

A Computational Implementation of Internally Headed Relative Clause Constructions

Jong-Bok Kim¹, Peter Sells², and Jaehyung Yang³

¹ School of English, Kyung Hee University, Seoul, Korea 130-701
jongbok@khu.ac.kr

² Dept. of Linguistics, Stanford University, USA
sells@stanford.edu

³ School of Computer Engineering, Kangnam University, Kyunggi, 446-702, Korea
jhyang@kangnam.ac.kr

Abstract. The so-called Internally Headed Relative Clause (IHRC) construction found in the head-final languages Korean and Japanese has received little attention from computational perspectives even though it is frequently found in both text and speech. This is partly because there have been no grammars precise enough to allow deep processing of the construction's syntactic and semantic properties. This paper shows that the typed feature structure grammar HPSG (together with the semantic representations of Minimal Recursion Semantics) offers a computationally feasible and useful way of deep-parsing the construction in question.

1 Introduction

In terms of truth conditions, there is no clear difference between a (Korean) IHRC (internally head relative clause) like (1)a and and EHRC (externally headed relative clause) like (1)b.¹

- (1) a. Tom-un [sakwa-ka cayngpan-wi-ey iss-nun kes]-ul mekessta
Tom-TOP apple-NOM tray-TOP-LOC exist-PNE KES-ACC ate
'Tom ate an apple, which was on the tray.'
b. Tom-un [___ cayngpan-wi-ey iss-nun sakwa]-ul mekessta.
Tom-TOP tray-TOP-LOC exist-PNE apple-ACC ate
'Tom ate an apple that was on the tray.'

Both describe an event in which an apple is on the tray, and Tom's eating it.²

Yet, there exist several intriguing differences between the two constructions. One crucial difference between the IHRC and EHRC comes from the fact that

¹ We thank anonymous reviewers for their helpful comments and suggestions. This work was supported by the Korea Research Foundation Grant funded by the Korean Government (KRF-2005-042-A00056).

² The following is the abbreviations used for glosses and feature attributes in this paper: ACC (ACCUSATIVE), COMP (COMPLEMENTIZER), LOC (LOCATIVE), NOM (NOMINATIVE), PNE (PRENOMINAL), TOP (TOPIC), etc.

the semantic object of *mekessta* ‘ate’ in the IHRC example (1)a is the NP *sakwa* ‘apple’ buried inside the embedded clause. It is thus the subject of the embedded clause that serves as the semantic argument of the main predicate ([1], [2]).

In the analysis of such IHRCs, the central questions thus involve (a) the key syntactic properties, (b) the association of the internal head of the IHRC clause with the matrix predicate so that the head can function as its semantic argument, and (c) the differences between the IHRC and EHRC. This paper provides a constraint-based analysis within the framework of HPSG (Head-driven Phrase Structure Grammar) and implements it in the existing HPSG grammar for Korean using the LKB (Linguistic Building Knowledge) system to check the computational feasibility of the proposed analysis.³

2 Implementing an Analysis

2.1 Syntactic Aspects of the IHRC

One main morphological property of the IHRC construction is shown in (2)b: the embedded clausal predicate should be in the adnominal present form of *(n)un*, followed by the so-called bound noun *kes*. This clearly contrasts with the EHRC example (2)a, in which the predicate can have any of the three different markers of tense information:⁴

- (2) a. Tom-i _i ilk-nun/un/ul chayk_i
 Tom-NOM read-PRES.PNE/PST.PNE/FUT.PNE book
 ‘the book that Tom reads/read/will read’
- b. Tom-un [sakwa-ka cayngpan-wi-ey **iss-nun/*ul** kes]-ul mekessta
 Tom-TOP apple-NOM tray-TOP-LOC exist-PNE KES-ACC ate
 ‘Tom ate an apple, which was (lit. ‘is’) on the tray.’

In traditional Korean grammar, *kes* in the IHRC is called a ‘dependent noun’, in that it always requires either a modifying determiner or clause, even in a non-IHRC usage:

- (3) a. *(i/ku/ce) kes ‘*(this/that) thing’
 b. *(nay-ka mek-un) kes ‘the thing (*that I ate)’

This close syntactic relation between the clause and the noun *kes* can also be found in the fact that unlike canonical nouns, it must combine with a preceding adnominal clause:

- (4) Na-nun *(kangto-ka unhayng-eyse nao-nun) kes-ul capassta
 I-TOP robber-NOM bank-from come-out-PNE KES-ACC caught
 ‘I arrested the robber who was coming out of the bank.’

³ The LKB, freely available with open source (<http://lingo.stanford.edu>), is a grammar and lexicon development environment for use with constraint-based linguistic formalisms such as HPSG. cf. [3].

⁴ These three pronominal markers in the EHRC extend their meanings to denote aspects when combined with (preceding) tense suffixes.

These examples show that the pronoun *kes* selects an adnominal clause as its complement, and that the IHRC requires a specific inflected form of its predicate.

Then, what is the relationship between the whole IHRC clause including *kes* and the matrix verb? To relate the matrix verb with this construction with an ‘internal semantic head’, it was assumed in transformational grammar that it was necessary to introduce an empty category such as *pro* to the right of the adnominal clause, on the assumption that the IHRC is an adjunct clause (Jhang 1991). However, there is ample evidence showing that the clause is a direct syntactic nominal complement of the matrix predicate. One strong argument against an adjunct treatment centers on the passivization of the IHRC clause. As shown in (5), an object IHRC clause can be promoted to the subject of the sentence.

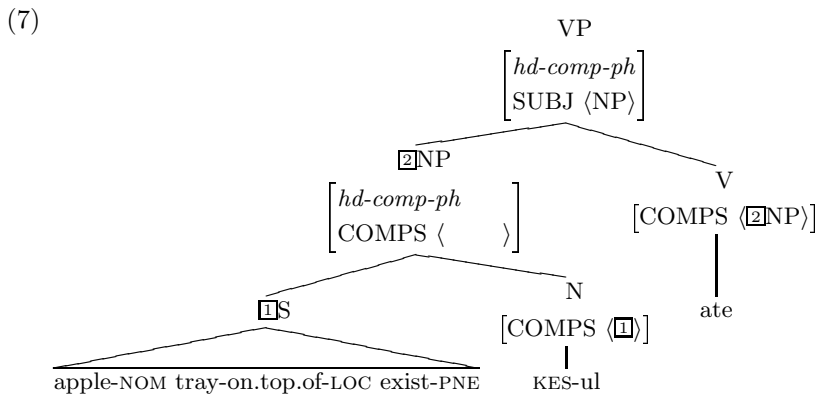
- (5) [Tom-i talli-nun kes]-i Mary-eyeuyhayse caphiessta
 Tom-NOM run-PNE KES-NOM Mary-by be.caught
 ‘Tom, who was running, was caught by Mary.’

Another fact concerning the status of the IHRC comes from stacking: whereas more than one EHRC clause can be stacked, only one IHRC clause is possible:

- (6) a. *kyongchal-i [kangto-ka unhayng-eyse nao-nun]
 police-NOM [robber-NOM bank-from come.out-PNE]
 [ton-ul hwumchi-n] kes-ul chephohayssta
 money-ACC steal-PNE KES-ACC arrested
 ‘(int.) The police arrested a thief coming out of the bank, stealing money.’
 b. kyongchal-i [__ unhayng-eyse nao-nun]
 police-NOM [bank-from come.out-PNE]
 [ton-ul hwumchi-n] kangto-lul chephohayssta
 money-ACC steal-PNE robber-ACC arrested
 ‘(int.) The police arrested a thief coming out of the bank, stealing money.’

This contrast implies that the adnominal clause which is the IHRC has the canonical properties of a complement clause.

Based on these observations, we assume the structure (7) for the internal and external structure of the IHRC in (1)a:



As represented in the tree, *kes* combines with its complement clause, forming a *hd-comp-ph* (*head-complement-ph*). This resulting NP also functions as the complement of the matrix verb *ate*.

2.2 Semantic Aspects of the IHRC and Related Constructions

One thing to note is that IHRCs are syntactically very similar to DPCs (direct perception constructions). IHRCs and DPCs both function as the syntactic argument of a matrix predicate. However, in the IHRC (8)a, the internal argument *John* within the embedded clause functions as the semantic argument of ‘caught’. Meanwhile, in (8)b it is the whole embedded clausal complement that functions as its semantic argument:

- (8) a. Mary-nun [John-i talli-nun kes]-ul **capassta**.
 Mary-TOP John-NOM run-PNE KES-ACC caught
 ‘Mary caught John who was running.’
- b. Mary-nun [John-i talli-nun kes]-ul **poassta**.
 Mary-TOP John-NOM run-PNE KES-ACC saw
 ‘Mary saw John running.’

The only difference between (8)a and (8)b is the matrix predicate, which correlates with the meaning difference. When the matrix predicate is an action verb such as *capta* ‘catch’, *chepohata* ‘arrest’, or *mekta* ‘eat’ as in (8)a, we obtain an entity reading for the clausal complement. But as in (8)b we will have only an event reading when the matrix predicate is a type of perception verb such as *po-ta* ‘see’, *al-ta* ‘know’, and *kiekhata* ‘remember’.

The key point in our analysis is thus that the interpretation of *kes* is dependent upon the type of matrix predicate. Hence the lexical entries in our grammar involve not only syntax but also semantics. For example, the verb *cap-ta* ‘catch’ in (9) lexically requires its object to refer to a *ref-ind* (referential-index) whereas the verb *po-ta* ‘see’ in (10) selects an object complement whose index is *indiv-ind* (individual index) whose subtypes include *ref-ind* and *event-ind*, indicating that its object can be either a referential individual or an event.⁵

⁵ The meaning representations adopted here involve Minimal Recursion Semantics (MRS), developed by [4]. This is a framework of computational semantics designed to enable semantic composition using only the unification of type feature structures. The value of the attribute SEM(ANTICS) we used here represents simplified MRS, though it originally includes HOOK, RELS, and HCONS. The feature HOOK represents externally visible attributes of the atomic predications in RELS (RELATIONS). The value of LTOP is the local top handle, the handle of the relations with the widest scope within the constituent. The value of XARG is linked to the external argument of the predicate. See [4] and [5] for the exact function(s) of each attribute. We suppress irrelevant features.

This grammar in which lexical information interacts with the other syntactic components ensures that the perception verb *saw* combines with an NP projected from (11)a whereas the action verb *caught* with an NP projected from (11)b. Otherwise, the resulting structure will not satisfy the selectional restrictions of the predicates.

Incorporating this into our Korean grammar,⁷ we implemented this analysis in the LKB and obtained the following two parsed trees and MRSs for the two examples:

The screenshot shows a window titled "존이 달리는 것을 잡았다" Simple MRS Display. The MRS data is as follows:

```

mrs
LTOP h1 h
INDEX e2 e
RELS
  [pro_rel] [named_rel] [proper_q_rel] [run_rel] [kes_rel] [prpstn_m_rel] [catch_rel]
  [LEL h3 h] [LEL h5 h] [LEL h7 h] [LEL h10 h] [LEL h12 h] [LEL h1 h] [LEL h14 h]
  [ARGO u4 u] [ARGO x8 x] [ARGO x6 h] [ARGO e11 e] [ARGO x6 h] [ARGO e2 e] [ARGO e15 e]
  [CARG john] [RSTR h8 h] [BODY h9 h] [ARG1 x6 h] [MARG h13 h] [ARG1 u4 u] [ARG2 x6 h]
HCONS
  [qeq] [qeq]
  [HARG h8] [HARG h13]
  [LARG h5] [LARG h14]
    
```

The parse tree shows a root S node branching into NP-ACC and VP. NP-ACC branches into S and N. The inner S branches into NP-NOM and V. NP-NOM branches into N (존이) and V (달리는). The outer V branches into N (것을) and V (잡았다).

The screenshot shows a window titled "존이 달리는 것을 보았다" Simple MRS Display. The MRS data is as follows:

```

mrs
LTOP h1 h
INDEX e2 e
RELS
  [pro_rel] [named_rel] [proper_q_rel] [run_rel] [kes_rel] [prpstn_m_rel] [see_rel]
  [LEL h3 h] [LEL h5 h] [LEL h7 h] [LEL h10 h] [LEL h12 h] [LEL h1 h] [LEL h14 h]
  [ARGO u4 u] [ARGO x8 x] [ARGO x6 h] [ARGO e11 e] [ARGO e11 h] [ARGO e2 e] [ARGO e15 e]
  [CARG john] [RSTR h8 h] [BODY h9 h] [ARG1 x6 h] [MARG h13 h] [ARG1 u4 u] [ARG2 e11 h]
HCONS
  [qeq] [qeq]
  [HARG h8] [HARG h13]
  [LARG h5] [LARG h14]
    
```

The parse tree is identical to the first screenshot, but the final verb is "보았다" (saw) instead of "잡았다" (caught).

⁷ The current Korean Resource Grammar has 394 type definitions, 36 grammar rules, 77 inflectional rules, 1100 lexical entries, and 2100 test-suite sentences, and aims to expand its coverage on real-life data.

Leaving aside the irrelevant parts, we can see that the two have the identical syntactic structures but different semantics. In the former, the ARG0 value of *kes* is identified with the *named_rel* (for ‘John’) but in the latter it is identified with *run_rel*.

The analysis thus provides a clean account of the complementary distribution of the IHRC and the DPC. That is, according to our analysis, we obtain an entity reading when the index value of *kes* is identified with that of the external argument. Meanwhile, we have an event reading when the index value is structure-shared with that of the adnominal S. This analysis thus correctly predicts that there exist no cases where the two readings are available simultaneously.

One of the welcome predictions that this analysis brings is that the canonical antecedent of the pronoun *kes* is the external argument:

- (12) [haksayng-i aktang-ul cha-nun kes-ul] capassta
 student-NOM rascal-ACC kick-PNE KES-ACC caught
 ‘(I) caught a student, who was then kicking a rascal.’

Even though one can catch either a student or a rascal, the semantic object of the verb ‘catch’ is not the object but the external argument *haksayng* (attested by our implementation but not included here because of limits on space).

3 Discussion and Conclusion

The analysis we have presented so far, part of the typed-feature structure grammar HPSG for Korean aiming at working with real-world data, has been implemented into LKB (Linguistic Knowledge Building System) to test its performance and feasibility.

We first inspected the Sejong Treebank Corpus (33,953 sentences) and identified 4,610 sentences with [S[FORM *nun*] + *kes*]. Of these, we inspected the 518 ACC marked examples, but found only 3 IHRC examples. Another 154 examples used *kes* in a cleft construction, and 361 as direct perception examples. Among these, we selected canonical types of the IHRC constructions to check if the grammar can parse them both in terms of syntax and semantics. As we have shown in section 2.2, the grammar is quite successful in picking up the appropriate semantic head from the IHRC. Of course, issues remain of extending the coverage of our grammar to parse more real-life data and further identifying other constructional types of *kes*, such as cleft usages.

Any grammar, aiming for real world application, needs to provide a correct syntax from which we can build semantic representations in compositional ways. In addition, these semantic representations must be rich enough to capture compositional as well as constructional meanings. In this respect, the analysis we have sketched here seems to be promising in the sense that it provides appropriate semantic representations for the IHRC and DPC in a compositional way, suitable for applications requiring deep natural language understanding.

References

1. Kim, Y.B.: Relevance in internally headed relative clauses in Korean. *Lingua* **112** (2002) 541–559
2. Chung, C., Kim, J.B.: Differences between externally and internally headed relative clause constructions. In: *Proceedings of HPSG 2002*, CSLI Publications (2003) 3–25
3. Copestake, A.: *Implementing Typed Feature Structure Grammars*. CSLI Publications, Stanford (2002)
4. Copestake, A., Flickenger, D., Sag, I., Pollard, C.: Minimal recursion semantics: An introduction. Manuscript (2003)
5. Bender, E.M., Flickinger, D.P., Oepen, S.: The grammar matrix: An open-source starter-kit for the rapid development of cross-linguistically consistent broad-coverage precision grammars. In Carroll, J., Oostdijk, N., Sutcliffe, R., eds.: *Proceedings of the Workshop on Grammar Engineering and Evaluation at the 19th International Conference on Computational Linguistics*, Taipei, Taiwan (2002) 8–14