

# Unifying Parsing and Reflective Printing for Fully Disambiguated Grammars

Zirun Zhu · Hsiang-Shang Ko ·  
Yongzhe Zhang · Pedro Martins ·  
João Saraiva · Zhenjiang Hu

Received: date / Accepted: date

**Abstract** Language designers usually need to implement parsers and printers. Despite being two closely related programs, in practice they are often designed separately, and then need to be revised and kept consistent as the language evolves. It will be more convenient if the parser and printer can be unified and developed in a single program, with their consistency guaranteed automatically. Furthermore, in certain scenarios (like showing compiler optimisation results to the programmer), it is desirable to have a more powerful *reflective* printer that, when an abstract syntax tree corresponding to a piece of program text is modified, can propagate the modification to the program text while preserving layouts, comments, and syntactic sugar.

To address these needs, we propose a domain-specific language BIYACC, whose programs denote both a parser and a reflective printer for a fully disambiguated context-free grammar. BIYACC is based on the theory of bidirectional transformations, which helps to guarantee by construction that the generated pairs of parsers and reflective printers are consistent. Handling grammatical ambiguity is particularly challenging: We propose an approach based on generalised parsing and disambiguation filters, which produce all the parse results and (try to) select the only correct one in the

---

Zirun Zhu · Yongzhe Zhang (E-mail: {zhu, yz915}@nii.ac.jp)

National Institute of Informatics

Department of Informatics, The Graduate University for Advanced Studies, SOKENDAI

Hsiang-Shang Ko (E-mail: joshko@iis.sinica.edu.tw)

Institute of Information Science, Academia Sinica

Pedro Martins (E-mail: pribeiro@uci.edu)

University of California, Irvine

João Saraiva (E-mail: jas@di.uminho.pt)

HASLab/INESC TEC

Department of Informatics, University of Minho

Zhenjiang Hu (E-mail: huzj@pku.edu.cn)

Department of Computer Science and Technology, Peking University

Key Laboratory of High Confidence Software Technology, Ministry of Education of China

National Institute of Informatics

parsing direction; the filters are carefully bidirectionalised so that they also work in the printing direction and do not break the consistency between the parsers and reflective printers. We show that BIYACC is capable of facilitating many tasks such as Pombrio and Krishnamurthi’s ‘resugaring’, simple refactoring, and language evolution.

**Keywords** Asymmetric Lenses · Disambiguation Filters · Bidirectional Transformations · Domain-Specific Languages · Parsing · Reflective Printing

## 1 Introduction

Whenever we come up with a new programming language, as the front-end part of the system we need to design and implement a parser and a printer to convert between program text and an internal representation. A piece of program text, while conforming to a *concrete syntax* specification, is a flat string that can be easily edited by the programmer. The parser extracts the tree structure from such a string to a concrete syntax tree (CST), and converts it to an *abstract syntax* tree (AST), which is a more structured and simplified representation and is easier for the back-end to manipulate. On the other hand, a printer converts an AST back to a piece of program text, which can be understood by the user of the system; this is useful for debugging the system, or reporting internal information to the user.

Parsers and printers do conversions in opposite directions and are closely related—for example, the program text printed from an AST should be parsed to the same tree. It is certainly far from being economical to write parsers and printers separately: The parser and printer need to be revised from time to time as the language evolves, and each time we must revise the parser and printer and also keep them consistent with each other, which is a time-consuming and error-prone task. In response to this problem, many domain-specific languages [6, 7, 13, 37, 44, 53] have been proposed, in which the user can describe both a parser and a printer in a single program.

Despite their advantages, these domain-specific languages cannot deal with synchronisation between program text and ASTs. Let us look at a concrete example in Figure 1: The original program text is an arithmetic expression, containing a negation, a comment, and parentheses (one pair of which is redundant). It is first parsed to an AST (supposing that addition is left-associative) where the negation is desugared to a subtraction, parentheses are implicitly represented by the tree structure, and the comment is thrown away. Suppose that the AST is optimised by replacing `Add (Num 1) (Num 1)` with a constant `Num 2`. The user may want to observe the optimisation made by the compiler, but the AST is an internal representation not exposed to the user, so a natural idea is to propagate the changes on the AST back to the program text to make it easy for the user to check where the changes are. With a conventional printer, however, the printed result will likely mislead the programmer into thinking that the negation is replaced by a subtraction by the compiler; also, since the comment is not preserved, it will be harder for the programmer to compare the updated and original versions of the text. The problem illustrated here has also been investigated in many other practical scenarios where the parser and printer are used as a bridge between the system and the user, for example,

Original program text:	Printed result from a <i>conventional</i> printer:
-a /* a is the variable denoting... */	(0 - a) * (2 + a)
* (1 + 1 + (a))	Printed result from our <i>reflective</i> printer:
Abstract syntax tree:	-a /* a is the variable denoting... */
Mul (Sub (Num 0) (Var "a"))	* (2 + (a))
(Add (Add (Num 1) (Num 1)) (Var "a"))	
Optimised abstract syntax tree:	
Mul (Sub (Num 0) (Var "a"))	
(Add (Num 2) (Var "a"))	

**Fig. 1** Comparison between conventional printing and reflective printing.

- in bug reporting [51], where a piece of program text is parsed to its AST to be checked but error messages should be displayed for the program text;
- in code refactoring [18], where instead of directly modifying a piece of program text, most refactoring tools will first parse the program text into its AST, perform code refactoring on the AST, and regenerate new program text; and
- in language-based editors, as introduced by Reps [45, 46], where the user needs to interact with different printed representations of the same underlying AST.

To address the problem, we propose a domain-specific language BIYACC, which enables the user to describe both a parser and a *reflective* printer for a fully disambiguated context-free grammar (CFG) in a single program. Different from a conventional printer, a reflective printer takes a piece of program text and an AST, which is usually slightly modified from the AST corresponding to the original program text, and propagates the modification back to the program text. Meanwhile the comments (and layouts) in the unmodified parts of the program text are all preserved. This can be seen clearly from the result of using our reflective printer on the above arithmetic expression example in Figure 1. It is worth noting that reflective printing is a generalisation of the conventional notion of printing, because a reflective printer can accept an AST and an *empty* piece of program text, in which case it will behave just like a conventional printer, producing a new piece of program text depending on the AST only.

From a BIYACC program we can generate a parser and a reflective printer; in addition, we want to guarantee that the two generated components are *consistent* with each other. Specifically, given a pair of parser *parse* and reflective printer *print*, we want to ensure two (inverse-like) consistency properties: Firstly, a piece of program text *s* printed from an abstract syntax tree *t* should be parsed to the same tree *t*, i.e.<sup>1</sup>

$$\text{parse}(\text{print } s \ t) = t. \quad (1)$$

Secondly, updating a piece of program text *s* with an AST parsed from *s* should leave *s* unmodified (including formatting details like parentheses and whitespaces), i.e.

$$\text{print } s (\text{parse } s) = s. \quad (2)$$

<sup>1</sup> We assume basic knowledge about functional programming languages and their notations, in particular HASKELL [5, 34]. In HASKELL, an argument of function application does not need to be enclosed in (round) parentheses, i.e. we write  $f \ x$  instead of  $f(x)$ ; type variables are implicitly universally quantified, i.e.  $f :: a \rightarrow b \rightarrow a$  is the same as  $f :: \forall a \ b. a \rightarrow b \rightarrow a$  where  $::$  means *has type*. Additionally, we omit universal quantification for free variables in an equation; for instance,  $\text{parse}(\text{print } s \ t) = t$  is in fact  $\forall s \ t. \text{parse}(\text{print } s \ t) = t$ .

These two properties are inspired by the theory of *bidirectional transformations* [19], in particular *lenses* [17], and are guaranteed by construction for all BiYACC programs.

An online tool that implements the approach described in the paper can be accessed at <http://www.prg.nii.ac.jp/project/biyacc.html>. The webpage also contains the test cases used in the paper. The structure of the paper is as follows: We start with an overview of BiYACC in Section 2, explaining how to describe in a single program both a parser and a reflective printer for synchronising program text and its abstract syntax representation. After reviewing some background on bidirectional transformations in Section 3, in particular the bidirectional programming language BiGUL [22, 27, 28], we first give the semantics of a basic version of BiYACC that handles unambiguous grammars by compiling it to BiGUL in Section 4, guaranteeing the properties (1) and (2) by construction. Then, inspired by the research on *generalised parsing* [50] and *disambiguation filters* [26], in Section 5 we revise the basic BiYACC architecture to allow the use of ambiguous grammars and disambiguation directives while still retaining the above-mentioned properties. We present a case study in Section 6, showing that BiYACC is capable of describing TIGER [4], which shares many similarities with fully-fledged languages. We demonstrate that BiYACC can handle syntactic sugar, partially subsume Pombrio and Krishnamurthi’s ‘resugaring’ [42, 43], and facilitate language evolution. In Section 7, we present detailed related work including comparison with other systems. Contributions are summarised in Section 8.

This is the extended version of our previous work *Parsing and Reflective Printing, Bidirectionally* presented at SLE’16 [55], and the differences are mainly as follows: (i) We propose the notion of bidirectionalised filters and integrate them into BiYACC for handling grammatical ambiguity (Section 5); the related work section is also updated accordingly. (ii) We restructure the narration for introducing the basic BiYACC system and in particular elaborate on the isomorphism between program text and CSTs. (iii) We present the definitions and theorems in a more formal way, and complete their proofs. (iv) We make several other revisions such as renewing the figures for introducing the BiYACC system and the syntax of BiYACC programs.

Throughout this paper, we typeset general definitions and properties in *math style* and specific examples in *code style*.

## 2 A First Look at BiYACC

We first give an overview of BiYACC by going through the BiYACC program shown in Figure 2, which deals with the arithmetic expression example given in Section 1. This program consists of definitions of the abstract syntax, concrete syntax, directives, and actions for reflectively printing ASTs to CSTs; we will introduce them in order.

### 2.1 Syntax Definitions

*Abstract Syntax.* The abstract syntax part, which starts with the keyword `#Abstract`, is just one or more definitions of HASKELL data types. In our example, the abstract syntax

```

1  #Abstract
2  data Arith = Num Int
3      | Var String
4      | Add Arith Arith
5      | Sub Arith Arith
6      | Mul Arith Arith
7      | Div Arith Arith
9  #Concrete
10 Expr -> Expr '+' Term
11      | Expr '-' Term
12      | Term ;
13
14 Term -> Term '*' Factor
15      | Term '/' Factor
16      | Factor ;
17
18 Factor -> '-' Factor
19         | Numeric
20         | Identifier
21         | '(' Expr ')' ;
22
23 #Directives
24 LineComment: "//" ;
25 BlockComment: "/*" "*/" ;
26
27 #Actions
28 Arith +> Expr
29   Add x y +> [x +> Expr] '+' [y +> Term];
30   Sub x y +> [x +> Expr] '-' [y +> Term];
31   e +> [e +> Term];
32
33 Arith +> Term
34   Mul x y +> [x +> Term] '*' [y +> Factor];
35   Div x y +> [x +> Term] '/' [y +> Factor];
36   e +> [e +> Factor];
37
38 Arith +> Factor
39   Sub (Num 0) y +> '-' [y +> Factor];
40   Num i +> [i +> Numeric];
41   Var n +> [n +> Identifier];
42   e +> '(' [e +> Expr] ')';
43   ;

```

**Fig. 2** A BiYACC program for the expression example.

is defined in lines 2–7 by a single data type `Arith` whose elements are constructed from constants and arithmetic operators. Different constructors—namely `Num`, `Var`, `Add`, `Sub`, `Mul`, and `Div`—are used to construct different kinds of expressions.

*Concrete Syntax.* The concrete syntax part, beginning with the keyword `#Concrete`, is defined by a context-free grammar. For our expression example, in lines 10–21 we use a standard unambiguous grammatical structure to encode operator precedence and order of association, involving three nonterminal symbols `Expr`, `Term`, and `Factor`: An `Expr` can produce a left-sided tree of `Terms`, each of which can in turn produce a left-sided tree of `Factors`. To produce right-sided trees or operators of lower precedence under those with higher precedence, the only way is to reach for the last production rule `Factor -> '(' Expr ')'`, resulting in parentheses in the produced program text. There are also predefined nonterminals *Numeric* and *Identifier*, which produce numerals and identifiers respectively.

*Directives.* The `#Directives` part defines the syntax of comments and disambiguation directives. For example, line 23 shows that the syntax for single line comments is `"//"`<sup>2</sup>, while line 24 states that `"/*"` and `"*/"` are respectively the beginning mark and ending mark for block comments. Since the grammar for arithmetic expressions is unambiguous, there is no need to give any disambiguation directive for this example (whereas the ambiguous version of the grammar in Figure 6 needs to be augmented with a few such directives).

<sup>2</sup> While single quotation marks are for characters, double quotation marks are for strings. For simplicity, the user can always use double quotation marks.

## 2.2 Printing Actions

The main part of a BIYACC program starts with the keyword `#Actions` and describes how to update a CST with an AST. For our expression example, the actions are defined in lines 27–42 in Figure 2. Before explaining the actions, we should first say that program text is identified with CSTs when programming BIYACC actions: Conceptually, whenever we write a piece of program text, we are actually describing a CST rather than just a sequence of characters. We will expound on this identification of program text with CSTs in Section 4.2 in detail.

The `#Actions` part consists of groups of actions, and each group begins with a ‘type declaration’ of the form *HsType* ‘+>’ *Nonterminal* stating that the actions in this group specify updates on CSTs generated from *Nonterminal* using ASTs of type *HsType*. Informally, given an AST and a CST, the semantics of an action is to perform pattern matching simultaneously on both trees, and then use components of the AST to update corresponding parts of the CST, possibly recursively. (The syntax ‘+>’ suggests that information from the left-hand side is embedded into the right-hand side.) Usually the nonterminals in a right-hand side pattern are overlaid with update instructions, which are also denoted by ‘+>’.

Let us look at a specific action—the first one for the expression example, at line 28 of Figure 2:

```
Add x y +> [x +> Expr] '+' [y +> Term];
```

The AST-side pattern `Add x y` is just a HASKELL pattern; as for the CST-side pattern, the main intention is to refer to the production rule `Expr -> Expr '+' Term` and use it to match those CSTs produced by this rule—since the action belongs to the group `Arith +> Expr`, the part ‘`Expr ->`’ of the production rule can be inferred and thus is not included in the CST-side pattern. Finally we overlay ‘`x +>`’ and ‘`y +>`’ on the nonterminal symbols `Expr` and `Term` to indicate that, after the simultaneous pattern matching succeeds, the subtrees `x` and `y` of the AST are respectively used to update the left and right subtrees of the CST.

Having explained what an action means, we can now explain the semantics of the entire program. Given an AST and a CST as input, first a group (of actions) is chosen according to the types of the trees. Then the actions in the group are tried in order, from top to bottom, by performing simultaneous pattern matching on both trees. If pattern matching for an action succeeds, the updating operations specified by the action is executed; otherwise the next action is tried. Execution of the program ends when the matched action specifies either no updating operations or only updates to primitive data types such as `Numeric`. BIYACC’s most interesting behaviour shows up when all actions in the chosen group fail to match—in this case a suitable CST will be created. The specific approach adopted by BIYACC is to perform pattern matching on the AST only and choose the first matched action. A suitable CST conforming to the CST-side pattern is then created, and after that the whole group of actions is tried again. This time the pattern matching will succeed at the action used to create the CST, and the program will be able to make further progress. For instance, assuming that the source is `1 * 2` while the view is `Add (Num 1) (Num 2)`, a new source skeleton representing `_ + _` will be created and the `_` part will be updated recursively later. We will elaborate more on this in Section 4.

*Deep Patterns.* By using deep patterns, we can write actions that establish nontrivial relationships between CSTs and ASTs. For example, the action at line 38 of Figure 2 associates abstract subtraction expressions whose left operand is zero with concrete negated expressions; this action is the key to preserving negated expressions in the CST. For an example of a more complex CST-side pattern: Suppose that we want to write a pattern that matches those CSTs produced by the rule `Factor -> '-' Factor`, where the inner nonterminal `Factor` produces a further `'-' Factor` using the same rule. This pattern is written by overlaying the production rule on the first nonterminal `Factor` (an additional pair of parentheses is required for the expanded nonterminal): `'-' (Factor -> '-' Factor)`. More examples involving this kind of deep patterns can be found in Section 6.

*Layout and Comment Preservation.* The reflective printer generated by BIYACC is capable of preserving layouts and comments, but, perhaps mysteriously, in Figure 2 there is no clue as to how layouts and comments are preserved. This is because we decide to hide layout preservation from the user, so that the more important logic of abstract and concrete syntax synchronisation is not cluttered with layout preserving instructions. Our approach is fairly simplistic: We store layout information following each terminal in an additional field in the CST implicitly, and treat comments in the same way as layouts. During the printing stage, if the pattern matching on an action succeeds, the layouts and comments after the terminals shown in the right-hand side of that action are preserved; on the other hand, layouts and comments are dropped when a CST is created in the situation where pattern matching fails for all actions in a group. The layouts and comments before the first terminal are always kept during the printing.

*Parsing Semantics.* So far we have been describing the reflective printing semantics of the BIYACC program, but we may also work out its parsing semantics intuitively by interpreting the actions from right to left, converting the production rules to the corresponding constructors. (This might remind the reader of the usual YACC [23] actions.) In fact, this paper will not define the parsing semantics formally, because the parsing semantics is completely determined by the reflective printing semantics: If the actions are written with the intention of establishing some relation between the CSTs and ASTs, then BIYACC will be able to derive the only well-behaved parser, which respects that relation. We will explain how this is achieved in the next section.

### 3 Foundation of BIYACC: Putback-based Bidirectional Programming

From a BIYACC program, in addition to generating a parser and a printer, we also need to guarantee that the two generated programs are consistent with each other, i.e. satisfy the properties (1) and (2) stated in Section 1. It is possible to implement the *print* and *parse* semantics separately in an ad hoc way, but verifying the two consistency properties takes extra effort. The implementation we present, however, is systematic and guarantees consistency by construction, thanks to the well-developed theory of *bidirectional transformations* (BXs for short), in particular *lenses* [17]. We will give

a brief introduction to BXs below; for a comprehensive treatment, the readers are referred to the lecture notes for the 2016 Oxford Summer School on Bidirectional Transformations [19].

### 3.1 Parsing and Printing as Lenses

The *parse* and *print* semantics of BIYACC programs are potentially *partial*—for example, if the actions in a BIYACC program do not cover all possible forms of program text and abstract syntax trees, *parse* and *print* will fail for those uncovered inputs. Thus we should take partiality into account when choosing a BX framework in which to model *parse* and *print*. The framework we use in this paper is an explicitly partial version [32, 40] of asymmetric lenses [17].

**Definition 1 (Lenses)** A *lens* between a *source* type  $S$  and a *view* type  $V$  is a pair of functions

$$\begin{aligned} \text{get} &:: S \rightarrow \text{Maybe } V \\ \text{put} &:: S \rightarrow V \rightarrow \text{Maybe } S \end{aligned}$$

satisfying the *well-behavedness* laws:

$$\begin{aligned} \text{put } s \ v = \text{Just } s' &\Rightarrow \text{get } s' = \text{Just } v && (\text{PUTGET}) \\ \text{get } s = \text{Just } v &\Rightarrow \text{put } s \ v = \text{Just } s && (\text{GETPUT}) \end{aligned}$$

Intuitively, a *get* function extracts a part of a source of interest to the user as a view, and a *put* function takes a source and a view and produces an updated source incorporating information from the view. Partiality is explicitly represented by making the functions return Maybe-values: a *get* or *put* function returns  $\text{Just } r$  where  $r$  is the result, or  $\text{Nothing}$  if the input is not in the domain. The PUTGET law enforces that *put* must embed all information of the view into the updated source, so the view can be recovered from the source by *get*, while the GETPUT law prohibits *put* from performing unnecessary updates by requiring that putting back a view directly extracted from a source by *get* must produce the same, unmodified source.

The *parse* and *print* semantics of a BIYACC program will be the pair of functions *get* and *put* in a lens, required by definition to satisfy the two well-behavedness laws, which are exactly the consistency properties (1) and (2) reformulated in a partial setting:

**Definition 2 (The Partial Version of Consistency Properties)**

$$\begin{aligned} \text{print } s \ t = \text{Just } s' &\Rightarrow \text{parse } s' = \text{Just } t \\ \text{parse } s = \text{Just } t &\Rightarrow \text{print } s \ t = \text{Just } s \end{aligned}$$



### 3.2 Putback-based Bidirectional Programming in BiGUL

Having rephrased parsing and printing in terms of lenses, we can now construct consistent pairs of parsers and printers using *bidirectional programming* techniques, in which the programmer writes a single program to denote the two directions of a lens. Specifically, BiYACC programs are compiled to the *putback-based* bidirectional programming language BiGUL [28]. It has been formally verified in Agda [39] that BiGUL programs always denote well-behaved lenses, and BiGUL has been ported to HASKELL as an embedded DSL library [22]. BiGUL is putback-based, meaning that a BiGUL program describes a *put* function, but—since BiGUL is bidirectional—can also be executed as the corresponding *get* function. The advantage of putback-based bidirectional programming lies in the fact that, given a *put* function, there is at most one *get* function that forms a (well-behaved) lens with this *put* function [16]. That is, once we describe a *put* function as a BiGUL program, the *get* semantics of the program is completely determined by its *put* semantics. We can therefore focus solely on the printing (*put*) behaviour, leaving the parsing (*get*) behaviour only implicitly (but unambiguously) specified. How the programmer can effectively work with this paradigm has been more formally explained in terms of a Hoare-style logic for BiGUL [27].

Compilation of BiYACC to BiGUL (Section 4) uses only three BiGUL operations, which we briefly introduce here; more details can be found in the lecture notes on BiGUL programming [22]. A BiGUL program has type  $\text{BiGUL } s \ v$ , where  $s$  and  $v$  are respectively the source and view types.

*Replace.* The simplest BiGUL operation we use is

$$\text{Replace} :: \text{BiGUL } s \ s$$

which discards the original source and returns the view—which has the same type as the source—as the updated source. That is, the *put* semantics of *Replace* is the function  $\lambda \ s \ v \rightarrow \text{Just } v$ .

*Update.* The next operation *update* is more complex, and is implemented with the help of Template Haskell [49]. The general form of the operation is

$$\$(\text{update } [p] \ \text{spat} \ [] \ [p] \ \text{vpat} \ [] \ [d] \ bs \ []) :: \text{BiGUL } s \ v.$$

This operation decomposes the source and view by pattern matching with the patterns *spat* and *vpat* respectively, pairs the source and view components as specified by the patterns (see below), and performs further BiGUL operations listed in *bs* on the source–view pairs; the way to determine which source and view components are paired and which operation is performed on a pair is by looking for the same names in the three arguments. For example, the *update* operation

$$\$(\text{update } [p] \ (x, \_) \ [] \ [p] \ x \ [] \ [d] \ x = \text{Replace} \ [])$$

matches the source with a tuple pattern  $(x, \_)$  and the view with a variable pattern  $x$ , so that the first component of the source tuple is related with the whole view; during the

update, the first component of the source is replaced by the whole view, as indicated by the operation  $x = \text{Replace}$ . (The part marked by underscore ( $\_$ ) simply means that it will be skipped during the update.) Given a source  $(1, 2)$  and a view  $3$ , the operation will produce  $(3, 2)$  as the updated source. In general, any (type-correct) BiGUL program can be used in the list of further updates, not just the primitive  $\text{Replace}$ .

*Case.* The most complex operation we use is *Case* for doing case analysis on the source and view:

$$\text{Case} :: [\text{Branch } s \ v] \rightarrow \text{BiGUL } s \ v.$$

*Case* takes a list of branches, of which there are two kinds: *normal* branches and *adaptive* branches. For a normal branch, we should specify a main condition using a source pattern *spat* and a view pattern *vpat*, and an exit condition using a source pattern *spat'*:

$$\$(\text{normalSV } [p \mid \text{spat}] \ [p \mid \text{vpat}] \ [p \mid \text{spat'}]) :: \text{BiGUL } s \ v \rightarrow \text{Branch } s \ v.$$

An adaptive branch, on the other hand, only needs a main condition:

$$\$(\text{adaptiveSV } [p \mid \text{spat}] \ [p \mid \text{vpat}]) :: (s \rightarrow v \rightarrow s) \rightarrow \text{BiGUL } s \ v.$$

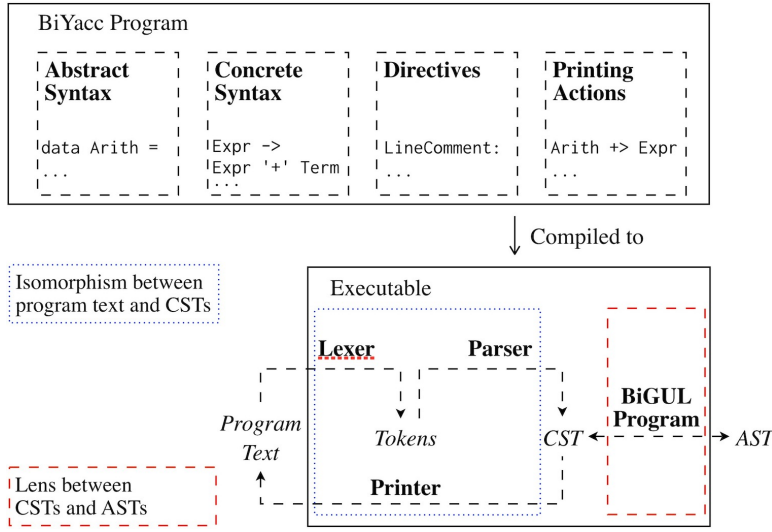
Their semantics in the *put* direction are as follows: A branch is applicable when the source and view respectively match *spat* and *vpat* in its main condition. Execution of a *Case* chooses the first applicable branch from the list of branches, and continues with that branch. When the applicable branch is a normal branch, the associated BiGUL operation is performed, and the updated source should satisfy the exit condition *spat'* (or otherwise execution fails); when the applicable branch is an adaptive branch, the associated function is applied to the source and view to compute an adapted source, and the whole *Case* is rerun on the adapted source and the view; it must go into a normal branch this time, otherwise the execution fails. Think of an adaptive branch as bringing a source that is too mismatched with the view to a suitable shape—for example, when the source is a subtraction while the view is an addition, which are by no means in correspondence, we must adapt the source to an addition—so that a normal branch that deals with sources and views in some sort of correspondence can take over. This adaptation mechanism is used by BiYACC to print an AST when the source program text is too different from the AST or even nonexistent at all.

#### 4 The Basic BiYACC

In this section we expound on a basic version of BiYACC that handles only unambiguous grammars. (Section 5 will present extensions for dealing with ambiguous grammars with disambiguation.) The architecture is illustrated in Figure 3, where a BiYACC program

$$\text{'\#Abstract' decls '\#Concrete' pgs '\#Directives' drctvs '\#Actions' ags}, \quad (3)$$

consisting of abstract syntax, concrete syntax, directives, and printing actions, as formally defined in Figure 4, is compiled into a few HASKELL source files and then



**Fig. 3** Architecture of BiYACC.

into an executable (by a HASKELL compiler) for converting between program text and ASTs. Specifically:

- The abstract syntax part (*decls* for HASKELL data type declarations) is already valid HASKELL code and is (almost) directly used as the definitions of AST data types.
- The concrete syntax part (*pgs* for production groups) is translated to definitions of CST data types (whose elements are representations of how a string is produced using the production rules), and also used to generate the pair of *concrete* parser (including a lexer) and printer for the conversion between program text and CSTs. This pair of concrete parser and printer can be shown to form a (partial) *isomorphism* (which will be defined in Section 4.1). This part will be explained in Section 4.2.
- The directives part (*drcvts* for directives) is used in the lexer for recognising single-line and multi-line comments.
- The printing actions part (*ags* for action groups) is translated to a BiGUL program (which is a lens, see Definition 1) for handling (the semantic part of) parsing and reflective printing between CSTs and ASTs. This part will be explained in Section 4.3.

The whole executable is a well-behaved lens since it is the composition of an isomorphism and a lens. We will start from a recap of this fact.

#### 4.1 Composition of Isomorphisms and Lenses

First we give the definition of (partial) isomorphisms.

**Definition 3 (Isomorphism)** A (partial) *isomorphism* between two types  $A$  and  $B$  is a pair of functions

$$\begin{aligned} to &:: A \rightarrow \text{Maybe } B \\ from &:: B \rightarrow \text{Maybe } A \end{aligned}$$

such that the *inverse properties* hold:

$$to\ a = \text{Just } b \iff from\ b = \text{Just } a.$$

**Definition 4 (Composition of Isomorphism and Lenses)** Given an isomorphism ( $to$  and  $from$ ) between  $A$  and  $B$  and a lens ( $get$  and  $put$ ) between  $B$  and  $C$ , we can compose them to form a new lens between  $A$  and  $C$ , whose components  $get'$  and  $put'$  are defined by

$$\begin{aligned} get' &:: A \rightarrow \text{Maybe } C \\ get'\ a &= to\ a \gg= get \\ put' &:: A \rightarrow C \rightarrow \text{Maybe } A \\ put'\ a\ c &= to\ a \gg= \lambda b \rightarrow put\ b\ c \gg= from \end{aligned}$$

where

$$\begin{aligned} (\gg=) &:: \text{Maybe } a \rightarrow (a \rightarrow \text{Maybe } b) \rightarrow \text{Maybe } b \\ \text{Just } x \gg= f &= f\ x \\ \text{Nothing} \gg= f &= \text{Nothing}. \end{aligned}$$

This is specialised from the standard definition of lens composition [17]—an isomorphism can be lifted to a lens (with  $get\ s = to\ s$  and  $put\ s\ v = from\ v$ ), which can then be composed with another lens to give rise to a new lens. We thus have the following lemma.

**Lemma 1** *Any lens resulted from the composition in Definition 4 is well-behaved.*

Therefore the whole BiYACC executable is a well-behaved lens, given that the concrete parser and printer form an isomorphism (Theorem 1) and the BiGUL program is a well-behaved lens (Theorem 2), which we will see next.

## 4.2 The Concrete Parsing and Printing Isomorphism

In this subsection, we describe the generation of CST data types and concrete printers (Section 4.2.1), the generation of concrete parsers (Section 4.2.2), and finally the inverse properties satisfied by the concrete parsers and printers (Section 4.2.3).

### 4.2.1 Generating CST Data Types and Concrete Printers

The production rules in a context-free grammar dictate how to produce strings from nonterminals, and a CST can be regarded as encoding one particular way of producing a string using the production rules. In BiYACC, we represent CSTs starting from a nonterminal  $nt$  as an automatically generated HASKELL data type named  $nt$ , whose constructors represent the production rules for  $nt$ . For each of these data types, we

```

Program ::= '#Abstract' HsDeclarations
         ['#Concrete' ProductionGroup+]
         '#Directives' CommentSyntaxDecl Disambiguation
         '#Actions' ActionGroup+
         ['#OtherFilters' OtherFilters]

ProductionGroup ::= Nonterminal '->' ProductionBody+ {'|'} ';'
ProductionBody ::= '[' Constructor ']' Symbol+ [{'#' Bracket '#'}]
Symbol ::= Primitive | Terminal | Nonterminal
Constructor ::= Nonterminal
CommentSyntaxDecl ::= 'LineComment:' String ';' 'BlockComment:' String ';'
Disambiguation ::= [Priority] [Associativity]
ActionGroup ::= HsType '+>' Nonterminal
               Action+ ';'
Action ::= HsPattern '+>' Update+ ';'
Update ::= Symbol
         | '[' HsVariable '+>' UpdateCondition ']'
         | '(' Nonterminal '->' Update+ ')'
UpdateCondition ::= Symbol
                 | '(' Nonterminal '->' UpdateCondition+ ')'

```

**Fig. 4** Syntax of BiYACC programs. (Nonterminals with prefix *Hs* denote HASKELL entities and follow the HASKELL syntax; the notation  $nt^+\{sep\}$  denotes a nonempty sequence of the same nonterminal  $nt$  separated by  $sep$ . Optional elements are enclosed in a pair of square brackets. The parts relating to disambiguation and filters will be explained in Section 5.)

also generate a printing function which takes a CST as input and produces a string as dictated by the production rules in the CST.

For instance, in Figure 2, the group of production rules from the nonterminal *Factor* (lines 18–21) is translated to the following HASKELL data type and concrete printing function:

```

data Factor = Factor1 String Factor
            | Factor2 (String, String)
            | Factor3 (String, String)
            | Factor4 String Expr String
            | FactorNull

cprtFactor :: Factor -> String
cprtFactor (Factor1 s1 factor1)    = "-" ++ s1 ++ cprtFactor factor1
cprtFactor (Factor2 (numeric, s1)) = numeric ++ s1
cprtFactor (Factor3 (identifier, s1)) = identifier ++ s1
cprtFactor (Factor4 s1 expr1 s2)   = "(" ++ s1 ++ cprtExpr expr1 ++ ")" ++ s2
cprtFactor FactorNull              = ""

```

where *Factor1* ... *Factor4* are constructors corresponding to the four production rules, and *FactorNull* represents an empty CST of type *Factor* and is used as the default value whenever we want to create new program text depending on the view only. As an example, *Factor1* represents the production rule *Factor*  $\rightarrow$  '-' *Factor*, and its *String* field stores the whitespaces appearing after a negation sign in the program text. The

Factor3 case makes a call to `cpriExpr :: Expr -> String`, which is the printing function generated for the nonterminal `Expr`.

Following this idea, we define the translation from production rule groups ( $pgs$  in formula (3)) to datatype definitions by source-to-source compilation rules:

$$\begin{aligned} \llbracket pgs \rrbracket_{ProductionGroup} &= \langle \llbracket pg \rrbracket_{ProductionGroup} \mid pg \in pgs \rangle \\ \llbracket nt \text{ '}' \rightarrow \text{' bodies'} \rrbracket_{ProductionGroup} &= \\ \text{'data' } nt \text{ '}' &= \langle \text{CON}(nt, body) \langle \text{FIELD}(s) \mid s \in body \rangle \text{'|'} \mid body \in bodies \rangle \\ &\quad \text{NULLCON}(nt) . \end{aligned}$$

Compilation rules of this kind will also be used later, so we introduce the notation here: Compilation rules are denoted by semantic brackets ( $\llbracket \cdot \rrbracket$ ), and refer to some auxiliary functions, whose names are in SMALL CAPS. A nonterminal in subscript gives the ‘type’ of the argument or metavariable before it. The angle bracket notation  $\langle f e \mid e \in es \rangle$  denotes the generation of a list of entities of the form  $f e$  for each element  $e$  in the list  $es$ , in the order of their appearance in  $es$ . The auxiliary function  $\text{CON}(nt, body)$  retrieves the constructor for a production rule. The fields of a constructor are generated from the right-hand side of the corresponding production rule in the way described by the auxiliary function `FIELD`—nonterminals that are not primitives are left unchanged (using their names for data types), primitives are stored in the `String` type<sup>3</sup>, terminal symbols are dropped, and an additional `String` field is added for each terminal and primitive for storing layout information (whitespaces and comments) appearing after the terminal or primitive in the program text. The last step is to insert an additional empty constructor, whose name is denoted by  $\text{NULLCON}(nt)$ .

#### 4.2.2 Generating Concrete Lexers and Parsers

The implementation of the concrete parser, which turns program text into CSTs, is further divided into two phases: lexing and parsing. In both phases, the layout information (whitespaces and comments) is automatically preserved, which makes the CSTs isomorphic to the program text.

*Lexer.* Apart from handling the terminal symbols appearing in a grammar, the lexer automatically derived by `BiYACC` can also recognise several kinds of literals, including integers, strings, and identifiers, respectively produced by the nonterminals `Numeric`, `String`, and `Identifier`. For now, the forms of these literals are predefined, but we take this as a step towards a lexerless grammar, in which strings produced by nonterminals can be specified in terms of regular expressions. Furthermore, whitespaces and comments are carefully handled in the derived lexer, so they can be completely stored in CSTs and correctly recovered to the program text in printing. This feature of `BiYACC`, which we explain below, makes layout preservation transparent to the programmer.

An assumption of `BiYACC` is that whitespaces are only regarded as separators between other tokens. (Although there exist some languages such as `HASKELL`

<sup>3</sup> The reason for storing primitives in the `String` type is because `String` is the most precise representation that will not cause the loss of any information. For instance, this is useful for retaining the leading zeros of an integer such as 073. Storing 073 as `Integer` will cause the loss of the leading zero.

and PYTHON where indentation does affect the meaning of a program, there are workarounds, e.g. writing a preprocessing program to insert explicit separators.) Usually, token separators are thrown away in the lexing phase, but since we want to keep layout information in CSTs, which are built by the parser, the lexer should leave the separators intact and pass them to the parser. The specific approach taken by BIYACC is wrapping a lexeme and the whitespaces following it into a single token. Beginning whitespaces are treated separately from lexing and parsing, and are always preserved. And in this prototype implementation, comments are also regarded as whitespaces.

*Parser.* The concrete parser is used to generate a CST from a list of tokens according to the production rules in the grammar. Our parser is built using the parser generator HAPPY [33], which takes a BNF specification of a grammar with semantic actions and produces a HASKELL module containing a parser function. The grammar we feed into HAPPY is still essentially the one specified in a BIYACC program, but in addition to parsing and constructing CSTs, the HAPPY actions also transfer the whitespaces wrapped in tokens to corresponding places in the CSTs. For example, the production rules for *Factor* in the expression example, as shown on the left below, are translated to the HAPPY specification on the right:

$\begin{array}{l} \text{Factor} \\ \rightarrow \text{'-'} \text{Factor} \\   \text{Numeric} \\   \text{Identifier} \\   \text{'(' Expr ')'}; \end{array}$	$\rightsquigarrow$	$\begin{array}{l} \text{Factor} \\ : \text{token1 Factor} \quad \{ \text{Factor1 \$1 \$2} \} \\   \text{tokenNumeric} \quad \{ \text{Factor2 \$1} \} \\   \text{tokenIdentifier} \quad \{ \text{Factor3 \$1} \} \\   \text{token2 Expr token3} \{ \text{Factor4 \$1 \$2 \$3} \} . \end{array}$
---	--------------------	--

We use the first expansion (*token1 Factor*) to explain how whitespaces are transferred: The generated HAPPY token *token1* matches a ‘-’ token produced by the lexer, and extracts the whitespaces wrapped in the ‘-’ token; these whitespaces are bound to *\$1*, which is placed into the first field of *Factor1* by the associated HASKELL action.

#### 4.2.3 Inverse Properties

Now we give the types of the concrete printer and parser generated from a BIYACC program and show that they form an isomorphism. Let the type CST be the set of all the CSTs defined by the grammar of a BIYACC program; by default it is the source type (nonterminal) of the first group of actions in the *#Actions* part. We have seen in Section 4.2.1 how to generate its datatype definition and a concrete printing function

$$cprint :: \text{CST} \rightarrow \text{String} .$$

On the other hand, from the grammar we directly use a parser generator to generate a concrete parsing function

$$cparse :: \text{String} \rightarrow \text{Maybe CST} ,$$

which is Maybe-valued since a piece of input text may be invalid. This *cparse* function is one direction of the isomorphism in the executable, while the other direction is

$$\text{Just} \circ cprint :: \text{CST} \rightarrow \text{Maybe String} .$$

Below we show that the inverse properties amount to the requirements that the generated parser is ‘correct’ and the grammar is unambiguous.

Since our concrete parsers are generated by the parser generator HAPPY [33], we need to assume that they satisfy some essential properties, for we cannot control the generation process and verify those properties.

**Definition 5 (Parser Correctness)** A parser  $cparse$  is *correct* with respect to a printer  $cprint$  exactly when

$$cparse\ text = Just\ cst \quad \Rightarrow \quad cprint\ cst = text \quad (4)$$

$$cprint\ cst = text \quad \Rightarrow \quad \exists\ cst'.\ cparse\ text = Just\ cst' . \quad (5)$$

To see what (4) means, recall that our CSTs, as described in Section 4.2.1, encode precisely the derivation trees, with the CST constructors representing the production rules used, and  $cprint$  traverses the CSTs and follows the encoded production rules to produce the derived program text. Now consider what  $cparse$  is supposed to do: It should take a piece of program text and find a derivation tree for it, i.e. a CST which  $cprints$  to that piece of program text. This statement is exactly (4). In other words, (4) is the functional specification of parsing, which is satisfied if the parser generator we use behaves correctly. Also it is reasonable to expect that a parser will be able to successfully parse any valid program text, and this is exactly (5).

We also need to make an assumption about concrete printers: recall that in this section we assume that the grammar is unambiguous, and this amounts to injectivity of  $cprint$ —for any piece of program text there is at most one CST that prints to it.

With these assumptions, we can now establish the isomorphism (which is rather straightforward).

**Theorem 1 (Inverse Properties)** *If a parser  $cparse$  is correct with respect to an injective printer  $cprint$ , then  $cparse$  and  $Just \circ cprint$  form an isomorphism, that is,*

$$cparse\ text = Just\ cst \quad \Leftrightarrow \quad (Just \circ cprint)\ cst = Just\ text .$$

*Proof* The left-to-right direction is immediate since the right-hand side is equivalent to  $cprint\ cst = text$ , and the whole implication is precisely (4). For the right-to-left direction, again the antecedent is equivalent to  $cprint\ cst = text$ , and we can invoke (5) to obtain  $cparse\ (text) = Just\ cst'$  for some  $cst'$ . This is already close to our goal—what remains to be shown is that  $cst'$  is exactly  $cst$ , which is indeed the case because

$$\begin{aligned} & cparse\ text = Just\ cst' \\ \Rightarrow & \{ \text{antecedent} \} \\ & cparse\ (cprint\ cst) = Just\ cst' \\ \Rightarrow & \{ (4) \} \\ & cprint\ cst' = cprint\ cst \\ \Rightarrow & \{ cprint\ \text{is injective} \} \\ & cst' = cst . \end{aligned}$$

□



```

[[vt '+>' st acts]]ActionGroup =
  PROG(vt, st) '::' 'BiGUL' st vt
  PROG(vt, st) '='
  'Case' '[' <[[a]]ActionN,vt,st ' ' | a ∈ acts> <[[a]]ActionA,st | a ∈ acts>{' ','}' ']'

[[vpat '+>' updates]]ActionN,vt,st =
  '$(normalSV'
  '[' SRCCOND(ERSVARS('[' st '->' updates '])')Update '[' ']' '[' vpat '[' ']'
  '[' SRCCOND(ERSVARS('[' st '->' updates '])')Update '[' '])'
  '$(update' '[' REMOVEAS(vpat) '[' ']'
  '[' SRCPAT('[' st '->' updates '])')Update '[' ']'
  '[' d '[' <[[u]]Updatevt,vpat | u ∈ updates> '[' '])'

[['[' var '+>' ucPrimitive ']]']Updatevt,vpat = var '=' Replace;
[['[' var '+>' ucNonterminal ']]']Updatevt,vpat = var '=' PROG(VARTYPE(vt, vpat, var), uc) ';';
[['[' var '+>' '(' nt '->' ... ')' ']]']Updatevt,vpat = [['[' var '+>' nt '[' ']]']Updatevt,vpat
[['(' ... '->' updates ')']Updatevt,vpat = <[[u]]Updatevt,vpat ' ' | u ∈ updates>
[[symbol]Updatevt,vpat = ''

[[vpat '+>' updates]]ActionA,st = '$(adaptiveSV' '[' _ '[' ']' '[' vpat '[' '])'
  '(\_ _ ->' DEFAULTEXPR(ERSVARS('[' st '->' updates ']))

FIELD(nt)Nonterminal = nt
FIELD(t)Terminal = 'String'
FIELD(p)Primitive = '(' p ' ', String)'

ERSVARS('[' var '+>' uc '])')Update = uc
ERSVARS('[' nt '->' updates ')')Update = '(' nt '->' <ERSVARS(u) | u ∈ updates> ')'
ERSVARS(symbol)Update = symbol

SRCCOND('[' nt '->' uconds ')')UpdateCondition = '(' CON(nt, <CONDHEAD(uc) | uc ∈ uconds>)
  <SRCCOND(uc) | uc ∈ uconds> ')'

SRCCOND(symbol)UpdateCondition = '_'

CONDHEAD('[' nt '->' ... ')')UpdateCondition = nt
CONDHEAD(symbol)UpdateCondition = symbol

SRCPAT('[' var '+>' ucPrimitive '])')Update = '(' var ' ', _)'
SRCPAT('[' var '+>' ucNonterminal '])')Update = var
SRCPAT('[' nt '->' updates ')')Update = '(' CON(nt, <CONDHEAD(uc) | uc ∈ ERSVARS(updates)>)
  <SRCPAT(u) | u ∈ updates> ')'

SRCPAT(symbol)Symbol = '_'

DEFAULTEXPR(symbol)Primitive = '(undefined, " ")'
DEFAULTEXPR(symbol)Nonterminal = NULLCON(symbol)
DEFAULTEXPR(symbol)Terminal = '" "'
DEFAULTEXPR('[' nt '->' uconds ')')UpdateCondition = CON(nt, <CONDHEAD(uc) | uc ∈ uconds>)
  <DEFAULTEXPR(uc) | uc ∈ uconds>

```

**Fig. 5** Semantics of BiYACC programs (as BiGUL programs).

### 4.3 Generating the BiGUL Lens

The source-to-source compilation from the actions part of a BiYACC program to a BiGUL program (i.e. lens) is shown in Figure 5. Additional arguments to the semantic bracket are typeset in superscript, and the notation  $\langle \dots \mid \dots \in \dots \rangle \{s\}$  means inserting  $s$  between the elements of the list.

*Action Groups.* Each group of actions is translated into a small BIGUL program, whose name is determined by the view type  $vt$  and source type  $st$  and denoted by  $\text{PROG}(vt, st)$ . The BIGUL program has one single `Case` statement, and each action is translated into two branches in this `Case` statement, one normal and the other adaptive. All the adaptive branches are gathered in the second half of the `Case` statement, so that the normal branches will be tried first. For example, the third group of type `Arith +> Factor` is compiled to

```
bigulArithFactor :: BiGUL Factor Arith
bigulArithFactor =
  Case [ ... -- normal branches
        ... -- adaptive branches
      ] .
```

*Normal Branches.* We said in Section 2 that the semantics of an action is to perform pattern matching on both the source and view, and then update parts of the source with parts of the view. This semantics is implemented with a normal branch: The source and view patterns are compiled to the main condition, and, together with the updates overlaid on the source pattern, also to an update operation. For example, the first action in the `Arith-Factor` group

```
Sub (Num 0) y +> '-' (y +> Factor)
```

is compiled to

```
$(normalSV [p| (Factor1 _ _) |] [p| Sub (Num 0) y |] [p| (Factor1 _ _) |])
$(update [p| Sub (Num 0) y |] [p| (Factor1 _ y) |] [d| y = bigulArithFactor; |]) .
```

When the CST is a `Factor1` and the AST matches `Sub (Num 0) y`, we enter this branch, decompose the source and view by pattern matching, and use the view's right subtree  $y$  to update the second field of the source while skipping the first field (which stores whitespaces); the name of the BIGUL program for performing the update is determined by the type of the smaller source  $y$  (deduced by `VARTYPE`) and that of the smaller view.

*Adaptive Branches.* When all actions in a group fail to match, we should adapt the source into a proper shape to correspond to the view. This is done by generating adaptive branches from the actions during compilation. For example, besides a normal branch, the first action in the `Arith-Factor` group `Sub (Num 0) y +> '-' (y +> Factor)` is also compiled to

```
$(adaptiveSV [p| _ |] [p| Sub (Num 0) _ |]) (\ _ _ -> Factor1 " " FactorNull) .
```

Since the source pattern of the main condition (of the adaptive branch) is a wildcard, the branch is always applicable if the view matches `Sub (Num 0) _`. The body of the adaptation function is generated by the auxiliary function `DEFAULTEXPR`, which creates a skeletal value—here `Factor1 " " FactorNull` represents a negation skeleton - whose value is not (resursively) created yet—that matches the source pattern. These adaptive branches are placed at the end of an action group and tried only if no normal branches are applicable so that unnecessary adaptation will never be performed.

*Entry Point.* The entry point of the program is chosen to be the BIGUL program compiled from the first group of actions. This corresponds to our assumption that the initial input concrete and abstract syntax trees are of the types specified for the first action group. (It is rather simple so the rules are not shown in the figure.) For the expression example, we generate a definition

$$\text{entrance} = \text{bigulArithExpr}$$

which is invoked in the `main` program.

*Well-behavedness.* Since BIGUL programs always denote well-behaved lenses, a fact which has been formally verified [39], we get the following theorem for free.

**Theorem 2 (Well-behavedness)** *The BIGUL program generated from a BIYACC program is a lens; that is, it satisfies the well-behavedness laws in Definition 1 with  $cst$  substituted for the source  $s$  and  $ast$  for the view  $v$ :*

$$\begin{aligned} \text{put } cst \text{ } ast = \text{Just } cst' &\Rightarrow \text{get } cst' = \text{Just } ast \\ \text{get } cst = \text{Just } ast &\Rightarrow \text{put } cst \text{ } ast = \text{Just } cst . \end{aligned}$$

## 5 Handling Grammatical Ambiguity

In Section 4, we have described the basic version of BIYACC, about which there is an important assumption (stated in Theorem 1) that grammars have to be unambiguous. Having this assumption can be rather inconvenient in practice, however, as ambiguous grammars (with disambiguation directives) are often preferred since they are considered more natural and human-friendly than their unambiguous versions [2, 26]. Therefore the purpose of this section is to revise the architecture of basic BIYACC to allow the use of ambiguous grammars and disambiguation directives. This is in fact a long-standing problem: tools designed for building parser and printer pairs usually do not support such functionality (Section 7.1).

For example, consider the ambiguous grammar (with disambiguation directives) and printing actions in Figure 6, which we will refer to throughout this section. Note that the parenthesis structure is dropped when converting a CST to its AST (as stated by the last printing action of `Arith +> Expr`). The grammar is converted to CST data types and constructors as in Section 4.2.1, but here we explicitly give names such as `Plus` and `Times` to production rules, and these names (instead of automatically generated ones) are used for constructors in CSTs. Compared with this grammar, the unambiguous one shown in Figure 2 is less intuitive as it uses different nonterminals to resolve the ambiguity regarding operator precedence and associativity.

In this section, we explain the problem brought by ambiguous grammars (Section 5.1) and address it (Section 5.2) using *generalised parsing* and *bidirectionalised filters* (bi-filters for short). Then we extend BIYACC with bi-filters (Section 5.3) while still retaining the well-behavedness. To program with bi-filters easily, we provide compositional bi-filter directives (Section 5.4) which compile to priority and associativity bi-filters. Power users can also define their own bi-filters (Section 5.5), and we illustrate this by writing a bi-filter that solves the (in)famous dangling-else problem.

```

#Concrete
Expr -> [Plus]      Expr '+' Expr
      | [Minus]     Expr '-' Expr
      | [Times]     Expr '*' Expr
      | [Division]  Expr '/' Expr
      | [Paren]     '(' Expr ')'
      | [Lit]       Numeric
      ;

#Directives
Priority:
Times > Plus    ;
Times > Minus   ;
Division > Plus ;
Division > Minus ;

Associativity:
Left: Plus, Minus, Times, Division ;

#Actions
Arith +> Expr
Add x y +> [x +> Expr] '+' [y +> Expr] ;
Sub x y +> [x +> Expr] '-' [y +> Expr] ;
Mul x y +> [x +> Expr] '*' [y +> Expr] ;
Div x y +> [x +> Expr] '/' [y +> Expr] ;
Num i   +> [i +> Numeric]          ;
e       +> '(' [e +> Expr] ')'      ;
;;

```

**Fig. 6** Arithmetic expressions defined by an ambiguous grammar and the corresponding printing actions. (For simplicity, the variable and negation productions are omitted.)

### 5.1 Problems with Ambiguous Grammars

Consider the original architecture of BiYACC in Figure 3, which we want to (and basically will) retain while adapting it to support ambiguous grammars. The first component (of the executable) we should adapt is  $cparse :: \text{String} \rightarrow \text{Maybe CST}$ , the (concrete) parsing direction of the isomorphism: since there can be multiple CSTs corresponding to the same program text,  $cparse$  needs to choose one of them as the result. Disambiguation directives [23] were invented to describe how to make this choice. For example, with respect to the grammar in Figure 6, text  $1 + 2 * 3$  will have either of the two CSTs<sup>4</sup>:

$$\begin{aligned} cst_1 &= \#Plus\ 1\ (Times\ 2\ 3) \\ cst_2 &= \#Times\ (Plus\ 1\ 2)\ 3 \end{aligned}$$

depending on the precedence of addition and multiplication. Conventionally, we can use the YACC-style disambiguation directives `%left '+'`; `%left '*'`; to specify that multiplication has higher precedence over addition, and instruct the parser to choose  $cst_1$ .

However, merely adapting  $cparse$  with disambiguation behaviour is not enough, since the isomorphism (Theorem 1), in particular its right to left direction (which is simplified as  $cparse\ (cprint\ cst) = Just\ cst$ ), cannot be established when an ambiguous grammar is used—in the example above,  $cparse\ (cprint\ cst_2) = Just\ cst_1 \neq Just\ cst_2$ . This is because the image of  $cparse$  is strictly smaller than the domain of  $cprint$ : if we start from any CST not in the image of  $cparse$ , we will never be able to get back to the same CST through  $cprint$  and then  $cparse$ . This tells us that, to retain the isomorphism, the domain of  $cprint$  should not be the whole CST but only the image of  $cparse$ ,

<sup>4</sup> For simplicity, we use  $\#$  to annotate type-incorrect CSTs in which fields for layouts (and comments) and unimportant constructors such as `Lit` are omitted.

i.e. the set of *valid* CSTs (as defined by the disambiguation directives), which we denote by  $\text{CST}_F$  (for reasons that will be made clear in Section 5.3).

Now that the right-hand side domain of the isomorphism is restricted to  $\text{CST}_F$ , the source of the lens should be restricted to this set as well. For  $\text{get} :: \text{CST} \rightarrow \text{Maybe AST}$  we need to restrict its domain, which is easy; for  $\text{put} :: \text{CST} \rightarrow \text{AST} \rightarrow \text{Maybe CST}$  we should revise its type to  $\text{CST}_F \rightarrow \text{AST} \rightarrow \text{Maybe CST}_F$ , meaning that *put* should now guarantee that the CSTs it produces are valid, which is nontrivial. For example, consider the result of  $\text{put } \text{cst } \text{ast}$  where  $\text{ast} = \text{Mul } (\text{Add } (\text{Num } 1) (\text{Num } 2)) (\text{Num } 3)$  and  $\text{cst}$  is some arbitrary tree. A natural choice is  $\text{cst}_2$ , which, however, is excluded from  $\text{CST}_F$  by disambiguation. A possible solution could be making *put* refuse to produce a result from  $\text{ast}$ , but this is unsatisfactory since  $\text{ast}$  is perfectly valid and should not be ignored by *put*. A more satisfactory way is creating a CST with *proper parentheses*, like  $\text{cst}_3 = \# \text{Times } (\text{Paren } (\text{Plus } 1\ 2))\ 3$ . But it is not clear in what cases parentheses need to be added, in what cases they need not, and in what cases they cannot.

We are now led to a fundamental problem: generally, *put* strategies for producing valid CSTs should be inferred from the disambiguation directives, but the semantics of YACC disambiguation directives are defined over the implementation of YACC's underlying LR parsing algorithm with a stack [3, 23], and therefore it is nontrivial to invent a dual semantics in the *put* direction. To have a simple and clear semantics of the disambiguation process, we turn away from YACC's traditional approach and opt for an alternative approach based on *generalised parsing* with *disambiguation filters* [9, 26], whose semantics can be specified implementation-independently. Based on this simple and clear semantics, we will be able to devise ways to amend *put* to produce only valid CSTs, and formally state the conditions under which the executable generated by the revised BiYACC is well-behaved.

## 5.2 Generalised Parsing and Bidirectionalised Filters

The idea of generalised parsing is for a parser to produce all possible CSTs corresponding to its input program text instead of choosing only one CST (possibly prematurely) [14, 47, 50, 54], and works naturally with ambiguous grammars. In practice, a generalised parser can be generated using, e.g., HAPPY's GLR mode [33], and we will assume that given a grammar we can obtain a generalised parser:

$$\text{cgparse} :: \text{String} \rightarrow [\text{CST}] .$$

The result of *cgparse* is a list of CSTs. We do not need to wrap the result type in *Maybe*—if *cgparse* fails, an empty list is returned. And we should note that, while the result is a list, what we really mean is a set (commonly represented as a list in HASKELL) since we do not care about the order of the output CSTs and do not allow duplicates.

With generalised parsing, program text is first parsed to all the possible CSTs; disambiguation then becomes an extremely simple concept: removing CSTs that the user does not want. One possible semantics of disambiguation may be a function  $\text{judge} :: \text{Tree} \rightarrow \text{Bool}$ ; during disambiguation, this function is applied to all candidate CSTs, and a candidate *cst* is removed if *judge cst* returns *False*, or kept otherwise. We

call these functions *disambiguation filters* (‘filters’ for short).<sup>5</sup> For example, to state that top-level addition is left-associative, we can use the following filter<sup>6</sup> to reject right-sided trees:

```
plusJudge :: Expr -> Bool
plusJudge (#Plus _ (Plus _ _)) = False
plusJudge _ = True .
```

This simple and clean semantics of disambiguation is then amenable to ‘bidirectionalisation’, which we do next.

Note that, unlike YACC’s disambiguation directives, which assign precedence and associativity to individual tokens and *implicitly* exclude ‘some’ CSTs, in `plusJudge` above we *explicitly* ban incorrect CSTs through pattern matching. Having described which CSTs are incorrect, we can further specify what to do with incorrect CSTs in the printing direction. Whenever a CST ‘in a bad shape’, i.e. rejected by a filter like `plusJudge`, is produced, we can repair it so that it becomes ‘in a good shape’:

```
plusRepair :: Expr -> Expr
plusRepair (#Plus t1 (Plus t2 t3)) = #Plus t1 (Paren (Plus t2 t3))
plusRepair t = t .
```

The above function states that whenever a `Plus` is another `Plus`’s right child, there must be a parenthesis structure `Paren` in between. Observant readers might have found that the trees processed by `plusJudge` and `plusRepair` have the same pattern. We can therefore pair the two functions and make a *bidirectionalised filter* (‘bi-filters’ for short):

```
plusLAssoc :: Expr -> (Expr, Bool)
plusLAssoc (#Plus t1 (Plus t2 t3)) = (#Plus t1 (Paren (Plus t2 t3)), False)
plusLAssoc t = (t, True) .
```

But there is still some redundancy in the definition of `plusLAssoc`, for when the input tree is correct we always return the same input tree; this can be further optimised:

```
plusLAssoc' :: Expr -> Maybe Expr
plusLAssoc' (#Plus t1 (Plus t2 t3)) = Just (#Plus t1 (Paren (Plus t2 t3)))
plusLAssoc' _ = Nothing .
```

Generalising the example above, we arrive at the definition of bi-filters.

**Definition 6 (Bidirectionalised Filters)** A *bidirectionalised filter*  $F$  working on trees of type  $t$  is a function of type  $\text{BiFilter } t$  defined by

$$\text{type BiFilter } t = t \rightarrow \text{Maybe } t$$

satisfying

$$\text{repair } F \, t = t' \quad \Rightarrow \quad \text{judge } F \, t' = \text{True} \quad (\text{RepairJudge})$$

<sup>5</sup> The general type for disambiguation filters is  $[t] \rightarrow [t]$ , which allows comparison among a list of CSTs. However, since in this paper we only consider *property filters* defined in terms of predicates (on a single tree), it is sufficient to use the simplified type  $t \rightarrow \text{Bool}$ . See Section 7.2.

<sup>6</sup> This is not a very realistic filter, although it sufficiently demonstrates the use of filters and removes ambiguity in simplest cases like  $1 + 2 * 3$ . In general, the filter should be *complete* (Definition 9) so that ambiguity is fully removed from the grammar.

where the two directions *repair* and *judge* are defined by

$$\begin{aligned} \text{repair} &:: \text{BiFilter } t \rightarrow (t \rightarrow t) \\ \text{repair } F \ t &= \text{case } F \ t \text{ of} \\ &\quad \text{Nothing} \rightarrow t \\ &\quad \text{Just } t' \rightarrow t' \\ \text{judge} &:: \text{BiFilter } t \rightarrow (t \rightarrow \text{Bool}) \\ \text{judge } F \ t &= \text{case } F \ t \text{ of} \\ &\quad \text{Nothing} \rightarrow \text{True} \\ &\quad \text{Just } \_ \rightarrow \text{False} . \end{aligned}$$

The functions *repair* and *judge* accept a bi-filter and return respectively the specialised *repair* and *judge* functions for that bi-filter. For clarity, we let  $\text{repair}_F$  denote *repair* *F* and let  $\text{judge}_F$  denote *judge* *F*. The bi-filter law *RepairJudge* dictates that  $\text{repair}_F$  should transform its input tree into a state accepted by  $\text{judge}_F$ . The reader may wonder why there is not a dual *JudgeRepair* law saying that if a tree is already of an allowed form justified by  $\text{judge}_F$ , then  $\text{repair}_F$  should leave it unchanged. In fact this is always satisfied according to the definitions of *judge* and *repair*, so we formulate it as a lemma.

**Lemma 2 (JudgeRepair)** *Any bi-filter  $F$  satisfies the JudgeRepair property:*

$$\text{judge}_F \ t = \text{True} \quad \Rightarrow \quad \text{repair}_F \ t = t .$$

*Proof* From  $\text{judge}_F \ t = \text{True}$  we deduce  $F \ t = \text{Nothing}$ , which implies  $\text{repair}_F \ t = t$ .  $\square$

In the next section, we will describe how to fit generalised parsers and bi-filters into the architecture of BIYACC. To let bi-filters work with the lens between CSTs and ASTs, we require a further property characterising the interaction between the repairing direction of a bi-filter and the get direction of a lens.

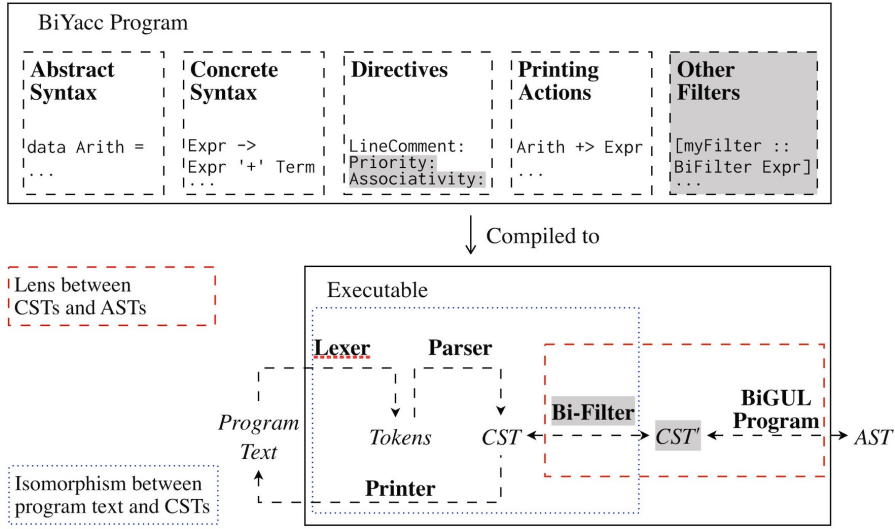
**Definition 7 (PassThrough)** A bi-filter *F* satisfies the *PassThrough* property with respect to a function *get* exactly when

$$\text{get} \circ \text{repair}_F = \text{get} .$$

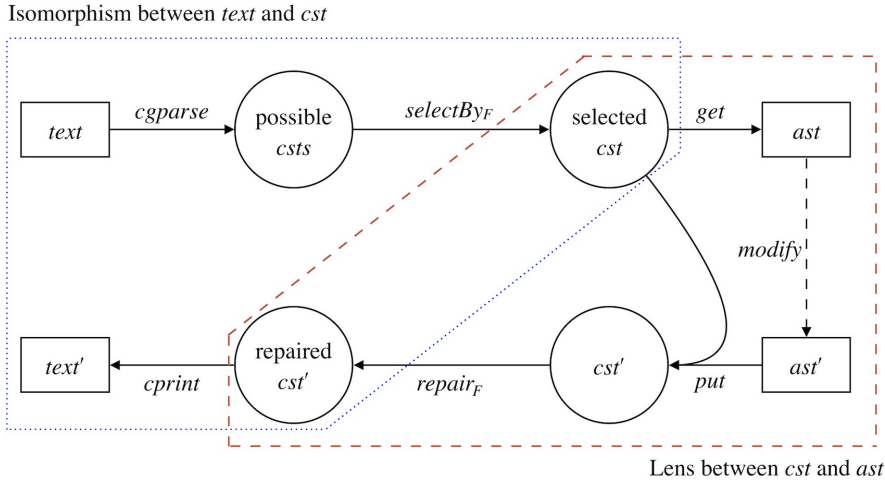
If we think of a *get* function as mapping CSTs to their semantics (in our case ASTs), then the *PassThrough* property is a reasonable requirement since it guarantees that the repaired CST will have the same semantics as before (since it is converted to the same AST). This property will be essential for establishing the well-behavedness of the executable generated by the revised BIYACC.

### 5.3 The New BIYACC System for Ambiguous Grammars

As depicted in Figure 7, the executable generated by the new BIYACC system is still the composition of an isomorphism and a lens, which is the structure we have tried to retain. To precisely identify the changes in several generated components (in the



**Fig. 7** New architecture of BiYACC. (New components are in light grey.)



**Fig. 8** A schematic diagram showing how parsing and printing work with a bi-filter.

executable file) and demonstrate how parsing and printing work with a bi-filter, we present Figure 8 and will use this one instead. In the new system, we will still use the *get* and *put* transformations generated from printing actions and the concrete printer *cprint* from grammars, while the concrete parser *cparse* is replaced with a generalised parser *cgparse*. Additionally, the `#Directives` and `#OtherFilters` parts will be used to generate a bi-filter *F*, whose *judge<sub>F</sub>* (used in the *selectBy<sub>F</sub>* function in Figure 8) and *repair<sub>F</sub>* components are integrated into the isomorphism and lens parts respectively, so that the right-hand side domain of the isomorphism and the source of the lens become



$\text{CST}_F$ , the set of valid CSTs:

$$\text{CST}_F = \{ \text{cst} \in \text{CST} \mid \text{judge}_F \text{ cst} = \text{True} \} .$$

Next, we introduce the (new) isomorphism and lens parts, and prove their inverse properties and well-behavedness respectively.

### 5.3.1 The Revised Isomorphism between Program Text and CSTs

Let us first consider the isomorphism part between `String` and  $\text{CST}_F$ , which is enclosed within the blue dotted lines in Figure 8 and consists of *cprint*, *cgparse*, and *selectBy<sub>F</sub>*:

```

cprint :: CST → String
cgparse :: String → [CST]
selectByF :: [CST] → Maybe CSTF
selectByF csts = case selectBy judgeF csts of
    [cst] → Just cst
    _     → Nothing
selectBy :: (a → Bool) → [a] → [a]
selectBy p [] = []
selectBy p (x : xs) | p x = x : selectBy p xs
selectBy p (x : xs) | otherwise = selectBy p xs .

```

In the parsing direction, first *cgparse* produces all the CSTs; then *selectBy<sub>F</sub>* utilises a function *selectBy* and a predicate *judge<sub>F</sub>* to (try to) select the only correct *cst*; if there is no correct CST or more than one correct CST, `Nothing` is returned. The function *selectBy*, which selects from the input list exactly the elements satisfying the given predicate, is named *filter* in HASKELL's standard libraries but renamed here to avoid confusion. In the printing direction, we still use *cprint* to flatten a (correct) CST back to program text. Formally, constructed from *cgparse* and *cprint*, the two directions of the isomorphism are

```

cparseF :: String → Maybe CSTF
cparseF = selectByF ∘ cgparse
cprintF :: CSTF → Maybe String
cprintF = Just ∘ cprint .

```

We are eager to give the revised version of the inverse properties (Theorem 3) and their proofs, which, however, depend on two assumptions about generalised parsers and bi-filters. So let us present them in order.

**Definition 8 (Generalised Parser Correctness)** A generalised parser *cgparse* is *correct* with respect to a printer *cprint* exactly when

$$\text{cgparse } \text{text} = \{ \text{cst} \in \text{CST} \mid \text{cprint } \text{cst} = \text{text} \} .$$

This is exactly Definition 3.7 of Klint and Visser [26]. We remind the reader again that we use sets and lists interchangeably for the parsing results.

**Definition 9 (Bi-Filter Completeness)** A bi-filter  $F$  is *complete* with respect to a printer  $cprint$  exactly when

$$text \in \text{Img } cprint \Rightarrow |\{cst \in \text{CST}_F \mid cprint \text{ } cst = text\}| = 1.$$

( $\text{Img } f = \{y \mid \exists x. fx = y\}$  is the image of the function  $f$ .)

This is revised from Definition 4.3 of Klint and Visser [26], where they require that filters select exactly one CST and reject all the others. Since it is undecidable to judge whether a given context-free grammar is ambiguous [10], we cannot tell whether a (bi-)filter (for the full CFG) is complete, either. But still, some checks can be performed in simple cases, as stated in Section 7.

The following two lemmas connect our two assumptions, Definitions 8 and 9, with the definitions of  $cparse_F$  and  $cprint_F$ .

**Lemma 3** Given  $cparse_F$  and  $cprint_F$  where  $cgparse$  is correct and  $F$  is complete with respect to  $cprint$ , we have

$$text \in \text{Img } cprint \Rightarrow \exists cst \in \text{CST}_F. cparse_F \text{ } text = \text{Just } cst \wedge cprint \text{ } cst = text.$$

*Proof* We reason:

$$\begin{aligned} & selectBy_F (cgparse \text{ } text) \\ = & \{ \text{Definition of } selectBy_F \} \\ & \text{case } selectBy \text{ } judge_F (cgparse \text{ } text) \text{ of } \{ [cst] \rightarrow \text{Just } cst; \_ \rightarrow \text{Nothing} \} \\ = & \{ \text{Generalised Parser Correctness} \} \\ & \text{case } selectBy \text{ } judge_F \{ cst \in \text{CST} \mid cprint \text{ } cst = text \} \text{ of } \\ & \quad \{ [cst] \rightarrow \text{Just } cst; \_ \rightarrow \text{Nothing} \} \\ = & \{ selectBy \text{ } judge_F \text{ only selects correct CSTs regarding } F \} \\ & \text{case } \{ cst \in \text{CST}_F \mid cprint \text{ } cst = text \} \text{ of } \{ [cst] \rightarrow \text{Just } cst; \_ \rightarrow \text{Nothing} \} \\ = & \{ \text{Bi-Filter Completeness, } \exists cst' \text{ s.t. } \{ cst \in \text{CST}_F \mid cprint \text{ } cst = text \} = [cst'] \} \\ & \text{case } [cst'] \text{ of } \{ [cst] \rightarrow \text{Just } cst; \_ \rightarrow \text{Nothing} \} \\ = & \{ \text{Definition of case} \} \\ & \text{Just } cst. \end{aligned}$$

Moreover,  $cst$  satisfies  $cprint \text{ } cst = text$ , since the latter is the comprehension condition of the set from which  $cst$  is chosen, and therefore  $cprint_F \text{ } cst = \text{Just } text$ .  $\square$

**Lemma 4 (Printer Injectivity)** If  $F$  is a complete bi-filter, then  $cprint_F$  is injective.

*Proof* Assume that  $cst, cst' \in \text{CST}_F$  and  $cprint \text{ } cst = cprint \text{ } cst' = text$  for some  $text$ ; that is, both  $cst$  and  $cst'$  are in the set  $P = \{cst \in \text{CST}_F \mid cprint \text{ } cst = text\}$ . Since  $text \in \text{Img } cprint$ , by the completeness of  $F$  we have  $|P| = 1$ , and hence  $cst = cst'$ .  $\square$

We can now prove a generalised version of Theorem 1 for ambiguous grammars.

**Theorem 3 (Inverse Properties with Bi-Filters)** Given  $cparse_F$  and  $cprint_F$  where  $cgparse$  is correct and  $F$  is complete, we have the following:

$$cparse_F \text{ } text = \text{Just } cst \Rightarrow cprint_F \text{ } cst = \text{Just } text \quad (6)$$

$$cprint_F \text{ } cst = \text{Just } text \Rightarrow cparse_F \text{ } text = \text{Just } cst. \quad (7)$$

*Proof* For (6): Let  $\text{Just } cst = \text{selectBy}_F (\text{cgparse } text)$ . According to the definition of  $\text{selectBy}_F$ , we have  $cst \in \text{cgparse } text$ . By Generalised Parser Correctness  $\text{cprint } cst = text$ , and therefore  $\text{cprint}_F cst = \text{Just } text$ .

For (7): The antecedent implies  $\text{cprint } cst = text$ . By Lemma 3, we have  $\text{cparse}_F text = \text{Just } cst'$  for some  $cst' \in \text{CST}_F$  such that  $\text{cprint}_F cst' = \text{Just } text = \text{cprint}_F cst$ . By Lemma 4 we know  $cst' = cst$ , and thus  $\text{cparse}_F text = \text{Just } cst$ .  $\square$

### 5.3.2 The Revised Lens between CSTs and ASTs

Recall that the `#Action` part of a BiYACC program produces a lens (BiGUL program) consisting of a pair of well-behaved *get* and *put* functions:

$$\begin{aligned} \text{get} &:: \text{CST} \rightarrow \text{Maybe AST} \\ \text{put} &:: \text{CST} \rightarrow \text{AST} \rightarrow \text{Maybe CST} . \end{aligned}$$

To work with a bi-filter  $F$ , in particular its  $\text{repair}_F$  component, they need to be adapted to  $\text{get}_F$  and  $\text{put}_F$ , which accept only valid CSTs:

$$\begin{aligned} \text{get}_F &:: \text{CST}_F \rightarrow \text{Maybe AST} \\ \text{get}_F &= \text{get} \\ \text{put}_F &:: \text{CST}_F \rightarrow \text{AST} \rightarrow \text{Maybe CST}_F \\ \text{put}_F \text{ cst ast} &= \text{fmap } \text{repair}_F (\text{put } \text{cst } \text{ast}) \end{aligned}$$

where *fmap* is a standard HASKELL library function defined (for Maybe) by

$$\begin{aligned} \text{fmap} &:: (a \rightarrow b) \rightarrow \text{Maybe } a \rightarrow \text{Maybe } b \\ \text{fmap } f \text{ Nothing} &= \text{Nothing} \\ \text{fmap } f (\text{Just } x) &= \text{Just } (f x) . \end{aligned}$$

We will need a lemma about *fmap*, which can be straightforwardly proved by a case analysis.

**Lemma 5** *If  $\text{fmap } f \text{ mx} = \text{Just } y$ , then there exists  $x$  such that  $\text{mx} = \text{Just } x$  and  $f x = y$ .*

Now we prove that  $\text{get}_F$  and  $\text{put}_F$  are well-behaved, which is a generalisation of Theorem 2 for ambiguous grammars.

**Theorem 4 (Well-behavedness with Bi-Filters)** *Given a complete bi-filter  $F$  and a well-behaved lens consisting of *get* and *put*, if *get* and  $F$  additionally satisfy *PassThrough*, then the  $\text{get}_F$  and  $\text{put}_F$  functions with respect to  $F$  are also well-behaved:*

$$\text{put}_F \text{ cst ast} = \text{Just } cst' \quad \Rightarrow \quad \text{get}_F \text{ cst}' = \text{Just } \text{ast} \quad (8)$$

$$\text{get}_F \text{ cst} = \text{Just } \text{ast} \quad \Rightarrow \quad \text{put}_F \text{ cst ast} = \text{Just } \text{cst} . \quad (9)$$

*Proof* For (8): The antecedent expands to  $\text{fmap } \text{repair}_F (\text{put } \text{cst } \text{ast}) = \text{Just } cst'$ , which, by Lemma 5, implies  $\text{put } \text{cst } \text{ast} = \text{Just } cst''$  for some  $cst''$  such that  $\text{repair}_F \text{ cst}''$

$= cst'$ . Now we reason:

$$\begin{aligned}
& get_F cst' \\
= & \{ \text{Definition of } get_F \text{ and } cst \in CST_F \} \\
& get cst' \\
= & \{ \text{Definition of } cst' \} \\
& get (repair_F cst'') \\
= & \{ \text{PassThrough} \} \\
& get cst'' \\
= & \{ \text{PutGet} \} \\
& Just ast .
\end{aligned}$$

For (9):

$$\begin{aligned}
& put_F cst ast \\
= & \{ \text{Definition of } put_F \} \\
& fmap repair_F (put cst ast) \\
= & \{ \text{GetPut} \} \\
& fmap repair_F (Just cst) \\
= & \{ \text{Definition of } fmap \} \\
& Just (repair_F cst) \\
= & \{ \text{Since } cst \in CST_F, judge_F cst = \text{True}. \text{ By JudgeRepair} \} \\
& Just cst .
\end{aligned}$$

□

## 5.4 Bi-Filter Directives

Until now, we have only considered working with a single bi-filter, but this is without loss of generality because we can provide a bi-filter composition operator (Section 5.4.1) so that we can build large bi-filters from small ones. This is a suitable semantic foundation for introducing YACC-like directives for specifying priority and associativity into BiYACC (Section 5.4.2), since we can give these directives a bi-filter semantics and interpret a collection of directives as the composition of their corresponding bi-filters. We will also discuss some properties related to this composition (Section 5.4.3).

### 5.4.1 Bi-Filter Composition

We start by defining bi-filter composition, with the intention of making the net effect of applying a sequence of bi-filters one by one the same as applying their composite. Although the intention is better captured by Lemma 6, which describes the *repair* and *judge* behaviour of a composite bi-filter in terms of the component bi-filters, we give the definition of bi-filter composition first.

**Definition 10 (Bi-Filter Composition)** The composition of two bi-filters is defined by

$$\begin{aligned}
 (\triangleleft) &:: (t \rightarrow \text{Maybe } t) \rightarrow (t \rightarrow \text{Maybe } t) \rightarrow (t \rightarrow \text{Maybe } t) \\
 (j \triangleleft i) t &= \text{case } i t \text{ of} \\
 &\quad \text{Nothing} \rightarrow j t \\
 &\quad \text{Just } t' \rightarrow \text{case } j t' \text{ of} \\
 &\quad \quad \text{Nothing} \rightarrow \text{Just } t' \\
 &\quad \quad \text{Just } t'' \rightarrow \text{Just } t'' .
 \end{aligned}$$

When applying a composite bi-filter  $j \triangleleft i$  to a tree  $t$ , if  $t$  is correct with respect to  $i$  (i.e.  $i t = \text{Nothing}$ ), we directly pass the original tree  $t$  to  $j$ ; otherwise  $t$  is repaired by  $i$ , yielding  $t'$ , and we continue to use  $j$  to repair  $t'$ . Note that if  $j t' = \text{Nothing}$ , we return the tree  $t'$  instead of  $\text{Nothing}$ .

**Lemma 6** For a composite bi-filter  $j \triangleleft i$ , the following two equations hold:

$$\begin{aligned}
 \text{repair } (j \triangleleft i) t &= (\text{repair}_j \circ \text{repair}_i) t \\
 \text{judge } (j \triangleleft i) t &= \text{judge}_j t \wedge \text{judge}_i t .
 \end{aligned}$$

*Proof* By the definition of bi-filter composition.  $\square$

Composition of bi-filters should still be a bi-filter and satisfy `RepairJudge` and `PassThrough`. This is not always the case though—to achieve this, we need some additional constraint on the component bi-filters, as formulated below.

**Definition 11** Let  $i$  and  $j$  be bi-filters. We say that  $j$  *respects*  $i$  exactly when

$$\text{judge}_i t = \text{True} \quad \Rightarrow \quad \text{judge}_i (\text{repair}_j t) = \text{True} .$$

If  $j$  respects  $i$ , then a later applied  $\text{repair}_j$  will never break what may already be repaired by a previous  $\text{repair}_i$ . Thus in this case we can safely compose  $j$  after  $i$ . This is proved as the following theorem.

**Theorem 5** Let  $i$  and  $j$  be bi-filters (satisfying `RepairJudge` and `PassThrough`). If  $j$  respects  $i$ , then  $j \triangleleft i$  also satisfy `RepairJudge` and `PassThrough`.

*Proof* For `RepairJudge`, we reason:

$$\begin{aligned}
 &\text{judge } (j \triangleleft i) (\text{repair } (j \triangleleft i) t) \\
 = &\quad \{ \text{Lemma 6} \} \\
 &\text{judge } (j \triangleleft i) (\text{repair}_j (\text{repair}_i t)) \\
 = &\quad \{ \text{Lemma 6} \} \\
 &\text{judge}_j (\text{repair}_j (\text{repair}_i t)) \wedge \text{judge}_i (\text{repair}_j (\text{repair}_i t)) \\
 = &\quad \{ \text{RepairJudge of } j \} \\
 &\text{True} \wedge \text{judge}_i (\text{repair}_j (\text{repair}_i t)) \\
 = &\quad \{ \text{judge}_i (\text{repair}_i t') = \text{True}; j \text{ respects } i \} \\
 &\text{True} \wedge \text{True} \\
 = &\quad \text{True} .
 \end{aligned}$$

And for PassThrough:

$$\begin{aligned}
 & \text{get}(\text{repair}(j \triangleleft i) t) \\
 = & \quad \{ \text{Lemma 6} \} \\
 & \text{get}(\text{repair}_j(\text{repair}_i t)) \\
 = & \quad \{ \text{PassThrough of } j \} \\
 & \text{get}(\text{repair}_i t) \\
 = & \quad \{ \text{PassThrough of } i \} \\
 & \text{get } t .
 \end{aligned}$$

□

#### 5.4.2 Priority and Associativity Directives

To relieve the burden of writing bi-filters manually and guaranteeing respect among bi-filters being composed, we provide some directives for constructing bi-filters dealing with priority<sup>7</sup> and associativity, which are generally comparable to YACC’s conventional disambiguation directives. The bi-filter directives in a BIYACC program can be thought of as specifying ‘production priority tables’, analogous to the operator precedence tables of, for example, the C programming language [24] (chapter *Expressions*) and HASKELL [34] (page 51). The main differences (in terms of the parsing direction) are as follows:

- For bi-filters, priority can be assigned independently of associativity and vice versa, while the YACC-style approach does not permit so—by design, when the YACC directives (%left, %right, and %nonassoc) are used on multiple tokens, they necessarily specify both the precedence and associativity of those tokens.
- For bi-filters, priority and associativity directives may be used to specify more than one production priority tables, making it possible to put unrelated operators in different tables and avoid (unnecessarily) specifying the relationship between them. It is impossible to do so with the YACC-style approach, for its concise syntax only allows a single operator precedence table.

(The bi-filter semantics of) our bi-filter directives repair CSTs violating priority and associativity constraints by adding parentheses—for example, if the production of addition expressions in Figure 6 is left-associative, then we can repair  $\#_{\text{Plus } 1}(\text{Plus } 2 \ 3)$  by adding parentheses around the right subtree, yielding  $\#_{\text{Plus } 1}(\text{Paren}(\text{Plus } 2 \ 3))$ , provided that the grammar has a production of parentheses annotated with the *bracket attribute* [8, 53]:

```
Expr -> ...
      | [Paren] '(' Expr ')' {# Bracket #} .
```

It instructs our bi-filter directives to use this production when parentheses need to be added. Internally, from the production and bracket attribute annotation, a type class AddParen and corresponding instances for each data type generated from concrete syntax (Expr for this example) are automatically created:

<sup>7</sup> The YACC-style approach adopts the word *precedence* [23] while the filter-based approaches tend to use the word *priority* [9, 26]. We follow the traditions and use either word depending on the context.

```

class AddParen t where
  canAddPar :: t -> Bool
  addPar    :: t -> t

```

where `canAddPar` tells whether a CST can be wrapped in a parenthesis structure and `addPar` adds that structure if it is possible or behaves as an identity function otherwise. This makes it possible to automatically generate bi-filters to repair incorrect CSTs (and help the user to define their own bi-filters more easily—see Section 5.5).

In order for bi-filter directives to work correctly, the user should notice the following requirements: (i) Directives shall not mention the parenthesis production annotated with bracket attribute so that they *respect* each other and work properly (as introduced in Definition 11). (ii) Suppose that the parenthesis production is  $NT \rightarrow \alpha NT_R \beta$  where  $\alpha$  and  $\beta$  denote a sequence of terminals and  $NT_R$  is a possibly different nonterminal from  $NT$  (on the right-hand side of the production)—for instance, `Expr -> '(' Expr ')'` above—there shall be exactly one printing action defined for the parenthesis production in the form of  $v \mapsto \alpha[v \mapsto NT_R]\beta$  for the `PassThrough` property to hold: for any CST, the (added) parenthesis structure will all be dropped through the conversion to its AST.

Next we introduce our priority and associativity directives and their bi-filter semantics. From a directive, we first generate a bi-filter that checks and repairs only the top of a tree; this bi-filter is then lifted to check and repair all the subtrees in a tree. In the following we will give the semantics of the directives in terms of the generation of the top-level bi-filters, and then discuss the lifted bi-filters and other important properties they satisfy in Section 5.4.3.

### Priority Directives

A priority directive defines relative priority between two productions; it removes (in the parsing direction) or repairs (in the printing direction) CSTs in which a node of lower priority is a direct child of the node of higher priority. For instance, we can define that (the production of) multiplication has higher priority than (the production of) addition for the grammar in Figure 6 by writing

```
Expr -> Expr '*' Expr > Expr -> Expr '+' Expr ; or just Times > Plus ; .
```

The directive first produces the following top-level bi-filter:<sup>8</sup>

```

fTimesPlusPrio (Times t1 t2 t3) =
  case or [match t1 p, match t2 p, match t3 p, False] of
    False -> Nothing
    True  -> Just (Times (if match t1 p then addPar t1 else t1)
                        (if match t2 p then addPar t2 else t2)
                        (if match t3 p then addPar t3 else t3))
  where p = Plus undefined undefined .

```

<sup>8</sup> Although terminals such as `'*'` and `'+'` are uniquely determined by constructors and not explicitly included in the CSTs, there are fields in CSTs for holding whitespaces after them. Thus `Times` still has three subtrees. Also, for simplicity, the bi-filter `fTimesPlusPrio` attempts to repair the whitespace subtree `t2` even though the repair can never happen since `t2` cannot match `p`.

We first check whether any of the subtrees  $t_1$ ,  $t_2$ , and  $t_3$  violates the priority constraint, i.e. having `Plus` as its top-level constructor—this is checked by the `match` function, which compares the top-level constructors of its two arguments. The resulting boolean values are aggregated using the list version of logical disjunction  $or :: [Bool] \rightarrow Bool$ . If there is any incorrect part, we repair it by inserting a parenthesis structure using `addPar`.

In general, the syntax of priority directives is

$$\begin{aligned} Priority &::= \text{'Priority:' } PDirective^+ \\ PDirective &::= ProdOrCons \text{'>'} ProdOrCons \text{';' } \\ &\quad | ProdOrCons \text{'<'} ProdOrCons \text{';' } \\ ProdOrCons &::= Prod \mid Constructor \\ Prod &::= Nonterminal \text{'->'} Symbol^+ \end{aligned}$$

where *Constructor* and *Symbol* are already defined in Figure 4; for each priority declaration, we can use either *productions* or their names (i.e. constructors).

If the user declares that a production  $NT_1 \rightarrow RHS_1$  has higher priority than another production  $NT_2 \rightarrow RHS_2$ , the following priority bi-filter will be generated:

```

TOPRIOFILTER[(RHS1, NT1, RHS2, NT2)] =
  'f'conRHS1 conRHS2 'Prio' '(conRHS1 FILLVARS(RHS1))' =
    'case or [' 'match' t 'p,' | t ∈ FILLVARS(RHS1)] 'False' ']' of
      'False -> Nothing'
      'True -> Just (conRHS1 (REPAIR(t) | t ∈ FILLVARS(RHS1)))'
    'where p = ' CON(NT2, RHS2) FILLUNDEFINED(RHS2)
  'f'conRHS1 conRHS2 'Prio' '_' '=' 'Nothing'
REPAIR(t) = '(if match' t 'p' 'then addPar' t 'else' t)'
conRHS1 = CON(NT1, RHS1)
conRHS2 = CON(NT2, RHS2) .

```

`CON` looks up constructor names for input productions (divided into nonterminals and right-hand sides); `FILLVARS(nt)` generates variable names for each terminal and nonterminal in *nt* (here  $RHS_1$ ); `FILLUNDEFINED` is similar to `FILLVARS` but it produces undefined values instead. If productions are referred to using their constructors, we can simply look up the nonterminals and right-hand sides and use the same code generation strategy.

*Transitive Closures.* In the same way as conventional YACC-style approaches, the priority directives are considered transitive. For instance,

```

Expr -> Expr '*' Expr > Expr -> Expr '+' Expr ;
Expr -> Expr '+' Expr > Expr -> Expr '&' Expr ;

```

implies that `Expr -> Expr '*' Expr > Expr -> Expr '&' Expr ;`. The feature is important in practice since it greatly reduces the amount of routine code the user needs to write (for large grammars).



### Associativity Directives

Associativity directives assign (left- or right-) associativity to productions. A left-associativity directive bans (or repairs, in the printing direction) CSTs having the pattern in which a parent and its right-most subtree are both left-associative, if the (relative) priority between the parent and the subtree is not defined; a right-associativity directive works symmetrically.

As an example, we can declare that both addition and subtraction are left-associative (for the grammar in Figure 6) by writing

```
Left: Expr -> Expr '+' Expr, Expr -> Expr '-' Expr;
```

or just `Left: Plus, Minus;`. Since the relative priority between `Plus` and `Minus` is not defined, we generate top-level bi-filters for all the four possible pairs formed out of `Plus` and `Minus`:

```
fPlusPlusLAssoc (Plus t1 t2 t3) =
  if match t3 p then Just (Plus t1 t2 (addPar t3)) else Nothing
  where p = Plus undefined undefined undefined
fPlusPlusLAssoc _ = Nothing

fMinusMinusLAssoc (Minus t1 t2 t3) =
  if match t3 p then Just (Minus t1 t2 (addPar t3)) else Nothing
  where p = Minus undefined undefined undefined
fMinusMinusLAssoc _ = Nothing

fPlusMinusLAssoc (Plus t1 t2 t3) =
  if match t3 p then Just (Plus t1 t2 (addPar t3)) else Nothing
  where p = Minus undefined undefined undefined
fPlusMinusLAssoc _ = Nothing

fMinusPlusLAssoc (Minus t1 t2 t3) =
  if match t3 p then Just (Minus t1 t2 (addPar t3)) else Nothing
  where p = Plus undefined undefined undefined
fMinusPlusLAssoc _ = Nothing .
```

For instance, `fPlusPlusLAssoc` accepts  $\sharp_{\text{Plus}} (\text{Plus } 1 \ 2) \ 3$  but not  $\sharp_{\text{Plus } 1} (\text{Plus } 2 \ 3)$ , which is repaired to  $\sharp_{\text{Plus } 1} (\text{Paren } (\text{Plus } 2 \ 3))$ .

Generally, the syntax of associativity directives is

$$\begin{aligned} \text{Associativity} &::= \text{'Associativity:' } \text{LeftAssoc } \text{RightAssoc} \\ \text{LeftAssoc} &::= \text{'Left:' } \text{ProdOrCons}^+ \{', '\} \text{' ;' } \\ \text{RightAssoc} &::= \text{'Right:' } \text{ProdOrCons}^+ \{', '\} \text{' ;' } . \end{aligned}$$

Now we explain the generation of (top-level) bi-filters from associativity directives. We will consider only left-associativity directives, as right-associativity directives are symmetric. For every pair of left-associative productions whose relative priority is not defined—including cases where the two productions are the same—we generate a bi-filter to repair CSTs whose top uses the first production and whose right-most child uses the second production. Let  $NT_1 \rightarrow \alpha_1 NT_{1R}$  and  $NT_2 \rightarrow \alpha_2 NT_{2R}$  be two such productions, where  $\alpha_1$  ( $\alpha_2$ ) matches a sequence of arbitrary symbols of any length

and  $NT_{1R}$  ( $NT_{2R}$ ) is the right-most symbol and must be a nonterminal. (If it is not a nonterminal, it is meaningless to discuss associativity.) The generated bi-filter is

```

toLASSOCFILTER[ $\alpha_1 NT_{1R}, NT_1, \alpha_2 NT_{2R}, NT_2$ ] =
  'f' conRHS1 conRHS2 'LAssoc' '(' conRHS1 FILLVARS( $\alpha_1 NT_{1R}$ ) ')' '='
  'if match ' ntrVar ' p'
    'then Just ( conRHS1 FILLVARS( $\alpha_1$ ) (addPar' ntrVar '))'
    'else Nothing'
  'where p = ' conRHS2 FILLUNDEFINED( $\alpha_2 NT_{2R}$ )
  'f' conRHS1 'LAssoc' '_' '=' 'Nothing'
conRHS1 = CON( $NT_1, \alpha_1 NT_{1R}$ )
conRHS2 = CON( $NT_2, \alpha_2 NT_{2R}$ )
ntrVar = FILLVARFROM(LENGTH( $\alpha_1$ ),  $NT_{1R}$ ) .

```

Functions CON, FILLUNDEFINED, and FILLVAR have the same behaviour as before; FILLVARFROM (which is a variation of FILLVARS) generates variable names for each terminal and nonterminal in its argument with suffix integers counting from a given number to avoid name clashing.

*Handling Injective Productions.* Sometimes the grammar may contain injective productions (also called chain productions) [9], which have only a single nonterminal on their right-hand side, like  $\text{InfE} \rightarrow [\text{FromE}] \text{Exp}$ . When we use it to define a grammar

```

InfE -> [FromE] Exp
Exp  -> [Plus] InfE '+' InfE
      | [Times] InfE '*' InfE ,

```

program text  $1 + 2 * 3$  will be parsed to two CSTs, namely  $\text{cst}_1 = \#Plus (\text{FromE } 1) (\text{FromE } (\text{Times } 2 \ 3))$  and  $\text{cst}_2 = \#Times (\text{FromE } (\text{Plus } 1 \ 2) (\text{FromE } 3))$ , and we want to spot  $\text{cst}_2$  and discard it using the priority directive  $\text{Times} > \text{Plus}$ . If handled naively, the bi-filter generated from the directive would only remove CSTs having pattern  $\text{Times} (\text{Plus } \_ \_)$  (and two other similar ones), but  $\text{cst}_2$  would not match the pattern due to the presence of the  $\text{FromE}$  node between  $\text{Times}$  and  $\text{Plus}$ . We made some effort in the implementation to make the *match* function ignore the nodes corresponding to injective productions ( $\text{FromE}$  in this case).

#### 5.4.3 Properties of the Generated Bi-Filters

We discuss some properties of the bi-filters generated from our priority and associativity directives, to justify that it is safe to use these bi-filters without disrupting the well-behavedness of the whole system. Specifically:

- The generated top-level bi-filters satisfy *RepairJudge*, and it is easy to write actions to make them satisfy *PassThrough*.
- The bi-filters lifted from the top-level bi-filters still satisfy *RepairJudge* and *PassThrough*.
- The lifted bi-filters are *commutative*, which not only implies that all such bi-filters respect each other and can be composed in any order, but also guarantees that we do not have to worry about the order of composition since it does not affect the behaviour.

We will give only high-level, even informal, arguments for these properties, since, due to the generic nature of the definitions of these bi-filters (in terms of *Scrap Your Boilerplate* [30]), to give formal proofs we would have to introduce rather complex machinery (e.g., datatype-generic induction), which would be tedious and distracting.

*Top-level bi-filters.* The fact that the generated top-level bi-filters satisfy `RepairJudge` can be derived from the requirement that the directives do not mention the parenthesis production. Because of the requirement, in the generated bi-filters, repairing is always triggered by matching a non-parenthesis production, and after that repairing will not be triggered again because a parenthesis production will have been added. For example, in the bi-filter `fTimesPlusPrio` (in Section 5.4.2), with `match t1 p`, `match t2 p`, and `match t3 p` we check whether `t1`, `t2`, and `t3` has `Plus` as the top-level production, which is different from the parenthesis production `Paren`; if any of the matching succeeds, say `t1`, then `addPar t1` will add `Paren` at the top of `t1`, and `match (addPar t1) p` is guaranteed to be `False`, so the subsequent invocation of `judge fTimesPlusPrio` will return `True`. For `PassThrough`, since all the top-level bi-filters do is add parenthesis productions, we can simply make sure that appearances of the parenthesis production are ignored by `get`, i.e. `get (addPar s) = get s` for all `s`; this, by well-behavedness, is the same as making `put` (printing actions) skip over parentheses. For example, for the grammar in Figure 6, we should write `t +> '(' [t +> Expr] ')'` as the only printing action mentioning parentheses, which means that `put (Paren s) t = fmap Paren (put s t)` for all `s` and `t`. Then the following reasoning implies that `get (Paren s) = get s` for all `s`:

$$\begin{aligned}
& \text{get (Paren s)} = \text{Just t} \\
& \Leftrightarrow \{ \Rightarrow \text{ by GetPut and } \Leftarrow \text{ by PutGet} \} \\
& \quad \text{put (Paren s) t} = \text{Just (Paren s)} \\
& \Leftrightarrow \{ \text{By the above statement: } \text{put (Paren s) t} = \text{fmap Paren (put s t)} \} \\
& \quad \text{fmap Paren (put s t)} = \text{Just (Paren s)} \\
& \Leftrightarrow \{ \text{Lemma 5 and the definition of fmap} \} \\
& \quad \text{put s t} = \text{Just s} \\
& \Leftrightarrow \{ \Rightarrow \text{ by PutGet and } \Leftarrow \text{ by GetPut} \} \\
& \quad \text{get s} = \text{Just t}
\end{aligned}$$

for all `s` and `t`.

*Lifted bi-filters.* The lifted bi-filters apply the top-level bi-filters to all the subtrees in a CST in a bottom-up order. Formally, we can define, datatype-generically, a lifted bi-filter as a composition of top-level bi-filters, and use datatype-generic induction to prove that there is suitable respect among the top-level bi-filters being composed, and that the lifted bi-filter satisfies `RepairJudge` and `PassThrough` if the top-level ones do. But here we provide only an intuitive argument. What the lifted bi-filters do is find all prohibited pairs of adjoining productions and separate all the pairs by adding parenthesis productions. For `RepairJudge`, since all prohibited pairs are eliminated after repairing, there will be nothing left to be repaired in the resulting CST, which will therefore be deemed valid. For `PassThrough`, the intuition is the same as that for the top-level bi-filters.

*Commutativity.* Composite bi-filters  $i \triangleleft j$  and  $j \triangleleft i$  may have different behaviour, so in general we need to know the order of composition to figure out the exact behaviour of a composite bi-filter. This can be difficult when using our bi-filter directives, since a lot of bi-filters are implicitly generated from the directives, and it is not straightforward to specify the order in which all the explicitly and implicitly generated bi-filters are composed. Fortunately we do not need to do so, for all the bi-filters generated from the directives are *commutative*, meaning that the order of composition does not affect the behaviour.

**Definition 12 (Bi-Filter Commutativity)** Two bi-filters  $i$  and  $j$  are *commutative* exactly when

$$\text{repair}_i \circ \text{repair}_j = \text{repair}_j \circ \text{repair}_i.$$

By Lemma 6, this implies  $\text{repair}(i \triangleleft j) = \text{repair}(j \triangleleft i)$ . Note that  $\text{judge}(i \triangleleft j) = \text{judge}(j \triangleleft i)$  by definition, so we do not need to require this in the definition of commutativity.

An important fact is that commutativity is stronger than respect, so it is always safe to compose commutative bi-filters.

**Lemma 7** *Commutative bi-filters respect each other.*

*Proof* Given commutative bi-filters  $i$  and  $j$ , we show that  $j$  respects  $i$ . Suppose that  $\text{judge}_i t = \text{True}$  for a given tree  $t$ . Then

$$\begin{aligned} & \text{judge}_i(\text{repair}_j t) \\ &= \{ \text{repair}_i t = t, \text{ since } \text{judge}_i t = \text{True} \} \\ & \quad \text{judge}_i(\text{repair}_j(\text{repair}_i t)) \\ &= \{ i \text{ and } j \text{ are commutative} \} \\ & \quad \text{judge}_i(\text{repair}_i(\text{repair}_j t)) \\ &= \{ \text{RepairJudge} \} \\ & \quad \text{True}. \end{aligned}$$

It follows by symmetry that  $i$  respects  $j$  as well.  $\square$

Now let us consider why any two different lifted bi-filters are commutative. (Commutativity is immediate if the two bi-filters are the same.) There are two key facts that lead to commutativity: (i) repairing does not introduce more prohibited pairs of productions, and (ii) the prohibited pairs of adjoining productions checked and repaired by the two bi-filters are necessarily different. Therefore the two bi-filters always repair different parts of a tree, and can repair the tree in any order without changing the final result. Fact (i) is, again, due to the requirement that the directives do not mention the parenthesis production, which is the only thing we add to a tree when repairing it. Fact (ii) can be verified by a careful case analysis. For example, we might be worried about the situation where a left-associative directive looks for production  $Q$  used at the right-most position under production  $P$ , while a priority directive also similarly looks for  $Q$  used under  $P$ , but the two directives cannot coexist in the first place since the first directive implies  $P$  and  $Q$  have no relative priority whereas the second one implies  $Q$  has lower priority than  $P$ .

## 5.5 Manually Written Bi-Filters

There are some other ambiguities that our directives cannot eliminate. In these cases, the user can define their own bi-filters and put them in the `#OtherFilters` part in a BiYACC program as shown in Figure 4. The syntax is

$$\begin{aligned} \text{OtherFilters} &::= \text{'[' } HsFunDecl^+ \text{'\{' , '\}' } \text{'\}' } HsCode \\ HsFunDecl &::= HsFunName \text{' :: BiFilter ' Nonterminal .} \end{aligned}$$

That is, this part of the program begins with a list of declarations of the names and types of the user-defined bi-filters, whose HASKELL definitions are then given below.

Now we demonstrate how to manually write a bi-filter by resolving the ambiguity brought by the dangling else problem. But before that, let us briefly review the problem, which arises, for example, in the following grammar:

```
Exp -> [ITE] 'if' Exp 'then' Exp 'else' Exp
      | [IT]  'if' Exp 'then' Exp .
```

With respect to this grammar, the program text `if a then if x then y else z` can be recognised as either `if a then (if x then y else z)` or `if a then (if x then y) else z`. To resolve the ambiguity, usually we prefer the ‘nearest match’ strategy (which is adopted by Pascal, C, and Java): `else` should match its nearest `then`, so that `if a then (if x then y else z)` is the only correct interpretation.

The user may think that the problem can be solved by a priority (bi-)filter `ITE > IT;`, in the hope that the production ‘if-then-else’ binds tighter than the production ‘if-then’. Unfortunately, this is incorrect as pointed out by Klint and Visser [26], because the corresponding (bi-)filter incorrectly rules out the pattern `#ITE _ _ (IT _ _)`, which prints to unambiguous text, e.g., `if a then b else if x then y`. In fact, the (dangling else) problem is tougher than one might think and cannot be solved by any (bi-)filter performing pattern matching with a fixed depth [26].

Klint and Visser [26] proposed an idea to disambiguate the dangling-else grammar: Let Greek letters  $\alpha, \beta, \dots$  match a sequence of symbols of any length. Then the program text `if  $\alpha$  then  $\beta$  else  $\gamma$`  should be banned if the right spine of  $\beta$  contains any `if  $\psi$  then  $\omega$` , as shown in the paper [26]. With the full power of (bi-)filters, which are fully-fledged HASKELL functions, we can implement this solution in the following bi-filter:

```
fCond (ITE c1 e1 e2) = case checkRightSpine e1 of
  True  -> Nothing
  False -> Just (ITE c1 (addPar e1) e2)

-- collect the names of the constructors in the right spine and
-- check if the collected constructors contain "IT"
checkRightSpine t = ... .
```

This bi-filter is commutative with the bi-filters generated from our directives, since it (i) only searches for non-parenthesis productions that are not declared in any other directives, and (ii) inserts only a parenthesis production when repairing incorrect CSTs. The reader may find the code of `checkRightSpine` in more detail in Figure 10.

## 6 Case Studies

The design of BIYACC may look simplistic and make the reader wonder how much it can describe. In fact, BIYACC can already handle real-world language features. For example, Kinoshita and Nakano [25] adopted BIYACC as part of their system for synchronising COQ functions and ~~their~~ corresponding OCAML programs. In this section, we demonstrate BIYACC with a medium-size case study: we use BIYACC to build a pair of parser and reflective printer for the TIGER language [4] and demonstrate some of their uses.

### 6.1 The TIGER Language

TIGER is a statically typed imperative language first introduced in Appel's textbook on compiler construction [4]. Since TIGER's purpose of design is pedagogical, it is not too complex and yet covers many important language features including conditionals, loops, variable declarations and assignments, and function definitions and calls. TIGER is therefore a good case study with which we can test the potential of our BX-based approach to constructing parsers and reflective printers. Some of these features can be seen in this TIGER program:

```
function foo() =
  (for i := 0 to 10
    do (print(if i < 5 then "smaller"
              else "bigger");
        print("\n")))
```

To give a sense of TIGER's complexity, it takes a grammar with 81 production rules to specify TIGER's syntax, while for C89 and C99 it takes respectively 183 and 237 rules without any disambiguation declarations (based on Kernighan and Ritchie [24] and the draft version of 1999 ISO C standard, excluding the preprocessing part). The difference is basically due to the fact that C has more primitive types and various kinds of assignment statements.

Excerpts of the abstract and concrete syntax of TIGER are shown in Figure 9. The abstract syntax is largely the same as the original one defined in Appel's textbook (page 98); as for the concrete syntax, Appel does not specify the whole grammar in detail, so we use a version slightly adapted from Hirzel and Rose's lecture notes [21]. Concretely, we add a parenthesis production to the grammar (and discard it when converting CSTs to ASTs, so that the PassThrough property could be satisfied), since TIGER's original grammar has no parenthesis production and an expression within round parentheses is regarded as a singleton expression sequence. This modification also makes it necessary to change the enclosing brackets for expression sequences from round brackets `()` to curly brackets `{}`, which helps (LALR(1) parsers) to distinguish a singleton expression sequence from an expression within parentheses. There is also another slight change in the definition of ASTs for handling a feature not supported by the current BIYACC: the AST constructors `TFunctionDec` and `TTypeDec` take a single function or type declaration instead of a list of adjacent declarations (for representing mutual recursion) as in Appel [4], since we cannot handle the synchronisation between a list of lists (in ASTs) and a list (in CSTs) with BIYACC's current syntax.

```

#Abstract
type TSymbol = String
data BBool = TT | FF
data MMaybe a = NN | JJ a
data Tuple a b = Tuple a b
data List a = Nil | Cons a (List a)

data TExp = TString String | TInt Int | TNilExp | TCond TExp TExp (MMaybe TExp)
         | TLet (List TDec) TExp | TOp TExp TOper TExp | TExpSeq (List TExp) | ...

data TOper = TPlusOp | TMinusOp | ... | TEqOp | TNeqOp | ...

data TDec = TVarDec TSymbol BBool (MMaybe TSymbol) TExp
         | TTypeDec (Tuple TSymbol TTy) | TFunctionDec TFundec

data TFundec = TFundec TSymbol (List TFieldDec) (MMaybe TSymbol) TExp
...

#Concrete
Exp -> LetExp | ArrExp | IfThen | IfThenElse | Prmtv
     | ForExp | RecExp | WhileExp | Assignment | 'break' ;

VarDec -> 'var' Identifier ':' Exp
        | 'var' Identifier ':' Identifier ':' Exp ;

LValue -> Identifier | OtherLValue ;
OtherLValue -> LValue '.' Identifier
             | Identifier '[' Exp ']' | OtherLValue '[' Exp ']' ;

SeqExp -> '{' ExpSeq '}' ;
ExpSeq -> Exp ';' ExpSeq | Exp ;

Prmtv -> [Paren] '(' Exp ')' {# Bracket #} | CallExp | SeqExp | ...
        | [Or] Prmtv '|' Prmtv | [And] Prmtv '&' Prmtv
        | [Plus] Prmtv '+' Prmtv | [Times] Prmtv '*' Prmtv | ...
        | [Neg] '-' Prmtv | Numeric | String | LValue | 'nil' ;

IfThenElse -> [ITE] 'if' Exp 'then' Exp 'else' Exp ;
IfThen -> [IT] 'if' Exp 'then' Exp ;
...

```

**Fig. 9** An excerpt of TIGER’s abstract and concrete syntax. (Here we define our own BBool type and MMaybe type to avoid name clashing with HASKELL’s built-in ones.)

Following Hirzel and Rose’s specification [21], the disambiguation directives for TIGER are shown in Figure 10; for instance, we define multiplication to be left-associative. The directives also include a concrete treatment for the dangling else problem, which is usually ‘not solved’ when using a YACC-like (LA)LR parser generator to implement parsers: rather than resolving the grammatical ambiguity, we often rely on the default behaviour of the parser generator—preferring shift.

We have successfully tested our BiYACC program for TIGER on all the sample programs provided on the homepage of Appel’s book<sup>9</sup>, including a merge sort implementation and an eight-queen solver, and there is no problem parsing and printing them with well-behavedness guaranteed. In the following subsections, we will present some

<sup>9</sup> <https://www.cs.princeton.edu/~appel/modern/testcases/>

```

#Directives
Priority:
Times > Plus ;
And > Or ; ...

Associativity:
Left: Times, Plus, And ... ;
Right: Assign, ... ;

#OtherFilters
[ fDanglingElse :: BiFilter IfThenElse ]

fDanglingElse (ITE t1 exp1 t2 exp2 t3 exp3) =
  case checkRightSpine exp2 of
    True  -> Nothing
    False -> Just (ITE t1 exp1 t2 (addPar exp2) t3 exp3)

checkRightSpine t = let spineStrs = getRSpineCons t
                     in  and $ map (\str -> str /= "IT") spineStrs

class GetRSpineCons t where
  getRSpineCons :: t -> [String]

instance GetRSpineCons IfThenElse where
  getRSpineCons (ITE _ _ _ _ r) = ["ITE"] ++ getRSpineCons r

instance GetRSpineCons IfThen where
  getRSpineCons (IT _ _ _ r) = ["IT"] ++ getRSpineCons r

instance GetRSpineCons LetExp where
  getRSpineCons (LetExp1 _ _ _ _ _) = ["LetExp1"]
...

```

**Fig. 10** An excerpt of the disambiguation directives for TIGER. (A type class `GetRSpineCons` is defined and implemented for collecting the constructors on the right spine of a given tree. Function `getRSpineCons` is recursively invoked for CSTs whose right-most subtree is (parsed from) a nonterminal.)

printing strategies described in the BIYACC program to demonstrate what BIYACC, in particular reflective printing, can achieve.

## 6.2 Syntactic Sugar and Resugaring

We start with a simple example about syntactic sugar, which is pervasive in programming languages and lets the programmer use some features in an alternative (usually conceptually higher-level) syntax. For instance, TIGER represents boolean values *false* and *true* respectively as zero and nonzero integers, and the logical operators `&` ('and') and `|` ('or') are converted to a conditional structure in the abstract syntax: `e1 & e2` is desugared and parsed to `TCond e1 e2 (TInt 0)` and `e1 | e2` to `TCond e1 (TInt 1) e2`. The printing actions for them in BIYACC are:

```

TExp +> Prmtv
TCond e1 (TInt 1) (JJ e2) +> [e1 +> Prmtv] '|' [e2 +> Prmtv];
TCond e1 e2 (JJ (TInt 0)) +> [e1 +> Prmtv] '&' [e2 +> Prmtv]; .

```

A conventional printer which takes only the AST as input cannot reliably determine whether an abstract expression should be printed to the basic form or the sugared form, whereas a reflective printer can make the correct decision by inspecting the CST.



The idea of *resugaring* [42] is to print evaluation sequences in a core language in terms of a surface syntax. Here we show that, without any extension, BIYACC is already capable of propagating some AST changes that result from evaluation back to the concrete syntax, subsuming a part of Pombrio and Krishnamurthi’s work [42, 43].

We borrow their example of resugaring evaluation sequences for the logical operators ‘or’ and ‘not’, but recast the example in TIGER. The ‘or’ operator has been defined as syntactic sugar in Section 6.2. For the ‘not’ operator, which TIGER lacks, we introduce ‘~’, represented by `TNot` in the abstract syntax. Now consider the source expression

$$\sim 1 \mid \sim 0 ,$$

which is parsed to

$$\text{TCond} (\text{TNot} (\text{TInt } 1)) (\text{TInt } 1) (\text{JJ} (\text{TNot} (\text{TInt } 0))) .$$

A typical evaluator will produce the following evaluation sequence given the above AST:

$$\begin{aligned} & \text{TCond} (\text{TNot} (\text{TInt } 1)) (\text{TInt } 1) (\text{JJ} (\text{TNot} (\text{TInt } 0))) \\ \rightarrow & \text{TCond} (\text{TInt } 0) (\text{TInt } 1) (\text{JJ} (\text{TNot} (\text{TInt } 0))) \\ \rightarrow & \text{TNot} (\text{TInt } 0) \\ \rightarrow & \text{TInt } 1 . \end{aligned}$$

If we perform reflective printing after every evaluation step using BIYACC, we will get the following evaluation sequence on the source:

$$\sim 1 \mid \sim 0 \quad \rightarrow \quad 0 \mid \sim 0 \quad \rightarrow \quad \sim 0 \quad \rightarrow \quad 1 .$$

Due to the PUTGET property, parsing these concrete terms will yield the corresponding abstract terms in the abstract evaluation sequence, and this is exactly Pombrio and Krishnamurthi’s ‘emulation’ property, which they have to prove for their system. For BIYACC, however, the emulation property holds by construction, since BIYACC programs are always well-behaved. Another difference is that we do not need to insert additional information (such as tags) into an AST for recording which surface syntax structure a node comes from. One advantage of our approach is that we keep the abstract syntax pure, so that other tools—the evaluator in particular—can process the abstract syntax without being modified, whereas in Pombrio and Krishnamurthi’s approach, the evaluator has to be adapted to work on an enriched abstract syntax.

### 6.3 Language Evolution

When a language evolves, some new features of the language (e.g. the `foreach` loops introduced in Java 5 [20]) can be implemented by desugaring to some existing features (e.g. ordinary `for` loops), so that the compiler back-end and abstract syntax definition do not need to be extended to handle the new features. As a consequence, all the engineering work about optimising transformations or refactoring [18] that has been developed for the abstract syntax remains valid.

Consider a kind of ‘generalised-if’ expression allowing more than two cases, resembling the alternative construct in Dijkstra’s guarded command language [12]. We extend TIGER’s concrete syntax with the following production rules:

```

Exp    -> ... | Guard | ... ;      CaseBs -> CaseB CaseBs | CaseB ;
Guard  -> 'guard' CaseBs 'end';    CaseB  -> LValue '=' Numeric '->' Exp ; .

```

For simplicity, we restrict the predicate produced by `CaseB` to the form `LValue '=' Numeric`, but in general **the Numeric part** can be any expression computing an integer. The reflective printing actions for this new construct can still be written within `BiYACC`, but require much deeper pattern matching:

```

TExp +> Guard
  TCond (TOp (TVar lv) TEqOp (TInt i)) e1 Nothing +>
    'guard' (CaseBs -> (CaseB -> [lv +> LValue] '=' [i +> Numeric] '->' [e1 +> Exp])
      ) 'end';
  TCond (TOp (TVar lv) TEqOp (TInt i)) e1 (J if2@(TCond _ _)) +>
    'guard' (CaseBs -> (CaseB -> [lv +> LValue] '=' [i +> Numeric] '->' [e1 +> Exp])
      [if2 +> CaseBs]
    ) 'end';
;;
TExp +> CaseBs
  TCond (TOp (TVar lv) TEqOp (TInt i)) e1 Nothing +>
    (CaseB -> [lv +> LValue] '=' [i +> Numeric] '->' [e1 +> Exp]);
  TCond (TOp (TVar lv) TEqOp (TInt i)) e1 (J if2@(TCond _ _)) +>
    (CaseB -> [lv +> LValue] '=' [i +> Numeric] '->' [e1 +> Exp])
    [if2 +> CaseBs];
;; .

```

Although being a little complex, these printing actions are in fact fairly straightforward: The first group of type `Tiger +> Guard` handles the enclosing `guard`–`end` pairs, distinguishes between single- and multi-branch cases, and delegates the latter case to the second group, which prints a list of branches recursively.

This is all we have to do—the corresponding parser is automatically derived and guaranteed to be consistent. Now `guard` expressions are desugared to nested `if` expressions in parsing and preserved in printing, and we can also resugar evaluation sequences on the ASTs to program text. For instance, the following `guard` expression

```

guard  choice = 1 -> 4
      choice = 2 -> 8
      choice = 3 -> 16 end

```

is parsed to

```

TCond (TOp (TVar (TSV "c")) TEqOp (TInt 1)) (TInt 4) (JJ
  (TCond (TOp (TVar (TSV "c")) TEqOp (TInt 2)) (TInt 8) (JJ
    (TCond (TOp (TVar (TSV "c")) TEqOp (TInt 3)) (TInt 16) NN))))

```

where `TSimpleVar` is shortened to `TSV`, and `choice` is shortened to `c`. Suppose that the value of the variable `choice` is 2. The evaluation sequence on the AST will then be:

```

TCond (TOp (TVar (TSV "c")) TEqOp (TInt 1)) (TInt 4) (JJ
  (TCond (TOp (TVar (TSV "c")) TEqOp (TInt 2)) (TInt 8) (JJ
    (TCond (TOp (TVar (TSV "c")) TEqOp (TInt 3)) (TInt 16) NN))))
→ TCond (TInt 0) (TInt 4) (JJ
  (TCond (TOp (TVar (TSV "c")) TEqOp (TInt 2)) (TInt 8) (JJ
    (TCond (TOp (TVar (TSV "c")) TEqOp (TInt 3)) (TInt 16) NN))))
→ TCond (TOp (TVar (TSV "c")) TEqOp (TInt 2)) (TInt 8) (JJ
  (TCond (TOp (TVar (TSV "c")) TEqOp (TInt 3)) (TInt 16) NN))
→ TCond (TInt 1) (TInt 8) (JJ
  (TCond (TOp (TVar (TSV "c")) TEqOp (TInt 3)) (TInt 16) NN))
→ TInt 8 .

```

And the reflected evaluation sequence on the concrete expression will be:

```

guard choice = 1 -> 4
      choice = 2 -> 8
      choice = 3 -> 16 end
↗
→ guard choice = 2 -> 8
      choice = 3 -> 16 end
↗
→ 8 .

```

Reflective printing fails for the first and third steps (the program text becomes an if-then-else expression if we do printing at these steps), but this behaviour in fact conforms to Pombrio and Krishnamurthi’s ‘abstraction’ property, which demands that core evaluation steps that make sense only in the core language must not be propagated to the surface. In our example, the first and third steps in the  $\tau_{\text{Cond}}$ -sequence evaluate the condition to a constant, but conditions in guard expressions are restricted to a specific form and cannot be a constant; evaluation of guard expressions thus has to proceed in bigger steps, throwing away or going into a branch in each step, which corresponds to two steps for  $\tau_{\text{Cond}}$ .

The reader may have noticed that, after the guard expression is reduced to two branches, the layout of the second branch is disrupted; this is because the second branch is in fact printed from scratch. In current BIYACC, the printing from an AST to a CST is accomplished by recursively performing pattern matching on both tree structures. This approach naturally comes with the disadvantage that the matching is mainly decided by the position of the nodes in the AST and CST. Consequently, a minor structural change on the AST may completely disrupt the matching between the AST and the CST.

## 6.4 Other Potential Applications

We conclude this section by shortly discussing several other potential applications. In general, (current) BIYACC can easily and reliably propagate AST changes that have local effect such as replacing part of an AST with a simpler tree, without destroying the layouts and comments of unaffected code. Thus it would not be surprising that BIYACC can also propagate (i) simplification-like optimisations such as constant folding and constant propagation and (ii) some code refactoring transformations such as variable renaming. All these functionalities are achieved for free by one ‘general-purpose’ BIYACC program, which does not need to be tailored for each application.

## 7 Related Work

### 7.1 Unifying Parsing and Printing

Much research has been devoted to describing parsers and printers in a single program. For example, both Rendel and Ostermann [44] and Matsuda and Wang [36, 37] adopt

a combinator-based approach<sup>10</sup> (whereas we use a generator-based approach), where small components are glued together to yield more sophisticated behaviour, and can guarantee properties similar to Theorem 1 with *cst* replaced by *ast* in the equations. (Let us call the variant version Theorem 1', since it will be used quite often later.) In Rendel and Ostermann's system (called 'invertible syntax descriptions', which we shorten to ISDs henceforth), both the parsing and printing semantics are predefined in the combinators and consistency is guaranteed by their partial isomorphisms, whereas in Matsuda and Wang's system (called FLIPPR), the combinators describing pretty printing are translated by a semantic-preserving transformation to a core syntax, which is further processed by their grammar-based inversion system [38] to realise the parsing semantics. Brabrand et al. [7] present a tool XSugar that handles bijections between the XML syntax (representation) and any other syntax (representation) for the same language, guaranteeing that the syntax transformation is reversible. However, the essential factor that distinguishes our system from others is that the printer produced from a BIYACC program is reflective and can deal with synchronisation.

Although the above-mentioned systems are tailored for unifying parsing and printing, there are design differences. An ISD is more like a parser, while FLIPPR lets the user describe a printer: To handle operator priorities, for example, the user of ISDs will assign priorities to different operators, consume parentheses, and use combinators such as `chainl` to handle left recursion in parsing, while the user of FLIPPR will produce necessary parentheses according to the operator priorities. For basic BIYACC (that deals with unambiguous grammars only), the user defines a concrete syntax that has a hierarchical structure (e.g., `Expr`, `Term`, and `Factor`) to express operator priority, and write printing strategies to produce (preserve) necessary parentheses. The user of XSugar will also likely need to use such a hierarchical structure.

It is interesting to note that the part producing parentheses in FLIPPR essentially corresponds to the hierarchical structure of grammars. For example, to handle arithmetic expressions in FLIPPR, we can write:

```
ppr' i (Minus x y) =
  parensIf (i >= 6) $ group $
    ppr 5 x <> nest 2
      (line' <> text "-" <> space' <> ppr 6 y); .
```

FLIPPR will automatically expand the definition and derive a group of `ppr_i` functions indexed by the priority integer *i*, corresponding to the hierarchical grammar structure. In other words, there is no need to specify the concrete grammar, which is already implicitly embedded in the printer program. This makes FLIPPR programs neat and concise. Following this idea, BIYACC programs can also be made more concise: In a BIYACC program, the user is allowed to omit the production rules in the concrete syntax part (or omit the whole concrete syntax part), and they will be automatically generated by extracting the terminals and nonterminals in the right-hand sides of all actions. However, if these production rules are supplied, BIYACC will perform some sanity checks: It will make sure that, in an action group, the user has covered all of the production rules of the nonterminal appearing in the 'type declaration', and never uses undefined production rules.

<sup>10</sup> Although they use different implementation techniques, we will not dive into them in our related work. See Matsuda and Wang's related work for a comparison [36].

Just like basic BiYACC, all of the systems described above (aim to) handle unambiguous grammars only. Theoretically, when the user-defined grammar (or the derived grammar) is ambiguous, ISDs' partial isomorphism could guarantee Theorem 1' by returning Nothing on ambiguous input; FLIPPR's (own) Theorem 1 is comparable to Theorem 1' by taking all the language constructs which may cause non-injective printing into account. However, according to the paper, FLIPPR's Theorem 1 appears to only consider nondeterministic printing based on prettiness (layouts). Since the discussion on ambiguous grammars has not been presented in their papers, we tested their implementation and the behaviour is as follows: Neither ISDs nor FLIPPR will notify the user that the (derived) grammar is ambiguous at compile time. For ISDs, the right-to-left direction of our Theorem 1' will fail, while for FLIPPR, both directions will fail. (They never promise to handle ambiguous grammars, though.) In contrast, Brabrand et al. [7] give a detailed discussion about ambiguity detection, and XSugar statically checks if the transformations are 'reversible'. If any ambiguity in the program is detected, XSugar will notify the user of the precise location where ambiguity arises. In BiYACC, the ambiguity detection of the input grammar is performed by the employed parser generator (currently HAPPY), and the result is reported at compile time; if no warning is reported, the well-behavedness is always guaranteed. Note that the ambiguity detection can produce false positives: warnings only mean that the grammar is not LALR(1) but does not necessarily mean that the grammar is ambiguous—ambiguity detection is undecidable for the full CFG [10].

Here we also briefly discuss ambiguity detection for the filter approaches: Priority and associativity (bi-)filters can be applied to (LA)LR parse tables to resolve (shift/reduce) conflicts [9, 26, 52, 53], and thus the completeness for simple (bi-)filters (see Definition 9) on LALR(1) grammars can be statically checked. However, our implementation does not support it, for bi-filter directives are more general, as stated in the beginning of Section 5.4.2, and therefore cannot be transformed to the underlying parser generator's YACC-style directives. Finding a way to directly apply priority and associativity bi-filters to parse tables (generated by HAPPY) is left as future work.

Finally, we compare BiYACC with an industrial tool, AUGEAS, which provides the user with a local configuration API that converts configuration data into a rose tree representation [31]. Similar to BiYACC, AUGEAS also uses the idea of state-based asymmetric lenses so that its *parse* and *print* functions satisfy well-behavedness and it tries to preserve comments and layouts when printing the tree representation back. However, since the purpose of AUGEAS and BiYACC are different, the differences between the tools are also noticeable: (i) AUGEAS works for regular grammars while BiYACC works for (unambiguous) context-free grammars. (ii) AUGEAS uses a combinator-based approach while BiYACC adopts a generator-based approach. (iii) AUGEAS works more like a simple parser that stops after constructing CSTs: in the parsing direction, AUGEAS unambiguously separates strings into sub-strings, turn sub-strings into tokens, and use tokens to build the corresponding tree; but since each lens combinator (of AUGEAS) has its predefined strategy to turn its acceptable strings into the tree representation, the corresponding tree will be determined once the input string and the lens combinators for parsing the string are given; AUGEAS does not provide a functionality to further transform a tree. On the other hand, BiYACC first turns a string into its isomorphic CST (fully determined the input string and the

grammar description) and finally converts the CST to its AST in accordance with the algebraic data types defined by the user; that is, the relation between a string (CST) and its AST is not predetermined but can be adjusted by the user (through printing actions).

## 7.2 Generalised Parsing, Disambiguation, and Filters

The grammar of a programming language is usually designed to be unambiguous. Various parser-dependent disambiguation methods such as grammar transformation [29] and parse table conflicts elimination [23] have been developed to guide the parser to produce a single correct CST [26]. On the other hand, natural languages that are inherently ambiguous usually require their parsing algorithms to produce all the possible CSTs; this requirement gives rise to algorithms such as Earley [14] and generalised LR [50] (GLR for short). Although these parsing algorithms produce all the possible CSTs, both their time complexity and space complexity are reasonable. For instance, GLR runs in cubic time in the worst situation and in linear time if the grammar is ‘almost unambiguous’ [48].

The idea to relate generalised parsing with parser-independent disambiguation for programming languages is proposed by Klint and Visser [26]. They proposed two classes of filters, *property filters* (defined in terms of predicates on a single tree) and *comparison filters* (defined in terms of relations among trees), but we only adapt and bidirectionalise predicate filters in this paper. One difficulty lies in the fact that it is unclear how to define *repair* for comparison filters, as they generally select better trees rather than absolutely correct ones—in the printing direction, since *put* only produces a single CST, we do not know whether this CST needs repairing or not (for there is no other CST to compare). This is also one of the most important problems for our future work.

Parser-independent disambiguation (for handling priority and associativity conflicts) can also be found in LaLonde and des Rivieres’s [29] and Aasa’s [1] work. At first glance, our *repair* function is quite similar to LaLonde and des Rivieres’s post-parse *tree transformations* that bring a CST into an *expression tree*, on whose nodes additional restrictions of priority and associativity are imposed. To be simple (but not completely precise), a CST’s corresponding expression tree is obtained by first dropping all the nodes constructed from injective productions<sup>11</sup> (note that parentheses nodes are still kept) and then use a *precedence-introducing tree transformation* to reshape the result. The transformation will do ‘repairing’ by rotating all the adjacent nodes of the tree where priority or associativity constraint is violated. By contrast, our *repair* function is simpler and only introduces parentheses in places where the *judge* function returns False. In short, their tree transformations are a kind of parser-independent disambiguation which does not require generalised parsing; however, those tree transformations are (almost) not applicable in the printing direction if well-behavedness is taken into consideration (due to the rotation of CSTs). Furthermore, it

<sup>11</sup> An injective production, or a chain production, is one whose right-hand side is a single nonterminal; for instance,  $E \rightarrow N$ .

is not clear whether their approach can be generalised to handle other types of conflicts rather than the ones caused by priority and associativity.

There is much research on how to handle ambiguity in the parsing direction as discussed above; conversely, little research is conducted for ‘handling ambiguity in the printing direction’ and we find only one paper [8] that describes how to produce correct program text regarding priority and associativity, which is also one of the bases of our work. We extend their work [8] by allowing the bracket attribute to work with injective productions such as  $E \rightarrow T$ ;  $T \rightarrow F$ ;  $F \rightarrow '(' E ')'$  {# Bracket #};. (The previous work seems to only support the bracket attribute in the form of  $E \rightarrow '(' E ')'$  {# Bracket #};; whether the nonterminal  $E$  on the left-hand side and right-hand side can be different is not made clear.)

Finally, we compare our approach with the conventional ones in general. In history, a printer is believed to be much simpler than a parser and is usually developed independently (of its corresponding parser). While a few printers choose to produce parentheses at every occasion naively, most of them take disambiguation information (for example, from the language’s operator precedence table) into account and try to produce necessary parentheses only. However, as the YACC-style conventional disambiguation [23] is parser-dependent, this parentheses-adding technique is also printer-dependent. As the post-parse disambiguation increases the modularity of the (front-end of the) compiler [29], we believe that our post-print parentheses-adding increases the modularity once again. Additionally, the unification of disambiguation for both parsing and printing makes it possible for us to impose bi-filter laws, which further makes it possible to guarantee the well-behavedness of the whole system.

### 7.3 Comparison with a Get-based Approach

Our work is theoretically based on asymmetric lenses [17] of bidirectional transformations [11, 19], particularly taking inspiration from the recent progress on putback-based bidirectional programming [15, 27, 28, 40, 41]. As explained in Section 3, the purpose of bidirectional programming is to relieve the burden of thinking bidirectionally—the programmer writes a program in only one direction, and a program in the other direction is derived automatically. We call a language *get-based* when programs written in the language denote *get* functions, and call a language *putback-based* when its programs denote *put* functions. In the context of parsing and reflecting printing, the get-based approach lets the programmer describe a parser, whereas the putback-based approach lets the programmer describe a printer. Below we discuss in more depth how the putback-based methodology affects BiYACC’s design by comparing BiYACC with a closely related, get-based system.

Martins et al. [35] introduces an attribute grammar-based BX system for defining transformations between two representations of languages (two grammars). The utilisation is similar to BiYACC: The programmer defines both grammars and a set of rules specifying a *forward* transformation (i.e. *get*), with a backward transformation (i.e. *put*) being automatically generated. For example, the BiYACC actions in lines

28–30 of Figure 2 can be expressed in Martins et al.’s system as

$$\begin{aligned} get_A^E(plus(x, '+', y)) &\rightarrow add(get_A^E(x), get_A^T(y)) \\ get_A^E(minus(x, '-', y)) &\rightarrow sub(get_A^E(x), get_A^T(y)) \\ get_A^E(et(e)) &\rightarrow get_A^T(e) \end{aligned}$$

which describes how to convert certain forms of CSTs to corresponding ASTs. The similarity is evident, and raises the question as to how get-based and putback-based approaches differ in the context of parsing and reflective printing.

The difference lies in the fact that, with a get-based system, certain decisions on the backward transformation are, by design, permanently encoded in the bidirectionalisation system and cannot be controlled by the user, whereas a putback-based system can give the user fuller control. For example, when no source is given and more than one rules can be applied, Martins et al.’s system chooses, by design, the one that creates the most specialised version. This might or might not be ideal for the user of the system. For example: Suppose that we port to Martins et al.’s system the BIYACC action that relates TIGER’s concrete ‘&’ operator with a specialised abstract *if* expression in Section 6.2, coexisting with a more general rule that maps a concrete *if* expression to an abstract *if* expression. Then printing the AST  $\tau_{Cond}(\text{TSV "a"}) (\text{TSV "b"}) \emptyset$  from scratch will and can only produce *a & b*, as dictated by the system’s hard-wired printing logic. By contrast, the user of BIYACC can easily choose to print the AST from scratch as *a & b* or *if a then b else  $\emptyset$*  by suitably ordering the printing actions.

This difference is somewhat subtle, and one might argue that Martins et al.’s design simply went one step too far—if their system had been designed to respect the rule ordering as specified by the user, as opposed to always choosing the most specialised rule, the system would have given its user the same flexibility as BIYACC. Interestingly, whether to let user-specified rule/action ordering affect the system’s behaviour is, in this case, exactly the line between get-based and putback-based design. The user of Martins et al.’s system writes rules to specify a forward transformation, whose semantics is the same regardless of how the rules are ordered, and thus it would be unpleasantly surprising if the rule ordering turned out to affect the system’s behaviour. By contrast, the user of BIYACC only needs to think in one direction about the printing behaviour, for which it is natural to consider how the actions should be ordered when an AST has many corresponding CSTs; the parsing behaviour will then be automatically and uniquely determined. In short, relevance of action ordering is incompatible with get-based design, but is a natural consequence of putback-based thinking.

## 8 Conclusion

We conclude the paper by summarising our contributions:

- We have presented the design and implementation of BIYACC, with which the programmer can describe both a parser and a reflective printer for a fully disambiguated context-free grammar in a single program. Our solution guarantees the partial version of the consistency properties (Definition 2) by construction.



- We proposed the notion of bi-filters, which enables BIYACC to disambiguate ambiguous grammars while still respecting the consistency properties. This is the main new contribution compared to the previous SLE'16 version [55].
- We have demonstrated that BIYACC can support various tasks of language engineering, from traditional constructions of basic machinery such as printers and parsers to more complex tasks such as resugaring, simple refactoring, and language evolution.

**Acknowledgements** We thank the reviewers and the editor for their selflessness and effort spent on reviewing our paper, a quite long one. With their help, the readability of the paper is much improved, especially regarding how several case studies are structured, how theorems for the basic BIYACC and theorems for the extended version handling ambiguous grammars are related, and how look-alike notions are 'disambiguated'. This work is partially supported by the Japan Society for the Promotion of Science (JSPS) Grant-in-Aid for Scientific Research (S) No. 17H06099; in particular, most of the second author's contributions were made when he worked at the National Institute of Informatics and funded by the grant.

## References

1. Aasa A (1995) Precedences in specifications and implementations of programming languages. In: *Selected Papers of the Symposium on Programming Language Implementation and Logic Programming*, Elsevier Science Publishers B. V., Amsterdam, The Netherlands, PLILP '91, pp 3–26, URL <http://dl.acm.org/citation.cfm?id=203429.203431>
2. Afroozeh A, Izmaylova A (2015) Faster, practical GLL parsing. In: Franke B (ed) *Compiler Construction*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 89–108, URL [https://doi.org/10.1007/978-3-662-46663-6\\_5](https://doi.org/10.1007/978-3-662-46663-6_5)
3. Aho AV, Johnson SC, Ullman JD (1975) Deterministic parsing of ambiguous grammars. *Communications of the ACM* 18(8):441–452
4. Appel AW (1998) *Modern Compiler Implementation in ML*, First Edition. Cambridge University Press, New York, NY, USA
5. Bird R (2014) *Thinking Functionally with Haskell*. Cambridge University Press, Cambridge, UK, URL <https://doi.org/10.1017/CB09781316092415>
6. Boulton R (1966) Syn: A single language for specifying abstract syntax trees, lexical analysis, parsing and pretty-printing. Tech. Rep. Number 390, Computer Laboratory, University of Cambridge
7. Brabrand C, Møller A, Schwartzbach MI (2008) Dual syntax for XML languages. *Information Systems* 33(4-5):385–406, URL <https://doi.org/10.1016/j.is.2008.01.006>
8. van den Brand M, Visser E (1996) Generation of formatters for context-free languages. *ACM Transactions on Software Engineering and Methodology* 5(1):1–41, URL <https://doi.org/10.1145/226155.226156>
9. van den Brand MGJ, Scheerder J, Vinju JJ, Visser E (2002) Disambiguation filters for scannerless generalized LR parsers. In: *Proceedings of the 11th International Conference on Compiler Construction*, Springer-Verlag, London, UK, CC '02, pp 143–158, URL [https://doi.org/10.1007/3-540-45937-5\\_12](https://doi.org/10.1007/3-540-45937-5_12)
10. Cantor DG (1962) On the ambiguity problem of Backus systems. *Journal of the ACM* 9(4):477–479
11. Czarnecki K, Foster JN, Hu Z, Lämmel R, Schürr A, Terwilliger JF (2009) Bidirectional transformations: A cross-discipline perspective. In: *Proceedings of the 2nd International Conference on Theory and Practice of Model Transformations*, Springer-Verlag, Berlin, Heidelberg, ICMT '09, pp 260–283, URL [https://doi.org/10.1007/978-3-642-02408-5\\_19](https://doi.org/10.1007/978-3-642-02408-5_19)
12. Dijkstra EW (1975) Guarded commands, nondeterminacy and formal derivation of programs. *Communications of the ACM* 18(8):453–457, URL <https://doi.org/10.1145/360933.360975>
13. Duregård J, Jansson P (2011) Embedded parser generators. In: *Proceedings of the 4th ACM Symposium on Haskell*, ACM, New York, NY, USA, Haskell '11, pp 107–117, URL <https://doi.org/10.1145/2034675.2034689>
14. Earley J (1970) An efficient context-free parsing algorithm. *Communications of the ACM* 13(2):94–102, URL <https://doi.org/10.1145/362007.362035>
15. Fischer S, Hu Z, Pacheco H (2015) The essence of bidirectional programming. *Science China Information Sciences* 58(5):1–21

16. Foster JN (2009) Bidirectional programming languages. PhD thesis, University of Pennsylvania
17. Foster JN, Greenwald MB, Moore JT, Pierce BC, Schmitt A (2007) Combinators for bidirectional tree transformations: A linguistic approach to the view-update problem. *ACM Transactions on Programming Language Systems* 29(3), URL <https://doi.org/10.1145/1232420.1232424>
18. Fowler M, Beck K (1999) Refactoring: Improving the Design of Existing Code. Addison-Wesley Professional, Boston, MA, USA
19. Gibbons J, Stevens P (eds) (2018) International Summer School on Bidirectional Transformations (Oxford, UK, 25–29 July 2016), *Lecture Notes in Computer Science*, vol 9715. Springer
20. Gosling J, Joy B, Steele G (2006) The Java language specification, third edition. URL <https://docs.oracle.com/javase/specs/>
21. Hirzel M, Rose KH (2013) Tiger language specification. URL <https://cs.nyu.edu/courses/fall13/CSCI-GA.2130-001/tiger-spec.pdf>
22. Hu Z, Ko HS (2018) Principles and practice of bidirectional programming in BiGUL. In: Gibbons J, Stevens P (eds) Bidirectional Transformations: International Summer School, Oxford, UK, July 25–29, 2016, Tutorial Lectures, Springer International Publishing, Cham, pp 100–150, URL [https://doi.org/10.1007/978-3-319-79108-1\\_4](https://doi.org/10.1007/978-3-319-79108-1_4)
23. Johnson SC (1975) Yacc: Yet another compiler-compiler. AT&T Bell Laboratories Technical Reports (AT&T Bell Laboratories Murray Hill, New Jersey 07974) (32)
24. Kernighan BW, Ritchie DM (1988) The C Programming Language. Prentice Hall Press, Upper Saddle River, NJ, USA
25. Kinoshita D, Nakano K (2017) Bidirectional certified programming. In: Eramo R, Johnson M (eds) Proceedings of the 6th International Workshop on Bidirectional Transformations co-located with The European Joint Conferences on Theory and Practice of Software (ETAPS 2017), CEUR Workshop Proceedings, Uppsala, Sweden, vol 1827, pp 31–38
26. Klint P, Visser E (1994) Using filters for the disambiguation of context-free grammars. In: Pighizzini G, Pietro PS (eds) Proceedings of the ASMICS Workshop on Parsing Theory, University of Milan, Italy, Milano, Italy, pp 1–20
27. Ko HS, Hu Z (2018) An axiomatic basis for bidirectional programming. Proceedings of the ACM on Programming Languages 2(POPL):41:1–41:29, URL <https://doi.org/10.1145/3158129>
28. Ko HS, Zan T, Hu Z (2016) BiGUL: A formally verified core language for putback-based bidirectional programming. In: Proceedings of the 2016 ACM SIGPLAN Workshop on Partial Evaluation and Program Manipulation, ACM, New York, NY, USA, PEPM '16, pp 61–72, URL <https://doi.org/10.1145/2847538.2847544>
29. LaLonde WR, des Rivieres J (1981) Handling operator precedence in arithmetic expressions with tree transformations. *ACM Transactions on Programming Language Systems* 3(1):83–103, URL <https://doi.org/10.1145/357121.357127>
30. Lämmel R, Jones SP (2003) Scrap your boilerplate: A practical design pattern for generic programming. In: Proceedings of the 2003 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation, ACM,

- New York, NY, USA, TLDI '03, pp 26–37, URL <https://doi.org/10.1145/604174.604179>
31. Lutterkort D (2008) Augeas—a configuration API. In: Proceedings of the Ottawa Linux Symposium, Ottawa, Canada, pp 47–56
  32. Macedo N, Pacheco H, Cunha A, Oliveira JN (2013) Composing least-change lenses. Proceedings of the Second International Workshop on Bidirectional Transformations 57:1–19, URL <https://doi.org/10.14279/tuj.eceasst.57.868>
  33. Marlow S, Gill A (2001) The parser generator for Haskell. URL <https://www.haskell.org/happy/>
  34. Marlow S, et al. (2010) Haskell 2010 language report. URL <https://www.haskell.org/onlinereport/haskell2010/>
  35. Martins P, Saraiva J, Fernandes JP, Van Wyk E (2014) Generating attribute grammar-based bidirectional transformations from rewrite rules. In: Proceedings of the ACM SIGPLAN 2014 Workshop on Partial Evaluation and Program Manipulation, ACM, New York, NY, USA, PEPM '14, pp 63–70, URL <https://doi.org/10.1145/2543728.2543745>
  36. Matsuda K, Wang M (2018) Embedding invertible languages with binders: A case of the FliPpr language. In: Proceedings of the 11th ACM SIGPLAN International Symposium on Haskell, ACM, New York, NY, USA, Haskell 2018, pp 158–171, URL <http://doi.org/10.1145/3242744.3242758>
  37. Matsuda K, Wang M (2018) FliPpr: A system for deriving parsers from pretty-printers. New Generation Computing 36(3):173–202, URL <https://doi.org/10.1007/s00354-018-0033-7>
  38. Matsuda K, Mu SC, Hu Z, Takeichi M (2010) A grammar-based approach to invertible programs. In: Gordon AD (ed) Proceedings of the 19th European Conference on Programming Languages and Systems, Springer-Verlag Berlin Heidelberg, Berlin, Heidelberg, no. 20 in ESOP'10, pp 448–467, DOI [https://doi.org/10.1007/978-3-642-11957-6\\_24](https://doi.org/10.1007/978-3-642-11957-6_24)
  39. Norell U (2007) Towards a practical programming language based on dependent type theory. PhD thesis, Chalmers University of Technology
  40. Pacheco H, Hu Z, Fischer S (2014) Monadic combinators for “putback” style bidirectional programming. In: Proceedings of the ACM SIGPLAN 2014 Workshop on Partial Evaluation and Program Manipulation, ACM, New York, NY, USA, PEPM '14, pp 39–50, URL <https://doi.org/10.1145/2543728.2543737>
  41. Pacheco H, Zan T, Hu Z (2014) BiFluX: A bidirectional functional update language for XML. In: Proceedings of the 16th International Symposium on Principles and Practice of Declarative Programming, ACM, New York, NY, USA, PPDP '14, pp 147–158, URL <https://doi.org/10.1145/2643135.2643141>
  42. Pombrio J, Krishnamurthi S (2014) Resugaring: Lifting evaluation sequences through syntactic sugar. In: Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation, ACM, New York, NY, USA, no. 6 in PLDI '14, pp 361–371, URL <https://doi.org/10.1145/2594291.2594319>
  43. Pombrio J, Krishnamurthi S (2015) Hygienic resugaring of compositional desugaring. In: Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming, ACM, New York, NY, USA, no. 13 in ICFP 2015, pp

- 75–87, URL <https://doi.org/10.1145/2784731.2784755>
44. Rendel T, Ostermann K (2010) Invertible syntax descriptions: Unifying parsing and pretty printing. In: *Proceedings of the Third ACM Haskell Symposium on Haskell*, ACM, New York, NY, USA, Haskell '10, pp 1–12, URL <https://doi.org/10.1145/1863523.1863525>
  45. Reps T, Teitelbaum T (1984) The synthesizer generator. In: *Proceedings of the First ACM SIGSOFT/SIGPLAN Software Engineering Symposium on Practical Software Development Environments*, ACM, New York, NY, USA, SDE 1, pp 42–48, URL <https://doi.org/10.1145/800020.808247>
  46. Reps T, Teitelbaum T, Demers A (1983) Incremental context-dependent analysis for language-based editors. *ACM Transactions on Programming Language Systems* 5(3):449–477, URL <https://doi.org/10.1145/2166.357218>
  47. Scott E, Johnstone A (2010) GLL parsing. *Electronic Notes in Theoretical Computer Science* 253(7):177–189, URL <https://doi.org/10.1016/j.entcs.2010.08.041>
  48. Scott E, Johnstone A, Economopoulos R (2007) BRNGLR: A cubic tomita-style glr parsing algorithm. *Acta Informatica* 44(6):427–461, URL <https://doi.org/10.1007/s00236-007-0054-z>
  49. Sheard T, Jones SP (2002) Template meta-programming for Haskell. In: *Proceedings of the 2002 ACM SIGPLAN Workshop on Haskell*, ACM, New York, NY, USA, Haskell '02, pp 1–16, URL <https://doi.org/10.1145/581690.581691>
  50. Tomita M (1985) An efficient context-free parsing algorithm for natural languages. In: *Proceedings of the 9th International Joint Conference on Artificial Intelligence - Volume 2*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, IJCAI'85, pp 756–764, URL <http://dl.acm.org/citation.cfm?id=1623611.1623625>
  51. Traver VJ (2010) On compiler error messages: What they say and what they mean. *Advances in Human-Computer Interaction* 2010:3:1–3:26, URL <https://doi.org/10.1155/2010/602570>
  52. Visser E (1997) A case study in optimizing parsing schemata by disambiguation filters. In: *International Workshop on Parsing Technology (IWPT 1997)*, Massachusetts Institute of Technology, Boston, USA, pp 210–224
  53. Visser E (1997) Syntax definition for language prototyping. PhD thesis, University of Amsterdam
  54. Younger DH (1967) Recognition and parsing of context-free languages in time  $n^3$ . *Information and Control* 10(2):189 – 208, URL <http://www.sciencedirect.com/science/article/pii/S00199586780007X>
  55. Zhu Z, Zhang Y, Ko HS, Martins P, Saraiva J, Hu Z (2016) Parsing and reflective printing, bidirectionally. In: *Proceedings of the 2016 ACM SIGPLAN International Conference on Software Language Engineering*, ACM, New York, NY, USA, SLE 2016, pp 2–14, URL <https://doi.org/10.1145/2997364.2997369>