

Bar-Hillel Theorem Mechanization in Coq

Closure of Context-Free Languages Under Intersection With Regular Languages is Mechanized in Coq

Sergey Bozhko

Doctoral Student

Max Planck Institute for Software
Systems (MPI-SWS)
Saarbrücken, Germany
sbozhko@mpi-sws.com

Leyla Khatbullina

Student

Department of Computer Science
and Technology
St.Petersburg Electrotechnical
University “LETI”
St.Petersburg, Russia
leila.xr@gmail.com

Semyon Grigorev

Associate Professor

St.Petersburg State University
St.Petersburg, Russia
s.v.grigoriev@spbu.ru
Researcher
JetBrains Research
St.Petersburg, Russia
semen.grigorev@jetbrains.com

Abstract

Formal language theory has a deep connection with such areas as static code analysis, graph database querying, formal verification, and compressed data processing. Many application problems can be formulated in terms of languages intersection. The Bar-Hillel theorem states that context-free languages are closed under intersection with a regular set. This theorem has a constructive proof and thus provides a formal justification of correctness of the algorithms for applications mentioned above. Mechanization of the Bar-Hillel theorem, therefore, is both a fundamental result of formal language theory and a basis for the certified implementation of the algorithms for applications. In this work, we present the mechanized proof of the Bar-Hillel theorem in Coq. We generalize results of Gert Smolka and Jana Hofmann and use them as the base for our work.

Keywords Formal languages, Coq, Bar-Hillel Theorem, Closure, Intersection, Regular Language, Context-free Language

1 Introduction

Formal language theory has deep connection with different areas such as static code analysis [25, 29, 34, 35, 37–39], graph database querying [19, 20, 23, 40], formal verification [?], and others. One of the most frequent uses is to formulate a problem in terms of languages intersection. In verification, one language can serve as a model of a program and another language describe undesirable behaviors. When intersection of these two languages is not empty, one can conclude that the program is incorrect. Usually, the only concern is the decidability of the languages intersection emptiness problem. But in some cases a constructive representation of the intersection may prove useful. This is the case,

for example, when the intersection of the languages models graph querying: a language produced by intersection is a query result and to be able to process it, one needs the appropriate representation of the intersection result.

Let us consider several applications starting with the user input validation. The problem is to check if the input provided by the user is correct with respect to some validation template such as a regular expression for e-mail validation. User input can be represented as a one word language. The intersection of such a language with the language specifying the validation template is either empty or contains the only string: the user input. If the intersection is empty, then the input should be rejected.

Checking that a program is syntactically correct is another example. The AST for the program (or lack thereof) is just a constructive representation of the intersection of the one word language (the program) and the programming language itself.

Graph database regular querying serves as an example of the intersection of two regular languages [1, 2, 23]. Next and one of the most comprehensive cases with decidable emptiness problem is an intersection of a regular language with a context-free language. This case is relevant for program analysis [35, 37, 38], graph analysis [17, 20, 40], context-free compressed data processing [26], and verification [?]. The constructive intersection representation in these applications is helpful for further analysis.

Intersection of some classes of languages is not generally decidable. For example, intersection of the linear conjunctive and the regular languages, used in the static code analysis [39], is undecidable while multiple context-free languages is closed under intersection with regular languages and emptiness problem for MCFLs is decidable [?]. Is it possible to express any useful properties in terms of regular and multiple context-free languages intersection? This question is beyond the scope of this paper but provides a good reason

Conference’17, July 2017, Washington, DC, USA

2018. ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

for future research in this area. In this paper we focus on the intersection of regular and context-free languages.

Some applications mentioned above require certifications. For verification this requirement is evident. For databases it is necessary to reason about security aspects and, thus, we should create certified solutions for query executing. Certified parsing may be critical for Web [?], as well as certified regular expressions for input validation. As a result, there is a big number of papers focusing on regular expressions mechanization and certification [?], and a number on certified parsers [?]. On the other hand, mechanization (formalization) is important by itself as theoretical results mechanization and verification, and there is a lot of work done on formal languages theory mechanization [?]. Also it is desirable to have a base to reason about parsing algorithms and other problems on languages intersection.

Context-free languages are closed under intersection with regular languages. It is stated as the Bar-Hillel theorem [3] which provides a constructive proof and a construction for the resulting language description. We believe that the mechanization of the Bar-Hillel theorem is a good starting point for certified application development and since it is one of fundamental theorems, it is an important part of formal language theory mechanization. And this work aims to provide such mechanization in Coq.

Our current work is a first step: we provide mechanization of theoretical results on context-free and regular languages intersection. We choose the result of Smolka and !!! on context-free languages mechanization [?] as a base for our work. The main contribution of this paper may be summarized as follows.

- We provide the constructive proof of the Bar-Hillel theorem in Coq.
- We generalize the CFL results of Smolka: terminals is abstract types....
- All code is published on GitHub: https://github.com/YaccConstructor/YC_in_Coq.

This work is organized as follows. In the section ?? we formulate Bar-Hillel theorem and provide the sketch of its proof. The next part is a brief discussion of the Chomsky normal form in section ?. After that we describe our solution in the section ?. This description is split into steps with respect to provided sketch and contains basic definitions, Smolka results generalization, handling of trivial cases, and steps summarization as a final proof. Finally, we discuss related works in the section ?? and conclude with the discussion of the presented work and possible directions for future research in the section ?.

2 Bar-Hillel Theorem

In this section we provide the Bar-Hillel theorem and sketch the proof which we use as base of our work. Also we provide

some additional lemmas which are used in the proof of the main theorem.

Lemma 2.1. *If L is a context free language and $\varepsilon \notin L$ then there is a grammar in Chomsky Normal Form that generates L .*

Lemma 2.2. *If $L \neq \emptyset$ and L is regular then L is the union of regular language A_1, \dots, A_n where each A_i is accepted by a DFA with exactly one final state.*

Theorem 2.3 (Bar-Hillel Theorem). *If L_1 is a context free language and L_2 is a regular language then $L_1 \cap L_2$ is context free.*

Sketch of the proof.

1. By lemma 2.1 we can assume that there is a context-free grammar G_{CNF} in Chomsky normal form, such that $L(G_{CNF}) = L_1$
2. By lemma 2.2 we can assume that there is a set of regular languages $\{A_1 \dots A_n\}$ where each A_i is recognized by a DFA with exactly one final state and $L_2 = A_1 \cup \dots \cup A_n$
3. For each A_i we can explicitly define a (?) grammar of the intersection: $L(G_{CNF}) \cap A_i$
4. Finally, we join them together with the (?) operation of union

3 The Chomsky Normal Form

The important part of our proof is that any context-free language can be described with grammar in Chomsky Normal Form (CNF) or, equally, any context-free grammar can be converted to the grammar in CNF which specifies the same language. Let us recall the definition of CNF and the algorithm for CFG to CNF conversion.

Definition 3.1 (Chomsky Normal Form). Context-free grammar is in CNF if:

- start nonterminal does not occur in the right-hand side of rules,
- all rules are of the form: $N_i \rightarrow t_i$, $N_i \rightarrow N_j N_k$ or $S \rightarrow \varepsilon$ where N_i, N_j, N_k are nonterminals, t_i is a terminal and S is the start nonterminal.

Transformation algorithm has the following steps.

1. Eliminate the start symbol from the right-hand sides of the rules.
2. Eliminate rules with nonsolary terminals.
3. Eliminate rules which right-hand side contains more than two nonterminals.
4. Delete ε -rules.
5. Eliminate unit rules.

As far as Bar-Hillel theorem operates with arbitrary context-free languages, but the proof requires grammar in nCNF, it is necessary to implement a certified algorithm for conversion of arbitrary CFG to CNF. We wanted to reuse existing

proof of conversion of arbitrary context-free grammar to CNF. We chose Smolka version which proves conversion to CNF correctness in the following way.

CNF!!!

Listing 1. TODO

4 Bar-Hillel Theorem Mechanization in Coq

In this section we describe in detail all the fundamental parts of the proof. Also in this section, we briefly describe motivation to use the chosen definitions. In addition, we discuss the advantages and disadvantages of using of third-party proofs.

Overall goal of this section is to provide step-by-step algorithm of constructing the CNF grammar of the intersection of two languages. Final formulation of the obtained theorem can be found in the last subsection.

All code are published on GitHub ¹.

4.1 Smolka's code generalization

In this section, we describe the exact steps taken to use the proof of TODO:Smolka's theorem in the proof of this article's theorem.

A substantial part of this proof relies on the work of TODO:Smc. From this work(?) many definitions and theorems were taken. Namely, the definition of a grammar, definitions of a derivation in grammar, some auxiliary lemmas about the decidability of properties of grammar/derivation, we also use the theorem that states that there always exists the transformation from context-free grammar to grammar in Chomsky Normal Form (CNF).

However, this proof had one major flaw that we needed to fix. One could define a terminal symbol as in inductive type over natural numbers[TODO].

```
Inductive ter : Type := | T : nat -> ter.
```

Listing 2. TODO

That is how it was done in TODO:Smolka. However for purposes of our proof, we need to consider nonterminals over the alphabet of triples. Therefore, it was decided to add polymorphism over the target alphabet. Namely, let Tt and Vt be types with decidable relation of equality, then we can define the types of terminal and nonterminal over alphabets Tt and Vt respectively as follows (???):

The proof of Smolka has a clear structure, therefore only part of the proof where the use of natural numbers was essential has become incorrect. One of the grammar transformations (namely deletion of long rules) requires the creation

¹https://github.com/YaccConstructor/YC_in_Coq

```
Inductive ter : Type := | T : Tt -> ter.
```

```
Inductive var : Type := | V : Vt -> var.
```

Listing 3. TODO

of many new non-terminals. In the original proof for this purpose, the maximum over non-terminals included in the grammar was used. However, it is impossible for an arbitrary type.

To tackle this problem we introduce an additional assumption on alphabet types for terminals and nonterminals. We require an existence of the bijection between natural numbers and alphabet of terminals as well as nonterminals.

Another difficulty is that the original work defines grammar as a list of rules (without a distinct starting nonterminal). Thus, in order to define the language that is defined by a grammar, one needs to specify the grammar and a starting terminal. This leads to the fact that the theorem about the equivalence of a CF grammar and the corresponding CNF grammar isn't formulated in the most general way, namely it guarantees equivalence only for non-empty words.

```
Lemma language_normal_form
  (G:grammar) (A: var) (u: word):
  u <> [] ->
  (language G A u <->
   language (normalize G) A u).
```

Listing 4. TODO, CHECK

Changes in the definition of grammar or language would lead to significant code corrections. However, the question of whether the empty word is derivable is decidable for both the CF grammar and the DFA. Therefore, it is possible to simply consider two cases (1) when the empty word is derivable in the grammar and (2) when the empty word is not derivable.

4.2 Part ... derivation and so on

In this section, we introduce the basic definitions used in the article.

We define a symbol is either a terminal or a nonterminal.

```
Inductive symbol : Type :=
  | Ts : ter -> symbol
  | Vs : var -> symbol.
```

Listing 5. TODO

Next we define a word and a phrase as lists of terminals and symbols respectively.

The notion of nonterminal doesn't make sense for DFA, but in order to construct derivation in grammar we need

Definition word := list ter.
Definition phrase := list symbol.

Listing 6. TODO

to use nonterminal in intermediate states. For phrases, we introduce a predicate that defines whenever a phrase consists of only terminals. And if so, the phrase it can be safely converted to the corresponding word.

We inheriting the definition of CFG from [Smplka] paper. Rule is defined as a pair of a nonterminal and a list of symbols. Grammar is a list of rules.

Inductive rule : Type :=
 | R : var -> phrase -> rule.

Definition grammar := list rule.

Listing 7. TODO

An important step towards the definition of a language (?) governed (formed?)(?) by a grammar is the definition of derivability. Having $der(G, A, p)$ — means that phrase p is derivable in grammar G starting from(?) nonterminal A .

Inductive der (G : grammar)
 (A : var) : phrase -> Prop :=
 | vDer : der G A [Vs A]
 | rDer l : (R A l) e1 G -> der G A l
 | replN B u w v :
 der G A (u ++ [Vs B] ++ w) ->
 der G B v -> der G A (u ++ v ++ w).

Listing 8. TODO

Proof of TODO requires grammar to be in CNF. We used statement that every grammar in convertible into CNF from TODO:Smolka work.

... ..

4.3 General scheme of the proof

General scheme of our proof is based on constructive proof presented by [?]. In the following subsections the main steps of the proof are presented. Overall, we will adhere to the following plan.

1. First we consider trivial case, when DFA has no state (TODO: del this?)
2. Every CF language can be converted to CNF
3. Every DFA can be presented as an union of DFAs with single final state
4. Intersecting grammar in CNF with DFA with one final state
5. Proving than union of CF languages is CF language

4.4 Part one: trivial case

(TODO: del?)

4.5 Part two: regular language and automata

In this section we describe definitions of DFA and DFA with exactly one final state, we also present function that converts any DFA to a set of DFA with one final state and lemma that states this split is well-defined(?).

We assume that regular language by definition is described by DFA. As the definition of an DFA, we have chosen a general definition, which does not impose any restrictions on the type of input symbols and the number of states. Thus, in our case, the DFA is a 5-tuple, (1) a state type, (2) a type of input symbols, (3) a start state, (4) a transition function, and (5) a list of final states.

Context {State T: Type}.
Record dfa: Type :=
 mkDfa {
 start: State;
 final: list State;
 next: State -> (@ter T) -> State;
 }.

Listing 9. TODO

Next we define a function that would evaluate the final state of the automaton if it starts from state s and receives a word w .

Fixpoint final_state
 (next_d: dfa_rule)
 (s: State)
 (w: word): State :=
 match w with
 | nil => s
 | h :: t => final_state next_d
 (next_d s h)
 t
 end.

Listing 10. TODO

We say that the automaton accepts a word w being in state s if the function $[final_state_sw]$ ends in one of the final states. Finally, we say that an automaton accepts a word w , if the DFA starts from the initial state and ends in one of the final states.

In order to define the DFA with exactly one final state, it is necessary to replace the list of final states by one final state in the definition of an(?) ordinary DFA. The definitions of "accepts" and "dfa_language" vary slightly.

Similarly, we can define functions $s_accepts$ and $s_dfa_language$ for sDFA. Since in this case, there is only one final state, to


```

441 Record s_dfa : Type :=
442   s_mkDfa {
443     s_start: State;
444     s_final: State;
445     s_next: State -> (@ter T) -> State;
446   }.

```

Listing 11. TODO

define function *s_accepts* it is enough to check the state in which the automaton stopped with the finite state. The function *s_dfa_language* repeats the function *dfa_language*, except that the function must use *s_accepts* instead of *accepts*.

Now we can define a function that converts an ordinary DFA into a set of DFAs with exactly one final state. Let *d* be a dfa. Then the list of its final states is known. For each such state, one can construct a copy of the original dfa, but with one current final state.

```

461 Fixpoint split_dfa_list
462   (st_d : State)
463   (next_d : dfa_rule)
464   (f_list : list State): list (s_dfa) :=
465   match f_list with
466   | nil => nil
467   | h :: t => (s_mkDfa st_d h next_d)
468               :: split_dfa_list st_d next_d t
469   end.

```

```

471 Definition split_dfa (d: dfa) :=
472   split_dfa_list (start d) (next d) (final d).

```

Listing 12. TODO

```

477 Lemma correct_split:
478   forall dfa w,
479     dfa_language dfa w <=>
480     exists sdfa,
481       In sdfa (split_dfa dfa) /\
482       s_dfa_language sdfa w.

```

Listing 13. TODO

We prove theorem that the function of splitting preserves the language.

Theorem 4.1. *Let dfa be an arbitrary dfa and w be a word. Then the fact that dfa accepts w implies that there exists a single-state dfa s_dfa, such that s_dfa ∈ split_dfa(dfa). And vice versa, For any s_dfa ∈ split_dfa(dfa) the fact that s_dfa accepts a word w implies that dfa also accepts w.*

Proof. Let us divide the proof into two parts. (1) Suppose *dfa* accepts *w*. Then we prove that there exists a single-state dfa *s_dfa*, such that *s_dfa* ∈ *split_dfa(dfa)*. Let *finals* be the set of final states of *dfa*. We carry out the proof by induction on *finals*. Base step: *finals* = [::]. Trivial by contradiction (DFA with no final state cannot accept a word). Induction step: *finals* = *a* :: *old_finals* and the statement holds for *old_finals*. Since *dfa* accepts *w*, it either ends up in *a*, or in one of the state from *old_finals*. If *dfa* ends up in *a*, then we simply choose an automaton with the final state that is equal to *a*. Such an automaton exists, since now the list of final states also contains *a*. On the other hand, if *dfa* ends up in one of the state from *old_finals*, then we can apply induction hypothesis.

(2) Similarly for the opposite direction. Assume that there exists an automaton with exactly one final state from *split_dfa(dfa)* that accepts *w*. Then we prove that *dfa* also accepts *w*. Let *finals* be the set of final states of *dfa*. We carry out the proof by induction on *finals*. Base step: *finals* = [::]. Trivial by contradiction. Induction step: *finals* = *a* :: *old_finals* and the statement holds for *old_finals*. We know that one of the DFAs from *split_dfa(dfa)* accepts *w*, its final state either is equal to *a*, or lies in *old_finals*. If the final state is equal to *a*, then *dfa* also ends up in state *a*. On the other hand, if final state lies in *old_finals*, then we can apply induction hypothesis.

4.6 Part :: Chomsky induction

In this section, we introduce the notion of Chomsky induction.

Naturally many statements about properties of language's words can be proved by induction over derivation structure. Unfortunately, grammar can derive phrase as an intermediate step, but DFA supposed to work only with words, so we can't simply apply induction over derivation structure. To tackle this problem we create custom induction principle for grammars in CNF.

As one might notice, TODO

The main point is that if we have a grammar in CNF, we can always divide the word into two parts, each of which is derived only from one nonterminal. Note that if we naively take a step back, we can get nonterminal in the middle of the word. Such a situation will not make any sense for DFA.

With induction we always work with subtrees that describes some part of word. Here is a picture of subtree describing intuition behind the Chomsky induction.

TODO: add picture

TODO: add Lemma derivability_backward_step.

(TODO: lemma) More formally: Let *G* be a grammar in CNF. Consider an arbitrary nonterminal *N* ∈ *G* and phrase which consists only on terminals *w*. If *w* is derivable from *N* and |*w*| ≥ 2, then there exist (TODO:s) two nonterminals *N*₁, *N*₂ and subphrases of *w* — *w*₁, *w*₂ such that: *N* →

$N_1N_2 \in G$, $der(N_1, w_1)$, $der(N_2, w_2)$, $|w_1| \geq 1$, $|w_2| \geq 1$ and $w_1 + w_2 = w$.

(TODO: fix) **Proof.** The proof heavily uses the fact that grammar G is in Chomsky Normal Form. We apply the hypothesis "syntactic analysis is possible". After application, we get the fact that word w is either an RHS of a rule of grammar G , or there is a phrase phr , such that (1) word w is derivable from phrase phr and (2) there exists a non-terminal N such that $N \rightarrow phr$ in G . The first case we finish with the proof by contradiction since the grammar is in CNF and there might be only a single terminal in an RHS (by assumption we have $|w| \geq 2$). On the other hand, if there is an intermediate phrase that was obtained by applying a rule, then the phrase has form N_1N_2 , since it also derived by rule in normal form. Finally, now we need to prove that both of this nonterminals has a non-empty contribution to word w . This is also true since it is impossible to derive empty word in CNF grammar (see ...).

(TODO: lemma) Let G be a grammar in CNF. And P be a predicate on nonterminals and phrases (i.e. $P : var \rightarrow phrase \rightarrow Prop$). Let's also assume that the following two hypotheses are satisfied: (1) for every terminal production (i.e. in the form $N \rightarrow a$) of grammar G , $P(r, [Tsr])$ holds and (2) for every $N, N_1, N_2 \in G$ and two phrases that consist only of terminals w_1, w_2 , if $P(N_1, w_1)$, $P(N_2, w_2)$, $der(G, N_1, w_1)$ and $der(G, N_2, w_2)$ then $P(N, w_1 + w_2)$. Then for any non-terminal N and any phrase consisting only of terminals w , the fact that w is derivable from N implies $P(N, w)$.

(TODO: fix) **Proof.** Let n be an upper bound of the length of word w . We carry out the proof by induction on n . Base case: $n = 0$. Proof by contradiction. $|w| \leq 0$ implies that w is empty. But an empty word cannot be derived in CNF grammar (see ...). Induction step: $|w| \leq n + 1$. This fact is equivalent to the following: $|w| = n + 1$ or $|w| < n$. In case of $|w| < n$ we use the induction hypothesis. Next we consider two new cases, either $|w| = 1$, or $1 < |w| = n + 1$. In the first case, it is clear that this is possible only if there is a production $N \rightarrow w$, which means you can apply assumption (1). If the word is longer than 1, then we apply the previous lemma, after that we can conclude that $\exists w_1w_2, w = w_1 + w_2$. After that, one need to apply assumption (2). We subgoals that are guaranteed by the lemma And for shorter words, we apply the induction hypothesis.

TODO: add some text

4.7 Part ..: intersection

Since we already have lemmas about the transformation of a grammar to CNF and the transformation a DFA to a DFA with exactly one state, further we assume that we have (1) DFA with exactly one final state — dfa and (2) grammar in CNF — G . In this section, we describe the proof of the lemma that states that for any grammar in CNF and any automaton with exactly one state there is the intersection grammar.

4.7.1 Function

Next we present adaptation of the algorithm given in [].

Let G_{INT} be the grammar of intersection. In G_{INT} nonterminals presented as triples $(from \times var \times to)$ where $from$ and to are states of dfa , and var is a nonterminal of(in?) G .

Since G is a grammar in CNF, it has only two type of productions: (1) $N \rightarrow a$ and (2) $N \rightarrow N_1N_2$, where N, N_1, N_2 are nonterminals and a is a terminal.

For every production $N \rightarrow N_1N_2$ in G we generate a set of productions of the form $(from, N, to) \rightarrow (from, N_1, m)(m, N_2, to)$ where: $from, m, to$ — goes through all dfa states.

```
Definition convert_nonterm_rule_2
  (r r1 r2: _)
  (state1 state2 : _) :=
  map (fun s3 => R (V (s1, r, s3))
      [Vs (V (s1, r1, s2));
       Vs (V (s2, r2, s3))])
  list_of_states.
```

```
Definition convert_nonterm_rule_1
  (r r1 r2: _)
  (s1 : _) :=
  flat_map (convert_nonterm_rule_2 r r1 r2 s1)
  list_of_states.
```

```
Definition convert_nonterm_rule (r r1 r2: _) :=
  flat_map (convert_nonterm_rule_1 r r1 r2)
  list_of_states.
```

Listing 14. TODO

For every production of the form $N \rightarrow a$ we add a set of productions $(from, N, (dfa_step(from, a))) \rightarrow a$ where: $from$ — goes through all dfa states and $dfa_step(from, a)$ is the state in which the dfa appears after receiving terminal a in state $from$.

```
Definition convert_terminal_rule
  (next: _)
  (r: _)
  (t: _): list TripleRule :=
  map (fun s1 => R (V (s1, r, next s1 t))
      [Ts t])
  list_of_states.
```

Listing 15. TODO

Next we join the functions above to get a generic function that works for both types of productions. Note that since the grammar is in CNF, (?) the third alternative can never be the case.

Note that at this point we do not have any manipulations with starting rules. Nevertheless(?), the hypothesis of the

```

661 Definition convert_rule (next: _) (r: _ ) :=
662   match r with
663   | R r [Vs r1; Vs r2] =>
664     convert_nonterm_rule r r1 r2
665   | R r [Ts t] =>
666     convert_terminal_rule next r t
667   | _ => [] (* Never called *)
668   end.
669
670 Definition convert_rules
671   (rules: list rule) (next: _) : list rule :=
672   flat_map (convert_rule next) rules.
673
674 (* Maps grammar and s_dfa to grammar over triples *)
675 Definition convert_grammar grammar s_dfa :=
676   convert_rules grammar (s_next s_dfa).
677

```

Listing 16. TODO

uniqueness of the final state of the DFA, will help us unambiguously introduce the starting nonterminal of the grammar of intersection.

4.7.2 Correctness

In this subsection we present a high-level description of the proof about correctness of the intersection function.

In the interest of clarity of exposition, we skip some auxiliary lemmas, such as (TODO:fix) "we can get the initial grammar from the grammar of intersection by projecting the triples back to terminals/nonterminals". Also note that the grammar after the conversion remains in CFN. Since the transformation of rules does not change the structure of the rules, but only replaces one(??!!) terminals and nonterminals with others.

Next we prove the two main lemmas. Namely, the derivability in the initial grammar and the *s_dfa* implies the derivability in the grammar of intersection. And the other way around, the derivability in the grammar of intersection implies the derivability in the initial grammar and the *s_dfa*.

Let *G* be a grammar in CNF. In order to use Chomsky Induction we also assume that syntactic analysis is possible.

Theorem 4.2. *Let s_dfa be an arbitrary DFA, let r be a non-terminal of grammar G, let from and to be two states of the DFA. We also pick an arbitrary word — w. If in grammar G it is possible to derive w out of r and starting from the state from when w is received, the s_dfa ends up in state to, then word w is also derivable in grammar (convert_rules G next) from the nonterminal (V (from, r, to)).*

Proof. TODO. It would be logical to use induction on the derivation structure in grammar *G*. But as it was discussed earlier, this is not the case, otherwise we will get a phrase

(list of terminals and nonterminals) instead of a word. Therefore we should use another way to use induction. For grammar in chomsky normal form it is possible (see the chomsky_induction section). Roughly speaking, we can split the word into two subwords, can be derived from some nonterminal.

Let's apply chomsky induction principle with the following predicate *P*:

$$P := \text{fun } r \text{ phr} \Rightarrow \forall (\text{next} : \text{dfa_rule})(\text{fromto} : \text{DfaState}),$$

$$\text{final_statenext from}(to_word \text{ phr}) = to \rightarrow$$

$$\text{der}(\text{convert_rules G next})(V(\text{from}, r, to)) \text{ phr}.$$

Basically, predicate *P* is the property that we are trying to prove in the theorem. Chomsky Induction has 3 assumptions. (1) The phrase to which *P* is applied should consist of only non-terminals. We consider only words in this theorem, therefore after conversion of the word to the phrase, no terminals can appear in it. So, we do not violate this assumption. Moreover, there is a base of induction (2) in the form of a property for a terminal rule and (3) an induction step for a non-terminal rule. Both statements can be proved by induction on the number of rules in the grammar *G* in combination with a simple calculation of the functions *convert_terminal_rule* and *convert_nonterm_rule* for terminal and non-terminal rules, respectively.

On the other side. Now we need to prove the theorems of the form "if it is derivable in the grammar of triples, then it is derivable in the automaton and in ordinary grammar"

We start with the DFA. TODO: Theorem. Let *from* and *to* be states of the automaton, *var* be an arbitrary non-terminal grammar of *G*. We prove that If a word *w* is derived from the non-terminal (*from*, *var*, *to*) in the grammar (*convert_rules G*), then the automaton, starting from the state *from* at the input *w* stops in state *to*. Proof. Like last time, we use the principle of Chomsky Induction. We apply the induction with the following parameter *P*:

$$(P := \text{fun } tr_n \text{ onphr} \Rightarrow \text{final_statenext}(\text{first3tr_non})(to_word \text{ phr}) = \text{third3tr_non})$$

Note that in this case, one need to use projections from triple-non-terminals to "single" non-terminals. Induction is carried out in the grammar of triples, but the property says(??) about the automaton. However(?), this property can also be expressed in terms of triples-non-terminals. But in general, as last time, *P* is the statement we want to prove. After applying the induction principle, it remains only to prove the fidelity of the assumptions. First of all, since *w* is a word, converting it into a phrase does not add any non-terminal. Next, one need to show that both assumption of the induction principle hold. TODO: Ill finish it tomorrow

Further we prove the theorem for grammar.

Theorem. Let *from* and *to* be the states of the automaton, let *var* be an arbitrary non-terminal of grammar *G*. We prove that if a word *w* is derivable from non-terminal (*from*, *var*, *to*) in the grammar (*convert_rules G*), then *w* is also derivable

in grammar G from nonterminal var . Proof. We again prove the theorem using Chomsky induction with the following predicate P :

$$(P := \text{funrphr} \Rightarrow \text{derG}(\text{snd3}(\text{unVarr}))(phr))$$

Again, note that induction is carried out in the grammar of triples, but the property is talking about the "unit" grammar, therefore we use projections from triples to non-triples.

This can be done by an "inverse" calculation of functions `convert_terminal_rule` and `convert_nonterm_rule` functions. TODO:

Well, in the end one need to combine both theorems to get full equivalence. On this, the correctness of the intersection is proved.

4.8 Part :: union

After the previous step, we have a list of grammars of CF languages, in this section, we provide a function by which we construct a grammar of the union of languages.

For this, we need nonterminals from every language to be from different nonintersecting sets. To achieve this we add labels to nonterminals. Thus, each grammar of the union would have its own unique ID number, all nonterminals within one grammar will have the same ID which coincides with the ID of a grammar. In addition, it is necessary to introduce a new starting nonterminal of the union.

```
Inductive labeled_Vt : Type :=
| start : labeled_Vt
| lV : nat -> Vt -> labeled_Vt.
```

```
Definition label_var (label: nat)
(v: @var Vt): @var
labeled_Vt :=
V (lV label v).
```

Listing 17. TODO

Construction of new grammar is quite simple. The function that constructs the union grammar takes a list of grammars, then, it (1) splits the list into head $[h]$ and tail $[tl]$, (2) labels $[length\ tl]$ to h , (3) adds a new rule from the start nonterminal of the union to the start nonterminal of the grammar $[h]$, finally (4) the function is recursively called on the tail $[tl]$ of the list.

4.8.1 Equivalence proof

In this section, we prove that function `grammar_union` constructs a correct grammar of union language indeed. Namely, we prove the following theorem.

Theorem 4.3. *Let grammars be a sequence of pairs of starting nonterminals and grammars. Then for any word w , the fact that w belongs to language of union is equivalent to the*

```
Definition label_grammar label grammar := ...
```

```
Definition label_grammar_and_add_start_rule
label
grammar :=
let '(st, gr) := grammar in
(R (V start) [Vs (V (lV label st))])
:: label_grammar label gr.
```

```
Fixpoint grammar_union
(grammars : seq (@var Vt * (@grammar Tt Vt)))
: @grammar
Tt
labeled_Vt :=
match grammars with
| [] => []
| (g::t) =>
label_grammar_and_add_start_rule
(length t)
g ++ (grammar_union t)
end.
```

Listing 18. TODO

fact that there exists a grammar $(st, gr) \in grammars$ such that w belongs to language generated by (st, gr) .

```
Variable grammars: seq (var * grammar).
```

```
Theorem correct_union:
forall word,
language (grammar_union grammars)
(V (start Vt)) (to_phrase word) <=>
exists s_l,
language (snd s_l) (fst s_l)
(to_phrase word) /\
In s_l grammars.
```

Listing 19. TODO

Proof of theorem 4.3. Since the statement is formulated as an equivalence, we divide the proof into two parts:

1. If w belongs to the union language, then w belongs to one of the initial language.

From an auxiliary lemma, we know that either (1) the phrase is equal to the starting nonterminal or (2) there exists a grammar G from the grammars-list such that TODO Let us prove that this is the grammar we are interested in. Since we consider the word, it cannot be a starting non-terminal. So this might be only the second case. According to another lemma, if the output doesn't start from the starting non-terminal, then it cannot appear in this derivation. So all the rules that

use the starting non-terminal can be safely (for this derivation) removed from the grammar. There is a lemma that says, for a derivation that starts from a nonterminal labeled with x , can not contain any nonterminals with a label other than x . Therefore, for this derivation, one can ignore the rules with other(?) labels. These two grammars are identical, but one of them is labeled and the other is not. It is clear that if there exists a bijection between nonterminals the set of derivable words doesn't change.

2. If w belongs to one of the initial language, then w belongs to the union language.

In this case, one explicitly specify the corresponding derivation in the union-grammar. The labeling function is arranged in such a way that knowing the place of a certain grammar in the list of grammars, one can calculate the exact number that will be assigned to this grammar as a label. After that proof can be finished in two steps. Step 1: one need to apply the rule from the starting non-terminal of the union-grammar to the starting non-terminal of the initial grammar. And step 2: one should use the fact that derivation in initial grammar and labeled grammar are equivalent.

4.9 Part N: taking all parts together

Theorem 4.4. *For any two decidable types Terminal and Non-terminal for type of terminals and nonterminals correspondingly. If there exists bijection from Nonterminal to \mathbb{N} and syntactic analysis in the sense of definition TODO is possible, then for any DFA dfa that accepts Terminal and any grammar G , there exists the grammar of intersection G_{int} .*

Proof. Let NG be the grammar in CNF obtained after applying the algorithm of TODO:Smolka. Let $sdfas$ be the list of DFAs with exactly one final state obtained after splitting dfa . Since we now have NG in CNF and the list of DFAs with one state, we can compute a list of their intersections. I.e. we intersect each of the $sdfa$ of the list with grammar NG . Next, we find the union of the languages.

Next, we divide the proof into two branches. We check whether the empty word is derivable in dfa and in grammar G . (1) If so, we add one more rule to the language of the union ($S \rightarrow \epsilon$). (2) If not, we add nothing.

Now for the cases (1) and (2) we will prove that G_{int} is the grammar of the intersection indeed. That is, if a word is derivable in dfa and G , then it must also be derivable in G_{int} . And vice versa.

For branch (1) we carry out the proof in 2 stages.

a) Consider the case when w is an empty word. By assumption, we already know that the empty word is derivable in the DFA and in the grammar. But we also know that we have added a rule from the starting terminal to the empty word to the grammar of intersection. So, if w is an empty word it derivable in both cases. b) Let's now prove for the case when w is a non-empty word. We consistently modify

the premises and conclusion using theorems about equivalences. First we prove that the fact that w derivable in G and dfa implies the fact that w is also derivable in G_{int} . We can safely remove the rule $S \rightarrow \epsilon$, since we apply unification only to grammars in CNF (any grammar of intersection is in CNF), the epsilon rule cannot be used anywhere except the initial step. We know, since the word is accepted by dfa , then there is a DFA with one final state $sdfa$, which also accepts this word. So, we can safely replace dfa with $sdfa$. For grammar G there is an equivalent grammar NG in CNF. And since word w is not empty, we maintain equivalence. Let INT be a grammar of the intersection of $sdfa$ and NG . We can use theorem TODO to prove that it is a grammar of intersection of $sdfa$ and G indeed. But by the construction, such a grammar belongs to the union of languages G_{int} . this finishes inclusion of G and $sdfa$ to G_{int} . In the other direction: the fact that w is derivable in G_{int} implies that w is derivable in grammar G and is accepted by DFA dfa .

Grammar G_{int} consists of a union of the empty language and list languages of the intersection of some sDFA and grammar in CNF. We can safely remove an empty grammar since w is not an empty word and any grammar in the list of languages of intersection is in CNF. We know that since the word is accepted by G_{int} grammar, there is a grammar from the union of grammars in which this word is derivable. But this grammar is a grammar of intersection of some sDFA and grammar in CNF. But now we can use TODO in order to prove equivalence.

The second case is when the empty word is not derivable either in G or in dfa . For an empty word, one needs to prove that it is not derivable in the grammar of the intersection. And indeed. None of the grammar from the union has the rule $S \rightarrow \epsilon$. Next, one have to repeat what is discussed above, but without the additional steps about the empty language.

5 Related Works

There is a big number of works in mechanization of different parts of formal languages theory and certified implementations of parsing algorithms and algorithms for graph data base querying. These works use different tools, such as Coq, Agda, Isabelle/HOL, and aimed to different problems such as theory mechanization or executable algorithm certification. We discuss only small part which is close enough to the scope of this work.

5.1 Formal Language Theory in Coq

Huge amount of work was done by Ruy de Queiroz who formalize different parts of formal language theory, such as pumping lemma [33], context-free grammar simplification [30] and closure properties [32] in Coq. The work on

closure properties contains mechanization of such properties as closure under union, Kleene star, but it does not contain mechanization of intersection with regular language. All these results are summarized in [31].

Gert Smolka et.al. also provide big set of works on regular and context-free languages formalization in Coq [11, 12, 21, 22]. The work [21] contains certified transformation of arbitrary context-free grammar to Chomsky normal form which is required for our proof of the Bar-Hillel theorem. Initially we hope to use these both parts because Bar-Hillel theorem is about both context-free and regular languages, and it was the reason to choose results of Gert Smolka as base for our work. But works on regular languages and on context-free languages are independent and we face with problems of reusing and integration and in current proof we use only results on context-free languages.

5.2 Formal Language Theory in Other Languages

In the parallel with works in Coq there exist works on formal languages mechanization in other languages and tools such as Agda or Isabelle/HOL.

First part is works of Denis Firsov who implements in Agda some parts of formal language theory and parsing algorithms. CYK parsing algorithm [13, 15] and Chomsky Normal Form [16], and some results on regular languages [14].

Another part is formal language theory mechanization in Isabelle/HOL [4, 6, 7] by Aditi Barthwal and Michael Norrish. This work contains basic definitions and big number of theoretical results, such as Chomsky normal form and Greibach normal form for context-free grammars. As an application of mechanized theory authors provide certified implementation of SLR parsing algorithm [5].

5.3 Certified Algorithms

Additionally we want to mention some works on certified applied algorithms based on formal language theory. Certification are required in different areas for various reasons and it is a reason to work on theory mechanization.

The first area where languages intersection may be applied is language constrained path querying in structured data (for example in graphs or XML). There exist works on certification of core of XQuery [10]. XQuery is a W3C standard for path querying in XML, extended for graph querying. Another result is a work on certified Regular Datalog querying in Coq [9]. Inspired by these results, our work may be a base for certified context-free path querying algorithm.

Another area which is fastly growth is certified parsers and parser generators based on different algorithms for different language classes [8, 18, 24].

6 Conclusion

We present mechanized in Coq proof of the Bar-Hillel theorem — the fundamental theorem on closure of context-free

languages under intersection with regular set. By this we increase mechanized part of formal language theory and provide a base for reasoning about many applicative algorithms which are based on languages intersection. Also we generalize results of Gert Smolka and Jana Hofmann: generalized terminal alphabet. It makes previously existing results more flexible and ease for reusing. All results are published at GitHub and equipped with automatically generated documentation.

The first open question, and seems that very important question, is integration of our results with other results on formal languages theory mechanization in Coq. There are two independent sets of results in this area: works of Ruy de Queiroz and works of Gert Smolka. We use part of Smolka's results in our work, but even here we do not use existing results on regular languages. We think that theory mechanization should be unified and results should be generalized. We think that these and other related questions should be discussed in community.

One of direction of future research is mechanization of practical algorithms which are just implementation of Bar-Hillel theorem. For example, context-free path querying algorithm, based on CYK [20, 40] or even on GLL [36] parsing algorithm [17]. Final target here is certified algorithm for context-free constrained path querying for graph databases.

Yet another direction is mechanization of other problems on language intersection which can be useful for applications. For example, intersection of two context-free grammars one of which describes finite language [27, 28]. It may be useful for compressed data processing [?] or speech recognition [?]. And, of course, all these works should be done on the common base of mechanized theoretical results.

Acknowledgments

We are grateful to Ekaterina Verbitskaia for discussion, careful reading, and pointing out some mistakes. The research was supported by the Russian Science Foundation grant No. 18-11-00100 and by the grant from JetBrains Research.

References

- [1] Serge Abiteboul and Victor Vianu. 1999. Regular Path Queries with Constraints. *J. Comput. System Sci.* 58, 3 (1999), 428 – 452. <https://doi.org/10.1006/jcss.1999.1627>
- [2] Faisal Alkhateeb. 2008. *Querying RDF(S) with Regular Expressions*. Theses. Université Joseph-Fourier - Grenoble I. <https://tel.archives-ouvertes.fr/tel-00293206>
- [3] Yehoshua Bar-Hillel, Micha Perles, and Eli Shamir. 1961. On formal properties of simple phrase structure grammars. *Sprachtypologie und Universalienforschung* 14 (1961), 143–172.
- [4] Aditi Barthwal. 2010. *A formalisation of the theory of context-free languages in higher order logic*. Ph.D. Dissertation. College of Engineering & Computer Science, The Australian National University.
- [5] Aditi Barthwal and Michael Norrish. 2009. Verified, Executable Parsing. In *Programming Languages and Systems*, Giuseppe Castagna (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 160–174.

- [6] Aditi Barthwal and Michael Norrish. 2010. A formalisation of the normal forms of context-free grammars in HOL4. In *International Workshop on Computer Science Logic*. Springer, 95–109.
- [7] Aditi Barthwal and Michael Norrish. 2010. Mechanisation of PDA and Grammar Equivalence for Context-Free Languages. In *Logic, Language, Information and Computation*, Anuj Dawar and Ruy de Queiroz (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 125–135.
- [8] Jean-Philippe Bernardy and Patrik Jansson. 2016. Certified Context-Free Parsing: A formalisation of Valiant's Algorithm in Agda. *arXiv preprint arXiv:1601.07724* (2016).
- [9] Angela Bonifati, Stefania Dumbrava, and Emilio Jesús Gallego Arias. 2018. Certified Graph View Maintenance with Regular Datalog. *arXiv preprint arXiv:1804.10565* (2018).
- [10] James Cheney and Christian Urban. 2011. Mechanizing the Metatheory of mini-XQuery. In *Certified Programs and Proofs*, Jean-Pierre Jouannaud and Zhong Shao (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 280–295.
- [11] Christian Doczkal, Jan-Oliver Kaiser, and Gert Smolka. 2013. A constructive theory of regular languages in Coq. In *International Conference on Certified Programs and Proofs*. Springer, 82–97.
- [12] Christian Doczkal and Gert Smolka. 2018. Regular Language Representations in the Constructive Type Theory of Coq. *Journal of Automated Reasoning* 61, 1 (01 Jun 2018), 521–553. <https://doi.org/10.1007/s10817-018-9460-x>
- [13] Denis Firsov. 2016. Certification of Context-Free Grammar Algorithms. (2016).
- [14] Denis Firsov and Tarmo Uustalu. 2013. Certified Parsing of Regular Languages. In *Certified Programs and Proofs*, Georges Gonthier and Michael Norrish (Eds.). Springer International Publishing, Cham, 98–113.
- [15] Denis Firsov and Tarmo Uustalu. 2014. Certified CYK parsing of context-free languages. *Journal of Logical and Algebraic Methods in Programming* 83, 5-6 (2014), 459–468.
- [16] Denis Firsov and Tarmo Uustalu. 2015. Certified normalization of context-free grammars. In *Proceedings of the 2015 Conference on Certified Programs and Proofs*. ACM, 167–174.
- [17] Semyon Grigorev and Anastasiya Ragozina. 2016. Context-Free Path Querying with Structural Representation of Result. *arXiv preprint arXiv:1612.08872* (2016).
- [18] Jason Gross and Adam Chlipala. 2015. Parsing Parses A Pearl of (Dependently Typed) Programming and Proof. 11.
- [19] J. Hellings. 2014. Conjunctive context-free path queries. (2014).
- [20] Jelle Hellings. 2015. Querying for Paths in Graphs using Context-Free Path Queries. *arXiv preprint arXiv:1502.02242* (2015).
- [21] Jana Hofmann. 2016. Verified Algorithms for Context-Free Grammars in Coq.
- [22] Jan-Oliver Kaiser. 2012. *Constructive Formalization of Regular Languages*. Ph.D. Dissertation. Saarland University.
- [23] André Koschmieder and Ulf Leser. 2012. Regular path queries on large graphs. In *International Conference on Scientific and Statistical Database Management*. Springer, 177–194.
- [24] Raul Lopes, Rodrigo Ribeiro, and Carlos Camarão. 2016. Certified Derivative-Based Parsing of Regular Expressions. In *Programming Languages*, Fernando Castor and Yu David Liu (Eds.). Springer International Publishing, Cham, 95–109.
- [25] Yi Lu, Lei Shang, Xinwei Xie, and Jingling Xue. 2013. An incremental points-to analysis with CFL-reachability. In *International Conference on Compiler Construction*. Springer, 61–81.
- [26] Sebastian Maneth and Fabian Peternek. 2018. Grammar-based graph compression. *Information Systems* 76 (2018), 19 – 45. <https://doi.org/10.1016/j.is.2018.03.002>
- [27] Mark-Jan Nederhof and Giorgio Satta. 2002. Parsing non-recursive context-free grammars. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*. Association for Computational Linguistics, 112–119.
- [28] Mark-Jan Nederhof and Giorgio Satta. 2004. The language intersection problem for non-recursive context-free grammars. *Information and Computation* 192, 2 (2004), 172–184.
- [29] Polyvios Pratikakis, Jeffrey S Foster, and Michael Hicks. 2006. Existential label flow inference via CFL reachability. In *International Static Analysis Symposium*. Springer, 88–106.
- [30] Marcus VM Ramos and Ruy JGB de Queiroz. 2015. Formalization of simplification for context-free grammars. *arXiv preprint arXiv:1509.02032* (2015).
- [31] Marcus Vinícius Midená Ramos, Ruy JGB de Queiroz, Nelma Moreira, and José Carlos Bacelar Almeida. 2016. On the Formalization of Some Results of Context-Free Language Theory. In *International Workshop on Logic, Language, Information, and Computation*. Springer, 338–357.
- [32] Marcus Vinícius Midená Ramos and Ruy J. G. B. de Queiroz. 2015. Formalization of closure properties for context-free grammars. *CoRR abs/1506.03428* (2015). [arXiv:1506.03428](http://arxiv.org/abs/1506.03428) <http://arxiv.org/abs/1506.03428>
- [33] Marcus Vinícius Midená Ramos, Ruy J. G. B. de Queiroz, Nelma Moreira, and José Carlos Bacelar Almeida. 2015. Formalization of the pumping lemma for context-free languages. *CoRR abs/1510.04748* (2015). [arXiv:1510.04748](http://arxiv.org/abs/1510.04748) <http://arxiv.org/abs/1510.04748>
- [34] Jakob Rehof and Manuel Fähndrich. 2001. Type-base flow analysis: from polymorphic subtyping to CFL-reachability. *ACM SIGPLAN Notices* 36, 3 (2001), 54–66.
- [35] Thomas Reps, Susan Horwitz, and Mooly Sagiv. 1995. Precise Interprocedural Dataflow Analysis via Graph Reachability. In *Proceedings of the 22Nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '95)*. ACM, New York, NY, USA, 49–61. <https://doi.org/10.1145/199448.199462>
- [36] Elizabeth Scott and Adrian Johnstone. 2010. GLL parsing. *Electronic Notes in Theoretical Computer Science* 253, 7 (2010), 177–189.
- [37] Dimitrios Vardoulakis and Olin Shivers. 2010. CFA2: A Context-free Approach to Control-flow Analysis. In *Proceedings of the 19th European Conference on Programming Languages and Systems (ESOP'10)*. Springer-Verlag, Berlin, Heidelberg, 570–589. https://doi.org/10.1007/978-3-642-11957-6_30
- [38] Dacong Yan, Guoqing Xu, and Atanas Rountev. 2011. Demand-driven Context-sensitive Alias Analysis for Java. In *Proceedings of the 2011 International Symposium on Software Testing and Analysis (ISSTA '11)*. ACM, New York, NY, USA, 155–165. <https://doi.org/10.1145/2001420.2001440>
- [39] Qirun Zhang and Zhendong Su. 2017. Context-sensitive Data-dependence Analysis via Linear Conjunctive Language Reachability. *SIGPLAN Not.* 52, 1 (Jan. 2017), 344–358. <https://doi.org/10.1145/3093333.3009848>
- [40] Xiaowang Zhang, Zhiyong Feng, Xin Wang, Guozheng Rao, and Wenrui Wu. 2016. Context-Free Path Queries on RDF Graphs. In *The Semantic Web – ISWC 2016*, Paul Groth, Elena Simperl, Alasdair Gray, Marta Sabou, Markus Krötzsch, Freddy Lecue, Fabian Flöck, and Yolanda Gil (Eds.). Springer International Publishing, Cham, 632–648.