Formalization of context-free language theory

Marcus V. M. Ramos

Centro de Informática, UFPE, Recife, Brazil

Ruy J. G. B. de Queiroz

Centro de Informática, UFPE, Recife, Brazil

Nelma Moreira

Departamento de Ciência de Computadores, Faculdade de Ciências, Universidade do Porto, Porto, Portugal

José Carlos Bacelar Almeida

Departamento de Informática, Universidade do Minho, Braga, Portugal

Abstract

Context-free language theory is a subject of high importance in computer language processing technology as well as in formal language theory. This paper presents a formalization, using the Coq proof assistant, of fundamental results related to context-free grammars and languages. These include closure properties (union, concatenation and Kleene star), grammar simplification (elimination of useless symbols inaccessible symbols, empty rules and unit rules) and the existence of a Chomsky Normal Form for context-free grammars.

Keywords: Context-free language theory, context-free languages, context-free grammars, closure properties, grammar simplification, Chomsky Normal Form, formalization, formal mathematics, proof assistant, interactive proof systems, Coq.

1. Introduction

The formalization of context-free language theory is a key to the certification of compilers and programs, as well as to the development of new languages and tools for certified programming.

The objective of this work is to formalize a substantial part of context-free language theory in the Coq proof assistant, making it possible to reason about

Email addresses: mvmr@cin.ufpe.br (Marcus V. M. Ramos), ruy@cin.ufpe.br (Ruy J. G. B. de Queiroz), nam@dcc.fc.up.pt (Nelma Moreira), jba@di.uminho.pt (José Carlos Bacelar Almeida)

it in a fully checked environment, with all the related advantages. Initially, however, the focus has been restricted to context-free grammars and associated results. Pushdown automata and their relation to context-free grammars are not considered.

In order to follow this paper, the reader is required to have basic knowledge of Coq and of context-free language theory. For the beginner, the recommended starting point for Coq is the book by Bertot and Castéran [4]. Background on context-free language theory can be found in [13] or [11], among others.

The general idea of formalizing context-free language theory in the Coq proof assistant is discussed in Section 2. A library that contains fundamental results on context-grammars, and supports the whole formalization, is briefly presented in Section 3. In Section 4 we give an overview of the work, by explaining the common approach adopted in its different parts. Specific results related to the formalization of closure properties of context-free languages, grammar simplification and Chomsky Normal Form are presented, respectively, in Sections 5, 6 and 7. Section 8 discusses related work by various other researchers and final conclusions are presented in Section 9.

As far as the authors are aware of, this is the first comprehensive formalization of important results of context-free language theory in the Coq proof assistant. Previous publications of the authors presented the formalization of closure properties for context-free grammars (in an earlier version) [9] and of simplification for context-free grammars [10]. All the definitions and proof scripts discussed in this paper were written in plain Coq and are available for download at:

https://github.com/mvmramos/chomsky

2. Basic Definitions

In this section we present how the main concepts and objects of contextfree language theory were defined in our formalization in Coq. They are used throughout the work.

2.1. Grammars

Context-free grammars were represented in Coq very closely to the usual algebraic definition $G=(V,\Sigma,P,S)$, where V is the vocabulary of G (it includes all non-terminal and terminal symbols), Σ is the set of terminal symbols (used in the construction of the sentences of the language generated by the grammar), $N=V\setminus\Sigma$ is the set of non-terminal symbols (representing different sentence abstractions), P is the set of rules and $S\in N$ is the start symbol (also called initial or root symbol). Rules have the form $\alpha\to\beta$, with $\alpha\in N$ and $\beta\in V^*$.

Basic definitions in Coq are presented below. The N and Σ sets are represented separately from G (respectively, by types non_terminal and terminal). Notations sf (sentential form) and sentence represent lists, possibly empty, of respectively terminal and non-terminal symbols and terminal only symbols.

```
Variables non_terminal terminal: Type.

Notation sf := (list (non_terminal + terminal)).

Notation sentence := (list terminal).

Notation nlist:= (list non_terminal).
```

The record representation cfg has been used for G. The definition states that cfg is a new type and contains three components. The first is the start_symbol of the grammar (a non-terminal symbol) and the second is rules, that represents the rules of the grammar. Rules are propositions (represented in Coq by Prop) that take as arguments a non-terminal symbol and a (possibly empty) list of non-terminal and terminal symbols (corresponding, respectively, to the left and right-hand side of a rule).

The predicate rules_finite_def assures that the set of rules of the grammar is finite by proving that the length of right-hand side of every rule is equal or less than a given value, and also that both left and right-hand side of the rules are built from finite sets of, respectively, non-terminal and terminal symbols (represented here by lists).

```
Definition rules_finite_def
   (non_terminal terminal : Type)
   (ss: non_terminal)
   (rules: non_terminal \rightarrow sf \rightarrow Prop)
   (n: nat)
   (ntl: list non_terminal)
   (tl: list terminal) :=
In ss ntl \wedge
(∀ left: non_terminal,
 ∀ right: list (non_terminal + terminal),
 \texttt{rules left right} \to
 \texttt{length} \; \texttt{right} \leq \texttt{n} \; \land \;
 In left ntl \wedge
 (\forall s: non\_terminal, In (inl s) right \rightarrow In s ntl) \land
 (\forall s: terminal, In (inr s) right \rightarrow In s tl)).
Record cfg (non_terminal terminal : Type): Type:= {
start_symbol: non_terminal;
rules: non_terminal \rightarrow sf \rightarrow Prop;
rules_finite:
   \exists n: nat,
   ∃ ntl: nlist.
   ∃ tl: tlist,
   rules_finite_def start_symbol rules n ntl tl }.
```

The decision of representing rules as propositions has the consequence that it will not allow for direct extraction of executable code from the formalization. It would surely be desirable, however, to be able to obtain certified algorithms for the operations described in this article. The alternative, in this case, would be to represent rules as a member of type list (non_terminal * sf) instead. This, however, would have changed the declarative approach of the present work into the algorithmic approach, by creating functions that generate new grammars

with the desired properties. The purely logical approach was considered more appealing, since its maps directly from the textbooks, and thus was selected as the choice for the present formalization. In any case, it does not affect the objectives listed in Section 1 and can be adapted in the future in order to allow for code extraction, although this should demand a considerable effort in the creation and proof of program-related scripts.

The example below represents the grammar

$$G = (\{S', A, B, a, b\}, \{a, b\}, \{S' \rightarrow aS', S' \rightarrow b\}, S')$$

that generates language a^*b :

```
Inductive nt1: Type:= \mid S' \mid A \mid B.

Inductive t1: Type:= \mid a \mid b.

Inductive rs1: nt1 \rightarrow list (nt1 + t1) \rightarrow Prop:=
r1: rs1 S' [inr a; inl S']
\mid r2: rs1 S' [inr b].

Definition g1: cfg nt1 t1:= {\mid start_symbol:= S'; rules:= rs1; rules_finite:= rs1_finite \mid}.
```

The term rs1_finite (the proof that the set of rules of g1 is finite) is not presented here, but can be easily constructed and is available from the link provided in Section 1.

2.2. Derivations

Another fundamental concept used in this formalization is the idea of *derivation*: a grammar g *derives* a string s2 from a string s1 if there exists a series of rules in g that, when applied to s1, eventually result in s2. A direct derivation (i.e, the application of a single rule) is represented by $s_1 \Rightarrow s_2$, and the reflexive and transitive closure of this relation (i.e, the application of zero or more rules) is represented by $s_1 \Rightarrow s_2$. An inductive predicate definition of this concept in Coq (derives) uses two constructors:

```
Inductive derives
  (non_terminal terminal : Type)
  (g : cfg non_terminal terminal
```

The constructors of this definition (derives_refl and derives_step) are the axioms of our theory. Constructor derives_refl asserts that every sentential form s can be derived from s itself. Constructor derives_step states that if a sentential form that contains the left-hand side of a rule is derived by a grammar, then the grammar derives the sentential form with the left-hand side replaced by the right-hand side of the same rule. This case corresponds to the application of a rule in a direct derivation step.

A grammar generates a string if this string can be derived from its start symbol. Finally, a grammar produces a sentence if it can be generated from its start symbol.

```
Definition generates (g: cfg) (s: sf): Prop:= derives g [inl (start_symbol g)] s.

Definition produces (g: cfg) (s: sentence): Prop:= generates g (map terminal_lift s).
```

Function terminal_lift converts a terminal symbol into an ordered pair of type (non_terminal + terminal). With these definitions, it has been possible to prove various lemmas about grammars and derivations, and also operations on grammars, all of which were useful when proving the main theorems of this article.

As an example, the lemma that states that G produces the string aab (that is, that $aab \in L(G)$) is represented as:

```
Lemma g1_produces_aab: produces g1 [a; a; b].
```

The proof of this lemma can be easily constructed and relates directly to the derivations in $S \Rightarrow aS \Rightarrow aaS \Rightarrow aab$, however in reverse order because of the way that derives is defined.

2.3. Languages

A language is a set of strings over a given alphabet. It is also useful to define the language generated by a grammar as the set of terminal strings generated by the grammar through derivations:

$$L(G) = \{ w \, | \, S \Rightarrow_g^* w \}$$

From the formalization point of view, we have defined language as a function that is parametrized over a certain type (representing the set of terminal symbols), takes a string built from elements of this type and returns a proposition asserting that the string belongs to the language:

```
Definition lang (terminal: Type):= sentence \rightarrow Prop.
```

The language generated by a grammar is then a function whose return value is the predicate produces presented earlier:

```
Definition lang_of_g (g: cfg): lang := fun w: sentence \Rightarrow produces g w.
```

Finally, we are able to define the equality of two languages, and also the fact that a certain language is a context-free language (which, in this case, means that there exists a context-free grammar that generates this language):

```
Definition lang_eq (1 k: lang) :=
\forall w, 1 w \leftrightarrow k w.
Infix "==" := lang_eq (at level 80).
Definition cfl (terminal: Type) (1: lang terminal): Prop:=
∃ non_terminal: Type,
∃ g: cfg non_terminal terminal,
1 == lang_of_g g.
   The symbol == is a notation for lang_eq.
   Two grammars g_1 (with start symbol S_1) and g_2 (with start symbol S_2)
are equivalent (denoted g_1 \equiv g_2) if they generate the same language, that is,
\forall s, (S_1 \Rightarrow_{q_1}^* s) \leftrightarrow (S_2 \Rightarrow_{q_2}^* s). This is represented in our formalization in Coq
by the predicate g_equiv:
Definition g_equiv
(non_terminal1 non_terminal2 terminal : Type)
(g1: cfg non_terminal1 terminal)
(g2: cfg non_terminal2 terminal): Prop:=
\forall s: sentence,
produces g1 s \leftrightarrow produces g2 s.
```

3. Generic CFG Library

The definitions presented in the previous section allowed the construction of a generic library of fundamental lemmas on context-free grammars. This library was later used in the formalization of the specific results discussed next. It includes, among others, the formalization of the following statements and the corresponding proofs:

```
• \forall g, s_1, s_2, s_3, (s_1 \Rightarrow_a^* s_2) \rightarrow (s_2 \Rightarrow_a^* s_3) \rightarrow (s_1 \Rightarrow_a^* s_3)
```

•
$$\forall g, s_1, s_2, s, s', (s_1 \Rightarrow_q^* s_2) \rightarrow (s \cdot s_1 \cdot s' \Rightarrow_q^* s \cdot s_2 \cdot s')$$

•
$$\forall g, s_1, s_2, s_3, s_4, (s_1 \Rightarrow_q^* s_2) \to (s_3 \Rightarrow_q^* s_4) \to (s_1 \cdot s_3 \Rightarrow_q^* s_2 \cdot s_4)$$

•
$$\forall g, s_1, s_2, s_3,$$

 $(s_1 \cdot s_2 \Rightarrow_g^* s_3) \to \exists s_1', s_2' \mid (s_3 = s_1' \cdot s_2') \land (s_1 \Rightarrow_g^* s_1') \land (s_2 \Rightarrow_g^* s_2')$

•
$$\forall g, s_1, s_2, n, w, (s_1 \cdot n \cdot s_2 \Rightarrow_a^* w) \rightarrow \exists w' \mid (n \Rightarrow_a^* w')$$

•
$$\forall g, n, w, (n \Rightarrow_q^* w) \rightarrow (n \rightarrow_q w) \lor (\exists right \mid n \rightarrow_q right \land right \Rightarrow_q^* w)$$

•
$$\forall g_1, g_2, g_3, (g_1 \equiv g_2) \land (g_2 \equiv g_3) \rightarrow (g_1 \equiv g_3)$$

where $s, s', s_1, s'_1, s'_2, s_2, s_3, s_4$ and right are sentential forms, n is a non-terminal symbol, w is a sentence and g_1, g_2 and g_3 are context-free grammars.

Also, additional notions of derivations were defined in order to simplify the proofs. This is the case, for example, where a reduction (the opposite of a derivation) or an induction on the number of derivation steps are required. For these cases, we have defined *derives2* (which reduces a sentential form according to a rule) and *derives6* (which controls the number of rules applied in the derivation):

```
Inductive derives2
   (non_terminal terminal : Type)
   (g: cfg non_terminal terminal)
   : sf \rightarrow sf \rightarrow Prop :=
    derives2_refl:
       \forall s: sf,
       derives2 g s s
    | derives2_step :
       ∀ (s1 s2 s3 : sf)
       ∀ (left: non_terminal)
       \forall (right : sf),
       derives2 g (s1 ++right ++s2) s3 \rightarrow
       \texttt{rules g left right} \rightarrow
       derives2 g (s1 ++inl left :: s2) s3
Inductive derives6
   (non_terminal terminal : Type)
   (g: cfg non_terminal terminal)
   : \ \mathtt{nat} \to \mathtt{sf} \to \mathtt{sf} \to \mathtt{Prop} :=
    derives6_0:
       \forall s: sf,
       derives6 g 0 s s
    derives6_sum :
       ∀ (left: non_terminal)
       \forall (\texttt{right} : \texttt{sf})
       \forall (i : nat)
       ∀ (s1 s2 s3 : sf),
       \mathtt{rules}\ \mathtt{g}\ \mathtt{left}\ \mathtt{right} \to
       derives6 g i (s1 ++right ++s2) s3 \rightarrow
       derives6 g (S i) (s1 ++[inl left] ++s2) s3
```

For the first case, we proved derives $g \ s_1 \ s_2 \leftrightarrow derives 2 \ g \ s_1 \ s_2$. For the second case, we proved derives $g \ s_1 \ s_2 \leftrightarrow \exists n, derives 6 \ g \ n \ s_1 \ s_2$

4. Methodology

This formalization is essentially about context-free grammar manipulation. That is, about the definition of a new grammar from a previous one (or two), such that it satisfies some very specific properties. This is exactly the case when we define new grammars that generate the union, concatenation, closure

(Kleene star) of given input grammar(s). Also, when we create new grammars that exclude empty rules, unit rules, useless symbols and inaccessible symbols from the original ones. Finally, when we propose a new grammar based on some other grammar, that satisfies a specific normalization standard, the Chomsky Normal Form.

For all these cases, the following approach has been adopted:

- 1. Depending on the case, define a new type of non-terminal symbols; this will be important, for example, when we want to guarantee that the start symbol of the grammar does not appear in the right-hand side of any rule or when we have to construct new non-terminals from the existing ones;
- Inductively define the rules of the new grammar, in a way that allows the
 construction of the proofs that the resulting grammar has the required
 properties; these new rules will likely make use of the new non-terminal
 symbols described above;
- 3. Define the new grammar by using the new non-terminal symbols and the new rules; define the new start symbol (which might be a new symbol or an existing one) and build a proof of the finiteness of the set of rules for this new grammar;
- 4. State and prove all the lemmas and theorems that will assert that the newly defined grammar has the desired properties;
- 5. Consolidate the results within the same scope and finally with the previously obtained results.

In the following sections, this approach will be explored with further detail for each main result achieved in this work.

5. Closure Properties

After context-free grammars and derivations were defined, and the generic CFG library was built, the basic operations of concatenation, union and closure for context-free grammars were described in a rather straightforward way. These operations provide, as their name suggests, new context-free grammars that generate, respectively, the concatenation, the union and the closure of the language(s) generated by the input grammar(s).

5.1. Union

Given two arbitrary context-free grammars g_1 and g_2 , the following definitions are used to construct g_3 such that $L(g_3) = L(g_1) \cup L(g_2)$ (that is, the language generated by g_3 is the union of the languages generated by g_1 and g_2).

The first definition below (g_uni_nt) represents the type of the non-terminal symbols of the union grammar, created from the non-terminal symbols of the source grammars (respectively, non_terminal1 and non_terminal2). Initially, the non-terminals of the source grammars are mapped to non-terminals of the union grammar. Second, there is the need to add a new and unique non-terminal symbol (Start_uni), which will be the start symbol of the union grammar.

The functions <code>g_uni_sf_lift1</code> and <code>g_uni_sf_lift2</code> simply map sentential forms from, respectively, the first or the second grammar, and produce sentential forms of the union grammar. This will be useful when defining the rules of the union grammar.

```
Inductive g_uni_nt (non_terminal_1 non_terminal_2 : Type): Type:=
 Start_uni
 {\tt Transf1\_uni\_nt: non\_terminal\_1 \rightarrow g\_uni\_nt}
| Transf2_uni_nt: non_terminal_2 \rightarrow g_uni_nt.
Notation sf1:= (list (non_terminal_1 + terminal)).
Notation sf2:= (list (non_terminal_2 + terminal)).
Notation sfu:= (list (g_uni_nt + terminal)).
Definition g_uni_sf_lift1 (c: non_terminal_1 + terminal)
: g_uni_nt + terminal:=
  match c with
  | inl nt ⇒ inl (Transf1_uni_nt nt)
  | inr t \Rightarrow inr t
  end.
Definition g_uni_sf_lift2 (c: non_terminal_2 + terminal)
: g_uni_nt + terminal:=
  match c with
  | inl nt \Rightarrow inl (Transf2_uni_nt nt)
  | inr t \Rightarrow inr t
  end.
```

The rules of the union grammar are represented by the inductive definition <code>g_uni_rules</code>. Constructors <code>Start1_uni</code> and <code>Start2_uni</code> state that two new rules are added to the union grammar: respectively the rule that maps the new start symbol to the start symbol of the first grammar, and the rule that does the same for the second grammar. Then, constructors <code>Lift1_uni</code> and <code>Lift2_uni</code> simply map rules of first (resp. second) grammar into rules of the union grammar.

```
Inductive g_uni_rules
(non_terminal_1 non_terminal_2 terminal : Type)
(g1: cfg non_terminal_1 terminal)
(g2: cfg non_terminal_2 terminal)
: g_uni_nt \rightarrow sfu \rightarrow Prop :=
| Start1_uni:
    g_uni_rules g1 g2 Start_uni [inl (Transf1_uni_nt (start_symbol g1))]
| Start2_uni:
    g_uni_rules g1 g2 Start_uni [inl (Transf2_uni_nt (start_symbol g2))]
| Lift1_uni:
    \forall nt: non_terminal_1,
    \forall s: sf1,
    rules g1 nt s \rightarrow
    g_uni_rules g1 g2 (Transf1_uni_nt nt) (map g_uni_sf_lift1 s)
| Lift2_uni:
```

```
\label{eq:continuous_state} \begin{array}{l} \forall \ nt: \ non\_terminal\_2, \\ \forall \ s: \ sf2, \\ \\ rules \ g2 \ nt \ s \rightarrow \\ \\ g\_uni\_rules \ g1 \ g2 \ (Transf2\_uni\_nt \ nt) \ (map \ g\_uni\_sf\_lift2 \ s). \end{array}
```

Finally, g_uni describes how to create a union grammar from two arbitrary source grammars. It uses the previous definitions to give values to each of the components of a new grammar definition.

Definition g_uni (non_terminal_1 non_terminal_2 terminal : Type) (g1: cfg non_terminal_1 terminal) (g2: cfg non_terminal_2 terminal) : (cfg g_uni_nt terminal):= {| start_symbol:= Start_uni; rules:= g_uni_rules g1 g2; rules_finite:= g_uni_finite g1 g2 |}.

Similar definitions were created to represent the concatenation of any two grammars and the closure of a grammar.

5.2. Concatenation

Given two arbitrary context-free grammars g_1 and g_2 , the following definitions are used to construct g_3 such that $L(g_3) = L(g_1) \cdot L(g_2)$ (that is, the language generated by g_3 is the concatenation of the languages generated by g_1 and g_2).

```
Inductive g_cat_nt (non_terminal_1 non_terminal_2 terminal : Type): Type:=
 Start_cat
 {\tt Transf1\_cat\_nt: non\_terminal\_1} \to {\tt g\_cat\_nt}
| Transf2_cat_nt: non_terminal_2 \rightarrow g_cat_nt.
Notation sf1:= (list (non_terminal_1 + terminal)).
Notation sf2:= (list (non_terminal_2 + terminal)).
Notation sfc := (list (g_cat_nt + terminal)).
Definition g_cat_sf_lift1 (c: non_terminal_1 + terminal):
g_cat_nt + terminal:=
  match c with
  | inl nt ⇒ inl (Transf1_cat_nt nt)
  | inr t \Rightarrow inr t
  end.
Definition g_cat_sf_lift2 (c: non_terminal_2 + terminal):
g_cat_nt + terminal:=
 match c with
  | inl nt ⇒ inl (Transf2_cat_nt nt)
  | inr t \Rightarrow inr t
  end.
Inductive g_cat_rules
```

```
(non_terminal_1 non_terminal_2 terminal : Type)
(g1: cfg non_terminal_1 terminal)
(g2: cfg non_terminal_2 terminal)
: \ \mathtt{g\_cat\_nt} \to \mathtt{sfc} \to \mathtt{Prop} :=
New_cat:
  g_cat_rules g1 g2 Start_cat
   ([inl (Transf1_cat_nt (start_symbol g1))]++
   [inl (Transf2_cat_nt (start_symbol g2))])
| Lift1_cat:
  \forall nt s,
  rules g1 nt s \rightarrow
   g_cat_rules g1 g2 (Transf1_cat_nt nt) (map g_cat_sf_lift1 s)
| Lift2_cat:
   \forall nt s,
  rules g2 nt s \rightarrow
  g_cat_rules g1 g2 (Transf2_cat_nt nt) (map g_cat_sf_lift2 s).
Definition g_cat
(non_terminal_1 non_terminal_2 terminal : Type)
(g1: cfg non_terminal_1 terminal)
(g2: cfg non_terminal_2 terminal)
: (cfg g_cat_nt terminal):=
   {| start_symbol:= Start_cat;
      rules:= g_cat_rules g1 g2;
      rules_finite:= g_cat_finite g1 g2 |}.
```

In this case, the new grammar (g_cat g1 g2) is built in such a way that it has all the rules of g1 and g2, plus a new rule that maps the new start symbol (Start_cat) to the concatenation of the start symbols of the argument grammars (respectively start_symbol g1 and start_symbol g2).

5.3. Kleene Star

Given an arbitrary context-free grammar g_1 , the following definitions are used to construct g_2 such that $L(g_2) = (L(g_1))^*$ (that is, the language generated by g_2 is the reflexive and transitive concatenation (Kleene star) of the language generated by g_1).

```
Notation sfc:= (list (g_clo_nt + terminal)).

Inductive g_clo_nt (non_terminal : Type): Type := | Start_clo : g_clo_nt | Transf_clo_nt : non_terminal \rightarrow g_clo_nt.

Definition g_clo_sf_lift (c: non_terminal + terminal): g_clo_nt + terminal:= match c with | inl nt \Rightarrow inl (Transf_clo_nt nt) | inr t \Rightarrow inr t end.
```

```
Inductive g_clo_rules
(non_terminal terminal : Type)
(g: cfg non_terminal terminal)
: \ \mathtt{g\_clo\_nt} \to \mathtt{sfc} \to \mathtt{Prop} :=
New1_clo:
   g_clo_rules g Start_clo ([inl Start_clo] ++
   [inl (Transf_clo_nt (start_symbol g))])
 New2_clo:
   g_clo_rules g Start_clo []
 Lift_clo:
   \forall nt: non_terminal,
   \forall s: sf,
   \mathtt{rules}\;\mathtt{g}\;\mathtt{nt}\;\mathtt{s}\to
   g_clo_rules g (Transf_clo_nt nt) (map g_clo_sf_lift s).
Definition g_clo (g: cfg non_terminal terminal):
(non_terminal terminal : Type)
(g: cfg g_clo_nt terminal):=
 {| start_symbol:= Start_clo;
    rules:= g_clo_rules g;
    rules_finite:= g_clo_finite g | }.
```

In this case, the new grammar $(g_clo\ g)$ is built in such a way that it has all the rules of g, plus two new rules: one that generates the empty string directly from the start symbol of $g_clo\ g$, and another one that allows for the arbitrary concatenation of strings generated by g.

5.4. Correctness and Completeness

Although simple in their structure, it must be proved that the definitions g_uni, g_cat and g_clo always produce the correct result. In other words, these definitions must be "certified", which is one of the main goals of formalization. In order to accomplish this, we must first state the theorems that capture the expected semantics of these definitions. Finally, we have to derive proofs of the correctness of these theorems.

This can be done with a pair of theorems for each grammar definition: the first relates the output to the inputs, and the other one does the converse, providing assumptions about the inputs once an output is generated. This is necessary in order to guarantee that the definitions do only what one would expect, and no more.

In what follows, an informal statement is presented right before the corresponding Coq theorem. This is intended to abstract over the necessary mappings that occur in the Coq terms, due to the different types of non-terminals, sentential forms etc involved.

For concatenation, the following Coq statement expresses the result we want to prove (considering that g_3 is the concatenation of g_1 and g_2 and g_3 , g_4 and g_5 are, respectively, the start symbols of g_3 , g_4 and g_5 :

$$\forall g_1 g_2, s_1, s_2, (S_1 \Rightarrow_{q_1}^* s_1) \land (S_2 \Rightarrow_{q_2}^* s_2) \rightarrow (S_3 \Rightarrow_{q_3}^* s_1 s_2)$$

```
Theorem g_cat_correct:  \begin{tabular}{ll} \hline $\forall$ g1: cfg non\_terminal\_1 terminal, \\ \hline $\forall$ g2: cfg non\_terminal\_2 terminal, \\ \hline $\forall$ s1: sf1, \\ \hline $\forall$ s2: sf2, \\ \hline $generates g1 s1 \land generates g2 s2 \rightarrow \\ \hline $generates (g\_cat g1 g2) ((map g\_cat\_sf\_lift1 s1)++(map g\_cat\_sf\_lift2 s2)). \\ \hline \end{tabular}
```

The above theorem states that if context-free grammars g1 and g2 generate, respectively, strings s1 and s2, then the concatenation of these two grammars, according to the proposed algorithm, generates the concatenation of string s1 with string s2. As mentioned before, the above theorem alone does not guarantee that g_cat will not produce outputs other than the concatenation of its input strings. This idea is captured by the following complementary theorem:

$$\forall s_3, (S_3 \Rightarrow_{g_3}^* s_3) \to \exists s_1, s_2 \mid (s_3 = s_1 \cdot s_2) \land (S_1 \Rightarrow_{g_1}^* s_1) \land (S_2 \Rightarrow_{g_2}^* s_2)$$

For the converse of concatenation, the following Coq statement expresses the result we want to prove:

```
Theorem g_cat_correct_inv:  \begin{tabular}{ll} \hline $\forall$ g1: cfg non\_terminal\_1 terminal, \\ $\forall$ g2: cfg non\_terminal\_2 terminal, \\ $\forall$ s: sfc, \\ generates (g_cat g1 g2) s \rightarrow \\ s = [inl (start\_symbol (g_cat g1 g2))] \lor \\ \hline $\exists$ s1: sf1, \\ \hline $\exists$ s2: sf2, \\ s = (map g_cat\_sf_lift1 s1) + + (map g_cat\_sf_lift2 s2) \land \\ generates g1 s1 \land generates g2 s2. \\ \hline \end{tabular}
```

The idea here is to express that, if a string is generated by <code>g_cat</code>, then it must only result from the concatenation of strings generated by the grammars combined by the definition. Together, these two theorems represent the semantics of the context-free grammar concatenation operation. The same ideas have been applied to the statement and proof of the following theorems, relative to the union and closure operations.

For union, we need to prove (considering that g_3 is the union of g_1 and g_2 and S_3, S_1 and S_2 are, respectively, the start symbols of g_3, g_1 and g_2):

$$\forall g_1, g_2, s_1, s_2, (S_1 \Rightarrow_{g_1}^* s_1 \to S_3 \Rightarrow_{g_3}^* s_1) \land (S_2 \Rightarrow_{g_2}^* s_2 \to S_3 \Rightarrow_{g_3}^* s_2)$$

which translates in Coq into:

```
Theorem g_uni_correct:  \forall \ g1: \ cfg \ non\_terminal\_1 \ terminal, \\ \forall \ g2: \ cfg \ non\_terminal\_2 \ terminal, \\ \forall \ s1: \ sf1, \\ \forall \ s2: \ sf2, \\ (generates \ g1 \ s1 \rightarrow generates \ (g_uni \ g1 \ g2) \ (map \ g_uni\_sf_lift1 \ s1)) \\ \land \\ (generates \ g2 \ s2 \rightarrow generates \ (g_uni \ g1 \ g2) \ (map \ g_uni\_sf_lift2 \ s2)).
```

For the converse of union we have:

$$\forall s_3, (S_3 \Rightarrow_{q_3}^* s_3) \to (S_1 \Rightarrow_{q_1}^* s_3) \lor (S_2 \Rightarrow_{q_2}^* s_3)$$

```
Theorem g_uni_correct_inv:  \forall \ g1: \ cfg \ non\_terminal\_1 \ terminal, \\ \forall \ g2: \ cfg \ non\_terminal\_2 \ terminal, \\ \forall \ s: \ sfu, \\ generates \ (g_uni \ g1 \ g2) \ s \rightarrow \\ (s=[inl \ (start\_symbol \ (g_uni \ g1 \ g2))]) \lor \\ (\exists \ s1: \ sf1, \ (s=(map \ g_uni\_sf\_lift1 \ s1) \land generates \ g1 \ s1)) \lor \\ (\exists \ s2: \ sf2, \ (s=(map \ g_uni\_sf\_lift2 \ s2) \land generates \ g2 \ s2)).
```

For closure, we have (considering that g_2 is the Kleene star of g_1 and S_2 and S_1 are, respectively, the start symbols of g_2 and g_1):

$$\forall g_1, s_1, s_2, (S_2 \Rightarrow_{q_2}^* \epsilon) \land ((S_2 \Rightarrow_{q_2}^* s_2) \land (S_1 \Rightarrow_{q_1}^* s_1) \rightarrow S_2 \Rightarrow_{q_2}^* s_2 \cdot s_1)$$

which translates in Coq into:

```
Theorem g_clo_correct:  \forall \ g: \ cfg \ non\_terminal \ terminal, \\ \forall \ s: \ sf, \\ \forall \ s': \ sfc, \\ generates \ (g\_clo \ g) \ nil \ \land \ (generates \ (g\_clo \ g) \ s' \ \land \ generates \ g \ s \ \rightarrow \\ generates \ (g\_clo \ g) \ (s'++ \ map \ g\_clo\_sf\_lift \ s)).
```

Finally:

$$\forall s_2, (S_2 \Rightarrow_{q_2}^* s_2) \to (s_2 = \epsilon) \lor (\exists s_1, s_2' \mid (s_2 = s_2' \cdot s_1) \land (S_2 \Rightarrow_{q_2}^* s_2') \land (S_1 \Rightarrow_{q_1}^* s_1))$$

```
Theorem g_clo_correct_inv:  \forall \ g: \ cfg \ non\_terminal \ terminal, \\ \forall \ s: \ sfc, \\ generates \ (g\_clo \ g) \ s \rightarrow \\ (s=[]) \ \lor \\ (s=[inl \ (start\_symbol \ (g\_clo \ g))]) \ \lor \\ (\exists \ s': \ sfc, \\ \exists \ s'': \ sf, \\ generates \ (g\_clo \ g) \ s' \ \land \ generates \ g \ s'' \ \land \ s=s' ++map \ g\_clo\_sf\_lift \ s'').
```

The proofs of all the six main theorems have been completed (g_uni_correct and g_uni_correct_inv for union, g_cat_correct and g_cat_correct_inv for concatenation and g_clo_correct and g_clo_correct_inv for closure). Most of them were obtained through induction over the predicate derives or one of its variants.

5.5. Closure over Languages

The previous results were all formulated over grammars, and it is desirable to obtain equivalent versions using languages instead. Thus, we have defined the union, concatenation and closure of arbitrary languages as follows:

```
Inductive l_uni (terminal : Type) (11 12: lang terminal): lang terminal:= | l_uni_11: \forall s: sentence, l1 s \rightarrow l_uni l1 l2 s | l_uni_12: \forall s: sentence, l2 s \rightarrow l_uni l1 l2 s.

Inductive l_cat (terminal : Type) (l1 l2: lang terminal): lang terminal:= | l_cat_app: \forall s1 s2: sentence, l1 s1 \rightarrow l2 s2 \rightarrow l_cat l1 l2 (s1 ++s2).

Inductive l_clo (terminal : Type) (l: lang terminal): lang terminal:= | l_clo_nil: l_clo 1 [| l_clo_app: \forall s1 s2: sentence, (l_clo 1) s1 \rightarrow 1 s2 \rightarrow l_clo 1 (s1 ++s2).
```

With these definitions, it is immediate to prove that the operations of union, concatenation and closure are correct, including the converse versions. However, it remains to be proved that the newly generated languages are also context-free, which leads to the following theorems:

```
Theorem l_uni_is_cfl: \forall \ 11\ 12:\ lang\ terminal, \\ cfl\ 11 \rightarrow cfl\ 12 \rightarrow cfl\ (l_uni\ 11\ 12). Theorem l_cat_is_cfl: \forall \ 11\ 12:\ lang\ terminal, \\ cfl\ 11 \rightarrow cfl\ 12 \rightarrow cfl\ (l_cat\ 11\ 12). Theorem l_clo_is_cfl: \forall \ l:\ lang\ terminal, \\ cfl\ 1 \rightarrow cfl\ (l_clo\ l).
```

In all cases, the proofs obtained rely on (i) the existence of context-free grammars that generated the original languages, a direct consequence of the definition of cfl and (ii) the results that were previously proved for context-free grammars.

6. Simplification

The definition of a context-free grammar, and also the operations defined in the previous section, allow for the inclusion of symbols and rules that might not contribute to the language being generated. Besides that, context-free grammars might also contain rules that can be substituted by equivalent smaller and simpler ones. Unit rules, for example, do not expand sentential forms (instead, they just rename the symbols in them) and empty rules can cause them to contract. Although the appropriate use of these features can be important for human communication in some situations, this is not the general case, since it leads to grammars that have more symbols and rules than necessary, making difficult its comprehension and manipulation. Thus, simplification is an important operation on context-free grammars.

Let G be a context-free grammar, L(G) the language generated by this grammar and ϵ the empty string. Different authors use different terminology when

presenting simplification results for context-free grammars. In what follows, we adopt the terminology and definitions of [13].

Context-free grammar simplification comprises the manipulation of rules and symbols, as described below:

- 1. An empty rule $r \in P$ is a rule whose right-hand side β is empty (e.g. $X \to \epsilon$). We formalize that for all G, there exists G' such that L(G) = L(G') and G' has no empty rules, except for a single rule $S \to \epsilon$ if $\epsilon \in L(G)$; in this case, S (the initial symbol of G') does not appear on the right-hand side of any rule in G';
- 2. A unit rule $r \in P$ is a rule whose right-hand side β contains a single non-terminal symbol (e.g. $X \to Y$). We formalize that for all G, there exists G' such that L(G) = L(G') and G' has no unit rules;
- 3. $s \in V$ is useful ([13], p. 116) if it is possible to derive a sentence from it using the rules of the grammar. Otherwise, s is called an useless symbol. A useful symbol s is one such that $s \Rightarrow^* \omega$, with $\omega \in \Sigma^*$. Naturally, this definition concerns mainly non-terminals, as terminals are trivially useful. We formalize that, for all G such that $L(G) \neq \emptyset$, there exists G' such that L(G) = L(G') and G' has no useless symbols;
- 4. $s \in V$ is accessible ([13], p. 119) if it is part of at least one string generated from the root symbol of the grammar. Otherwise, it is called an inaccessible symbol. An accessible symbol s is one such that $S \Rightarrow^* \alpha s \beta$, with $\alpha, \beta \in V^*$. We formalize that for all G, there exists G' such that L(G) = L(G') and G' has no inaccessible symbols.

Finally, we formalize a unification result: that for all G, if G is non-empty, then there exists G' such that L(G) = L(G') and G' has no empty rules (except for one, if G generates the empty string), no unit rules, no useless symbols, no inaccessible symbols and the start symbol of G' does not appear on the right-hand side of any other rule of G'.

In all these four cases and the five grammars that are discussed next (namely g_emp, g_emp', g_unit, g_use and g_acc), the proof of rules_finite is based on the proof of the correspondent predicate for the argument grammar. Thus, all new grammars satisfy the cfg specification and are finite as well.

6.1. Empty rules

Result (1) is achieved in two steps. First, the idea of a *nullable* ([13], p. 107) symbol was represented by the definition empty:

```
Definition empty
(g: cfg terminal _) (s: non_terminal + terminal): Prop:= derives g [s] [].
```

Notation sf' represents a sentential form that is constructed with elements of non_terminal' and terminal. Definition symbol_lift maps a pair of type (non_terminal + terminal) into a pair of type (non_terminal' + terminal) by replacing each non_terminal with the corresponding non_terminal':

```
Inductive non_terminal': Type:=
| Lift_nt: non_terminal → non_terminal'
| New_ss.

Notation sf' := (list (non_terminal' + terminal)).

Definition symbol_lift
(s: non_terminal + terminal): non_terminal' + terminal:=
match s with
| inr t ⇒ inr t
| inl n ⇒ inl (Lift_nt n)
end.
```

With these, a new grammar <code>g_emp</code> <code>g</code> has been created, such that the language generated by it matches the language generated by the original grammar (g), except for the empty string. Predicate <code>g_emp_rules</code> states that every nonempty rule of <code>g</code> is also a rule of <code>g_emp</code> <code>g</code>, and also adds new rules to <code>g_emp</code> <code>g</code> where every possible combination of nullable non-terminal symbols that appears on the right-hand side of a rule of <code>g</code> is removed, as long as the resulting right-hand side is not empty. Finally, it adds a rule that maps a new symbol, the start symbol of the new grammar (<code>New_ss</code>), to the start symbol of the original grammar. For this reason, the new type <code>non_terminal</code>, has been defined. The motivation for introducing a new start symbol at this point is to be able to prove that the start symbol does not appear in the right-hand side of any rule of the new grammar, a result that will be important in future developments.

```
Inductive g_emp_rules
(non_terminal terminal : Type)
(g: cfg non_terminal terminal)
: non\_terminal' \rightarrow sf' \rightarrow Prop :=
| Lift_direct :
    ∀ left: non_terminal,
    ∀ right: sf,
    right \neq [] \rightarrow rules g left right \rightarrow
    g_emp_rules g (Lift_nt left) (map symbol_lift right)
Lift_indirect:
    ∀ left: non_terminal,
    ∀ right: sf,
    g_emp_rules g (Lift_nt left) (map symbol_lift right)\rightarrow
    ∀ s1 s2: sf,
    \forall s: non_terminal,
    right = s1 ++(inl s) :: s2 \rightarrow
    \texttt{empty g (inl s)} \rightarrow
    s1 ++s2 \neq [] \rightarrow
    g_emp_rules g (Lift_nt left) (map symbol_lift (s1 ++s2))
| Lift_start_emp:
    g_emp_rules g New_ss [inl (Lift_nt (start_symbol g))].
Definition g_emp
(non_terminal terminal : Type)
```

```
(g: cfg non_terminal terminal)
: cfg non_terminal' terminal :=
    {| start_symbol:= New_ss;
        rules:= g_emp_rules g;
        rules_finite:= g_emp_finite g |}.
```

Suppose, for example, that X, A, B, C are non-terminals, of which A, B and C are nullable, a, b and c are terminals and $X \to aAbBcC$ is a rule of g. Then, the above definitions assert that $X \to aAbBcC$ is a rule of g_emp g, and also:

- $X \rightarrow aAbBc$;
- $X \rightarrow abBcC$;
- $X \rightarrow aAbcC$;
- $X \to aAbc$;
- $X \rightarrow abBc$;
- $X \rightarrow abcC$;
- $X \to abc$.

Observe that grammar g_emp g does not generate the empty string. The second step, thus, was to define g_emp'g, such that g_emp'g generates the empty string if g generates the empty string. This was done by stating that every rule from g_emp g is also a rule of g_emp'g and also by adding a new rule that allow g_emp'g to generate the empty string directly if necessary.

```
Inductive g_emp'_rules
(non_terminal terminal : Type)
(g: cfg non_terminal terminal)
: non_terminal' non_terminal \rightarrow sf' \rightarrow Prop :=
Lift_all:
   ∀ left: non_terminal'_,
  ∀ right: sf',
  rules (g_emp g) left right \rightarrow g_emp'_rules g left right
| Lift_empty:
   empty g (inl (start_symbol g)) \rightarrow g_emp'_rules g (start_symbol (g_emp g)) \|.
Definition g_emp'
(non_terminal terminal : Type)
(g: cfg non_terminal terminal)
: cfg (non_terminal'_) terminal :=
  {| start_symbol:= New_ss _;
    rules:= g_emp'_rules g;
    rules_finite:= g_emp'_finite g |}.
```

Note that the generation of the empty string by $g_{-}emp'$ g depends on g generating the empty string.

The proof of the correctness of these definitions is achieved through the following theorem:

```
Theorem g_emp'_correct:  \forall \ g: \ cfg \ non\_terminal \ terminal, \\ g_equiv \ (g_emp' \ g) \ g \ \land \\ (generates_empty \ g \rightarrow has\_one\_empty\_rule \ (g_emp' \ g)) \ \land \\ (\sim \ generates\_empty \ g \rightarrow has\_no\_empty\_rules \ (g_emp' \ g)) \ \land \\ start\_symbol\_not\_in\_rhs \ (g_emp' \ g).
```

Four auxiliary predicates have been used in this statement: g_equiv (introduced in Section 2.3) for two context-free grammars that generate the same language, generates_empty for a grammar whose language includes the empty string, has_one_empty_rule for a grammar that has an empty rule whose left-hand side is the initial symbol, and all other rules are not empty and has_no_empty_rules for a grammar that has no empty rules at all.

The definition of g_equiv, when applied to this theorem, yields:

```
\forall s: sentence, produces (g_emp'g) s \leftrightarrow produces g s.
```

For the \rightarrow part, the strategy is to prove that for every rule $left \rightarrow_{g_emp'} right$, either $left \rightarrow_{g} right$ is a rule of g or $left \Rightarrow_{g}^{*} right$. For the \leftarrow part, the strategy is a more complicated one, and involves induction over the number of derivation steps in g.

6.2. Unit rules

For result (2), definition unit expresses the relation between any two non-terminal symbols X and Y, and is true when $X \Rightarrow^* Y$ ([13], p. 114).

```
Inductive unit
(terminal non_terminal : Type)
(g: cfg terminal non_terminal)
(a: non_terminal)
: non_terminal → Prop:=
| unit_rule:
    ∀ (b: non_terminal),
    rules g a [inl b] → unit g a b
| unit_trans:
    ∀ b c: non_terminal,
    unit g a b → unit g b c → unit g a c.
```

Grammar g_unit g represents the grammar that is equivalent to g, except that the unit rules of the latter have been substituted by others, non-unit rules, that produce the same results in terms of the generated language. The idea is that g_unit g has all non-unit rules of g, plus new rules that are created by anticipating the possible application of unit rules in g, as informed by g_unit.

```
Inductive g_unit_rules
(terminal non_terminal : Type)
(g: cfg non_terminal terminal)
: non_terminal \rightarrow sf \rightarrow Prop :=
| Lift_direct':
    \forall left: non_terminal,
```

```
∀ right: sf,
    (\forall r: non\_terminal, right \neq [inl r]) \rightarrow
    \texttt{rules} \ \texttt{g} \ \texttt{left} \ \texttt{right} \rightarrow
    g_unit_rules g left right
Lift_indirect':
    \forall a b: non_terminal,
    unit g a b \rightarrow
    \forall right: sf,
    \mathtt{rules}\;\mathtt{g}\;\mathtt{b}\;\mathtt{right}\;\to\;
    (\forall \ \mathtt{c:} \ \mathtt{non\_terminal}, \ \mathtt{right} \neq [\mathtt{inl} \ \mathtt{c}]) \rightarrow
    g_unit_rules g a right.
Definition g_unit
(terminal non_terminal : Type)
(g: cfg non_terminal terminal)
: cfg non_terminal terminal :=
  {| start_symbol:= start_symbol g;
      rules:= g_unit_rules g;
      rules_finite:= g_unit_finite g | }.
```

Finally, the correctness of g_unit comes from the following theorem:

```
Theorem g_unit_correct:

∀ g: cfg non_terminal terminal,

g_equiv (g_unit g) g ∧ has_no_unit_rules (g_unit g).
```

The predicate has_no_unit_rules states that the argument grammar has no unit rules at all.

Similar to the previous case, for the \rightarrow part of the g_equiv (g_unit g) g proof, the strategy adopted is to prove that for every rule $left \rightarrow_{g_unit} right$ of (g_unit g), either $left \rightarrow_{g} right$ is a rule of g or $left \Rightarrow_{g}^{*} right$. For the \leftarrow part, the strategy is also a more complicated one, and involves induction over a predicate that is equivalent to derives (derives3), but generates the sentence directly without considering the application of a sequence of rules, which allows one to abstract the application of unit rules in g.

6.3. Useless symbols

For result (3), the idea of a useful symbol is captured by the definition useful:

```
Definition useful
(terminal non_terminal : Type)
(g: cfg non_terminal terminal)
(s: non_terminal + terminal): Prop:=
match s with
| inr t ⇒ True
| inl n ⇒ ∃ s: sentence, derives g [inl n] (map term_lift s)
end
```

The removal of useless symbols comprises, first, the identification of useless symbols in the grammar and, second, the elimination of the rules that use them.

Definition g_use_rules selects, from the original grammar, only the rules that do not contain useless symbols. The new grammar, without useless symbols, can then be defined as in g_use:

```
Inductive g_use_rules
(terminal non_terminal: Type)
(g: cfg non_terminal terminal)
: non\_terminal \rightarrow sf \rightarrow Prop :=
Lift_use:
   ∀ left: non_terminal,
   ∀ right: sf,
   rules g left right \rightarrow
   useful g (inl left) \rightarrow
   (\forall s: non\_terminal + terminal, In s right \rightarrow useful g s) \rightarrow
   g_use_rules g left right.
Definition g_use
(terminal non_terminal: Type)
(g: cfg non_terminal terminal)
: cfg non_terminal terminal:=
  {| start_symbol:= start_symbol g;
     rules:= g_use_rules g;
     rules_finite:= g_use_finite g | }.
```

The g_use definition, of course, can only be used if the language generated by the original grammar is not empty, that is, if the start symbol of the original grammar is useful. If it were useless then it would be impossible to assign a root to the grammar and the language would be empty. The correctness of the useless symbol elimination operation can be certified by proving theorem g_use_correct, which states that every context-free grammar whose start symbol is useful generates a language that can also be generated by an equivalent context-free grammar whose symbols are all useful.

```
Theorem g_use_correct:  \forall \ g: \ cfg \ non\_terminal \ terminal, \\ non\_empty \ g \rightarrow g\_equiv \ (g\_use \ g) \ g \land has\_no\_useless\_symbols \ (g\_use \ g).
```

The predicates non_empty, and has_no_useless_symbols used above assert, respectively, that grammar g generates a language that contains at least one string (which in turn may or may not be empty) and the grammar has no useless symbols at all.

The \rightarrow part of the g_equiv proof is straightforward, since every rule of g_use is also a rule of g. For the converse, it is necessary to show that every symbol used a the derivation of g is useful, and thus the rules used in this derivation also appear in g_use.

6.4. Inaccessible symbols

Result (4) is similar to the previous case, and definition accessible has been used to represent accessible symbols in context-free grammars.

```
Definition accessible
(terminal non_terminal : Type)
(g : cfg non_terminal terminal)
(s: non_terminal + terminal): Prop:=

∃ s1 s2: sf, derives g [inl (start_symbol g)] (s1++s::s2).
```

Definition g_acc_rules selects, from the original grammar, only the rules that do not contain inaccessible symbols. Definition g_acc represents a grammar whose inaccessible symbols have been removed:

```
Inductive g_acc_rules
(terminal non_terminal : Type)
(g : cfg non_terminal terminal)
: non_terminal \rightarrow sf \rightarrow Prop :=
| Lift_acc : \forall left: non_terminal,
    \forall right \right accessible g (inl left) \rightarrow g_acc_rules g left right.

Definition g_acc
(terminal non_terminal : Type)
(g : cfg non_terminal terminal)
: cfg non_terminal terminal :=
{| start_symbol:= start_symbol g;
    rules:= g_acc_rules g;
    rules_finite:= g_acc_finite g |}.
```

The correctness of the inaccessible symbol elimination operation can be certified by proving theorem <code>g_acc_correct</code>, which states that every context-free grammar generates a language that can also be generated by an equivalent context-free grammar where symbols are all accessible.

```
Theorem g_acc_correct: \forall g: cfg non_terminal terminal, g_equiv (g_acc g) g \land has_no_inaccessible_symbols (g_acc g).
```

In a way similar to has_no_useless_symbols, the absence of inaccessible symbols in a grammar is expressed by predicate has_no_inaccessible_symbols used above.

Similar to the previous case, the \rightarrow part of the <code>g_equiv</code> proof is also straightforward, since every rule of <code>g_acc</code> is also a rule of <code>g</code>. For the converse, it is necessary to show that every symbol used in the derivation of <code>g</code> is accessible, and thus the rules used in this derivation also appear in <code>g_acc</code>.

6.5. Unification

If one wants to obtain a new grammar simultaneously free of empty and unit rules, and of useless and inaccessible symbols, it is not enough to consider the previous independent results: it is necessary to establish a suitable order to apply these simplifications, in order to guarantee that the final result satisfies all desired conditions. Then, it is necessary to prove that the claims do hold.

For the order, we should start with (i) the elimination of empty rules, followed by (ii) the elimination of unit rules. The reason for this is that (i) might introduce new unit rules in the grammar, and (ii) will surely not introduce empty rules, as long as original grammar is free of them (except for $S \to \epsilon$, in which case S, the initial symbol of the grammar, must not appear on the right-hand side of any rule). Then, elimination of useless and inaccessible symbols (in either order) is the right thing to do, since they only remove rules from the original grammar (which is specially important because they do not introduce new empty or unit rules).

The formalization of this result is captured in the following theorem, which represents the main result of this section:

```
Theorem g_simpl:  \forall \ g: \ cfg \ non\_terminal \ terminal, \\ non\_empty g \rightarrow \\ \exists \ g': \ cfg \ (non\_terminal' \ non\_terminal) \ terminal, \\ g\_equiv \ g' g \land \\ has\_no\_inaccessible\_symbols \ g' \land \\ has\_no\_useless\_symbols \ g' \land \\ (generates\_empty g \rightarrow has\_one\_empty\_rule \ g') \land \\ (\sim \ generates\_empty g \rightarrow has\_no\_empty\_rules \ g') \land \\ has\_no\_unit\_rules \ g' \land \\ start\_symbol\_not\_in\_rhs \ g'.
```

Hypothesis non_empty g is necessary in order to allow for the elimination of useless symbols. The predicate start_symbol_not_in_rhs states that the start symbol does not appear in the right-hand side of any rule of the argument grammar.

The proof of <code>g_simpl</code> demands auxiliary lemmas to prove that the characteristics of the initial transformations are preserved by the following ones. For example, that all of the unit rules elimination, useless symbol elimination and inaccessible symbol elimination operations preserve the characteristics of the empty rules elimination operation.

7. Chomsky Normal Form

The Chomsky Normal Form (CNF) theorem asserts:

$$\forall G = (V, \Sigma, P, S), \ \exists G' = (V', \Sigma, P', S') \mid$$

$$L(G) = L(G') \land \forall (\alpha \to_{G'} \beta) \in P', (\beta \in \Sigma) \lor (\beta \in N \cdot N)$$

That is, every context-free grammar can be converted to an equivalent one whose rules have only one terminal symbol or two non-terminal symbols in the right-hand side. Naturally, this is valid only if G does not generate the empty string. If this is the case, then the grammar that has this format, plus a single rule $S' \to_G \epsilon$, is also considered to be in the Chomsky Normal Form, and generates the original language, including the empty string. It can also

be assured that in either case the start symbol of G' does not appear on the right-hand side of any rule of G'.

The existence of a CNF can be used for a variety of purposes, including to prove that there is an algorithm to decide whether an arbitrary context-free language accepts an arbitrary string, and to test if a language is not context-free (using the Pumping Lemma for context-free languages, which can be proved with the help of CNF grammars).

The idea of mapping G into G' consists of creating a finite number of new non-terminal symbols and new rules, in the following ways:

- 1. For every terminal symbol σ that appears in the right-hand side of a rule $r = \alpha \to_G \beta_1 \cdot \sigma \cdot \beta_2$ of G, create a new non-terminal symbol $[\sigma]$, a new rule $[\sigma] \to_{G'} \sigma$ and substitute σ for $[\sigma]$ in r;
- 2. For every rule $r = \alpha \to_G N_1 N_2 \cdots N_k$ of G, where N_i are all non-terminals, create a new set of non-terminals and a new set of rules such that:

$$\begin{array}{ccc} \alpha & \rightarrow_{G'} & N_1[N_2\cdots N_k], \\ [N_2\cdots N_k] & \rightarrow_{G'} & N_2[N_3\cdots N_k], \\ & & \ddots & \\ [N_{k-2}N_{k-1}N_k] & \rightarrow_{G'} & N_{k-2}[N_{k-1}N_k], \\ [N_{k-1}N_k] & \rightarrow_{G'} & N_{k-1}N_k \end{array}$$

Case (i) substitutes all terminal for non-terminal symbols. Case (ii) splits rules that have three or more non-terminal symbols on the right-hand side by a set of rules that has only two non-terminal symbols in the right-and side. Both changes preserve the language of the original grammar.

As an example, consider $G = (\{S', X, Y, Z, a, b, c\}, \{a, b, c\}, P, S')$ with P equal to:

$$\begin{cases} S' & \rightarrow & XYZd, \\ X & \rightarrow & a, \\ Y & \rightarrow & b, \\ Z & \rightarrow & c, \end{cases}$$

The CNF grammar G', equivalent to G, would then be the one with the following set of rules:

$$\begin{cases} S' & \rightarrow & X[YZd], \\ [YZd] & \rightarrow & Y[Zd], \\ [Zd] & \rightarrow & Z[d], \\ [d] & \rightarrow & d, \\ X & \rightarrow & a, \\ Y & \rightarrow & b, \\ Z & \rightarrow & c, \end{cases}$$

These ideas are captured by the following definitions. The non-terminals of the new grammar g_cnf g are represented by the type non_terminal'. Its elements are associated with sentential forms of g via the constructor Lift_r:

```
Inductive non_terminal' (non_terminal terminal : Type): Type:=
| Lift_r: sf \rightarrow non_terminal'.
Notation sf':= (list (non_terminal' + terminal)).
Notation term_lift:= ((terminal_lift non_terminal) terminal).
   The function symbol_lift, presented below, maps sentential forms of g into
sentential forms of g_cnf g:
Definition symbol_lift (s: non_terminal + terminal)
: non_terminal' + terminal:=
match s with
| inr t \Rightarrow inr t
| inl n ⇒ inl (Lift_r [inl n])
   The rules of g_cnf g and g_cnf g itself are defined as:
Inductive g_cnf_rules
(non_terminal terminal : Type)
(g: cfg non_terminal terminal)
: non_terminal' \rightarrow sf' \rightarrow Prop:=
Lift_cnf_t:
   \forall t: terminal,
   \forall left: non_terminal,
  \forall s1 s2: sf,
   rules g left (s1++[inr t]++s2) \rightarrow
   g_cnf_rules g (Lift_r [inr t]) [inr t]
Lift_cnf_1:
   \forall left: non_terminal,
   \forall t: terminal,
   \texttt{rules g left [inr t]} \to
   g_cnf_rules g (Lift_r [inl left]) [inr t]
| Lift_cnf_2:
   ∀ left: non_terminal,
  \forall s1 s2: symbol,
   \forall beta: sf.
   \texttt{rules g left} \; (\texttt{s1} :: \; \texttt{s2} :: \; \texttt{beta}) \to
   g_cnf_rules g (Lift_r [inl left])
   [inl (Lift_r [s1]); inl (Lift_r (s2 :: beta))]
| Lift_cnf_3:
  \forall left: sf,
   \forall s1 s2 s3: symbol,
   \forall beta: sf,
   g_cnf_rules g (Lift_r left)
   [inl (Lift_r [s1]); inl (Lift_r (s2 :: s3 :: beta))] \rightarrow
   g_cnf_rules g (Lift_r (s2 :: s3 :: beta))
   [inl (Lift_r [s2]); inl (Lift_r (s3 :: beta))].
```

Next, we prove that g_cnf g is equivalent to g. It should be noted, however, that the set of rules defined above do not generate the empty string. If this is the case, the definitions below define a new grammar g_cnf , that adds a new rule that generates the empty string:

```
Inductive g_cnf'_rules
(non_terminal terminal : Type)
(g: cfg non_terminal terminal)
: \mathtt{non\_terminal'} \to \mathtt{sf'} \to \mathtt{Prop} :=
Lift_cnf'_all:
   \forall left: non_terminal',
   ∀ right: sf',
   g\_cnf\_rules g left right \rightarrow
   g_cnf'_rules g left right
| Lift_cnf'_new:
   g_cnf'_rules g (start_symbol (g_cnf g)) [].
Definition g_cnf'
(non_terminal terminal : Type)
(g: cfg non_terminal terminal)
: cfg non_terminal' terminal:=
  {| start_symbol:= start_symbol (g_cnf g);
     rules:= g_cnf'_rules g;
     rules_finite:= g_cnf'_finite g |}.
    The statement of the CNF theorem can then be presented as:
Theorem g_cnf_final:
∀ g: cfg non_terminal terminal,
(produces_empty g \lor \sim produces_empty g) \land
(\texttt{produces\_non\_empty} \ \texttt{g} \ \lor \sim \texttt{produces\_non\_empty} \ \texttt{g}) \ \to \\
\exists g': cfg non_terminal' terminal,
{\tt g\_equiv} \ {\tt g'} \ {\tt g} \ \land \\
(is_cnf g' ∨ is_cnf_with_empty_rule g').
```

The predicates used above assert that the argument grammar:

- produces the empty string (produces_empty);
- does not produce the empty string (produces_non_empty);
- is in the Chomsky Normal Form (is_cnf);

• is in the Chomsky Normal Form and has a single empty rule with the start symbol in the left-hand side (is_cnf_with_empty_rule).

The proof of this theorem requires, among other things, that the original grammar is first simplified according to the results discussed in the previous section. For the \leftarrow part of g_equiv, the strategy adopted is to prove that for every rule $left \rightarrow right$ of (g), either $left \rightarrow right$ is a rule of g_cnf g or $left \Rightarrow^* right$ in g_cnf g.

For the \to part, that is, $(s_1 \Rightarrow_{g_cnfg}^* s_2) \to (s_1 \Rightarrow_g^* s_2)$, it is enough to note that the sentential forms of g are embedded in the sentential forms of g_cnf g, specifically in the arguments of the constructor Lift_r of non_terminal'. Thus, a simple extraction mechanism allows the implication to be proved by induction on the structure of the sentential form s_1 .

Using the previous example, suppose we have: $X[YZd] \Rightarrow_{g_cnfg}^* abcd$, which would be represented in our formalization as:

```
derives (g_cnf g) [inl X] ++[inl (Lift_r ([inl Y; inl Z; inr d]))]
(map (·symbol_lift _ _) (map term_lift [inr a; inr b; inr c; inr d]))
```

The extraction mechanism, applied to this case, would yield:

derives g [inl X; inl Y; inl Z; inr d] (map term_lift [inr a; inr b; inr c; inr d]) which is exactly the expected result $(XYZd\Rightarrow_q^*abcd)$.

8. Related Work

Context-free language theory formalization is a relatively new area of research, with some results already obtained with a diversity of proof assistants, including Coq, HOL4 and Agda. Most of the effort started in 2010 and have been devoted to the certification and validation of parser generators. Examples of this are the works of Koprowski and Binsztok (using Coq, [8]), Ridge (using HOL4, [12]), Jourdan, Pottier and Leroy (using Coq, [7]) and, more recently, Firsov and Uustalu (in Agda, [5]).

On the more theoretical side, on which the present work should be considered, Norrish and Barthwal published on general context-free language theory formalization using the using HOL4 proof assistant [1, 2, 3], including the existence of normal forms for grammars, pushdown automata and closure properties. Recently, Firsov and Uustalu proved the existence of a Chomsky Normal Form grammar for every general context-free grammar, using the Agda proof assistant [6].

It can thus be noted that so far apparently no formalization has been done in Coq for results not related directly to parsing and parser verification (except in HOL4 and Agda), and that this constitutes an important motivation for the present work, mainly due to the increasing usage and importance of Coq in different areas and communities. Specifically, the formalization done by Norrish and Barthwal in HOL4 is quite comprehensive and extends our work with the Greibach Normal Form and pushdown automata and its relation to context-free

grammars. It does not include, however, a proof of either the decidability of the membership problem or the Pumping Lemma for context-free languages, which are objectives of the present work. The formalization by Firsov and Uustalu in Agda comprises basically the existence of a Chomsky Normal Form, and formalizes the elimination of empty and unit rules, but not elimination of useless and inaccessible symbols.

9. Conclusions

All important objects related with context-free grammars have been properly represented and different grammar manipulation strategies were formalized. Proofs of their correctness were successfully constructed. The proofs of all lemmas and theorems presented in this article have been formalized in Coq and comprise approximately 18,000 lines of scripts. This number can be explained for the following reasons:

- 1. The style adopted for writing the scripts: for the sake of clarity, each tactic is placed in its own line, despite the possibility of combining several tactics in the same line. Also, bullets (for structuring the code) were used as much as possible and the sequence tactical (using the semicolon symbol) was avoided at all. This duplicates parts of the code but has the advantage of keeping the static structure of the script related to its dynamic behaviour, which favors legibility and maintenance.
- 2. The formalization includes not only the main theorems described here, but also an extensive library of other fundamental and auxiliary lemmas on context-free grammars and derivations, which have been used to obtain the main results presented here, were used in the previously obtained results and will be used in future developments.
 - The results presented in this paper are fundamental to context-free language theory. They create an adequate framework in which to pursue further results, including a proof of the decidability of the membership problem and a proof of the Pumping Lemma for context-free languages.

Bibliography

References

- [1] Barthwal, A., Norrish, M., August 2010. A formalisation of the normal forms of context-free grammars in HOL4. In: Dawar, A., Veith, H. (Eds.), Computer Science Logic, 24th International Workshop, CSL 2010, 19th Annual Conference of the EACSL, Brno, Czech Republic, August 23–27, 2010. Proceedings. Vol. 6247 of Lecture Notes in Computer Science. Springer, pp. 95–109.
- [2] Barthwal, A., Norrish, M., 2010. Mechanisation of PDA and grammar equivalence for context-free languages. In: Dawar, A., de Queiroz, R. J.

- G. B. (Eds.), Logic, Language, Information and Computation, 17th International Workshop, WoLLIC 2010. Vol. 6188 of Lecture Notes in Computer Science. pp. 125–135.
- [3] Barthwal, A., Norrish, M., 2014. A mechanisation of some context-free language theory in HOL4. Journal of Computer and System Sciences (WoL-LIC 2010 Special Issue, A. Dawar and R. de Queiroz, eds.) 80 (2), 346 – 362. URL http://www.sciencedirect.com/science/article/pii/S0022000013000925
- [4] Bertot, Y., Castéran, P., 2004. Interactive Theorem Proving and Program Development. Springer.
- [5] Firsov, D., Uustalu, T., 2014. Certified {CYK} parsing of context-free languages. Journal of Logical and Algebraic Methods in Programming 83 (5-6), 459 468, the 24th Nordic Workshop on Programming Theory (NWPT 2012).
 URL http://www.sciencedirect.com/science/article/pii/S2352220814000601
- [6] Firsov, D., Uustalu, T., 2015. Certified normalization of context-free grammars. In: Proceedings of the 2015 Conference on Certified Programs and Proofs. CPP '15. ACM, New York, NY, USA, pp. 167–174. URL http://doi.acm.org/10.1145/2676724.2693177
- [7] Jourdan, J.-H., Pottier, F., Leroy, X., 2012. Validating LR(1) parsers. In: Proceedings of the 21st European Conference on Programming Languages and Systems. ESOP'12. Springer-Verlag, Berlin, Heidelberg, pp. 397–416. URL http://dx.doi.org/10.1007/978-3-642-28869-2_20
- [8] Koprowski, A., Binsztok, H., 2010. TRX: A formally verified parser interpreter. In: Proceedings of the 19th European Conference on Programming Languages and Systems. ESOP'10. Springer-Verlag, Berlin, Heidelberg, pp. 345–365.
 - URL http://dx.doi.org/10.1007/978-3-642-11957-6_19
- [9] Ramos, M. V. M., de Queiroz, R. J. G. B., Jun. 2015. Formalization of closure properties for context-free grammars. ArXiv e-prints. URL http://arxiv.org/abs/1506.03428
- [10] Ramos, M. V. M., de Queiroz, R. J. G. B., 2015. Formalization of simplification for context-free grammars. In: Preliminary Proceedings of the 10th Workshop on Logical and Semantic Frameworks, with Applications. LSFA'15.
- [11] Ramos, M. V. M., Neto, J. J., Vega, I. S., 2009. Linguagens Formais: Teoria Modelagem e Implementação. Bookman.
- [12] Ridge, T., 2011. Simple, functional, sound and complete parsing for all context-free grammars. In: CPP. pp. 103–118.
- [13] Sudkamp, T. A., 2006. Languages and Machines, 3rd Edition. Addison-Wesley.