here, the two semantic arguments of *part-of_rel*, ARG0 and ARG1, have identical values with the ARG0 value of *card_rel* (x14) and *man_rel* (x4), respectively.

Meanwhile, as given in the second parsing result Figure 4, the FQ *sey myeng-i* ‘three CL-NOM’ modifies the simple VP *chakhata-ko* ‘honest-COMP’ only. Since the VP that the FQ modifies has only its SUBJ unsaturated, the SUBJ is the only possible antecedent. The output MRS reflects this raising property: The ARG0 value of *part-of_rel* identified with that of *card_rel* whereas its ARG1 value is identified with the ARG0 value of *woman_rel*. Our system thus correctly links the NOM-marked FQ with the ACC-marked antecedent even though they have different case values.

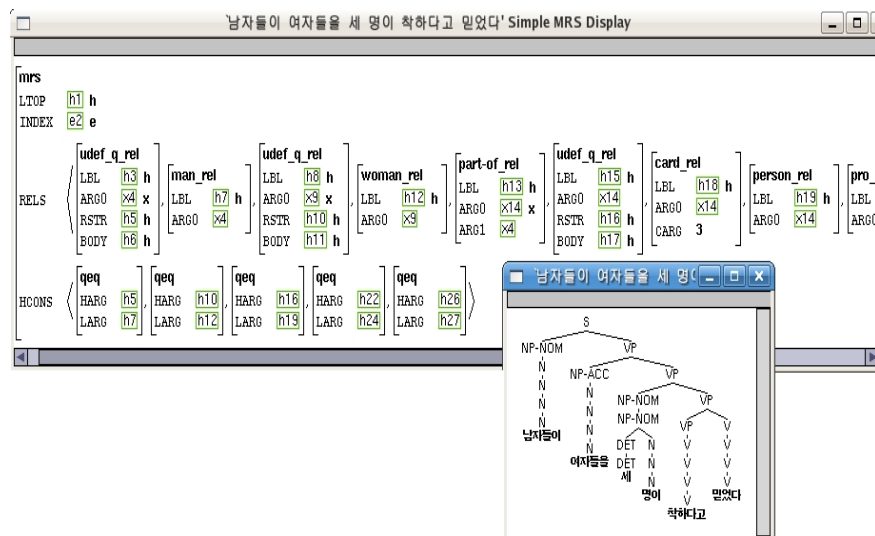


Figure 4: Parsed Tree and MRS for ‘Men (NOM) thought there are three (NOM) women (ACC) are honest.’

The grammar we have built within the typed-feature structure system and well-defined constraints, eventually aiming at working with real-world data, has been implemented in the HPSG for Korean (cf. Kim (2004), Kim and Yang (2004)). We have shown that the grammar can parse the appropriate syntactic and semantic aspects of the FQ constructions. The test results provide a promising indication that the grammar, built upon the typed feature structure system, is efficient enough to build semantic representations for the simple as well as complex FQ constructions.

References

- Bender, Emily M. and Melanie Siegel. 2004. Implementing the Syntax of Japanese Numeral Classifiers. *Proceedings of IJCNLP-04*, Hainan Island, China.
- Bond, Francis and Kyoung-Hee Paik. 2000. Reusing an ontology to generate numeral classifiers. In *Coling 2000*, Saarbrücken, Germany.
- Copestake, Ann. 2002. *Implementing Typed Feature Structure Grammars*. CSLI Publications.
- Copestake, Ann, Dan Flickinger, Carl Pollard, and Ivan A. Sag. 2006. Minimal Recursion Semantics: an Introduction. *Research on Language and Computation* 3.4: 281–332.
- Kang, Beom-mo. 2002. Categories and meanings of Korean floating quantifiers-with some reference to Japanese. *Journal of East Asian Linguistics* 11, 375–398.

Kim, Jong-Bok. 2004. *Korean Phrase Structure Grammar* (In Korean). Hankook Publishing.

Kim, Jong-Bok, Jaehyung Yang. 2004. Projections from Morphology to Syntax in the Korean Resource Grammar: Implementing Typed Feature Structures, *Lecture Notes in Computer Science*, Vol.2945, pp.13-24, Springer-Verlag, 2004.2.