

## Bar-Hillel Theorem Mechanization in Coq

Sergey Bozhko, Leyla Khatbullina, **Semyon Grigorev**

JetBrains Research, Programming Languages and Tools Lab  
Saint Petersburg University

July 05, 2019

- Automatization of checking of the proofs correctness

# Automated Theorem Proving

- Automatization of checking of the proofs correctness
- Also a way to create correct-by-construction algorithms
  - ▶ Coq proof assistant
    - ★ Based on calculus of inductive constructions
    - ★ Supports extraction of certified programs to executable programming languages

## Goals:

- Check nontrivial proofs
- Ensure correctness of algorithms
  - ▶ Parsing algorithms
  - ▶ Algorithms over regular expressions
  - ▶ Algorithms over finite automata

# The Bar-Hillel Theorem

## Theorem (Bar-Hillel)

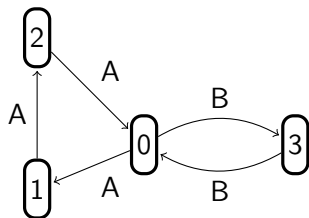
*If  $L_1$  is a context-free language and  $L_2$  is a regular language, then  $L_1 \cap L_2$  is context-free.*

# Context-Free Path Quierying (CFPQ)

Navigation through an edge-labelled graph

# Context-Free Path Quierying (CFPQ)

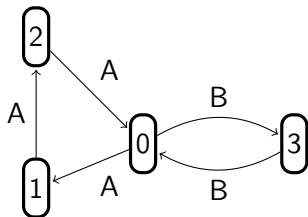
Navigation through an edge-labelled graph



# Context-Free Path Quierying (CFPQ)

Navigation through an edge-labelled graph

- Are there paths in graph, which form well-balanced sequences over A and B?

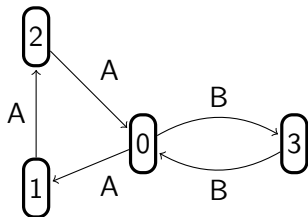




# Context-Free Path Quierying (CFPQ)

Navigation through an edge-labelled graph

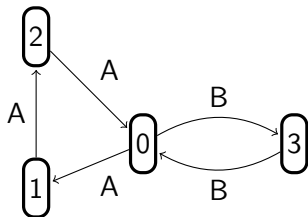
- Are there paths in graph, which form well-balanced sequences over A and B?
- Find all paths, such that they form a word in the Dyck language over A and B



# Context-Free Path Quierying (CFPQ)

Navigation through an edge-labelled graph

- Are there paths in graph, which form well-balanced sequences over A and B?
- Find all paths, such that they form a word in the Dyck language over A and B



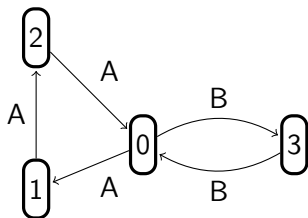
Paths filter (query):

$$s \rightarrow A s B s \mid \varepsilon$$

# Context-Free Path Quierying (CFPQ)

Navigation through an edge-labelled graph

- Are there paths in graph, which form well-balanced sequences over A and B?
- Find all paths, such that they form a word in the Dyck language over A and B



Paths filter (query):

$$s \rightarrow A s B s \mid \varepsilon$$

Answer:

- $2 \xrightarrow{A} 0 \xrightarrow{B} 3$
- $1 \xrightarrow{A} 2 \xrightarrow{A} 0 \xrightarrow{B} 3 \xrightarrow{B} 0$
- ...

- $\mathbb{G} = (\Sigma, N, P, S)$  — context-free grammar
  - ▶  $L(\mathbb{G}) = \{w \mid S \Rightarrow^* w\}$

- $\mathbb{G} = (\Sigma, N, P, S)$  — context-free grammar
  - ▶  $L(\mathbb{G}) = \{w \mid S \Rightarrow^* w\}$
- $G = (V, E, L)$  — directed graph
  - ▶  $v \xrightarrow{I} u \in E$
  - ▶  $L \subseteq \Sigma$

- $\mathbb{G} = (\Sigma, N, P, S)$  — context-free grammar
  - ▶  $L(\mathbb{G}) = \{w \mid S \Rightarrow^* w\}$
- $G = (V, E, L)$  — directed graph
  - ▶  $v \xrightarrow{l} u \in E$
  - ▶  $L \subseteq \Sigma$
- $\omega(\pi) = \omega(v_0 \xrightarrow{l_0} v_1 \xrightarrow{l_1} \cdots \xrightarrow{l_{n-2}} v_{n-1} \xrightarrow{l_{n-1}} v_n) = l_0 l_1 \cdots l_{n-1}$

- $\mathbb{G} = (\Sigma, N, P, S)$  — context-free grammar
  - ▶  $L(\mathbb{G}) = \{w \mid S \Rightarrow^* w\}$
- $G = (V, E, L)$  — directed graph
  - ▶  $v \xrightarrow{l} u \in E$
  - ▶  $L \subseteq \Sigma$
- $\omega(\pi) = \omega(v_0 \xrightarrow{l_0} v_1 \xrightarrow{l_1} \cdots \xrightarrow{l_{n-2}} v_{n-1} \xrightarrow{l_{n-1}} v_n) = l_0 l_1 \cdots l_{n-1}$
- $R = \{(n, m) \mid \exists n\pi m, \text{ such that } \omega(\pi) \in L(\mathbb{G})\}$

- $\mathbb{G} = (\Sigma, N, P, S)$  — context-free grammar
  - ▶  $L(\mathbb{G}) = \{w \mid S \Rightarrow^* w\}$
- $G = (V, E, L)$  — directed graph
  - ▶  $v \xrightarrow{l} u \in E$
  - ▶  $L \subseteq \Sigma$
- $\omega(\pi) = \omega(v_0 \xrightarrow{l_0} v_1 \xrightarrow{l_1} \dots \xrightarrow{l_{n-2}} v_{n-1} \xrightarrow{l_{n-1}} v_n) = l_0 l_1 \dots l_{n-1}$
- $R = \{(n, m) \mid \exists n\pi m, \text{ such that } \omega(\pi) \in L(\mathbb{G})\}$
- $P = \{\pi \mid \pi \text{ is a path in } G, \text{ such that } \omega(\pi) \in L(\mathbb{G})\}$



- Graph database querying
  - ▶ Mihalis Yannakakis, “Graph-theoretic methods in database theory” (1990)
  - ▶ X. Zhang et al, “Context-free path queries on RDF graphs” (2016)

- Graph database querying
  - ▶ Mihalis Yannakakis, “Graph-theoretic methods in database theory” (1990)
  - ▶ X. Zhang et al, “Context-free path queries on RDF graphs” (2016)
- Static code analysis
  - ▶ Thomas Reps. “Program Analysis via Graph Reachability” (1997)
  - ▶ Andrei Marian Dan et al, “Finding Fix Locations for CFL-Reachability Analyses via Minimum Cuts” (2017)

# Sketch of the Proof<sup>1</sup>

## Theorem (Bar-Hillel)

*If  $L_1$  is a context-free language and  $L_2$  is a regular language, then  $L_1 \cap L_2$  is context-free.*

- 1 Assume that there is a context-free grammar  $\mathbb{G}_{CNF}$  in Chomsky Normal Form, such that  $L(\mathbb{G}_{CNF}) = L_1$

---

<sup>1</sup>Richard Beigel and William Gasarch

# Sketch of the Proof<sup>1</sup>

## Theorem (Bar-Hillel)

*If  $L_1$  is a context-free language and  $L_2$  is a regular language, then  $L_1 \cap L_2$  is context-free.*

- 1 Assume that there is a context-free grammar  $\mathbb{G}_{CNF}$  in Chomsky Normal Form, such that  $L(\mathbb{G}_{CNF}) = L_1$
- 2 Assume that there is a set of regular languages  $\{A_1 \dots A_n\}$  where each  $A_i$  is recognized by a DFA with precisely one final state and  $L_2 = A_1 \cup \dots \cup A_n$

---

<sup>1</sup>Richard Beigel and William Gasarch

# Sketch of the Proof<sup>1</sup>

## Theorem (Bar-Hillel)

*If  $L_1$  is a context-free language and  $L_2$  is a regular language, then  $L_1 \cap L_2$  is context-free.*

- ① Assume that there is a context-free grammar  $\mathbb{G}_{CNF}$  in Chomsky Normal Form, such that  $L(\mathbb{G}_{CNF}) = L_1$
- ② Assume that there is a set of regular languages  $\{A_1 \dots A_n\}$  where each  $A_i$  is recognized by a DFA with precisely one final state and  $L_2 = A_1 \cup \dots \cup A_n$ 
  - ▶ If  $L \neq \emptyset$  and  $L$  is regular then  $L$  is the union of regular languages  $A_1, \dots, A_n$  where each  $A_i$  is accepted by a DFA with a single final state

---

<sup>1</sup>Richard Beigel and William Gasarch

# Sketch of the Proof<sup>1</sup>

## Theorem (Bar-Hillel)

*If  $L_1$  is a context-free language and  $L_2$  is a regular language, then  $L_1 \cap L_2$  is context-free.*

- ❶ Assume that there is a context-free grammar  $\mathbb{G}_{CNF}$  in Chomsky Normal Form, such that  $L(\mathbb{G}_{CNF}) = L_1$
- ❷ Assume that there is a set of regular languages  $\{A_1 \dots A_n\}$  where each  $A_i$  is recognized by a DFA with precisely one final state and  $L_2 = A_1 \cup \dots \cup A_n$ 
  - ▶ If  $L \neq \emptyset$  and  $L$  is regular then  $L$  is the union of regular languages  $A_1, \dots, A_n$  where each  $A_i$  is accepted by a DFA with a single final state
- ❸ For each  $A_i$  we can explicitly define a grammar of the intersection:  
 $L(\mathbb{G}_{CNF}) \cap A_i$

---

<sup>1</sup>Richard Beigel and William Gasarch

# Sketch of the Proof<sup>1</sup>

## Theorem (Bar-Hillel)

*If  $L_1$  is a context-free language and  $L_2$  is a regular language, then  $L_1 \cap L_2$  is context-free.*

- ❶ Assume that there is a context-free grammar  $\mathbb{G}_{CNF}$  in Chomsky Normal Form, such that  $L(\mathbb{G}_{CNF}) = L_1$
- ❷ Assume that there is a set of regular languages  $\{A_1 \dots A_n\}$  where each  $A_i$  is recognized by a DFA with precisely one final state and  $L_2 = A_1 \cup \dots \cup A_n$ 
  - ▶ If  $L \neq \emptyset$  and  $L$  is regular then  $L$  is the union of regular languages  $A_1, \dots, A_n$  where each  $A_i$  is accepted by a DFA with a single final state
- ❸ For each  $A_i$  we can explicitly define a grammar of the intersection:  
 $L(\mathbb{G}_{CNF}) \cap A_i$
- ❹ Finally, join them together with the operation of the union

---

<sup>1</sup>Richard Beigel and William Gasarch

# Hofmann's Results Generalization

Jana Hofmann provides mechanization of some theorems for context-free languages in Coq

- Basic definitions: terminal, nonterminal, grammar, word, ...



# Hofmann's Results Generalization

Jana Hofmann provides mechanization of some theorems for context-free languages in Coq

- Basic definitions: terminal, nonterminal, grammar, word, ...
- **Context-Free grammar to the Chomsky Normal Form conversion**

# Hofmann's Results Generalization

Jana Hofmann provides mechanization of some theorems for context-free languages in Coq

- Basic definitions: terminal, nonterminal, grammar, word, ...
- **Context-Free grammar to the Chomsky Normal Form conversion**

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

Jana Hofmann

# Hofmann's Results Generalization

Jana Hofmann provides mechanization of some theorems for context-free languages in Coq

- Basic definitions: terminal, nonterminal, grammar, word, ...
- **Context-Free grammar to the Chomsky Normal Form conversion**

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

Jana Hofmann

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

We needed an arbitrary type for terminals and nonterminals!

# Hofmann's Results Generalization

Jana Hofmann provides mechanization of some theorems for context-free languages in Coq

- Basic definitions: terminal, nonterminal, grammar, word, ...
- **Context-Free grammar to the Chomsky Normal Form conversion**

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

Jana Hofmann

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

We needed an arbitrary type for terminals and nonterminals!

We had to carefully refactor everything...

# DFA Splitting

If  $L \neq \emptyset$  and  $L$  is regular, then  $L$  is the union of regular languages  $A_1, \dots, A_n$  where each  $A_i$  is accepted by a DFA with precisely one final state

# DFA Splitting

If  $L \neq \emptyset$  and  $L$  is regular, then  $L$  is the union of regular languages  $A_1, \dots, A_n$  where each  $A_i$  is accepted by a DFA with precisely one final state

**Lemma** `correct_split`:

```
forall dfa w,  
  dfa_language dfa w <->  
  exists sdfa,  
    In sdfa (split_dfa dfa) /\ s_dfa_language sdfa w.
```

# Chomsky Induction

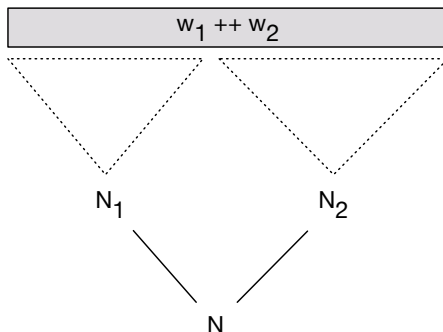
## Lemma

*Let  $\mathbb{G}$  be a grammar in CNF. Consider an arbitrary nonterminal  $N \in \mathbb{G}$  and phrase which consists only of terminals  $w$ . If  $w$  is derivable from  $N$  ( $\text{der}(\mathbb{G}, N, w)$ ) and  $|w| \geq 2$ , then there exists two nonterminals  $N_1, N_2$  and two phrases  $w_1, w_2$  such that:  $N \rightarrow N_1 N_2 \in \mathbb{G}$ ,  $\text{der}(\mathbb{G}, N_1, w_1)$ ,  $\text{der}(\mathbb{G}, N_2, w_2)$ ,  $|w_1| \geq 1$ ,  $|w_2| \geq 1$  and  $w_1 ++ w_2 = w$ .*

# Chomsky Induction

## Lemma

Let  $\mathbb{G}$  be a grammar in CNF. Consider an arbitrary nonterminal  $N \in \mathbb{G}$  and phrase which consists only of terminals  $w$ . If  $w$  is derivable from  $N$  ( $\text{der}(\mathbb{G}, N, w)$ ) and  $|w| \geq 2$ , then there exists two nonterminals  $N_1, N_2$  and two phrases  $w_1, w_2$  such that:  $N \rightarrow N_1 N_2 \in \mathbb{G}$ ,  $\text{der}(\mathbb{G}, N_1, w_1)$ ,  $\text{der}(\mathbb{G}, N_2, w_2)$ ,  $|w_1| \geq 1$ ,  $|w_2| \geq 1$  and  $w_1 ++ w_2 = w$ .





# Chomsky Induction in Coq

```
Definition syntactic_analysis_is_possible :=  
forall (G : grammar) (A : var) (w : phrase),  
  der G A w -> (R A w \in G)  
    \/  
    (exists rhs, R A rhs \in G /\ derf G rhs w).
```

```
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.
```

# The Final Theorem

## Theorem

*For any two decidable types  $\mathbf{T}t$  and  $\mathbf{N}t$  for types of terminals and nonterminals correspondingly. If there exists a bijection from  $\mathbf{N}t$  to  $\mathbb{N}$  and syntactic analysis is possible (in the sense of our definition), then for any DFA  $\mathbf{dfa}$  and any context-free grammar  $\mathbb{G}$ , there exists the context-free grammar  $\mathbb{G}_{INT}$ , such that  $L(\mathbb{G}_{INT}) = L(\mathbb{G}) \cap L(\mathbf{dfa})$ .*

# The Final Theorem in Coq

```
Theorem grammar_of_intersection_exists:  
  exists  
    (NewNonterminal: Type)  
    (IntersectionGrammar: @grammar Terminal NewNonterminal)  
    St,  
  forall word,  
    dfa_language dfa word /\ language G S (to_phrase word)  
    <->  
    language IntersectionGrammar St (to_phrase word).
```

# Conclusion

- We present mechanized in Coq proof of the Bar-Hillel theorem on the closure of context-free languages under intersection with the regular languages

# Conclusion

- We present mechanized in Coq proof of the Bar-Hillel theorem on the closure of context-free languages under intersection with the regular languages
- We generalize the results of Jana Hofmann and Gert Smolka
  - ▶ The definition of the terminal and nonterminal alphabets in context-free grammar were made generic
  - ▶ All related definitions and theorems were adjusted to work with the updated definition

# Conclusion

- We present mechanized in Coq proof of the Bar-Hillel theorem on the closure of context-free languages under intersection with the regular languages
- We generalize the results of Jana Hofmann and Gert Smolka
  - ▶ The definition of the terminal and nonterminal alphabets in context-free grammar were made generic
  - ▶ All related definitions and theorems were adjusted to work with the updated definition
- All results are published at GitHub and are equipped with automatically generated documentation

- Ruy J. G. B. de Queiroz vs Jana Hifmann
  - ▶ We use results of Jana Hofman
  - ▶ Results of Ruy J. G. B. de Queiroz seem more mature
  - ▶ Is it possible to create one “true” solution in this area?
    - ★ Is our grammar-based proof better then PDA-based one in all contexts?



- Ruy J. G. B. de Queiroz vs Jana Hifmann
  - ▶ We use results of Jana Hofman
  - ▶ Results of Ruy J. G. B. de Queiroz seem more mature
  - ▶ Is it possible to create one “true” solution in this area?
    - ★ Is our grammar-based proof better then PDA-based one in all contexts?
- Mechanization of practical algorithms which are just implementation of the Bar-Hillel theorem
  - ▶ Context-free path querying algorithm, based on CYK or even on GLL parsing algorithm
  - ▶ Certified algorithm for context-free constrained path querying for graph databases

# Contact Information

- Semyon Grigorev:
  - ▶ s.v.grigoriev@spbu.ru
  - ▶ Semen.Grigorev@jetbrains.com
- Sergey Bozhko:
  - ▶ Max Planck Institute for Software Systems (MPI-SWS), Saarbrücken, Germany
  - ▶ sbozhko@mpi-sws.com
- Leyla Khatbullina:
  - ▶ St.Petersburg Electrotechnical University “LETI”, St.Petersburg, Russia
  - ▶ leila.xr@gmail.com
- Sources: [https://github.com/YaccConstructor/YC\\_in\\_Coq](https://github.com/YaccConstructor/YC_in_Coq)

Thanks!