# A Swedish Grammar Checker

Johan Carlberger Rickard Domeij Viggo Kann Ola Knutsson\* Royal Institute of Technology, Stockholm

This article describes the construction and performance of Granska – a surface-oriented system for grammar checking of Swedish text. With the use of carefully constructed error detection rules, the system can detect and suggest corrections for a number of grammatical errors in Swedish texts. In this article, we specifically focus on how erroneously split compounds and noun phrase disagreement are handled in the rules.

The system combines probabilistic and rule-based methods to achieve high efficiency and robustness. This is a necessary prerequisite for a grammar checker that will be used in real time in direct interaction with users. We hope to show that the Granska system with higher efficiency can achieve the same or better results than systems that use rule-based parsing alone.

Parts of this work were presented at Nodalida-99 (Domeij et al., 1999).

**Keywords**: grammar checking, part-of-speech tagging, word tagging, error detection rules, optimization, hidden Markov models.

# **1** Introduction

Grammar checking is one of the most widely used tools within language engineering. Spelling, grammar and style checking for English has been an integrated part of common word processors for some years now. However, most such programs are strictly commercial, and therefore there exists no documentation of the algorithms and rules used. An exception is the rule-based system by Vosse (Vosse, 1994).

For smaller languages, such as Swedish, advanced tools have been lacking. Recently, the first grammar checker for Swedish, developed by the Finnish company Lingsoft, was launched in Word 2000 (Arppe, 1999). This grammar checker is based on the Swedish constraint grammar SWECG.

In this article, another grammar checker for Swedish is presented. This grammar checker, called GRANSKA, has been developed at KTH for about four years. We will first present the structure of GRANSKA, and then in more detail describe four important parts of the system: the part-of-speech tagging module, the construction of error detection rules, the algorithms for rule matching, and the generation of error corrections. Finally, we describe the performance of the tagging, error detection and NP recognition.

# 2 The Granska System

GRANSKA is a hybrid system that uses surface grammar rules to check grammatical constructions in Swedish. The system combines probabilistic and rule-based methods to achieve high efficiency and robustness. This is a necessary prerequisite for a grammar checker that runs in real time in direct interaction with users (Kukich, 1992). Using so called error rules, the system

<sup>\*</sup> Nada, Department of Numerical Analysis and Computer Science, KTH, SE-100 44 Stockholm. E-mail: {jfc,domeij,viggo,knutsson}@nada.kth.se.



#### Figure 1

An overview of the GRANSKA system.

can detect a number of Swedish grammar problems and suggest corrections for them.

In Figure 1, the modular structure of the system is presented. First, in the tokenizer, potential words and special characters are recognized as such. In the next step, a tagger is used to assign disambiguated part of speech and inflectional form information to each word. The tagged text is then sent to the error rule matching component where error rules are matched with the text in order to search for specified grammatical problems. The error rule component also generates error corrections and instructional information about detected problems that are presented to the user in a graphical interface. The system also contains a spelling detection and correction module that can handle Swedish compounds (Domeij, Hollman, and Kann, 1994; Kann et al., 1999). The spelling detection module can also be used from inside the error rules, e.g. for checking split compound errors.

The GRANSKA system is implemented in C++ under UNIX, and it can be tested using a simple web interface<sup>1</sup>. There is ongoing work for designing two graphical interfaces for Windows, an add-in in Word and a stand-alone grammar checking editor that can be used interactively during writing. The system will be used as a research tool for studying usability aspects with real users.

## **3 Part-of-Speech Tagging**

In POS tagging of a text, each word and punctuation mark in the text is assigned a morphosyntactic tag. We have designed and implemented a tagger based on a second order Hidden Markov Model (Charniak et al., 1993; Charniak, 1996). Given a sequence of words  $w_{1..n}$ , the model finds the most probable sequence of tags  $t_{1..n}$  according to the equation

$$\mathcal{T}(w_{1..n}) = \arg\max_{t_{1..n}} \prod_{i=1}^{n} P(t_i | t_{i-2}, t_{i-1}) P(w_i | t_i).$$
(1)

<sup>1</sup> See the web page of the project: http://www.nada.kth.se/theory/projects/granska/

Estimations of the two probabilities in this equation are based on the interpolation of relative counts of sequences of 1, 2 and 3 tags and word-tag pairs extracted from a large tagged corpus.

For unknown words, we use a statistical morphological analysis adequate for Swedish and other moderately inflecting languages. This analysis is based on relative counts of observed tags for word types ending with the same 1 to 5 letters. This captures both inflections (tense *-ade* in *hämtade* (fetched)) and derivations (nounification *-ning* in *hämtning* (pick-up)).

We also perform an analysis that finds the last word form of compounds, which are common in Swedish. The possible tags of the last word form indicate possible tags (and probability estimation) for an unknown compound word. These two analyses are heuristically combined to get estimations of  $P(w_i|t_i)$ , which enables unknown words to work in the model. This method combines morphological information for unknown words with contextual information of surrounding words, and resulted in a tagger that tags 98 % of known and 93 % of unknown words correctly (Carlberger and Kann, 1999).

## **4 Error Rules**

The error rule component uses carefully constructed error rules to process the tagged text in search for grammatical errors. Since the Markov model also disambiguates and tags morphosyntactically deviant words with only one tag, there is normally no need for further disambiguation in the error rules in order to detect an error. An example of an agreement error is **en** *litet hus* (a small house), where the determiner *en* (a) does not agree with the adjective *liten* (small) and the noun *hus* (house) in gender. The strategy differs from most rule-based systems which often use a complete grammar in combination with relaxation techniques to detect morphosyntactical deviations (see for example (Vosse, 1994; Sågvall Hein, 1998)).

The error rules of GRANSKA are expressed in a rule language that we developed especially for this project. It is partly object-oriented and has a syntax resembling Java or C++. An error rule in GRANSKA that can detect the agreement error in *en liten hus*, is shown in Rule 1 below.

#### Rule 1

```
cong22@incongruence {
   X(wordcl=dt),
   Y(wordcl=jj)*,
   Z(wordcl=nn & (gender!=X.gender | num!=X.num | spec!=X.spec))
-->
   mark(X Y Z)
   corr(X.form(gender:=Z.gender, num:=Z.num, spec:=Z.spec))
   info("The determiner" X.text "does not agree with the noun" Z.text)
   action(scrutinizing)
}
```

Rule 1 has two parts separated with an arrow. The first part contains a matching condition. The second part specifies the action that is triggered when the matching condition is fulfilled. In the example, the action is triggered when a determiner is found followed by a noun (optionally preceded by one or more (\*) attributes) that differs (!=) in gender, number or (|) species from the determiner. Each line in the first part contains an expression that must evaluate to true in a matching rule. This expression may be a general Java expression and may refer to values (matching texts, word classes, or features) of the earlier parts of the rule.

The action part of the rule first (in the mark statement) specifies that the erroneous phrase should be marked in the text. Then (in the corr statement) a function is used to generate a new inflection of the article from the lexicon, one that agrees with the noun. This correction suggestion (in the example **ett** *litet hus*) is presented to the user together with a diagnostic comment (in the info statement) describing the error.

In most cases, the tagger succeeds in choosing the correct tag for the deviant word on probabilistic grounds (in the example *en* is correctly analyzed as an indefinite, singular and common gender article by the tagger). However, since errors are statistically rare compared to grammatical constructions, the tagger can sometimes choose the wrong tag for a morphosyntactically deviant form. In such cases, when the tagger is known to make mistakes, the error rules can be used in re-tagging the sentence to correct the tagging mistake. Thus, a combination of probabilistic and rule-based methods is used even during basic word tag disambiguation.

We use *help rules* (sub-routines) to define phrase types that can be used as context conditions in the error rules. In rule 2 below, two help rules are used in detecting agreement errors in predicative position. The help rules specify that copula should be preceded by an NP followed by one or more (+) PPs.

#### Rule 2

```
pred2@predicative {
   T(wordcl!=pp),
   (NP)(),
   (PP)()+
   X(wordcl=vb & vbt=kop),
   Y(wordcl=jj & (gender!=NP.gender | num!=NP.num)),
   Z(wordcl!=jj & wordcl!=nn)
-->
   mark(all)
   corr(if NP.spec=def then
        Y.form(gender:=NP.gender, num:=NP.num, spec:=ind) else
        Y.form(gender:=NP.gender, num:=NP.num) end)
   info("The noun phrase" NP.text "does not agree with the adjective" Y.text)
   action(scrutinizing)
}
NP@ {
   X(wordcl=dt)?,
   Y(wordcl=jj)*,
   Z(wordcl=nn)
-->
   action(help, gender:=Z.gender, num:=Z.num, spec:=Z.spec, case:=Z.case)
}
PP@ {
   X(wordcl=pp),
   (NP)()
-->
   action(help, gender:=NP.gender, num:=NP.num, spec:=NP.spec, case:=NP.case)
}
```

The help rules in the example are specified in the subroutines NP@ and PP@ which define noun phrases and prepositional phrases respectively. These subroutines are called from the higher level rule for predicative agreement (pred2@predicative). Note that the help rule PP@ uses the other help rule NP@ to define the prepositional phrase.

The help rules make the analysis approach that of a phrase structure grammar. Help rules make it possible for the system to perform a local phrase analysis selectively, without parsing other parts of the sentence that are not needed in the detection of the targeted error type. Thus, by calibrating the level of analysis that is needed for the case at hand, the system obtains high efficiency.

Above, we have shown how agreement errors are handled in the system. Another frequently occurring error type is erroneously split compounds. In contrast to English, a Swedish compound is regularly formed as one word, so split compounds are regarded as ungrammatical. So far, we

have mainly focussed on erroneously split compounds of the type noun+noun which stands for about 70 % of the various types (Domeij, Knutsson, and Öhrman, 1999).

Detection of erroneously split compounds where the first part cannot stand alone is trivial. This is done by listing those first parts in the lexicon and classifying them so that an error rule can be made to search the text for such a first part in combination with any other noun, as for example in *pojk byxor* where *pojk* is the first part form of *pojke* (boy) which is combined with *byxor* (trousers).

In other cases, when both parts have the form of full words, the strategy for detecting erroneously split compounds makes use of the fact that the first noun, unlike the last, must be uninflected (indefinite and singular). Since the combination uninflected noun followed by any noun is an unusual syntactical combination in grammatically correct sentences, it can be used to find candidates for split compound errors. Other contextual cues are also used before checking the candidate against a spell checker for compound recognition. If the spell checker recognizes the compound as such, the two nouns in the text are marked as a split compound and the corrected compound is given as a suggestion alternative.

Many errors can be difficult to detect because of ambiguities that are irresolvable on contextual grounds only. One example is *en exekverings enhet* (an execution unit). The first noun *exekvering* belongs to a group of nouns that take an *-s* when compounded with another noun *(exekvering-s+enhet)*. When the compound is erroneously split, the form of the first noun coincides with the genitive form (an execution's unit), which has the same syntactical distribution as the error construction and therefore cannot be distinguished from the split compound case.

#### 5 Rule Matching

Presently, there are about 200 scrutinizing rules and 50 help rules in GRANSKA. Each rule may be matched at any position (i.e. word) in the text, and there may even exist several matchings of a rule with the same starting position and of different length. The rule matcher tries to match rules from left to right, evaluating the expression of each token in the left hand side of the rule, and stopping as soon it finds out that the rule cannot be matched.

The rule language allows the operators \* (zero or more), + (one or more) and ? (zero or one) for tokens, and *i* (or) between rules. Together with the possibility of writing help rules, this makes the rule language a general regular language. Unfortunately, it is not possible to construct a standard linear time regular expression analyzer (Aho, Sethi, and Ullman, 1986), because the expressions in the rules may (and often do) use values of other tokens. However, the structure of the rules is in practice not very complicated. Therefore, a simple recursive matching algorithm is used and works well.

It is inefficient to try to match each error rule at each position in the text. We therefore perform a statistical optimization, where each rule is analyzed in advance. For each position in the rule, the possible matching words and taggings are computed. In fact, the possible tag bigrams for each pair of positions are computed. Then, using statistics on word and tag bigram relative frequencies, the position of the rule that is least probable to match a Swedish text is determined. This means that this rule is checked by the matcher *only* at the positions in the text where the words or tag bigrams of this least probable position in the rule occur. For example, a noun phrase disagreement rule may require a plural adjective followed by a singular noun in order to match. Such tag combinations are rare, and with this optimization approach, only the small portion of word sequences in a text containing this tag combination will be inspected by this rule.

It is important to note that this optimization is fully automatic. The program itself detects the optimal positions in each rule and stores two tables representing this information on disk. The first table describes, for each tag bigram, which rules that should be checked when that tag bigram occurs in the text. The second table contains the words appearing in the rules and describes, for each word, which rules that should be checked when that text. With the current set of error rules in GRANSKA the rule matching performs six times faster with optimization than without. Furthermore, due to the optimization, it is almost free (with respect to performance) to add many new rules as long as they contain some uncommon word or tag bigram.

## **6** Lexicons and Word Form Generation

The lexicon of the system is derived from the tagged Stockholm-Umeå Corpus (SUC) in addition to morphological information from various sources.

The grammar rules require the functionality to generate alternate inflection forms of any given word. Instead of having a lexicon containing all more or less common forms of each base form word, we use inflection rules to derive word forms from a base form. This approach has two advantages. Firstly, all inflectional forms of a word can be derived as long as its base form is known, and thus a smaller lexicon can be used. Secondly, unknown compound words can inherit the inflection rule of its last word form constituent, which enables corrections of unknown compound words.

## 7 Evaluation and Ranking of Error Corrections

It is often the case that an error rule matching generates more than one correction alternative. There are several reasons for this: different syntactic features may be applicable when a word form is changed, a base form may have more than one applicable inflection rule, and an error rule may have more than one correction field. These alternative sentences are first scrutinized and then ranked before being suggested to the user.

As the error rules are applied locally and not to an entire clause, sentence or paragraph, there will inevitably be false alarms. Therefore, each corrected sentence generated from an error rule matching is scrutinized with all other error rules in order to determine if another error was introduced. In such cases, the correction alternative is discarded.

If one of the correction alternatives is identical to the original sentence, it indicates not that the original sentence was erroneous, but that it was incorrectly tagged. For example, the noun *verktyg* (tool) has the same spelling in singular and plural. If the tagger tags *verktyg* as a plural noun in *Ett mycket bra verktyg* (A very good tool), a noun phrase disagreement error rule will correct the phrase to *Ett mycket bra verktyg*, where the only difference is the tag of the last word. Thus, when a corrected sentence identical to the original sentence is generated, the entire error matching is regarded as a false alarm.

These two approaches of discarding correction alternatives have indeed shown to increase precision more than they decrease recall.

There is another benefit from scrutinizing the sentences generated from error rules. The probability given by the tagging equation is a suitable measure for ranking these sentences, so that the sentence with most "common" words and syntactic structure is given as first alternative. We believe that it is important for a spell and grammar checker to suggest reasonable corrections. A spell or grammar checker that suggests a non-grammatical correction will lose in confidence from the user. The notions of trust and credibility have received increased attention in recent research about human-computer interaction. It applies not only to language support systems, but to all systems providing information and services to a human user. A recent overview is presented in (Fogg and Tseng, 1999).

If a sentence has a great proportion of unknown words, it makes little sense to apply grammar and spell checking rules to it, since it is probably a non-Swedish sentence. Instead, such a sentence is either ignored, marked as suspect in its entirety, or scrutinized anyway, according to the user's preference.

## **8** Results and Further Work

The tagging module has a processing speed of more than 20 000 words per second on a SUN Sparc station Ultra 5. In a previously unseen text, 97 % of the words are correctly tagged, a good result in an international comparing. Unknown words are correctly tagged in 93 % of the cases. The whole system (with about 20 rule categories of about 250 error rules) processes about 3 500 words per second, tagging included. The numbers are hard to compare to those of other systems, since they are seldom reported, but we believe that we have achieved a comparably high performance.

We are still working to optimize the system and improve the error rules. Preliminary tests with the error rules show that we can hope for a recall rate over 50 % and a precision rate over 90 % for agreement errors and erroneously split compounds. Still, it is unrealistic to hope for full recall and precision. Therefore, we think it is important to develop a user friendly and instructive graphical interface and test the program on users in practice to study usability aspects as well as the effects on writing and writing ability.

We have not yet performed a full comparison between the recall and precision of GRANSKA and Lingsoft's Swedish grammar checker Grammatifix (Arppe, 1999; Birn, 1999). Lingsoft has put much more effort in improving the precision of the system than we have had possibility to do yet. But, on the other hand, GRANSKA may detect split compounds, something that is not done by Grammatifix.

Most of the help rules of GRANSKA are trying to recognize different kinds of Swedish noun phrases. An evaluation of the first tentative NP recognition rules (Johansson, 2000) found both precision and recall to be about 80 %. The evaluation task was to recognize the best, comparing to human annotation, NP candidates without post attributes, in a text. This is a bit different and more difficult task than the ordinary usage of the help rules in GRANSKA. However, we have found that by refining and extending the help rules, we increase both recall and precision of grammar checking. We also hope that help rules for clause boundary recognition could even more increase recall and precision, especially for split compounds. Therefore, we have experimented with rules based on Ejerhed's clause boundary recognition algorithm (Ejerhed, 1999). By applying the error rules for split compounds only on for example clauses without ditransitive verbs, GRANSKA can avoid false alarms and still detect errors in another clause within the same sentence. The preliminary results so far are promising.

# Acknowledgments

The work has been funded by the Swedish research councils TFR, HSFR and Nutek. Project leader of the project has been Prof. Kerstin Severinson-Eklundh.

Språkdata at Göteborg University and the Swedish Academy let us use Svenska Akademiens ordlista as a source for words in GRANSKA. Prof. Eva Ejerhed and Prof. Gunnel Källgren let us use SUC.

#### References

- Aho, A., R. Sethi, and J. Ullman. 1986. Compilers: Principles, Techniques and Tools. Addison-Wesley, Reading, Mass.
- Arppe, A. 1999. Developing a grammar checker for Swedish. In Proc. 12th Nordic Conference in Computational Linguistics, Nodalida-99.
- Birn, J. 1999. Detecting grammar errors with Lingsoft's Swedish grammar checker. In Proc. 12th Nordic Conference in Computational Linguistics, Nodalida-99.

Carlberger, J. and V. Kann. 1999. Implementing an efficient part-of-speech tagger. *Software–Practice and Experience*, 29(9):815–832.

- Charniak, E. 1996. *Statistical language learning*. MIT Press, Cambridge, Massachusetts.
- Charniak, E., C. Hendrickson, N. Jacobson, and M. Perkowitz. 1993. Equations for part-of-speech tagging. In *11th National Conf. Artificial Intelligence*, pages 784–789.

Domeij, R., J. Hollman, and V. Kann. 1994. Detection of spelling errors in Swedish not using a word list en clair. J. Quantitative Linguistics, 1:195–201.

- Domeij, R., O. Knutsson, J. Carlberger, and V. Kann. 1999. Granska – an efficient hybrid system for Swedish grammar checking. In Proc. 12th Nordic Conference in Computational Linguistics, Nodalida-99.
- Domeij, R., O. Knutsson, and L. Öhrman. 1999. Inkongruens och felaktigt särskrivna sammansättningar – en beskrivning av två feltyper och möjligheten att detektera felen automatiskt (Incongruence and erroneously split compounds), in Swedish. In *Svenskans beskrivning-99*.
- Ejerhed, E. 1999. Finite state segmentation of discourse into clauses. In A. Kornai, editor, *Extended Finite State Models of Language*. Cambridge University Press, chapter 13.
- Fogg, B. J. and H. Tseng. 1999. The elements of computer credibility. In CHI-99, Human

*Factors in Computing Systems*, pages 80–87, Pittsburgh, PA. ACM Press.

- Johansson, V. 2000. NP-detektion, utvärdering och förslag till förbättringar av Granskas NP-regler (NP detection, evaluation and proposals of improvements of the NP rules of Granska), in Swedish. Technical report, Department of Linguistics, Stockholm University, Stockholm.
- Kann, V., R. Domeij, J. Hollman, and M. Tillenius. 1999. Implementation aspects and applications of a spelling correction algorithm. In R. Koehler, L. Uhlirova, and G. Wimmer, editors, *Text as a Linguistic Paradigm: Levels, Constituents, Constructs. Festschrift in honour of Ludek Hrebicek.* Universitat Verlag, Trier, Germany. To appear.
- Kukich, K. 1992. Techniques for automatically correcting words in text. ACM Computing Surveys, 24(4):377–439.
- Sågvall Hein, A. 1998. A chart-based framework for grammar checking. In Proc. 11th Nordic Conference in Computational Linguistics, Nodalida-98.
- Vosse, T. 1994. The Word Connection. Grammar-Based Spelling Error Correction in Dutch. Enschede: Neslia Paniculata. ISBN 90-75296-01-0.