

Implementing the Syntax of Japanese Numeral Classifiers

Emily M. Bender¹ and Melanie Siegel²

¹ University of Washington, Department of Linguistics,
Box 354340, Seattle WA 98195-4340
ebender@u.washington.edu

² Saarland University, Computational Linguistics,
PF 15 11 50, D-66041 Saarbrücken
siegel@dfki.uni-sb.de

Abstract. While the sortal constraints associated with Japanese numeral classifiers are well-studied, less attention has been paid to the details of their syntax. We describe an analysis implemented within a broad-coverage HPSG that handles an intricate set of numeral classifier construction types and compositionally relates each to an appropriate semantic representation, using Minimal Recursion Semantics.

1 Introduction

Much attention has been paid to the semantic aspects of Japanese numeral classifiers, in particular, the semantic constraints governing which classifiers co-occur with which nouns [1, 2]. Here, we focus on the syntax of numeral classifiers: How they combine with number names to create numeral classifier phrases, how they modify head nouns, and how they can occur as stand-alone NPs. We find that there is both broad similarity and differences in detail across different types of numeral classifiers in their syntactic and semantic behavior. We present semantic representations for two types and describe how they can be constructed compositionally in an implemented broad-coverage HPSG [3] for Japanese.

The grammar of Japanese in question is JACY¹, originally developed as part of the *Verbmobil* project [4] to handle spoken Japanese, and then extended to handle informal written Japanese (email text; [5]) and newspaper text. Recently, it has been adapted to be consistent with the LinGO Grammar Matrix [6].

2 Types of Numeral Classifiers

[7] divide Japanese numeral classifiers into five major classes: *sortal*, *event*, *mensural*, *group* and *taxonomic*, and several subclasses. The classes and subclasses can be differentiated according to the semantic relationship between the classifiers and the nouns they modify, on two levels: First, what properties of the

¹ <http://www.dfki.uni-sb.de/~siegel/grammar-download/JACY-grammar.html>

The above examples illustrate the contexts with a sortal numeral classifier, but mensural numeral classifiers can also appear both as modifiers (3a) and as NPs in their own right (3b):

- (3) a. ni kiro no ringo wo katta
 2 NumCl (kg) GEN apple ACC bought
 ‘(I) bought two kilograms of apples.’
 b. ni kiro wo katta
 2 NumCl (kg) ACC bought
 ‘(I) bought two kilograms.’

NumCIPs serving as NPs can also appear as modifiers of other nouns:

- (4) a. san nin no deai wa 80 nen haru
 3 NumCl GEN meeting TOP 80 year spring
 ‘The three’s meeting was in the spring of ’80.’

As a result, tokens following the syntactic pattern of (2b) and (3a) are systematically ambiguous, although the non-anaphoric reading tends to be preferred.

Certain mensural classifiers can be followed by the word *han* ‘half’:

- (5) ni kiro han
 two kg half
 ‘two and a half kilograms’

In order to build their semantic representations compositionally, we make the numeral classifier (here, *kiro*) the head of the whole expression. *Kiro* can then orchestrate the semantic composition of the two dependents as well as the composition of the whole expression with the noun it modifies (see §6 below).

4 Data: Distribution

We used ChaSen [9] to segment and tag 10,000 paragraphs of the Mainichi Shinbun 2002 corpus. Of the resulting 490,202 words, 11,515 (2.35%) were tagged as numeral classifiers. 4,543 of those were potentially time/date expressions, leaving 6,972 numeral classifiers, or 1.42% of the words. 203 orthographically distinct numeral classifiers occur in the corpus. The most frequent is *nin* (the numeral classifier for people) which occurs 1,675 times.

We sampled 100 sentences tagged as containing numeral classifiers to examine the distribution of the constructions outlined in §3. These sentences contained a total of 159 numeral classifier phrases and the vast majority (128) were stand-alone NPs. This contrasts with Downing’s study [8] of 500 examples from modern works of fiction and spoken texts, where most of the occurrences are not anaphoric. Furthermore, while our sample contains no examples of the floated variety, Downing’s contains 96. The discrepancy probably arises because Downing only included sortal numeral classifiers, and not any other type. Another possible contributing factor is the effect of genre. In future work we hope to study the distribution of both the types of classifiers and the constructions involving them in the Hinoki treebank [10].

5 Semantic Representations

One of our main goals in implementing a syntactic analysis of numeral classifiers is to compositionally construct semantic representations, and in particular, Minimal Recursion Semantics (MRS) representations [11, 12]. Abstracting away from handle constraints (the representation of scope), illocutionary force, tense/aspect, and the unexpressed subject, the representation we build for (2b,c) is as in (6).

(6) `_cat_n_rel(x)`, `undef_rel(x)`, `card_rel(x, "2")`, `_raise_v_rel(z,x)`

This can be read as follows: A relation of raising holds between z (the unexpressed subject), and x . x denotes a cat entity, and is bound by an underspecified quantifier (`undef_rel`) as there is no explicit determiner. x is also an argument of a `card_rel` (short for ‘cardinal.relation’), whose other argument is the constant value 2, meaning that there are in fact two cats being referred to.

For anaphoric numeral classifiers (2a), the representation contains an underspecified `noun_relation`, to be resolved in further processing.

(7) `noun_relation(x)`, `undef_rel(x)`, `card_rel(x, "2")`, `_raise_v_rel(z,x)`

Mensural classifiers have somewhat more elaborated semantic representations, which we treat as similar to English measure NPs [13]. On this analysis, the NumCIP denotes the extent of some dimension or property of the modified N. This dimension or property is represented with an underspecified relation (`unspec_adj_rel`), and a `degree_rel` relates the measured amount to the underspecified adjective relation. The underspecified adjective relation modifies the N in the usual way. This is illustrated in (8), the semantic representation for (3a).

(8) `_kilogram_n_rel(x)`, `undef_rel(x)`, `card_rel(x, "2")`,
`degree_rel(unspec_adj_rel, x)`, `unspec_adj_rel(y)`, `_apple_n_rel(y)`,
`undef_rel(y)`, `_buy_v_rel(z,y)`

When mensural NumCIPs are used anaphorically (3b), the element modified by the `_unspec_adj_rel` is an underspecified `noun_relation`, analogously to the case of sortal NumCIPs used anaphorically:

(9) `_kilogram_n_rel(x)`, `undef_rel(x)`, `card_rel(x, "2")`,
`degree_rel(unspec_adj_rel, x)`, `unspec_adj_rel(y)`, `noun_relation(y)`,
`undef_rel(y)`, `_buy_v_rel(z,y)`

6 Implementing an Analysis

Our analysis consists of: (1) a lexical type hierarchy cross-classifying numeral classifiers along three dimensions (Fig. 1), (2) a special lexical entry for *no* for linking NumCIPs with nouns, (3) a unary-branching phrase structure rules for promoting NumCIPs to nominal constituents.

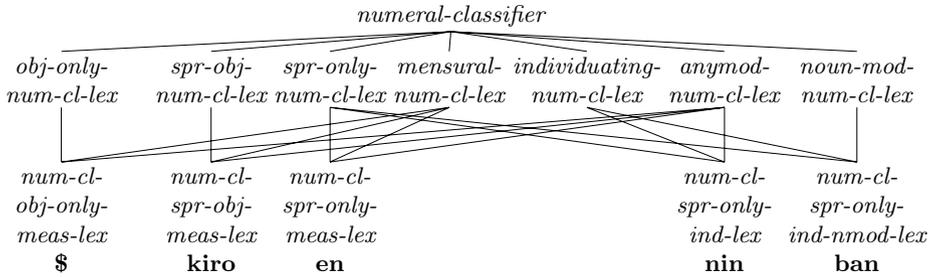


Fig. 1. Type hierarchy under *numeral-classifier*.

6.1 Lexical Types

Fig. 1 shows the lexical types for numeral classifiers, which are cross-classified along three dimensions: semantic relationship to the modified noun (*individuating* or *mensural*), modificational possibilities (NPs or PPs: *anymod*/NPs: *noun-mod*), and relationship to the number name (number name precedes: *spr-only*, number name precedes but may take *han*: *spr-obj*, number name follows: *obj-only*). Not all the possibilities in this space are instantiated (e.g., we have found no sortal classifiers which can take *han*), but we leave open the possibility that we may find in future work examples that fill in the range of possibilities. In this section, we treat each of the types in turn.

The constraint in (10) ensures that all numeral classifiers have the head type *num-cl_head*, as required by the unary phrase structure rule discussed in §6.3 below. Furthermore, it identifies two key pieces of semantic information made available for further composition, the INDEX and LTOP (local top handle) of the modified element, with the numeral classifier’s own INDEX and LTOP, as these are intersective modifiers [6]. The constraints on the type *num-cl_head* (not shown here) ensure that numeral classifiers can modify only saturated NPs or PPs (i.e., NPs marked with a case postposition *wo* or *ga*), and that they only combine via intersective head-modifier rules.

(10) *numeral-classifier* :=

$$\left[\begin{array}{l} \dots \text{CAT.HEAD} \\ \dots \text{CONT.HOOK} \end{array} \left[\begin{array}{l} \text{MOD} \left\langle \left[\begin{array}{l} \dots \text{INDEX} \\ \dots \text{LTOP} \end{array} \right] \begin{array}{l} \boxed{1} \\ \boxed{2} \end{array} \right\rangle \\ \text{INDEX} \quad \boxed{1} \\ \text{LTOP} \quad \quad \boxed{2} \end{array} \right. \right]$$

The constraints on the types *spr-only-num-cl-lex*, *obj-only-num-cl-lex* and *spr-obj-num-cl-lex* account for the position of the numeral classifier with respect to the number name and for the potential presence of *han*. Both the number name (a phrase of head type *int_head*) and *han* (given the distinguished head value *han_head*) are treated as dependents of the numeral classifier expression,

but variously as specifiers or complements according to the type of the numeral classifier. In the JACY grammar, specifiers immediately precede their heads, while complements are not required to do so and can even follow their heads (in rare cases). Given all this, in the ordinary case (*spr-only-num-cl-lex*), we treat the number name as the specifier of the numeral classifier. The other two cases involve numeral classifiers taking complements: with no specifier, in the case of pre-number unit expressions like the symbol \$ (*obj-only-num-cl-lex*) and both a number-name specifier and the complement *han* in the case of unit expressions appearing with *han* (*spr-obj-num-cl-lex*). Finally, the type *spr-obj-num-cl-lex* does some semantic work as well, providing the **plus_rel** which relates the value of the number name to the “ $\frac{1}{2}$ ” contributed by *han*, and identifying the ARG1 of the **plus_rel** with the XARG of the SPR and COMPS so that they will all share an index argument (eventually the index of the modified noun for sortal classifiers and of the measure noun relation for mensural classifiers).

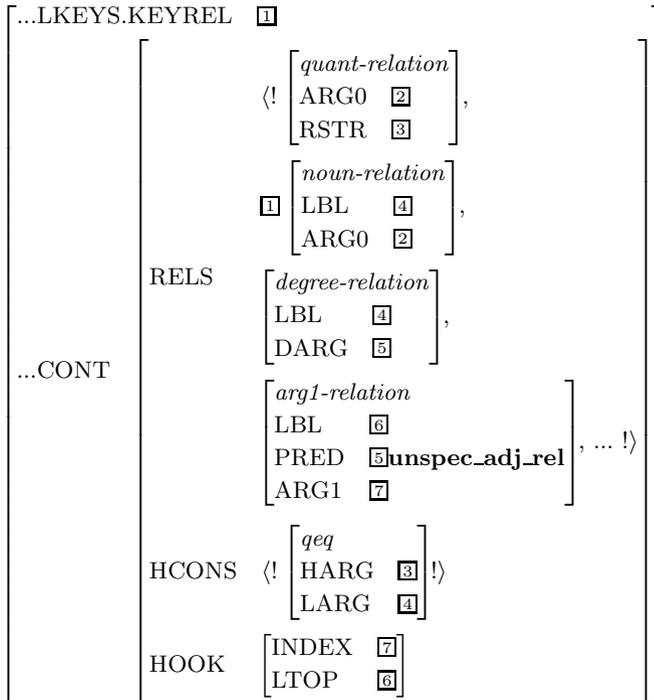
(11) *spr-obj-num-cl-lex* :=

$$\left[\begin{array}{l} \dots \text{VAL} \\ \dots \text{RELS} \end{array} \left[\begin{array}{l} \text{SUBJ } \textit{null} \\ \text{OBJ } \left[\begin{array}{l} \dots \text{CAT.HEAD } \textit{han_head} \\ \dots \text{CONT.HOOK } \left[\begin{array}{l} \text{LTOP } \boxed{1} \\ \text{XARG } \boxed{2} \end{array} \right] \end{array} \right] \\ \text{SPR } \left[\begin{array}{l} \dots \text{CAT.HEAD } \textit{int_head} \\ \dots \text{CONT.HOOK } \left[\begin{array}{l} \text{LTOP } \boxed{3} \\ \text{XARG } \boxed{2} \end{array} \right] \end{array} \right] \end{array} \right] \left[\begin{array}{l} \textit{plus-relation} \\ \text{ARG1 } \boxed{2} \\ \text{TERM1 } \boxed{3} \\ \text{TERM2 } \boxed{1} \end{array} \right] \! \} \right]$$

In the second dimension of the cross-classification, *anymod-num-cl-lex* and *noun-mod-num-cl-lex* constrain what the numeral classifier may modify, via the MOD value. Prenominal numeral classifiers are linked to the head noun with *no*, which mediates the modifier-modifiee relationship (see (2) and §6.2). However, numeral classifiers can appear after the noun (2c), modifying it directly. Some numeral classifiers can also ‘float’ outside the NP, either immediately after the case postposition or to the position before the verb (2d). While we leave the latter kind of float to future work (see §7), we handle the former by allowing most numeral classifiers to appear as post-head modifiers of PPs. Thus *noun-mod-num-cl-lex* further constrains the HEAD value of the element on the MOD list to be *noun_head*, but *anymod-num-cl-lex* leaves it as inherited (*noun-or-case-p_head*). This type does, however, constrain the modifier to show up after the head ([POSTHEAD *right*]), and further constrains the modified head to be [NUCL *nucl_plus*], in order to rule out vacuous attachment ambiguities between numeral classifiers attaching to the right left-attaching modifiers of the same NP.

The final dimension of the classification captures the semantic differences between sortal and mensural numeral classifiers. The sortal numeral classifiers contribute no semantic content of their own (represented with empty RELS and HCONS lists). In contrast, mensural numeral classifiers contribute quite a bit of semantic information, and therefore have quite rich RELS and HCONS values. As shown in (12), the *noun-relation* is identified with the lexical key relation value (LKEYS.KEYREL) so that specific lexical entries of this type can easily further specify it (e.g., *kīro* constrains its PRED to be **kilogram_n_rel**). The type also makes reference to the HOOK value so that the INDEX and LTOP (also the INDEX and LTOP of the modified noun, see (10)) can be identified with the appropriate values inside the RELS list. The length of the RELS list is left unbounded, because some mensural classifiers also inherit from *spr-obj-num-cl-lex*, and therefore must be able to add the **plus_rel** to the list.

(12) *mensural-num-cl-lex* :=



The types in the bottom part of the hierarchy in Fig. 1 join the dimensions of classification. They also do a little semantic work, making the INDEX and LTOP of the modified noun available to their number name argument, and, in the case of subtypes of *mensural-num-cl-lex*, they constrain the final length of the RELS list, as appropriate.

6.2 The Linker *no*

We posit a special lexical entry for *no* which mediates the relationship between NumCIPs and the nouns they modify. In addition to the constraints that it shares with other entries for *no* and other modifier-heading postpositions, this special *no* is subject to constraints that specify that *no* makes no semantic contribution, that it takes a NumCIP as a complement, and that the element on the MOD list of *no* shares its local top handle and index with the element on the MOD list of the NumCIP (i.e., that *no* effectively inherits its complement's MOD possibility). Even though (most) numeral classifiers can either modify NPs or PPs, all entries for *no* are independently constrained to only modify NPs, and only as pre-head modifiers.

6.3 Unary-Branching Phrase Structure Rule

We treat NumCIPs serving as nominal constituents by means of an exocentric unary-branching rule. This rule specifies that the mother is a noun subcategorized for a determiner specifier (these constraints are expressed on *noun_sc*), while the daughter is a numeral classifier phrase whose valence is saturated. Furthermore, it contributes (via its C-CONT, or constructional content feature) an underspecified *noun-relation* which serves as the thing (semantically) modified by the numeral classifier phrase. The reentrancies required to represent this modification are implemented via the LTOP and INDEX features.

(13) *nominal-numcl-rule-type* :=

$$\left[\begin{array}{l} \dots \text{CAT} \left[\begin{array}{l} \text{HEAD} \quad \textit{ordinary_noun_head} \\ \text{VAL} \quad \quad \textit{noun_sc} \end{array} \right] \\ \\ \text{C-CONT} \left[\begin{array}{l} \text{HOOK} \quad \left[\begin{array}{l} \text{LTOP} \quad \boxed{1} \\ \text{INDEX} \quad \boxed{2} \end{array} \right] \\ \text{RELS} \quad \langle ! \left[\begin{array}{l} \textit{noun-relation} \\ \text{LBL} \quad \boxed{1} \\ \text{ARG0} \quad \boxed{2} \end{array} \right] ! \rangle \end{array} \right] \\ \\ \text{ARGS} \left\langle \left[\begin{array}{l} \dots \text{CAT} \left[\begin{array}{l} \text{HEAD} \quad \textit{num-cl_head} \\ \text{VAL} \quad \quad \textit{saturated} \end{array} \right] \\ \dots \text{CONT.HOOK} \left[\begin{array}{l} \text{LTOP} \quad \boxed{1} \\ \text{INDEX} \quad \boxed{2} \end{array} \right] \end{array} \right] \right\rangle \end{array} \right]$$

This rule works for both sortal and mensural NumCIPs, as both are expecting to modify a noun.

7 Future Work

We have not yet implemented an analysis of pre-verbal floated NumCIPs, but we sketch one here. The key is that NumCIPs are treated as simple modifiers, not

quantifiers. Therefore, they can attach syntactically to the verb, but semantically to one of its arguments. In our HPSG analysis, the verb will have unsaturated valence features, making the indices of its arguments ‘visible’ to any modifiers attaching to it.

There appear to be constraints on which arguments can ‘launch’ floated quantifiers, although their exact nature is as yet unclear. Proposals include: only nominals marked with the case particles *ga* or *wo* [14], only subjects or direct objects [15], or c-command-based constraints [16]. While there are exceptions to all of these generalizations, [8] notes that the vast majority of actually occurring cases satisfy all of them, and further that it is primarily *intransitive* subjects which participate in the construction.

These observations will help considerably in reducing the ambiguity inherent in introducing an analysis of floated NumCIPs. We could constrain floated NumCIPs to only modify intransitive verbs (semantically modifying the subject) or transitive verbs (semantically modifying the object). Some ambiguity will remain, however, as the pre-verbal and post-nominal positions often coincide.

Also missing from our analysis are the sortal constraints imposed by classifiers on the nouns they modify. In future work, we hope to merge this analysis with an implementation of the sortal constraints, such as that of [2]. We believe that such a merger would be extremely useful: First, the sortal constraints could be used to narrow down the possible referents of anaphoric uses of NumCIPs. Second, sortal constraints could reduce ambiguity in NumCIP+*no*+N strings, whenever they could rule out the ordinary numeral classifier use, leaving the anaphoric interpretation (see (4) above). Third, sortal constraints will be crucial in generation [2]. Without them, we would propose an additional string for each sortal classifier whenever a **card_rel** appears in the input semantics, most of which would in fact be unacceptable. Implementing sortal constraints could be simpler for generation than for parsing, since we wouldn’t need to deal with varying inventories or metaphorical extensions.

8 Conclusion

Precision grammars require compositional semantics. We have described an approach to the syntax of Japanese numeral classifiers which allows us to build semantic representations for strings containing these prevalent elements — representations suitable for applications requiring natural language understanding, such as (semantic) machine translation and automated email response.

Acknowledgements

This research was carried out as part a joint R&D effort between YY Technologies and DFKI, and we are grateful to both for the opportunity. We would also like to thank Francis Bond, Dan Flickinger, Stephan Oepen, Atsuko Shimada and Tim Baldwin for helpful feedback in the process of developing and implementing this analysis and Setsuko Shirai for grammaticality judgments. This research was partly supported by the EU project DeepThought IST-2001-37836.

References

1. Matsumoto, Y.: Japanese numeral classifiers: A study of semantic categories and lexical organization. *Linguistics* **31** (1993) 667–713
2. Bond, F., Paik, K.H.: Reusing an ontology to generate numeral classifiers. In: *Coling 2000*, Saarbrücken, Germany (2000)
3. Pollard, C., Sag, I.A.: *Head-Driven Phrase Structure Grammar*. U of Chicago Press, Chicago (1994)
4. Siegel, M.: HPSG analysis of Japanese. In Wahlster, W., ed.: *Verbmobil: Foundations of Speech-to-Speech Translation*. Springer, Berlin (2000)
5. Siegel, M., Bender, E.M.: Efficient deep processing of Japanese. In: *Proceedings of the 3rd Workshop on Asian Language Resources and Standardization*, Coling 2002, Taipei (2002)
6. Bender, E.M., Flickinger, D., Oepen, S.: The Grammar Matrix: An open-source starter-kit for the rapid development of cross-linguistically consistent broad-coverage precision grammars. In: *Proceedings of the Workshop on Grammar Engineering and Evaluation*, Coling 2002, Taipei (2002) 8–14
7. Paik, K., Bond, F.: Spatial representation and shape classifiers in Japanese and Korean. In Beaver, D.I., Casillas Martínez, L.D., Clark, B.Z., Kaufmann, S., eds.: *The Construction of Meaning*. CSLI Publications, Stanford CA (2002) 163–180
8. Downing, P.: *Numeral Classifier Systems: The Case of Japanese*. John Benjamins, Philadelphia (1996)
9. Asahara, M., Matsumoto, Y.: Extended models and tools for high-performance part-of-speech tagger. In: *Coling 2000*, Saarbrücken, Germany (2000)
10. Bond, F., Fujita, S., Hashimoto, C., Kasahara, K., Nariyama, S., Nichols, E., Ohtani, A., Tanaka, T., Amano, S.: The Hinoki Treebank: A treebank for text understanding. In: *Proceedings of the IJC-NLP-2004*, Springer-Verlag (2004) this volume.
11. Copestake, A., Flickinger, D.P., Sag, I.A., Pollard, C.: Minimal Recursion Semantics. An introduction. Under review. (2003)
12. Copestake, A., Lascarides, A., Flickinger, D.: An algebra for semantic construction in constraint-based grammars. In: *ACL 2001*, Toulouse, France (2001)
13. Flickinger, D., Bond, F.: A two-rule analysis of measure noun phrases. In Müller, S., ed.: *Proceedings of the 10th International Conference on Head-Driven Phrase Structure Grammar*, Stanford CA, CSLI Publications (2003) 111–121
14. Shibatani, M.: *Nihongo no Bunseki*. Tasishuukan, Tokyo (1978)
15. Inoue, K.: *Nihongo no Bunpou Housoku*. Tasishuukan, Tokyo (1978)
16. Miyagawa, S.: *Structure and Case Marking in Japanese*. Academic Press, New York (1989)