



# 1 Unifying Parsing and Reflective Printing for Fully 2 Disambiguated Grammars

3 Zirun Zhu<sup>1</sup> · Hsiang-Shang Ko<sup>2</sup> · Yongzhe Zhang<sup>1</sup> · Pedro Martins<sup>3</sup> ·  
4 João Saraiva<sup>4</sup> · Zhenjiang Hu<sup>5</sup>

5 Received: 31 January 2019 / Accepted: 16 December 2019  
6 © Ohmsha, Ltd. and Springer Japan KK, part of Springer Nature 2020

## 7 Abstract

8 Language designers usually need to implement parsers and printers. Despite being  
9 two closely related programs, in practice they are often designed separately, and  
10 then need to be revised and kept consistent as the language evolves. It will be more  
11 convenient if the parser and printer can be unified and developed in a single pro-  
12 gram, with their consistency guaranteed automatically. Furthermore, in certain  
13 scenarios (like showing compiler optimisation results to the programmer), it is  
14 desirable to have a more powerful reflective printer that, when an abstract syntax  
15 tree corresponding to a piece of program text is modified, can propagate the  
16 modification to the program text while preserving layouts, comments, and syntactic  
17 sugar. To address these needs, we propose a domain-specific language BiYACC,  
18 whose programs denote both a parser and a reflective printer for a fully disam-  
19 biguated context-free grammar. BiYACC is based on the theory of bidirectional  
20 transformations, which helps to guarantee by construction that the generated pairs of  
21 parsers and reflective printers are consistent. Handling grammatical ambiguity is  
22 particularly challenging: we propose an approach based on generalised parsing and  
23 disambiguation filters, which produce all the parse results and (try to) select the only  
24 correct one in the parsing direction; the filters are carefully bidirectionalised so that  
25 they also work in the printing direction and do not break the consistency between  
26 the parsers and reflective printers. We show that BiYACC is capable of facilitating  
27 many tasks such as Pombrio and Krishnamurthi's 'resugaring', simple refactoring,  
28 and language evolution.

29  
30 **Keywords** Asymmetric lenses · Disambiguation filters · Bidirectional  
31 transformations · Domain-specific languages · Parsing · Reflective  
32 printing  
33

34

**Electronic supplementary material** The online version of this article (<https://doi.org/10.1007/s00354-019-00082-y>) contains supplementary material, which is available to authorized users.

Extended author information available on the last page of the article

## 35 Introduction

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

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

55 Despite their advantages, these domain-specific languages cannot deal with  
 56 synchronisation between program text and ASTs. Let us look at a concrete example  
 57 in Fig. 1: the original program text is an arithmetic expression, containing a  
 58 negation, a comment, and parentheses (one pair of which is redundant). It is first  
 59 parsed to an AST (supposing that addition is left-associative) where the negation is  
 60 desugared to a subtraction, parentheses are implicitly represented by the tree  
 61 structure, and the comment is thrown away. Suppose that the AST is optimised by  
 62 replacing `Add (Num 1) (Num 1)` with a constant `Num 2`. The user may want to  
 63 observe the optimisation made by the compiler, but the AST is an internal  
 64 representation not exposed to the user, so a natural idea is to propagate the changes  
 65 on the AST back to the program text to make it easy for the user to check where the  
 66 changes are. With a conventional printer, however, the printed result will likely  
 67 mislead the programmer into thinking that the negation is replaced by a subtraction  
 68 by the compiler; also, since the comment is not preserved, it will be harder for the

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

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

69 programmer to compare the updated and original versions of the text. The problem  
 70 illustrated here has also been investigated in many other practical scenarios where  
 71 the parser and printer are used as a bridge between the system and the user, for  
 72 example,

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

80  
 81 To address the problem, we propose a domain-specific language BiYACC, which  
 82 enables the user to describe both a parser and a reflective printer for a fully  
 83 disambiguated context-free grammar (CFG) in a single program. Different from a  
 84 conventional printer, a reflective printer takes a piece of program text and an AST,  
 85 which is usually slightly modified from the AST corresponding to the original  
 86 program text, and propagates the modification back to the program text. Meanwhile  
 87 the comments (and layouts) in the unmodified parts of the program text are all  
 88 preserved. This can be seen clearly from the result of using our reflective printer on  
 89 the above arithmetic expression example in Fig. 1. It is worth noting that reflective  
 90 printing is a generalisation of the conventional notion of printing, because a  
 91 reflective printer can accept an AST and an empty piece of program text, in which  
 92 case it will behave just like a conventional printer, producing a new piece of  
 93 program text depending on the AST only.

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

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

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

$$print s (parse s) = s . \quad (2)$$

103

IFL01 <sup>1</sup> We assume basic knowledge about functional programming languages and their notations, in particular  
 IFL02 HASKELL [5, 34]. In HASKELL, an argument of function application does not need to be enclosed in (round  
 IFL03 parentheses, i.e. we write  $fx$  instead of  $f(x)$ ; type variables are implicitly universally quantified,  
 IFL04 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  
 IFL05 universal quantification for free variables in an equation; for instance,  $parse (print s t) = t$  is in fact  
 IFL06  $\forall s t. parse (print s t) = t$ .

104 These two properties are inspired by the theory of bidirectional transformations  
 105 [19], in particular lenses [17], and are guaranteed by construction for all BiYACC  
 106 programs.

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

126 This is the extended version of our previous work Parsing and Reflective  
 127 Printing, Bidirectionally presented at SLE'16 [55], and the differences are mainly as  
 128 follows: (1) we propose the notion of bidirectionalised filters and integrate them into  
 129 BiYACC for handling grammatical ambiguity (Sect. 5); the related work section is  
 130 also updated accordingly. (2) We restructure the narration for introducing the basic  
 131 BiYACC system and in particular elaborate on the isomorphism between program  
 132 text and CSTs. (3) We present the definitions and theorems in a more formal way,  
 133 and complete their proofs. (4) We make several other revisions such as renewing the  
 134 figures for introducing the BiYACC system and the syntax of BiYACC programs.

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

## 137 A First Look at BiYACC

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

## 142 Syntax Definitions

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

```

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
8
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

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

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

159 *Directives* The *#Directives* part defines the syntax of comments and disambiguation  
 160 directives. For example, line 23 shows that the syntax for single line comments  
 161 is  $//$ ,<sup>2</sup> while line 24 states that  $/*$  and  $*/$  are, respectively, the beginning mark  
 162 and ending mark for block comments. Since the grammar for arithmetic expressions  
 163 is unambiguous, there is no need to give any disambiguation directive for this  
 164 example (whereas the ambiguous version of the grammar in Fig. 6 needs to be  
 165 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.

166 **Printing Actions**

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

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

183 Let us look at a specific action—the first one for the expression example, at line  
 184 28 of Fig. 2:

$$\text{Add } x \ y \ +> \ [x \ +> \ \text{Expr}] \ '+' \ [y \ +> \ \text{Term}];$$

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

194 Having explained what an action means, we can now explain the semantics of the  
 195 entire program. Given an AST and a CST as input, first a group (of actions) is  
 196 chosen according to the types of the trees. Then, the actions in the group are tried in  
 197 order, from top to bottom, by performing simultaneous pattern matching on both  
 198 trees. If pattern matching for an action succeeds, the updating operations specified  
 199 by the action is executed, otherwise the next action is tried. Execution of the  
 200 program ends when the matched action specifies either no updating operations or  
 201 only updates to primitive data types such as Numeric. BiYACC’s most interesting  
 202 behaviour shows up when all actions in the chosen group fail to match—in this case  
 203 a suitable CST will be created. The specific approach adopted by BiYACC is to  
 204 perform pattern matching on the AST only and choose the first matched action. A  
 205 suitable CST conforming to the CST-side pattern is then created, and after that the  
 206 whole group of actions is tried again. This time the pattern matching will succeed at  
 207 the action used to create the CST, and the program will be able to make further  
 208 progress. For instance, assuming that the source is  $1 * 2$  while the view is Add (Num



209 1) (Num 2), a new source skeleton representing  $- + -$  will be created and the  $-$  part  
 210 will be updated recursively later. We will elaborate more on this in Sect. 4.

211 *Deep patterns* Using deep patterns, we can write actions that establish nontrivial  
 212 relationships between CSTs and ASTs. For example, the action at line 38 of Fig. 2  
 213 associates abstract subtraction expressions whose left operand is zero with concrete  
 214 negated expressions; this action is the key to preserving negated expressions in the  
 215 CST. For an example of a more complex CST-side pattern: suppose that we want to  
 216 write a pattern that matches those CSTs produced by the rule  $\text{Factor} \rightarrow \text{'-' Factor}$ ,  
 217 where the inner nonterminal  $\text{Factor}$  produces a further  $\text{'-' Factor}$  using the same rule.  
 218 This pattern is written by overlaying the production rule on the first nonterminal  
 219  $\text{Factor}$  (an additional pair of parentheses is required for the expanded nonterminal):  
 220  $\text{'-' (Factor} \rightarrow \text{'-' Factor)}$ . More examples involving this kind of deep patterns can  
 221 be found in Sect. 6.

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

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

## 245 Foundation of BiYACC: Putback-Based Bidirectional Programming

246 From a BiYACC program, in addition to generating a parser and a printer, we also  
 247 need to guarantee that the two generated programs are consistent with each other,  
 248 i.e. satisfy the properties (1) and (2) stated in Sect. 1. It is possible to implement the  
 249 *print* and *parse* semantics separately in an ad hoc way, but verifying the two  
 250 consistency properties takes extra effort. The implementation we present, however,  
 251 is systematic and guarantees consistency by construction, thanks to the well-

252 developed theory of bidirectional transformations (BXs for short), in particular  
 253 lenses [17]. We will give a brief introduction to BXs below; for a comprehensive  
 254 treatment, the readers are referred to the lecture notes for the 2016 Oxford Summer  
 255 School on Bidirectional Transformations [19].

## 256 Parsing and Printing as Lenses

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

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

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

266 satisfying the well-behavedness laws:

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

268  
 269 Intuitively, a *get* function extracts a part of a source of interest to the user as a  
 270 view, and a *put* function takes a source and a view and produces an updated source  
 271 incorporating information from the view. Partiality is explicitly represented by  
 272 making the functions return *Maybe* values: a *get* or *put* function returns *Just*  $r$  where  
 273  $r$  is the result, or *Nothing* if the input is not in the domain. The PUTGET law enforces  
 274 that *put* must embed all information of the view into the updated source, so the view  
 275 can be recovered from the source by *get*, while the GETPUT law prohibits *put* from  
 276 performing unnecessary updates by requiring that putting back a view directly  
 277 extracted from a source by *get* must produce the same, unmodified source.

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

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

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

283

284 **Putback-Based Bidirectional Programming in BiGUL**

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

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

306 *Replace* The simplest BiGUL operation we use is

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

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

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

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

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

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

321 matches the source with a tuple pattern  $(x, \_)$  and the view with a variable pattern  $x$ ,  
 322 so that the first component of the source tuple is related with the whole view; during  
 323 the update, the first component of the source is replaced by the whole view, as  
 324 indicated by the operation  $x = \text{Replace}$ . (The part marked by underscore  $(\_)$  simply



325 means that it will be skipped during the update.) Given a source (1,2) and a view 3,  
 326 the operation will produce (3,2) as the updated source. In general, any (type-correct)  
 327 BiGUL program can be used in the list of further updates, not just the primitive  
 328 **Replace**.

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

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

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

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

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

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

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

## 355 The Basic BiYACC

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

'#Abstract' decls '#Concrete' pgs '#Directives' drctvs '#Actions' ags , (3)

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

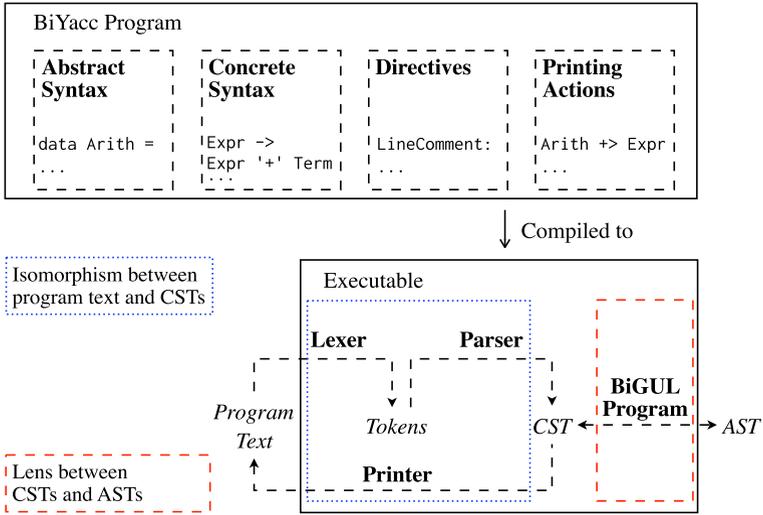


Fig. 3 Architecture of BiYACC

```

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 Sect. 5)

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

- 365 • The abstract syntax part (*decls* for HASKELL data type declarations) is already  
366 valid HASKELL code and is (almost) directly used as the definitions of AST data  
367 types.
- 368 • The concrete syntax part (*pgs* for production groups) is translated to definitions  
369 of CST data types (whose elements are representations of how a string is  
370 produced using the production rules), and also used to generate the pair of  
371 concrete parser (including a lexer) and printer for the conversion between  
372 program text and CSTs. This pair of concrete parser and printer can be shown to  
373 form an (partial) isomorphism (which will be defined in Sect. 4.1). This part will  
374 be explained in Sect. 4.2.
- 375 • The directives part (*drctvs* for directives) is used in the lexer for recognising  
376 single line and multi-line comments.
- 377 • The printing actions part (*ags* for action groups) is translated to a BIGUL  
378 program (which is a lens, see Definition 1) for handling (the semantic part of)  
379 parsing and reflective printing between CSTs and ASTs. This part will be  
380 explained in Sect. 4.3.

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

### 383 Composition of Isomorphisms and Lenses

384 First, we give the definition of (partial) isomorphisms.

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

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

388 such that the inverse properties hold:

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

389

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

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

395 where

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

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

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

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

## 406 The Concrete Parsing and Printing Isomorphism

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

## 410 Generating CST Data Types and Concrete Printers

411 The production rules in a context-free grammar dictate how to produce strings from  
412 nonterminals, and a CST can be regarded as encoding one particular way of  
413 producing a string using the production rules. In BiYACC, we represent CSTs starting  
414 from a nonterminal  $nt$  as an automatically generated HASKELL data type named  $nt$ ,  
415 whose constructors represent the production rules for  $nt$ . For each of these data  
416 types, we also generate a printing function which takes a CST as input and produces  
417 a string as dictated by the production rules in the CST.

418 For instance, in Fig. 2, the group of production rules from the nonterminal Factor  
419 (lines 18–21) is translated to the following HASKELL data type and concrete printing  
420 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 = ""
```

421 where Factor1 ... Factor4 are constructors corresponding to the four production rules,  
 422 and FactorNull represents an empty CST of type *Factor* and is used as the default  
 423 value whenever we want to create new program text depending on the view only. As  
 424 an example, Factor1 represents the production rule  $\text{Factor} \rightarrow \text{'-'} \text{Factor}$ , and its  
 425 String field stores the whitespaces appearing after a negation sign in the program text.  
 426 The Factor3 case makes a call to `cprtExpr::Expr → String`, which is the printing  
 427 function generated for the nonterminal Expr.

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

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

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

## 446 Generating Concrete Lexers and Parsers

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

451 *Lexer* Apart from handling the terminal symbols appearing in a grammar, the  
 452 lexer automatically derived by BiYACC can also recognise several kinds of literals,  
 453 including integers, strings, and identifiers, respectively, produced by the nontermi-

3FL01 <sup>3</sup> The reason for storing primitives in the String type is because String is the most precise representation  
 3FL02 that will not cause the loss of any information. For instance, this is useful for retaining the leading zeros of  
 3FL03 an integer such as 073. Storing 073 as Integer will cause the loss of the leading zero.

456 nals Numeric, String, and Identifier. For now, the forms of these literals are  
 457 predefined, but we take this as a step towards a lexerless grammar, in which strings  
 458 produced by nonterminals can be specified in terms of regular expressions.  
 459 Furthermore, whitespaces and comments are carefully handled in the derived lexer,  
 460 so they can be completely stored in CSTs and correctly recovered to the program  
 461 text in printing. This feature of BiYACC, which we explain below, makes layout  
 462 preservation transparent to the programmer.

463 An assumption of BiYACC is that whitespaces are only regarded as separators  
 464 between other tokens. (Although there exist some languages such as HASKELL and  
 465 PYTHON where indentation does affect the meaning of a program, there are  
 466 workarounds, e.g. writing a preprocessing program to insert explicit separators.)  
 467 Usually, token separators are thrown away in the lexing phase, but since we want to  
 468 keep layout information in CSTs, which are built by the parser, the lexer should  
 469 leave the separators intact and pass them to the parser. The specific approach taken  
 470 by BiYACC is wrapping a lexeme and the whitespaces following it into a single  
 471 token. Beginning whitespaces are treated separately from lexing and parsing, and  
 472 are always preserved. And in this prototype implementation, comments are also  
 473 regarded as whitespaces.

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

<pre>Factor -&gt; '-' Factor      Numeric      Identifier      '(' Expr ')';</pre>	<pre>Factor : token1 Factor { Factor1 \$1 \$2 }   tokenNumeric { Factor2 \$1 }   tokenIdentifier { Factor3 \$1 }   token2 Expr token3 { Factor4 \$1 \$2 \$3 } .</pre>
--	---

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

## 489 Inverse Properties

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



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

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

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

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

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

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

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

508 **Definition 5** (Parser correctness) A parser *cparse* is correct with respect to a printer  
509 *cprint* exactly when

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

$$511 \quad cprint \text{ cst} = \text{text} \quad \Rightarrow \quad \exists cst'. \text{ cparse } \text{ text} = \text{Just } cst' . \quad (5)$$

512  
513 To see what (4) means, recall that our CSTs, as described in Sect. 4.2.1, encode  
514 precisely the derivation trees, with the CST constructors representing the production  
515 rules used, and *cprint* traverses the CSTs and follows the encoded production rules  
516 to produce the derived program text. Now consider what *cparse* is supposed to do: it  
517 should take a piece of program text and find a derivation tree for it, i.e. a CST which  
518 *prints* to that piece of program text. This statement is exactly (4). In other words,  
519 (4) is the functional specification of parsing, which is satisfied if the parser generator  
520 we use behaves correctly. Also it is reasonable to expect that a parser will be able to  
521 successfully parse any valid program text, and this is exactly (5).

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

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

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

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

529



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

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

537 □

## 538 Generating the BiGUL Lens

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

543 *Action groups* Each group of actions is translated into a small BiGUL program,  
 544 whose name is determined by the view type  $vt$  and source type  $st$  and denoted by  
 545  $PROG(vt, st)$ . The BiGUL program has one single Case statement, and each action is  
 546 translated into two branches in this Case statement, one normal and the other  
 547 adaptive. All the adaptive branches are gathered in the second half of the Case  
 548 statement, so that the normal branches will be tried first. For example, the third  
 549 group of type  $Arith \rightarrow Factor$  is compiled to

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

552 *Normal branches* We said in Sect. 2 that the semantics of an action is to perform  
 553 pattern matching on both the source and view, and then update parts of the source  
 554 with parts of the view. This semantics is implemented with a normal branch: the  
 555 source and view patterns are compiled to the main condition, and, together with the  
 556 updates overlaid on the source pattern, also to an update operation. For example, the  
 557 first action in the  $Arith \rightarrow Factor$  group

```

[[vt '+>' st acts]]ActionGroup =
  PROG(vt, st) ':: 'BiGUL' st vt
  PROG(vt, st) '=
    'Case' '[' <[[a]]N,vt,stAction ' ; ' | a ∈ acts > <[[a]]A,stAction | a ∈ acts > { ' ; ' } ' ]'
[[vpat '+>' updates]]N,vt,stAction =
  '$(normalSV
    '[p]' SRCCOND(ERSVARS('[' st '->' updates '])Update) ' | ]' '[p]' vpat ' | ]'
    '[p]' SRCCOND(ERSVARS('[' st '->' updates '])Update) ' | ]'
    '$(update' '[p]' REMOVEAS(vpat) ' | ]'
      '[p]' SRCPAT('[' st '->' updates '])Update ' | ]'
      '[d]' <[[u]]vt,vpatUpdate | u ∈ updates > ' | ]'
  )'
[[ '[' var '+>' ucPrimitive ' ] ] ]vt,vpatUpdate = var '= Replace;'
[[ '[' var '+>' ucNonterminal ' ] ] ]vt,vpatUpdate = var '= PROG(VARTYPE(vt, vpat, var), uc) ' ; '
[[ '[' var '+>' ' ( ' nt '->' ... ' ) ' ] ] ]vt,vpatUpdate = [[ '[' var '+>' nt ' ] ] ]vt,vpatUpdate
[[ ' ( ' ... '->' updates ' ) ' ] ]vt,vpatUpdate = <[[u]]vt,vpatUpdate ; ' | u ∈ updates >
[[symbol]]vt,vpatUpdate = ''
[[vpat '+>' updates]]A,stAction = '$(adaptiveSV' '[p] _ | ]' '[p]' 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)

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

559 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; | ] ) .

```

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

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

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

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

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

```
entrance = bigulArithExpr
```

587 which is invoked in the main program.

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

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

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

593

## 594 Handling Grammatical Ambiguity

595 In Sect. 4, we have described the basic version of BiYACC, about which there is an  
 596 important assumption (stated in Theorem 1) that grammars have to be unambigu-  
 597 ous. Having this assumption can be rather inconvenient in practice, however, as

```

#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)

598 ambiguous grammars (with disambiguation directives) are often preferred since  
 599 they are considered more natural and human friendly than their unambiguous  
 600 versions [2, 26]. Therefore, the purpose of this section is to revise the architecture of  
 601 basic BiYACC to allow the use of ambiguous grammars and disambiguation  
 602 directives. This is in fact a long-standing problem: tools designed for building parser  
 603 and printer pairs usually do not support such functionality (Sect. 7.1).

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

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

## 620 Problems with Ambiguous Grammars

621 Consider the original architecture of BiYACC in Fig. 3, which we want to (and  
 622 basically will) retain while adapting it to support ambiguous grammars. The first

623 component (of the executable) we should adapt is  $cparse :: \text{String} \rightarrow \text{Maybe CST}$ ,  
 624 the (concrete) parsing direction of the isomorphism: since there can be multiple  
 625 CSTs corresponding to the same program text,  $cparse$  needs to choose one of them  
 626 as the result. Disambiguation directives [23] were invented to describe how to make  
 627 this choice. For example, with respect to the grammar in Fig. 6, text  $1 + 2 * 3$  will  
 628 have either of the two CSTs<sup>4</sup>:

$$\begin{aligned} \text{cst}_1 &= \text{Plus } 1 \text{ (Times } 2 \text{ )} \\ \text{cst}_2 &= \text{Times (Plus } 1 \text{ ) } 3 \end{aligned}$$

630 depending on the precedence of addition and multiplication. Conventionally, we can  
 631 use the YACC-style disambiguation directives `%left '+'`; `%left '**'`; to specify that  
 632 multiplication has higher precedence over addition, and instruct the parser to choose  
 633  $\text{cst}_1$ .

634 However, merely adapting  $cparse$  with disambiguation behaviour is not enough,  
 635 since the isomorphism (Theorem 1), in particular its right to left direction (which is  
 636 simplified as  $cparse (cprint \text{ cst}) = \text{Just } \text{cst}$ ) cannot be established when  
 637 an ambiguous grammar is used—in the example above,  $cparse (cprint \text{ cst}_2) = \text{Just}$   
 638  $\text{cst}_1 \neq \text{Just } \text{cst}_2$ . This is because the image of  $cparse$  is strictly smaller than the  
 639 domain of  $cprint$ : if we start from any CST not in the image of  $cparse$ , we will never  
 640 be able to get back to the same CST through  $cprint$  and then  $cparse$ . This tells us  
 641 that, to retain the isomorphism, the domain of  $cprint$  should not be the whole CST  
 642 but only the image of  $cparse$ , i.e. the set of valid CSTs (as defined by the  
 643 disambiguation directives), which we denote by  $\text{CST}_F$  (for reasons that will be made  
 644 clear in Sect. 5.3).

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

658 We are now led to a fundamental problem: generally,  $put$  strategies for producing  
 659 valid CSTs should be inferred from the disambiguation directives, but the semantics  
 660 of YACC disambiguation directives are defined over the implementation of YACC's  
 661 underlying LR parsing algorithm with a stack [3, 23], and therefore it is nontrivial to  
 662 invent a dual semantics in the  $put$  direction. To have a simple and clear semantics of

4FL01 <sup>4</sup> For simplicity, we use  $\#$  to annotate type-incorrect CSTs in which fields for layouts (and comments) and  
 4FL02 unimportant constructors such as `Lit` are omitted.

663 the disambiguation process, we turn away from YACC's traditional approach and opt  
 664 for an alternative approach based on generalised parsing with disambiguation filters  
 665 [9, 26], whose semantics can be specified implementation independently. Based on  
 666 this simple and clear semantics, we will be able to devise ways to amend *put* to  
 667 produce only valid CSTs, and formally state the conditions under which the  
 668 executable generated by the revised BiYACC is well behaved.

## 669 Generalised Parsing and Bidirectionalised Filters

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

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

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

681 With generalised parsing, program text is first parsed to all the possible CSTs;  
 682 disambiguation then becomes an extremely simple concept: removing CSTs that the  
 683 user does not want. One possible semantics of disambiguation may be a function  
 684 *judge* :: *Tree* → *Bool*; during disambiguation, this function is applied to all  
 685 candidate CSTs, and a candidate *cst* is removed if *judge cst* returns *False*, or  
 686 kept otherwise. We call these functions disambiguation filters ('filters' for short).<sup>5</sup>  
 687 For example, to state that top-level addition is left-associative, we can use the  
 688 following filter<sup>6</sup> to reject right-sided trees:  
 689

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

691 This simple and clean semantics of disambiguation is then amenable to  
 692 'bidirectionalisation', which we do next.

693 Note that, unlike YACC's disambiguation directives, which assign precedence and  
 694 associativity to individual tokens and implicitly exclude 'some' CSTs, in *plusJudge*  
 695 above we explicitly ban incorrect CSTs through pattern matching. Having described  
 696 which CSTs are incorrect, we can further specify what to do with incorrect CSTs in

5FL01 <sup>5</sup> The general type for disambiguation filters is  $[t] \rightarrow [t]$ , which allows comparison among a list of CSTs.  
 5FL02 However, since in this paper we only consider property filters defined in terms of predicates (on a single  
 5FL03 tree), it is sufficient to use the simplified type  $t \rightarrow \text{Bool}$ . See Sect. 7.2.

6FL01 <sup>6</sup> This is not a very realistic filter, although it sufficiently demonstrates the use of filters and removes  
 6FL02 ambiguity in simplest cases like  $1 + 2 * 3$ . In general, the filter should be complete (Definition 9) so that  
 6FL03 ambiguity is fully removed from the grammar.

697 the printing direction. Whenever a CST ‘in a bad shape’, i.e. rejected by a filter like  
 698 `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 .
```

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

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

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

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

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

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

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

714 satisfying

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

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

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

717

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

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

$$judge_F\ t = \text{True} \Rightarrow repair_F\ t = t.$$

727

728 **Proof** From  $judge_F\ t = \text{True}$  we deduce  $F\ t = \text{Nothing}$ , which implies  $repair_F\ t = t$ .  
 729 □

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

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

$$get \circ repair_F = get.$$

736

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

## 742 The New BiYACC System for Ambiguous Grammars

743 As depicted in Fig. 7, the executable generated by the new BiYACC system is still the  
 744 composition of an isomorphism and a lens, which is the structure we have tried to  
 745 retain. To precisely identify the changes in several generated components (in the  
 746 executable file) and demonstrate how parsing and printing work with a bi-filter, we  
 747 present Fig. 8 and will use this one instead. In the new system, we will still use the  
 748 *get* and *put* transformations generated from printing actions and the concrete printer  
 749 *cprint* from grammars, while the concrete parser *cparse* is replaced with a  
 750 generalised parser *cgpars*. Additionally, the #Directives and #OtherFilters parts will  
 751 be used to generate a bi-filter  $F$ , whose  $judge_F$  (used in the  $selectBy_F$  function in  
 752 Fig. 8) and  $repair_F$  components are integrated into the isomorphism and lens parts

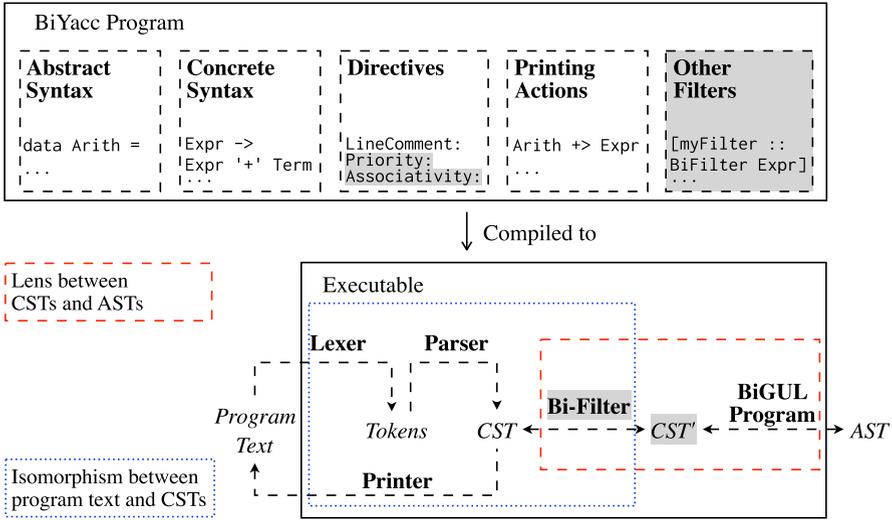


Fig. 7 New architecture of BiYACC (new components are in light grey)

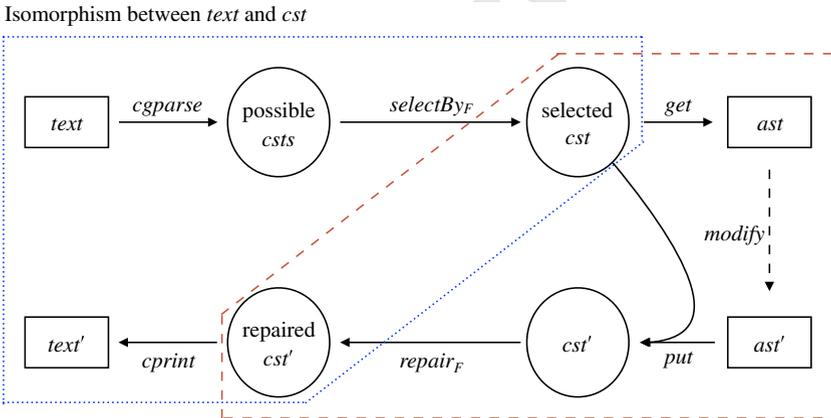


Fig. 8 A schematic diagram showing how parsing and printing work with a bi-filter

753 respectively, so that the right-hand side domain of the isomorphism and the source  
 754 of the lens become  $CST_F$ , the set of valid CSTs:

$$CST_F = \{ cst \in CST \mid judge_F\ cst = True \} .$$

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

758 **The Revised Isomorphism between Program Text and CSTs**

759 Let us first consider the isomorphism part between `String` and `CSTF`, which is  
 760 enclosed within the blue dotted lines in Fig. 8 and consists of `cprint`, `cgparse`, and  
 761 `selectByF`:

$$\begin{aligned}
 cprint &:: \text{CST} \rightarrow \text{String} \\
 cgparse &:: \text{String} \rightarrow [\text{CST}] \\
 selectBy_F &:: [\text{CST}] \rightarrow \text{Maybe } \text{CST}_F \\
 selectBy_F \text{ csts} &= \text{case } selectBy \text{ judge}_F \text{ csts of} \\
 &\quad [cst] \rightarrow \text{Just } cst \\
 &\quad \_ \rightarrow \text{Nothing} \\
 selectBy &:: (a \rightarrow \text{Bool}) \rightarrow [a] \rightarrow [a] \\
 selectBy \text{ p} [] &= [] \\
 selectBy \text{ p} (x : xs) | \text{px} = x &: selectBy \text{ p} \text{ xs} \\
 selectBy \text{ p} (x : xs) | \text{otherwise} &= selectBy \text{ p} \text{ xs} .
 \end{aligned}$$

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

$$\begin{aligned}
 cparse_F &:: \text{String} \rightarrow \text{Maybe } \text{CST}_F \\
 cparse_F &= selectBy_F \circ cgparse \\
 cprint_F &:: \text{CST}_F \rightarrow \text{Maybe } \text{String} \\
 cprint_F &= \text{Just} \circ cprint .
 \end{aligned}$$

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

775 **Definition 8** (Generalised parser correctness) A generalised parser `cgparse` is  
 776 correct with respect to a printer `cprint` exactly when

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

777

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



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

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

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

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

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

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

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

794 **Proof** We reason:  
 795

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

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

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

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

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

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

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

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

809

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

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

### 817 The Revised Lens between CSTs and ASTs

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

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

821 To work with a bi-filter  $F$ , in particular its  $repair_F$  component, they need to be  
822 adapted to  $get_F$  and  $put_F$ , which accept only valid CSTs:

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

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

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

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

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

830 Now we prove that  $get_F$  and  $put_F$  are well-behaved, which is a generalisation of  
831 Theorem 2 for ambiguous grammars.

832 **Theorem 4** (*Well-behavedness with bi-filters*) Given a complete bi-filter  $F$  and a  
 833 well-behaved lens consisting of  $get$  and  $put$ , if  $get$  and  $F$  additionally satisfy  
 834 *PassThrough*, then the  $get_F$  and  $put_F$  functions with respect to  $F$  are also well-  
 835 behaved:

$$put_F\ cst\ ast = \text{Just}\ cst' \Rightarrow get_F\ cst' = \text{Just}\ ast \quad (8)$$

837  $get_F\ cst = \text{Just}\ ast \Rightarrow put_F\ cst\ ast = \text{Just}\ cst . \quad (9)$

838  
 839 **Proof** For (8): the antecedent expands to  $fmap\ repair_F\ (put\ cst\ ast) = \text{Just}\ cst'$ ,  
 840 which, by Lemma 5, implies  $put\ cst\ ast = \text{Just}\ cst''$  for some  $cst''$  such that  
 841  $repair_F\ cst'' = cst'$ . Now we reason:

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

843 For (9):

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

845 □

846 **Bi-Filter Directives**

847 Until now, we have only considered working with a single bi-filter, but this is  
 848 without loss of generality because we can provide a bi-filter composition operator  
 849 (Sect. 5.4.1) so that we can build large bi-filters from small ones. This is a

850 suitable semantic foundation for introducing YACC-like directives for specifying  
 851 priority and associativity into BiYACC (Sect. 5.4.2), since we can give these  
 852 directives a bi-filter semantics and interpret a collection of directives as the  
 853 composition of their corresponding bi-filters. We will also discuss some properties  
 854 related to this composition (Sect. 5.4.3).

## 855 Bi-Filter Composition

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

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

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

862

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

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

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

868 **Proof** By the definition of bi-filter composition. □

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

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

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

875

876 If  $j$  respects  $i$ , then a later applied  $repair_j$  will never break what may already be  
 877 repaired by a previous  $repair_i$ . Thus in this case we can safely compose  $j$  after  
 878  $i$ . This is proved as the following theorem.

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

881 **Proof** For *RepairJudge*, we reason:

$$\begin{aligned}
 & judge(j \triangleleft i)(repair (j \triangleleft i) t) \\
 = & \{ \text{Lemma 6} \} \\
 & judge(j \triangleleft i)(repair_j(repair_i t)) \\
 = & \{ \text{Lemma 6} \} \\
 & judge_j(repair_j(repair_i t)) \wedge judge_i(repair_j(repair_i t)) \\
 = & \{ \text{RepairJudge of } j \} \\
 & \text{True} \wedge judge_i(repair_j(repair_i t)) \\
 = & \{ judge_i(repair_i t) = \text{True}; j \text{ respects } i \} \\
 & \text{True} \wedge \text{True} \\
 = & \text{True} .
 \end{aligned}$$

883 And for *PassThrough*:

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

885 □

## 886 Priority and Associativity Directives

887 To relieve the burden of writing bi-filters manually and guaranteeing respect among  
 888 bi-filters being composed, we provide some directives for constructing bi-filters  
 889 dealing with priority<sup>7</sup> and associativity, which are generally comparable to YACC's  
 890 conventional disambiguation directives. The bi-filter directives in a BiYACC program  
 891 can be thought of as specifying 'production priority tables', analogous to the  
 892 operator precedence tables of, for example, the C programming language [24]

7FL01 <sup>7</sup> The YACC-style approach adopts the word *precedence* [23] while the filter-based approaches tend to use  
 7FL02 the word *priority* [9, 26]. We follow the traditions and use either word depending on the context.

893 (chapter *Expressions*) and HASKELL [34] (page 51). The main differences (in terms of  
894 the parsing direction) are as follows:

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

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

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

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

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

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

922 In order for bi-filter directives to work correctly, the user should notice the  
923 following requirements: (1) directives shall not mention the parenthesis production  
924 annotated with bracket attribute so that they respect each other and work properly  
925 (as introduced in Definition 11). (2) Suppose that the parenthesis production is  
926  $NT \rightarrow \alpha NT_R \beta$  where  $\alpha$  and  $\beta$  denote a sequence of terminals and  $NT_R$  is a possibly  
927 different nonterminal from  $NT$  (on the right-hand side of the production)—for  
928 instance, Expr  $\rightarrow$  '(' Expr ')' above—there shall be exactly one printing action  
929 defined for the parenthesis production in the form of  $v \rightarrow \alpha[v \rightarrow NT_R]\beta$  for the  
930 PassThrough property to hold: for any CST, the (added) parenthesis structure will  
931 all be dropped through the conversion to its AST.



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

## 938 Priority Directives

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

$$\text{Expr} \rightarrow \text{Expr} \text{'*'} \text{Expr} > \text{Expr} \rightarrow \text{Expr} \text{'+'} \text{Expr} ;$$

or just Times > Plus; .

945 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 undefined .
```

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

953 In general, the syntax of priority directives is

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

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

8FL01 <sup>8</sup> Although terminals such as “\*” and “+” are uniquely determined by constructors and not explicitly  
 8FL02 included in the CSTs, there are fields in CSTs for holding whitespaces after them. Thus Times still has  
 8FL03 three subtrees. Also, for simplicity, the bi-filter fTimesPlusPrio attempts to repair the whitespace  
 8FL04 subtree  $t_2$  even though the repair can never happen since  $t_2$  cannot match  $p$ .



957 If the user declares that a production  $NT_1 \rightarrow RHS_1$  has higher priority than  
 958 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) .
  
```

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

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

$$\text{Expr} \rightarrow \text{Expr} \text{ '*' } \text{Expr} > \text{Expr} \rightarrow \text{Expr} \text{ '+' } \text{Expr};$$

$$\text{Expr} \rightarrow \text{Expr} \text{ '+' } \text{Expr} > \text{Expr} \rightarrow \text{Expr} \text{ '&' } \text{Expr};$$

969 implies that  $\text{Expr} \rightarrow \text{Expr} \text{ '*' } \text{Expr} > \text{Expr} \rightarrow \text{Expr} \text{ '&' } \text{Expr}$ . The feature is important in  
 970 practice since it greatly reduces the amount of routine code the user needs to write  
 971 (for large grammars).

## 972 Associativity Directives

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

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

$$\text{Left:Expr} \rightarrow \text{Expr} \text{ '+' } \text{Expr}, \text{Expr} \rightarrow \text{Expr} \text{ '-' } \text{Expr};$$

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

984

```

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 .

```

986 For instance, `fPlusPlusLAssoc` accepts `#Plus (Plus 1 2) 3` but not  
 987 `#Plus 1 (Plus 2 3)`, which is repaired to `#Plus 1 (Paren (Plus 2 3))`.  
 988 Generally, the syntax of associativity directives is

*Associativity ::= 'Associativity: LeftAssoc RightAssoc*  
*LeftAssoc ::= 'Left: ProdOrCons<sup>+</sup>{','};'*  
*RightAssoc ::= 'Right: ProdOrCons<sup>+</sup>{','};'* .

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

```

TO_LASSOC_FILTER[ $\alpha_1 NT_{1R}, NT_1, \alpha_2 NT_{2R}, NT_2$ ] =
  '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 = FILLVARSFROM(LENGTH( $\alpha_1$ ),  $NT_{1R}$ ) .

```

1001

1002 Functions `CON`, `FILLUNDEFINED`, and `FILLVAR` have the same behaviour as before;  
 1003 `FILLVARS—FROM` (which is a variation of `FILLVARS`) generates variable names for each  
 1004 terminal and nonterminal in its argument with suffix integers counting from a given  
 1005 number to avoid name clashing.

1006 *Handling injective productions* Sometimes the grammar may contain injective  
 1007 productions (also called chain productions) [9], which have only a single  
 1008 nonterminal on their right-hand side, like  $\text{InfE} \rightarrow [\text{FromE}] \text{Exp}$ . When we use it  
 1009 to define a grammar

$$\begin{aligned} \text{InfE} &\rightarrow [\text{FromE}] \text{Exp} \\ \text{Exp} &\rightarrow [\text{Plus}] \text{InfE} '+' \text{InfE} \\ &\quad | [\text{Times}] \text{InfE} '*' \text{InfE} , \end{aligned}$$

1011 program text  $1 + 2 * 3$  will be parsed to two CSTs, namely  $\text{cst}_1 = \#_{\text{Plus}} (\text{FromE } 1)$  and  
 1012  $\text{cst}_2 = \#_{\text{Times}} (\text{FromE } (\text{Plus } 1 \ 2) (\text{FromE } 3))$ , and we want to spot  $(\text{FromE } (\text{Times } 2 \ 3))$   
 1013 and discard it using the priority directive `Times > Plus`. If handled naively, the bi-  
 1014 filter generated from the directive would only remove CSTs having pattern  
 1015  $\text{Times} (\text{Plus } \_ \_)$  (and two other similar ones), but  $\text{cst}_2$  would not match the pattern  
 1016 due to the presence of the `FromE` node between `Times` and `Plus`. We made some  
 1017 effort in the implementation to make the match function ignore the nodes  
 1018 corresponding to injective productions (`FromE` in this case).

## 1019 Properties of the Generated Bi-Filters

1020 We discuss some properties of the bi-filters generated from our priority and  
 1021 associativity directives, to justify that it is safe to use these bi-filters without  
 1022 disrupting the well-behavedness of the whole system. Specifically:

- 1023 • The generated top-level bi-filters satisfy [RepairJudge](#), and it is easy to write  
 1024 actions to make them satisfy [PassThrough](#).
- 1025 • The bi-filters lifted from the top-level bi-filters still satisfy [RepairJudge](#) and  
 1026 [PassThrough](#).
- 1027 • The lifted bi-filters are commutative, which not only implies that all such bi-  
 1028 filters respect each other and can be composed in any order, but also guarantees  
 1029 that we do not have to worry about the order of composition since it does not  
 1030 affect the behaviour.

1031 We will give only high-level, even informal, arguments for these properties, since,  
 1032 due to the generic nature of the definitions of these bi-filters (in terms of *Scrap Your  
 1033 Boilerplate* [30]), to give formal proofs we would have to introduce rather complex  
 1034 machinery (e.g. datatype-generic induction), which would be tedious and  
 1035 distracting.

1036 *Top-level bi-filters* The fact that the generated top-level bi-filters satisfy  
 1037 [RepairJudge](#) can be derived from the requirement that the directives do not  
 1038 mention the parenthesis production. Because of the requirement, in the generated bi-  
 1039 filters, repairing is always triggered by matching a non-parenthesis production, and  
 1040 after that repairing will not be triggered again because a parenthesis production will  
 1041 have been added. For example, in the bi-filter `fTimesPlusPrio` (in Sect. 5.4.2), with



1042 match  $t_1$   $p$ , match  $t_2$   $p$ , and match  $t_3$   $p$  we check whether  $t_1$ ,  $t_2$ , and  $t_3$  has Plus as the  
 1043 top-level production, which is different from the parenthesis production Paren; if  
 1044 any of the matching succeeds, say  $t_1$ , then  $\text{addPar } t_1$  will add Paren at the top of  $t_1$ ,  
 1045 and  $\text{match}(\text{addPar } t_1) p$  is guaranteed to be False, so the subsequent invocation of  
 1046  $\text{judge } f\text{TimesPlusPrio}$  will return True. For [PassThrough](#), since all the top-level bi-  
 1047 filters do is add parenthesis productions, we can simply make sure that appearances  
 1048 of the parenthesis production are ignored by  $\text{get}$ , i.e.  $\text{get}(\text{addPar } s) = \text{get } s$  for  
 1049 all  $s$ ; this, by well-behavedness, is the same as making  $\text{put}$  (printing actions) skip  
 1050 over parentheses. For example, for the grammar in Figure 6, we should write  
 1051  $t \mapsto \text{'(' [t} \mapsto \text{Expr] ')}$  as the only printing action mentioning parentheses, which  
 1052 means that  $\text{put}(\text{Paren } s) t = \text{fmap } \text{Paren}(\text{put } s t)$  for all  $s$  and  $t$ . Then the  
 1053 following reasoning implies that  $\text{get}(\text{Paren } s) = \text{get } s$  for all  $s$ :

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

1055 for all  $s$  and  $t$ .

1056 *Lifted bi-filters* The lifted bi-filters apply the top-level bi-filters to all the subtrees  
 1057 in a CST in a bottom-up order. Formally, we can define, datatype-generically, a  
 1058 lifted bi-filter as a composition of top-level bi-filters, and use datatype-generic  
 1059 induction to prove that there is suitable respect among the top-level bi-filters being  
 1060 composed, and that the lifted bi-filter satisfies [RepairJudge](#) and [PassThrough](#) if the  
 1061 top-level ones do. But here we provide only an intuitive argument. What the lifted  
 1062 bi-filters do is find all prohibited pairs of adjoining productions and separate all the  
 1063 pairs by adding parenthesis productions. For [RepairJudge](#), since all prohibited pairs  
 1064 are eliminated after repairing, there will be nothing left to be repaired in the  
 1065 resulting CST, which will therefore be deemed valid. For [PassThrough](#), the intuition  
 1066 is the same as that for the top-level bi-filters.

1067 *Commutativity* Composite bi-filters  $i \triangleleft j$  and  $j \triangleleft i$  may have different behaviours,  
 1068 so in general we need to know the order of composition to figure out the exact  
 1069 behaviour of a composite bi-filter. This can be difficult when using our bi-filter  
 1070 directives, since a lot of bi-filters are implicitly generated from the directives, and it  
 1071 is not straightforward to specify the order in which all the explicitly and implicitly  
 1072 generated bi-filters are composed. Fortunately, we do not need to do so, for all the  
 1073 bi-filters generated from the directives are commutative, meaning that the order of  
 1074 composition does not affect the behaviour.

1075 **Definition 12** (Bi-filter commutativity) Two bi-filters  $i$  and  $j$  are commutative  
 1076 exactly when

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

1079 By Lemma 6, this implies  $\text{repair}(i \triangleleft j) = \text{repair}(j \triangleleft i)$ . Note that  
 1080  $\text{judge}(i \triangleleft j) = \text{judge}(j \triangleleft i)$  by definition, so we do not need to require this in  
 1081 the definition of commutativity.  
 1082

1083 An important fact is that commutativity is stronger than respect, so it is always  
 1084 safe to compose commutative bi-filters.

1085 **Lemma 7** *Commutative bi-filters respect each other.*

1086 **Proof** Given commutative bi-filters  $i$  and  $j$ , we show that  $j$  respects  $i$ . Suppose that  
 1087  $\text{judge}_i t = \text{True}$  for a given tree  $t$ . Then

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

1089 It follows by symmetry that  $i$  respects  $j$  as well. □

1090 Now let us consider why any two different lifted bi-filters are commutative.  
 1091 (Commutativity is immediate if the two bi-filters are the same.) There are two key  
 1092 facts that lead to commutativity: (1) repairing does not introduce more prohibited  
 1093 pairs of productions, and (2) the prohibited pairs of adjoining productions checked  
 1094 and repaired by the two bi-filters are necessarily different. Therefore the two bi-  
 1095 filters always repair different parts of a tree, and can repair the tree in any order  
 1096 without changing the final result. Fact (1) is, again, due to the requirement that the  
 1097 directives do not mention the parenthesis production, which is the only thing we add  
 1098 to a tree when repairing it. Fact (2) can be verified by a careful case analysis. For  
 1099 example, we might be worried about the situation where a left-associative directive  
 1100 looks for production  $Q$  used at the right-most position under production  $P$ , while a  
 1101 priority directive also similarly looks for  $Q$  used under  $P$ , but the two directives  
 1102 cannot coexist in the first place since the first directive implies  $P$  and  $Q$  have no  
 1103 relative priority whereas the second one implies  $Q$  has lower priority than  $P$ .

## 1104 Manually Written Bi-Filters

1105 There are some other ambiguities that our directives cannot eliminate. In these  
 1106 cases, the user can define their own bi-filters and put them in the #OtherFilters part in  
 1107 a BiYACC program as shown in Fig. 4. The syntax is

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

1109 That is, this part of the program begins with a list of declarations of the names and  
1110 types of the user-defined bi-filters, whose HASKELL definitions are then given below.

1111 Now we demonstrate how to manually write a bi-filter by resolving the ambiguity  
1112 brought by the dangling else problem. But before that, let us briefly review the  
1113 problem, which arises, for example, in the following grammar:  
1114

$$\begin{aligned} \text{Exp} &\rightarrow \text{[ITE] 'if' Exp 'then' Exp 'else' Exp} \\ &\quad | \text{[IT] 'if' Exp 'then' Exp .} \end{aligned}$$

1116 With respect to this grammar, the program text if a then if x then y else z can be  
1117 recognised as either if a then (if x then y else z) or if a then (if x then y) else z. To  
1118 resolve the ambiguity, usually we prefer the ‘nearest match’ strategy (which is  
1119 adopted by Pascal, C, and Java): else should match its nearest then, so that if a then  
1120 (if x then y else z) is the only correct interpretation.

1121 The user may think that the problem can be solved by a priority (bi-)filter  
1122  $\text{ITE} > \text{IT}$ ; in the hope that the production ‘if-then-else’ binds tighter than the  
1123 production ‘if-then’. Unfortunately, this is incorrect as pointed out by Klint and  
1124 Visser [26], because the corresponding (bi-)filter incorrectly rules out the pattern  
1125  $\#_{\text{ITE}} \_ \_ (\text{IT} \_ \_)$ , which prints to unambiguous text, e.g. if a then b else if x then y. In  
1126 fact, the (dangling else) problem is tougher than one might think and cannot be  
1127 solved by any (bi-)filter performing pattern matching with a fixed depth [26].

1128 Klint and Visser [26] proposed an idea to disambiguate the dangling-else  
1129 grammar: let Greek letters  $\alpha, \beta, \dots$  match a sequence of symbols of any length. Then  
1130 the program text if  $\alpha$  then  $\beta$  else  $\gamma$  should be banned if the right spine of  $\beta$  contains  
1131 any if  $\psi$  then  $\omega$ , as shown in the paper [26]. With the full power of (bi-)filters, which  
1132 are fully fledged HASKELL functions, we can implement this solution in the following  
1133 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 = ... .
```

1135 This bi-filter is commutative with the bi-filters generated from our directives,  
1136 since it (1) only searches for non-parenthesis productions that are not declared in  
1137 any other directives, and (2) inserts only a parenthesis production when repairing



1138 incorrect CSTs. The reader may find the code of `checkRightSpine` in more detail in  
 1139 Fig. 10.

## 1140 Case Studies

1141 The design of `BiYACC` may look simplistic and make the reader wonder how much it  
 1142 can describe. In fact, `BiYACC` can already handle real-world language features. For  
 1143 example, Kinoshita and Nakano [25] adopted `BiYACC` as part of their system for  
 1144 synchronising `Coq` functions and corresponding `Ocaml` programs. In this section,  
 1145 we demonstrate `BiYACC` with a medium-size case study: we use `BiYACC` to build a  
 1146 pair of parser and reflective printer for the `Tiger` language [4] and demonstrate some  
 1147 of their uses.

## 1148 The `Tiger` Language

1149 `Tiger` is a statically typed imperative language first introduced in Appel's textbook  
 1150 on compiler construction [4]. Since `Tiger`'s purpose of design is pedagogical, it is  
 1151 not too complex and yet covers many important language features including  
 1152 conditionals, loops, variable declarations and assignments, and function definitions  
 1153 and calls. `Tiger` is therefore a good case study with which we can test the potential  
 1154 of our `BX`-based approach to constructing parsers and reflective printers. Some of  
 1155 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")))
```

1157 To give a sense of `Tiger`'s complexity, it takes a grammar with 81 production  
 1158 rules to specify `Tiger`'s syntax, while for `C89` and `C99` it takes, respectively, 183  
 1159 and 237 rules without any disambiguation declarations (based on Kernighan and  
 1160 Ritchie [24] and the draft version of 1999 ISO C standard, excluding the  
 1161 preprocessing part). The difference is basically due to the fact that `C` has more  
 1162 primitive types and various kinds of assignment statements.

1163 Excerpts of the abstract and concrete syntax of `Tiger` are shown in Fig. 9. The  
 1164 abstract syntax is largely the same as the original one defined in Appel's textbook  
 1165 (page 98); as for the concrete syntax, Appel does not specify the whole grammar in  
 1166 detail, so we use a version slightly adapted from Hirzel and Rose's lecture notes  
 1167 [21]. Concretely, we add a parenthesis production to the grammar (and discard it  
 1168 when converting CSTs to ASTs, so that the `PassThrough` property could be  
 1169 satisfied), since `Tiger`'s original grammar has no parenthesis production and an  
 1170 expression within round parentheses is regarded as a singleton expression sequence.  
 1171 This modification also makes it necessary to change the enclosing brackets for  
 1172 expression sequences from round brackets `()` to curly brackets `{}`, which helps  
 1173 (LALR(1) parsers) to distinguish a singleton expression sequence from an  
 1174 expression within parentheses. There is also another slight change in the definition

```

#Abstract
type TSymbol = String
data Tuple a b = Tuple a b
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)

1175 of ASTs for handling a feature not supported by the current BiYACC: the AST  
 1176 constructors TFunctionDec and TTypeDec take a single function or type declaration  
 1177 instead of a list of adjacent declarations (for representing mutual recursion) as in  
 1178 Appel [4], since we cannot handle the synchronisation between a list of lists (in  
 1179 ASTs) and a list (in CSTs) with BiYACC's current syntax.

1180 Following Hirzel and Rose's specification [21], the disambiguation directives for  
 1181 TIGER are shown in Fig. 10; for instance, we define multiplication to be left-  
 1182 associative. The directives also include a concrete treatment for the dangling else  
 1183 problem, which is usually 'not solved' when using a YACC-like (LA)LR parser  
 1184 generator to implement parsers: rather than resolving the grammatical ambiguity,  
 1185 we often rely on the default behaviour of the parser generator—preferring shift.



```

#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)

1186 We have successfully tested our BiYACC program for TIGER on all the sample  
 1187 programs provided on the homepage of Appel's book,<sup>9</sup> including a merge sort  
 1188 implementation and an eight-queen solver, and there is no problem parsing and  
 1189 printing them with well-behavedness guaranteed. In the following subsections, we  
 1190 will present some printing strategies described in the BiYACC program to  
 1191 demonstrate what BiYACC, in particular reflective printing, can achieve.

## 1192 Syntactic Sugar and Resugaring

1193 We start with a simple example about syntactic sugar, which is pervasive in  
 1194 programming languages and lets the programmer use some features in an alternative  
 1195 (usually conceptually higher-level) syntax. For instance, TIGER represents boolean  
 1196 values *false* and *true*, respectively, as zero and nonzero integers, and the logical  
 1197 operators `&` ('and') and `|` ('or') are converted to a conditional structure in the  
 1198 abstract syntax: `e1 & e2` is desugared and parsed to `TCond e1 e2 (TInt 0)` and `e1 | e2`  
 1199 to `TCond e1 (TInt 1) e2`. The printing actions for them in BiYACC are:

9FL01 <sup>9</sup> <https://www.cs.princeton.edu/~appel/modern/testcases/>.

```

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

1201 A conventional printer which takes only the AST as input cannot reliably  
 1202 determine whether an abstract expression should be printed to the basic form or the  
 1203 sugared form, whereas a reflective printer can make the correct decision by  
 1204 inspecting the CST.

1205 The idea of resugaring [42] is to print evaluation sequences in a core language in  
 1206 terms of a surface syntax. Here we show that, without any extension, BiYACC is already  
 1207 capable of propagating some AST changes that result from evaluation back to the  
 1208 concrete syntax, subsuming a part of Pombrio and Krishnamurthi’s work [42, 43].

1209 We borrow their example of resugaring evaluation sequences for the logical  
 1210 operators ‘or’ and ‘not’, but recast the example in TIGER. The ‘or’ operator has been  
 1211 defined as syntactic sugar in Section 6.2. For the ‘not’ operator, which TIGER lacks,  
 1212 we introduce ‘~’, represented by TNot in the abstract syntax. Now consider the  
 1213 source expression

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

1215 which is parsed to

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

1217 A typical evaluator will produce the following evaluation sequence given the above  
 1218 AST:

```

TCond (TNot (TInt 1)) (TInt 1) (JJ (TNot (TInt 0)))
→ TCond (TInt 0) (TInt 1) (JJ (TNot (TInt 0)))
→ TNot (TInt 0)
→ TInt 1 .
    
```

1220 If we perform reflective printing after every evaluation step using BiYACC, we  
 1221 will get the following evaluation sequence on the source:  
 1222

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

1224 Due to the PUTGET property, parsing these concrete terms will yield the  
 1225 corresponding abstract terms in the abstract evaluation sequence, and this is exactly  
 1226 Pombrio and Krishnamurthi’s ‘emulation’ property, which they have to prove for their  
 1227 system. For BiYACC, however, the emulation property holds by construction, since  
 1228 BiYACC programs are always well-behaved. Another difference is that we do not need  
 1229 to insert additional information (such as tags) into an AST for recording which surface  
 1230 syntax structure a node comes from. One advantage of our approach is that we keep the  
 1231 abstract syntax pure, so that other tools—the evaluator in particular—can process the  
 1232 abstract syntax without being modified, whereas in Pombrio and Krishnamurthi’s  
 1233 approach, the evaluator has to be adapted to work on an enriched abstract syntax.

1234 **Language Evolution**

1235 When a language evolves, some new features of the language (e.g. the foreach loops  
 1236 introduced in Java 5 [20]) can be implemented by desugaring to some existing  
 1237 features (e.g. ordinary for loops), so that the compiler back-end and abstract syntax  
 1238 definition do not need to be extended to handle the new features. As a consequence,  
 1239 all the engineering work about optimising transformations or refactoring [18] that  
 1240 has been developed for the abstract syntax remains valid.

1241 Consider a kind of ‘generalised-if’ expression allowing more than two cases,  
 1242 resembling the alternative construct in Dijkstra’s guarded command language [12].  
 1243 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 ; .
```

1246 For simplicity, we restrict the predicate produced by CaseB to the form LValue ‘=’  
 1247 Numeric, but in general the Numeric part can be any expression computing an  
 1248 integer. The reflective printing actions for this new construct can still be written  
 1249 within B1YACC, but require much deeper pattern matching:  
 1250

```
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];
;; .
```

1252 Although being a little complex, these printing actions are in fact fairly  
 1253 straightforward: The first group of type Tiger+> Guard handles the enclosing  
 1254 guard–end pairs, distinguishes between single- and multi-branch cases, and delegates  
 1255 the latter case to the second group, which prints a list of branches recursively.

1256 This is all we have to do—the corresponding parser is automatically derived and  
 1257 guaranteed to be consistent. Now guard expressions are desugared to nested if  
 1258 expressions in parsing and preserved in printing, and we can also resugar evaluation  
 1259 sequences on the ASTs to program text. For instance, the following guard expression

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

1261 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))))
```

1263 where TSimpleVar is shortened to TSV, and choice is shortened to c. Suppose that  
 1264 the value of the variable choice is 2. The evaluation sequence on the AST will then  
 1265 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 .
```

1269 And the reflected evaluation sequence on the concrete expression will be:  
 1270  
 1271

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

1273 Reflective printing fails for the first and third steps (the program text becomes an  
 1274 if-then-else expression if we do printing at these steps), but this behaviour in fact  
 1275 conforms to Pombrio and Krishnamurthi's 'abstraction' property, which demands  
 1276 that core evaluation steps that make sense only in the core language must not be  
 1277 propagated to the surface. In our example, the first and third steps in the TCond-  
 1278 sequence evaluate the condition to a constant, but conditions in guard expressions  
 1279 are restricted to a specific form and cannot be a constant; evaluation of guard  
 1280 expressions thus has to proceed in bigger steps, throwing away or going into a  
 1281 branch in each step, which corresponds to two steps for TCond.

1282 The reader may have noticed that, after the guard expression is reduced to two  
 1283 branches, the layout of the second branch is disrupted; this is because the second  
 1284 branch is in fact printed from scratch. In current BiYACC, the printing from an AST  
 1285 to a CST is accomplished by recursively performing pattern matching on both tree  
 1286 structures. This approach naturally comes with the disadvantage that the matching is  
 1287 mainly decided by the position of the nodes in the AST and CST. Consequently, a  
 1288 minor structural change on the AST may completely disrupt the matching between  
 1289 the AST and the CST.

## 1290 Other Potential Applications

1291 We conclude this section by shortly discussing several other potential applications.  
 1292 In general, (current) BiYACC can easily and reliably propagate AST changes that  
 1293 have local effect such as replacing part of an AST with a simpler tree, without  
 1294 destroying the layouts and comments of unaffected code. Thus it would not be  
 1295 surprising that BiYACC can also propagate (1) simplification-like optimisations such  
 1296 as constant folding and constant propagation and (2) some code refactoring  
 1297 transformations such as variable renaming. All these functionalities are achieved for  
 1298 free by one ‘general-purpose’ BiYACC program, which does not need to be tailored  
 1299 for each application.

## 1300 Related Work

### 1301 Unifying Parsing and Printing

1302 Much research has been devoted to describing parsers and printers in a single  
 1303 program. For example, both Rendel and Ostermann [44] and Matsuda and Wang  
 1304 [36, 37] adopt a combinator-based approach<sup>10</sup> (whereas we use a generator-based  
 1305 approach), where small components are glued together to yield more sophisticated  
 1306 behaviour, and can guarantee properties similar to Theorem 1 with *cst* replaced by  
 1307 *ast* in the equations. (Let us call the variant version Theorem 1', since it will be used  
 1308 quite often later.) In Rendel and Ostermann's system (called ‘invertible syntax  
 1309 descriptions’, which we shorten to ISDs henceforth), both the parsing and printing  
 1310 semantics are predefined in the combinators and consistency is guaranteed by their  
 1311 partial isomorphisms, whereas in Matsuda and Wang's system (called FLIPPR), the  
 1312 combinators describing pretty printing are translated by a semantic-preserving  
 1313 transformation to a core syntax, which is further processed by their grammar-based  
 1314 inversion system [38] to realise the parsing semantics. Brabrand et al. [7] present a  
 1315 tool XSugar that handles bijections between the XML syntax (representation) and  
 1316 any other syntax (representation) for the same language, guaranteeing that the  
 1317 syntax transformation is reversible. However, the essential factor that distinguishes  
 1318 our system from others is that the printer produced from a BiYACC program is  
 1319 reflective and can deal with synchronisation.

1320 Although the above-mentioned systems are tailored for unifying parsing and  
 1321 printing, there are design differences. An ISD is more like a parser, while FLIPPR lets  
 1322 the user describe a printer: To handle operator priorities, for example, the user of  
 1323 ISDs will assign priorities to different operators, consume parentheses, and use  
 1324 combinators such as *chainl* to handle left recursion in parsing, while the user of  
 1325 FLIPPR will produce necessary parentheses according to the operator priorities. For  
 1326 basic BiYACC (that deals with unambiguous grammars only), the user defines a  
 1327 concrete syntax that has a hierarchical structure (e.g. *Expr*, *Term*, and *Factor*) to  
 1328 express operator priority, and write printing strategies to produce (preserve)

10FL01 <sup>10</sup> Although they use different implementation techniques, we will not dive into them in our related work.  
 10FL02 See Matsuda and Wang's related work for a comparison [36].

1329 necessary parentheses. The user of XSugar will also likely need to use such a  
1330 hierarchical structure.

1331 It is interesting to note that the part producing parentheses in FLIPPR essentially  
1332 corresponds to the hierarchical structure of grammars. For example, to handle  
1333 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); .
```

1336 FLIPPR will automatically expand the definition and derive a group of ppr<sub>i</sub>  
1337 functions indexed by the priority integer *i*, corresponding to the hierarchical  
1338 grammar structure. In other words, there is no need to specify the concrete grammar,  
1339 which is already implicitly embedded in the printer program. This makes FLIPPR  
1340 programs neat and concise. Following this idea, BiYACC programs can also be made  
1341 more concise: in a BiYACC program, the user is allowed to omit the production rules  
1342 in the concrete syntax part (or omit the whole concrete syntax part), and they will be  
1343 automatically generated by extracting the terminals and nonterminals in the right-  
1344 hand sides of all actions. However, if these production rules are supplied, BiYACC  
1345 will perform some sanity checks: it will make sure that, in an action group, the user  
1346 has covered all of the production rules of the nonterminal appearing in the ‘type  
1347 declaration’, and never uses undefined production rules.

1348 Just like basic BiYACC, all of the systems described above (aim to) handle  
1349 unambiguous grammars only. Theoretically, when the user-defined grammar (or the  
1350 derived grammar) is ambiguous, ISDs’ partial isomorphism could guarantee  
1351 Theorem 1’ by returning Nothing on ambiguous input; FLIPPR’s (own) Theorem 1 is  
1352 comparable to Theorem 1’ by taking all the language constructs which may cause  
1353 non-injective printing into account. However, according to the paper, FLIPPR’s  
1354 Theorem 1 appears to only consider nondeterministic printing based on prettiness  
1355 (layouts). Since the discussion on ambiguous grammars has not been presented in  
1356 their papers, we tested their implementation and the behaviour is as follows: neither  
1357 ISDs nor FLIPPR will notify the user that the (derived) grammar is ambiguous at  
1358 compile time. For ISDs, the right-to-left direction of our Theorem 1’ will fail, while  
1359 for FLIPPR, both directions will fail. (They never promise to handle ambiguous  
1360 grammars, though.) In contrast, Brabrand et al. [7] give a detailed discussion about  
1361 ambiguity detection, and XSugar statically checks if the transformations are  
1362 ‘reversible’. If any ambiguity in the program is detected, XSugar will notify the user  
1363 of the precise location where ambiguity arises. In BiYACC, the ambiguity detection  
1364 of the input grammar is performed by the employed parser generator (currently  
1365 HAPPY), and the result is reported at compile time; if no warning is reported, the  
1366 well-behavedness is always guaranteed. Note that the ambiguity detection can  
1367 produce false positives: warnings only mean that the grammar is not LALR(1) but  
1368 does not necessarily mean that the grammar is ambiguous—ambiguity detection is  
1369 undecidable for the full CFG [10].

1370 Here we also briefly discuss ambiguity detection for the filter approaches: priority  
1371 and associativity (bi-)filters can be applied to (LA)LR parse tables to resolve (shift/

1372 reduce) conflicts [9, 26, 52, 53], and thus the completeness for simple (bi-)filters  
 1373 (see Definition 9) on LALR(1) grammars can be statically checked. However, our  
 1374 implementation does not support it, for bi-filter directives are more general, as  
 1375 stated in the beginning of Sect. 5.4.2, and therefore cannot be transformed to the  
 1376 underlying parser generator's YACC-style directives. Finding a way to directly apply  
 1377 priority and associativity bi-filters to parse tables (generated by HAPPY) is left as  
 1378 future work.

1379 Finally, we compare BiYACC with an industrial tool, AUGEAS, which provides the  
 1380 user with a local configuration API that converts configuration data into a rose tree  
 1381 representation [31]. Similar to BiYACC, AUGEAS also uses the idea of state-based  
 1382 asymmetric lenses so that its *parse* and *print* functions satisfy well-behavedness and  
 1383 it tries to preserve comments and layouts when printing the tree representation back.  
 1384 However, since the purpose of AUGEAS and BiYACC is different, the differences  
 1385 between the tools are also noticeable: (1) AUGEAS works for regular grammars while  
 1386 BiYACC works for (unambiguous) context-free grammars. (2) AUGEAS uses a  
 1387 combinator-based approach while BiYACC adopts a generator-based approach.  
 1388 (3) AUGEAS works more like a simple parser that stops after constructing CSTs: in  
 1389 the parsing direction, AUGEAS unambiguously separates strings into sub-strings, turn  
 1390 sub-strings into tokens, and use tokens to build the corresponding tree; but since  
 1391 each lens combinator (of AUGEAS) has its predefined strategy to turn its  
 1392 acceptable strings into the tree representation, the corresponding tree will be  
 1393 determined once the input string and the lens combinators for parsing the string are  
 1394 given; AUGEAS does not provide a functionality to further transform a tree. On the  
 1395 other hand, BiYACC first turns a string into its isomorphic CST (fully determined the  
 1396 input string and the grammar description) and finally converts the CST to its AST in  
 1397 accordance with the algebraic data types defined by the user; that is, the relation  
 1398 between a string (CST) and its AST is not predetermined but can be adjusted by the  
 1399 user (through printing actions).

## 1400 Generalised Parsing, Disambiguation, and Filters

1401 The grammar of a programming language is usually designed to be unambiguous.  
 1402 Various parser-dependent disambiguation methods such as grammar transformation  
 1403 [29] and parse table conflicts elimination [23] have been developed to guide the  
 1404 parser to produce a single correct CST [26]. On the other hand, natural languages  
 1405 that are inherently ambiguous usually require their parsing algorithms to produce all  
 1406 the possible CSTs; this requirement gives rise to algorithms such as Earley [14] and  
 1407 generalised LR [50] (GLR for short). Although these parsing algorithms produce all  
 1408 the possible CSTs, both their time complexity and space complexity are reasonable.  
 1409 For instance, GLR runs in cubic time in the worst situation and in linear time if the  
 1410 grammar is 'almost unambiguous' [48].

1411 The idea to relate generalised parsing with parser-independent disambiguation  
 1412 for programming languages is proposed by Klint and Visser [26]. They proposed  
 1413 two classes of filters, property filters (defined in terms of predicates on a single tree)  
 1414 and comparison filters (defined in terms of relations among trees), but we only adapt  
 1415 and bidirectionalise predicate filters in this paper. One difficulty lies in the fact that

1416 it is unclear how to define *repair* for comparison filters, as they generally select  
 1417 better trees rather than absolutely correct ones— in the printing direction, since *put*  
 1418 only produces a single CST, we do not know whether this CST needs repairing or  
 1419 not (for there is no other CST to compare). This is also one of the most important  
 1420 problems for our future work.

1421 Parser-independent disambiguation (for handling priority and associativity  
 1422 conflicts) can also be found in LaLonde and des Rivieres's [29] and Aasa's [1]  
 1423 work. At first glance, our *repair* function is quite similar to LaLonde and des  
 1424 Rivieres's post-parse tree transformations that bring a CST into an expression tree,  
 1425 on whose nodes additional restrictions of priority and associativity are imposed. To  
 1426 be simple (but not completely precise), a CST's corresponding expression tree is  
 1427 obtained by first dropping all the nodes constructed from injective productions<sup>11</sup>  
 1428 (note that parentheses nodes are still kept) and then use a precedence-introducing  
 1429 tree transformation to reshape the result. The transformation will do 'repairing' by  
 1430 rotating all the adjacent nodes of the tree where priority or associativity constraint is  
 1431 violated. By contrast, our *repair* function is simpler and only introduces parentheses  
 1432 in places where the *judge* function returns `False`. In short, their tree transformations  
 1433 are a kind of parser-independent disambiguation which does not require generalised  
 1434 parsing; however, those tree transformations are (almost) not applicable in the  
 1435 printing direction if well-behavedness is taken into consideration (due to the rotation  
 1436 of CSTs). Furthermore, it is not clear whether their approach can be generalised to  
 1437 handle other types of conflicts rather than the ones caused by priority and  
 1438 associativity.

1439 There is much research on how to handle ambiguity in the parsing direction as  
 1440 discussed above; conversely, little research is conducted for 'handling ambiguity in  
 1441 the printing direction' and we find only one paper [8] that describes how to produce  
 1442 correct program text regarding priority and associativity, which is also one of the bases  
 1443 of our work. We extend their work [8] by allowing the bracket attribute to work with  
 1444 injective productions such as  $E \rightarrow T; T \rightarrow F; F \rightarrow (' E ') \# \text{Bracket} \#;$ . (The previous  
 1445 work seems to only support the bracket attribute in the form of  
 1446  $E \rightarrow (' E ') \# \text{Bracket} \#;$ ; whether the nonterminal  $E$  on the left-hand side and right-  
 1447 hand side can be different is not made clear.)

1448 Finally, we compare our approach with the conventional ones in general. In  
 1449 history, a printer is believed to be much simpler than a parser and is usually  
 1450 developed independently (of its corresponding parser). While a few printers choose  
 1451 to produce parentheses at every occasion naively, most of them take disambiguation  
 1452 information (for example, from the language's operator precedence table) into  
 1453 account and try to produce necessary parentheses only. However, as the  $Y_{ACC}$ -style  
 1454 conventional disambiguation [23] is parser-dependent, this parentheses-adding  
 1455 technique is also printer-dependent. As the post-parse disambiguation increases the  
 1456 modularity of the (front-end of the) compiler [29], we believe that our post-print  
 1457 parentheses-adding increases the modularity once again. Additionally, the unifica-  
 1458 tion of disambiguation for both parsing and printing makes it possible for us to

11FL01 <sup>11</sup> An injective production, or a chain production, is one whose right-hand side is a single nonterminal;  
 11FL02 for instance,  $E \rightarrow N$ .

1459 impose bi-filter laws, which further makes it possible to guarantee the well-  
1460 behavedness of the whole system.

## 1461 Comparison with a Get-Based Approach

1462 Our work is theoretically based on asymmetric lenses [17] of bidirectional  
1463 transformations [11, 19], particularly taking inspiration from the recent progress on  
1464 putback-based bidirectional programming [15, 27, 28, 40, 41]. As explained in  
1465 Sect. 3, the purpose of bidirectional programming is to relieve the burden of  
1466 thinking bidirectionally—the programmer writes a program in only one direction,  
1467 and a program in the other direction is derived automatically. We call a language  
1468 get-based when programs written in the language denote *get* functions, and call a  
1469 language putback-based when its programs denote *put* functions. In the context of  
1470 parsing and reflecting printing, the get-based approach lets the programmer describe  
1471 a parser, whereas the putback-based approach lets the programmer describe a  
1472 printer. Below we discuss in more depth how the putback-based methodology  
1473 affects BiYACC’s design by comparing BiYACC with a closely related, get-based  
1474 system.

1475 Martins et al. [35] introduces an attribute grammar-based BX system for defining  
1476 transformations between two representations of languages (two grammars). The  
1477 utilisation is similar to BiYACC: The programmer defines both grammars and a set of  
1478 rules specifying a *forward* transformation (i.e. *get*), with a backward transformation  
1479 (i.e. *put*) being automatically generated. For example, the BiYACC actions in lines  
1480 28–30 of Fig. 2 can be expressed in Martins et al.’s system as

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

1482 which describes how to convert certain forms of CSTs to corresponding ASTs. The  
1483 similarity is evident, and raises the question as to how get-based and putback-based  
1484 approaches differ in the context of parsing and reflective printing.

1485 The difference lies in the fact that, with a get-based system, certain decisions on  
1486 the backward transformation are, by design, permanently encoded in the bidirec-  
1487 tionalisation system and cannot be controlled by the user, whereas a putback-based  
1488 system can give the user fuller control. For example, when no source is given and  
1489 more than one rule can be applied, Martins et al.’s system chooses, by design, the  
1490 one that creates the most specialised version. This might or might not be ideal for  
1491 the user of the system. For example: suppose that we port to Martins et al.’s system  
1492 the BiYACC action that relates TIGER’s concrete ‘&’ operator with a specialised  
1493 abstract if expression in Sect. 6.2, coexisting with a more general rule that maps a  
1494 concrete if expression to an abstract if expression. Then printing the AST TCond  
1495 (TSV “a”) (TSV “b”) 0 from scratch will and can only produce a & b, as dictated by the  
1496 system’s hard-wired printing logic. By contrast, the user of BiYACC can easily  
1497 choose to print the AST from scratch as a & b or if a then b else 0 by suitably  
1498 ordering the printing actions.

1499 This difference is somewhat subtle, and one might argue that Martins et al.'s  
 1500 design simply went one step too far—if their system had been designed to respect  
 1501 the rule ordering as specified by the user, as opposed to always choosing the most  
 1502 specialised rule, the system would have given its user the same flexibility as  
 1503 BiYACC. Interestingly, whether to let user-specified rule/action ordering affect the  
 1504 system's behaviour is, in this case, exactly the line between get-based and putback-  
 1505 based design. The user of Martins et al.'s system writes rules to specify a forward  
 1506 transformation, whose semantics is the same regardless of how the rules are ordered,  
 1507 and thus it would be unpleasantly surprising if the rule ordering turned out to affect  
 1508 the system's behaviour. By contrast, the user of BiYACC only needs to think in one  
 1509 direction about the printing behaviour, for which it is natural to consider how the  
 1510 actions should be ordered when an AST has many corresponding CSTs; the parsing  
 1511 behaviour will then be automatically and uniquely determined. In short, relevance of  
 1512 action ordering is incompatible with get-based design, but is a natural consequence  
 1513 of putback-based thinking.

## 1514 Conclusion

1515 We conclude the paper by summarising our contributions:

- 1516 • We have presented the design and implementation of BiYACC, with which the  
 1517 programmer can describe both a parser and a reflective printer for a fully  
 1518 disambiguated context-free grammar in a single program. Our solution  
 1519 guarantees the partial version of the consistency properties (Definition 2) by  
 1520 construction.
- 1521 • We proposed the notion of bi-filters, which enables BiYACC to disambiguate  
 1522 ambiguous grammars while still respecting the consistency properties. This is  
 1523 the main new contribution compared to the previous SLE'16 version [55].
- 1524 • We have demonstrated that BiYACC can support various tasks of language  
 1525 engineering, from traditional constructions of basic machinery such as printers  
 1526 and parsers to more complex tasks such as resugaring, simple refactoring, and  
 1527 language evolution.

1528 **Acknowledgements** We thank the reviewers and the editor for their selflessness and effort spent on  
 1529 reviewing our paper, a quite long one. With their help, the readability of the paper is much improved,  
 1530 especially regarding how several case studies are structured, how theorems for the basic BiYACC and  
 1531 theorems for the extended version handling ambiguous grammars are related, and how look-alike notions  
 1532 are 'disambiguated'. This work is partially supported by the Japan Society for the Promotion of Science  
 1533 (JSPS) Grant-in-Aid for Scientific Research (S) No. 17H06099; in particular, most of the second author's  
 1534 contributions were made when he worked at the National Institute of Informatics and funded by the Grant.  
 1535

## 1536 References

- 1537 1. Aasa, A.: Precedences in specifications and implementations of programming languages. In: Selected  
 1538 Papers of the Symposium on Programming Language Implementation and Logic Programming,  
 1539 Elsevier Science Publishers B. V., Amsterdam, PLILP '91, pp. 3–26. <http://dl.acm.org/citation.cfm?id=203429.203431> (1995)

- 1541 2. Afrozeh, A., Izmaylova, A.: Faster, practical GLL parsing. In: Franke, B. (ed.) *Compiler Construction*, pp. 89–108. Springer, Berlin (2015). [https://doi.org/10.1007/978-3-662-46663-6\\_5](https://doi.org/10.1007/978-3-662-46663-6_5)
- 1542 3. Aho, A.V., Johnson, S.C., Ullman, J.D.: Deterministic parsing of ambiguous grammars. *Commun. ACM* **18**(8), 441–452 (1975)
- 1543 4. Appel, A.W.: *Modern Compiler Implementation in ML*, 1st edn. Cambridge University Press, New York (1998)
- 1544 5. Bird, R.: *Thinking Functionally with Haskell*. Cambridge University Press, Cambridge (2014). <https://doi.org/10.1017/CBO9781316092415>
- 1545 6. Boulton, R.: Syn: a single language for specifying abstract syntax trees, lexical analysis, parsing and pretty-printing. Tech. Rep. Number 390, Computer Laboratory, University of Cambridge (1966)
- 1546 7. Brabrand, C., Møller, A., Schwartzbach, M.I.: Dual syntax for XML languages. *Inf. Syst.* **33**(4–5), 385–406 (2008). <https://doi.org/10.1016/j.is.2008.01.006>
- 1547 8. van den Brand, M., Visser, E.: Generation of formatters for context-free languages. *ACM Trans. Softw. Eng. Methodol.* **5**(1), 1–41 (1996). <https://doi.org/10.1145/226155.226156>
- 1548 9. van den Brand, M.G.J., Scheerder, J., Vinju, J.J., Visser, E.: Disambiguation filters for scannerless generalized LR parsers. In: *Proceedings of the 11th International Conference on Compiler Construction*, Springer, London, UK, CC '02, pp. 143–158. [https://doi.org/10.1007/3-540-45937-5\\_12](https://doi.org/10.1007/3-540-45937-5_12) (2002)
- 1549 10. Cantor, D.G.: On the ambiguity problem of Backus systems. *J. ACM* **9**(4), 477–479 (1962)
- 1550 11. Czarnecki, K., Foster, J.N., Hu, Z., Lämmel, R., Schürr, A., Terwilliger, J.F.: Bidirectional transformations: a cross-discipline perspective. In: *Proceedings of the 2nd International Conference on Theory and Practice of Model Transformations*, Springer, Berlin, ICMT '09, pp. 260–283. [https://doi.org/10.1007/978-3-642-02408-5\\_19](https://doi.org/10.1007/978-3-642-02408-5_19) (2009)
- 1551 12. Dijkstra, E.W.: Guarded commands, nondeterminacy and formal derivation of programs. *Commun. ACM* **18**(8), 453–457 (1975). <https://doi.org/10.1145/360933.360975>
- 1552 13. Duregård, J., Jansson, P.: Embedded parser generators. In: *Proceedings of the 4th ACM Symposium on Haskell*, ACM, New York, NY, USA, Haskell '11, pp. 107–117. <https://doi.org/10.1145/2034675.2034689> (2011)
- 1553 14. Earley, J.: An efficient context-free parsing algorithm. *Commun. ACM* **13**(2), 94–102 (1970). <https://doi.org/10.1145/362007.362035>
- 1554 15. Fischer, S., Hu, Z., Pacheco, H.: The essence of bidirectional programming. *Sci. China Inf. Sci.* **58**(5), 1–21 (2015)
- 1555 16. Foster, J.N.: *Bidirectional programming languages*. PhD thesis, University of Pennsylvania (2009)
- 1556 17. Foster, J.N., Greenwald, M.B., Moore, J.T., Pierce, B.C., Schmitt, A.: Combinators for bidirectional tree transformations: a linguistic approach to the view-update problem. *ACM Trans. Program. Lang. Syst.* **29**, 3 (2007). <https://doi.org/10.1145/1232420.1232424>
- 1557 18. Fowler, M., Beck, K.: *Refactoring: Improving the Design of Existing Code*. Addison-Wesley Professional, Boston (1999)
- 1558 19. Gibbons, J., Stevens, P.: *International Summer School on Bidirectional Transformations* (Oxford, UK, 25–29 July 2016). *Lecture Notes in Computer Science*, vol. 9715. Springer, Berlin (2018)
- 1559 20. Gosling, J., Joy, B., Steele, G.: *The Java Language Specification*, 3rd ed (2006). <https://docs.oracle.com/javase/specs/>
- 1560 21. Hirzel, M., Rose, K.H.: *Tiger language specification* (2013). <https://cs.nyu.edu/courses/fall13/CSCI-GA.2130-001/tiger-spec.pdf>
- 1561 22. Hu, Z., Ko, H.S.: 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*, pp. 100–150. Springer, Cham (2018). [https://doi.org/10.1007/978-3-319-79108-1\\_4](https://doi.org/10.1007/978-3-319-79108-1_4)
- 1562 23. Johnson, S.C.: Yacc: Yet another compiler-compiler. AT&T Bell Laboratories Technical Reports (AT&T Bell Laboratories Murray Hill, New Jersey 07974). p. 32 (1975)
- 1563 24. Kernighan, B.W., Ritchie, D.M.: *The C Programming Language*. Prentice Hall Press, Upper Saddle River (1988)
- 1564 25. Kinoshita, D., Nakano, K.: 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 (2017)



- 1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652
26. Klint, P., Visser, E.: Using filters for the disambiguation of context-free grammars. In: Pighizzini, G., Pietro, P.S. (eds) Proceedings of the ASMICS Workshop on Parsing Theory, University of Milan, Italy, Milano, Italy, pp. 1–20 (1994)
  27. Ko, H.S., Hu, Z.: An axiomatic basis for bidirectional programming. *Proc. ACM Program. Lang.* **2**(POPL), 41:1–41:29 (2018). <https://doi.org/10.1145/3158129>
  28. Ko, H.S., Zan, T., Hu, Z.: 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 (2016). <https://doi.org/10.1145/2847538.2847544>
  29. LaLonde, W.R., des Rivieres, J.: Handling operator precedence in arithmetic expressions with tree transformations. *ACM Trans. Program. Lang. Syst.* **3**(1), 83–103 (1981). <https://doi.org/10.1145/357121.357127>
  30. Lämmel, R., Jones, S.P.: 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 (2003). <https://doi.org/10.1145/604174.604179>
  31. Lutterkort, D.: Augeas—a configuration API. In: Proceedings of the Ottawa Linux Symposium, Ottawa, Canada, pp. 47–56 (2008)
  32. Macedo, N., Pacheco, H., Cunha, A., Oliveira, J.N.: Composing least-change lenses. *Proc. Sec. Int. Workshop Bidirect. Transform.* **57**, 1–19 (2013). <https://doi.org/10.14279/tuj.eceasst.57.868>
  33. Marlow, S., Gill, A.: The parser generator for Haskell. <https://www.haskell.org/happy/> (2001)
  34. Marlow, S., et al.: Haskell 2010 language report. <https://www.haskell.org/onlinereport/haskell2010/> (2010)
  35. Martins, P., Saraiva, J., Fernandes, J.P., Van Wyk, E.: 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 (2014). <https://doi.org/10.1145/2543728.2543745>
  36. Matsuda, K., Wang, M.: 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 (2018a). <https://doi.org/10.1145/3242744.3242758>
  37. Matsuda, K., Wang, M.: FliPpr: a system for deriving parsers from pretty-printers. *New Gener. Comput.* **36**(3), 173–202 (2018b). <https://doi.org/10.1007/s00354-018-0033-7>
  38. Matsuda, K., Mu, S.C., Hu, Z., Takeichi, M.: A grammar-based approach to invertible programs. In: Gordon, A.D. (ed) Proceedings of the 19th European Conference on Programming Languages and Systems, Springer, Berlin, no. 20 in ESOP'10, pp. 448–467 (2010). [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.: Towards a practical programming language based on dependent type theory. PhD thesis, Chalmers University of Technology (2007)
  40. Pacheco, H., Hu, Z., Fischer, S.: 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 (2014a). <https://doi.org/10.1145/2543728.2543737>
  41. Pacheco, H., Zan, T., Hu, Z.: 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 (2014b). <https://doi.org/10.1145/2643135.2643141>
  42. Pombrio, J., Krishnamurthi, S.: 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 (2014). <https://doi.org/10.1145/2594291.2594319>
  43. Pombrio, J., Krishnamurthi, S.: 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 (2015). <https://doi.org/10.1145/2784731.2784755>
  44. Rendel, T., Ostermann, K.: 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 (2010). <https://doi.org/10.1145/1863523.1863525>

- 1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684
45. Reps, T., Teitelbaum, T.: 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 (1984). <https://doi.org/10.1145/800020.808247>
  46. Reps, T., Teitelbaum, T., Demers, A.: Incremental context-dependent analysis for language-based editors. *ACM Trans. Program. Lang. Syst.* **5**(3), 449–477 (1983). <https://doi.org/10.1145/2166.357218>
  47. Scott, E., Johnstone, A.: GLL parsing. *Electron. Notes Theor. Comput. Sci.* **253**(7), 177–189 (2010). <https://doi.org/10.1016/j.entcs.2010.08.041>
  48. Scott, E., Johnstone, A., Economopoulos, R.: BRNGLR: a cubic tomita-style GLR parsing algorithm. *Acta Inform.* **44**(6), 427–461 (2007). <https://doi.org/10.1007/s00236-007-0054-z>
  49. Sheard, T., Jones, S.P.: 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 (2002). <https://doi.org/10.1145/581690.581691>
  50. Tomita, M.: 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 (1985). <http://dl.acm.org/citation.cfm?id=1623611.1623625>
  51. Traver, V.J.: On compiler error messages: what they say and what they mean. *Ad. Hum. Comput. Interact.* **2010**, 3:1–3:26 (2010). <https://doi.org/10.1155/2010/602570>
  52. Visser, E.: A case study in optimizing parsing schemata by disambiguation filters. *International Workshop on Parsing Technology (IWPT 1997)*, pp. 210–224. Massachusetts Institute of Technology, Boston, USA (1997a)
  53. Visser, E.: Syntax definition for language prototyping. PhD thesis, University of Amsterdam (1997b)
  54. Younger, D.H.: Recognition and parsing of context-free languages in time  $n^3$ . *Inf. Control* **10**(2), 189–208 (1967)
  55. Zhu, Z., Zhang, Y., Ko, H.S., Martins, P., Saraiva, J., Hu, Z.: 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. <https://doi.org/10.1145/2997364.2997369> (2016)
- Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## 1685 Affiliations

1686 Zirun Zhu<sup>1</sup>  · Hsiang-Shang Ko<sup>2</sup>  · Yongzhe Zhang<sup>1</sup>  · Pedro Martins<sup>3</sup> ·  
1687 João Saraiva<sup>4</sup>  · Zhenjiang Hu<sup>5</sup> 

1688 ✉ Zirun Zhu  
1689 zhu@nii.ac.jp

1690 Hsiang-Shang Ko  
1691 joshko@iis.sinica.edu.tw

1692 Yongzhe Zhang  
1693 zyz915@nii.ac.jp

1694 Pedro Martins  
1695 pribeiro@uci.edu

1696 João Saraiva  
1697 jas@di.uminho.pt

1698 Zhenjiang Hu  
1699 huzj@pku.edu.cn

- 1700 <sup>1</sup> National Institute of Informatics, Chiyoda, Tokyo, Japan  
1702 <sup>2</sup> Institute of Information Science, Academia Sinica, Taipei, Taiwan  
1703 <sup>3</sup> University of California, Irvine, USA  
1704 <sup>4</sup> Department of Informatics, University of Minho, Braga, Portugal  
1705 <sup>5</sup> Department of Computer Science and Technology, Peking University, Beijing, China  
1706

REVISED PROOF