# Robust error correction in a generalized LR parser

## STEPHEN LEAKE, retired

 Error correction in a parser is important when the parser is used in an interactive development environment (IDE); the source code is often not syntactically correct. Several extensions to work by McKenzie, Yeatman, and De Vere [15] are described, including 'Minimal\_Complete', which quickly provides the missing end of any grammar production. The algorithm has been in production use in Emacs ada-mode for over a year; metrics show that the algorithm is useful.<sup>1</sup>

CCS Concepts: • Software and its engineering  $\rightarrow$  Parsers; Translator writing systems and compiler generators; • Humancentered computing  $\rightarrow$  User interface programming; • Computing methodologies  $\rightarrow$  Shared memory algorithms.

### ACM Reference Format:

Stephen Leake. 2020. Robust error correction in a generalized LR parser. 1, 1 (July 2020), 16 pages. https://doi.org/10.1145/nnnnnnnnnnnnn

## 1 INTRODUCTION

Emacs Ada mode has used an LR parser to support indentation, syntax highlighting, and navigation since 2013 [13]. However, the parser did not provide error correction, so indentation was often confusing when the syntax was incorrect, as it usually is in an interactive editing environment. This motivated the search for error correction algorithms.

Grune and Jacobs [11] provides a thorough overview of error correction algorithms in LR and LL parsers. Of those, the work by McKenzie, Yeatman, and De Vere [15] provides the foundation for the current work.

The McKenzie algorithm works by exploring the parse table (or Deterministic Parsing Automata (DPA) as Mckenzie et al. [15] calls it) at the error point, finding tokens to insert. It also tries deleting tokens following the error point. Each possible solution, together with the parse stack at the error point, forms a *configuration*. Each configuration also has a cost, determined by what tokens are inserted and deleted. At each step in the algorithm, new configurations are generated from the current error point. Then the minimum cost configuration is checked to see if it succeeds; if not, more configurations are generated.

There are several situations where the McKenzie algorithm can take a long time, or is inefficient. For example, consider this Ada code:

procedure Example\_01

## is begin

Msg : constant String;

begin

Put\_Line (Msg);

end;

<sup>1</sup>This work was supported in part by Eurocontrol

Author's address: Stephen Leake, stephen\_leake@stephe-leake.orgretired.

Unpublished working draft. Not for distribution.

© 2020 Convright held by the owner/author(s)

Manuscript submitted to ACM

2020-07-24 19:10. Page 1 of 1-16. Manuscript submitted to ACM

53 There is an extra 'begin' immediately after 'is' (a common occurance while editing code). However, the error is not 54 detected until ':' after 'Msg', which can only occur in declarations, not statements (in Ada, declarations occur between 55 'in' and 'begin'; statements between 'begin' and 'end'). 56 To fix this, the McKenzie algorithm must insert 57 58 '; end; begin IDENTIFIER', or delete 'constant String; begin' and then insert ';'. A much better solution would be 59 to "push back" 'begin Msg', and then delete 'begin'. 60 A harder problem is when there are several missing "end"s, because the user is typing a nested statement: 61 62 63 begin 64 if A then 65 66 B ; 67 Ation. if C then 68 69 loop 70 Do\_Something; 71 end; 72 73 74 Fig. 1. missing many 'end's 75

Here we are missing several tokens: 'end loop; end if; end if;'. McKenzie will theoretically find the solution that inserts all of these, but along the way it will try inserting every possible statement as well, wasting a lot of time, and in practice hitting a time-out limit. This article introduces the 'Minimal\_Complete' algorithm, which quickly finds the minimal number of tokens to insert to complete the trailing statement in situations like this.

This article also introduces several other new operations for the core McKenzie algorithm to try at the error point, and adapts the algorithm to work with a generalized parser.

Unless quoting from another source, we use the notation described by DeRemer and Pennello [6] or Aho, Sethi, and Ullman [3].

## 2 THE PARSING CONTEXT

## 2.1 Generalized parsing

We use a generalized LR parser Tomita [16], to tolerate conflicts in the parse table. This allows using an Ada grammar that is close to the one given in the ISO Ada language standard appendix P [9], which is not LR(1). Using that grammar as faithfully as possible ensures we are implementing the correct language, and simplifies updating the parser to a new language version.

This means the main parser may have several parallel parsers executing when an error is encountered. The McKenzie error correction algorithm is enhanced to maintain state information for each parser. When more than one solution is 98 99 found with the same cost, additional parsers are created to use them.

100 The parser builds a syntax tree; to handle parallel parsers, we use a branched tree similar to Tomita's sub-tree 101 sharing. The syntax tree will contain "virtual tokens" that are inserted by error correction. 102

The entire input token sequence is kept in memory, to allow arbitrary push back.

104 Manuscript submitted to ACM

2020-07-24 19:10. Page 2 of 1-16.

2

76 77

78

79

80 81

82

83

84

85 86

87 88

> 89 90

> 91

92

93 94

95

96

97

## 2.2 Partial parse

105 106

107

108 109

110

111

112

113 114

115

116

117

118 119

120

121 122 123

124

125 126 127

128

129

130

131 132

133

134 135 136

137

138 139

140

141 142 143

144

145 146

147

148

In order to handle very large files, the parser supports "partial parse"; parsing only part of a file.

Note that this is not the same as "incremental parse", where an existing syntax tree from a previous parse is modified based on a text edit. Integrating this error correction algorithm with incremental parse is the subject of future work.

In order to minimize the amount of text passed from the editor to the parser, we modify the grammar to allow smaller chunks of code to be accepted as a complete parse. In Ada, this means adding 'declaration' and 'statement' to 'compilation unit'.

When the file length is greater than a threshold, Emacs invokes a partial parse whenever a parse is needed. It first uses a regular expression search to find a reasonable start point, then finds a possible matching end that includes the requested parse position, and passes that region to the parser.

In Ada, the search for a start point finds a block begin, or the point after a block end. If a block begin was found, the matching end is looked for; otherwise, the requested parse point is the end point; 'Minimal Complete' will provide Warning war the missing tokens.

# **3 EXTENSIONS TO MCKENZIE**

We define the following operations that are tried in each McKenzie step:

- 'push\_back'
- 'undo\_reduce'
- 'Try\_Insert\_Quote'
- 'Minimal\_Complete'
- 'Language\_Matching\_Begin\_Tokens'
- 'Language\_Fixes'

#### 3.1 push\_back

'push\_back' pops the top parse stack item, and moves the input stream pointer back to the first terminal contained by that item. We call the point in the input stream at which insert and delete is done the "edit point"; it may not be an error point. 'push\_back' moves the edit point.

### 3.2 undo reduce

'undo\_reduce' undoes the reduce that produced the top stack item (which must be a nonterminal), replacing the top stack item by the sequence of stack items just before the reduction, without moving the edit point. This requires a syntax tree that records the shifts and reductions done during the parse (a "concrete" syntax tree). This operation serves two purposes;

149 150 151

152

153 154

155 156

- (1) It allows a subsequenct 'push\_back' to push back fewer tokens.
- (2) It allows token insertions that would otherwise be forbidden by the grammar.
- To illustrate the second point, consider:

2020-07-24 19:10. Page 3 of 1-16.

```
Stephen Leake
```

```
157
       procedure Example_2
158
       i s
159
160
           I : Integer;
161
       begin
162
           procedure Put_Top_10
163
164
           is begin
165
            . . .
166
           end Put_Top_10;
167
168
       begin
169
       end Example_2;
170
171
172
                                                              Fig. 2
173
174
175
         There is an extra 'begin' after 'I : Integer'. The error is detected at 'procedure'; at that point, the parse stack
176
       looks like (top is to the left):
177
       245 : BEGIN, 208 : declarative_part, 159 : IS, 36 : subprogram_specification, 0 :
178
179
         Here the numbers label the states, terminals are in uppercase, nonterminals in lowercase.
180
         The grammar productions relevant to this example are:
181
182
       declarative_part <= declarations |__;
183
       declarations <= declarations declaration | declaration ;</pre>
184
       declaration <= subprogram_declaration | ... ;</pre>
185
186
         Fixing the parse error starts by 'push_back BEGIN, delete BEGIN', leaving ' declarative_part ' on top of the stack.
187
       'procedure' is the next token, which is illegal in state 208; it starts a 'subprogram_declaration', but we've already
188
189
       "closed" the declaration section by reducing to ' declarative_part '. We could do 'push_back declarative_part ', but
190
       that moves the edit point to before the object declaration for 'I', where there is no error and nothing helpful to insert
191
       or delete. Instead, 'undo_reduce' leaves the stack as:
192
193
       137 : declarations, 159 : IS, 36 : subprogram_specification, 0 :
194
       and now inserting 'procedure' is legal.
195
196
         We do not maintain a syntax tree for the parsing done during error correction; doing that proved to be much too slow.
197
       Therefore the 'undo_reduce' operation can only be applied to configurations where the top stack item was produced
198
       by the main parse, so it has a valid syntax tree entry. An exception is when the nonterminal is empty; that is easy to
199
       undo.
200
201
202
       3.3 Try_Insert_Quote
203
       Missing string quotes cause problems for the McKenzie algorithm. Consider the code:
204
205
      A : String := Now is the time for all good men";
206
207
208
       Manuscript submitted to ACM
                                                                                                 2020-07-24 19:10. Page 4 of 1-16.
```

There is a missing quote before 'Now'. In Ada, strings cannot cross newline, and the lexer handles the error by inserting a virtual quote just before the existing one. Then the parser sees a list of identifiers followed by an empty string literal.

The McKenzie algorithm would have to delete all the identifiers one by one, with a cost for each.

The 'Try\_Insert\_Quote' operation attempts to find a better place to insert the string quote, depending on the relative placement of the unbalanced quote and the parse error.

• If the parse error is at the unbalanced quote, assume the unbalanced quote is the intended closing quote, and insert the opening quote one non-empty token before it. Example:

A := "for\_all" & good ";

We are in the process of splitting a string across lines; we just added '"\_&', but are missing the '"' before 'good'. This solution inserts that missing quote.

• If the parse error is after the unbalanced quote, assume the unbalanced quote is the intended opening quote, and insert the closing quote at the line end. Example:

A := "for\_all" & "good

The missing '"' should be after 'good'. This solution inserts that missing quote.

• If the parse error is before the unbalanced quote, assume the unbalanced quote is the intended closing quote, and insert the opening quote in several places (generate one new configuration for each):

– before the error token. Example:

A := for all good ";

The missing '"' should be before '**for**'. The parse error is at '**all**'; this solution inserts the missing quote before '**all**', which is almost right.

- One non-empty token before the unbalanced quote. Same example, but this inserts the '"' before 'all', which
  is correct.
- If there is a string literal on the parse stack, assume the closing quote of that string literal is new (or extra), push back thru the token containing that string literal, and extend the string literal to the unbalanced quote.
   Example:

A := "for\_all" good";

The '"' after 'all' is extra. The parse error is at 'good'; this solution in effect deletes the extra quote.

Note that the search for a string literal on the parse stack must take into account nonterminals that may contain a string literal. The set of nonterminals that may contain a string literal is provided as a function call written by the grammar author; it is not computed from the grammar because it should not include higher level nonterminals that are not likely to be contained in a string; for Ada, it stops at expression.

Since the lexer recognizes string literals, we cannot actually insert an unbalanced string quote; we actually delete all tokens between the inserted quote and the unbalanced quote, which matches what the lexer would have returned. This is still much better than the original McKenzie algorithm, because we do all the deletions in one step, with one low cost.

260 2020-07-24 19:10. Page 5 of 1-16.

### 3.4 Minimal Complete 261

The 'Minimal Complete' strategy is to insert the minimum number of tokens to finish the current grammar production. 263 This is supported by precomputing a set of 'Minimal\_Complete\_Action' for each parser state. 264

Section 4 gives the algorithm for computing the 'Minimal\_Complete\_Actions'; here we describe how they are used in error correction.

Each 'Minimal Complete Action' is either "insert this terminal token" or "reduce to this nonterminal token".

Consider the code in figure 1. When parsing this code, an error is detected at the final ';'; the parser is expecting 'end loop;'. The kernel of that state has one production:

loop\_statement <= LOOP sequence\_of\_statements END ^ LOOP</pre>

identifier\_opt SEMICOLON

where the caret (^) shows the parse point (also known as "dot" in an LR(1) item). In this case, it is easy to see that 'Minimal\_Complete\_Action' must be "insert 'loop' ". Similarly, after parsing 'loop', 'Minimal\_Complete\_Action' is "reduce to ' identifier opt ' ", and then "insert ';' ". After that, we will be completing the inner 'if then' statement, and then the outer 'if then' statement. At that point, the existing 'end ;' is legal.

To handle cases where only part of a nonterminal production needs to be inserted, the check to see if the original error token is now legal must be performed after each token is inserted; this means that 'Minimal\_Complete' inserts at most one token for each cycle of the underlying McKenzie algorithm.

In order to prefer the 'Minimal\_Complete' solution over others, we give it a negative cost. In figure 1, the grammar is a small subset of Ada, the default insert and delete cost is 4, delete 'begin' is cost 1, delete 'end', ';' are cost 2. With 'Minimal\_Complete' cost 0, -1, or -2, no solution for figure 1 is found, with an enqueue limit of 120,000. With 'Minimal Complete' cost -3, the desired solution is found with cost 9, after enqueueing 6804 configurations and check-288 ing 1051. 289

As usual, there is a trade off here; sometimes other solutions would be better than 'Minimal\_Complete'. For example, consider:

## for I in 1 to Result\_Length loop

end loop; 294

Here 'to' should be '.. ' (this is an actual error the author typed late one night). The error is detected at 'to', so the 296 297 desired solution is

298 'delete IDENTIFIER, insert ...'. With the 'Minimal\_Complete' cost set to 0, this solution is found, along with 19 other 299 solutions with the same cost; a few of these change the code to: 300

301 in to.Result\_Length loop 302

```
in 1 .. to * Result Length loop
303
```

```
304
     in 1 .. to/Result_Length loop
```

Finding these takes a relatively long time; 1273 configurations were enqueued and 218 checked.

307 Setting 'Minimal\_Complete' cost to -1 finds similar solutions, but more quickly (341 enqueued, 56 checked); 'insert ...' 308 is done by 'Minimal\_Complete', so it is cheaper, and more expensive solutions are not checked. 309

Setting 'Minimal Complete' cost to -3 (as required for figure 1) finds one cost 3 solution very quickly (67 enqueued, 310 311 10 checked); it changes the code to:

312 Manuscript submitted to ACM

2020-07-24 19:10. Page 6 of 1-16.

262

265

266

267

268

269 270

271

272

273 274

275

276

277 278

279

280

281

282 283

284

285

286

287

290

291 292

293

295

305

```
for I in 1 .. to loop Result_Length; loop
313
314
       end loop;
315
316
      which causes another syntax error later in the code (missing 'end loop;'). Here three tokens were inserted by 'Minimal_Complete';
317
      '.. loop ;'.
318
         In the production Ada parser, -3 proves to be a good compromise for
319
      'Minimal Complete' cost.
320
321
         In some cases, the action required for 'Minimal_Complete' is reduce, not shift. For example:
322
       case Current Token is
323
324
       = +Right_Paren_ID then
325
           Matching_Begin_Token := +Left_Paren_ID;
326
327
       else
328
           Matching_Begin_Token := Invalid_Token_ID;
329
       end if;
330
331
      Here the user is in the middle of converting an 'if' statement to a 'case' statement. The relevant grammar productions
332
      are:
333
334
       if_statement <=</pre>
335
           IF expression_opt THEN sequence_of_statements
336
           elsif_statement_list ELSE sequence_of_statements
337
338
           END IF SEMICOLON
339
      case_statement <=</pre>
340
           CASE expression_opt IS case_statement_alternative_list
341
           END CASE SEMICOLON
342
343
      An error is detected at '='; 'Minimal_Complete' inserts
344
      'when NUMERIC LITERAL =>', then the original McKenzie algorithm inserts
345
      'if NUMERIC_LITERAL', which makes the remaining code legal, terminating one error correction session. Then pars-
346
347
      ing proceeds to then end of the input, where another error is enountered; missing 'end case;'. At that point, the parse
348
      state kernel has one production:
349
       if_statement <= IF expression_opt THEN sequence_of_statements
350
351
          ELSE sequence_of_statements END IF SEMICOLON ^;
352
353
      Minimal_Complete_Action => if_statement
354
355
      Since the 'Minimal_Complete_Action' token is a nonterminal, the action is reduce, not shift. That leads to several more
356
       states where the
357
      'Minimal_Complete_Action' is reduce:
358
359
      compound_statement <= if_statement ^</pre>
360
      Minimal_Complete_Action => compound_statement
361
362
363
      statement <= compound_statement ^</pre>
364
      2020-07-24 19:10. Page 7 of 1-16.
                                                                                               Manuscript submitted to ACM
```

```
365
       Minimal_Complete_Action => statement
366
367
368
369
         and finally arrives at the state:
370
       case_statement <= CASE expression_opt IS</pre>
371
          case_statement_alternative_list ^ END CASE SEMICOLON
372
373
       case_statement_alternative_list <=</pre>
374
          case_statement_alternative_list ^ case_statement_alternative
375
376
       Minimal_Complete_Action => END
377
378
       which inserts 'end'. 'Minimal_Complete' does all required reductions, and one insertion, in one McKenzie step (similar
379
       to McKenzie insert).
380
         If there is more than one production in the kernel for a parse state that gives the minimum length for that state,
381
382
       there is more than one
383
      'Minimal_Complete_Action' for that state. In that case, we look at the length after the parse point for each production
384
       in the kernel for the state that the shift or reduce goes to; if one of the actions results in a minimum length, that action
385
       is chosen. If more than one action gives the minimum length, all are kept. Therefore 'Minimal_Complete' maintains a
386
       queue of configurations with an action to check. Each may produce a new configuration for the next McKenzie step.
387
388
         Sometimes the minimal length production cannot be computed at compile time. Consider the Java code fragment:
389
390
       { B = UpdateText (A }
391
392
393
                                                             Fig. 3
394
395
396
         This is missing '); ' after 'A'. The error is detected at '}'. At that point, the kernel is:
397
       LambdaExpression <= Identifier ^ MINUS_GREATER Identifier
398
       LeftHandSide <= Identifier ^
399
       ClassType <= Identifier ^
400
401
       MethodInvocation <= Identifier ^ LEFT_PAREN ArgumentList RIGHT_PAREN
402
403
       Minimal_Complete_Action => (MINUS_GREATER, LeftHandSide, ClassType)
404
405
       Here two actions are reduce, and therefore there are no tokens in these productions after the parse point, so we have
406
       to look at the next state to find the minimal length production. Reducing to 'LeftHandSide' goes to the state:
407
408
       Assignment <= LeftHandSide ^ EQUAL Expression
409
      where there are two tokens needed to complete the production.
410
         Reducing to 'ClassType' goes to the state:
411
412
       PostfixExpression <= ClassType ^</pre>
413
       ClassType <= ClassType ^ DOT Identifier
414
415
       Once again, we reduce 'ClassType' to the next state, and the next, until we find a shift. This ends in the state:
```

```
416 Manuscript submitted to ACM
```

2020-07-24 19:10. Page 8 of 1-16.

Robust error correction in a generalized LR parser

```
MethodInvocation <= Identifier LEFT_PAREN ArgumentList ^ RIGHT_PAREN
417
418
       ArgumentList <= ArgumentList ^ COMMA Expression
419
       which has 1 token after the parse point; this is the minimal length. All of these reductions, and the final shift, are done
420
421
       in one McKenzie step. The function that computes the length after parse point is recursive, and explores all paths that
422
       might be minimal.
423
         Sometimes 'Minimal_Complete' results in an ambiguous parse. Consider the following Java code:
424
425
       { UpdateText (A) }
426
427
      In the subset of Java we are using for this example, this is not a complete statement. The relevant productions are:
428
       StatementExpression
                                    <= PostIncrementExpression | PostDecrementExpression ;</pre>
429
       PostIncrementExpression <= PostfixExpression '++' ;</pre>
430
431
       PostDecrementExpression <= PostfixExpression '--' ;</pre>
432
       PostfixExpression
                                    <=
433
          ClassType | MethodInvocation | PostIncrementExpression
434
        | PostDecrementExpression ;
435
436
                                   <= Identifier | ClassType '.' Identifier ;
       ClassType
437
      To complete the code fragment, we have to insert either '++' or '--'. The error is detected at the closing brace; the
438
       state is:
439
440
       PostIncrementExpression <= PostfixExpression ^ PLUS_PLUS
441
      PostDecrementExpression <= PostfixExpression ^ MINUS_MINUS
442
       Minimal_Complete_Action => (PLUS_PLUS, MINUS_MINUS)
443
444
      There are two actions, neither reduce, so we enqueue both. The next McKenzie step inserts ';' in each, completing the
445
       error recovery session with two solutions. In the main parser, both solutions parse to end of input. In the absence of
446
       error correction multiple parallel parsers getting to end of input is an error (ambiguous parse), but since error correction
447
448
       often produces multiple solutions, we don't report an error; the parser with an error solution that has minimum cost
449
       and minimum recover ops length is chosen to be the final parse.
450
         In general, we can rely on cost to exclude cycles in error recovery, but this is not true for 'Minimal Complete', since
451
       it has negative cost. We will see in section 4 that minimal actions that might lead to a cycle are dropped at compile
452
453
       time, so there is no need for detecting cycles at runtime.
454
455
       3.5 Matching_Begin
456
       Consider the Ada code fragment:
457
458
       procedure ... end;
459
460
461
           end if;
462
       end Foo;
463
464
      This is a result of cut and paste. The code is missing
465
       'procedure Foo is begin if expression then' before 'end if'.
466
467
         The error is detected at the 'end' in 'end if'. 'Minimal_Complete' is no help here; the parse stack looks like:
468
       2020-07-24 19:10. Page 9 of 1-16.
                                                                                                   Manuscript submitted to ACM
```

34 : subprogram\_body, 0 :

'Minimal\_Action' in state 34 is reduce to 'compilation\_unit', which is then reduced to ' compilation\_unit\_list ', which
is the grammar start symbol, which has no minimal complete action. So 'Minimal\_Complete' has nothing useful to
insert.

To the grammar author, the solution is obvious. We capture that knowledge by having the grammar author provide a function

'Language\_Matching\_Begin\_Tokens', which takes the current parse stack and the next three input tokens as input, and
 returns a list of tokens to try inserting at the edit point.

The 'Language\_Matching\_Begin\_Tokens' function is similar to the table

E(A, a) given by Fischer, Milton, and Quiring in [7] for an LL parser:

$$E(A, a) \equiv x \mid A \Rightarrow^* xay, \operatorname{Cost}(x) \text{ is minimized}$$
(1)

However, we are using an LR parser, so that table is not directly applicable. We could try to compute a table that gives the minimal sequence of tokens to insert starting in any state *s*, and allowing any token *a* as the next token:

490

491

492

493 494

498 499

500

480

481 482

483 484

485

 $M(s, a) \equiv y | A \Rightarrow^{*} xyaz, A \in Kernel(s), Cost(x) \text{ is minimized}$   $\tag{2}$ 

where x is the prefix of the production A in state s. That is the table that the McKenzie algorithm computes on the fly; it is not worth precomputing. It is worth providing 'Language\_Matching\_Begin\_Tokens' to shortcut some common situations.

In Ada, three tokens are required to determine the proper match for 'end'; consider:

```
495 case is end case;
```

<sup>496</sup> loop end loop;

block\_1 : **begin end** Block\_1;

```
package Parent_1.Child_1 is begin end Parent_1.Child_1;
```

In the case and loop statements, the three tokens starting with 'end' are 'end loop ;' and 'end case ;'; here we only need two tokens to determine that the matching begin is 'case' or 'loop'. To distinguish between the named block statement and the package declaration, we need three tokens; in a named block statement the name must be a simple identifier, with no dots, so the third token must be ';'.

In Ada, 'Language\_Matching\_Begin\_Tokens' returns the proper statement start token for 'end ... ', and for 'then, else, elsif, excepting
 For 'when', it returns 'case IDENTIFIER is', which assumes a partial case statement is more common than a partial
 exception handler. For any other error token, it returns an empty token list; no guess is better than a bad guess.

'Language\_Matching\_Begin\_Tokens' also returns a flag

<sup>511</sup> 'Forbid\_Minimal\_Complete', which is True when 'Minimal\_Complete' would be harmful. In Ada, this is set True when <sup>512</sup> the error point is after '**end**' in one of the '**end**' sequences above; it is better to push back '**end**'.

514

510

## 515 3.6 Language\_Fixes

To take advantage of the redundant block name information in Ada, we provide a general hook 'Language\_Fixes'; it takes as input a configuration, the parse table, and the syntax tree and token stream for one parser. It enqueues new configurations to test. Consider:

520 Manuscript submitted to ACM

2020-07-24 19:10. Page 10 of 1-16.

10

Robust error correction in a generalized LR parser

```
521 procedure Proc_1
522 is begin
```

```
Block_1:
begin
```

end Proc\_1;

524 525

526

527 528

529

530

531 532

533

534 535

536

537

538

539 540

541

542

543 544

545

546

547

548 549

550

551 552

553

554 555

556

557 558

559

560

561 562

563

564 565

566

567 568

569

570

Here 'end Block\_1;' is missing. The grammar rules for Ada do not require the start and end block names to match; that is checked later in the compilation process. To take advantage of it for error correction, we add that check as a parse-time action in the grammar declaration:

```
block_statement <=</pre>
```

```
block_label_opt BEGIN handled_sequence_of_statements END
identifier_opt SEMICOLON
%()%
%(return Match_Names
      (Lexer, Descriptor, Tokens, 1, 5, End_Name_Optional);)%
```

Here the first '%()%' gives the post-parse action, which is run after the parse is complete and the syntax tree is available; Emacs uses this action to compute indent, navigation, and highlight. The second '%()%' gives the in-parse action, which is run when the production is reduced, both in the main parse and during error correction.

```
Here 'Match_Names' will return a status of 'Match_Names_Error' if the names do not match, with the error point after the final ';', and the production not reduced (so it is possible to edit the token sequence without an 'undo_reduce'). 'Language_Fixes' then tries to determine the best fix based on the name information.
```

In this example, the error is reported after 'end Proc\_1;'.

```
'Language_Fixes' finds the matching 'procedure Proc_1' on the parse stack, and inserts 'end Block_1;' before 'end Proc_1;'. Consider:
```

```
package Parent_1.Child_1
```

is begin

```
begin
```

```
end Parent_1.Child_1;
```

Here we might expect 'Match\_Names' will fail with 'Extra\_Name\_Error', but instead we get a parse error on '.' in 'Parent\_1.Child\_1'; block names must be simple identifiers. Since this is similar to 'Match\_Names\_Error', 'Language\_Fixes' handles it, finding the matching name; the fix is to insert '**end**;' before

```
'end Parent 1.Child 1;'.
```

Consider:

```
Block_1 :
```

# begin

```
if A then
```

```
null;
```

**end** Block\_1;

```
572 2020-07-24 19:10. Page 11 of 1-16.
```

Here we get a syntax error on the final 'Block_1'; ' <b>if</b> ' is expected. Again, 'Language_Fixes' handles it, finding metabing name	the
More complex patterns can be recognized in 'Language_Fixes'; see the production code for the Ada parser.	
4 COMPUTING MINIMAL_COMPLETE	
4.1 Grammar recursions	
We first compute the recursions in the grammar, so we can exclude cycles from 'Minimal_Complete'. A 'recursic	n' is
the result of a cycle in the grammar productions; for example:	
association_list <=	
association_list COMMA association_opt	
association_opt	
The first right hand side (RHS) is direct left recursive; the second is not recursive. Recursion can also be indirect; consider:	
name <=	
IDENTIFIER	
selected_component	
selected_component <= name DOT IDENTIFIER	
Together, 'name, selected_component' are indirect recursive.	
We can form a graph representing the grammar by taking the nonterminals as the graph vertices, and if there	is a
production $A \Rightarrow xBy$ then there is a directed edge from the A to B. Then recursion is represented by a cycle ir	the
graph.	
In a useful grammar, every recursion must have a non-recursive RHS in one of the nonterminals, to terminate	the
recursion.	
We use Johnson's algorithm [12] to find the cycles in the graph. However, that algorithm does not apply to ' $n$	ulti-
graphs', which have more than one edge connecting any two nodes. Real grammars can be multigraphs, so we	first
filter out all such edges, and add them back after we compute the cycles.	
Some languages have a lot of recursion, so it can be prohibitive to compute the exact set of cycles in the g	aph.
For example, Java SE 12 (using the grammar given in chapter 19 of the language reference manual [8]) takes too	ong
to compute. In that case, we only compute the strongly connected components (SCCs) (which is one of the step	s in

Johnson [12]), and use that as the recursion. This gives more recursion than necessary, and thus reduces the number of minimal actions computed, but still gives good performance in error correction. For Ada 2012, computing the full recursion is fast, but for Ada 2020 draft 25 [10], it is prohibitively slow. 

Once we have the cycles or SCCs, we add a Boolean 'recursive' flag to the items in the state kernels; true if the production is in a cycle or SCC and the recursive token is the first token (ie the production is left recursive, direct or indirect), false otherwise. 

Manuscript submitted to ACM

2020-07-24 19:10. Page 12 of 1-16.

## 4.2 Other preliminaries

We also need the minimal terminal token sequence for each production. This is the same as S(A) from Fischer et al. [7]:

$$S(A) \equiv x \in V_t^* \mid A \Rightarrow^* x, \operatorname{Cost}(x) \text{ is minimized}$$
(3)

where the cost of each token is 1. Including individual token costs at this point would make it very difficult to assign useful costs; we only use token costs in the run-time portion of the algorithm.

Next we need 'Minimal\_Terminal\_First (A)', which is the first token in S(A), or the invalid token  $\xi$  if S(A) is empty. Finally we need  $Nullable(A, \omega)$ , which gives the production that reduces A to  $\epsilon$ :

$$Nullable(A,\omega) \equiv (B,\psi) \in P \mid (A \Longrightarrow B\gamma, B \Longrightarrow^* \epsilon) \ else \ \xi \tag{4}$$

We say a nonterminal token 'A' is "nullable" if  $Nullable(A, \omega)$  returns  $\xi$  for some  $\omega$ .

## 4.3 Minimal Complete Action

For each state, we consider each item in the kernel. There are several possible cases, listed in priority order:

- (0) There is only one item in the kernel. Recursion is ignored, because any other McKenzie operation also has no choice here. The minimal action is given by 'Compute\_Action (Dot)' (see figure 4), where 'Dot' is the token after dot in the kernel item.
- (1) Dot is at the end of the production, or all tokens following dot are nullable; the actual production length must be computed at runtime. If Dot is at the end of the production, recursion is ignored again because any other McKenzie operation also has no choice here; the minimal action is reduce to the item LHS. If Dot is not at the end of the production, we cannot ignore recursion; the null token might be in a recursion cycle.
- (2) The item is left recursive; there is no minimal action.
- (3) There is no recursion, and dot is not at the end of the production. If the number of tokens after dot is minimal within the kernel, include 'Compute\_Action (Dot)' in the minimal actions.

```
657
     function Compute_Action (Token)
658
         if Token in Terminals then
659
660
            return Action (State, Token);
661
         else
662
            if Minimal Terminal First (Token) = Invalid Token then
663
664
                return (Reduce 0 tokens to Token);
665
            else
666
667
                return Action (State, Minimal_Terminal_First (Token));
668
            end if;
669
         end if;
670
671
     end Compute_Action;
672
```

'Action (State, Token)' returns the parse action for the token in the state.

676 2020-07-24 19:10. Page 13 of 1-16.

### 5 IMPLEMENTATION

The algorithm presented here is implemented in the WisiToken parser [14], and used in Emacs ada-mode.

McKenzie uses a bit map to prune redundant (but presumably higher cost) configurations from the queue; instead, we use a Fibonacci min heap (Cormen, Leiserson, and Rivest Cormen et al. [5] chapter 19), so finding the minimal cost configuration is asymptotically free.

Checking each configuration to see if it is a solution, and generating new configurations from it, is independent of all other configurations, so the underlying McKenzie algorithm is easily parallelized. We use one Ada protected object to store the min heap of configurations to check, and one Ada task per processor to check and generate configurations. However, that results in only 40% speedup with 8 processors; more work is needed in this area.

The configure data structure is optimized for speed; it has a bounded parse stack of 70 items, and a bounded vector of 80 recover operations; creating a new configuration on the min heap requires only one fixed length memory allocation. Hitting either of those limits is not at all likely in a low cost solution, so we simply drop any configurations that do so.

We use an LR1 parse table, not LALR. The extra information retained by the LR1 parse table is helpful in error recovery. Using the same stress-test case as the parallel task speed test, with the LALR parser there are 4 parsers active when recovery is entered, they enqueue a total of 1\_094\_252 configurations (two hit the enqueue limit of 500\_000), and check a total of 48\_035 configurations. With the LR1 parser, there are two parsers active when recovery is entered, they enqueue 558\_102 configurations (one hit the enqueue limit; a gain of almost 50%) and check 47\_530. On the other ant in reco hand, due in part to using multiple tasks, the total time spent in recovery is about the same between the LALR and LR1 parsers; 1.10 seconds in this case.

# 6 RESULTS

This error correction algorithm has been in production use in Emacs ada-mode since November 2018. The parser can be configured to output information about each error recovery; that is summarized in figures 5 and 6 for one month of the author's use. The enqueue limit is set at 58 000; there is a noticeable delay when that limit is hit. The maximimum enqueue value in the table is higher because multiple tasks are used in error recovery; each task is allowed to finish its current operation before aborting due to the limit.

This shows that the minimal complete operation is used in a large majority of cases, along with insert.

Manuscript submitted to ACM

2020-07-24 19:10. Page 14 of 1-16.

8		
	count	percent
fail enqueue limit	1_126	3%
ignore_error	2_130	1%
language_fix	11_138	6%
minimal_complete	151_977	78%
matching_begin	2_848	1%
push_back	13_694	7%
undo_reduce	7_491	4%
insert	141_629	72%
delete	18_557	9%
string_quote	2_615	1%
	mean	max
enqueue	796.3	60_147
check	61.8	5496
		P
		4: .   £:
. o. Error recovery stati	istics for par	tial file pars
	count	percent
fail enqueue limit	16	0%
ignore_error	850	5%
		Y

Fig. 5. Error recovery statistics for full file parsing

rig. o. Error recovery statistics for partial file parsi	Fig. 6.	Error recover	y statistics for	partial file	parsing
--	---------	---------------	------------------	--------------	---------

	count	percent
fail enqueue limit	16	0%
ignore_error	850	5%
language_fix	356	2%
minimal_complete	14_579	77%
matching_begin	2_355	13%
push_back	1_019	5%
undo_reduce	544	3%
insert	6_076	32%
delete	1_266	7%
string_quote	11	0%
	mean	max
enqueue	922.6	49_286
check	78.5	4201

> To compare our error correction to other parsers, we use a set of 59 Ada source files with known errors - this is the set of files used to test Emacs indentation etc. We write code to use the parsers to output the corrected string of tokens found for each file. Then we difference that string from the nominal correct string, using the 'diff ' program. The length of the diff output is then a measure of error correction quality (shorter is better).

2020-07-24 19:10. Page 15 of 1-16.

libadalang is a parser provided by AdaCore [2], used in their GNAT Studio IDE (although not yet for indentation; see AdaCore [1]). The results of comparing wisitoken with libadalang are given in figure 7. Here "perfect files" is the count of files where the corrected token stream is the same as the nominally correct token stream; "better files" means the diff is shorter than the other algorithm. Overall, WisiToken does a better job, although there are 7 files where libadalang does a better job.

	total	perfect	better
	diff size	files	files
wisitoken	13_097	23	49
libadalang	27_935	4	7

Tree-sitter [4] is an incremental parser, designed for use in IDEs. However, it cannot cope with the full Ada grammar, so we used the Ada subset grammar to compare error correction. The test files were edited to conform to the subset grammar; three became meaningless after that and were dropped. The results are in figure 8

ig. 8. error compa	arison metri	ics with sub	set Ada g
	total	perfect	better
	diff size	files	files
wisitoken	20_450	12	51
tree-sitter	47 047	1	4

## REFERENCES

- [1] AdaCore. Adacore gnat studio, 2020. URL https://www.adacore.com/gnatpro/toolsuite/gnatstudio.
- 813
   [2] AdaCore. Adacore libadalang, 2020. URL https://github.com/AdaCore/libadalang.
- [3] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. Compilers: Principles, Techniques, and Tools. Addison-Wesley, 1985.
- [4] Max Brunsfeld. tree-sitter home page, 2020. URL https://github.com/tree-sitter/tree-sitter.
  - [5] Thomas H. Cormen, Charles Eric. Leiserson, Ronald L. Rivest, and Clifford Stein. Introduction to Algorithms. MIT Press, third edition, 2009.
- [6] Frank Deremer and Thomas Pennello. Efficient computation of lalr(1) look-ahead sets. ACM Transactions on Programming Languages and Systems
   (TOPLAS), 4(4):615-649, 1982. doi: 10.1145/69622.357187.
  - [7] C. N. Fischer, D. R. Milton, and S. B. Quiring. Efficient ll(1) error correction and recovery using only insertions. Acta Informatica, 13(2):141–154, 1980. doi: 10.1007/bf00263990.
  - [8] James Gosling, Bill Joy, Guy Steele, Gilad Bracha, Alex Buckley, and Daniel Smith. Java language grammar, 2019. URL https://docs.oracle.com/javase/specs/jls/se12/html/index.html.
  - [9] Ada Working Group. Ada 2012 language standard, 2012. URL http://ada-auth.org/standards/ada12.html.
- [5] Ada Working Group. Ada 2020 anguage standard, 2012. URL http://ada/aduti.org/standards/ada/2x.html.
   [10] Ada Working Group. Ada 2020 draft 25, 2020. URL http://www.ada-auth.org/standards/ada/2x.html.
- [11] Dick Grune and Ceriel J.H Jacobs. *Parsing Techniques; A Practical Guide*. Springer Science, second edition, 2008.
- [12] Donald B. Johnson. Finding all the elementary circuits of a directed graph. SIAM Journal on Computing, 4(1):77–84, 1975. doi: 10.1137/0204007.
- 826 [13] Stephen Leake. Ada mode homepage, 2020. URL http://www.nongnu.org/ada-mode/NEWS-ada-mode.text.
- [14] Stephen Leake. Wisitoken home page, 2020. URL https://stephe-leake.org/ada/wisitoken.html.
- [15] Bruce J. Mckenzie, Corey Yeatman, and Lorraine De Vere. Error repair in shift-reduce parsers. ACM Transactions on Programming Languages and
   Systems (TOPLAS), 17(4):672–689, 1995. doi: 10.1145/210184.210193.
- 830 [16] Masaru Tomita. Efficient parsing for natural language. Kluwer Academic Publishers, 1986.

832 Manuscript submitted to ACM

2020-07-24 19:10. Page 16 of 1-16.