Automating Grammar Comparison

Ravichandhran Madhavan

EPFL, Switzerland ravi.kandhadai@epfl.ch

Mikaël Mayer EPFL, Switzerland mikael.mayer@epfl.ch

Viktor Kuncak * EPFL, Switzerland viktor.kuncak@epfl.ch

Sumit Gulwani

Microsoft Research, USA sumitg@microsoft.com



Abstract

We consider from a practical perspective the problem of checking equivalence of context-free grammars. We present techniques for proving equivalence, as well as techniques for finding counter-examples that establish non-equivalence. Among the key building blocks of our approach is a novel algorithm for efficiently enumerating and sampling words and parse trees from arbitrary context-free grammars; the algorithm supports polynomial time random access to words belonging to the grammar. Furthermore, we propose an algorithm for proving equivalence of context-free grammars that is complete for LL grammars, yet can be invoked on any context-free grammar, including ambiguous grammars.

Our techniques successfully find discrepancies between different syntax specifications of several real-world languages, and are capable of detecting fine-grained incremental modifications performed on grammars. Our evaluation shows that our tool improves significantly on the existing available state of the art tools. In addition, we used these algorithms to develop an online tutoring system for grammars that we then used in an undergraduate course on computer language processing. On questions involving grammar constructions, our system was able to automatically evaluate the correctness of 95% of the solutions submitted by students: it disproved 74% of cases and proved 21% of them.

This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive version was published in the following publication:

OOPSLA'15, October 25–30, 2015, Pittsburgh, PA, USA ACM. 978-1-4503-3689-5/15/10... http://dx.doi.org/10.1145/2814270.2814304 *Categories and Subject Descriptors* D.3.1 [*Programming Languages*]: Formal Definitions and Theory—Syntax; D.3.4 [*Programming Languages*]: Processors—Parsing, Compilers, Interpreters; K.3.2 [*Computers and Education*]: Computer and Information Science Education; D.2.5 [*Software Engineering*]: Testing and Debugging

General Terms Languages, Theory, Verification

Keywords Context-free grammars, equivalence, counterexamples, proof system, tutoring system

1. Introduction

Context-free grammars are pervasively used in verification and compilation, both for building input parsers and as foundation of algorithms for model checking, program analysis, and testing. They also play an important pedagogical role in introducing fundamentals of formal language theory, and are an integral part of undergraduate computer science education. Despite their importance, and despite decades of theoretical advances, practical tools that can check semantic properties of grammars are still scarce, except for specific tasks such as parsing.

In this paper, we develop practical techniques for checking equivalence of context-free grammars. Our techniques can find counter-examples that disprove equivalence, and can prove that two context-free grammars are equivalent, much like a software model checker. Our approaches are motivated by two applications: (a) comparing real-world grammars, such as those used in production compilers, (b) automating tutoring and evaluation of context-free grammars. These applications are interesting and challenging for a number of reasons.

Programming Language Grammars. Much of the frontends of modern compilers and interpreters are automatically or manually derived from grammar-based descriptions of programming languages, more so with integrated language support for domain specific languages. When two compilers

^{*} This work is supported in part by the European Research Council (ERC) Project *Implicit Programming* and Swiss National Science Foundation Grant *Constraint Solving Infrastructure for Program Analysis.*

$$S \to S + S \mid S * S \mid ID \qquad \begin{array}{c} S \to ID \ E \\ E \to +S \mid * S \mid \epsilon \end{array}$$

Figure 1. Grammars recognizing simple arithmetic expressions. An example proven equivalent by our tool.

$$\begin{array}{ll} S \to A \Rightarrow S \mid Int \\ A \to Int, A \mid Int \end{array} \qquad \begin{array}{ll} S \to Int \; G \\ G \to \Rightarrow Int \; G \mid , Int \; A \mid \epsilon \\ A \to, Int \; A \mid \Rightarrow Int \; G \end{array}$$

Figure 2. Grammars defining well-formed function signatures over *Int*. An example proven equivalent by our tool.

or other language tools are built according to two different reference grammars, knowing how they differ in the programs they support is essential. Our experiments show that two grammars for the same language almost always differ, even if they aim to implement the same standard. For instance, we found using our tool that two high quality standard Java grammars (namely, a Java grammar written for Antlr v4 parser generator [1] and the Java language specification [2]) disagree on more than 50% of words that are randomly sampled from them.

Even though the differences need not correspond to incompatibilities between the compilers that use these grammars (since their handling could have been intentionally delegated to type checkers and other backend phases), the sheer volume of these discrepancies does raise serious concerns. In fact, most deviations found by our tool are not purposeful. Moreover, in the case of dynamic languages like Javascript, where parsing may happen at run-time, differences between parsers can produce diverging run-time behaviors. Our experiments show that (incorrect) expressions like "++ RegExp - this" discovered by our tool while comparing Javascript grammars result in different behaviors on Firefox and Internet Explorer browsers, when the expressions are wrapped inside functions and executed using the eval construct (see section 5).

Besides detecting incompatibilities, comparing real-world grammars can help identify portions of the grammars that are overly permissive. For instance, many real-world Java grammars generate programs like "enum ID implements char { ID }", "private public class ID" (which were discovered while comparing them with other grammars). These imprecisions can be eliminated with little effort without compromising parsing efficiency.

Furthermore, often grammars are rewritten extensively to make them acceptable by parser generators, which is laborious and error prone. Parser generators have become increasingly permissive over the years to mitigate this problem. However, there still remains considerable overhead in this process, and there is a general need for tools that pin-point subtle changes in the modified versions (documented in works such as [33]). It is almost always impossible to spot differences between large real-world grammars through manual

(a)
$$\begin{array}{cc} S \to A \Rightarrow Int \mid Int \\ A \to S, Int \mid Int \end{array}$$
 (b) $\begin{array}{cc} S \to A \Rightarrow S \mid Int \\ A \to S, S \mid Int \end{array}$

Figure 3. Grammars subtly different from the grammars shown in Fig. 2. The grammar on the left does not accept "Int \Rightarrow Int \Rightarrow Int".

scanning, because the grammars typically appear similar, and even use the same name for many non-terminals. A challenge this paper addresses is developing techniques that scales to real-world grammars, which have hundreds of non-terminals and productions.

Grammars in Teaching. An equally compelling application of grammar comparison arises from the importance of context-free grammar in computer science education. Assignments involving context-free grammars are hard to grade and provide feedback on because of the great variation in the possible solutions arising due to the succinctness of the grammars. The complexity is further aggravated when the solutions are required to belong to subclasses like LL(1). For example, Figures 1 and 2 show two pairs of grammars that are equivalent. The grammars shown on the right are LL(1)grammars, and are reference solutions. The grammars shown on the left are intuitive solutions that a student comes up with initially. Proving equivalence of these pairs of grammars is challenging because they do not have any similarity in their structure, but recognize the same language. On the other hand, Figure 3 shows two grammars (written by students) that subtly differ from the grammars of Fig. 2. The smallest counter-example for the grammar shown in Fig. 3(a) is the string "Int \Rightarrow Int \Rightarrow Int". We invite the readers to identify a counter-example that differentiates the grammar of Fig. 3(b) from those of Fig. 2.

In our experience, a practical system that can prove that a student's solution is correct and provide a counter-example if it is not can greatly aid tutoring of context-free grammars. The state of the art for providing feedback on programming assignments is to use test cases (though there has been recent work on generating repair based feedback [29]). We bring the same fundamentals to context-free grammar education. Furthermore, we exploit the large, yet under-utilized, theoretical research on decision procedures for equivalence of context-free grammars to develop a practical algorithm that can prove the correctness of solutions provided by the students.

Overview and Contributions. At the core of our system is a fast approach for enumerating words and parse trees of an arbitrary context-free grammar, which supports exhaustive enumeration as well as random sampling of parse trees and words. These features are supported by an efficient *polynomial time* random access operation that constructs a unique parse tree for any given natural number index. We construct a scalable counter-example detection algorithm by integrating our enumerators with a state-of-the-art parsing technique [25].

We develop and implement an algorithm for proving equivalence by extending decision procedures for subclasses of deterministic context-free grammars to arbitrary (possibly ambiguous) context-free grammars, while preserving soundness. We make the algorithm practical by performing numerous optimizations, and use concrete examples to guide the proof exploration. We are not aware of any existing system that supports both proving as well as disproving of equivalence of context-free grammars. The following are our main contributions:

- We present an enumerator for generating parse trees of arbitrary context-free grammars that supports the following operations: 1) a *polynomial time* random access operation lookup(i, l) that given an index *i* returns the unique parse tree generating a word of length *l*, corresponding to the index, and 2) sample(n, l) that generates *n* uniformly random samples from the parse trees of the grammar generating words of length *l* (Section 2).
- We use the enumerators to discover counter-examples for equivalence of context-free grammars (Section 3).
- We integrate and extend the algorithms of Korenjak and Hopcroft [15], Olshansky and Pnueli [24], and Harrison et al. [12], for proving equivalence of LL context-free grammars to arbitrary context-free grammars. Our extensions are sound but incomplete. We show using experiments that the algorithm is effective on many grammars that lie outside the classes with known decision procedures (Section 4).
- We evaluate the counter-example detection algorithm on 10 real-world grammars describing the syntax of 5 mainstream programming languages. The algorithm discovers deep, fine-grained errors, by finding counter-examples with an average length of 35, detecting almost 3 times more errors than a state-of-the-art approach (Section 5).
- We implement and evaluate an online tutoring system for context-free grammars. Our system is able to decide the veracity of 95% of the submissions, detecting counter-examples in 74% of the submissions, and proving correctness of 21% of the submissions (Section 6).

2. Enumeration of Parse Trees and Words

A key ingredient of our approach for finding counterexamples is enumeration of words and parse trees belonging to a context-free grammar. Enumeration is also used in optimizing and improving the scope of our grammar equivalence proof engine. We model our enumerators as functions from natural numbers to objects that are enumerated (which are parse trees or words), as opposed to viewing them as iterators for a sequence of objects as is typical in programming language theory. The enumerators we propose are *bijective* functions from natural numbers to parse trees in which the image and pre-image of any given value is efficiently computable in *polynomial time* (formalized in Theorem 1). The functions are partial if the set that is enumerated is finite. Using bijective functions to construct enumerators has many advantages, for example, it immediately provides a way of sampling elements from the given set. It also ensures that there is no duplication during enumeration. Additionally, the algorithm we present here can be configured to enumerate parse trees that generate words having a desired length.

Notations. A context-free grammar is a quadruple (N, Σ, P, S) , where N is a set of non-terminals, Σ is a set of terminals, $P \subseteq N \times (N \cup \Sigma)^*$ is a finite set of productions and $S \in N$ is the start symbol. A parse tree belonging to a grammar is a labelled tree with internal nodes labelled by non-terminals, and leaves labelled by terminals, that has the property that if an internal node labelled A has children labelled C_1, \dots, C_n (from left to right) then $A \to C_1 \dots C_n$. Let \mathcal{T} denote the set of parse trees belonging to a grammar.

We refer to sequences of terminals and non-terminals belonging to $(N \cup \Sigma)^*$ as *sentential forms* of the grammar. If a sentential form has only terminals, we refer to it as a *word*, and also sometimes as a *string*. We adopt the usual convention of using greek characters α, β etc. to represent sentential forms and upper-case latin characters to represent non-terminals. We use lower-case latin characters a, b, c etc. to represent terminals and w, x, y etc. to denote words. We say that α derives β in one step, denoted $\alpha \Rightarrow \beta$, iff β is obtained by replacing the left most non-terminal in α by one of its right-hand-sides. Let \Rightarrow^* denote the reflexive transitive closure of \Rightarrow . The language recognized by a sentential form α , denoted $L(\alpha)$, is the set of words that can be derived from α i.e, $L(\alpha) = \{w \in \Sigma^* \mid \alpha \Rightarrow^* w\}$. We introduce more notations as they are needed.

2.1 Constructing Random Access Enumerators

We use $\operatorname{Enum}[\alpha] : \mathbb{N} \to \mathcal{T}^*$ to denote an enumerator for a sentential form α of the input grammar. The enumerators are *partial functions* from natural numbers to tuples of parse trees of the grammar, one rooted at every symbol in the sentential form. For brevity, we refer to the tuple as parse trees of sentential forms. We define $\operatorname{Enum}[\alpha]$ recursively following the structure of the grammar as explained in the sequel.

For a terminal a belonging to a grammar, $\operatorname{Enum}[a]$ is defined as $\{0 \rightarrow leaf(a)\}$. That is, the enumerator for a terminal a maps the first index to a parse tree with a single leaf node containing a and is undefined for every other index. We now describe an enumerator for a non-terminal. Consider for a moment the non-terminal S of the grammar shown in Fig. 4. The parse trees rooted at S are constructed out of the parse trees that belong to the non-terminal A and the sentential form BA. Assume that we have enumerators defined for Aand BA, namely $\operatorname{Enum}[A]$ and $\operatorname{Enum}[BA]$ that are functions from natural numbers to parse trees (a pair of them in the case of BA). Our algorithm constructs an enumerator for Scompositionally using the enumerators for A and BA.

Figure 4. An example grammar and illustrations of the Enum functions for the symbols of the grammar. Choose and π are defined in Fig. 6 and Appendix A, respectively.

$$\tau(\alpha) = \prod_{i=0}^{n-1} \tau(M_i), \text{ where } \alpha = M_0 \cdots M_{n-1}, n > 1$$

$$\tau(A) = \sum_{i=0}^{n-1} \tau(\alpha_i), \text{ where } A \to \alpha_0 \mid \cdots \mid \alpha_{n-1}$$

$$\tau(a) = 1, \text{ where } a \in \Sigma$$

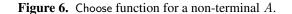
Figure 5. Equations defining the number of parse trees of sentential forms. #t is the greatest fix-point of the equations.

Recall that we consider enumerators as bijective functions from natural numbers. So, given an index i we need to define a unique parse tree of S corresponding to i (provided i is within the number of parse trees rooted at S). To associate a parse tree of S to an index i, we first need to identify a right-handside α of S and select a parse tree t of the right-hand-side. To determine a parse tree t of the right-hand-side α , it suffices to determine the index of t in the enumerator for α . Hence, we define a function $\text{Choose}[A] : \mathbb{N} \to ((N \cup \Sigma)^* \times \mathbb{N})$ for every non-terminal A, that takes an index and returns a right-hand-side of A, and an index for accessing an element of the right-hand-side. We define the enumerator for a nonterminal as: $\text{Enum}[A](i) = node(A, \text{Enum}[\alpha](j))$, where $(\alpha, j) = \text{Choose}[A](i)$. That is, as a node labelled A and having the tuple $\text{Enum}[\alpha](j)$ as children.

In the simplest case, if A has n right-hand-sides $\alpha_0, \alpha_2, \cdots, \alpha_{n-1}$, the choose function $\mathsf{Choose}[A](i)$ could be defined as $(\alpha_{i\%n}, |i/n|)$. This definition, besides being simple, also ensures a fair usage of the right-hand-sides of A by mapping successive indices to different right-hand-sides, which ensures that any sequence of enumeration of the words belonging to a non-terminal alternates over the right-handsides of the non-terminal. However, this definition is well defined only when every right-hand-side of A has unbounded number of parse trees. For instance, consider the non-terminal A shown in Fig. 4. It has two right-hand-sides a and aS of which a has only a single parse tree. Defining Choose[A] as $(\alpha_{i\%2}, |i/2|)$ is incorrect as, for example, Enum[A](2) maps to Enum[a](1), which is not defined. Therefore, we extend the above function so that it takes into account the number of parse trees belonging to the right-hand-sides, which is denoted using $\#t(\alpha)$.

It is fairly straightforward to compute the number of parse trees of non-terminals and right-hand-sides in a grammar. For

$$\begin{aligned} \mathsf{Choose}[A](i) &= \\ \mathsf{let} \ A \to \alpha_0 \mid \cdots \mid \alpha_{n-1} \ \mathsf{s.t.} \\ &\forall 1 \leq m < n. \ \# t(\alpha_{m-1}) \leq \# t(\alpha_m) \ \mathsf{in} \\ \mathsf{let} \ b_0 &= 0, \ \mathsf{and} \ b_1, \cdots, b_n = \# t(\alpha_0), \cdots, \# t(\alpha_{n-1}) \ \mathsf{in} \\ \mathsf{let} \ \forall 0 \leq m \leq n. \ i_m = b_m (n-m+1) + \sum_{i=0}^{m-1} b_i \ \mathsf{in} \\ \mathsf{let} \ k \ \mathsf{be} \ \mathsf{such} \ \mathsf{that} \ 0 \leq k \leq n-1 \ \mathsf{and} \ i_k \leq i < i_{k+1} \ \mathsf{in} \\ \mathsf{let} \ q = \lfloor (i-i_k)/(n-k) \rfloor \ \mathsf{and} \ r = (i-i_k)\%(n-k) \ \mathsf{in} \\ (\alpha_{k+r}, b_k + q) \end{aligned}$$



completeness, we show a formal definition in Fig. 5. Consider the recursive equations shown in Fig. 5 of the form $\tau = F(\tau)$ that define a function F from the set $(N \cup \Sigma)^* \to (\mathbb{N} \cup \{\infty\})$ to itself. Let \leq^{\sharp} be the relation \leq lifted to the domain of F, i.e, for any $\tau_1, \tau_2, \tau_1 \leq^{\sharp} \tau_2$ iff $\forall \alpha.\tau_1(\alpha) \leq \tau_2(\alpha)$. The number of parse trees #t is the greatest fix-point of F with respect to the ordering \leq^{\sharp} . Note that the greatest fix-point can be computed iteratively starting from the initial value $\lambda x.\infty$. As shown in the equations, the number of (tuples of) parse trees of a sentential form is the product of the number of parse trees of a non-terminal is the sum of the number of parse trees of its right-hand-sides, and the number of parse trees of a terminal is one. Note that if the grammar has cycles, #t could be infinite for some sentential forms.

Fig.6 defines a Choose function, explained below, that can handle right-hand-sides with a finite number of parse trees. The definition guarantees that whenever Choose returns a pair (α, i) , *i* is less than $\#t(\alpha)$, which ensures that $\text{Enum}[\alpha]$ is defined for *i*. In Fig.6, the right-hand-sides of the nonterminal $A: \alpha_0, \dots, \alpha_{n-1}$, are sorted in ascending order of the number of parse trees belonging to them. We define b_0 as zero and use b_1, \dots, b_n to denote the number of parse trees of the right-hand-sides. (Note that $\#t(\alpha_i)$ is given by b_{i+1} .) The index i_m is the smallest index (of Enum[A]) at which the m^{th} right-hand-side α_m becomes undefined, which is determined using the number of parse trees of each righthand-side as shown. Given an index *i*, Choose[A](i) first determines the right-hand-sides that need to be skipped i.e, whose enumerators are not defined for the index *i*, by finding a *k* such that $i_k \leq i < i_{k+1}$. It then chooses a right-handside (namely α_j) from the remaining n - k right-hand-sides whose enumerators are defined for the index *i*, and computes the index to enumerate from the chosen right-hand-side.

Most of the computations performed by Fig. 6 – such as computing the number of parse trees of the right-handsides (and hence b_0, \dots, b_n) and the indices i_0, \dots, i_m , and sorting the right-hand-sides of non-terminals by their number of parse trees – need to performed once per grammar. Therefore, for each index *i*, the Choose function may only have to scan through the right-hand-sides to determine the value of *k*, and perform simple arithmetic operations to compute *q* and *r*.

Note that the Choose function degenerates to the simple definition $(\alpha_{i\%n}, \lfloor i/n \rfloor)$ presented earlier when #t is unbounded for every right-hand-side of N. The function also preserves fairness by mapping successive indices to different right-hand-sides of the non-terminals. For instance, in the case of the non-terminal A shown in Fig. 4, the Choose function maps index 0 to (a, 0), index 1 to (aS, 0), but index 2 is mapped to (aS, 1) as a has only one parse tree, i.e, #t(a) = 1.

We now describe the enumerator for a sentential form α with more than one symbol. Let $\alpha = M_1 M_2 \cdots M_m$. The tuples of parse trees belonging to the sentential form is the cartesian product of the parse trees of M_1, \dots, M_m . However, eagerly computing the cartesian product is impossible for most realistic grammars because it is either unbounded or untractably large. Nevertheless, we are interested only in accessing a tuple at a given index *i*. Hence, it suffices to determine for every symbol M_i , the parse tree t_i that is used to construct the tuple at index i. The tree t_j can be determined if we know its index in $Enum[M_i]$. Therefore, it suffices to define a bijective function $\pi: \mathbb{N} \to \mathbb{N}^m$ that maps a natural number (the index of $Enum[\alpha]$) to a point in an *m*dimensional space of natural numbers. The j^{th} component of $\pi(i)$ is the index of Enum $[M_i]$ that corresponds to the j^{th} parse tree of the tuple. In other words, $Enum[M_1 \cdots M_m](i)$ could be defined as $(\operatorname{Enum}[M_1](i_1), \cdots, \operatorname{Enum}[M_m](i_m)),$ where i_j is the j^{th} component of $\pi(i)$.

When *m* is two, the function π reduces to an *inverse* pairing function that is a bijection from natural numbers to pairs of natural numbers. Our algorithm uses only an inverse pairing function as we normalize the right-hand-sides of the productions in the grammar to have at most two symbols. We use the well known *Cantor's inverse pairing function* [26]. But, this function assumes that the two dimension space is unbounded in both directions, and hence cannot be employed directly when the number of parse trees generated by the symbols in the sentential form are bounded. We extend the inverse paring functions to two dimensional spaces that are bounded in one or both the directions. The extended functions take three arguments, the index that is to be mapped, and

Prgm	\rightarrow	<pre>import QName ; ClassDef</pre>
QName	\rightarrow	ID ID . QName
ClassDef	\rightarrow	<pre>class { Body }</pre>

Figure 7. A grammar snippet illustrating the need to bound the length of the generated words during enumeration.

the sizes of the x and y dimensions (or infinity if they are unbounded). We present a formal definition of the functions in Appendix A.

Using the extended Cantor's inverse pairing function π we define the enumerator for a sentential form with two symbols as: $\operatorname{Enum}[M_1M_2](i) = (\operatorname{Enum}[M_1](i_1), \operatorname{Enum}[M_2](i_2))$, where $(i_1, i_2) = \pi(i, \#t(M_1), \#t(M_2))$.

Termination of Random Access. Later in section 2.3 we present a bound on the running time of the algorithm, but now we briefly discuss termination. If A is a recursive nonterminal e.g, if it has a production of the form $A \to \alpha A\beta$, the enumerator for N may recursively invoke itself, either directly or through other enumerators. However, for every index other than 0 and 1, the recursive invocations will always be passed a strictly smaller index. This follows from the definition of the Choose and the inverse pairing functions used by our algorithm. (In the case of the inverse pairing function, if $\pi(i) = (j, k)$, j and k are strictly smaller than i for all i > 1). For indices 0 and 1 the recursive invocations may happen with the same index. However, this will not result in non-termination if the following properties are ensured: (a) for every non-terminal, the right-hand-side chosen for index 0 is the first production in the shortest derivation starting from the non-terminal and ending at a word. (b) There are no unproductive non-terminals (which are non-terminals that do not generate any word) in the input grammar.

From Parse Trees to Words. We obtain enumerators for words using the enumerators for parse trees by mapping the enumerated parse trees to words. However, when the input grammar is ambiguous, the resulting enumerators are no longer bijective mappings from indices to words. The number of indices that map to a word is equal to the number of parse trees of the word.

2.2 Enumerating Fixed Length Words

The enumeration algorithm we have described so far is agnostic to the lengths of the enumerated words. As a consequence, the algorithm may generate undesirably long words, and in fact may also favour the enumeration of long words over shorter ones. Fig. 7 shows a snippet from the Java grammar that results in this behavior.

In the Fig. 7, the productions of the non-terminal Body are not shown for brevity. It generates all syntactically correct bodies allowed for a class in a Java program. Consider the enumeration of the words (or parse trees) belonging to the non-terminal Prgm starting from index 1. A fair enumeration strategy, such as ours, will try to generate almost equal num-

Figure 8. (a) An example grammar. (b) the result of restricting the grammar shown in (a) to words of size 3.

$$\llbracket N \rrbracket_l = \llbracket N \to \alpha_1 \rrbracket_l \cup \dots \cup \llbracket N \to \alpha_n \rrbracket_l,$$

where $N \to \alpha_1, | \dots | \alpha_n$
$$\llbracket N \to a \rrbracket_l = \begin{cases} \{N_l \to a\} & \text{if } l = 1\\ \emptyset & \text{otherwise} \end{cases}$$

$$\llbracket N \to AB \rrbracket_l = \bigcup_{i=1}^{l-1} (\{N_l \to A_i B_{l-i}\} \cup \llbracket A \rrbracket_i \cup \llbracket B \rrbracket_{l-i})$$

Figure 9. Transforming non-terminals and productions of a grammar to a new set of non-terminals and productions that generate only words of length l.

ber of words from the non-terminals QName and ClassDef. However, the lengths of the words generated for the same index differ significantly between the non-terminals. For instance, the word generated from the non-terminal QName at an index *i* has length i + 1. On the other hand, the lengths of the words generated from the non-terminal ClassDef grow slowly relative to their indices, since it has many right-handsides, and each right-hand-side is in turn composed of nonterminals having many alternatives. In essence, the words generated for the non-terminal Prgm will have long import declarations followed by very short class definitions.

Moreover, this also results in reduced coverage of rules since the enumeration heavily reuses productions of QName, but fails to explore many alternatives reachable through ClassDef. We address this issue by extending the enumeration algorithm so that it generates only parse trees of words having a specified length. We accomplish this by transforming the input grammar in such way that it produces only strings that are of the required length, and use the transformed grammar in enumeration. The idea behind the transformation is quite standard. For instance, previous works [14], [19] on theoretical algorithms for random sampling of unambiguous grammars also resort to a similar approach. However, what is unique to our algorithm is using the transformation to construct bijective enumerators while guaranteeing random access property for all words of the specified length.

Fig. 8 illustrates this transformation on an example, which is explained in detail below. For explanatory purposes, assume that the input grammar is in *Chomsky's Normal Form* (CNF) [16] which ensures that every right-hand-side of the grammar is either a terminal or has two non-terminals.

Fig. 9 formally defines the transformation. For every nonterminal N of the input grammar, the transformation creates a non-terminal N_l that generates only those words of N that have a length l. The productions of N_l are obtained by transforming the productions of N. For every production of the form $N \rightarrow a$, where a is a terminal, the transformation creates a production $N_l \rightarrow a$ if l = 1. For every production of the form $N \rightarrow AB$ that has two non-terminals on the right-hand-side, the transformation considers every possible way in which a word of size l can be split between the two non-terminals, and creates a production of the form $N_l \rightarrow A_i B_{l-i}$ for each possible split (i, l-i). Additionally, the transformation recursively produces rules for the nonterminals A_i and B_{l-i} . The transformed grammar may have unproductive non-terminals and rules that do not generate any word (like the non-terminal B_2 and rule $S_3 \rightarrow B_2 S_1$ of Fig. 9(b)), and hence may have to be simplified. Observe that the transformer grammar is acyclic and generates only a finite number of parse trees.

This transformation increases the sizes of the right-handsides by a factor of l, in the worst case. For efficiency reasons, we construct the productions of the transformed grammar on demand, when it is required during the enumeration of a parse tree. Therefore, the functions Enum and Choose take the length of the word to be enumerated as an additional parameter.

Sampling Parse Trees and Words. Having constructed enumerators with the above characteristics, it is straightforward to sample parse trees and words of a non-terminal N having a given length l. We sample numbers in the interval [0, #t(N)-1] uniformly at random, and lookup the parse tree or word at the sampled index using the enumerators. Since we have a bijection from numbers in the range [0, #t(N)-1] to parse trees of N, this approach guarantees a uniform random sampling of parse trees. However, sampling of words is guaranteed to be uniform only if the grammar is unambiguous. In general, the probability of choosing a word w of length l in a sample of size s is equal to $\frac{t \times s}{\#t(N)}$, where t is the number of parse trees of the word w.

2.3 Running Time of Random Access

We derive an upper bound on the time taken by the random access operation by bounding the number of recursive invocations of Enum, and the time spent between two successive invocations of Enum. The number of recursive invocations performed by Enum for generating a word of length l is equal to the number of nodes and edges in the parse tree of the word, which is O(l) for a grammar in CNF [16]. The time spent between two successive invocations of Enum is dominated by the Choose operation, whose running time is bounded by $O(r \cdot l \cdot |i|^2)$ since it performs a linear scan over the right-hand-sides of a non-terminal performing basic arithmetic operations over the given index i. Therefore, we have the following theorem. **Theorem 1.** Let G be a grammar in Chomsky's Normal Form. Let r denote the largest number of right-hand-sides of a non-terminal. Let i be an index, and l the length of the word to be generated. For any non-terminal A belonging to the grammar, the time taken by Enum[A](i,l) is upper bounded by $O(r \cdot |i|^2 \cdot l^2)$, provided the number of parse trees generated by each non-terminal and right-hand-side is precomputed.

Notice that the time taken for random access is polynomial in the size of the input grammar, the number of bits in the index *i* (which is $O(\log i)$), and the size of the generated word *l*. We now briefly discuss the complexity of computing the number of parse trees (#*t*). For a grammar in CNF, the number of parse trees that generate a word of length *l* is $O(r^{2^l})$ in the worst case. (For unambiguous grammars, it is $O(c^l)$, where *c* is the number of terminals.) Thus, computing the number of parse trees could, in principle, be expensive. However, in practice, the number of parse trees, in spite of being large, is efficiently computable.

Prior theoretical works on uniform random sampling (such as [19]) for context-free grammars assume that the input grammar is a constant, and that the arithmetic operations take constant time. (Our approach matches the best known running time $O(l \log l)$ under these assumptions). But, this assumption is quite restrictive in the real-world. For example, the Java 7 grammar has 430 non-terminals and 2447 rules when normalized to CNF, and the number of parse trees increases rapidly with the length of the generated word. In fact, for length 50, it is a 84 digit number (in base 10). Using numbers as big as these in computation introduces significant overhead which cannot be discounted. Our enumerators offer quite some flexibility in sampling by supporting random access. For example, we can sample only from a restricted range instead of using the entire space of parse trees. Since we ensure a fair usage of rules while mapping rules to indices, restricting the sizes of indices still provides a good coverage of rules. In fact, our implementation exposes a parameter for limiting the size of indices, which we found useful in practice.

3. Counter-Example Detection

We apply the enumerators described in previous sections to find counter-examples for equivalence of two context-free grammars. We sample words (of length within a predefined range) from one grammar and check if they are accepted by the other and vice versa. Bounding the length of words greatly aids in reducing the parsing overhead especially while using generic parsers.

In section 5, we present detailed results about the efficiency and accuracy of counter-example detection. The results show that the implementation is able to enumerate and parse millions of words within a few minutes on large realworld grammars. In the sequel, we present an overview of the parsers used by the tool. Parsing. We use a suite of parsers consisting of CYK parser [16], Antlr v4 parser [1] (in compiler and interpreter modes), and LL(1) [3] parser, and employ them selectively depending on the context. For instance, for testing large programming language grammars for equivalence, we compile the grammars to parsers (at runtime) using Antlr v4, which uses adaptive LL(*) parsing algorithm [25], and use the parsers to check if the generated words are accepted by the grammar. The CYK parsing algorithm we implement in our tool, is a top-down, non-recursive algorithm that memoizes the parsing information computed for the substrings of the word being parsed (using a trie data structure), and reuses the information on encountering the same substring again, during a parse of the same or another word. Though the top-down evaluation introduces some overheads compared to the conventional dynamic programming approach, it improves the performance of the CYK parser by orders of magnitude when used in batch mode to parse a collection of words using the same grammar.

We mostly rely on the optimized CYK parser for checking the correctness of students' solutions. We find that quite often the solutions provided by students are convoluted, and are tricky to parse using specialized parsers. For instance, for a grammar with productions $S \rightarrow a \mid B$ and $B \rightarrow aaBb \mid aB \mid \epsilon$, the performance of the Antlr v4 parser degenerates severely with the length of word that is parsed.

4. Proving Equivalence

Our approach for proving equivalence is based on the algorithms proposed by Korenjak and Hopcroft [15], and extended by Olshansky and Pnueli [24] and Harrison et al. [12]. This family of algorithms is attractive because it works directly on context-free grammars without requiring conversions to other representation like push-down automata. Moreover, they come with strong completeness guarantees. Korenjak and Hopcroft [15] introduce a decision procedure for checking equivalence of simple deterministic grammars, which are LL(1) grammars in Griebach Normal Form (GNF). Olshansky and Pnueli [24] extend this algorithm to LL(k) grammars in GNF, while Harrison et al. [12] extend the algorithm in another direction, namely to decide equivalence of deterministic GNF grammars one of which is LL(1). (A grammar is in GNF iff the right-hand-side of every production starts with a terminal. A subtle point to note is that an LL(k) grammar with epsilon productions may become LL(k+1) when expressed in GNF [28].)

Our approach extends the work of Olshansky and Pnueli [24] by incorporating several aspects of the algorithm of Harrison et al. [12]. The resulting algorithm is applicable to arbitrary context-free grammars, but at the same time is complete for LL grammars. (Our implementation is complete only for LL(2) grammars since we limit the lookahead for efficiency reasons.) Furthermore, we perform several extensions to the algorithm that improves its precision and also performance in practice. In particular, we extend the approach

(-)	S	\rightarrow	aT	(1-)	P	\rightarrow	aR
(a)	T	\rightarrow	aT $aTb \mid b$	(b)	R	\rightarrow	$aR \\ abb \mid aRb \mid b$

Figure 10. GNF grammars for the language $a^n b^n$.

to handle inclusion relations, which provides an alternative way of establishing equivalence when the equivalence query is not directly provable. We also introduce transformations that use concrete examples to dynamically refine the queries during the course of the algorithm. Our experiments show that the algorithm succeeds in 82% of the cases that passed all test cases, even on queries involving ambiguous grammars (see section 5).

Algorithm. We use the grammars shown in Fig. 10 for the language $a^n b^n$ as a running example. Observe that the grammar shown on the right is ambiguous – it has two parse trees for *aabb*. We formalize the verification algorithm as a proof system that uses the inference rules shown in Fig. 11. We later discuss an extension to the algorithm that augments the rules with a fixed *lookahead* distance k. Fig. 12 illustrates the algorithm on our running example.

In the sequel, we make the following assumptions: (a) the input grammars have the same set of terminals, (b) the grammars do not have any epsilon productions, and (c) the non-terminals belonging to grammars are unique. In our implementation, if an input grammar has epsilon productions, we remove them using the standard transformations [16]. However, in the case of LL(1) grammars we use the specialized, but expensive algorithm introduced by Rosenkrantz and Stearns [28] that preserves the LL property. This ensures that the algorithm is complete for arbitrary LL(1) grammars including those not in GNF.

Derivatives. We express our algorithm using the notion of a *derivative* of a sentential form which is defined as follows. A derivative $d : \Sigma^* \times (N \cup \Sigma)^* \to 2^{(N \cup \Sigma)^*}$ is a function that given a (non-empty) word w and a sentential form α , computes the sentential forms β that remain immediately after deriving w from α . We define derivatives using left most derivations.

$$d(w,\alpha) = \begin{cases} \{\beta\} & \text{if } \alpha = w\beta \\ \{\beta \mid \exists x \in \Sigma^*, A \in N, \gamma \in (N \cup \Sigma)^* \\ \alpha \Rightarrow^* xA\gamma \Rightarrow w\beta \land |x| < |w| \} \text{ otherwise} \end{cases}$$

For example, given the grammar with rules $A \rightarrow a$ and $B \rightarrow bB \mid b, d(aab, AaB) = \{\epsilon, B\}$, and $d(b, AaB) = \emptyset$. We refer to $d(w, \alpha)$ as a derivative of α with respect to w. Though derivatives are defined for any grammar, they are more efficiently computable when a grammar is normalized to GNF. For a grammar in GNF every production of the grammar starts with a terminal. Hence, a word w that is derivable from a sentential form α would be derivable in at most |w| steps,

We lift the derivative operation to a set of sentential forms as: $\hat{d}(w, \alpha) = \bigcup_{\alpha \in \alpha} d(w, \alpha)$.

Inference Rules. We consider two types of relations between sets of sentential forms: equivalence (\equiv) and inclusion (\subseteq). A relation $\alpha \equiv \beta$ (or $\alpha \subseteq \beta$) holds if the set of words generated by the sentential forms in α i.e, $\bigcup_{\alpha \in \alpha} L(\alpha)$ is equal to (or included in) the set of words generated by the sentential forms in β i.e, $\bigcup_{\beta \in \beta} L(\beta)$. Though we are only interested in proving equivalence of sentential forms, our algorithm sometimes uses inclusion relations in the intermediate steps to establish equivalence. As a consequence, the approach can also be used to prove inclusion of grammars. However, the rules do not guarantee completeness for inclusion queries. The rules shown Fig. 11 use judgements of the form $\mathcal{C} \vdash \alpha \subseteq \beta$, where \mathcal{C} is a set of relations that can be assumed to hold when deciding the truth of $\alpha \subseteq \beta$. Every inference rule shown Fig. 11 provides a set of judgements, given by the antecedents, which, when established, guarantees that the consequent holds. In other words, the antecedents provide a sufficient condition for the consequent. (Sometimes they are also necessary conditions.)

Consider the illustration shown in Fig. 12. Our goal is to establish that the start symbols of the two grammars are equivalent under an empty context, i.e, $\emptyset \vdash [S \equiv P]$. We prove this by finding a derivation for the judgement using the inference rules. In Fig. 12, the relations that are added to the context are marked with \dagger . At any step in the derivation, we can assume that every relation that is marked in the preceding steps leading to the current step hold.

BRANCH Rule. Initially, we apply the BRANCH rule to $[S \equiv P]$. The rule asserts that a relation α op β (where op denotes \subseteq or \equiv) holds in a context C if for every alphabet a, the derivatives of α and β with respect to a are related under the same operation. The correctness of this part is obvious: if two sentential forms are equivalent, the sentential forms that remain after deriving the first character ought to be equivalent. Additionally, the BRANCH rule allows the relation α op β to be considered as valid when proving the antecedents. This is because the BRANCH rule also incorporates inductive reasoning. To prove α op β , the rule hypothesises that the relation holds for all words with length smaller than k, and attempts to establish that the relation holds for words of length k. It suffices for the antecedents to hold for all words of length less than k since we peel off the first character from the words generated by α and β by computing the derivative. Therefore, during the proof of the antecedents if the relation α op β is encountered again then we know that it needs to hold only for words of length less than k, which holds by hypothesis.

An equivalent contrapositive argument is that, if the relation α op β has a counter-example then the antecedents will have a strictly smaller counter-example. However, when α op β is encountered during the proof of the antecedents it need not be explored any further because it would not lead to the smallest counter-example. Harrison et al. [12] refer to this property, wherein the counter-examples of the newly created

$$\begin{array}{c} \frac{\mathsf{BRANCH}}{\forall a \in \Sigma. \ \mathcal{C} \cup \{\alpha \text{ op } \beta\} \vdash \hat{d}(a, \alpha) \text{ op } \hat{d}(a, \beta)}{\mathcal{C} \vdash \alpha \text{ op } \beta} & \frac{\mathsf{INCLUSION}}{\mathcal{C} \vdash \alpha \text{ op } \beta} & \frac{\mathcal{C} \vdash \beta \subseteq \alpha}{\mathcal{C} \vdash \alpha \equiv \beta} & \frac{\mathsf{DIST}}{\mathcal{C} \vdash \{\alpha_i\} \subseteq \beta} \\ & \frac{\mathsf{SPLIT}}{\mathcal{C} \vdash \{\alpha \text{ op } \beta\}} & \frac{\mathsf{SPLIT}}{\mathsf{C} \vdash \{\alpha_i\} = \beta} \\ & \frac{\mathsf{SPLIT}}{\mathsf{C} \vdash \{\alpha\} = \beta} & \frac{\mathsf{SPLIT}}{\mathsf{C} \vdash \{\alpha\} = \beta} \\ & \frac{\mathsf{INDUCT}}{\mathsf{C} \vdash \{\alpha\} = \alpha \text{ op } \beta} & \frac{\mathcal{C}' = \mathcal{C} \cup \{\{A\gamma\} \text{ op } \Psi\} & \mathcal{C}' \vdash \{\gamma\} \text{ op } \hat{d}(x, \Psi) = \bigcup_{i=1}^{m} \rho_i \beta_i & \forall 0 \leq i \leq m. \ |\beta_i| > 0 \\ & \frac{\mathsf{C}' = \mathcal{C} \cup \{\{A\gamma\} \text{ op } \Psi\} & \mathcal{C}' \vdash \{\gamma\} \text{ op } \hat{d}(x, \Psi) & \forall 0 \leq i \leq m. \ |\beta_i| > 0 \\ & \mathcal{C} \vdash \{A\gamma\} \text{ op } \{\psi\} & \mathcal{C} \vdash \{A\gamma\} \text{ op } \{\psi\} & \mathcal{C} \vdash \{A\gamma\} \text{ op } \{\psi\} \\ & \mathcal{C} \vdash \{A\gamma\} \text{ op } \Psi \\ & \mathcal{C} \vdash \{A\gamma\} \text{ op } \Psi \end{array}$$

Figure 11. Basic inference rules of the verification algorithm. In the figure, $op \in \{\equiv, \subseteq\}$, ||A|| is the set of shortest words derivable from A, and $rel_1 \Rightarrow rel_2$ is a syntactic implication check that holds if rel_1 is stronger than rel_2 .

$$\begin{split} ^1[S \equiv P]^{\dagger} & \xrightarrow{\text{BRANCH}} \ ^2[T \equiv R]^{\dagger} & \xrightarrow{\text{BRANCH}} \ ^3[Tb \equiv Rb \cup bb] \wedge ^4[b \equiv b] \\ ^4[b \equiv b] & \xrightarrow{\text{BRANCH}} \ ^5[\epsilon \equiv \epsilon] & \xrightarrow{\text{EPSILON}} \ [\emptyset \equiv \emptyset] & \xrightarrow{\text{EMPTY}} \text{ proved} \\ ^3[Tb \equiv Rb \cup bb] & \xrightarrow{\text{INCLUSION}} \ ^6[Tb \subseteq Rb \cup bb] \wedge ^7[Rb \cup bb \subseteq Tb] \\ ^6[Tb \subseteq Rb \cup bb] & \xrightarrow{\text{TESTCASES}} \ ^8[Tb \subseteq Rb]^{\dagger} & \xrightarrow{\text{SPLIT}} \ ^9[b \subseteq b] \wedge ^{10}[T \subseteq R] & \xrightarrow{\text{INDUCT}} \ ^{11}[b \subseteq b] & \xrightarrow{\star} \text{ proved} \\ ^7[Rb \cup bb \subseteq Tb] & \xrightarrow{\text{DIST}} \ ^{12}[Rb \subseteq Tb] \wedge ^{13}[bb \subseteq Tb] \\ ^{13}[bb \subseteq Tb]^{\dagger} & \xrightarrow{\text{BRANCH}} \ ^{14}[b \subseteq b] \wedge ^{15}[\emptyset \subseteq Tbb] & \xrightarrow{\star} \text{ proved} \\ ^{12}[Rb \subseteq Tb]^{\dagger} & \xrightarrow{\text{SPLIT}} \ ^{16}[b \subseteq b] \wedge ^{17}[R \subseteq T] & \xrightarrow{\text{INDUCT}} \ ^{18}[b \subseteq b] & \xrightarrow{\star} \text{ proved} \end{split}$$

Figure 12. Illustration of application of the rules on the running example. A star (*) denotes application of one or more rules. Curly braces around singleton sets are omitted.

relations (antecedents) are strictly smaller than the counterexamples of the input relation (consequent) when they exist, as *monotonicity*. In our system, the only other monotonic rule is SPLIT.

Applying the BRANCH rule to $[S \equiv P]$ produces the relation $[T \equiv R]$ for the terminal *a* since *T* and *R* are derivatives of *S* and *P* w.r.t *a*, and produces the empty relation $[\emptyset \equiv \emptyset]$ for terminal *b*. The empty relation trivially holds, as asserted by rule *Empty*, and hence is not shown.

Equivalence to Inclusion. The INCLUSION rule reduces equivalence relations to pairs of inclusion relations, (e.g. see relation 3 in Fig. 12). The DIST rule simplifies the inclusion relations by distributing the inclusion operation over the left-hand-sides, as illustrated on the relation 7. These rules ensure that every relation generated during the algorithm is normalized to the form $\{\alpha\} \equiv \{\beta\}$, or $\{\alpha\} \subseteq \beta$.

The TESTCASES rule applies to a relation of the form $\{\alpha\} \subseteq \beta$. It samples a predefined set of words from α and searches for a strict subset of β that accepts all the samples. On finding such a subset *S*, it construct a stronger relation $\{\alpha\} \subseteq S$ that implies the input relation. For instance, the rule

reduces the relation 6: $[Tb \subseteq Rb \cup bb]$ to $[Tb \subseteq Rb]$ using a set of examples. This rule uses an enumerator to sample words from sentential forms and a parser to check if the sample words are accepted by the sentential forms. In our implementation, we use a CYK parser extended for parsing sentential forms to check if the sample words are accepted by the sentential forms.

The TESTCASES rule, despite being incomplete, is useful in practice. It makes the tool converge faster, and also helps in proving more queries by making other rules applicable. For instance, removing smaller sentential forms from a union may make the SPLIT rule (described shortly) applicable. (Fig. 18 of section 5 sheds light on the usefulness of the TESTCASES rule in practice.)

INDUCT Rule. The INDUCT rule asserts that all relations implied by the context hold. The implication check only uses syntactic equality of the sentential forms. In particular, for equality relations $\alpha \equiv \beta$, we check if the context contains the same relation or $\beta \equiv \alpha$. For inclusion relations of the form $\alpha \subseteq \beta$, we check if the context contains an equivalence relation between α and β or an inclusion relation of the form

 $\alpha \subseteq S$, where S has fewer sentential forms than β . For instance, in the derivation shown in Fig. 12, the relations 10 and 17 are implied by the relation 2: $[T \equiv R]$, added to the context in step 2 during the application of BRANCH rule, and hence are considered valid.

SPLIT Rule. The main purpose of the SPLIT Rule is to prevent the sentential forms in the relations from becoming excessively long. The key idea behind the rule is to *split* the sentential forms that are compared (say $[A\gamma \equiv \beta\delta]$) into smaller chunks that are piece-wise equivalent e.g. as $[A\rho \equiv \beta]$, and $[\gamma \equiv \rho\delta]$ (where ρ is a sentential form derived from β), while preserving completeness under some restrictions. It identifies the point to split by deriving a shortest word of A from the other side of the relation.

We apply this rule only to a relation whose left-handside is a singleton (since all relations will be reduced to this form). Let r_1 be the relation $\{A\gamma\}$ op Ψ (with non-empty γ). Let x be one of the shortest words derived from A, denoted ||A||. The SPLIT rule requires that every sentential form in Ψ can be split into $\psi_i\beta_i$ such that ψ_i can derive x, and β_i is non-empty. (However, this requirement can be relaxed as described shortly.) This implies that the derivative of Ψ w.r.t the word x will preserve the suffix β_i . That is, $\hat{d}(x, \Psi)$ will be of the form $\bigcup_i \rho_i\beta_i$, where ρ_i is a union of sentential forms corresponding to the derivative of ψ_i .

Under the above conditions, the rule asserts that if γ op $\hat{d}(x, \Psi)$, and, for all i, $A\rho_i$ op ψ_i holds, then so does r_1 . Furthermore, the rule allows assuming r_1 while proving the antecedents. The requirement that all β_i s are non-empty ensures the monotonicity of the rule. (Note that β_i s cannot derive the empty string as there are no epsilon productions.) If it is not possible to split Ψ in a way that all β_i s are non-empty, the rule is still applicable but r_1 cannot be added to the context, since the rule may not be monotonic.

The soundness of this assertion is relatively easy to establish. For all *i*, $A\rho_i$ op ψ_i implies $A\rho_i\beta_i$ op $\psi_i\beta_i$ (since we are concatenating the left- and the right-handsides with the same sentential form). This entails that $\bigcup_i A\rho_i\beta_i$ op $\bigcup_i \psi_i\beta_i$. We are also given that γ and $\bigcup_i \rho_i\beta_i$ (which is $\hat{d}(x, \Psi)$) are related by op, where op $\in \{\equiv, \subseteq\}$. Substituting $\bigcup_i \rho_i\beta_i$ with γ yields $A\gamma$ op $\bigcup_i \psi_i\beta_i$, which is the relation r_1 . Hence, the antecedents imply the consequent. However, the converse does not necessarily hold. It holds (for equivalence) when the grammars satisfy the *suffix property*: $\alpha\beta \equiv \gamma\beta \Rightarrow \alpha \equiv \gamma$, and are *strict* deterministic [12] (which includes LL(1) grammars).

In the illustration shown in Fig. 12, the split rule is applied on relations 8 and 12. Consider the relation 8: $[Tb \subseteq Rb]$. The shortest word derivable from T is b (see Fig. 14). Since, d(b, Rb) = b, we can deduce that ψ_1 is R, β_1 is b, and $\rho_1 = \{\epsilon\}$ (which are the sentential forms that remain after deriving b from R). The new relations created by the SPLIT rule are γ op d(b, Rb), and $A\rho_1$ op ψ_1 , which correspond to $[b \subseteq b]$ and $[T \subseteq R]$. Note that without the application of the SPLIT rule, the relation $[Tb \subseteq Rb]$ will gradually grow with the application of BRANCH rule and lead to non-termination.

Application Strategy and Termination Checks. In order to preserve termination and completeness of the algorithm for LL grammars, we adopt a specific strategy for applying the rules. We use the INCLUSION rule to convert an equivalence relation to inclusion relations only when at least one of the operands of the relation has size greater than one. Such cases will not arise if both the grammars are LL(1) (or LL(k) when the rules are augmented with a lookahead distance of k). We prevent the sentential forms from growing beyond a threshold by applying SPLIT rule whenever the threshold is reached. We prioritize the application of rules EMPTY, INDUCT, and TESTCASES that simplify the relations over the BRANCH rule. Note that TESTCASES rule, which is incomplete, will not apply to LL(1) grammars since inclusion relations will not be created during its proof exploration.

We use a set of filters to identify relations that are false and to terminate the algorithm. An important filter is the *Length* filter, which checks for every equivalence query $\{\alpha\} \equiv \{\beta\}$, whether the length of the left sentential form α is larger than the length of the shortest word that can be generated by β , and vice versa. If this check fails, one of the sentential forms cannot generate the shortest word of the other and the relation does not hold. (Recall that the input grammar do not have epsilon productions.) We also enumerate words from the sentential forms contained in the relations to detect counterexamples that violate the relation. This helps in quickly aborting the search and reporting a failure, especially for inclusion relations.

The algorithm described above reduces to the algorithm of Korenjak and Hopcroft [15] for LL(1) grammars that are in GNF. Hence, our algorithm is a decision procedure for LL(1) grammars in GNF. Our algorithm may not terminate for grammars outside this class, since the sentential forms in an inclusion relation can grow arbitrarily long. In our implementation, we abort the algorithm and return failure if the algorithm exhausts the memory resources or exceeds a parametrizable time limit (fixed as 10s in our experiments).

4.1 Incorporating Lookahead

In general, the SPLIT rule described above is incomplete. Recall that given a relation $[A\gamma \equiv \psi]$, the rule computes a word x of shortest length derivable from A, and equates γ with the derivative of ψ with respect to x. This is because, since $A\gamma \Rightarrow^* x\gamma$, the rule optimistically hypothesises that γ and $d(x, \psi)$ are equivalent. However, if there are other derivations of the form $A\gamma \Rightarrow^* x\beta$ where $\beta \neq \gamma$, equating γ alone with $d(x, \psi)$ could be incomplete. In the case of LL(1) grammars there is at most one derivation for a sentential form starting with a prefix x from a non-terminal A (if x is non-empty), which entails completeness of the SPLIT rule. Interestingly, this property also holds for LL(k) grammars provided we know the string w of length at least k - 1 that would follow x. In the sequel, we briefly describe the extensions we perform to the proof rules shown in Fig. 11 to exploit this property. Our extensions are based on the algorithm of Olshansky and Pnueli [24]. In essence, the extensions statically enumerate all strings of length smaller than k, referred to as lookahead strings, and use them to create fine-grained relations between sentential forms.

We perform two major extensions to the relations and sentential forms: (a) We qualify every relation with a (possibly empty) word x, which restricts the relations to only words having x as a prefix. For instance, $\alpha \equiv_x \beta$ holds iff α and β generate the same set of words having the prefix x. (b) We introduce two special types of sentential forms: prefix restricted sentential forms (PRS), and grouped variables. A PRS is of the form $[x, \alpha]$ where x is a word and α is a sentential form. It allows only those derivations of α that will lead to a word having x as the prefix. A grouped variable is a disjoint union of two or more PRS that have different prefixes. A grouped variable allows all derivations that are possible through its individual members, akin to a union of sentential forms. PRS and grouped variables are formally defined by Olshansky and Pnueli [24]. They can be treated as any other sentential form. For example, they can be concatenated with other sentential forms, used in derivative computation and so on.

We extend the definition of a derivative $d(w, \alpha)$ so that it additionally accepts a string x and refines the result of $d(w, \alpha)$ to include only those sentential forms that can derive the string x. That is, $d(w, x, \alpha) = \{\beta \mid \beta \in d(w, \alpha) \land$ $(\exists \gamma \text{ s.t. } \beta \Rightarrow^* x\gamma)\}$. We refer to this parameter x as a lookahead as it is not consumed by the derivative but is used to select the sentential forms. We denote using \hat{d} the operation d lifted to a set of sentential forms.

We adapt the BRANCH and SPLIT rules shown in Fig. 11 to the extended domain of relations and sentential forms. (Other rules in Fig. 11 do not require any extensions.) We now discuss the extended branch rule. For brevity, we present the extended SPLIT rule in Appendix B.

BRANCHEXT.

$$\frac{x = aw}{\forall b \in \Sigma. \ \mathcal{C} \cup \{ \boldsymbol{\alpha} \text{ op}_x \ \boldsymbol{\beta} \} \vdash \hat{d}(a, wb, \boldsymbol{\alpha}) \text{ op}_{wb} \ \hat{d}(a, wb, \boldsymbol{\beta})}{\mathcal{C} \vdash \boldsymbol{\alpha} \text{ op}_x \ \boldsymbol{\beta}}$$

Similar to BRANCH rule, the BRANCHEXT rule removes the first character a (of the words considered by the relation) from the sentential forms in α and β . However, unlike the BRANCH rule that compares all the sentential forms left behind after deriving the first character, the BRANCHEXT rule looks ahead at the string wb that follows the character ato choose the sentential forms that have to be compared. Note that the derivative operation only returns the sentential forms that can derive the lookahead string wb.

Given a lookahead distance k, and two grammars with start symbols S_1 and S_2 , we begin the algorithm with the initial set of relations $S_1 \equiv_{w_{k-1}} S_2$, where w_{k-1} is a word of length $\leq k-1$. The grammars are equivalent if every relation

Query	# Ctr. Exs.	# Samples	RA time
$c1 \equiv c2$	82	227	1.0ms
$p1 \equiv p2$	417	1053	0.2ms
$js1 \equiv js2$	75	150	0.8ms
$j1 \equiv j2$	133	240	1.5ms
$v1 \equiv v2$	41	52	2.7ms

Figure 14. Counter-examples found in 1min when comparing grammars of the same programming language. The column *RA time* denotes the average time taken for one random access.

is proven using the inference rules. In our implementation, we fix the lookahead distance as 2. Our implementation reduces to the algorithm of Olshansky and Pnueli [24] when the input grammars are LL(2) GNF grammars, and hence is complete for LL(2) grammars in GNF.

5. Experimental Results

We developed a grammar analysis system based on the algorithms presented in this paper, using the Scala programming language. All the experiments discussed in this section were performed on a 64 bit OpenJDK1.7 VM running on Ubuntu Linux operating system, executing on a server with two 3.5 GHz Intel Xeon processors, and 128GB memory.

5.1 Evaluations with Programming Language Grammars

Fig. 13 presents details about the benchmarks used in the evaluation. The column *Language* shows the list of programming languages chosen for evaluation. For each language we considered at least two grammars which are denoted using the names shown in column *B*. For each language we hand-picked grammars that cover almost all features of the language and are expected to be identical. For example, in the case of Javascript and VHDL, the grammars we choose were supposed to implement the same standard, namely ECMA standard and VHDL-93. In some cases, the grammars even use identical names for many non-terminals, which, however, our algorithm does not attempt to exploit. The column *Size* shows the number of non-terminals and productions in each grammar when expressed in standard BNF form. The column *Source* shows the source of the grammars.

Comparing Real-world Grammars. As an initial experiment, we compared the grammars belonging to the same programming language for equivalence. We ran the counter-example detector for 1 minute on each pair of grammars, fixing the maximum length of the word that is enumerated as 50. Fig. 14 show the results of this experiment. The column *Ctr.Exs* shows the number of counter-examples that were found in 1min, and the column *Samples* shows the number of samples generated during counter-example detection.

The column *RA time* shows the average time taken for accessing one word (of length between 1 and 50) uniformly at

Language	В	Size	Source
C 2011	C 2011 Cl		Antlr v4
C 2011	<i>c</i> 2	(75 269)	www.quut.com/c/ANSI-C-grammar-y.html
Pascal	pl	(177 79)	ftp://ftp.iecc.com/pub/file/pascal-grammar
Tascal	p2 (148		Antlr v3
JavaScript	jsl	(128 336)	www-archive.mozilla.org/js/language/grammar14.html
js2 ((124 278)	Antlr v4
Java 7	Java 7 jl (256 530)		docs.oracle.com/javase/specs/jls/se7/html/jls-18.html
Java /	j2	(229 490)	Antlr v4
VHDL	vl	(286 587)	tams-www.informatik.uni-hamburg.de/vhdl/vhdl.html
VIIDL	v HDL v2 (475 9-		Antlr v4

Figure 13. Benchmarks, their sizes as pairs of number of non-terminals and productions, and their sources. Antlr v4 and Antlr v3 denote the repositories: github.com/antlr/grammars-v4/ and www.antlr3.org/grammar/.

try {
$eval("var k = function() \{ ++ /ab*/ - this \}");$
false;
} catch(err) {
true;
}

Figure 15. A Javascript program, created using a counterexample discovered by our tool, that returns true in Firefox/Chrome browsers, and false in Internet Explorer.

random. The results show that the operation is quite efficient, taking only a few milliseconds across all benchmarks.

Interestingly, as shown by the results, the grammars have many differences even when they implement the same standard. In many cases, more than 40% of the sampled words are counter-examples. Manually inspecting a few counterexamples revealed that this is mostly due to rules that are more permissive than they ought to be. For instance, the string "enum ID implements char { ID }" is generated by j2 (Antlr v4 Java grammar), but is not accepted by j1 [2]. The counter-examples were mostly distinct but sometimes strings that have many parse trees occurred more than once.

From Counter-examples to Incompatibilities. Focusing on Javascript, we studied the usefulness of the counterexamples found by our tool in discovering incompatibilities between Javascript interpreters (embedded within browsers). In this experiment, we automatically converted every counterexample found by our tool (in one minute) while comparing Javascript grammars to a valid Javascript expression, wrapped the expression inside a function, and passed the function as a string to the eval construct. For example, for the counter-example "++ RegExp - this", we generate the program eval("var k = function(){ ++ /ab*/ this }"). This program when executed may either assign a function value to the variable k if the body of the function is parsed correctly, or throw an exception if the body has parse errors¹. We executed the code snippets on Mozilla Firefox (version 38), Google Chrome (version 43) and Microsoft Internet Explorer (version 11) browsers. On five counter-examples, the code snippet threw an exception (either ParseError or ReferenceError) in Firefox and Chrome browsers, but terminated normally in Internet Explorer assigning a function value for k. Exploiting this we created a pure Javascript program shown in Fig. 15 that returns true in Firefox/Chrome browsers and false in Internet Explorer. This can potentially be used as a hidden identification mechanism in adaptive attacks on browser vulnerabilities.

In essence, the experiment highlights that in dynamic languages that supports constructs like eval, where parsing may happen at run-time, differences in parsers will likely manifest as differences in run-time behaviors.

Discovering Injected Errors. In this experiment, we evaluate the effectiveness of our tool on grammars that have comparatively fewer, and subtle counter-examples. Since grammars obtained from independent sources are likely to have many differences, in order to obtain pairs of grammars that almost recognize the same language, we resort to automatic, controlled tweaking of our benchmarks. We introduce 3 types of errors as explained below. (Let G_m denote the modified grammar and G the original grammar).

- Type 1 Errors. We construct G_m by removing one production of G chosen at random. In this case, $L(G_m) \subseteq L(G)$.
- Type 2 Errors. We create G_m by choosing (at random) one production of G having at least two non-terminals, and removing (at random) one non-terminal from the right-hand-side. In this case, neither $L(G_m)$ nor L(G) has to necessarily include the other.
- Type 3 Errors. We construct G_m as follows. We randomly choose one production of the grammar, say P, having at least two non-terminals, and also choose one non-terminal of the right-hand-side, say N. We then create a copy (say N') of the non-terminal N that has every production of N except one (determined at random). We replace N by N' in the production P.

¹www.ecma-international.org/ecma-262/5.1/#sec-15.1.2.1

The above error model has some correspondence to practical scenarios. For instance, most grammars we found do not enforce that the condition of an if statement should be a boolean valued expression. They have productions like $S \rightarrow$ if E then E else $E \mid \cdots$, and $E \rightarrow E + E \mid E \geq E \mid \cdots$. Enforcing this property requires one or more type 3 fixes, as we need to create a copy E' of E that do not have some productions of E, and use E' in place of E to define an if condition.

We avoid injecting errors that can be discovered through small counter-examples using the following heuristic. We repeat the random error injection process until the modified grammar agrees with the original grammar on the number of parse trees (the function #t defined in Fig. 5) generating words of length ≤ 15 . This ensures that the minimum counterexample, if it exists, is at least 15 tokens long. We relax this bound to 10 and 7 for C and JavaScript grammars, respectively, since the approach failed to produce errors that satisfy larger bounds within reasonable time limits. We also ensured that the same error is not reintroduced. It is to be noted that the counter-example detection algorithm is not aware of the similarities between the input grammars, neither does it attempt to discover such similarities.

For each benchmark b and error type t, we create 10 defective versions of b each containing one error of type t. In total, we create 300 defective grammars. In each case, we query the tool for the equivalence of the erroneous and the original versions, with a time out of 15 minutes. Fig. 16 shows the results of this experiment. We categorize the results based on the type of the error that was injected. For now consider only the sub-columns labelled *ours*.

The column *Disproved* shows the number of queries disproved, i.e, the cases where the defective grammar was identified to be not equivalent to the original version. (The maximum possible value for this columns is 10.) Note that for this experiment we ran our tool only until it finds one counter-example. The column *Avg.Time/query* shows the average time taken by the tool on queries where it found a counter-example. The column *Avg.Ctr.Size* shows the average length of the counter-example discovered by the tool. The last row of the table summaries the results by showing the total number of queries disproved, average time taken to disprove a query, and the average length of a counter-example.

The results show that the tool was successful in disproving all queries except 3 for *Type 1 Errors*, and 92 out of 100 queries for *Type 2 Errors*, within a few seconds. For *Type 3 Errors*, which are quite subtle, the tool succeeded in finding counter-examples for 73 out of 100 queries taking at most 200s. It timed out after 15 min in the remaining cases. We found that the tool generated millions of words before timing out on a query, across all benchmarks,.

To put these results in perspective, we now present a comparison with the approach proposed in Axelsson et al.

	Type 1 Errors						
B	Disproved		Avg.Tin	ne/query	Avg.Ctr.Size		
	our	cfga	our	cfga	our	cfga	
<i>c</i> 1	10	7	12.7s	396.7s	29.1	10.0	
<i>c</i> 2	10	4	13.8s	325.0s	30.3	10.3	
<i>p</i> 1	10	7	6.8s	127.8s	39.3	15.0	
<i>p</i> 2	10	5	6.8s	329.2s	43.2	16.2	
js1	10	0	10.9s	-	32.2	-	
js2	10	9	9.6s	190.9s	31.2	8.1	
j1	8	0	14.5s	-	41.1	-	
j2	9	0	14.3s	-	32.1	-	
v1	10	1	16.9s	810.4s	39.3	15.0	
v2	10	0	23.4s	-	39.0	-	
			Type 2 E	rrors			
<i>c</i> 1	9	3	13.3s	319.1s	33.8	10.0	
<i>c</i> 2	10	6	9.1s	300.7s	35.6	10.3	
<i>p</i> 1	10	5	6.2s	358.5s	41.3	16.0	
p2	10	5	7.9s	229.8s	40.0	15.8	
js1	10	0	12.3s	-	33.8	-	
js2	7	8	15.3s	52.8s	31.4	7.4	
j1	7	0	16.3s	-	33.9	-	
j2	9	0	15.1s	-	38.1	-	
v1	10	2	16.4s	729.2s	43.7	15.0	
v2	10	0	58.0s	-	35.8	-	
			Type 3 E	rrors			
<i>c</i> 1	5	4	37.2s	413.6s	17.8	10.3	
<i>c</i> 2	6	5	131.3s	361.2s	30.3	10.0	
<i>p</i> 1	10	3	11.0s	272.5s	34.8	15.0	
<i>p</i> 2	10	5	7.5s	526.8s	34.8	15.8	
js1	5	0	198.6s	-	28.2	-	
js2	5	2	34.0s	79.3s	33.2	7.5	
j1	8	0	25.7s	-	35.4	-	
j2	6	0	24.8s	-	36.3	-	
v1	9	0	17.7s	-	38.6	-	
v2	9	0	54.6s	-	37.3	-	
	262	81	28.1s	342.6s	35.0	12.2	

Figure 16. Identification of automatically injected errors, using our tool (*our*) and the implementation of Axelsson et al. [4] (*cfga*).

[4], which is also used in the more recent work of Creuss and Godoy [7].

Comparison with Previous SAT Solving Approach. The approach proposed in [4] finds counter-examples for equivalence by constructing a propositional formula that is unsatisfiable iff the input grammars are equivalent upto a bound l, i.e, they accept (or reject) the same set of words of length $\leq l$. The approach uses an incremental SAT solver to obtain a satisfying assignment of the formula, which corresponds to a counter-example for equivalence. We ran their tool cfgAnalyzer on the same set of equivalence queries constructed by automatically injecting errors in our benchmarks as described

earlier, with the same time out of 15 minutes. We present the results obtained using their tool in Fig. 16 adjacent to our results, under the sub-column cfga. The cfgAnalyzer tool was run in its default mode, wherein the bound l on the length of the words is incremented in unit steps starting from 1 until a counter-example is found. (We tried adjusting the start length to 15 and greater, and also tried varying the length increments, but they resulted in worse behaviour. This is probably because of incremental solving which may benefit starting from 1.)

The results show that our tool out performs cfgAnalyzer by a huge margin on these benchmarks. When aggregated over all benchmarks, our tool disproves 3 times more queries than cfgAnalyzer. Observe that on Java, VHDL and the first Javascript (js1) benchmarks, cfgAnalyzer timed out on almost all queries. In general, we found that the performance of cfgAnalyzer degrades with the length of the counterexamples, and with the sizes of the grammars. On the other hand, as highlighted by the results in Fig. 16, our tool discovers large counter-examples within seconds.

To the credit of cfgAnalyzer, in cases where it terminates, it finds the shortest counter-example (as a consequence of running it in the default mode). This, however, is not a limitation of our tool, since we can progressively search for smaller counter-examples by narrowing the range of the possible word lengths after discovering a counter-example.

6. Tutoring System for Context-free Grammars

We implemented an online grammar tutoring system available at grammar.epfl.ch using our tool. The tutoring system offers three types of exercises: (a) constructing (LL(1) as well as arbitrary) context-free grammars from English descriptions, (b) converting a given context-free grammar to normal forms like CNF and GNF, and (c) writing left most derivations for automatically generated strings belonging to a grammar. Each class of exercise had about 20 problems each with varying levels of difficulty.

For exercises (a) and (b), the system automatically checks the correctness of the grammars submitted by the users by comparing them to a pre-loaded reference solution for the question. The following are the possible outcomes of checking a solution: (i) the *shortest* counter-example that was found within the time limits, or (ii) a message that the grammar has been proved correct, or (iii) a message that the grammar passed all test cases but was not proved to be correct.

The system also supports checking LL(1) property and ambiguity of grammars. Moreover, it also has experimental support for generating hints (a feature outside the scope of this paper). The system offers a very intuitive syntax for writing grammars, and also supports EBNF form that permits using regular expressions in right-hand-sides of productions.

Query	Refuted	Proved	Unprvd.	time/query
1395	1042	289	64	107ms
(100%)	(74.6%)	(20.7%)	(4.6%)	

Figure 17. Summary of evaluating students' solutions.

Query	Proved	Time	LL1	LL2	Amb		
353	289	410ms	7	56	101		
100%	81.9%		2%	15.9%	28.6%		
w/o TESTCASES rule							
353	280	630ms	7	56	94		

Figure 18. Evaluation of the verification algorithm on students' solutions.

6.1 Evaluations of the Algorithms in the Context of a Tutoring System

We used our tutoring system in a 3rd year undergraduate course on computer language processing. We summarize the results of this study in Fig. 17. The column Queries shows the total number of *distinct* equivalence queries that the tool was run on. The system refuted 1042 queries by finding counter-examples. (It was configured to enumerate at most 1000 words of length 1 to 11). Among the 353 submissions for which no counter-example was found, the tool proved the correctness of 289 submissions. For 64 submissions, the tool was neither able to find a counter-example nor able to prove correctness. In essence, the tool was to able to decide the veracity of 95% of the submissions, and was incomplete on the remaining 5% (in which cases we report that the student's solution is possibly correct). The grammars submitted by students on average had around 3 non-terminals and 6 productions (the maximum was 9 non-terminals and 43 productions). Moreover, at least 370 of the submissions were ambiguous. We now present detailed results on the effectiveness of the verification algorithm, which is, to our knowledge, a unique feature of our grammar tutoring system.

Evaluation of the Verification Algorithm. Our tutoring system uses the verification algorithm described in section 4 to establish the correctness of the submissions for which no counter-examples are found within the given time limit and sample size. In our evaluation, there are 353 such submissions. The first row of Fig. 18 shows the results of using the algorithm, with all of its features enabled, on the 353 submissions. We used a time out of 10s per query. The system proved almost 82% of the queries taking on average less than half a second per query (as shown by column *Time*). The remaining columns further classify the queries that were verified based on the nature of the grammars that are compared.

The column LL1 shows the number of queries in which the grammars that are compared are LL(1) when normalized to GNF. The algorithm of [15] is applicable only to these cases. The results show that only a meager 2% of the queries belong this category. This is expected since even LL(1) grammars may become non-LL(1) when epsilon productions are eliminated [28] (which is required by the verification algorithm).

The column *LL2* shows the number of queries in which the grammars compared are LL(2) but not LL(1), after conversion to GNF. About 16% of the queries belong this category. This class is interesting because the algorithm of [24] is complete for these cases. (Although the algorithm of [12] is also applicable, it seldom succeeds for these queries since it uses no lookhead.) A vast majority (72%) of the queries that are proven involved at least one grammar that is not LL(2). In fact, about 28% of the queries involved ambiguous grammars. (Neither [24] nor [15] is directly applicable to this class of grammars, and [12] is likely to be ineffective.) This indicates that without our extensions a vast majority of the queries may remain unproven. We are not aware of any existing algorithm that can prove equivalence queries involving ambiguous grammars.

We also measure the impact of the TESTCASES inference rule, which uses concrete examples to refine inclusion relations (see section 4). The second row of Fig. 18 shows the results of running the verification algorithm without this rule. Overall, the number of queries proven decreases by 9 when this rule is disabled. The impact is mostly on queries involving ambiguous grammars. Moreover, the verifier is slower in this case as shown by the increase in the average time per query. It also timed out on 25 queries after 10s. This is due to the increase in the number and sizes of *relations* created during the verification algorithm. We measured a two fold increase in the average number of sentential forms contained in a relation.

7. Related Work

Grammar Analysis Systems. Axelsson et al. [4] present a constraint based approach for checking bounded properties of context-free grammars including equivalence and ambiguity. In section 5 we presented a comparison of our counterexample detection algorithm with this work, which shows that our approach does better especially when the counterexamples and grammars are large. Creus and Godoy [7] present RACSO an online judge for context-free grammars. RACSO integrates many strategies for counter-example detection including the approach of Axelsson et al. [4]. We differ from this work in many aspects. For instance, our enumerators support random access and uniform sampling, scale to large programming language grammars generating millions of strings within seconds. Our system can additionally prove equivalence of grammars. (An empirical comparison with this work was not possible since their interface restricts the sizes of grammars that can be used while creating problems, by requiring that non-terminals have to be upper case characters.)

Decision Procedures for Equivalence. Decision procedures for restricted classes of context-free grammars have been extensively researched [15], [24], [12], [28], [32], [23], [5], [31]. For brevity we only highlight a few important works. Korenjak and Hopcroft [15], and later Bastien et al. [5] developed algorithms for deciding equivalence of *simple grammars*. Rosenkrantz and Streans [28] introduced LL(k) grammars and showed that their equivalence is decidable. Later, Olshansky and Pnueli [24] proposed a direct algorithm for deciding equivalence of LL(k) grammars. Nijholt [23] presented a similar result for LL-regular grammars, which properly contain LL(k) grammars. Decision procedures for several proper subclasses of deterministic grammars were studied by Harrison et al. [12], and Valiant [32]. Sénizergues [31] showed that equivalence of arbitrary deterministic grammars is decidable.

We are not aware of any practical applications of these algorithms. We extend the algorithms of Olshansky and Pnueli [24], and Harrison et al. [12] to a sound but incomplete approach for proving equivalence of arbitrary grammars, and use it to power a grammar tutoring system.

Uniform Sampling of Words. Hickey and Cohen [14], and Mairson [19] present algorithms for sampling words from unambiguous grammars uniformly at random (u.a.r). Gore et al. [10] develop a subexponential time algorithm for sampling words from (possibly ambiguous) grammars, where the probability of generating a word varies from uniform by a factor $1 + \epsilon$, $\epsilon \in (0, 1)$. Bertoni et al. [6] present an algorithm for sampling from a finitely ambiguous grammar in polynomial time.

Our approach has a comparable running time for sampling a word u.a.r, and is not restricted to uniform random sampling. We are not aware of any implementations of these related works.

Enumeration in the Context of Testing. Grammar-based software testing approaches (such as [27], [22], [30], [21], [13], [18], [20], [9], [11]) generate strings belonging to grammars describing the structure of the input, and use them to test softwares like refactoring engines and compilers. In contrast to our objective, there the focus is on generating strings from grammars satisfying complex semantic properties, such as data-structure invariants, type correctness etc., that will expose bugs in the software under test.

Purdom [27], and Malloy [21] present specialized algorithms for generating small number of strings that result in semantically correct test cases useful for detecting bugs. Maurer [22], Sirer and Bershad [30], and Guo and Qiu [11] propose approaches for stochastic enumeration of strings from probabilistic grammars where productions are weighted by probabilities. The probabilities are either manually provided or dynamically adjusted during enumeration. A difference compared to our approach is that they do not sample words by restricting their length (which is hard in the presence of semantic properties), but control the frequency with which the productions are used.

Hennessy [13], and Lämmel and Schulte [18] explore various criteria for covering the productions of the grammar that can be beneficial in discovering bugs in softwares. Majumdar and Xu [20], and Godefroid et al. [9] propose approaches for selectively generating strings from a grammar that will exercise a path in the program under test using symbolic execution.

Daniel et al [8], and Kuraj and Kuncak [17] present generic approaches for constructing enumerators for arbitrary structures, by way of enumerator combinators. They allow combining simple enumerators using a set of *combinators* (such as union and product) to produce more complex enumerators. These approaches (Kuraj and Kuncak [17] in particular) were an inspiration for our enumeration algorithm, which is specialized for grammars, and provides more functionalities like polynomial time random access, and uniform random sampling.

8. Conclusions

We present scalable algorithms for enumerating and sampling words (and parse trees) belonging to context-free grammars, using bijective functions from natural numbers to parse trees that provide random access to words belonging to a grammar in polynomial time. Our experiments show that the enumerators are effective in finding discrepancies between large, real world grammars meant to describe the same language, as well as for unraveling deep, subtle differences between grammars, outperforming the available state of the art. We also show that the counter-examples serve as good test cases for discovering incompatibilities between interpreters, especially for languages like Javascript.

We also develop a practical system for proving equivalence of arbitrary context-free grammars building on top of prior theoretical research. We built a grammar tutoring system, available at grammar.epfl.ch, using our algorithms. Our evaluations show that the system is able to decide the correctness of 95% of the submissions, proving over 80% of grammars that pass all test cases. To our knowledge, this is the first tutoring system for grammars that can certify the correctness of the solutions. This opens up the possibility of using our tool in massive open online courses to introduce grammars to large populations of students.

References

- [1] Antlr version 4. http://www.antlr.org/.
- [2] Java 7 language specification. http://docs.oracle.com/ javase/specs/jls/se7/html/jls-18.html.
- [3] A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Princiles, Techniques, and Tools*. Addison-Wesley, 1986. ISBN 0-201-10088-6.
- [4] R. Axelsson, K. Heljanko, and M. Lange. Analyzing context-free grammars using an incremental SAT solver. In *Automata, Languages and Programming, ICALP*, pages 410–422, 2008. URL http://dx.doi.org/10.1007/ 978-3-540-70583-3_34.

- [5] C. Bastien, J. Czyzowicz, W. Fraczak, and W. Rytter. Prime normal form and equivalence of simple grammars. *Theor. Comput. Sci.*, 363(2):124–134, 2006.
- [6] A. Bertoni, M. Goldwurm, and M. Santini. Random generation and approximate counting of ambiguously described combinatorial structures. In STACS 2000, pages 567–580. 2000.
- [7] C. Creus and G. Godoy. Automatic evaluation of context-free grammars (system description). In *Rewriting and Typed Lambda Calculi RTA-TLCA*, pages 139–148, 2014. URL http://dx.doi.org/10.1007/978-3-319-08918-8_10.
- [8] B. Daniel, D. Dig, K. Garcia, and D. Marinov. Automated testing of refactoring engines. In *Foundations of Software Engineering*, pages 185–194, 2007.
- [9] P. Godefroid, A. Kiezun, and M. Y. Levin. Grammar-based whitebox fuzzing. In *Programming Language Design and Implementation*, pages 206–215, 2008.
- [10] V. Gore, M. Jerrum, S. Kannan, Z. Sweedyk, and S. R. Mahaney. A quasi-polynomial-time algorithm for sampling words from a context-free language. *Inf. Comput.*, 134(1):59–74, 1997.
- [11] H. Guo and Z. Qiu. Automatic grammar-based test generation. In *Testing Software and Systems ICTSS*, pages 17–32, 2013.
- [12] M. A. Harrison, I. M. Havel, and A. Yehudai. On equivalence of grammars through transformation trees. *Theor. Comput. Sci.*, 9:173–205, 1979.
- [13] M. Hennessy. An analysis of rule coverage as a criterion in generating minimal test suites for grammar-based software. In *Automated Software Engineering*, pages 104–113, 2005.
- [14] T. J. Hickey and J. Cohen. Uniform random generation of strings in a context-free language. SIAM J. Comput., 12(4): 645–655, 1983.
- [15] A. J. Korenjak and J. E. Hopcroft. Simple deterministic languages. In *Symposium on Switching and Automata Theory* (*Swat*), pages 36–46, 1966.
- [16] D. Kozen. Automata and computability. Undergraduate texts in computer science. Springer, 1997. ISBN 978-0-387-94907-9.
- [17] I. Kuraj and V. Kuncak. Scife: Scala framework for efficient enumeration of data structures with invariants. In *Scala Workshop*, pages 45–49, 2014.
- [18] R. Lämmel and W. Schulte. Controllable combinatorial coverage in grammar-based testing. In *Testing of Communicating Systems, TestCom*, pages 19–38, 2006.
- [19] H. G. Mairson. Generating words in a context-free language uniformly at random. *Inf. Process. Lett.*, 49(2):95–99, 1994.
- [20] R. Majumdar and R. Xu. Directed test generation using symbolic grammars. In *Automated Software Engineering*, pages 553–556, 2007.
- [21] B. A. Malloy. An interpretation of purdom's algorithm for automatic generation of test cases. In *International Conference* on Computer and Information Science, pages 3–5, 2001.
- [22] P. M. Maurer. Generating test data with enhanced context-free grammars. *IEEE Software*, 7(4):50–55, 1990.
- [23] A. Nijholt. The equivalence problem for LL- and LR-regular grammars. pages 149–161, 1982.

- [24] T. Olshansky and A. Pnueli. A direct algorithm for checking equivalence of LL(k) grammars. *Theor. Comput. Sci.*, 4(3): 321–349, 1977.
- [25] T. Parr, S. Harwell, and K. Fisher. Adaptive LL(*) parsing: the power of dynamic analysis. In *Object Oriented Programming Systems Languages & Applications, OOPSLA*, pages 579–598, 2014.
- [26] S. Pigeon. Pairing function. http://mathworld.wolfram. com/PairingFunction.html.
- [27] P. Purdom. A sentence generator for testing parsers. BIT Numerical Mathematics, pages 366–375, 1972.
- [28] D. J. Rosenkrantz and R. E. Stearns. Properties of deterministic top down grammars. In *Symposium on Theory of Computing STOC*, pages 165–180, 1969.
- [29] R. Singh, S. Gulwani, and A. Solar-Lezama. Automated feedback generation for introductory programming assignments. In *Programming Language Design and Implementation PLDI*, pages 15–26, 2013.
- [30] E. G. Sirer and B. N. Bershad. Using production grammars in software testing. In *Domain-Specific Languages DSL*, pages 1–13, 1999.
- [31] G. Sénizergues. L(a)=l(b)? decidability results from complete formal systems. *Theoretical Computer Science*, 251(1–2):1 – 166, 2001.
- [32] L. G. Valiant. Decision procedures for families of deterministic pushdown automata. Technical report, University of Warwick, Coventry, UK, 1973.
- [33] A. Warth, J. R. Douglass, and T. D. Millstein. Packrat parsers can support left recursion. In *Symposium on Partial Evaluation* and Semantics-based Program Manipulation, PEPM, pages 103–110, 2008.

A. Cantor's Inverse Pairing Functions for Bounded and Unbounded Domains

The basic inverse pairing function π that maps a natural number in one dimensional space to a number in two dimensional space that is unbounded along both directions [26]. $\pi(z) = (x, y)$, where x and y are defined as follows:

$$x = w - y$$
$$y = z - t$$
$$(t, w) = simple(z)$$

where, simple(z) = (t, w) is a function defined as follows:

$$t = \frac{w(w+1)}{2}$$
$$w = \left\lfloor \frac{\lfloor \sqrt{8z+1} \rfloor - 1}{2} \right\rfloor$$

We extend the Cantor's inverse pairing function to two dimensional spaces bounded along one or both directions. The inverse pairing function π takes three arguments: the number z that has to be mapped, and the bounds of the x and y dimensions x_b and y_b (which could be ∞). x_b is the (inclusive) bound on the x-axis i.e, $\forall x.x \leq x_b$, and y_b is the (exclusive) bound on the y-axis i.e, $\forall y.y < y_b$. We define $\pi(z, x_b, y_b) = (x, y)$, where

$$\begin{aligned} x &= w - y \\ y &= z - t \end{aligned}$$

$$(t, w) = \begin{cases} bskip(z) & \text{if } z \ge z_b \\ xskip(z) & \text{if } z_x \le z < z_b \\ yskip(z) & \text{if } z_y \le z < z_b \\ simple(z) & \text{Otherwise} \end{cases}$$

where, z_x , z_y and z_b are indices at which the bounds along the x or y or both directions are crossed, respectively. The values are defined as follows:

$$\begin{aligned} z_y &= \frac{y_b(y_b+1)}{2} \\ z_x &= \frac{(x_b+1)(x_b+2)}{2} \\ z_b &= \begin{cases} y_b(x_b-y_b+1) + z_y & \text{if } x_b > y_b - 1 \\ (x_b+1)(y_b-x_b-1) + z_x & \text{if } y_b - 1 > x_b \\ z_y & \text{Otherwise} \end{cases} \end{aligned}$$

We define xskip(z) as (t, w), where t and w are defined as follows:

$$t = \frac{2wx_b - x_b^2 + x_b}{2}$$
$$w = \left\lfloor \frac{2z + x_b^2 + x_b}{2(x_b + 1)} \right\rfloor$$

Define yskip(z) as (t, w), where

l

$$t = \frac{2wy_b - y_b^2 + y_b}{2}$$
$$v = \left\lfloor \frac{2z + y_b^2 - y_b}{2y_b} \right\rfloor$$

Define bskip(z) as (t, w), where

$$t = \frac{(2w_b - 1)w - w^2 - s_b + w_b}{2}$$
$$w = \left\lfloor \frac{r - \left\lceil \sqrt{r^2 - 8z - 4s_b + 4y_b - 4x_b} \right\rceil}{2} \right\rfloor$$
$$r = 2w_b + 1$$
$$w_b = x_b + y_b$$
$$s_b = x_b^2 + y_b^2$$

The above definitions use only integer arithmetic operations. (Note that we always compute floor or ceil of divisions and square roots). These operations take at most quadratic time on the sizes of the inputs, and can be optimized even further. Moreover, many multiplications and divisions are by powers of 2, and hence can be implemented using bit shift operations. SplitExt

$$\begin{aligned} x \in ||A|| \qquad \Psi = \bigcup_{i=1}^{m} \alpha_i \delta_i \beta_i \qquad \forall w \in \Theta_{k-1}(\gamma). \ d(x, w, \alpha_i \delta_i \beta_i) = \boldsymbol{\rho}_i^w \delta_i \beta_i \qquad \forall 0 \le i \le m. \ |\delta_i| > 0, |\beta_i| > 0 \\ \mathcal{C}' = \mathcal{C} \cup \{\{A\gamma\} \text{ op}_z \ \Psi\} \\ \forall w \in \Theta_{k-1}(\gamma). \ \mathcal{C}' \vdash \{\gamma\} \text{ op}_w \ \hat{d}(x, w, \Psi) \qquad \forall 0 \le i \le m. \ \mathcal{C}' \vdash \left(\bigcup_{w \in \Theta_{k-1}(\gamma)} A[\![w, \boldsymbol{\rho}_i^w \delta_i]\!]\right) \text{ op}_z \ \{\alpha_i \delta_i\} \\ \hline \mathcal{C} \vdash \{A\gamma\} \text{ op}_z \ \Psi \end{aligned}$$

Figure 19. Extended Split Rule. $\Theta_{k-1}(\gamma)$ is the set of all words of length at most k-1 derivable from the sentential form γ . ||A|| is the set of shortest words derivable from A. $[w, \alpha]$ denotes a prefix restricted sentential form defined by Olshansky and Pnueli [24].

B. Extended Split Rule

Fig. 19 shows the extended split rule that incorporates a finite amount of lookahead. We assume that the lookahead distance k is at least 2. We define $\Theta_{k-1}(\gamma)$ as the set of all words of length at most k - 1 derivable from a sentential form γ . That is, $\Theta_{k-1}(\gamma) = \{w \mid w \in L(\gamma) \land |w| \le k - 1\}$. Recall the definition of SPLIT shown in Fig. 11. The SPLITEXT rule has a similar structure but it creates more constrained relations by specializing the SPLIT rule for every possible lookahead string in $\Theta_{k-1}(\gamma)$.

Let x be one of the shortest words derivable from A. Akin to the SPLIT rule, SPLITEXT rule applies only when Ψ can be expressed as a union of sentential forms ψ_i each of which has a non-empty suffix of length at least two (represented using $\delta_i\beta_i$) that is preserved by its derivative with respect to x. The rule computes the derivative of ψ_i w.r.t x looking ahead at the strings $w \in \Theta_{k-1}(\gamma)$ (which are possible strings that may follow x). We denote using ρ_i^w the derivative of (the prefix of) ψ_i w.r.t x when the lookahead string is w.

The relations shown in the antecedents of the SPLITEXT rule are straightforward extensions of the antecedents asserted by the SPLIT rule that take into account the lookahead strings in $\Theta_{k-1}(\gamma)$. For instance, the antecedent relations on the left assert that for any lookahead string w, the strings generated by γ and $\hat{d}(x, w, \Psi)$ starting with the prefix w should be related by op. The sentential form $[\![w, \rho_i^w \delta_i]\!]$ used in the antecedent relations on the right denotes a prefix restricted sentential form [24] that generates only those strings of $\rho_i^w \delta_i$ having the prefix w. The extended split rule is applicable to any arbitrary grammar. However, for LL(k) grammars we replace the union in the antecedent relations shown on the right by a grouped variable [24] that denotes a disjoint union, in order to emulate the algorithm of Olshansky and Pnueli [24].