# Embedded Controlled Languages

Aarne Ranta

Department of Computer Science and Engineering
Chalmers University of Technology and University of Gothenburg

**Abstract.** Inspired by embedded programming languages, an embedded CNL (controlled natural language) is a proper fragment of an entire natural language (its host language), but it has a parser that recognizes the entire host language. This makes it possible to process out-of-CNL input and give useful feedback to users, instead of just reporting syntax errors. This extended abstract explains the main concepts of embedded CNL implementation in GF (Grammatical Framework), with examples from machine translation and some other ongoing work.

**Keywords**: controlled language, domain-specific language, embedded language, Grammatical Framework, machine translation

## 1 Introduction

A controlled natural language (CNL) is a strictly defined fragment of a natural language [1]. As fragments of natural languages, CNLs are analogous to **embedded domain-specific languages**, which are fragments of general purpose programming languages [2]. Such languages have been introduced as an alternative to traditional **domain-specific languages** (DSL), which have their own syntax and semantics, and require therefore a specific learning effort. An embedded DSL is a part of a general-purpose programming language, the **host language**, and is therefore readily usable by programmers who already know the host language. At its simplest, an embedded DSL is little more than a **library** in the host language. Using the library helps programmers to write compact, efficient, and correct code in the intended domain. But whenever the library does not provide all functionalities wanted, the programmer can leave its straight-jacket and use the host language directly, of course at her own risk.

Embedding a language fragment in the full language presupposes that a grammar of the full language is available. In the case of natural languages, this is by no means a trivial matter. On the contrary, it is widely acknowledged that "all grammars leak", which means that any formal grammar defining a natural language is bound to be either incomplete or overgenerating. As a consequence, defining CNLs *formally* as subsets of natural languages looks problematic.

However, *if* a grammar of the host language exists, then it is useful to define the CNL as an embedded language. It enables us to build systems that provide, at the same time, the rigour of controlled languages and the comfort of graceful degradation. The user of the system can be guided to stay inside the controlled language, but she will also be understood, at least to some extent, if she goes outside it.

In this extended abstract, we will outline some recent work on building controlled languages in the embedded fashion. Our focus will be on multilingual systems, where the CNL yields high-quality translation and the host language yields browsing quality. But the same structure should be useful for other applications of CNLs as well, such as query languages.

In Section 2, we will summarize how CNLs are traditionally defined by using GF, Grammatical Framework. In Section 3, we will show how they can be converted to embedded CNLs. In Section 4, we summarize some on-going work and suggest some more applications.

## 2   Defining Controlled Languages in GF

GF [3] is a grammar formalism based on a distinction between **abstract syntax** and **concrete syntax**. The abstract syntax is a system of **trees**. The concrete syntax is a reversible mapping from trees to **strings** and **records**, reminiscent of **feature structures** in unification-based grammar formalisms. The separation between abstract and concrete syntax makes GF grammars **multilingual**, since one and the same abstract syntax can be equipped with many concrete syntaxes. The abstract syntax is then usable as an **interlingua**, which enables **translation** via **parsing** the source language string into a tree followed by the **linearization** of the tree into the target language.

As an example, consider a predicate expressing the age of a person. The abstract syntax rule is

```
fun aged : Person -> Numeral -> Fact
```

defining a **function** with the name `aged`, whose **type** is a function type of two arguments, of types `Person` and `Numeral`, yielding a value of type `Fact`. A simple concrete syntax rule for English is

```
lin aged p n = p ++ "is" ++ n ++ "years old"
```

stating that a function application (`aged p n`) is **linearized** to the string where the linearization of `p` is concatenated (`++`) with the string `"is"`, the linearization of `n`, and the string `"years old"`. The corresponding rule for French is

```
lin aged p n = p ++ "a" ++ n ++ "ans"
```

producing sentences literally equivalent to *p has n years*. Thus the concrete syntax allows the production of different syntactic structures in different languages, while the abstract syntax form (`aged p n`) remains the same, and stands in a compositional relation to the linearizations.

GF is widely used for defining CNLs; [4,5,6,7,8] are some examples. Much of its power comes from the **resource grammar libraries** (RGL), which are general purpose grammars enabling GF grammar writers to delegate the "low-level" linguistic details to generic library code [9,10]. If we dig deeper into the concrete syntax of the `aged` predicate, we will find lots of such details to cope with: number agreement (*one year*

vs. *five years*), subject-verb agreement: (*I am*, *you are*, *she is*), word order (*you are* in declaratives vs. *are you* in questions), etc; French generally poses harder problems than English. The use of feature structures instead of plain strings does make it possible to express all this compositionally in GF. But these details can make the grammar prohibitively difficult to write, especially since controlled language designers are not always theoretical linguists but experts in the various domains of application of CNLs. The RGL addresses this problem by providing general-purpose functions such as

```
mkCl : NP -> VP -> Cl
```

which builds a clause (`Cl`) from a noun phrase (`NP`) and a verb phrase (`VP`) and takes care of all details of agreement and word order. The CNL linearization rules can be written as combinations of such functions. The English rule, in full generality, comes out as compact as

```
lin aged p n = mkCl p (mkVP (mkAP (mkNP n year_N) old_A))
```

and the French,

```
lin aged p n = mkCl p (mkVP avoir_V2 (mkNP n an_N))
```

If a function `quest` is added to turn facts into questions, the linearization is in both languages just a simple RGL function call:

```
lin quest fact = mkQS fact
```

The API (application programmer's interface) is the same for all languages in the RGL (currently 29), but the low-level details that it hides are language-dependent.

## 3    Embedding a Controlled Language in the Host Language

The standard practice in GF is to define CNLs by using the RGL rather than low-level hand-crafted linearization rules. In addition to saving effort, this practice guarantees that the CNL is a valid fragment of a natural language. The reason is that the RGL is designed only to allow grammatically correct constructions.

But what is missing in a CNL built in this way is the rest of the host language—the part that is not covered by the RGL rule combinations actually used in the CNL. In general, a random sentence has a probability close to zero to be recognized by the parser. The standard solution is to guide the user by a predictive editor [11,4,12]. But there are many situations in which this does not work so well: for instance, with speech input, or when processing text in batch mode. Then it can be more appropriate to include not only the CNL but also the host language in the system. The system as a whole will then be similar to an embedded DSL with its general-purpose host language.

The easiest way to combine a CNL with a host language is to introduce a new start category `S` with two productions: one using the CNL start category as its argument type, the other using the host language start category:

```
fun UseCNL  : S_CNL  -> S
fun UseHost : S_Host -> S
```

The CNL trees can be given priority by biasing the weights of these functions in probabilistic GF parsing [13]. The CNL parse tree will then appear as the first alternative, whenever it can be found. Since the system sees where the tree comes from, it can give feedback to the user, for instance by using colours: green colour for CNL trees, yellow for host language trees, and red for unanalysed input, as shown (in greyscale in the printed version) in Figure 1.

Since the RGL is designed to guarantee grammaticality, and since all grammars leak, the RGL does not cover any language entirely. But if the purpose is wide-coverage parsing, we can relax this strictness. An easy way to do this is to extend the grammar with **chunking**. A chunk can be built from almost any RGL category: sentences, noun phrases, nouns, adjectives, etc:

```
fun ChunkS  : S_Host -> Chunk
fun ChunkNP : NP     -> Chunk
fun ChunkN  : N      -> Chunk
```

The top-level grammar has a production that recognizes lists of chunks (`[Chunk]`):

```
fun UseChunks : [Chunk] -> S
```

It is relatively easy to make the chunking grammar **robust**, in the sense that it returns some analysis for any combination of words. If the input has out-of-dictionary words, they can be dealt with by named entity recognition and morphological guessing. Weights can be set in such a way that longer chunks are favoured. For instance, *this old city* should be analyzed as one NP chunk rather than a determiner chunk followed by an adjective chunk and a noun chunk. The user can be given feedback, not only by a red colour indicating that chunks are used, but also by showing the chunk boundaries.

A further step of integration between CNL and the host language is obtained if CNL sentences are treated as chunks:

```
fun ChunkCNL : S_CNL -> Chunk
```

In the resulting trees, one can use different colours for different chunks inside one and the same sentence.

But since both the CNL and the host language are defined in terms of the same RGL structures, one can take the integration much further, by *intertwining* the CNL and host language rules. Consider again the CNL predicate

```
fun aged : Person -> Numeral -> Fact
```

where `Person` and `Fact` are CNL categories and `Numeral` is an RGL category. One can generalize the use of this predicate and other ones by introducing coercions from RGL categories to CNL categories used as arguments, and from CNL categories used as values to RGL categories (making sure that no cycles are created):

```
fun np2person : NP -> Person
fun fact2cl   : Fact -> Cl
```

The effect of this is seen when we analyse the English sentence

<p align="center"><em>John does not believe that the queen is sixty-five years old</em></p>

The resulting tree is

```
mkS negativePol (mkCl John believe_VS (fact2Cl
     (aged (np2person (mkNP the_Det queen_N)) (mkNumeral "65"))))
```

where those parts that belong to the CNL are boldfaced. Thus the predicate `aged` is from the CNL, but uses as argument *the queen*, which is not in the CNL. The resulting `Fact` is used as a complement to the verb *believe*, which requires an RGL clause. The resulting French translation is

<p align="center"><em>John ne croit pas que la reine</em> <strong>ait soixante-cinq ans</strong></p>

which correctly renders the `aged` idiom defined by the CNL, even though its subject is not in the CNL, and even though the negated main verb puts it into the subjunctive mood, which might never occur in the CNL itself.
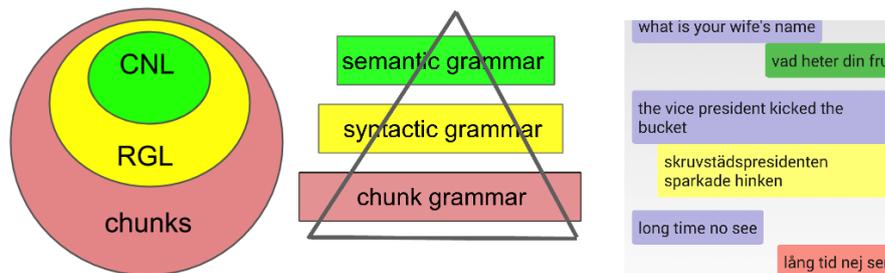


**Fig. 1.** From left: CNL embedded in general purpose RGL embedded in a chunk grammar; the corresponding levels in the Vauquois triangle; a mobile translation application showing the level of confidence in colours (from top: semantic translation from CNL, syntactic translation from RGL, chunk-based translation).

## 4   Work in Progress

The translation example with the `aged` predicate shows that embedded CNL functions introduce semantic structures in translation. This has been exploited in the wide-coverage GF-based translation system [14] [1] where three levels are distinguished by

---

[1] "The Human Language Compiler", `http://grammaticalframework.org/demos/app.html`

colours: green for the CNL (the MOLTO phrasebook [6]), yellow for the RGL syntax, and red for chunks (Figure 1). These levels correspond to levels in the **Vauquois triangle**, where translation systems are divided on the basis of whether they use a semantic interlingua, syntactic transfer, or word-to-word transfer [15]. The effect of the layered structure with an embedded CNL is that these levels coexist in one and the same system. The system is currently implemented for 11 languages. Its architecture permits easily changing the CNL to some other one and also including several CNLs in a single system. At least two experiments are in progress: one with a mathematical grammar library [5], another with the multilingual version of Attempto Controlled English [8].

In addition to translation, embedded CNLs could be used in query systems, where "green" parse trees are semantically well-formed queries and other colours are interpreted as key phrases and key words. It can also be interesting to match out-of-CNL trees with CNL trees and try to find the closest semantically complete interpretations. Yet another application of embedded CNLs would be to implement languages of the "Simplified English" type, which are not defined by formal grammars but by restrictions posed on the full language [1]. Such a language could be parsed by using the host language grammar together with a procedure that checks on the tree level how the restrictions are followed.

### Acknowledgement

## References

1. Kuhn, T.: A Survey and Classification of Controlled Natural Languages. Computational Linguistics **40**(1) (2014) 121–170
2. Hudak, P.: Building domain-specific embedded languages. ACM Computing Surveys (CSUR) **28**(4es) (1996) 196
3. Ranta, A.: Grammatical Framework: Programming with Multilingual Grammars. CSLI Publications, Stanford (2011)
4. Ranta, A., Angelov, K.: Implementing Controlled Languages in GF. In: Proceedings of CNL-2009, Athens. Volume 5972 of LNCS. (2010) 82–101
5. Saludes, J., Xambo, S.: The GF mathematics library. In: THedu'11. (2011)
6. Ranta, A., Détrez, G., Enache, R.: Controlled language for everyday use: the molto phrasebook. In: CNL 2012: Controlled Natural Language. Volume 7175 of LNCS/LNAI. (2012)
7. Davis, B., Enache, R., Van Grondelle, J., Pretorius, L.: Multilingual Verbalisation of Modular Ontologies using GF and lemon. In: Controlled Natural Language. Springer (2012) 167–184
8. Kaljurand, K., Kuhn, T.: A multilingual semantic wiki based on Attempto Controlled English and Grammatical Framework. In: The Semantic Web: Semantics and Big Data. Springer (2013) 427–441
9. Ranta, A.: The GF Resource Grammar Library. Linguistics in Language Technology **2** (2009)
10. Ranta, A.: Grammars as Software Libraries. In Bertot, Y., Huet, G., Lévy, J.J., Plotkin, G., eds.: From Semantics to Computer Science. Essays in Honour of Gilles Kahn, Cambridge University Press (2009) 281–308

11. Khegai, J., Nordström, B., Ranta, A.: Multilingual Syntax Editing in GF. In Gelbukh, A., ed.: Intelligent Text Processing and Computational Linguistics (CICLing-2003), Mexico City, February 2003. Volume 2588 of LNCS., Springer-Verlag (2003) 453–464

12. Kuhn, T.: Codeco: A Practical Notation for Controlled English Grammars in Predictive Editors. In: Controlled Natural Language (CNL), Springer (2012)

13. Angelov, K., Ljunglöf, P.: Fast statistical parsing with parallel multiple context-free grammars. In: Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, Gothenburg, Sweden, Association for Computational Linguistics (April 2014) 368–376

14. Angelov, K., Bringert, B., Ranta, A.: Speech-enabled hybrid multilingual translation for mobile devices. In: Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics, Gothenburg, Sweden, Association for Computational Linguistics (April 2014) 41–44

15. Vauquois, B.: A survey of formal grammars and algorithms for recognition and transformation in mechanical translation. In: IFIP Congress (2). (1968) 1114–1122