



## Bounded seas



Jan Kurš\*, Mircea Lungu, Rathesan Iyadurai, Oscar Nierstrasz

Software Composition Group, University of Bern, Switzerland<sup>1</sup>

### ARTICLE INFO

#### Article history:

Received 1 June 2015

Accepted 7 August 2015

Available online 24 August 2015

#### Keywords:

Semi-parsing

Island parsing

Parsing expression grammars

### ABSTRACT

Imprecise manipulation of source code (semi-parsing) is useful for tasks such as robust parsing, error recovery, lexical analysis, and rapid development of parsers for data extraction. An island grammar precisely defines only a subset of a language syntax (islands), while the rest of the syntax (water) is defined imprecisely.

Usually water is defined as the negation of islands. Albeit simple, such a definition of water is naïve and impedes composition of islands. When developing an island grammar, sooner or later a language engineer has to create water tailored to each individual island. Such an approach is fragile, because water can change with any change of a grammar. It is time-consuming, because water is defined manually by an engineer and not automatically. Finally, an island surrounded by water cannot be reused because water has to be defined for every grammar individually.

In this paper we propose a new technique of island parsing — bounded seas. Bounded seas are composable, robust, reusable and easy to use because island-specific water is created automatically. Our work focuses on applications of island parsing to data extraction from source code. We have integrated bounded seas into a parser combinator framework as a demonstration of their composability and reusability.

© 2015 Elsevier Ltd. All rights reserved.

## 1. Introduction

Island grammars [1] offer a way to parse input without complete knowledge of the target grammar. They are especially useful for extracting selected information from source files, reverse engineering and similar applications. The approach assumes that only a subset of the language syntax is known or of interest (the islands), while the rest of the syntax is undefined (the water). During parsing, any unrecognized input (water) is skipped until an island is found.

A common misconception is that water should consume everything until some island is detected. Rules for such water are easy to define, but they cause composability problems. Consider a parser where local variables are defined as islands within a method body. Now suppose a method declaring no local variables is followed by one that does. In this case the water might consume the end of the first method as well as the start of the second method until a variable declaration is found. The method variables from the second method will then be improperly assigned to the first one.

In practice, language engineers define many small islands to guide the parsing process. However it is difficult to define such islands in a robust way so that they function correctly in multiple contexts. As a consequence they are neither reusable nor composable.

\* Corresponding author.

E-mail address: [kurs@iam.unibe.ch](mailto:kurs@iam.unibe.ch) (J. Kurš).

<sup>1</sup> <http://scg.unibe.ch>

To prevent our variable declaring island from skipping to another method, we have to make its water stop at most at the end of a method. In general, we have to analyze and update each particular island's water, depending on its context. Yet island-specific water is fragile, hard to define and it is not reusable. It is fragile, because it requires re-evaluation by a language engineer after any change in a grammar. It is hard to define, because it requires the engineer's time for detailed analysis of a grammar. It is not reusable, because island-specific water depends on rules following the island, thus it is tailored to the context in which the island is used — it is not general.

In this paper we propose a new technique for island parsing: *bounded seas* [2]. Bounded seas are composable, reusable, robust and easy to use. The key idea of bounded seas is that specialized water is defined for each particular island (depending on the context of the island) so that an island can be embedded into any rule. To achieve such composability, water is not allowed to consume any input that would be consumed by a following rule.

To prevent fragility and to improve reusability, we compute water automatically, without user interaction. To prove feasibility, we integrate bounded seas into Petit Parser [3], a PEG-based [4] (see Appendix A) parser combinator [5] framework.

In addition to our previous work [2] we evaluate the usability of bounded seas in two case studies, we present a performance study, and we provide more details about the implementation. The contributions of the paper are:

- the definition of bounded seas, a composable, reusable, robust and easy method of island parsing;
- a formalization of bounded seas for PEGs;
- an implementation of bounded seas in a PEG-based parser combinator framework; and
- case studies of semi-parsing of Java and Ruby using bounded seas.

*Structure:* Section 2 motivates this work by presenting the limitations of island grammars with an example. Section 3 presents our solution to overcoming these limitations by introducing bounded seas. Section 4 introduces a sea operator for PEGs, which creates a bounded sea from an arbitrary PEG expression. Section 5 presents our implementation of bounded seas in PetitParser. Section 6 discusses the applicability of bounded seas to GLL parsers, design decisions and some limitations of bounded seas. Section 7 analyzes how well bounded seas perform compare to other island parsers. Section 8 analyzes usability of bounded seas for context-sensitive grammars, particularly for indentation-sensitive grammars. Section 9 surveys other semi-parsing techniques and highlights similarities and differences between them and bounded seas. Finally, Section 10 concludes this paper with a summary of the contributions.

## 2. Motivating example

Let us consider the source code in Listing 1 written in some proprietary object-oriented language. We don't have a grammar specification for the code, because the parser was written using *ad hoc* techniques, and we do not have access to its implementation. Let us suppose that our task is to extract class and method names. Classes may be contained within other classes and we need to keep track of which class each method belongs to.

**Listing 1.** Source code of the `Shape` class in a proprietary language.

```
class Shape
  Color color;

  method getColor {
    return color;
  }
  int uid = UIDGenerator.newUID;
endclass
```

### 2.1. Why not use regular expressions?

To extract a flat list of method names, we could use regular expressions. We need, however, to keep track of the nesting of classes and methods within classes. Regular expressions are only capable of keeping track of finite state, so are formally too weak to analyze our input. To deal with nested structures, we need at least a context-free parser.

Modern implementations of regular expression frameworks can parse more than regular languages (e.g., using recursive patterns<sup>2</sup>). Such powerful frameworks can handle our rather simple task. However regular expressions are not meant to specify complex grammars since they tend to be hard to maintain when the complexity of the grammar grows.

<sup>2</sup> <http://perldoc.perl.org/perlre.html>

## 2.2. A naïve island grammar

To write a parser, we need a grammar. Because the grammar can easily consist of a hundred rules (e.g.,  $\approx 80$  for Python,  $\approx 180$  for Java) and since we are only interested in specific parts of the grammar, we define an island grammar as a PEG (see [Appendix A](#)) with fewer than ten rules as in [Listing 2](#). We initially assume that each class body contains just one method.<sup>3</sup> Since we are interested in extracting method names, we define the `method` rule as an island inside of the `methodWater` rule which surrounds it with water. The `methodWater` rule is defined imprecisely: water skips everything until the string “method” is found.

We also define the `block` rule, which consumes an open curly brace and then skips everything until the closing curly brace is found.

**Listing 2.** Our first island grammar.

```

start      ← class
class      ← 'class' id classBody 'endclass'
classBody  ← methodWater

methodWater ← (!'method' .)* method (!'endclass' .)*

method     ← 'method' id block
block      ← '{' (!'}' .)* '}'

id         ← letter (letter / number)*
letter     ← 'a' / 'b' / 'c' / ...
number     ← '1' / '2' / '3' / ...

```

The `methodWater` rule in the grammar in [Listing 2](#) uses a naïve definition of water. It will work as long as we do not complicate the grammar.

### 2.2.1. Composability problems

Suppose that in order to allow multiple classes in a single file we modify the start rule to allow repetition (`start ← class*`). Parsing the input in [Listing 3](#) should fail because `Shape` does not contain a method. The result, however, no matter whether we use PEG or CFG, is only one class – `Shape` (instead of `Shape` and `Circle`) – with a method `getDiameter`, which is wrong. We see that our water is too greedy here, trying to find a `method` at any cost and ignoring the ‘endclass’ and the `Circle` definition.

**Listing 3.** Source code of `Shape` and `Circle` classes.

```

class Shape
    int uid = UIDGenerator.newUID;
endclass

class Circle
    int diameter;

    method getDiameter {
        return diameter;
    }
endclass

```

<sup>3</sup> We use an almost standard PEG formalism for grammar definitions (see [Appendix A](#)). A terminal is quoted ‘terminal’, a non-terminal is not quoted nonterminal, a sequence is a concatenation of expressions, prioritized choice is marked as /, repetition as \*, a not-predicate as !, and . stands for any character.

Things do not get better when we allow multiple repetitions of `methodWater` within `classBody` (`classBody ← methodWater*`). The parser will stay confused, and, depending on the technology (CFG, PEG), the result will be either ambiguous (CFG) or incorrect (PEG).

The language engineer has to use either (a) disambiguation rules and filters [6,7] to filter out unwanted results of CFGs; or (b) predicates to prevent the incorrect decisions of CFGs and PEGs. Since predicates are applicable for both technologies (CFGs and PEGs), we focus on this approach.

### 2.3. An advanced island grammar

To make the `methodWater` rule composable we must make it possible for it to be embedded into optional (?) or repetition (+, \*) rules. We consequently define the grammar as in Listing 4. This new definition can properly parse multiple classes in a file with an arbitrary number of methods in a class. We achieve composability by forbidding the water to go beyond the 'endclass' keyword and by forbidding the water to consume any method definition.

**Listing 4.** Complete and final island grammar.

```

start      ← class*
class      ← 'class' id classBody 'endclass'
classBody  ← (methodWater)*

methodWater ← (!'method' !'endclass'.)*
            method
            (!'method' !'endclass'.)*

method     ← 'method' id block
block      ← '{'
            (
              (!'}' !'{' .)*
              block
              (!'}' !'{' .)*
            )*
            '}'

id         ← letter (letter / number)*
letter     ← 'a' / 'b' / 'c' ...
number     ← '1' / '2' / '3' ...

```

One can see that the syntactic predicates in the `methodWater` are more complicated. They have been inferred from the rest of the grammar by analyzing which tokens can appear after the `method` island. In case we decide to allow for nested classes, i.e., if we extend the rule `classBody` to:

```
classBody ← (methodWater / classWater)*
```

we have to revise the predicates of `methodWater` to add `!'class'`, and we have to find the proper predicates for the `classWater` rule.

#### 2.3.1. Ease of use, robustness, and reusability problems

The limitations of defining `methodWater` and `classWater` by hand illustrate the general problems of semi-parsing [8,9] with island grammars:

1. Water rules are hard to define correctly because they require the entire grammar to be analyzed.
2. The definition of water is fragile because predicates need to be re-evaluated after any change in a grammar.
3. Finally, the water rules are tailored just for a specific grammar and cannot be reused in another grammar with different rules.

### 3. Bounded seas

#### 3.1. The sea operator in a nutshell

We have shown that water must be tailored both to the island within the water and to the surroundings of the water (e.g., `methodWater` in Listing 4). In this paper, we define a *bounded sea* to be *an island surrounded by context-aware water*.

To automate the definition of bounded seas we introduce a new operator for building tolerant grammars: *the sea operator*. We use the notation  $\sim \text{island} \sim$  to create sea from `island`, which can be a terminal or non-terminal. Instead of having to produce complex definitions of sea, a language engineer can use the sea operator which will do the hard work. Listing 5 shows how the grammar of Listing 4 can be defined using the sea operator.

**Listing 5.** Island Grammar from Listing 4 rewritten with the sea operator.

```
class      ← 'class' id classBody 'endclass'
classBody ← methodSea*

methodSea ← ~method~
method    ← 'method' id block

block     ← '{' ~(block / ε)* '}'

id        ← letter (letter / number)*
letter    ← 'a' / 'b' / 'c' ...
number    ← '1' / '2' / '3' ...
```

A rule defined with the sea operator (e.g.,  $\sim \text{method} \sim$ ) maintains the composability property of the advanced grammar since by applying the sea operator we search for the island in a restricted scope. Moreover, such a rule is reusable, robust, and simple to define.

Bounded seas are based on two ideas:

1. *Water never consumes any input from the right context of the bounded sea, i.e.,* any input that can appear after the bounded sea. This is very different from the water of “traditional” island grammars, where water is not guaranteed to not consume a part of a valid input (cf. Section 2.2.1). The water of bounded seas is unambiguous, thus improving composability.
2. *Everything is fully automated.* The sea is created using the sea operator  $\sim \text{island} \sim$ . Once the sea is placed in the grammar, the grammar is analyzed and appropriate water is created without user interaction. This way the sea can be placed in any grammar. In case the grammar is changed, the water is recomputed automatically. Automatic water computation eases grammar definition, and ensures robustness and reusability of rules.

Bounded seas can be integrated into a parser combinator framework, a highly modular framework for building a parser from other composable parsers [10]. The fact that a bounded sea can be implemented as a parser combinator demonstrates its composability and flexibility.

#### 3.2. The sea boundary

Ideally water should never consume any input that can appear after a bounded sea, i.e., it should never consume an input from its right context. We will call the right context the *boundary* of a sea. The right context of the sea consists of the inputs accepted by parsing expressions that appear after the island. In the case of  $A \leftarrow \sim 'a' \sim (B/C)$ , the right context of  $\sim 'a' \sim$  is any input accepted either by `B` or by `C`.

Being aware of the boundary, a tolerant parser can search for methods in a class without the risk that other classes will interfere. Bounded seas would correctly parse the input in Listing 3 because water of a method sea would not be allowed to consume `endclass`, which is a boundary of the `methodSea`.

The island-specific water has to stop in two cases: first, when an island is reached; second, when a boundary is reached. If a boundary is reached before an island is found, the sea fails. The fact that sea can fail implies that sea can be embedded into optional or repetition expressions without ambiguous results. For example, we can define the superclass specification

**Table 1**

The seas A and B recognize different inputs depending on the context.

	Rule	Input	Result	
1	$R1 \leftarrow A$	"...a..b.."	A recognizes "...a..b.."	
2	$R1 \leftarrow A$	"...a..c.."	A recognizes "...a..c.."	
3	$R2 \leftarrow B$	"...a..b.."	B recognizes "...a..b.."	
4	$R2 \leftarrow B$	"...a..c.."	B fails	
5	$R3 \leftarrow A'b'$	"...a..b.."	A recognizes "...a.."	'b' recognizes 'b'
6	$R3 \leftarrow A'b'$	"...a..c.."	A recognizes "...a..b.."	'b' fails
7	$R4 \leftarrow A'c'$	"...a..b.."	A recognizes "...a..b.."	'c' fails
8	$R4 \leftarrow A'c'$	"...a..c.."	A recognizes "...a.."	'c' recognizes 'c'
9	$R5 \leftarrow A B$	"...a..b.."	A recognizes "...a.."	B recognizes "...b.."
10	$R5 \leftarrow A B$	"...a..c.."	A recognizes "...a..c.."	B fails

as an optional island:

```
~classDef~ ~superclassSpec~? classBody 'endclass'
```

If `superclassSpec` is not present for the particular class, it will simply fail upon reaching `classBody` instead of searching for `superclassSpec` further and further. The same holds for repetitions.

```
classBody ← ~method~*
```

This rule will consume only methods until it reaches "`endclass`" in the input string, since `endclass` is in the boundary of `~method~`, so methods in another class cannot be inadvertently consumed.

We first define bounded seas generally, and subsequently provide a PEG-specific definition.

**Definition 1** (*Bounded sea*). A bounded sea consists of a sequence of three parsing phases:

1. *Before-water*: Consume input until an island or the right context appears. Fail the whole sea if we hit the right context. Continue if we hit an island.
2. *Island*: Consume an island.
3. *After-water*: Consume input until the right context is reached.

### 3.3. The context sensitivity of bounded seas

In order to preserve the unambiguity of water in bounded seas, they need to be context-sensitive. A bounded sea recognizes different substrings of an input depending on what surrounds the sea. There are two cases where context-sensitivity emerges:

1. A bounded sea recognizes different input depending on what immediately follows the sea.
2. A bounded sea recognizes different input depending on what immediately precedes the sea.

Let us demonstrate context sensitivity of bounded seas using rules from Listing 6 and two inputs, "...a..b.." and "...a..c..". On its own, A recognizes any input with 'a' and B recognizes any input with 'b' (see rows 1–4 in Table 1), because they are not bounded by anything.

**Listing 6.** Rules for demonstrating context-sensitive behavior.

```
A ← ~'a'~
B ← ~'b'~

R1 ← A      R2 ← B
R3 ← A 'b'  R4 ← A 'c'
R5 ← A B
```

However, when the two islands are not alone, their boundary can differ, depending on the context. The right context of  $A$  is 'b' in  $R_3$ , and the right context of  $A$  is 'c' in  $R_4$ . Therefore  $A$  consumes different substrings of input depending whether it is called from  $R_3$  or  $R_4$  (see rows 5–8 in Table 1).

A more complex case of context-sensitivity, which we call the *overlapping sea problem*, arises when one sea is immediately followed by another. Consider, for example, rule  $R_5$ , where the sea  $A$  has as its right context  $B$ , which is also a sea. Note that the before-water of  $B$  should consume anything up to its island 'b' or its own right context, *including the island of its preceding sea*  $A$ . Now, the before-water of  $A$  should consume anything up to either its island 'a' or its right context  $B$ . But the very search for the right context will now consume the island we are looking for, since  $B$ 's before-water will consume 'a'! We must therefore take special care to avoid a “shipwreck” in the case of overlapping seas by disabling the before-water of the second sea. Therefore  $B$  recognizes “...a..b..” when called from  $R_2$  and “b..” when called from  $R_5$  (see rows 3 and 9 in Table 1). For the detailed example of the  $\sim a \sim \sim b \sim$  sequence, see B.3.

#### 4. Bounded seas in parsing expression grammars

Starting from the standard definition of PEGs (see Appendix A), we now show how to add the sea operator to PEGs while avoiding the overlapping sea problem. To define the sea operator, we first need the following two abstractions:

1. **The water operator** consumes uninteresting input. Water ( $\approx$ ) is a new PEG prefix operator that takes as its argument an expression that specifies when the water ends. We discuss this in detail in Section 4.1.
2. **The NEXT function** approximates the boundary of a sea. Intuitively,  $NEXT(e)$  returns the set of expressions<sup>4</sup> that can appear directly after a particular expression  $e$ . The details of the NEXT function are given in Section 4.2.

**Definition 2** (Sea operator). Given the definitions of  $\approx$  and NEXT, we define the sea operator as follows:  $\sim e \sim$  is a sequence expression

$$\begin{array}{c} \approx(e / next_1 / next_2 / \dots next_n) \\ e \\ \approx(next_1 / next_2 / \dots next_n) \end{array}$$

where  $next_i \in NEXT(e)$  for  $i = 1..n$  and  $n = |NEXT(e)|$ .

That is, the before-water consumes everything up to the island or the boundary, and the after-water consumes everything up to the boundary.

##### 4.1. The Water operator

The purpose of a water expression is to consume uninteresting input. Water consumes input until it encounters the expression specified in its argument (i.e., the *boundary*). We must, however, take care to avoid the overlapping sea problem.

If two seas overlap (one sea is followed by another), the right boundary of the first sea starts with the second sea. Yet it should only start with the island of the second sea as illustrated in Section 3.3. In order to do so, the second sea has to simply disable its before-water.

We detect overlapping seas as follows: if sea  $s_2$  is invoked from the water of another sea  $s_1$ , it means that the water of  $s_1$  is testing for its boundary  $s_2$  and thus  $s_2$  has to disable its before-water. To distinguish between nested seas (e.g., ' $\sim x' \sim island \sim 'x' \sim$ ') and overlapping seas (e.g., ' $\sim 'x' \sim \sim 'y' \sim$ '), we test the position where this sea was invoked. In case of nested seas the positions differ, and in case of overlapping seas they are the same.

**Definition 3** (Extended semantics of PEGs). In order to detect overlapping seas and to compute the NEXT set, we extend the original semantics of a PEG  $G = \{N, T, R, e_s\}$  (see Definition 8 in Appendix A) with a stack of invoked expressions and their positions. For standard PEG operators there is no change except that an explicit stack  $S$  is maintained. We define a relation  $\Rightarrow$  from tuples of the form  $(x, S)$  to the output  $o$ , where  $x \in T^*$  is an input string to be recognized,  $S$  is a stack of tuples  $(e, p)$ , where  $e$  is a parsing expression and  $p \geq 0$  is a position, and  $o \in T^* \cup \{f\}$  indicates the result of a recognition attempt. The distinguished symbol  $f \notin T$  indicates failure. Function  $len(x)$  returns the length of an input  $x$ . Function  $(e, p):S$  denotes a stack with tuple  $(e, p)$  on the top and stack  $S$  below.  $S$  is initialized with the pair  $(e_s, 0)$ .

We define  $\Rightarrow$  inductively as follows (without any semantic changes for standard PEG operators)<sup>5</sup>:

$$\text{Empty: } \frac{x \in T^*}{(x, (e, p):S) \Rightarrow e}$$

<sup>4</sup> The NEXT function is modeled after FOLLOW sets from parsing theory, except that instead of returning a set of tokens, it returns a set of parsers.

<sup>5</sup> Note that in these rules  $p$  is implicitly defined as the current position in the input.

$$\begin{aligned}
\text{Terminal (success case): } & \frac{a \in T \quad x \in T^*}{(ax, (a, p): S) \Rightarrow a} \\
\text{Terminal (failure case): } & \frac{a \neq b \quad (a, \epsilon, S) \Rightarrow f}{(bx, (a, p): S) \Rightarrow f} \\
\text{Nonterminal: } & \frac{A \leftarrow e \in R \quad (x, (e, p): S) \Rightarrow o}{(x, (A, p): S) \Rightarrow o} \\
& \quad (x_1 x_2 y, (e_1, p): S) \Rightarrow x_1 \\
\text{Sequence (success case): } & \frac{(x_2 y, (e_2, p + \text{len}(x_1)): S) \Rightarrow x_2}{(x_1 x_2 y, (e_1 e_2, p): S) \Rightarrow x_1 x_2} \\
\text{Sequence (failure case): } & \frac{(x, (e_1, p): S) \Rightarrow f}{(x, (e_1 e_2, p): S) \Rightarrow f} \\
& \quad (xy, (e_1, p): S) \Rightarrow x \\
\text{Sequence (failure case 2): } & \frac{(y, (e_2, p + \text{len}(x)): S) \Rightarrow f}{(xy, (e_1 e_2, p): S) \Rightarrow f} \\
\text{Alternation (case 1): } & \frac{(xy, (e_1, p): S) \Rightarrow x}{(x, (e_1 / e_2, p): S) \Rightarrow x} \\
& \quad (x, (e_1, p): S) \Rightarrow f \\
\text{Alternation (case 2): } & \frac{(x, (e_2, p): S) \Rightarrow o}{(x, (e_1 / e_2, p): S) \Rightarrow o} \\
& \quad (x_1 x_2 y, (e, p): S) \Rightarrow x_1 \\
\text{Repetitions (repetition case): } & \frac{(x_2, (e*, p + \text{len}(x_1)): S) \Rightarrow x_2}{(x_1 x_2 y, (e*, p): S) \Rightarrow x_1 x_2} \\
\text{Repetitions (termination case): } & \frac{(x, (e, p): S) \Rightarrow f}{(x, (e*, p): S) \Rightarrow \epsilon} \\
\text{Not predicate (case 1): } & \frac{(xy, (e, p): S) \Rightarrow x}{(xy, (!e, p): S) \Rightarrow f} \\
\text{Not predicate (case 2): } & \frac{(xy, (e, p): S) \Rightarrow f}{(xy, (!e, p): S) \Rightarrow \epsilon}
\end{aligned}$$

A detailed example can be found in [B.3](#).

**Definition 4** (*Water operator*). With the extended semantics of PEGs we can define a prefix **water operator**  $\approx$ . It searches for a boundary and consumes input until it reaches a boundary. If the water starts a boundary of another sea, it stops immediately. Function  $\text{seasOverlap}(S, p_1)$  returns true if there is a pair  $(\approx e, p_2)$  on a stack  $S$  where  $p_1 = p_2$  and  $e$  is any parsing expression and returns false otherwise.  $x \in T^*, y \in T^*, z \in T^*$ .

$$\begin{aligned}
\text{Overlapping seas: } & \frac{\text{seasOverlap}(S, p)}{(x, (\approx e, p): S) = \epsilon} \\
& \quad (yz, (e, p): S) \Rightarrow y \\
& \quad (x', (e, p + \text{len}(x')): (\approx e, p + \text{len}(x')): S) \Rightarrow f \\
\text{Boundary found: } & \frac{\forall x = x' x'' x'''}{(xyz, (\approx e, p): S) = x}
\end{aligned}$$

In case of *directly nested seas* (e.g.,  $\sim \sim \text{island} \sim \sim$ ) we obtain the same behavior as with  $\sim \text{island} \sim$ . The function  $\text{seasOverlap}$  returns true in case a sea is directly invoked from another sea without consuming any input. Applying the rule *Overlapping seas* from [Definition 4](#), water of the inner sea is eliminated and the boundary is the same for the both seas. Therefore  $\sim \sim \text{island} \sim \sim$  is equivalent to  $\sim \text{island} \sim$ .

#### 4.2. The NEXT function

Any input that can appear after the sea forms a boundary of a sea. The NEXT function returns a set of expressions that can appear directly after a particular expression.

Consider the grammar in the example from [Listing 7](#). The `code` rule is defined in such a way that it accepts an arbitrary number of class and structure islands in the beginning (classes and structures can be in any order) and there is a main method at the end. Intuitively, another class island, a structure island or a main method can appear after a class island.



The NEXT set approximates the boundary. Its expressions recognize prefixes of the boundary and not necessarily the whole boundary. The reason for using NEXT is the limited backtracking ability of PEGs. PEGs are not capable of taking globally correct decisions because they are not able to revert choices that have already been taken.<sup>6</sup>

**Listing 7.** Definition of code that consists of classes and structures followed by main method.

```
code      ← (∼class∼/∼struct∼)* mainMethod
class     ← 'class' ID classBody
struct    ← 'struct' ID sbody
mainMethod ← 'public' 'method' 'main' block

classBody ← ...
sbody     ← ...
block     ← ...
ID        ← ...
```

For practical reasons, elements of NEXT cannot accept an empty string. For example, an optional expression is not a suitable approximation of a boundary, because it succeeds for any input. Consider a simple expression  $\sim e \sim 'a'? 'b'$ . The  $'a'?$  can appear after the  $'island'$  but  $'b'$  as well if  $'a'$  fails. Therefore NEXT has to return  $'a'? 'b'$ , not just  $'a'?$ .

We will use *abstract simulation* [4] in order to recognize an expression that accepts an empty string.

**Definition 5** (*Abstract simulation*). We define a relation  $\rightarrow$  consisting of pairs  $(e, o)$ , where  $e$  is an expression and  $o \in \{0, 1, f\}$ . If  $e \rightarrow 0$ , then  $e$  can succeed on some input while consuming no input. If  $e \rightarrow 1$ , then  $e$  can succeed on some input while consuming at least one terminal. If  $e \rightarrow f$ , then  $e$  may fail on some input. We will use variable  $s$  to represent a  $\rightarrow$  outcome of either 0 or 1. We will define the simulation relation  $\rightarrow$  as follows:

1.  $e \rightarrow 0$ .
2.  $t \rightarrow 1$ ,  $t \in T$ .
3.  $t \rightarrow f$ ,  $t \in T$ .
4.  $A \rightarrow o$  if  $e \rightarrow o$  and  $A \leftarrow e$  is a rule of the grammar  $G$ .
5.  $e_1 e_2 \rightarrow 0$  if  $e_1 \rightarrow 0$  and  $e_2 \rightarrow 0$ .  
 $e_1 e_2 \rightarrow 1$  if  $e_1 \rightarrow 1$  and  $e_2 \rightarrow s$ .  
 $e_1 e_2 \rightarrow 1$  if  $e_1 \rightarrow s$  and  $e_2 \rightarrow 1$ .
6.  $e_1 e_2 \rightarrow f$  if  $e_1 \rightarrow f$
7.  $e_1 e_2 \rightarrow f$  if  $e_1 \rightarrow s$  and  $e_2 \rightarrow f$ .
8. (a)  $e_1 / e_2 \rightarrow 0$  if  $e_1 \rightarrow 0$   
(b)  $e_1 / e_2 \rightarrow 1$  if  $e_1 \rightarrow 1$
9.  $e_1 / e_2 \rightarrow o$  if  $e_1 \rightarrow f$  and  $e_2 \rightarrow o$ .
10.  $e^* \rightarrow 1$  if  $e \rightarrow 1$
11.  $e^* \rightarrow 0$  if  $e \rightarrow f$
12.  $!e \rightarrow f$  if  $e \rightarrow s$
13.  $!e \rightarrow 0$  if  $e \rightarrow f$

Because this relation does not depend on the input string, and there are a finite number of expressions in a grammar, we can compute this relation over any grammar [4]. An example of abstract simulation can be found in Appendix B.1.

**Definition 6** (*NEXT*). Let  $S$  be a stack of (*expression, position*) pairs representing positions and invoked parsing expressions, where  $\top(S)$  pops an element from the stack  $S$  returning a stack  $S'$  without the top element,  $s_n, s_{n-1}, \dots, s_2, s_1$  are expressions on the stack  $S$  (top of the stack is to the left, bottom to the right),  $\$$  is a special symbol signaling end of input, and  $E_1 \times E_2$  is a product of two sets of parsing expressions,  $E_1$  and  $E_2$ , such that  $E_1 \times E_2 = \{e_i e_j | e_i \in E_1, e_j \in E_2\}$ , we define  $NEXT(S)$  as a set of expressions such that:

- if  $s_n = e_1$  and  $s_{n-1} = e_1 e_2$  and  $e_2 \rightarrow 0$  then  $NEXT(S) = \{e_2\}$
- if  $s_n = e_1$  and  $s_{n-1} = e_1 e_2$  and  $e_2 \rightarrow 0$  then  $NEXT(S) = \{e_2\} \times NEXT(\top(S))$

<sup>6</sup> See for example: <http://www.webcitation.org/6YrGmNAi7>.

- if  $s_n = e_1$  and  $s_{n-1} = e_1 e_2$  then  $NEXT(S) = \{e_2\}$
- if  $s_n = e_2$  and  $s_{n-1} = e_1 e_2$  then  $NEXT(S) = NEXT(\overline{\wedge}(S))$
- if  $s_n = e_1$  or  $s_n = e_2$  and  $s_{n-1} = e_1 / e_2$  then  $NEXT(S) = NEXT(\overline{\wedge}(S))$
- if  $s_n = e$  and  $s_{n-1} = e*$  then  $NEXT(S) = e \cup NEXT(\overline{\wedge}(S))$
- if  $s_n = e$  and  $s_{n-1} = !e$  then  $NEXT(S) = \{\}$
- if  $s_n = e \in N$  then  $NEXT(S) = NEXT(\overline{\wedge}(S))$
- if  $n=0$  (stack is empty) then  $NEXT(S) = \{\$\}$

An example of NEXT computation can be found in Appendix B.2.

## 5. Implementation

As a validation of bounded sea composability and reusability we report on an implementation of bounded seas in the PetitParser framework.<sup>7</sup> The bounded sea extension of PetitParser is part of Moose – a platform for software and data analysis.<sup>8</sup>

### 5.1. PetitParser internals

PetitParser [3,11] is a PEG-based parser combinator [5] framework utilizing scannerless parsing [12] and packrat parsing [13]. Implementations of PetitParser exist for Pharo Smalltalk<sup>9</sup> (the version we extended), Java<sup>10</sup> and Dart.<sup>11</sup>

PetitParser combinators are subclasses of the `PPParser` class, which defines an abstract method `parse: anInput`. If parsing fails, `PPFailure` is returned, otherwise a result is returned. For example, the `PPSequence` combinator is subclass of `PPParser`, having two extra instance variables referring to two parsers that should be in sequence as you can see in Listing 8. The method `parse: anInput` is implemented as shown in Listing 9. The method returns a failure if either of the two parsers fails, and returns both results in an array if they both succeed.

**Listing 8.** `PPSequence` has two instance variables, `p1` and `p2`.

```
PPParser subclass: #PPSequence
    "Sequence of two parsers, p1 and p2"
    instanceVariables: 'p1 p2'.
```

**Listing 9.** Implementation of `PPSequence > > parse:` in PetitParser.

```
PPSequence>>parse: anInputStream
| result1 result2 |
result1 ← p1 parse: anInputStream.
result1 ifFailure: [ ↑ result1 ].
result2 ← p2 parse: anInputStream.
result2 ifFailure: [ ↑ result2 ].

"return array with both results"
↑ { result1 . result2 }
```

### 5.2. Implementation of BoundedSeas in PetitParser

To support bounded seas, we changed the interface of the `parse: anInput` method to `parse: aPPContext`. `PPContext` is an object that provides access to the stack of invoked expressions. `PPContext` as well implements the interface of the `InputStream` so that it can be used as `InputStream`. In order to manage the stack of invoked expressions, a parser is dispatched via `PPContext>>invoked:` and a value is returned via `PPContext>>return:` or `PPContext>>fail:`.

<sup>7</sup> <http://scg.unibe.ch/research/IslandParsing/CLSS2015>

<sup>8</sup> <http://moosetechnology.org>

<sup>9</sup> <http://smalltalkhub.com/#!/~Moose/PetitParser>

<sup>10</sup> <https://github.com/petitparser/java-petitparser>

<sup>11</sup> <https://github.com/petitparser/dart-petitparser>

PPBoundedSea is defined as in Listing 10. Even though a bounded sea consists of a sequence of three parsers, it has only one instance variable `island`, before-water and after-water being created dynamically depending on the state of PPContext. The `parse:` method of PPBoundedSea is in Listing 11. The three phases of `parse:` correspond to the phases in Definition 1. In order to detect an overlapping sea, there is a check in `parseBeforeWater:` (see Listing 12).

**Listing 10.** PPBoundedSea has only one instance variable `island`, before and after-water are created dynamically, depending on the state of the PPContext.

```
PPParser subclass: #PPBoundedSea
  instanceVariables: 'island'.
```

**Listing 11.** Implementation of a `parse:` method in PPBoundedSea. The three phases correspond to the phases in Definition 1.

```
PPBoundedSea>>parse: aPPContext
| result1 result2 result3 |
aPPContext invoked: self.
"Phase One"
result1 ← self parseBeforeWater: aPPContext.
result1 ifFailure: [
  ↑ aPPContext fail: 'boundary or island not found'
].

"Phase Two"
result2 ← island parse: aPPContext
result2 ifFailure: [
  ↑ aPPContext fail: 'island not found'
]

"Phase Three"
result3 ← self parseAfterWater: aPPContext.
result3 ifFailure: [
  ↑ aPPContext fail: 'boundary not found'
].

↑ aPPContext return: { result1 . result2 . result3 }
```

**Listing 12.** Implementation of a `beforeWater:` method in PPBoundedSea.

```
PPBoundedSea>>parseBeforeWater: aPPContext
| next |
"Catch Overlapping Seas Problem"
aPPContext seasOverlap ifTrue: [
  ↑ nil
].
next ← aPPContext next.
↑ self goUpTo: island / next.
```

PPContext manages the parsing expression invocation stack, computes the next set and detects the overlapping seas. Thanks to the fact that the method invocation stack can be accessed in the Pharo environment, PPContext can reuse the method invocation stack to access the invoked expressions. Because the method invocation stack does not contain the invoked position, PPContext manages this separately and only for PPBoundedSea parsers (see Listing 13). Overlapping seas can be detected trivially (see Listing 14). The NEXT function implementation follows straightforwardly from recursive Definition 6 (see Listing 15).

**Listing 13.** Implementation of an `invoked:` method and `return:` method in PPContext.

```

PPContext>>invoked: parser
  self assert: parser isBoundedSea.
  self invokedPositions push: self position.

PPContext>>return: parser
  self assert: parser isBoundedSea.
  self invokedPositions pop.

```

**Listing 14.** Implementation of a seasOverlap method in PPContext.

```

PPContext>>seasOverlap
  ↑ self invokedPositions top ==
  self invokedPositions secondTop

```

**Listing 15.** Fragment of a next method in PPContext.

```

PPContext>>next
| stack |
stack ← self expressionStackFrom: thisContext.
↑ self next: stack into: Set new

PPContext>>next:stack into: set
  "first sequence case: e1 on top, e1e2 second top,
  e2 isNotNullable then NEXT = {e2}"
  (stack secondTop isSequence and:
  [ stack secondTop first == stack top ] and:
  [ stack secondTop second isNotNullable not ]) ifTrue: [
    set add: stack secondTop second.
    ↑ set
  ]
  ...
  "repetition case"
  (stack secondTop) isRepetition ifTrue: [
    set add: stack pop.
    ↑ self next: stack into: set
  ]

```

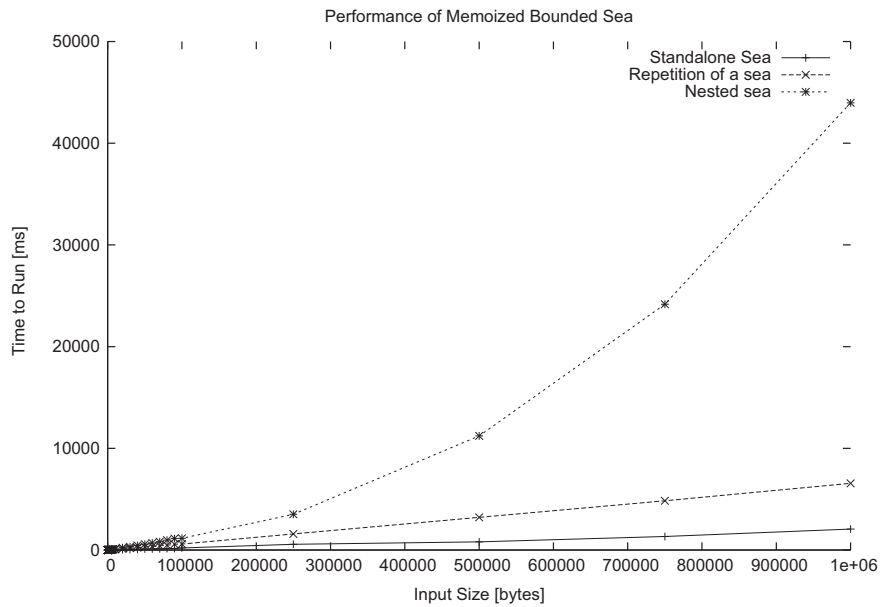
### 5.3. Performance

In this section we briefly report on the performance of bounded seas. We focus on the time complexity of the three different placements of a sea: standalone seas, repetition of a sea and a nested sea.

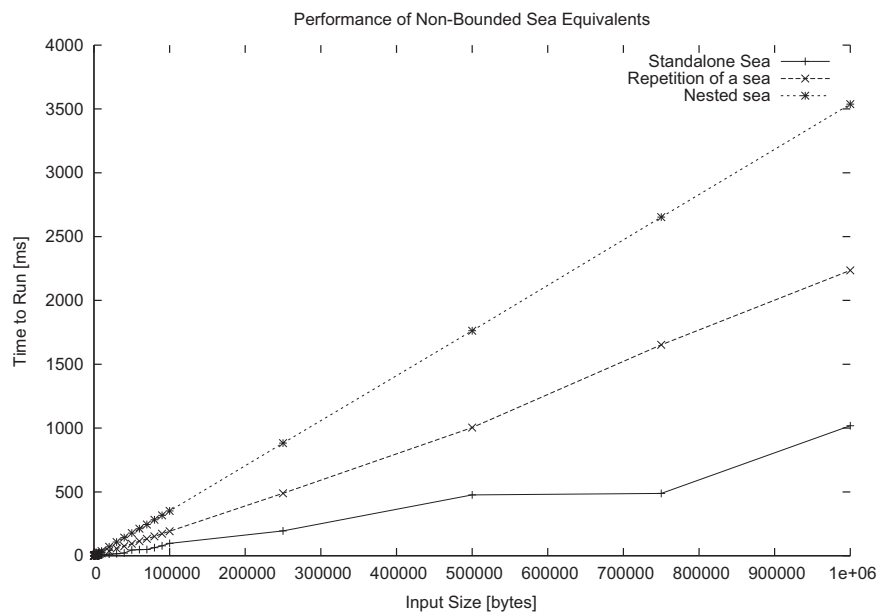
We performed measurements on the following parsers and inputs:

1. *Stand-alone sea*  $\sim 'a' \sim$  searches for the island "a" in an input. An input consists of randomly generated string of dots . (representing water) and a single character "a" at a random position.
2. *Repetition of a sea*  $\sim 'a' \sim +$  searches for sequences of islands "a" in an input. An input consists of a randomly generated string of dots . (for water) and island characters "a", e.g., "..a.....a....aa..".
3. *Nested sea*  $\text{block} \leftarrow \sim \{ 'block+ / \sim e \sim ' \} \sim +$  searches for sequences of nested blocks in an input. An input consists of block starting with "{" and ending with "}". A block contains a possibly empty sequence of other blocks, e.g., "{...{...}{...}}".

Fig. 1 shows that the time complexity is linear compared to the input size for a stand-alone sea and a repetition of a sea. For the nested sea, we measured an exponential complexity. All of the measured parsers used a memoized version [13,14] of bounded seas.



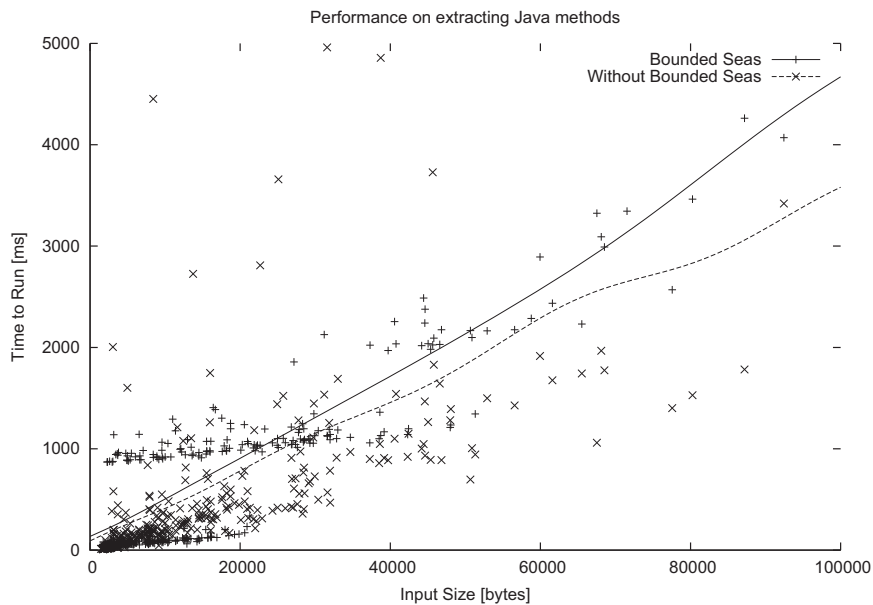
**Fig. 1.** The performance comparison of memoized bounded seas for a stand-alone sea  $\sim 'a' \sim$ , a repetition of a sea  $\sim 'a' \sim +$  and for nested sea on randomly generated inputs.



**Fig. 2.** The performance comparison of non-bounded sea equivalents for a stand-alone sea  $\sim 'a' \sim$ , a repetition of a seas  $\sim 'a' \sim +$  and for a nested sea on randomly generated inputs.

Fig. 2 reports on equivalents of bounded seas implemented without using the bounded sea operator. The complexity of these parsers is linear. We see that there is room for improvement and this still remains an open issue.

In the case studies we performed (Sections 7 and 8), bounded seas showed performance comparable to the non-bounded seas versions. We assume that the typical size of the input and parser complexity used in the case studies is below the threshold where the exponential complexity manifests itself. To support our assumption, Fig. 3 compares a bounded sea parser and a non-bounded sea equivalent. Both parsers extract Java methods from Java standard library files (details about extracting Java methods are provided in Section 7) with a comparable time performance.



**Fig. 3.** The performance comparison of a bounded sea parser and its non-bounded sea equivalent on extracting Java methods from approximately four hundred files from the Java standard library.

## 6. Discussion

In this section we discuss some implementation decisions of bounded seas as well as an implementation of bounded seas for generalized LL grammars and parser combinator libraries. We also discuss why the NEXT set is needed instead of the simpler FOLLOW sets.

### 6.1. Bounded seas as meta-syntactic sugar

Bounded seas can be implemented in two ways: as a meta-syntactic sugar or as a parser extension. In this work we take the latter approach. There are several reasons why we did not choose a grammar transformation that transforms a sea expression into a standard PEG expression.

First, it is easier to implement the NEXT function and detect overlapping seas during parsing than detecting overlapping seas statically and transforming the boundary in such a way that the overlapping seas problem cannot arise. Second, PetitParser is a very agile framework where a parser can be updated simply by changing an object reference at any time. Furthermore, the graph of parser combinators corresponds exactly to the grammar, which makes PetitParser easy to understand and debug. Grammar transformation would add an extra level of complexity into the implementation and it would complicate comprehension and debugging.

### 6.2. Generalized LL parsing

In this paper we have discussed bounded seas for PEGs. However, the essence of bounded seas is not in the grammar formalism used but in the fact that water is specific for each island and it is computed automatically from a stack of invoked expressions. We argue that bounded islands are useful for Context Free Grammars (CFGs) [15] as well.

The key difference between PEGs and CFGs is that CFGs may return ambiguous results whereas PEGs cannot. Implementing an island grammar as a CFG may lead to ambiguous results even though only one of the results is desired. The undesired, remaining results are present only because of vaguely defined water. This is problematic since it is hard to decide which of the results is the correct one.

Bounded seas eliminate ambiguities by adopting a more precise definition of water. Water of a bounded sea never consumes any input that might be valid in a given parsing context. Even though we define a bounded sea with an island 'y' and we run such a rule on the input "xyz", the water of the bounded sea consumes only "x", never "xyz", thus avoiding ambiguities.

Generalized LL Parsing [16] can handle any CFG, allows all the choices of CFGs to be explored in parallel, and, in case of ambiguity returns all possible results. Bounded seas can be implemented in a GLL parser because their top-down nature allows for a stack of parsing expressions and they support syntactic predicates used in a boundary.

### 6.3. Integration with monadic parser libraries

To compute  $NEXT(e_1)$  in a sequence  $e_1e_2$  we need to know what  $e_2$  is. However in some cases, e.g., in monadic parser combinator [5] libraries,  $e_2$  could be a closure. For example, when parsing the HTTP header containing a value indicating the length of a content,<sup>12</sup> we might read that value and use it to create the parser that reads the content itself (i.e., by length-times reading a character):

```
length >>= \length -> applyNTimes length readChar
```

Now, if we want to use a bounded sea to extract the length, i.e.,

```
~length~ >>= \length -> applyNTimes length readChar
```

we cannot determine the boundary of  $\sim\text{length}\sim$ , because it depends on the result of the `length`. As a consequence we can only use bounded seas in a sequence  $e_1e_2$  if we can compute  $e_2$  before parsing the  $e_1$ .

Bounded seas do not allow for context-sensitive dependencies between an island and its border, but for one exception: when a sea is bounded by another sea, we disable water if another water is already invoked at the same position.

### 6.4. FOLLOW vs. NEXT

The NEXT function introduces extra complexity into bounded seas, even though it resembles the FOLLOW function from LL parsing theory [17, pp. 235–361]. The key difference between FOLLOW and NEXT is that the former returns only terminals, while the latter returns parsing expressions.

Why is it not sufficient to use the well-known FOLLOW sets instead of the more complicated NEXT function? The reason is that the right context (boundary) of a sea is in general an  $LL(k), k \geq 1$  language, and a simple FOLLOW set is not usually sufficient to recognize the boundary.

As an example, consider the grammar from Listing 7. The boundary of `class` is  $NEXT(\text{class}) = \{\sim\text{class}\sim, \sim\text{struct}\sim, \text{mainMethod}\}$ . Suppose that instead we take as the boundary of `class` its FOLLOW set, i.e.,  $FOLLOW(\text{class}) = \{\text{'class'}, \text{'struct'}, \text{'public'}\}$ . If there are other elements in the input that start with `'public'` (e.g., `"public int i=0;"`), they will be indistinguishable from the `mainMethod` and the water of bounded seas would finish in an invalid position.

Bounded seas are supposed to work only with a skeleton of an original grammar with as little information as possible. Therefore, information about other input that can interfere with a boundary (e.g., `"public int i=0;"`) is not usually available. If bounded seas are provided with a baseline grammar this would not be problem as the techniques described by Klusener and Lämmel [18] can then be applied.

## 7. Java parser case study

The goal of this case study is to demonstrate the suitability of bounded seas for extracting data from Java sources without any baseline grammar provided. First we focus on a simpler task without considering nested classes. Because bounded seas target extensibility we subsequently investigate the effort required to extend the parser with nested classes.

We compare four kinds of Java parsers and we measure how well can they extract classes and their methods from a Java source code.<sup>13</sup>

1. **PetitJava** is an open-source Java parser using PetitParser [3] provided by the Moose analysis platform community [19]. We used version 159.<sup>14</sup>
2. **Naïve Island Parser** is an island parser with water defined simply as the negation of the island we are searching for. The sea rules in this parser can be reused, because they do not consider their surroundings and they are grammar-independent. The sea rules are defined in a simple form: consume input until an island is found, then consume an island.
3. **Advanced Island Parser** is a more complex version of the naïve island parser. The water is more complicated to prevent the most frequent failures of island parsers. The sea rules in this parser are hard-wired to the grammar and cannot be reused. The sea rules are customized for a particular islands.
4. **Island Parser with Bounded Seas** is an island parser written using bounded seas. The sea rules were defined using the sea operator.

<sup>12</sup> [https://en.wikipedia.org/wiki/List\\_of\\_HTTP\\_header\\_fields](https://en.wikipedia.org/wiki/List_of_HTTP_header_fields)

<sup>13</sup> The case study and instructions can be found at the following prepared web-page: <http://scg.unibe.ch/research/IslandParsing/CLSS2015>.

<sup>14</sup> <http://smalltalkhub.com/#!/~Moose/PetitJava/>

The PetitJava parser parses Java 6 code. All the island parsers (island, advanced and bounded) are very similar, with approximately 20 rules per each. PetitJava itself contains over 200 rules. The island parsers were designed to extract classes and the methods that belong to them. None of the parsers was optimized to provide a better performance.

We compare the three island parsers (almost identical in a structure) written by the first author. It is very likely that the advanced island parser can be modified to achieve better precision and better performance, but at the cost of considerable engineering work. We demonstrate that naïve water rules do not work and that the advanced version of water is needed. We further show that with bounded seas we can obtain high precision and performance without needing to define an advanced island parser. Finally, we show that extending an island parser is a highly demanding task, unless bounded seas are used.

*Test data:* For our case study we randomly selected 50 files ( $N$ ) containing 50 classes from the JDK 6 library. These 50 classes contain 81 nested classes and a total of 1380 methods  $M$ . We extract the reference data using the Verveine<sup>15</sup> parser.

Each parser returns a set  $m$  of fully qualified method names,<sup>16</sup> some of which are true positives  $m_{tp}$ . If a parser fails, an error is returned and the set of all errors is  $e$ . Failure is treated as though no classes or methods were found. We measure precision  $P = |m_{tp}|/|m|$ , recall  $R = |m_{tp}|/|M|$ , error rate  $err = |e|/N$  and time per file  $t = t_{total}/N$ .

### 7.1. Without nested classes

First of all, we evaluate our parsers on extracting method names without considering the nested classes and their methods. We can easily skip the nested classes by defining properly paired blocks starting with '{' and ending with '}' and ignoring everything inside.

*Results:* As we see in Table 2, PetitJava parser provides perfect precision, but recall is poor because of the high error rate.<sup>17</sup> On the other hand, the error rate of all island parsers (island, advanced and bounded) is very low,<sup>18</sup> but precision and recall are not perfect, even though they are relatively good. Amongst the imprecise parsers, the Bounded parser provides the best precision and recall.

### 7.2. With nested classes

In this step, we extend our island parsers to include nested classes and their methods. We do this by making a single change, where we extend the `classBody` rule from this<sup>19</sup>:

```
classBody ← '{' method island * '}'
```

to this:

```
classBody ← '{' (method / class) island * '}'
```

*Results:* As we see in Table 3 the PetitJava parser performs as in the previous case. Yet the imprecise parsers (Island, Advanced) start to struggle. Their error rate has increased and recall has dropped dramatically. The errors were mostly due either to parsing timeouts (when parsing took more than ten seconds per file) or various parsing errors. On the other hand, the Bounded parser maintains high precision and recall, zero error rate, and improves time per file slightly.

In Table 3 we also measured the Advanced' parser, which made use of refined rules for water to take into account the grammar changes.<sup>20</sup> This improved recall, parsing time and the error rate. We would, however, need to invest even more effort to reach the quality of the Bounded parser.

## 8. Ruby parser case study

The standard approach to recognize the structure of the input is to track all language elements that affect structure, as we did in the Java case study where we defined a rule for blocks. As it turns out, almost anything can affect the structure of a Ruby program. For this reason, we turned to indentation as it turns out to be a good proxy for structure [20]. In this case study we focus on using bounded seas to extract the structure of a Ruby program by exploiting indentation.<sup>21</sup>

<sup>15</sup> <https://gforge.inria.fr/projects/verveinej>

<sup>16</sup> <http://docs.oracle.com/javase/specs/jls/se7/html/jls-6.html#jls-6.7>

<sup>17</sup> The PetitJava failures are due to bugs in the grammar specification.

<sup>18</sup> Failures of the imprecise parsers are due to parsing timeout (set to 10 s).

<sup>19</sup> `island` here creates either an island, an advanced island or a bounded sea depending on the parser we use.

<sup>20</sup> We investigated the reasons for failures and added an extra boundary to `classBody`.

<sup>21</sup> The case study and instructions can be found at the prepared web-page: <http://scg.unibe.ch/research/IslandParsing/CLSS2015>.



### 8.1. The dangling end problem

Ruby poses interesting parsing challenges even for imprecise parsers. The biggest problem we faced is the *dangling end problem*: Normally a control structure like an if statement terminates with `end`. However there is also an *if modifier*, as in `return error if check?`, which does not require an `end`.

Such modifiers pose problems for parsing. There exist numerous such modifiers in Ruby,<sup>22</sup> which resemble conditional blocks, but have a different syntax. From the perspective of an imprecise parser, it is hard to distinguish between a modifier, loops and conditional blocks.

Ruby structures (such as classes, methods, blocks) end with the `'end'` keyword (see Listing 16). To capture the structure of Ruby code, we need to define rules for these structural elements, including conditional blocks and others, such as loops, do blocks, and brace pairs.

**Listing 16.** Example of a Ruby code.

```
class Shape
  def draw
    if (x > 0)
      do_something()
    end
  end
end
```

Ruby modifiers are not paired with any `'end'` as we can see in Listing 17. If we incorrectly pair `'end'` we change the structure of a program. Unfortunately, it is hard to recognize when `'if'` belongs to a modifier and when to a conditional block, unless we specify a complete grammar to recognize all the constructs that it could possibly modify.

**Listing 17.** Example of Ruby code where `'if'` is not paired with any `'end'`.

```
class Shape
  def draw
    return error if check?
    if (x > 0)
      do_something
    end
  end
end
```

### 8.2. Indentation

It is known that indentation is a good proxy for structure in programming languages [20]. We can exploit this fact to define a context-sensitive parser that uses both indentation and bounded seas to recognize modifiers. From the perspective of indentation, modifiers look like loops or conditional blocks with a single line scope.

Because PetitParser produces scannerless parsers [12] and it doesn't use any preprocessing (*i.e.*, tokenizing), an indentation-sensitive parser is context-sensitive since the question whether code is indented or not depends on the results of previously invoked parsers.

Inspired by Landin's offside rule [21], indentation in PetitParser uses a stack of indentation levels and adds extra layout-oriented parsing expressions (*e.g.*, `align`, `inOffside`). These expressions consult the stack and the current indentation level to verify that the input complies with the given layout criteria [22].

Although we have seen earlier that bounded seas are incompatible with monadic parser libraries where a boundary may depend on what has been parsed earlier, indentation parsing is a special case that does not interfere with bounded seas. As we shall see, the use of indentation-sensitive parsers simplifies the implementation of the parsers and even improves the overall performance.

<sup>22</sup> <http://docs.huihoo.com/ruby/ruby-man-1.4/syntax.html#if-mod>, <http://docs.huihoo.com/ruby/ruby-man-1.4/syntax.html#unless-mod>, <http://docs.huihoo.com/ruby/ruby-man-1.4/syntax.html#while-mod>, <http://docs.huihoo.com/ruby/ruby-man-1.4/syntax.html#until-mod>

**Table 2**

Precision, recall error rate and time of the four tested parsers without considering nested classes.

Parser	Precision	Recall	Time (ms)	Error rate
PetitJava	1.00	0.71	308	0.28
Island	0.87	0.90	1225	0.04
Advanced	0.92	0.90	1336	0.04
Bounded	0.96	1.00	941	0.00

**Table 3**

Precision, recall, time and error rate including nested classes.

Parser	Precision	Recall	Time (ms)	Error rate
PetitJava	1.00	0.67	299	0.28
Island	0.87	0.54	934	0.12
Advanced	0.94	0.32	1734	0.34
Advanced'	0.91	0.68	847	0.03
Bounded	0.97	0.99	627	0.00

We define a context-sensitive grammar that recognizes modules, classes, methods and class methods in Ruby code by utilizing indentation and bounded seas. The scope of a class or method extends as far as code appears to the right of the class or method declaration (*i.e.*, in the onside position). The `class` definition is in [Listing 18](#).

**Listing 18.** Indentation sensitive definition of a Ruby class.

```
class ← setOffsideLine, 'class' identifier
      ~(class / method)~ onside *
      unsetOffsideLine
```

### 8.3. Parsing results

In this section we report on the complexity, performance, precision and recall of three parsers: a classical *island parser* (46 grammar rules, 9K characters), a *bounded parser* that does not utilize indentation (41 rules, 8.5K characters), and an *indent bounded parser* that utilizes indentation (27 rules, 4K characters).

The island parser and the bounded parser are almost identical. The sea parser uses bounded seas, while the island parser uses manually defined islands and water. From the number of methods, we can see that indentation simplifies the implementation. The bounded parser and the island parser must implement additional rules to recognize the dangling end.

The bounded parser shows its flexibility here. For example, the method definition in the bounded grammar does not require `arguments` ([Listing 19](#)) contrary to the method definition in the island grammar ([Listing 20](#)).

**Listing 19.** Method definition in the bounded grammar.

```
method ← 'def' name arguments primary* 'end'
primary ← (! (comment / keyword / modifier / ...)
           #any)*
           (method / class)
arguments ← ...
```

**Listing 20.** Method definition in island grammar.

```
methodDef ← 'def' name primary* 'end'
primary ← ~method / class~
```

To measure precision and recall, we used *jruby-parser*<sup>23</sup> as a reference parser. We compared the structure (modules, classes, methods and class methods) of Ruby code as detected by *jruby-parser* with the structure detected by our parsers. We describe the structure as a set of methods where each method is prepended with a path consisting of other methods, classes and modules depending on the location of the method in a code, similar to Java's fully qualified names. For example:

```
<module>graphics.<class>Shape.<method>draw
```

refers to a method `draw` defined in the class `Shape`. The `Shape` belongs to the `graphics` module. On the other hand:

```
<class>Shape.<class>Renderer.<class-method>Instance
```

refers to the class-side method `Instance` of the `Shape`'s inner class `Renderer`.

**Test data:** We performed our study on a sample of  $N=100$  files of six popular projects on Github: Rails,<sup>24</sup> Discourse,<sup>25</sup> Diaspora,<sup>26</sup> Cucumber,<sup>27</sup> Valgrant<sup>28</sup> and Typhoeus.<sup>29</sup> The sampled files contain a total of 520 methods.

Parsers return a set of fully qualified methods  $m$ , where some of them are true positives  $m_{tp}$ . If a parser fails, an error is returned. The set of all errors is  $e$ . We measure precision  $P = |m_{tp}|/|m|$ , recall  $R = |m_{tp}|/|M|$ , error rate  $err = |e|/N$  and time per file  $t = t_{total}/N$ . Failure is treated as though no classes or methods are found.

Table 4 shows precision and recall are rather high in all of the cases. The island parser has perfect precision, but recall is not perfect due to some failures. The bounded parser has worse precision, because it did not fail for one of the inputs, but misplaced the methods into the wrong module. The indent bounded parser can parse any of the files with very high precision and recall. It misplaced only one<sup>30</sup> of all the methods.

As we have seen, the island parser contains 46 rules, the bounded parser 41, and the indent parser 27. This shows that both bounded seas and indentation help to reduce the complexity of the Ruby grammar. Bounded seas perform better than traditional islands. The indentation parser is even better than the bounded parser, because fewer rules are needed to determine the boundaries.

## 9. Related work

**Agile parsing:** Agile parsing [8] is a recent paradigm for source analysis and reverse engineering tools. In agile parsing the effective grammar used by a particular tool is a combination of two parts: the standard base grammar for the input language, and a set of explicit grammar overrides that modify the parse to support the task at hand. There are several agile parsing idioms: (i) *rule abstraction* (grammar rules can be parametrized); (ii) *grammar specialization* (grammar rules can be specialized based on the semantic needs); (iii) *grammar categorization* (to deal with context-free ambiguities); (iv) *union of grammars* (to unify multiple grammars); (v) *markup* (to match and mark chunks of interest); (vi) *semi-parsing* (to define islands and lakes); and (vi) *data structure grammars* (separate grammars that hold auxiliary data structures).

The semi-parsing idiom [8] uses the *not* predicate to prevent water from consuming islands. This approach is the same as that taken by bounded seas. Contrary to the semi-parsing idiom, bounded seas are able to infer the predicates on their own. The agile parsing idioms are based on a transformation of a well-defined baseline grammar, whereas bounded seas do not expect such a well-defined grammar and must infer the predicates only from the available *skeleton*.

**Island grammars:** Island grammars were proposed by Moonen [1] as a method of semi-parsing to deal with irregularities in the artifacts that are typical for the reverse engineering domain. Island grammars make use of a special syntactic rule called *water* that can accept any input. Water is annotated with a special keyword `avoid` that will ensure that water will be accepted only if there is no other rule that can be applied.

Contrary to Moonen, we propose boundaries (based on the NEXT function) that limit the scope in which water can be applied. Because each island has a different boundary, our solution does not use the single water rule; instead our water is tailored to each particular island.

**Non-greedy rules:** Non-greedy operators are well-known from regular expressions introduced in Perl.<sup>31</sup> `??`, `*?`, and `+?` are non-greedy versions of `?`, `*` and `+`, which match as little of a string as possible while preserving the overall match. The backtracking algorithm admits a simple implementation of non-greedy operators: try the shorter match before the longer

<sup>23</sup> <https://github.com/jruby/jruby-parser>

<sup>24</sup> <https://github.com/rails/rails>

<sup>25</sup> <https://github.com/discourse/discourse>

<sup>26</sup> <https://github.com/diaspora/diaspora>

<sup>27</sup> <https://github.com/cucumber/cucumber>

<sup>28</sup> <https://github.com/mitchellh/vagrant>

<sup>29</sup> <https://github.com/typhoeus/typhoeus>

<sup>30</sup> If a method declaration with a modifier follows an inner class defined on a single line, the method with the modifier is incorrectly assigned to the inner class.

<sup>31</sup> <http://perldoc.perl.org/perlre.html>

**Table 4**

Precision, recall error rate and time of compared parsers.

Parser	Precision	Recall	Time (ms)	Error rate
Island parser	1.00	0.96	495	0.03
Bounded parser	0.97	0.96	283	0.01
Indent bounded parser	0.99	0.99	203	0.00

one. For example, in a standard backtracking implementation,  $e?$  first tries using  $e$  and then tries not using it;  $e??$  uses the other order.<sup>32</sup>

Non-greedy operators are also available in ANTLR as parser operators. A non-greedy parser matches the shortest sequence of tokens that preserves a successful parse for a valid input sentence. Contrary to regular expressions, a non-greedy parser never makes a decision that will ultimately cause valid input to fail later on during the parse. The central idea is to match the shortest sequence of tokens that preserves a successful parse for a valid input sentence.<sup>33</sup>

Bounded seas are distinct from non-greedy rules in two ways. First, bounded seas do not require globally correct decisions, since they are not available in traditional PEGs. Though PEGs can backtrack while choosing between alternatives, once the choice is made it cannot be changed, thus making a globally correct decision impossible. In order to realize non-greedy repetitions, PEGs feature predicates, which have to be specified by an engineer (as illustrated in Section 2). Bounded seas remove the burden of predicates from a language engineer by computing the NEXT set automatically.

Second, bounded seas target transparent composability. A language engineer can treat a bounded sea like any other PEG rule without bothering about its implementation. For example, the following grammar can be easily modified by changing the body to  $\text{body} \leftarrow \text{sea}^*$ ,  $\text{body} \leftarrow \text{sea}?$  or  $\text{body} \leftarrow \text{sea}? \text{sea}?$ .

```
start ← ('begin' body 'end')*
body  ← sea
sea   ← ~sea~
```

If we define  $\text{sea}$  using lazy repetition  $*?$ , the normal  $\text{sea}$  can be defined as

```
start ← ('begin' body 'end')*
body  ← sea
sea   ← .*? 'body' .*?
```

the optional version as

```
start ← ('begin' body)*
body  ← sea
sea   ← .*? ('body' | 'end')
```

the repetition version as

```
start ← ('begin' body 'end')*
body  ← sea
sea   ← ('body' | 'end')*?
```

and the sequence of two optional seas as

```
start ← ('begin' body)*
body  ← sea1
sea1  ← .*? ('body1' sea2 | 'body2' 'end' | 'end')
sea2  ← .*? ('body2' 'end' | 'end')
```

**Noise skipping parsing:** GLR\* is a noise-skipping parsing algorithm for context-free grammars able to parse any input sentence by ignoring unrecognizable parts of the sentence [23]. The parser nondeterministically skips some words in a

<sup>32</sup> <https://swtch.com/~rsc/regexp/regexp1.html>

<sup>33</sup> <https://theantlruguy.atlassian.net/wiki/display/ANTLR4/Wildcard+Operator+and+Nongreedy+Subrules>

sentence and returns the parse with fewest skipped words. The parser is a modification of Generalized LR (Tomita) parsing algorithm [24].

The GLR\* application domain is parsing of spontaneous speech. Contrary to bounded seas, GLR\* itself decides what is noise (water in our case) and where it is. In the case of bounded seas the positions of the noise (water) are explicitly defined.

*Fuzzy parsing:* The term fuzzy parser was coined for Sniff [25], a commercial C++ IDE that uses a hand-made top-down parser. Sniff can process incomplete programs or programs with errors by focusing on symbol declarations (classes, members, functions, variables) and ignoring function bodies. In linguistics or natural language processing [26], the notion of fuzzy parsing corresponds to an algorithm that recognizes fuzzy languages.

The semi-formal definition of a fuzzy parser was introduced by Koppler [27]. Fuzzy parsers recognize only parts of a language by means of an unstructured set of rules. Compared with whole-language parsers, a fuzzy parser remains idle until its scanner encounters an anchor in the input or reaches the end of the input. Thereafter the parser behaves like a normal parser.

*Skeleton grammars:* Skeleton grammars [18] address the issue of false positives and false negatives when performing tolerant parsing by inferring a tolerant (skeleton) grammar from a precise baseline grammar.

Our approach tackles the same problem as skeleton grammars: improving the precision of island grammars. They both maintain the composability property and both can be automated. Skeleton grammars use the standard first and sets known from standard parsing theory [17, pp. 235–361] for synchronization with the baseline grammar.

Bounded seas do not require a precise baseline grammar and they have to find point of synchronization based only on the main grammar itself. Therefore the main grammar has to contain all the relevant information (e.g., when extracting classes and methods with bounded seas block definitions are essential to place methods properly). Because the main grammar of bounded seas is typically far from complete, bounded seas use the NEXT set (instead of first and follow) to reach the required precision. If bounded seas are provided with the baseline grammar, the boundaries can be computed from the baseline.

*Bridge parsing:* Bridge parsing is a novel, lightweight recovery algorithm that complements existing recovery techniques [28]. Bridge parsing extends an island grammar with the notion of bridges and reefs. Islands denote tokens that open or close scopes. Reefs are attributed tokens and they add information (e.g., indentation) to nearby islands. Islands and reefs are created in a tokenizing phase. Bridges connect matching opening and closing islands in a bridge-building phase. The corresponding islands are searched with the help of reefs (e.g., indentation can be used to find matching brackets). If some islands are not connected (e.g., if the opening or closing scope island is missing), the bridge repair phase tries to repair them with the help of information from reefs.

The focus of bounded seas is on data extraction rather than on error recovery and bounded seas are missing advanced error-recovery techniques available in the bridge parsing. Bounded seas are meant to be used on valid inputs without errors. If an erroneous chunk appears, bounded seas skip such a chunk until a valid chunk is found. To our best knowledge, techniques used in bridge parsing are complementary to bounded seas and might help improve precision of bounded seas on erroneous inputs.

*Permissive grammars:* The main idea of permissive grammars [29,30] is to derive a permissive grammar from a standard grammar. Such a permissive grammar accepts programs with minor errors (missing brackets, etc.). A permissive grammar is also a normal grammar and can be tweaked by the language engineer. Using a specialized version of the GLR algorithm, both syntactically correct and incorrect programs can be efficiently parsed using these grammars [29].

Contrary to bounded seas, which target the area of rapid data extraction, permissive grammars are supposed to help IDE developers with interactive parsing and error recovery as the user is writing a program. Similar to bounded seas, permissive grammars extend the concept of island grammars and use water for error recovery. Even though bounded seas can be used to skip over noise in an input, bounded seas handle missing or misspelled input simply by ignoring the whole erroneous chunk until a valid chunk is found. Permissive grammars try to find the best way to fix an erroneous chunk (and not only skip over it).

## 10. Conclusion

In this paper we have presented bounded seas — composable, reusable, robust and easy to use islands. Contrary to the traditional approach of island parsing, bounded seas compute the scope within which water can consume the input. We have extended the semantics of PEGs to implement useful and practical bounded seas. Boundaries are computed by a NEXT function, inspired by the follow function from standard parsing theory. The automation of the process that creates the bounded sea ensures that bounded seas are easy to use and are not error-prone. Bounded seas as presented in this work are context-sensitive.

As a validation of the composability and reusability of bounded seas, we have presented an implementation of bounded seas as a parser combinator in the PetitParser framework. Furthermore we have presented two case studies applying bounded sea parsers to extracting method names from Java and Ruby code, and we have compared these parsers to conventional parsers based on a precise grammar and based on island grammars. We show that bounded seas provide both good precision and performance.

## Acknowledgments

We gratefully acknowledge the financial support of the Swiss National Science Foundation for the project “Agile Software Assessment” (SNSF project no. 200020-144126/1, January 1, 2013 – December 30, 2015).

We also thank the anonymous referees for their invaluable comments.

## Appendix A. Parsing expression grammars

PEGs were first introduced by Ford [4] and the formalism is closely related to top-down parsing. PEGs are syntactically similar to CFGs [15], but they have different semantics. The main semantic difference is that the choice operator in PEG is ordered – it selects the first successful match – while the choice operator in CFG is ambiguous. PEGs are composed using the operators in Table A1.

**Definition 7** (PEG definition). We use the standard definition as suggested by Ford [4]. A *parsing expression grammar* (PEG) is a 4-tuple  $G = \{N, T, R, e_s\}$ , where  $N$  is a set of nonterminals,  $T$  is a set of terminals,  $R$  is a set of rules,  $e_s$  is a start expression.  $N \cap T = \emptyset$ . Each  $r \in R$  is a pair  $(A, e)$ , which we write  $A \leftarrow e$ , where  $A \in N$ ,  $e$  is a parsing expression. Parsing expressions are defined inductively. If  $e$ ,  $e_1$  and  $e_2$  are parsing expressions, then so is:

- $\epsilon$ , the empty string
- $a$ , any terminal where  $a \in T$
- $A$ , any nonterminal where  $A \in N$
- $e_1 e_2$ , a sequence
- $e_1 / e_2$ , a prioritized choice
- $e^*$ , zero or more repetitions
- $!e$  a not-predicate

The following operators are syntactic sugar:

- *Any character*:  $\cdot$  is character class containing all letters
- *Character class*:  $[a_1, a_2, \dots, a_n]$  character class is  $a_1 / a_2 / \dots / a_n$
- *Optional expression*:  $e?$  is  $e_d / \epsilon$ , where  $e_d$  is desugaring of  $e$
- *One-or-more repetitions*:  $e+$  is  $e_d e_d^*$ , where  $e_d$  is desugaring of  $e$
- *And-predicate*:  $\&e$  is  $!(e_d)$ , where  $e_d$  is desugaring of  $e$

We will use text in quotation marks to refer to terminals e.g., 'a', 'b', 'class'. We will use identifiers `A`, `B`, `C`, `class` or `method` to refer to nonterminals. We will use  $e$  or indexed  $e$ :  $e_1$ ,  $e_2$ , ... to refer to parsing expressions.

**Definition 8** (PEG semantics). To formalize the semantics of a grammar  $G = \{N, T, R, e_s\}$ , we define a relation  $\Rightarrow$  from pairs of the form  $(e, x)$  to the output  $o$ , where  $e$  is a parsing expression,  $x \in T^*$  is an input string to be recognized and  $o \in T^* \cup \{f\}$  indicates the result of a recognition attempt. The distinguished symbol  $f \notin T$  indicates failure.

**Empty:**  $\frac{x \in T^*}{(\epsilon, x) \Rightarrow \epsilon}$

**Terminal (success case):**  $\frac{a \in T, x \in T^*}{(a, ax) \Rightarrow a}$

**Terminal (failure case):**  $\frac{a \neq b, (a, \epsilon) \Rightarrow f}{(a, bx) \Rightarrow f}$

**Nonterminal:**  $\frac{A \leftarrow e \in R \ (e, x) \Rightarrow o}{(A, x) \Rightarrow o}$

**Table A1**  
Operators for constructing parsing expressions.

Operator	Description
' '	Literal string
[]	Character class
.	Any character
(e)	Grouping
e?	Optional
e*	Zero-or-more repetitions of e
e+	One-or-more repetitions of e
&e	And-predicate, does not consume input
!e	Not-predicate, does not consume input
e <sub>1</sub> e <sub>2</sub>	Sequence
e <sub>1</sub> / e <sub>2</sub>	Prioritized choice

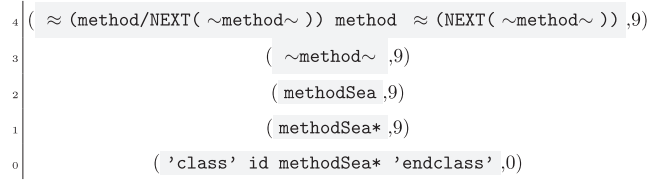


Fig. B1. State of a stack after parsing “class Foo” in the input “class Foo endclass”.

$$\begin{aligned}
 & (e_1, x_1 x_2 y) \Rightarrow x_1 \\
 & (e_2, x_2 y) \Rightarrow x_2 \\
 \text{Sequence (success case): } & \frac{(e_1 e_2, x_1 x_2 y) \Rightarrow x_1 x_2}{(e_1, x_1 y) \Rightarrow x_1} \\
 \text{Sequence (failure case 1): } & \frac{(e_1, x) \Rightarrow f}{(e_1 e_2, x) \Rightarrow f} \\
 \text{Sequence (failure case 2): } & \frac{(e_1, x_1 y) \Rightarrow x_1 \quad (e_2, y) \Rightarrow f}{(e_1 e_2, x_1 y) \Rightarrow f} \\
 \text{Alternation (case 1): } & \frac{(e_1, xy) \Rightarrow x}{(e_1/e_2, x) \Rightarrow x} \\
 \text{Alternation (case 2): } & \frac{(e_1, x) \Rightarrow f \quad (e_2, x) \Rightarrow o}{(e_1/e_2, x) \Rightarrow o} \\
 & (e, x_1 x_2 y) \Rightarrow x_1 \\
 & (e^*, x_2) \Rightarrow x_2 \\
 \text{Repetitions (repetition case): } & \frac{(e^*, x_1 x_2 y) \Rightarrow x_1 x_2}{(e, x_1 y) \Rightarrow x_1} \\
 \text{Repetitions (termination case): } & \frac{(e, x) \Rightarrow f}{(e^*, x) \Rightarrow \epsilon} \\
 \text{Not predicate (case 1): } & \frac{(e, xy) \Rightarrow x}{(!e, xy) \Rightarrow f} \\
 \text{Not predicate (case 2): } & \frac{(e, xy) \Rightarrow f}{(!e, xy) \Rightarrow \epsilon}
 \end{aligned}$$

## Appendix B. Examples

### B.1. Example of abstract simulation

Let us compute the abstract simulation (see [Definition 5](#)) for the following grammar:

```

S      ← E1 E2
E1     ← 'a' / ε
E2     ← 'b' 'c'

```

Because of the recursive nature of the definition, we will compute—for terminals first and we will infer the—for more complex expressions once we have computed—for the simpler ones:

- ‘a’→1 (rule 2), same for ‘b’ and ‘c’
- ‘a’→f (rule 3), same for ‘b’ and ‘c’
- ε→0 (rule 10)
- E1→0 (rule 9)
- E2→1 (rule 5)
- E2→f (rule 6)
- S→1 (rule 5)
- S→f (rule 7)

### B.2. Example of NEXT computation

Let us compute *NEXT* of the *method* island defined in the island grammar in [Listing 5](#). Let us suppose we have already parsed “class Foo” in the input “class Foo endclass”. The stack now looks as shown below in [Fig. B1](#).

Nonterminal:

$$\begin{array}{l}
 1. \quad S \leftarrow \sim a \sim b \sim \in R \\
 2. \quad \begin{array}{c} \dots a \dots b \dots \end{array} \quad \begin{array}{c} 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} o = \\ \dots a \dots b \dots \end{array}
 \end{array}$$


---


$$\begin{array}{c} \dots a \dots b \dots \end{array} \quad \begin{array}{c} 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} o = \\ \dots a \dots b \dots \end{array}$$

**Fig. B2.** The inference rule for nonterminal.

Sequence I (success case):

$$\begin{array}{l}
 1. \quad \begin{array}{c} \dots a \dots b \dots \end{array} \quad \begin{array}{c} 2 \mid ( \sim a \sim , 0 ) \\ 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} x_1 = \\ \dots a \dots \end{array} \\
 2. \quad \begin{array}{c} b \dots \end{array} \quad \begin{array}{c} 2 \mid ( \sim b \sim , 5 ) \\ 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} x_2 = \\ b \dots \end{array}
 \end{array}$$


---


$$\begin{array}{c} \dots a \dots b \dots \end{array} \quad \begin{array}{c} 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} x_1 x_2 = \\ \dots a \dots b \dots \end{array}$$

**Fig. B3.** The inference rule for sequence.

Rewrite according to the Definition 2:

$$\begin{array}{l}
 1. \quad \begin{array}{c} \dots a \dots b \dots \end{array} \quad \begin{array}{c} 2 \mid ( \sim a \sim , 0 ) \\ 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} o = \\ \dots a \dots \end{array}
 \end{array}$$


---


$$\begin{array}{c} \dots a \dots b \dots \end{array} \quad \begin{array}{c} 2 \mid ( \approx (a/NEXT(\sim a \sim)) a \approx (NEXT(\sim a \sim)) , 0 ) \\ 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} o = \\ \dots a \dots \end{array}$$

**Fig. B4.** Rewrite rule according to Definition 2.

Sea Sequence I (success case):

$$\begin{array}{l}
 1. \quad \begin{array}{c} \dots a \dots b \dots \end{array} \quad \begin{array}{c} 3 \mid ( \approx (a/NEXT(\sim a \sim)) , 0 ) \\ 2 \mid ( \approx (a/NEXT(\sim a \sim)) a \approx (NEXT(\sim a \sim)) , 0 ) \\ 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} y_1 = \\ \dots \end{array} \\
 2. \quad \begin{array}{c} a \dots b \dots \end{array} \quad \begin{array}{c} 3 \mid ( a , 2 ) \\ 2 \mid ( \approx (a/NEXT(\sim a \sim)) a \approx (NEXT(\sim a \sim)) , 0 ) \\ 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} y_2 = \\ a \end{array} \\
 3. \quad \begin{array}{c} \dots b \dots \end{array} \quad \begin{array}{c} 3 \mid ( \approx (NEXT(\sim a \sim)) , 3 ) \\ 2 \mid ( \approx (a/NEXT(\sim a \sim)) a \approx (NEXT(\sim a \sim)) , 0 ) \\ 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} y_3 = \\ \dots \end{array}
 \end{array}$$


---


$$\begin{array}{c} \dots a \dots b \dots \end{array} \quad \begin{array}{c} 2 \mid ( \approx (a/NEXT(\sim a \sim)) a \approx (NEXT(\sim a \sim)) , 0 ) \\ 1 \mid ( \sim a \sim b \sim , 0 ) \\ 0 \mid ( S , 0 ) \end{array} \quad \Rightarrow \quad \begin{array}{c} y_1 y_2 y_3 = \\ \dots a \dots \end{array}$$

**Fig. B5.** The inference rule for sequence.



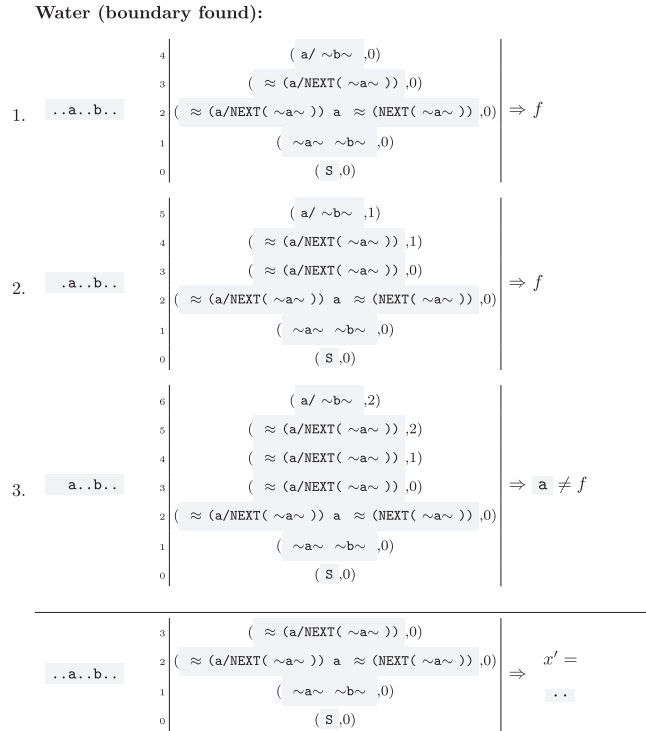


Fig. B6. The inference rule for water.

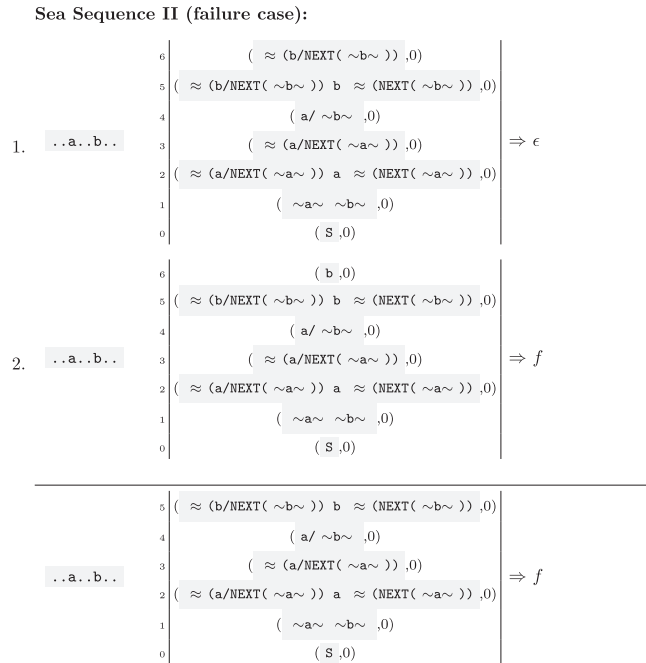


Fig. B7. The inference rule for sequence (failure case).

To parse  $\sim \text{method} \sim$  we need to compute  $NEXT(\sim \text{method} \sim)$ . We do this in the following steps.<sup>34</sup>

1. Initialize:  $NEXT(\text{methodSea}) = \{\}, n = 2$

<sup>34</sup> To simplify, we start from stack position 2, because  $NEXT(\sim \text{method} \sim)$  (stack position 3) is trivially  $NEXT(\text{methodSea})$  (stack position 2).

Water (Overlapping Seas Case):

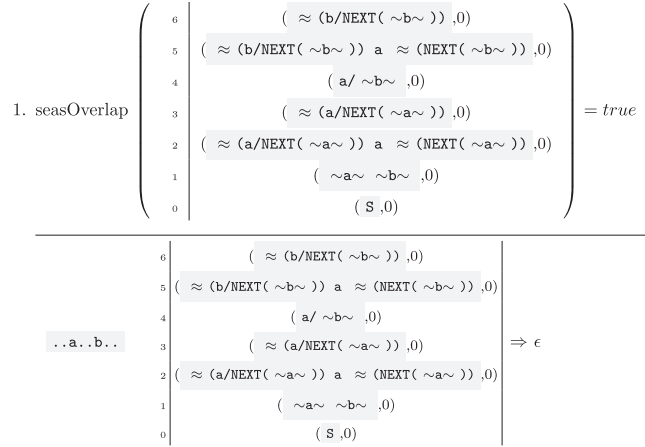


Fig. B8. The inference rule for overlapping seas.

## 2. Check stack:

$s_n = s_2 = \text{methodSea}$  and

$s_{n-1} = s_1 = e^*$ , where  $e = \text{methodSea}$

## 3. Apply rule for $e^*$ : $\text{NEXT}(\text{methodSea}) = \{\text{methodSea}\} \cup \text{NEXT}(\text{methodSea}^*)$

(a) Call:  $\text{NEXT}(\text{methodSea}^*)$

(b) Initialize:  $\text{NEXT}(\text{methodSea}^*) = \{\}$ ,  $n = 1$

(c) Check stack:

$s_n = s_1 = \text{methodSea}^*$  and

$s_{n-1} = s_0 = e_1 e_2 e_3 e_4$ ,  $e_3 = \text{methodSea}^*$ ,  $e_4 = \text{'endclass'}$

(d) Apply the rule for sequence, where  $e_4 \neq 0$ :  $\text{NEXT}(\text{methodSea}^*) = \{\text{'endclass'}\}$

(e) Return:  $\text{NEXT}(\text{methodSea}^*) = \{\text{'endclass'}\}$

## 4. Return: $\text{NEXT}(\text{methodSea}) = \{\text{methodSea 'endclass'}\}$

### B.3. PEG example

Let us go through the grammar  $S \leftarrow \sim a \sim \sim b \sim$  using “..a..b..” as an input. As we see in Fig. B2, the stack is initialized with  $(S, 0)$  and the whole result is “..a..b..”, because it is a result of nonterminal expansion  $S \leftarrow \sim a \sim \sim b \sim$ . The sequence on the top is straightforward, as  $\sim a \sim$  consumes “..a..” and  $\sim b \sim$  consumes “b..”, and the result is then “..a..b..” (see Fig. B3).

In order to get result of  $\sim a \sim$  invoked in position 0, we first follow Definition 2 (see Fig. B4). It is a sequence of three parsers (generalization from the sequence of two to the sequence of three is straightforward). In Fig. B5 we see that before-water consumes “..”, the island itself consumes the desired “a” and another “..” is consumed by after-water.

Let us investigate what happens in before-water of  $\sim a \sim$ . First of all, we need to determine the  $\text{NEXT}(\sim a \sim)$ . In this case it is  $\sim b \sim$  (see B.2 for more complex example). Once we know the boundary, before-water tries to find the island  $a$  or its boundary  $\sim b \sim$  at positions 0 and 1 until it finds the island at the position 2 (see Fig. B6). We return a substring of all the positions for which we failed, i.e., “..”

**Overlapping seas:** The interesting question is, why does  $\sim b \sim$  fail in position 0? We already explained the problem with overlapping seas in Section 3.3, and now we show the computation formally. First of all, we rewrite the sea on top of the stack according to Definition 2. The new sequence on top of the stack fails because before-water returns  $\epsilon$  and there is no  $b$  at position 0 (see Fig. B7).

The before-water of  $\sim b \sim$  returns  $\epsilon$ , because of the overlapping seas case. It analyzes the stack and notices the before-water of  $\sim a \sim$  invoked on the position 0 (using the `seasOverlap` function) and returns  $\epsilon$  (see Fig. B8).

If there is no case of overlapping seas in the grammar, the before-water of  $\sim b \sim$  consumes “..a..” contrary to the correct parse  $\epsilon$  (see Fig. B8). This means that the before-water of  $\sim a \sim$  (see Fig. B6) would be  $x' = \epsilon$ . This would then fail the whole  $\sim a \sim$  and consequently the whole  $\sim a \sim \sim b \sim$ .

## References

- [1] Moonen L. Generating robust parser using island grammars. In: Burd E, Aiken P, Koschke R. editors. In: Proceedings of eighth working conference on reverse engineering (WCRE 2001). IEEE Computer Society; Wiley; 2001. p. 13–22. <http://dx.doi.org/10.1109/WCRE.2001.957806>. URL: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.13.1027>.

- [2] Kurš J, Lungu M, Nierstrasz O. Bounded seas: island parsing without shipwrecks. In: Combemale B, Pearce D, Barais O, Vinju J, editors. Software language engineering. Lecture notes in computer science, vol. 8706. Cham: Springer International Publishing; 2014. p. 62–81. [http://dx.doi.org/10.1007/978-3-319-11245-9\\_4](http://dx.doi.org/10.1007/978-3-319-11245-9_4) URL: (<http://scg.unibe.ch/archive/papers/Kurs14b-BoundedSeas.pdf>).
- [3] Renggli L, Ducasse S, Gîrba T, Nierstrasz O. Practical dynamic grammars for dynamic languages. In: 4th workshop on dynamic languages and applications (DYLA 2010), Malaga, Spain, 2010. p. 1–4. URL: (<http://scg.unibe.ch/archive/papers/Reng10cDynamicGrammars.pdf>).
- [4] Ford B. Parsing expression grammars: a recognition-based syntactic foundation. In: POPL '04: Proceedings of the 31st ACM SIGPLAN-SIGACT symposium on principles of programming languages. New York, NY, USA: ACM; 2004. p. 111–22. <http://dx.doi.org/10.1145/964001.964011>. URL: (<http://pdos.csail.mit.edu/~baford/packrat/pop104/peg-pop104.pdf>).
- [5] Hutton G, Meijer E. Monadic parser combinators. Technical report NOTTCS-TR-96-4. Department of Computer Science, University of Nottingham; 1996. URL: ([citeseer.ist.psu.edu/hutton96monadic.html](http://citeseer.ist.psu.edu/hutton96monadic.html)) (<http://eprints.nottingham.ac.uk/237/1/monparsing.pdf>).
- [6] Klint P, Visser E. Using filters for the disambiguation of context-free grammars. In: Proceedings of ASMICS workshop on parsing theory, 1994. p. 1–20.
- [7] van den Brand M, Scheerder J, Vinju JJ, Visser E. Disambiguation filters for scannerless generalized LR parsers. In: Horspool N, editor. Compiler construction (CC'02). Lecture notes in computer science, vol. 2304. Grenoble, France: Springer-Verlag; 2002. p. 143–58 URL: (<http://www.cs.uu.nl/people/visser/ftp/BSVV02.pdf>).
- [8] Dean TR, Cordy JR, Malton AJ, Schneider KA. Agile parsing in TXL. Autom Softw Eng 2003;10(4):311–36 URL: ([http://research.cs.queensu.ca/~cordy/Papers/JASE\\_AP.pdf](http://research.cs.queensu.ca/~cordy/Papers/JASE_AP.pdf)).
- [9] Zaytsev V. Formal foundations for semi-parsing. In: Software maintenance, reengineering and reverse engineering (CSMR-WCRE). IEEE conference on 2014 software evolution week, 2014. p. 313–17. <http://dx.doi.org/10.1109/CSMR-WCRE.2014.6747184>. URL: (<http://grammarware.net/text/2014/semiparsing.pdf>).
- [10] Frost R, Launchbury J. Constructing natural language interpreters in a lazy functional language. Comput J 1989;32(2):108–21. <http://dx.doi.org/10.1093/comjnl/32.2.108> URL: (<https://courses.cit.cornell.edu/ling4424/frost-launchbury.pdf>).
- [11] Kurš J, Larcheveque G, Renggli L, Berge A, Cassou D, Ducasse S, et al. PetitParser: building modular parsers. In: Deep into pharo, square bracket associates, 2013. p. 36. URL: (<http://scg.unibe.ch/archive/papers/Kurs13a-PetitParser.pdf>).
- [12] Visser E. Scannerless generalized-LR parsing. Technical report P9707. Programming Research Group, University of Amsterdam; July 1997. URL: (<http://www.cs.uu.nl/people/visser/ftp/P9707.ps.gz>).
- [13] Ford B. Packrat parsing: simple, powerful, lazy, linear time, functional pearl. In: ICFP 02: proceedings of the seventh ACM SIGPLAN international conference on functional programming, vol. 37/9. New York, NY, USA: ACM; 2002. p. 36–47. <http://dx.doi.org/10.1145/583852.581483>. URL: (<http://pdos.csail.mit.edu/~baford/packrat/icfp02/packrat-icfp02.pdf>).
- [14] Ford B. Packrat parsing: a practical linear-time algorithm with backtracking (Master's thesis). Massachusetts Institute of Technology; 2002. URL: (<http://pdos.csail.mit.edu/~baford/packrat/thesis/>) (<http://pdos.csail.mit.edu/~baford/packrat/thesis/thesis.pdf>).
- [15] Chomsky N. Three models for the description of language. IRE Trans Inf Theory 1956;2:113–24 URL: (<http://www.chomsky.info/articles/195609-.pdf>).
- [16] Scott E, Johnstone A. GLL parsing. Electron Notes Theor Comput Sci 2010;253(7):177–89. <http://dx.doi.org/10.1016/j.entcs.2010.08.041>.
- [17] Grune D, Jacobs CJ. Parsing techniques — a practical guide. New York: Springer; 2008 URL: (<http://www.cs.vu.nl/~dick/PT2Ed.html>).
- [18] Klusener S, Lämmel R. Deriving tolerant grammars from a base-line grammar. In: Proceedings of the international conference on software maintenance (ICSM 2003). Wiley; IEEE Computer Society; 2003. p. 179–88. <http://dx.doi.org/10.1109/ICSM.2003.1235420>.
- [19] Nierstrasz O, Ducasse S, Gîrba T. The story of Moose: an agile reengineering environment. In: Proceedings of the European software engineering conference (ESEC/FSE'05). New York, NY, USA: ACM Press; 2005. p. 1–10, invited paper. <http://dx.doi.org/10.1145/1095430.1081707>. URL: (<http://scg.unibe.ch/archive/papers/Nier05cStoryOfMoose.pdf>).
- [20] Hindle A, Godfrey MW, Holt RC. Reading beside the lines: indentation as a proxy for complexity metrics. In: ICPC '08: Proceedings of the 2008 IEEE international conference on program comprehension. Washington, DC, USA: IEEE Computer Society; 2008. p. 133–42. <http://dx.doi.org/10.1109/ICPC.2008.13>. URL: (<http://swag.uwaterloo.ca/~ahindle/pubs/hindle08icpc.pdf>).
- [21] Landin P. The next 700 programming languages. Commun ACM 1966;9(3):157–66. <http://dx.doi.org/10.1145/365230.365257> URL: (<http://www.cs.utah.edu/~eide/compilers/old/papers/p157-landin.pdf>).
- [22] Givi AS. Layout sensitive parsing in the PetitParser framework (Bachelor's thesis). University of Bern; October 2013. URL: (<http://scg.unibe.ch/archive/projects/Sade13a.pdf>).
- [23] Lavie A, Tomita M. GLR\* — an efficient noise-skipping parsing algorithm for context free grammars. In: Proceedings of the third international workshop on parsing technologies, 1993. p. 123–34.
- [24] Tomita M. Efficient parsing for natural language: a fast algorithm for practical systems. Norwell, MA, USA: Kluwer Academic Publishers; 1985.
- [25] Bischofberger WR. Sniff: a pragmatic approach to a C++ programming environment. In: C++ conference, 1992. p. 67–82. URL: (<http://citeseer.nj.nec.com/bischofberger92sniff.html>).
- [26] Asveld P. A fuzzy approach to erroneous inputs in context-free language recognition. In: Proceedings of the fourth international workshop on parsing technologies IWPT'95. Prague, Czech Republic: Institute of Formal and Applied Linguistics. Charles University; 1995. p. 14–25. URL: (<http://doc.utwente.nl/64694/>).
- [27] Koppler R. A systematic approach to fuzzy parsing. Softw: Pract Exp 1997;27(6):637–49. [http://dx.doi.org/10.1002/\(SICI\)1097-024X\(199706\)27:6<637::AID-SPE99>3.0.CO;2-3](http://dx.doi.org/10.1002/(SICI)1097-024X(199706)27:6<637::AID-SPE99>3.0.CO;2-3) URL: (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.4.3198&rep=rep1&type=pdf>).
- [28] Nilsson-Nyman E, Ekman T, Hedin G. Practical scope recovery using bridge parsing. In: Gašević D, Lämmel R, Van Wyk E, editors. Software language engineering. Lecture notes in computer science, vol. 5452. Berlin, Heidelberg: Springer; 2009. p. 95–113. [http://dx.doi.org/10.1007/978-3-642-00434-6\\_7](http://dx.doi.org/10.1007/978-3-642-00434-6_7).
- [29] Kats LCL, de Jonge M, Nilsson-Nyman E, Visser E. Providing rapid feedback in generated modular language environments. Adding error recovery to scannerless generalized-LR parsing. In: Leavens GT, editor. Proceedings of the 24th ACM SIGPLAN conference on object-oriented programming, systems, languages, and applications (OOPSLA 2009). ACM SIGPLAN notices. New York, NY, USA: ACM Press; 2009.
- [30] de Jonge M, Kats LCL, Soderberg E, Visser E. Natural and flexible error recovery for generated modular language environments. ACM Trans Programm Lang Syst 2012; 34 (4), article no. 15, 50 p. <http://dx.doi.org/10.1145/2400676.2400678>.